

Towards User-level Differential Privacy at Scale

Krishna Pillutla

Google Research -> IIT Madras

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

AHA, FOUND THEM!

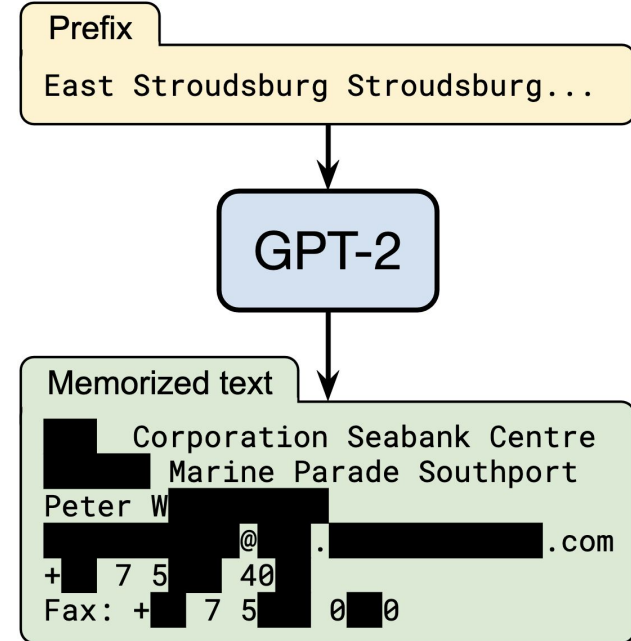


WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data

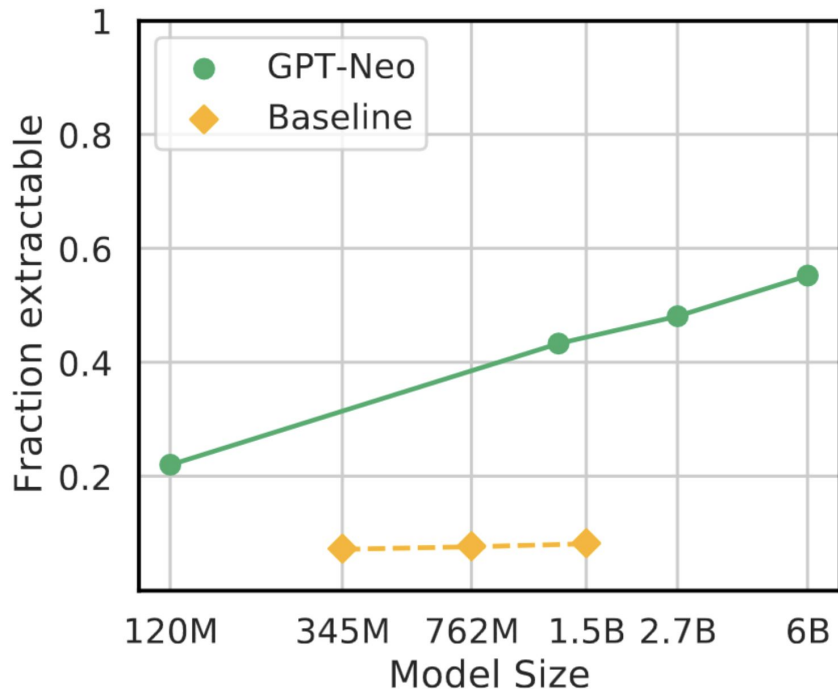


WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

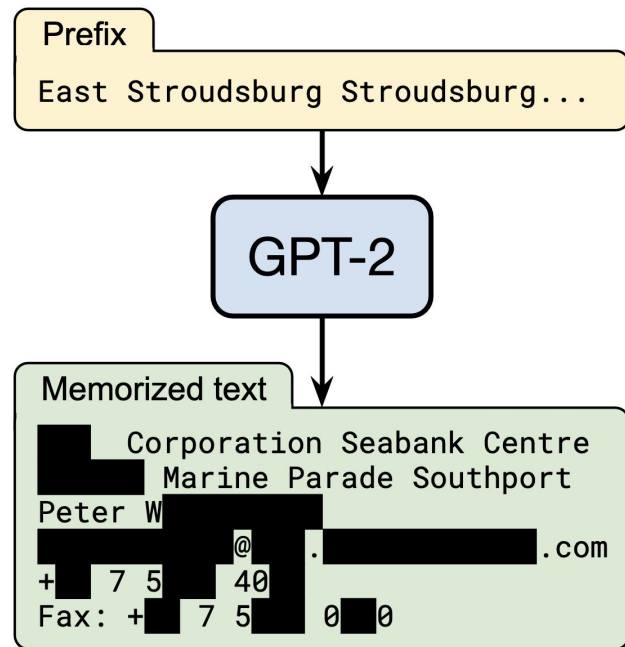


Carlini et al. (USENIX Security 2021)

Models leak information about their training data *reliably*



Carlini et al. (ICLR 2023)



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

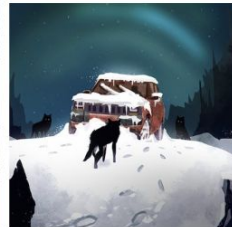
🇺🇸 New York University

goldblum@nyu.edu

Generation



LAION-A Match



Differential privacy (DP)

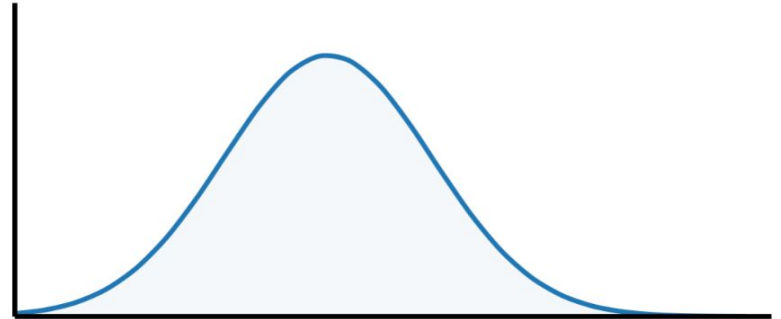
Dataset



Randomized
Algorithm



Output Distribution
(e.g. over models)



Differential privacy (DP)

Dataset



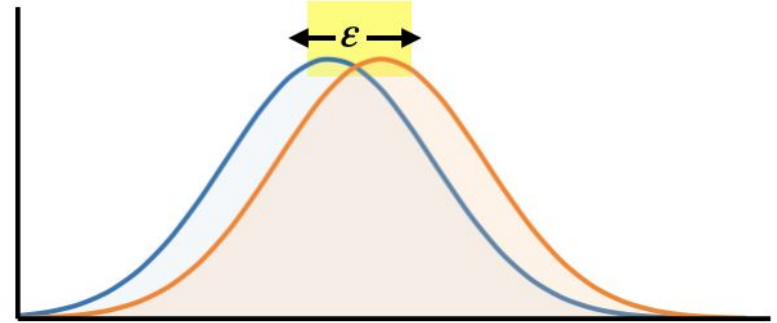
+



Randomized
Algorithm



Output Distribution
(e.g. over models)



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

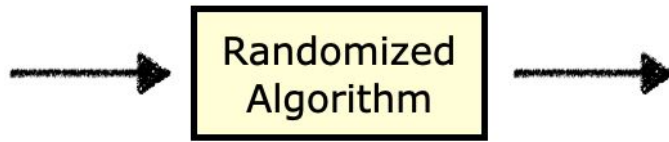
 Unit of data
= **example**

Example-level Differential privacy (DP)

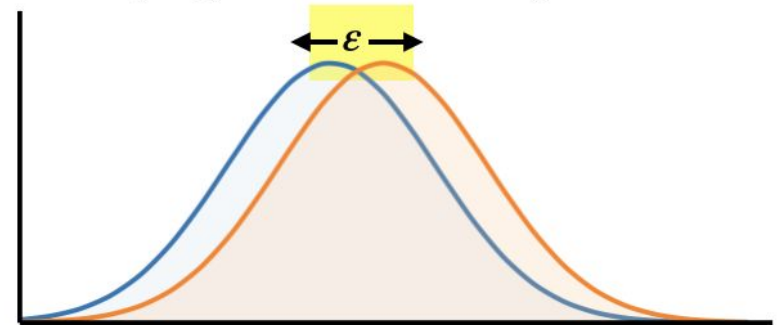
Dataset



+

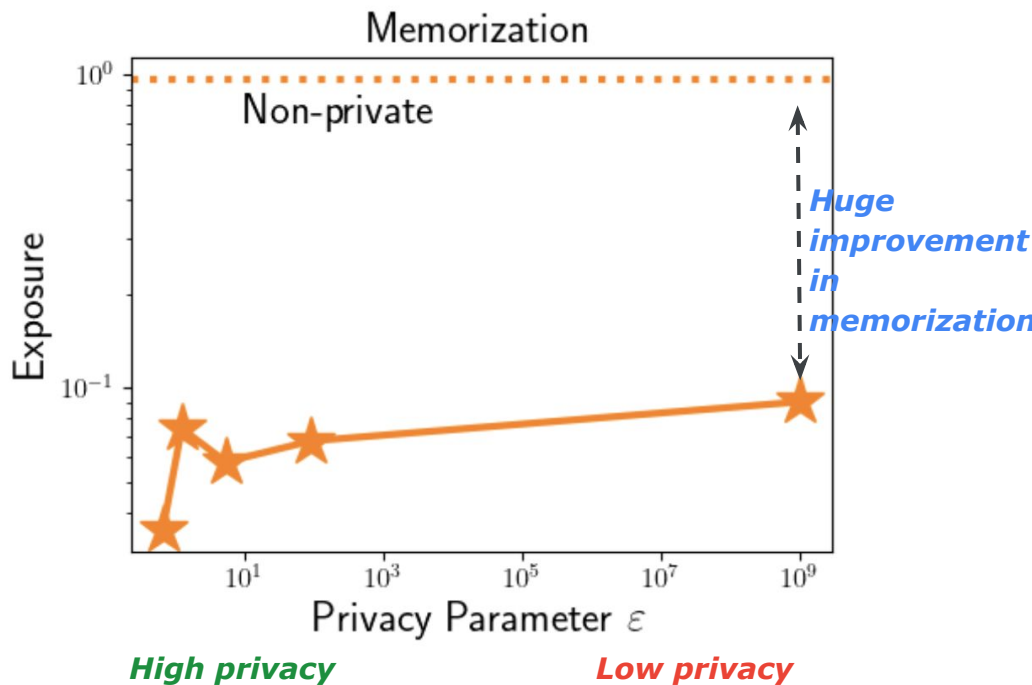
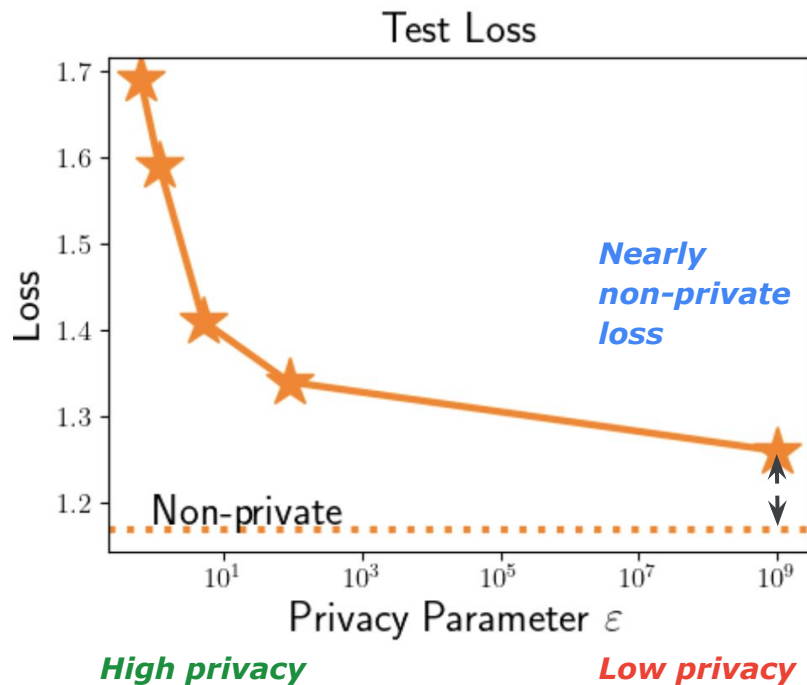


Output Distribution
(e.g. over models)



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 eliminates memorization

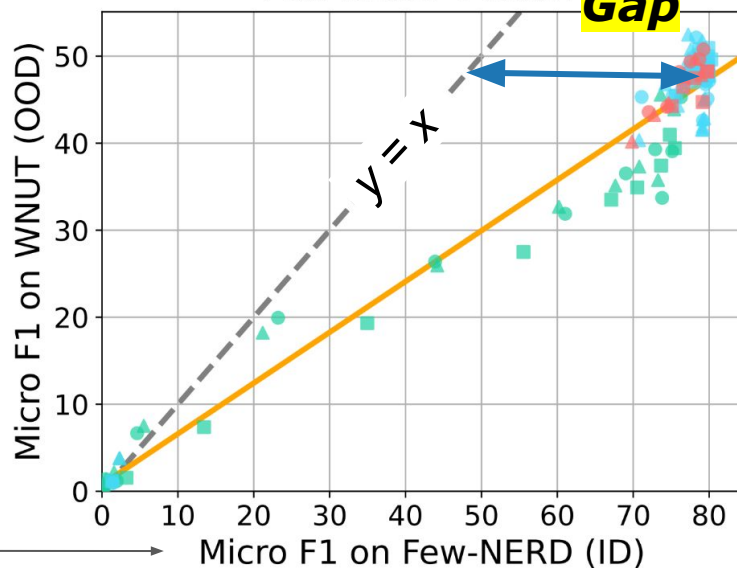


Which data do we use to train/finetune/align these models?

Robustness

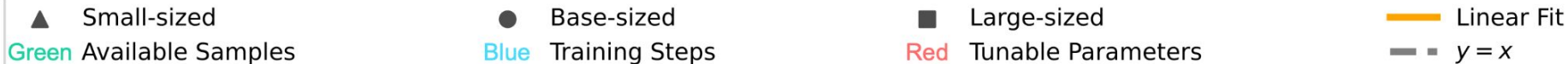
Gap

Few-NERD → WNUT



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

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

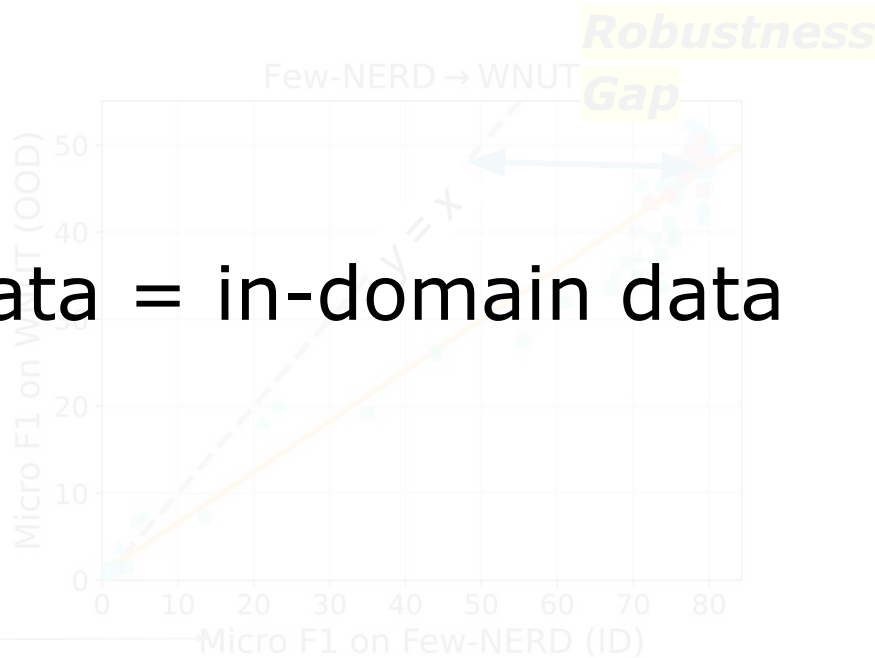


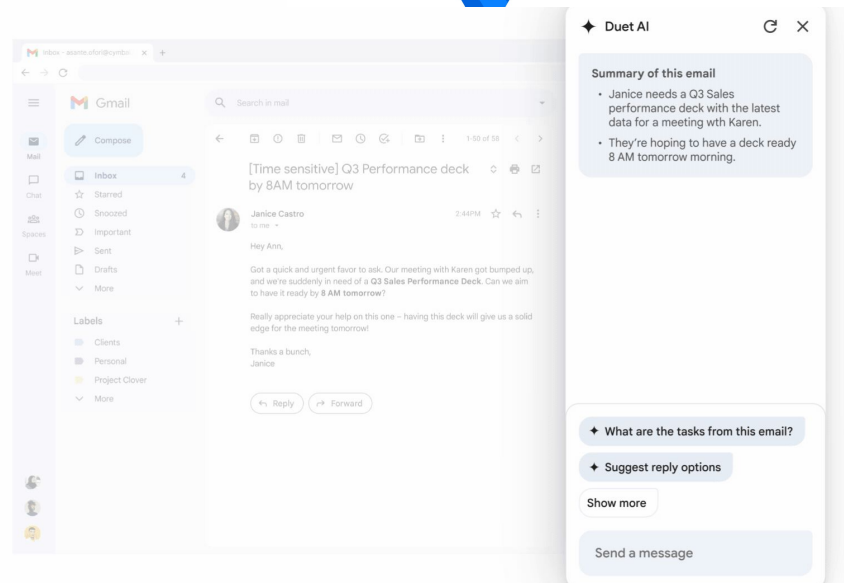
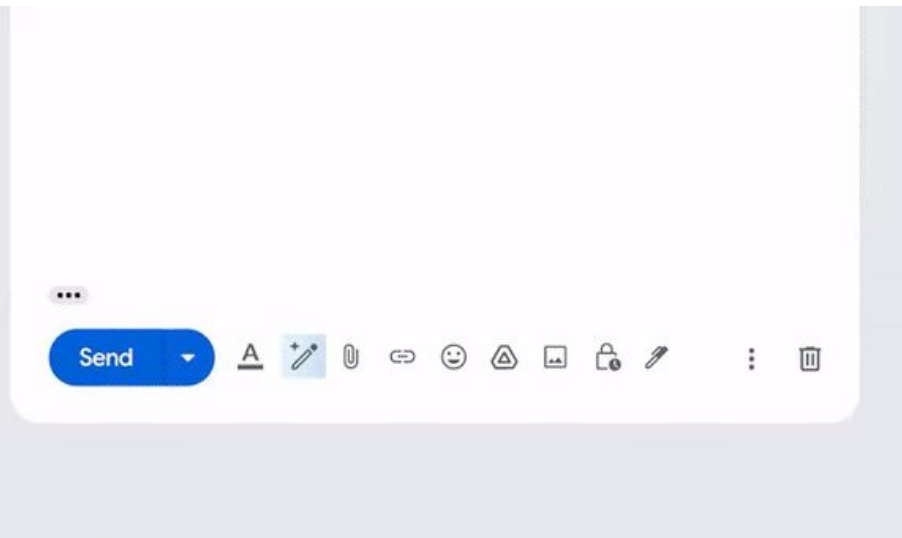
Which data do we use to train/finetune/align these models?

Best training data = in-domain data

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

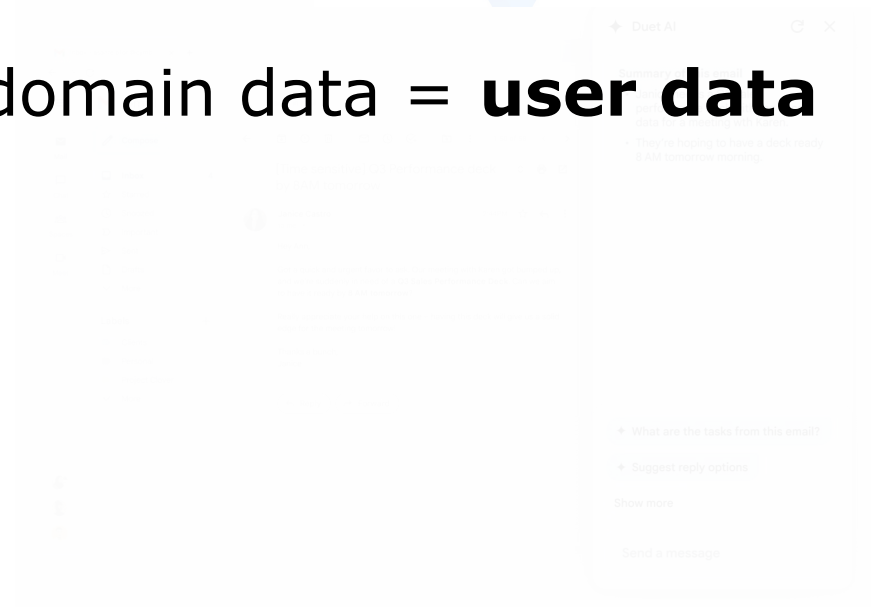
Test on **training distribution**
(in-domain / ID)







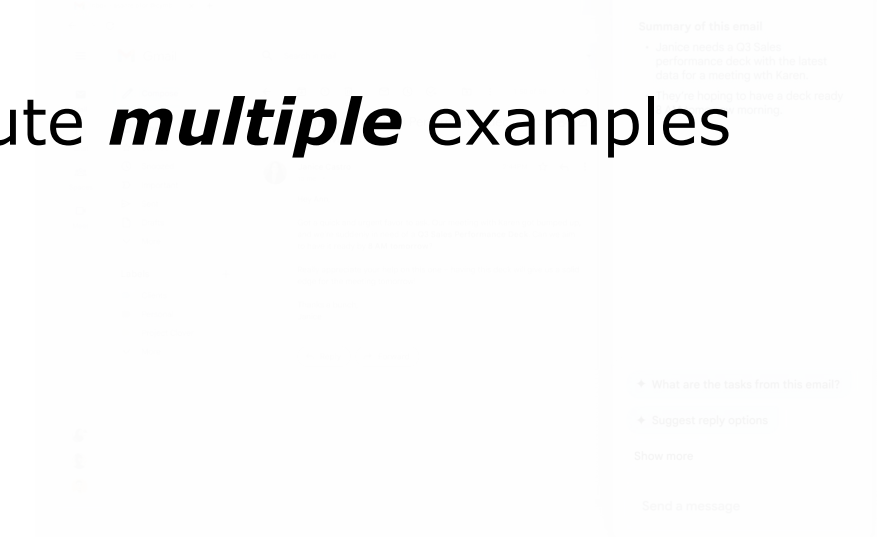
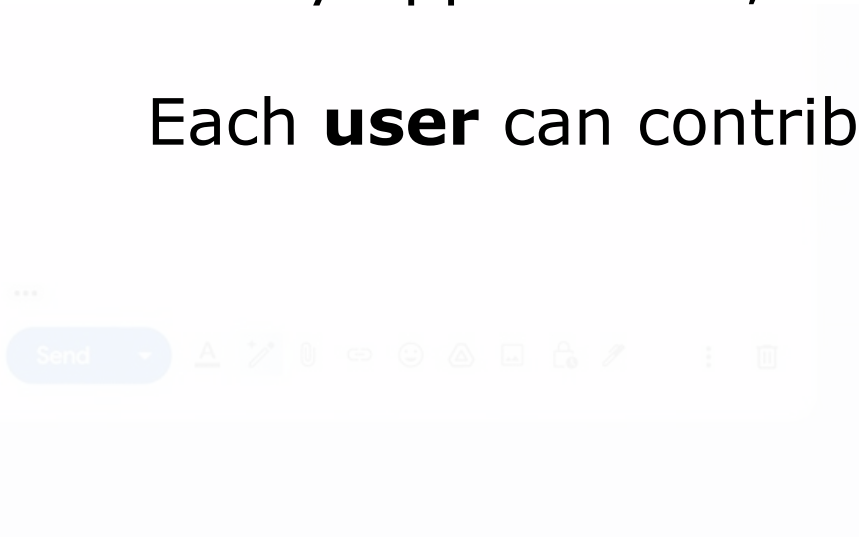
For many applications, in-domain data = **user data**





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

Each **user** can contribute *multiple* examples



Unit of data
= **example**

Example-level Differential privacy (DP)

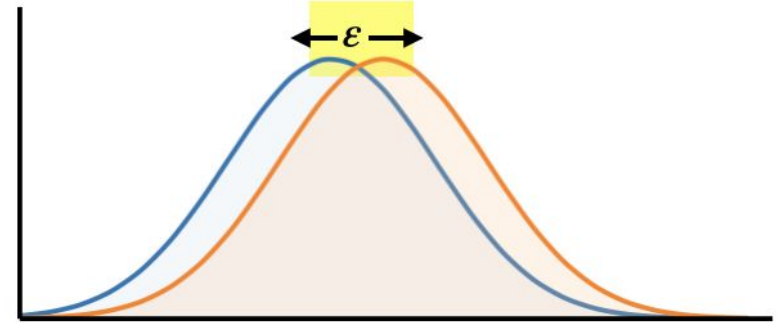
Dataset



+



Output Distribution
(e.g. over models)

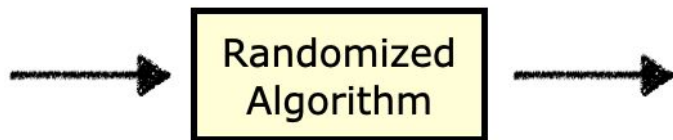


A randomized algorithm is ϵ -**differentially private** if the addition of **one example** does not alter its output distribution by more than ϵ

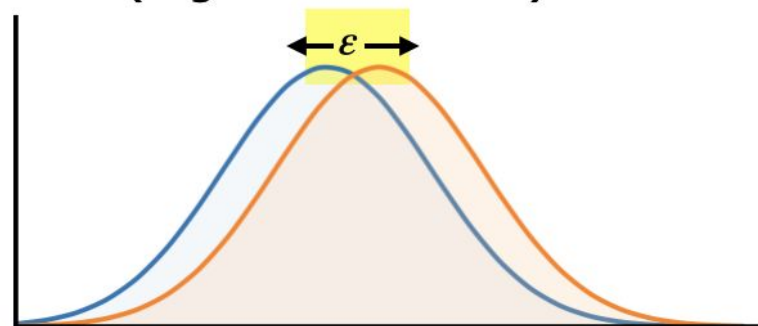
User Example-level Differential privacy (DP)

Unit of data
= **user**

Dataset



Output Distribution
(e.g. over models)



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?

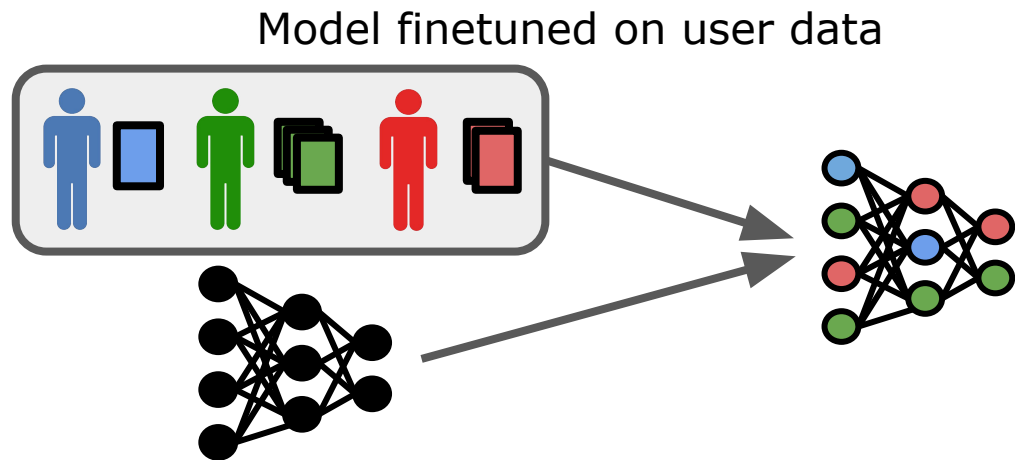
Why do we need user-level DP?

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

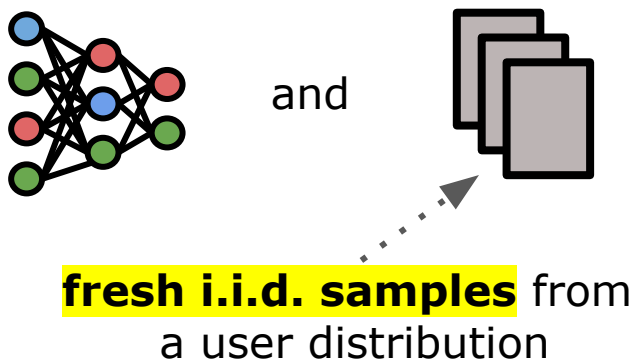
Nikhil Kandpal, **KP**, Alina Oprea, Peter Kairouz, Chris Choquette-Choo, Zheng Xu.
Submitted (2024)



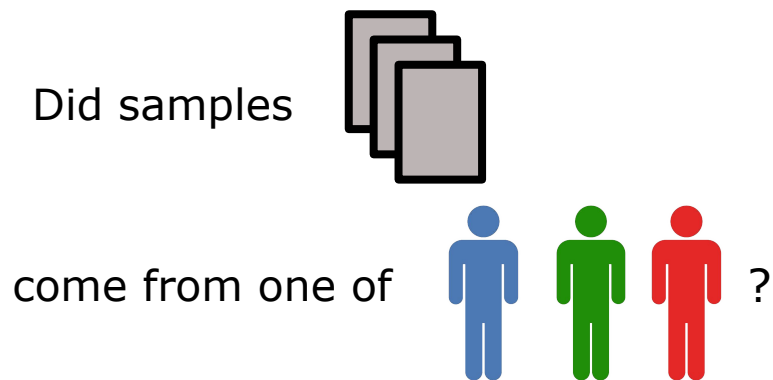
User Inference Attack



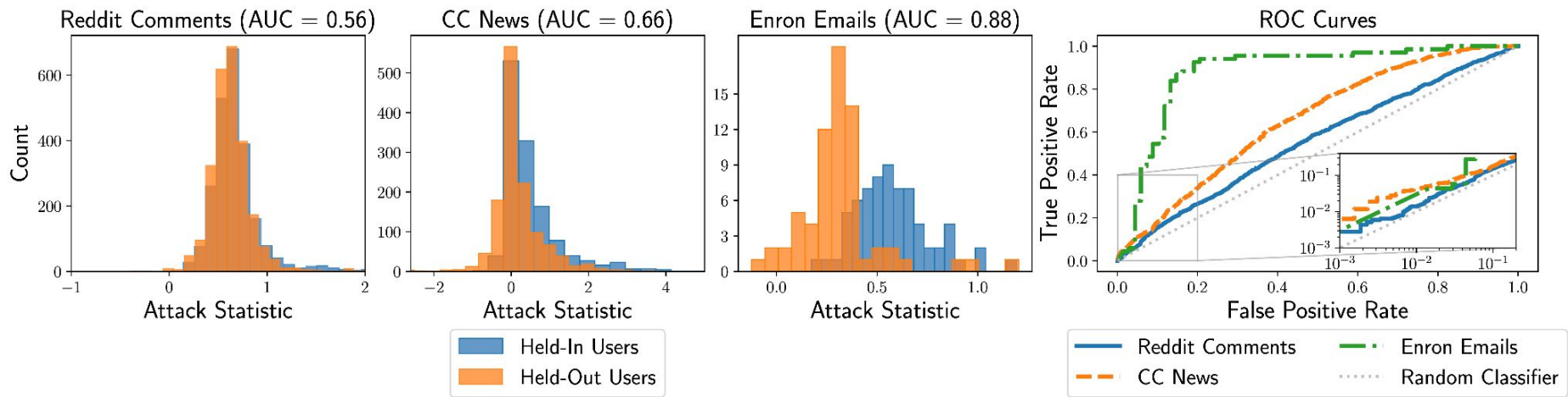
Attacker Has:



Attacker Wants to Infer:



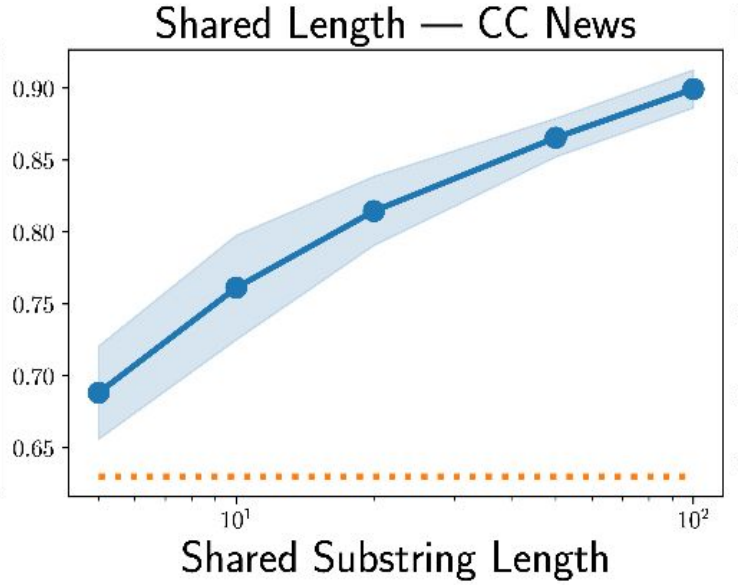
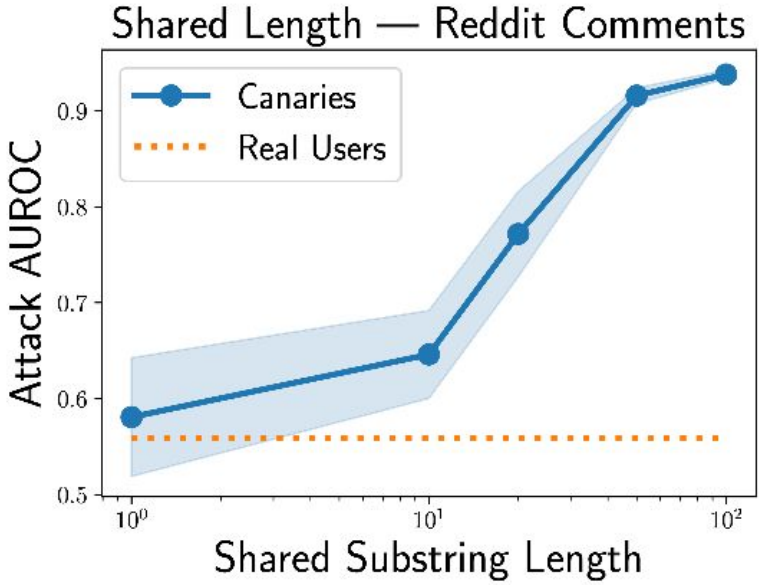
User inference is effective when #users is small and data per user is large



More fine-tuning samples per user

More users

Short common phrases can exacerbate user inference



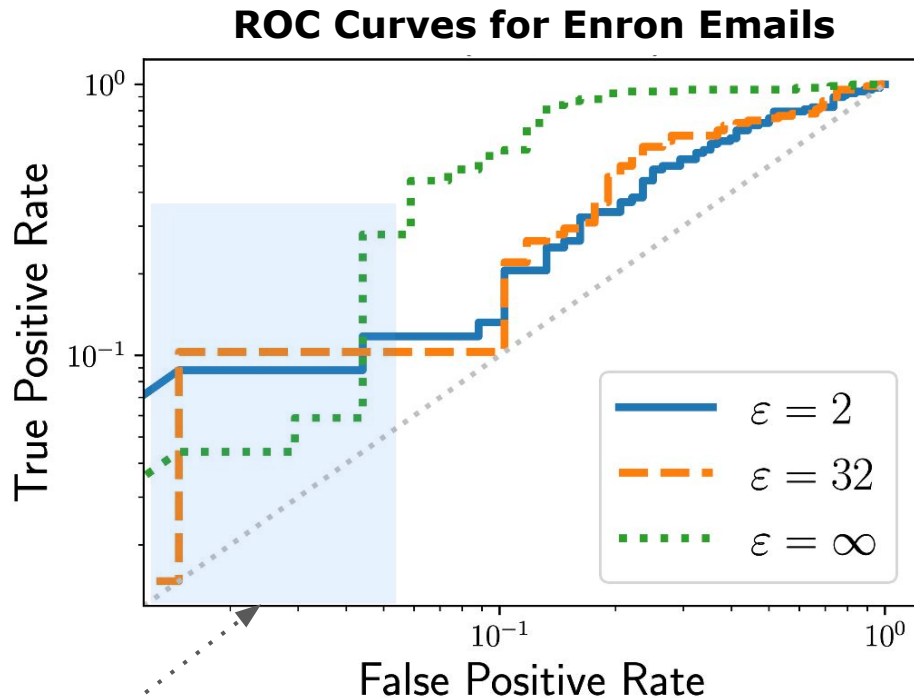
Example-level DP offers limited mitigation

AUROC:

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

Utility:

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

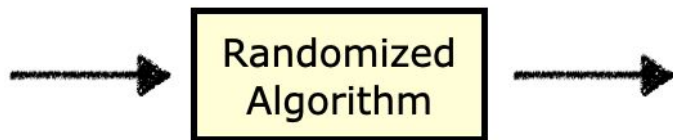


Example-level DP does not help here

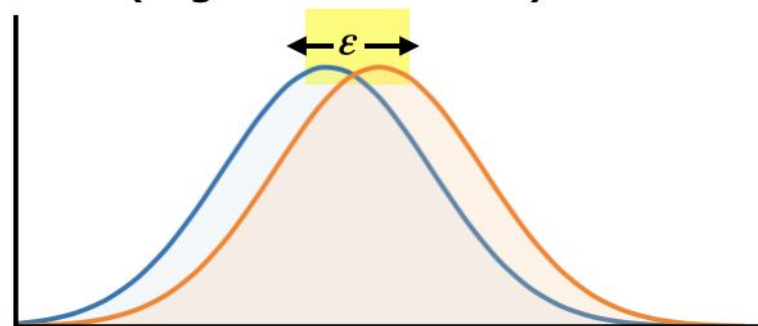
User Example-level Differential privacy (DP)

Unit of data
= **user**

Dataset



Output Distribution
(e.g. over models)



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:

(Anti-) correlated noise provably beats independent noise

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

Independent noise	$\Theta(d)$
Correlated noise	$\tilde{O}(d_{\text{eff}})$
Lower bound	$\Omega(d_{\text{eff}})$

Outline: how do we realize user-level DP?

Learning algorithms:

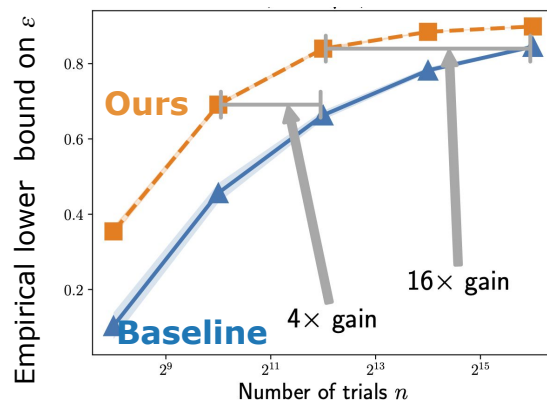
(Anti-) correlated noise provably beats independent noise

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

<i>Independent noise</i>	$\Theta(d)$
<i>Correlated noise</i>	$\tilde{O}(d_{\text{eff}})$
<i>Lower bound</i>	$\Omega(d_{\text{eff}})$

Auditing:

Randomness makes the audit more computationally efficient



Part 1: How do we learn with user-level DP?

(Anti-)correlated noise **provably** beats independent noise

ICLR 2024



**Chris
Choquette-Choo***



**Dj
Dvijotham***



**Krishna
Pillutla***



Arun
Ganesh



Thomas
Steinke



Abhradeep
Thakurta

*Equal contribution, $\alpha\beta$ -order

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

The diagram illustrates the DP-SGD update equation: $\theta_{t+1} = \theta_t - \eta (g_t + z_t)$. Three callout boxes provide details: the top-left box explains that the stochastic gradient g_t is clipped to $\|g_t\| \leq G$ per-example; the top-right box identifies z_t as independent Gaussian noise; and the bottom-left box identifies η as the learning rate.

Stochastic gradient
clipped to $\|g_t\| \leq G$
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**?

The diagram illustrates the DP-FedAvg update equation: $\theta_{t+1} = \theta_t - \eta (g_t + z_t)$. Three callout boxes provide details: 'Stochastic gradient clipped to $\|g\| \leq G$ per-user' points to g_t ; 'Independent Gaussian noise' points to z_t ; and 'Learning rate' points to η .

Stochastic gradient clipped to $\|g\| \leq G$ **per-user**

Independent
Gaussian noise

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

Learning rate

DP-SGD: DP Training with ***Independent*** Noise

For ρ -zCDP, take
noise variance = $\frac{G^2}{2\rho}$

(G = gradient clip norm)

Independent
Gaussian noise

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

DP-FTRL: DP Training with ***Correlated*** Noise

Correlated
Gaussian noise
(z_t i.i.d. Gaussian)

$$\theta_{t+1} = \theta_t - \eta \left(g_t + \sum_{\tau=0}^t \beta_{t,\tau} z_{t-\tau} \right)$$

DP-FTRL: DP Training with ***Correlated*** Noise

For Q -zCDP, take
noise variance =

$$\frac{G^2}{2\rho} \max_t \|[B^{-1}]_{:,t}\|_2^2$$

sensitivity

$$B = \begin{pmatrix} \beta_{0,0} & 0 & 0 & \dots \\ \beta_{1,0} & \beta_{1,1} & 0 & \dots \\ \beta_{2,0} & \beta_{2,1} & \beta_{2,2} & \dots \\ \vdots & & & \end{pmatrix}$$

Correlated
Gaussian noise
(z_t i.i.d. Gaussian)

$$\theta_{t+1} = \theta_t - \eta \left(g_t + \sum_{\tau=0}^t \beta_{t,\tau} z_{t-\tau} \right)$$

Production Training

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

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

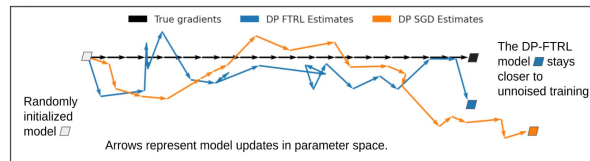
Federated Learning with Formal Differential Privacy Guarantees

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.

Do we use *independent* or *correlated* noise?



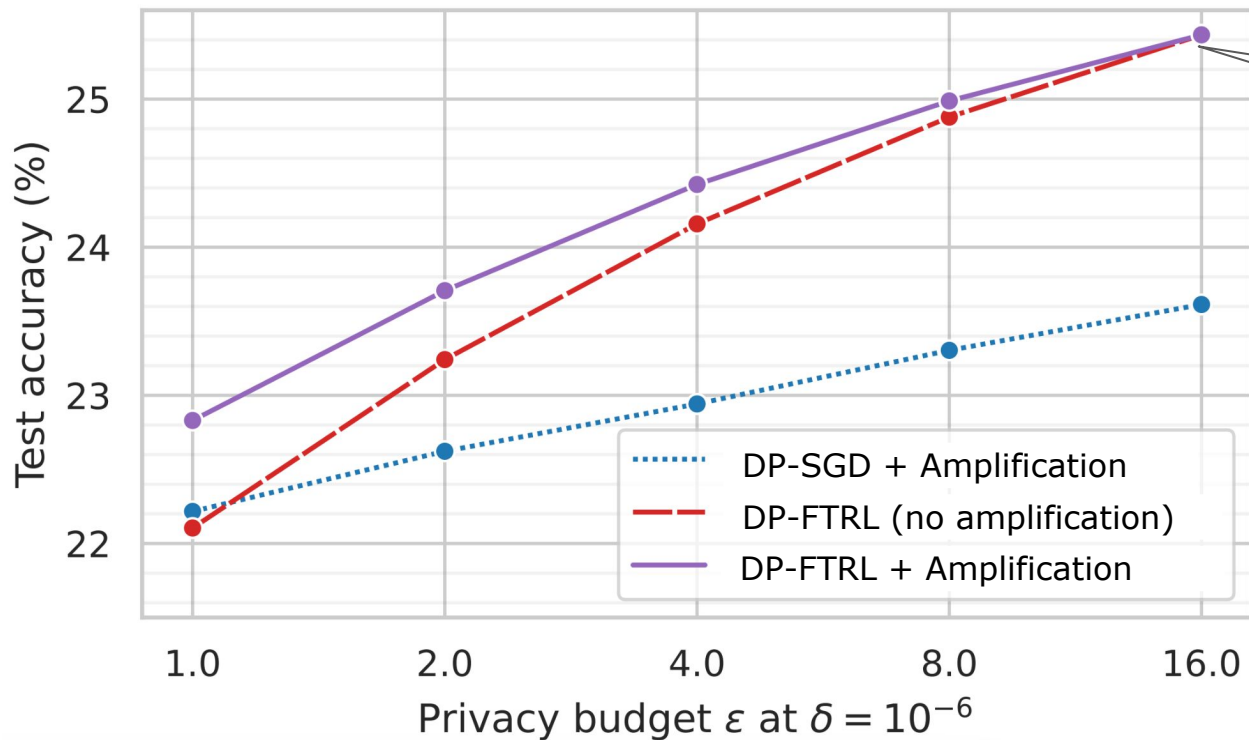
DP-SGD



DP-FTRL

Prior work: (Empirically) correlated noise outperforms independent noise

Experiment:
User-level DP with StackOverflow



DP-FTRL (+ amplification) uniformly beats **DP-SGD**

Our goal: a *provable* gap between DP-SGD & DP-FTRL

DP-FTRL vs. DP-SGD: Previous Theory

For convex & G -Lipschitz losses

DP-SGD	$\frac{Gd^{1/4}}{\sqrt{\rho T}}$
DP-FTRL	$\frac{Gd^{1/4}}{\sqrt{\rho^2 T}}$

ρ : privacy level (zCDP)

d : dimension

T : #iterations

Kairouz, McMahan, Song, Thakkar, Thakurta, Xu. **Practical and Private (Deep) Learning without Sampling or Shuffling**. ICML 2021.

Setting and Simplifications

The diagram shows the equation $\min_{\theta} [F(\theta) = \mathbb{E}_{x \sim P} [f(\theta; x)]]$ with three callout boxes. A box labeled 'Model parameters' points to the θ in the minimization. A box labeled 'Data' points to the x in the expectation. A box labeled 'Loss function' points to the $f(\theta; x)$ term.

$$\min_{\theta} [F(\theta) = \mathbb{E}_{x \sim P} [f(\theta; x)]]$$

Streaming setting: Suppose we draw a fresh data point $x_t \sim P$ in each iteration t (i.e. only 1 epoch)

Toeplitz noise correlations: $\beta_{t,\tau} = \beta_\tau$

$$\theta_{t+1} = \theta_t - \eta \left(g_t + \sum_{\tau=0}^t \beta_\tau z_{t-\tau} \right)$$

$$B = \begin{pmatrix} \beta_{0,0} & & & & \\ \beta_{0,1} & \beta_{1,0} & & & \\ \beta_{0,2} & \beta_{1,1} & \beta_{2,0} & \cdots & \\ \vdots & & & & \end{pmatrix} \longrightarrow B = \begin{pmatrix} \beta_0 & & & & \\ \beta_1 & \beta_0 & & & \\ \beta_2 & \beta_1 & \beta_0 & \cdots & \\ \vdots & & & & \end{pmatrix}$$

Computationally: store $O(T)$ coefficients instead of $O(T^2)$

Asymptotics: Iterates converge to a stationary distribution as $t \rightarrow \infty$

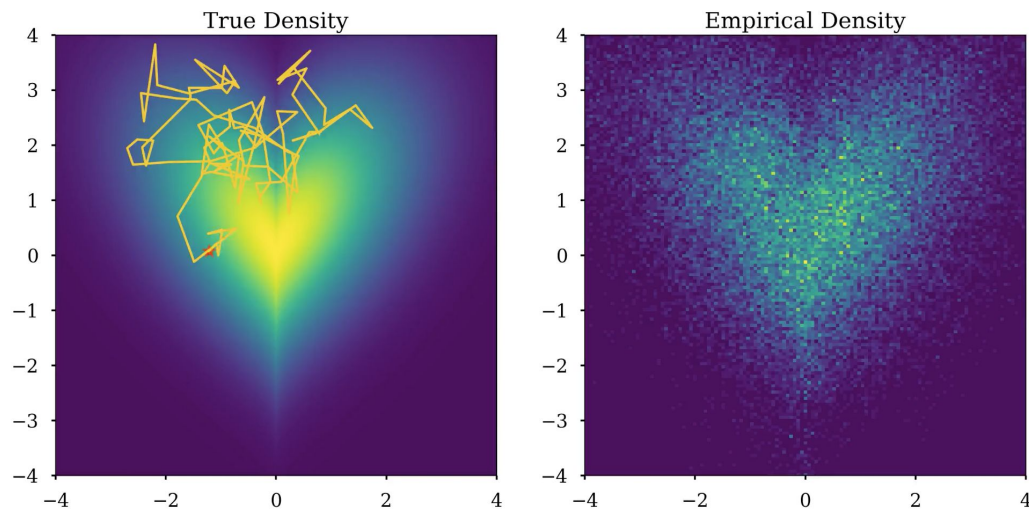


Image credit:
[Abdul Fatir Ansari](#)

Asymptotics: Iterates converge to a stationary distribution as $t \rightarrow \infty$

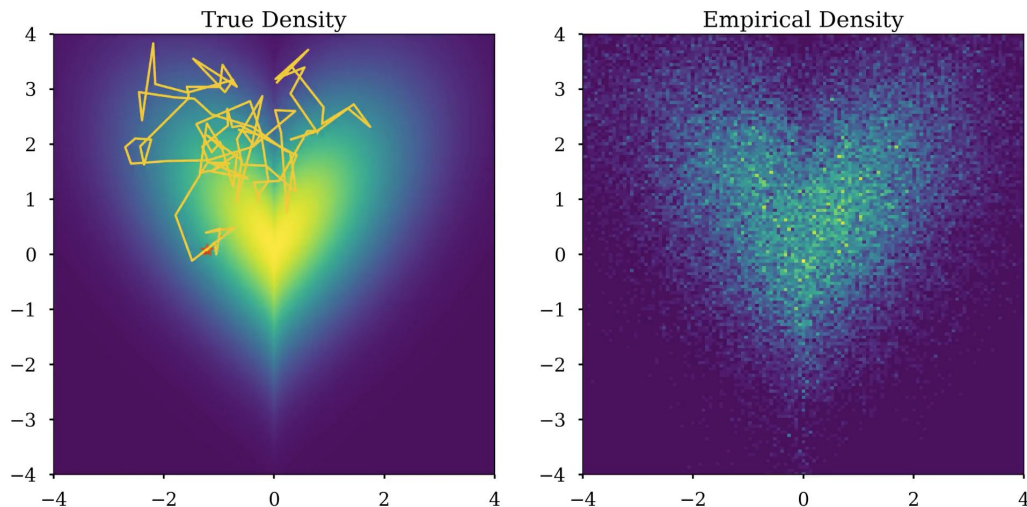


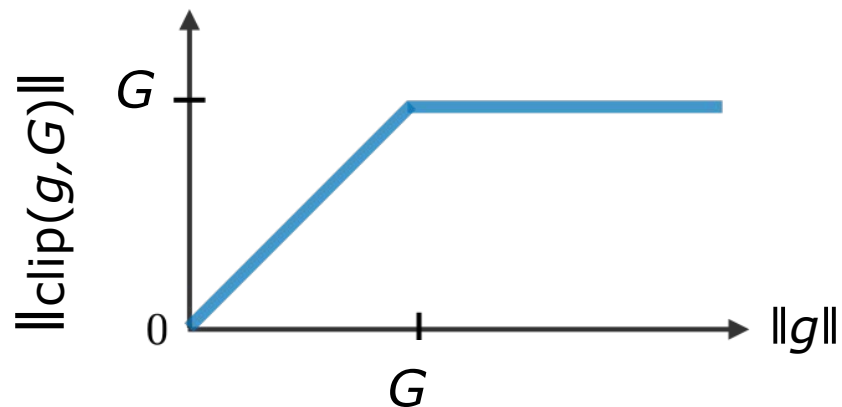
Image credit:
[Abdul Fatir Ansari](#)

**Asymptotic
error**

$$F_{\infty}(\beta) = \lim_{t \rightarrow \infty} \mathbb{E} [F(\theta_t) - F(\theta_{\star})]$$

Asymptotics at a fixed learning rate $\eta > 0$

Noisy-SGD/Noisy-FTRL: DP-SGD/DP-FTRL without clipping



Lets us study the noise dynamics of the algorithms
(do not satisfy DP guarantees)

Mean estimation in 1 dimension

$$\min_{\theta} [F(\theta) = \mathbb{E}_{x \sim P} (\theta - x)^2]$$

Data distribution
s.t. $|x| \leq 1$

Solve with stochastic optimization problem
with DP-SGD/DP-FTRL

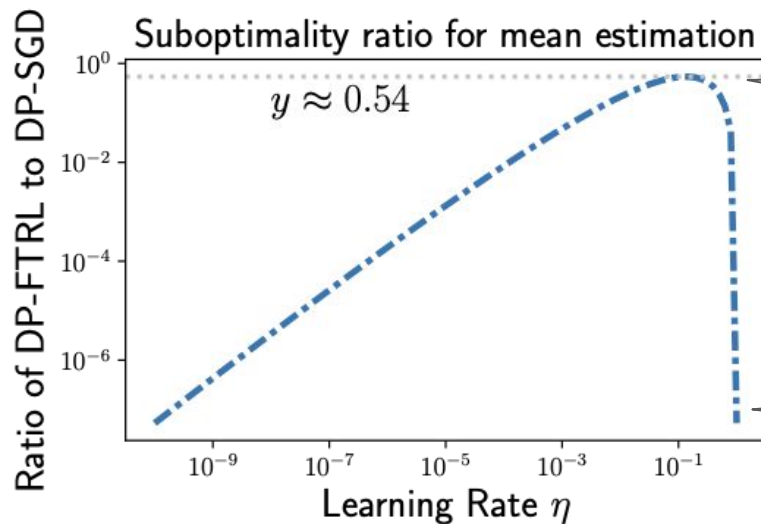
Mean estimation in 1 dimension

Informal Theorem: The asymptotic error of a ϱ -zCDP sequence is

Independent noise (DP-SGD)	$F_\infty(\beta^{\text{sgd}}) = \rho^{-1}\eta$
Correlated noise (DP-FTRL)	$\inf_{\beta} F_\infty(\beta) = F_\infty(\beta^\star) = \rho^{-1}\eta^2 \log^2 \frac{1}{\eta}$

η : learning rate (constant and non-zero)

ϱ : privacy level



DP-FTRL is always better than DP-SGD

DP-FTRL is significantly better at $\eta \rightarrow 0$ or $\eta \rightarrow 1$

Closed form correlations for mean estimation

Proposition: The correlations $\beta_0^* = 1$, $\beta_t^* = -t^{-3/2}(1 - \eta)^t$ attain the optimal error

$$\inf_{\beta} F_{\infty}(\beta) = F_{\infty}(\beta^*) = \rho^{-1} \eta^2 \log^2 \frac{1}{\eta}$$

Closed form correlations for mean estimation

Proposition: The correlations $\beta_0^* = 1$, $\beta_t^* = -t^{-3/2}(1 - \eta)^t$ attain the optimal error

$$\inf_{\beta} F_{\infty}(\beta) = F_{\infty}(\beta^*) = \rho^{-1} \eta^2 \log^2 \frac{1}{\eta}$$

ν -DP-FTRL

For general problems, use $\beta_0 = 1$, $\beta_t = -t^{-3/2}(1 - \nu)^t$

and tune the parameter ν

Linear regression

$$\min_{\theta} [F(\theta) = \mathbb{E}(y - \langle \theta, x \rangle)^2]$$

where $x \sim \mathcal{N}(0, H)$

H is also the
Hessian of the
objective

Linear regression

$$\min_{\theta} [F(\theta) = \mathbb{E}(y - \langle \theta, x \rangle)^2]$$

where $x \sim \mathcal{N}(0, H)$

Well-specified
linear model

$$y|x \sim \mathcal{N}(x^\top \theta_*, \sigma^2)$$

Informal Theorem: The asymptotic error is

Independent noise (Noisy-SGD)	=	$d \rho^{-1} \eta$
Correlated noise (v -Noisy-FTRL)	\leq	$d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$
Lower bound for any algorithm	\geq	$d_{\text{eff}} \rho^{-1} \eta^2$

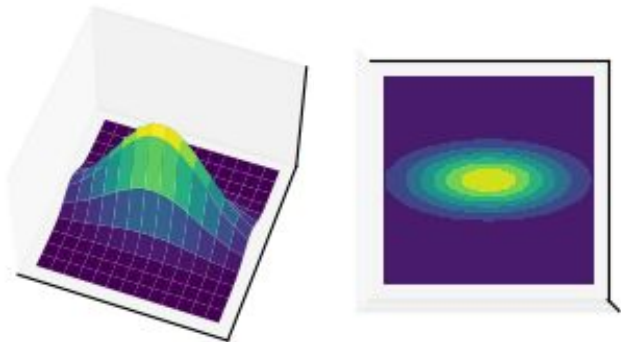
Improve **dimension d** to
problem-dependent
effective dimension d_{eff}

Effective dimension

$$d_{\text{eff}} = \text{Tr}(H) / \|H\|_2 \leq d$$

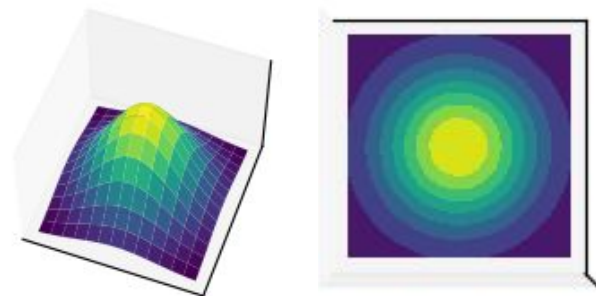
Low effective dimension

$$\lambda_1 = 1, \lambda_2 = \dots = \lambda_d = 1/d$$



High effective dimension

$$\lambda_1 = \lambda_2 = \dots = \lambda_d = 1$$



Closely connected to **numerical/stable rank**

SAMPLING FROM LARGE MATRICES: AN APPROACH THROUGH GEOMETRIC FUNCTIONAL ANALYSIS

MARK RUDELSON AND ROMAN VERSHYNIN

Remark 1.3 (Numerical rank). The numerical rank $r = r(A) = \|A\|_F^2 / \|A\|_2^2$ in Theorem 1.1 is a relaxation of the exact notion of rank. Indeed, one always has $r(A) \leq \text{rank}(A)$. But as opposed to the exact rank, the numerical rank is stable under small perturbations of the matrix A . In particular, the numerical rank of A tends to be low when A is close to a low rank matrix, or when A is sufficiently sparse.

$$d_{\text{eff}} = \text{srank}(H^{1/2})$$

The stable rank appears in:

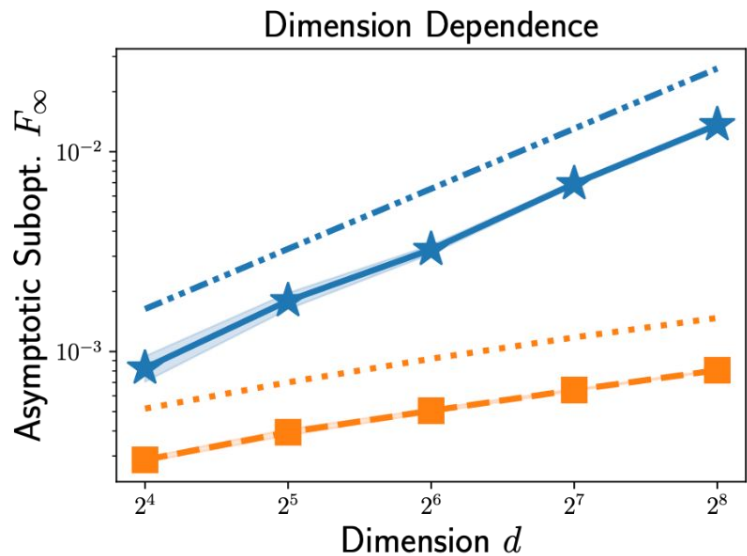
- Numerical linear algebra (e.g. randomized matrix multiplications) [Tropp (2014), Cohen-Nelson-Woodruff (2015)]
- Matrix concentration [Hsu-Kakade-Zhang (2012), Minsker (2017)]
- ...

Informal Theorem: The asymptotic error is

Independent noise (Noisy-SGD)	=	$d \rho^{-1} \eta$
Correlated noise (v -Noisy-FTRL)	\leq	$d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$
Lower bound for any algorithm	\geq	$d_{\text{eff}} \rho^{-1} \eta^2$

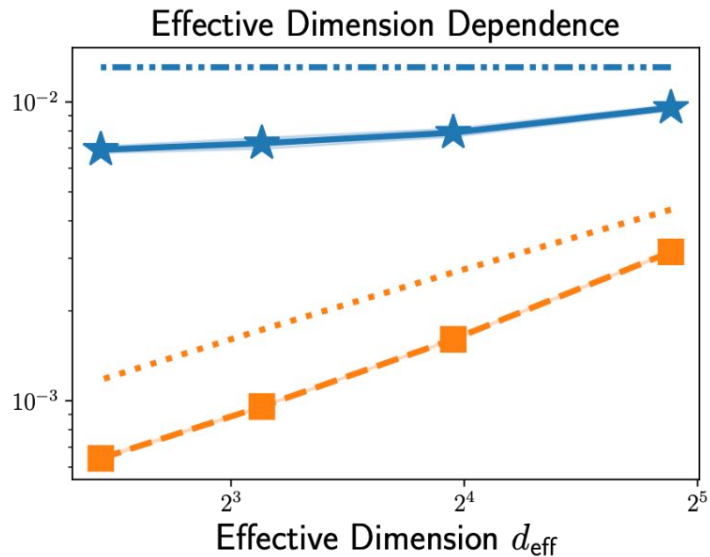
Improve **dimension d** to
problem-dependent
effective dimension d_{eff}

Linear regression: theory predicts simulations



Noisy-SGD
scales with d

Noisy-FTRL
scales with d_{eff}

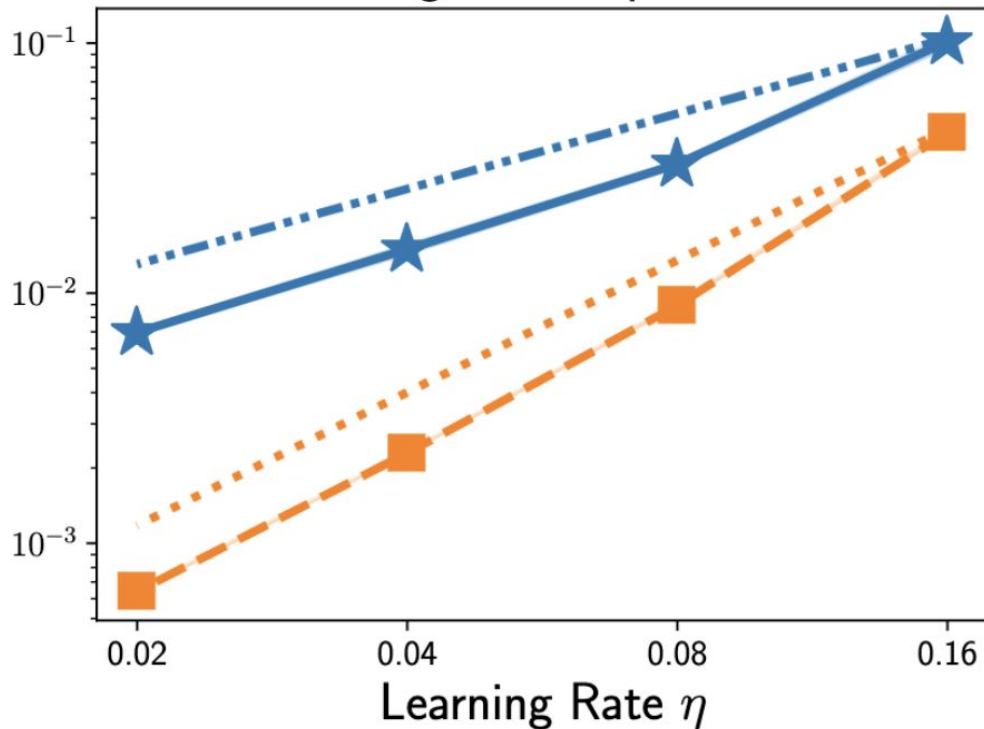


Informal Theorem: The asymptotic error for $0 < \eta < 1$ is

Independent noise (Noisy-SGD)	$= d \rho^{-1} \eta$
Correlated noise (v -Noisy-FTRL)	$\leq d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$
Lower bound for any algorithm	$\geq d_{\text{eff}} \rho^{-1} \eta^2$

*Improved dependence on
the learning rate η*

Learning Rate Dependence



Noisy-SGD scales as η

ν -Noisy-FTRL
scales as η^2

Noisy-FTRL \gg **Noisy-SGD** at small η

Finite-time rates with DP: Linear Regression

Independent noise (DP-SGD)	$\frac{1}{\rho T} + \frac{1}{T}$
Correlated noise (ν -DP-FTRL)	$\frac{1}{\rho T^2} + \frac{1}{T}$

Privacy error

T : number of iterations

ρ : privacy level

η : learning rate is optimized

Proof sketch for Mean Estimation

Updates are not Markovian (key for all stochastic gradient proofs)

Our approach: Analysis the Fourier domain

Letting $\delta_t = \theta_t - \theta_*$, the DP-FTRL update can be written as

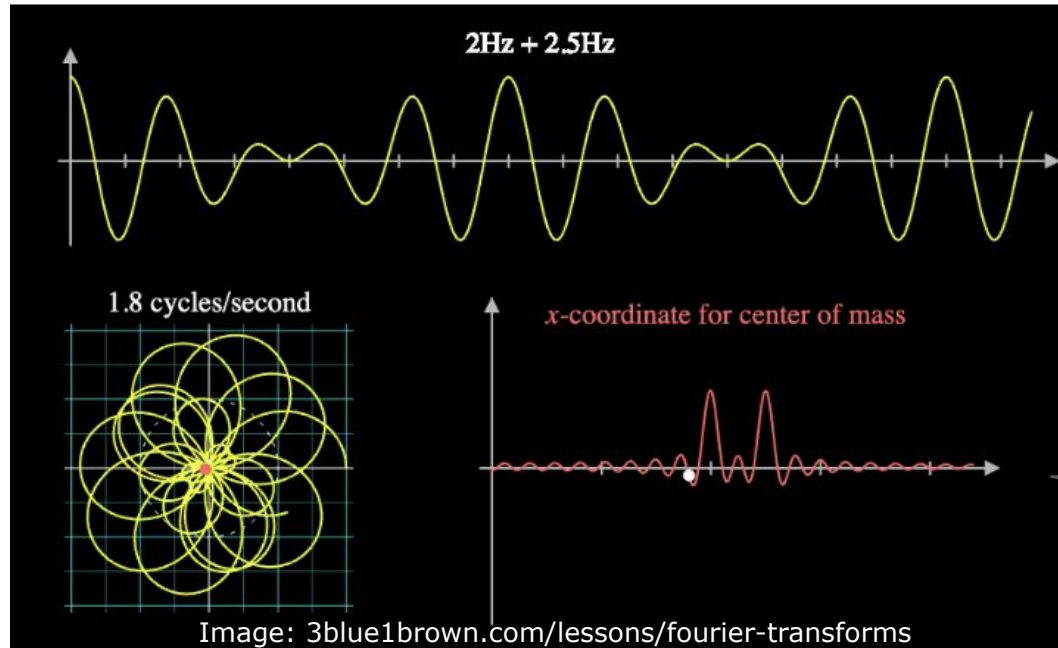
Linear
Time-Invariant
(LTI) system

$$\delta_{t+1} = (1 - \eta)\delta_t - \eta \sum_{\tau=0}^t \beta_{\tau} z_{t-\tau}$$

Convolution of the
noise

Fourier analysis can give the stationary variance of δ_t in terms of the **discrete-time Fourier transform** $B(\omega) = \sum_{t=0}^{\infty} \beta_t e^{i\omega t}$ of the convolution weights β

Frequency



Time-domain description

Frequency-domain description

Letting $\delta_t = \theta_t - \theta_*$, the DP-FTRL update can be written as

Linear
Time-Invariant
(LTI) system

$$\delta_{t+1} = (1 - \eta)\delta_t - \eta \sum_{\tau=0}^t \beta_{\tau} z_{t-\tau}$$

Convolution of the
noise

The stationary variance of δ_t can be given as

$$\lim_{t \rightarrow \infty} \mathbb{E}[\delta_t^2] = \frac{\eta^2}{2\pi} \left(\int_{-\pi}^{\pi} \frac{|B(\omega)|^2}{|1 - \eta e^{i\omega}|^2} d\omega \right) \mathbb{E}[z_t^2]$$

$$\lim_{t \rightarrow \infty} \mathbb{E}[\delta_t^2] = \frac{\eta^2}{2\pi} \left(\int_{-\pi}^{\pi} \frac{|B(\omega)|^2}{|1 - \eta - e^{i\omega}|^2} d\omega \right) \mathbb{E}[z_t^2]$$

sensitivity

For ρ -zCDP, take

$$\begin{aligned} \mathbb{E}[z_t^2] &= \frac{1}{2\rho} \max_t \|[B^{-1}]_{:,t}\|_2^2 \\ &= \frac{1}{2\rho} \int_{-\pi}^{\pi} \frac{d\omega}{2\pi|B(\omega)|^2} \end{aligned}$$

$$B = \begin{pmatrix} \beta_0 & & & \\ \beta_1 & \beta_0 & & \\ \beta_2 & \beta_1 & \beta_0 & \cdots \\ \vdots & & & \end{pmatrix}$$

$$\lim_{t \rightarrow \infty} \mathbb{E}[\delta_t^2] = \frac{\eta^2}{2\pi} \left(\int_{-\pi}^{\pi} \frac{|B(\omega)|^2}{|1 - \eta - e^{i\omega}|^2} d\omega \right) \mathbb{E}[z_t^2]$$

Requires $|B(\omega)|$
small

For ρ -zCDP, take

$$\mathbb{E}[z_t^2] = \frac{1}{2\rho} \max_t \|[B^{-1}]_{:,t}\|_2^2$$

$$= \frac{1}{2\rho} \int_{-\pi}^{\pi} \frac{d\omega}{2\pi|B(\omega)|^2}$$

sensitivity

$$B = \begin{pmatrix} \beta_0 & & & \\ \beta_1 & \beta_0 & & \\ \beta_2 & \beta_1 & \beta_0 & \cdots \\ \vdots & & & \end{pmatrix}$$

Requires $|B(\omega)|$
large

$$\lim_{t \rightarrow \infty} \mathbb{E}[\delta_t^2] = \frac{\eta^2}{2\pi} \left(\int_{-\pi}^{\pi} \frac{|B(\omega)|^2}{|1 - \eta - e^{i\omega}|^2} d\omega \right) \mathbb{E}[z_t^2]$$

Requires $|B(\omega)|$
small

For ρ -zCDP, take $\mathbb{E}[z_t^2] = \frac{1}{2\rho} \max_t \|[B^{-1}]_{:,t}\|_2^2$

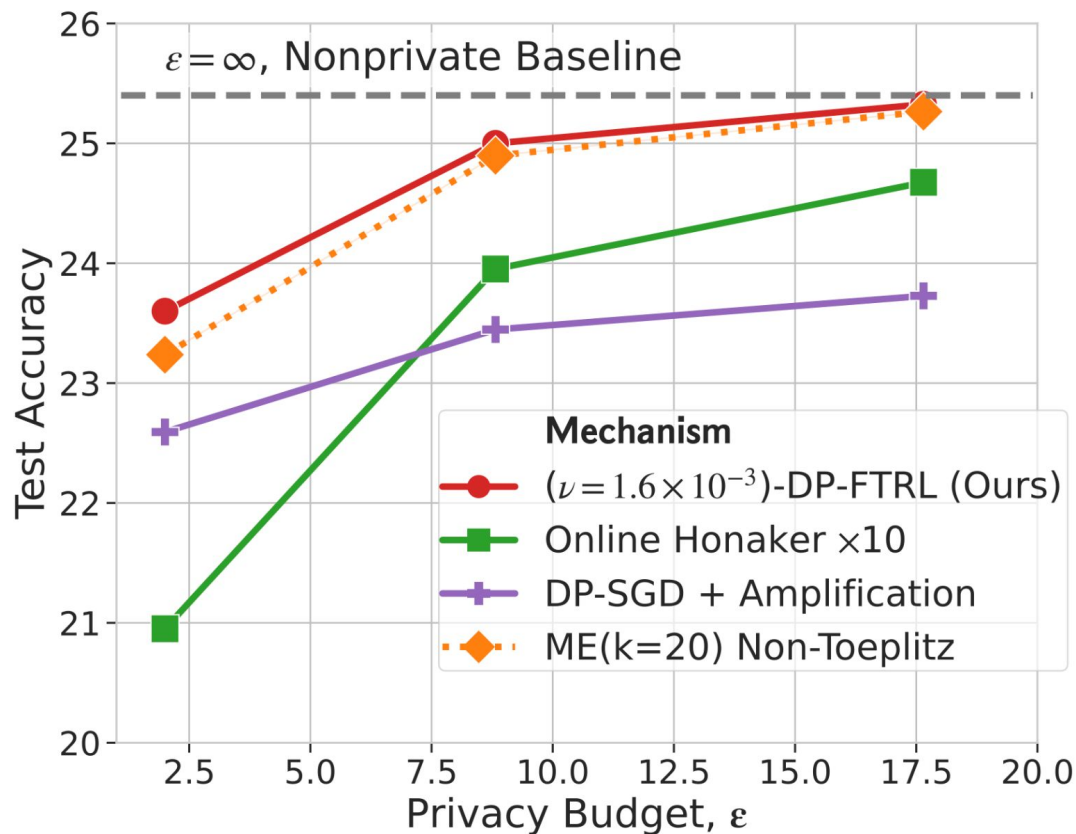
$$= \frac{1}{2\rho} \int_{-\pi}^{\pi} \frac{d\omega}{2\pi |B(\omega)|^2}$$

$$B = \begin{pmatrix} \beta_0 & & & \\ \beta_1 & \beta_0 & & \\ \beta_2 & \beta_1 & \beta_0 & \dots \\ \vdots & & & \end{pmatrix}$$

Requires $|B(\omega)|$
large

Optimizing for $|B(\omega)|$ gives the theorem

Language modeling with Stack Overflow | User-level DP



Ours
matches
SoTA!

Image classification with CIFAR-10 | Example-level DP

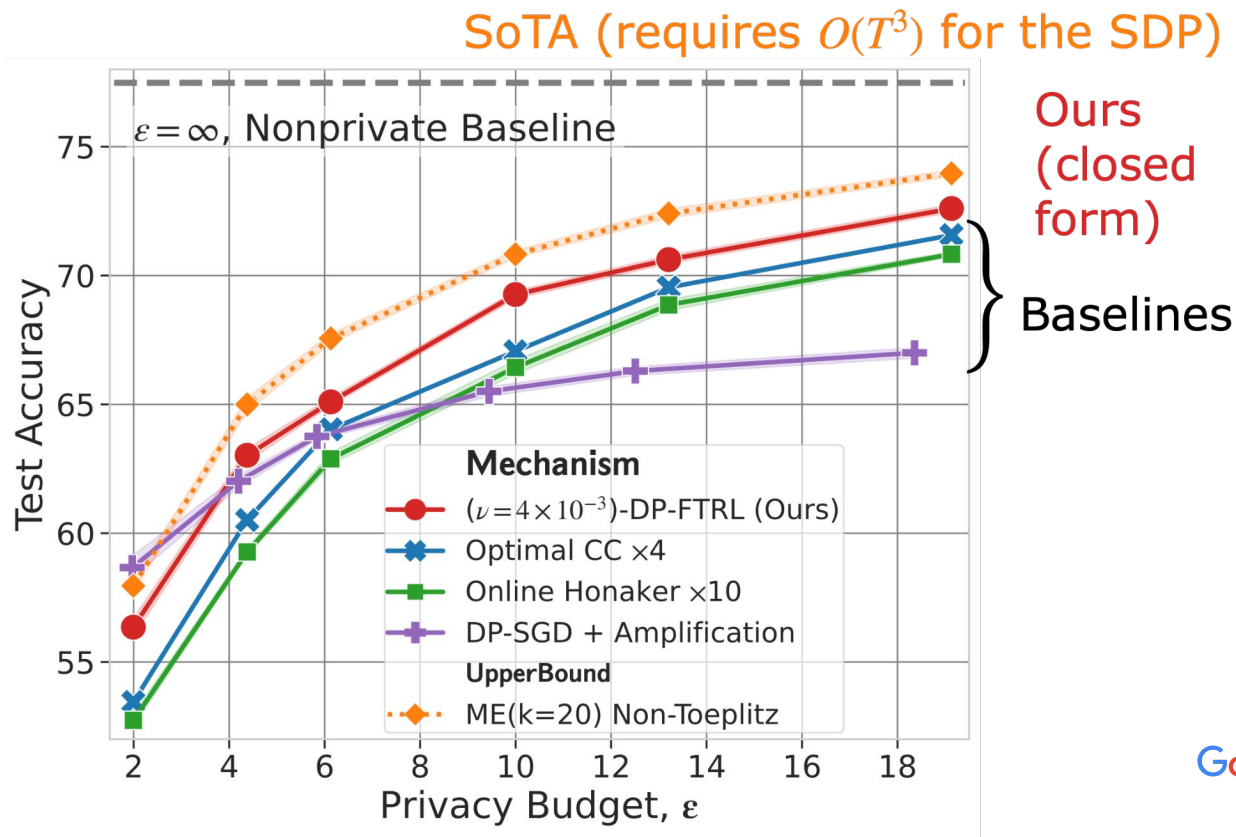
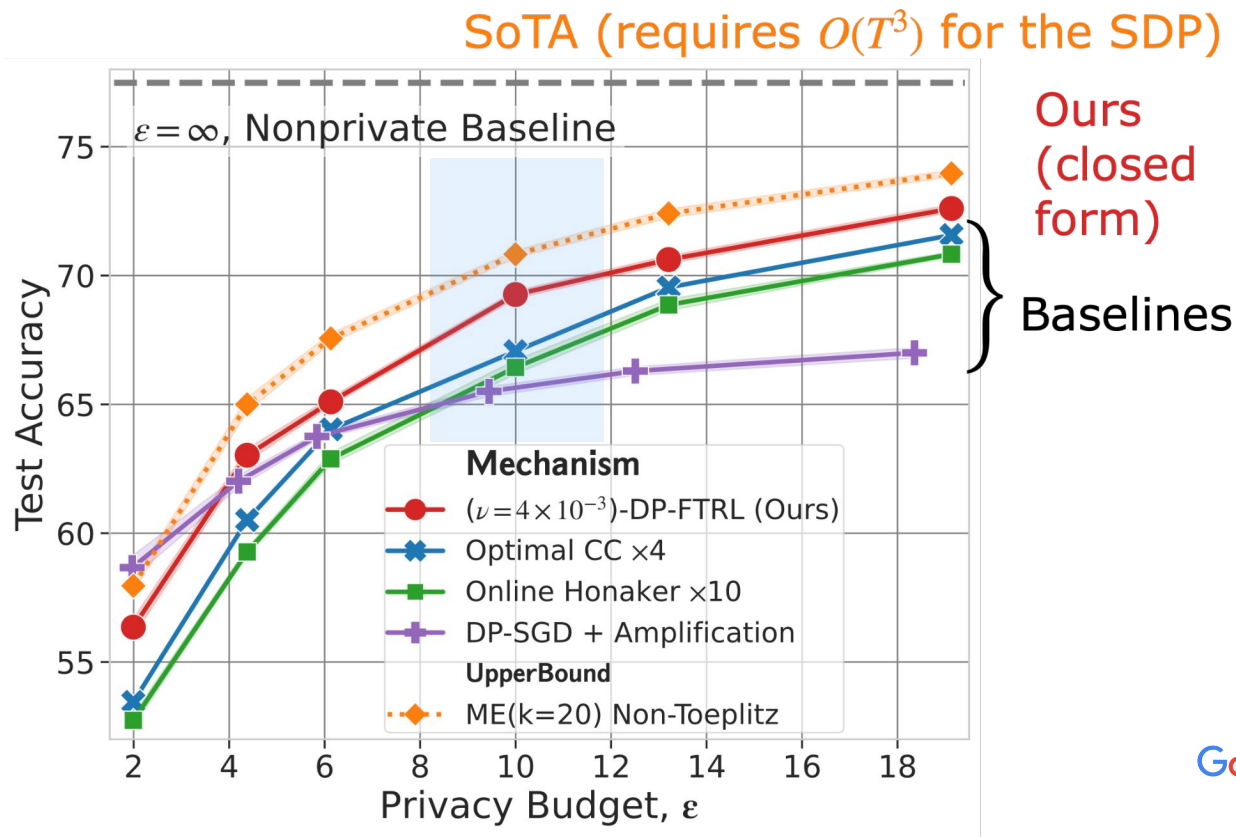


Image classification with CIFAR-10 | Example-level DP



Computational cost

- **SoTA**: cubic complexity to generate the β 's
- **Ours**: linear complexity (closed form)
 - nearly matches SoTA empirically

Summary

- Correlated noise is **provably** better
- Depends on effective dimension instead of dimension
- Matches lower bounds

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

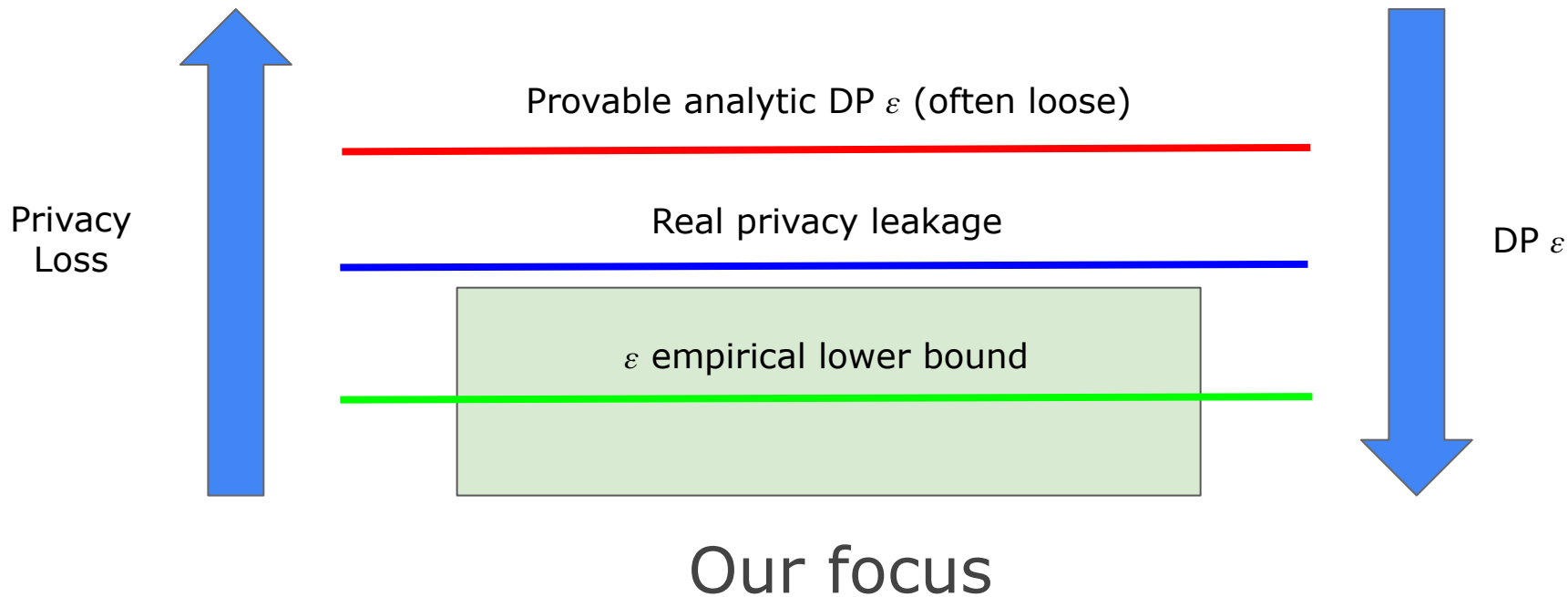


Alina Oprea



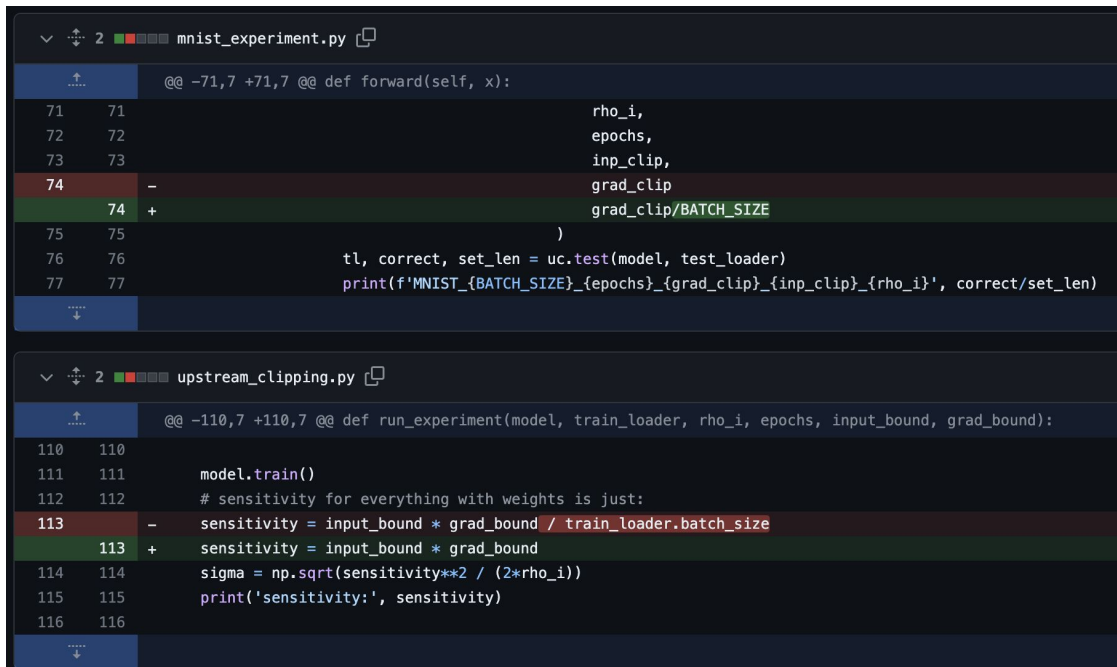
Sewoong Oh

Empirical privacy auditing



Why empirical privacy auditing?

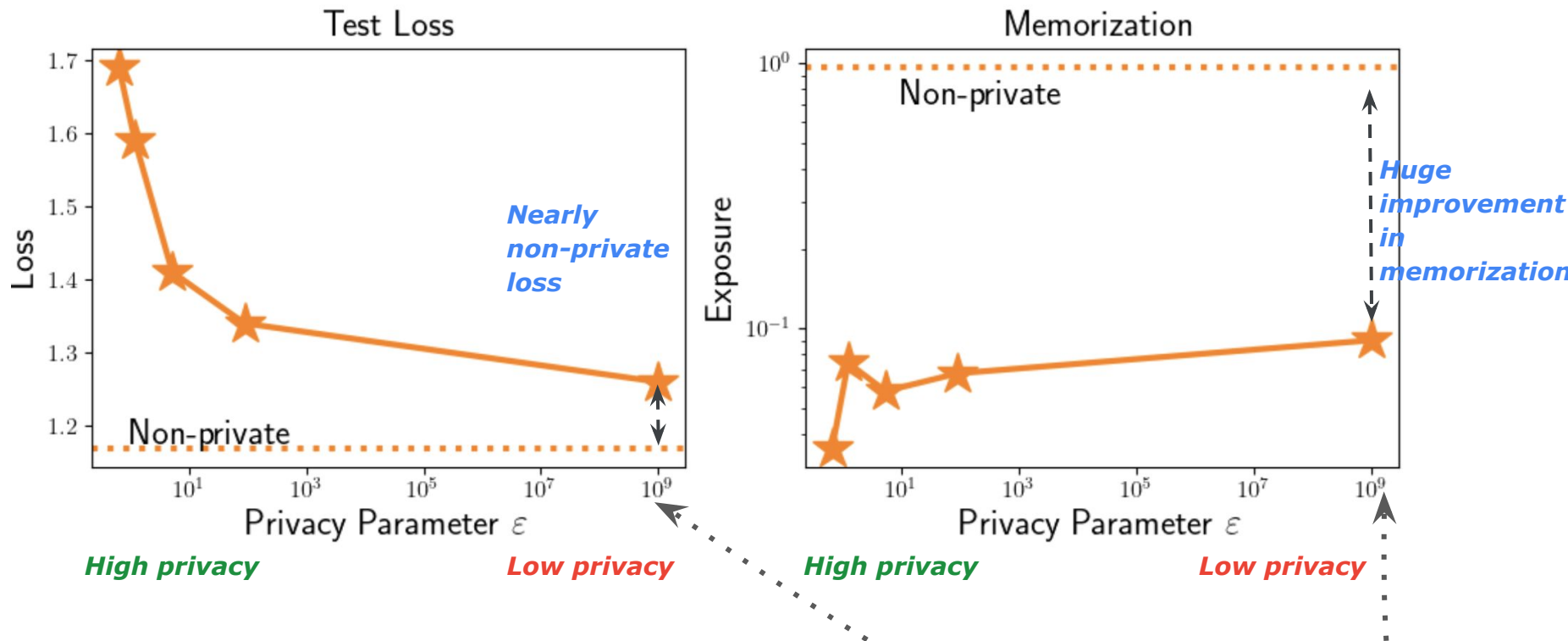
To verify that we actually provide the guarantee we claim
(no bugs in proofs/implementation)



```
mnist_experiment.py
@@ -71,7 +71,7 @@ def forward(self, x):
71 71         rho_i,
72 72         epochs,
73 73         inp_clip,
74 74         grad_clip
74 74 +         grad_clip/BATCH_SIZE
75 75     )
76 76     tl, correct, set_len = uc.test(model, test_loader)
77 77     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     model.train()
111 111     # sensitivity for everything with weights is just:
112 112
113 113 -     sensitivity = input_bound * grad_bound / train_loader.batch_size
113 113 +     sensitivity = input_bound * grad_bound
114 114     sigma = np.sqrt(sensitivity**2 / (2*rho_i))
115 115     print('sensitivity:', sensitivity)
116 116
```

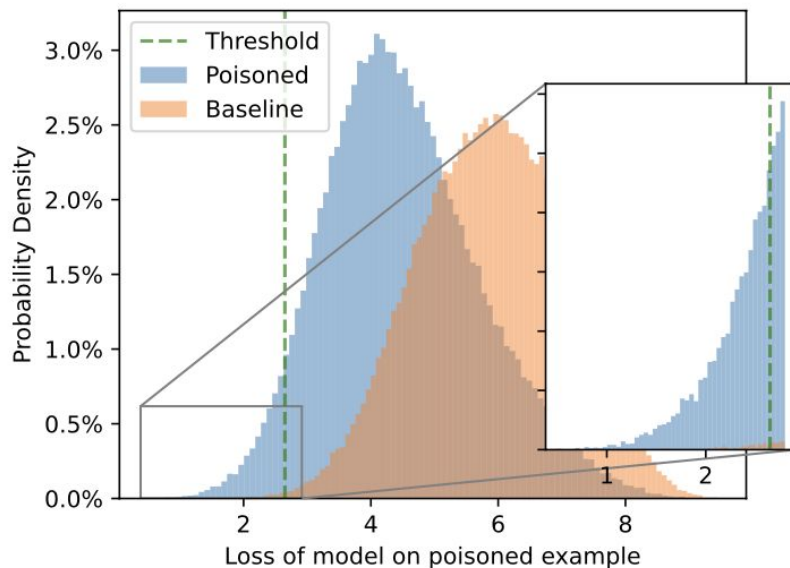
Gap between DP guarantees and empirical behavior: Memorization



Privacy guarantee is vacuous at this ϵ !

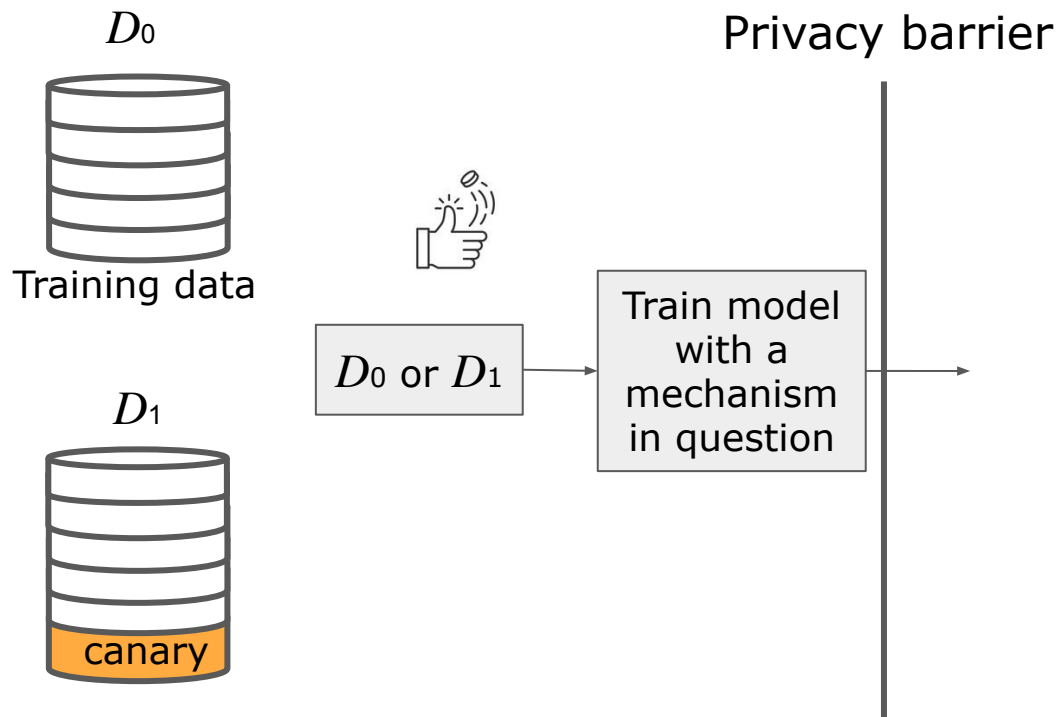
Empirical Privacy Auditing requires **many samples**

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

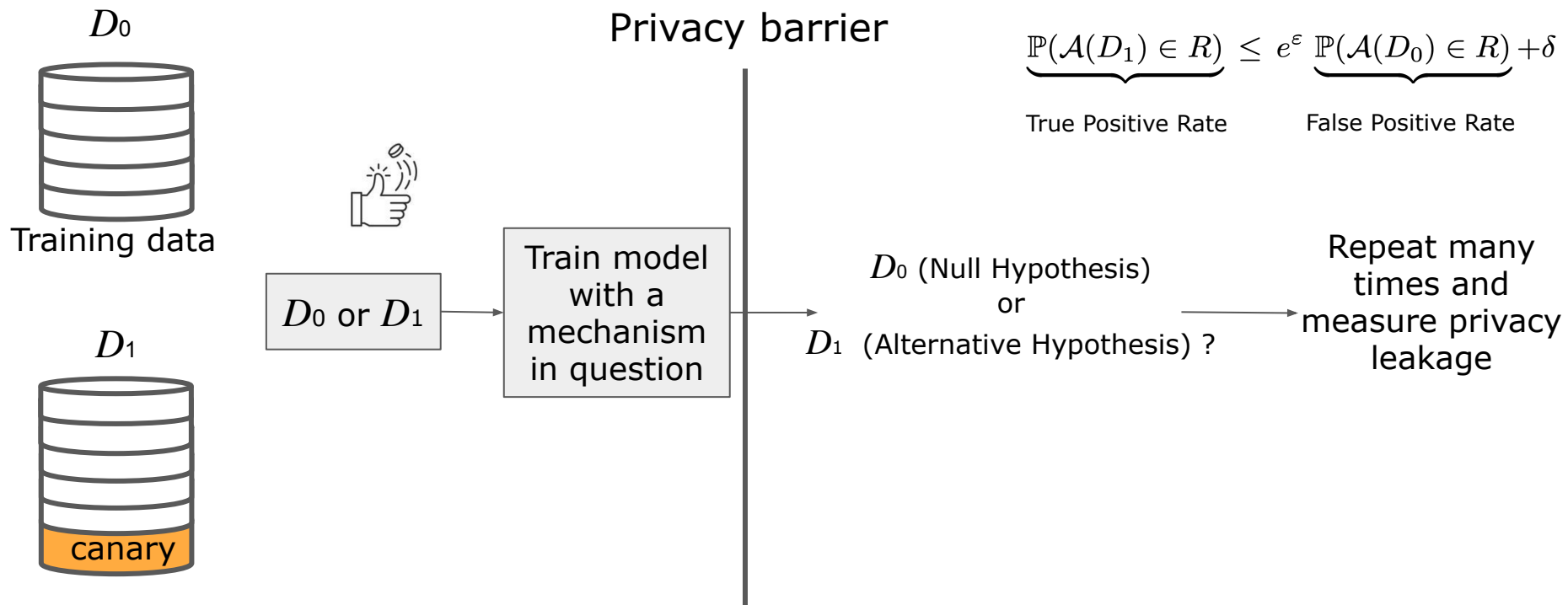


Our goal: make empirical privacy auditing
more *sample-efficient*

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

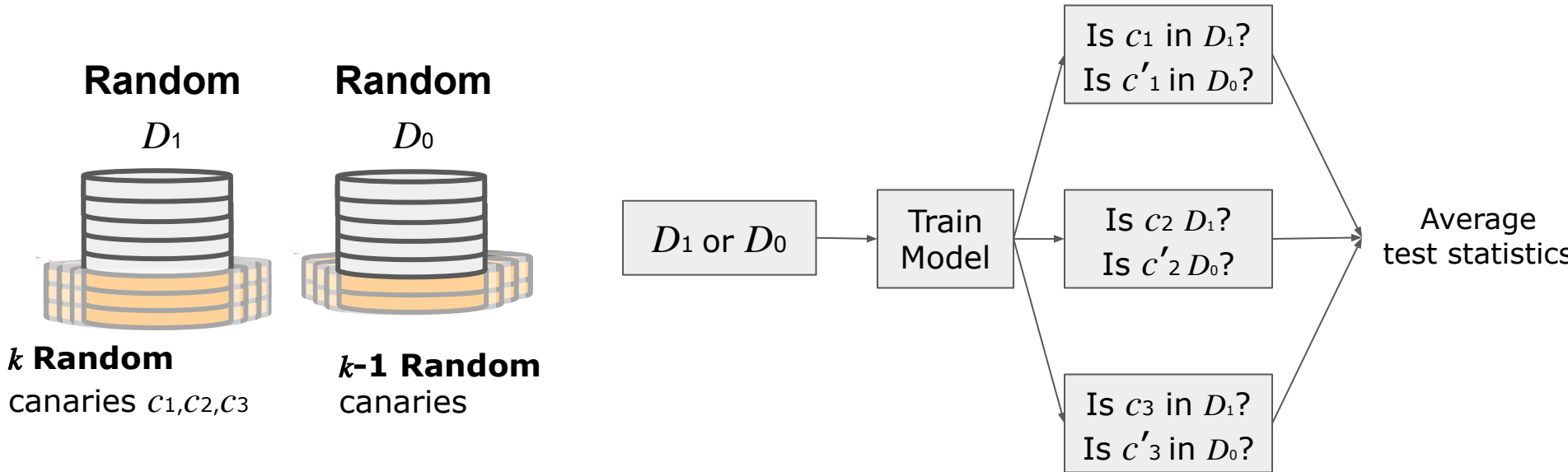
$$\begin{aligned} & \text{Actual TPR/FPR} \\ & \uparrow \\ \varepsilon & \geq \log \left(\frac{\text{TPR} - \delta}{\text{FPR}} \right) \\ & \geq \log \left(\frac{\hat{\text{TPR}}_n - \frac{c}{\sqrt{n}} - \delta}{\hat{\text{FPR}}_n + \frac{c}{\sqrt{n}}} \right) \\ & \downarrow \\ & \text{Empirical TPR/FPR}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

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)

Average test statistics

Is c_3 in D_1 ?
Is c'_3 in D_0 ?

k Random canaries c_1, c_2, c_3

$k-1$ Random canaries

Novel higher-order confidence interval

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

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

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

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

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

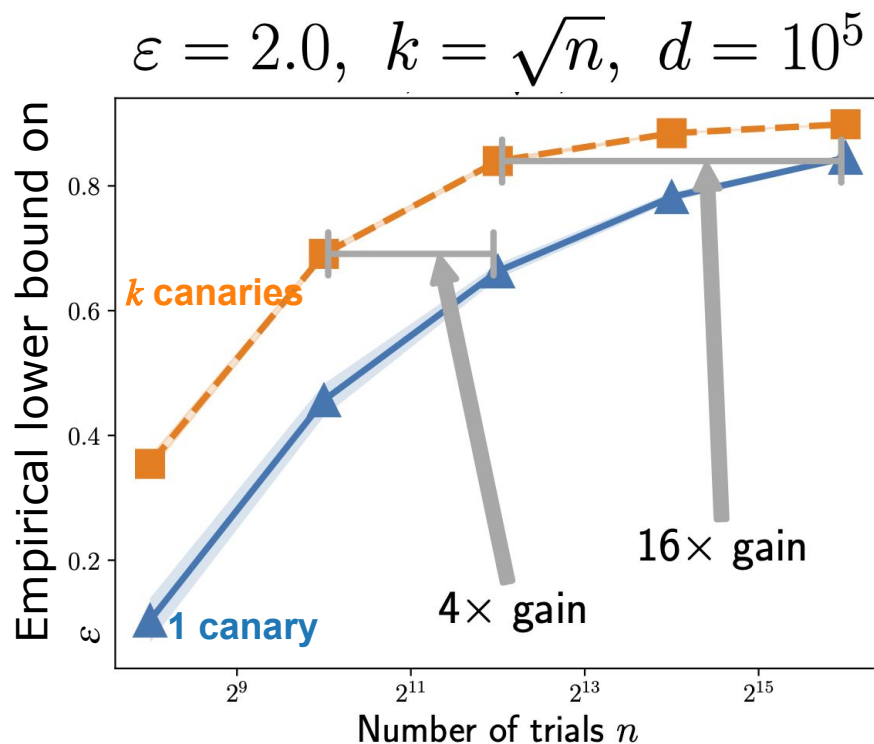
c' - Data-dependent constant

**Lower variance =>
Tighter confidence intervals**

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

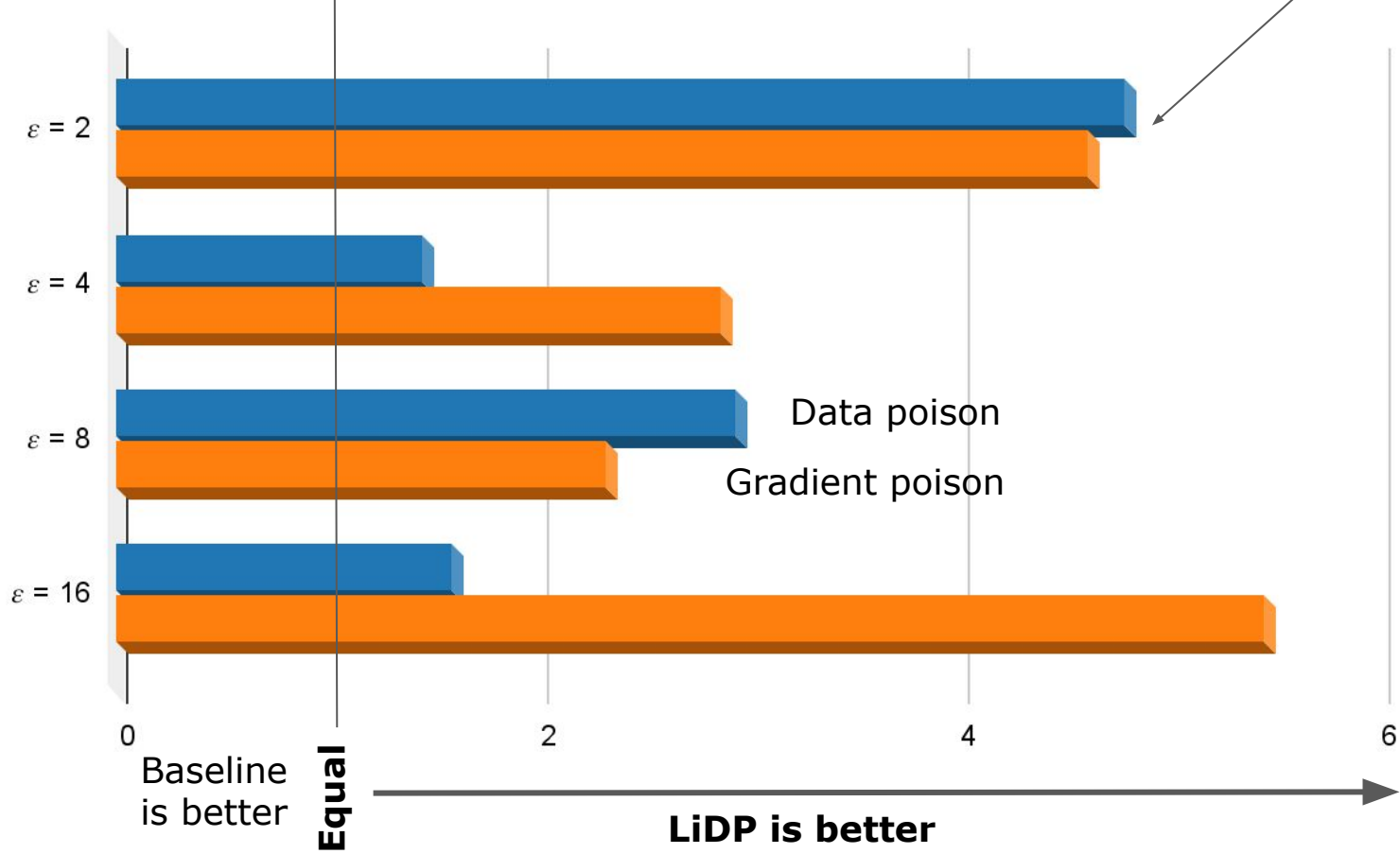
Proof of concept with Gaussian mechanisms

- 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

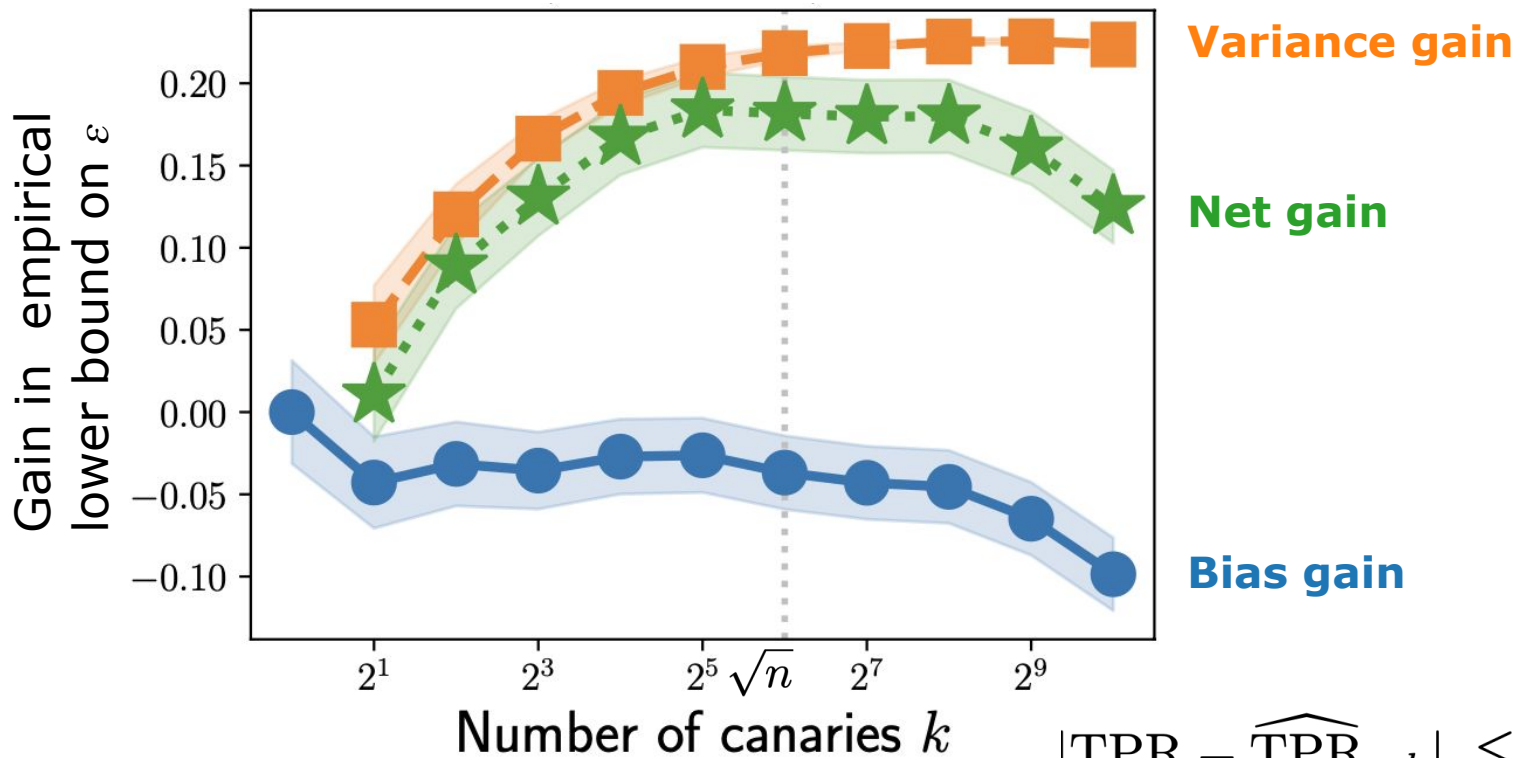
Thomas Steinke*
Google DeepMind
steinke@google.com

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

Other highlights: large-scale group-stratified datasets

Dataset Grouper

Library for creating group-structured datasets.

- **Scalable:** can handle millions of clients ✓
- **Flexible:** any custom partition function on any TFDS/HuggingFace dataset ✓
- **Platform-agnostic:** works with TF, PyTorch, JAX, NumPy, ... ✓

Zach Charles*, Nicole Mitchell*, *KP**,
Michael Reneer, Zach Garrett.
NeurIPS D&B 2023



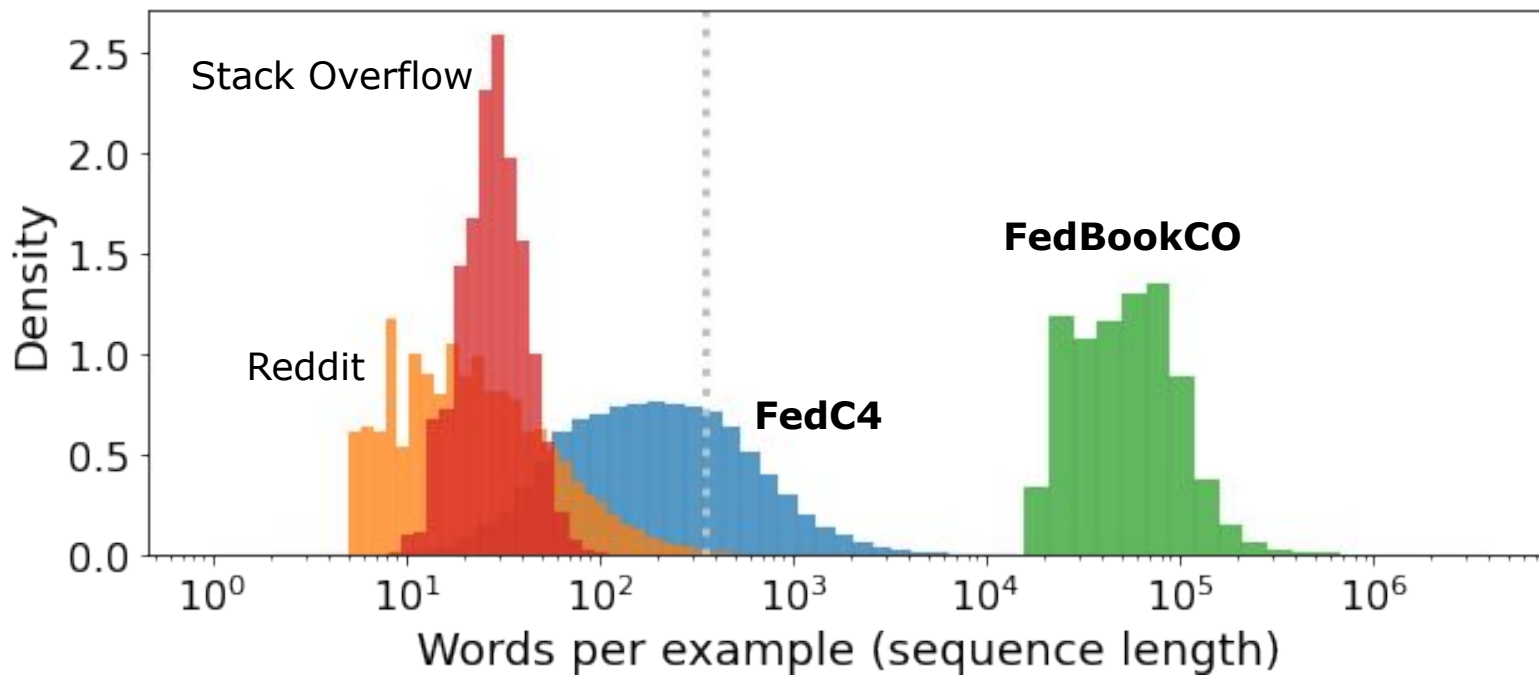
New federated LLM datasets: longer sequences

Largest previous datasets:

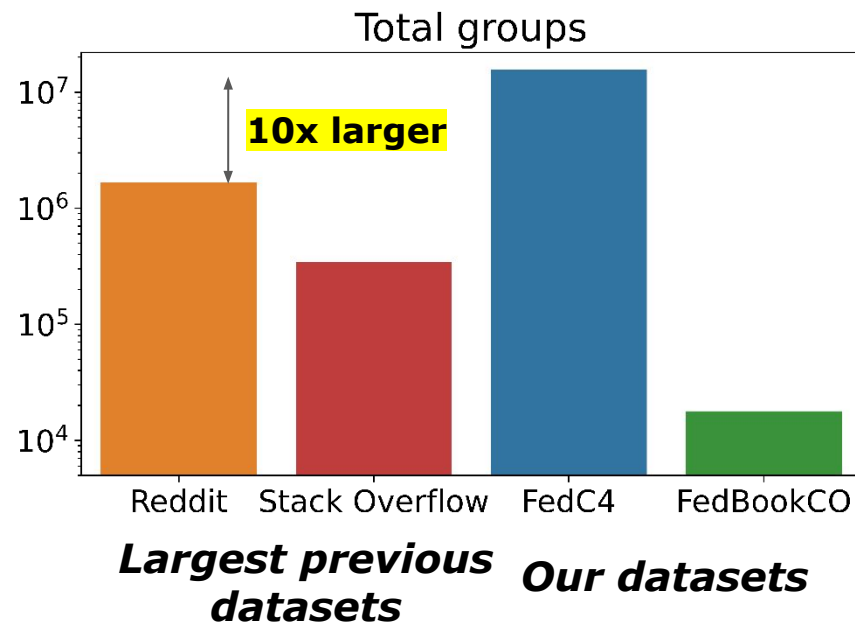
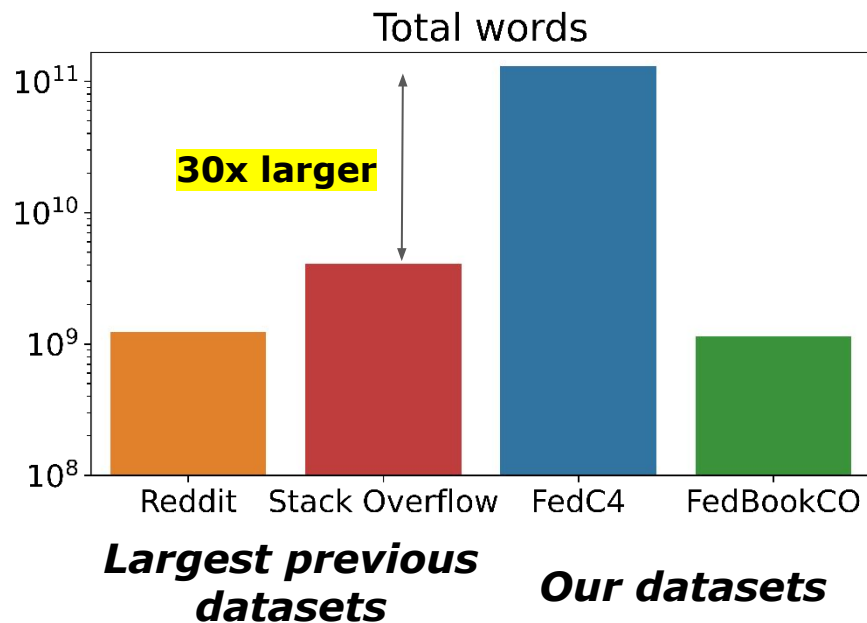
Reddit, Stack Overflow

*Typical
sequence
length of LLMs*

Our datasets:
FedC4, FedBookCO



New federated LLM datasets: more words & groups



Thank you!

