
Automatic Clipping: Differentially Private Deep Learning Made Easier and Stronger

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Per-example gradient clipping is a key algorithmic step that enables practical
2 differentially private (DP) training for deep learning models. The choice of clipping
3 threshold R , however, is shown to be vital for achieving high accuracy under
4 DP. We propose an easy-to-use replacement, called AutoClipping, that eliminates
5 the need to tune R for any DP optimizers, including DP-SGD, DP-Adam, DP-
6 LAMB and many others. The automatic variants are as private and computationally
7 efficient as existing DP optimizers, but require no DP-specific hyperparameters
8 and thus make DP training as amenable as the standard non-private training. We
9 give a rigorous convergence analysis of automatic DP-SGD in the non-convex
10 setting, which shows that it can enjoy an asymptotic convergence rate that matches
11 the standard SGD, under a symmetric noise assumption of the per-sample gradi-
12 ents. We also demonstrate on various language and vision tasks that automatic
13 clipping outperforms or matches the state-of-the-art, and can be easily employed
14 with minimal changes to existing codebases.

15 1 Introduction

16 Deep learning has achieved impressive progress in a wide range of computer vision and natural
17 language processing tasks. These successes are made available, in part, by the collection of large
18 datasets, sometimes containing sensitive private information of individual data points (e.g., chest scan
19 images, DNA sequences). Prior works have illustrated that deep learning models pose severe privacy
20 risks to individual subjects in the training data and are susceptible to various practical attacks. For
21 example, machine learning services such as Google Prediction API and Amazon Machine Learning
22 can leak membership information from the purchase records [58]; if one feeds the GPT2 language
23 model with some specific prefix, the model will autocomplete texts that contain the full name, phone
24 number, email address, etc., from the training data that it memorizes [11].

25 Differential privacy (DP) [21, 23, 22] is a formal definition of privacy that has been shown to prevent
26 the aforementioned privacy risks in deep learning [1]. On a high level, the key difference between the
27 DP deep learning and the regular one is whether the gradient is privately released. In other words,
28 while the standard optimizers update on the summed gradient $\sum_i \mathbf{g}_i$, and DP optimizers update on
29 the *private gradient*:

$$\text{DP Optimizer}(\{\mathbf{g}_i\}_{i=1}^B) = \text{Optimizer}(\overbrace{\sum_i \mathbf{g}_i \cdot \text{Clip}(\|\mathbf{g}_i\|; R) + \sigma R \cdot \mathcal{N}(0, \mathbf{I})}^{\text{private gradient}}) \quad (1.1)$$

$$\text{Standard Optimizer}(\{\mathbf{g}_i\}_{i=1}^B) = \text{Optimizer}(\sum_i \mathbf{g}_i) \quad (1.2)$$

30 Here $\mathbf{g}_i \in \mathbb{R}^d$ is the per-sample gradient of loss l_i , \mathcal{N} is the standard normal random variable, σ is the
31 noise multiplier, and R is the clipping threshold. The clipping function $\text{Clip} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is defined
32 such that $\|\mathbf{g}_i \cdot \text{Clip}(\mathbf{g}_i; R)\| \leq R$. For instance, the DP-SGD in [1] on batch B_t is

$$\text{DP-SGD}_{\text{Abadi}} : \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left(\sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} \min \left(R / \left\| \frac{\partial l_i}{\partial \mathbf{w}_t} \right\|, 1 \right) + \sigma R \cdot \mathcal{N}(0, \mathbf{I}) \right) \quad (1.3)$$

In comparison to the regular training (1.2), two additional DP-specific hyperparameters R and σ need to be determined in DP learning (1.1). On the one hand, setting the noise multiplier σ is easy and can be derived analytically prior to the training. Whenever the privacy budget (ϵ, δ) is determined, one can apply off-the-shelf privacy accounting tools in Section 2.1 to determine σ , based on the subsampling probability p (e.g. expected batch size over sample size) and the number of iterations T :

$$\text{privacy_accountant}(\sigma, p, T; \delta) = \epsilon$$

On the other hand, the choice of clipping threshold R is crucial to the performance of DP models, yet the hyperparameter tuning is much labor-intensive. Recent advances of DP deep learning on ImageNet [34] and on E2E datasets [37], using ResNet18 and GPT2 respectively, illustrate that the performance is very sensitive to R . We have reproduced their results in Figure 1. Observe that on ImageNet, ResNet18 can drop from the highest 45% accuracy to 31% if R is chosen 2 times larger, and to 0.1% if R is chosen 4 times larger. Similar drastic drop can also be observed in [34, Figure 3] even if the noise multiplier $\sigma = 0$. Unlike the noise multiplier σ , the clipping threshold R cannot be inferred from the privacy budget (ϵ, δ) and have to be tuned. Consequently, DP training necessarily requires a 2D grid search for (R, η) , like the lower plot of Figure 1, whereas the regular training only requires an easy 1D grid search for η . Even worse, the difficulty of tuning a per-layer clipping threshold vector [43], i.e. one clipping threshold for one layer, may increase exponentially as the number of layers increases.

To save the effort of tuning R , previous researches have proposed different approaches. In [3, 51, 25], researchers advocate to use data-adaptive information to select R , such as a specified quantile of the gradient norm distribution. These adaptive clipping methods can be a little ad-hoc: they often replace the the need to tune R by the need to tune one or more new hyperparameters, e.g. the quantile to use and the ratio to split the privacy budget between the quantile decision and the gradient perturbation. Another approach used by the practitioners is to replace an expensive 2D grid search by multiple cheaper 1D grid searches. For example, the researchers propose, in [34, Section 3.3] to fine-tune η with non-DP SGD, fix η and sweep over various values of the clipping threshold R with DP-SGD, then further fix R and do one more grid search on η . However, tuning R formally in a data-dependent way (e.g. through cross-validation) introduces additional privacy loss [48], and most existing empirical work does not privately conduct hyperparameter tuning.

We take a completely different route by proposing a new clipping principle that removes R , instead of coming up with methods to find the appropriate R . We term our method as *automatic clipping* (*AutoClipping*) and we term the versions of DP optimizers using it as *automatic DP optimizers*.

We summarize our contributions as follows.

1. We propose the automatic clipping in (4.1) that expunge the clipping threshold from general DP optimizers, allowing DP learning to be as amenable as regular learning.
2. We show that automatic DP optimizers are as private and efficient as existing DP optimizers.
3. We show that automatic DP-SGD converges in the non-convex setting, at the same asymptotic convergence rate as the standard SGD. Our theoretical analysis successfully explains the training behaviors in previous empirical works.
4. We demonstrate the superiority of automatic clipping on a variety of vision and language tasks, especially with large models including ResNet, RoBERTa and GPT2.
5. In Appendix K, we include simple code snippets that demonstrate how easy it is to switch from Abadi’s clipping to our automatic clipping in popular codebases, that implement DP optimizers for deep learning, e.g. Opacus and ObjAX.

2 Preliminaries

2.1 Differential Privacy

We consider the (ϵ, δ) -DP in Definition 2.1, where smaller (ϵ, δ) means stronger privacy guarantee.

Definition 2.1 ([22]). A randomized algorithm M is (ϵ, δ) -differentially private (DP) if for any two neighboring¹ datasets S, S' , and for any event E ,

$$\mathbb{P}[M(S) \in E] \leq e^\epsilon \mathbb{P}[M(S') \in E] + \delta. \quad (2.1)$$

¹ S' is a neighbor of S if one can obtain S' by adding or removing one data point from S .

76 In words, DP restricts the influence of an arbitrary sample, so that the information contributed by such
 77 sample is limited and less vulnerable to privacy attacks. In deep learning, DP is generally achieved by
 78 applying the *subsampling Gaussian mechanism* to privatize the minibatch gradients during training.

79 As illustrated in Equation (1.1), the subsampled Gaussian mechanism involves (1) Sampling a
 80 minibatch by including each data point iid with probability p (2) per-sample gradient clipping to
 81 bound the l_2 norm sensitivity at R and (3) adding independent Gaussian noise proportional to the
 82 sensitivity R and σ , which is derived from the privacy loss ϵ . This can be realized by leveraging a
 83 variety of modern privacy accounting tools, such as those based on Renyi DP (or moments accountant)
 84 [1, 45, 63], Privacy Loss distribution (Fourier accountants) [33, 27, 70], or Gaussian DP [16, 7].

85 2.2 Differentially Private optimizers with general clipping operations

86 Privately released stochastic gradients (through the Gaussian mechanism) can be used to instantiate
 87 various off-the-shelf optimizers, which gives rise to DP-SGD in (1.3), DP-HeavyBall, DP-AdaGrad,
 88 DP-Adam, DP-FedAvg, DP-FedSGD [43], etc. To improve the performance of DP optimizers,
 89 previous researches can be classified into two categories.

90 The first category, where the majority of researches lie in, works with Abadi’s clipping and focuses
 91 on better design of R . To name a few examples, one can adaptively design R_t for each iteration t
 92 [3, 51, 25], or design the per-layer clipping threshold vector $\mathbf{R} \in \mathbb{R}^L$ for L layers [1, 43] so as to
 93 apply a different clipping threshold for each layer.

94 Much fewer works fall into the second category that proposes new clipping method. In fact, any
 95 function $\text{Clip} : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying $\|\text{Clip}(\mathbf{g}) \cdot \mathbf{g}\| \leq R$ can serve as a valid clipping function besides
 96 Abadi’s. For instance, the global clipping [9] proposes $\text{Clip}_{\text{global}}(\mathbf{g}) := \mathbb{I}(\|\mathbf{g}\| < Z) \cdot R/Z$ to correct
 97 the bias in the private gradient and alleviate the mis-calibration issue of DP classifiers. Our automatic
 98 clipping also belongs to this category. We note that different clipping methods work orthogonally to
 99 optimizers, network architectures and gradient norm computation (see Section 7).

100 3 Motivation

101 3.1 Small clipping threshold works best

102 One intriguing observation that we can make about the re-
 103 cent studies on DP learning with large models is that the
 104 state-of-the-art (SOTA) results are often achieved with very
 105 small clipping threshold R . This observation is consistent
 106 in both vision and language tasks. In [37], GPT2 (over 800
 107 million parameters) and RoBERTa models (over 400 mil-
 108 lions parameters) achieve the best results under DP on QNLI,
 109 MNLI, SST-2, QQP, E2E, and DART datasets, with each per-
 110 sample gradient clipped to length $R = 0.1$. In [34, 14, 44],
 111 ResNets and Vision Transformers achieve the best DP re-
 112 sults on ImageNet with $R = 1$; in [60], the best DP results
 113 on CIFAR10 use $R = 0.1$ with ResNeXt-29 and SimCLRv2
 114 [12]. The effectiveness of small clipping threshold together
 115 with proper learning rate is depicted in Figure 1.

116 Intuitively, smaller R implies that the Abadi’s clipping (3.1)
 117 happens, which means $\min(R/\|\mathbf{g}_i\|, 1) = R/\|\mathbf{g}_i\|$. Given
 118 that the clipping threshold R is so small compared to the
 119 number of parameters in large neural networks, and that
 120 strong DP is guaranteed when the number of training iter-
 121 ations is small (i.e. $\|\mathbf{g}_i\|$ has not converged to small values
 122 yet), we expect and empirically observe that the clipping
 123 happens on a large proportion of per-sample gradients at
 124 all iterations. For instance, we find in the GPT2 genera-
 125 tion experiments in [37] that 100% of per-sample gradi-
 126 ents are clipped at all iterations; in classification tasks such
 127 as QQP/QNLI/MNLI, the percentage of clipping is about
 128 20 ~ 60% on average (more details in Appendix H.1).

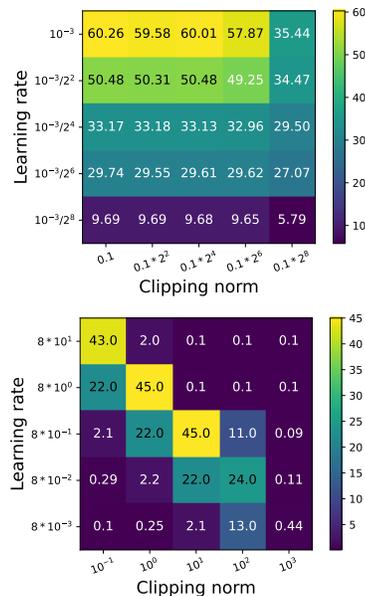


Figure 1: Ablation study of clipping threshold and learning rate that achieves SOTA results. Upper: BLEU score of GPT2 on E2E dataset, adapted from [37], trained with DP-AdamW. Lower: Test accuracy of ResNet18 on ImageNet dataset, adapted from [34], trained with DP-SGD with momentum.

129 **3.2 Per-sample gradient normalization as new clipping: a similarity viewpoint**

130 In the small clipping threshold regime, we can approximately view

$$\text{Clip}_{\text{Abadi}}(\mathbf{g}_i; R) = \min(R/\|\mathbf{g}_i\|, 1) \approx R/\|\mathbf{g}_i\| =: \text{Clip}_{\text{AUTO-V}}(\mathbf{g}_i; R) \quad (3.1)$$

131 and thus derive a novel private gradient $\sum_i R \frac{\mathbf{g}_i}{\|\mathbf{g}_i\|} + \sigma R \cdot \mathcal{N}(0, \mathbf{I})$. Here AUTO-V stands for the
 132 vanilla automatic clipping, which essentially performs the gradient normalization on each per-sample
 133 gradient. As a specific example, we can write the R -dependent automatic DP-SGD as

$$R\text{-dependent DP-SGD}_{\text{AUTO-V}} : \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left(\sum_{i \in B_t} R \frac{\partial l_i}{\partial \mathbf{w}_t} / \|\frac{\partial l_i}{\partial \mathbf{w}_t}\| + \sigma R \cdot \mathcal{N}(0, \mathbf{I}) \right) \quad (3.2)$$

134 We may view our AUTO-V clipping as to maximize the dot-product similarity between the clipped
 135 gradient and the regular gradient, a commonly used similarity measure, e.g. in the attention block in
 136 transformers [61]. Suppose we want

$$\max_{C_i} \left\langle \sum_i C_i \mathbf{g}_i, \sum_j \mathbf{g}_j \right\rangle \quad \text{s.t. } 0 \leq C_i \leq R/\|\mathbf{g}_i\|$$

137 Note that the constraint is a sufficient condition for clipping, as discussed in Section 2.2. It is not
 138 hard to see that the optimal clipping factor is

$$C_i = \begin{cases} R/\|\mathbf{g}_i\| & \text{if } \langle \mathbf{g}_i, \sum_j \mathbf{g}_j \rangle > 0 \\ 0 & \text{if } \langle \mathbf{g}_i, \sum_j \mathbf{g}_j \rangle \leq 0 \end{cases}$$

139 If the per-sample gradients are indeed concentrated in the sense $\forall i, \langle \mathbf{g}_i, \sum_j \mathbf{g}_j \rangle \geq 0$, then AUTO-V
 140 is the optimal per-sample gradient clipping. We compare with Abadi’s clipping in Figure 2, where
 the dot-product similarity is significantly magnified by our AUTO-V clipping.

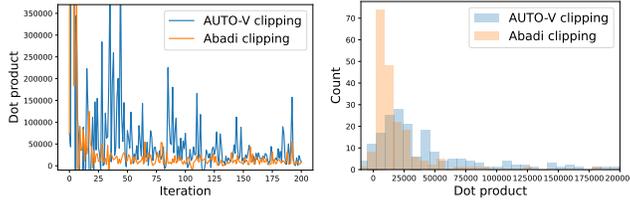


Figure 2: RoBERTa-base with DP-Adam ($\epsilon = 3$) on SST2 dataset, as in Section 6.2.

142 **3.3 Stability constant breaks scale-invariance and remains stationary**

143 One potential drawback of AUTO-V clipping is that all
 144 gradients lose their magnitudes information completely,
 145 since $\|\mathbf{g}_i \cdot \text{Clip}_{\text{AUTO-V}}(\mathbf{g}_i; R)\| = R, \forall i$. This scale-
 146 invariance in AUTO-V and partially in Abadi’s clipping
 147 (when $\|\mathbf{g}_i\| > R$) leads to the "lazy region" issue: the
 148 parameters will not be updated by DP-GD even if the
 149 true gradients are non-zero. In Figure 3, we illustrate
 150 in a logistic regression² that AUTO-V and Abadi’s clip-
 151 ping have zero clipped gradient for the trainable parameter
 152 $\theta \in [-2, 2]$, as the per-sample gradients from two classes cancel each other.

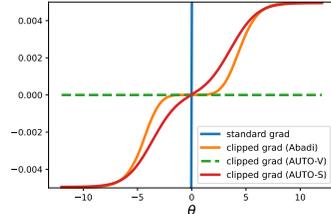


Figure 3: Gradient (scalar) at each θ .

153 Another benefit of γ is to remain stationary as $\mathbf{g}_i \rightarrow 0$, i.e. making the clipped gradient $C_i \mathbf{g}_i \rightarrow \mathbf{g}_i/\gamma$
 154 small rather than having a magnitude R in AUTO-V. We elaborate this point in Section 4.3.

155 To preserve the magnitude information and thus escape the lazy region, we propose the AUTO-S
 156 clipping, with a positive stability constant γ :

$$\text{Clip}_{\text{AUTO-S}}(\mathbf{g}_i; R) := R/(\|\mathbf{g}_i\| + \gamma) \quad (3.3)$$

157 We visualize in Figure 4 that AUTO-S allows larger per-sample gradients to have larger magnitudes
 158 after the clipping, while still allowing smaller gradients to vanish after “clipping”. This is critical in
 159 our convergence analysis and allows DP-SGD_{AUTO-S} (but not DP-SGD_{AUTO-V}) to converge to zero
 160 gradient norms in Section 5.

²The settings are in Appendix F, where the lazy region issues also emerge in the mean estimation problem. We note that the lazy region is also discussed in [13, Example 2].

161 4 Automatic DP Training

162 One may wonder why our clipping (3.1)(3.3) is automatic at all, if the hyperparameter R is still
 163 present and there is an additional parameter γ to choose. It turns out that any constant choice of $R > 0$
 164 is equivalent to choosing $R = 1$, and common deep learning optimizers are insensitive to the choice
 165 of γ (e.g. for any $\gamma > 0$, we show that the gradient norm converges to zero at the same asymptotic
 166 rate in Theorem 4; see also the ablation study in Figure 14). Consequently, we set $\gamma = 0.01$ as the
 167 default. Specifically, let us redefine the R -independent clipping function:

$$\text{Clip}_{\text{AUTO-S}}(\mathbf{g}_i) := 1/(\|\mathbf{g}_i\| + \gamma). \quad (4.1)$$

168 With this clipping, we can design automatic DP optimizers similar to (1.1):

$$\begin{aligned} \text{Automatic DP Optimizer}(\{\mathbf{g}_i\}_{i=1}^B) &= \text{Optimizer}(\hat{\mathbf{g}}_t) \\ \text{where } \hat{\mathbf{g}}_t &:= \sum_{i \in B_t} \frac{\mathbf{g}_{t,i}}{\|\mathbf{g}_{t,i}\| + \gamma} + \sigma \cdot \mathcal{N}(0, \mathbf{I}) \end{aligned} \quad (4.2)$$

169 Clearly, the new private gradient $\hat{\mathbf{g}}_t$ from our automatic clipping is R -independent, in contrast to the
 170 one used in (1.1). A concrete example (in the case of $\gamma = 0$) that is comparable to (3.2) will be

$$R\text{-independent DP-SGD}_{\text{AUTO-V}}: \quad \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \left(\sum_{i \in B_t} \frac{\partial l_i}{\partial \mathbf{w}_t} / \left\| \frac{\partial l_i}{\partial \mathbf{w}_t} \right\| + \sigma \cdot \mathcal{N}(0, \mathbf{I}) \right) \quad (4.3)$$

171 Leveraging the private gradient $\hat{\mathbf{g}}_t$ in (4.2), we can train DP neural networks without tuning DP-
 172 specific hyperparameters R and σ , as demonstrated in Algorithm 1.

Algorithm 1 Automatic Deep Learning with DP

Parameters: initial weights \mathbf{w}_0 , learning rate η_t , sampling probability p , number of iterations T .

- 1: Find σ such that $\epsilon_{\text{Accountant}}(\delta, \sigma, p, T) \leq \epsilon$ from any privacy accountant.
 - 173 2: **for** iteration $t = 1, \dots, T$ **do**
 - 3: Sample a batch B_t by including each data point iid with probability p ($\mathbb{E}[\text{BatchSize}] = pn$).
 - 4: Apply automatic clipping to per-sample gradients $\{\mathbf{g}_i\}_{i \in B_t}$: $\hat{\mathbf{g}}_i = \mathbf{g}_i / (\|\mathbf{g}_i\|_2 + 0.01)$.
 - 5: Add Gaussian noise to the sum of clipped gradients: $\hat{\mathbf{g}} = \sum_i \hat{\mathbf{g}}_i + \sigma \cdot \mathcal{N}(0, \mathbf{I})$.
 - 6: Update \mathbf{w}_t by any optimizer on the private gradient $\hat{\mathbf{g}}$ with learning rate η_t .
-

We will elaborate two distinct reasons in each section for the following statement:

$$\boxed{\text{DP Optimizer}_{\text{Abadi}} \approx R\text{-dependent DP Optimizer}_{\text{AUTO}} \equiv R\text{-independent DP Optimizer}_{\text{AUTO}}}$$

174 which reduces the hyperparameter tuning of DP training to that of the regular training, i.e. only on
 175 learning rate, weight decay, etc. The significant save in the tuning effort is illustrated in Figure 15.

176 4.1 Non-adaptive optimizer couples clipping threshold with learning rate

With R -dependent automatic clipping, DP-SGD becomes

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

177 We can view $\eta_{\text{effective}} \equiv \eta R$ as a whole: increasing R has the same effect as increasing η , which
 178 explains the diagonal pattern in Figure 1(lower plot) where DP-SGD_{Abadi} is applied with small
 179 clipping threshold³. We extend to general non-adaptive optimizers in Theorem 1, with proof in
 180 Appendix B.1⁴.

181 **Theorem 1.** *Non-adaptive R -dependent automatic DP optimizers (including SGD, Heavyball[52]*
 182 *and NAG[47]), with learning rate η and weight decay λ , is equivalent to R -independent automatic*
 183 *DP optimizers, with learning rate $\eta' = \eta R$ and weight decay $\lambda' = \lambda/R$.*

³When we further consider weight decay in automatic clipping (included in Theorem 1), increasing R is no longer equivalent to increasing η , as η also couples with the weight decay constant λ .

⁴This coupling of η and R is also partially observed in [14] through a reparameterization trick of Abadi's clipping. Unlike AUTO-S/V, their coupling is not strict (e.g. doubling R is not equivalent to doubling η in their Figure 8, thus necessitating tuning both (η, R)), and the relationship to weight decay was not discussed.

184 **4.2 Adaptive optimizer can be insensitive to clipping threshold**

Adaptive automatic DP optimizers are different than the non-adaptive ones, as the clipping threshold cancels out instead of being coupled with learning rate. To see this, we scrutinize DP-Adam_{Abadi} (which is similar to DP-Adam_{AUTO-V}) in Figure 1 (upper plot), where columns to the left are almost identical. Further evidence is observed in [44, Table 5] that shrinking R has zero effect on LAMB. We now give a simple explanation using AdaGrad [19]:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{\mathbf{g}_t}{\sqrt{G_t}}$$

where $\mathbf{g}_t = \sum_i \mathbf{g}_{t,i}$ is the gradient sum and $G_t = \sum_{\tau < t} \mathbf{g}_\tau^2$ is sum of gradient square by Hadamard product over the past iterations. In R -dependent DP-AdaGrad_{AUTO-V}, the private gradient is $R\hat{\mathbf{g}}_t$ in place of the standard gradient sum \mathbf{g}_t , and $\hat{G}_t = R^2 \sum_{\tau < t} \hat{\mathbf{g}}_\tau^2$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{R\hat{\mathbf{g}}_t}{\sqrt{\hat{G}_t}} = \mathbf{w}_t - \eta \frac{\hat{\mathbf{g}}_t}{\sqrt{\sum_{\tau < t} (\hat{\mathbf{g}}_\tau)^2}}.$$

185 We generalize to the general adaptive optimizers in Theorem 2, with proof in Appendix B.2.

186 **Theorem 2.** Adaptive R -dependent automatic DP optimizers (including AdaGrad[19], AdaDelta[69],
187 AdaMax/Adam[31], NAdam[17], RAdam[39], LARS[65], LAMB[66]), with learning rate η and
188 weight decay λ is equivalent to R -independent automatic DP optimizers with learning rate η and
189 weight decay $\lambda' = \lambda/R$. With decoupled weight decay[42], R -dependent automatic DP-AdamW is
190 equivalent to R -independent automatic DP-AdamW with the same η and λ .

191 Similarly, we demonstrate the automatic DP optimizers with per-layer clipping style in Appendix B.3.

192 **4.3 Automatic clipping guarantees the same level of privacy while maximizes utility**

193 In Theorem 3 (proved in Appendix A), we show that
194 the new private gradient $\hat{\mathbf{g}}_t$ in (4.2) has the same level of
195 privacy guarantee as the existing one in (1.1), since
196 the global sensitivity remains the same (see Figure 4).
197 We note that as long as $\gamma > 0$, the magnitude information
198 of per-sample gradients is preserved by AUTO-
199 S, in the sense that $\|\mathbf{g}_i\| > \|\mathbf{g}_j\| \iff \|C_i \mathbf{g}_i\| >$
200 $\|C_j \mathbf{g}_j\|$, whereas this can be violated in both the
201 AUTO-V and Abadi’s clipping (as depicted by the
202 flat curve in Figure 4 when $\|\mathbf{g}_i\| > 1$).

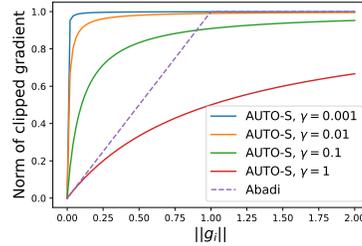


Figure 4: Gradient norms before and after being clipped by different methods at $R = 1$.

203 Additionally, note that when γ is small, almost all data
204 points “max out” the signal relative to the amount of noise we add. To say it differently, for the same
205 amount of noise, AUTO-S with small γ allows more signal to be pushed through a differentially
206 private channel. Towards the end of the training, i.e., at the limit when $\|\mathbf{g}_i\| \rightarrow 0$ for all i , then we
207 have $\sum_i \frac{\mathbf{g}_i}{\|\mathbf{g}_i\| + \gamma} \rightarrow \frac{1}{\gamma} \sum_i \mathbf{g}_i$. In words, the clipped gradients become closer to the standard SGD,
208 thus do not suffer from the instability of AUTO-V.

209 **Theorem 3.** Under the noise multiplier σ , number of iterations T , subsampling probability B/n ,
210 DP optimizers using AUTO-V or AUTO-S clipping satisfy $(\epsilon_{\text{Accountant}}(\delta, \sigma, B/n, T), \delta)$ -DP, where
211 $\epsilon_{\text{Accountant}}$ is any valid privacy accountant for DP-SGD under Abadi’s clipping.

212 **5 Convergence analysis of DP-SGD with automatic clipping**

213 **5.1 Convergence theory of DP-SGD to stationary points**

214 We highlight that automatic clipping can be more amenable to analysis than Abadi’s clipping in [13],
215 since we no longer need to decide whether each per-sample gradient is clipped.

216 To analyze the convergence of automatic DP-SGD (4.2) in the non-convex setting, we follow the
217 standard assumptions in the SGD literature [24, 2, 6], with one additional symmetry assumption on
218 the gradient noise.

219 **Assumption 5.1** (Lower bound of loss). For all \mathbf{w} and some constant \mathcal{L}_* , we have $\mathcal{L}(\mathbf{w}) \geq \mathcal{L}_*$.

220 **Assumption 5.2** (Smoothness). Let $\mathbf{g}(\mathbf{w})$ denote the gradient of the objective $\mathcal{L}(\mathbf{w})$. Then $\forall \mathbf{w}, \mathbf{v}$,
221 there is an non-negative constant L such that

$$\mathcal{L}(\mathbf{v}) - [\mathcal{L}(\mathbf{w}) + \mathbf{g}(\mathbf{w})^\top (\mathbf{v} - \mathbf{w})] \leq \frac{L}{2} \|\mathbf{w} - \mathbf{v}\|^2. \quad (5.1)$$

Assumption 5.3 (Gradient noise). The per-sample gradient noise $\tilde{\mathbf{g}}_{t,i} - \mathbf{g}_t$ is i.i.d. from some distribution such that

$$\mathbb{E}(\tilde{\mathbf{g}}_{t,i} - \mathbf{g}_t) = 0, \mathbb{E}\|\tilde{\mathbf{g}}_{t,i} - \mathbf{g}_t\|^2 \leq \xi^2,$$

and $\tilde{\mathbf{g}}_{t,i}$ is centrally symmetric⁵ about \mathbf{g}_t in distribution:

$$\tilde{\mathbf{g}}_{t,i} - \mathbf{g}_t \stackrel{\mathcal{D}}{=} \mathbf{g}_t - \tilde{\mathbf{g}}_{t,i}.$$

222 We show in Theorem 4 that DP-SGD with AUTO-S clipping allows the true gradient norm to converge
223 to zero, but not so with AUTO-V clipping. We leave the proof in Appendix C.1.

224 **Theorem 4.** Under Assumption 5.1, 5.2, 5.3, running DP-SGD with automatic clipping for T
225 iterations and setting the learning rate $\eta \propto 1/\sqrt{T}$ give

$$\min_{0 \leq t \leq T} \mathbb{E}(\|\mathbf{g}_t\|) \leq \mathcal{G} \left(\frac{4}{\sqrt{T}} \sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L \left(1 + \frac{\sigma^2 d}{B^2}\right)}; \xi, \gamma \right) := \min_{r > 0} \frac{\xi}{r} + \mathcal{F}(\dots; r, \xi, \gamma). \quad (5.2)$$

226 Here \dots represents the first argument of \mathcal{G} , and \mathcal{G} is increasing and positive. As $T \rightarrow \infty$, we have
227 $\min_t \mathbb{E}(\|\mathbf{g}_t\|) = O(T^{-1/4})$ for AUTO-S, the same rate as the standard SGD given in Theorem 9.

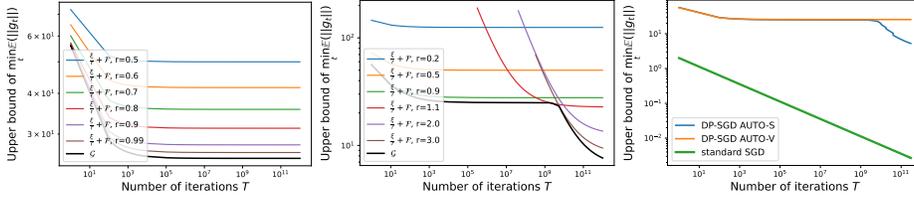


Figure 5: Left: DP-SGD with AUTO-V clipping. Middle: DP-SGD with AUTO-S clipping. Right: Log-log plot of convergence rate in comparison to standard SGD. Here $\xi = 25, \gamma = 0.01$, and the $O(1/\sqrt{T})$ term is set to 10 for DP-SGD and to 2 for standard SGD.

228 **Remark 5.4.** In Theorem 4, the upper bound takes an implicit form of $\mathcal{G}(\cdot; \xi, \gamma)$ because it is a
229 lower envelope of functions $\frac{\xi}{r} + \mathcal{F}(\cdot; r, \xi, \gamma)$ over all possible $r > 0$, whose forms are detailed in
230 Theorem 6. Notice that \mathcal{G} results only from the clipping operation, not from the noise addition.

231 **Remark 5.5.** We show in Theorem 6 and demonstrate in Figure 5 that the upper bound (5.2) is always
232 larger than ξ with AUTO-V ($\gamma = 0$), and can only be reduced to zero with AUTO-S ($\gamma > 0$). We
233 provide real data evidence in Figure 13 that strictly positive γ reduces the gradient norm significantly.

234 5.2 Analysis of factors affecting the convergence

235 We now analyze the many factors that affect the convergence in Theorem 4, from a unified viewpoint
236 of both the convergence and the privacy.

237 We start with the stability constant γ and the learning rate η_t , both only affect the convergence not
238 the privacy. We empirically observe in Figure 7 that small γ benefits the convergence at initial
239 iterations (when the privacy guarantee is strong) but larger γ converges faster asymptotically. For η_t ,
240 the optimal is in fact the minimizer of the hyperbola in (C.4), that is unique and tunable.

241 Next, we focus on the hyperparameters that affect both convergence and privacy: the batch size B ,
242 the noise multiplier σ , and the number of iterations T . These hyperparameters have to be considered
243 along the privacy-accuracy tradeoff, not just from a convergence perspective.

244 Recall that given a fixed privacy budget (ϵ, δ) , we rely on modern privacy accountant for computing
245 the appropriate combinations of parameter σ, T, B . The exact expression of the bound as a function of
246 (ϵ, δ) is somewhat messy. For this reason, we illustrate our analysis in terms of the surrogate parameter
247 μ for μ -GDP [16]. [7] showed that DP-SGD's privacy guarantee asymptotically converges to μ -GDP
248 (as $T \rightarrow \infty$) with $\mu = \frac{B}{n} \sqrt{T}(e^{1/\sigma^2} - 1)$. μ -GDP implies (ϵ, δ) -DP with $\epsilon = \mu^2 + \mu\sqrt{2 \log(1/\delta)}$ ⁶.
249 We can alternatively leverage ρ -tCDP [10] for similar conclusions, using ρ in place of μ^2 in (5.3).

⁵The symmetry assumption has been empirically verified in [13, Figure 3]. For theoretical analysis, it can be extended to mirror symmetry about the hyperplane normal to \mathbf{g}_t , that is $\{\mathbf{v} : \mathbf{g}_t^\top \mathbf{v} = 0\}$.

⁶More precisely, μ -GDP is equivalent to an entire family of (ϵ, δ) -DP for any $\epsilon > 0$ and $\delta = \Phi(\mu/2 - \epsilon/\mu) - e^\epsilon \Phi(-\mu/2 - \epsilon/\mu)$ where Φ is the standard Gaussian CDF.

250 **Theorem 5.** Under Assumption 5.1, 5.2, 5.3, fixing the asymptotic $\mu(\epsilon, \delta)$ -GDP parameter, running
 251 DP-SGD with automatic clipping for T iterations and setting the learning rate $\eta \propto 1/\sqrt{T}$ give

$$\min_{0 \leq t \leq T} \mathbb{E}(\|g_t\|) \leq \mathcal{G} \left(4\sqrt{(\mathcal{L}_0 - \mathcal{L}_*)L} \left(\frac{1}{T} + \frac{d}{\mu^2 n^2} + O\left(\frac{1}{B^2 T}\right) \right); \xi, \gamma \right) \quad (5.3)$$

252 To show that our analysis matches the training behaviors observed in SOTA empirical work [37, 34,
 253 14, 60, 44, 68], we minimize the first argument of \mathcal{G} in (5.3), denoted as $X(B, T, \mu, d, L, \mathcal{L}_0)$.

- 254 1. **[Train longer with larger noise]** Fixing the expected batch size B , we see that X is decreasing
 255 in T . Hence larger T and consequently larger σ are preferred.
- 256 2. **[Larger batch size helps]** Fixing number of iterations T or epochs $E = BT/n$, we see that X is
 257 decreasing in B . Hence larger B and consequently larger σ are preferred.
- 258 3. **[Pretraining is critical]** Pretraining can boost the accuracy of DP learning through a much smaller
 259 initial loss \mathcal{L}_0 and from a smooth (small L) and flat (small ξ , c.f. Figure 7(left)) initialization.
- 260 4. **[Learning rate needs tuning]** The optimal learning rate by minimizing (C.4) is $\sqrt{\frac{(\mathcal{L}_0 - \mathcal{L}_*)\mu^2 n^2}{L(\mu^2 n^2 + dT)}}$.
 261 This indicates that one should use larger learning rate for smaller model d , weaker privacy (larger
 262 μ or small ϵ), or smaller iteration budget T . Interestingly, the optimal choice of learning rate is
 263 independent to (expected) batch-size B .

264 6 Experiments

265 We evaluate our automatic DP training on image classification, sentence classification, and table-to-
 266 text generation tasks. Detailed settings including hyperparameters can be found in Appendix G.

267 6.1 Image classification

268 For MNIST/FashionMNIST, we use the same setup as in [49, 60, 57] with a simple CNN. For
 269 CIFAR10, we use the same setup as in [60] with pretrained SimCLRv2 [12]. For ImageNette, a
 270 10-class sub-task of ImageNet [15], we use the same setup as in [32] without the learning rate decay.
 271 For CelebA [41], the real human face dataset, we train ResNet9 [28] with group normalization to
 272 replace the batch normalization. Notice that CelebA contains high-resolution (178x218) images, each
 273 with 40 labels. We consider CelebA for either multi-class classification on one label, e.g. ‘Smiling’
 274 and ‘Male’, or for multi-label/multi-task problem to learn all labels simultaneously.

Task	Model	(ϵ, δ)	Accuracy %		
			Abadi’s clipping	AUTO-S clipping	non-DP ($\epsilon = \infty$)
MNIST	4-layer CNN	(3, 1e-5)	98.04 ± 0.09	98.15 ± 0.07	99.11 ± 0.07
FashionMNIST	4-layer CNN	(3, 1e-5)	86.04 ± 0.26	86.36 ± 0.18	89.57 ± 0.13
CIFAR10 pretrained	SimCLRv2	(2, 1e-5)	92.44 ± 0.13	92.70 ± 0.02	94.42 ± 0.01
ImageNette	ResNet9	(8, 1e-4)	60.29 ± 0.53	60.71 ± 0.48	71.11 ± 0.37
CelebA [Smiling]	ResNet9	(8, 5e-6)	90.75 ± 0.11	91.08 ± 0.08	92.61 ± 0.20
CelebA [Male]	ResNet9	(8, 5e-6)	95.54 ± 0.14	95.70 ± 0.07	97.90 ± 0.04
CelebA Multi-label	ResNet9	(3, 5e-6)	86.81 ± 0.03	87.05 ± 0.01	90.30 ± 0.02
CelebA Multi-label	ResNet9	(8, 5e-6)	87.52 ± 0.15	87.58 ± 0.04	90.30 ± 0.02

Table 1: Average test accuracy and 95% confidence interval on image tasks over 5 runs.

275 In Table 1, we observe that AUTO-S clipping outperforms existing clipping in all datasets with statisti-
 276 cal significance. Interestingly, the standard deviation from different runs is smaller for automatic
 277 DP optimizers, indicating better reproducibility and stability. We additionally experiment 40 binary
 278 classification problems on CelebA with respect to each label, and observe that the mean accuracy
 279 further improves to 91.63% at $\epsilon = 8$ for AUTO-S (see Appendix J).

280 6.2 Sentence classification

281 On five benchmark language datasets (MNLI(m/mm)[64], QQP[30], QNLI[55], SST2[59]), we
 282 compare our automatic DP training with reparameterized gradient perturbation (RGP, [68]) and
 283 full-parameter finetuning (full, [37]) using RoBERTa models [40]. These methods use the same
 284 experimental setup. For language models, our automatic training is based on the codebase of [37]⁷.

⁷See <https://github.com/lxuechen/private-transformers> and the detailed modification in Appendix K.3.

Method	$\epsilon = 3$				$\epsilon = 8$				$\epsilon = \infty$ (non-DP)			
	MNLI	QQP	QNLI	SST2	MNLI	QQP	QNLI	SST2	MNLI	QQP	QNLI	SST2
RGP [68]	-	-	-	-	80.5/79.6	85.5	87.2	91.6	83.6/83.2	89.3	91.3	92.9
full [37]	82.45/82.99	85.56	87.42	91.86	83.20/83.46	86.08	87.94	92.09				
full AUTO-V	81.21/82.03	84.72	86.56	91.86	82.18/82.64	86.23	87.24	92.09	85.91/86.14	87.34	91.40	94.49
full AUTO-S	83.22/83.21	85.76	86.91	92.32	83.82/83.55	86.58	87.85	92.43				

Table 2: Test accuracy on language tasks with RoBERTa-base (12 blocks, 163 million parameters).

Method	$\epsilon = 3$				$\epsilon = 8$				$\epsilon = \infty$ (non-DP)			
	MNLI	QQP	QNLI	SST2	MNLI	QQP	QNLI	SST2	MNLI	QQP	QNLI	SST2
RGP [68]	-	-	-	-	86.1/86.0	86.7	90.0	93.0	-	-	-	-
full [37]	86.43 /86.46	86.43	90.76	93.04	87.02/ 87.26	87.47	91.10	93.81				
full AUTO-V	85.33/85.61	86.61	89.99	93.12	85.91/86.10	86.86	90.55	93.35	90.33/90.03	87.90	93.61	96.21
full AUTO-S	86.27/ 86.67	86.76	91.01	93.92	87.07 /87.16	87.47	91.45	94.61				

Table 3: Test accuracy on language tasks with RoBERTa-large (24 blocks, 407 million parameters).

285 In Table 2 and Table 3, we note that full parameter finetuning with AUTO-S outperforms or at least
 286 matches SOTA on all tasks. We use *exactly the same* hyperparameters as in [37].

287 6.3 Table-to-text generation

288 We compare our automatic DP training with a variety of fine-tuning methods, for table-to-text
 289 generation task on E2E dataset [20], where the goal is to generate texts about different aspects of a
 290 restaurant’s data. We measure the success on this task by BLEU, ROUGE-L (in Table 4), METEOR,
 291 NIST, CIDEr (extended in Table 7), with higher value meaning better model quality.

Metric	DP guarantee	GPT2 large	GPT2 medium	GPT2							
		full	full	full	full	full	LoRA	RGP	prefix	top2	retrain
		AUTO-S	AUTO-S	AUTO-S	AUTO-V	[37]	[29]	[68]	[36]	[37]	[37]
BLEU	$\epsilon = 3$	64.180	63.850	61.340	61.519	61.519	58.153	58.482	47.772	25.920	15.457
	$\epsilon = 8$	64.640	64.220	63.600	63.189	63.189	63.389	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	$\epsilon = 3$	67.857	67.071	65.872	65.670	65.670	65.773	65.560	58.964	44.536	35.240
	$\epsilon = 8$	68.968	67.533	67.073	66.429	66.429	67.525	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

Table 4: Test performance on E2E dataset with GPT2. Additional performance measures are included in Table 7. The best two GPT2 models for each row are marked in bold.

292 Competitive methods include low-rank adaption (LoRA), prefix-tuning (prefix), RGP, only fine-tuning
 293 the top 2 Transformer blocks (top2), and training from scratch (retrain), as were recorded in [37].
 294 Again, we use the *exactly the same* hyperparameters as in [37]. For GPT2 (163 million parameters),
 295 GPT2 medium (406 million), and GPT2 large (838 million), Table 4 shows that AUTO-S is scalable
 296 with stronger performance on larger models. Our automatic full-parameter finetuning has the best
 297 overall performance. Additionally, we highlight that AUTO-S and methods like LoRA are not
 298 mutually exclusive and can be combined to yield strong performance, since AUTO-S modifies the
 299 optimizers and LoRA modifies the architecture.

300 7 Discussion

301 In this work, we proposed AutoClipping as a drop-in replacement to the standard per-example
 302 clipping differentially private training. This is the first technique that eliminate the need to tune the
 303 clipping threshold R , thus making DP deep learning as easy as regular learning. Our AUTO-S method
 304 enjoys both theoretical guarantee of convergence in non-convex problems (under various conditions),
 305 and strong empirical performance that advances the state-of-the-art (SOTA) of DP learning on both
 306 computer vision and language tasks.

307 We are excited about the future of automatic DP training, especially along with other working
 308 techniques. Notably, our automatic clipping applies compatibly with general optimizers (e.g. [8, 18]),
 309 clipping styles (all-layer or per-layer), architecture modifications (e.g. LoRA, RGP, prefix), and data
 310 augmentation (e.g. adversarial training [26] and multiple augmentation [14]). Thus, we expect to
 311 achieve comparable results to all SOTA in a lightweight fashion.

312 References

- 313 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,
314 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*
315 *conference on computer and communications security*, pages 308–318, 2016.
- 316 [2] Zeyuan Allen-Zhu. Natasha 2: Faster non-convex optimization than sgd. *Advances in neural*
317 *information processing systems*, 31, 2018.
- 318 [3] Galen Andrew, Om Thakkar, Brendan McMahan, and Swaroop Ramaswamy. Differentially
319 private learning with adaptive clipping. *Advances in Neural Information Processing Systems*,
320 34, 2021.
- 321 [4] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with
322 improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic*
323 *and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages
324 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics.
- 325 [5] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar.
326 signsgd: Compressed optimisation for non-convex problems. In *International Conference on*
327 *Machine Learning*, pages 560–569. PMLR, 2018.
- 328 [6] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine
329 learning. *Siam Review*, 60(2):223–311, 2018.
- 330 [7] Zhiqi Bu, Jinshuo Dong, Qi Long, and Weijie J Su. Deep learning with gaussian differential
331 privacy. *Harvard data science review*, 2020(23), 2020.
- 332 [8] Zhiqi Bu, Sivakanth Gopi, Janardhan Kulkarni, Yin Tat Lee, Hanwen Shen, and Uthaiapon
333 Tantipongpipat. Fast and memory efficient differentially private-sgd via jl projections. *Advances*
334 *in Neural Information Processing Systems*, 34, 2021.
- 335 [9] Zhiqi Bu, Hua Wang, and Qi Long. On the convergence and calibration of deep learning with
336 differential privacy. *arXiv preprint arXiv:2106.07830*, 2021.
- 337 [10] Mark Bun, Cynthia Dwork, Guy N Rothblum, and Thomas Steinke. Composable and versatile
338 privacy via truncated cdp. In *Proceedings of the 50th Annual ACM SIGACT Symposium on*
339 *Theory of Computing*, pages 74–86, 2018.
- 340 [11] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine
341 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training
342 data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*,
343 pages 2633–2650, 2021.
- 344 [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
345 for contrastive learning of visual representations. In *International conference on machine*
346 *learning*, pages 1597–1607. PMLR, 2020.
- 347 [13] Xiangyi Chen, Steven Z Wu, and Mingyi Hong. Understanding gradient clipping in private sgd:
348 A geometric perspective. *Advances in Neural Information Processing Systems*, 33:13773–13782,
349 2020.
- 350 [14] Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlocking
351 high-accuracy differentially private image classification through scale. *arXiv preprint*
352 *arXiv:2204.13650*, 2022.
- 353 [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
354 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
355 *recognition*, pages 248–255. Ieee, 2009.
- 356 [16] Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *Journal of the Royal*
357 *Statistical Society Series B*, 84(1):3–37, 2022.
- 358 [17] Timothy Dozat. Incorporating nesterov momentum into adam. 2016.
- 359 [18] Jian Du and Haitao Mi. Dp-fp: Differentially private forward propagation for large models.
360 *arXiv preprint arXiv:2112.14430*, 2021.
- 361 [19] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
362 and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.

- 363 [20] Ondrej Dusek, Jekaterina Novikova, and Verena Rieser. Evaluating the State-of-the-Art of
364 End-to-End Natural Language Generation: The E2E NLG Challenge. *Computer Speech &*
365 *Language*, 59:123–156, January 2020.
- 366 [21] Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory*
367 *and applications of models of computation*, pages 1–19. Springer, 2008.
- 368 [22] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to
369 sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284.
370 Springer, 2006.
- 371 [23] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found.*
372 *Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- 373 [24] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex
374 stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- 375 [25] Aditya Golatkar, Alessandro Achille, Yu-Xiang Wang, Aaron Roth, Michael Kearns, and
376 Stefano Soatto. Mixed differential privacy in computer vision. In *Proceedings of the IEEE/CVF*
377 *Conference on Computer Vision and Pattern Recognition*, pages 8376–8386, 2022.
- 378 [26] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
379 examples. In *International Conference on Learning Representations*, 2015.
- 380 [27] Sivakanth Gopi, Yin Tat Lee, and Lukas Wutschitz. Numerical composition of differential
381 privacy. *Advances in Neural Information Processing Systems*, 34, 2021.
- 382 [28] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
383 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
384 pages 770–778, 2016.
- 385 [29] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu
386 Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference*
387 *on Learning Representations*, 2021.
- 388 [30] Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. First quora dataset release: Question pairs,
389 2017.
- 390 [31] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International*
391 *Conference on Learning Representations*, 12 2014.
- 392 [32] Helena Klause, Alexander Ziller, Daniel Rueckert, Kerstin Hammernik, and Georgios Kaissis.
393 Differentially private training of residual networks with scale normalisation. *arXiv preprint*
394 *arXiv:2203.00324*, 2022.
- 395 [33] Antti Koskela, Joonas Jälkö, and Antti Honkela. Computing tight differential privacy guarantees
396 using fft. In *International Conference on Artificial Intelligence and Statistics*, pages 2560–2569.
397 PMLR, 2020.
- 398 [34] Alexey Kurakin, Steve Chien, Shuang Song, Roxana Geambasu, Andreas Terzis, and Abhradeep
399 Thakurta. Toward training at imagenet scale with differential privacy. *arXiv preprint*
400 *arXiv:2201.12328*, 2022.
- 401 [35] Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von
402 Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison,
403 Mario Sasko, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor
404 Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger,
405 Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault
406 Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. Datasets: A
407 community library for natural language processing. In *Proceedings of the 2021 Conference on*
408 *Empirical Methods in Natural Language Processing: System Demonstrations*, pages 175–184,
409 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational
410 Linguistics.
- 411 [36] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
412 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*
413 *and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long*
414 *Papers)*, pages 4582–4597, 2021.

- 415 [37] Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large language mod-
416 els can be strong differentially private learners. In *International Conference on Learning*
417 *Representations*, 2021.
- 418 [38] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summariza-*
419 *tion Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational
420 Linguistics.
- 421 [39] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and
422 Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International*
423 *Conference on Learning Representations*, 2019.
- 424 [40] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
425 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
426 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 427 [41] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
428 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 429 [42] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International*
430 *Conference on Learning Representations*, 2018.
- 431 [43] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially
432 private recurrent language models. In *International Conference on Learning Representations*,
433 2018.
- 434 [44] Harsh Mehta, Abhradeep Thakurta, Alexey Kurakin, and Ashok Cutkosky. Large scale transfer
435 learning for differentially private image classification. *arXiv preprint arXiv:2205.02973*, 2022.
- 436 [45] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations*
437 *symposium (CSF)*, pages 263–275. IEEE, 2017.
- 438 [46] Diganta Misra. Mish: A self regularized non-monotonic activation function. *BMVC 2020*, 2019.
- 439 [47] Yurii E Nesterov. A method for solving the convex programming problem with convergence
440 rate $o(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- 441 [48] Nicolas Papernot and Thomas Steinke. Hyperparameter tuning with renyi differential privacy.
442 In *International Conference on Learning Representations*, 2021.
- 443 [49] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Úlfar Erlingsson.
444 Tempered sigmoid activations for deep learning with differential privacy. In *Proceedings of the*
445 *AAAI Conference on Artificial Intelligence*, volume 35, pages 9312–9321, 2021.
- 446 [50] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic
447 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for*
448 *Computational Linguistics, ACL '02*, page 311–318, USA, 2002. Association for Computational
449 Linguistics.
- 450 [51] Venkatadheeraj Pichapati, Ananda Theertha Suresh, Felix X Yu, Sashank J Reddi, and Sanjiv
451 Kumar. Adacclip: Adaptive clipping for private sgd. *arXiv preprint arXiv:1908.07643*, 2019.
- 452 [52] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr*
453 *computational mathematics and mathematical physics*, 4(5):1–17, 1964.
- 454 [53] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging.
455 *SIAM journal on control and optimization*, 30(4):838–855, 1992.
- 456 [54] Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Micro-batch training with
457 batch-channel normalization and weight standardization. *arXiv preprint arXiv:1903.10520*,
458 2019.
- 459 [55] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+
460 questions for machine comprehension of text. In *EMNLP*, 2016.
- 461 [56] Seyed Omid Sadjadi, Timothee Kheyrkhah, Audrey Tong, Craig S Greenberg, Douglas A
462 Reynolds, Elliot Singer, Lisa P Mason, Jaime Hernandez-Cordero, et al. The 2017 nist language
463 recognition evaluation. In *Odyssey*, pages 82–89, 2018.
- 464 [57] Ali Shahin Shamsabadi and Nicolas Papernot. Losing less: A loss for differentially private deep
465 learning. 2021.

- 466 [58] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference
467 attacks against machine learning models. In *2017 IEEE symposium on security and privacy*
468 (*SP*), pages 3–18. IEEE, 2017.
- 469 [59] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng,
470 and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment
471 treebank. In *Proceedings of the 2013 conference on empirical methods in natural language*
472 *processing*, pages 1631–1642, 2013.
- 473 [60] Florian Tramer and Dan Boneh. Differentially private learning needs better features (or much
474 more data). In *International Conference on Learning Representations*, 2020.
- 475 [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
476 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
477 *processing systems*, 30, 2017.
- 478 [62] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
479 description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern*
480 *recognition*, pages 4566–4575, 2015.
- 481 [63] Yu-Xiang Wang, Borja Balle, and Shiva Prasad Kasiviswanathan. Subsampled rényi differ-
482 ential privacy and analytical moments accountant. In *International Conference on Artificial*
483 *Intelligence and Statistics*, pages 1226–1235. PMLR, 2019.
- 484 [64] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus
485 for sentence understanding through inference. In *Proceedings of the 2018 Conference of the*
486 *North American Chapter of the Association for Computational Linguistics: Human Language*
487 *Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational
488 Linguistics, 2018.
- 489 [65] Yang You, Igor Gitman, and Boris Ginsburg. Scaling sgd batch size to 32k for imagenet training.
490 *arXiv preprint arXiv:1708.03888*, 6(12):6, 2017.
- 491 [66] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan
492 Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep
493 learning: Training bert in 76 minutes. In *International Conference on Learning Representations*,
494 2020.
- 495 [67] Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad,
496 Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, Graham Cormode,
497 and Ilya Mironov. Opacus: User-friendly differential privacy library in PyTorch. *arXiv preprint*
498 *arXiv:2109.12298*, 2021.
- 499 [68] Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. Large scale private learning
500 via low-rank reparametrization. In *International Conference on Machine Learning*, pages
501 12208–12218. PMLR, 2021.
- 502 [69] Matthew D Zeiler. Adadelat: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*,
503 2012.
- 504 [70] Yuqing Zhu, Jinshuo Dong, and Yu-Xiang Wang. Optimal accounting of differential privacy
505 via characteristic function. In *International Conference on Artificial Intelligence and Statistics*,
506 pages 4782–4817. PMLR, 2022.

507 Checklist

- 508 1. For all authors...
- 509 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
510 contributions and scope? [Yes]
- 511 (b) Did you describe the limitations of your work? [N/A]
- 512 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 513 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?
514 [Yes]
- 515 2. If you are including theoretical results...

- 516 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assumptions
517 5.1-5.3.
- 518 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.
- 519 3. If you ran experiments...
- 520 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
521 results (either in the supplemental material or as a URL)? [Yes] We include our code in
522 Appendix K, and cite all datasets.
- 523 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
524 chosen)? [Yes] See all details in Appendix G.
- 525 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
526 multiple times)? [Yes] See Table 1.
- 527 (d) Did you include the total amount of compute and the type of resources used (e.g., type of
528 GPUs, internal cluster, or cloud provider)? [N/A] This work is unrelated to computational
529 efficiency.
- 530 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 531 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite all codebases,
532 datasets, models clearly.
- 533 (b) Did you mention the license of the assets? [N/A] Experiments are conducted on public data.
- 534 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We
535 include code snippet in Appendix K and will include a URL to release the full codebase.
- 536 (d) Did you discuss whether and how consent was obtained from people whose data you're
537 using/curating? [N/A]
- 538 (e) Did you discuss whether the data you are using/curating contains personally identifiable
539 information or offensive content? [N/A]
- 540 5. If you used crowdsourcing or conducted research with human subjects...
- 541 (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
542 [N/A]
- 543 (b) Did you describe any potential participant risks, with links to Institutional Review Board
544 (IRB) approvals, if applicable? [N/A]
- 545 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
546 participant compensation? [N/A]