

## A Differential Privacy Metric

The  $(\epsilon, \delta)$ -Differential Privacy is the most widely used privacy metric which is defined as

**Definition A.1** (Differential Privacy (Dwork et al. 2006b; Dwork 2006; Dwork et al. 2006a)). Let  $\mathcal{D}^n$  be space of datasets including  $n$  samples. A randomized algorithm  $M : \mathcal{D}^n \rightarrow \mathcal{Z}$  is  $(\epsilon, \delta)$ -DP if for all adjacent<sup>4</sup>  $D, D' \in \mathcal{D}^n$  and for all subset of outputs  $S \subseteq \mathcal{Z}$ , we have

$$P(M(D) \in S) \leq e^\epsilon P(M(D') \in S) + \delta.$$

If  $\delta = 0$ , then  $M$  achieves pure DP which is stronger than the case  $\delta > 0$ . To achieve better properties described in Section 3.1, we give a privacy guarantee by a new privacy metric and its proofs here.

A variant of DP is called  $\omega$ -truncated  $\rho$ -Concentrated Differential Privacy (Bun et al. 2018) as follows.

**Definition A.2** ( $(\rho, \omega)$ -tCDP). Let  $\omega > 1$  and  $\rho > 0$ . A randomized algorithm  $M : \mathcal{D}^n \rightarrow \mathbb{R}$  satisfied  $(\rho, \omega)$ -tCDP if, for all adjacent inputs  $d, d' \in \mathcal{D}^n$ ,

$$D_\alpha(M(d) \| M(d')) \leq \rho\alpha, \forall \alpha \in (1, \omega)$$

where  $D_\alpha(\cdot \| \cdot)$  denotes the Rényi divergence (Rényi 1961) of order  $\alpha$ .

The  $(\rho, \omega)$ -tCDP provides nice properties for private learning as discussed in Section 3.1. The general usage can be found in the Algorithms 1 and 2. However, the privacy parameters,  $(\rho, \omega)$ , have to be carefully initialized and chosen to satisfy their constraint.

In this section, we first summarize the properties of  $(\rho, a)$ -ctCDP and then provide proofs based on the theoretic results of  $(\rho, \omega)$ -tCDP. Also, we show how to use the properties of the ctCDP to initialize the privacy parameters and add essential constraints in detail.

### A.1 Properties of $(\rho, a)$ -ctCDP

We suggest an  $a$ -constrained truncated  $\rho$ -CDP (or  $(\rho, a)$ -ctCDP) based on the  $(\rho, \omega)$ -tCDP to provide a minimal example with tight composition bound, linear composition function and a simple Gaussian mechanism. It is defined in Definition 3.1. The proposal originates from the scenario where the privacy budget is constrained as a constant, typically as  $(\epsilon, \delta)$ -DP, and a series of properties are required for scalable gradient perturbation (Algorithm 1). Other than ctCDP, more complicated privacy metric can be used on demand. Here, we first summarize the theoretic results of ctCDP with proofs in the following sections. Note that most of the results are identical to tCDP by fixing  $\omega$  as  $\omega_a$ . The major difference lies on the Lemma A.5 where the noising mechanism is assumed to be derived from zCDP, resulting in a different constraint on  $\rho$ .

**Lemma A.1** (Relation to  $(\epsilon, \delta)$ -DP). Suppose a randomized algorithm  $M$  satisfies  $(\rho_a, a)$ -ctCDP with  $a = \log(1/\delta)/\epsilon$  and  $\rho_a$  defined in Definition 3.1. Then  $M$  satisfies  $(\epsilon, \delta)$ -DP.

<sup>4</sup>By adjacent, we mean the two datasets have only one different entry at most.

**Lemma A.2** (Relation to tCDP). Suppose a randomized algorithm  $M$  satisfies  $(\rho, a)$ -ctCDP. Then  $M$  satisfies  $(\rho, \omega_a)$ -tCDP where  $\omega_a$  is defined in Definition 3.1.

**Definition A.3** (Sensitivity). Given  $f : \mathcal{D}^n \rightarrow \mathbb{R}^d$ , the  $L_2$  sensitivity of  $f$ , denoted by  $\Delta_2(f)$  (or simply  $\Delta$ ) is defined as:

$$\Delta_2(f) = \max_{D, D'} \|f(D) - f(D')\|_2,$$

where  $D, D'$  are adjacent datasets.

**Lemma A.3.** The sensitivity of  $\|\nabla_t\|_2^2$  is  $\Delta_2(\|\nabla_t\|_2^2) = (2|D_t| - 1)C_g^2$ .

*Proof.* By definition,

$$\begin{aligned} & \max_{D, D'} \left\| \|\nabla_t\|_2^2 - \|\nabla'_t\|_2^2 \right\|_2 \\ &= \max_{D, D'} \left\| \left\| \sum_i \nabla_{t,i} \right\|_2^2 - \left\| \sum_{i \neq m} \nabla_{t,i} \right\|_2^2 \right\|_2 \end{aligned}$$

where we assuming the  $m$ -th sample is eliminated from computing the gradient  $\nabla'_t$ . Also, we have

$$\begin{aligned} \left\| \sum_i \nabla_{t,i} \right\|_2^2 &= \sum_{i,j} \nabla_{t,i}^\top \nabla_{t,j}, \\ \left\| \sum_i \nabla'_{t,i} \right\|_2^2 &= \sum_{i,j} \nabla_{t,i}^\top \nabla_{t,j} - 2 \sum_j \nabla_{t,m}^\top \nabla_{t,j} + \nabla_{t,m}^\top \nabla_{t,m}. \end{aligned}$$

Thus,

$$\begin{aligned} & \left\| \sum_i \nabla_{t,i} \right\|_2^2 - \left\| \sum_i \nabla'_{t,i} \right\|_2^2 \\ &= (2 \sum_j \nabla_{t,j} - \nabla_{t,m})^\top \nabla_{t,m} \end{aligned}$$

whose norm is lower than or equal to  $(2|D_t| - 1)C_g^2$ . Thus,  $\Delta_2(\|\nabla_t\|_2^2) = (2|D_t| - 1)C_g^2$ .  $\square$

**Lemma A.4** (Composition & Post-processing). Let two mechanisms be  $M : \mathcal{D}^n \rightarrow \mathcal{Y}$  and  $M' : \mathcal{D}^n \times \mathcal{Y} \rightarrow \mathcal{Z}$ . Suppose  $M$  satisfies  $(\rho_1, a)$ -ctCDP and  $M'(\cdot, y)$  satisfies  $(\rho_2, a)$ -ctCDP for  $\forall y \in \mathcal{Y}$ . Then, mechanism  $M'' : \mathcal{D}^n \rightarrow \mathcal{Z}$  (defined by  $M''(x) = M'(x, M(x))$ ) satisfies  $(\rho_1 + \rho_2, a)$ -ctCDP.

The Lemma A.4 is directly given by tCDP when we transform ctCDP to tCDP with fixed  $\omega_a$  (see Lemma A.2).

Instead of deriving a canonical noise mechanism for the ctCDP, we directly use the Gaussian mechanism theorem of zCDP and amplify its privacy cost in the form of ctCDP by subsampling.

**Lemma A.5** (Gaussian mechanism for  $\rho$ -zCDP (Bun and Steinke 2016)). Let  $f : \mathcal{D}^n \rightarrow \mathcal{Z}$  have sensitivity  $\Delta$ . Define a randomized algorithm  $M : \mathcal{D}^n \rightarrow \mathcal{Z}$  by

$$M(x) \leftarrow f(x) + \mathcal{N}(0, \sigma^2).$$

Then  $M$  satisfies  $\frac{\Delta^2}{2\sigma^2}$ -zCDP.

**Lemma A.6** ( $(\rho, a)$ -ctCDP from  $\rho$ -zCDP by privacy amplification through subsampling). Let  $\omega_a = (1+a) + \sqrt{a(a+1)}$ . Let  $\rho, q \in (0, 0.1]$ ,  $n, N \in \mathcal{Y}$  with  $q = n/N$  (sampling rate) and satisfy

$$\log(1/q) \geq 3\rho(2 + \log_2(1/\rho)), \quad (9)$$

$$0 \leq \rho < \min \left\{ \frac{\log(1/q)}{4\omega_a}, \frac{\rho_a}{13q^2} \right\}. \quad (10)$$

Let  $M : \mathcal{D}^n \rightarrow \mathbb{R}$  satisfy  $\rho$ -zCDP. Define the mechanism  $M_s : \mathcal{D}^N \rightarrow \mathcal{Y}$  by  $M_s(x) = M(x_S)$ , where  $x_S \in \mathcal{D}^N$  is the restriction of  $x \in \mathcal{D}^N$  to the entries specified by a uniformly random subset  $S \subseteq [N]$  with  $|S| = n$ .

The algorithm  $M_s : \mathcal{D}^N \rightarrow \mathcal{Y}$  satisfies  $(13q^2\rho, a)$ -ctCDP.

Remarkably, the Lemma A.6 assumes the sub-routine is  $\rho$ -zCDP. Since the subsampling-based privacy amplification happens after the noise mechanism, it is natural to use the Gaussian mechanism (Lemma A.5) which results in a zCDP privacy cost.

These privacy guarantees are derived from  $(\rho, \omega)$ -tCDP by constraining the range of  $\rho$  and  $\omega$  where  $\omega$  is simply a constant. By sacrificing the flexibility of the privacy parameters, we can get a single parameter metric which is simple in notation. Moreover, by fixing the  $\omega$ , we will be able to update the privacy parameters by gradient descent in meta-learning<sup>5</sup>.

In this paper, we apply the  $(\rho, a)$ -ctCDP to the gradient perturbation scenario (Algorithm 1) where the total privacy cost  $\rho_{tot}$  is constrained by  $(\epsilon, \delta)$ -DP. A completed pipeline of private learning includes initialization of the parameters based on the given privacy budget and the protected learning (Algorithm 1). For the convenience of implementation, we provide the detailed steps of the initialization in Algorithm 3 as supplementary to private learning. The algorithm is based on the ctCDP with subsampling. If subsampling is not used, only the step budget needs to be modified as  $\rho_0 = \rho'_0$  with the upper bound  $\rho_{ub} = \rho_a$ .

## A.2 Budget constraint from standard DP

Our motivation for proposing ctCDP is that existing private learning methods are typically compared by performance under a given privacy budget. Therefore, we first introduce a budget constraint using the standard DP, which is translated as bounds on the tCDP parameters.

**Theorem A.1** (Transformation of tCDP to DP). Suppose  $M$  satisfy  $(\rho, \omega)$ -tCDP. Then, for all  $\delta > 0$  and all  $1 < \alpha \leq \omega$ ,  $M$  satisfies  $(\epsilon, \delta)$ -DP with

$$\epsilon = \begin{cases} \rho + 2\sqrt{\rho \log(1/\delta)}, & \log(1/\delta) \leq (\omega - 1)^2\rho \\ \rho\omega + \frac{\log(1/\delta)}{\omega - 1}, & \log(1/\delta) \geq (\omega - 1)^2\rho \end{cases}$$

When  $\epsilon$  and  $\delta$  are fixed, we want to maximize the available budget  $\rho$  and fix  $\omega_a$ . We consider  $\log(1/\delta) \geq (\omega - 1)^2\rho$  to find the upper bound of  $\rho$  when  $\omega$  can be maximized, as well. First, we solve the linear function of  $\rho$ , i.e.,  $\epsilon = \rho\omega + \frac{\log(1/\delta)}{\omega - 1}$ , given some  $\omega$ . Let

$$a \triangleq \log(1/\delta)/\epsilon$$

<sup>5</sup>Note that the composition of  $\omega$  in  $(\rho, \omega)$ -tCDP is not continuously differentiable if all sub-mechanisms have varying  $\omega$ .

## Algorithm 3 Private learning initialization.

$(\epsilon, \delta)$ -DP, the number of learning iterations ( $T$ ) estimated from non-private algorithms.

- 1: Transformation from DP to ctCDP (Lemma A.1):  $a \leftarrow \log(1/\delta)/\epsilon$  and compute  $\rho_a$  and  $\omega_a$  by Definition 3.1
- 2:  $\rho_{tot} \leftarrow \rho_a$ .
- 3: Estimate step budget by uniformly decomposing  $\rho_{tot}$  into  $T$  steps:  $\rho'_0 \leftarrow \frac{\rho_{tot}}{T}$  (Lemma A.4).
- 4: Estimate a batch sampling rate  $q$ , e.g.,  $q \leftarrow (\sqrt{|D|} + 10)/|D|$ .
- 5:  $\rho_0 \leftarrow \frac{\rho'_0}{13q^2}$  (subsampling by Lemma A.6)
- 6: If  $\rho_0$  and  $q$  do not satisfy Eqs. (9) and (10), re-estimate  $q$  by choosing the smaller solution to

$$\log(1/q) = 3 \frac{\rho_{tot}}{13q^2} (2 - \log_2(\frac{\rho'_0}{13q^2}))$$

$$\frac{\rho'_0}{13q^2} = \frac{\log(1/q)}{4\omega_a}.$$

Then re-compute  $\rho_0$  using the new  $q$ .

- 7: Get the upper bound of step budgets:  $\rho_{ub}$  (Eq. (10)).
- 8: Compute noise scale by Lemma A.5:
- 9: estimated step noise:  $\sigma_g = \Delta/\sqrt{2\rho_0}$
- 10: step noise lower bound:  $\sigma_{min} = \Delta/\sqrt{2\rho_{ub}}$
- 11: Compute batch size  $|D_t| \equiv \lfloor qn \rfloor$ .
- 12: Output:  $\rho_{tot}, \sigma_g, \sigma_{min}, q, |D_t|$

and denote the solution as

$$\rho_{\text{tCDP}} = \frac{\epsilon}{\omega} \left( 1 - \frac{\log(1/\delta)/\epsilon}{\omega - 1} \right) = \frac{\epsilon(\omega - (a + 1))}{\omega(\omega - 1)}. \quad (11)$$

Now, we substitute the  $\rho_{\text{tCDP}}$  into  $\log(1/\delta) \geq (\omega - 1)^2\rho$  to obtain

$$a\omega \geq (\omega - 1)(\omega - (a + 1))$$

$$\Rightarrow (1 + a) - \sqrt{a(a + 1)} \leq \omega \leq (1 + a) + \sqrt{a(a + 1)}$$

By the definition of tCDP,  $\rho_{\text{tCDP}} \geq 0$  and  $\omega > \alpha > 1$ . Thus,

$$\omega > 1 \text{ and } \omega \geq a + 1.$$

Because  $a + 1 > (1 + a) - \sqrt{a(a + 1)}$  and  $a > 0$ , the only solution to  $\log(1/\delta) = (\omega - 1)^2\rho$  is the upper bound. Now, we denote the upper bound of  $\omega$  as

$$\omega_a \triangleq (1 + a) + \sqrt{a(a + 1)}$$

and substitute it into Eq. (11) to get

$$\rho_a \triangleq \epsilon \frac{\sqrt{a(a + 1)}}{((1 + a) + \sqrt{a(a + 1)})(a + \sqrt{a(a + 1)})} \quad (12)$$

which is also the solution of  $\epsilon = \rho + 2\sqrt{\rho \log(1/\delta)}$ .

Until now, we finish the proof of the bounds in Definition 3.1. In addition, it is straightforward to get that  $(\rho_a, \omega_a)$ -tCDP is equivalent to  $(\rho_a, a)$ -ctCDP. Thus, we have the Lemma A.1. We can further extend it for  $\rho < \rho_a$ .

**Lemma A.7** ( $(\rho, a)$ -ctCDP to  $(\epsilon, \delta)$ -DP). Suppose a randomized algorithm  $M$  satisfy  $(\rho, a)$ -ctCDP for  $0 < \rho \leq \rho_a$  and  $a > 0$ , then  $M$  satisfies  $(\epsilon, \delta)$ -tCDP with

$$\epsilon = \frac{(\omega_a - 1)\omega_a}{\omega_a - a - 1}\rho, \delta = \exp(-a\epsilon).$$

### 801 A.3 Noise mechanism

802 The canonical noise for tCDP is a Gaussian noise reshaped  
803 by a sinh function. We restate the theorem by rearranging the  
804 variables.

**Theorem A.2** (Sinh-Normal mechanism for  $(\rho, \omega)$ -tCDP (Bun et al. 2018)). Let  $f : \mathcal{D}^n \rightarrow \mathcal{Z}$  has sensitivity  $\Delta$ . Let  $\rho = \frac{8\Delta^2}{\sigma^2}$ ,  $\omega$  satisfy  $\frac{1}{4\omega^2} \leq \rho < 16$  and  $A = 8\Delta\omega$ . Define a randomized algorithm  $M : \mathcal{D}^n \rightarrow \mathcal{Z}$  by

$$M(x) \leftarrow f(x) + A \operatorname{arsinh}\left(\frac{1}{A}\mathcal{N}(0, \sigma^2)\right).$$

805 Then  $M$  satisfies  $(\rho, \omega)$ -tCDP.

806 If  $\omega \rightarrow \infty$ , then  $A \rightarrow \infty$  in which case the  $(\rho, \omega =$   
807  $\infty)$ -tCDP is just  $\rho$ -zCDP and the Sinh-Normal distribution  
808 degrades as the normal distribution. However, due to the  
809 truncation of  $\omega$ , the privacy cost, i.e.,  $\frac{8\Delta^2}{\sigma^2}$ , is not as optimal  
810 as  $\rho$ -zCDP. Therefore, we use the noise mechanism of  $\rho$ -  
811 zCDP (Lemma A.5) when  $\omega \rightarrow \infty$ .

### 812 A.4 Privacy amplification by subsampling

813 In stochastic gradient descent, a batch of data subsampled  
814 from the whole dataset is used to update models. It is criti-  
815 cal for implementing scalable learning algorithms. Because  
816 of the randomness of subsampling, it provably reduce the  
817 privacy cost. Technically, there are two ways to subsample  
818 the batch. One is random sampling without replacement or  
819 reshuffling (RF) which is widely used in the non-private deep  
820 learning. Yu et al. (2019) proved the composed privacy cost  
821 is the maximum of batch costs in one RF epoch. Numerically,  
822 each batch is  $q\rho$ -zCDP if the full batch cost is  $\rho$  and batch  
823 sample rate is  $q$ . In this case, the dynamic budget alloca-  
824 tion for batches within one epoch is always worse than the  
825 uniform schedule.

826 The other strategy is the random sampling with replace-  
827 ment (RS), for example, in SGD-MA for the private deep  
828 learning (Abadi et al. 2016). Compared to RF, RS injects  
829 more randomness and therefore scale down the privacy cost  
830 more (Yu et al. 2019), for example, a  $q^2$  factor in the MA.  
831 The lack of privacy amplification for RS motivates the de-  
832 velopment of extensions. Both tCDP (Bun et al. 2018) and  
833 the modified zCDP (Yu et al. 2019) spot the issue theoretic-  
834 ally and provide similar solutions by truncating the order of  
835 Rényi divergence. A privacy amplification of tCDP is given  
836 Theorem A.3.

837 **Theorem A.3** (Privacy amplification by subsampling for  
838  $(\rho, \omega)$ -tCDP). Let  $\rho, q \in (0, 0.1]$  and positive integers  $n, N$   
839 with  $q = n/N$  and  $\log(1/q) \geq 3\rho(2 + \log_2(1/\rho))$ . Let  $M :$   
840  $\mathcal{D}^n \rightarrow \mathbb{R}$  satisfy  $(\rho, \omega')$ -tCDP for  $\omega' \geq \frac{1}{2\rho} \log(1/q) \geq 3$ .  
841 Define the mechanism  $M_s : \mathcal{D}^N \rightarrow \mathcal{Y}$  by  $M_s(x) = M(x_S)$ ,

842 where  $x_S \in \mathcal{D}^N$  is the restriction of  $x \in \mathcal{D}^N$  to the en-  
843 tries specified by a uniformly random subset  $S \subseteq [N]$  with  
844  $|S| = n$ .

845 The algorithm  $M_s : \mathcal{D}^N \rightarrow \mathcal{Y}$  satisfies  $(13q^2\rho, \omega'')$ -tCDP  
846 for  $\omega'' = \frac{\log(1/q)}{4\rho}$ .

847 In comparison, the modified zCDP does not have a strict  
848 theoretic proof of the scale factor of the privacy cost but  
849 empirically shows that  $q^2\rho$  works for a wide range of  $\rho$ . Here,  
850 we use tCDP to derive the range of privacy parameters.

851 In Theorem A.3,  $\omega''$  is a variable depending on the  $\rho$  rather  
852 than  $\omega'$ . Thus, we let  $\omega' \rightarrow \infty$  to degrade  $(\rho, \omega')$ -tCDP as  
853  $\rho$ -zCDP when  $\rho \in (0, 16)$ .

Recall our target is to simplify the tCDP by eliminat-  
ing  $\omega$ . Because the subsampled mechanism also satisfies  
 $(13q^2\rho, \omega_a)$ -tCDP if  $\omega_a \leq \omega''$ , we constrain  $\rho$  as

$$\rho \leq \min \left\{ \frac{\log(1/q)}{4\omega_a}, \frac{\log(1/q)}{6}, \frac{\rho_a}{13q^2} \right\} \quad (13)$$

854 where  $\frac{\rho_a}{13q^2}$  comes from the constraint of  $(\rho, a)$ -ctCDP on  
855  $13q^2\rho$ . Typically, when  $\omega_a > 1.5$ , the  $\log(1/q)/6$  can be  
856 ignored. Because 1.5 is too small to reach for  $\omega_a$  in practice,  
857 we may assume it is satisfied generally.

## 858 B Methodology supplementaries

### 859 B.1 Model-based private learning

860 Here, we provide the formal statement and proof of Theo-  
861 rem 4.1.

**Theorem B.1** (Privacy guarantee of model-based gradient  
862 descent). Suppose a gradient-based algorithm Algorithm 1  
863 is protected by Algorithm 2 and  $\sigma(\cdot)$  and  $\pi(\cdot)$  are crafted  
864 fully independently from the private data. The output of the  
865 algorithm, i.e.,  $\theta_T$  (assuming the loops stop at step  $T$ ), is  
866  $\hat{\rho}$ -ctCDP where  $\hat{\rho} \leq \rho_{\text{tot}}$ , if  $f_C(\cdot)$ ,  $f_S(\cdot)$  and  $\rho(\cdot)$  are defined  
867 based on ctCDP properties (Lemmas A.4 to A.6). 868

869 *Proof.* For brevity, we omit the  $a$  in notations. Denote  
870 the sub-routine defined in Algorithm 2 is  $g_t, \rho_t, z_{t+1} =$   
871  $M_t(\nabla_t, z_t)$  where  $z_t$  denotes the hidden states. Then each  
872 iteration of private learning in Algorithm 1 can be abstracted  
873 as  $\theta_{t+1}, z_{t+1} = A_t(\theta_t, g_t, z_t, \rho_t)$ . Because of the linear com-  
874 position, Lemma A.4, the condition  $f_C(\rho_{1:t}) > \rho_{\text{tot}}$  can be  
875 justified by  $\rho_{\text{residual}} > 0$  where  $\rho_{\text{residual}} \leftarrow \rho_{\text{residual}} - \rho_t$ .

876 By rearranging variables, without changing the meaning  
877 of the mappings, we can write the iteration as  $\theta_{t+1}, z_{t+1} =$   
878  $A_t(M_t(\theta_t, z_t))$  where  $d$  denotes the private batch data. Sup-  
879 pose  $(\theta_t, z_t)$  is  $\hat{\rho}_t$ -ctCDP w.r.t. the dataset and the mapping  
880  $M_t(\cdot, \cdot)$  is  $\rho_t$ -ctCDP w.r.t. the dataset. Thus, according to  
881 Lemma A.4,  $\theta_{t+1}, z_{t+1} = A_t(M_t(\theta_t, z_t))$  is  $\hat{\rho}_{t+1}$ -ctCDP  
882 where  $\hat{\rho}_{t+1} = \hat{\rho}_t + \rho_t$ .

Next, we show  $M_t(\cdot, \cdot)$  is  $\rho_t$ -ctCDP for  $t \leq T$  and some  
 $\rho_t < \infty$ . According to Lemma A.5, the noised gradient is  
 $1/2\sigma_t^2$ -zCDP and the noised gradient norm is  $1/2\sigma_g^2$ -zCDP  
(note its sensitivity is proved by Lemma A.3). Further using  
the Lemma A.6, we can compute

$$\rho_t = 13q^2 \left( \frac{1}{2\sigma_t^2} + \frac{1}{2\sigma_g^2} \right) < \infty$$

883 if  $\sigma_t$  and  $\sigma_g$  are non-zero.

884 Now we show  $A_1(M_1(\theta_1, z_1))$  is  $\rho_1$ -ctCDP. Typically,  
 885  $\theta_1, z_1$  are randomly initialized or constantly zero which  
 886 are independent from the dataset. Therefore,  $(\theta_1, z_1)$  is  
 887 0-ctCDP. By Lemma A.4, because  $M_1(\cdot, \cdot)$  is  $\rho_1$ -ctCDP,  
 888  $A_1(M_1(\theta_1, z_1))$  is  $\rho_1$ -ctCDP.

In summary, the output of model-based private learning,  
 i.e.,  $\theta_T = A_T(M_T(\theta_T, z_T))$  (omitting  $z_{T+1}$ ), is  $\hat{\rho}_T$ -ctCDP  
 where

$$\hat{\rho}_T = \hat{\rho}_{T-1} + \rho_T = \sum_{i=1}^T \rho_i \leq \rho_{\text{tot}}.$$

889 □

## 890 B.2 Augmented Lagrangian algorithm

891 Given  $\mu_0 > 0$ , tolerance  $\tau_0 > 0$  (Nocedal and Wright 1999)  
 892 (Chapter 17), starting point  $\sigma_0$  and  $\lambda^0$ , the variables are itera-  
 893 tively updated:

1. Line search  $s$  such that  $\sigma^+$  is an approximate minimizer  
 of  $\mathcal{L}_{\text{aug}}$  (the gradient norm is less than  $\tau_k$ ):

$$\sigma^+ = \sigma - s \left[ \nabla_{\sigma} F(T, \sigma_T) + \frac{dh}{d\sigma}(z - h(\sigma)/\mu) \right] \quad (14)$$

- 894 2. If the final convergence criteria satisfied, stop with approx-  
 895 imate solution  $\sigma$
3. Update Lagrange multiplier:

$$z^+ = z - h(\sigma^+)/\mu \quad (15)$$

- 896 4. Choose new penalty parameter  $\mu^+ \in (0, \mu)$ .

897 where  $s$  is the step size and we let  $\sigma$  be a vec-  
 898 tor  $[\sigma_1, \dots, \sigma_T]^\top$  or constant scalar. The update on  $\sigma^+$   
 899 can be replaced by another line search, i.e.,  $\sigma^+ =$   
 900  $\arg \min_{\sigma'} \mathcal{L}_{\text{arg}}(\sigma'(s))$  where  $\sigma'(s)$  is given by Eq. (14). In  
 901 practice, we want to avoid the second time of unrolling  $\sigma$   
 902 because it is required in Eq. (15). To fix this issue, we proceed  
 903 with steps 3, 4, first and then finally perform step 1.

## 904 B.3 Analysis of the gradients

A generic gradient descent method can be summarized as a  
 set of sequential updates on the parameter  $\theta$ , i.e.,

$$\begin{aligned} \theta_T &= \theta_1 + \sum_{t=1}^{T-1} g_t \\ &= \theta_1 + \sum_{t=1}^{T-1} \pi(\nabla_t + \sigma_t \nu_t), \quad \nu_t \sim \mathcal{N}(0, I). \end{aligned}$$

Assume  $\frac{\partial \sigma_t}{\partial \sigma_{t-1}} = 0$  and  $\frac{\partial g_t}{\partial \sigma_{t-1}} = 0$ . Therefore, we can com-  
 pute the gradient w.r.t.  $\sigma_t$  as

$$\begin{aligned} \frac{\partial f_T}{\partial \sigma_t} &= \frac{\partial \tilde{\nabla}_t}{\partial \sigma_t} \frac{\partial g_t}{\partial \tilde{\nabla}_t} \frac{\partial f_T}{\partial \theta_T} \\ &= \nu_t^\top \frac{\partial \pi(\tilde{\nabla}_t)}{\partial \tilde{\nabla}_t} \nabla_T \end{aligned} \quad (16)$$

$$\begin{aligned} &= \frac{1}{\sigma_t} (\tilde{\nabla}_t - \nabla_t)^\top \frac{\partial \pi(\tilde{\nabla}_t)}{\partial \tilde{\nabla}_t} \nabla_T \\ &\approx \frac{1}{\sigma_t} (\pi(\tilde{\nabla}_t) - \pi(\nabla_t))^\top \nabla_T \end{aligned} \quad (17)$$

where  $\tilde{\nabla}_t = \nabla_t + \sigma_t \nu_t$  and we approximate in the last  
 term by Taylor expansion. Taking expectation, we can see  
 from Eq. (16) that the gradient is related to the covariance  
 between the noise  $\nu_t$  and the final gradient  $\nabla_T$ . Intuitively,  
 if  $\frac{\partial \sigma_t}{\partial \sigma_{t-1}} = 0$  and  $\frac{\partial g_t}{\partial \sigma_{t-1}} = 0$ , the gradient updates on  $\sigma_t$  will  
 increase a lot when  $t \ll T$ . Together with the observation in  
 Eq. (17), it can be witnessed that the scheduler is decided by  
 the denoising effect of the  $\pi$ .

In our implementation, we use an RNN to model the  $\pi$  and  
 $\sigma$  which could greatly denoise the updates according to its  
 memory. Furthermore, in Appendix B.6, we give an exact  
 bound of the utility which closely relate the scheduler and  
 projector together.

## B.4 Optimality Stability

The optimality of learning scheduler is achieved by vanishing  
 the gradient in Eq. (3). If both of the two terms in Eq. (3) are  
 zero, the second term will be quite unstable since the  $f_t$  and  
 $\hat{F}$  are random variables. Especially when the expectation of  
 the gradient is estimated by only few samples, the instability  
 will be a major issue. Here, we focus on the second term  
 and analyze the probability for it to be zero. Assume the  
 optimization of  $f$  has converged and the expectation of the  
 gradient is estimated by one sample.

**Non-batch algorithms.** An optimal case for the non-batch  
 algorithm is  $f_t = f_{t'}$  for all  $t, t' \in \mathcal{T}$  where  $\mathcal{I}_t \neq 0$  and  
 $\mathcal{I}_{t'} \neq 0$ . By convergence, we assume the expected loss values  
 cannot be decreased any more since the private updates are  
 totally governed by the noise instead of gradients. Therefore,  
 we can further assume the losses are identically independently  
 distributed with variance  $\text{Var}[f_i] \equiv \sigma_f^2$ . According to the  
 Chebyshev's inequality, we have, for a constant  $\xi$ ,

$$P(|f_{t'} - f_t| > \xi) \leq \frac{2\sigma_f^2}{\xi^2}$$

which does not vanish since  $\sigma_f$  is non-zero due to private  
 noise. That means only using one optimization process can  
 barely reach the zero gradient condition.

**Batch algorithms.** The condition  $f_t = f_{t'}$  can be easily  
 extended to the batch case, i.e.,

$$\frac{1}{|\mathcal{B}_j|} \sum_{t \in \mathcal{B}_j} f_t = \frac{1}{|\mathcal{B}_i|} \sum_{t \in \mathcal{B}_i} f_t \quad (18)$$



for all  $i, j \in \mathcal{T}_B$  where  $\mathcal{I}_i \neq 0$  and  $\mathcal{I}_j \neq 0$ . Therefore,

$$P \left( \left| \frac{1}{|\mathcal{B}_i|} \sum_{t \in \mathcal{B}_j} f_t - \frac{1}{|\mathcal{B}_i|} \sum_{t \in \mathcal{B}_i} f_t \right| > \xi \right) \leq \frac{2\sigma_f^2}{\xi^2 |\mathcal{B}_i|}$$

931 which has a smaller failure probability if  $|\mathcal{B}_i| > 1$ . In other  
932 words, the batch algorithm is stabler.

### 933 B.5 Batch Augmented Lagrangian Algorithm

With the basic augmented objective Eq. (2), we can extend it to the batch case, i.e.,

$$L_b^{\text{aug}}(\sigma; r_b) = \mathbb{E}[\bar{F}_i] - z_b h_b + \frac{\|h_b\|_2^2}{2\mu_b}$$

where  $\bar{F}_i$  is the batch-averaged loss defined in Eq. (5). After decomposing the budget constraint into batches, the augmented objective on the whole optimization process has to be replaced by

$$\mathcal{L}_b^{\text{aug}}(\sigma, r; \rho_{\text{tot}}) = \sum_{b=1}^B L_b^{\text{aug}}(\sigma; r_b) - z_r h_r + \frac{\|h_r\|_2^2}{2\mu_r},$$

934 which constrains the  $b$ -th-batch privacy cost by  $r_b$  and the  
935 overall cost by  $\rho_{\text{tot}}$ . Note that  $z_b$  will be gradually reduced  
936 to zero when the batch constraint is getting tighter. With the  
937 Lagrangian multiplier  $z_b$ , the batch scheduler will be allowed  
938 to fetch needed budget slightly ignoring the constraint  $r_b$ .  
939 Therefore, we can define  $\hat{\rho}_b = f_C(\{\rho(\sigma_t)\}_{t \in \mathcal{B}_b}) - z_b \mu_b$  to  
940 be the batch privacy cost supplemented by the Lagrangian  
941 multiplier.

With the constraint decomposition, we can update the  $\sigma$  only using one batch loss  $L_b^{\text{aug}}(\sigma; r_b)$  independently if  $r_b$  is fixed. Then, we update  $r$  by optimizing Eq. (6), i.e.,

$$\mathcal{L}^{\text{aug}}(r; \rho_{\text{tot}}) = \sum_{b=1}^B \frac{\|\hat{\rho}_b - r_b\|_2^2}{2\mu_b} - z_r h_r + \frac{\|h_r\|_2^2}{2\mu_r}$$

942 We conceptually illustrate the enforcement of constraints  
943 between batch budget  $r_b$  and the global budget  $\rho_{\text{tot}}$  in Fig. 5.  
944 The global budget allocation will be enforced to align the  
945 total budget. However, the relation between the batch cost  
946 and the batch budget are bi-directions. When we optimize  
947 w.r.t.  $r$ , the  $r_b$  will also be encouraged to align the  $\hat{\rho}_b$ . When  
948 we optimize w.r.t.  $\rho_t$  in batch  $b$ , the batch privacy cost will  
949 be enforced to match the budget  $r_b$ .

950 In the unconstrained batch algorithm, we will use an one-  
951 pass fashion to update the meta-model, i.e., the optimizer.  
952 That means we do not need to store any batch data (except  
953 for the meta-model) that has been used, which could greatly  
954 reduce the space complexity. However, for budget constrained  
955 L2P, we cannot directly drop the used batches, since the batch  
956 state is essential to check if we need to adjust the constraint  
957 to fulfill the budget requirement. In Fig. 5, the dependency is  
958 represented by the interaction between the batch privacy cost  
959 and the constraint.

960 Rather than a one-pass method, we suggest a two-pass  
961 way to update the parameters. First, we update the meta-  
962 models with one pass. Second, by fixing the meta-models,  
963 we unroll the protected learning, update and store  $\hat{\rho}_b$ ,  $z_b$  and  
964  $\mu_b$ . With the recorded data, we minimize  $\mathcal{L}^{\text{aug}}(r; \rho_{\text{tot}})$  w.r.t.  $r$   
965 and update corresponding AL variables.

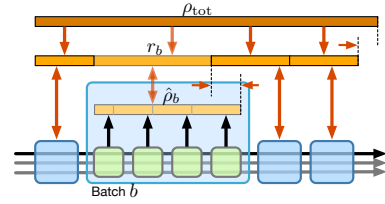


Figure 5: Illustration of the privacy budget constraint in batches. We strip the output of the protectors to be only the privacy costs which are composed into the batch cost  $\hat{\rho}_b$  supplemented by the Lagrangian multiplier. The red arrows between  $r_b$  and  $\hat{\rho}_b$  are the interaction during L2P optimization where they will both be forced to match the value of the other.

### B.6 Optimality analysis

In this paper, we cast the searching for the optimal scheduler as a learning problem. Formally, we minimize the objective:

$$\mathbb{E}[\hat{F}(\sigma)] = \mathbb{E} \left[ \frac{\sum_{t=1}^T \mathcal{I}_t f_t}{\sum_{t=1}^T \mathcal{I}_t} \right].$$

We first assume  $\mathcal{I}_t$  is a general weight function varying by  $t$ . Recall the gradient of  $\hat{F}(\sigma)$  w.r.t.  $\sigma$  (Eq. (3)) is

$$\partial \hat{F}(\sigma) = \frac{\sum_{t=1}^T \mathcal{I}_t \partial f_t}{\sum_{t=1}^T \mathcal{I}_t} + \frac{\sum_{t=1}^T (f_t - \hat{F}) \partial \mathcal{I}_t}{\sum_{t=1}^T \mathcal{I}_t}. \quad (19)$$

By vanishing the gradient, we can get the optimal condition of the L2P objective.

We restrict the non-zero range of  $\mathcal{I}_t$  within  $[T_\rho - 1, T_\rho]$ . Define the weight  $\alpha = \mathcal{I}_{T_\rho-1} / (\mathcal{I}_{T_\rho-1} + \mathcal{I}_{T_\rho})$ .

970 We first summarize the major results in Theorem B.2, with  
971 which we can get an approximated convergence guarantee  
972 to some (local or global) solution of our objective  $f$ . It gives  
973 us the insight that the final gradient norm is upper bounded  
974 by the covariance between the accumulated noise variables  
975 (transformed by  $\frac{\partial g_t}{\partial \nabla_t}$ ) and the final gradient. The upper bound  
976 will be improved when we train the projector on a fixed  
977 scheduler. In short words, to improve the utility, the projector  
978 training have to denoise the protected updates which reduces  
979 the covariance between the  $\nu_t$  and the  $g_t$ . And the covariance  
980 between  $\nu_t$  and  $f_{t+1}$  is reduced meanwhile. As a result, we  
981 will see  $\mathbb{C}_{T_\rho-1}$  approaching zero.  
982

**Definition B.1** ( $L$ -smooth function). A differentiable objective function  $f: \Theta \times \mathcal{X} \rightarrow \mathbb{R}$  is  $L$ -smooth over  $\theta \in \Theta$  with respect to the norm  $\|\cdot\|$  if for any  $x \in \mathcal{X}$  and  $\theta_1, \theta_2 \in \Theta$ , we have:

$$\|\nabla f(\theta_1, x) - \nabla f(\theta_2, x)\|_* \leq L \|\theta_1 - \theta_2\|,$$

where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ . If  $\|\cdot\|$  is  $l_2$ -norm, this yields

$$f(\theta_1, x) - f(\theta_2, x) \leq \nabla^\top f(\theta_2, x)(\theta_1 - \theta_2) + \frac{L}{2} \|\theta_1 - \theta_2\|^2.$$

**Theorem B.2** (Utility bound of a stationary L2P protector). Suppose  $f$  is  $L$ -smooth and  $\sigma$  is independent from the noise

variables  $\nu_t$ . If the  $\sigma$  is a stationary point of the constrained optimization problem, the following is satisfied:

$$\mathbb{E} \|\nabla_{T_\rho-1}\|^2 \leq \sigma_{\zeta_{T_\rho-1}}^2 p - \frac{2L}{C_\sigma} (\alpha \mathbb{C}_{T_\rho-1} + (1-\alpha) \mathbb{C}_{T_\rho}), \quad (20)$$

983 where  $C_\sigma$  is a constant depending on the scheduler  $\sigma$   
 984 (Eq. (24)),  $\mathbb{C}_T$  represents the covariance between the noise  
 985 and the true gradient (Eq. (32)) and  $\sigma_{\zeta_{T_\rho-1}}$  is the upper  
 986 bound of reduced noise variance (Eq. (27)).

*Proof.* Generally, we assume the  $\mathcal{I}_t$  be a function of  $h_t$  whose gradient is

$$\partial \mathcal{I}_t = \frac{\partial \mathcal{I}_t}{\partial \sigma} = \frac{\partial \mathcal{I}_t}{\partial h_t} \frac{\partial h_t}{\partial \sigma}.$$

Let us first look into  $\frac{\partial h_t}{\partial \sigma}$  which is

$$\begin{aligned} \frac{\partial h_t}{\partial \sigma} &= \sum_{\tau=1}^t \frac{\partial \rho_\tau}{\partial \sigma} \frac{\partial \sigma_\tau}{\partial \sigma} = - \sum_{\tau=1}^t \frac{\Delta^2}{\sigma_\tau^3} \frac{\partial \sigma_\tau}{\partial \sigma} = - \sum_{\tau=1}^t \frac{2\rho_\tau}{\sigma_\tau} \frac{\partial \sigma_\tau}{\partial \sigma} \\ \frac{\partial h_t}{\partial \sigma} &= \frac{\partial h_{t-1}}{\partial \sigma} - \frac{2\rho_t}{\sigma_t} \frac{\partial \sigma_t}{\partial \sigma} \end{aligned}$$

987 if the ctCDP is utilized. Because  $\mathcal{I}_t$  is always centered around  
 988 some real value  $t_0$  for  $h_{t_0} = 0$  (by continuous approxima-  
 989 tion), we may assume  $\frac{\partial \mathcal{I}_t}{\partial h_t}$  of different signs on the different  
 990 sides of  $t_0$ .

Let the gradient be zero and rearrange the variables.

$$\frac{1}{Z_T} \sum_{t=1}^T \mathcal{I}_t \partial f_t = - \frac{1}{Z_T} \sum_{t=1}^T (f_t - \hat{F}) \partial \mathcal{I}_t \quad (21)$$

991 where  $Z_T = \sum_{t=1}^T \mathcal{I}_t$ . Let  $T_\rho$  be the integer such that  $h_{T_\rho} =$   
 992  $\xi > 0$  and  $h_{T_\rho-1} = \xi' = \xi - \rho_T < 0$ . Denote the left-hand-  
 993 side and right-hand-side of Eq. (21) as lhs and rhs.

We restrict the non-zero range of  $\mathcal{I}_t$  within  $[T_\rho - 1, T_\rho]$ . Then,

$$\begin{aligned} \mathcal{I}_{T_\rho} &= 1 - \gamma h_{T_\rho} = 1 - \gamma \xi \\ \mathcal{I}_{T_\rho-1} &= 1 + \gamma h_{T_\rho-1} = 1 + \gamma \xi' = 1 + \gamma \xi - \gamma \rho_{T_\rho} \end{aligned}$$

whose summation is  $2 + \gamma(\xi' - \xi) = 2 - \gamma \rho_{T_\rho}$  and gradients are:

$$\partial \mathcal{I}_{T_\rho} = -\gamma \frac{\partial h_{T_\rho}}{\partial \sigma} = \gamma \sum_{\tau=1}^{T_\rho} \frac{2\rho_\tau}{\sigma_\tau} \frac{\partial \sigma_\tau}{\partial \sigma} \quad (22)$$

$$\partial \mathcal{I}_{T_\rho-1} = \gamma \frac{\partial h_{T_\rho-1}}{\partial \sigma} = -\gamma \sum_{\tau=1}^{T_\rho-1} \frac{2\rho_\tau}{\sigma_\tau} \frac{\partial \sigma_\tau}{\partial \sigma}. \quad (23)$$

Since  $\alpha = \mathcal{I}_{T_\rho-1} / (\mathcal{I}_{T_\rho-1} + \mathcal{I}_{T_\rho})$ ,

$$\begin{aligned} \text{rhs} &= - \frac{1}{Z_T} \sum_{t=1}^T (f_t - \hat{F}) \partial \mathcal{I}_t \\ &= - \frac{1}{Z_T} [(1-\alpha) \partial \mathcal{I}_{T_\rho-1} - \alpha \partial \mathcal{I}_{T_\rho}] (f_{T_\rho-1} - f_{T_\rho}) \end{aligned}$$

If  $\sigma$  is independent from the noise variables  $\nu_t$ , e.g., uniform schedule, the coefficient is a constant, i.e.,

$$C_\sigma = - \frac{1}{Z_T} [(1-\alpha) \partial \mathcal{I}_{T_\rho-1} - \alpha \partial \mathcal{I}_{T_\rho}] \quad (24)$$

based on which we can get the expectation,

$$\mathbb{E}[\text{rhs}] = C_\sigma \mathbb{E}[f_{T_\rho-1} - f_{T_\rho}]. \quad (25)$$

If  $f$  is  $L$ -smooth (Definition B.1),

$$\mathbb{E}[f_{T_\rho} - f_{T_\rho-1}] \leq \mathbb{E}[\nabla_{T_\rho-1}^\top g_{T_\rho-1} + \frac{L}{2} \|g_{T_\rho-1}\|^2].$$

Define  $\zeta_T \triangleq Lg_T + \nabla_T$ . Since  $\pi$  is just a variant of the SGD, the direction of  $g_T$  should be opposite to the  $\nabla_T$ . Therefore,

$$\zeta_t = L(\pi(\nabla_t + \sigma_t \nu_t) - \frac{\nabla_t}{L}), \quad (26)$$

which represents the difference between the projected updates ( $g_T$ ) and the gradient descent update with the step size  $\frac{1}{L}$ . A rational guess is that the  $\zeta_{T_\rho-1}$  is the residual noise noise after the denoising operation,  $\pi$ . Thus, it is rational to assume the  $\mathbb{E} \|\zeta_{T_\rho-1}\|^2$  is bounded as

$$\mathbb{E} \|\zeta_{T_\rho-1}\|^2 \leq \sigma_{\zeta_{T_\rho-1}}^2 p \quad (27)$$

for some parameter  $\sigma_{\zeta_{T_\rho-1}}$  depending on the scheduler where  $p$  is the dimension of  $\theta$ . Then

$$g_t = -\frac{1}{L} \nabla_t + \frac{1}{L} \zeta_t,$$

which leads to

$$\begin{aligned} \mathbb{E}[f_{T_\rho} - f_{T_\rho-1}] &\leq -\frac{1}{2L} \mathbb{E} \|\nabla_{T_\rho-1}\|^2 + \frac{1}{2L} \mathbb{E} \|\zeta_{T_\rho-1}\|^2 \\ &\leq -\frac{1}{2L} \mathbb{E} \|\nabla_{T_\rho-1}\|^2 + \frac{\sigma_{\zeta_{T_\rho-1}}^2 p}{2L}. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E} \|\nabla_{T_\rho-1}\|^2 &\leq \sigma_{\zeta_{T_\rho-1}}^2 p + 2L \mathbb{E}[f_{T_\rho-1} - f_{T_\rho}] \\ &= \sigma_{\zeta_{T_\rho-1}}^2 p + 2L \mathbb{E}[\text{rhs}] / C_\sigma. \end{aligned} \quad (28)$$

Thus, we complete the discussion of the rhs.

To the left-hand-side of Eq. (19), we first calculate the derivatives of the loss functions,

$$\partial f_T = \left( \sum_{t=1}^{T-1} \frac{\partial \sigma_t}{\partial \sigma} \nu_t^\top \frac{\partial g_t}{\partial \tilde{\nabla}_t} \right) \nabla_T. \quad (29)$$

Define a random variable as

$$V_T^\top \triangleq \sum_{t=1}^T \frac{\partial \sigma_t}{\partial \sigma} \nu_t^\top \frac{\partial (-g_t)}{\partial \tilde{\nabla}_t}, \quad (30)$$

where the negative sign is added because  $g_t$  is usually the opposite to the  $\tilde{\nabla}_t$ , for example,  $g_t \propto -\tilde{\nabla}_t$  in SGD. Now we substitute Eqs. (29) and (30) into Eq. (19) to obtain

$$\text{lhs} = -\alpha V_{T_\rho-2}^\top \nabla_{T_\rho-1} - (1-\alpha) V_{T_\rho-1}^\top \nabla_{T_\rho}.$$

For brevity, we rewrite the expectation as

$$\mathbb{E}[\text{lhs}] = -\alpha \mathbb{C}_{T_{\rho-1}} - (1-\alpha) \mathbb{C}_{T_{\rho}}, \quad (31)$$

where we define

$$\mathbb{C}_t = \mathbb{E}[V_{t-1}^{\top} \nabla_t] = \sum_{i=1}^p \text{Cov}(V_{t-1,i}, \nabla_{t,i}), \quad (32)$$

995 where we utilize  $\mathbb{E}V_{T-1} = 0$  because  $\nu_t$  is an i.i.d. Gaussian  
 996 random vector. Therefore,  $\mathbb{C}_T$  represents the covariance between  
 997 two vectors and will be zero only when the two vectors  
 998 are uncorrelated. Combining Eqs. (25), (28) and (31), we can  
 999 get Eq. (20).  $\square$

**Uniform scheduler.** We assume the  $\sigma_t \equiv \sigma$  where the scheduler degrades as a constant  $\sigma$ . Therefore, with Eqs. (22) and (23), we have

$$\begin{aligned} T_{\rho} &= \lceil \rho_{\text{tot}} / \rho \rceil, \quad \xi = \rho T_{\rho} - \rho_{\text{tot}} \\ \mathcal{I}_{T_{\rho}} &= 1 - \gamma \xi, \quad \mathcal{I}_{T_{\rho-1}} = 1 + \gamma \xi - \gamma \rho \\ \partial \mathcal{I}_{T_{\rho}} &= \gamma T_{\rho} \frac{2\rho}{\sigma}, \quad \partial \mathcal{I}_{T_{\rho-1}} = -\gamma (T_{\rho} - 1) \frac{2\rho}{\sigma} \end{aligned}$$

where  $\rho = \Delta^2 / 2\sigma^2$ . In addition, we have

$$\begin{aligned} f_{T_{\rho}} - \hat{F} &= \alpha (f_{T_{\rho}} - f_{T_{\rho-1}}) \\ f_{T_{\rho-1}} - \hat{F} &= -(1-\alpha) (f_{T_{\rho}} - f_{T_{\rho-1}}). \end{aligned}$$

Substitute what we have into Eq. (21) giving

$$\begin{aligned} \text{rhs} &= \frac{1}{Z_T} [\alpha T_{\rho} + (1-\alpha)(T_{\rho} - 1)] \gamma \frac{2\rho}{\sigma} (f_{T_{\rho-1}} - f_{T_{\rho}}) \\ &= C_{\sigma} (f_{T_{\rho-1}} - f_{T_{\rho}}), \end{aligned}$$

where we update the constant  $C_{\sigma}$  from Eq. (24) as

$$C_{\sigma} = \frac{T_{\rho} - (1-\alpha) \frac{2\rho\gamma}{\sigma}}{2 - \gamma\rho}.$$

1000 If  $\gamma \in (0, 2/\rho)$ , then  $C_{\sigma} > 0$ . In our implementation, the  
 1001 condition always holds since  $\gamma = 1/\rho_{\text{tot}} < 2/\rho$ .

Taking the expectations of lhs and rhs, we have:

$$-\alpha \mathbb{C}_{T_{\rho-1}} - (1-\alpha) \mathbb{C}_{T_{\rho}} = C_{\sigma} \mathbb{E}[f_{T_{\rho-1}} - f_{T_{\rho}}], \quad (33)$$

where  $\mathbb{C}_T$  is given by substituting  $\frac{\partial \sigma_t}{\partial \sigma} = 1$  into Eq. (30) and its definition, i.e.,

$$\mathbb{C}_T = \mathbb{E} \left[ \sum_{t=1}^{T-1} \nu_t^{\top} \frac{\partial(-g_t)}{\partial \nabla_t} \nabla_T \right].$$

**Analysis of batch algorithm.** Recall the objective for the batch algorithm (Eq. (8)) is

$$\mathcal{L}^{\text{aug}}(r) = \sum_{b=1}^B \frac{1}{2\mu_b} \|\hat{\rho}_b - r_b\|_2^2 + \frac{\sum_{b \in \mathcal{T}_B} \mathcal{I}_b \bar{F}_i}{\sum_{b \in \mathcal{T}_B} \mathcal{I}_b}$$

where  $Z_B = \sum_{b \in \mathcal{T}_B} \mathcal{I}_b$  and we use  $\mathcal{I}_b$  to denote  $\mathcal{I}(h_b)$ . Generally, we assume the batch budget is scheduled by parameterized model  $r(\cdot)$  or  $r$  for simplicity, e.g., LSTMs. In

addition, assume  $\Delta B = 1$ ,  $\alpha = \mathcal{I}_{B_{\rho-1}} / \mathcal{I}_{B_{\rho-1}} + \mathcal{I}_{B_{\rho}}$ . Vanishing  $\frac{\partial \mathcal{L}^{\text{aug}}}{\partial r}$  causes

$$\sum_{b=1}^B \frac{1}{\mu_b} (r_b - \hat{\rho}_b) \frac{\partial r_b}{\partial r} = - \frac{\sum_{b \in \mathcal{T}_B} (\bar{F}_b - \hat{F}) \partial \mathcal{I}_b}{Z_B}, \quad (34)$$

where  $Z_B = \sum_{b \in \mathcal{T}_B} \mathcal{I}_b$  and we define the notation  $\bar{x}_t^b = \frac{1}{m} \sum_{t \in \mathcal{B}_b} x_t$  for any variables  $x_t$  related to the step  $t$ . In addition, we need to make the gradient of batch objective be zero, i.e.,

$$0 = \frac{1}{m} \sum_{t \in \mathcal{B}_b} \frac{\partial f_t}{\partial \sigma} + \frac{1}{m\mu_b} (\hat{\rho}_b - r_b) \sum_{t \in \mathcal{B}_b} \frac{\partial \rho_t}{\partial \sigma}, \quad (35)$$

where we let the batch size,  $|\mathcal{B}_b|$ , be  $m$  for any  $b$ . If the equalities hold in Eqs. (34) and (35), we can extend non-batch utility bound, Theorem B.2, to the batch version in Theorem B.3. Compared to the non-batch result, the batch utility bound is extended by the average of steps in batches. For example,  $\mathbb{C}$  is replaced by  $\bar{\mathbb{C}}$ . 1002 1003 1004 1005 1006 1007

**Theorem B.3** (Utility bound of batch L2P protector). *If  $f$  is  $L$ -smooth and  $\sigma$  and  $r$  are independent from the noise variables  $\nu_t$ , then we have:*

$$\mathbb{E} \|\nabla_t\|^2^{B_{\rho-1}} \leq p \sigma_{\zeta_t}^{2B_{\rho-1}} - \frac{2L}{C_r} \sum_{b=1}^B \frac{\bar{\mathbb{C}}_t^b}{C_{\sigma,t}^{1-b}} \frac{\partial r_b}{\partial r}, \quad (36)$$

where  $C_{\sigma,b}$  is a constant depending on the scheduler  $\sigma$  and the batch  $b$  (Eq. (42)),  $C_r$  is a constant depending on the batch scheduler  $r$ ,  $\mathbb{C}_t$  represents the covariance between the noise and the true gradient (Eq. (32)), and  $\sigma_{\zeta_t}$  is the upper bound of reduced noise variance (Eq. (27)). 1008 1009 1010 1011 1012

*Proof.* We can easily get the derivative  $\partial \mathcal{I}_b$  based on Eqs. (22) and (23):

$$\partial \mathcal{I}_{B_{\rho}} = \gamma \sum_{b=1}^{B_{\rho}} \frac{\partial r_b}{\partial r}, \quad \partial \mathcal{I}_{B_{\rho-1}} = -\gamma \sum_{\tau=1}^{B_{\rho-1}} \frac{\partial r_{\tau}}{\partial r}. \quad (37)$$

Still, we use lhs and rhs to denote the two sides of the Eq. (34). From the non-batch analysis, we can extend Eq. (25) as

$$\mathbb{E}[\text{rhs}] = C_r \mathbb{E}[\bar{F}_{B_{\rho-1}} - \bar{F}_{B_{\rho}}] \quad (38)$$

$$C_r = -\frac{1}{Z_B} [(1-\alpha) \partial \mathcal{I}_{B_{\rho-1}} - \alpha \partial \mathcal{I}_{B_{\rho}}] \quad (39)$$

where we still assume  $r(\cdot)$  is independent from the private noise which makes  $C_r$  constant. 1013 1014

Consider the case when the  $f$  is  $L$ -smooth. Thus,

$$\mathbb{E}[f_t - f_{t-m}] \leq \mathbb{E}[\nabla_{t-m}^{\top} g_{t-m} + \frac{L}{2} \|g_{t-m}\|^2],$$

for all  $t$  in  $\mathcal{B}_{B_{\rho}}$ . Averaging over  $t$ , we get:

$$\begin{aligned} \mathbb{E}[\bar{F}_{B_{\rho}} - \bar{F}_{B_{\rho-1}}] &\leq \mathbb{E} \left[ \frac{1}{m} \sum_{t \in \mathcal{B}_{B_{\rho-1}}} \nabla_t^{\top} g_t + \frac{L}{2} \|g_t\|^2 \right] \\ &\leq -\frac{1}{2L} \frac{\overline{\mathbb{E} \|\nabla_t\|^2}^{B_{\rho-1}}}{\mathbb{E} \|\nabla_t\|^2} + \frac{p \sigma_{\zeta_t}^{2B_{\rho-1}}}{2L}, \end{aligned}$$

where we make use of  $\zeta_t$  defined in Eq. (26) and its bound in Eq. (27). Combine this with Eqs. (34) and (38) to get

$$\mathbb{E}\|\nabla_t\|^{2B\rho-1} \leq \frac{1}{p\sigma_{\zeta_t}^{2B\rho-1}} + \frac{2L}{C_r} \sum_{b=1}^B \frac{1}{\mu_b} (r_b - \hat{\rho}_b) \frac{\partial r_b}{\partial r}. \quad (40)$$

To find the value of  $\frac{1}{\mu_b}(r_b - \hat{\rho}_b)$ , we need to use Eq. (35) which gives:

$$\frac{1}{\mu_b}(r_b - \hat{\rho}_b) = \frac{1}{mC'_{\sigma,b}} \sum_{t \in \mathcal{B}_b} \frac{\partial f_t}{\partial \sigma}, \quad (41)$$

$$\overline{C'_{\sigma,t}}^b = \frac{1}{m} \sum_{t \in \mathcal{B}_b} \frac{\partial \rho_t}{\partial \sigma}. \quad (42)$$

According to Eqs. (29) and (30), it can be attained that

$$\mathbb{E}\left[\frac{1}{\mu_b}(r_b - \hat{\rho}_b)\right] = -\frac{\overline{C'_t}^b}{C'_{\sigma,t}} \quad (43)$$

where we modify Eqs. (30) and (32) as

$$C'_t = \mathbb{E}[V_{t-1}^b \top \nabla_t] = \sum_{i=1}^p \text{Cov}(V_{t-1,i}^b, \nabla_{t,i})$$

$$V_{T,i}^b = \sum_{t \in \mathcal{B}_b, t \leq T} \frac{\partial \sigma_t}{\partial \sigma} \nu_t \top \frac{\partial(-g_t)}{\partial \nabla_t}, T \in \mathcal{B}_b.$$

Substituting it into Eq. (40), we can get Eq. (36). This thus completes the proof.  $\square$

## B.7 Implementation details

In this section, we present the implementation details for the the projector and scheduler models. We use the Long-Short Term Memory (LSTM) networks as the backbone models.

**Constrain Noise-Scale Prediction.** To stabilize the L2P training, we explicitly constrain the range of the noise scale by using a Sigmoid activation in the scheduler. In addition, assuming the sigmoid output of the LSTM is  $y$ , we scale the output as

$$\sigma_{\min} + 2(\sigma_g - \sigma_{\min})y,$$

which is constrained in  $(\sigma_{\min}, 2\sigma_g)$ . The  $\sigma_{\text{lb}}$  is the lower bound of noise scale which derived from the upper bound of privacy budget, e.g.,  $\rho_a$  for  $(\rho, a)$ -ctCDP. The  $\sigma_g$  is estimated by uniformly scheduling budgets. Generally, we will expect the predicted  $\sigma$  is centered around  $\sigma_g$  and is not too large, e.g., larger than  $2\sigma_g$ , which will violate the utility greatly. With the constraint, the noise prediction will not fluctuate significantly.

**Coordinate-wise LSTM.** Following the implementation in (Andrychowicz et al. 2016), we share the parameters of LSTM for all optimized parameters. Therefore, a small LSTM can work for optimizing large-scale neural networks.

**Incremental Pre-training.** Training an L2P model from scratch may suffer from a great amount of DP noise such that no useful information can be learned. For simple tasks, pre-training without noise can mitigate this noise gap since

it could avoid some random optimization exploration at the beginning. For complicated tasks, e.g., deep neural networks or large-scale models, the gap between L2L models and high-privacy L2P models can still be huge. The DP noise is added without considering the scale of the model. Specifically, when the size of model parameters increases and the scale of their every coordinate decreases meanwhile, the DP noise will not change if the clipping norm is fixed. Thus, the noise is relatively amplified. Especially for deep models, the small coordinates may greatly affect the model performance and thus deep models are more sensitive to DP noise. Therefore, neither a scratch nor an L2L model could be robust enough as an initialization for the L2P model. Instead, we suggest an incremental pre-training in which the privacy scale  $\epsilon$  will incrementally increase from 0.

## C Additional experiments

### C.1 Quadratic optimization

**Setup.** To show the optimality of L2P training, we compare different algorithms by non-privately tuning them. Formally, given a fixed size of privacy budget, we train or tune a private optimizer on the quadratic optimization problems:

$$\min_{\theta} f(\theta) = \sum_{i=1}^{60} \|W_i \theta - y_i\|_2^2 + 0.001 \|\theta\|_2^2,$$

with random constants  $W_i \in \mathbb{R}^{2 \times 2}$  and  $y_i \in \mathbb{R}^2$  for  $i \in \{1, \dots, 60\}$ . We note that the tuning/meta-training is non-private such that we can see if the L2P can converge to the best private optimizer on the auxiliary datasets in comparison to baselines.

L2P-Proj (L2P with only projector) and L2P models are trained independently. Hence we can see the effect of adaptive perturbation. All optimizers are only tested on identical  $W$ ,  $y$  and initial variables. The L2P and L2P-Proj are trained with normally randomized  $W$  and  $y$  for 200 epochs after they are pre-trained without noise in the same way and the best model are selected with the lowest loss when their privacy budgets are used up in validation. The iteration numbers for SGD-Adv and L2P-Proj are chosen in range  $\{10, 20, 30, 40, 50, 60\}$  which are enough for convergence of such quadratic problems. The step size is chosen from 0.001 to 0.02 with 20 choices for SGD-Adv, while AGD uses the line search in the same range with 20 choices.

**Results.** In Fig. 8, four optimization methods are compared at the same  $(0.05, 10^{-8})$ -DP. As shown in Fig. 8, the proposed L2P converges to the zone close to but not exactly at the noise-free optimal solution. The optimization algorithms stop before reaching the optimal, because of the imposed budget constraint. Recall that the model at the optimal solution may leak sensitive information. We see that L2P guides the optimization toward the optimal by adjusting the update directions. More importantly, L2P-Proj reduces the noise magnitude, uses more step budget but converges in less steps. Because L2P-Proj has omitted no budget scheduler, it stop in a different spot. In comparison, the SGD-Adv algorithm randomly walks in a rather large region. Though AGD reduces variances relatively, it barely finds the correct optimization direction.



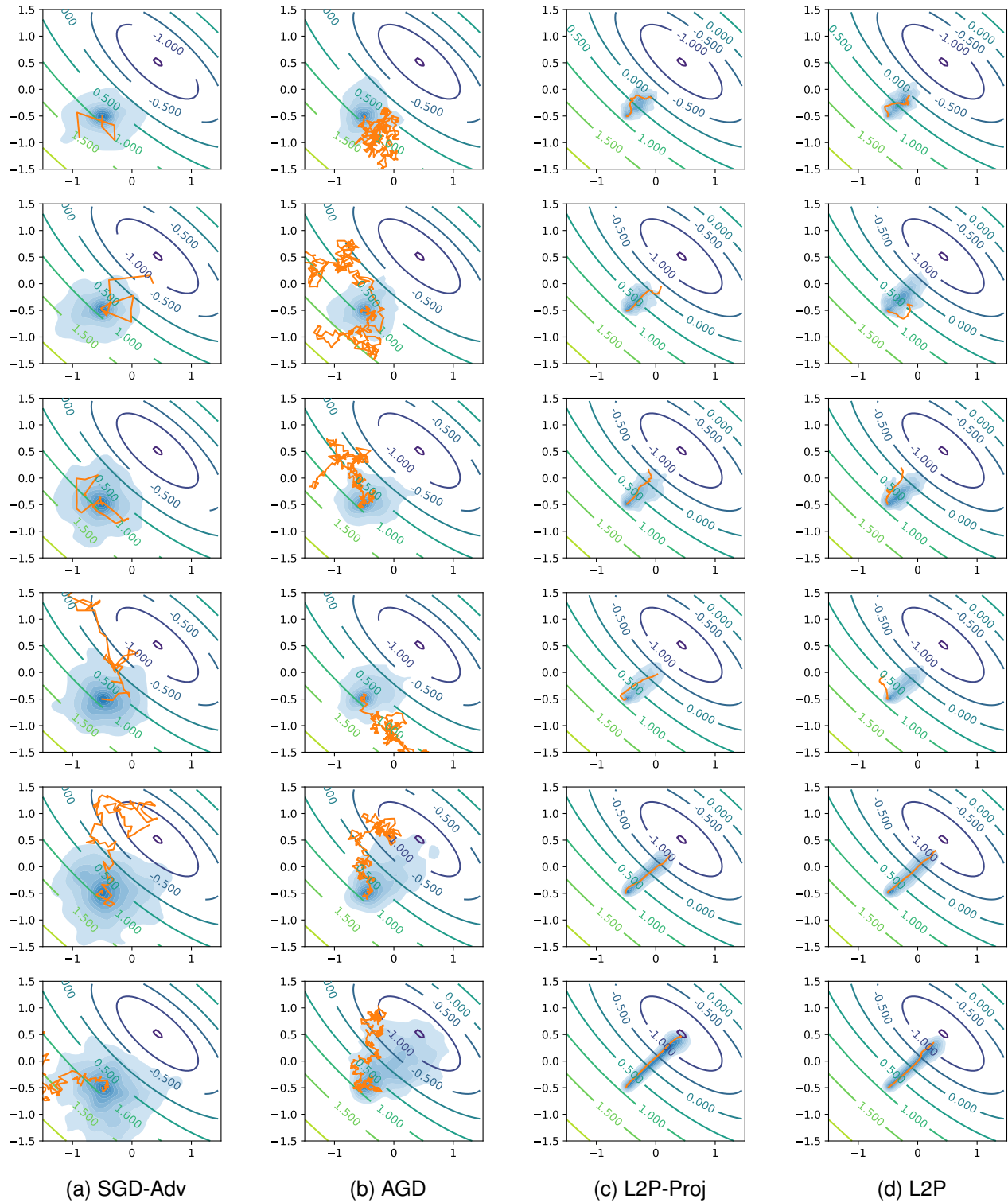


Figure 6: Comparison of the convergence  $(\epsilon, 10^{-8})$ -DP with  $\epsilon$  varying as 0.05, 0.1, 0.2, 0.4, 0.8, 1.6 from top to bottom.

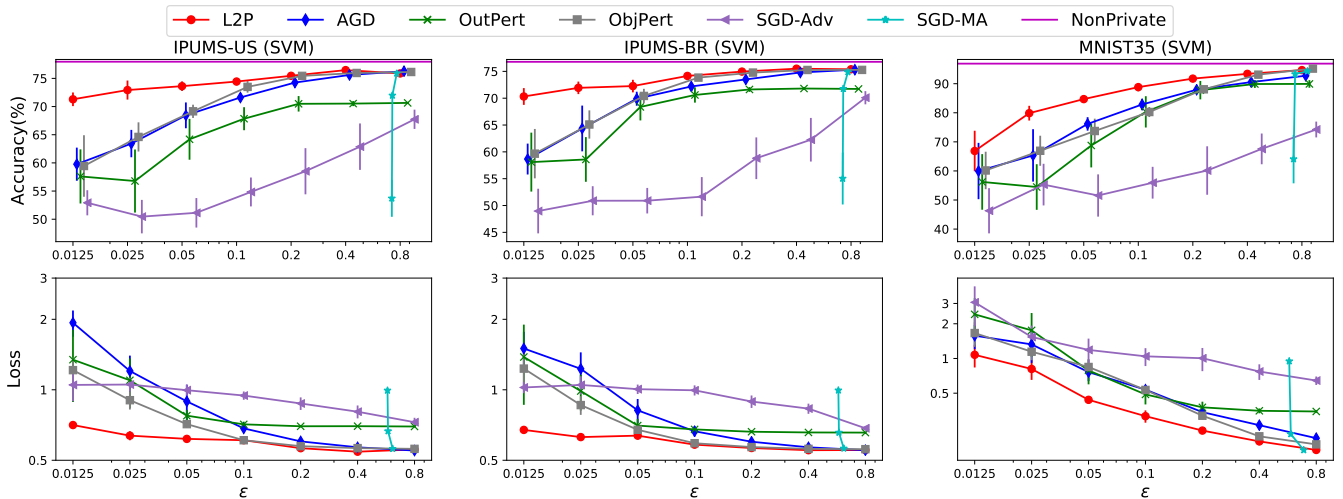


Figure 7: Test performance (top) and training loss values (bottom) by varying  $\epsilon$  of SVMs classifiers on IPUMS and MNIST datasets. The error bar presents the size of standard deviations. For better visualization of error bars, some virtual horizontal offsets are added to every point.

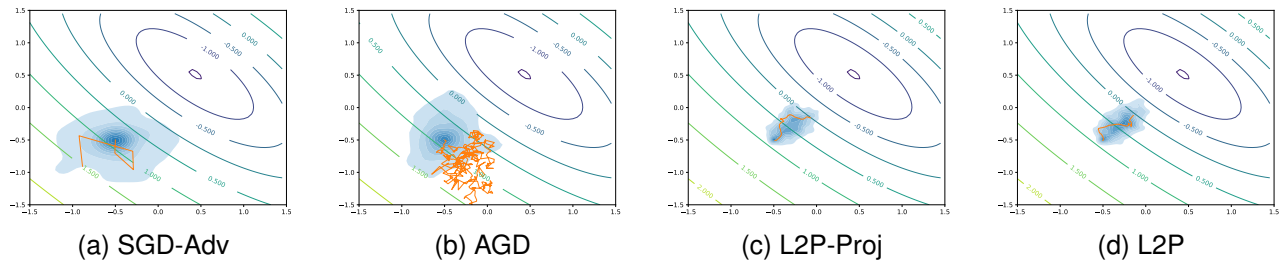


Figure 8: Comparisons of  $(0.05, 10^{-8})$ -DP algorithms on a quadratic problem. Solid contour lines illustrate the loss values. The trajectory distribution of 100 repeated optimizations are shown in blue shadowed contours. Sampled trajectories are plotted in orange.

1087 Additional quadratic optimization results for different  $\epsilon$  are  
 1088 shown in Fig. 6. Because the quadratic problem uses very few  
 1089 data and its gradients are in small scale, the optimization will  
 1090 be very sensitive to the noise. In this case, SGD-Adv rarely  
 1091 find the proper directions to go. In contrast, adaptive DP  
 1092 algorithms perform better. L2P-Proj behaves similarly to the  
 1093 L2P. However, when  $\epsilon$  gets smaller, L2P is capable to use the  
 1094 budget more efficiently such that it can converge to a better  
 1095 position. Meanwhile, L2P-Proj cannot adaptively adjust its  
 1096 step budgets which make the execution length shorter. AGD  
 1097 shows some ability to correct the noised directions but it fails  
 1098 when the privacy constraint is higher.

## 1099 C.2 Experiments of generalization to different 1100 distributions

1101 In this section, we provide additional experiments for evalu-  
 1102 ating the generalization ability of L2P.

1103 **Experiments of SVMs.** The results are reported in Fig. 7.  
 1104 The results are similar to the Logistic.

## 1105 C.3 Classification on MNIST35 datasets with 1106 non-convex objectives and varying $\epsilon$

1107 In addition to convex objectives, we also evaluate our models  
 1108 on a popular non-convex model, neural networks. The  
 1109 evaluated network includes two layers of 20 and 2 units (for  
 1110 binary classification), respectively. The layers are connected  
 1111 with sigmoid activations. The loss is computed by the cross-  
 1112 entropy function.

1113 Different from logistic and SVM models, the patterns of  
 1114 optimizing a neural network could be hard to learn for L2P.  
 1115 The first issue is that the relative magnitude of noise w.r.t. the  
 1116 gradient coordinates is enlarged when the size of the gradient  
 1117 increases. For MNIST35 images of  $28 \times 28 = 784$  pixels  
 1118 and a network with 20 units in the first layer, the number  
 1119 of connection weights could be  $20 \times 784$  which is 20 times  
 1120 of an SVM model. Since a constant L2 sensitivity, e.g., 2,  
 1121 is expected, the gradient norm will be less than 2, which  
 1122 makes each coordinate much smaller while the number of  
 1123 coordinates increases. Meanwhile, the scale of noise will  
 1124 not change for each coordinate, which means it increases in  
 1125 a relative way. This issue makes private learning methods  
 1126 hard to achieve the same utility performance under the same

1127 privacy requirement. As a result, we adjust the clipping norm  
 1128 to 2 which can slightly reduce the noise scale.

1129 The second issue has been discussed in L2L (Andrychow-  
 1130 icz et al. 2016). The optimization of the L2L or the L2P  
 1131 projector will encounter numerical issues for that the weights  
 1132 in different layers of the optimized networks have differ-  
 1133 ent magnitude. Since the optimization often focuses on the  
 1134 large elements, it will cause the optimizer barely updated  
 1135 and thus no convergence can be witnessed. As suggested by  
 1136 Andrychowicz et al. (Andrychowicz et al. 2016), either scal-  
 1137 ing the gradient by logarithm and sign mapping or scaling  
 1138 the output of the L2L by a constant can mitigate the issue.  
 1139 The former method augments the insignificant information  
 1140 in the gradients, while the second one resembles the learning  
 1141 rate such that the predicted updates will not fluctuate too  
 1142 much. To avoid that the gradient values are overwhelmed  
 1143 by noise and the L2P model absorb useless information, we  
 1144 recommend the latter method and use the scaling constant as  
 1145 0.05.

1146 In addition, optimizing small gradient coordinates with  
 1147 noise could be challenging. Incremental pre-training intro-  
 1148 duced in Appendix B.7 could reduce the hardness step by step  
 1149 through the step could be flexible. For example, when trained  
 1150 with  $\epsilon = 0.2$ , the L2P should be initialized by L2P without  
 1151 noise. When trained with  $\epsilon = 0.05$ , the L2P should be ini-  
 1152 tialized by L2P with  $\epsilon = 0.2$  instead. Other experimental  
 1153 settings follow the same principles in previous ones.

1154 For a more precise trade-off between utility and privacy  
 1155 losses, a tuning of the privacy loss coefficient is necessary. A  
 1156 recommended range of the coefficient is {500, 1000, 5000}.  
 1157 Using the training data to monitor the convergence curves  
 1158 will be helpful for choosing a proper coefficient. Meanwhile,  
 1159 the estimated number of iterations which determines the ini-  
 1160 tial privacy cost should be selected in {30, 60, 100, 400}. For  
 1161 a small  $\epsilon$ , a small iteration number will be more helpful for  
 1162 the convergence.

1163 In the last column of Fig. 3, we compare the DP  $\epsilon$  against  
 1164 the utility metrics, accuracy and loss, on the MNIST35  
 1165 dataset. The ObjPert and OutPert are excluded since it is  
 1166 designed for convex problems only<sup>6</sup>. Because some methods  
 1167 cannot converge in optimizing the network due to above-  
 1168 mentioned computation difficulties with a large noise, we  
 1169 adjust the range of  $\epsilon$ . It can be seen that for similar low pri-  
 1170 vacy conditions, L2P can train models with higher accuracy  
 1171 in most cases. The  $\epsilon$  of SGD-MA increases slowly after 0.8  
 1172 and its left boundary is given when the the number of iter-  
 1173 ation is 1, for which only a narrow range is available for  
 1174 presentation.

1175 Notably, when  $\epsilon > 0.6$ , the performance of SGD-MA is  
 1176 better than other methods except the L2P and AGD, which  
 1177 is quite distinct from previous experimental results. Because  
 1178 SGD-MA is originally designed for optimizing deep mod-  
 1179 els (Abadi et al. 2016), the moment accountant method is  
 1180 used for calculating the privacy level  $\epsilon$  is more suitable for

<sup>6</sup>Though the OutPert is claimed to be capable for nonconvex  
 problems with SGD algorithm (Zhang et al. 2017), the algorithm re-  
 quires a constant of  $\beta$ -smoothness which can not be easily obtained  
 or designed for neural networks.

Table 1: Space and time complexity on different batch sizes.  
 #unroll represents the number of unrolled steps in L2P train-  
 ing. ‘full’ means a full unrolling in one batch while 40 and  
 20 denotes the sizes of mini-batches.

#unroll	Memory (Mb)			Epoch time (sec)		
	full	40	20	full	40	20
200	474	282	250	19	21	21
400	730	282	250	41	90	89
600	1242	282	250	57	102	173
1000	1843	282	250	164	188	190
1200	2066	282	250	172	248	266
1600	2266	282	250	418	357	350

1181 mini-batch optimization. In other words, the noise scale in-  
 1182 creases slower by  $\epsilon$  using SGD-MA. Since the L2P uses the  
 1183 same batch privacy estimation, it is rational to see the L2P  
 1184 could share the benefit in optimization. When  $\epsilon \geq 0.8$ , SGD-  
 1185 MA outperforms other methods. It is because the moment  
 1186 accountant of privacy costs can lead to a tighter bound of  
 1187 compositions than  $\rho$ -zCDP used by L2P and AGD when  $\epsilon$   
 1188 increases. But  $\rho$ -zCDP can provide a more convenient and  
 1189 efficient way to compute the privacy cost explicitly. Moment  
 1190 accountant has to compute the privacy cost by iterating over  
 1191 the moment orders which is relatively slow. Though L2P does  
 1192 not outperform in accuracies when  $\epsilon > 0.8$ , it has obviously  
 1193 lower training losses. It means L2P can optimize the losses  
 1194 better within less iterations, which might be local optimal,  
 1195 though.

#### 1196 C.4 Scalability

1197 When extending the meta-training of L2P from non-  
 1198 constrained optimization to the constrained one, a critical  
 1199 issue is the scalability of the algorithm. Here we compare the  
 1200 time and space complexity of the batch and non-batch L2P  
 1201 algorithms to give a view of the issue.

1202 **Setup.** The memory usage is measured by the GPU mem-  
 1203 ory through the ‘nvidia-smi’ command on a Ubuntu 16 sys-  
 1204 tem with a TITAN X GPU and CUDA 10.1 driver. The pro-  
 1205 gram is written using TensorFlow 1.15<sup>7</sup> and allocates memory  
 1206 on need. The time is measured by the process time of one  
 1207 epoch averaged on 100 epochs. We use a 4-layer MLP and  
 1208 MNIST2 dataset for demonstration of budget-constrained op-  
 1209 timization of schedulers. Because the memory usage grows  
 1210 nonlinearly due to the TensorFlow allocation, it is slightly  
 1211 more (around 20 to 50 Mb) than the true value while the trend  
 1212 is not affected.

1213 **Results.** We empirically show the time and space com-  
 1214 plexity versus the unrolling length in Table 1. We see that  
 1215 the memory size increased quickly using the full batch while  
 1216 mini-batch does not need extra memory. Instead, mini-batch  
 1217 trade the memory with higher but acceptable time complex-  
 1218 ity. Experiments for larger networks (e.g., 128 layers) are  
 1219 included in the supplementary. Experimental results suggest  
 1220 that when a longer unrolling and larger network (e.g., 1000

<sup>7</sup><https://www.tensorflow.org/>

1221 steps for 128-layer network) are needed, slowly increasing  
1222 the batch size will be beneficial to fit the algorithm into a  
1223 limited GPU memory.

## 1224 D Discussion

1225 **a) Query Efficiency.** In comparison to AGD (Lee and Kifer  
1226 2018), the proposed L2P requires fewer times in querying the  
1227 datasets to obtain the private model, because that AGD needs  
1228 to query the dataset at each regret. To see this, we assume  
1229 that the unit privacy cost of one query is  $\epsilon$ . At each iteration,  
1230 AGD conducts two queries, including one query for objective  
1231 and the other for the gradient, if no regrets occur. Once one  
1232 regret occurs, at least one additional query is required. On the  
1233 other hand, in L2P the regret query happens in the training  
1234 process of protector on auxiliary learning tasks: when a bad  
1235 gradient  $g_t$  causes a lower loss at  $L_t$ , the effects will be back-  
1236 propagated to the LSTM cells, as shown in Fig. 1. Since  
1237 there are no privacy concerns in training the protector, the  
1238 back-propagation is more accurate than the random objective  
1239 queries used in AGD.

1240 **b) From Noised Model Training to Optimizer Training.**  
1241 In many learning algorithms, the noise-injected training, e.g.,  
1242 dropout training (Wager, Wang, and Liang 2013), has shown  
1243 to be a useful way to improve the robustness or generaliza-  
1244 tion of an algorithm. Especially if there are infinitely many  
1245 additional noised samples for training, the classification per-  
1246 formance can be improved against specific noise test envi-  
1247 ronment and both in linear space (Maaten et al. 2013) and in  
1248 nonlinear one (Hong, Chen, and Lin 2018). A critical differ-  
1249 ence between traditional noise-gradient-based DP algorithm  
1250 and noised training is the number of noised samples in noised  
1251 training or gradients in DP<sup>8</sup>. Because the constraint of pri-  
1252 vacy budget, the allowed training step is limited. In other  
1253 words, the number of noised gradients is far away from in-  
1254 finity. Thus, the DP training can only result in a degraded  
1255 model.

1256 Since, in DP, the noised component is the gradient which  
1257 is the input of an optimizer, we propose to improve the opti-  
1258 mizer by training it with noise. It is a direct extension of the  
1259 noised training except that we also train the noise variance  
1260 which is related to the privacy budget.

1261 **c) The Denoising Effect of Utility Projector** The projector  
1262 in L2P is a denoising post-processing step which does *not*  
1263 expose the original data, though. The guarantee is given in  
1264 Lemma A.4. Denoising is not new in this area which has been  
1265 studied in different directions. Recently, Balle and Wang en-  
1266 hanced the one-time query utility on Gaussian mechanism  
1267 by calibration and statical denoising (Balle and Wang 2018).  
1268 They proved that a scaling factor on the query result could  
1269 lead to a smaller expected distance between the private output  
1270 and the original one. Though their method is the analytic noi-  
1271 sing mechanism, it lacks necessary precise composition theory  
1272 for multiple queries (e.g., a learning algorithm) in compari-  
1273 son to their baseline moment accountants (Abadi et al. 2016).  
1274 Earlier, Barak et al. (Barak et al. 2007) and Hay et al. (Hay  
1275 et al. 2009) show that accurate estimation can be achieved

<sup>8</sup>Since noised samples can lead to noised gradients, we put them  
in approximately equivalent position here.

1276 by enforcing table releases and graph degree sequences to be  
1277 consistent. Karwa et al. make use of the knowledge of the  
1278 noise distribution to efficiently infer a DP graph. In addition,  
1279 the idea integrating prior into the Bayesian inference from pri-  
1280 vate outputs is formulated in (Williams and Mcsherry 2010).  
1281 Bernstein et al. use Expectation-Maximization to denoise the  
1282 parameter of a class of probabilistic graphical model (Bern-  
1283 stein et al. 2017). When a target solution is sparse, it is also  
1284 possible to project linear regression model to a known  $l_1$ -ball  
1285 which improves the resultant error.

1286 Among these work, Balle and Wang’s work (Balle and  
1287 Wang 2018) and Lee and Kifer’s work (Lee and Kifer 2018)  
1288 is the first to adaptively perturb the outputs. Balle and Wang  
1289 chose to scale the outputs with a factor adapted to the size  
1290 of private outputs. This idea is also reflected in our adaptive  
1291 perturbation where the step noise variance is adjusted accord-  
1292 ing to the private gradient norm. Differently, the variance is  
1293 adaptively calibrated according to an additional query to an  
1294 alternative objective. Also, this is leveraged in our method  
1295 while the objective query happens in auxiliary training before  
1296 a private execution.

1297 **d) Protecting L2P Training Data.** When there are very dif-  
1298 ficult learning tasks and hard to identify public auxiliary  
1299 learning tasks, one may want to use some private data for  
1300 auxiliary learning, which may cause privacy concerns when  
1301 using protector in the sensitive learning. In such a case, the  
1302 training of L2P protector should also be done in a private  
1303 learning setting, e.g., perturbing the gradients or objective  
1304 functions through classical privacy-preserving algorithms.

1305 **e) Choice of Auxiliary Tasks.** The L2P protector as well  
1306 as the learning-to-learn (Andrychowicz et al. 2016) are in  
1307 fact performing transfer learning methods that gain gradient  
1308 knowledge from auxiliary tasks and apply to a target learning  
1309 task, with and without privacy consideration respectively. We  
1310 see from our experiments that even though arbitrary choices  
1311 of auxiliary tasks can deliver promising protectors, more  
1312 relevant ones can further bring significant performance gains.  
1313 This points out an important direction for future work, i.e.,  
1314 how to quantify the task relatedness so we can use high-  
1315 performance protectors for a given learning task.

1316 **f) The availability of a public auxiliary dataset similar  
1317 enough to the private one.**

1318 Prior than our paper, public dataset has been suggested for  
1319 tuning hyper-parameters of private learning algorithms (Wu  
1320 et al. 2017). However, they did not state how to access the  
1321 public data and the affect of using different auxiliary datasets.  
1322 Our method extend the setting for practical purpose. In prac-  
1323 tice, choosing public auxiliary dataset may not be a trivial  
1324 work which greatly affect the performance. Here, we show the  
1325 affects in experiments and with some primitive criteria, we  
1326 can select useful auxiliary dataset easily. More complicated  
1327 methods could be developed based on our primitive settings.  
1328 For example, use cross validation to verify the effectiveness  
1329 of the auxiliary datasets and extract more non-private infor-  
1330 mation from the target private datasets for accurate auxiliary  
1331 dataset selection.

1332 More reasons can support the usage of auxiliary datasets in  
1333 private learning. First off, the availability of auxiliary datasets  
1334 is the main assumption of this paper and however this is a



1335 rather common assumption used by other lines of work, such  
1336 as learning-to-learn (L2L), where the learning trajectories  
1337 from other tasks are leveraged. **2)** Secondly, for most learn-  
1338 ing tasks in real-world there *are* similar publicly available  
1339 datasets, such as electronic medical records or computer vi-  
1340 sion tasks, on which we can construct auxiliary learning  
1341 tasks. **3)** Moreover, the proposed L2P framework is learning  
1342 momentum experiences from other optimization problems,  
1343 instead of heavily relying on similar datasets, we therefore  
1344 can leverage a wide spectrum of auxiliary optimization tasks  
1345 *of the same class*. For example, a quadratic programming  
1346 (QP) task may benefit from optimization procedures of many  
1347 other QP, sometimes even a random QP problem of the same  
1348 size according to our empirical study (Fig. 2). **4)** To evaluate  
1349 the influence of the choice of auxiliary datasets, an experi-  
1350 ment comparing different subsets of MNIST classes is con-  
1351 ducted in Fig. 4. The experiment is constructed to simulate  
1352 the scenario that both the auxiliary and protected datasets are  
1353 used for binary classification task with same losses. It turns  
1354 out that visually similar class sets, e.g.,  $\{4, 6\}$  (auxiliary) to  
1355  $\{3, 5\}$  (protected), yields better accuracies while less similar  
1356 ones still show performance above the best baseline.