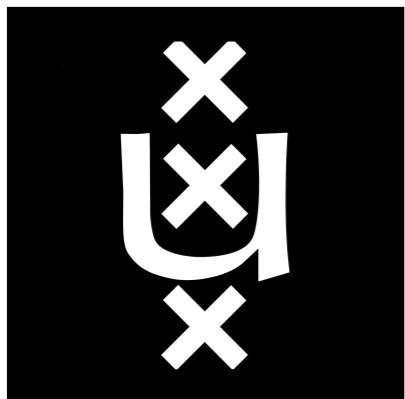
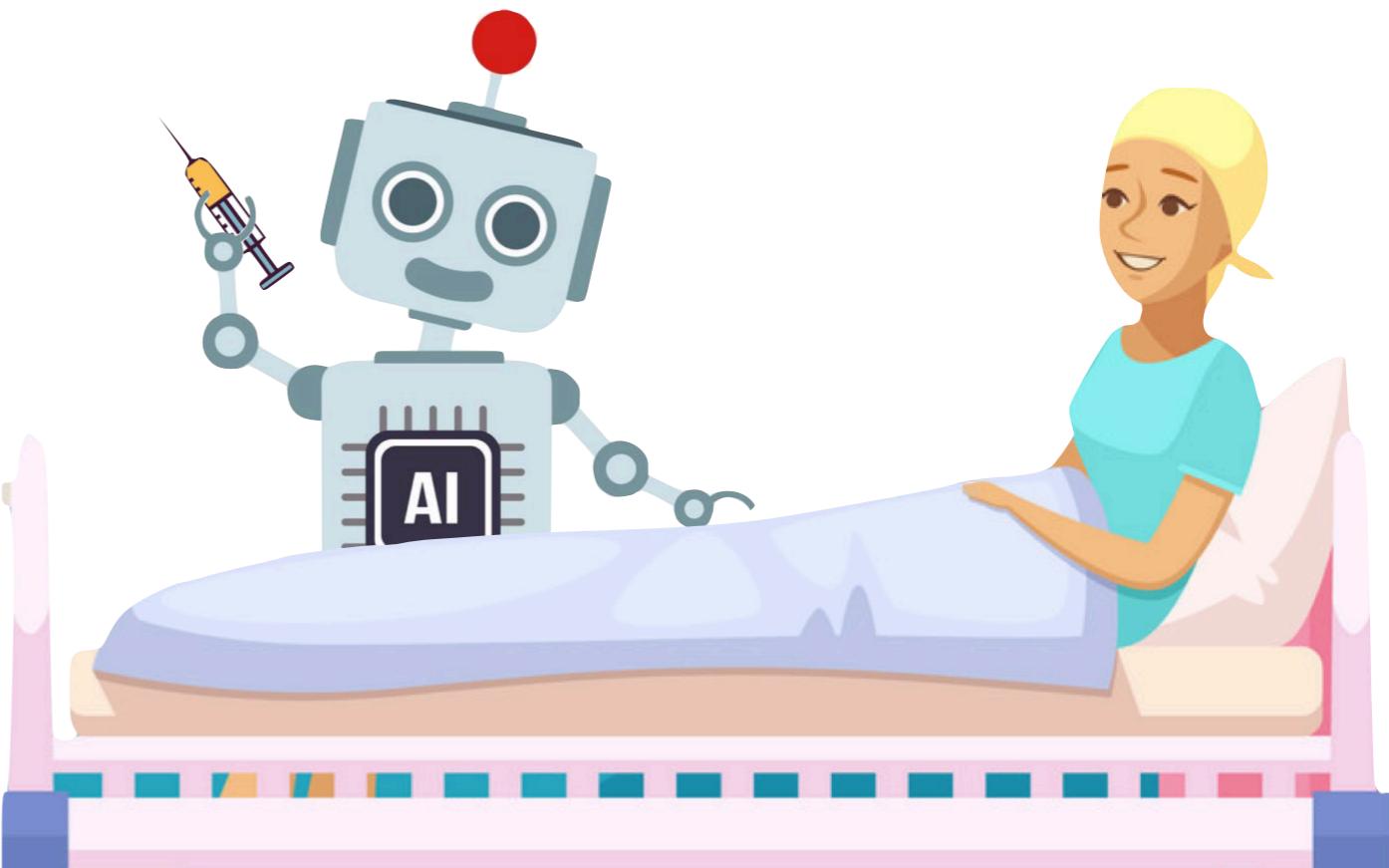
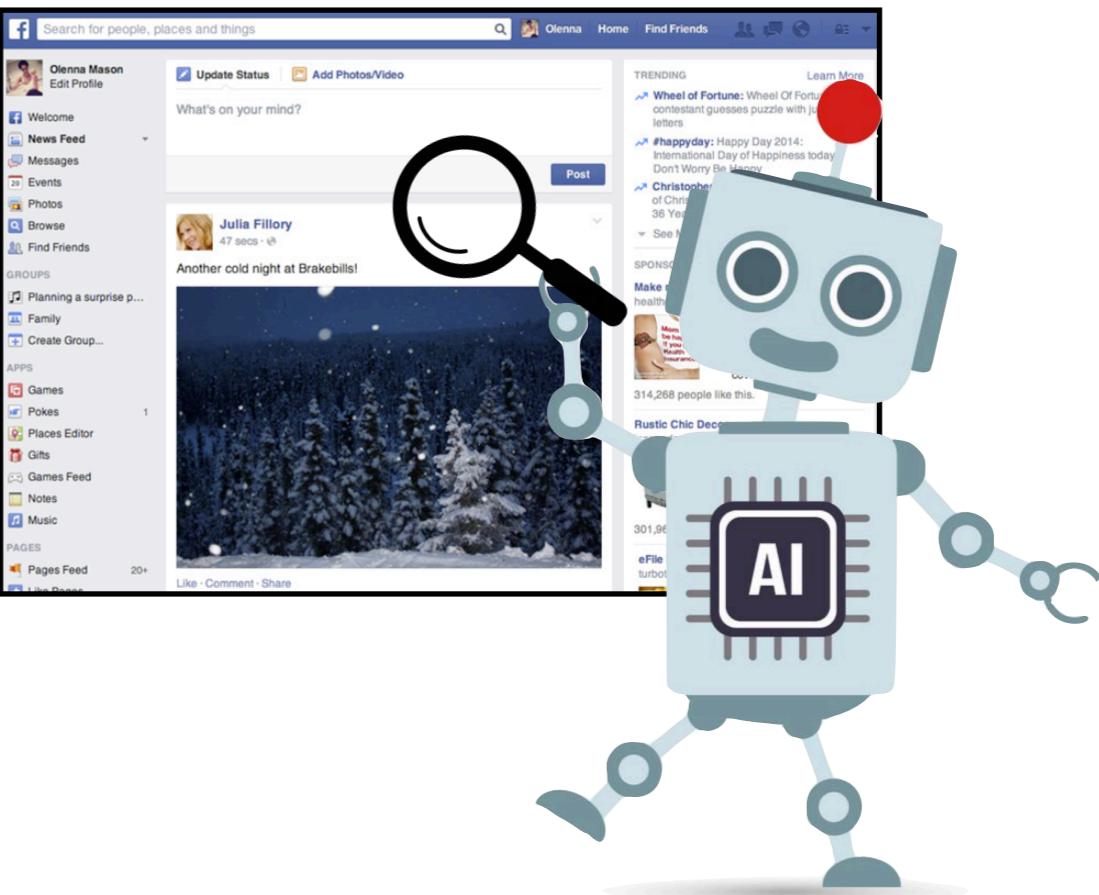
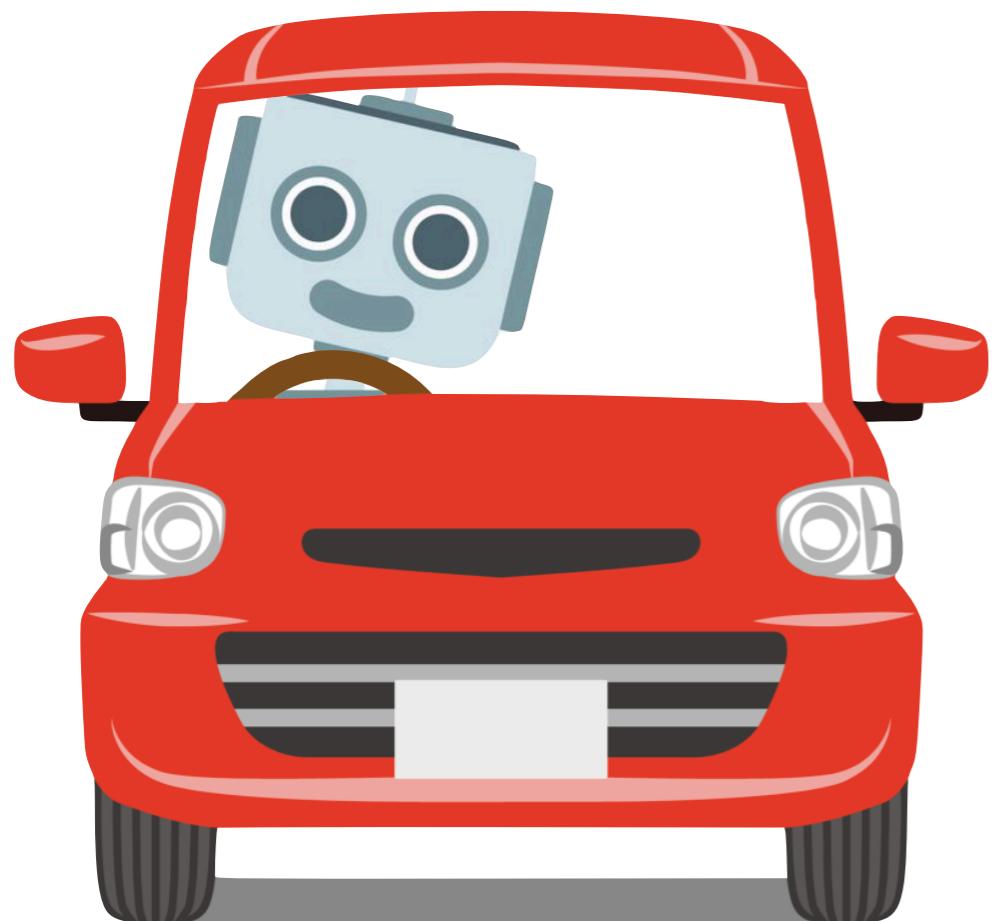


Learning to Defer to One, Multiple, or a Population of Expert(s)

Eric Nalisnick

University of Amsterdam





Los Angeles Times

A Tesla mystery: Why didn't auto-braking stop these crashes?

The image shows a scene of multiple vehicle collisions on a highway. In the foreground, a dark-colored Tesla sedan has crashed into the side of a white and blue emergency response vehicle, which appears to be a California State Trooper's car. Other vehicles are visible in the background, and several people are standing near the accident site.

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)

[Facebook](#) [Twitter](#) [LinkedIn](#) [Print](#) [Email](#)



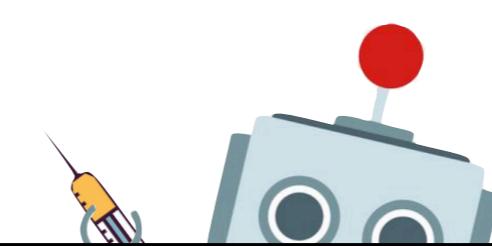
Amhara militia members ride in the back of a truck towards a fight with the Tigray People's Liberation Front. Reuters/Tiksa Negeri

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





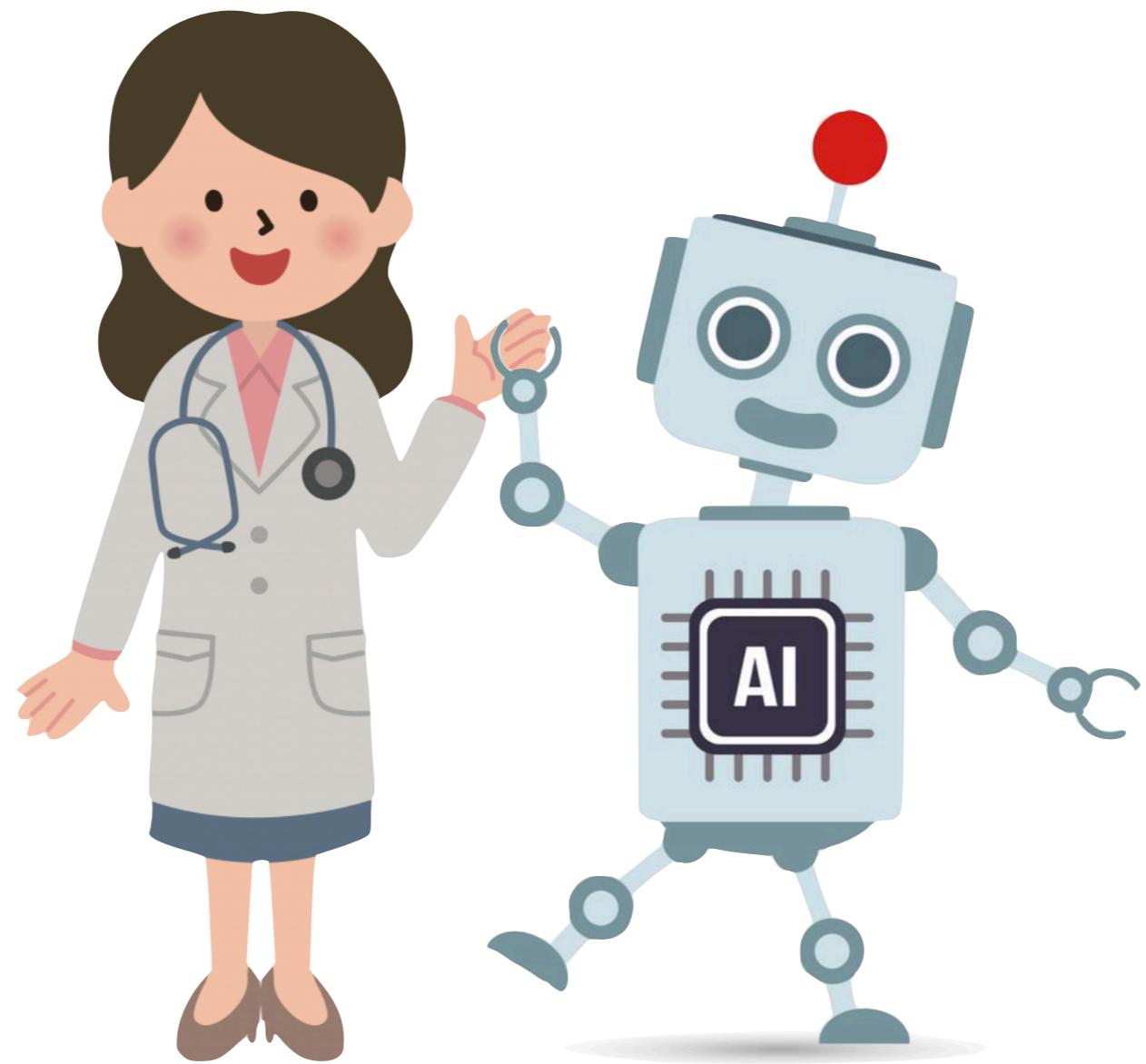
Medscape Tuesday, December 13, 2022

NEWS & PERSPECTIVE DRUGS & DISEASES CME & EDUCATION ACADEMY VIDEO DECISION POINT

News > Medscape Medical News > Conference News > CHEST 2022

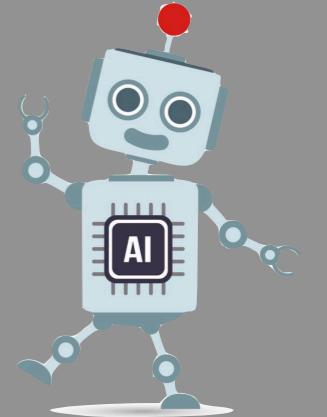
Sepsis Predictor Tool Falls Short in Emergency Setting

Heidi Splete
October 17, 2022



human-AI collaboration

input
features



classifier



expert

learning to defer (to an expert)

input
features



allocation
mechanism

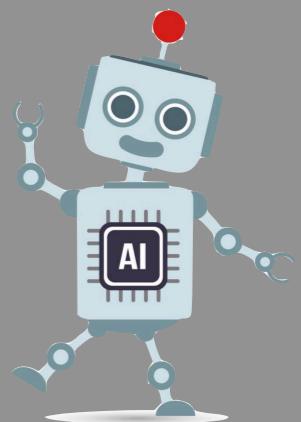
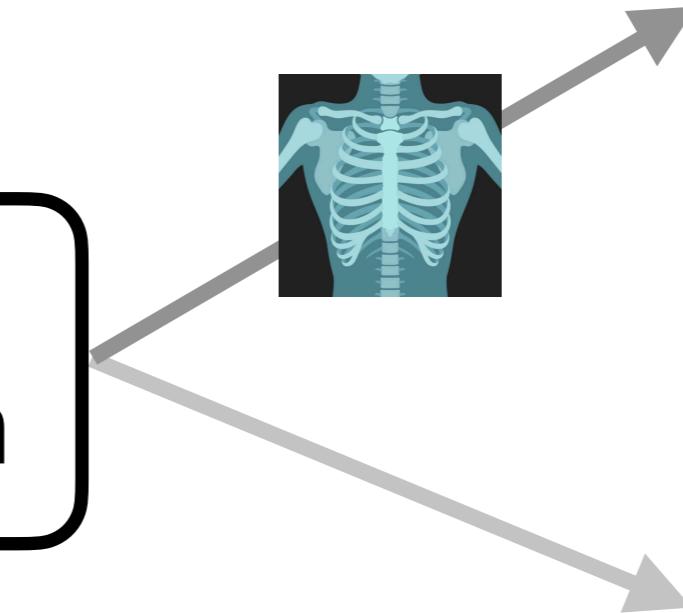


learning to defer (to an expert)

input
features



allocation
mechanism



classifier



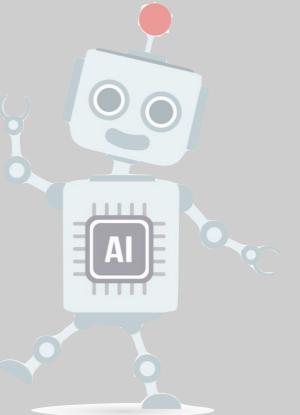
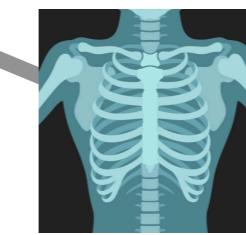
expert

learning to defer (to an expert)

input
features



allocation
mechanism



classifier



expert

learning to defer (to an expert)

input
features

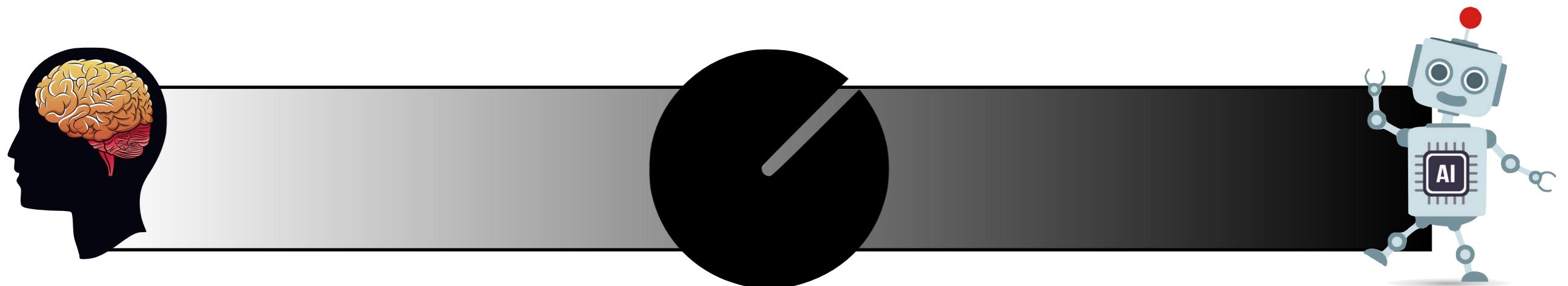


allocation
mechanism



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

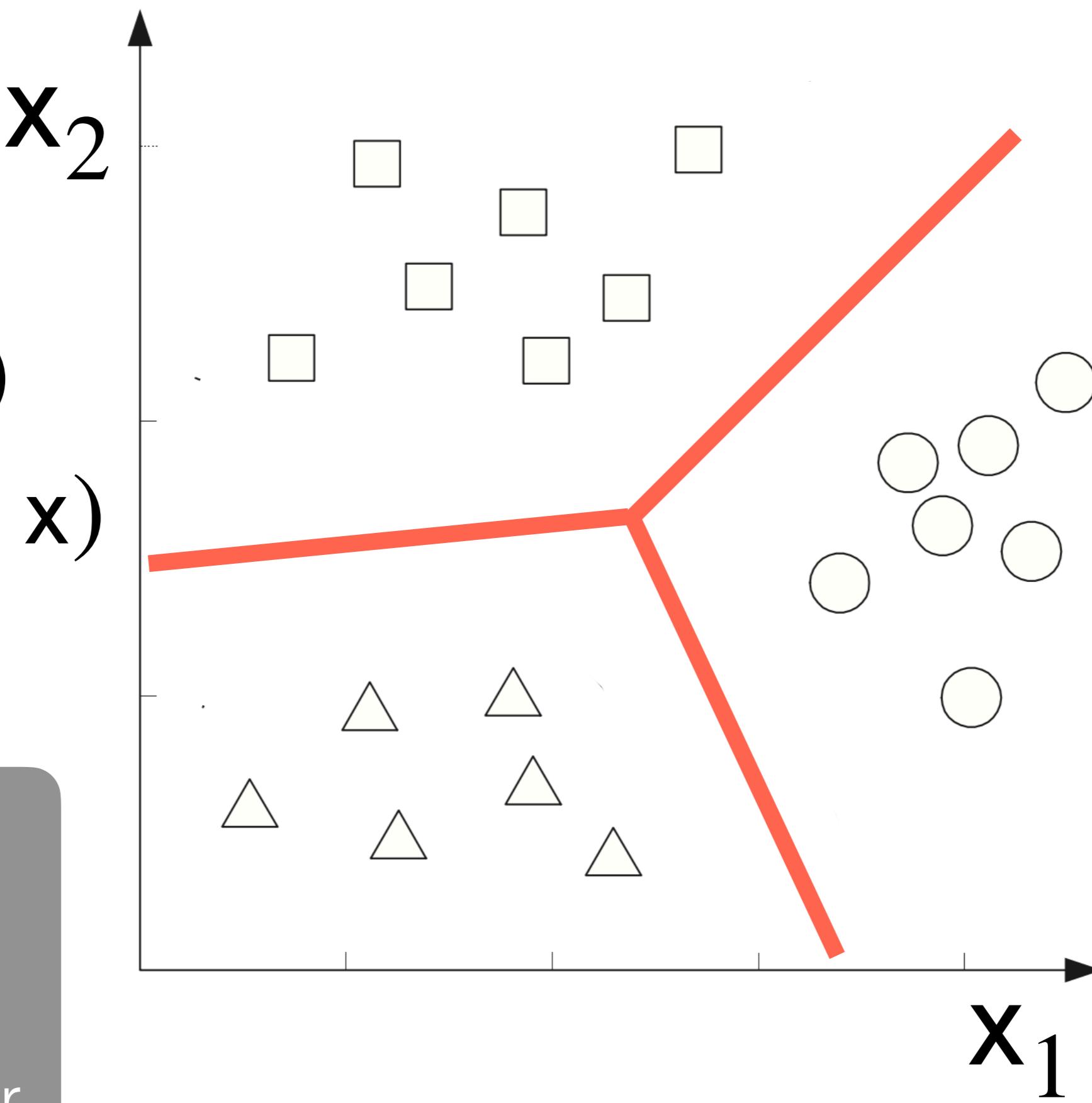
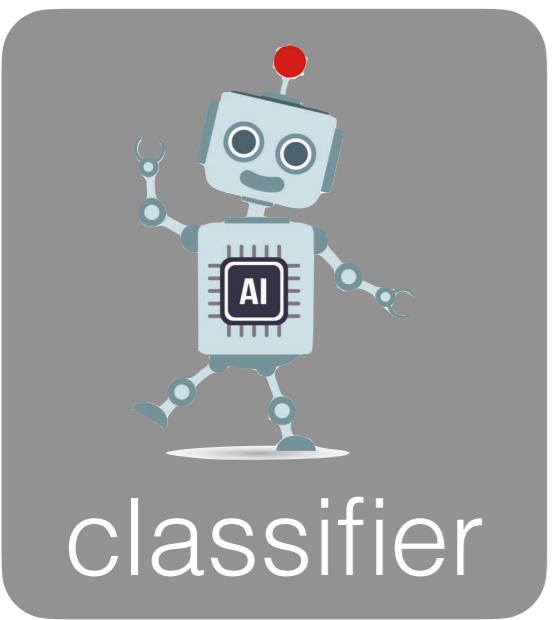
input
features



allocation
mechanism



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



input
features



allocation
mechanism

classifier



expert

input
features

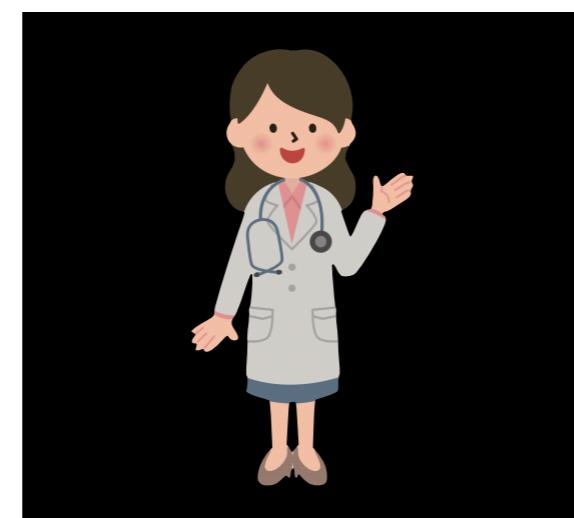


allocation
mechanism



input
features

X



m
prediction

(black box)



input
features



allocation
mechanism

???



classifier

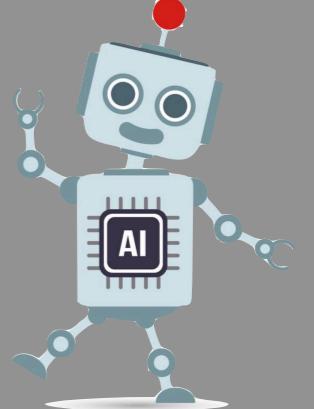


expert

input
features



allocation
mechanism



classifier

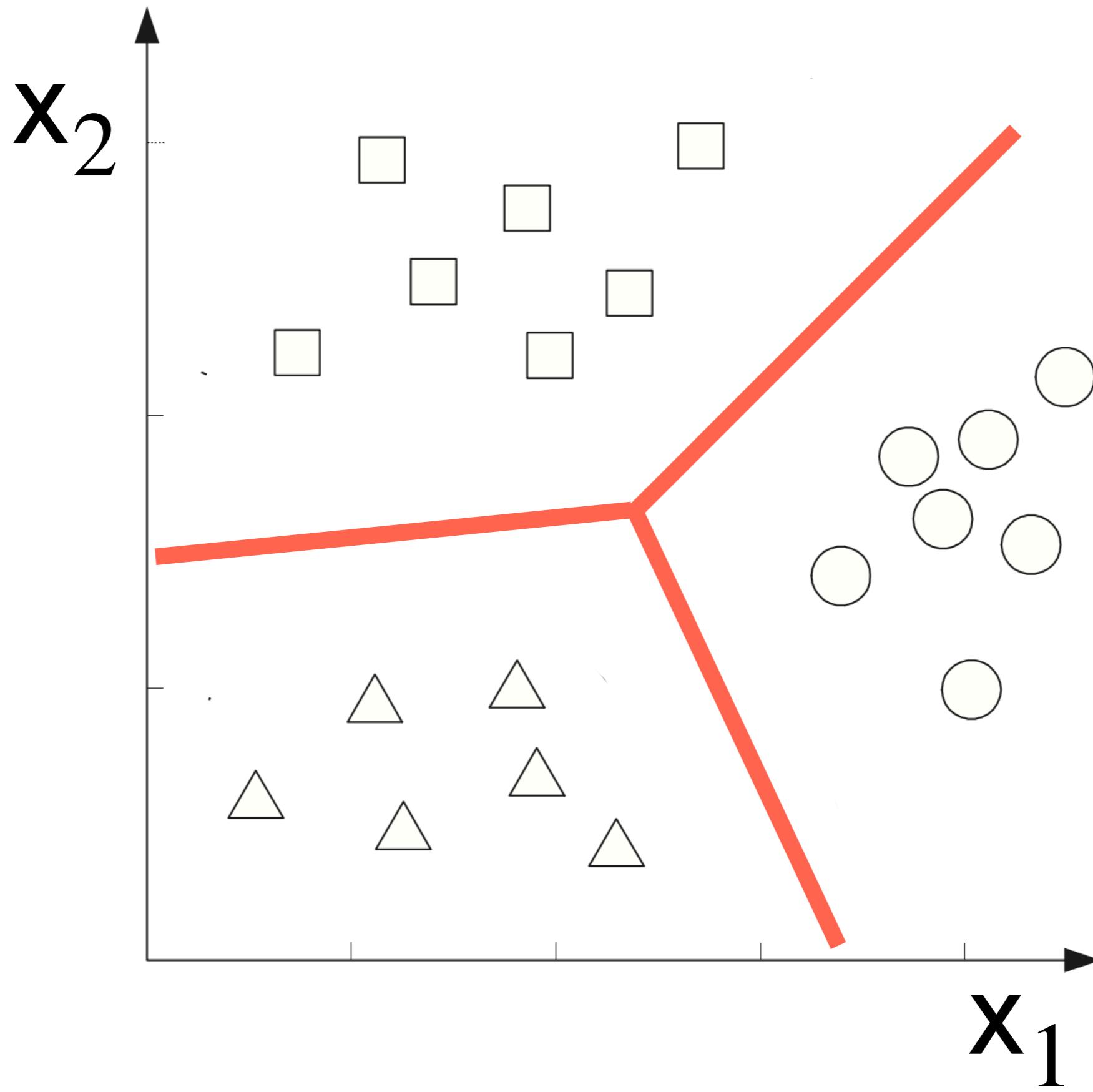


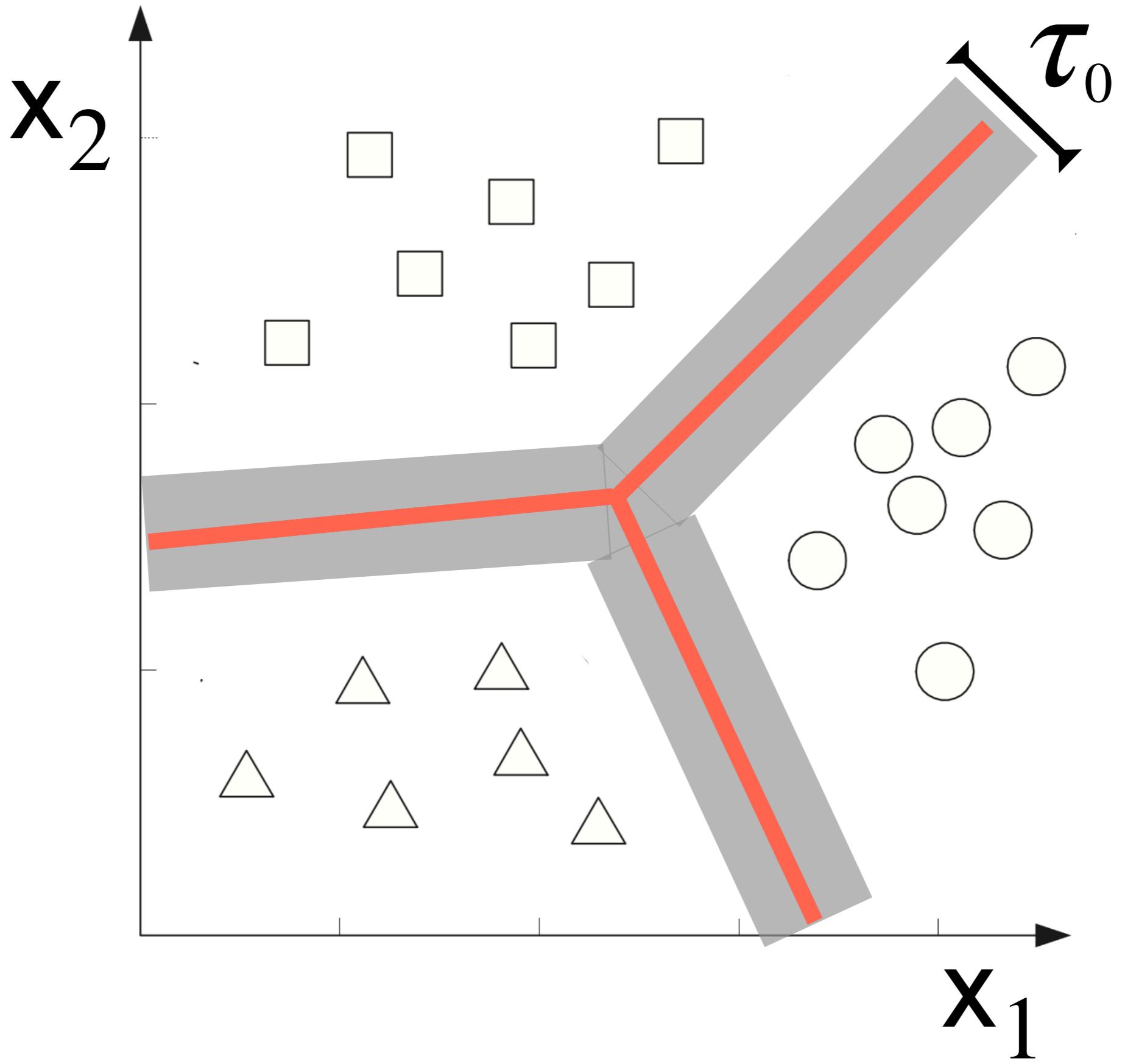
expert

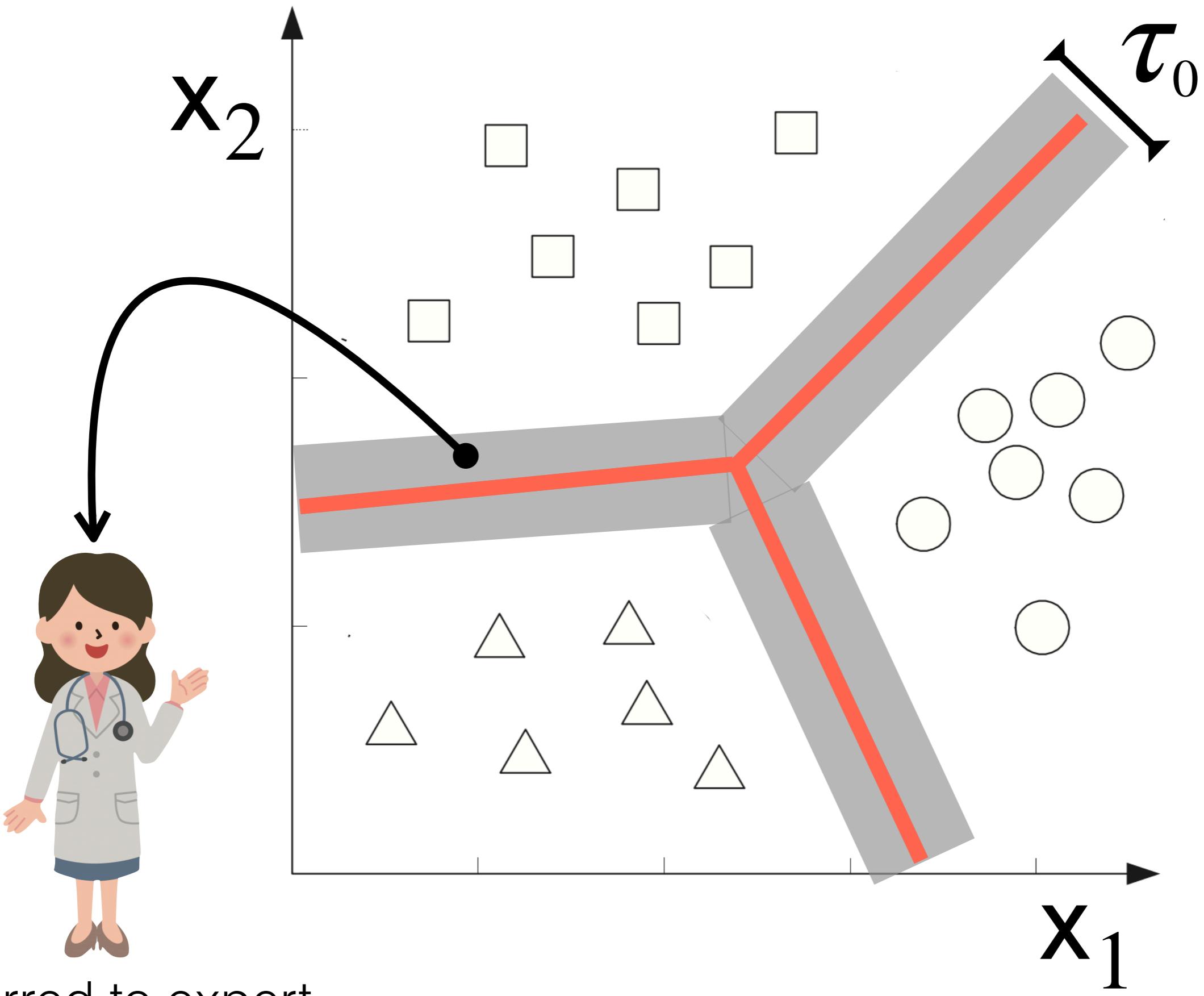
defer to expert if...

$$\max_y p(y|x) \leq \tau_0$$

(constant)







deferred to expert

input
features



allocation
mechanism



defer to expert if...

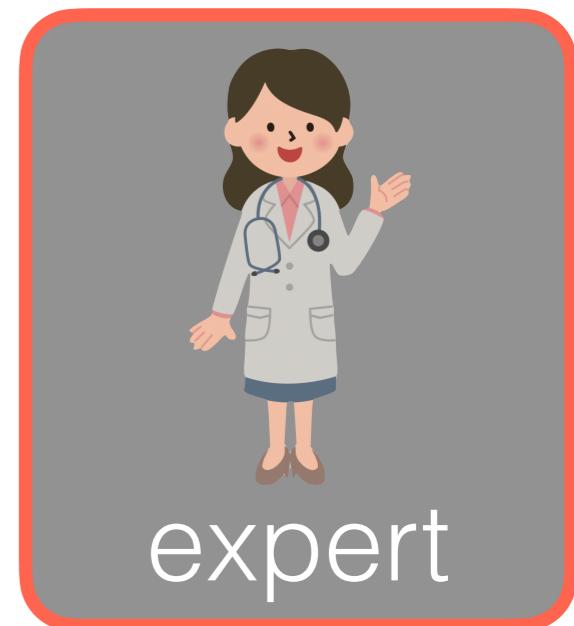
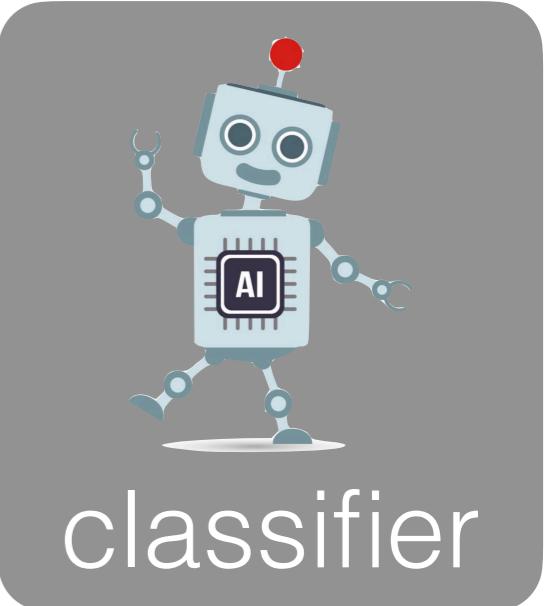
$$\max_y p(y|x) \leq \tau_0 \quad (\text{constant})$$

problem?

input
features



allocation
mechanism



defer to expert if...

$$\max_y p(y|x) \leq \tau_0 \quad (\text{constant})$$

the expert's
knowledge is
not considered!

input
features



allocation
mechanism



defer to expert if...

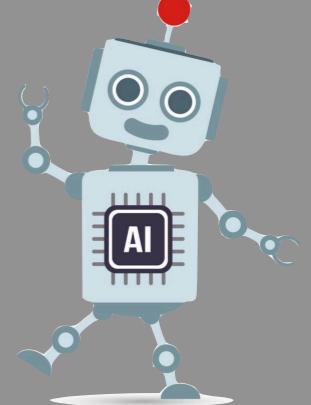
$$\max_y p(y|x) \leq \tau \left(\text{[skeleton icon]}, \text{[doctor icon]} \right)$$

input
features



allocation
mechanism

L_{0-1}



classifier



expert

defer to expert if...

$$\max_y p(y|x) \leq \tau \left(\text{[skeleton icon]}, \text{[doctor icon]} \right),$$

input
features



allocation
mechanism

L_{0-1}



defer to expert if...

$$\max_y p(y|x) \leq \tau \left(\text{[skeleton icon]}, \text{[doctor icon]} \right)$$

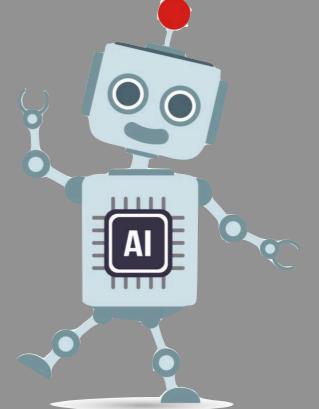
input
features



allocation
mechanism

L_{0-1}

L_{0-1}



classifier



expert

Bayes optimal deferral rule:

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

y

probability that the expert is correct

softmax implementation

[Mozannar & Sontag, 2020]

training data

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

softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

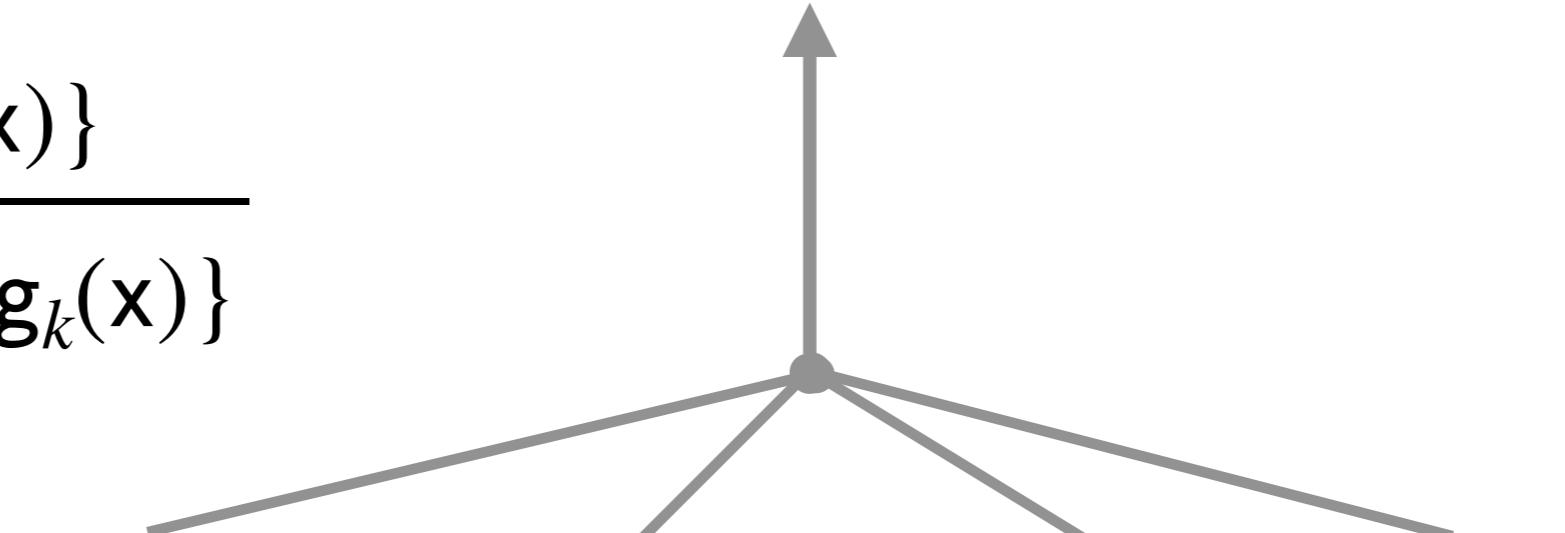
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$$\begin{array}{c} h_1(\mathbf{x}) \quad \cdots \quad h_k(\mathbf{x}) \quad \cdots \quad h_K(\mathbf{x}) \quad h_\perp(\mathbf{x}) \\ \hline \end{array}$$

\mathbf{x}

\mathbf{g}_θ

$$\begin{array}{c} g_1(\mathbf{x}) \quad \cdots \quad g_k(\mathbf{x}) \quad \cdots \quad g_K(\mathbf{x}) \quad g_\perp(\mathbf{x}) \\ \hline \end{array}$$



softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

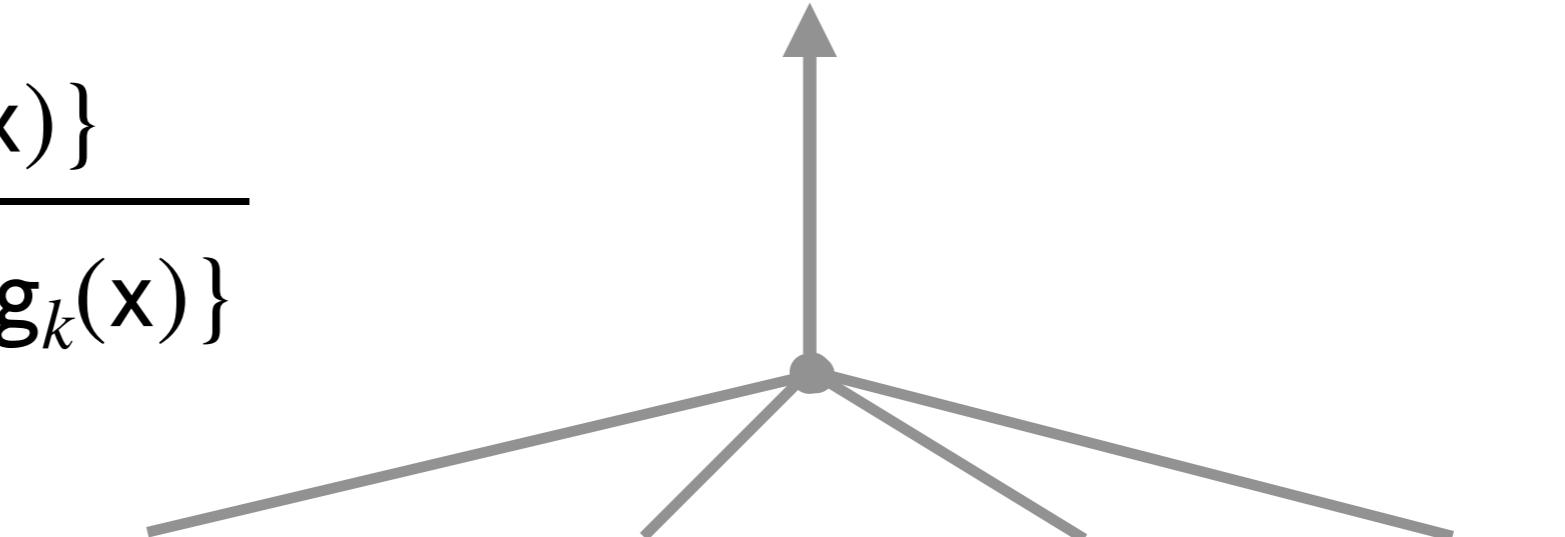
$$h_i(\mathbf{x}) = \frac{\exp\{g_i(\mathbf{x})\}}{\sum_{k=1}^{K+1} \exp\{g_k(\mathbf{x})\}}$$

$$\begin{array}{c} h_1(\mathbf{x}) \quad \cdots \quad h_k(\mathbf{x}) \quad \cdots \quad h_K(\mathbf{x}) \quad h_{\perp}(\mathbf{x}) \\ \hline \end{array}$$

\mathbf{x}

\mathbf{g}_{θ}

$$\begin{array}{c} g_1(\mathbf{x}) \quad \cdots \quad g_k(\mathbf{x}) \quad \cdots \quad g_K(\mathbf{x}) \quad g_{\perp}(\mathbf{x}) \\ \hline \end{array}$$



softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

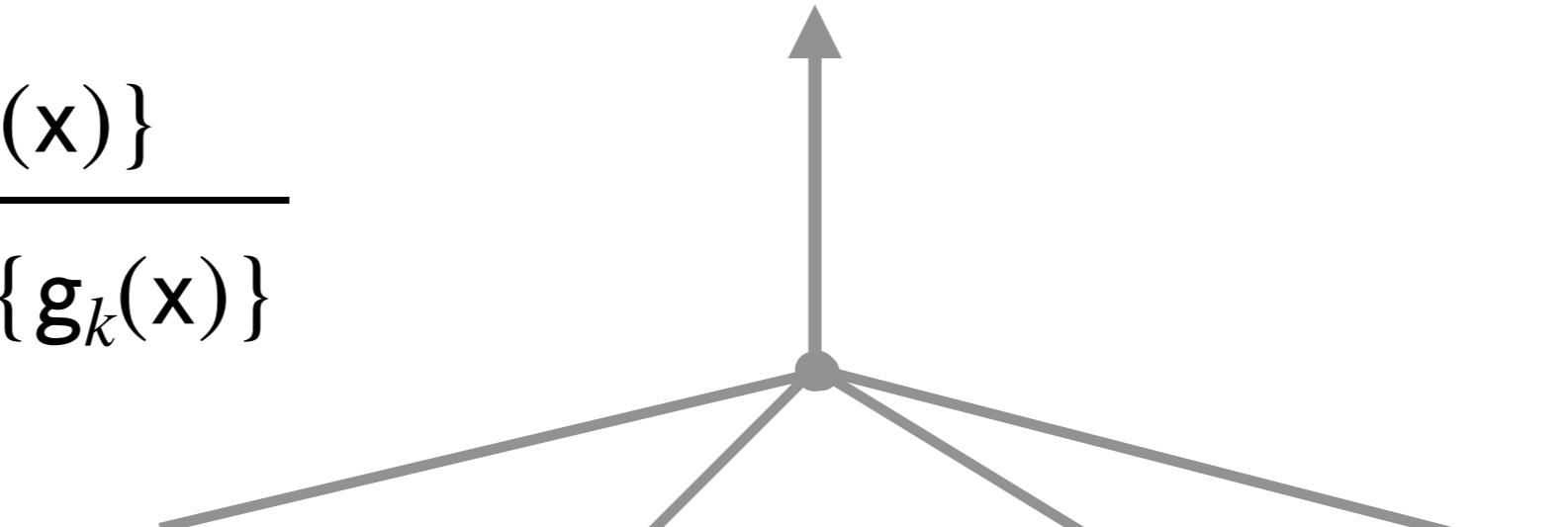
$$h_{\perp}(x) = \frac{\exp\{g_{\perp}(x)\}}{\sum_{k=1}^{K+1} \exp\{g_k(x)\}}$$

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

x

$$g_{\theta}$$

$g_1(x)$	\dots	$g_k(x)$	\dots	$g_K(x)$	$g_{\perp}(x)$
----------	---------	----------	---------	----------	----------------



softmax implementation

[Mozannar & Sontag, 2020]

training data

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

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

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

model

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

loss function

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

softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

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

loss function

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

softmax implementation

[Mozannar & Sontag, 2020]

training data

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

model

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

loss function

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

input
features



allocation
mechanism



defer to expert if...

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

- ⊗ 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

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

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

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

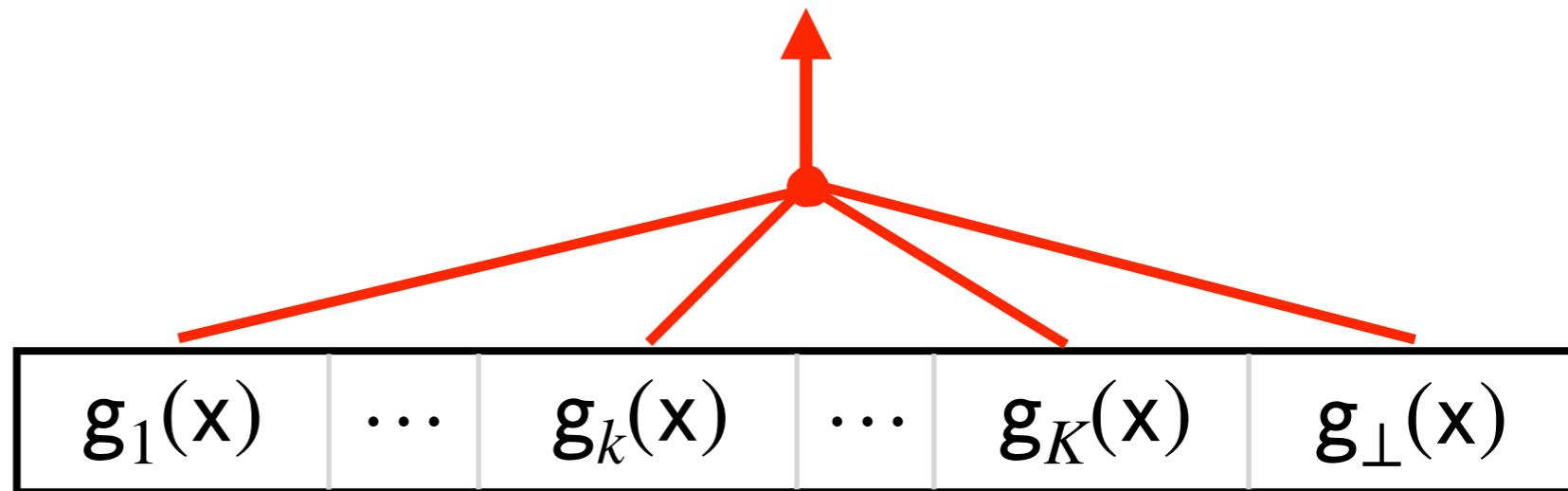
How well can the softmax-based system estimate expert correctness?

$$\hat{p}(m = y | x) \cancel{\approx} P(m = y | x)$$

degenerate
parameterization

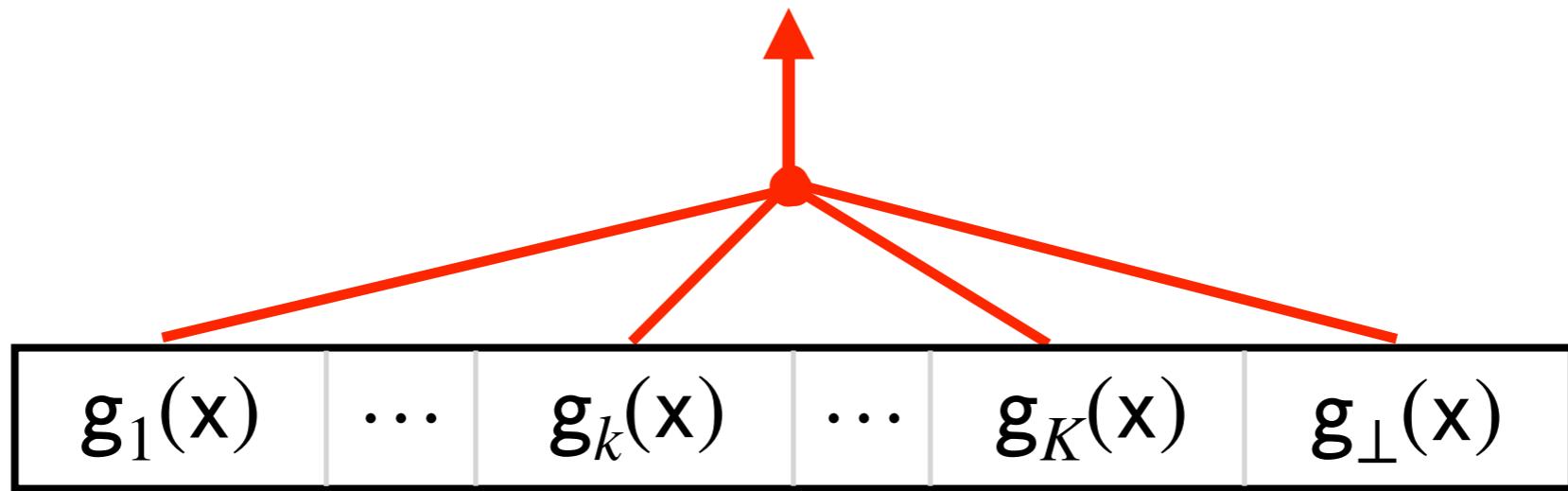
[Proposition 3.1]

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

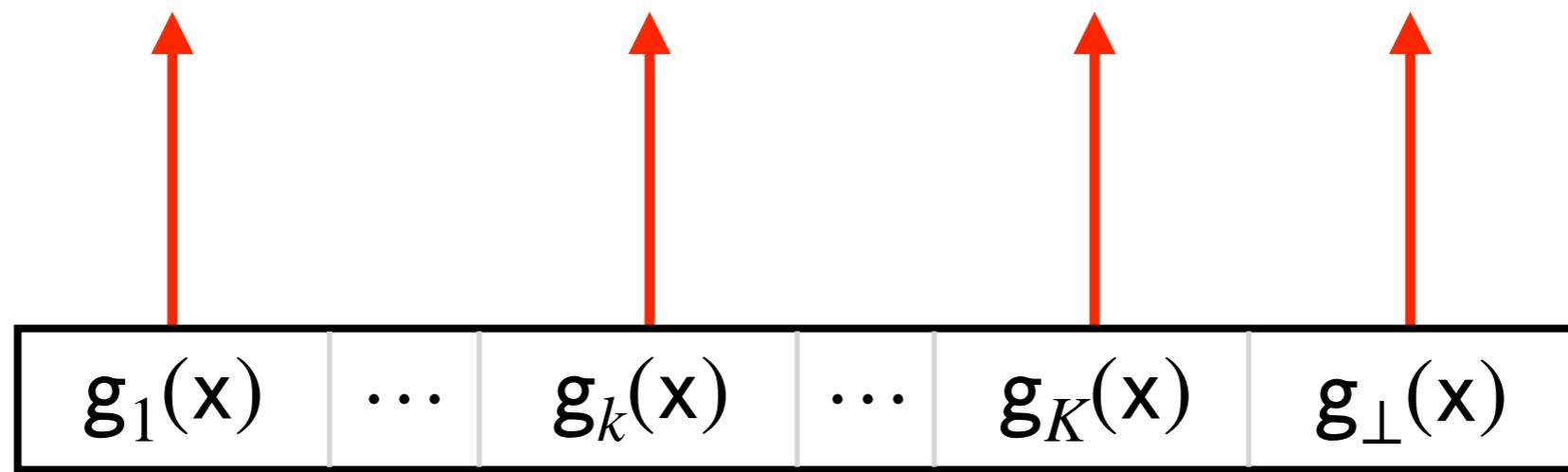


our solution: switch to a one-vs-all parameterization

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

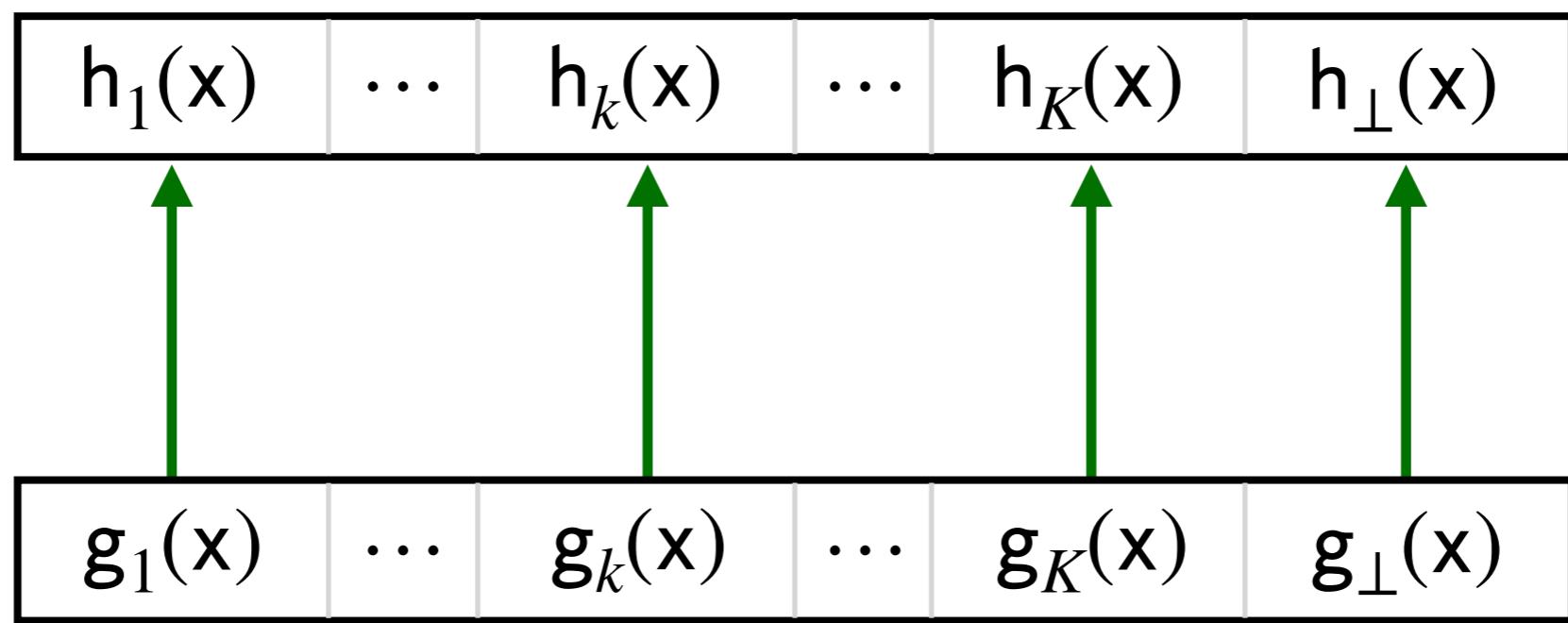


our solution: switch to a
one-vs-all parameterization



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

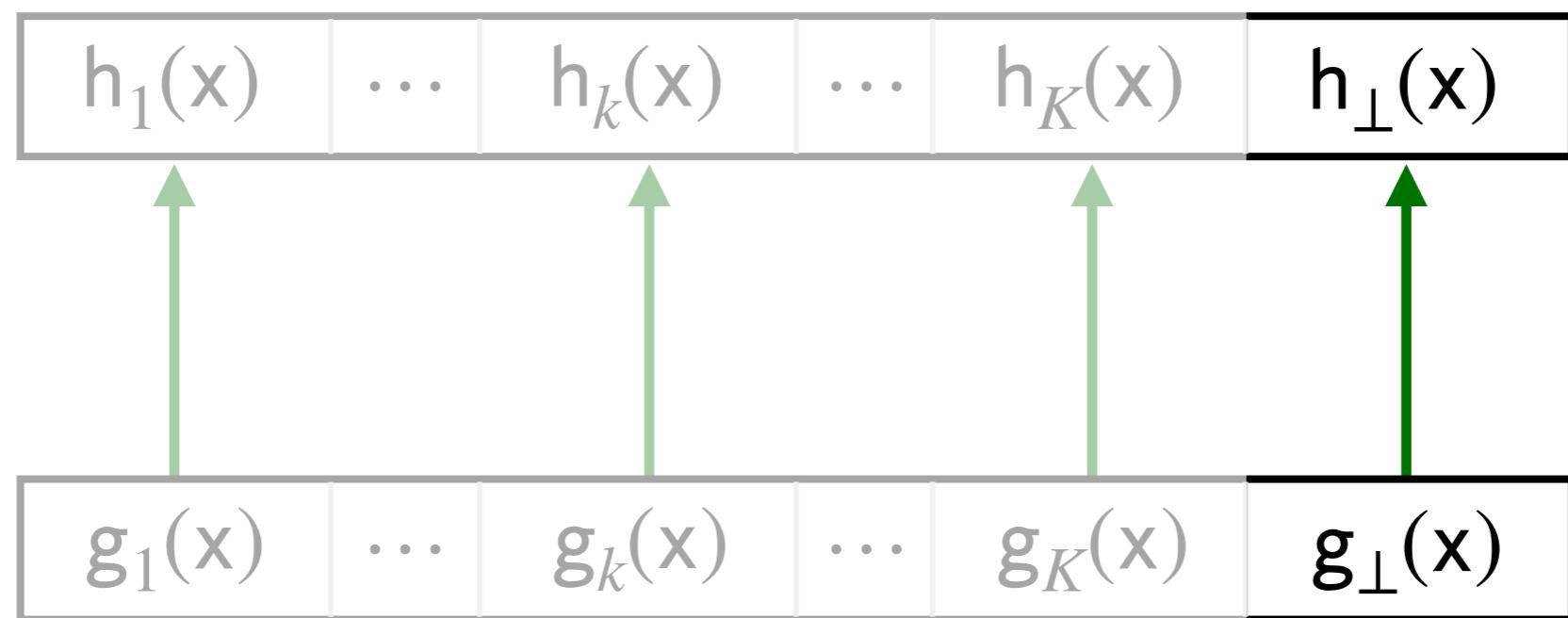
$$h_i(x) = \frac{1}{1 + \exp \{-g_i(x)\}}$$



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

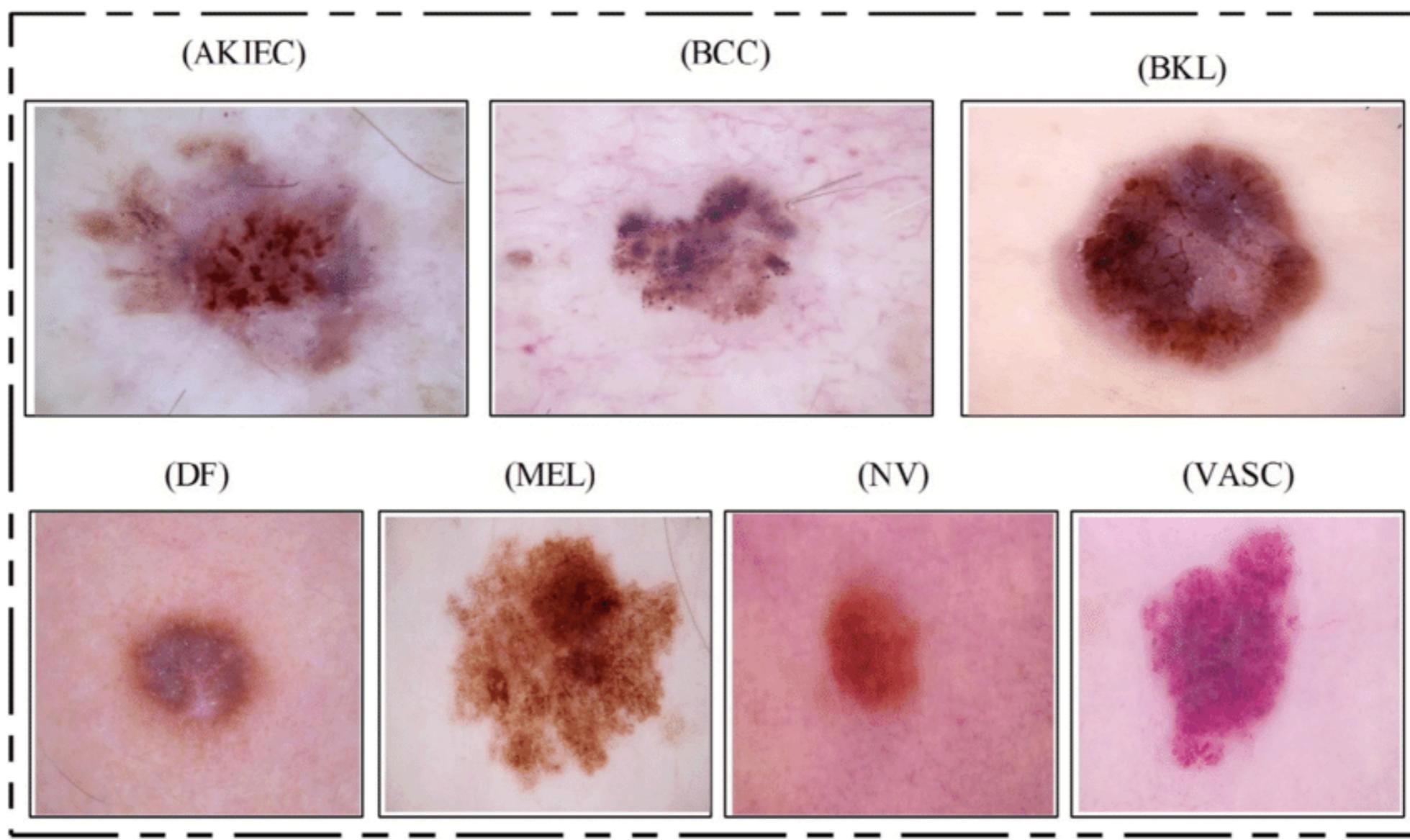
[Theorem 4.1]

$$h_{\perp}^*(x) = P(m = y | x)$$



estimating expert correctness

skin lesion diagnosis



estimating expert correctness

\hat{p}

distance: \hat{p} vs P

softmax

26.7 ± 1.8

one-vs-all
(ours)

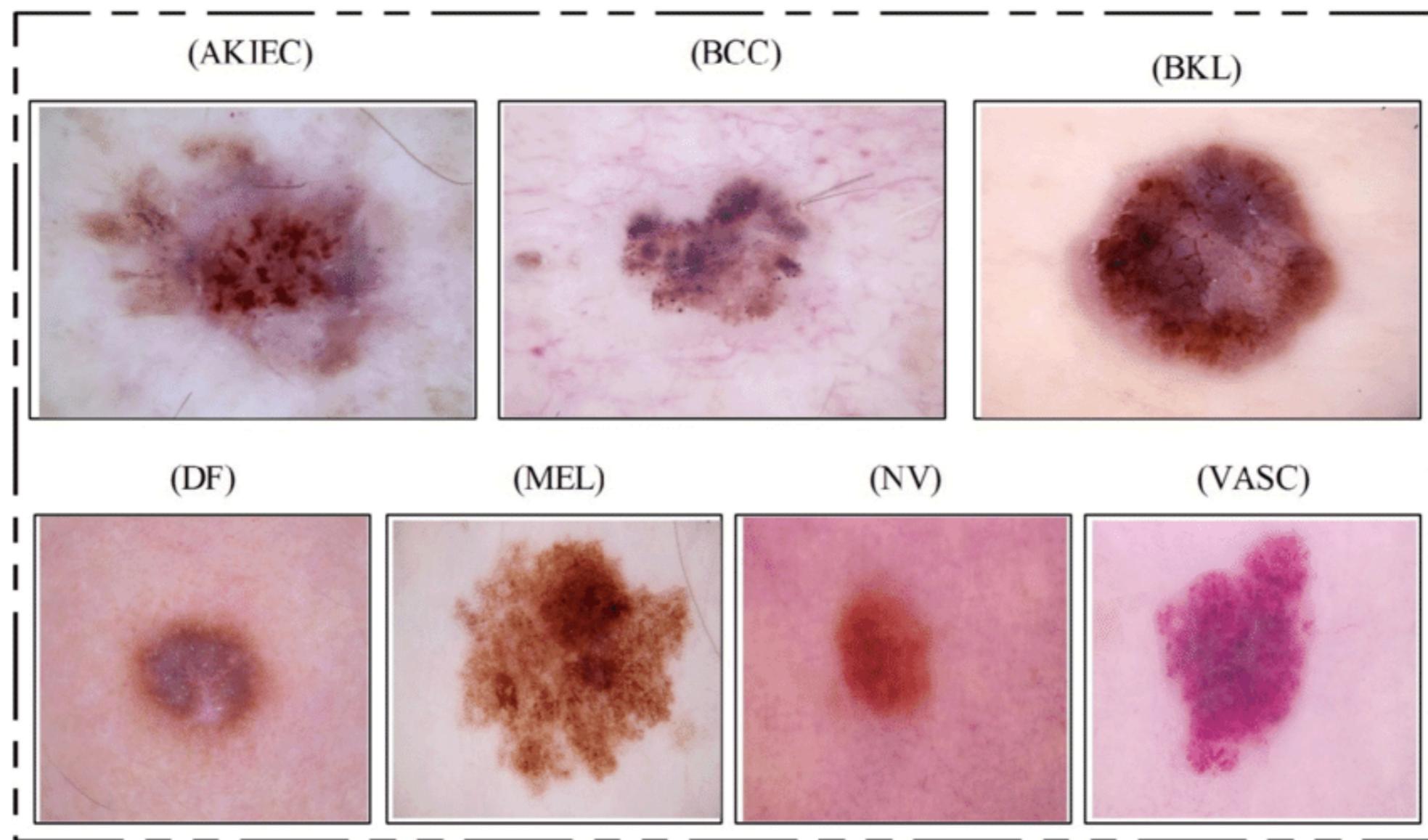
8.0 ± 1.0

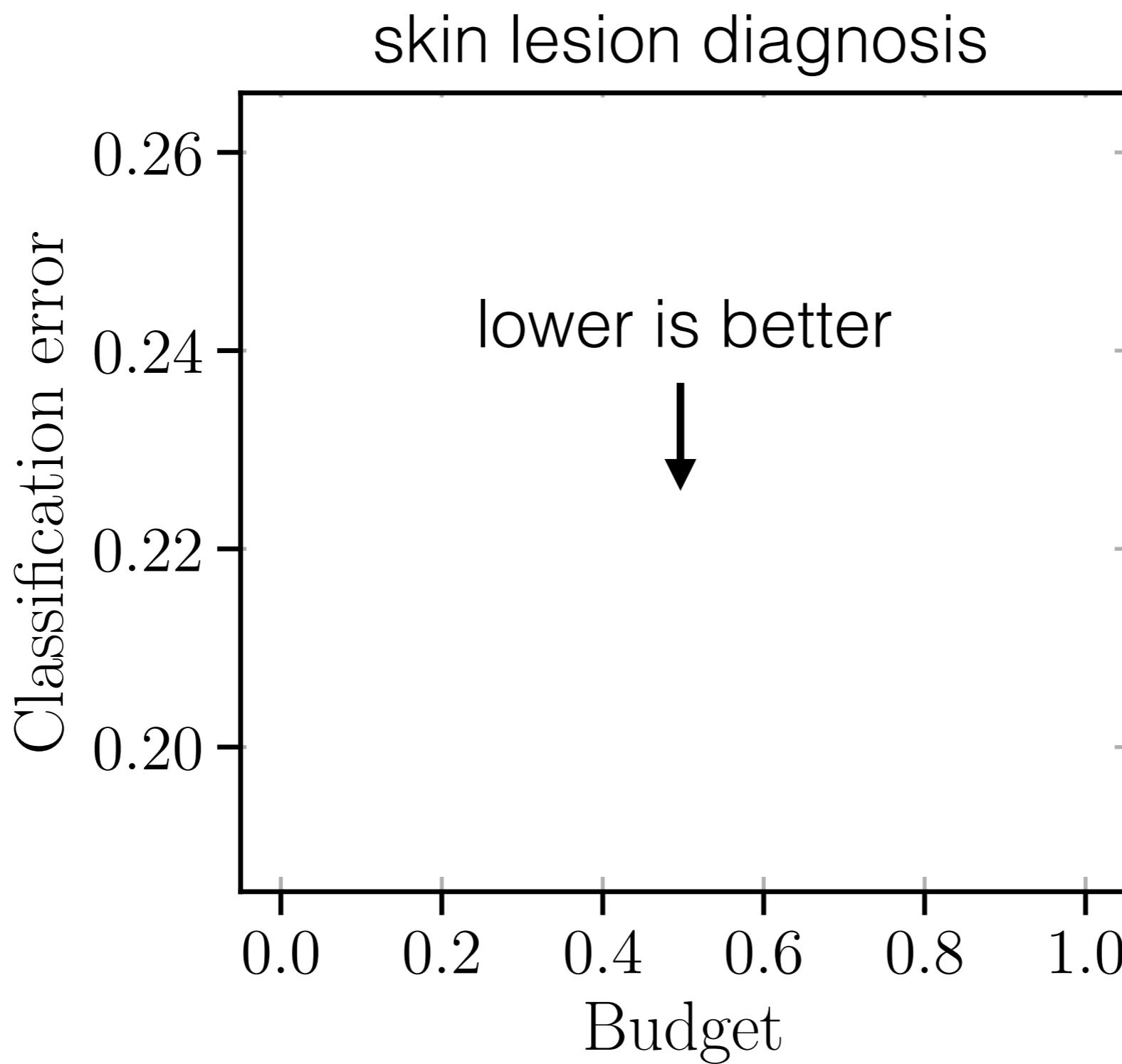
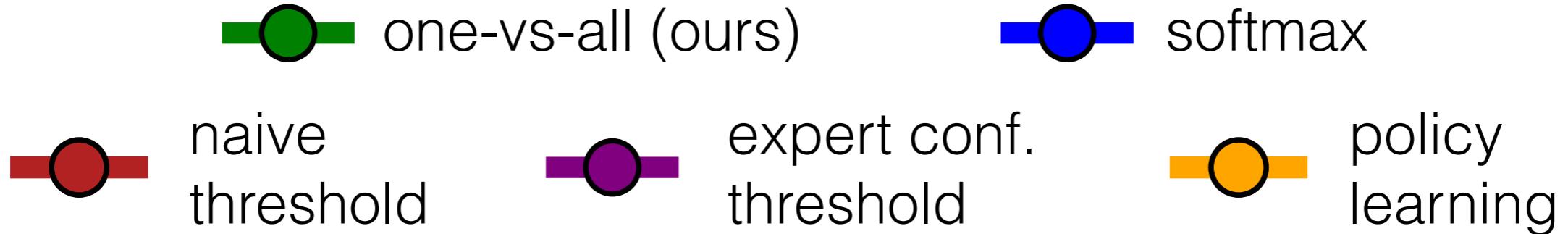


But does one-vs-all
result in more accurate models?



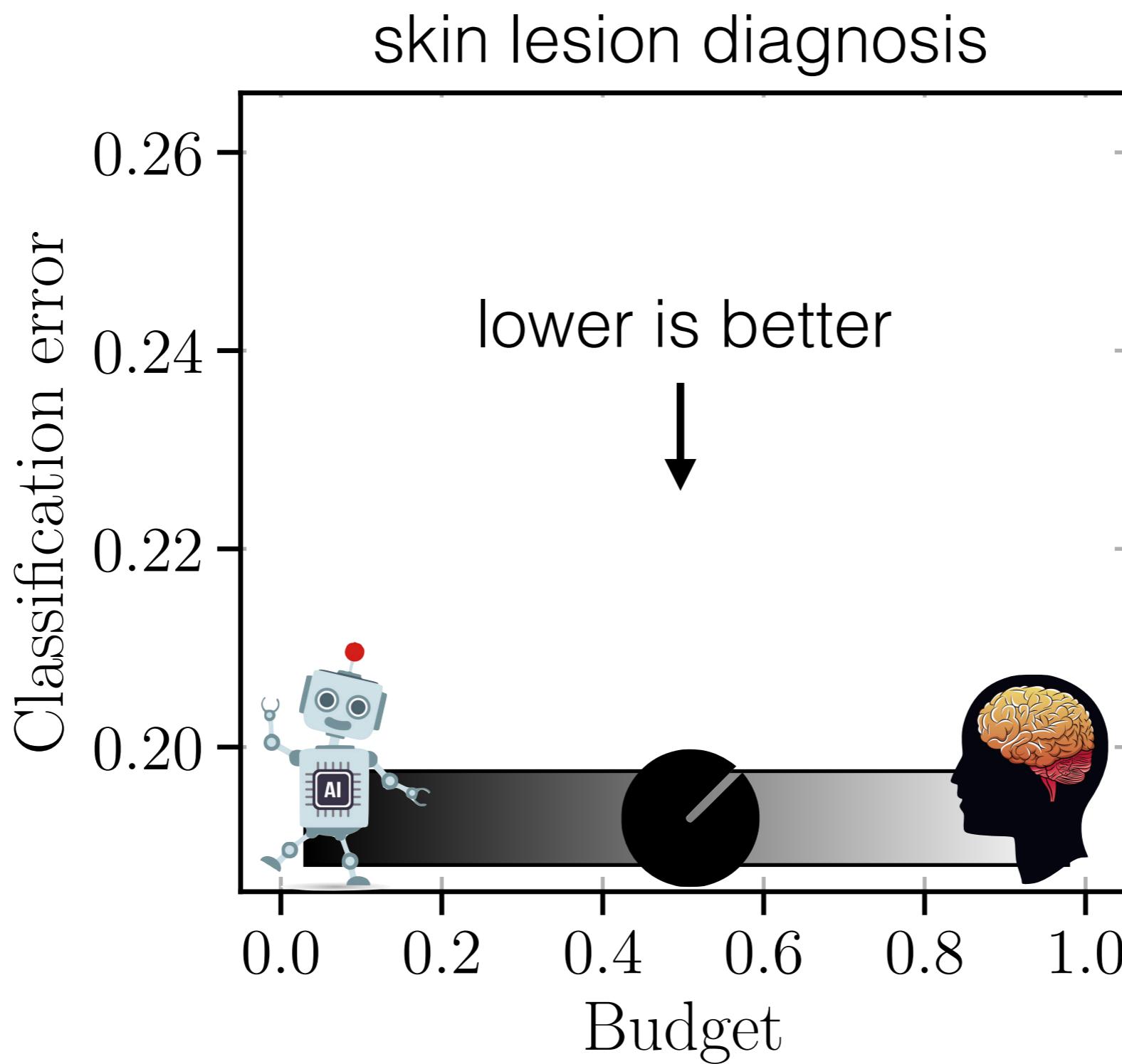
skin lesion diagnosis

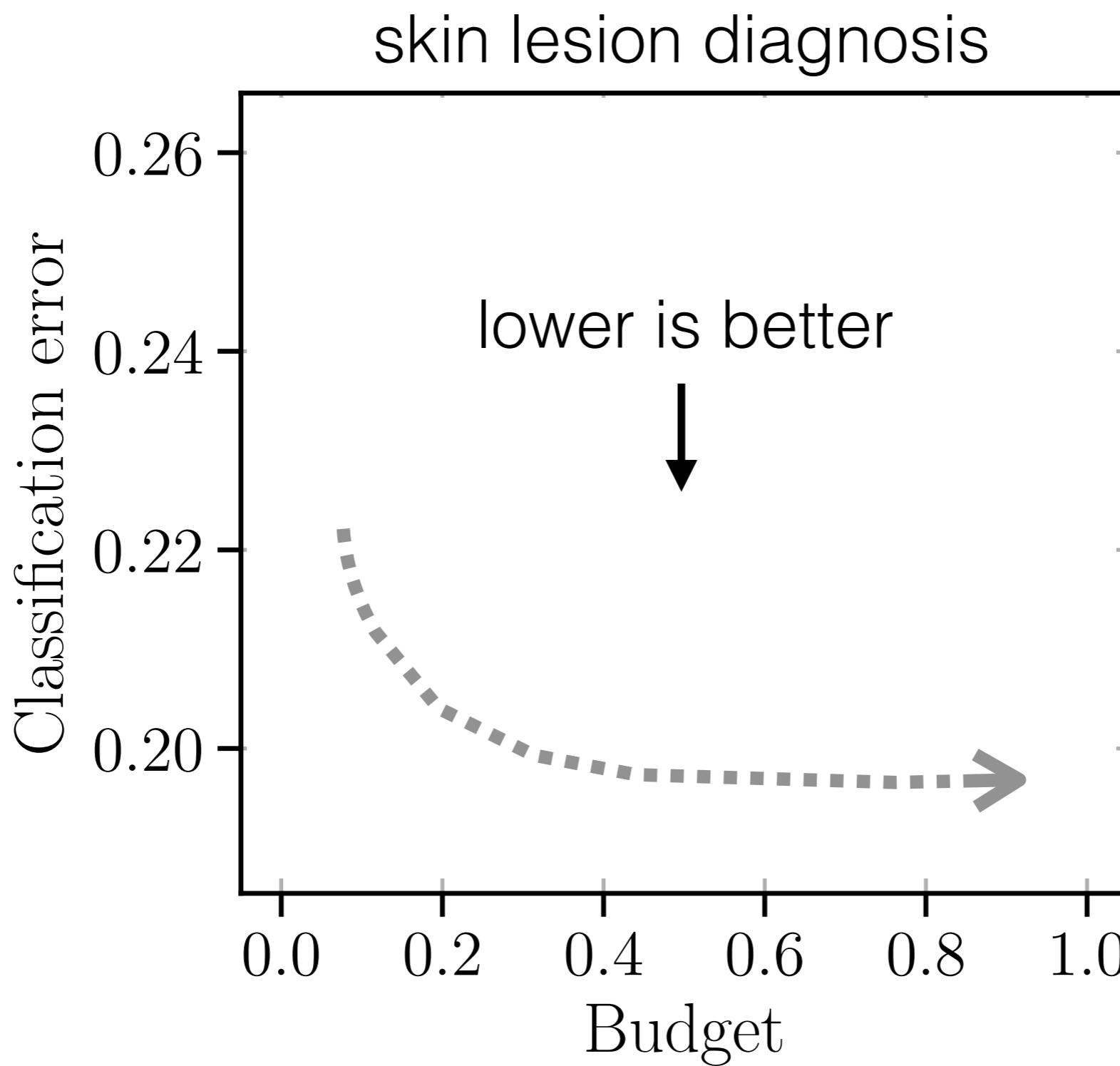




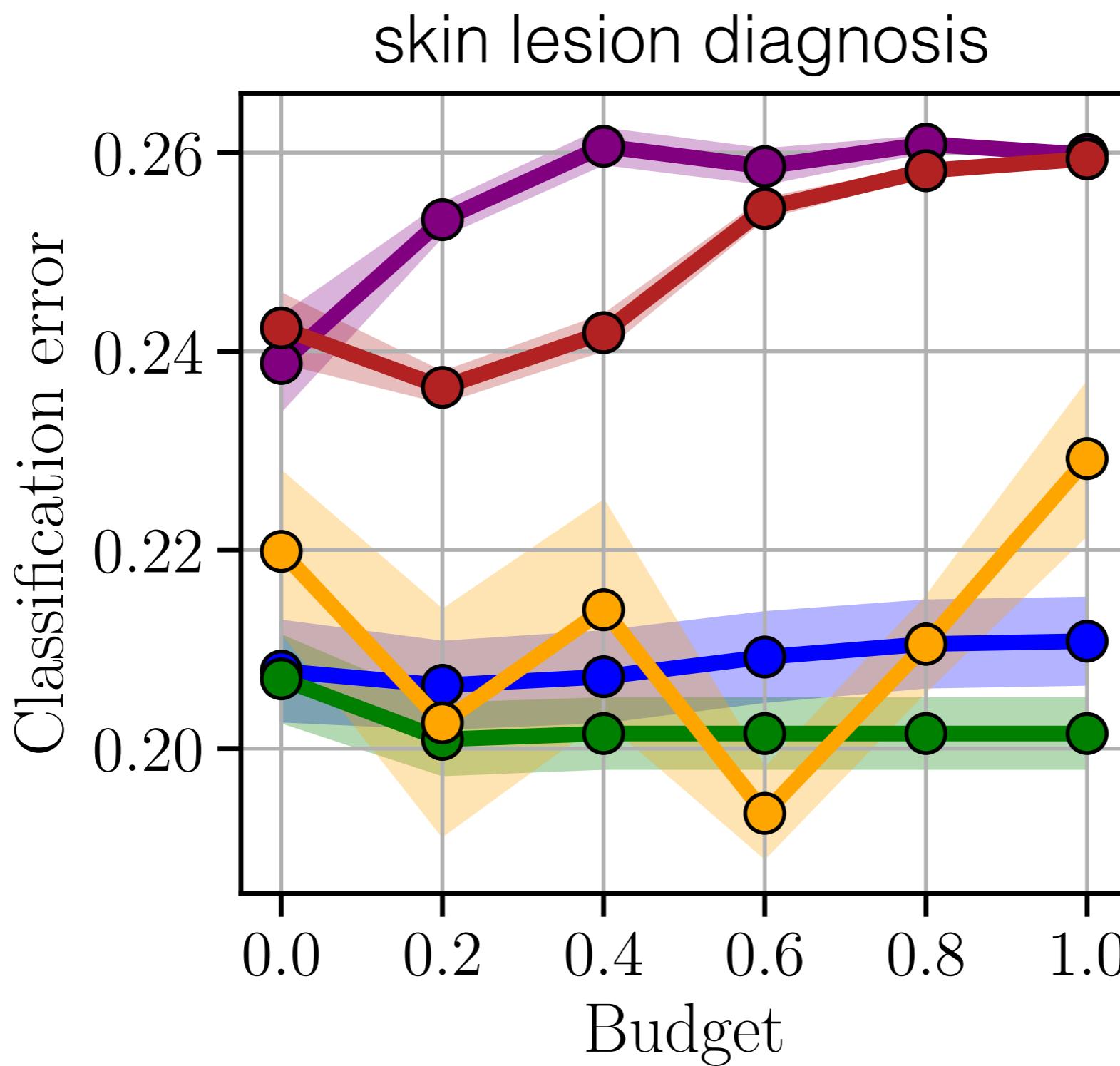
Legend:

- one-vs-all (ours) (Green)
- softmax (Blue)
- naive threshold (Red)
- expert conf. threshold (Purple)
- policy learning (Yellow)



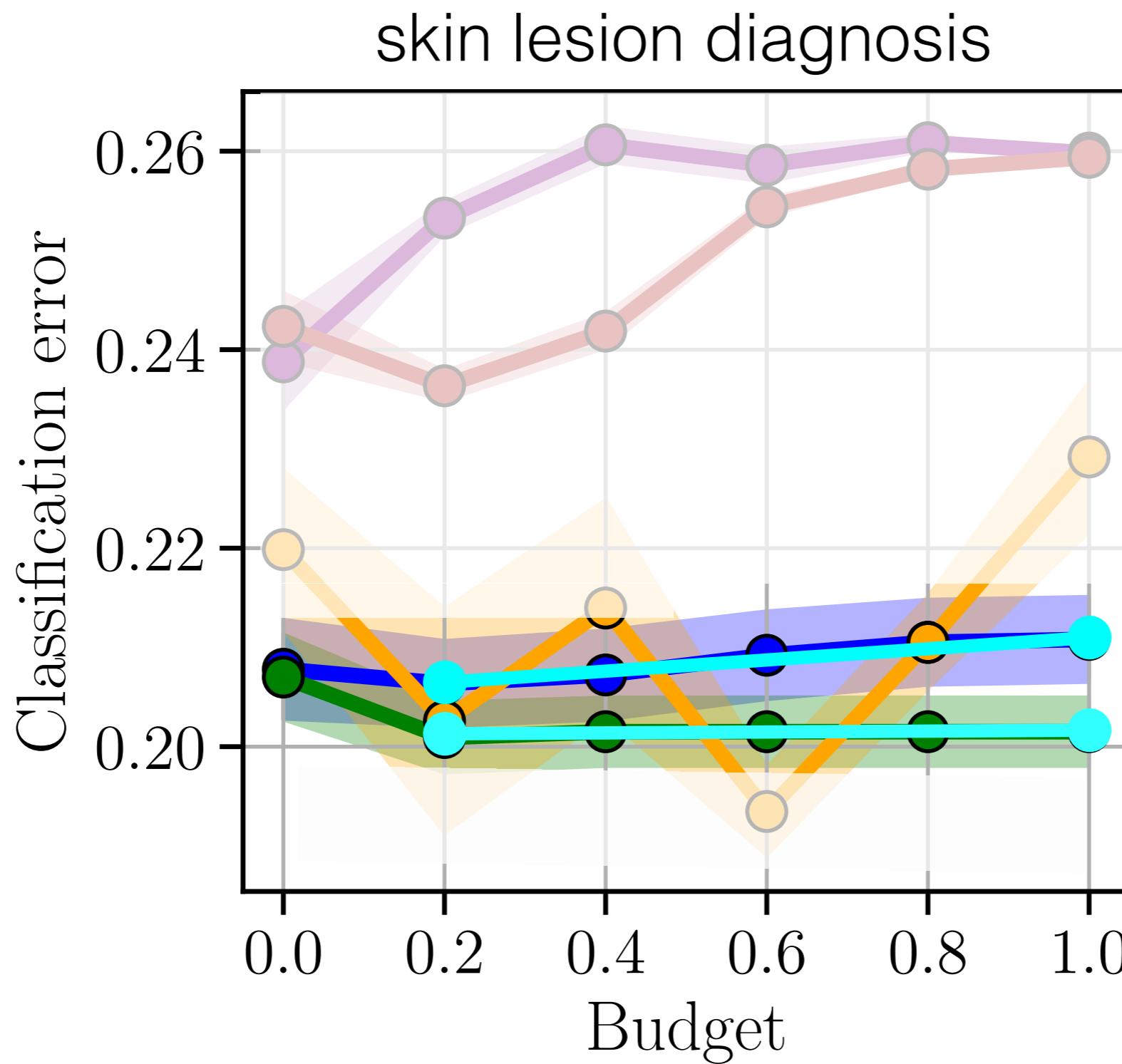


● one-vs-all (ours) ● softmax
● naive threshold ● expert conf. threshold ● policy learning



one-vs-all (ours) softmax

naive threshold expert conf. threshold policy learning

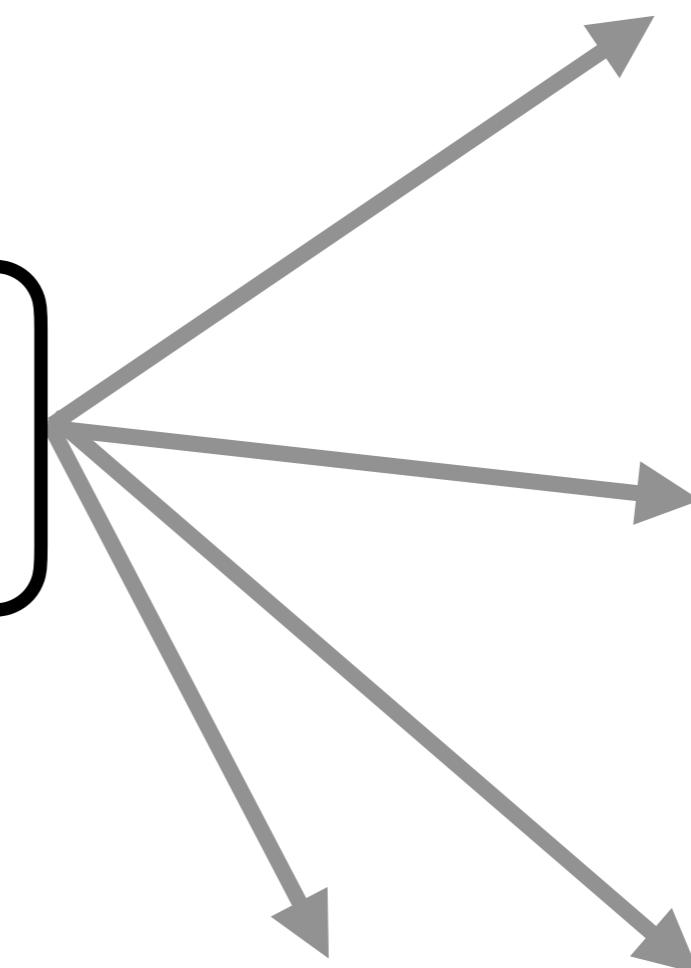


- ⊗ single expert
 - ⊗ softmax surrogate loss
 - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
 - ⊗ surrogate losses
 - ⊗ conformal sets of experts
- ⊗ population of experts
 - ⊗ surrogate losses
 - ⊗ meta-learning a rejector

input
features



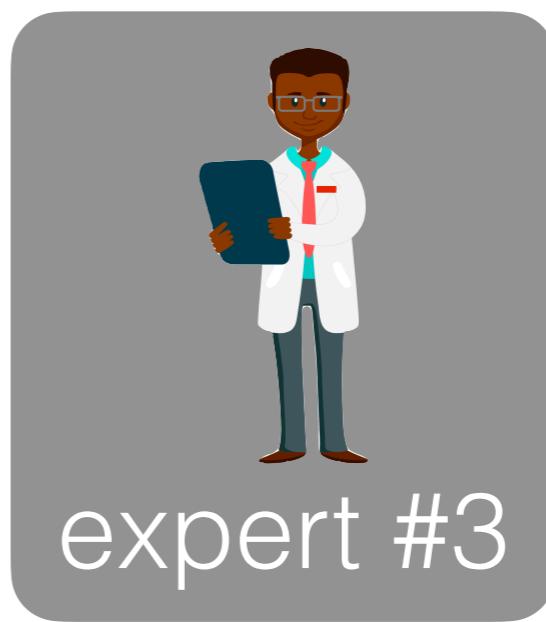
allocation
mechanism



classifier



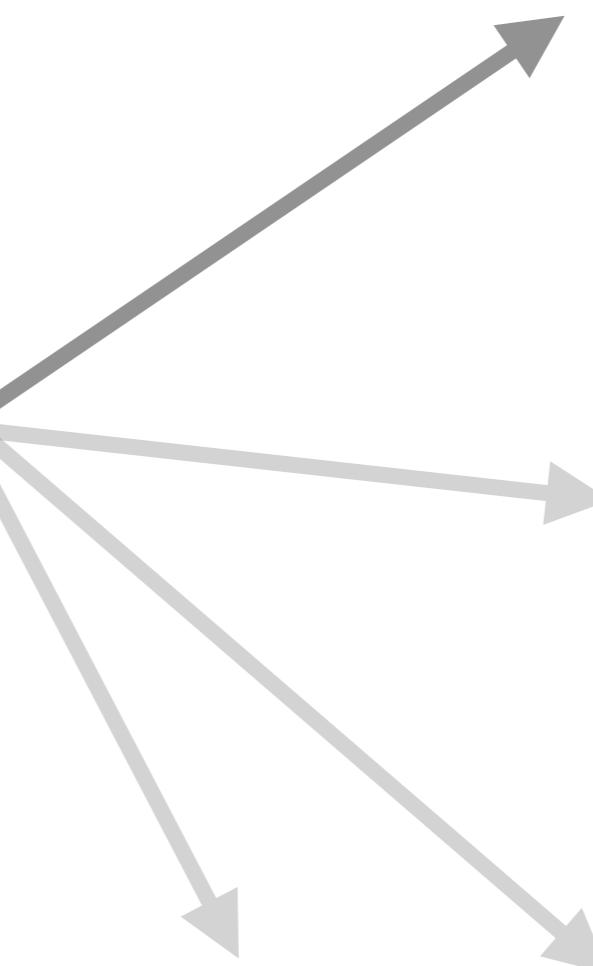
expert #1



input
features



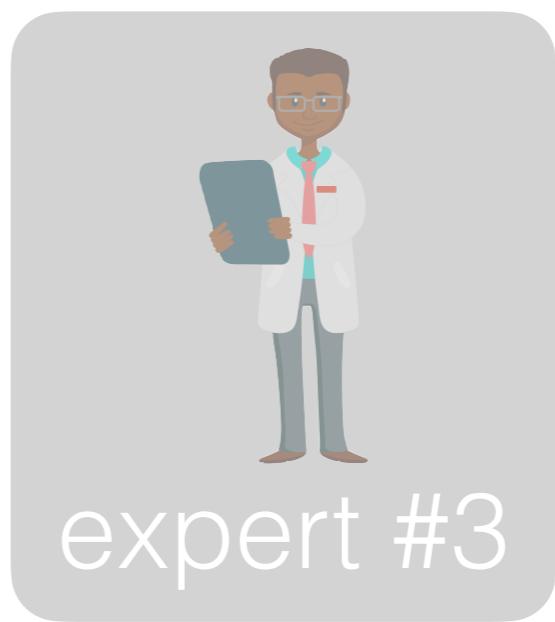
allocation
mechanism



classifier



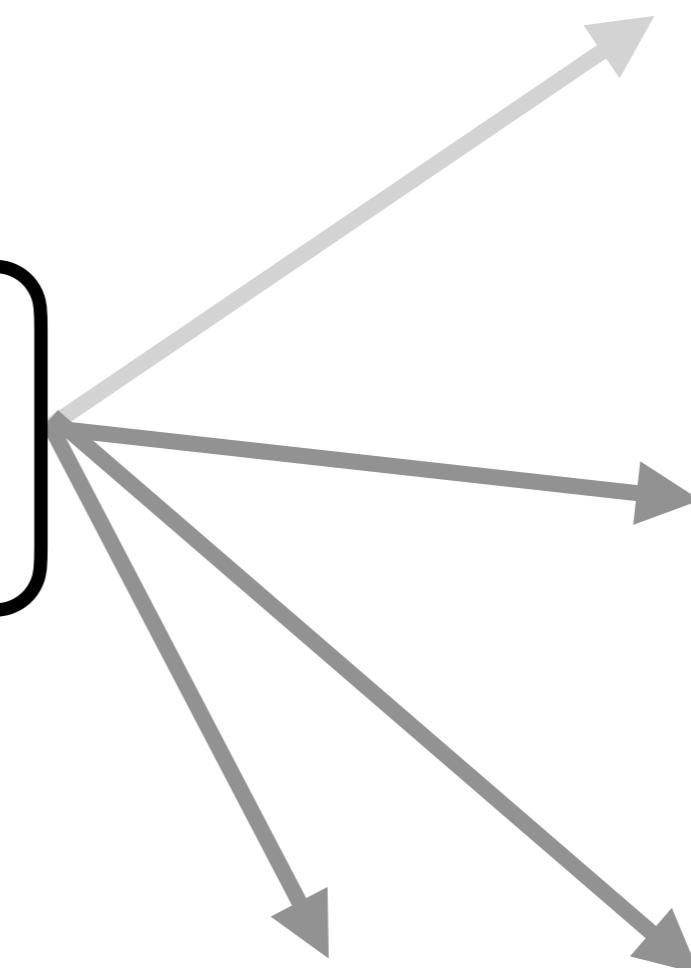
expert #1



input
features



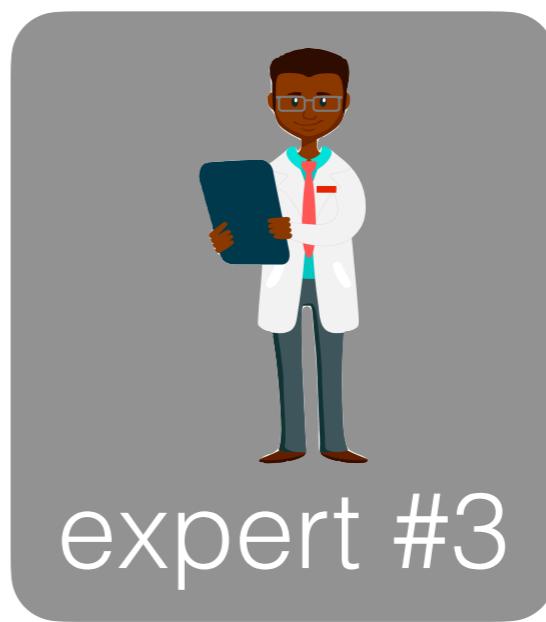
allocation
mechanism



classifier



expert #1

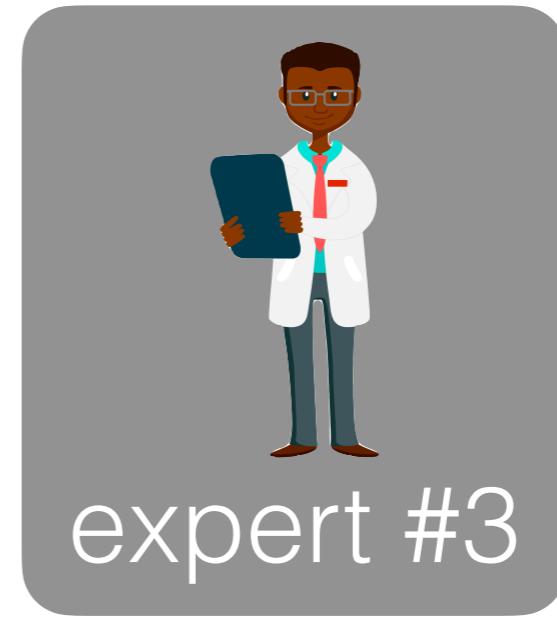


input
features



allocation
mechanism

expert #3



expert #2

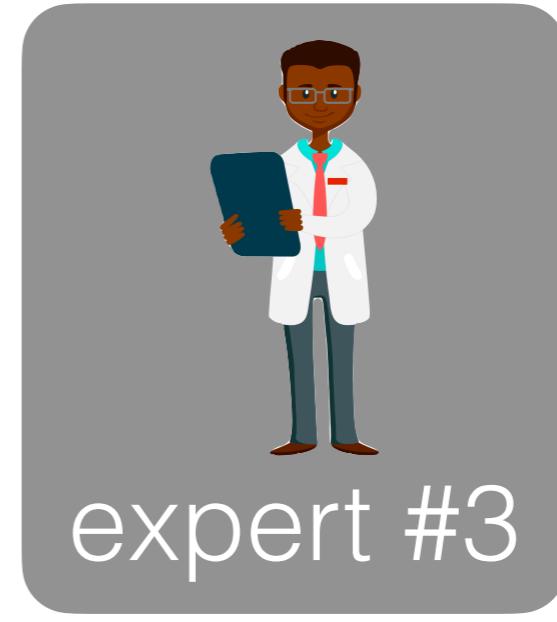


input
features

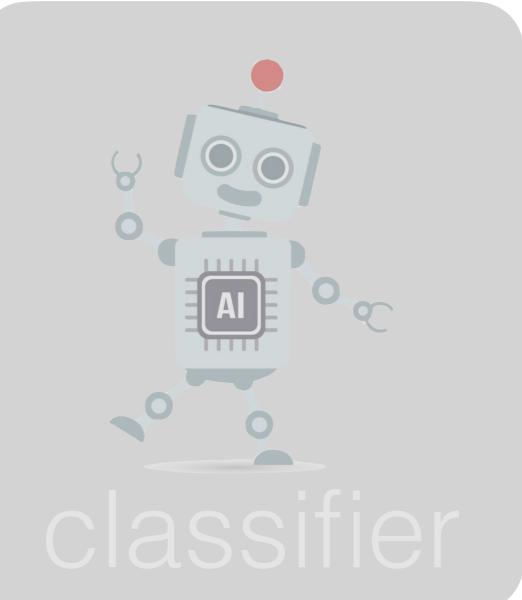


allocation
mechanism

expert #3



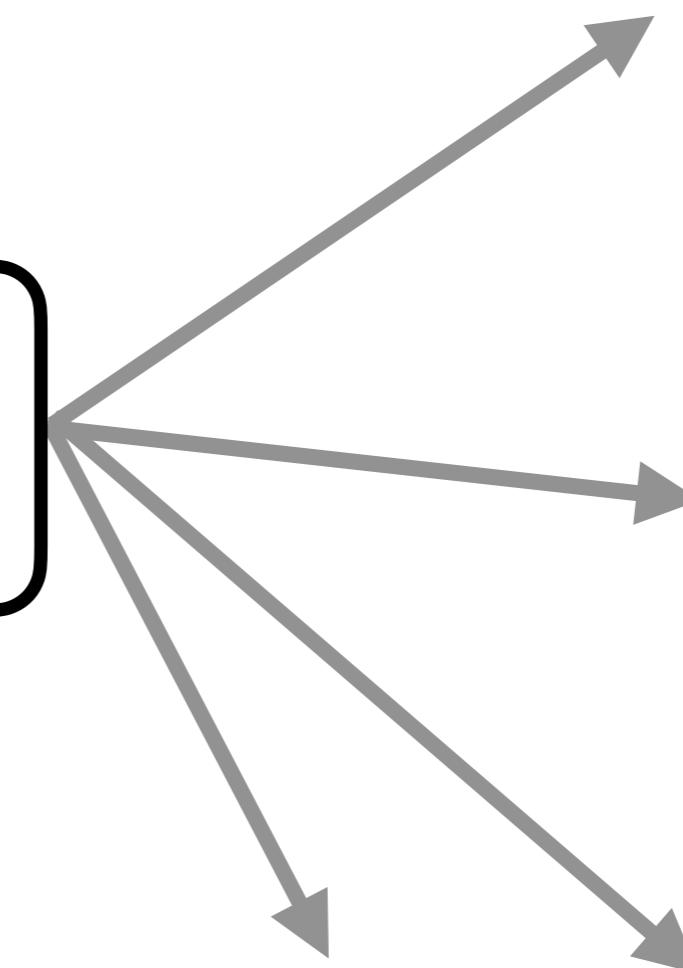
expert #2



input
features



allocation
mechanism

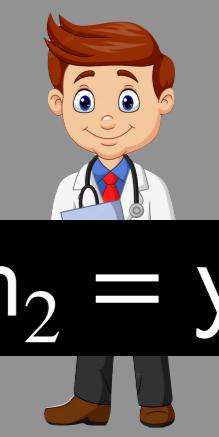


classifier



$$P(m_1 = y | x)$$

expert #1



$$P(m_3 = y | x)$$

expert #3

$$P(m_2 = y | x)$$

expert #2



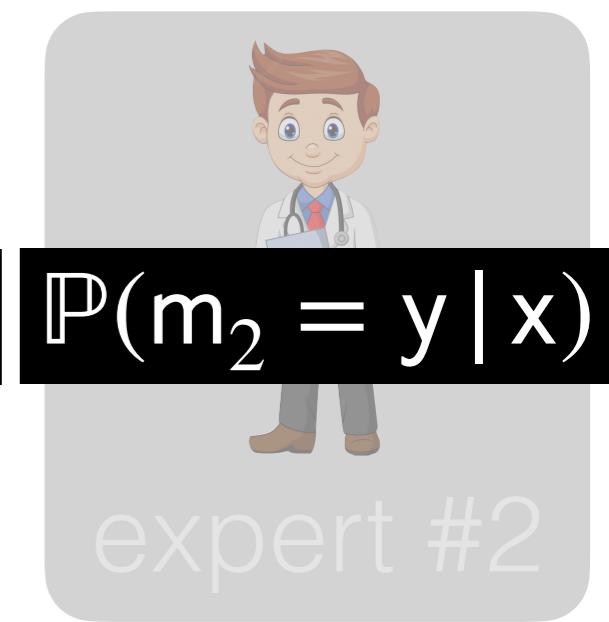
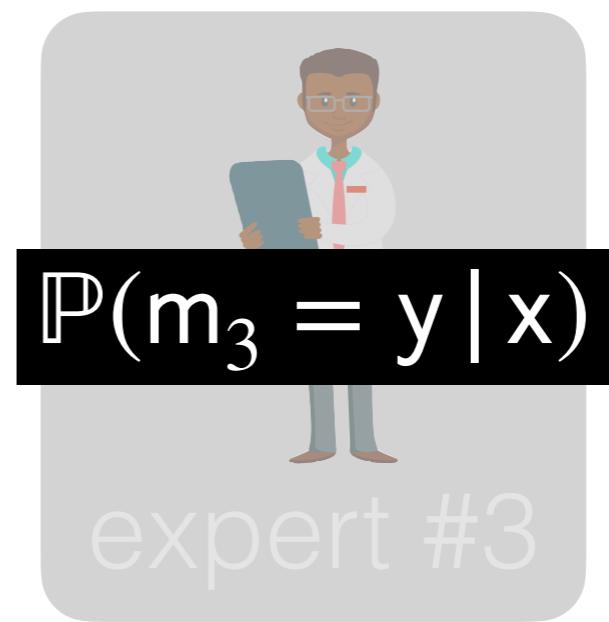
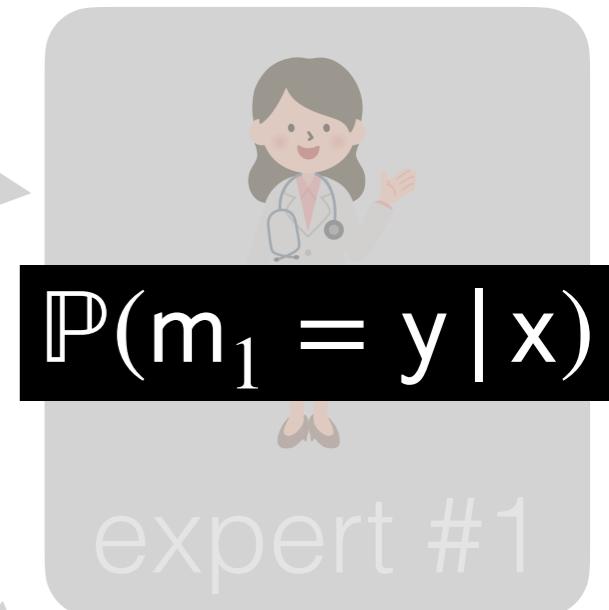
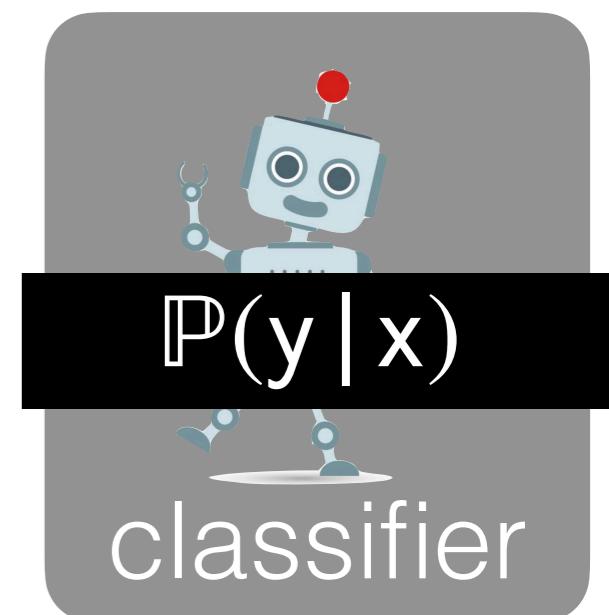
input
features



allocation
mechanism

use classifier if...

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



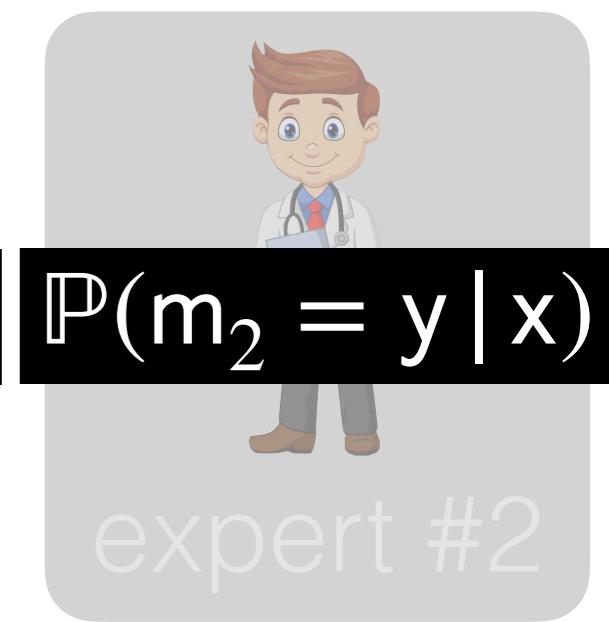
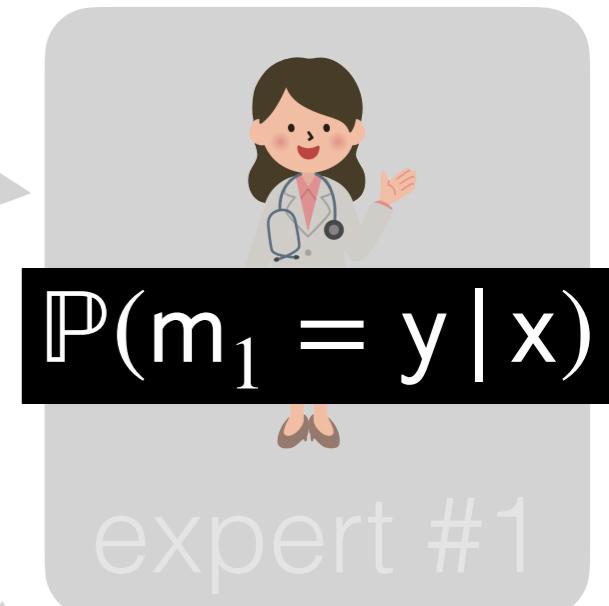
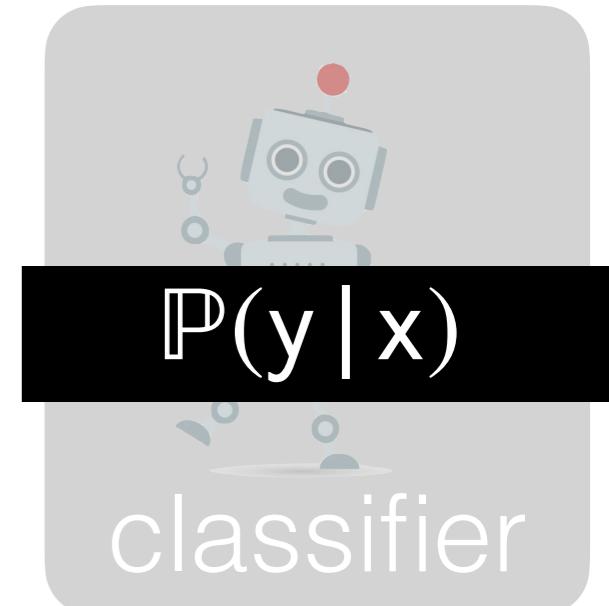
input
features



allocation
mechanism

else, pick best expert:

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



multi-expert implementation

training data

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

multi-expert implementation

training data

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

model

$h_1(\mathbf{x})$	\dots	$h_k(\mathbf{x})$	\dots	$h_K(\mathbf{x})$	$h_{\perp,1}(\mathbf{x})$	\dots	$h_{\perp,J}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	---------------------------	---------	---------------------------

K classes

J experts

$g_1(\mathbf{x})$	\dots	$g_k(\mathbf{x})$	\dots	$g_K(\mathbf{x})$	$g_{\perp,1}(\mathbf{x})$	\dots	$g_{\perp,J}(\mathbf{x})$
-------------------	---------	-------------------	---------	-------------------	---------------------------	---------	---------------------------

multi-expert implementation

training data

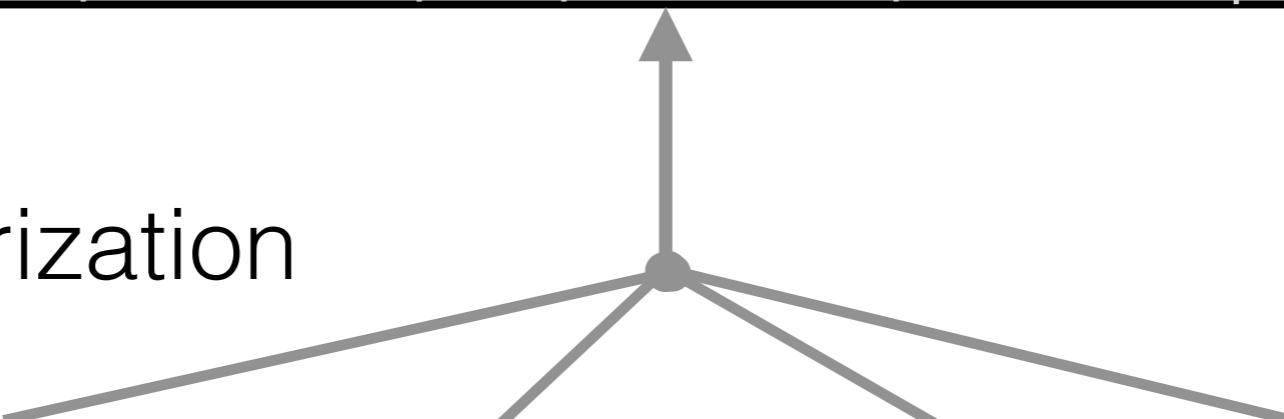
$$\mathcal{D} = \left\{ \mathbf{x}_n, \mathbf{y}_n, \mathbf{m}_{n,1}, \dots, \mathbf{m}_{n,J} \right\}_{n=1}^N$$

model

$$\boxed{h_1(\mathbf{x}) \mid \dots \mid h_k(\mathbf{x}) \mid \dots \mid h_K(\mathbf{x}) \mid h_{\perp,1}(\mathbf{x}) \mid \dots \mid h_{\perp,J}(\mathbf{x})}$$

softmax
parametrization

$$\boxed{g_1(\mathbf{x}) \mid \dots \mid g_k(\mathbf{x}) \mid \dots \mid g_K(\mathbf{x}) \mid g_{\perp,1}(\mathbf{x}) \mid \dots \mid g_{\perp,J}(\mathbf{x})}$$



multi-expert implementation

training data

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

model

$$\boxed{h_1(\mathbf{x}) \mid \dots \mid h_k(\mathbf{x}) \mid \dots \mid h_K(\mathbf{x}) \mid h_{\perp,1}(\mathbf{x}) \mid \dots \mid h_{\perp,J}(\mathbf{x})}$$

one-vs-all
parameterization

$$\boxed{g_1(\mathbf{x}) \mid \dots \mid g_k(\mathbf{x}) \mid \dots \mid g_K(\mathbf{x}) \mid g_{\perp,1}(\mathbf{x}) \mid \dots \mid g_{\perp,J}(\mathbf{x})}$$

multi-expert implementation

training data

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

model

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

multi-expert implementation

training data

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

model

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

softmax loss function

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

multi-expert implementation

training data

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

model

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

softmax loss function

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

multi-expert implementation

training data

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

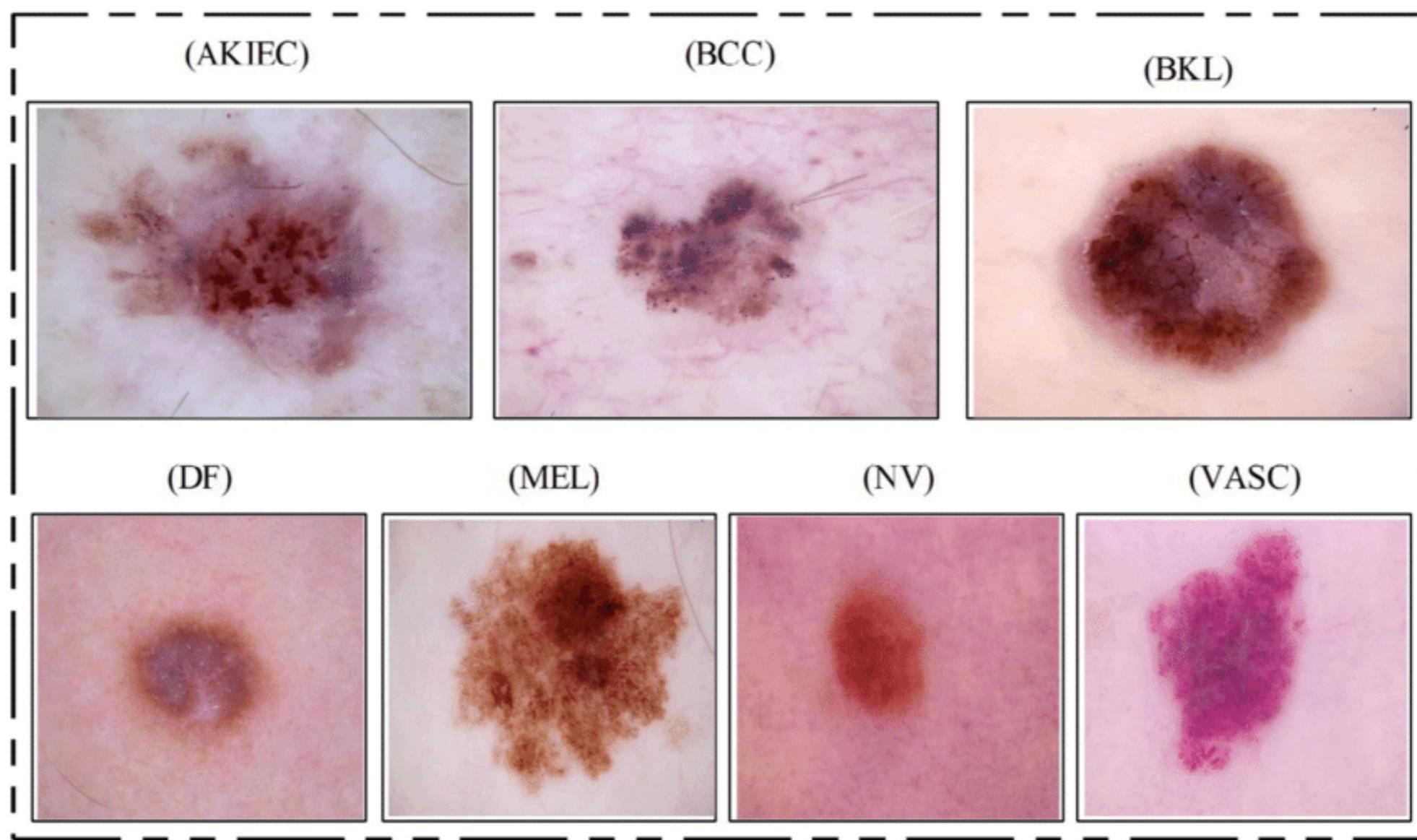
model

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

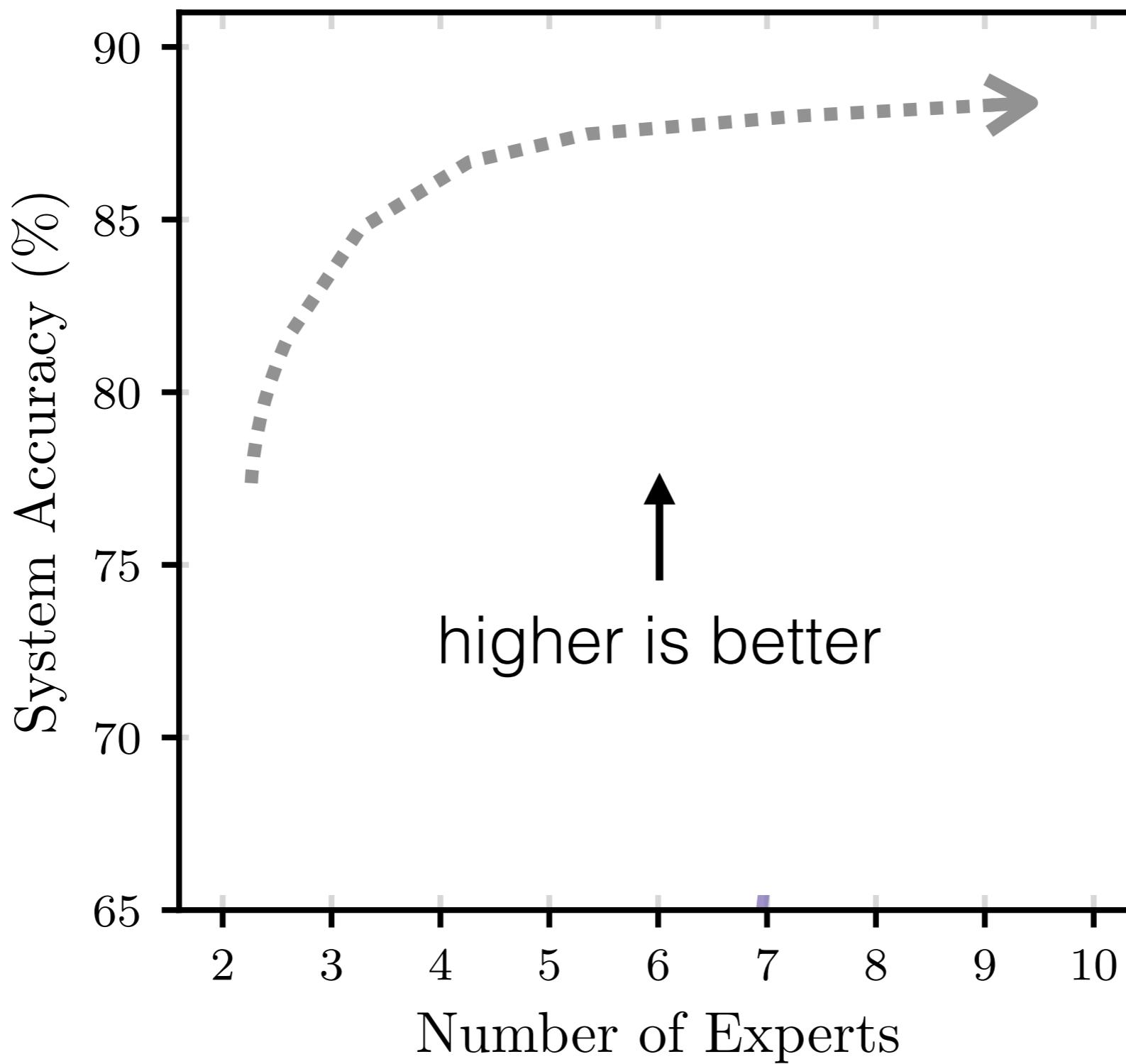
softmax loss function

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

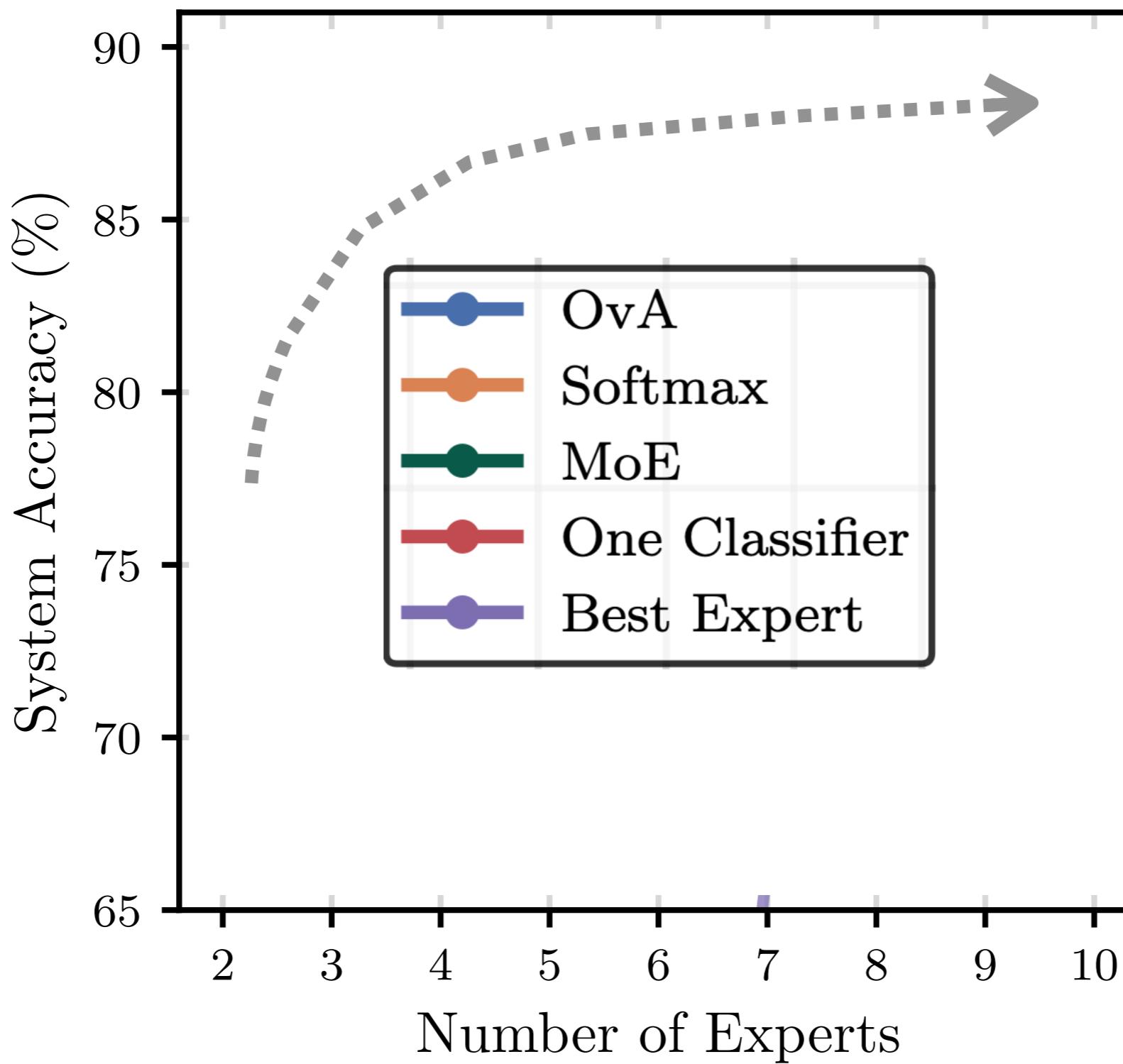
skin lesion diagnosis



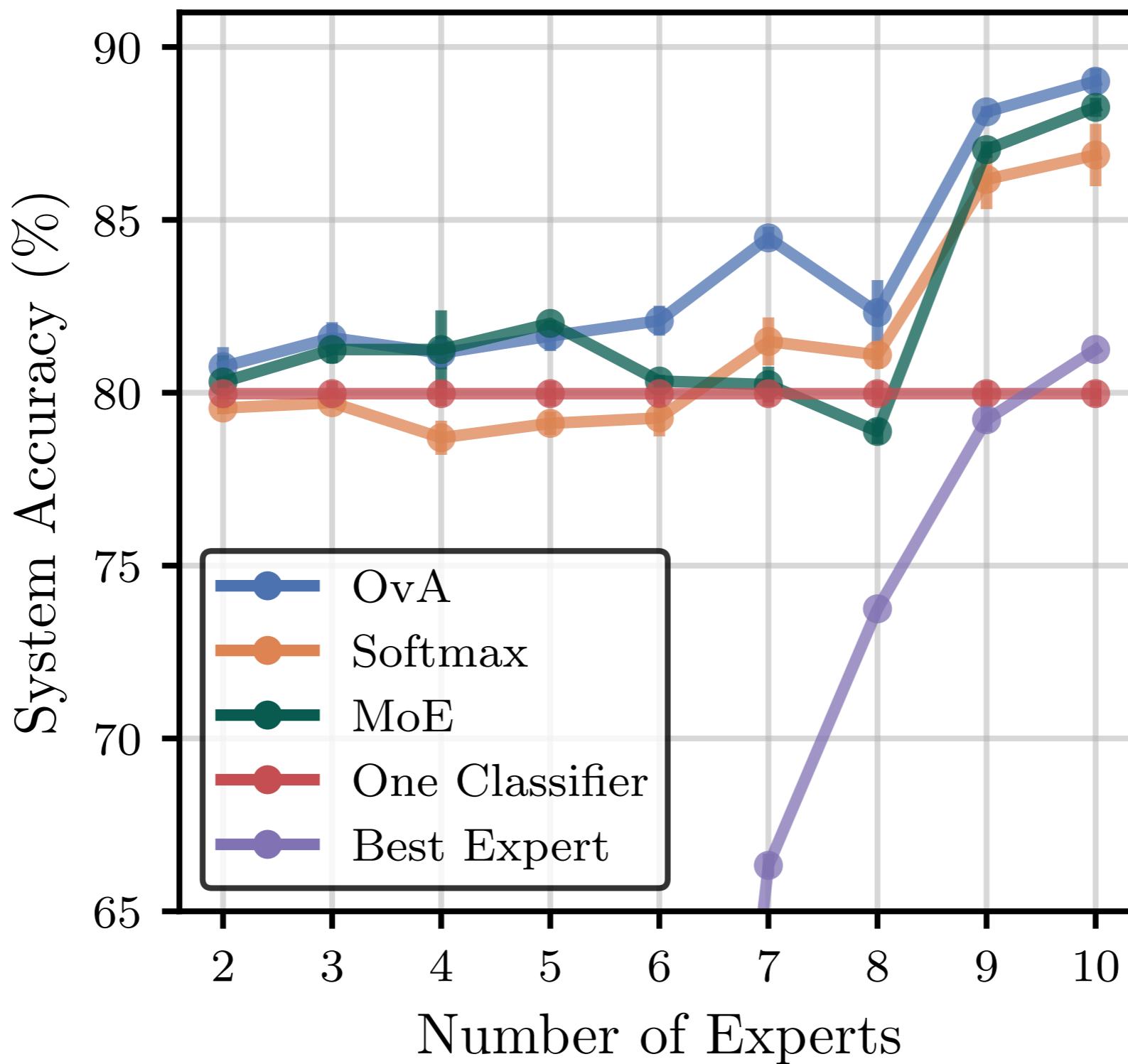
skin lesion diagnosis



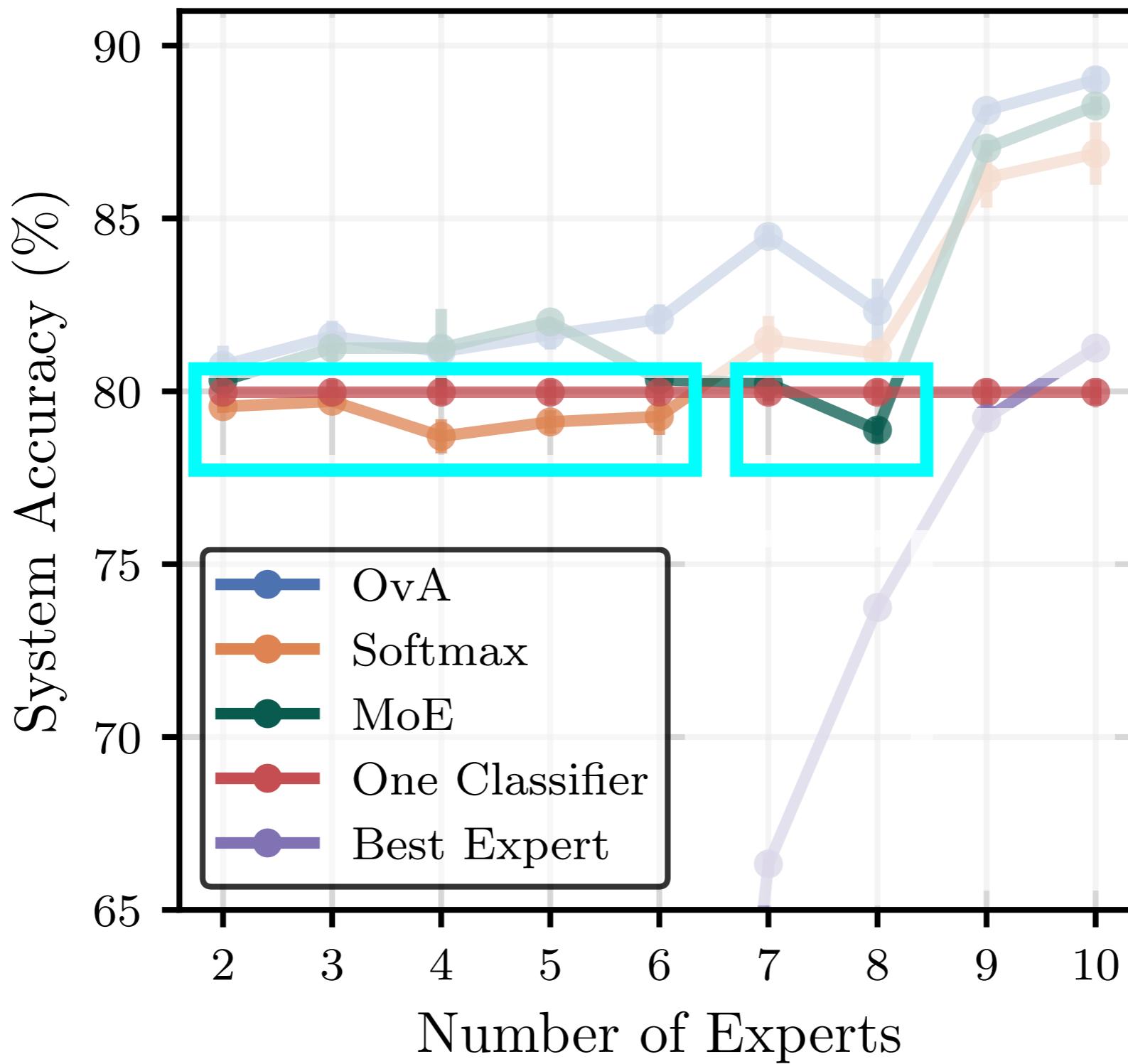
skin lesion diagnosis



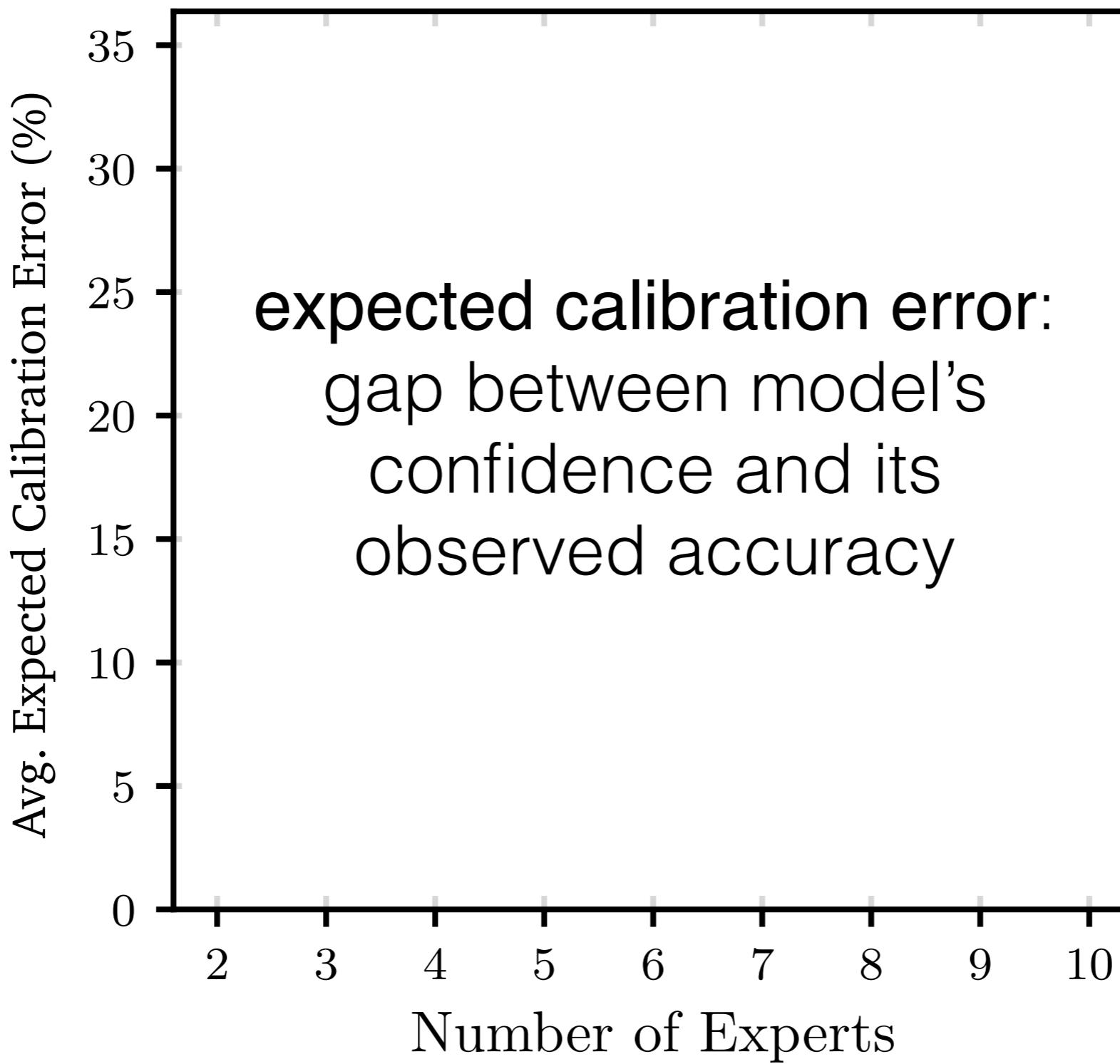
skin lesion diagnosis



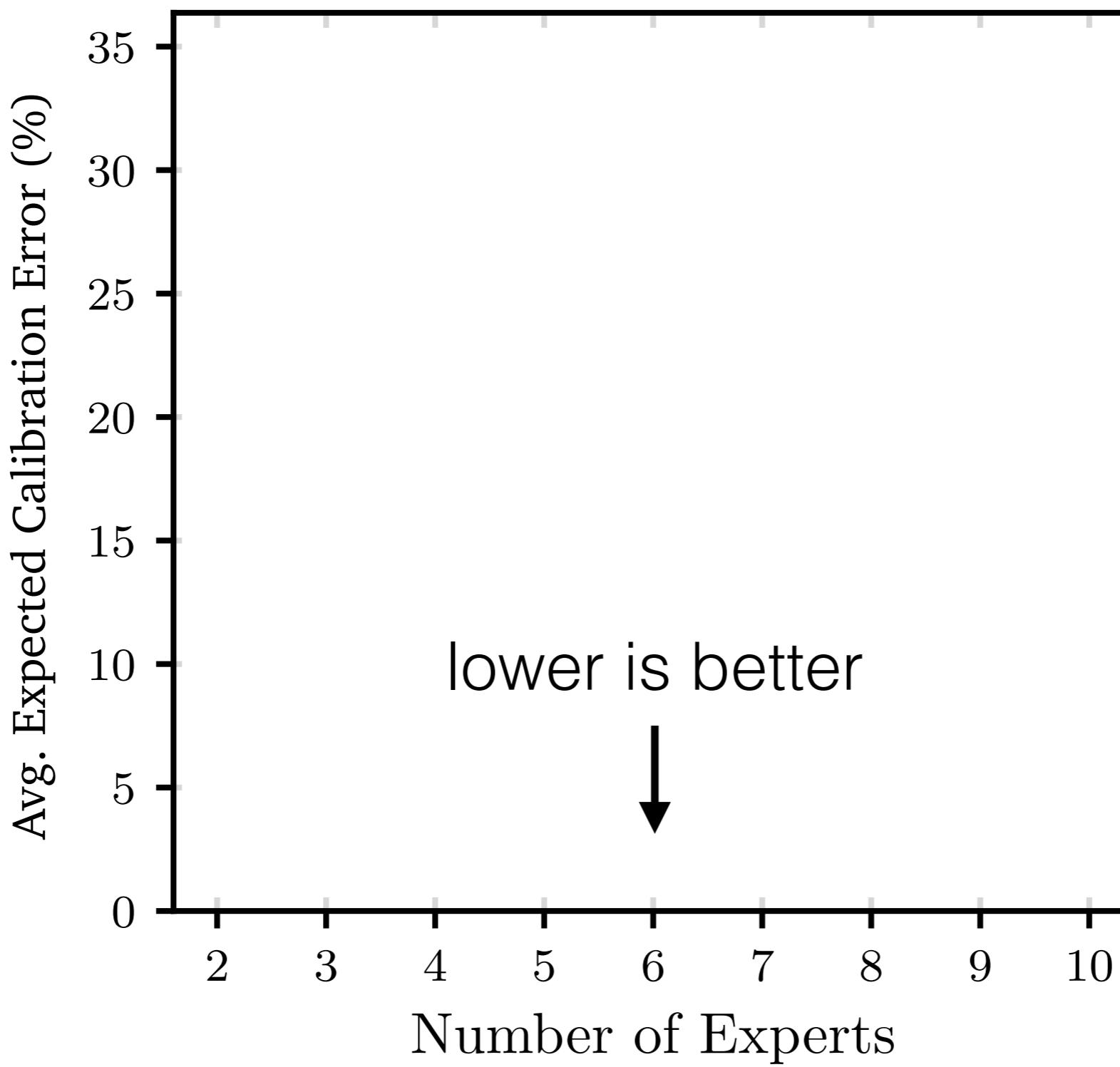
skin lesion diagnosis



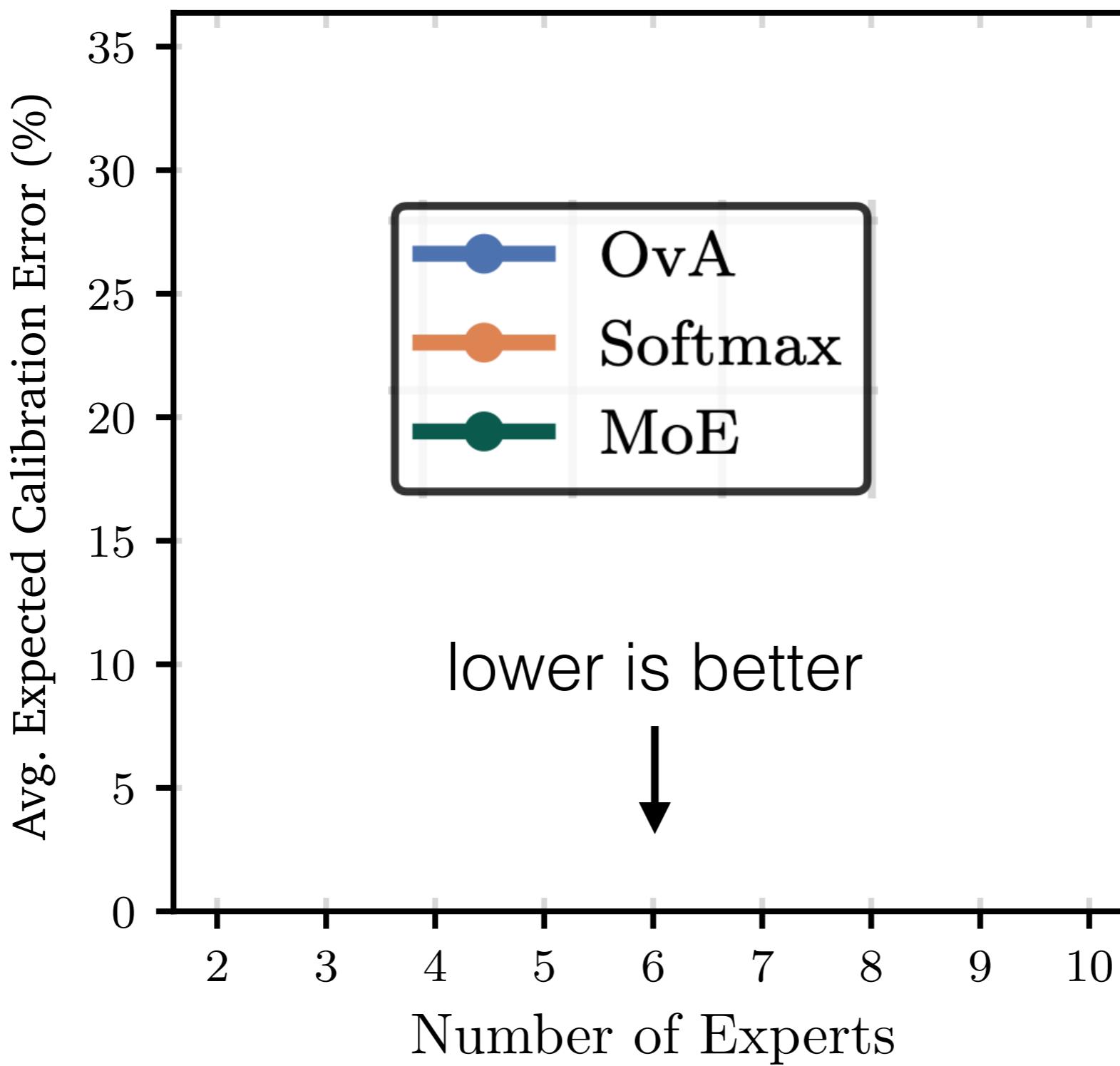
skin lesion diagnosis



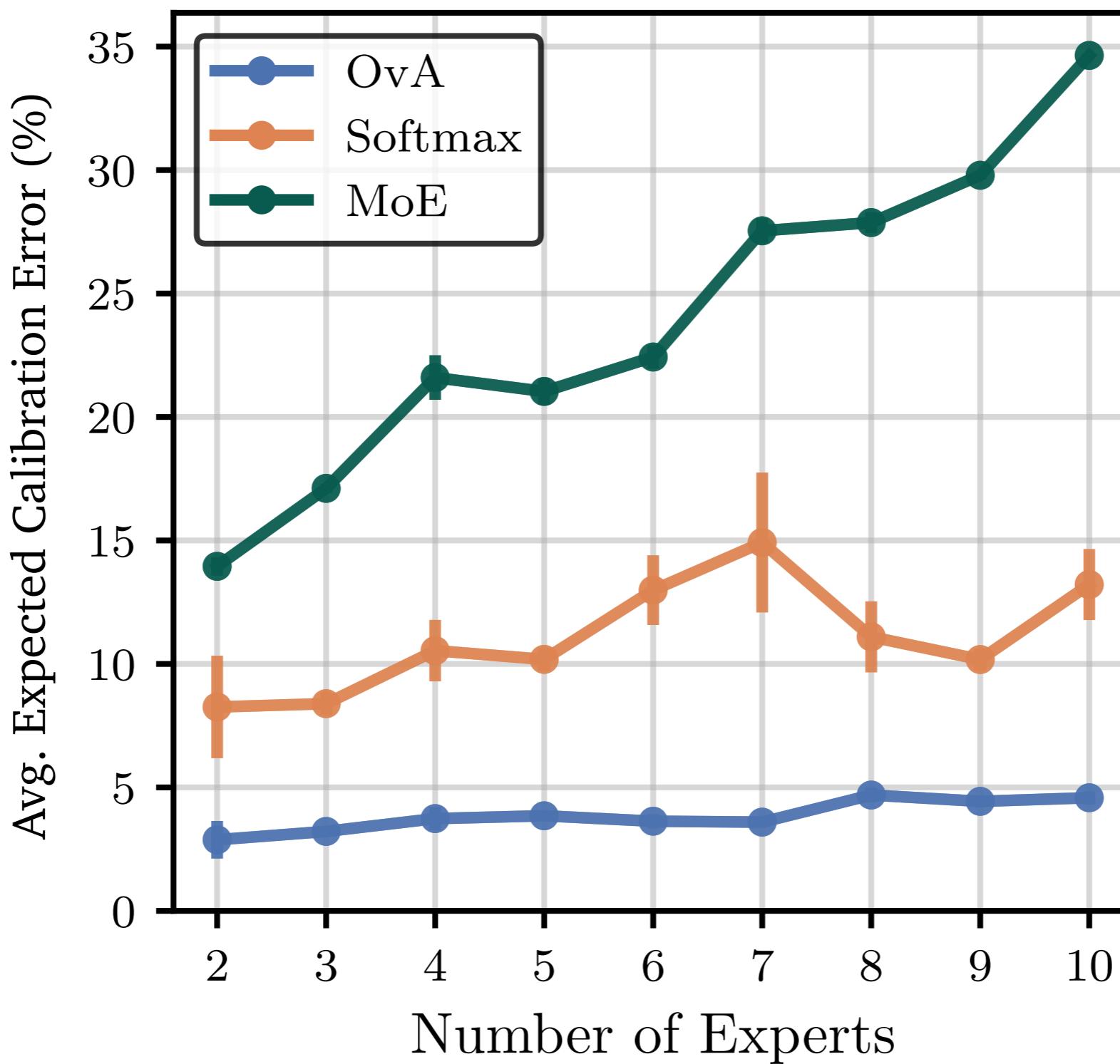
skin lesion diagnosis



skin lesion diagnosis



skin lesion diagnosis



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

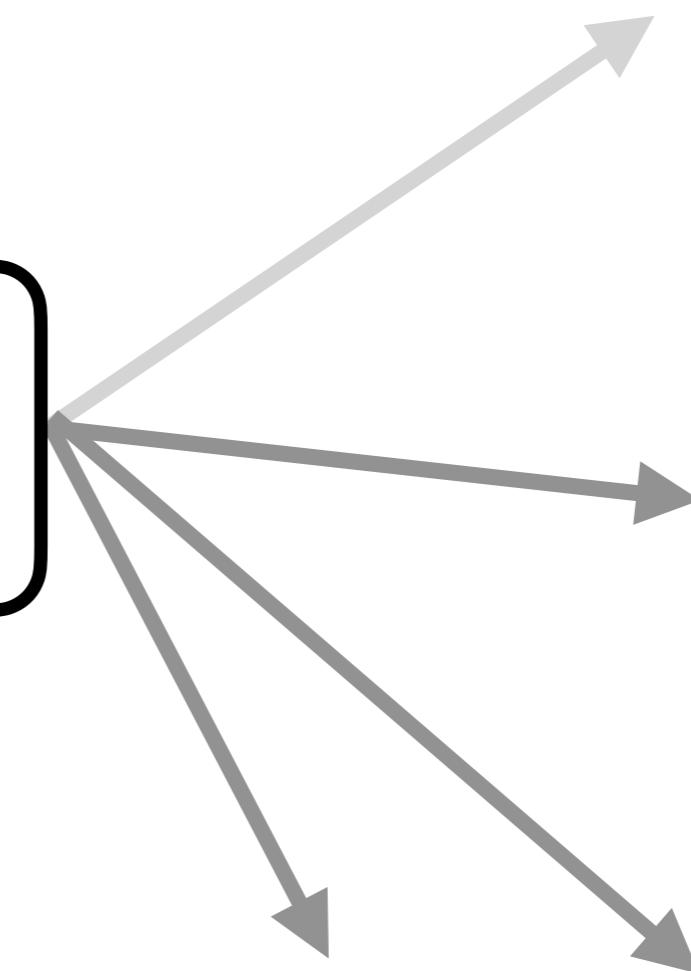
- ⊗ multiple experts
 - ⊗ surrogate losses
 - ⊗ conformal sets of experts

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

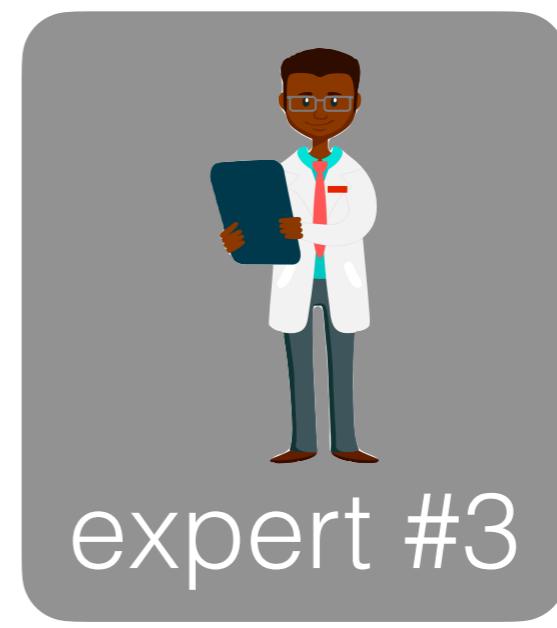
input
features



allocation
mechanism



expert #3



expert #2



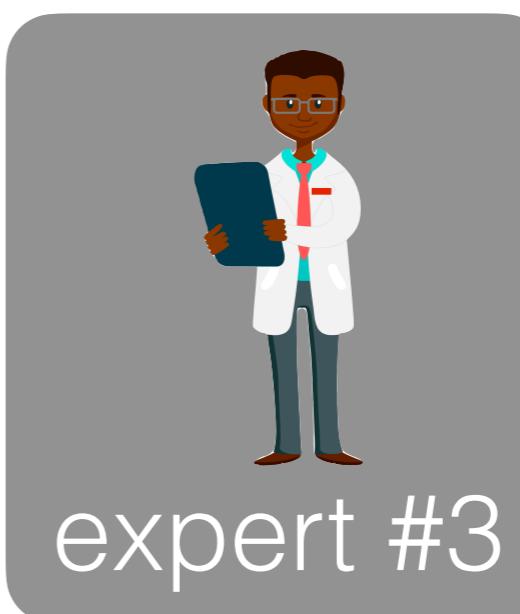
expert #1



input
features



allocation
mechanism



input
features



allocation
mechanism

expert #3



expert #2



expert #1



conformal inference

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

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

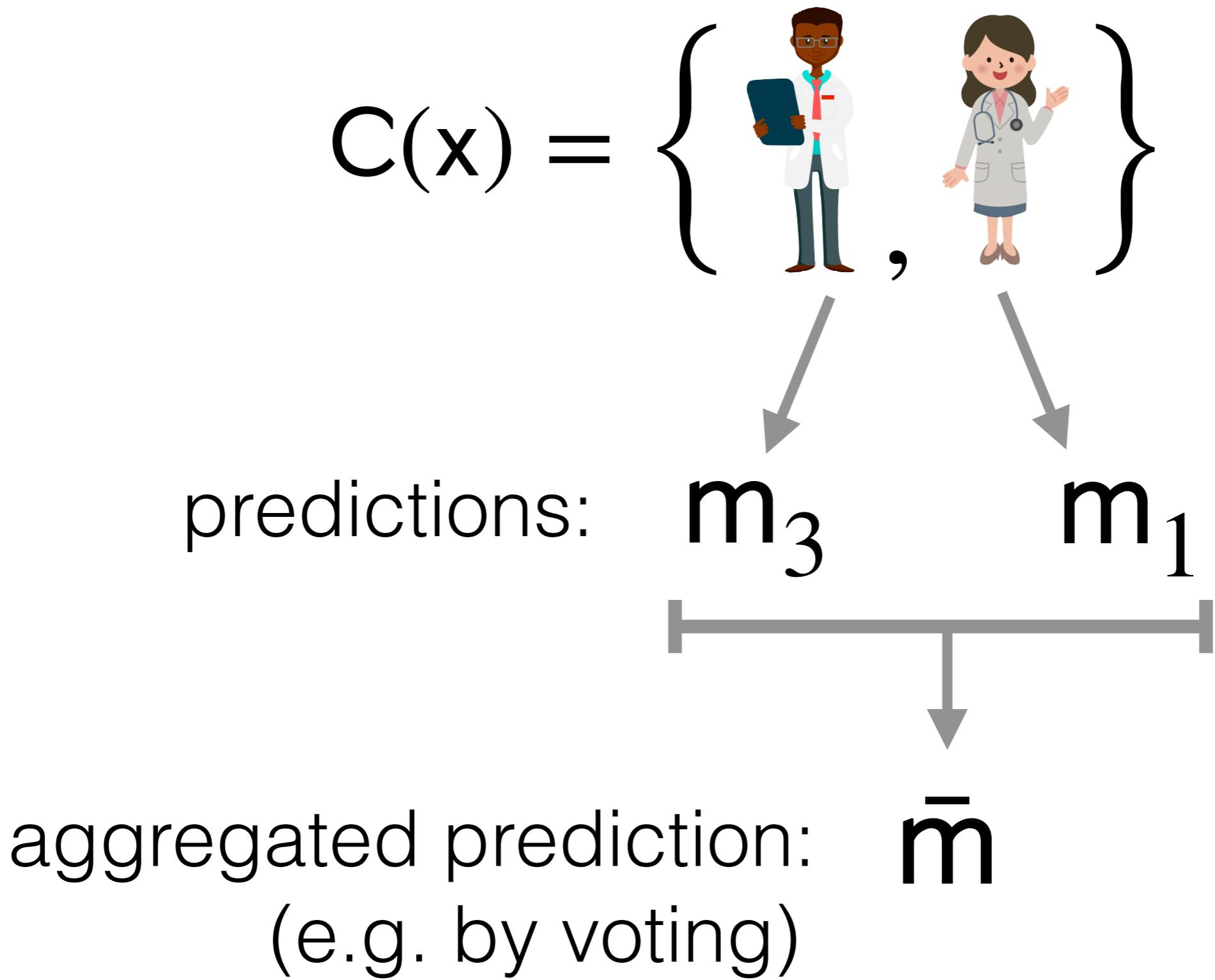
construct a confidence set of experts:

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



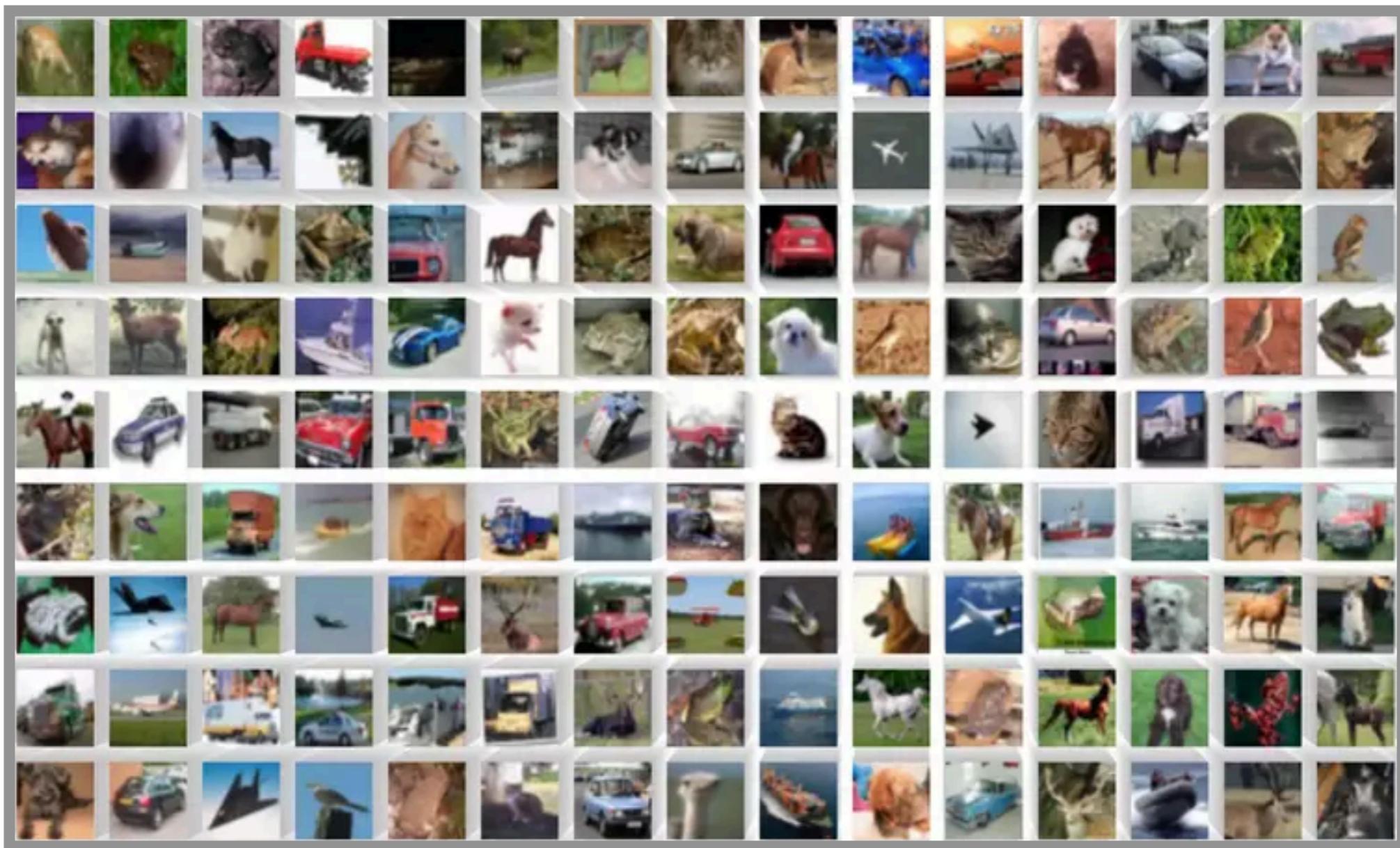
team of experts: adaptive in size and membership

conformal inference: ensembling

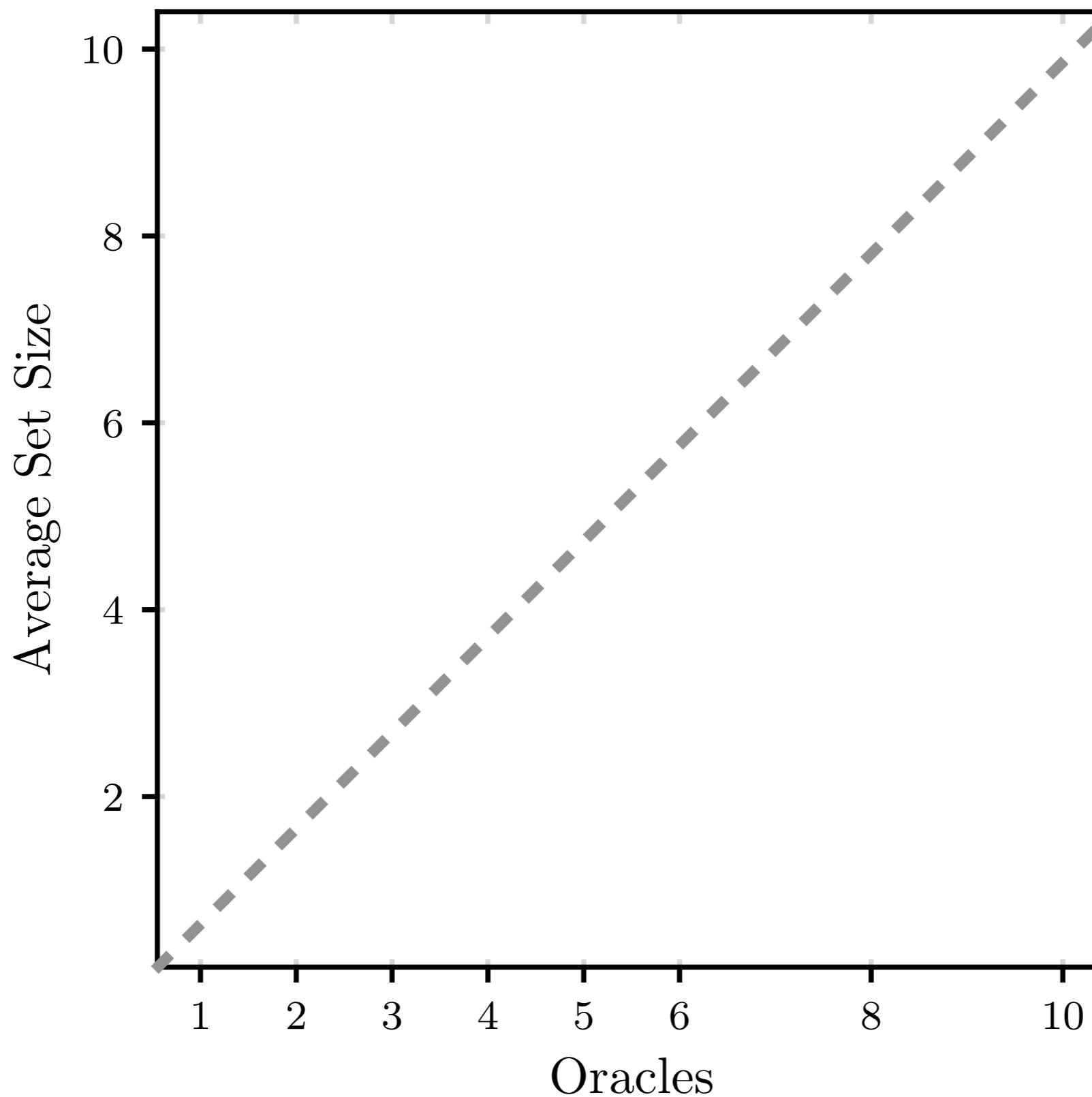


expert selection

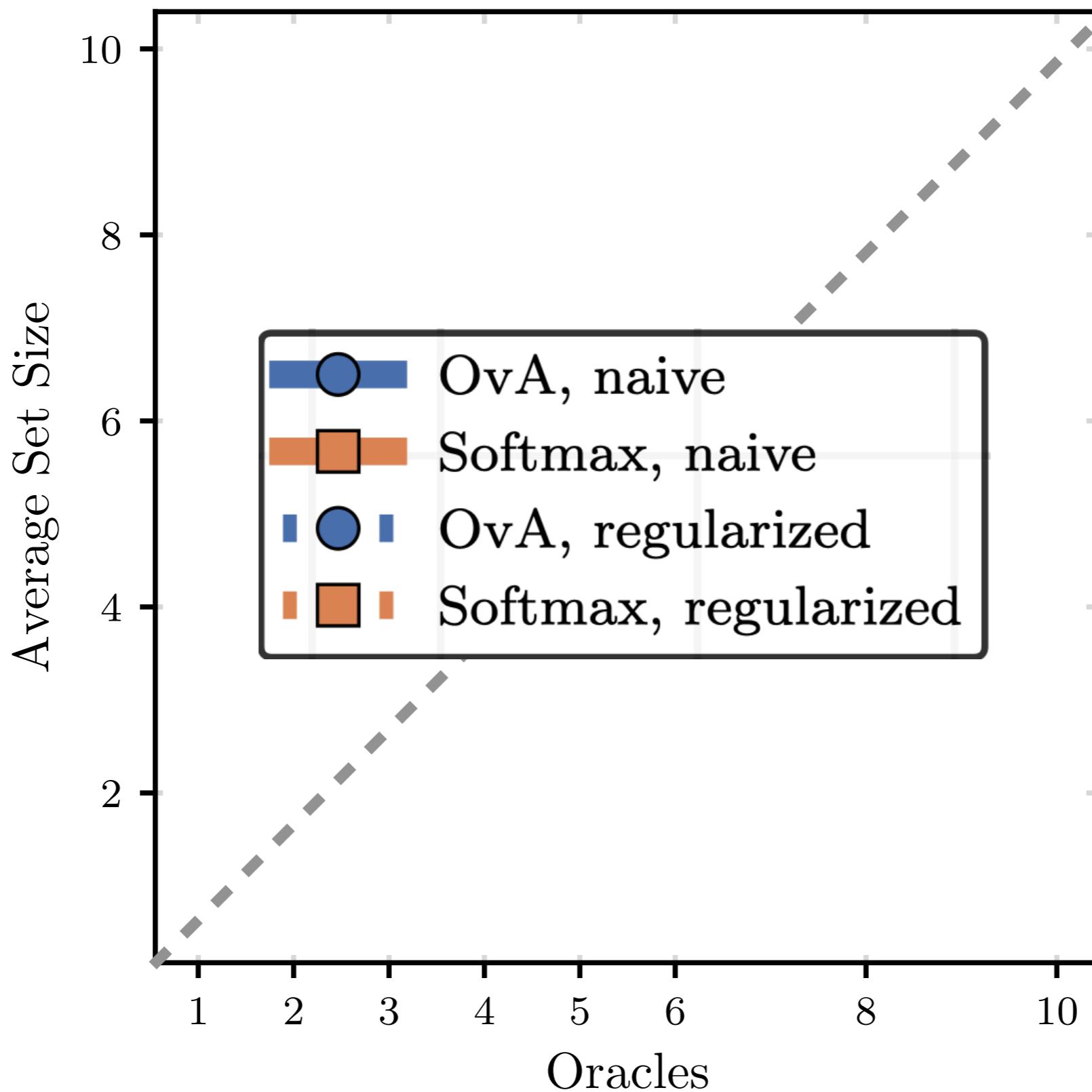
CIFAR-10



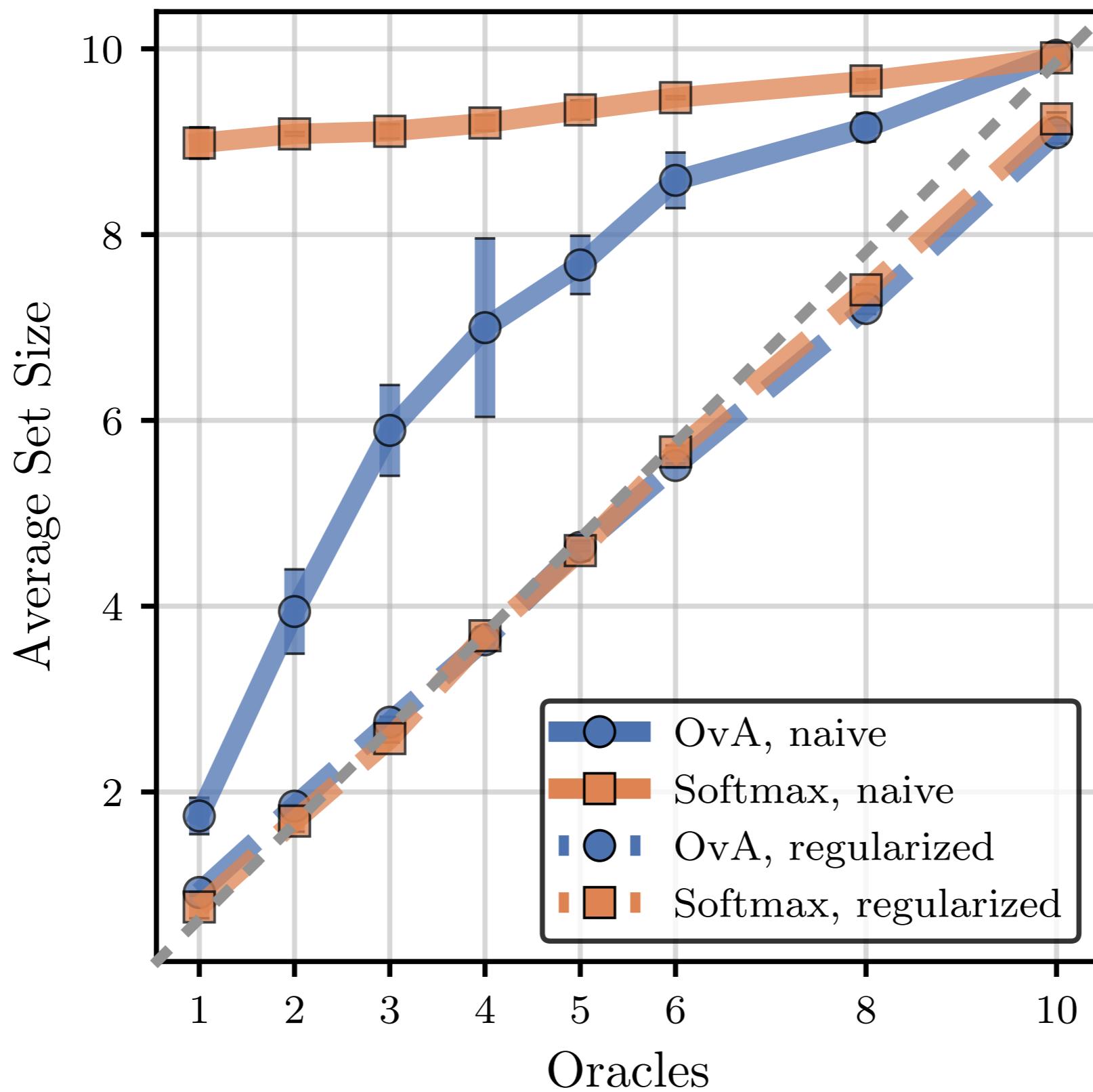
expert selection



expert selection



expert selection

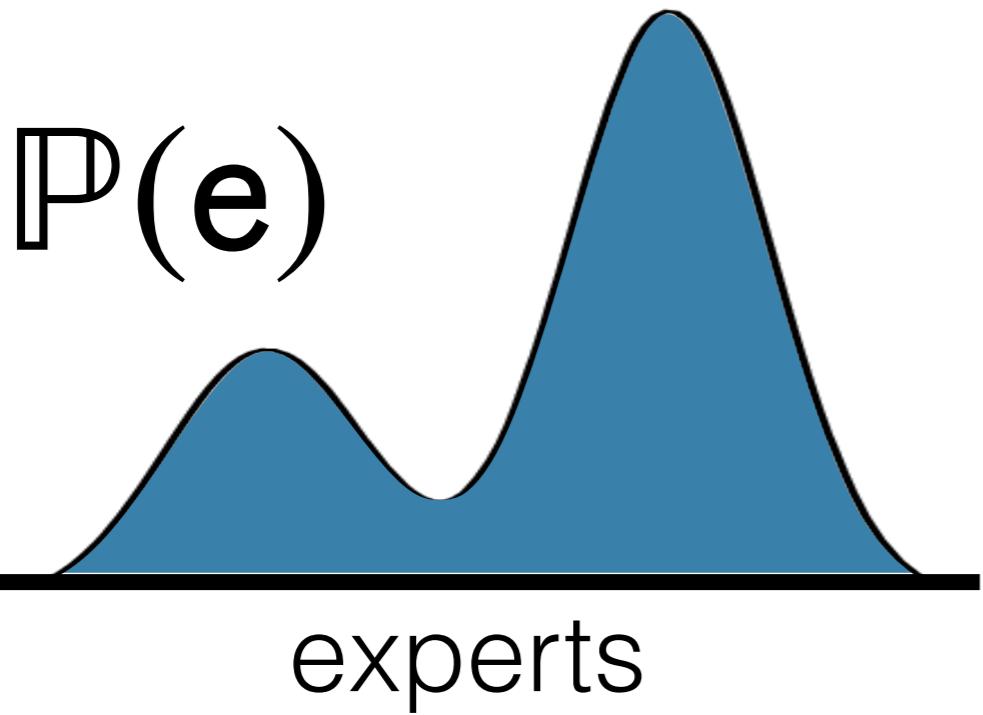


- ⊗ single expert
 - ⊗ softmax surrogate loss
 - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
 - ⊗ surrogate losses
 - ⊗ conformal sets of experts
- ⊗ population of experts
 - ⊗ surrogate losses
 - ⊗ meta-learning a rejector

input
features



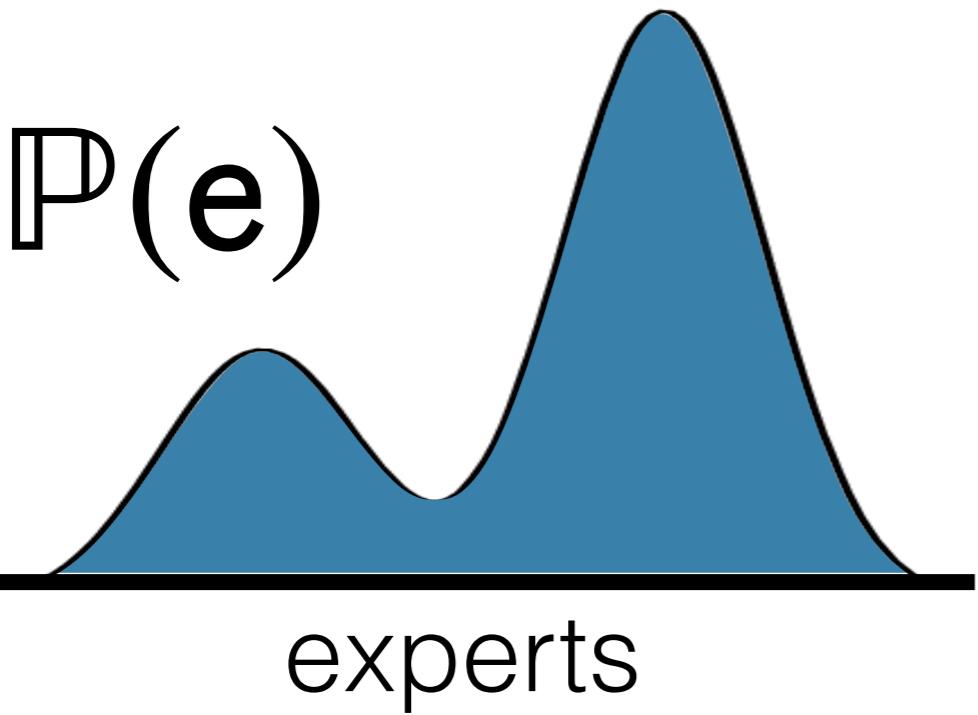
allocation
mechanism



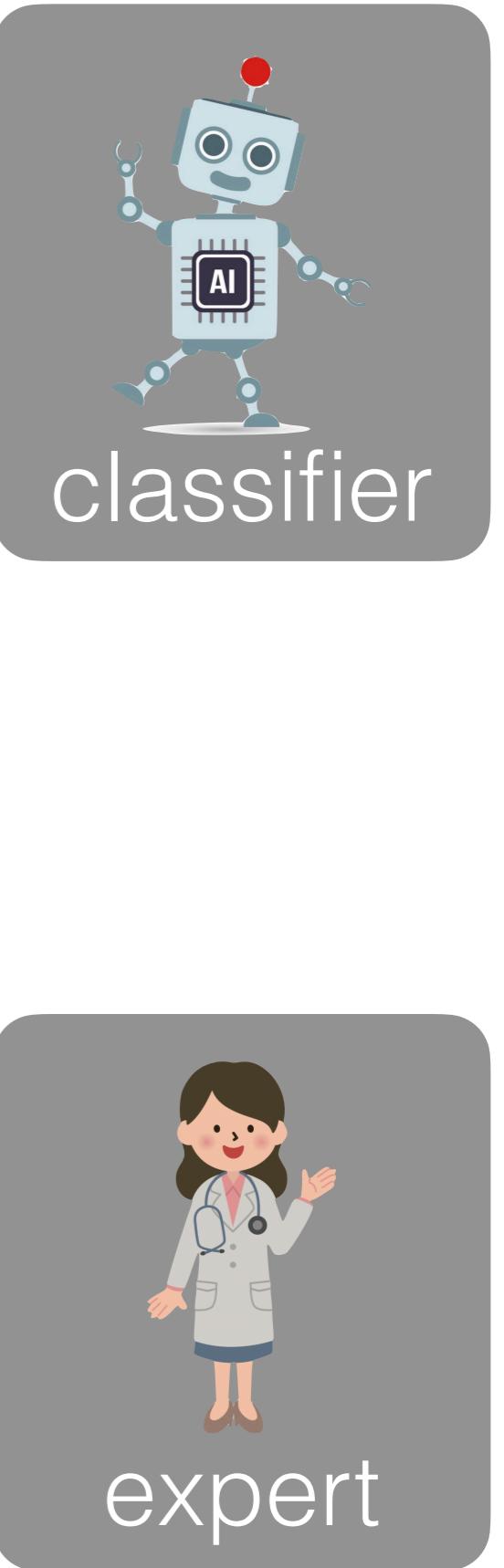
input
features



allocation
mechanism



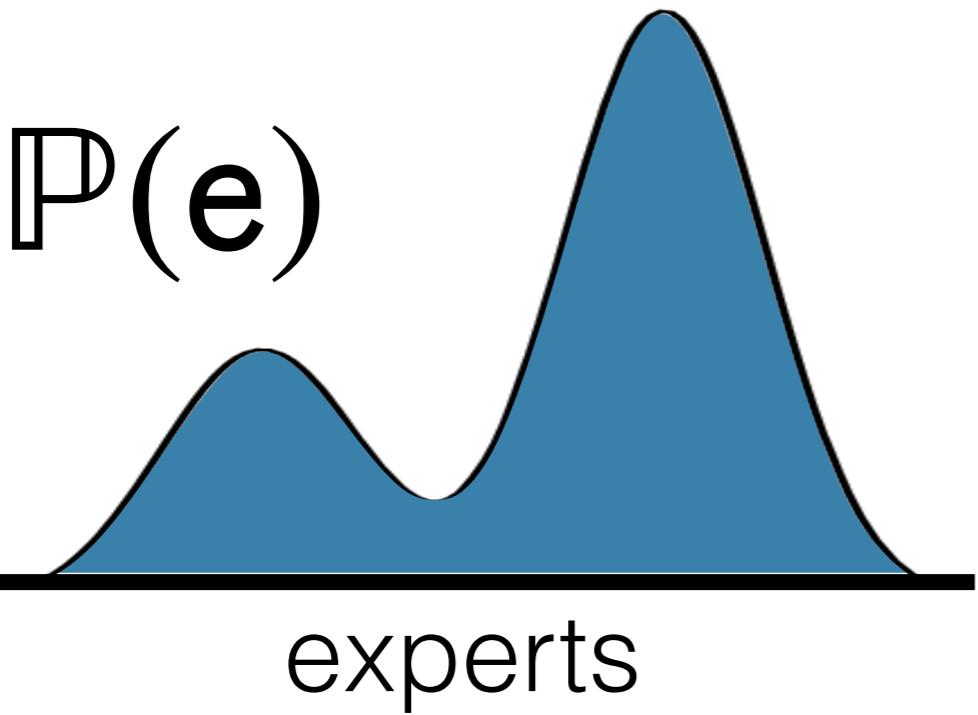
$e \sim P(e)$



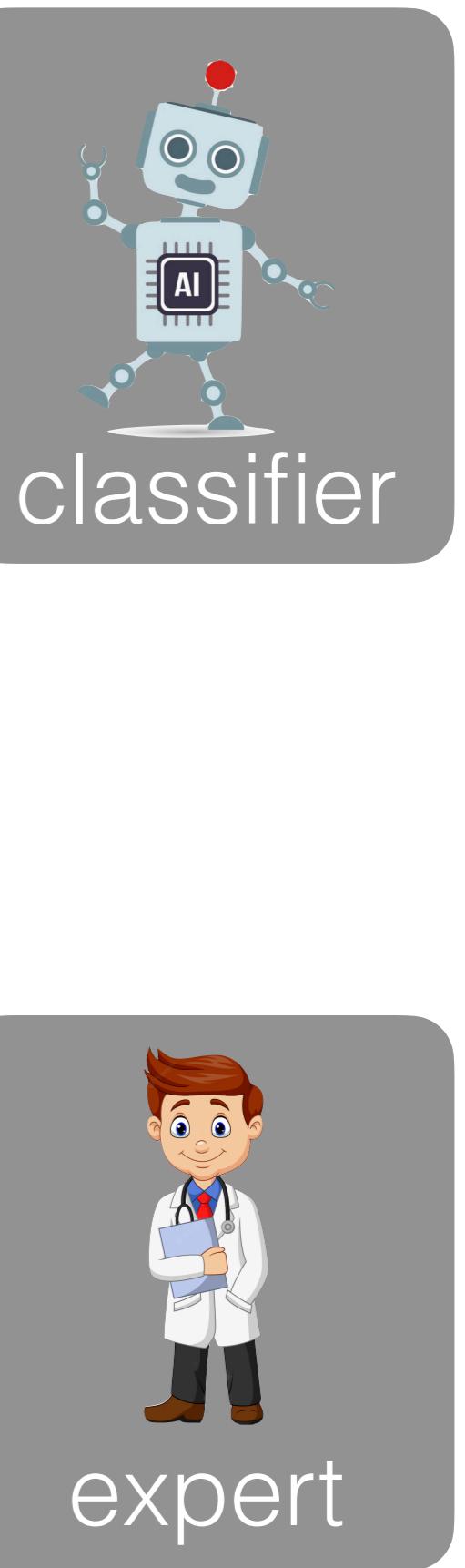
input
features



allocation
mechanism



$$e \sim P(e)$$



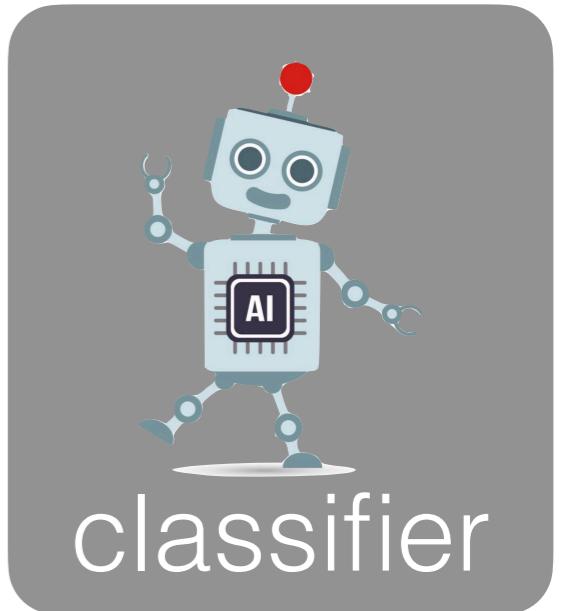
input
features



allocation
mechanism

$$e \sim P(e)$$

$$m \sim P(m | e)$$



input
features



allocation
mechanism

L_{0-1}

classifier

L_{0-1}

?

expert

Bayes optimal deferral rule:

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

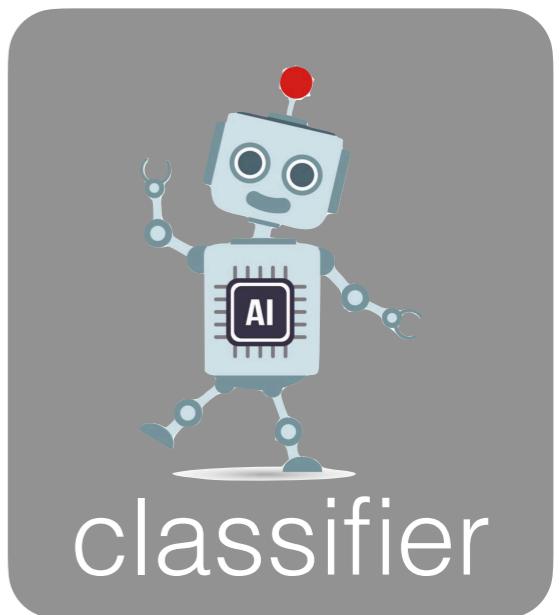
input
features



allocation
mechanism

defer to expert if...

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



input
features

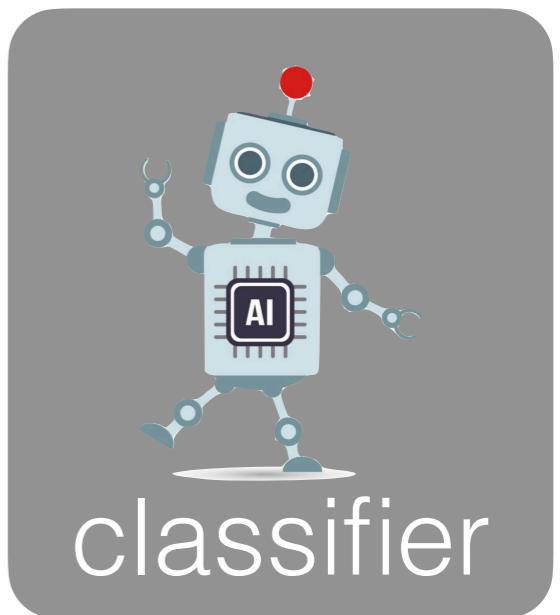


allocation
mechanism

defer to expert if...

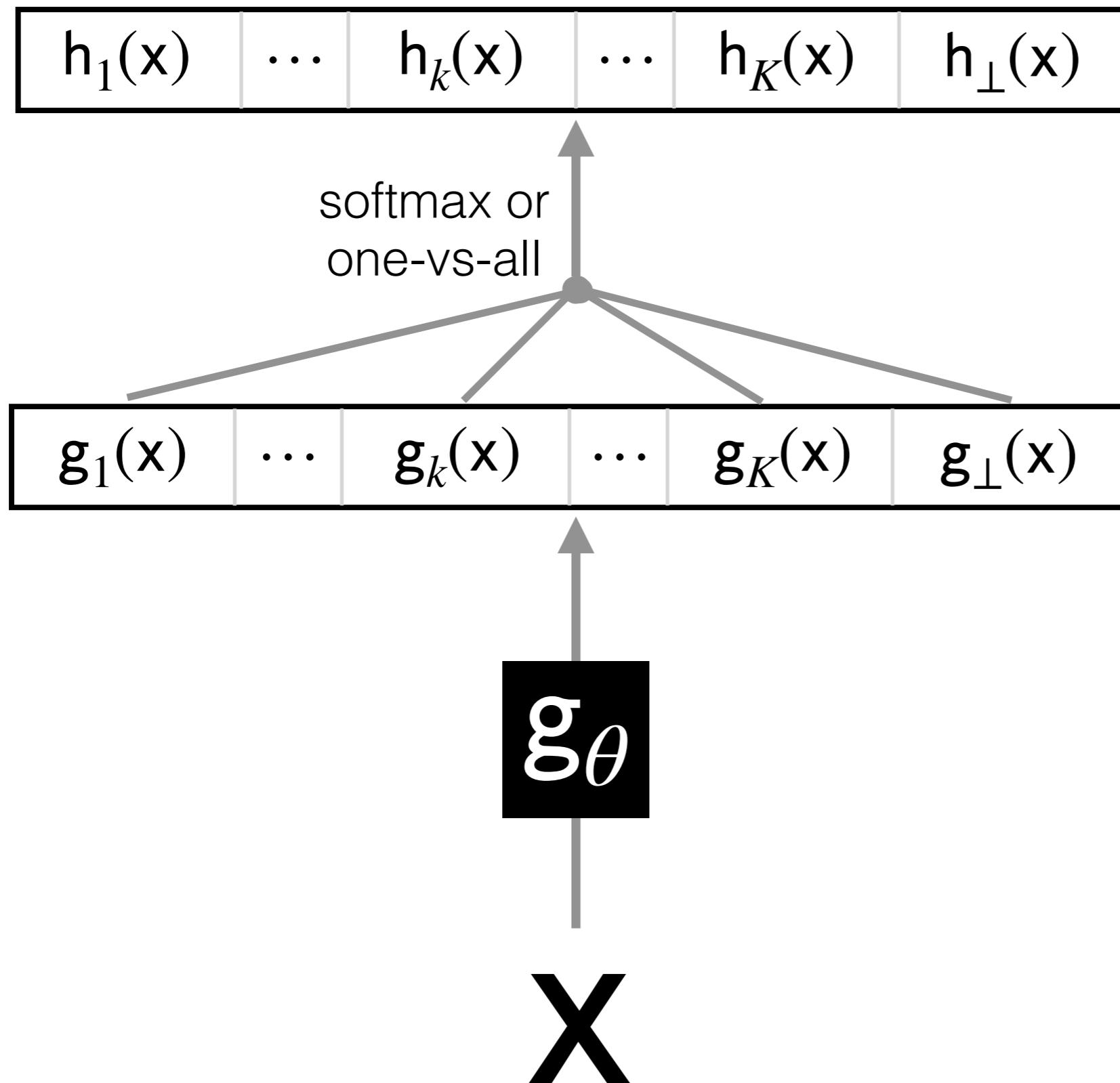
$$\max_{y \in [1, K]} h_y(x) \leq h_\perp(x, e)$$

?

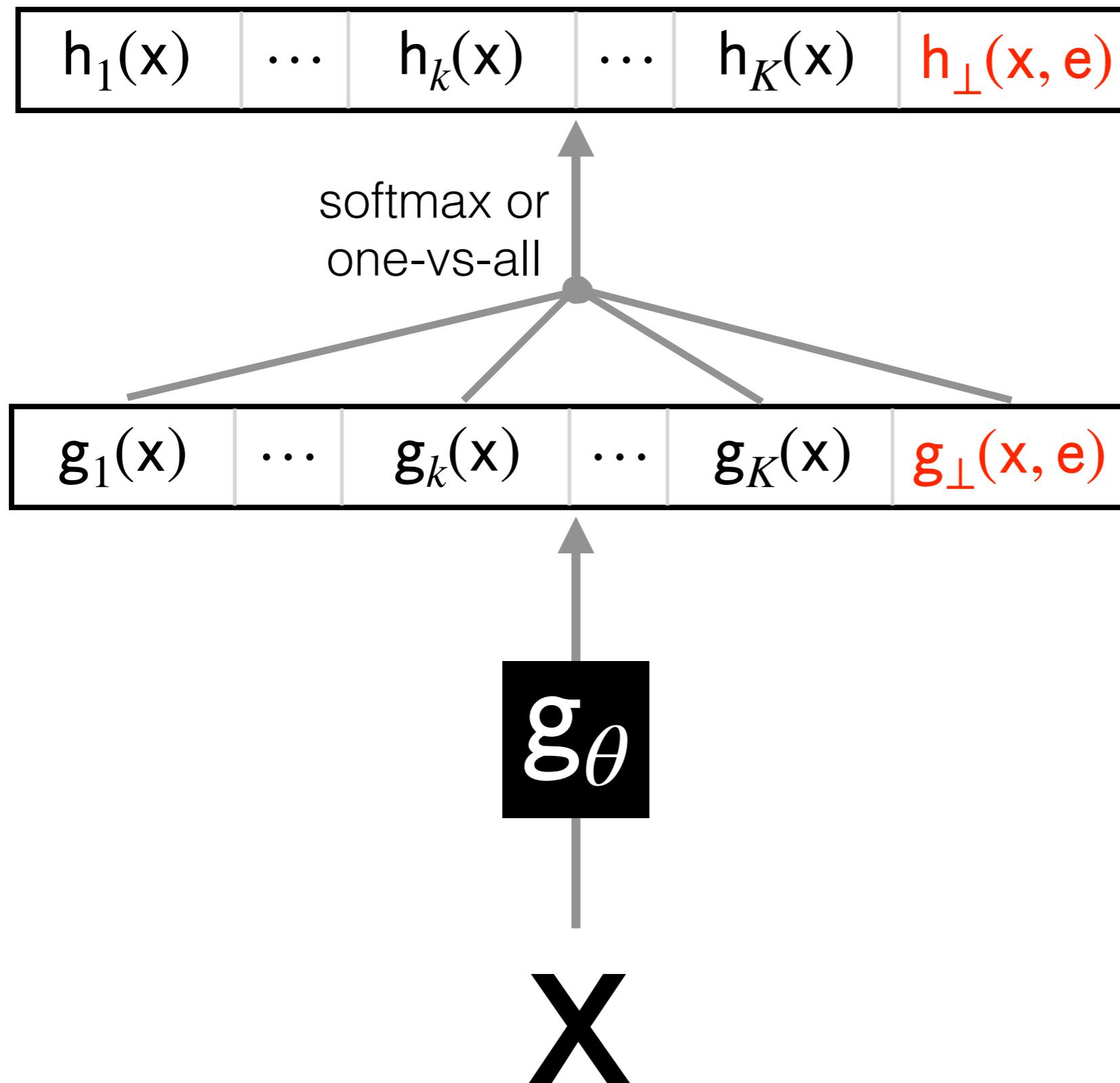


- ⊗ single expert
 - ⊗ softmax surrogate loss
 - ⊗ improving calibration via one-vs-all
- ⊗ multiple experts
 - ⊗ surrogate losses
 - ⊗ conformal sets of experts
- ⊗ population of experts
 - ⊗ surrogate losses
 - ⊗ meta-learning a rejector

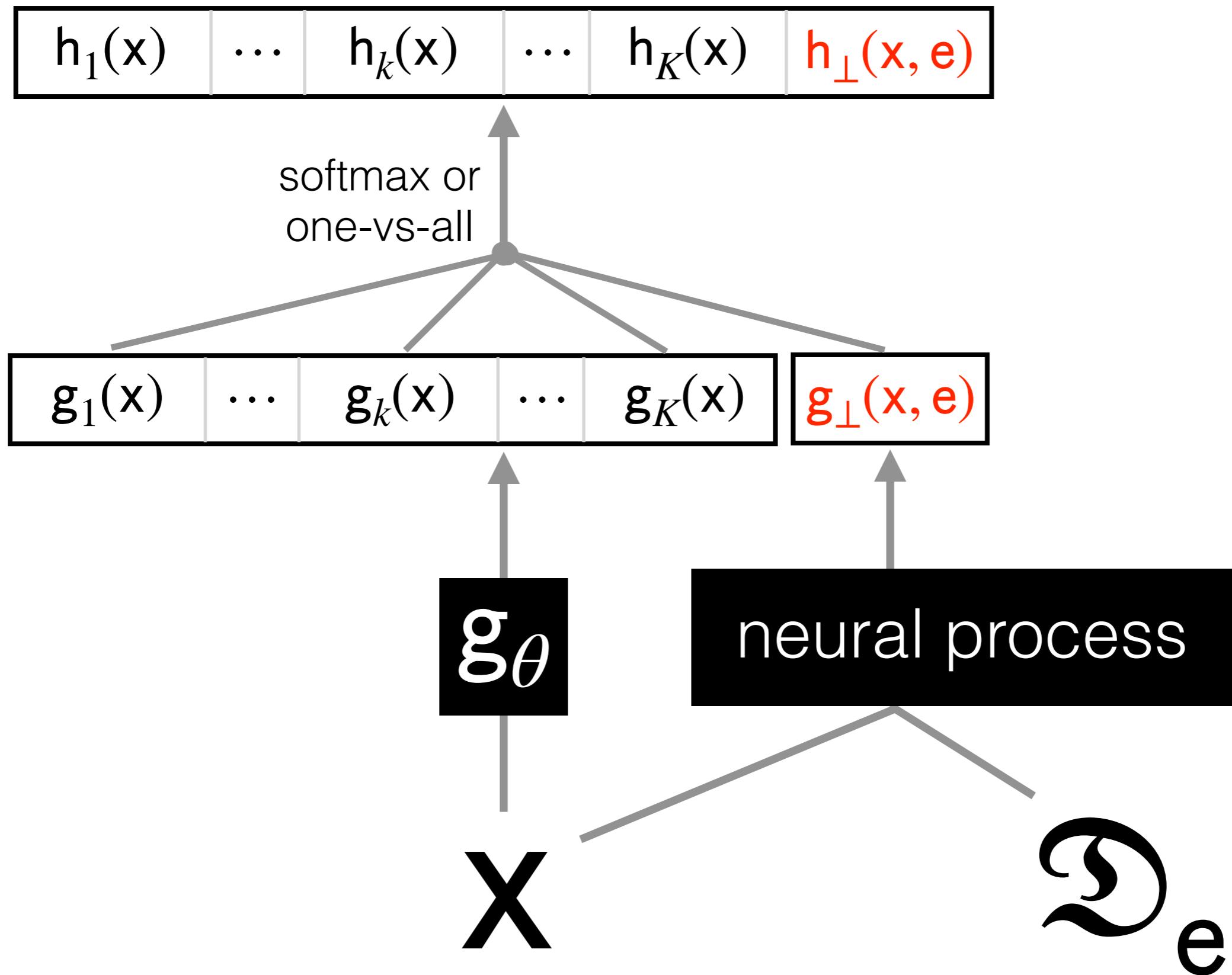
meta-learning implementation



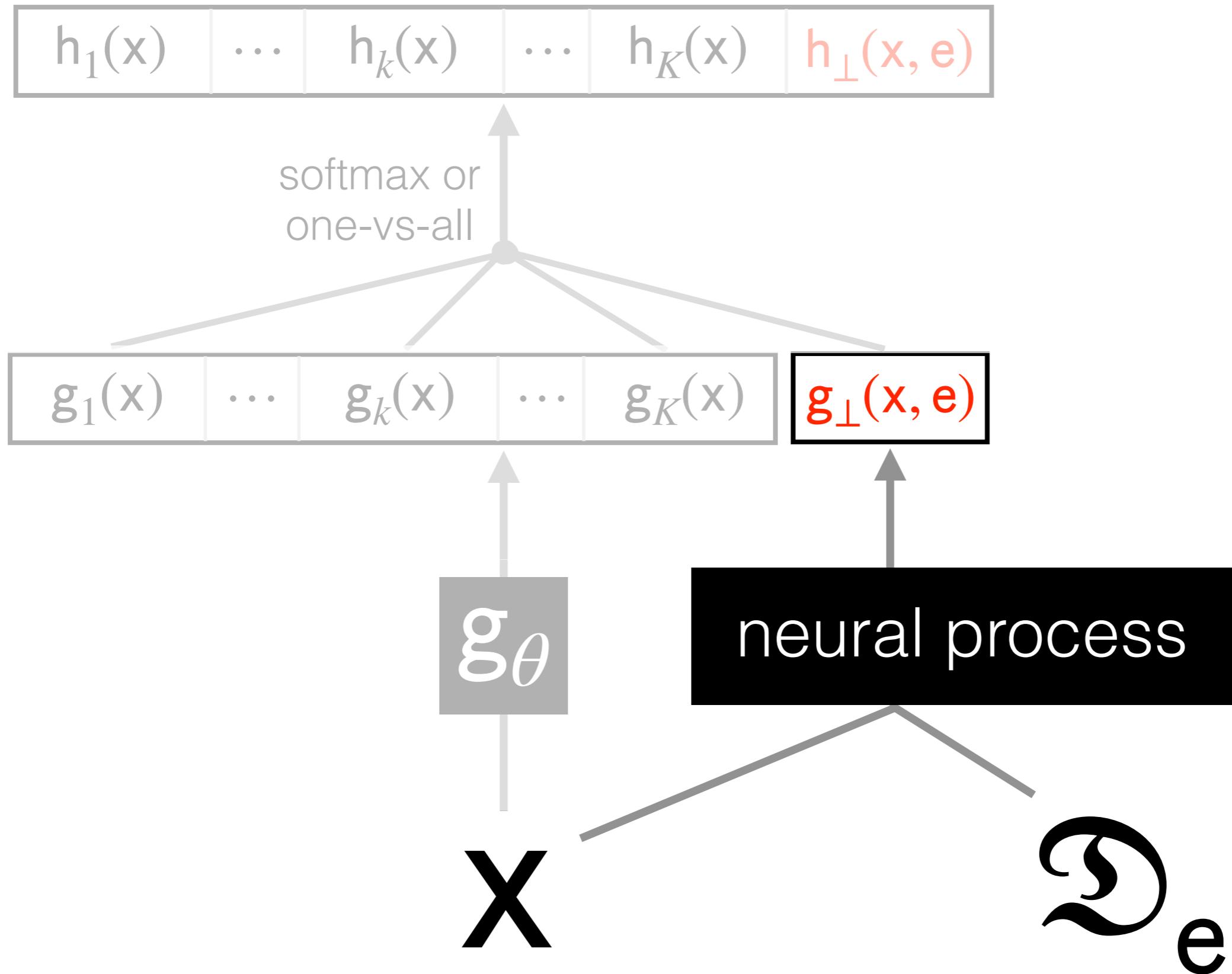
meta-learning implementation



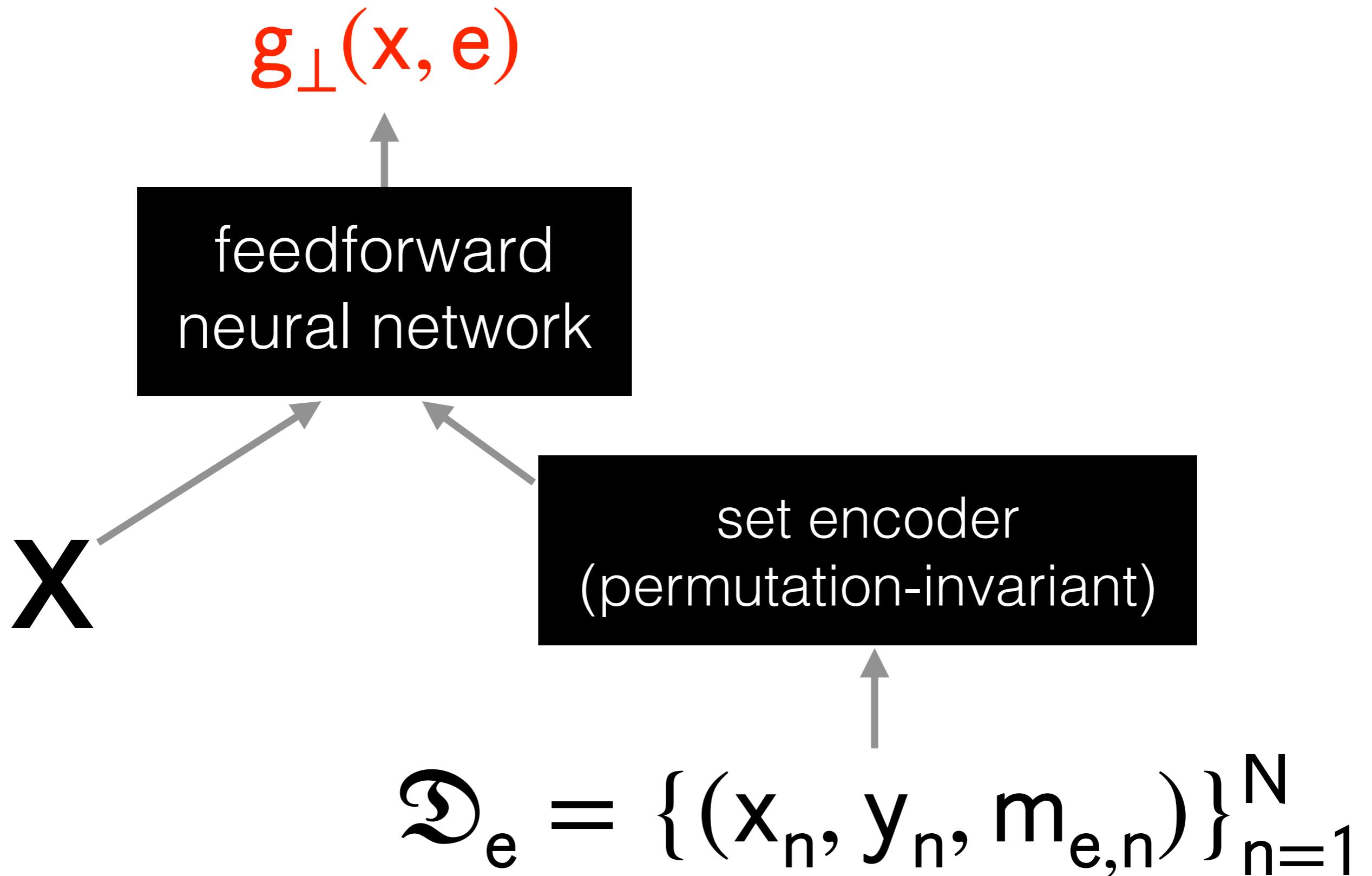
meta-learning implementation



meta-learning implementation



neural process rejector



- ⊗ 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



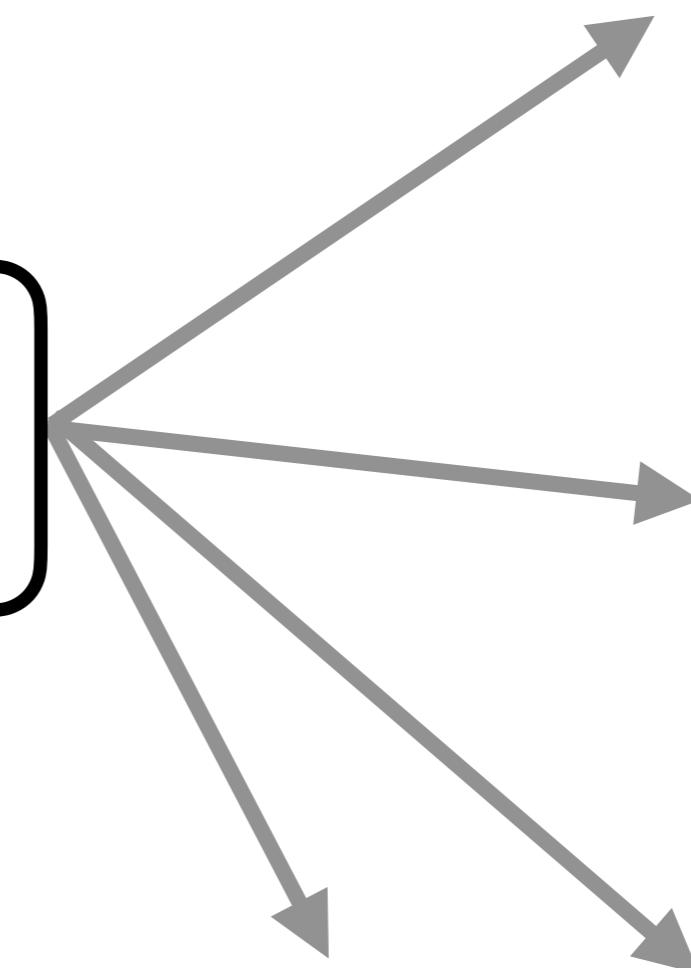
allocation
mechanism



input
features



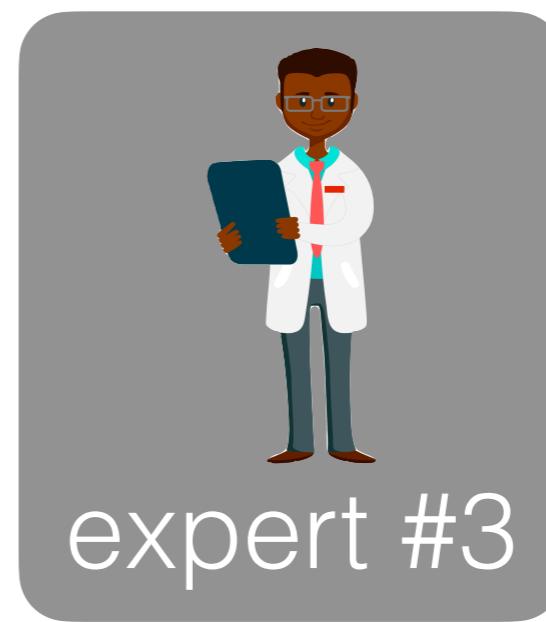
allocation
mechanism



classifier



expert #1



expert #3



expert #2

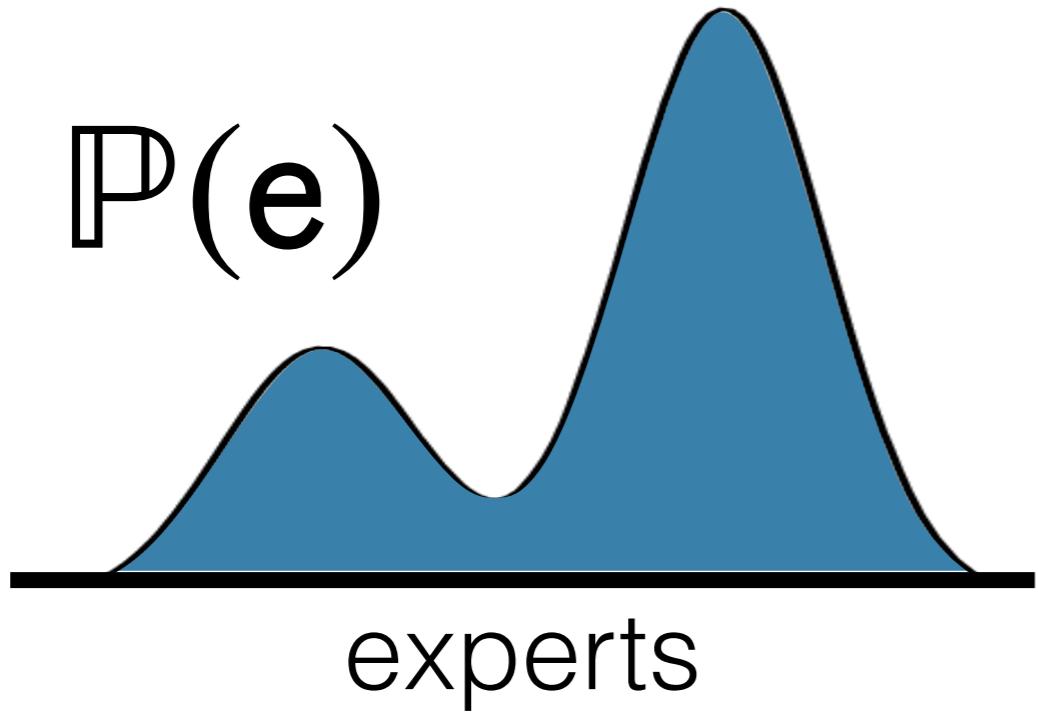
input
features



allocation
mechanism

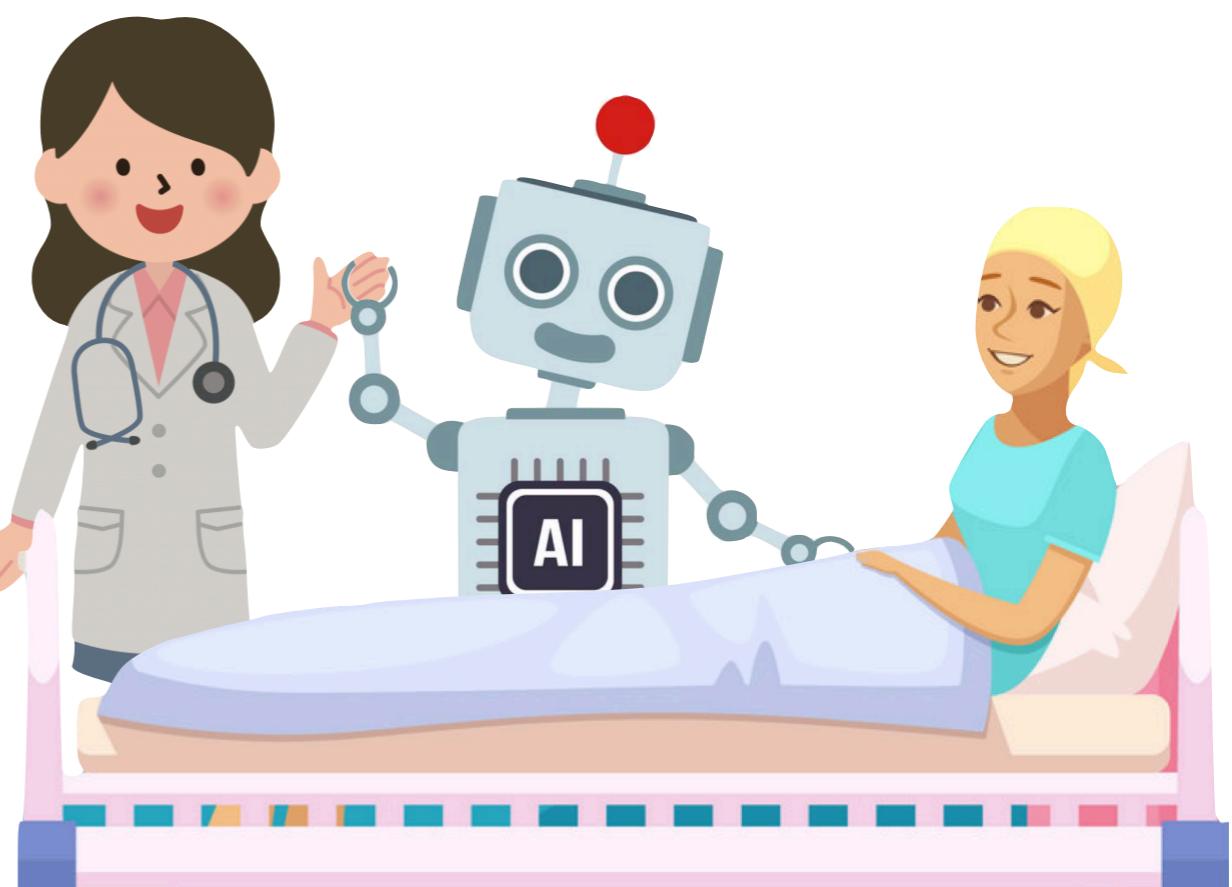
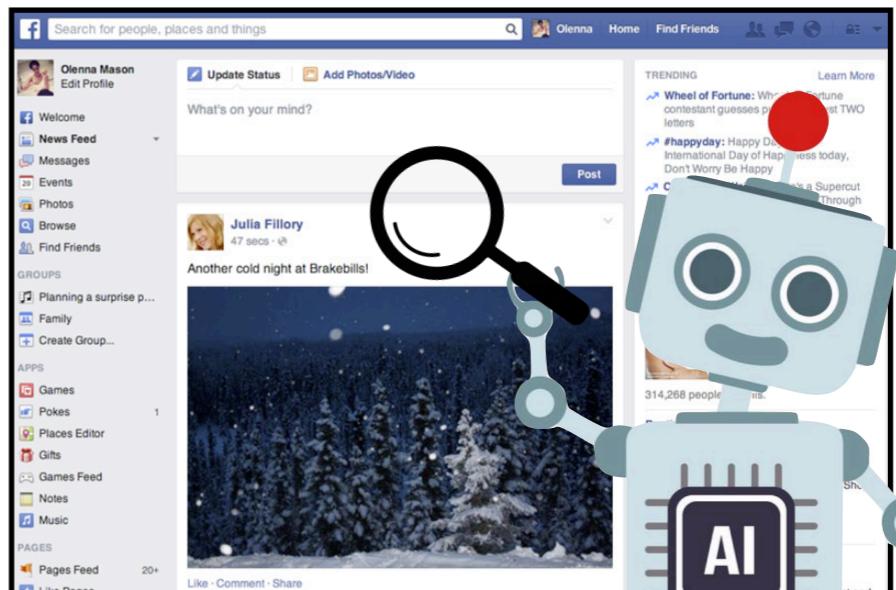
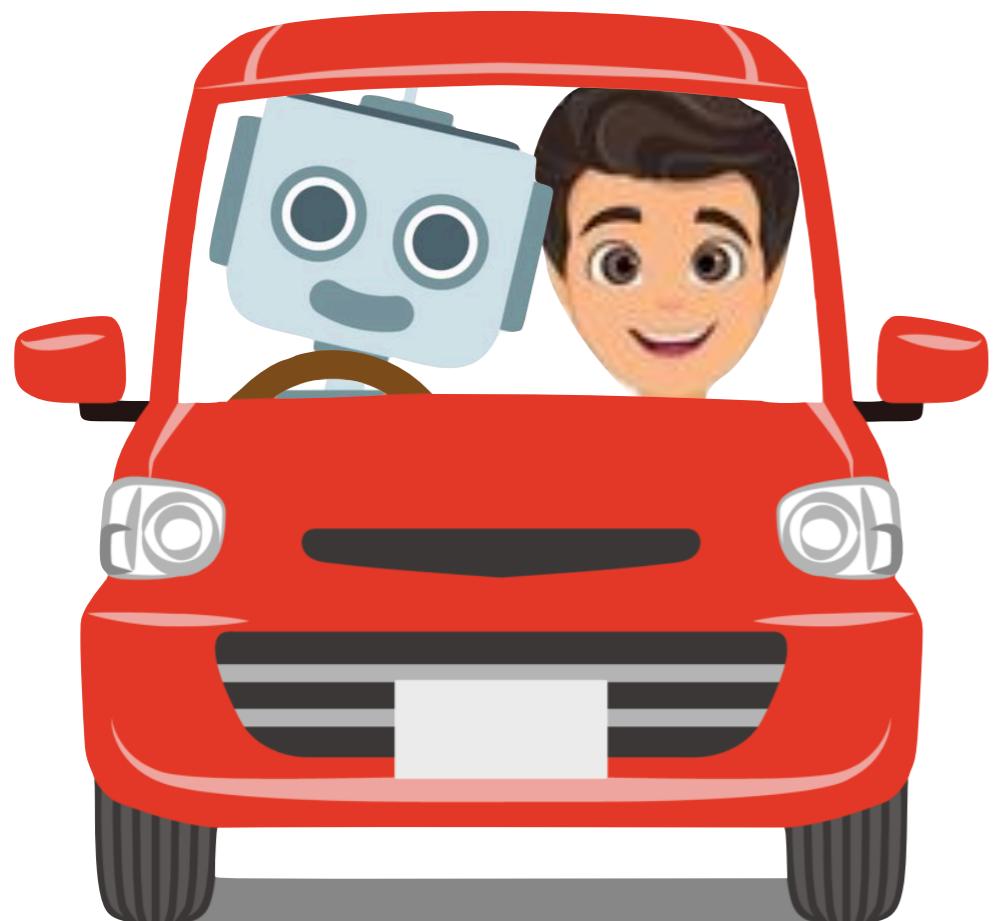


$P(e)$



$e \sim P(e)$





papers & code



funding provided by



co-authors



Rajeev
Verma



Daniel
Barrejón



Dharmesh
Tailor



Putra
Manggala



Aditya
Patra