# Large models are impossible to regulate.
## CHANGE MY MIND

Augustin Godinot (WIDE)
Université de Rennes, IRISA, INRIA, PEReN

Université de Rennes · Inria · UMR IRISA · GOUVERNEMENT · PEReN Pôle d'Expertise de la Régulation Numérique

## 🏛 What regulators ask for…



- ▶ Digital Services Act (**DSA**): Large platforms induce risks for society, they have to implement risk mitigation meechanisms.
- ▶ Digital Markets Act (**DMA**): Large platforms have a lot of power, we must avoid power imbalance.
- ▶ Artificial Intelligence Act (**AI Act**): limit the use of some algorithms.

## 🕵 ML audit you said ?

- ▶ **Input space** $\mathcal{X}$. *Example: The space of all possible* $1000 \times 1000$ *images.*
- ▶ **Hypothesis** $h : \mathcal{X} \to \{0,1\}$. *Example: a deep neural network.*
- ▶ **Hypothesis class** $\mathcal{H} \subset \{0,1\}^{\mathcal{X}}$. *Example: all the ResNet models with* 50 *blocks.*

Audit a parity metric

$$\mu(h, S) = \mathbb{P}\big(h(X) = 1 \,|\, X \in S, E\big) - \mathbb{P}\big(h(X) = 1 \,|\, X \in S, \overline{E}\big)$$

Example: make sure that in average, men are not advantaged compared to women by a resume screening algorithm.
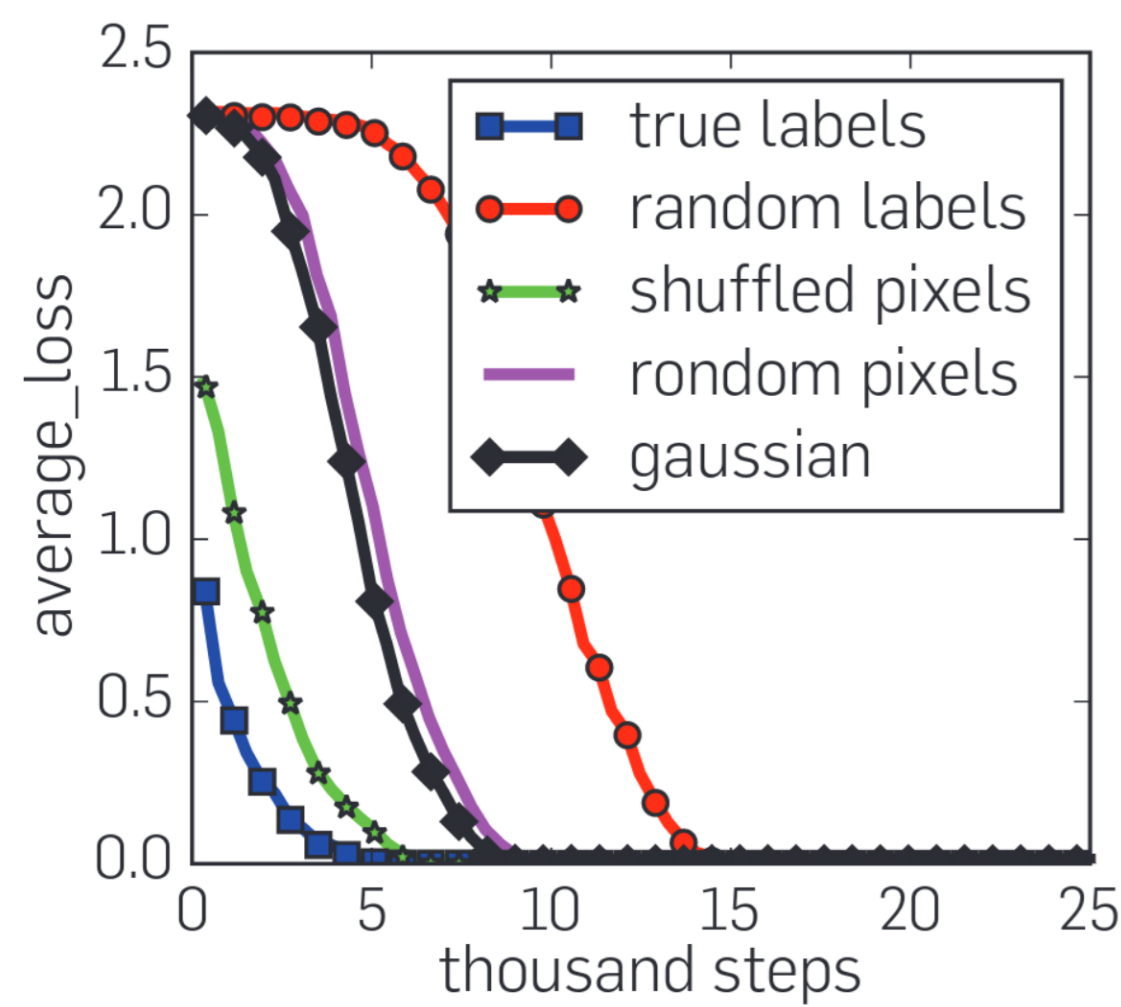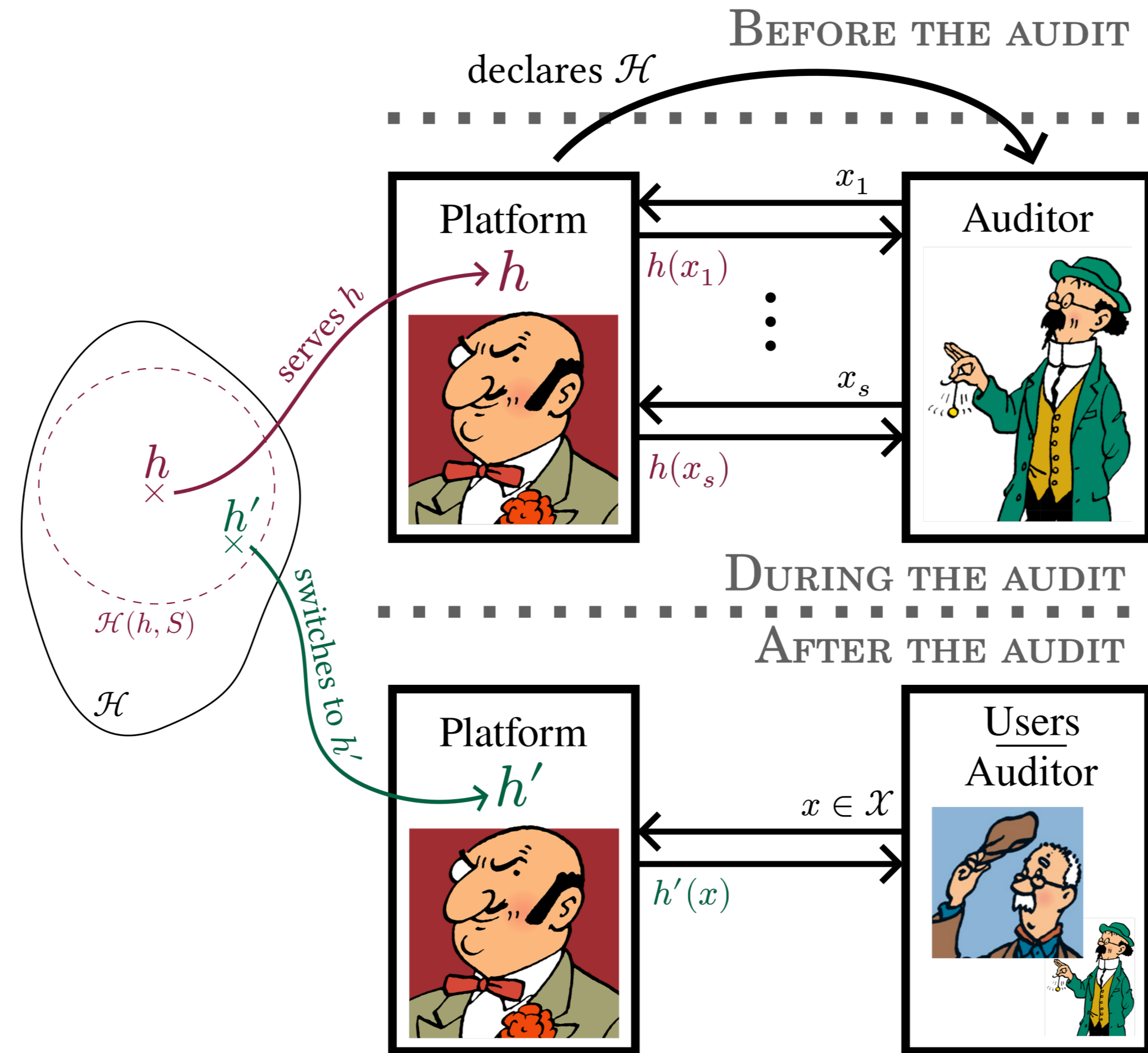
## 🤖 Large Machine Learning models



Figure 1: The training loss of an Inception model trained on CI-FAR10. After enough steps, the loss reaches $0$ *even when trained on random labels.*

Taken from *Understanding deep learning requires rethinking generalization* (Zang et al, CACM 2021)

- ▶ Current ML models can reach **billions of parameters**.
- ▶ Current ML models can **overfit the train data** and have **good generalization** properties.
- ▶ Some explanation attempts: **benign overfitting** and **double descent**.

## ⚠️ Threat model



BEFORE THE AUDIT

declares $\mathcal{H}$

Platform $h$ — serves $h$

$x_1$ ... $x_s$

$h(x_1)$ ... $h(x_s)$

Auditor

$h$ ×  $h'$ ×

$\mathcal{H}(h, S)$

$\mathcal{H}$

switches to $h'$

DURING THE AUDIT

AFTER THE AUDIT

Platform $h'$

$x \in \mathcal{X}$

$h'(x)$

Users Auditor

## 📏 Mesuring the effect of potential manipulations



$$\mathcal{H}(S, h) = \{h' \in \mathcal{H} : \forall x \in S, h'(x) = h(x)\}$$

$$\mathrm{diam}_\mu \;\blacksquare\; = \max_{h' \in \blacksquare} |\mu(h') - \mu(h)|$$

$\mu$-diameter · Version space

## 🔍 Impossibility theorem

### Definition 2: Benign overfitting on $c$

$\mathcal{H}$ exhibits benign overfitting with respect to $c$ iif fhere exists $d_0 \in \mathbb{N}_*$ and $\varepsilon \in [0, 1)$ such that $\forall d \le d_0, S \in \mathcal{X}, \sigma \in \{0,1\}^d$,

$$\exists h \in \mathcal{H}, \begin{cases} \forall x_i \in S, h(x_i) = \sigma_i \text{ (fits any train set)} \\ \mathbb{P}\big(h(X) = c(X) \,\big|\, X \in \overline{S}\big) = 1 - \varepsilon \text{ (low error)} \end{cases}$$

### Theorem: Better than random? No can do.

If $\mathcal{H}$ exhibits benign overfitting with respect to the sensitive attribute, then,

$$\forall S, |S| = |S_{\mathrm{random}}|, \quad \mathrm{diam}_\mu(h, S) = \mathrm{diam}_\mu(h, S_{\mathrm{random}})$$

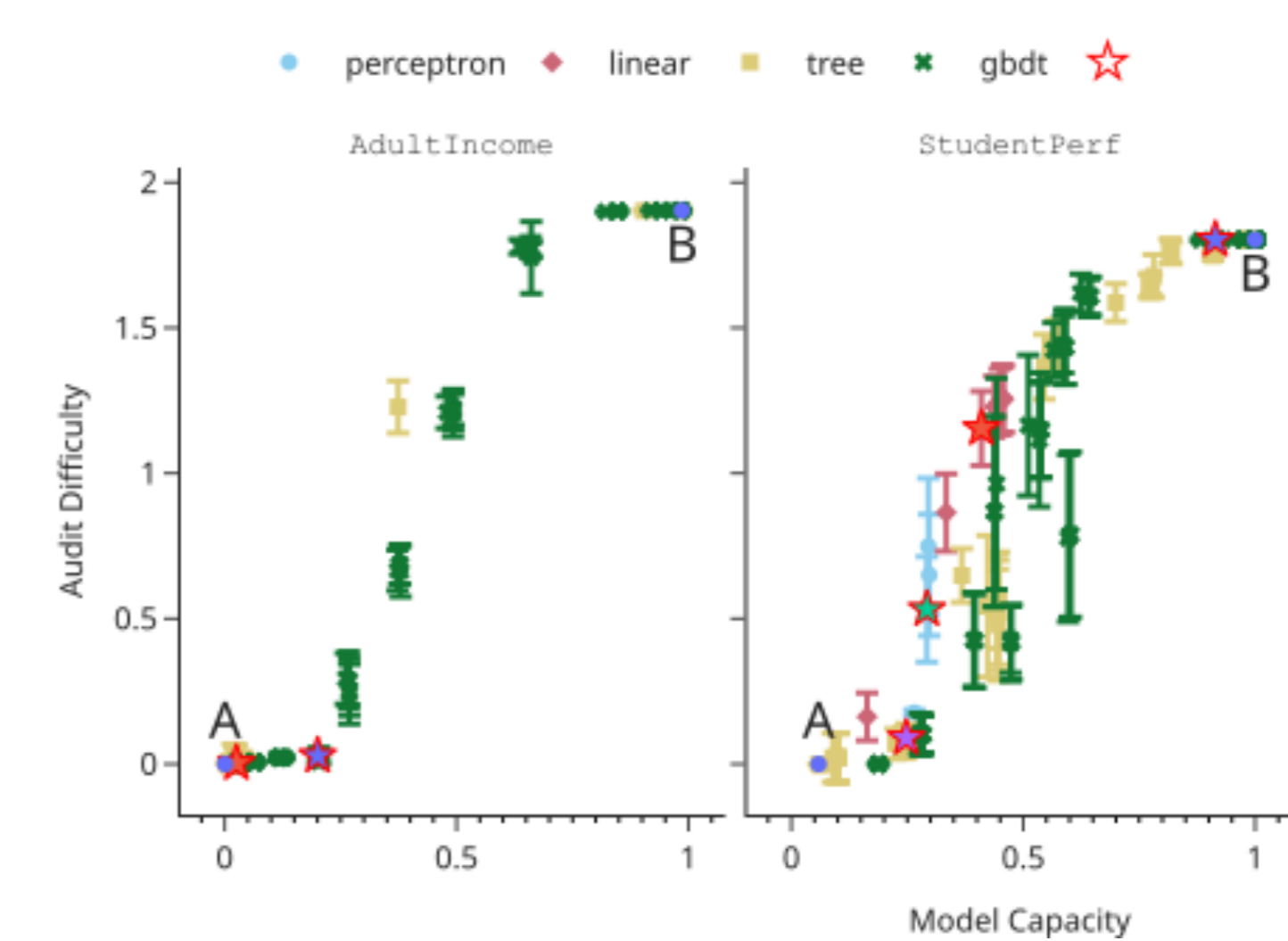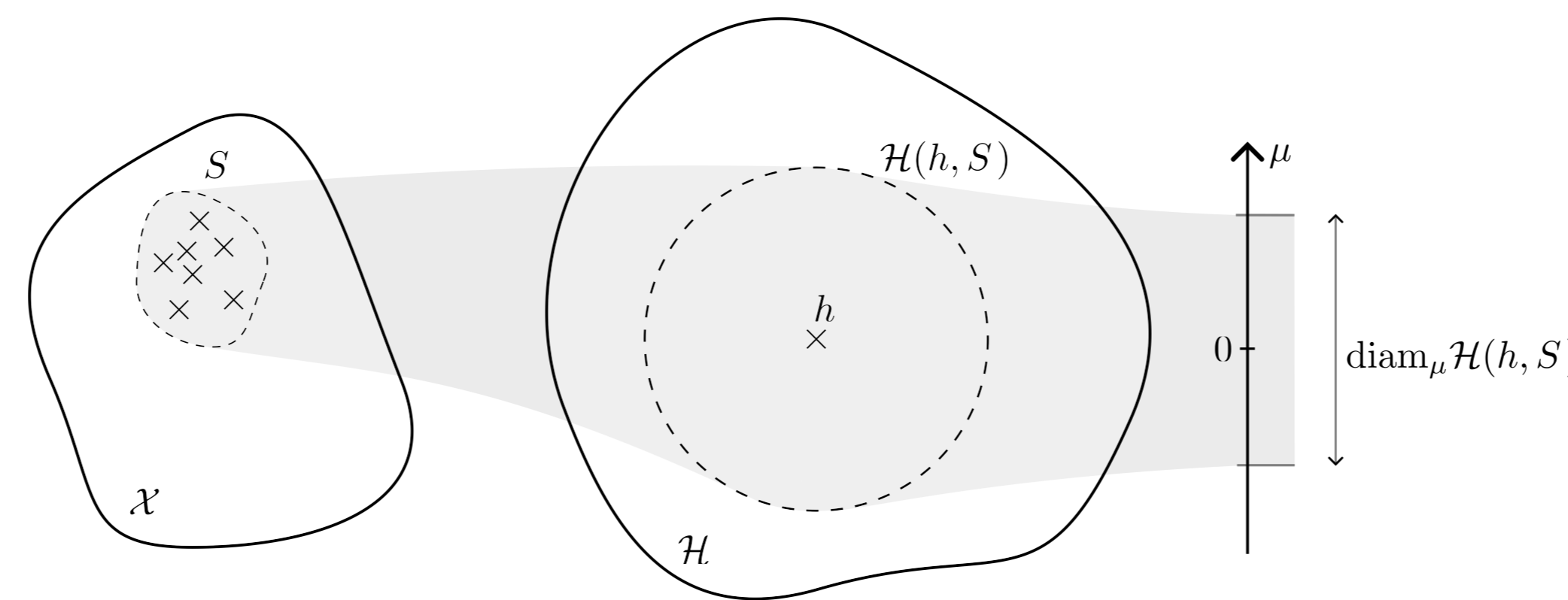## 📊 And in practice ?



perceptron · linear · tree · gbdt

Figure 2: The value of the $\mu$-diameter with respect to the Rademacher complexity of the hypothesis class. Informally: hypothesis class = fixed architeecture + hyperparameters.



Figure 3: What is the accuracy cost for a platform to evade an audit? Not much. Let $\mathcal{F} = \big(\mathcal{H}_1, ..., \mathcal{H}_f\big)$ be a family of hypothesis classes. *Example: all the decision trees with varying maximum depth.*

- ▶ $\mathcal{H}^* \in \mathcal{F}$ with best test accuracy.
- ▶ $\mathcal{H}_{\mathrm{evade}} \in \mathcal{F}$ with largest $\mu$-diameter.

$$\mathrm{CostOfExhaustion}(\mathcal{F}) = \mathrm{Accuracy}(\mathcal{H}^*) - \mathrm{Accuracy}(\mathcal{H}_{\mathrm{evade}})$$