# Gaussian Process Optimization with Adaptive Sketching: Scalable and No Regret

**Daniele Calandriello**                                      DANIELE.CALANDRIELLO@IIT.IT
*LCSL - Istituto Italiano di Tecnologia, Genova, Italy & MIT, Cambridge, USA*

**Luigi Carratino**                                          LUIGI.CARRATINO@DIBRIS.UNIGE.IT
*UNIGE - Università degli Studi di Genova, Genova, Italy*

**Alessandro Lazaric**                                       LAZARIC@FB.COM
*FAIR - Facebook AI Research, Paris, France*

**Michal Valko**                                             MICHAL.VALKO@INRIA.FR
*INRIA Lille - Nord Europe, SequeL team, Lille, France*

**Lorenzo Rosasco**                                          LROSASCO@MIT.EDU
*LCSL - Istituto Italiano di Tecnologia, Genova, Italy & MIT, Cambridge, USA*
*UNIGE - Università degli Studi di Genova, Genova, Italy*

**Editors:** Alina Beygelzimer and Daniel Hsu

## Abstract

Gaussian processes (GP) are a stochastic processes, used as Bayesian approach for the optimization of black-box functions. Despite their effectiveness in simple problems, GP-based algorithms hardly scale to high-dimensional functions, as their per-iteration time and space cost is at least *quadratic* in the number of dimensions $d$ and iterations $t$. Given a set of $A$ alternatives to choose from, the overall runtime $\mathcal{O}(t^3 A)$ is prohibitive. In this paper, we introduce BKB (*budgeted kernelized bandit*), a new approximate GP algorithm for optimization under bandit feedback that achieves near-optimal regret (and hence near-optimal convergence rate) with near-constant per-iteration complexity and remarkably no assumption on the input space or covariance of the GP.

We combine a kernelized linear bandit algorithm (GP-UCB) leverage score sampling as a randomized matrix sketching and prove that selecting inducing points based on their posterior variance gives an accurate low-rank approximation of the GP, preserving variance estimates and confidence intervals. As a consequence, BKB does not suffer from *variance starvation*, an important problem faced by many previous sparse GP approximations. Moreover, we show that our procedure selects at most $\widetilde{\mathcal{O}}(d_{\text{eff}})$ points, where $d_{\text{eff}}$ is the *effective* dimension of the explored space, which is typically much smaller than both $d$ and $t$. This greatly reduces the dimensionality of the problem, thus leading to a $\mathcal{O}(T A d_{\text{eff}}^2)$ runtime and $\mathcal{O}(A d_{\text{eff}})$ space complexity.

**Keywords:** *sparse Gaussian process optimization*; *kernelized linear bandits*; *regret*; *sketching*; *Bayesian optimization*; *black-box optimization*; *variance starvation*

## 1. Introduction

Efficiently selecting the best alternative out of a set of alternatives is important in sequential decision making, with practical applications ranging from recommender systems (Li et al., 2010) to experimental design (Robbins, 1952). It is also the main focus of the research in bandits (Lattimore and Szepesvári, 2019) and Bayesian optimization (Mockus, 1989; Pelikan, 2005; Snoek et al., 2012),

---

Extended abstract. Full version appears as https://arxiv.org/abs/1903.05594.

that study optimization under bandit feedback. In this setting, a learning algorithm sequentially interacts with a reward or utility function $f$. Over $T$ interactions, the algorithm chooses a point $\mathbf{x}_t$ and it has only access to a noisy black-box evaluation of $f$ at $\mathbf{x}_t$. The goal of the algorithm is to minimize the cumulative regret, which compares the reward accumulated at the points selected over time, $\sum_t f(\mathbf{x}_t)$, to the reward obtained by repeatedly selecting the optimum of the function, i.e., $T \max_x f(x)$. In this paper take the Gaussian process optimization approach. In particular, we study the GP-UCB algorithm first introduced by Srinivas et al. (2010).

Starting from a Gaussian process prior over $f$, GP-UCB alternates between evaluating the function,x and using the evaluations to build a posterior of $f$. This posterior is composed by a mean function $\mu$ that estimates the value of $f$, and a variance function $\sigma$ that captures the uncertainty $\mu$. These two quantities are combined in a single upper confidence bound (UCB) that drives the selection of the evaluation points, and trades off between evaluating high-reward points (*exploitation*) and testing possibly sub-optimal points to reduce the uncertainty on the function (*exploration*). The performance of GP-UCB has been studied by Srinivas et al. (2010); Valko et al. (2013); Chowdhury and Gopalan (2017) to show that GP-UCB iprovably achieves low regret both in a Bayesian and non-Bayesian setting. However, the main limiting factor to its applicability is its computational cost. When choosing between $A$ alternatives, GP-UCB requires $\Omega(At^2)$ time/space to select each new point, and therefore does not scale to complex functions. Several approximations of GP-UCB have been suggested (Quinonero-Candela et al., 2007; Liu et al., 2018) and we review them next:

**Inducing points:** The GP can be restricted to lie in the range of a small subset of inducing points. The subset should cover the space well for accuracy, but also be as small as possible for efficiency. Methods referred to as *sparse GPs*, have been proposed to select the inducing points and an approximation based on the subset. Popular instances of this approach are the subset of regressors (SoR, Wahba, 1990) and the deterministic training conditional (DTC, Seeger et al., 2003). While these methods are simple to interpret and efficient, they do not come with regret guarantees. Moreover, when the subset does not cover the space well, they suffer from *variance starvation* (Wang et al., 2018), as they underestimate the variance of points far away from the inducing points.

**Random Fourier features:** Another approach is to use explicit feature expansions to approximate the GP covariance function, and embed the points in a low-dimensional space, usually exploiting some variation of Fourier expansions (Rahimi and Recht, 2007). Among these methods, Mutný and Krause (2018) recently showed that discretizing the posterior on a fine grid of quadrature Fourier features (QFF) incurs a negligible approximation error. It is sufficient to prove that the maximum of the approximate posterior can be efficiently found and that it is accurate enough that Thompson sampling with quadratures provably achieves low regret. However this approach does not extend to non-stationary (or translation invariant) kernels and although its dependence on $t$ is small, the approximation and posterior maximization procedure scales exponentially with the input dimension.

**Variational inference:** The next approach replaces the true GP likelihood with a variational approximation that can be optimized efficiently. Although recent methods provide guarantees on the approximate posterior mean and variance (Huggins et al., 2019), these guarantees only apply to GP regression and not to the harder optimization setting.

**Linear case:** In the linear bandits, there are methods that reduce the complexity of algorithms such as LinUCB have. Kuzborskij et al. (2019) uses the *frequent directions* (FD, Ghashami et al., 2016) to project the design matrix data to a smaller subspace. Unfortunately, the size of the subspace has to be specified in advance and when the size is not sufficiently large, the method suffers linear regret. Prior

to the FD approach, CBRAP (Yu et al., 2017) used random projections instead, but faced similar issues. This turns out to be a fundamental weakness of all approaches that do not adapt the *actual* size of the space defined by the sequence of points selected by the learning algorithm. Indeed, Ghosh et al. (2017) showed a lower bound for the linear case that shows that as soon as one single arm does not abide by the linear model we can suffer linear regret.

### 1.1. Contributions

In this paper, we show a way to adapt the size of the projected space online and devise the BKB (budgeted kernel bandit) algorithm achieving near-optimal regret with a computational complexity drastically smaller than GP-UCB. This is achieved *without assumptions on the complexity* of the input or on the kernel function. BKB leverages several well-known tools: a DTC approximation of the posterior variance, based on inducing points, and a confidence interval construction based on state-of-the-art self-normalized concentration inequalities (Abbasi-Yadkori et al., 2011). It also introduces two novel tools: a selection strategy to select inducing points based on ridge leverage score (RLS) sampling (Alaoui and Mahoney, 2015) that is provably accurate, and an approximate confidence interval that is not only nearly as accurate as the one of GP-UCB, but also efficient.

Ridge leverage score sampling was introduced for randomized kernel matrix approximation (Alaoui and Mahoney, 2015). In the context of GPs, RLS correspond to the posterior variance of a point, which allows adapting algorithms and their guarantees from the RLS sampling to the GP setting. This solves two problems in sparse GPs and linear bandit approximations. First, BKB constructs estimates of the variance that are provably accurate, i.e., it does not suffer from variance starvation, which results in provably accurate confidence bounds as well. The only method with comparable guarantees, Thompson sampling with quadrature FF (Mutný and Krause, 2018), only applies to stationary kernels, and only in extremely low dimension. Moreover our approximation guarantees are qualitatively different since they do not require a corresponding *uniform* approximation bound on the GP. Second, BKB adaptively chooses the size of the inducing point set based on the *effective dimension* $d_{\text{eff}}$ of the problem. This is crucial to achieve low regret, since fixed approximation schemes may suffer linear regret. Moreover, in a problem with $A$ arms, using a set of $\mathcal{O}(d_{\text{eff}})$ inducing points results in an algorithm with $\mathcal{O}(Ad_{\text{eff}}^2)$ per-step runtime and $\mathcal{O}(Ad_{\text{eff}})$ space, a significant improvement over the $\mathcal{O}(At^2)$ time and $\mathcal{O}(At)$ space cost of GP-UCB.

Finally, while in our work we only address kernelized (GP) bandits, our work could be extended to more complex online learning problems, such as to recent advances in kernelized reinforcement learning (Chowdhury and Gopalan, 2019). Moreover, inducing point methods have clear interpretability and our analysis provides insight both from a bandit and Bayesian optimization perspective, making it applicable to a large amount of downstream tasks.

## 2. Background

**Notation.** We use lower-case letters $a$ for scalars, lower-case bold letters $\mathbf{a}$ for vectors, and upper-case bold letters $\mathbf{A}$ for matrices and operators, where $[\mathbf{A}]_{ij}$ denotes its element $(i, j)$. We denote by $\|\mathbf{x}\|_{\mathbf{A}}^2 \triangleq \mathbf{x}^\intercal \mathbf{A}\mathbf{x}$, the norm with metric $\mathbf{A}$, and $\|\mathbf{x}\| \triangleq \|\mathbf{x}\|_{\mathbf{I}}$ with $\mathbf{I}$ being the identity. Finally, we denote the first $T$ integers as $[T] \triangleq \{1, \ldots, T\}$.

**Online optimization under bandit feedback.** Let $f : \mathcal{A} \to \mathbb{R}$ be a reward function that we wish to optimize over a set of decisions $\mathcal{A}$, also called actions or arms. For simplicity, we assume that $\mathcal{A} \triangleq \{\mathbf{x}_i\}_{i=1}^A$ is a fixed finite set of $A$ vectors in $\mathbb{R}^d$. We discuss how to relax these assumptions in

Section 5. In optimization under bandit feedback, a learner aims to optimize $f$ through a sequence of interactions. At each step $t \in [T]$, the learner **(1)** chooses an arm $\mathbf{x}_t \in \mathcal{A}$, **(2)** receives reward $y_t \triangleq f(\mathbf{x}_t) + \eta_t$, where $\eta_t$ is a zero-mean noise, **(3)** updates its model of the problem.

The goal of the learner is to minimize its cumulative regret $R_T \triangleq \sum_{t=1}^{T} f(\mathbf{x}_\star) - f(\mathbf{x}_t)$ w.r.t. the best[1] $\mathbf{x}_\star$, where $\mathbf{x}_\star \triangleq \arg\max_{\mathbf{x}_i \in \mathcal{A}} f(\mathbf{x}_i)$. In particular, the objective of a *no-regret* algorithm is to have $R_T/T$ go to zero as $T$ grows as fast as possible. Recall that the regret is related to the convergence rate and the optimization performance. In fact, let $\overline{\mathbf{x}}_T$ be an arm chosen at random from the sequence of arms $(\mathbf{x}_1, \ldots, \mathbf{x}_T)$ selected by the learner, then $f(\mathbf{x}_\star) - \mathbb{E}[f(\overline{\mathbf{x}}_T)] \leq R_T/T$.

### Gaussian process optimization and GP-UCB

GP-UCB is popular no-regret algorithm for optimization under bandit feedback and was introduced by Srinivas et al. (2010) for Gaussian process optimization. We first give the formal definition of a Gaussian process (Rasmussen and Williams, 2006), and then briefly present GP-UCB.

A *Gaussian process* $\mathrm{GP}(\mu, k)$ is a generalization of the Gaussian distribution to a space of functions and it is defined by a mean function $\mu(\cdot)$ and covariance function $k(\cdot, \cdot)$. We consider zero-mean $\mathrm{GP}(0, k)$ priors and bounded covariance $k(\mathbf{x}_i, \mathbf{x}_i) \leq \kappa^2$ for all $\mathbf{x}_i \in \mathcal{A}$. An important property of Gaussian processes is that if we combine a prior $f \sim \mathrm{GP}(0, k)$ and assume that the observation noise is zero-mean Gaussian (i.e., $\eta_t \sim \mathcal{N}(0, \xi^2)$), then the posterior distribution of $f$ conditioned on a set of observations $\{(\mathbf{x}_s, y_s)\}_{s=1}^{t}$ is also a GP. More precisely, if $\mathbf{X}_t \triangleq [\mathbf{x}_1, \ldots, \mathbf{x}_t]^{\mathsf{T}} \in \mathbb{R}^{t \times d}$ is the matrix with all arms selected so far and $\mathbf{y}_t \triangleq [y_1, \ldots, y_t]^{\mathsf{T}}$ the corresponding observations, then the posterior is still a GP and the mean and variance of the function at a test point $\mathbf{x}$ are defined as

$$\mu_t(\mathbf{x} \mid \mathbf{X}_t, \mathbf{y}_t) = \mathbf{k}_t(\mathbf{x})^{\mathsf{T}}(\mathbf{K}_t + \lambda\mathbf{I})^{-1}\mathbf{y}_t, \tag{1}$$

$$\sigma_t^2(\mathbf{x} \mid \mathbf{X}_t) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_t(\mathbf{x})^{\mathsf{T}}(\mathbf{K}_t + \lambda\mathbf{I})^{-1}\mathbf{k}_t(\mathbf{x}), \tag{2}$$

where $\lambda \triangleq \xi^2$, $\mathbf{K}_t \in \mathbb{R}^{t \times t}$ is the matrix $[\mathbf{K}_t]_{i,j} \triangleq k(\mathbf{x}_i, \mathbf{x}_j)$ constructed from all pairs $\mathbf{x}_i, \mathbf{x}_j$ in $\mathbf{X}_t$, and $\mathbf{k}_t(\mathbf{x}) \triangleq [k(\mathbf{x}_1, \mathbf{x}), \ldots, k(\mathbf{x}_t, \mathbf{x})]^{\mathsf{T}}$. Notice that $\mathbf{k}_t(\mathbf{x})$ can be seen as an *embedding* of an arm $\mathbf{x}$ represented using by the arms $\mathbf{x}_1, \ldots, \mathbf{x}_T$ observed so far.

The GP-UCB *algorithm* is a uses a Gaussian process $\mathrm{GP}(0, k)$ as a prior for $f$. Inspired by the optimism-in-face-of-uncertainty principle, at each time step $t$, GP-UCB uses the posterior GP to compute the mean and variance of an arm $\mathbf{x}_i$ and obtain the score

$$u_t(\mathbf{x}_i) \triangleq \mu_t(\mathbf{x}_i) + \beta_t\sigma_t(\mathbf{x}_i), \tag{3}$$

where we use the short-hand notation $\mu_t(\cdot) \triangleq \mu(\cdot \mid \mathbf{X}_t, \mathbf{y}_t)$ and $\sigma_t(\cdot) \triangleq \sigma(\cdot \mid \mathbf{X}_t)$. Finally, GP-UCB chooses the maximizer $\mathbf{x}_{t+1} \triangleq \arg\max_{\mathbf{x}_i \in \mathcal{A}} u_t(\mathbf{x}_i)$ as the next arm to evaluate. According to the score $u_t$, an arm $\mathbf{x}$ is likely to be selected if it has high mean reward $\mu_t$ *or* high variance $\sigma_t$, i.e., its estimated reward $\mu_t(\mathbf{x})$ is very uncertain. As a result, selecting the arm $\mathbf{x}_{t+1}$ with the largest score trades off between collecting (estimated) large reward (*exploitation*) and improving the accuracy of the posterior (*exploration*). The parameter $\beta_t$ balances between these two objectives and must be properly tuned to guarantee low regret. Srinivas et al. (2010) proposes different approaches for tuning $\beta_t$, depending on the assumptions on $f$ and $\mathcal{A}$.

While GP-UCB is interpretable, simple to implement and provably achieves low regret, it is computationally expensive. In particular, computing $\sigma_t(\mathbf{x})$ has a complexity at least $\Omega(t^2)$.

---

1. We assume a consistent and arbitrary tie-breaking strategy.

Multiplying this complexity by $T$ iterations and $A$ arms results in an overall $\mathcal{O}(AT^3)$ cost, which does not scale to large number of iterations $T$.

## 3. Budgeted Kernel Bandits

In this section, we introduce the BKB (*budgeted kernel bandit*) algorithm, a novel efficient approximation of GP-UCB, and we provide guarantees for its computational complexity. The analysis in Section 4 shows that BKB can be tuned to significantly reduce the complexity of GP-UCB with a negligible impact on the regret. We begin by introducing the two major contributions of this section: an approximation of the GP-UCB scores supported only by a small subset $\mathcal{S}_t$ of *inducing points*, and a method to *incrementally and adaptively* construct an accurate subset $\mathcal{S}_t$.

### 3.1. The algorithm

The main computational bottleneck for getting the scores in Equation 3 is due to the fact that after $t$ steps, the posterior GP is supported on *all* $t$ previously seen arms. As a consequence, evaluating Equations 1 and 2 requires computing a $t$ dimensional vector $\mathbf{k}_t(\mathbf{x})$ and $t \times t$ matrix $\mathbf{K}_t$ respectively. To avoid this dependency we restrict both $\mathbf{k}_t$ and $\mathbf{K}_t$ to be supported on a *subset $\mathcal{S}_t$* of $m$ arms. This approach is a case of the sparse Gaussian process approximation (Quinonero-Candela et al., 2007), or equivalently, linear bandits constrained to a subspace (Kuzborskij et al., 2019).

**Approximated GP-UCB scores.** Consider a subset of arm $\mathcal{S}_t \triangleq \{\mathbf{x}_i\}_{i=1}^m$ and let $\mathbf{X}_{\mathcal{S}_t} \in \mathbb{R}^{m \times d}$ be the matrix with all arms in $\mathcal{S}_t$ as rows. Let $\mathbf{K}_{\mathcal{S}_t} \in \mathbb{R}^{m \times m}$ be the matrix constructed by evaluating the covariance $k$ between any two pairs of arms in $\mathcal{S}_t$ and $\mathbf{k}_{\mathcal{S}_t}(\mathbf{x}) \triangleq [k(\mathbf{x}_1, \mathbf{x}), \ldots, k(\mathbf{x}_m, \mathbf{x})]^\mathsf{T}$. The Nyström embedding $\mathbf{z}_t(\cdot)$ associated with subset $\mathcal{S}_t$ is defined as the mapping[2]

$$\mathbf{z}_t(\cdot) \triangleq \left(\mathbf{K}_{\mathcal{S}_t}^{1/2}\right)^{+} \mathbf{k}_{\mathcal{S}_t}(\cdot) : \mathbb{R}^d \to \mathbb{R}^m,$$

where $(\cdot)^{+}$ indicates the pseudo-inverse. We denote with $\mathbf{Z}_t(\mathbf{X}_t) \triangleq [\mathbf{z}_t(\mathbf{x}_1), \ldots, \mathbf{z}_t(\mathbf{x}_t)]^\mathsf{T} \in \mathbb{R}^{t \times m}$ the associated matrix of points and we define $\mathbf{V}_t \triangleq \mathbf{Z}_t(\mathbf{X}_t)^\mathsf{T} \mathbf{Z}_t(\mathbf{X}_t) + \lambda \mathbf{I}$. Then, we approximate the posterior mean, variance, and UCB for the value of the function at $\mathbf{x}_i$ as

$$\widetilde{\mu}_t(\mathbf{x}_i) \triangleq \mathbf{z}_t(\mathbf{x}_i)^\mathsf{T} \mathbf{V}_t^{-1} \mathbf{Z}_t(\mathbf{x}_i)^\mathsf{T} \mathbf{y}_t,$$

$$\widetilde{\sigma}_t^2(\mathbf{x}_i) \triangleq \frac{1}{\lambda}\Big(k(\mathbf{x}_i, \mathbf{x}_i) - \mathbf{z}_t(\mathbf{x}_i)^\mathsf{T} \mathbf{Z}_t(\mathbf{X}_t)^\mathsf{T} \mathbf{Z}_t(\mathbf{X}_t) \mathbf{V}_t^{-1} \mathbf{z}_t(\mathbf{x}_i)\Big),$$

$$\widetilde{u}_t(\mathbf{x}_i) \triangleq \widetilde{\mu}_t(\mathbf{x}_i) + \widetilde{\beta}_t \widetilde{\sigma}_t(\mathbf{x}_i), \tag{4}$$

where $\widetilde{\beta}_t$ is appropriately tuned to achieve small regret in the theoretical analysis of Section 4. Finally, at each time step $t$, BKB selects arm $\widetilde{\mathbf{x}}_{t+1} = \arg\max_{\mathbf{x}_i \in \mathcal{A}} \widetilde{u}_t(\mathbf{x}_i)$.

Notice that in general, $\widetilde{\mu}_t$ and $\widetilde{\sigma}_t$ do *not* correspond to any GP posterior. In fact, if we were simply replacing the $k(\mathbf{x}_i, \mathbf{x}_i)$ in the expression of $\widetilde{\sigma}_t^2(\mathbf{x}_i)$ by its value in the Nyström embedding, i.e., $\mathbf{z}_t(\mathbf{x}_i)^\mathsf{T}\mathbf{z}_t(\mathbf{x}_i)$, then we would recover a sparse GP approximation known as the *subset of regressors*. Using $\mathbf{z}_t(\mathbf{x}_i)^\mathsf{T}\mathbf{z}_t(\mathbf{x}_i)$ is known to cause *variance starvation*, as it can severely underestimate the variance of a test point $\mathbf{x}_i$ when it is far from the points in $\mathcal{S}_t$. Our formulation of $\widetilde{\sigma}_t$ is known in Bayesian world as the *deterministic training conditional* (DTC), where it is used as a heuristic to

---

2. Recall that in the exact version, $\mathbf{k}_t(\mathbf{x})$ can be seen as an embedding of any arm $\mathbf{x}$ into the space induced by all the $t$ arms selected so far, i.e. using all selected points as inducing points.

prevent variance starvation. However, DTC does *not* correspond to a GP since it violates consistency (Quinonero-Candela et al., 2007). In this work, we justify this approach rigorously and analyze it.

**Choosing the inducing points.** A critical for the low complexity of BKB while still controlling the regret is to carefully choose the inducing points to include in the subset $\mathcal{S}_t$. As the complexity of computing $\widetilde{u}_t$ scales with the size $m$ of $\mathcal{S}_t$, a smaller set gives a faster algorithm. Conversely, the difference between $\widetilde{\mu}_t$ and $\widetilde{\sigma}_t$ and their exact counterparts depends on the accuracy of the embedding $\mathbf{z}_t$, which increases with the size of $\mathcal{S}_t$. Moreover, even for a fixed $m$, the quality of the embedding greatly depends on *which* inducing points are included. For instance, selecting the same arm as inducing point twice, or two co-linear arms, does not improve accuracy as the embedding space does not change. Finally, we need to take into account two important aspects of sequential optimization when choosing $\mathcal{S}_t$. First, we need to focus our approximation more on regions of $\mathcal{A}$ that are relevant to the objective (i.e., high-reward arms). Second, as these regions change over time, we need to keep adapting the composition and size of $\mathcal{S}_t$. With these objectives in mind, we choose to construct $\mathcal{S}_t$ by randomly subsampling only out of the set of arms $\widetilde{\mathbf{X}}_t$ evaluated so far. In particular, arms are selected for inclusion in $\mathcal{S}_t$ with a probability proportional to their posterior variance $\sigma_t$ at step $t$. We report the selection procedure in Algorithm 1 with the complete BKB algorithm.

We initialize $\mathcal{S}_1 \triangleq \{\widetilde{\mathbf{x}}_1\}$ by selecting an arm uniformly at random. At each step $t$, after selecting $\widetilde{\mathbf{x}}_{t+1}$, we must regenerate $\mathcal{S}_t$ to reflect the changes in $\widetilde{\mathbf{X}}_{t+1}$. Ideally, we would sample each arm in $\widetilde{\mathbf{X}}_{t+1}$ proportionally to $\sigma_{t+1}^2$, but this would be too computationally expensive. Therefore, we first approximate $\sigma_{t+1}^2$ with $\sigma_t^2$. This is equivalent to ignoring the last arm and does not significantly impact the accuracy. We can then replace $\sigma_t^2$ with $\widetilde{\sigma}_t^2$ that was already efficiently computed when constructing Equation 4. Finally, given a parameter $\overline{q} \geq 1$, we set our approximate inclusion probability as $\widetilde{p}_{t+1,i} \triangleq \overline{q}\widetilde{\sigma}_t^2(\widetilde{\mathbf{x}}_s)$. The $\overline{q}$ parameter is used to increase the inclusion

---

**Algorithm 1:** BKB

**Data:** Arm set $\mathcal{A}$, $q$, $\{\beta_t\}_{t=1}^T$
**Result:** Arm choices $\mathcal{D}_T \leftarrow \{(\widetilde{\mathbf{x}}_t, y_t)\}$
 Select uniformly at random $\mathbf{x}_1$ and observe $y_1$;
 Initialize $\mathcal{S}_1 \leftarrow \{\mathbf{x}_1\}$;
**for** $t = \{1, \ldots, T-1\}$ **do**
  Compute $\widetilde{\mu}_t(\mathbf{x}_i)$ and $\widetilde{\sigma}_t^2(\mathbf{x}_i)$ for all $\mathbf{x}_i \in \mathcal{A}$;
  Select $\widetilde{\mathbf{x}}_{t+1} \leftarrow \arg\max_{\mathbf{x}_i \in \mathcal{A}} \widetilde{u}_t(\mathbf{x}_i)$ (Eq. 4);
  **for** $i = \{1, \ldots, t+1\}$ **do**
   Set $\widetilde{p}_{t+1,i} \leftarrow \overline{q} \cdot \widetilde{\sigma}_t^2(\widetilde{\mathbf{x}}_i)$;
   Draw $q_{t+1,i} \sim \text{Bernoulli}(\widetilde{p}_{t+1,i})$;
   **If** $q_{t+1} = 1$ **then** include $\widetilde{\mathbf{x}}_i$ in $\mathcal{S}_{t+1}$;
  **end**
**end**

---

probability in order to boost the overall success probability of the approximation procedure at the expense of a small increase in the size of $\mathcal{S}_{t+1}$. Given $\widetilde{p}_{t+1,i}$, we start from an empty $\mathcal{S}_{t+1}$ and iterate over all $\widetilde{\mathbf{x}}_i$ for $i \in [t+1]$ drawing $q_{t+1,i}$ from a Bernoulli distribution with probability $\widetilde{p}_{t+1,i}$. If $q_{t+1,i} = 1$, $\widetilde{\mathbf{x}}_i$ is included in $\mathcal{S}_{t+1}$.

Notice that while constructing $\mathcal{S}_t$ based on $\sigma_t^2$ is a common heuristic for sparse GPs, it has not been yet rigorously justified. In the next section, we show that this approach is equivalent to $\lambda$-ridge leverage score (RLS) sampling (Alaoui and Mahoney, 2015), a well studied tool in randomized linear algebra. We leverage the results from this field to prove both accuracy and efficiency guarantees.

## 3.2. Complexity analysis

Let $m_t \triangleq |\mathcal{S}_t|$ be the size of the set $\mathcal{S}_t$ at step $t$. At each step, we first compute the embedding $\mathbf{z}_t(\mathbf{x}_i)$ of all arms in $\mathcal{O}(Am_t^2 + m_t^3)$ time, which corresponds to one inversion of $\mathbf{K}_{\mathcal{S}_t}^{1/2}$ and the matrix-vector product specific to each arm. We then rebuild the matrix $\mathbf{V}_t$ from scratch using all
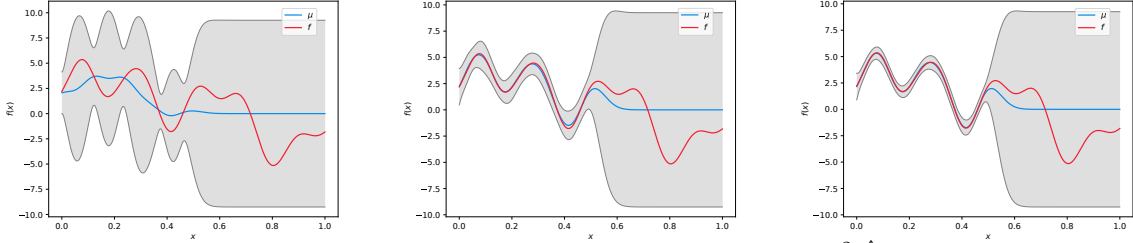
Fig. 1: We simulate a GP on $[0, 1] \in \mathbb{R}$ using Gaussian kernel with bandwidth $\sigma^2 \triangleq 100$. We draw $f$ from the GP and give to BKB $t \in \{6, 63, 215\}$ evaluations sampled uniformly in $[0, 0.5]$. We plot $f$ and $\widetilde{\mu}_t \pm 3\widetilde{\sigma}_t$.

the arms observed so far. In general, it is sufficient to record the counters of the arms pulled so far, rather than the full list of arms, so that $\mathbf{V}_t$ can be constructed in $\mathcal{O}(\min\{t, A\}m_t^2)$ time. Then, the inverse $\mathbf{V}_t^{-1}$ is computed in $\mathcal{O}(m_t^3)$ time. We can now efficiently compute $\widetilde{\mu}_t$, $\widetilde{\sigma}_t$, and $\widetilde{u}_t$ for all arms in $\mathcal{O}(Am_t^2)$ time reusing the embeddings and $\mathbf{V}_t^{-1}$. Finally, computing all $q_{t+1,i}$s and $\mathcal{S}_{t+1}$ takes $\mathcal{O}(\min\{t + 1, A\})$ time using the estimated variances $\widetilde{\sigma}_t^2$. As a result, the per-step complexity is of order $\mathcal{O}\big((A + \min\{t, A\})m_T^2\big)$.[3] Space-wise, we only need to store the embedded arms and $\mathbf{V}_t$ matrix, which takes at most $\mathcal{O}(Am_T)$ space.

**The size of $\mathcal{S}_T$.** The size $m_t$ of $\mathcal{S}_t$ can be expressed using the $q_{t,i}$ r.v. as the sum $m_t \triangleq \sum_{i=1}^{t} q_{t,i}$. In order to provide a bound on the total number of inducing points, which directly determines the computational complexity of BKB, we go through three major steps.

The first is to show that w.h.p., $m_t$ is close to the sum $\sum_{i=1}^{t} \widetilde{p}_{t,i} = \sum_{i=1}^{t} \overline{q}\widetilde{\sigma}_t^2(\widetilde{\mathbf{x}}_i)$, i.e., close to the sum of the probabilities we used to sample each $q_{t,i}$. However, the different $q_{t,i}$ are *not independent* and each $\widetilde{p}_{t,i}$ is itself a r.v. Nonetheless all $q_{t,i}$ are conditionally independent given the previous $t - 1$ steps, and this is sufficient to obtain the result.

The second and a more complex step is to guarantee that the random sum $\sum_{i=1}^{t} \widetilde{\sigma}_t^2(\widetilde{\mathbf{x}}_i)$ is close to $\sum_{i=1}^{t} \sigma_t^2(\widetilde{\mathbf{x}}_i)$ and, at a lower level, that each individual estimate $\widetilde{\sigma}_t^2(\cdot)$ is close to $\sigma_t^2(\cdot)$. To achieve this we exploit the connection between ridge leverage scores and posterior variance $\sigma_t^2$. In particular, we show that the variance estimator $\widetilde{\sigma}_t^2(\cdot)$ used by BKB is a variation of the RLS estimator of Calandriello et al. (2017a).

The first two steps lead to $m_t \approx \sum_{i=1}^{t} \sigma_i^2(\widetilde{\mathbf{x}}_i)$, for which we need to derive a more explicit bound. In the GP analyses, this quantity is bounded using the maximal information gain $\gamma_T$ after $T$ rounds. For this, let $\mathbf{X}_{\mathcal{A}} \in \mathbb{R}^{A \times d}$ be the matrix with all arms as rows, $\mathcal{D}$ a subset of these rows, potentially with duplicates, and $\mathbf{K}_{\mathcal{D}}$ the associated kernel matrix. Then, Srinivas et al. (2010) define

$$\gamma_T \triangleq \max_{\mathcal{D} \subset \mathcal{A}: |\mathcal{D}| = T} \frac{1}{2} \log \det(\mathbf{K}_{\mathcal{D}}/\lambda + \mathbf{I}),$$

and show that $\sum_{i=1}^{t} \sigma_i^2(\widetilde{\mathbf{x}}_i) \leq \gamma_t$, and that $\gamma_T$ itself can be bounded for specific $\mathcal{A}$ and kernel functions, e.g., $\gamma_T \leq \mathcal{O}(\log(T)^{d+1})$ for Gaussian kernels. Using the equivalence between RLS and posterior variance $\sigma_t^2$, we can also relate the posterior variance $\sigma_t^2(\widetilde{\mathbf{x}}_i)$ of the evaluated arms to the GP's *effective dimension* $d_{\text{eff}}$ or degrees of freedom

$$d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_T) \triangleq \sum_{i=1}^{t} \sigma_t^2(\widetilde{\mathbf{x}}_i) = \text{Tr}(\mathbf{K}_T(\mathbf{K}_T + \lambda\mathbf{I})^{-1}), \tag{5}$$

using the following inequality by Calandriello et al. (2017b),

$$\log \det(\mathbf{K}_T/\lambda + \mathbf{I}) \leq \text{Tr}(\mathbf{K}_T(\mathbf{K}_T + \lambda\mathbf{I})^{-1})\Big(1 + \log\Big(\frac{\|\mathbf{K}_T\|}{\lambda} + 1\Big)\Big). \tag{6}$$

---

3. Notice that $m_t \leq \min\{t, a\}$ and thus the complexity term $\mathcal{O}(m_t^3)$ is absorbed by the other terms.

We use both RLS and $d_{\text{eff}}$ to describe BKB's selection. We now give the main result of this section.

**Theorem 1** *For a desired $0 < \varepsilon < 1$, $0 < \delta < 1$, let $\alpha \triangleq (1 + \varepsilon)/(1 - \varepsilon)$. If we run* BKB *with $\overline{q} \geq 6\alpha \log(4T/\delta)/\varepsilon^2$, then with probability $1 - \delta$, for all $t \in [T]$ and for all $\mathbf{x} \in \mathcal{A}$, we have*

$$\sigma_t^2(\mathbf{x})/\alpha \leq \widetilde{\sigma}_t^2(\mathbf{x}) \leq \alpha\sigma_t^2(\mathbf{x}) \qquad \text{and} \qquad |\mathcal{S}_t| \leq 3(1 + \kappa^2/\lambda)\alpha\overline{q}d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_t).$$

**Computational complexity.** We already showed that BKB's implementation with Nyström embedding requires $\mathcal{O}(T(A + \min\{t, A\})m_T^3)$ time and $\mathcal{O}(Am_T)$ space. Combining this with Theorem 1 and the bound $m_T \leq \widetilde{\mathcal{O}}(d_{\text{eff}})$, we obtain a $\widetilde{\mathcal{O}}(TAd_{\text{eff}}^2 + \min\{t, A\})d_{\text{eff}}^3)$ time complexity. Whenever $d_{\text{eff}} \ll T$ and $T \ll A$, this is essentially a quadratic $\mathcal{O}(T^2)$ runtime, a large improvement over the quartic $\mathcal{O}(T^4) \leq \mathcal{O}(T^3A)$ runtime of GP-UCB.

**Tuning $\overline{q}$.** Note that although $\overline{q}$ must satisfy the condition of Theorem 1 for the result to hold, it is quite robust to uncertainty on the desired horizon $T$. In particular, the bound holds for *any $\varepsilon > 0$*, and even if we continue updating $\mathcal{S}_T$ after the $T$-th step, the bound still holds by implicitly increasing the parameter $\varepsilon$. Alternatively, after the $T$-th iteration the user can suspend the algorithm, increase $\overline{q}$ to suit the new desired horizon, and rerun only the subset selection on the arms selected so far.

**Avoiding variance starvation.** Another important consequence of Theorem 1 is that BKB's variance estimate is always close to the exact one up to a small constant factor. To the best of our knowledge, it makes BKB the first efficient and general GP algorithm that provably avoids variance starvation, which can be caused by two sources of error. The first source is the degeneracy, i.e., low-rankness of the GP approximation which causes the estimate to grow over-confident when the number of observed points grows and exceeds the degrees of freedom of the GP. BKB *adaptively chooses its degrees of freedom* as the size of $\mathcal{S}_t$ scales with the effective dimension. The second source of error arises when a point is far away from $\mathcal{S}_t$. Our use of a DTC variance estimator avoids under-estimation before we update the subset $\mathcal{S}_t$. Afterward, we can use guarantees on the quality of $\mathcal{S}_t$ to guarantee that we do not over-estimate the variance too much, exploiting a similar approach used to guarantee accuracy in RLS estimation. Both problems, and BKB's accuracy, are highlighted in Figure 1 using a benchmark experiment proposed by Wang et al. (2018).

**Incremental dictionary update.** At each step $t$, BKB recomputes the dictionary $\mathcal{S}_{t+1}$ from scratch by sampling each of the arms pulled so far with a suitable probability $\widetilde{p}_{t+1,i}$. A more efficient variant would be to build $\mathcal{S}_{t+1}$ by adding the new point $\mathbf{x}_{t+1}$ with probability $\widetilde{p}_{t+1,t+1}$ and including the points in $\mathcal{S}_t$ with probability $\widetilde{p}_{t+1,i}/\widetilde{p}_{t,i}$. This strategy is used in the streaming setting to avoid storing all points observed so far and update the dictionary (see Calandriello et al., 2017a). Nonetheless, the stream of points, although arbitrary, is assumed to be generated *independently* from the dictionary itself. On the other hand, in our bandit setting, the points $\widetilde{\mathbf{x}}_1, \widetilde{\mathbf{x}}_2, \dots$ are actually chosen by the learner depending on the dictionaries built over time, thus building a strong dependency between the stream of points and the dictionary itself. How to analyze such dependency and whether the accuracy of the inducing points is preserved in this case remains as an open question.

## 4. Regret Analysis

We are now ready to present the second main contribution of this paper, a bound on the regret achieved by BKB. To prove our result we additionally assume that the reward function $f$ has a bounded norm, i.e., $\|f\|_{\mathcal{H}}^2 \triangleq \langle f, f \rangle < \infty$. We use an upper-bound $\|f\|_{\mathcal{H}} \leq F$ to properly tune $\widetilde{\beta}_t$ to the range of the rewards. If $F$ is not known in advance, standard guess-and-double techniques apply.

**Theorem 2** *Assume $\|f\|_{\mathcal{H}} \leq F < \infty$. For any desired $0 < \varepsilon < 1$, $0 < \delta < 1$, $0 < \lambda$, let $\alpha \triangleq (1+\varepsilon)/(1-\varepsilon)$ and $\overline{q} \geq 6\alpha \log(4T/\delta)/\varepsilon^2$. If we run BKB with*

$$\widetilde{\beta}_t \triangleq 2\xi \sqrt{\alpha \log(\kappa^2 t) \left( \sum_{s=1}^t \widetilde{\sigma}_t^2(\widetilde{\mathbf{x}}_s) \right) + \log(1/\delta)} + \left( 1 + \tfrac{1}{\sqrt{1-\varepsilon}} \right) \sqrt{\lambda} F,$$

*then, with probability of at least $1 - \delta$, the regret $R_T$ of BKB is bounded as*

$$R_T \leq 2(2\alpha)^{3/2}\sqrt{T}\left( \xi d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_T) \log(\kappa^2 T) + \sqrt{\lambda F^2 d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_T) \log(\kappa^2 T)} + \xi \log(1/\delta) \right).$$

Theorem 2 shows that BKB achieves exactly the same regret as (exact) GP-UCB up to small $\alpha$ constant and $\log(\kappa^2 T)$ multiplicative factor.[4] For instance, setting $\varepsilon = 1/2$ results in a bound only $3 \log T$ times larger than the of of GP-UCB. At the same time, the choice $\varepsilon = 1/2$ only accounts for a constant factor 12 in the per-step computational complexity, which is still dramatically reduced from $t^2 A$ to $d_{\text{eff}}^2 A$. Note also that even if we send $\varepsilon$ to 0, in the worst case we will include all arms selected so far, i.e., $\mathcal{S}_t = \{\widetilde{\mathbf{X}}_t\}$. Therefore, even in this case BKB's runtime does not grow unbounded, but BKB transforms back into exact GP-UCB. Moreover, we show that $d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_T) \leq \log \det(\mathbf{K}_T/\lambda + \mathbf{I})$, as in Proposition 5 in the appendix, so any bound on $\log \det(\mathbf{K}_T/\lambda + \mathbf{I})$ available for GP-UCB applies directly to BKB. This means that up to an extra $\log T$ factor, we match GP-UCB's $\widetilde{\mathcal{O}}(\log(T)^{2d})$ rate for the Gaussian kernel, $\widetilde{\mathcal{O}}(T^{\frac{1}{2}\frac{2\nu+3d^2}{2\nu+d^2}})$ rate for the Matérn kernel, and $\widetilde{\mathcal{O}}(d\sqrt{T})$ for the linear kernel. While these bounds are not minimax optimal, they closely follow the lower bounds derived in Scarlett et al. (2017). On the other hand, in the case of linear kernel (i.e., the linear bandits) we nearly match the lower bound of Dani et al. (2008).

Another interesting aspect of BKB is that computing the trade-off parameter $\widetilde{\beta}_t$ can be done efficiently. Previous methods bounded this quantity with a loose (deterministic) upper bound, e.g., $\mathcal{O}(\log(T)^d)$ for Gaussian kernels, to avoid the large cost of computing $\log \det(\mathbf{K}_T/\lambda + \mathbf{I})$. In our $\widetilde{\beta}_t$, we bound the $\log \det$ by $d_{\text{eff}}$, which is then bounded by $\sum_{s=1}^t \widetilde{\sigma}_t^2(\mathbf{x}_s)$, see Thm. 1, where all $\widetilde{\sigma}_t^2$s are already efficiently computed at each step. While this is up to $\log t$ larger than the exact $\log \det$, it is *data adaptive* and much smaller than the known worst case upper bounds.

It it crucial, that our regret guarantee is achieved without requiring an *increasing accuracy* in our approximation. One would expect that to obtain a sublinear regret the error induced by the approximation should decrease as $1/T$. Instead, in BKB, the constants $\varepsilon$ and $\lambda$ that govern the accuracy level are fixed and thus it is not possible to guarantee that $\widetilde{\mu}_t$ will ever get close to $\mu_t$ everywhere. Adaptivity is the key: we can afford the same approximation level at every step because accuracy is actually increased only on a specific part of the arm set. For example, if a suboptimal arm is selected too often due to bad approximation, it will be eventually included in $\mathcal{S}_t$. After the inclusion, the approximation accuracy in the region of the suboptimal arm increases, and it would not be selected anymore. As the set of inducing points is updated *fast enough*, the impact of inaccurate approximations is limited over time, thus preventing large regret to accumulate. Note that this is a significant divergence from existing results. In particular approximation bounds that are uniformly accurate for all $\mathbf{x}_i \in \mathcal{A}$, such as those obtained with quadrature FF (Mutný and Krause, 2018), rely on packing arguments. Due to the nature of packing, this usually causes the runtime or regret to scale

---

4. Here we derive a *frequentist* regret bound and thus we compare with the result of Chowdhury and Gopalan (2017) rather than the original *Bayesian* analysis of Srinivas et al. (2010).

exponentially wiht the input dimension $d$, and requires kernel $k$ to have a specific structure, e.g., to be statonary. Our new analysis avoids both of these problems.

Finally, we point out that the adaptivity of BKB allows drawing an interesting connection between learning and computational complexity. In fact, both the regret and the computation of BKB scale with the log-determinant and the effective dimension of $\mathbf{K}_T$, which is related to the effective dimension of the sequence of arms selected over time. As a result, if the problem is difficult as.a learning problem, then BKB automatically adapts the set $\mathcal{S}_t$ by including many more inducing points to guarantee the level of accuracy needed to solve the problem. Conversely, if the problem is simple, then BKB greatly reduces the size of $\mathcal{S}_t$ and achieves the derived level of accuracy.

### 4.1. Proof sketch

We build on the GP-UCB analysis of Chowdhury and Gopalan (2017). Their analysis relies on a confidence interval formulation of GP-UCB that is more conveniently expressed using an explicit feature-based representation of the GP. For any GP with covariance $k$, there is a corresponding RKHS $\mathcal{H}$ with $k$ as its kernel function. Furthermore, any kernel function $k$ is associated to a non-linear feature map $\phi(\cdot) : \mathbb{R}^d \to \mathcal{H}$ such that $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}')^\mathsf{T}\phi(\mathbf{x}')$. As a result, any reward function $f \in \mathcal{H}$ can be written as $f(\mathbf{x}) = \phi(\mathbf{x})^\mathsf{T}\mathbf{w}_\star$, where $\mathbf{w}_\star \in \mathcal{H}$.

**Confidence-interval view of GP-UCB.** Let $\mathbf{\Phi}(\mathbf{X}_t) \triangleq [\phi(\mathbf{x}_1), \ldots, \phi(\mathbf{x}_t)]^\mathsf{T}$ be the matrix $\mathbf{X}_t$ after the application of $\phi(\cdot)$ to each row. We can then define the regularized design matrix as $\mathbf{A}_t \triangleq \mathbf{\Phi}(\mathbf{X}_t)^\mathsf{T}\mathbf{\Phi}(\mathbf{X}_t) + \lambda\mathbf{I}$, and then computer regularized least-squares estimate as

$$\widehat{\mathbf{w}}_t \triangleq \arg\min_{\mathbf{w}\in\mathcal{H}} \sum_{i=1}^{t} (y_i - \phi(\mathbf{x}_i)^\mathsf{T}\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2 = \mathbf{A}_t^{-1}\mathbf{\Phi}(\mathbf{X}_t)^\mathsf{T}\mathbf{y}_t.$$

We define the *confidence interval* $C_t$ as the ellipsoid induced by $\mathbf{A}_t$ with center $\widehat{\mathbf{w}}_t$ and radius $\beta_t$

$$C_t \triangleq \{\mathbf{w} : \|\mathbf{w} - \widehat{\mathbf{w}}_t\|_{\mathbf{A}_t} \le \beta_t\}, \qquad \beta_t \triangleq \lambda^{1/2}F + R\sqrt{2(\log\det(\mathbf{A}_t/\lambda) + \log(1/\delta))}, \quad (7)$$

where the radius $\beta_t$ is such that $\mathbf{w}_\star \in C_t$ w.h.p. (Chowdhury and Gopalan, 2017). Finally, using Lagrange multipliers we reformulate the GP-UCB scores as

$$u_t(\mathbf{x}_i) = \max_{\mathbf{w}\in C_t} \phi(\mathbf{x}_i)^\mathsf{T}\mathbf{w} = \overbrace{\phi(\mathbf{x}_i)^\mathsf{T}\widehat{\mathbf{w}}_t}^{\mu_t(\mathbf{x}_i)} + \beta_t \overbrace{\sqrt{\phi(\mathbf{x}_i)^\mathsf{T}\mathbf{A}_t^{-1}\phi(\mathbf{x}_i)}}^{\sigma_t(\mathbf{x}_i)}. \quad (8)$$

**Approximating the confidence ellipsoid.** Consider subset of arm $\mathcal{S}_t = \{\mathbf{x}_i\}_{i=1}^{m}$ chosen by BKB at each step and denote by $\mathbf{X}_{\mathcal{S}_t} \in \mathbb{R}^{m\times d}$ the matrix with all arms in $\mathcal{S}_t$ as rows. Let $\widetilde{\mathcal{H}}_t \triangleq \mathrm{Im}(\mathbf{\Phi}(\mathbf{X}_{\mathcal{S}_t}))$ be the smaller $m$-rank RKHS spanned by $\mathbf{\Phi}(\mathbf{X}_{\mathcal{S}_t})$; and by $\mathbf{P}_t$ the symmetric orthogonal projection operator on $\widetilde{\mathcal{H}}_t$. We then define an *approximate* feature map $\widetilde{\phi}_t(\cdot) \triangleq \mathbf{P}_t\phi(\cdot) : \mathbb{R}^d \to \widetilde{\mathcal{H}}_t$ and associated approximations of $\mathbf{A}_t$ and $\widehat{\mathbf{w}}_t$ as

$$\widetilde{\mathbf{A}}_t \triangleq \widetilde{\mathbf{\Phi}}_t(\mathbf{X}_t)^\mathsf{T}\widetilde{\mathbf{\Phi}}_t(\mathbf{X}_t) + \lambda\mathbf{I}, \quad (9)$$

$$\widetilde{\mathbf{w}}_t \triangleq \arg\min_{\mathbf{w}\in\mathcal{H}} \sum_{i=1}^{t} (y_i - \widetilde{\phi}(\mathbf{x}_i)^\mathsf{T}\mathbf{w})^2 + \lambda\|\mathbf{w}\|_2^2 = \widetilde{\mathbf{A}}_t^{-1}\widetilde{\mathbf{\Phi}}_t(\mathbf{X}_t)^\mathsf{T}\mathbf{y}_t. \quad (10)$$

This leads to an approximate confidence ellipsoid $\widetilde{C}_t \triangleq \{\mathbf{w} : \|\mathbf{w} - \widetilde{\mathbf{w}}_t\|_{\widetilde{\mathbf{A}}_t} \le \widetilde{\beta}_t\}$. A subtle element in these definitions is that while $\widetilde{\mathbf{\Phi}}_t(\mathbf{X}_t)^\mathsf{T}\widetilde{\mathbf{\Phi}}_t(\mathbf{X}_t)$ and $\widetilde{\mathbf{w}}_t$ are now *restricted* to $\widetilde{\mathcal{H}}_t$, the identity operator $\lambda\mathbf{I}$ in the regularization of $\widetilde{\mathbf{A}}_t$ still *acts over the whole* $\mathcal{H}$, and therefore $\widetilde{\mathbf{A}}_t$ does not belong

to $\widetilde{\mathcal{H}}_t$ and remains full-rank and invertible. This immediately leads to the usage of $k(\mathbf{x}_i, \mathbf{x}_i)$ in the definition of $\widetilde{\sigma}$ in Eq. 4, instead of the its approximate version using the Nyström embedding.

**Bounding the regret.** To find an appropriate $\widetilde{\beta}_t$, we follow an approach similar to the of Abbasi-Yadkori et al. (2011). Exploiting the relationship $y_t = \widetilde{\phi}(\widetilde{\mathbf{x}}_t)^\intercal \mathbf{w}_\star + \eta_t$, we bound

$$\|\mathbf{w}_\star - \widetilde{\mathbf{w}}_t\|_{\widetilde{\mathbf{A}}_t}^2 \leq \overbrace{\lambda^{1/2}\|\mathbf{w}_\star\|}^{(a)} + \overbrace{\|\widetilde{\boldsymbol{\Phi}}_t(\mathbf{X}_t)\boldsymbol{\eta}_t\|_{\widetilde{\mathbf{A}}_t^{-1}}}^{(b)} + \overbrace{\|\boldsymbol{\Phi}(\mathbf{X}_t)^\intercal\|_{\mathbf{I}-\mathbf{P}_t} \cdot \|\mathbf{w}_\star\|}^{(c)}.$$

Both $(a)$ and $(b)$ are present in GP-UCB's and OFUL's analysis. First, term $(a)$ is due to the bias introduced in the least-square estimator $\widetilde{\mathbf{w}}_t$ by the regularization $\lambda$. Then, term $(b)$ is due to the noisy rewards. Note that the same term $(b)$ appears in GP-UCB's analysis as $\|\boldsymbol{\Phi}(\mathbf{X}_t)\boldsymbol{\eta}_t\|_{\mathbf{A}_t^{-1}}$ and it is bounded by $\log\det(\mathbf{A}_t/\lambda)$ using self-normalizing concentration inequalities (Chowdhury and Gopalan, 2017). However, our $\|\widetilde{\boldsymbol{\Phi}}_t(\mathbf{X}_t)\boldsymbol{\eta}_t\|_{\widetilde{\mathbf{A}}_t^{-1}}$ is a more complex object, since the projection $\mathbf{P}_t$ contained in $\widetilde{\boldsymbol{\Phi}}_t(\mathbf{X}_t) \triangleq \mathbf{P}_t\boldsymbol{\Phi}(\mathbf{X}_t)$ depends on the whole process up to time time $t$, and therefore $\widetilde{\boldsymbol{\Phi}}_t(\mathbf{X}_t)$ also depends on the whole process, losing its martingale structure. To avoid this, we use Sylvester's identity and the projection operator $\mathbf{P}_t$ to bound

$$\log\det(\widetilde{\mathbf{A}}_t/\lambda) = \log\det\left(\tfrac{\boldsymbol{\Phi}(\mathbf{X}_t)\mathbf{P}_t\boldsymbol{\Phi}(\mathbf{X}_t)^\intercal}{\lambda} + \mathbf{I}\right) \leq \log\det\left(\tfrac{\boldsymbol{\Phi}(\mathbf{X}_t)\boldsymbol{\Phi}(\mathbf{X}_t)^\intercal}{\lambda} + \mathbf{I}\right) = \log\det(\mathbf{A}_t/\lambda).$$

In other words, restricting the problem to $\widetilde{\mathcal{H}}_t$ acts as a regularization and reduces the variance of the martingale. Unfortunately, $\log\det(\mathbf{A}_t/\lambda)$ is too expensive to compute, so we first bound it with $d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_t)\log(\kappa^2 t)$, and then we bound $d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_t) \leq \alpha\sum_{s=1}^t \widetilde{\sigma}_t^2(\mathbf{x}_s)$, Theorem 1, which can be computed efficiently. Finally, a new bias term $(c)$ appears. Combining Theorem 1 with the results of Calandriello and Rosasco (2018) for projection $\mathbf{P}_t$ obtained using RLSs sampling, we show that

$$\mathbf{I} - \mathbf{P} \preceq \lambda\mathbf{A}_t^{-1}/(1-\varepsilon).$$

The combination of $(a)$, $(b)$, and $(c)$ leads to the definition of $\widetilde{\beta}_t$ and the final regret bound as $R_T \leq \sqrt{\widetilde{\beta}_T}\sqrt{\sum_{t=1}^T \boldsymbol{\phi}(\mathbf{x}_t)^\intercal \widetilde{\mathbf{A}}_t^{-1}\boldsymbol{\phi}(\mathbf{x}_t)}$. To conclude the proof, we bound $\sum_{t=1}^T \boldsymbol{\phi}(\mathbf{x}_t)^\intercal \widetilde{\mathbf{A}}_t^{-1}\boldsymbol{\phi}(\mathbf{x}_t)$ with the following corollary of Theorem 1.

**Corollary 3** *Under the same conditions as Theorem 2, for all $t \in T$, we have $\mathbf{A}_t/\alpha \preceq \widetilde{\mathbf{A}}_t \preceq \alpha\mathbf{A}_t$.*

**Remarks.** The novel bound $\|\boldsymbol{\Phi}(\mathbf{X}_t)^\intercal\|_{\mathbf{I}-\mathbf{P}_t} \leq \frac{\lambda}{1-\varepsilon}\|\boldsymbol{\Phi}(\mathbf{X}_t)^\intercal\|_{\mathbf{A}_t^{-1}}$ has a crucial role in controlling the bias due to the projection $\mathbf{P}_t$. Note that the second term measures the error with the same metric $\mathbf{A}_t^{-1}$ used by the variance martingale. In other words, the bias introduced by BKB's approximation can be seen as a *self-normalizing* bias. It is larger along directions that have been sampled less frequently, and smaller along directions correlated with arms selected often (e.g., the optimal arm).

Our analysis bears some similarity with the one recently and independently developed by Kuzborskij et al. (2019). Nonetheless, our proof improves their result along two dimensions. First, we consider the more general (and challenging) GP optimization setting. Second, *we do not fix* the rank of our approximation in advance. While their analysis also exploits a self-normalized bias argument, this applies only to the $k$ largest components. If the problem has an effective dimension larger than $k$, their radius and regret is linear. In BKB we use our adaptive sampling scheme to include all necessary directions and to achieve the same regret rate as exact GP-UCB.

## 5. Discussion

As the prior work in Bayesian optimization is vast, we do not compare to alternative GP acquisition functions, such as GP-EI or GP-PI, and only focus on approximation techniques with theoretical

guarantees. Similarly, we exclude scalable variational inference based methods, even when their approximate posterior is provably accurate such as pF-DTC (Huggins et al., 2019), since they only provide guarantees for GP regression and not for the more difficult optimization setting. We also do not discuss SUPKERNELUCB (Valko et al., 2013), which has a tighter analysis than GP-UCB, since the algorithm does not work well in practice.

**Infinite arm sets.** Looking at the proof of Theorem 1, the guarantees on $\widetilde{u}_t$ hold for any $\mathcal{H}$, and in Theorem 2, we only require that the maximum $\widetilde{\mathbf{x}}_{t+1} \triangleq \arg\max_{\mathbf{x} \in \mathcal{A}} \max_{\mathbf{w} \in \widetilde{C}_t} \phi(\mathbf{x})^{\mathsf{T}} \mathbf{w}$ is returned. Therefore, the accuracy and regret guarantees also hold also for an infinite set of arms $\mathcal{A}$. However, the search over $\mathcal{A}$ can be difficult. In the general case, maximization of a GP posterior is an NP-hard problem, with algorithms that often scale exponentially with the input dimension $d$ and are not practical. We treated the easier case of finite sets, where the enumeration is sufficient. Note that this automatically introduces an $\Omega(A)$ runtime dependency, which could be removed if the user provides an efficient method to solve the maximization problem on a specific infinite set $\mathcal{A}$. As an example, Mutný and Krause (2018) prove that a GP posterior approximated using QFF can be optimized efficiently in low dimensions and we expect similar results hold for BKB and low *effective* dimension. Finally, note that recomputing a new set $\mathcal{S}_t$ still requires $\min\{A, t\} d_{\text{eff}}^2$ at each step. This is a bottleneck in BKB independent from the arm selection, and must be addressed in future.

**Linear bandit with matrix sketching.** Our analysis is related to the ones of CBRAP (Yu et al., 2017) and SOFUL (Kuzborskij et al., 2019). CBRAP uses Gaussian projections to embed all arms in a lower dimensional space. Unfortunately, their approach must either use an embedded space at least $\Omega(T)$ large, which can making it slower than the exact OFUL,or it incurs regret w.h.p. Another approach for linear bandits by Kuzborskij et al. (2019), SOFUL, uses frequent directions (Ghashami et al., 2016), a method similar to incremental PCA, to embed the arms into $\mathbb{R}^m$, where $m$ is *fixed* in advance. To compare, we distinguish between SOFUL-UCB and SOFUL-TS, a variant based on Thompson sampling. SOFUL-UCB achieves a $\widetilde{\mathcal{O}}(TAm^2)$ runtime and $\widetilde{\mathcal{O}}((1+\varepsilon_m)^{3/2}(d+m)\sqrt{T})$ regret, where $\varepsilon_m$ is the sum of the $d - m$ smallest eigenvalues of $\mathbf{A}_T$. However, notice that if the tail do not decrease quickly, this algorithm also suffers linear regret and no adaptive way to tune $m$ is known. On the same task BKB achieves a $\widetilde{\mathcal{O}}(d\sqrt{T})$ regret, since it adaptively chooses the size of the embedding. Computationally, directly instantiating BKB to use a linear kernel would achieve a $\widetilde{\mathcal{O}}(TAm_t^2)$ runtime, matching Kuzborskij et al. (2019)'s. Compared to SOFUL-TS, BKB achieves better regret, but is potentially slower. Since Thompson sampling does not need to compute all confidence intervals, but solves a simpler optimization problem, SOFUL-TS requires only $\widetilde{\mathcal{O}}(TAm)$ time against BKB's $\widetilde{\mathcal{O}}(TAm_t^2)$. It is unknown if a variant of BKB can match this complexity.

**Approximate GP with RFF.** RFF methods are popular to transform GP optimization to finite-dimensional problems and allow for scalability. Unfortunately, GP-UCB with standard RFF is not low-regret as RFF are suffer from variance starvation (Wang et al., 2018) and unfeasibly large RFF embeddings are necessary to prevent it. Recently, Mutný and Krause (2018) proposed an alternative based on QFF, a specialized approach for stationary kernels. They achieve the same regret rate as GP-UCB and BKB, with a near-optimal $\mathcal{O}(TA\log(T)^{d+1})$ runtime and give additional variations based on Thompson sampling and exact posterior maximization. However, quadrature-based approaches apply only to stationary kernels and require to $\varepsilon$-cover $\mathcal{A}$, hence they cannot escape an exponential dependence on $d$. Conversely, BKB applies to *any* kernel and while not specifically designed for this task it achieves a close $\widetilde{\mathcal{O}}(TA\log(T)^{3(d+1)})$ runtime. Moreover, in practice the size of $\mathcal{S}_T$ is less than exponential in $d$.

# References

Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Neural Information Processing Systems*, 2011.

Ahmed El Alaoui and Michael W. Mahoney. Fast randomized kernel methods with statistical guarantees. In *Neural Information Processing Systems*, 2015.

Daniele Calandriello and Lorenzo Rosasco. Statistical and computational trade-offs in kernel k-means. In *Neural Information Processing Systems*, 2018.

Daniele Calandriello, Alessandro Lazaric, and Michal Valko. Distributed adaptive sampling for kernel matrix approximation. In *International Conference on Artificial Intelligence and Statistics*, 2017a.

Daniele Calandriello, Alessandro Lazaric, and Michal Valko. Second-order kernel online convex optimization with adaptive sketching. In *International Conference on Machine Learning*, 2017b.

Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, 2017.

Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized Markov decision processes. In *International Conference on Artificial Intelligence and Statistics*, 2019.

Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback. In *Conference on Learning Theory*, 2008.

Mina Ghashami, Edo Liberty, Jeff M Phillips, and David P. Woodruff. Frequent directions: Simple and deterministic matrix sketching. *The SIAM Journal of Computing*, pages 1–28, 2016.

Avishek Ghosh, Sayak Ray Chowdhury, and Aditya Gopalan. Misspecified linear bandits. In *AAAI Conference on Artificial Intelligence*, 2017.

Elad Hazan, Adam Tauman Kalai, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. In *Conference on Learning Theory*, 2006.

Jonathan H. Huggins, Trevor Campbell, Mikołaj Kasprzak, and Tamara Broderick. Scalable Gaussian process inference with finite-data mean and variance guarantees. In *International Conference on Artificial Intelligence and Statistics*, 2019.

Ilja Kuzborskij, Leonardo Cella, and Nicolò Cesa-Bianchi. Efficient linear bandits through matrix sketching. In *International Conference on Artificial Intelligence and Statistics*, 2019.

Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. 2019.

Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to personalized news article recommendation. *International World Wide Web Conference*, 2010.

Haitao Liu, Yew-Soon Ong, Xiaobo Shen, and Jianfei Cai. When Gaussian process meets big data: A Review of scalable GPs. Technical report, 2018.

Jonas Mockus. *Global optimization and the Bayesian approach*. 1989.

Mojmír Mutný and Andreas Krause. Efficient high-dimensional Bayesian optimization with additivity and quadrature Fourier features. In *Neural Information Processing Systems*, 2018.

Martin Pelikan. Hierarchical Bayesian optimization algorithm. In *Studies in Fuzziness and Soft Computing*, pages 105–129. 2005.

Joaquin Quinonero-Candela, Carl Edward Rasmussen, and Christopher K. I. Williams. Approximation methods for gaussian process regression. *Large-scale kernel machines*, pages 203–224, 2007.

Ali Rahimi and Ben Recht. Random features for large-scale kernel machines. In *Neural Information Processing Systems*, 2007.

Carl Edward. Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. MIT Press, 2006.

Herbert Robbins. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematics Society*, 58:527–535, 1952.

Jonathan Scarlett, Ilija Bogunovic, and Volkan Cevher. Lower bounds on regret for noisy Gaussian process bandit optimization. In *Conference on Learning Theory*, 2017.

Matthias Seeger, Christopher Williams, and Neil Lawrence. Fast forward selection to speed up sparse Gaussian process regression. In *International Conference on Artificial Intelligence and Statistics*, 2003.

Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms. In *Neural Information Processing Systems*, 2012.

Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *International Conference on Machine Learning*, 2010.

Joel Aaron Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends in Machine Learning*, 8(1-2):1–230, 2015.

Michal Valko, Nathan Korda, Rémi Munos, Ilias Flaounas, and Nelo Cristianini. Finite-time analysis of kernelised contextual bandits. In *Uncertainty in Artificial Intelligence*, 2013.

Grace Wahba. *Spline models for observational data*. Society for Industrial and Applied Mathematics, 1990.

Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched Large-scale Bayesian Optimization in High-dimensional Spaces. In *International Conference on Artificial Intelligence and Statistics*, 2018.

Xiaotian Yu, Michael R. Lyu, and Irwin King. CBRAP: Contextual bandits with random projection. In *AAAI Conference on Artificial Intelligence*, 2017.

## Appendix A. Relaxing assumptions

In our derivations, we make several assumptions. While some are necessary, others can be relaxed.
**Assumptions on the noise.** Throughout the paper, we assume that the noise $\eta_t$ is i.i.d. Gaussian. Since Chowdhury and Gopalan's results hold for any $\xi$-sub-Gussian noise that is measurable based with respect to the prior observations, this assumption can be easily relaxed.
**Assumptions on the arms.** So far we considered a set of arms that is $(a)$ in $\mathbb{R}^d$, $(b)$ fixed for all $t$, and $(c)$ finite. Relaxing $(a)$ is easy, since we do not make any assumption beyond boundedness on the kernel function $k$ and there are many bounded kernel function for non-Euclidean spaces, e.g., strings or graphs. Relaxing $(b)$ is trivial, we just need to embed the changing arm sets as they are provided, and store and re-embed previously selected arms as necessary. The per-step time complexity will now depend on the size of the set of arms available at each step. Relaxing $(c)$ is straightforward from a theoretical perspective, but has varying computational consequences. In particular, looking at the proof of Theorem 1, the guarantees on $\widetilde{u}_t$ hold for all $\mathcal{H}$ and in Theorem 2, we only require that the maximum $\widetilde{\mathbf{x}}_{t+1} \triangleq \arg\max_{\mathbf{x} \in \mathcal{A}} \max_{\mathbf{w} \in \widetilde{C}_t} \phi(\mathbf{x})^\mathsf{T} \mathbf{w}$ is returned. Therefore, at least from the regret point of view, everything holds also for infinite $\mathcal{A}$. However, while the inner maximization over $\widetilde{C}_t$ can be solved in closed form for a fixed $\mathbf{x}$, the same cannot be said of the search over $\mathcal{A}$. If the designer can provide an efficient method to solve the maximization problem on an infinite $\mathcal{A}$, e.g., linear bandit optimization over compact subsets or $\mathbb{R}^d$, then all BKB guarantees apply.

## Appendix B. Properties of the posterior variance

For simplicity and completeness we provide known statements regarding the posterior variance $\sigma_t^2(\cdot)$. While most of these hold for generic RLS, we will adapt them to our notation.

**Proposition 4 (Calandriello et al., 2017a)** *For the posterior variance, we have that*

$$\frac{1}{\kappa^2/\lambda + 1}\sigma_{t-1}^2(\widetilde{\mathbf{x}}_t) \leq \frac{1}{\sigma_{t-1}^2(\widetilde{\mathbf{x}}_t) + 1}\sigma_{t-1}^2(\widetilde{\mathbf{x}}_t) \leq \sigma_t^2(\widetilde{\mathbf{x}}_t) \leq \sigma_{t-1}^2(\widetilde{\mathbf{x}}_t).$$

**Proof** The leftmost inequality follows from $\kappa^2/\lambda \geq \sigma_0^2(x)$ and $\sigma_a^2(x) \geq \sigma_b^2(x), \forall a \leq b$, the others are are by Calandriello et al., 2017a. ∎

**Proposition 5 (Hazan et al., 2006; Calandriello et al., 2017b)** *The effective dimension $d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)$ is upperbounded as*

$$\begin{aligned}
d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T) &\triangleq \mathrm{Tr}(\mathbf{K}_T(\mathbf{K}_T + \lambda\mathbf{I})^{-1}) = \sum\nolimits_{t=1}^{T} \sigma_T^2(\widetilde{\mathbf{x}}_t) \\
&\overset{(1)}{\leq} \sum\nolimits_{t=1}^{T} \sigma_t^2(\widetilde{\mathbf{x}}_t) \\
&\overset{(2)}{\leq} \log\det(\mathbf{K}_T/\lambda + \mathbf{I}) \\
&\overset{(3)}{\leq} \mathrm{Tr}(\mathbf{K}_T(\mathbf{K}_T + \lambda\mathbf{I})^{-1})\Big(1 + \log\Big(\frac{\|\mathbf{K}_T\|}{\lambda} + 1\Big)\Big).
\end{aligned}$$

**Proof** Inequality (1) is due to Proposition 4, inequality (2) is due to Hazan et al. (2006), and Inequality (3) is due to Calandriello et al. (2017b). ∎

# Appendix C. Proof of Theorem 1

Let $B_t$ be the unfavorable event where the guarantees of Theorem 1 do not hold. Our goal is to prove that $B_t$ happens at most with probability $\delta$ uniformly for all $t \in [T]$.

## C.1. Notation

In the following we refer to $\mathbf{\Phi}(\widetilde{\mathbf{X}}_t)$ as $\mathbf{\Phi}_t$, $\widetilde{\mathbf{\Phi}}(\widetilde{\mathbf{X}}_t)$ as $\widetilde{\mathbf{\Phi}}_t$ and $\phi(\widetilde{\mathbf{x}}_t)$ as $\phi_t$. When the subscript is clear from the context, we omit it. Since we leverage several results of Calandriello et al. (2017b), we start with some additional notation. First we extend our notation for the subset $\mathcal{S}_t$ to include a possible reweighing of the inducing points. We denote with $\mathcal{S}_t \triangleq \{(\phi_j, s_j)\}_{j=1}^{m_t}$, a *weighted* subset, i.e., a weighted *dictionary*, of columns from $\mathbf{\Phi}_t$, with positive weights $s_j > 0$ that must be appropriately chosen. Now, denote with $i_j \in [t]$, the index of the sample $\phi_j$ as a column in $\mathbf{\Phi}_t$. Using a standard approach (Alaoui and Mahoney, 2015), we choose $s_j \triangleq 1/\sqrt{\widetilde{p}_{t,i_j}}$, where $\widetilde{p}_{t,i} \triangleq \overline{q}\widetilde{\sigma}_{t-1}^2(\widetilde{\mathbf{x}}_i)$ is the probability[5] used by Algorithm 1 when sampling $\phi_{i_j}$ from $\mathbf{\Phi}_t$.

Let $\mathbf{S}_t \in \mathbb{R}^{t \times t}$ be the diagonal matrix with $q_{t,i}/\sqrt{\widetilde{p}_{t,i}}$ on the diagonal, where $q_{t,i}$ are the $\{0,1\}$ random variables selected by Algorithm 1. Then, we can see that

$$\sum_{j=1}^{m_t} \frac{1}{\widetilde{p}_{t,i_j}} \phi_{i_j} \phi_{i_j}^\mathsf{T} = \sum_{i=1}^t \frac{q_{t,i}}{\widetilde{p}_{t,i}} \phi_i \phi_i^\mathsf{T} = \mathbf{\Phi}_t \mathbf{S}_t \mathbf{S}_t^\mathsf{T} \mathbf{\Phi}_t^\mathsf{T}. \tag{11}$$

Calandriello et al. (2017a) define $\mathcal{S}_t$ to be an $\varepsilon$-accurate dictionary of $\mathbf{\Phi}_t$ if it satisfies

$$(1-\varepsilon)\mathbf{\Phi}_t\mathbf{\Phi}_t^\mathsf{T} - \varepsilon\lambda\mathbf{I} \preceq \mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^\mathsf{T}\mathbf{\Phi}_t^\mathsf{T} \preceq (1+\varepsilon)\mathbf{\Phi}_t\mathbf{\Phi}_t^\mathsf{T} + \varepsilon\lambda\mathbf{I}. \tag{12}$$

We can also now fully define the projection operator at time $t$ (see Section 4.1 for more details) as

$$\mathbf{P}_t \triangleq \mathbf{\Phi}_t\mathbf{S}_t(\mathbf{S}_t^\mathsf{T}\mathbf{\Phi}_t^\mathsf{T}\mathbf{\Phi}_t\mathbf{S}_t)^+\mathbf{S}_t^\mathsf{T}\mathbf{\Phi}_t^\mathsf{T},$$

which is the projection matrix spanned by the dictionary.

## C.2. Event decomposition

We decompose Theorem 1 into an accuracy part, i.e., $\mathcal{S}_t$ must induce accurate $\widetilde{\sigma}_t$, and an efficiency part, i.e., $m_t \leq d_{\text{eff}}(t)$. We also the accuracy of $\widetilde{\sigma}_t$ to the definition of $\varepsilon$-accuracy.

**Lemma 6** *Let $\alpha \triangleq \frac{1+\varepsilon}{1-\varepsilon}$. If $\mathcal{S}_t$ is $\varepsilon$-accurate w.r.t. $\mathbf{\Phi}_t$, then*

$$\mathbf{A}_t/\alpha \preceq \widetilde{\mathbf{A}}_t \preceq \alpha\mathbf{A}_t \quad and \quad \sigma_t^2(\mathbf{x})/\alpha \leq \min\{\widetilde{\sigma}_t^2(\mathbf{x}), 1\} \leq \alpha\sigma_t^2(\mathbf{x}) \text{ for all } \mathbf{x} \in \mathcal{A}.$$

**Proof** Inverting the bound in Equation 12 and using the fact that $\mathbf{P}_t\mathbf{\Phi}_t\mathbf{S}_t = \mathbf{\Phi}_t\mathbf{S}_t$, we get

$$\mathbf{P}_t\mathbf{\Phi}_t\mathbf{\Phi}_t^\mathsf{T}\mathbf{P}_t \preceq \frac{1}{1-\varepsilon}(\mathbf{P}_t\mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^\mathsf{T}\mathbf{\Phi}_t^\mathsf{T}\mathbf{P}_t + \varepsilon\lambda\mathbf{P}_t) \preceq \frac{1}{1-\varepsilon}(\mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^\mathsf{T}\mathbf{\Phi}_t^\mathsf{T} + \varepsilon\lambda\mathbf{P}_t)$$

$$\preceq \frac{1}{1-\varepsilon}((1+\varepsilon)\mathbf{\Phi}_t\mathbf{\Phi}_t^\mathsf{T} + \varepsilon\lambda\mathbf{I} + \varepsilon\lambda\mathbf{P}_t) \preceq \frac{1+\varepsilon}{1-\varepsilon}\left(\mathbf{\Phi}_t\mathbf{\Phi}_t^\mathsf{T} + \frac{2\varepsilon}{1+\varepsilon}\lambda\mathbf{I}\right).$$

---

5. Note that $\widetilde{p}_{t,i}$ might be larger than 1, but with a small abuse of notation and without the loss of generality we still refer to it as a probability.

Repeating the same process for the other side, we obtain

$$\frac{1-\varepsilon}{1+\varepsilon}\left(\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}} - \frac{2\varepsilon}{1-\varepsilon}\lambda\mathbf{I}\right) \preceq \mathbf{P}_t\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}}\mathbf{P}_t \preceq \frac{1+\varepsilon}{1-\varepsilon}\left(\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}} + \frac{2\varepsilon}{1+\varepsilon}\lambda\mathbf{I}\right).$$

Applying the above to $\widetilde{\mathbf{A}}_t$, we get

$$\widetilde{\mathbf{A}}_t = \mathbf{P}_t\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}}\mathbf{P}_t + \lambda\mathbf{I} \succeq \frac{1-\varepsilon}{1+\varepsilon}\left(\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}} - \frac{2\varepsilon}{1-\varepsilon}\lambda\mathbf{I}\right) + \lambda\mathbf{I} = \frac{1-\varepsilon}{1+\varepsilon}(\boldsymbol{\Phi}_t\boldsymbol{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I}) = \frac{1-\varepsilon}{1+\varepsilon}\mathbf{A}_t,$$

which can again be applied on the other side to obtain our result. To prove the accuracy of the approximate posterior variance $\widetilde{\sigma}_t^2(\mathbf{x}_i)$ we simply apply the definition to get

$$\frac{1-\varepsilon}{1+\varepsilon}\overbrace{\boldsymbol{\phi}_i^{\mathsf{T}}\mathbf{A}_t\boldsymbol{\phi}_i}^{\sigma_t^2(\mathbf{x}_i)} \preceq \overbrace{\boldsymbol{\phi}_i^{\mathsf{T}}\widetilde{\mathbf{A}}_t\boldsymbol{\phi}_i}^{\widetilde{\sigma}_t^2(\mathbf{x}_i)} \preceq \frac{1+\varepsilon}{1-\varepsilon}\overbrace{\boldsymbol{\phi}_i^{\mathsf{T}}\mathbf{A}_t\boldsymbol{\phi}_i}^{\sigma_t^2(\mathbf{x}_i)}.$$

$\blacksquare$

Using Lemma 6, we decompose our unfavorable event $B_t \triangleq A_t \cup E_t$, where $A_t$ is the event where $\mathcal{S}_t$ is not $\varepsilon$-accurate w.r.t. $\boldsymbol{\Phi}_t$ and $E_t$ is the event where $m_t$ is much larger than $d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_t)$. We now further decompose the event $A_t$ as

$$A_t = (A_t \cap A_{t-1}) \cup (A_t \cap A_{t-1}^{\complement})$$
$$\subseteq A_{t-1} \cup (A_t \cap A_{t-1}^{\complement}) = A_0 \cup \left(\bigcup_{s=1}^{t}(A_s \cap A_{s-1}^{\complement})\right) = \bigcup_{s=1}^{t}(A_s \cap A_{s-1}^{\complement}),$$

where $A_0$ is the empty event since $\boldsymbol{\Phi}_0$ is empty and it is well approximated by the empty $\mathcal{S}_0$. Moreover, we simplify a part of the expression by noting

$$B_t = A_t \cup E_t = A_t \cup (E_t \cap A_{t-1}^{\complement}) \cup (E_t \cap A_{t-1}) \subseteq A_t \cup A_{t-1} \cup (E_t \cap A_{t-1}^{\complement}),$$

which will help us when bounding the event $E_t$, where we will directly act as if $A_t$ does not hold. Putting it all together, we get

$$\bigcup_{t=1}^{T} B_t = \bigcup_{t=1}^{T}(A_t \cup E_t) \subseteq \bigcup_{t=1}^{T}\left(A_t \cup A_{t-1} \cup (E_t \cap A_{t-1}^{\complement})\right)$$
$$= \left(\bigcup_{t=1}^{T} A_t\right) \cup \left(\bigcup_{t=1}^{T}(E_t \cap A_{t-1}^{\complement})\right) = \left(\bigcup_{t=1}^{T} A_t\right) \cup \left(\bigcup_{t=1}^{T}(E_t \cap A_{t-1}^{\complement})\right)$$
$$\subseteq \left(\bigcup_{t=1}^{T}\left(\bigcup_{s=1}^{t}(A_s \cap A_{s-1}^{\complement})\right)\right) \cup \left(\bigcup_{t=1}^{T}(E_t \cap A_{t-1}^{\complement})\right)$$
$$= \left(\bigcup_{t=1}^{T}(A_t \cap A_{t-1}^{\complement})\right) \cup \left(\bigcup_{t=1}^{T}(E_t \cap A_{t-1}^{\complement})\right).$$

## C.3. Bounding $\Pr(A_t \cap A_{t-1}^{\mathsf{C}})$

We now bound the probability of event $A_t \cap A_{t-1}^{\mathsf{C}}$. In our first step, we formally define $A_t$ using Equation 12. In particular, we rewrite the $\varepsilon$-accuracy condition as

$$(1-\varepsilon)\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} - \varepsilon\lambda\mathbf{I} \preceq \mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^{\mathsf{T}}\mathbf{\Phi}_t^{\mathsf{T}} \preceq (1+\varepsilon)\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \varepsilon\lambda\mathbf{I}$$

$$\iff -\varepsilon(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I}) \preceq \mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^{\mathsf{T}}\mathbf{\Phi}_t^{\mathsf{T}} - \mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} \preceq \varepsilon(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})$$

$$\iff -\varepsilon\mathbf{I} \preceq (\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-1/2}(\mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^{\mathsf{T}}\mathbf{\Phi}_t^{\mathsf{T}} - \mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}})(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-1/2} \preceq \varepsilon\mathbf{I}$$

$$\iff \|(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-1/2}(\mathbf{\Phi}_t\mathbf{S}_t\mathbf{S}_t^{\mathsf{T}}\mathbf{\Phi}_t^{\mathsf{T}} - \mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}})(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-1/2}\| \leq \varepsilon,$$

where $\|\cdot\|$ is the spectral norm. We now focus on the last reformulation and frame it as a random matrix concentration question in RKHS $\mathcal{H}$. Let $\psi_{t,i} \triangleq (\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-\frac{1}{2}}\phi_i$ and $\mathbf{\Psi}_t \triangleq \mathbf{\Phi}_t(\mathbf{\Phi}_t^{\mathsf{T}}\mathbf{\Phi}_t + \lambda\mathbf{I})^{-\frac{1}{2}} = [\psi_{t,1}, \ldots, \psi_{t,t}]^{\mathsf{T}}$, and define the operator $\mathbf{G}_{t,i} \triangleq \left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}$. Then we rewrite $\varepsilon$-accuracy as

$$\left\|(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{\Phi}_t(\mathbf{S}_t\mathbf{S}_t^{\mathsf{T}} - \mathbf{I})\mathbf{\Phi}_t^{\mathsf{T}}(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-\frac{1}{2}}\right\| = \left\|\sum_{i=1}^{t}\left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\right\| = \left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \leq \varepsilon,$$

and the event $A_t$ as the event where $\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon$, Note that this reformulation exploits the fact that $q_{t,i} = 0$ encodes the column that are not selected in $\mathcal{S}_t$ (see Equation 11). To study this random object, we begin by defining the filtration $\mathcal{F}_t \triangleq \{q_{s,i}, \eta_s\}_{s=1}^{t}$ at time $t$ containing all the randomness coming from the construction of the various $\mathcal{S}_s$ and the noise on the function $\eta_t$. In particular, note that the $\{0,1\}$ r.v. $q_{t,i}$ used by Algorithm 1 are not necessarily Bernoulli r.v.s, since the probability $\widetilde{p}_{t,i}$ used to select $0$ or $1$ is itself random. However, they become well defined Bernoulli when conditioned on $\mathcal{F}_{t-1}$. Let $\mathbb{I}\{\cdot\}$ indicates the indicator function of an event. We have that

$$\Pr(A_t \cap A_{t-1}^{\mathsf{C}}) = \Pr\left(\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon \cap \left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon\right)$$

$$= \mathop{\mathbb{E}}_{\mathcal{F}_t}\left[\mathbb{I}\left\{\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon \cap \left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon\right\}\right]$$

$$= \mathop{\mathbb{E}}_{\mathcal{F}_{t-1}}\left[\mathop{\mathbb{E}}_{\eta_t, \{q_{t,i}\}}\left[\mathbb{I}\left\{\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon \cap \left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon\right\}\,\middle|\,\mathcal{F}_{t-1}\right]\right]$$

$$= \mathop{\mathbb{E}}_{\mathcal{F}_{t-1}}\left[\mathop{\mathbb{E}}_{\{q_{t,i}\}}\left[\mathbb{I}\left\{\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon \cap \left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon\right\}\,\middle|\,\mathcal{F}_{t-1}\right]\right],$$

where the last passage is due to the fact that $\mathbf{G}_{t,i}$ is independent from $\eta_t$. Next, notice that conditioned on $\mathcal{F}_{t-1}$, the event $A_{t-1}^{\mathsf{C}}$ becomes deterministic, and we can restrict our expectations to the outcomes where $\left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon$,

$$\Pr(A_t \cap A_{t-1}^{\mathsf{C}}) = \mathop{\mathbb{E}}_{\mathcal{F}_{t-1}:\left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\|\leq\varepsilon}\left[\mathop{\mathbb{E}}_{\{q_{t,i}\}}\left[\mathbb{I}\left\{\left\|\sum_{i=1}^{t}\mathbf{G}_{t,i}\right\| \geq \varepsilon\right\}\,\middle|\,\mathcal{F}_{t-1}\right]\right].$$

Moreover, conditioned on $\mathcal{F}_{t-1}$ all the $q_{t,i}s$ become independent r.v., and we are able to use the following result of Tropp (2015).

**Proposition 7** *Let* $\mathbf{G}_1, \ldots, \mathbf{G}_n$ *be a sequence of independent self-adjoint random operators such that* $\mathbb{E}[\mathbf{G}_i] = 0$ *and* $\|\mathbf{G}_i\| \leq R$ *a.s. Denote* $\sigma^2 = \left\|\sum_{i=1}^{t} \mathbb{E}[\mathbf{G}_i^2]\right\|$. *Then, for any* $\varepsilon \geq 0$,

$$\Pr\left(\left\|\sum_{i=1}^{t} \mathbf{G}_i\right\| \geq \varepsilon\right) \leq 4t \exp\left(\frac{\varepsilon^2/2}{\sigma^2 + R\varepsilon/3}\right).$$

We begin by computing the mean of $\mathbf{G}_{t,i}$,

$$
\begin{aligned}
\mathbb{E}_{q_{t,i}}[\mathbf{G}_{t,i} \mid \mathcal{F}_{t-1}] &= \mathbb{E}_{q_{t,i}}\left[\left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \,\middle|\, \mathcal{F}_{t-1}\right] \\
&= \left(\frac{\mathbb{E}_{q_{t,i}}[q_{t,i} \mid \mathcal{F}_{t-1}]}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}} = \left(\frac{\widetilde{p}_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}} = \mathbf{0},
\end{aligned}
$$

where we use the fact that $\widetilde{p}_{t,i}$ is fixed conditioned on $\mathcal{F}_{t-1}$ and it is the (conditional) expectation of $q_{t,i}$. Since $\mathbf{G}$ is zero-mean, we can use Proposition 7. First, we find $R$ and for that, we upper bound

$$\|\mathbf{G}_{t,i}\| = \left\|\left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\right\| \leq \left|\left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)\right|\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\| \leq \frac{1}{\widetilde{p}_{t,i}}\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\|.$$

Note that due to the definition of $\psi_{t,i}$,

$$\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\| = \psi_{t,i}^{\mathsf{T}}\psi_{t,i} = \phi_i^{\mathsf{T}}(\mathbf{\Phi}_t\mathbf{\Phi}_t^{\mathsf{T}} + \lambda\mathbf{I})^{-1}\phi_i = \sigma_t^2(\widetilde{\mathbf{x}}_i).$$

Moreover, we are only considering outcomes of $\mathcal{F}_{t-1}$ where $\left\|\sum_{i=1}^{t}\mathbf{G}_{t-1,i}\right\| \leq \varepsilon$, which implies that $\mathcal{S}_{t-1}$ is $\varepsilon$-accurate, and by Lemma 6 we have that $\widetilde{\sigma}_{t-1}(\widetilde{\mathbf{x}}_i) \geq \sigma_{t-1}(\widetilde{\mathbf{x}}_i)/\alpha$. Finally, due to Proposition 4, we have $\sigma_{t-1}(\widetilde{\mathbf{x}}_i) \geq \sigma_t(\widetilde{\mathbf{x}}_i)$. Putting this all together we can bound

$$\frac{1}{\widetilde{p}_{t,i}}\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\| = \frac{1}{\overline{q}\widetilde{\sigma}_{t-1}(\widetilde{\mathbf{x}}_i)}\sigma_t(\widetilde{\mathbf{x}}_i) \leq \frac{\alpha}{\overline{q}} \triangleq R.$$

For the variance term, we expand

$$
\begin{aligned}
\sum_{i=1}^{t}\mathbb{E}_{q_{t,i}}[\mathbf{G}_{t,i}^2 \mid \mathcal{F}_{t-1}] &= \sum_{i=1}^{t}\mathbb{E}_{q_{t,i}}\left[\left(\frac{q_{t,i}}{\widetilde{p}_{t,i}} - 1\right)^2 \,\middle|\, \mathcal{F}_{t-1}\right]\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \\
&= \sum_{i=1}^{t}\left(\mathbb{E}_{q_{t,i}}\left[\frac{q_{t,i}^2}{\widetilde{p}_{t,i}^2} \,\middle|\, \mathcal{F}_{t-1}\right] - \mathbb{E}_{q_{t,i}}\left[2\frac{q_{t,i}}{\widetilde{p}_{t,i}} \,\middle|\, \mathcal{F}_{t-1}\right] + 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \\
&= \sum_{i=1}^{t}\left(\mathbb{E}_{q_{t,i}}\left[\frac{q_{t,i}}{\widetilde{p}_{t,i}^2} \,\middle|\, \mathcal{F}_{t-1}\right] - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\psi_{t,i}\psi_{t,i}^{\mathsf{T}} = \sum_{i=1}^{t}\left(\mathbb{E}_{q_{t,i}}\left[\frac{q_{t,i}}{\widetilde{p}_{t,i}^2} \,\middle|\, \mathcal{F}_{t-1}\right] - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \\
&= \sum_{i=1}^{t}\left(\frac{1}{\widetilde{p}_{t,i}} - 1\right)\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \preceq \sum_{i=1}^{t}\frac{1}{\widetilde{p}_{t,i}}\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\|\psi_{t,i}\psi_{t,i}^{\mathsf{T}} \preceq \sum_{i=1}^{t}R\psi_{t,i}\psi_{t,i}^{\mathsf{T}},
\end{aligned}
$$

where we used the fact that $q_{t,i}^2 = q_{t,i}$ and $\mathbb{E}_{q_{t,i}}[q_{t,i}|\mathcal{F}_{t-1}] = \widetilde{p}_{t,i}$. We can now bound this quantity as

$$\left\|\sum_{i=1}^{t}\mathbb{E}_{q_{t,i}}[\mathbf{G}_{t,i}^2 \mid \mathcal{F}_{t-1}]\right\| \leq \left\|\sum_{i=1}^{t}R\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\right\| = R\left\|\sum_{i=1}^{t}\psi_{t,i}\psi_{t,i}^{\mathsf{T}}\right\| = R\|\mathbf{\Psi}_t^{\mathsf{T}}\mathbf{\Psi}_t\| \leq R \triangleq \sigma^2.$$

Therefore, we have $\sigma^2 = R$ and $R = 1/\overline{q}$. Now, applying Proposition 7 and a union bound we conclude the proof.

## C.4. Bounding $\Pr(E_t \cap A_{t-1}^{\complement})$

We will use the following concentration for independent Bernoulli random variables.

**Proposition 8 ([Calandriello et al., 2017a], App. D.4)**  *Let $\{q_s\}_{s=1}^t$ be independent Bernoulli random variables, each with success probability $p_s$, and let $d = \sum_{s=1}^t p_s \geq 1$ be their sum. Then,[6]*

$$\mathbb{P}\left(\sum_{s=1}^t q_s \geq 3d\right) \leq \exp\{-3d(3d - (\log(3d) + 1))\} \leq \exp\{-2d\}.$$

We now rigorously define event $E_t$ as the event where

$$\sum_{i=1}^t q_{t,i} \geq 3\alpha(1 + \kappa^2/\lambda)\log(t/\delta)\sum_{i=1}^t \sigma_t^2(\widetilde{\mathbf{x}}_i) = 3\alpha(1 + \kappa^2/\lambda)d_{\text{eff}}(\lambda, \widetilde{\mathbf{X}}_t)\log(t/\delta).$$

Once again, we use conditioning and in particular,

$$\Pr(E_t \cap A_t^{\complement}) = \mathop{\mathbb{E}}_{\mathcal{F}_{t-1}:\left\|\sum_{i=1}^t \mathbf{G}_{t-1,i}\right\| \leq \varepsilon}\left[\mathop{\mathbb{E}}_{\{q_{t,i}\}}\left[\mathbb{I}\left\{\sum_{i=1}^t q_{t,i} \geq 3\alpha(1+\kappa^2/\lambda)\log(t/\delta)\sum_{i=1}^t \sigma_t^2(\widetilde{\mathbf{x}}_i)\right\}\bigg| \mathcal{F}_{t-1}\right]\right].$$

Conditioned on $\mathcal{F}_{t-1}$ the r.v. $q_{t,i}$ becomes independent Bernoulli with probability $\widetilde{p}_{t,i} \triangleq \overline{q}\widetilde{\sigma}_{t-1}(\widetilde{\mathbf{x}}_i)$. Since we restrict the outcomes to $A_{t-1}^{\complement}$, we can exploit Lemma 6 and the guarantees of $\varepsilon$-accuracy to bound $\widetilde{p}_{t,i} \leq \alpha\sigma_{t-1}^2(\widetilde{\mathbf{x}}_i)$. Then, we use Proposition 4 to bound $\sigma_{t-1}^2(\widetilde{\mathbf{x}}_i) \leq (1 + \kappa^2/\lambda)\sigma_t^2(\widetilde{\mathbf{x}}_i)$. Therefore, $q_{t,i}$ are conditionally independent Bernoulli with probability at most $\overline{q}(1 + \kappa^2/\lambda)\sigma_t^2(\widetilde{\mathbf{x}}_i)$. Applying a simple stochastic dominance argument and Proposition 8 gets the needed statement.

## Appendix D. Proof of Theorem 2

Following [Abbasi-Yadkori et al.] (2011), we divide the proof in two parts, first bounding the approximate confidence ellipsoid, and then bounding the regret.

### D.1. Bounding the confidence ellipsoid

We begin by proving an intermediate result regarding the confidence ellipsoid.

**Theorem 9**  *Under the same assumptions as Theorem 2 with probability at least $1 - \delta$ and for all $t \geq 0$, $\mathbf{w}_\star$ lies in the set*

$$\widetilde{C}_t \triangleq \left\{\mathbf{w} : \|\mathbf{w} - \widetilde{\mathbf{w}}_t\|_{\widetilde{\mathbf{A}}_t} \leq \widetilde{\beta}_t\right\}$$

*with*

$$\widetilde{\beta}_t \triangleq 2\xi\sqrt{\alpha\log(\kappa^2 t)\left(\sum_{s=1}^t \widetilde{\sigma}_t^2(\mathbf{x}_s)\right) + \log\left(\frac{1}{\delta}\right)} + \left(1 + \frac{1}{\sqrt{1-\varepsilon}}\right)\sqrt{\lambda}F.$$

---

6. This is a simple variant of the Chernoff bound where the Bernoulli random variables are not identically distributed.

**Proof** For simplicity, we omit the subscript $t$. We begin by noticing that

$$
\begin{aligned}
(\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\widetilde{\mathbf{w}} - \mathbf{w}_\star) &= (\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal \mathbf{y} - \mathbf{w}_\star) \\
&= (\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal (\boldsymbol{\Phi}\mathbf{w}_\star + \eta - \mathbf{w}_\star) \\
&= (\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\underbrace{\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{w}_\star - \mathbf{w}_\star}_{\text{bias}}) + (\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}^{1/2} \underbrace{\widetilde{\mathbf{A}}^{-1/2}\widetilde{\boldsymbol{\Phi}}^\intercal \eta}_{\text{variance}}.
\end{aligned}
$$

**Bounding the bias.** We first focus on the first term, which is difficult to analyze due to the mismatch $\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}$. We have that

$$
\begin{aligned}
\widetilde{\mathbf{A}}(\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{w}_\star - \mathbf{w}_\star) &= \widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{w}_\star - \widetilde{\boldsymbol{\Phi}}^\intercal \widetilde{\boldsymbol{\Phi}}\mathbf{w}_\star - \lambda\mathbf{w}_\star \\
&= \widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star + \widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{P}\mathbf{w}_\star - \widetilde{\boldsymbol{\Phi}}^\intercal \widetilde{\boldsymbol{\Phi}}_\star^{\mathbf{w}} - \lambda\mathbf{w}_\star \\
&= \widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star - \lambda\mathbf{w}_\star.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
(\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{w}_\star - \mathbf{w}_\star) &= (\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star - \lambda(\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \mathbf{w}_\star \\
&\leq \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}} \left( \|\widetilde{\mathbf{A}}^{-1/2}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star\| + \lambda\|\mathbf{w}^*\|_{\widetilde{\mathbf{A}}^{-1}} \right) \\
&\leq \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}} \left( \|\widetilde{\mathbf{A}}^{-1/2}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star\| + \tfrac{\lambda}{\sqrt{\lambda}}\|\mathbf{w}^*\| \right).
\end{aligned}
$$

Then, we have that

$$
\begin{aligned}
\|\widetilde{\mathbf{A}}^{-1/2}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\mathbf{w}_\star\| &\leq \|\widetilde{\mathbf{A}}^{-1/2}\widetilde{\boldsymbol{\Phi}}^\intercal\| \|\boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\| \|\mathbf{w}_\star\| \\
&\leq \sqrt{\lambda_{\max}(\widetilde{\boldsymbol{\Phi}}\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal)}\sqrt{\lambda_{\max}(\boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})^2\boldsymbol{\Phi}^\intercal)}\|\mathbf{w}_\star\|.
\end{aligned}
$$

It is easy to see that

$$
\lambda_{\max}(\widetilde{\boldsymbol{\Phi}}\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal) = \lambda_{\max}(\widetilde{\boldsymbol{\Phi}}(\widetilde{\boldsymbol{\Phi}}^\intercal \widetilde{\boldsymbol{\Phi}} + \lambda\mathbf{I})^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal) \leq 1.
$$

To bound the other term we use the following result by Calandriello and Rosasco (2018).

**Proposition 10** *If $\mathcal{S}_t$ is $\varepsilon$-accurate w.r.t. $\boldsymbol{\Phi}_t$, then*

$$
\mathbf{I} - \mathbf{P}_t \preceq \mathbf{I} - \boldsymbol{\Phi}_t \mathbf{S}_t (\mathbf{S}_t^\intercal \boldsymbol{\Phi}_t^\intercal \boldsymbol{\Phi}_t \mathbf{S}_t + \lambda\mathbf{I})^{-1}\mathbf{S}_t^\intercal \boldsymbol{\Phi}_t^\intercal \preceq \frac{\lambda}{1 - \varepsilon}(\boldsymbol{\Phi}_t \boldsymbol{\Phi}_t^\intercal + \lambda\mathbf{I})^{-1}.
$$

Since from Theorem 1, we have that $\mathcal{S}_t$ is $\varepsilon$-accurate, by Theorem 10, we have that

$$
\boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})^2 \boldsymbol{\Phi}^\intercal = \boldsymbol{\Phi}(\mathbf{I} - \mathbf{P})\boldsymbol{\Phi}^\intercal \preceq \frac{\lambda}{1 - \varepsilon}\boldsymbol{\Phi}(\boldsymbol{\Phi}^\intercal \boldsymbol{\Phi} + \lambda\mathbf{I})^{-1}\boldsymbol{\Phi}^\intercal \preceq \frac{\lambda}{1 - \varepsilon}\mathbf{I}.
$$

Putting it all together, we obtain

$$
(\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\intercal \widetilde{\mathbf{A}}(\widetilde{\mathbf{A}}^{-1}\widetilde{\boldsymbol{\Phi}}^\intercal \boldsymbol{\Phi}\mathbf{w}_\star - \mathbf{w}_\star) \leq \left( 1 + \frac{1}{\sqrt{1 - \varepsilon}} \right) \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}}\sqrt{\lambda}\|\mathbf{w}_\star\|.
$$

**Bounding the variance.** We use the the following self-normalized martingale concentration inequality by Abbasi-Yadkori et al. (2011). It can be trivially extended to RKHSs in the case of finite sets such as our $\mathcal{A}$. Note that if the reader is interested in infinite sets, Chowdhury and Gopalan (2017) provide a generalization with slightly worse constants.

**Proposition 11 ([Abbasi-Yadkori et al., 2011](#))** *Let $\{\mathcal{F}_t\}_{t=0}^{\infty}$ be a filtration, let $\{\eta_t\}_{t=1}^{\infty}$ be a real-valued stochastic process such that $\eta_t$ is $\mathcal{F}_t$-measurable and zero-mean $\xi$-subgaussian; let $\{\mathbf{\Phi}_t\}_{t=1}^{\infty}$ be an $\mathcal{H}$-valued stochastic process such that $\mathbf{\Phi}_t$ is $\mathcal{F}_{t-1}$-measurable, and let $\mathbf{I}$ be the identity operator on $\mathcal{H}$. For any $t \geq 1$, define*

$$\mathbf{A}_t = \mathbf{\Phi}_t^\mathsf{T}\mathbf{\Phi}_t + \lambda\mathbf{I} \quad and \quad \mathbf{V}_t = \mathbf{\Phi}_t^\mathsf{T}\eta_t.$$

*Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t \geq 0$,*

$$\|\mathbf{V}_t\|_{\mathbf{A}_t^{-1}}^2 \leq 2\xi^2 \log\left(\frac{\det(\mathbf{A}_t/\lambda)}{\delta}\right).$$

Recalling the definition of $\alpha \geq 1$ from Theorem [1](#), we reformulate

$$
\begin{aligned}
(\widetilde{\mathbf{w}} - \mathbf{w}_\star)^\mathsf{T}\widetilde{\mathbf{A}}^{1/2}\widetilde{\mathbf{A}}^{-1/2}\widetilde{\mathbf{\Phi}}\eta &\leq \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}}\|\widetilde{\mathbf{\Phi}}\eta\|_{\widetilde{\mathbf{A}}^{-1}} \\
&= \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}}\|\widetilde{\mathbf{\Phi}}^\mathsf{T}\eta\|_{(\widetilde{\mathbf{\Phi}}^\mathsf{T}\widetilde{\mathbf{\Phi}}+\lambda\mathbf{I})^{-1}} \\
&= \|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}}\|\widetilde{\mathbf{\Phi}}^\mathsf{T}\eta/\lambda\|_{(\widetilde{\mathbf{\Phi}}^\mathsf{T}\widetilde{\mathbf{\Phi}}/\lambda+\mathbf{I})^{-1}}.
\end{aligned}
$$

We now make a remark that requires temporal notation. Note that we cannot directly apply Theorem [11](#) to $\widetilde{\mathbf{\Phi}}_t\eta_t = \mathbf{P}_t\mathbf{\Phi}_t\eta_t$. In particular, for $s < t$ we have that $\widetilde{\mathbf{\Phi}}_s\eta_s = \mathbf{P}_t\mathbf{\Phi}_s\eta_s$ is not $\mathcal{F}_{s-1}$ measurable, since $\mathbf{P}_t$ depends on all randomness up to time $t$. However, since $\mathbf{P}_t$ is always a projection matrix we know that the variance of the projected process is bounded by the variance of the original process, in particular,

$$
\begin{aligned}
\|\widetilde{\mathbf{\Phi}}^\mathsf{T}\eta/\lambda\|_{(\widetilde{\mathbf{\Phi}}^\mathsf{T}\widetilde{\mathbf{\Phi}}/\lambda+\mathbf{I})^{-1}} &= \sqrt{\eta^\mathsf{T}\widetilde{\mathbf{\Phi}}(\widetilde{\mathbf{\Phi}}^\mathsf{T}\widetilde{\mathbf{\Phi}}/\lambda+\mathbf{I})^{-1}\widetilde{\mathbf{\Phi}}^\mathsf{T}\eta/\lambda} = \sqrt{\eta^\mathsf{T}\widetilde{\mathbf{\Phi}}\widetilde{\mathbf{\Phi}}^\mathsf{T}(\widetilde{\mathbf{\Phi}}\widetilde{\mathbf{\Phi}}^\mathsf{T}/\lambda+\mathbf{I})^{-1}\eta/\lambda} \\
&\overset{(a)}{=} \sqrt{\eta^\mathsf{T}(\mathbf{I} - \lambda(\widetilde{\mathbf{\Phi}}\widetilde{\mathbf{\Phi}}^\mathsf{T}/\lambda+\mathbf{I})^{-1})\eta/\lambda} = \sqrt{\eta^\mathsf{T}(\mathbf{I} - \lambda(\mathbf{\Phi}\mathbf{P}\mathbf{\Phi}^\mathsf{T}/\lambda+\mathbf{I})^{-1})\eta/\lambda} \\
&\overset{(b)}{\leq} \sqrt{\eta^\mathsf{T}(\mathbf{I} - \lambda(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T}/\lambda+\mathbf{I})^{-1})\eta/\lambda} \overset{(c)}{=} \|\mathbf{\Phi}^\mathsf{T}\eta/\lambda\|_{(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}/\lambda+\mathbf{I})^{-1}},
\end{aligned}
$$

where in $(a)$ we added and subtracted $\lambda\mathbf{I}$ from $\widetilde{\mathbf{\Phi}}\widetilde{\mathbf{\Phi}}^\mathsf{T}$, in $(b)$ we used the fact that $\|\mathbf{P}\| \leq 1$ for all projection matrices, and in $(c)$ we reversed the reformulation from $(a)$. We can finally use Theorem [11](#) to obtain

$$
\begin{aligned}
\|\mathbf{\Phi}^\mathsf{T}\eta/\lambda\|_{(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}/\lambda+\mathbf{I})^{-1}} &\leq \sqrt{2\xi^2 \log(\mathrm{Det}(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}/\lambda+\mathbf{I})/\delta)} \\
&= \sqrt{2\xi^2 \log(\mathrm{Det}(\mathbf{A}/\lambda)/\delta)}.
\end{aligned}
$$

While above is a valid bound on the radius of the confidence interval, it is still not satisfactory. In particular, we can use Sylvester's identity to reformulate

$$\log\det(\mathbf{A}/\lambda) = \log\det(\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}/\lambda + \mathbf{I}) = \log\det(\mathbf{\Phi}\mathbf{\Phi}^\mathsf{T}/\lambda + \mathbf{I}) = \log\det(\mathbf{K}/\lambda + \mathbf{I}).$$

Computing the radius would require constructing the matrix $\mathbf{K} \in \mathbb{R}^{t \times t}$ and this is way too expensive. Instead, we obtain a cheap but still a small enough upper bound as follows,

$$
\begin{aligned}
\log \det(\mathbf{K}_t/\lambda + \mathbf{I}) &\leq \operatorname{Tr}(\mathbf{K}_t(\mathbf{K}_t + \lambda \mathbf{I})^{-1})(1 + \log(\|\mathbf{K}_t\| + 1)) \\
&\leq \operatorname{Tr}(\mathbf{K}_t(\mathbf{K}_t + \lambda \mathbf{I})^{-1})(1 + \log(\operatorname{Tr} \mathbf{K}_t + 1)) \\
&\leq \operatorname{Tr}(\mathbf{K}_t(\mathbf{K}_t + \lambda \mathbf{I})^{-1})(1 + \log(\kappa^2 t + 1)) \\
&= (1 + \log(\kappa^2 t + 1)) \sum_{s=1}^{t} \sigma_t^2(\mathbf{x}_s) \\
&\leq \alpha(1 + \log(\kappa^2 t + 1)) \sum_{s=1}^{t} \widetilde{\sigma}_t^2(\mathbf{x}_s) \\
&\leq 2\alpha \log(\kappa^2 t) \sum_{s=1}^{t} \widetilde{\sigma}_t^2(\mathbf{x}_s),
\end{aligned}
$$

where $\widetilde{\sigma}_t^2(\mathbf{x}_s)$ can be computed efficiently and it is actually already done by the algorithm at every step! Putting it all together, we get that

$$
\begin{aligned}
\|\widetilde{\mathbf{w}} - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}} &\leq 2\xi \sqrt{\alpha \log(\kappa^2 t) \left( \sum_{s=1}^{t} \widetilde{\sigma}_t^2(\mathbf{x}_s) \right) + \log(1/\delta)} + \left(1 + \frac{1}{\sqrt{1 - \varepsilon}}\right)\sqrt{\lambda}\|\mathbf{w}_\star\| \\
&\leq 2\xi \sqrt{\alpha \log(\kappa^2 t) \left( \sum_{s=1}^{t} \widetilde{\sigma}_t^2(\mathbf{x}_s) \right) + \log(1/\delta)} + \left(1 + \frac{1}{\sqrt{1 - \varepsilon}}\right)\sqrt{\lambda}F \triangleq \widetilde{\beta}_t.
\end{aligned}
$$

∎

### D.2. Bounding the regret

The regret analysis is straightforward. Assume that $\mathbf{w}_\star \in \widetilde{C}_t$ is satisfied (i.e., the event from Theorem 9 holds) and remember that by the definition, $\phi_t \triangleq \arg\max_{\mathbf{x}_i in \mathcal{A}} \max_{\mathbf{w} \in \widetilde{\mathbf{C}}_t} \phi_i^\mathsf{T}\mathbf{w}$. We also define $\overline{\mathbf{w}}_{t,i} \triangleq \arg\max_{\mathbf{w} \in \widetilde{\mathbf{C}}_t} \phi_i^\mathsf{T}\mathbf{w}$ as the auxiliary vector which encodes the optimistic behaviour of the algorithm. With a slight abuse of notation, we also use $\star$ as a subscript to indicate the (unknown) index of the optimal arm, so that $\overline{\mathbf{w}}_{t,\star} \triangleq \arg\max_{\mathbf{w} \in \widetilde{\mathbf{C}}_t} \phi_\star^\mathsf{T}\mathbf{w}$. Since $\mathbf{w}_\star \in \widetilde{C}_t$, we have that

$$
\phi_t^\mathsf{T}\overline{\mathbf{w}}_{t,t} \geq \phi_\star^\mathsf{T}\overline{\mathbf{w}}_{t,\star} \geq \phi_\star^\mathsf{T}\mathbf{w}_*.
$$

We can now bound the instantaneous regret $r_t$ as

$$
\begin{aligned}
r_t = \phi_\star^\mathsf{T}\mathbf{w}_\star - \phi_t^\mathsf{T}\mathbf{w}_\star &\leq \phi_t^\mathsf{T}\overline{\mathbf{w}}_{t,t} - \phi_t^\mathsf{T}\mathbf{w}_\star \\
&= \phi_t^\mathsf{T}(\overline{\mathbf{w}}_{t,t} - \widehat{\mathbf{w}}_t) + \phi_t^\mathsf{T}(\widehat{\mathbf{w}}_t - \mathbf{w}_\star) \\
&= \phi_t^\mathsf{T}\widetilde{\mathbf{A}}_t^{-1/2}\widetilde{\mathbf{A}}_t^{1/2}(\overline{\mathbf{w}}_{t,t} - \widehat{\mathbf{w}}_t) + \phi_t^\mathsf{T}\widetilde{\mathbf{A}}_{t-1}^{-1/2}\widetilde{\mathbf{A}}_t^{1/2}(\widehat{\mathbf{w}}_t - \mathbf{w}_\star) \\
&\leq \sqrt{\phi_t^\mathsf{T}\widetilde{\mathbf{A}}_t^{-1}\phi_t} \left( \|\overline{\mathbf{w}}_{t,t} - \widehat{\mathbf{w}}_t\|_{\widetilde{\mathbf{A}}_t} + \|\widehat{\mathbf{w}}_t - \mathbf{w}_\star\|_{\widetilde{\mathbf{A}}_t} \right) \\
&\leq 2\widetilde{\beta}_t \sqrt{\phi_t^\mathsf{T}\widetilde{\mathbf{A}}_t^{-1}\phi_t}.
\end{aligned}
$$

Summing over $t$ and taking the max over $\widetilde{\beta}_t$, we get

$$R_t \leq 2\widetilde{\beta}_T \sum_{t=1}^{T} \sqrt{\phi_t^\intercal \widetilde{\mathbf{A}}_t^{-1} \phi_t} \leq 2\widetilde{\beta}_T \sqrt{T} \sqrt{\sum_{t=1}^{T} \phi_t^\intercal \widetilde{\mathbf{A}}_t^{-1} \phi_t} \leq 2\widetilde{\beta}_T \sqrt{T} \sqrt{\alpha \sum_{t=1}^{T} \phi_t^\intercal \mathbf{A}_t^{-1} \phi_t}.$$

We can now use once again Proposition 5 to obtain

$$R_T \leq 2\widetilde{\beta}_T \sqrt{\alpha T \sum_{t=1}^{T} \phi_t^\intercal \mathbf{A}_t^{-1} \phi_t} = 2\widetilde{\beta}_T \sqrt{\alpha T \sum_{t=1}^{T} \sigma_t^2(\widetilde{\mathbf{x}}_t)} \leq 2\widetilde{\beta}_T \sqrt{2\alpha T d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T) \log(\kappa^2 T)}.$$

We can also further upper bound $\widetilde{\beta}_T$ as

$$\widetilde{\beta}_T = 2\xi \sqrt{\alpha \log(\kappa^2 T)\left(\sum_{s=1}^{T} \widetilde{\sigma}_t^2(\mathbf{x}_s)\right) + \log(1/\delta)} + \left(1 + \frac{1}{\sqrt{1-\varepsilon}}\right)\sqrt{\lambda}F$$

$$\leq 2\xi \sqrt{\alpha^2 \log(\kappa^2 T)\left(\sum_{s=1}^{T} \sigma_t^2(\mathbf{x}_s)\right) + \log(1/\delta)} + \left(1 + \frac{1}{\sqrt{1-\varepsilon}}\right)\sqrt{\lambda}F$$

$$\leq 2\xi\alpha\sqrt{d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T) + \log(1/\delta)} + \left(1 + \frac{1}{\sqrt{1-\varepsilon}}\right)\sqrt{\lambda}F.$$

Putting it together, we obtain

$$R_T \leq 2\left(2\xi\alpha\sqrt{d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T) + \log(1/\delta)}\right)\sqrt{2\alpha T d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T)}$$

$$+ 2\left(\left(1 + \frac{1}{\sqrt{1-\varepsilon}}\right)\sqrt{\lambda}F\right)\sqrt{2\alpha T d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T)}$$

$$\leq 2\xi(2\alpha)^{3/2}\left(d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T) + \log(1/\delta)\right) + 2\left(2\sqrt{\alpha}\sqrt{\lambda}F\right)\sqrt{2\alpha T d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T)}$$

$$\leq 2(2\alpha)^{3/2}\left(\sqrt{T}\xi d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T) + \sqrt{T}\log(1/\delta) + \sqrt{T\lambda F^2 d_{\mathrm{eff}}(\lambda, \widetilde{\mathbf{X}}_T)\log(\kappa^2 T)}\right).$$