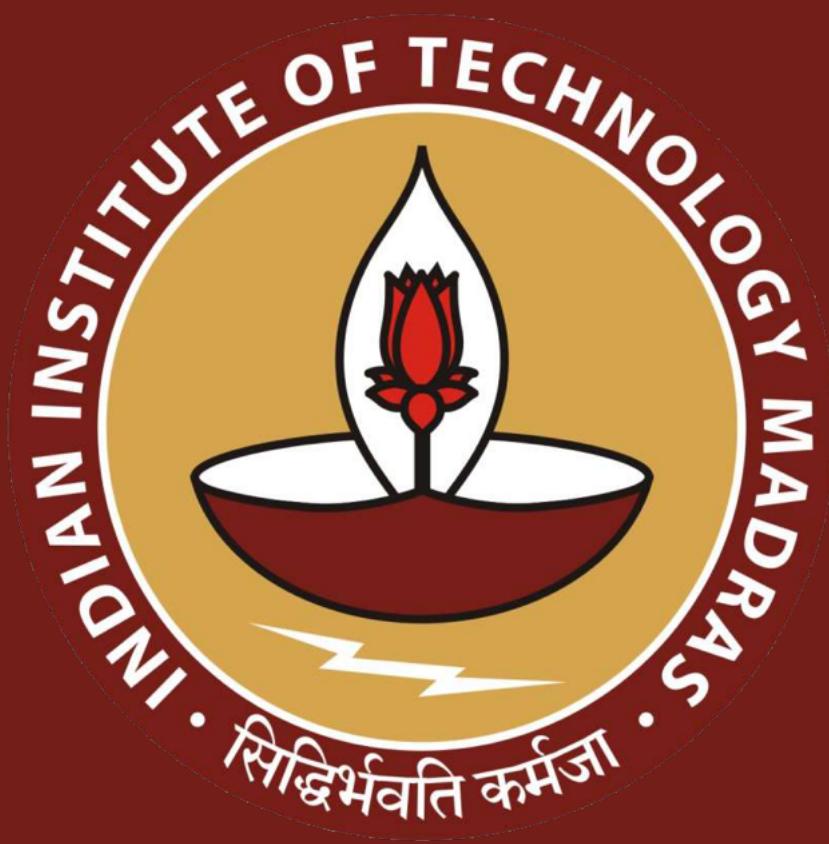


Towards user-level differential privacy at scale

Krishna Pillutla

Wadhwani School of Data Science & AI,
IIT Madras



LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



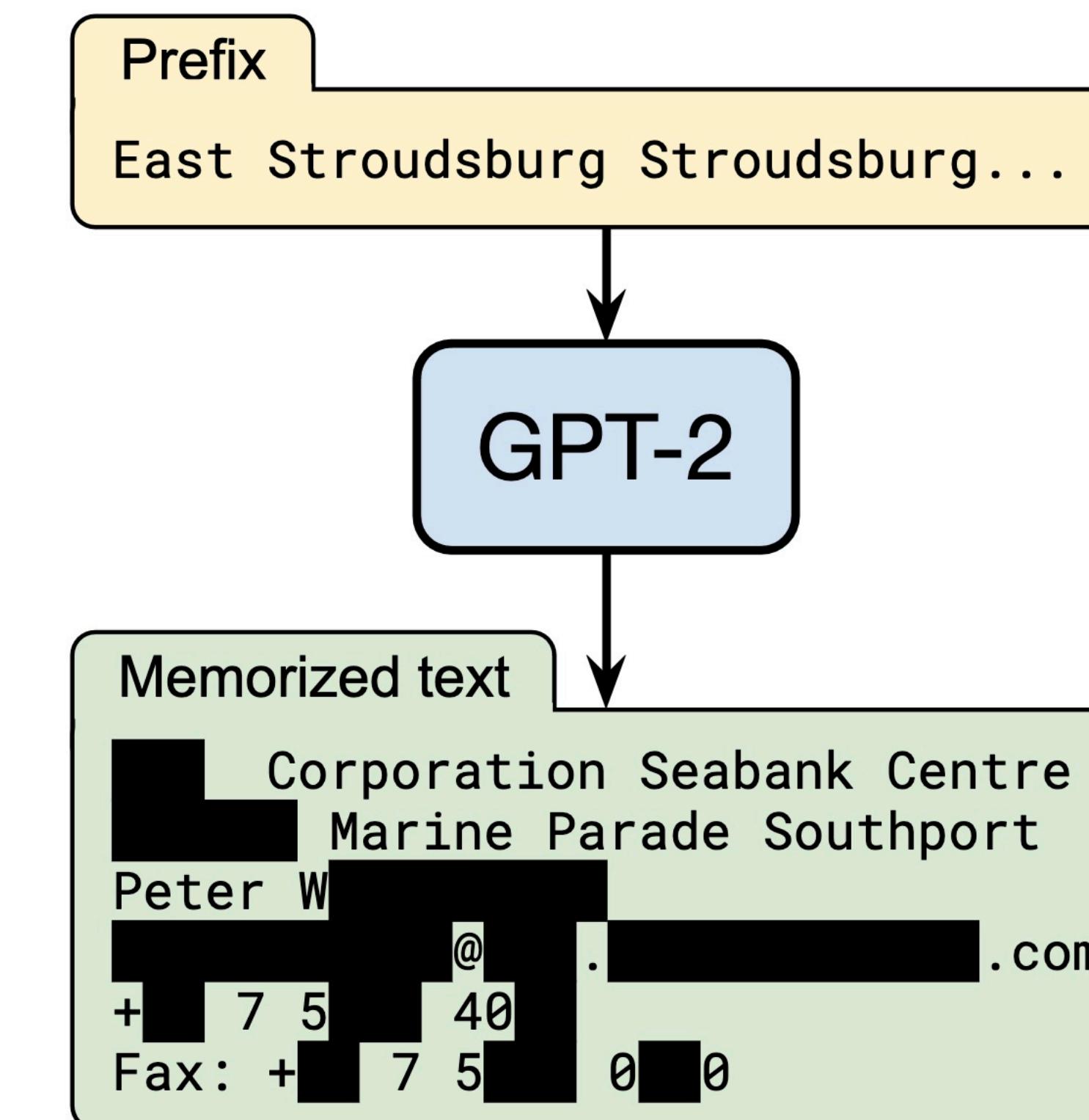
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



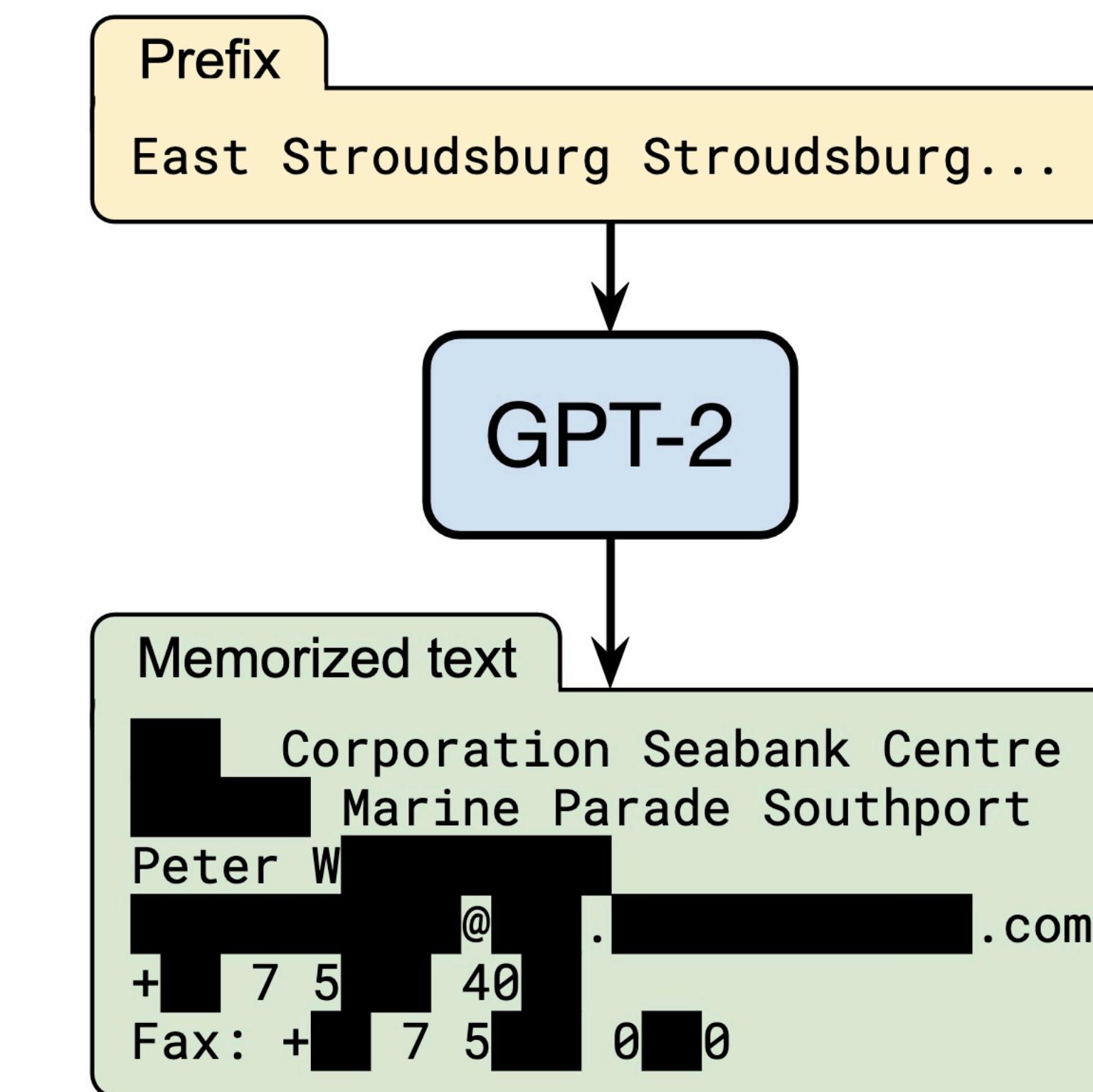
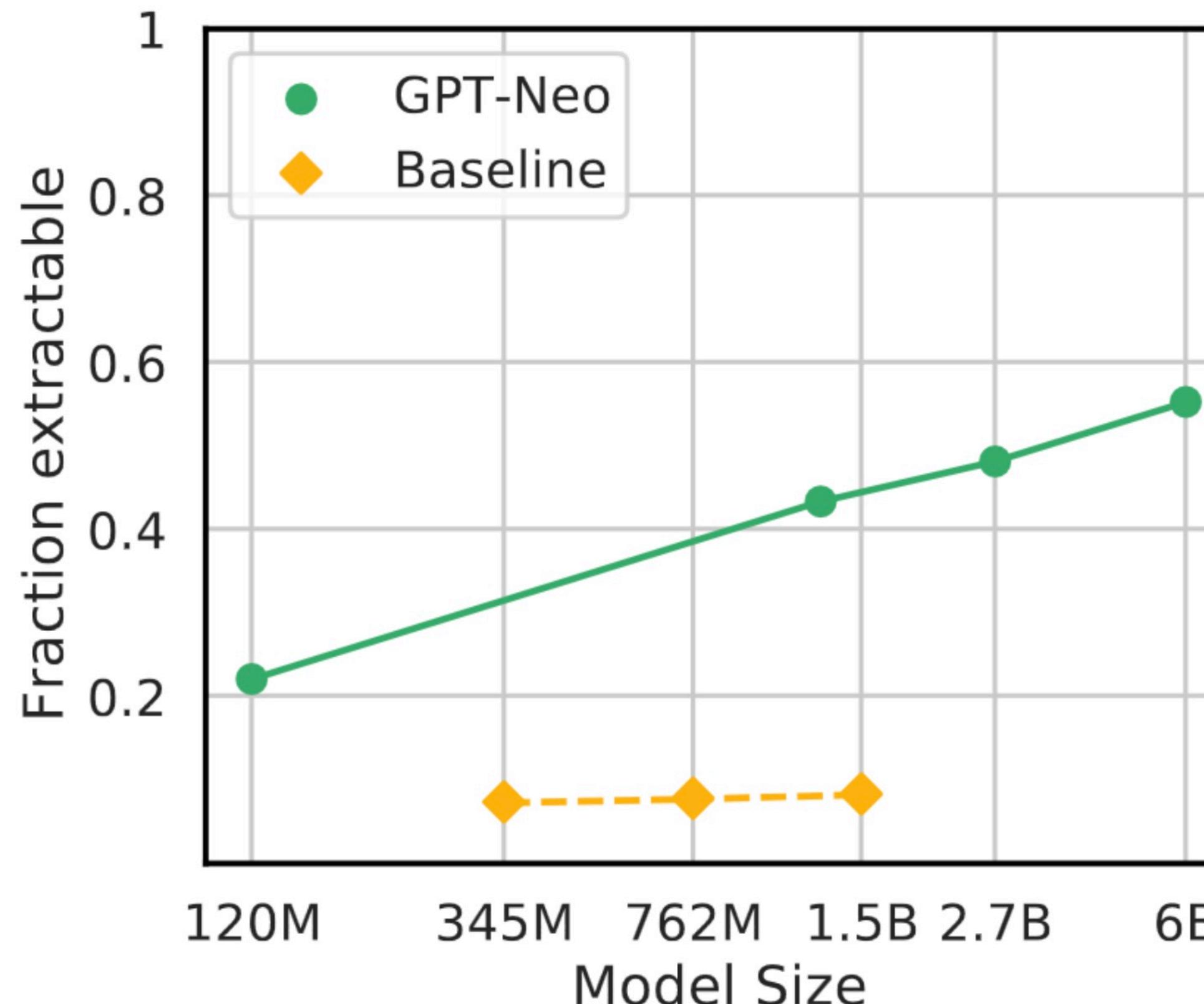
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data



Carlini et al. (USENIX Security 2021)

Models leak information about their training data *reliably*



Carlini et al. (USENIX Security 2021)

Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models

Gowthami Somepalli , Vasu Singla , Micah Goldblum , Jonas Geiping , Tom Goldstein 

 University of Maryland, College Park

{gowthami, vsingla, jgeiping, tomg}@cs.umd.edu

 New York University

goldblum@nyu.edu

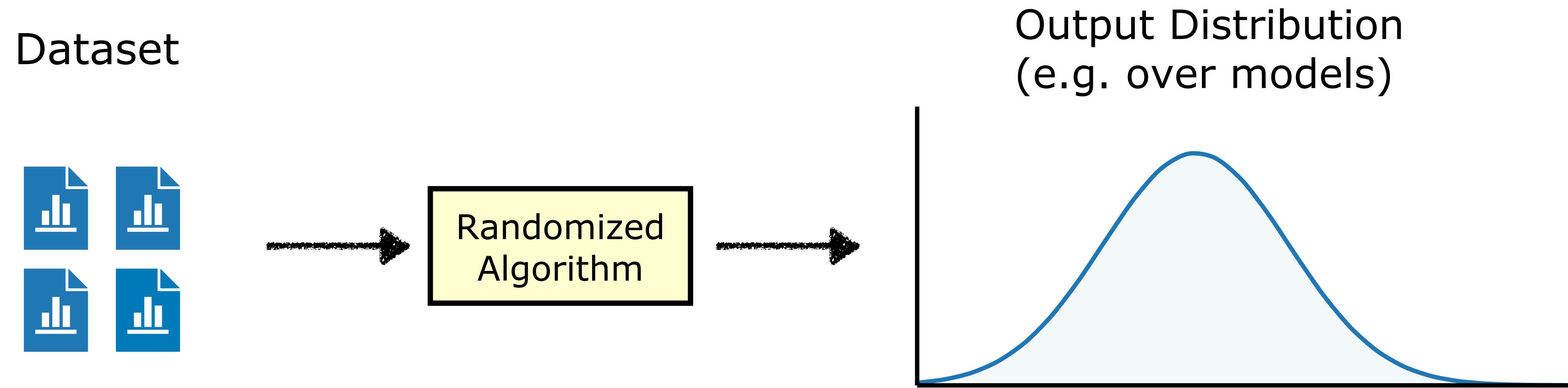
Generation



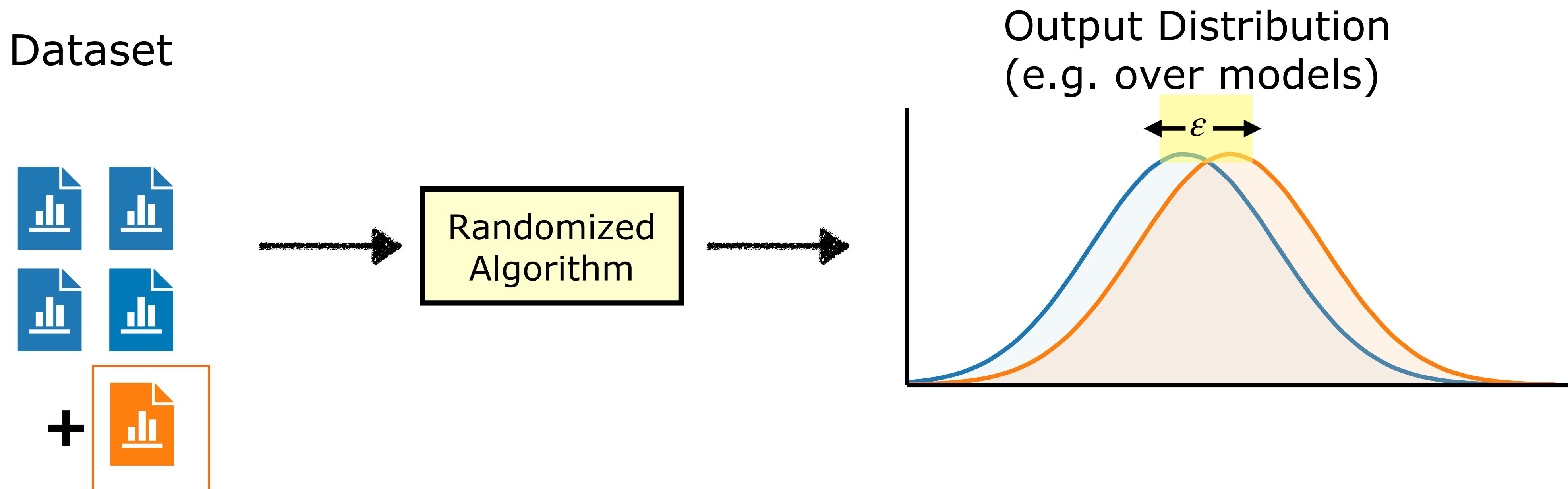
LAION-A Match



Differential privacy (DP)



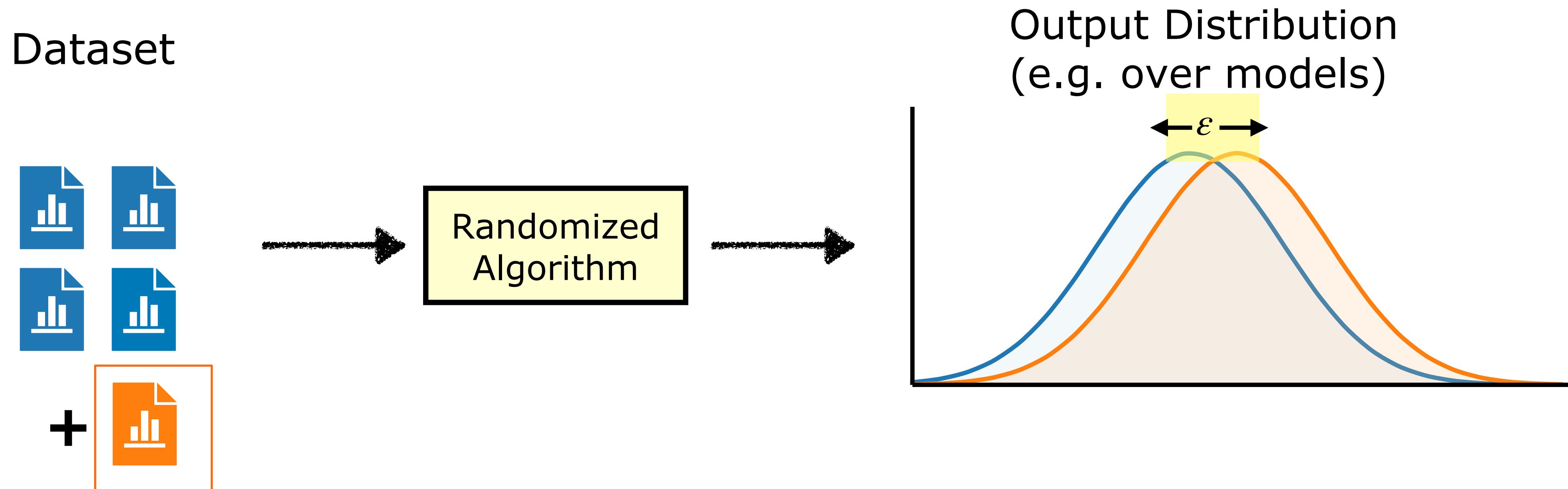
Differential privacy (DP)



A randomized algorithm is **ε -differentially private** if the addition of **one unit of data** does not alter its output distribution by more than ε

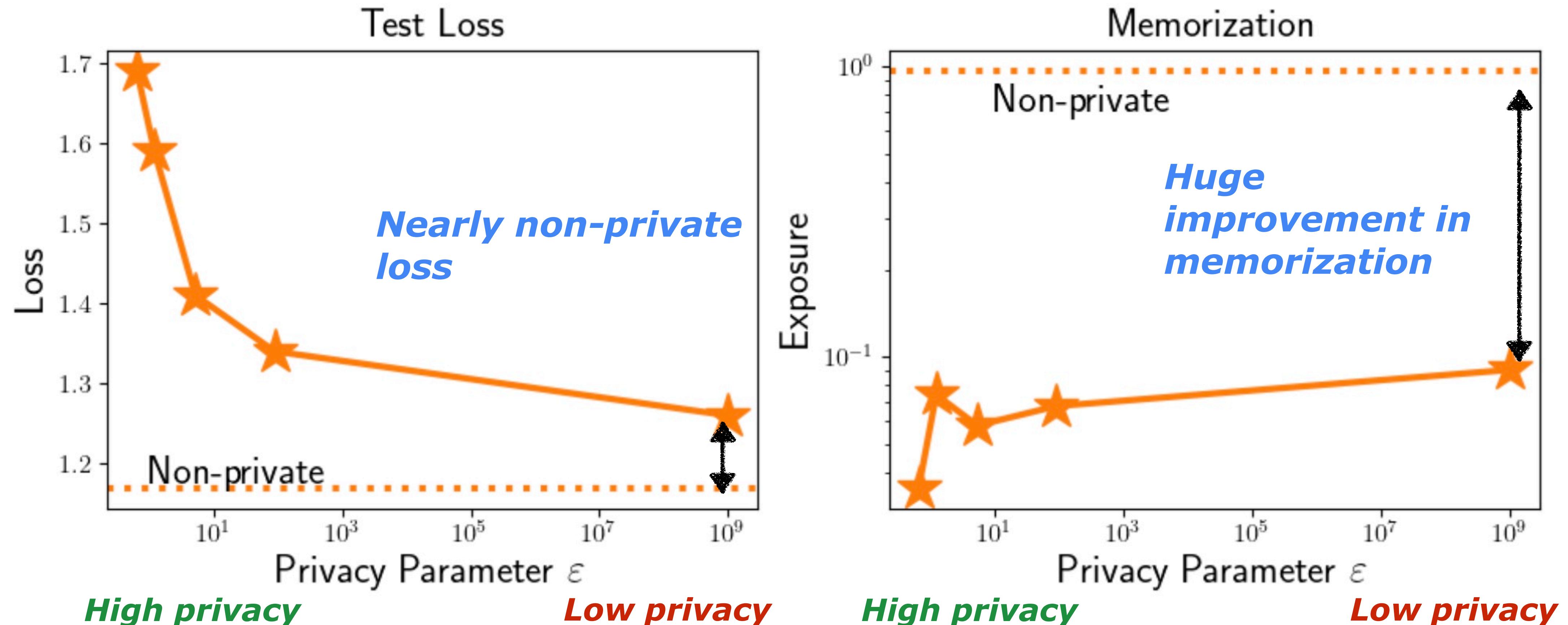
Example-level Differential privacy (DP)

 Unit of data
= **example**



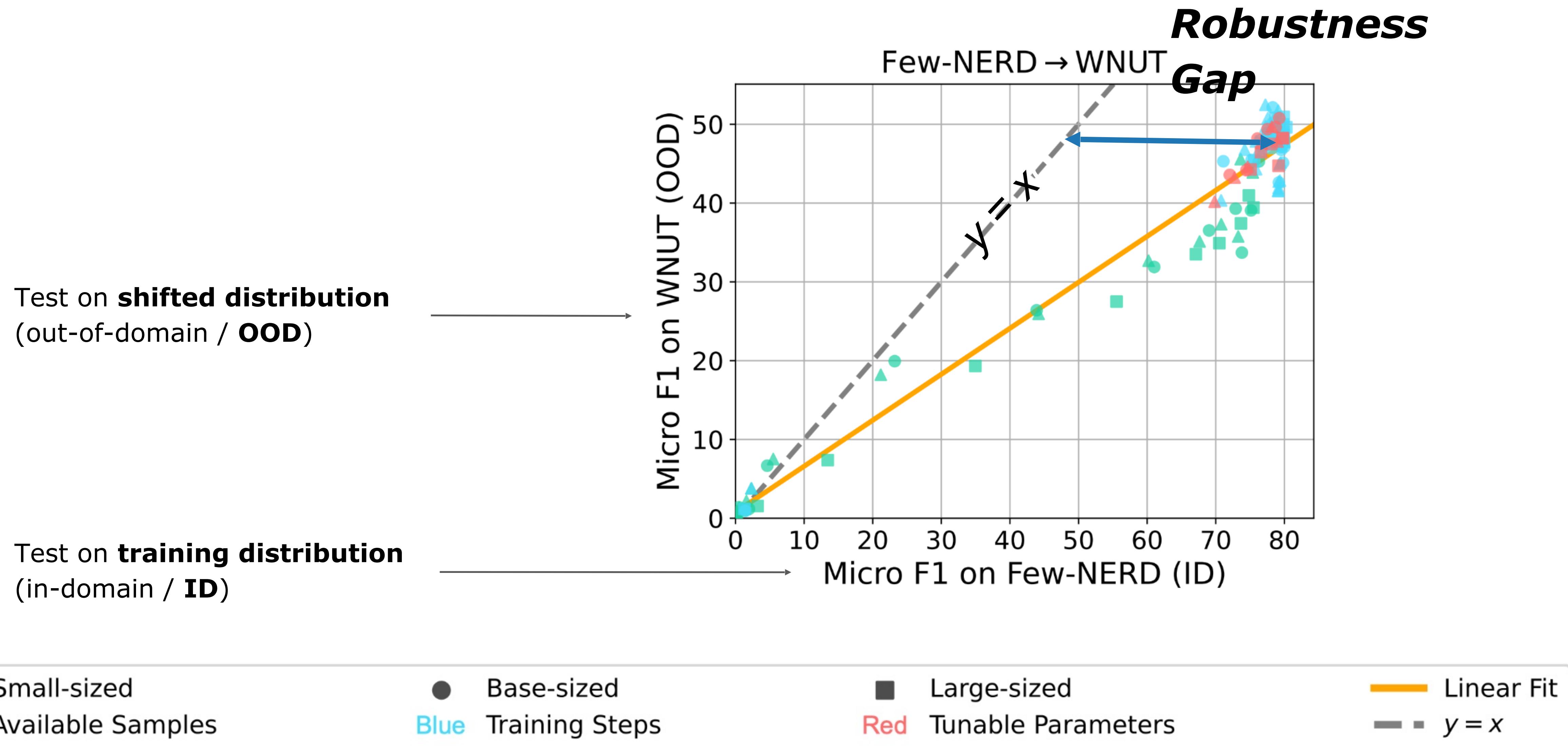
A randomized algorithm is **ε -differentially private** if the addition of **one example** does not alter its output distribution by more than ε

Differential privacy nearly eliminates memorization



Carlini, Liu, Erlingsson, Kos, Song. **The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks**. USENIX Security 2019.

Which data do we use to train/finetune/align these models?



Which data do we use to train/finetune/align these models?

Best training data = in-domain data

Test on **shifted distribution**
(out-of-domain / OOD)

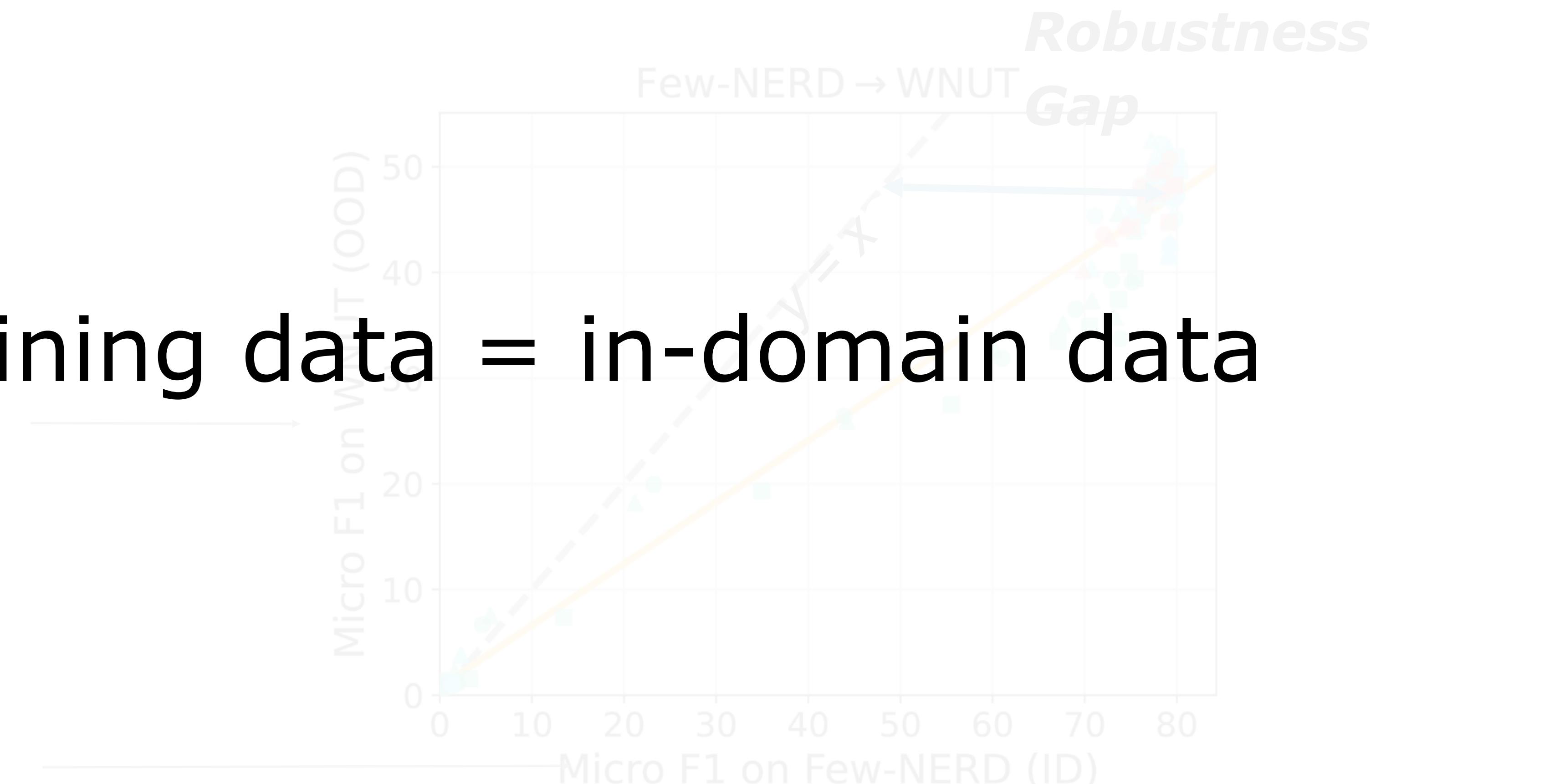
Test on **training distribution**
(in-domain / ID)

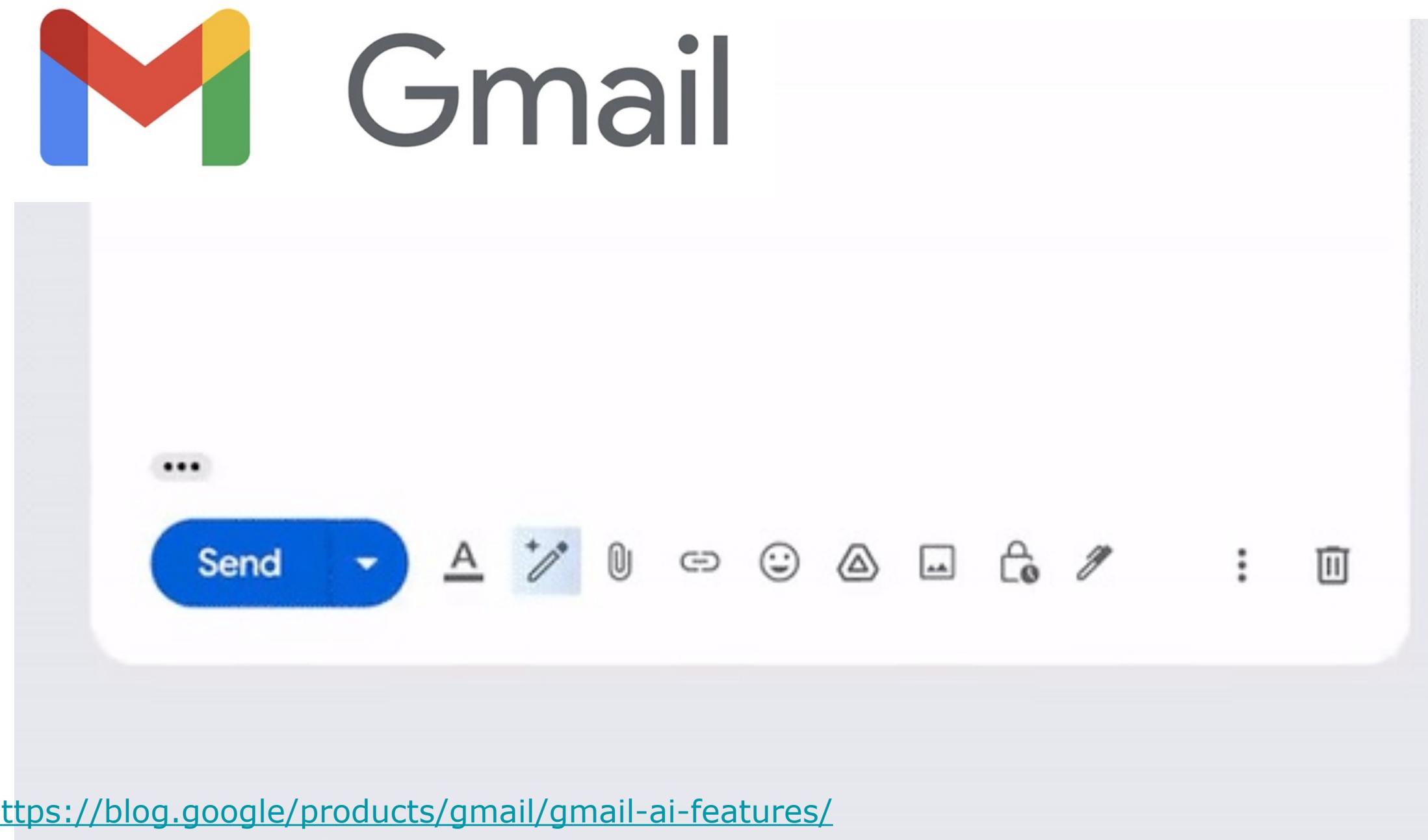
▲ Small-sized
Green Available Samples

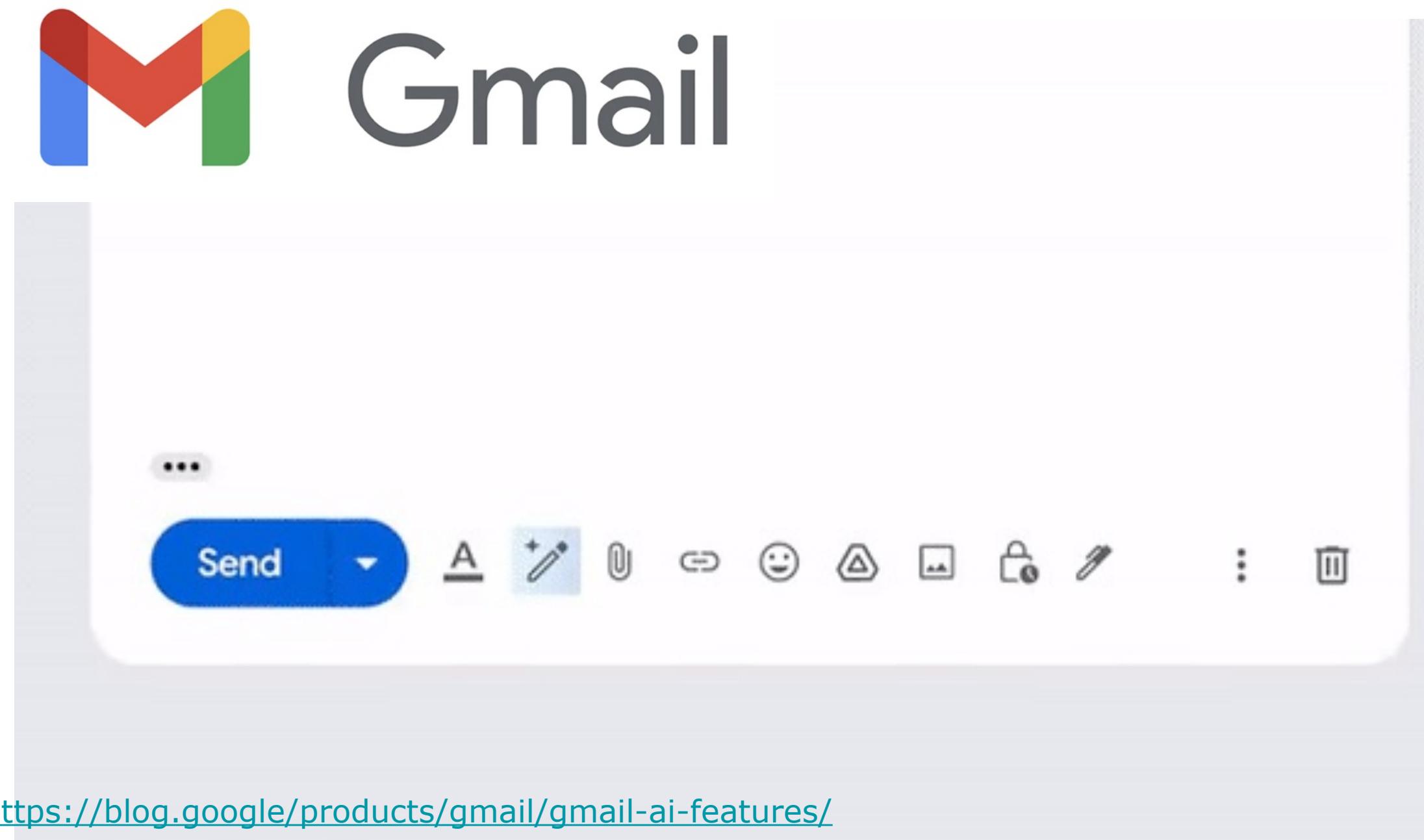
● Base-sized
Blue Training Steps

■ Large-sized
Red Tunable Parameters

— Linear Fit
— $y = x$



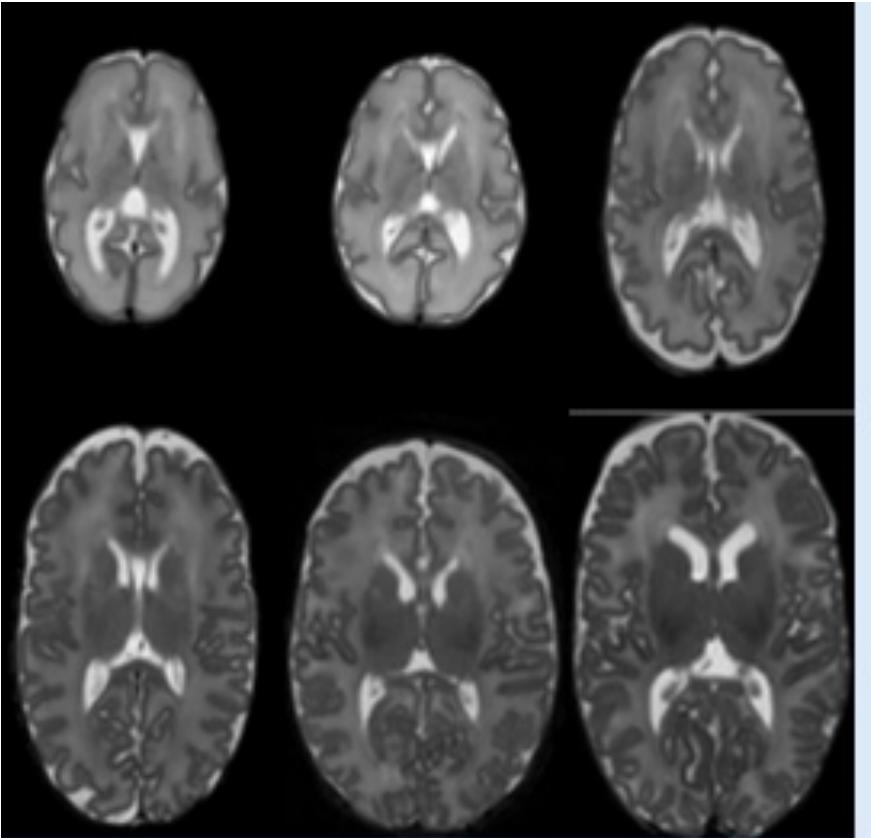
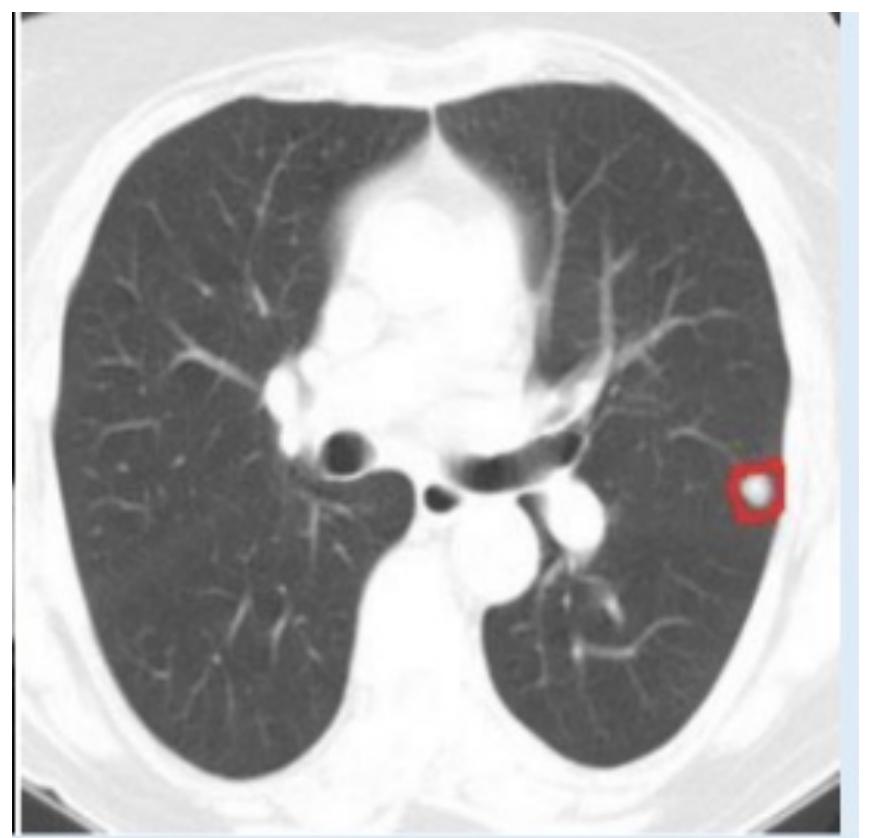
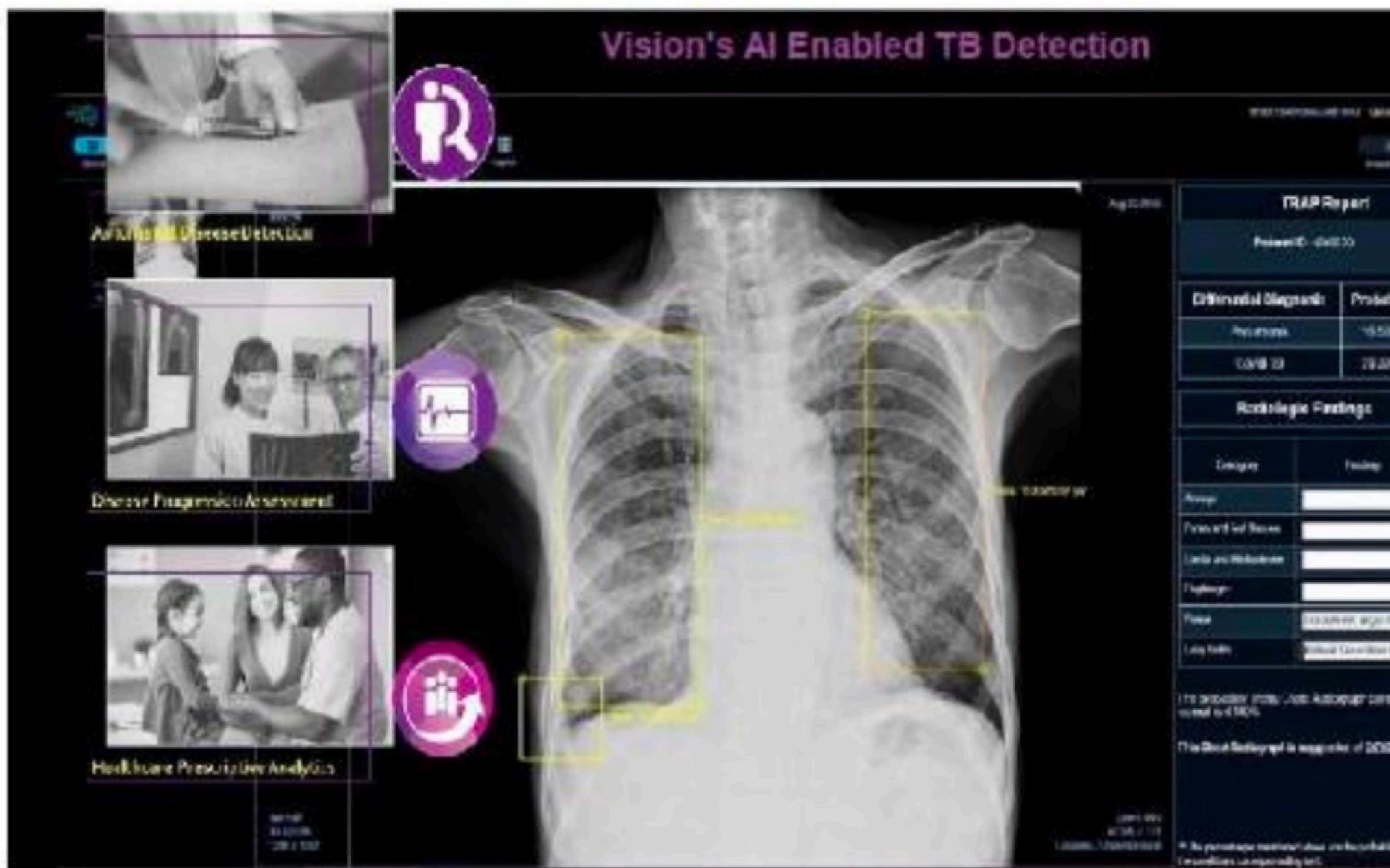


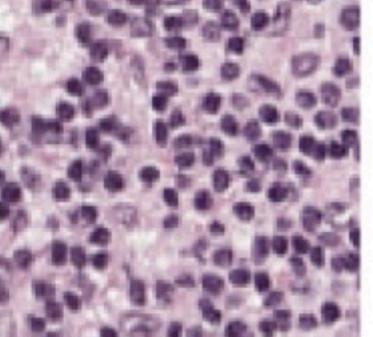
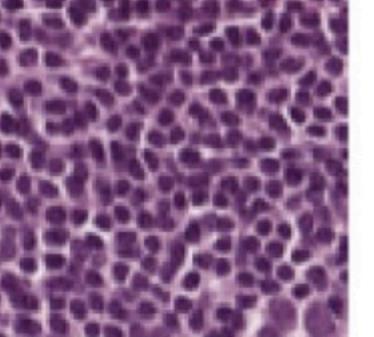
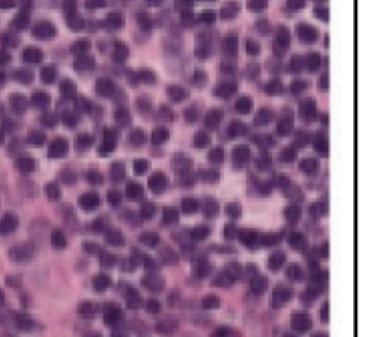
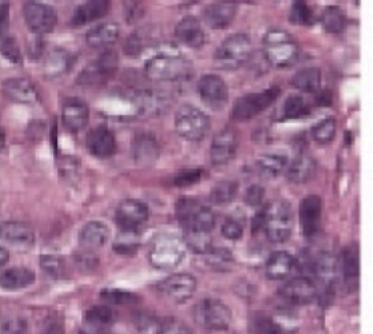
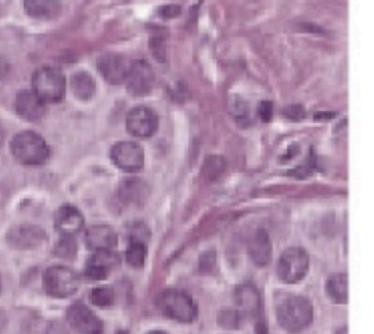
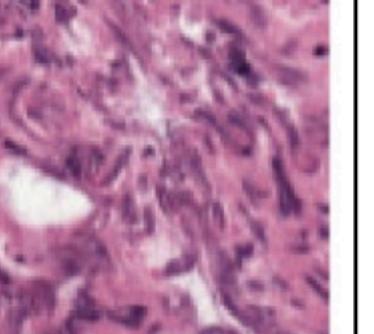
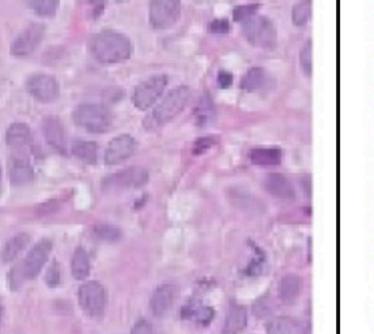
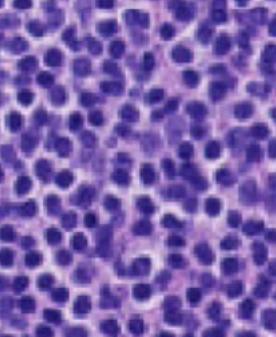
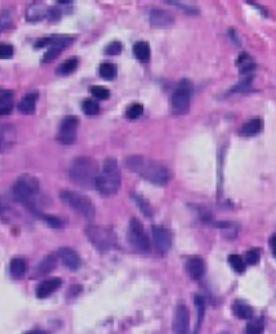


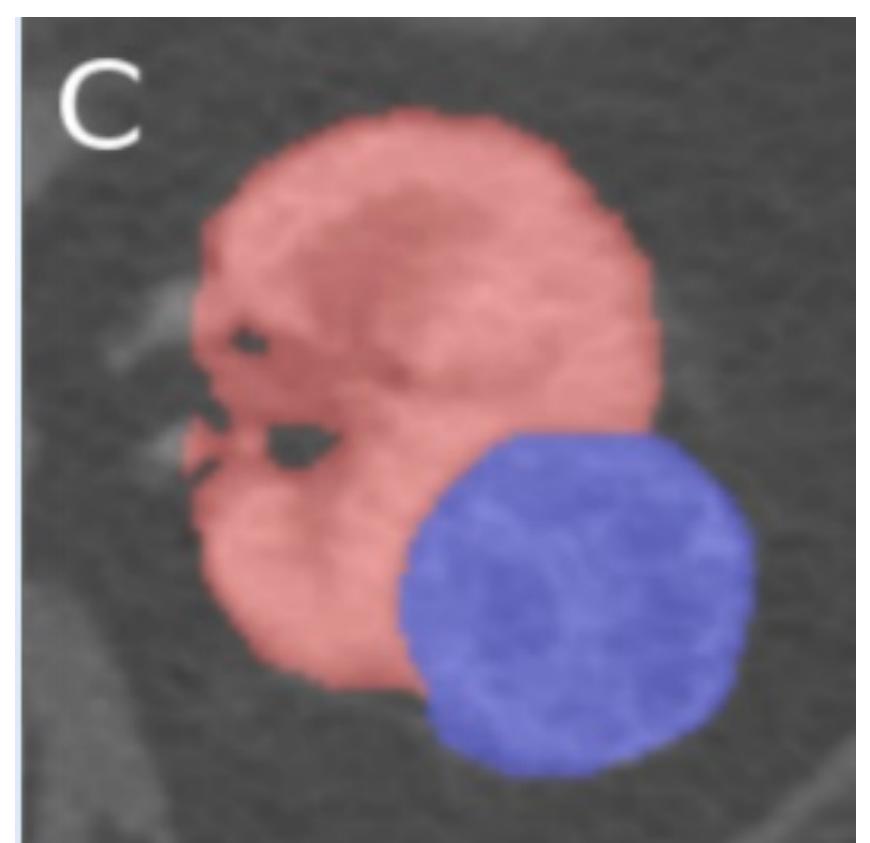
<https://blog.google/products/gmail/gmail-ai-features/>



For many applications, in-domain data = **user data**



Train			Val (OOD)	Test (OOD)
$d = \text{Hospital 1}$	$d = \text{Hospital 2}$	$d = \text{Hospital 3}$	$d = \text{Hospital 4}$	$d = \text{Hospital 5}$
				
$y = \text{Normal}$				
$y = \text{Tumor}$				



Digital Health Laws and Regulations India 2024

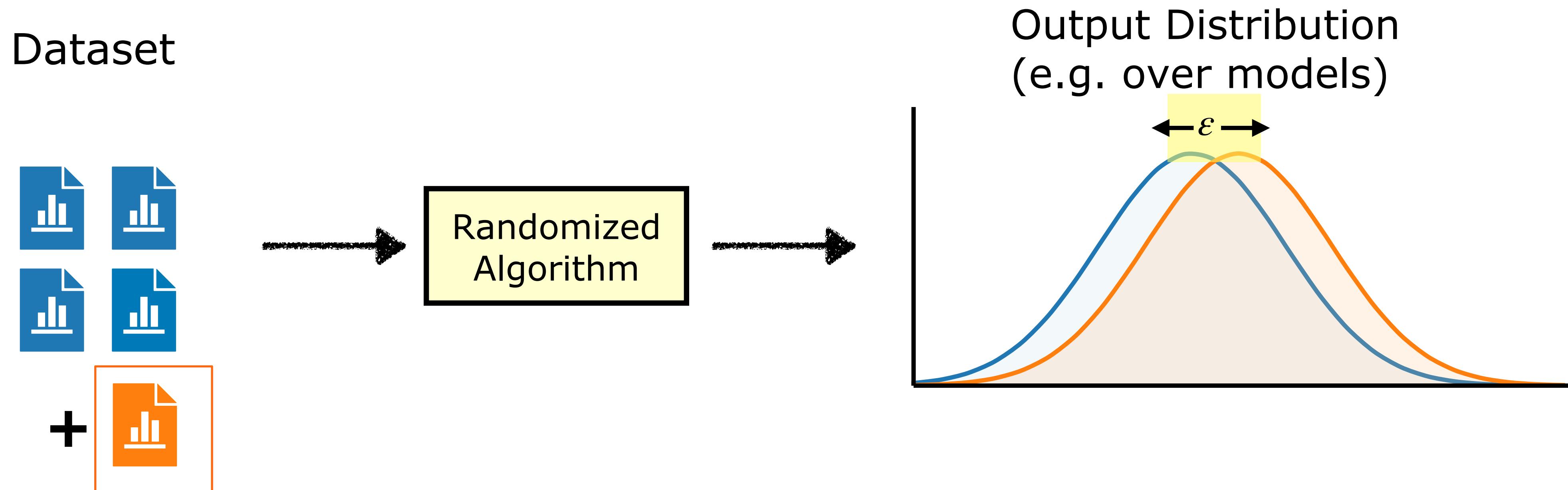


Figure 4: The CAMELYON17-WILDS dataset comprises tissue patches from different hospitals. The goal is to accurately predict the presence of tumor tissue in patches taken from hospitals that are not in the training set. In this figure, each column contains two patches, one of normal tissue and the other of tumor tissue from the same slide.

For many applications, in-domain data = **user data**

Example-level Differential privacy (DP)

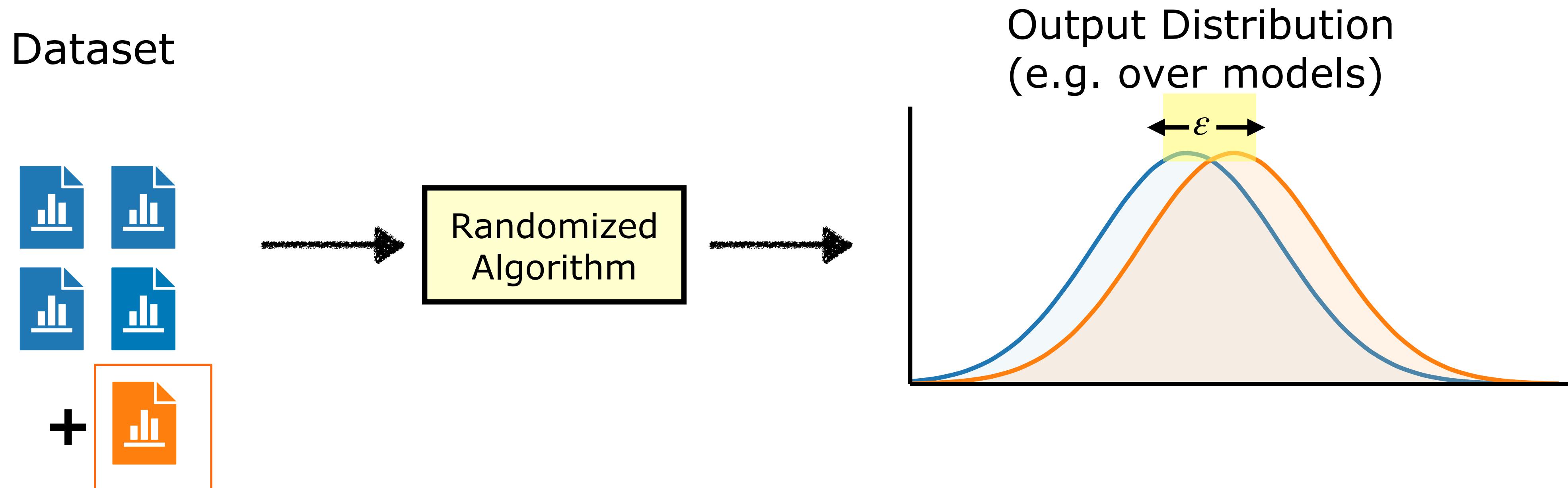
 Unit of data
= **example**



A randomized algorithm is **ε -differentially private** if the addition of **one example** does not alter its output distribution by more than ε

User Example-level Differential privacy (DP)

 Unit of data
= user



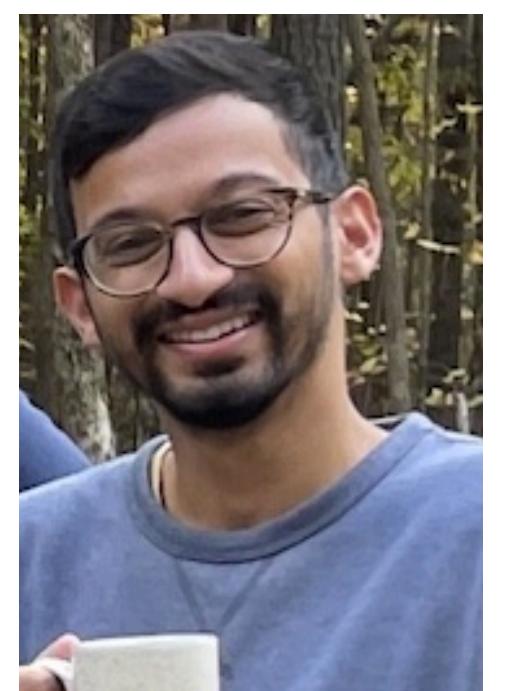
A randomized algorithm is **ε -differentially private** if the addition of **one user's data** does not alter its output distribution by more than ε

Why do we need user-level DP?

Why do we need user-level DP?

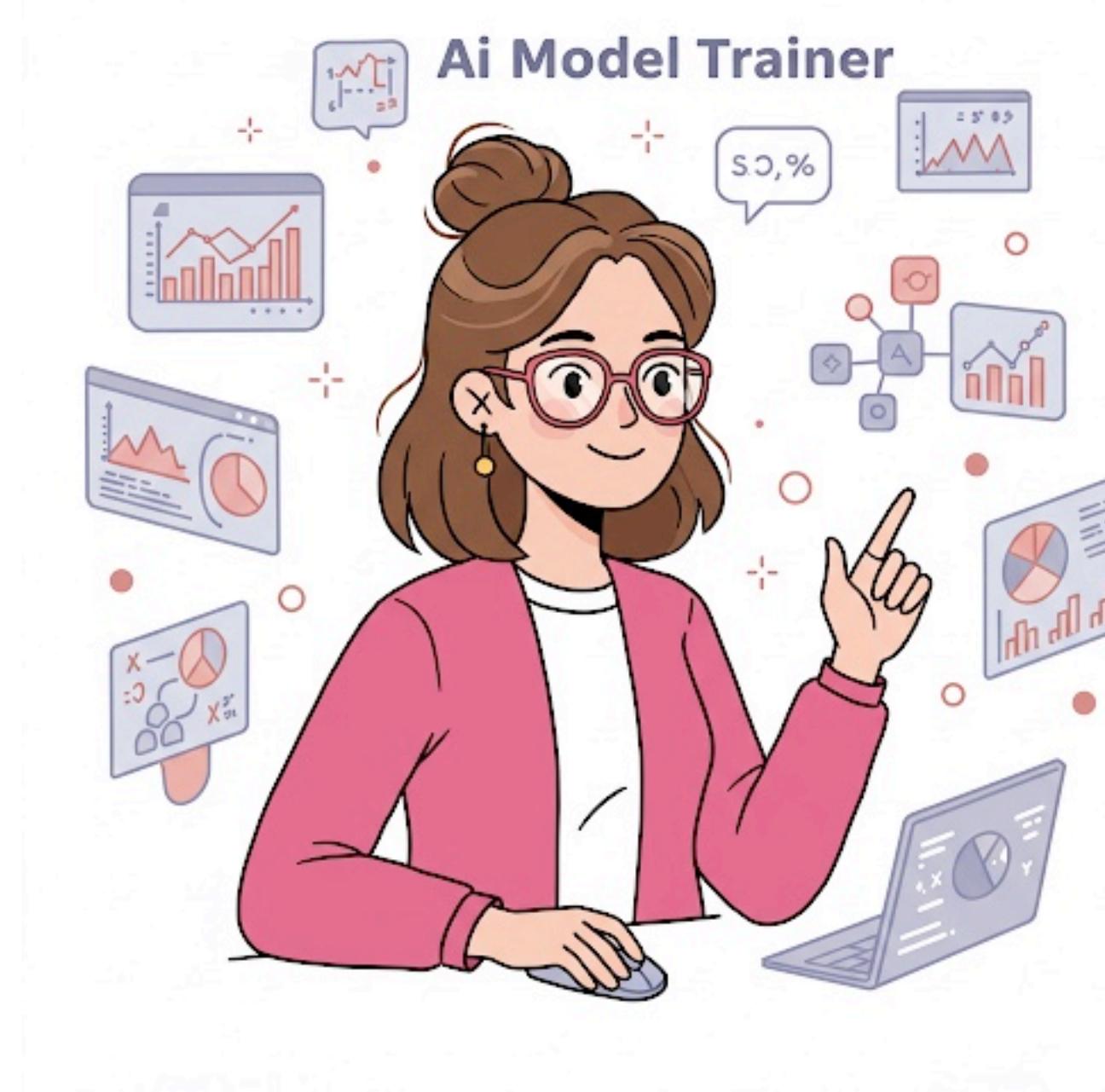
*Standard LLM finetuning pipelines are susceptible to
user inference attacks!*

Nikhil Kandpal, **P.**, Alina Oprea, Peter Kairouz, Chris Choquette-Choo, Zheng Xu.
EMNLP (2024) Oral



User Inference Attack

You:
Train AI model

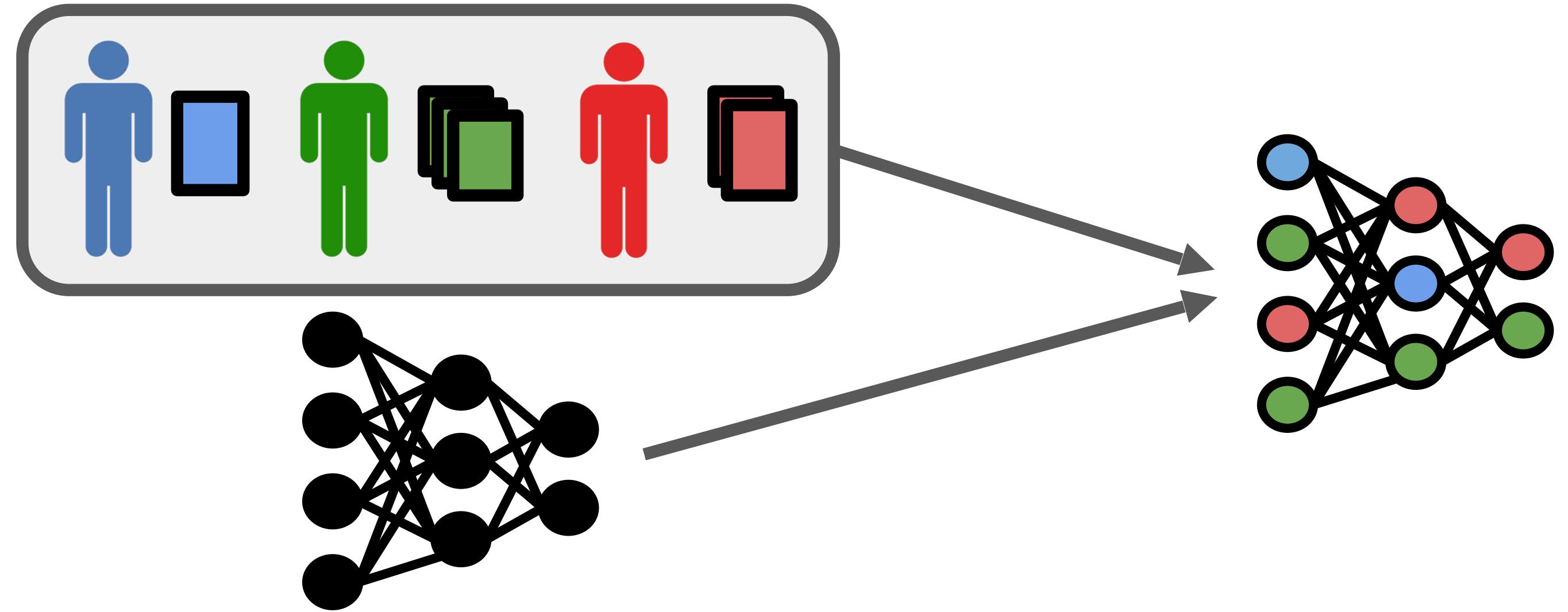


Adversary:
“Attack” the model



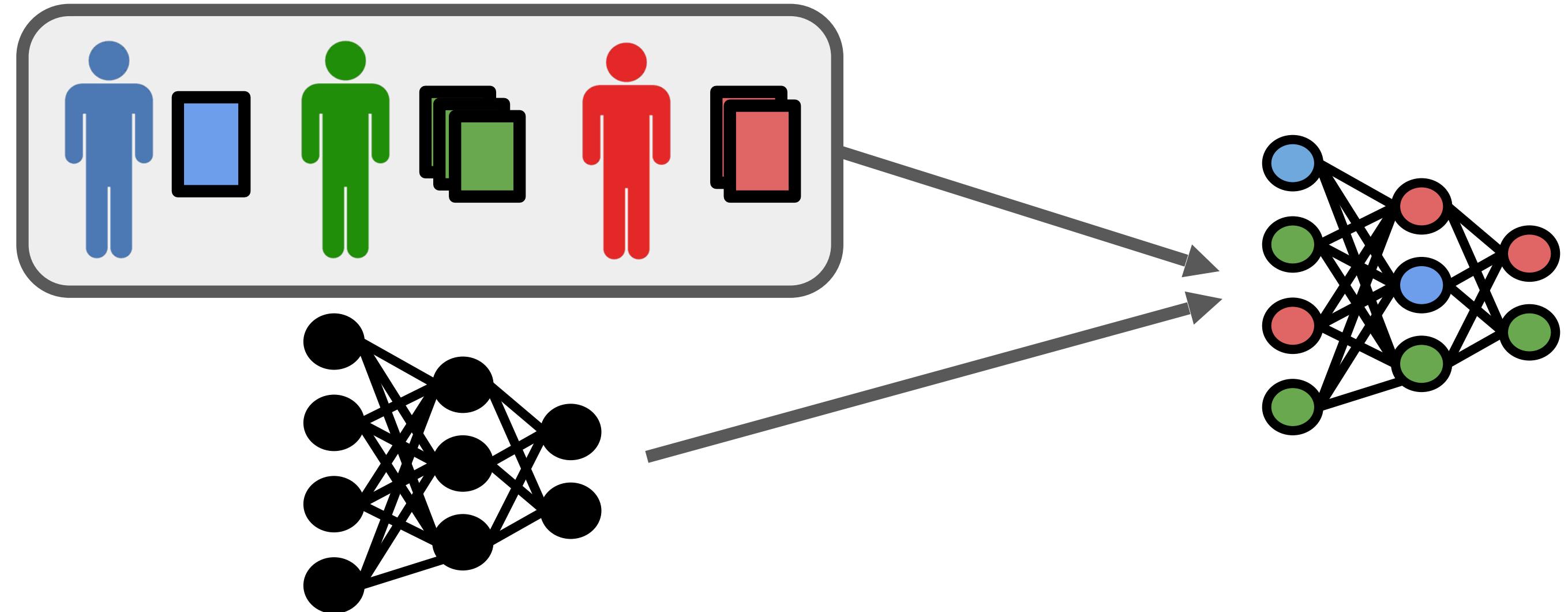
User Inference Attack

Model fine-tuned on user data

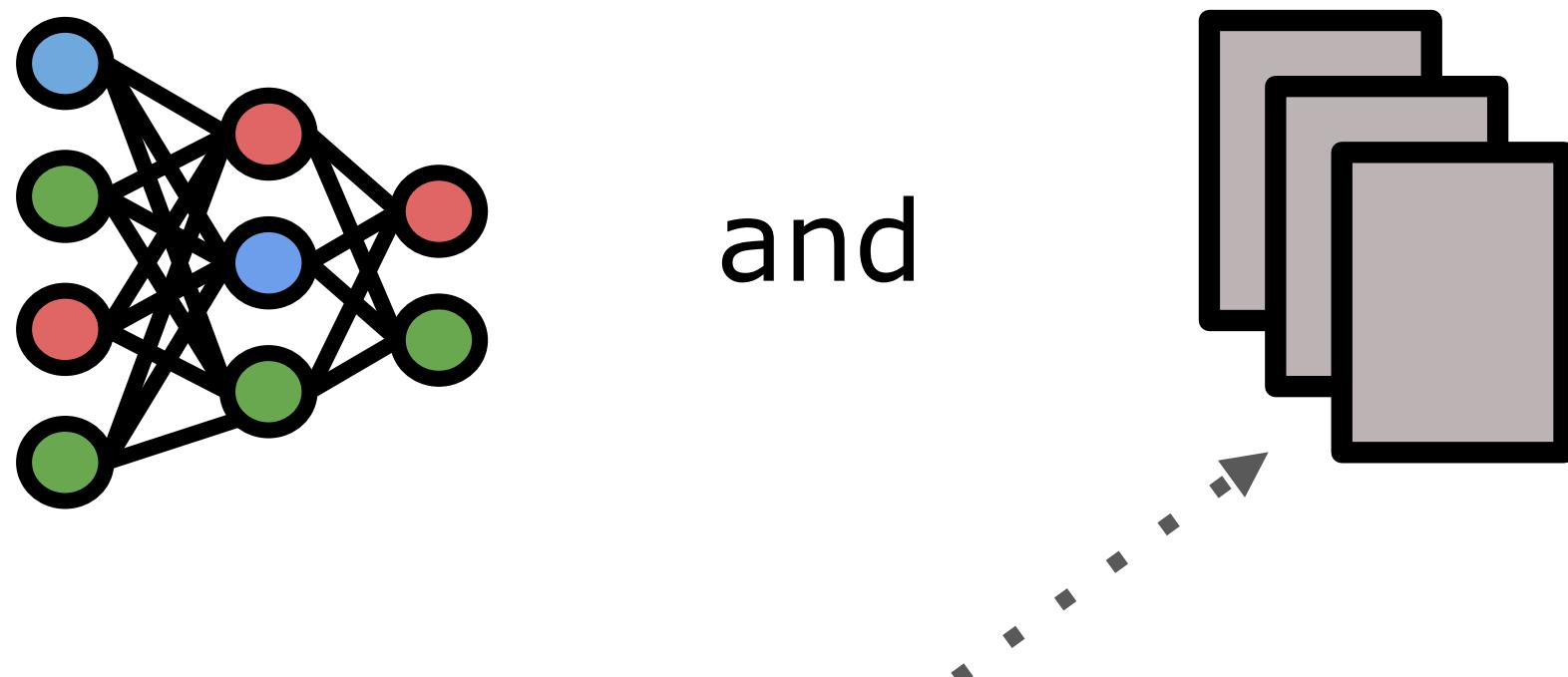


User Inference Attack

Model fine-tuned on user data



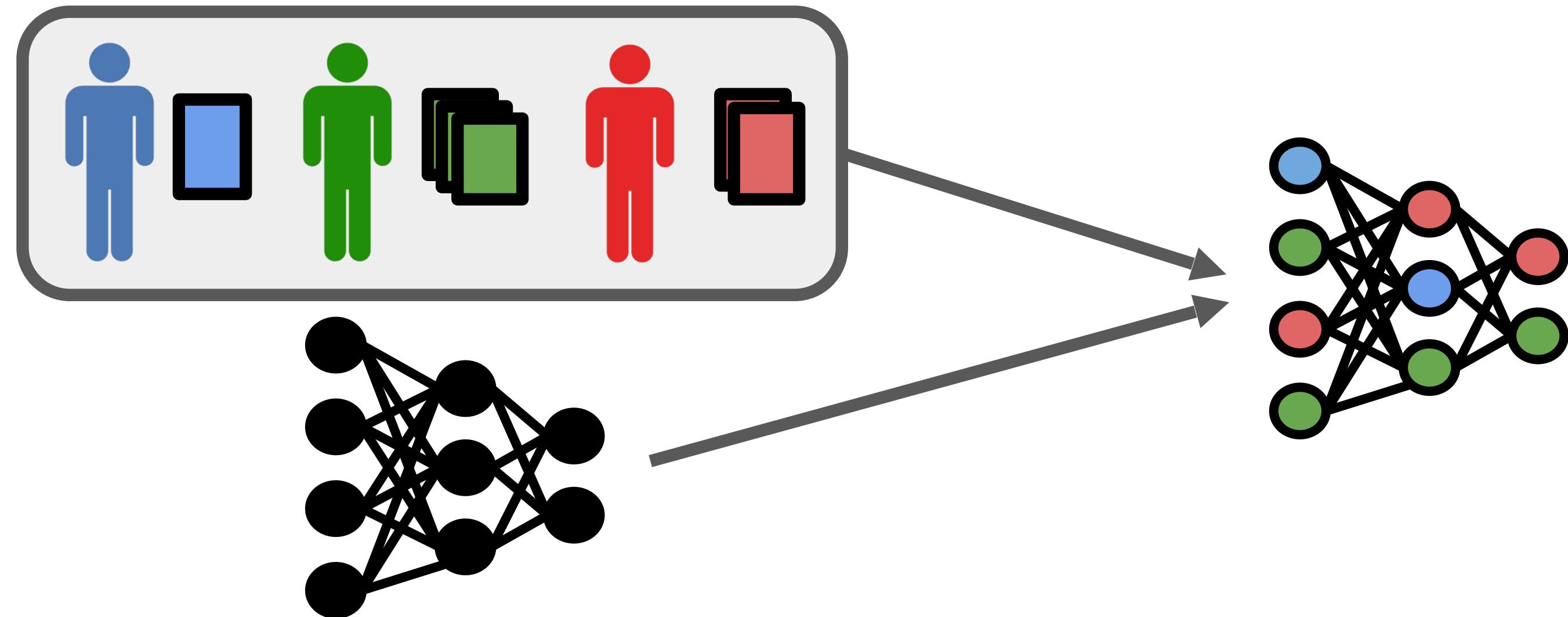
Adversary Has:



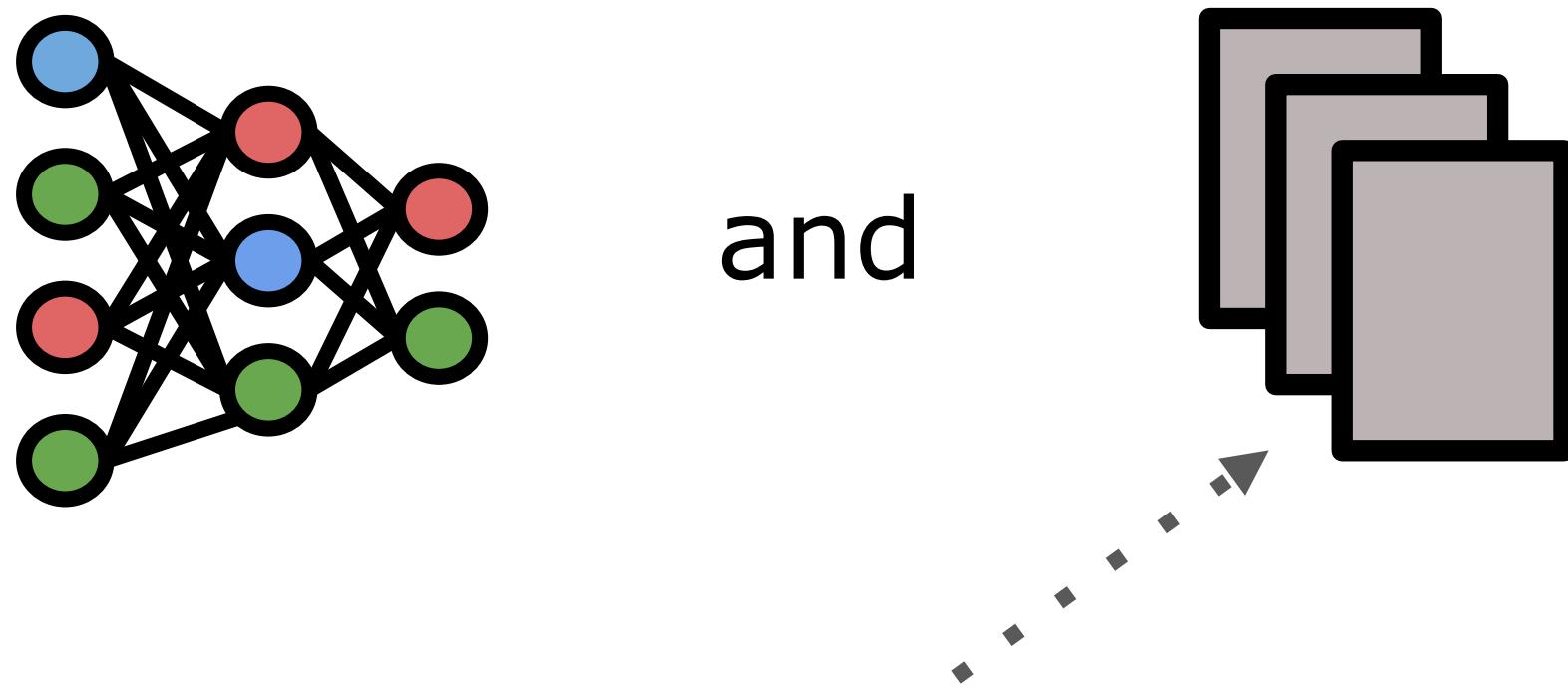
fresh i.i.d. samples from
a user distribution

User Inference Attack

Model fine-tuned on user data



Adversary Has:

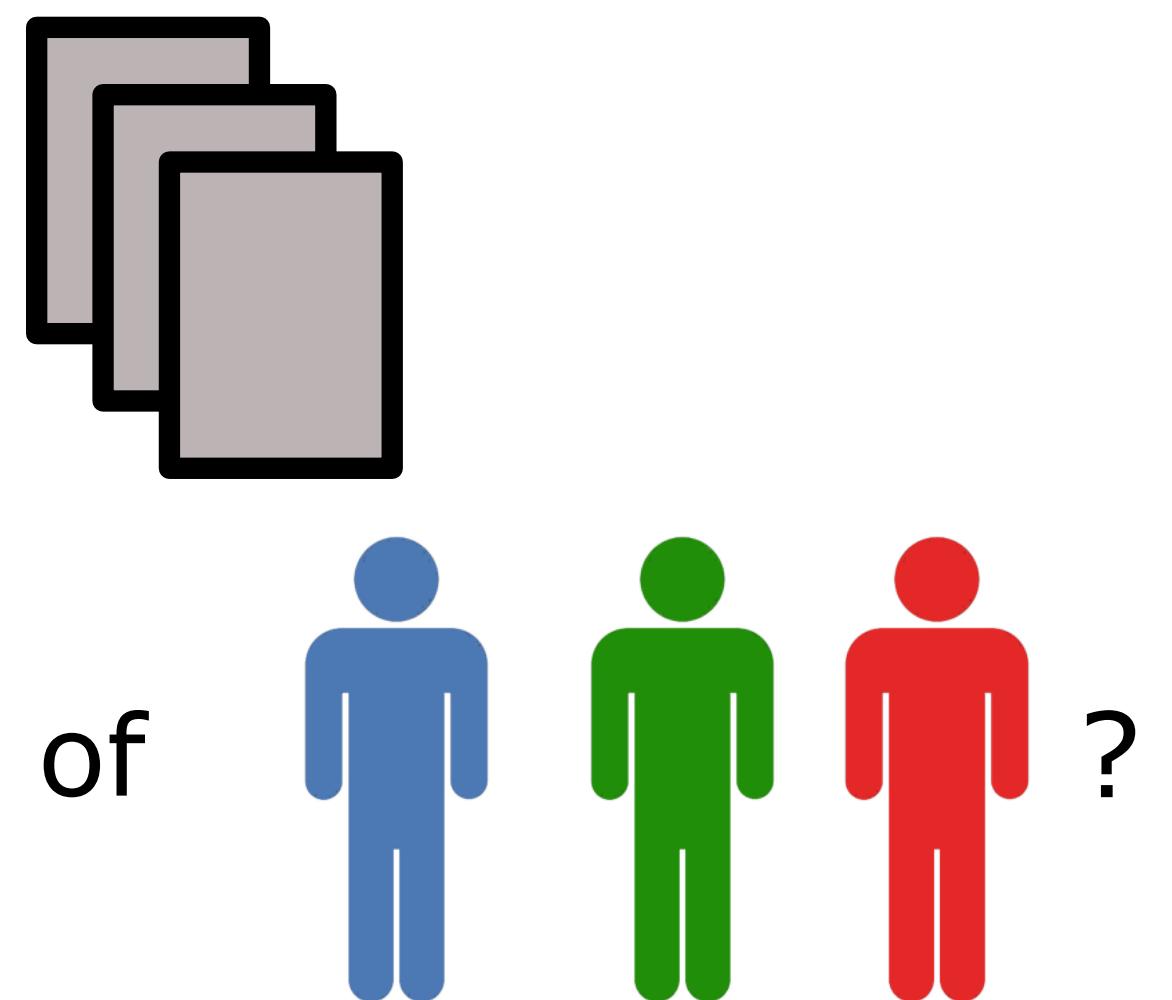


fresh i.i.d. samples from
a user distribution

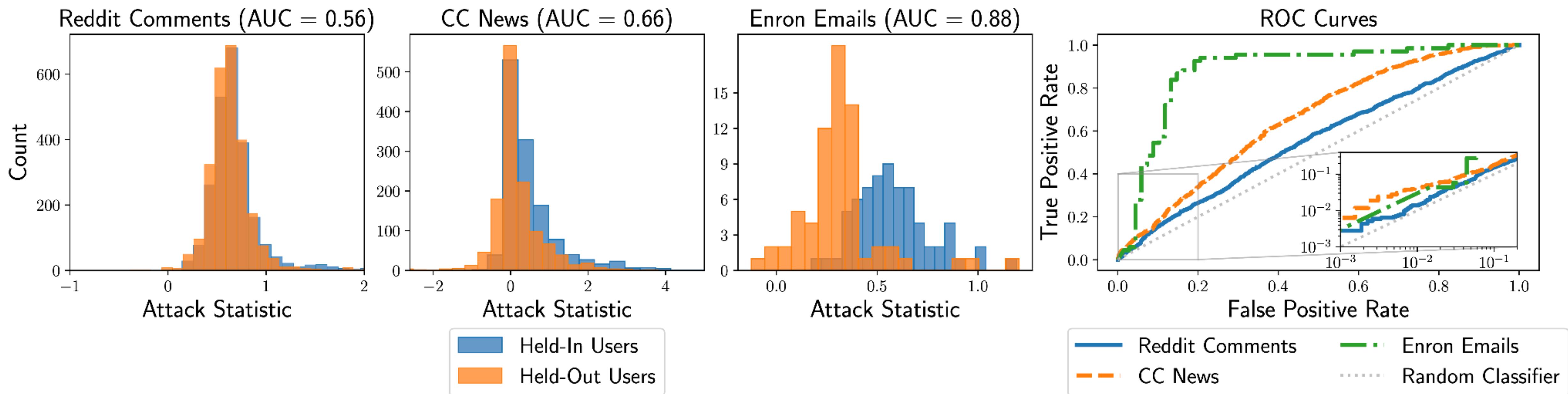
Adversary Wants to Infer:

Did samples

come from one of



User inference is effective when
#users is small and data per user is large

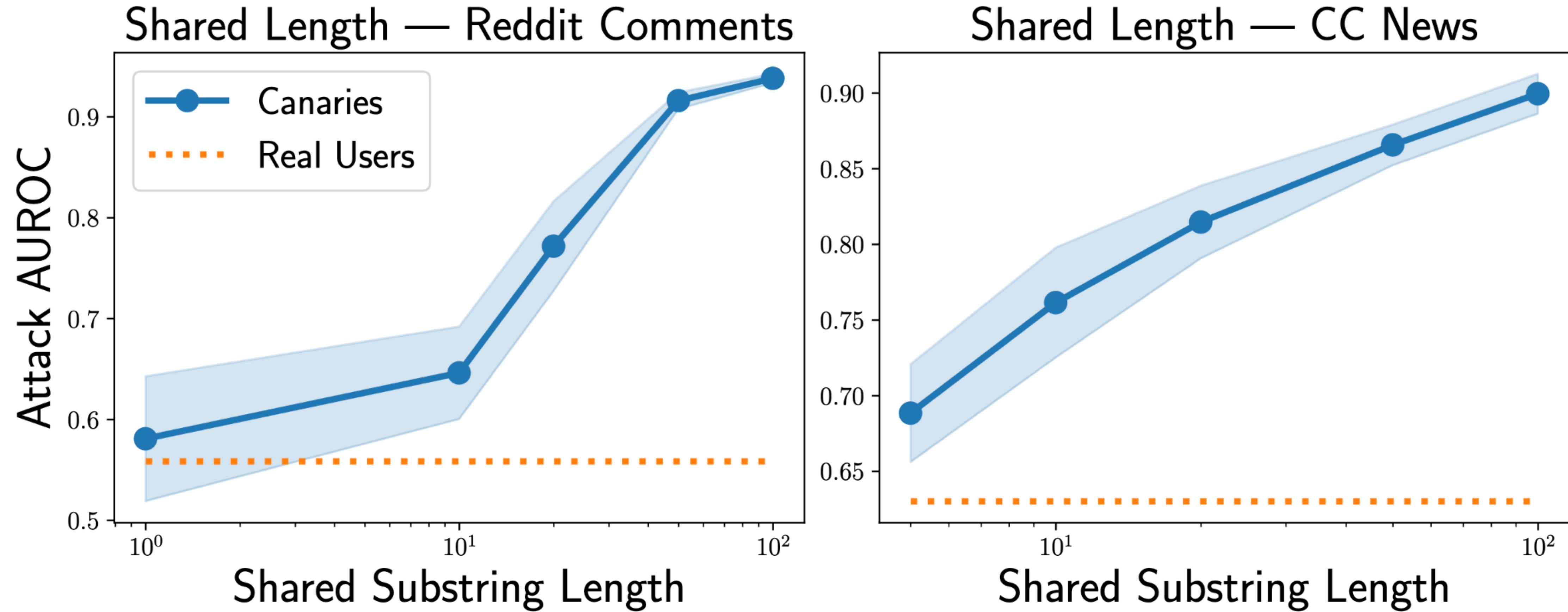


More fine-tuning samples per user



More users

Short common phrases can exacerbate user inference



Example-level DP offers limited mitigation

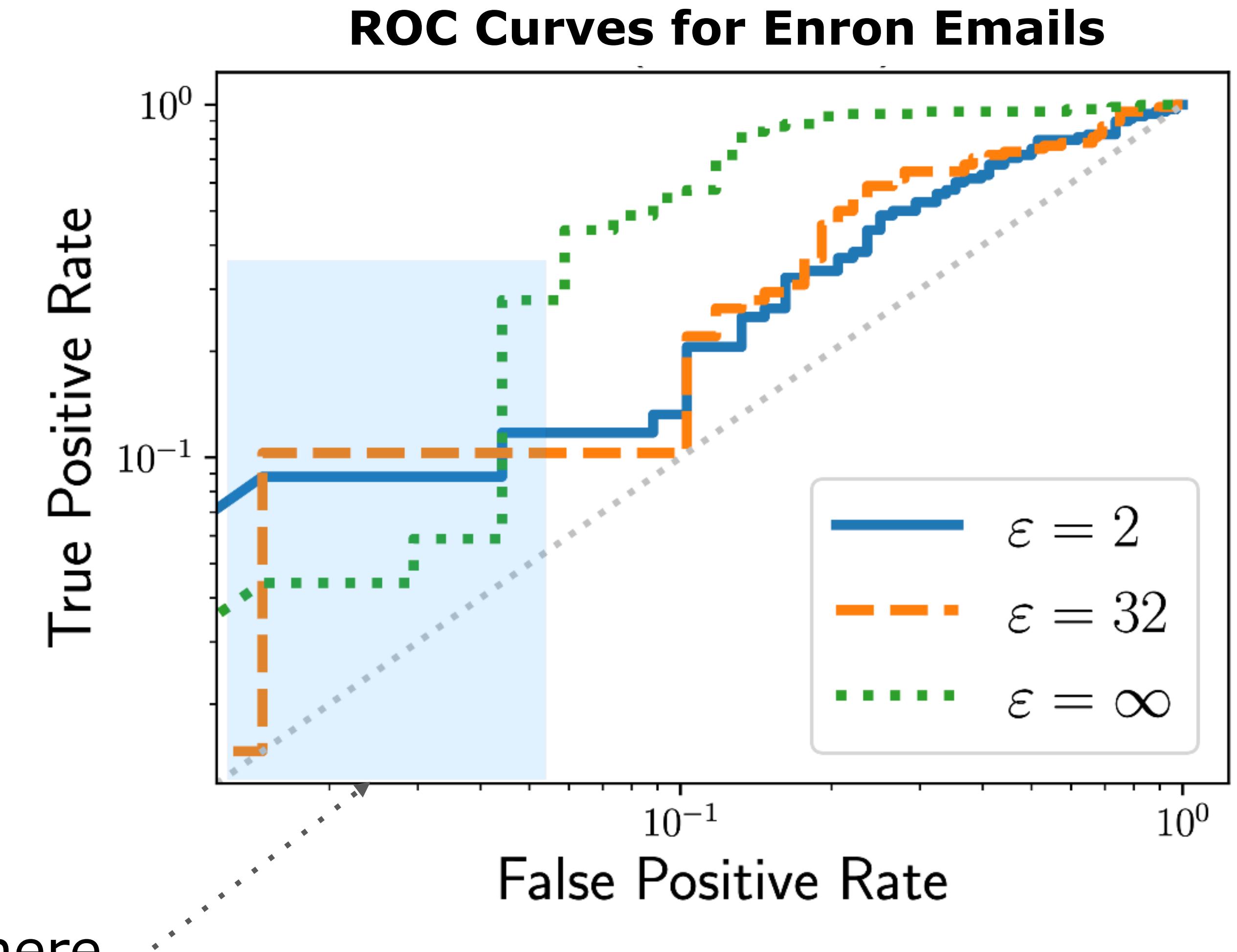
AUROC:

- non-private: 88%
- $\epsilon = 32$: 70%

Utility:

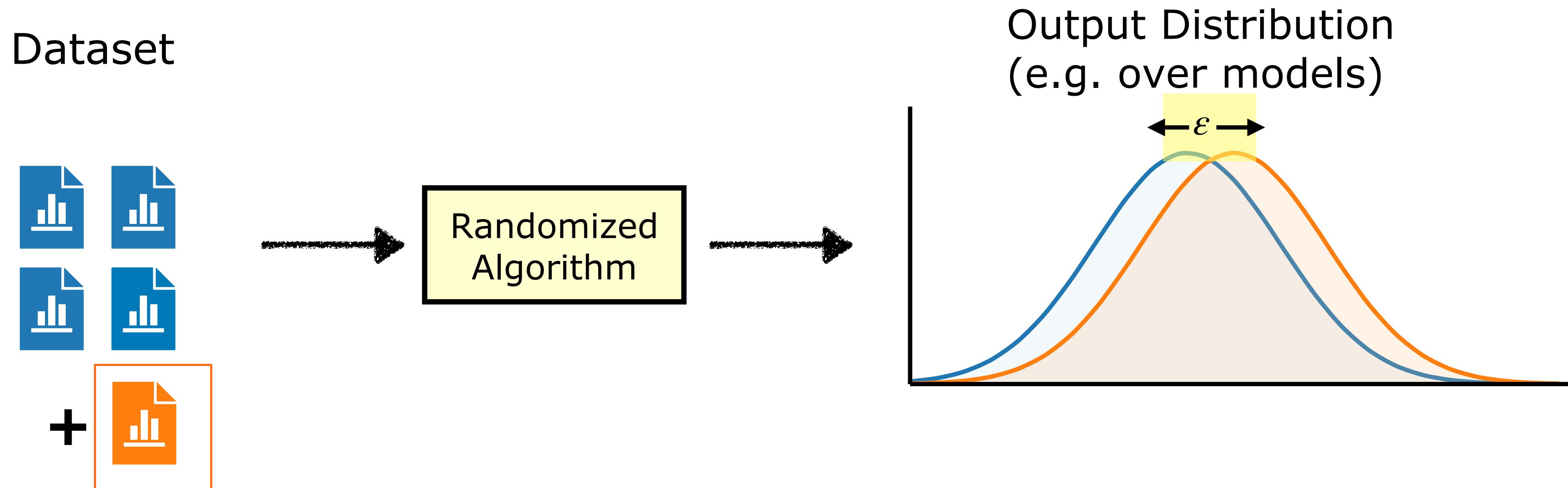
- DP model reaches what the private model achieves in 1/3 epoch

Example-level DP does not help here



User Example-level Differential privacy (DP)

 Unit of data
= user



A randomized algorithm is **ε -differentially private** if the addition of **one user's data** does not alter its output distribution by more than ε

How do we realize user-level DP?

Outline: how do we realize user-level DP?

Learning algorithms:

Improve the runtime of SoTA correlated noise algorithms from $O(n^2)$ to $O(n \log^2 n)$ at the same performance

(n = number of steps)

Noise	Error	Time / step
Independent	$\Theta(\sqrt{n})$	$O(1)$
Optimal	$\frac{\log(n)}{\pi}$	$O(n)$
Correlated		
Ours	$\frac{\log(n)}{\pi} + c$	$O(\log^2(n/c))$

Outline: how do we realize user-level DP?

Learning algorithms:

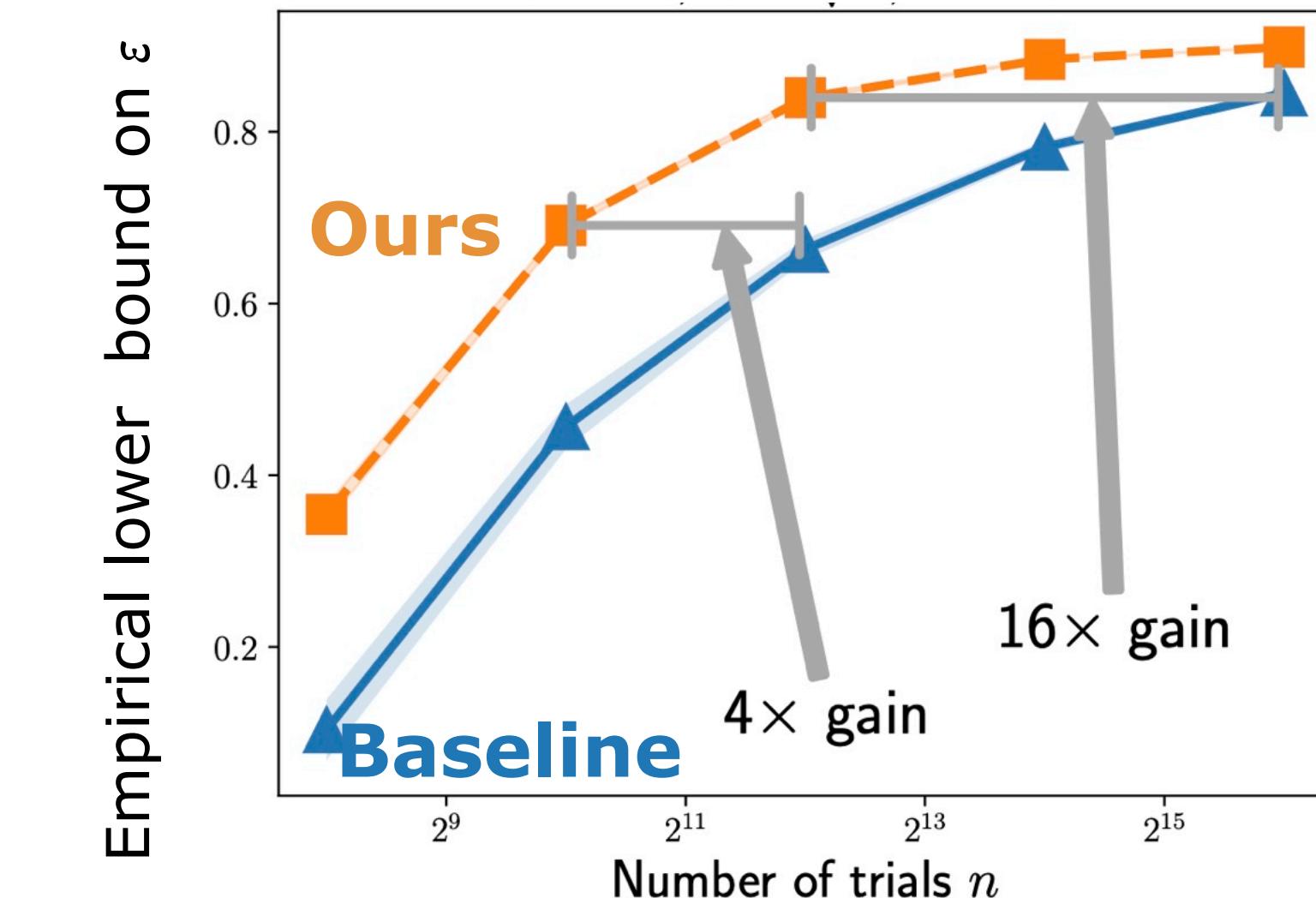
Improve the runtime of SoTA correlated noise algorithms from $O(n^2)$ to $O(n \log^2 n)$ at the same performance

(n = number of steps)

Noise	Error	Time / step
Independent	$\Theta(\sqrt{n})$	$O(1)$
Optimal	$\frac{\log(n)}{\pi}$	$O(n)$
Correlated		
Ours	$\frac{\log(n)}{\pi} + c$	$O(\log^2(n/c))$

Auditing:

Randomness makes the audit more computationally efficient



Part 1: Faster learning algorithms

FOCS (2024)



Dj
Dvijotham



Brendan
McMahan



Krishna
Pillutla



Thomas
Steinke



Abhradeep
Thakurta

DP-SGD: How do we train models with **example-level DP**?

Stochastic gradient
clipped to $\|g\|_2 \leq 1$
per-example

Independent
Gaussian noise

$$\theta_{t+1} = \theta_t - \eta (g_t + z_t)$$

Learning
rate

DP-FedAvg: How do we train models with **user-level DP**?

Stochastic gradient
clipped to $\|g\|_2 \leq 1$
per-user

Independent
Gaussian noise

$$\theta_{t+1} = \theta_t - \eta (g_t + z_t)$$

Learning
rate

DP-SGD: DP Training with *Independent* Noise

Independent
Gaussian noise

$$\theta_{t+1} = \theta_t - \eta (g_t + z_t)$$

DP-FTRL: DP Training with *Correlated* Noise

The diagram illustrates the DP-FTRL update rule:

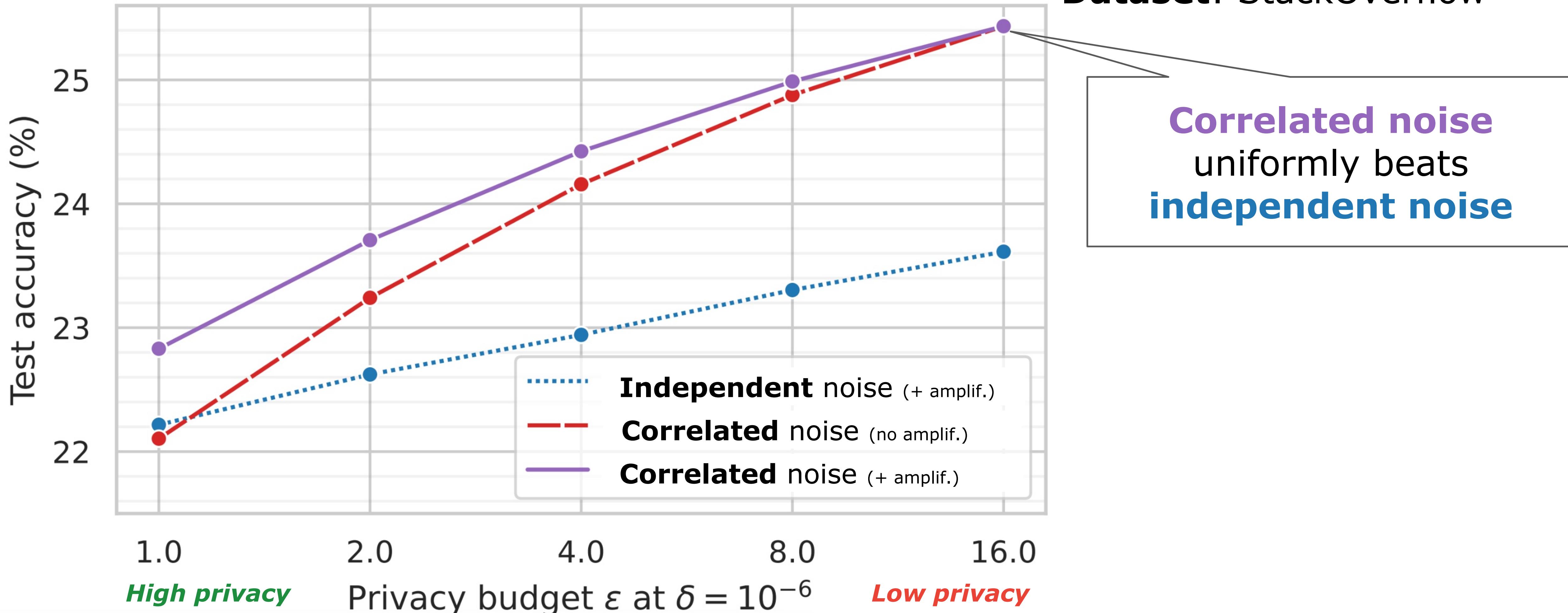
$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau} \right)$$

A callout box with a yellow background and black border contains the text: **(Anti-)correlated Gaussian noise** (z_t i.i.d. Gaussian). A grey arrow points from this box to the term z_t in the update rule.

Prior work: (Empirically) correlated noise outperforms independent noise

Experiment: user-level DP + language modeling

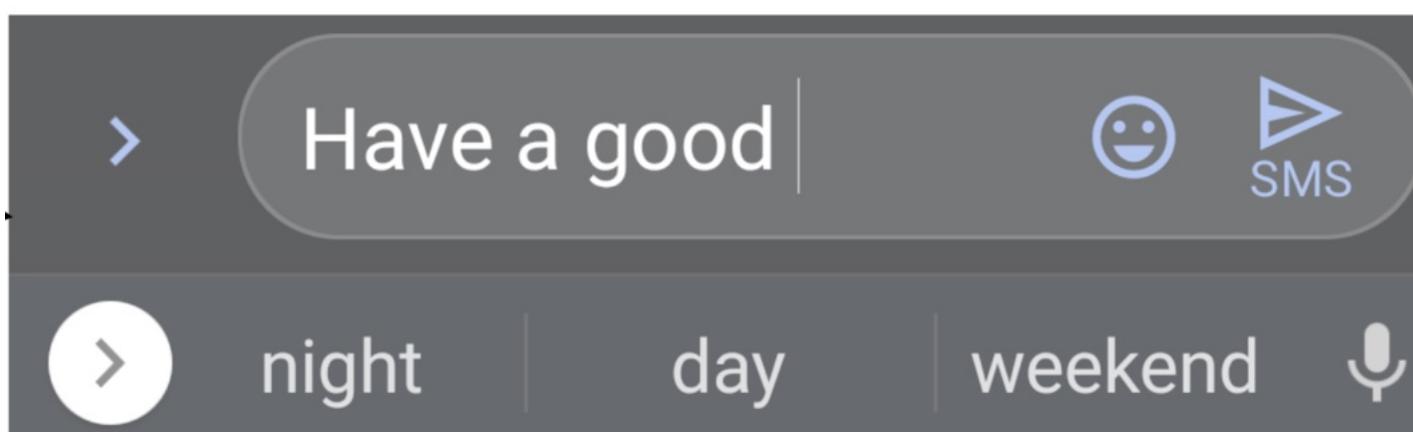
Dataset: StackOverflow



Production Training

"the first production neural network trained directly on user data announced with a formal DP guarantee."

- [Google AI Blog post](#), Feb 2022



The latest from Google Research

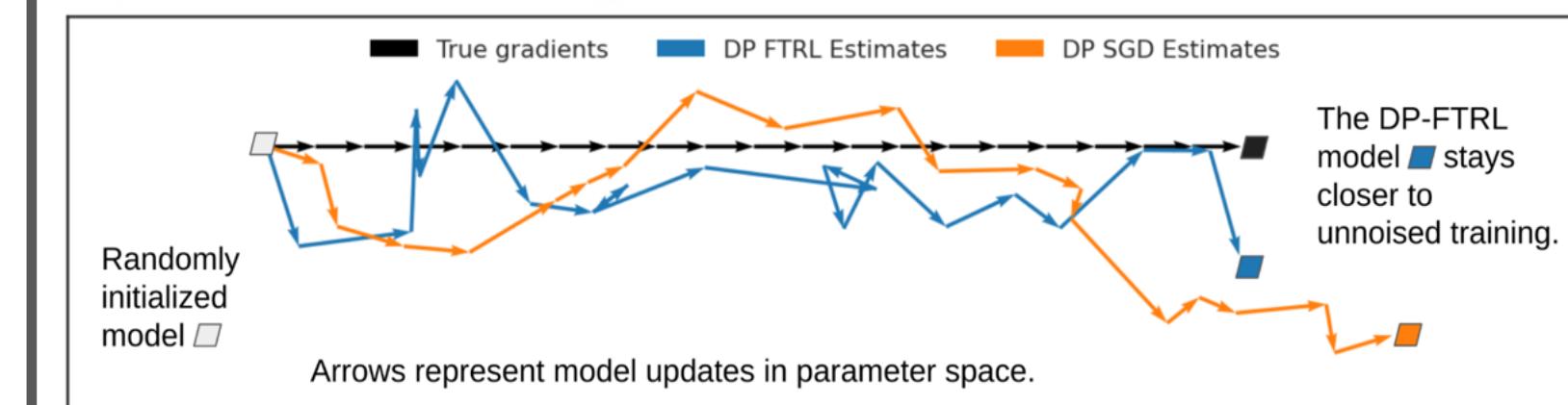
Federated Learning with Formal Differential Privacy Guarantees

Monday, February 28, 2022

Posted by Brendan McMahan and Abhradeep Thakurta, Research Scientists, Google Research

In 2017, Google [introduced federated learning \(FL\)](#), an approach that enables mobile devices to collaboratively train machine learning (ML) models while keeping the raw training data on each user's device, decoupling the ability to do ML from the need to store the data in the cloud. Since its introduction, Google has continued to [actively engage in FL research](#) and deployed FL to power many features in [Gboard](#), including next word prediction, emoji suggestion and out-of-vocabulary word discovery. Federated learning is improving the "Hey Google" detection models in Assistant, [suggesting replies](#) in Google Messages, [predicting text selections](#), and more.

While FL allows ML without raw data collection, [differential privacy \(DP\)](#) provides a quantifiable measure of data anonymization, and when applied to ML can address concerns about models memorizing sensitive user data. This too has been a top research priority, and has yielded one of the first production uses of DP for analytics with [RAPPOR](#) in 2014, [our open-source DP library](#), [Pipeline DP](#), and [TensorFlow Privacy](#).



Data Minimization and Anonymization in Federated Learning

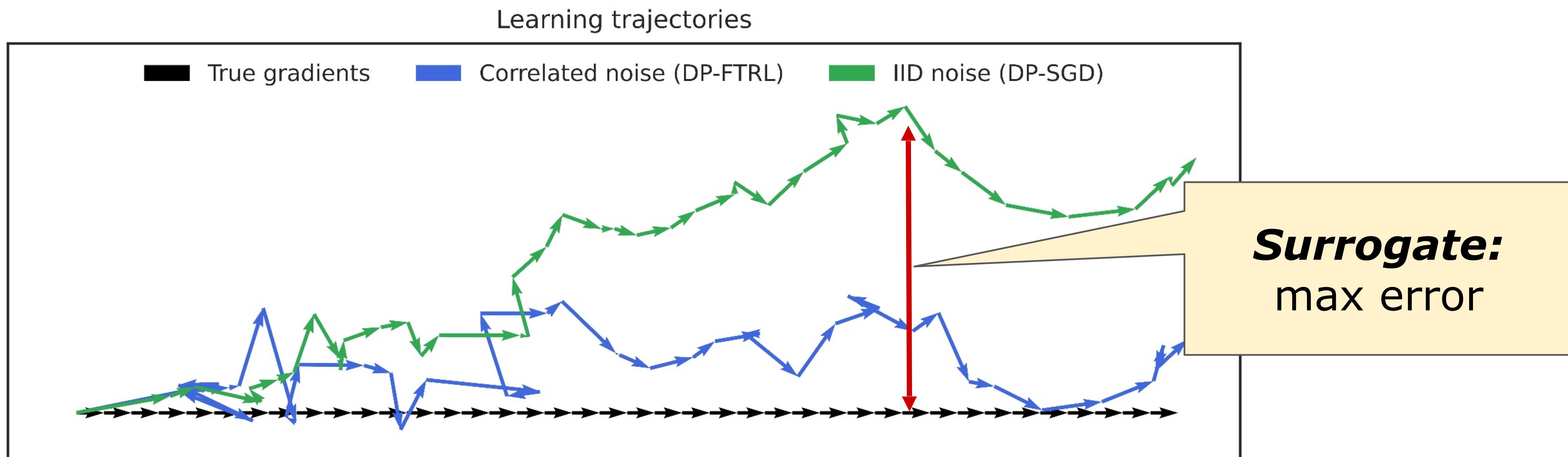
Along with fundamentals like transparency and consent, the [privacy principles of data minimization and anonymization](#) are important in ML applications that involve sensitive data.

How do we find the noise coefficients?

How do we find the noise coefficients?

Current Approach:

Find the noise coefficients to ***minimize the cumulative noise*** added to the learning trajectory (such that a given DP constraint is satisfied)



How do we find the noise coefficients?

$$\theta_{t+1} = \theta_t - \eta \left(g_t + \underbrace{z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau}}_{=:w_t} \right)$$

Find the noise coefficients β_t to **minimize the max error** (i.e. cumulative noise added to the learning trajectory):

**Surrogate
Objective**

$$\mathcal{E}(\beta)^2 = \max_{t \leq n} \mathbb{E}_{z_\tau \sim \mathcal{N}(0, \sigma^2 I)} \left\| \sum_{\tau=0}^t w_\tau \right\|_2^2$$

where the variance σ^2 is chosen so that θ_t 's satisfy a given DP constraint

Toeplitz mechanism: optimal max error

$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \boxed{\beta_\tau} z_{t-\tau} \right)$$

Theorem

[Fichtenberger, Henzinger, Upadhyay (ICML '23); Dvijotham, McMahan, **P.**, Steinke, Thakurta (FOCS '24)]

For any number n of steps, the optimal max error is obtained by coefficients $\beta_t^* = t^{-3/2}$ and satisfies the bounds

$$\mathcal{E}(\beta^*) = \frac{\log n}{\pi} + \text{constant}$$

Toeplitz mechanism: optimal max error

$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \boxed{\beta_\tau} z_{t-\tau} \right)$$

Theorem

[Fichtenberger, Henzinger, Upadhyay (ICML '23); Dvijotham, McMahan, **P.**, Steinke, Thakurta (FOCS '24)]

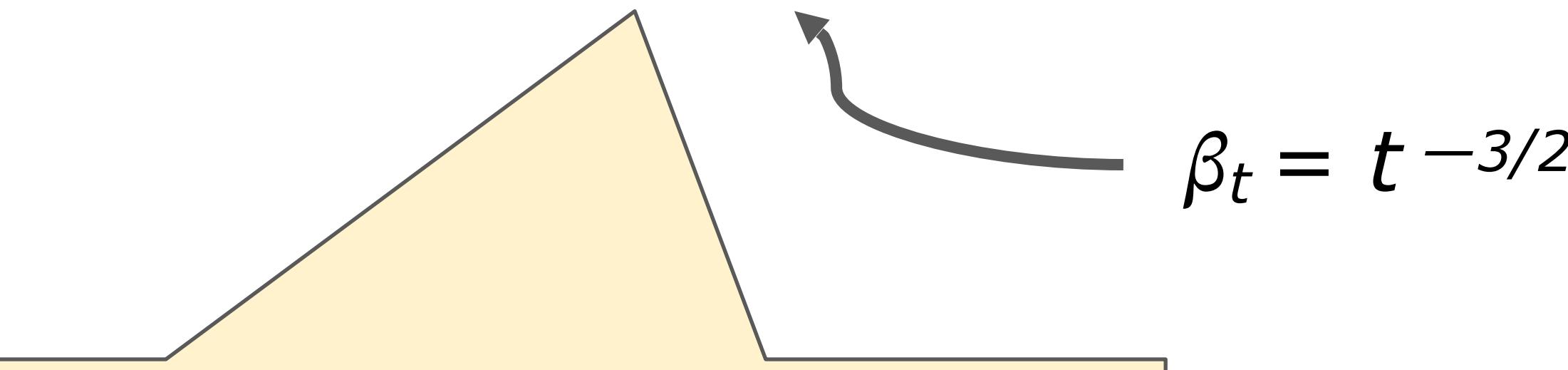
For any number n of steps, the optimal max error is obtained by coefficients $\beta_t^* = t^{-3/2}$ and satisfies the bounds

$$\mathcal{E}(\beta^*) = \frac{\log n}{\pi} + \text{constant}$$

$$\mathcal{E}(\beta^{\text{SGD}}) = \Theta(\sqrt{n})$$

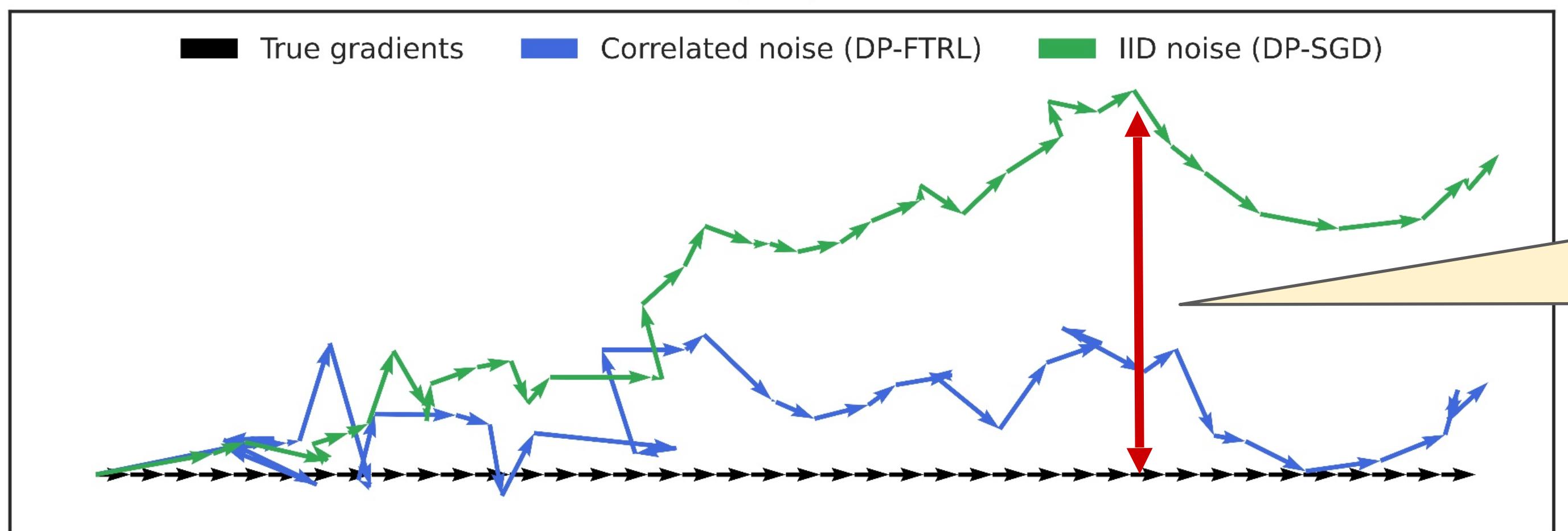
Exponential
improvement over
independent noise

Our challenge: running time

$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau} \right)$$

$$\beta_t = t^{-3/2}$$

Quadratic time complexity:
Noise generation requires $O(t)$ time in iteration t

Learning trajectories



Max Error

Noise generation time
(in iteration t)

Independent noise

$$\Theta(\sqrt{n})$$

$$O(\dim)$$

Optimal correlated noise

$$\frac{\log n}{\pi} + c$$

$$O(t \cdot \dim)$$

A first attempt: the banded mechanism

Set $\beta_t = 0$ for $t > b$

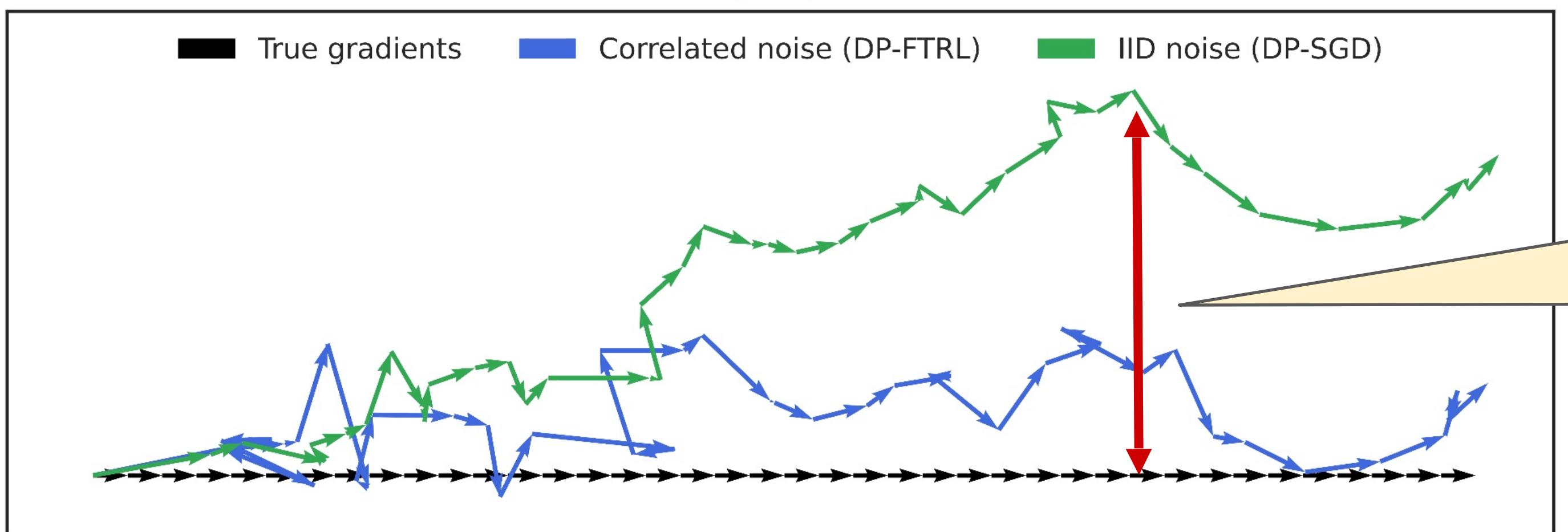
Then, we only have to sum b terms in

$$\sum_{\tau=1}^b \beta_\tau z_{t-\tau}$$

Linear complexity:

Noise generation requires $O(b)$ time in each iteration

Learning trajectories



Max Error

Noise generation time
(in iteration t)

Independent noise

$$\Theta(\sqrt{n})$$

$$O(\dim)$$

Optimal correlated noise

$$\frac{\log n}{\pi} + c$$

$$O(t \cdot \dim)$$

b -Banded

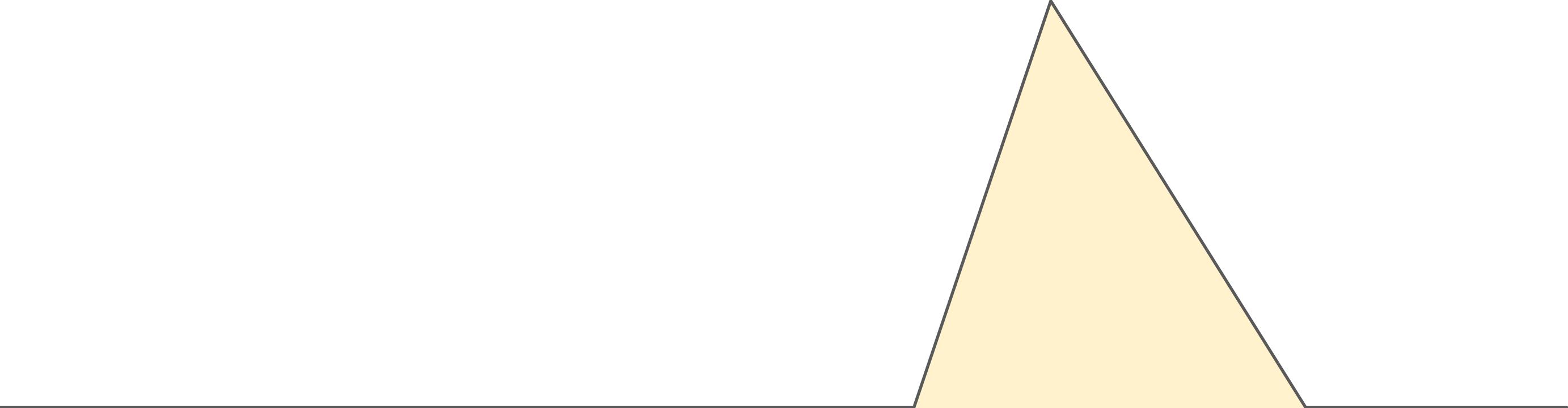
$$O\left((\sqrt{n/b} - 1) \log b\right)$$

$$O(b \cdot \dim)$$

Our approach: Intuition

Consider an exponentially decaying sequence $\beta_t = \alpha \lambda^{t-1}$.

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$
using the recurrence $w_{t+1} = \alpha z_t + \lambda w_{t-1}$



Linear complexity:

Noise generation requires $O(\dim)$ time in each iteration

Our approach: Intuition

Consider an exponentially decaying sequence $\beta_t = \alpha \lambda^{t-1}$.

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$
using the recurrence $w_{t+1} = \alpha z_t + \lambda w_{t-1}$



Linear complexity:

Noise generation requires $O(\dim)$ time in each iteration

Our approach: Intuition

Consider sums of exponentials:

$$\beta_t = \alpha_1 \lambda_1^{t-1} + \alpha_2 \lambda_2^{t-1}$$

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$

using

Our approach: Intuition

Consider sums of exponentials:

$$\beta_t = \alpha_1 \lambda_1^{t-1} + \alpha_2 \lambda_2^{t-1}$$

The diagram illustrates the decomposition of β_t into two components. A blue square labeled β_t' and an orange square labeled β_t'' are positioned above a larger blue rectangle containing the term $\alpha_1 \lambda_1^{t-1}$. A larger orange rectangle contains the term $\alpha_2 \lambda_2^{t-1}$. Arrows point from both the blue and orange squares down to their respective components in the main equation.

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$ using

Our approach: Intuition

Consider sums of exponentials:

$$\beta_t = \alpha_1 \lambda_1^{t-1} + \alpha_2 \lambda_2^{t-1}$$

The diagram illustrates the decomposition of β_t into two components. It shows a blue box containing $\alpha_1 \lambda_1^{t-1}$ and an orange box containing $\alpha_2 \lambda_2^{t-1}$. Above these boxes, a blue box labeled β_t' and an orange box labeled β_t'' are shown. Arrows point from the blue box to the blue box and from the orange box to the orange box, indicating they are separate entities.

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$

using

$$w'_t = \sum_{\tau=1}^t \beta'_\tau z_{t-\tau}$$

$$w''_t = \sum_{\tau=1}^t \beta''_\tau z_{t-\tau}$$

$$w_t = w'_t + w''_t$$

Our approach: Intuition

Consider sums of exponentials:

$$\beta_t = \sum_{t=1}^t \alpha_1 \lambda_1^{t-1} + \alpha_2 \lambda_2^{t-1}$$

The diagram illustrates the decomposition of β_t into two components: β_t' (blue box) and β_t'' (orange box). Arrows point from the terms $\alpha_1 \lambda_1^{t-1}$ and $\alpha_2 \lambda_2^{t-1}$ in the sum to their respective boxes. The sum $\beta_t' + \beta_t''$ is shown above the boxes.

Then, we can compute the correlated noise $w_t = \sum_{\tau=1}^t \beta_\tau z_{t-\tau}$

using

$$w'_t = \sum_{\tau=1}^t \beta'_\tau z_{t-\tau}$$

$$w''_t = \sum_{\tau=1}^t \beta''_\tau z_{t-\tau}$$

$$w_t = w'_t + w''_t$$

Linear time +
space

Our approach: **Buffered Linear Toeplitz (BLT) Mechanism**

Approximate the optimal noise coefficients with d exponentials as

$$\beta_t \approx \beta'_t = \sum_{i=1}^d \alpha_i \lambda_i^{t-1}$$

Time & space complexity:
 $O(d \times \text{dimension})$

Our approach: **Buffered Linear Toeplitz (BLT) Mechanism**

	Max Error	Noise generation time (in iteration t)
Independent noise	$\Theta(\sqrt{n})$	$O(\dim)$
Optimal correlated noise	$\frac{\log n}{\pi} + c$	$O(t \cdot \dim)$
b-Banded	$O\left((\sqrt{n/b} - 1) \log b\right)$	$O(b \cdot \dim)$
BLT of degree d	??????	$O(d \cdot \dim)$

Our approach: Buffered Linear Toeplitz (BLT) Mechanism

	Max Error	Noise generation time (in iteration t)
Independent noise	$\Theta(\sqrt{n})$	$O(\dim)$
Optimal correlated noise	$\frac{\log n}{\pi} + c$	$O(t \cdot \dim)$
b-Banded	$O\left((\sqrt{n/b} - 1) \log b\right)$	$O(b \cdot \dim)$
BLT of degree d	??????	$O(d \cdot \dim)$

Approximation Theory!!

From sequences to functions

$$r(x) = 1 - \beta_1 x - \beta_2 x^2 - \dots$$



From sequences to functions

$$r(x) = 1 - \beta_1 x - \beta_2 x^2 - \dots$$



Taylor expansion around $x = 0$

From sequences to functions

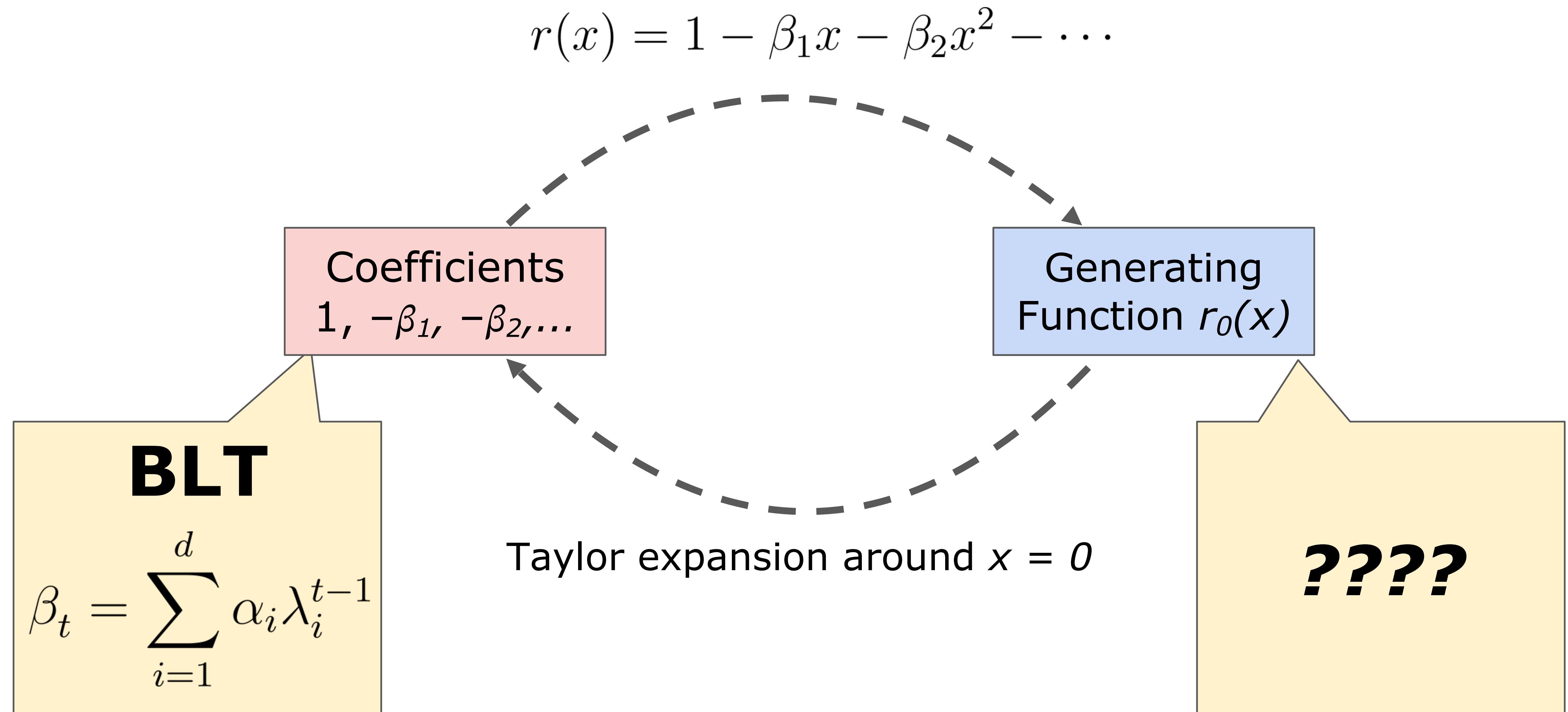
$$r(x) = 1 - \beta_1 x - \beta_2 x^2 - \dots$$



Taylor expansion around $x = 0$

Coefficients $\beta_t = \Theta(t^{-3/2}) \Leftrightarrow$ generating function $r_0(x) = (1 - x)^{1/2}$

From sequences to functions



BLT generating functions

Manuel Kauers
Peter Paule

The Concrete Tetrahedron

Symbolic Sums, Recurrence Equations,
Generating Functions, Asymptotic Estimates

Theorem (Informal):

The following properties are equivalent:

- β 's are a (complex) BLT sequence:
$$\beta_t = \sum_{i=1}^d \alpha_i \lambda_i^{t-1}$$
- Its generating function $r(x)$ is a **rational function** of degree d
- β 's satisfy a linear recurrence
$$\beta_t = \sum_{i=1}^d q_i \beta_{t-i}$$

From functions to efficient noise generation



Theorem [Dvijotham, McMahan, P., Steinke, Thakurta 2024]

The max error of a sequence (β_t) with generating function $r(x)$ is

$$\mathcal{E}(\beta) \leq \frac{\log n}{\pi} + O(n \cdot \text{err}(r))$$

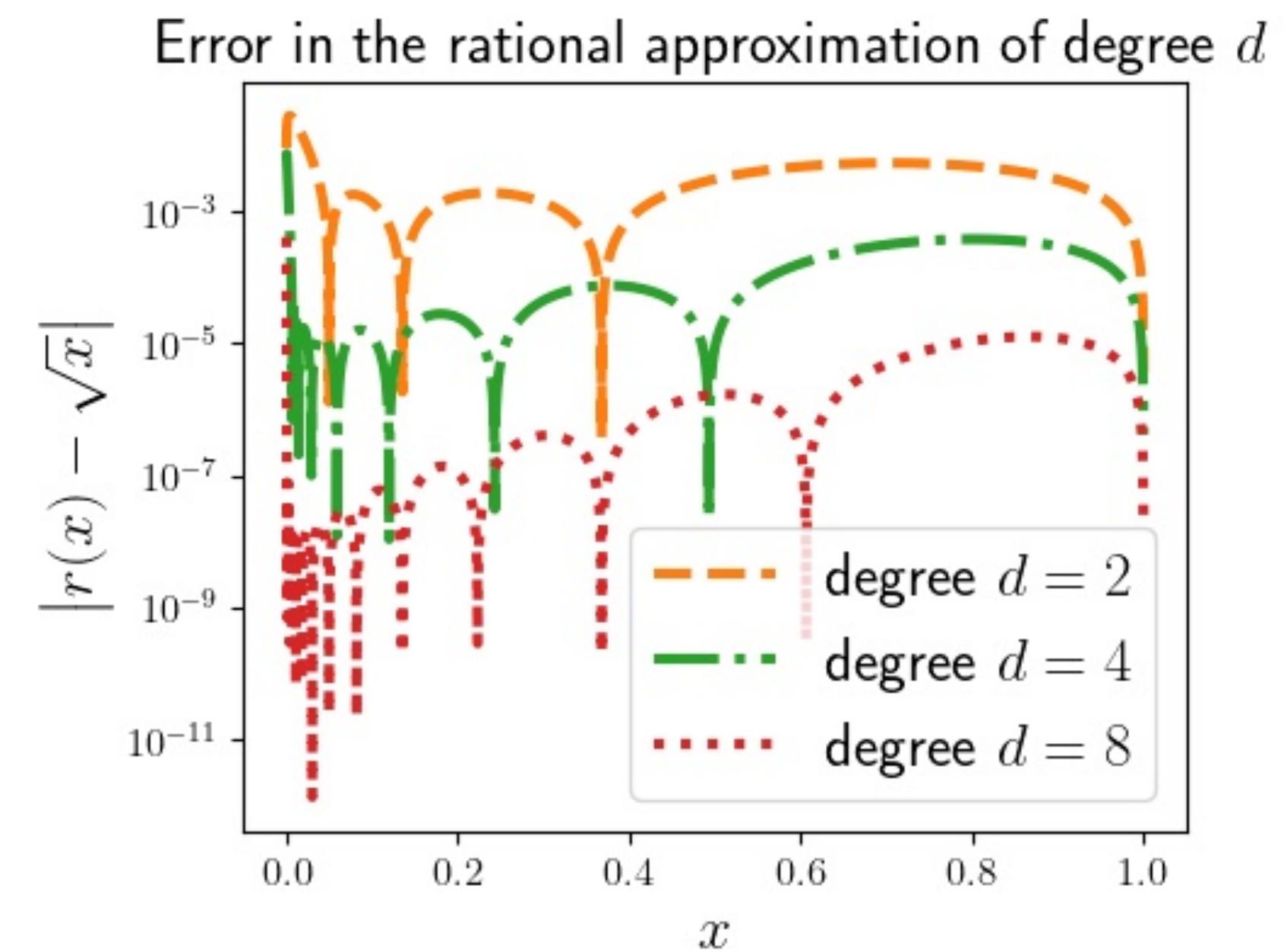
where $\text{err}(r)$ quantifies the **approximation quality**

$$\text{err}(r) = \max_{x \in \mathbb{C} : |x|=1-n^{-1}} |r(x) - \sqrt{1-x}|$$

There exists a degree- d rational function that satisfies the tight approximation bound:

$$\sup_{x \in [0,1]} |r(x) - \sqrt{x}| \leq 3 \cdot \exp(-\sqrt{d}).$$

Newman. **Rational approximation to $|x|$.** Michigan Math. J. (1964)



where $\text{err}(r)$ quantifies the **approximation quality**

$$\text{err}(r) = \max_{x \in \mathbb{C} : |x|=1-n^{-1}} |r(x) - \sqrt{1-x}|$$

Our approach: Buffered Linear Toeplitz (BLT) Mechanism

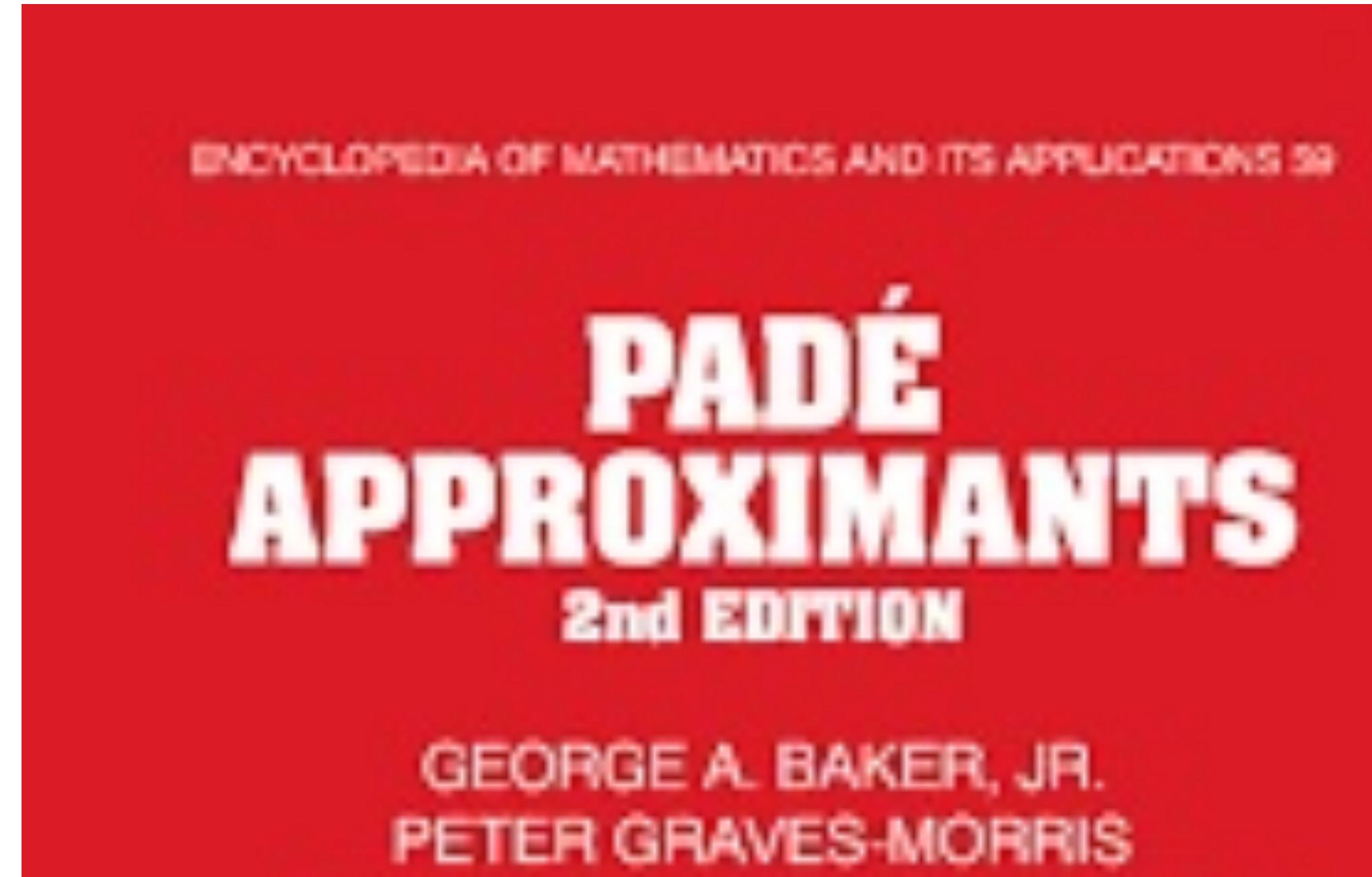
	Max Error	Noise generation time (in iteration t)
Independent noise	$\Theta(\sqrt{n})$	$O(\dim)$
Optimal correlated noise	$\frac{\log n}{\pi} + c$	$O(t \cdot \dim)$
b-Banded	$O\left((\sqrt{n/b} - 1) \log b\right)$	$O(b \cdot \dim)$
BLT of degree d	$\frac{\log n}{\pi} + O(n \cdot \exp(-\sqrt{d}))$	$O(d \cdot \dim)$

Suffices to take $d=O(\log^2 n)$!

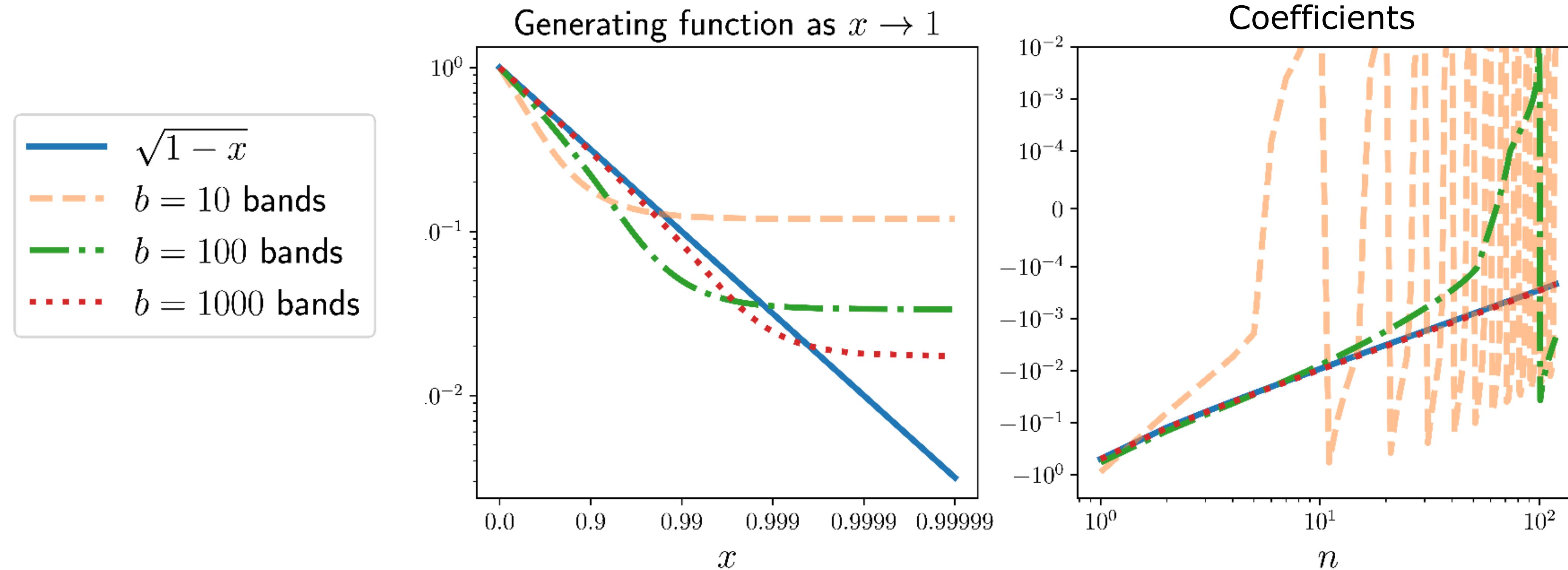
Key difference: approximation quality

Banded: Set $\beta_t = 0$ for $t > b \Rightarrow$ polynomial approximation

BLT:
$$\beta'_t = \sum_{i=1}^d \alpha_i \lambda_i^{t-1} \Rightarrow$$
 rational approximation

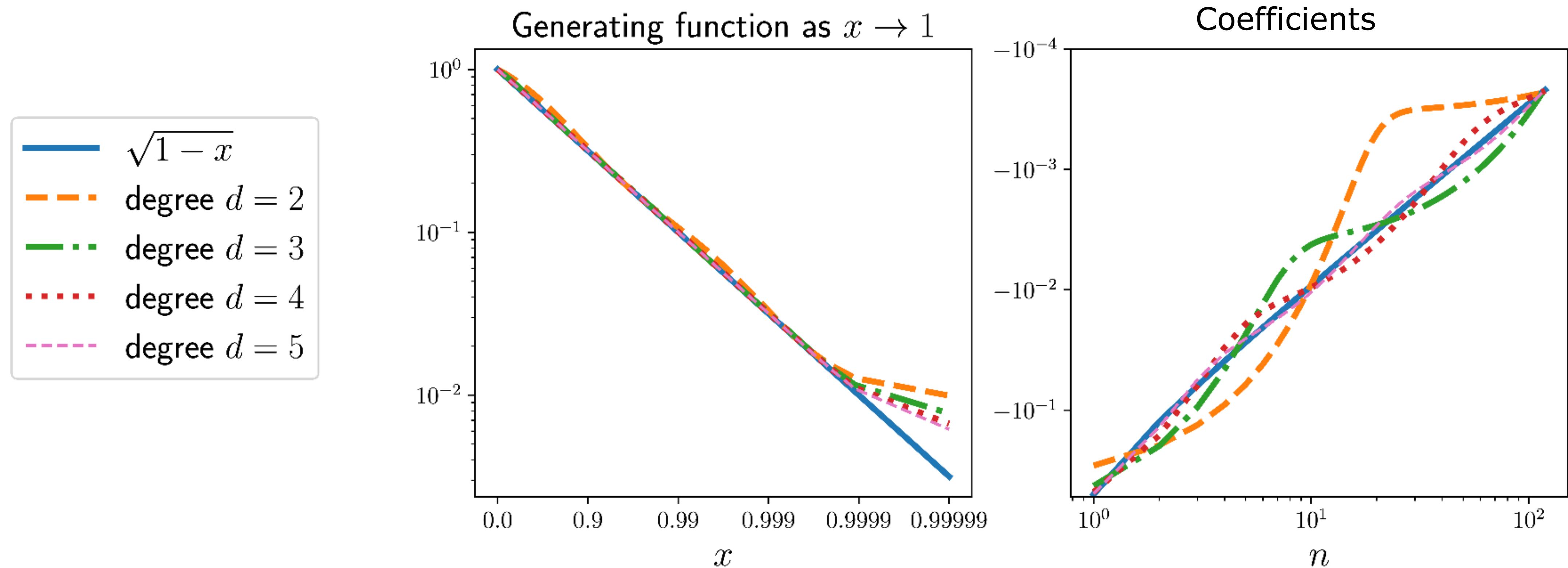


Approximation quality: banded mechanism



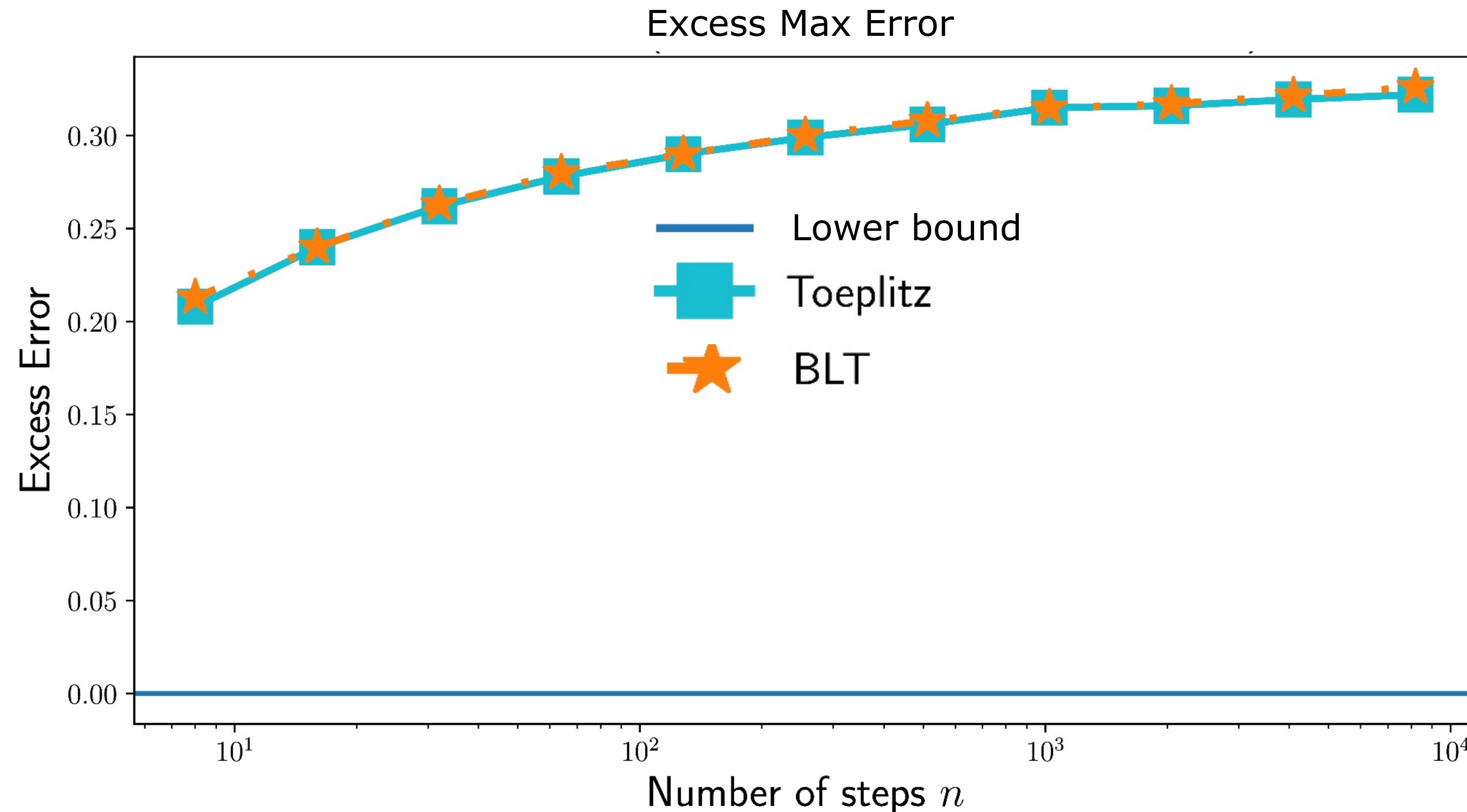
Note: here, we use a polynomial approximation to $1 / (1 - x)^{1/2}$ rather than $(1 - x)^{1/2}$

Approximation quality: BLT mechanism

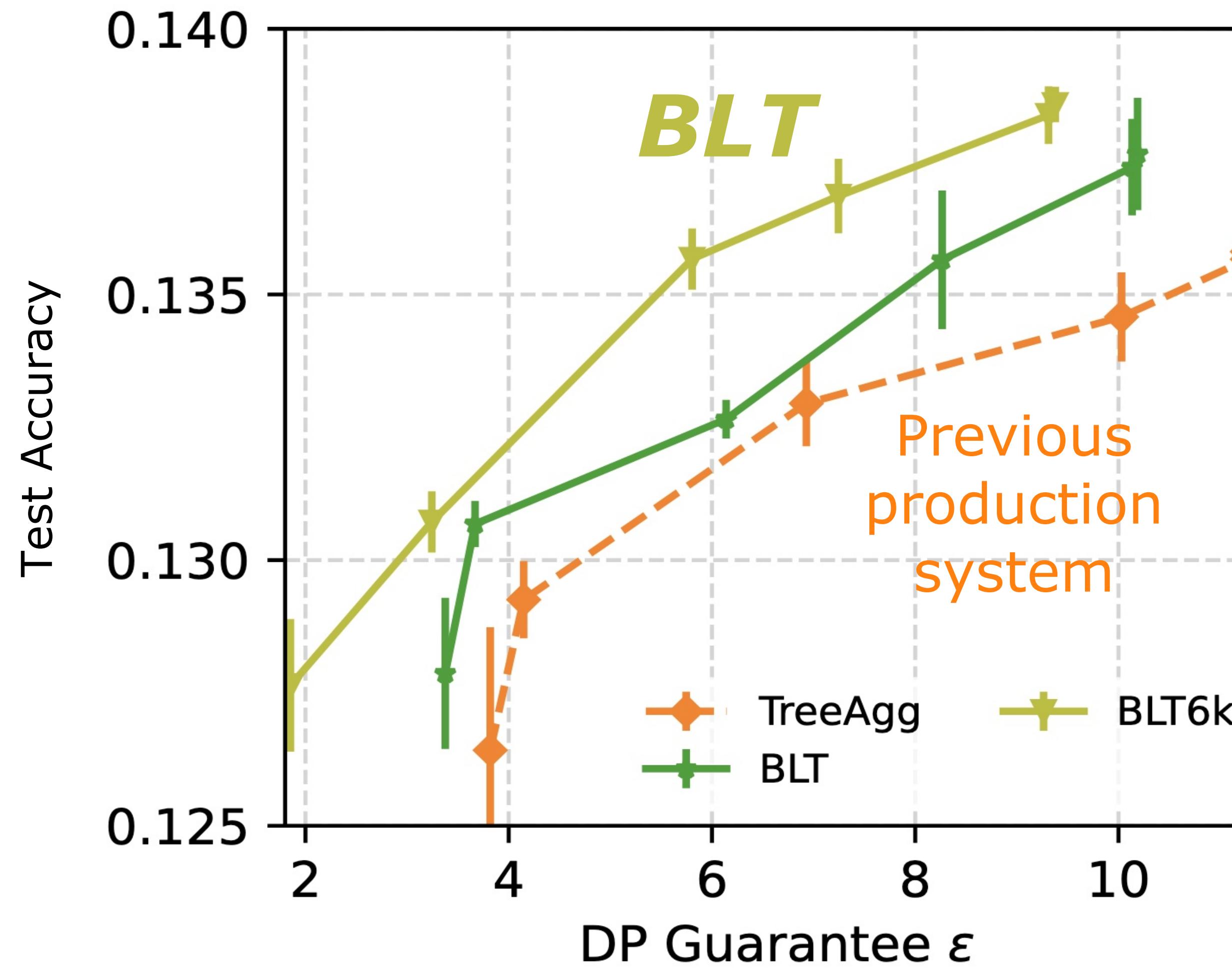


Note: BLT approximation of $1 / (1 - x)^{1/2}$ \Leftrightarrow BLT approximation of $(1 - x)^{1/2}$

Empirical Results



Practical Impact: Google's production language model (Portuguese)



Plot: McMahan, Xu, Zhang (2024)

Aside: Theoretical evidence in favour of correlated noise for *learning problems*

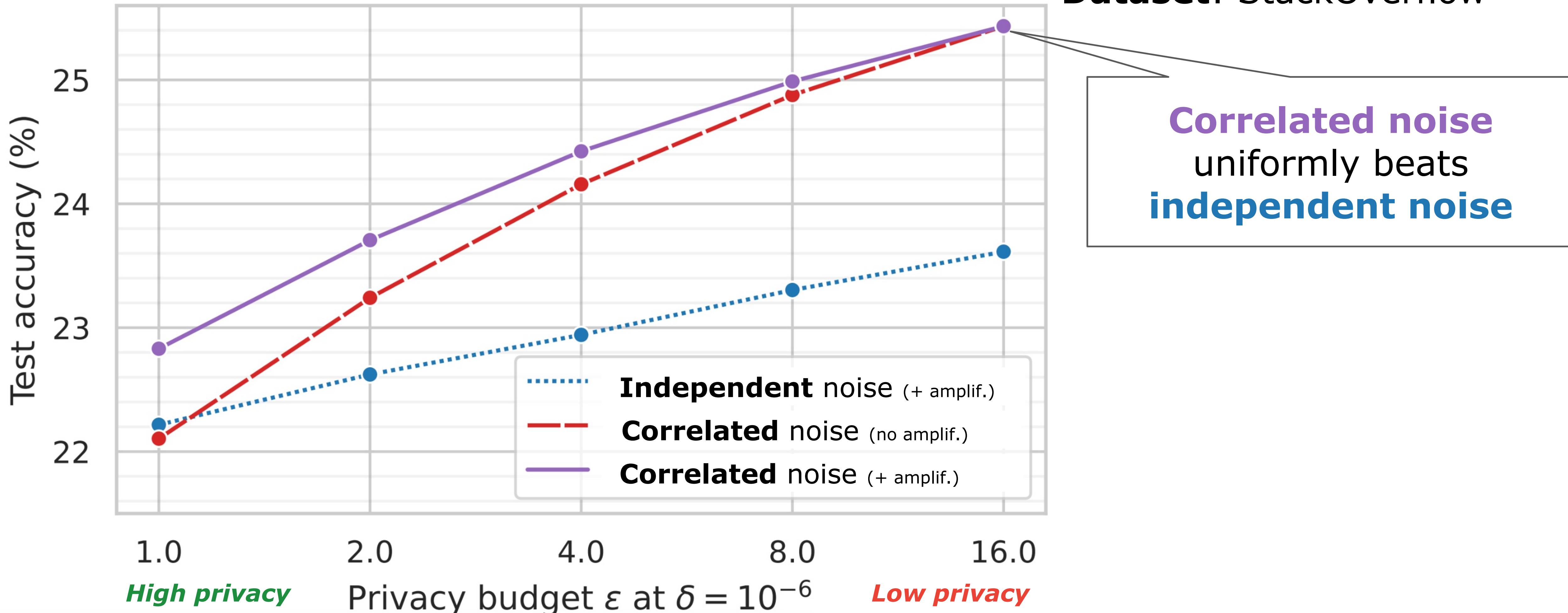
Choquette-Choo*, Dvijotham*, **P.***, Ganesh, Steinke, Thakurta.

Correlated Noise Provably Beats Independent Noise for Differentially Private Learning.
ICLR (2024)

Prior work: (Empirically) correlated noise outperforms independent noise

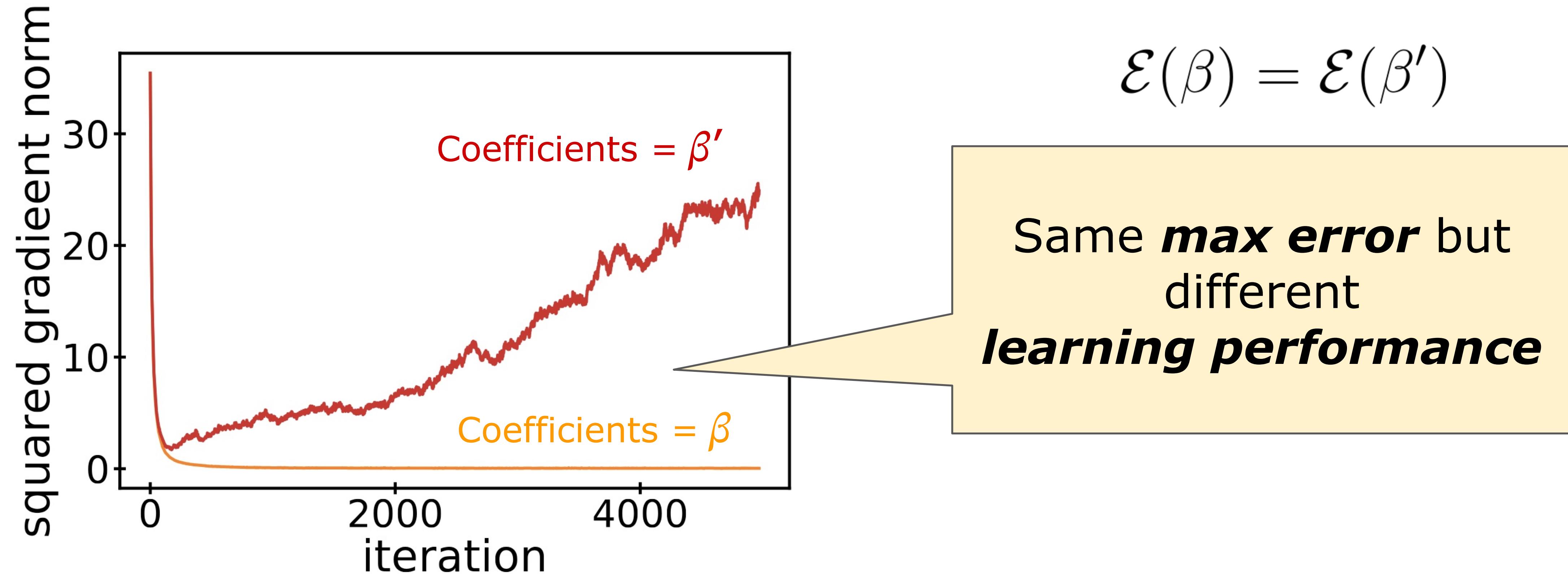
Experiment: user-level DP + language modeling

Dataset: StackOverflow



Is correlated noise provably better for learning problems?

The surrogate objective is not related to the learning objective



Our result: Correlated noise *is* provably better for learning problems

(Anti-) correlated noise provably beats independent noise

For linear regression, **dimension d** improves to problem-dependent **effective dimension d_{eff}**

Independent noise

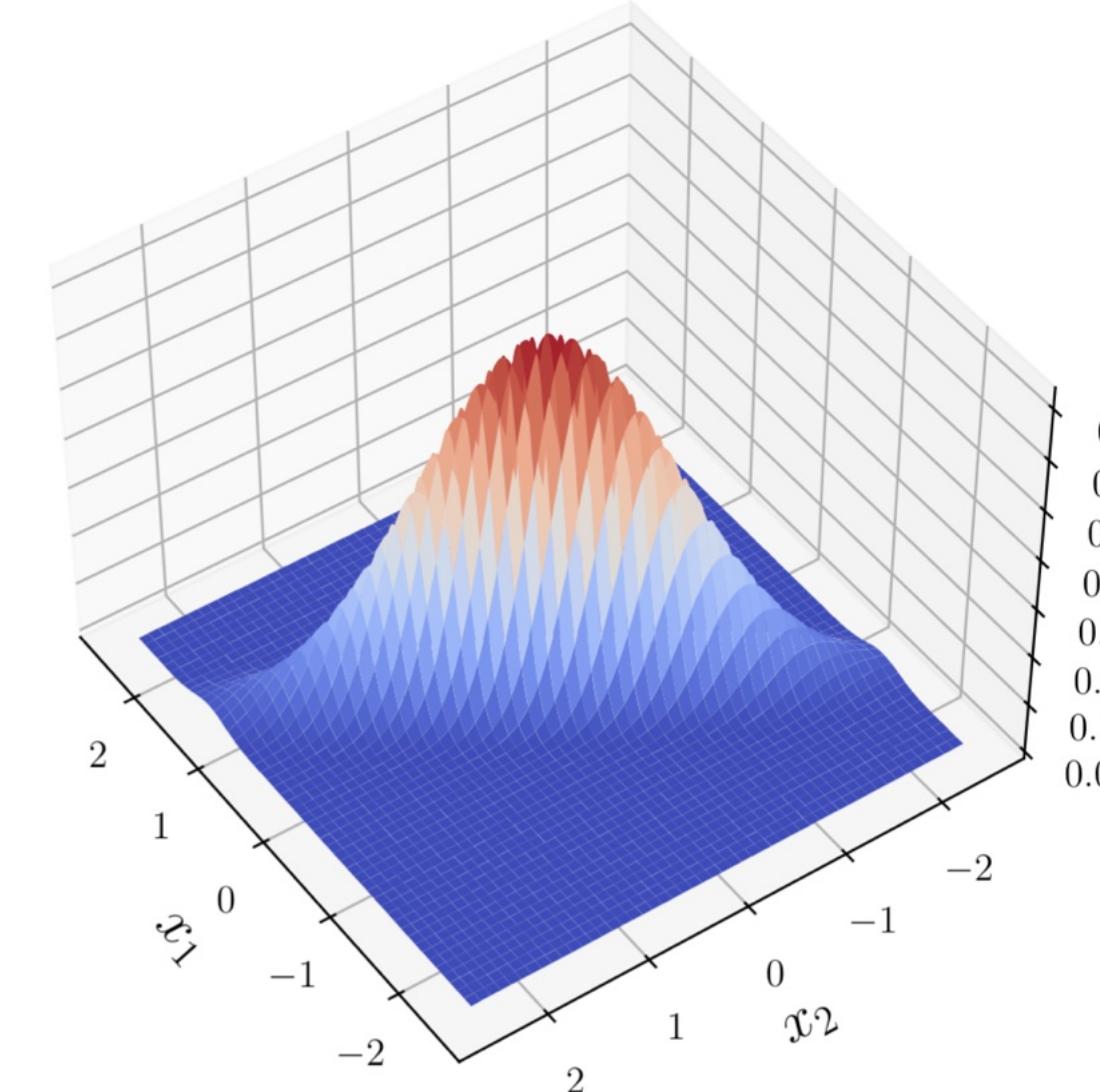
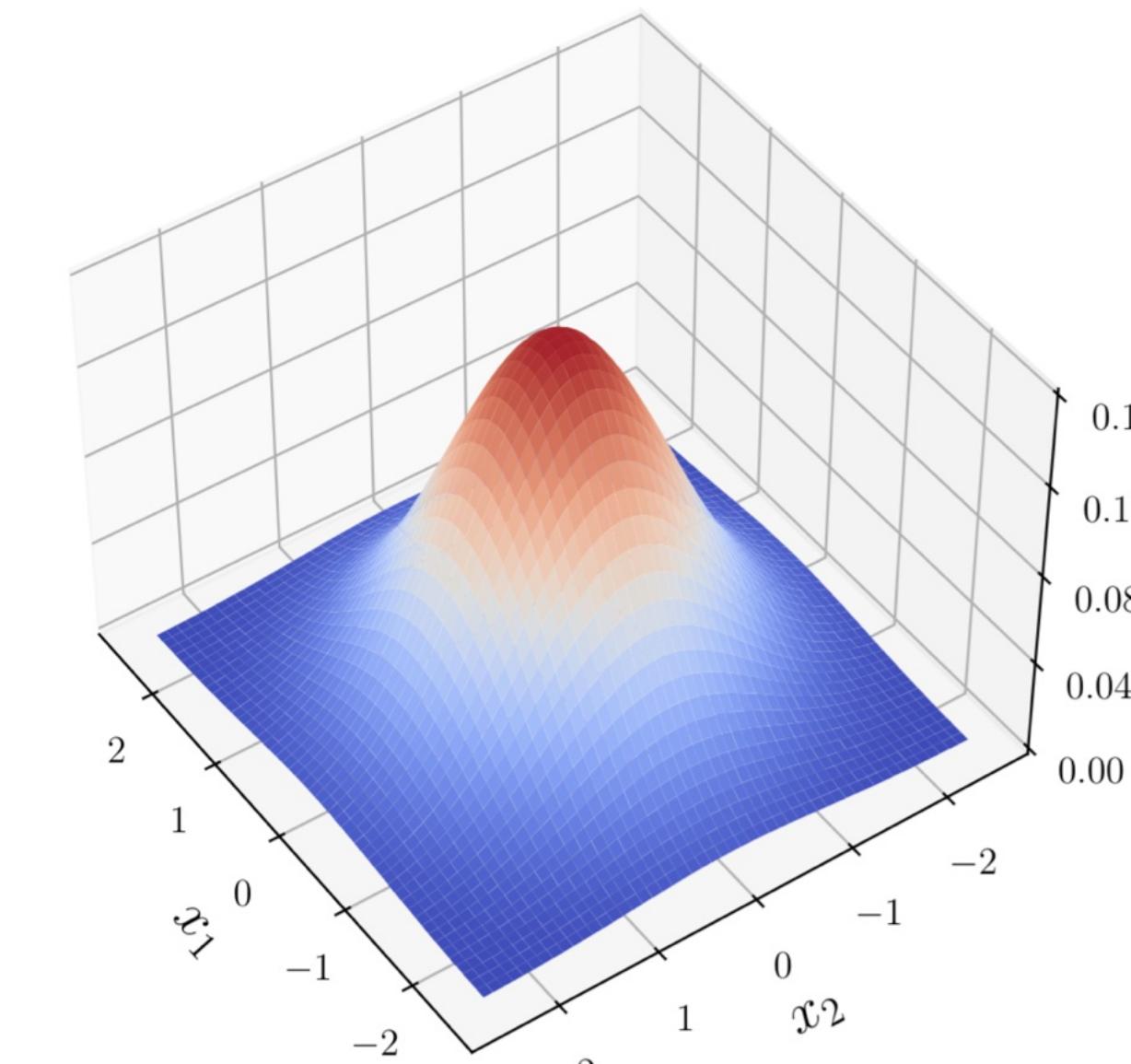
Correlated noise

Lower bound

$\Theta(d)$

$\tilde{O}(d_{\text{eff}})$

$\Omega(d_{\text{eff}})$



High
effective
dimension

Low
effective
dimension

Aside 2: Fine-tuning LLMs with user-level DP

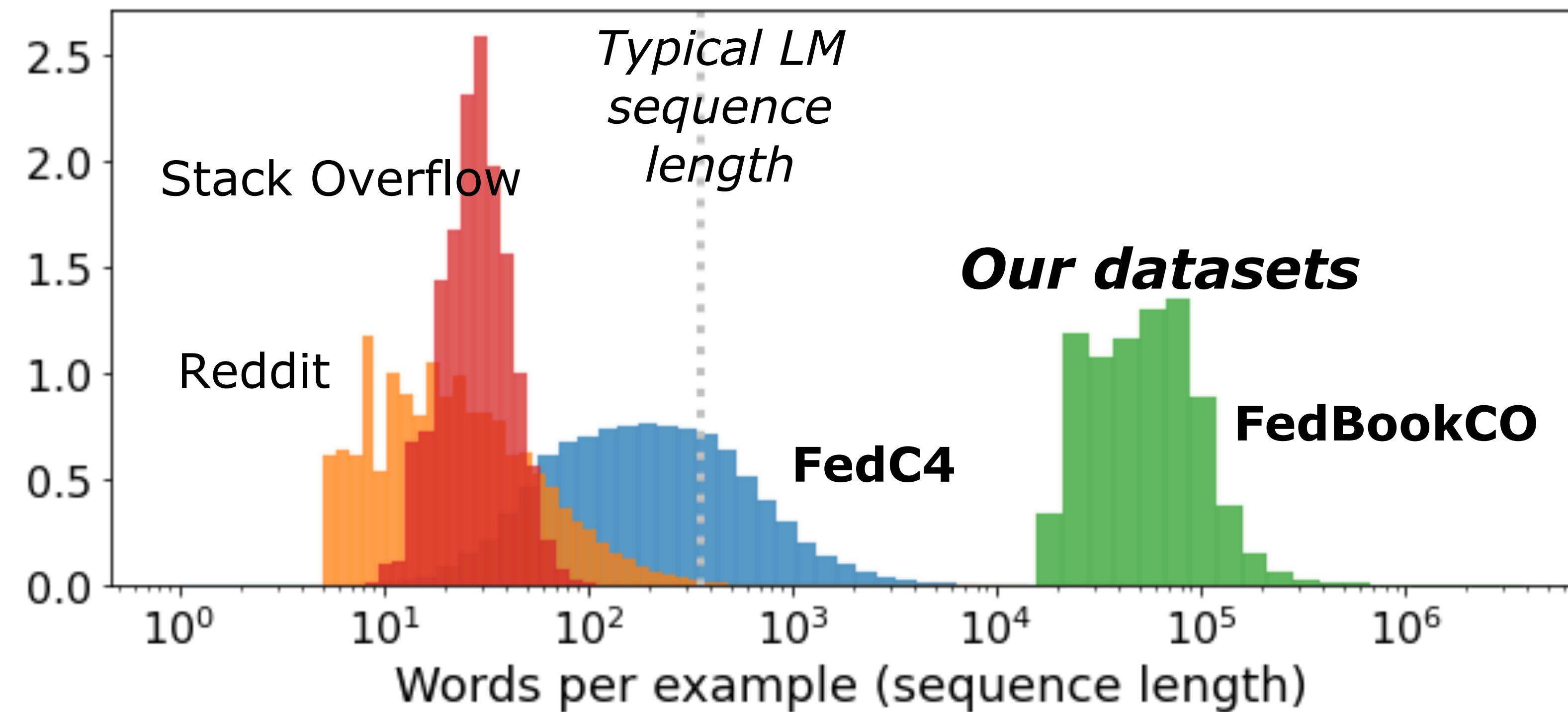
Scaling up user-level DP to LLMs (on a budget)

NeurIPS D&B 2023

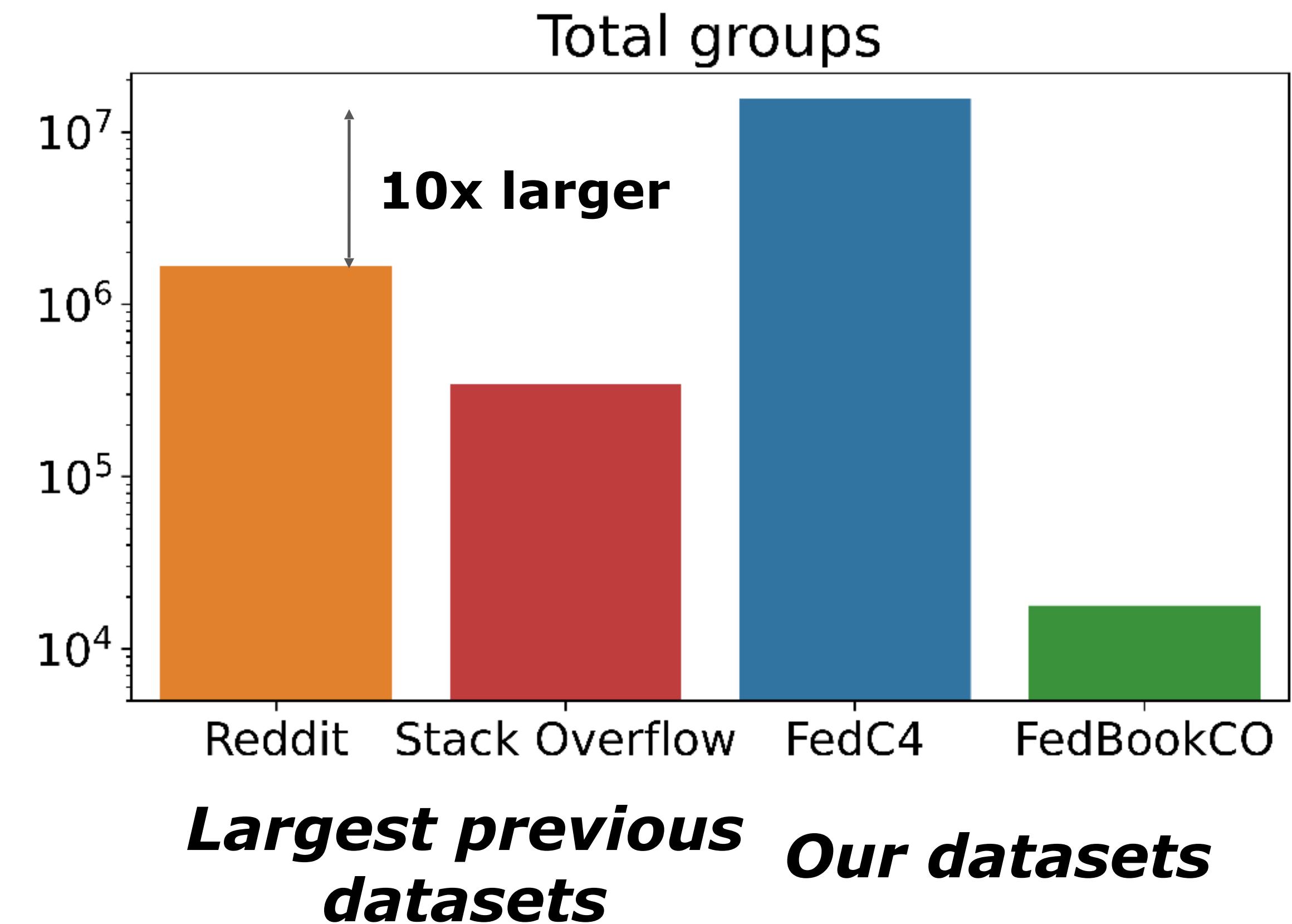
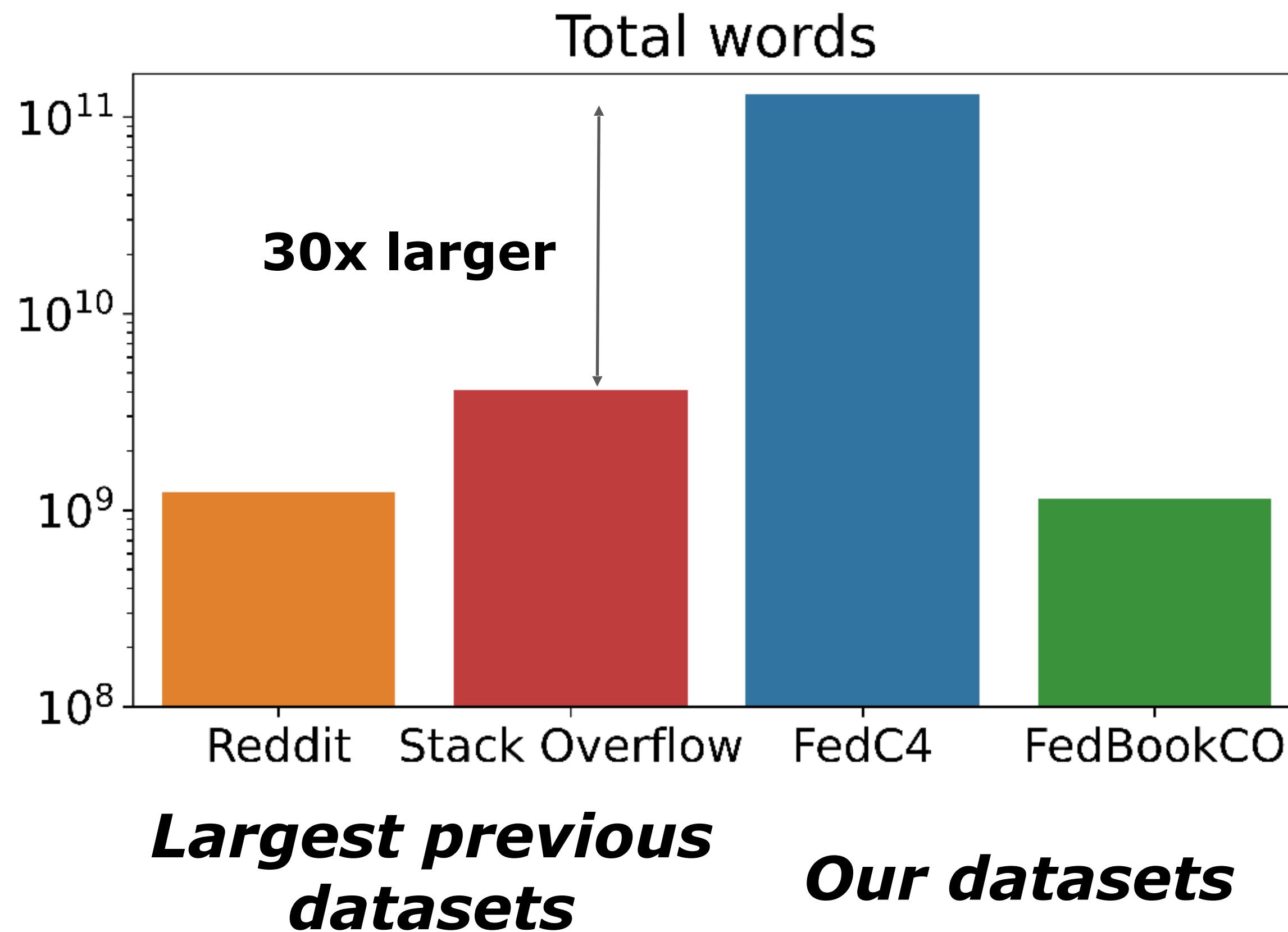
SaTML 2025

- First user-level DP benchmarks for LLMs
- Training with **$O(0.5B)$ params** and **$O(100K)$ users**

***Largest
previous
datasets***



More words & groups than any previous benchmarks



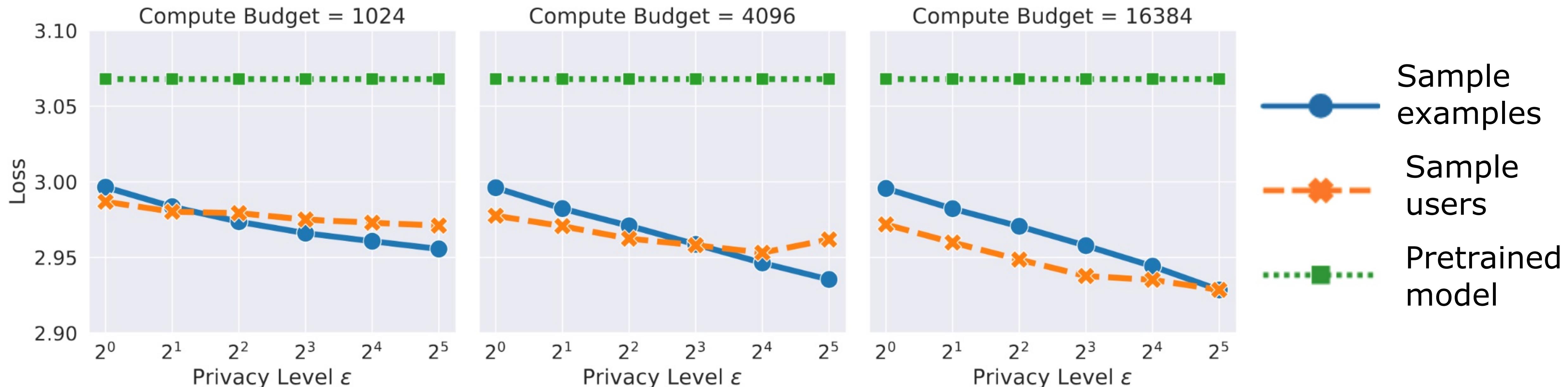
Scaling up user-level DP to LLMs (on a budget)

NeurIPS D&B 2023

SaTML 2025

Scaling user-level DP to LLMs (on a budget)
with independent noise:

- First user-level DP benchmarks for LLMs
- Training with **$O(0.5B)$ params** and **$O(100K)$ users**



Coming soon: Monograph and tutorial on correlated noise mechanisms!

Contents

1	Introduction and Background	4
1.1	Introduction to Differential Privacy	5
1.2	Problem Statement: DP Estimation of Weighted Prefix Sums	9
1.3	Correlated Noise Mechanisms	13
1.4	Why Correlated Noise Mechanisms?	16
1.5	Design Space and Detailed Outline of the Monograph	23
1.6	Some Technical Considerations*	27
1.7	Chapter Notes	29
1.8	Bibliographic Notes	31
2	Correlated Noise Mechanisms for Streaming Prefix Sums	34
2.1	Design Considerations	35
2.2	Dense Mechanism	39
2.3	Toeplitz Mechanism	40
2.4	Banded Toeplitz Mechanism	44
2.5	Buffered Linear Toeplitz (BLT) Mechanism	49
2.6	Tree Aggregation*	54
2.7	Empirical Comparison of the Mechanisms	58
2.8	Other Error Metrics*	60
2.9	An Approximation Theory Viewpoint*	63
2.10	Bibliographic Notes	69
3	Correlated Noise Mechanisms for Machine Learning	74
3.1	Motivation	75
3.2	Learning Problems as Weighted Prefix Sums	78
3.3	Multi-Epoch Correlated Noise Mechanisms	80
3.4	Simulations	99
3.5	Learning Guarantees for Correlated Noise Mechanisms*	99
3.6	Proofs of Multi-Epoch Sensitivity*	106
3.7	Privacy Amplification by Sampling*	107
3.8	Bibliographic Notes	109
4	Implementation Details and Practical Recommendations	112
4.1	Numerical Mechanism Optimization	113
4.2	Optimizing the Dense Mechanism	115
4.3	Optimizing Parameterized Mechanisms	123
4.4	Open-Source Software	130
4.5	Choosing a Correlated Noise Mechanism	130
4.6	Bibliographic Notes	133
5	Challenges and Open Questions	135
5.1	Directions Forward for Practice	135
5.2	Directions Forward for Theory	139
	References	140
	Appendices	148
	Common Notions of Differential Privacy	149
A.1	Zero-Concentrated DP (zCDP)	149
A.2	Approximate DP	150
	Review of Linear Algebra	151
A.3	Induced Matrix norms	152
A.4	Matrix Decompositions	153
A.5	Toeplitz Matrices	154

Open Problem: Continuous time limits

$$\theta_{t+1} = \theta_t - \eta \left(g_t + \boxed{z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau}} \right)$$

Proceedings of Machine Learning Research vol 195:1–44, 2023

36th Annual Conference on Learning Theory

Universality of Langevin Diffusion for Private Optimization, with Applications to Sampling from Rashomon Sets

Arun Ganesh
Google Research

Abhradeep Thakurta
Google DeepMind

Jalaj Upadhyay
Rutgers University

ARUNGANESH@GOOGLE.COM

ATHAKURTA@GOOGLE.COM

JALAJ.UPADHYAY@RUTGERS.EDU

Precise analysis
(better rates)

Algorithm design

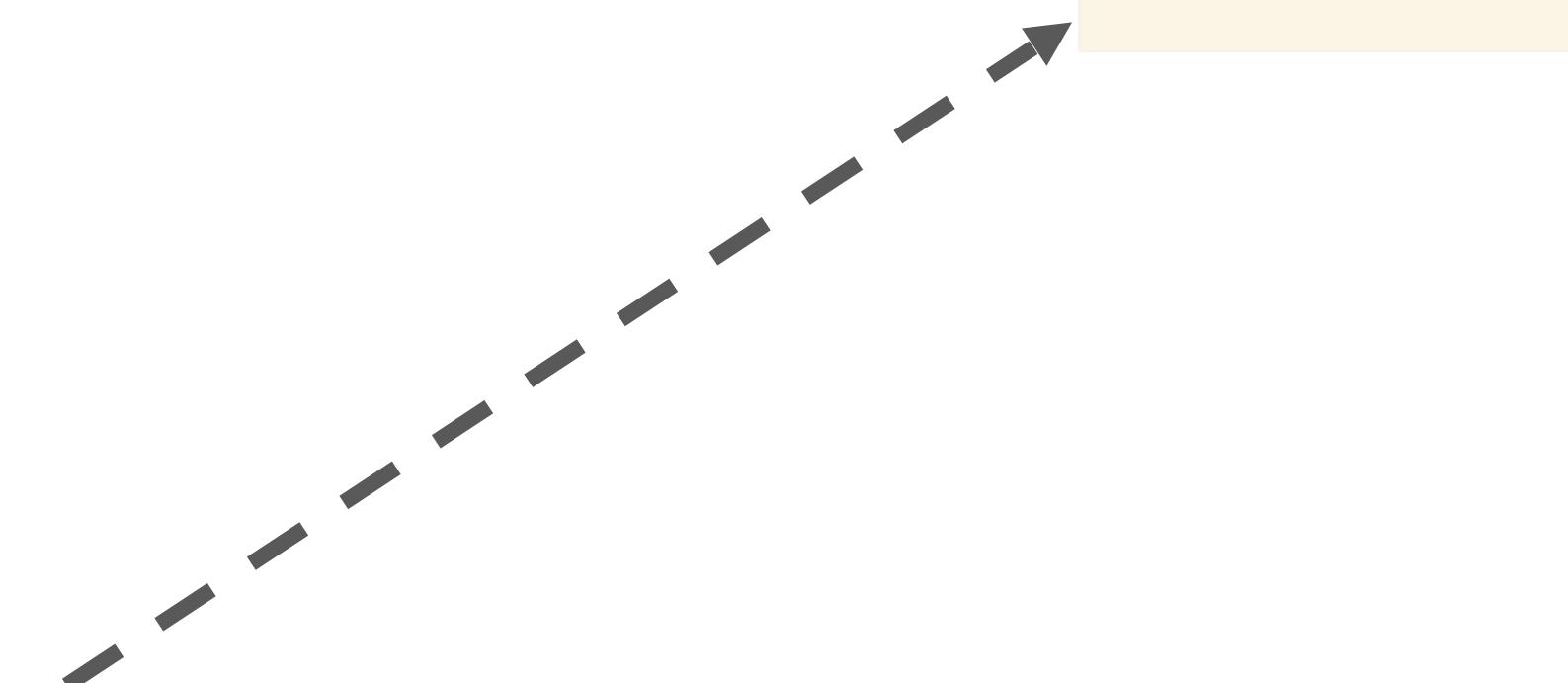
Open Problem: Adaptive Gradient Algorithms

SGD update (without noise)

$$\theta_t - \theta_0 = - \sum_{\tau=0}^{t-1} g_\tau$$

Adam update (without noise)

$$\begin{aligned} v_t &= (1 - \beta_1)v_{t-1} + \beta_1 g_t \\ s_t &= (1 - \beta_2)s_{t-1} + \beta_2 g_t^2 \\ \theta_{t+1} &= \theta_t - \eta \frac{v_t}{\sqrt{s_t} + \delta} \end{aligned}$$



Non-linear functions of the injected noise

Part 2: How audit user-level DP?

Unleashing the power of randomness in auditing DP

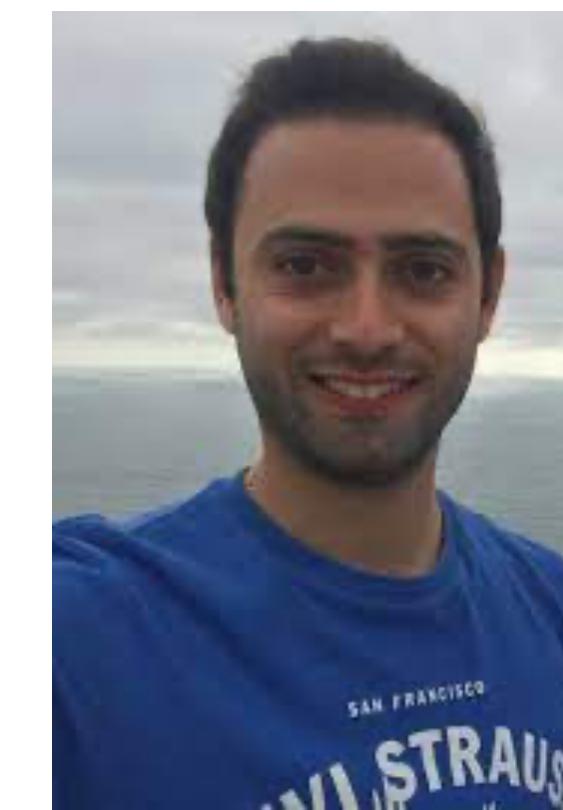
NeurIPS 2023



Krishna Pillutla



Galen Andrew



Peter Kairouz



Brendan McMahan

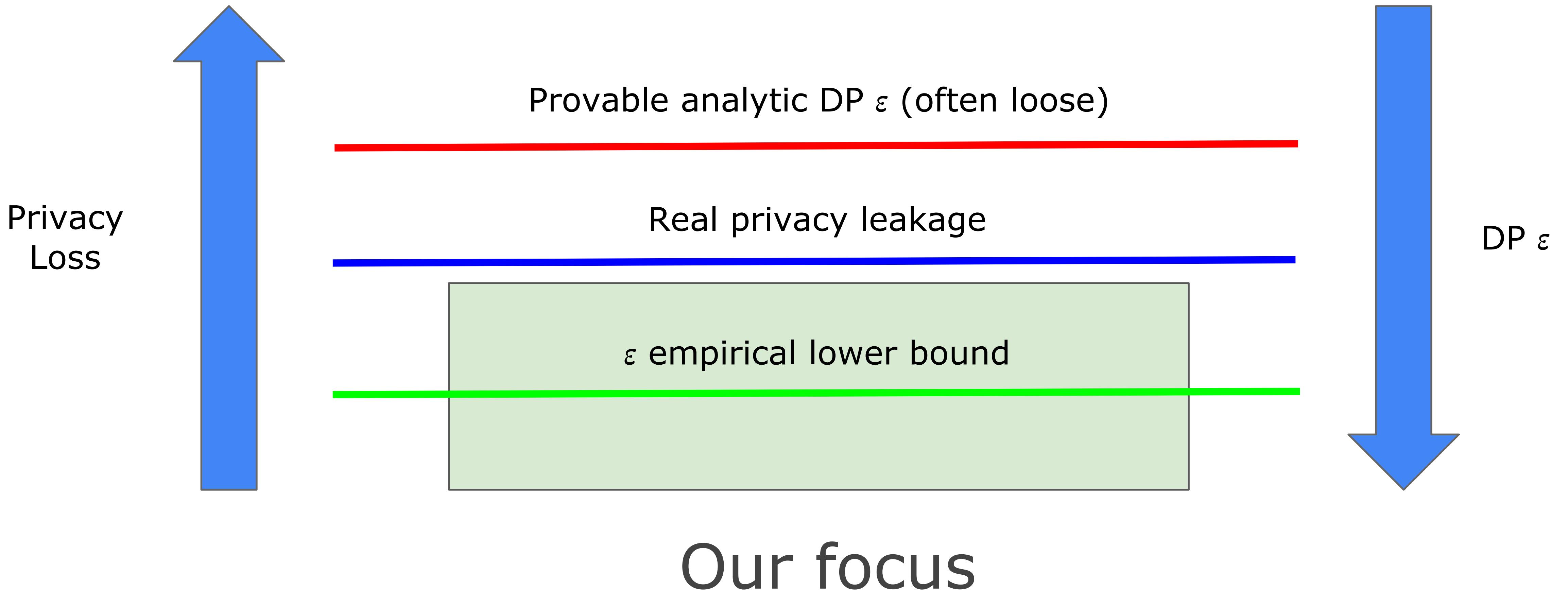


Alina Oprea



Sewoong Oh

Empirical privacy auditing



Why empirical privacy auditing?

To verify that we actually provide the guarantee we claim
(no bugs in proofs/implementation)

The image shows a code editor interface with two tabs open, displaying the contents of `mnist_experiment.py` and `upstream_clipping.py`. The code is presented in a diff format, where changes from the upstream version to the local version are highlighted.

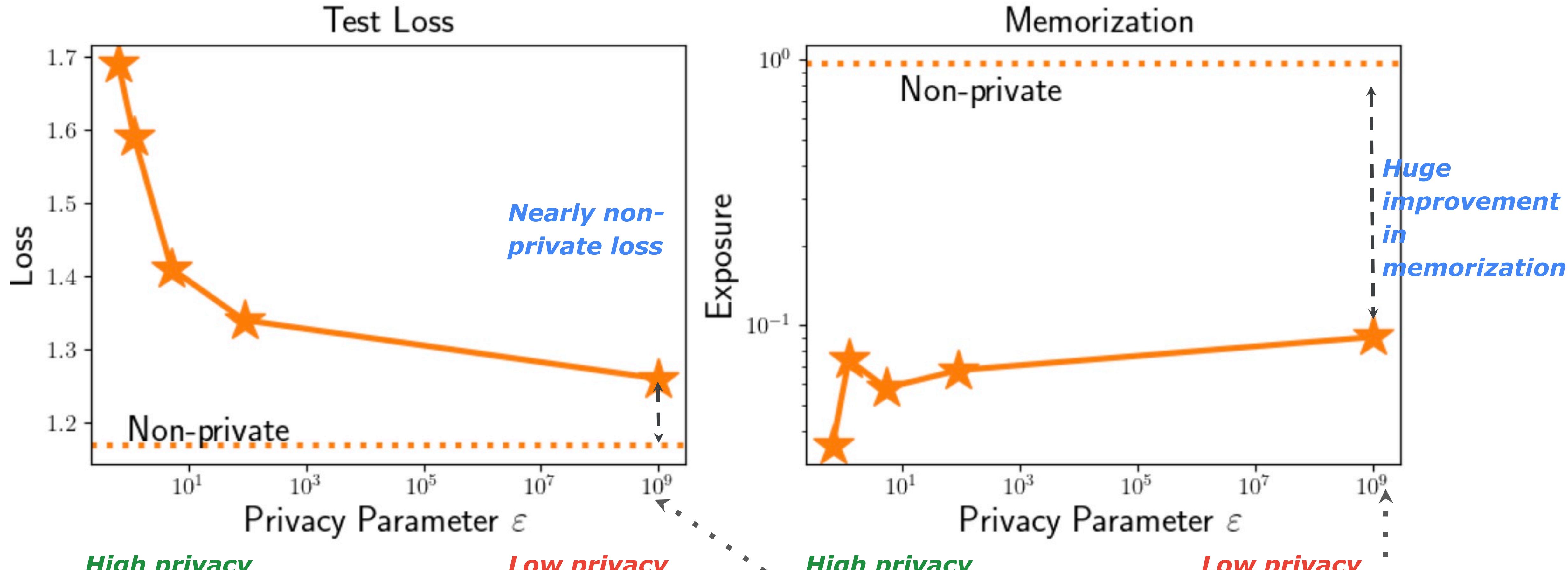
mnist_experiment.py:

```
@@ -71,7 +71,7 @@ def forward(self, x):
    71      71
    72      72
    73      73
  74      - rho_i,
    74      + epochs,
    75      75         inp_clip,
    76      76         grad_clip
    77      77         grad_clip/BATCH_SIZE
)
tl, correct, set_len = uc.test(model, test_loader)
print(f'MNIST_{BATCH_SIZE}_{epochs}_{grad_clip}_{inp_clip}_{rho_i}', correct/set_len)
```

upstream_clipping.py:

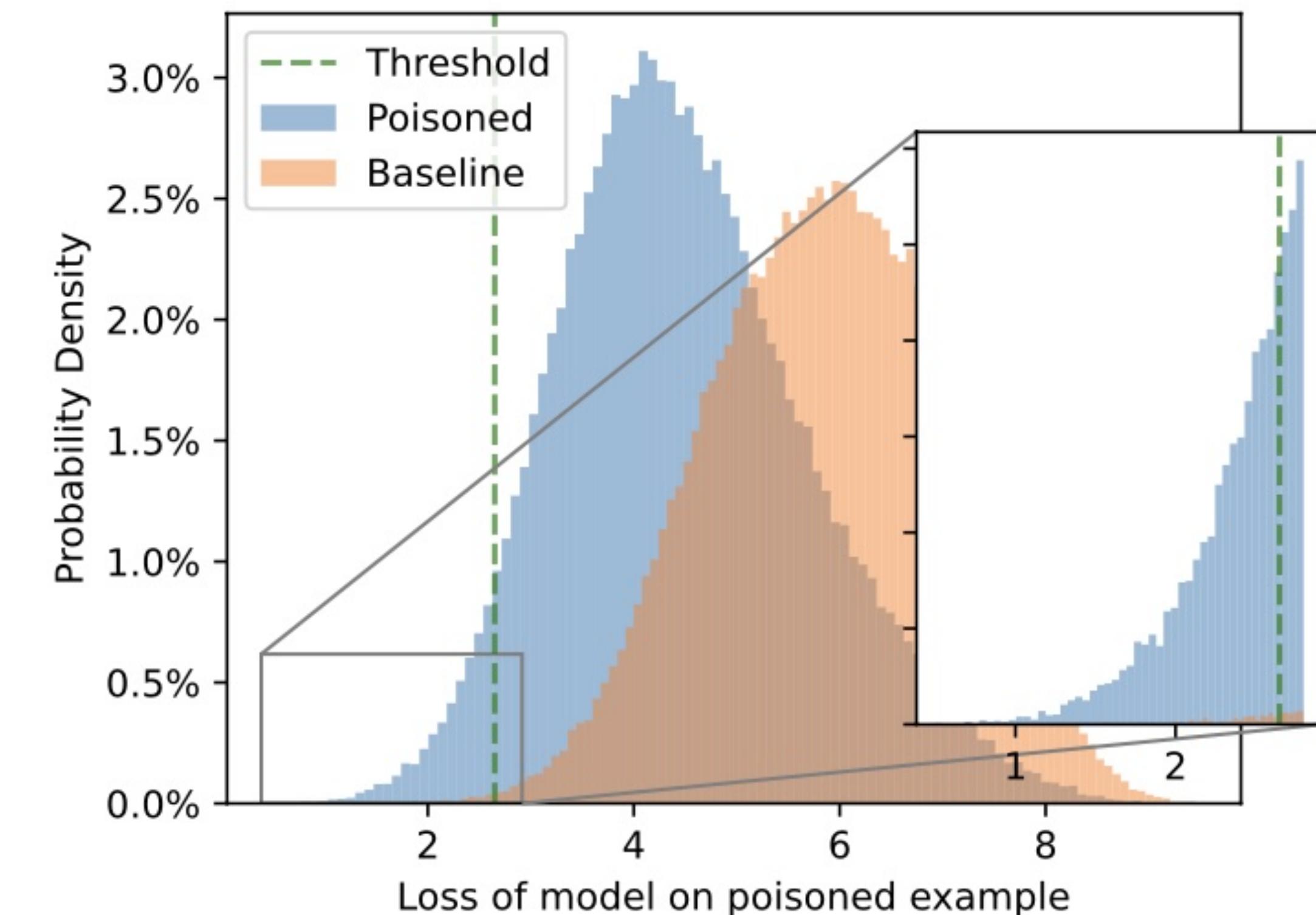
```
@@ -110,7 +110,7 @@ def run_experiment(model, train_loader, rho_i, epochs, input_bound, grad_bound):
 110      110
 111      111          model.train()
 112      112          # sensitivity for everything with weights is just:
 113      - sensitvity = input_bound * grad_bound / train_loader.batch_size
 113      + sensitvity = input_bound * grad_bound
 114      114          sigma = np.sqrt(sensitivity**2 / (2*rho_i))
 115      115          print('sensitivity:', sensitivity)
 116      116
```

Gap between DP guarantees and empirical behavior: Memorization



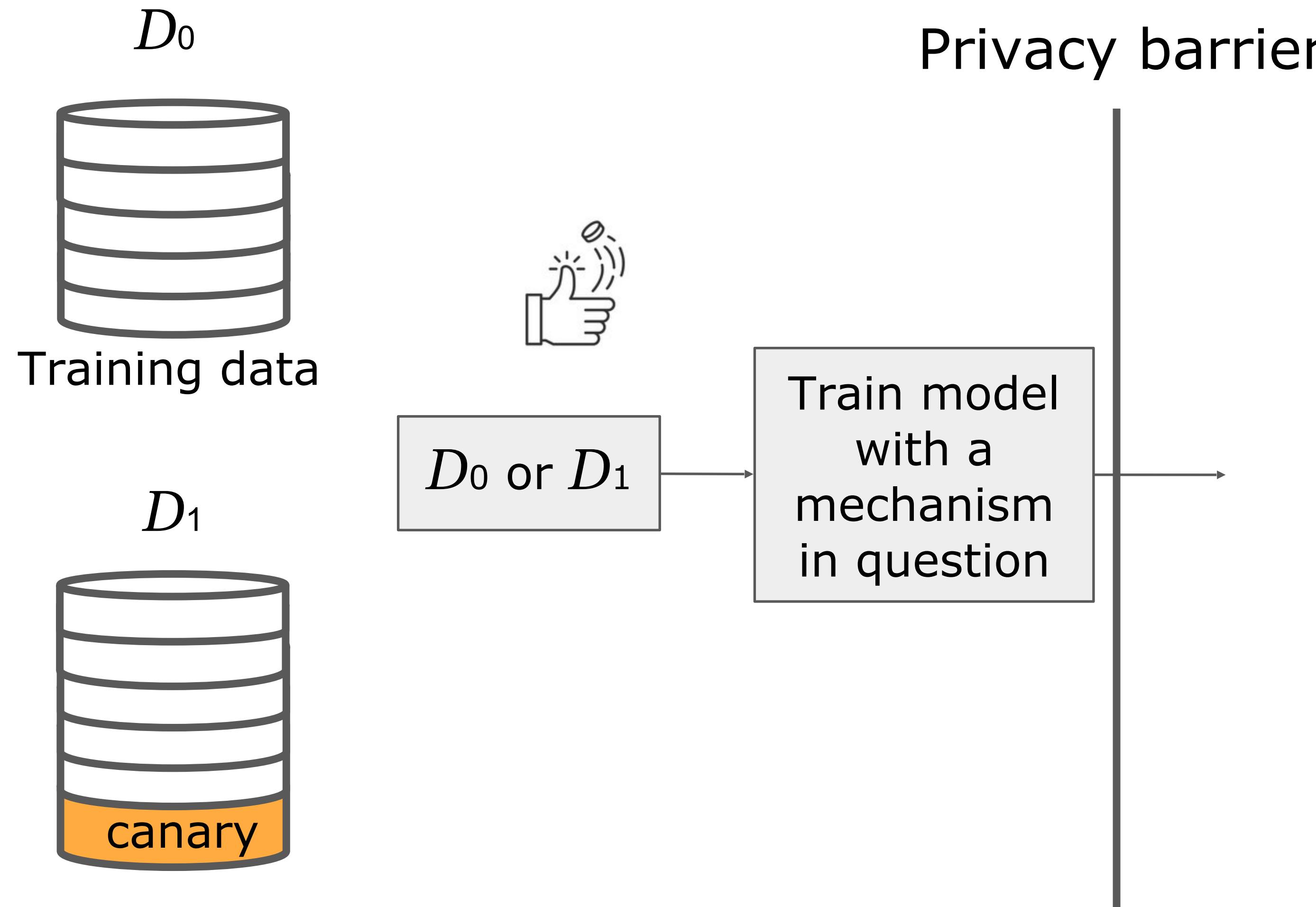
Empirical Privacy Auditing requires **many samples**

- Trained w/ $(0.21, 10^{-5})$ -DP
but empirically $\epsilon > 2.79$ with
confidence $1 - 10^{-8} \Rightarrow \text{bug in}$
implementation
- This required training
 $n=200,000$ models

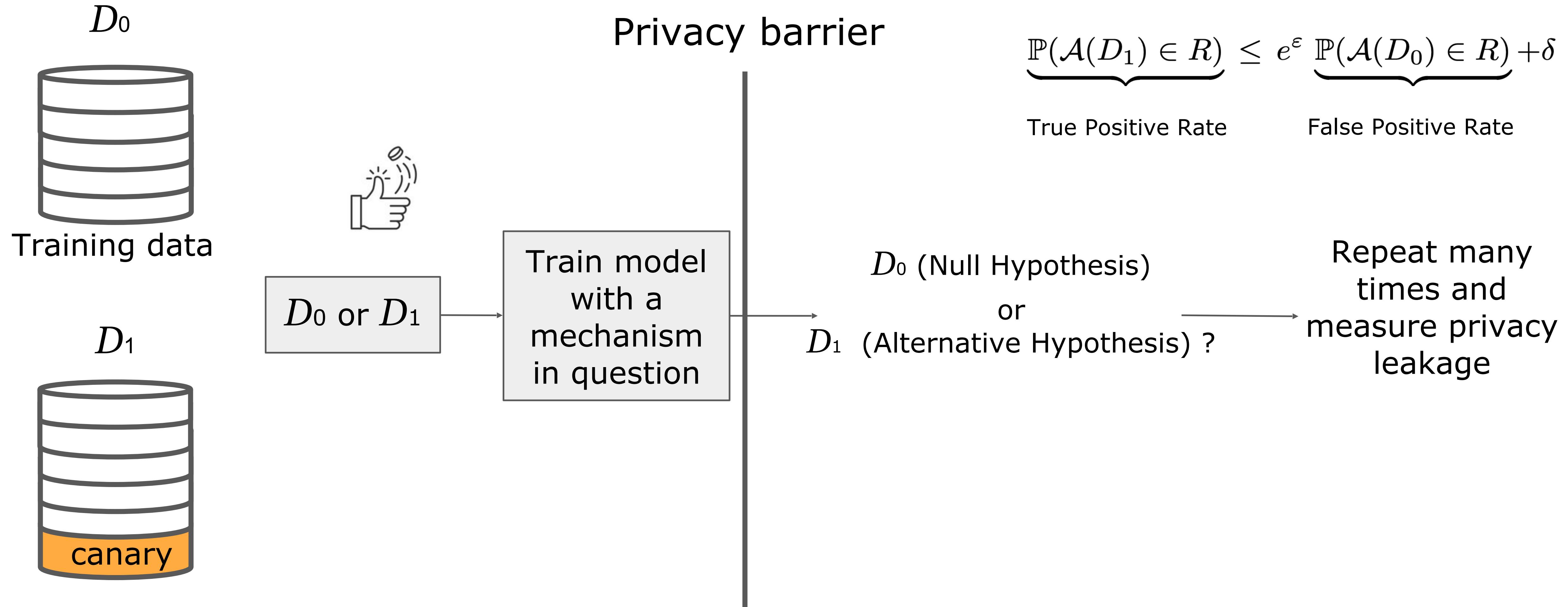


Our goal: make empirical privacy auditing
more *sample-efficient*

Standard approaches for auditing privacy: **binary hypothesis testing**



Standard approaches for auditing privacy: **binary hypothesis testing**



Bottleneck: Bernoulli confidence intervals

- Confidence intervals based on n trials

$$\text{TPR} \approx \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{Guess } i \text{ correct})}_{\text{Empirical TPR/FPR}} + \sqrt{\frac{\text{Variance}}{n}}$$

Actual
TPR/FPR

Empirical TPR/
FPR

Sample size n needs to be large
for good estimates

Actual TPR/
FPR

$$\begin{aligned}\varepsilon &\geq \log\left(\frac{\text{TPR} - \delta}{\text{FPR}}\right) \\ &\geq \log\left(\frac{\hat{\text{TPR}}_n - \frac{c}{\sqrt{n}} - \delta}{\hat{\text{FPR}}_n + \frac{c}{\sqrt{n}}}\right)\end{aligned}$$

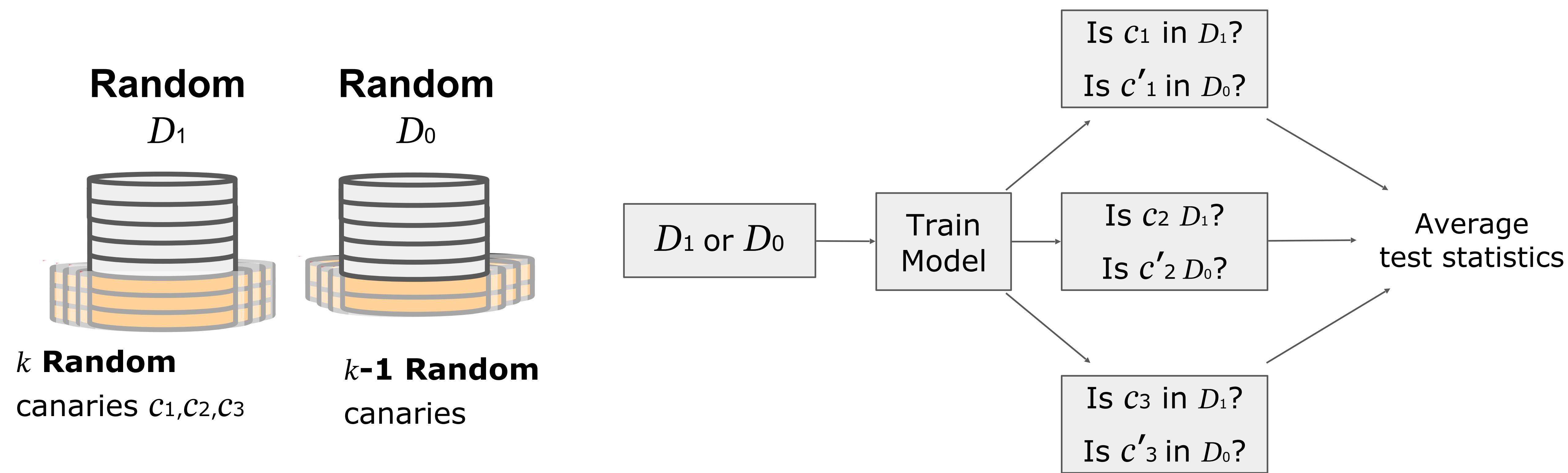
Empirical TPR/
FPR

Our approach: leverage randomness

- **Lifted DP:** Equivalent notion of DP with randomized datasets
- Multiple randomized hypothesis tests
- **Adaptive confidence intervals** capitalizing on low correlations

Multiple hypothesis tests for auditing Lifted DP

- **Leave-One Out** construction with **i.i.d. random canaries**



Multiple hypothesis tests for auditing Lifted DP

If the statistics are independent \Rightarrow better confidence intervals

Unfortunately, they are **dependent**
(but highly uncorrelated)

κ

Random

canaries c_1, c_2, c_3

$\kappa - 1$ Random

canaries

Is c_3 in D_1 ?

Is c'_3 in D_0 ?



Novel higher-order confidence interval

- 2nd-order confidence interval using empirical correlations between two tests

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

- Ideally, when **correlation**= $O(1/k)$, the confidence interval improves as

$$|TPR - \widehat{TPR}_{n,k}| \lesssim \sqrt{\frac{1}{nk}} + \frac{1}{n^{3/4}}$$

Takeaway: Reduces variance from randomness in trials

Standard approach:

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

c - Universal constant

c' - Data-dependent constant

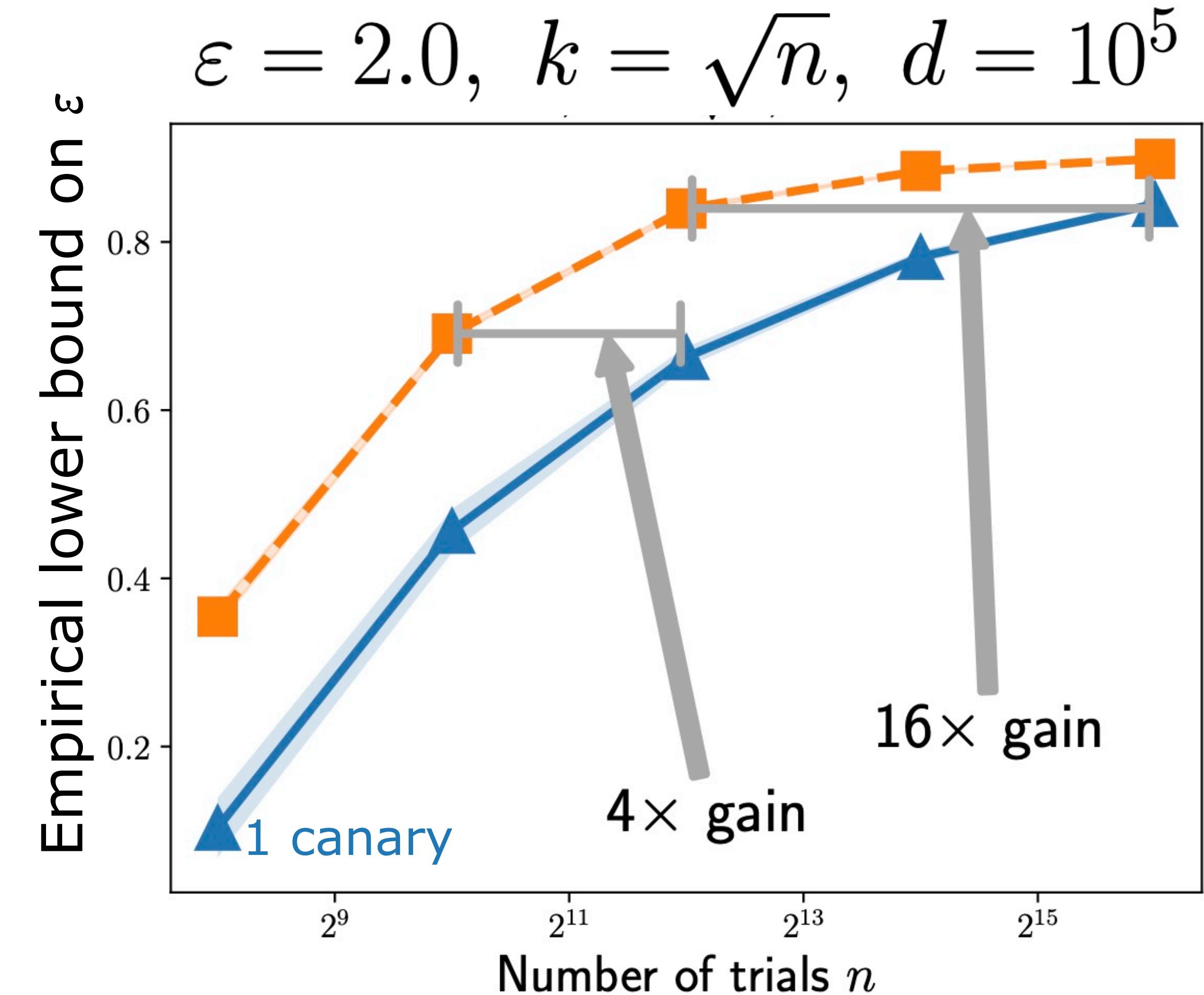
**Lower variance =>
Tighter confidence intervals**

Our approach:

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

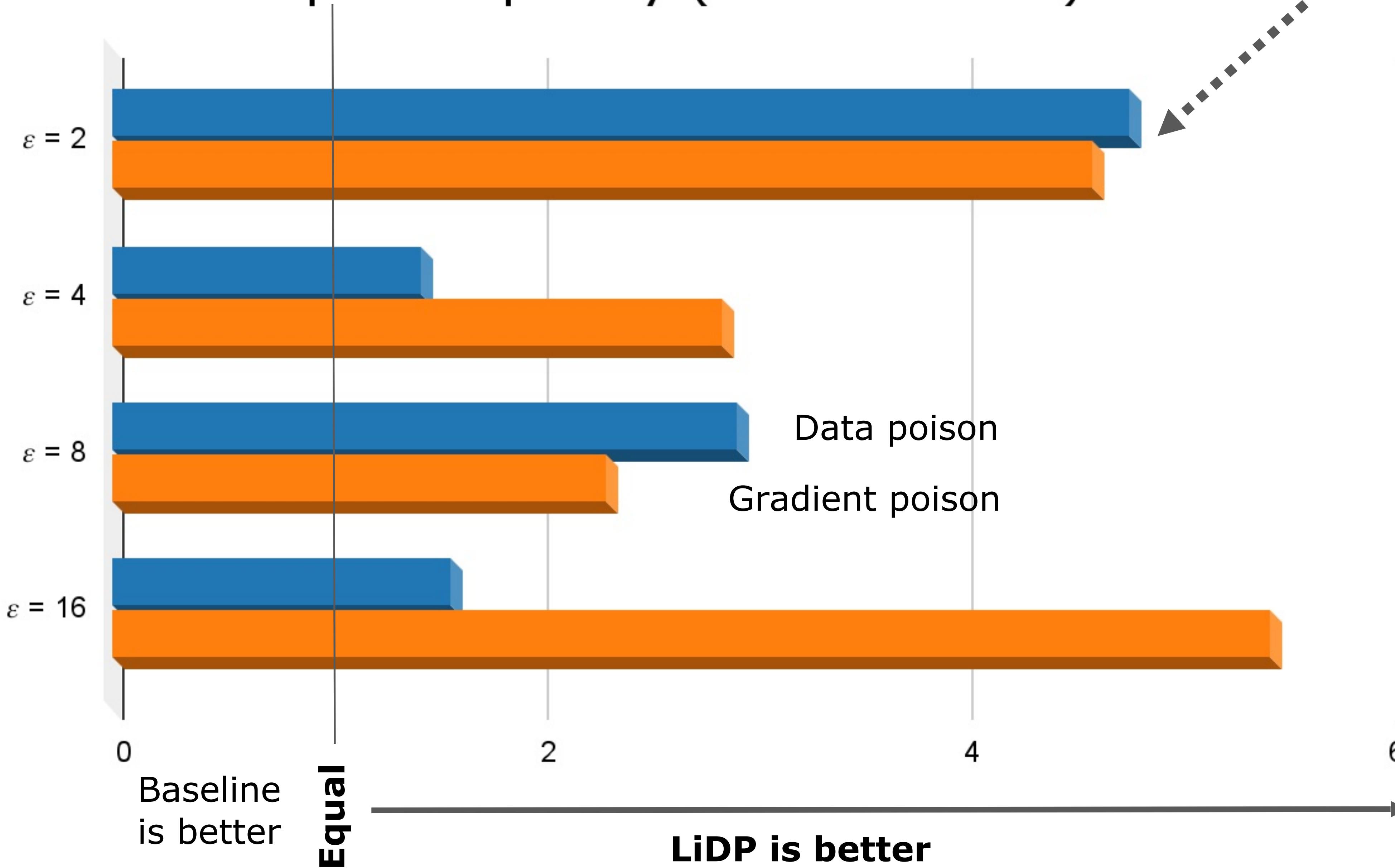
Proof of concept with Gaussian mechanisms

- Sum query with sensitivity 1
- Gaussian mechanism
- k **canaries** uniformly random on the sphere
- **Test statistic** is inner product



Gain in sample complexity (FashionMNIST)

Suffices to train **200 models**
instead of 1000 models



Privacy Auditing with One (1) Training Run

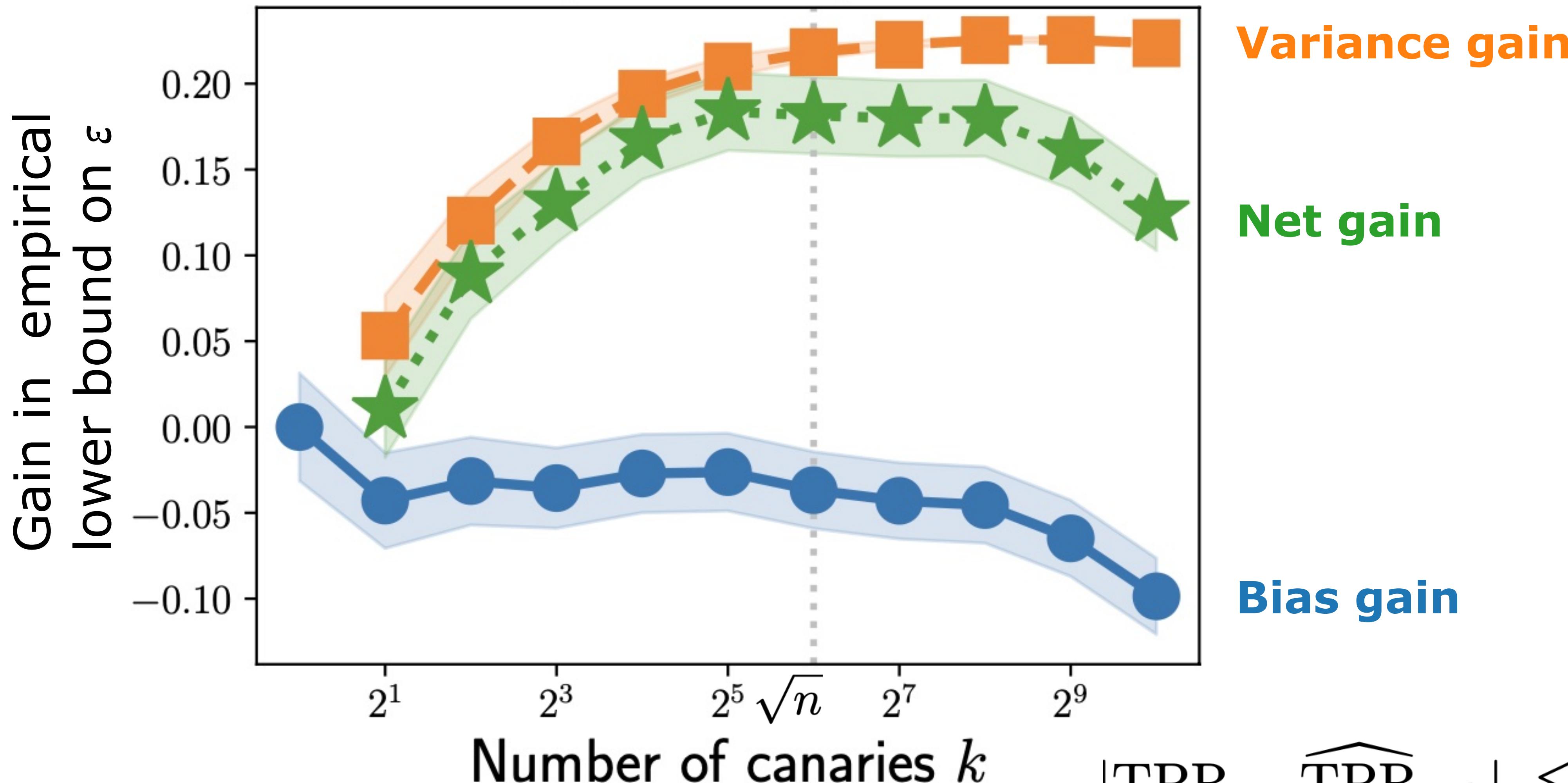
Thomas Steinke*
Google DeepMind
steinke@google.com

Milad Nasr*
Google DeepMind
srxzr@google.com

Matthew Jagielski*
Google DeepMind
jagielski@google.com

Bias-variance tradeoff in the number of canaries k

$$\varepsilon = 4.0, n = 4096, d = 10^4$$

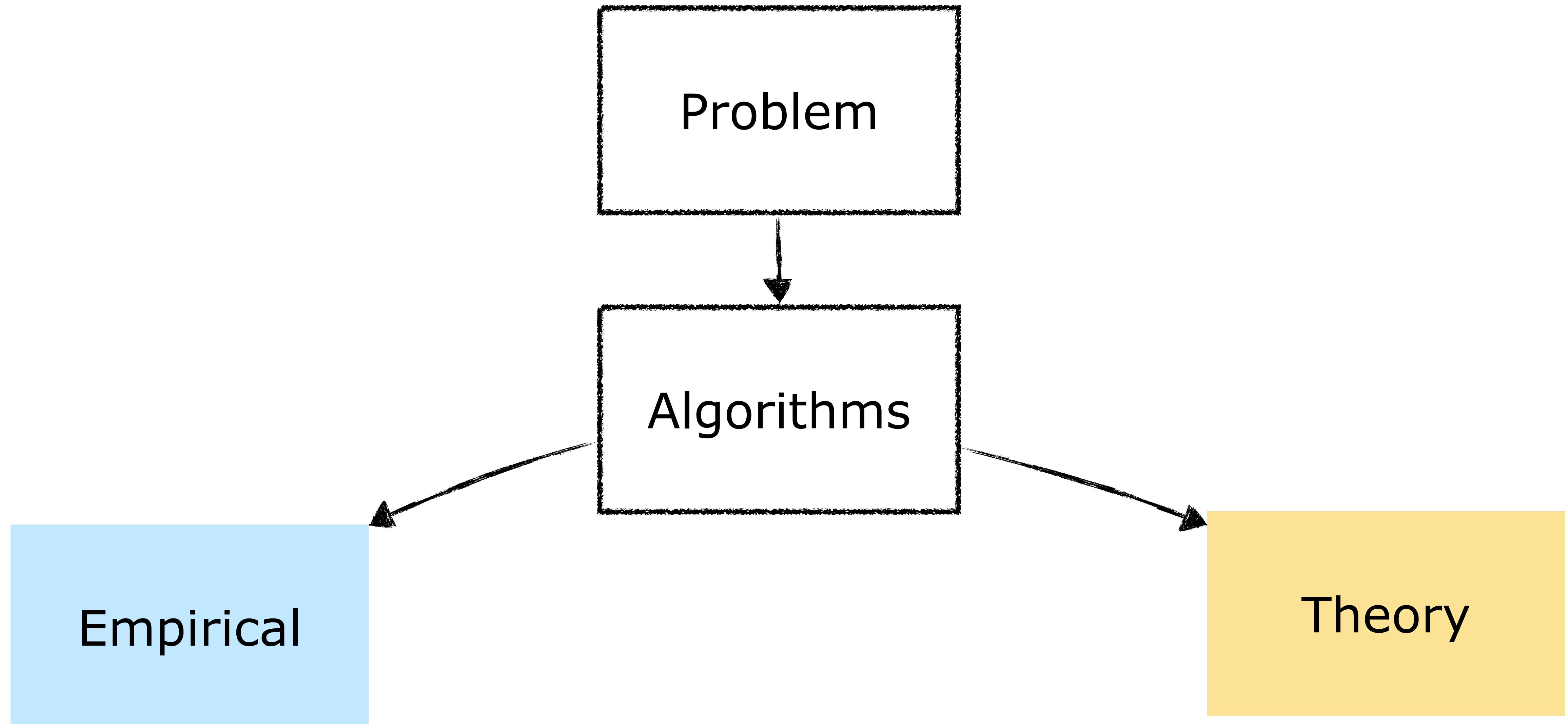


$$|\text{TPR} - \widehat{\text{TPR}}_{n,k}| \lesssim \sqrt{\frac{1}{nk}} + \frac{1}{n^{3/4}}$$

Summary

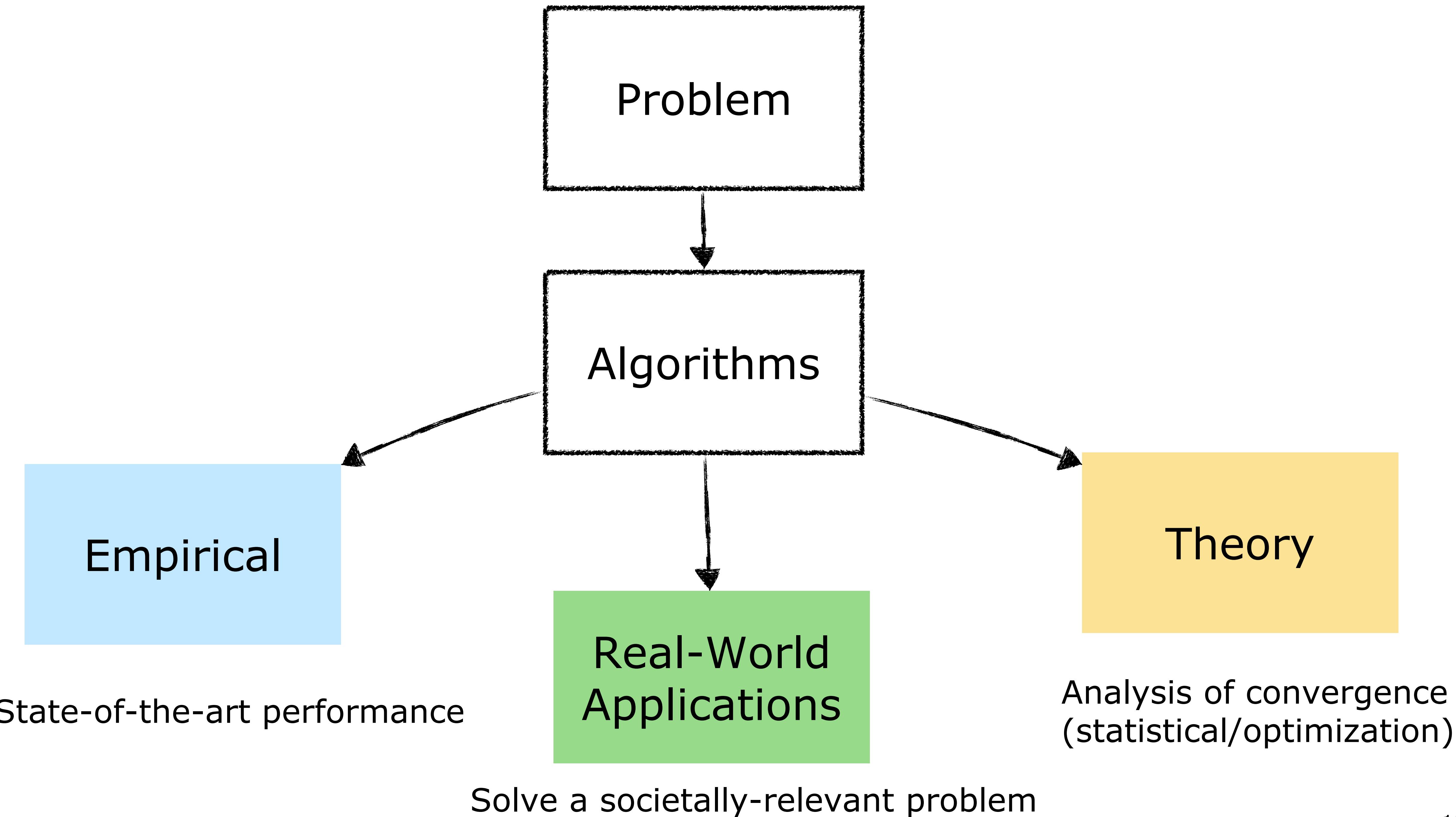
- **Auditing Lifted DP** (equivalent to usual DP) using multiple **i.i.d. random canaries** to improve sample dependence of the confidence intervals
- Can integrate with existing recipes for designing canaries

Ongoing Projects



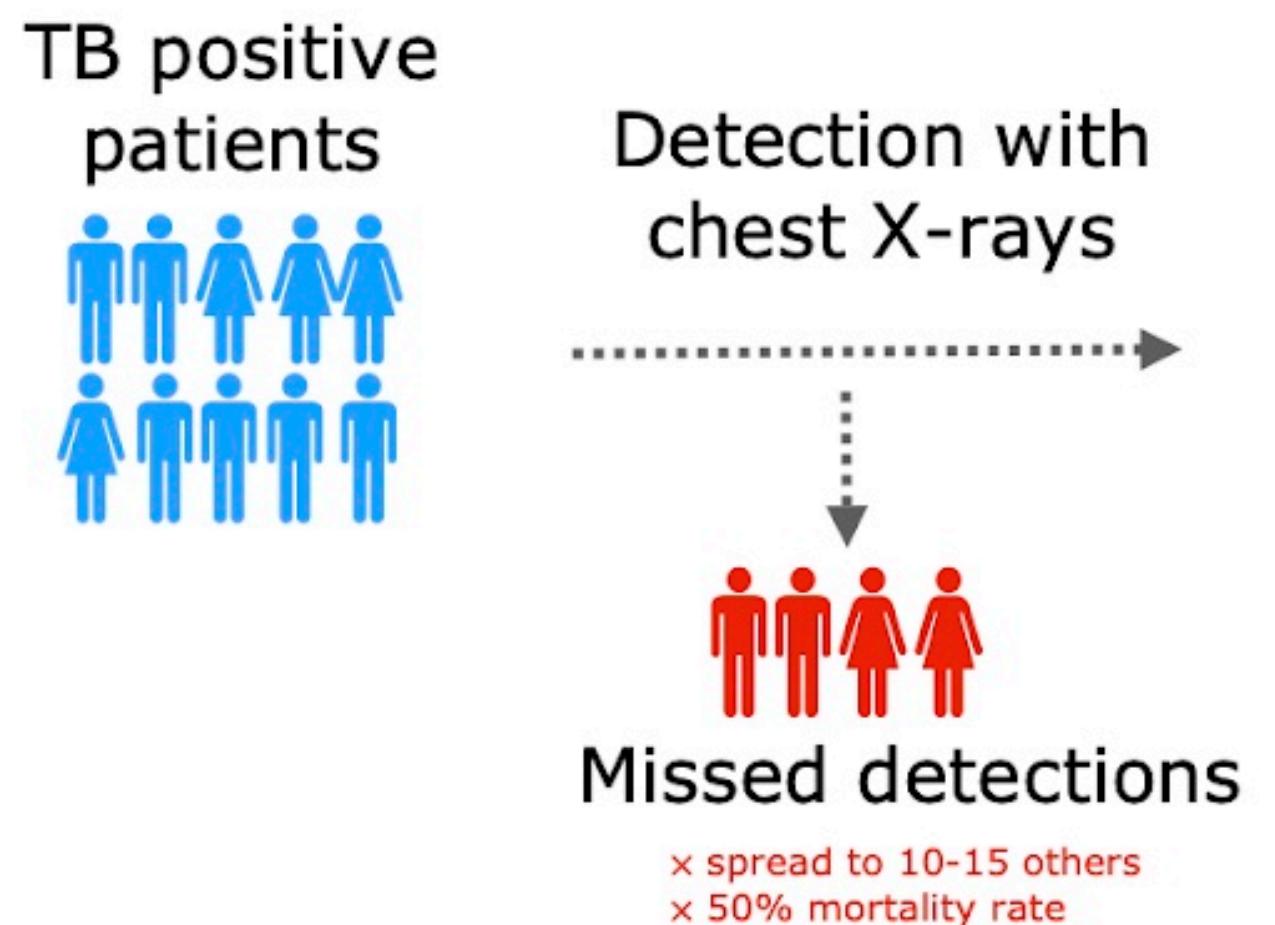
State-of-the-art performance

Analysis of convergence
(statistical/optimization)



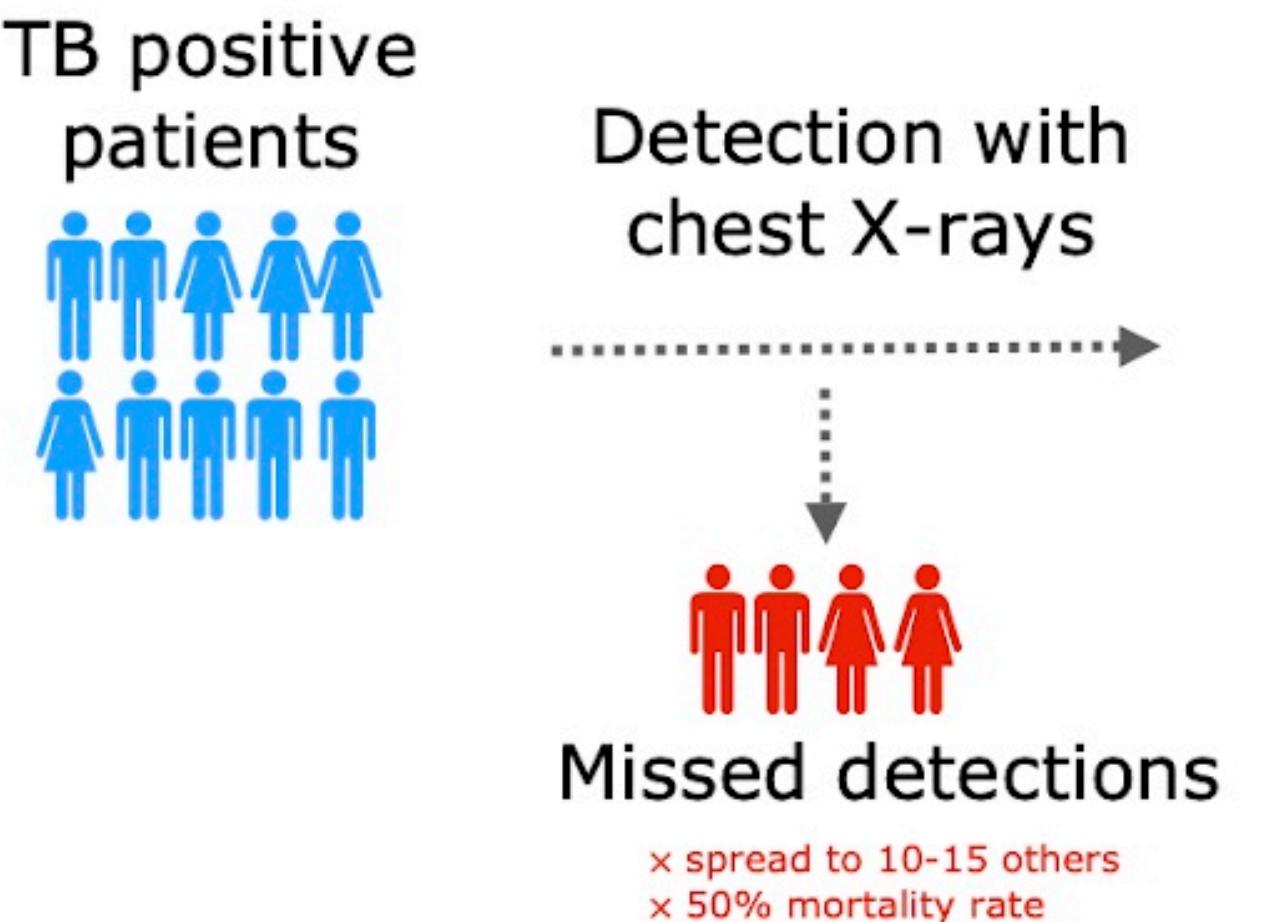
TB Detection with Privacy-Preserving AI

Usual Approach

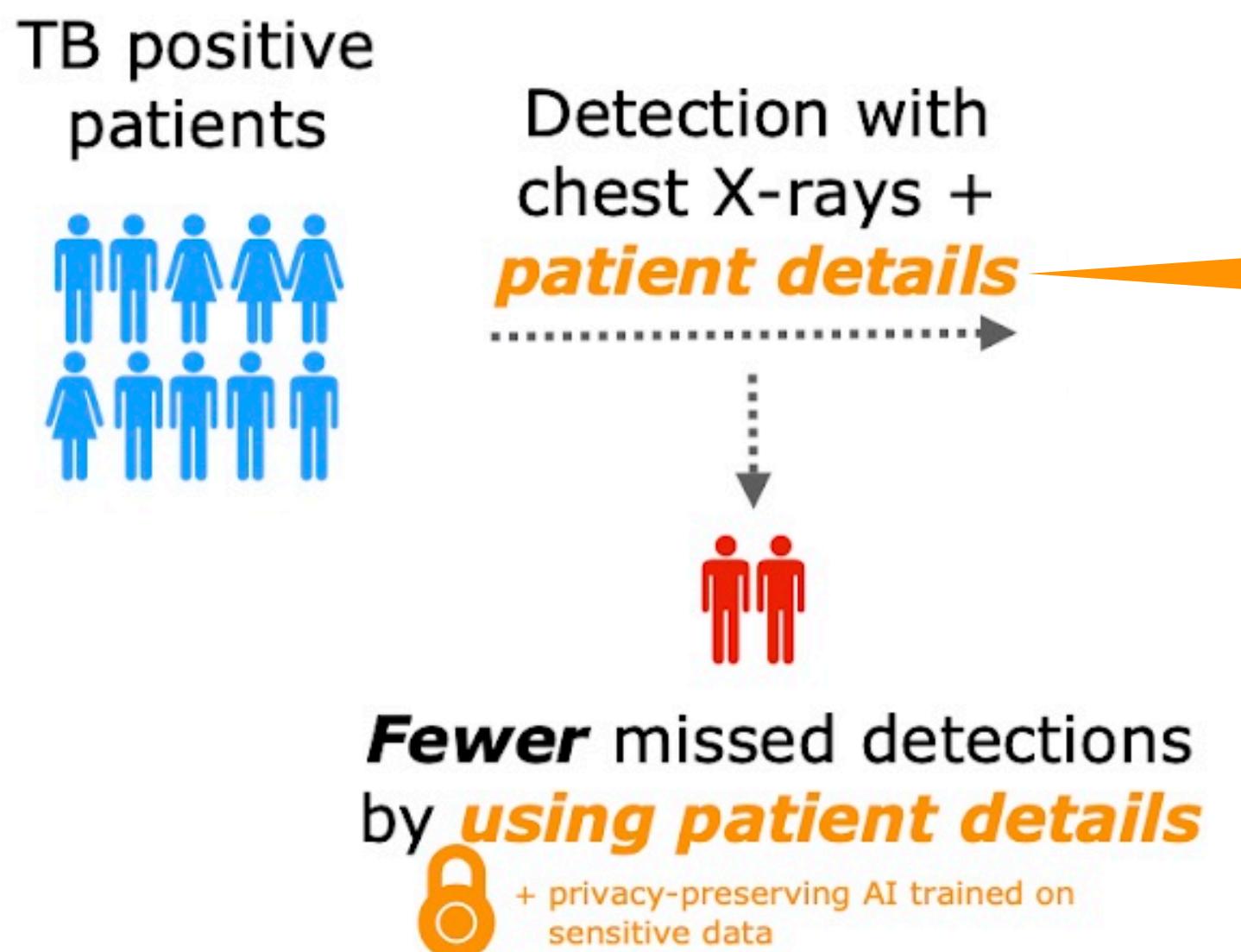


TB Detection with Privacy-Preserving AI

Usual Approach



Our Approach

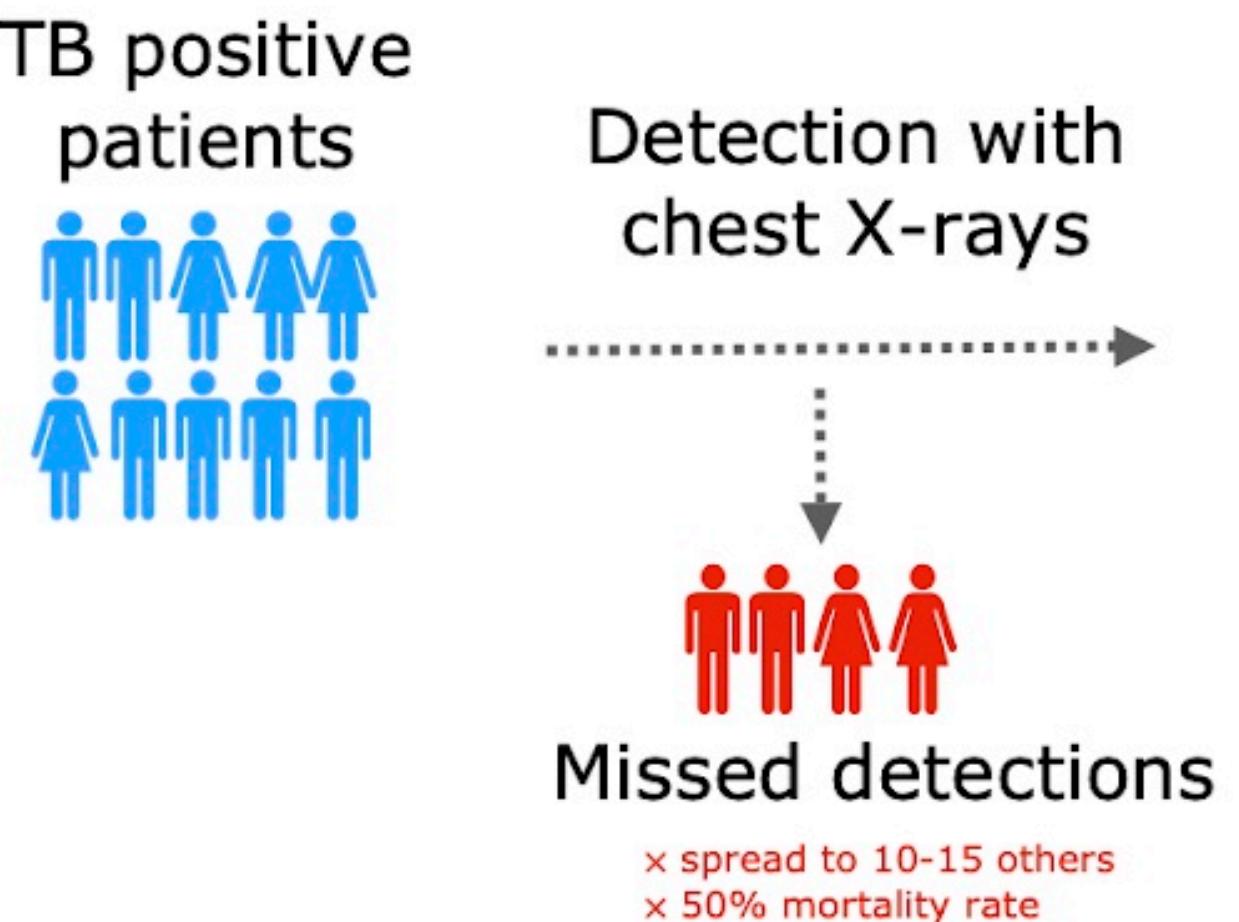


Privacy-sensitive!

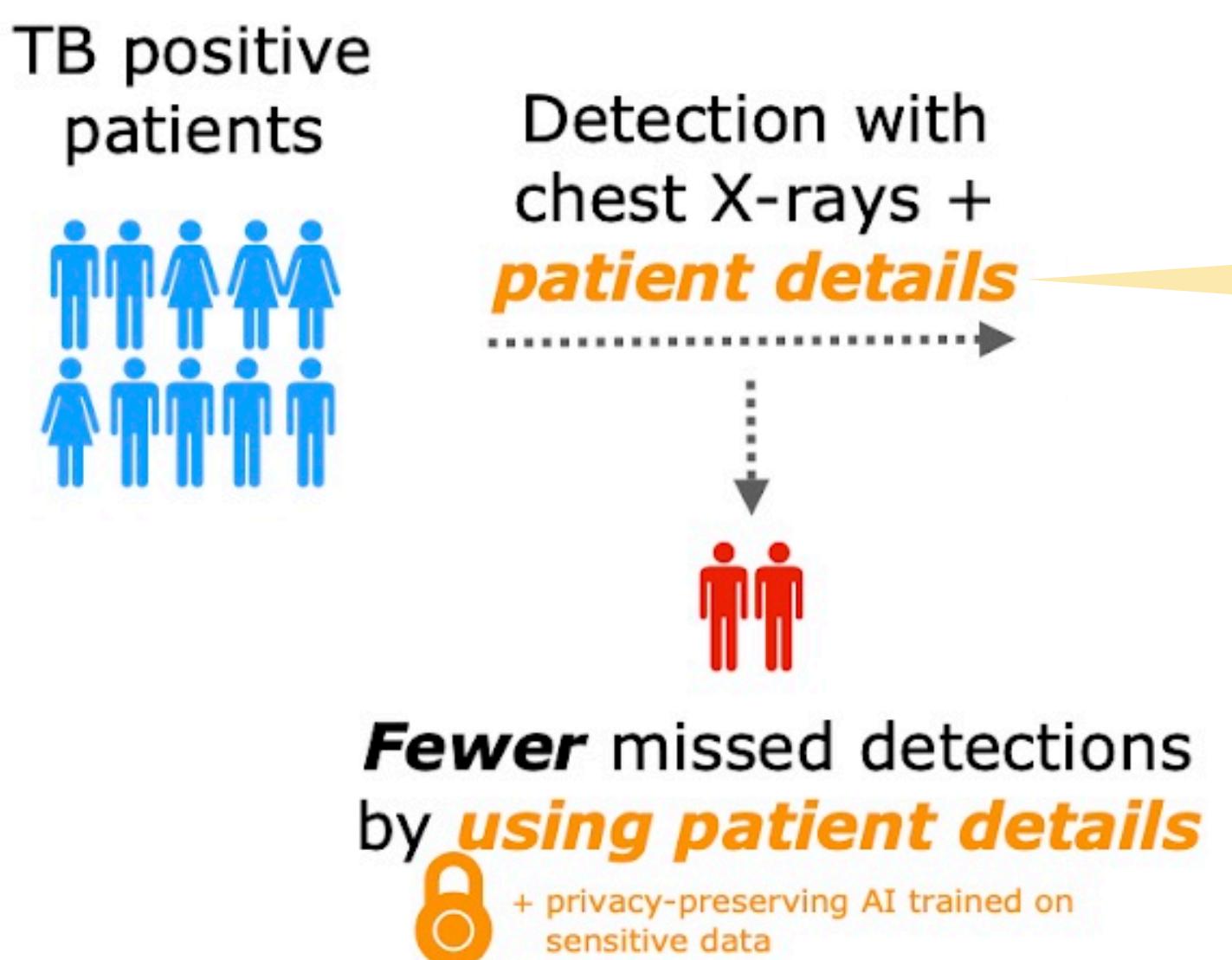


TB Detection with Privacy-Preserving AI

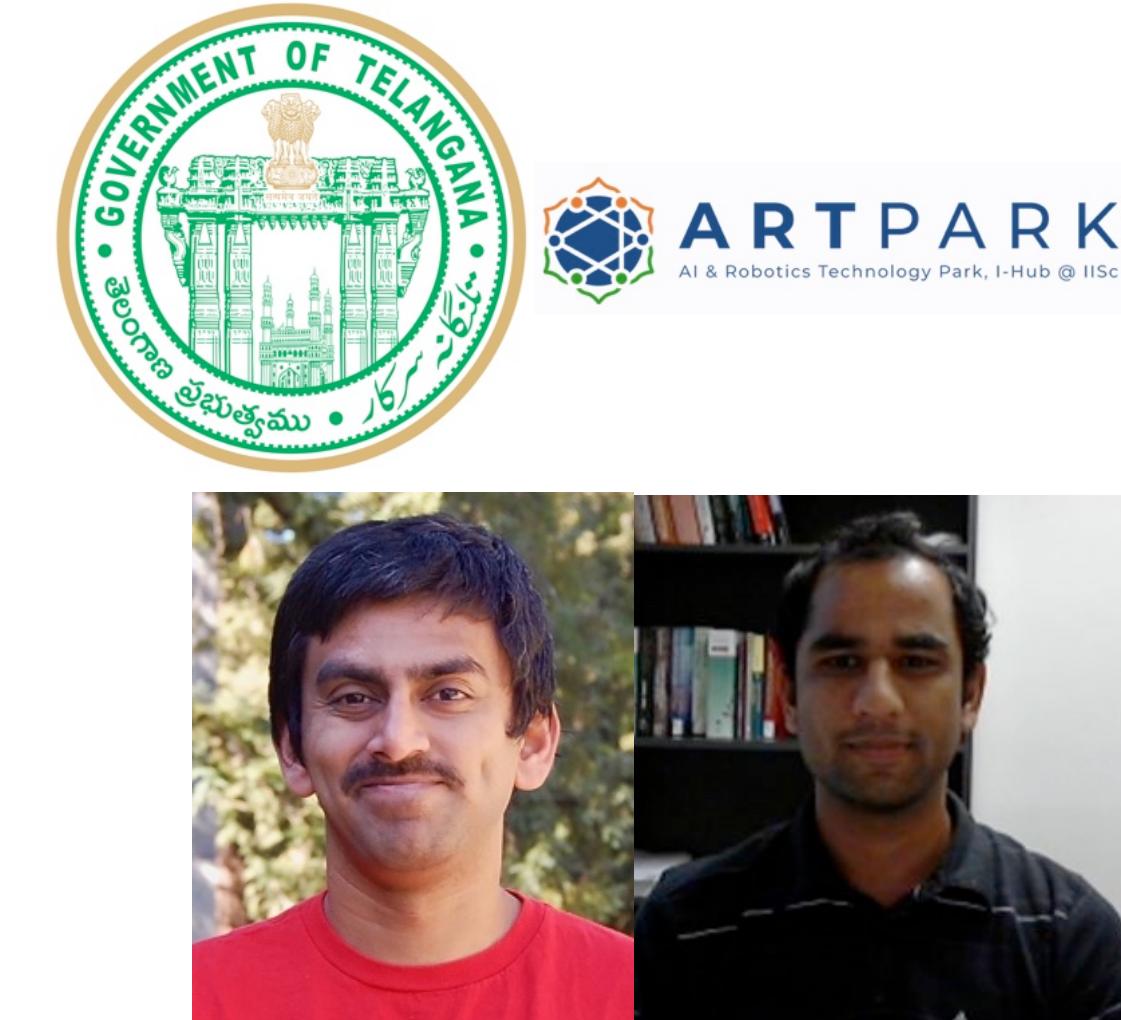
Usual Approach



Our Approach

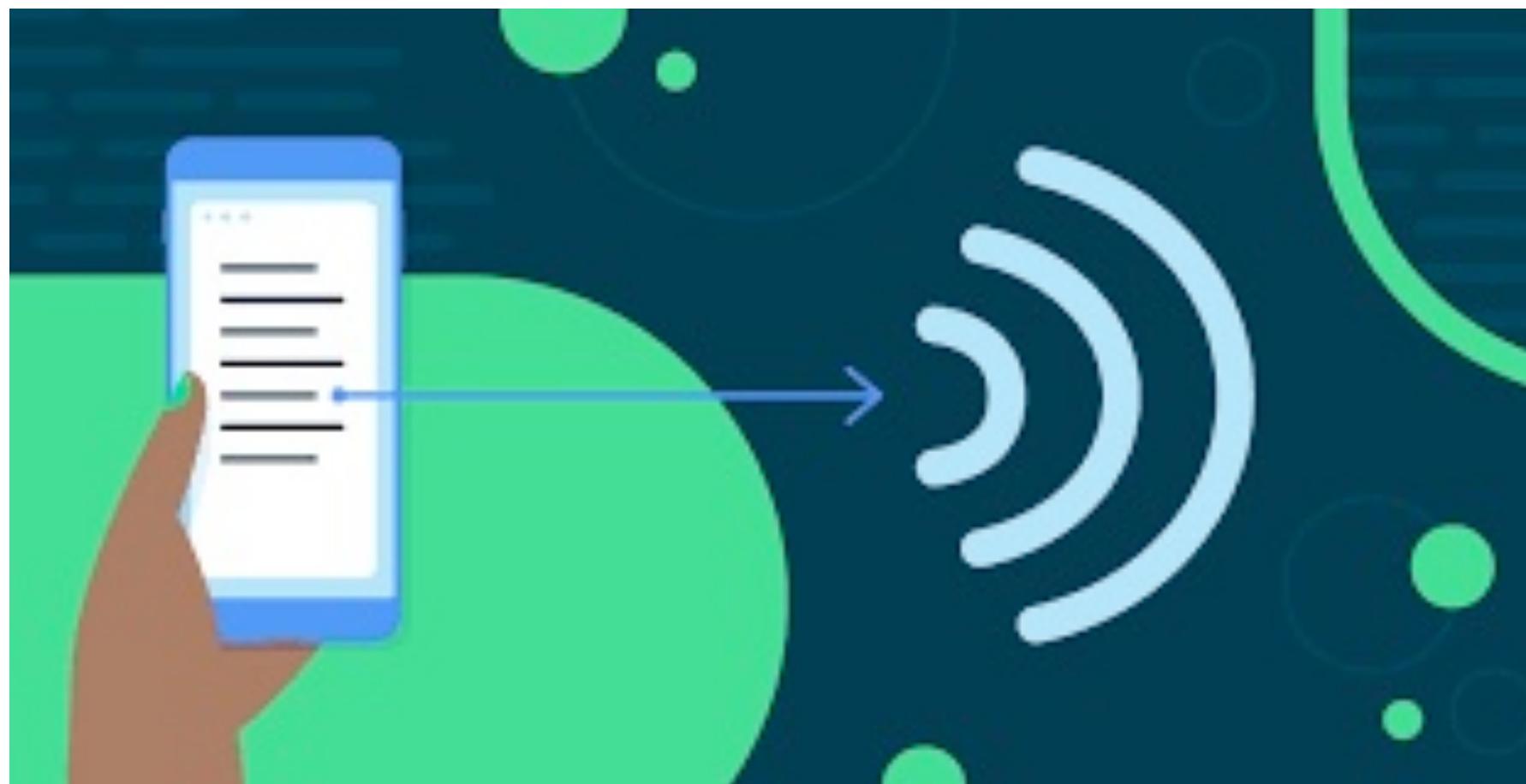


Technical Problem:
Mixed public-private
multi-modal learning



Privacy-sensitive synthetic data generation

Text-to-speech data augmentation



Financial fraud detection



Review of Gen AI Models for Financial Risk Management

👤 Satyadhar Joshi

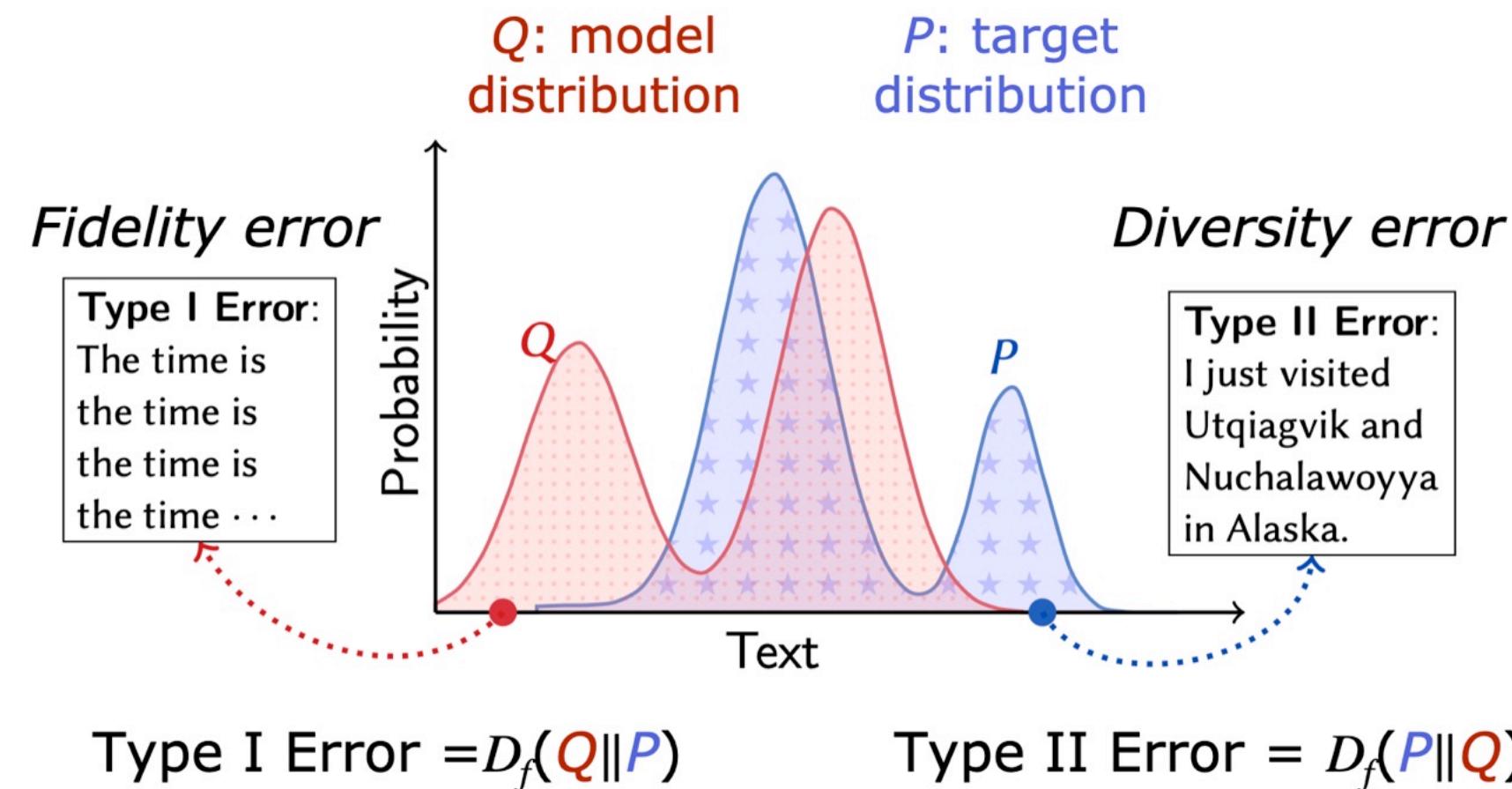
🏛️ BOFA Jersey City, USA

Useful to address data imbalances, bias

Evaluating synthetic data from LLMs/Gen AI

Project Page: <https://krishnap25.github.io/mauve-overview/>

- **Pillutla***, Liu*, Thickstun, Welleck, Swayamdipta, Zellers, Oh, Choi, Harchaoui. **MAUVE Scores for Generative Models: Theory and Practice.** *Journal of Machine Learning Research* (2023).
- **Pillutla**, Swayamdipta, Zellers, Thickstun, Welleck, Choi, Harchaoui. **MAUVE: Measuring the Gap Between Neural Text and Human Text using Divergence Frontiers.** *NeurIPS* (2021). **Outstanding Paper Award.**
- Liu, **Pillutla**, Welleck, Oh, Choi, Harchaoui. **Divergence Frontiers for Generative Models: Sample Complexity, Quantization Effects, and Frontier Integrals.** *NeurIPS* (2021).



mauve-text 0.4.0

pip install mauve-text

```
import mauve

# call mauve.compute_mauve using raw text on GPU 0; each
out = mauve.compute_mauve(p_text=p_text, q_text=q_text,
print(out.mauve) # prints 0.9917
```



Evaluate

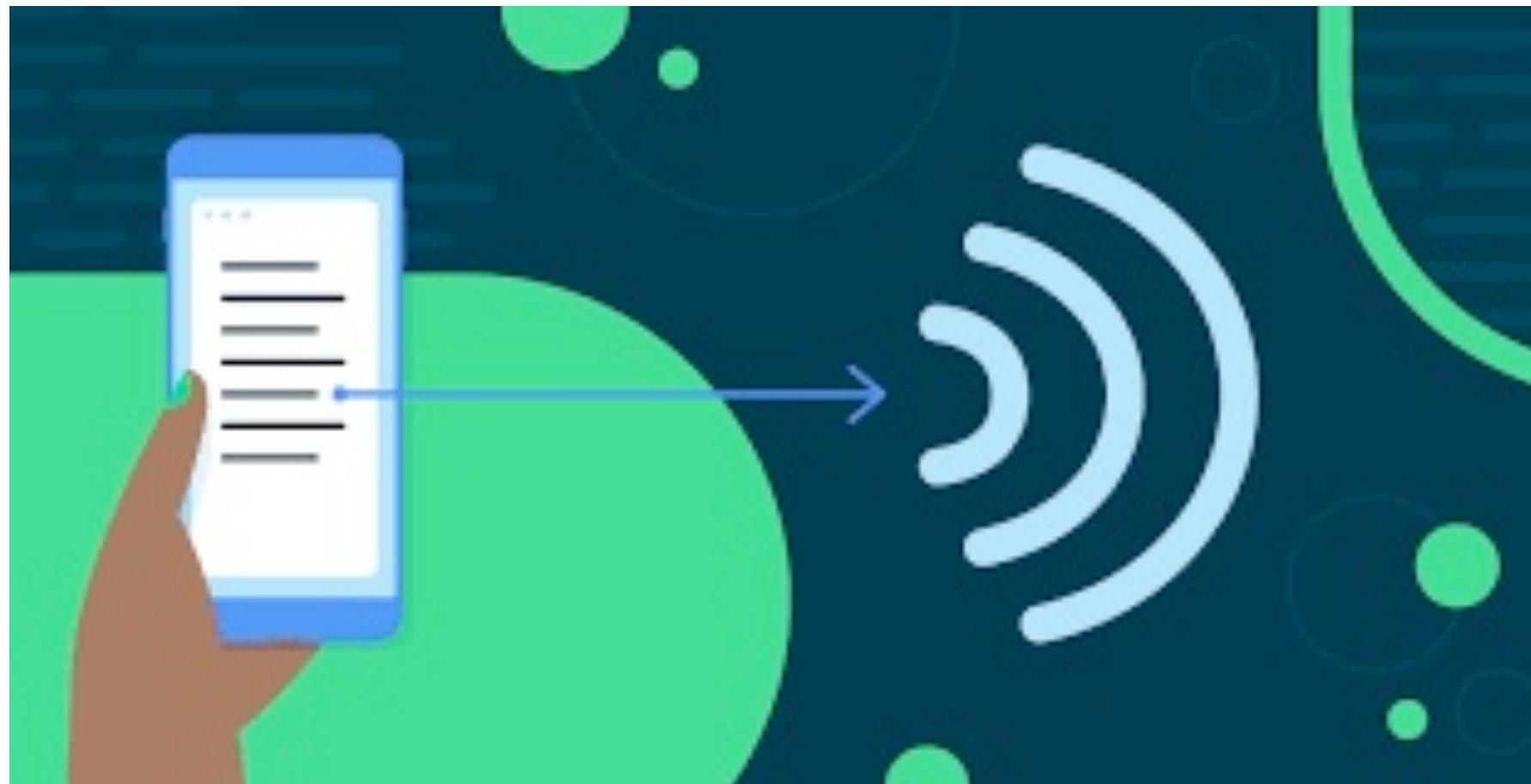
Outstanding Paper



NeurIPS 2021

Privacy-sensitive synthetic data generation

Text-to-speech data augmentation



Financial fraud detection



Review of Gen AI Models for Financial Risk Management

👤 Satyadhar Joshi

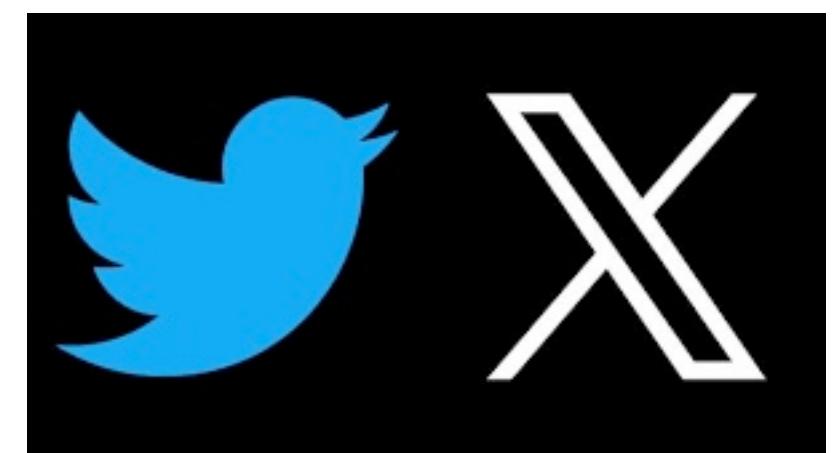
🏛️ BOFA Jersey City, USA

Useful to address data imbalances, bias

Thank you!



<https://krishnap25.github.io>



@KrishnaPillutla



<https://www.linkedin.com/in/krishna-pillutla-a0b1b2a8/>



Advertisement: MS/PhD Openings in my group at IIT Madras

- Areas of interest in ML/AI:
 - Privacy-preserving AI
 - Making (generative) AI more robust
 - Applications in healthcare + public good
- Flavour:
 - Theoretical foundations +
 - State of the art empirical performance +
 - Real-world applications