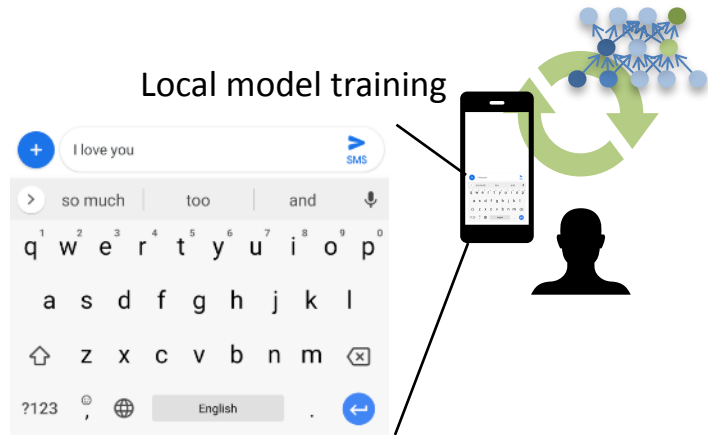# Dynamic Policies on Differentially Private Learning
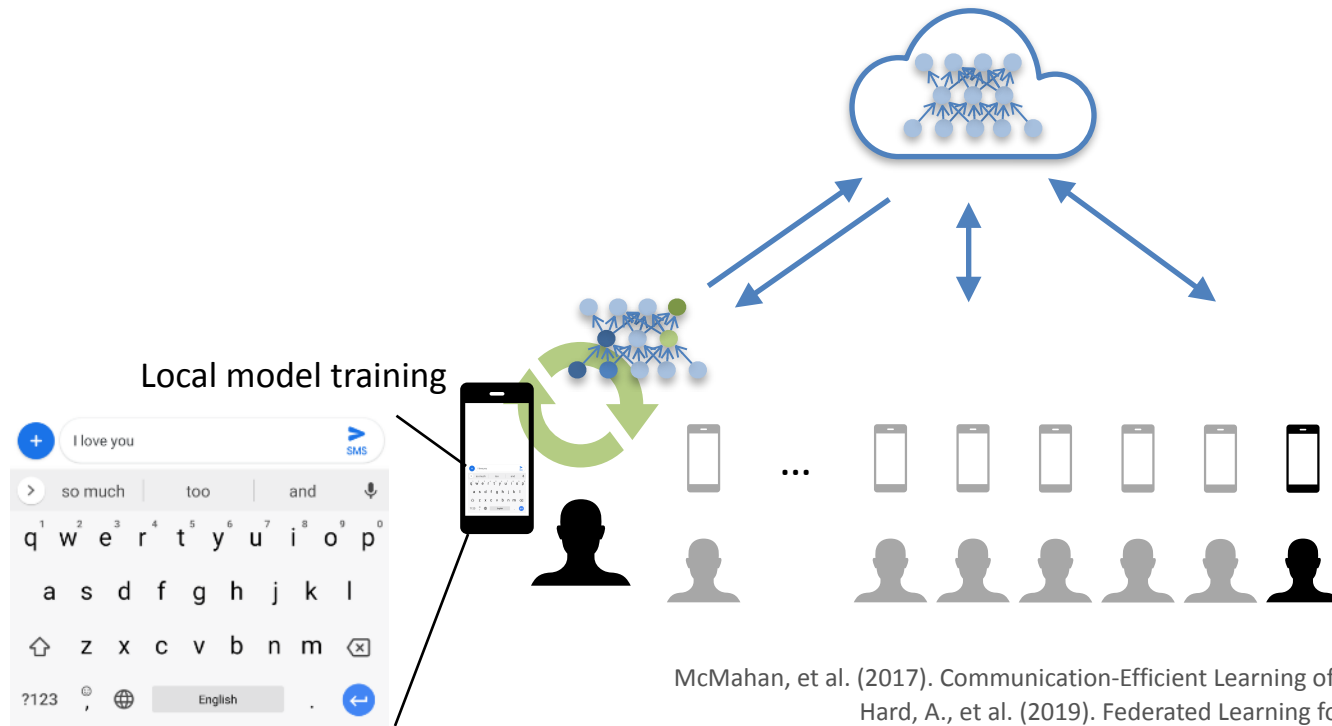
Junyuan Hong

Michigan State University
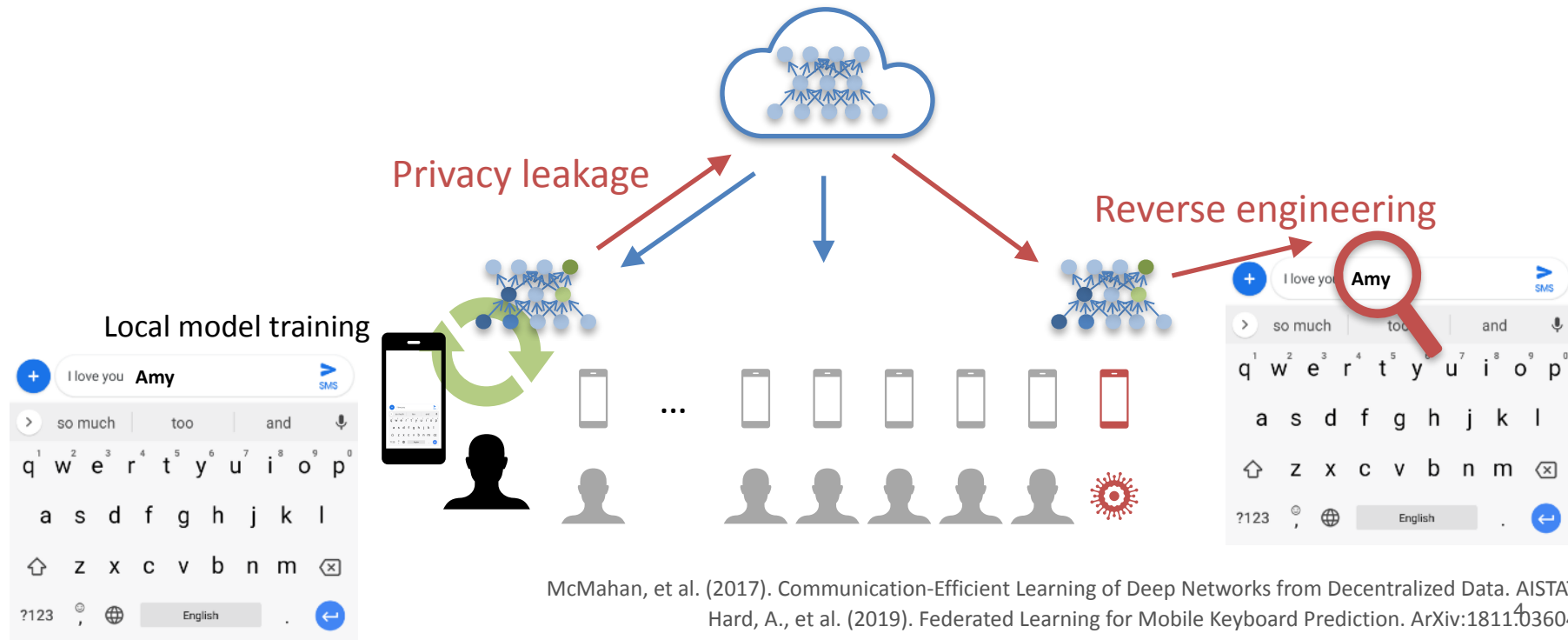
# Machine Learning in Our Life

Local model training

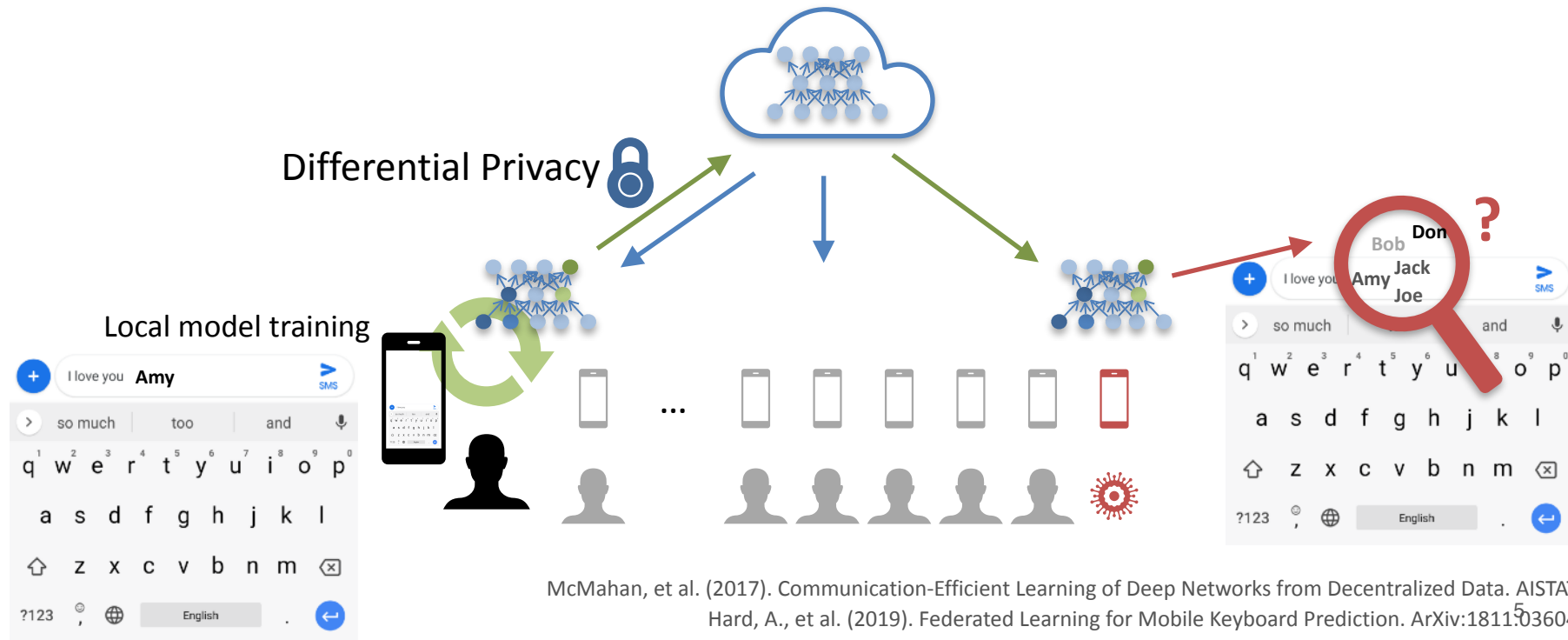# Federated Learning



Local model training

McMahan, et al. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. AISTAT
Hard, A., et al. (2019). Federated Learning for Mobile Keyboard Prediction. ArXiv:1811.03604

# Federated Learning



Privacy leakage

Reverse engineering

Local model training

McMahan, et al. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. AISTAT
Hard, A., et al. (2019). Federated Learning for Mobile Keyboard Prediction. ArXiv:1811.03604

4

# Federated Learning



Differential Privacy

Local model training

McMahan, et al. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. AISTAT
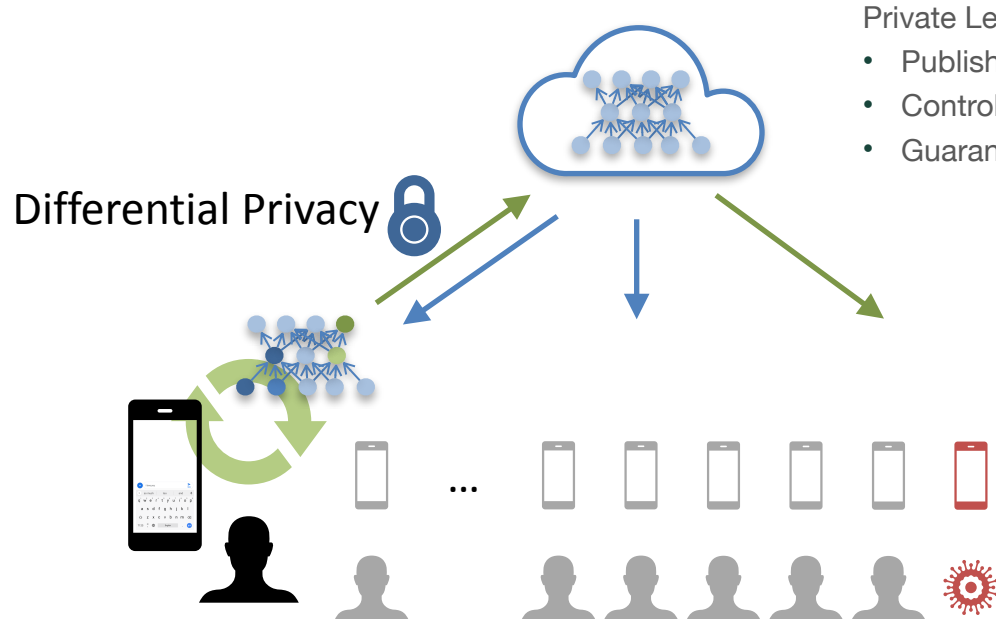Hard, A., et al. (2019). Federated Learning for Mobile Keyboard Prediction. ArXiv:1811.03604

5

# Federated Learning

Differential Privacy

Private Learning:
- Publish knowledge (model) rather than data
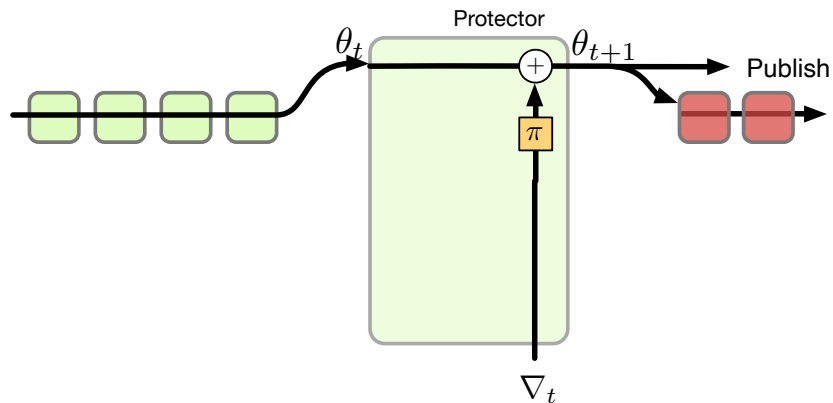- Control the privacy loss
- Guarantee the convergence

...

# Private Learning

Algorithm

Convergence theory and dynamic policy
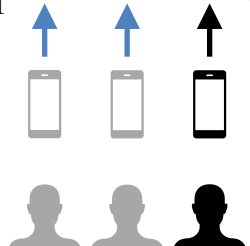
# Learning by Gradient Descent

$\rho$ privacy measure

$\pi$ projection: AdaGrad, etc

Protector

$\theta_t$ $\theta_{t+1}$ Publish

$\pi$

$\nabla_t$

$$\nabla_t = \frac{1}{N} \sum_{n=1}^{N} \boxed{\nabla f(\theta; x_n)}$$

Private sample
(to protect)

# Privacy attack

- **2019Grad**: Deep Leakage from Gradients, Zhu et al.: $x = \arg\min_{x} \|\nabla f(x) - \nabla_t\|^2$

- **2017MIA**: Membership Inference Attacks, Shokri et al.: $P(x \in D_{\text{train}}) = h(f(x; \theta))$ where $h()$ is a trained attack.

- **2017GAN**: Info Leakage from Collaborative Deep Learning, Hitaj et al. 2017: $x = G(z)$ where $z = \max_{z} f(G(z); \theta)$

- **2015MI**: Model Inversion, Fredrikson et al.: $x = \arg\max_{x} f(x)$ (statistical model)

prediction

2015MI

2019Grad $\nabla_t$

$D$

2017GAN

2017MIA

# Quantify privacy

If privacy cost is over a budget, we stop and publish model



$$\cdots + \rho_{t-3} + \rho_{t-2} + \rho_{t-1} + \rho_t \longrightarrow \text{Privacy cost}$$

$\theta_t \qquad \theta_{t+1} \qquad \text{Publish}$

$\cdots \quad \nabla_{t-3} \quad \nabla_{t-2} \quad \nabla_{t-1} \quad \nabla_t$

# Quantify privacy: Differential Privacy (DP)

# Differential Privacy

# Differential Privacy



**Privacy loss at** $y$    $Z(y) \triangleq \log \left( \dfrac{p(\mathscr{A}(D) = y)}{p(\mathscr{A}(D') = y)} \right)$    **where** $y \sim \mathscr{A}(D)$ **and** $D, D'$ **are adjacent (differing at one sample)**

13

# Differential Privacy

$$P(Z > t)$$



**Privacy loss at** $y$

$$Z(y) \triangleq \log \left( \frac{p(\mathscr{A}(D) = y)}{p(\mathscr{A}(D') = y)} \right)$$

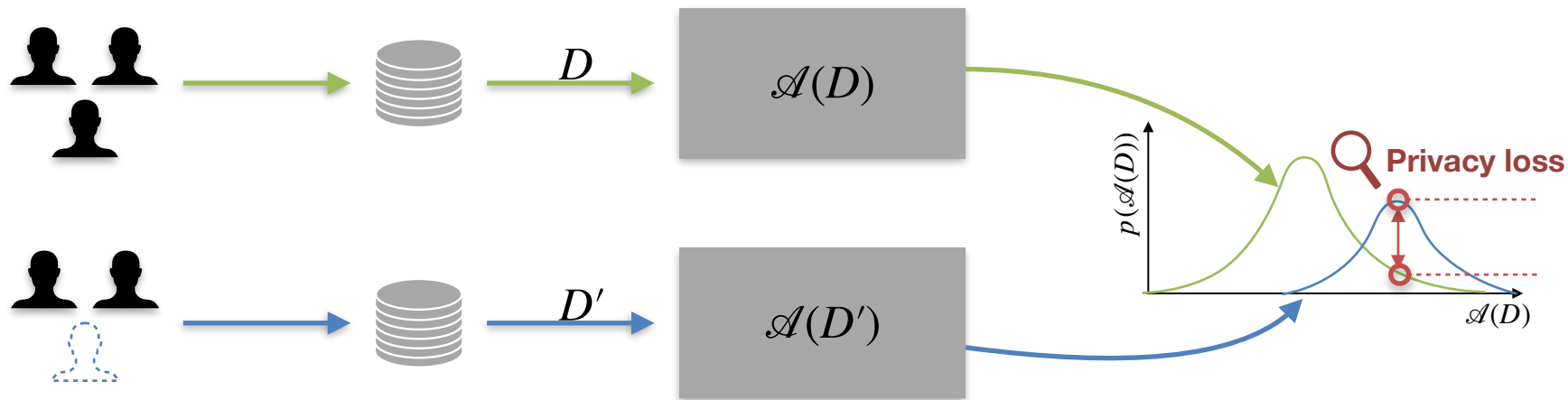**where** $y \sim \mathscr{A}(D)$

$\mathscr{A}$ is $\epsilon$-DP: $Z \leq \epsilon$ or $P(Z > \epsilon) = 0$

$\mathscr{A}$ is $(\epsilon, \delta)$-DP: $P(Z > \epsilon) = \delta$

$\mathscr{A}$ is $\rho$-zCDP: $P(Z > t + \rho) \leq e^{-t^2/(4\rho)}$ for $t \geq 0$

$\mathscr{A}$ is $(\rho, \omega)$-tCDP: $P(Z > t + \rho) \leq e^{-t^2/(4\rho)}$ for $t \in [0, 2\rho(\omega - 1)]$

$$P(Z > t + \rho) \leq e^{(\omega-1)^2\rho} \cdot e^{-(\omega-1))t} \text{ for } t \in (2\rho(\omega - 1), \infty)$$

14

Bun, M., Dwork, C., et al. (2018). Composable and Versatile Privacy via Truncated CDP. *STOC*

# Quantify privacy: Accumulate privacy loss

Compose **dynamic** privacy parameter



LEMMA 3.5. *(Composition) Suppose two mechanisms* $M, M'$ : $\mathcal{D}^n \to \mathbb{R}^d$ *satisfy* $\rho_1$-*zCDP and* $\rho_2$-*zCDP, then their composition satisfies* $(\rho_1 + \rho_2)$-*zCDP.*

Note: zCDP allows $\rho_1$ and $\rho_2$ to be different, but DP does not. For DP, an additional privacy cost has to be paid.

Bun, M., & Steinke, T. (2016). Concentrated Differential Privacy: Simplifications, Extensions, and Lower Bounds, TOC
Dwork, C., & Rothblum, G. N. (2016). Concentrated Differential Privacy. ArXiv:1603.01887
Rogers, et al. (2016). Privacy Odometers and Filters: Pay-as-you-Go Composition. NeurIPS

15

# Quantify privacy

$\pi$ is **deterministic**
which is non-private



$\boxed{\rho}$ privacy measure

$\boxed{\pi}$ projection: AdaGrad, etc

$$\nabla_t = \frac{1}{N} \sum_{n=1}^{N} \boxed{\nabla f(\theta; x_n)}$$

Private sample
(to protect)

16

# Privatize Gradients



Privacy loss $\rho_t$-zCDP

Protector

$\theta_t$     $\theta_{t+1}$     Publish

$\rho$    $\pi$

$\mathcal{N}$

$\sigma$

$C$

$\nabla_t$

$\rho$   privacy measure

$\pi$   projection: AdaGrad, etc

$\sigma$   **noise schedule**

$\mathcal{N}$   **noise distribution**

$C$   **sensitivity constraint**

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \dots, \nabla_t^{(n)}]$, residual privacy budget $R_t$
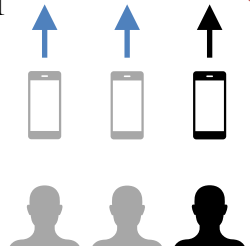
1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\|\}$      ▷ Sensitivity constraint

2: $\rho_t \leftarrow 1/\sigma_t^2$      ▷ Budget request

3: **if** $\rho_t < R_t$ **then**

4:     $R_{t+1} \leftarrow R_t - \rho_t$      Cost some privacy budget

5:     $g_t \leftarrow \nabla_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$      ▷ Privacy noise

6:     **return** $\eta_t g_t, R_{t+1}$      ▷ Utility projection

7: **else**

8:     Terminate

# Privatize Gradients



Privacy loss $\rho_t$-zCDP

Protector

$\theta_t$ $\qquad$ $\theta_{t+1}$ $\qquad$ Publish

$\rho$ $\quad$ $\pi$

$\mathcal{N}$

$\sigma$

$C$

$\nabla_t$

$\rho$ privacy measure

$\pi$ projection: AdaGrad, etc

$\sigma$ noise schedule

$\mathcal{N}$ noise distribution

$C$ **sensitivity constraint**

LEMMA 3.1 ($L_2$ SENSITIVITY). *Given mapping from a n-element dataset domain to d-dimensional real space $f : \mathcal{D}^n \to \mathbb{R}^d$, the $L_2$ sensitivity of $f$, denoted by $\Delta_2(f)$ is defined as:*

$$\Delta_2(f) = \max_{D, D'} \left\| f(D) - f(D') \right\|_2 ,$$

*where $D, D'$ are adjacent datasets.*

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\|\}$ $\qquad$ ▷ Sensitivity constraint

Control the influence of a sample

2: $\rho_t \leftarrow 1/\sigma_t^2$ $\qquad$ ▷ Budget request
3: **if** $\rho_t < R_t$ **then**
4: $\quad$ $R_{t+1} \leftarrow R_t - \rho_t$
5: $\quad$ $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$ $\qquad$ ▷ Privacy noise
6: $\quad$ **return** $\eta_t g_t, R_{t+1}$ $\qquad$ ▷ Utility projection
7: **else**
8: $\quad$ Terminate

18

# Differentially Private Learning



Privacy loss

$\rho_t$-zCDP

Protector

$\theta_t$    $\theta_{t+1}$    Publish

$\rho$   $\pi$

$\mathcal{N}$

$\sigma$

$C$

$\nabla_t$

$\boxed{\rho}$ privacy measure

$\boxed{\pi}$ projection: AdaGrad, etc

$\boxed{\sigma}$ **noise schedule**

$\boxed{\mathcal{N}}$ **noise distribution**

$\boxed{C}$ sensitivity control

LEMMA 3.4. *The Gaussian mechanism, which returns* $\boxed{f(D)}$ $+ \sigma v$ *satisfies* $\Delta_2(f)^2/(2\sigma^2)$-*zCDP.*
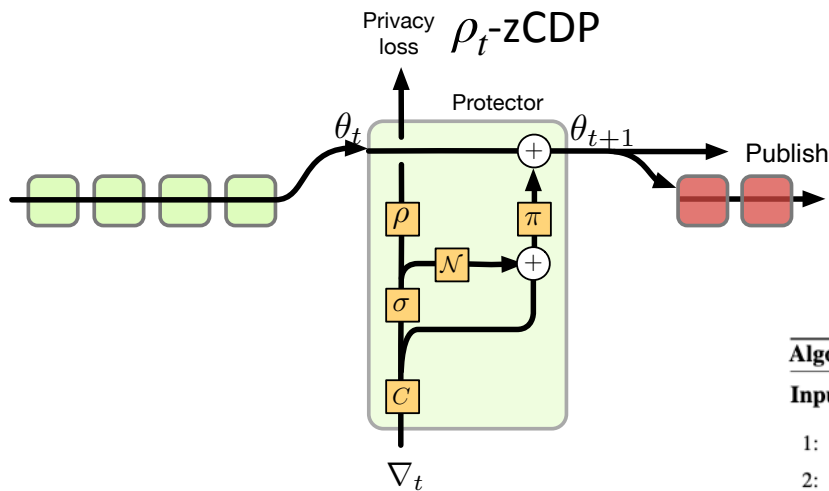
A deterministic function

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\|\}$      ▷ Sensitivity constraint

2: $\rho_t \leftarrow 1/\sigma_t^2$      ▷ Budget request

3: **if** $\rho_t < R_t$ **then**

4:      $R_{t+1} \leftarrow R_t - \rho_t$

5:      $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$      ▷ Privacy noise

6:      **return** $\eta_t g_t, R_{t+1}$      ▷ Utility projection

7: **else**

8:      Terminate

# Differentially Private Learning



Privacy loss $\rho_t$-zCDP

Protector

$\theta_t$    $\theta_{t+1}$    Publish

$\rho$   $\pi$

$\mathcal{N}$

$\sigma$

$C$

$\nabla_t$

$\boxed{\rho}$ privacy measure

$\boxed{\pi}$ projection: AdaGrad, etc

$\boxed{\sigma}$ noise schedule

$\boxed{\mathcal{N}}$ noise distribution

$\boxed{C}$ sensitivity control

If gradients are a stochastic mini-batch, e.g., sampled by q-probability, the privacy cost is $\propto q^2 \rho$ for DP metric, e.g, tCDP.

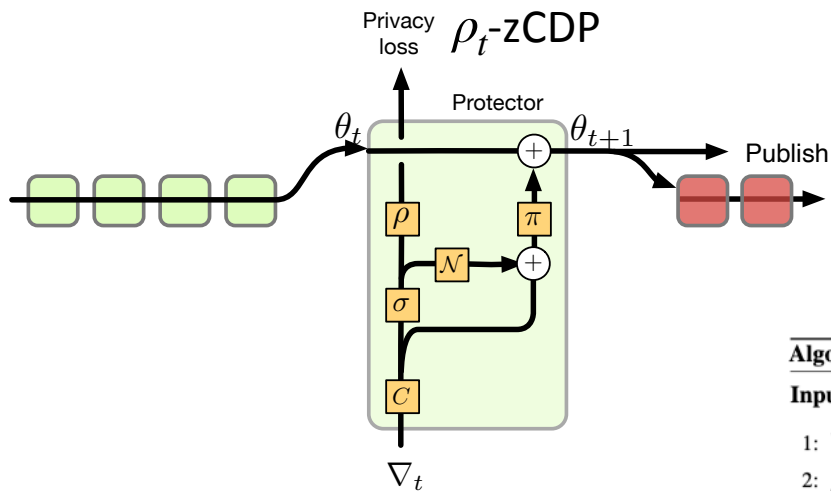**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$ residual privacy budget $R_t$

1:   $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\| \}$    ▷ Sensitivity constraint

2:   $\rho_t \leftarrow 1/\sigma_t^2$    ▷ Budget request

3:   **if** $\rho_t < R_t$ **then**

4:      $R_{t+1} \leftarrow R_t - \rho_t$

5:      $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$    ▷ Privacy noise

6:      **return** $\eta_t g_t, R_{t+1}$    ▷ Utility projection
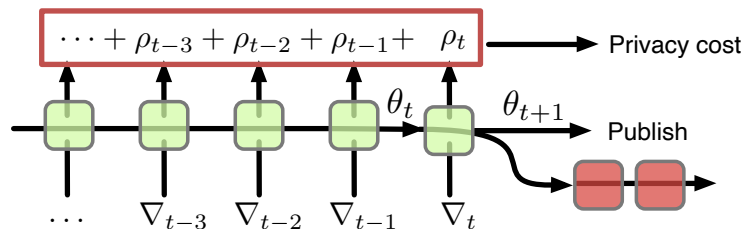
7:   **else**

8:      Terminate

# Privatize Gradients



$\rho_t$-zCDP

Privacy loss

Protector

$\theta_t$    $\theta_{t+1}$

Publish

$\rho$    privacy measure

$\pi$    projection: AdaGrad, etc

$\sigma$    **noise schedule**

$\mathcal{N}$    **noise distribution**

$C$    **sensitivity constraint**

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t/\left\|\nabla_t^{(n)}\right\|\}$        ▷ Sensitivity constraint
2: $\rho_t \leftarrow 1/\sigma_t^2$        ▷ Budget request
3: **if** $\rho_t < R_t$ **then**
4:      $R_{t+1} \leftarrow R_t - \rho_t$
5:      $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$        ▷ Privacy noise
6:      **return** $\eta_t g_t, R_{t+1}$        ▷ Utility projection
7: **else**
8:      Terminate

# Differentially Private Learning



$$\cdots + \rho_{t-3} + \rho_{t-2} + \rho_{t-1} + \rho_t \longrightarrow \text{Privacy cost}$$

$\theta_t$ $\theta_{t+1}$ Publish

$\nabla_{t-3}$ $\nabla_{t-2}$ $\nabla_{t-1}$ $\nabla_t$

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^N \nabla_t^{(n)} \min\{1, C_t/\left\|\nabla_t^{(n)}\right\|\}$      ▷ Sensitivity constraint
2: $\rho_t \leftarrow 1/\sigma_t^2$      ▷ Budget request
3: **if** $\rho_t < R_t$ **then**
4:      $R_{t+1} \leftarrow R_t - \rho_t$
5:      $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t/N, \nu_t \sim \mathcal{N}(0, I)$      ▷ Privacy noise
6:      **return** $\eta_t g_t, R_{t+1}$      ▷ Utility projection
7: **else**
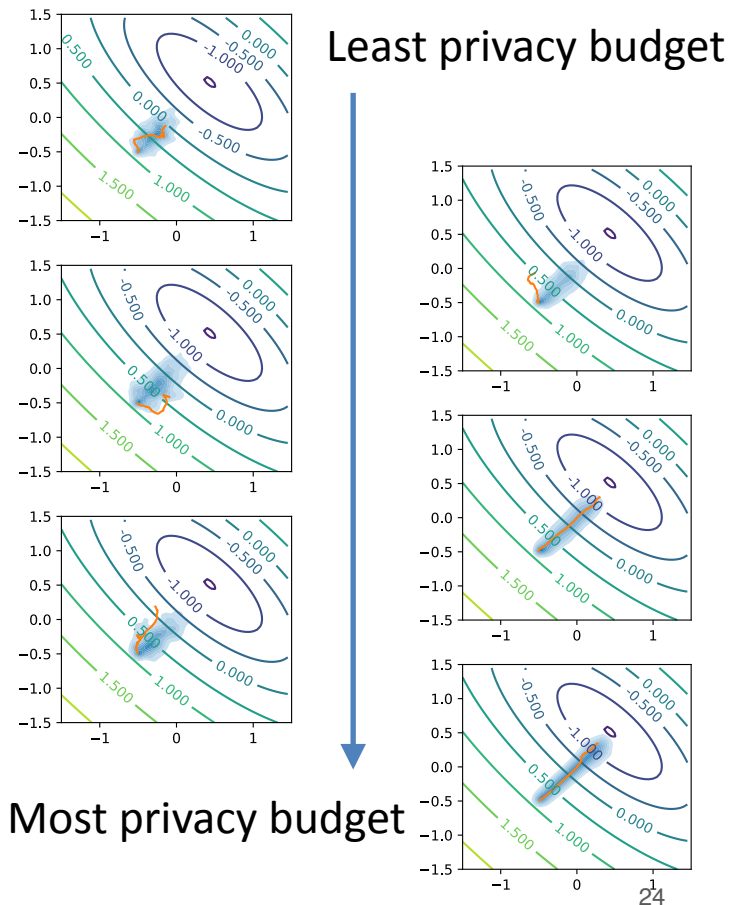8:      Terminate

# Private Learning

Algorithm

**Convergence theory and dynamic policy**

# Does private learning converge?

- Not converge to the optimal
  - Finite iteration
  - Noise
- Improve the final iterate loss given a privacy budget:

$$\mathrm{EER} = \mathbb{E}_\nu[f(\theta_{T+1})] - f(\theta^*)$$

  - The upper bound of EER

Least privacy budget

Most privacy budget



24

# Why study convergence upper bound?

- Bound the worst case.

- Find a way to speed up optimization algorithm

- To study the impact of privacy operations, e.g., noise magnitude, clipping norm, etc.

- To compare different algorithms: convergence rate

# Assumptions

- $G$-Lipschitz continuous loss,

  $\left\| f(x) - f(x') \right\| \leq G\|x - x'\| \Leftrightarrow \|f'(x)\| \leq G$ if $f$ is differentiable.

- $M$-Lipschitz continuous gradient or $M$-smooth loss:

  $\left\| \nabla f(x) - \nabla f(x') \right\| \leq M\|x - x'\|$

- $\mu$-Polyak-Lojasiewicz (PL) condition $< \mu$-strongly convex

  $\left\| \nabla f(\theta) \right\|^2 \geq 2\mu(f(\theta) - f(\theta*))$

# Convergence

**Algorithm 1** Privatizing gradients

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \ldots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\|\}$      ▷ Sensitivity constraint
2: $\rho_t \leftarrow 1/\sigma_t^2$      ▷ Budget request
3: **if** $\rho_t < R_t$ **then**
4:      $R_{t+1} \leftarrow R_t - \rho_t$
5:      $g_t \leftarrow \tilde{\nabla}_t + C_t \sigma_t \nu_t / N, \nu_t \sim \mathcal{N}(0, I)$      ▷ Privacy noise
6:      **return** $\eta_t g_t, R_{t+1}$      ▷ Utility projection
7: **else**
8:      Terminate

**Theorem 3.2.** *Let $\alpha$, $\kappa$ and $\gamma$ be defined in Eq. (5), and $\eta_t = \frac{1}{M}$. Suppose $f(\theta; x_i)$ is G-Lipschitz M-smooth and satisfies the Polyak-Lojasiewicz condition. If $C_t \leq G$, then clipping does not take place, i.e., $\tilde{\nabla}_t = \nabla_t$ and the following holds:*

$$\text{EER} = \mathbb{E}_\nu[f(\theta_{T+1})] - f(\theta^*) \leq \left( \gamma^T + R \sum_{t=1}^{T} q_t \sigma_t^2 \right) (f(\theta_1) - f(\theta^*)), \quad (6)$$

$$\text{where } q_t \triangleq \gamma^{T-t} \alpha_t. \quad (7)$$
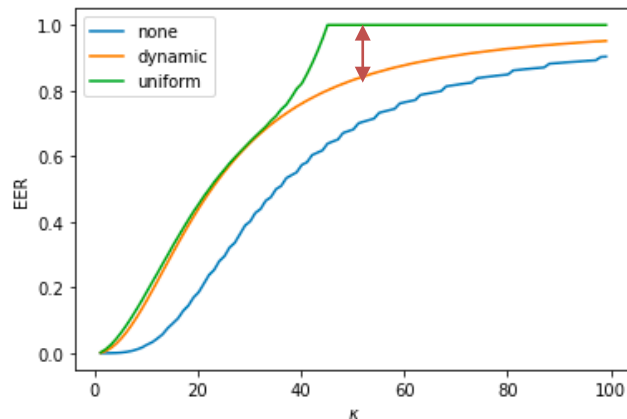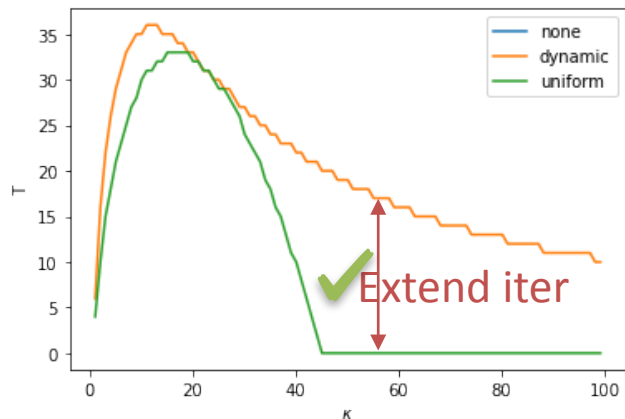
$$\alpha_t \triangleq \frac{MD}{2R} \left( \frac{\eta_t C_t}{N} \right)^2 \frac{1}{f(\theta_1) - f(\theta^*)} > 0, \ \kappa \triangleq \frac{M}{\mu} \geq 1, \text{ and } \gamma \triangleq 1 - \frac{1}{\kappa} \in [0, 1). \quad (5)$$

27

# Convergence

**Theorem 3.2.** *Let $\alpha$, $\kappa$ and $\gamma$ be defined in Eq. (5), and $\eta_t = \frac{1}{M}$. Suppose $f(\theta; x_i)$ is G-Lipschitz M-smooth and satisfies the Polyak-Lojasiewicz condition. If $\tilde{C}_t \leq G$, then clipping does not take place, i.e., $\tilde{\nabla}_t = \nabla_t$ and the following holds:*

$$\text{EER} \leq \left( \gamma^T + R \sum_{t=1}^{T} q_t \sigma_t^2 \right) (f(\theta_1) - f(\theta^*)), \tag{6}$$

$$\text{where } q_t \triangleq \gamma^{T-t} \alpha_t. \tag{7}$$

Finite iteration — Noise impact

- Schedule noise to
  - Extend iteration T
  - Reduce the effect of noise

28

# **Convergence**

**Theorem 3.2.** Let $\alpha$, $\kappa$ and $\gamma$ be defined in Eq. (5), and $\eta_t = \frac{1}{M}$. Suppose $f(\theta; x_i)$ is G-Lipschitz M-smooth and satisfies the Polyak-Lojasiewicz condition. If $\tilde{C}_t \le G$, then clipping does not take place, i.e., $\tilde{\nabla}_t = \nabla_t$ and the following holds:

$$\text{EER} \le \left( \gamma^T + R \sum_{t=1}^{T} q_t \sigma_t^2 \right) (f(\theta_1) - f(\theta^*)), \tag{6}$$

$$\text{where } q_t \triangleq \gamma^{T-t} \alpha_i. \tag{7}$$

Influence of noise

**Lemma 3.1** (Dynamic schedule). Suppose $\sigma_t$ satisfy $\sum_{t=1}^{T} \sigma^{-2} = R$. Given a positive sequence $\{q_t\}$, the following equation holds

✓Reduce noise impact $\quad \min_{\sigma} R \sum_{t=1}^{T} q_t \sigma_t^2 = \left( \sum_{t=1}^{T} \sqrt{q_t} \right)^2, \text{ when } \sigma_t = \sqrt{\frac{1}{R} \sum_{i=1}^{T} \sqrt{\frac{q_i}{q_t}}}. \tag{10}$

How much improvement can we achieve?

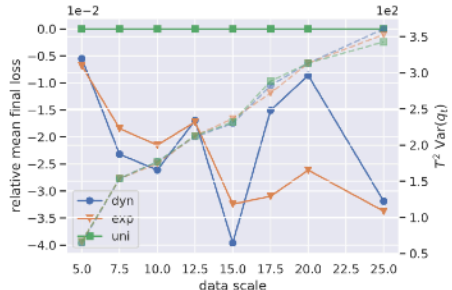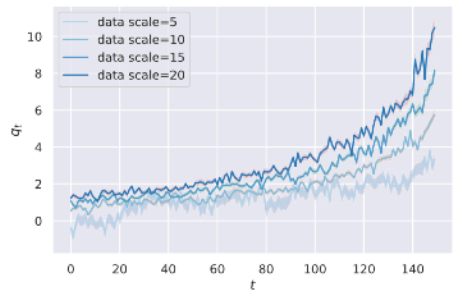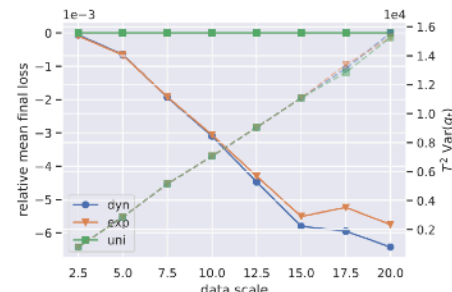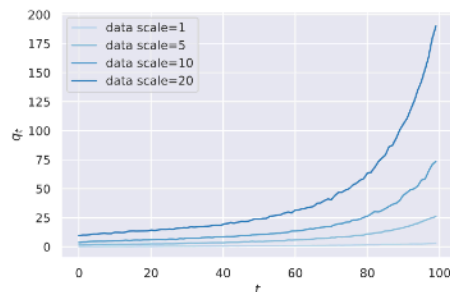# Advantage of dynamic schedule on optimal upper bound



✅ Extend iter



stable when the loss curvature is sharp

# Advantage of dynamic schedule

- Empirically check the $q_t$

$$\text{EER} \leq \left( \gamma^T + R \sum_{t=1}^{T} q_t \sigma_t^2 \right) (f(\theta_1) - f(\theta^*)),$$
$$where \; q_t \triangleq \gamma^{T-t} \alpha_t.$$



31

# Further reduce the noise by momentum

**Algorithm 2** Privatizing gradients with debiased momentum

**Input**: Private gradient $\nabla_t$ summed from $[\nabla_t^{(1)}, \dots, \nabla_t^{(n)}]$, residual privacy budget $R_t$

1: $\tilde{\nabla}_t \leftarrow \frac{1}{N} \sum_{n=1}^{N} \nabla_t^{(n)} \min\{1, C_t / \left\| \nabla_t^{(n)} \right\|\}$      ▷ Sensitivity constraint

2: $\rho_t \leftarrow 1/\sigma_t^2$      ▷ Budget request

3: **if** $\rho_t < R_t$ **then**

4:      $R_{t+1} \leftarrow R_t - \rho_t$

5:      $g_t \leftarrow \tilde{\nabla}_t + \nu_t, \; \nu_t \sim \mathcal{N}(0, (C_t \sigma_t / N)^2 I)$      ▷ Privacy noise

6:      $v_{t+1} = \beta v_t + (1 - \beta) g_t, \; v_1 = 0$

7:      $\hat{v}_{t+1} = v_{t+1}/(1 - \beta^t)$

8:      **return** $\eta_t \hat{v}_{t+1}, R_{t+1}$      ▷ Utility projection

9: **else**

10:      Terminate

# Further reduce the noise by momentum



**Theorem 3.4** (Convergence under PL condition). *Suppose $f(\theta; x_i)$ is $M$-smooth, $G$-Lipschitz and satisfies the Polyak-Lojasiewicz condition. Let $\eta_t = \eta_0$. If $C_t \geq G$ which implies $\tilde{\nabla}_t = \nabla_t$ (clipping does not take place), then the following holds:*

$$\text{EER} \leq \gamma^T(f(\theta_1) - f(\theta^*)) + \frac{2\eta_0 D}{N^2}\underbrace{\sum_{t=1}^{T} q_t (C_t \sigma_t)^2}_{noise\ varinace} + \eta_0 \zeta \underbrace{\sum_{t=1}^{T} \gamma^{T-t}\|v_{t+1}\|^2}_{momentum\ effect} \quad (16)$$
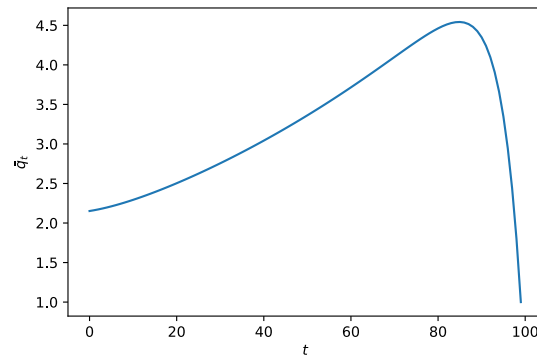
$$\text{where } q_t = \frac{\beta^{2(T-t+1)} - \gamma^{T-t+1}}{\beta^2 - \gamma}, \ \gamma = 1 - \eta_0\mu, \ \zeta = \frac{4M^2\beta\gamma}{(\gamma-\beta)^2(1-\beta)^3}\eta_0^2 + \frac{1}{2}M\eta_0 - 1. \quad (17)$$

*Especially, when $\eta_0 \leq \frac{\beta(1-\beta)^3}{8M}\left[\sqrt{\frac{1}{4} + \frac{16}{\beta(1-\beta)^3}} - 1\right]$, the noise variance dominates the bound, i.e.,*

$$\text{EER} = \mathcal{O}\left(\frac{2\eta_0 D}{N^2}\sum_{t=1}^{T} q_t (C_t \sigma_t)^2\right).$$
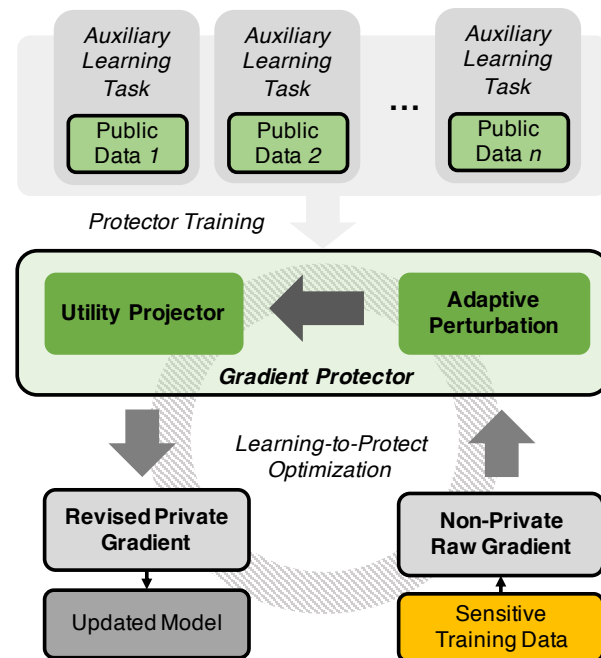
A negative term if $\eta_0$ is small.
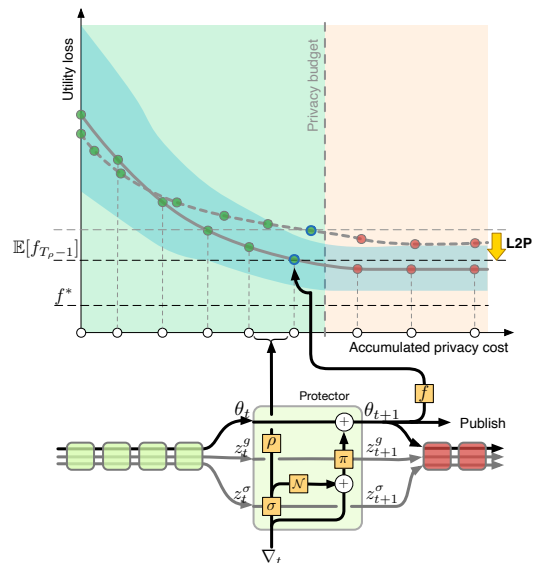
The GD noise

# Beyond dynamic noise magnitude

- Learning to protect: Transfer the dynamic policies learned from auxiliary tasks to private task.
- AdaClip (Pichapati et al. 2019): Adaptively clipping the gradients
- Dynamic batch size (Feldman et al., 2019, STOC): Increase the batch size to improve non-convex convergence bound.



34

# Beyond dynamic noise magnitude

- **Learning to protect**: Transfer the dynamic policies learned from auxiliary tasks to private task.

- AdaClip (Pichapati et al. 2019): Adaptively clipping the gradients

- Dynamic batch size (Feldman et al., 2019, STOC): Increase the batch size to improve non-convex convergence bound.



$$\min_{\pi,\sigma,T} \mathbb{E}\left[\tilde{F}(\sigma,\pi,T)\right], \text{ s.t. } h_T(\sigma;\rho_{\text{tot}}) = 0$$

# Beyond dynamic noise magnitude

- **Learning to protect**: Transfer the dynamic policies learned from auxiliary tasks to private task.

- AdaClip (Pichapati et al. 2019): Adaptively clipping the gradients

- Dynamic batch size (Feldman et al., 2019, STOC): Increase the batch size to improve non-convex convergence bound.
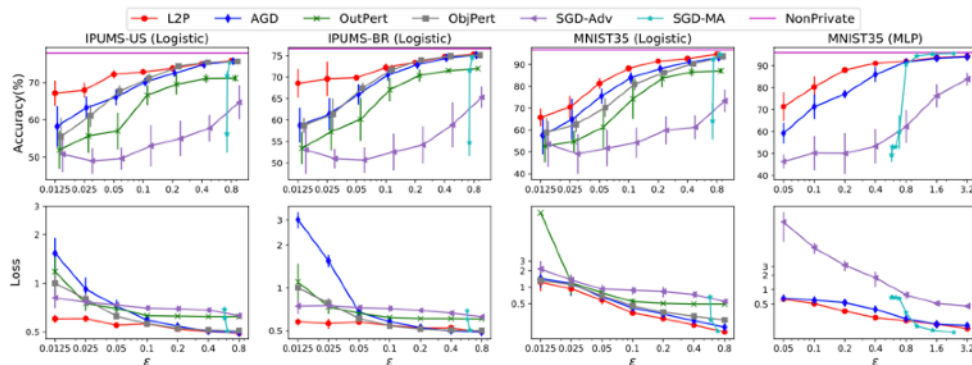


Figure 2: Test performance (top) and training loss values (bottom) by varying $\epsilon$ of logistic and MLP classifiers on IPUMS and MNIST35 datasets. The error bar presents the size of standard deviations. For better visualization, some horizontal offsets are added to every point.

$$\min_{\pi,\sigma,T} \mathbb{E}\left[\tilde{F}(\sigma,\pi,T)\right], \text{ s.t. } h_T(\sigma;\rho_{\text{tot}}) = 0$$

36

# Thank you for your time!