

The Saddle-Point Accountant for Differential Privacy

Shahab Asoodeh

Joint work with Wael Alghamdi, Felipe Gomez, Flavio Calmon (Harvard University),
Oliver Kosut, Lalitha Sankar (ASU)

October 19, 2022



زن_زندگی_آزادی

#Women-Life-Freedom

Joint work with



Wael Alghamdi
(Harvard)



Flavio Calmon
(Harvard)



Felipe Gomez
(Harvard)



Oliver Kosut
(ASU)



Lalitha Sankar
(ASU)

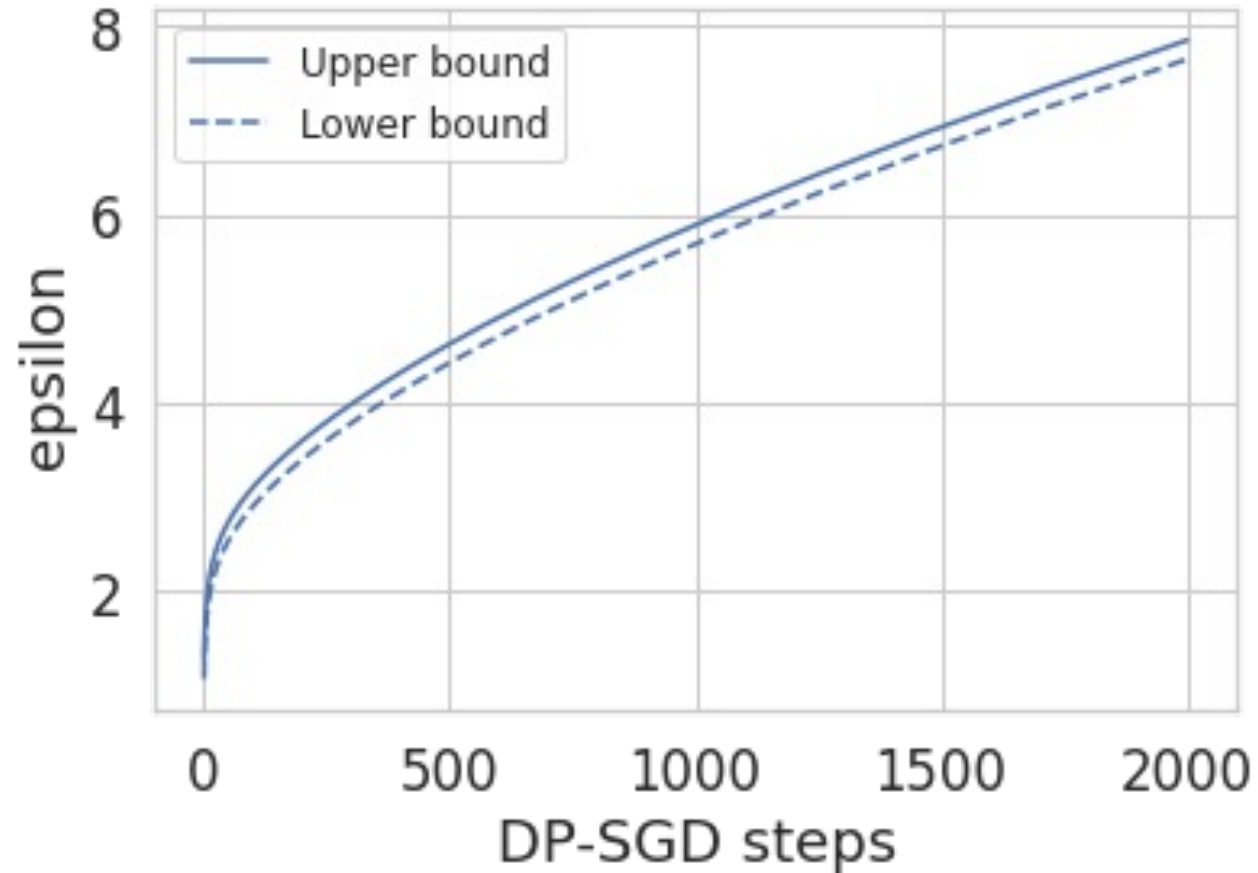
State-of-the-art composition

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$



Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

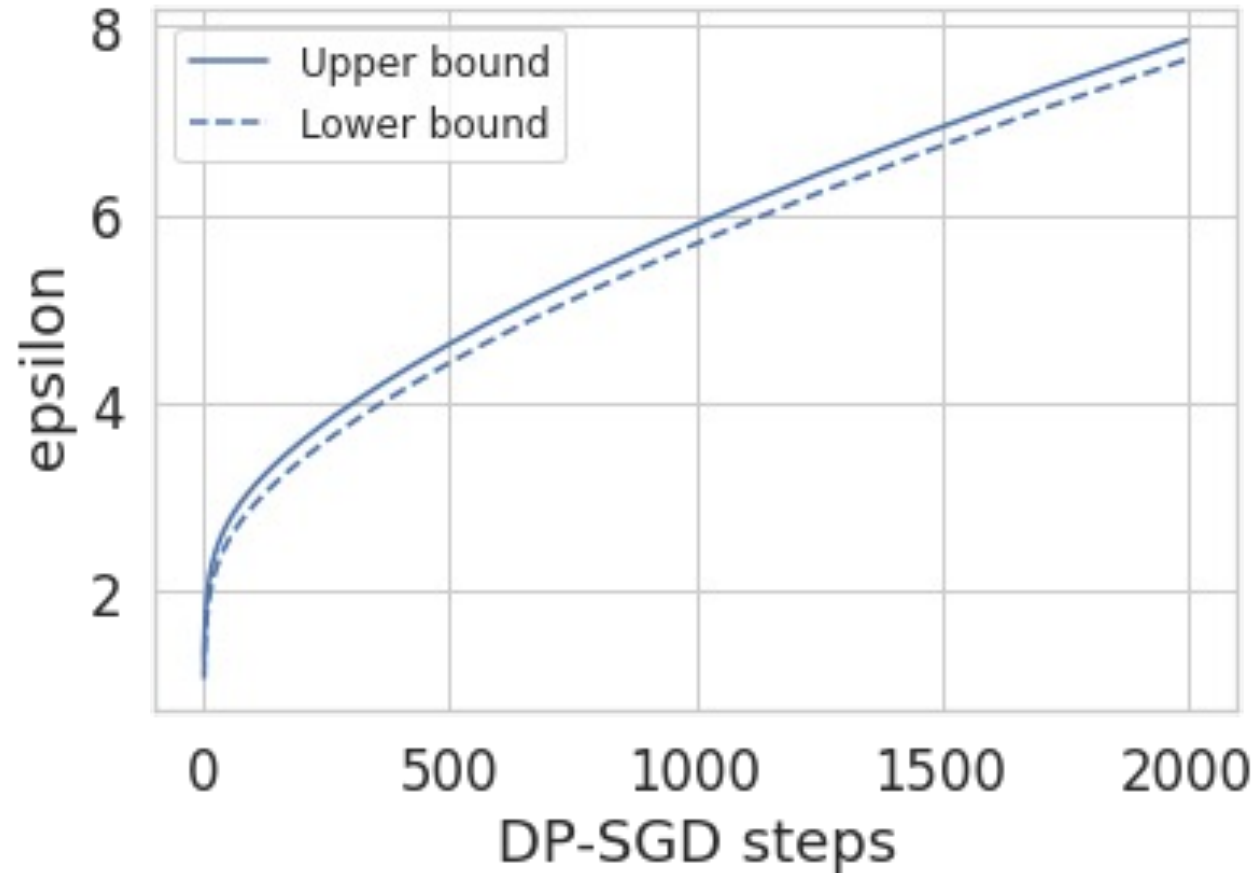
State-of-the-art composition

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$



runtime complexity

$$O(\sqrt{n} \log n)$$

Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

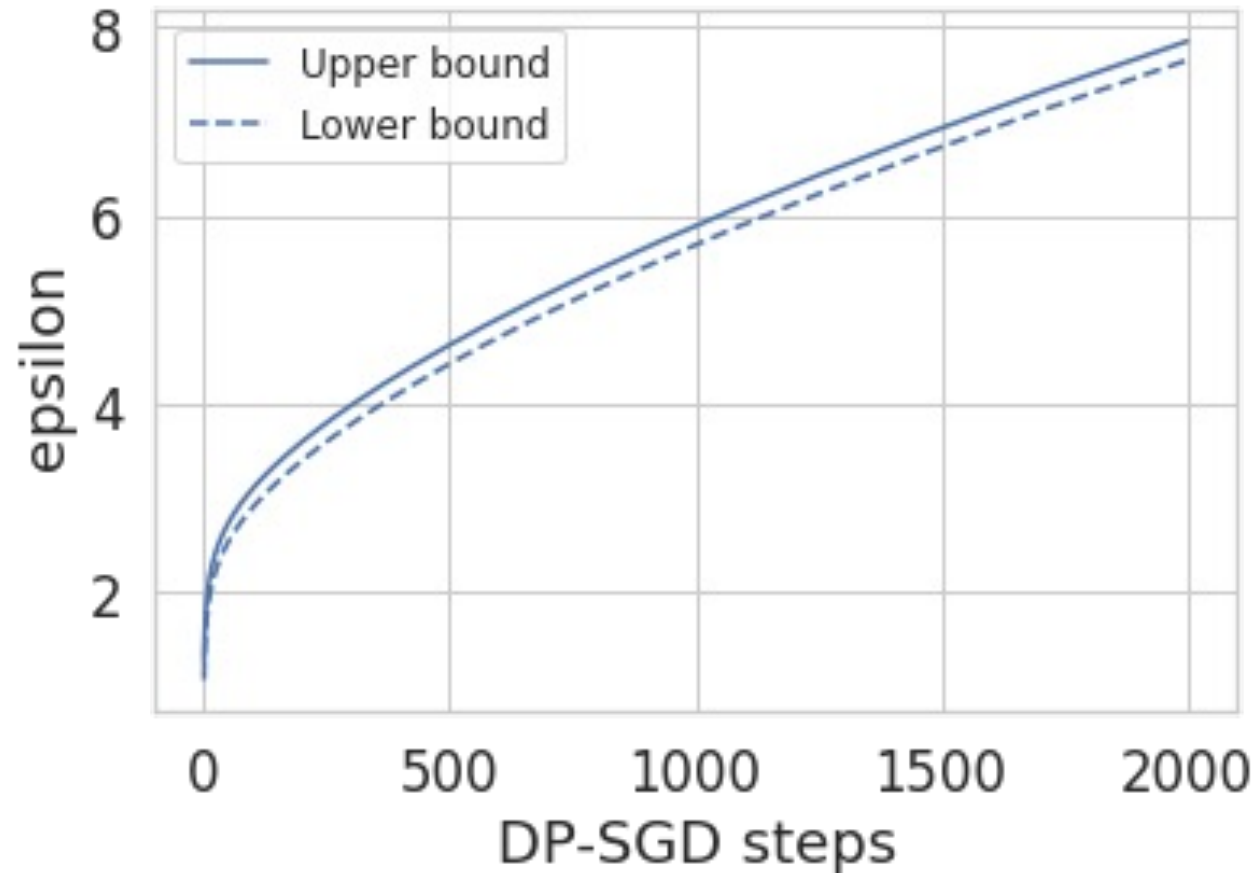
State-of-the-art composition

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$



runtime complexity

$$O(\sqrt{n} \log n)$$



$$O(\text{polylog}(n))$$

Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

Ghazi, Kamath, Kumar, and Manurangsi, Faster Privacy Accounting via Evolving Discretization, ICML 2022

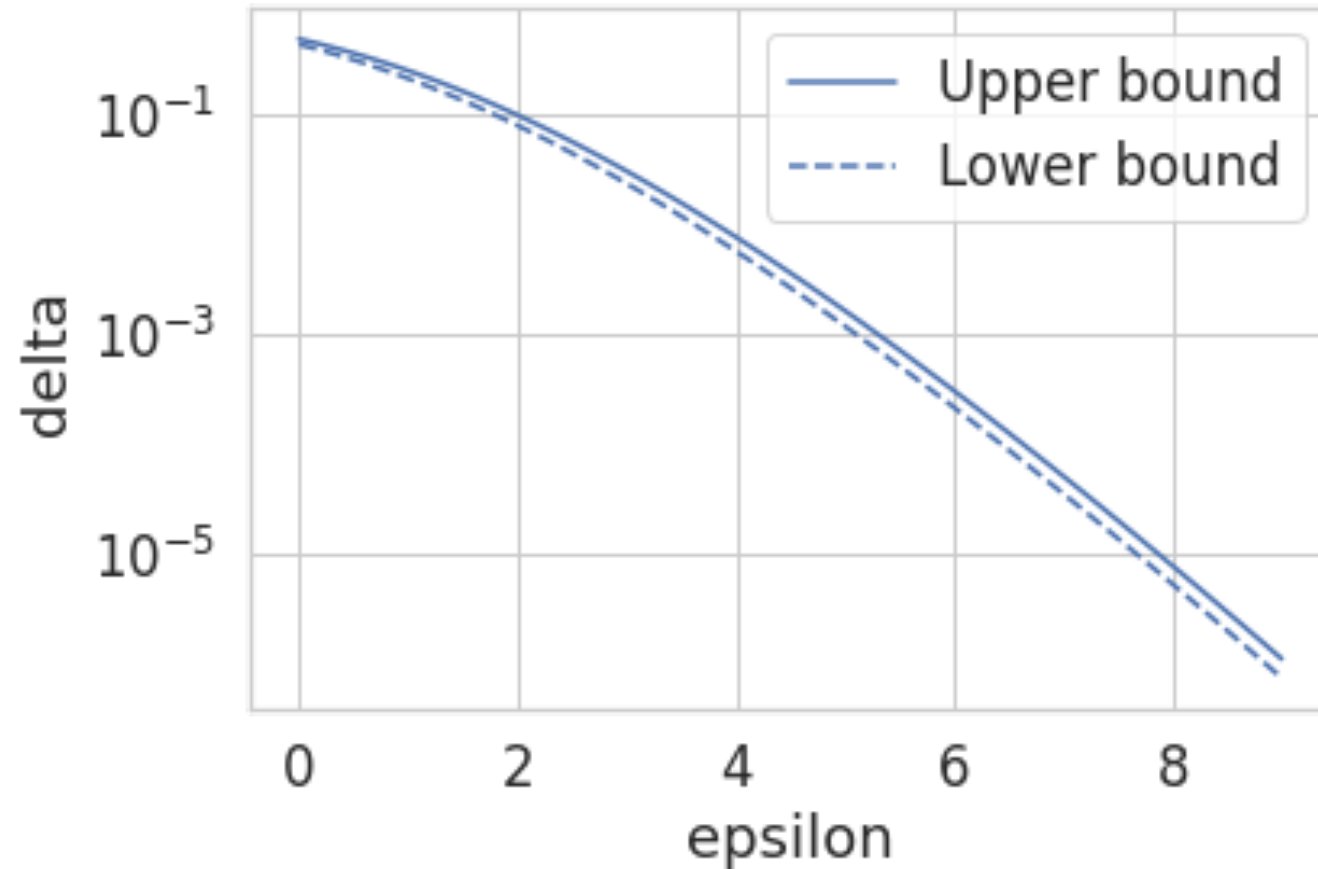
State-of-the-art composition

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$n = 2000$$



Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

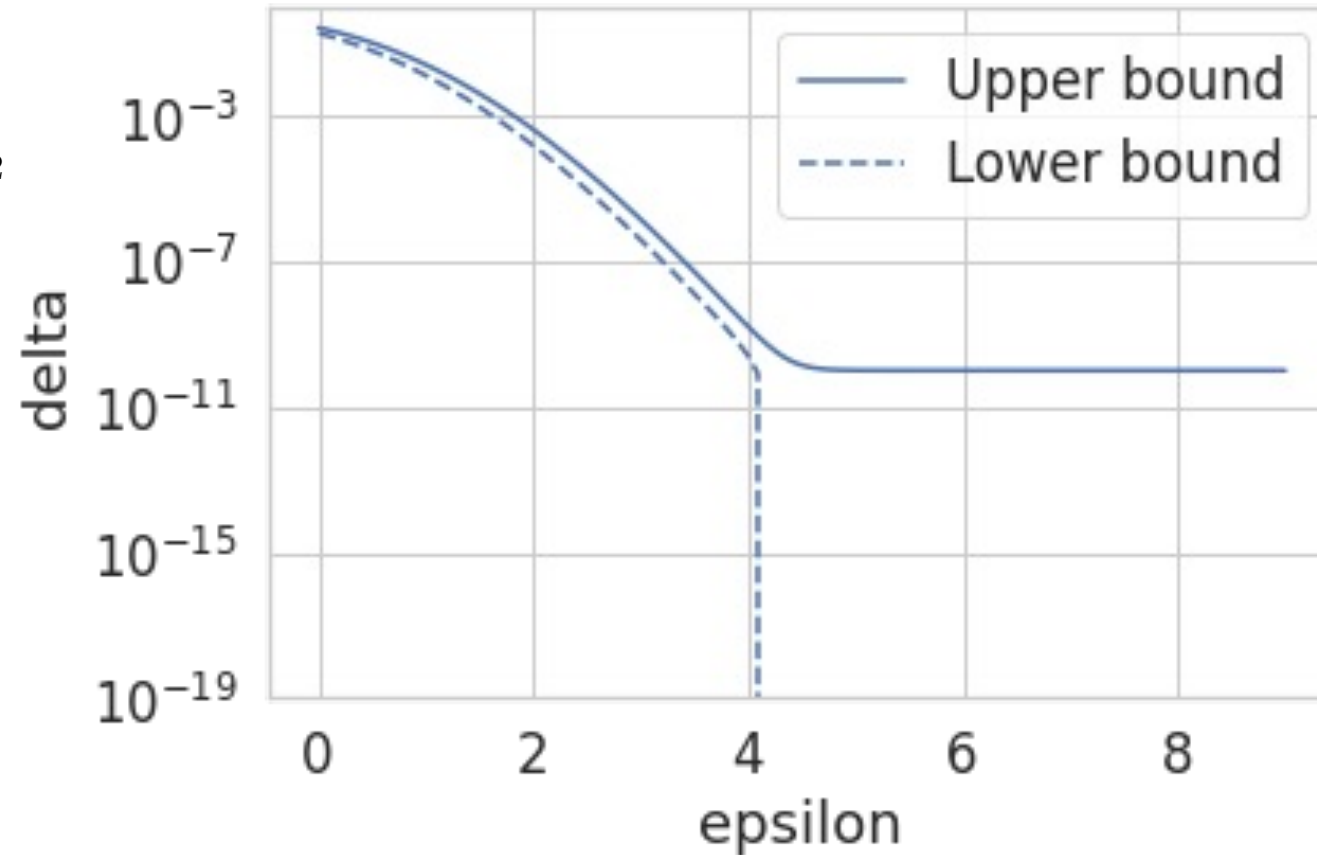
State-of-the-art composition

DP-SGD:

$$\sigma = 1$$

$$\text{subsampling} = 10^{-2}$$

$$n = 2000$$



Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

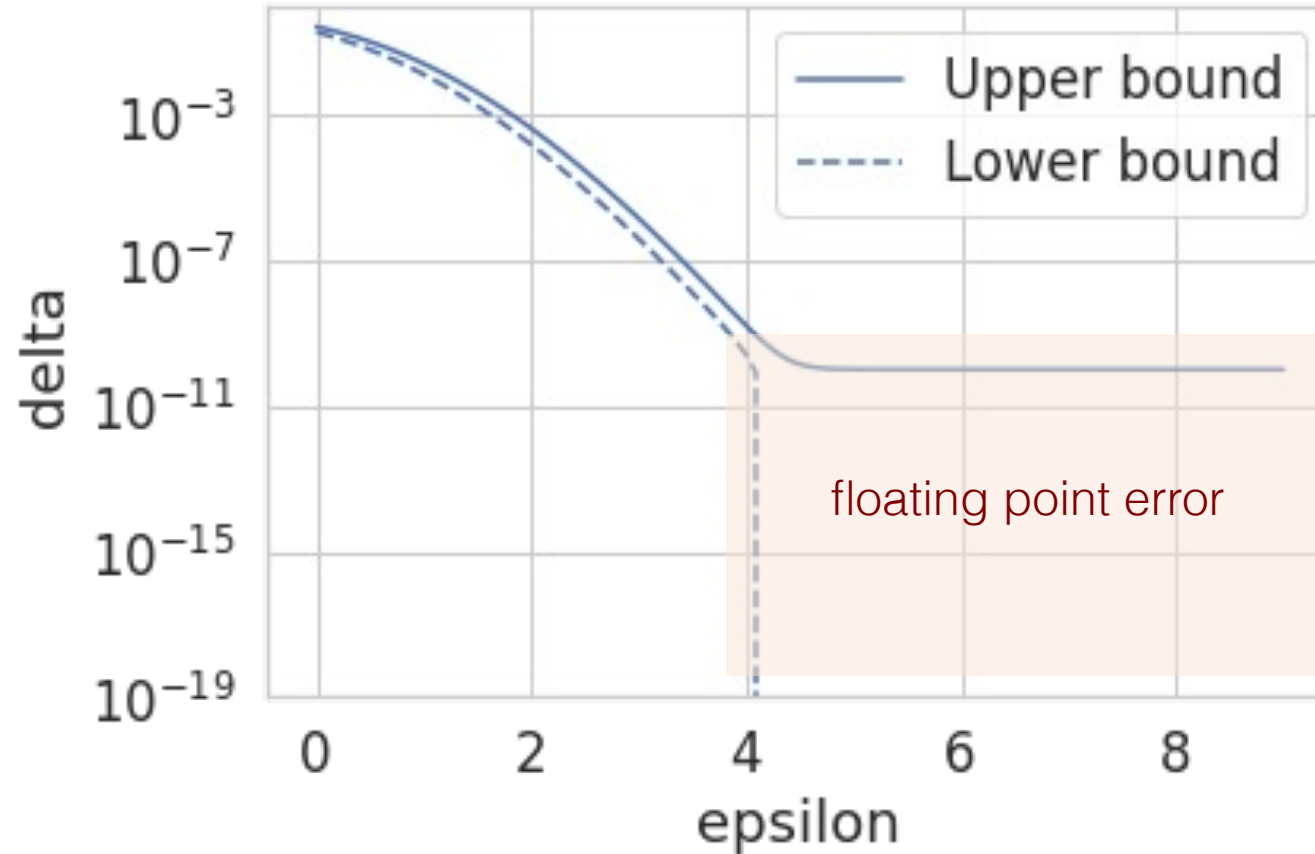
State-of-the-art composition

DP-SGD:

$$\sigma = 1$$

$$\text{subsampling} = 10^{-2}$$

$$n = 2000$$



Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

Today's talk: Develop DP numerical composition using saddle-point approximation:

- Runtime complexity independent of # composition
- Works for all epsilon and delta

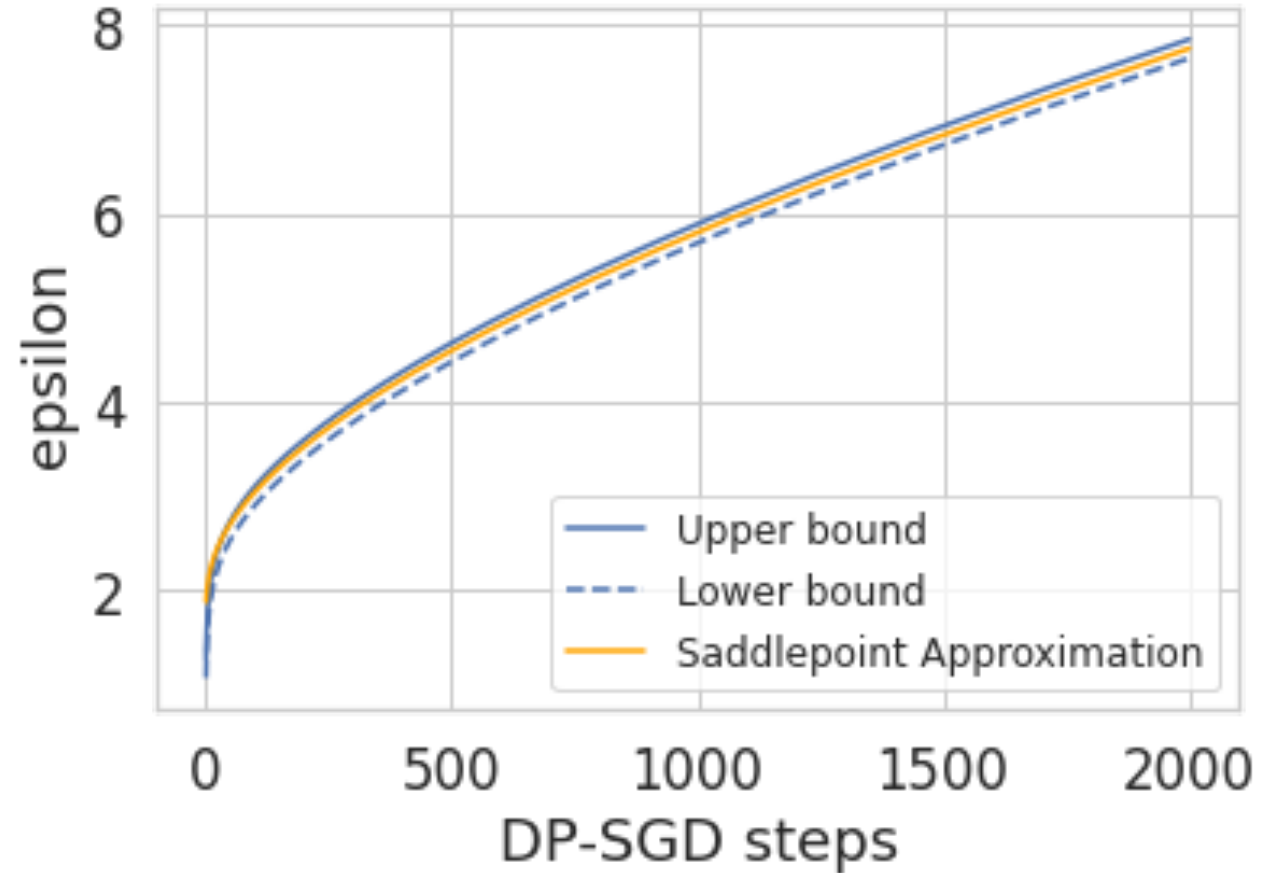
Saddle-point vs. state-of-the-art

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$



runtime complexity
 $O(1)$

Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

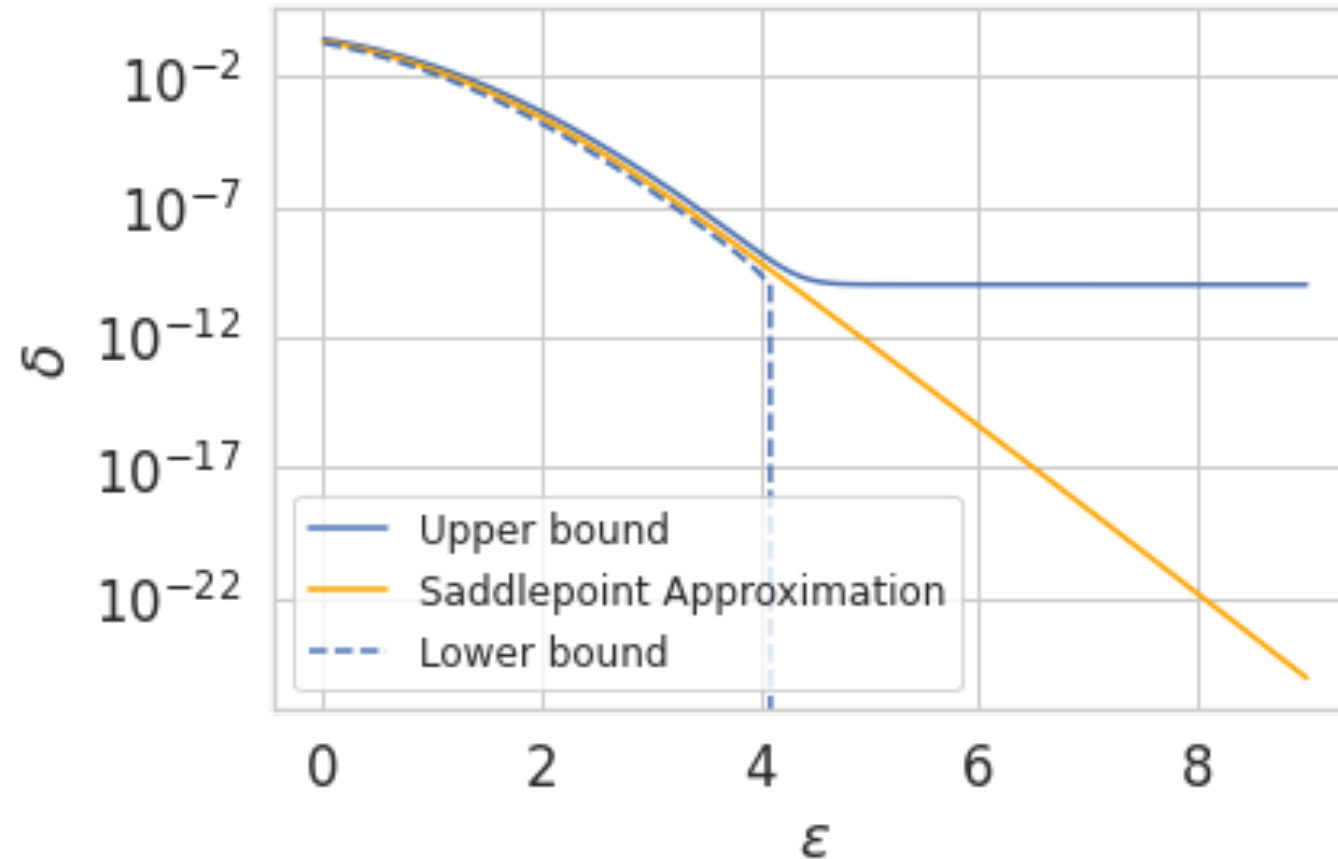
Saddle-point vs. state-of-the-art

DP-SGD:

$$\sigma = 1$$

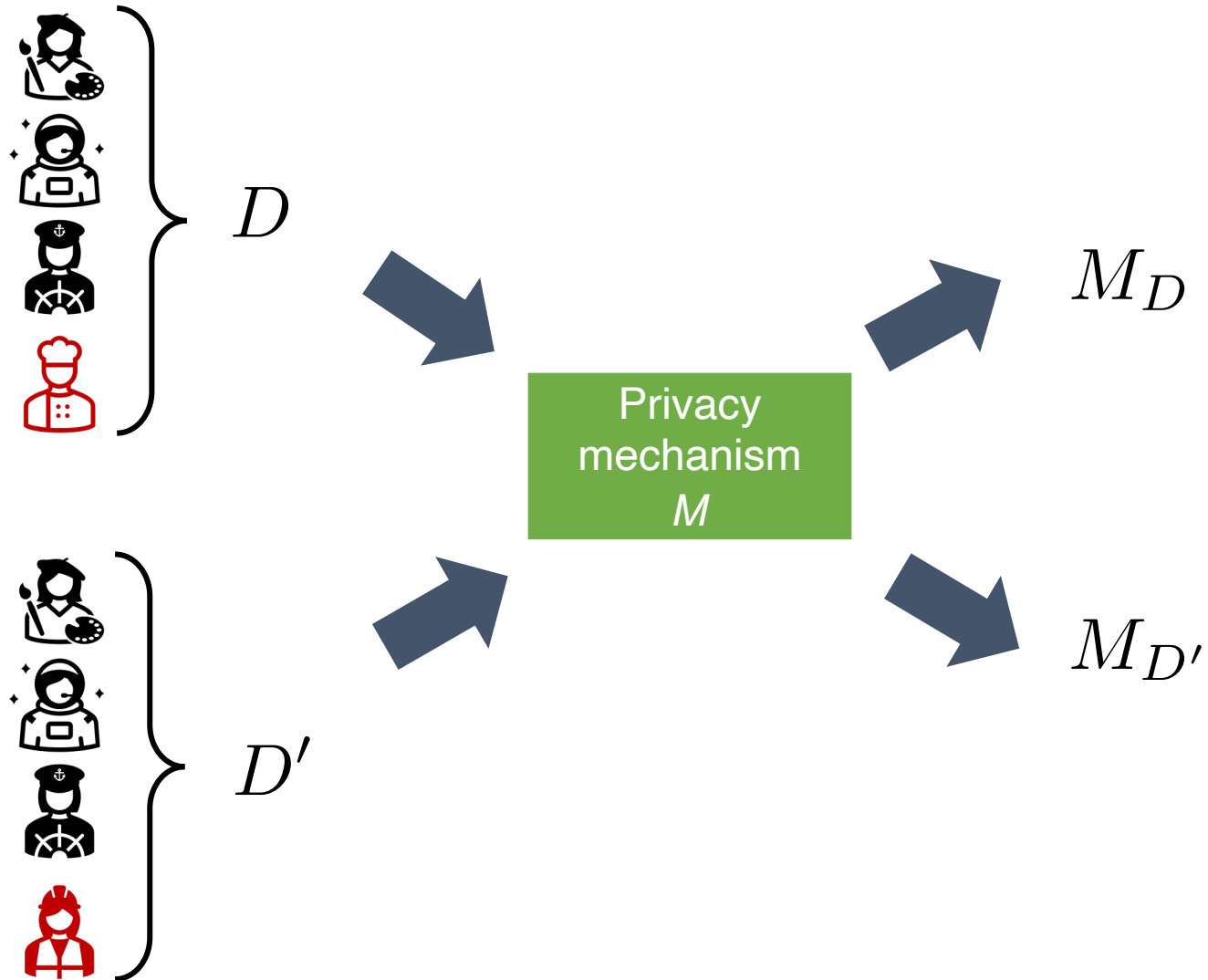
$$\text{subsampling} = 10^{-2}$$

$$n = 2000$$



Gopi, Lee, and Wutschitz, Numerical Composition of Differential Privacy, NeurIPS 2021

Differential privacy



M is (ϵ, δ) -DP if $\forall D \sim D'$

$$\sup_{\text{subset } A} [M_D(A) - e^\epsilon M_{D'}(A)] \leq \delta$$

Hockey-stick divergence
 $E_\epsilon(M_D || M_{D'})$

M is (ϵ, δ) -DP if $\forall D \sim D'$

$$E_\epsilon(M_D || M_{D'}) \leq \delta$$

Dominating distribution, PLRV

A pair (P, Q) is said to dominate M if

$$\sup_{D \sim D'} \mathbb{E}_\varepsilon(M_D \| M_{D'}) \leq \mathbb{E}_\varepsilon(P \| Q)$$

or tightly dominate M if equality is achieved for all ε .

$\delta(\varepsilon) \triangleq$ smallest δ such that M is (ε, δ) -DP

privacy curve

If (P, Q) tightly dominates M , then

$$\delta(\varepsilon) = \mathbb{E}_\varepsilon(P \| Q) = \mathbb{E} \left[\left(1 - e^{\varepsilon - L} \right)_+ \right]$$

where

$$L = \log \frac{dP}{dQ}(X) \quad \text{with } X \sim P$$

privacy loss random variable

Dominating distribution, PLRV

A pair (P, Q) is said to dominate M if

$$\sup_{D \sim D'} \mathbb{E}_\varepsilon(M_D \| M_{D'}) \leq \mathbb{E}_\varepsilon(P \| Q)$$

or tightly dominate M if equality is achieved for all ε .

$\delta(\varepsilon) \triangleq$ smallest δ such that M is (ε, δ) -DP

privacy curve

If (P, Q) ~~tightly~~ dominates M , then

$$\delta(\varepsilon) \stackrel{\leq}{\neq} \mathbb{E}_\varepsilon(P \| Q) = \mathbb{E} \left[\left(1 - e^{\varepsilon - L} \right)_+ \right]$$

where

$$L = \log \frac{dP}{dQ}(X) \quad \text{with } X \sim P$$

privacy loss random variable

DP-SGD

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

 Take a random sample L_t with sampling probability L/N

Compute gradient

 For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

DP-SGD

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

 Take a random sample L_t with sampling probability L/N

Compute gradient

 For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

Tightly dominating distributions for each iteration:

$$P = p\mathcal{N}(0, \sigma^2 C^2) + (1 - p)\mathcal{N}(C, \sigma^2 C^2) \quad Q = \mathcal{N}(0, \sigma^2 C^2)$$

DP-SGD

Algorithm 1 Differentially private SGD (Outline)

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

 Take a random sample L_t with sampling probability L/N

Compute gradient

 For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

Clip gradient

$\bar{\mathbf{g}}_t(x_i) \leftarrow \mathbf{g}_t(x_i) / \max(1, \frac{\|\mathbf{g}_t(x_i)\|_2}{C})$

Add noise

$\tilde{\mathbf{g}}_t \leftarrow \frac{1}{L} (\sum_i \bar{\mathbf{g}}_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 \mathbf{I}))$

Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

Tightly dominating distributions for each iteration:

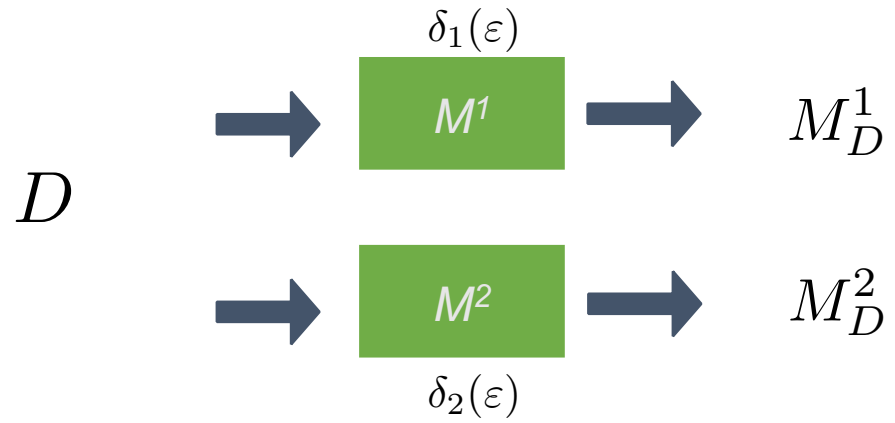
$$P = p\mathcal{N}(0, \sigma^2 C^2) + (1 - p)\mathcal{N}(C, \sigma^2 C^2) \quad Q = \mathcal{N}(0, \sigma^2 C^2)$$



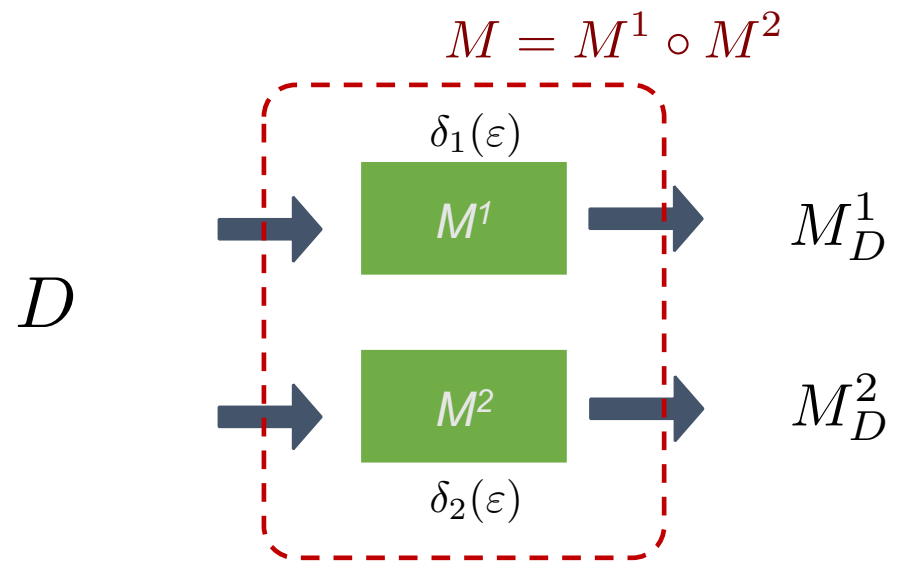
$$\delta(\epsilon) = \mathbb{E}_\epsilon(P \| Q) = \mathbb{E} \left[(1 - e^{\epsilon - L})_+ \right]$$

$$L = \log \left(1 - p + p \cdot e^{\frac{C(2X - C)}{2\sigma^2}} \right), \quad X \sim P$$

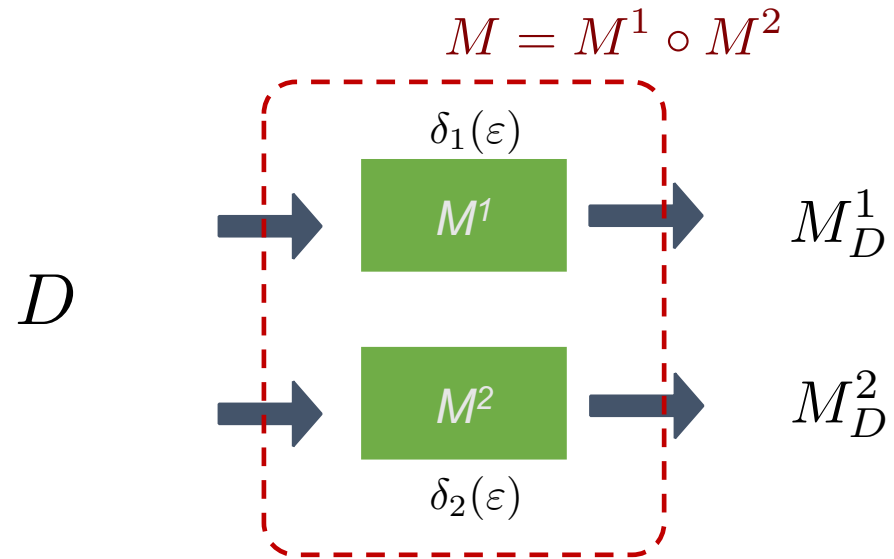
Composition of DP



Composition of DP

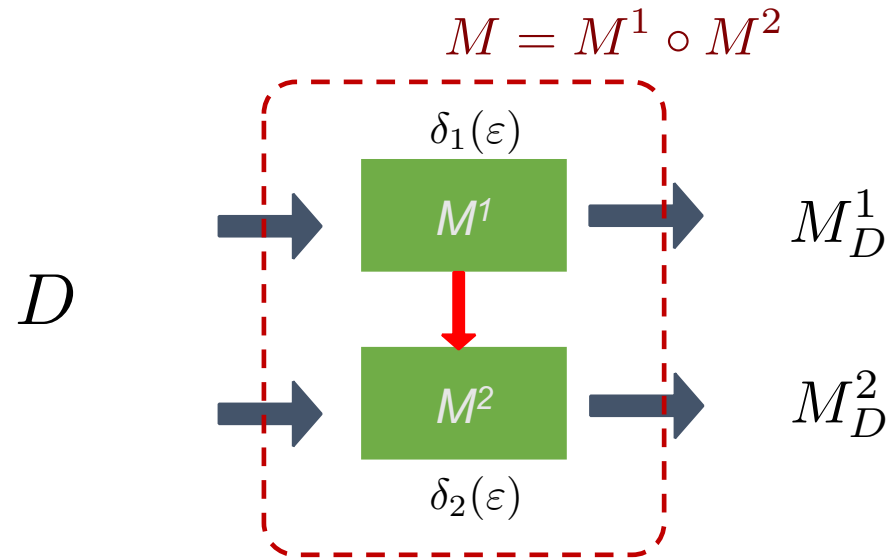


Composition of DP



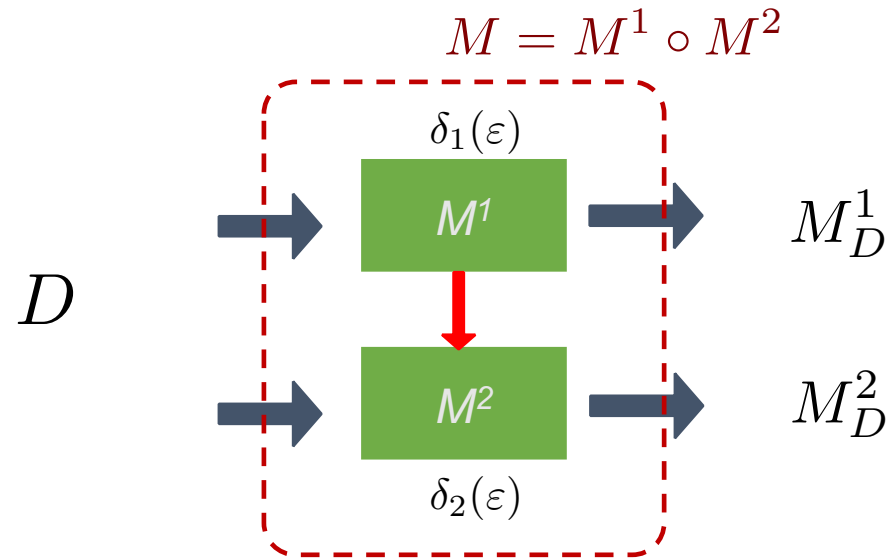
What is the privacy curve $\delta(\epsilon)$ of M ?

Composition of DP



What is the privacy curve $\delta(\epsilon)$ of M ?

Composition of DP

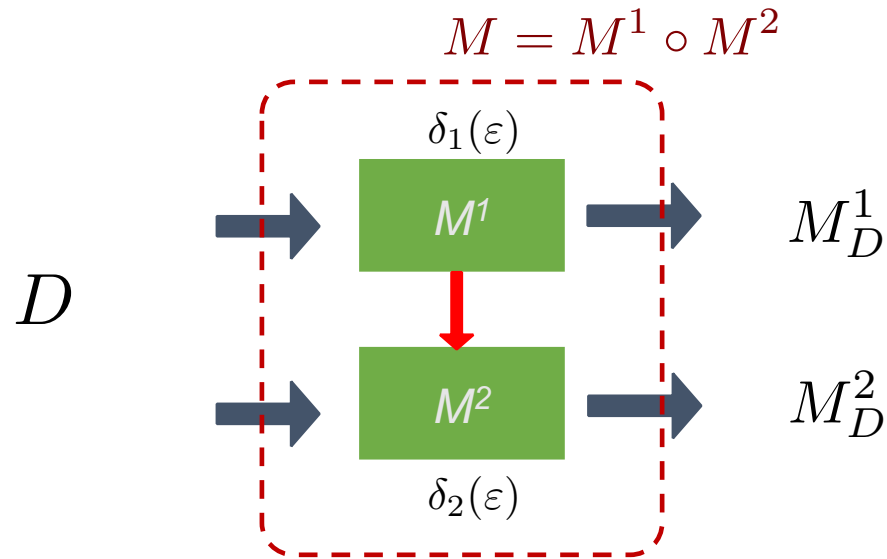


What is the privacy curve $\delta(\epsilon)$ of M ?

(P^1, Q^1) tightly dominates M^1

(P^2, Q^2) tightly dominates M^2

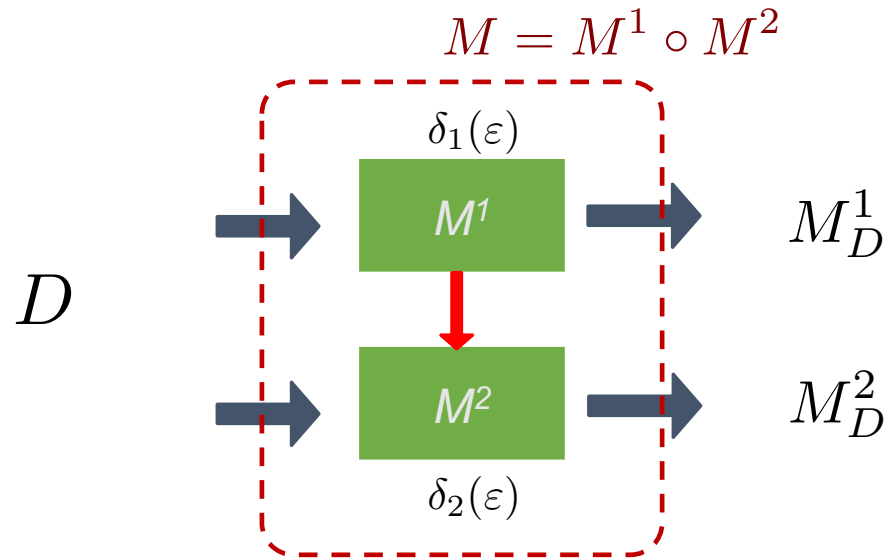
Composition of DP



What is the privacy curve $\delta(\epsilon)$ of M ?

(P^1, Q^1) tightly dominates M^1 }
 (P^2, Q^2) tightly dominates M^2 } $\implies (P^1 \times P^2, Q^1 \times Q^2)$ dominates M

Composition of DP

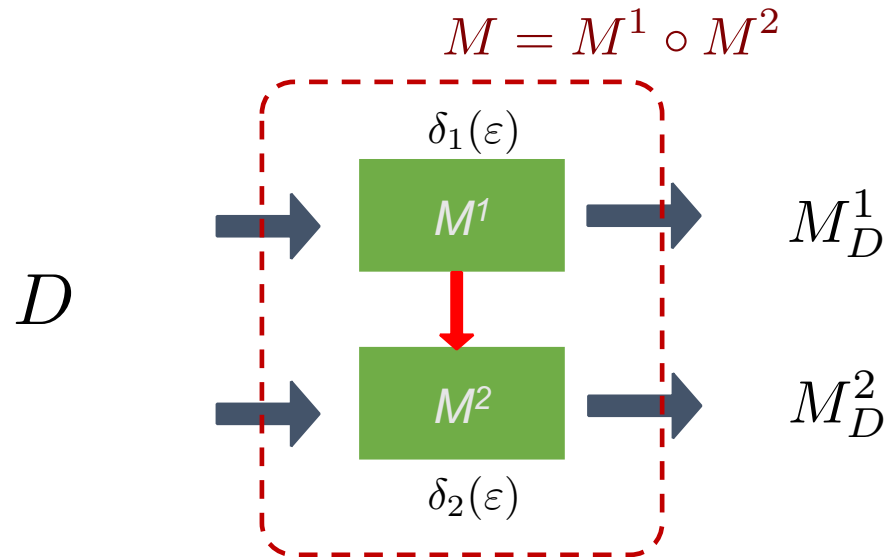


What is the privacy curve $\delta(\epsilon)$ of M ?

(P^1, Q^1) tightly dominates M^1 }
 (P^2, Q^2) tightly dominates M^2 } $\implies (P^1 \times P^2, Q^1 \times Q^2)$ dominates M

$$\delta(\epsilon) \leq \mathbb{E}_\epsilon(P^1 \times P^2 \| Q^1 \times Q^2)$$

Composition of DP



What is the privacy curve $\delta(\epsilon)$ of M ?

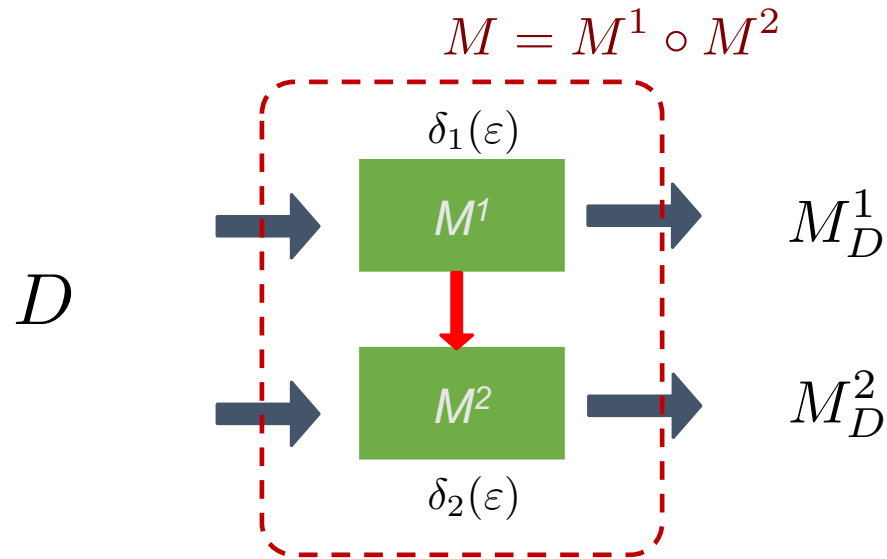
(P^1, Q^1) tightly dominates M^1
 (P^2, Q^2) tightly dominates M^2

$\implies (P^1 \times P^2, Q^1 \times Q^2)$ dominates M

\Downarrow

$$\delta(\epsilon) \leq \mathbf{E}_\epsilon(P^1 \times P^2 \| Q^1 \times Q^2) = \mathbb{E} \left[\left(1 - e^{\epsilon - (L_1 + L_2)}\right)_+ \right]$$

Composition of DP



What is the privacy curve $\delta(\epsilon)$ of M ?

$$\left. \begin{array}{l} (P^1, Q^1) \text{ tightly dominates } M^1 \\ (P^2, Q^2) \text{ tightly dominates } M^2 \end{array} \right\} \implies (P^1 \times P^2, Q^1 \times Q^2) \text{ dominates } M$$

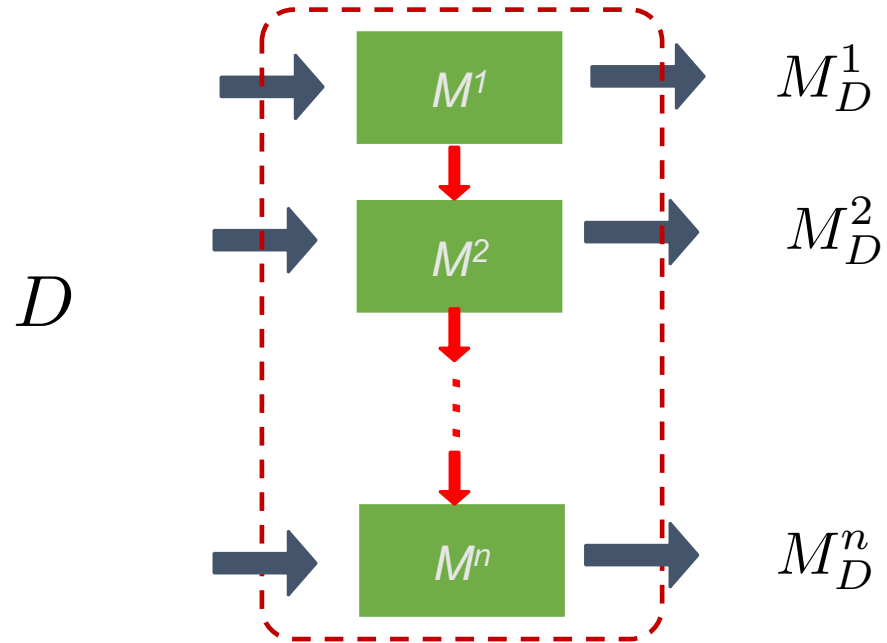
$$\Downarrow$$

$$\delta(\epsilon) \leq \mathbb{E}_\epsilon(P^1 \times P^2 \| Q^1 \times Q^2) = \mathbb{E} \left[\left(1 - e^{\epsilon - (L_1 + L_2)} \right)_+ \right]$$

\swarrow PLRV for M^1 \swarrow PLRV for M^2

Composition of DP

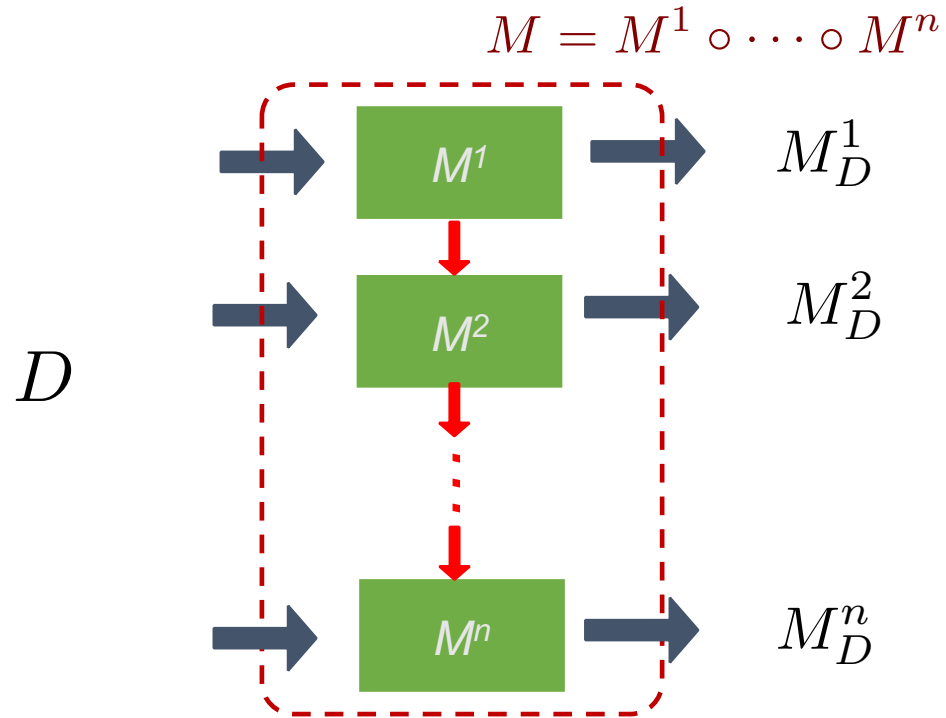
$$M = M^1 \circ \dots \circ M^n$$



What is the privacy curve $\delta(\varepsilon)$ of M ?

(P^i, Q^i) tightly dominates $M^i \implies (P^1 \times \dots \times P^n, Q^1 \times \dots \times Q^n)$ dominates M

Composition of DP



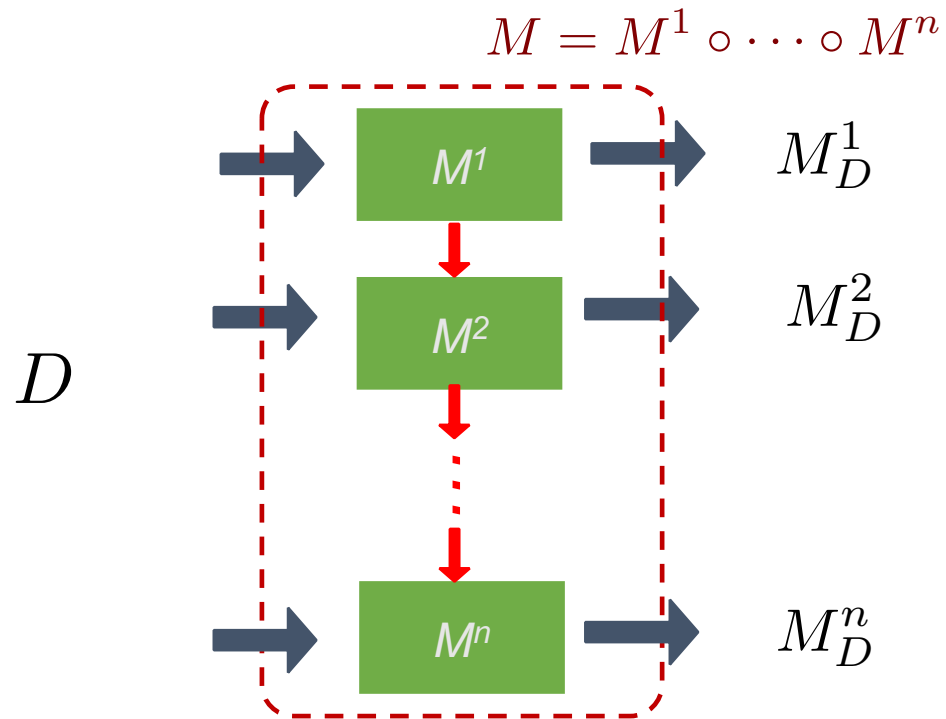
What is the privacy curve $\delta(\varepsilon)$ of M ?

(P^i, Q^i) tightly dominates $M^i \implies (P^1 \times \dots \times P^n, Q^1 \times \dots \times Q^n)$ dominates M

\Downarrow

$$\delta(\varepsilon) \leq \mathbf{E}_\varepsilon(P^1 \times \dots \times P^n \parallel Q^1 \times \dots \times Q^n)$$

Composition of DP



What is the privacy curve $\delta(\varepsilon)$ of M ?

(P^i, Q^i) tightly dominates $M^i \implies (P^1 \times \dots \times P^n, Q^1 \times \dots \times Q^n)$ dominates M

\Downarrow

$$\delta(\varepsilon) \leq \mathbf{E}_\varepsilon(P^1 \times \dots \times P^n \| Q^1 \times \dots \times Q^n) = \mathbb{E} \left[(1 - e^{\varepsilon - (L_1 + \dots + L_n)})_+ \right]$$

Composition Results

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - (L_1 + \dots + L_n)} \right)_+ \right]$$

Composition Results

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - (L_1 + \dots + L_n)} \right)_+ \right]$$

- Moments accountant: [Abadi et al'16], [Mironov'17]
- Central limit theorem: [Dong et al'19], [Sommer et al'19]
- Fast Fourier transform: [Koskela et al'20], [Koskela and Honkela'20], [Koskela et al'21], [Gopi et al'21], [Ghazi et al'22]
- Characteristic function: [Zhu et al'22]
- Piece-wise linearization of HS divergence: [Doroshenko et al'22]

Composition Results

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - (L_1 + \dots + L_n)} \right)_+ \right]$$

- Moments accountant: [Abadi et al'16], [Mironov'17]
- Central limit theorem: [Dong et al'19], [Sommer et al'19]
- Fast Fourier transform: [Koskela et al'20], [Koskela and Honkela'20], [Koskela et al'21], [Gopi et al'21], [Ghazi et al'22]
- Characteristic function: [Zhu et al'22]
- Piece-wise linearization of HS divergence: [Doroshenko et al'22]

Saddle-point accountant

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - (L_1 + \dots + L_n)} \right)_+ \right]$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right]$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell$$

Saddle-point accountant


$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} e^{t\ell} f_L(\ell) d\ell\end{aligned}$$

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} e^{t\ell} f_L(\ell) d\ell \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} \frac{e^{t\ell} f_L(\ell)}{\mathbb{E}[e^{tL}]} d\ell\end{aligned}$$

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L}\right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ f_L(l) dl \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} e^{tl} f_L(l) dl \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} \frac{e^{tl} f_L(l)}{\mathbb{E}[e^{tL}]} dl\end{aligned}$$

MGF of L 

Saddle-point accountant

$$\delta(\varepsilon) \leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ f_L(l) dl$$

$$= \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} e^{tl} f_L(l) dl$$

MGF of L

$$= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} \frac{e^{tl} f_L(l)}{\mathbb{E}[e^{tL}]} dl$$

pdf of \tilde{L} exponentially tilted of L

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ f_L(l) dl \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} e^{tl} f_L(l) dl \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} \frac{e^{tl} f_L(l)}{\mathbb{E}[e^{tL}]} dl \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} e^{-tl} (1 - e^{\varepsilon - l})_+ f_{\tilde{L}}(l) dl\end{aligned}$$

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} e^{t\ell} f_L(\ell) d\ell \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} \frac{e^{t\ell} f_L(\ell)}{\mathbb{E}[e^{tL}]} d\ell \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} e^{-t\ell} (1 - e^{\varepsilon - \ell})_+ f_{\tilde{L}}(\ell) d\ell\end{aligned}$$

Plancherel's Theorem:
$$\int_{\mathbb{R}} f(x)g(x)dx = \frac{1}{2\pi} \int_{\mathbb{R}} \mathcal{F}(f)(\omega)\overline{\mathcal{F}(g)(\omega)}d\omega$$

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ f_L(l) dl \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} e^{tl} f_L(l) dl \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - l})_+ e^{-tl} \frac{e^{tl} f_L(l)}{\mathbb{E}[e^{tL}]} dl \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} e^{-tl} (1 - e^{\varepsilon - l})_+ f_{\tilde{L}}(l) dl \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_{\varepsilon}(t+is)} ds\end{aligned}$$

Saddle-point accountant

$$\begin{aligned}\delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell \\ &= \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} e^{t\ell} f_L(\ell) d\ell \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} \frac{e^{t\ell} f_L(\ell)}{\mathbb{E}[e^{tL}]} d\ell \\ &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} e^{-t\ell} (1 - e^{\varepsilon - \ell})_+ f_{\tilde{L}}(\ell) d\ell \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds\end{aligned}$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Saddle-point accountant

$$\begin{aligned}
 \delta(\varepsilon) &\leq \mathbb{E} \left[\left(1 - e^{\varepsilon - \overbrace{(L_1 + \dots + L_n)}^L} \right)_+ \right] = \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ f_L(\ell) d\ell \\
 &= \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} e^{t\ell} f_L(\ell) d\ell \\
 &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} (1 - e^{\varepsilon - \ell})_+ e^{-t\ell} \frac{e^{t\ell} f_L(\ell)}{\mathbb{E}[e^{tL}]} d\ell \\
 &= \mathbb{E}[e^{tL}] \int_{\mathbb{R}} e^{-t\ell} (1 - e^{\varepsilon - \ell})_+ f_{\tilde{L}}(\ell) d\ell \\
 &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds
 \end{aligned}$$

cumulant generating function (CGF) of L

$$K_L(z) = \log \mathbb{E}[e^{zL}]$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε :

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*) \implies$$

Along the real line, F_ε is minimized at $z = t_*$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*) \implies$$

Along the real line, F_ε is minimized at $z = t_*$

Parallel to imaginary axis, F_ε is maximized at $z = t_*$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

Second approximation:

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation: Vanilla saddle-point approximation

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation: Vanilla saddle-point approximation

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$



$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds \approx \frac{e^{F_\varepsilon(t_*)}}{\sqrt{2\pi F''_\varepsilon(t_*)}}$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$



$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds \approx \frac{e^{F_\varepsilon(t_*)}}{\sqrt{2\pi F''_\varepsilon(t_*)}} = \frac{e^{K_L(t_*) - \varepsilon t_*}}{\sqrt{2\pi [t_*^2(1+t_*)^2 K''_L(t_*) + t_*^2 + (t_* + 1)^2]}}$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$



$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds \approx e^{K_L(t_*) - \varepsilon t_*} \mathbb{E} \left[e^{t_*(\varepsilon - Z)} (1 - e^{\varepsilon - Z})_+ \right]$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$



$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds \approx e^{K_L(t_*) - \varepsilon t_*} \mathbb{E} \left[e^{t_*(\varepsilon - Z)} (1 - e^{\varepsilon - Z})_+ \right]$$

Saddle-point accountant

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$

Take t to be **saddle-point** of F_ε : Unique t_* satisfying $F'_\varepsilon(t_*) = 0$, i.e.,

$$K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_* + 1}$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F''_\varepsilon(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K''_L(t_*)$$



$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds \approx e^{K_L(t_*) - \varepsilon t_*} \mathbb{E} \left[e^{t_*(\varepsilon - Z)} (1 - e^{\varepsilon - Z})_+ \right]$$

$$Z \sim \mathcal{N}(K'_L(t_*), K''_L(t_*))$$

Saddle-point accountant: Algorithm

Input: Tightly dominating pairs $\{(P_i, Q_i)\}_{i=1}^n$ for mechanisms $\{M_i\}_{i=1}^n$ and ε

- Compute (numerically estimate) $K_{L_i}(t) = \log \mathbb{E}[e^{tL_i}]$
- $K_L(t) = \sum_i K_{L_i}(t)$ (since L_i are independent)
- Find saddle-point t_* by solving $K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_*+1}$

Outputs:

Saddle-point accountant: Algorithm

Input: Tightly dominating pairs $\{(P_i, Q_i)\}_{i=1}^n$ for mechanisms $\{M_i\}_{i=1}^n$ and ε

- Compute (numerically estimate) $K_{L_i}(t) = \log \mathbb{E}[e^{tL_i}]$
- $K_L(t) = \sum_i K_{L_i}(t)$ (since L_i are independent)
- Find saddle-point t_* by solving $K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_*+1}$

Outputs:

$$\hat{\delta}_1(\varepsilon) = \frac{e^{K_L(t_*) - \varepsilon t_*}}{\sqrt{2\pi [t_*^2(1+t_*)^2 K''_L(t_*) + t_*^2 + (t_*+1)^2]}}$$

$$\hat{\delta}_2(\varepsilon) = e^{K_L(t_*) - \varepsilon t_*} \mathbb{E} \left[e^{t_*(\varepsilon - Z)} (1 - e^{\varepsilon - Z})_+ \right]$$

Saddle-point accountant: Algorithm

Input: Tightly dominating pairs $\{(P_i, Q_i)\}_{i=1}^n$ for mechanisms $\{M_i\}_{i=1}^n$ and ε

- Compute (numerically estimate) $K_{L_i}(t) = \log \mathbb{E}[e^{tL_i}]$
- $K_L(t) = \sum_i K_{L_i}(t)$ (since L_i are independent)
- Find saddle-point t_* by solving $K'_L(t_*) = \varepsilon + \frac{1}{t_*} + \frac{1}{t_*+1}$

Outputs:

$$\hat{\delta}_1(\varepsilon) = \frac{e^{K_L(t_*) - \varepsilon t_*}}{\sqrt{2\pi [t_*^2(1+t_*)^2 K''_L(t_*) + t_*^2 + (t_*+1)^2]}}$$

$$\hat{\delta}_2(\varepsilon) = e^{K_L(t_*) - \varepsilon t_*} \mathbb{E} \left[e^{t_*(\varepsilon - Z)} (1 - e^{\varepsilon - Z})_+ \right]$$

Both are “corrected” versions of moments accountant

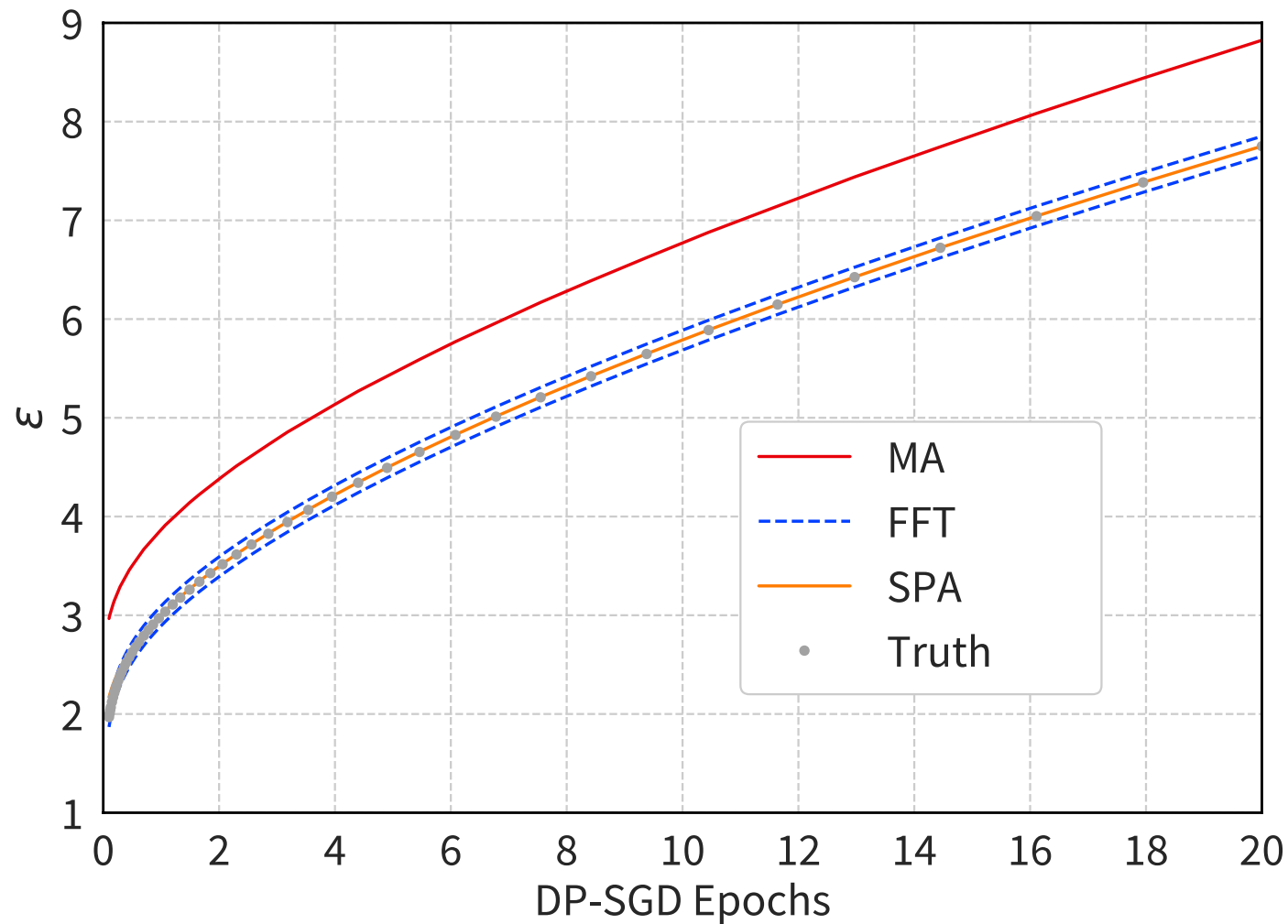
Numerical experiments

DP-SGD:

$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$



Numerical experiments

DP-SGD:

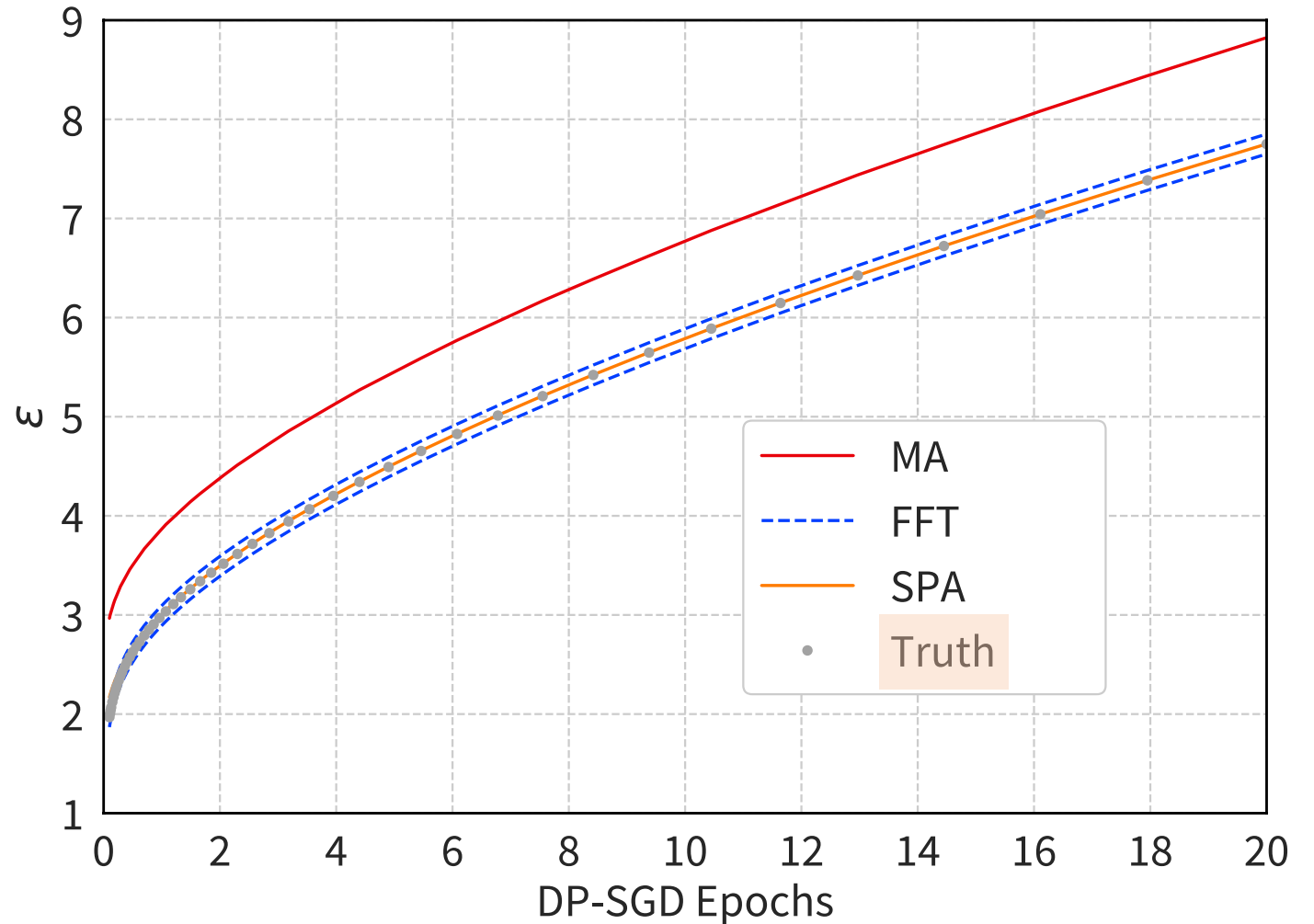
$$\sigma = 0.65$$

$$\text{subsampling} = 10^{-2}$$

$$\delta = 10^{-5}$$

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

$$F_\varepsilon(z) \triangleq K_L(z) - z\varepsilon - \log z - \log(1+z)$$



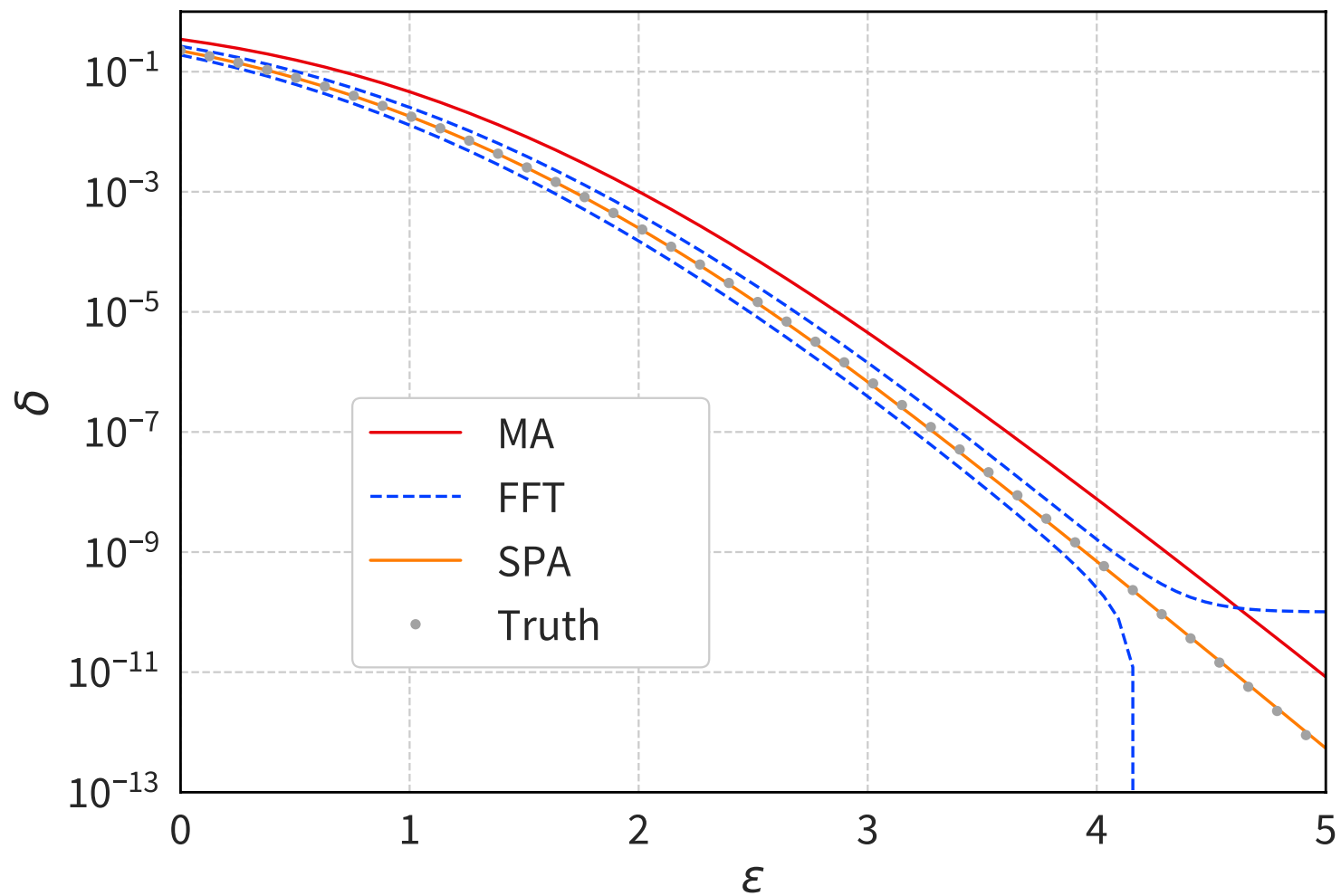
Numerical experiments

DP-SGD:

$$\sigma = 1$$

$$\text{subsampling} = 10^{-2}$$

$$n = 2000$$



Error analysis

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F_\varepsilon''(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K_L''(t_*)$$

Error analysis

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F_\varepsilon''(t_*)$$

vanilla saddle-point approximation

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K_L''(t_*)$$

Edgeworth expansion

Error analysis

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F_\varepsilon''(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K_L''(t_*)$$

For any $\varepsilon \geq 0$

$$\left| \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds - \hat{\delta}_2(\varepsilon) \right| \leq e^{K_L(t_*) - \varepsilon t_*} \left(\frac{t_*}{1 + t_*} \right)^{t_*} \frac{P_{t_*}}{K_L''(t_*)^{3/2}}$$

Error analysis

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

First approximation:


$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F_\varepsilon''(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K_L''(t_*)$$

For any $\varepsilon \geq 0$

$$\left| \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds - \hat{\delta}_2(\varepsilon) \right| \leq e^{K_L(t_*) - \varepsilon t_*} \left(\frac{t_*}{1 + t_*} \right)^{t_*} \frac{P_{t_*}}{K_L''(t_*)^{3/2}}$$

$$\sum_{i=1}^n \mathbb{E}[|\tilde{L}_i - \mathbb{E}[\tilde{L}_i]|^3]$$


Error analysis

$$\delta(\varepsilon) \leq \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds$$

First approximation:

$$F_\varepsilon(z) \approx F_\varepsilon(t_*) + \frac{1}{2}(z - t_*)^2 F_\varepsilon''(t_*)$$

Second approximation:

$$K_L(z) \approx K_L(t_*) + \frac{1}{2}(z - t_*)^2 K_L''(t_*)$$

For any $\varepsilon = \mathbb{E}[L] + b \cdot \text{var}(L)$, we have

$$\left| \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{F_\varepsilon(t+is)} ds - \hat{\delta}_2(\varepsilon) \right| \leq \frac{C}{\sqrt{n}}$$

Summary

- Accountant algorithm comparable with state-of-the-art
- Runtime complexity independent of # composition
- Works for all epsilon and delta

Summary

- Accountant algorithm comparable with state-of-the-art
- Runtime complexity independent of # composition
- Works for all epsilon and delta

Ongoing work

- Tightly dominating distribution of composed mechanisms
- Tighter error analysis
- Efficient estimator for CGF

Summary

- Accountant algorithm comparable with state-of-the-art
- Runtime complexity independent of # composition
- Works for all epsilon and delta

Ongoing work

- Tightly dominating distribution of composed mechanisms
- Tighter error analysis
- Efficient estimator for CGF

Pre-print available here!

