# Towards Anytime Uncertainty Estimation in Early-Exit Neural Networks
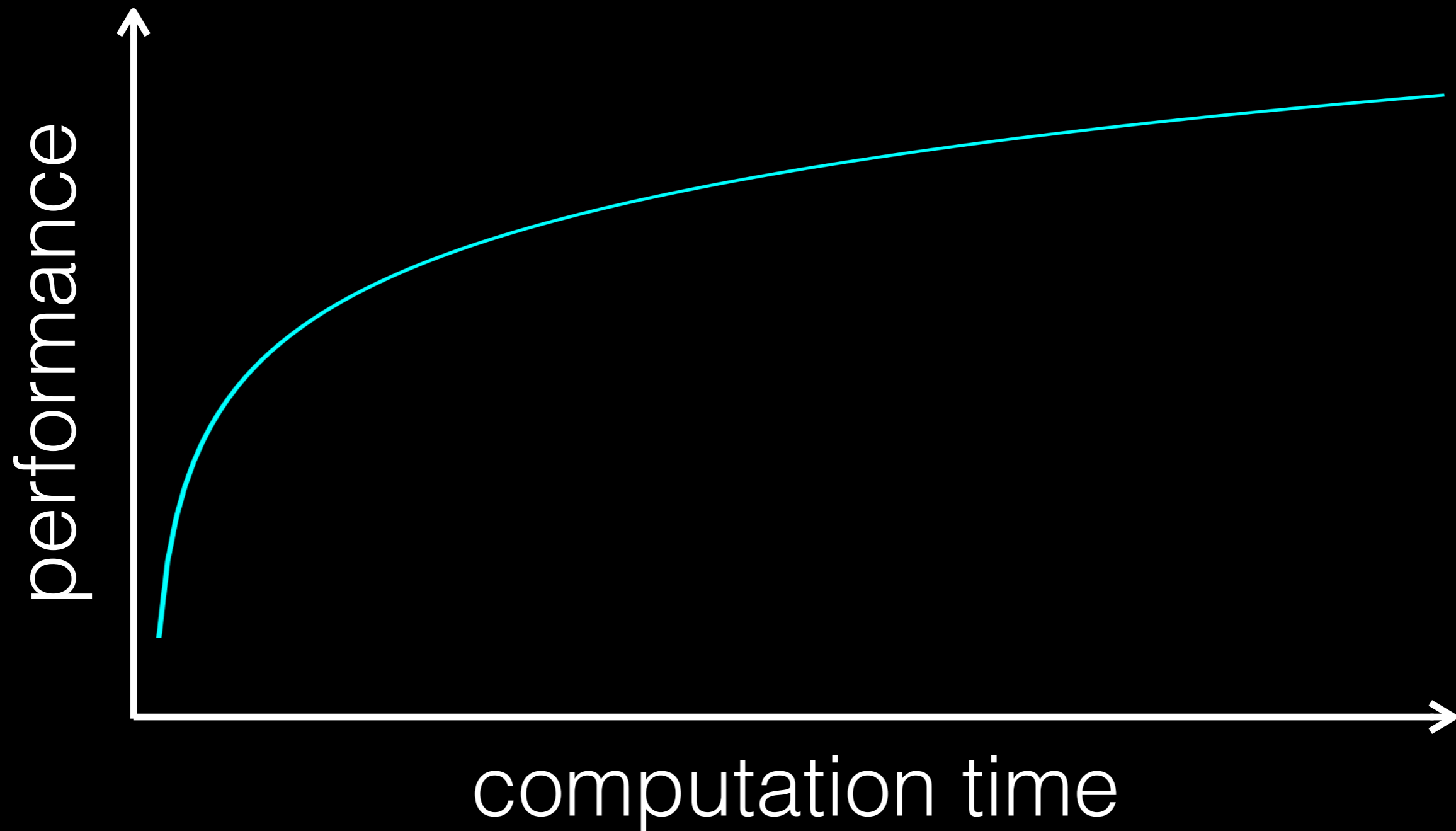
## Eric Nalisnick

University of Amsterdam
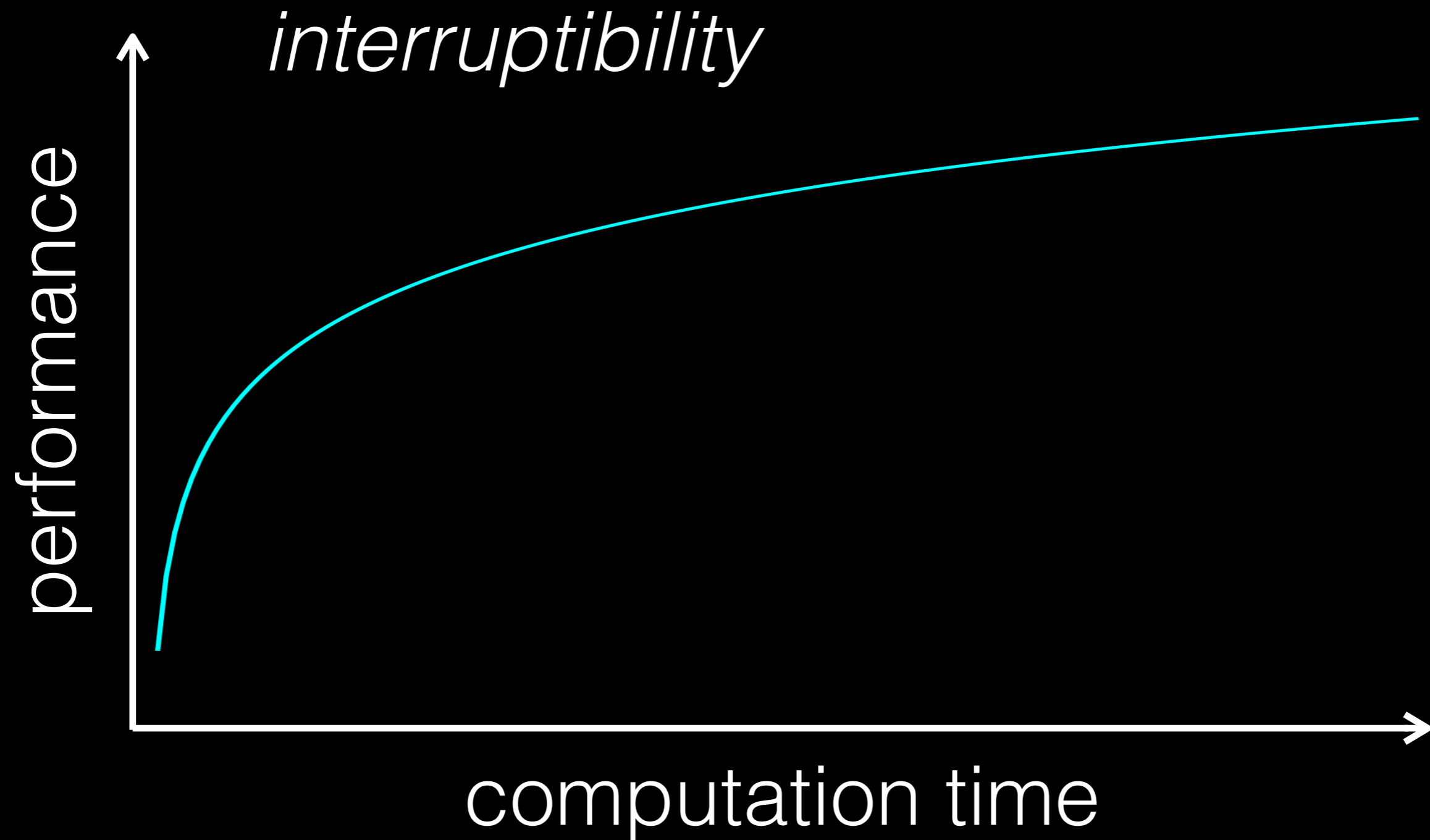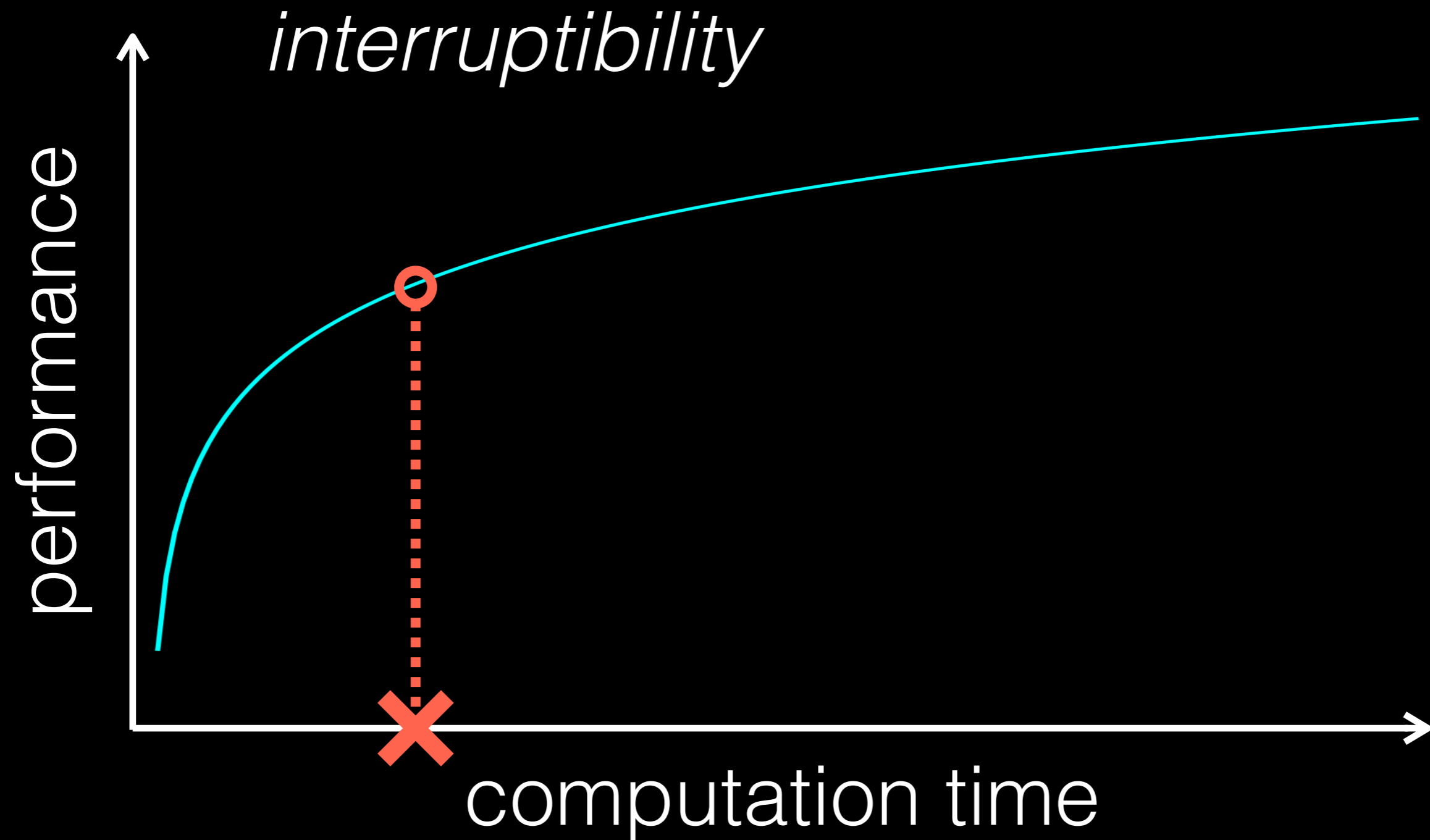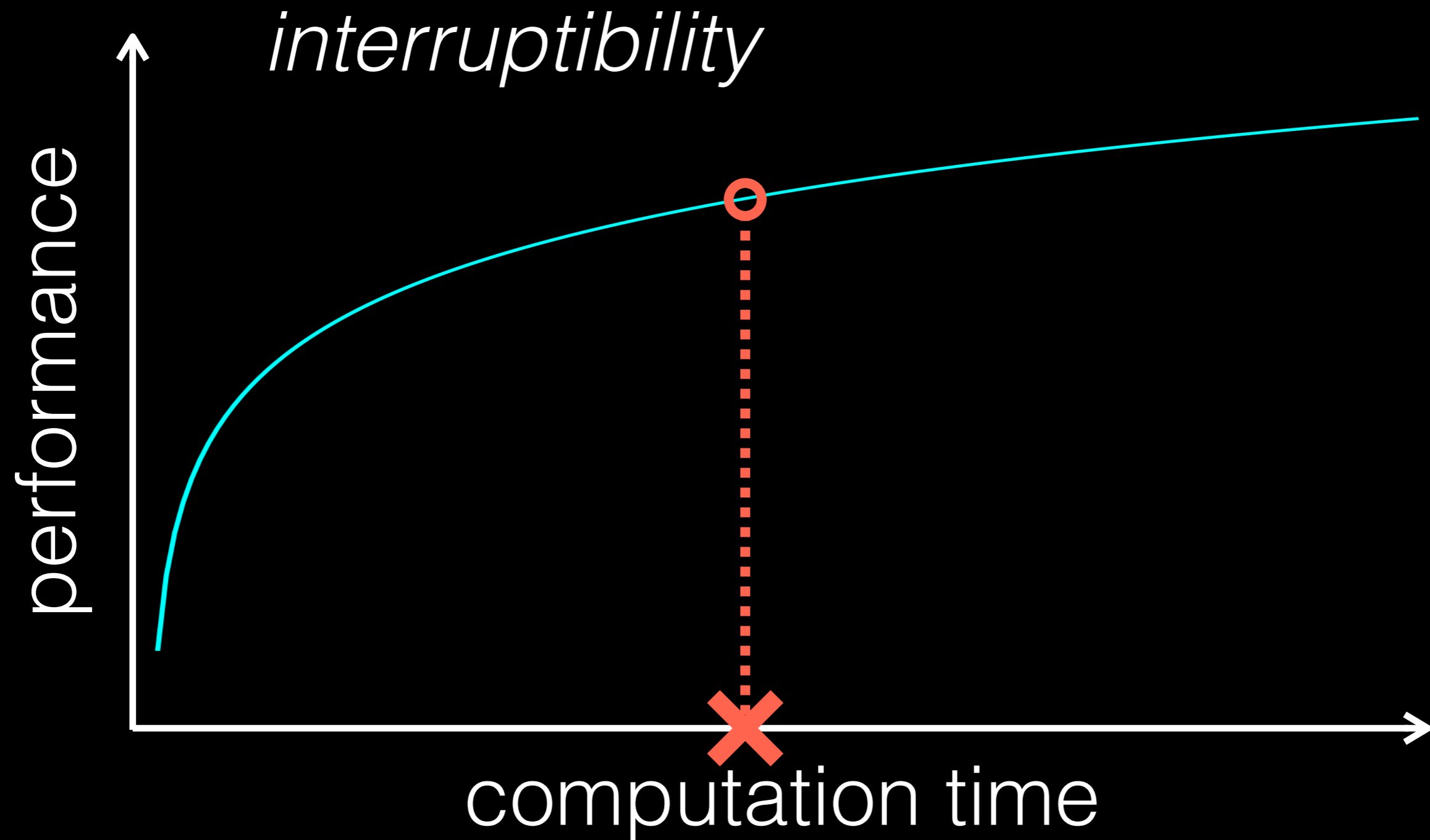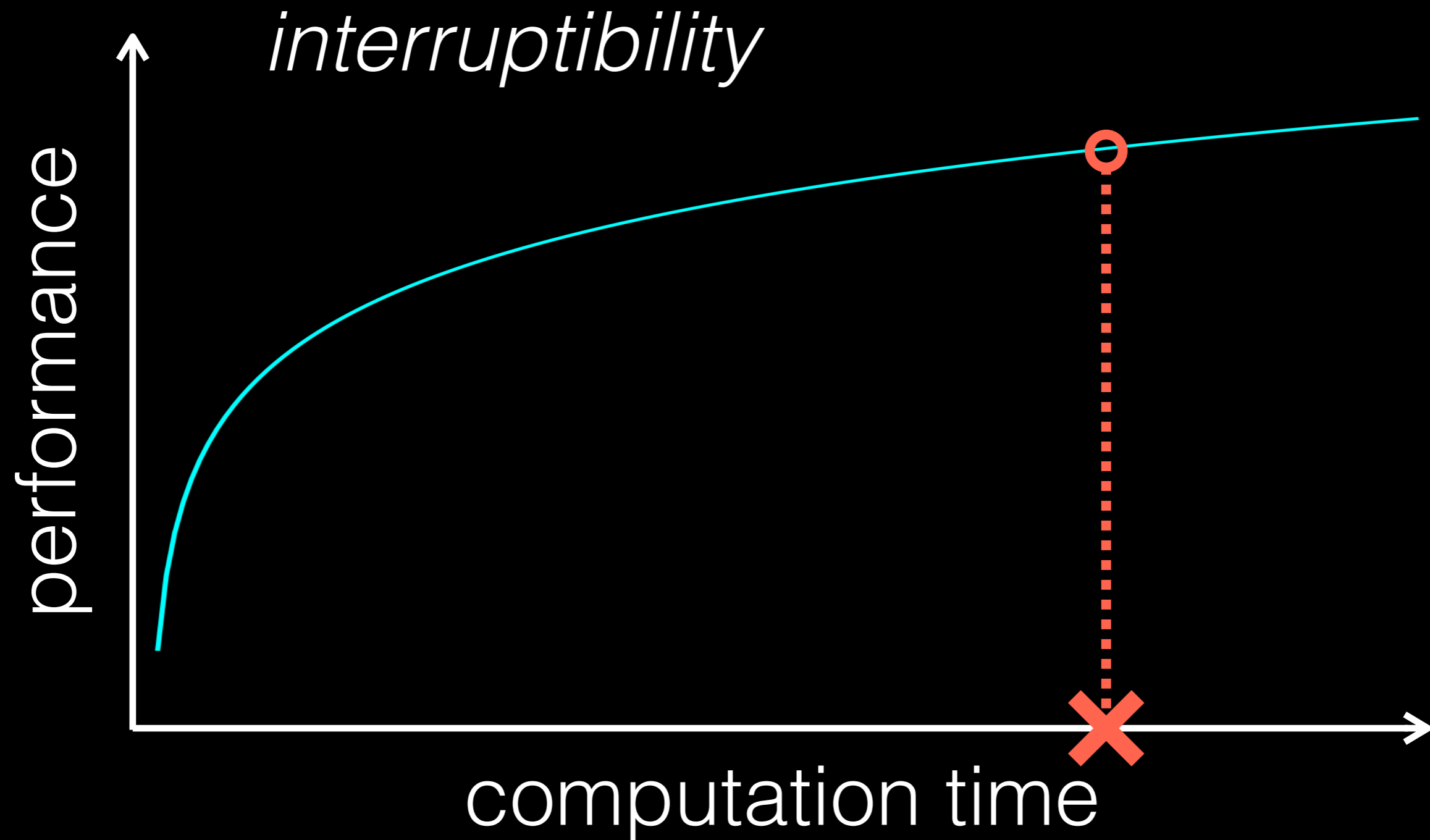
# Anytime Models

# Anytime Models

# Anytime Models

# Anytime Models

*interruptibility*

performance

computation time

# Anytime Models

# Anytime Models

# Early-Exit Neural Networks

$$\ell\left(\theta_{1:E}\right) = -\sum_{e=1}^{E} \log p_e(y \mid x, \theta_{1:e})$$

# Early-Exit Neural Networks

- ⊗ interruptibility?

- ⊗ monotonicity?

- ⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

⊗ interruptibility?

⊗ monotonicity?

⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

⊗ interruptibility ✅

⊗ monotonicity?

⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

- ⊗ interruptibility ✅

- ⊗ monotonicity?

- ⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Multi-Scale Dense Net: CIFAR-100

# Multi-Scale Dense Net: CIFAR-100

Anytime Models

*marginal* *monotonicity*

average performance

computation time

# Early-Exit Neural Networks

⊗ interruptibility ✅

⊗ monotonicity? ✅

⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Multi-Scale Dense Net: CIFAR-100

Multi-Scale Dense Net: CIFAR-100

Multi-Scale Dense Net: CIFAR-100

# Multi-Scale Dense Net: CIFAR-100

# Multi-Scale Dense Net: CIFAR-100

# Multi-Scale Dense Net: Overthinking

**Overthinking**: having the correct prediction but then switching to a wrong prediction.

[Kaya et al., ICML 2019]

$$\Delta = (\text{test error at final exit}) -$$
$$(\text{test error if exited at correct prediction})$$

$$\Delta(\text{CIFAR} - 100) = \sim 14\,\%$$

$$\Delta(\text{ImageNet}) = \sim 9\,\%$$

# Early-Exit Neural Networks

⊗ interruptibility ✅

⊗ monotonicity ✅

⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

⊗ interruptibility ✅

⊗ monotonicity ✅

⊗ diminishing returns?

[Zilberstein, AI Magazine 1996]

Multi-Scale Dense Net: CIFAR-100

# Early-Exit Neural Networks

⊗ interruptibility ✅

⊗ monotonicity ✅

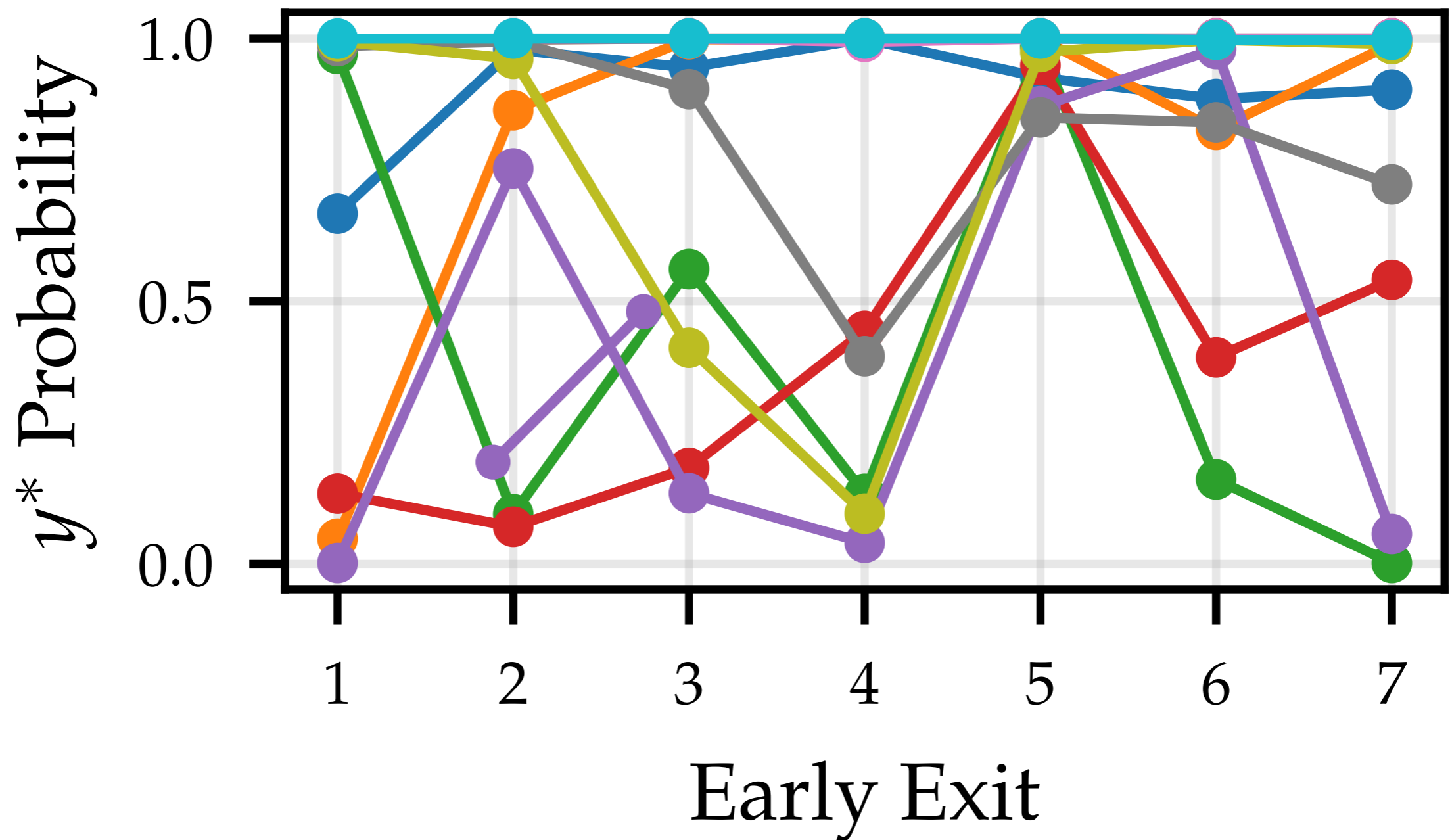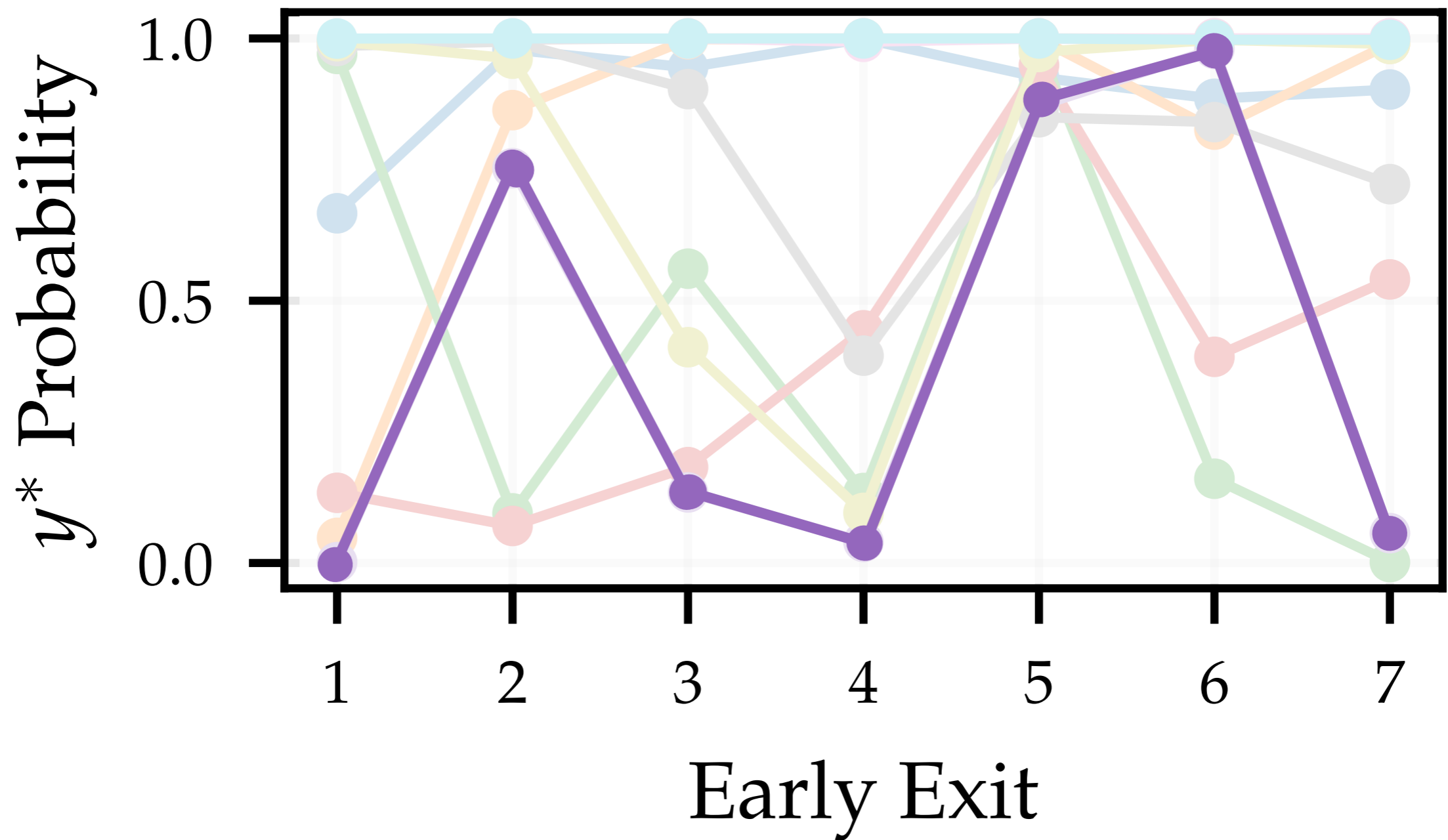⊗ diminishing returns ✅

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

⊗  interruptibility  ✅

⊗  monotonicity  ✅

⊗  diminishing returns  ✅

only marginally

[Zilberstein, AI Magazine 1996]

# Early-Exit Neural Networks

⊗ interruptibility ✔

⊗ monotonicity ✔

⊗ diminishing returns ✔

[Zilberstein, AI Magazine 1996]

# A simple, post-hoc method for encouraging conditional monotonicity

Metod Jazbec

James U. Allingham

Dan Zhang

NEURAL INFORMATION PROCESSING SYSTEMS

2023

# Idea: combine the early-exits via a product of experts



experts

$p_1(y|x)$    $p_2(y|x)$    $p_3(y|x)$

$x \rightarrow h_1 \rightarrow h_2 \rightarrow h_3$

# Idea: combine the early-exits via a product of experts

$p_{1:2}(y \mid x)$

$p_1(y \mid x) \quad p_2(y \mid x) \quad p_3(y \mid x)$

$x \rightarrow h_1 \rightarrow h_2 \rightarrow h_3$

# Idea: combine the early-exits via a product of experts

$$p_{1:2}(y \mid x) = \frac{p_1(y \mid x) \cdot p_2(y \mid x)}{\sum_{y'} p_1(y' \mid x) \cdot p_2(y' \mid x)}$$

$p_{1:2}(y \mid x)$

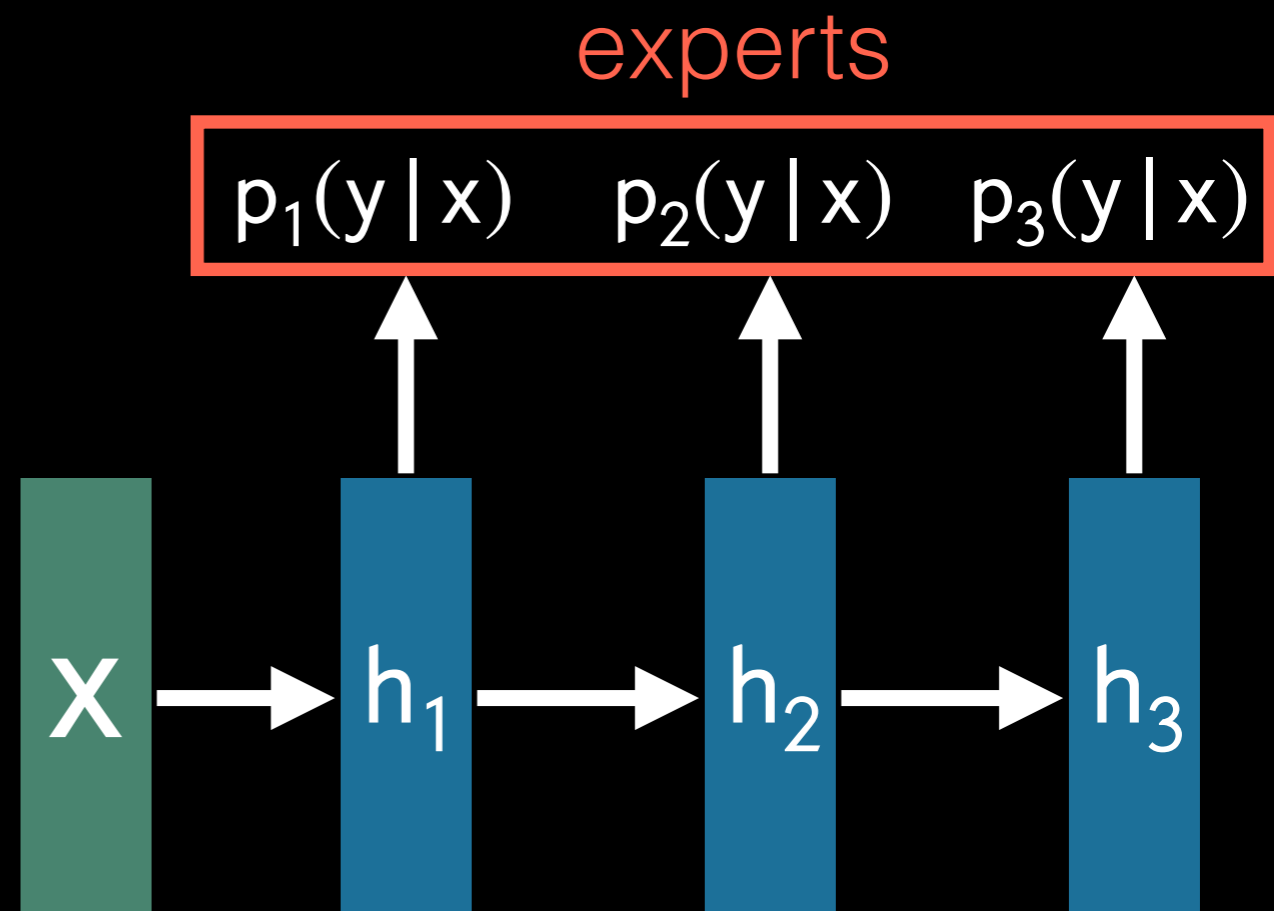$p_1(y \mid x) \quad p_2(y \mid x) \quad p_3(y \mid x)$

$x \rightarrow h_1 \rightarrow h_2 \rightarrow h_3$

# Idea: combine the early-exits via a product of experts

# Idea: combine the early-exits via a product of experts

$$p_{1:3}(y \mid x)$$
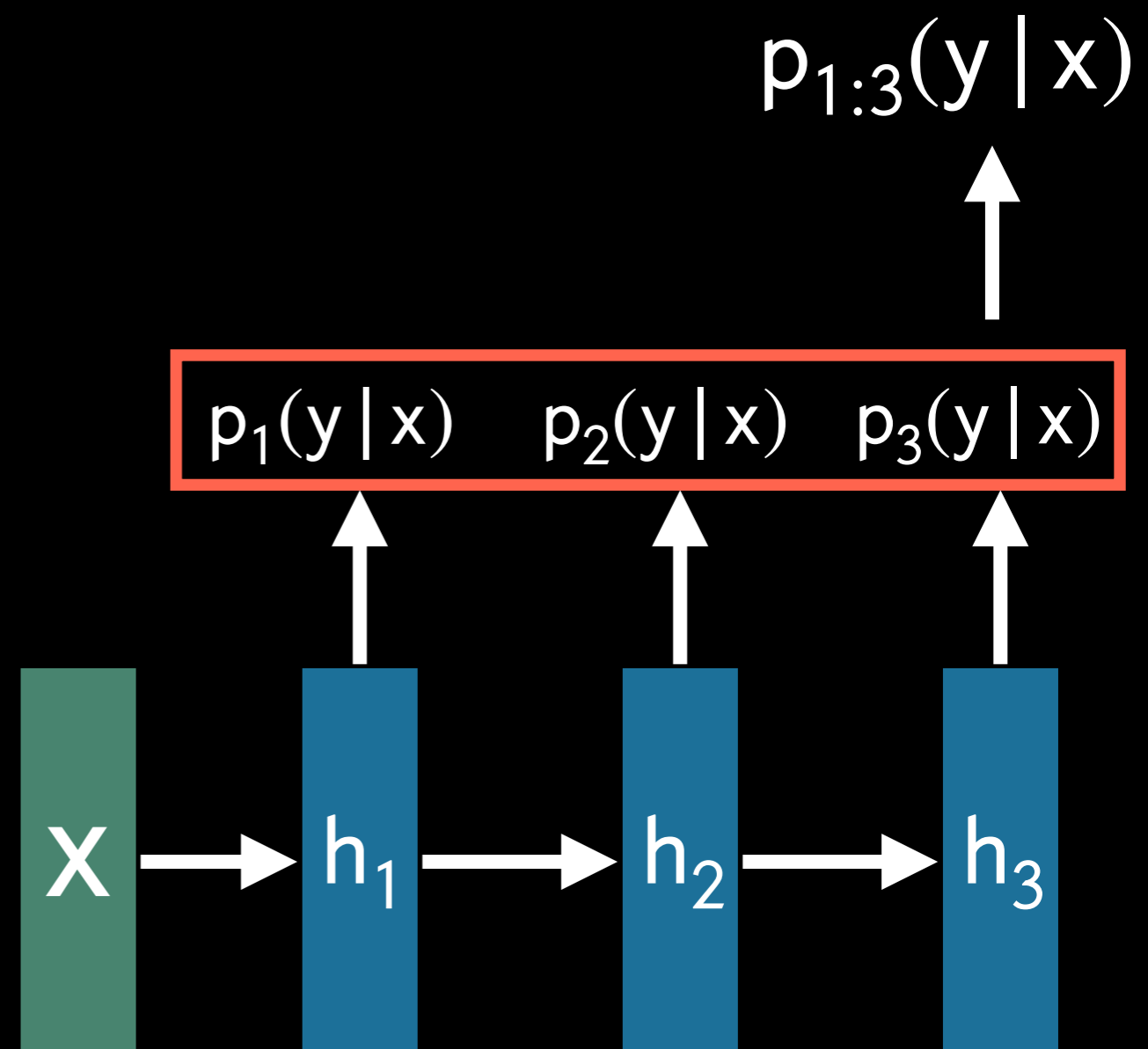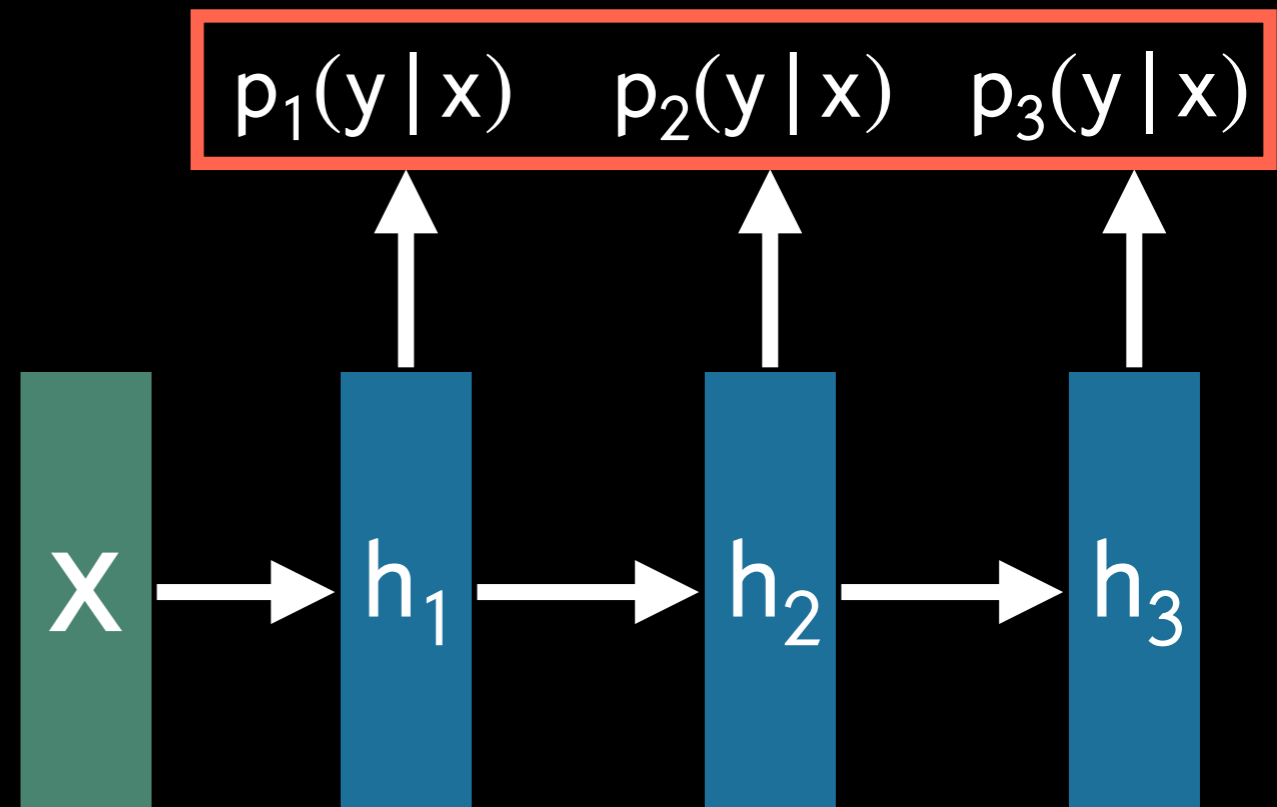
$$p_1(y \mid x) \quad p_2(y \mid x) \quad p_3(y \mid x)$$
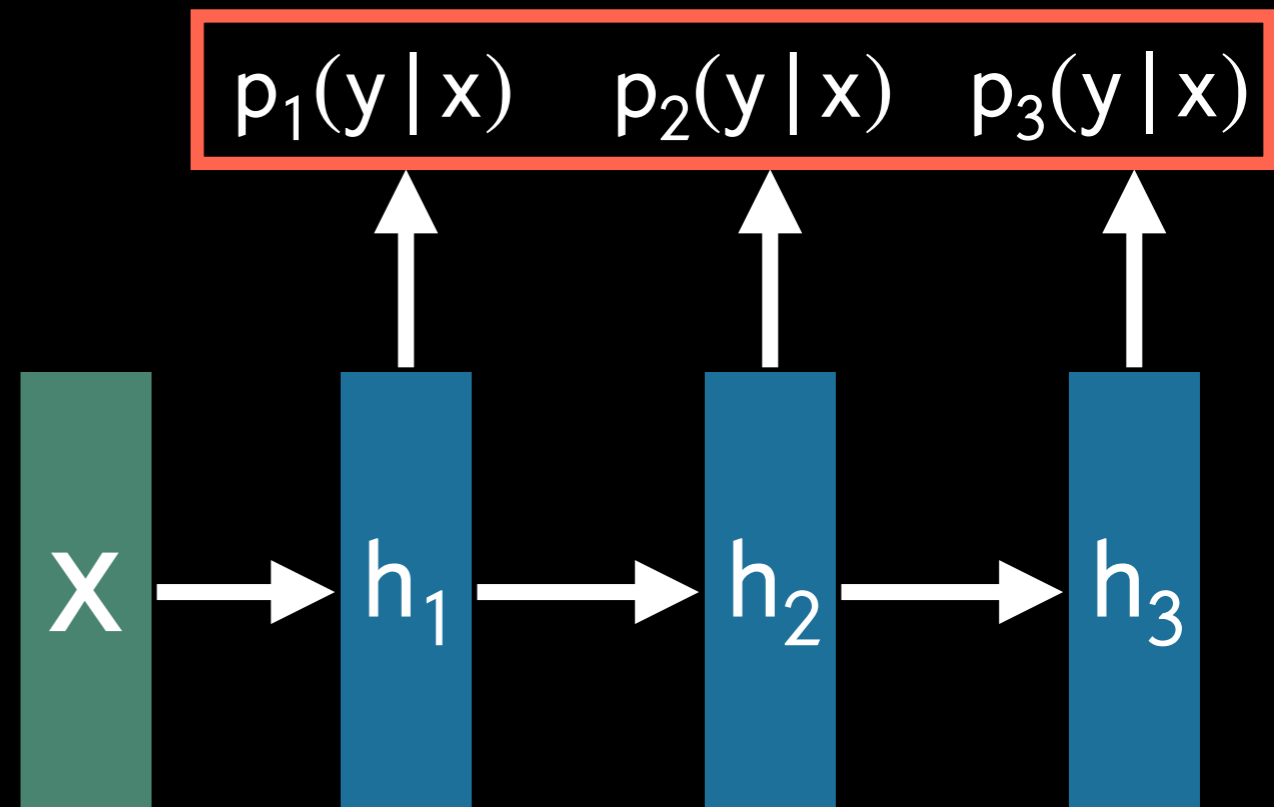
x → $h_1$ → $h_2$ → $h_3$

# Idea: combine the early-exits via a product of experts

$$p_{1:3}(y \mid x) = \frac{p_1(y \mid x) \cdot p_2(y \mid x) \cdot p_3(y \mid x)}{\sum_{y'} p_1(y' \mid x) \cdot p_2(y' \mid x) \cdot p_3(y' \mid x)}$$

$p_{1:3}(y \mid x)$

# Idea: combine the early-exits via a product of experts

$$p_{1:e}(y \mid x) = \frac{\prod_{j=1}^{e} p_j(y \mid x)}{\sum_{y'} \prod_{j=1}^{e} p_j(y' \mid x)}$$
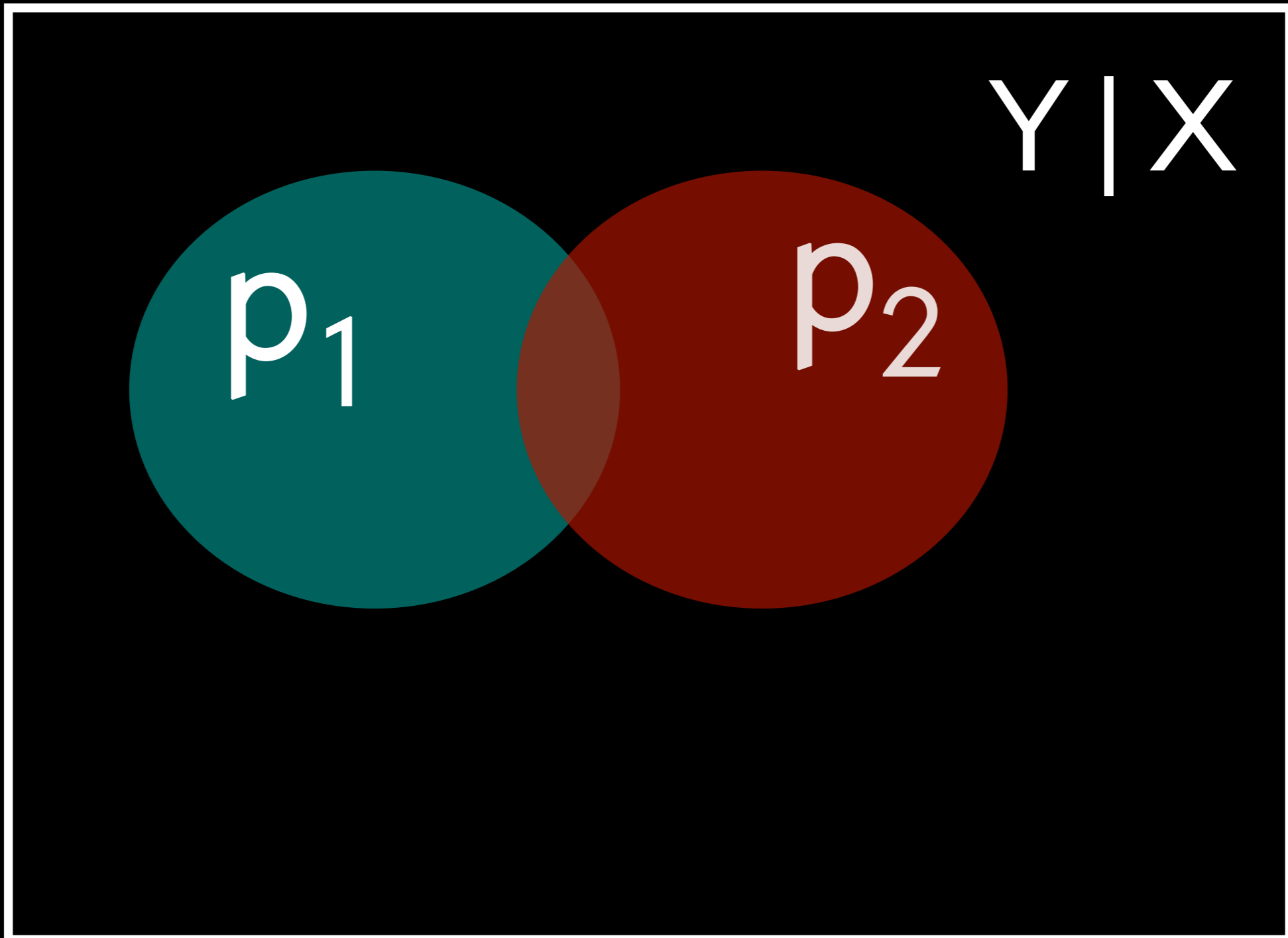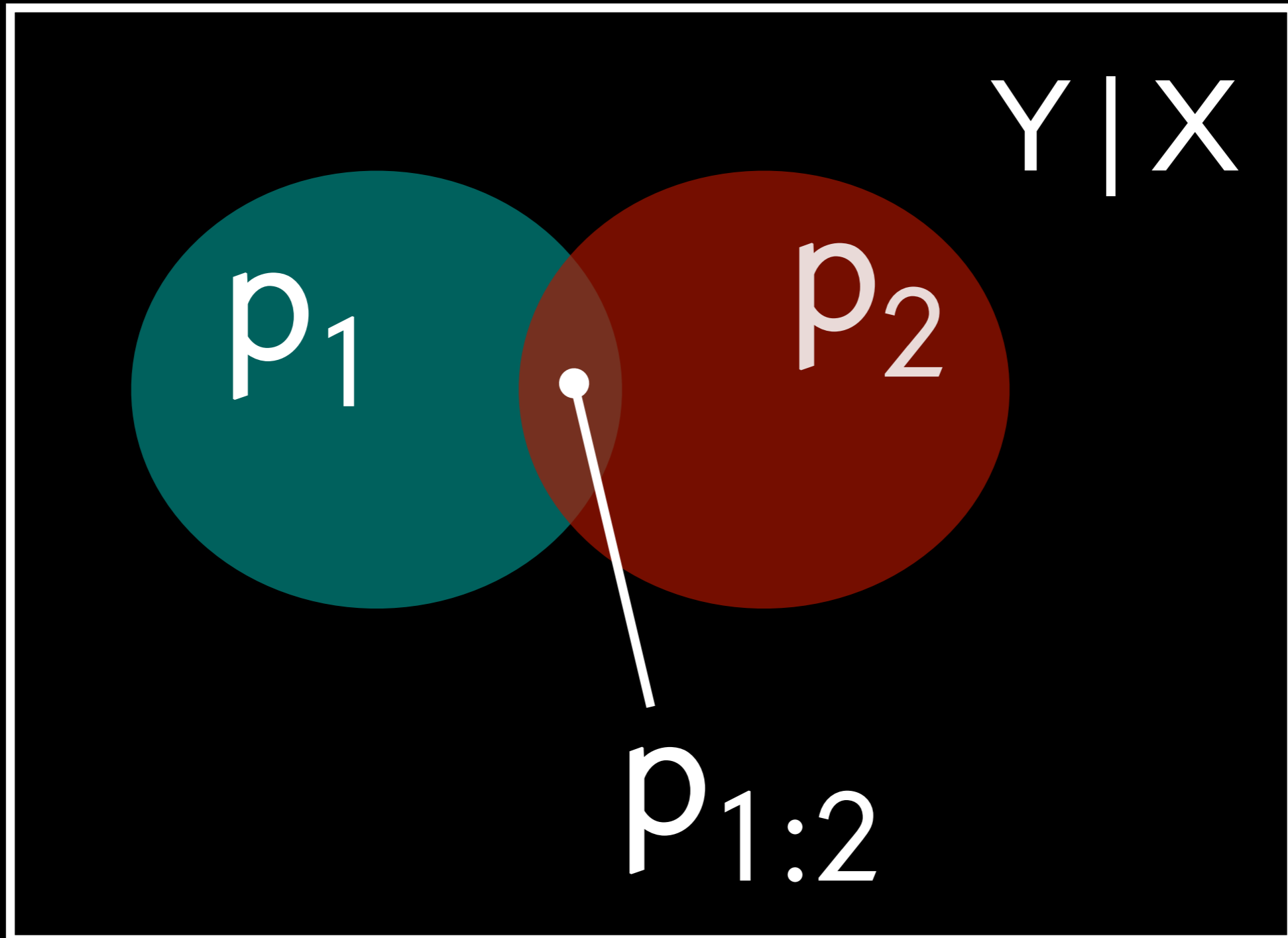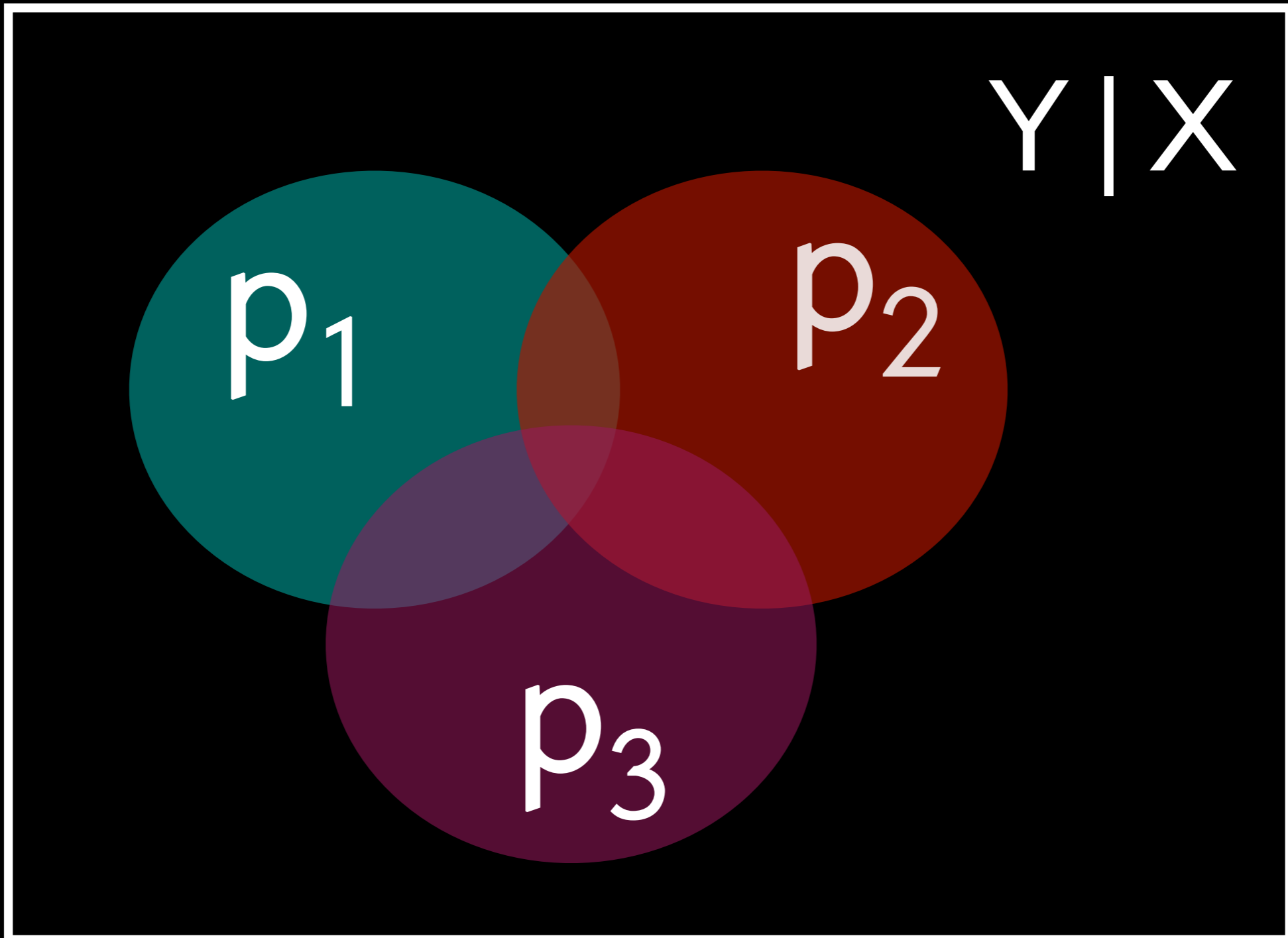
# Idea: combine the early-exits via a product of experts

$$Y \mid X$$

Idea: combine the early-exits via a product of experts
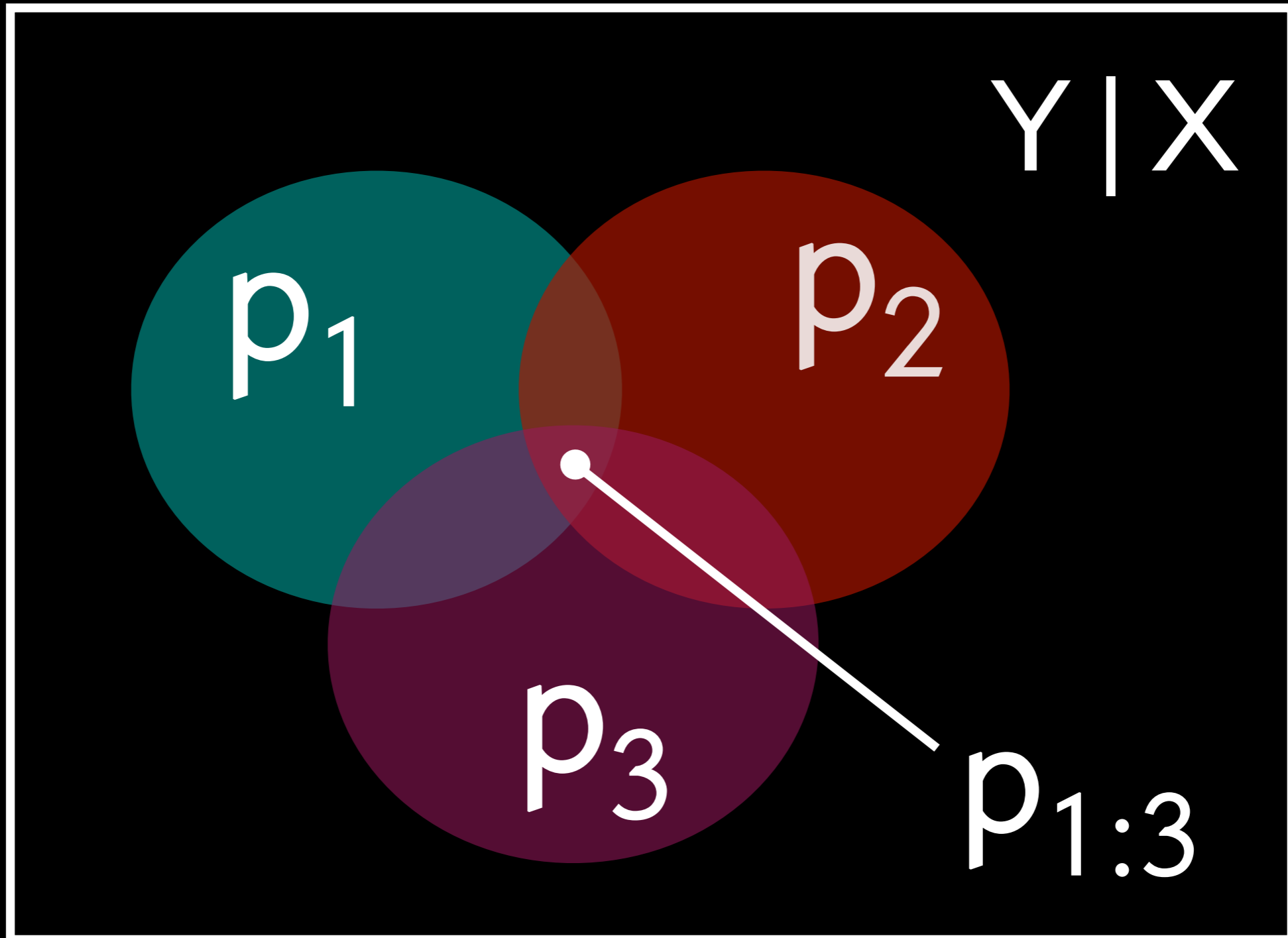
Idea: combine the early-exits via a product of experts

$p_1$ $p_2$ $p_{1:2}$ $Y|X$

Idea: combine the early-exits via a product of experts

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.
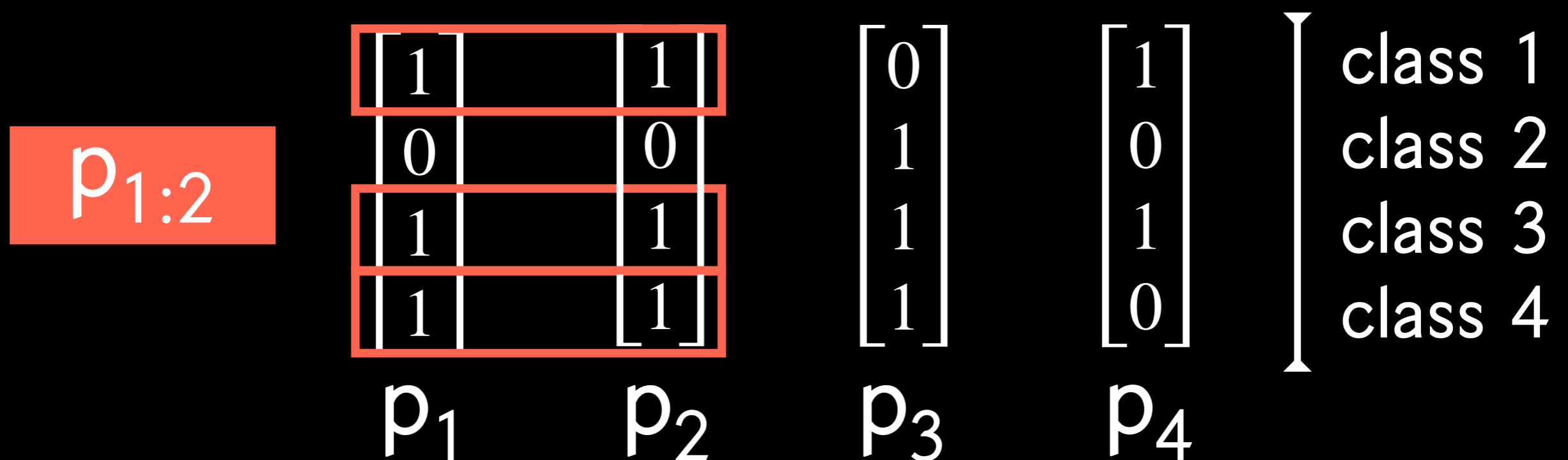
**Ideal case**: binary one-vs-rest.

$$
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{matrix} \text{class 1} \\ \text{class 2} \\ \text{class 3} \\ \text{class 4} \end{matrix}
$$

$p_1 \qquad p_2 \qquad p_3 \qquad p_4$

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.
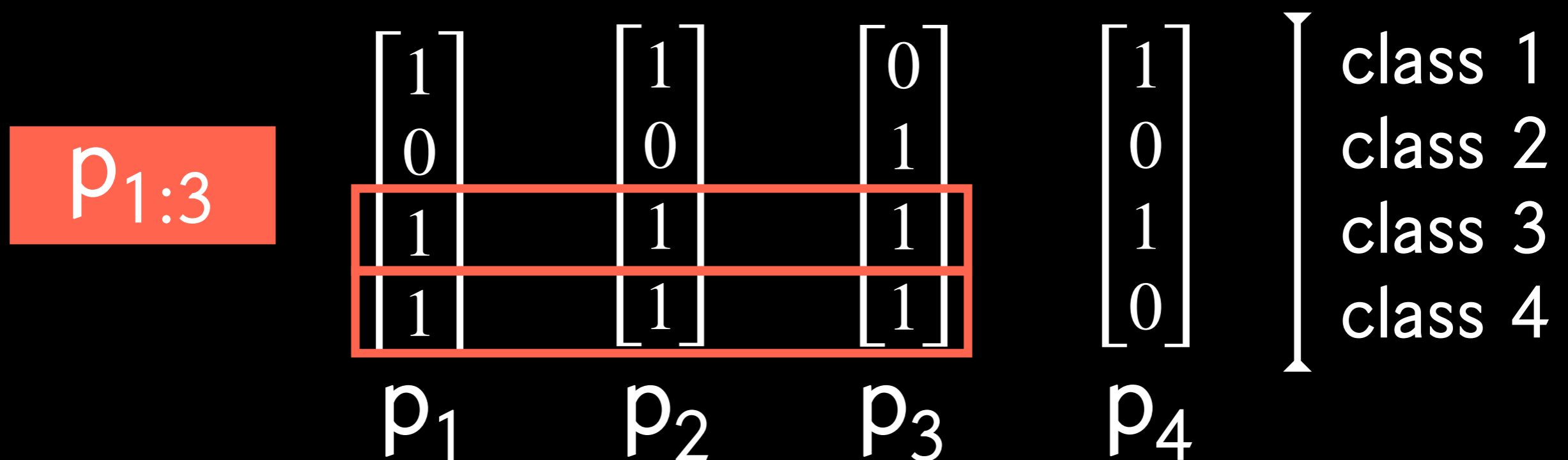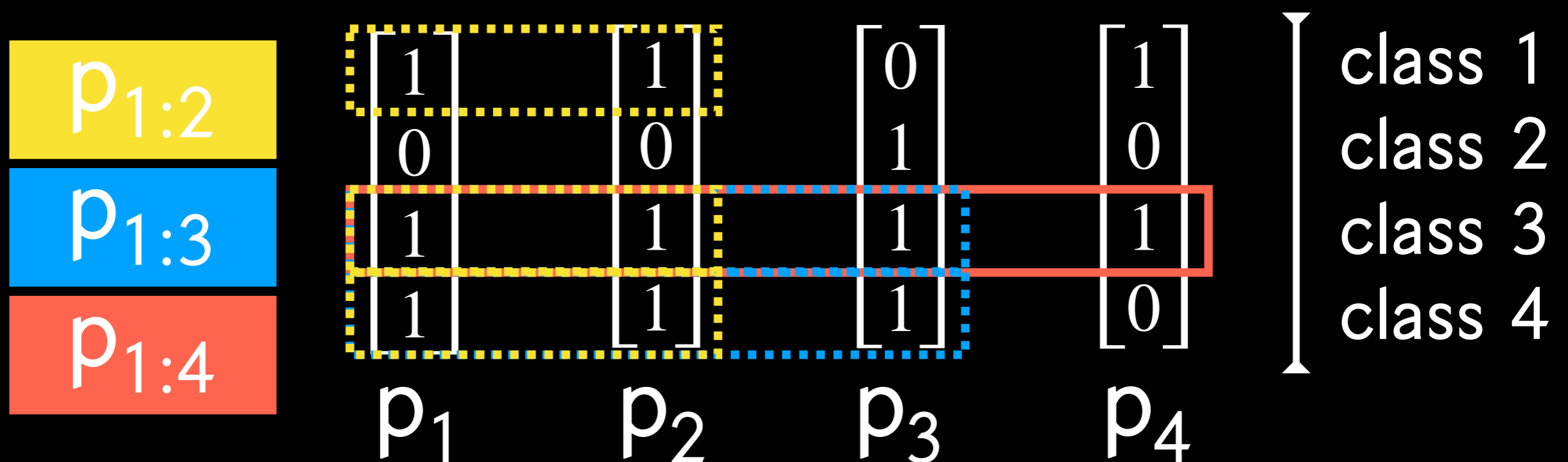
**Ideal case**: binary one-vs-rest.

$$
p_{1:2} \quad
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}
\quad
\begin{matrix} \text{class 1} \\ \text{class 2} \\ \text{class 3} \\ \text{class 4} \end{matrix}
$$

$$
p_1 \qquad p_2 \qquad p_3 \qquad p_4
$$

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.

**Ideal case**: binary one-vs-rest.

$$p_{1:3} \qquad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \qquad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} \qquad \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \qquad \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \qquad \begin{matrix} \text{class 1} \\ \text{class 2} \\ \text{class 3} \\ \text{class 4} \end{matrix}$$

$$p_1 \qquad p_2 \qquad p_3 \qquad p_4$$

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.

**Ideal case**: binary one-vs-rest.

$$
p_{1:4} \quad
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}
\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}
\quad
\begin{array}{l} \text{class 1} \\ \text{class 2} \\ \text{class 3} \\ \text{class 4} \end{array}
$$

$$p_1 \qquad p_2 \qquad p_3 \qquad p_4$$

# Idea: combine the early-exits via a product of experts

**One catch**: exit distributions must have finite (or quickly decaying) support to bound influence of (e+1)th expert.

**Ideal case**: binary one-vs-rest.

# Implementation with ReLUs

$$p_{1:e}(y \mid x) = \frac{\prod_{j=1}^{e} \max\left(0,\, f_{j,y}(x)\right)}{\sum_{y'} \prod_{j=1}^{e} \max\left(0,\, f_{j,y'}(x)\right)}$$

$f_{j,y}(x)$ is logit for yth class at jth exit

# Implementation with ReLUs

$$p_{1:e}(y \mid x) = \frac{\prod_{j=1}^{e} \max \left( 0, \ f_{j,y}(x) \right)}{\sum_{y'} \prod_{j=1}^{e} \max \left( 0, \ f_{j,y'}(x) \right)}$$

Clipping logits controls deviation from perfect monotonicity.

We apply this transformation post-hoc!

# Monotonicity: CIFAR-100

# Monotonicity: CIFAR-100

# Accuracy: CIFAR-100 & ImageNet



CIFAR-100

ImageNet

Test Accuracy (↑)

Early Exit

Early Exit

baseline
PA (ours)
CA

MSDNet
IMTA

# Accuracy: CIFAR-100 & ImageNet

# Overthinking: CIFAR-100 & ImageNet

# Overthinking: CIFAR-100 & ImageNet



*Doesn't mean that overall accuracy is improved by this amount since our model makes more mistakes at intermediate exits.

# Ensuring consistency across exits in predictive uncertainty estimates



Metod Jazbec

Dan Zhang

Patrick Forré

Stephan Mandt

# Anytime Models

performance

computation time

Anytime Uncertainty

# Anytime Uncertainty Estimation

We want nested, non-increasing prediction intervals across exits.

consistency:
$$C_1(x) \subseteq C_2(x) \subseteq C_3(x)$$

# Anytime-Valid Confidence Sequences

We construct an *anytime-valid confidence sequence* across the exits.

$$\mathbb{P}\left(\forall t,\ y^* \in C_t(x)\right) \geq 1 - \alpha$$

*Due to approximations, we can only hope to achieve this for large datasets (and if $y^*$ is from the training distribution).

[Robbins, AMS 1970]

# Anytime-Valid Confidence Sequences

Derived from the following predictive-likelihood martingale:

$$R_t(y) = \prod_{e=1}^{t} \frac{p_e(y \mid x, \mathfrak{D})}{p_e(y \mid x, \hat{\theta}_e)} \qquad \hat{\theta}_e \sim p\left(\theta_e \mid x, \mathfrak{D}\right)$$

# Anytime-Valid Confidence Sequences

Derived from the following predictive-likelihood martingale:

$$R_t(y) = \prod_{e=1}^{t} \frac{p_e(y \mid x, \mathfrak{D})}{p_e(y \mid x, \hat{\theta}_e)} \qquad \hat{\theta}_e \sim p\left(\theta_e \mid x, \mathfrak{D}\right)$$

Construct set at time t as:

$$C_t(x) = \left\{ y \in Y \mid R_t(y) \leq 1/\alpha \right\}$$

# Regression Simulation



t = 1

t = 5

t = 15

x

Regression Simulation

y

t = 1

y

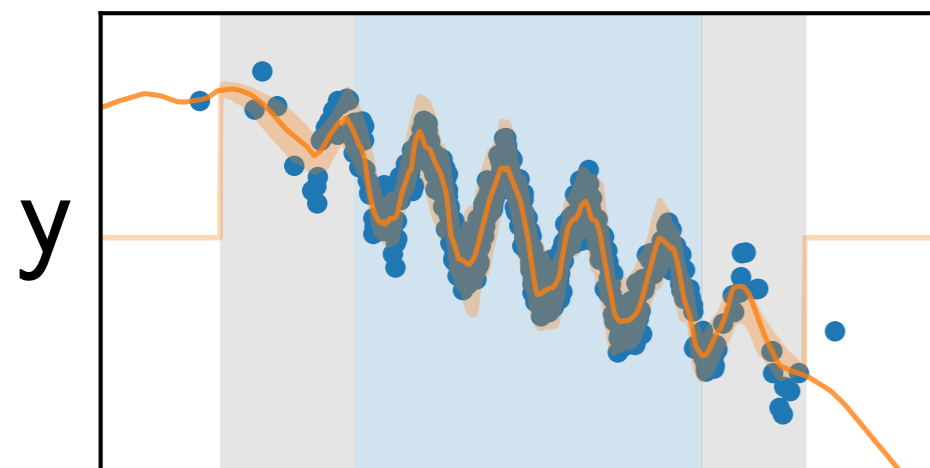t = 5

y

t = 15

Intersection
Current
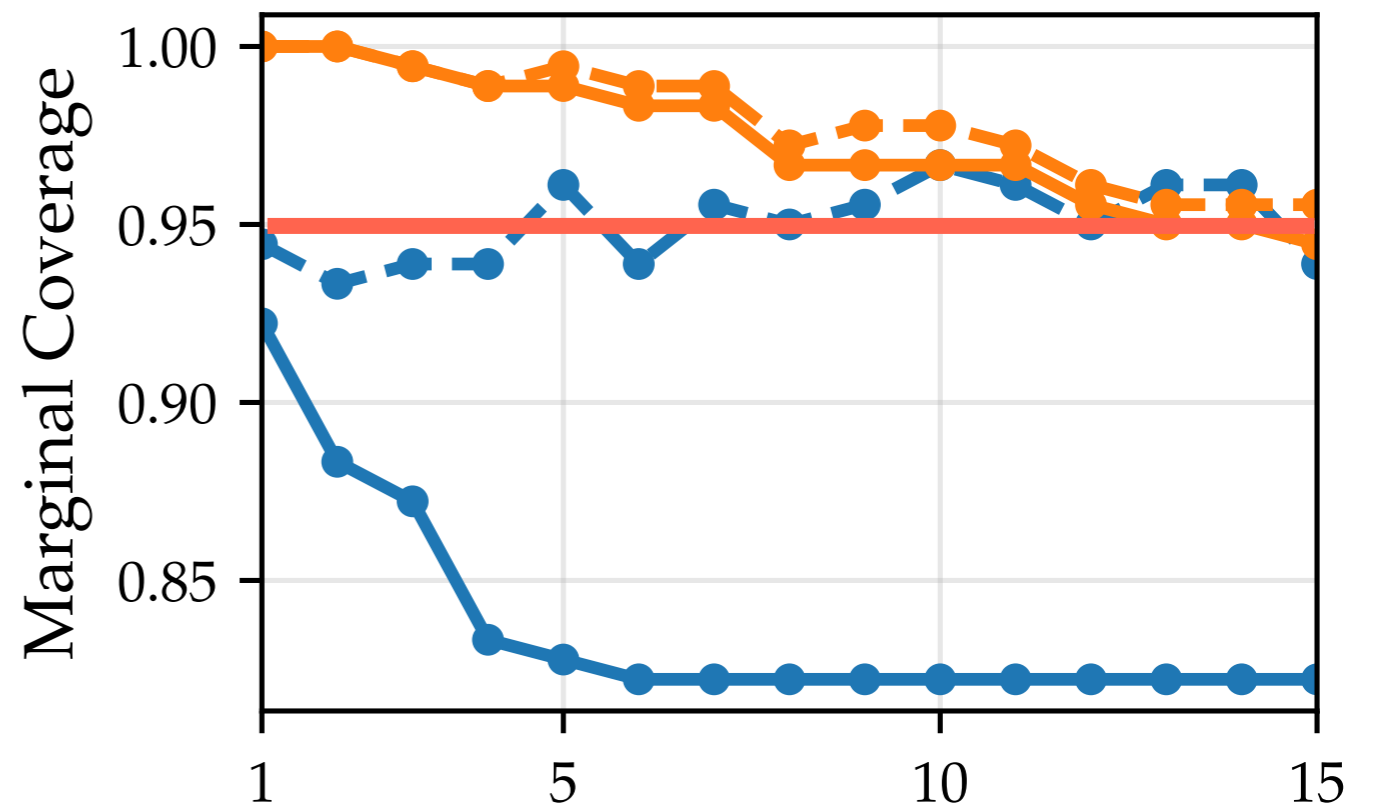
EENN-Bayes
EENN-AVCS

X

102

Regression Simulation

Regression Simulation

Regression Simulation

Regression Simulation

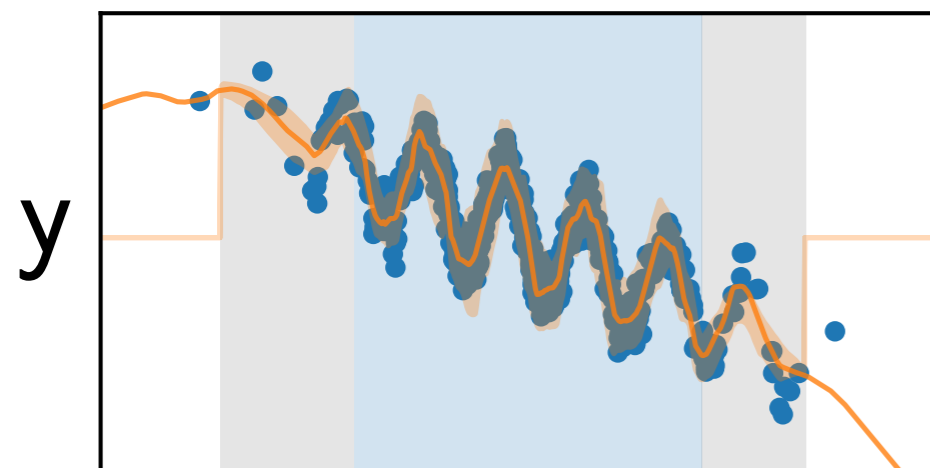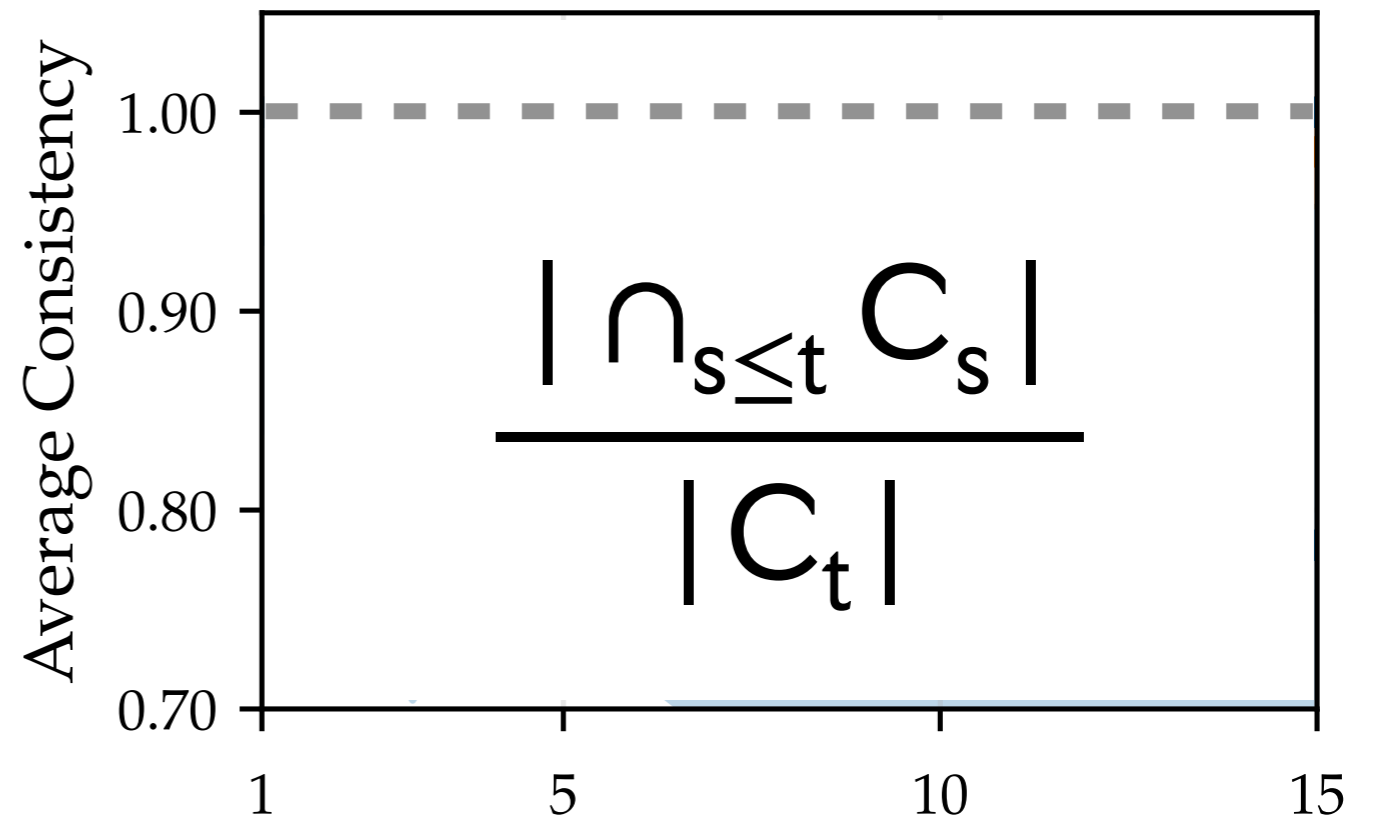Regression Simulation

$t = 1$

$t = 5$

$t = 15$

y

y

y

X

Average Consistency

$$\frac{|\cap_{s \leq t} C_s|}{|C_t|}$$

1.00

0.90

0.80

0.70

1     5     10     15

Time / Early-Exit

Intersection

Current

EENN-Bayes

EENN-AVCS

107
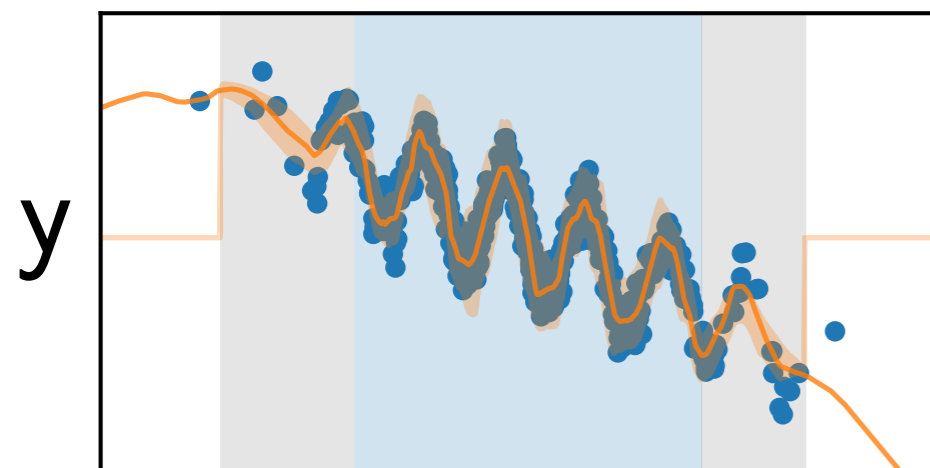
Regression Simulation
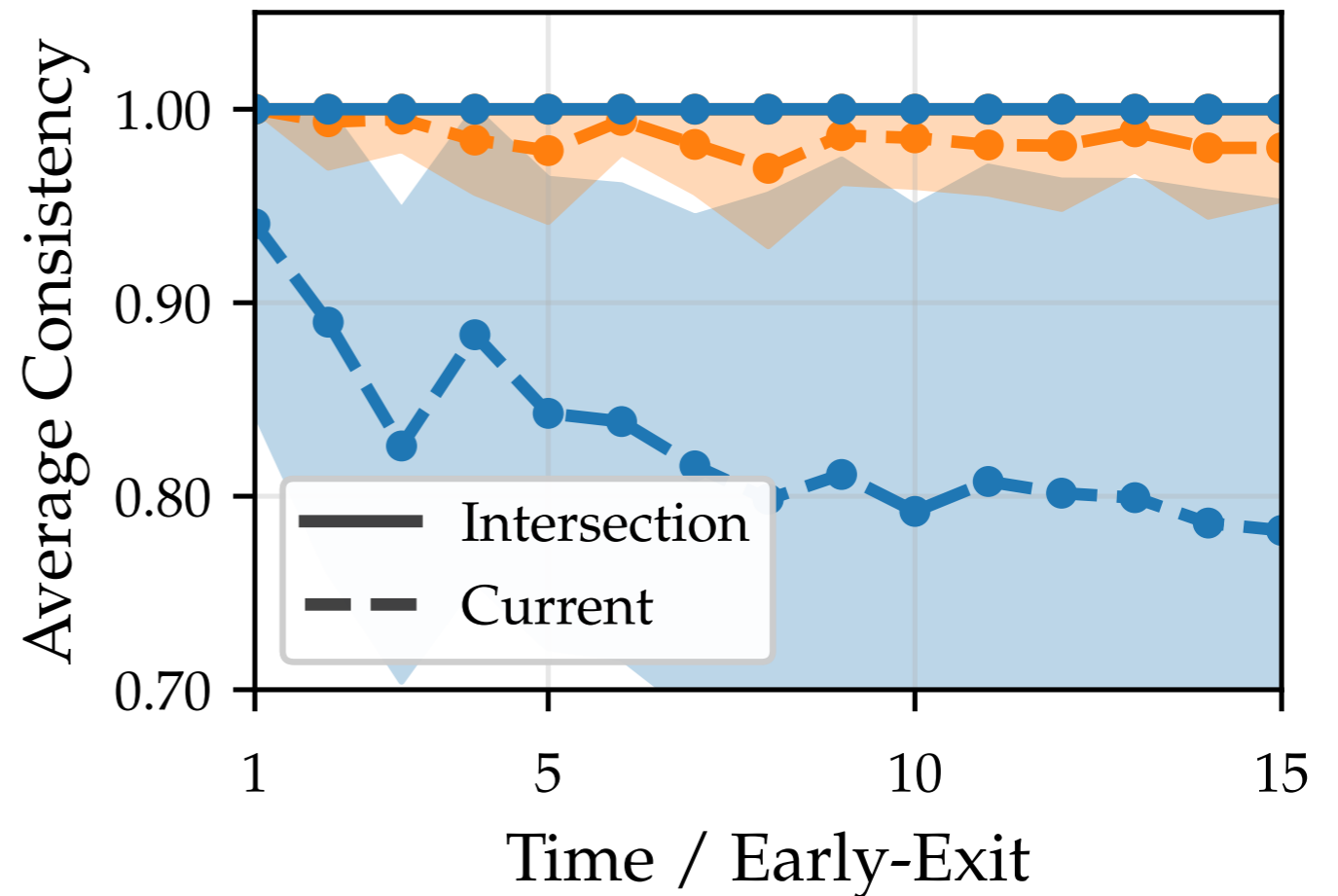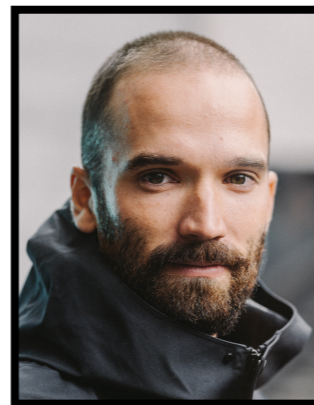
# Summary

⊗ Early-exit neural networks have mostly marginal anytime properties (and overthink)

⊗ We give them better conditional monotonicity via a product ensemble.

---

⊗ Also want consistency in predictive uncertainty across exits.

⊗ We enforce this with anytime-valid confidence sequences.

# Thank you!  Questions?

paper





Metod
Jazbec



James U.
Allingham

UvA - BOSCH
DELTA LAB



Dan
Zhang



Patrick
Forré



Stephan
Mandt