

Multi-Epoch Matrix Factorization Mechanisms for Private Machine Learning

Christopher A. Choquette-Choo

H. Brendan McMahan, Keith Rush, Abhradeep Thakurta



Google DeepMind

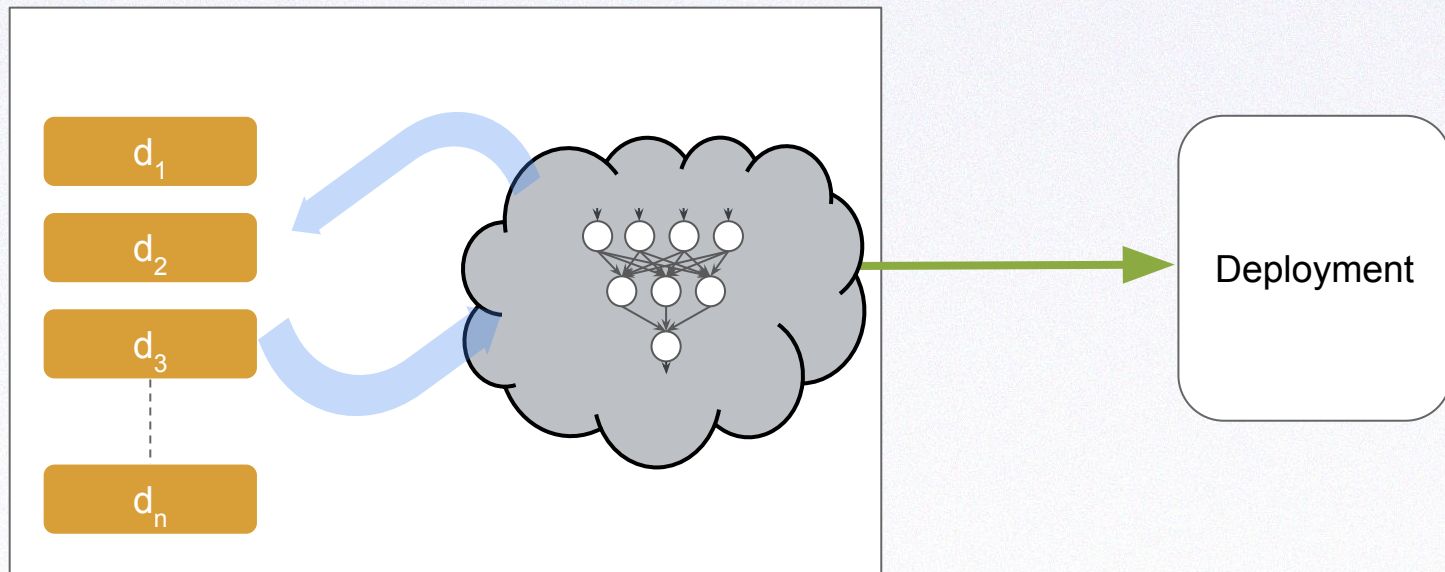
Google Research



www.christopherchoquette.com

@chris_choquette

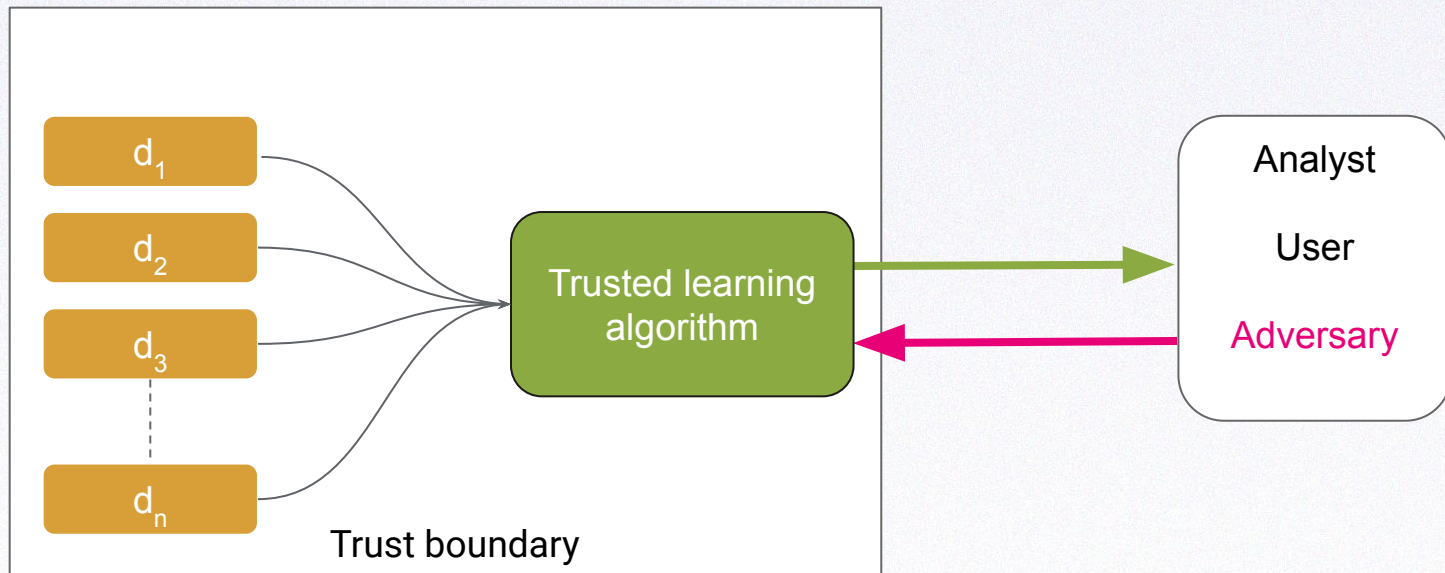
Standard Central Machine Learning Setup



Learn a **high utility** model via stochastic gradient descent (SGD).



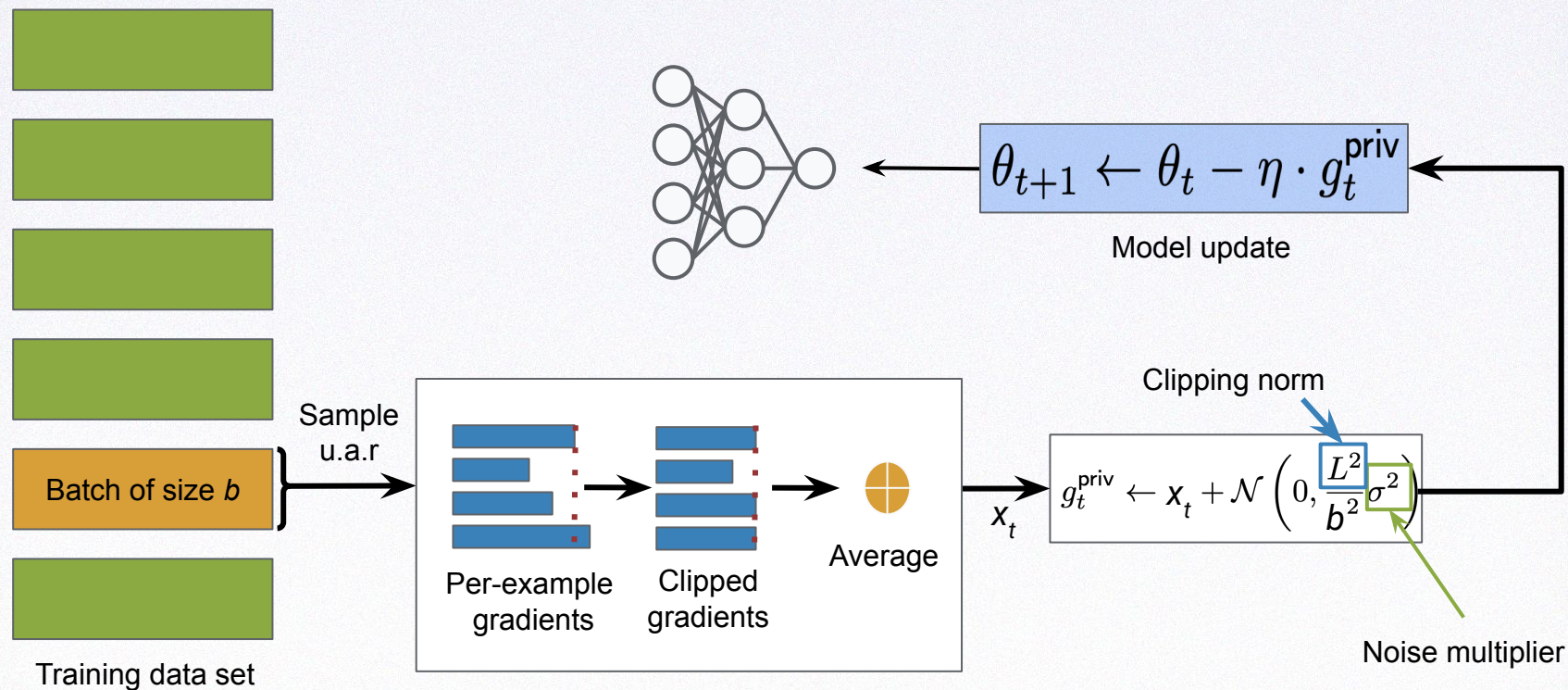
Our Goal: To do this with Differential Privacy (DP)



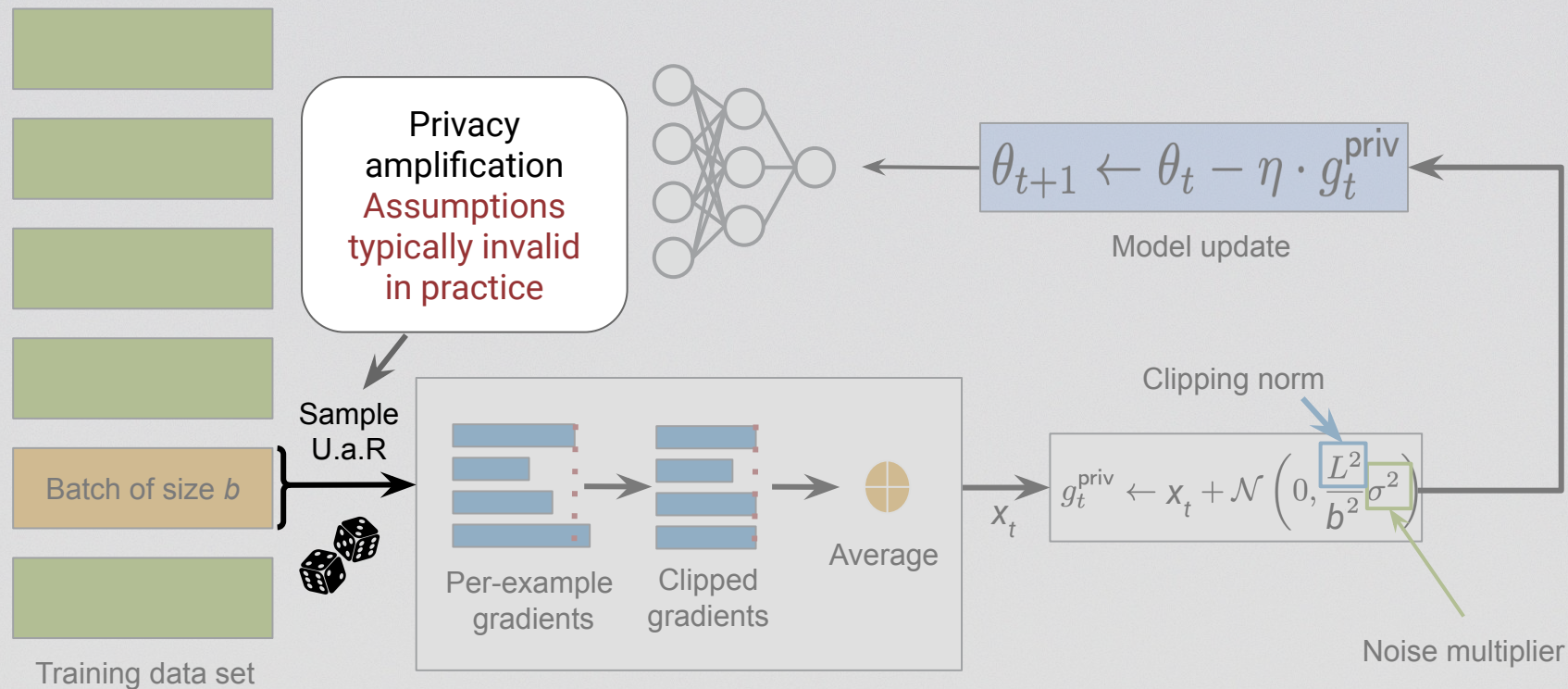
How can we maximize **utility** and **privacy** while minimizing **computation**



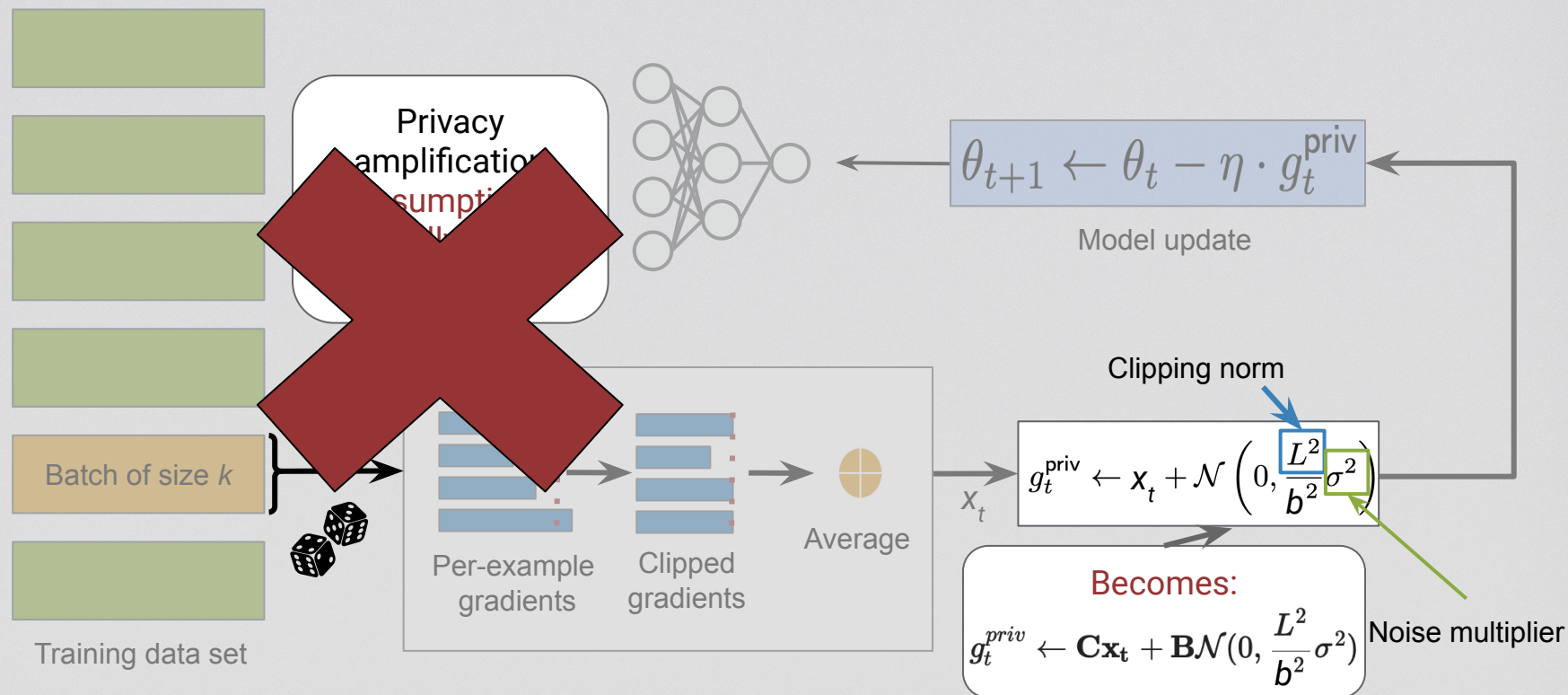
DP-SGD: Most Common Trusted Learning Algorithm



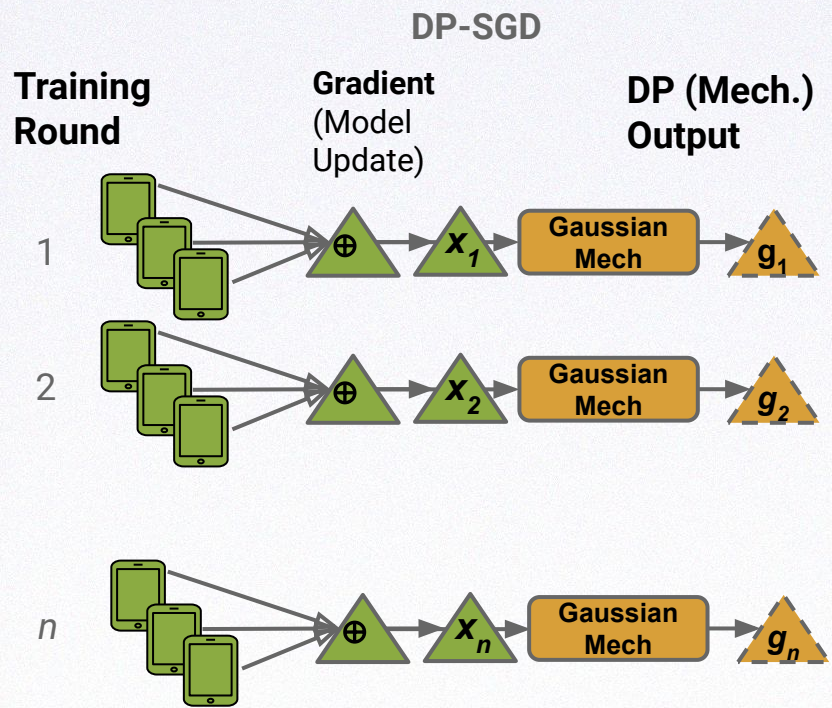
DP-SGD: Most Common Trusted Learning Algorithm



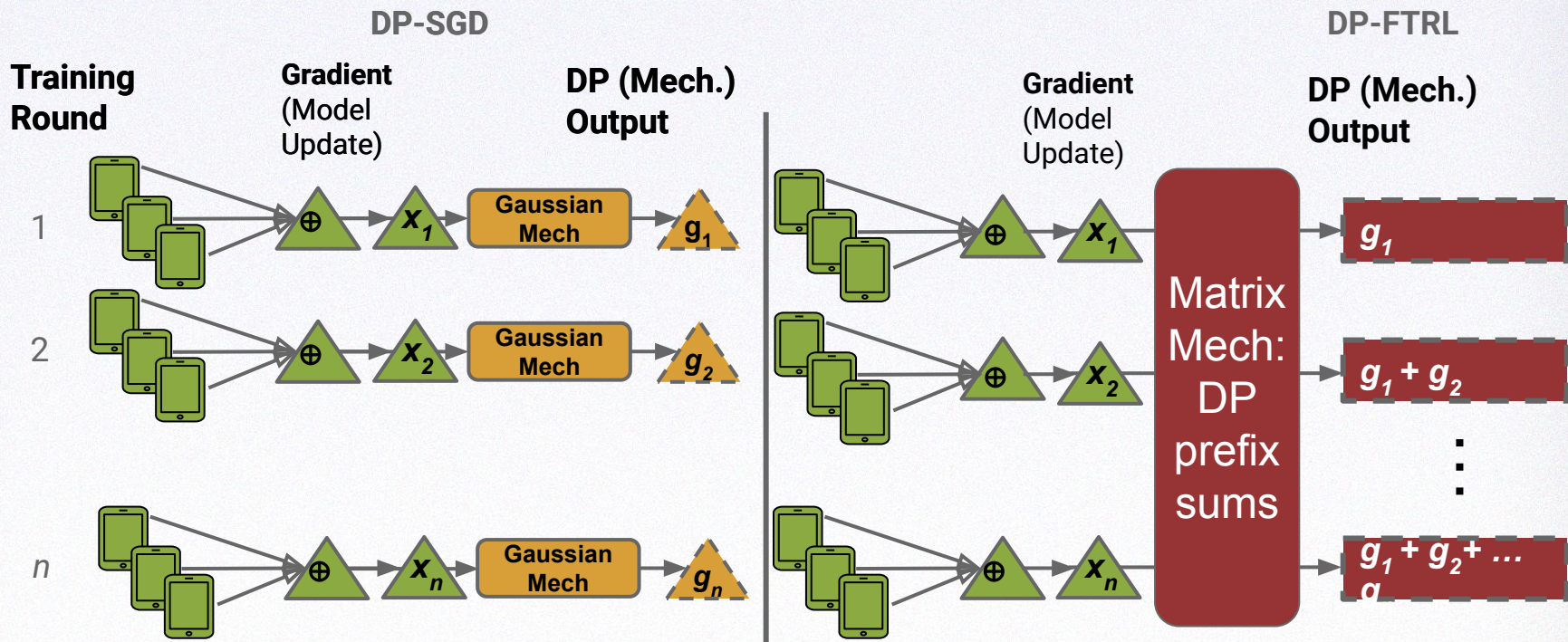
DP-SGD: Most Common Trusted Learning Algorithm



DP-SGD to DP-Follow the Regularized Leader (FTRL)



DP-SGD to DP-Follow the Regularized Leader (FTRL)



DP-FTRL Formulation and Matrix Factorization

Stream of **adaptively** chosen data vectors (a.k.a. gradients)

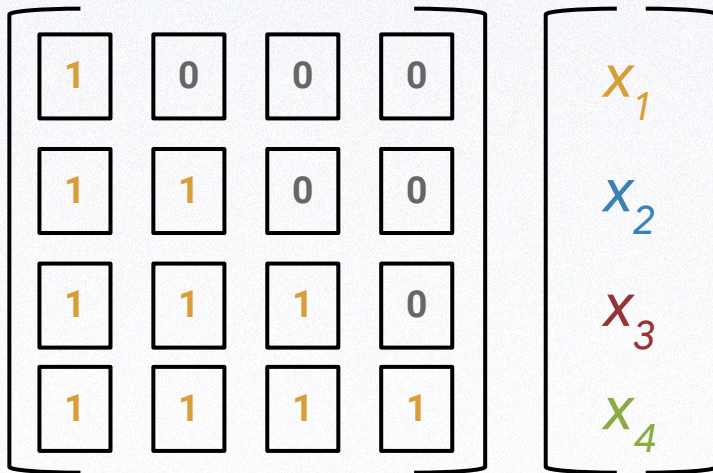
x_1 x_2 x_3 ... x_n

Query matrix: A

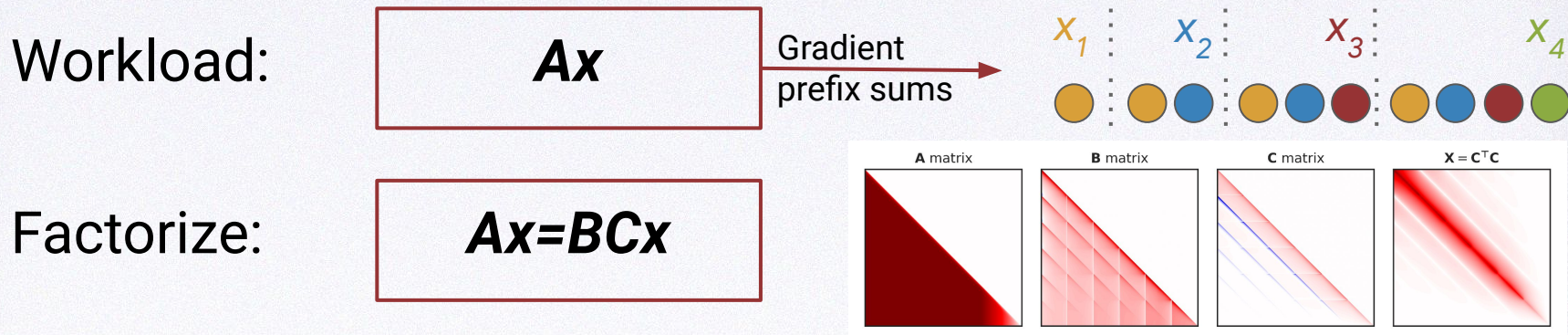
Data vector: x

Prefix sum as mat.-vec. product :

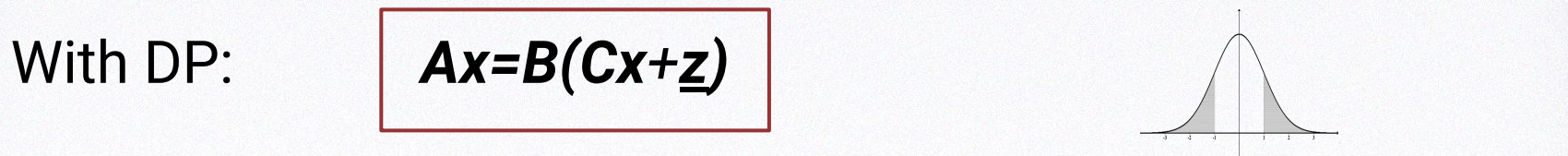
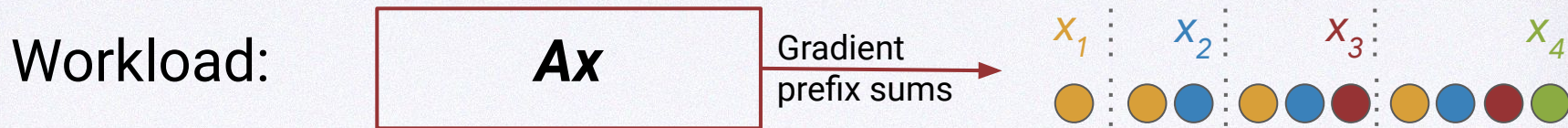
x_1 x_1+x_2 $x_1+x_2+x_3$ $x_1+x_2+x_3+x_4$



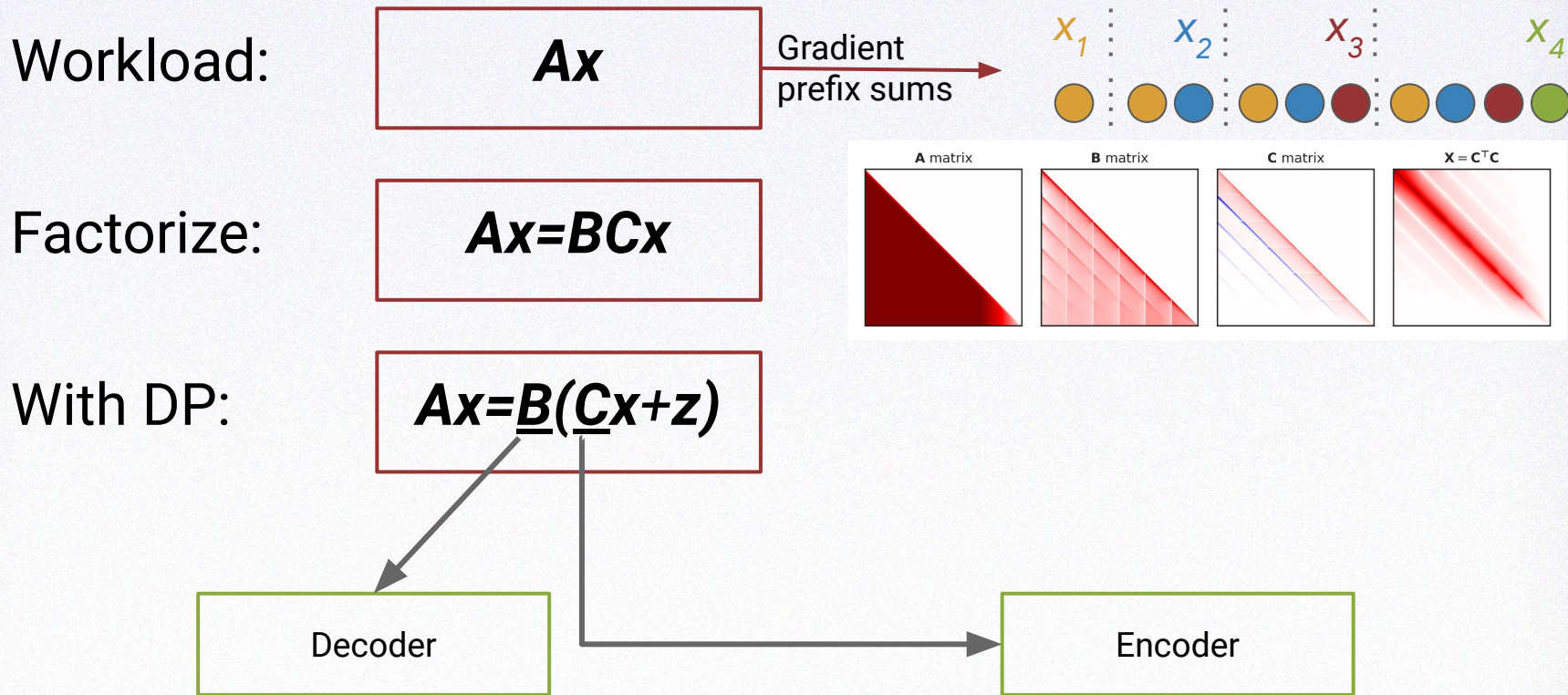
DP-FTRL Formulation and Matrix Factorization



DP-FTRL Formulation and Matrix Factorization



DP-FTRL Formulation and Matrix Factorization



What are B and C ?

Objective: factorize $A=BC$ such that...

$\|Bz\|_2$ is minimized

Decoder

$Cx+z$ satisfies DP

Encoder



What are \mathbf{B} and \mathbf{C} ?

Objective: factorize $\mathbf{A}=\mathbf{BC}$ such that...

$$\min_{\mathbf{C}} \|\mathbf{A}\mathbf{C}^\dagger\|_F^2 \text{ where our loss is } \text{sens}_{\mathcal{D}}^2(\mathbf{C}) \|\mathbf{B}\|_F^2$$
$$\text{sens}(\mathbf{C}) = 1; \quad \text{as } \forall \alpha, \mathcal{L}(\mathbf{B}, \mathbf{C}) = \mathcal{L}(\alpha\mathbf{B}, \mathbf{C}/\alpha)$$

$\|\mathbf{Bz}\|_2$ is minimized

Decoder

$\mathbf{Cx}+\mathbf{z}$ satisfies DP

Encoder



What are B and C ?

Objective: factorize $A=BC$ such that...

$$\min_{\mathbf{C}} \|\mathbf{A}\mathbf{C}^\dagger\|_F^2 \text{ where our loss is } \text{sens}_{\mathcal{D}}^2(\mathbf{C}) \|\mathbf{B}\|_F^2$$

$\text{sens}(\mathbf{C}) = 1;$ as $\forall \alpha, \mathcal{L}(\mathbf{B}, \mathbf{C}) = \mathcal{L}(\alpha\mathbf{B}, \mathbf{C}/\alpha)$

$\|\mathbf{Bz}\|_2$ is minimized

Decoder

$\mathbf{Cx}+\mathbf{z}$ satisfies DP

Encoder



Bound sens(**Cx**)

Desiderata

1. **Tight:** Leave nothing on the table.
2. **Efficient:** Or optimization of the factorization will be slow.
3. **Practical:** Can be implemented easily to not break assumptions/guarantees.

Challenges

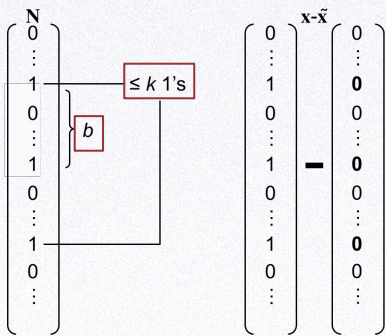
1. Sensitivity for scalar contributions does not always extend to vector contributions.
2. In general, sensitivity is NP-hard to compute.



Bound sens($\mathbf{C}\mathbf{x}$)

define $\mathcal{D}_{\Pi}^d = \{\mathbf{x} - \tilde{\mathbf{x}} \mid (\mathbf{x}, \tilde{\mathbf{x}}) \in \mathcal{N}\}$, set of all participation patterns.

$$\text{sens}_{\mathcal{D}}(\mathbf{C}) = \sup_{(\mathbf{x}, \tilde{\mathbf{x}}) \in \mathcal{N}} \|\mathbf{C}\mathbf{x} - \mathbf{C}\tilde{\mathbf{x}}\|_F = \sup_{\mathbf{u} \in \mathcal{D}} \|\mathbf{C}\mathbf{u}\|_F$$



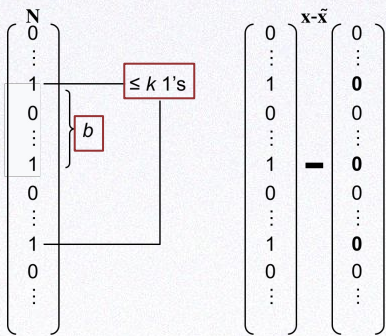
- Formalize participation patterns
 - When users are allowed to participate.
 - Enforced via data pipelines.
- Require a notion of “neighbouring streams”
 - Zero-out contributions of an example/user
 - The delta must still be a valid participation pattern



Bound sens(Cx)

define $\mathcal{D}_{\Pi}^d = \{\mathbf{x} - \tilde{\mathbf{x}} \mid (\mathbf{x}, \tilde{\mathbf{x}}) \in \mathcal{N}\}$, set of all participation patterns.

$$\text{sens}_{\mathcal{D}}(\mathbf{C}) = \sup_{(\mathbf{x}, \tilde{\mathbf{x}}) \in \mathcal{N}} \|\mathbf{C}\mathbf{x} - \mathbf{C}\tilde{\mathbf{x}}\|_F = \sup_{\mathbf{u} \in \mathcal{D}} \|\mathbf{C}\mathbf{u}\|_F$$



- Formalize participation patterns
 - When users are allowed to participate.
 - Enforced via data pipelines.
- Require a notion of “neighbouring streams”
 - Zero-out contributions of an example/user
 - The delta must still be a valid participation pattern
- Identify and propose $(\underline{k}, \underline{b})$ -participation
 - Participate at most \underline{k} times, exactly \underline{b} steps apart.

Practical

- Shuffle just once prior to training.

Efficient

- Requires only $\mathcal{O}(k^2b)$ time.

Tight

- Bounds exactly the k contributions per data point.



What are \mathbf{B} and \mathbf{C} ?

Objective: factorize $\mathbf{A}=\mathbf{BC}$ such that...

$$\min_{\mathbf{C}} \|\mathbf{A}\mathbf{C}^\dagger\|_F^2 \text{ where our loss is } \text{sens}_{\mathcal{D}}^2(\mathbf{C}) \|\mathbf{B}\|_F^2$$
$$\text{sens}(\mathbf{C}) = 1; \quad \text{as } \forall \alpha, \mathcal{L}(\mathbf{B}, \mathbf{C}) = \mathcal{L}(\alpha\mathbf{B}, \mathbf{C}/\alpha)$$

$\|\mathbf{Bz}\|_2$ is minimized

Decoder

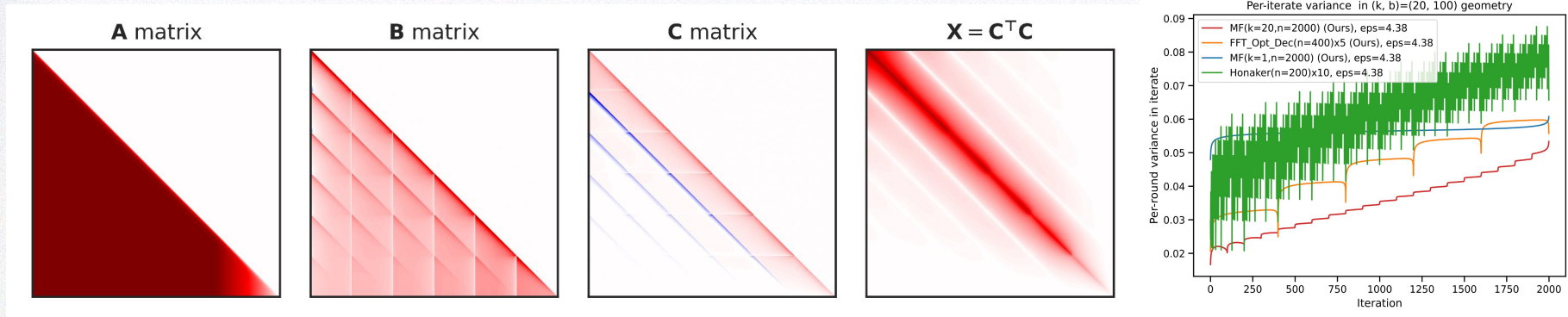
$\mathbf{Cx}+\mathbf{z}$ satisfies DP

Encoder



$\|Bz\|_2$ is minimized

- We define a mathematical program for minimizing $\|Bz\|_2$.
- We solve for the lagrangian and dual functions in closed-form.
- We solve for the gradient of the dual.
 - Allows us to directly leverage prior fixed-point iteration optimization algorithms to optimize B, C .
- Noise added is lower than all prior algorithms.

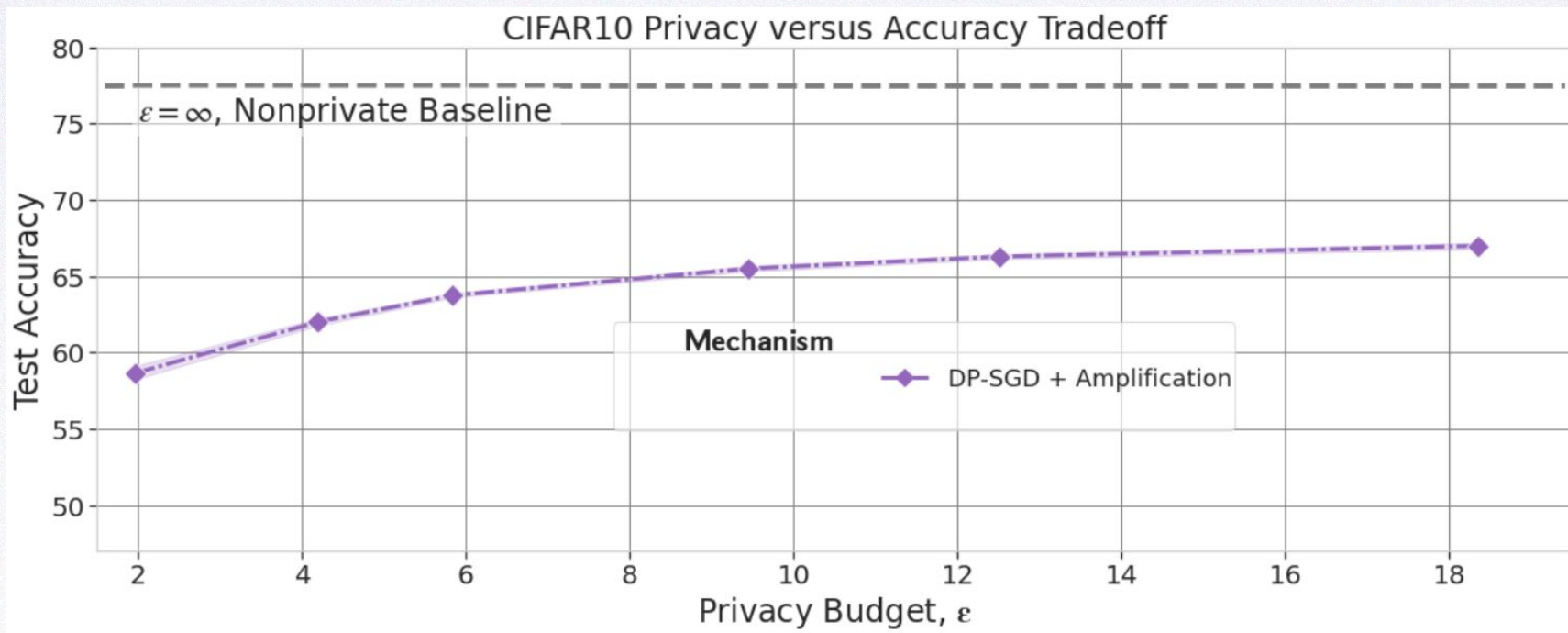


Empirical Details

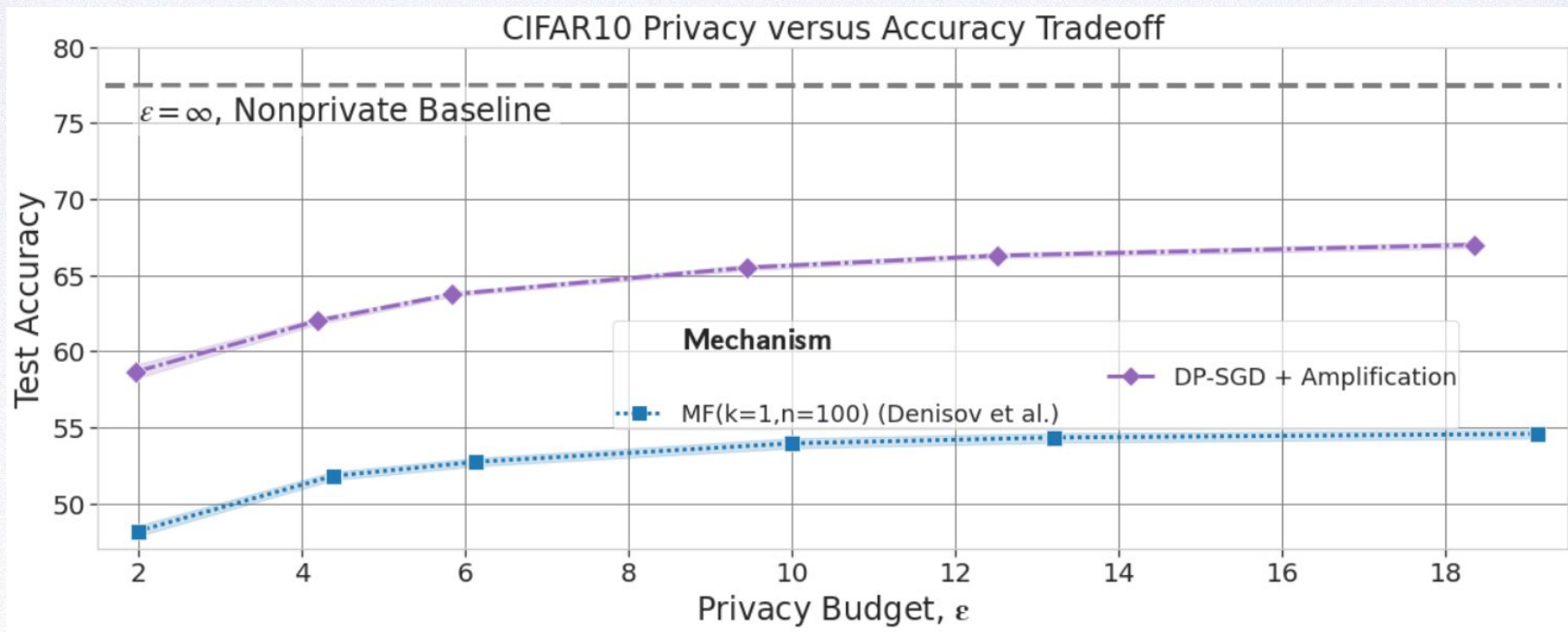
- Focus on CIFAR-10 (see paper for Stack Overflow next work prediction)
- Train for $k=20$ epochs, this gives $b=100$
- Factorize \mathbf{B}, \mathbf{C} matrices for (20,100)-participation
 - We also propose two other methods: FFT and stamping. See the paper for more details!
- Scale as $\alpha \mathbf{B}, 1/\alpha \mathbf{C}$ so that $\text{sens}(\mathbf{C}) = 1$
- Choose noise standard deviation σ so that a single Gaussian event achieves (ϵ, δ) -DP.
- Generate per-step noise as $(\mathbf{A}\mathbf{C}^{-1})_{[:,i]} \mathcal{N}(0, \sigma)$



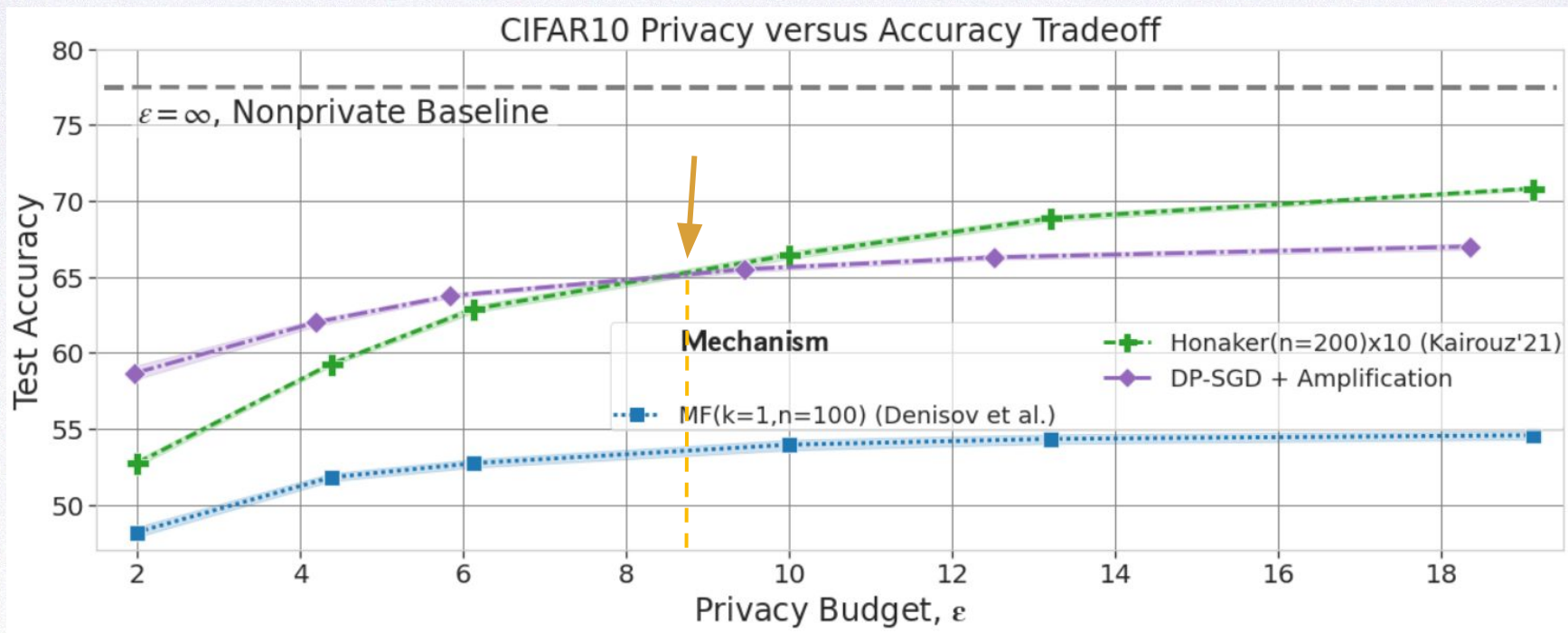
CIFAR-10



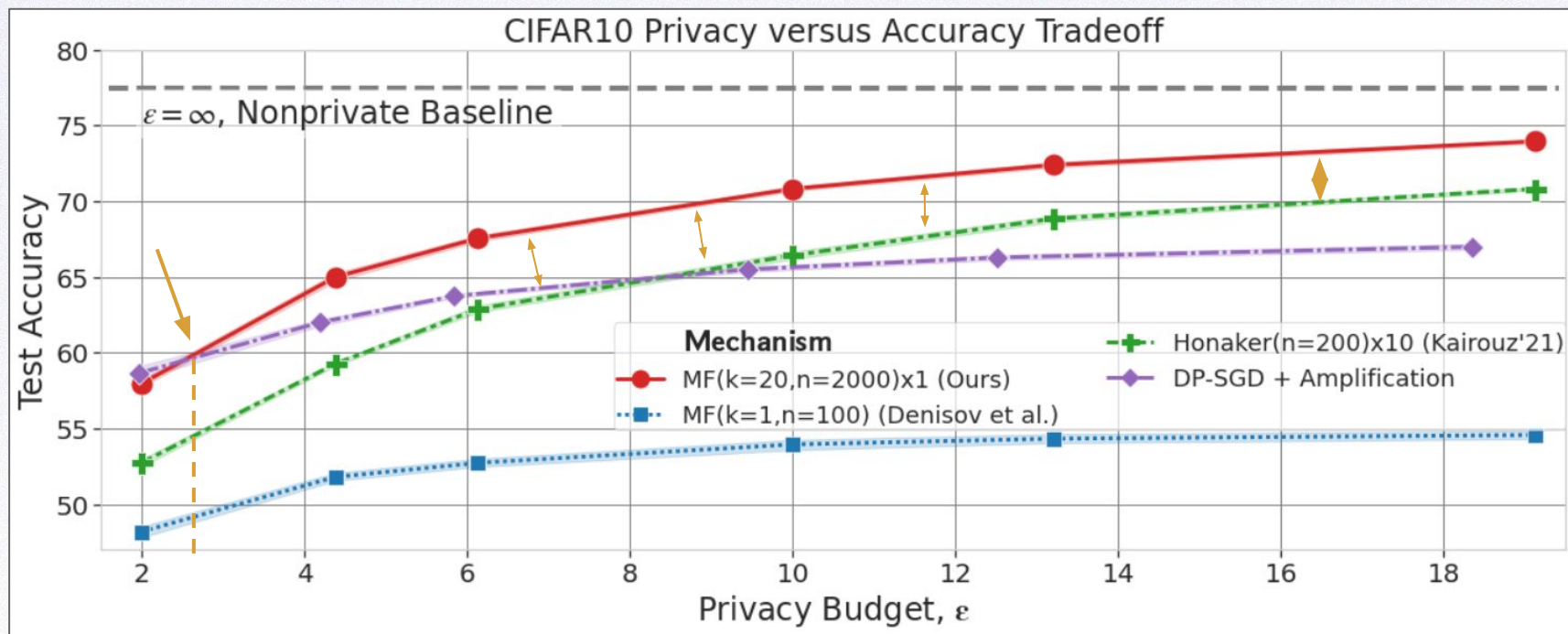
CIFAR-10



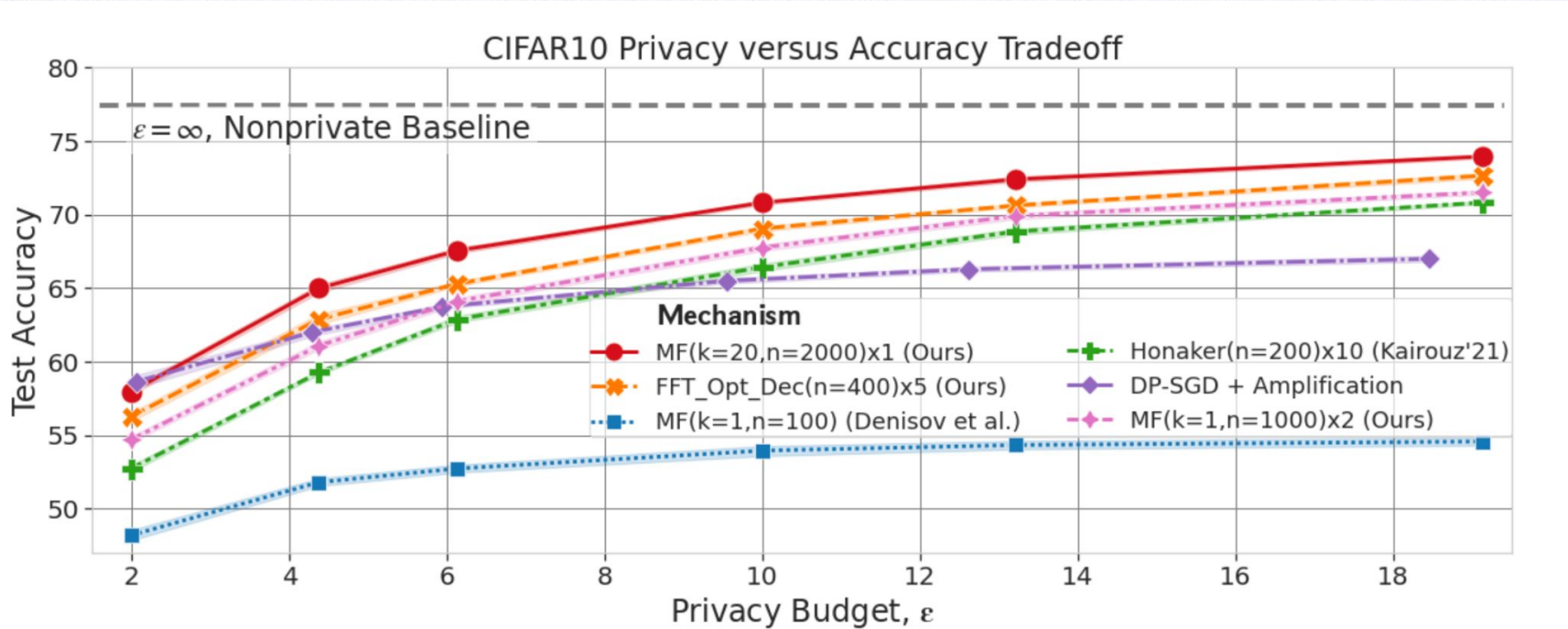
CIFAR-10



CIFAR-10



CIFAR-10



Conclusion

- New differentially private mechanisms.
- Achieve state-of-the-art privacy-utility tradeoffs (~5% points above DP-SGD).
- Realistic assumptions that can be implemented in practice.
- Applicable with tricks used for DP-SGD (no changes other than minimizing noise)
- Computationally practical

www.christopherchoquette.com

@chris_choquette

