



**PEReN**  
Pôle d'Expertise de la  
Régulation Numérique

# Manipulation-proof auditing

*Under manipulations, are there models harder to audit?*

PrivSec seminar · April. 7th 2024



Augustin Godinot



Erwan Le Merrer



Gilles Tredan



Camilla Penzo



François Taïani

# A first example

Qty: 1 ▾

**\$204.60** + Free Shipping  
In stock. Sold by **-MOTOPILEX-**

 Add to Cart

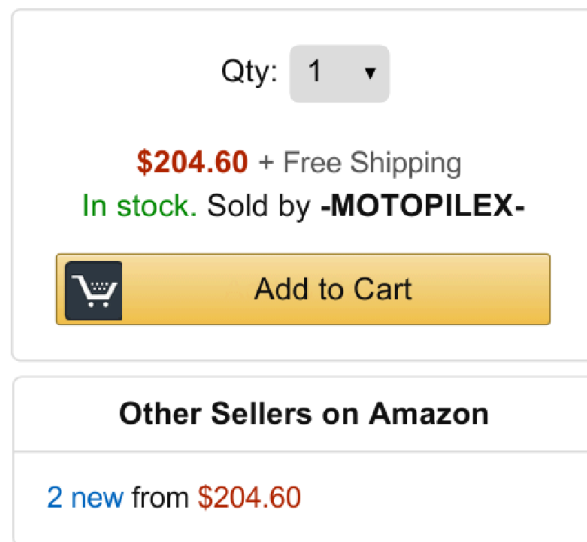
**Other Sellers on Amazon**

2 new from **\$204.60**

**Metric** Demographic parity  
between amazon and the other  
sellers




# A first example



Qty: 1 ▾

**\$204.60** + Free Shipping

In stock. Sold by **-MOTOPILEX-**

 Add to Cart

**Other Sellers on Amazon**

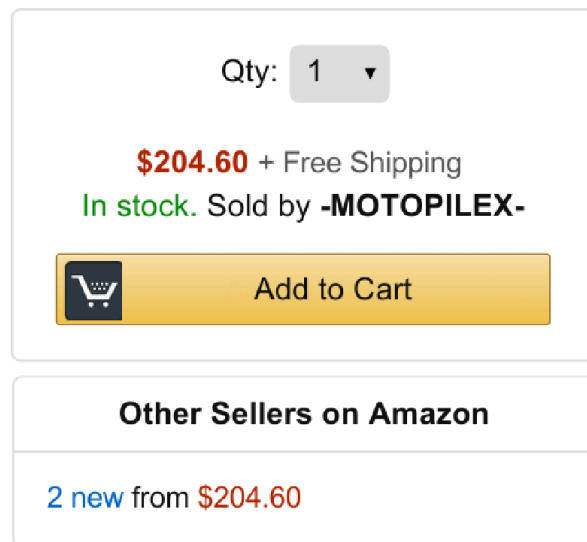
2 new from **\$204.60**

**Metric** Demographic parity  
between amazon and the other  
sellers

**Audit queries** Top- $k$  best selling  
products




# A first example



Qty: 1 ▾

**\$204.60** + Free Shipping

In stock. Sold by **-MOTOPILEX-**

 Add to Cart

**Other Sellers on Amazon**

2 new from **\$204.60**

**Metric** Demographic parity  
between amazon and the other  
sellers

**Audit queries** Top- $k$  best selling  
products

**Data collection** shameless scraping



# In this talk

## Context

How are audits currently conducted?

## Framework

What do we mean by robust auditing: manipulation-proofness.



# In this talk

## Context

How are audits currently conducted?

## Framework

What do we mean by robust auditing: manipulation-proofness.

## A theoretical peek

**Our contributions**

Large models cannot be audited more efficiently than by random sampling.

## Empirical study

In practice, the cost to evade a black-box audit is mild.



# In this talk

## Context

How are audits currently conducted?

## Framework

What do we mean by robust auditing: manipulation-proofness.

## A theoretical peek

**Our contributions**

Large models cannot be audited more efficiently than by random sampling.

## Empirical study

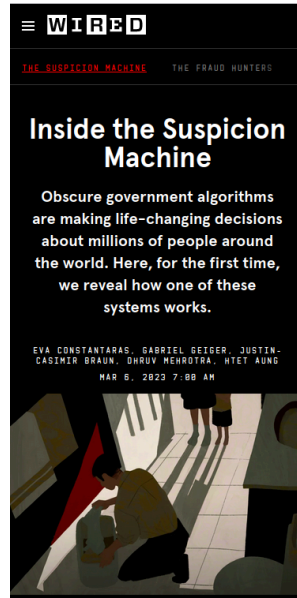
In practice, the cost to evade a black-box audit is mild.

## Concluding remarks

The implications for AI regulation.



# Algorithm audits



Qty: 1 ▼

**\$204.60** + Free Shipping  
In stock. Sold by **MOTOPILEX-**

 Add to Cart

**Other Sellers on Amazon**

2 new from **\$204.60**




**HIRING PLATFORM**

**Fast. Fair. Flexible.**  
Finally, hiring technology that works how you want it to.

HireVue is a talent experience platform designed to automate workflows and make scaling hiring easy. Improve how you engage, screen and hire talent with text recruiting, assessments, and video interviewing software.

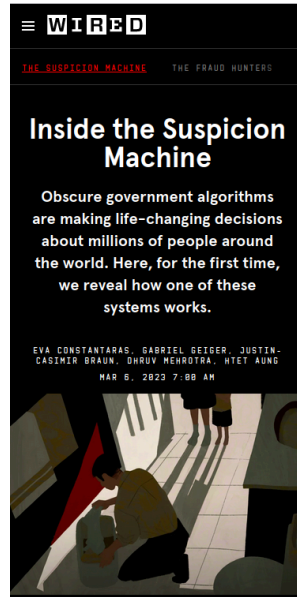
*Hirevue claims it is "Fast. Fair. Flexible."*

## Context

-  Framework
-  A theoretical peek
-  Empirical study
- Concluding remarks




# Algorithm audits



Qty: 1 ▼

**\$204.60** + Free Shipping  
In stock. Sold by **-MOTOPILEX-**

 Add to Cart

**Other Sellers on Amazon**

2 new from **\$204.60**




**HIRING PLATFORM**

**Fast. Fair. Flexible.**  
Finally, hiring technology that works how you want it to.

HireVue is a talent experience platform designed to automate workflows and make scaling hiring easy. Improve how you engage, screen and hire talent with text recruiting, assessments, and video interviewing software.

*Hirevue claims it is "Fast. Fair. Flexible."*

## Context

-  Framework
-  A theoretical peek
-  Empirical study
- Concluding remarks

A screenshot of the European Parliament News website. The header includes the European Parliament logo and the word "News". Below the header is a navigation bar with links for "Headlines", "Press room", "Agenda", "FAQ", and "Election Press Kit". The main content area shows a headline: "EU AI Act: first regulation on artificial intelligence" with a sub-headline "intelligence". It also includes the text "Society Updated: 14-06-2023 - 14:06" and "Created: 08-06-2023 - 11:40".

J. Dastin, L. Chen, A. Mislove, and C. Wilson, , J. Larson, S. Mattu, L. Kirchner, and J. Angwin, Rédaction



# Audit steps

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks

## Choosing the metric

- ▶ FairML book [6]
- ▶ Political implications of the metric [7]
- ▶ Data minimization [8]
- ▶ Privacy auditing [9]

## Choosing the queries

- ▶ Classical random sampling [10]
- ▶ Crafted datasets
- ▶ Active learning [11]
- ▶ Fairness by betting [12]

## Data collection

- ▶ Do we get explanations? [13], [14]
- ▶ Do we have access to private API? [15]
- ▶ What if the platform lies? [11]



# Audit steps

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks

## Choosing the metric

- ▶ FairML book [6]
- ▶ Political implications of the metric [7]
- ▶ Data minimization [8]
- ▶ Privacy auditing [9]

## Choosing the queries

- ▶ Classical random sampling [10]
- ▶ Crafted datasets
- ▶ Active learning [11]
- ▶ Fairness by betting [12]

## Data collection

- ▶ Do we get explanations? [13], [14]
- ▶ Do we have access to private API? [15]
- ▶ What if the platform lies? [11] **this talk**



**my PhD topic:** how to design audit methods that are not easily gamed by platforms?

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks



**my PhD topic:** how to design audit methods that are not easily gamed by platforms?

- ▶ detecting lies, robust auditing, robust estimation  $\Leftrightarrow$  comparing the observations to a prior.

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks



**my PhD topic:** how to design audit methods that are not easily gamed by platforms?

- ▶ detecting lies, robust auditing, robust estimation  $\Leftrightarrow$  comparing the observations to a prior.
- ▶ which prior to choose and how does it impact the audit process in practice?  
 $\Rightarrow$  This talk: one choice of prior

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks



**my PhD topic:** how to design audit methods that are not easily gamed by platforms?

- ▶ detecting lies, robust auditing, robust estimation  $\Leftrightarrow$  comparing the observations to a prior.
- ▶ which prior to choose and how does it impact the audit process in practice?  
 $\Rightarrow$  This talk: one choice of prior
- ▶ how to verify that the prior holds ?  
 $\Rightarrow$  Model change detection

## Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks



# Manipulation- proof auditing

*Threat model*

Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{€} | X \in S, \text{amazon}) - \mathbb{P}(\text{€} | X \in S, \text{france})$

**Hypothesis space**

$\mathcal{H} \subset \{0, 1\}^x, h \in \mathcal{H}$

[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# Manipulation-proof auditing

*Threat model*

Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{€} | X \in S, \text{a}) - \mathbb{P}(\text{€} | X \in S, \text{f})$

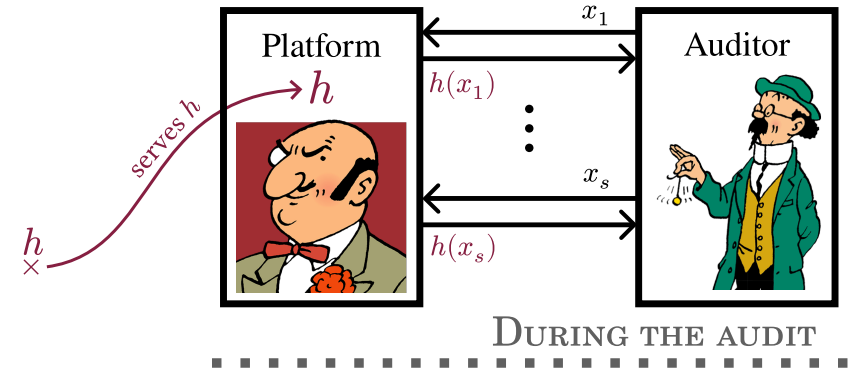
**Hypothesis space**

$\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}, h \in \mathcal{H}$

**Audit set**

$S = (x_1, \dots, x_s) \subset \mathcal{X}$

$Y = (h(x_1), \dots, h(x_s))$



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# Manipulation-proof auditing

*Threat model*

Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{€} | X \in S, \text{🇺🇸}) - \mathbb{P}(\text{€} | X \in S, \text{🇫🇷})$

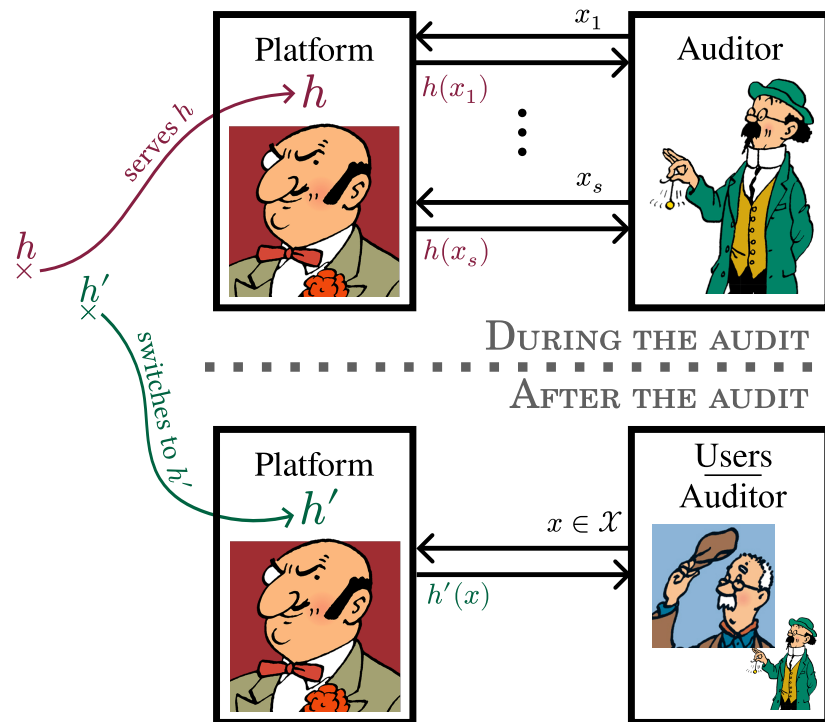
**Hypothesis space**

$$\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}, h \in \mathcal{H}$$

**Audit set**

$$S = (x_1, \dots, x_s) \subset \mathcal{X}$$

$$Y = (h(x_1), \dots, h(x_s))$$



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# Manipulation-proof auditing

Threat model

Context

Framework

A theoretical peek

Empirical study

Concluding remarks

$$\text{Audit metric } \mu(h, S) = \mathbb{P}(\text{€} \mid X \in S, \text{🇺🇸}) - \mathbb{P}(\text{€} \mid X \in S, \text{🇫🇷})$$

Hypothesis space

$$\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}, h \in \mathcal{H}$$

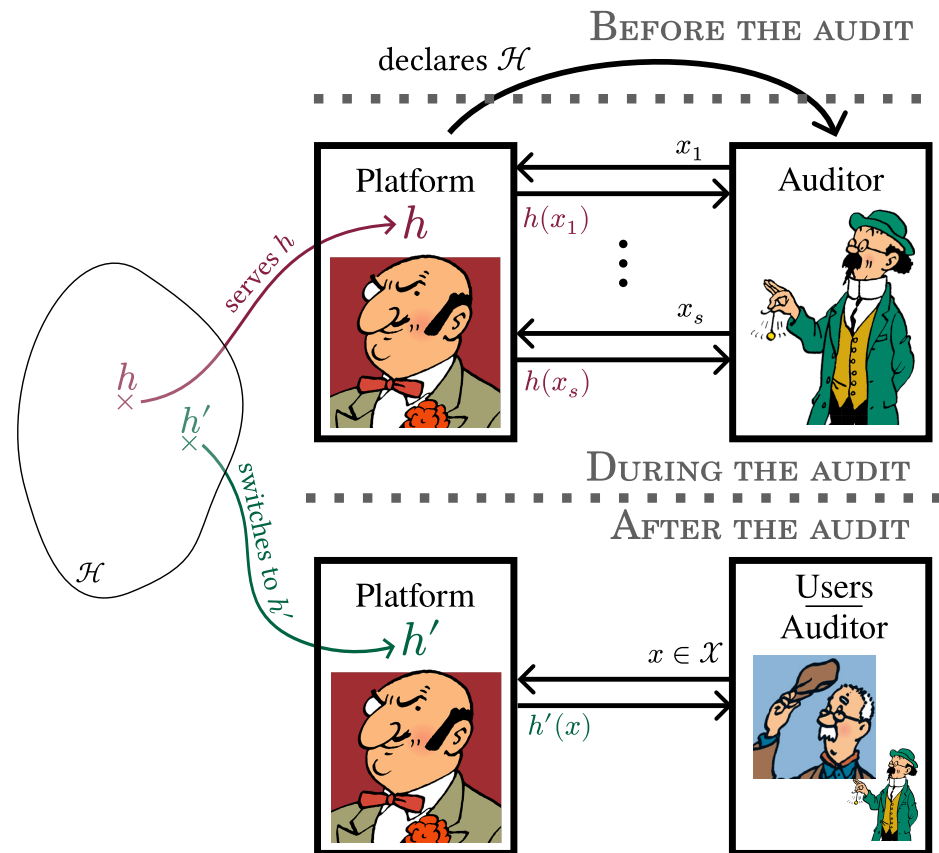
Audit set

$$S = (x_1, \dots, x_s) \subset \mathcal{X}$$

$$Y = (h(x_1), \dots, h(x_s))$$

Assumptions

1. *Auditor prior*:  $\mathcal{H}$  is known



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# Manipulation-proof auditing

Threat model

Context

Framework

A theoretical peek

Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{€} | X \in S, \text{a}) - \mathbb{P}(\text{€} | X \in S, \text{f})$

**Hypothesis space**

$$\mathcal{H} \subset \{0, 1\}^{\mathcal{X}}, h \in \mathcal{H}$$

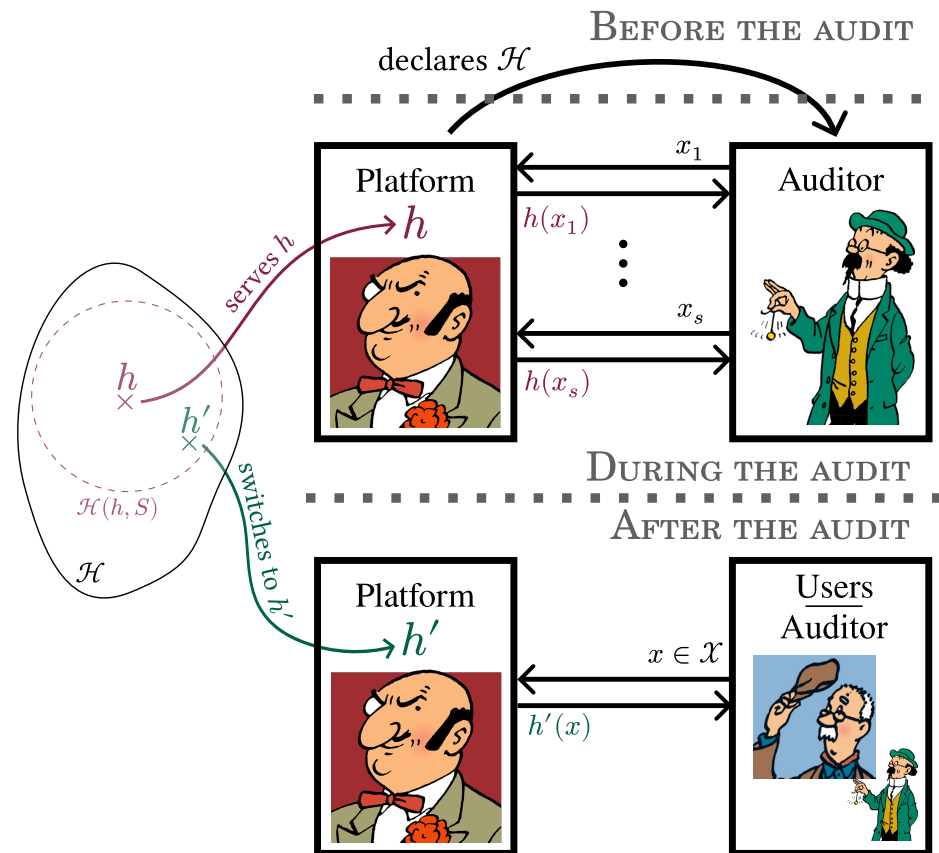
**Audit set**

$$S = (x_1, \dots, x_s) \subset \mathcal{X}$$

$$Y = (h(x_1), \dots, h(x_s))$$

**Assumptions**

1. *Auditor prior*:  $\mathcal{H}$  is known
2. *Self-consistency*: once platform reveals its labeling of  $x$ , cannot change it.



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# Manipulation-proof auditing

*Audit manipulability*

Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

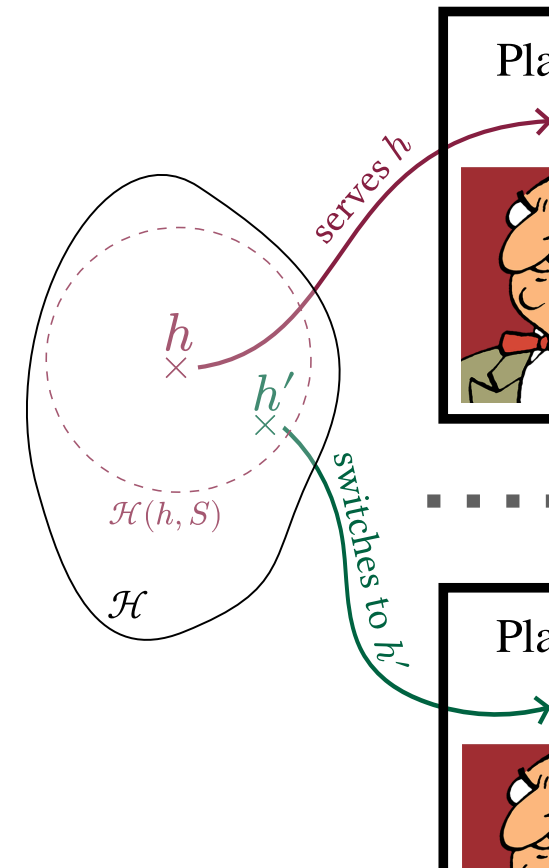
How much can  $\mu$  change after the audit?

**Measured** during  
the audit  
 $\mu(h)$

**True value** after  
the audit  
 $\mu(h')$

**Audit manipulability**

$$|\mu(h) - \mu(h')| \leq ?$$



# Manipulation-proof auditing

*Audit manipulability*

Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

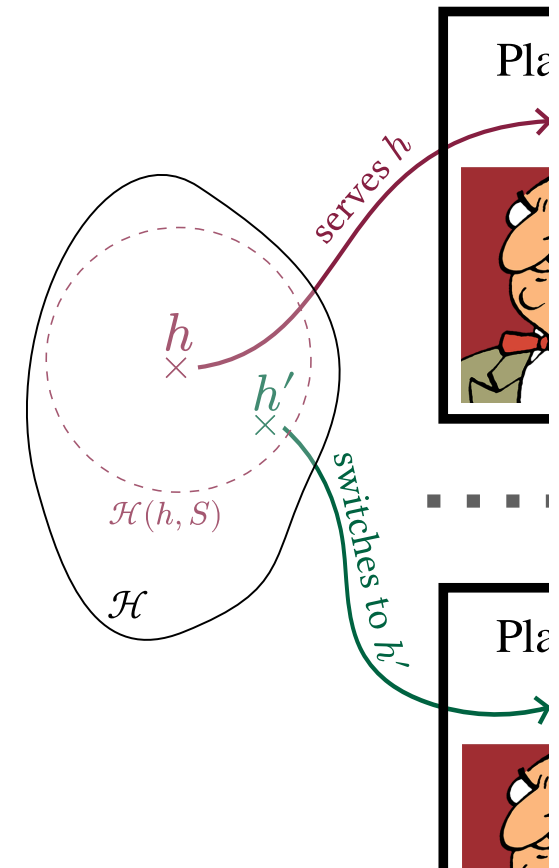
How much can  $\mu$  change after the audit?

**Measured** during  
the audit  
 $\mu(h)$

**True value** after  
the audit  
 $\mu(h')$

**Audit manipulability**

$$|\mu(h) - \mu(h')| \leq \text{diam}_\mu \mathcal{H}(S, Y)$$



# Manipulation-proof auditing

## Audit manipulability

Context

 Framework

 A theoretical peek

 Empirical study

Concluding remarks

How much can  $\mu$  change after the audit?

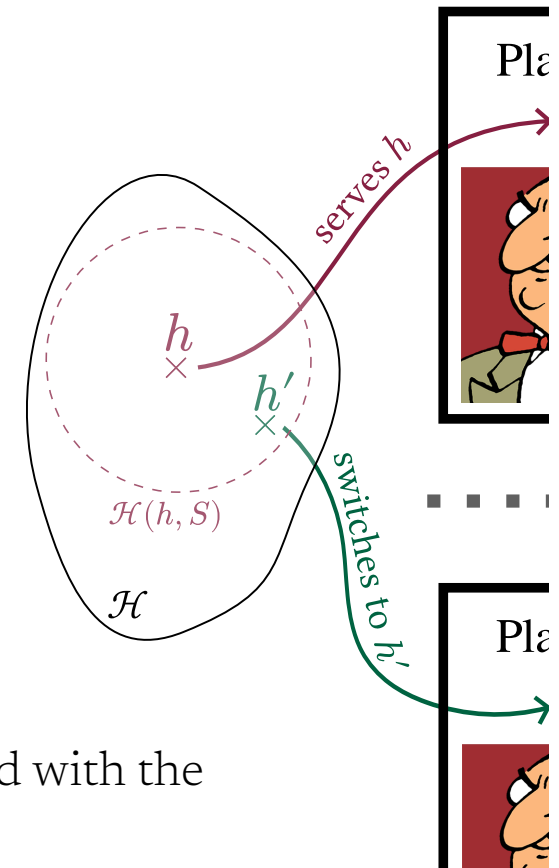
**Measured** during  
the audit  
 $\mu(h)$

**True value** after  
the audit  
 $\mu(h')$

**Audit manipulability**

$$|\mu(h) - \mu(h')| \leq \text{diam}_\mu \mathcal{H}(S, Y)$$

diameter of the version space (equipped with the pseudo metric  $h, h' \mapsto |\mu(h) - \mu(h')|$ )



# Manipulation- proof auditing

Version

Space 🌟

Context

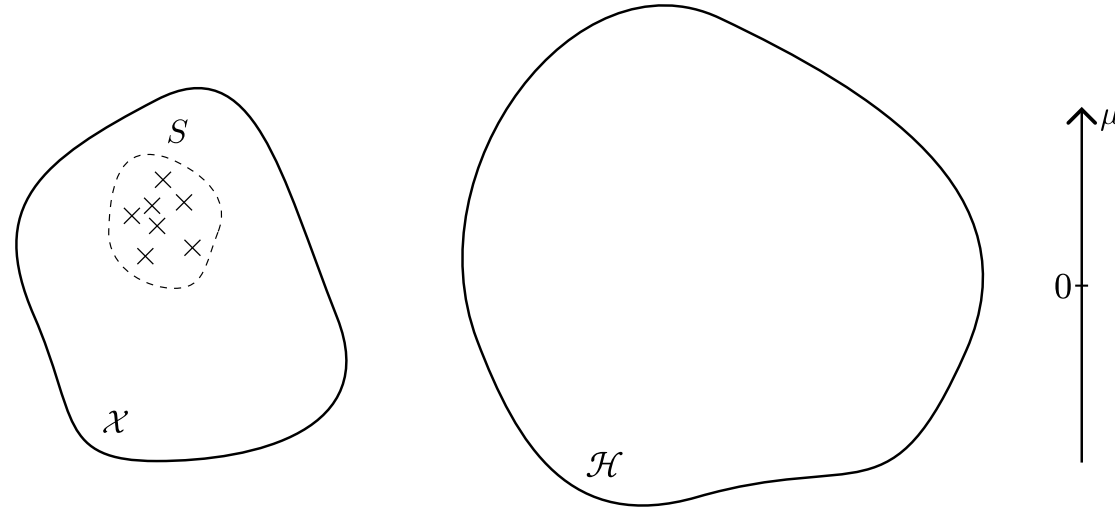
🔨 **Framework**

🔍 A theoretical peek

📊 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷}) - \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷})$



# Manipulation- proof auditing

Version

Space 🌟

Context

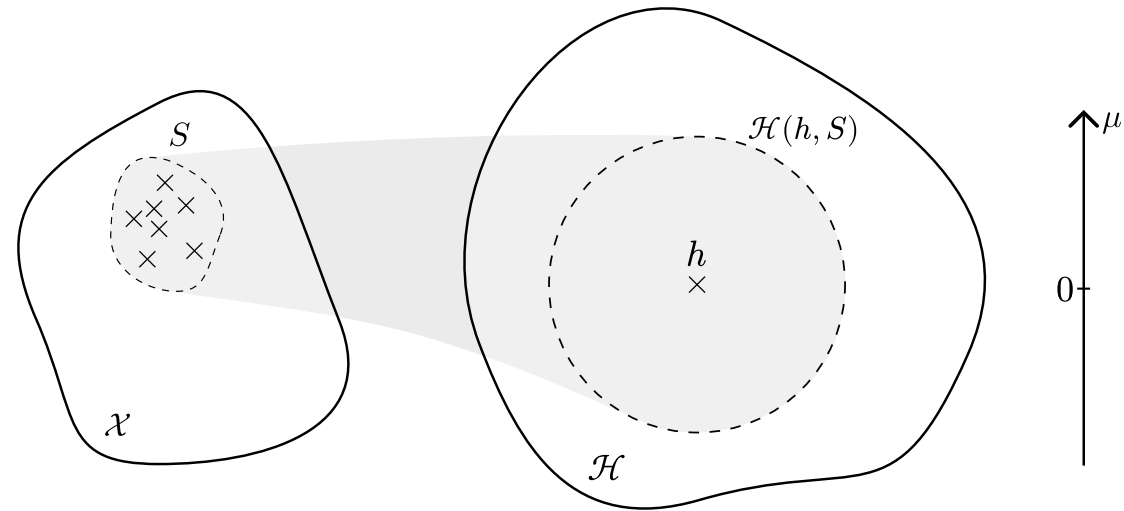
🔨 **Framework**

🔍 A theoretical peek

📊 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷}) - \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷})$



# Manipulation-proof auditing

Version

Space 🌟

Context

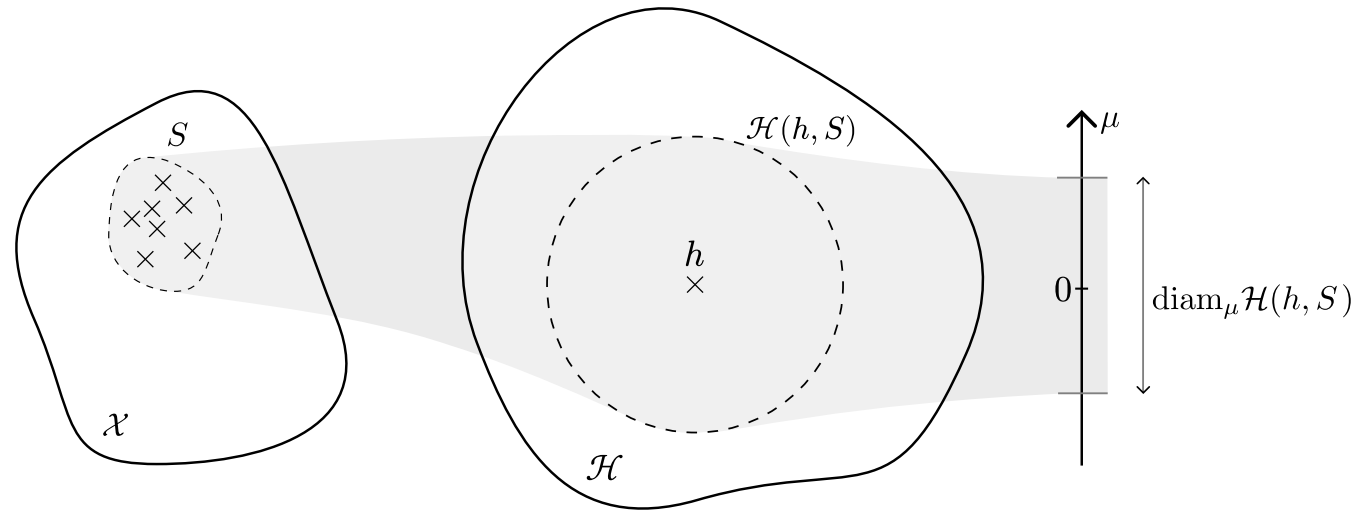
🔧 Framework

🔍 A theoretical peek

📊 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{🇺🇸} \mid X \in S, \text{🗣️}) - \mathbb{P}(\text{🇺🇸} \mid X \in S, \text{🇫🇷})$



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

$$\text{diam}_\mu \mathcal{H}(S, h) = \max_{h' \in \mathcal{H}(S, h)} |\mu(h') - \mu(h)|$$



# Manipulation-proof auditing

Version Space 

Context

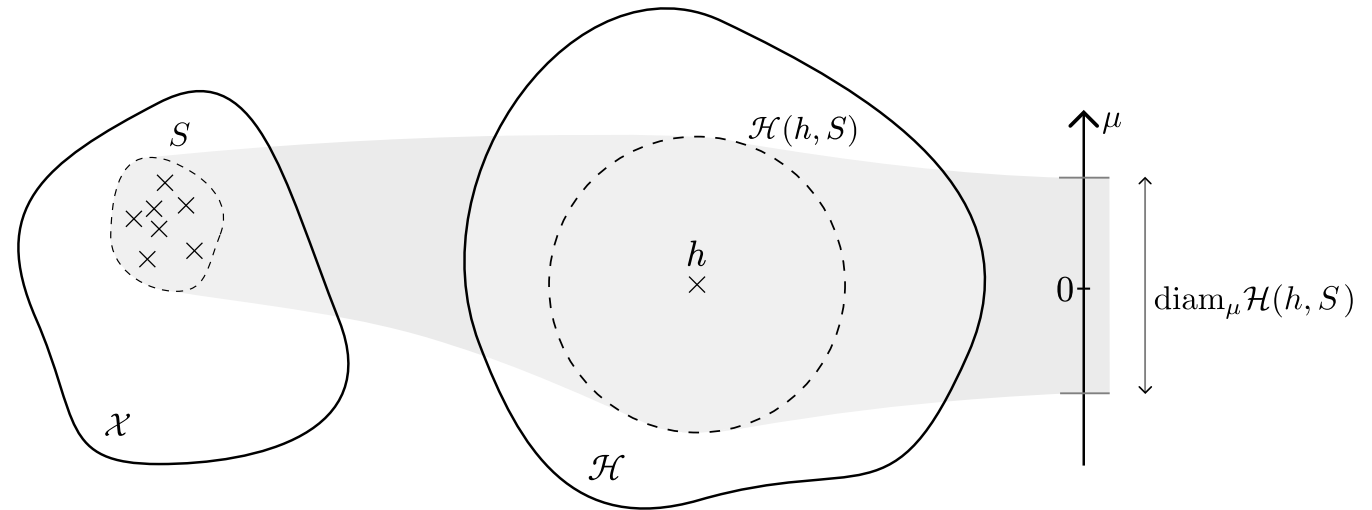
 Framework

 A theoretical peek

 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇯🇵}) - \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷})$



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

$$\text{diam}_\mu \mathcal{H}(S, h) = \max_{h' \in \mathcal{H}(S, h)} |\mu(h') - \mu(h)|$$

Version space



# Manipulation-proof auditing

Version Space 

Context

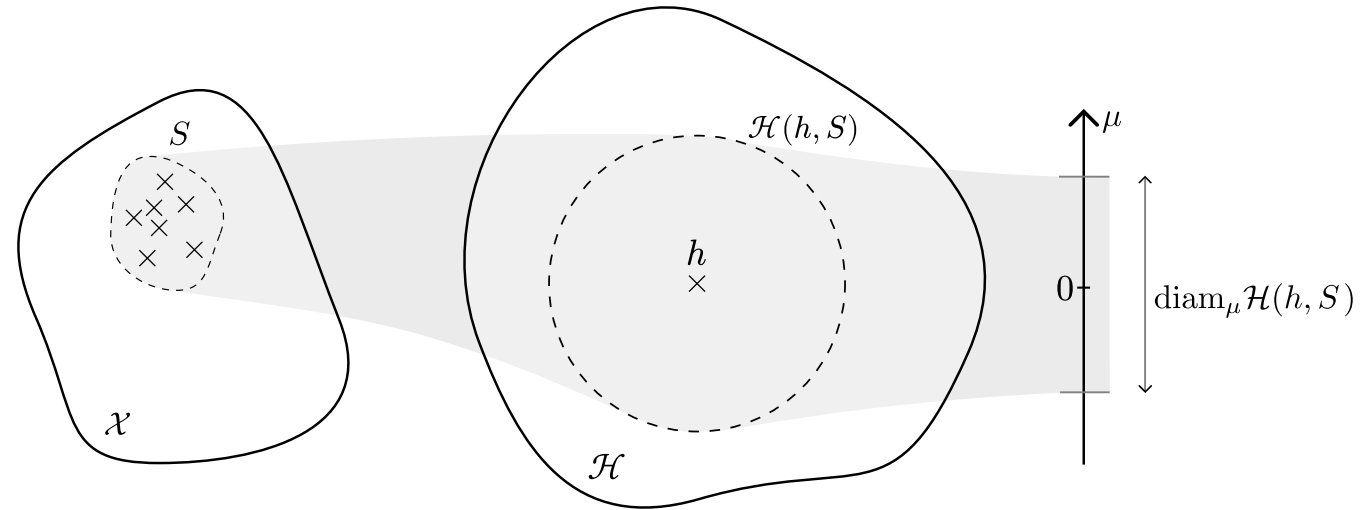
 Framework

 A theoretical peek

 Empirical study

Concluding remarks

**Audit metric**  $\mu(h, S) = \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇯🇵}) - \mathbb{P}(\text{🇺🇸} | X \in S, \text{🇫🇷})$



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

$$\text{diam}_\mu \mathcal{H}(S, h) = \max_{h' \in \mathcal{H}(S, h)} |\mu(h') - \mu(h)|$$

$\mu$ -diameter

Version space



# State of the Art

## *Active auditing algorithms*

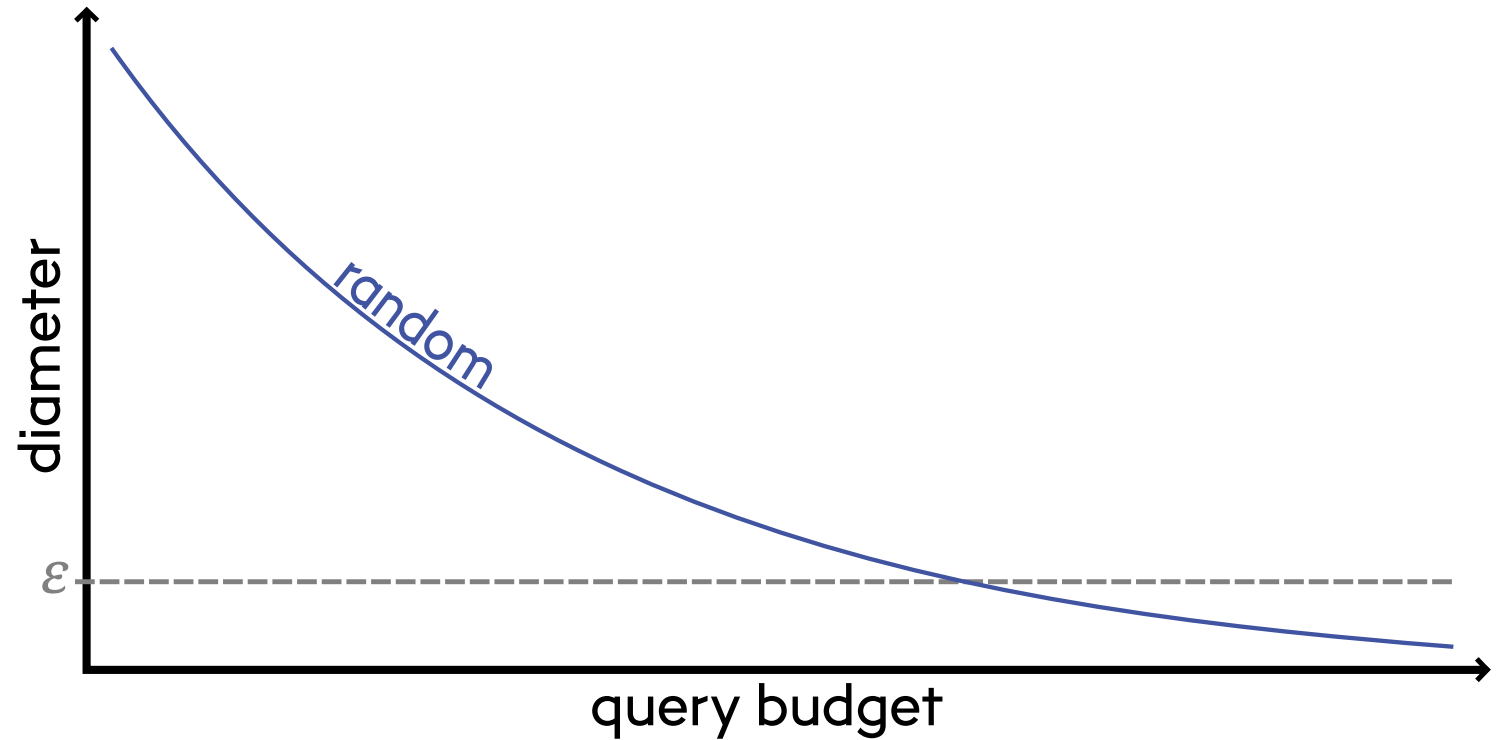
Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# State of the Art

## *Active auditing algorithms*

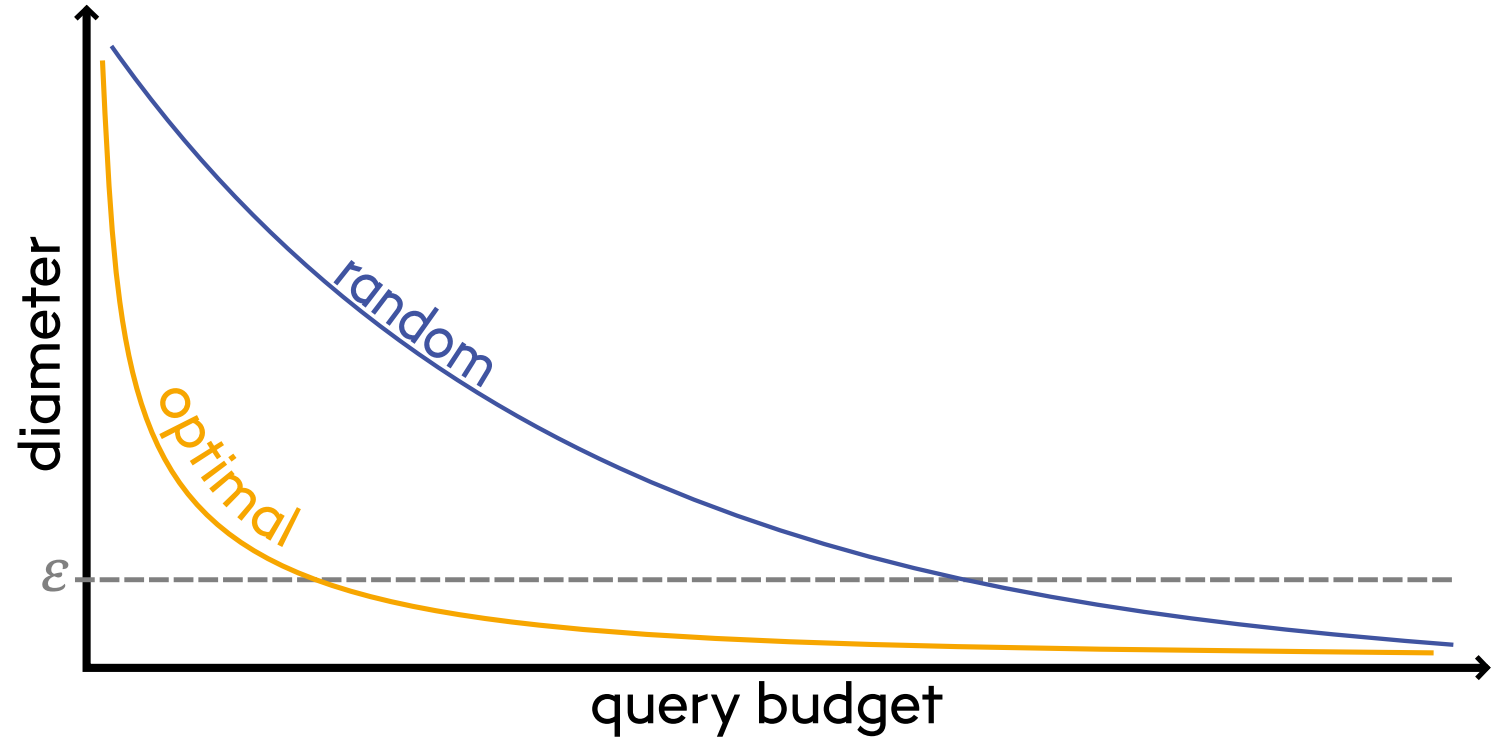
Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# State of the Art

## Active auditing algorithms

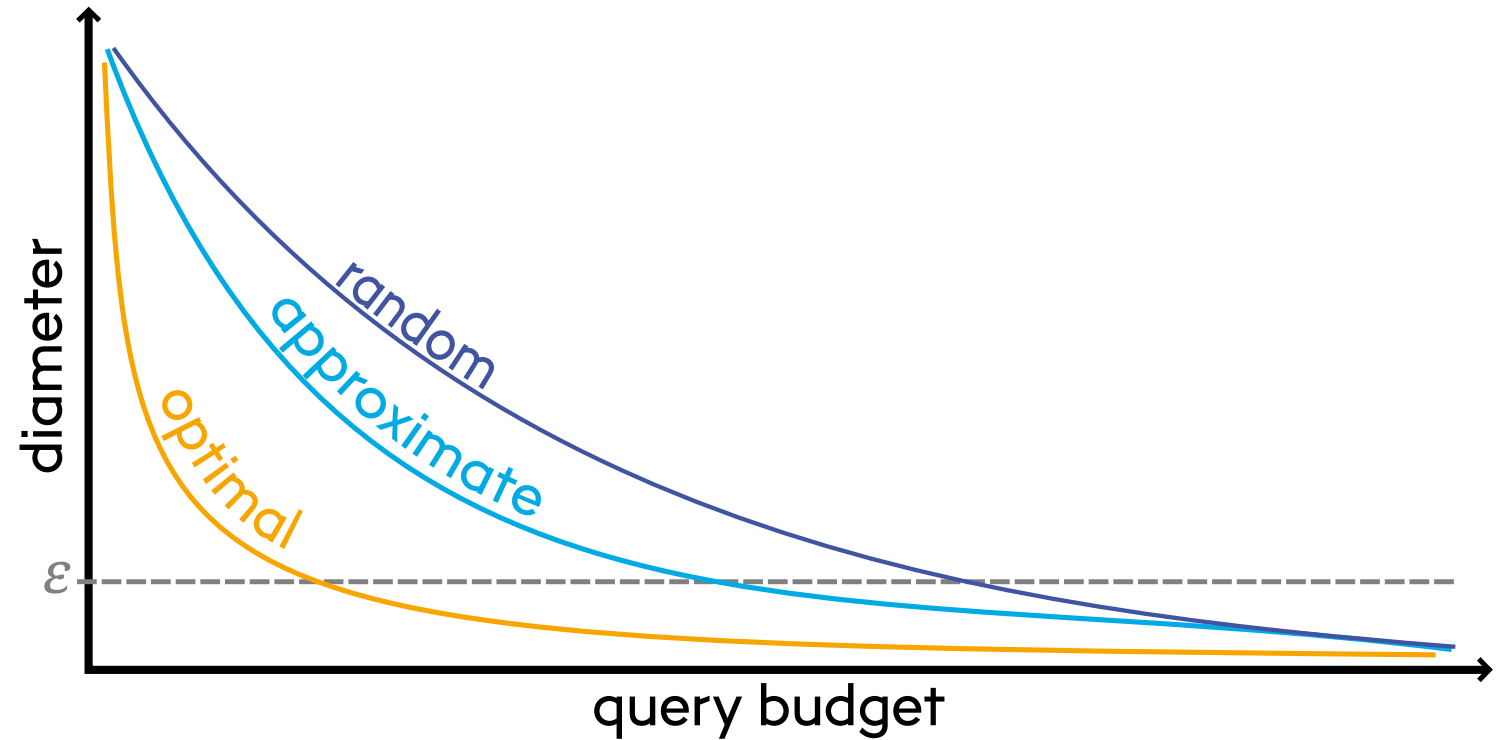
Context

 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# State of the Art

## Active auditing algorithms

Context

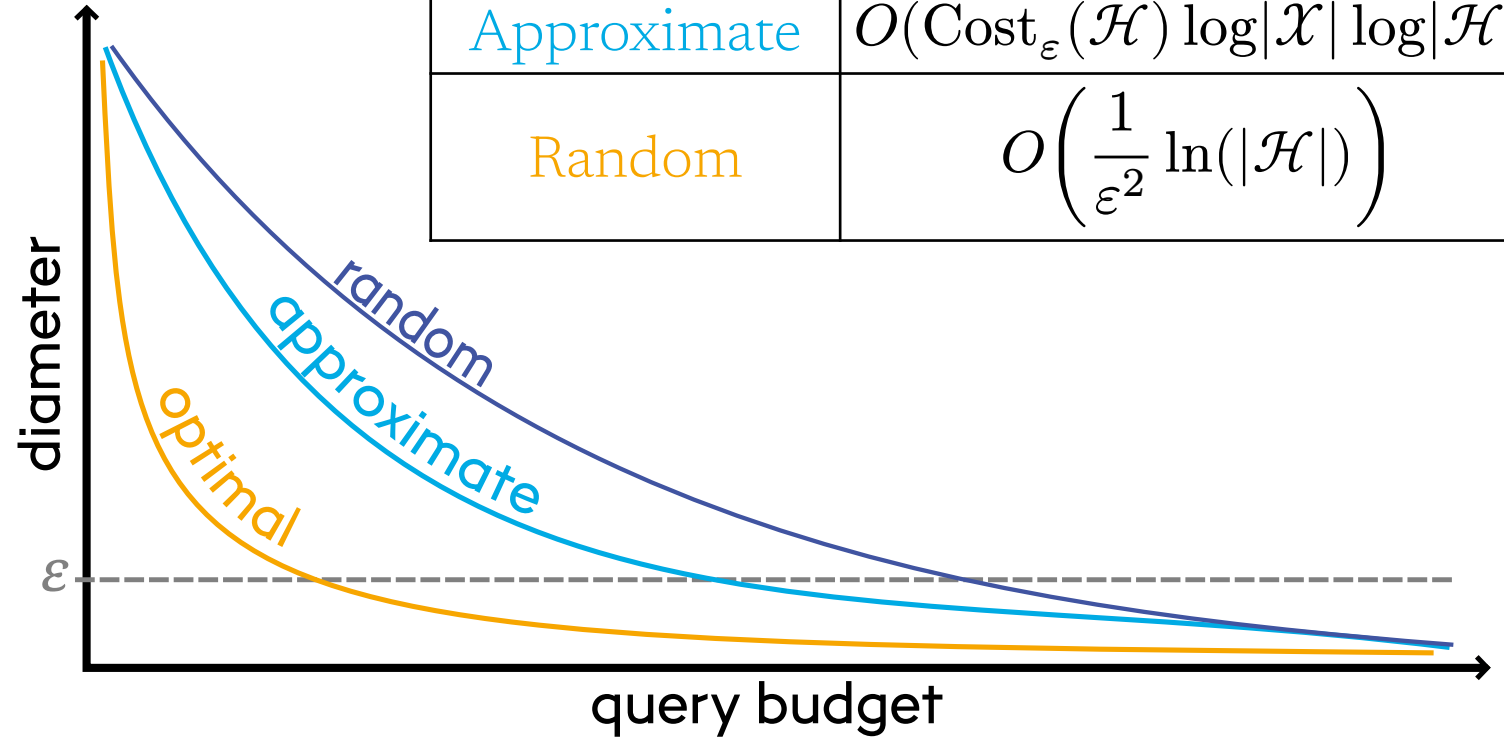
 **Framework**

 A theoretical peek

 Empirical study

Concluding remarks

Audit method	Query complexity
Optimal	$\text{Cost}_\varepsilon(\mathcal{H})$
Approximate	$O(\text{Cost}_\varepsilon(\mathcal{H}) \log \mathcal{X}  \log \mathcal{H} )$
Random	$O\left(\frac{1}{\varepsilon^2} \ln( \mathcal{H} )\right)$



[11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# State of the Art

## Active auditing algorithms

Context

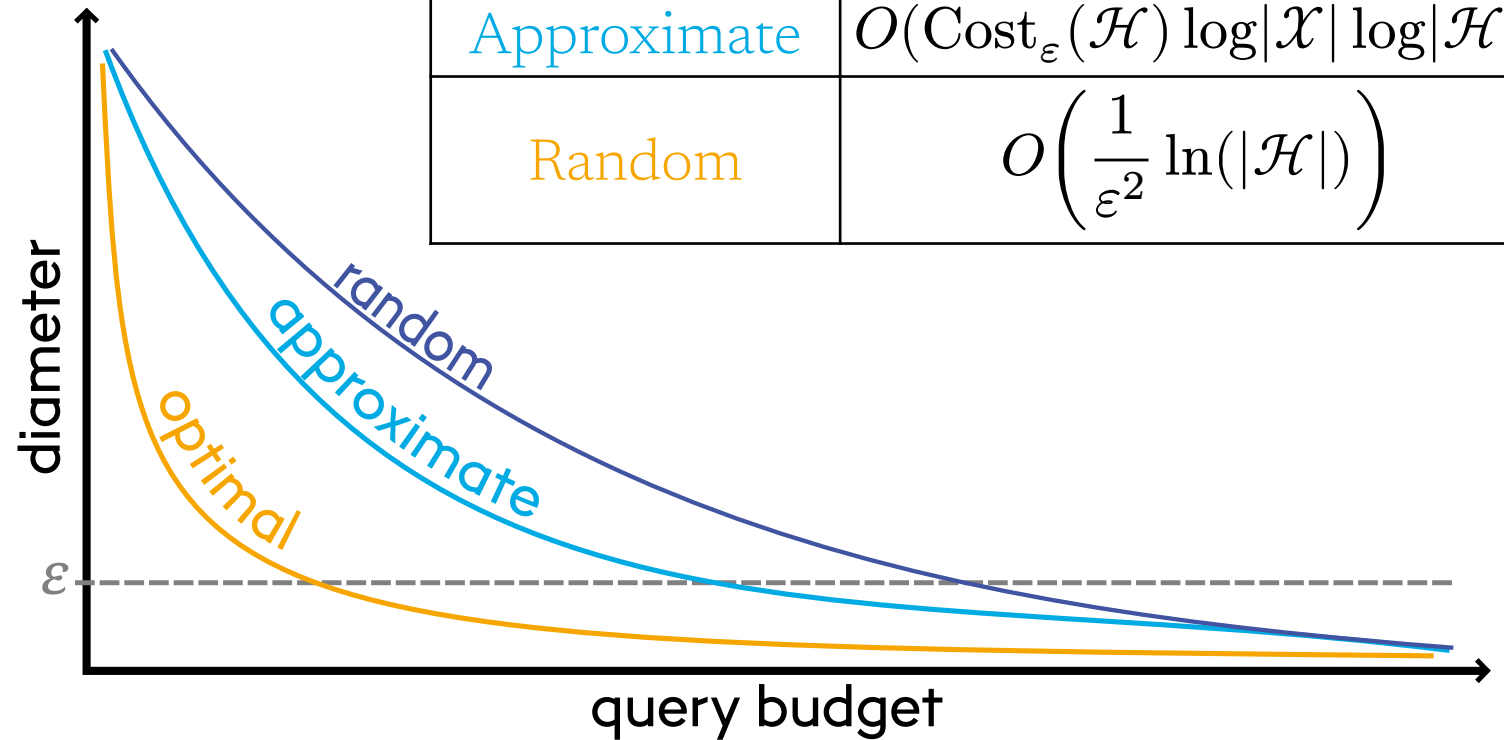
 Framework

 A theoretical peek

 Empirical study

Concluding remarks

Audit method	Query complexity
Optimal	$\text{Cost}_\varepsilon(\mathcal{H})$
Approximate	$O(\text{Cost}_\varepsilon(\mathcal{H}) \log \mathcal{X}  \log \mathcal{H} )$
Random	$O\left(\frac{1}{\varepsilon^2} \ln( \mathcal{H} )\right)$



**Problem:** Computational cost depends on  $\mathcal{H}$ !

[11] T. Yan and C. Zhang, "Active Fairness Auditing," in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



# State of the Art

## Active auditing algorithms

Context

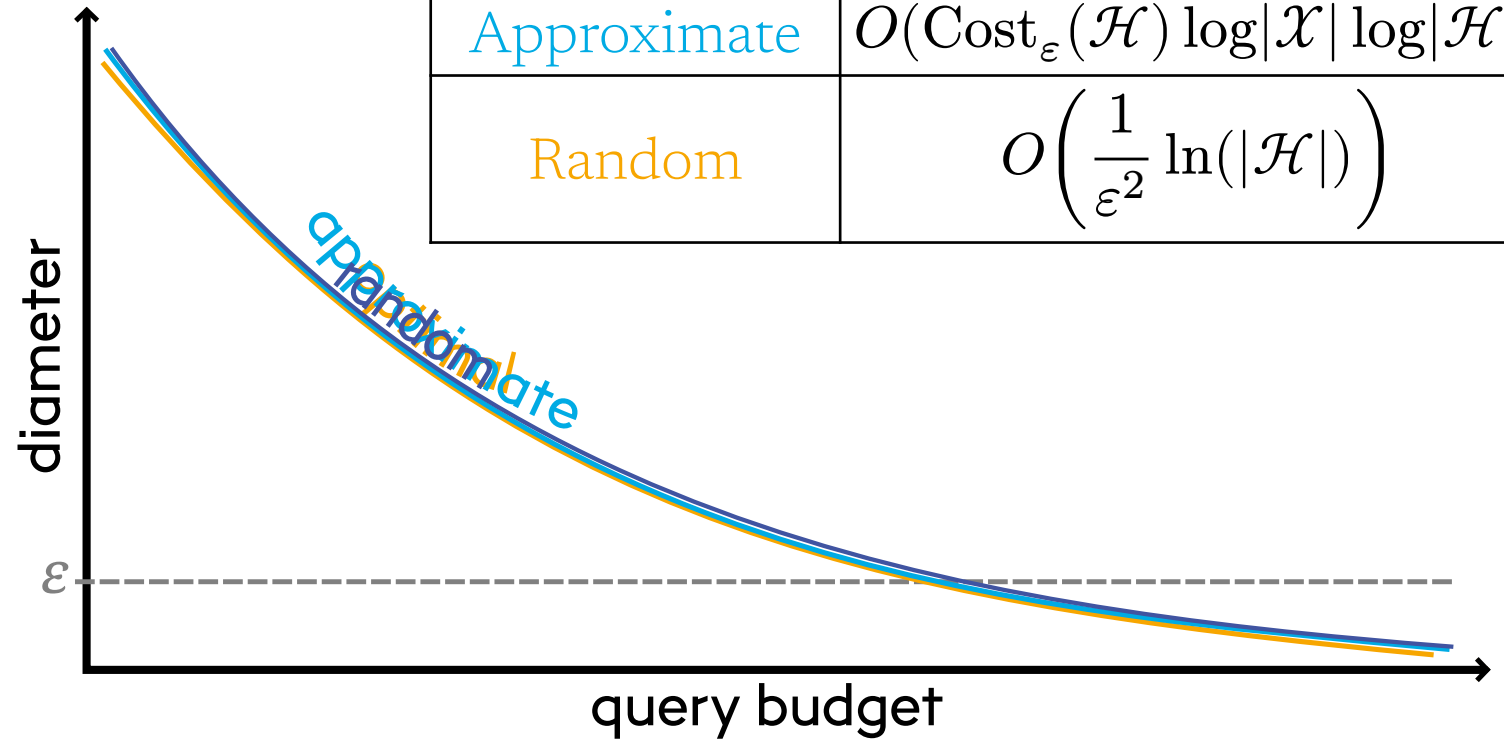
 Framework

 A theoretical peek

 Empirical study

Concluding remarks

Audit method	Query complexity
Optimal	$\text{Cost}_\varepsilon(\mathcal{H})$
Approximate	$O(\text{Cost}_\varepsilon(\mathcal{H}) \log \mathcal{X}  \log \mathcal{H} )$
Random	$O\left(\frac{1}{\varepsilon^2} \ln( \mathcal{H} )\right)$



**Problem:** Computational cost depends on  $\mathcal{H}$ !

[11] T. Yan and C. Zhang, "Active Fairness Auditing," in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.



## Research questions

**RQ1**  $\exists \mathcal{H}$  such that  $\left\{ \begin{array}{l} \text{Complexity}(\mathcal{H}, \text{random audit}) \\ = \\ \text{Complexity}(\mathcal{H}, \text{optimal audit}) \end{array} \right. ?$

**RQ2** Do these  $\mathcal{H}$  exist in practice ?

# A simple case

## *Shattering hypothesis class*

Context

 Framework

 **A theoretical peek**

 Empirical study

Concluding remarks

### Theorem 1: No need to aim

If  $\mathcal{H} = \{0, 1\}^x$ , then

$$\begin{aligned} \text{diam}_{\mu} \mathcal{H}(h^*, S) &= 2 - (\mathbb{P}(X \in S \mid X_A = 1) \\ &\quad + \mathbb{P}(X \in S \mid X_A = 0)) \end{aligned}$$



# A simple case

## *Shattering hypothesis class*

Context

 Framework

 **A theoretical peek**

 Empirical study

Concluding remarks

### Theorem 1: No need to aim

If  $\mathcal{H} = \{0, 1\}^x$ , then

$$\text{diam}_{\mu} \mathcal{H}(h^*, S) = 2 - (\mathbb{P}(X \in S \mid X_A = 1) + \mathbb{P}(X \in S \mid X_A = 0))$$

### **Proof intuition:**

1. Split the value of the  $\mu$ -diameter on  $S$  and  $\bar{S}$
2. Construct the “optimal” hypotheses  $h^{\uparrow}$  and  $h^{\downarrow}$
3. Express the result as a function of  $\mathbb{P}(X \in S \mid X_A = 0 \text{ or } 1)$



# A more refined case

## *Dictionary models*

Context

 Framework

 **A theoretical peek**

 Empirical study

Concluding remarks

### Theorem 2: Little Robert (informal)

Let  $d \in \{0, 1\}^x$  be a dictionary of memory  $m$ .



# A more refined case

## *Dictionary models*

Context

 Framework

 **A theoretical peek**

 Empirical study

Concluding remarks

## Theorem 2: Little Robert (informal)

Let  $d \in \{0, 1\}^x$  be a dictionary of memory  $m$ .

Then, for  $m$  large enough,  $\text{diam}_\mu \mathcal{H}(S, Y)$  is a function of

- ▶ sensitive groups  $\mathbb{P}(X \in S \mid X_A = 1, 2)$  ( $\searrow$ )
- ▶ the dictionary's size  $m$  ( $\nearrow$ )



# A more refined case

## Dictionary models

Context

Framework

**A theoretical peek**

Empirical study

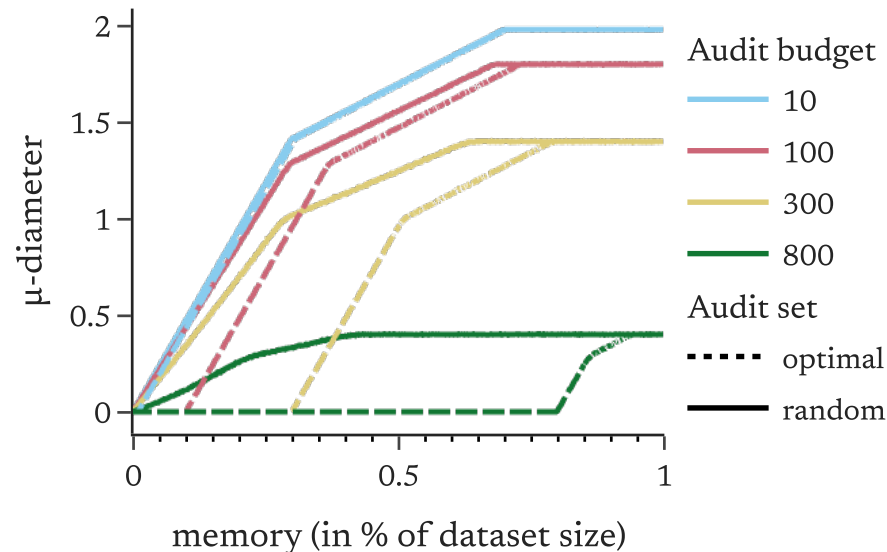
Concluding remarks

## Theorem 2: Little Robert (informal)

Let  $d \in \{0, 1\}^x$  be a dictionary of memory  $m$ .

Then, for  $m$  large enough,  $\text{diam}_\mu \mathcal{H}(S, Y)$  is a function of

- ▶ sensitive groups  $\mathbb{P}(X \in S \mid X_A = 1, 2)$  ( $\searrow$ )
- ▶ the dictionary's size  $m$  ( $\nearrow$ )



# Benign overfitting

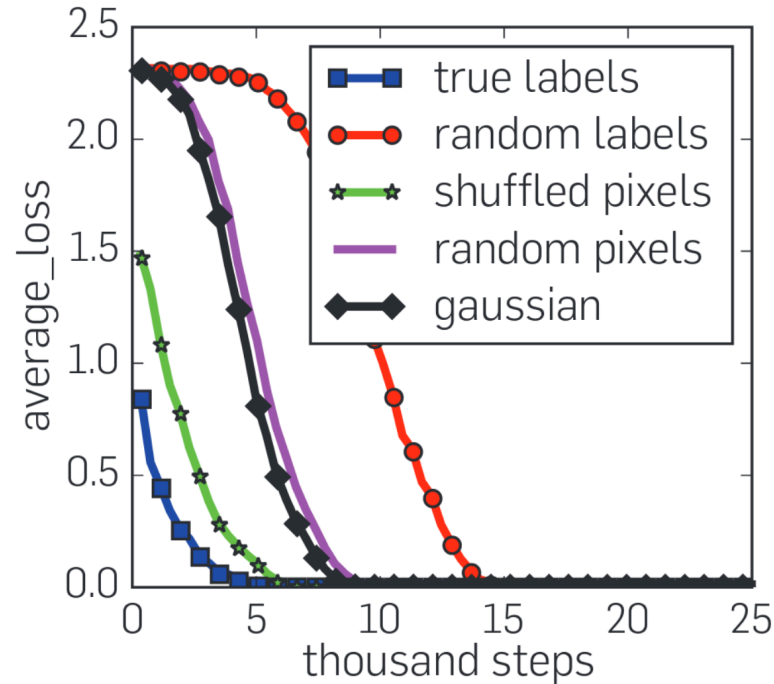
Context

Framework

A theoretical peek

Empirical study

Concluding remarks



(a) Learning curves

## Train loss of an Inception net on CIFAR10

**Taken** from [16] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding Deep Learning (Still) Requires Rethinking Generalization,” *Communications of the ACM*, vol. 64, no. 3, pp. 107–115, Feb. 2021, doi: 10.1145/34446776.



# Benign overfitting

Context

Framework

A theoretical peek

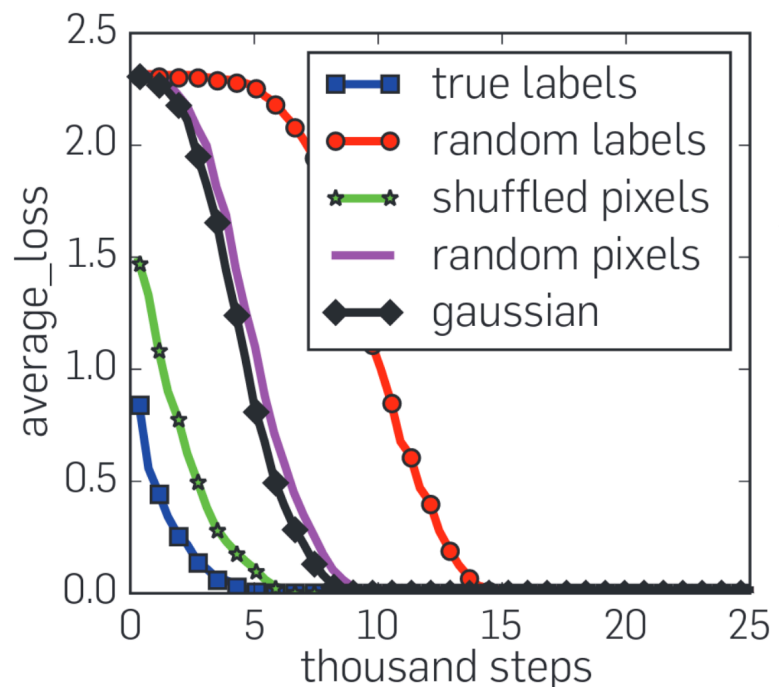
Empirical study

Concluding remarks

## Definition 2: Benign overfitting on $c$

$\mathcal{H}$  exhibits benign overfitting iff  $\exists h^* \in \mathcal{H}$  such that

1.  $\forall D \subset \mathcal{X}, |D| \leq d_0, \text{error}(h, D) = 0$
2.  $\text{error}(h^*, \mathcal{X}) \leq \varepsilon$



## Train loss of an Inception net on CIFAR10

**Taken** from [16] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding Deep Learning (Still) Requires Rethinking Generalization,” *Communications of the ACM*, vol. 64, no. 3, pp. 107–115, Feb. 2021, doi: 10.1145/3446776.



# Benign overfitting

*and audit difficulty*

Context

 Framework

 **A theoretical peek**

 Empirical study

Concluding remarks

## Corollary 1: Large models are difficult to audit

If  $\mathcal{H}$  exhibits benign overfitting with respect to the sensitive attribute, then (for  $\varepsilon < \frac{1}{2}$ ),  $\text{diam}_\mu(h, S)$  is a function of

- ▶ sensitive groups  $\mathbb{P}(X \in S \mid X_A = 1, 2)$  ( $\searrow$ )
- ▶ the model's error-rate  $\varepsilon$  ( $\searrow$ )



## Research questions

**RQ1**  $\exists \mathcal{H}$  such that  $\begin{cases} \text{Complexity}(\mathcal{H}, \text{random audit}) \\ = \\ \text{Complexity}(\mathcal{H}, \text{optimal audit}) \end{cases} ?$

$\Rightarrow$  **Yes:** *models with high capacity*

**RQ2** Do these  $\mathcal{H}$  exist in practice ?

# Metrics

Context

 Framework

 A theoretical peek

 **Empirical study**

Concluding remarks

## Simulated models

**Model family  $\mathcal{F}$**

e.g. decision trees

● perceptron   ◆ linear   ■ tree   ✖ gbdt   ★  $\mathcal{H}_{\text{opt}}$

**Hypothesis class  $\mathcal{H} \in \mathcal{F}$**

e.g. trees with `max_depth = 2`

## Metrics

- ▶  $\text{AuditManipulability}(\mathcal{H})$
- ▶  $\text{ModelCapacity}(\mathcal{H})$



# Metrics

Context

 Framework

 A theoretical peek

 **Empirical study**

Concluding remarks

## Simulated models

**Model family  $\mathcal{F}$**

e.g. decision trees

● perceptron   ◆ linear   ■ tree   ✖ gbdt   ★  $\mathcal{H}_{\text{opt}}$

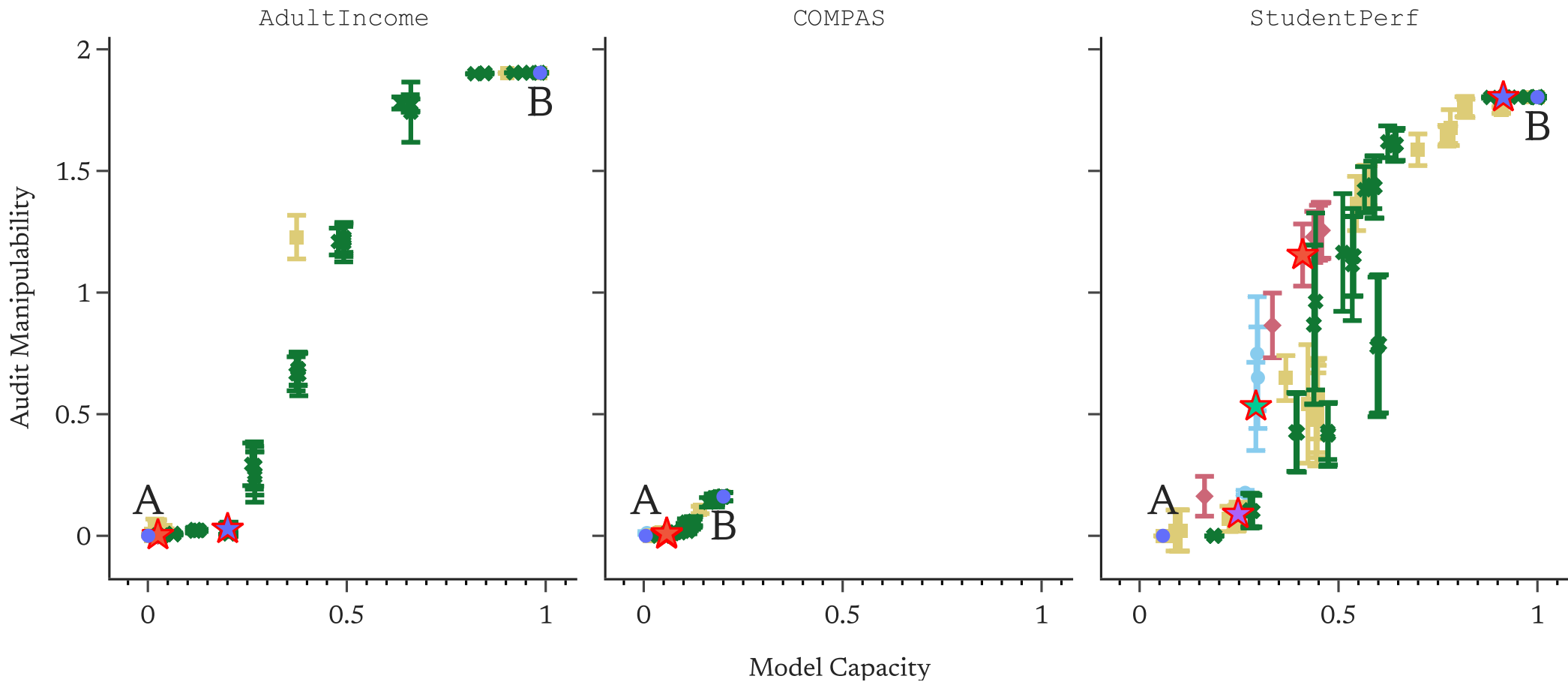
**Hypothesis class  $\mathcal{H} \in \mathcal{F}$**

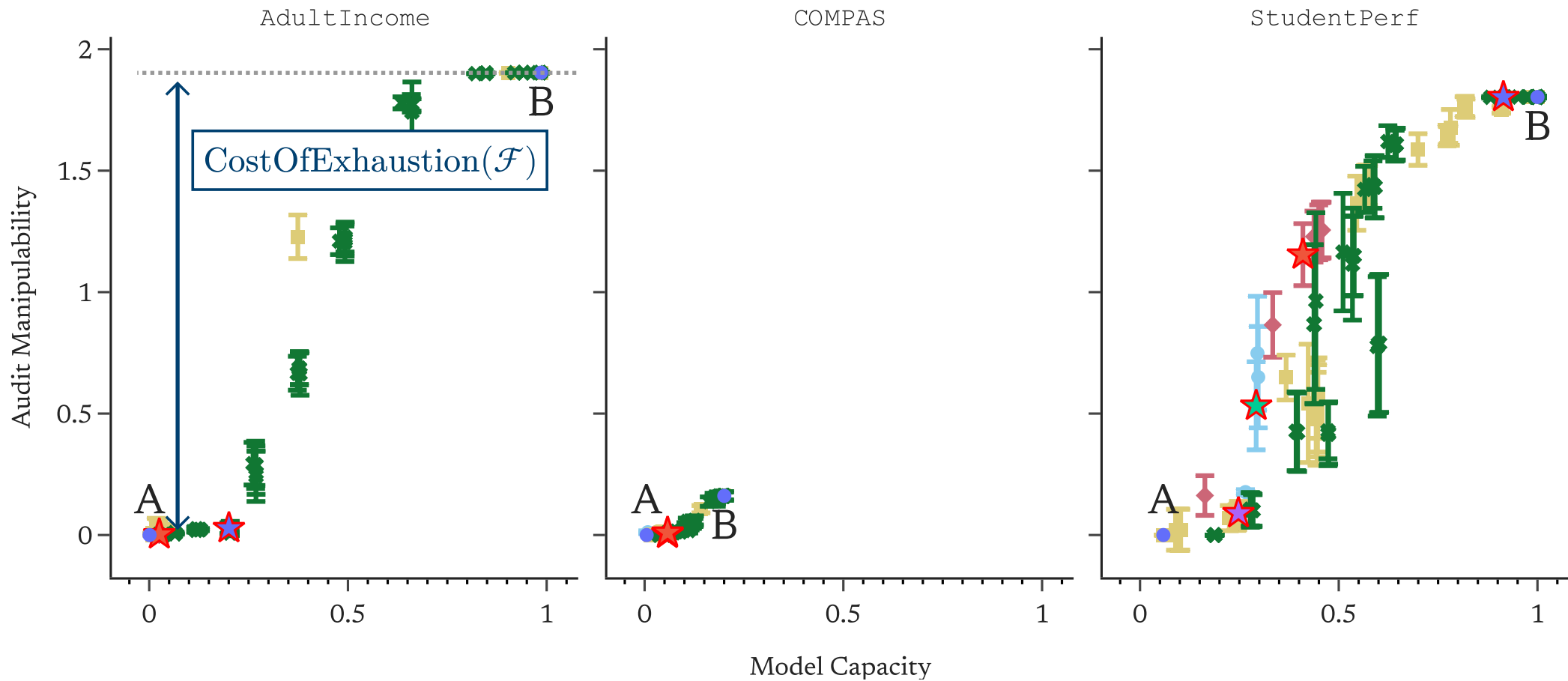
e.g. trees with `max_depth = 2`

## Metrics

- ▶  $\text{AuditManipulability}(\mathcal{H}) = \mathbb{E}_{S \sim \mathcal{X}^m} [\text{diam}_\mu(h^*, S)]$
- ▶  $\text{ModelCapacity}(\mathcal{H}) = \mathbb{E}_{D \sim \mathcal{X}^r} [\text{Rademacher}(\mathcal{H}, D)]$







# Cost of exhaustion

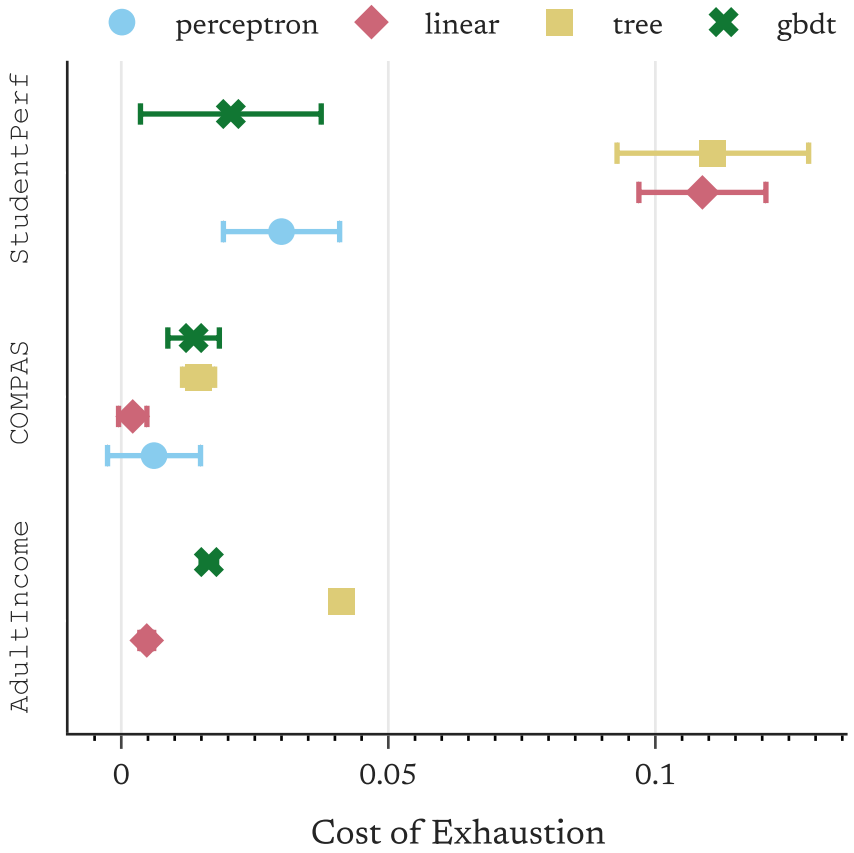
Context

Framework

A theoretical peek

**Empirical study**

Concluding remarks



# Conclusion

It seems [...] a platform could always *game the system* [...] *without sacrificing a lot of accuracy* of the model learnt.

– Anonymous reviewer

## Paper link



## Conclusions

- ▶ Robustness needs a prior
- ▶ Prior on the model  
⇒ guarantees depend on the capacity of the model

## Implications for AI regulation

- ▶ Need more access to the model
- ▶ And/Or anonymous auditor

Context

 Framework

 A theoretical peek

 Empirical study

**Concluding remarks**



Thanks for your attention !



Announcing a new team in SaTML town:

**ARTISHAU**

**ART**ificial Intelligence: **S**ecurity,  
**T**rut**H**fulness and **AU**dit

Centre Inria de l'Université de Rennes,  
France

# Bibliography

- [1] J. Dastin, “Amazon Scraps Secret AI Recruiting Tool That Showed Bias against Women,” *Reuters*, Oct. 2018, Accessed: Mar. 06, 2023. [Online]. Available: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>
- [2] L. Chen, A. Mislove, and C. Wilson, “An Empirical Analysis of Algorithmic Pricing on Amazon Marketplace,” in *Proceedings of the 25th International Conference on World Wide Web*, in WWW '16. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee, Apr. 2016, pp. 1339–1349. doi: 10.1145/2872427.2883089.
- [3] “EU AI Act: First Regulation on Artificial Intelligence | News | European Parliament.” Accessed: Jun. 21, 2023. [Online]. Available: <https://www.europarl.europa.eu/news/en/headlines/society/20230601STO93804/eu-ai-act-first-regulation-on-artificial-intelligence>
- [4] J. Larson, S. Mattu, L. Kirchner, and J. Angwin, “How We Analyzed the COMPAS Recidivism Algorithm,” *ProPublica*, May 2016, Accessed: Mar. 06, 2023. [Online]. Available: <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
- [5] Rédaction, “Numérique : que sont le DMA et le DSA, les règlements européens qui visent à réguler internet ?” Accessed: Jun. 21, 2023. [Online]. Available: <https://www.touteleurope.eu/societe/numerique-que-sont-le-dma-et-le-dsa-les-reglements-europeens-qui-veulent-reguler-internet/>
- [6] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities*. fairmlbook.org, 2019.
- [7] Arvind Narayanan, “Tutorial: 21 Fairness Definitions and Their Politics.” Accessed: Oct. 12, 2023. [Online]. Available: <https://www.youtube.com/watch?v=jIXIuYdnyyk>
- [8] B. Rastegarpanah, K. Gummedi, and M. Crovella, “Auditing Black-Box Prediction Models for Data Minimization Compliance,” in *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 2021, pp. 20621–20632. Accessed: Nov. 02, 2022. [Online]. Available: <https://proceedings.neurips.cc/paper/2021/hash/ac6b3cce8c74b2e23688c3e45532e2a7-Abstract.html>
- [9] F. Lu *et al.*, “A General Framework for Auditing Differentially Private Machine Learning,” presented at the Advances in Neural Information Processing Systems, Dec. 2022, pp. 4165–4176. Accessed: Aug. 16, 2023. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2022/hash/1add3bbdbc20c403a383482a665eb5a4-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/1add3bbdbc20c403a383482a665eb5a4-Abstract-Conference.html)
- [10] J. Bandy, “Problematic Machine Behavior: A Systematic Literature Review of Algorithm Audits,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 5, no. CSCW1, pp. 1–34, Apr. 2021, doi: 10.1145/3449148.
- [11] T. Yan and C. Zhang, “Active Fairness Auditing,” in *Proceedings of the 39th International Conference on Machine Learning*, PMLR, Jun. 2022.
- [12] B. Chugg, S. Cortes-Gomez, B. Wilder, and A. Ramdas, “Auditing Fairness by Betting.” 2023.
- [13] C. Yadav, M. Moshkovitz, and K. Chaudhuri, “XAudit : A Theoretical Look at Auditing with Explanations.” Accessed: Sep. 18, 2023. [Online]. Available: <http://arxiv.org/abs/2206.04740>
- [14] A. Shahin Shamsabadi *et al.*, “Washing The Unwashable : On The (Im)possibility of Fairwashing Detection,” in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., Curran Associates, Inc., 2022, pp. 14170–14182. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/5b84864ff8474fd742c66f219b2eaac1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/5b84864ff8474fd742c66f219b2eaac1-Paper-Conference.pdf)
- [15] J. G. Bourrée, E. L. Merrer, G. Tredan, and B. Rottembourg, “On the relevance of APIs facing fairwashed audits,” *arXiv preprint arXiv:2305.13883*, 2023.
- [16] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding Deep Learning (Still) Requires Rethinking Generalization,” *Communications of the ACM*, vol. 64, no. 3, pp. 107–115, Feb. 2021, doi: 10.1145/3446776.