

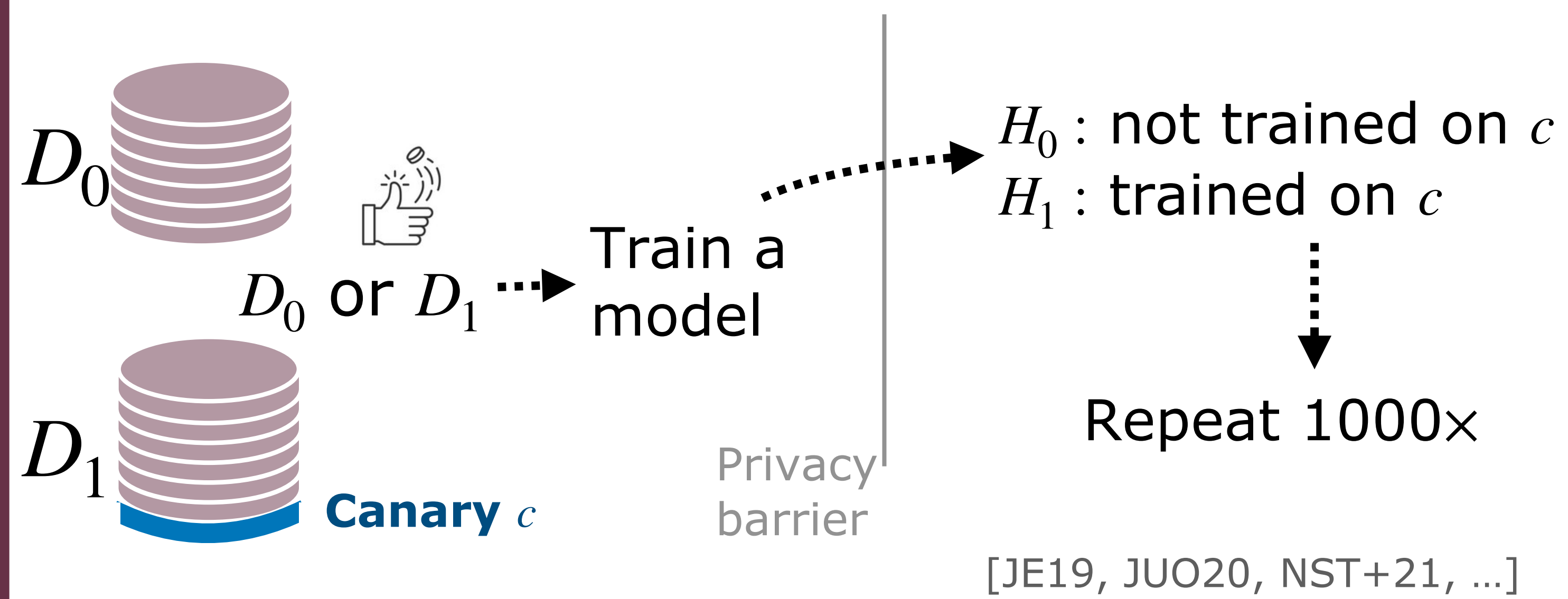
Unleashing the Power of Randomization in Auditing Differentially Private ML

Krishna Pillutla, Galen Andrew, Peter Kairouz, H. Brendan McMahan, Alina Oprea, Sewoong Oh



Auditing DP guarantees in ML

The standard approach: binary hypothesis testing



The key components of this approach

Step 1: DP definition

For all neighboring D_0, D_1 and R , we have

$$\mathbb{P}(\mathcal{A}(D_1) \in R) \leq e^\epsilon \mathbb{P}(\mathcal{A}(D_0) \in R) + \delta \quad (1)$$

True positive rate

False positive rate

Step 2: Set up the hypothesis test

Take $D_0 = \text{dataset}$, $D_1 = D_0 \cup \{\text{canary}\}$ and the test statistic as $R = \{\theta : \text{Loss}(\text{canary}; \theta) \leq \tau\}$

Step 3: Bernoulli Confidence intervals

Run n trials (each trial = one model training run)

$$\text{TPR} \approx \frac{1}{n} \sum_{i=1}^n \text{Loss}_i(D_1) \leq \tau \pm \sqrt{\frac{\text{variance}}{n}} \quad (2)$$

True rate

Empirical rate

Overall, (1) + (2) \Rightarrow

$$\epsilon \geq \log \left(\frac{\text{TPR} - \delta}{\text{FPR}} \right) \geq \log \left(\frac{\widehat{\text{TPR}}_n - \frac{1}{\sqrt{n}} - \delta}{\widehat{\text{FPR}}_n + \frac{1}{\sqrt{n}}} \right)$$

Problem: the $1/\sqrt{n}$ term requires n large

How do we solve this? Add multiple canaries

Key: Avoid group privacy with *randomization*

Auditing Lifted DP

Step 1: Lifted DP (LiDP) definition

Def: \mathcal{A} is (ϵ, δ) -LiDP if for all random $(D_0, D_1, R) \sim \mathcal{P}$ independent of \mathcal{A} s.t. D_0, D_1 are neighboring, we have

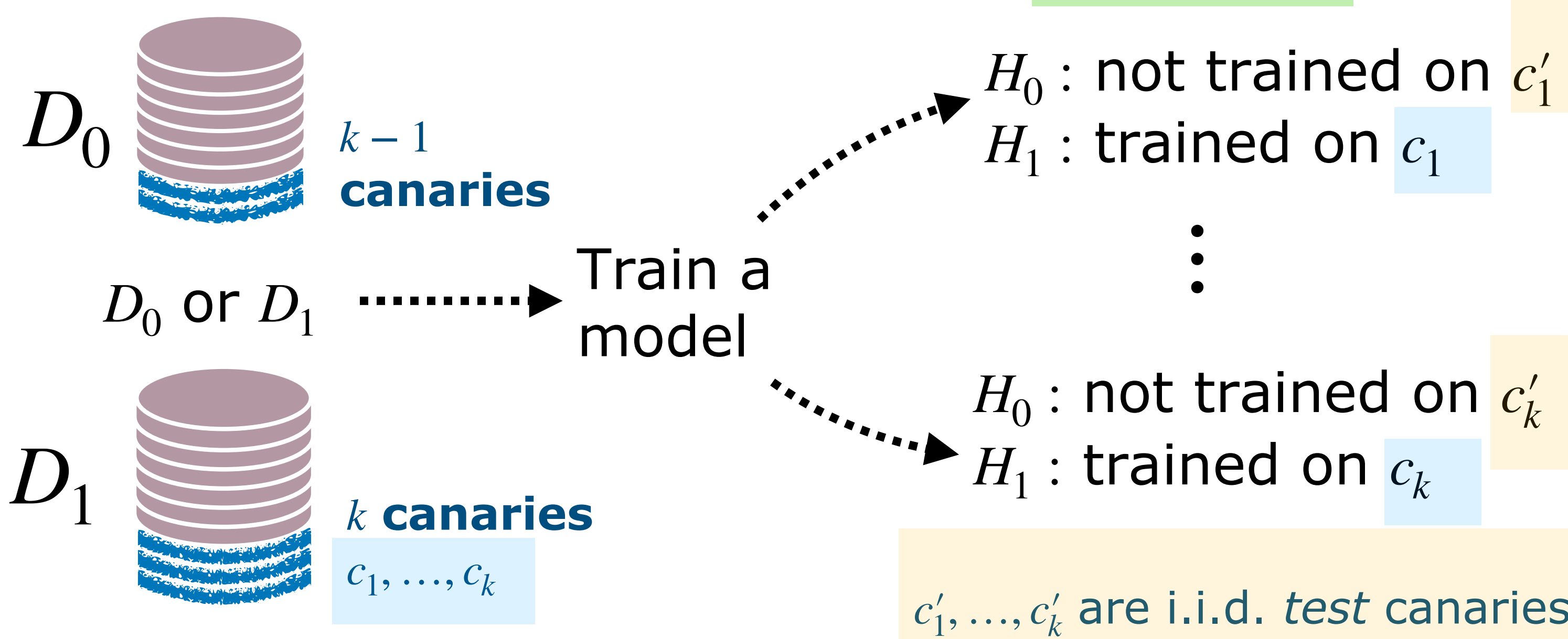
$$\mathbb{P}(\mathcal{A}(D_1) \in R) \leq e^\epsilon \mathbb{P}(\mathcal{A}(D_0) \in R) + \delta \quad (3)$$

Theorem: \mathcal{A} is (ϵ, δ) -DP $\Leftrightarrow \mathcal{A}$ is (ϵ, δ) -LiDP

Consequence: We can have *random* canaries!

Step 2: Randomized hypothesis tests

Test for k vs. $k-1$ canaries that are drawn i.i.d. from P



Consequence: Get k statistics from each trial

Step 3: Adaptive higher-order confidence

Challenge: the statistics are *correlated* (not i.i.d.)

We derive novel CIs using *empirical* correlations!

$$\left| \text{TPR} - \widehat{\text{TPR}}_{n,k} \right| \leq \sqrt{\frac{1}{n} \left(\text{corr.} + \frac{1}{k} + \sqrt{\frac{4\text{th moment}}{n}} \right)} \quad (4)$$

If $\text{corr.} = O(1/k)$, **improvement:** (3) + (4) \Rightarrow

$$\epsilon \geq \log \left(\frac{\text{TPR} - \delta}{\text{FPR}} \right) \geq \log \left(\frac{\widehat{\text{TPR}}_{n,k} - \frac{1}{\sqrt{nk}} - \frac{M_4}{n^{3/4}} - \delta}{\widehat{\text{FPR}}_{n,k} + \frac{1}{\sqrt{nk}} + \frac{M_4}{n^{3/4}}} \right)$$

Experiments

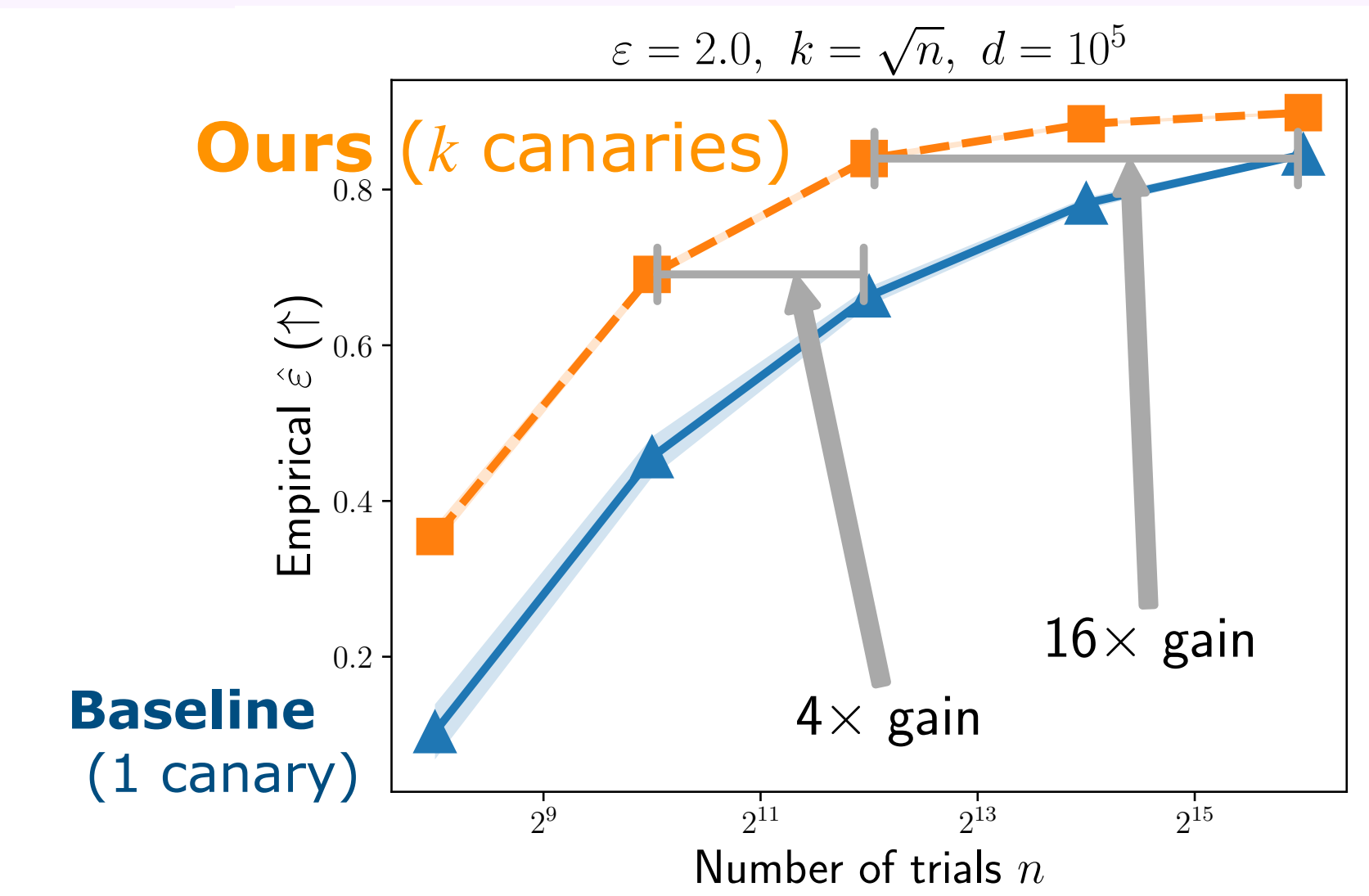
Auditing a Gaussian mechanism

Setup:

- Sum query
- Canaries: uniform over unit sphere
- Test: inner product

Result:

4-16 \times gain in sample complexity

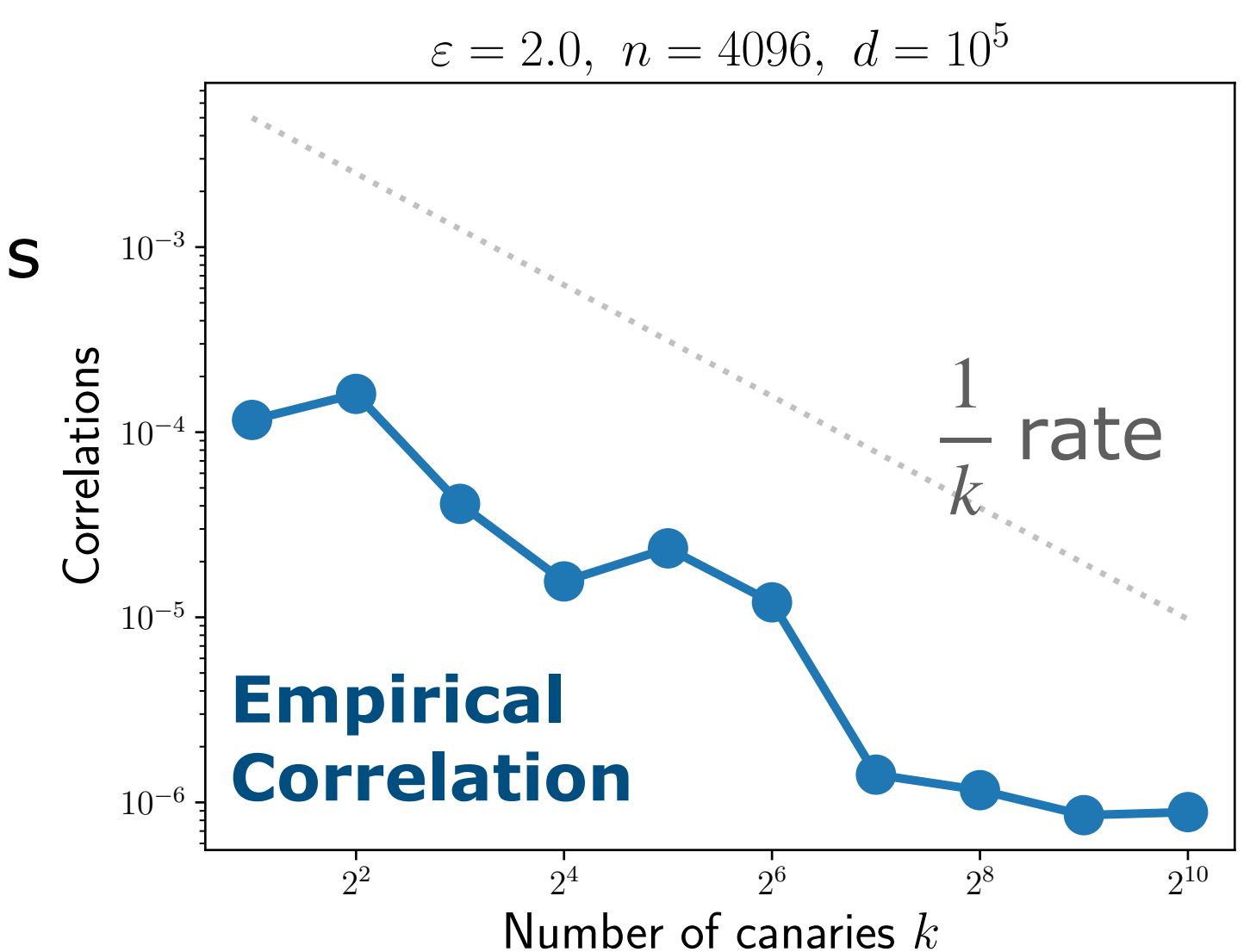


Analysis:

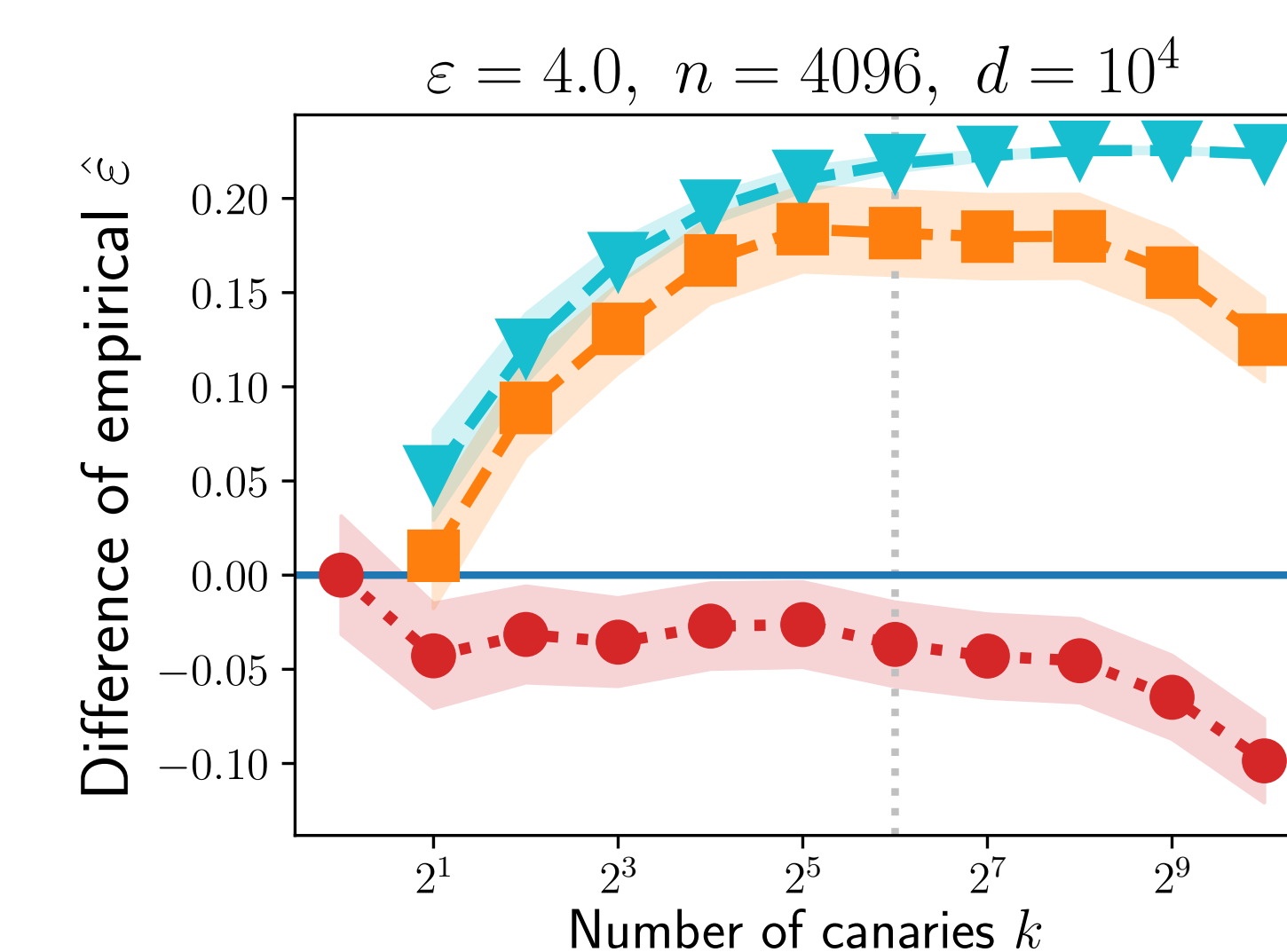
Empirical canary correlations are small, so LiDP auditing gives large wins.

Practical Guidance:

Multiple canaries should be "orthogonal"



Bias-Variance Tradeoff of LiDP:



Variance reduction of LiDP (width of the confidence interval)

Net benefit of LiDP (Balancing bias and variance)

Bias of LiDP (no higher-order estimators)

Experiments: FashionMNIST + MLP model

Gain in sample complexity from LiDP auditing

