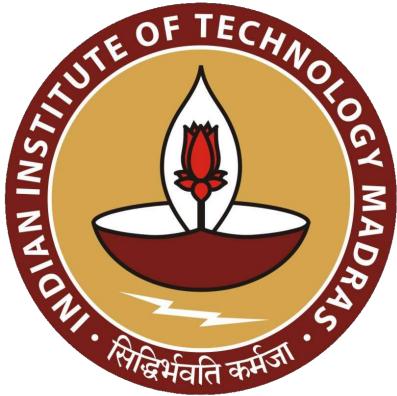


Near-Optimal Private Learning with Correlated Noise Mechanisms

Krishna Pillutla

Jan. 12 2025 @ CMI (BIRS Workshop)



Joint work with



Chris
Choquette-Choo



Dj
Dvijotham



Brendan
McMahan



Arun
Ganesh



Thomas
Steinke



Abhradeep
Thakurta

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

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

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

Dvijotham, McMahan, **P.**, Steinke, Thakurta.

Efficient and Near-Optimal Noise Generation for Streaming Differential Privacy.
FOCS (2024).

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

AHA, FOUND THEM!



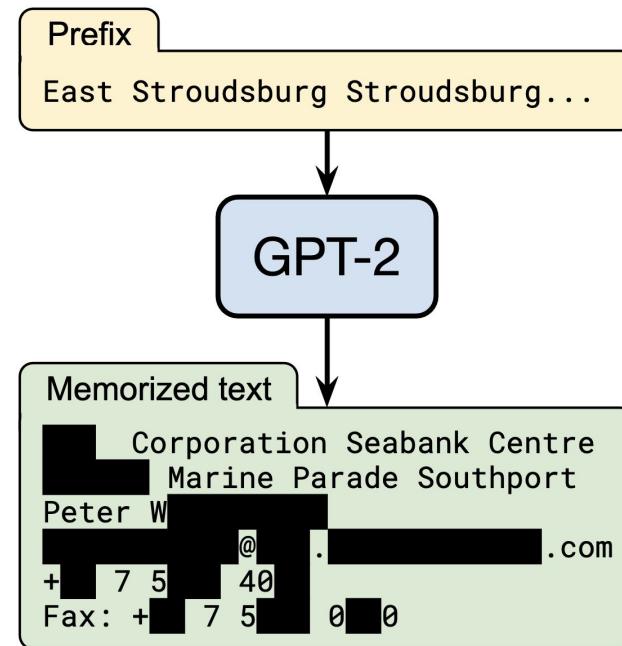
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

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



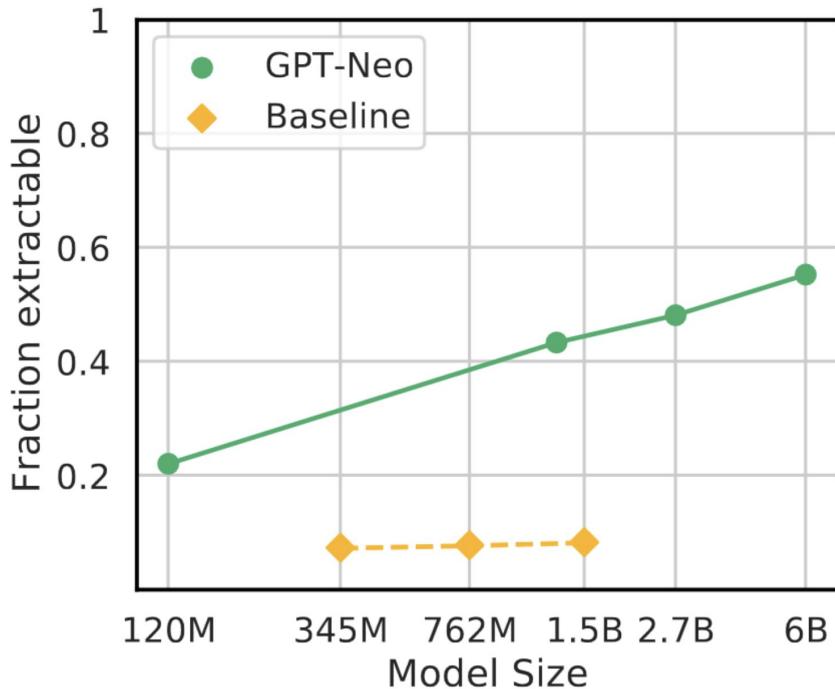
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data

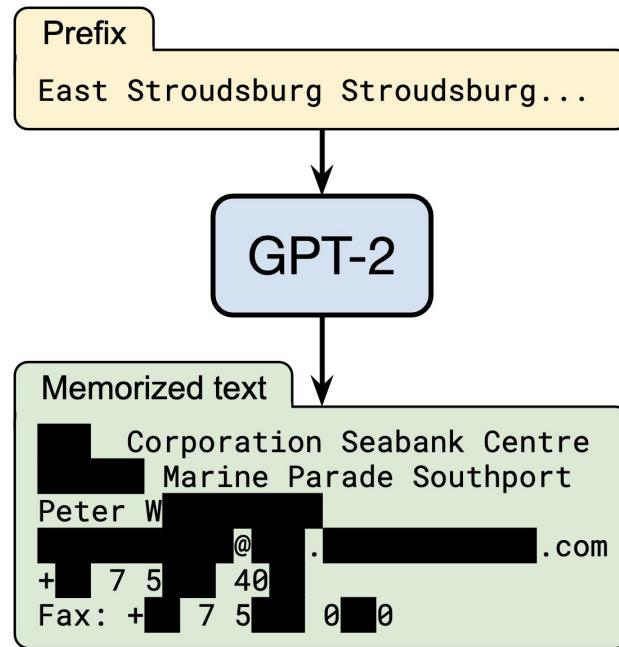


Carlini et al. (USENIX Security 2021)

Models leak information about their training data ***reliably***



Carlini et al. (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)

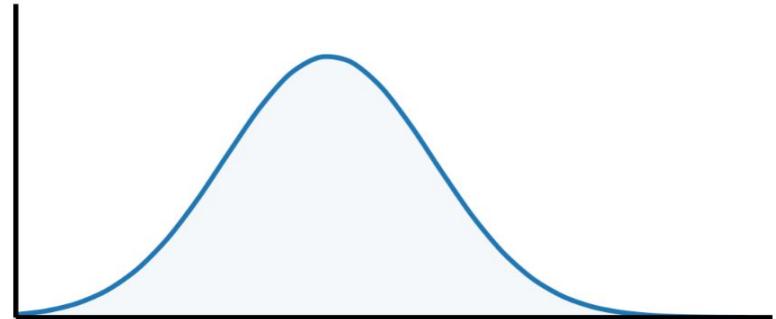
Dataset



Randomized
Algorithm

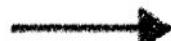
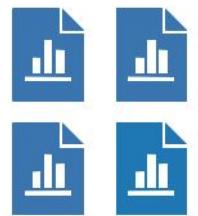


Output Distribution
(e.g. over models)



Differential privacy (DP)

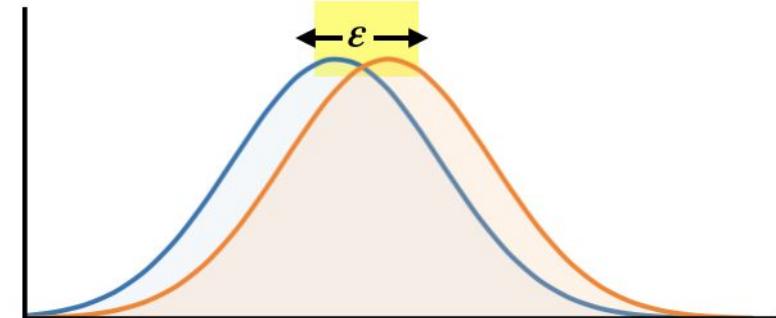
Dataset



Randomized
Algorithm



Output Distribution
(e.g. over models)

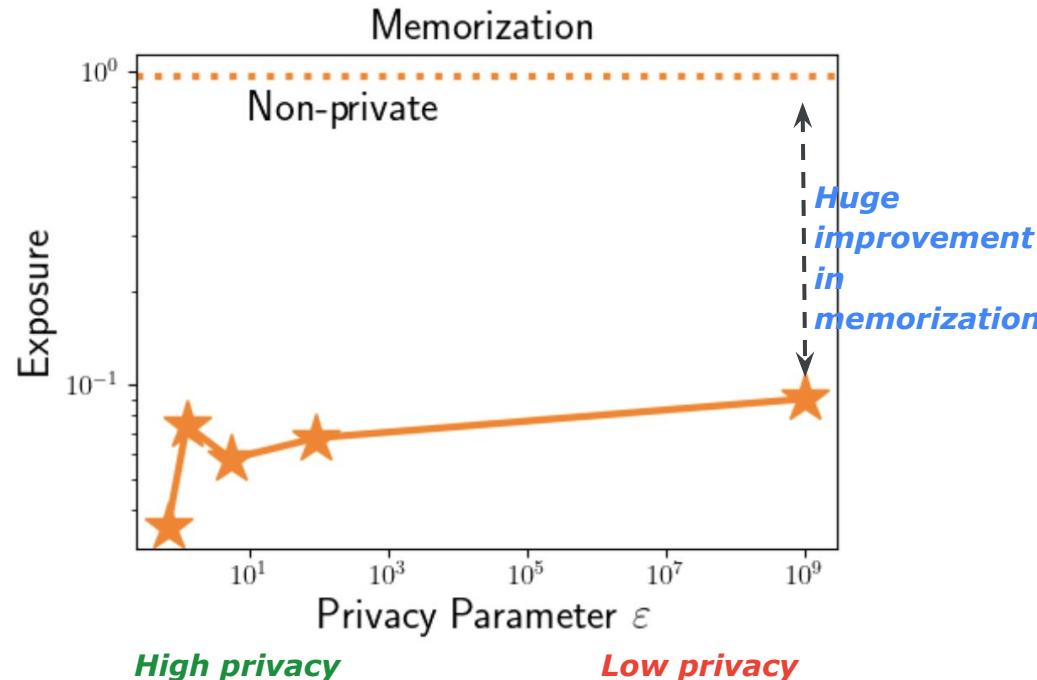
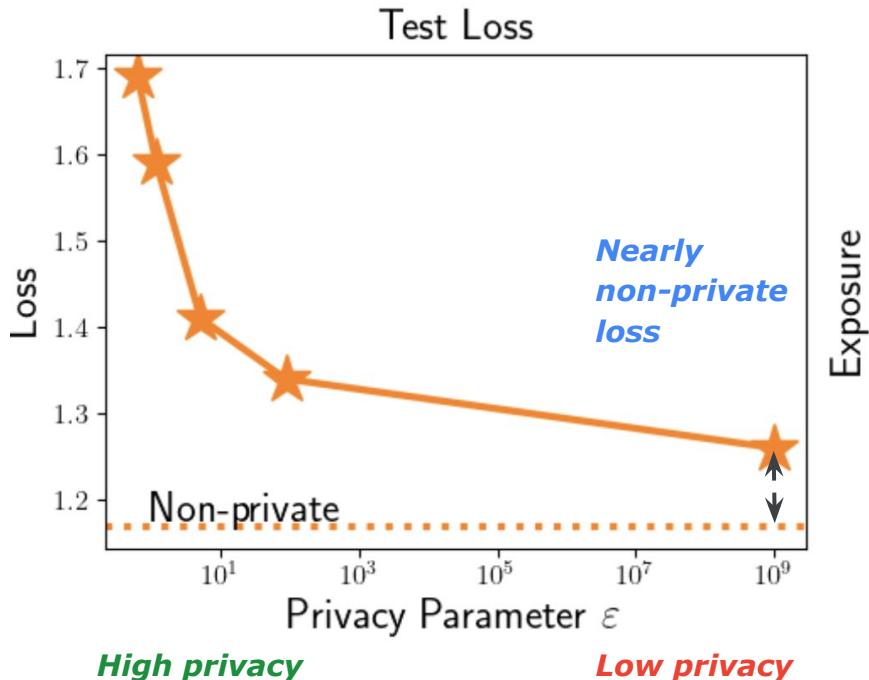


+

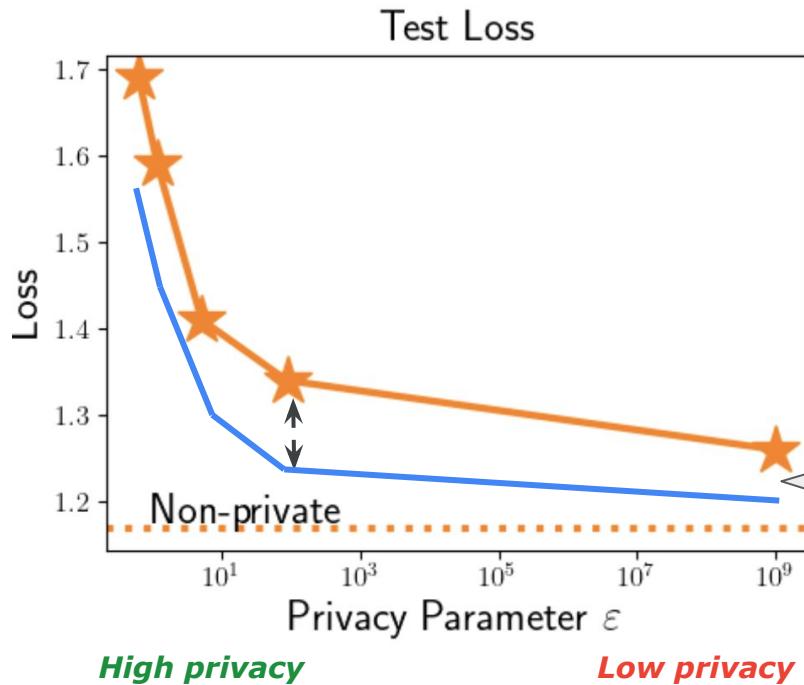


A randomized algorithm is **ϵ -differentially private** if the addition of **one unit of data** does not alter its output distribution by more than ϵ .

Differential privacy eliminates memorization



Goal: Better privacy-utility trade-offs



Goal: Achieve *smaller loss* at each privacy level

How do we train models with DP?

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

Model parameters

Loss function

Data

DP-SGD: How do we train models with DP?

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

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

Independent Gaussian noise

Learning rate

The diagram illustrates the DP-SGD update rule. It features a central equation $\theta_{t+1} = \theta_t - \eta (g_t + z_t)$. Three callout boxes point to specific components: one to the term g_t with the text "Stochastic gradient clipped to $\|g\|_2 \leq 1$ per-example", another to the term z_t with the text "**Independent** Gaussian noise", and a third to the learning rate term η with the text "Learning rate".

DP-FTRL: DP Training with *Correlated* Noise

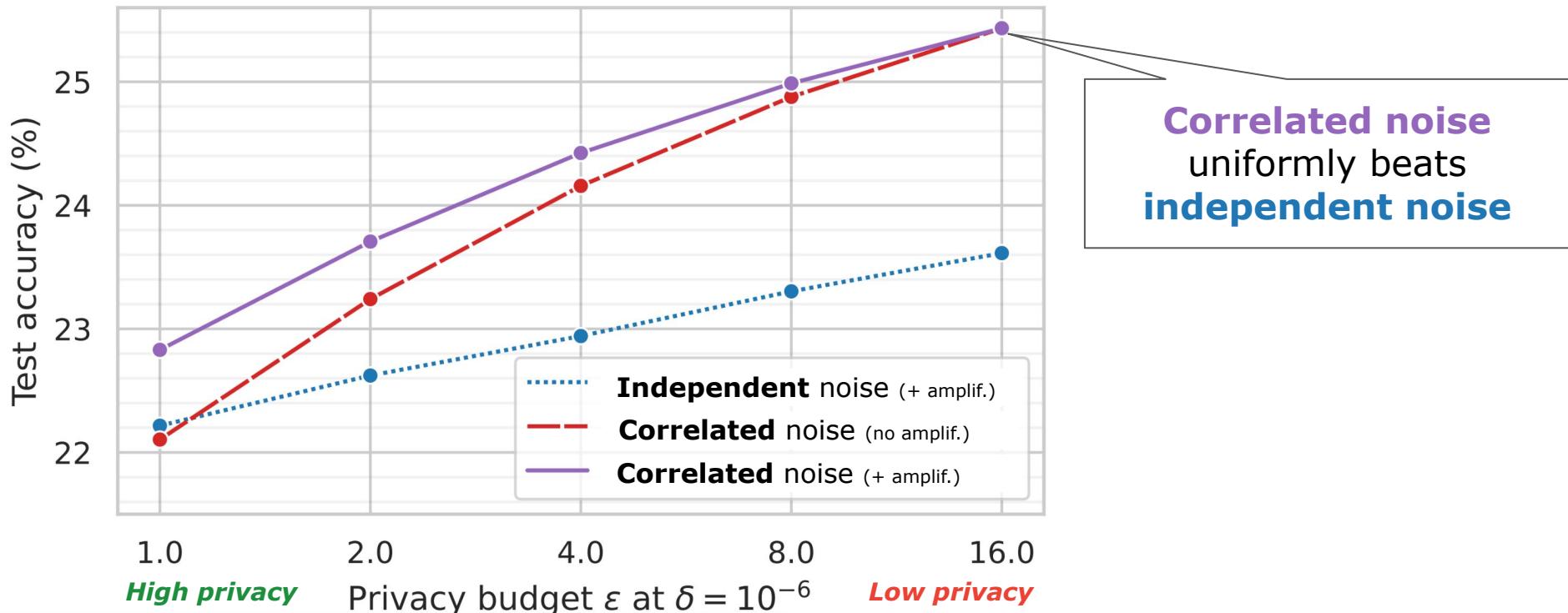
The diagram illustrates the addition of correlated Gaussian noise to the FTRL update rule. A yellow box labeled '(Anti-)correlated Gaussian noise (z_t i.i.d. Gaussian)' has a downward-pointing arrow pointing to the term z_t in the update equation.

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

Kairouz, McMahan, Song, Thakkar, Thakurta, Xu. **Practical and Private (Deep) Learning without Sampling or Shuffling**. ICML 2021.
Denisov, McMahan, Rush, Smith, Thakurta. **Improved Differential Privacy for SGD via Optimal Private Linear Operators on Adaptive Streams**. NeurIPS 2022.

Prior work: (Empirically) correlated noise outperforms independent noise

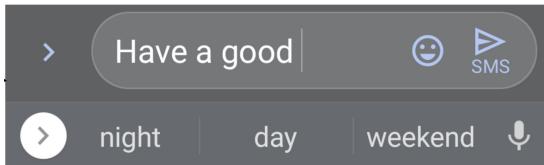
Experiment: DP language modeling
Dataset: StackOverflow



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

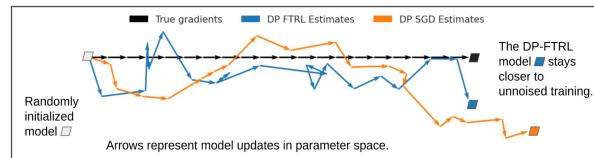
Federated Learning with Formal Differential Privacy Guarantees

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

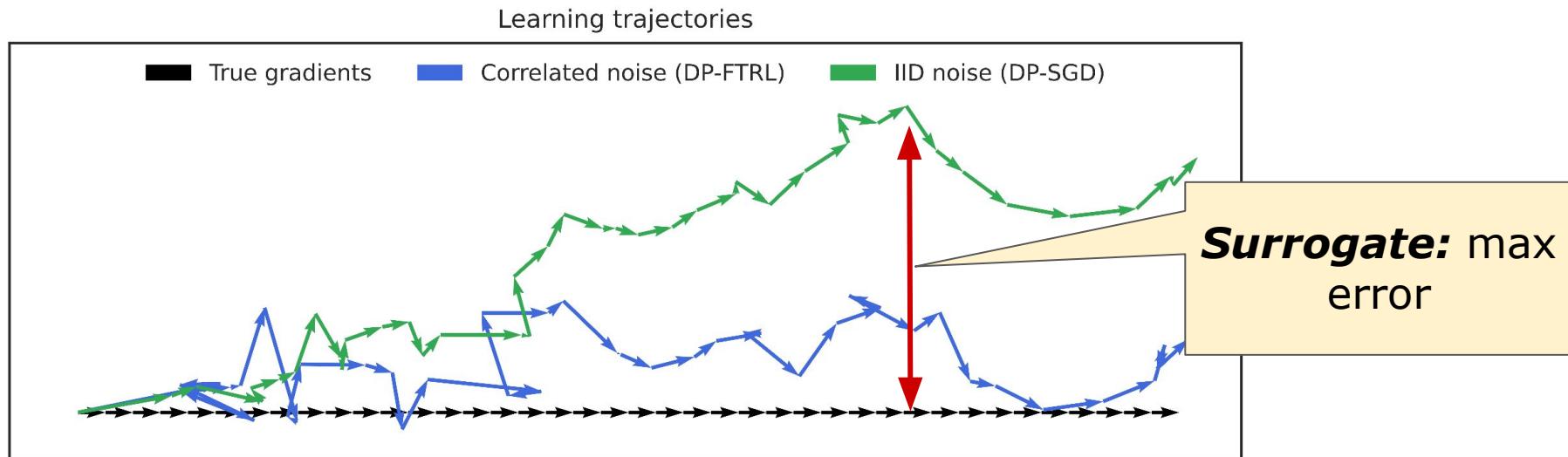
Along with fundamentals like transparency and consent, the [privacy principles of data minimization and anonymization](#) are important in ML applications that involve sensitive data.

How do we find the noise coefficients?

How do we find the noise coefficients?

Current Approach:

Find the noise coefficients β_t 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?

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

$\underbrace{\phantom{g_t + z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau}}}_{=:w_t}$

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

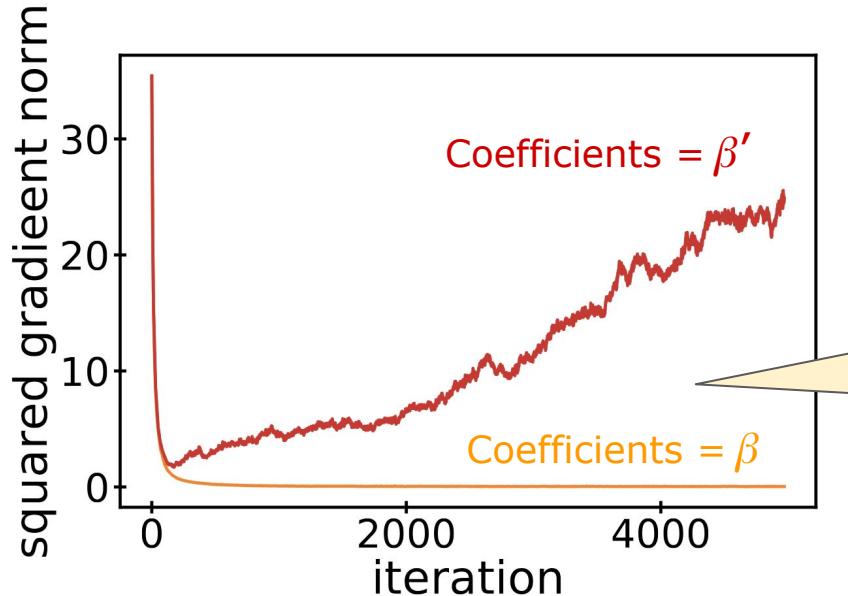
**Surrogate
Objective**

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

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

Part 1: Is correlated noise provably better for learning problems?

The surrogate objective is not related to the learning objective



$$\mathcal{E}(\beta) = \mathcal{E}(\beta')$$

Same **max error** but
different
learning performance

Part 1: 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

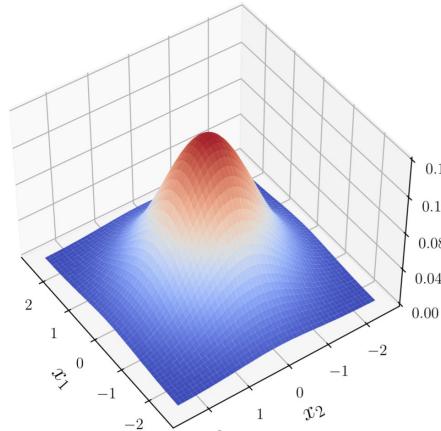
$$\Theta(d)$$

Correlated noise

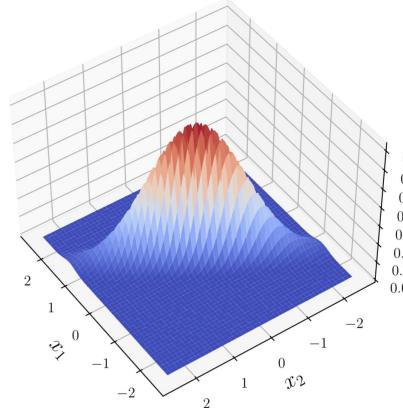
$$\tilde{\Theta}(d_{\text{eff}})$$

Lower bound

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



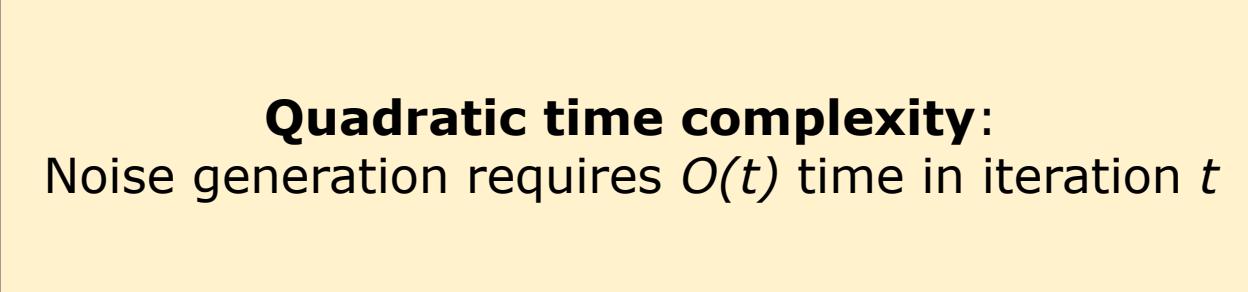
High
effective
dimension



Low
effective
dimension

Part 2: Noise generation time complexity

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



Quadratic time complexity:

Noise generation requires $O(t)$ time in iteration t

Part 2: Near-optimal noise generation time complexity

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

Our approach: Approximate the noise coefficients as

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

Per-iteration time:
 $O(d \times \text{dimension})$

Part 2: Near-optimal noise generation time complexity

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

Our approach: Approximate the noise coefficients as

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

Per-iteration time:
 $O(d \times \text{dimension})$

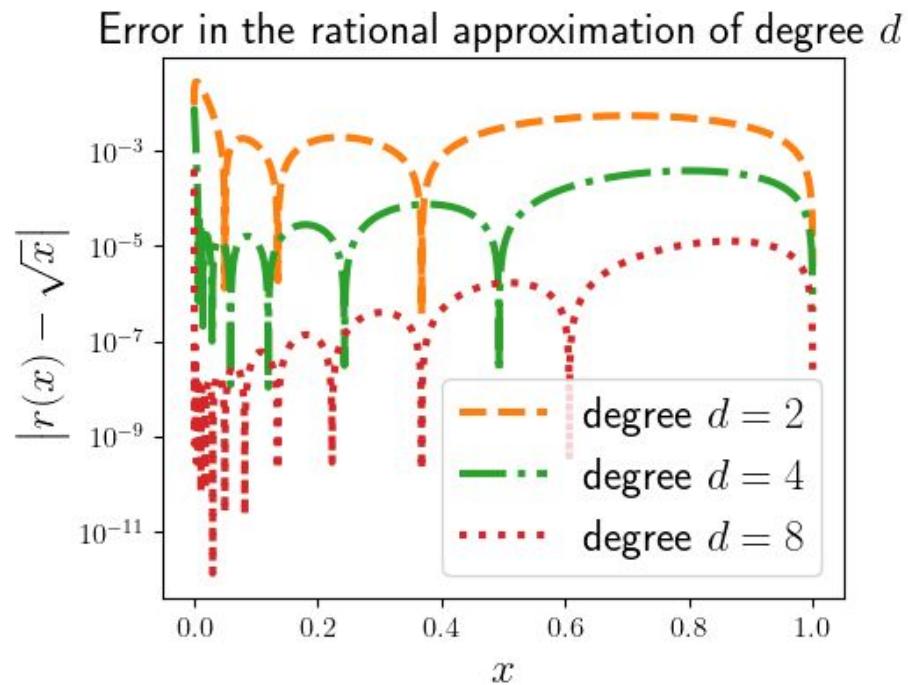
Error: If $d = O(\log^2(n/c))$, then the error is $\mathcal{E}(\beta') \leq \mathcal{E}(\beta) + c$
(n = Number of steps)

Part 2: Near-optimal noise generation time complexity

Key insight: Approximation theory

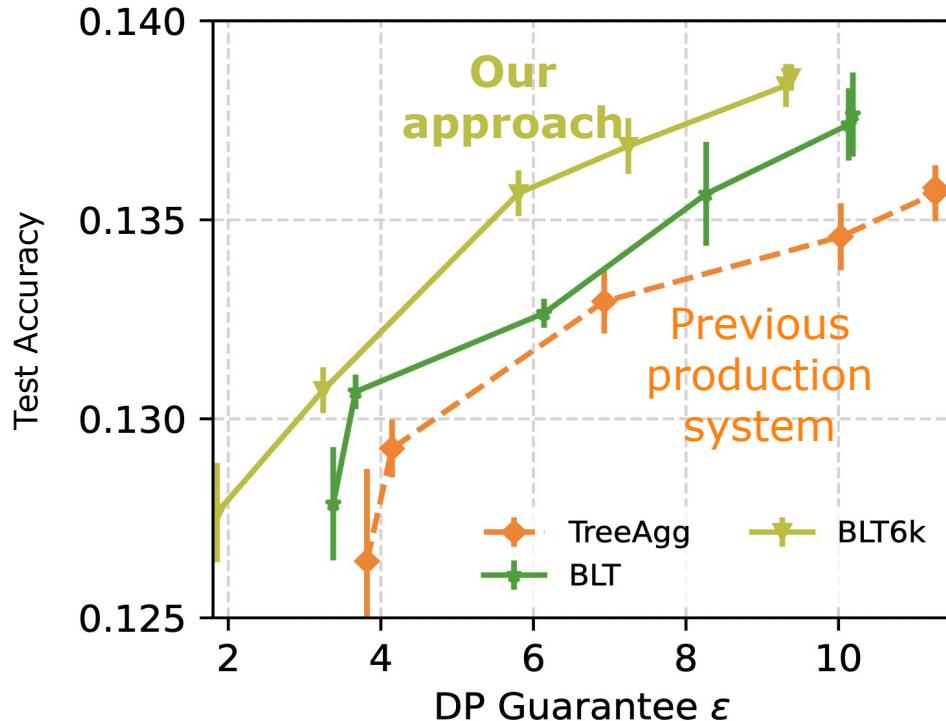
There exists a rational function $r(x)$ of degree- d that satisfies the approximation:

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



Practical Impact:

Google's production language model (Portuguese)

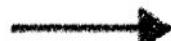


Plot: McMahan, Xu, Zhang (2024)

Background

Differential privacy (DP)

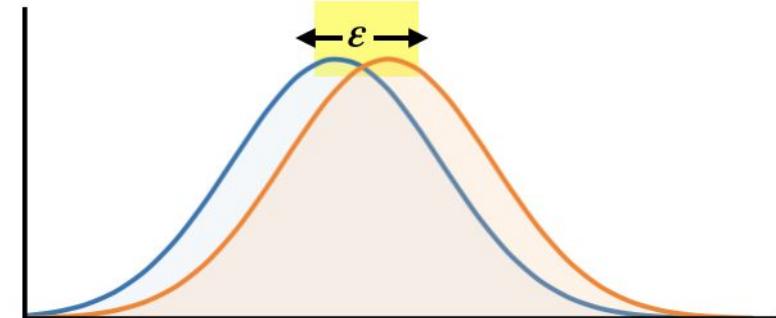
Dataset



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

ϱ -Zero-Concentrated DP (ϱ -zCDP)

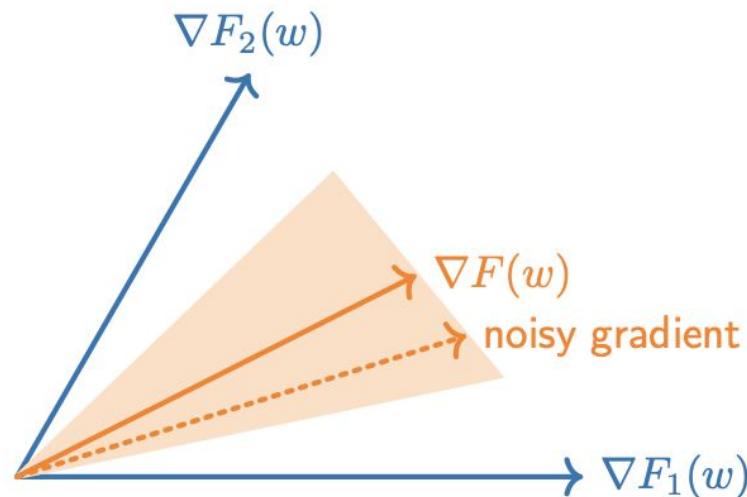
For all $0 < \alpha < \infty$, we have

$$D_\alpha \left(\mathcal{A} \left(\begin{array}{c} \text{file} \\ \text{file} \\ \text{file} \\ \text{file} \end{array} \right) \middle\| \mathcal{A} \left(\begin{array}{c} \text{file} \\ \text{file} \\ \text{file} \\ \text{file} \\ + \\ \text{file} \end{array} \right) \right) \leq \rho\alpha$$

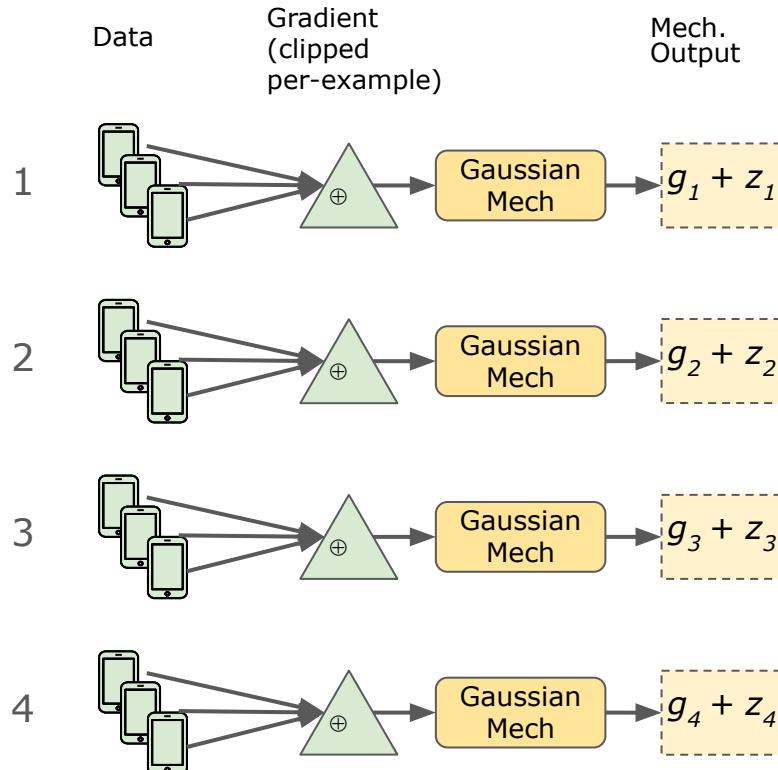
↑
Rényi α -divergence

DP-SGD / Independent noise

Primitive: private mean estimation of minibatch (clipped) gradients in each iteration



DP-SGD adds independent noise in each iteration

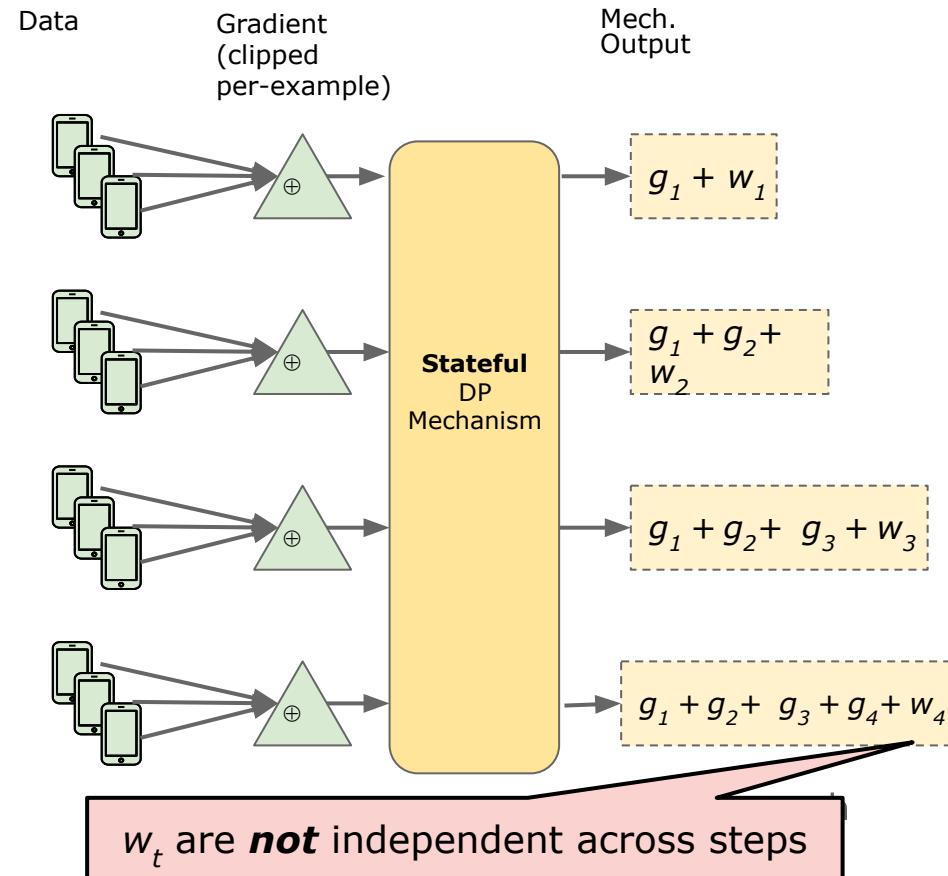


For ϱ -zCDP, take
 $\text{Var}(z_t) = 1/(2\varrho)$

DP-FTRL: privatize prefix sums of gradients

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

SGD update (without noise)



DP-FTRL: privatize prefix sums of gradients

For ϱ -zCDP, take

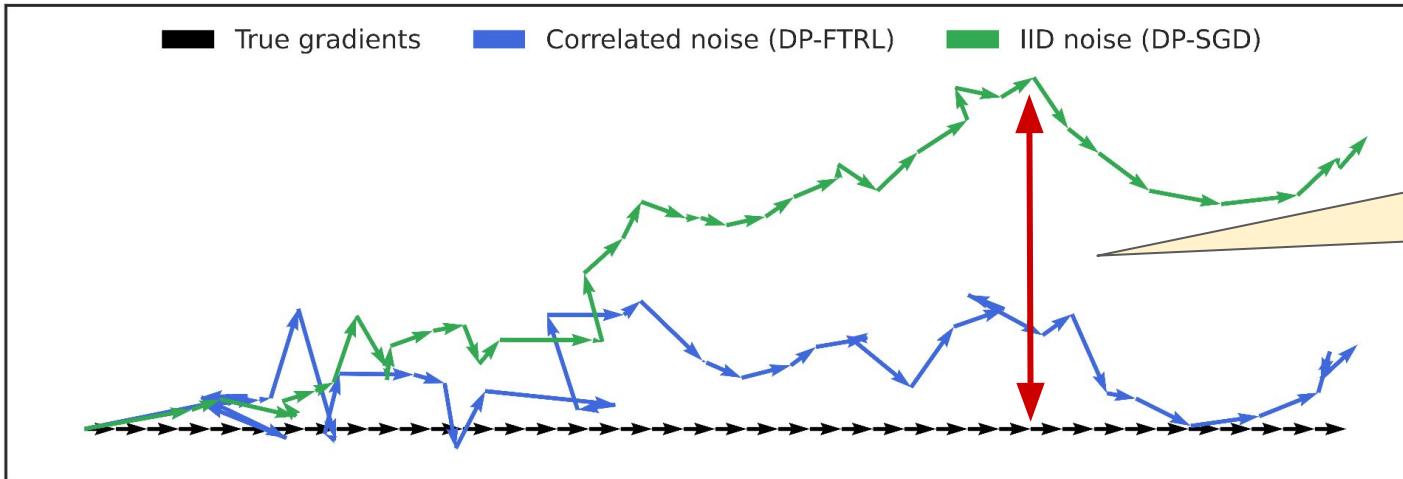
$$\text{Var}(z_t) = \frac{1}{2\rho} \left(\max_{t \leq n} \| (B^{-1})[:, t] \|_2^2 \right)$$

$$B = \begin{pmatrix} 1 & & & & \\ -\beta_1 & 1 & & & \\ -\beta_2 & -\beta_1 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \\ -\beta_{n-1} & -\beta_{n-2} & \cdots & \cdots & 1 \end{pmatrix} \quad \text{sensitivity}$$

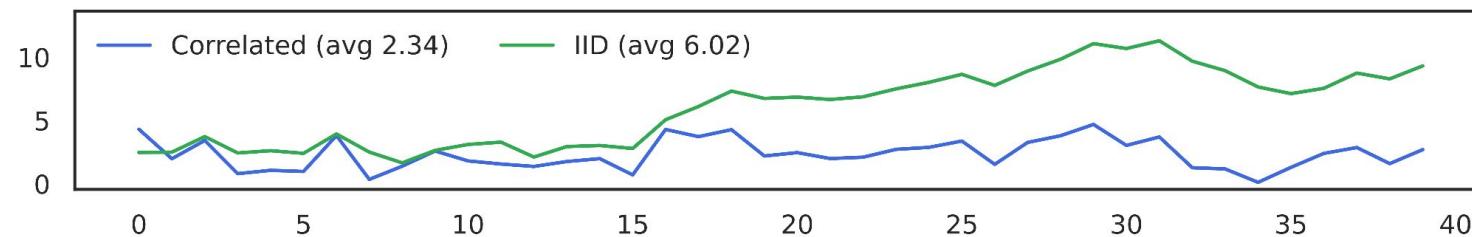
$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau} \right)$$

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

Learning trajectories



Prefix sum error



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

Part 1: Learning guarantees

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

ICLR 2024

DP-FTRL vs. DP-SGD: Previous Theory

For convex & G -Lipschitz losses

Independent Noise

$$\frac{Gd^{1/4}}{\sqrt{\rho T}}$$

Correlated Noise

$$\frac{Gd^{1/4}}{\sqrt{\rho^2 T}}$$

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

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

Model parameters

Loss function

Data

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

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

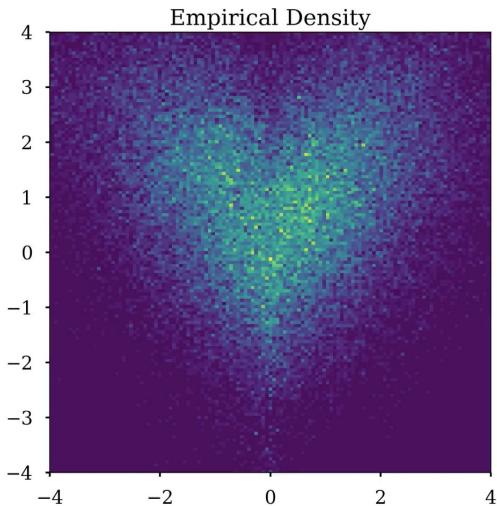
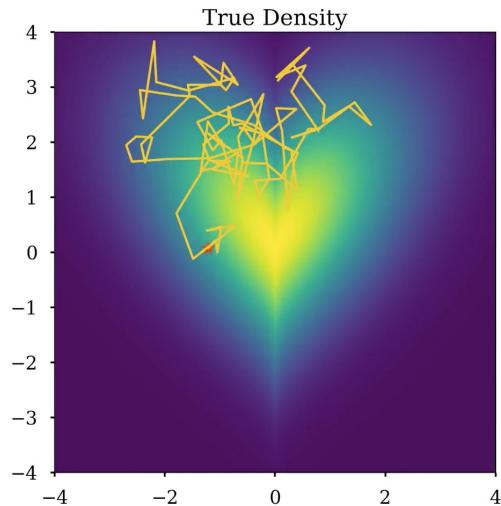


Image credit:
[Abdul Fatir Ansari](#)

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

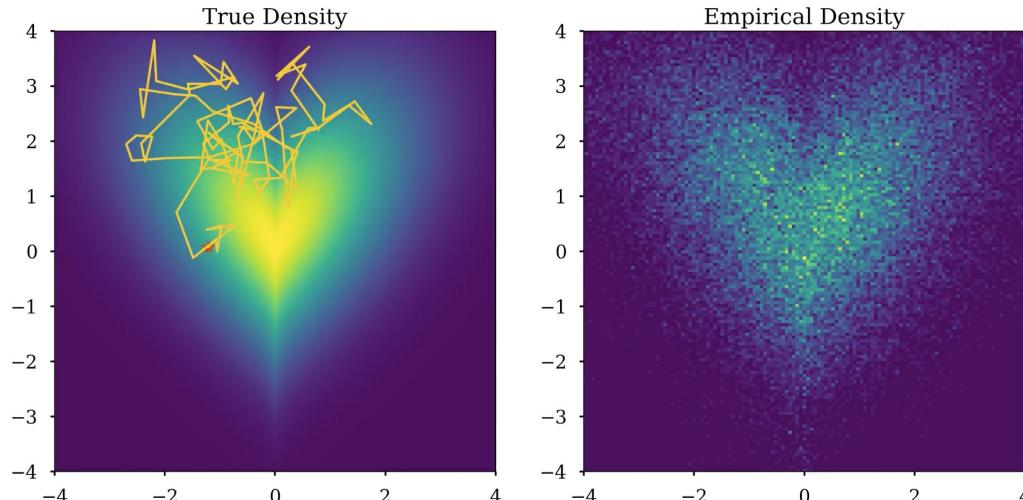


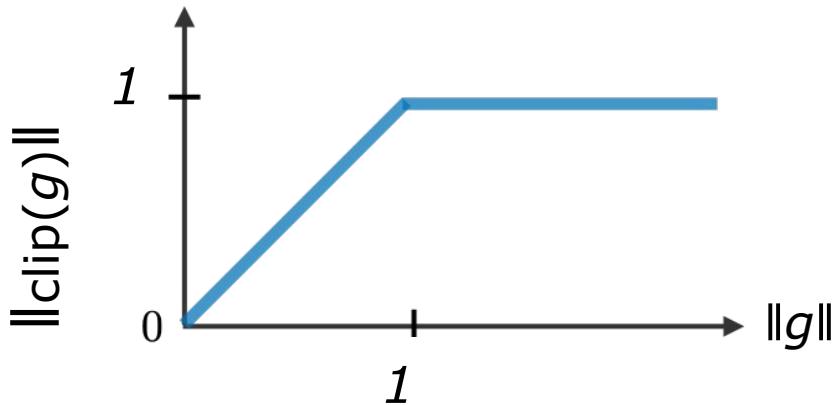
Image credit:
[Abdul Fatir Ansari](#)

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)

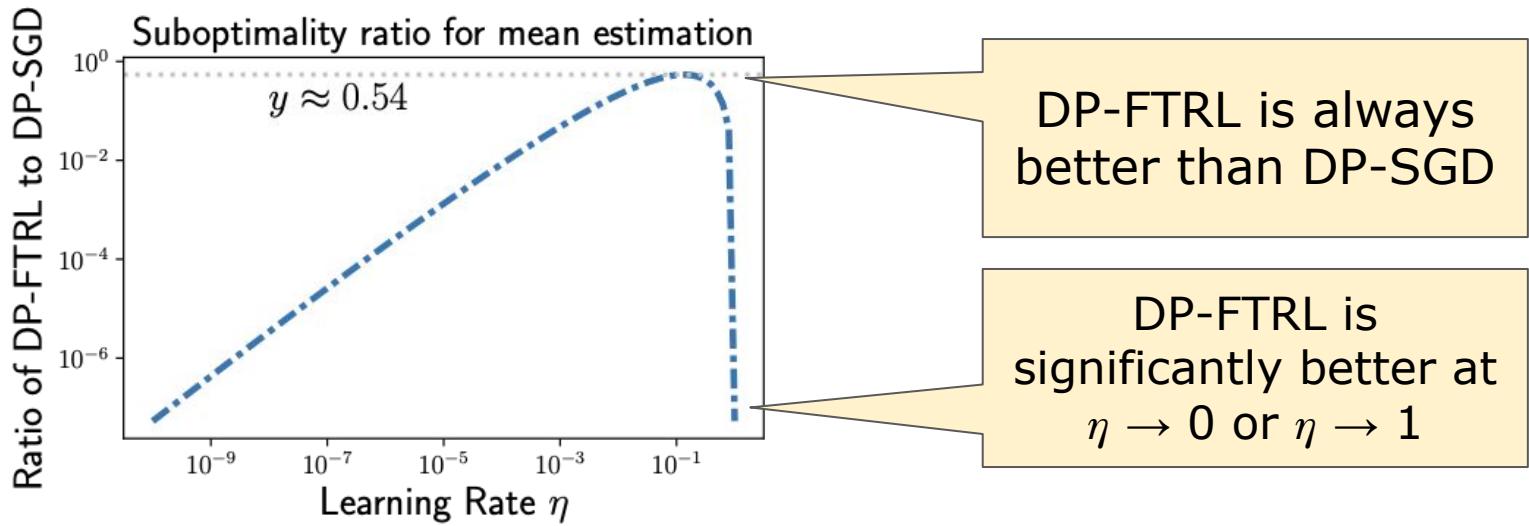
$$F_\infty(\beta^{\text{sgd}}) = \rho^{-1}\eta$$

Correlated noise (DP-FTRL)

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

η : constant learning rate in $(0, 1)$

ϱ : privacy level



Closed form correlations for mean estimation

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

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

Closed form correlations for mean estimation

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

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

ν -DP-FTRL

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

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 (ν -Noisy-FTRL)

$$\leq d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$$

Lower bound for any algorithm

$$\geq d_{\text{eff}} \rho^{-1} \eta^2$$

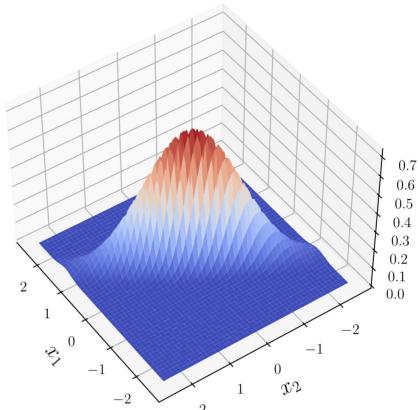
Improve **dimension d** to
problem-dependent
effective dimension d_{eff}

Effective dimension

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

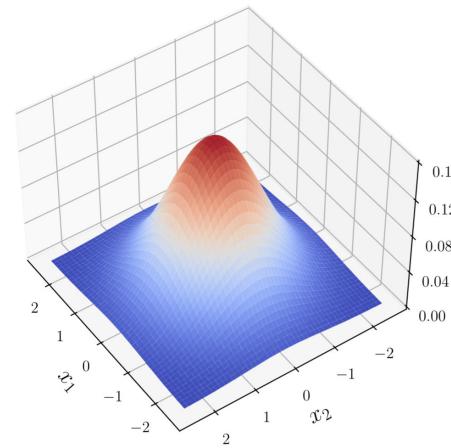
Low effective dimension

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



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.

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

[Rudelson & Vershynin (J. ACM 2007)]

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 (ν -Noisy-FTRL)

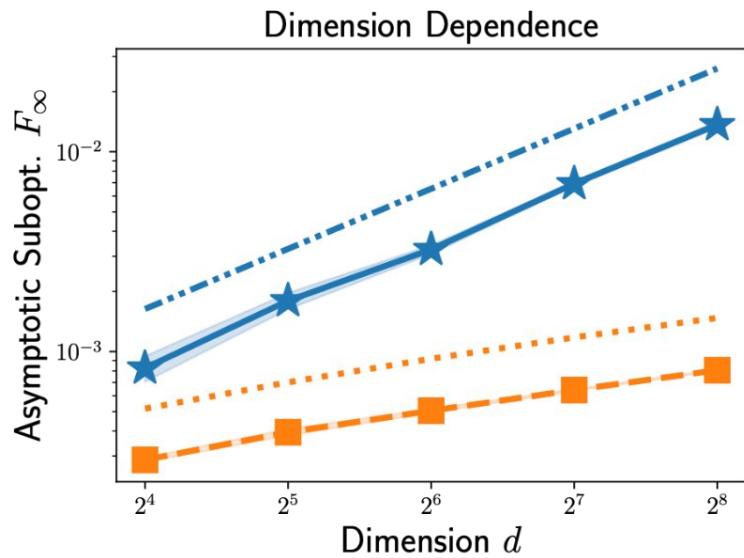
$$\leq d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$$

Lower bound for any algorithm

$$\geq d_{\text{eff}} \rho^{-1} \eta^2$$

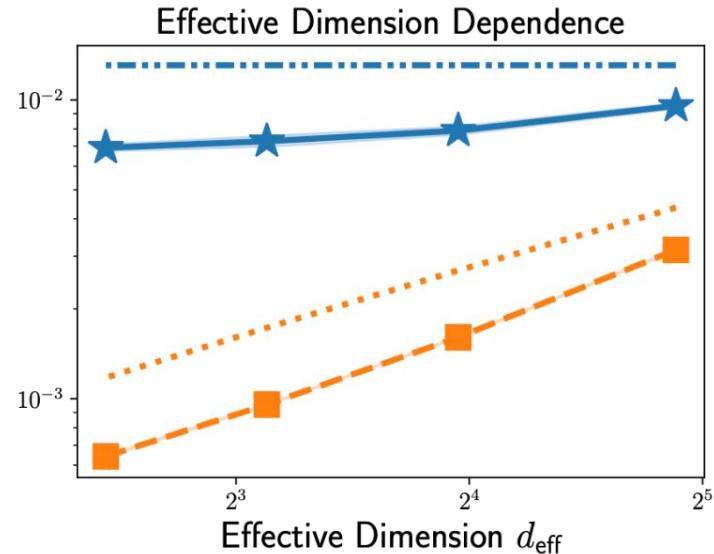
Improve *dimension d* to
problem-dependent
effective dimension d_{eff}

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 (ν -Noisy-FTRL)

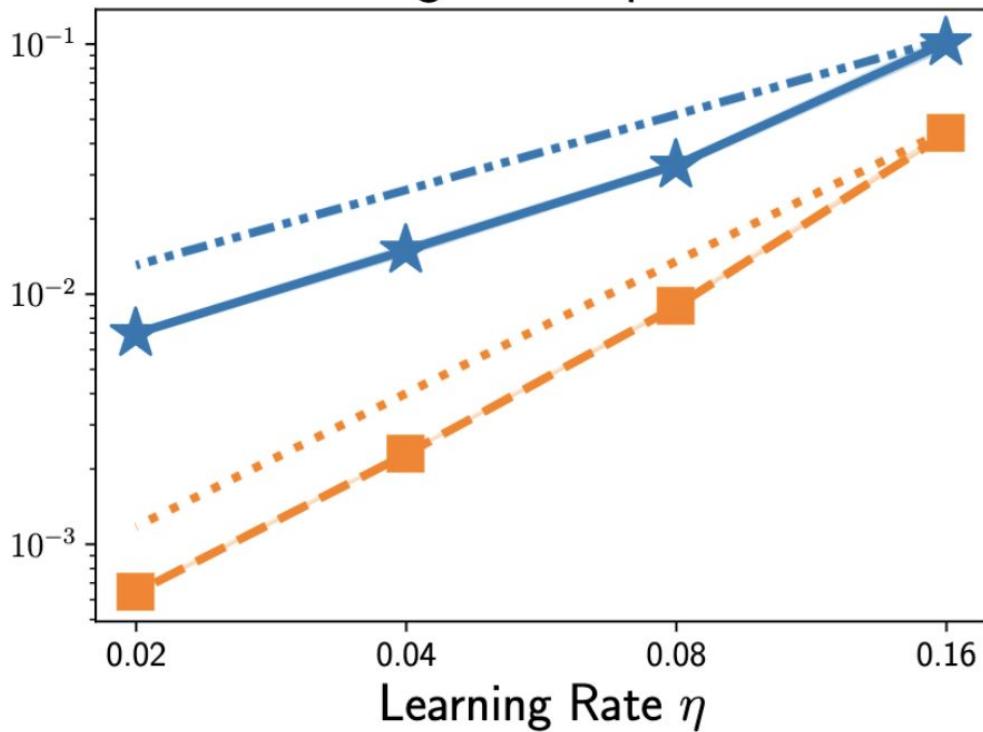
$$\leq d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \left(\frac{1}{\eta \mu} \right)$$

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

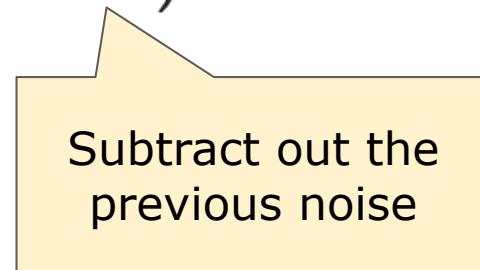
Noisy-FTRL \gg Noisy-SGD at small η

Anticorrelated Noise Injection for Improved Generalization

Antonio Orvieto ^{*1} Hans Kersting ^{*2} Frank Proske ³ Francis Bach ² Aurelien Lucchi ⁴

Anti-PGD [Orvieto et al. (ICML '22)] corresponds to $\beta_1 = 1$

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



Subtract out the
previous noise

Anticorrelated Noise Injection for Improved Generalization

Antonio Orvieto ^{*1} Hans Kersting ^{*2} Frank Proske ³ Francis Bach ² Aurelien Lucchi ⁴

Anti-PGD [Orvieto et al. (ICML '22)] corresponds to $\beta_1 = 1$

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

Asymptotic error = ∞ (as sensitivity scales of $O(t)$ for t iterations)

Anti-PGD can be adapted for DP by damping: take $\beta_1 = \nu$ ($0 < \nu < 1$)

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

$$\text{Asymptotic error} = \sqrt{d d_{\text{eff}}} \rho^{-1} \eta^{3/2}$$

Geometric mean of
Noisy-SGD and
lower bound

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

Usual stochastic gradient proof patterns do not work:

No Markovian/martingale structure in the noise

Our approach: Analysis the *Fourier* domain

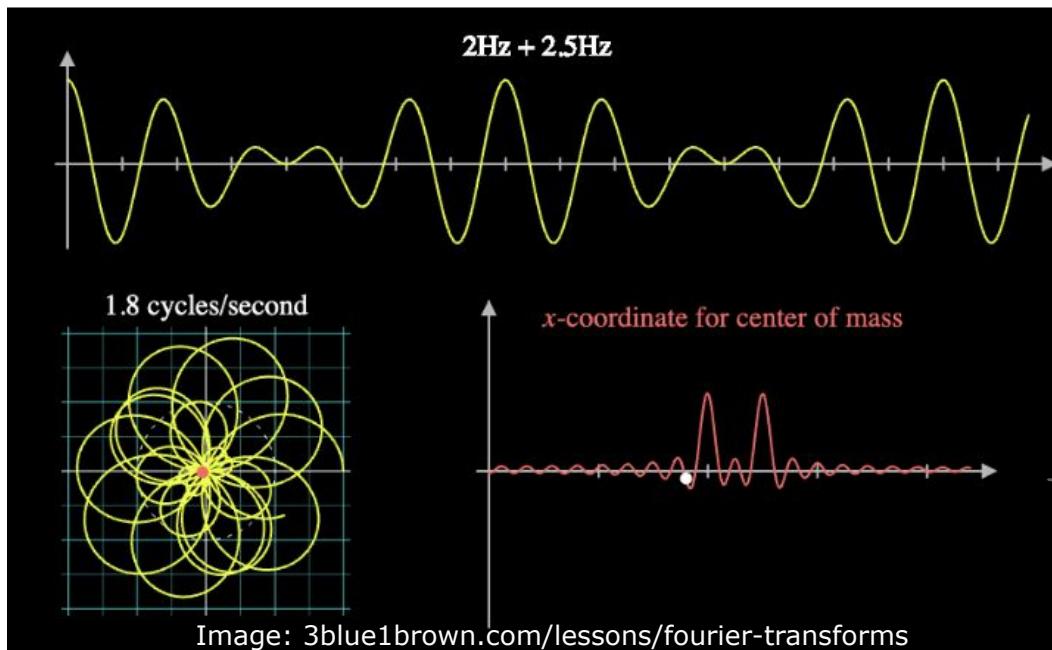
Letting $\delta_t = \theta_t - \theta_*$, the DP-FTRL update can be written as

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



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

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

$$B = \begin{pmatrix} 1 & & & & \\ -\beta_1 & 1 & & & \\ -\beta_2 & -\beta_1 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \\ -\beta_{n-1} & -\beta_{n-2} & \cdots & \cdots & 1 \end{pmatrix}$$

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

Requires $|B(\omega)|$
small

sensitivity

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

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

Requires $|B(\omega)|$
large

$$B = \begin{pmatrix} 1 & & & & \\ -\beta_1 & 1 & & & \\ -\beta_2 & -\beta_1 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \\ -\beta_{n-1} & -\beta_{n-2} & \cdots & \cdots & 1 \end{pmatrix}$$

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

Requires $|B(\omega)|$
small

sensitivity

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

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

Requires $|B(\omega)|$
large

$$B = \begin{pmatrix} 1 & & & & \\ -\beta_1 & 1 & & & \\ -\beta_2 & -\beta_1 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \\ -\beta_{n-1} & -\beta_{n-2} & \cdots & \cdots & 1 \end{pmatrix}$$

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

For linear regression:

$$\boldsymbol{\theta}'_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\boldsymbol{\theta}'_t + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_{\tau} \mathbf{w}_{t-\tau}. \quad (25)$$

Multiplicative
noise

$$\boldsymbol{\theta}'_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\boldsymbol{\theta}'_t + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_\tau \mathbf{w}_{t-\tau}. \quad (25)$$

Decomposition:

$$\begin{aligned} \boldsymbol{\theta}_{t+1}^{(0)} &= (\mathbf{I} - \eta \mathbf{H})\boldsymbol{\theta}_t^{(0)} + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_\tau \mathbf{w}_{t-\tau}, \\ \boldsymbol{\theta}_{t+1}^{(r)} &= (\mathbf{I} - \eta \mathbf{H})\boldsymbol{\theta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r-1)} \text{ for } r > 0, \\ \boldsymbol{\delta}_{t+1}^{(r)} &= (\mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\delta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r)}. \end{aligned}$$

$$\boldsymbol{\theta}'_t = \sum_{r=0}^m \boldsymbol{\theta}_t^{(r)} + \boldsymbol{\delta}_t^{(m)}.$$

Aguech, Moulines, Priouret. **On a Perturbation Approach for the Analysis of Stochastic Tracking Algorithms.** SIAM J. Control. Optim., 2000
 Bach and Moulines. **Non-Strongly-Convex Smooth Stochastic Approximation with Convergence Rate $O(1/n)$.** NeurIPS 2013.

$$\boldsymbol{\theta}'_{t+1} = (\mathbf{I} - \eta(\mathbf{x}_t \otimes \mathbf{x}_t))\boldsymbol{\theta}'_t + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_{\tau} \mathbf{w}_{t-\tau}. \quad (25)$$

Decomposition:

$$\begin{aligned} \boldsymbol{\theta}_{t+1}^{(0)} &= (\mathbf{I} - \eta \mathbf{H})\boldsymbol{\theta}_t^{(0)} + \eta \xi_t \mathbf{x}_t - \eta \sum_{\tau=0}^{\infty} \beta_{\tau} \mathbf{w}_{t-\tau}, \\ \boldsymbol{\theta}_{t+1}^{(r)} &= (\mathbf{I} - \eta \mathbf{H})\boldsymbol{\theta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r-1)} \text{ for } r > 0, \\ \boldsymbol{\delta}_{t+1}^{(r)} &= (\mathbf{I} - \eta \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\delta}_t^{(r)} + \eta(\mathbf{H} - \mathbf{x}_t \otimes \mathbf{x}_t)\boldsymbol{\theta}_t^{(r)}. \end{aligned} \quad \boldsymbol{\theta}'_t = \sum_{r=0}^m \boldsymbol{\theta}_t^{(r)} + \boldsymbol{\delta}_t^{(m)}.$$

Aguech, Moulines, Priouret. **On a Perturbation Approach for the Analysis of Stochastic Tracking Algorithms.** SIAM J. Control. Optim., 2000
 Bach and Moulines. **Non-Strongly-Convex Smooth Stochastic Approximation with Convergence Rate $O(1/n)$.** NeurIPS 2013.

Key idea: $\mathbb{E} [\boldsymbol{\delta}_0^{(m)} \otimes \boldsymbol{\delta}_0^{(m)}] \rightarrow \mathbf{0}$ as $m \rightarrow \infty$.

Thus, $\|\boldsymbol{\theta}'_t\| \leq \sum_{r=0}^{\infty} \|\boldsymbol{\theta}_t^{(r)}\|$

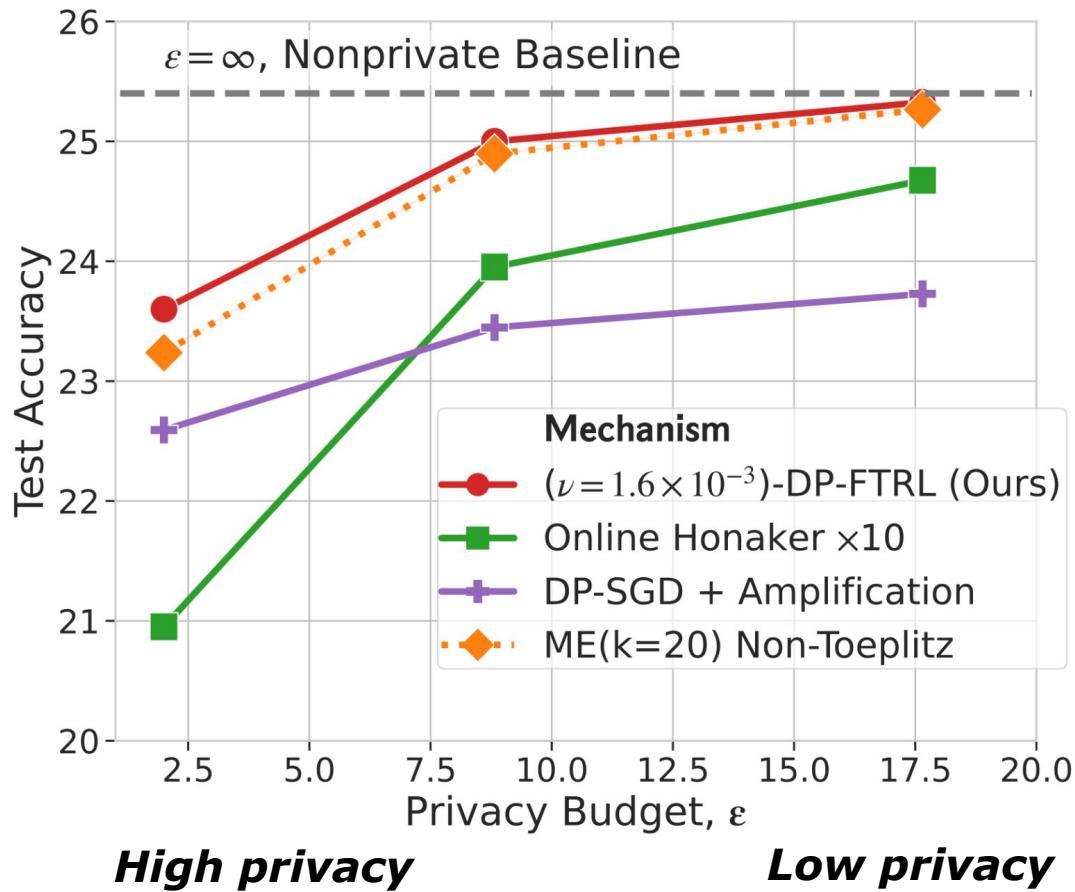
Experiments

ν -DP-FTRL

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

and tune the parameter ν

Language modeling with Stack Overflow | User-level DP

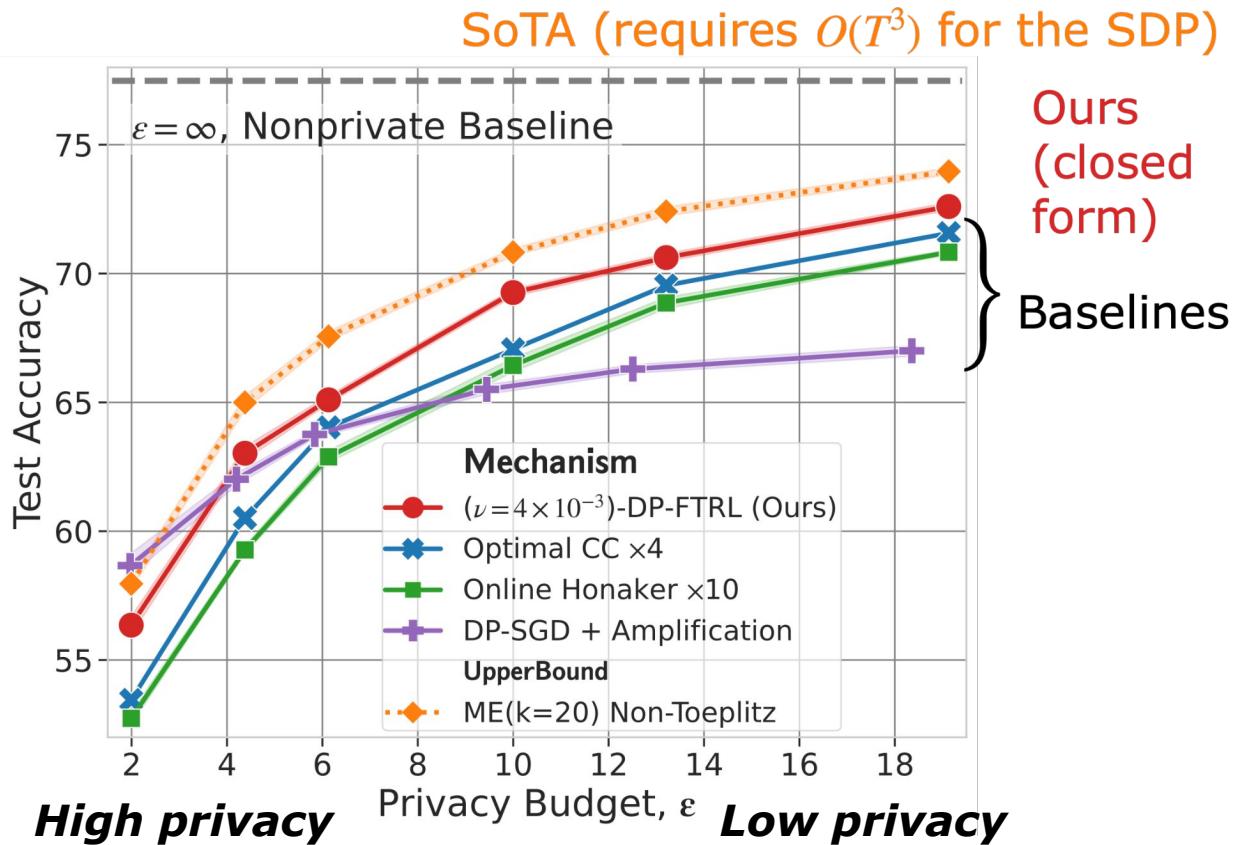


Ours
matches
SoTA!

High privacy

Low privacy

Image classification with CIFAR-10 | Example-level DP



Part 2: Efficient noise generation

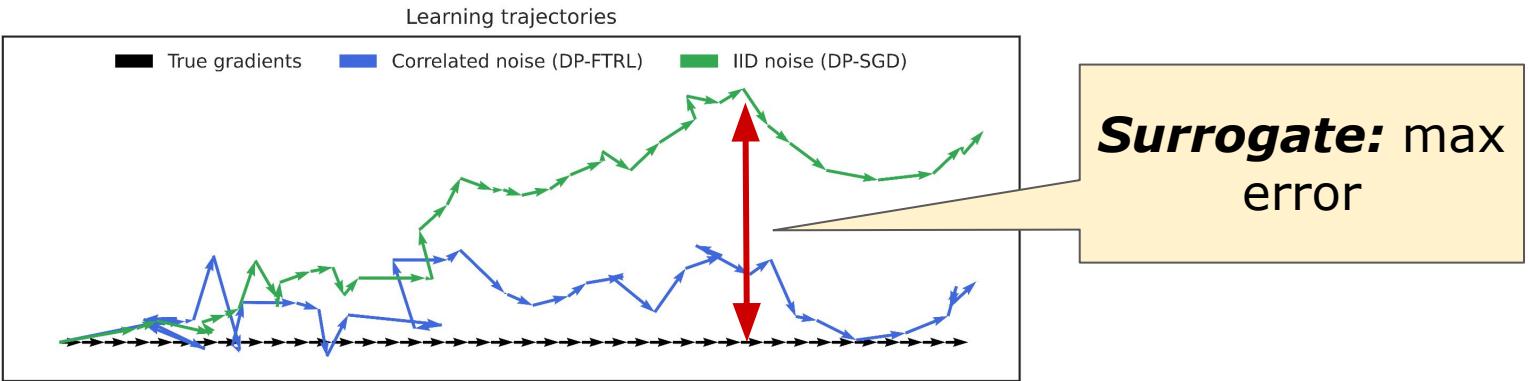
With near-optimal privacy-utility trade-offs

FOCS 2024

$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t - \sum_{\tau=1}^t \beta_\tau z_{t-\tau} \right)$$
$$\beta_t = t^{-3/2}$$

Quadratic time complexity:

Noise generation requires $O(t)$ time in iteration t



Max Error

Independent noise

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

Noise generation time
(in iteration t)

$$O(\dim)$$

Optimal correlated noise

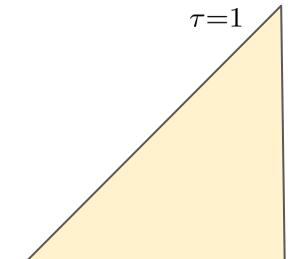
$$\frac{\log n}{\pi} + c$$

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

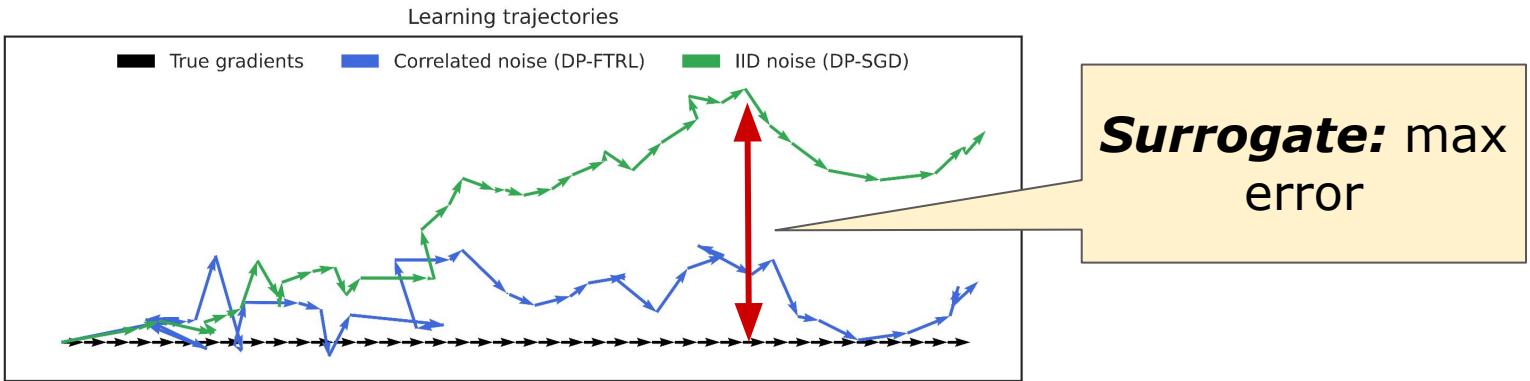
A first attempt: the banded mechanism

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

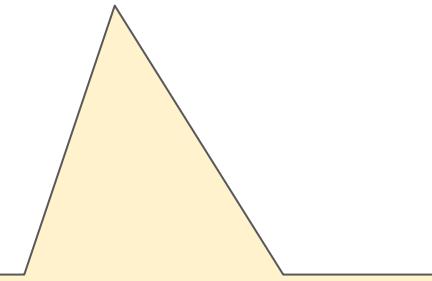


	Max Error	Noise generation time (in iteration t)
Independent noise	$\Theta(\sqrt{n})$	$O(\dim)$
Optimal correlated noise	$\frac{\log n}{\pi} + c$	$O(t \cdot \dim)$
b-Banded	$O\left((\sqrt{n/b} - 1) \log b\right)$	$O(b \cdot \dim)$

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 $w_{t+1} = \alpha z_t + \lambda w_{t-1}$



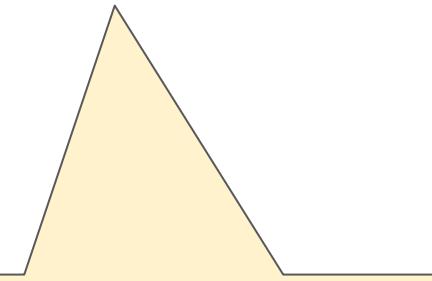
Linear complexity:

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 the recurrences

$$s_{t+1}^{(1)} = z_t + \lambda_1 z_{t-1} + \dots = s_t^{(1)} + \lambda_1 z_t$$

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 the recurrences

$$s_{t+1}^{(1)} = z_t + \lambda_1 z_{t-1} + \dots = s_t^{(1)} + \lambda_1 z_t$$

$$s_{t+1}^{(2)} = z_t + \lambda_2 z_{t-1} + \dots = s_t^{(2)} + \lambda_2 z_t$$

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 the recurrences

$$s_{t+1}^{(1)} = z_t + \lambda_1 z_{t-1} + \dots = s_t^{(1)} + \lambda_1 z_t$$

$$s_{t+1}^{(2)} = z_t + \lambda_2 z_{t-1} + \dots = s_t^{(2)} + \lambda_2 z_t$$

$$w_{t+1} = \alpha_1 s_{t+1}^{(1)} + \alpha_2 s_{t+1}^{(2)}$$

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 the recurrences

$$s_{t+1}^{(1)} = z_t + \lambda_1 z_{t-1} + \dots = s_t^{(1)} + \lambda_1 z_t$$

$$s_{t+1}^{(2)} = z_t + \lambda_2 z_{t-1} + \dots = s_t^{(2)} + \lambda_2 z_t$$

$$w_{t+1} = \alpha_1 s_{t+1}^{(1)} + \alpha_2 s_{t+1}^{(2)}$$

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	$\Theta(\sqrt{n})$	$O(\text{dim})$
Optimal correlated noise	$\frac{\log n}{\pi} + c$	$O(t \cdot \text{dim})$
b-Banded	$O\left((\sqrt{n/b} - 1) \log b\right)$	$O(b \cdot \text{dim})$
BLT of degree d	??????	$O(d \cdot \text{dim})$

Approximation Theory!!

From sequences to functions

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

Coefficients
 $1, -\beta_1, -\beta_2, \dots$

Generating
Function $r_0(x)$



From sequences to functions

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

Coefficients
 $1, -\beta_1, -\beta_2, \dots$

Generating
Function $r_0(x)$

Taylor expansion around $x = 0$

From sequences to functions

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

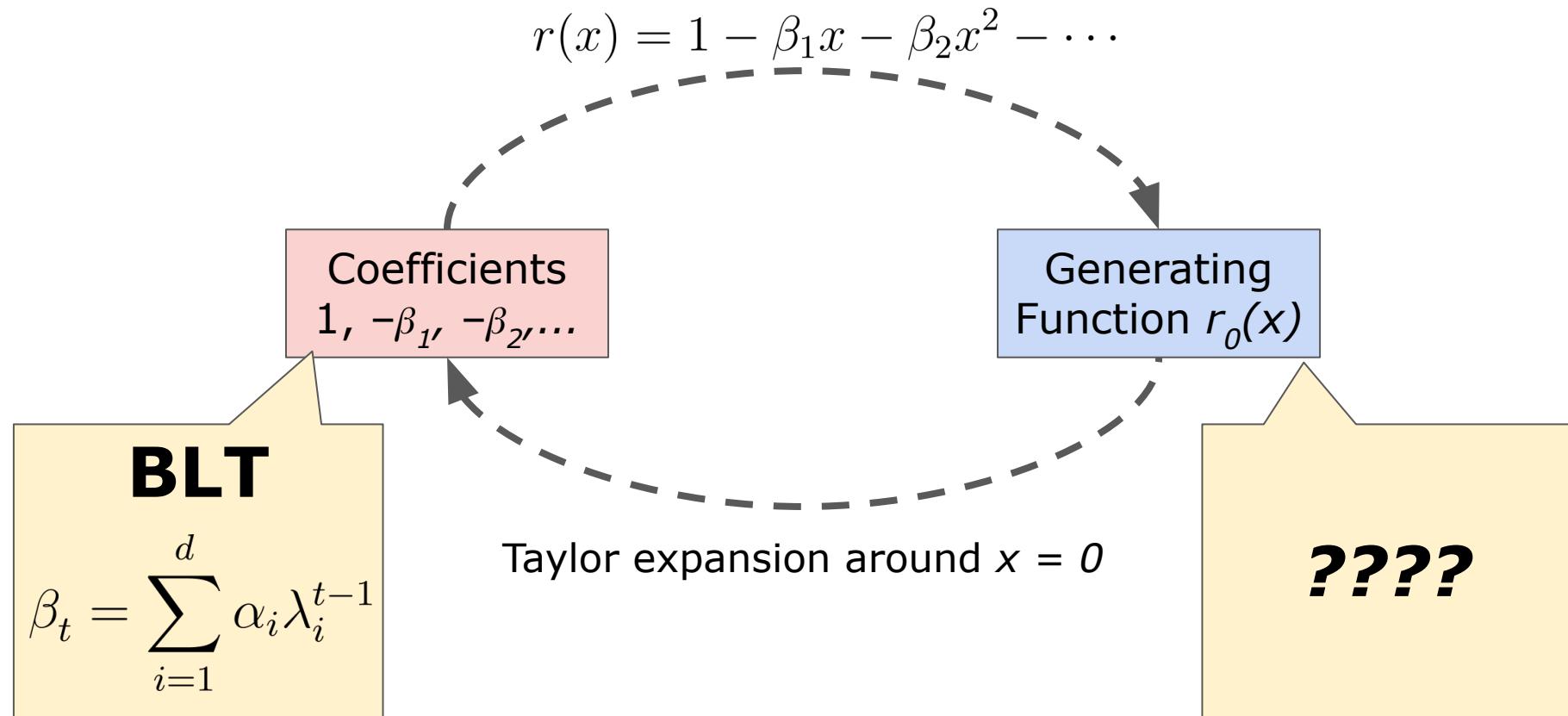
Coefficients
 $1, -\beta_1, -\beta_2, \dots$

Generating
Function $r_0(x)$

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:

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

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

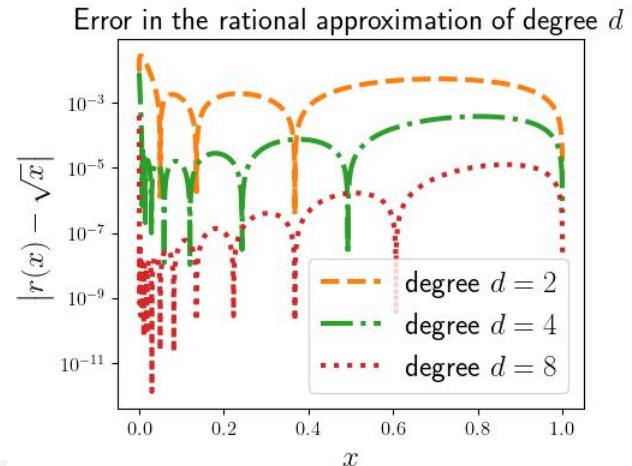
where $\text{err}(r)$ quantifies the **approximation quality**

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

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

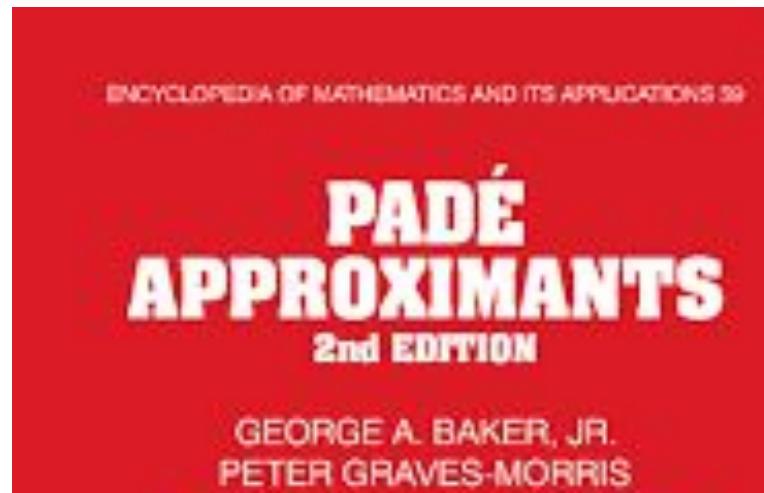
Our approach: **B**uffered **L**inear **T**oeplitz (BLT) Mechanism

	Error	Noise generation time (n : # iterations)
Independent	$\Theta(\sqrt{n})$	$O(\dim)$
Optimal correlated	$\frac{\log n}{\pi} + c$	$O(n \cdot \dim)$
Ours	$\frac{\log n}{\pi} + c$	$O(\log^2(n) \cdot \dim)$

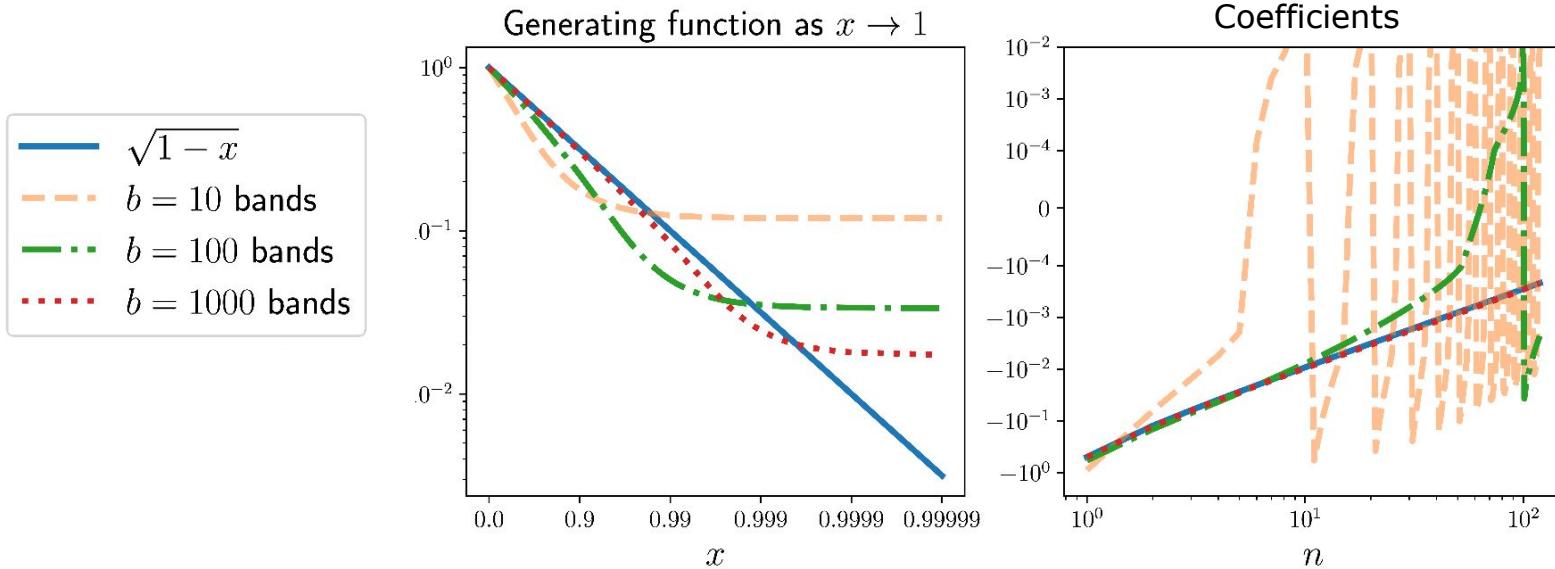
Key difference: approximation quality

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

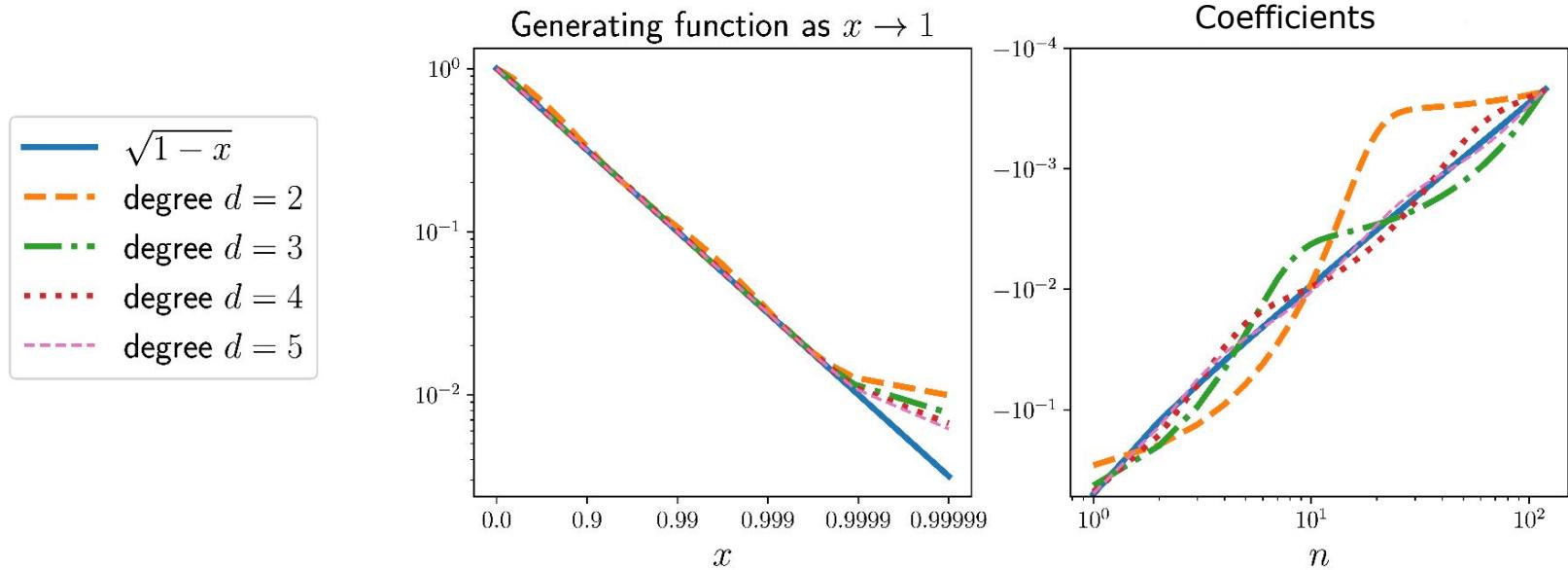


Approximation quality: banded mechanism



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

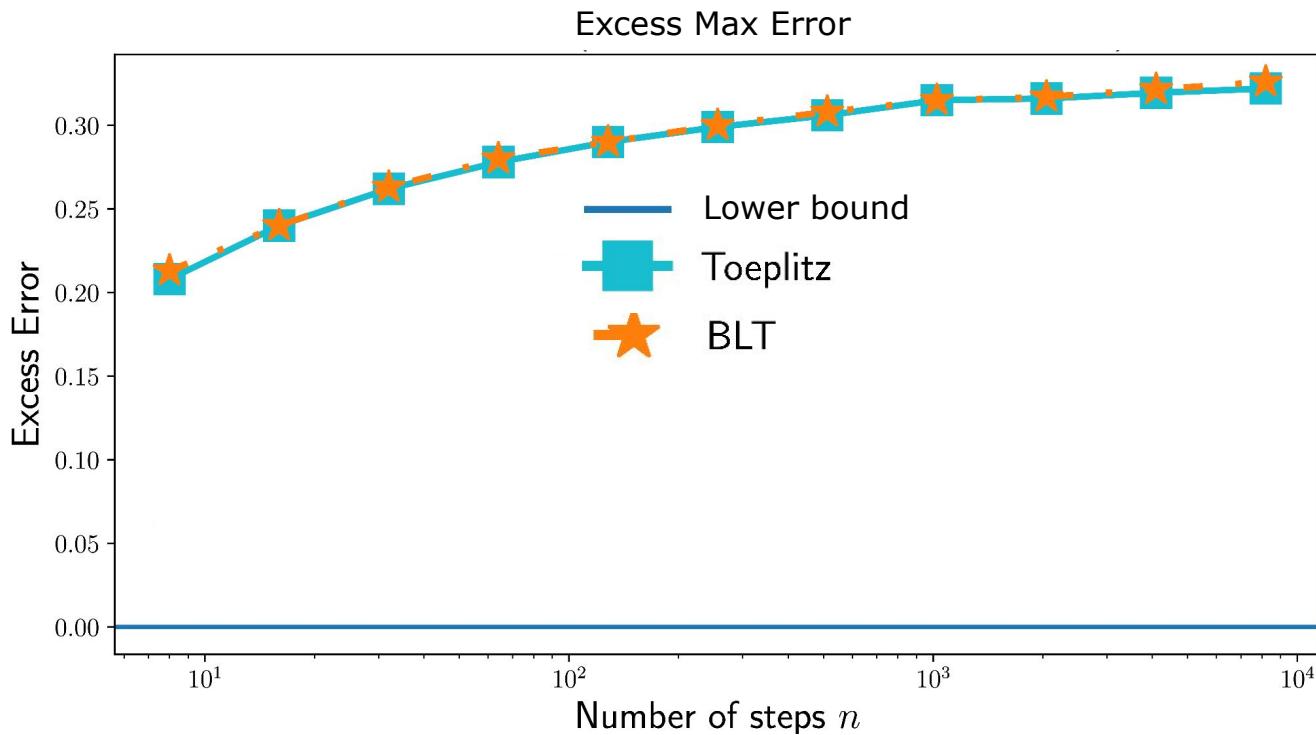
Approximation quality: BLT mechanism



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

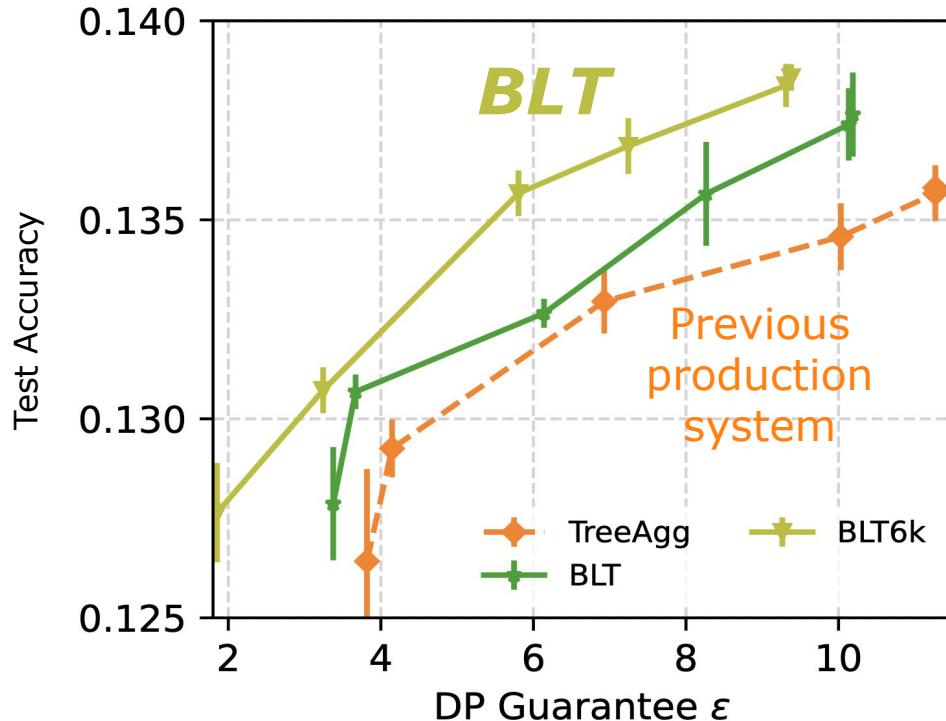
[McMahan and **P.** (2025)]

Empirical Results



Practical Impact:

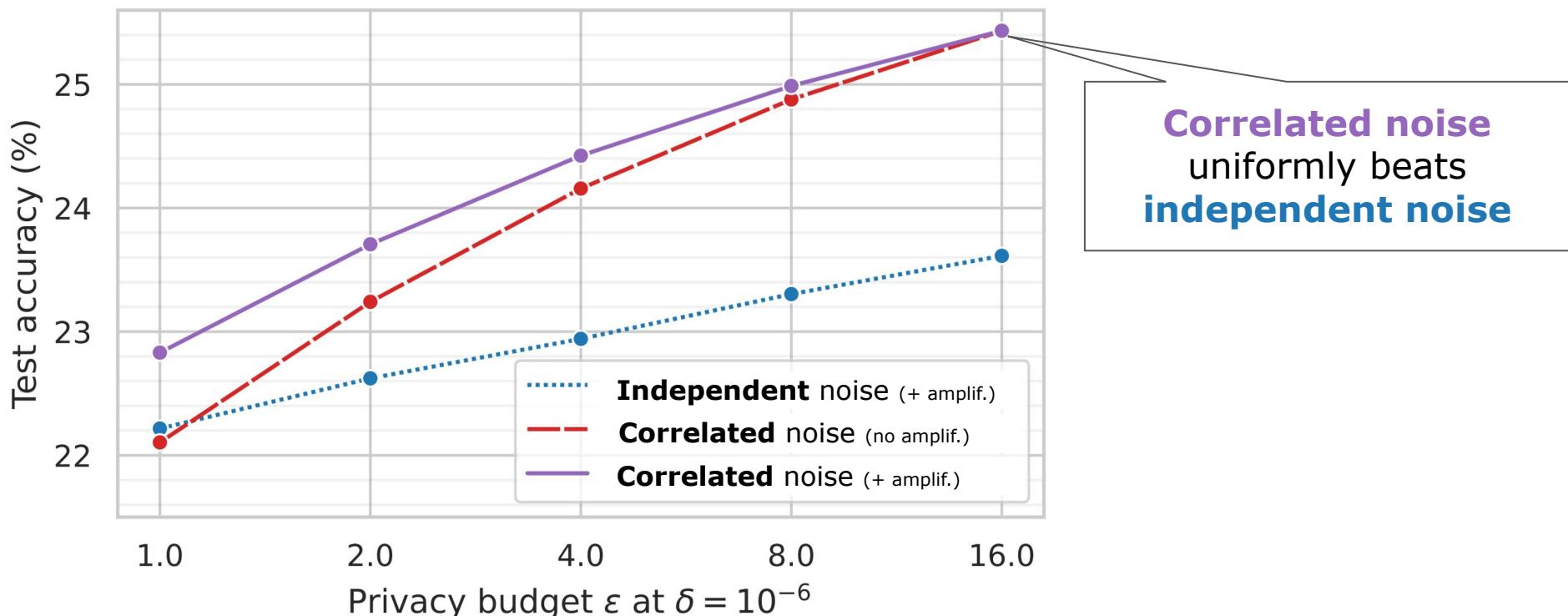
Google's production language model (Portuguese)



Plot: McMahan, Xu, Zhang (2024)

Conclusion and Open Problems

Goal: Better privacy-utility trade-offs

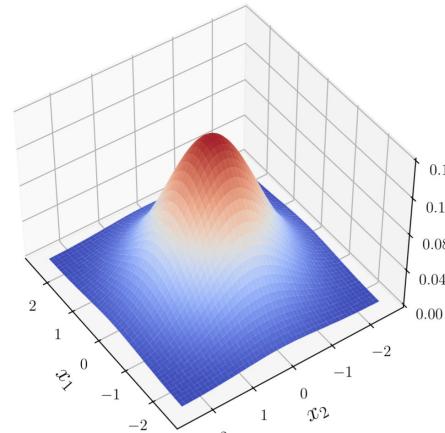


Part 1: Correlated noise algorithms **are** provably better

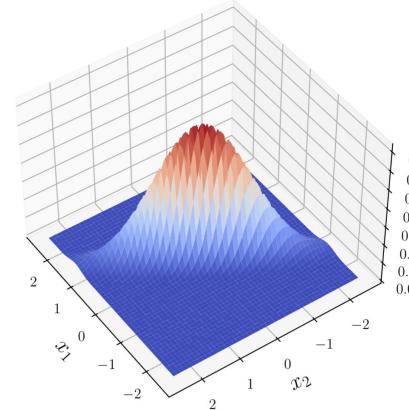
(Anti-) correlated noise provably beats independent noise

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

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



High
effective
dimension



Low
effective
dimension

Part 2: Near-optimal noise generation time complexity

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

Our approach: Approximate the noise coefficients as

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

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

Error: If $d = O(\log^2(n/c))$, then the error is $\mathcal{E}(\beta') \leq \mathcal{E}(\beta) + c$
(n = Number of steps)

Coming soon: Survey/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

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

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

ARUNGANESH@GOOGLE.COM

Abhradeep Thakurta
Google DeepMind

ATHAKURTA@GOOGLE.COM

Jalaj Upadhyay
Rutgers University

JALAJ.UPADHYAY@RUTGERS.EDU

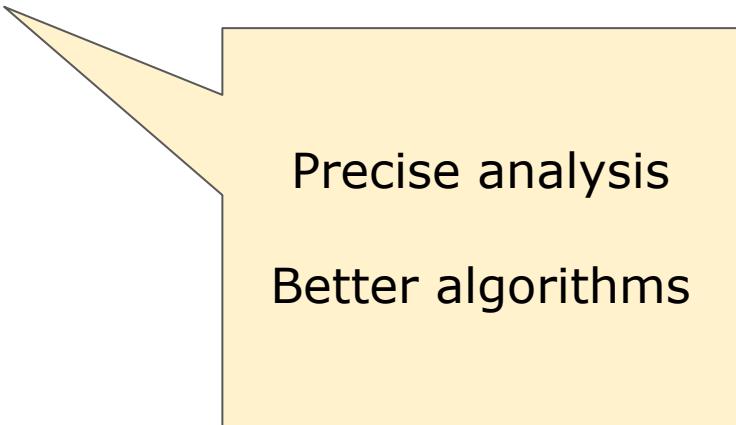
Precise analysis
(better rates)

Algorithm design

Open Problem: Multi-epoch Learning Guarantees

Assumptions in this talk:

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



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



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

Thank you!

Future Work

Theory

- Averaged iterate analysis + precise finite time bounds
- Analysis for non-Toeplitz systems

Ruppert. **Efficient Estimations from a Slowly Convergent Robbins-Monro Process.** 1998

Polyak and Juditsky. **Acceleration of Stochastic Approximation by Averaging.** SIAM J Control Optim. (1992)

Jain, Kakade, Kidambi, Netrapalli, Sidford. **Parallelizing Stochastic Gradient Descent for Least Squares Regression: Mini-batching, Averaging, and Model Misspecification.** JMLR (2018).

DP-FTRL: privatize prefix sums of gradients

SGD update (without noise)

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

$$\begin{pmatrix} g_0 \\ g_0 + g_1 \\ \vdots \\ g_0 + \cdots + g_{n-1} \end{pmatrix} = \boxed{\begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ \vdots & \vdots & \ddots & \\ 1 & 1 & \cdots & 1 \end{pmatrix}} \begin{pmatrix} g_0 \\ g_1 \\ \vdots \\ g_{n-1} \end{pmatrix}$$

Call this matrix as A

Gradient Descent with Linearly Correlated Noise: Theory and Applications to Differential Privacy

NeurIPS 2023

Anastasia Koloskova*
EPFL, Switzerland

Ryan McKenna
Google Research

Zachary Charles
Google Research

Keith Rush
Google Research

Brendan McMahan
Google Research

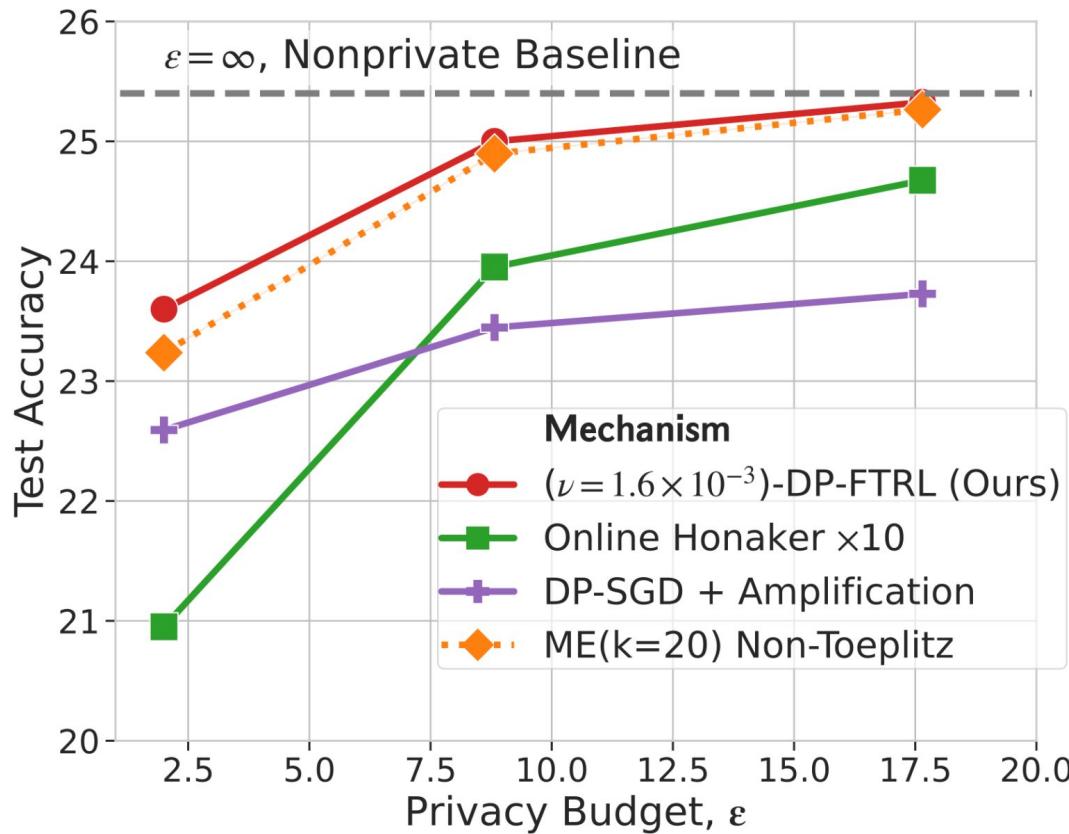
Theorem 4.7 (convex). *Under Assumptions 4.1, 4.2, and 4.3, if $\gamma \leq 1/4L$ and $\tau = \tilde{\Theta}(1/\gamma L)$, then (7) produces iterates with average error $(T+1)^{-1} \sum_{t=0}^T \mathbb{E}[f(\mathbf{x}_t) - f^*]$ upper bounded by*

$$\tilde{\mathcal{O}}\left(\frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{\gamma T} + \frac{\sigma^2}{TL\tau} \times \left[\frac{1}{\tau} \sum_{t=1}^T \left\| \mathbf{b}_t - \mathbf{b}_{\lfloor \frac{t}{\tau} \rfloor \tau} \right\|^2 + \sum_{t=0 \bmod \tau}^{1 \leq t \leq T} \left\| \mathbf{b}_t - \mathbf{b}_{t-\tau} \right\|^2 + \left\| \mathbf{b}_{\lfloor \frac{T}{\tau} \rfloor \tau} \right\|^2 \right] \right).$$

Improved analysis of DP-FTRL

No provable gap between DP-SGD & DP-FTRL (same as previous)

Empirical results for private language modeling



Ours
matches
SoTA!