

DPZero: Private Fine-Tuning of Language Models without Backpropagation



Liang Zhang

(ETH Zurich & Max Planck Institute)



Bingcong Li

(ETH Zurich)



Kiran Thekumparampil

(Amazon)



Sewoong Oh

(University of Washington)

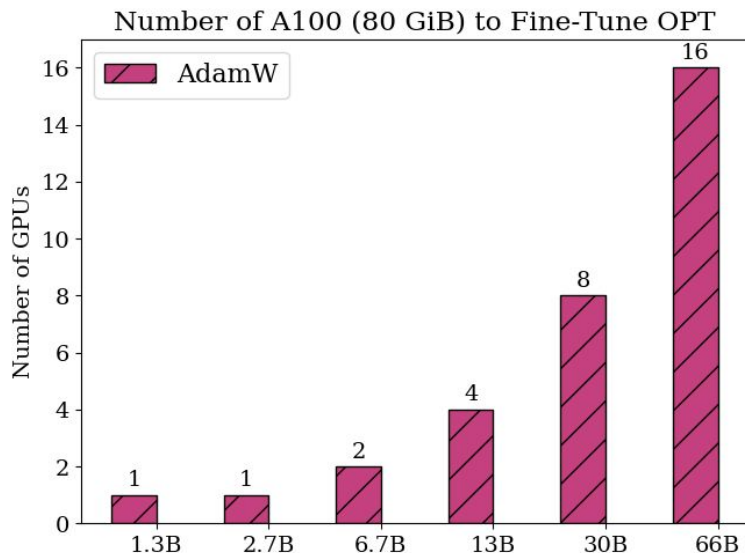


Niao He

(ETH Zurich)

Bottleneck in Fine-Tuning LLMs

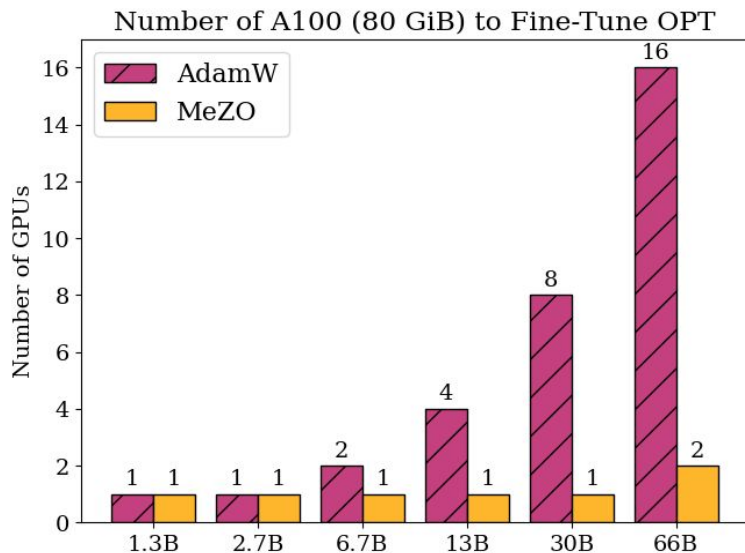
- **Memory** becomes a bottleneck when fine-tuning billion-sized LLMs



- Backpropagation heavy in memory

Bottleneck in Fine-Tuning LLMs

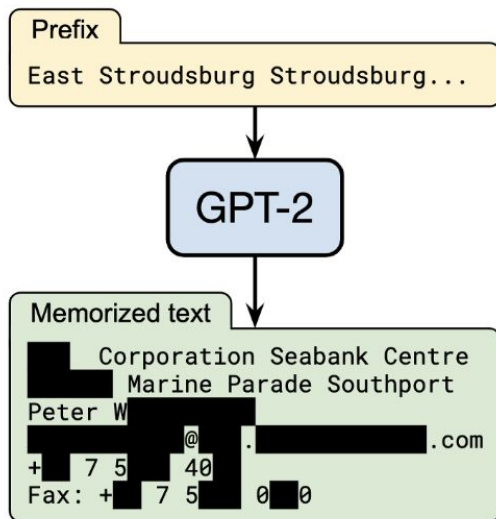
- **Memory** becomes a bottleneck when fine-tuning billion-sized LLMs



- Backpropagation heavy in memory
- MeZO [1]: **zeroth-order** methods with **only forward passes**

Bottleneck in Fine-Tuning LLMs

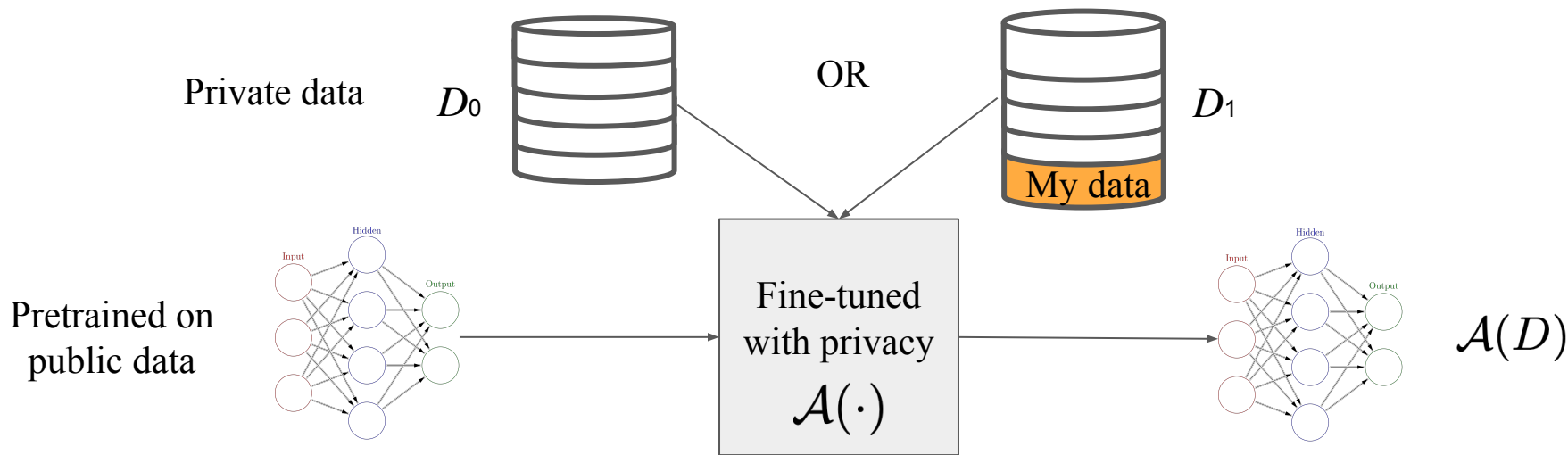
- High quality data are often **private**



- LLMs memorize and reveal sensitive data !!!
- Fine-tuning LLMs with **differential privacy (DP)**

Bottleneck in Fine-Tuning LLMs

- High quality data are often **private** → fine-tuning with (ϵ, δ) -Differential Privacy



$$\mathbb{P}(\mathcal{A}(D_1) \in R) \leq e^\epsilon \mathbb{P}(\mathcal{A}(D_0) \in R) + \delta$$

Private and Zeroth-Order Optimization for LLMs Fine-Tuning

Gap Between Theory and Practice

$$\text{minimize}_x F_S(x) := \frac{1}{n} \sum_{i=1}^n f(x; \xi_i)$$

Theory

Convergence rates of both

- 👉 **private** first-order optimization
 - 👉 non-private **zeroth-order** optimization
- depend on the **dimension**

Practice

They perform well on LLMs fine-tuning with
dimension scaling to **billions**

Non-Private Zeroth-Order Optimization

$$x_{t+1} \leftarrow x_t - \alpha \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda}}_{\text{0-th order gradient estimate}} u_t \right)$$

Non-Private Zeroth-Order Optimization

$$x_{t+1} \leftarrow x_t - \alpha \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda}}_{\text{0-th order gradient estimate}} u_t \right)$$

Theory

- **Unbiased** when $\lambda \rightarrow 0$ & $E[uu^T] = I_d$

$$\|\nabla F_S(x)\|^2 \lesssim \frac{d}{T}$$

Non-Private Zeroth-Order Optimization

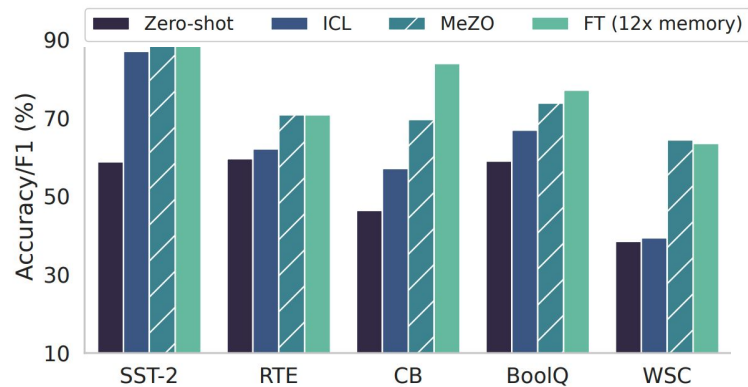
$$x_{t+1} \leftarrow x_t - \alpha \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} u_t \right)}_{\text{0-th order gradient estimate}}$$

Theory

- **Unbiased** when $\lambda \rightarrow 0$ & $E[uu^T] = I_d$

$$\|\nabla F_S(x)\|^2 \lesssim \frac{d}{T}$$

Practice



OPT-13B

Private First-Order Optimization: DP-GD

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C(\nabla f(x_t; \xi_i)) + \frac{C}{n} z_t \right)$$

Private First-Order Optimization: DP-GD

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C(\nabla f(x_t; \xi_i)) + \frac{C}{n} z_t \right)$$

Theory

- (ϵ, δ) -differential privacy with

$$z_t \sim \mathcal{N}(0, (4\sqrt{2T \log(1.25/\delta)}/\epsilon)^2 \mathbf{I}_{d \times d})$$

$$\|\nabla F_S(x)\|^2 \lesssim \frac{\sqrt{d \log(1/\delta)}}{n \epsilon}$$

Private First-Order Optimization: DP-GD

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C(\nabla f(x_t; \xi_i)) + \frac{C}{n} z_t \right)$$

Theory

- (ϵ, δ) -differential privacy with

$$z_t \sim \mathcal{N}(0, (4\sqrt{2T \log(1.25/\delta)}/\epsilon)^2 \mathbf{I}_{d \times d})$$

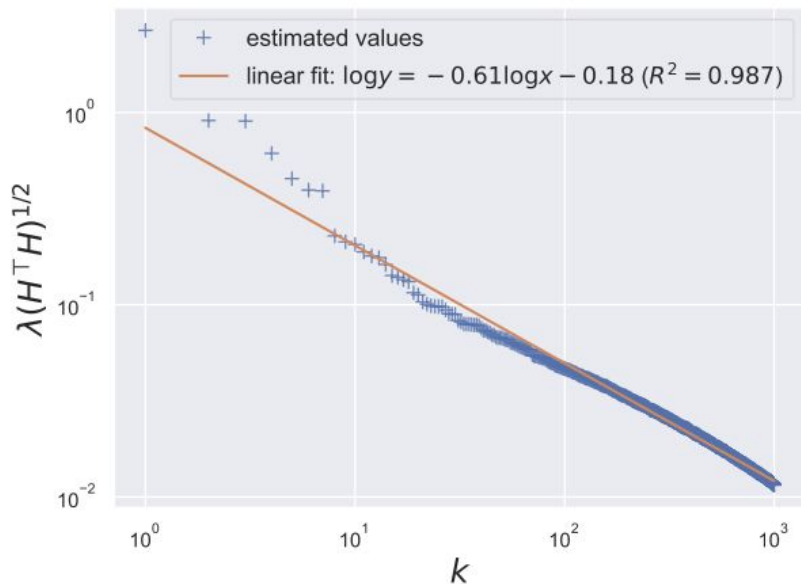
$$\|\nabla F_S(x)\|^2 \lesssim \frac{\sqrt{d \log(1/\delta)}}{n \epsilon}$$

Practice

Model	BLEU	BLEU (DP)	Drop
GPT2-(345M)	47.1	42.0	5.1
GPT2-(774M)	47.5	43.1	4.4
GPT2-(1.5B)	48.1	43.8	4.3

$(\epsilon = 6.8, \delta = 1\text{e-}5)$ on DART

Low-Dimensional Structure



Effective Rank

$$-H \preceq \nabla^2 F_S(x) \preceq H$$

$$\text{Tr}(H) \leq r \|H\|_2$$

- Not necessarily low rank
- Recover smoothness when $r = d$

Dimension-Independent
Private and Zeroth-Order Optimization
under Low Effective Rank Structure

First Attempt: DPGD-0th

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\underbrace{\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda}}_{\text{0-th order gradient estimate}} u_t \right) + \frac{C}{n} z_t \right)$$

First Attempt: DPGD-0th

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\underbrace{\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda}}_{\text{0-th order gradient estimate}} u_t \right) + \frac{C}{n} z_t \right)$$

$f(x; \xi_i)$ is L -Lipschitz and ℓ -smooth

- Without effective rank

$$\mathbb{E} [\|\nabla F_S(x_\tau)\|^2] \lesssim ((F_S(x_0) - F_S^*)\ell + L^2) \frac{d\sqrt{d \log(1/\delta)}}{n\varepsilon}$$

- With effective rank

$$\mathbb{E} [\|\nabla F_S(x_\tau)\|^2] \lesssim ((F_S(x_0) - F_S^*)\ell + L^2) \frac{d\sqrt{r \log(1/\delta)}}{n\varepsilon}$$

First Attempt: DPGD-0th

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\underbrace{\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda}}_{\text{0-th order gradient estimate}} u_t \right) + \frac{C}{n} z_t \right)$$

$f(x; \xi_i)$ is L -Lipschitz and ℓ -smooth

$$C \sim Ld$$

- Without effective rank

$$\mathbb{E} [\|\nabla F_S(x_\tau)\|^2] \lesssim ((F_S(x_0) - F_S^*)\ell + L^2) \frac{d\sqrt{d \log(1/\delta)}}{n\varepsilon}$$

- With effective rank

$$\mathbb{E} [\|\nabla F_S(x_\tau)\|^2] \lesssim ((F_S(x_0) - F_S^*)\ell + L^2) \frac{d\sqrt{r \log(1/\delta)}}{n\varepsilon}$$

DPZero: Dimension-Independent

- No need to add noise to the update direction

$$x_{t+1} \leftarrow x_t - \alpha \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right)}_{\text{(approx.) directional derivative}} + \underbrace{\frac{C}{n} z_t}_{\text{scalar noise}} \right) u_t$$

DPZero: Dimension-Independent

- No need to add noise to the update direction

$$x_{t+1} \leftarrow x_t - \alpha \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right)}_{\text{(approx.) directional derivative}} + \underbrace{\frac{C}{n} z_t}_{\text{scalar noise}} \right) u_t$$

(approx.) directional derivative

scalar noise

$$\simeq \langle \nabla f(x_t; \xi_i), u_t \rangle \simeq \begin{cases} \sqrt{d} L & \text{worst-case} \\ L & \text{w.h.p} \end{cases}$$

$$C \sim L$$

DPZero: Dimension-Independent

- No need to add noise to the update direction

$$x_{t+1} \leftarrow x_t - \alpha \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right)}_{\text{(approx.) directional derivative}} + \underbrace{\frac{C}{n} z_t}_{\text{scalar noise}} \right) u_t$$

$$\simeq \langle \nabla f(x_t; \xi_i), u_t \rangle \simeq \begin{cases} \sqrt{d} L & \text{worst-case} \\ L & \text{w.h.p} \end{cases}$$

- Under effective rank structure

$$C \sim L$$

$$\mathbb{E}[\|\nabla F_S(x_\tau)\|^2] \lesssim ((F_S(x_0) - F_S^*)\ell + L^2) \frac{\sqrt{r \log(1/\delta)}}{n\varepsilon}$$

DPZero

Algorithm 3 DPZERO

Input: Dataset $S = \{\xi_1, \dots, \xi_n\}$, initialization $x_0 \in \mathbb{R}^d$, number of iterations T , stepsize $\alpha > 0$, smoothing parameter $\lambda > 0$, clipping threshold $C > 0$, privacy parameters $\varepsilon > 0, \delta \in (0, 1)$.

1: **for** $t = 0, 1, \dots, T - 1$ **do**

2: Sample u_t uniformly at random from the Euclidean sphere $\sqrt{d} \mathbb{S}^{d-1}$.

3: Sample $z_t \sim \mathcal{N}(0, \sigma^2)$ with variance $\sigma = 4\sqrt{2T \log(e + (\varepsilon/\delta))}/\varepsilon$, and

$$x_{t+1} \leftarrow x_t - \alpha \left(\frac{1}{n} \sum_{i=1}^n \text{clip}_C \left(\frac{f(x_t + \lambda u_t; \xi_i) - f(x_t - \lambda u_t; \xi_i)}{2\lambda} \right) + \frac{C}{n} z_t \right) u_t.$$

Output: x_τ for τ sampled uniformly at random from $\{0, 1, \dots, T - 1\}$.

Summary of Results

- **Effective rank:** $-H \preceq \nabla^2 F_S(x) \preceq H$

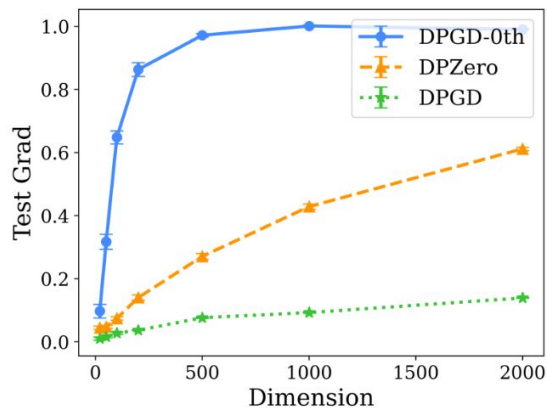
$$\text{Tr}(H) \leq r \|H\|_2$$

	without effective rank	with effective rank r
DPGD-0th	$d\sqrt{d}$	$d\sqrt{r}$
DPZERO	$(\log d)\sqrt{d}$	$(\log d)\sqrt{r}$
DP-GD	\sqrt{d}	\sqrt{r}

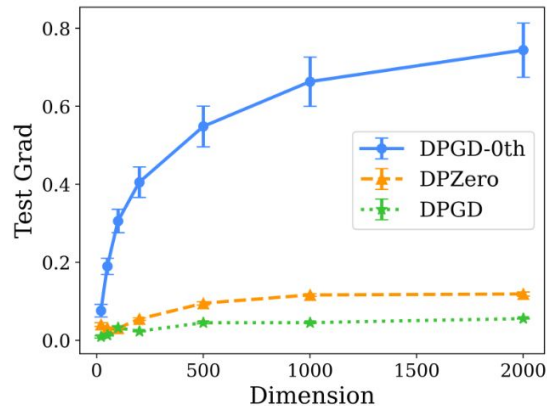
Experiments on Synthetic Examples

- $(\varepsilon = 2, \delta = 10^{-6})$ - DP on a quadratic

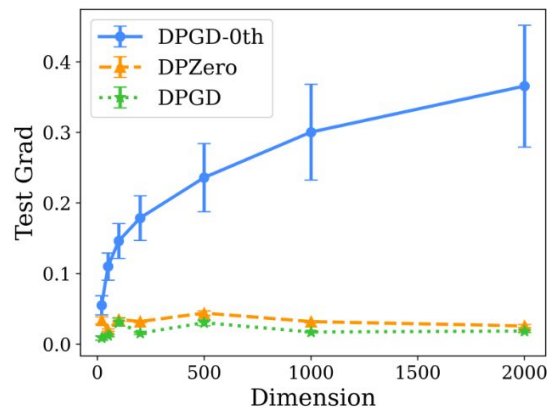
$$\min_{x \in \mathbb{R}^d} F_S(x) = \frac{1}{2n} \sum_{i=1}^n (x - x_i)^\top A (x - x_i).$$



(a) $\text{Tr}(A) = \mathcal{O}(d)$.

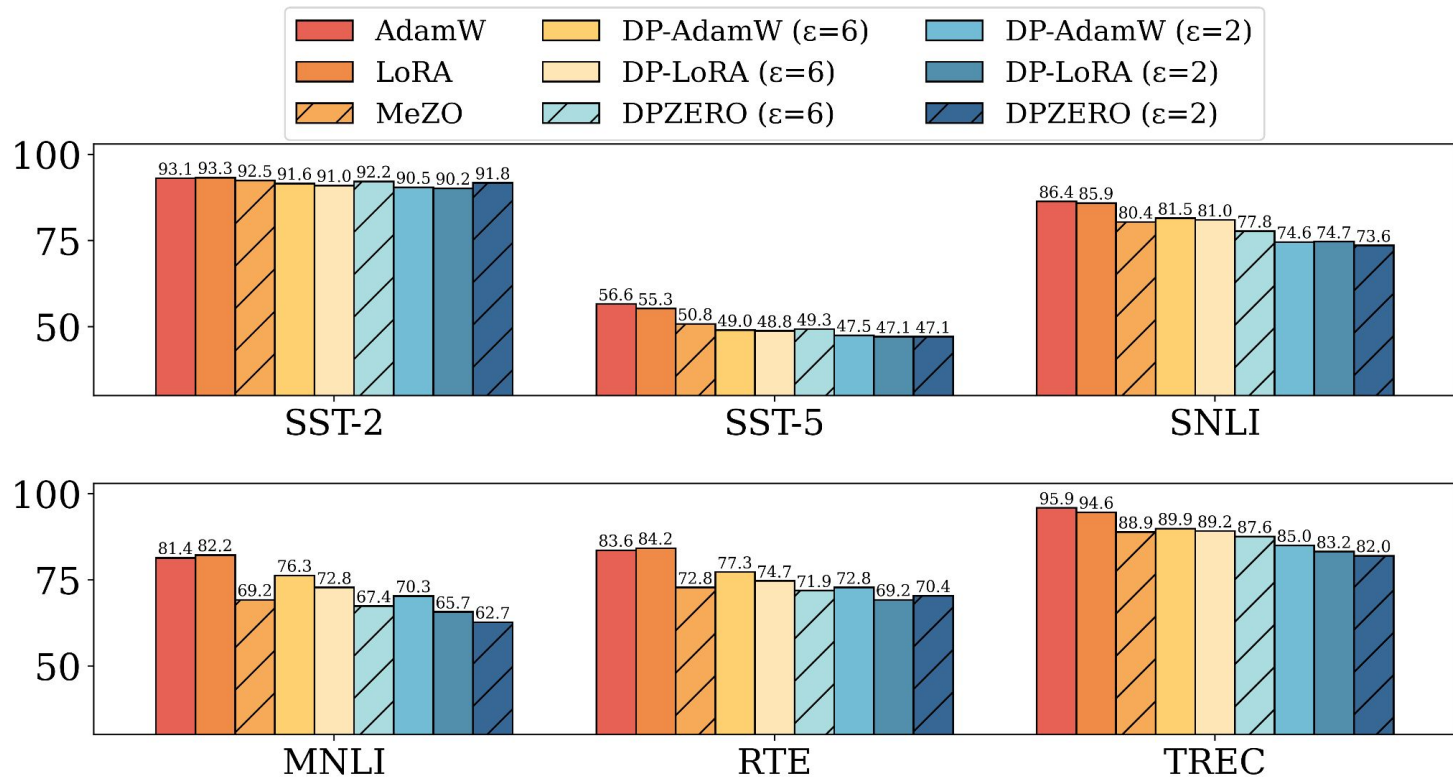


(b) $\text{Tr}(A) = \mathcal{O}(\sqrt{d})$.



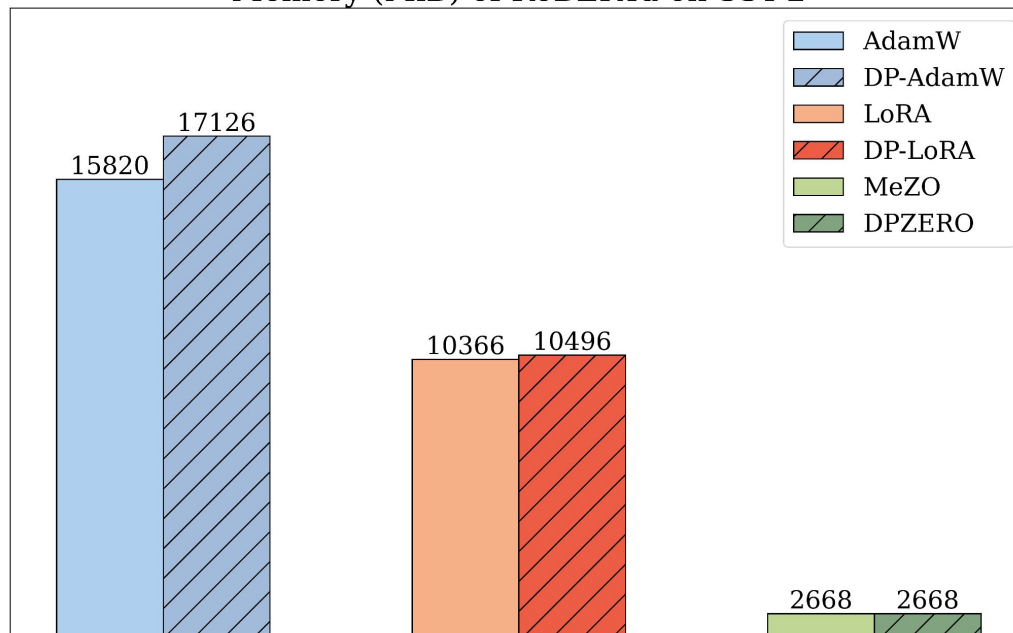
(c) $\text{Tr}(A) = \mathcal{O}(\log d)$.

Private Fine-Tuning RoBERTa (355M)



Private Fine-Tuning RoBERTa (355M)

Memory (MiB) of RoBERTa on SST-2



- **DPZero**: Nearly **no additional cost**
 - ✓ Efficient **per-sample loss clipping** (v.s. **per-sample gradient clipping** in DP-GD)

Private Fine-Tuning OPT

