

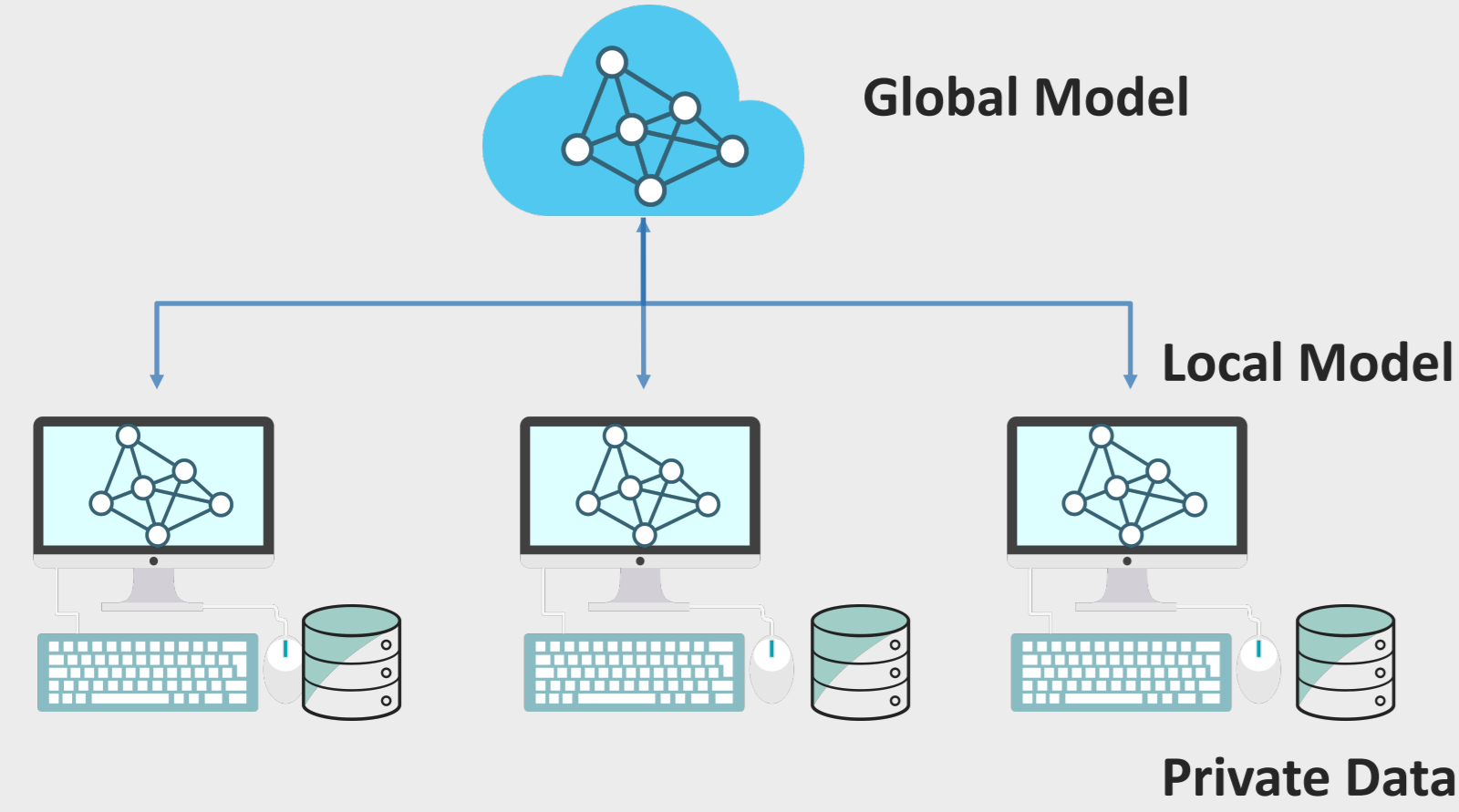
Federated Learning with Local Differential Privacy: Trade-offs between Privacy, Utility, and Communication

Muah Kim¹, Onur Günlü¹, and Rafael F. Schaefer²

¹Information Theory and Applications Chair, TU Berlin, Germany, ²Chair of Communications Engineering and Security, Universität Siegen, Germany
{muah.kim, guenlue}@tu-berlin.de, {rafael.schaefer}@uni-siegen.de

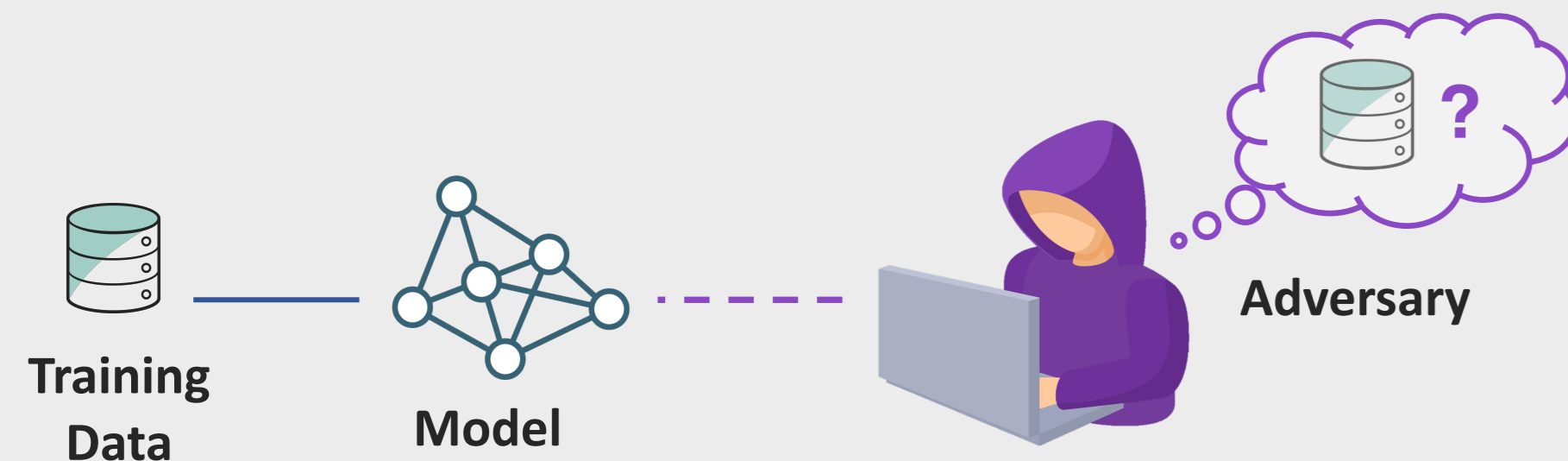
Motivation – Privacy of Machine Learning

1. Federated Learning (FL)



The global model is achieved by merging locally trained models.
⇒ More **private** than a centralized model: users' data remain at users.
⇒ It costs **iterative communication** for every weight updates.

2. Privacy Preserving Scheme



The training data can be partially recovered by using the model.
⇒ **Privacy preserving scheme** is required,
⇒ while minimizing learning **performance degradation** from it.

FL + privacy preserving scheme

Privacy, Utility, Transmission Rate should be jointly considered.
Privacy after T iterations should be tightly accounted.

Main Contribution

- We analyze the bounds of **privacy, utility, and transmission rate** of an FL model with stochastic gradient descent (SGD) algorithm and Gaussian mechanism.
- The **trade-offs** between three metrics are observed:
Privacy \uparrow – Utility \downarrow – Transmission Rate \uparrow
- The trade-offs are improved by adopting an enhanced privacy accounting method over many iterations. [1]
- A **generalized** FL model is assumed: heterogeneity of users, variable query sensitivity of the privacy mechanism, parameterized gradient norm clipping threshold, etc.

System Modeling

Algorithm: FL-SGD with Gaussian Mechanism

Input: User datasets $\{\mathcal{D}_k\}_{k=1}^K$, data sampling rates $\{q_k\}_{k=1}^K$, sampled datasets $\{\mathcal{J}_k^{(t)}\}_{k=1}^K$, total sampled dataset $\mathcal{J}^{(t)} = \cup_{k=1}^K \mathcal{J}_k^{(t)}$, loss function $\mathcal{L}_k(\mathbf{w}^{(t)}, \mathcal{J}_k^{(t)}) = \frac{1}{|\mathcal{J}_k^{(t)}|} \sum_{x \in \mathcal{J}_k^{(t)}} \ell(\mathbf{w}^{(t)}, x)$.
Parameters: learning rate η_t , noise scale $\{\sigma_k\}_{k=1}^K$, clipping norm threshold C .

Initialize $\mathbf{w}^{(0)}$ randomly
for $t \in [0: T - 1]$ **do** Training T iterations in total
for $k \in [K]$ Training user k 's local model

Download $\mathbf{w}^{(t)}$

Sample $\mathcal{J}_k^{(t)}$ from \mathcal{D}_k with probability q_k

Compute gradient

$$\mathbf{g}_k^{(t)}(\mathbf{w}^{(t)}, \mathcal{J}_k^{(t)}) \leftarrow \nabla_{\mathbf{w}} \mathcal{L}_k(\mathbf{w}^{(t)}, \mathcal{J}_k^{(t)})$$

Gradient Norm Clipping

$$\bar{\mathbf{g}}_k^{(t)} = \frac{\mathbf{g}_k^{(t)}}{\max(1, \|\mathbf{g}_k^{(t)}\|/C)}$$

Add Gaussian Noise

$$\tilde{\mathbf{g}}_k^{(t)} = \mathcal{M}_k(\bar{\mathbf{g}}_k^{(t)}) = \bar{\mathbf{g}}_k^{(t)} + \mathcal{N}(\mathbf{0}, C^2 \sigma_k^2 \mathbf{I}_d)$$

Upload $\tilde{\mathbf{g}}_k^{(t)}$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \cdot \sum_{k=1}^K \frac{|\mathcal{J}_k^{(t)}|}{|\mathcal{J}^{(t)}|} \tilde{\mathbf{g}}_k^{(t)}$$

Output: $\mathbf{w}^{(T)}$

Performance Metrics

Local Differential Privacy (LDP) of User k

Definition. (ϵ_k, δ_k) – **Local Differential Privacy (LDP)**
Gaussian mechanism \mathcal{M}_k is (ϵ_k, δ_k) –LDP w.r.t. dataset \mathcal{D}_k , if \forall two **neighboring** datasets $\mathcal{D} \sim \mathcal{D}' \subseteq \mathcal{D}_k$ and $\forall \mathcal{S} \subseteq \text{Range}(\mathcal{M}_k)$
$$\Pr[\mathcal{M}_k(\bar{\mathbf{g}}_k(\mathcal{D})) \in \mathcal{S}] \leq e^{\epsilon_k} \cdot \Pr[\mathcal{M}_k(\bar{\mathbf{g}}_k(\mathcal{D}')) \in \mathcal{S}] + \delta_k.$$

Gaussian Mechanism for (ϵ_k, δ_k) –LDP

A Gaussian mechanism $\mathcal{M}_k(\bar{\mathbf{g}}_k^{(t)}) = \bar{\mathbf{g}}_k^{(t)} + \mathcal{N}(\mathbf{0}, C^2 \sigma_k^2 \mathbf{I}_d)$ satisfies (ϵ_k, δ_k) –LDP
if $\delta_k \geq \frac{4}{5} \exp(-C \sigma_k \epsilon_k / 2)$.

Global Utility is defined by the multiplicative inverse of the convergence rate, i.e.,

$$\mathcal{U}(T) = \frac{1}{\mathbb{E}[\mathcal{L}(\mathbf{w}^{(T)}, \mathcal{J}^{(T)}) - \mathcal{L}(\mathbf{w}^*, \cup_{k=1}^K \mathcal{D}_k)]}$$

$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}^{(t)}, \cup_{k=1}^K \mathcal{D}_k)$: optimal weight vector

Transmission rate is defined by the **differential entropy** of a noisy gradient:
$$R_{tr,k} = h(\tilde{\mathbf{g}}_k^{(t)}).$$

Theoretical Analysis

Theoretical Bounds on the Metrics

Theorem 1. For $\epsilon_k > 2 \log(\delta_k^{-1}) \max\{\delta_k, \frac{1}{\sigma_k^2 \ln \frac{1}{q_k \sigma_k}}\}$,

$q_k < \frac{1}{16\sigma_k}$, and $\sigma_k \geq 1$, user k 's Gaussian mechanism \mathcal{M}_k is (ϵ_k, δ_k) –LDP after T iterations if

$$\sigma_k^2 \geq \frac{4q_k^2 T}{1 - q_k} \left[\frac{2}{\epsilon_k^2} \log \frac{1}{\delta_k} + \frac{1}{\epsilon_k} - \frac{2}{\epsilon_k^2} (\log(2 \log \delta_k^{-1}) + 1 - \log \epsilon_k) \right] + \mathcal{O}\left(\frac{\log^2(\log \delta_k^{-1})}{\log \delta_k^{-1}}\right).$$

For a μ – smooth and λ – strongly convex loss $\mathcal{L}(\mathbf{w}; \mathcal{S})$ w.r.t. a d – dimensional weight vector $\mathbf{w} \in \mathbb{R}^d$ given an arbitrary subset \mathcal{S} of $\mathcal{D} = \cup_{k=1}^K \mathcal{D}_k$, i.e., $\mathcal{S} \subseteq \mathcal{D}$, and for a learning rate $\eta_t = \frac{G}{C\lambda t}$,

$$\mathcal{U}(T) \geq \frac{\lambda^2 T}{\mu G^2} \min\left\{\frac{1}{2}, \frac{1}{d\sigma^2}\right\}$$

where $\sigma^2 = \frac{\sum_{k=1}^K (|\mathcal{D}_k| q_k \sigma_k^2)}{(\sum_{k=1}^K |\mathcal{D}_k| q_k)^2}$, and G is the maximum norm of the gradient. The transmission rate $R_{tr,k}$ of user k with Gaussian noise $\mathcal{N}(\mathbf{0}, C^2 \sigma_k^2 \mathbf{I}_d)$ satisfies

$$R_{tr,k} \leq d \log_2 \left(\frac{2\pi e C^2 \sigma_k}{\sqrt{d}} \right) \text{ (bits per gradient).}$$

Theoretical Analysis (cont'd)

- We set the number of iterations T and obtain the range of σ_k^2 that can achieve a target LDP level (ϵ_k, δ_k) .
- The noise variance σ_k^2 connects the target privacy (ϵ_k, δ_k) with the bounds of utility $\mathcal{U}(T)$ and transmission rate $R_{tr,k}$:
- Privacy \uparrow & $\epsilon_k \downarrow$ – $\sigma_k \uparrow$ – Utility \downarrow – Transmission Rate \uparrow .

⇒ **Trade-off relationship**

Comparison with other Privacy Accounting Methods

Moment Accountant (MA)

$$\sigma_k^2 \geq \frac{4q_k^2 T}{1 - q_k} \left[\frac{2}{\epsilon_k^2} \log \frac{1}{\delta_k} + \frac{1}{\epsilon_k} + \mathcal{O}(\log \delta_k^{-1}) \right]$$

Advanced Composition 1 (AC1)

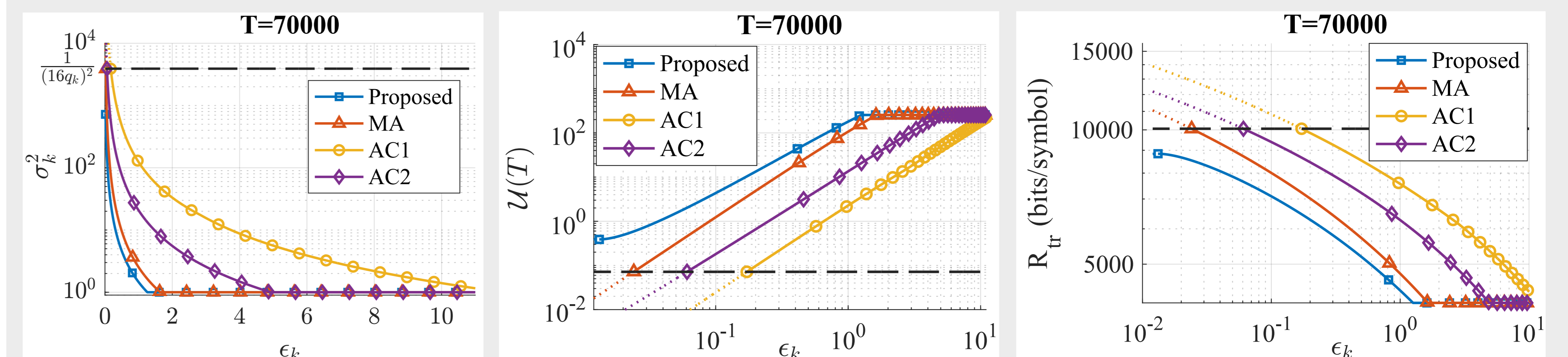
$$\sigma_k^2 \geq \frac{4q_k^2}{1 - q_k} \cdot \frac{2}{\epsilon_0} \log \left(\frac{4}{5\delta_0} \right)$$

Advanced Composition 2 (AC2)

$$\sigma_k^2 \geq \frac{4q_k^2}{1 - q_k} \cdot \frac{8T \log \left(e + \frac{\epsilon_k}{\delta_k} \right)}{\epsilon_k^2}$$

- These three methods are commonly used for privacy analysis.
- The bound of σ_k^2 in **Theorem 1.** is smaller than that of **(MA)**.
- (MA)** is shown to outperform **(AC1)** and **(AC2)**. [2]
- ⇒ The **proposed** bound of σ_k^2 is expected to be the smallest.

Interpretation by Simulations



- Parameters: $K = 100$, $\delta_k = 10^{-4}$, $q_k = 10^{-3}$ for all $k = 1, 2, \dots, 100$, $d = 10^4$, $\mu = 1$, $\lambda = 1$, $C = 1$, and $G = 5$.
- The bounds of σ_k^2 , **utility** $\mathcal{U}(T)$ and **transmission rate** $R_{tr,k}$ are plotted for varying ϵ_k by using σ_k^2 as a parameter.
- Trade-offs** between three metrics are observed: the utility increases and transmission rate decreases as ϵ_k grows, i.e., the target privacy becomes weaker.
- The proposed scheme accounts the differential privacy **the tightest** compared to using the other privacy accounting methods, which results in a smaller noise level σ_k^2 to achieve the same (ϵ_k, δ_k) –LDP after $T = 70,000$ iterations and, accordingly, a greater guaranteed value of utility and a smaller worst case transmission rate. ⇒ **Better Trade-offs**
- Due to the condition $q_k < \frac{1}{16\sigma_k}$, the domain of ϵ_k is restricted, and it restricts the ranges by the black dashed lines in the graphs. By the condition $\sigma_k \geq 1$, the curves show saturating behavior when σ_k reaches to 1 as ϵ_k increases.