# DIFFERENTIALLY PRIVATE DEEP MODEL-BASED REIN FORCEMENT LEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

Paper under double-blind review

#### ABSTRACT

We address private deep offline reinforcement learning (RL), where the goal is to train a policy on standard control tasks that is differentially private (DP) with respect to individual trajectories in the dataset. To achieve this, we introduce PRIMORL, a model-based RL algorithm with formal differential privacy guarantees. PRIMORL first learns an ensemble of trajectory-level DP models of the environment from offline data. It then optimizes a policy on the penalized private model, without any further interaction with the system or access to the dataset. In addition to offering strong theoretical foundations, we demonstrate empirically that PRIMORL enables the training of private RL agents on offline continuous control tasks with deep function approximations, whereas current methods are limited to simpler tabular and linear Markov Decision Processes (MDPs). We furthermore outline the tradeoffs involved in achieving privacy in this setting.

#### 1 INTRODUCTION

Despite Reinforcement Learning's (RL) notable advancements in various tasks, there have been 026 many obstacles to its adoption for the control of real systems in the industry. In particular, online 027 interaction with the system may be impractical or hazardous in real-world scenarios. Offline RL (Levine et al., 2020) refers to the set of methods enabling the training of control agents from static 029 datasets. While this paradigm shows promise for real-world applications, its deployment is not without concerns. Many studies have warned of the risk of privacy leakage when deploying machine 031 learning models, as these models can memorize part of the training data. For instance, Rigaki & Garcia (2020) review the proliferation of sophisticated privacy attacks. Of the various attack types, 033 membership inference attacks (MIAs) (Shokri et al., 2017) stand out as the most prevalent. In these 034 attacks, the adversary, with access to a black-box model trainer, attempts to predict whether a specific data point was part of the model's training data. Unfortunately, RL is no exception to these threats. In recent work, Gomrokchi et al. (2023) exploit the temporal correlation of RL samples to perform 036 powerful membership inference attacks using convolutional neural classifiers. More precisely, they 037 demonstrate that given access to the output policy, an adversary can learn to infer the presence of a specific trajectory — which is the result of a sequence of interactions between a user and the system — in the training dataset with great accuracy. The threat of powerful MIAs is particularly concerning 040 in reinforcement learning, where a trajectory can unveil sensitive user information. For instance, 041 when using RL to train autonomous vehicles (Kiran et al., 2022), we need to collect a large number of 042 trips that may disclose locations and driving habits. Similarly, a browsing journey collected to train a 043 personalized recommendation engine may contain sensitive information about the user's behavior 044 (Zheng et al., 2018). In healthcare, RL's potential for personalized treatment recommendation (Liu et al., 2022) underscores the need to safeguard patients' treatment and health history.

A large body of work has focused on protecting against privacy leakages in machine learning.
 Differential Privacy (DP), which allows learning models without exposing sensitive information about any particular user in the training dataset, has emerged as the gold standard. While successfully applied in various domains, such as neural network training (Abadi et al., 2016) and multi-armed bandits (Tossou & Dimitrakakis, 2016), extending differential privacy to reinforcement learning poses challenges. In particular, the many ways of collecting data and the correlated nature of training samples resulting from online interactions make it difficult to come up with a universal and meaningful DP definition in this setting, despite several attempts such as local and joint DP. In addition to its practical significance, the offline RL setting arguably offers a more natural framework for privacy

compared to the classic online setting. An offline RL method can indeed be seen as a black-box randomized algorithm h taking in as input a fixed dataset  $\mathcal{D}$ , partitioned in trajectories, and outputting a policy  $\hat{\pi}$ . An adversary with access to  $\hat{\pi}$  may successfully learn to infer the membership of a specific trajectory in  $\mathcal{D}$ , which can, as emphasized before, reveal sensitive user information. Hence, similarly to Qiao & Wang (2023a), we use the following informal DP definition for offline RL, which we refer to as *trajectory-level differential privacy* (TDP): adding or removing a single trajectory from the input dataset of an offline RL algorithm must not impact significantly the distribution of the output policy.

061 While reinforcement learning encounters the same privacy challenges as other areas of machine 062 learning, no existing work has proposed a private RL method that matches the versatility, scalability, 063 and empirical effectiveness of DP-SGD (Abadi et al., 2016) for supervised learning. Indeed, existing 064 research largely remains theoretical and demonstrates limited practical applicability. In the online setting, numerous private algorithms have been developed (e.g., Vietri et al. (2020), Garcelon et al. 065 (2021), Qiao & Wang (2023b)) but their scope remains restricted to tabular and linear Markov 066 Decision Processes (MDPs) with finite horizon. Qiao & Wang (2023a) have proposed the first private 067 algorithms for offline RL, building on value iteration methods, but they present similar limitations. 068 These approaches cannot intrinsically scale to the problems typically encountered in deep RL, leaving 069 a huge gap between the current private RL literature and real-world applications. This work addresses this gap by introducing the first deep RL method with provable privacy guarantees. In contrast to 071 previous work, our method is applicable to general MDPs with continuous state and action spaces 072 and deals with the classic  $\gamma$ -discounted setting, paving the way for enhanced applications of private 073 RL in complex, risk-sensitive scenarios.

074

075 **Contributions.** While the current differentially private RL literature is mainly theoretical and has 076 limited practical relevance, this work is the first attempt to tackle deep RL problems with differential 077 privacy guarantees. Specifically, we address the offline setting under the well-founded concept of trajectory-level privacy, and introduce a model-based approach named PRIMORL. Protecting entire trajectories, rather than individual examples, precludes the use of vanilla optimizers such 079 as DP-SGD. Additionally, the standard approach of using bootstrap ensembles to handle model 080 uncertainty presents an extra challenge in the private setting, as the ensemble size directly impacts 081 the privacy budget. A key contribution of this work is therefore the introduction of a training method for model ensembles that ensures differential privacy at the trajectory level and effectively controls 083 the privacy budget. We also provide a theoretical analysis of how private training influences model 084 reliability. We then perform policy optimization under the resulting pessimistic private model and 085 prove the formal privacy guarantees of the resulting policy. We show empirically that PRIMORL can train private policies with competitive privacy-performance trade-offs on standard continuous control 087 benchmarks, demonstrating the potential of our approach.

088 089

090

## 2 RELATED WORK

091 Offline RL (Levine et al., 2020; Prudencio et al., 2022) focuses on training agents without further 092 interactions with the system, making it essential in scenarios where data collection is impractical (Singh et al., 2022; Liu et al., 2020; Kiran et al., 2022). Model-based RL (Moerland et al., 2023) 094 can further reduce costs or safety risks by using a learned environment model to simulate beyond 095 the collected data and improve sample efficiency (Chua et al., 2018). Argenson & Dulac-Arnold 096 (2021) demonstrate that model-based offline planning, where the model is trained on a static dataset, 097 performs well in robotic tasks. However, offline RL faces challenges like distribution shift (Fujimoto 098 et al., 2019), where the limited coverage of the dataset can lead to inaccuracies in unexplored stateaction regions, affecting performance. Methods like MOPO (Yu et al., 2020), MOREL (Kidambi 099 et al., 2020), and COUNT-MORL (Kim & Oh, 2023) address this by penalizing rewards based on 100 model uncertainty, achieving strong results on offline benchmarks. Still, key design choices in offline 101 MBRL require further exploration, as highlighted by Lu et al. (2022). 102

On the other hand, Differential Privacy (DP), established by Dwork (2006), has become the standard
for privacy protection. Recent research has focused on improving the privacy-utility trade-off, with
relaxations of DP and advanced composition tools enabling tighter privacy analyses (Dwork et al.,
2010; Dwork & Rothblum, 2016; Bun & Steinke, 2016; Mironov, 2017a). Notably, DP-SGD (Abadi
et al., 2016) has facilitated the development of private deep learning algorithms, despite ongoing
practical challenges (Ponomareva et al., 2023). Concurrently, sophisticated attack strategies have

108 underscored the necessity for robust DP algorithms (Rigaki & Garcia, 2020). Recent studies have shown that reinforcement learning (RL) is also vulnerable to privacy threats (Pan et al., 2019; Prakash 110 et al., 2022; Gomrokchi et al., 2023). As RL is increasingly applied in personalized services (den 111 Hengst et al., 2020), the need for privacy-preserving training techniques is critical. Although DP 112 has been successfully extended to multi-armed bandits (Tossou & Dimitrakakis, 2016; Basu et al., 2019), existing RL algorithms (e.g., Vietri et al. (2020), Zhou (2022), Qiao & Wang (2023b)) with 113 formal DP guarantees mainly apply to episodic tabular or linear MDPs and lack empirical validation 114 beyond basic simulations. Moreover, private offline RL remains underexplored. Only Qiao & Wang 115 (2023a) have proposed DP offline algorithms, which, while theoretically strong, are also restricted to 116 finite-horizon tabular and linear MDPs. Consequently, no existing work has introduced DP methods 117 that can handle deep RL environments in the infinite-horizon discounted setting, a critical step toward 118 deploying private RL algorithms in real-world applications. With this work, we aim to fill this gap by 119 proposing a differentially private, deep model-based RL method for the offline setting. 120

121

123

124

145

147

154

155

#### **3** PRELIMINARIES

#### 3.1 OFFLINE MODEL-BASED REINFORCEMENT LEARNING

125 We consider an infinite-horizon discounted MDP, that is a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma, \rho_0)$  where  $\mathcal{S}$  and 126  $\mathcal{A}$  are respectively the state and action spaces,  $P: \mathcal{S} \times \mathcal{A} \longrightarrow \Delta(\mathcal{S})$  the transition dynamics (where 127  $\Delta(\mathcal{X})$  denotes the space of probability distributions over  $\mathcal{X}$ ),  $r: \mathcal{S} \times \mathcal{A} \longrightarrow [0, 1]$  the reward function, 128  $\gamma \in [0,1)$  a discount factor and  $\rho_0 \in \Delta(S)$  the initial state distribution. The dynamics satisfy the 129 Markov property, *i.e.*, the next state s' only depends on current state and action. The goal is to learn 130 a policy  $\pi: \mathcal{S} \longrightarrow \Delta(\mathcal{A})$  maximizing the expected discounted return  $\eta_{\mathcal{M}}(\pi) := \mathbb{E}_{\tau \sim \pi, \mathcal{M}}[R(\tau)]$ , where  $R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$ . The expectation is taken w.r.t. the trajectories  $\tau = ((s_t, a_t, r_t))_{t \geq 0}$ 131 132 generated by  $\pi$  in the MDP  $\mathcal{M}$ , *i.e.*,  $s_0 \sim \rho_0$ ,  $s_{t+1} \sim P(\cdot|s_t, a_t)$  and  $a_t \sim \pi(\cdot|s_t)$ . 133

In offline RL, we assume access to a dataset of K trajectories  $\mathcal{D}_K = (\tau_k)_{k=1}^K$ , where each  $\tau_k =$ 134  $(s_t^{(k)}, a_t^{(k)}, r_t^{(k)})_{t>0}$  has been collected with an unknown behavioral policy  $\pi^B$ .  $\tau_k$  can be seen as the 135 result of the interaction of a user  $u_k$  with the environment. The objective is then to learn a policy  $\hat{\pi}$ 136 from  $\mathcal{D}_K$  (without any further interaction with the environment) which performs as best as possible 137 in  $\mathcal{M}$ . To achieve this goal, we consider a model-based approach. In this context, we learn estimates 138 of both the transition dynamics and the reward function, denoted  $\hat{P}$  and  $\hat{r}$  respectively, from the 139 offline dataset  $\mathcal{D}_K$ . This results in an estimate of the MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \dot{P}, \hat{r}, \gamma, \rho_0)$ . We can then use 140 the model  $\hat{\mathcal{M}}$  as a simulator of the environment to learn a policy  $\hat{\pi}_{\hat{\mathcal{M}}}$ , without further access to the 141 dataset or interactions with the real environment modeled by  $\mathcal{M}$ . Note that if the policy  $\hat{\pi}_{\hat{\mathcal{M}}}$  is trained 142 143 to maximize the expected discounted return in the MDP model  $\mathcal{M}$ , *i.e.*,  $\hat{\pi}_{\mathcal{M}} \in \operatorname{argmax} \eta_{\mathcal{M}}(\pi)$ , we eventually want to evaluate the policy in the true environment  $\mathcal{M}$ , that is using  $\eta_{\mathcal{M}}$ . 144

#### 146 3.2 DIFFERENTIAL PRIVACY

When learning patterns from a dataset, differential privacy (Dwork, 2006) protects against the leakage of sensitive information in the data by ensuring that the output of the algorithm does not change significantly when adding or removing a data point, as formally stated in Definition 3.1.

**Definition 3.1.**  $(\epsilon, \delta)$ -differential privacy. Given  $\epsilon > 0, \delta \in [0, 1)$ , a mechanism h (i.e., a randomized function of the data) is  $(\epsilon, \delta)$ -DP if for any pair of datasets D, D' that differ in at most one element (referred to as neighboring datasets, and denoted d(D, D') = 1), and any subset  $\mathcal{E}$  in h's range:

$$\mathbb{P}(h(D) \in \mathcal{E}) \le e^{\epsilon} \cdot \mathbb{P}(h(D') \in \mathcal{E}) + \delta .$$

In particular,  $\epsilon$  controls the strength of the privacy guarantees, decreasing as  $\epsilon$  grows. To achieve ( $\epsilon, \delta$ )-DP, the standard approach is to add a zero-mean random noise to the output of the (non-private) function f, whose magnitude  $\sigma$  scales with  $\Delta_{\ell}(f)/\epsilon$ , where  $\Delta_{\ell}(f) := \max_{d(D,D')=1} ||f(D) - f(D')||_{\ell}$ is the sensitivity of f. One of the most used DP mechanisms is the *Gaussian mechanism*, which provably guarantees ( $\epsilon, \delta$ )-DP for  $\epsilon, \delta \in (0, 1)$  by adding random noise from a Gaussian distribution with magnitude  $\sigma = \epsilon^{-1} \sqrt{2 \log (1.25/\delta)} \cdot \Delta_2(f)$ . From such simple mechanisms, we can derive



4 DIFFERENTIALLY PRIVATE MODEL-BASED OFFLINE REINFORCEMENT LEARNING

189 190 191

192

193

194

195

187

188

We now describe our model-based approach for learning differentially private RL agents from offline data, which we call PRIMORL (for <u>Private Model-Based Offline RL</u>). After defining trajectory-level differential privacy (TDP) in offline RL (Section 4.1), we address the learning of a private model from offline data (Section 4.2). Finally, we demonstrate how we optimize a policy under the private model (Section 4.3). Exploiting the *post-processing* property of DP, we show that ensuring model privacy alone is enough to achieve a private policy. Figure 1 provides a high-level description of PRIMORL.

196 197

199

4.1 TRAJECTORY-LEVEL PRIVACY IN OFFLINE REINFORCEMENT LEARNING

In supervised learning, differential privacy is typically applied at the *example-level*, under the assumption that the examples in the dataset are independent. If this assumption is already questionable in the supervised setting, it certainly does not hold in RL where the transitions  $(s_t, a_t, r_t)$  are obviously correlated. Several works, for instance Liu et al. (2016), have demonstrated that data correlation degrades privacy guarantees in the traditional *per-example* setting. It thus appears that protecting individual transitions is insufficient in RL, calling instead for data protection at the *trajectory level*.

We introduce the following formal definition for trajectory-level differential privacy (TDP) in offline RL which protects whole trajectories. It states that the learned policy is roughly the same for two offline datasets D, D' where D' is obtained by adding or removing one full trajectory from D. It can be seen as a reformulation of the definition used in Qiao & Wang (2023a), which is the first work to tackle differential privacy in this setting.

211 212 213 214 Definition 4.1.  $(\epsilon, \delta)$ -*TDP*. Let h be an offline RL algorithm, that takes as input an offline dataset and outputs a policy. Given  $\epsilon > 0$  and  $\delta \in (0, 1)$ , h is  $(\epsilon, \delta)$ -TDP if for any trajectory-neighboring datasets  $\mathcal{D}_K, \mathcal{D}_{K \setminus \{k\}}$ , and any subset of policies  $\Pi$ :

$$\mathbb{P}\left(h(\mathcal{D}_K) \in \Pi\right) \le e^{\epsilon} \cdot \mathbb{P}\left(h(\mathcal{D}_{K \setminus \{k\}}) \in \Pi\right) + \delta$$

# 4.2 MODEL LEARNING WITH DIFFERENTIAL PRIVACY

Following previous work (Yu et al., 2020; Kidambi et al., 2020), we jointly model the transition dynamics  $\hat{P}$  and reward  $\hat{r}$  with a Gaussian distribution  $\hat{M}$  conditioned on the current state and action. Its mean and covariance are parameterized with neural networks  $\theta = (\phi, \psi)$ :

221 222

$$\hat{M}_{\theta}\left(\Delta_t^{t+1}(s), r_t | s_t, a_t\right) = \mathcal{N}\left(\mu_{\phi}(s_t, a_t), \Sigma_{\psi}(s_t, a_t)\right) \ .$$

To carry out uncertainty estimation (see Section 4.3), we train an ensemble of N models  $\hat{M}_{\theta_i}$ ,  $i \in [\![1, N]\!]$ , all sharing the same architecture. The core aspect of PRIMORL, as illustrated in Figure 1, is therefore to learn a trajectory-level DP dynamics model ensemble.

This poses two major challenges. First, the traditional approach to privatize neural networks, DP-SGD, is designed for example-level privacy and is unsuitable for guaranteeing TDP. Moreover, since the training of all the models in the ensemble consumes the same dataset  $\mathcal{D}_K$ , we must deal with the dependence of the privacy budget on the ensemble size N. A key contribution of our work, developed in Section 4.2.1, is thus to introduce a training method that 1) guarantees privacy at the trajectory level and 2) efficiently manages the privacy budget across an ensemble of models.

#### 233

235

### 4.2.1 TRAJECTORY-LEVEL DP TRAINING FOR MODEL ENSEMBLES

As DP-SGD ensures per-example privacy by clipping each per-example gradient, limiting the 236 contribution of each data point to the final model, the key to achieving trajectory-level privacy is to 237 compute and clip per-trajectory updates. Therefore, our training method partitions the dataset by 238 trajectories, *i.e.*,  $\mathcal{D}_K = \bigcup_{k=1}^K \{\tau_k\}$ , computes independent updates from each trajectory's data, and 239 bounds the  $L_2$ -norm of each update before aggregation. This idea has been developed in McMahan 240 et al. (2017) to achieve user-level privacy when training recurrent language models. Building on 241 prior training algorithms from federated learning, they introduce DP-FEDAVG, which leverages 242 privacy amplification by sub-sampling to achieve competitive privacy-utility trade-offs in language 243 modeling. To address the unique privacy challenges of our task, we build on this approach and adapt 244 it to ensembles of dynamics models. We present the resulting training procedure in Algorithm 1. 245

The core idea behind TDP MODEL ENSEMBLE TRAINING is to draw, at each iteration t, a random 246 subset  $\mathcal{U}_t$  of the K trajectories (line 2) using Poisson sampling. Each trajectory is drawn with 247 probability q, resulting in an expected qK trajectories being selected per step. The sampling ratio q 248 plays a critical role in determining the strength of privacy guarantees. Specifically, a smaller q reduces 249 the likelihood of any given trajectory being included in an update, thereby limiting its influence on the 250 final model — this forms the basis of privacy amplification by sub-sampling. However, in the offline 251 RL setting, where trajectory data is highly correlated, q must remain large enough to ensure that the 252 model update incorporates a sufficiently diverse set of trajectories. Interestingly, while the theoretical 253 analysis of common private deep learning methods like DP-SGD relies on Poisson sampling, most implementations actually use fixed-size batches with shuffling in practice, in order to overcome the 254 computational challenges due to batches of varying size. This can lead to significant underestimation 255 of the actual privacy leakage, as pointed out in Chua et al. (2018). Our implementation, however, 256 does indeed use Poisson sampling, allowing us to compute correct theoretical privacy guarantees. 257

For each trajectory  $\tau_k \in \mathcal{U}_t$ , the clipped gradients  $\{\Delta_{i,k}^{\text{clipped}}(t)\}$  are then computed from  $\tau_k$ 's data 258 259 only (line 3 to 7). During this step, we perform multiple local updates on the same trajectory's data, 260 leveraging larger global updates without incurring more privacy leakage. This is made possible 261 because the global model is updated with clipped gradients only. We later introduce ensemble-adapted 262 clipping strategies to control the privacy budget over model ensembles, ensuring that the sensitivity of the ensemble gradient  $\Delta_k^{\text{clipped}}(t) = \left(\Delta_{i,k}^{\text{clipped}}(t)\right)_{i=1}^N$  is bounded by C. We then compute an unbiased 263 264 estimator of the subset gradient average whose sensitivity is bounded by C/qK (line 8). We can then 265 apply the Gaussian mechanism with magnitude  $\sigma = zC/qK$ , where z controls the strength of the 266 privacy guarantee  $\epsilon$ , and update the ensemble model  $\theta(t) = (\theta_i(t))_{i=1}^N$  with noisy gradient (line 9): 267 avg ( .) . . . . . ( ( ... 268

$$\theta(t+1) \longleftarrow \theta(t) + \Delta^{\operatorname{avg}}(t) + \mathcal{N}\left(0_{Nd}, \sigma^2 I_{Nd}\right) \quad .$$

270

Algorithm 1 TDP MODEL ENSEMBLE TRAINING

1: for each iteration  $t \in [0, T-1]$  do 2:  $\mathcal{U}_t \leftarrow (\text{sample with replacement trajectories from } \mathcal{D}_K \text{ with prob. } q)$ 

3: for each trajectory  $\tau_k \in \mathcal{U}_t$  do

4:

Clone current models  $\{\theta_i^{\text{start}}\}_{i=1}^N \leftarrow \{\theta_i(t)\}_{i=1}^N$  $\{\theta_{i,k}\}_{i=1}^N \leftarrow \text{ENSCLIPGD}\left(\tau_k, \{\theta_i^{\text{start}}\}_{i=1}^N; C, \text{local epochs } E, \text{batch size } B\right)$  $\Delta_{i,k}^{\text{clipped}}(t) \leftarrow \theta_{i,k} - \theta_i^{\text{start}}, i = 1, ..., N$ 5:

277 278 279 6:

7:

8:

9:

10: end for

275

276

281

283 284 285

#### 4.2.2 PRIVACY GUARANTEES FOR THE MODEL

 $\Delta_i^{\text{avg}}(t) = \frac{\sum_{k \in \mathcal{U}_t} \Delta_{i,k}^{\text{clipped}}(t)}{qK}, \ i = 1, ..., N$ 

 $\theta(t+1) \leftarrow \theta(t) + \Delta^{\operatorname{avg}}(t) + \mathcal{N}\left(0_{Nd}, \left(\frac{zC}{qK}\right)^2 I_{Nd}\right)$ 

287 We can now derive formal privacy guarantees for a model trained using Algorithm 1. A key challenge 288 in our setting arises from training an ensemble of N models for uncertainty estimation, all using the 289 same dataset  $\mathcal{D}_K$ . Treating each model independently, with separate clipping and noise addition, would be inefficient and significantly increase the privacy budget by composition. This could be 290 mitigated by limiting the ensemble size, but at the cost of performance, as shown in Lu et al. (2022). 291

292 To address this challenge, we process all the gradients of the model ensemble simultaneously and 293 distribute the global clipping norm C across all models, on the same principle as the per-layer clipping 294 used in McMahan et al. (2017). Denoting  $\Delta_{i,\ell}$  the gradient of layer  $\ell$  for model i, we propose and experiment with two ensemble clipping strategies: Flat Ensemble Clipping, which clips the whole 295 model gradient  $\Delta_i = (\Delta_{i,\ell})_{\ell=1}^L$  with  $C_i = C/\sqrt{N}$ ; and **Per Layer Ensemble Clipping**, which clips 296

297 298

per-layer gradients  $\Delta_{i,\ell}$  with  $C_{i,\ell} = C/\sqrt{N \times L}$ , so that  $C = \sqrt{\sum_{i=1}^{N} C_i^2} = \sqrt{\sum_{i=1}^{N} \sum_{\ell=1}^{L} C_{i,\ell}^2}$ . For both strategies, we verify that that  $\Delta_k^{\text{clipped}} = \left(\Delta_{i,k}^{\text{clipped}}\right)_{i=1}^{K}$  has sensitivity bounded by C (see Theorem 4.2's proof in appendix), and that the contribution of a given trajectory to the *model ensemble* is constrained by C (see 299 300 301 is appropriately limited. Ensemble clipping eliminates the linear dependence of the privacy budget 302 on the number of models. However, it does not entirely remove the negative impact of increasing N. 303 Indeed, for a given noise level, a larger N requires a smaller clipping threshold  $C_i$  or  $C_{i,\ell}$ , which can 304 degrade model convergence by losing too much information from the original gradient. Nevertheless, the clipping threshold scales with the square root of N, mitigating the impact to some extent. 305

306 We now formally derive the privacy guarantees for an ensemble of models trained with Algorithm 1. 307 Mapping users in federating learning to trajectories in offline RL, we can directly adapt Theorem 1 308 from McMahan et al. (2018) to state that, with the sensitivity of clipped gradients  $\Delta_{i,k}^{\text{clipped}}$  effectively 309 bounded by C, the moments accounting method from Abadi et al. (2016) computes correctly the 310 privacy loss of Algorithm 1 at trajectory-level for the noise multiplier  $z = \sigma/\mathbb{C}$  with  $\mathbb{C} = C/qK$ . 311 We can therefore use the moments accountant to compute, given  $\delta \in (0, 1)$ ,  $z > 0, q \in (0, 1)$  and 312  $T \in \mathbb{N}$ , the total privacy budget  $\epsilon$  spent by Algorithm 1, and obtain  $(\epsilon, \delta)$ -TDP guarantees for our 313 dynamics model, as stated in Theorem 4.2 (full proof in appendix).

314 **Theorem 4.2.**  $(\epsilon, \delta)$ -TDP guarantees for dynamics model. Given  $\delta \in (0, 1)$ , noise multiplier z, 315 sampling ratio q and number of training iterations T, let  $\epsilon := \epsilon^{MA}(z, q, T, \delta)$  be the privacy budget 316 computed by the moments accounting method from (Abadi et al. (2016), more details in Section H.6). 317 The dynamics model output by Algorithm 1 is  $(\epsilon, \delta)$ -TDP.

318 319

#### 4.3 POLICY OPTIMIZATION UNDER A PRIVATE MODEL 320

Now that we learned a private model  $\hat{M}$  from offline data, we use it as a simulator of the environment 321 to learn a private policy  $\hat{\pi}$  with a model-based policy optimization approach. The use of a private 322 model and the privacy constraints on the end policy introduce additional challenges compared to the 323 non-private case, as demonstrated in Section 4.3.1. We study solutions to mitigate the detrimental

324 effects of private training on policy performance in Section 4.3.2, before deriving formal privacy 325 guarantees for a policy learned under a private model in Section 4.3.3. 326

#### 4.3.1 IMPACT OF PRIVACY ON POLICY OPTIMIZATION

It is first essential to examine the complexities of policy optimization in model-based offline RL and assess whether they are amplified in the private setting. A major challenge in model-based offline RL 330 is to handle the discrepancy between the true and the learned dynamics when optimizing the policy. 331 Indeed, model inaccuracies cause errors in policy evaluation that may be exploited, resulting in poor 332 performance in the real environment. According to the Simulation Lemma (Kearns & Singh, 2002; 333 Xu et al., 2020), the value evaluation error of a policy  $\pi$  in model-based RL can be decomposed 334 into a model error term and a policy distribution shift term. Formally, denoting  $\rho_P^{\pi^B}$  the state-action 335 discounted occupancy measure of the data-collection policy  $\pi^B$  under the true MDP, if the model error 336 is bounded as  $\mathbb{E}_{(s,a)\sim\rho_P^{\pi B}}\left[D_{KL}\left(P(\cdot|s,a)\|\hat{P}(\cdot|s,a)\right)\right] \leq \varepsilon_m$  and the distribution shift is bounded as  $\max_s D_{KL}\left(\pi(\cdot|s)\|\pi^B(\cdot|s)\right) \leq \varepsilon_{\pi}$ , then the value evaluation error of  $\pi$  is bounded as: 337 338 339

$$|\hat{V}^{\pi} - V^{\pi}| \le \frac{\sqrt{2\gamma}}{(1-\gamma)^2} \sqrt{\epsilon_m} + \frac{2\sqrt{2}}{(1-\gamma)^2} \sqrt{\varepsilon_\pi} \quad , \tag{1}$$

342 where  $\hat{V}^{\pi}$  and  $V^{\pi}$  denote the value of  $\pi$  under the learned and the true dynamics, respectively. 343 Controlling this quantity for an arbitrary  $\pi$  is crucial in our setting, as it ensures that the learned 344 MDP is a reasonable simulator of the true environment. Moreover, (1) directly implies a bound on 345 the sub-optimality gap, since  $|V^{\star} - V^{\hat{\pi}}| \leq 2 \sup_{\pi} |\hat{V}^{\pi} - V^{\pi}|$ . Under some assumptions regarding 346 the model loss function, Proposition 4.3 states the model error term in terms of the size  $N_{\mathcal{D}}$  of the 347 dataset.

348 Proposition 4.3. Value evaluation error in non-private offline MBRL. Let the model loss function 349 be L-Lipschitz and  $\Delta$ -strongly convex, and assumptions from the simulation lemma hold. There is a 350 stochastic convex optimization algorithm for learning the model and a constant M such that, with 351 probability at least  $1 - \alpha$ , and for sufficiently large  $N_{\mathcal{D}}$ , the value evaluation error of  $\pi$  is bounded 352 as:

$$|\hat{V}^{\pi} - V^{\pi}| \le \frac{\sqrt{2\gamma}}{(1-\gamma)^2} \cdot M \cdot \frac{L \log^{1/2}(N_{\mathcal{D}}/\alpha)}{\sqrt{\Delta N_{\mathcal{D}}}} + \frac{2\sqrt{2\gamma}}{(1-\gamma)^2} \sqrt{\varepsilon_{\pi}}$$

355 When we learn the model with differential privacy, we disrupt model convergence because of gradient 356 clipping and noise. This likely results in a less accurate dynamics model (although it may help prevent 357 overfitting in some cases) and increased value evaluation error. Intuitively, DP training impacts model 358 error in (1) as a direct result of gradient perturbations: Bassily et al. (2014), in particular, shows that 359 noisy gradient descent (GD) has increased excess risk compared to non-private GD. In the simpler 360 case where the model is trained with a vanilla DP noisy GD algorithm, Proposition 4.4 states the 361 value evaluation error under the private model.

Proposition 4.4. Value evaluation error in private offline MBRL. Let assumptions from Proposition 4.3 hold. If the model is learned with  $(\epsilon, \delta)$ -DP gradient descent, then, with probability at least  $1 - \alpha$ , there is a constant M' such that for large enough  $N_{\mathcal{D}}$ , the value evaluation error of  $\pi$ :

$$|\hat{V}_{DP}^{\pi} - V^{\pi}| \leq \frac{\sqrt{2}\gamma}{(1-\gamma)^2} \cdot M' \cdot \frac{Ld^{1/4}\log(N_{\mathcal{D}}/\delta) \cdot \operatorname{poly}\log(1/\alpha)}{\sqrt{\Delta N_{\mathcal{D}}\epsilon\alpha}} + \frac{2\sqrt{2}\gamma}{(1-\gamma)^2}\sqrt{\varepsilon_{\pi}}$$

where  $V_{\rm DP}^{\pi}$  is the value of  $\pi$  under the privately learned dynamics. Comparing the value evaluation 369 errors in Propositions 4.3 and 4.4 (both proven in appendix), we observe how DP training may degrade 370 performance in MBRL. The private bound has an explicit dependence on the problem dimension dwhich is not present in the non-private bound, and the  $\sqrt{\epsilon}$  factor in the denominator shows that the 372 error will degrade with strong privacy guarantees. On the other hand, the distribution shift term does 373 not depend on the learned dynamics and is therefore not affected by private training.

374 375

376

371

327

328

340 341

353 354

362

364

366 367 368

#### 4.3.2 MITIGATING PRIVATE MODEL UNCERTAINTY

In Section 4.3.1, we showed that private model training impacts the reliability of our model for 377 evaluating policies due to an increased dynamics error, which can lead to misjudging the quality of a 378 policy in the true environment. In the non-private case, this is typically handled by penalizing the 379 reward with a measure of the uncertainty of the model, denoted  $u: S \times A \to \mathbb{R}_+$ . Therefore, if the 380 model is believed to be unreliable at a given state-action pair (s, a) (*i.e.*, large u(s, a)), the possibly 381 over-estimated reward will be corrected as:

$$\tilde{r}(s,a) = \hat{r}(s,a) - \lambda \cdot u(s,a) \quad , \tag{2}$$

where  $\lambda$  is an hyperparameter. The policy is then optimized under the resulting pessimistic MDP 384  $\tilde{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \hat{P}, \tilde{r}, \gamma, \rho_0)$ . MOPO (Yu et al., 2020), MOREL (Kidambi et al., 2020) and more 385 recently COUNT-MORL (Kim & Oh, 2023) achieve impressive results on traditional offline RL 386 benchmarks with this approach, using different heuristics to estimate model uncertainty. 387

388 As suggested by the simulation lemma, the valuation error can depend on both model error and 389 distribution shift. However, we demonstrated that private training affects only model error. Interest-390 ingly, Lu et al. (2022), which studies design choices in offline model-based RL and the properties of various uncertainty estimators, finds that the uncertainty measures proposed in the literature are 391 more strongly correlated with model error than with distribution shift. Based on this, we believe 392 that existing uncertainty measures are well-suited to mitigate the diminished reliability of the model 393 under private training, as they will effectively capture the increased error. In particular, we consider 394 the maximum aleatoric uncertainty  $u_{MA}(s,a) = \max_{i \in [1,N]} \|\Sigma_{\psi_i}(s,a)\|_F$  (Yu et al., 2020) and the 395 maximum pairwise difference  $u_{\text{MPD}}(s, a) = \max_{i, j \in [1, N]} \|\mu_{\phi_i}(s, a) - \mu_{\phi_j}(s, a)\|_2$  (Kidambi et al., 396 2020). We compare both estimators (see Table 10 in the appendix) and find that neither is consistently 397 superior. However, we observe that the choice of estimator can affect performance on a specific 398 task. In addition, it seems reasonable to moderately increase the reward penalty  $\lambda$  compared to the 399 non-private case to take into account the greater uncertainty. 400

401 4.3.3 PRIVATE POLICY OPTIMIZATION

Given a choice of uncertainty estimator  $u \in \{u_{MA}, u_{MPD}\}$ , we now consider optimizing the policy 403 within the pessimistic private MDP  $\tilde{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \hat{P}, \tilde{r}_u, \gamma, \rho_0)$ , with  $\tilde{r}_u = \hat{r}(s, a) - \lambda \cdot u(s, a)$ . We 404 use Soft Actor-Critic (SAC, Haarnoja et al. (2018))<sup>1</sup>, a classic off-policy algorithm with entropy 405 regularization, to learn the policy from  $\mathcal{M}$ , in line with existing approaches in the offline MBRL 406 literature. Offline model-based methods typically mix real offline data from  $\mathcal{D}_K$  with model data 407 during policy learning (in MOPO, for instance, each batch contains 5% of real data). Here, however, 408 we learn the policy exclusively from model data to avoid incurring privacy loss beyond what is needed 409 to train the model, and thus control the privacy guarantees. Algorithm 4 in appendix provides a 410 pseudo-code for SAC policy optimization in the pessimistic private MDP. Using the post-processing 411 property of DP, we can now state in Theorem 4.5 that, given the  $(\epsilon, \delta)$ -TDP model  $\hat{M} = (\hat{P}, \hat{r})$ 412 learned as described in Section 4.2, the policy learned with Algorithm 4 under  $\tilde{M}$  is also  $(\epsilon, \delta)$ -TDP. 413 The full proof of this theorem is provided in appendix. 414

**Theorem 4.5.**  $(\epsilon, \delta)$ -TDP guarantees for PRIMORL. Given an  $(\epsilon, \delta)$ -TDP model  $(\hat{P}, \hat{r})$  learned 415 with Algorithm 1, the policy obtained with private policy optimization (Algorithm 4) within the 416 417 pessimistic model  $(\hat{P}, \hat{r} - \lambda u)$  is  $(\epsilon, \delta)$ -TDP. 418

419 420

421

423

427

428

402

382

5 **EXPERIMENTS** 

We empirically assess PRIMORL in three continuous control tasks: CARTPOLE-BALANCE and 422 CARTPOLE-SWINGUP from the DeepMind Control Suite (Tassa et al., 2018) as well as PENDULUM from OpenAI's Gym (Brockman et al., 2016). We also conduct experiments on HALFCHEETAH 424 (Wawrzynski, 2009), which we present in appendix (Section J). For simplicity, we refer to CARTPOLE-425 BALANCE and CARTPOLE-SWINGUP as BALANCE and SWINGUP. 426

5.1 EXPERIMENTAL SETTING

429 Following common practice, we evaluate the offline policies by running them in the real environment 430 We aim to assess the policy's performance degradation when varying the privacy level, as DP training 431

<sup>&</sup>lt;sup>1</sup>This could be any model-based policy optimization or planning algorithm that does not use offline data.



Figure 2: Learning curves on PENDULUM (left), BALANCE (middle) and SWINGUP (right).

may negatively affect it. We consider MOPO as our non-private baseline. For PRIMORL, we consider different configurations outlined in Table 3. The NO PRIVACY variant, without noise (z = 0), isolates the impact of trajectory-level model training on performance. The two private variants ( $\epsilon < \infty$ ) PRIMORL LOW and PRIMORL HIGH correspond to different noise multipliers. We discuss the choice of privacy parameters, as well as training hyperparameters and implementation details in appendix (Section H).

As the existing SWINGUP offline benchmark from Gülçehre et al. (2020) is very small (K = 40), and DP training of ML models typically requires significantly more data compared to non-private training (see, for instance, Ponomareva et al. (2023), and our discussion in appendix, Section L), we build our own dataset with 30k trajectories (*i.e.*, 30M steps). We follow the same approach for BALANCE and PENDULUM for which we are not aware of any existing offline benchmark. Data collection, detailed in appendix (Section D), follows the philosophy of standard benchmarks like D4RL (Fu et al., 2020).

454 455

456

439

440 441 442

443

444

445

446

447

#### 5.2 MAIN RESULTS

457 We present results on BALANCE, SWINGUP and PENDULUM for PRIMORL and baselines in Table 1 458 and Figure 2. Both report policy performance in the real MDP as the mean episodic return over 10 459 episodes per SAC training epoch. Average performance and 95% confidence intervals are computed 460 by re-training the model and the policy from scratch on at least 5 random seeds to assess the stability 461 of the full training process. We also report the corresponding theoretical upper bound on  $\epsilon$ , as 462 computed from the hyperparameters z, q, T and  $\delta$  using the moments accountant  $\epsilon^{MA}(z, q, T, \delta)$  (see 463 Table 2 and discussion in Section H.6 for further explanations regarding the moments accountant).

464 These results show a well-expected trade-off: performance tends to degrade with stronger privacy 465 guarantees (*i.e.*, smaller  $\epsilon$ 's), as the model training gets perturbed with higher levels of noise. Moreover, private model training makes the policy performance less stable over several runs, which 466 is also expected since differential privacy adds another source of randomness during training. We 467 notice that noise is not the sole factor that negatively impacts performance, as suggested by the gap 468 between MOPO and PRIMORL NO PRIVACY: gradient clipping and trajectory-level training also 469 contribute to performance degradation. In some cases, a small amount of DP noise might actually be 470 beneficial, acting as a kind of regularization, as in SWINGUP and PENDULUM. Moreover, experiments 471 on HALFCHEETAH (Section J) show that PRIMORL performs worse in higher-dimensional tasks. 472 This could be expected based on the theoretical analysis led in Section 4.3.1, as DP training adds 473 a dependence on the dimension d of the task in the valuation gap. Despite this trade-off, private 474 agents trained with PRIMORL remain competitive with MOPO for  $\epsilon$  in the 10<sup>1</sup> to 10<sup>2</sup> range. For 475 PENDULUM, we plot policy performance against  $\epsilon$  in Figure 3, and observe even no performance 476 degradation until  $\epsilon$  reaches the 1 to 10 range. Although algorithms from Qiao & Wang (2023a) are not 477 suited for direct comparison on the same tasks, we argue that our empirical results are significantly stronger. Indeed, converting  $\rho$ -zero-concentrated DP guarantees into standard ( $\epsilon, \delta$ )-DP guarantees 478 for clarity and fair comparison, we observe that PRIMORL achieves comparable privacy-performance 479 trade-offs, but on much more complex environments (more details in appendix, Section F). 480

While the privacy budgets  $\epsilon$  from Table 1 do not correspond to strong theoretical privacy guarantees, we must consider the worst-case nature of the differential privacy definition, along with its very strong assumptions on the adversary side. In offline RL especially, the definition of DP assumes the adversary only has to discriminate between two precise neighboring datasets D and  $D' = D \cup \{\tau\}$ as well as the release of all gradients, whereas in practice the adversary faces the much harder task of reconstructing a high-dimensional trajectory based on the output policy and limited side information

4	8	6
4	8	7
4	8	8

502

504

Table 1: Results for PENDULUM,	BALANCE and SWINGUP.
--------------------------------	----------------------

			PENDUL	UM		CARTPO	DLE-BALANCE		CARTPOLE-SWINGUP
Method		ε	Retu	RN	ε		Return	ε	RETURN
МОРО		$\infty$	795.9	$\pm 6.5$	$\infty$	97	$6.3 \pm 26.8$	$\infty$	$804.9\pm89.6$
PRIMORL NO I	PRIV.	$\infty$	$810.4 \pm 27.4$	5 (101.8%)	$\infty$	947.5	± 68.3 ( <b>97.1%</b> )	$\infty$	$774.1 \pm 81.7$ (96.17%)
PRIMORL LOW	r	22.3	$817.4 \pm 21.5$	7 (102.7%)	85.0	815.8	± 97.2 (83.6%)	94.2	$772.4 \pm 73.9$ (95.96%)
PRIMORL HIGH	н	5.1	$778.9 \pm 53.$	.5 ( <b>97.9</b> %)	8.2	$758.2 \pm$	187.2 (77.7%)	17.0	698.3 ± 57.5 ( <b>86.75</b> %
								Performanc	$e = f(\varepsilon)$ on Pendulum
								Performanc	$e = f(\varepsilon)$ on Pendulum
		z	T	q	δ	ε	3000 T J	Performanc	e = f(s) on Pendulum
Pendulum	Low	z 0.35	$T$ $7.10^3$	$q$ $10^{-3}$	$\delta$ $10^{-5}$	ε 22.3	1000 800	Performanc	e = f(ε) on Pendulum
Pendulum	Low High	z 0.35 0.52	T 7.10 <sup>3</sup> 7.10 <sup>3</sup>	q $10^{-3}$ $10^{-3}$	$\delta$ 10 <sup>-5</sup> 10 <sup>-5</sup>	$\epsilon$ 22.3 5.1	1000 800 400	Performanc	e = f(c) on Pendulum
Pendulum Balance	Low High Low	z 0.35 0.52 0.25	$\frac{T}{7.10^{3}}$ 7.10 <sup>3</sup> 7.10 <sup>3</sup>	$\frac{q}{10^{-3}}\\ 10^{-3}\\ 10^{-3}$	$\delta$ $10^{-5}$ $10^{-5}$ $10^{-5}$	$\epsilon$ 22.3 5.1 85.0	000 000 beform	Performanc	e = f(c) on Pendulum
Pendulum Balance	Low High Low High	z 0.35 0.52 0.25 0.45	$\begin{array}{c} T \\ \hline 7.10^3 \\ 7.10^3 \\ 7.10^3 \\ 7.10^3 \end{array}$	$\begin{array}{c} q \\ 10^{-3} \\ 10^{-3} \\ 10^{-3} \\ 10^{-3} \end{array}$	$\delta$ $10^{-5}$ $10^{-5}$ $10^{-5}$ $10^{-5}$	<ul> <li>ε</li> <li>22.3</li> <li>5.1</li> <li>85.0</li> <li>8.2</li> </ul>	Ebisodo a termina (1990)	Performanc	e = f(ε) on Pendulum
Pendulum Balance Swingup	Low High Low High Low	z 0.35 0.52 0.25 0.45 0.25	$\begin{array}{c} T \\ 7.10^{3} \\ 7.10^{3} \\ 7.10^{3} \\ 7.10^{3} \\ 10.10^{3} \end{array}$	$\begin{array}{c} q \\ 10^{-3} \\ 10^{-3} \\ 10^{-3} \\ 10^{-3} \\ 10^{-3} \end{array}$	$\frac{\delta}{10^{-5}}\\ 10^{-5}\\ 10^{-5}\\ 10^{-5}\\ 10^{-5}$	<ul> <li> <i>ϵ</i> 22.3 5.1 85.0 8.2 94.2 </li> </ul>	600 900 900 900 900 900	Performanc	e = f(¢) on Pendulum

Table 2: Hyperparameters and computation of the theo- Figure 3: Policy performance on PENretical privacy budget  $\epsilon := \epsilon^{MA}(z, q, T, \delta)$ 

DULUM as a function of  $\epsilon$ 

505 only. Therefore, backed by recent work on empirical privacy auditing (e.g., Carlini et al. (2019); 506 Ponomareva et al. (2022)), we argue that such  $\epsilon$ 's can provide adequate privacy protection in practical 507 offline RL applications. According to Ponomareva et al. (2023),  $\epsilon \lesssim 10$  is actually a realistic and 508 widely used goal in private deep learning applications. We discuss this matter more in depth in 509 appendix (Section G). We also point out that achieving a strong privacy-utility trade-off in offline RL requires access to datasets with a very large number of trajectories and that current benchmarks, with 510 datasets of only dozens to thousands of trajectories, are insufficient for studying privacy effectively. 511 In contrast, other fields often use datasets containing millions of users (between  $10^6$  to  $10^9$  users in 512 McMahan et al. (2018)) to ensure robust privacy guarantees, which would be very costly to study in 513 offline RL. In appendix (Section L), we provide evidence that increasing dataset size improves the 514 privacy-performance trade-off, demonstrating even greater potential for PRIMORL. 515

516 517

518

#### 6 DISCUSSION

519 While existing DP RL methods are limited to tabular and linear finite-horizon MDPs, we are the 520 first to address deep offline RL with privacy guarantees in the infinite-horizon discounted setting, and propose a model-based approach named PRIMORL. We empirically show that PRIMORL 521 is capable of learning trajectory-level private, neural-based policies in standard control tasks with 522 only limited performance cost, achieving a new standard in differentially private RL. Although the 523 reported privacy budgets are typically considered too large to stand as formal DP guarantees, we 524 argue based on recent studies on practical DP that they can offer satisfying privacy protection in 525 practice, especially considering the worst-case nature of DP which can yield too pessimistic privacy 526 budgets. Empirical evaluation of the robustness of our algorithm against privacy attacks, for which 527 a rigorous and standardized benchmark has to be developed, will thus be an important research 528 direction for future work. We further point out that our approach has the potential for achieving 529 greater privacy-utility trade-offs given access to large enough offline datasets, hence calling for new 530 benchmarks in the increasingly important field of private offline RL.

531 With the aim of shifting the paradigm in how private RL is approached — from predominantly 532 theoretical research to practical algorithms — this work sets the stage for future efforts to scale 533 to higher-dimensional problems. We identify several promising research avenues. First, we may 534 consider limiting the number of real trajectories used during training to leverage privacy amplification through sub-sampling, for example, by using data augmentation techniques. As our theoretical 536 analysis highlights the impact of dimensionality on private model error, another promising direction is learning compact representations of high-dimensional inputs and performing planning directly in the latent space, as explored by Jiang et al. (2023). We leave these avenues for future work. Overall, 538 we believe our work represents a significant step toward the much-needed deployment of private RL methods in practical applications.

#### 540 REFERENCES 541

585

542 543 544 545	Martín Abadi, Andy Chu, Ian J. Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In <i>Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security</i> , pp. 308–318. ACM, 2016. URL https://doi.org/10.1145/2976749.2978318.
546 547 548	Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. In 9th International Conference on Learning Representations, ICLR, 2021. URL https://openreview.net/forum?id=OMNB1G5xzd4.
549 550 551 552	Borja Balle, Gilles Barthe, and Marco Gaboardi. Privacy amplification by subsampling: Tight analyses via couplings and divergences. In <i>Proceedings of NeurIPS</i> , 2018. URL https://proceedings.neurips.cc/paper/2018/hash/3b5020bb891119b9f5130f1fea9bd773-Abstract.html.
553 554 555 556 557	Raef Bassily, Adam D. Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In <i>55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014</i> , pp. 464–473. IEEE Computer Society, 2014. URL https://doi.org/10.1109/FOCS.2014.56.
558 559 560	Debabrota Basu, Christos Dimitrakakis, and Aristide C. Y. Tossou. Differential privacy for multi- armed bandits: What is it and what is its cost? <i>CoRR</i> , abs/1905.12298, 2019. URL http: //arxiv.org/abs/1905.12298.
561 562 563	Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. <i>CoRR</i> , abs/1606.01540, 2016. URL http://arxiv.org/abs/1606.01540.
564 565 566 567 568 569	Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In <i>Theory of Cryptography - 14th International Conference, TCC 2016-B, Beijing, China, October 31 - November 3, 2016, Proceedings, Part I, volume 9985 of Lecture Notes in Computer Science</i> , pp. 635–658, 2016. URL https://doi.org/10.1007/978-3-662-53641-4_24.
570 571 572 573 574	Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In Nadia Heninger and Patrick Traynor (eds.), 28th USENIX Security Symposium, USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019, pp. 267–284. USENIX Association, 2019. URL https://www. usenix.org/conference/usenixsecurity19/presentation/carlini.
575 576 577	Sayak Ray Chowdhury and Xingyu Zhou. Differentially Private Regret Minimization in Episodic Markov Decision Processes, December 2021. URL http://arxiv.org/abs/2112.10599. arXiv:2112.10599 [cs, math].
578 579 580 581	Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforce- ment learning in a handful of trials using probabilistic dynamics models. In <i>Proceedings</i> of <i>NeurIPS</i> , 2018. URL https://proceedings.neurips.cc/paper/2018/hash/ 3de568f8597b94bda53149c7d7f5958c-Abstract.html.
583 584	Chris Cundy, Rishi Desai, and Stefano Ermon. Privacy-constrained policies via mutual information regularized policy gradients. In <i>International Conference on Artificial Intelligence and Statistics</i> , volume 238 of <i>Proceedings of Machine Learning Research</i> , pp. 2809–2817. PMLR, 2024. URL

586 Floris den Hengst, Eoin Grua, Ali el Hassouni, and Mark Hoogendoorn. Reinforcement learning 587 588 for personalization: A systematic literature review. Data Science, 3:1-41, 04 2020. doi: 10.3233/ DS-200028. 589

https://proceedings.mlr.press/v238/j-cundy24a.html.

- 590 Cynthia Dwork. Differential Privacy. In Proceedings of ICALP, 2006. URL https://www. 591 microsoft.com/en-us/research/publication/differential-privacy/. 592
- Cynthia Dwork and Guy N. Rothblum. Concentrated differential privacy. CoRR, abs/1603.01887, 593 2016. URL http://arxiv.org/abs/1603.01887.

594

595 Annual IEEE Symposium on Foundations of Computer Science, FOCS, pp. 51–60. IEEE Computer 596 Society, 2010. URL https://doi.org/10.1109/FOCS.2010.12. 597 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: datasets for 598 deep data-driven reinforcement learning. CoRR, abs/2004.07219, 2020. URL https://arxiv. org/abs/2004.07219. 600 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without 601 exploration. In Proceedings of the 36th International Conference on Machine Learning, ICML 602 2019, volume 97 of Proceedings of Machine Learning Research, pp. 2052–2062. PMLR, 2019. 603 URL http://proceedings.mlr.press/v97/fujimoto19a.html. 604 605 Evrard Garcelon, Vianney Perchet, Ciara Pike-Burke, and Matteo Pirotta. Local differ-606 ential privacy for regret minimization in reinforcement learning. In Proceedings of 607 *NeurIPS*, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/ 580760fb5def6e2ca8eaf601236d5b08-Abstract.html. 608 609 Maziar Gomrokchi, Susan Amin, Hossein Aboutalebi, Alexander Wong, and Doina Precup. Mem-610 bership inference attacks against temporally correlated data in deep reinforcement learning. 611 IEEE Access, 11:42796-42808, 2023. URL https://doi.org/10.1109/ACCESS.2023. 612 3270860. 613 Çaglar Gülçehre, Ziyu Wang, Alexander Novikov, Thomas Paine, Sergio Gómez Colmenarejo, Konrad 614 Zolna, Rishabh Agarwal, Josh Merel, Daniel J. Mankowitz, Cosmin Paduraru, Gabriel Dulac-615 Arnold, Jerry Li, Mohammad Norouzi, Matthew Hoffman, Nicolas Heess, and Nando de Freitas. 616 RL unplugged: A collection of benchmarks for offline reinforcement learning. In Proceedings 617 of NeurIPS, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 618 51200d29d1fc15f5a71c1dab4bb54f7c-Abstract.html. 619 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy 620 maximum entropy deep reinforcement learning with a stochastic actor. In Proceedings of the 621 35th International Conference on Machine Learning, ICML 2018, volume 80 of Proceedings 622 of Machine Learning Research, pp. 1856–1865. PMLR, 2018. URL http://proceedings. 623 mlr.press/v80/haarnoja18b.html. 624 Zhengyao Jiang, Tianjun Zhang, Michael Janner, Yueying Li, Tim Rocktäschel, Edward Grefenstette, 625 and Yuandong Tian. Efficient planning in a compact latent action space. In The Eleventh 626 International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. 627 OpenReview.net, 2023. URL https://openreview.net/forum?id=cA77NrVEuqn. 628 629 Yilin Kang, Jian Li, Yong Liu, and Weiping Wang. Data heterogeneity differential privacy: From theory to algorithm. In Computational Science - ICCS 2023 - 23rd International Conference, 630 Prague, Czech Republic, July 3-5, 2023, Proceedings, Part I, volume 14073 of Lecture Notes 631 in Computer Science, pp. 119-133. Springer, 2023. URL https://doi.org/10.1007/ 632 978-3-031-35995-8\_9. 633 634 Michael J. Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Mach.* 635 Learn., 49(2-3):209-232, 2002. URL https://doi.org/10.1023/A:1017984413808. 636 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. 637 MOReL: Model-based offline reinforcement learning. In Proceedings of NeurIPS, 638 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 639 f7efa4f864ae9b88d43527f4b14f750f-Abstract.html. 640 Byeongchan Kim and Min Hwan Oh. Model-based offline reinforcement learning with count-641 based conservatism. In International Conference on Machine Learning, ICML 2023, volume 642 202 of Proceedings of Machine Learning Research, pp. 16728–16746. PMLR, 2023. URL 643 https://proceedings.mlr.press/v202/kim23q.html. 644 B. Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Kumar 645 Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. 646 IEEE Trans. Intell. Transp. Syst., 23(6):4909-4926, 2022. URL https://doi.org/10. 647

Cynthia Dwork, Guy N. Rothblum, and Salil P. Vadhan. Boosting and differential privacy. In 51th

12

1109/TITS.2021.3054625.

648 649 650	Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. <i>CoRR</i> , abs/2005.01643, 2020. URL https://arxiv.org/abs/2005.01643.
651 652 653 654	Chonghua Liao, Jiafan He, and Quanquan Gu. Locally differentially private reinforcement learning for linear mixture markov decision processes. <i>CoRR</i> , abs/2110.10133, 2021. URL https://arxiv.org/abs/2110.10133.
655 656 657 658	Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In <i>4th International Conference on Learning Representations, ICLR 2016</i> , 2016. URL http://arxiv.org/abs/1509.02971.
659 660 661 662 663	Changchang Liu, Supriyo Chakraborty, and Prateek Mittal. Dependence makes you vulnberable: Differential privacy under dependent tuples. In 23rd Annual Network and Distributed System Security Symposium, NDSS 2016, San Diego, California, USA, February 21-24, 2016. The Inter- net Society, 2016. URL https://www.princeton.edu/~pmittal/publications/ ddp-ndss16.pdf.
665 666	Mingyang Liu, Xiaotong Shen, and Wei Pan. Deep reinforcement learning for personalized treatment recommendation. <i>Statistics in Medicine</i> , 41, 06 2022.
667 668 669	Siqi Liu, Kay Choong See, Kee Yuan Ngiam, Leo Anthony Celi, Xingzhi Sun, and Mengling Feng. Reinforcement learning for clinical decision support in critical care: Comprehensive review. <i>J Med Internet Res</i> , 22(7):e18477, Jul 2020. URL https://www.jmir.org/2020/7/e18477.
670 671 672 673 674	Cong Lu, Philip J. Ball, Jack Parker-Holder, Michael A. Osborne, and Stephen J. Roberts. Revisiting design choices in offline model based reinforcement learning. In <i>The Tenth International Conference on Learning Representations, ICLR 2022</i> , 2022. URL https://openreview.net/forum?id=zz9hXVhf40.
675 676 677	Paul Luyo, Evrard Garcelon, Alessandro Lazaric, and Matteo Pirotta. Differentially private explo- ration in reinforcement learning with linear representation, 2021. URL https://arxiv.org/ abs/2112.01585.
678 679 680 681 682	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Proceedings</i> of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, volume 54 of <i>Proceedings of Machine Learning Research</i> , pp. 1273–1282. PMLR, 2017. URL http://proceedings.mlr.press/v54/mcmahan17a.html.
683 684 685 686	<ul> <li>H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. In 6th International Conference on Learning Representations, ICLR, 2018. URL https://openreview.net/forum?id=BJ0hF1Z0b.</li> </ul>
687 688 689	Ilya Mironov. Rényi differential privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF). IEEE, aug 2017a. URL https://doi.org/10.1109%2Fcsf.2017. 11.
690 691 692	Ilya Mironov. Renyi Differential Privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF), pp. 263–275, 2017b. URL http://arxiv.org/abs/1702.07476. arXiv:1702.07476 [cs].
693 694 695 696	Thomas M. Moerland, Joost Broekens, Aske Plaat, and Catholijn M. Jonker. Model-based re- inforcement learning: A survey. <i>Found. Trends Mach. Learn.</i> , 16(1):1–118, 2023. URL https://doi.org/10.1561/220000086.
697 698 699 700 701	Dung Daniel T. Ngo, Giuseppe Vietri, and Steven Wu. Improved regret for differentially private exploration in linear MDP. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), <i>International Conference on Machine Learning, ICML 2022</i> , 17-23 July 2022, Baltimore, Maryland, USA, volume 162 of Proceedings of Machine Learning Research, pp. 16529–16552. PMLR, 2022. URL https://proceedings.mlr.press/ v162/ngo22a.html.

702 Xinlei Pan, Weiyao Wang, Xiaoshuai Zhang, Bo Li, Jinfeng Yi, and Dawn Song. How You Act Tells 703 a Lot: Privacy-Leaking Attack on Deep Reinforcement Learning. Reinforcement Learning, 2019. 704 Natalia Ponomareva, Jasmijn Bastings, and Sergei Vassilvitskii. Training text-to-text transformers 705 with privacy guarantees. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), 706 Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 707 2022, pp. 2182–2193. Association for Computational Linguistics, 2022. doi: 10.18653/v1/2022. 708 findings-acl.171. URL https://doi.org/10.18653/v1/2022.findings-acl.171. 709 710 Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H. Brendan 711 McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Thakurta. How to DP-fy ML: 712 A Practical Guide to Machine Learning with Differential Privacy, March 2023. URL http: //arxiv.org/abs/2303.00654. arXiv:2303.00654 [cs, stat]. 713 714 Kritika Prakash, Fiza Husain, Praveen Paruchuri, and Sujit Gujar. How Private Is Your RL Policy? 715 An Inverse RL Based Analysis Framework. Proceedings of the AAAI Conference on Artificial 716 Intelligence, 36(7):8009–8016, June 2022. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v36i7. 717 20772. URL https://ojs.aaai.org/index.php/AAAI/article/view/20772. 718 719 Rafael Figueiredo Prudencio, Marcos R. O. A. Máximo, and Esther Luna Colombini. A survey on offline reinforcement learning: Taxonomy, review, and open problems. CoRR, abs/2203.01387, 720 2022. URL https://doi.org/10.48550/arXiv.2203.01387. 721 722 Dan Qiao and Yu-Xiang Wang. Offline reinforcement learning with differential privacy. In Proceed-723 ings of NeurIPS, 2023a. URL http://papers.nips.cc/paper\_files/paper/2023/ 724 hash/claaf7c3f306fe94f77236dc0756d771-Abstract-Conference.html. 725 Dan Qiao and Yu-Xiang Wang. Near-optimal differentially private reinforcement learning. In 726 Francisco J. R. Ruiz, Jennifer G. Dy, and Jan-Willem van de Meent (eds.), International Conference 727 on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia, Spain, 728 volume 206 of Proceedings of Machine Learning Research, pp. 9914–9940. PMLR, 2023b. URL 729 https://proceedings.mlr.press/v206/giao23a.html. 730 731 Maria Rigaki and Sebastian Garcia. A survey of privacy attacks in machine learning. CoRR, 732 abs/2007.07646, 2020. URL https://arxiv.org/abs/2007.07646. 733 Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Stochastic convex 734 optimization. In COLT 2009, 2009. URL http://www.cs.mcgill.ca/%7Ecolt2009/ 735 papers/018.pdf#page=1. 736 737 R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine 738 learning models. In 2017 IEEE Symposium on Security and Privacy (SP), pp. 3–18. IEEE Computer 739 Society, 2017. URL https://doi.ieeecomputersociety.org/10.1109/SP.2017. 740 41. 741 Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. Reinforcement learning in robotic applications: 742 a comprehensive survey. Artif. Intell. Rev., 55(2):945-990, 2022. URL https://doi.org/ 743 10.1007/s10462-021-09997-9. 744 745 Richard S. Sutton and Andrew G. Barto. Reinforcement learning - an introduction. Adaptive 746 computation and machine learning. MIT Press, 1998. ISBN 978-0-262-19398-6. URL https: //www.worldcat.org/oclc/37293240. 747 748 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, 749 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy P. Lillicrap, and Martin A. Riedmiller. 750 Deepmind control suite. CoRR, abs/1801.00690, 2018. URL http://arxiv.org/abs/ 751 1801.00690. 752 753 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2012, 754 Vilamoura, Algarve, Portugal, October 7-12, 2012, pp. 5026–5033. IEEE, 2012. URL https: 755 //doi.org/10.1109/IROS.2012.6386109.

 Aristide Tossou and Christos Dimitrakakis. Algorithms for Differentially Private Multi-Armed Bandits. In Proceedings of AAAI, 2016. URL https://aaai.org/papers/ 212-algorithms-for-differentially-private-multi-armed-bandits/.

- Giuseppe Vietri, Borja Balle, Akshay Krishnamurthy, and Zhiwei Steven Wu. Private reinforcement learning with PAC and regret guarantees. In *Proceedings of the 37th International Conference* on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pp. 9754–9764. PMLR, 2020. URL http://proceedings. mlr.press/v119/vietri20a.html.
- Baoxiang Wang and Nidhi Hegde. Privacy-preserving q-learning with functional noise in continuous spaces. In Advances in Neural Information Processing Systems, volume 32, 2019. URL https://proceedings.neurips.cc/paper\_files/paper/2019/file/6646b06b90bd13dabc11ddba01270d23-Paper.pdf.
- Pawel Wawrzynski. A cat-like robot real-time learning to run. In Adaptive and Natural Computing Algorithms, 9th International Conference, ICANNGA 2009, Kuopio, Finland, April 23-25, 2009, Revised Selected Papers, volume 5495 of Lecture Notes in Computer Science, pp. 380–390.
   Springer, 2009. URL https://doi.org/10.1007/978-3-642-04921-7\_39.
- Tian Xu, Ziniu Li, and Yang Yu. Error bounds of imitating policies and environments. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/b5c01503041b70d41d80e3dbe31bbd8c-Abstract.html.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y. Zou, Sergey Levine, Chelsea
   Finn, and Tengyu Ma. MOPO: model-based offline policy optimization. In *Proceedings* of *NeurIPS*, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ a322852ce0df73e204b7e67cbbef0d0a-Abstract.html.
- Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and
  Zhenhui Li. DRN: A deep reinforcement learning framework for news recommendation. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*, pp. 167–176. ACM, 2018. URL https://doi.org/10.1145/
  3178876.3185994.
- Xingyu Zhou. Differentially Private Reinforcement Learning with Linear Function Approximation.
   *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 6(1):1–27, 2022.
   URL https://dl.acm.org/doi/10.1145/3508028.

792 793 794

807 808 809

## <sup>810</sup> A PROOFS

811 812

816

850

851

852

854

855

856

857 858 859

**Theorem 4.2.**  $(\epsilon, \delta)$ -TDP guarantees for dynamics model. Given  $\delta \in (0, 1)$ , noise multiplier z, sampling ratio q and number of training iterations T, let  $\epsilon := \epsilon^{MA}(z, q, T, \delta)$  be the privacy budget computed by the moments accounting method from (Abadi et al. (2016), more details in Section H.6). The dynamics model output by Algorithm 1 is  $(\epsilon, \delta)$ -TDP.

817 *Proof.* Theorem 1 from McMahan et al. (2018) shows that the moments accounting method from 818 Abadi et al. (2016) computes correctly the privacy loss of DP-FEDAVG at user-level for the noise 819 multiplier  $z = \sigma/\mathbb{C}$  with  $\mathbb{C} = C/qK$  if, for each user  $u_k$ , the clipped gradient  $\Delta_k^{\text{clipped}}$  computed 820 from  $u_k$ 's data has sensitivity bounded by C (referred to as **Condition 1**). With TDP MODEL 821 ENSEMBLE TRAINING, we train the model ensemble as a single big model: at each training iteration, 822 the same input batch is processed forward by all models in a single pass, a single loss is computed 823 for the ensemble, and the parameters are then updated in a single backward pass. The ensemble 824 of models can therefore be seen as a concatenation of all individual models, equivalent to a larger model  $\theta = (\theta_i)_{i=1}^N$ . We can therefore extend this theorem by mapping users in federating learning to 825 trajectories in offline RL, as long as **Condition 1** holds for every trajectory  $\tau_k$ . 826

Since we use **ensemble clipping**, we verify that, for trajectory  $\tau_k$ , the ensemble gradient  $\Delta_k^{\text{CLIPPED}} = \begin{pmatrix} \Delta_{i,k}^{\text{CLIPPED}} \end{pmatrix}$  has sensitivity bounded by C. With **flat ensemble clipping**, the gradient of each model  $i \in [\![1, N]\!]$  is clipped by a factor  $C_i = \frac{C}{\sqrt{N}}$  (see Algorithm 3). By construction,  $\Delta_{i,k}^{\text{CLIPPED}}$  has sensitivity bounded by  $C_i$ , *i.e.*,  $\max_{d(D,D')=1} \|\Delta_{i,k}^{\text{CLIPPED}}(D) - \Delta_{i,k}^{\text{CLIPPED}}(D')\|_2 \le C_i$ . Therefore, for two neighboring datasets D and D':

$$\begin{split} \|\Delta_k^{\text{CLIPPED}}(D) - \Delta_k^{\text{CLIPPED}}(D')\|_2 &= \|\left(\Delta_k^{\text{CLIPPED}}(D) - \Delta_k^{\text{CLIPPED}}(D')\right)_{i=1}^N\|_2 \\ &= \sqrt{\sum_{i=1}^N \|\Delta_{i,k}^{\text{CLIPPED}}(D) - \Delta_{i,k}^{\text{CLIPPED}}(D')\|_2^2} \\ &\leq \sqrt{\sum_{i=1}^N C_i^2} \\ &= \sqrt{\sum_{i=1}^N \frac{C^2}{N}} \\ &= C \ . \end{split}$$

847 This implies  $\max_{d(D,D')=1} \|\Delta_k^{\text{CLIPPED}}(D) - \Delta_k^{\text{CLIPPED}}(D')\|_2 \le C$ :  $\Delta_k^{\text{CLIPPED}}$  has sensitivity bounded 848 by *C*. We can derive the same proof for **per-layer ensemble clipping**. Therefore, Theorem 1 from 849 McMahan et al. (2018) holds for TDP MODEL ENSEMBLE TRAINING.

We can therefore use the moments accountant  $\epsilon^{MA}$  to compute, given z > 0,  $\delta \in (0, 1)$ ,  $q \in (0, 1)$ and  $T \in \mathbb{N}$ , the total privacy budget  $\epsilon$  spent by Algorithm 1, *i.e.*,  $\epsilon = \epsilon^{MA}(z, q, T, \delta)$ .

The dyanmics model output by Algorithm 1 is therefore  $(\epsilon, \delta)$ -TDP.

**Theorem 4.5.**  $(\epsilon, \delta)$ -TDP guarantees for PRIMORL. Given an  $(\epsilon, \delta)$ -TDP model  $(\hat{P}, \hat{r})$  learned with Algorithm 1, the policy obtained with private policy optimization (Algorithm 4) within the pessimistic model  $(\hat{P}, \hat{r} - \lambda u)$  is  $(\epsilon, \delta)$ -TDP.

860 *Proof.* First, we establish that the pessimistic MDP  $\tilde{M}$  is private for  $u \in \{u_{MA}, u_{MPD}\}$ . By Theo-861 rem 4.2, both the mean estimators  $\{\mu_{\phi_i}\}_{i=1}^N$  the covariance estimators  $\{\Sigma_{\psi_i}\}_{i=1}^N$  are private. There-862 fore, both uncertainty estimators  $u_{MA}(s, a) = \|\Sigma_{\psi_i}(s, a)\|_F$  and  $u_{MPD}(s, a) = \max_{i,j} \|f_{\phi_i} - f_{\phi_j}\|_2$ , 863 as data-independent transformations of the above quantities, are also private thanks to the postprocessing property of DP. Therefore, the pessimistic model  $\tilde{M}$  remains  $(\epsilon, \delta)$ -TDP. Now, we can think of SAC model-based policy optimization (Algorithm 4) as an abstract, randomized function  $h_{\Pi}$ , that takes as input  $\hat{M}$  and outputs as policy  $\hat{\pi}$ . Furthermore, let  $h_M$  denote the mechanism that takes as input the private offline dataset  $\mathcal{D}_K$  and outputs the private pessimistic model  $\hat{M}$ , and which is  $(\epsilon, \delta)$ -TDP following 4.2. We observe that  $h = h_{\Pi} \circ h_M$ , where h is the global offline RL algorithm which is the object of Definition 4.1. Since SAC only uses data from the model, as stated in Section 4.3.3,  $h_{\Pi}$  is independent of the private offline data  $\mathcal{D}_K$ . In other words,  $h_{\Pi}$  is a data-independent transformation of the private mechanism  $h_M$ . Thanks again to the post-processing property of differential privacy, h is also  $(\epsilon, \delta)$ -TDP.

872

879 880

888 889 890

891 892 893

894 895

896

897

899 900

905 906 907

912 913 914

917

873 We now prove the following two propositions:

**Proposition 4.3.** Value evaluation error in non-private offline MBRL. Let the model loss function be L-Lipschitz and  $\Delta$ -strongly convex, and assumptions from the simulation lemma hold. There is a stochastic convex optimization algorithm for learning the model and a constant M such that, with probability at least  $1 - \alpha$ , and for sufficiently large  $N_D$ , the value evaluation error of  $\pi$  is bounded as:

$$|\hat{V}^{\pi} - V^{\pi}| \le \frac{\sqrt{2}\gamma}{(1-\gamma)^2} \cdot M \cdot \frac{L \log^{1/2}(N_{\mathcal{D}}/\alpha)}{\sqrt{\Delta N_{\mathcal{D}}}} + \frac{2\sqrt{2}\gamma}{(1-\gamma)^2} \sqrt{\varepsilon_{\pi}} .$$

**Proposition 4.4.** Value evaluation error in private offline MBRL. Let assumptions from Proposition 4.3 hold. If the model is learned with  $(\epsilon, \delta)$ -DP gradient descent, then, with probability at least  $1 - \alpha$ , there is a constant M' such that for large enough  $N_D$ , the value evaluation error of  $\pi$ :

$$|\hat{V}_{DP}^{\pi} - V^{\pi}| \leq \frac{\sqrt{2}\gamma}{(1-\gamma)^2} \cdot M' \cdot \frac{Ld^{1/4}\log(N_{\mathcal{D}}/\delta) \cdot \operatorname{poly}\log(1/\alpha)}{\sqrt{\Delta N_{\mathcal{D}}\epsilon\alpha}} + \frac{2\sqrt{2}\gamma}{(1-\gamma)^2}\sqrt{\varepsilon_{\pi}}$$

*Proof.* Let  $\mathcal{F}$  denote the function class of the model. The model is estimated by maximizing the likelihood of the data  $\mathcal{D}_K = (s_i, a_i, s'_i)_{i=1}^N$ , which is collected by an unknown behavioral policy  $\pi^B$ . This is equivalent to minimizing the negative log-likelihood. The population risk of the estimated model  $\hat{P}$  obtained with DP-SGD, is therefore:

$$\mathcal{L}(\hat{P}) = \mathbb{E}_{(s,a)\sim\rho_P^{\pi^B}, s'\sim P(\cdot|s,a)} \left[ -\log \hat{P}(s'|s,a) \right] ,$$

where  $\rho_P^{\pi^B}$  is the (normalized) state-action occupancy measure under policy  $\pi^B$  and dynamics P. Let us further assume that the true model P belongs to the function class  $\mathcal{F}$ , and that  $P \in \operatorname{argmin}_{P' \in \mathcal{F}} \mathcal{L}(P')$ . We can therefore write the excess population risk of the model estimator  $\hat{P}$  as:

$$\mathcal{L}(\hat{P}) - \mathcal{L}(P) = \mathbb{E}_{(s,a) \sim \rho_P^{\pi^B}, s' \sim P(\cdot|s,a)} \left[ \frac{\log P(s'|s,a)}{\log \hat{P}(s'|s,a)} \right]$$

But, denoting  $D_{KL}(A, B)$  the Kullback-Leibler divergence between distributions A, B:

$$D_{\mathrm{KL}}\left(P(s,a),\hat{P}(s,a)\right) = \mathbb{E}_{s' \sim P(\cdot|s,a)}\left[\frac{\log P(s'|s,a)}{\log \hat{P}(s'|s,a)}\right]$$

We can therefore rewrite the above excess population risk as:

$$\mathcal{L}(\hat{P}) - \mathcal{L}(P) = \mathbb{E}_{(s,a) \sim \rho_P^{\pi^B}} \left[ D_{\mathrm{KL}} \left( P(s,a), \hat{P}(s,a) \right) \right] \quad . \tag{3}$$

If the objective function  $\mathcal{L}$  is *L*-Lipschitz and  $\Delta$ -strongly convex, Bassily et al. (2014) shows (Theorem F.2) that, given  $N_{\mathcal{D}}$  data points, a noisy gradient descent algorithm with  $(\epsilon, \delta)$ -DP guarantees satisfies, with probability at least  $1 - \alpha$ :

$$\mathcal{L}(\hat{P}) - \mathcal{L}(P) = \mathcal{O}\left(\frac{L^2\sqrt{d}\log^2(N_{\mathcal{D}}/\delta) \cdot \operatorname{poly}\log(1/\alpha)}{\Delta N_{\mathcal{D}}\epsilon\alpha}\right) \quad . \tag{4}$$

In the non-private case, Shalev-Shwartz et al. (2009) provides the following bound under the same assumptions:

$$\mathcal{L}(\hat{P}) - \mathcal{L}(P) = \mathcal{O}\left(\frac{L^2 \log(N_{\mathcal{D}}/\alpha)}{\Delta N_{\mathcal{D}}}\right) \quad .$$
(5)

On the other hand, we have from the Simulation Lemma (Kearns & Singh, 2002; Xu et al., 2020) that for a MDP  $\mathcal{M}$  with reward upper bounded by  $r_{\text{max}} = 1$  and dynamics P, a behavioral policy  $\pi^B$  and a learned transition model  $\hat{P}$  with:

$$\mathbb{E}_{(s,a)\sim\rho_P^{\pi^B}}\left[D_{\mathrm{KL}}\left(P(s,a),\hat{P}(s,a)\right)\right] \leq \varepsilon_M \quad , \tag{6}$$

which by 3 is equivalent to:

$$\mathcal{L}(\hat{P}) - \mathcal{L}(P) \le \varepsilon_M \quad , \tag{7}$$

Let  $\pi$  be an arbitrary policy. If the divergence between  $\pi$  and the behavioral policy is bounded:

$$\max_{s} D_{\mathrm{KL}}\left(\pi(\cdot|s), \pi^{B}(\cdot|s)\right) \le \varepsilon_{\pi} \quad , \tag{8}$$

930 then the value evaluation error of  $\pi$  is bounded as:

$$|\hat{V}^{\pi} - V^{\pi}| \le \frac{\sqrt{2\gamma}}{(1-\gamma)^2} \sqrt{\varepsilon_M} + \frac{2\sqrt{2\gamma}}{(1-\gamma)^2} \sqrt{\varepsilon_\pi} \quad . \tag{9}$$

Since f(x) = O(g(x)) implies  $\sqrt{f(x)} = O(\sqrt{g(x)})^2$ , we note that we can replace  $\sqrt{\varepsilon_M}$  in the model term of the right-hand side of 9 by the (square root of) the bounds from 4 and 5 in the private case and in the non-private case, respectively.

938 This result holds for any policy  $\pi$  verifying 8. In particular, if:

$$\max_{s} D_{\mathrm{KL}}\left(\hat{\pi}(\cdot|s), \pi^{B}(\cdot|s)\right) \le \varepsilon_{\hat{\pi}} \quad , \tag{10}$$

then:

$$|\hat{V}^{\hat{\pi}} - V^{\hat{\pi}}| \le \frac{\sqrt{2}\gamma}{(1-\gamma)^2}\sqrt{\varepsilon_M} + \frac{2\sqrt{2}\gamma}{(1-\gamma)^2}\sqrt{\varepsilon_{\hat{\pi}}} \quad .$$
(11)

_			

<sup>2</sup>Indeed, for 
$$f(x)$$
 positive, for any  $x \ge x_0$ ,  $|f(x)| = f(x) \le M' \times g(x)$ , then, for any  $x \ge x_0$ ,  $\sqrt{f(x)} = |\sqrt{f(x)}| \le \sqrt{M'} \times \sqrt{g(x)} = M \times \sqrt{g(x)}$ 

#### 972 B RELATED WORK (EXTENDED) 973

# 974 B.1 MODEL-BASED OFFLINE REINFORCEMENT LEARNING

976 Unlike classical RL (Sutton & Barto, 1998) which is online in nature, offline RL (Levine et al., 977 2020; Prudencio et al., 2022) aims at learning and controlling autonomous agents without further 978 interactions with the system. This approach is preferred or even unavoidable in situations where data collection is impractical (see for instance Singh et al. (2022); Liu et al. (2020); Kiran et al. (2022)). 979 980 Model-based RL (Moerland et al., 2023) can also help when data collection is expensive or unsafe as a good model of the environment can generalize beyond in-distribution trajectories and allow 981 simulations. Moreover, model-based RL has been shown to be generally more sample efficient than 982 model-free RL (Chua et al., 2018). Argenson & Dulac-Arnold (2021) also show that model-based 983 offline planning, where the model is learned offline on a static dataset and subsequently used for 984 control without further accessing the system, is a viable approach to control agents on robotic-like 985 tasks with good performance. Unfortunately, the offline setting comes with its own major challenges. 986 In particular, when the data is entirely collected beforehand, we are confronted to the problem of 987 distribution shift (Fujimoto et al., 2019): as the logging policy used to collect the training dataset only 988 covers a limited (and potentially small) region of the state-action space, the model can only be trusted 989 in this region, and may be highly inaccurate in other parts of the space. This can lead to a severe 990 decrease in the performance of classic RL methods, particularly in the model-based setting where the acting agent may exploit these inaccuracies in the model, causing large gap between performances 991 in the true and the learned environment. MOPO (Yu et al., 2020) and MOREL (Kidambi et al., 992 2020), and more recently COUNT-MORL (Kim & Oh, 2023) have effectively tackled this issue 993 by penalizing the reward proportionally to the model's uncertainty, achieving impressive results on 994 popular offline benchmarks. Nonetheless, there remain many areas for improvement, as highlighted 995 by Lu et al. (2022), which extensively study and challenge key design choices in offline MBRL 996 algorithms.

997 998 999

#### B.2 PRIVACY IN REINFORCEMENT LEARNING

1000 Differential Privacy (DP), first formalized in Dwork (2006), has become the gold standard in terms 1001 of privacy protection. Over the recent years, the design of algorithms with better privacy-utility 1002 trade-offs has been a major line of research. In particular, relaxations of differential privacy and more 1003 advanced composition tools have allowed tighter analysis of privacy bounds (Dwork et al., 2010; 1004 Dwork & Rothblum, 2016; Bun & Steinke, 2016; Mironov, 2017a). Leveraging these advances, the 1005 introduction of DP-SGD (Abadi et al., 2016) has allowed to design private deep learning algorithms, paving the way towards a wider adoption of DP in real-world settings, although the practicalities of differential privacy remain challenging (Ponomareva et al., 2023). In parallel to the theoretical 1007 analysis of privacy, many works have focused on designing more and more sophisticated attacks, 1008 justifying further the need to design DP algorithms ((Rigaki & Garcia, 2020)). 1009

1010 Recent works on RL-specific attacks (Pan et al., 2019; Prakash et al., 2022; Gomrokchi et al., 2023) 1011 have demonstrated that reinforcement learning (RL) is no more immune to privacy threats. With RL being increasingly used to provide personalized services (den Hengst et al., 2020), which may 1012 expose sensitive user data, developing privacy-preserving techniques for training policies has become 1013 crucial. Shortly after DP was successfully extended to multi-armed bandits (Tossou & Dimitrakakis, 1014 2016; Basu et al., 2019), a substantial body of work (e.g., Vietri et al. (2020); Garcelon et al. (2021); 1015 Liao et al. (2021); Luyo et al. (2021); Chowdhury & Zhou (2021); Zhou (2022); Ngo et al. (2022); 1016 Qiao & Wang (2023b)) addressed privacy in online RL, extending definitions from bandits. However, 1017 relying on count-based and UCB-like methods, current RL algorithms with formal DP guarantees 1018 are essentially limited to episodic tabular or linear MDPs, and have not been assessed empirically 1019 beyond simple numerical simulations. However, current RL algorithms with formal DP guarantees 1020 are essentially limited to episodic tabular or linear MDPs, and have not been assessed empirically 1021 beyond simple numerical simulations. Few works have proposed private RL methods for more 1022 general problems, however with significant limitations or in different contexts. Wang & Hegde (2019) 1023 tackle continuous state spaces by adding functional noise to Q-Learning, but the approach is restricted to unidimensional states and focuses on protecting reward information. Recently, Cundy et al. (2024) 1024 addressed high-dimensional control and robotic tasks; however, they consider a specific notion of 1025 privacy that protects sensitive state variables based on a mutual information framework.

Despite the relevance of the setting for real-world RL deployments, private offline RL has received comparatively less attention. To date, only Qiao & Wang (2023a) have proposed DP offline algorithms, building on non-private value iteration methods. While their approach lays the groundwork for private offline RL and offers strong theoretical guarantees, it remains limited to episodic tabular and linear MDPs. Consequently, no existing work has introduced DP methods that can handle deep RL environments in the infinite-horizon discounted setting, a critical step toward deploying private RL algorithms in real-world applications. With this work, we aim to fill this gap by proposing a differentially private, deep model-based RL method for the offline setting. 

## <sup>1080</sup> C PRESENTATION OF THE TASKS

## D DATA COLLECTION

To collect our offline dataset for CARTPOLE and PENDULUM, we used DDPG (Lillicrap et al., 2016), a model-free RL algorithm for continuous action spaces. We ran 600 independent runs of 50,000 steps each for CARTPOLE-BALANCE, 150 independent runs of 200,000 steps each for CARTPOLE-SWINGUP, and 6 independent runs of 1M steps each for PENDULUM. We collect all training episodes to ensure a correct mix between random, medium and expert episodes (similar to *replay* datasets in Fu et al. (2020)).

#### E BASELINES

## Table 3: PRIMORL configurations.

VARIANT	TRAJECTORY-LEVEL ENS. TRAINING	CLIP	NOISE	
NO CLIP	1	X	X	
NO PRIVACY	$\checkmark$	1	X	e
Low, High	✓	1	1	e

11101111The first two baselines, PRIMORL NO CLIP and PRIMORL NO PRIVACY are not private ( $\epsilon < \infty$ )1112but allow us to isolate the impact of trajectory-level model ensemble training (without clipping and1113noise addition) and clipping on policy performance. We do not report results for PRIMORL NO1114CLIP for CARTPOLE and PENDULUM as we found that the model optimized with TDP MODEL1115ENSEMBLE TRAINING diverges without clipping.

#### F COMPARISON TO EXISTING METHODS

The closest and only comparable work in offline DPRL is Qiao & Wang (2023a). In Table 4, we compare the characteristics of PRIMORL with their algorithm DP-VAPVI (since their other algorithm is only for tabular MDPs and obviously does not compare). This comparison highlights that PRIMORL and DP-VAPVI are designed for very distinct settings. We cannot efficiently implement DP-VAPVI on our benchmark, in particular because of the continuous action spaces and the fact that it explicitly relies on a finite, relatively small horizon H (not only the number of statistics to maintain and privatize depends on H but the amount of noise needed to privatize each statistic also grows linearly with H). 

Although their scope is limited to finite, tabular and linear MDPs and their algorithms are not suited for direct comparison on the same benchmarks, we provide below a side-by-side comparison of our respective empirical results, with the aim of re-contextualizing our results within the current state of the literature.

First, we compare the complexity of the benchmark tasks considered here and the evaluation environment used in Qiao & Wang (2023a). Qiao & Wang (2023a) evaluate their algorithms on an episodic synthetic linear MDP with 2 states and 100 actions, and horizon H = 20. On the other hand, we

consider standard control tasks with multi-dimensional continuous state and action spaces. Moreover,

#### Table 4: Comparison between PRIMORL (Ours) and DP-VAPVI (Qiao & Wang)

	PRIMORL (Ours)	DP-VAPVI (Qiao & Wang)
MODEL-BASED	✓	1
SETTING	$\gamma$ -discounted infinite horizon	Finite horizon H
FUNCTION APPROXIMATION	General function approximation, including NN	Linear function approximation with known feature $\phi(s,a)$
SPACES	Continuous ${\mathcal S}$ and ${\mathcal A}$	Could theoretically handle continuous actions, b arg max <sub>π<sub>h</sub></sub> (. s) $\langle \hat{Q}_h(s, a), \pi_h(\cdot s) \rangle$ is impractic to compute for large or infinite $\mathcal{A}$
Model type	Global model ensemble (step independent), set of weights $\theta = \{\theta_i\}_i = 1^N$ with $\theta_i \in \mathbb{R}^d$	Step-dependent model represented by $5H$ statistics
PRIVACY BUDGET	Scales with training hyperparameters $q, T, N$ (indirectly)	Scales with horizon $H$ , a parameter of the problem

1147 1148

1134

1135

our tasks have long horizons and high frequency, which makes them impractical to represent in the episodic setting, justifying the use of the  $\gamma$ -discounted infinite-horizon setting.

We then compare the privacy-performance trade-offs achieved by Qiao & Wang (2023a) and PRI-1151 MORL. In Qiao & Wang (2023a), they do not mention explicitly the privacy budgets  $\epsilon$ , but instead 1152 mention the zero-concentrated differential privacy (z-CDP) parameter  $\rho$ . For clarity and fair com-1153 parison, we convert the z-CDP guarantee into a DP guarantee. For this, we use Proposition 1.3 1154 from Bun & Steinke (2016): if a mechanism is  $\rho$ -zero-concentrated DP, then for any  $\delta > 0$  it is 1155  $(\epsilon, \delta)$ -DP, with  $\epsilon = \rho + 2\sqrt{\rho \log(1/\delta)}$ . As they evaluate their algorithms for a dataset size up to 1156 1000, we consider two values of  $\delta \in \{1/100, 1/1000\}$ . Table 5 shows the results for the various 1157 parameters  $\rho$  mentioned in Figure 1 from Qiao & Wang (2023a). We observe Qiao & Wang (2023a) 1158 also considers the low privacy regime with  $\rho = 25$  yielding  $\epsilon$  close to 50, which is comparable to our 1159 low privacy variant. They indeed consider  $\epsilon$  close to 1 with  $\rho = 0.1$ , but the cost is a 2 to 3 times 1160 worse utility. Other configurations proposed are closed in privacy budgets to what we consider in our 1161 paper. Overall, our work achieves comparable privacy-utility trade-offs than Qiao & Wang (2023a), 1162 but on significantly more complex tasks.

- 1163 1164
- 1165 1166

Table 5: Results from Qiao & Wang (2023a), converted from z-CDP

Z-CDP GUARANTEE $\rho$	DP $\epsilon$ for $\delta=10^{-1}$	DP $\epsilon$ for $\delta=10^{-3}$
25	40.2	51.3
5	11.8	16.8
1	4.0	6.26
0.1	1.1	1.8

1171 1172 1173

1174

#### G DISCUSSION ON THE $\epsilon$ parameter

1175 1176 As the privacy budgets  $\epsilon$ 's presented in our experimental results do not provide strong theoretical DP guarantees, we would like to further discuss the implications of such privacy budgets in practice.

1178 First, we point out that such  $\epsilon$  values are comparable to existing work. In particular, as pointed out in 1179 Section F, Qiao & Wang (2023a) achieves similar privacy-performance trade-offs and also consider 1180 the "low privacy regime" with  $\epsilon$ 's approaching 50 for their best-performing variant. We argue that 1181 studying different privacy regimes allows us to clearly highlight the trade-offs between privacy and 1182 performance.

1183 Moreover, in light of recent literature on achieving differential privacy in practical deep learning 1184 (Carlini et al., 2019; Ponomareva et al., 2022; 2023), we argue that these  $\epsilon$  values may offer an 1185 adequate level of privacy in real-world applications. Ponomareva et al. (2023) states  $\epsilon \leq 10$  as a 1186 realistic and widely used goal in DP deep learning and a "sweet spot" where it is possible to preserve 1187 acceptable utility for complex ML models. Moreover, these studies point out the overly restrictive 1187 assumptions on the adversary side, which may yield unnecessarily pessimistic privacy bounds. In



For all tasks, the model is approximated with a deep neural network with SWISH activation functions and decaying weights. Models take as input a concatenation of the current state s and the taken action a and predict the difference between the next state s' and the current state s along with the reward r. Table 7 provides further implementation details.

The code repository for PRIMORL is provided as part of the supplementary material and will be made public upon acceptance. For MOPO, we use the official implementation from https: //github.com/tianheyu927/mopo, as well as the PyTorch re-implementation from https: //github.com/junming-yang/mopo. Our implementation of PRIMORL, which mainly uses PyTorch, is also based on these codebases. To collect the datasets, we use DDPG implementation from https://github.com/schatty/DDPG-pytorch.

1241 Model training with TDP MODEL ENSEMBLE TRAINING is parallelized over 16 CPUs using JobLib, while SAC training is conducted over a single Nvidia Tesla P100 GPU.

1242	Table 7:	Table 7: Implementation details						
1243		r · · · · ·						
1244								
1245		CARTPOLE	Pendulum	HalfCheetah				
1246	MODEL INPUT DIMENSION	6	4	23				
1247	MODEL OUTPUT DIMENSION	6	4	18				
1248	MODEL HIDDEN LAYERS	2	2	4				
1249	NEURONS PER LAYER	128	64	200				
1250	WEIGHT DECAY	<b>√</b>	<b>√</b>	1				
1251	ACTIVATION FUNCTIONS ENSEMBLE SIZE $N$	SWISH 5	SWISH 3	SWISH 7				
1050								

#### Table 8: Training and Hyperparameters details

	CARTPOLE	Pendulum	HALFCHEETAH
Test set size	$1\% \times K$	$1\% \times K$	$10\% \times K$
EARLY STOPPING	✓ patience = $10$	✓ PATIENCE = $10$	✓ PATIENCE = $5$
SAMPLING RATIO $q$	$10^{-3}$	$10^{-3}$	$10^{-2}$
Model local epochs $E$	1	1	1
Model batch size $B$	16	16	16
Model LR $\eta$	$10^{-3}$	$10^{-3}$	$10^{-3}$
CLIPPING STRATEGY	FLAT	PER-LAYER	PER-LAYER
SAC LR	$3.10^{-4}$	$3.10^{-4}$	$3.10^{-4}$
Rollout length $H$	20	30	5
Reward penalty $\lambda$	2.0	2.0	1.0
Auto- $\alpha$	✓	✓	1
Target entropy $H$	-3	-3	-3
UNCERTAINTY ESTIMATE	$u_{\mathrm{MPD}}$ (Bal.), $u_{\mathrm{MA}}$ (Swi.)	$u_{ m MPD}$	$u_{ m MA}$

1270 1271

1253

1255 1256 1257

1259

#### 1272

1273 H.3 1274

1275 Before model training, we split the offline dataset into two parts: a train set used to train the model, and a test set used to track model performance. We consider the test set public so that this operation does 1276 not involve additional privacy leakage. The split is made by episode (instead of by transitions), so that 1277 the test set contains 1% of the episodes for CARTPOLE and PENDULUM and 20% for HALFCHEETAH. 1278 To tune the clipping norm, we set z = 0 and progressively decreased C until it started to adversely 1279 affect performance provided the best results. Moreover, we set the sampling ratio so that a few dozen 1280 episodes are randomly selected at each step, which proved to work best in our experiments, which 1281 correspond to  $q = 10^{-3}$  for CARTPOLE and PENDULUM. The model is trained until convergence 1282 using *early stopping*. Test set prediction error is used to track model improvement. For SAC training, 1283 the real-to-model ratio  $r_{real}$  is zero, meaning that SAC is trained using only simulated data from the 1284 model, and does not access any data from the offline dataset. Training details are provided in Table 8. 1285

1286

# 1287 H.4 HYPERPARAMETERS

TRAINING DETAILS

The model is trained using TDP MODEL ENSEMBLE TRAINING with learning rate  $\eta = 10^{-3}$ , batch size B = 16, and number of local epochs E = 1.

The policy is optimized within the model using Soft Actor-Critic with rollout, with rollout length and penalty depending on the task. We use a learning rate of  $3.10^{-4}$  for both the actor and the critic. For entropy regularization, we use auto- $\alpha$  with target entropy H = -3.

1295 Hyperparameters are summarized in Table 8. We do not report the privacy loss resulting from hyperparameter tuning, although we recognize its importance in real-world applications.

#### 1296 H.5 PRIVACY PARAMETERS 1297

1298 In Table 1, we provide the privacy budgets  $\epsilon$  computed with the moments accountant method from 1299 Abadi et al. (2016). We use the DP accounting tools from Google's Differential Privacy library, available on GitHub. Privacy budget are computed for  $\delta = 10^{-5}$ , *i.e.* less than  $K^{-1}$  as recommended 1300 in the literature. It also depends on the noise multiplier z, the number of training round T and the 1301 sampling ratio q. Since we use early stopping and the different training runs have different durations, 1302 we use the average number of training rounds in the privacy budget computations. 1303

1304 For PENDULUM, we use z = 0.35 and z = 0.52 for PRIMORL LOW and PRIMORL HIGH, 1305 respectively. For CARTPOLE-BALANCE, we use z = 0.25 and z = 0.45 for PRIMORL LOW and PRIMORL HIGH, respectively. For CARTPOLE-SWINGUP, we use z = 0.25 and z = 0.38for PRIMORL LOW and PRIMORL HIGH, respectively. The value for PRIMORL HIGH is chosen by incrementally increasing z until policy performance drops below acceptable levels. The corresponding  $\epsilon$  is therefore roughly the best privacy budget we can obtain while keeping acceptable 1309 policy performance. The value for PRIMORL LOW is chosen arbitrarily to provide a weaker level of 1310 privacy that typically yields higher policy performance, illustrating the trade-off between the strength 1311 of the privacy guarantee and the performance. Table 9 summarizes the parameters used to compute 1312  $\epsilon^{\text{MA}}(z,q,T,\delta)$  in our experiments. 1313

Table 9: Training and privacy parameters used to compute  $\epsilon^{MA}(z, q, T, \delta)$ .

		z	T	q	δ	$\epsilon$
Pendulum	Low	0.35	$7.10^{3}$	$10^{-3}$	$10^{-5}$	22.3
	High	0.52	$7.10^{3}$	$10^{-3}$	$10^{-5}$	<b>5.1</b>
BALANCE	Low	0.25	$7.10^{3}$	$10^{-3}$	$10^{-5}$	85.0
	High	0.45	$7.10^{3}$	$10^{-3}$	$10^{-5}$	8.2
SWINGUP	Low	0.25	$10.10^{3}$	$10^{-3}$	$10^{-5}$	94.2
	High	0.38	$10.10^{3}$	$10^{-3}$	$10^{-5}$	17.0

#### 1327 H.6 COMPUTING $\epsilon$ : THE MOMENTS ACCOUNTANT 1328

Theorem 1 from McMahan et al. (2018) allows us to compute the privacy guarantees 1329  $(\epsilon^{MA}(z,q,T,\delta),\delta)$  of Algorithm 1 using the Moments Accountant from Abadi et al. (2016). To com-1330 pute  $\epsilon^{MA}(z, q, T, \delta)$  in our experiments, we use the DP accounting tools from Google's Differential 1331 Privacy library, which provides an improved version of the moments accountant based on Rényi 1332 Differential Privacy (RDP) Mironov (2017b). Since the computations of the RDP accountant are 1333 quite involved while the underlying principles are the same, we rather present the original moments 1334 accounting method based on Section 3.2 from Abadi et al. (2016). 1335

1336 By taking into account the DP noise distribution, the moments accountant allows to get a tighter 1337 bound on the total privacy leakage compared to the standard strong composition theorem. Using an  $(\epsilon, \delta)$ -DP mechanism at each gradient step, Algorithm 1 with T training steps and a sampling ratio q 1338 is  $(\mathcal{O}(q\epsilon\sqrt{T}), \delta)$ -DP by the moments accountant. For comparison, the strong composition theorem 1339 1340 would yield  $\left(\mathcal{O}(q\epsilon_{\lambda}/T\log(1/\delta)), Tq\delta\right)$ . 12/1

ts of the privacy loss random variable. We denote  $\mathcal{M}_{\sigma^2}$  the Gaussian mechanism at each training step t, which is characterized by the magnitude  $\sigma^2 := \sigma^2(z, q, T, \delta)$  of the Gaussian noise. The privacy loss for  $\mathcal{M}_{\sigma^2}$  at output o is 1345 defined as follows:  $c(o; \sigma^2, D, D') = \log \frac{\mathbb{P}(\mathcal{M}_{\sigma^2}(D) = o)}{\mathbb{P}(\mathcal{M}_{\sigma^2}(D') = o)} ,$ 1946

1326

1348

where D, D' are neighboring datasets. It quantifies the privacy leakage for the specific output o1349 taking into account the randomness of the algorithm. The  $\lambda$ -th moment  $\alpha(\lambda; D, D')$  is defined as the logarithm of the moment generating function:

$$\alpha_{\mathcal{M}_{\sigma^2}}(\lambda) = \max_{D,D'} \log \mathbb{E}_{o \sim \mathcal{M}_{\sigma^2}(D)} \left[ \exp(\lambda c(o; \mathcal{M}_{\sigma^2}, D, D')) \right]$$

To bound  $\alpha_{\mathcal{M}_{\sigma^2}}(\lambda)$  for a Gaussian mechanism of scale  $\sigma^2$ , Abadi et al. (2016) show that, denoting  $\mu_x$  the p.d.f. of  $\mathcal{N}(x, \sigma^2)$  and  $\mu = (1 - q)\mu_0 + q\mu_1$ , it suffices to estimate  $\alpha(\lambda) = \log \max(E_1, E_2)$ with:

 $E_1 = \mathbb{E}_{z \sim \mu_0} \left[ (\mu_0(z)/\mu(z))^{\lambda} \right]$ 

 $E_2 = \mathbb{E}_{z \sim \mu} \left[ (\mu(z)/\mu_0(z))^{\lambda} \right] .$ 

1360 Implementations of the moments accountant typically use numerical integration to estimate  $\alpha(\lambda)$ .

To compute  $\epsilon^{MA}(z, q, T, \delta)$ , a bound on the total privacy loss of Algorithm 1, it then suffices to compute a bound on  $\alpha_{\mathcal{M}_{\sigma^2}}(\lambda)$  at each step and sum over all steps. Since we cannot compute a bound for all  $\lambda$ , we need to specify as input a discrete list  $\Lambda = \{\lambda_1, ..., \lambda_S\}$  of moments to bound, and select the  $\lambda$  yielding the best privacy budget. Abadi et al. (2016) find that it usually suffices to compute  $\alpha(\lambda)$  for  $\lambda \leq 32$  (see Section 4).

1368 H.7 COMPUTATIONAL RESOURCES

We perform training on a single machine with 64 CPUs and 6 Tesla P100 GPUs with 16GB RAM each. The full training of a single policy, from model learning to policy optimization, takes several hours.



Figure 5: Learning curves for the SAC policy on HALFCHEETAH (*right*). Policy performance (episodic return) is evaluated in the true MDP at the end of each training epoch, over 10 evaluation episodes with different random seeds.

1423 I ADDITIONAL EXPERIMENTS

1422

1424

1444

1425Table 10: Ablation study investigating the impact of using different clipping methods (*flat clipping*1426(FC) and *per-layer clipping* (LC)) and different uncertainty estimates ( $u_{MA}$  and  $u_{MPD}$ ). For instance,1427MA + LC means the model has been train with *per-layer clipping* and the policy optimized under1428uncertainty estimate  $u_{MA}$ . We report policy performance as the mean episodic return over 101429evaluation episodes, averaged over the last 10 epochs of policy optimization. Average performance1430and 95% confidence intervals are computed on at least 3 seeds.

	Z	MA + LC	MPD + LC	MA + FC	MPD + FC
PENDULUM	0.35 0.52	$\begin{array}{c} 734.5 \pm 28.9 \\ 750.1 \pm 75.97 \end{array}$	$\begin{array}{c} \mathbf{817.4 \pm 21.7} \\ \mathbf{778.9 \pm 53.5} \end{array}$	$\begin{array}{c} 638.9 \pm 91.7 \\ 612.4 \pm 76.8 \end{array}$	$\begin{array}{c} 757.9 \pm 76.5 \\ 723.0 \pm 47.0 \end{array}$
BALANCE	0.25 0.45	$\begin{array}{c} 785.4 \pm 77.9 \\ 749.0 \pm 115.6 \end{array}$	$792 \pm 85.8 \\ 738.5 \pm 85.7$	$\begin{array}{c} {\bf 819.9 \pm 66.1} \\ {\bf 722.7 \pm 105.7} \end{array}$	$815.8 \pm 97.2$ <b>758.2</b> $\pm$ <b>187.2</b>
SWINGUP	0.25 0.38	$\begin{array}{c} 711.4 \pm 90.6 \\ 536.1 \pm 112.8 \end{array}$	$704.2 \pm 46.2 \\ 575.4 \pm 89.7$	$\begin{array}{c} 772.4 \pm 73.9 \\ \textbf{698.3} \pm \textbf{57.5} \end{array}$	$\begin{array}{c} {\bf 777.1 \pm 35.0} \\ {\bf 590.8 \pm 34.8} \end{array}$

In Table 10, we compare performance depending on the clipping method used to train the model (*flat* or *per-layer clipping*) and the uncertainty estimator used to optimize the policy ( $u_{MPD}$  or  $u_{MA}$ ). We can see that no clipping method or uncertainty estimator is significantly superior overall, but these choices may impact privacy performance for a specific task. Preliminary results on this ablation study led us to choose, for each task, the clipping strategy and the uncertainty estimate stated in Table 8.

# 1445 J EXPERIMENTS ON HALFCHEETAH

1447 We conduct experiments on the MEDIUM-EXPERT dataset (K = 2,003) from the classic D4RL 1448 benchmark (Fu et al., 2020). Experimental results are reported in Figure 5 and Table 11 (in appendix), 1449 using C = 15.0 and  $q = 10^{-2}$ .

1450 If PRIMORL can train competitive policies with small enough noise levels - a tiny amount of 1451 noise like  $z = 10^{-4}$  proving even beneficial, possibly acting as a kind of regularization —, we 1452 were not able to obtain reasonable  $\epsilon$ 's. Indeed, a noise multiplier as small as  $z = 10^{-3}$  is enough 1453 to cause a significant decline in performance. HALFCHEETAH thus appears a significantly harder 1454 tasks than CARTPOLE and PENDULUM. It is not surprising as HALFCHEETAH is higher-dimensional, and the theoretical analysis led in Section 4.3.1 showed that the dimension d of the problem could 1455 negatively impact the performance of the policy. However, we point out that the size of the dataset 1456 for HALFCHEETAH is very limited, and argue that larger datasets with substantially more episodes 1457 would translate into competitive privacy-performance trade-offs, as we develop in Section L.

1458 Table 11: Results for HALFCHEETAH MEDIUM-EXPERT. RETURN is the return of the SAC policy 1459 evaluated over 10 episodes at the end of each training epoch, averaged over the last 20 epochs.

Method	z	Return
МОРО	0.0	$10931\pm1326$
PRIMORL NO CLIP PRIMORL NO NOISE	$\begin{array}{c} 0.0\\ 0.0\end{array}$	$\begin{array}{c} 7062 \pm 2230 \\ 8792 \pm 2053 \end{array}$
PriMORL	$z = 1.10^{-4}$ z = 1.10^{-3}	$9729 \pm 2018 \\ 3697 \pm 1465$

1468 1469 1470

1471

#### Κ ALGORITHMS

1472 Algorithm 2 is the fully detailed pseudo-code for PRIMORL. Algorithm 3 details the clipping 1473 method used in TDP MODEL ENSEMBLE TRAINING. Algorithm 4 is the pseudo-code for SAC 1474 policy optimization on the pessimistic private model. This pseudo-code is based on https:// 1475 spinningup.openai.com/en/latest/algorithms/sac.html

1476 1477 Algorithm 2 Model Training with TDP MODEL ENSEMBLE TRAINING 1478 1: Input: offline dataset  $\mathcal{D}_K$ , sampling ratio  $q \in (0,1)$ , noise multiplier  $z \ge 0$ , clipping norm 1479 C > 0, local epochs E, batch size B, learning rate  $\eta$ 1480 2: **Output:** private model  $M_{\theta}$ 1481 3: Initialize model parameters  $\theta_0$ 1482 4: for each iteration  $t \in [0, T-1]$  do  $\mathcal{U}_t \leftarrow (\text{sample with replacement trajectories from } \mathcal{D}_K \text{ with prob. } q)$ 1483 5: 6: for each trajectory  $\tau_k \in \mathcal{U}_t$  do 1484 Clone current models  $\{\theta_i^{\text{start}}\}_{i=1}^N \leftarrow \{\theta_i(t)\}_{i=1}^N$ 7: 1485  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{\mathrm{start}} := \left(\boldsymbol{\theta}^{\mathrm{start}}\right)_{i=1}^{N}$ 1486 8: for each local epoch  $i \in [1, E]$  do 1487 9:  $\mathcal{B} \leftarrow (\tau_k$ 's data split into size *B* batches) 10: 1488  $\left. \begin{array}{l} \left. \begin{array}{l} \text{end} \left\{ \boldsymbol{\theta}_{k} \in \mathcal{B} \text{ do} \\ \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \boldsymbol{\eta} \nabla \mathcal{L}(\boldsymbol{\theta}; b) \\ \boldsymbol{\theta} \leftarrow \boldsymbol{\theta}^{\text{start}} + \text{ENSEMBLECLIP}(\boldsymbol{\theta} - \boldsymbol{\theta}^{\text{start}}; C) \end{array} \right\} \text{ENSCLIPGD} \left( \tau_{k}, \left\{ \boldsymbol{\theta}_{i}^{\text{start}} \right\}_{i=1}^{N}; C, E, B \right) \\ \text{ad for} \end{array} \right.$ 11: for each batch  $b \in \mathcal{B}$  do 1489 12: 1490 13: 1491 end for 14: 1492 end for 15:  $\Delta_{t,k}^{\text{clipped}} \leftarrow \theta - \theta^{\text{start}}$ 1493 16: 1494 end for 17: 1495  $\Delta^{\mathrm{avg}}_i(t) = \tfrac{\sum_{k \in \mathcal{U}_t} \Delta^{\mathrm{clipped}}_{i,k}(t)}{qK}$ 18: 1496  $\theta(t+1) \leftarrow \theta(t) + \Delta^{\operatorname{avg}}(t) + \mathcal{N}\left(0_{N_d}, \left(\frac{zC}{qK}\right)^2 I_{N_d}\right)$ 1497 19: 1498 20: end for 1499 1500 1501 1502 Algorithm 3 Ensemble Clipping (ENSEMBLECLIP) 1503 1: Input: ensemble size N, number of model layers L, unclipped gradient  $\Delta = \{\Delta_{i,\ell}\}_{i,\ell=1}^{N,L}$ 1504 clipping norm C2: **Output:** clipped gradient  $\Delta^{\text{clipped}}$ 3:  $\Delta_i \leftarrow (\Delta_{i,\ell})_{\ell=1}^L$ ,  $C_i = \frac{C}{\sqrt{N}}$ 1507

```
4:
                                                                                      \Delta_i^{\text{clipped}} \leftarrow \frac{\Delta_i}{\max\left(1, \frac{\|\Delta_i\|_2}{C_i}\right)}, \ j = 1, ..., m.
1509
1510
1511
```

Alg	orithm 4 Private Model-Based Optimization with SAC
1:	<b>Input:</b> private model $\hat{M} = (\hat{P}, \hat{r})$ , empty replay buffer $\mathcal{B}$ , uncertainty estimator $u \in$
2	$\{u_{\text{MA}}, u_{\text{MPD}}\}$
2:	<b>Output:</b> private policy $\pi^{22}$
5. ⊿∙	for enoch $e \in [1, E]$ do
+. 5∙	while episode is not terminated do
<i>6</i> :	Observe state s and select action $a \sim \pi_{\mathcal{E}}(\cdot s)$
7.	Execute a in the pessimistic MDP $\tilde{M}$ and observe next state $s' \sim \hat{P}( s a)$ reward
<i>.</i>	$r \sim \hat{r}(s, a) - \lambda u(s, a)$ and done signal d
8:	Store $(s, a, r, s', d)$ in replay buffer $\mathcal{B}$
9:	if time to update then
10:	Sample a batch of transitions $B = \{(s, a, r, s', d)\}$ from buffer $\mathcal{B}$
11:	Compute targets for Q-functions:
	$y(r, s', d) = r + \gamma(1 - d) \left( \min_{i=1,2} Q_{\omega_{\text{targ},i}}(s', \tilde{a}') - \alpha \log \pi_{\xi}(\tilde{a}' s') \right),  \tilde{a}' \sim \pi_{\xi}(\cdot s')  .$
12:	Update O-functions by one step of gradient descent using:
	$ abla_{\omega_i} rac{1}{ B } \sum_{(s,a,r,s',d) \in B} \left( Q_{\omega_i}(s,a) - y(r,s',d) \right)^2, \ \  ext{for} \ i = 1,2.$
13:	Update policy by one step of gradient ascent using:
	$\nabla = \frac{1}{2} \sum \left( \min \left( 0, \left( \frac{\pi}{2}, \left( \frac{\pi}{2} \right) \right) - 1 + \frac{\pi}{2} - \left( \frac{\pi}{2}, \left( \frac{\pi}{2} \right) \right) \right) - \frac{\pi}{2} - \left( \frac{\pi}{2} \right) - \frac{\pi}{2} + \frac{\pi}$
	$ v_{\xi} \frac{ B }{ B } \sum_{z \in B} \left( \min_{i=1,2} Q_{\omega_i}(s, a_{\zeta}(s)) - \alpha \log \pi_{\xi}(a_{\zeta}(s) s) \right),  a_{\zeta}(s) \sim \pi_{\xi}(\cdot s). $
	$s \in B$
14:	Update target networks with:
	$\omega_{\text{targ},i} \leftarrow \rho \omega_{\text{targ},i} + (1-\rho) \omega_i$ , for $i = 1, 2$ .
	······································
	ena II and while
15:	chu white
15: 16: 17:	Evaluate $\pi_{\star}$ is the true environment M
15: 16: 17: 18:	Evaluate $\pi_{\xi}$ is the true environment $\mathcal{M}$ .

# <sup>1566</sup> L THE PRICE OF PRIVACY IN OFFLINE RL

In this section, we provide theoretical and practical arguments to further justify the need for (much)
 larger datasets in order to achieve competitive privacy trade-offs in offline RL, as pointed out in
 (Section 5).

1571 Why does privacy benefit so much from large datasets? From a theoretical perspective, it stems 1572 from two facts: 1)  $\epsilon$  scales with the sampling ratio q (privacy amplification by subsampling), and 1573 2) noise magnitude  $\sigma$  is inversely proportional to  $\mathbb{E}[|\mathcal{U}_t|] = qK$ . Clearly, the privacy-performance 1574 trade-off would benefit from both small q (reducing  $\epsilon$ ) and large qK (reducing noise levels and thus 1575 improving performance), which are conflicting objectives for a fixed K. However, if we consider 1576 using larger datasets of size  $K' \gg K$ , it becomes possible to find a K' large enough so that we can use  $q' \ll q$  and  $q'K' \gg qK$ , achieving both much stronger privacy and better performance. We can 1577 1578 even argue that for a given privacy budget  $\epsilon$  (obtained for a given q) and an unlimited capacity to increase K, we could virtually tend to zero noise levels and achieve optimal performance. Therefore, 1579 PRIMORL, already capable of producing good policies with significant noise levels and  $\epsilon$ , has the 1580 potential to achieve stronger privacy guarantees provided access to large enough datasets. 1581

1582 An aspect that deserves further development is the iterative aspect of the used training methods and its effect on privacy. Differential privacy being a worst-case definition, it assumes that all intermediate models are released during training. Although the practicality of this hypothesis is debatable, it 1585 definitely impacts privacy: privacy loss is incurred at each training iteration (corresponding to a gradient step on the global model in DP-SGD and TDP MODEL ENSEMBLE TRAINING) and privacy 1586 budget, therefore, scales with the number of iterations T. Consequently, limiting the number of 1587 iterations is even more crucial with DP training than with non-private training. Training a model on the 1588 kind of tasks we considered nonetheless requires a lot of iterations to reach convergence (empirically, 1589 thousands of iterations for CARTPOLE and tens of thousands of iterations for HALFCHEETAH), and 1590 the privacy budget suffers unavoidably. 1591

However, one way to circumvent this is to leverage privacy amplification by subsampling. Indeed, as 1592 McMahan et al. (2017) observe, the additional privacy loss incurred by additional training iterations 1593 becomes negligible when the sampling ratio q is small enough, which is a direct effect of privacy 1594 amplification by subsampling. We discussed above how increasing dataset size K allowed to decrease 1595 both sampling ratio q and noise levels. Therefore, by increasing the size of the dataset, we also greatly 1596 reduce the impact of the number of training iterations, likely promoting model convergence. This 1597 further reinforces the need for large datasets in offline RL in order to study privacy. As an example, 1598 McMahan et al. (2018) consider datasets with  $10^6$  to  $10^9$  users to train DP recurrent language models, and this is arguably the main reason why they achieve formal strong privacy guarantees. For comparison, the classical RL UNPLUGGED and D4RL benchmarks provide datasets with  $K \approx 10^{1}$ to  $K \approx 10^3$  datasets. Achieving the privacy-performance trade-offs demonstrated in Section 5 would not have been possible without the collection of large datasets. Moreover, datasets orders of magnitude larger would be required to attain formal, strong privacy guarantees, such as  $\epsilon < 1$ . 1603 While conducting experiments in deep offline RL with such extensive datasets demands substantial 1604 computational resources, we argue that scenarios involving access to datasets with a vast number of trajectories are reflective of real-world situations. For this reason, we consider this case worthy of 1606 thorough investigation.

Figure 6 illustrates this point in another way. Given  $\epsilon \in \{10^{-4}, 10^{-3}, 10^{-2}\}$ , we plot for a range of sampling ratio q the maximum number of iterations T that is allowed so that the total privacy loss does not exceed  $\epsilon$ , as a function of the noise multiplier z. We can see how decreasing q makes it well easier to train a private model: dividing q by 10, we "gain" roughly 10 times more iterations across all noise levels.

1613

1614

1615

1616

1617

1618



Figure 6: Maximum number of iterations T so that the privacy loss does not exceed  $\epsilon$ , as function of the noise multiplier z.

#### M BROADER IMPACTS

As recent advances in the field have moved reinforcement learning closer to widespread real-world application, from healthcare to autonomous driving, and as many works have shown that it is no more immune to privacy attacks than any other area in machine learning, it has become crucial to design algorithmic techniques that protect user privacy. In this paper, we contribute to this endeavor by introducing a new approach to privacy in offline RL, tackling more complex control problems and thus paving the way towards real-world private reinforcement learning. We firmly believe in the importance of pushing the boundaries of this research field and are hopeful that this work will contribute to practical advancements in achieving trustworthy machine learning.