

Differentially Private Stochastic Linear Bandits: (Almost) for Free

Osama A. Hanna* Antonious M. Girgis* Christina Fragouli Suhas Diggavi

University of California, Los Angeles, USA
Email:{ohanna, amgirgis, christina.fragouli, suhasdiggavi}@ucla.edu

Abstract

In this paper, we propose differentially private algorithms for the problem of stochastic linear bandits in the central, local and shuffled models. In the central model, we achieve almost the same regret as the optimal non-private algorithms, which means we get privacy for free. In particular, we achieve a regret of $\tilde{O}(\sqrt{T} + \frac{1}{\epsilon})$ matching the known lower bound for private linear bandits, while the best previously known algorithm achieves $\tilde{O}(\frac{1}{\epsilon}\sqrt{T})$. In the local case, we achieve a regret of $\tilde{O}(\frac{1}{\epsilon}\sqrt{T})$ which matches the non-private regret for constant ϵ , but suffers a regret penalty when ϵ is small. In the shuffled model, we also achieve regret of $\tilde{O}(\sqrt{T} + \frac{1}{\epsilon})$ while the best previously known algorithm suffers a regret of $\tilde{O}(\frac{1}{\epsilon}T^{3/5})$. Our numerical evaluation validates our theoretical results.

1 Introduction

Stochastic linear bandits offer a sequential decision framework where a learner interacts with an environment over rounds, and decides what is the optimal (from a potentially infinite set) action to play so as to achieve the best possible reward (minimize her regret). In particular, at each round, the learner may take into account all past rewards and actions to decide the next action to play, and in return receive a new reward. This model has been widely adopted both in theory but also in a number of applications, including recommendation systems, health, online education, and resource allocation [1–4]. Motivated by the fact that many of these applications are privacy-sensitive, in this paper we explore what is the performance in terms of regret we can achieve, if we are constrained to use a privacy-preserving stochastic linear bandit algorithm.

In particular, in this paper we aim to design algorithms that preserve the privacy of the rewards, from an adversary that can observe all actions that the learner plays. We assume that the learner is connected through a secure communication channel with clients, who play the requested actions. For example, the central learner may make restaurant recommendations to mobile devices, may regulate the operation of on-body sensors in senior living communities, may decide what educational exercises to provide to students, or what jobs to allocate to workers. The actions the clients play - what restaurant is visited, which sensor is activated, what is the exercise solved, what is the job performed - may be naturally visible especially in public environments. What we care to protect are the rewards, that may capture private information, such as personal preferences in recommendation systems, health indices in online health, performance in online education, and income gained in resource allocation. Our goal is to design algorithms that preserve the privacy of the rewards, while still (almost) achieve the same regret as the traditional algorithms that do not take privacy into consideration.

We do so for three different setups, depicted in Figure 1, in each case measuring the privacy using Differential Privacy (DP) measures [5, 6]. In the **central DP model**, the learner is a trusted server. The server employs a DP mechanism on aggregates of the reward realizations she collects, to ensure that the actions do not reveal information on the rewards. In the **local DP model**, the learner is an

*The first and second authors made equal contribution.

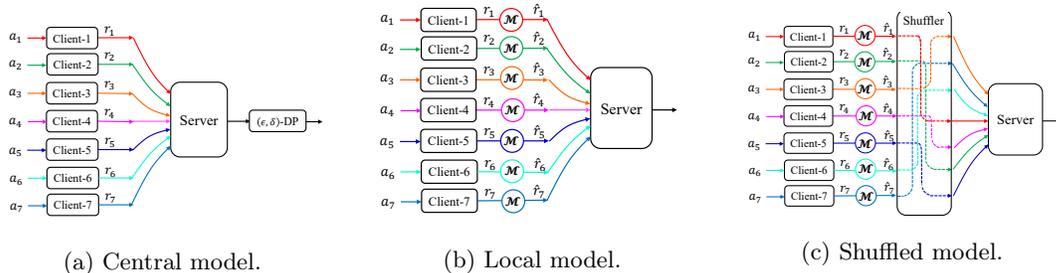


Figure 1: In case (a) the server is trusted, and we ensure that the publicly observable actions maintain privacy of the rewards. In (b) and (c) we maintain privacy from an untrusted server.

untrusted server. The clients provide privatized rewards to the server, who then uses this noisy input to decide her next actions. In the **shuffled model**, the learner is still an untrusted server, but now a trusted node, that can act as a relay in the communication between the clients and the server, serves as a shuffler, and can randomly permute the privatized rewards before making them available to the server. A shuffler offers a privacy-amplification mechanism that has recently become popular in the literature, as it is easy to implement (simply takes a set of inputs and randomly permutes them), and may enable better privacy-regret performance [7–11].

Our main contributions are as follows.

- For the **central DP model**, we design an algorithm that guarantees ϵ -DP (see Definition 1 in Section 2) and achieves regret that matches existing lower bounds. In particular, over T rounds, it achieves regret $R_T = O\left(\sqrt{T \log T} + \frac{\log^2 T}{\epsilon}\right)$ w.h.p., which is optimal within a $\log T$ factor: a lower bound of $O(\sqrt{T})$ is proven in [12] for non-private linear bandits, while a lower bound of $O(\frac{\log T}{\epsilon})$ is shown in [13] for ϵ -DP linear bandits. Note that for $\epsilon \approx 1$ (perhaps the most common case) the dominant term $O(\sqrt{T \log T})$ matches the regret of the best known algorithms for the non-private case (eg., LinUCB [12, 14]), and hence, we get privacy for free.
- For the **local DP model**, we design an algorithm that guarantees ϵ_0 -LDP (see Definition 2 in Section 2) and achieves regret $R_T = O\left(\sqrt{T \log(T)}/\epsilon_0\right)$ w.h.p.; this regret matches the non-private regret for constant ϵ_0 , but suffers a regret penalty when ϵ_0 is small. Although our algorithm does not improve the regret order as compared to the best-known algorithm for private (contextual) linear bandits in [15], it offers an alternative approach that serves as a foundation for the shuffled case.
- For the **shuffled model**, we leverage the help of a trusted shuffler to ensure both that the output of each client satisfies ϵ_0 -LDP and that the output of the secure shuffler satisfies ϵ -DP requirements. Our algorithm achieves regret $R_T = O\left(\sqrt{T \log(T)} + \frac{\log(T)}{\epsilon}\right)$ w.h.p. that matches the regret of the best non-private algorithms, same as the central model. Furthermore, our algorithm outperforms the best known algorithm for private (contextual) linear bandits in [16, 17] that use shuffling.

Our results are summarized in Table I, where we also provide known results in the literature (see also discussion next). To the best of our knowledge, in all three cases, our algorithms achieve the best currently known results for private linear bandits, significantly improving from the previously best known results in the case of the central and shuffled model, and closely matching in some cases existing lower bounds.

Our Work vs. Related Work. Differential Privacy (DP) algorithms have been proposed for the generic multi-armed bandits (MAB) problems [18–20], yet these algorithms would not work well for linear bandits, as linear bandits allow for an infinite set of actions while generic MAB have a regret that increases with the number of actions. Closer to ours is work on DP for contextual linear bandits [13, 15, 16, 21]; indeed, linear bandits can be viewed as (a special case of) contextual linear bandit setup with a single context. The work in [13] considers contextual linear bandits with DP in a centralized setting and propose an algorithm that achieves a regret of $\tilde{O}(\sqrt{T}/\epsilon)$. This does not match the best known lower bound for the centralized setting of $\Omega(\sqrt{T} + \log(T)/\epsilon)$ [13]. Our work achieves the lower bound of $\Omega(\sqrt{T} + \log(T)/\epsilon)$ up to logarithmic factors for the special case of stochastic linear bandits. The work in [21] considers contextual linear bandits with LDP, where the contexts can be adversarial. The work proposes an algorithm that achieves a regret of $\tilde{O}(T^{3/4}/\epsilon_0)$ and conjectures that the regret is optimal up to a logarithmic factor. The

Algorithm	Regret Bound	Context	Privacy Model	
			Central DP	Local DP
Central DP [13]	$\tilde{O}\left(\frac{\sqrt{T}}{\epsilon}\right)$	Adversarial	(ϵ, δ)	N/A
LDP [21]	$\tilde{O}\left(\frac{T^{3/4}}{\epsilon_0}\right)$	Adversarial	$(\epsilon = \epsilon_0, \delta)$	(ϵ_0, δ)
LDP+shuffling [16]	$\tilde{O}\left(\frac{T^{2/3}}{\epsilon^{1/3}}\right)$	Adversarial	(ϵ, δ)	$(\epsilon_0 = \epsilon^{2/3}T^{1/6}, \delta)$
LDP [15]	$\tilde{O}\left(\frac{\sqrt{T}}{\epsilon_0}\right)$	Stochastic	$(\epsilon = \epsilon_0, \delta)$	(ϵ_0, δ)
Central DP (Theorem 1)	$\tilde{O}\left(\sqrt{T} + \frac{1}{\epsilon}\right)$	Free	$(\epsilon, 0)$	N/A
LDP (Theorem 2)	$\tilde{O}\left(\frac{\sqrt{T}}{\epsilon_0}\right)$	Free	$(\epsilon = \epsilon_0, 0)$	$(\epsilon_0, 0)$
LDP+shuffling(Theorem 3)	$\tilde{O}\left(\sqrt{T} + \frac{1}{\epsilon}\right)$	Free	(ϵ, δ)	$(\epsilon_0 = \epsilon T^{1/4}, 0)$

Table 1: Upper part: known results. Lower part: our results. The \tilde{O} notation hides the dependencies on the dimension d , privacy parameter δ and log factors.

authors in [15] consider a special case, where the contexts are generated from a distribution, and propose a method that achieves a regret of $\tilde{O}(\sqrt{T}/\epsilon_0)$ under certain assumptions on the context distribution. Our algorithm for the local model achieves the same regret order using an alternative method. The works in [16,17] consider contextual linear bandits in the shuffled model where the best known algorithm achieves a regret of $\tilde{O}(T^{3/5})$. Our proposed algorithms achieve a regret of $\tilde{O}(\sqrt{T} + 1/\epsilon)$, matching the information theoretic lower bound in [13], for stochastic linear bandits in the shuffled model. A summary of the best results for DP contextual linear bandits and our results is presented in Table 1.

Paper organization. We present the problem formulation in Section 2. We design and analyze privacy-preserving linear bandit algorithms for the central model in Section 3, for the local model in Section 4 and for the shuffled model in Section 5. We provide numerical results in Section 6.

2 Notation and Problem Formulation

Stochastic linear bandits. In stochastic linear bandits a learner interacts with clients over T rounds by taking a sequence of decisions and receiving rewards. In particular, at each round $t \in [T]$, the learner plays an action a_t from a set $\mathcal{A} \subset \mathbb{R}^d$ and receives a reward $r_t \in \mathbb{R}$. The reward r_t is a noisy linear function of the action, i.e., $r_t = \langle \theta_*, a_t \rangle + \eta_t$, where $\langle \cdot \rangle$ denotes inner product, η_t is an independent zero-mean noise and $\theta_* \in \mathbb{R}^d$ is an unknown parameter vector. The goal of the learner is to minimize the total regret over the T rounds, which is calculated as:

$$R_T = T \max_{a \in \mathcal{A}} \langle \theta_*, a \rangle - \sum_{t=1}^T \langle \theta_*, a_t \rangle. \quad (1)$$

The regret captures the difference between the reward for the optimal action and the rewards for the actions chosen by the learner. The basic approach in all algorithms is to play actions that enable the learner to learn θ_* well enough to identify a (near) optimal action. The best known algorithms (for example, LinUCB [12, 14]) achieve a regret of order $O(\sqrt{T \log T})$, which is the best we can hope for (matches existing lower bounds [12]).

In this paper, we make the following standard assumptions (see, e.g., [13, 14]).

Assumption 1. We consider stochastic linear bandits with:

1. *Sub-gaussian noise:* $\mathbb{E}[\eta_{t+1} | \mathcal{F}_t] = 0$ and $\mathbb{E}[\exp(\lambda \eta_{t+1}) | \mathcal{F}_t] \leq \exp(\frac{\lambda^2}{2}) \forall \lambda \in \mathbb{R}$, where $\mathcal{F}_t = \sigma(a_1, r_1, \dots, a_t, r_t)$ is the σ -field summarizing the information available before round t .
2. *Bounded actions:* $\|a\|_2 \leq 1 \forall a \in \mathcal{A}$.
3. *Bounded unknown parameter:* $\|\theta_*\|_2 \leq 1$.
4. *Bounded rewards:* $|r_t| \leq 1$.

Privacy Goal and Measures. Our goal in this paper is to achieve the minimum possible regret in (1) while preserving privacy of the rewards $\{r_t\}_{t \in [T]}$ (as discussed in Section 1 the rewards can represent sensitive information of the clients). To measure privacy, we use the popular central and local differential privacy definitions that we provide for completeness next. For simplicity, we assume that a different client plays each action (e.g., visits a recommended restaurant).

Differential Privacy (DP). We say that two sequences of rewards $\mathcal{R} = (r_1, \dots, r_T)$ and $\mathcal{R}' = (r'_1, \dots, r'_T)$ are neighboring if they differ in a single reward, i.e., there is a round $t \in [T]$ such that $r_t \neq r'_t$, but $r_j = r'_j$ for all $j \neq t$. To preserve privacy, we use a randomized mechanism \mathcal{M} designed for stochastic linear bandits, that observes rewards and outputs publicly observable actions.

Definition 1. ([5, 6]): A randomized mechanism \mathcal{M} for stochastic linear bandits is said to be (ϵ, δ) Differentially Private ((ϵ, δ) -DP) if for any two neighboring sequences of rewards $\mathcal{R} = (r_1, \dots, r_T)$ and $\mathcal{R}' = (r'_1, \dots, r'_T)$, and any subset of outputs $\mathcal{O} \subset \mathcal{A}^T$, \mathcal{M} satisfies:

$$\Pr[\mathcal{M}(\mathcal{R}) \in \mathcal{O}] \leq e^\epsilon \Pr[\mathcal{M}(\mathcal{R}') \in \mathcal{O}] + \delta. \quad (2)$$

When $\delta = 0$, we say that the mechanism \mathcal{M} is pure differentially private (ϵ -DP). The DP mechanisms maintain that the distribution on the output of the mechanism does not significantly change when replacing a single client with reward r_t with another client with reward r'_t . Thus, the adversary observing the output of the DP mechanism does not infer the clients rewards.

Local Differential Privacy (LDP). If the central learner is untrusted, we need a local private mechanism \mathcal{M} whose output is all the information available to the central learner. We denote the range of the output of the local mechanism by \mathcal{Z} .

Definition 2. ([22]) A randomized mechanism $\mathcal{M} : [-1, 1] \rightarrow \mathcal{Z}$ is said to be (ϵ_0, δ_0) Local Differentially Private ((ϵ_0, δ_0) -LDP) if for any rewards r_t and r'_t , and any subset of outputs $\mathcal{O} \subset \mathcal{Z}$, the algorithm \mathcal{M} satisfies:

$$\Pr[\mathcal{M}(r_t) \in \mathcal{O}] \leq e^{\epsilon_0} \Pr[\mathcal{M}(r'_t) \in \mathcal{O}] + \delta_0. \quad (3)$$

Similar to the DP definition, we say that \mathcal{M} is pure locally differentially private (ϵ_0 -LDP) when $\delta_0 = 0$. Observe that the input of the LDP mechanism is a single reward, and hence, each client preserves privacy of her observed reward r_t , even if the adversary knows what is the action she plays and observes a function of her reward.

System Model. We consider three different models for private stochastic linear bandits. In all three cases, our setup is that of a learner, who asks clients to play publicly observable actions, and collects the resulting rewards using a secure communication channel (see Figure 1). The models differ on whether the learner is a trusted or untrusted server, and whether a shuffler is available or not. A shuffler simply performs a random permutation on its input.

1) Central DP model: The learner is a **trusted server** who can collect the clients' rewards and take actions. Thus, the trusted server can apply a DP mechanism (see Definition 1) to preserve the privacy of the collected rewards against any adversary observing the actions of the clients.

2) LDP model: The learner is an **untrusted server**. Hence, each client needs to privatize her own reward by applying an LDP mechanism (see Definition 2) before sending it to the untrusted server. The server takes decisions on next actions using the collected privatized rewards.

3) Shuffled model: Similar to the LDP model, the learner is an **untrusted server**. However, we consider that there exists a **trusted shuffler** that collects the LDP responses of the clients and randomly permutes them before passing them to the server, see Figure 1.

3 Stochastic Linear Bandits with central DP

In this section we consider the case where the learner is a trusted server. We present an algorithm that offers ϵ -DP (see Definition 1) for stochastic linear bandits, with no regret penalty: we achieve the same order regret performance as the best algorithms that operate under no privacy considerations.

Main Idea. Our algorithm follows the structure of elimination algorithms: it runs in batches, where in each batch i we maintain a “good set of actions” \mathcal{A}_i , that almost surely contain the optimal one, and

gradually eliminate sub-optimal actions, shrinking the sets \mathcal{A}_i as i increases. As is fairly standard in elimination algorithms, in our case as well, during batch i , the learner plays actions in \mathcal{A}_i , calculates an updated estimate $\hat{\theta}_i$ of the unknown parameter vector θ_* , and eliminates from \mathcal{A}_i actions if their estimated reward is $2\gamma_i$ from the estimated reward of the arm that appears to be best, where γ_i is the confidence of the reward estimates.

Our new idea, that enables to make our algorithm offer ϵ -DP, is at a high level as follows. **If by playing a smaller number of distinct actions we are able to identify the optimal action, we need to overall add a smaller amount of noise to guarantee privacy than if we play a larger number of distinct actions.** Indeed, if an action a is played for n_a times, the learner, to estimate θ_* , only needs to use the sum of these n_a rewards. To offer ϵ -DP we can perturb this sum by adding independent Laplacian noise ($\text{Lap}(\frac{1}{\epsilon})$); clearly, the smaller the number of distinct actions we play, the smaller the overall amount of noise we need to add. Thus our algorithm, at each batch iteration i , plays actions from a carefully selected subset of \mathcal{A}_i , of cardinality as small as possible. The technical question we address is, starting from a continuous action space \mathcal{A} , how to select at each batch iteration a small cardinality subset that maintains the ability to identify the optimal action.

We next describe the steps in implementing this idea. Recall that our actions come from a set $\mathcal{A} \subseteq \mathbb{R}^d$, and we assume they are bounded, namely, $\|a\|_2 \leq 1, \forall a \in \mathcal{A}$ (see Assumptions 1 in Section 2).

1. Our first step is to reduce the continuous action space to a discrete action space problem. To do so, we finely discretize \mathcal{A} to create what we call a ζ -net, a discrete set of actions $\mathcal{N}_\zeta \subseteq \mathcal{A}$ such that distances are approximately preserved. Namely, for any $a \in \mathcal{A}$, there is some $a' \in \mathcal{N}_\zeta$ with $\|a' - a\|_2 \leq \zeta$. Lemma 1, proved in [23, Cor. 4.2.13], states that we can always find such a discrete set with cardinality at most $(\frac{3}{\zeta})^d + d$. As a result, all the “good sets” \mathcal{A}_i will also be discrete.

Lemma 1. (*ζ -net for \mathcal{A} [23]*) *For any set $\mathcal{A} \subseteq \{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\}$ that spans \mathbb{R}^d , there is a set $\mathcal{N}_\zeta \subseteq \mathcal{A}$ (zeta-net) with cardinality at most $(\frac{3}{\zeta})^d + d$ such that \mathcal{N}_ζ spans \mathbb{R}^d , and for any $a \in \mathcal{A}$, there is some $a' \in \mathcal{N}_\zeta$ with $\|a' - a\|_2 \leq \zeta$.*

2. We introduce the use of a **core set** \mathcal{C}_i , a subset of the actions of the set of “good actions” \mathcal{A}_i . During batch i , **the learner only plays actions in \mathcal{C}_i , each with some probability $\pi_i(a)$.** Lemma 2, proved in [Ch.21] [24], states that if \mathcal{A}_i spans some space \mathbf{R}^k , we can find a core set of size at most Bk (with B a constant) and an associated probability distribution π , so that, playing actions only from \mathcal{C}_i enables to calculate a good estimate of $\langle a, \theta_* \rangle$ for each $a \in \mathcal{A}_i$.

Lemma 2. (*Core set for \mathcal{A} [24]*) *For any finite set of actions $\mathcal{A} \subset \{x \in \mathbf{R}^d \mid \|x\|_2 \leq 1\}$ that spans \mathbb{R}^d , there is a subset \mathcal{C} of size at most Bd that spans \mathbb{R}^d , where B is a constant, and a distribution π on \mathcal{C} such that for any $a \in \mathcal{A}$*

$$a^\top \left(\sum_{\alpha \in \mathcal{C}} \pi(\alpha) \alpha \alpha^\top \right)^{-1} a \leq 2d. \quad (4)$$

Moreover, \mathcal{C} and π can be found in polynomial time in d .

3. To preserve the privacy of rewards, we **perturb the sum rewards of each action by adding Laplace noise.** Adding noise affects the confidence of the reward estimates γ (step 5 in Algorithm 1 shows that γ increases as ϵ decreases), and thus delays the elimination of bad actions and increases the regret by an additive term of $\tilde{O}(\frac{1}{\epsilon})$. Replacing a possibly large set \mathcal{A}_i with the smaller core set \mathcal{C}_i effectively decreases the cumulative noise affecting the estimate of θ_* .

Remark. The computation of \mathcal{C}, π can be formulated as a convex optimization problem with many efficient approximation algorithms available. One example is the Franke-Wolfe algorithm [24, 25] that starts with an initialization of π, π_0 and updates it according to

$$\pi_{i+1}(a) = (1 - \nu_i) \pi_i(a) + \nu_i \mathbf{1}(a = \arg \max_{\alpha \in \mathcal{A}} \|\alpha\|_{V(\pi_i)}^2), \quad (5)$$

where $V(\pi_i) = \sum_{\alpha \in \mathcal{A}} \pi_i(\alpha) \alpha \alpha^\top$ and ν_i is a step size. If π_0 is chosen to be the uniform distribution over \mathcal{A} , then in $O(d \log \log |\mathcal{A}|)$ iterations we can find a π , and $\mathcal{C} = \{a \in \mathcal{A} \mid \pi(a) \neq 0\}$ that satisfy (4). Using a more sophisticated initialization, the dependence on $|\mathcal{A}|$ in the number of iterations can be eliminated entirely and the core set is guaranteed to have size of $O(d)$.

Algorithm Pseudo-Code. Our algorithm pseudo-code in Algorithm 1, starts by initializing the good action set \mathcal{A}_1 to be an $\frac{1}{T}$ -net of \mathcal{A} according to Lemma 1. Then, the algorithm operates in batches that grow exponentially in length, where the length of batch i is approximately q^i and $q = (2T)^{1/\log T}$. In each batch i , we construct the core set \mathcal{C}_i and the associated distribution π_i \mathcal{A}_i as per Lemma 2. Each action in \mathcal{C}_i is pulled $n_{ia} = \lceil \pi(a)q^i \rceil$ times, where the length of batch i is $n_i = \sum_{a \in \mathcal{C}_i} n_{ia}$. To preserve privacy, the sum of the rewards of each action is perturbed with $\text{Lap}(\frac{1}{\epsilon})$ noise. The learner uses these privatized sum rewards to compute the least squares estimate of θ_* , $\hat{\theta}_i$. At the end of batch i the learner eliminates from \mathcal{A}_i the actions with estimated mean reward, $\langle a, \hat{\theta}_i \rangle$, that fail to be within $2\gamma_i$ from the action that appears to be best, where γ_i is our confidence in the mean estimates. After the iteration $i = \log T - 1$ is completed, the learner simply plays the action that appears to be best.

Algorithm 1 ϵ -DP algorithm for stochastic linear bandits: central model

- 1: Input: set of actions \mathcal{A} , time horizon T , and privacy parameter ϵ .
 - 2: Let \mathcal{A}_1 be a ζ -net for \mathcal{A} as in Lemma 1, with $\zeta = \frac{1}{T}$.
 - 3: $q \leftarrow (2T)^{1/\log T}$.
 - 4: **for** $i = 1 : \log(T) - 1$ **do**
 - 5: $\gamma_i \leftarrow \sqrt{\frac{4d}{q^i} \log(4|\mathcal{A}_i|T^2)} + \frac{2Bd^2 + 2d \log(4|\mathcal{A}_i|T^2)}{\epsilon q^i}$.
 - 6: For $\mathcal{A}_i \subseteq \mathbf{R}^m$, $m \leq d$, let \mathcal{C}_i be a core set of size at most Bm as in Lemma 2 and π_i the associated distribution.
 - 7: Pull each action $a \in \mathcal{C}_i$, $n_{ia} = \lceil \pi_i(a)q^i \rceil$ times to get rewards $r_{ia}^{(1)}, \dots, r_{ia}^{(n_{ia})}$.
 - 8: $\bar{r}_{ia} \leftarrow \sum_{k=1}^{n_{ia}} r_{ia}^{(k)}$, $\hat{r}_{ia} \leftarrow \bar{r}_{ia} + z_{ia} \forall a \in \mathcal{C}_i$, where z_{ia} is an independent noise that follows $\text{Lap}(\frac{1}{\epsilon})$.
 - 9: $V \leftarrow \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top$, $\hat{\theta}_i \leftarrow V^{-1} \sum_{a \in \mathcal{C}_i} \hat{r}_{ia} a$.
 - 10: $\mathcal{A}_{i+1} \leftarrow \{a \in \mathcal{A}_i \mid \langle a, \hat{\theta}_i \rangle \geq \max_{\alpha \in \mathcal{A}} \langle \alpha, \hat{\theta}_i \rangle - 2\gamma_i\}$
 - 11: Play action $\arg \max_{\alpha \in \mathcal{A}_{\log(T)-1}} \langle \alpha, \hat{\theta}_{\log(T)-1} \rangle$ for the remaining time.
-

Algorithm Performance. We next prove that Algorithm 1 is ϵ -DP and provide a bound on its regret.

Theorem 1. *Algorithm 1 is ϵ -differentially private. Moreover, it achieves a regret*

$$R_T \leq C \left(\sqrt{T \log T} + \frac{\log^2 T}{\epsilon} \right), \quad (6)$$

with probability at least $1 - \frac{1}{T}$, where C is a constant that does not depend on ϵ, T .

Proof Outline. The privacy result follows from the Laplace mechanism [6]. To bound the regret, we first argue that with probability at least $1 - \frac{1}{T}$, and for all i and all $a \in \mathcal{A}_i$, we have that $|\langle a, \hat{\theta}_i \rangle - \langle a, \hat{\theta}_* \rangle| \leq \gamma_i$. Conditioned on this event, an action with gap Δ_a is eliminated when, or before, $\gamma_i < \Delta_a/2$. Hence, all actions in batch i have gap that is at most $4\gamma_i$. The regret bound follows by summing $4\gamma_i n_i$ for all batches. The complete proof is provided in Appendix A. \square

Remark 1. We note that the high probability bound in Theorem 1 implies a bound in expectation

$$\mathbb{E}[R_T] \leq C \left(\sqrt{T \log T} + \frac{\log^2 T}{\epsilon} \right). \quad (7)$$

This is because the regret is trivially $O(T)$ and the algorithm fails with probability $\frac{1}{T}$, which overall contributes $O(1)$ to the expectation.

Remark 2. The regret in Theorem 1 is optimal up to $\log T$ factor; a lower bound of $O(\sqrt{T})$ is proven in [12] for the non-private case, while a lower bound of $\frac{\log T}{\epsilon}$ is shown in [13] for the private case.

Remark 3. We observe that the privacy parameter ϵ is typically ≈ 1 . In this case, the dominating term in the regret in (11) is $O(\sqrt{T \log T})$ which matches the regret of the best known algorithm for the non-private case (see LinUCB in [12, 14]), and hence, we get privacy for free.

¹We note that $e \leq q \leq e^2$.

4 Stochastic Linear Bandits with LDP

In this section, the learner is an untrusted server, and thus we design a linear bandit algorithm (Algorithm 2) that operates under LDP constraints.

Main Idea. As in Algorithm 1, we here also utilize a core set of actions; the difference is that, since the server is untrusted, each client privatizes her own reward before providing it to the server. Our algorithm offers an alternative approach to [15] that achieves the same regret, while using operation in batches, which may in some applications be more implementation-friendly, and also forms a foundation for the Algorithm 3 we discuss in the next section.

Algorithm 2 ϵ_0 -LDP algorithm for stochastic linear bandits: local model

- 1: Input: set of actions \mathcal{A} , time horizon T , and privacy parameter ϵ_0 .
 - 2: Let \mathcal{A}_1 be a ζ -net for \mathcal{A} as in Lemma 1, with $\zeta = \frac{1}{T}$.
 - 3: $q \leftarrow (2T)^{1/\log T}$.
 - 4: **for** $i = 1 : \log(T) - 1$ **do**
 - 5: **Client side:**
 - 6: Receive action a from the server. Play action a and receive a reward r .
 - 7: Send $\hat{r} = r + \text{Lap}(\frac{1}{\epsilon_0})$.
 - 8: **Server side:**
 - 9: Let \mathcal{C}_i be a core set for \mathcal{A}_i as in Lemma 2 with distribution π_i , and $n_{ia} = \lceil \pi_i(a)q^i \rceil$.
 - 10: Send each action $a \in \mathcal{C}_i$ to a set of n_{ia} clients to get rewards $\hat{r}_{ia}^{(1)}, \dots, \hat{r}_{ia}^{(n_{ia})}$.
 - 11: $n_i \leftarrow \sum_{a \in \mathcal{C}_i} n_{ia}$.
 - 12: $\gamma_i \leftarrow \sqrt{\frac{4d}{q^i} \log(4|\mathcal{A}_i|T^2)} + \frac{2d}{q^i \epsilon_0} \sqrt{n_i \log(4|\mathcal{A}_i|T^2)}$.
 - 13: $\hat{r}_{ia} \leftarrow \sum_{k=1}^{n_j} \hat{r}_{ia}^{(k)} \forall a \in \mathcal{C}_i$.
 - 14: $V \leftarrow \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top$, $\hat{\theta}_i \leftarrow V^{-1} \sum_{a \in \mathcal{C}_i} \hat{r}_{ia} a$.
 - 15: $\mathcal{A}_{i+1} \leftarrow \{a \in \mathcal{A}_i \mid \langle a, \hat{\theta}_i \rangle \geq \max_{\alpha \in \mathcal{A}_i} \langle \alpha, \hat{\theta}_i \rangle - 2\gamma_i\}$.
 - 16: Play action $\arg \max_{\alpha \in \mathcal{A}_{\log(T)-1}} \langle \alpha, \hat{\theta}_{\log(T)-1} \rangle$ for the remaining time.
-

Algorithm Pseudocode. Algorithm 2 operates like Algorithm 1, except for the addition of $\text{Lap}(1/\epsilon_0)$ noise for each reward individually as opposed to adding $\text{Lap}(1/\epsilon)$ to the sum of the rewards of each arm in the central model. The value of γ_i is adjusted to account for this change. **Algorithm Performance.** The following Theorem 2 presents the privacy-regret tradeoffs of the LDP stochastic bandits Algorithm 2. The proof is deferred to Appendix B and follows the same main steps as the proof of Theorem 1, but with the modified values of γ_i .

Theorem 2. *Algorithm 2 is ϵ_0 -LDP. Moreover, it achieves a regret*

$$R_T \leq C \left(1 + \frac{1}{\epsilon_0}\right) \left(d\sqrt{dT \log T}\right), \quad (8)$$

with probability at least $1 - \frac{1}{T}$, where C is a constant that does not depend on ϵ_0 and T .

Remark 4. Since the regret is trivially bounded by $O(T)$ when Algorithm 2 fails, which happens with probability $\frac{1}{T}$, we can upper bound the expected regret as

$$\mathbb{E}[R_T] \leq C \left(1 + \frac{1}{\epsilon_0}\right) \left(d\sqrt{dT \log T}\right). \quad (9)$$

Remark 5. When $\epsilon_0 > 1$, the regret R_T would be $\mathcal{O}\left(\sqrt{T} \log(T)\right)$ that matches the non-private case. However, the constants of the regret convergence are larger than that of the non-private case.

Remark 6. (Comparison to the central (ϵ, δ) -DP model.) Observe that when $\epsilon_0 < 1$, the dominating term in the regret bound is $R_T = \mathcal{O}\left(\frac{T \log(T)}{\epsilon_0}\right)$. In other words, we obtain the regret of the non-private

case divided by the LDP parameter ϵ_0 . In contrast, the central DP parameter ϵ appears as an additive term in the regret of the central model. This difference is because, in the local model noise is added on every reward, while in the central model directly on the reward aggregates; thus the noise variance of the aggregate rewards and the confidence parameter γ_i increases in the local model. In the high privacy regimes; for example, assume that $\epsilon_0 = \mathcal{O}\left(\frac{1}{T^\alpha}\right)$ for some $0 < \alpha \leq \frac{1}{2}$, we get a regret R_T of order $\mathcal{O}\left(T^{\frac{1}{2}+\alpha}\right)$ that becomes linear function of T as $\epsilon_0 \rightarrow \frac{1}{\sqrt{T}}$.

5 Stochastic Linear Bandits in the Shuffled Model

Algorithm 3 DP algorithm for stochastic linear bandits: shuffled model

- 1: Input: set of actions \mathcal{A} , time horizon T , and privacy parameters (ϵ, δ) .
 - 2: Let \mathcal{A}_1 be a ζ -net for \mathcal{A} as in Lemma 1, with $\zeta = \frac{1}{T}$.
 - 3: $q \leftarrow (2T)^{1/\log T}$.
 - 4: **for** $i = 1 : \log(T) - 1$ **do**
 - 5: **Client side:**
 - 6: Receive action a and the value n_i from the shuffler.
 - 7: Play action a and receive a reward r .
 - 8: $\epsilon_0^{(i)} \leftarrow f_{n_i, \delta}^{-1}(\epsilon)$
 - 9: Send $\hat{r} = r + \text{Lap}\left(\frac{1}{\epsilon_0^{(i)}}\right)$ to the shuffler.
 - 10: **Shuffler:**
 - 11: Let \mathcal{C}_i be a core set for \mathcal{A}_i as in Lemma 2 with distribution π_i .
 - 12: Let $n_{ia} = \lceil \pi_i(a)q^i \rceil, n_i \leftarrow \sum_{a \in \mathcal{C}_i} n_{ia}$.
 - 13: Let $\mathcal{A}_{\mathcal{C}_i} = \cup_{a \in \mathcal{C}_i} \{a\}_{l=1}^{n_{ia}}$ be a set of n_i actions where action $a \in \mathcal{C}_i$ is repeated n_{ia} times.
 - 14: Let a_1, \dots, a_{n_i} be an enumeration of $\mathcal{A}_{\mathcal{C}_i}$.
 - 15: Send action $a_{\pi(j)}$ and the value n_i to client $j, j = 1, \dots, n_i$, where π is a random permutation of $1, \dots, n_i$.
 - 16: Receive the action-reward pairs $\{(a_1, \hat{r}_{ia_1}), \dots, (a_{n_i}, \hat{r}_{ia_{n_i}})\}$, and send them to the server.
 - 17: **Server side:**
 - 18: Receive the action-reward pairs from the shuffler.
 - 19: $\gamma_i \leftarrow \sqrt{\frac{4d}{q^i} \log(4|\mathcal{A}_i|T^2) + \frac{2d}{q^i \epsilon_0^{(i)}} \sqrt{n_i \log(4|\mathcal{A}_i|T^2)}}$.
 - 20: $\hat{r}_{ia} \leftarrow \sum_{k=1}^{n_j} \hat{r}_{ia}^{(1)} \forall a \in \mathcal{C}_i$.
 - 21: $V \leftarrow \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top, \hat{\theta}_i \leftarrow V^{-1} \sum_{a \in \mathcal{C}_i} \hat{r}_{ia} a$.
 - 22: $\mathcal{A}_{i+1} \leftarrow \{a \in \mathcal{A}_i \mid \langle a, \hat{\theta}_i \rangle \geq \max_{\alpha \in \mathcal{A}} \langle \alpha, \hat{\theta}_i \rangle - 2\gamma_i\}$.
 - 23: Play action $\arg \max_{\alpha \in \mathcal{A}_{i \log(T)-1}} \langle \alpha, \hat{\theta}_{i \log(T)-1} \rangle$ for the remaining time.
-

In this section, we consider the case of an untrusted server and a trusted shuffler. We propose Algorithm 3 that (almost) achieves the same order regret as the best non-private algorithms.

Main idea. To use shuffling, we need to use an algorithm that operates over batches of actions, so as to be able to shuffle them. The use of a core set is critical to enable a selection of actions that lead to a good estimate for θ_\star . For example, if the original set \mathcal{A} contains a large number of actions along one direction in the space, but only a few actions along other directions, then pulling each action in \mathcal{A} once will not result in a good estimate of θ_\star . Use of the core set and the associated distribution π will balance such asymmetries and enable to explore multiple directions of the space for a sufficient number of times to acquire a good estimate of θ_\star .

Accordingly, we follow the same approach as in Algorithm 2 with two changes: we use a shuffler (in a manner tailored to bandits) to realize privacy amplification gains, and we adjust the amount of Laplace noise we add in each batch, depending on the batch size.

We use the trusted shuffler as follows. The actions to be played in the i th batch are shuffled by the trusted shuffler at the beginning of the batch. The shuffler asks clients to play actions in the shuffled

order. Then, at the end of the batch, the shuffler reverses the shuffling operation, associates every action with its observed LDP reward, and conveys it to the untrusted learner.²

We adjust the amount of added Laplace noise per batch as follows. To offer privacy guarantees, we want to add noise to the rewards so that the output of the shuffler is (ϵ, δ) -DP for each batch $i \in [\log(T)]$. This implies that the entire algorithm will be (ϵ, δ) -DP, since we assume that each client contributes at only one of the batches. The privacy amplification of the shuffling depends on the size of the batch (see e.g. [10, Theorem 1]); thus the larger the batch size, the less noise needs to be added to the rewards of the clients. To ensure that the output of batch i is (ϵ, δ) -DP, it is sufficient to add to each reward noise $\text{Lap}(\frac{1}{\epsilon_0^{(i)}})$, where $\epsilon_0^{(i)} \leftarrow f_{n_i, \delta}^{-1}(\epsilon)$, and n_i is the size of batch i . The function $f_{n, \delta} : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ captures privacy amplification via shuffling [10] and is defined as follows

$$f_{n, \delta}(\epsilon_0) = \log \left(1 + \frac{e^{\epsilon_0} - 1}{e^{\epsilon_0} + 1} \left(\frac{8\sqrt{e^{\epsilon_0} \log(4/\delta)}}{\sqrt{n}} + \frac{8e^{\epsilon_0}}{n} \right) \right). \quad (10)$$

Since the noise added to the rewards varies for each batch i , we modify the confidence bounds, γ_i , to reflect this. The pseudo-code is provided in Algorithm 3.

Algorithm Performance. The following theorem proves that Algorithm 3 is (ϵ, δ) -DP and provides an upper bound on its regret that matches the information theoretic lower bound for $\epsilon = \tilde{O}(\frac{1}{\sqrt{T}})$.

Theorem 3. *Algorithm 3 is (ϵ, δ) -differentially private. Moreover, for $\epsilon = O(\sqrt{\frac{\log(1/\delta)}{T}})$ it achieves a regret*

$$R_T \leq C \left(\sqrt{T \log T} + \frac{\sqrt{\log(1/\delta)} \log^{3/2} T}{\epsilon} \right), \quad (11)$$

with probability at least $1 - \frac{1}{T}$, where C is a constant that does not depend on ϵ and T .

Proof Outline. The proof of Theorem 3 is deferred to Appendix C. The privacy guarantee is proved by reducing the scheme to one that shuffles the rewards but does not shuffle the corresponding actions and using results from [10]. The regret analysis follows similar ideas as in Theorem 1 and Theorem 2.

Remark 7. Note that if we had the shuffler to simply permute the collected rewards of the clients (and not the actions) we would get no privacy gains in some cases. For example, consider the case where all actions to be pulled in a batch are unique and the action pulled by each client is known to the central learner (e.g., for MAB algorithms where the policy is a deterministic function of the history), then the learner can undo the shuffling using the action associated with each shuffled reward.

Remark 8. Algorithm 3 almost achieves the same order regret as the best non-private algorithms. Indeed, Theorem 3 proves that Algorithm 3 achieves a regret that matches the regret of the central DP Algorithm 1 for the high privacy regimes $\epsilon = O(\sqrt{\log(1/\delta)/T})$. For the low privacy regime $\epsilon > 1$, the shuffling does not offer privacy gains, $\epsilon_0^{(i)} \approx \epsilon$ for all $i \in [\log(T)]$ and the regret of Algorithm 3 is similar to the regret of Algorithm 2 of the local DP model. However, for the low privacy regime the local DP model also achieves the same regret as non-private algorithms up to constant factors (see Remark 5). Hence in both cases, Algorithm 3 achieves the same order regret as Algorithm 1 which almost matches the regret of non-private algorithms.

Remark 9. Algorithm's 3 improved regret performance over Algorithm 2 is thanks to the smaller amount of noise added to rewards. In particular, the noise added in Step 9 of Algorithm 3 has variance $\frac{2}{\epsilon_0^{(i)2}} \approx \frac{2}{n_i \epsilon^2}$ for small ϵ .

6 Numerical Results

We here present indicative results on the performance of our proposed Algorithms 1, 2 and 3; additional details and numerical evaluation plots are provided in Appendix D. We consider synthetic data generated

²We assume that the server cannot directly observe which action is played by which client, for instance due to geographical separation.

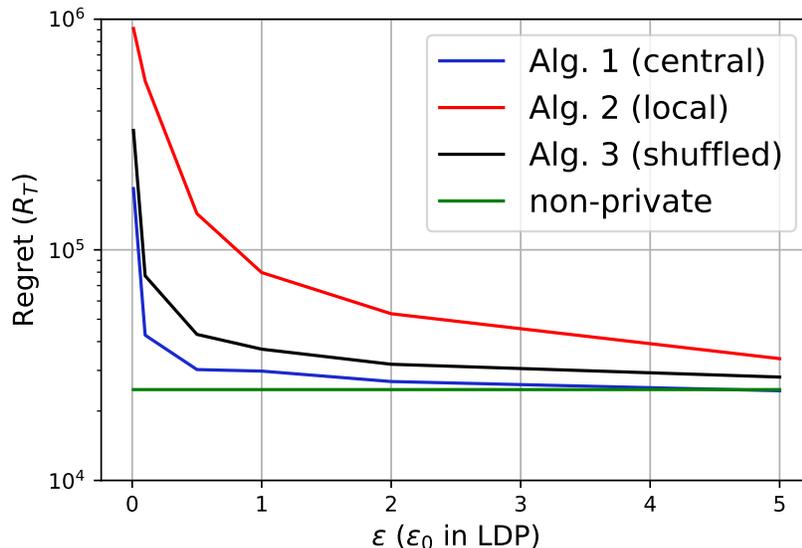


Figure 2: Regret-privacy trade-offs for stochastic linear bandits algorithms.

as follows. The set of actions \mathcal{A} contains $K = 10$ actions, where each action $a \in \mathcal{A}$ is a $d = 2$ -dimensional vector. The actions $a \in \mathcal{A}$ and the optimal parameter θ_* are generated uniformly at random from the unit ball $\mathcal{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$ (a similar method is considered in [15]). Figure 2 plots the total regret R_T over an horizon $T = 10^6$ as a function of the privacy budget (ϵ or ϵ_0). Figure 2 shows that the regret achieved by all three algorithms, Algorithm 1 (central model), Algorithm 2 (local model), and Algorithm 3 (shuffled model) converges to the regret of non-private stochastic linear bandit algorithms [24, Ch. 22] as $\epsilon \rightarrow \infty$ ($\epsilon_0 \rightarrow \infty$), albeit at different rates. As predicted from the theoretical analysis, Algorithms 1 (central) and 3 (shuffled) offer privacy (almost) for free, closely following the non-private regret.

A Regret and Privacy Analysis of The Central DP Model (Proof of Theorem 1)

A.1 Privacy Analysis

We first show that Algorithm 1 is ϵ -DP. Let $\bar{r}_i = [\bar{r}_{ia_1}, \dots, \bar{r}_{ia_{|C_i|}}]$, $\hat{r}_i = [\hat{r}_{ia_1}, \dots, \hat{r}_{ia_{|C_i|}}] = \bar{r}_i + z_i$, $z_i = [z_{ia_1}, \dots, z_{ia_{|C_i|}}]$, where $a_1, \dots, a_{|C_i|}$ is an enumeration of the elements of C_i . We construct the concatenated reward vector denoted by $\bar{r} = [\bar{r}_1, \dots, \bar{r}_{\log(T)-1}]$, and let $\hat{r} = [\hat{r}_1, \dots, \hat{r}_{\log(T)-1}] = \bar{r} + z$, $z = [z_1, \dots, z_{\log(T)-1}]$.

Now consider two neighboring sequence of rewards $\mathcal{R}, \mathcal{R}'$, that only differ in r_k, r'_k , with corresponding concatenated reward vectors \bar{r}, \bar{r}' . We notice that each reward in \mathcal{R} appears once in \bar{r} , and similarly for \mathcal{R}', \bar{r}' . Thus, we get:

$$\|\bar{r} - \bar{r}'\|_1 \leq \max_{r_k, r'_k} |r_k - r'_k| \leq 1, \quad (12)$$

where the last inequality follows from Assumption 1 with bounded rewards $|r_k| \leq 1$. Then, from [5, Theorem 3.6], \hat{r} is ϵ -DP. We notice that the output of Algorithm 1 depends on r_1, \dots, r_T only through \hat{r} . Hence, by post processing, Algorithm 1 is ϵ -DP.

A.2 Regret Analysis

We next prove the regret bound of Algorithm 1 for stochastic linear bandits.

Our analysis follows the known confidence bound technique in [26] by designing confidence intervals (in step 5) that take into consideration the privacy effect.

Let $K = (3T)^d$ be the size of the $\frac{1}{T}$ -net set $\mathcal{N}_{1/T}$ from Lemma 1. We first bound the following regret:

$$\tilde{R}_T = T \max_{a \in \mathcal{N}_{1/T}} \langle a, \theta_* \rangle - \sum_{t=1}^T \langle a_t, \theta_* \rangle, \quad (13)$$

where $a_1, a_2, \dots, a_T \in \mathcal{N}_{1/T}$. We then bound the regret R_T by showing that we only lose a constant term when we choose actions from $\mathcal{N}_{1/T}$ instead of the bigger set \mathcal{A} .

We start with a set of actions $\mathcal{A}_0 = \mathcal{N}_{1/T}$ with cardinality $|\mathcal{A}_0| = K$. Furthermore, we have $|\mathcal{A}_i| \leq |\mathcal{A}_{i-1}|$, and hence, we get $|\mathcal{A}_i| \leq K$ for all $i \in [\log(T)]$.

For given batch $i \in [\log(T)]$, let \mathcal{C}_i be the core set of \mathcal{A}_i that has at most Bd actions. At the i th batch, each action $a \in \mathcal{C}_i$ is picked n_{ia} times, where $n_{ia} = \lceil \pi_i(a)q^i \rceil$. Let \mathcal{G} be the good event $\left\{ \left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| < \gamma_i \ \forall i \in [\log T] \ \forall a \in \mathcal{A}_i \right\}$. Lemma 3 shows that the event \mathcal{G} holds with probability at least $1 - \frac{1}{T}$. In the remaining part of the proof, we condition on the event \mathcal{G} .

We first show that the best action $a_* = \arg \max_{a \in \mathcal{N}_{1/T}} \langle a, \theta_* \rangle$ will not be eliminated at any batch $i \in [\log T]$; this is because the elimination criterion will not hold for the optimal action a_* :

$$\langle a, \hat{\theta}_i \rangle - \langle a_*, \hat{\theta}_i \rangle < (\langle a, \theta_* \rangle + \gamma_i) - (\langle a_*, \theta_* \rangle - \gamma_i) \leq 2\gamma_i \quad \forall a \in \mathcal{A}_i \ \forall i \in [\log T]. \quad (14)$$

For each sub-optimal action $a \in \mathcal{A}_0$ with $\Delta_a = \langle a_* - a, \theta_* \rangle$, let i be the smallest integer for which $\gamma_i < \frac{\Delta_a}{4}$. From the triangle inequality, we get that

$$\langle a_*, \hat{\theta}_i \rangle - \langle a, \hat{\theta}_i \rangle \geq (\langle a_*, \hat{\theta}_i \rangle - \gamma_i) - (\langle a, \hat{\theta}_i \rangle + \gamma_i) = \Delta_a - 2\gamma_i > 2\gamma_i. \quad (15)$$

This implies that a will be eliminated before the beginning of batch $i + 1$. Hence, each action $a \in \mathcal{A}_{i+1}$ at batch $i + 1$ has a gap at most $4\gamma_i$. Let $n_i = \sum_{a \in \mathcal{C}_i} n_{ia} \leq Bd + q^i$ denote the total number of rounds at the i -th batch. Note that the number of batches is upper bounded by $\log T$ since $\sum_{i=1}^{\log T} q^i \geq T$. When $q^i < Bd$, the regret can be bounded by $2Bd$, and when $q^i \geq Bd$, we bound $n_i \leq 2q^i$. Thus, there is universal constants C', C such that the total regret in (13) can be bounded as

$$\begin{aligned} \tilde{R}_T &\leq 2Bd \log(T) + \sum_{i=1}^{\log T} 4n_i \gamma_{i-1} & (16) \\ &\leq 2Bd \log(T) + \sum_{i=1}^{\log T} 8q^i \left(\sqrt{\frac{4d}{q^{i-1}} \log(4KT^2)} + \frac{2Bd^2 + 2d \log(4KT^2)}{\epsilon q^{i-1}} \right) \\ &\leq C' \left(d \log(T) + d \sqrt{\log T} \sum_{i=1}^{\log T} q^{(i-1)/2} + \frac{d^2 \log^2 T}{\epsilon} \right) q \\ &\stackrel{(a)}{\leq} C' q \left(d \log(T) + d \sqrt{\log T} q^{\log T/2} + \frac{d^2 \log^2 T}{\epsilon} \right) \\ &\stackrel{(b)}{\leq} C' q \left(d \log(T) + d \sqrt{T \log T} + \frac{d^2 \log^2 T}{\epsilon} \right) \\ &\stackrel{(c)}{\leq} C \left(d \sqrt{T \log T} + \frac{d^2 \log^2 T}{\epsilon} \right), & (17) \end{aligned}$$

where step (a) follows from the sum of a geometric series and $q > 1$, step (b) uses $q = (2T)^{1/\log T}$, and step (c) follows from the facts $q \leq e^2$, $\log T = O(\sqrt{T})$.

Hence, with probability at least $1 - \frac{1}{T}$ the regret in (13) is bounded as

$$\tilde{R}_T \leq C \left(d \sqrt{T \log T} + \frac{d^2 \log^2 T}{\epsilon} \right). \quad (18)$$

Next, we bound the exact regret R_T . Observe that the first step in our Algorithm is to use the finite $\frac{1}{T}$ -net set $\mathcal{N}_{1/T}$ of actions. Thus, for any round $t \in [T]$ and any action $a \in \mathcal{A}$, there exists an action

$a' \in \mathcal{N}_{1/T}$ such that $\|a - a'\| \leq \frac{1}{T}$. As a result, we get $\langle a, \theta_* \rangle - \langle a', \theta_* \rangle \leq \|a - a'\| \|\theta_*\| \leq \frac{1}{T}$, where $\|\theta_*\| \leq 1$. Hence, there is a universal constant C such that we can bound the regret R_T as

$$\begin{aligned} R_T &= T \max_{a \in \mathcal{A}} \langle a, \theta_* \rangle - \sum_{t=1}^T \langle a_t, \theta_* \rangle \\ &= \left[T \max_{a \in \mathcal{A}} \langle a, \theta_* \rangle - T \max_{a' \in \mathcal{N}_{1/T}} \langle a', \theta_* \rangle \right] + \left[T \max_{a' \in \mathcal{N}_{1/T}} \langle a', \theta_* \rangle - \sum_{t=1}^T \langle a_t, \theta_* \rangle \right] \\ &\leq T \frac{1}{T} + \tilde{R}_T \\ &= 1 + \tilde{R}_T. \end{aligned} \quad (19)$$

Hence, with probability at least $1 - \frac{1}{T}$ the regret R_T is bounded as

$$R_T \leq C \left(d\sqrt{T \log T} + \frac{d^2 \log^2 T}{\epsilon} \right). \quad (20)$$

This concludes the proof of Theorem 1.

Lemma 3. *Let $\hat{\theta}_i$ be the least square estimate of θ_* at the end of the i th batch of Algorithm 1. Then, we have that*

$$\Pr \left[\left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| > \gamma_i \quad \forall i \in [T] \forall a \in \mathcal{A}_i \right] \leq \frac{1}{T}, \quad (21)$$

where $\gamma_i = \sqrt{\frac{4d}{q^i} \log(4KT^2)} + \frac{2Bd^2 + 2d \log(4KT^2)}{\epsilon q^i}$.

Proof. Let $\hat{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} \hat{r}_{ia} a$ be the private estimate of θ_* and $\bar{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} \bar{r}_{ia} a$ be the non-private estimate of θ_* as $\{\bar{r}_{ia}\}$ are the non-private rewards, where $V_i = \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top$. From [Chapter 21, Eqn 21.1], for each $a \in \mathcal{A}_i$, we get:

$$\Pr \left[\langle a, \bar{\theta}_i - \theta_* \rangle \geq \sqrt{2 \|a\|_{V_i^{-1}}^2 \log \left(\frac{1}{\beta} \right)} \right] \leq \beta, \quad (22)$$

where $\beta \in (0, 1)$ and $\|a\|_{V_i^{-1}}^2 = a^\top V_i^{-1} a$. Let $V_i(\pi_i) = \sum_{a \in \mathcal{C}_i} \pi_i(a) a a^\top$ and hence we have

$$V_i = \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top \geq q^i \sum_{a \in \mathcal{C}_i} \pi_i(a) a a^\top = q^i V_i(\pi_i). \quad (23)$$

Observe that for any symmetric random variable x if $\Pr[x \geq t] \leq \beta$, then $\Pr[|x| \geq t] = \Pr[x \geq t] + \Pr[-x \geq t] \leq 2\beta$. Thus, from lemma 2, we have $\|a\|_{V_i^{-1}}^2 = \frac{1}{q^i} a^\top V_i(\pi_i)^{-1} a \leq \frac{2d}{q^i}$ for each $a \in \mathcal{A}_i$. By setting $\beta = \frac{1}{4KT^2}$ and $\|a\|_{V_i^{-1}}^2 \leq \frac{2d}{q^i}$ for each $a \in \mathcal{A}_i$ in (22), we get that:

$$\Pr \left[\left| \langle a, \bar{\theta}_i - \theta_* \rangle \right| \geq \sqrt{\frac{4d}{q^i} \log(4KT^2)} \right] \leq \frac{1}{2KT^2}, \quad (24)$$

for each $a \in \mathcal{A}_i$. Now, we compute the effect of the privacy in estimating θ_* by bounding difference $\langle a, \bar{\theta}_i - \hat{\theta}_i \rangle$. Observe that $\hat{r}_{ia} = \bar{r}_{ia} + z_{ia}$, where $z_{ia} \sim \text{Lap}(\frac{1}{\epsilon})$, and hence, we can write $\hat{\theta}_i - \bar{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} z_{ia} a$. Thus, for any $\alpha \in \mathcal{A}_i$, we have that:

$$\langle \alpha, \hat{\theta}_i - \bar{\theta}_i \rangle = \sum_{a \in \mathcal{C}_i} \alpha^\top V_i^{-1} a z_{ia}, \quad (25)$$

where $\alpha^\top V_i^{-1} a \leq \max_{b \in \mathcal{A}_i} \|b\|_{V_i^{-1}}^2 \leq \frac{2d}{q^i}$ for each $a \in \mathcal{C}_i$ that holds from the fact that V_i is positive semi-definite. From Lemma 4 presented at the end of the section, by setting $b = \epsilon$, $n = Bd$, $c = \frac{2d}{q^i}$, and $t = 2\frac{Bd^2}{\epsilon q^i} + \frac{2d \log(4KT^2)}{\epsilon q^i}$, we get that:

$$\Pr \left[\left| \langle a, \bar{\theta}_i - \hat{\theta}_i \rangle \right| \geq 2\frac{Bd^2}{\epsilon q^i} + \frac{2d \log(4KT^2)}{\epsilon q^i} \right] \leq \frac{1}{2KT^2}, \quad (26)$$

Then, by the union bound and triangle inequality we have that

$$\Pr \left[\left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| > \gamma_i \quad \forall i \in [l \log T] \forall a \in \mathcal{A}_i \right] \leq \frac{1}{T}, \quad (27)$$

where $\gamma_i = \sqrt{\frac{4d}{q^i} \log(4KT^2)} + \frac{2Bd^2 + 2d \log(4KT^2)}{\epsilon q^i}$. This concludes the proof of Lemma 3. \blacksquare

Lemma 4. *Let $x_i = l_i z_i$ for $i \in [n]$, where $z_i \sim \text{Lap}(1/b)$ and l_i is constant such that $|l_i| \leq c$. Let $\bar{x} = \sum_{i=1}^n x_i$. We have that*

$$\Pr[\bar{x} \geq t] \leq \begin{cases} \exp\left(-\frac{t^2 b^2}{2nc^2}\right) & \text{if } t \leq \frac{nc}{b} \\ \exp\left(\frac{n}{2} - \frac{b}{c}t\right) & \text{if } t > \frac{nc}{b} \end{cases} \quad (28)$$

Proof. The proof follows from the concentration results of the Laplace distribution (e.g., see). We have that

$$\begin{aligned} \Pr[\bar{x} \geq t] &= \Pr[\exp(\lambda \bar{x}) \geq e^{\lambda t}] && \forall \lambda \geq 0 \\ &\stackrel{(a)}{\leq} \frac{\mathbb{E}[\exp(\lambda \bar{x})]}{e^{\lambda t}} \\ &\stackrel{(b)}{=} \frac{\prod_{i=1}^n \mathbb{E}[e^{\lambda x_i}]}{e^{\lambda t}} \\ &\stackrel{(c)}{\leq} \frac{\prod_{i=1}^n e^{\lambda^2 \frac{l_i^2}{2b^2}}}{e^{\lambda t}} && \forall 0 \leq \lambda \leq \frac{b}{c} \\ &\stackrel{(d)}{\leq} \frac{e^{\lambda^2 n \frac{c^2}{2b^2}}}{e^{\lambda t}} && \forall 0 \leq \lambda \leq \frac{b}{c} \end{aligned} \quad (29)$$

where step (a) follows from Markov's inequality and step (b) follows from the fact that z_1, \dots, z_n are independent Laplace random variables. Step (c) follows from the fact that z_i is sub-exponential random variable with proxy $\frac{l_i^2}{2b^2}$. Step (d) follows from the fact that $|l_i| \leq c$. By choosing $\lambda = \frac{tb^2}{nc^2}$ when $t < \frac{nc}{b}$ and $\lambda = \frac{b}{c}$ when $t > \frac{nc}{b}$, we get that

$$\Pr[\bar{x} \geq t] \leq \begin{cases} \exp\left(-\frac{t^2 b^2}{2nc^2}\right) & \text{if } t \leq \frac{nc}{b} \\ \exp\left(\frac{n}{2} - \frac{b}{c}t\right) & \text{if } t > \frac{nc}{b} \end{cases}, \quad (30)$$

This completes the proof of Lemma 4. \blacksquare

B Regret and Privacy Analysis of The local DP Model (Proof of Theorem 2)

B.1 Privacy Analysis

The privacy proof is straightforward. For any client, since the reward is bounded by $|r| \leq 1$, the output $\hat{r} = r + \text{Lap}(1/\epsilon_0)$ is ϵ_0 -LDP from [5, Theorem 3.6].

B.2 Regret Analysis

We next prove the regret bound of Algorithm 2 for stochastic linear bandits with LDP. Our proof is similar to the proofs of the central DP Algorithm presented in Section A.2.

Let \tilde{R}_T be the regret defined in (13). Let \mathcal{G} be the good event $\left\{ \left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| < \gamma_i \ \forall i \in [\log T] \forall a \in \mathcal{A}_i \right\}$. Lemma 5 shows that the event \mathcal{G} holds with probability at least $1 - \frac{1}{T}$. In the remaining part of the proof we condition on the event \mathcal{G} . When $q^i < \max\{Bd, 2\log(4KT^2)\}$, the regret can be bounded by $\max\{Bd, 2\log(4KT^2)\}$, and when $q^i \geq \max\{Bd, 2\log(4KT^2)\}$, we bound $n_i \leq 2q^i$, and hence,

$$\gamma_i \leq \sqrt{\frac{4d}{q^i} \log(4KT^2)} + \frac{2d}{\epsilon_0} \sqrt{\frac{\log(4KT^2)}{q^i}} \leq \left(1 + \frac{1}{\epsilon_0}\right) 2d \sqrt{\frac{\log(4KT^2)}{q^i}}.$$

By following similar steps as in the central DP, we can show that there is universal constants C', C such that the total regret in (13) can be bounded as

$$\begin{aligned} \tilde{R}_T &\leq (Bd + 2\log(4KT^2)) \log(T) + \sum_{i=1}^{\log T} 4n_i \gamma_{i-1} \\ &\leq (Bd + 2\log(4KT^2)) \log(T) + \left(1 + \frac{1}{\epsilon_0}\right) 2d \sum_{i=1}^{\log T} 8q^i \sqrt{\frac{1}{q^{i-1}} \log(4KT^2)} \\ &\leq C' \left(1 + \frac{1}{\epsilon_0}\right) \left(d\sqrt{d} \log^2(T) + d\sqrt{d \log T} \sum_{i=1}^{\log T} q^{(i-1)/2} \right) q \\ &\stackrel{(a)}{\leq} C' \left(1 + \frac{1}{\epsilon_0}\right) q \left(d\sqrt{d} \log^2(T) + d\sqrt{d \log T} q^{\log T/2} \right) \\ &\stackrel{(b)}{\leq} C' \left(1 + \frac{1}{\epsilon_0}\right) q \left(d\sqrt{d} \log^2(T) + d\sqrt{dT \log T} \right) \\ &\stackrel{(c)}{\leq} C \left(1 + \frac{1}{\epsilon_0}\right) \left(d\sqrt{dT \log T} \right), \end{aligned} \tag{31}$$

where step (a) follows from the sum of a geometric series and $q > 1$, step (b) uses $q = (2T)^{1/\log T}$, and step (c) follows from the facts $q \leq e^2$, $\log^2 T = O(\sqrt{T})$.

Hence, following similar steps as in the proof of the central DP algorithm, with probability at least $1 - \frac{1}{T}$ the regret is bounded as

$$R_T \leq \tilde{R}_T + 1 \leq C \left(1 + \frac{1}{\epsilon_0}\right) \left(d\sqrt{dT \log T} \right). \tag{32}$$

Lemma 5. *Let $\hat{\theta}_i$ be the least square estimate of θ_* at the end of the i th batch of Algorithm 2. Then, we have that*

$$\Pr \left[\left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| > \gamma_i \ \forall i \in [\log T] \forall a \in \mathcal{A}_i \right] \leq \frac{1}{T}, \tag{33}$$

where $\gamma_i = \sqrt{\frac{4d}{q^i} \log(4KT^2)} + \frac{2d}{q^i \epsilon_0} \sqrt{n_i \log(4KT^2)}$.

Proof. Let $\hat{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} \hat{r}_{ia} a$ be the private estimate of θ_* and $\bar{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} \bar{r}_{ia} a$ be the non-private estimate of θ_* as $\{\bar{r}_{ia}\}$ are the non-private rewards, where $V_i = \sum_{a \in \mathcal{C}_i} n_{ia} a a^\top$ and $\hat{r}_{ia} = \sum_{j=1}^{n_{ia}} \hat{r}_{ia}^{(j)}$. Similar to the central DP in Section 3, we have that

$$\Pr \left[\left| \langle a, \bar{\theta}_i - \theta_* \rangle \right| \geq \sqrt{\frac{4d}{q^i} \log(4KT^2)} \right] \leq \frac{1}{2KT^2}, \tag{34}$$

for each $a \in \mathcal{A}_i$. Now, we compute the effect of the LDP in estimating θ_* by bounding difference $\langle a, \bar{\theta}_i - \hat{\theta}_i \rangle$. Observe that $\hat{r}_{ia} = \sum_{j=1}^{n_{ia}} \hat{r}_{ia}^{(j)} = \bar{r}_{ia} + z_{ia}$, where $\bar{r}_{ia} = \sum_{j=1}^{n_{ia}} r_{ia}^{(j)}$ and $z_{ia} = \sum_{j=1}^{n_{ia}} z_{ia}^{(j)}$, where

$z_{ia}^{(j)} \sim \text{Lap}(\frac{1}{\epsilon_0})$. Hence, we can write $\hat{\theta}_i - \bar{\theta}_i = V_i^{-1} \sum_{a \in \mathcal{C}_i} z_{ia} a$. Thus, for any $\alpha \in \mathcal{A}_i$, we have that:

$$\langle \alpha, \hat{\theta}_i - \bar{\theta}_i \rangle = \sum_{a \in \mathcal{C}_i} \sum_{j=1}^{n_{ia}} \alpha^\top V_i^{-1} a z_{ia}^{(j)}, \quad (35)$$

where $\alpha^\top V_i^{-1} a \leq \max_{b \in \mathcal{A}_i} \|b\|_{V_i^{-1}}^2 \leq \frac{2d}{q^i}$ for each $a \in \mathcal{C}_i$ that holds from the fact that V_i is positive semi-definite. From Lemma 4 presented in Section 3, by setting $b = \epsilon_0$, $n = n_i$, $c = \frac{2d}{q^i}$, and $t = \frac{2d}{q^i \epsilon_0} \sqrt{n_i \log(4KT^2)}$, we get that:

$$\Pr \left[\left| \langle a, \bar{\theta}_i - \hat{\theta}_i \rangle \right| \geq \frac{2d}{q^i \epsilon_0} \sqrt{n_i \log(4KT^2)} \right] \leq \frac{1}{2KT^2}, \quad (36)$$

Then, by the union bound and triangle inequality we have that

$$\Pr \left[\left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| > \gamma_i \quad \forall i \in [\log T] \forall a \in \mathcal{A}_i \right] \leq \frac{1}{T}, \quad (37)$$

where $\gamma_i = \sqrt{\frac{4d}{q^i} \log(4KT^2)} + \frac{2d}{q^i \epsilon_0} \sqrt{n_i \log(4KT^2)}$. This concludes the proof of Lemma 5. \blacksquare

C Regret and Privacy Analysis of The Shuffled Model (Proof of Theorem 3)

C.1 Privacy Analysis

We note that the data of each user j can be represented as $\cup_{a \in \mathcal{C}_i} \{(a, r_a^{(j)})\}$. We observe that our scheme is equivalent to performing the following steps

- Each user $j \in [n_i]$ sends its data $\mathcal{D}_j = \cup_{a \in \mathcal{C}_i} \{(a, r_a^{(j)})\}$ to the shuffler.
- The shuffler randomly permutes the sets $\mathcal{D}_1, \dots, \mathcal{D}_{n_i}$ to get $\mathcal{D}_{\pi(1)}, \dots, \mathcal{D}_{\pi(n_i)}$.
- The shuffler reveals n_i action reward pairs $(a_1, \hat{r}_{ia_1}), \dots, (a_{n_i}, \hat{r}_{ia_{n_i}})$, where $(a_j, \hat{r}_{ia_j}) \in \mathcal{D}_{\pi(j)}$, and \hat{r}_{ia_j} is the LDP version of r_{ia_j} ($\hat{r}_{ia_j} = r_{ia_j} + \text{Lap}(\frac{1}{\epsilon_0^{(i)}})$).

Hence, we shuffle the data, then feed it to an LDP mechanism with LDP parameter $\epsilon_0^{(i)}$ (as proved in Theorem 2). As a result, it follows from [10] that the output of the shuffler is (ϵ_i, δ) -DP where

$$\epsilon_i = \log \left(1 + \frac{e^{\epsilon_0^{(i)}} - 1}{e^{\epsilon_0^{(i)}} + 1} \left(\frac{8\sqrt{e^{\epsilon_0^{(i)}} \log(4/\delta)}}{\sqrt{n_i}} + \frac{8e^{\epsilon_0^{(i)}}}{n_i} \right) \right). \quad (38)$$

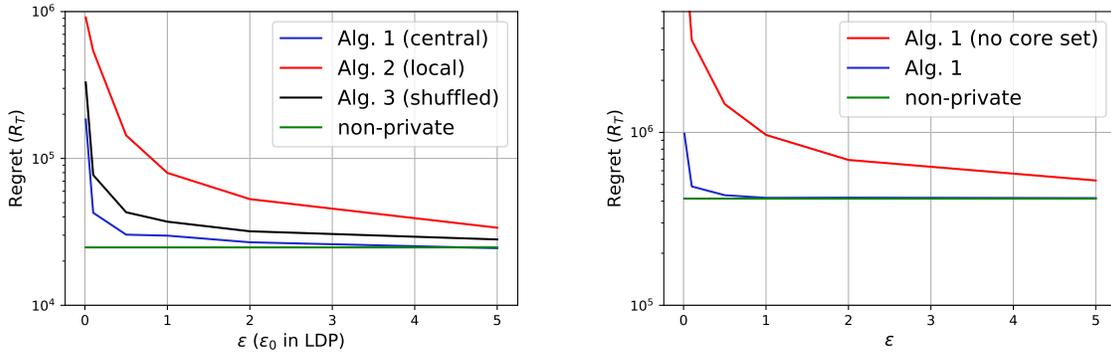
By the choice of $\epsilon_0^{(i)}$ as an inverse of the function $f_{n_i, \delta}$, we have that $\epsilon_i = \epsilon$ for all $i \in [\log T]$.

We observe that for any neighboring datasets D, D' , there is only one user data that is different between D, D' . That user appears in exactly one batch. It follows that Algorithm 3 is (ϵ, δ) -DP.

C.2 Regret Analysis

We next prove the regret bound of Algorithm 3 for stochastic linear bandits in the shuffled model. Our proof is similar to the proofs of the LDP Algorithm presented in Section B.2.

Let \tilde{R}_T be the regret defined in (13). Let \mathcal{G} be the good event $\left\{ \left| \langle a, \hat{\theta}_i - \theta_* \rangle \right| < \gamma_i \quad \forall i \in [\log T] \forall a \in \mathcal{A}_i \right\}$. Lemma 5 shows that the event \mathcal{G} holds with probability at least $1 - \frac{1}{T}$. In the remaining part of the proof



(a) Central, local and shuffled models, $K = 10, T = 10^6$. (b) Effect of core set size, $K = 1000, T = 10^7$.

Figure 3: Regret-privacy trade-offs for stochastic linear bandits algorithms.

we condition on the event \mathcal{G} . When $q^i < Bd$, the regret can be bounded by Bd . By following similar steps as in the central DP, we can show that there is universal constants C' such that the total regret in (13) can be bounded as

$$\begin{aligned}
\tilde{R}_T &\leq Bd \log(T) + \sum_{i=1}^{\log T} 4n_i \gamma_{i-1} \\
&\stackrel{(a)}{\leq} Bd \log(T) + \sum_{i=1}^{\log T} 8q^i \sqrt{\frac{4d}{q^{i-1}} \log(4KT^2)} + C' \frac{2d}{\epsilon} \sum_{i=1}^{\log T} 8q \sqrt{\log(4KT^2) \log(1/\delta)} \\
&\leq C \left(d\sqrt{T \log T} + \frac{(d \log T)^{3/2} \sqrt{\log(1/\delta)}}{\epsilon} \right), \tag{39}
\end{aligned}$$

where step (a) follows from the fact that from the privacy analysis, when $\epsilon_0^{(i)} \leq 1$, we get that $\epsilon = O(\epsilon_0^{(i)} \sqrt{\frac{\log(1/\delta)}{n_i}})$.

Hence, following similar steps as in the proof of the central DP algorithm, with probability at least $1 - \frac{1}{T}$ the regret is bounded as

$$R_T \leq \tilde{R}_T + 1 \leq C \left(d\sqrt{T \log T} + \frac{(d \log T)^{3/2} \sqrt{\log(1/\delta)}}{\epsilon} \right). \tag{40}$$

D Additional Numerical Results

Data Generation. We generate synthetic data generated as follows. The set of actions \mathcal{A} contains K actions, where each action $a \in \mathcal{A}$ is a $d = 2$ -dimensional vector. The actions $a \in \mathcal{A}$ and the optimal parameter θ_* are generated uniformly at random from the unit sphere $\mathcal{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$. Figure 4 plots the total regret R_T over an horizon $T = 10^6$ as a function of the privacy budget (ϵ or ϵ_0 in case of LDP mechanisms). Figure 3 plots the total regret R_T over an horizon T as a function of the privacy budget (ϵ or ϵ_0 in case of LDP mechanisms).

Usefulness of Core Set. In Figure 3b, we explore potential benefits on the performance of Algorithm 1 that use of the core set can offer. We consider $K = 1000$ and $T = 10^7$, and plot the regret of Algorithm 1 for two cases: (i) when we use a core set of size 2-3 actions, similar to the dimension of our space (labeled as Alg. 1), and (ii) when no core set is used, and instead the good set of actions of the batched algorithm is the whole action set (labeled as Alg. 1 no-core-set). We find that, as expected from our theoretical

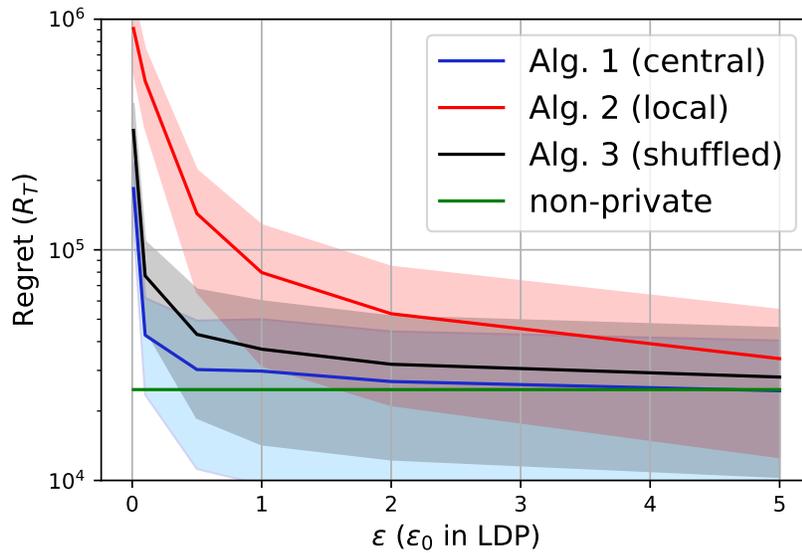


Figure 4: Regret-privacy trade-offs for stochastic linear bandits algorithms with $T = 10^6$.

analysis, using a core set enables to achieve performance very close to that of a non-private batched algorithm that adds no noise. In contrast, using (and adding noise to) the entire action space significantly degrades the performance.

We here present indicative results on the performance of our proposed Algorithms 1, 2 and 3.

Comparison of Algorithms 1, 2 and 3. In Figure 3a, we compare the regret of the proposed algorithms in the central, local and shuffled models using $K = 10, T = 10^6$. We observe that all algorithms converge to the regret of non-private stochastic linear bandit algorithms [24] as $\epsilon \rightarrow \infty$ ($\epsilon_0 \rightarrow \infty$), albeit at different rates. As predicted from the theoretical analysis, Algorithms 1 (central) and 3 (shuffled) offer privacy (almost) for free, closely following the non-private regret.

References

- [1] J. Mary, R. Gaudel, and P. Preux, “Bandits and recommender systems,” in *International Workshop on Machine Learning, Optimization and Big Data*. Springer, 2015, pp. 325–336.
- [2] D. Bouneffouf, I. Rish, and G. A. Cecchi, “Bandit models of human behavior: Reward processing in mental disorders,” in *International Conference on Artificial General Intelligence*. Springer, 2017, pp. 237–248.
- [3] A. N. Rafferty, H. Ying, and J. J. Williams, “Bandit assignment for educational experiments: Benefits to students versus statistical power,” in *International Conference on Artificial Intelligence in Education*. Springer, 2018, pp. 286–290.
- [4] D. Bouneffouf and I. Rish, “A survey on practical applications of multi-armed and contextual bandits,” *arXiv preprint arXiv:1904.10040*, 2019.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. D. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Theory of Cryptography Conference (TCC)*, 2006, pp. 265–284.
- [6] C. Dwork and A. Roth, “The algorithmic foundations of differential privacy,” *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [7] A. Cheu, A. D. Smith, J. Ullman, D. Zeber, and M. Zhilyaev, “Distributed differential privacy via shuffling,” in *Advances in Cryptology - EUROCRYPT 2019*, vol. 11476. Springer, 2019, pp. 375–403.
- [8] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, “Amplification by shuffling: From local to central differential privacy via anonymity,” in *SODA*. SIAM, 2019, pp. 2468–2479.
- [9] B. Balle, J. Bell, A. Gascón, and K. Nissim, “The privacy blanket of the shuffle model,” in *Annual International Cryptology Conference*. Springer, 2019, pp. 638–667.
- [10] V. Feldman, A. McMillan, and K. Talwar, “Hiding among the clones: A simple and nearly optimal analysis of privacy amplification by shuffling,” in *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE, 2022, pp. 954–964.
- [11] A. M. Girgis, D. Data, S. Diggavi, A. T. Suresh, and P. Kairouz, “On the renyi differential privacy of the shuffle model,” in *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021, pp. 2321–2341.
- [12] P. Rusmevichientong and J. N. Tsitsiklis, “Linearly parameterized bandits,” *Mathematics of Operations Research*, vol. 35, no. 2, pp. 395–411, 2010.
- [13] R. Shariff and O. Sheffet, “Differentially private contextual linear bandits,” vol. 31, 2018.
- [14] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári, “Improved algorithms for linear stochastic bandits,” *Advances in neural information processing systems*, vol. 24, 2011.
- [15] Y. Han, Z. Liang, Y. Wang, and J. Zhang, “Generalized linear bandits with local differential privacy,” vol. 34, 2021.
- [16] E. Garcelon, K. Chaudhuri, V. Perchet, and M. Pirotta, “Privacy amplification via shuffling for linear contextual bandits,” in *International Conference on Algorithmic Learning Theory*. PMLR, 2022, pp. 381–407.
- [17] S. R. Chowdhury and X. Zhou, “Shuffle private linear contextual bandits,” *arXiv preprint arXiv:2202.05567*, 2022.
- [18] T. Sajed and O. Sheffet, “An optimal private stochastic-mab algorithm based on optimal private stopping rule,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5579–5588.

- [19] W. Ren, X. Zhou, J. Liu, and N. B. Shroff, “Multi-armed bandits with local differential privacy,” *arXiv preprint arXiv:2007.03121*, 2020.
- [20] J. Tenenbaum, H. Kaplan, Y. Mansour, and U. Stemmer, “Differentially private multi-armed bandits in the shuffle model,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [21] K. Zheng, T. Cai, W. Huang, Z. Li, and L. Wang, “Locally differentially private (contextual) bandits learning,” vol. 33, 2020, pp. 12 300–12 310.
- [22] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, “What can we learn privately?” *SIAM Journal on Computing*, vol. 40, no. 3, pp. 793–826, 2011.
- [23] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*. Cambridge university press, 2018, vol. 47.
- [24] T. Lattimore and C. Szepesvári, *Bandit algorithms*. Cambridge University Press, 2020.
- [25] M. Frank and P. Wolfe, “An algorithm for quadratic programming,” *Naval research logistics quarterly*, vol. 3, no. 1-2, pp. 95–110, 1956.
- [26] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, “Gambling in a rigged casino: The adversarial multi-armed bandit problem,” in *Proceedings of IEEE 36th annual foundations of computer science*. IEEE, 1995, pp. 322–331.