ONLINE DIFFERENTIAL PRIVACY BAYESIAN OPTI-MIZATION WITH SLICED WASSERSTEIN COMPRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

The increasing prevalence of streaming data and rising privacy concerns pose significant challenges for traditional Bayesian optimization (BO), which is often ill-suited for real-time, privacy-aware learning. In this paper, we propose a novel online locally differentially private BO framework that enables zero-order optimization with rigorous privacy guarantees in dynamic environments. Specifically, we develop a one-pass Gaussian process compression algorithm based on the sliced Wasserstein distance, which effectively addresses the challenges of kernel matrix scalability, memory efficiency, and numerical stability under streaming updates. We further establish a systematic non-asymptotic convergence analysis to characterize the privacy–utility trade-off of the proposed estimators. Extensive experiments on both simulated and real-world datasets demonstrate that our method consistently delivers accurate, stable, and privacy-preserving results without sacrificing efficiency.

1 Introduction

Bayesian optimization (BO) (Močkus, 1974; Jones et al., 1998) is a sample-efficient framework widely used for the global optimization of expensive, non-convex, or black-box functions, with applications in hyperparameter tuning, robotics, and scientific discovery (Snoek et al., 2012; Berkenkamp et al., 2023). In particular, BO iteratively selects query points using a probabilistic surrogate model and balances exploration and exploitation through the predictive mean and uncertainty, often achieving high-performance solutions with relatively few evaluations. To date, BO has been extensively studied, leading to numerous methodological advances, including local descent strategies (Müller et al., 2021; Nguyen et al., 2022), mixed-space optimization techniques (Neiswanger et al., 2022), scalable acquisition via Monte Carlo methods (Balandat et al., 2020), and extensions to iterative and bilevel problems (Fu et al., 2024), supported by theoretical analyses of high-dimensional Gaussian processes (Hvarfner et al., 2024). Furthermore, practical robustness has been enhanced through improved constraint handling (Nguyen et al., 2024), contextual uncertainty modeling (Tay et al., 2024), and meta-learning strategies for rapid adaptation (Ravi & Beatson, 2019).

Building on this line of work, several methods have sought to accelerate convergence by incorporating gradient information via finite differences or kernel-based estimation (Wu et al., 2017; Eriksson et al., 2019). For example, Müller et al. (2021) reformulated BO as an approximate gradient descent procedure, a formulation later extended by the gradient information BO framework (Wu et al., 2023), which reduces gradient uncertainty and guarantees convergence to low-gradient regions in reproducing kernel Hilbert spaces (RKHS). More recently, Sopa et al. (2025) adapted these methods to tackle high-dimensional problems. Nonetheless, the aforementioned BO methods remain predicated on static datasets and are not designed for streaming environments, thereby limiting their applicability in dynamic and continually evolving settings, such as autonomous systems or large-scale monitoring, where data are generated at rates that make batch learning and reprocessing infeasible.

The growing demand for real-time decision-making in streaming data environments has elevated online learning to a central paradigm, with stochastic gradient descent (SGD) serving as its primary optimization tool (Robbins & Monro, 1951; Bottou, 2010). Recent advances have extended SGD beyond classical settings to a variety of estimation settings, including online learning (Su & Zhu, 2023; Xie et al., 2025), contextual bandits (Ding et al., 2021), and high dimensional infer-

ence tasks (Han et al., 2024). Yet these methods remain rooted in the frequentist paradigm and rely heavily on heuristic exploration, and depend on gradient access, which constrains data efficiency and often results in slow convergence in complex, non-convex functions (Ruder, 2016). By contrast, BO does not require gradient information and provides a principled framework for balancing exploration and exploitation, thereby enabling more sample-efficient optimization in such settings (Jones et al., 1998). From a Bayesian standpoint, online learning has largely been investigated in sequential decision-making contexts, such as hyperparameter tuning (Snoek et al., 2012), black-box optimization (Frazier, 2018), and sequential hypothesis testing (She et al., 2021), but these methods typically emphasize decision efficiency over functional exploration and often lack expressive input—output modeling beyond classification. Consequently, they are ill-suited for streaming environments, where adaptive and sample-efficient exploration of the response surface is essential, highlighting the need for a scalable BO framework explicitly designed for online settings.

On the other hand, the increasing complexity and scale of data amplify the challenges of safeguarding individual privacy and sustaining public trust, particularly in applications that involve sensitive user information, such as financial transactions in banking or location data from mobile applications. Differential Privacy (DP) (Dwork, 2006; Dwork et al., 2014), one of the most widely adopted frameworks for privacy-preserving data analysis, provides a rigorous guarantees the output of a computation does not reveal sensitive information about any individual in the dataset. DP is typically implemented under two models: central DP (CDP), where a trusted server injects noise into aggregated data (Ponomareva et al., 2023), and local DP (LDP), where users privatize their data before sharing, thereby removing the need for a trusted server (Duchi et al., 2018; Lowy & Razaviyayn, 2023; Duchi & Ruan, 2024). Although substantial advances in both paradigms, most existing methods continue to be developed within the frequentist framework.

Recently, increasing attention has been devoted to privacy-preserving estimation in BO under the CDP framework. Early work by Heikkilä et al. (2017) proposed a distributed DP-Bayesian learning method that leverages secure multi-party aggregation and Gaussian mechanisms for efficient privacy-preserving inference. Subsequently, Dimitrakakis et al. (2017) introduced a Bayesian DP framework based on posterior sampling, establishing sensitivity bounds for arbitrary data metrics. Building on this foundation, Triastcyn & Faltings (2020) incorporated distributional information to provide more practical privacy guarantees, while Zhang & Zhang (2023) further advanced the line of research by designing an exact and efficient DP Metropolis-Hastings algorithm. In parallel, Li et al. (2023) investigated DP synthetic data generation using Bayesian networks and established statistical accuracy guarantees for marginal-based methods. Makhija et al. (2024) developed a federated Bayesian learning framework that trains personalized models across clients with rigorous DP guarantees, and Chew et al. (2025) introduced a risk-weighted pseudo-posterior distribution to address imbalanced data in DP deep learning. More recently, Sopa et al. (2025) proposed a DP gradient-informed BO method for high-dimensional problems with exponential convergence guarantees. Despite these advances, existing methods are primarily designed for batch learning and typically assume a trusted data curator. To the best of our knowledge, no scalable and statistically rigorous method has yet been developed for online BO under the LDP framework. This gap naturally motivates the following fundamental question:

Is it possible to develop an <u>online</u>, <u>gradient-free</u>, <u>Bayesian optimization framework that provides rigorous <u>LDP</u> guarantees without sacrificing statistical efficiency?</u>

The main goal of this paper is to address the question outlined above. To this end, we propose a fully online LDP framework for real-time BO. Specifically, we introduce a novel one-pass, online, gradient-free LDP-BO algorithm that integrates a Sliced Wasserstein Compression (SWC) strategy, which enables efficient kernel compression to control memory growth while simultaneously ensuring privacy-preserving learning in streaming data environments. An overview of the proposed framework is provided in Figure 1. The key contributions of this work are summarized as follows:

Online LDP Bayesian estimation framework: Our framework provides rigorous periteration LDP guarantees for BO in an online setting, thereby enabling privacy-preserving real-time estimation and addressing a key limitation of existing methods that typically require access to the entire dataset in dynamic environments. By constructing a surrogate model, we further develop a zeroth-order optimizer that eliminates the need for gradient



Figure 1: Flowchart of the proposed online privacy-preserving Bayesian framework. Data is processed sequentially, and privacy-preserving estimates are obtained using the LDP-BO algorithm. During this process, the kernel dictionary is compressed via the sliced Wasserstein distance to control memory growth.

information, making the framework well-suited for complex objective functions with non-differentiable points or discontinuities.

- Efficient compression algorithm: We propose an efficient compression algorithm based on the Sliced Wasserstein distance to manage the kernel dictionary in streaming data environments. The algorithm reduces memory overhead while preserving numerical stability, and we establish that the kernel dictionary size remains uniformly bounded, ensuring efficient BO without loss of model fidelity. Moreover, the proposed algorithm achieves $\mathcal{O}(1)$ time and space complexity per iteration. By eliminating the need to store or re-access historical data, our method avoids the $\mathcal{O}(t^3)$ computational cost and $\mathcal{O}(t)$ memory requirements inherent standard BO and inducing point-based batch methods.
- Non-asymptotic analysis: We establish non-asymptotic convergence rates for our estimator under decaying step sizes, addressing both strongly convex losses and the more general smooth (but not necessarily convex) losses. The rates depend explicitly on the sample size, privacy budget, and BO compression error. Specifically, in the strongly convex setting, the estimation error achieves the same order as that of SGD, whereas under smoothness alone we provide guarantees of convergence to stationary points. Notably, our method achieves SGD-like convergence behavior without requiring access to exact gradients at any stage of the optimization process.

2 Problem formulation

In this paper, we consider an online learning framework in which independent and identically distributed (i.i.d.) observations $\{z_i\}_{i=1}^t$ with $t \geq 1$, arrive sequentially, where each $z_i = (x_i^\top, y_i)^\top$ consists of a covariate vector $x_i \in \mathbb{R}^p$ and a response $y_i \in \mathbb{R}$, jointly drawn from an underlying distribution \mathcal{F} . Specifically, we consider the following optimization problem:

$$\boldsymbol{\theta}^{\star} = \operatorname{argmin}_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \left(f(\boldsymbol{\theta}) := E_{\boldsymbol{z} \sim \mathcal{P}_{\boldsymbol{z}}} \left[\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{z}) \right] = \int \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{z}) d\mathcal{P}_{\boldsymbol{z}}(\boldsymbol{z}) \right), \tag{1}$$

where $\mathcal{L}(\theta, z)$ denotes a pre-specified loss function with respect to θ and z is a random variable from the distribution $\mathcal{P}z$.

We aim to estimate an unknown parameter θ^* from streaming data within the BO framework, where observations are received sequentially over time. The BO framework adopts a Gaussian process (GP) as a probabilistic surrogate model. By placing a GP prior with a twice-differentiable kernel K, the objective function f can be efficiently approximated without explicit gradient computations. Given a collection of points $\mathcal{D} = \{\theta_i\}_{i=1}^t$, the posterior distribution $f \mid \mathcal{D} \sim \mathrm{GP}(m_{\mathcal{D}}, K_{\mathcal{D}})$ yields closed-form estimates, while the gradient process $\nabla f \mid \mathcal{D}$ (Müller et al., 2021)

$$\nabla f(\boldsymbol{\theta}) \mid \mathcal{D} \sim N\left(\nabla m_{\mathcal{D}}(\boldsymbol{\theta}), \nabla^2 K_{\mathcal{D}}(\boldsymbol{\theta}, \boldsymbol{\theta})\right),$$
 (2)

where

$$\nabla m_{\mathcal{D}}(\boldsymbol{\theta}) = \nabla m(\boldsymbol{\theta}) + \nabla K(\boldsymbol{\theta}, \mathcal{D}) K(\mathcal{D}, \mathcal{D})^{-1} (f(\mathcal{D}) - m(\mathcal{D})),$$

$$\nabla^{2} K_{\mathcal{D}}(\boldsymbol{\theta}, \boldsymbol{\theta}) = \nabla^{2} K(\boldsymbol{\theta}, \boldsymbol{\theta}) - \nabla K(\boldsymbol{\theta}, \mathcal{D}) K(\mathcal{D}, \mathcal{D})^{-1} \nabla K(\mathcal{D}, \boldsymbol{\theta}).$$

This procedure only depends on zeroth-order function evaluations, thereby eliminating the need for explicit gradient calculations. Since the true distribution \mathcal{P}_z is unknown, the expected risk $f(\theta)$ is intractable and is instead approximated by the empirical loss $\mathcal{L}(\theta,z)$ based on observed data. For simplicity, we assume throughout this work that the prior mean function is zero, i.e., $m(\cdot) \equiv 0$.

Unfortunately, the standard BO framework suffers from two major limitations: (1) it does not scale to online learning, as the storage requirement for \mathcal{D} grows unbounded as new data arrive sequentially, and (2) it is vulnerable to privacy breaches because repeated data queries during the optimization process may leak sensitive information, such as medical records (Liu et al., 2024) or consumer data (Hard et al., 2018). (Additional preliminaries on LDP are provided in Appendix A.1) To address these challenges, we propose GP-based BO framework to a privacy-preserving online setting that achieves computationally efficient estimation with reduced time and space complexity, while simultaneously providing rigorous individual-level privacy guarantees.

3 METHODOLOGY

In this section, we propose the online locally privacy-preserving estimation within the BO framework to the minimization problem (1).

3.1 ONLINE LOCALLY DIFFERENTIALLY PRIVATE BAYESIAN OPTIMIZATION

We first leverage BO to approximate the gradient of the underlying function defined in (1) through the gradient of a surrogate model. In particular, at each iteration, the BO procedure selects query points that minimize an acquisition function, thereby maximizing information gain in the optimization process (see Wu et al. (2023) for further details). In line with Müller et al. (2021), this paper adopts gradient information as the acquisition function, which is defined as

$$GI(\boldsymbol{\xi}; \mathcal{D}, \boldsymbol{\theta}) = Tr(\nabla^2 K_{\mathcal{D} \cup \boldsymbol{\xi}}(\boldsymbol{\theta}, \boldsymbol{\theta})), \tag{3}$$

where ξ denotes a candidate point in the parameter space Θ . This strategy minimizes the trace of the Hessian of the kernel, thereby reducing the uncertainty of gradient estimates. Furthermore, since the kernel K is smooth and Θ is compact, the acquisition function $\mathrm{GI}(\xi;\mathcal{D},\theta)$ is uniformly bounded above by a constant L (Wu et al., 2023).

At each iteration, the candidate point ξ is obtained by optimizing $\mathrm{GI}(\xi;\mathcal{D},\theta)$ and subsequently incorporated into the kernel dictionary \mathcal{D} . In streaming settings with infinitely arriving data, however, the kernel dictionary would grow unbounded as iterations proceed, which fundamentally limits the applicability of BO in online learning. To overcome this issue, we propose a compression algorithm, i.e., SWC, based on the sliced Wasserstein distance to efficiently compress \mathcal{D} (see Section 3.2 for details). This algorithm guarantees that the size of the kernel dictionary remains bounded independently of t, while ensuring that the compressed probability distribution converges to the domain of the true probability distribution.

Using the BO surrogate model, we then obtain the approximate gradient at iteration t as

$$\widehat{\nabla \mathcal{L}}_t = \boldsymbol{\mu}_{\mathcal{D}_{t-1}} = \nabla K(\hat{\boldsymbol{\theta}}_{t-1}, \mathcal{D}_{t-1}) K(\mathcal{D}_{t-1}, \mathcal{D}_{t-1})^{-1} \mathcal{L}(\hat{\boldsymbol{\theta}}_{t-1}, \boldsymbol{z}_t). \tag{4}$$

This formulation enables iterative updates without requiring storage of historical raw data or direct access to the gradient of the objective function. Upon receiving the t-th sample $z_t = (\boldsymbol{x}_t^\top, y_t)^\top$, the parameter estimate is updated via

$$\hat{\boldsymbol{\theta}}_t = \hat{\boldsymbol{\theta}}_{t-1} - \eta_t \widehat{\nabla \mathcal{L}}_t,$$

where η_t denotes the step size at iteration t. Throughout the procedure, only the estimator $\hat{\theta}_{t-1}$ and the kernel dictionary \mathcal{D}_{t-1} are required, thereby ensuring greater flexibility and substantially reduced memory usage.

However, while the above procedure enables efficient online estimation, it does not inherently safe-guard sensitive information. In streaming environments, where each newly arriving observation may expose individual data, privacy protection is indispensable. Unlike traditional centralized approaches to DP (Sopa et al., 2025), which inject noise into the entire algorithm in a post-hoc manner,

our framework embeds privacy protection directly into each iteration. This design eliminates the reliance on a trusted data curator and achieves LDP by ensuring that data are privatized at the source before any aggregation occurs. To enforce rigorous LDP guarantees, we first clip the approximate gradient to a fixed bound B>0, i.e.,

$$g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) = \boldsymbol{\mu}_{\mathcal{D}_{t-1}} \cdot \min \left\{ 1, \frac{B}{\|\boldsymbol{\mu}_{\mathcal{D}_{t-1}}\|} \right\},$$

and then perturb it with noise drawn from a suitable distribution to ensure privacy. Common choices include Gaussian, Laplace, or more sophisticated mechanisms (Dwork et al., 2014; Dong et al., 2022). In this work, we adopt the Gaussian mechanism primarily for illustrative purposes, owing to its analytical simplicity. Nevertheless, our proposed framework is general and can be easily extended to other noise distributions. Let ω_t denote Gaussian noise with mean zero and covariance matrix $2(2B/\varepsilon_t)^2\log(1.25/\delta_t)\mathbf{I}_p$, where (ε_t,δ_t) is the privacy budget allocated to the t-th iteration. The proposed private estimator is initialized at $\hat{\boldsymbol{\theta}}_0 = \tilde{\boldsymbol{\theta}}_0 = \mathbf{0}_p$ and updated as

$$\hat{\boldsymbol{\theta}}_t = \hat{\boldsymbol{\theta}}_{t-1} - \eta_t \{ g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) + w_t \}, \quad \tilde{\boldsymbol{\theta}}_t = \{ (t-1)\tilde{\boldsymbol{\theta}}_{t-1} + \hat{\boldsymbol{\theta}}_t \} / t.$$
 (5)

Notably, the optimization of the acquisition function, the SWC compression, and the posterior mean evaluation depend only on the kernel K, the compressed dictionary \mathcal{D}_{t-2} , and the previous parameter estimate $\hat{\theta}_{t-1}$, making the proposed method well-suited to streaming environments. The proposed LDP-BO procedure is summarized in Algorithm 1.

Algorithm 1 Online Locally Differentially Private Bayesian Optimization Algorithm (LDP-BO).

1: **Input**: User-defined loss function $\mathcal{L}(\cdot, \mathbf{z})$, a clipping bound B > 0, learning rates $\{\eta_t\}_{t \geq 1}$, privacy parameters $\{(\varepsilon_t, \delta_t)\}_{t \geq 1}$, and a compression budget $\kappa > 0$.

2: **Initialize**: Non-data-dependent parameters $\hat{\theta}_0 = \tilde{\theta}_0 = \mathbf{0}_p$, and evaluation set $\mathcal{D}_{-1} = \emptyset$.

3: **for** $t = 1, 2, \dots$ **do**

4: Collect a new data point $z_t = (x_t^\top, y_t)^\top$.

5: Select the candidate point $\boldsymbol{\xi} = \arg\min_{\boldsymbol{\xi}} \operatorname{GI}(\boldsymbol{\xi}; \mathcal{D}_{t-2}, \hat{\boldsymbol{\theta}}_{t-1}).$

6: Update the compressed dictionary via SWC Algorithm 2 $\mathcal{D}_{t-1} = \text{SWC}(\mathcal{D}_{t-2}, \boldsymbol{\xi})$.

7: Evaluate the loss function at $\mathcal{L}(\hat{\theta}_{t-1}, z_t)$ at point z_t .

8: Compute the posterior mean $\mu_{\mathcal{D}_{t-1}}$ by (4).

9: Clip the gradient to obtain $g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) = \boldsymbol{\mu}_{\mathcal{D}_{t-1}} \cdot \min\left\{1, \frac{B}{\|\boldsymbol{\mu}_{\mathcal{D}_{t-1}}\|}\right\}$.

10: Perform the noisy gradient descent step and update $\hat{\theta}_t$ and $\tilde{\theta}_t$ by (5).

11: **end for** 12: **Output**: $\tilde{\theta}_t$.

By the post-processing property A.4 of LDP, we establish the following privacy guarantee for Algo-

rithm 1.
Theorem

Theorem 3.1. Given an initial estimate $\hat{\theta}_0 \in \mathbb{R}^p$, consider the iterates $\{\hat{\theta}_t\}_{t\geq 1}$ defined in Algorithm 1. Then the final output $\tilde{\theta}_t$ satisfies $(\max\{\varepsilon_1,\ldots,\varepsilon_t\},\max\{\delta_1,\ldots,\delta_t\})$ -LDP.

Theorem 3.1 guarantees that each update of the proposed LDP-BO algorithm satisfies $(\max\{\varepsilon_1,\ldots,\varepsilon_t\},\max\{\delta_1,\ldots,\delta_t\})$ -LDP by introducing Gaussian noise calibrated to the sensitivity of the gradient. This mechanism safeguards the privacy of every individual sample at each iteration while eliminating the need to store raw data. The analysis for time-varying privacy parameters (ε_t,δ_t) proceeds analogously to that of the constant- (ε,δ) case. Hence, for clarity of exposition, we focus on a fixed privacy level (ε,δ) in the subsequent discussion.

3.2 SLICED WASSERSTEIN COMPRESSION

As discussed above, a major challenge in streaming data settings is the unbounded growth of the kernel dictionary as new points are continuously arrived. To address this issue, we develop an SWC strategy that controls the growth of the dictionary while preserving the statistical fidelity of the surrogate model. Specifically, in Algorithm 1, whenever a candidate point ξ is selected by (3), the posterior distribution $\rho_{\tilde{\mathcal{D}}_t}$ is updated according to (2), where $\tilde{\mathcal{D}}_t = \mathcal{D}_{t-1} \cup \xi$. To ensure

computational efficiency, the enlarged dictionary $\tilde{\mathcal{D}}_t$ is subsequently compressed using the Sliced Wasserstein (SW) distance, which quantifies discrepancies between probability distributions through their one-dimensional projections (see Bonneel et al. (2015) for details).

Our primary goal is to guarantee that the compressed dictionary \mathcal{D}_t satisfies

$$SW_2(\rho_{\mathcal{D}_t}, \rho_{\tilde{\mathcal{D}}_t}) < \kappa,$$

for a prescribed budget parameter κ , where ρ denotes the posterior density. We define the model order M_t as the column dimension of the compressed kernel dictionary \mathcal{D}_t . This compression step ensures that $M_t \leq M_{t-1} + 1$, thereby keeping the dictionary size bounded over time. The detailed SWC procedure is provided in Algorithm 2.

Algorithm 2 Sliced Wasserstein Compression (SWC).

```
1: Input: Previous dictionary \mathcal{D}_{t-1}, new acquisition point \boldsymbol{\xi} and a compression budget \kappa > 0.
 2: Initialize: \mathcal{D}_t = \mathcal{D}_{t-1} \cup \boldsymbol{\xi} and index set \mathcal{I} = \{1, \dots, M_t\}.
      while \mathcal{I} \neq \emptyset do
            for j \in \mathcal{I} do
 4:
                   Compute Sliced Wasserstein distance \eta_i = SW_2(\rho_{\mathcal{D}_{-i}}, \rho_{\tilde{\mathcal{D}}_i}).
 5:
 6:
            Identify index with minimal distance j^* = \arg\min_{i \in \mathcal{I}} \eta_i.
 7:
 8:
            if \eta_{i^*} > \kappa then break
 9:
                   \mathcal{I} = \mathcal{I} \setminus \{j^*\}, \mathcal{D}_t = \tilde{\mathcal{D}}_{\mathcal{I}}.
10:
11:
            end if
12: end while
13: Output: Compressed dictionary \mathcal{D}_t such that SW_2(\rho_{\mathcal{D}_t}, \rho_{\tilde{\mathcal{D}}_t}) \leq \kappa.
```

To ensure that the posterior distribution produced by Algorithm 2 converges to a stationary region, we impose the following assumption.

Assumption 3.2. For any c > 0, let ρ_t denote the true posterior density, and define the events: $\psi_t = \{ \operatorname{SW}_2(\rho_t, \rho_{t-1}) < c \mid \mathcal{D}_t \}$, $\tilde{\psi}_t = \{ \operatorname{SW}_2(\rho_{\mathcal{D}_t}, \rho_{\mathcal{D}_{t-1}}) < c \mid \mathcal{D}_t \}$. We assume that compression does not increase the probability of divergence relative to the original model, i.e., $\mathbb{P}\{\psi_t\} \geq \mathbb{P}\{\tilde{\psi}_t\}$.

Assumption 3.2 is mild, as the likelihood of the true posterior is at least as large as that of the sparse GP, a condition also adopted in Koppel et al. (2021). In our analysis, Assumption 3.2 serves as the Bayesian analogue of the nonexpansiveness property of projection operators. This property is essential for establishing an upper bound on the error introduced by kernel dictionary compression.

Theorem 3.3. For the compression process in Algorithm 2, the model order M_t of each posterior ρ_{D_t} is uniformly bounded as

$$M_t \leq \mathcal{O}\left(\frac{1}{\kappa}\right)^p$$
 for all t .

Theorem 3.3 establishes that, in the streaming setting, the kernel dictionary size in our BO framework remains uniformly bounded, with dependence only on the compression budget κ and the input dimension p. By operating directly on one-dimensional sample projections, the proposed method circumvents explicit density estimation and thereby mitigates sensitivity to both ambient dimensionality and discretization errors (Kolouri et al., 2015).

4 THEORETICAL PROPERTIES

In this section, we investigate the finite-sample properties of the proposed estimator. Firstly, we establish theoretical guarantees for the estimator produced by Algorithm 1 under the strongly convex loss. In order to obtain the convergence property, we also need the following assumptions.

Assumption 4.1. There exists a $B < \infty$ such that all $t \ge 1, \theta \in \Theta$, we have $\|\nabla \mathcal{L}(\theta, z_t)\| \le B$.

Assumption 4.2. For all $t \geq 1$, $\mathcal{L}(:, \mathbf{z}_t) \in \mathcal{H} = RKHS(K)$, where K is the kernel used in Algorithm 1. Moreover, there exists a constant $C_{\mathcal{X}} < \infty$ such that for all t, $\|\mathcal{L}(:, \mathbf{z}_t)\|_{\mathcal{H}} \leq C_{\mathcal{X}}$

Assumption 4.3. Assume that the objective function $f(\theta)$ is differentiable, ζ -smoothness, and λ -strongly convex, in the sense

$$(i) \quad f(\boldsymbol{\theta}_1) - f(\boldsymbol{\theta}_2) \leq \langle \nabla f(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \frac{\zeta}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2, \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \boldsymbol{\Theta} \subseteq \mathbb{R}^p,$$

$$(ii) \quad f(\boldsymbol{\theta}_1) - f(\boldsymbol{\theta}_2) \geq \langle \nabla f(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle + \frac{\lambda}{2} \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|^2, \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \boldsymbol{\Theta} \subseteq \mathbb{R}^p.$$

Assumption 4.1 ensures that the sensitivity of the gradient is uniformly bounded, a condition frequently imposed in LDP optimization to control the amount of noise required for privacy see, e.g., Song et al. (2013); Avella-Medina et al. (2023). Assumption 4.2 requires the target function to lie within the kernel-induced space, enabling convergence and estimation bounds under standard regularity conditions (Sopa et al., 2025). Assumption 4.3 imposes strong convexity and smoothness on the loss function, which are standard conditions for the convergence analysis of (stochastic) gradient optimization methods. Similar conditions can be found in Vaswani et al. (2022); Zhu et al. (2023).

Recall that $\hat{\theta}_t$ is the estimate obtained at the t-th iteration of the proposed LDP-BP Algorithm 1 under (ε, δ) -LDP, while θ^* denotes the true parameter value. The theorem below provides a non-asymptotic bound on the mean squared error of the estimate at iteration t.

Theorem 4.4 ((ε, δ) -**LDP**). Under Assumptions 4.1-4.3, there exist some positive constants $t_0 \in \mathbb{N}$ and c_p that depends on the dimension p, such that for $t \geq t_0$, $\hat{\Delta}_t = \hat{\theta}_t - \theta^*$ satisfies

$$E(\|\hat{\boldsymbol{\Delta}}_t\|_2^2) \lesssim t^{-\alpha} \{ (\eta c_p B^2 \log(1.25/\delta)/(\lambda \varepsilon^2) + \eta (L + p\kappa + 2B^2)/\lambda + \|\hat{\boldsymbol{\Delta}}_0\|_2^2 \},$$

when the step-size is chosen to be $\eta_t = \eta t^{-\alpha}$ with $\eta > 0$ and $1/2 < \alpha < 1$.

Theorem 4.4 establishes that the mean squared error $E(\|\hat{\Delta}_t\|_2^2)$ converges at rate $\mathcal{O}(t^{-\alpha})$ under the step size $\eta_t = \eta t^{-\alpha}$. The bound consists of three components: the privacy-induced noise term $B^2 \log(1.25/\delta)/(\lambda \varepsilon^2)$, the compression error $L + p\kappa$, and the error from the initial estimate. Notably, L can be made arbitrarily small by minimizing the acquisition function over p+1 points (Wu et al., 2023). Furthermore, as the compression budget $\kappa \to 0$, the rate coincides with that of Xie et al. (2025). Unlike their result, which requires a restrictive assumption on the conditional covariance of gradient noise, our analysis avoids this condition, thereby providing broader applicability.

In practice, however, many loss functions are neither strongly convex nor convex. Although non-convexity rules out guarantees of global optimality, our analysis relies only on the weaker assumption of ζ -smoothness, under which we establish convergence to an approximate stationary point. In non-convex settings with multiple local minima, convergence is typically analyzed through gradient norms rather than parameter estimates (Garrigos & Gower, 2023).

Theorem 4.5. Under Assumption 4.1, 4.2 and 4.3 (i), there exist some positive constants c', when the step-size is chosen to be $\eta_t = \eta t^{-\alpha}$ with $\eta > 0$ and $1/2 < \alpha < 1$, it follows that for every $t \ge 1$

$$\min_{1 \le i \le t} E \|\nabla f(\hat{\boldsymbol{\theta}}_i)\|^2 \le c' \frac{(f(\hat{\boldsymbol{\theta}}_0) - f(\boldsymbol{\theta}^*)) + \zeta(L + p\kappa + B^2) + pB^2/\varepsilon^2 \log(1.25/\delta)}{t^{1-\alpha}}.$$

Theorem 4.5 establishes an $O(t^{-(1-\alpha)})$ convergence rate of the gradient norm under a step size $\eta_t = \eta t^{-\alpha}$ in ζ -smooth optimization without assuming strong convexity. With a fixed step size, the rate reduces to the classical $O(t^{-1/2})$ result (Fang et al., 2023; Bu et al., 2023). The weaker ζ -smoothness assumption still enables meaningful gradient-based analysis, and by controlling the BO approximation error, our method achieves rates comparable to classical non-convex optimization (Garrigos & Gower, 2023). Notably, our guarantees avoid restrictive conditions such as fixing the Lipschitz constant to a specific value (e.g., 1), as required in prior work (Béthune et al., 2023).

In contrast to Theorem 4.4, which relies on strong convexity to establish a convergence rate for parameter estimation, the lack of convexity precludes direct control over the parameter error, thereby presenting a fundamental challenge. To address this, Theorem 4.5 leverages recursive moment bounds on the gradients and averaging techniques, yielding a convergence rate in gradient norm and guaranteeing convergence to an approximate stationary point. These findings align with existing literature (Stich, 2019; Garrigos & Gower, 2023): strong convexity enables rapid parameter recovery, whereas the general analysis guarantees convergence to stationarity in non-convex settings.

5 EXPERIMENTS

We assess the finite-sample performance of our method on two synthetic datasets and one real-world dataset, comparing it with LDP-SGD (Xie et al., 2025) in the parametric case and with a non-private deep neural network (Schmidhuber, 2015) in the nonparametric case. We compare the estimates of the coefficients based on 100 simulation replications. Details about the data generating process can be found in Appendix E.

Example 5.1 (Parametric Models). We evaluate the proposed LDP-BO algorithm on synthetic data under three regression settings: linear, logistic, and ReLU. We generate $T=20{,}000$ i.i.d. samples with features $x_t \sim N(0, \mathbf{I}_p)$ and true parameters $\theta=\mathbf{1}_p$, considering dimensions $p\in 2,5$. The compressed budget is set to $\kappa\in 0.1,0.2$, and the privacy budget is either fixed at $\varepsilon\in 1,2$ or randomly drawn from U(1,2) per iteration, with $\delta=0.2$. For comparison, we include LDP-SGD (Xie et al., 2025), as well as non-private BO and SGD as benchmarks.

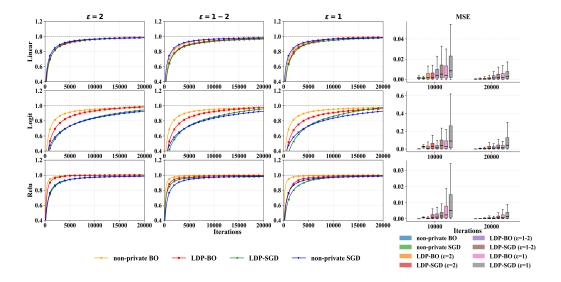


Figure 2: Evolution of the first-dimension coefficient estimate (true value = 1) and MSEs over iterations for linear, logistic, and ReLU models (rows) in Example 5.1. Columns correspond to privacy budgets $\varepsilon=2, \varepsilon\sim U(1,2)$, and $\varepsilon=1$, and Boxplots of coefficient MSEs.

The first three columns of Figure 2 shows the evolution of the average of the first-dimension coefficient estimate (true value = 1) over iterations. For simple models (linear), LDP-BO and LDP-SGD closely track their non-private counterparts, while in more complex models (logistic, ReLU), BO-based methods outperform SGD across all privacy levels. The last column of Figure 2 reports Mean-Squared Errors (MSEs) of the estimates, calculated as $\text{MSE} = \sum_{j=1}^p \text{MSE}_j/p = \sum_{j=1}^p \sum_{i=1}^t (\hat{\theta}_{i,j} - \theta_j)^2/(tp)$, where LDP-BO consistently achieves lower error and variability than LDP-SGD, especially in complex settings. Under strong privacy ($\varepsilon = 1$), LDP-BO converges faster and attains accuracy comparable to non-private BO and SGD. These results highlight LDP-BO's modeling advantage in nonlinear problems, mitigating the utility loss common in gradient-based methods. Results for p = 5 and varying compression budgets, reported in Appendix E, are similar.

Example 5.2 (Nonparametric Models). In this example, we evaluate our approach under nonparametric settings using the Sine and Friedman functions. A Gaussian process regression model is employed to estimate the unknown function, with kernel parameters optimized via our proposed LDP-BO framework (see Appendix E for details). We compare its utility against a non-private deep neural network (denoted as DNN) (Schmidhuber, 2015) trained incrementally with one data point per iteration.

We generate $T=10{,}000$ i.i.d. samples with features $\boldsymbol{x}_t \sim U(-1,1)$. For the Sine function, $y_t=\sin(2\pi x_t)+\varepsilon_t$; for the Friedman function, $y_t=\sin(\pi x_{1t}x_{2t})+(x_{3t}-0.5)^2+x_{4t}+x_{5t}+\varepsilon_t$, where $\varepsilon_t \sim \mathcal{N}(0,0.1^2)$. We set the compression budget to $\kappa=0.1$, the privacy budget to $(\varepsilon,\delta)=(1,0.2)$

 and B=2. We report the MSE of averaged estimators at sample sizes n=2000,5000, and 10000, and provide function fitting plots at n=10000 using 100 randomly generated test points.

Figure 3 presents the prediction errors (calculated as $\operatorname{error}_t = \frac{t-1}{t}\operatorname{error}_{t-1} + \frac{1}{t}(y_t - \hat{y}_t)^2$ in the online setting) and function fitting results for the proposed LDP-BO method and the DNN baseline. The LDP-BO method consistently outperforms the non-private DNN, even under privacy constraints. The boxplots show that LDP-BO achieves lower variance and fewer outliers, indicating greater stability and robustness across trials. The fitted curves further demonstrate that LDP-BO closely tracks the true function, capturing both global trends and fine-scale structure—particularly in high-value regions critical for optimization. In contrast, the DNN exhibits larger deviations and unstable oscillations, reflecting weaker generalization and poorly calibrated uncertainty.

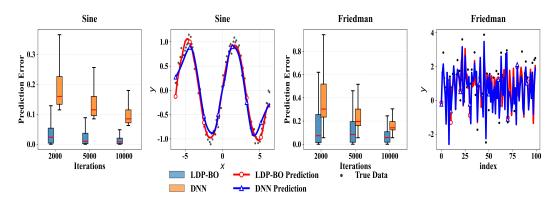


Figure 3: Predcition errors and function fitting plots of the proposed LDP-BO and DNN methods in Example 5.2.

Example 5.3 (Real Data Analysis). In this example, we apply LDP-SGD and the proposed LDP-BO to a real Uber Fares Dataset ¹, which comprises approximately 21,000 historical trip records collected between 2014 and 2015 in New York City. The selected features include distance, hour of day, day of week and passenger count; see Appendix E for full preprocessing details. These predictors, which collectively capture spatial, temporal, and demand-related determinants of Uber fare variations, have been similarly employed in prior studies (Khandelwal et al., 2021; Silveira-Santos et al., 2023; Huynh et al., 2025). The response is chosen to be the fare.

We adopt a Gaussian regression framework with a 4-dimensional parameter space for possible complex relationships. Among privacy-preserving methods, LDP-SGD applied to a linear model is the only one supporting both LDP and online parameter estimation; thus, we use it as a baseline for comparing prediction error across methods.

Figure 4 compares the performance of LDP-BO and LDP-SGD under $(\epsilon, \delta) = (1, 0.2)$ across sample sizes of 5000, 10000, and 20000 in terms of the prediction error. It show that LDP-BO consistently outperforms LDP-SGD across all metrics, achieving lower prediction error and exhibiting narrower interquartile ranges as sample size increases. This trend indicates reduced estimation variance and improved stability for LDP-BO.

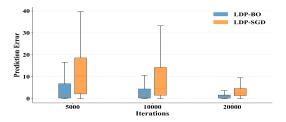


Figure 4: Fare prediction errors of LDP-BO and LDP-SGD in Example 5.3.

https://www.kaggle.com/datasets/yasserh/uber-fares-dataset

AUTHOR CONTRIBUTIONS

488 If you'd like to, you may include a section for author contributions as is done in many journals. This is optional and at the discretion of the authors.

491 ACKNOWLEDGMENTS

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper.

ETHICS STATEMENT

Our research strictly adheres to the ICLR Code of Ethics requirements in all aspects.

REPRODUCIBILITY STATEMENT

Algorithms 1-2, Section 5 and Appendix E have provided detailed information to ensure the reproduction of core results. We provide open access to the code with sufficient instructions, as described in supplemental material.

REFERENCES

- Martin Anthony and Peter L Bartlett. *Neural network learning: Theoretical foundations*. Cambridge University Press, 2009.
- Marco Avella-Medina, Casey Bradshaw, and Po-Ling Loh. Differentially private inference via noisy optimization. *The Annals of Statistics*, 51(5):2067–2092, 2023.
- Maximilian Balandat, Brian Karrer, Daniel Jiang, Samuel Daulton, Ben Letham, Andrew G Wilson, and Eytan Bakshy. Botorch: A framework for efficient monte-carlo bayesian optimization. *Advances in Neural Information Processing Systems*, 33:21524–21538, 2020.
- Felix Berkenkamp, Andreas Krause, and Angela P Schoellig. Bayesian optimization with safety constraints: safe and automatic parameter tuning in robotics. *Machine Learning*, 112(10):3713–3747, 2023.
- Louis Béthune, Thomas Massena, Thibaut Boissin, Yannick Prudent, Corentin Friedrich, Franck Mamalet, Aurélien Bellet, Mathieu Serrurier, and David Vigouroux. Dp-sgd without clipping: the lipschitz neural network way. *arXiv preprint arXiv:2305.16202*, 2023.
- Nicolas Bonneel, Julien Rabin, Gabriel Peyré, and Hanspeter Pfister. Sliced and radon wasserstein barycenters of measures. *Journal of Mathematical Imaging and Vision*, 51(1):22–45, 2015.
- Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010: 19th International Conference on Computational StatisticsParis France, August 22-27, 2010 Keynote, Invited and Contributed Papers*, pp. 177–186. Springer, 2010.
- Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- Zhiqi Bu, Yu-Xiang Wang, Sheng Zha, and George Karypis. Automatic clipping: Differentially private deep learning made easier and stronger. *Advances in Neural Information Processing Systems*, 36:41727–41764, 2023.
- Xi Chen, Jason D. Lee, Xin T. Tong, and Yichen Zhang. Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics*, 48(1):251 273, 2020.
- Robert Chew, Matthew R Williams, Elan A Segarra, Alexander J Preiss, Amanda Konet, and Terrance D Savitsky. Bayesian pseudo posterior mechanism for differentially private machine learning. arXiv preprint arXiv:2503.21528, 2025.
 - Taeryon Choi and Mark J Schervish. On posterior consistency in nonparametric regression problems. *Journal of Multivariate Analysis*, 98(10):1969–1987, 2007.

- Christos Dimitrakakis, Blaine Nelson, Zuhe Zhang, Aikaterini Mitrokotsa, and Benjamin IP Rubinstein. Differential privacy for bayesian inference through posterior sampling. *Journal of Machine Learning Research*, 18(11):1–39, 2017.
 - Qin Ding, Cho-Jui Hsieh, and James Sharpnack. An efficient algorithm for generalized linear bandit: Online stochastic gradient descent and thompson sampling. In *International Conference on Artificial Intelligence and Statistics*, pp. 1585–1593. PMLR, 2021.
 - Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(1):3–37, 2022.
 - John C Duchi and Feng Ruan. The right complexity measure in locally private estimation: It is not the fisher information. *The Annals of Statistics*, 52(1):1–51, 2024.
 - John C Duchi, Michael I Jordan, and Martin J Wainwright. Minimax optimal procedures for locally private estimation. *Journal of the American Statistical Association*, 113(521):182–201, 2018.
 - Cynthia Dwork. Differential privacy. In *International colloquium on automata, languages, and programming*, pp. 1–12. Springer, 2006.
 - Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
 - Yaakov Engel, Shie Mannor, and Ron Meir. The kernel recursive least-squares algorithm. *IEEE Transactions on Signal Processing*, 52(8):2275–2285, 2004.
 - David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. Scalable global optimization via local bayesian optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Huang Fang, Xiaoyun Li, Chenglin Fan, and Ping Li. Improved convergence of differential private sgd with gradient clipping. In *The Eleventh International Conference on Learning Representations*, 2023.
 - Peter I Frazier. A tutorial on bayesian optimization. arXiv preprint arXiv:1807.02811, 2018.
 - Shi Fu, Fengxiang He, Xinmei Tian, and Dacheng Tao. Convergence of bayesian bilevel optimization. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Guillaume Garrigos and Robert M Gower. Handbook of convergence theorems for (stochastic) gradient methods. *arXiv preprint arXiv:2301.11235*, 2023.
 - Ruijian Han, Lan Luo, Yuanyuan Lin, and Jian Huang. Online inference with debiased stochastic gradient descent. *Biometrika*, 111(1):93–108, 2024.
 - Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
 - Mikko Heikkilä, Eemil Lagerspetz, Samuel Kaski, Kana Shimizu, Sasu Tarkoma, and Antti Honkela. Differentially private bayesian learning on distributed data. *Advances in Neural Information Processing Systems*, 30, 2017.
 - Tuyet Ngoc Thi Huynh, Huu Dat Bui, Tuyet Nam Thi Nguyen, and Tan Dat Trinh. Enhancing prediction of ride-hailing fares using advanced deep learning techniques. *New Trends in Computer Sciences*, 3(1):64–82, 2025.
 - Carl Hvarfner, Erik Orm Hellsten, and Luigi Nardi. Vanilla bayesian optimization performs great in high dimensions. *In International Conference on Machine Learning*, pp. 20793–20817, 2024.
 - Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998.
 - Kunal Khandelwal, Atharva Sawarkar, and Swati Hira. A novel approach for fare prediction using machine learning techniques. *International Journal of Next-Generation Computing*, 12(5), 2021.

- Soheil Kolouri, Se Rim Park, and Gustavo K Rohde. The radon cumulative distribution transform and its application to image classification. *IEEE Transactions on Image Processing*, 25(2):920–934, 2015.
 - Alec Koppel, Hrusikesha Pradhan, and Ketan Rajawat. Consistent online gaussian process regression without the sample complexity bottleneck. *Statistics and Computing*, 31(6):76, 2021.
 - Ximing Li, Chendi Wang, and Guang Cheng. Statistical theory of differentially private marginal-based data synthesis algorithms. *arXiv* preprint arXiv:2301.08844, 2023.
 - WeiKang Liu, Yanchun Zhang, Hong Yang, and Qinxue Meng. A survey on differential privacy for medical data analysis. *Annals of Data Science*, 11(2):733–747, 2024.
 - Andrew Lowy and Meisam Razaviyayn. Private federated learning without a trusted server: Optimal algorithms for convex losses. In *The Eleventh International Conference on Learning Representations*, 2023.
 - Disha Makhija, Joydeep Ghosh, and Nhat Ho. A bayesian approach for personalized federated learning in heterogeneous settings. *Advances in Neural Information Processing Systems*, 37: 102428–102455, 2024.
 - Ilya Mironov. Rényi differential privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF), pp. 263–275. IEEE, 2017.
 - Jonas Močkus. On bayesian methods for seeking the extremum. In *IFIP Technical Conference on Optimization Techniques*, pp. 400–404. Springer, 1974.
 - Sarah Müller, Alexander von Rohr, and Sebastian Trimpe. Local policy search with bayesian optimization. *Advances in Neural Information Processing Systems*, 34:20708–20720, 2021.
 - Willie Neiswanger, Lantao Yu, Shengjia Zhao, Chenlin Meng, and Stefano Ermon. Generalizing bayesian optimization with decision-theoretic entropies. *Advances in Neural Information Processing Systems*, 35:21016–21029, 2022.
 - Quan Nguyen, Kaiwen Wu, Jacob Gardner, and Roman Garnett. Local bayesian optimization via maximizing probability of descent. *Advances in Neural Information Processing Systems*, 35: 13190–13202, 2022.
 - Quoc Phong Nguyen, Wan Theng Ruth Chew, Le Song, Bryan Kian Hsiang Low, and Patrick Jaillet. Optimistic bayesian optimization with unknown constraints. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Natalia Ponomareva, Hussein Hazimeh, Alex Kurakin, Zheng Xu, Carson Denison, H Brendan McMahan, Sergei Vassilvitskii, Steve Chien, and Abhradeep Guha Thakurta. How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research*, 77:1113–1201, 2023.
 - Sachin Ravi and Alex Beatson. Amortized bayesian meta-learning. In *International Conference on Learning Representations*, 2019.
 - Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pp. 400–407, 1951.
 - Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
 - Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015.
 - Xueyuan She, Saurabh Dash, and Saibal Mukhopadhyay. Sequence approximation using feedforward spiking neural network for spatiotemporal learning: Theory and optimization methods. In *International Conference on Learning Representations*, 2021.

- Tulio Silveira-Santos, Anestis Papanikolaou, Thais Rangel, and Jose Manuel Vassallo. Understanding and predicting ride-hailing fares in madrid: A combination of supervised and unsupervised techniques. *Applied Sciences*, 13(8):5147, 2023.
 - Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. *Advances in Neural Information Irocessing Systems*, 25, 2012.
 - Shuang Song, Kamalika Chaudhuri, and Anand D Sarwate. Stochastic gradient descent with differentially private updates. In 2013 IEEE Global Conference on Signal and Information Processing, pp. 245–248. IEEE, 2013.
 - Getoar Sopa, Juraj Marusic, Marco Avella-Medina, and John P Cunningham. Scalable differentially private bayesian optimization. *arXiv preprint arXiv:2502.06044*, 2025.
 - Sebastian U Stich. Unified optimal analysis of the (stochastic) gradient method. *arXiv preprint* arXiv:1907.04232, 2019.
 - Weijie J Su and Yuancheng Zhu. Higrad: Uncertainty quantification for online learning and stochastic approximation. *Journal of Machine Learning Research*, 24(124):1–53, 2023.
 - Sebastian Shenghong Tay, Chuan-Sheng Foo, Daisuke Urano, Richalynn Leong, and Bryan Kian Hsiang Low. A unified framework for bayesian optimization under contextual uncertainty. In *The Twelfth International Conference on Learning Representations*, 2024.
 - Aleksei Triastcyn and Boi Faltings. Bayesian differential privacy for machine learning. In *International Conference on Machine Learning*, pp. 9583–9592. PMLR, 2020.
 - Sharan Vaswani, Benjamin Dubois-Taine, and Reza Babanezhad. Towards noise-adaptive, problem-adaptive (accelerated) stochastic gradient descent. In *International Conference on Machine Learning*, pp. 22015–22059. PMLR, 2022.
 - Cédric Villani. Topics in optimal transportation, volume 58. American Mathematical Soc., 2021.
 - Jian Wu, Matthias Poloczek, Andrew G Wilson, and Peter Frazier. Bayesian optimization with gradients. *Advances in Neural Information Processing Systems*, 30, 2017.
 - Kaiwen Wu, Kyurae Kim, Roman Garnett, and Jacob Gardner. The behavior and convergence of local bayesian optimization. *Advances in Neural Information Processing Systems*, 36:73497–73523, 2023.
 - Jinhan Xie, Enze Shi, Bei Jiang, Linglong Kong, and Xuming He. Online differentially private inference in stochastic gradient descent. *arXiv preprint arXiv:2505.08227*, 2025.
 - Xingxing Xiong, Shubo Liu, Dan Li, Zhaohui Cai, and Xiaoguang Niu. A comprehensive survey on local differential privacy. *Security and Communication Networks*, 2020(1):8829523, 2020.
 - Wanrong Zhang and Ruqi Zhang. Dp-fast mh: Private, fast, and accurate metropolis-hastings for large-scale bayesian inference. In *International Conference on Machine Learning*, pp. 41847–41860. PMLR, 2023.
 - Wanrong Zhu, Xi Chen, and Wei Biao Wu. Online covariance matrix estimation in stochastic gradient descent. *Journal of the American Statistical Association*, 118(541):393–404, 2023.

A BACKGROUND ON LDP AND SLICED WASSERSTEIN DISTANCE

A.1 DIFFERENTIAL PRIVACY

In this section, we begin with the basic concepts and properties of Local Differential Privacy (LDP), Rényi Differential Privacy (RDP) and Gaussian Differential Privacy (GDP). The intuition underlying LDP is that a randomized algorithm produces outputs that are statistically similar, even when a single individual's information in the dataset is modified or removed, thereby ensuring the protection of individual privacy. The formal definition of LDP is presented below.

Definition A.1. $((\epsilon, \delta)\text{-}LDP \ (Xiong \ et \ al., 2020))$ Let \mathcal{X} be the sample space for an individual data, a randomized algorithm $\mathcal{A}: \mathcal{X} \to \mathbb{R}$ is $(\epsilon, \delta)\text{-}LDP$ if and only if for any pair of input single values $z, z' \in \mathcal{X}$ and for any $S \subseteq \mathbb{R}$, the inequality below holds

$$\mathbb{P}\left(\mathcal{A}(\boldsymbol{z}) \in S\right) \leq e^{\varepsilon} \cdot \mathbb{P}\left(\mathcal{A}(\boldsymbol{z}') \in S\right) + \delta.$$

In contrast to CDP, LDP imposes a stricter requirement in which each individual perturbs their data locally before submission. This design eliminates the need for a trusted data curator and is particularly well suited to streaming environments, where data are continuously generated and transmitted. To formalize the guarantee, we introduce the notion of sensitivity, which quantifies the maximum change in an algorithm's output resulting from the modification of a single data entry.

Definition A.2. For any deterministic function $g: \mathcal{X} \to \mathbb{R}$ and any pair of input single values $z, z' \in \mathcal{X}$, the ℓ_p -sensitivity of g is defined as

$$\Delta_p(g) = \sup_{\boldsymbol{z}, \boldsymbol{z}' \in \mathcal{X}} \|g(\boldsymbol{z}) - g(\boldsymbol{z}')\|_p.$$

Among various LDP mechanisms, we introduce the following Gaussian mechanism for illustrative purposes, as it facilitates clear exposition.

Definition A.3. (The Gaussian Mechanism (Dwork, 2006)) Let $g: \mathcal{X} \to \mathbb{R}$ be a deterministic function with $\Delta_2(g) < \infty$. For $\mathbf{w} \in \mathbb{R}$ with coordinates w_1, w_2, \cdots, w_p be i.i.d samples drawn from $N(0, 2(\Delta_2(g)/\varepsilon)^2 \log(1.25/\delta))$, $g(\mathbf{z}) + \mathbf{w}$ is (ε, δ) -LDP.

The post-processing and parallel composition properties are fundamental to LDP, enabling complex algorithms to be systematically constructed from simpler components.

Proposition A.4. (Post-processing Property for LDP (Xiong et al., 2020)) Let A be an (ε, δ) -LDP algorithm and g be an arbitrary function which takes A(z) as input, then g(A(z)) is also (ε, δ) -LDP.

Proposition A.5. (Parallel Composition for LDP (Xiong et al., 2020)) Suppose n mechanisms $\{A_1, \ldots, A_n\}$ satisfy $(\varepsilon_i, \delta_i)$ -LDP, respectively, and are computed on disjoint subsets of data, then a mechanism formed by $(A_1(z_1), \ldots, A_n(z_n))$ satisfies $(\max(\varepsilon_i), \max(\delta_i))$ -LDP.

As an alternative to standard LDP, RDP was introduced by Mironov (2017) as a generalization of LDP based on Rényi divergence, providing a more structured and flexible framework for privacy accounting. RDP quantifies privacy loss through the Rényi divergence of order q>1 between the output distributions of an algorithm on adjacent datasets. For two probability distributions P and Q, the Rényi divergence of order q is defined as

$$D_q(P||Q) = \frac{1}{q-1} \log E_Q \left\{ \left(\frac{P}{Q}\right)^{q-1} \right\},\,$$

whenever the expectation exists. This divergence provides a smooth and fine-grained measure of dissimilarity that depends on the order q, thereby enabling more precise tracking of cumulative privacy loss under composition compared to the standard (ε, δ) -LDP framework. Formally, RDP is defined as follows:

Definition A.6. (RDP, Mironov (2017)). Let A be a randomized algorithm, and let z and z' be two adjacent datasets. For any real number $\alpha > 1$, the algorithm A satisfies (q, ε) -RDP if

$$D_q(\mathcal{A}(z) \| \mathcal{A}(z')) \le \varepsilon,$$

where A(z) denotes the distribution of the output of A on data z.

Building on this hypothesis testing framework, Dong et al. (2022) introduced GDP, a privacy notion with a natural statistical interpretation: determining whether an individual's data is included in a dataset is at least as difficult as distinguishing between N(0,1) and $N(\mu,1)$ based on a single observation, for some $\mu > 0$. Formally, GDP is defined as follows:

Definition A.7. (GDP, Dong et al. (2022) Let A be a randomized algorithm.

- 1. A satisfies f-DP if, for any α -level test of H_0 , the power function $\beta(\alpha)$ satisfies $\beta(\alpha) \leq 1 f(\alpha)$, where f is convex, continuous, non-increasing, and $f(\alpha) \leq 1 \alpha$ for all $\alpha \in [0, 1]$.
- 2. A satisfies μ -GDP if it is f-DP with $f(\alpha) \ge \Phi(\Phi^{-1}(1-\alpha)-\mu)$ for all $\alpha \in [0,1]$, where $\Phi(\cdot)$ denotes the standard normal CDF.

A.2 SLICED WASSERSTEIN DISTANCE

Definition A.8. (Wasserstein Distance (Villani, 2021)) The Wasserstein distance $W_p(u, \nu)$ quantifies the optimal transport cost between two probability distributions u and ν , defined as the minimal expected cost required to redistribute mass from u to ν . For univariate distributions, it admits the closed-form

$$W_p(u,\nu) = \left(\int_{\mathcal{X}} \left| x - F_{\nu}^{-1}(F_u(x)) \right|^p du(x) \right)^{1/p} = \left(\int_0^1 \left| F_u^{-1}(t) - F_{\nu}^{-1}(t) \right|^p dt \right)^{1/p},$$

where $F(\cdot)$ denotes the cumulative distribution function (CDF). In particular, if $u = N(m_1, \sigma_1^2)$ and $\nu = N(m_2, \sigma_2^2)$, are univariate Gaussian distributions, their 2-Wasserstein distance admits the analytic form $W_2(u, \nu) = \sqrt{(m_1 - m_2)^2 + (\sigma_1 - \sigma_2)^2}$.

Definition A.9. (Sliced Wasserstein (SW) Distance (Bonneel et al., 2015)) The Sliced Wasserstein distance generalizes the Wasserstein distance to higher dimensions via Radon transforms. Specifically, it projects multivariate distributions onto one-dimensional subspaces determined by directions $\theta \in \mathbb{S}^{p-1}$, computes the Wasserstein distance between these projections, and then averages across directions:

$$SW_p(u,\nu) = \left(\int_{\boldsymbol{\theta} \in \mathbb{S}^{p-1}} W_p^p(\mathcal{R}u_{\boldsymbol{\theta}}, \mathcal{R}\nu_{\boldsymbol{\theta}}) d\boldsymbol{\theta}\right)^{1/p}.$$

In practice, the SW distance is typically approximated using Monte Carlo sampling over m random directions: $SW_p(u,\nu) \approx \{\sum_{l=1}^m W_p^p \left(\mathcal{R}u_{\theta},\mathcal{R}\nu_{\theta}\right)/m\}^{1/p}$. For our experiments, we used a value of m=100.

B ADDITIONAL COROLLARIES

In this section, we additionally present two corollaries that provide non-asymptotic error bounds for the LDP-BO algorithm under specific privacy definitions.

Corollary B.1 ((q, ε) -**RDP**). Suppose the conditions of Theorem 4.4 hold. Under (q, ε) -Rényi Differential Privacy (RDP), where noise $\omega_t = B\sqrt{q/(2\varepsilon)} \cdot N(0, \mathbf{I}_p)$ is added at each iteration in Algorithm 1, the expected estimation error satisfies

$$E(\|\hat{\boldsymbol{\Delta}}_t\|_2^2) \lesssim t^{-\alpha} \{ (\eta c_p B^2 q/(2\lambda \varepsilon) + \eta (L + p\kappa + 2B^2)/\lambda + \|\hat{\boldsymbol{\Delta}}_0\|_2^2 \}.$$

Corollary B.2 (μ -GDP). Suppose the conditions of Theorem 4.4 hold. Under μ -Gaussian Differential Privacy (GDP), where noise $\omega_t = \frac{2B}{\mu} \cdot N(0, \mathbf{I}_p)$ is added at each iteration in Algorithm 1, the expected estimation error satisfies

$$E(\|\hat{\Delta}_t\|_2^2) \lesssim t^{-\alpha} \{ (\eta c_p B^2 / (\lambda \mu^2) + \eta (L + p\kappa + 2B^2) / \lambda + \|\hat{\Delta}_0\|_2^2 \}.$$

Corollaries B.1–B.2 present the expected estimation error under two specific privacy definitions, (q,ε) -RDP and μ -GDP. The bounds follow the same structure as Theorem 4.5, with identical second and third components, while the first component varies by privacy definition. Specifically, Corollary B.1 shows that (α,ε) -Rényi DP improves the bound from $\mathcal{O}\left(t^{-\alpha}\cdot B^2\log(1/\delta)/(\lambda\varepsilon^2)\right)$ to $\mathcal{O}\left(t^{-\alpha}\cdot B^2\alpha/(\lambda\varepsilon)\right)$, whereas Corollary B.2 demonstrates that μ -Gaussian DP yields a bound of order $\mathcal{O}\left(t^{-\alpha}\cdot B^2/(\lambda\mu^2)\right)$.

C SUPPORTING LEMMAS

Lemma C.1. Let ρ_0 be the corresponding true population posterior distribution. Suppose the following conditions hold:

- i. For any measurable subset $A \subseteq [0,1]^p$ with Lebesgue measure $\lambda(A) \geq (K_p t)^{-1}$, where $K_p \in (0,1]$ is a constant, A contains at least one sample point θ_t .
- ii. For all $t \geq 1$, the kernel matrix is positive definite $K_t \succ 0$.
- 818
 819
 820
 iii. The covariance kernel is translation-invariant, taking the form $K(\theta, \theta') = K(\beta \| \theta \theta' \|)$ for some scale parameter $\beta > 0$.
 - iv. There exist constants $\delta \in (0, 1/2)$ and $b_1, b_2 > 0$ such that for all $t \ge 1$, $\mathbb{P}_{\Pi} \left(\beta > t^{\delta} \right) < b_1 e^{-b_2 t}$, where \mathbb{P}_{Π} denotes the probability under the Gaussian prior Π for β .
 - Then, the posterior distribution without compression ρ_t is asymptotically consistent, i.e. for every c > 0,

$$\mathbb{P}\left(\mathrm{SW}_2(\rho_t, \rho_0) < c \mid \mathcal{D}_t\right) \to 1 \quad (a.s.).$$

Proof of Lemma C.1. The results of this lemma are well established, with detailed proofs provided in Theorem 6 of Choi & Schervish (2007).

Lemma C.2. Assuming the regularity conditions specified in Lemma C.1, which guarantee the well-behaved geometry of the target distribution, Algorithm 2 achieves κ -approximate convergence under the SW metric. Specifically, for any c > 0

$$\lim_{t \to \infty} \mathbb{P}\left\{ SW_2(\rho_{\mathcal{D}_t}, \rho_{\mathcal{D}_{t-1}}) < c + \kappa \mid \mathcal{D}_t \right\} = 1.$$

Proof of Lemma C.2. Using triangle inequality, we obtain

$$\mathrm{SW}_2(\rho_{\mathcal{D}_t}, \rho_{\mathcal{D}_{t-1}}) \leq \mathrm{SW}_2(\rho_{\mathcal{D}_t}, \rho_{\tilde{\mathcal{D}}_t}) + \mathrm{SW}_2(\rho_{\tilde{\mathcal{D}}_t}, \rho_{\mathcal{D}_{t-1}}),$$

The first term corresponds exactly to the stopping criterion in Algorithm 2, and is therefore bounded above by κ . Consequently, following the argument of Koppel et al. (2021), we have the following containment relationship for any c' > 0:

$$\{\operatorname{SW}_{2}(\rho_{\mathcal{D}_{t}}, \rho_{\mathcal{D}_{t-1}}) < c'\} \subset \{\operatorname{SW}_{2}(\rho_{\mathcal{D}_{t}}, \rho_{\tilde{\mathcal{D}}_{t}}) + \operatorname{SW}_{2}(\rho_{\tilde{\mathcal{D}}_{t}}, \rho_{\mathcal{D}_{t-1}}) < c'\}$$
$$\subset \{\operatorname{SW}_{2}(\rho_{\tilde{\mathcal{D}}_{t}}, \rho_{\mathcal{D}_{t-1}}) + \kappa < c'\}.$$

Taking prior probability with respect to Π , it follows that

$$\mathbb{P}_{\Pi}\{\mathrm{SW}_{2}(\rho_{\mathcal{D}_{t}}, \rho_{\mathcal{D}_{t-1}}) < c'\} \leq \mathbb{P}_{\Pi}\{\mathrm{SW}_{2}(\rho_{\mathcal{D}_{t}}, \rho_{\tilde{\mathcal{D}}_{t}}) + \mathrm{SW}_{2}(\rho_{\tilde{\mathcal{D}}_{t}}, \rho_{\mathcal{D}_{t-1}}) < c'\} \\
\leq \mathbb{P}_{\Pi}\{\mathrm{SW}_{2}(\rho_{\tilde{\mathcal{D}}_{t}}, \rho_{\mathcal{D}_{t-1}}) + \kappa < c'\} \\
\leq \mathbb{P}_{\Pi}\{\mathrm{SW}_{2}(\rho_{\tilde{\mathcal{D}}_{t}}, \rho_{\mathcal{D}_{t-1}}) < c' - \kappa\}$$

By Assumption 3.2, which states that $\mathbb{P}_{\Pi}\{\psi_t\} \geq \mathbb{P}_{\Pi}\{\tilde{\psi}_t\}$, we have

$$\mathbb{P}_{\Pi}\{SW_2(\rho_{\tilde{\mathcal{D}}_t}, \rho_{\mathcal{D}_{t-1}}) < c' - \kappa\} \le \mathbb{P}_{\Pi}\{SW_2(\rho_t, \rho_{t-1}) < c' - \kappa\}$$

By Lemma C.1 the supremum of the probability of the right-hand side of tends 1 as $t\to\infty$ for $c=c'-\kappa>0$. Therefore

$$\lim_{t \to \infty} \sup \mathbb{P}_{\Pi} \{ SW_2(\rho_{\mathcal{D}_t}, \rho_{\mathcal{D}_{t-1}}) < c' \} = 1.$$

Exploiting the continuity of both the GP posterior and the SW metric, we conclude that the above limit exists. Substituting $c'=c+\kappa$, Lemma C.2 follows.

Lemma C.3. For a vector $v \in \mathbb{R}^p$, define the projection operator $\Pi_B(v) = v \cdot \min\{1, \frac{B}{\|v\|}\}$, which projects v onto the Euclidean ball $B_B(0)$ of radius B centered at the origin. Under Assumption 4.1, we have, $\forall t \geq 1$,

$$\|\Pi_B(\boldsymbol{\mu}_{\mathcal{D}_t}) - \nabla \mathcal{L}(\boldsymbol{\theta}_t, z_t)\| \le \|\boldsymbol{\mu}_{\mathcal{D}_t} - \nabla \mathcal{L}(\boldsymbol{\theta}_t, z_t)\|.$$

Proof of Lemma C.3. Notice that $\Pi_B(x) = \arg\min_{x' \in B_B(0)} \|x - x'\|$, that is, $\Pi_B(x)$ is the Euclidean projection of x onto the ball $B_B(0)$. Now, let $y \in B_B(0)$. Since $B_B(0)$ is convex, for any $0 < \eta < 1$, the convex combination $z := \eta y + (1 - \eta)\Pi_B(x) = \Pi_B(x) + \eta(y - \Pi_B(x))$, also belongs to $B_B(0)$, i.e., $z \in B_B(0)$.

We then obtain

$$||x - \Pi_B(x)||^2 \le ||x - z||^2 = ||x - \Pi_B(x) - \eta(y - \Pi_B(x))||^2$$

$$= ||x - \Pi_B(x)||^2 + \eta^2 ||y - \Pi_B(x)||^2 - 2\eta\langle x - \Pi_B(x), y - \Pi_B(x)\rangle,$$
(6)

where the inequality follows from the definition of $\Pi_B(x)$ as the closest point in $B_B(0)$ to x. Thus, we have

$$\langle x - \Pi_B(x), \Pi_B(x) - y \rangle + \frac{\eta}{2} ||y - \Pi_B(x)||^2 \ge 0.$$

As $0 < \eta < 1$ is arbitrary, we obtain

$$\langle x - \Pi_B(x), \Pi_B(x) - y \rangle = \lim_{n \to 0^+} \langle x - \Pi_B(x), \Pi_B(x) - y \rangle + \frac{\eta}{2} ||y - \Pi_B(x)||^2 \ge 0$$

for all $y \in B_B(0)$. Using inequality (6), we can further derive the following bound:

$$\|\boldsymbol{\mu}_{\mathcal{D}_{t}} - \nabla \mathcal{L}(\boldsymbol{\theta}_{t}, z_{t})\|^{2} = \|\boldsymbol{\mu}_{\mathcal{D}_{t}} - \Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}) + \Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}) - \nabla \mathcal{L}(\boldsymbol{\theta}_{t}, z_{t})\|^{2}$$

$$= \|\boldsymbol{\mu}_{\mathcal{D}_{t}} - \Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}})\|^{2} + \|\Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}) - \nabla \mathcal{L}(\boldsymbol{\theta}_{t}, z_{t})\|^{2}$$

$$+ 2 \langle \boldsymbol{\mu}_{\mathcal{D}_{t}} - \Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}), \Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}) - \nabla \mathcal{L}(\boldsymbol{\theta}_{t}, z_{t})\rangle$$

$$> \|\Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}) - \nabla \mathcal{L}(\boldsymbol{\theta}_{t}, z_{t})\|^{2},$$

where the final inequality follows from the fact that both the first and last terms on the right-hand side of (6) are nonnegative, since by Assumption 4.1 we have $\nabla \mathcal{L}(\theta_t, z_t) \in B_B(0)$.

Lemma C.4. Assume Assumption 4.1 and Assumption 4.2 hold. let $\theta \in \Theta$ and let \mathcal{D} denote a set containing points θ . Denote $g(\theta_t) = \Pi_B(\nabla K(\theta_t, \mathcal{D}_t)K(\mathcal{D}_t, \mathcal{D}_t)^{-1}f(\theta_t))$. Then, there exists some constant $c_1 > 0$ such that

$$\|\nabla f(\boldsymbol{\theta}_t) - g(\boldsymbol{\theta}_t)\|^2 \le c_1(L + p\kappa).$$

Proof of Lemma C.4. Combining Assumption 4.2 with Lemma C.3 of Wu et al. (2023), we obtain

$$\|\nabla f(\boldsymbol{\theta}_t) - g(\boldsymbol{\theta}_t)\|^2 \le \|\nabla f(\boldsymbol{\theta}_t) - \nabla K(\boldsymbol{\theta}_t, \mathcal{D}_t) K(\mathcal{D}_t, \mathcal{D}_t)^{-1} f(\boldsymbol{\theta}_t, z_t)\|^2 \le C_{\mathcal{X}} Tr(\nabla^2 K_{\mathcal{D}_t}(\boldsymbol{\theta}_t, \boldsymbol{\theta}_t)),$$

Since D_t is obtained by compressing $\tilde{D}_t = D_{t-1} \cup \boldsymbol{\xi}$, we then have

$$SW_2(\rho_{\mathcal{D}_t}, \rho_{\tilde{\mathcal{D}}_t}) \leq \kappa.$$

Using the expression of the Sliced Wasserstein distance for multivariate normal distributions, it follows that

$$SW_{2}^{2}(\rho_{\mathcal{D}_{t}}, \rho_{\tilde{\mathcal{D}}_{t}})$$

$$= E_{\boldsymbol{\theta} \sim \mathcal{U}(\mathbb{S}^{p-1})} \left[(\boldsymbol{\theta}^{\top} (\boldsymbol{\mu}_{t+1}|_{\mathcal{D}_{t}} - \boldsymbol{\mu}_{t+1}|_{\tilde{\mathcal{D}}_{t}}))^{2} + \left(\sqrt{\boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\mathcal{D}_{t}} \boldsymbol{\theta}} - \sqrt{\boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\tilde{\mathcal{D}}_{t}} \boldsymbol{\theta}} \right)^{2} \right]$$

$$\leq \kappa^{2}.$$

This implies $E_{\theta \sim \mathcal{U}(\mathbb{S}^{p-1})} \{ (\sqrt{\theta^{\top} \Sigma_{t+1}}|_{\mathcal{D}_t} \theta - \sqrt{\theta^{\top} \Sigma_{t+1}}|_{\tilde{\mathcal{D}}_t} \theta)^2 \} \leq \kappa^2$. Notice that θ is the projection on the unit sphere. We then have $E_{\theta \sim \mathcal{U}(\mathbb{S}^{p-1})} [\theta^{\top} \Sigma \theta] = \frac{1}{2} tr(\Sigma)$. Therefore, we obtain

$$tr(\Sigma_{t+1}|_{\mathcal{D}_t}) - tr(\Sigma_{t+1}|_{\tilde{\mathcal{D}}_t}) = p \cdot E_{\boldsymbol{\theta} \sim \mathcal{U}(\mathbb{S}^{p-1})} \left[\boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\mathcal{D}_t} \boldsymbol{\theta} - \boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\tilde{\mathcal{D}}_t} \right].$$

Hence.

$$\boldsymbol{\theta}^{\top}(\Sigma_{t+1}|_{\mathcal{D}_t} - \Sigma_{t+1}|_{\tilde{\mathcal{D}}_t})\boldsymbol{\theta} = \left(\sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\mathcal{D}_t}\boldsymbol{\theta}} + \sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\tilde{\mathcal{D}}_t}\boldsymbol{\theta}}\right)\left(\sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\mathcal{D}_t}\boldsymbol{\theta}} - \sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\tilde{\mathcal{D}}_t}\boldsymbol{\theta}}\right)$$

Without loss of generality, assume the operator (spectral) norms of $\sqrt{\boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\mathcal{D}_t} \boldsymbol{\theta}}$ and $\sqrt{\boldsymbol{\theta}^{\top} \Sigma_{t+1}|_{\tilde{\mathcal{D}}_t} \boldsymbol{\theta}}$ are uniformly bounded by C. We then have

$$\boldsymbol{\theta}^{\top}(\Sigma_{t+1}|_{\mathcal{D}_t} - \Sigma_{t+1}|_{\tilde{\mathcal{D}}_t})\boldsymbol{\theta} \leq 2C\left(\sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\mathcal{D}_t}\boldsymbol{\theta}} - \sqrt{\boldsymbol{\theta}^{\top}\Sigma_{t+1}|_{\tilde{\mathcal{D}}_t}\boldsymbol{\theta}}\right)$$

Therefore, we obtain

$$tr(\Sigma_{t+1}|_{\mathcal{D}_t}) - tr(\Sigma_{t+1}|_{\tilde{\mathcal{D}}_t}) \le 2Cp\kappa.$$

As established in the discussion of BO (Wu et al., 2023), there exists some constant L > 0 such that

$$tr(\Sigma_{t+1}|_{\tilde{\mathcal{D}}_{t}}) = tr(\nabla^{2}K_{D\cup\mathbf{z}}(\boldsymbol{\theta},\boldsymbol{\theta})) \leq L.$$

Consequently, we obtain that, for some constant $c_1 > 0$,

$$\|\nabla \mathcal{L}(\boldsymbol{\theta}_t, z_t) - \boldsymbol{\mu}_{\mathcal{D}_t}\|^2 \le c_1(L + p\kappa).$$

П

Lemma C.5. Let $g_t(\theta_t)$ be defined as in Algorithm 1. Under Assumptions 4.1 and 4.2, there exists some constant $c_1 > 0$ such that

$$||g_t(\boldsymbol{\theta}_t) - g(\boldsymbol{\theta}_t)||^2 \le 2B^2.$$

Proof of Lemma C.5. Using Lemma C.3, the effect of the projection operator Π_B can be removed from the analysis. Consequently, we obtain

$$\|g_{t}(\boldsymbol{\theta}_{t}) - g(\boldsymbol{\theta}_{t})\|^{2} = \|\Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}(z_{t})) - \Pi_{B}(\nabla K(\boldsymbol{\theta}_{t}, \mathcal{D}_{t})K(\mathcal{D}_{t}, \mathcal{D}_{t})^{-1}f(\boldsymbol{\theta}_{t}))\|^{2}$$

$$\leq \|\Pi_{B}(\boldsymbol{\mu}_{\mathcal{D}_{t}}(z_{t}))\|^{2} + \|\Pi_{B}(\nabla K(\boldsymbol{\theta}_{t}, \mathcal{D}_{t})K(\mathcal{D}_{t}, \mathcal{D}_{t})^{-1}f(\boldsymbol{\theta}_{t}))\|^{2}$$

$$\leq B^{2} + B^{2}$$

$$\leq 2B^{2}.$$

Lemma C.6. (1) Suppose that $f: \mathbb{R}^p \to \mathbb{R}$ is a λ -strongly convex function, we have

$$\langle \nabla f(\boldsymbol{\theta}_1) - \nabla f(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle \ge \lambda \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2^2, \quad \forall \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^p,$$

and if f is twice-differentiable, then $\nabla^2 f(\boldsymbol{\theta}) \succeq \lambda I$, $\forall \boldsymbol{\theta} \in \mathbb{R}^p$.

(2) Suppose that $f: \mathbb{R}^p \to \mathbb{R}$ is a convex and ζ -smooth function, we have for any $\theta_1, \theta_2 \in \mathbb{R}^p$,

$$\|\nabla f(\boldsymbol{\theta}_1) - \nabla f(\boldsymbol{\theta}_2)\|_2^2 < \zeta \langle \nabla f(\boldsymbol{\theta}_1) - \nabla f(\boldsymbol{\theta}_2), \boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 \rangle,$$

and

$$\|\nabla f(\boldsymbol{\theta}_1) - \nabla f(\boldsymbol{\theta}_2)\|_2 \le \zeta \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_2.$$

If f is twice-differentiable, then $\nabla^2 f(\boldsymbol{\theta}) \preceq \zeta I$, $\forall \boldsymbol{\theta} \in \mathbb{R}^p$.

Proof of Lemma C.6. The results of this lemma are standard and can be found in the convex optimization literature; see, for example, Boyd & Vandenberghe (2004) for detailed proofs. □

D ALL TECHNIQUE PROOFS

Proof of Theorem 3.1. Consider two neighboring data points z_t and z_t' for $t \ge 1$, differing in exactly one entry, i.e., $d_H(z_t, z_t') = 1$. Recall that

$$\mu_{t-1} = \nabla K(\boldsymbol{\theta}_{t-1}, \mathcal{D}) K(\mathcal{D}, \mathcal{D})^{-1} \mathcal{L}(\mathcal{D}, \boldsymbol{z}_t),$$

$$\tilde{\mu}_{t-1} = \nabla K(\boldsymbol{\theta}_{t-1}, \mathcal{D}) K(\mathcal{D}, \mathcal{D})^{-1} \mathcal{L}(\mathcal{D}, \boldsymbol{z}_t').$$

and

$$g_t = \boldsymbol{\mu}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\boldsymbol{\mu}_{t-1}\|} \right\}, \tilde{g}_t = \tilde{\boldsymbol{\mu}}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\tilde{\boldsymbol{\mu}}_{t-1}\|} \right\}.$$

It follows that the global sensitivity of the estimated gradient at time t is

$$||g_{t} - \tilde{g}_{t}|| = \left\| \boldsymbol{\mu}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\boldsymbol{\mu}_{t-1}\|} \right\} - \tilde{\boldsymbol{\mu}}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\tilde{\boldsymbol{\mu}}_{t-1}\|} \right\} \right\|$$

$$\leq \left(\left\| \boldsymbol{\mu}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\boldsymbol{\mu}_{t-1}\|} \right\} \right\| + \left\| \tilde{\boldsymbol{\mu}}_{t-1} \cdot \min \left\{ 1, \frac{B}{\|\tilde{\boldsymbol{\mu}}_{t-1}\|} \right\} \right\| \right)$$

$$\leq B + B = 2B.$$

Hence, by adding noise sampled from $\mathcal{N}\left(0,2(2B/\varepsilon_t)^2\log(1.25/\delta_t)\mathbf{I}_p\right)$ at each iteration, the gradient update is guaranteed to satisfy (ε_t,δ_t) -LDP. Moreover, by the parallel composition property of DP, the cumulative output $\tilde{\boldsymbol{\theta}}_t$ produced by Algorithm 1 satisfies $(\max\{\varepsilon_1,\ldots,\varepsilon_t\},\max\{\delta_1,\ldots,\delta_t\})$ -LDP.

Without loss of generality, we assume that the first iteration of Algorithm 1 satisfies $(\varepsilon_1, \delta_1)$ -LDP. Since the initial estimate $\hat{\theta}_0$ is deterministic, it follows directly that $\hat{\theta}_1$ also satisfies $(\varepsilon_1, \delta_1)$ -LDP. At the second iteration, $\hat{\theta}_2$, depends on both the privatized output $\hat{\theta}_1$ and the disjoint sample z_2 . It follows from Proposition A.4 that the two-fold composed algorithm $(\hat{\theta}_1, \hat{\theta}_2)$ satisfies $(\max\{\varepsilon_1, \varepsilon_2\}, \max\{\delta_1, \delta_2\})$ -LDP when the samples z_1 and z_2 are disjoint. Iteratively applying this argument, we conclude that after t iterations the entire sequence of updates satisfies $(\max\{\varepsilon_1, \ldots, \varepsilon_t\}, \max\{\delta_1, \ldots, \delta_t\})$ -LDP. By the post-processing property, both $\hat{\theta}_t$ and its averaged version $\tilde{\theta}_t$ inherit the same privacy guarantees.

Proof of Theorem 3.3. Our proof builds upon the framework of Koppel et al. (2021), which depends on the Hellinger distance, but here we adapt the analysis to the Sliced Wasserstein distance. Let $\rho_{\mathcal{D}_t}$ denote the posterior distribution at iteration t, where \mathcal{D}_t is a dictionary of size M_t . When a new sample θ_t is incorporated at iteration t+1, the dictionary is augmented to $\tilde{\mathcal{D}}_{t+1} = [\mathcal{D}_t; \theta_t]$, increasing its size to M_t+1 . The stopping criterion for Algorithm 2 is violated whenever

$$\min_{j=1,\dots,M_t+1} \eta_j \le \kappa. \tag{7}$$

Notice that (7) provides a lower bound on the approximation error η_{M_t+1} incurred by removing the newly added point θ_t . In particular, if $\eta_{M_t+1} \leq \kappa$, then the criterion in (7) is satisfied, and the model order remains unchanged. Consequently, η_{M_t+1} can serve as a proxy for η_j for all $j=1,\ldots,M_t+1$.

For the case of the Sliced Wasserstein distance between multivariate Gaussian distributions, the approximation error η_{M_t+1} depends only on the changes in the mean vector and covariance matrix induced by incorporating the new sample. θ_t . Specifically,

$$\eta_{M_t+1} \propto (\boldsymbol{\mu}_{t+1}|_{\mathcal{D}_t} - \boldsymbol{\mu}_{\mathcal{D}_t}, \, \boldsymbol{\Sigma}_{t+1}|_{\mathcal{D}_t} - \boldsymbol{\Sigma}_{\mathcal{D}_t}),$$

where $\mu_{t+1}|_{\mathcal{D}_t}$ and $\Sigma_{t+1}|_{\mathcal{D}_t}$ denote the mean and covariance conditioned on the dictionary \mathcal{D}_t , respectively, and $\mu_{\mathcal{D}_t}$, $\Sigma_{\mathcal{D}_t}$ are the corresponding quantities without θ_t .

Although there is no closed-form expression directly linking these mean and covariance differences to the Sliced Wasserstein distance, one can interpret the problem geometrically in terms of the Hilbert subspace defined by the current dictionary, $\mathcal{H}_{\mathcal{D}_t} := \mathrm{span}\{K(\mathcal{D}_j,\cdot)\}_{j=1}^{M_t}$. In particular, the approximation quality is governed by the distance between the kernel evaluation at the new point $K(\theta_t,\cdot)$ and the subspace $\mathcal{H}_{\mathcal{D}_t}$. Intuitively, if this distance is small, the new point contributes little additional information and can be safely excluded without degrading the fidelity of the surrogate model, thereby satisfying the compression criterion. The approximation quality is then determined by the distance from the kernel evaluation at the new point to the current dictionary's Hilbert subspace:

$$\operatorname{dist}(K(\boldsymbol{\theta}_t,\cdot),\mathcal{H}_{\mathcal{D}_t}) := \min_{\mathbf{v} \in \mathbb{R}^{M_t}} \left\| K(\boldsymbol{\theta}_t,\cdot) - \mathbf{v}^\top \boldsymbol{\nu}_{\mathcal{D}_t}(\cdot) \right\|_{\mathcal{H}},$$

where $\mathcal{H}_{\mathcal{D}_t} := \operatorname{span}\{K(\mathcal{D}_j,\cdot)\}_{j=1}^{M_t}$ denotes the subspace spanned by the kernel functions in the current dictionary.

Therefore, if there exists some constant c>0 such that $\mathrm{dist}(K(\theta_t,\cdot),\mathcal{H}_{\mathcal{D}_t})\leq c$, then there exists some $\kappa>0$ for which $\eta_{M_t+1}\leq \kappa$. This ensures that the approximation error remains sufficiently small, and hence the model order does not increase. Since θ lies in a compact set and K

is continuous, the range of the kernel embedding $\phi(\theta) := K(\theta, \cdot)$ is compact (Engel et al., 2004). Consequently, the number of balls of radius c required to cover $\phi(\theta)$ is finite and determined by the covering number of $\phi(\theta)$ at scale c (Anthony & Bartlett, 2009).

In particular, there exists a finite constant M^{∞} such that, if $M_t = M^{\infty}$, then $\operatorname{dist}(K(\boldsymbol{\theta}_t, \cdot), \mathcal{H}_{\mathcal{D}_t}) \leq c$, and consequently $\eta_{M_t+1} \leq \kappa$. Therefore, $M_t \leq M^{\infty}$ for all t. As shown by Engel et al. (2004), for a Lipschitz continuous Mercer kernel defined on a compact domain $\boldsymbol{\theta} \subset \mathbb{R}^p$, the covering number satisfies

$$M \le \mathcal{O}\left(\frac{1}{\kappa}\right)^p$$
.

We have completed the proof of this theorem.

Proof of Theorem 4.4. Recall that

$$\hat{\boldsymbol{\theta}}_t = \hat{\boldsymbol{\theta}}_{t-1} - \eta_t \big(g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) + \omega_t \big).$$

Define the shifted functions

$$\tilde{g}_{t-1}(\boldsymbol{\Delta}) = g_{t-1}(\boldsymbol{\Delta} + \boldsymbol{\theta}^*), \quad \tilde{g}(\boldsymbol{\Delta}) = g(\boldsymbol{\Delta} + \boldsymbol{\theta}^*), \quad \tilde{f}(\boldsymbol{\Delta}) = f(\boldsymbol{\Delta} + \boldsymbol{\theta}^*),$$

which correspond to a change of variables centered at the true parameter θ^* . We then have

$$\begin{split} \hat{\boldsymbol{\Delta}}_t &= \hat{\boldsymbol{\Delta}}_{t-1} - \eta_t g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) + \eta_t \omega_t \\ &= \hat{\boldsymbol{\Delta}}_{t-1} - \eta_t \nabla \tilde{f}(\hat{\boldsymbol{\Delta}}_{t-1}) + \eta_t \{ \nabla \tilde{f}(\hat{\boldsymbol{\Delta}}_{t-1}) - \tilde{g}(\hat{\boldsymbol{\Delta}}_{t-1}) \} \\ &+ \eta_t \{ \tilde{g}(\hat{\boldsymbol{\Delta}}_{t-1}) - \tilde{g}_{t-1}(\hat{\boldsymbol{\Delta}}_{t-1}) \} + \eta_t \omega_t \\ &= \hat{\boldsymbol{\Delta}}_{t-1} - \eta_t \nabla \tilde{f}(\hat{\boldsymbol{\Delta}}_{t-1}) + \eta_t \xi_{1t} + \eta_t \xi_{2t} + \eta_t \omega_t, \end{split}$$

where $\xi_{1t} = \nabla \tilde{f}(\hat{\Delta}_{t-1}) - \tilde{g}(\hat{\Delta}_{t-1}), \quad \xi_{2t} = \tilde{g}(\hat{\Delta}_{t-1}) - \tilde{g}_{t-1}(\hat{\Delta}_{t-1}).$

Therefore,

$$\|\hat{\Delta}_{t}\|_{2}^{2} = \|\hat{\Delta}_{t-1}\|_{2}^{2} - 2\eta_{t} \left\langle \hat{\Delta}_{t-1}, \nabla \tilde{f}(\hat{\Delta}_{t-1}) - \xi_{1t} - \xi_{2t} - \omega_{t} \right\rangle + \eta_{t}^{2} \left\| \nabla \tilde{f}(\hat{\Delta}_{t-1}) - \xi_{1t} - \xi_{2t} - \omega_{t} \right\|_{2}^{2}.$$
(8)

Notice that $E[\omega_t] = 0$, the expectation of gradient estimate $\nabla f(\hat{\theta}_{t-1})$ is $g(\hat{\theta}_{t-1})$, and $g_{t-1}(\hat{\theta}_{t-1}) - g(\hat{\theta}_{t-1})$ is a transformation of the martingale difference sequence $\nabla \mathcal{L}(\hat{\theta}_{t-1}, \mathbf{z}_t) - \nabla f(\hat{\theta}_{t-1})$. This implies that

$$E\left[\left\langle \hat{\Delta}_{t-1}, \xi_{1t} + \xi_{2t} + \omega_t \right\rangle\right] = 0.$$

Meanwhile, applying Lemma C.6(i) to the pair $(\theta^*, \hat{\theta}_{t-1})$, we obtain

$$\langle \nabla \tilde{f}(\hat{\boldsymbol{\Delta}}_{t-1}), \hat{\boldsymbol{\Delta}}_{t-1} \rangle \geq \tilde{f}(\hat{\boldsymbol{\Delta}}_{t-1}) + \frac{\lambda}{2} \|\hat{\boldsymbol{\Delta}}_{t-1}\|_2^2 \geq \frac{\lambda}{2} \|\hat{\boldsymbol{\Delta}}_{t-1}\|_2^2.$$

Using the upper equations above, we obtain

$$E\{2\eta_{t}\langle \hat{\Delta}_{t-1}, \nabla \tilde{f}(\hat{\Delta}_{t-1}) - \xi_{1t} - \xi_{2t} - \omega_{t} \rangle\} \ge \frac{\lambda}{2} \|\hat{\delta}_{t-1}\|_{2}^{2}.$$
(9)

Applying Lemma C.6(ii) to the pair $(\theta^*, \hat{\theta}_{t-1})$, we obtain the gradient norm bound $\|\nabla \tilde{f}(\hat{\Delta}_{t-1})\|_2 \leq \zeta \|\hat{\Delta}_{t-1}\|_2$. In addition, Lemma C.4 and Lemma C.5 jointly provide explicit upper bounds on the second moments of the stochastic error terms: $E(\|\xi_{1t}\|_2^2) \leq c_1(L + p\kappa)$ and $E(\|\xi_{2t}\|_2^2) \leq 2B^2$.

Using Young's inequality, we then have

$$E\{\|\nabla f(\hat{\Delta}_{t-1}) - \xi_{1t} - \xi_{2t} - \omega_t\|_2^2\}$$

$$\leq 4\|\nabla f(\hat{\Delta}_{t-1})\|_2^2 + 4E(\|\xi_{1t}\|_2^2) + 4E(\|\xi_{2t}\|_2^2) + 4E\|\omega_t\|_2^2$$

$$\leq 4\zeta^2\|\hat{\Delta}_{t-1}\|_2^2 + 8B^2 + 4c_1(L + p\kappa) + 32pB^2/\varepsilon^2\log(1.25/\delta).$$
(10)

Replacing the appropriate terms in (8) with (9) and (10), we have

$$E(\|\hat{\Delta}_t\|_2^2) \le (1 - \lambda \eta_t + c'\eta_t^2) \|\hat{\Delta}_{t-1}\|_2^2 + cp\eta_t^2 B^2 / \varepsilon^2 \log(1.25/\delta) + 4\eta_t^2 (c_1(L + p\kappa) + 2B^2).$$

Therefore, there exists some positive constant a_p depending on the dimension p such that

$$E(\|\hat{\boldsymbol{\Delta}}_t\|_2^2) \le (1 - \lambda \eta_t + a_p^2 \eta_t^2) \|\hat{\boldsymbol{\Delta}}_{t-1}\|_2^2 + a_p \eta_t^2 B^2 / \varepsilon^2 \log(1.25/\delta) + 4\eta_t^2 (c_1(L + p\kappa) + 2B^2),$$

Define $t_0 = \min\{t : \lambda \ge 2a_p^2\eta_t, \lambda\eta_t t \ge 8\alpha \log t\}$. Then, for any $t \ge t_0$ and some constant $b_p = O(a_p)$, the equation simplifies to

$$E(\|\hat{\Delta}_t\|_2^2) \le (1 - \lambda \eta_t/2) \|\hat{\Delta}_{t-1}\|_2^2 + b_p \eta_t^2 B^2 / \varepsilon^2 \log(1.25/\delta) + 4\eta_t^2 (c_1(L + p\kappa) + 2B^2),$$

Note that $\exp(-t\lambda\eta_t/4) \le \exp(-\lambda\eta t^{1-\alpha}/4) \le t^{-2\alpha} \le t^{-\alpha}$ for $t \ge 2t_0$. Therefore, using the same arguments as in Chen et al. (2020), for $t \ge 2t_0$, we have

$$E(\|\hat{\Delta}_{t}\|_{2}^{2}) \leq \exp(-t\lambda\eta_{t}/4)E\|\hat{\Delta}_{t/2}\|_{2}^{2} + 2b_{p}\eta_{t/2}B^{2}\log(1.25/\delta)/(\lambda\varepsilon^{2}) + 8\eta_{t/2}^{2}(c_{1}(L+p\kappa) + 2B^{2})$$

$$\leq \exp(-t\lambda\eta_{t}/4)(E\|\hat{\Delta}_{n_{0}}\|_{2}^{2} + 2b_{p}\eta_{n_{0}}B^{2}\log(1.25/\delta)/(\lambda\varepsilon^{2})$$

$$+ 8\eta_{n_{0}}(c_{1}(L+p\kappa) + 2B^{2})/\lambda) + 2b_{p}\eta(t/2)^{-\alpha}B^{2}\log(1.25/\delta)/(\lambda\varepsilon^{2})$$

$$+ 8\eta(t/2)^{-\alpha}(c_{1}(L+p\kappa) + 2B^{2})/\lambda$$

$$\leq \exp(-t\lambda\eta_{t}/4)\{c(1+\|\hat{\Delta}_{0}\|_{2}^{2}) + 2b_{p}\eta_{n_{0}}B^{2}\log(1.25/\delta)/(\lambda\varepsilon^{2})$$

$$+ 8\eta_{n_{0}}(c_{1}(L+p\kappa) + 2B^{2})/\lambda\} + 2b_{p}\eta(t/2)^{-\alpha}B^{2}\log(1.25/\delta)/(\lambda\varepsilon^{2})$$

$$+ 8\eta(t/2)^{-\alpha}(c_{1}(L+p\kappa) + 2B^{2})/\lambda$$

$$\leq c't^{-\alpha}\{\|\hat{\Delta}_{0}\|_{2}^{2} + c''b_{p}\eta B^{2}\log(1.25/\Delta)/(\lambda\varepsilon^{2}) + \eta(L+p\kappa + 2B^{2})/\lambda\}.$$

Proof of Theorem 4.5. Recall that $\hat{\boldsymbol{\theta}}_t = \hat{\boldsymbol{\theta}}_{t-1} - \eta_t(g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) + \omega_t)$. By Assumption 4.3, we have $f(\hat{\boldsymbol{\theta}}_t) \leq f(\hat{\boldsymbol{\theta}}_{t-1}) + \langle \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \hat{\boldsymbol{\theta}}_t - \hat{\boldsymbol{\theta}}_{t-1} \rangle + \frac{\zeta}{2} \|\hat{\boldsymbol{\theta}}_t - \hat{\boldsymbol{\theta}}_{t-1}\|^2.$

Thus, substituting the step sizes, we obtain

$$\begin{split} f(\hat{\theta}_{t}) &\leq f(\hat{\theta}_{t-1}) - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), g_{t-1}(\hat{\theta}_{t-1}) + \omega_{t} \rangle + \frac{\zeta \eta_{t}^{2}}{2} \|g_{t-1}(\hat{\theta}_{t-1}) + \omega_{t}\|^{2} \\ &= f(\hat{\theta}_{t-1}) - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), g_{t-1}(\hat{\theta}_{t-1}) \rangle - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), \omega_{t} \rangle \\ &+ \frac{\zeta \eta_{t}^{2}}{2} \left(\|g_{t-1}(\hat{\theta}_{t-1})\|^{2} + \|\omega_{t}\|^{2} + 2\langle g_{t-1}(\hat{\theta}_{t-1}), \omega_{t} \rangle \right) \\ &\leq f(\hat{\theta}_{t-1}) - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), g_{t-1}(\hat{\theta}_{t-1}) - \nabla f(\hat{\theta}_{t-1}) + \nabla f(\hat{\theta}_{t-1}) \rangle - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), \omega_{t} \rangle \\ &+ \frac{\zeta \eta_{t}^{2}}{2} \left(\|g_{t-1}(\hat{\theta}_{t-1})\|^{2} + 8pB^{2}/\varepsilon^{2} \log(1.25/\delta) + 2\langle g_{t-1}(\hat{\theta}_{t-1}), \omega_{t} \rangle \right) \\ &\leq f(\hat{\theta}_{t-1}) - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), g_{t-1}(\hat{\theta}_{t-1}) - \nabla f(\hat{\theta}_{t-1}) \rangle - \eta_{t} \|\nabla f(\hat{\theta}_{t-1})\|^{2} - \eta_{t} \langle \nabla f(\hat{\theta}_{t-1}), \omega_{t} \rangle \\ &+ \frac{\zeta \eta_{t}^{2}}{2} \left(\|\nabla f(\hat{\theta}_{t-1})\|^{2} + \|g_{t-1}(\hat{\theta}_{t-1}) - \nabla f(\hat{\theta}_{t-1})\|^{2} + 2\langle \nabla f(\hat{\theta}_{t-1}), g_{t-1}(\hat{\theta}_{t-1}) - \nabla f(\hat{\theta}_{t-1}) \rangle \right) \\ &+ \frac{\zeta \eta_{t}^{2}}{2} \left(8pB^{2}/\varepsilon^{2} \log(1.25/\delta) + 2\langle g_{t-1}(\hat{\theta}_{t-1}), \omega_{t} \rangle \right) \\ &\leq f(\hat{\theta}_{t-1}) - \frac{\eta_{t}}{2} \|\nabla f(\hat{\theta}_{t-1})\|^{2} + \eta_{t} \langle \nabla g_{t-1}(\hat{\theta}_{t-1}) - f(\hat{\theta}_{t-1}), \omega_{t} \rangle \\ &+ \frac{\zeta \eta_{t}^{2}}{2} \left(\|g_{t-1}(\hat{\theta}_{t-1}) - \nabla f(\hat{\theta}_{t-1})\|^{2} + 8pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right), \end{split}$$

where the first inequality follows from ζ -smoothness and the last inequality holds due to $\eta_t \leq \frac{1}{\zeta}$. The result is obtained by rearranging terms.

$$\frac{\eta_t}{2} \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^2 \le f(\hat{\boldsymbol{\theta}}_{t-1}) - f(\hat{\boldsymbol{\theta}}_t) + \eta_t \langle g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \omega_t \rangle
+ \frac{\zeta \eta_t^2}{2} \left(\|g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^2 + 8pB^2 / \varepsilon^2 \log(1.25/\delta) \right).$$

Summing the inequalities over t = 1, ..., T, we have

$$\sum_{t=1}^{T} \eta_{t} \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2} \leq 2(f(\hat{\boldsymbol{\theta}}_{0}) - f(\hat{\boldsymbol{\theta}}_{T-1})) + \sum_{t=1}^{T} 2\eta_{t} \langle g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \omega_{t} \rangle
+ \sum_{t=1}^{T} 2\zeta \eta_{t}^{2} \left(\|g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2} + 16pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right)
\leq 2(f(\hat{\boldsymbol{\theta}}_{0}) - f(\boldsymbol{\theta}^{*})) + \sum_{t=1}^{T} 2\eta_{t} \langle g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \omega_{t} \rangle
+ \sum_{t=1}^{T} 2\zeta \eta_{t}^{2} \left(\|g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2} + 16pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right).$$
(11)

Dividing both sides by $\sum_{t=1}^{T} \eta_t$ yields

$$\begin{split} \frac{\sum_{t=1}^{T} \eta_{t} \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2}}{\sum_{t=1}^{T} \eta_{t}} &\leq \frac{2(f(\hat{\boldsymbol{\theta}}_{0}) - f(\boldsymbol{\theta}^{\star}))}{\sum_{t=1}^{T} \eta_{t}} + \frac{\sum_{t=1}^{T} 2\eta_{t} \langle g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \omega_{t} \rangle}{\sum_{t=1}^{T} \eta_{t}} \\ &+ \frac{\sum_{t=1}^{T} 2\zeta \eta_{t}^{2} \left(\|g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2} + 16pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right)}{\sum_{t=1}^{T} \eta_{t}} \end{split}$$

Note that $E(\omega_t) = 0$, the expectation of gradient estimate $\nabla f(\hat{\theta}_{t-1})$ is $g(\hat{\theta}_{t-1})$, and $g_{t-1}(\hat{\theta}_{t-1}) - g(\hat{\theta}_{t-1})$ is a transformation of the martingale difference sequence $\nabla \mathcal{L}(\hat{\theta}_{t-1}, \mathbf{z}_t) - \nabla f(\hat{\theta}_{t-1})$, implying

$$E(\langle g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1}), \omega_t \rangle) = 0.$$

Furthermore,

$$||g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})||^2 \le ||g_{t-1}(\hat{\boldsymbol{\theta}}_{t-1}) - g(\hat{\boldsymbol{\theta}}_{t-1})||^2 + ||g(\hat{\boldsymbol{\theta}}_{t-1}) - \nabla f(\hat{\boldsymbol{\theta}}_{t-1})||^2$$

$$\le 2B^2 + c_1(L + p\kappa).$$

Taking the expectation with respect to these terms and substituting into (11), we obtain

$$\frac{\sum_{t=1}^{T} \eta_{t} E \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2}}{\sum_{t=1}^{T} \eta_{t}} \leq \frac{2(f(\hat{\boldsymbol{\theta}}_{0}) - f(\boldsymbol{\theta}^{\star}))}{\sum_{t=1}^{T} \eta_{t}} + \frac{\sum_{t=1}^{T} 2\zeta \eta_{t}^{2} \left((c_{1}(L + p\kappa) + 2B^{2}) + 16pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right)}{\sum_{t=1}^{T} \eta_{t}}$$

We then obtain

$$\min_{1 \le t \le T} E \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^{2} \le \frac{2(f(\hat{\boldsymbol{\theta}}_{0}) - f(\boldsymbol{\theta}^{*}))}{\sum_{t=1}^{T} \eta_{t}} + \frac{\sum_{t=1}^{T} 2\zeta \eta_{t}^{2} \left((c_{1}(L + p\kappa) + 2B^{2}) + 16pB^{2}/\varepsilon^{2} \log(1.25/\delta) \right)}{\sum_{t=1}^{T} \eta_{t}}.$$

Recall that $\eta_t = \eta_0 t^{-\alpha}$. Following the integral bounding technique in Garrigos & Gower (2023), there exist constants c_2 and c_3 such that $\sum_{t=1}^T \eta_t = \eta_0 \sum_{t=1}^T t^{-\alpha} \le c_2 T^{1-\alpha}$ and $\sum_{t=1}^T \eta_t^2 = \eta_0 \sum_{t=1}^T t^{-2\alpha} \le c_3$. Therefore, the inequality simplifies to

$$\min_{1 \le t \le T} E \|\nabla f(\hat{\boldsymbol{\theta}}_{t-1})\|^2 \le c' \frac{(f(\hat{\boldsymbol{\theta}}_0) - f(\boldsymbol{\theta}^{\star})) + \zeta((L + p\kappa) + B^2) + pB^2/\varepsilon^2 \log(1.25/\delta)}{T^{1-\alpha}}.$$

E ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide details of data generating processes and additional results in Section 5.

Example 5.1 (Continued). We evaluate the proposed algorithm and the competing methods under linear, logistic and ReLU regression models, respectively.

Linear regression. We sample T=20000 i.i.d. data points $\{(\boldsymbol{x}_t,y_t)\}_{t=1}^T$, where the covariates are drawn as $\boldsymbol{x}_t \sim N(0,\mathbf{I}_p)$, and the responses are generated according to

$$y_t = \boldsymbol{x}_t^{\top} \boldsymbol{\theta} + \varepsilon_t,$$

with true parameter vector $\boldsymbol{\theta} = \mathbf{1}_p$ and noise terms $\varepsilon_t \overset{\text{i.i.d.}}{\sim} N(0,1)$. We employ the Huber loss function ρ_c with threshold c=1, and incorporate gradient sensitivity control to ensure stability. The overall objective function is given by

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{T} \sum_{t=1}^{T} \rho_c \left(y_t - \boldsymbol{x}_t^{\top} \boldsymbol{\theta} \right) \cdot \min \left(1, \frac{2}{\|\boldsymbol{x}_t\|^2} \right).$$

This reweighting scheme effectively bounds the influence of high-magnitude gradients, serving as a form of implicit gradient clipping that enhances robustness during optimization.

Logistic regression. The feature vectors $\mathbf{x}_t \in \mathbb{R}^d$ are sampled independently from a standard normal distribution, $\mathbf{x}_t \sim N(0, \mathbf{I}_p)$. Binary labels $y_t \in \{-1, +1\}$ are generated according to the logistic model:

$$\mathbb{P}(y_t = 1 \mid \boldsymbol{x}_t) = \frac{1}{1 + \exp(-\boldsymbol{x}_t^{\top} \boldsymbol{\theta})},$$

where the true parameter vector $\theta = \mathbf{1}_p$ defines the underlying decision boundary. The learning objective is defined via the binary cross-entropy loss, which measures the discrepancy between the predicted probabilities and the true labels. Specifically, we minimize the following empirical risk:

$$\mathcal{L}(\boldsymbol{\theta}) = -\frac{1}{T} \sum_{t=1}^{T} \left[y_t \log(p_t) + (1 - y_t) \log(1 - p_t) \right] \cdot \min\left(1, \frac{2}{\|\boldsymbol{x}_t\|^2}\right),$$

where, $p_t = \mathbb{P}(y_t = 1 \mid x_t)$ represents the predicted probability of the positive class for sample t, given by the sigmoid function applied to the linear combination of features and parameters.

ReLU regression. We generate synthetic data $\{(x_t, y_t)\}_{t=1}^T$ according to the model:

$$y_t = \text{ReLU}(\boldsymbol{x}_t^{\top} \boldsymbol{\theta}),$$

with true parameter vector $\theta = \mathbf{1}_p$. The objective is to minimize the squared loss, which quantifies the discrepancy between the predicted values and the true responses. The empirical risk is thus defined as:

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{T} \sum_{t=1}^{T} \rho_c \left(y_t - \text{ReLU}(\boldsymbol{x}_t^{\top} \boldsymbol{\theta}) \right) \cdot \min \left(1, \frac{2}{\|\boldsymbol{x}_t\|^2} \right).$$

This setup allows us to evaluate how effectively each method can handle nonlinear transformations and non-continuous derivative functions, as introduced by the ReLU activation. By applying this nonlinearity, we test the robustness of various algorithms in approximating complex, discontinuous mappings while maintaining low prediction error.

Figure 5 presents additional results for p=5. The first three columns of Figure 5 illustrate the trajectory of the first-dimensional coefficient estimate (true value = 1) across iterations in the p=5 setting. For the linear model, both LDP-BO and LDP-SGD closely track their non-private counterparts. In nonlinear models (logistic and ReLU), however, BO-based methods consistently outperform SGD-based approaches under all privacy regimes. The last column of Figure 5 reports MSE of the parameter estimates, revealing that LDP-BO achieves consistently lower error and reduced variability compared to LDP-SGD in complex settings. Even under strong privacy constraints ($\varepsilon=1$), LDP-BO exhibits faster convergence and attains accuracy on par with non-private BO and SGD. These results underscore the modeling advantage of LDP-BO in handling nonlinear problems in

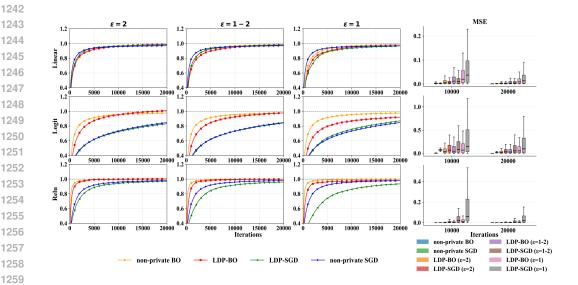


Figure 5: Left figure represents evolution of the first-dimension coefficient estimate (true value = 1) over iterations for linear, logistic, and ReLU models (rows) in Example 5.1. Columns correspond to privacy budgets $\varepsilon = 2$, $\varepsilon \sim U(1,2)$, and $\varepsilon = 1$. Right figure represents boxplots of coefficient MSEs across three models under different privacy budgets in Example 5.1.

moderate-dimensional (p = 5) scenarios, where it effectively mitigates the utility degradation often associated with gradient-based private optimization.

The compression budget strikes a balance between prediction time and prediction accuracy. A smaller compression budget retains more essential information, leading to improved results at the cost of increased computational time. Figure 6 further illustrates the impact of different compression budgets (0.1 and 0.2) on the performance of linear, logistic, and ReLU regression models under varying privacy budgets ($\varepsilon = 2$, $\varepsilon = U(1, 2)$, and $\varepsilon = 1$). Across all settings, a smaller compression budget (0.1, represented by red lines) consistently leads to better performance compared to a larger budget (0.2, represented by blue lines), as evidenced by faster convergence and higher final accuracy. This improvement is particularly pronounced in complex models such as logistic and ReLU regression, where the underlying data structure is more nonlinear and intricate. In these cases, a smaller compression budget helps preserve a greater amount of critical kernel information during the Bayesian optimization process, which is essential for accurately modeling complex decision boundaries. Therefore, tighter compression—achieved through a smaller budget—is especially beneficial in complex models, as it enables the algorithm to retain more informative data points, leading to more reliable and accurate parameter estimates. The results suggest that carefully controlling the compression budget is crucial for balancing efficiency and utility, with more complex problems generally requiring stricter (i.e., smaller) compression budgets to achieve optimal performance.

Example 5.2 (Continued). In this example, we perform LDP-BO with $(\epsilon, \delta) = (1, 0.2), \kappa = 0.1$ and B=2. The following is a detailed description of the models, including the Sine function and the Friedman function.

Sine function. We apply an exact Gaussian process regression model designed under privacy constraints. The model employs a constant mean function m(x) = 0 and a scaled radial basis function (RBF) covariance kernel:

$$K(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{ ext{output}}^2 \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2\ell^2}\right),$$

The kernel contains two trainable parameters: the length scale ℓ , which controls the smoothness of the function, and the output scale $\sigma_{
m output}$, which modulates the amplitude of the output. The model

1299

1301

1309

1315

1316

1317

1318 1319 1320

1321

1322 1323

1324 1325 1326

1327

1330

1331

1332

1333

1334

1335

1336 1337

1338

1339

1340 1341

1344

1347

1348

1349

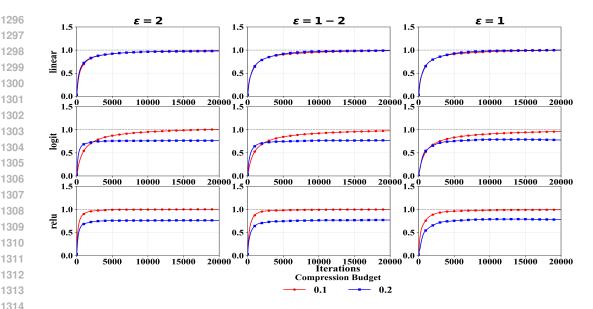


Figure 6: Results of experiments with different compression budget, where dimension p=5, and privacy budget $\delta = 0.2$. Each row corresponds to a different model: linear regression, logistic regression, and ReLU regression. Each column represents a different privacy budget $\varepsilon = 2, Unif(1,2), 1$, ordered from highest to lowest noise intensity.

is trained by minimizing the negative log marginal likelihood (NLL), which serves as our objective function:

$$\mathcal{L}(\boldsymbol{\theta}) = -\log p(y \mid \boldsymbol{x}, \boldsymbol{\theta}) = \frac{1}{2} y^{\top} K_y^{-1} y + \frac{1}{2} \log |K_y| + \frac{1}{2} \log(2\pi),$$

where $K_y = K + \sigma_{\text{noise}}^2 \mathbf{I}$ denotes the noise-perturbed covariance matrix. This loss function naturally balances data fit (first term) and model complexity (second term), providing a probabilistically principled measure of model adequacy. We set $\sigma_{\text{noise}}^2 = 10^{-4}$.

We optimize the parameters in log space to ensure positivity and improve numerical stability. The trainable parameter vector is thus $\theta = (\log \ell, \log \sigma_{\text{output}})$, making this a two-dimensional optimization problem. The actual kernel parameters are recovered via exponentiation: $\ell = \exp(\log \ell)$, $\sigma_{\rm output} = \exp(\log \sigma_{\rm output})$. This formulation enables efficient Bayesian optimization of the kernel parameters while providing a tractable and interpretable objective for privacy-preserving parameter optimization. The entire framework offers a rigorous foundation for adaptive, nonparametric regression under DP constraints.

Friedman function. We propose an adaptive Gaussian process GP regression framework employing automatic relevance determination (ARD) to handle multidimensional input spaces in sequential learning scenarios. The model utilizes a constant mean function and a scaled radial basis function (RBF) covariance kernel with ARD:

$$K(\boldsymbol{x}, \boldsymbol{x}') = \sigma_{ ext{output}}^2 \exp\left(-rac{1}{2} \sum_{j=1}^p rac{(x_j - x_j')^2}{\ell_j^2}
ight),$$

where each input dimension p has its own trainable length scale ℓ_i , allowing the model to automatically learn the relevance of each feature. The output scale σ_{output} can be either optimized or fixed to modulate function amplitude. In our simulations, we fixed it to 1.

The training objective minimizes the negative log marginal likelihood:

 $\mathcal{L}(\boldsymbol{\theta}) = -\log p(y \mid \boldsymbol{x}, \boldsymbol{\theta}) = \frac{1}{2} y^{\top} K_y^{-1} y + \frac{1}{2} \log |K_y| + \frac{1}{2} \log(2\pi),$

where $\theta = (\log \ell_1, \log \ell_2, \dots, \log \ell_p)$ represents the *p*-dimensional hyperparameter vector optimized in log space to ensure positivity and numerical stability. The ARD formulation enables automatic feature selection by assigning larger length scales to less relevant dimensions, effectively suppressing their contribution to the covariance function.

This approach provides a principled probabilistic framework for high-dimensional regression, with the optimization complexity scaling linearly with the input dimension p. The model maintains computational tractability through exact inference while offering interpretable insights into feature relevance through the learned length scales, making it particularly suitable for Bayesian optimization in parameterized spaces.

Example 5.3 (Continued). The data preprocessing pipeline starts with comprehensive cleaning to enhance data robustness. We remove records with invalid fare amounts, such as negative values or extreme outliers beyond predefined percentile thresholds, and handle missing values in key fields. Following this, feature engineering is conducted to extract meaningful signals from the raw data.

Original features such as passenger_count are retained to account for the impact of group travel on fare pricing. Spatial information is derived from the provided geographic coordinates: pickup_longitude and pickup_latitude (indicating where the trip began), along with dropoff_longitude and dropoff_latitude (marking the destination). From these, we compute the Manhattan distance between pickup and drop-off points—a more accurate proxy for actual travel distance in New York City's grid-like street layout than Euclidean distance.

Temporal patterns are captured by extracting features from the pickup_patetime field, including the hour of the day and day of the week, which help model variations in demand, traffic congestion, and surge pricing dynamics.

The final feature set combines cleaned original variables with these engineered spatial and temporal features, forming the input for downstream regression models designed to accurately predict fare amounts.

F THE USE OF LARGE LANGUAGE MODELS

In the preparation of this manuscript, we employed a large language model (LLM) to assist in the polishing and refinement of the writing. The model was used exclusively for improving linguistic expression, enhancing clarity, and ensuring consistency of terminology—tasks that contribute to the overall readability and academic tone of the document. All technical content, mathematical reasoning, and scientific conclusions remain entirely formulated by the authors. The use of LLM-assisted editing did not alter the theoretical contributions or empirical results presented in this work.