

Correlated Noise Provably Beats Independent Noise for Differentially Private Learning

Christopher Choquette-Choo*, Krishnamurthy Dvijotham*, Krishna Pillutla*,
Arun Ganesh, Thomas Steinke, Abhradeep Thakurta

Google DeepMind
Google Research



Differentially Private Learning

Summary

How to add **noise** to SGD for differentially private learning?

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

Main result: (anti-) correlated noise is **provably** better

Standard Approach: SGD with **independent** noise

$$\theta_{t+1} = \theta_t - \eta \left(g_t + z_t \right) \quad \text{(DP-SGD)}$$

Learning Rate Clipped gradient s.t. $\|g_t\|_2 \leq G$ **Independent** Gaussian noise

For streaming data, take $z_t \sim \mathcal{N}(0, G^2/2\rho)$ for a desired level of privacy ρ zero-concentrated DP

Recent Work: SGD with **correlated** noise

[Kairouz et al. (ICML '21), Denisov et al. (NeurIPS '22), ...]

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

Sensitivity **Correlated** Gaussian noise

$$\text{Variance}(z_t) = \frac{G^2}{2\rho} \max_t \left\| [B^{-1}]_{:,t} \right\|_2^2$$

required for the same privacy

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

Empirically: Correlated noise \gg independent noise
[Choquette-Choo et al. (NeurIPS '23)]

Compute: Solve a semi-definite program for B

This work: *first* clear separation in theory!

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

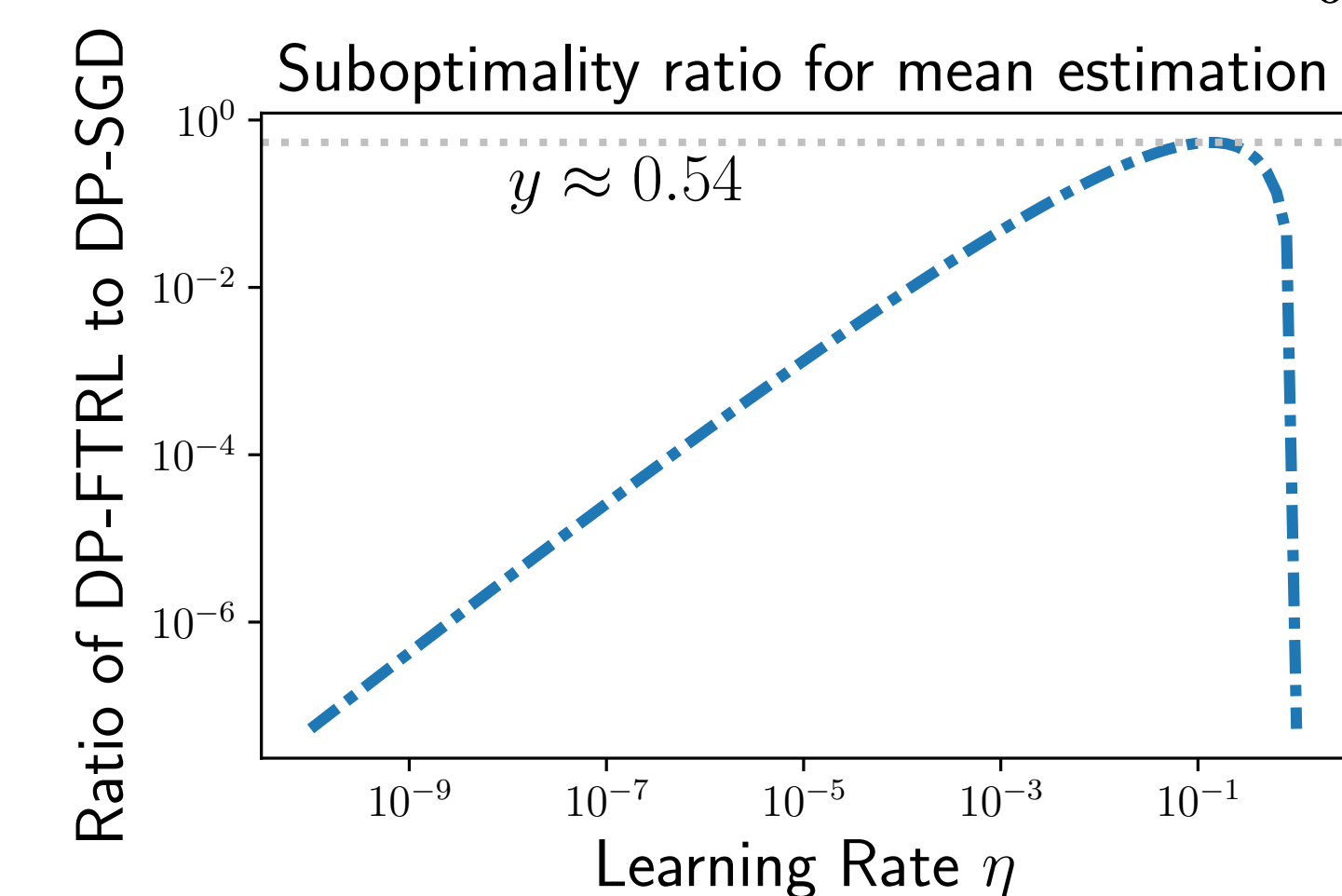
Theoretical Results

Mean Estimation in 1D with $f(\theta; x) = (\theta - x)^2$

Theorem: $F_{\infty}(\beta^{\text{sgd}}) = \rho^{-1} \eta$ **(DP-SGD)**

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

Optimal correlations are $\beta_0^{\star} = 1$, $\beta_t^{\star} = -t^{-3/2}(1-\eta)^t$.



DP-FTRL is always better

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

ν -DP-FTRL: use $\beta_t^{\nu} = -t^{-3/2}(1-\nu)^t$ with a tunable ν

Linear Regression

$f(\theta; x, y) = (y - \langle \theta, x \rangle)^2$ with $x \sim \mathcal{N}(0, H)$

Theorem: For algorithms without clipping, we have

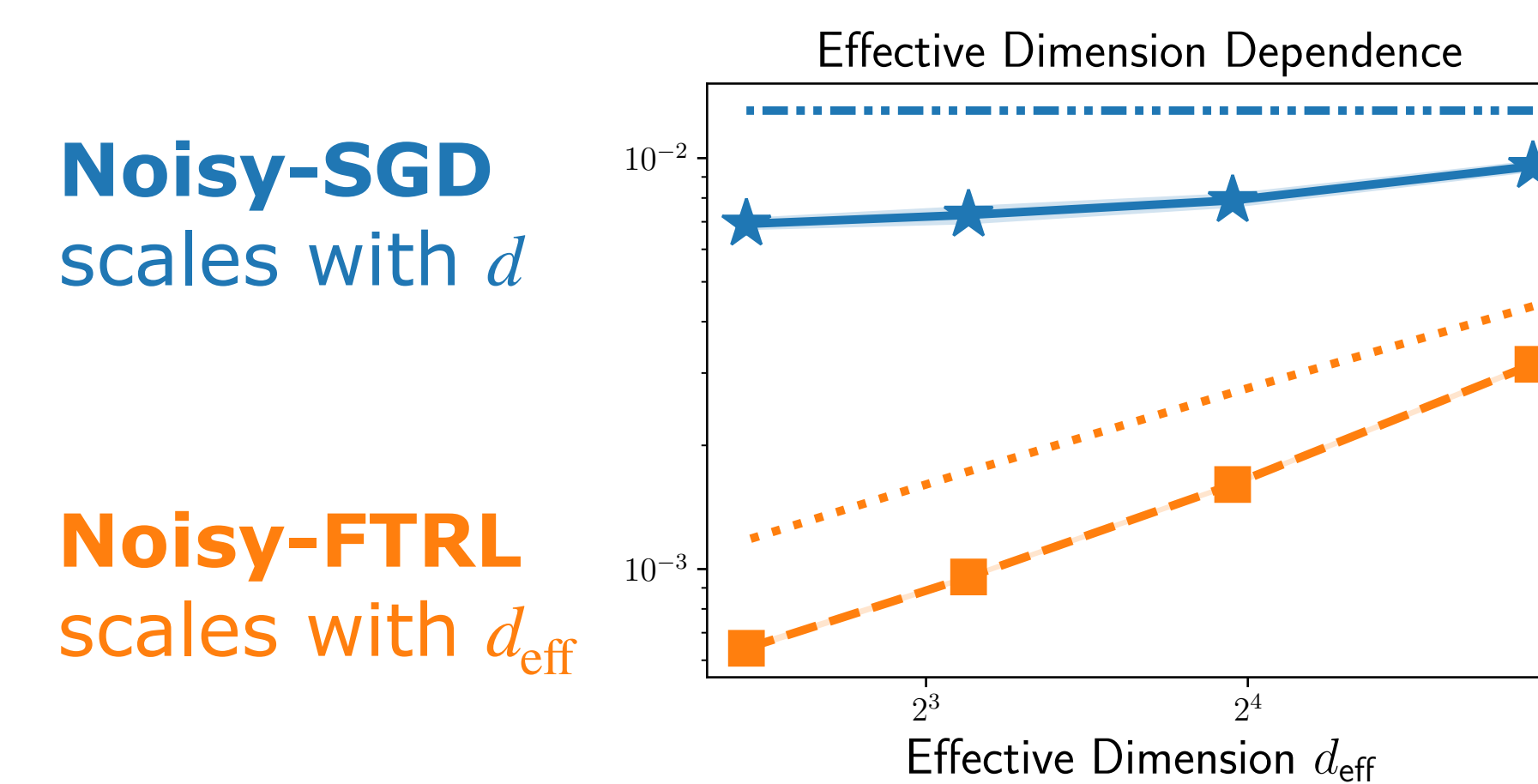
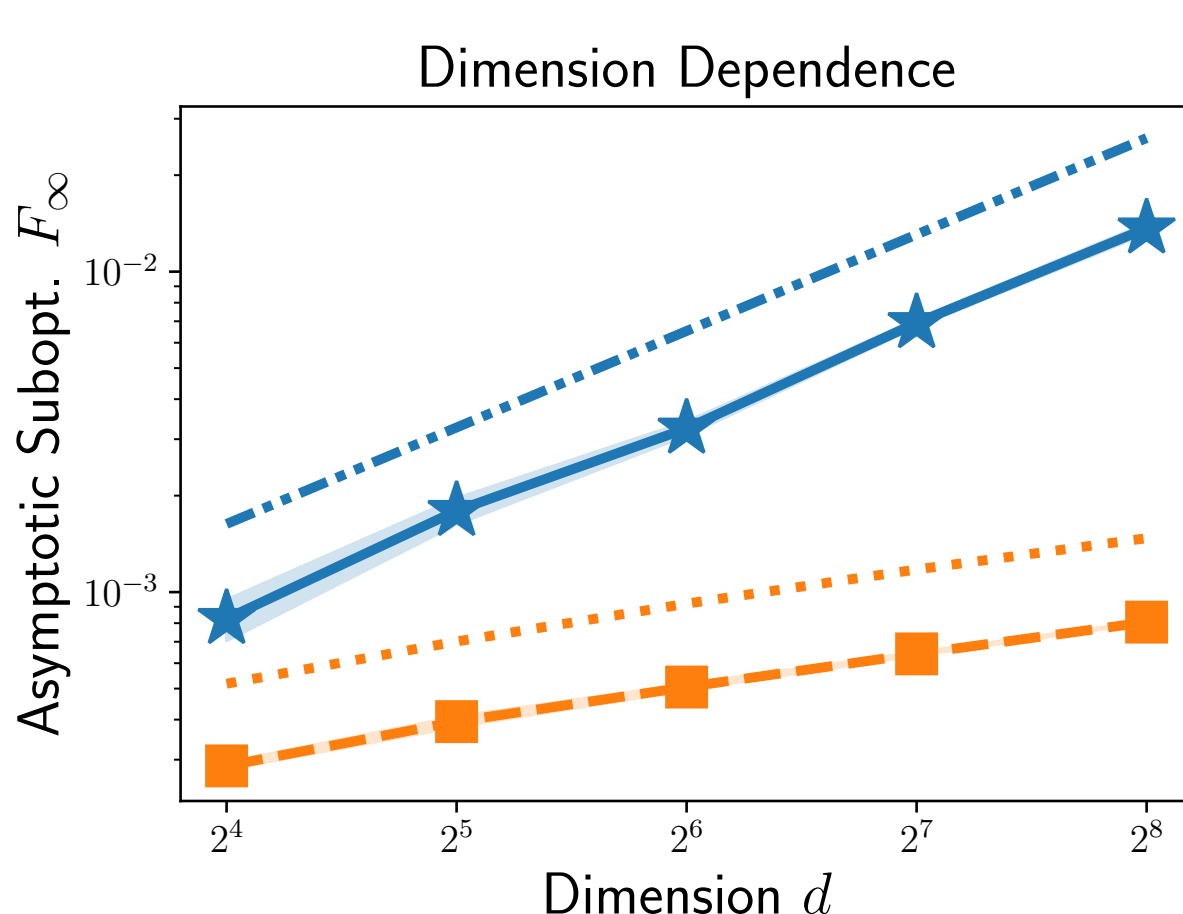
$F_{\infty}(\beta^{\text{sgd}}) = d \rho^{-1} \eta$ **(Noisy-SGD)**

$F_{\infty}(\beta^{\nu}) \leq d_{\text{eff}} \rho^{-1} \eta^2 \log^2 \frac{1}{\eta \mu}$ **(Noisy-FTRL)**

$\inf_{\beta} F_{\infty}(\beta) \geq d_{\text{eff}} \rho^{-1} \eta^2$ **(Lower bound)**

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

Near-optimal up to log factors

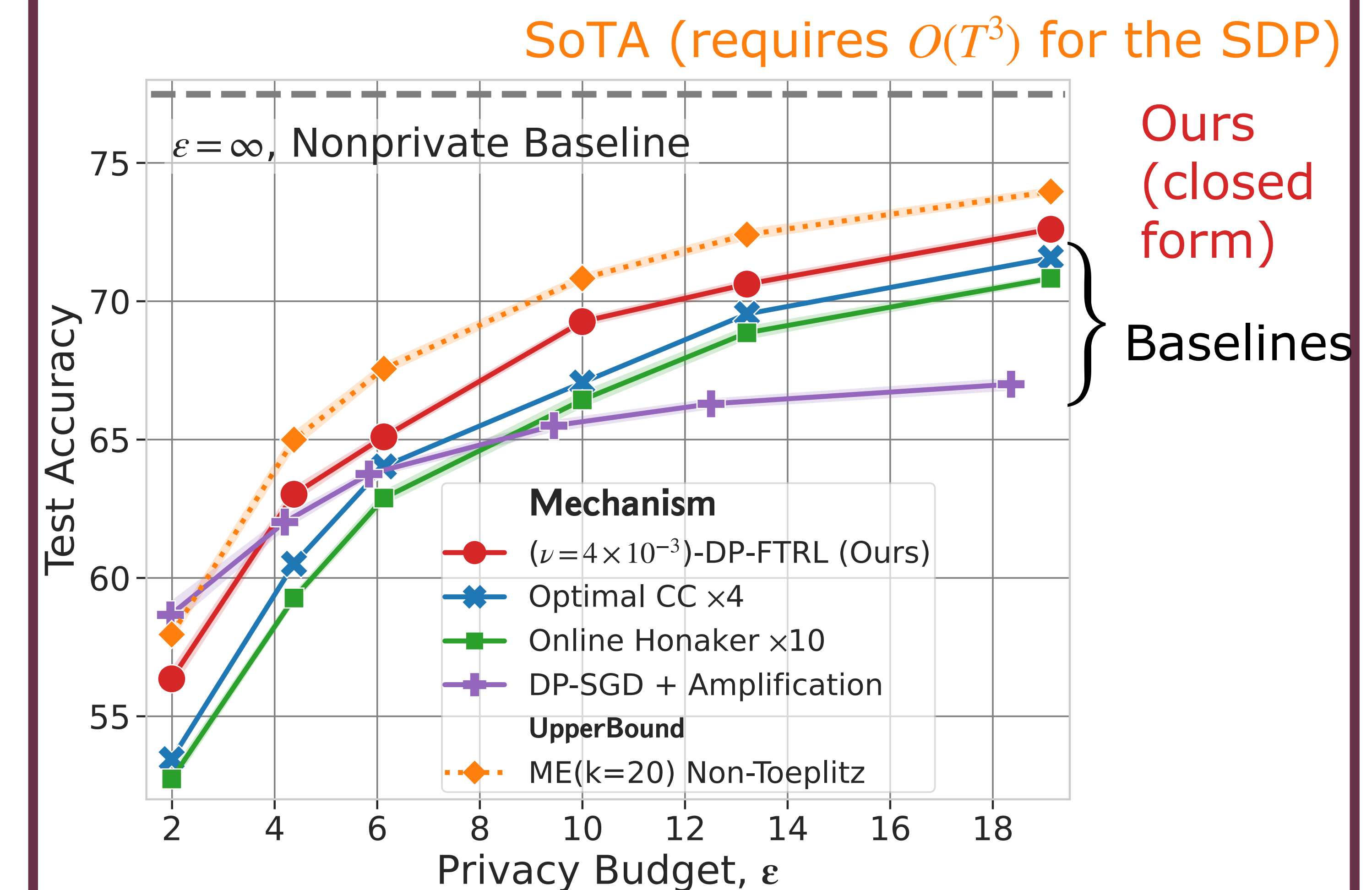


Challenge: Updates are no longer Markovian

Solution: Analysis in the Fourier domain

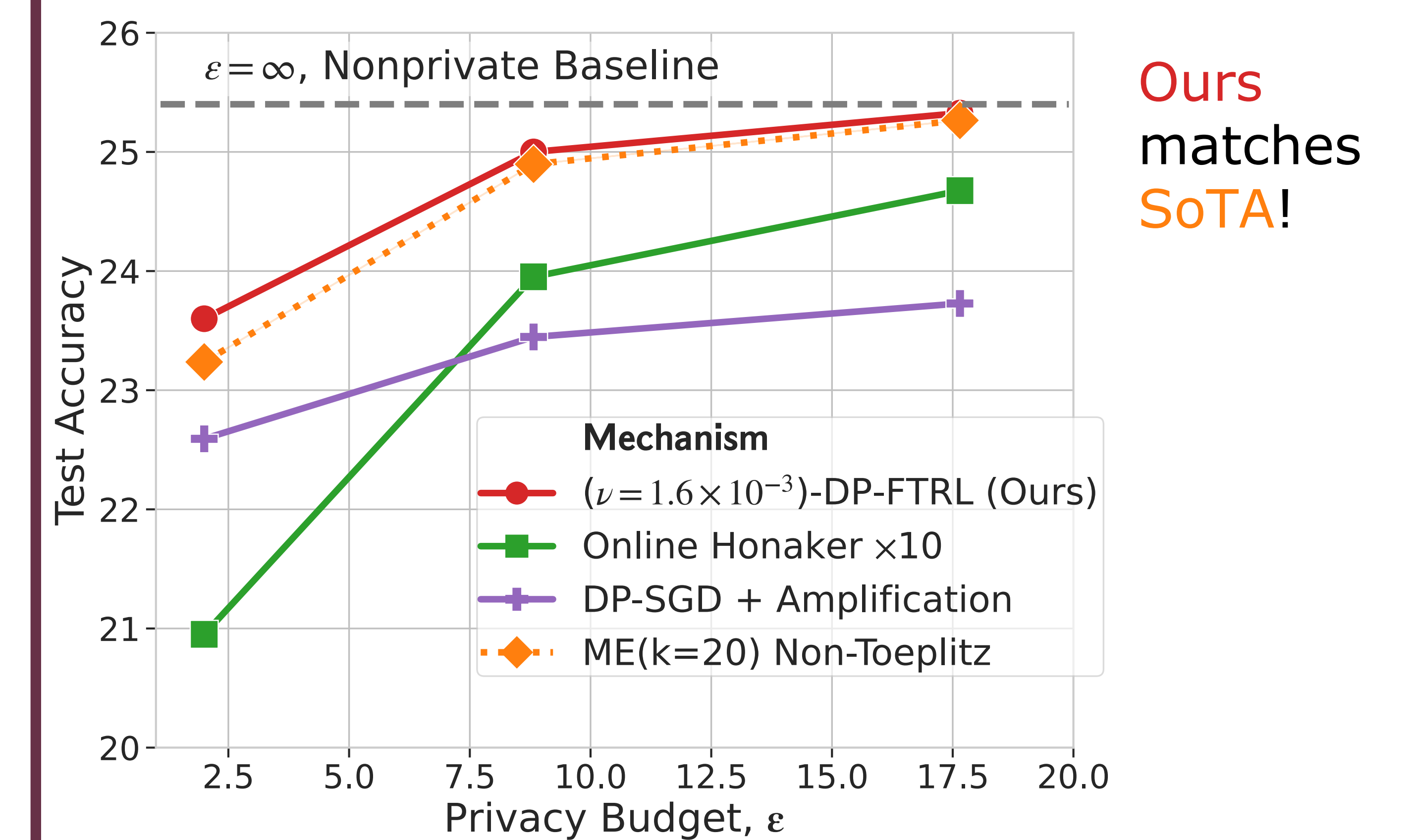
Experiments: Private Deep Learning

Image Classification (CIFAR-10)



Language Modeling (StackOverflow)

Federated learning + user-level DP



Extensions (see paper)

- Finite-time DP bounds for linear regression
- Bounds for general strongly convex functions



Arxiv link