

## A PROOF OF DIFFERENTIAL PRIVACY

*Proof of Theorem 3.* Define the  $\ell_2$  **sensitivity** of any function  $g$  to be  $\Delta g = \sup_{S, S'} \|g(S) - g(S')\|_2$  where the supreme is over all neighboring  $(S, S')$ . Then the **Gaussian mechanism**  $\hat{g}(S) = g(S) + \sigma \Delta g \cdot \mathcal{N}(0, \mathbf{I})$ .

$\sigma$  denotes the “Noise multiplier”, which corresponds to the noise-level when a Gaussian mechanism is applied to a query with sensitivity 1.

Observe that automatic clipping (AUTO-V and AUTO-S (4.1)) ensures the bounded global-sensitivity of the stochastic gradient as in Abadi’s clipping. Aligning the noise-multiplier (rather than the noise-level itself) ensures that the the noise-to-sensitivity ratio  $\frac{\sigma \Delta g}{\Delta g} = \sigma$  is fixed regardless of  $\Delta g$ . The Gaussian mechanism’s privacy guarantees are equivalent. Thus from the privacy accountant perspective, DP-SGD with both Abadi’s clipping and our autoclipping method can be equivalently represented as the adaptive composition of  $T$  Poisson sampled Gaussian Mechanism with sampling probability  $B/n$  and noise multiplier  $\sigma$ .  $\square$

## B PROOF OF AUTOMATICITY

### B.1 NON-ADAPTIVE DP OPTIMIZERS

*Proof of Theorem 1.* We prove Theorem 1 by showing that, DP-SGD using  $R$ -dependent AUTO-S with learning rate  $\eta$  and weight decay  $\lambda$  is equivalent to  $R$ -independent AUTO-S with learning rate  $\eta R$  and weight decay  $\lambda/R$ . We claim other non-adaptive optimizers such as HeavyBall and NAG can be easily shown in a similar manner.

Recall the standard SGD with weight decay is

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left( \sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} + \lambda \mathbf{w}_t \right)$$

Replacing the standard gradient  $\sum_i \frac{\partial l_i}{\partial \mathbf{w}_t}$  with the private gradient, we write the  $R$ -dependent case as

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \left( \sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} \cdot R / \left\| \frac{\partial l_i}{\partial \mathbf{w}_t} \right\|_2 + \sigma R \cdot \mathcal{N}(0, \mathbf{I}) + \lambda \mathbf{w}_t \right) \\ &= \mathbf{w}_t - \eta R \left( \sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} / \left\| \frac{\partial l_i}{\partial \mathbf{w}_t} \right\|_2 + \sigma \cdot \mathcal{N}(0, \mathbf{I}) \right) - \eta \lambda \mathbf{w}_t \end{aligned}$$

which is clearly equivalent to the  $R$ -independent case:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta' \left( \sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} / \left\| \frac{\partial l_i}{\partial \mathbf{w}_t} \right\|_2 + \sigma \cdot \mathcal{N}(0, \mathbf{I}) + \lambda' \mathbf{w}_t \right)$$

if we use  $\eta' = \eta R$  and  $\lambda' = \lambda/R$ .  $\square$

### B.2 ADAPTIVE DP OPTIMIZERS

*Proof of Theorem 2.* We prove Theorem 2 by showing that, DP-AdamW using  $R$ -dependent AUTO-S with learning rate  $\eta$  and weight decay  $\lambda$  is equivalent to  $R$ -independent AUTO-S with the same learning rate  $\eta$  and weight decay  $\lambda/R$ . This is the most complicated case. We claim other adaptive optimizers such as AdaDelta, Adam with weight decay (not AdamW), and NAdam can be easily shown in a similar manner.

Recall the standard AdamW is

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left( \frac{\mathbf{m}_t / (1 - \beta_1)}{\sqrt{\mathbf{v}_t / (1 - \beta_2)}} + \lambda \mathbf{w}_t \right)$$

where  $\beta_1, \beta_2$  are constants,  $\mathbf{g}_t := \sum_i \frac{\partial l_i}{\partial \mathbf{w}_t}$  is the standard gradient,

$$\begin{aligned}\mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \longrightarrow \mathbf{m}_t = \sum_{\tau} \beta_1^{t-\tau} (1 - \beta_1) \mathbf{g}_{\tau}, \\ \mathbf{v}_t &= \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2 \longrightarrow \mathbf{v}_t = \sum_{\tau} \beta_2^{t-\tau} (1 - \beta_2) \mathbf{g}_{\tau}^2.\end{aligned}$$

Replacing the standard gradient with the private gradient  $R\tilde{\mathbf{g}}_t := R(\sum_i \frac{\partial l_i}{\partial \mathbf{w}_t} / \|\frac{\partial l_i}{\partial \mathbf{w}_t}\|_2 + \sigma \cdot \mathcal{N}(0, \mathbf{I}))$ , we write the  $R$ -dependent DP-AdamW as

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left( \frac{\tilde{\mathbf{m}}_t / (1 - \beta_1)}{\sqrt{\tilde{\mathbf{v}}_t / (1 - \beta_2)}} + \lambda \mathbf{w}_t \right)$$

where

$$\begin{aligned}\tilde{\mathbf{m}}_t &= \beta_1 \tilde{\mathbf{m}}_{t-1} + (1 - \beta_1) R\tilde{\mathbf{g}}_t \longrightarrow \tilde{\mathbf{m}}_t = \sum_{\tau} \beta_1^{t-\tau} (1 - \beta_1) R\tilde{\mathbf{g}}_{\tau}, \\ \tilde{\mathbf{v}}_t &= \beta_2 \tilde{\mathbf{v}}_{t-1} + (1 - \beta_2) R^2 \tilde{\mathbf{g}}_t^2 \longrightarrow \tilde{\mathbf{v}}_t = \sum_{\tau} \beta_2^{t-\tau} (1 - \beta_2) R^2 \tilde{\mathbf{g}}_{\tau}^2.\end{aligned}$$

Clearly, the  $R$  factor in the numerator and denominator of  $\frac{\tilde{\mathbf{m}}_t / (1 - \beta_1)}{\sqrt{\tilde{\mathbf{v}}_t / (1 - \beta_2)}}$  cancel each other. Therefore we claim that the  $R$ -dependent DP-AdamW is in fact completely independent of  $R$ .  $\square$

### B.3 AUTOMATIC PER-LAYER CLIPPING

In some cases, the per-layer clipping is desired, where we use a clipping threshold vector  $\mathbf{R} = [R_1, \dots, R_L]$  and each layer uses a different clipping threshold. We claim that DP optimizers under automatic clipping works with the per-layer clipping when  $\mathbf{R}$  is tuned proportionally, e.g.  $\mathbf{R} = R \cdot [a_1, \dots, a_L]$ , but not entry-wise (see counter-example in Fact B.1). One special case is the *uniform per-layer clipping* when  $R_1 = \dots = R_L = R/\sqrt{L}$ . This is widely applied as only one norm  $R$  requires tuning, instead of  $L$  norms in  $\mathbf{R}$ , particularly in the case of deep models with hundreds of layers. The corresponding DP-SGD with AUTO-S in (3.3) gives

$$\mathbf{w}_{t+1}^{(l)} = \mathbf{w}_t^{(l)} - \eta \left( \sum_{i \in B_t} \frac{R}{\sqrt{L}} \frac{\mathbf{g}_{t,i}^{(l)}}{\|\mathbf{g}_{t,i}^{(l)}\| + \gamma} + \sigma R \cdot \mathcal{N}(0, \mathbf{I}) \right)$$

Here the superscript  $(l)$  is the layer index. Clearly  $R$  couples with the learning rate  $\eta$  and the same analysis as in Theorem 1 follows. The adaptive optimizers can be similarly analyzed from Theorem 2.

**Fact B.1.** Changing one clipping threshold in the clipping threshold vector  $\mathbf{R}$  (i.e. not proportionally) can break the coupling with learning rate.

*Proof of Fact B.1.* We prove by a counter-example of  $\mathbf{R}$  in  $\mathbb{R}^2$ . Consider DP-SGD with per-layer clipping thresholds  $(R_1, R_2) = (9, 12)$ :

$$\mathbf{w}_{t+1}^{(l)} = \mathbf{w}_t^{(l)} - \eta \left( \sum_{i \in B} \frac{R_l \mathbf{g}_{t,i,l}}{\|\mathbf{g}_{t,i,l}\|} + \sigma \sqrt{R_1^2 + R_2^2} \cdot \mathcal{N}(0, \mathbf{I}) \right)$$

Increasing  $R_1$  from 9 to 16 changes the update for the first layer

$$\eta \left( \sum_{i \in B} \frac{9 \mathbf{g}_{t,i,l}}{\|\mathbf{g}_{t,i,l}\|} + 15\sigma \cdot \mathcal{N}(0, 1) \right) \rightarrow \eta \left( \sum_{i \in B} \frac{16 \mathbf{g}_{t,i,l}}{\|\mathbf{g}_{t,i,l}\|} + 20\sigma \cdot \mathcal{N}(0, \mathbf{I}) \right)$$

The noise-to-signal ratio decreases from 5/3 to 5/4 for this layer, and increases from 5/4 to 5/3 for the second layer. This breaks the coupling with learning rate, since the coupling does not change the noise-to-signal ratio.  $\square$

## C MAIN RESULTS OF CONVERGENCE FOR DP-SGD WITH AUTOMATIC CLIPPING

### C.1 MAIN PROOF OF CONVERGENCE FOR DP-SGD (THE ENVELOPE VERSION)

*Proof of Theorem 4.* In this section, we prove two parts of Theorem 4.

The first part of Theorem 4 is the upper bound on  $\min_t \mathbb{E}(\|g_t\|)$ , which is a direct result following from Theorem 6, and we prove it in Appendix C.2.

**Theorem 6.** *Under Assumption 5.1, 5.2, 5.3, running DP-SGD with automatic clipping for  $T$  iterations gives*

$$\min_t \mathbb{E}(\|g_t\|) \leq \frac{\xi}{r} + \mathcal{F}\left(\frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left(1 + \frac{\sigma^2 d}{B^2}\right)}; r, \xi, \gamma\right) \quad (\text{C.1})$$

where

- for  $r < 1, \gamma = 0$  and  $\eta \propto 1/\sqrt{T}$ ,  $\mathcal{F}(x) = \frac{x}{\min_{0 < c < 1} f(c, r)}$  and  $f(c, r) := \frac{(1+rc)}{\sqrt{r^2+2rc+1}} + \frac{(1-rc)}{\sqrt{r^2-2rc+1}}$ ; for  $r \geq 1, \gamma = 0$  and  $\eta \propto 1/\sqrt{T}$ ,  $\mathcal{F}(x) = \infty$ ;
- for  $r \geq 1, \gamma > 0$  and  $\eta \propto 1/\sqrt{T}$ ,  $\mathcal{F}$  is the convex envelope of (C.8), and is strictly increasing.

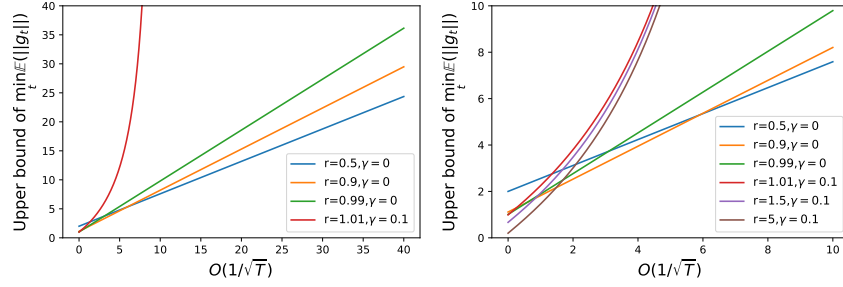


Figure 6: Visualization of upper bound  $\frac{\xi}{r} + \mathcal{F}\left(O(1/\sqrt{T}); r, \xi, \gamma\right)$  for gradient norm, with  $O(1/\sqrt{T})$  in (C.1). Here  $\xi = 1$ . The right plot is a zoom-in (with additional lines) of the left one.

Notice that, (C.1) holds for any  $r > 0$ . However, we have to consider an envelope curve over  $r$  in (C.1) to reduce the upper bound: with AUTO-V clipping ( $\gamma = 0$ ), the upper bound in (C.1) is always larger than  $\xi$  as  $r < 1$ ; we must use AUTO-S clipping ( $\gamma > 0$ ) to reduce the upper bound to zero, as can be seen from Figure 6. In fact, larger  $T$  needs larger  $r$  to reduce the upper bound.

All in all, we specifically focus on  $r \geq 1$  and  $\gamma > 0$ , which is the only scenario that (C.1) can converge to zero. This scenario is also where we prove the second part of Theorem 4.

The second part of Theorem 4 is the asymptotic convergence rate  $O(T^{-1/4})$  of DP-SGD, only possible under  $r \geq 1$  and  $\gamma > 0$ .

By (C.1) in Theorem 6, our upper bound  $\mathcal{G}$  from Theorem 4 can be simplified to

$$\min_{r>0} \frac{\xi}{r} + (\mathcal{M}^{-1})_{ccv} \left( \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left(1 + \frac{\sigma^2 d}{B^2}\right)}; r, \xi, \gamma \right)$$

where the function  $\mathcal{M}^{-1}$  is explicitly defined in (C.8) and the subscript *ccv* means the upper concave envelope. Clearly, as  $T \rightarrow \infty$ ,  $\mathcal{M}^{-1}(\frac{1}{\sqrt{T}}) \rightarrow 0$ . We will next show that the convergence rate of  $\mathcal{M}^{-1}$  is indeed  $O(\frac{1}{\sqrt{T}})$  and the minimization over  $r$  makes the overall convergence rate  $O(T^{-1/4})$ .

Starting from (C.8), we denote  $x = \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L(1 + \frac{\sigma^2 d}{B^2})}$  and write

$$\begin{aligned}
\mathcal{M}^{-1}(x; r, \xi, \gamma) &= \frac{-\frac{\xi}{r}\gamma + (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \gamma\sqrt{(\frac{\xi}{r})^2 + 2\xi x + 2\gamma x + x^2}}{2\gamma - (r^2 - 1)x} \\
&= \left( -\frac{\gamma\xi}{r} + (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \gamma\sqrt{(\frac{\xi}{r})^2 + 2\xi x + 2\gamma x + x^2} \right) \\
&\quad \cdot \frac{1 + \frac{r^2 - 1}{2\gamma}x + O(x^2)}{2\gamma} \\
&= \frac{1}{2\gamma} \left( -\frac{\gamma\xi}{r} + (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \frac{\gamma\xi}{r} \sqrt{1 + \frac{2(\xi + \gamma)r^2 x}{\xi^2} + O(x^2)} \right) \\
&\quad \cdot \left( 1 + \frac{r^2 - 1}{2\gamma}x + O(x^2) \right) \\
&= \frac{1}{2\gamma} \left( -\frac{\gamma\xi}{r} + (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \frac{\gamma\xi}{r} \left( 1 + \frac{(\xi + \gamma)r^2 x}{\xi^2} + O(x^2) \right) \right) \\
&\quad \cdot \left( 1 + \frac{r^2 - 1}{2\gamma}x + O(x^2) \right) \\
&= \frac{1}{2\gamma} \left( (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \frac{\gamma(\xi + \gamma)r x}{\xi} + O(x^2) \right) \cdot \left( 1 + \frac{r^2 - 1}{2\gamma}x + O(x^2) \right) \\
&= \frac{1}{2\gamma} \left( (r^2 - 1)\frac{\xi}{r} + r\gamma + \frac{\gamma(\xi + \gamma)r}{\xi} \right) \cdot x + O(x^2) \\
&= \frac{1}{2\gamma} \left( \frac{(\xi + \gamma)^2}{\xi} r - \frac{\xi}{r} \right) \cdot x + O(x^2)
\end{aligned}$$

Since  $\mathcal{M}^{-1}$  is asymptotically linear as  $x \rightarrow 0$ , we instead study

$$\min_{r>0} \frac{\xi}{r} + \mathcal{M}^{-1}(x; r, \xi, \gamma) \equiv \min_{r>0} \frac{\xi}{r} + \frac{1}{2\gamma} \left( \frac{(\xi + \gamma)^2}{\xi} r - \frac{\xi}{r} \right) \cdot x + O(x^2).$$

That is, ignoring the higher order term for the asymptotic analysis, the  $\mathcal{M}^{-1}$  part converges as  $O(x) = O(1/\sqrt{T})$ , and we visualize this in Figure 8.

Although DP-SGD converges faster than SGD, the former converges to  $\xi/r$  and the latter converges to 0. Thus, taking  $\xi/r$  into consideration, the objective reduces to a hyperbola

$$\frac{\left(\xi(1 - \frac{x}{2\gamma})\right)}{r} + \frac{x(\xi + \gamma)^2}{2\gamma\xi} \cdot r$$

whose minimum over  $r$  is obviously  $2\sqrt{\xi(1 - \frac{x}{2\gamma})\frac{x(\xi + \gamma)^2}{2\gamma\xi}} = O(\sqrt{x}) = O(T^{-1/4})$ .  $\square$

To give more details about the upper bound in (5.2), we demonstrate its dependence on  $\xi$  and  $\gamma$  in Figure 7.

## C.2 MAIN PROOF OF CONVERGENCE FOR DP-SGD (THE NON-ENVELOPE VERSION)

*Proof of Theorem 6.* Consider DP-SGD with AUTO-S clipping

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left( \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} + \sigma \mathcal{N}(0, \mathbf{I}) \right)$$

where  $\tilde{\mathbf{g}}_{t,i}$  is i.i.d. samples of  $\tilde{\mathbf{g}}_t$ , an unbiased estimate of  $\mathbf{g}_t$ , with a bounded variance as described in Assumption 5.3.

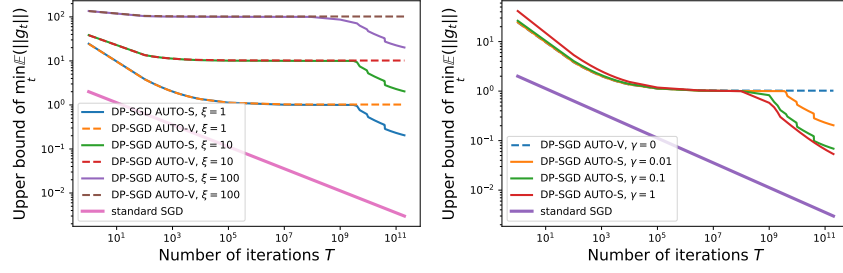


Figure 7: Dependence of the upper bound  $\mathcal{G}$  on  $\xi$  (left) and  $\gamma$  (right). Here the  $O(1/\sqrt{T})$  term is set to 10 and either  $\gamma = 0.01$  (left) or  $\xi = 1$  (right).

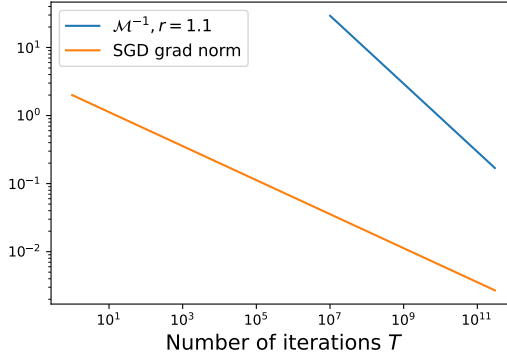


Figure 8: Convergence with respect to  $T$ . Same setting as Figure 5.

By Lipschitz smoothness in Assumption 5.2, and denoting  $Z = \mathcal{N}(0, \mathbf{I})$ , we have

$$\begin{aligned}
 \mathcal{L}_{t+1} - \mathcal{L}_t &\leq \mathbf{g}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_t) + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\
 &= -\eta \mathbf{g}_t^\top \left( \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} + \sigma Z \right) + \frac{L\eta^2}{2} \left\| \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} + \sigma Z \right\|^2 \\
 &\leq -\eta \mathbf{g}_t^\top \left( \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} + \sigma Z \right) \\
 &\quad + L\eta^2 \left( \left\| \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} \right\|^2 + \sigma^2 \|Z\|^2 \right)
 \end{aligned}$$

where the last inequality follows from Cauchy Schwartz.

Given the fact that  $\|\tilde{\mathbf{g}}_{t,i}/(\|\tilde{\mathbf{g}}_{t,i}\| + \gamma)\| \leq 1$ , the expected improvement at one iteration is

$$\begin{aligned}
 \mathbb{E}(\mathcal{L}_{t+1} - \mathcal{L}_t | \mathbf{w}_t) &\leq -\eta \mathbf{g}_t^\top \mathbb{E} \left( \sum_i \frac{\tilde{\mathbf{g}}_{t,i}}{\|\tilde{\mathbf{g}}_{t,i}\| + \gamma} \right) + L\eta^2 (B^2 + \sigma^2 d) \\
 &= -\eta B \mathbf{g}_t^\top \mathbb{E} \left( \frac{\tilde{\mathbf{g}}_t}{\|\tilde{\mathbf{g}}_t\| + \gamma} \right) + L\eta^2 (B^2 + \sigma^2 d)
 \end{aligned} \tag{C.2}$$

Now we want to lower bound  $\mathbf{g}_t^\top \mathbb{E} \left( \frac{\tilde{\mathbf{g}}_t}{\|\tilde{\mathbf{g}}_t\| + \gamma} \right)$  in (C.2).

Write  $\tilde{\mathbf{g}}_t = \mathbf{g}_t + \Delta_t$  where the gradient noise  $\Delta_t$  follows  $\mathbb{E}\Delta_t = 0, \mathbb{E}\|\Delta_t\| < \xi$  by Assumption 5.3. Then

$$\begin{aligned} \mathbf{g}_t^\top \mathbb{E} \left( \frac{\tilde{\mathbf{g}}_t}{\|\tilde{\mathbf{g}}_t\| + \gamma} \right) &= \mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} \right) \\ &= \frac{1}{2} \mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) + \frac{1}{2} \mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} \middle| \Delta_t \in H_- \right) \\ &= \frac{1}{2} \mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) + \frac{1}{2} \mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 - \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t - \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) \end{aligned}$$

where we use the hyperplane perpendicular to  $\mathbf{g}_t$  to divide the support of  $\Delta_t$  into two half-spaces:

$$H_+ := \{\mathbf{v} : \mathbf{g}_t^\top \mathbf{v} > 0\}, \quad H_- := \{\mathbf{v} : \mathbf{g}_t^\top \mathbf{v} < 0\}.$$

We use the symmetry assumption in Assumption 5.3 to get

$$\mathbb{P}(\Delta_t \in H_+) = \mathbb{P}(\Delta_t \in H_-) = \frac{1}{2}$$

and notice that  $\Delta_t \stackrel{D}{=} -\Delta_t$ , i.e., if  $\Delta_t \in H_+$ , then  $-\Delta_t \in H_-$  with the same distribution.

The next result further gives a lower bound for  $\mathbf{g}_t^\top \mathbb{E} \left( \frac{\tilde{\mathbf{g}}_t}{\|\tilde{\mathbf{g}}_t\| + \gamma} \right)$  using  $\|\mathbf{g}_t\|$ .

**Lemma C.1.**

$$\mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} + \frac{\|\mathbf{g}_t\|^2 - \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t - \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) \geq \min_{0 < c \leq 1} f(c, r; \frac{\gamma}{\|\mathbf{g}_t\|}) \cdot (\|\mathbf{g}_t\| - \xi/r)$$

for any  $r > 0$  and  $f(c, r; \Gamma) = \frac{(1+rc)}{\sqrt{r^2+2rc+1+\Gamma}} + \frac{(1-rc)}{\sqrt{r^2-2rc+1+\Gamma}}$ .

For the simplicity of notation, we denote the distance measure

$$\mathcal{M}(\|\mathbf{g}_t\| - \xi/r; r, \xi, \gamma) = \min_{0 < c \leq 1} f \left( c, r; \frac{\gamma}{\|\mathbf{g}_t\|} \right) \cdot (\|\mathbf{g}_t\| - \xi/r) \quad (\text{C.3})$$

and leave the fine-grained analysis (e.g. its explicit form in some scenarios) at the end of this section.

Using the lower bound from Lemma C.1, the expected improvement (C.2) becomes

$$\mathbb{E}(\mathcal{L}_{t+1} - \mathcal{L}_t | \mathbf{w}_t) \leq -\frac{\eta B}{2} \mathcal{M}(\|\mathbf{g}_t\| - \xi/r) + L\eta^2 B^2 \left( 1 + \frac{\sigma^2 d}{B^2} \right)$$

Now extend the expectation over randomness in the trajectory, and perform a telescoping sum over the iterations

$$\begin{aligned} \mathcal{L}_0 - \mathcal{L}_* &\geq \mathcal{L}_0 - \mathbb{E}\mathcal{L}_T = \sum_t \mathbb{E}(\mathcal{L}_t - \mathcal{L}_{t+1}) \\ &\geq \frac{\eta B}{2} \mathbb{E} \left( \sum_t \mathcal{M}(\|\mathbf{g}_t\| - \xi/r) \right) - TL\eta^2 B^2 \left( 1 + \frac{\sigma^2 d}{B^2} \right) \end{aligned}$$

Substituting  $\eta B = \eta_0/\sqrt{T}$  where  $\eta_0$  is a base learning rate, we have

$$2(\mathcal{L}_0 - \mathcal{L}_*) \geq \sqrt{T}\eta_0 \mathbb{E} \left( \frac{1}{T} \sum_t \mathcal{M}(\|\mathbf{g}_t\| - \xi/r) \right) - 2L\eta_0^2 \left( 1 + \frac{\sigma^2 d}{B^2} \right)$$

and finally

$$\mathbb{E} \left( \frac{1}{T} \sum_t \mathcal{M}(\|\mathbf{g}_t\| - \xi/r) \right) \leq \frac{1}{\sqrt{T}} \left[ \frac{2(\mathcal{L}_0 - \mathcal{L}_*)}{\eta_0} + 2L\eta_0 \left( 1 + \frac{\sigma^2 d}{B^2} \right) \right] \quad (\text{C.4})$$

With  $\eta_0$  chosen properly at  $\eta_0 = \sqrt{\frac{\mathcal{L}_0 - \mathcal{L}_*}{L(1 + \frac{\sigma^2 d}{B^2})}}$ , the hyperbola on the right hand side in (C.4) is minimized to  $4\sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L(1 + \frac{\sigma^2 d}{B^2})}$ , and we obtain

$$\mathbb{E} \left( \frac{1}{T} \sum_t \mathcal{M}(\|g_t\| - \xi/r) \right) \leq \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)}$$

Since the minimum of a sequence is smaller than the average, we have

$$\min_t \mathbb{E}(\mathcal{M}(\|g_t\| - \xi/r)) \leq \frac{1}{T} \sum_t \mathbb{E}(\mathcal{M}(\|g_t\| - \xi/r)) \leq \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)} \quad (\text{C.5})$$

We claim that  $\mathcal{M}$  may not be concave or convex. Therefore we use  $\mathcal{M}_{cvx}$  to denote its lower convex envelope, i.e. the largest convex function that is smaller than  $\mathcal{M}$ . Then by Jensen's inequality (C.5) becomes

$$\min_t \mathcal{M}_{cvx}(\mathbb{E}(\|g_t\| - \xi/r)) \leq \min_t \mathbb{E}(\mathcal{M}_{cvx}(\|g_t\| - \xi/r)) \leq \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)} \quad (\text{C.6})$$

It is obvious that  $\mathcal{M}_{cvx}$  is increasing as  $\mathcal{M}$  is increasing by Theorem 8. Hence,  $(\mathcal{M}_{cvx})^{-1}$  is also increasing, as the inverse of  $\mathcal{M}_{cvx}$ . We write (C.6) as

$$\min_t \mathbb{E}(\|g_t\| - \xi/r) \leq (\mathcal{M}_{cvx})^{-1} \left( \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)} \right)$$

and equivalently

$$\min_t \mathbb{E}(\|g_t\|) \leq \frac{\xi}{r} + (\mathcal{M}_{cvx})^{-1} \left( \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)} \right) \quad (\text{C.7})$$

Finally, we derive the explicit properties of  $\mathcal{M}(\|g_t\| - \xi/r)$  in Theorem 8. These properties allow us to further analyze on the convergence of  $\mathcal{M}(\|g_t\| - \xi/r)$ , based on AUTO-V and AUTO-S, respectively.

**1. DP-SGD with AUTO-V clipping.** By Theorem 8, we write

$$\mathcal{M}(x; r) = \min_{c \in (0, 1]} f(c, r; 0) \cdot x$$

This is a linear function and thus  $\mathcal{M}_{cvx} = \mathcal{M} = 1/\mathcal{M}_{cvx}^{-1}$ . As a result, we have

$$\min_t \mathbb{E}(\|g_t\|) \leq \frac{\xi}{r} + \frac{1}{\min_{c \in (0, 1]} f(c, r; 0)} \cdot \frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left( 1 + \frac{\sigma^2 d}{B^2} \right)}$$

We note here  $r$  plays an important role under AUTO-V clipping: when  $r < 1$ , we spend more iterations to converge to better and smaller gradient norm  $\xi/r$ ; when  $r \geq 1$ ,  $\min_c f(c, r; 0) = f(1, r; 0) = 0$  and it takes forever to converge. This is demonstrated in the left plot of Figure 5.

**2. DP-SGD with AUTO-S clipping.** By Theorem 8 and for  $r > 1$ , we write

$$\mathcal{M}(x; r, \xi, \gamma) = \left( \frac{\gamma}{(r-1)(x + \xi/r) + \gamma} - \frac{\gamma}{(r+1)(x + \xi/r) + \gamma} \right) \cdot x.$$

Notice that the inverse of a lower convex envelope is equivalent to the upper concave envelope (denoted by the subscript  $ccv$ ) of an inverse. Therefore we can derive  $(\mathcal{M}_{cvx})^{-1} = (\mathcal{M}^{-1})_{ccv}$  with the explicit form

$$\mathcal{M}^{-1}(x; r, \xi, \gamma) = \frac{-\frac{\xi}{r}\gamma + (r^2 - 1)\frac{\xi}{r}x + r\gamma x + \gamma\sqrt{(\frac{\xi}{r})^2 + 2\xi x + 2\gamma x + x^2}}{2\gamma - (r^2 - 1)x}. \quad (\text{C.8})$$

we can derive it based on  $r, \xi, \gamma$  and substitute back to (C.7).

Note that the domain of  $\mathcal{M}^{-1}$  (or the image of  $\mathcal{M}$ ) is  $[0, \frac{\gamma}{r-1} - \frac{\gamma}{r+1})$ .

In comparison to the AUTO-V clipping,  $\mathcal{M}^{-1}$  takes a much more complicated form, as depicted in the middle plot of Figure 5, where  $r > 1$  plays an important role for the gradient norm to converge to zero.  $\square$

### C.3 PROOF OF LEMMA C.1

*Proof of Lemma C.1.* We want to lower bound

$$\mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} + \frac{\|\mathbf{g}_t\|^2 - \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t - \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) \quad (\text{C.9})$$

To simplify the notation, we denote noise-to-signal ratio  $S := \frac{\|\Delta_t\|}{\|\mathbf{g}_t\|}$  and  $c := \cos \theta = \frac{\mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t\| \|\Delta_t\|}$ , with  $\theta$  be the random angle between  $\mathbf{g}_t$  and  $\Delta_t$ . Note that  $0 < c \leq 1$  when  $\Delta_t \in H_+$ .

The term inside the conditional expectation in (C.9) can be written as

$$\begin{aligned} & \frac{(1+Sc)\|\mathbf{g}_t\|^2}{\sqrt{S^2+2Sc+1}\|\mathbf{g}_t\|+\gamma} + \frac{(1-Sc)\|\mathbf{g}_t\|^2}{\sqrt{S^2-2Sc+1}\|\mathbf{g}_t\|+\gamma} \\ = & \|\mathbf{g}_t\| \left( \frac{(1+Sc)}{\sqrt{S^2+2Sc+1}+\gamma/\|\mathbf{g}_t\|} + \frac{(1-Sc)}{\sqrt{S^2-2Sc+1}+\gamma/\|\mathbf{g}_t\|} \right) \end{aligned}$$

Defining  $\Gamma = \gamma/\|\mathbf{g}_t\|$  and

$$f(c, S; \Gamma) := \frac{(1+Sc)}{\sqrt{S^2+2Sc+1}+\Gamma} + \frac{(1-Sc)}{\sqrt{S^2-2Sc+1}+\Gamma}, \quad (\text{C.10})$$

we turn the conditional expectation in (C.9) into

$$\mathbb{E} \left( \frac{\|\mathbf{g}_t\|^2 + \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t + \Delta_t\| + \gamma} + \frac{\|\mathbf{g}_t\|^2 - \mathbf{g}_t^\top \Delta_t}{\|\mathbf{g}_t - \Delta_t\| + \gamma} \middle| \Delta_t \in H_+ \right) = \|\mathbf{g}_t\| \mathbb{E}(f(c, S; \Gamma) | \Delta_t \in H_+) \quad (\text{C.11})$$

for which we want to lower bound  $f(c, S; \Gamma)$  over  $0 < c \leq 1, S > 0, \Gamma > 0$ . We use the next theorem to prepare some helpful properties. The proof can be found in Appendix E.1.

**Theorem 7.** *For  $f$  defined in (C.10), we have*

1.  $f(c, S; \Gamma)$  is strictly decreasing in  $S$  for all  $0 < c < 1$  and  $\Gamma > 0$ .
2. Consequently,  $\min_{c \in (0,1)} f(c, S; \Gamma)$  is strictly decreasing in  $S$ .
3.  $f(c, S; \Gamma)$  is strictly decreasing in  $c$  for all  $S > 1$  and  $\Gamma > 0$ .

We consider a thresholding ratio  $r > 0$  and we will focus on the regime that  $S < r$ . This  $r$  will turn out to measure the minimum gradient norm at convergence: informally speaking,  $\|\mathbf{g}_t\|$  converges to  $\xi/r$ .



By the law of total expectation, (C.11) can be relaxed as follows.

$$\begin{aligned}
& \|g_t\| \mathbb{E} \left( f(c, S; \Gamma) \middle| \Delta \in H_+ \right) \\
&= \|g_t\| \mathbb{E} \left( f(c, S; \Gamma) \middle| \Delta \in H_+, S < r \right) \mathbb{P}(r\|g_t\| > \|\Delta\| \middle| \Delta \in H_+) \\
&\quad + \|g_t\| \mathbb{E} \left( f(c, S; \Gamma) \middle| \Delta \in H_+, S > r \right) \mathbb{P}(r\|g_t\| < \|\Delta\| \middle| \Delta \in H_+) \\
&\geq \|g_t\| \mathbb{E} \left( f(c, S; \Gamma) \middle| \Delta \in H_+, S < r \right) \mathbb{P}(r\|g_t\| > \|\Delta\| \middle| \Delta \in H_+) \\
&\geq \|g_t\| \mathbb{E} \left( f(c, r; \Gamma) \middle| \Delta \in H_+, S < r \right) \mathbb{P}(r\|g_t\| > \|\Delta\| \middle| \Delta \in H_+) \\
&= \|g_t\| \mathbb{E} \left( f(c, r; \Gamma) \middle| \Delta \in H_+, S < r \right) \mathbb{P}(r\|g_t\| > \|\Delta\|) \\
&\geq \min_{c \in (0,1]} f(c, r; \Gamma) \cdot \underbrace{\|g_t\| \mathbb{P}(r\|g_t\| > \|\Delta\|)}_{\star}
\end{aligned} \tag{C.12}$$

where in the first inequality, the ignoring of last term is justified by  $f(c, S; \Gamma) \geq \min_{c \in (0,1]} f(c, S; \Gamma) \geq \min_{c \in (0,1]} f(c, \infty; \Gamma) = 0$ , from the monotonicity (second statement) in Theorem 7.

We first lower bound  $\star$  by applying the Markov's inequality:

$$\mathbb{P}(r\|g_t\| > \|\Delta_t\|) \geq 1 - \frac{\mathbb{E}\|\Delta_t\|}{r\|g_t\|}$$

and hence by Assumption 5.3,

$$\|g_t\| \mathbb{P}(r\|g_t\| > \|\Delta_t\|) \geq \|g_t\| - \mathbb{E}\|\Delta_t\|/r \geq \|g_t\| - \xi/r.$$

Finally, the conditional expectation of interest in (C.9) gives

$$\mathbb{E} \left( \frac{\|g_t\|^2 + g_t^\top \Delta_t}{\|g_t + \Delta_t\|} + \frac{\|g_t\|^2 - g_t^\top \Delta_t}{\|g_t - \Delta_t\|} \middle| \Delta_t \in H_+ \right) \geq \min_{0 < c \leq 1} f(c, r; \frac{\gamma}{\|g_t\|}) \cdot (\|g_t\| - \xi/r)$$

□

#### C.4 PROOF OF THEOREM 8

To derive some properties of  $\min_c f(c, r; \Gamma)$ , we need to compute separately for AUTO-V (without the stability constant,  $\Gamma = 0$ ) and for AUTO-S (with the stability constant,  $\Gamma > 0$ ), as shown in Theorem 8. As we will show, as the number of training iterations  $T \rightarrow \infty$ , DP-SGD with AUTO-V clipping can only compress  $\|g_t\|$  to  $\xi/r$  for  $r < 1$ . However, DP-SGD with AUTO-S clipping can compress  $\|g_t\|$  to  $\xi/r$  to any  $r > 1$ .

##### Theorem 8.

1. For  $0 < r < 1$  and  $\Gamma = 0$ , we have  $\min_{c \in (0,1]} f(c, r; 0) > 0$ . Then Equation (C.11) is lower bounded by

$$\min_{c \in (0,1]} f(c, r; 0) \cdot (\|g_t\| - \xi/r)$$

which is increasing in  $\|g_t\| - \xi/r$ .

2. For  $r \geq 1$  and  $\Gamma = 0$ , we have  $\min_{c \in (0,1]} f(c, r; \Gamma) = f(1, r; 0) = 0$ . In words, (C.9) has a trivial lower bound and Theorem 6 cannot compress  $\|g_t\|$  to  $\xi/r$ .
3. For  $r \geq 1$  and  $\Gamma > 0$ , we have  $\min_{c \in (0,1]} f(c, r; \Gamma) = f(1, r; \Gamma) = \left( \frac{\Gamma}{r+\Gamma-1} - \frac{\Gamma}{r+\Gamma+1} \right)$ . Then Equation (C.11) is lower bounded by

$$\left( \frac{\gamma}{(r-1)\|g_t\| + \gamma} - \frac{\gamma}{(r+1)\|g_t\| + \gamma} \right) \cdot (\|g_t\| - \xi/r)$$

which is increasing in  $\|g_t\| - \xi/r$ .

*Proof.* To prove statement 1, we use the second statement from Theorem 7 and show that  $\min_c f(c, r; 0) > \min_c f(c, \infty; 0) = 0$ . To prove statement 2 and 3, we use the third statement from Theorem 7 and see that  $\min_c f(c, r; \Gamma) = f(1, r; \Gamma)$  with an explicit formula.  $\square$

## D CONVERGENCE RATE OF STANDARD SGD

**Theorem 9.** *Under Assumption 5.1, 5.2, 5.3 (without the symmetry assumption), running the standard non-DP SGD for  $T$  iterations gives, for  $\eta \propto 1/\sqrt{T}$ ,*

$$\min_t \mathbb{E}(\|\mathbf{g}_t\|) \leq \frac{1}{T^{1/4}} \sqrt{2(\mathcal{L}_0 - \mathcal{L}_*)L + \frac{\xi^2}{B}}$$

*Proof of Theorem 9.* Consider the standard SGD

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{\sum_i \tilde{\mathbf{g}}_{t,i}}{B}$$

where  $\tilde{\mathbf{g}}_{t,i}$  is i.i.d. unbiased estimate of  $\mathbf{g}_t$ , with a bounded variance as described in Assumption 5.3.

By Lipschitz smoothness assumption in Assumption 5.2,

$$\mathcal{L}_{t+1} - \mathcal{L}_t \leq \mathbf{g}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_t) + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 = -\eta \mathbf{g}_t^\top \left( \sum_i \frac{1}{B} \tilde{\mathbf{g}}_{t,i} \right) + \frac{L\eta^2}{2} \left\| \sum_i \frac{1}{B} \tilde{\mathbf{g}}_{t,i} \right\|^2$$

The expected improvement at one iteration is

$$\begin{aligned} \mathbb{E}(\mathcal{L}_{t+1} - \mathcal{L}_t | \mathbf{w}_t) &\leq -\eta \mathbf{g}_t^\top \mathbb{E} \tilde{\mathbf{g}}_{t,i} + \frac{L\eta^2}{2} \mathbb{E} \left\| \sum_i \frac{1}{B} \tilde{\mathbf{g}}_{t,i} \right\|^2 \\ &\leq -\eta \|\mathbf{g}_t\|^2 + \frac{L\eta^2}{2} \left( \|\mathbf{g}_t\|^2 + \frac{\xi^2}{B} \right) \end{aligned} \tag{D.1}$$

Now we extend the expectation over randomness in the trajectory, and perform a telescoping sum over the iterations

$$\mathcal{L}_0 - \mathcal{L}_* \geq \mathcal{L}_0 - \mathbb{E} \mathcal{L}_T = \sum_t \mathbb{E}(\mathcal{L}_t - \mathcal{L}_{t+1}) \geq \left( \eta - \frac{L\eta^2}{2} \right) \mathbb{E} \left( \sum_t \|\mathbf{g}_t\|^2 \right) - \frac{TL\eta^2 \xi^2}{2B}$$

Notice that we do not need the symmetry assumption in Assumption 5.3 in the non-DP SGD analysis.

We apply the same learning rate as in Bernstein et al. (2018),  $\eta = \frac{1}{L\sqrt{T}}$ ,

$$2(\mathcal{L}_0 - \mathcal{L}_*) \geq \left( \frac{2}{L\sqrt{T}} - \frac{1}{LT} \right) \mathbb{E} \left( \sum_t \|\mathbf{g}_t\|^2 \right) - \frac{T\xi^2}{BLT} \geq \frac{\sqrt{T}}{L} \mathbb{E} \left( \frac{1}{T} \sum_t \|\mathbf{g}_t\|^2 \right) - \frac{\xi^2}{BL}$$

and finally

$$\min_t \mathbb{E}(\|\mathbf{g}_t\|^2) \leq \mathbb{E} \left( \frac{1}{T} \sum_t \|\mathbf{g}_t\|^2 \right) \leq \frac{1}{\sqrt{T}} \left[ 2(\mathcal{L}_0 - \mathcal{L}_*)L + \frac{\xi^2}{B} \right]$$

Using the Jensen's inequality, we can have

$$\min_t \mathbb{E}(\|\mathbf{g}_t\|) \leq \frac{1}{T^{1/4}} \sqrt{2(\mathcal{L}_0 - \mathcal{L}_*)L + \frac{\xi^2}{B}}$$

$\square$

## E AUXILIARY PROOFS

### E.1 PROOF OF THEOREM 7

*Proof.* We first show  $\frac{df(c, S; \Gamma)}{dS} < 0$  for all  $0 < c < 1, \Gamma > 0$  and  $S > 0$ , as visualized in the left plot of Figure 9. We can explicitly write down the derivative, by WolframAlpha

$$\frac{df(c, S; \Gamma)}{dS} = \frac{-(A\Gamma^2 + B\Gamma + C)}{\sqrt{S^2 - 2cS + 1}\sqrt{S^2 + 2cS + 1}(\Gamma + \sqrt{S^2 - 2cS + 1})^2(\Gamma + \sqrt{S^2 + 2cS + 1})^2} \quad (\text{E.1})$$

with

$$\begin{aligned} A(c, S) &= \sqrt{S^2 + 2cS + 1} (3c^2S - 2c(S^2 + 1) + S) + \sqrt{S^2 - 2cS + 1} (3c^2S + 2c(S^2 + 1) + S) \\ B(c, S) &= 4S \left[ (S^2 + 1)(1 - c^2) + c^2\sqrt{S^2 + 2cS + 1}\sqrt{S^2 - 2cS + 1} \right] \\ C(c, S) &= (1 - c^2)S \left[ (S^2 - 2cS + 1)^{3/2} + (S^2 + 2cS + 1)^{3/2} \right] \end{aligned}$$

It is obvious that, since  $c < 1$ ,

$$S^2 \pm 2cS + 1 > S^2 \pm 2cS + c^2 = (S \pm c)^2 \geq 0. \quad (\text{E.2})$$

From (E.2), the denominator in (E.1) is positive and it suffices to show  $A\Gamma^2 + B\Gamma + C > 0$  for all  $0 < c < 1$  and  $S > 0$ , in order to show  $\frac{df}{dS} < 0$ .

Also from (E.2), we can easily see  $B(c, S) > 0$  and  $C(c, S) > 0$ . We will show that  $A(c, S) > 0$  in Lemma E.1, after very heavy algebraic computation.

Now we can claim that  $A\Gamma^2 + B\Gamma + C > 0$  by Fact E.3, and complete the proof of the first statement.

To further see that  $\min_c f(c, S; \Gamma)$  is decreasing in  $S$ , let us denote  $c^*(x; \Gamma) := \arg \min_{c \in [0, 1]} f(c, x; \Gamma)$ . Then considering  $S < S'$ , we prove the second statement by observing

$$\min_c f(c, S; \Gamma) = f(c^*(S; \Gamma), S; \Gamma) > f(c^*(S; \Gamma), S'; \Gamma) \geq \min_c f(c, S'; \Gamma).$$

This statement is also visualized in the right plot of Figure 9.

We next show  $\frac{df(c, S; \Gamma)}{dc} < 0$  for all  $0 < c < 1, \Gamma > 0$  and  $S > 1$ . We can explicitly write down the derivative, by WolframAlpha

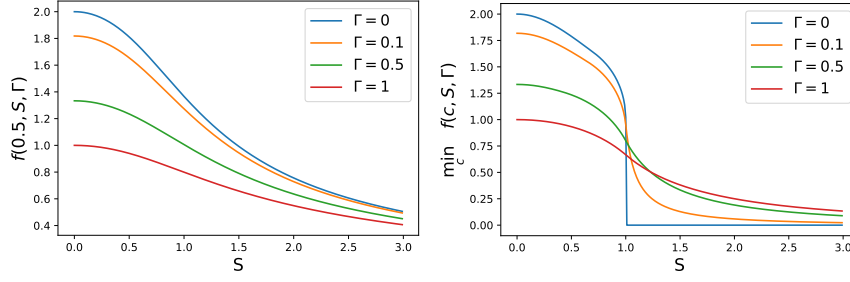
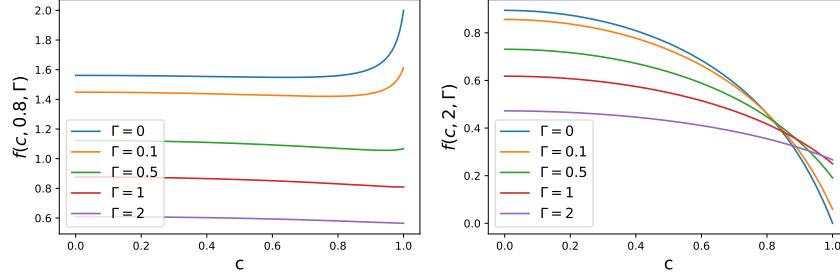
$$\frac{df(c, S; \Gamma)}{dc} = \frac{-S(A'\Gamma^2 + B'\Gamma + C')}{\sqrt{S^2 - 2cS + 1}\sqrt{S^2 + 2cS + 1}(\Gamma + \sqrt{S^2 - 2cS + 1})^2(\Gamma + \sqrt{S^2 + 2cS + 1})^2} \quad (\text{E.3})$$

with

$$\begin{aligned} A'(c, S) &= \left[ (S^2 + 3cS + 2)\sqrt{S^2 - 2cS + 1} - (S^2 - 3cS + 2)\sqrt{S^2 + 2cS + 1} \right] \\ B'(c, S) &= 4Sc \left[ \sqrt{S^2 + 2cS + 1}\sqrt{S^2 - 2cS + 1} + (S^2 - 1) \right] \\ C'(c, S) &= S \left[ (c + S)(S^2 - 2cS + 1)^{3/2} + (c - S)(S^2 + 2cS + 1)^{3/2} \right] \end{aligned}$$

Clearly  $B'(c, S) > 0$  and  $C'(c, S) > 0$ , since  $S^2 + 2cS + 1 > S^2 - 2cS + c^2 = (S - c)^2 \geq 0$ . And we will show  $A'(c, S) > 0$  in Lemma E.2, after some algebra.

We again claim that  $A'\Gamma^2 + B'\Gamma + C' > 0$  by Fact E.3, which guarantees that the numerator in (E.3) is negative and that  $\frac{df}{dc} < 0$ . This is visualized in Figure 10.  $\square$

Figure 9: Visualization of  $f(0.5, S, \Gamma)$  (left) and  $\min_{0 \leq c \leq 1} f(c, S, \Gamma)$  over  $S > 0$ .Figure 10: Visualization of  $f(c, 0.8, \Gamma)$  (left) and  $f(c, 2, \Gamma)$  over  $0 \leq c \leq 1$ .

## E.2 PROOF OF LEMMA E.1

**Lemma E.1.** For all  $0 < c < 1$  and  $S > 0$ ,

$$A := \sqrt{S^2 + 2cS + 1} (3c^2S - 2c(S^2 + 1) + S) + \sqrt{S^2 - 2cS + 1} (3c^2S + 2c(S^2 + 1) + S) > 0.$$

*Proof.* We prove by contradiction. Suppose

$$\sqrt{S^2 + 2cS + 1} (3c^2S - 2c(S^2 + 1) + S) + \sqrt{S^2 - 2cS + 1} (3c^2S + 2c(S^2 + 1) + S) < 0.$$

Then

$$0 < \sqrt{S^2 - 2cS + 1} (3c^2S + 2c(S^2 + 1) + S) < -\sqrt{S^2 + 2cS + 1} (3c^2S - 2c(S^2 + 1) + S).$$

where the first inequality comes from  $S^2 - 2cS + 1 > S^2 - 2cS + c^2 = (S - c)^2 \geq 0$ .

Squaring everything gives

$$(S^2 - 2cS + 1) (3c^2S + 2c(S^2 + 1) + S)^2 < (S^2 + 2cS + 1) (3c^2S - 2c(S^2 + 1) + S)^2.$$

Taking the difference gives

$$4cS(2 + 3S^2 - 9c^4S^2 + 2S^4 + 2c^2(1 - S^2 + S^4)) < 0$$

Given that  $c > 0, S > 0$ , we have

$$2 + 3S^2 - 9c^4S^2 + 2S^4 + 2c^2(1 - S^2 + S^4) < 0$$

Denoting  $X := S^2$  and viewing the above as a quadratic polynomial of  $X$ , we have

$$\underbrace{(2c^2 + 2)X^2 + (3 - 2c^2 - 9c^4)X + (2c^2 + 2)}_{\textcircled{1}} < 0$$

Using the closed-form minimizer of quadratic polynomial  $\textcircled{1}$ , after some heavy algebra, one can check the minimum of  $\textcircled{1}$  is

$$\frac{(1 + 3c^2)^2(1 - c^2)(7 + 9c^2)}{8(1 + c^2)}$$

which is clearly positive. Contradiction!  $\square$

### E.3 PROOF OF LEMMA E.2

**Lemma E.2.** For all  $0 < c < 1$  and  $S > 1$ ,

$$(S^2 + 3cS + 2)\sqrt{S^2 - 2cS + 1} - (S^2 - 3cS + 2)\sqrt{S^2 + 2cS + 1} > 0.$$

*Proof.* Notice that  $(S^2 + 3cS + 2) > S^2 + 2 > 0$  and  $\sqrt{S^2 \pm 2cS + 1} > 0$ . Therefore if  $S^2 - 3cS + 2 \leq 0$ , we are done.

Otherwise, we prove by contradiction and suppose

$$0 < (S^2 + 3cS + 2)\sqrt{S^2 - 2cS + 1} < (S^2 - 3cS + 2)\sqrt{S^2 + 2cS + 1}.$$

under the condition that  $S^2 - 3cS + 2 > 0$ .

Squaring everything gives

$$(S^2 + 3cS + 2)^2(S^2 - 2cS + 1) < (S^2 - 3cS + 2)^2(S^2 + 2cS + 1).$$

Taking the difference gives

$$cS(8 + 20S^2 - 36c^2S^2 + 8S^4) < 0$$

Given that  $c > 0, S > 0$ , we have

$$2 + 5S^2 - 9c^2S^2 + 2S^4 < 0$$

Denoting  $X := S^2$  and viewing the above as a quadratic polynomial of  $X$ , we have, for  $X > 1$ ,

$$\underbrace{2X^2 + (5 - 9c^2)X + 2}_{\textcircled{2}} < 0$$

The closed-form minimizer of quadratic polynomial  $\textcircled{2}$  is  $\frac{(9c^2-5)}{4}$ . Given that  $0 < c < 1$ , we must have  $-\frac{5}{4} < \frac{9c^2-5}{4} < 1$ . Hence the minimizer is not within the feasible domain  $(1, \infty)$  of  $X$ . Thus the minimum of  $\textcircled{2}$  is achieved with  $X = 1$  at  $9(1 - c^2)$ . This is positive. Contradiction!  $\square$

### E.4 PROOF OF FACT E.3

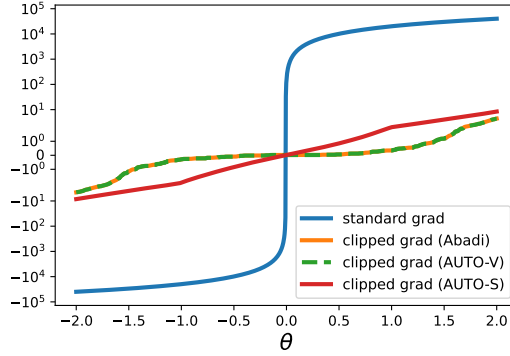
**Fact E.3.** For a quadratic polynomial  $Ax^2 + Bx + C$  with  $A, B, C > 0$ , the minimum value on the domain  $x \geq 0$  is  $C$ , at  $x = 0$ . Therefore  $Ax^2 + Bx + C > 0$ .

*Proof.* Since  $A > 0$ , the quadratic polynomial is convex and increasing on the domain  $x > -\frac{B}{2A}$ . Since  $B > 0$  as well, we know  $-\frac{B}{2A} < 0$  and hence the quadratic polynomial is strictly increasing on  $x > 0$ . Therefore the minimum value is achieved when  $x = 0$ , and we obtain  $Ax^2 + Bx + C \geq C > 0$  for all  $x \geq 0$ .  $\square$

## F EXAMPLES OF LAZY REGIONS

### F.1 BALANCED BINARY CLASSIFICATION

We describe the data generation in Section 3.3. The label is uniformly  $\pm 1$ , that is  $\mathbb{P}(y_i = +1) = \mathbb{P}(y_i = -1) = 0.5$ . We have 10000 positive and negative samples  $x_i \sim \mathcal{N}(y_i, 1)$ . We consider a logistic regression model  $\mathbb{P}(Y = y|x) = \mathbb{I}(y = 1) \cdot \text{Sigmoid}(x + \theta) + \mathbb{I}(y = -1) \cdot (1 - \text{Sigmoid}(x + \theta)) = \frac{1}{1 + e^{-y(\theta + x)}}$ , where  $\theta \in \mathbb{R}$  is the intercept. The gradient with respect to this only trainable parameter is  $\frac{\partial \mathcal{L}_i}{\partial \theta} = -y \left(1 - \frac{1}{1 + e^{-y(\theta + x)}}\right)$ . We set the clipping threshold  $R = 0.01$  and the stability constant  $\gamma = 0.01$ .

Figure 11: Scalar gradient  $\frac{\partial \mathcal{L}}{\partial \theta}$  at each  $\theta$ .

## F.2 MEAN ESTIMATION ON GAUSSIAN MIXTURE DATA

We also observe the lazy region issue in the mean estimation problem  $\min_{\theta} \frac{1}{2} \|\theta - x_i\|^2$ . Here  $\mathbb{P}(x_i \sim \mathcal{N}(4, 1)) = \mathbb{P}(x_i \sim \mathcal{N}(4, 1)) = 0.5$ . We have 10000 samples from each Gaussian distribution. The regular minimum is clearly  $\sum_i x_i \rightarrow 0$ , where the regular gradient and AUTO-S clipped gradient vanish. Yet both AUTO-V and Abadi’s clipping lose motivation to update the mean estimator on the interval  $(-1, 1)$ . We set the clipping threshold  $R = 0.01$  and the stability constant  $\gamma = 0.1$ .

## G EXPERIMENTS SETTINGS

### G.1 IMAGE CLASSIFICATION SETTINGS

We give the experiments settings for computer vision tasks in Table 1.

- **MNIST**: We use the network architecture from Papernot et al. (2021); Tramer & Boneh (2020); Shamsabadi & Papernot (2021), with 40 epochs, 512 batch size, 0.5 learning rate (or 0.005 non-DP learning rate), 0.1 clipping threshold, DP-SGD with 0.9 momentum, and without pretraining. This setting is the same as Tramer & Boneh (2020).
- **FashionMNIST**: We use the same network architecture as MNIST, with 40 epochs, 2048 batch size, 4 learning rate (or 0.04 non-DP learning rate), DP-SGD with 0.9 momentum, and without pretraining. This setting is the same as Tramer & Boneh (2020).
- **CIFAR10 pretrained**: We use the SimCLR model from Chen et al. (2020a)<sup>9</sup>, with 50 epochs, 1024 batch size, 4 learning rate (or 0.04 non-DP learning rate), 0.1 clipping threshold, and DP-SGD with 0.9 momentum. The SimCLR model is pretrained on unlabelled ImageNet dataset. After pretraining, we obtain a feature of dimension 4096 on which a linear classifier is trained privately. This setting is the same as Tramer & Boneh (2020).
- **ImageNette**: We use the ResNet9 (2.5 million parameters) with Mish activation function Misra (2019). We set 50 epochs, 1000 batch size, 0.0005 learning rate (or 0.000005 non-DP learning rate), 1.5 clipping threshold, and use DP-NAdam, without pretraining. This setting is the same as Klause et al. (2022) except we did not apply the learning rate decaying scheduler.
- **CelebA (Smiling and Male and Multi-label)** We use the same ResNet9 as above, with 10 epochs, 500 batch size, 0.001 DP learning rate (or 0.00001 non-DP learning rate), 0.1 clipping threshold, and use DP-Adam, without pretraining. We use the labels ‘Smiling’ and ‘Male’ for two binary classification tasks, with cross-entropy loss. For the multi-label task uses a scalar loss by summing up the 40 binary cross-entropy losses from each label.

We refer the code for MNIST, FashionMNIST, CIFAR10, CIFAR10 pretrained to <https://github.com/ftramer/Handcrafted-DP> by Tramer & Boneh (2020). ResNet9 can be found in <https://github.com/cbenitez81/Resnet9>.

<sup>9</sup>See implementation in <https://github.com/google-research/simclr>.

Throughout all experiments, we do not apply tricks such as random data augmentation (single or multiple times De et al. (2022)), weight standardization Qiao et al. (2019), or parameter averaging Polyak & Juditsky (1992).

## G.2 SENTENCE CLASSIFICATION SETTINGS

We experiment on five datasets in Table 2 and Table 3.

- **MNLI(m)** MNLI-matched, the matched validation and test splits from Multi-Genre Natural Language Inference Corpus.
- **MNLI(mm)** MNLI-mismatched, the matched validation and test splits from Multi-Genre Natural Language Inference Corpus.
- **QQP** The Quora Question Pairs2 dataset.
- **QNLI** The Stanford Question Answering dataset.
- **SST2** The Stanford Sentiment Treebank dataset.

The datasets are processed and loaded from Huggingface Lhoest et al. (2021), as described in <https://huggingface.co/datasets/glue>. We follow the same setup as Yu et al. (2021) and Li et al. (2021). We refer the interested readers to (Li et al., 2021, Appendix G,H,I,K,N) for more details.

We emphasize that our automatic clipping uses exactly the same hyperparameters as the Abadi’s clipping in Li et al. (2021), which is released in their Private-Transformers library <sup>10</sup>.

Dataset	MNLI(m/mm)	QQP	QNLI	SST2
Epoch	18	18	6	3
Batch size	6000	6000	2000	1000
clipping threshold $R$	0.1	0.1	0.1	0.1
DP learning rate	5e-4	5e-4	5e-4	5e-4
non-DP learning rate	5e-5	5e-5	5e-5	5e-5
learning rate decay	Yes	Yes	Yes	Yes
AdamW weight decay	0	0	0	0
Max sequence length	256	256	256	256

Table 5: Hyperparameters of automatic clipping and Abadi’s clipping, for sentence classification in Table 2 and Table 3, using either RoBERTa base or large.

Notice that we use DP learning rate 5e-4 across tasks for the  $R$ -dependent automatic DP-Adam, which is equivalent to  $R$ -independent automatic DP-Adam with the same learning rate. We demonstrate that the results are not sensitive to learning rates around the optimal choice. That is, the automatic clipping does not eliminate  $R$  at the cost of more difficult tuning of learning rate.

learning rate	1e-4	3e-4	5e-4	8e-4	1e-3
RoBERTa-base	93.92	94.38	94.49	94.72	93.35
RoBERTa-large	95.76	96.21	96.21	96.33	95.99

Table 6: SST2 accuracy with respect to learning rate.

## G.3 TABLE-TO-TEXT GENERATION SETTINGS

We experiment multiple GPT2 models on E2E dataset from Huggingface Lhoest et al. (2021) in Table 4. We follow the same setup as Li et al. (2021), and our automatic clipping uses exactly the

<sup>10</sup>See [https://github.com/lxuechen/private-transformers/blob/main/examples/classification/run\\_wrapper.py](https://github.com/lxuechen/private-transformers/blob/main/examples/classification/run_wrapper.py)

same hyperparameters as the Abadi’s clipping in Li et al. (2021), which is released in their Private-Transformer library <sup>11</sup>.

Model	GPT2	GPT2 medium	GPT2 large
Epoch	10	10	10
Batch size	1024	1024	1024
clipping threshold $R$	0.1	0.1	0.1
DP learning rate	2e-3	2e-3	2e-3
non-DP learning rate	2e-4	1e-4	1e-4
learning rate decay	No	No	No
AdamW weight decay	0.01	0.01	0.01
Max sequence length	100	100	100

Table 7: Hyperparameters of automatic clipping and Abadi’s clipping, for the E2E generation task in Table 4.

## H FIGURE ZOO

### H.1 FREQUENCY OF CLIPPING

We show that in all sentence classification tasks, Abadi’s clipping happens on a large proportion of per-sample gradients. This supports the similarity between Abadi’s clipping and AUTO-V in (3.1).

We note that for GPT2, GPT2 medium and GPT2 large, empirically in all iterations 100% of the per-sample gradients are clipped by the Abadi’s clipping, making the performance of Abadi’s clipping equivalent to AUTO-V clipping, as shown in Table 4.

### H.2 STABILITY CONSTANT HELPS AUTO CLIPPING REDUCE GRADIENT NORM

To corroborate our claim in Theorem 6, that the stability  $\gamma$  reduces the gradient norm, we plot the actual gradient norm by iteration.

### H.3 CHOICE OF STABILITY CONSTANT IS ROBUST

We claim in Theorem 6 that, as long as  $\gamma > 0$  in our automatic clipping, the asymptotic convergence rate of gradient norm is the same as that by standard non-private SGD. We plot the ablation study of learning rate and the stability constant  $\gamma$  to show that it is easy to set  $\gamma$ : in Table 2 and Table 3, we adopt learning rate 0.0005, under which a wide range of  $0.0001 < \gamma < 1$  gives similar accuracy. Note that the largest good  $\gamma$  is 1000 times bigger than the smallest good  $\gamma$ .

### H.4 AUTOMATIC CLIPPING AVOIDS ABLATION STUDY

We plot the ablation study of learning rate and clipping threshold in Abadi’s clipping below. This demonstrates that, AUTO-S clipping only requires 1D grid search to tune the learning rate, avoiding the expensive 2D grid search that is unfortunately necessary for the Abadi’s clipping. Hence our automatic clipping can save the tuning effort substantially.

## I FULL TABLE OF GPT2 GENERATION TASK ON E2E DATASET

This is the extended version of Table 4 on E2E dataset. The performance measures are BLEU Papineni et al. (2002), ROGUE-L Lin (2004), NIST Sadjadi et al. (2018), METEOR Banerjee & Lavie (2005), and CIDEr Vedantam et al. (2015) scores. Here  $\epsilon$  is accounted by RDP Mironov (2017), where  $\epsilon = 3$  corresponds to 2.68 if accounted by Gaussian DP Dong et al. (2022); Bu et al. (2020) or to 2.75 if accounted by numerical composition Gopi et al. (2021), and  $\epsilon = 8$  corresponds to 6.77 if accounted by Gaussian DP or to 7.27 if accounted by numerical composition.

<sup>11</sup>See <https://github.com/lxuechen/private-transformers/blob/main/examples/table2text/run.sh>



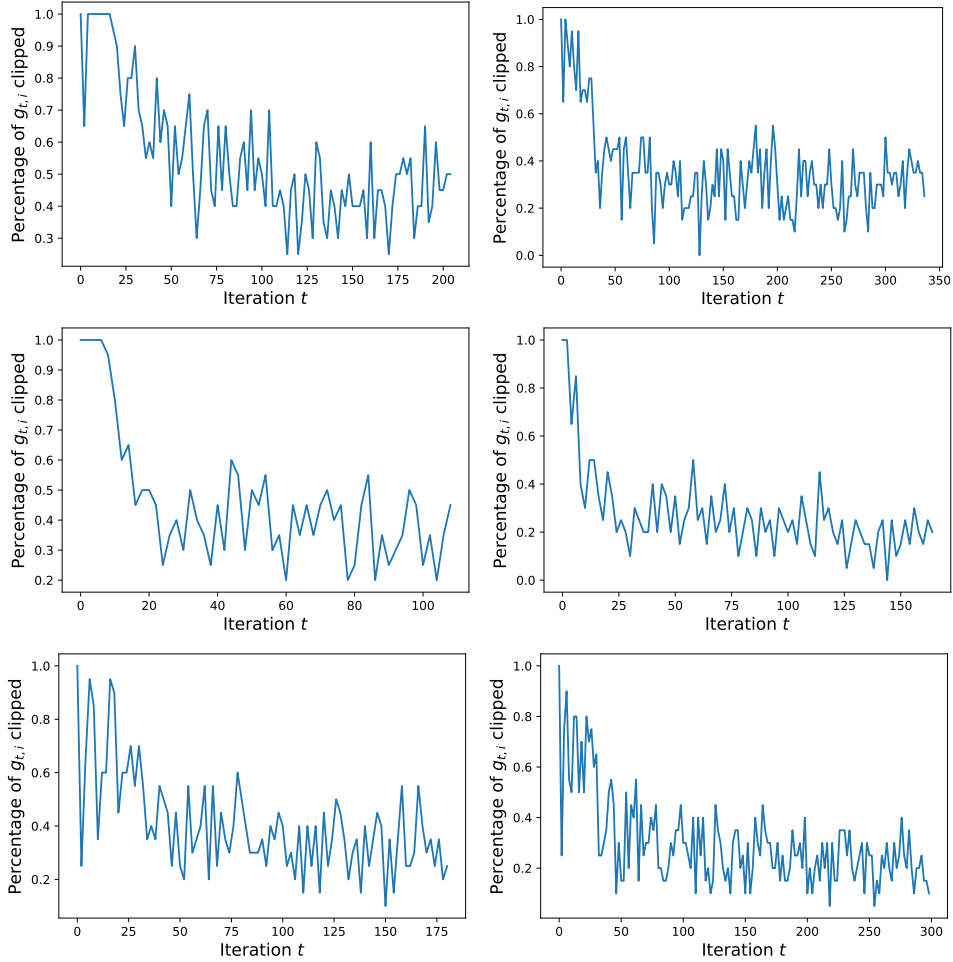


Figure 12: Percentage of clipped per-sample gradients when training with DP-Adam<sub>Abadi</sub> ( $\epsilon = 3$ ), as in Section 6.2. Left panel is RoBERTa-base and right panel is RoBERTa-large. Top row: MNLI. Middle row: QNLI. Bottom row: QQP.

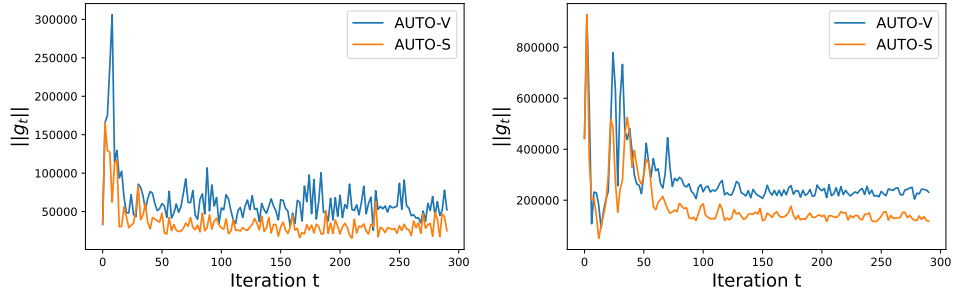


Figure 13: Gradient norm by different automatic clipping methods, on SST2 (left) and MNLI (right), trained with RoBERTa-base.

We observe that GPT2 (163 million parameters), GPT2-medium (406 million), and GPT2-large (838 million), Table 4 trained with our automatic clipping consistently perform better in comparison to other methods. In some cases, LoRA trained with Abadi’s clipping also demonstrates strong performance and it would be interesting to see how LoRA trained with the automatic clipping will behave.

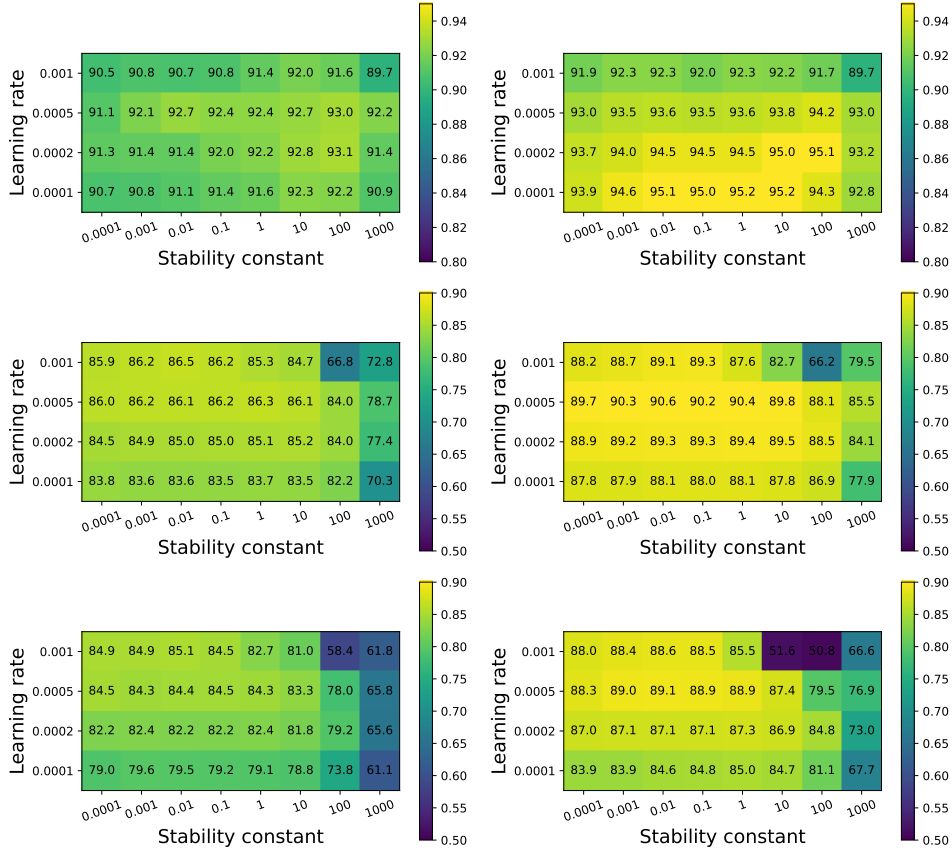


Figure 14: Test accuracy by different stability constant  $\gamma$  and learning rate  $\eta$  in automatic clipping ( $\epsilon = 3$ ). Upper row: SST2 for full 3 epochs. Middle row: QNLI for full 6 epochs. Lower row: QNLI for one epoch. Trained with RoBERTa-base (left) and RoBERTa-large (right).

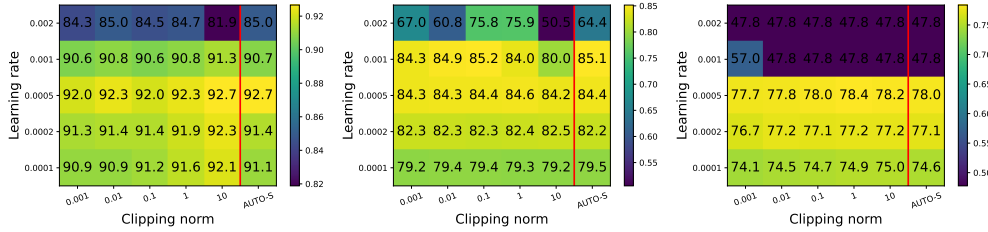


Figure 15: Test accuracy by different clipping threshold  $R$  in DP-Adam<sub>Abadi</sub> and learning rate  $\eta$ , on SST2 (left, 3 epochs) / QNLI (middle, 1 epoch) / MNLI (right, 1 epoch),  $\epsilon = 3$ , trained with RoBERTa-base.

## J FURTHER EXPERIMENTS ON CELEBA DATASET

In this section, we present a complete summary of accuracy results, with DP constraint or not, for the CelebA dataset. We do not apply any data-preprocessing. In the first experiment, we apply a single ResNet on the 40 labels as the multi-task/multi-label learning. In the second experiment, we apply one ResNet on one label. As expected, our automatic DP optimizers have comparable test accuracy to the Abadi’s DP optimizers, but we do not need to tune the clipping threshold for each individual task/label. We also notice that, learning different labels separately gives better accuracy than learning all labels together, though at the cost of heavier computational burden.

Metric	DP guarantee	GPT2 large	GPT2 medium	GPT2							
		full	full	full	full	full	full	full	full	full	full
		AUTO-S	AUTO-S	AUTO-S	AUTO-V	Li et al. (2021)	Hu et al. (2021)	Yu et al. (2021)	Li & Liang (2021)	Li et al. (2021)	Li et al. (2021)
BLEU	$\epsilon = 3$	<b>64.180</b>	<b>63.850</b>	<b>61.340</b>	<b>61.519</b>	<b>61.519</b>	58.153	58.482	47.772	25.920	15.457
	$\epsilon = 8$	<b>64.640</b>	<b>64.220</b>	<b>63.600</b>	63.189	63.189	<b>63.389</b>	58.455	49.263	26.885	24.247
	non-DP	66.840	68.500	69.463	69.463	69.463	69.682	68.328	68.845	65.752	65.731
ROGUE-L	$\epsilon = 3$	<b>67.857</b>	<b>67.071</b>	<b>65.872</b>	65.670	65.670	<b>65.773</b>	65.560	58.964	44.536	35.240
	$\epsilon = 8$	<b>68.968</b>	<b>67.533</b>	<b>67.073</b>	66.429	66.429	<b>67.525</b>	65.030	60.730	46.421	39.951
	non-DP	70.384	71.458	71.359	71.359	71.359	71.709	68.844	70.805	68.704	68.751
NIST	$\epsilon = 3$	<b>7.937</b>	<b>7.106</b>	<b>7.071</b>	<b>6.697</b>	<b>6.697</b>	5.463	5.775	5.249	1.510	0.376
	$\epsilon = 8$	<b>8.301</b>	<b>8.172</b>	<b>7.714</b>	7.444	7.444	<b>7.449</b>	6.276	5.525	1.547	1.01
	non-DP	8.730	8.628	8.780	8.780	8.780	8.822	8.722	8.722	8.418	8.286
METEOR	$\epsilon = 3$	<b>0.403</b>	<b>0.387</b>	<b>0.387</b>	<b>0.384</b>	<b>0.384</b>	0.370	0.331	0.363	0.197	0.113
	$\epsilon = 8$	<b>0.420</b>	<b>0.418</b>	<b>0.404</b>	0.400	0.400	<b>0.407</b>	0.349	0.364	0.207	0.145
	non-DP	0.460	0.449	0.461	0.461	0.461	0.463	0.456	0.445	0.443	0.429
CIDEr	$\epsilon = 3$	<b>2.008</b>	<b>1.754</b>	<b>1.801</b>	<b>1.761</b>	<b>1.761</b>	1.581	1.300	1.507	0.452	0.116
	$\epsilon = 8$	<b>2.163</b>	<b>2.081</b>	<b>1.938</b>	1.919	1.919	<b>1.948</b>	1.496	1.569	0.499	0.281
	non-DP	2.356	2.137	2.422	2.422	2.422	2.491	2.418	2.345	2.180	2.004

Table 8: Test performance on E2E dataset with GPT2. The best two GPT2 models for each row are marked in bold.

### J.1 MULTI-LABEL CLASSIFICATION

We apply ResNet9 as in Appendix G.1 on the multi-label classification task. I.e. the output layer has 40 neurons, each corresponding to one sigmoid cross-entropy loss, that are summed to a single loss and all labels are learnt jointly.

Index	Attributes	Abadi's $\epsilon = 3$	AUTO-S $\epsilon = 3$	Abadi's $\epsilon = 8$	AUTO-S $\epsilon = 8$	non-DP $\epsilon = \infty$
0	5 o Clock Shadow	90.64	90.99↑	90.81	91.28↑	93.33
1	Arched Eyebrows	75.15	76.31↑	76.84	77.11↑	81.52
2	Attractive	75.85	76.10↑	77.50	77.74↑	81.15
3	Bags Under Eyes	80.75	81.12↑	82.15	82.13↓	84.81
4	Bald	97.84	97.87↑	98.04	97.98↓	98.58
5	Bangs	92.71	92.68↓	93.46	93.55↑	95.50
6	Big Lips	67.51	67.78↑	68.34	68.44↑	71.33
7	Big Nose	78.01	80.23↑	76.69	80.59↑	83.54
8	Black Hair	81.92	80.95↓	83.33	83.28↓	88.55
9	Blond Hair	92.25	92.38↑	93.52	93.09↓	95.49
10	Blurry	94.91	94.82↓	95.08	94.90↓	95.78
11	Brown Hair	80.13	82.50↑	83.74	83.89↑	87.79
12	Bushy Eyebrows	88.06	88.23↑	89.72	88.80↓	92.19
13	Chubby	94.72	94.54↓	94.54	94.50↓	95.56
14	Double Chin	95.19	95.49↑	95.50	95.51↑	96.09
15	Eyeglasses	97.06	97.64↑	98.32	98.06↓	99.39
16	Goatee	95.68	95.45↓	95.84	95.87↑	97.06
17	Gray Hair	96.77	96.79↑	97.02	97.03↑	98.06
18	Heavy Makeup	84.96	85.70↑	87.58	87.29↓	90.76
19	High Cheekbones	81.46	81.42↓	82.62	82.72↑	86.62
20	Male	92.05	92.17↑	93.32	93.17↓	97.46
21	Mouth Slightly Open	86.20	86.32↑	87.84	88.48↑	93.07
22	Mustache	96.05	95.96↓	96.08	95.99↓	96.74
23	Narrow Eyes	84.90	84.78↓	85.14	85.18↑	86.98
24	No Beard	91.55	91.67↑	92.29	92.45↑	95.18
25	Oval Face	71.26	71.42↑	71.98	71.25↓	74.62
26	Pale Skin	96.09	96.04↓	96.15	96.17↑	96.93
27	Pointy Nose	70.34	72.11↑	72.23	73.01↑	75.68
28	Receding Hairline	91.53	91.37↓	91.75	91.74↓	92.87
29	Rosy Cheeks	93.26	93.02↓	93.56	93.35↓	94.86
30	Sideburns	96.16	96.09↓	96.27	96.46↑	97.44
31	Smiling	86.39	87.08↑	88.87	88.63↓	92.25
32	Straight Hair	76.20	77.95↑	78.78	78.52↓	80.66
33	Wavy Hair	70.30	71.79↑	73.58	73.19↓	79.15
34	Wearing Earrings	80.53	81.52↑	82.29	82.20↓	87.56
35	Wearing Hat	96.99	96.83↓	97.46	97.31↓	98.68
36	Wearing Lipstick	88.95	88.04↓	89.87	90.72↑	93.49
37	Wearing Necklace	84.59	85.83↑	85.93	85.42↓	86.61
38	Wearing Necktie	93.91	93.91↓	94.43	94.08↓	96.30
39	Young	81.35	81.21↓	82.18	82.52↑	87.18

Table 9: Accuracy on CelebA dataset with settings in Appendix G.1 from one run. The green arrow indicates AUTO-S is better than Abadi’s clipping under the same  $\epsilon$ ; the red arrow indicates otherwise; the black bar indicates the same accuracy.

## J.2 MULTIPLE BINARY CLASSIFICATION

For the second experiment, we apply ResNet9 on each label as a binary classification task. I.e. the output layer has 1 neuron and we run 40 different models for all labels separately.

## K CODE IMPLEMENTATION OF AUTOMATIC CLIPPING

Changing Abadi’s clipping to automatic clipping is easy in available codebases. One can set the clipping  $R = 1$  or any other constant, as explained in Theorem 1 and Theorem 2.

Index	Attributes	Abadi's Single $\epsilon = 8$	AUTO-S Single $\epsilon = 8$	Abadi's Multi $\epsilon = 8$	AUTO-S Multi $\epsilon = 8$	non-DP Multi $\epsilon = \infty$
0	5 o Clock Shadow	92.15	92.29↑	90.81	91.28↑	93.33
1	Arched Eyebrows	81.18	80.19↓	76.84	77.11↑	81.52
2	Attractive	79.31	79.79↑	77.50	77.74↑	81.15
3	Bags Under Eyes	83.52	83.48↓	82.15	82.13↓	84.81
4	Bald	97.89	97.88↓	98.04	97.98↓	98.58
5	Bangs	94.52	94.83↑	93.46	93.55↑	95.50
6	Big Lips	67.32	67.53↑	68.34	68.44↑	71.33
7	Big Nose	82.31	82.36↑	76.69	80.59↑	83.54
8	Black Hair	87.08	86.93↓	83.33	83.28↓	88.55
9	Blond Hair	94.29	94.73↑	93.52	93.09↓	95.49
10	Blurry	94.95	95.20↑	95.08	94.90↓	95.78
11	Brown Hair	87.41	87.19↓	83.74	83.89↑	87.79
12	Bushy Eyebrows	91.23	91.43↑	89.72	88.80↓	92.19
13	Chubby	94.70	94.70–	94.54	94.50↓	95.56
14	Double Chin	95.43	95.43–	95.50	95.51↑	96.09
15	Eyeglasses	98.88	99.14↑	98.32	98.06↓	99.39
16	Goatee	96.12	96.07↓	95.84	95.87↑	97.06
17	Gray Hair	97.48	97.34↓	97.02	97.03↑	98.06
18	Heavy Makeup	88.85	88.72↓	87.58	87.29↓	90.76
19	High Cheekbones	85.66	85.45↓	82.62	82.72↑	86.62
20	Male	95.42	95.70↑	95.53	93.17↓	97.46
21	Mouth Slightly Open	92.67	92.74↑	87.84	88.48↑	93.07
22	Mustache	96.13	96.13–	96.08	95.99↓	96.74
23	Narrow Eyes	85.13	85.13–	85.14	85.18↑	86.98
24	No Beard	94.26	94.58↑	92.29	92.45↑	95.18
25	Oval Face	70.77	73.05↑	71.98	71.25↓	74.62
26	Pale Skin	96.38	96.34↓	96.15	96.17↑	96.93
27	Pointy Nose	71.48	73.37↑	72.23	73.01↑	75.68
28	Receding Hairline	91.51	91.51–	91.75	91.74↓	92.87
29	Rosy Cheeks	93.26	93.35↑	93.56	93.35↓	94.86
30	Sideburns	96.46	96.34↓	96.27	96.46↑	97.44
31	Smiling	90.82	90.87↑	88.87	88.63↓	92.25
32	Straight Hair	79.01	79.01–	78.78	78.52↓	80.66
33	Wavy Hair	77.55	78.83↑	73.58	73.19↓	79.15
34	Wearing Earrings	87.33	87.50↑	82.29	82.20↓	87.56
35	Wearing Hat	98.04	98.11↑	97.46	97.31↓	98.68
36	Wearing Lipstick	92.05	90.46↓	89.87	90.72↑	93.49
37	Wearing Necklace	86.21	86.21–	85.93	85.42↓	86.61
38	Wearing Necktie	95.85	95.94↑	94.43	94.08↓	96.30
39	Young	85.19	84.12↓	82.18	82.52↑	87.18

Table 10: Accuracy on CelebA dataset with settings in Appendix G.1 from one run. ‘Single’ means each attribute is learned separately as a binary classification task. ‘Multi’ means all attributes are learned jointly as a multi-label classification task. The green arrow indicates AUTO-S is better than Abadi’s clipping under the same  $\epsilon$  and the same task; the red arrow indicates otherwise; the black bar indicates the same accuracy.

## K.1 OPACUS

For Opacus Yousefpour et al. (2021) version 1.1.2 (latest), we can implement the all-layer automatic clipping by changing Line 399-401 in <https://github.com/pytorch/opacus/blob/main/opacus/optimizers/optimizer.py> to

```
per_sample_clip_factor = self.max_grad_norm / (per_sample_norms + 0.01)
```

The per-layer automatic clipping requires changing Line 61-63 in <https://github.com/pytorch/opacus/blob/main/opacus/optimizers/perlayeroptimizer.py> to

```
per_sample_clip_factor = max_grad_norm / (per_sample_norms + 0.01)
```

For older version ( $< 1.0$ , e.g. 0.15) of Opacus, we can implement the all-layer automatic clipping by changing Line 223-225 in <https://github.com/pytorch/opacus/blob/v0.15.0/opacus/utils/clipping.py> to

```
per_sample_clip_factor = self.flat_value / (norms[0] + 0.01)
```

or implement the per-layer automatic clipping by changing Line 301-302 in <https://github.com/pytorch/opacus/blob/main/opacus/optimizers/perlayeroptimizer.py> to

```
per_sample_clip_factor = threshold / (norm + 0.01)
clipping_factor.append(per_sample_clip_factor)
```

## K.2 OBJAX

For ObJAX version 1.6.0 (latest), we can implement the automatic clipping in <https://github.com/google/objax/blob/master/objax/privacy/dpsgd/gradient.py> by changing Line 92 to

```
idivisor = self.l2_norm_clip / (total_grad_norm+0.01)
```

and changing Line 145 to

```
idivisor = self.l2_norm_clip/(grad_norms+0.01)
```

## K.3 PRIVATE-TRANSFORMERS

To reproduce our experiments for sentence classification and table-to-text generation, we modify the ‘private-transformers’ codebase of Li et al. (2021). The modification is in [https://github.com/lxuechen/private-transformers/blob/main/private\\_transformers/privacy\\_utils/privacy\\_engine.py](https://github.com/lxuechen/private-transformers/blob/main/private_transformers/privacy_utils/privacy_engine.py), by changing Line 349 to

```
return self.max_grad_norm / (norm_sample + 0.01)
```

and Line 510-512 to

```
coef_sample = self.max_grad_norm * scale / (norm_sample + 0.01)
```