

---

# DP-SGD Without Clipping: The Lipschitz Neural Network Way

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 State-of-the-art approaches for training Differentially Private (DP) Deep Neural  
2 Networks (DNN) faces difficulties to estimate tight bounds on the sensitivity of  
3 the network's layers, and instead rely on a process of per-sample gradient clipping.  
4 This clipping process not only biases the direction of gradients but also proves  
5 costly both in memory consumption and in computation. To provide sensitivity  
6 bounds and bypass the drawbacks of the clipping process, our theoretical analysis  
7 of Lipschitz constrained networks reveals an unexplored link between the Lipschitz  
8 constant with respect to their input and the one with respect to their parameters.  
9 By bounding the Lipschitz constant of each layer with respect to its parameters  
10 we guarantee DP training of these networks. This analysis not only allows the  
11 computation of the aforementioned sensitivities at scale but also provides leads  
12 on to how maximize the gradient-to-noise ratio for fixed privacy guarantees. To  
13 facilitate the application of Lipschitz networks and foster robust and certifiable  
14 learning under privacy guarantees, we provide a Python package that implements  
15 building blocks allowing the construction and private training of such networks.

## 16 1 Introduction

17 Machine learning relies more than ever on foundational models, and such practices raise questions  
18 about privacy. Differential privacy allows to develop methods for training models that preserve  
19 the privacy of individual data points in the training set. The field seeks to enable deep learning on  
20 sensitive data, while ensuring that models do not inadvertently memorize or reveal specific details  
21 about individual samples in their weights. This involves incorporating privacy-preserving mechanisms  
22 into the design of deep learning architectures and training algorithms, whose most popular example  
23 is Differentially Private Stochastic Gradient Descent (DP-SGD) [1]. One main drawback of classical  
24 DP-SGD methods is that they require costly per-sample backward processing and gradient clipping.  
25 In this paper, we offer a new method that unlocks fast differentially private training through the use  
26 of Lipschitz constrained neural networks. Additionally, this method offers new opportunities for  
27 practitioners that wish to easily "DP-fy" [2] the training procedure of a deep neural network.

28 **Differential privacy fundamentals.** Informally, differential privacy is a *definition* that quantifies how  
29 much the change of a single sample in a dataset affects the range of a stochastic function (here the DP  
30 training), called *mechanism* in this context. This quantity can be bounded in an inequality involving  
31 two parameters  $\epsilon$  and  $\delta$ . A mechanism fulfilling such inequality is said  $(\epsilon, \delta)$ -DP (see Definition 1).  
32 This definition is universally accepted as a strong guarantee against privacy leakages under various  
33 scenarii, including data aggregation or post-processing [3]. A popular rule of thumb suggests using  
34  $\epsilon \leq 10$  and  $\delta < \frac{1}{N}$  with  $N$  the number of records [2] for mild guarantees. In practice, most classic  
35 algorithmic procedures (called *queries* in this context) do not readily fulfill the definition for useful  
36 values of  $(\epsilon, \delta)$ , in particular the deterministic ones: randomization is mandatory. This randomization

---

```

model = DP_Sequential( # step 1: use DP_Sequential to build a model
    [
        # step 2: add Lipschitz layers of known sensitivity
        DP_BoundedInput(input_shape=(28, 28, 1), upper_bound=20.),
        DP_SpectralConv2D(filters=16, kernel_size=3, use_bias=False),
        DP_GroupSort(2),
        DP_Flatten(),
        DP_SpectralDense(10),
    ],
    noise_multiplier = 1.2, # step 3: choose DP parameters
    sampling_probability = batch_size / dataset_size,
) # step 4: compile the model, and choose any first order optimizer
model.compile(loss=DP_Crossentropy(), optimizer=Adam(1e-3))
model.fit( # step 5: train the model and measure the DP guarantees
    train_dataset, validation_data=val_dataset,
    epochs=num_epochs, callbacks=[DP_Accountant()]
)

```

---

Figure 1: An example of usage of our framework, illustrating how to create a small Lipschitz VGG and how to train it under  $(\epsilon, \delta)$ -DP guarantees while reporting  $(\epsilon, \delta)$  values.

37 comes at the expense of “utility”, i.e the usefulness of the output for downstream tasks [4]. The goal  
38 is then to strike a balance between privacy and utility, ensuring that the released information remains  
39 useful and informative for the intended purpose while minimizing the risk of privacy breaches. The  
40 privacy/utility trade-off yields a Pareto front, materialized by plotting  $\epsilon$  against a measurement of  
41 utility, such as validation accuracy for a classification task.

42 **Private gradient descent.** The SGD algorithm consists of a sequence of queries that (i) take the  
43 dataset in input, sample a minibatch from it, and return the gradient of the loss evaluated on the  
44 minibatch, before (ii) performing a descent step following the gradient direction. The sensitivity (see  
45 Definition 2) of SGD queries is proportional to the norm of the per-sample gradients. DP-SGD turns  
46 each query into a Gaussian mechanism by perturbing the gradients with a noise  $\zeta$ . The upper bound  
47 on gradient norms is generally unknown in advance, which leads practitioners to clip it to  $C > 0$ , in  
48 order to bound the sensitivity manually. This is problematic for several reasons: **1.** Hyper-parameter  
49 search on the broad-range clipping value  $C$  is required to train models with good privacy/utility trade-  
50 offs [5], **2.** The computation of per-sample gradients is expensive: DP-SGD is usually slower and  
51 consumes more memory than vanilla SGD, in particular for the large batch sizes often used in private  
52 training [6], **3.** Clipping the per-sample gradients biases their average [7]. This is problematic as the  
53 average direction is mainly driven by misclassified examples, that carry the most useful information  
54 for future progress.

55 **An unexplored approach: Lipschitz constrained networks.** We propose to train neural networks  
56 for which the parameter-wise gradients are provably and analytically bounded during the whole  
57 training procedure, in order to get rid of the clipping process. This allows for rapid training of models  
58 without a need for tedious hyper-parameter optimization.

59 The main reason why this approach has not been experimented much in the past is that upper bounding  
60 the gradient of neural networks is often intractable. However, by leveraging the literature of Lipschitz  
61 constrained networks [8], we show that these networks allows to estimate their gradient bound.  
62 This yields tight bounds on the sensitivity of SGD steps, making their transformation into Gaussian  
63 mechanisms inexpensive - hence the name **Clipless DP-SGD**.

64 Informally, the Lipschitz constant quantifies the rate at which the function’s output varies with respect  
65 to changes in its input. A Lipschitz constrained network is one in which its weights and activations  
66 are constrained such that it can only represent  $l$ -Lipschitz functions. In this work, we will focus our  
67 attention on feed-forward networks (refer to Definition 3). Note that the most common architectures,  
68 such as Convolutional Neural Networks (CNNs), Fully Connected Networks (FCNs), Residual  
69 Networks (ResNets), or patch-based classifiers (like MLP-Mixers), all fall under the category of  
70 feed-forward networks. We will also tackle the particular case of Gradient Norm Preserving (GNP)  
71 networks, a subset of Lipschitz networks that enjoy tighter bounds (see appendix).

72 **Contributions**

73 While the properties of Lipschitz constrained networks regarding their inputs are well explored, the  
 74 properties with respect to its parameters remain non-trivial. This work provides a first step to fill this  
 75 gap: our analysis shows that under appropriate architectural constraints, a  $l$ -Lipschitz network has a  
 76 tractable, finite Lipschitz constant with respect to its parameters. We prove that this Lipschitz constant  
 77 allows for easy estimation of the sensitivity of the gradient computation queries. The prerequisite and  
 78 details of the method to compute the sensitivities are explained in Section 2.

79 Our contributions are the following:

- 80 1. We extend the field of applications of Lipschitz constrained neural networks. So far the  
 81 literature focused on Lipschitzness with respect to the *inputs*: we extend the framework to  
 82 **compute the Lipschitzness with respect to the parameters**. This is exposed in Section 2.
- 83 2. We propose a **general framework to handle layer gradient steps as Gaussian mechanisms**  
 84 that depends on the loss and the model structure. Our framework covers widely used  
 85 architectures, including VGG and ResNets.
- 86 3. We show that SGD training of deep neural networks can be achieved **without gradient**  
 87 **clipping** using Lipschitz layers. This allows the use of larger networks and larger batch  
 88 sizes, as illustrated by our experiments in Section 4.
- 89 4. We establish connections between **Gradient Norm Preserving (GNP)** networks and **im-**  
 90 **proved privacy/utility trade-offs** (Section 3.1).
- 91 5. Finally, a **Python package**<sup>1</sup> companions the project, with pre-computed Lipschitz constant  
 92 and noise for each layer type, ready to be forked on any problem of interest (Section 3.2).

93 **1.1 Differential Privacy and Lipschitz Networks**

94 The definition of DP relies on the notion of neighboring datasets, i.e datasets that vary by at most one  
 95 example. We highlight below the central tools related to the field, inspired from [9].

96 **Definition 1** ( $(\epsilon, \delta)$ -Differential Privacy). *A labeled dataset  $\mathcal{D}$  is a finite collection of input/label*  
 97 *pairs  $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ . Two datasets  $\mathcal{D}$  and  $\mathcal{D}'$  are said to be neighboring*  
 98 *for the “replace-one” relation if they differ by at most one sample:  $\mathcal{D}' = \mathcal{D} \cup \{(x'_i, y'_i)\} \setminus \{(x_i, y_i)\}$ .*  
 99 *Let  $\epsilon$  and  $\delta$  be two non-negative scalars. A mechanism  $\mathcal{A}$  is  $(\epsilon, \delta)$ -DP if for any two neighboring*  
 100 *datasets  $\mathcal{D}$  and  $\mathcal{D}'$ , and for any  $S \subseteq \text{range}(\mathcal{A})$ :*

$$\mathbb{P}[\mathcal{A}(\mathcal{D}) \in S] \leq e^\epsilon \times \mathbb{P}[\mathcal{A}(\mathcal{D}') \in S] + \delta. \quad (1)$$

101 A cookbook to create a  $(\epsilon, \delta)$ -DP mechanism from a query is to compute its *sensitivity*  $\Delta$  (see  
 102 Definition 2), and to perturb its output by adding a Gaussian noise of predefined variance  $\zeta^2 = \Delta^2 \sigma^2$ ,  
 103 where the  $(\epsilon, \delta)$ -DP guarantees depends on  $\sigma$ . This yields what is called a *Gaussian mechanism* [3].

104 **Definition 2** ( $l_2$ -sensitivity). *Let  $\mathcal{M}$  be a query mapping from the space of the datasets to  $\mathbb{R}^p$ . Let  $\mathcal{N}$*   
 105 *be the set of all possible pairs of neighboring datasets  $\mathcal{D}, \mathcal{D}'$ . The  $l_2$  sensitivity of  $\mathcal{M}$  is defined by:*

$$\Delta(\mathcal{M}) = \max_{\mathcal{D}, \mathcal{D}' \in \mathcal{N}} \|\mathcal{M}(\mathcal{D}) - \mathcal{M}(\mathcal{D}')\|_2. \quad (2)$$

106 **Differentially Private SGD.** The classical algorithm keeps track of  $(\epsilon, \delta)$ -DP values with a *moments*  
 107 *accountant* [1] which allows to keep track of privacy guarantees at each epoch, by composing  
 108 different sub-mechanisms. For a dataset with  $N$  records and a batch size  $b$ , it relies on two parameters:  
 109 the sampling ratio  $p = \frac{b}{N}$  and the “noise multiplier”  $\sigma$  defined as the ratio between effective noise  
 110 strength  $\zeta$  and sensitivity  $\Delta$ . Bounds on gradient norm can be turned into bounds on sensitivity  
 111 of SGD queries. In “replace-one” policy for  $(\epsilon, \delta)$ -DP accounting, if the gradients are bounded by  
 112  $K > 0$ , the sensitivity of the gradients averaged on a minibatch of size  $b$  is  $\Delta = 2K/b$ .

113 Crucially, the algorithm requires a bound on  $\|\nabla_\theta \mathcal{L}(\hat{y}, y)\|_2 \leq K$ . The whole difficulty lies in  
 114 bounding tightly this value in advance for neural networks. Currently, gradient clipping serves as a  
 115 patch to circumvent the issue [1]. Unfortunately, clipping individual gradients in the batch is costly  
 116 and will bias the direction of their average, which may induce underfitting [7].

<sup>1</sup>Code and documentation are given as supplementary material during review process.

117 **Lipschitz constrained networks.** Our proposed solution comes from the observation that the norm  
 118 of the gradient and the Lipschitz constant are two sides of the same coin. The function  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$   
 119 is said  $l$ -Lipschitz for  $l_2$  norm if for every  $x, y \in \mathbb{R}^m$  we have  $\|f(x) - f(y)\|_2 \leq l\|x - y\|_2$ . Per  
 120 Rademacher’s theorem [10], its gradient is bounded:  $\|\nabla_x f\| \leq l$ . Reciprocally, continuous functions  
 121 gradient bounded by  $l$  are  $l$ -Lipschitz.

122 In Lipschitz networks, the literature has predominantly concentrated on investigating the control  
 123 of Lipschitzness with respect to the inputs (i.e bounding  $\nabla_x f$ ), primarily motivated by concerns  
 124 of robustness [11]. However, in this work, we will demonstrate that it is also possible to control  
 125 Lipschitzness with respect to parameters (i.e bounding  $\nabla_\theta f$ ), which is essential for ensuring privacy.  
 126 Our first contribution will point out the tight link that exists between those two quantities.

127 **Definition 3** (Lipschitz feed-forward neural network). *A feedforward neural network of depth  $D$ ,  
 128 with input space  $\mathcal{X} \subset \mathbb{R}^n$ , output space  $\mathcal{Y} \subset \mathbb{R}^K$  (e.g logits), and parameter space  $\Theta \subset \mathbb{R}^p$ , is a  
 129 parameterized function  $f : \Theta \times \mathcal{X} \rightarrow \mathcal{Y}$  defined by the sequential composition of layers  $f_d$ :*

$$f(\theta, x) := (f_D(\theta_d) \circ \dots \circ f_2(\theta_2) \circ f_1(\theta_1))(x). \quad (3)$$

130 *The parameters of the layers are denoted by  $\theta = (\theta_d)_{1 \leq d \leq D} \in \Theta$ . For affine layers, it corresponds  
 131 to bias and weight matrix  $\theta_d = (W_d, b_d)$ . For activation functions, there is no parameters:  $\theta_d = \emptyset$ .*

132 *Lipschitz networks are feed-forward networks, with the additionnal constraint that each  
 133 layer  $x_d \mapsto f_d(\theta_d, x_d) := y_d$  is  $l_d$ -Lipschitz for all  $\theta_d$ . Consequently, the function  $x \mapsto f(\theta, x)$   
 134 is  $l$ -Lipschitz with  $l = l_1 \times \dots \times l_d$  for all  $\theta \in \Theta$ .*

135 In practice, this is enforced by using activations with Lipschitz constant  $l_d$ , and by applying a con-  
 136 straint  $\Pi : \mathbb{R}^p \rightarrow \Theta$  on the weights of affine layers. This corresponds to spectrally normalized matrices  
 137 [12, 13], since for affine layers we have  $l_d = \|W_d\|_2 := \max_{\|x\| \leq 1} \|W_d x\|_2$  hence  $\Theta = \{\|W_d\| \leq l_q\}$ .

138 The seminal work of [8] proved that universal approximation in the set of  $l$ -Lipschitz functions was  
 139 achievable by this family of architectures. Concurrent approaches are based on regularization (like  
 140 in [14, 15, 16]) but they fail to produce formal guarantees. While they have primarily been studied in  
 141 the context of adversarial robustness [11, 17], recent works have revealed additional properties of  
 142 these networks, such as improved generalization [13, 18]. However, the properties of their parameter  
 143 gradient  $\nabla_\theta f(\theta, x)$  remain largely unexplored.

## 144 2 Clipless DP-SGD with $l$ -Lipschitz networks

145 Our framework consists of **1**, a method that computes the maximum gradient norm of a network with  
 146 respect to its parameters to obtain a *per-layer* sensitivity  $\Delta_d$ , **2**, a moments accountant that relies on  
 147 the per-layer sensitivities to compute  $(\epsilon, \delta)$ -DP guarantees. The method 1. is based on the recursive  
 148 formulation of the chain rule involved in backpropagation, while 2. keeps track of  $(\epsilon, \delta)$ -DP values  
 149 with RDP accounting. It requires some natural assumptions that we highlight below.

150 **Requirement 1** (Lipschitz loss.). *The loss function  $\hat{y} \mapsto \mathcal{L}(\hat{y}, y)$  must be  $L$ -Lipschitz with respect to  
 151 the logits  $\hat{y}$  for all ground truths  $y \in \mathcal{Y}$ . This is notably the case of Categorical Softmax-Crossentropy.*

152 The Lipschitz constants of common classification losses can be found in the appendix.

153 **Requirement 2** (Bounded input). *There exists  $X_0 > 0$  such that for all  $x \in \mathcal{X}$  we have  $\|x\| \leq X_0$ .*

154 While there exist numerous approaches for the parametrization of Lipschitz networks (e.g differenti-  
 155 able re-parametrization [19, 8], optimization over matrix manifolds [20] or projections [21]), our  
 156 framework only provides sensitivity bounds for projection-based algorithms (see appendix).

157 **Requirement 3** (Lipschitz projection). *The Lipschitz constraints must be enforced with a projection  
 158 operator  $\Pi : \mathbb{R}^p \rightarrow \Theta$ . This corresponds to Tensorflow [22] constraints and Pytorch [23] hooks.  
 159 Projection is a post-processing of private gradients: it induces no privacy leakage [3].*

160 To compute the per-layer sensitivities, our framework mimics the backpropagation algorithm, where  
 161 *Vector-Jacobian* products (VJP) are replaced by *Scalar-Scalar* products of element-wise bounds. For  
 162 an arbitrary layer  $x_d \mapsto f_d(\theta_d, x_d) := y_d$  the operation is sketched below:

$$\underbrace{\nabla_{x_d} \mathcal{L} := (\nabla_{y_d} \mathcal{L}) \frac{\partial f_d}{\partial x_d}}_{\text{Vector-Jacobian product: backpropagate gradients}} \implies \underbrace{\|\nabla_{x_d} \mathcal{L}\|_2 \leq \|\nabla_{y_d} \mathcal{L}\|_2 \times \left\| \frac{\partial f_d}{\partial x_d} \right\|_2}_{\text{Scalar-Scalar product: backpropagate bounds}}. \quad (4)$$

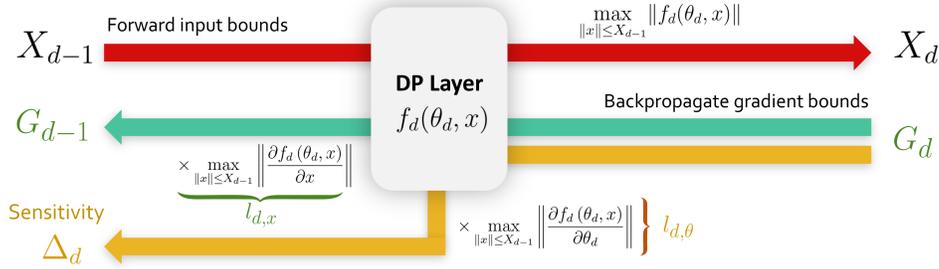


Figure 2: **Backpropagation for bounds**, Algorithm 1. Compute the per-layer sensitivity  $\Delta_d$ .

163 The notation  $\|\cdot\|_2$  must be understood as the spectral norm for Jacobian matrices, and the Euclidean  
 164 norm for gradient vectors. The scalar-scalar product is inexpensive. For Lipschitz layers the spectral  
 165 norm of the Jacobian  $\left\| \frac{\partial f}{\partial x} \right\|$  is kept constant during training with projection operator  $\Pi$ . The bound of the  
 166 gradient with respect to the parameters then takes a simple form:

$$\|\nabla_{\theta_d} \mathcal{L}\|_2 = \|\nabla_{y_d} \mathcal{L}\|_2 \times \left\| \frac{\partial f_d}{\partial \theta_d} \right\|_2. \quad (5)$$

167 Once again the operation is inexpensive. The upper bound  $\left\| \frac{\partial f}{\partial \theta} \right\|_2$  typically depends on the supremum  
 168 of  $\|x_d\|_2$ , that can also be analytically bounded, as exposed in the following section.

## 169 2.1 Backpropagation for bounds

170 The pseudo-code of **Clipless DP-SGD** is sketched in Algorithm 2. The algorithm avoids clipping by  
 171 computing a *per-layer* bound on the element-wise gradient norm. The computation of this *per-layer*  
 172 bound is described by Algorithm 1 (graphically explained in Figure 2). Crucially, it requires to  
 173 compute the spectral norm of the Jacobian of each layer with respect to input and parameters.

174 **Input bound propagation (line 2).** We compute  $X_d = \max_{\|x\| \leq X_{d-1}} \|f_d(x)\|_2$ . For activation  
 175 functions it depends on their range. For linear layers, it depends on the spectral norm of the operator  
 176 itself. This quantity can be computed with SVD or Power Iteration [24, 19], and constrained during  
 177 training using projection operator  $\Pi$ . In particular, it covers the case of convolutions, for which tight  
 178 bounds are known [25]. For affine layers, it additionally depends on the amplitude of the bias  $\|b_d\|$ .

179 **Remark 1** (Tighter bounds in literature.). *Although libraries such as Decomon [26] or*  
 180 *auto-LiRPA [27] provide tighter bounds  $X_d$  via linear relaxations [28, 29], our approach is ca-*  
 181 *capable of delivering practically tighter bounds than worst-case scenarios thanks to the projection*  
 182 *operator  $\Pi$ , while also being significantly less computationally expensive. Moreover, hybridizing our*  
 183 *method with scalable certification methods can be a path for future extensions.*

184 **Computing maximum gradient norm (line 6).** We bound the Jacobian  $\frac{\partial f_d(\theta_d, x)}{\partial \theta_d}$ . In neural networks,  
 185 the parameterized layers  $f(\theta, x)$  (fully connected, convolutions) are bilinear operators. Hence we  
 186 typically obtain bounds of the form:

$$\left\| \frac{\partial f_d(\theta_d, x)}{\partial \theta_d} \right\|_2 \leq K(f_d, \theta_d) \|x\|_2 \leq K(f_d, \theta_d) X_{d-1}, \quad (6)$$

187 where  $K(f_d, \Theta_d)$  is a constant that depends on the nature of the operator.  $X_{d-1}$  is obtained in line 2  
 188 with input bound propagation. Values of  $K(f_d, \theta_d)$  for popular layers are pre-computed in the library.

189 **Backpropagate cotangent vector bounds (line 7).** We bound the Jacobian  $\frac{\partial f_d(\theta_d, x)}{\partial x}$ . For activa-  
 190 tion functions this value can be hard-coded, while for affine layers it is the spectral norm of the linear  
 191 operator. Like before, this value is constrained with projection operator  $\Pi$ .

## 192 2.2 Privacy accounting for Clipless DP-SGD

193 Two strategies are available to keep track of  $(\epsilon, \delta)$  values as the training progresses, based on  
 194 accounting either a per-layer “local” sensitivity, either by aggregating them into a “global” sensitivity.

---

**Algorithm 1 Backpropagation for Bounds**( $f, X$ )

---

**Input:** Feed-forward architecture  $f(\theta, \cdot) = f_D(\theta_D, \cdot) \circ \dots \circ f_1(\theta_1, \cdot)$

**Input:** Weights  $\theta = (\theta_1, \theta_2, \dots, \theta_D)$ , input bound  $X_0$

- 1: **for all** layers  $1 \leq d \leq D$  **do**
  - 2:      $X_d \leftarrow \max_{\|x\| \leq X_{d-1}} \|f_d(\theta_d, x)\|_2$ . ▷ Input bounds propagation
  - 3: **end for**
  - 4:  $G \leftarrow L/b$ . ▷ Lipschitz constant of the loss for batchsize  $b$
  - 5: **for all** layers  $D \geq d \geq 1$  **do**
  - 6:      $\Delta_d \leftarrow G \max_{\|x\| \leq X_{d-1}} \left\| \frac{\partial f_d(\theta_d, x)}{\partial \theta_d} \right\|_2$ . ▷ Compute sensitivity from gradient norm
  - 7:      $G \leftarrow G \max_{\|x\| \leq X_{d-1}} \left\| \frac{\partial f_d(\theta_d, x)}{\partial x} \right\|_2 = G l_d$ . ▷ Backpropagate cotangent vector bounds
  - 8: **end for**
  - 9: **return** sensitivities  $\Delta_1, \Delta_2, \dots, \Delta_D$
- 

---

**Algorithm 2 Clipless DP-SGD with local** sensitivity accounting

---

**Input:** Feed-forward architecture  $f(\theta, \cdot) = f_D(\theta_D, \cdot) \circ \dots \circ f_1(\theta_1, \cdot)$

**Input:** Initial weights  $\theta = (\theta_1, \theta_1, \dots, \theta_D)$ , learning rate  $\eta$ , noise multiplier  $\sigma$ .

- 1: **repeat**
  - 2:      $\Delta_1, \Delta_2, \dots, \Delta_D \leftarrow$  **Backpropagation for Bounds**( $f, X$ ).
  - 3:     Update Moment Accountant state with **local** sensitivities  $\Delta_1, \Delta_2, \dots, \Delta_D$ .
  - 4:     Sample a batch  $\mathcal{B} = \{(x_1, y_1), (x_2, y_2), \dots, (x_b, y_b)\}$ .
  - 5:     Compute per-layer averaged gradient:  $g_d := \frac{1}{b} \sum_{i=1}^b \nabla_{\theta_d} \mathcal{L}(f(\theta, x_i), y_i)$ .
  - 6:     Sample local noise:  $\zeta_d \sim \mathcal{N}(0, \sigma \Delta_d)$ .
  - 7:     Perform noisified gradient step:  $\theta_d \leftarrow \theta_d - \eta(g_d + \zeta_d)$ .
  - 8:     Enforce Lipschitz constraint with projection:  $\theta_d \leftarrow \Pi(\theta_d)$ .
  - 9: **until** privacy budget  $(\epsilon, \delta)$ -DP budget has been reached.
- 

195 **The “global” strategy.** Illustrated in the appendix, this strategy simply aggregates the individual  
196 sensitivities  $\Delta_d$  of each layer to obtain the global sensitivity of the whole gradient vector  $\Delta =$   
197  $\sqrt{\sum_d \Delta_d^2}$ . The origin of the clipping-based version of this strategy can be traced back to [30]. With  
198 noise variance  $\sigma^2 \Delta^2$  we recover the accountant that comes with DP-SGD. It tends to overestimate  
199 the true sensitivity (in particular for deep networks), but its implementation is straightforward with  
200 existing tools.

201 **The “local” strategy.** Recall that we are able to characterize the sensitivity  $\Delta_d$  of every layer of  
202 the network. Hence, we can apply a different noise to each of the gradients. We dissect the whole  
203 training procedure in Figure 3. At same noise multiplier  $\sigma$ , it tends to produce a higher value of  $\epsilon$  per  
204 epoch than “global” strategy, but has the advantage over the latter to add smaller effective noise  $\zeta$  to  
205 each weight.

206 We rely on the autodp<sup>2</sup> library [32, 33, 34] as it uses the Renyi Differential Privacy (RDP) adaptive  
207 composition theorem [35, 36], that ensures tighter bounds than naive DP composition.

### 208 3 From theory to practice

209 Beyond the application of Algorithms 1 and 2, our framework provides numerous opportunities to  
210 enhance our understanding of prevalent techniques identified in the literature. An in-depth exploration  
211 of these is beyond the scope of this work, so we focus on giving insights on promising tracks based  
212 on our theoretical analysis. In particular, we discuss how the tightness of the bound provided by  
213 Algorithm 1 can be influenced by working on the architecture, the input pre-processing and the loss  
214 post-processing.

---

<sup>2</sup><https://github.com/yuxiangw/autodp> distributed under Apache License 2.0 licence.

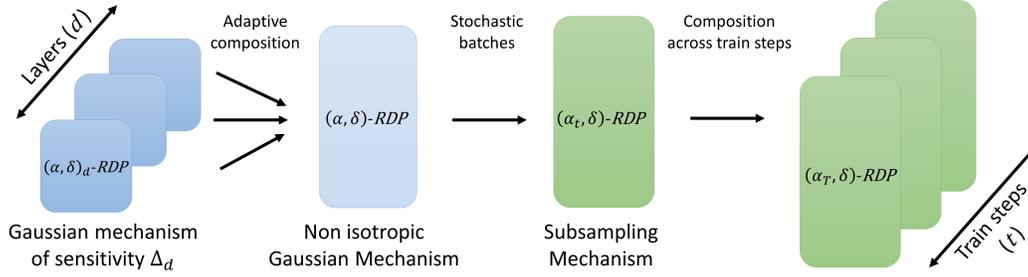


Figure 3: **Accountant for locally enforced differential privacy.** (i) The gradient query for each layer is turned into a Gaussian mechanism [9], (ii) their composition at the scale of the whole network is a non isotropic Gaussian mechanism, (iii) that benefits from amplification via sub-sampling [31], (iv) the train steps are composed over the course of training.

### 215 3.1 Gradient Norm Preserving networks

216 We can manually derive the bounds obtained from Algorithm 2 across diverse configurations. Below,  
 217 we conduct a sensitivity analysis on  $l$ -Lipschitz networks.

218 **Theorem (informal) 1. Gradient Norm of Lipschitz Networks.** *Assume that every layer  $f_d$  is*  
 219  *$K$ -Lipschitz, i.e  $l_1 = \dots = l_D = K$ . Assume that every bias is bounded by  $B$ . We further assume that*  
 220 *each activation is centered in zero (e.g ReLU, tanh, GroupSort). We recall that  $\theta = [\theta_1, \theta_2, \dots, \theta_D]$ .*  
 221 *Then the global upper bound of Algorithm 2 can be expanded analytically.*

222 **1. If  $K < 1$  we have:**  $\|\nabla_{\theta} \mathcal{L}(f(\theta, x), y)\|_2 = \mathcal{O}(L(K^D(X_0 + B) + 1))$ .

223 *Due to the  $K^D \ll 1$  term this corresponds to a vanishing gradient phenomenon [37]. The output of*  
 224 *the network is essentially independent of its input, and the training is nearly impossible.*

225 **2. If  $K > 1$  we have:**  $\|\nabla_{\theta} \mathcal{L}(f(\theta, x), y)\|_2 = \mathcal{O}(LK^D(X_0 + B))$ .

226 *Due to the  $K^D \gg 1$  term this corresponds to an exploding gradient phenomenon [38]. The upper*  
 227 *bound becomes vacuous for deep networks: the added noise  $\zeta$  is at risk of being too high.*

228 **3. If  $K = 1$  we have:**  $\|\nabla_{\theta} \mathcal{L}(f(\theta, x), y)\|_2 = \mathcal{O}\left(L\left(X_0 + \sqrt{D} + \sqrt{BX_0D} + BD^{3/2}\right)\right)$ ,

229 *which for linear layers without biases further simplify to  $\mathcal{O}(L(X_0 + \sqrt{D}))$ .*

230 The formal statement can be found in appendix. From Theorem 1 we see that most favorable bounds  
 231 are achieved by 1-Lipschitz neural networks with 1-Lipschitz layers. In classification tasks, they are  
 232 not less expressive than conventional networks [18]. Hence, this choice of architecture is not at the  
 233 expense of utility. Moreover an accuracy/robustness trade-off exists, determined by the choice of  
 234 loss function [18]. However, setting  $K = 1$  merely ensures that  $\|\nabla_x f\| \leq 1$ , and in the worst-case  
 235 scenario we have  $\|\nabla_x f\| < 1$  almost everywhere. This could result in a situation where the bound of  
 236 case 3 in Theorem 1 is not tight, leading to an underfitting regime as in case  $K < 1$ . With Gradient  
 237 Norm Preserving (GNP) networks [17], we expect to mitigate this issue.

238 **Controlling  $K$  with Gradient Norm Preserving (GNP) networks.** GNP networks are 1-Lipschitz  
 239 neural networks with the additional constraint that the Jacobian of layers consists of orthogonal  
 240 matrices. They fulfill the Eikonal equation  $\left\|\frac{\partial f_d(\theta_d, x_d)}{\partial x_d}\right\|_2 = 1$  for any intermediate activation  
 241  $f_d(\theta_d, x_d)$ . Without biases these networks are also norm preserving:  $\|f(\theta, x)\| = \|x\|$ .

242 As a consequence, the gradient of the loss with respect to the parameters is easily bounded by

$$\|\nabla_{\theta_d} \mathcal{L}\| = \|\nabla_{y_D} \mathcal{L}\| \times \left\|\frac{\partial f_d(\theta_d, x_d)}{\partial \theta_d}\right\|, \quad (7)$$

243 which for weight matrices  $W_d$  further simplifies to  $\|\nabla_{W_d} \mathcal{L}\| \leq \|\nabla_{y_D} \mathcal{L}\| \times \|f_{d-1}(\theta_{d-1}, x_{d-1})\|$ . We  
 244 see that this upper bound crucially depends on two terms that can be analyzed separately. On one  
 245 hand,  $\|f_{d-1}(\theta_{d-1}, x_{d-1})\|$  depends on the scale of the input. On the other,  $\|\nabla_{y_D} \mathcal{L}\|$  depends on the  
 246 loss, the predictions and the training stage. We show below how to intervene on these two quantities.

247 **Remark 2** (Implementation of GNP Networks). *In practice, GNP are parametrized with GroupSort*  
 248 *activation [8, 39], Householder activation [40], and orthogonal weight matrices [17, 41]. Strict*  
 249 *orthogonality is challenging to enforce, especially for convolutions for which it is still an active*  
 250 *research area (see [42, 43, 44, 45, 46] and references therein). Our line of work traces an additional*  
 251 *motivation for the development of GNP and the bounds will strengthen as the field progresses.*

252 **Controlling  $X_0$  with input pre-processing.** The weight gradient norm  $\|\nabla_{\theta_d} \mathcal{L}\|$  indirectly depends  
 253 on the norm of the inputs. This observation implies that the pre-processing of input data significantly  
 254 influences the bounding of sensitivity. Multiple strategies are available to keep the input’s norm under  
 255 control: projection onto the ball (“norm clipping”), or projection onto the sphere (“normalization”).  
 256 In the domain of natural images for instance, this result sheds light on the importance of color space  
 257 such as RGB, HSV, YIQ, YUV or Grayscale. These strategies are natively handled by our library.

258 **Controlling  $L$  with the hybrid approach, loss gradient clipping.** As training progresses, the  
 259 magnitude of  $\|\nabla_f \mathcal{L}\|$  tends to diminish when approaching a local minima, quickly falling below the  
 260 upper bound and diminishing the gradient norm to noise ratio. To circumvent the issue, the gradient  
 261 clipping strategy is still available in our framework. Crucially, instead of clipping the parameter  
 262 gradient  $\nabla_{\theta} \mathcal{L}$ , any intermediate gradient  $\nabla_{f_d} \mathcal{L}$  can be clipped during backpropagation. This can  
 263 be achieved with a special “clipping layer” that behaves like the identity function at the forward  
 264 pass, and clips the gradient during the backward pass. The resulting cotangent vector is not a true  
 265 gradient anymore, but rather a descent direction [47]. In vanilla DP-SGD the clipping is applied on  
 266 the batched gradient  $\nabla_{W_d} \mathcal{L}$  of size  $b \times h^2$  for matrix weight  $W_d \in \mathbb{R}^{h \times h}$  and clipping this vector can  
 267 cause memory issues or slowdowns [6]. In our case,  $\nabla_{y_D} \mathcal{L}$  is of size  $b \times h$  which reduces overhead.

### 268 3.2 Lip-dp library

269 To foster and spread accessibility, we provide an opensource tensorflow library for Clipless DP-SGD  
 270 training, named `lip-dp`. It provides an exposed Keras API for seamless usability. It is implemented  
 271 as a wrapper over the Lipschitz layers of `dee1-lip3` library [48]. Its usage is illustrated in Figure 1.

## 272 4 Experimental results

273 We validate our implementation with a speed benchmark against competing approaches, and we  
 274 present the privacy/utility Pareto front that can be obtained with GNP networks.

275 **Speed and memory consumption.** We bench-  
 276 marked the median runtime per epoch of vanilla  
 277 DP-SGD against the one of Clipless DP-SGD,  
 278 on a CNN architecture and its Lipschitz equiv-  
 279 alent respectively. The experiment was run on  
 280 a GPU with 48GB video memory. We compare  
 281 against the implementation of `tf_privacy`,  
 282 `opacus` and `optax`. In order to allow a fair com-  
 283 parison, when evaluating Opacus, we reported  
 284 the runtime with respect to the logical batch size,  
 285 while capping the physical batch size to avoid  
 286 Out Of Memory error (OOM). Although our lib-  
 287 rary does not implement logical batching yet,  
 288 it is fully compatible with this feature.

289 An advantage of projection  $\Pi$  over per-sample  
 290 gradient clipping is that the projection cost is  
 291 independent of the batch size. Fig 4 validates that our method scales much better than vanilla  
 292 DP-SGD, and is compatible with large batch sizes. It offers several advantages: firstly, a larger batch  
 293 size contributes to a decrease of the sensitivity  $\Delta \propto 1/b$ , which diminishes the ratio between noise  
 294 and gradient norm. Secondly, as the batch size  $b$  increases, the variance decreases at the parametric  
 295 rate  $\mathcal{O}(\sqrt{b})$  (as demonstrated in appendix), aligning with expectations. This observation does not  
 296 apply to DP-SGD: gradient clipping biases the direction of the average gradient, as noticed by [7].

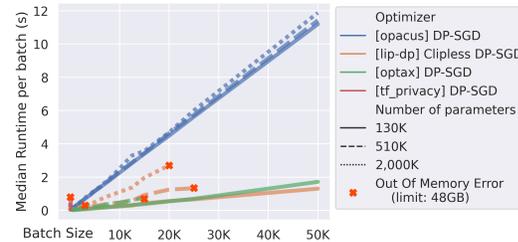


Figure 4: **Our approach outperforms concurrent frameworks in terms of runtime and memory:** we trained CNNs (ranging from 130K to 2M parameters) on CIFAR-10, and report the median batch processing time (including noise, and constraints application  $\Pi$  or gradient clipping).

<sup>3</sup><https://github.com/dee1-ai/dee1-lip> distributed under MIT License (MIT).

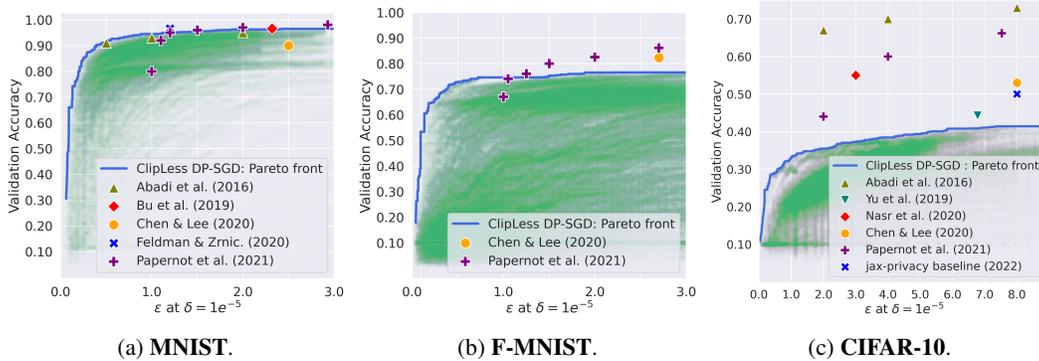


Figure 5: **Our framework paints a clearer picture of the privacy/utility trade-off.** We trained models in an "out of the box setting" (no pre-training, no data augmentation and no handcrafted features) on multiple tasks. While our results align with the baselines presented in other frameworks, we recognize the importance of domain-specific engineering. In this regard, we find the innovations introduced in [49, 50, 51] and references therein highly relevant. These advancements demonstrate compatibility with our framework and hold potential for future integration.

297 **Pareto front of privacy/utility trade-off.** We performed a search over a broad range of hyper-  
 298 parameters values to cover the Pareto front between utility and privacy. Results are reported in  
 299 Figure 5. We emphasize that our experiments did not use the elements behind the success of most  
 300 recent papers (pre-training, data preparation, or handcrafted feature are examples). Hence our  
 301 results are more representative of the typical performance that can be obtained in an "out of the  
 302 box" setting. Future endeavors or domain-specific engineering can enhance the performance even  
 303 further, but such improvements currently lie beyond the scope of our work. We also benchmarked  
 304 architectures inspired from VGG [52], Resnet [53] and MLP\_Mixers [54] see appendix for more  
 305 details. Following standard practices of the community [2], we used *sampling without replacement*  
 306 at each epoch (by shuffling examples), but we reported  $\epsilon$  assuming *Poisson sampling* to benefit from  
 307 privacy amplification [31]. We also ignore the privacy loss that may be induced by hyper-parameter  
 308 search, which is a limitation per recent studies [5], but is common practice.

309 **5 Limitations and future work**

310 Although this framework offers a novel approach to address differentially private training, it introduces  
 311 new challenges. We primary rely on GNP networks, where high performing architectures are  
 312 quite different from the usual CNN architectures. As emphasized in Remark 2, we anticipate that  
 313 progress in these areas would greatly enhance the effectiveness of our approach. Additionally, to  
 314 meet requirement 3, we rely on projections, necessitating additional efforts to incorporate recent  
 315 advancements associated with differentiable reparametrizations [42, 43]. It is worth noting that  
 316 our methodology is applicable to most layers. Another limitation of our approach is the accurate  
 317 computation of sensitivity  $\Delta$ , which is challenging due to the non-associativity of floating-point  
 318 arithmetic and its impact on numerical stability [55]. This challenge is exacerbated on GPUs, where  
 319 operations are inherently non-deterministic [56]. Finally, as mentioned in Remark 1, our propagation  
 320 bound method can be refined.

321 **6 Concluding remarks and broader impact**

322 Besides its main focus on differential privacy, our work provides **(1) a motivation to further develop**  
 323 **Gradient Norm Preserving architectures.** Furthermore, the development of networks with known  
 324 Lipschitz constant with respect to parameters is a question of independent interest, **(2) a useful tool for**  
 325 **the study of the optimization dynamics** in neural networks. Finally, Lipschitz networks are known  
 326 to enjoy certificates against adversarial attacks [17, 57], and from generalization guarantees [13],  
 327 without cost in accuracy [18]. We advocate for the spreading of their use in the context of robust and  
 328 certifiable learning.

## References

- 329
- 330 [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar,  
331 and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC*  
332 *conference on computer and communications security*, pages 308–318, 2016.
- 333 [2] Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H Brendan  
334 McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Thakurta. How to dp-fy ml: A  
335 practical guide to machine learning with differential privacy. *arXiv preprint arXiv:2303.00654*,  
336 2023.
- 337 [3] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to  
338 sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography*  
339 *Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284.  
340 Springer, 2006.
- 341 [4] Mário S Alvim, Miguel E Andrés, Konstantinos Chatzikokolakis, Pierpaolo Degano, and  
342 Catuscia Palamidessi. Differential privacy: on the trade-off between utility and information  
343 leakage. In *Formal Aspects of Security and Trust: 8th International Workshop, FAST 2011,*  
344 *Leuven, Belgium, September 12-14, 2011. Revised Selected Papers 8*, pages 39–54. Springer,  
345 2012.
- 346 [5] Nicolas Papernot and Thomas Steinke. Hyperparameter tuning with renyi differential privacy.  
347 In *International Conference on Learning Representations*, 2022.
- 348 [6] Jaewoo Lee and Daniel Kifer. Scaling up differentially private deep learning with fast per-  
349 example gradient clipping. *Proceedings on Privacy Enhancing Technologies*, 2021(1), 2021.
- 350 [7] Xiangyi Chen, Steven Z Wu, and Mingyi Hong. Understanding gradient clipping in private sgd:  
351 A geometric perspective. *Advances in Neural Information Processing Systems*, 33:13773–13782,  
352 2020.
- 353 [8] Cem Anil, James Lucas, and Roger Grosse. Sorting out lipschitz function approximation. In  
354 *International Conference on Machine Learning*, pages 291–301. PMLR, 2019.
- 355 [9] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Founda-*  
356 *tions and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- 357 [10] Leon Simon et al. *Lectures on geometric measure theory*. The Australian National University,  
358 Mathematical Sciences Institute, Centre . . . , 1983.
- 359 [11] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-  
360 fellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference*  
361 *on Learning Representations*, 2014.
- 362 [12] Yuichi Yoshida and Takeru Miyato. Spectral norm regularization for improving the generaliz-  
363 ability of deep learning. *arXiv preprint arXiv:1705.10941*, 2017.
- 364 [13] Peter L Bartlett, Dylan J Foster, and Matus Telgarsky. Spectrally-normalized margin bounds for  
365 neural networks. In *Proceedings of the 31st International Conference on Neural Information*  
366 *Processing Systems*, pages 6241–6250, 2017.
- 367 [14] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville.  
368 Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*,  
369 volume 30, pages 5767–5777. Curran Associates, Inc., 2017.
- 370 [15] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier.  
371 Parseval networks: Improving robustness to adversarial examples. In *International Conference*  
372 *on Machine Learning*, pages 854–863. PMLR, 2017.
- 373 [16] Henry Gouk, Eibe Frank, Bernhard Pfahringer, and Michael J Cree. Regularisation of neural  
374 networks by enforcing lipschitz continuity. *Machine Learning*, 110:393–416, 2021.

- 375 [17] Qiyang Li, Saminul Haque, Cem Anil, James Lucas, Roger B Grosse, and Jörn-Henrik Jacobsen.  
376 Preventing gradient attenuation in lipschitz constrained convolutional networks. In *Advances in*  
377 *Neural Information Processing Systems (NeurIPS)*, volume 32, Cambridge, MA, 2019. MIT  
378 Press.
- 379 [18] Louis Béthune, Thibaut Boissin, Mathieu Serrurier, Franck Mamalet, Corentin Friedrich, and  
380 Alberto Gonzalez Sanz. Pay attention to your loss : understanding misconceptions about  
381 lipschitz neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun  
382 Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 383 [19] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization  
384 for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- 385 [20] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix*  
386 *manifolds*. Princeton University Press, 2008.
- 387 [21] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial  
388 networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- 389 [22] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,  
390 Greg Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,  
391 Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser,  
392 Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray,  
393 Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul  
394 Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden,  
395 Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: Large-scale  
396 machine learning on heterogeneous distributed systems, 2015.
- 397 [23] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
398 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative  
399 style, high-performance deep learning library. *Advances in neural information processing*  
400 *systems*, 32, 2019.
- 401 [24] Lloyd N Trefethen and David Bau. *Numerical linear algebra*, volume 181. Siam, 2022.
- 402 [25] S Singla and S Feizi. Fantastic four: Differentiable bounds on singular values of convolution  
403 layers. In *International Conference on Learning Representations (ICLR)*, 2021.
- 404 [26] Airbus. Decomon. <https://github.com/airbus/decomon>, 2023.
- 405 [27] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya  
406 Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified  
407 robustness and beyond. *Advances in Neural Information Processing Systems*, 33:1129–1141,  
408 2020.
- 409 [28] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. An abstract domain for  
410 certifying neural networks. *Proceedings of the ACM on Programming Languages*, 3(POPL):1–  
411 30, 2019.
- 412 [29] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neu-  
413 ral network robustness certification with general activation functions. *Advances in neural*  
414 *information processing systems*, 31, 2018.
- 415 [30] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially  
416 private recurrent language models. In *International Conference on Learning Representations*,  
417 2018.
- 418 [31] Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy amplification by subsampling: Tight  
419 analyses via couplings and divergences. *Advances in Neural Information Processing Systems*,  
420 31, 2018.
- 421 [32] Yu-Xiang Wang, Borja Balle, and Shiva Prasad Kasiviswanathan. Subsampled rényi differential  
422 privacy and analytical moments accountant. In *The 22nd International Conference on Artificial*  
423 *Intelligence and Statistics*, pages 1226–1235. PMLR, 2019.

- 424 [33] Yuqing Zhu and Yu-Xiang Wang. Poission subsampled rényi differential privacy. In *International*  
425 *Conference on Machine Learning*, pages 7634–7642. PMLR, 2019.
- 426 [34] Yuqing Zhu and Yu-Xiang Wang. Improving sparse vector technique with renyi differential  
427 privacy. *Advances in Neural Information Processing Systems*, 33:20249–20258, 2020.
- 428 [35] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations*  
429 *symposium (CSF)*, pages 263–275. IEEE, 2017.
- 430 [36] Ilya Mironov, Kunal Talwar, and Li Zhang. R\`enyi differential privacy of the sampled gaussian  
431 mechanism. *arXiv preprint arXiv:1908.10530*, 2019.
- 432 [37] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent  
433 neural networks. In *International conference on machine learning*, pages 1310–1318. Pmlr,  
434 2013.
- 435 [38] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with  
436 gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.
- 437 [39] Ugo Tanielian and Gerard Biau. Approximating lipschitz continuous functions with groupsort  
438 neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages  
439 442–450. PMLR, 2021.
- 440 [40] Zakaria Mhammedi, Andrew Hellicar, Ashfaque Rahman, and James Bailey. Efficient orthogonal  
441 parametrisation of recurrent neural networks using householder reflections. In *International*  
442 *Conference on Machine Learning*, pages 2401–2409. PMLR, 2017.
- 443 [41] Shuai Li, Kui Jia, Yuxin Wen, Tongliang Liu, and Dacheng Tao. Orthogonal deep neural  
444 networks. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1352–1368,  
445 2019.
- 446 [42] Asher Trockman and J Zico Kolter. Orthogonalizing convolutional layers with the cayley  
447 transform. In *International Conference on Learning Representations*, 2021.
- 448 [43] Sahil Singla and Soheil Feizi. Skew orthogonal convolutions. In *International Conference on*  
449 *Machine Learning*, pages 9756–9766. PMLR, 2021.
- 450 [44] El Mehdi Achour, François Malgouyres, and Franck Mamalet. Existence, stability and scalability  
451 of orthogonal convolutional neural networks. *The Journal of Machine Learning Research*,  
452 23(1):15743–15798, 2022.
- 453 [45] Sahil Singla and Soheil Feizi. Improved techniques for deterministic l2 robustness. *arXiv*  
454 *preprint arXiv:2211.08453*, 2022.
- 455 [46] Xiaojun Xu, Linyi Li, and Bo Li. Lot: Layer-wise orthogonal training on improving l2 certified  
456 robustness. *arXiv preprint arXiv:2210.11620*, 2022.
- 457 [47] Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press,  
458 2004.
- 459 [48] Mathieu Serrurier, Franck Mamalet, Alberto González-Sanz, Thibaut Boissin, Jean-Michel  
460 Loubes, and Eustasio Del Barrio. Achieving robustness in classification using optimal transport  
461 with hinge regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
462 *and Pattern Recognition*, pages 505–514, 2021.
- 463 [49] Nicolas Papernot, Abhradeep Thakurta, Shuang Song, Steve Chien, and Úlfar Erlingsson.  
464 Tempered sigmoid activations for deep learning with differential privacy. In *Proceedings of the*  
465 *AAAI Conference on Artificial Intelligence*, volume 35, pages 9312–9321, 2021.
- 466 [50] Florian Tramèr and Dan Boneh. Differentially private learning needs better features (or much  
467 more data), 2021.
- 468 [51] Soham De, Leonard Berrada, Jamie Hayes, Samuel L Smith, and Borja Balle. Unlocking  
469 high-accuracy differentially private image classification through scale. *arXiv preprint*  
470 *arXiv:2204.13650*, 2022.

- 471 [52] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
472 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 473 [53] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
474 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
475 pages 770–778, 2016.
- 476 [54] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas  
477 Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-  
478 mixer: An all-mlp architecture for vision. *Advances in neural information processing systems*,  
479 34:24261–24272, 2021.
- 480 [55] David Goldberg. What every computer scientist should know about floating-point arithmetic.  
481 *ACM computing surveys (CSUR)*, 23(1):5–48, 1991.
- 482 [56] Hadi Jooybar, Wilson WL Fung, Mike O’Connor, Joseph Devietti, and Tor M Aamodt. Gpudet:  
483 a deterministic gpu architecture. In *Proceedings of the eighteenth international conference on*  
484 *Architectural support for programming languages and operating systems*, pages 1–12, 2013.
- 485 [57] Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas.  
486 Efficient and accurate estimation of lipschitz constants for deep neural networks. *Advances in*  
487 *Neural Information Processing Systems*, 32, 2019.