

# Communication-Privacy-Accuracy Tradeoffs under Distributed DP for FL

University of California Santa Barbara

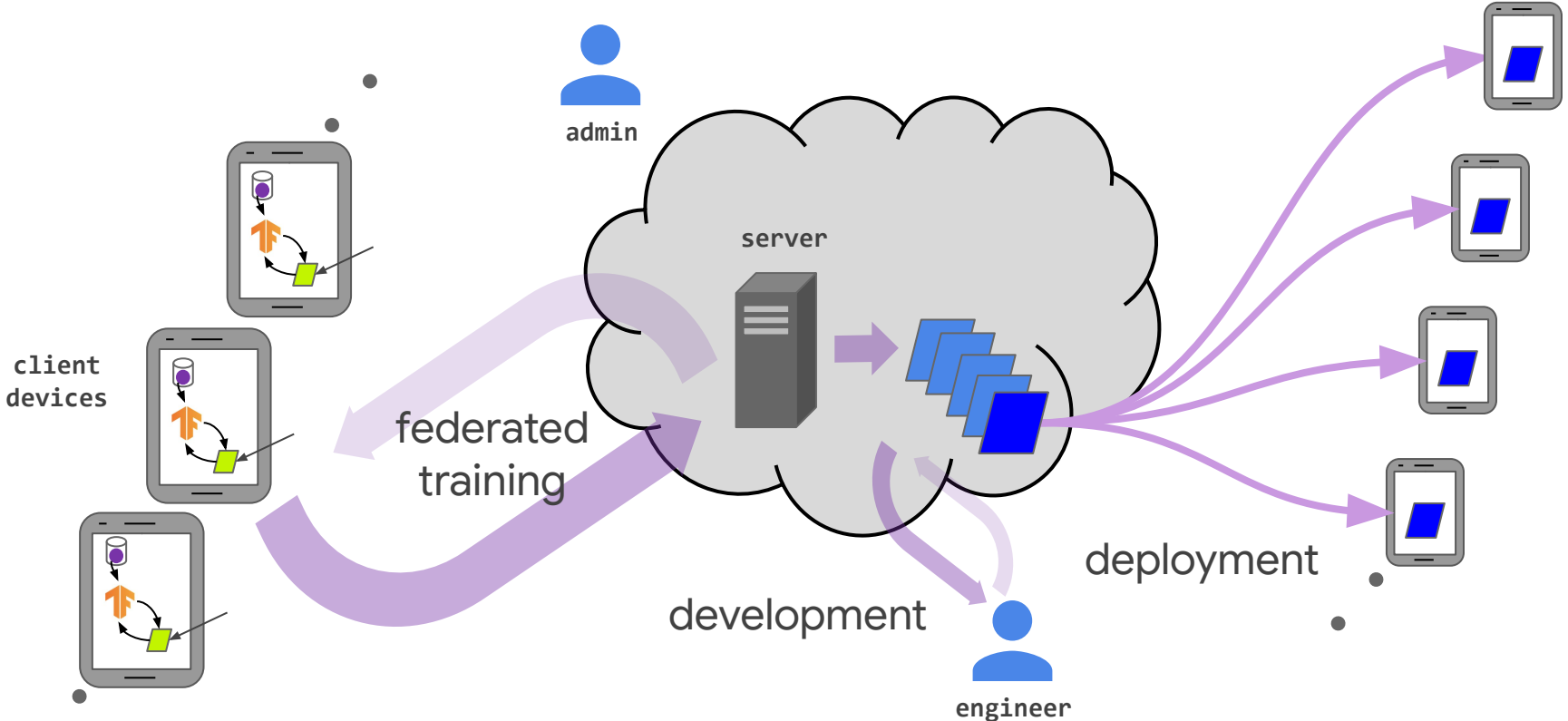


**Peter Kairouz**  
Research Scientist | Google

Twitter: @KairouzPeter

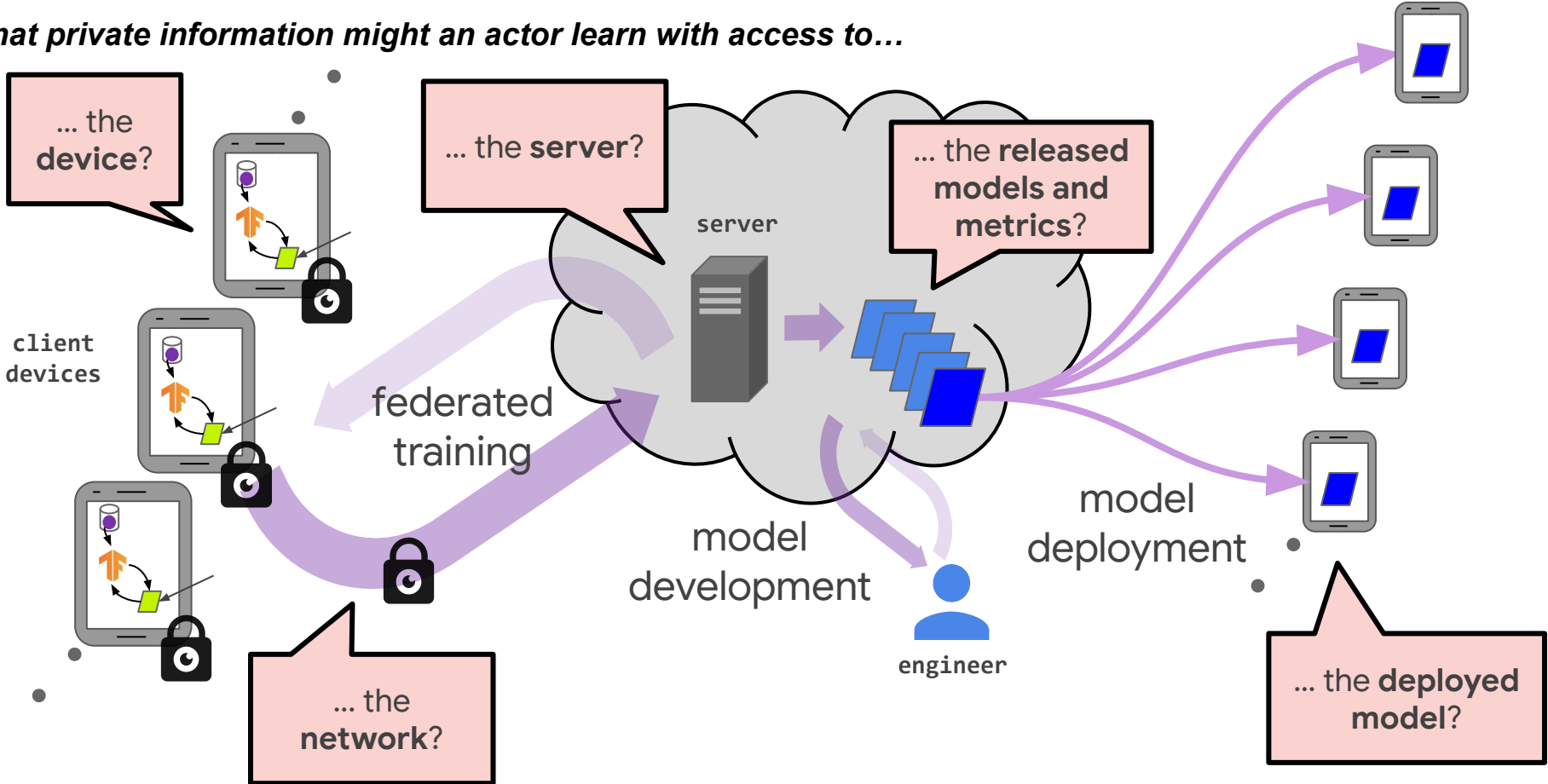
Presenting joint work with Ken Liu, Thomas Steinke, Naman Agarwal, and others

# Federated Learning

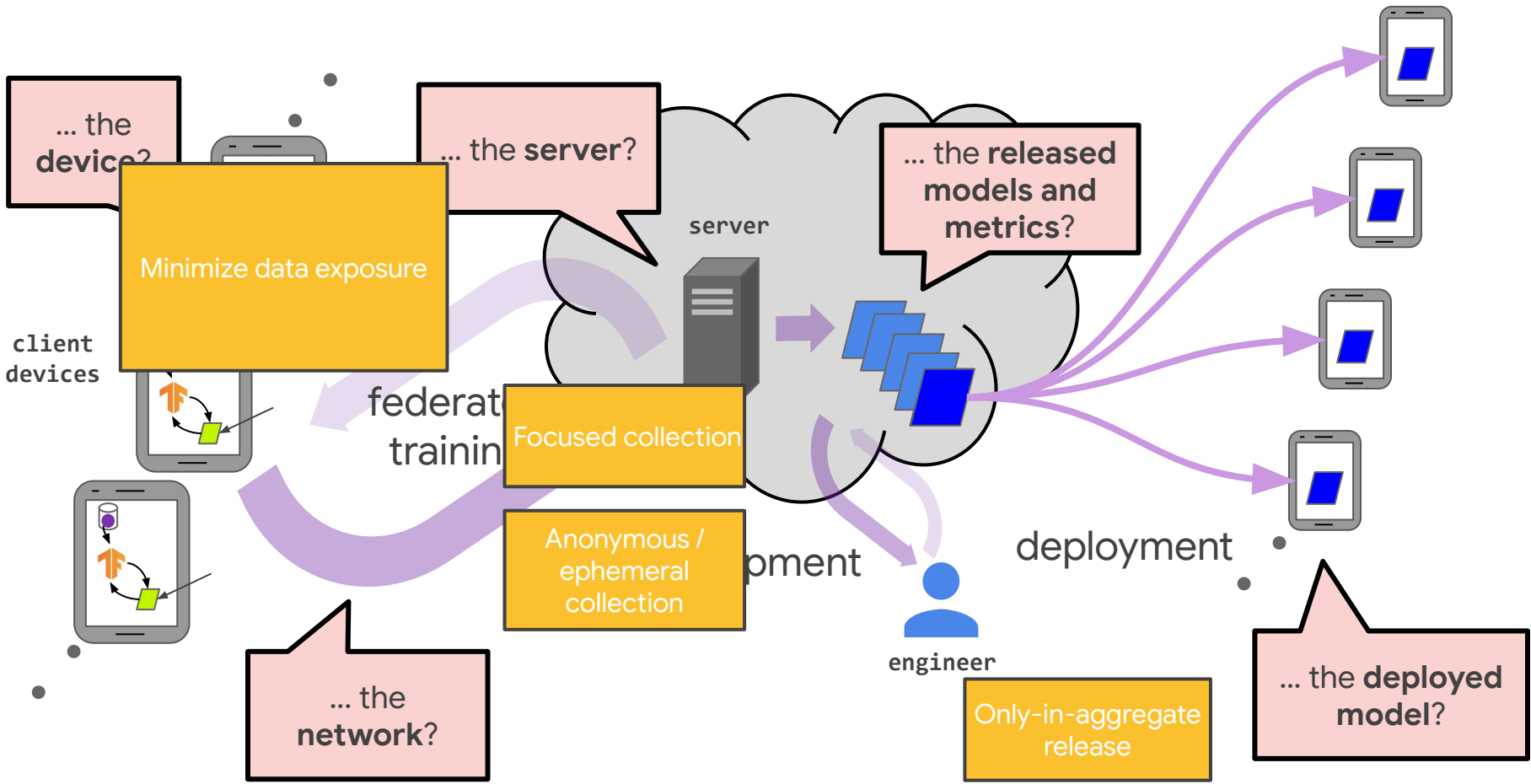


# Ensuring privacy of participating users

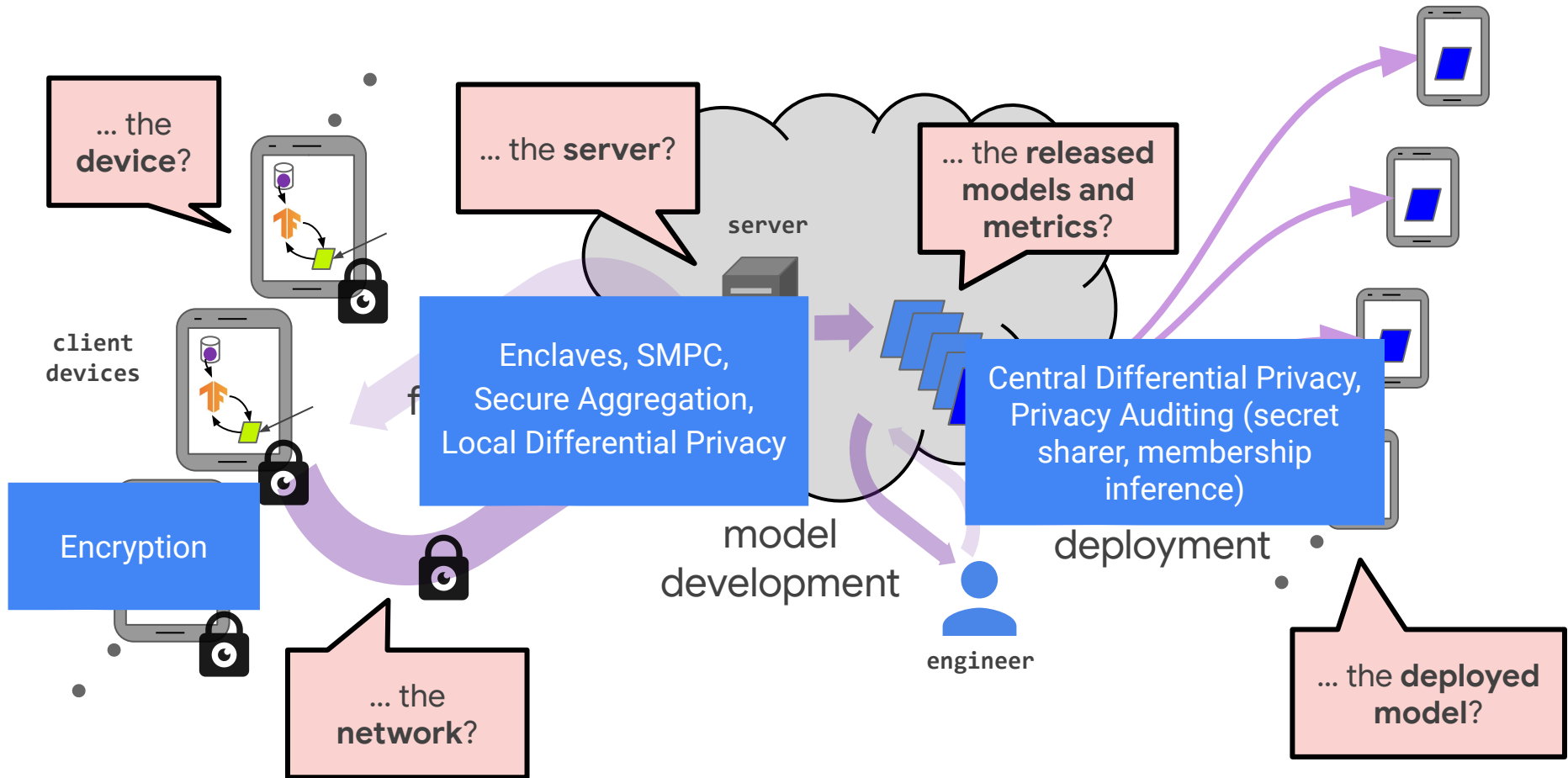
*What private information might an actor learn with access to...*



# Data minimization principles for FL

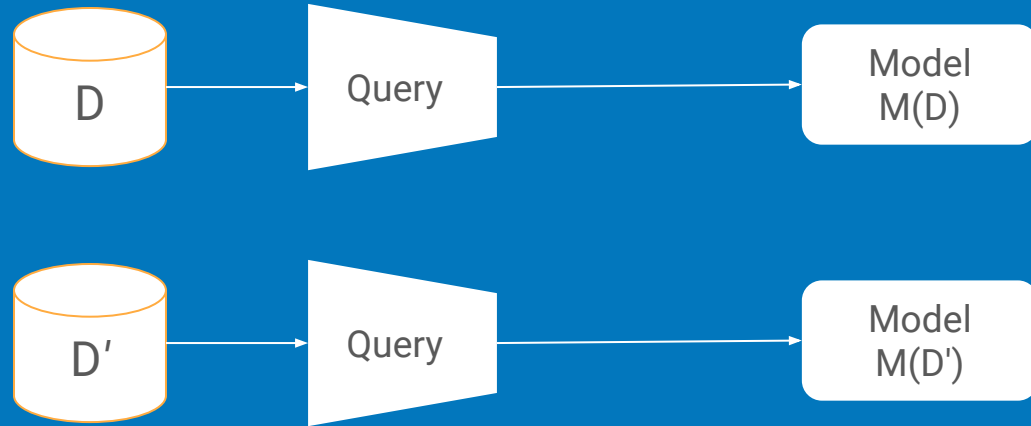


# Complementary privacy technologies



# Differential Privacy for FL

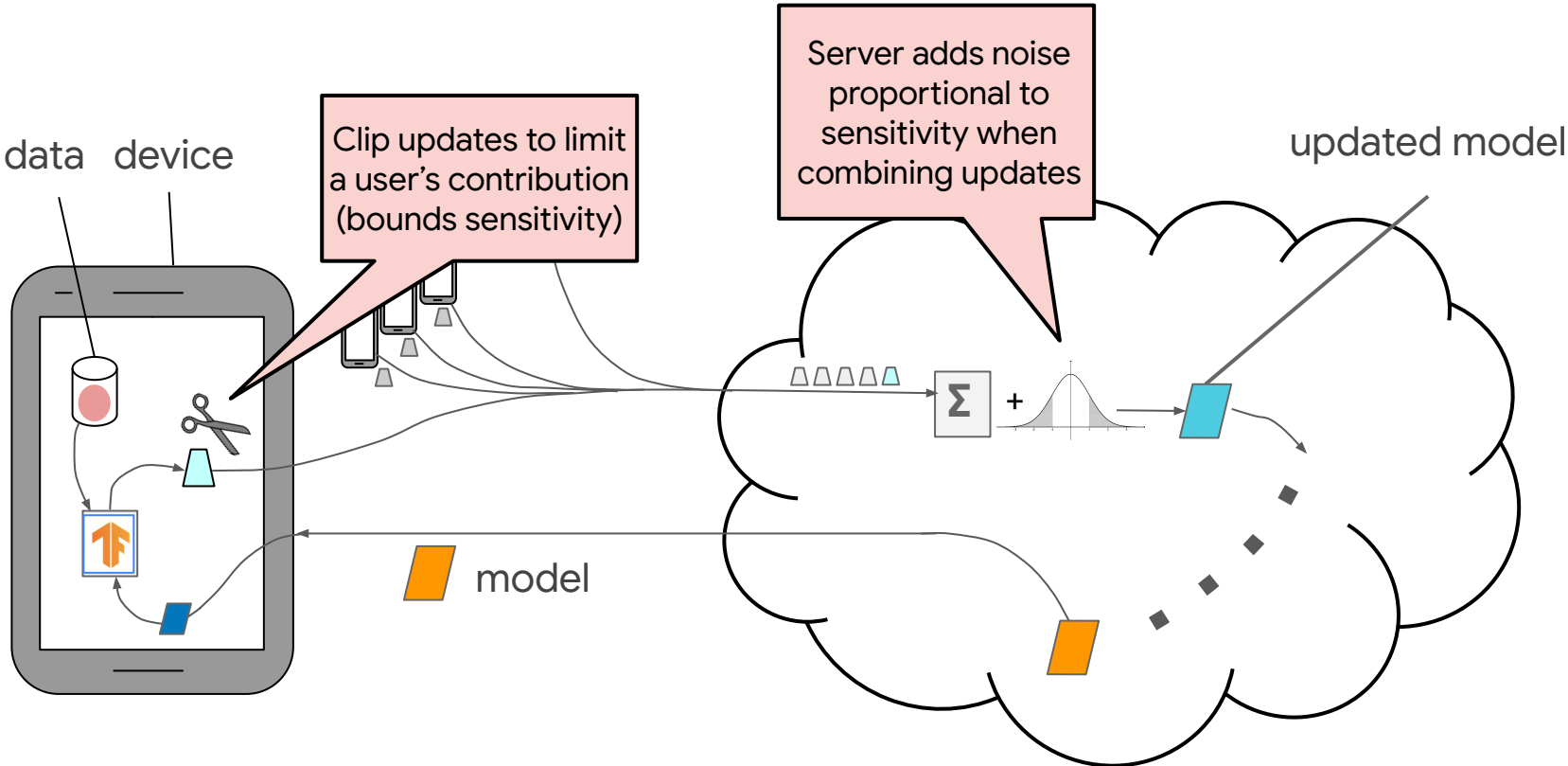
# Differential Privacy



**$(\epsilon, \delta)$ -Differential Privacy:** The distribution of the output  $M(D)$  (a trained model) on database (training dataset)  $D$  is **nearly the same** as  $M(D')$  for all adjacent databases  $D$  and  $D'$

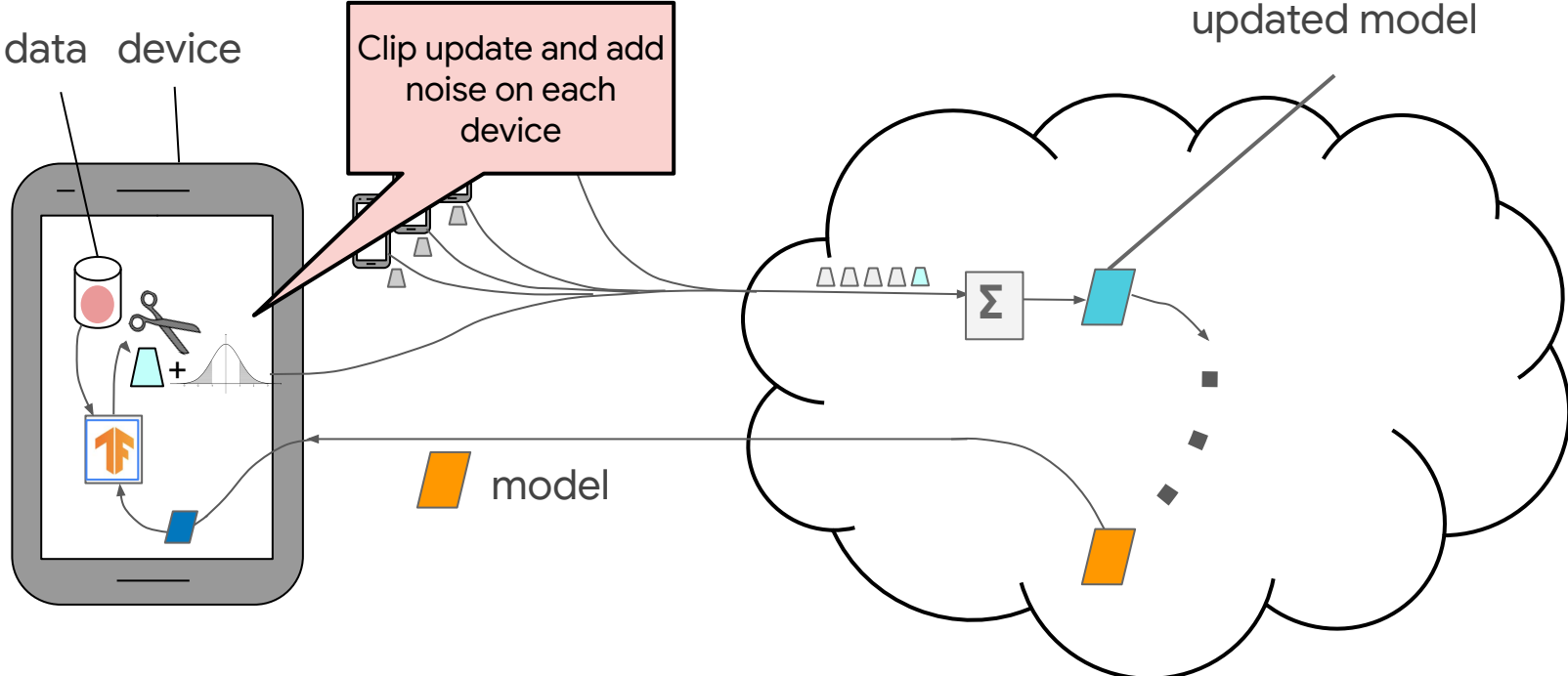
$$\forall S: \Pr[M(D) \in S] \leq \exp(\epsilon) \cdot \Pr[M(D') \in S] + \delta$$

# Centrally differentially private federated learning





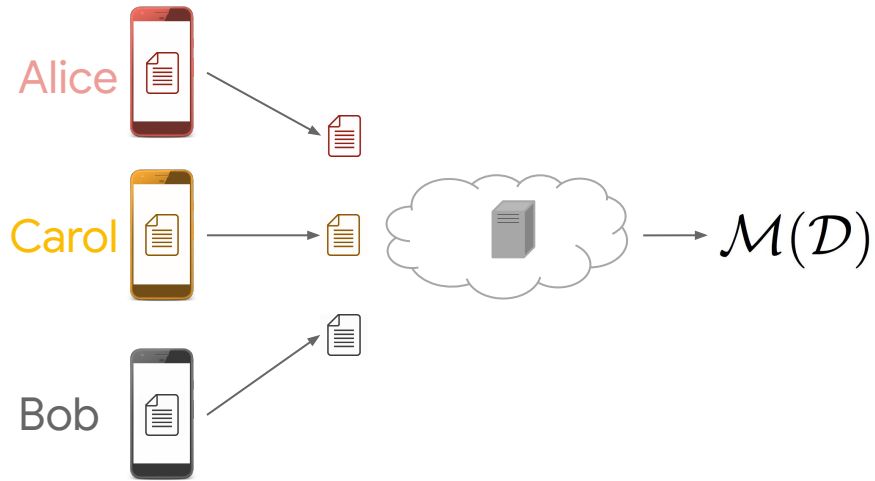
# Locally differentially private federated training



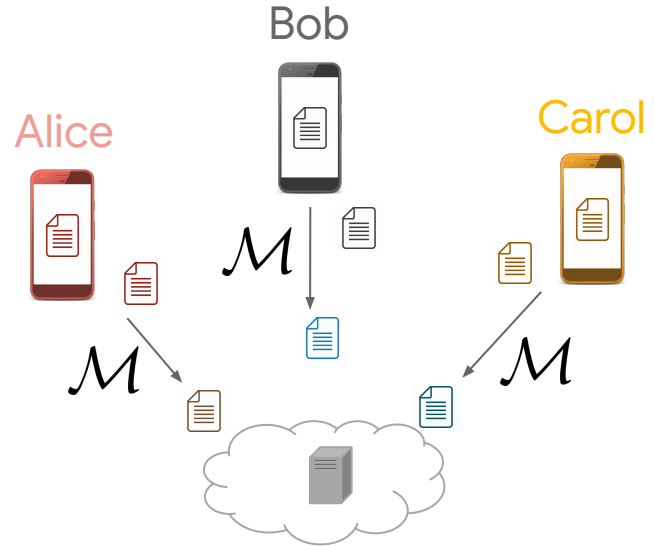
*Evfimievski, Alexandre, et al. Privacy preserving mining of association rules. Information Systems 2004*  
*Warner, Stanley L. Randomized response: A survey technique for eliminating evasive answer bias. JASA 1965*

# Can we combine the best of both worlds?

## Distributed Differential Privacy

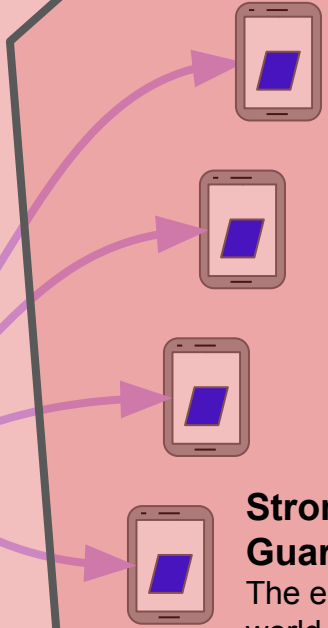
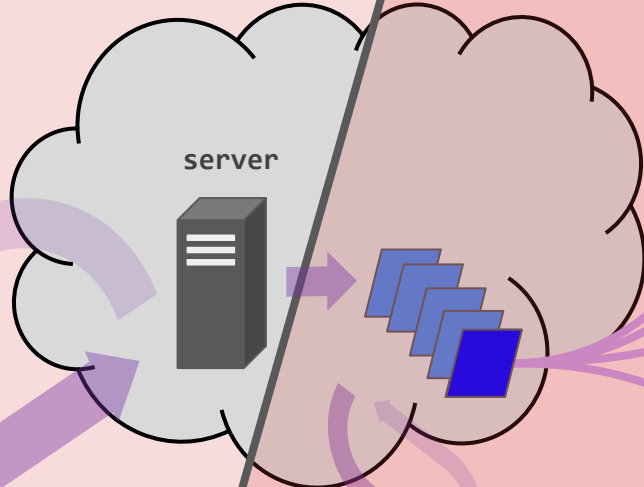
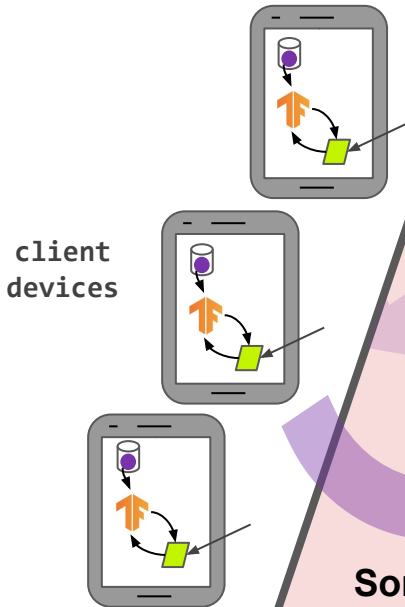


**Central DP: full trust in service provider**  
Higher utility at reasonable privacy levels



**Local DP: weaker trust assumptions**  
Utility often suffers

# Distributed DP



**Some DP guarantees**  
Very few (~10s) people have access to the server

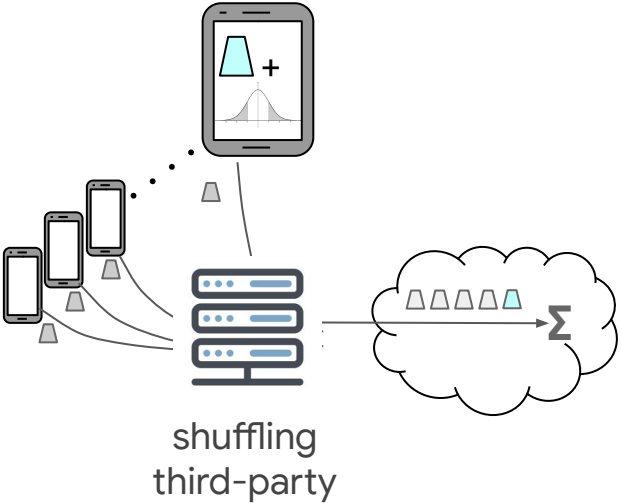
**Stronger Central DP Guarantee**  
More people (~1000s) have access to the model iterates

**Strongest Guarantee**  
The entire world will have access to deployed model



# Distributing trust for private aggregation

1 Trusted “third party”



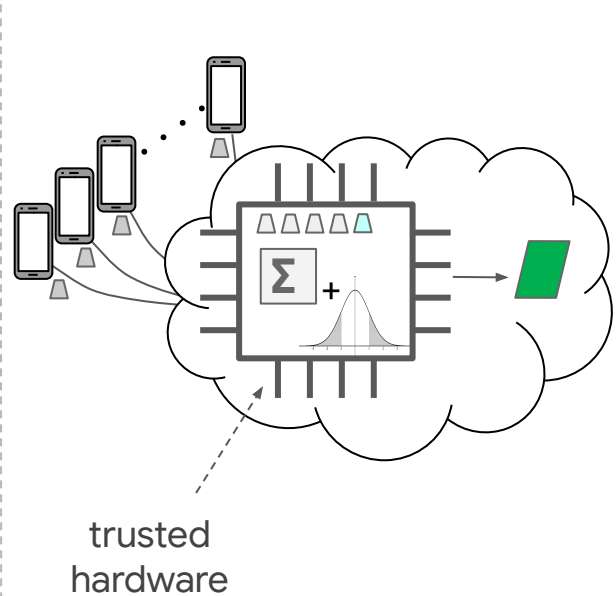
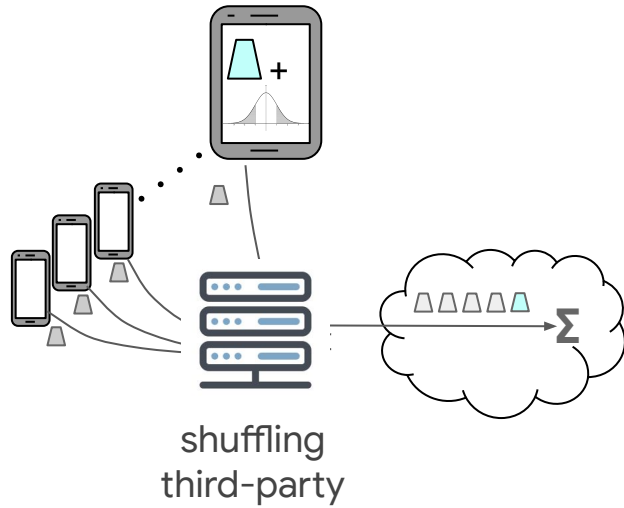
Andrea Bittau, et al. *Prochlo: Strong Privacy for Analytics in the Crowd*. SOSP 2017

Úlfar Erlingsson, et al. *Amplification by Shuffling: From Local to Central Differential Privacy via Anonymity*. SODA 2019

# Distributing trust for private aggregation

## 1 Trusted “third party”

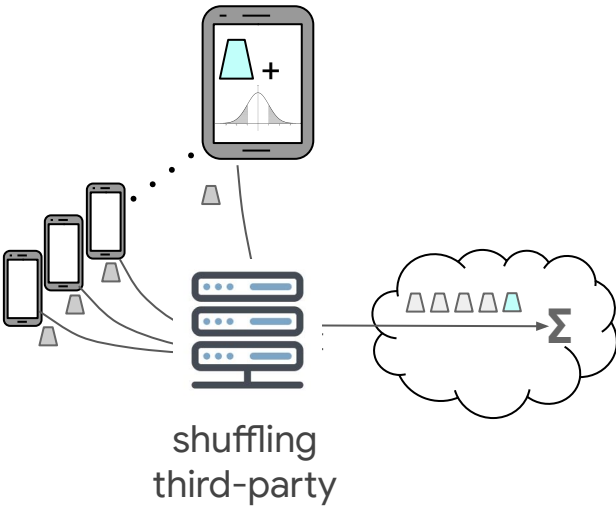
## 2 Trusted Execution Environments



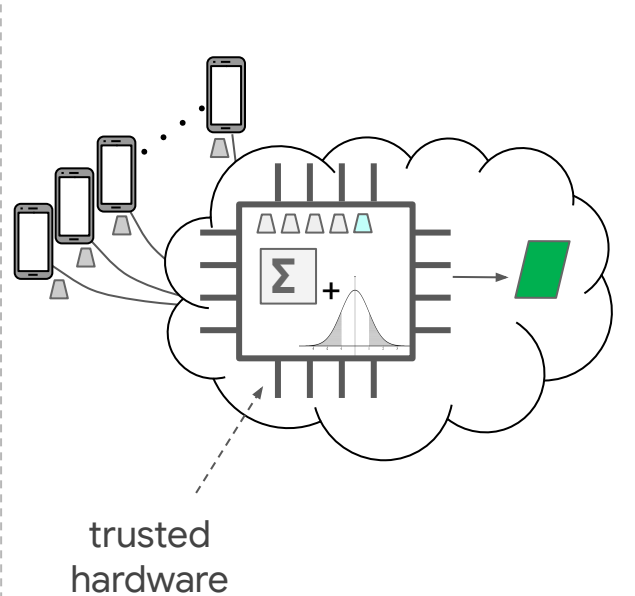
Andrea Bittau, et al. *Prochlo: Strong Privacy for Analytics in the Crowd*. SOSP 2017  
Úlfar Erlingsson, et al. *Amplification by Shuffling: From Local to Central Differential Privacy via Anonymity*. SODA 2019

# Distributing trust for private aggregation

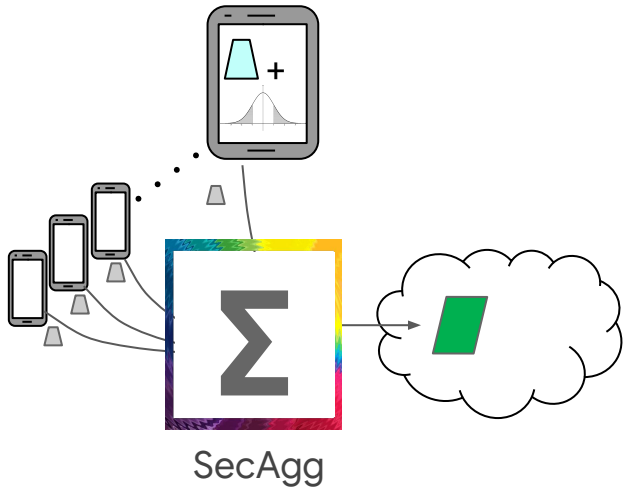
1 Trusted “third party”



2 Trusted Execution Environments



3 Trust via Cryptography

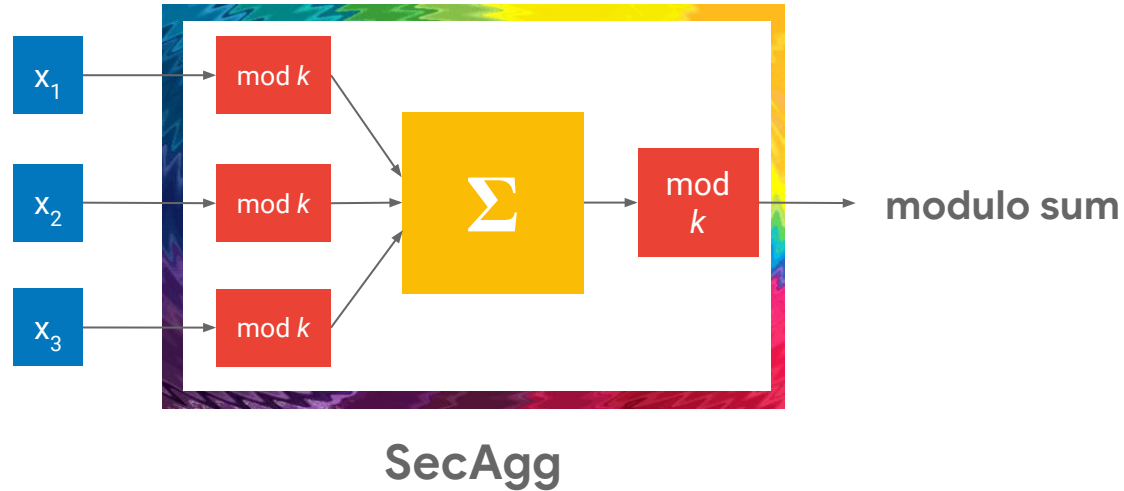


K. A. Bonawitz, et al. *Practical secure aggregation for privacy-preserving machine learning* CCS 2017  
J. Bell, et al. *Secure Single-Server Vector Aggregation with (Poly) Logarithmic Overhead* CCS 2020

**Secure Aggregation** allows a server to obtain the sum of high-dimensional vectors of client-held data in a way that ensures (cryptographically) that the server learns *just the sum*, and *no individual data whatsoever* \*.

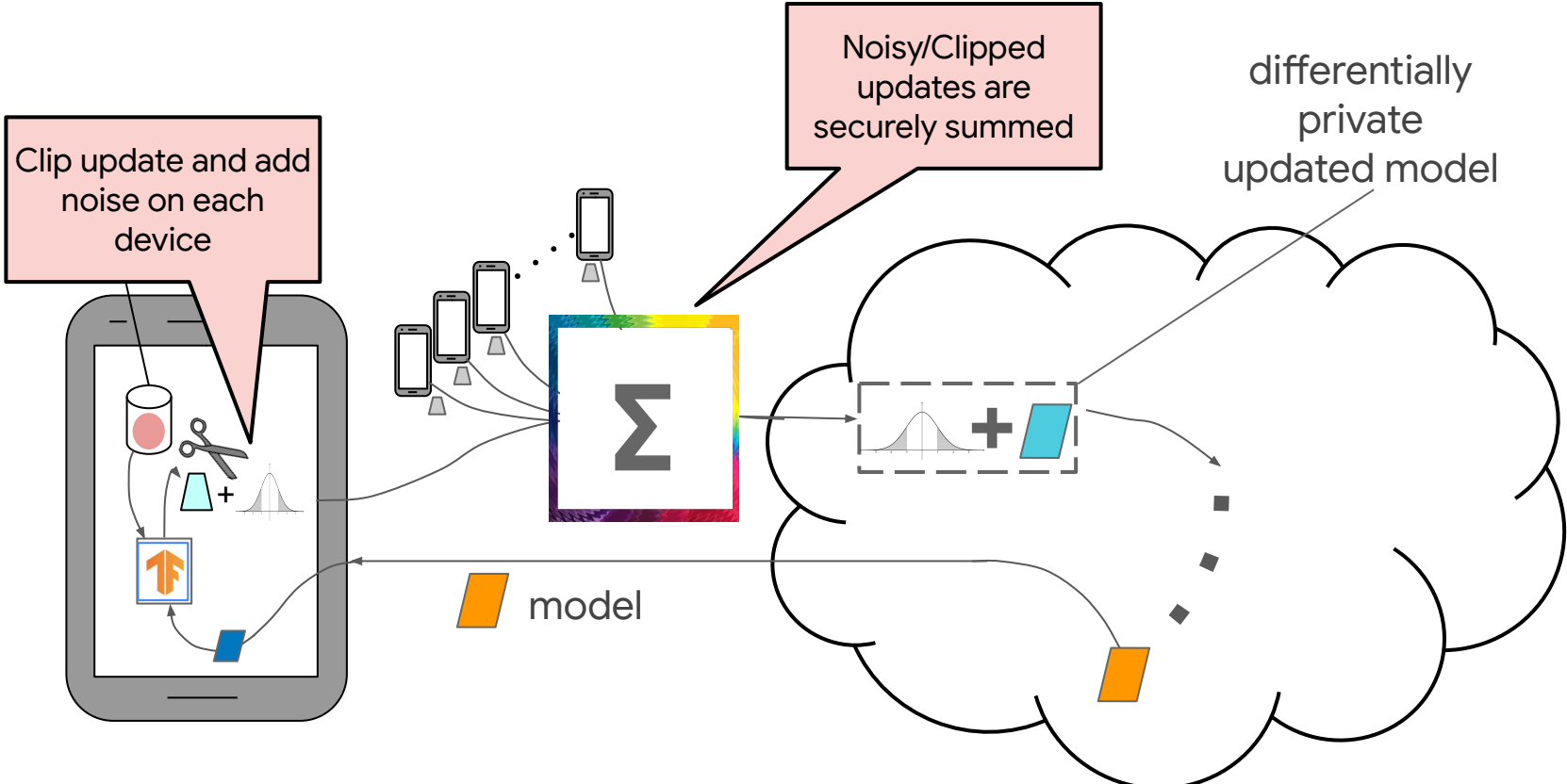
\* even if some users are malicious (and collude with the server), and some drop out.

# SecAgg: a closed box that performs integer modulo sums





# Why not just add continuous Gaussian noise?



# Challenges

- SecAgg operates on a finite group with integer modulo arithmetic
  - Need clever data discretization methods that do not inflate the sensitivity
  - Cannot use continuous mechanisms

# Challenges

- SecAgg operates on a finite group with integer modulo arithmetic
  - Need clever data discretization methods that do not inflate the sensitivity
  - Cannot use continuous mechanisms
- Communication efficiency is an important bottleneck in FL
  - Do not want privacy to depend on communication cost (SecAgg's group size)

# Challenges

- SecAgg operates on a finite group with integer modulo arithmetic
  - Need clever data discretization methods that do not inflate the sensitivity
  - Cannot use continuous mechanisms
- Communication efficiency is an important bottleneck in FL
  - Do not want privacy to depend on communication cost (SecAgg's group size)
- Many discrete mechanisms (e.g. k-RR) are not closed under summation
  - Analyzing sums of these mechanisms is difficult, especially in high dimensions

# Challenges

- SecAgg operates on a finite group with integer modulo arithmetic
  - Need clever data discretization methods that do not inflate the sensitivity
  - Cannot use continuous mechanisms
- Communication efficiency is an important bottleneck in FL
  - Do not want privacy to depend on communication cost (SecAgg's group size)
- Many discrete mechanisms (e.g. k-RR) are not closed under summation
  - Analyzing sums of these mechanisms is difficult, especially in high dimensions
- Discrete mechanisms with finite tails\* do not satisfy Rényi or concentrated DP
  - Avoids catastrophic privacy failures and allows for tight privacy accounting

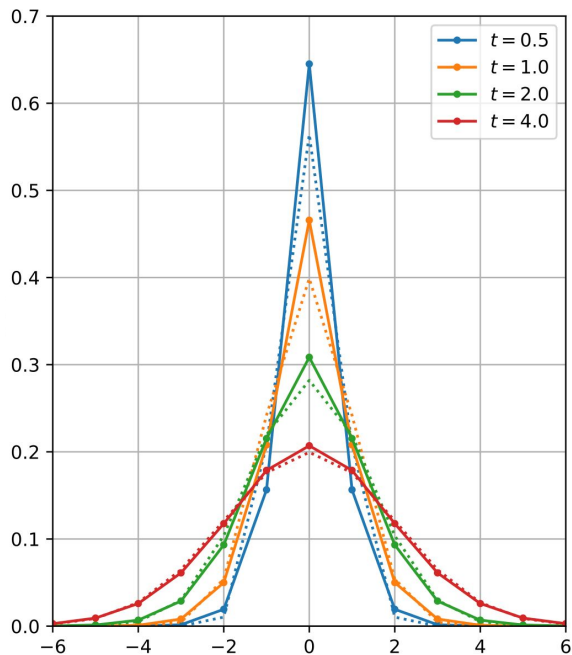
\*For example, the multi-dimensional binomial mechanism (Agarwal, et al. **cpSGD: Communication-efficient and differentially-private distributed SGD**. NeurIPS 2018.) does not achieve Rényi DP

# Challenges

- SecAgg operates on a finite group with integer modulo arithmetic
  - Need clever data discretization methods that do not inflate the sensitivity
  - Cannot use continuous mechanisms
- Communication efficiency is an important bottleneck in FL
  - Do not want privacy to depend on communication cost (SecAgg's group size)
- Many discrete mechanisms (e.g. k-RR) are not closed under summation
  - Analyzing sums of these mechanisms is difficult, especially in high dimensions
- Discrete mechanisms with finite tails\* do not satisfy Rényi or concentrated DP
  - Avoids catastrophic privacy failures and allows for tight privacy accounting
- Need mechanisms that can be sampled from exactly and efficiently

\*For example, the multi-dimensional binomial mechanism (Agarwal, et al. **cpSGD: Communication-efficient and differentially-private distributed SGD**. NeurIPS 2018.) does not achieve Rényi DP

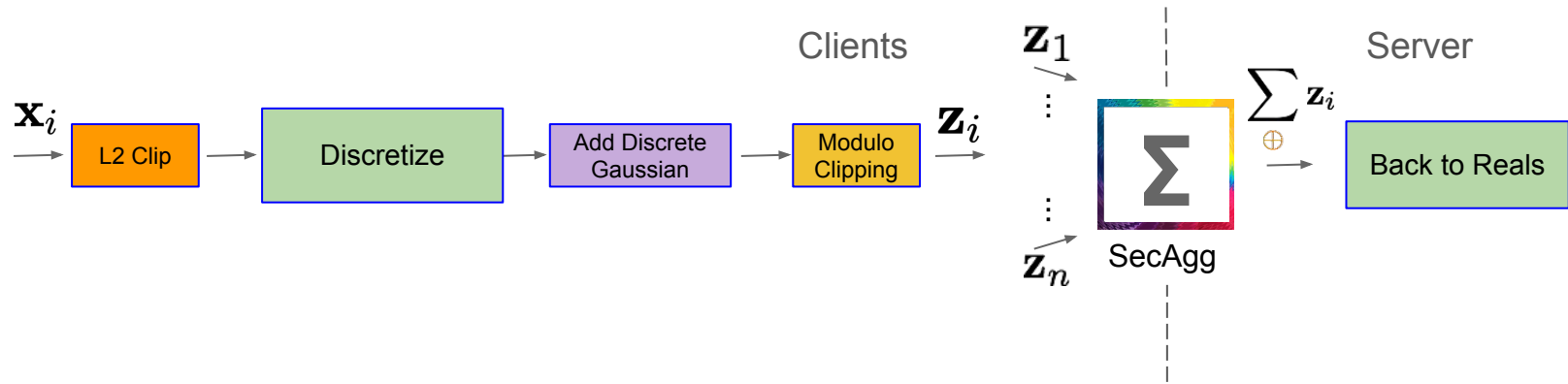
# The discrete Gaussian mechanism



$$\mathbb{P}_{X \leftarrow \mathcal{N}_{\mathbb{Z}}(\mu, t^2)} [X = n] = \frac{e^{-(n-\mu)^2/2t^2}}{\sum_{k \in \mathbb{Z}} e^{-(k-\mu)^2/2t^2}}$$

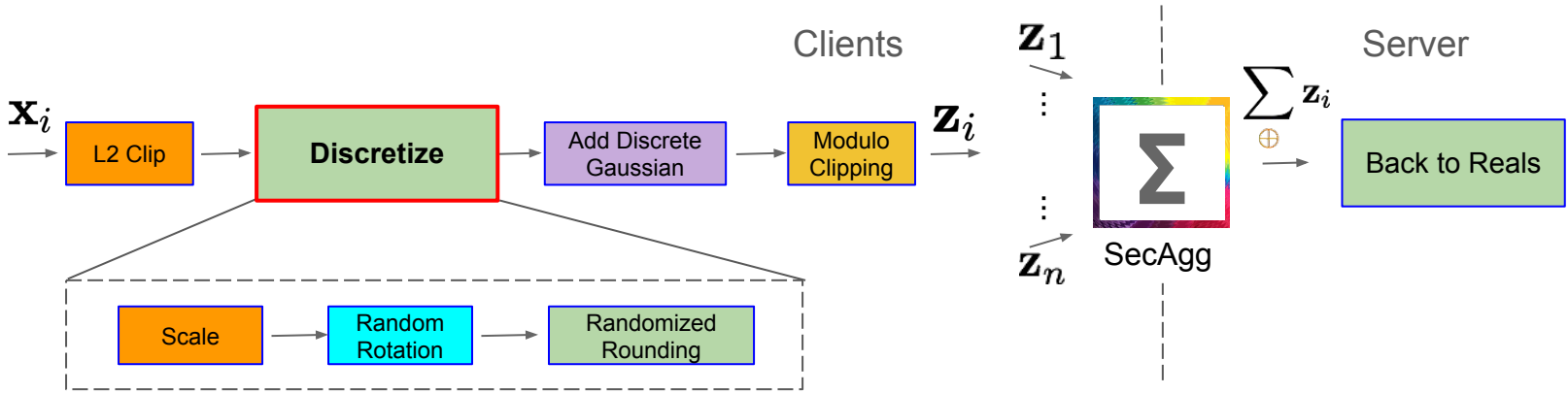
- **Discrete Gaussian:** discrete analog of continuous Gaussian ( $\neq$  rounding Gaussian to nearest ints)
  - Essentially **the same privacy-accuracy trade-off** as continuous Gaussian\*
  - zCDP / Rényi DP for tight compositions in learning contexts
- **Problem:** sums of discrete Gaussians  $\neq$  discrete Gaussians

# The Distributed Discrete Gaussian Mechanism



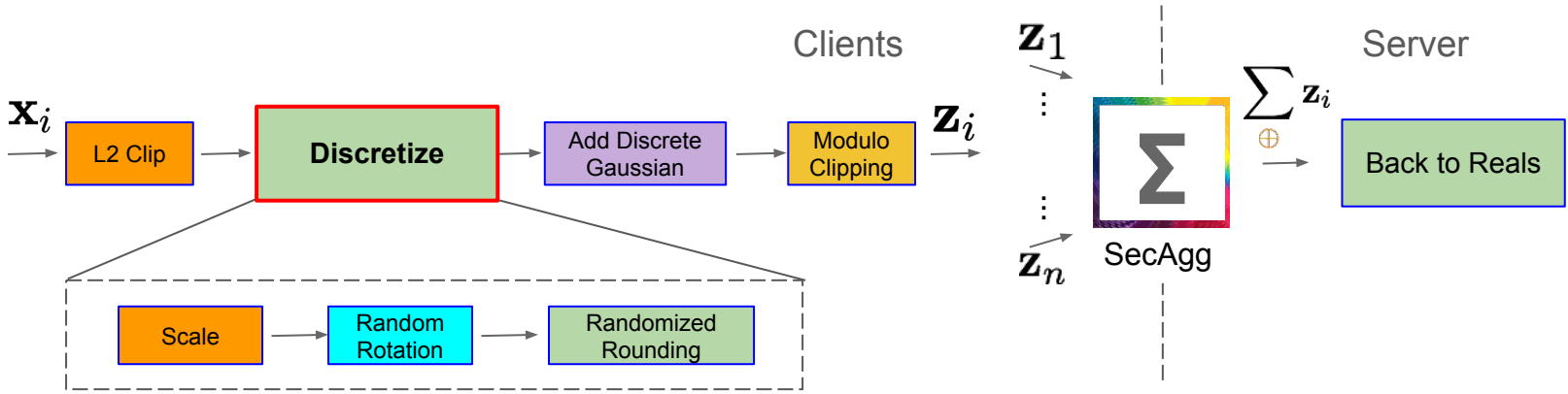


# Data Quantization



- **Scaling:** stretching the signal  $\rightarrow$  reduces quantization error
- **Random rotation:** “flatten” concentrated coordinates  $\rightarrow$  controls the L-inf norm
- **Randomized rounding:** values stochastically rounded to integers (unbiased)

# Data Quantization



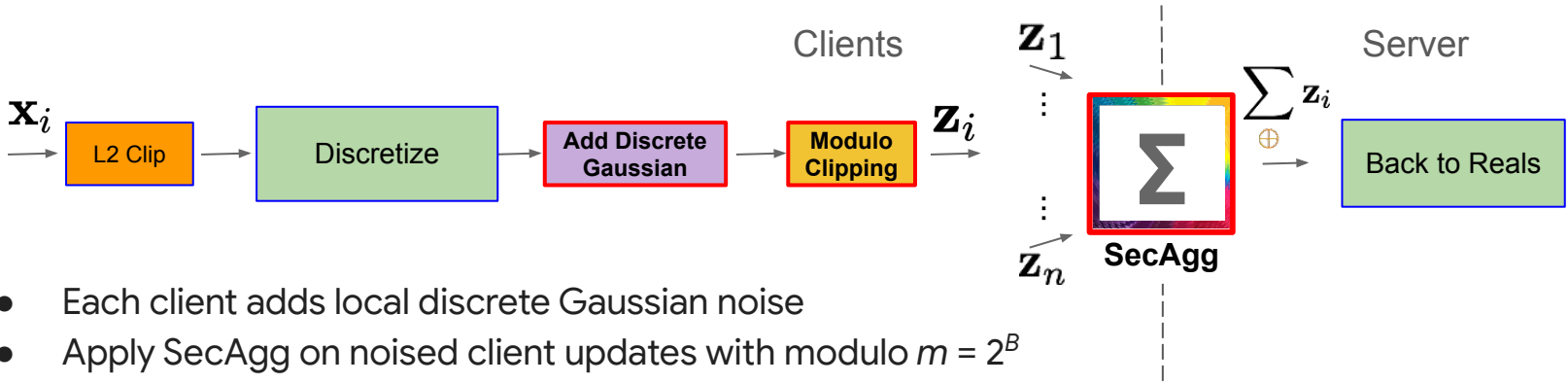
- **Scaling:** stretching the signal  $\rightarrow$  reduces quantization error
- **Random rotation:** “flatten” concentrated coordinates  $\rightarrow$  controls the L-inf norm
- **Randomized rounding:** values stochastically rounded to integers (unbiased)
- We can probabilistically bound the L2 norm growth from rounding (**helps reduce DP noise**):

**Proposition 22** (Properties of Randomized Rounding). *Let  $\beta \in [0, 1)$ ,  $\gamma > 0$ , and  $x \in \mathbb{R}^d$ .*

$$\Delta_2^2 := \min \left\{ \begin{aligned} &\|x\|_2^2 + \frac{1}{4}\gamma^2 d + \sqrt{2 \log(1/\beta)} \cdot \gamma \cdot \left( \|x\|_2 + \frac{1}{2}\gamma\sqrt{d} \right), \\ &\left( \|x\|_2 + \gamma\sqrt{d} \right)^2 \end{aligned} \right\}$$

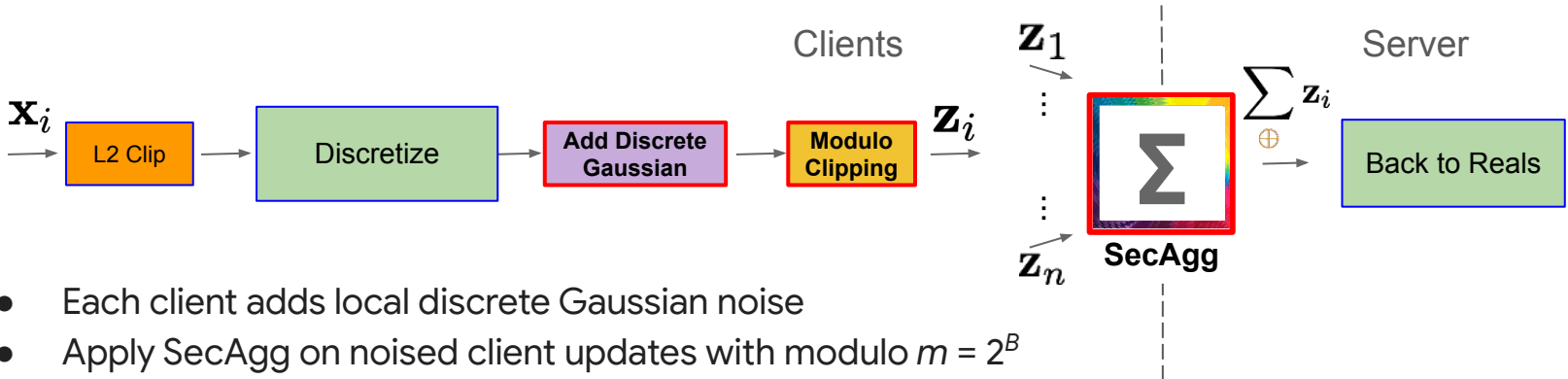
$d$ : client vector dim  
 $\gamma$ : rounding granularity;  
 inverse scaling factor  
 $\beta$ : rounding bias

# Local Noising & (Secure) Sums of Discrete Gaussians



- Each client adds local discrete Gaussian noise
- Apply SecAgg on noised client updates with modulo  $m = 2^B$

# Local Noising & (Secure) Sums of Discrete Gaussians



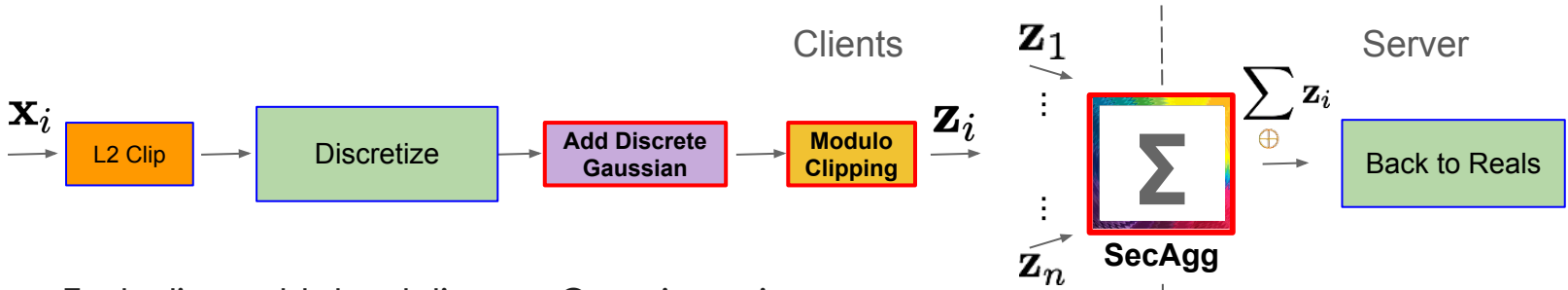
- Each client adds local discrete Gaussian noise
- Apply SecAgg on noised client updates with modulo  $m = 2^B$
- While sums of discrete Gaussians  $\neq$  discrete Gaussian, we show that they are **extremely close**:

**Theorem 11** (Convolution of two Discrete Gaussians). Let  $\sigma, \tau \geq \frac{1}{2}$ . Let  $X \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \sigma^2)$  and  $Y \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \tau^2)$  be independent. Let  $Z = X + Y$ . Let  $W \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \sigma^2 + \tau^2)$ . Then

$$D_{\pm\infty}(Z||W) = \sup_{z \in \mathbb{Z}} \left| \log \left( \frac{\mathbb{P}[Z = z]}{\mathbb{P}[W = z]} \right) \right| \leq 5 \cdot e^{-2\pi^2 / (1/\sigma^2 + 1/\tau^2)}.$$

- Exponentially small with larger variance;  $\leq 10^{-12}$  if  $\sigma_1^2 = \sigma_2^2 = 3$ .
- Noise is added on **quantized** client values, so  $\sigma^2$  is **scaled** and this is even smaller

# Local Noising & (Secure) Sums of Discrete Gaussians



- Each client adds local discrete Gaussian noise
- Apply SecAgg on noised client updates with modulo  $m = 2^B$
- While sums of discrete Gaussians  $\neq$  discrete Gaussian, we show that they are **extremely close**:

**Theorem 11** (Convolution of two Discrete Gaussians). Let  $\sigma, \tau \geq \frac{1}{2}$ . Let  $X \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \sigma^2)$  and  $Y \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \tau^2)$  be independent. Let  $Z = X + Y$ . Let  $W \leftarrow \mathcal{N}_{\mathbb{Z}}(0, \sigma^2 + \tau^2)$ . Then

$$D_{\pm\infty}(Z||W) = \sup_{z \in \mathbb{Z}} \left| \log \left( \frac{\mathbb{P}[Z = z]}{\mathbb{P}[W = z]} \right) \right| \leq 5 \cdot e^{-2\pi^2/(1/\sigma^2 + 1/\tau^2)}.$$

- Weak dependence on  $d$  (number of model params);
- 1st term same as central Gaussian/DGaussian

**Main Privacy Guarantees** with  $n$  clients:

$$\tau := 10 \cdot \sum_{k=1}^{n-1} e^{-2\pi^2 \frac{\sigma^2}{\gamma^2} \cdot \frac{k}{k+1}}$$

$$\epsilon := \min \left\{ \sqrt{\frac{\Delta_2^2}{n\sigma^2} + \frac{1}{2}\tau d}, \frac{\Delta_2}{\sqrt{n\sigma}} + \tau\sqrt{d} \right\}$$



Code

# Stack Overflow Next Word Prediction

- Next word prediction for question/answer sentences on StackOverflow.com with LSTMs
- $\sim 10^9$  sentences grouped by the  $N = 342477$  SO users/clients
- **Fig. 1: DDGauss matches continuous Gaussian as long as the bit-width  $B$  is sufficient**
- **Fig. 2: DDGauss scales (1000 clients per round) and works in low-noise (utility-first) settings**

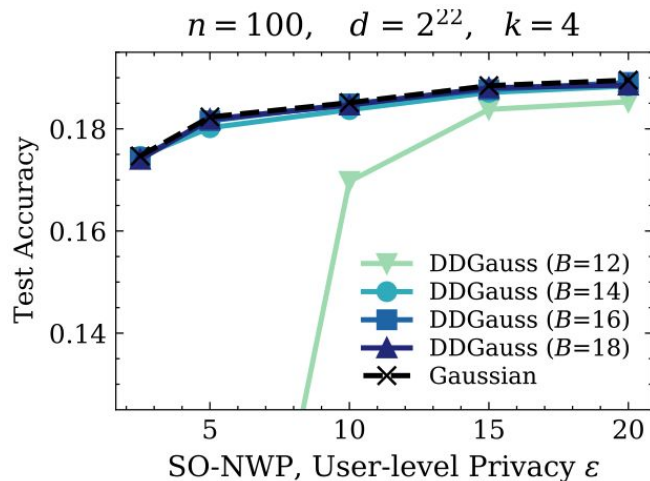


Fig. 1: Test acc with different  $\epsilon$  and  $B$  ( $n = 100$ ).

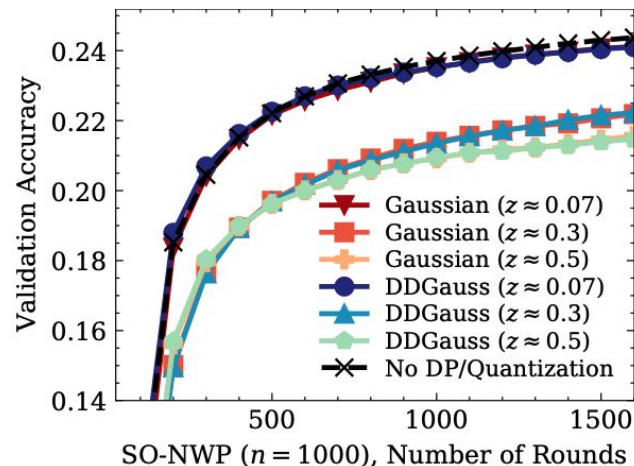
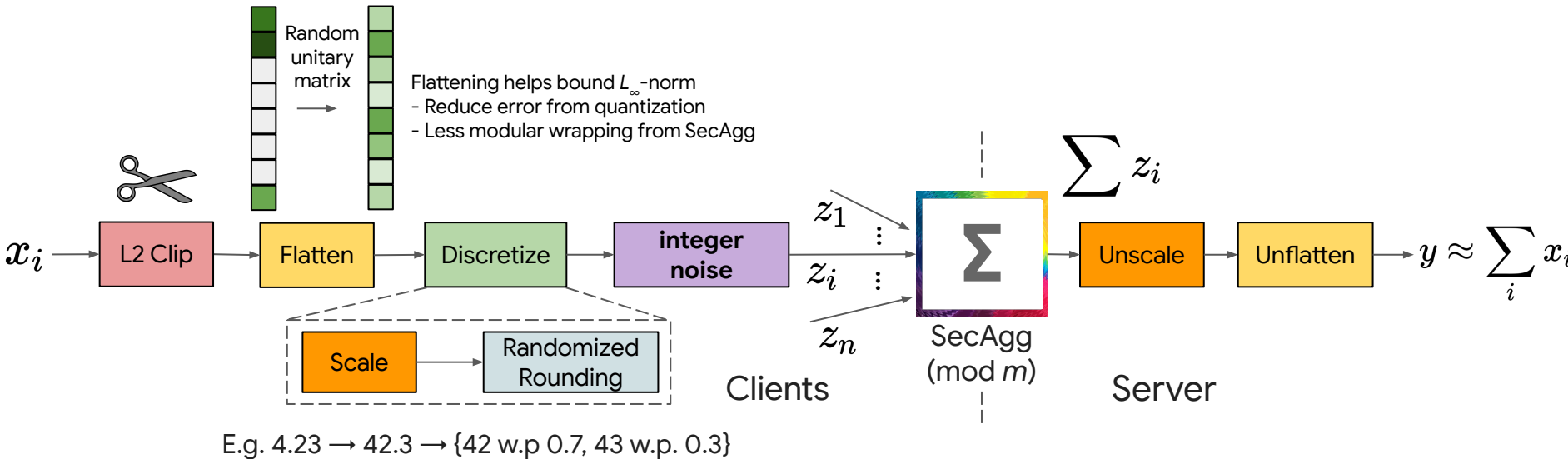


Fig. 2: Val acc with with  $n = 1000$  clients,  $B = 18$ .  
 $z$ : approximate noise multiplier aligned on  $\epsilon$ .

# Our end-to-end solution



# (Symmetric) Skellam Distribution

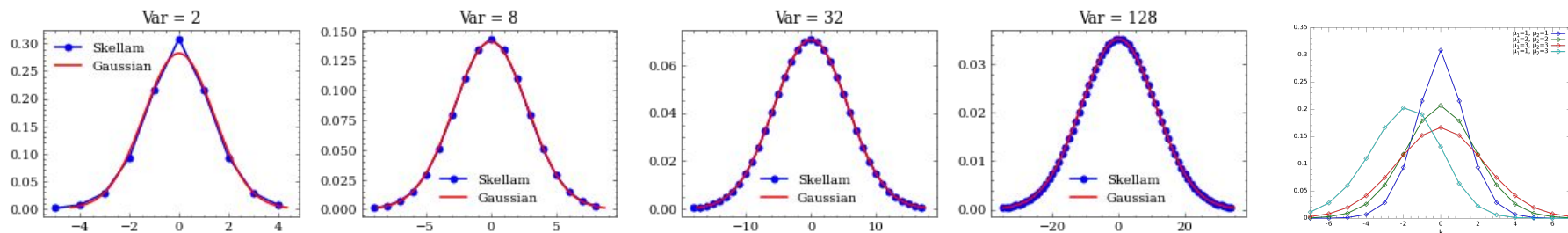
- Difference of two independent Poisson RVs. With mean  $\Delta$  and variance  $\mu$ ,

$I_k(z)$ : modified Bessel function of the first kind

$$X_i \sim \text{Sk}_{\Delta_i, \mu} \quad \text{with} \quad P(X_i = k) = e^{-\mu} I_{k-\Delta_i}(\mu)$$

- **Closed under summation**: easily switch between **central DP** & **distributed DP** (central vs local noise)
- **Easy to sample**: `np.random.poisson``
- Skellam gets closer to Gaussian as variance increases and we scale the output appropriately
- **Skellam Mechanism**: for an integer-valued query  $f(D)$ ,

$$\text{Sk}_{0, \mu}(f(D)) = f(D) + Z \quad \text{where} \quad Z \sim \text{Sk}_{0, \mu}$$





# (Distributed) Skellam

$$X \sim \text{Sk}_{\Delta, \mu}(X) \triangleq e^{-\mu} I_{k-\Delta}(\mu)$$

- **Main Rényi DP guarantee**

Gaussian RDP

L1 bound (after quantization)

$$\Delta_1 \leq \Delta_2 \cdot \min(\sqrt{d}, \Delta_2)$$

2nd term goes to 0 with larger variance (higher privacy)

For  $\ell_1, \ell_2$  sensitivities  $\Delta_1, \Delta_2$ , central variance  $\mu$ , and order  $\alpha > 1, \alpha \in \mathbb{Z}$ ,

$$\epsilon(\alpha) \leq \frac{\alpha \Delta_2^2}{2\mu} + \min\left(\frac{(2\alpha - 1)\Delta_2^2 + 6\Delta_1}{4\mu^2}, \frac{3\Delta_1}{2\mu}\right)$$

# (Distributed) Skellam

$$X \sim \text{Sk}_{\Delta, \mu}(X) \triangleq e^{-\mu} I_{k-\Delta}(\mu)$$

- Main Rényi DP guarantee**

Gaussian RDP

L1 bound (after quantization)

$$\Delta_1 \leq \Delta_2 \cdot \min(\sqrt{d}, \Delta_2)$$

2nd term goes to 0 with larger variance (higher privacy or large scaling)

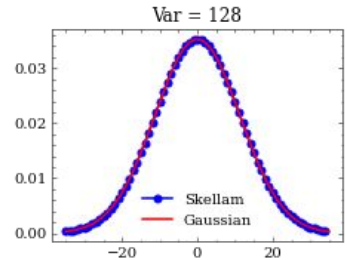
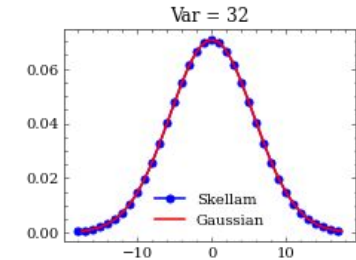
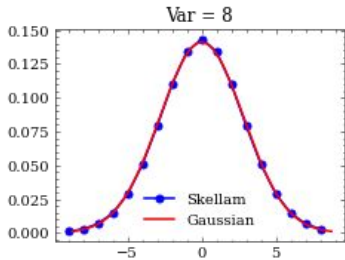
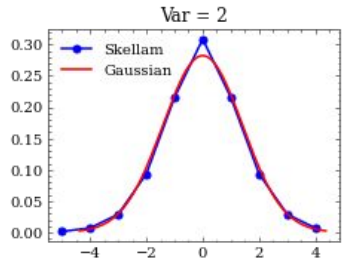
For  $\ell_1, \ell_2$  sensitivities  $\Delta_1, \Delta_2$ , central variance  $\mu$ , and order  $\alpha > 1, \alpha \in \mathbb{Z}$ ,

$$\epsilon(\alpha) \leq \frac{\alpha \Delta_2^2}{2\mu} + \min\left(\frac{(2\alpha - 1)\Delta_2^2 + 6\Delta_1}{4\mu^2}, \frac{3\Delta_1}{2\mu}\right)$$

- Effect of scaling** (scale both noise stddev and sensitivity)

**Corollary 4.1** (Scaled Skellam Mechanism). *With a scaling factor  $s \in \mathbb{R}$ , the multi-dimensional Skellam Mechanism is  $(\alpha, \epsilon)$ -RDP with*

$$\epsilon(\alpha) \leq \frac{\alpha \Delta_2^2}{2\mu} + \min\left(\frac{(2\alpha - 1)\Delta_2^2}{4s^2\mu^2} + \frac{3\Delta_1}{2s^3\mu^2}, \frac{3\Delta_1}{2s\mu}\right)$$

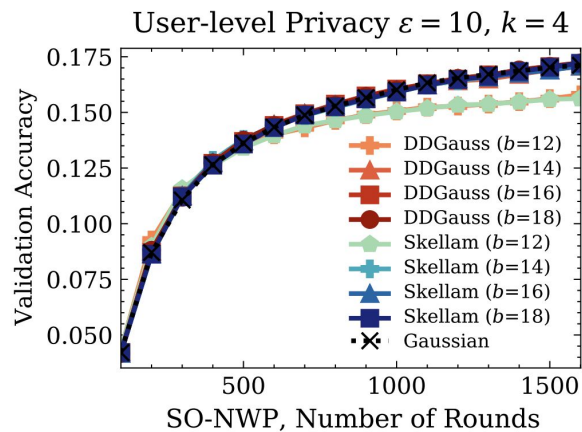
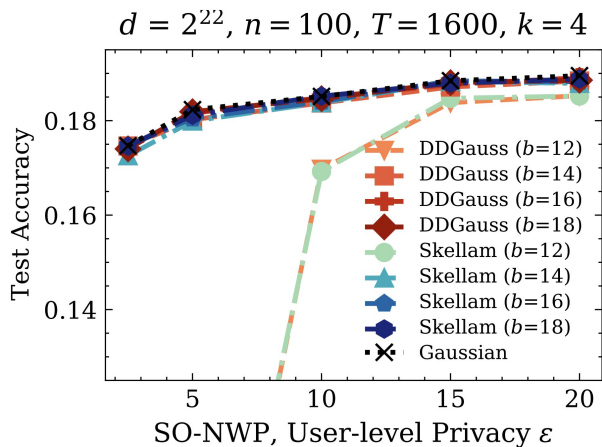


# Stack Overflow Next Word Prediction



Code

- Next word prediction for questions/answers sentences on StackOverflow.com with LSTMs
- $\sim 10^9$  sentences grouped by  $N = 342477$  users on Stack Overflow
- **Left:** Test acc across various privacy levels  $\epsilon$  and bit-widths  $b$
- **Right:** Validation acc across training rounds
- Skellam matches continuous Gaussian and distributed discrete Gaussian



# Better communication efficiency?

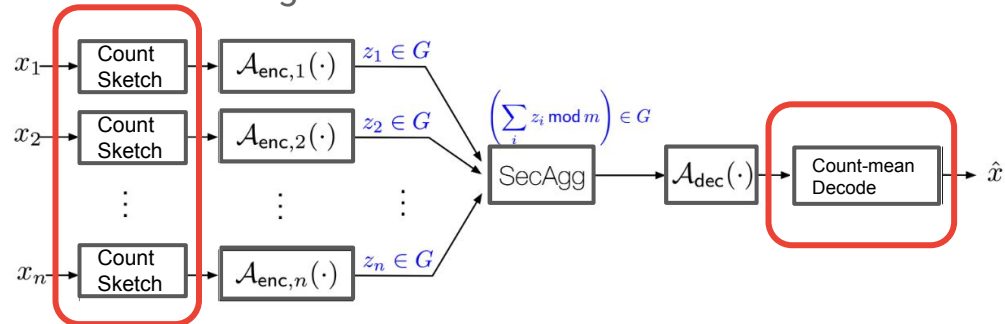
To achieve centralized error of  $O\left(\frac{d}{n^2\varepsilon^2}\right)$  each client must transmit  $\Omega\left(d\log\left(\frac{n^2\varepsilon^2}{d}\right)\right) = \tilde{\Omega}(d)$  bits.

# Better communication efficiency?

To achieve centralized error of  $O\left(\frac{d}{n^2\varepsilon^2}\right)$  each client must transmit  $\Omega\left(d \log\left(\frac{n^2\varepsilon^2}{d}\right)\right) = \tilde{\Omega}(d)$  bits.

In the worst-case, each client **cannot** transmit less than the entire gradient!

- But, gradients may be near-sparse! Is their sum?
- We can leverage this structure to compress each  $x_i$ !



- We will use a count-mean sketch: **efficient** and **linear** dimensionality reduction

Stich, Sebastian U., Jean-Baptiste Cordonnier, and Martin Jaggi. "Sparsified SGD with memory." arXiv preprint arXiv:1809.07599 (2018).

Barnes, Leighton Pate, et al. "rTop-k: A statistical estimation approach to distributed SGD." IEEE Journal on Selected Areas in Information Theory 1.3 (2020): 897-907.

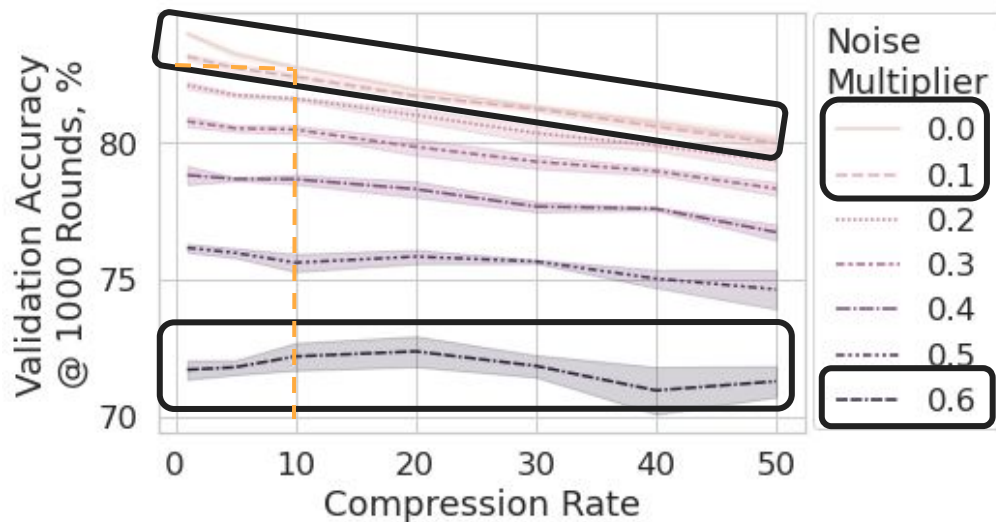
Rothchild, Daniel, et al. "Fetchsgd: Communication-efficient federated learning with sketching." International Conference on Machine Learning. PMLR, 2020.

Haddadpour, Farzin, et al. "Fedsketch: Communication-efficient and private federated learning via sketching." arXiv preprint arXiv:2008.04975 (2020).

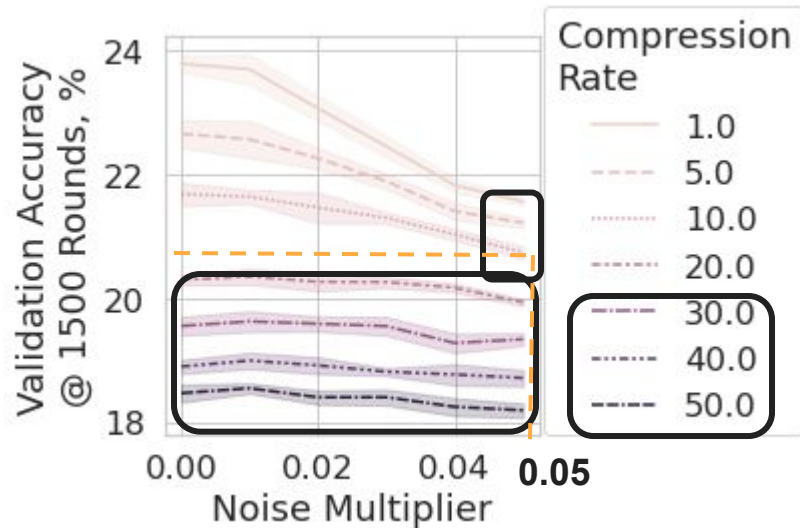
# Experiments

**Want:**  
Noise Multiplier ↑ ↔ ◆◆  
Accuracy ↑

## Federated EMNIST-62 @ 100 Clients



## Stack Overflow Next Word Prediction @ 100 Clients



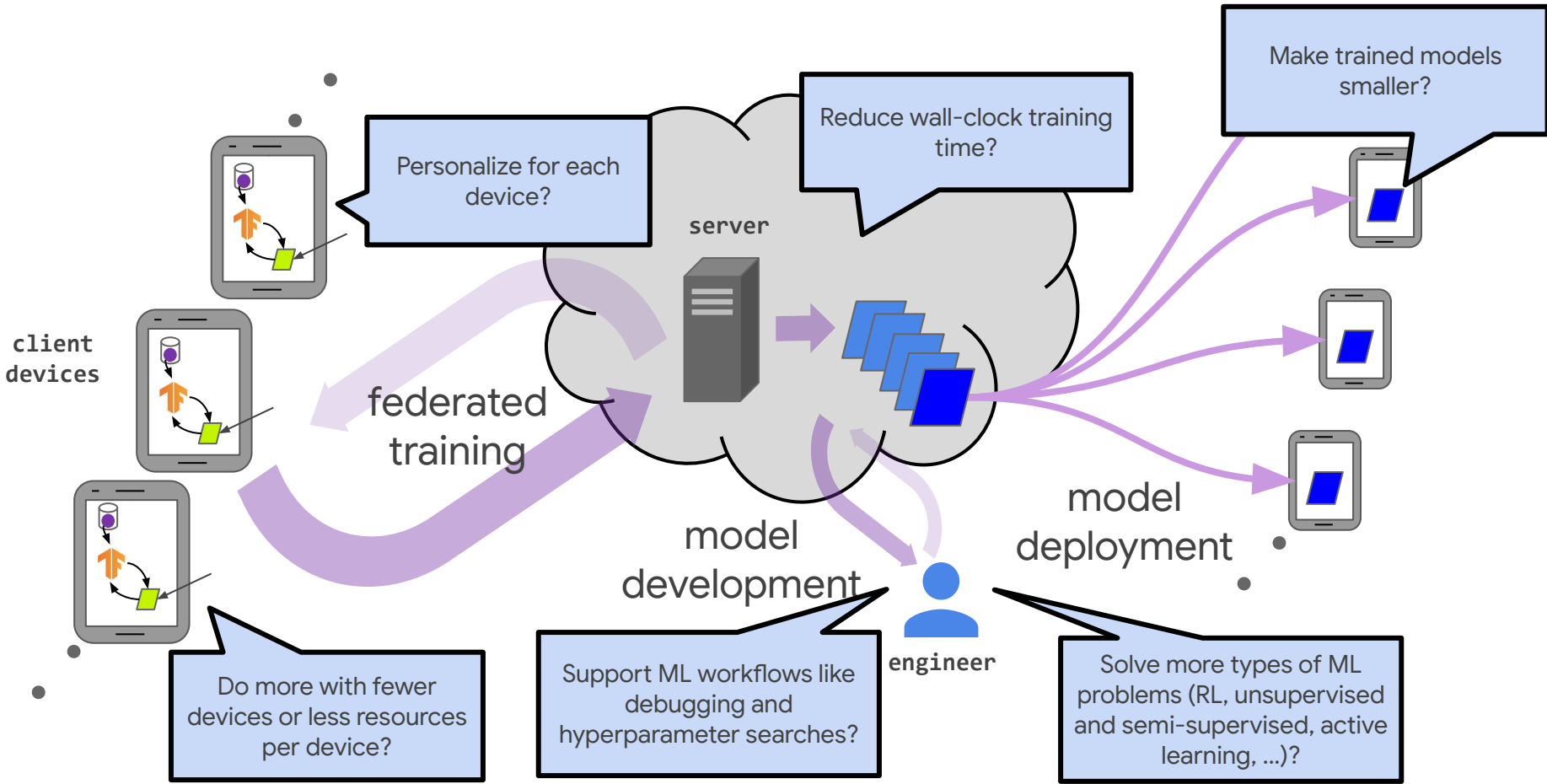
# Challenges & Opportunities

# Open technical challenges in privacy

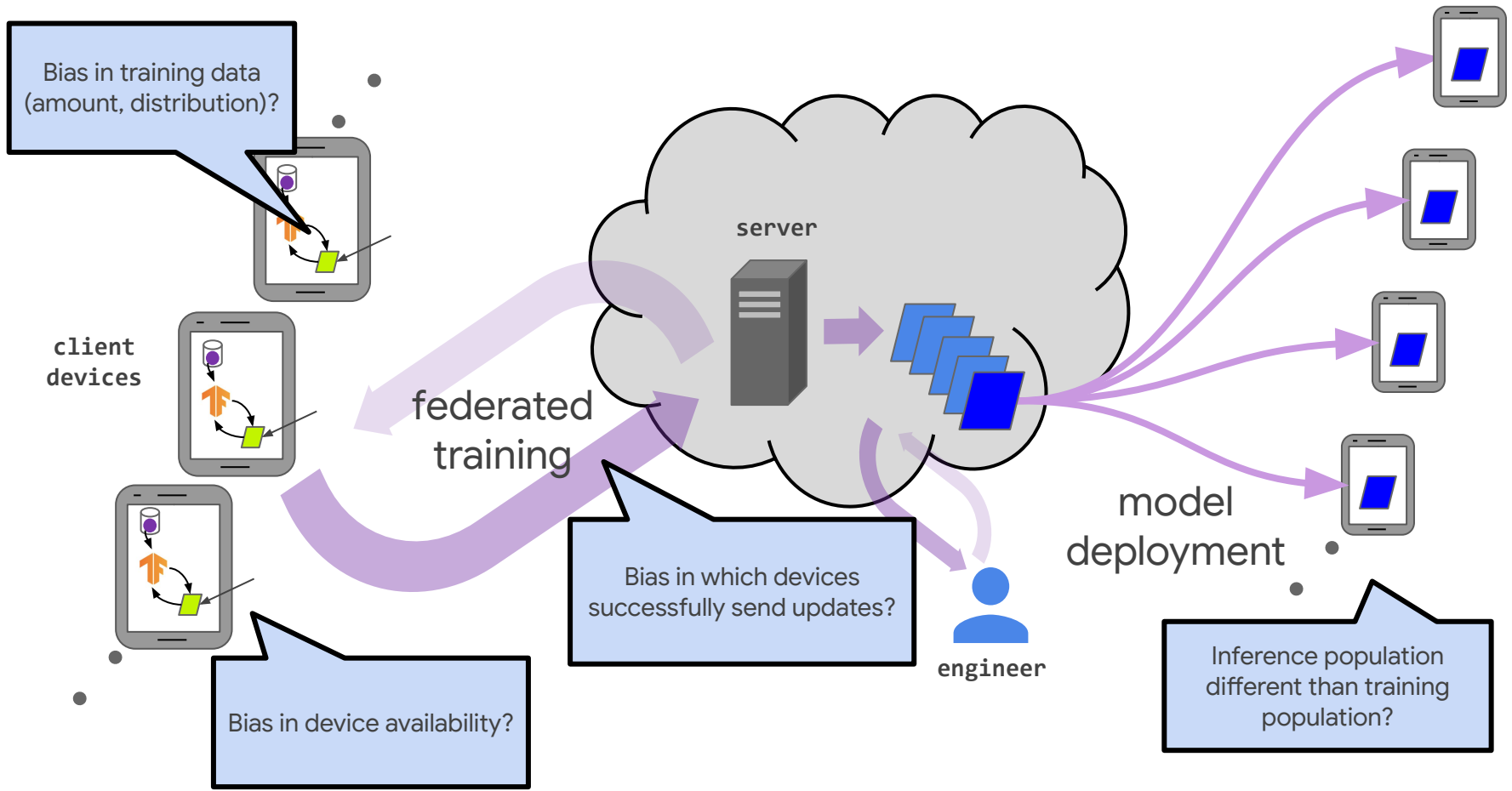
- Privacy is multifaceted
  - Need to better understand *privacy, communication, computation, accuracy, sparsity* tradeoffs
  - Tensions between privacy, robustness, and fairness are very interesting and remain underexplored – personalization may play an important role in easing the tensions
  - Cryptographic techniques will play a critical role in strengthening privacy
- Differential privacy provides an incredibly useful tool
  - But it often comes at a “hit” in accuracy
  - If we have to pay, we'd usually rather pay with more computation (not privacy or accuracy)
  - How to choose epsilon remains (and perhaps will always be) an open question
  - How to make sense of large-ish epsilons?
  - Model auditing techniques for measuring privacy loss (memorization) are complimentary
- Privacy budgeting and management systems are not available
  - Can scientists apply complex and repeated learning tasks on the same or similar datasets?
  - How do we efficiently track and quantify the privacy loss of a complex system?
- Public data is largely underutilized
  - Public data will play a key role in improving privacy-accuracy tradeoffs
  - How do we optimally combine public and private datasets during training?



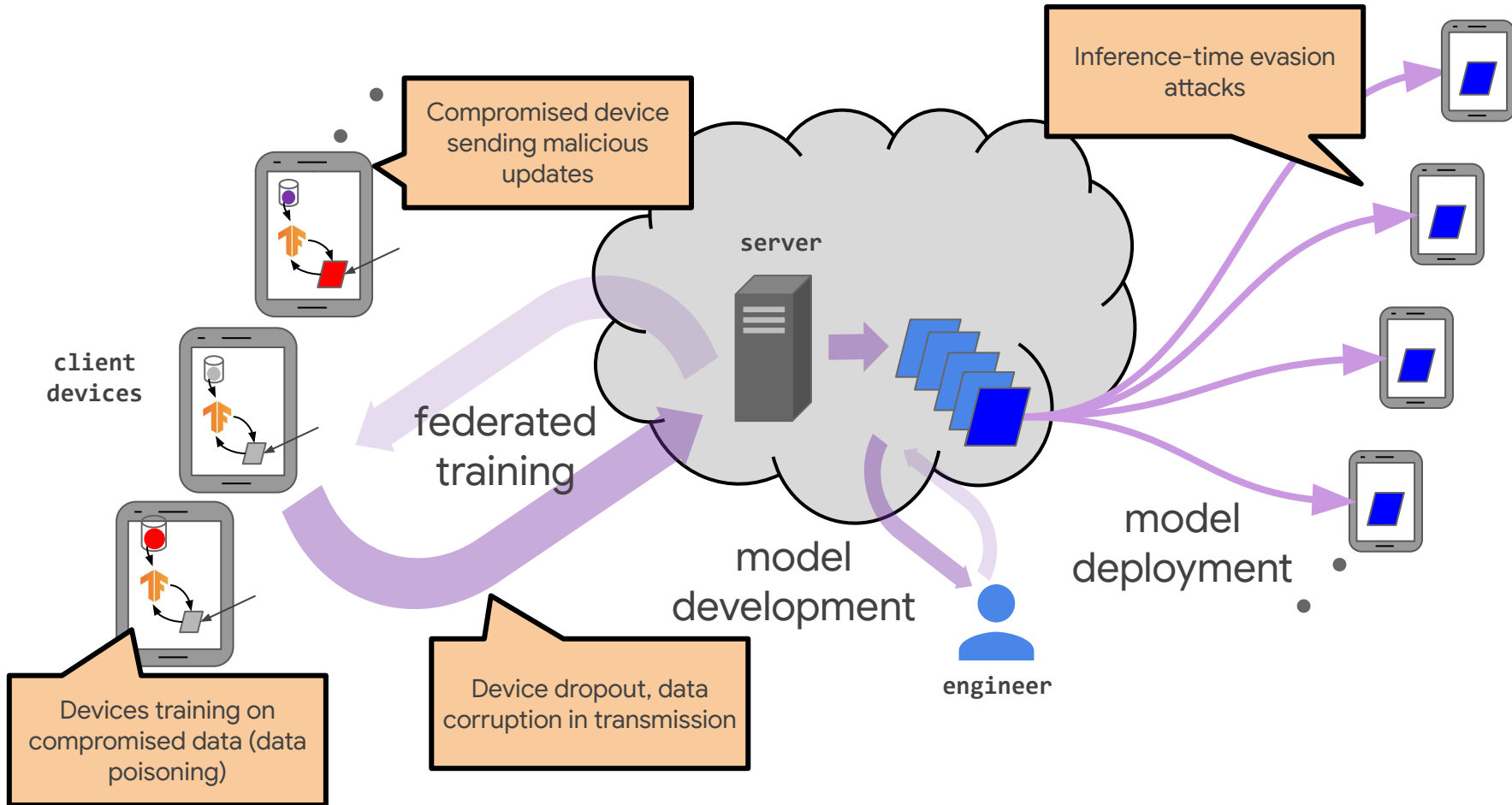
# Improving efficiency and effectiveness



# Ensuring fairness and addressing sources of bias



# Robustness to attacks and failures



## Advances and Open Problems in Federated Learning

Peter Kairouz<sup>7\*</sup> H. Brendan McMahan<sup>7\*</sup> Brendan Avent<sup>21</sup> Aurélien Bellet<sup>9</sup>  
Mehdi Bennis<sup>19</sup> Arjun Nitin Bhagoji<sup>13</sup> Kallista Bonawitz<sup>7</sup> Zachary Charles<sup>7</sup>  
Graham Cormode<sup>23</sup> Rachel Cummings<sup>6</sup> Rafael G.L. D'Oliveira<sup>14</sup>  
Hubert Eichner<sup>7</sup> Salim El Rouayheb<sup>14</sup> David Evans<sup>22</sup> Josh Gardner<sup>24</sup>  
Zachary Garrett<sup>7</sup> Adrià Gascón<sup>7</sup> Badih Ghazi<sup>7</sup> Phillip B. Gibbons<sup>2</sup>  
Marco Gruteser<sup>7,14</sup> Zaid Harchaoui<sup>24</sup> Chaoyang He<sup>21</sup> Lie He<sup>4</sup>  
Zhouyuan Huo<sup>20</sup> Ben Hutchinson<sup>7</sup> Justin Hsu<sup>25</sup> Martin Jaggi<sup>4</sup> Tara Javidi<sup>17</sup>  
Gauri Joshi<sup>2</sup> Mikhail Khodak<sup>2</sup> Jakub Konečný<sup>7</sup> Aleksandra Korolova<sup>21</sup>  
Farinaz Koushanfar<sup>17</sup> Sanmi Koyejo<sup>7,18</sup> Tancrede Lepoint<sup>7</sup> Yang Liu<sup>12</sup>  
Prateek Mittal<sup>13</sup> Mehryar Mohri<sup>7</sup> Richard Nock<sup>1</sup> Ayfer Özgür<sup>15</sup>  
Rasmus Pagh<sup>7,10</sup> Hang Qi<sup>7</sup> Daniel Ramage<sup>7</sup> Ramesh Raskar<sup>11</sup>  
Mariana Raykova<sup>7</sup> Dawn Song<sup>16</sup> Weikang Song<sup>7</sup> Sebastian U. Stich<sup>4</sup>  
Ziteng Sun<sup>3</sup> Ananda Theertha Suresh<sup>7</sup> Florian Tramèr<sup>15</sup> Praneeth Vepakomma<sup>11</sup>  
Jianyu Wang<sup>2</sup> Li Xiong<sup>5</sup> Zheng Xu<sup>7</sup> Qiang Yang<sup>8</sup> Felix X. Yu<sup>7</sup> Han Yu<sup>12</sup>  
Sen Zhao<sup>7</sup>

<sup>1</sup>Australian National University, <sup>2</sup>Carnegie Mellon University, <sup>3</sup>Cornell University,

<sup>4</sup>École Polytechnique Fédérale de Lausanne, <sup>5</sup>Emory University, <sup>6</sup>Georgia Institute of Technology,

<sup>7</sup>Google Research, <sup>8</sup>Hong Kong University of Science and Technology, <sup>9</sup>INRIA, <sup>10</sup>IT University of Copenhagen,

<sup>11</sup>Massachusetts Institute of Technology, <sup>12</sup>Nanyang Technological University, <sup>13</sup>Princeton University,

<sup>14</sup>Rutgers University, <sup>15</sup>Stanford University, <sup>16</sup>University of California Berkeley,

<sup>17</sup>University of California San Diego, <sup>18</sup>University of Illinois Urbana-Champaign, <sup>19</sup>University of Oulu,

<sup>20</sup>University of Pittsburgh, <sup>21</sup>University of Southern California, <sup>22</sup>University of Virginia,

<sup>23</sup>University of Warwick, <sup>24</sup>University of Washington, <sup>25</sup>University of Wisconsin-Madison

### Abstract

Federated learning (FL) is a machine learning setting where many clients (e.g. mobile devices or whole organizations) collaboratively train a model under the orchestration of a central server (e.g. service provider), while keeping the training data decentralized. FL embodies the principles of focused data collection and minimization, and can mitigate many of the systemic privacy risks and costs resulting from traditional, centralized machine learning and data science approaches. Motivated by the explosive growth in FL research, this paper discusses recent advances and presents an extensive collection of open problems and challenges.

## Advances and Open Problems in FL

59 authors from 25 top institutions

[arxiv.org/abs/1912.04977](https://arxiv.org/abs/1912.04977)

*Foundations and Trends in Machine Learning*





# A Field Guide to Federated Optimization

Jianyu Wang<sup>\*1</sup>, Zachary Charles<sup>\*3</sup>, Zheng Xu<sup>\*3</sup>, Gauri Joshi<sup>\*1</sup>, H. Brendan McMahan<sup>\*3</sup>, Blaise Agüera y Arcas<sup>3</sup>, Maruan Al-Shedivat<sup>1</sup>, Galen Andrew<sup>3</sup>, Salman Avestimehr<sup>13</sup>, Katharine Daly<sup>3</sup>, Deepesh Data<sup>9</sup>, Suhas Diggavi<sup>9</sup>, Hubert Eichner<sup>3</sup>, Advait Gadhikar<sup>1</sup>, Zachary Garrett<sup>3</sup>, Antonious M. Girgis<sup>9</sup>, Filip Hanzely<sup>8</sup>, Andrew Hard<sup>3</sup>, Chaoyang He<sup>13</sup>, Samuel Horváth<sup>4</sup>, Zhouyuan Huo<sup>3</sup>, Alex Ingerman<sup>3</sup>, Martin Jaggi<sup>2</sup>, Tara Javidi<sup>10</sup>, Peter Kairouz<sup>3</sup>, Satyen Kale<sup>3</sup>, Sai Praneeth Karimireddy<sup>2</sup>, Jakub Konečný<sup>3</sup>, Sanmi Koyejo<sup>11</sup>, Tian Li<sup>1</sup>, Luyang Liu<sup>3</sup>, Mehryar Mohri<sup>3</sup>, Hang Qi<sup>3</sup>, Sashank J. Reddi<sup>3</sup>, Peter Richtárik<sup>4</sup>, Karan Singhal<sup>3</sup>, Virginia Smith<sup>1</sup>, Mahdi Soltanolkotabi<sup>13</sup>, Weikang Song<sup>3</sup>, Ananda Theertha Suresh<sup>3</sup>, Sebastian U. Stich<sup>2</sup>, Ameet Talwalkar<sup>1</sup>, Hongyi Wang<sup>14</sup>, Blake Woodworth<sup>8</sup>, Shanshan Wu<sup>3</sup>, Felix X. Yu<sup>3</sup>, Honglin Yuan<sup>6</sup>, Manzil Zaheer<sup>3</sup>, Mi Zhang<sup>5</sup>, Tong Zhang<sup>3,7</sup>, Chunxiang Zheng<sup>3</sup>, Chen Zhu<sup>12</sup>, and Wennan Zhu<sup>3</sup>

<sup>1</sup>Carnegie Mellon University, <sup>2</sup>École Polytechnique Fédérale de Lausanne, <sup>3</sup>Google Research, <sup>4</sup>King Abdullah University of Science and Technology, <sup>5</sup>Michigan State University, <sup>6</sup>Stanford University, <sup>7</sup>The Hong Kong University of Science and Technology, <sup>8</sup>Toyota Technological Institute at Chicago, <sup>9</sup>University of California, Los Angeles, <sup>10</sup>University of California, San Diego, <sup>11</sup>University of Illinois Urbana-Champaign, <sup>12</sup>University of Maryland, College Park, <sup>13</sup>University of Southern California, <sup>14</sup>University of Wisconsin–Madison

## Abstract

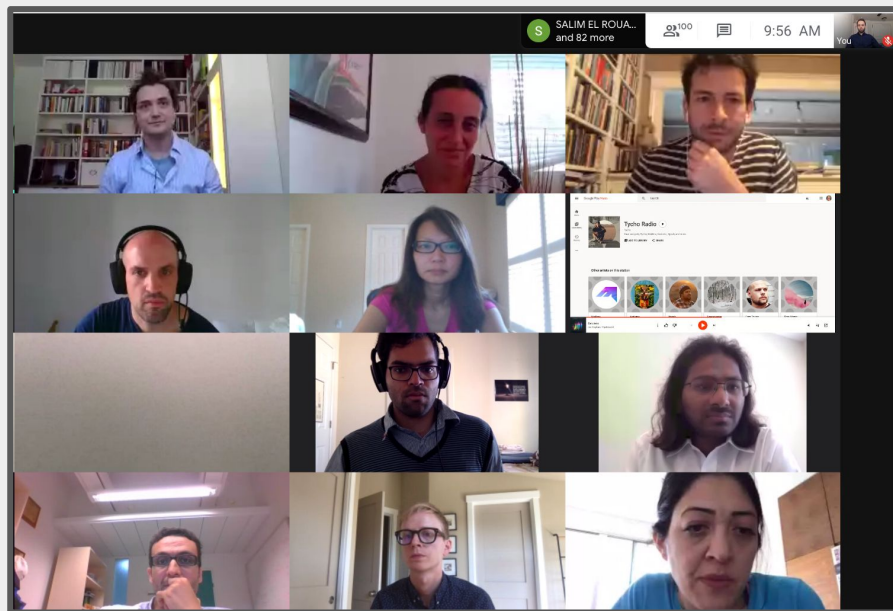
Federated learning and analytics are a distributed approach for collaboratively learning models (or statistics) from decentralized data, motivated by and designed for privacy protection. The distributed learning process can be formulated as solving federated optimization problems, which emphasize communication efficiency, data heterogeneity, compatibility with privacy and system requirements, and other constraints that are not primary considerations in other problem settings. This paper provides recommendations and guidelines on formulating, designing, evaluating and analyzing federated optimization algorithms through concrete examples and practical implementation, with a focus on conducting effective simulations to infer real-world performance. The goal of this work is not to survey the current literature, but to inspire researchers and practitioners to design federated learning algorithms that can be used in various practical applications.

# A Field Guide to Federated Optimization

53 authors from 14 top institutions

[arxiv.org/abs/2107.06917](https://arxiv.org/abs/2107.06917)

Tensorflow Federated Implementation



**Thank you for your time!**

Twitter: @KairouzPeter