



Calibration in Deep Learning: Theory & Practice

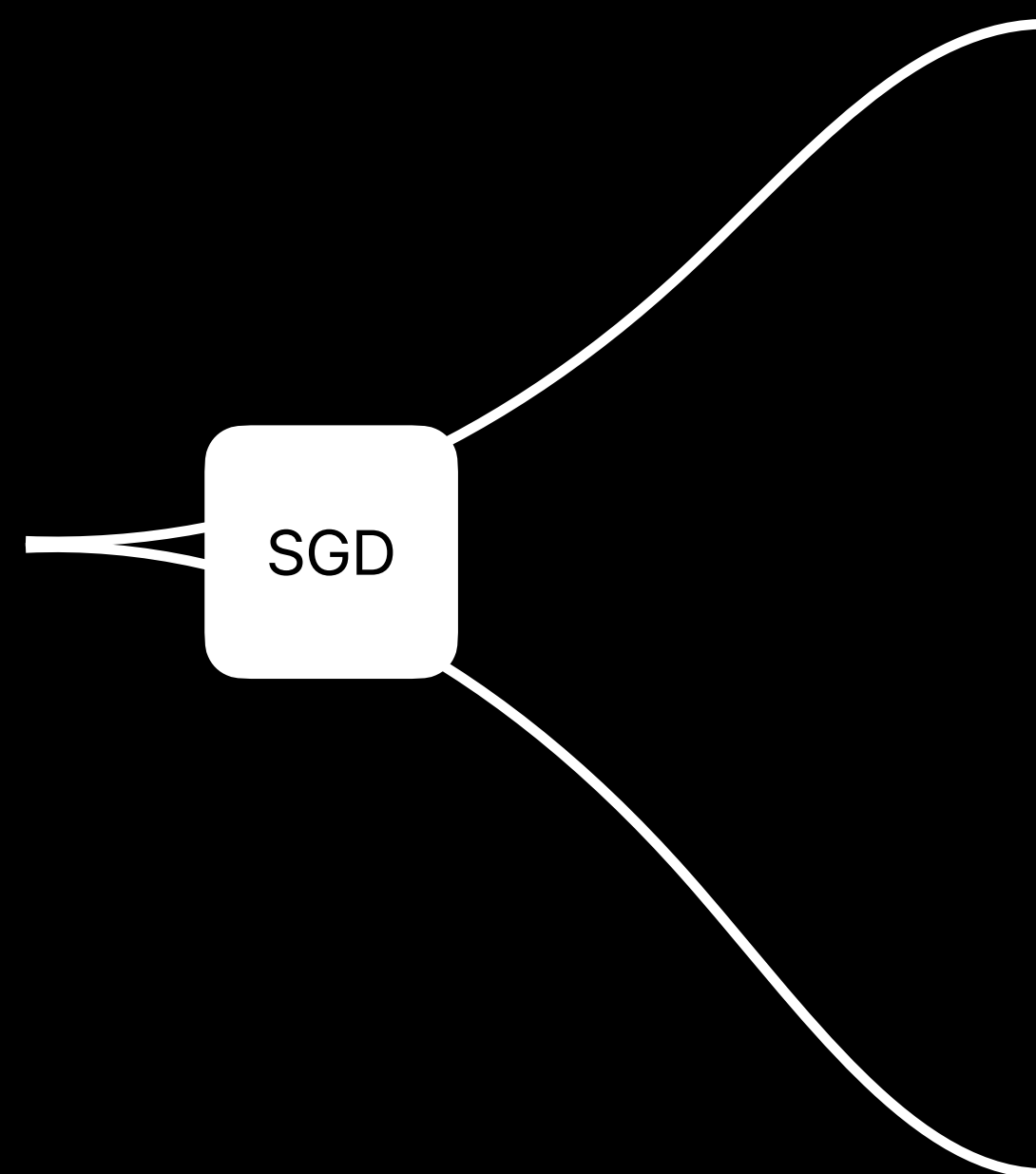
Preetum Nakkiran

Aspen Center for Physics | Apple Inc. | Feb 28 2023

Motivation

*“Why do we **get** more than we **asked for** in Deep Learning?”*

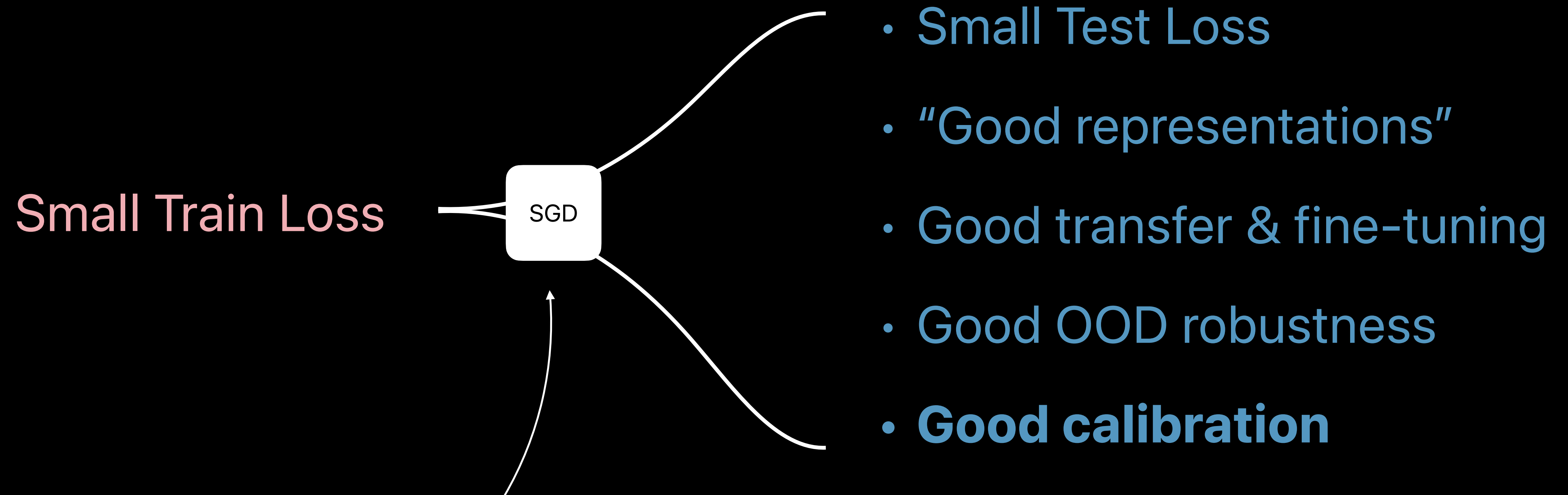
Small Train Loss



- Small Test Loss
- “Good representations”
- Good transfer & fine-tuning
- Good OOD robustness
- **Good calibration**

Motivation

*“Why do we **get** more than we **asked for** in Deep Learning?”*



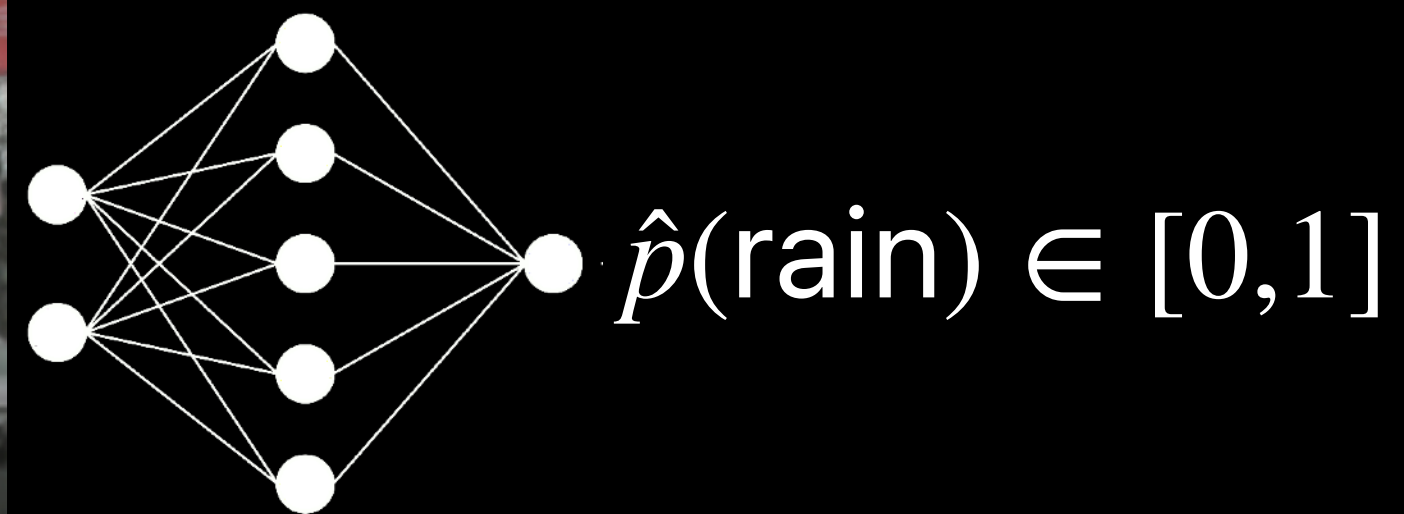
Goal: What’s important about this box?

What is Calibration?

Setting: Binary classification

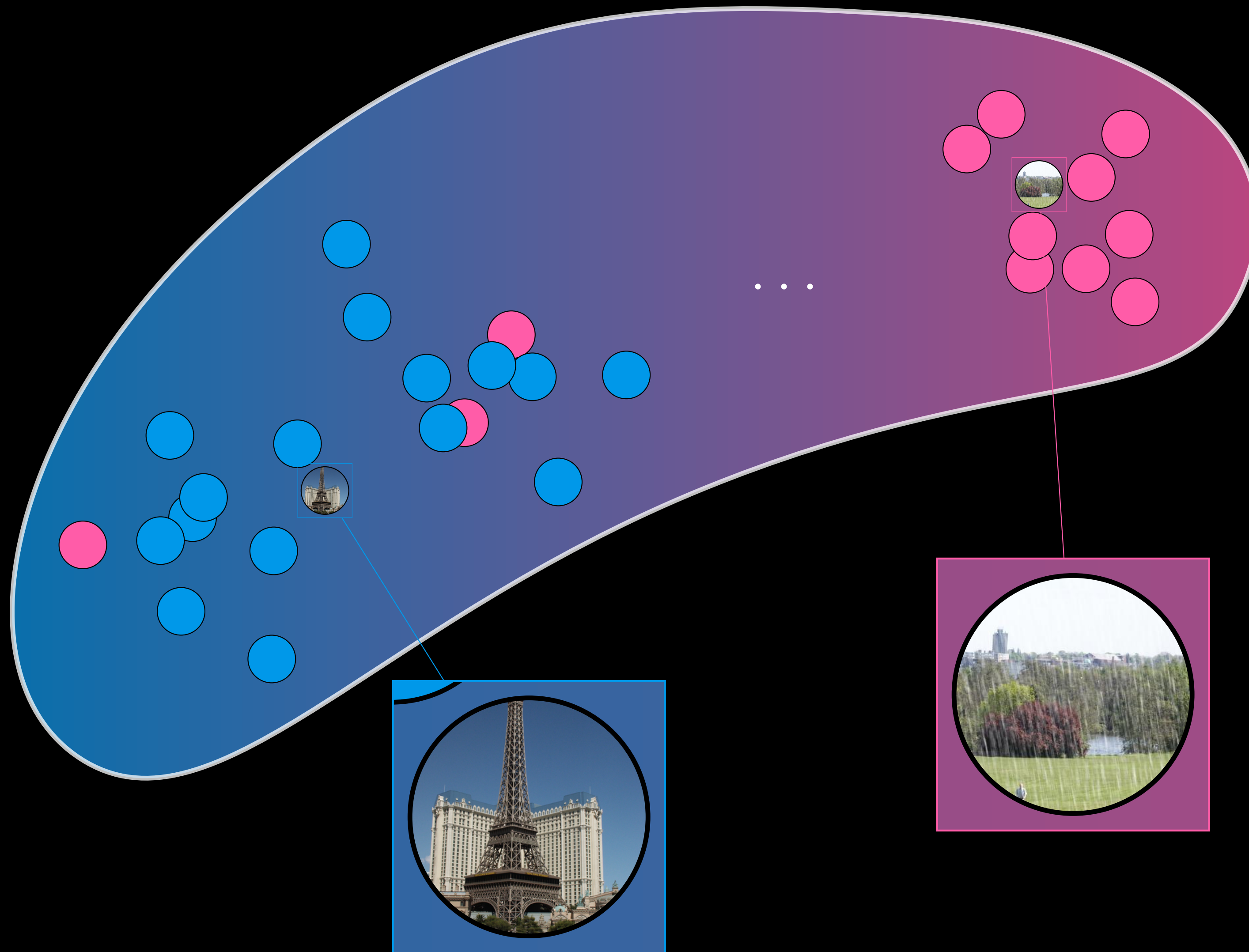
Test distribution D over $(x, y) \in \mathcal{X} \times \{0, 1\}$

Predictor $f: \mathcal{X} \rightarrow [0, 1]$ " $f(x)$ is confidence of $y(x) = 1$ "

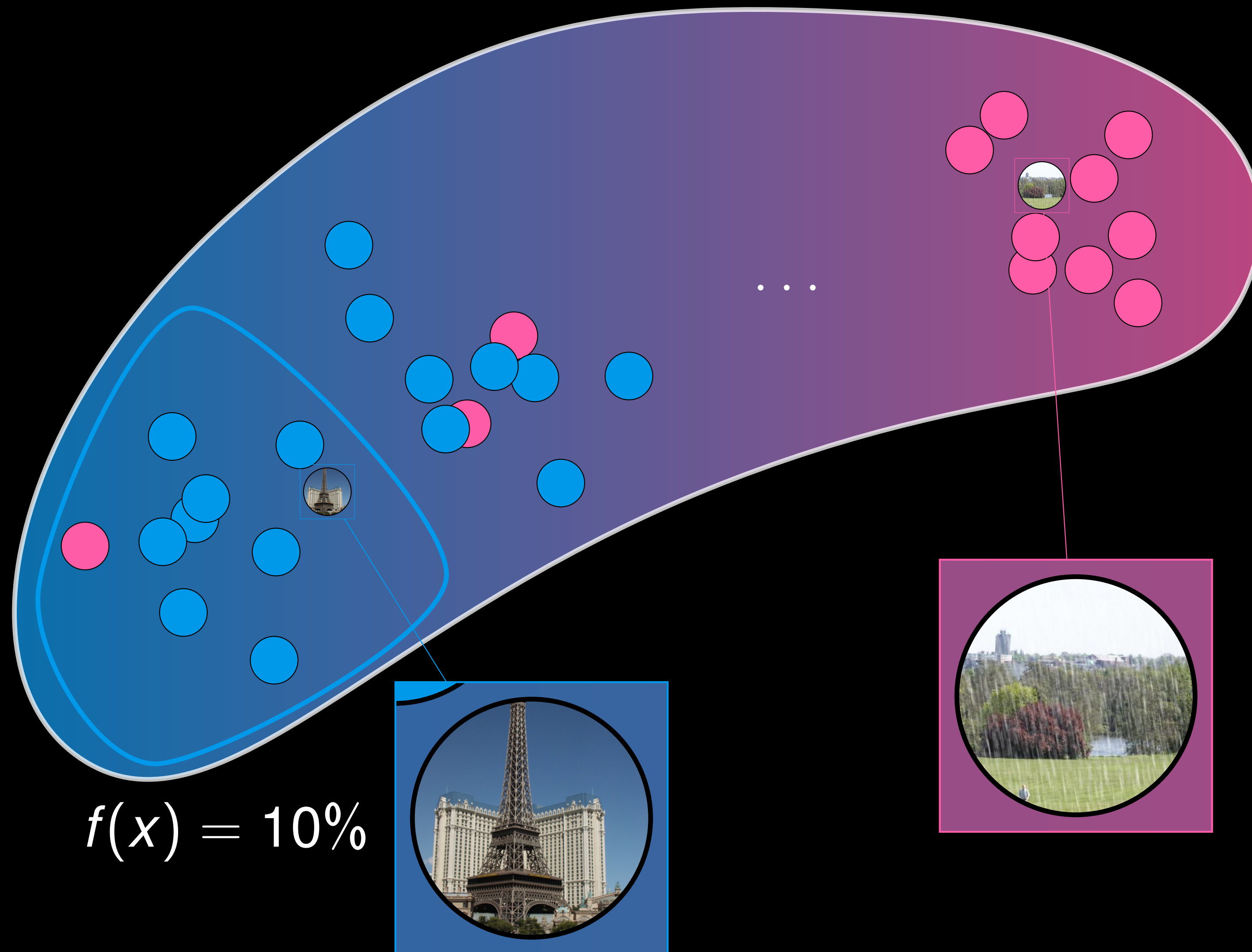


Perfect calibration is a property of the pair (f, D)

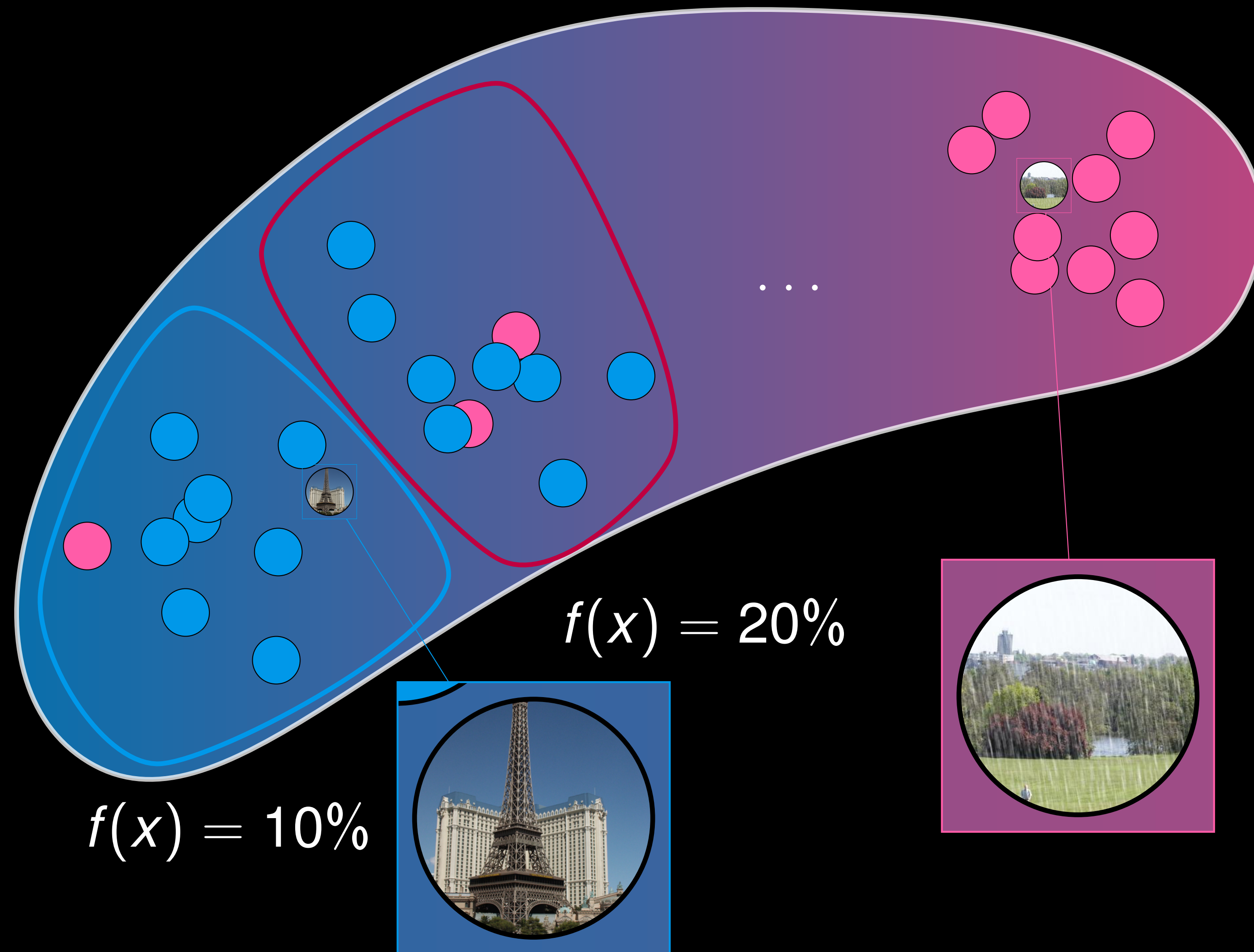
Calibration: Predictor f is perfectly calibrated w.r.t. distribution D if...



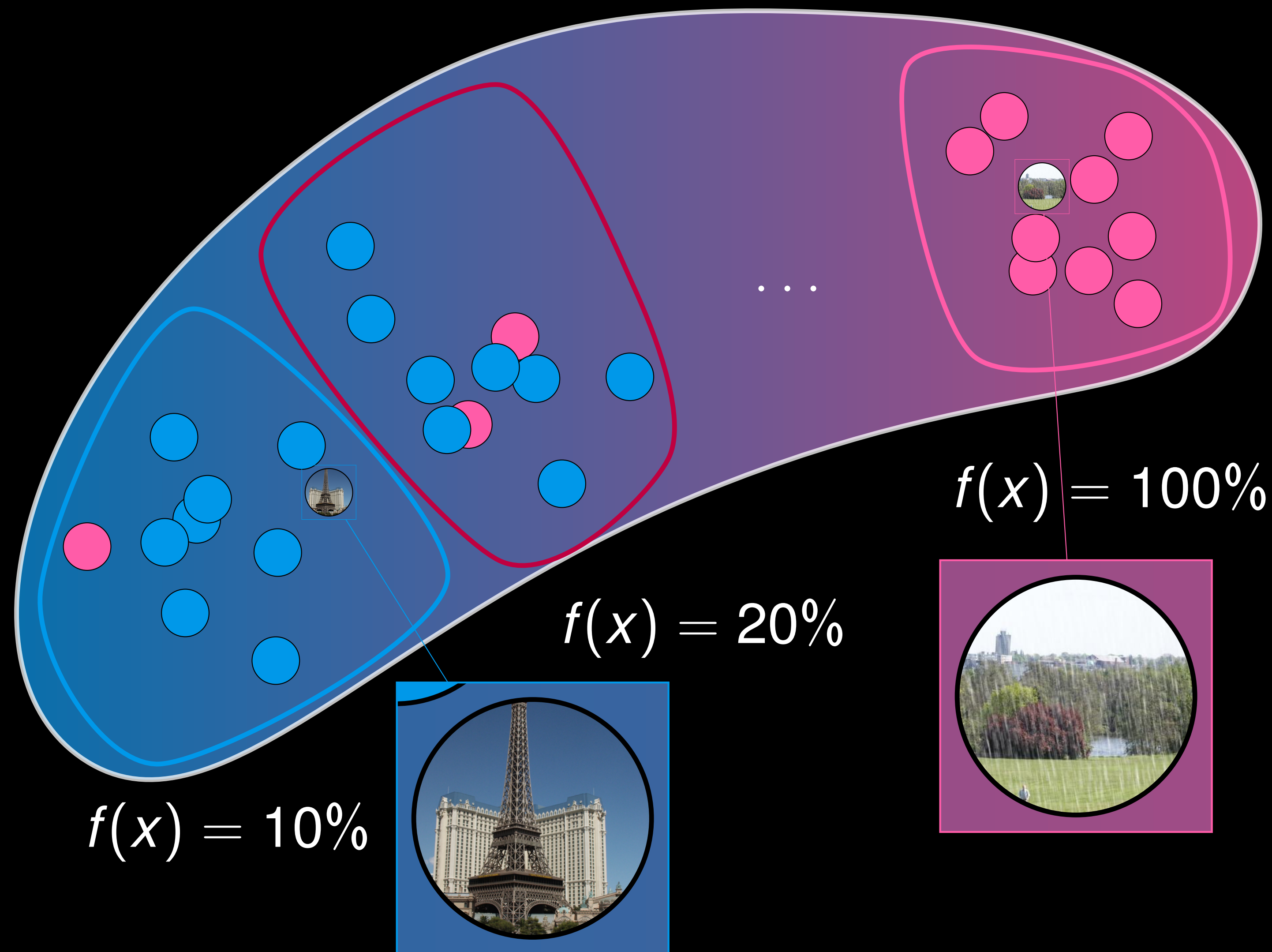
Calibration: Predictor f is perfectly calibrated w.r.t. distribution D if...



Calibration: Predictor f is perfectly calibrated w.r.t. distribution D if...



Calibration: Predictor f is perfectly calibrated w.r.t. distribution D if...



Perfect Calibration:

Predictor f is perfectly calibrated w.r.t. D if

$$\forall \ell \in [0,1] : \quad \mathbb{E}_{x,y \sim D}[y \mid f(x) = \ell] = \ell$$

Perfect Calibration:

Predictor f is perfectly calibrated w.r.t. D if

$$\forall \ell \in [0,1] : \quad \mathbb{E}_{x,y \sim D}[y \mid f(x) = \ell] = \ell$$

What's calibration good for?

1. **Interpretability:** $f(x)$ is a meaningful quantity, "confidence that $y = 1$ "
e.g. doctor informing patient of "80% probability of heart disease"
2. **Operational Uncertainty:** Systems downstream of $f(x)$ can behave differently on "high confidence" vs. "low confidence" inputs

$$\Pr[y=1 \mid f(x) = 0.5] = 0.5$$

$$\Pr[y=1 \mid f(x) = 0.8] = 0.8$$

$$\Pr[y=1 \mid f(x) = 1.0] = 1.0$$

Perfect Calibration:

Predictor f is perfectly calibrated w.r.t. D if

$$\forall \ell \in [0,1] : \quad \mathbb{E}_{x,y \sim D}[y \mid f(x) = \ell] = \ell$$

What's calibration good for?

1. **Interpretability:** $f(x)$ is a meaningful quantity, "confidence that $y = 1$ "
e.g. doctor informing patient of "80% probability of heart disease"
2. **Operational Uncertainty:** Systems downstream of $f(x)$ can behave differently on "high confidence" vs. "low confidence" inputs


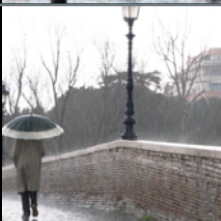

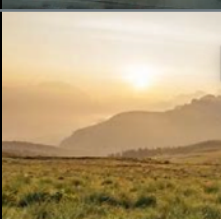

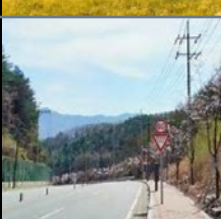
$$\Pr[y=1 \mid f(x) = 0.5] = 0.5$$

$$\Pr[y=1 \mid f(x) = 0.8] = 0.8$$

$$\Pr[y=1 \mid f(x) = 1.0] = 1.0$$

3. **Interesting:** "Models knows when it doesn't know" self-consistency

Calibration is **orthogonal** to accuracy

input x	ground-truth y	prediction $f(x)$
	1	0.5
	1	0.5
	1	0.5
	0	0.5
	0	0.5
	0	0.5

Goal

Understand when and why (DL) models are well-calibrated
(& what factors affect calibration)

This Talk

1. How to Define & Measure Miscalibration

"A Unifying Theory of Distance from Calibration." *STOC 2023*. *arXiv:2211.16886*. [Błasiok, Gopalan, Hu, **N.**]

2. Empirical Conjectures

"The Calibration Generalization Gap." *arXiv:2210.01964*. [Carrell, Mallinar, Lucas, **N.**]
+ upcoming work

3. Theory

"When Loss Minimization Yields Calibration." *In preparation*. [Błasiok, Gopalan, Hu, **N.**]

"Loss minimization yields multicalibration for large neural networks." *In submission*. [Kalai]

This Talk

1. How to Define & Measure Miscalibration

"A Unifying Theory of Distance from Calibration." *STOC 2023*. *arXiv:2211.16886*. [Błasiok, Gopalan, Hu, **N.**]

2. Empirical Conjectures

"The Calibration Generalization Gap." *arXiv:2210.01964*. [Carrell, Mallinar, Lucas, **N.**]
+ upcoming work

3. Theory

"When Loss Minimization Yields Calibration." *In preparation*. [Błasiok, Gopalan, Hu, **N.**]

"Loss minimization yields multicalibration for large neural networks." *In submission*. [Kalai]

Part 2.

Calibration of DNNs, Experimentally

Can we empirically characterize which DNNs have small calibration error?

Landscape of Calibration

- {Model, Data-distribution}



Landscape of Calibration



No single design choice determines calibration:

- Not just architecture
- Not just model size
- Not just test accuracy/loss

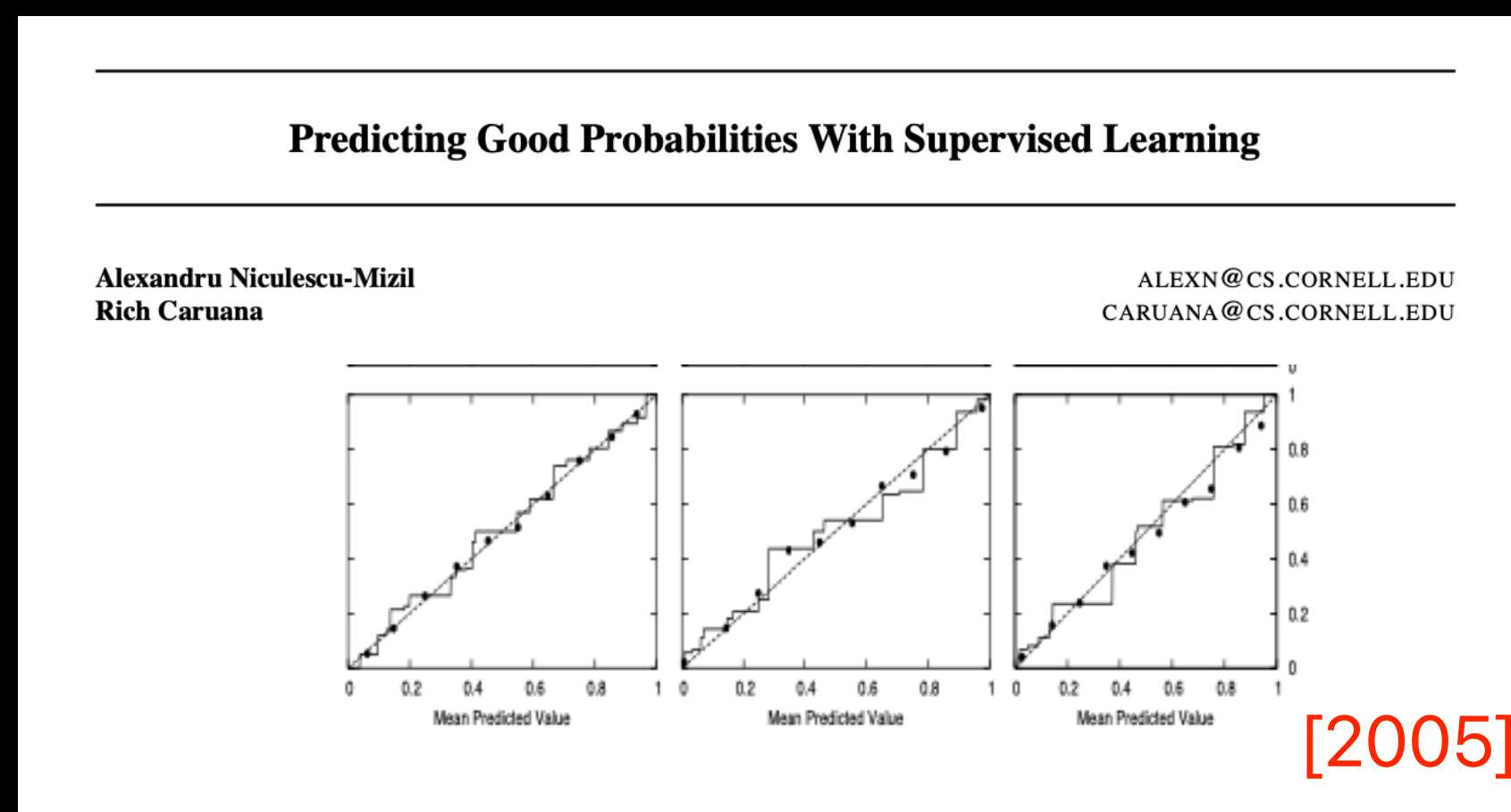
ex: small 3-layer MLP is well-calibrated on ImageNet

POORLY-
CALIBRATED

WELL-
CALIBRATED

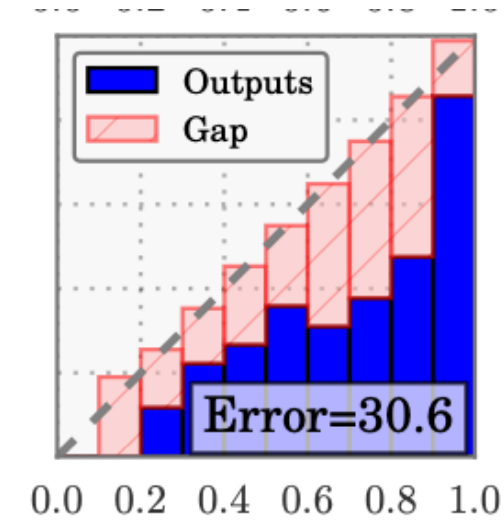
←
**POORLY-
CALIBRATED**

**WELL-
CALIBRATED** →



On Calibration of Modern Neural Networks

Chuan Guo^{*1} Geoff Pleiss^{*1} Yu Sun^{*1} Kilian Q. Weinberger¹



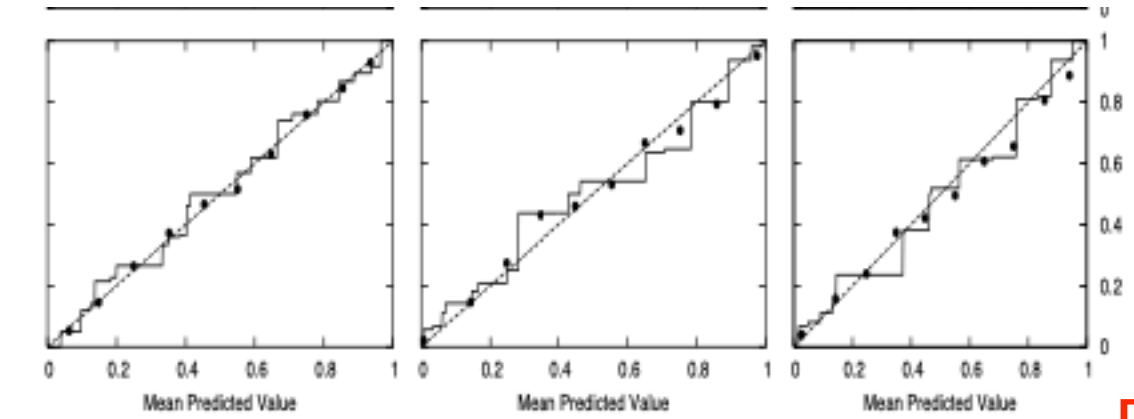
ResNet (2016)
CIFAR-100

[2017]

Predicting Good Probabilities With Supervised Learning

Alexandru Niculescu-Mizil
Rich Caruana

ALEXN@CS.CORNELL.EDU
CARUANA@CS.CORNELL.EDU



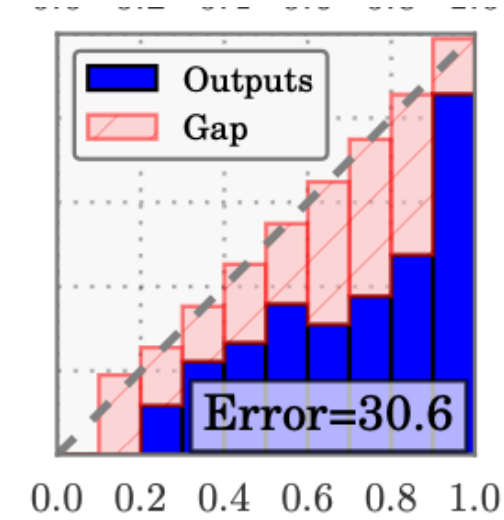
[2005]

←
POORLY-
CALIBRATED

→
WELL-
CALIBRATED

On Calibration of Modern Neural Networks

Chuan Guo^{*1} Geoff Pleiss^{*1} Yu Sun^{*1} Kilian Q. Weinberger¹



ResNet (2016)
CIFAR-100

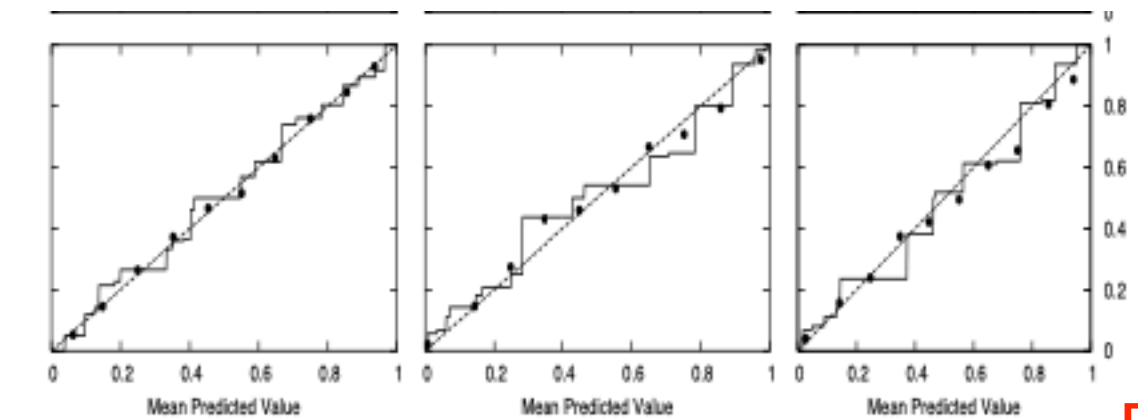
[2017]

←
**POORLY-
CALIBRATED**

Predicting Good Probabilities With Supervised Learning

Alexandru Niculescu-Mizil
Rich Caruana

ALEXN@CS.CORNELL.EDU
CARUANA@CS.CORNELL.EDU



[2005]

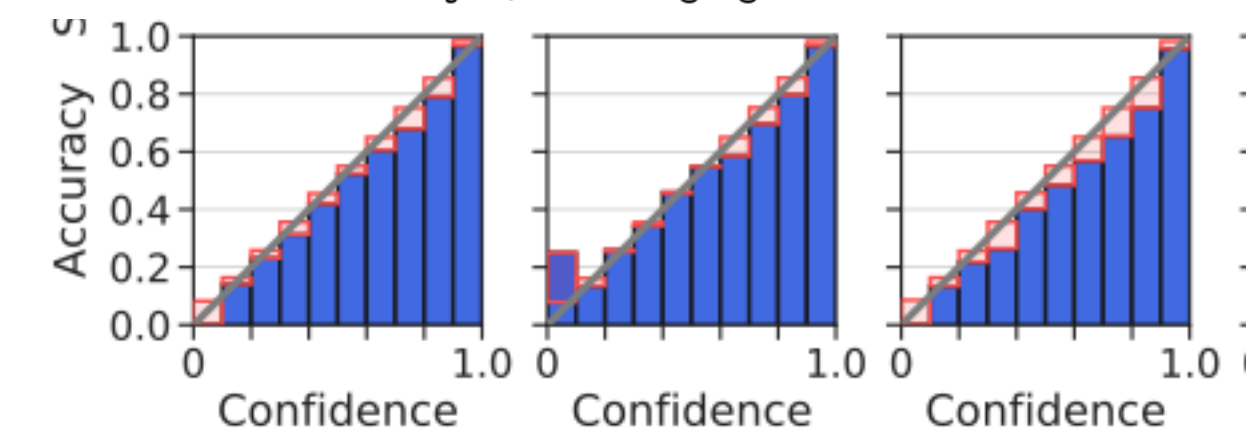
**WELL-
CALIBRATED**

Revisiting the Calibration of Modern Neural Networks

Matthias Minderer Josip Djolonga Rob Romijnders Frances Hubis
Xiaohua Zhai Neil Houlsby Dustin Tran Mario Lucic

Google Research, Brain Team

{mjlm, lucic}@google.com



[2021]

Proposal: Study test-calibration as we study test-error.

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Proposal: Study test-calibration as we study test-error.

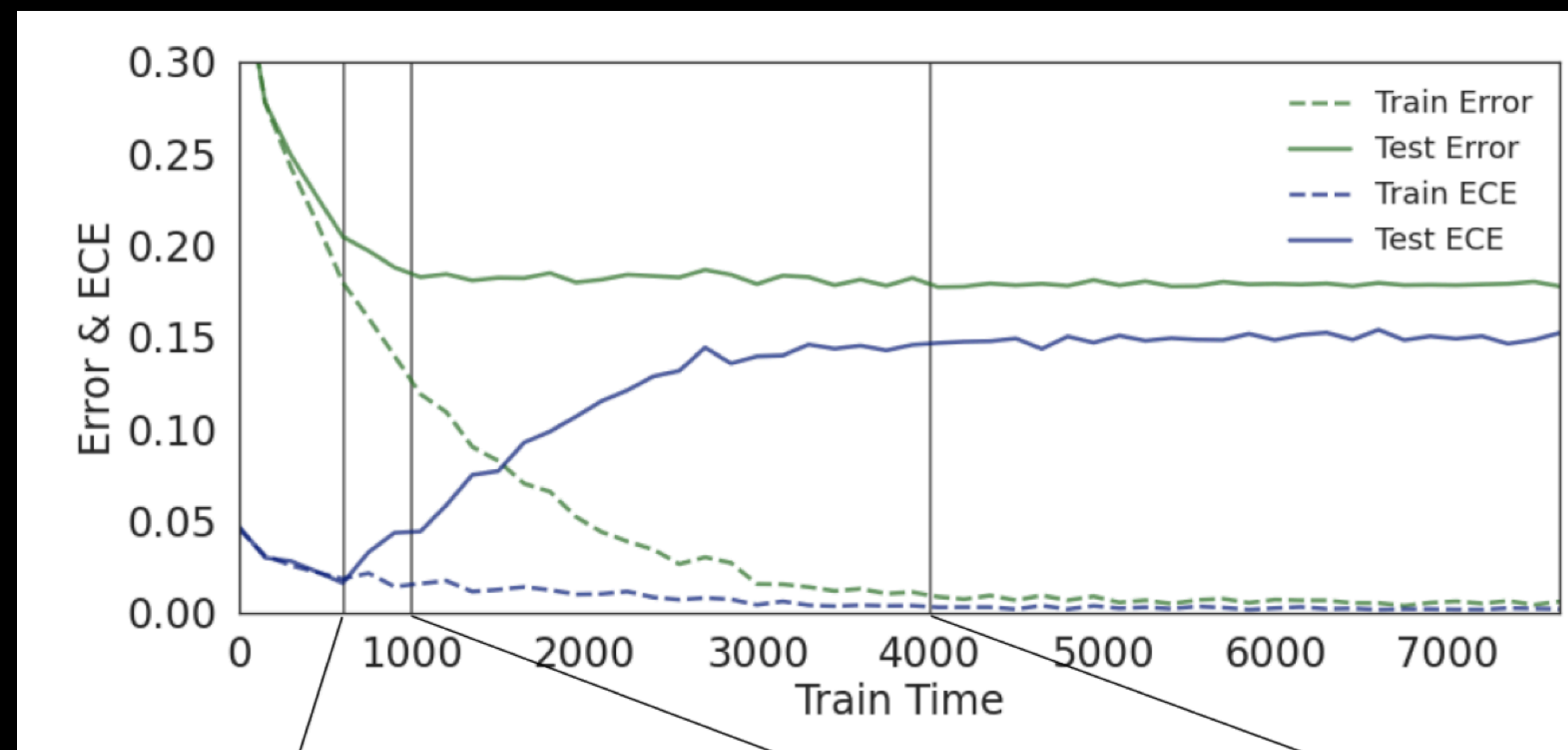
Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

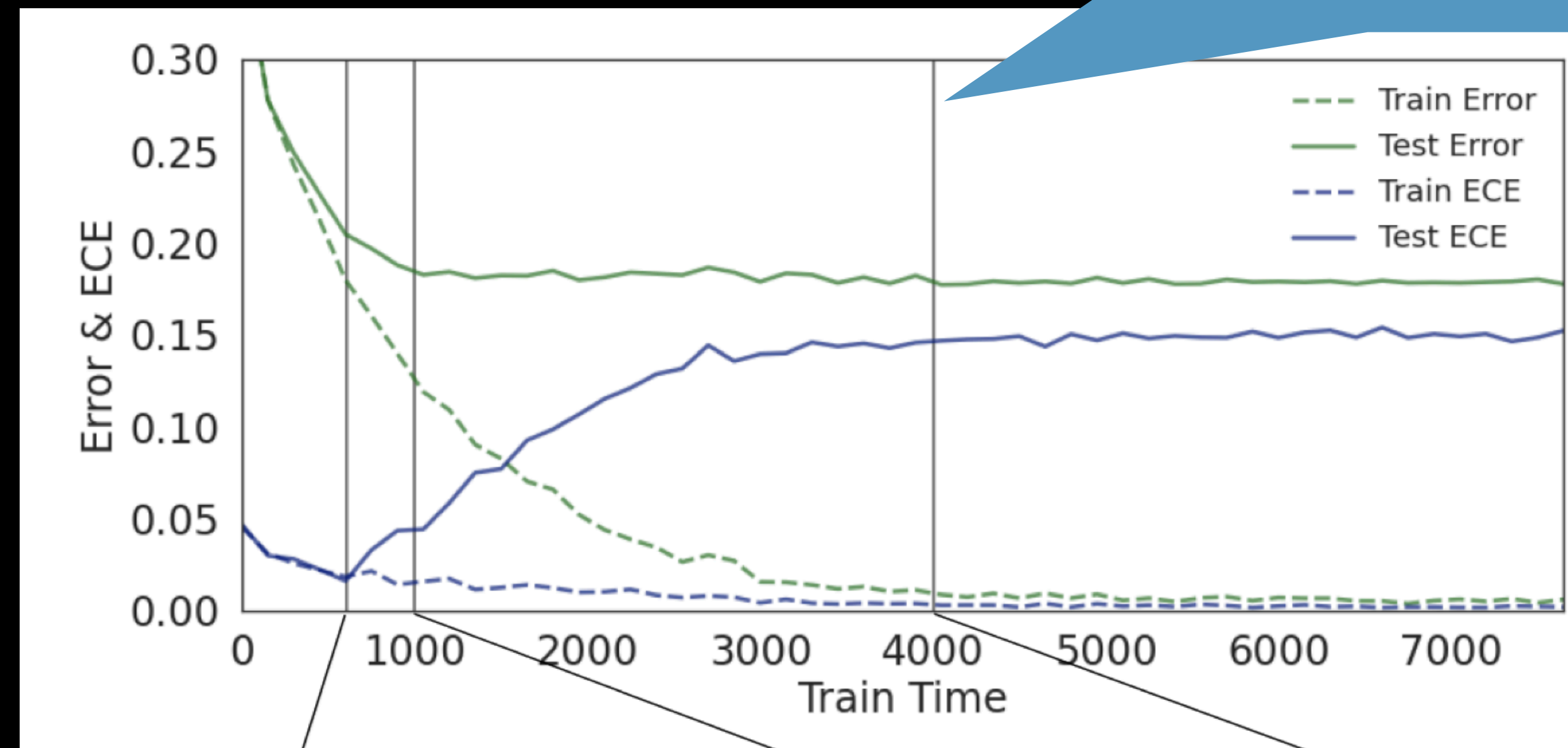
Trivial, BUT:

1. Suggests methodology: study each part separately
2. Insightful: parts are simpler than the whole

(higher=worse)



(higher=worse) ↑



End of training:
models are overconfident

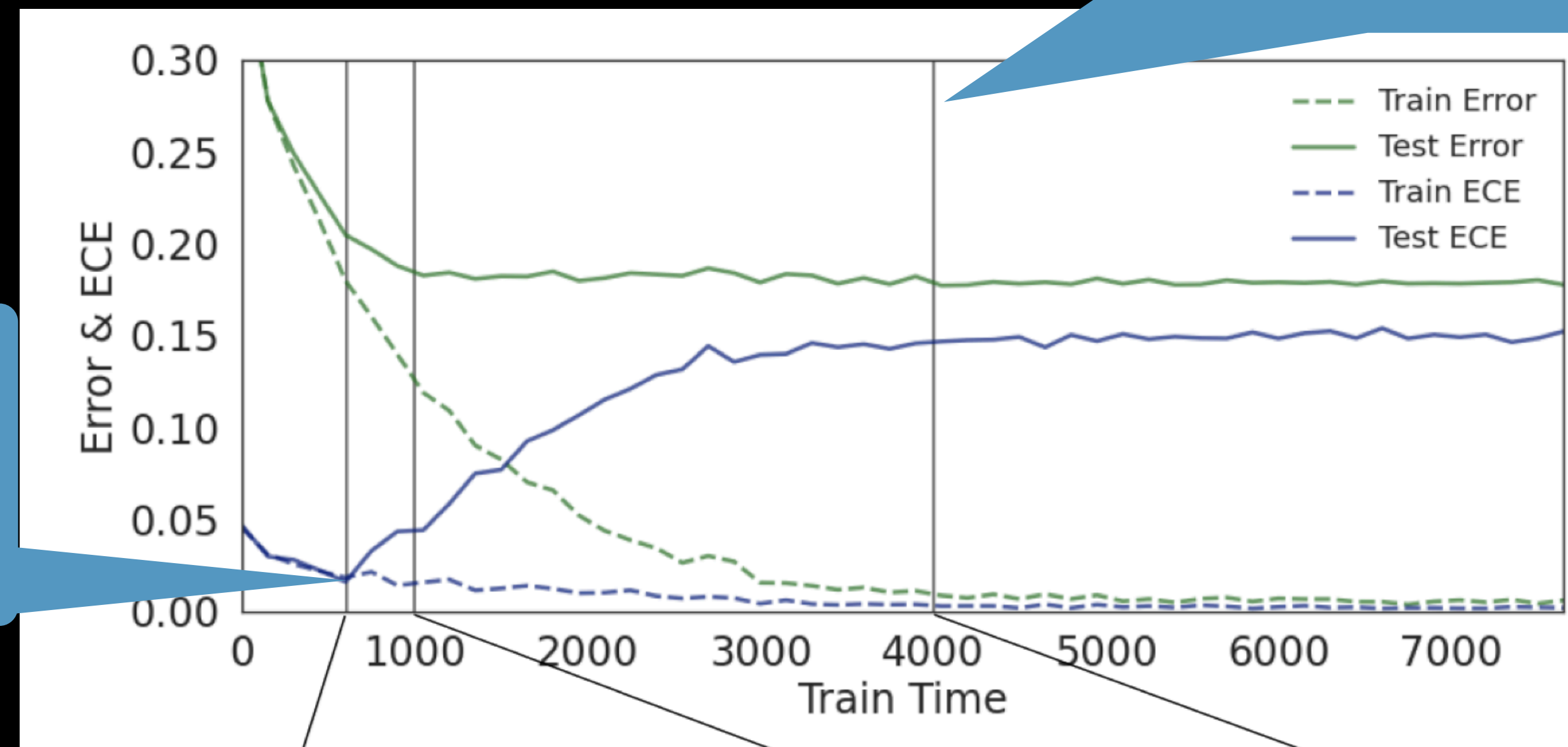
Train {Error,CE} ≈ 0

Test {Error,CE} $\gg 0$

(higher=worse) ↑

Throughout training:

Train CE ≈ 0



End of training:
models are overconfident

Train {Error,CE} ≈ 0

Test {Error,CE} $\gg 0$

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

* $\text{depth} \geq 2$, trained with proper scoring rule, no severe augmentations, ...

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

* $\text{depth} \geq 2$, trained with proper scoring rule, no severe augmentations, ...

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

Even when underfitting!

* depth ≥ 2 , trained with proper scoring rule, no severe augmentations, ...

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

Even when underfitting!

* depth ≥ 2 , trained with proper scoring rule, no severe augmentations, ...

Fundamental Decomposition

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration Error on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration Error on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}} \leq \text{Error Generalization Gap}$$

Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

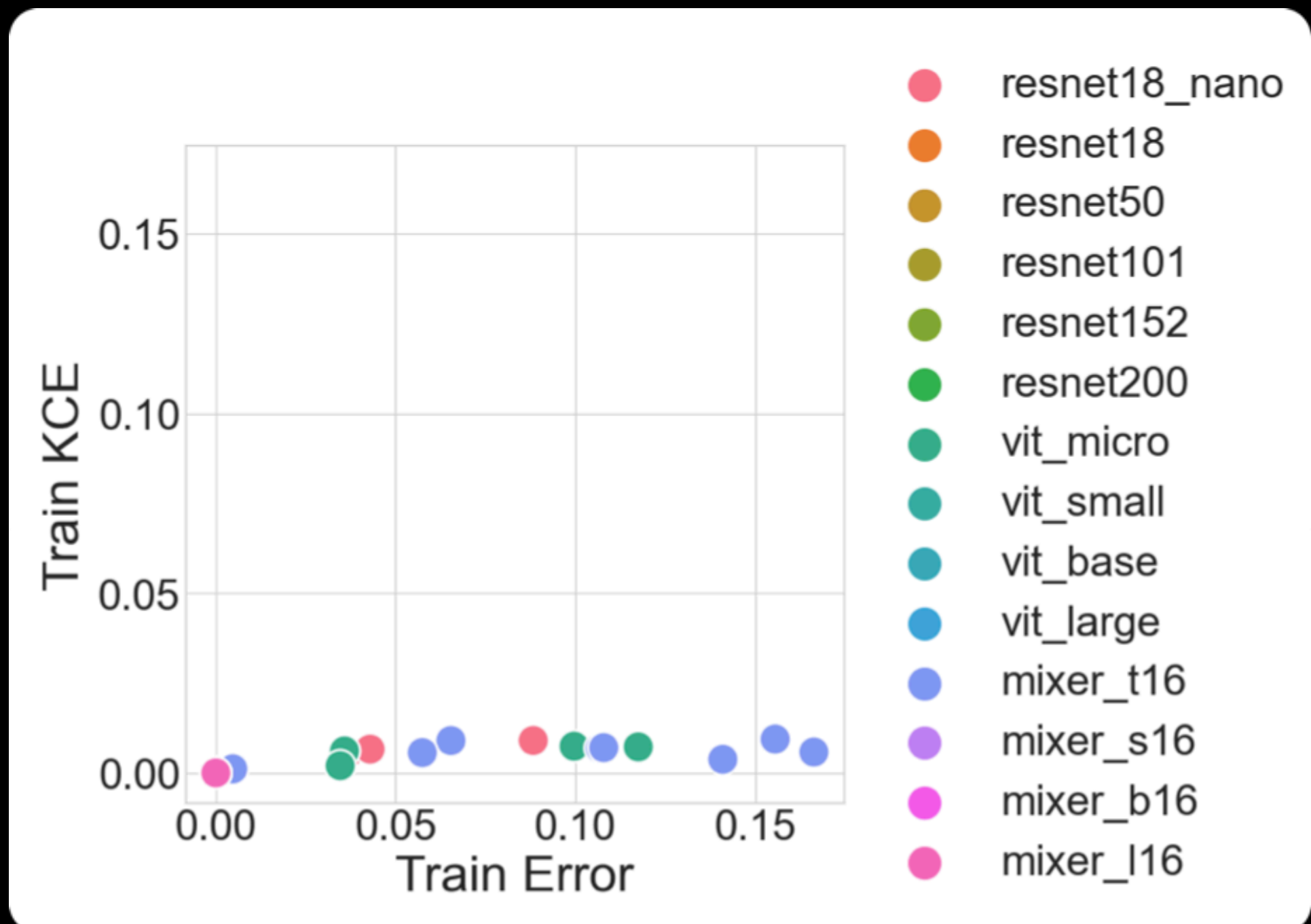
Even when underfitting!

* depth ≥ 2 , trained with proper scoring rule, no severe augmentations, ...

Empirical Claim 1

For almost all* DNNs

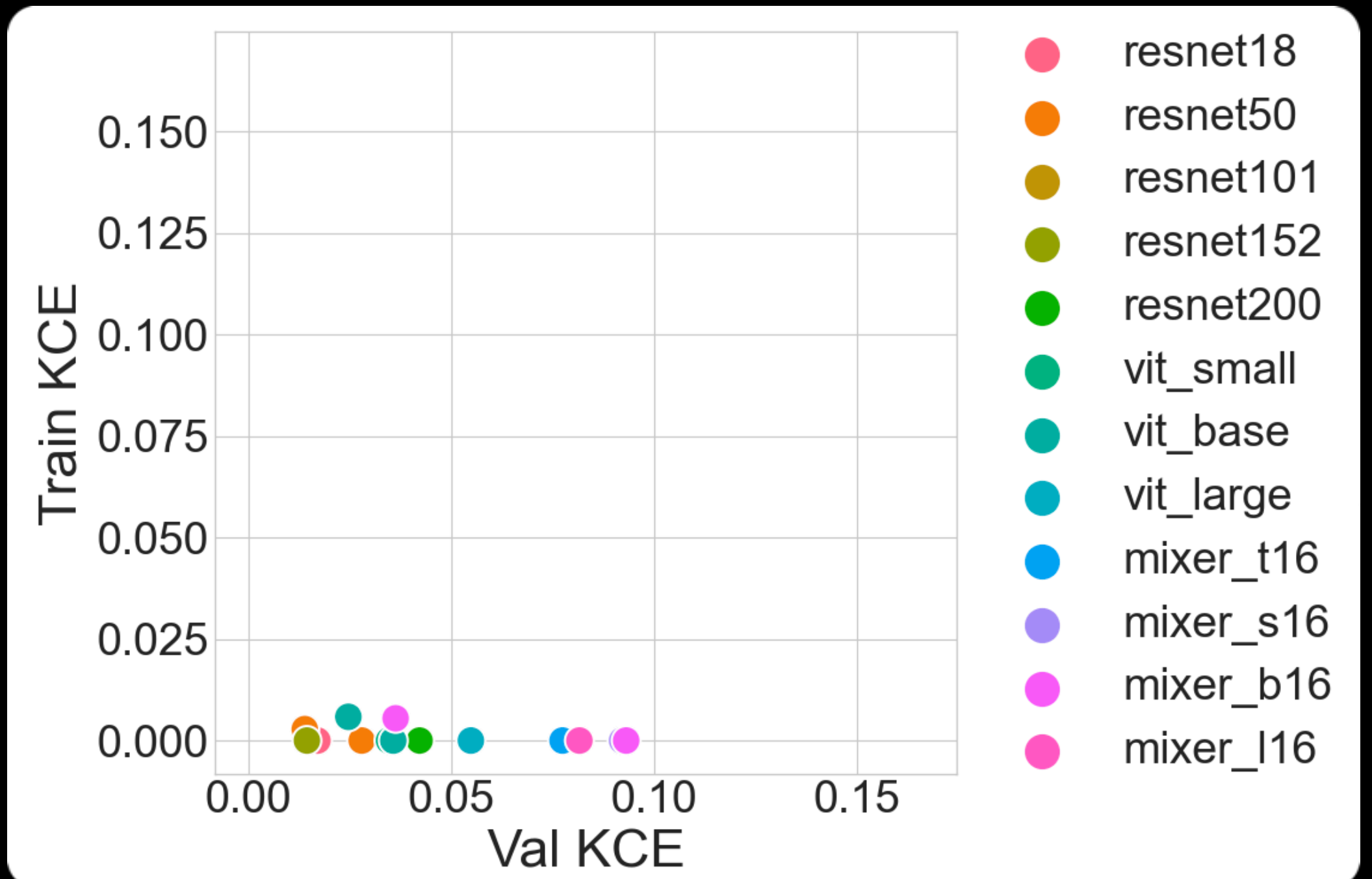
$$\mu_{\text{Train}} \approx 0$$



Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

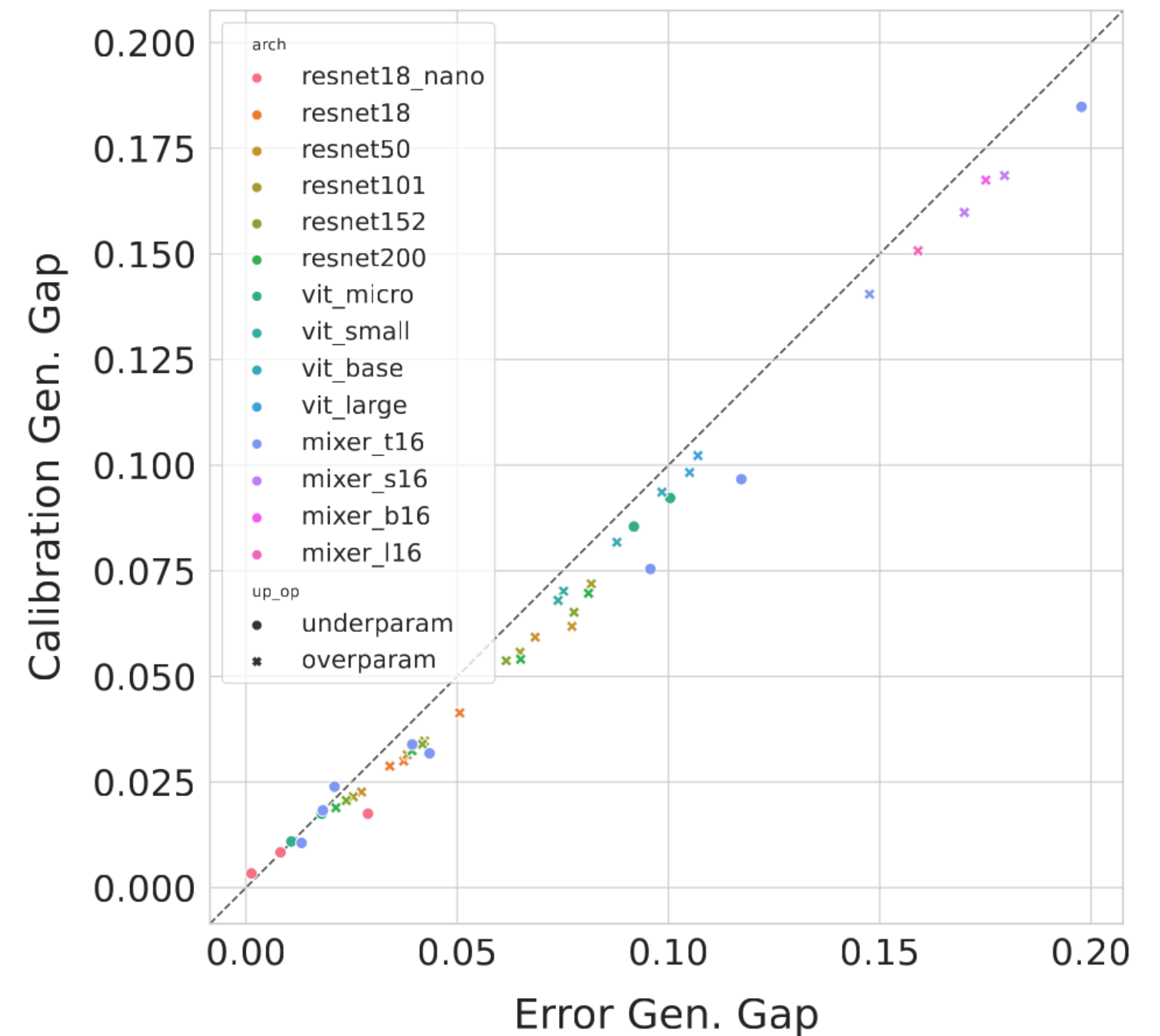


(binary-ImageNet)

Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$



(binary-ImageNet)

Takeaways

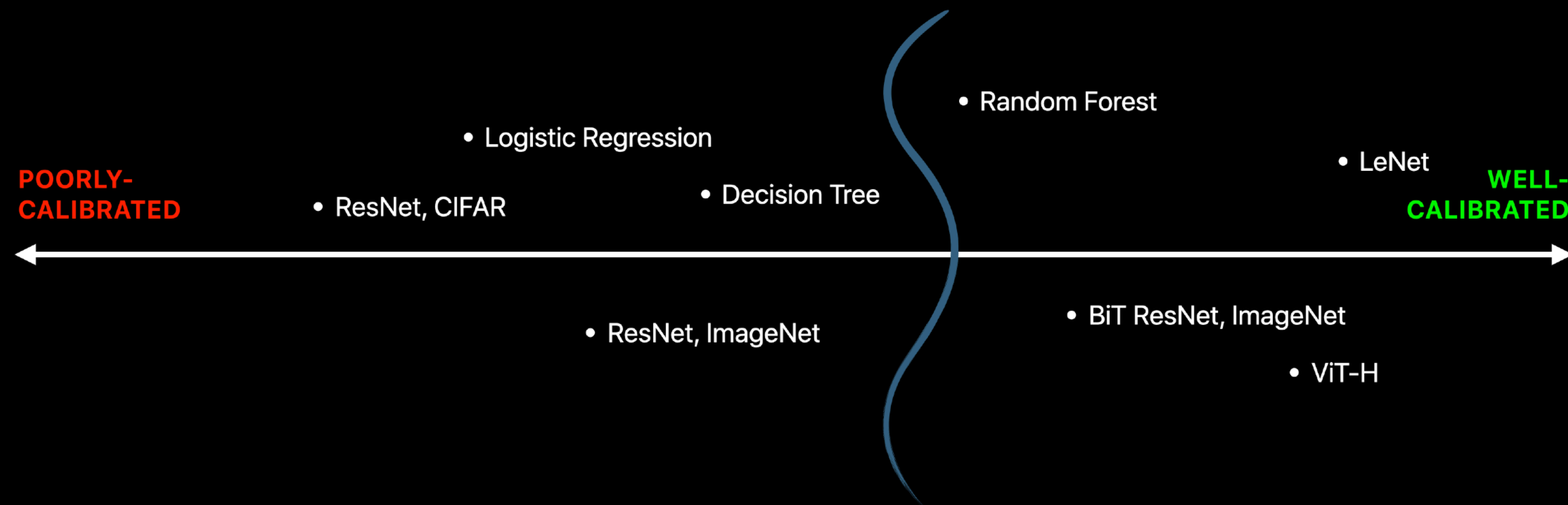
$$(\text{Test Calibration Error}) \approx |\text{Train Error} - \text{Test Error}|$$

“Models with small generalization-gap are typically well-calibrated”

Takeaways

$$(\text{Test Calibration Error}) \approx |\text{Train Error} - \text{Test Error}|$$

“Models with small generalization-gap are typically well-calibrated”



Takeaways

$$(\text{Test Calibration Error}) \approx |\text{Train Error} - \text{Test Error}|$$

“Models with small generalization-gap are typically well-calibrated”

The following are well-calibrated:

1. Small models, on large data-sets (e.g. large vision models)
2. All models trained for 1-epoch (e.g. LLMs)

Applications

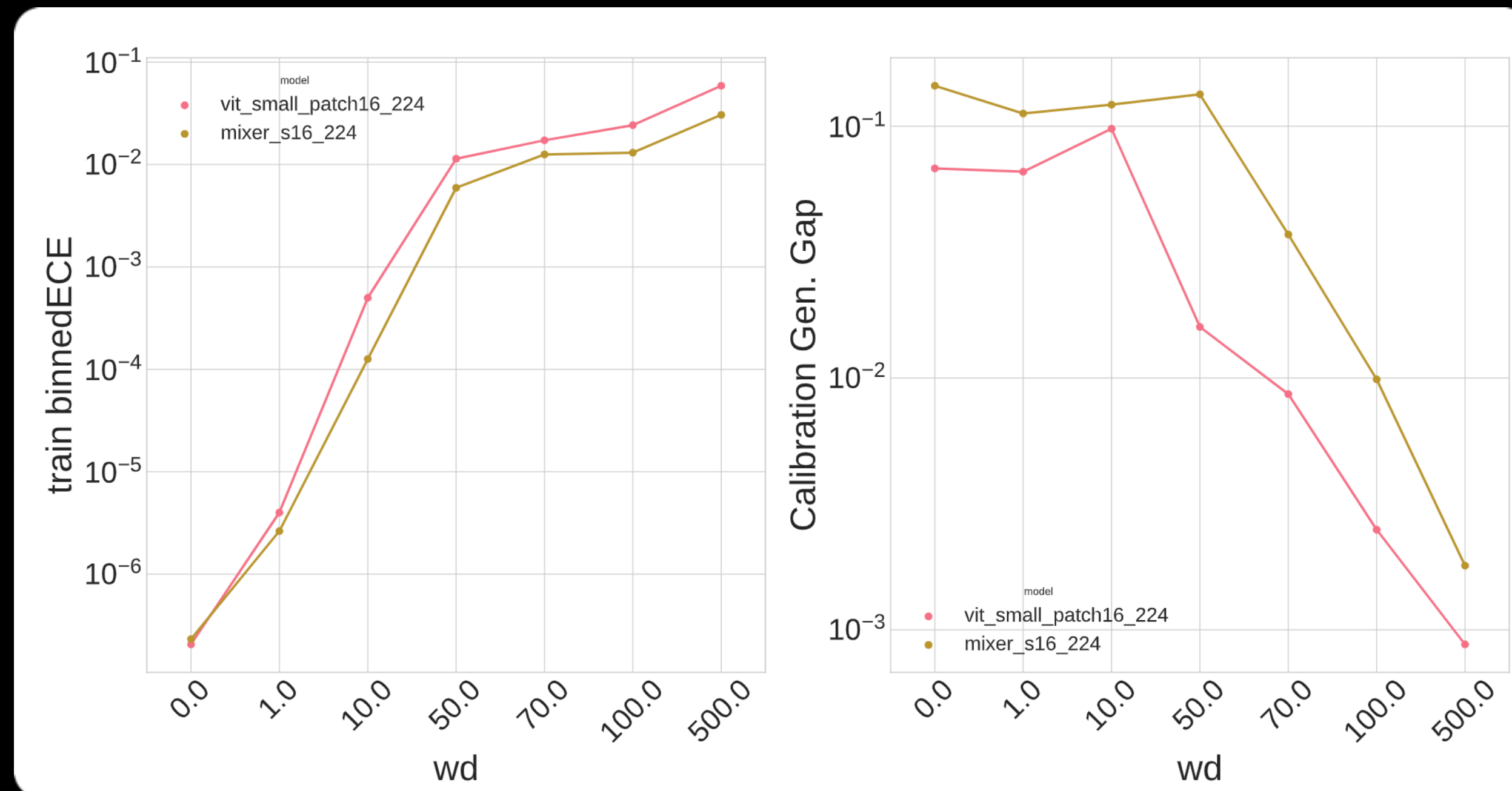
$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

For any intervention (changing the augmentation, regularizer, etc), study its effect on:

- (1) Train calibration
- (2) Generalization gap

Applications: Regularization Strength

$$\underbrace{\mu^{\text{Test}}}_{\text{Calibration on Test Set}} \leq \underbrace{\mu^{\text{Train}}}_{\text{Calibration on Train Set}} + \underbrace{|\mu^{\text{Test}} - \mu^{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$



Applications: Data Augmentation

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$

“Standard” data-augmentation (measure-preserving):

- Same TrainCE; Shrinks generalization gap

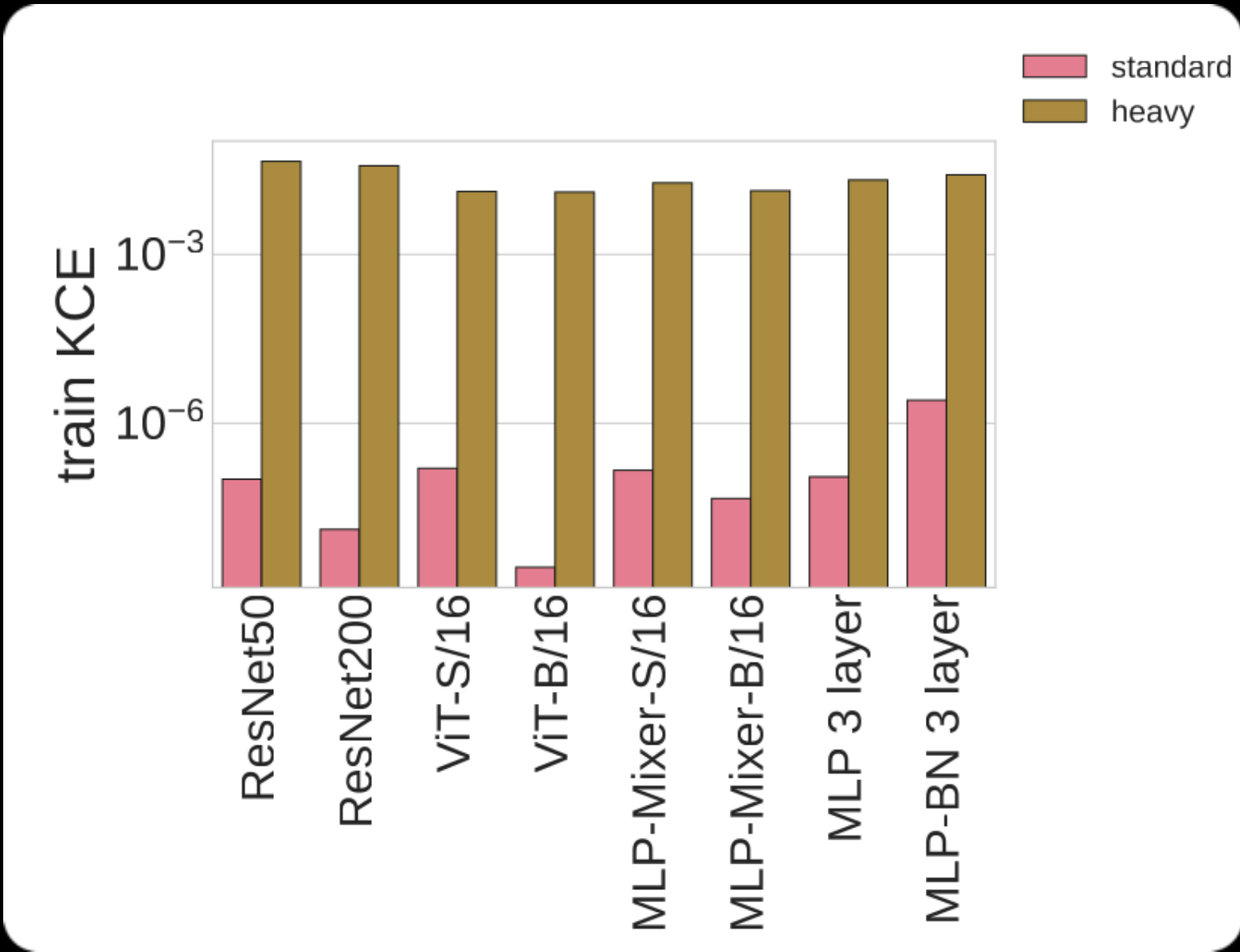
“Exotic” data-augmentation:

- *Increases TrainCE*; Shrinks generalization gap



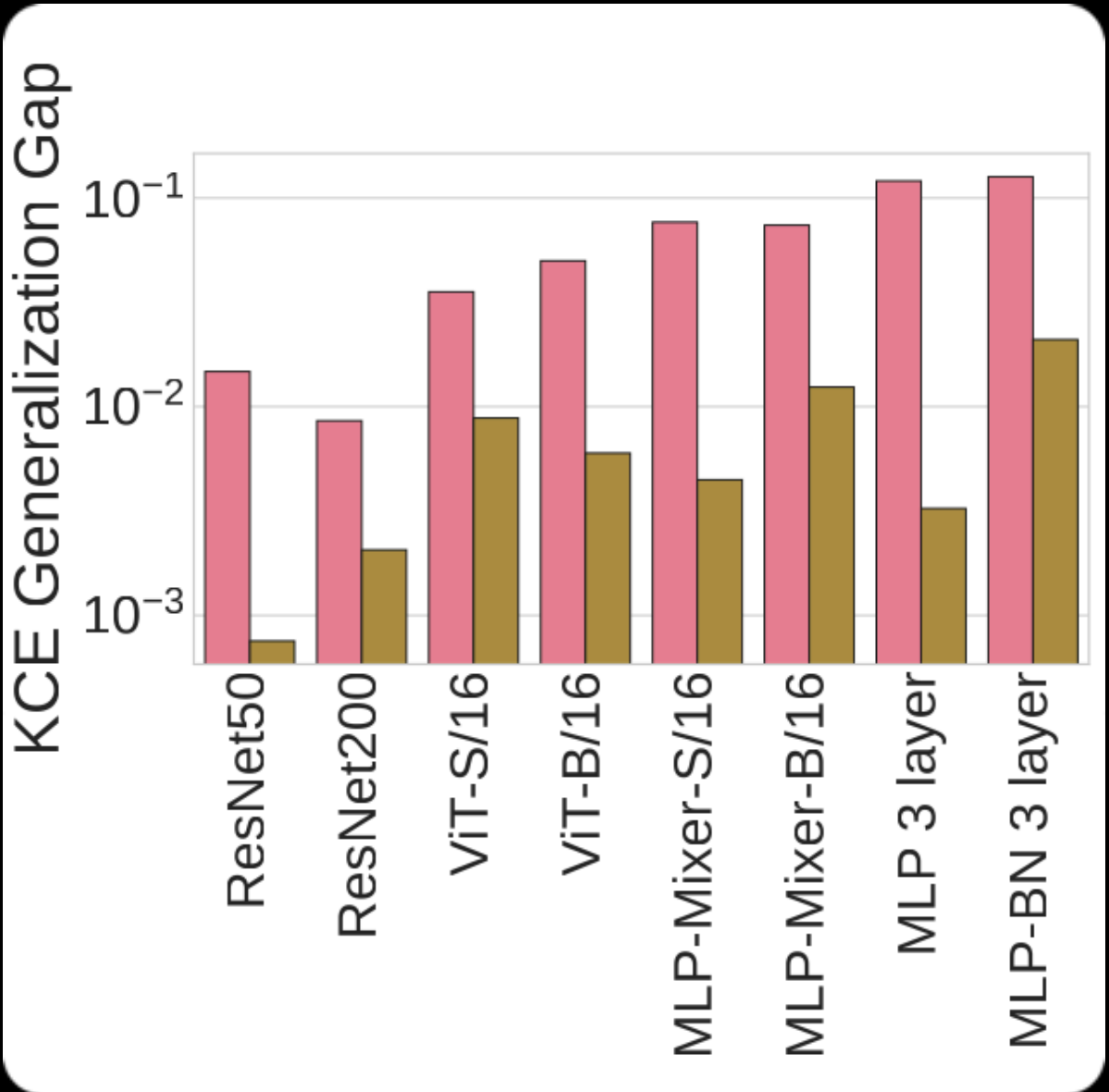
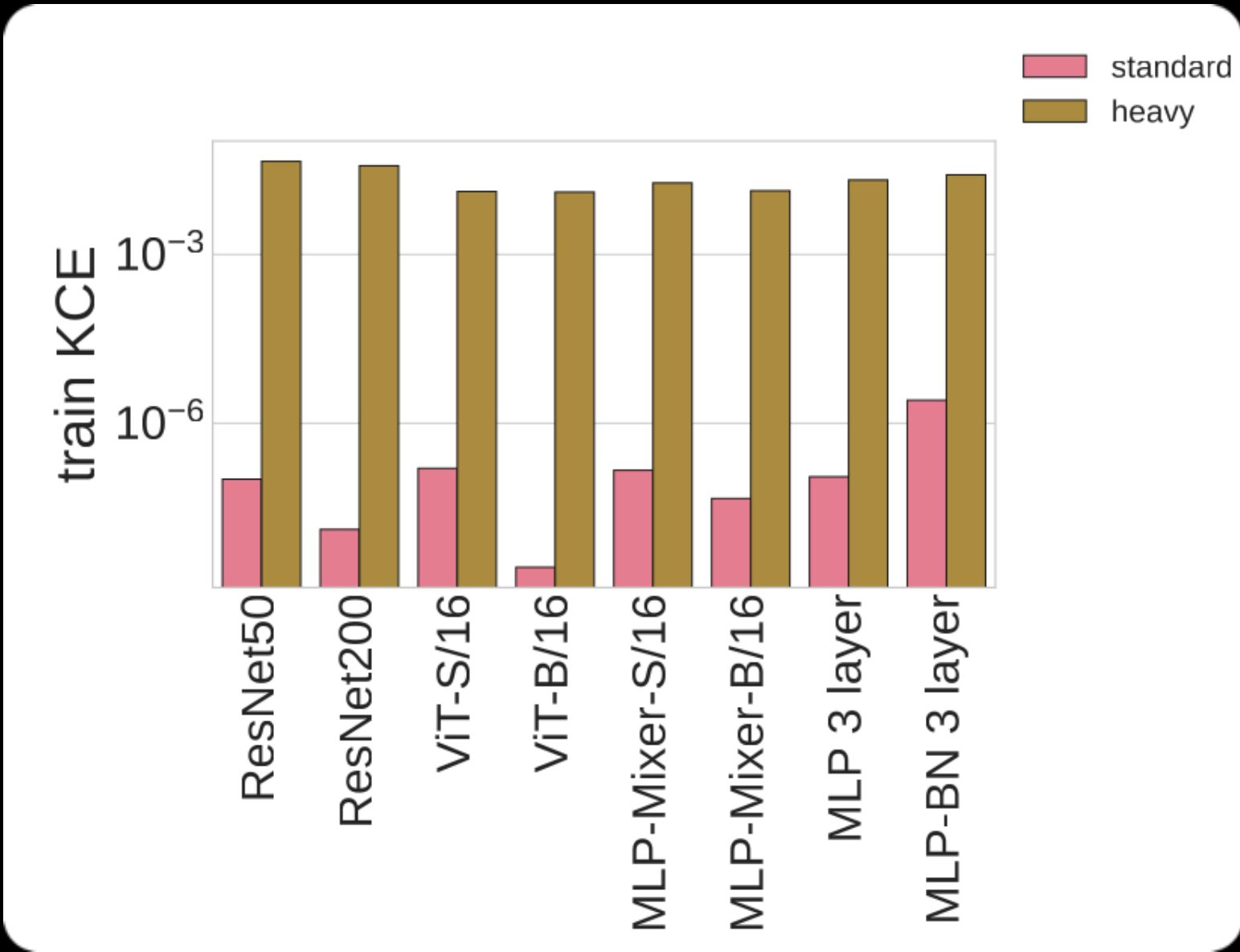
Applications: Data Augmentation

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$



Applications: Data Augmentation

$$\underbrace{\mu_{\text{Test}}}_{\text{Calibration on Test Set}} \leq \underbrace{\mu_{\text{Train}}}_{\text{Calibration on Train Set}} + \underbrace{|\mu_{\text{Test}} - \mu_{\text{Train}}|}_{\text{Calibration Generalization Gap}}$$



Part 3. Theory

When are Claims 1 & 2 provably true?

Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

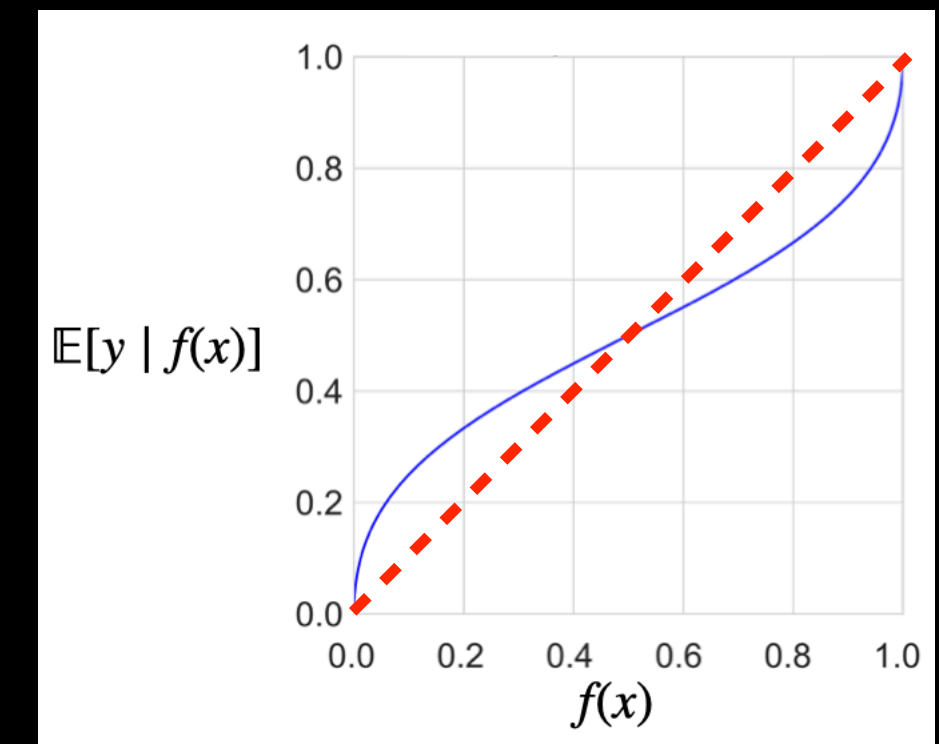
Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

Assumption 1: f is overconfident on test.

$$\forall \ell \in [0, 1] : \mathbb{E}_{\mathbb{D}_{\text{test}}} [\text{Acc}(f, y) \mid f = \ell] \leq \text{Conf}(v)$$



Empirical Claim 2

For almost all* DNNs

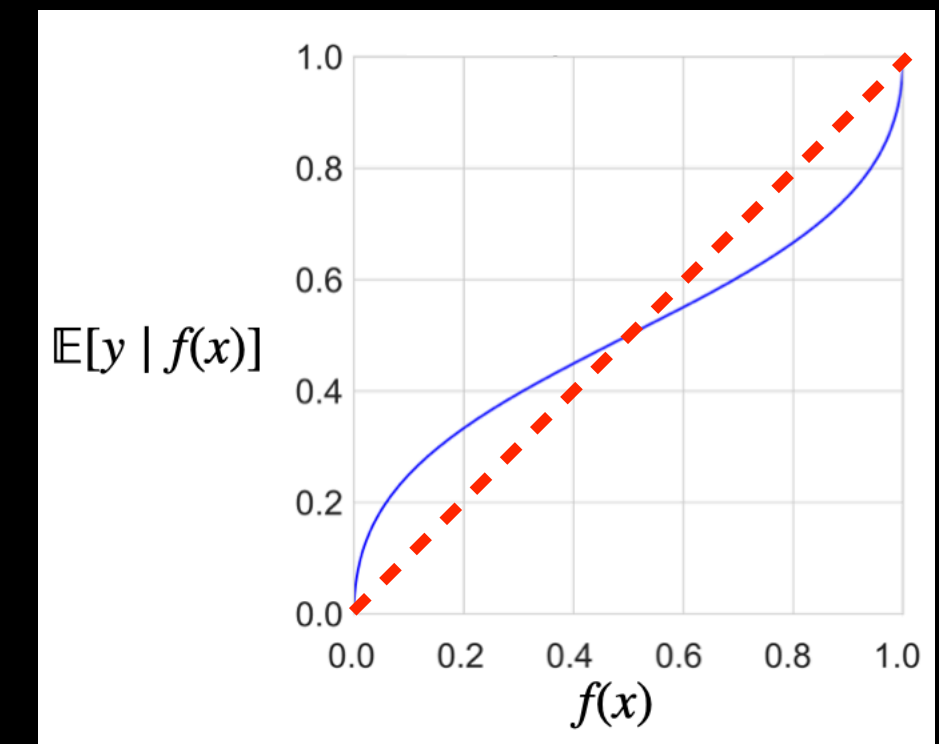
$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

Assumption 1: f is overconfident on test.

$$\forall \ell \in [0, 1] : \mathbb{E}_{\mathbb{D}_{\text{test}}} [\text{Acc}(f, y) \mid f = \ell] \leq \text{Conf}(v)$$

Assumption 2: f is more confident on train than on test.

$$\mathbb{E}_{\mathbb{D}_{\text{train}}} [\text{Conf}(f)] \geq \mathbb{E}_{\mathbb{D}_{\text{test}}} [\text{Conf}(f)]$$



Empirical Claim 2

For almost all* DNNs

$$|\mu_{\text{Test}} - \mu_{\text{Train}}| \leq |\text{TestError} - \text{TrainError}|$$

Assumption 1: f is overconfident on test.

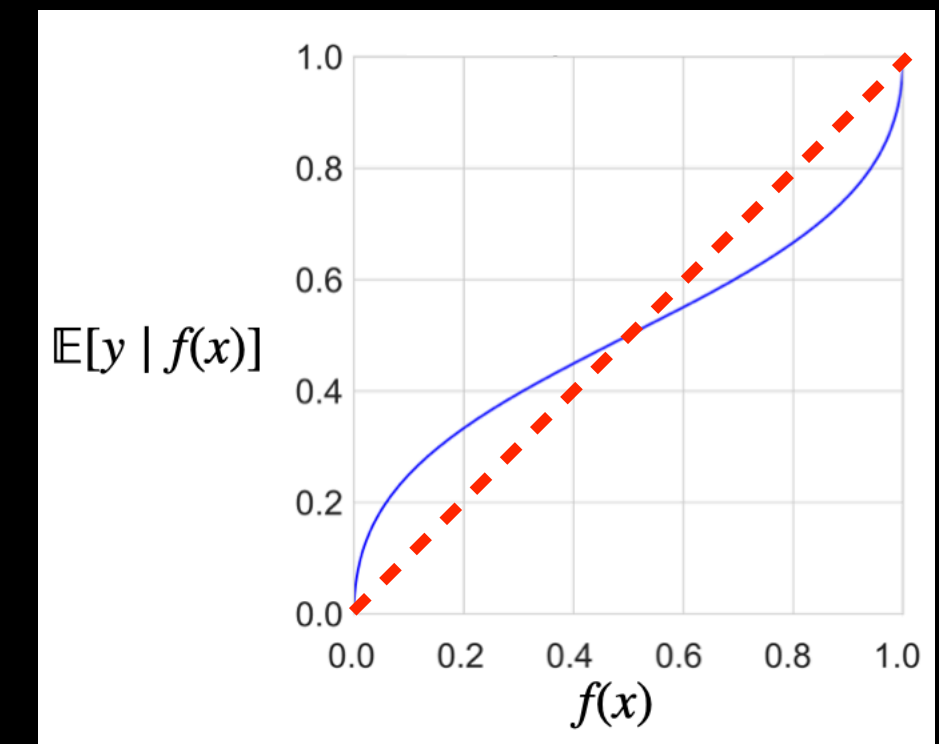
$$\forall \ell \in [0, 1] : \mathbb{E}_{\mathbb{D}_{\text{test}}} [\text{Acc}(f, y) \mid f = \ell] \leq \text{Conf}(v)$$

Assumption 2: f is more confident on train than on test.

$$\mathbb{E}_{\mathbb{D}_{\text{train}}} [\text{Conf}(f)] \geq \mathbb{E}_{\mathbb{D}_{\text{test}}} [\text{Conf}(f)]$$

Theorem 2 (Calibration Generalization Bound). *Under Assumptions 2 and 3, we have*

$$\text{ECE}(\mathbb{D}_{\text{test}}) - \text{ECE}(\mathbb{D}_{\text{train}}) \leq \text{Error}(\mathbb{D}_{\text{test}}) - \text{Error}(\mathbb{D}_{\text{train}})$$



Empirical Claim 1

For almost all* DNNs

$$\mu_{\text{Train}} \approx 0$$

Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

What we'd like to do:

Exactly minimize expected loss, over
all functions:

$$f^* = \operatorname{argmin}_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f(x), y)]$$

Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

What we'd like to do:

Exactly minimize expected loss, over
all functions:

$$f^* = \operatorname{argmin}_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f(x), y)]$$

$$\implies f^*(x) = p_{\mathcal{D}}(y|x)$$

perfectly calibrated

Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

What we'd like to do:

Exactly minimize expected loss, over *all functions*:

$$f^* = \operatorname{argmin}_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f(x), y)]$$

$\implies f^*(x) = p_{\mathcal{D}}(y|x)$
perfectly calibrated

What we actually do:

Run SGD* to *approximately minimize* expected loss, over *restricted family* $\{f_{\theta} : \theta \in \Theta\}$:

$$\tilde{f} = \operatorname{SGDmin}_{\theta \in \Theta} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

What we'd like to do:

Exactly minimize expected loss, over *all functions*:

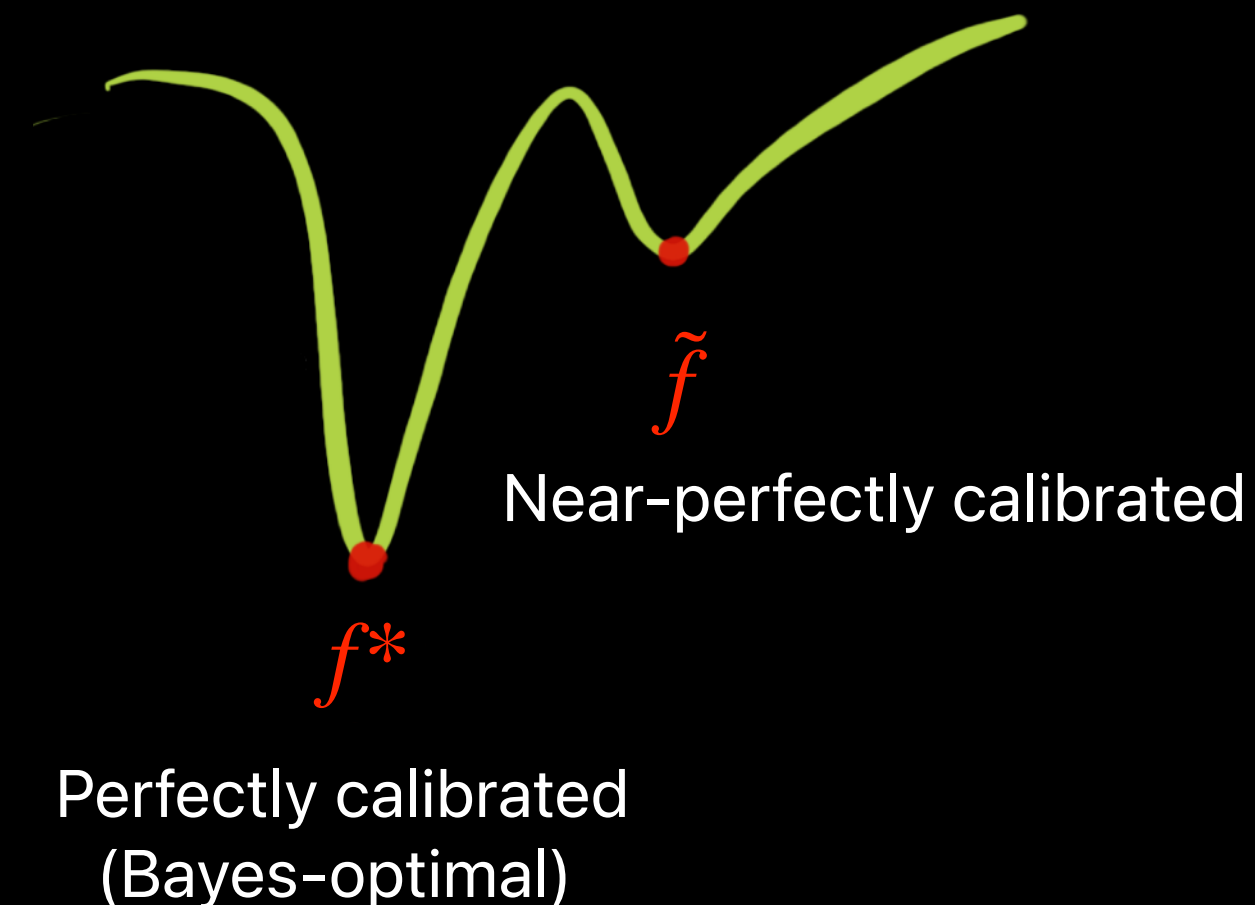
$$f^* = \operatorname{argmin}_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f(x), y)]$$

$\Rightarrow f^*(x) = p_{\mathcal{D}}(y|x)$
perfectly calibrated

What we actually do:

Run SGD* to *approximately minimize* expected loss, over *restricted family* $\{f_{\theta} : \theta \in \Theta\}$:

$$\tilde{f} = \operatorname{SGDmin}_{\theta \in \Theta} \mathbb{E}_{x,y \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$



Given: Distribution $D = \widehat{D} = \{(x_i, y_i)\}_{i \in [n]}$

What we'd like to do:

Exactly minimize expected loss, over *all functions*:

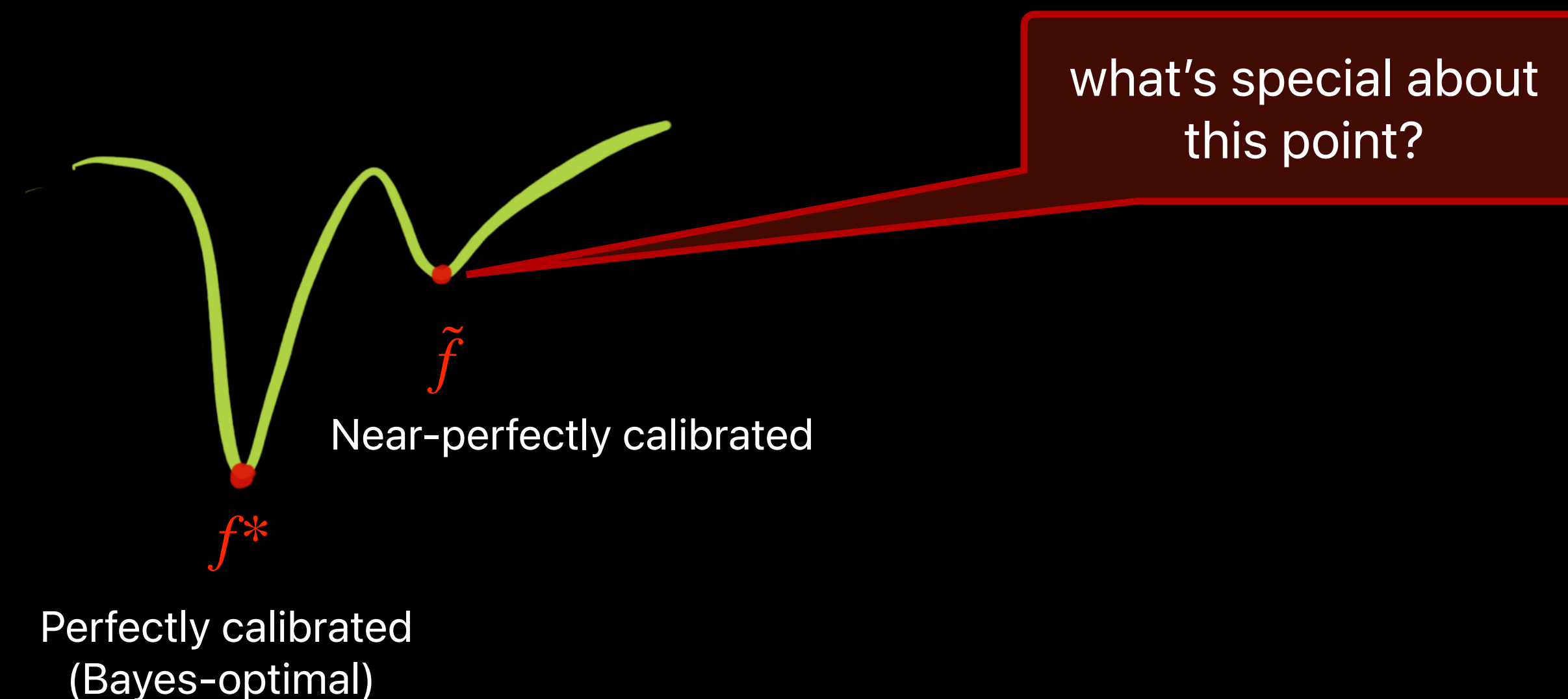
$$f^* = \operatorname{argmin}_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x,y \sim D} [\ell(f(x), y)]$$

$\Rightarrow f^*(x) = p_D(y|x)$
perfectly calibrated

What we actually do:

Run SGD* to *approximately minimize* expected loss, over *restricted family* $\{f_\theta : \theta \in \Theta\}$:

$$\tilde{f} = \operatorname{SGDmin}_{\theta \in \Theta} \mathbb{E}_{x,y \sim D} [\ell(f_\theta(x), y)]$$



In general, when does:

suboptimal loss-minimization \implies *near-optimal* calibration?

Toy Theorem

For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of $f : \mathcal{X} \rightarrow [0,1]$ on D cannot be improved by post-processing $\kappa : [0,1] \rightarrow [0,1]$

$$\forall \kappa : [0,1] \rightarrow [0,1], \quad \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

where $\mathcal{L}_D(f) := \mathbb{E}_{x,y \sim D}[\ell(f(x), y)]$ is the expected loss

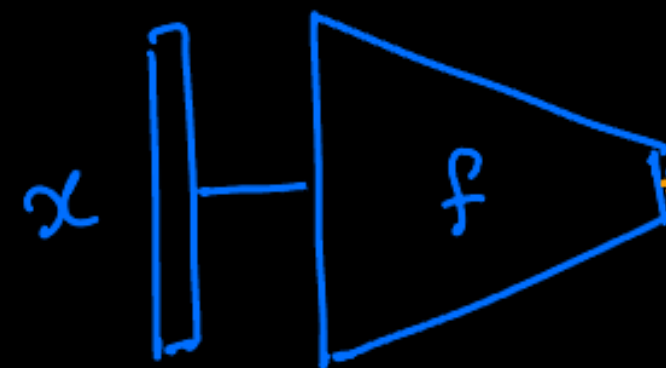
Toy Theorem

For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of $f : \mathcal{X} \rightarrow [0,1]$ on D cannot be improved by post-processing $\kappa : [0,1] \rightarrow [0,1]$

$$\forall \kappa : [0,1] \rightarrow [0,1], \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

where $\mathcal{L}_D(f) := \mathbb{E}_{x,y \sim D}[\ell(f(x), y)]$ is the expected loss



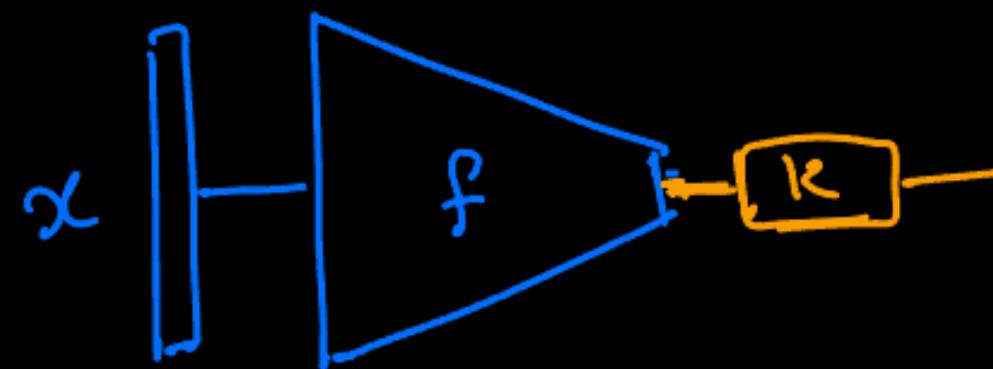
Toy Theorem

For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of $f : \mathcal{X} \rightarrow [0,1]$ on D cannot be improved by post-processing $\kappa : [0,1] \rightarrow [0,1]$

$$\forall \kappa : [0,1] \rightarrow [0,1], \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

where $\mathcal{L}_D(f) := \mathbb{E}_{x,y \sim D}[\ell(f(x), y)]$ is the expected loss



Toy Theorem

For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of f on D cannot be improved by post-processing:

$$\forall \kappa : \mathbb{R} \rightarrow \mathbb{R}, \quad \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

Suggestive properties:

1. Requires only “weak local-optimality”, not global optimality
2. Post-processing can be represented by adding a layer

Toy Theorem

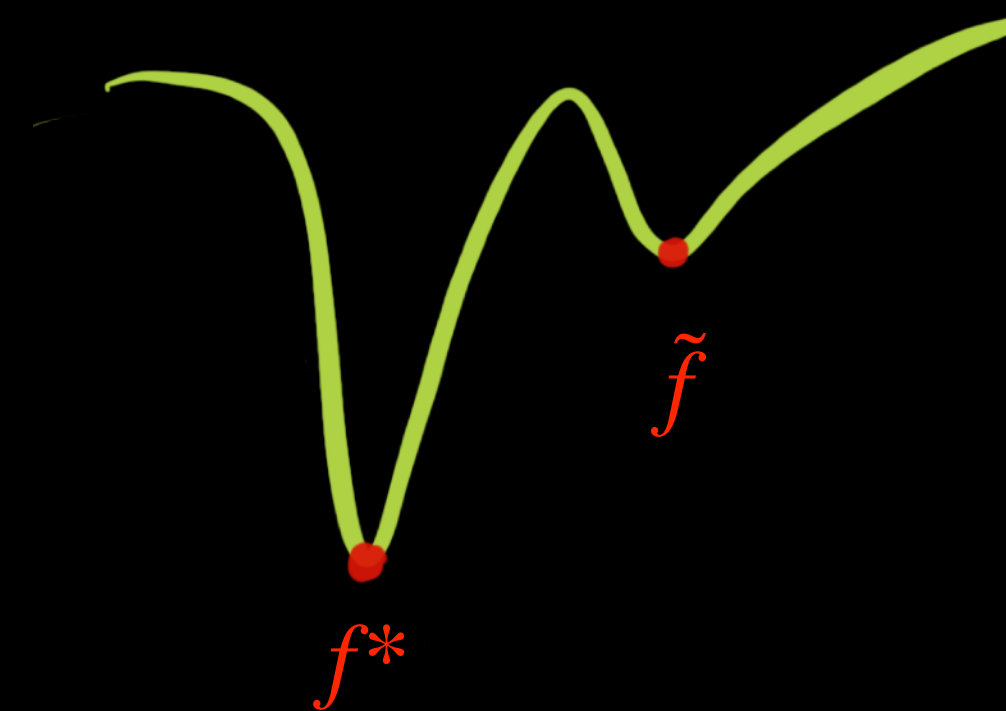
For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of f on D cannot be improved by post-processing:

$$\forall \kappa : \mathbb{R} \rightarrow \mathbb{R}, \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

Suggestive properties:

1. Requires only "weak local-optimality", not global optimality
2. Post-processing can be represented by adding a layer



Toy Theorem

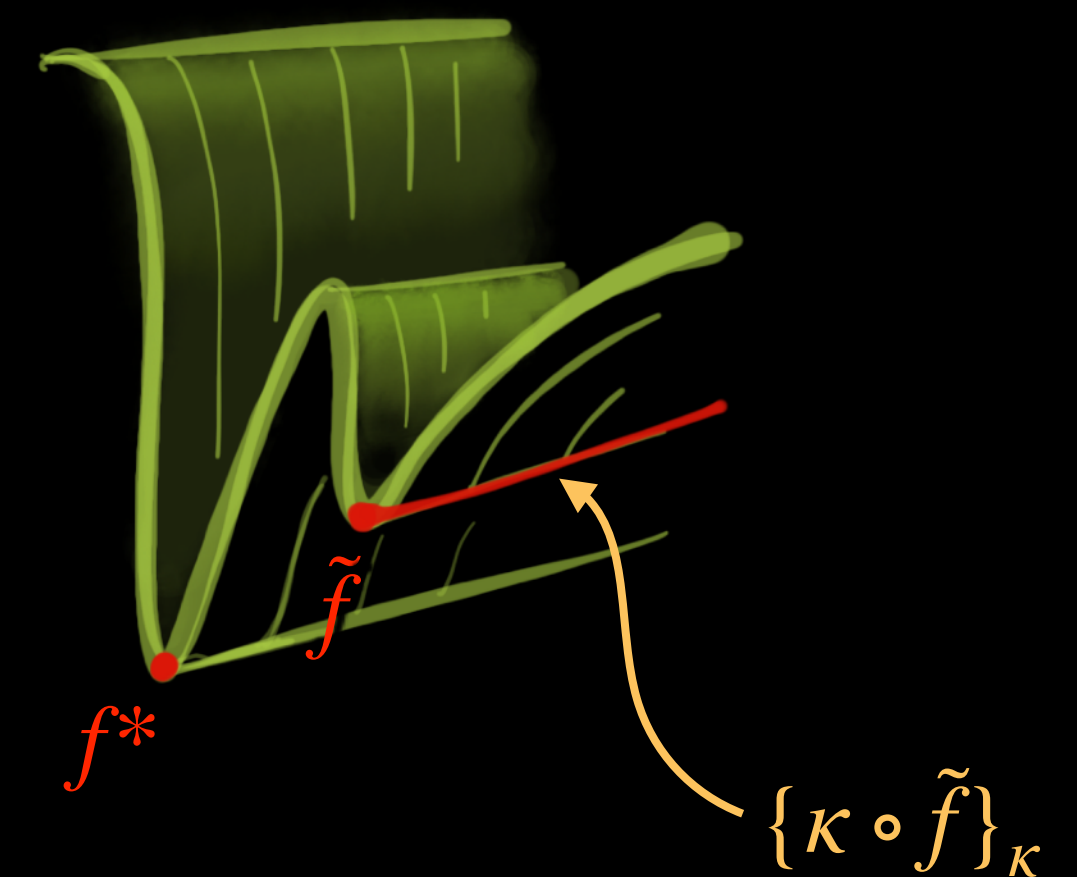
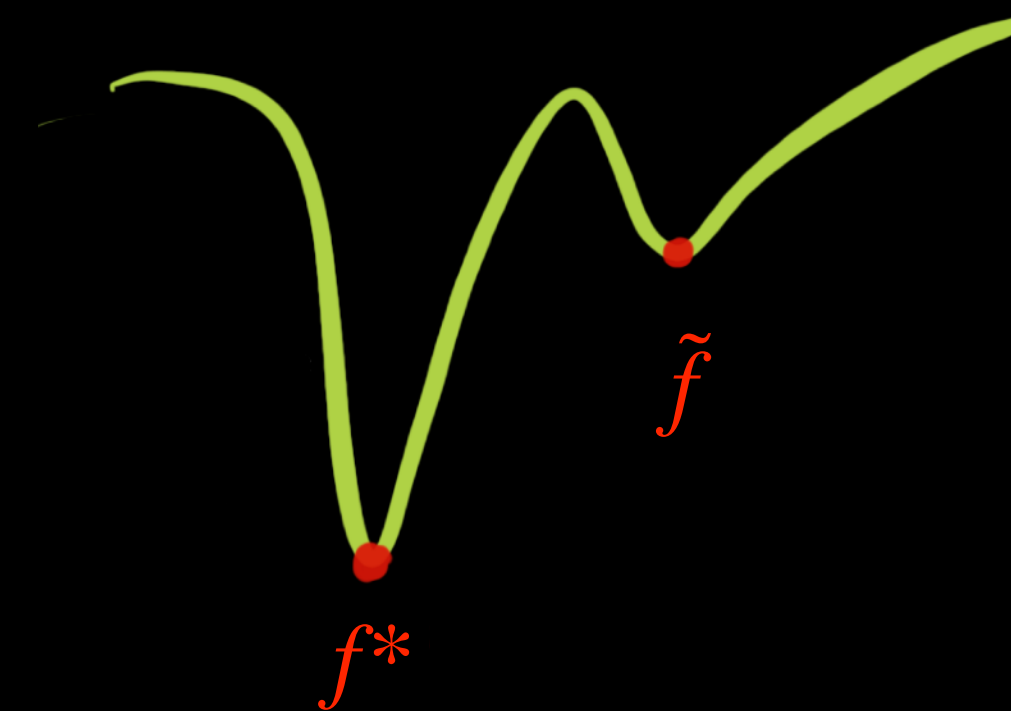
For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of f on D cannot be improved by post-processing:

$$\forall \kappa : \mathbb{R} \rightarrow \mathbb{R}, \quad \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

Suggestive properties:

1. Requires only "weak local-optimality", not global optimality
2. Post-processing can be represented by adding a layer



Toy Theorem

For all f, D , and proper loss ℓ , TFAE:

1. f is perfectly calibrated w.r.t. D
2. The loss of f on D cannot be improved by post-processing:

$$\forall \kappa : \mathbb{R} \rightarrow \mathbb{R}, \mathcal{L}_D(f) \leq \mathcal{L}_D(\kappa \circ f)$$

Problems:

1. Only characterizes perfect calibration
2. Requires composition with arbitrary functions
(not just "nice" ones that can be represented by NNs)

Toy Theorem

*" f is **perfectly calibrated** iff its loss can't be improved **at all** by post-processing **with an arbitrary function**"*

Toy Theorem

*" f is **perfectly calibrated** iff its loss can't be improved **at all** by post-processing **with an arbitrary function**"*

Dream Theorem

*" f is **close to calibrated** iff its loss can't be improved **much** by post-processing **with a smooth function**"*

Toy Theorem

*" f is **perfectly calibrated** iff its loss can't be improved **at all** by post-processing **with an arbitrary function**"*

Dream Theorem

*" f is **close to calibrated** iff its loss can't be improved **much** by post-processing **with a smooth function**"*

How to formalize **"close to"**? Calibration distance $dCE(f)$!

Toy Theorem

f is *perfectly calibrated* iff its $\text{loss}_{\text{[test]}}$ can't be improved *at all* by post-processing *with an arbitrary function*

Dream Theorem

f is *close to calibrated* iff its loss can't be improved *much* by post-processing *with a smooth function*

How to formalize "close to"? Calibration distance $d\text{CE}(f)$!

Dream Theorem

*" f is **close to** calibrated iff its loss can't be improved **much** by post-processing **with a smooth function**"*

Dream Theorem

*" f is **close to** calibrated iff its loss can't be improved **much** by post-processing **with a smooth function**"*

Theorem

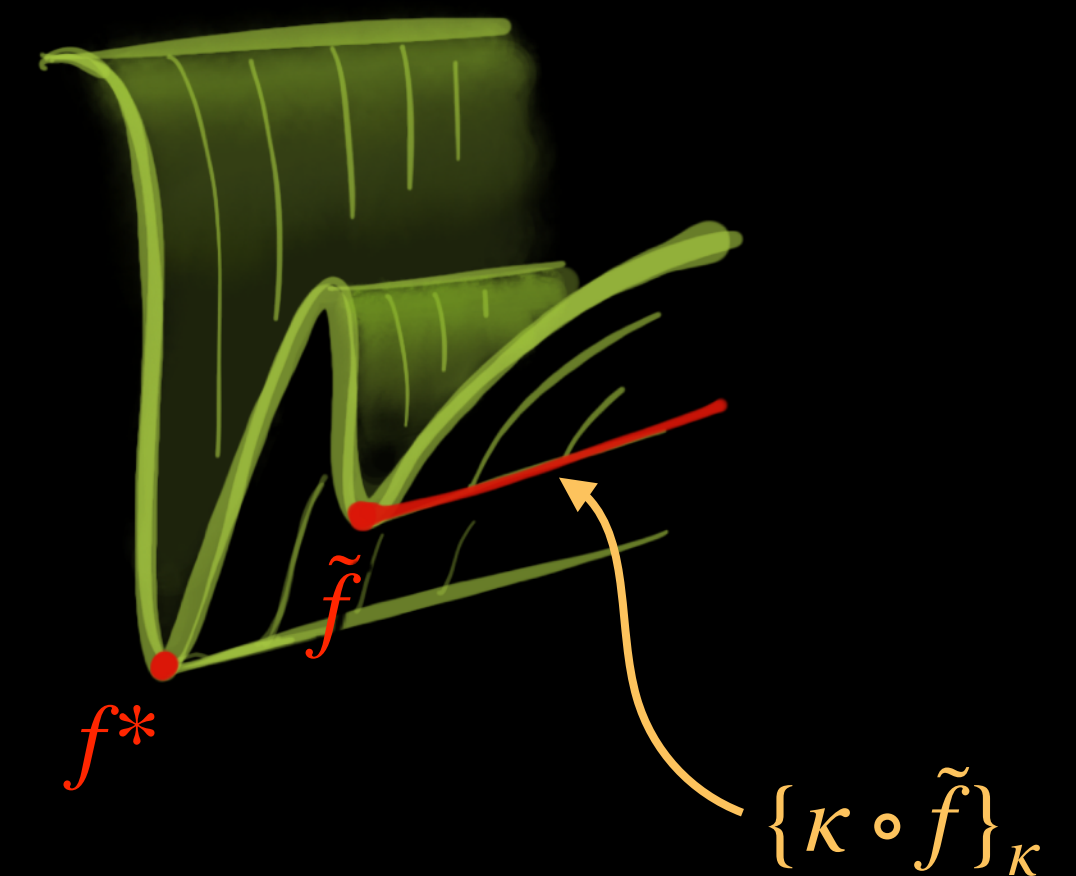
(distance from calibration) \sim poly(potential post-processing improvement)

Dream Theorem

*" f is **close to** calibrated iff its loss can't be improved **much** by post-processing **with a smooth function**"*

Theorem

(distance from calibration) \sim poly(potential post-processing improvement)



Dream Theorem

*" f is **close to** calibrated iff its loss can't be improved **much** by post-processing **with a smooth function**"*

Theorem

(distance from calibration) \sim poly(potential post-processing improvement)

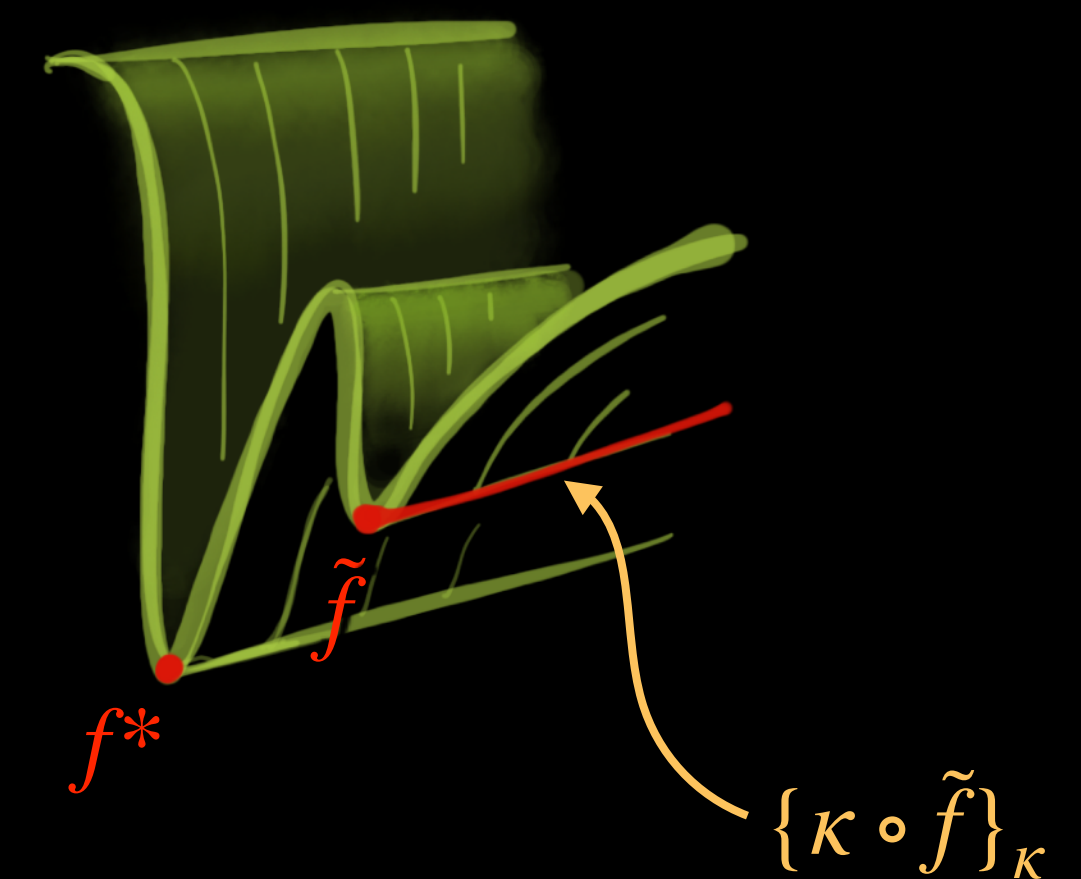
Theorem 1.3. *There exist constants $c_1, c_2 > 0$ such that for all predictors $f : \mathcal{X} \rightarrow [0, 1]$ and all distributions \mathcal{D} , the following holds.*

Let K denote the family of all post-processing functions $\kappa : [0, 1] \rightarrow [0, 1]$ such that the update function $\eta(f) = \kappa(f) - f$ is 1-Lipschitz. Define the "gap calibration error" of f as the maximum improvement in MSE loss via post-processings in K :

$$\text{gapCE}(f) = \text{MSE}_{\mathcal{D}}(f) - \min_{\kappa \in K} \text{MSE}_{\mathcal{D}}(\kappa \circ f).$$

Then, the maximum loss improvement (gapCE) polynomially bounds the distance from calibration (dCE):

$$c_1 \text{dCE}(f)^4 \leq \text{gapCE}(f) \leq c_2 \text{dCE}(f).$$

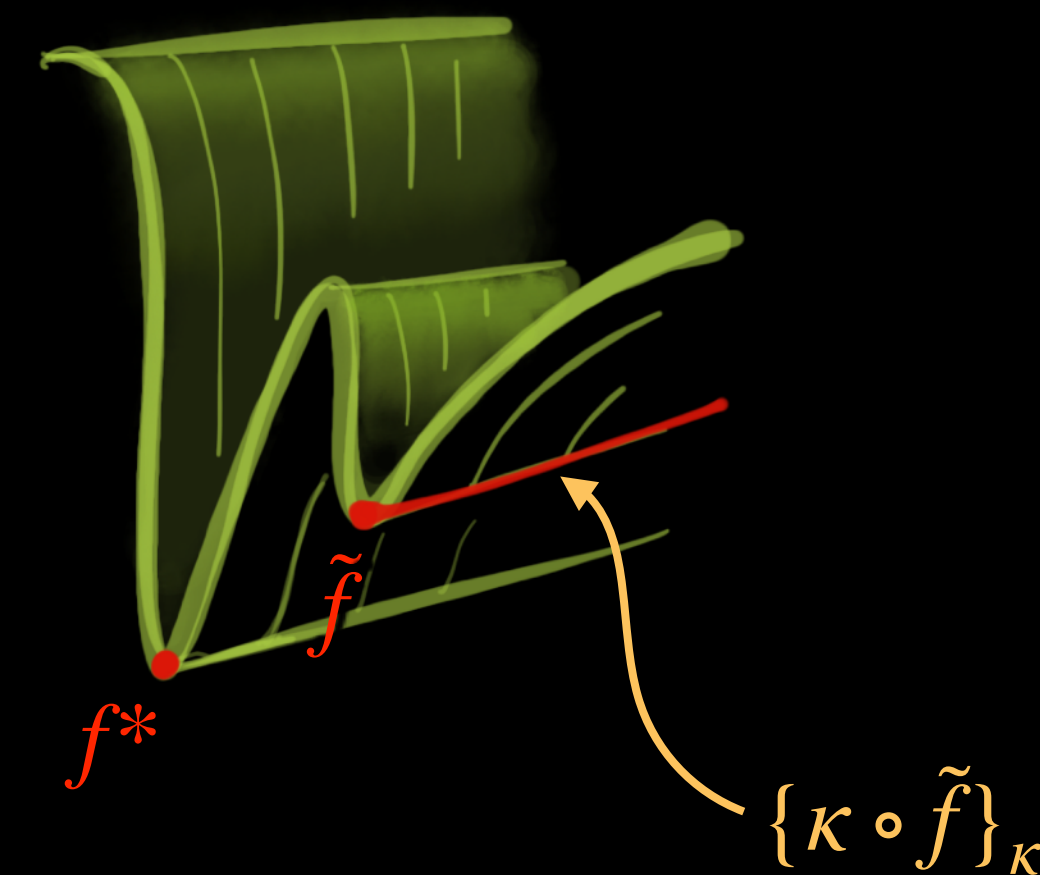
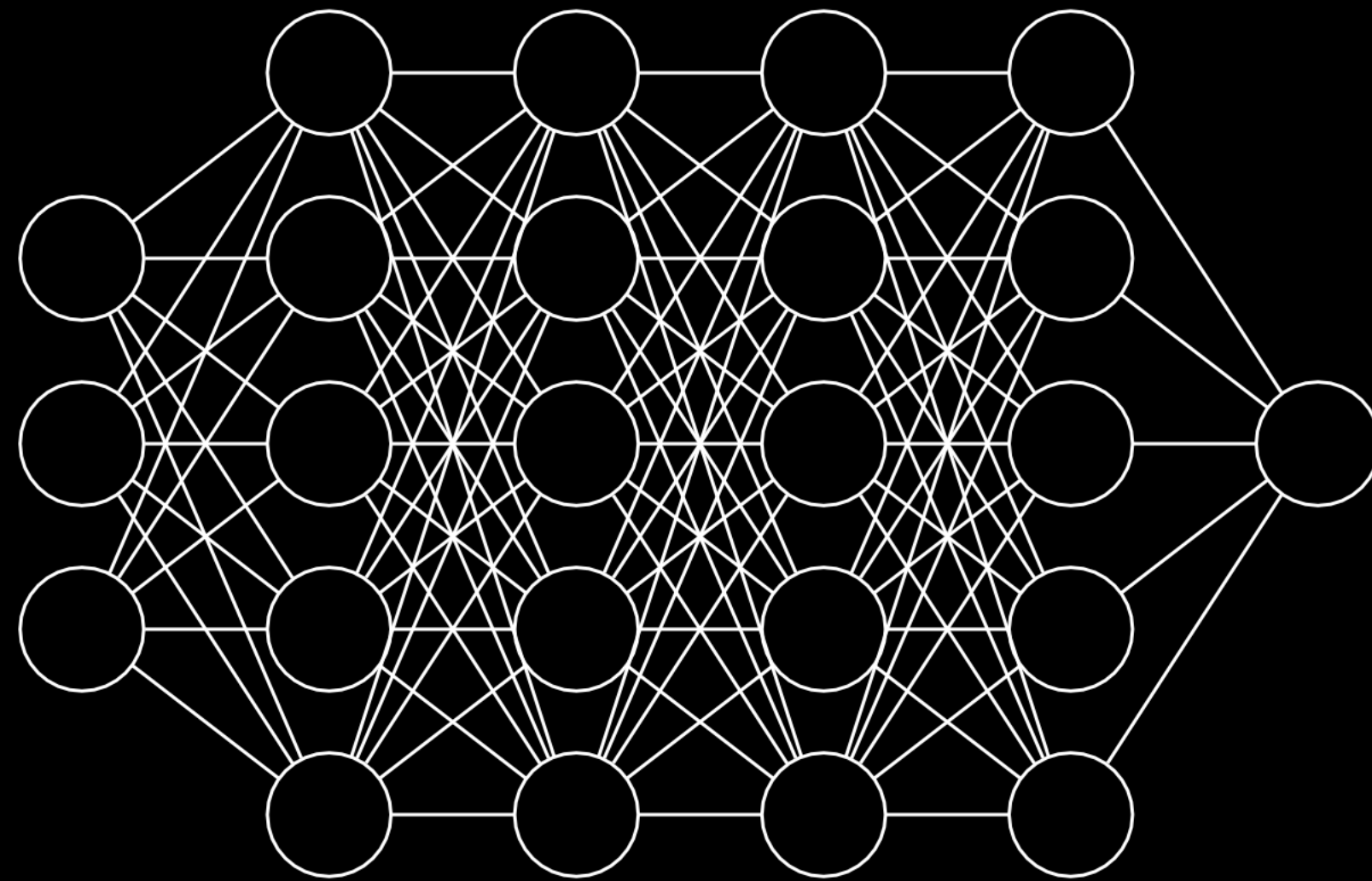


When is the “no loss improvement” condition satisfied?

1. (Algorithmic assumption):

If it were possible to improve loss via a simple post-processing, SGD would have done it already.

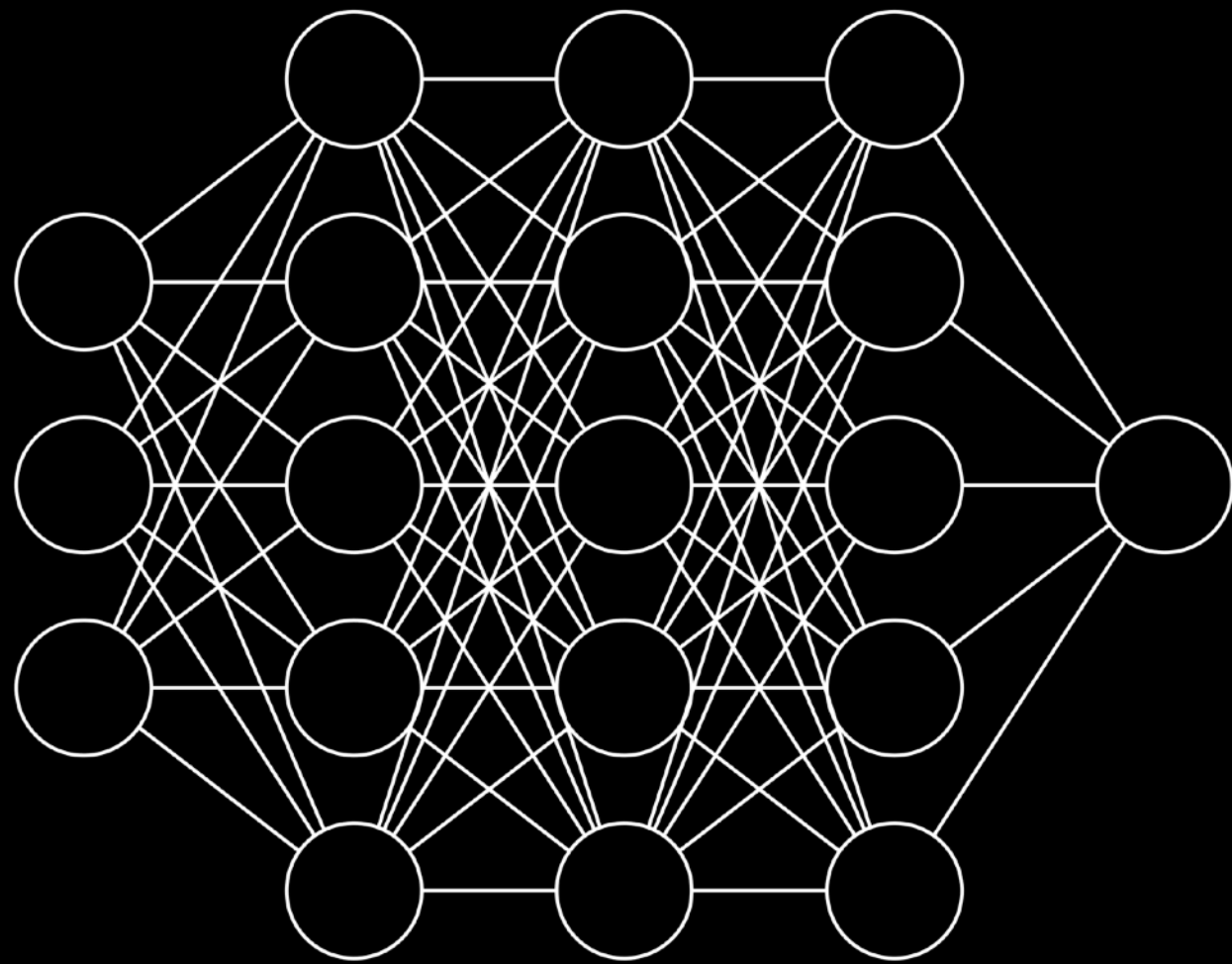
$f \mapsto \kappa \circ f$ is a “simple” update for SGD on deep nets



When is the “no loss improvement” condition satisfied?

2. **(Human-in-loop assumption):**

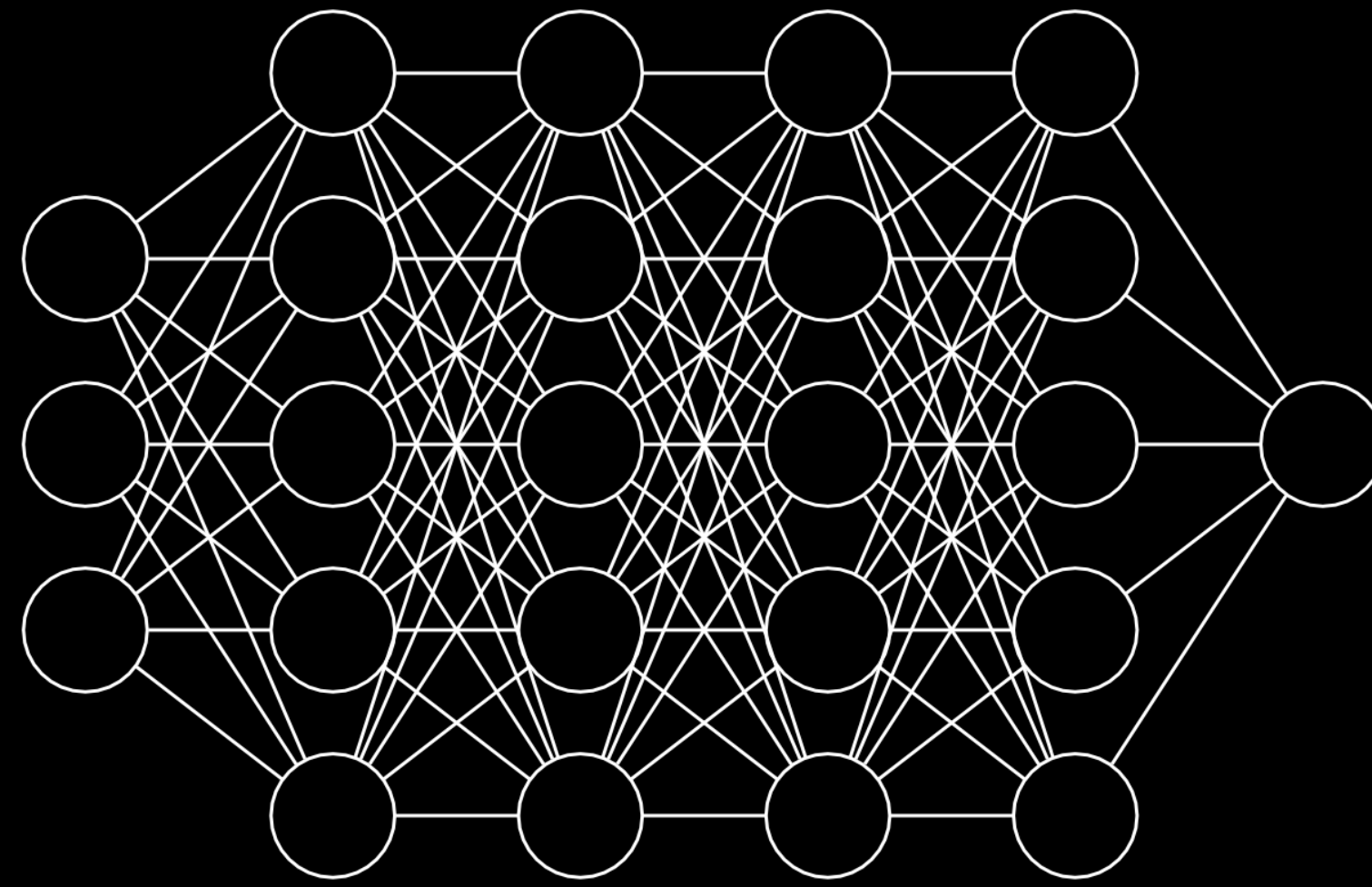
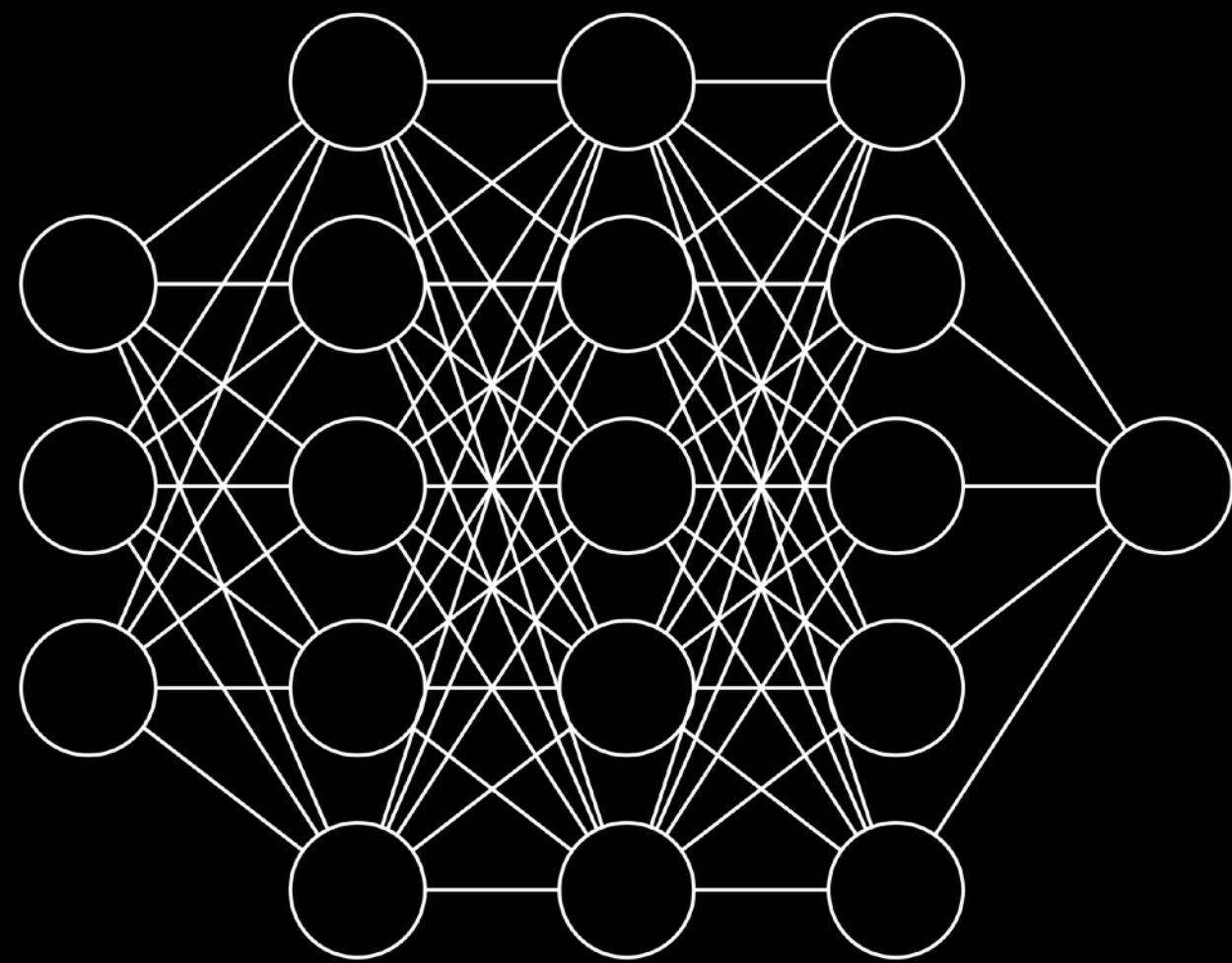
If it were possible to add a layer, train it optimally, and improve the loss, then the human trainer would have done it already



When is the “no loss improvement” condition satisfied?

2. **(Human-in-loop assumption):**

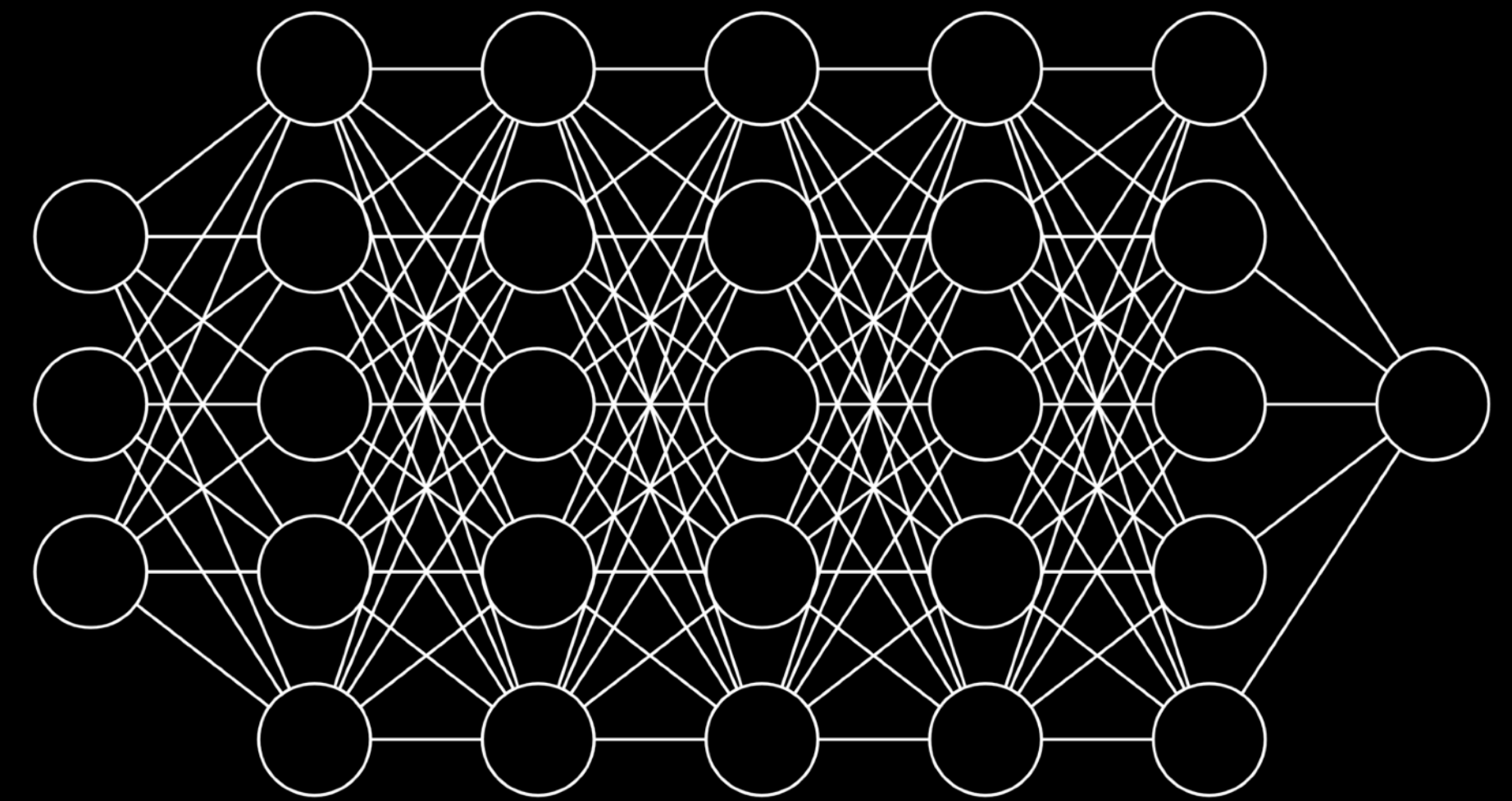
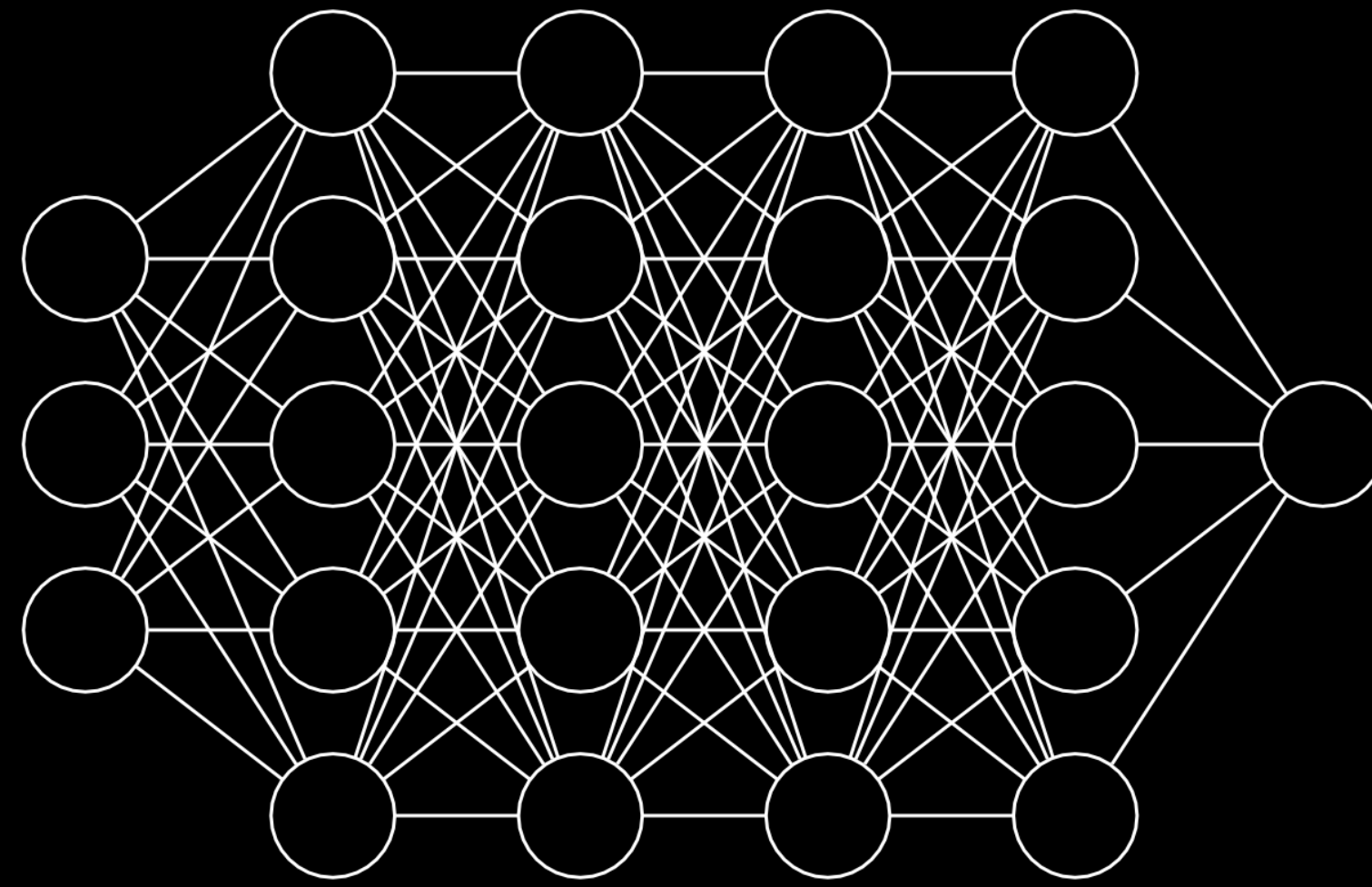
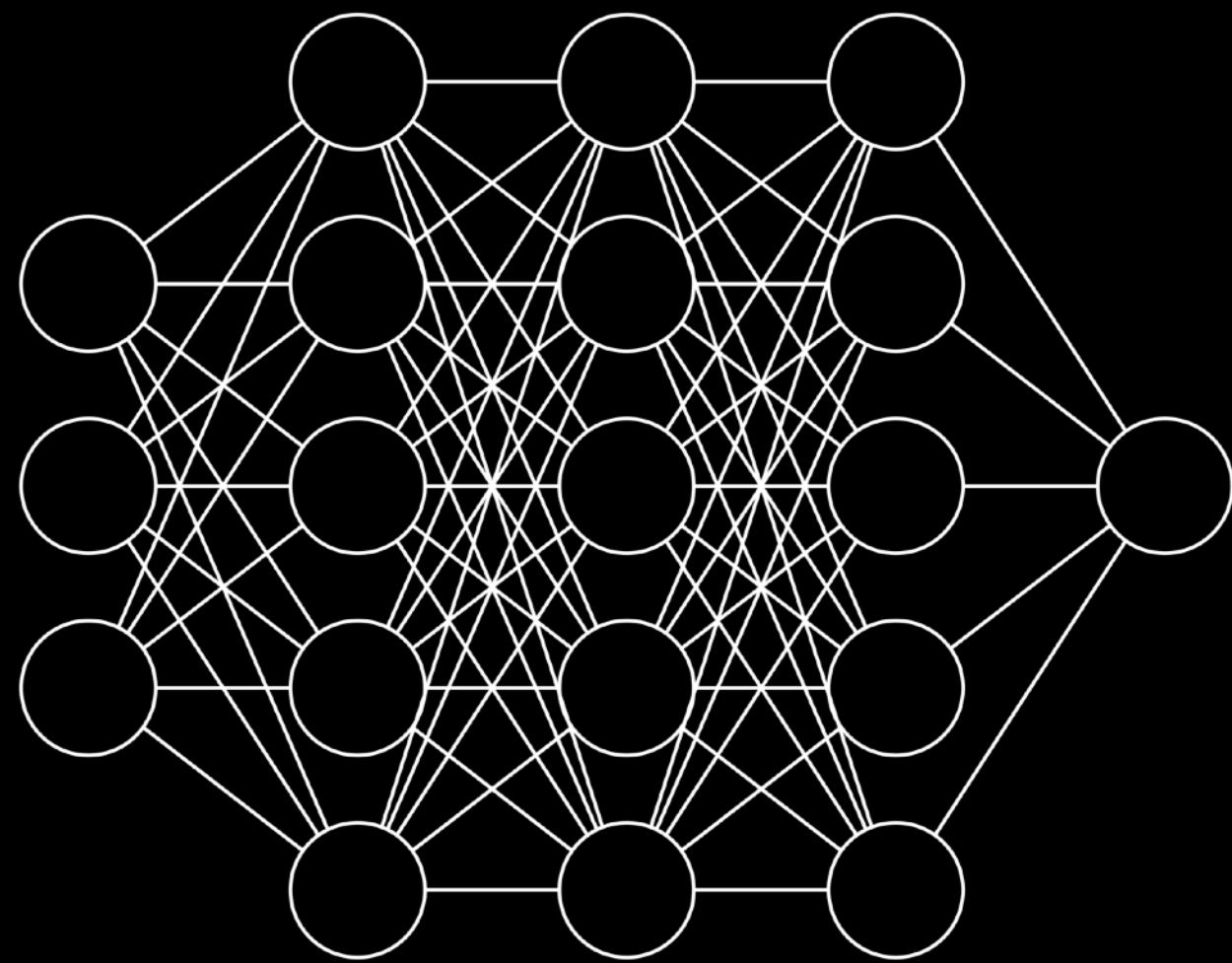
If it were possible to add a layer, train it optimally, and improve the loss, then the human trainer would have done it already



When is the “no loss improvement” condition satisfied?

2. (Human-in-loop assumption):

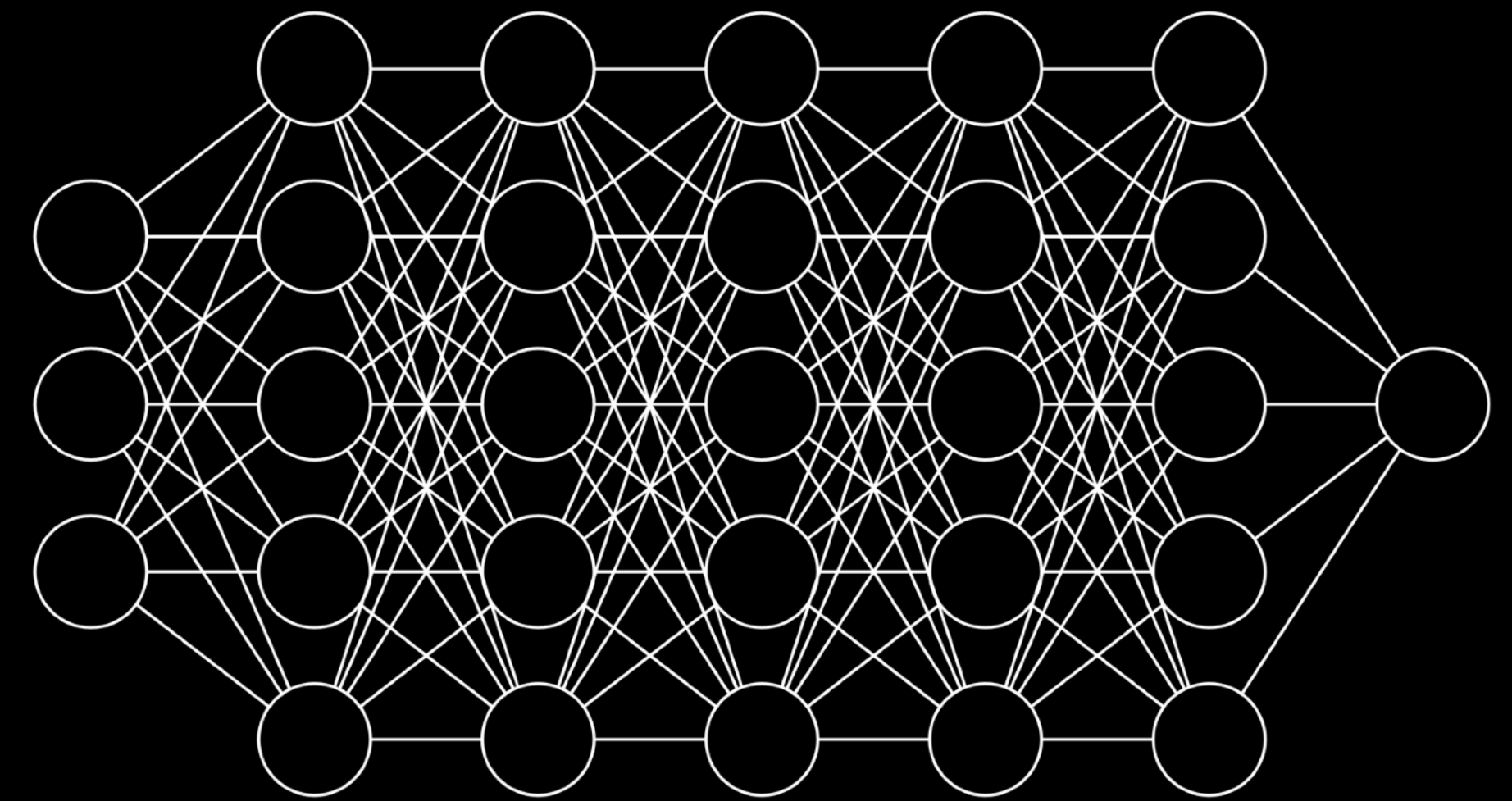
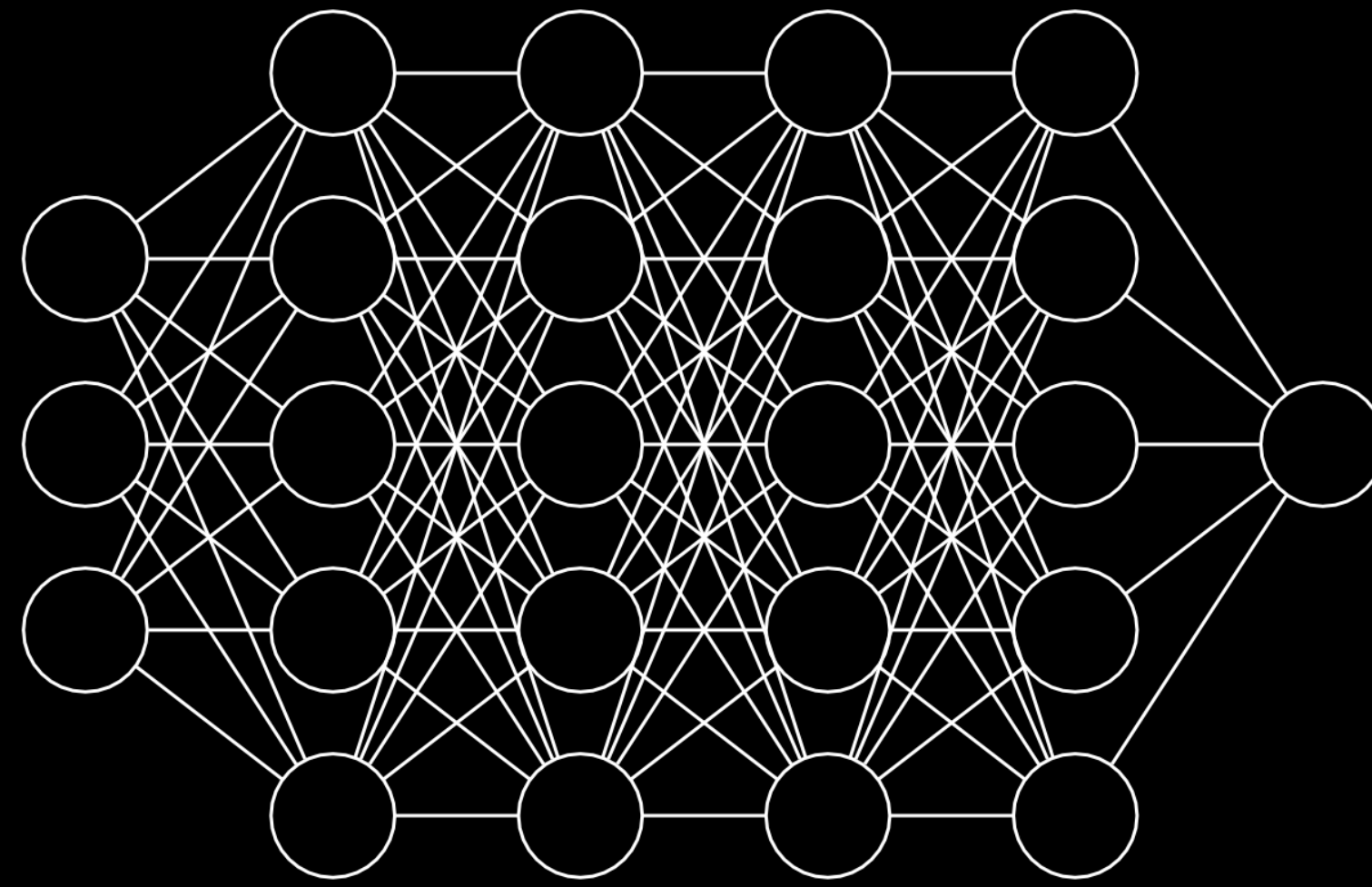
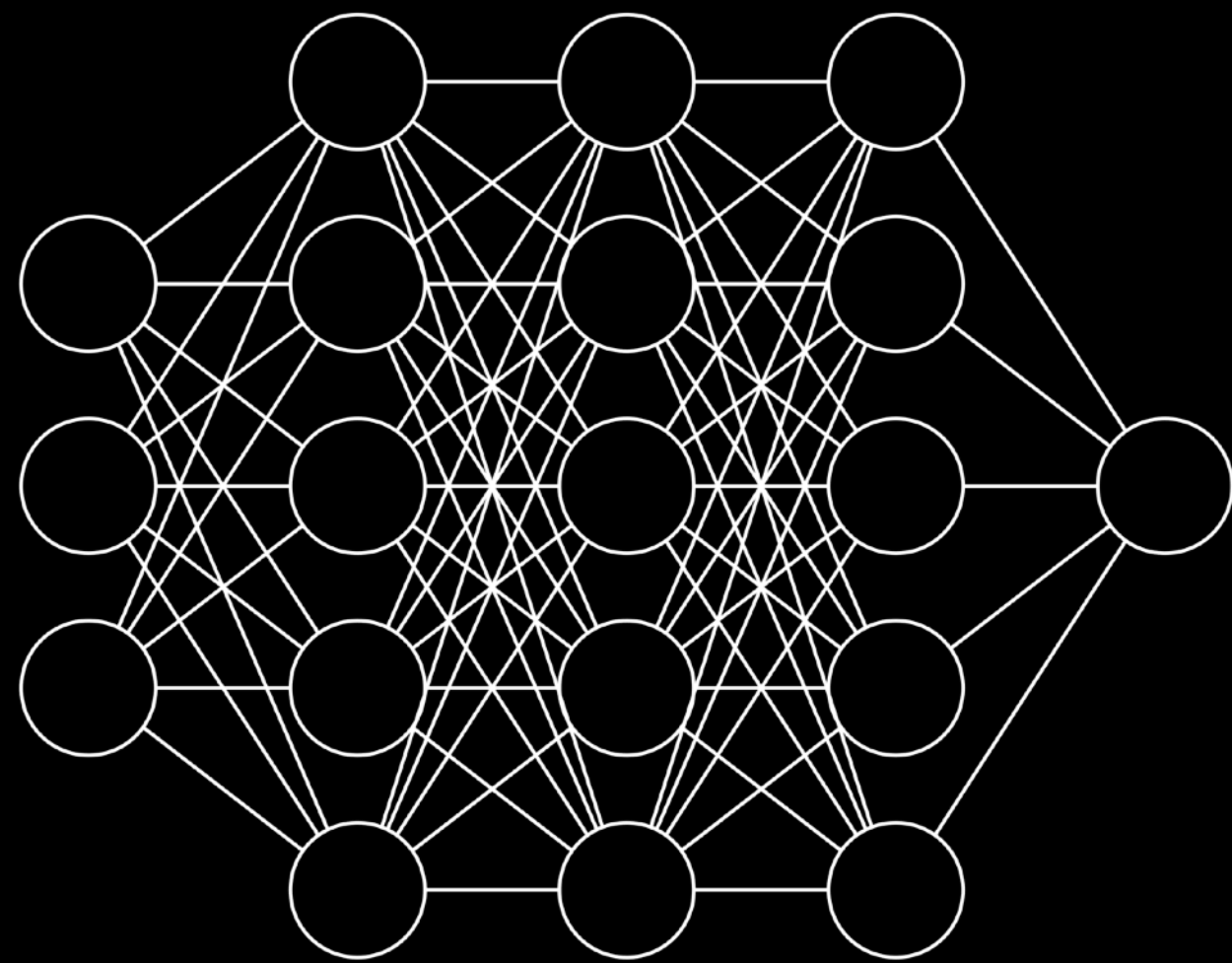
If it were possible to add a layer, train it optimally, and improve the loss, then the human trainer would have done it already



When is the “no loss improvement” condition satisfied?

2. (Human-in-loop assumption):

If it were possible to add a layer, train it optimally, and improve the loss, then the human trainer would have done it already \implies output of human is “nearly post-processing optimal”



When is the “no loss improvement” condition satisfied?

3. (Theory assumption):

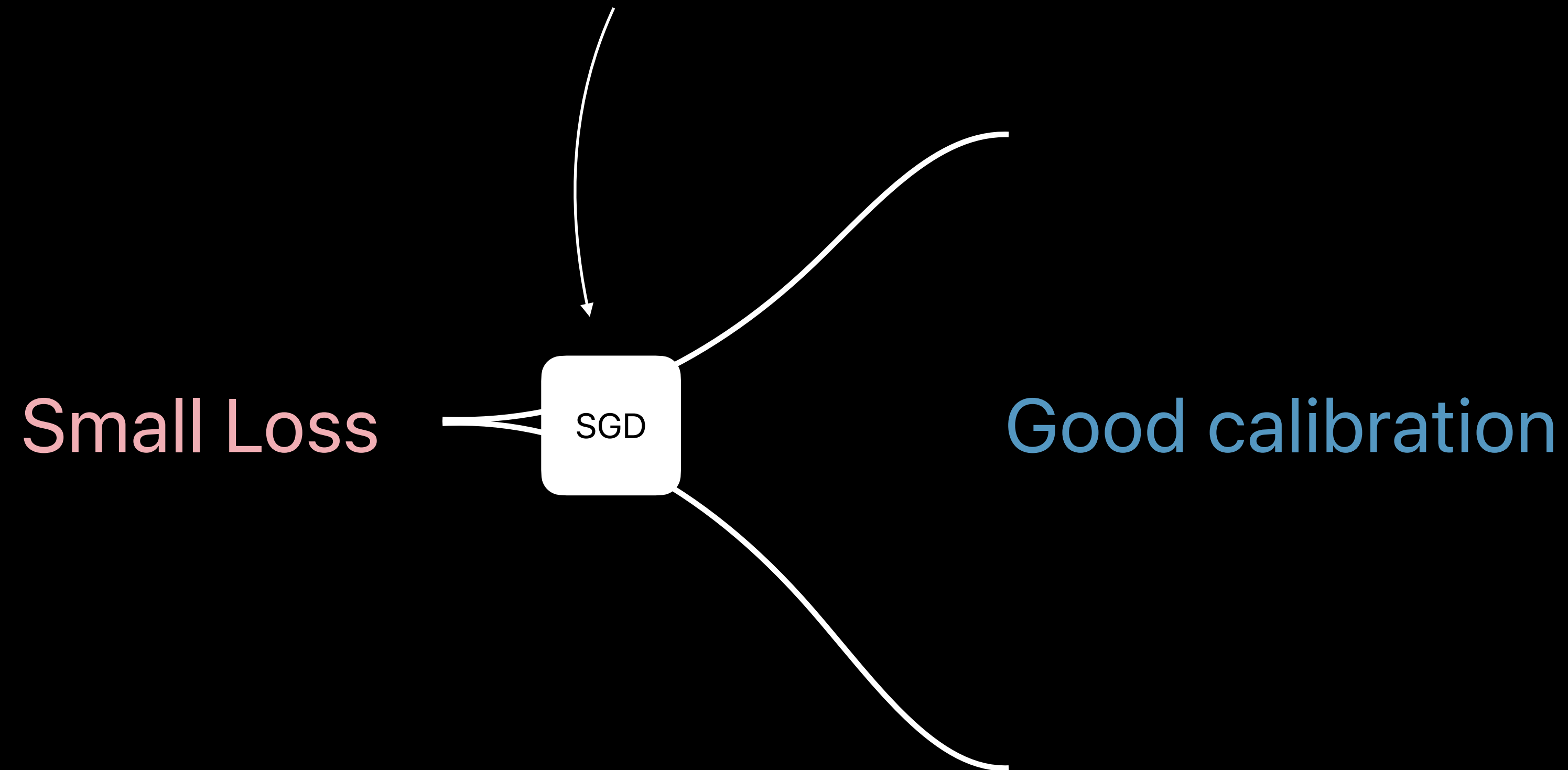
Structural risk minimization with any “well-behaved” complexity measure

$$\min_{f \in \mathcal{F}} \text{MSE}_{\mathcal{D}}(f) + \lambda \mu(f).$$

Implications

- Generic characterization of when (sub-optimal) **loss-minimization** yields (near-optimal) **calibration**
- Importance of **depth** for calibration
- Importance of **proper scoring rules** for calibration
- Non-Baysean reasons for calibration

Q: What's important about this box?



A: Output is (nearly) post-processing-optimal w.r.t. loss

Thanks!

In Collaboration With

Parikshit Gopalan
Apple

Vimal Thilak
Apple

Omid Saremi
Apple

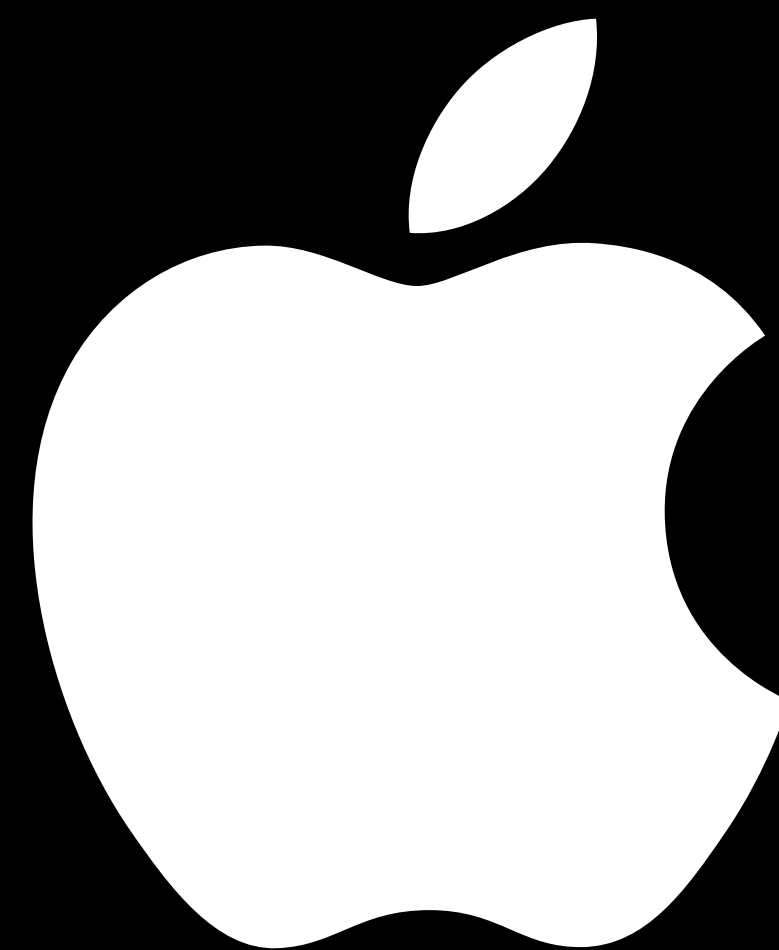
Joshua Suskind
Apple

Jarosław Błasiok
Columbia

Annabelle Carrell
Cambridge, Apple intern

Lunjia Hu
Stanford, Apple intern

Elan Rosenfeld
CMU, Apple intern



Defining “Almost All”

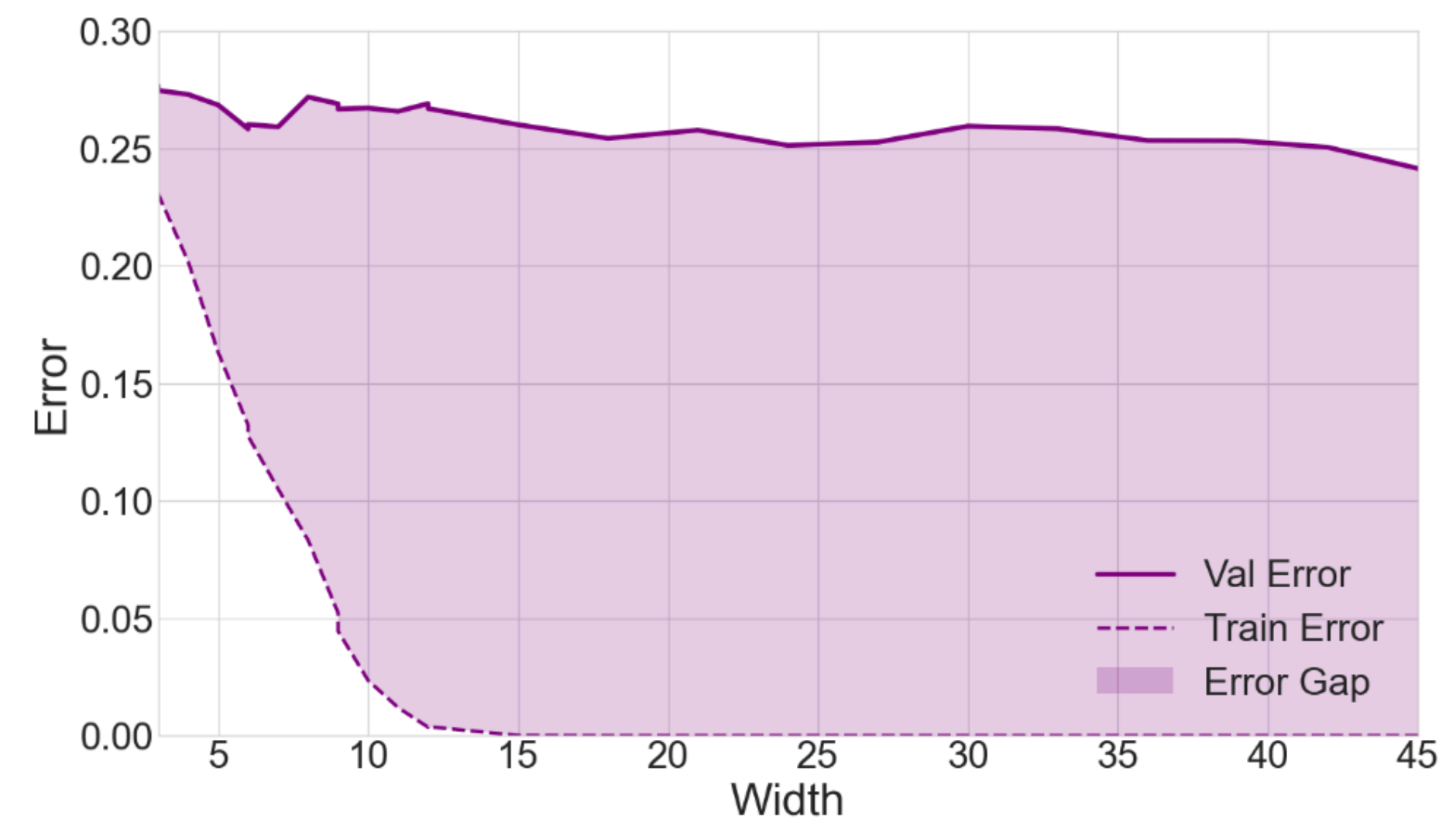
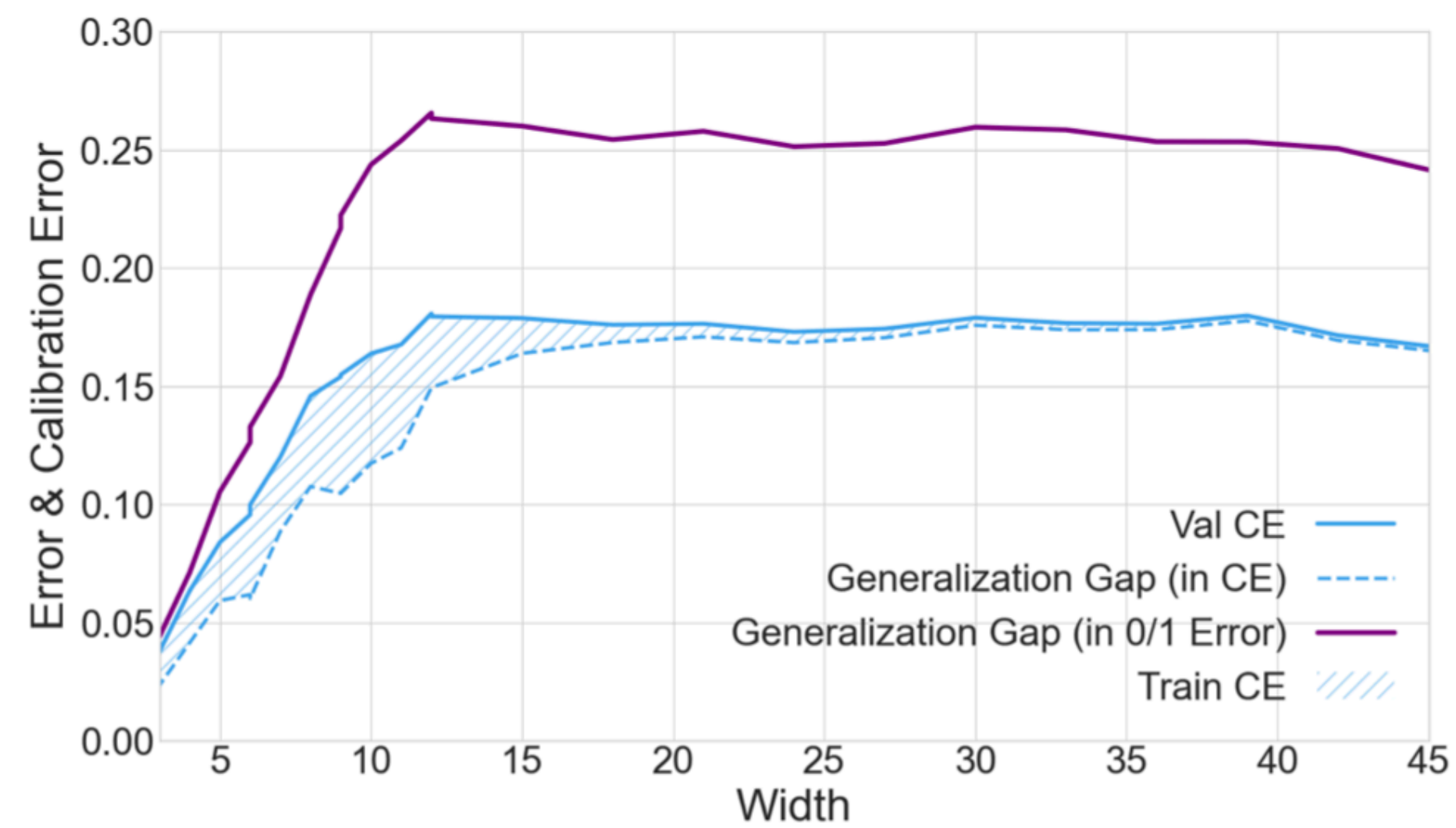
Definition-by-example:

Data distribution	Any
Architecture	Any* (MLP, ConvNet, Transformer,...)
Model depth	≥ 2
Model width	Any* (≥ 100)
Optimizer	Any* SGD-variant (SGD, Adam, ...)
Optimization steps	Any* (≥ 10 , after “warm-up” period)
Sample size	Any
Data-aug	None, or “standard” (measure-preserving)
Loss function	Any proper scoring rule (MSE, xent, ...)
Regularization	None, or very weak (e.g. wd=1-e4)

Empirical Claim 1:

For almost all* ML models

$$\mu_{\text{Train}} \approx 0$$

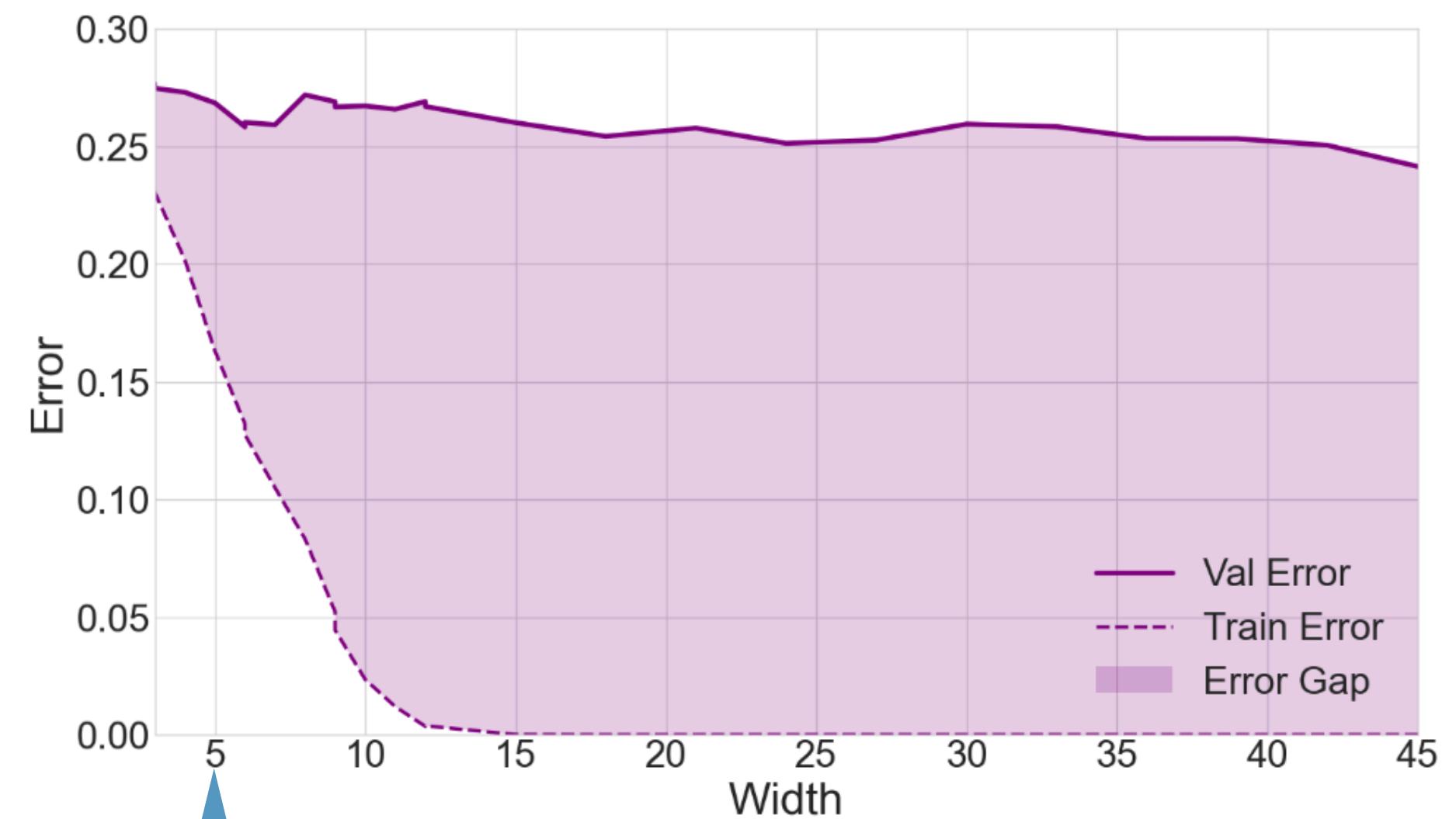
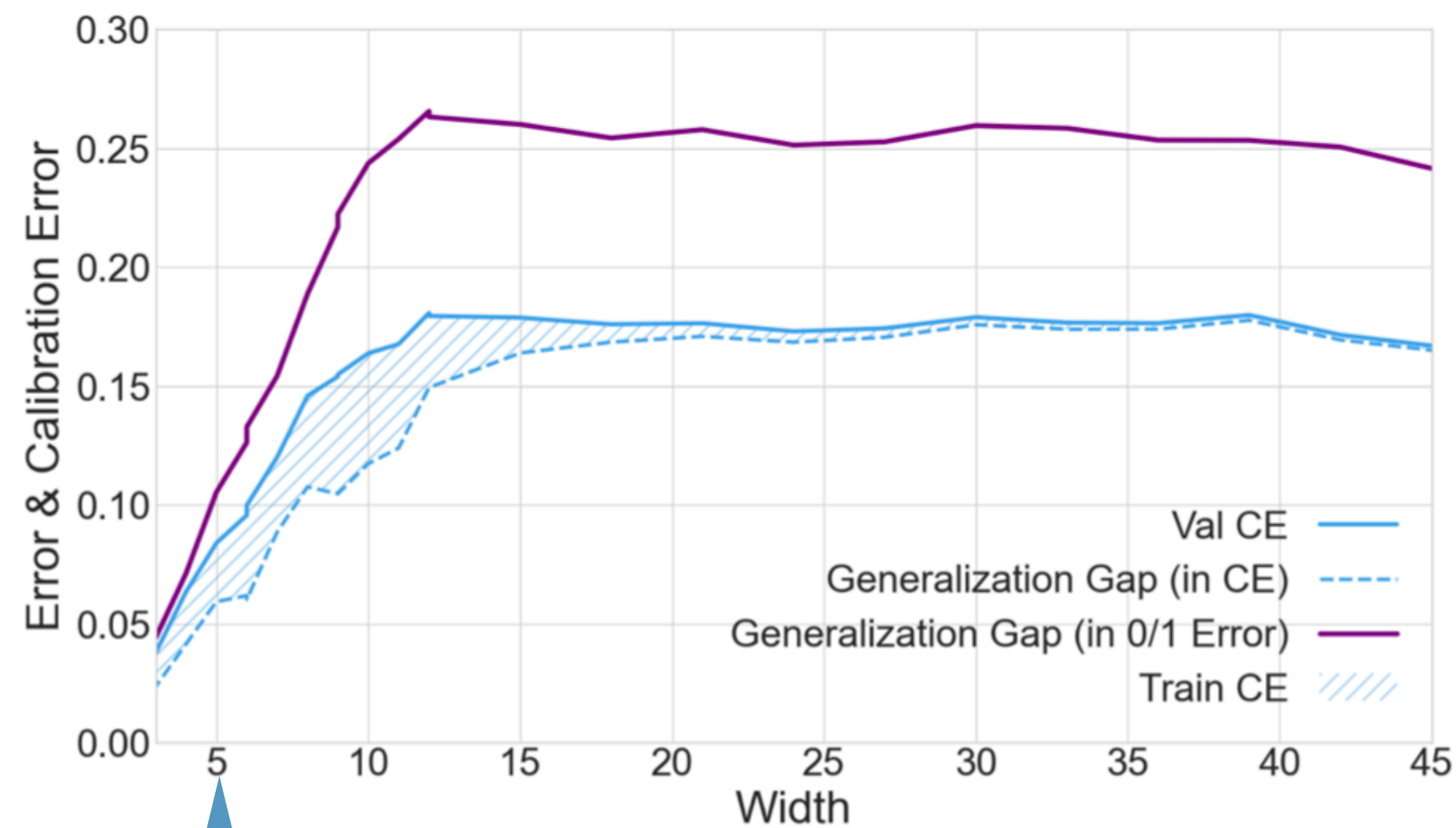


(ResNets on binary-CIFAR-10)

Empirical Claim 1:

For almost all* ML models

$$\mu_{\text{Train}} \approx 0$$



Surprising: small DNNs, with **high train error**, have good train calibration.

(ResNets on binary-CIFAR-10)

Part 1.

Measuring Miscalibration

Most models aren't *perfectly calibrated*.

How do we measure *degree-of-miscalibration*?

Summary: How to Measure Miscalibration [Błasiok, Gopalan, Hu, N. STOC 2023]

Most models aren't *perfectly calibrated*. How do we measure *degree-of-miscalibration*?

Desire:

Function $\mu_D(f) \in [0, \infty)$ that measures “degree of miscalibration”

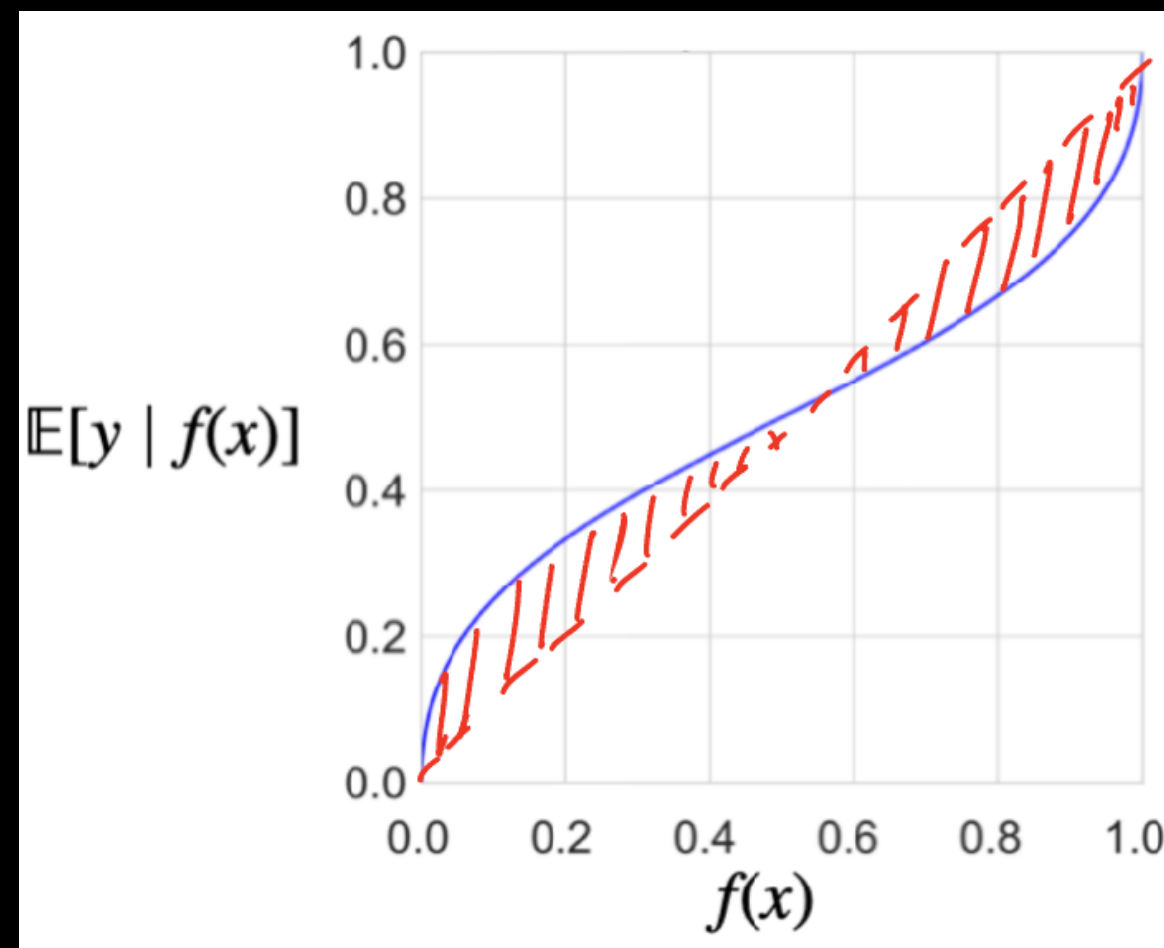
Summary: How to Measure Miscalibration [Błasiok, Gopalan, Hu, N. STOC 2023]

Most models aren't *perfectly calibrated*. How do we measure *degree-of-miscalibration*?

DON'T:

- Use “Expected Calibration Error (ECE)”

$$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y | f(x)] - f(x)|]$$



$\text{ECE}(f)$ is discontinuous in f !

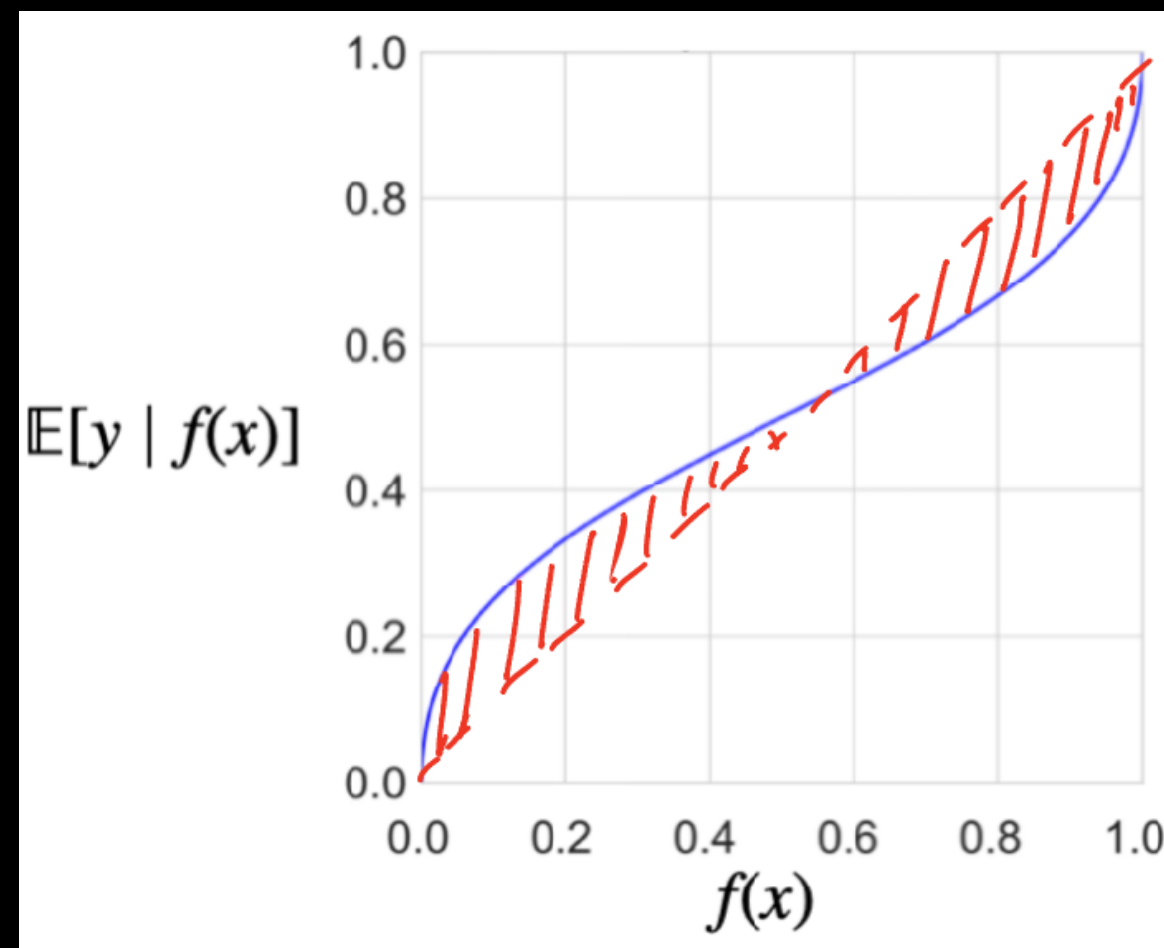
Summary: How to Measure Miscalibration [Błasiok, Gopalan, Hu, N. STOC 2023]

Most models aren't *perfectly calibrated*. How do we measure *degree-of-miscalibration*?

DON'T:

- Use "Expected Calibration Error (ECE)"

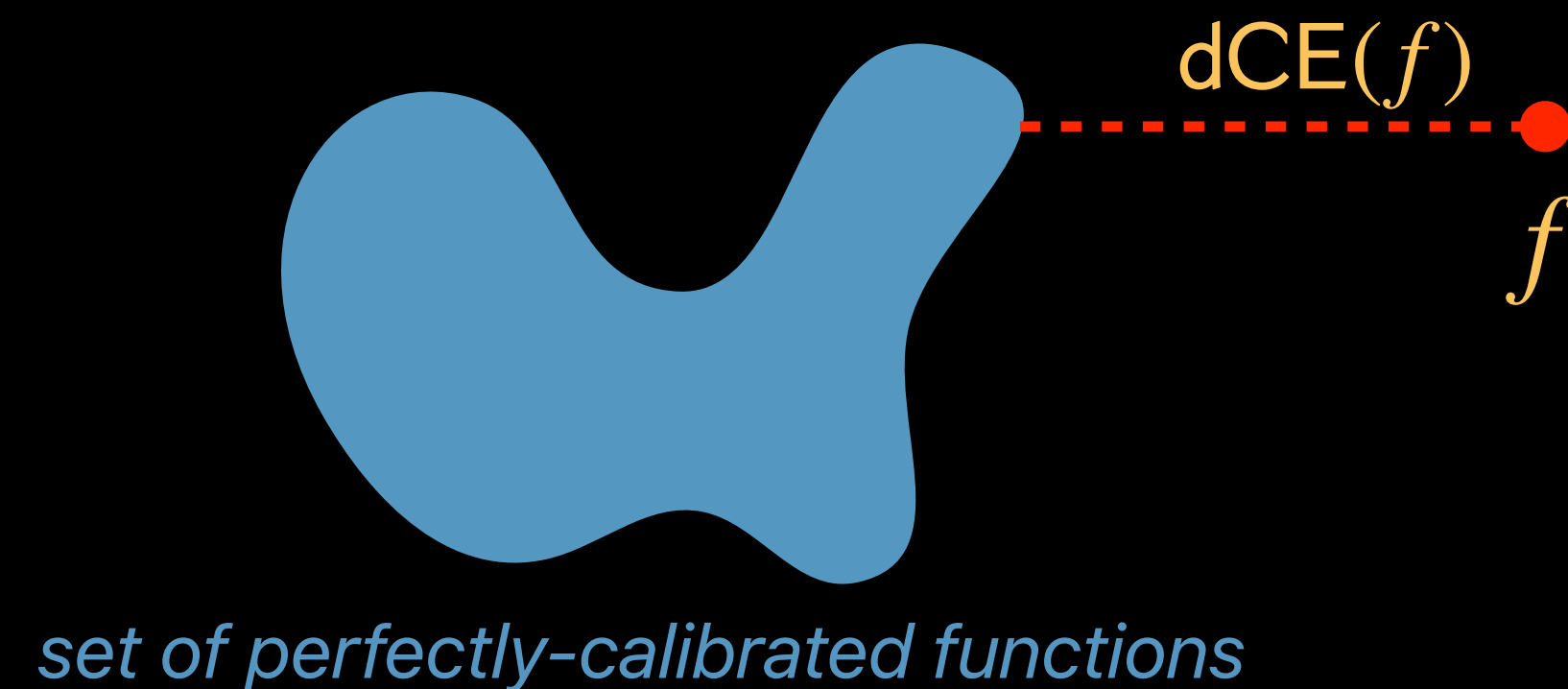
$$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y | f(x)] - f(x)|]$$



$\text{ECE}(f)$ is discontinuous in f !

DO:

- Use " ℓ_1 distance from perfect calibration"



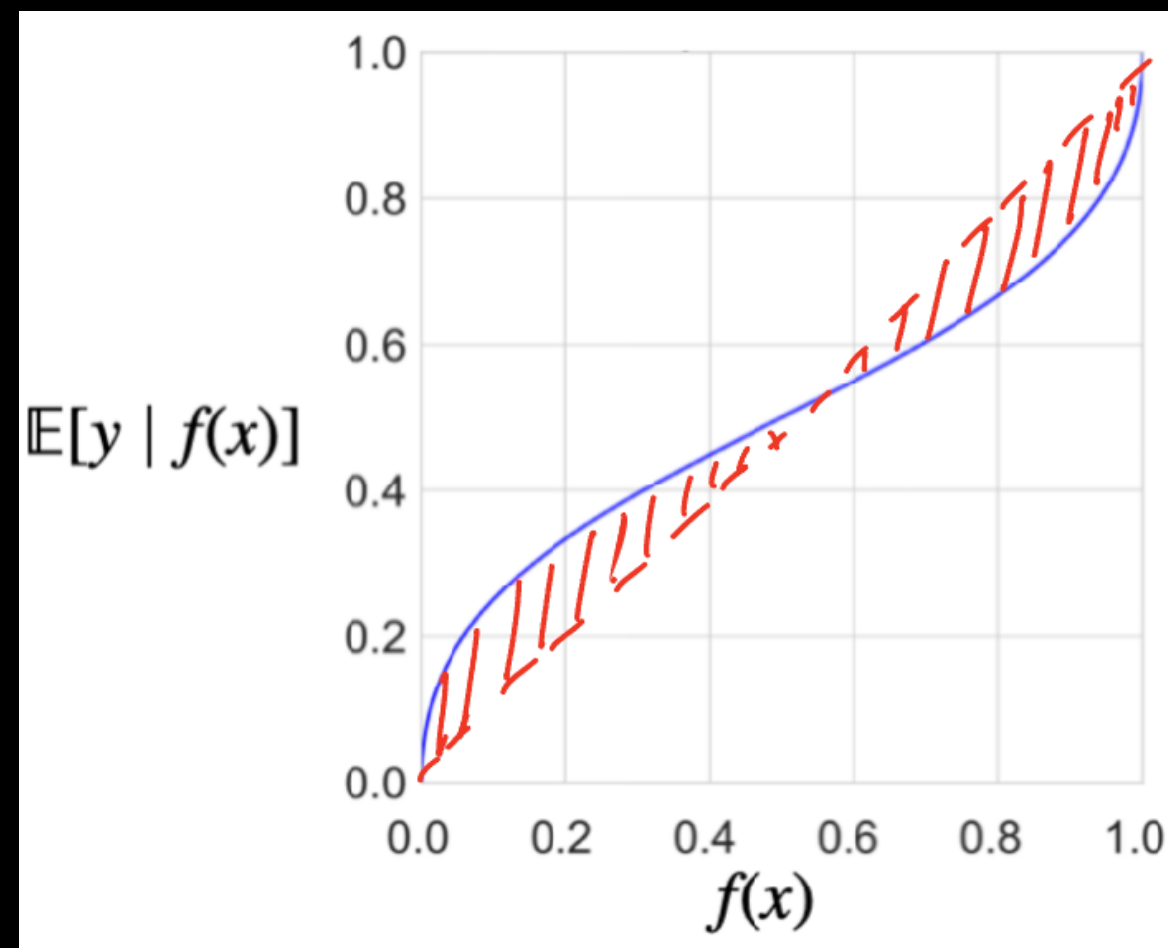
Summary: How to Measure Miscalibration [Błasiok, Gopalan, Hu, N. STOC 2023]

Most models aren't *perfectly calibrated*. How do we measure *degree-of-miscalibration*?

DON'T:

- Use "Expected Calibration Error (ECE)"

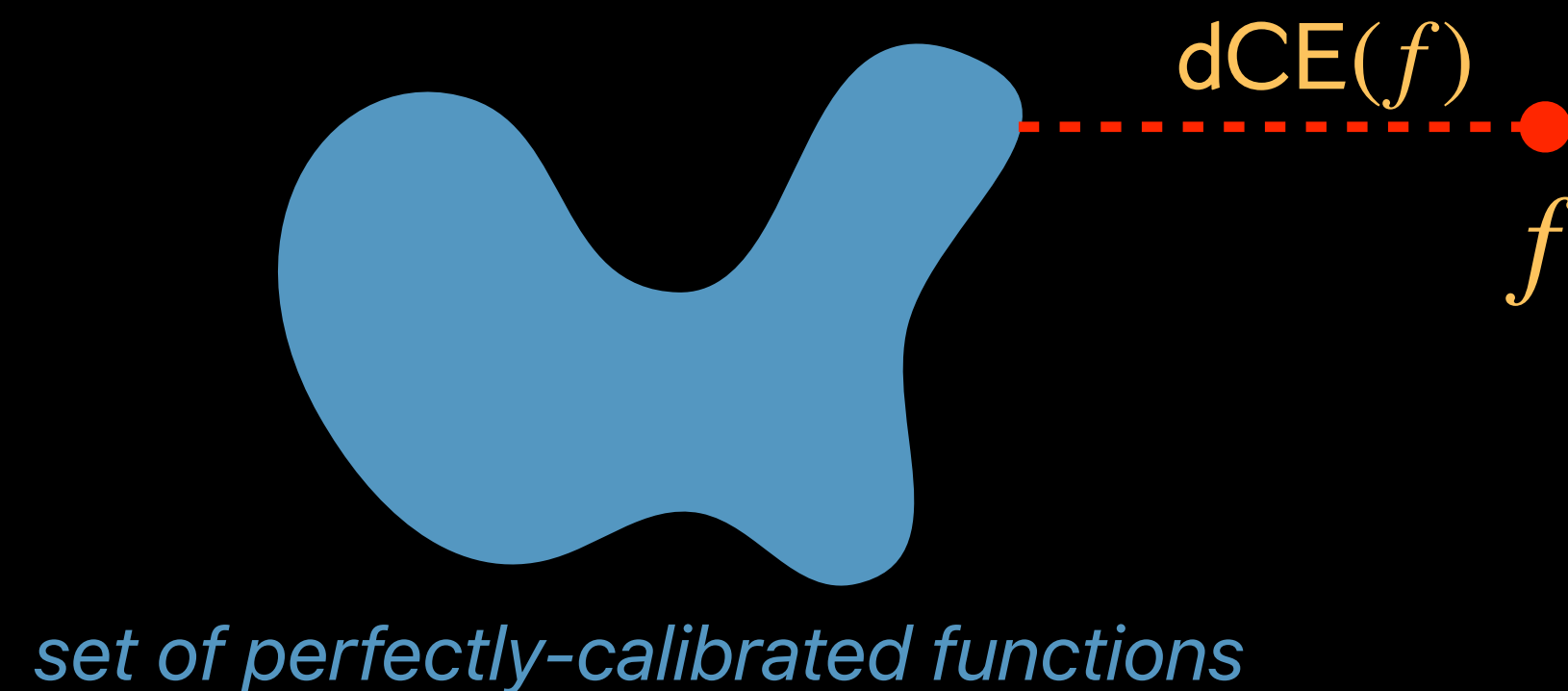
$$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y | f(x)] - f(x)|]$$



$\text{ECE}(f)$ is discontinuous in f !

DO:

- Use " ℓ_1 distance from perfect calibration"



- Estimate with a "consistent calibration metric" e.g. Kernel calibration error (kCE)

$$\text{kCE}_{\mathcal{D}}(f) := \sup_{w: \|w\|_K \leq 1} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$

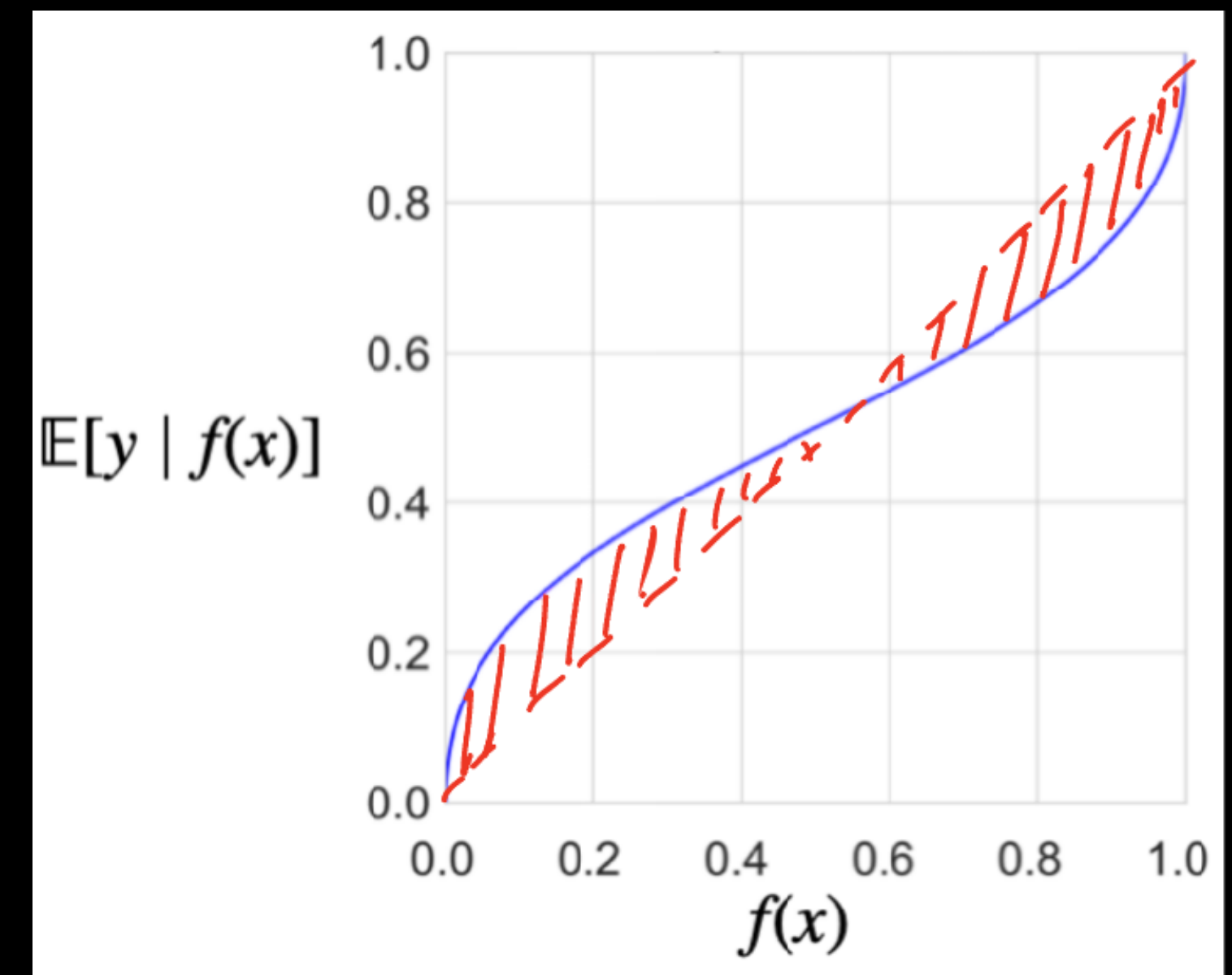
Measuring Miscalibration

Most models aren't *perfectly calibrated*.

How to measure **degree-of-miscalibration**?

Many proposed measures are problematic. Eg, ECE:

$$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y \mid f(x)] - f(x)|]$$





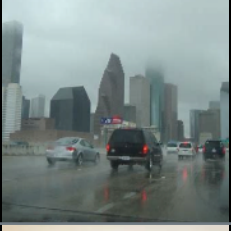



Problem: $\text{ECE}(f)$ is discontinuous in f

1. $\|f_1 - f_2\| \leq \varepsilon$

2. $\text{ECE}(f_1) - \text{ECE}(f_2) \geq 0.5 - \varepsilon$

Problem: $\text{ECE}(f)$ is discontinuous in f

- 1. $\|f_1 - f_2\| \leq \varepsilon$
- 2. $\text{ECE}(f_1) - \text{ECE}(f_2) \geq 0.5 - \varepsilon$



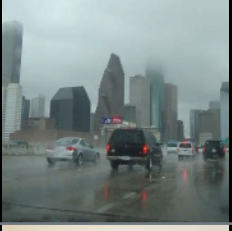



x	y	f ₁ (x)	f ₂ (x)
	1	0.5	0.5+ε
	1	0.5	0.5+ε
	1	0.5	0.5+ε
	0	0.5	0.5-ε
	0	0.5	0.5-ε
	0	0.5	0.5-ε

$\text{ECE}(f_1) = 0 \quad \text{ECE}(f_2) \approx 0.5$

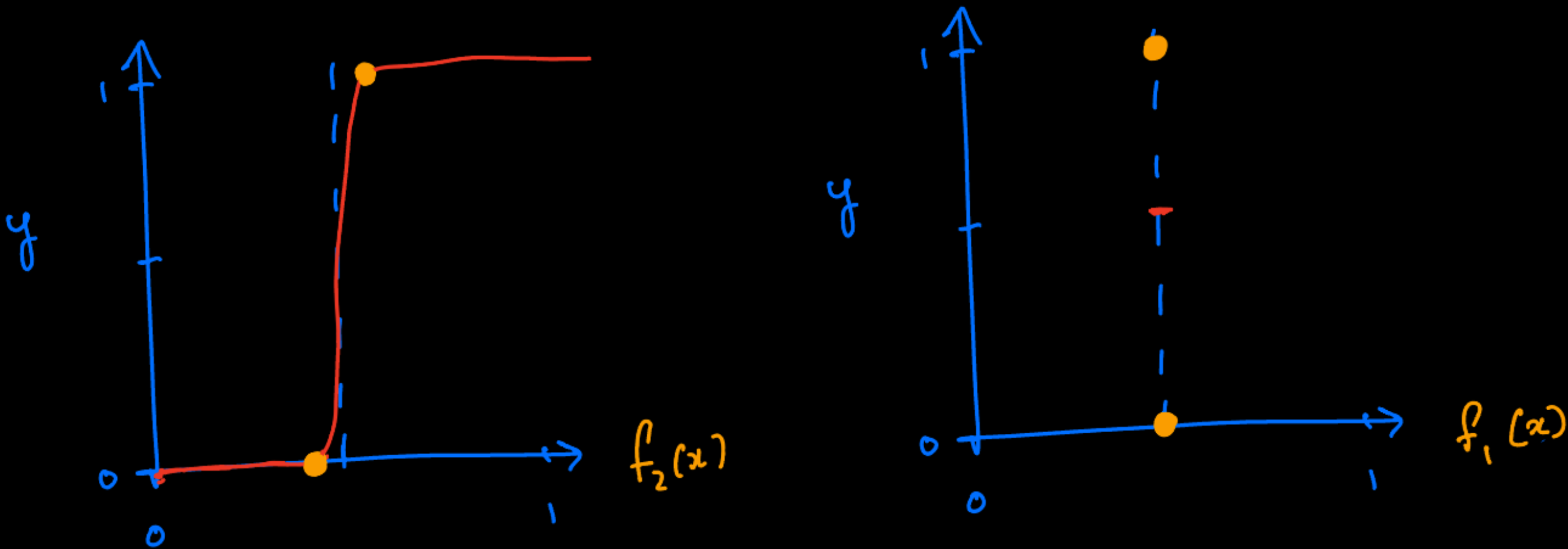
$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y \mid f(x)] - f(x)|]$

Problem: $\text{ECE}(f)$ is discontinuous in f

- 1. $\|f_1 - f_2\| \leq \varepsilon$
- 2. $\text{ECE}(f_1) - \text{ECE}(f_2) \geq 0.5 - \varepsilon$

x	y	f ₁ (x)	f ₂ (x)
	1	0.5	0.5+ε
	1	0.5	0.5+ε
	1	0.5	0.5+ε
	0	0.5	0.5-ε
	0	0.5	0.5-ε
	0	0.5	0.5-ε

$\text{ECE}(f_1) = 0$ $\text{ECE}(f_2) \approx 0.5$



$\text{ECE}(f) = \mathbb{E}[|\mathbb{E}[y \mid f(x)] - f(x)|]$

Axiomatic construction of **degree-of-miscalibration** $\mu_D(f)$?

Want $\mu(f) \in \mathbb{R}_{\geq 0}$ to satisfy:

1. Correctness:

$\mu(f) = 0 \iff f$ is perfectly calibrated

2. $\mu(f)$ is continuous in f

3. Can be estimated from samples

Axiomatic construction of **degree-of-miscalibration** $\mu_D(f)$?

Want $\mu(f) \in \mathbb{R}_{\geq 0}$ to satisfy:

1. Correctness:

$\mu(f) = 0 \iff f$ is perfectly calibrated

2. $\mu(f)$ is continuous in f

3. Can be estimated from samples

	Correctness	Continuity	Estimation
ECE	✓	✗	✗
Binned-ECE	✗	✗	✓
Brier	✗	✓	✓
NLL	✗	✓	✓
NCE	✗	✓	✓
kCE/MMCE	✓	✓	✓
smCE	✓	✓	✓

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

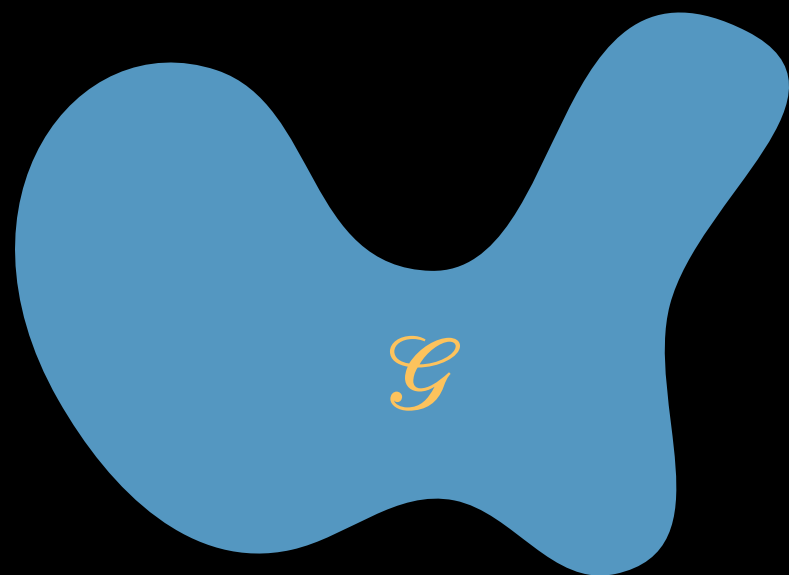
$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$



$\mathcal{G} :=$ set of perfectly-calibrated functions

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$



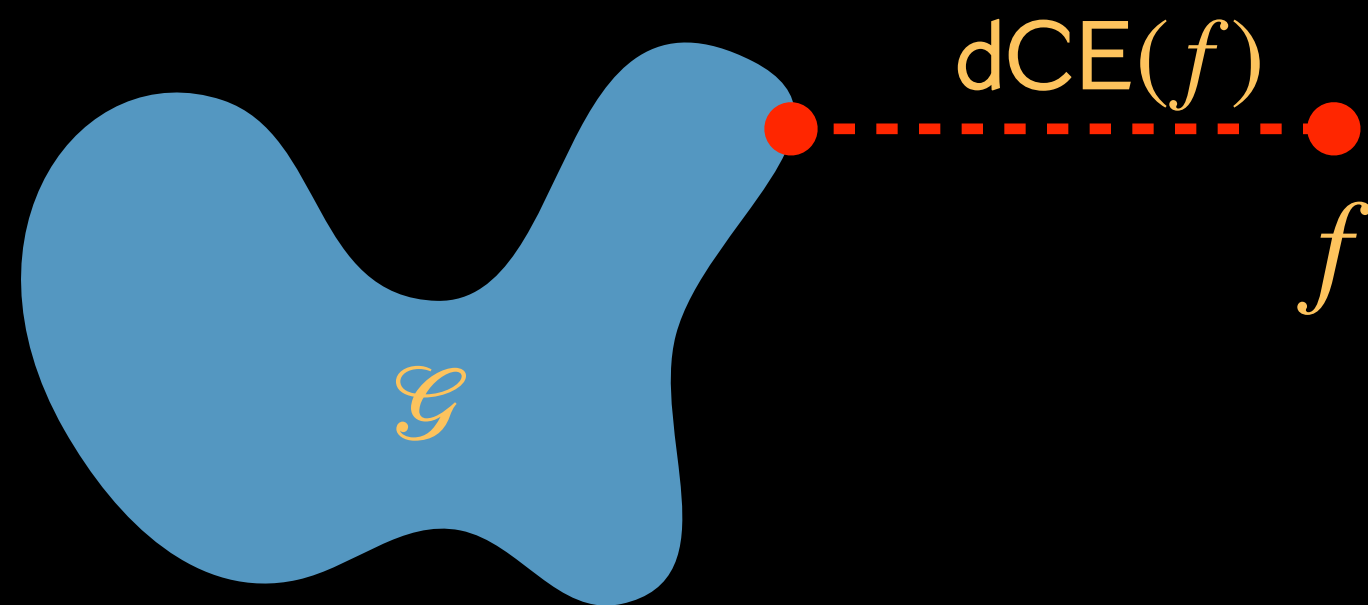
$\mathcal{G} :=$ set of perfectly-calibrated functions

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$



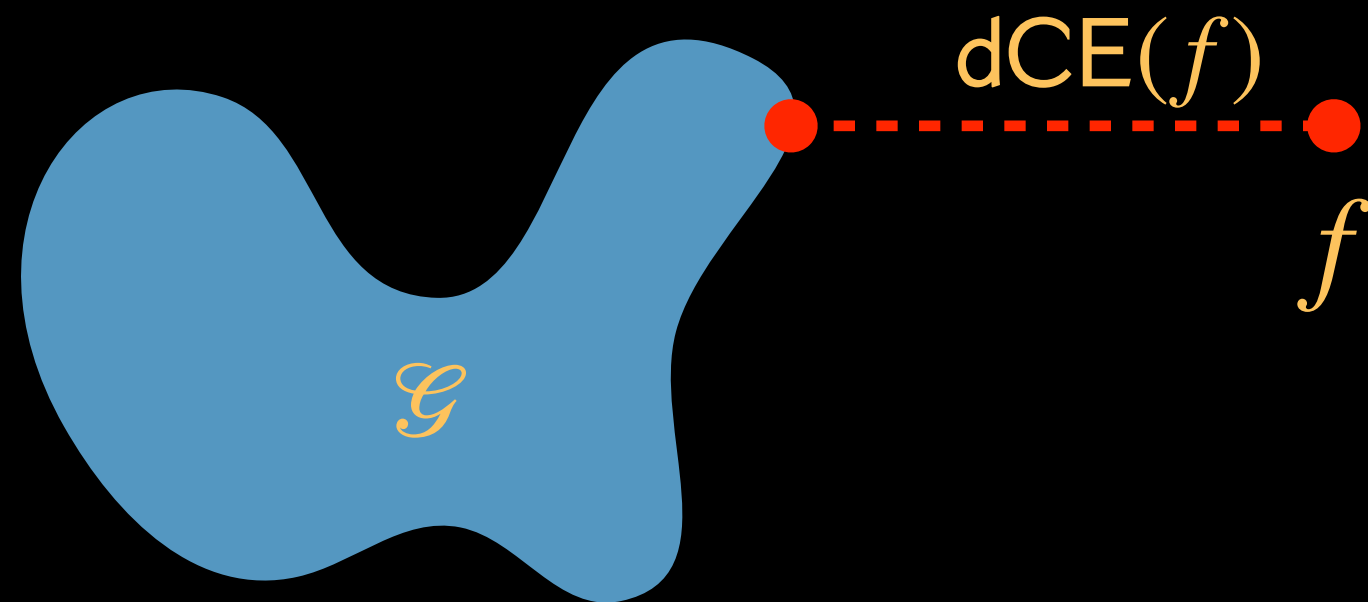
$\mathcal{G} :=$ set of perfectly-calibrated functions

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$



$\mathcal{G} :=$ set of perfectly-calibrated functions

The Calibration Distance

$$\text{dCE}(f) := \min_{g \in \mathcal{G}} d_1(f, g)$$

Correctness:

$$\mu(f) = 0 \iff f \text{ is perfectly calibrated}$$

Robust Correctness (informally):

$$\mu(f) \text{ is "close" to } 0 \iff f \text{ is "close" to perfectly calibrated}$$

Robust Correctness:

$$\text{dCE}(f)^\beta \leq \mu(f) \leq \text{dCE}(f)^\alpha$$

The Calibration Distance

$$\text{dCE}(f) := \min_{g \in \mathcal{G}} d_1(f, g)$$

Why not use $\text{dCE}(f)$ as calibration measure μ ?

Satisfies robust completeness:

$\mu(f)$ is "close" to 0 $\iff f$ is "close" to perfectly calibrated

The Calibration Distance

$$\text{dCE}(f) := \min_{g \in \mathcal{G}} d_1(f, g)$$

Why not use $\text{dCE}(f)$ as calibration measure μ ?

Satisfies robust completeness:

$\mu(f)$ is "close" to 0 $\iff f$ is "close" to perfectly calibrated

Q: How to estimate from samples $\{(f(x_i), y_i)\}$?

- Both info-theoretic, and computational issues...

The Calibration Distance

$$\text{dCE}(f) := \min_{g \in \mathcal{G}} d_1(f, g)$$

Unification

New metric $\text{dCE}(f)$ intimately related to existing metrics:

- $\text{kCE}(f)$: kernel calibration / MMCE [Kumar Sarawagi Jain 2018]
- $\text{smCE}(f)$: smooth calibration [Foster Hart 2018]
- $\text{intCE}(f)$: interval calibration

Unification

New metric $\text{dCE}(f)$ intimately related to existing metrics:

- $\text{kCE}(f)$: kernel calibration / MMCE [Kumar Sarawagi Jain 2018]
- $\text{smCE}(f)$: smooth calibration [Foster Hart 2018]
- $\text{intCE}(f)$: interval calibration

Theorem: For $\mu \in \{\text{kCE}, \text{smCE}, \text{intCE}\}$:

$$\mu^2 \leq \text{dCE} \leq \mu^{1/3}$$

Unification

New metric $\text{dCE}(f)$ intimately related to existing metrics:

- $\text{kCE}(f)$: kernel calibration / MMCE [Kumar Sarawagi Jain 2018]
- $\text{smCE}(f)$: smooth calibration [Foster Hart 2018]
- $\text{intCE}(f)$: interval calibration

Theorem: For $\mu \in \{\text{kCE}, \text{smCE}, \text{intCE}\}$:

$$\mu^2 \leq \text{dCE} \leq \mu^{1/3}$$

Takeaway:

1. Estimate dCE from samples
2. Prior metrics are related

Practical Takeaways

Measure calibration with either:

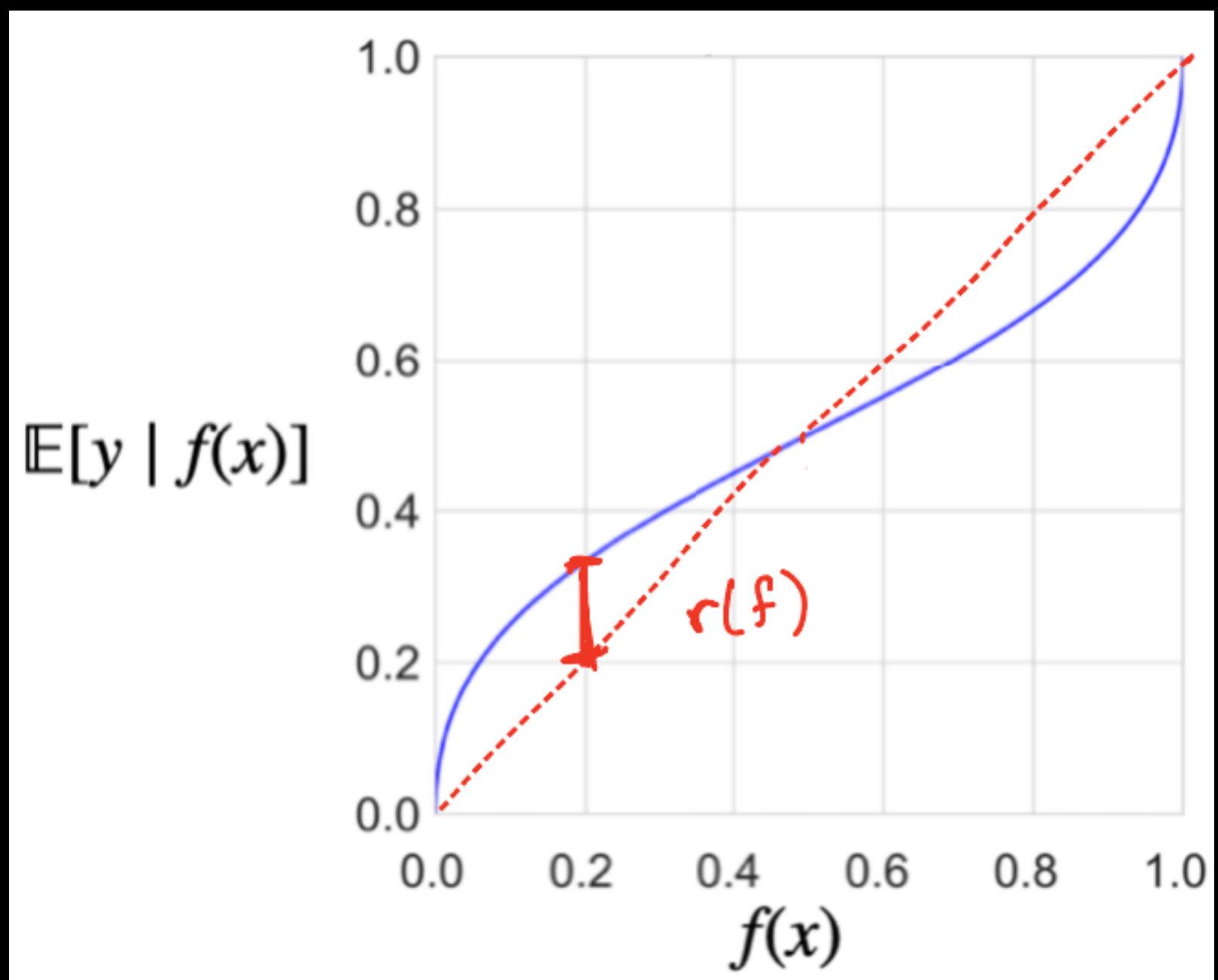
1. Kernel Calibration Error
2. Interval Calibration Error (modification of binnedECE)

Kernel Calibration Error

$$\text{ECE}_{\mathcal{D}}(f) := \sup_{w:[0,1] \rightarrow [-1,1]} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$

Kernel Calibration Error

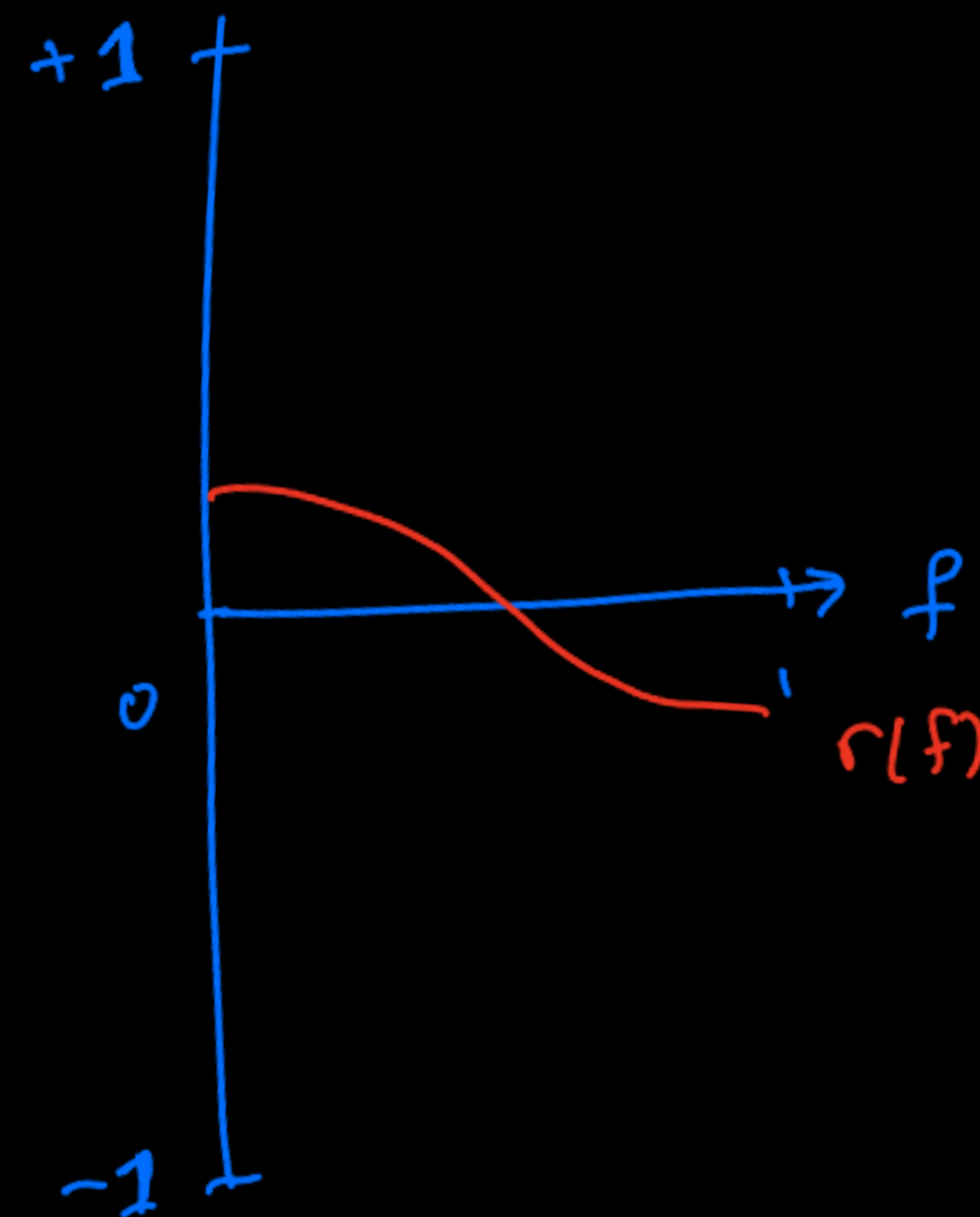
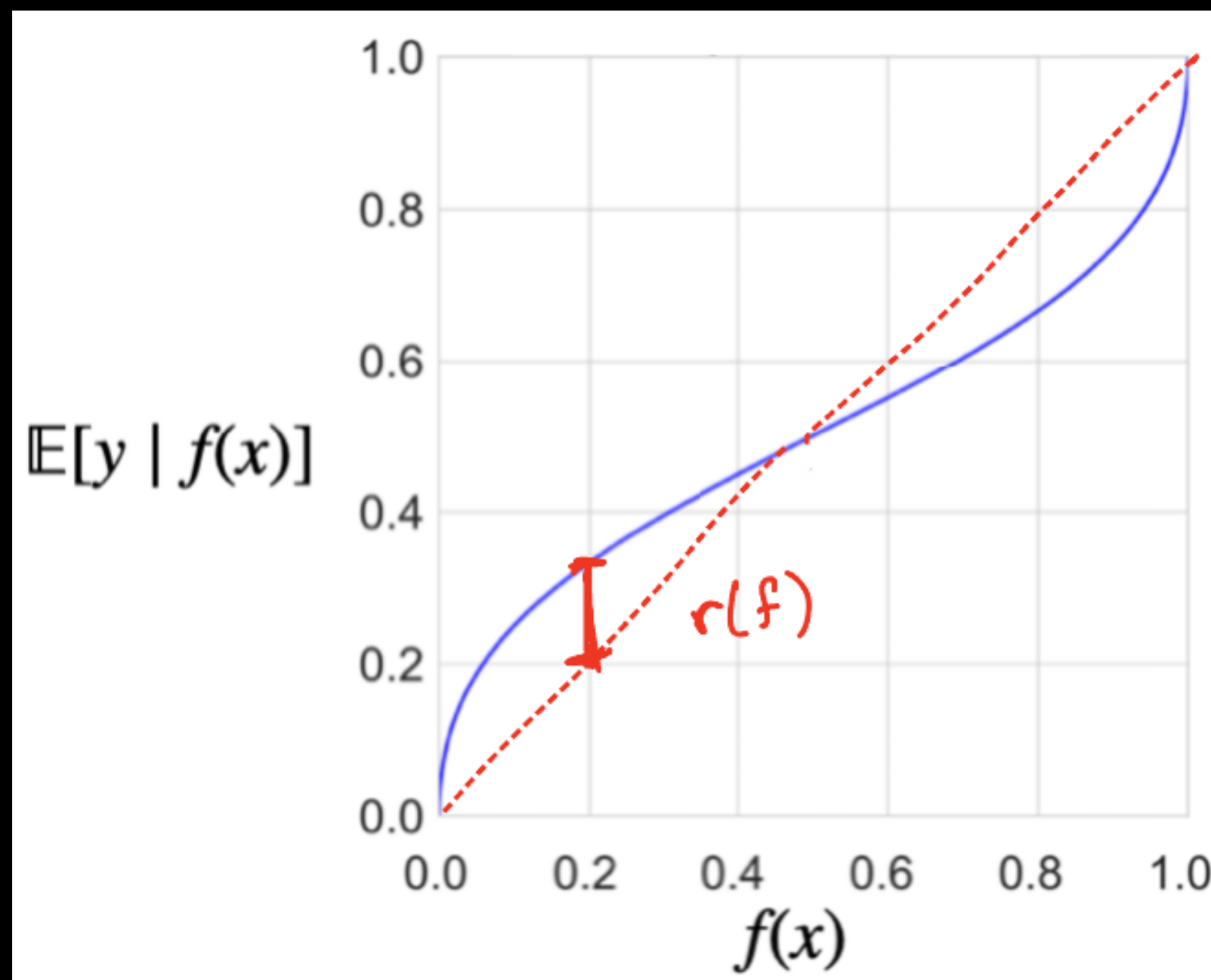
$$\text{ECE}_{\mathcal{D}}(f) := \sup_{w: [0,1] \rightarrow [-1,1]} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$



“Residual”: $r(f) := \mathbb{E}[y \mid f] - f$

Kernel Calibration Error

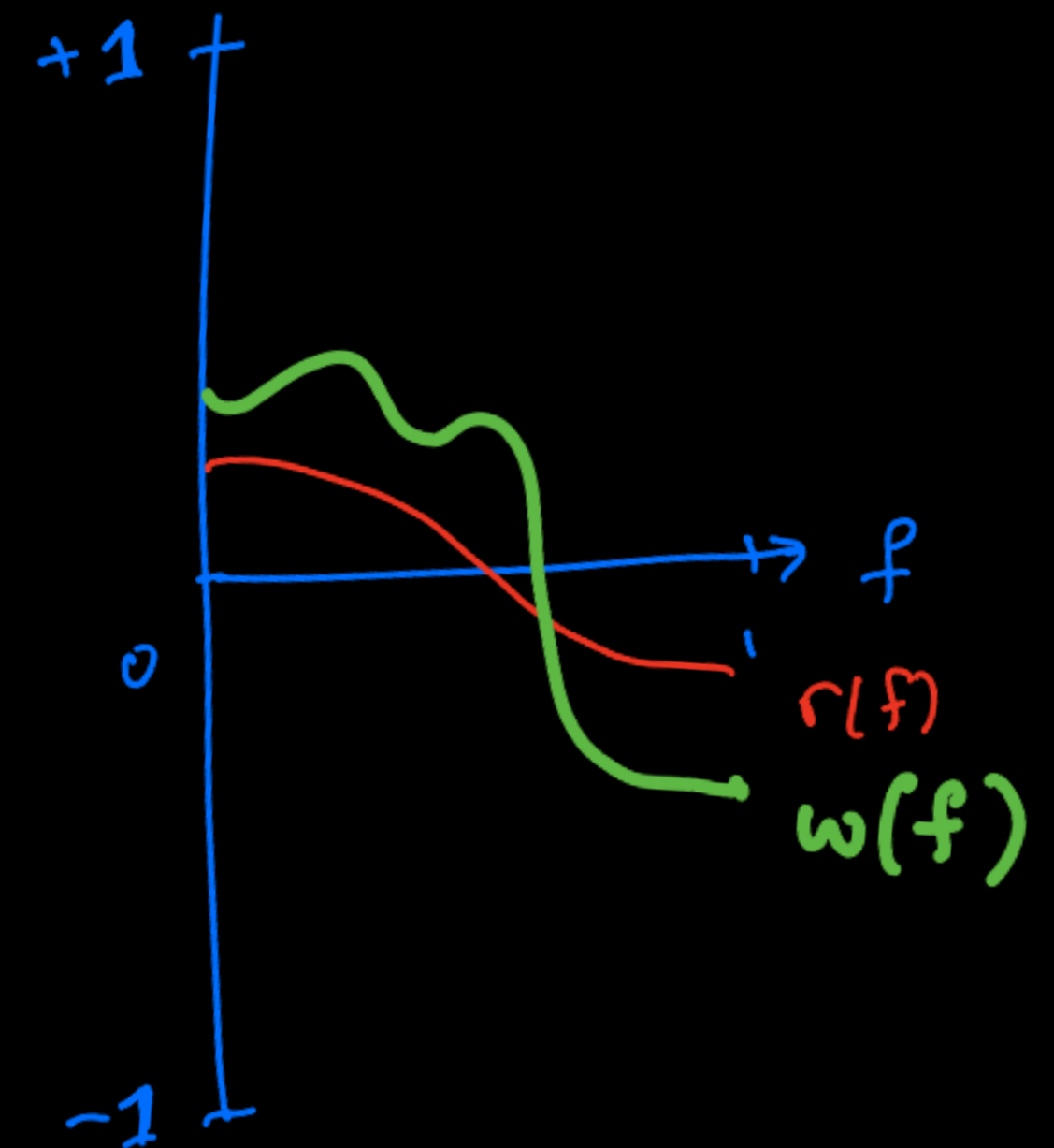
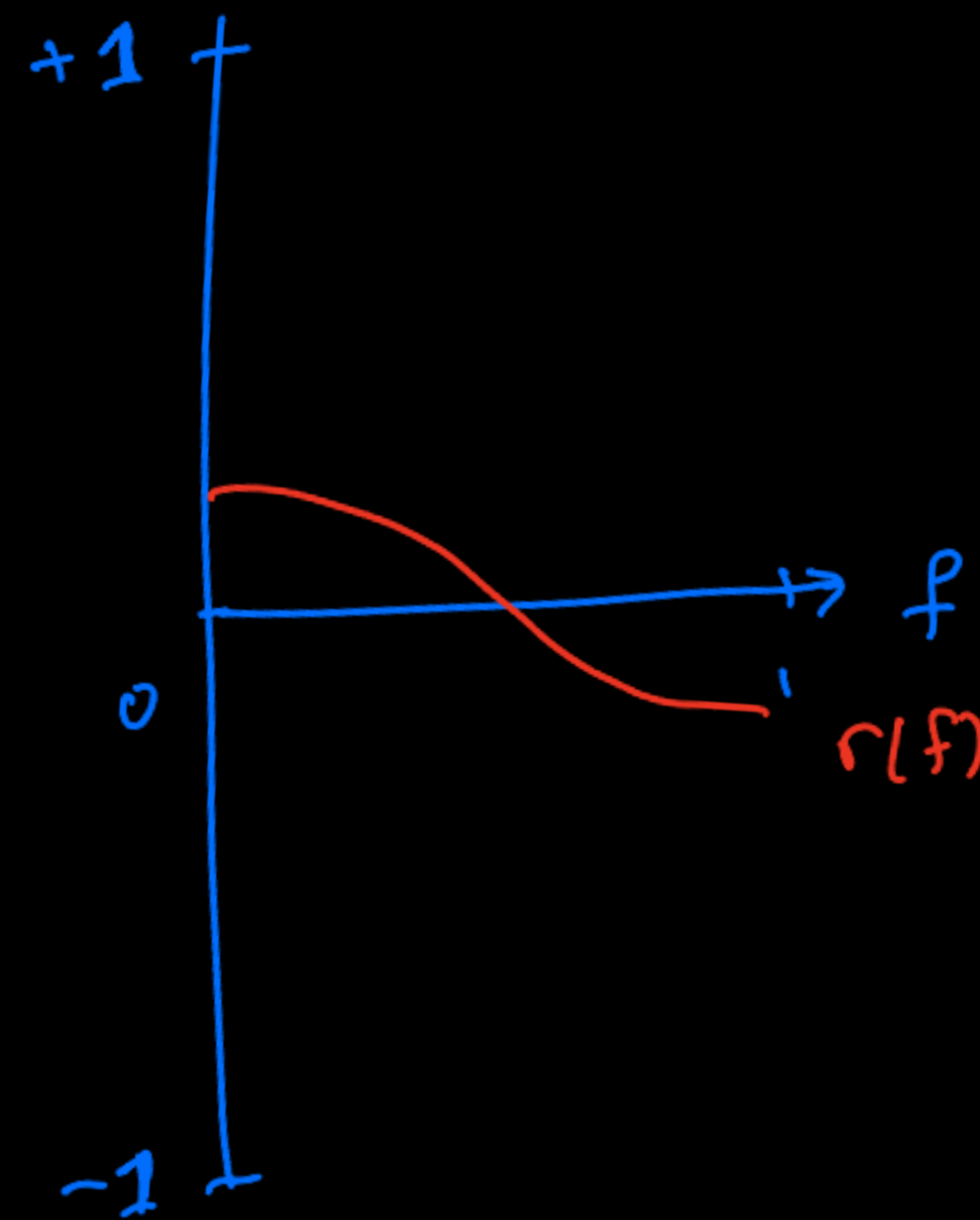
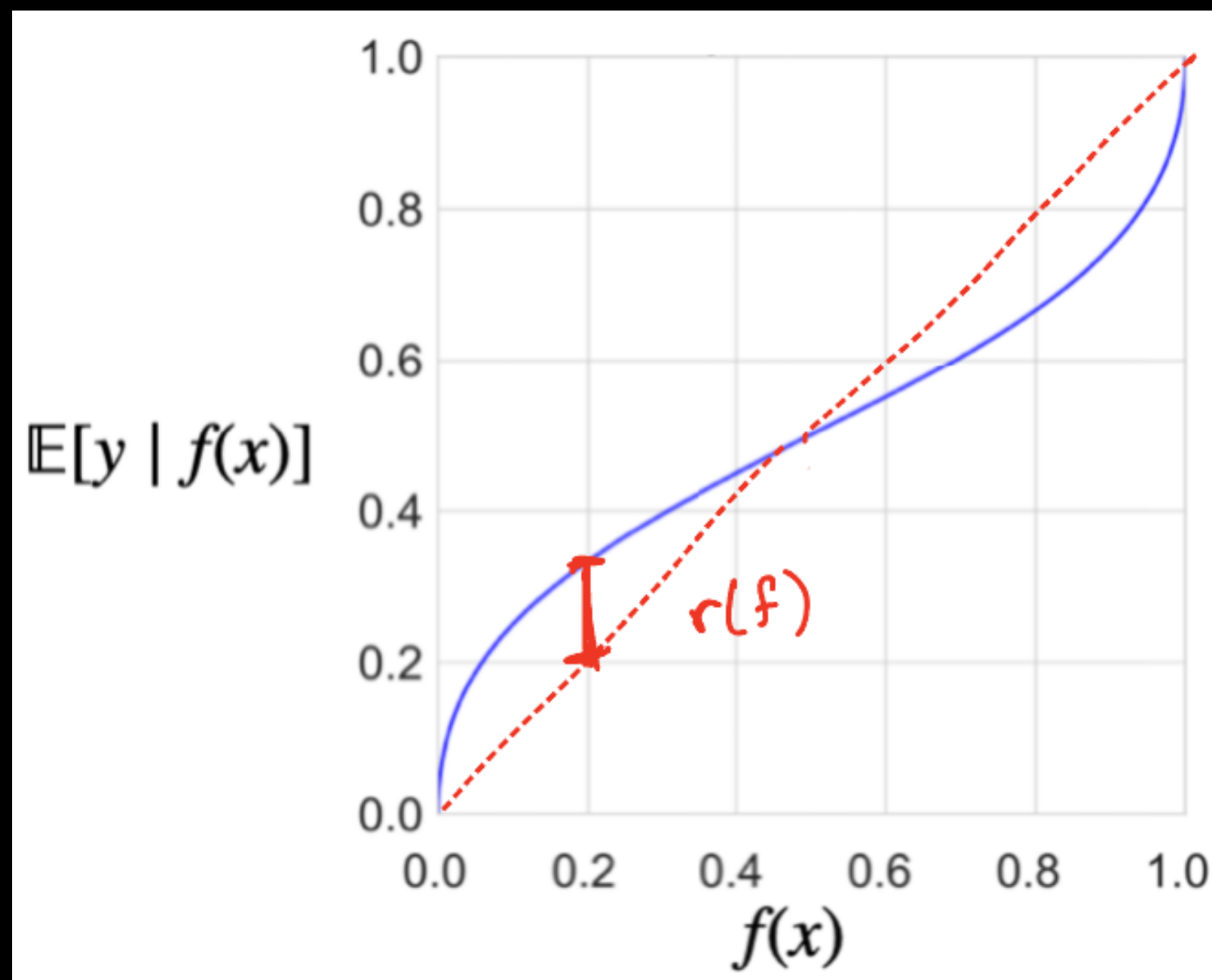
$$\text{ECE}_{\mathcal{D}}(f) := \sup_{w: [0,1] \rightarrow [-1,1]} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$



"Residual": $r(f) := \mathbb{E}[y | f] - f$

Kernel Calibration Error

$$\text{ECE}_{\mathcal{D}}(f) := \sup_{w: [0,1] \rightarrow [-1,1]} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$

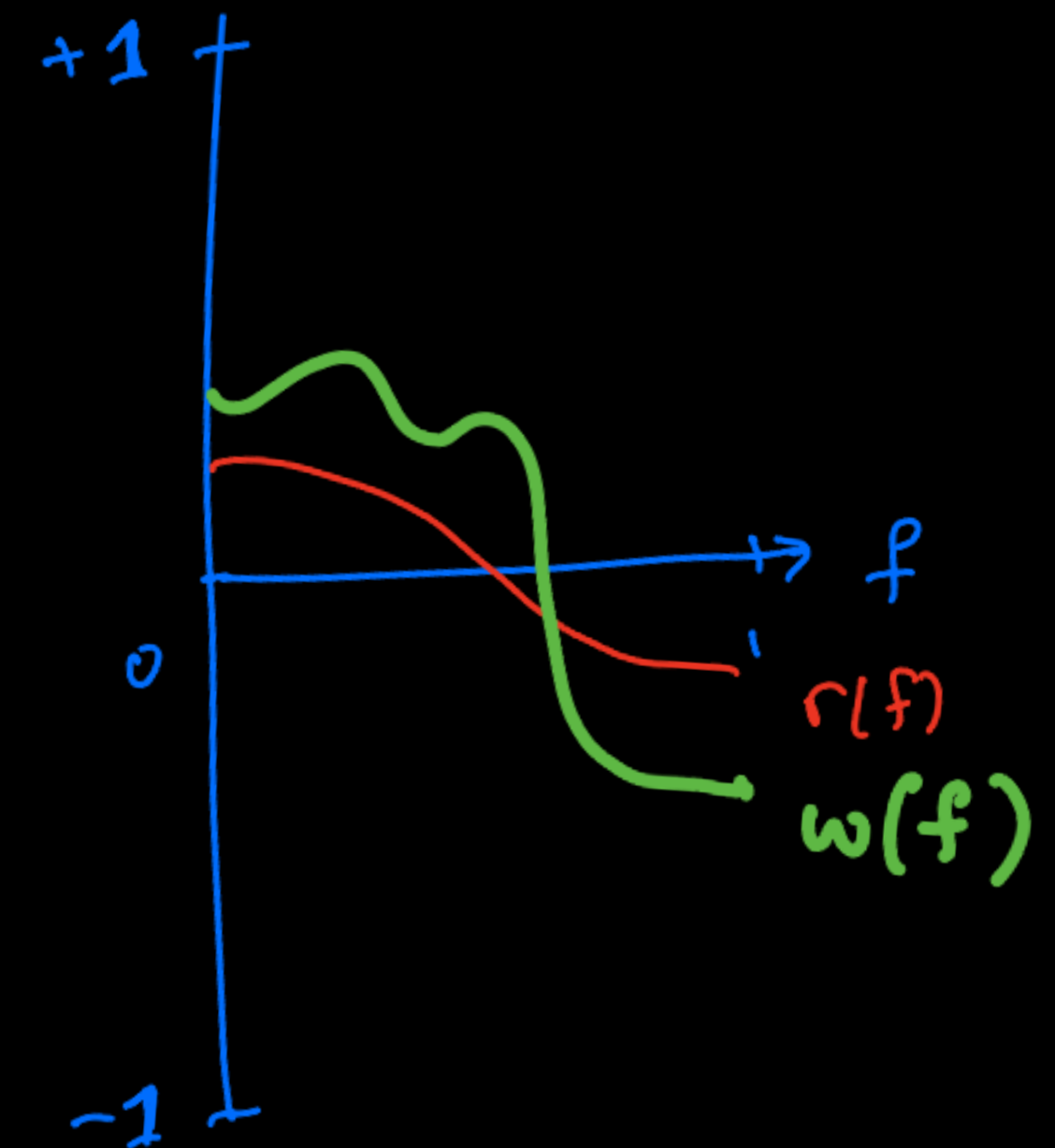
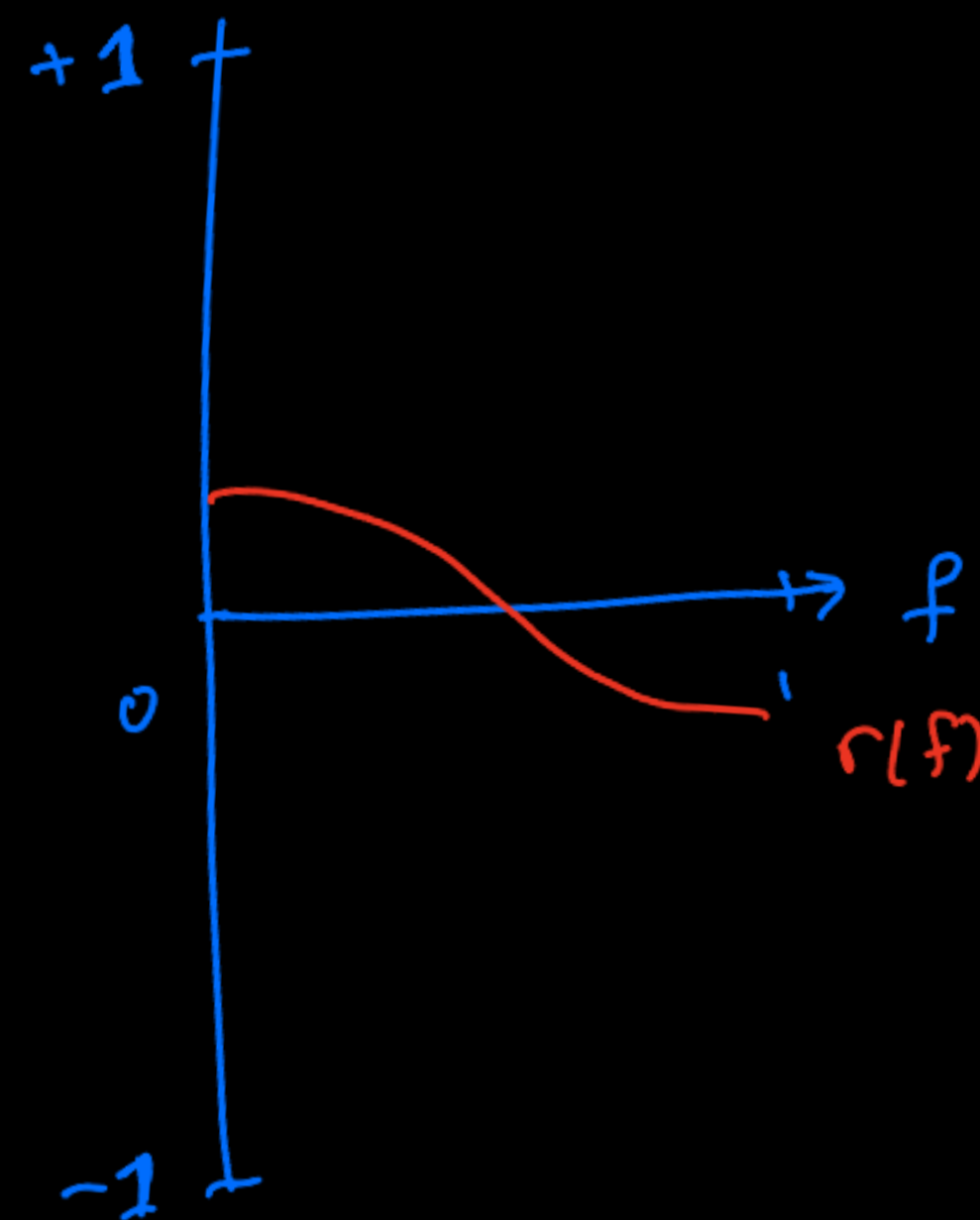
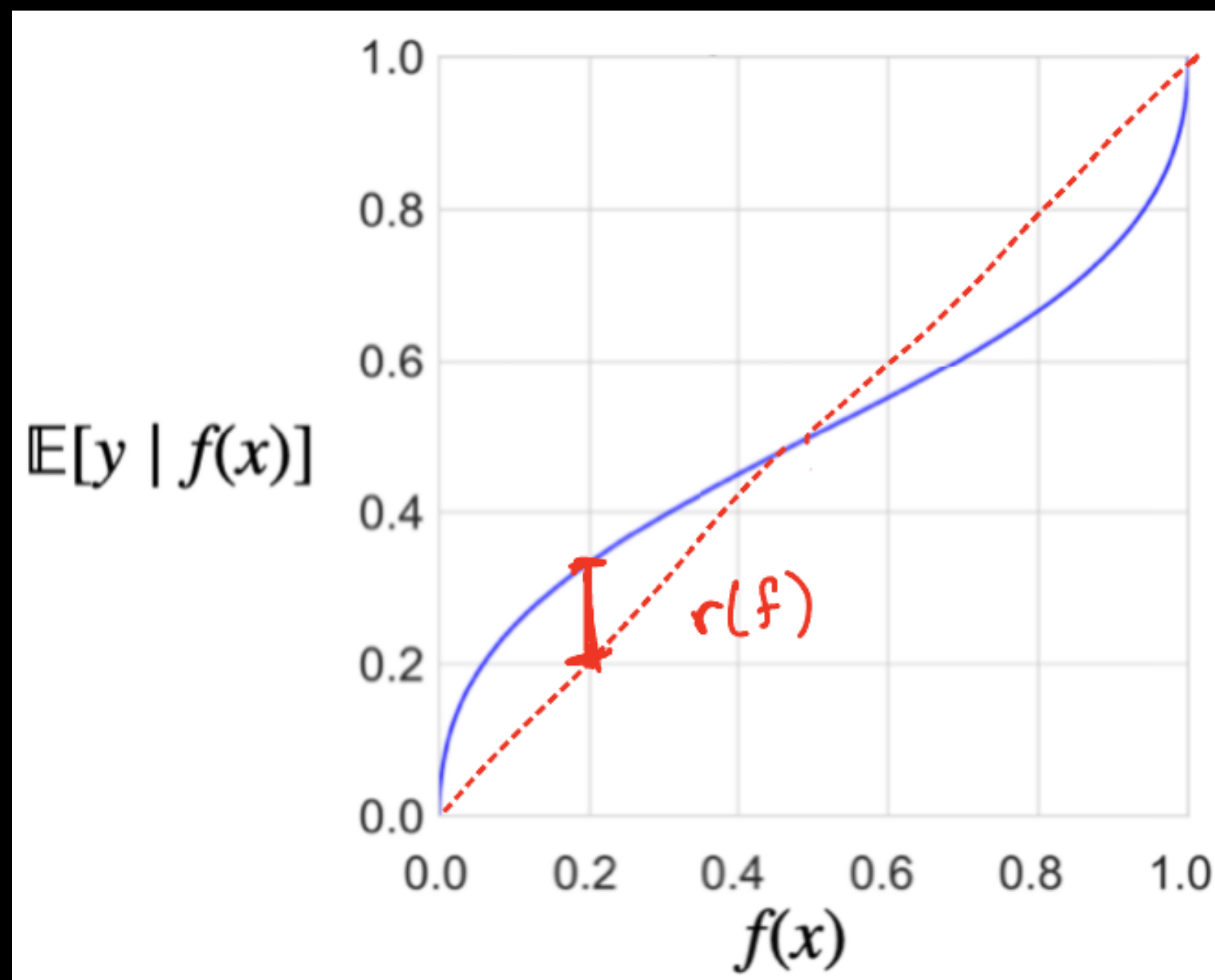


"Residual": $r(f) := \mathbb{E}[y | f] - f$

Kernel Calibration Error

$$\text{ECE}_{\mathcal{D}}(f) := \sup_{w: [0,1] \rightarrow [-1,1]} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$

$$\text{kCE}_{\mathcal{D}}(f) := \sup_{w: \|w\|_K \leq 1} \mathbb{E}_{(f,y) \sim \mathcal{D}_f} [w(f)(y - f)]$$



"Residual": $r(f) := \mathbb{E}[y | f] - f$

Kernel Calibration Error: Sample Estimation

Given: Samples $(f(x_i), y_i) =: (f_i, y_i)$

Residuals: (f_i, r_i) for $r_i := (y_i - f_i)$

Kernel Calibration Error: Sample Estimation

Given: Samples $(f(x_i), y_i) =: (f_i, y_i)$

Residuals: (f_i, r_i) for $r_i := (y_i - f_i)$

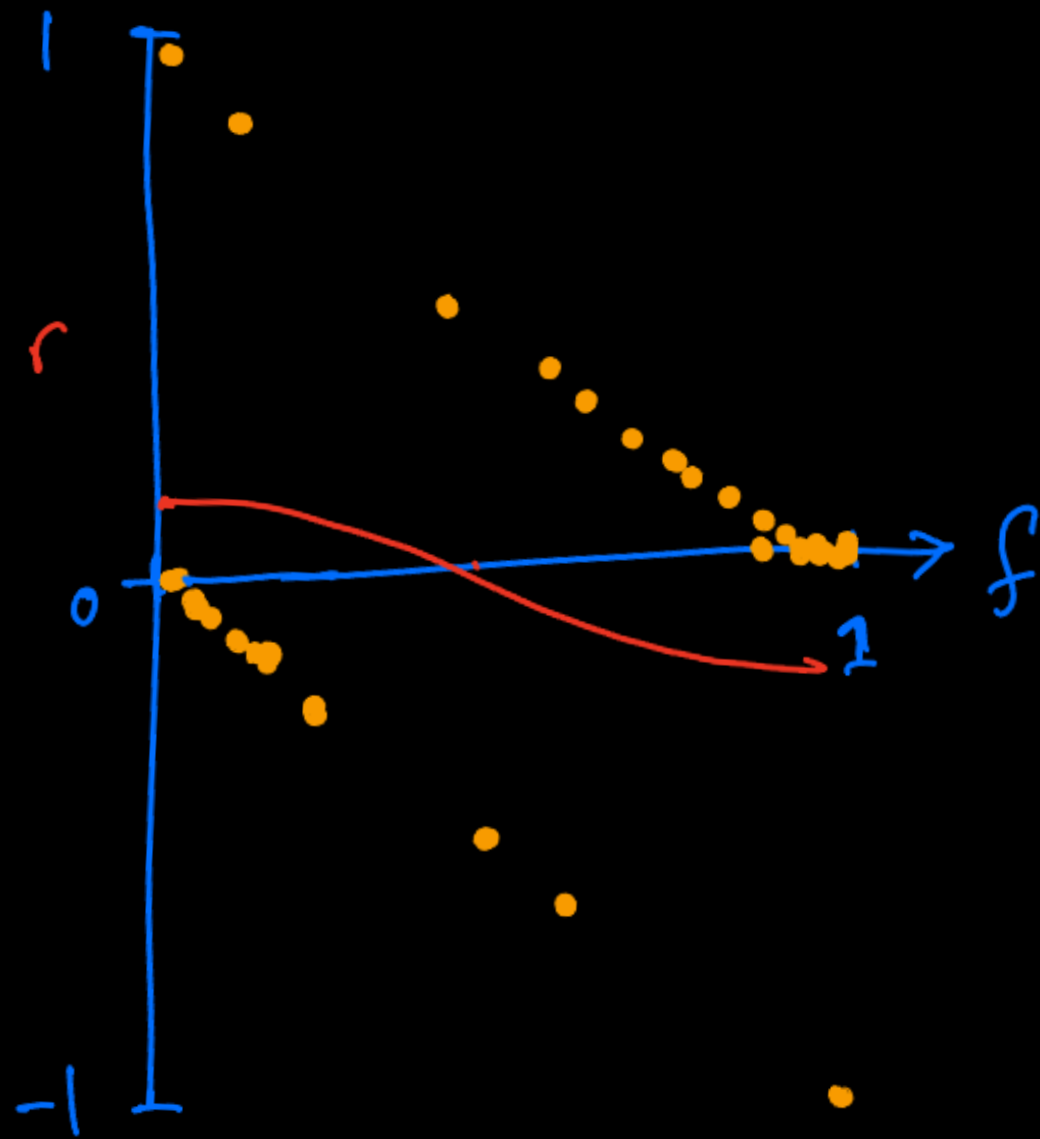
$$\widehat{\text{kCE}}_{\mathcal{D}}(f) = \sqrt{\frac{1}{n^2} \sum_{i,j} r_i r_j K(f_i, f_j)} = \|r\|_{K(f,f)}$$

Kernel Calibration Error: Sample Estimation

Given: Samples $(f(x_i), y_i) =: (f_i, y_i)$

Residuals: (f_i, r_i) for $r_i := (y_i - f_i)$

$$\widehat{\text{kCE}}_{\mathcal{D}}(f) = \sqrt{\frac{1}{n^2} \sum_{i,j} r_i r_j K(f_i, f_j)} = \|r\|_{K(f,f)}$$

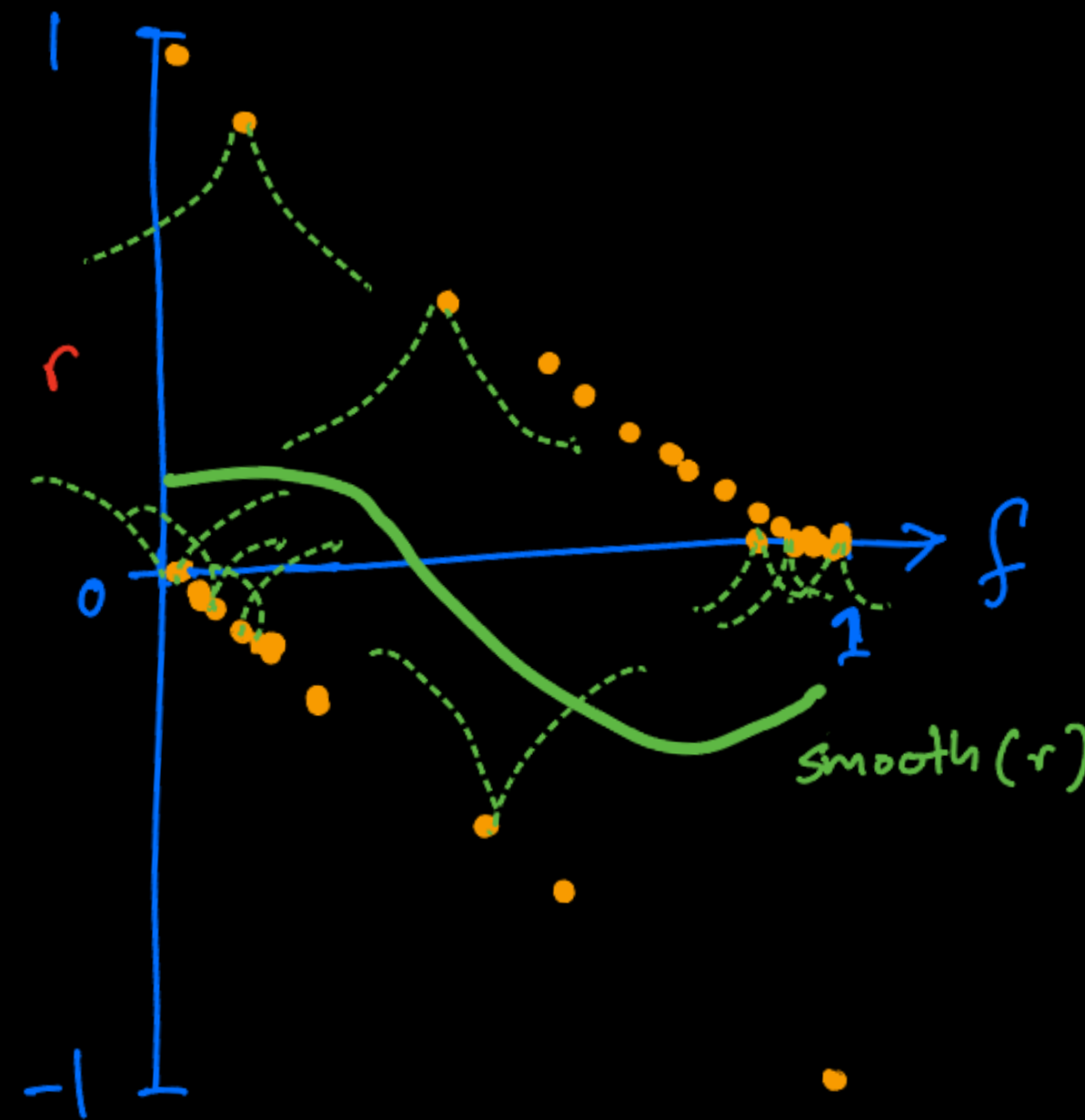
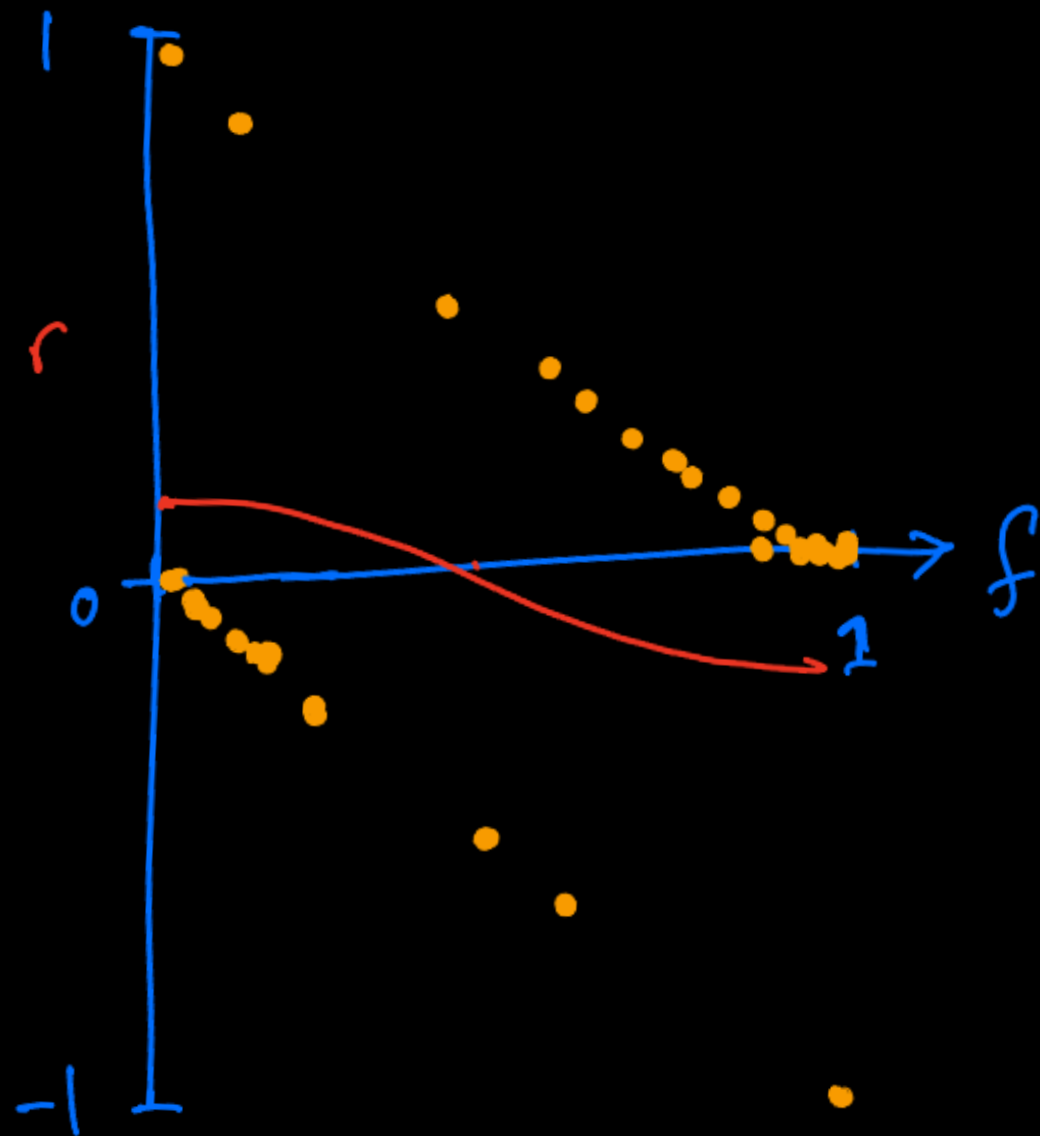


Kernel Calibration Error: Sample Estimation

Given: Samples $(f(x_i), y_i) =: (f_i, y_i)$

Residuals: (f_i, r_i) for $r_i := (y_i - f_i)$

$$\begin{aligned}\widehat{\text{kCE}}_{\mathcal{D}}(f) &= \sqrt{\frac{1}{n^2} \sum_{i,j} r_i r_j K(f_i, f_j)} = \|r\|_{K(f,f)} \\ &= \sqrt{\langle r, \text{smooth}_K(r) \rangle}\end{aligned}$$

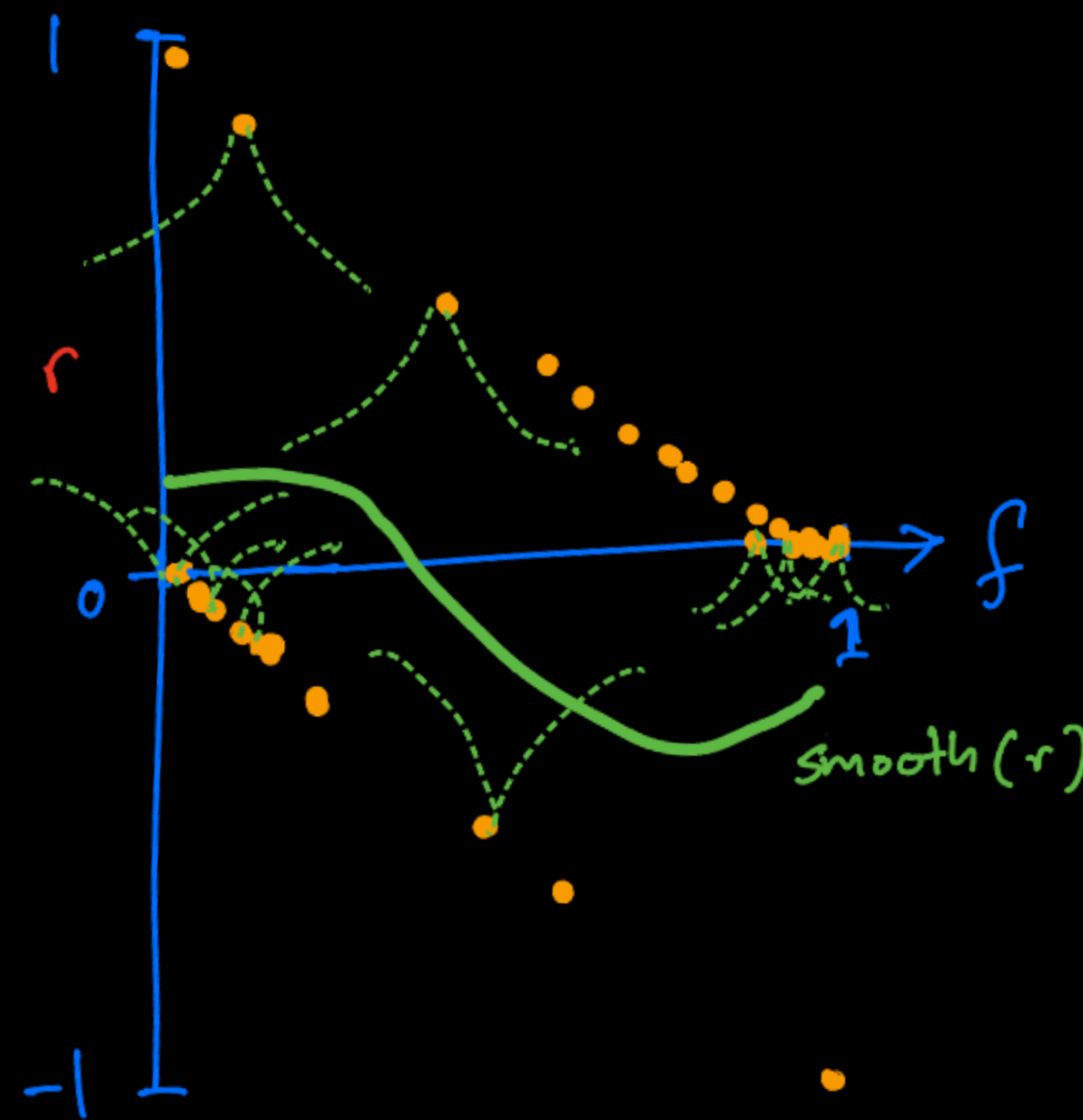
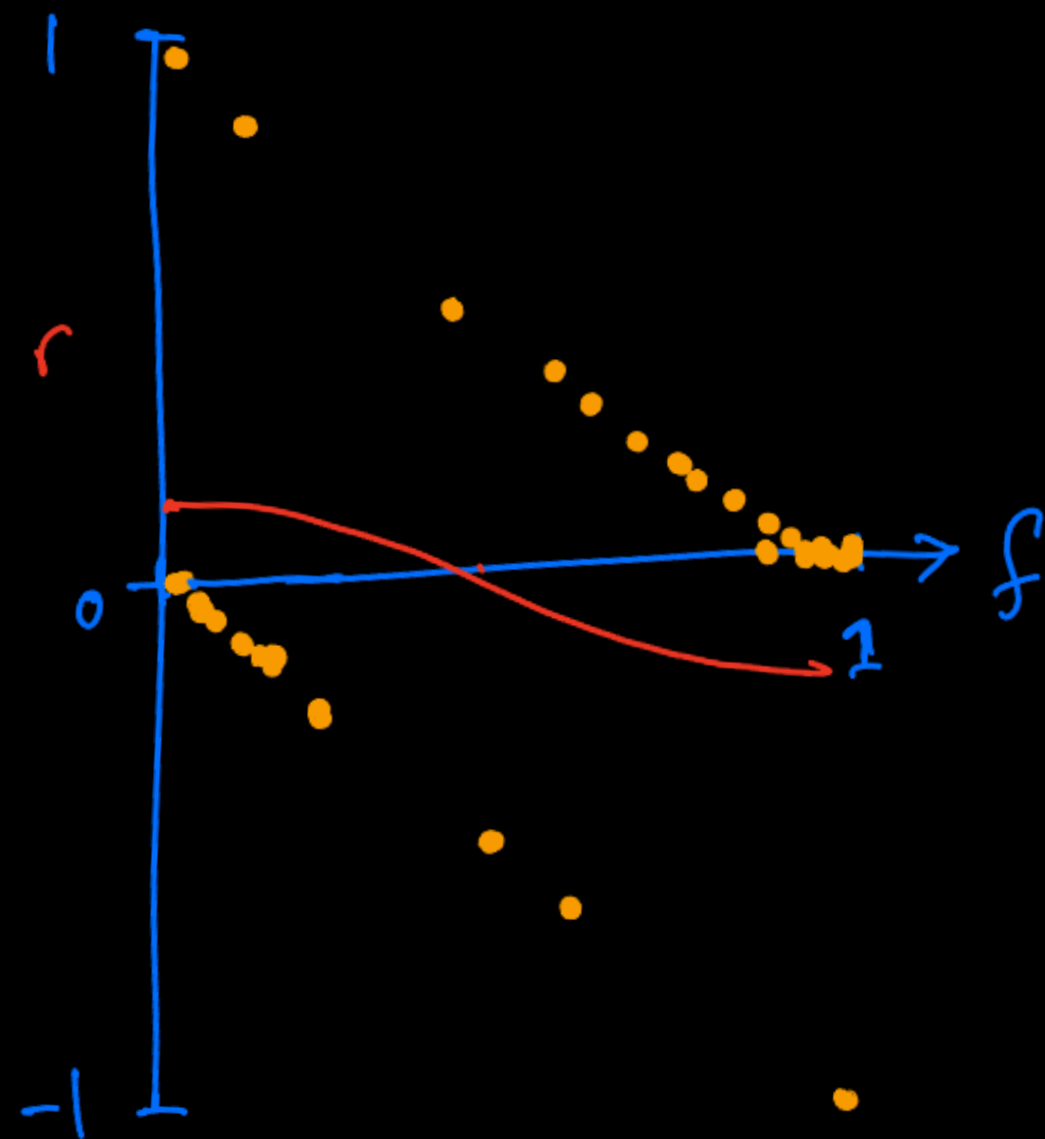


Kernel Calibration Error: Sample Estimation

Given: Samples $(f(x_i), y_i) =: (f_i, y_i)$

Residuals: (f_i, r_i) for $r_i := (y_i - f_i)$

$$\begin{aligned}\widehat{\text{kCE}}_{\mathcal{D}}(f) &= \sqrt{\frac{1}{n^2} \sum_{i,j} r_i r_j K(f_i, f_j)} = \|r\|_{K(f,f)} \\ &= \sqrt{\langle r, \text{smooth}_K(r) \rangle}\end{aligned}$$

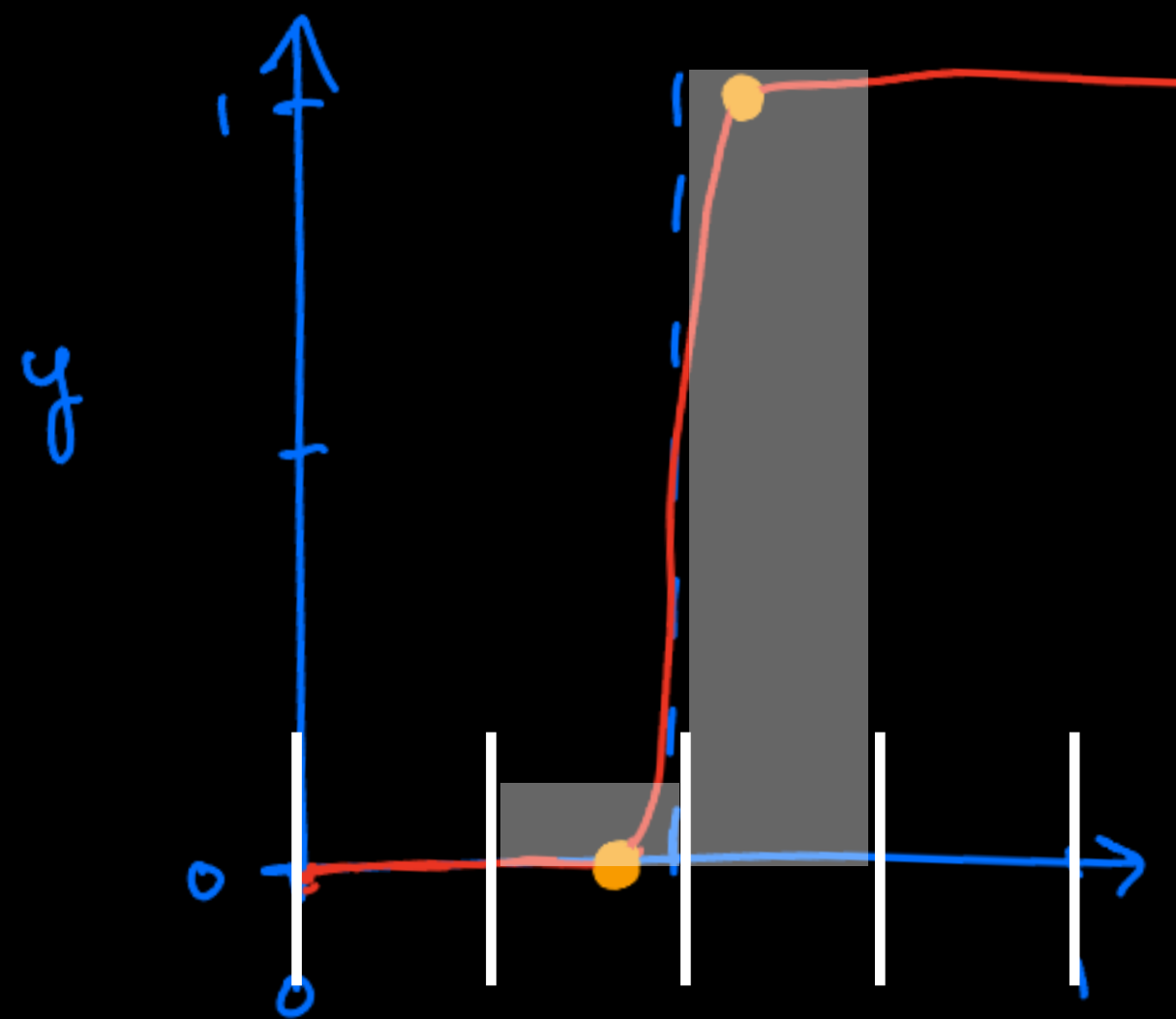
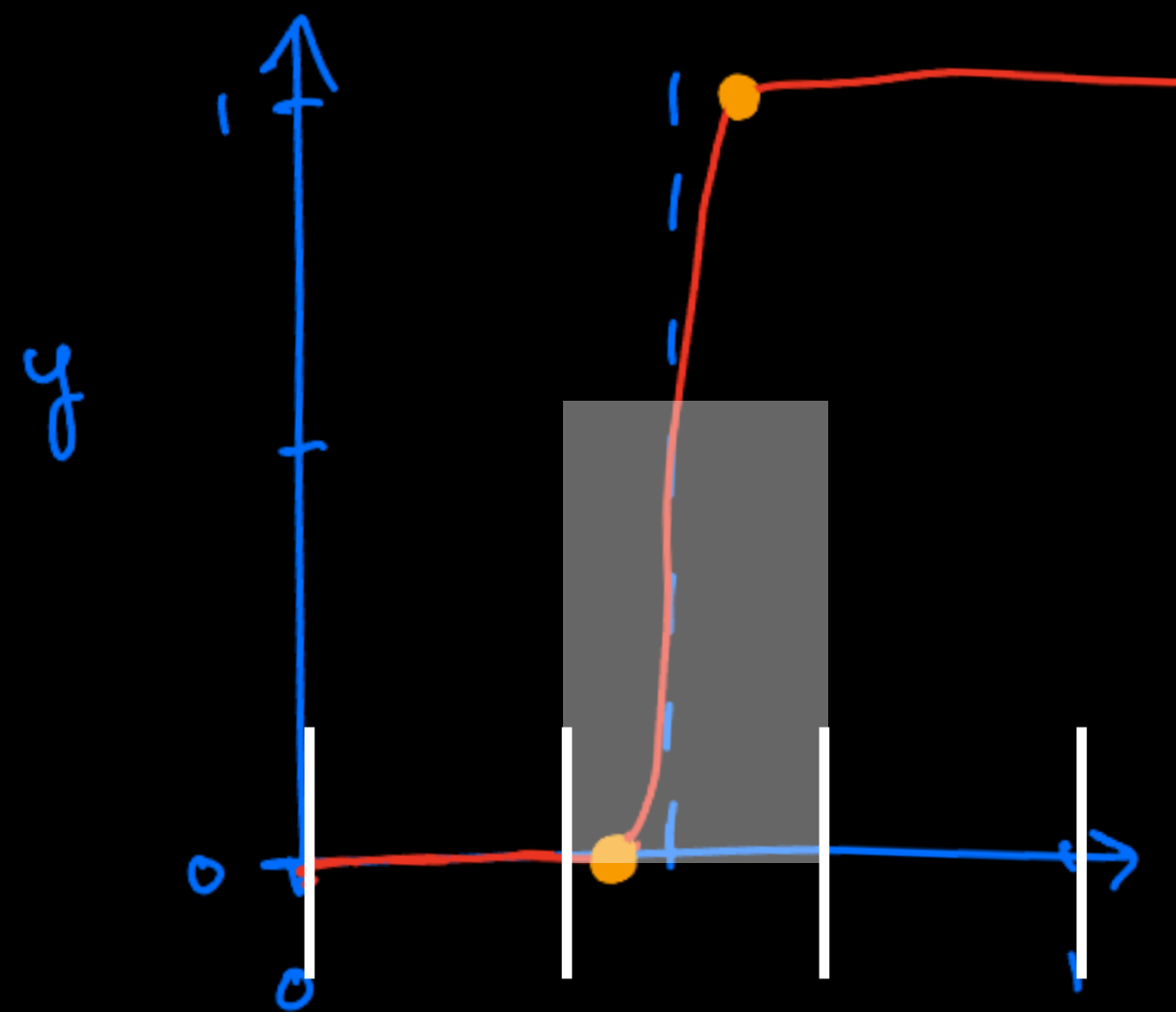


- Linear-time estimation: sub-sample $\Theta(n)$ terms
- Requires **Laplace** kernel

Interval Calibration

$$\text{binnedECE}(f, \mathcal{I}) := \text{ECE}(\text{round}_{\mathcal{I}}(f))$$

binnedECE: Unclear how to choose bins (any fixed choice violates continuity & correctness)



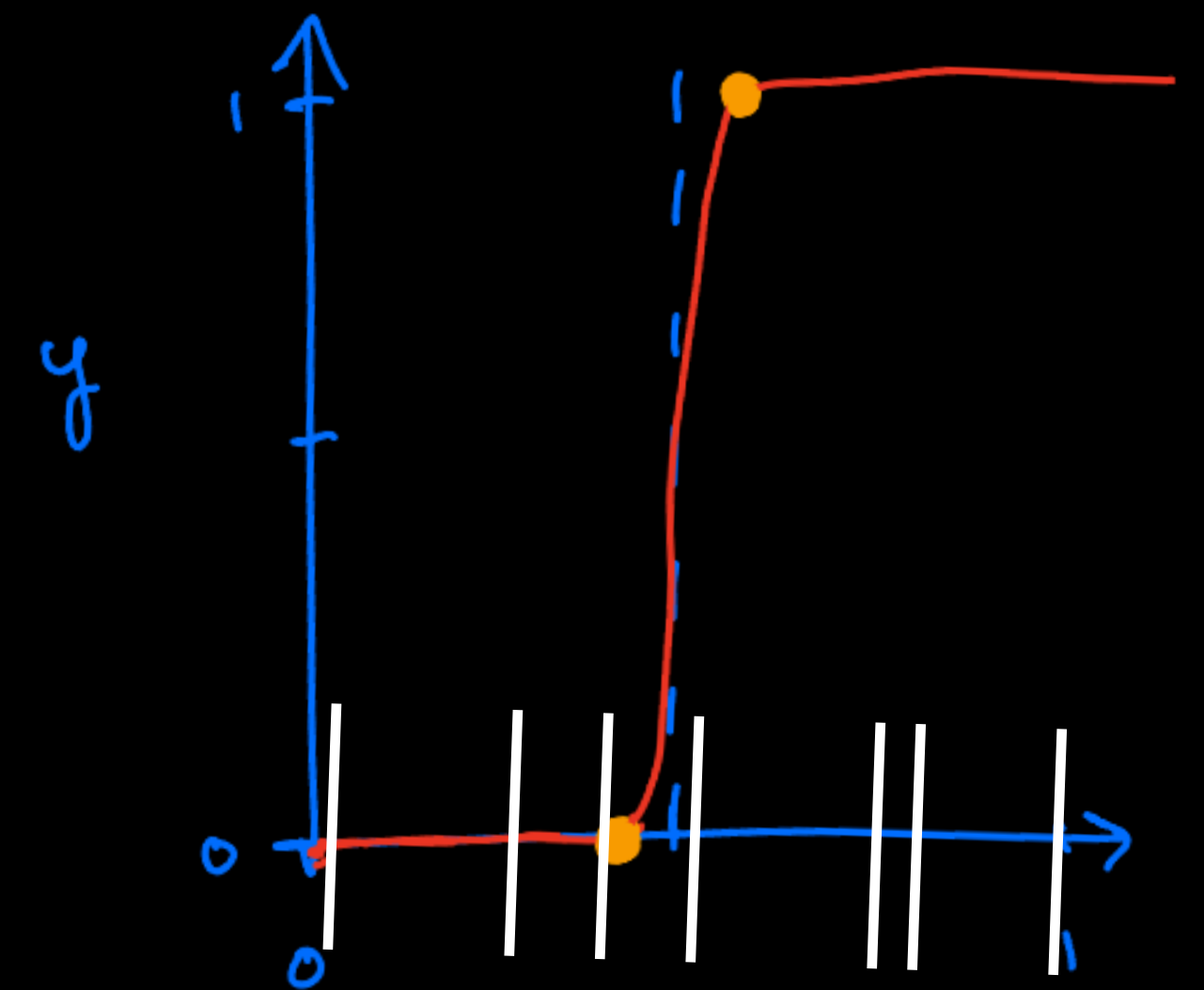
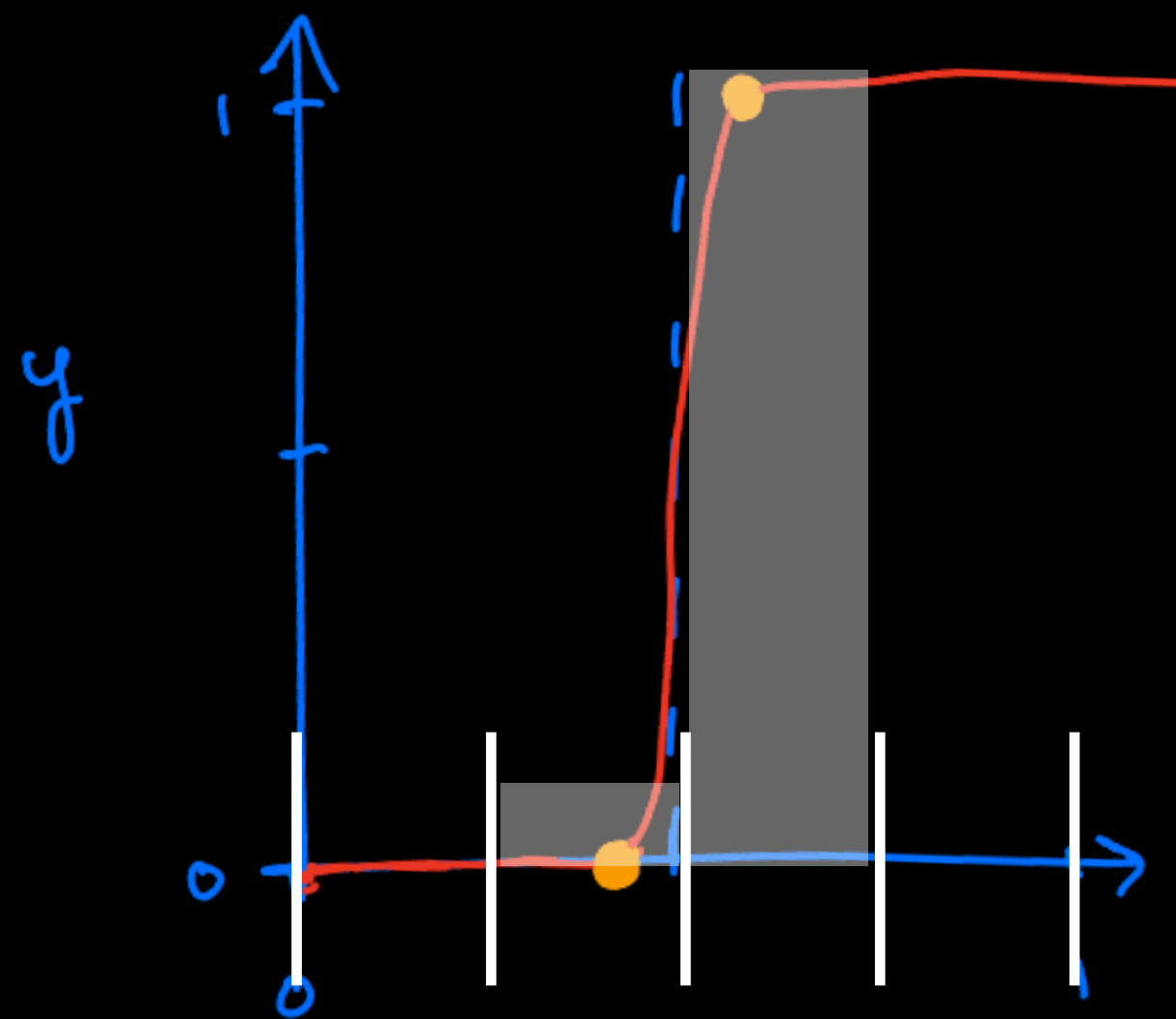
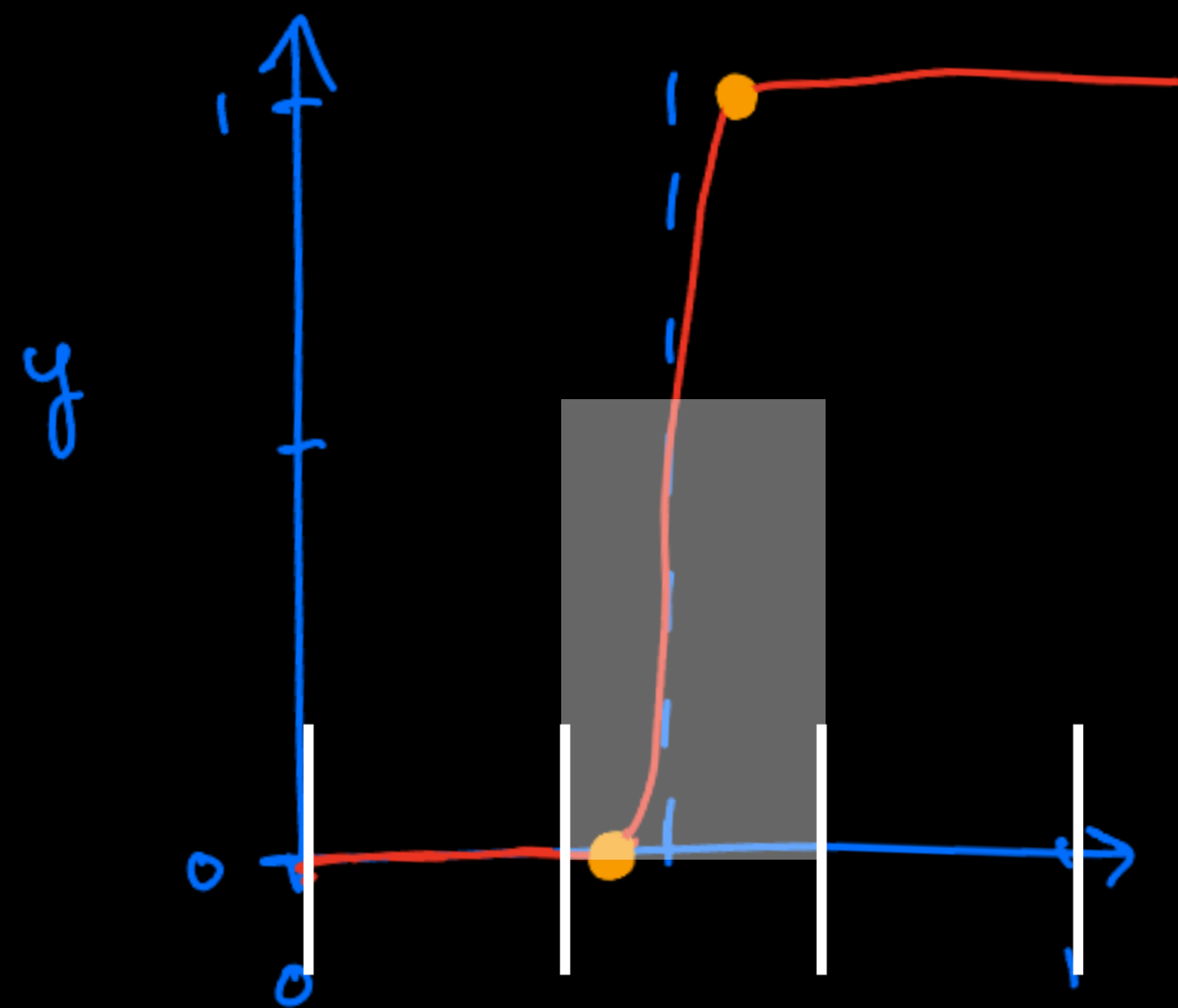
Interval Calibration

$$\text{binnedECE}(f, \mathcal{I}) := \text{ECE}(\text{round}_{\mathcal{I}}(f))$$

binnedECE: Unclear how to choose bins (any fixed choice violates continuity & correctness)

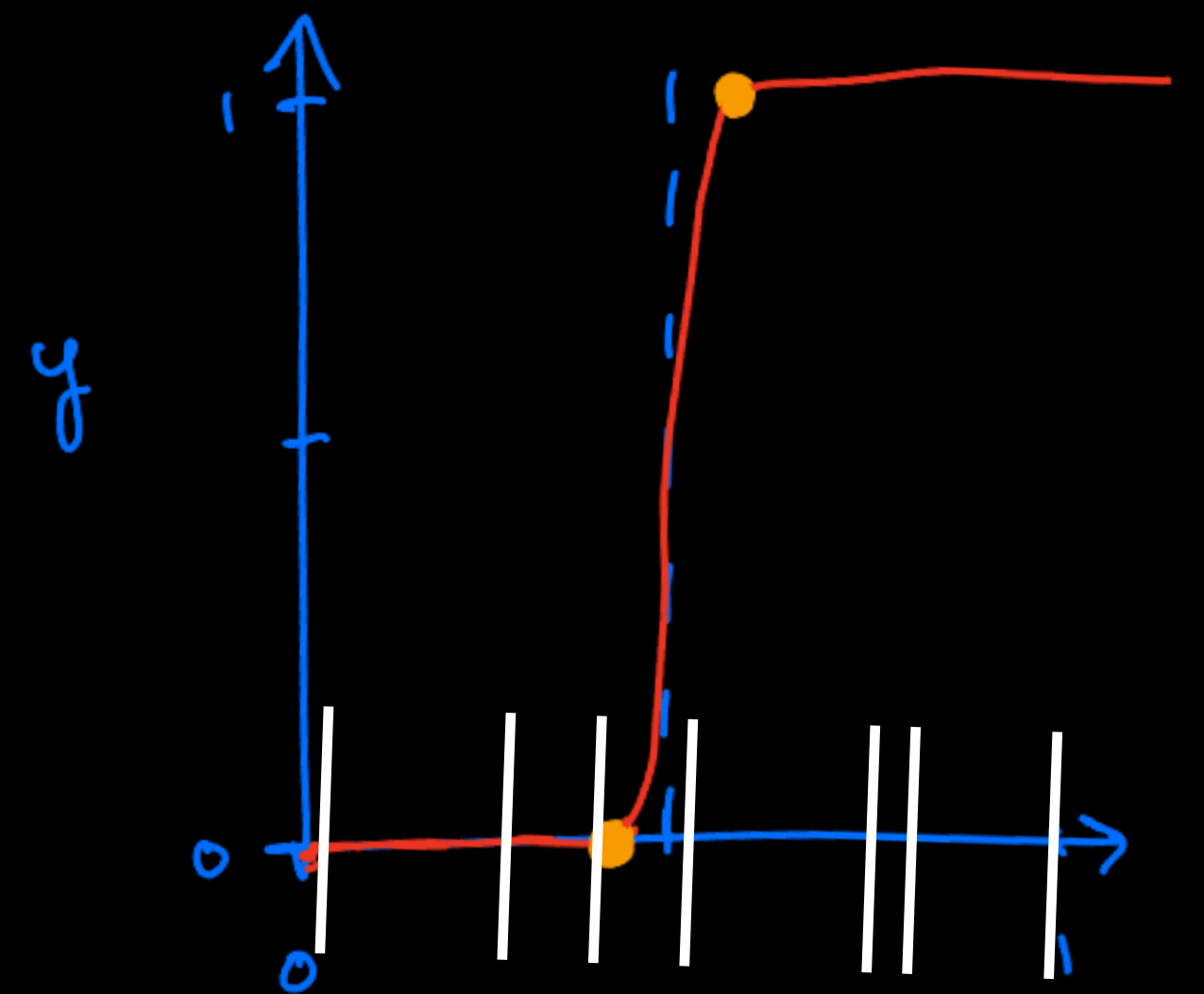
But, adding a “width regularizer” guarantees upper-bound. For all interval-partitions:

$$\text{dCE}(f) \leq \text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$$



Interval Calibration

$$\text{dCE}(f) \leq \text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$$

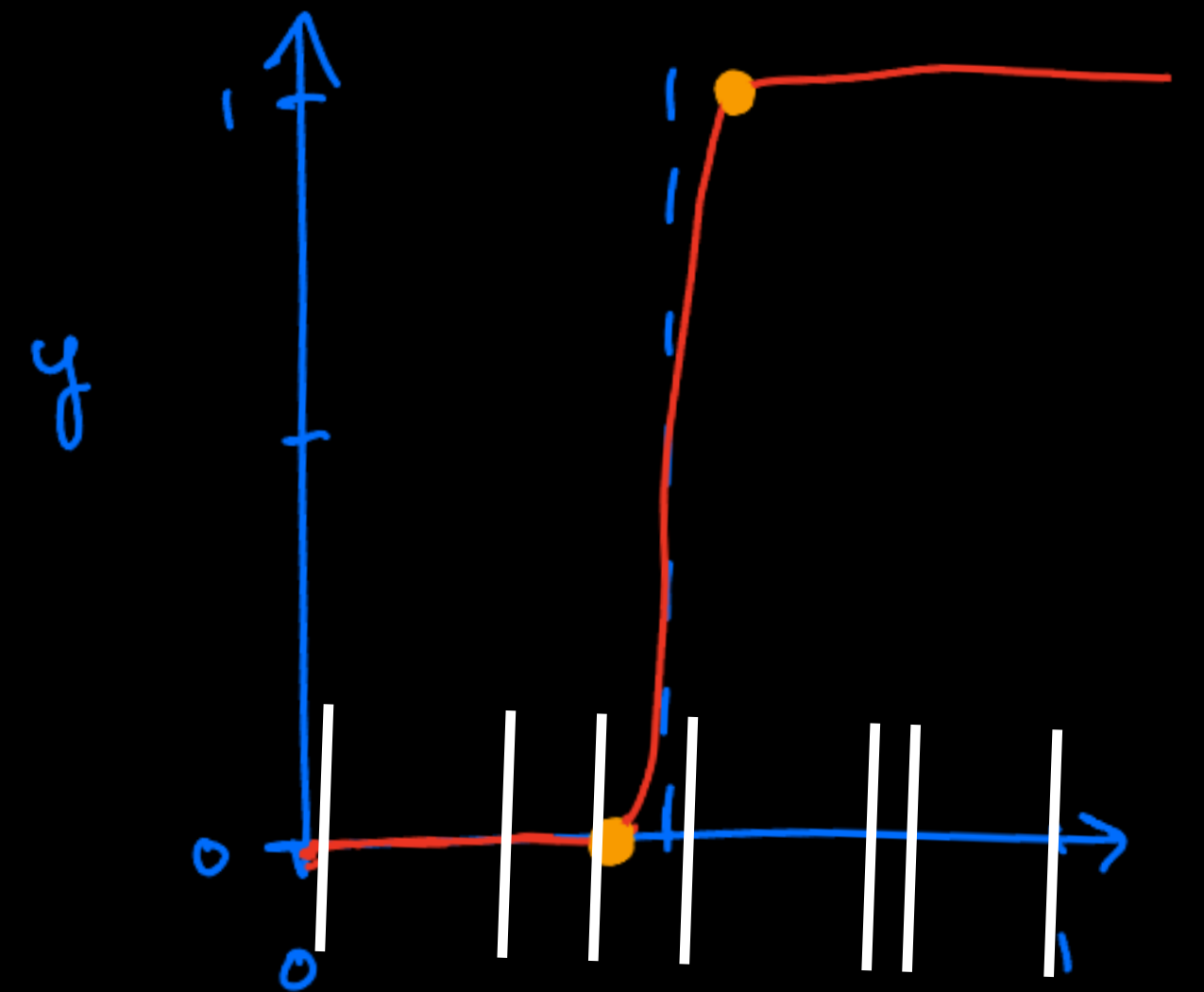


Interval Calibration

$$\text{dCE}(f) \leq \text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$$

Best-possible upper-bound:

$$\text{intCE}(f) := \inf_{\mathcal{I}: \text{Interval partition}} (\text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I}))$$



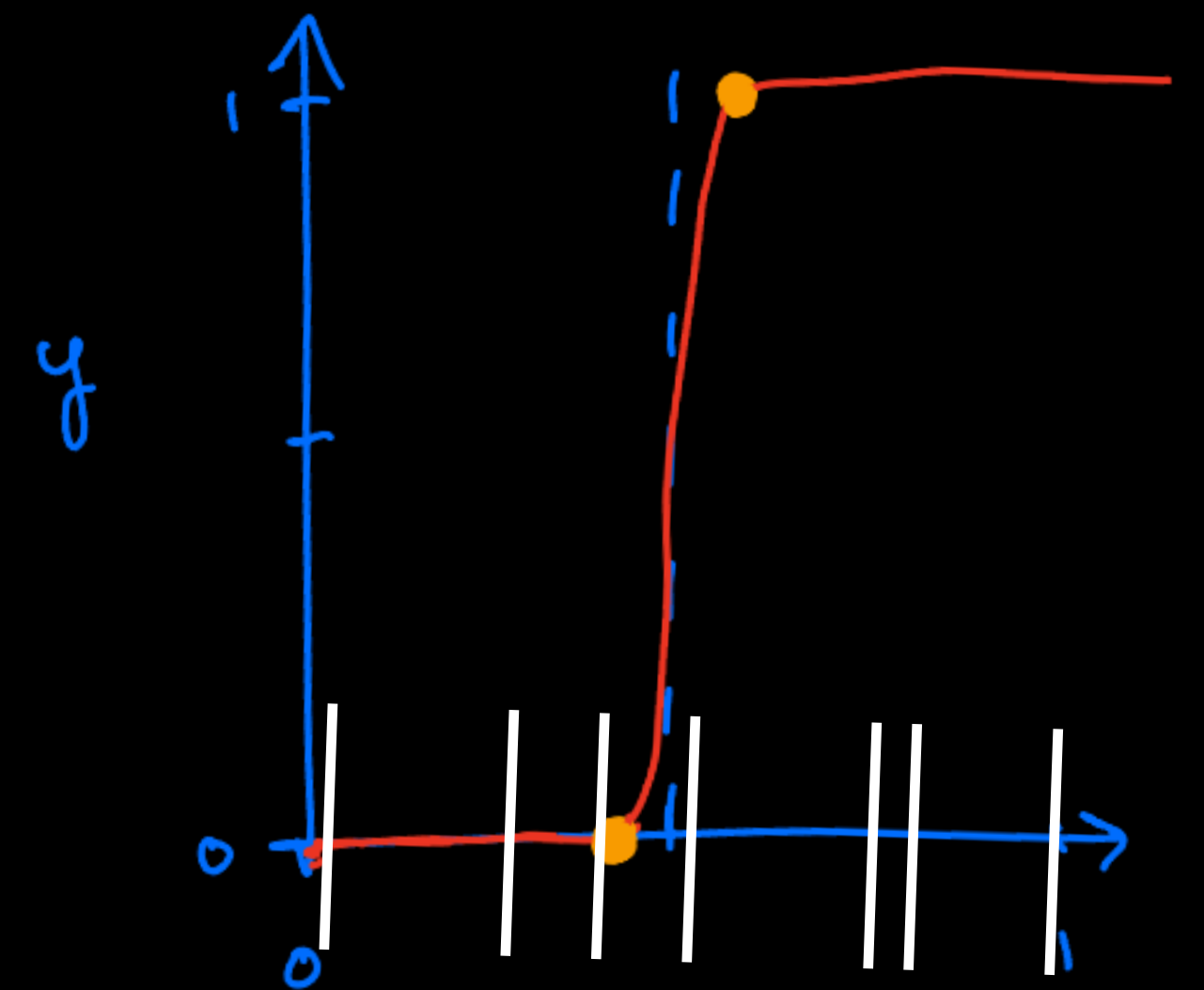
Interval Calibration

$$\text{dCE}(f) \leq \text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$$

Best-possible upper-bound:

$$\text{intCE}(f) := \inf_{\mathcal{I}: \text{Interval partition}} (\text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I}))$$

Can we get a lower-bound?



Interval Calibration

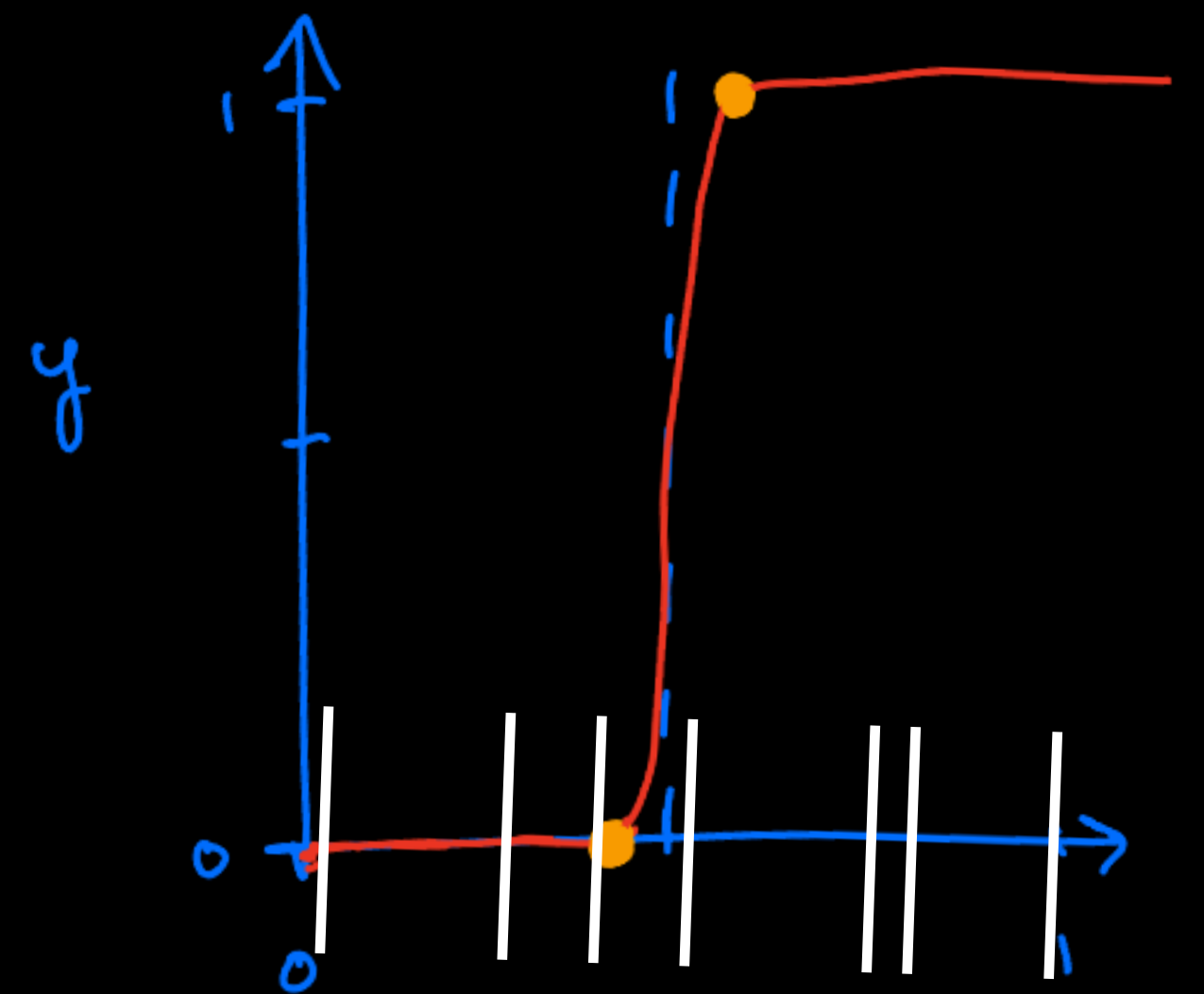
$$\text{dCE}(f) \leq \text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$$

Best-possible upper-bound:

$$\text{intCE}(f) := \inf_{\mathcal{I}: \text{Interval partition}} (\text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I}))$$

Can we get a lower-bound?

$$\frac{1}{16} \text{intCE}(f)^2 \leq \text{dCE}(f) \leq \text{intCE}(f)$$



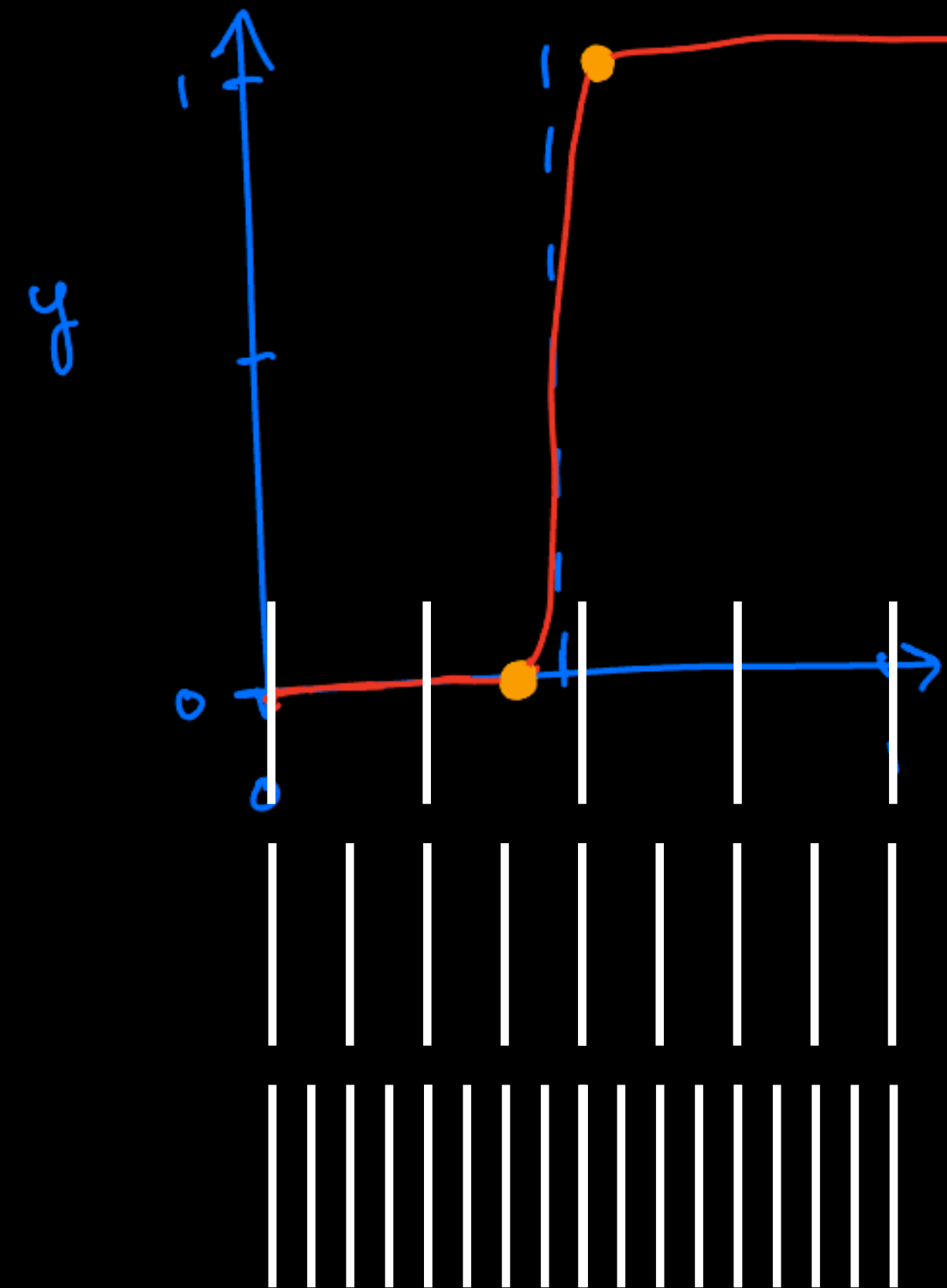
Interval Calibration

$$\text{intCE}(f) := \inf_{\mathcal{I}: \text{Interval partition}} (\text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I}))$$

Computationally, sufficient to minimize over $i \in \mathbb{N}$:

1. Construct regular intervals of width = 2^{-i}
2. Randomly shift intervals (together)
3. Compute $\text{binnedECE}(f, \mathcal{I}) + \text{width}(\mathcal{I})$

This gives same guarantees!



Practical Takeaways

Measure calibration with either:

1. Kernel Calibration Error

2. Interval Calibration Error

or, if you must use binnedECE, add max-interval-width “regularizer”

In Practice: $k\text{CE} \approx \text{binnedECE}$

