ON KERNEL RL WITHOUT OPTIMISTIC CLOSURE

Anonymous authors

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

018

019

021

024

025

026027028

029

031

032

033

037

038

040 041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

We study episodic reinforcement learning with a kernel (RKHS) structure on stateaction pairs. Previous optimistic analyses in this case either pay a data-dependent covering-number penalty that can grow with time and undermine no-regret guarantees, or it assumes a strong "optimistic closure" condition requiring all optimistic proxies to lie in a fixed state-RKHS ball. We take a different approach that removes the covering-number dependence without invoking optimistic closure. Our analysis builds a uniform confidence bound, derived via conditional mean embeddings, that holds simultaneously for all proxy value functions within a bounded state-RKHS class. We introduce **KOVI-Proj**, an optimistic value-iteration scheme that explicitly projects the optimistic proxy back into the state-RKHS ball at every step, ensuring that the uniform bound applies throughout the learning process. Under a restricted Bellman-embedding assumption (bounded conditional mean embeddings), KOVI-Proj enjoys a high-probability no-regret guarantee whose rate is governed by the task horizon and the kernel's information gain. When the optimal value function lies in the chosen state-RKHS ball (realizability), the regret is sublinear; in the agnostic case, an explicit approximation term reflects the best RKHS approximation error. Overall, this work provides a new pathway to no-regret kernel RL that is strictly weaker than optimistic closure and avoids covering-number penalties. Numerical experiments validate our claims.

1 Introduction

Kernel-based function approximation are known to provide an interesting link between linear models and the behavior of infinitely wide neural networks. However, obtaining sharp, no-regret (let alone order-optimal) guarantees in kernel-based reinforcement learning (RL) remains challenging (Vakili, 2024). Previous optimistic analyses in RKHSs typically follow one of the two approaches: (i) apply a union bound over a *data-dependent*, *evolving* class of optimistic value proxies, thereby incurring a covering-number penalty that can scale in the order of $\Omega(\sqrt{T})$ and spoil no-regret for common kernels (e.g., kernelized optimistic LSVI: Least Squares Value Iteration (Yang et al., 2020)); or (ii) assume a strong *optimistic closure* property stating that *every* optimistic proxy already contained in a fixed state-RKHS ball (as found in CME-based optimistic RL) (Chowdhury & Oliveira, 2023). The former is statistically loose; the latter is structurally strong and not obviously aligned with standard optimistic constructions.

We take a different approach based on uniform concentration without covering. The key observation is that for any V in a state RKHS \mathcal{H}_{ℓ} , the Bellman image $[P_h V]$ can be written as an inner product

$$[P_h V](z) = \langle \mu_h(z), V \rangle_{\mathcal{H}_{\ell}},$$

where $\mu_h:\mathcal{Z}\to\mathcal{H}_\ell$ is the conditional mean embedding (CME) of the next-state distribution (Muandet et al., 2017b; Song et al., 2013; Muandet et al., 2017a). When μ_h is contained in an appropriate vector-valued RKHS over \mathcal{Z} with bounded norm, the map $V\mapsto [P_hV]$ is a bounded linear operator from $(\mathcal{H}_\ell,\|\cdot\|_{\mathcal{H}_\ell})$ to $(\mathcal{H}_k,\|\cdot\|_{\mathcal{H}_k})$ (Carmeli et al., 2010). This viewpoint lets us control, via a single vector-valued regression problem, the Bellman images $[P_hV]$ for all V in the state ball $\{V:\|V\|_{\mathcal{H}_\ell}\leq B\}$ simultaneously, yielding a uniform kernel-ridge confidence bound with no data-dependent covering (leveraging information-gain / elliptical-potential tools standard in kernel bandits) (Chowdhury & Gopalan, 2017). Importantly, we enforce the bounded-norm property algorithmically by projecting the optimistic proxy value onto the state-RKHS ball each step. Experimental results show promise of the proposed method, and tends to shows lower regret than the baselines.

Constibutions of this paper is listed as follows:

- (1) **Restricted Bellman-embedding assumption (RBE).** We use a mild assumption under which the CME μ_h belongs to the vector-valued RKHS on \mathcal{Z} with kernel k I and norm at most U. This is strictly weaker than optimistic closure (which presumes *all* optimistic proxies already lie in a fixed state-RKHS ball), and it is natural under standard CME regularity (Muandet et al., 2017b; Carmeli et al., 2010).
- (2) Uniform confidence without covering. We prove a high-probability bound

$$\sup_{\|V\|_{\mathcal{H}_{\ell}} \le B} \ \left| [P_h V](z) - \hat{f}_{h,n}^{V}(z) \right| \ \le \ \beta_{h,n} \, \sigma_{h,n}(z),$$

holding for all z and all V in the ball, where $\widehat{f}_{h,n}^V$ is the kernel-ridge predictor trained on labels V(s') and $\sigma_{h,n}$ is the posterior standard deviation. The multiplier satisfies

$$\beta_{h,n} = BU + \frac{B\sigma}{\sqrt{\rho}}\sqrt{2\gamma(n,\rho) + 2\log(1/\delta)},$$

depending on the ball radius B, operator norm U, sub-Gaussian scale σ , ridge parameter $\rho > 0$, and (regularized) information gain $\gamma(n,\rho)$ -but it does not depend on any covering number of the proxy class (Chowdhury & Gopalan, 2017).

- (3) **KOVI-Proj: Kernel-Optimistic Value Iteration with projection.** We propose a practical optimistic method that (i) performs kernel-ridge backups to estimate $[P_hV]$, (ii) adds an uncertainty bonus $\beta_{h,n}\sigma_{h,n}$, and (iii) the method projects optimistic value proxy onto the state-RKHS ball (with clipping), thereby guaranteeing $\|V\|_{\mathcal{H}_\ell} \leq B$ and placing all proxies within the scope of the uniform bound above.
- (4) No-regret guarantee. Under the realizability $(V_h^* \in \mathcal{H}_{\ell}(B))$ assumption, the proposed KOVI-Proj method attains

$$R(T) = \tilde{\mathcal{O}}\left(H^2 B\left(U + \frac{\sigma}{\sqrt{\rho}} \sqrt{\gamma(HT, \rho)}\right) \sqrt{T \gamma(HT, \rho)}\right),$$

which is sublinear for kernels with sublinear information gain (e.g., Matérn/Squared-Exponential under regularization) Chowdhury & Gopalan (2017). In the agnostic case $(V_h^* \notin \mathcal{H}_\ell(B))$, we add an explicit approximation term of order $HT \, \varepsilon_B$, where $\varepsilon_B := \max_h \sup_{\|V\| \le B} \|V_h^* - V\|_{\infty}$, and we show how a slowly growing $B = B_T$ balances both terms to remain o(T).

2 PROBLEM SETUP AND ASSUMPTIONS

We consider an episodic MDP M = (S, A, H, P, r) with horizon $H \in \mathbb{N}$. Let $\mathcal{Z} = S \times A$. For step $h \in [H]$, transition kernel is $P_h(\cdot | z)$ on S is unknown. We take rewards $r_h : \mathcal{Z} \to [0, 1]$ to be known and deterministic for clarity; for a policy π and step h, we have

$$Q_h^{\pi}(z) = r_h(z) + \mathbb{E}_{s' \sim P_h(\cdot|z)}[V_{h+1}^{\pi}(s')], \qquad V_h^{\pi}(s) = \max_{a \in \mathcal{A}} Q_h^{\pi}(s, a), \qquad V_{H+1}^{\pi} \equiv 0$$

The per-episode regret will be measured against optimal value V_1^* as follows

$$R(T) := \sum_{t=1}^{T} (V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t})).$$

RKHS structure on \mathcal{Z} and on \mathcal{S} . Let $k: \mathcal{Z} \times \mathcal{Z} \to \mathbb{R}$ be a positive-definite kernel with RKHS $(\mathcal{H}_k, \|\cdot\|_{\mathcal{H}_k})$ and $k(z,z) \leq \kappa_k^2$. Let $\ell: \mathcal{S} \times \mathcal{S} \to \mathbb{R}$ be a positive-definite kernel with RKHS $(\mathcal{H}_\ell, \|\cdot\|_{\mathcal{H}_\ell})$ and let $\ell(s,s) \leq \kappa_\ell^2$. We plan to use kernel ridge regression (KRR) on \mathcal{Z} and consider proxy value functions $V: \mathcal{S} \to \mathbb{R}$ in \mathcal{H}_ℓ .

The following assumption is novel to our work, but is inspired by the conditional mean embedding literature (Muandet et al., 2017b) and the theory of vector-valued RKHSs (Carmeli et al., 2010).

¹The extension to unknown (possibly stochastic) rewards can be handled with an additional KRR estimator and a union bound; see the discussion section.

Assumption 2.1 (Restricted Bellman-embedding (RBE)). For each $h \in [H]$ assume that there exists a conditional mean embedding $\mu_h : \mathcal{Z} \to \mathcal{H}_\ell$ such that

$$[P_hV](z) := \mathbb{E}_{s' \sim P_h(\cdot|z)}[V(s')] = \langle \mu_h(z), V \rangle_{\mathcal{H}_{\ell}} \quad \text{for all } V \in \mathcal{H}_{\ell} \text{ and all } z \in \mathcal{Z} \quad (1)$$

and μ_h belongs to the vector-valued RKHS over \mathcal{Z} with operator-valued kernel $K(z,z') = k(z,z') I_{\mathcal{H}_{\ell}}$, with $\|\mu_h\|_{\mathcal{H}_k \otimes \mathcal{H}_{\ell}} \leq U$.

Remark 2.2. Assumption 2.1 is a standard conditional-mean-embedding (CME) property written in a vector-valued RKHS: $z \mapsto \mu_h(z)$ is an \mathcal{H}_ℓ -valued function whose inner product with V equals the Bellman image $[P_hV]^2$. This assumption is strictly weaker than optimistic closure (which would require that all optimistic proxies lie in a fixed state-RKHS ball in advance) and is implied by common regularity conditions under which CMEs exist with bounded norm.

Data model at step h. On observing a transition (z_i, s_i') , we define \mathcal{H}_{ℓ} -valued observation $\phi_i = \phi(s_i')$, where $\phi : \mathcal{S} \to \mathcal{H}_{\ell}$ is the canonical feature map of ℓ . Then we have

$$\mathbb{E}[\phi_i \mid z_i] = \mu_h(z_i), \qquad \varepsilon_i := \phi_i - \mu_h(z_i) \in \mathcal{H}_{\ell},$$

so $\{\varepsilon_i\}$ is a martingale-difference sequence in Hilbert space \mathcal{H}_{ℓ} . We assume $\|\varepsilon_i\|_{\mathcal{H}_{\ell}} \leq \kappa_{\ell}$ almost surely and that ε_i is σ -sub-Gaussian in \mathcal{H}_{ℓ} conditionally on the past. For any $V \in \mathcal{H}_{\ell}$, we define scalar labels

$$y_i^{(V)} := V(s_i') = \langle \phi_i, V \rangle_{\mathcal{H}_{\ell}} = \langle \mu_h(z_i), V \rangle_{\mathcal{H}_{\ell}} + \xi_i^{(V)}, \qquad \xi_i^{(V)} := \langle \varepsilon_i, V \rangle_{\mathcal{H}_{\ell}},$$

so that $\xi_i^{(V)}$ is conditionally sub-Gaussian with proxy variance proportional to $\|V\|_{\mathcal{H}_\ell}^2$ (and $|\xi_i^{(V)}| \le \kappa_\ell \|V\|_{\mathcal{H}_\ell}$ almost surely).

Kernel ridge predictors and variances. Given n observations at step h with design points $z_{1:n}$, Gram matrix $K_n = [k(z_i, z_j)]_{i,j=1}^n$, regularization $\rho > 0$, and labels $\boldsymbol{y}^{(V)} = [V(s_1'), \dots, V(s_n')]^\top$, we define

$$\widehat{f}_{h,n}^{V}(z) = k_n(z)^{\top} (K_n + \rho I)^{-1} \boldsymbol{y}^{(V)}, \qquad \sigma_{h,n}^{2}(z) = k(z,z) - k_n(z)^{\top} (K_n + \rho I)^{-1} k_n(z),$$
 (2)

where $k_n(z) = [k(z, z_1), \dots, k(z, z_n)]^{\top}$. We also use the (regularized) information gain Chowdhury & Gopalan (2017); Srinivas et al. (2010)

$$\gamma(n,\rho) := \frac{1}{2} \log \det (I + \rho^{-1} K_n)$$

Note that all quantities here carry a step index h, which we will suppress when it will be clear from context.

3 A UNIFORM CONFIDENCE BOUND FOR ALL V WITH $\|V\|_{\mathcal{H}_\ell} \leq B$

The next proposition will be a key algebraic identity: it trades uniformity over an *uncountable* class of scalar predictors for a single bound on a *vector-valued* kernel ridge estimator.

Proposition 3.1 (Scalar KRR = inner product with a vector-valued KRR). Fix a step h and data $\{(z_i, s_i')\}_{i=1}^n$. Define the \mathcal{H}_ℓ -valued (vector) KRR estimator

$$\widehat{\mu}_n(z) := \sum_{i=1}^n \alpha_i(z) \, \phi(s_i') \in \mathcal{H}_\ell, \qquad \alpha(z) := (K_n + \rho I)^{-1} k_n(z).$$

Then for every $V \in \mathcal{H}_{\ell}$ and $z \in \mathcal{Z}$, we have

$$\widehat{f}_{h,n}^{V}(z) = \langle \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}}$$

Proof. By equation 2, $\widehat{f}_{h,n}^{V}(z) = \sum_{i=1}^{n} \alpha_i(z) V(s_i') = \sum_{i=1}^{n} \alpha_i(z) \langle \phi(s_i'), V \rangle = \langle \sum_{i=1}^{n} \alpha_i(z) \phi(s_i'), V \rangle$. For detailed proof, see C.1 in Appendix.

²See, e.g., Muandet et al. (2017) for CMEs and Carmeli-De Vito-Toigo (2008) for vector-valued RKHS foundations.

Proposition 3.1 reduces uniform control over all scalar targets V to control of the *vector-valued* estimation error $\|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_\ell}$. The next lemma extends self-normalized kernel concentration to the vector-valued CME.

Lemma 3.2 (Vector-valued kernel ridge concentration). Suppose Assumption 3.2 holds, $k(z,z) \le \kappa_k^2$, $\ell(s,s) \le \kappa_\ell^2$. Let $\rho > 0$ and define $\sigma_{h,n}(\cdot)$ by equation 2. Then for any $\delta \in (0,1)$, with probability at least $1 - \delta$, simultaneously for all $z \in \mathcal{Z}$,

$$\|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}} \leq \left(\sqrt{\rho} \, U \, + \, \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right) \sigma_{h,n}(z).$$

Proof sketch; full details in Appendix E. Write the vector-valued regression as $\phi_i = \mu(z_i) + \varepsilon_i$, with $\varepsilon_i \in \mathcal{H}_\ell$ a martingale difference, conditionally σ -sub-Gaussian and $\|\varepsilon_i\|_{\mathcal{H}_\ell} \leq \kappa_\ell$ a.s. The KRR error decomposes as

$$\mu(z) - \widehat{\mu}_n(z) = \underbrace{\mu(z) - \prod_n \mu(z)}_{\text{bias}} - \underbrace{\Phi^{\top}(K_n + \rho I)^{-1}k_n(z)}_{\text{noise}},$$

where Π_n is the Tikhonov projector in the vector-valued RKHS induced by kI, and $\Phi: \mathbb{R}^n \to \mathcal{H}_\ell$ maps $b \mapsto \sum_i b_i \phi_i$. The bias is controlled by the standard RKHS interpolation inequality: $\|\mu(z) - \Pi_n \mu(z)\|_{\mathcal{H}_\ell} \le \sqrt{\rho} \|\mu\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell} \sigma_{h,n}(z) \le \sqrt{\rho} \, U \, \sigma_{h,n}(z)$. For the noise, a Hilbert-space self-normalized bound (made explicit in Appendix E) yields $\|\Phi^\top (K_n + \rho I)^{-1} k_n(z)\|_{\mathcal{H}_\ell} \le \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}} \, \sigma_{h,n}(z)$ with probability at least $1-\delta$. Summing the two contributions gives the claim.

Combining Proposition 3.1 with Lemma 3.2 yields the desired uniform scalar bound.

Theorem 3.3 (Uniform CI for all $||V||_{\mathcal{H}_{\ell}} \leq B$). Under the conditions of Lemma 3.2, for any B > 0 and $\delta \in (0,1)$, with probability at least $1 - \delta$, for all $V \in \mathcal{H}_{\ell}$ with $||V||_{\mathcal{H}_{\ell}} \leq B$ and all $z \in \mathcal{Z}$,

$$|[P_h V](z) - \widehat{f}_{h,n}^V(z)| \le \beta_{n,\delta} \, \sigma_{h,n}(z), \qquad \beta_{n,\delta} := B\left(\sqrt{\rho} \, U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right).$$

Proof. By equation 1 and Proposition 3.1, $[P_hV](z) - \widehat{f}_{h,n}^V(z) = \langle \mu(z) - \widehat{\mu}_n(z), V \rangle$; Cauchy-Schwarz and Lemma 3.2 complete the proof. Detailed proof in Appendix 3.3.

Remark 3.4 (Notational simplification used later). For simplicity in subsequent sections (e.g., in the algorithmic confidence radius and regret display), we may absorb the $\sqrt{\rho}$ factor into the constant by defining $U' := \sqrt{\rho} U$ and writing $\beta_{n,\delta} = B \left(U' + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho)} + 2\log(1/\delta) \right)$. We keep Lemma 3.2 in the explicit $\sqrt{\rho}$ form for clarity.

4 ALGORITHM: KOVI-PROJ (KERNEL-OPTIMISTIC VALUE ITERATION WITH PROJECTION)

We now describe our algorithn. We maintain a separate KRR model for each step $h \in [H]$. Let $\mathcal{D}_{h,t-1} = \{(z_{h,\tau},s_{h+1,\tau})\}_{\tau=1}^{n_{h,t-1}}$ denote the transitions collected so far at step h before episode t, with $n_{h,t-1} = |\mathcal{D}_{h,t-1}|$. At the start of episode t, set $V_{H+1,t} \equiv 0$ and perform a backward pass for $h = H, H-1, \ldots, 1$:

- 1. **Kernel-ridge backup.** Using equation 2 with design points $z_{h,1:n_{h,t-1}}$ and labels $y_{\tau} = V_{h+1,t}(s_{h+1,\tau})$, compute the predictor $\widehat{f}_{h,t}^{V_{h+1,t}}(\cdot)$ and its posterior deviation $\sigma_{h,t}(\cdot)$.
- 2. Confidence radius and optimism. Let $\delta \in (0,1)$ be the overall failure probability. Define the per-step confidence multiplier (cf. Theorem 3.3 and Remark 3.4)

$$\beta_{h,t} \ := \ B\Big(\sqrt{\rho}\, U \ + \ \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{2HT}{\delta}}\,\Big),$$

and form the optimistic action-value

$$\widetilde{Q}_{h,t}(z) := r_h(z) + \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \beta_{h,t} \, \sigma_{h,t}(z).$$
 (3)

3. Optimistic value and projection onto the state-RKHS ball. Let

$$\widetilde{V}_{h,t}(s) := \max_{a \in \mathcal{A}} \widetilde{Q}_{h,t}(s,a),$$

then obtain $V_{h,t}$ by projecting $\widetilde{V}_{h,t}$ onto $\{V \in \mathcal{H}_{\ell} : \|V\|_{\mathcal{H}_{\ell}} \leq B\}$ (with range clipping) under a reference measure ν on \mathcal{S} :

$$V_{h,t} \in \arg\min_{V \in \mathcal{H}_{\ell}} \left\{ \|V - \widetilde{V}_{h,t}\|_{L^{2}(\nu)} : \|V\|_{\mathcal{H}_{\ell}} \le B, \ 0 \le V \le H - h + 1 \right\}. \tag{4}$$

Interaction. Within episode t, act greedily with respect to $\widetilde{Q}_{h,t}$: pick $a_{h,t} \in \arg\max_{a \in \mathcal{A}} \widetilde{Q}_{h,t}(s_{h,t},a)$, observe $s_{h+1,t} \sim P_h(\cdot \mid s_{h,t},a_{h,t})$, and append $(z_{h,t},s_{h+1,t})$ to $\mathcal{D}_{h,t}$. Proceed to step h+1.

Projection in finite dimension (The QP form). In practice, we instantiate equation 4 via the representer theorem. Let $\{\bar{s}_j\}_{j=1}^{m_h}$ be a set of anchor states for step h (e.g., the distinct states observed at step h so far, optionally augmented by a cover of \mathcal{S}). Denote the Gram matrix $L_h = [\ell(\bar{s}_i, \bar{s}_j)]_{i,j}$ and the vector of target values $v_{h,t} = [\widetilde{V}_{h,t}(\bar{s}_j)]_{j=1}^{m_h}$. Seeking $V \in \mathcal{H}_\ell$ of the form $V(s) = \sum_{j=1}^{m_h} \alpha_j \, \ell(s, \bar{s}_j)$, the projection reduces to the convex quadratic program (although standard, a proof is in appendix I)

$$\min_{\alpha \in \mathbb{R}^{m_h}} \frac{1}{m_h} \| L_h \alpha - v_{h,t} \|_2^2 \quad \text{s.t.} \quad \alpha^\top L_h \alpha \le B^2, \qquad 0 \le (L_h \alpha)_j \le H - h + 1 \ \forall j. \quad (5)$$

The optimizer yields $V_{h,t}(s) = \sum_{j=1}^{m_h} \alpha_j \, \ell(s, \bar{s}_j)$. Problem equation 5 is a small QP solvable in $\tilde{\mathcal{O}}(m_h^3)$ time per step; in our experiments we take m_h to be the number of unique states observed at step h (with optional down-sampling).

Remarks.

- (i) The confidence radius $\beta_{h,t}$ incorporates a union bound over all (h,t) via the $\log(2HT/\delta)$ term, ensuring the uniform event of Theorem 3.3 holds jointly for all steps and episodes.
- (ii) The projection step guarantees $||V_{h,t}||_{\mathcal{H}_{\ell}} \leq B$ and range constraints; thus *every* optimistic proxy used by the algorithm lies in the state-RKHS ball, placing it within the scope of the uniform confidence bound without any data-dependent covering.
- (iii) Choice of ν in equation 4 can be the empirical state distribution at step h or an exploratory cover over \mathcal{S} ; the finite-dimensional form equation 5 corresponds to taking ν uniform over the anchor set.
- (iv) If rewards are unknown and/or stochastic, one can learn \hat{r}_h via a separate KRR with its own confidence band and add it to equation 3 (with a union bound across reward and transition estimators).
- (v) For notational simplicity one may absorb $\sqrt{\rho}$ into U (Remark 3.4) and write $\beta_{h,t} = B(U' + \frac{\sigma}{\sqrt{\rho}}\sqrt{2\gamma(n_{h,t-1},\rho)} + 2\log(2HT/\delta))$ with $U' := \sqrt{\rho}U$.

5 REGRET ANALYSIS

We state the main guarantee under Assumption 2.1. The proof follows the optimistic value-iteration template, combining (i) the uniform confidence event from Theorem 3.3 enforced by the projection step, (ii) a standard telescoping decomposition, and (iii) an elliptical-potential (information-gain) bound summed across steps.

Theorem 5.1 (No-regret under RBE). Suppose Assumption 2.1 holds for all $h \in [H]$, rewards lie in [0,1], and $k(z,z) \le \kappa_k^2$, $\ell(s,s) \le \kappa_\ell^2$. Let $\rho \in (0,1]$ and let $\gamma(\cdot,\rho)$ be the regularized information gain of k on $\mathcal Z$ as in equation 2. Run KOVI-Proj with ball radius B and failure probability $\delta \in (0,1/T]$. Then with probability at least $1-\delta$, after T episodes,

$$R(T) \; \leq \; \tilde{\mathcal{O}}\!\!\left(H^2\,B\!\left(\sqrt{\rho}\,U + \tfrac{\sigma}{\sqrt{\rho}}\sqrt{\gamma(HT,\rho)}\right)\sqrt{T\,\gamma(HT,\rho)}\right) \; + \; HT\,\varepsilon_B,$$

where $\varepsilon_B := \max_{h \in [H]} \inf_{\|V\|_{\mathcal{H}_{\ell}} \leq B} \|V_h^* - V\|_{\infty}$. In particular, under realizability $(V_h^* \in \mathcal{H}_{\ell}(B))$ for all h) we have

$$R(T) = \tilde{\mathcal{O}}\left(H^2 B\left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{\gamma(HT, \rho)}\right) \sqrt{T \gamma(HT, \rho)}\right),$$

which is o(T) whenever $\gamma(n, \rho) = o(n)$.

Proof sketch; full details in Appendix D. Good event. By Theorem 3.3 with a union bound over all steps and episodes (the $\log(2HT/\delta)$ term inside $\beta_{h,t}$ in equation 3), there exists an event \mathcal{G} of probability at least $1-\delta$ such that, for all h,t and all $z \in \mathcal{Z}$,

$$[P_h V_{h+1,t}](z) \le \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \beta_{h,t} \sigma_{h,t}(z)$$

where $\beta_{h,t}$ is as in Section 4. The projection step guarantees $\|V_{h,t}\|_{\mathcal{H}_{\ell}} \leq B$, hence every optimistic proxy used by the algorithm lies within the scope of the uniform bound.

Optimism and telescoping. Define $\widetilde{Q}_{h,t}$ by equation 3 and $\widetilde{V}_{h,t}(s) = \max_a \widetilde{Q}_{h,t}(s,a)$. On \mathcal{G} and up to the agnostic term ε_B , a standard dynamic-programming induction yields $Q_h^*(z) \leq \widetilde{Q}_{h,t}(z)$ and hence $V_h^*(s) \leq \widetilde{V}_{h,t}(s)$. Therefore the per-episode regret telescopes as follows

$$R(T) \leq \sum_{t=1}^{T} \sum_{h=1}^{H} \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t}) \right) + HT \varepsilon_B$$

$$= \sum_{t=1}^{T} \sum_{h=1}^{H} \beta_{h,t} \, \sigma_{h,t}(z_{h,t}) + HT \varepsilon_B$$

Elliptical potential across steps. Let $n_{h,T}$ be the total number of transitions observed at step h by time T (then $\sum_h n_{h,T} = HT$). For each fixed h, the standard GP/RKHS potential argument gives $\sum_{t=1}^T \sigma_{h,t}(z_{h,t}) \leq \sqrt{2 n_{h,T} \gamma(n_{h,T},\rho)}$. We sum over h and apply Cauchy-Schwarz, to obtain

$$\sum_{t=1}^{T} \sum_{h=1}^{H} \sigma_{h,t}(z_{h,t}) \leq \sum_{h=1}^{H} \sqrt{2 \, n_{h,T} \, \gamma(n_{h,T}, \rho)} \leq \sqrt{2 \, \left(\sum_{h} n_{h,T} \right) \left(\sum_{h} \gamma(n_{h,T}, \rho) \right)} \, = \sqrt{2 \, HT \, \Gamma_{T}},$$

where $\Gamma_T := \sum_{h=1}^H \gamma(n_{h,T},\rho)$. Since the per-step Gram matrices are disjoint, Γ_T equals the information gain of the block-diagonal kernel on the stacked design and satisfies $\Gamma_T \leq \gamma(HT,\rho)$. Hence $\sum_{t,h} \sigma_{h,t}(z_{h,t}) \leq \sqrt{2\,HT\,\gamma(HT,\rho)}$.

Putting it together. Use $\beta_{h,t} \leq \tilde{\mathcal{O}}(B(\sqrt{\rho}\,U + \frac{\sigma}{\sqrt{\rho}}\sqrt{\gamma(HT,\rho)}))$ uniformly over h,t, multiply by $\sum_{t,h} \sigma_{h,t}(z_{h,t})$, and absorb polylogarithms to obtain the stated bound.

Remark 5.2 (On the H-dependence). The H^2 factor arises from the optimistic LSVI-style decomposition and coarse bounding of stepwise contributions. We expect sharper analysis (e.g., refined Bellman-error coupling or variance decomposition) could improve this to $H^{3/2}$ or even H, but we leave this for future work.

6 DISCUSSION AND COMPARISONS

Versus covering-number analyses. Theorem 3.3 yields a confidence multiplier of the form

$$\beta_{n,\delta} = B\left(\sqrt{\rho}U + \frac{\sigma}{\sqrt{\rho}}\sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right),$$

without any covering-number factor over the evolving proxy class. Intuitively, by estimating the conditional mean embedding μ_h once, we control all Bellman images $[P_hV]$ for V in the ball $\{V: \|V\|_{\mathcal{H}_\ell} \leq B\}$ via Cauchy-Schwarz:

$$\sup_{\|V\|_{\mathcal{H}_{\ell}} \leq B} \left| [P_h V](z) - \widehat{f}_{h,n}^V(z) \right| \leq B \|\mu_h(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}} \lesssim \beta_{n,\delta} \, \sigma_{h,n}(z),$$

as formalized by Lemma 3.2. This directly replaces the union-bound-over-a-cover step used in earlier kernel-RL analyses.

Versus optimistic closure. We *do not* assume that every optimistic proxy automatically lies in a fixed state-RKHS ball. Instead we *enforce* it by an explicit projection step (Section 4). Analytically, this is sufficient: it is the set of *actual* proxies used by the algorithm that needs to lie inside the uniform-confidence event of Theorem 3.3. Thus the projection step plays the role that *optimistic closure* previously assumed.

When does RBE hold? Assumption 2.1 requires that the CME map $\mu_h: \mathcal{Z} \to \mathcal{H}_\ell$ belongs to the vector-valued RKHS with kernel kI and $\|\mu_h\| \leq U$. This is natural when: (i) $z \mapsto P_h(\cdot|z)$ varies smoothly in a kernel mean sense (e.g., Hölder or Lipschitz in the MMD induced by ℓ), (ii) ℓ is bounded and universal (e.g., RBF on compact \mathcal{S}), and (iii) k is bounded on \mathcal{Z} . In such cases, conditional mean embeddings exist and admit finite RKHS norm. The constant U estimates the operator norm of the Bellman map $V \mapsto [P_h V]$ from $(\mathcal{H}_\ell, \|\cdot\|_{\mathcal{H}_\ell})$ to $(\mathcal{H}_k, \|\cdot\|_{\mathcal{H}_k})$.

Computational considerations. The projection step reduces to the QP in equation 5 with complexity $\tilde{\mathcal{O}}(m_h^3)$ per step, where m_h is the number of anchor states. In practice, m_h can be taken as the distinct observed states at step h (optionally sub-sampled) or a small cover; this keeps the overhead modest relative to KRR updates on \mathcal{Z} .

Agnostic setting. When $V_h^* \notin \mathcal{H}_\ell(B)$, the only degradation is the explicit $HT \, \varepsilon_B$ term in Theorem 5.1. For universal kernels, $\varepsilon_B \to 0$ as $B \to \infty$; choosing $B = B_T$ to grow slowly (e.g., $B_T = \tilde{\mathcal{O}}(\sqrt{\log T})$) balances approximation and estimation so that R(T) = o(T) whenever $\gamma(HT, \rho) = o(HT)$.

Relation to kernel bandits (H=1). For H=1, KOVI-Proj specializes to a GP/KRR-UCB scheme where the uniform CME bound recovers the familiar information-gain control of regret. Our analysis is consistent with recent refined bounds for GP-UCB and shows how the CME perspective naturally extends to multi-step RL.

Limitations and possible improvements. Our current regret bound scales as H^2 , inherited from an optimistic LSVI-style decomposition. Tighter coupling of stepwise Bellman errors or a variance-aware decomposition could plausibly reduce this to $H^{3/2}$ or H. Extending RBE examples and verifying U for broader kernel/state-action families, and integrating unknown rewards with joint confidence control, are also natural next steps.

BROADER IMPACT

This work proposes a CME-based uniformization mechanism for kernel RL that removes an obstacle to no-regret guarantees while relaxing structural assumptions (no optimistic closure). Broader impacts include more reliable kernelized RL with principled uncertainty quantification; as always, care is warranted when deploying RL systems in safety-critical settings.

7 LLM USAGE

LLM was used for polishing texts to rephrase and correct grammar.

- 8 EXPERIMENTS
- 9 NUMERICAL EXPERIMENT: 1D DOUBLE-WELL (QCQP PROJECTION, ABSORBING GOAL)

Setup. We consider the classical quartic *double-well* in 1D with overdamped Langevin dynamics and additive control:

$$x_{t+1} = x_t + \Delta t \left(x_t - x_t^3 + u_t \right) + \sigma \varepsilon_t, \qquad u_t \in \{ -u_0, 0, +u_0 \}, \ \varepsilon_t \sim \mathcal{N}(0, 1).$$

Table 1: Double–Well (H = 40, T = 100): final cumulative regret (mean over seeds) and SEM.

Algorithm	Final Cum. Regret (↓)	SEM
KOVI-Proj	93.287	4.308
KOVI0	118.236	0.085
Kernel-LSVI-ε	118.301	0.007

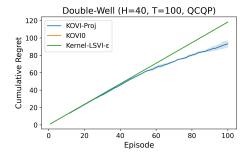
Episodes have horizon H=40. The goal set is an absorbing tube around x=+1 of radius τ ; the reward is one-shot hit+1 (upon first entry) minus a step penalty 0.01 each step, and the episode terminates on hit. This makes the benchmark V_1^* O(1) and aligned with the environment (details in the appendix).

Algorithms. We compare three methods: (i) **KOVI-Proj**, which performs backward optimistic value iteration with a kernel surrogate for $[P_hV]$ and *always* projects V_h by solving the QCQP

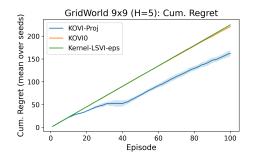
$$\min_{\alpha} \frac{1}{m} \|L\alpha - v_h\|_2^2 \quad \text{s.t.} \quad \alpha^{\top} L\alpha \leq B^2, \ 0 \leq (L\alpha)_j \leq H - h + 1,$$

(ii) **KOVI0**, the same but *without* the RKHS projection (ridge only), and (iii) **Kernel-LSVI-** ε , a non-optimistic KRR baseline with ε -greedy exploration. All methods warm-start with a few random episodes; we amortize planning with a plan-every-K schedule. We evaluate over K=100 episodes and three seeds.

Results. Figure 1a shows the mean cumulative regret (shaded: SEM) against episodes. KOVI-Proj learns substantially faster and attains markedly lower regret. Table 1 summarizes final cumulative regret (mean over seeds) and its SEM.



(a) Double–Well ($H=40,\,T=100$): mean cumulative regret vs. episode (QCQP projection always on). KOVI-Proj (blue) outperforms both the no–projection ablation KOVI0 (orange) and Kernel-LSVI- ε (green).



(b) GridWorld 9×9 (H = 5), 100 episodes: mean cumulative regret (shaded: SEM across seeds). Here KOVI-pro in blue has least growth in regret. Projection is QCQP.

Figure 1: Regret Plots for Double-well and Grid World.

Discussion. Three observations are consistent across seeds: (i) Level. KOVI-Proj lowers the final cumulative regret by about 21% relative to the non–projected optimistic ablation and the non-optimistic baseline. This reflects substantially higher hit probability of the absorbing goal within the H=40-step window. (ii) Rate. The slope of the regret curve is strictly smaller for KOVI-Proj across the training horizon, indicating faster value improvement per episode. (iii) Role of projection. Removing the RKHS ball + range constraints (KOVI0) collapses the optimism guarantee: the upper-confidence target \widetilde{Q}_h no longer reliably upper-bounds the Bellman image, leading to miscalibrated targets and markedly worse exploration. In contrast, the QCQP projection keeps value iterates within the feasible hypothesis set, preserving the UCB validity and translating into consistent goal-reaching behavior.

Table 2: GridWorld (9 \times 9, H = 5), 100 episodes summary metrics (mean over seeds).

Algorithm	Final Cum. Regret (↓)	Regret Slope / ep (\downarrow)	Mean Return (†)	SEM(Return)
KOVI-Proj	162.348	1.579	0.533	0.602
KOVI 0	220.928	2.221	-0.053	0.132
Kernel-LSVI- ε	224.968	2.248	-0.093	0.048

10 GRIDWORLD BENCHMARK

Environment. We use a 9×9 GridWorld (states $\{0,\dots,8\}^2$) with start at (0,0) and goal at (8,8). The horizon is H=5 per episode. Actions are $\{\text{up}, \text{right}, \text{down}, \text{left}\}$. With slip probability $p_{\text{slip}}=0.1$, the executed action is replaced uniformly at random. The reward is +1 upon entering the goal and -0.01 otherwise. We have RBF kernel over states with lengthscale $\ell=0.35$, product kernel for Q over state-action, KRR ridge $\lambda_Q=10^{-2}$ for Q, ridge $\lambda_V=10^{-3}$ for the ridge baseline, anchors placed on a stride-2 grid (m=25 anchors), UCB scale $\beta=0.8\sqrt{\log((mH+1)/\delta)}$ with $\delta=0.1$, and RKHS ball radius B=4.0 for the projection.

Algorithms. We compare (i) **KOVI-Proj** (QCQP projection for V_h enforcing $||V_h||_{\mathcal{H}_\ell} \leq B$ and $0 \leq V_h \leq H - h + 1$), (ii) **KOVI0** (same optimism, but V_h via ridge without constraints), and (iii) **Kernel-LSVI-** ε (non-optimistic KRR targets with ε -greedy; ε decays as in the code). At each episode, we perform a backward planning pass to update $\{V_h\}_{h=H}^1$ from replayed targets, then run one episode of interaction.

Metrics. We compute the optimal benchmark V_1^* by dynamic programming and report (i) mean cumulative regret over K=100 episodes, (ii) the least-squares per-episode regret slope, and (iii) mean return; all statistics are averaged over three seeds with SEM bands.

Results. Figure 1b shows the regret curves; Table 2 summarizes final numbers.

Discussion. KOVI-Proj substantially improves both level and rate of regret: its final cumulative regret is ≈ 162.3 versus ≈ 220.9 (KOVI0) and ≈ 225.0 (Kernel-LSVI- ε), corresponding to a relative reduction of $\sim 26\%$ against both baselines. The estimated regret slope drops from ≈ 2.22 –2.25 to ≈ 1.58 , indicating faster learning throughout training. In terms of return, KOVI-Proj achieves a positive average (≈ 0.53) while the baselines remain near the step-penalty floor (≈ -0.05 to -0.09), confirming that optimism together with the RKHS *projection* (norm ball + range constraints) materially helps the agent reach the goal within the short horizon despite slippage. The higher SEM for KOVI-Proj reflects mixed outcomes early on (goal reached vs. not reached) typical of sparse-reward exploration; this variance shrinks with longer runs or denser anchors.

REPRODUCIBILITY.

All the results are generated by python notebook code and it is attached in supplementary. Details of experiment are in paper and supplementary.

ETHICS STATEMENT

We have followed the ICLR Code of Ethics throughout this work. Our study does not have any ethical issue.

REFERENCES

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

492

493

494

495

496

497

498 499

500

501

502

504 505

506

507

508 509

510

511 512

513

514

515

516 517

518

519 520

521

522 523

524

525

527 528

529 530

531

532

534 535

536

538

- 486 Claudio Carmeli, Ernesto De Vito, Alessandro Toigo, and Veronica Umanità. Vector-valued reproducing kernel hilbert spaces and universality. Analysis and Applications, 8(1):19–61, 2010. 488
- Sattar R. Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *Proceedings of the* 489 34th International Conference on Machine Learning (ICML), volume 70, pp. 844–853. PMLR, 490 2017. 491
 - Sattar R. Chowdhury and Roberto I. Oliveira. Value function approximations via kernel embeddings for no-regret reinforcement learning. In Proceedings of the Asian Conference on Machine Learning (ACML). PMLR, 2023.
 - Crispin W. Gardiner. Stochastic Methods: A Handbook for the Natural and Social Sciences. Springer Series in Synergetics. Springer, Berlin, Heidelberg, 4 edition, 2009. doi: 10.1007/ 978-3-540-70713-4.
 - Peter Hänggi, Peter Talkner, and Michal Borkovec. Reaction-rate theory: fifty years after kramers. Reviews of Modern Physics, 62(2):251-341, 1990. doi: 10.1103/RevModPhys.62.251. Classic review using the quartic double-well potential $U(x) = ax^4 - bx^2$ as the canonical bistable system.
 - George Kimeldorf and Grace Wahba. Some results on tchebycheffian spline functions. Journal of Mathematical Analysis and Applications, 33(1):82–95, 1971.
 - Krikamol Muandet, Kenji Fukumizu, Bharath Sriperumbudur, and Bernhard Schölkopf. Kernel mean embedding of distributions: A review and beyond. Foundations and Trends in Machine Learning, 10(1-2):1-141, 2017a.
 - Krikamol Muandet, Kenji Fukumizu, Bharath K. Sriperumbudur, and Bernhard Schölkopf. Kernel mean embedding of distributions: A review and beyond. Foundations and Trends in Machine Learning, 10(1-2):1-141, 2017b.
 - Jonathan Scarlett and Ilija Bogunovic. Gaussian process bandits: A tutorial. Foundations and Trends in Machine Learning, 11(5-6):421–516, 2018. Survey framing γ_T and its determinant form $\frac{1}{2} \log \det(I + \rho^{-1} K_n)$.
 - Bernhard Schölkopf and Alexander J. Smola. Learning with Kernels. MIT Press, 2002.
 - Bernhard Schölkopf, Ralf Herbrich, and Alexander J. Smola. A generalized representer theorem. In COLT, 2001.
 - Le Song, Kenji Fukumizu, and Arthur Gretton. Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference. Journal of Machine Learning Research, 14:1415–1444, 2013.
 - Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias W. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In Proceedings of the 27th International Conference on Machine Learning (ICML), 2010. Introduces the GP information gain $\gamma_T = \max_{|A|=T} I(y_A; f_A)$ and uses $\frac{1}{2} \log \det(I + \rho^{-1}K_A)$.
 - Sattar Vakili. Open problem: Order-optimal regret bounds for kernel-based rl. COLT Open Problems, 2024.
 - Zhuoran Yang, Chi Jin, Zhaoran Wang, and Michael I. Jordan. Provably efficient reinforcement learning with kernel and neural function approximations. In Advances in Neural Information Processing Systems (NeurIPS), volume 33, 2020.

ENVIRONMENT AND IMPLEMENTATION DETAILS (DOUBLE-WELL)

Continuous-time model and discretization. We consider the standard overdamped Langevin dynamics in the quartic double-well potential

$$U(x) = \frac{1}{4}x^4 - \frac{1}{2}x^2, \qquad b(x) := -\nabla U(x) = x - x^3,$$

a canonical bistable system (Hänggi et al., 1990; Gardiner, 2009). With additive control u_t and thermal diffusion D>0, the SDE is

$$dX_t = (b(X_t) + u_t) dt + \sqrt{2D} dW_t.$$

We simulate via Euler–Maruyama with step $\Delta t > 0$:

$$X_{t+1} = X_t + \Delta t \left(b(X_t) + u_t \right) + \sigma \varepsilon_t, \qquad \sigma^2 := 2D \, \Delta t, \ \varepsilon_t \sim \mathcal{N}(0, 1).$$
 (6)

Finite-horizon MDP. Episodes have horizon H. The MDP $M = (S, A, P, r, H, \mu_1)$ is:

- State space: S = [-2, 2] (we clip draws from equation 6 to [-2, 2]).
- Action space: $A = \{-u_0, 0, +u_0\}$ (discrete pushes).
- Transitions: $X_{t+1} | X_t = x, A_t = a \sim \mathcal{N}(x + \Delta t (x x^3 + a), \sigma^2)$.
- Goal and absorption: The goal tube is $\mathcal{G} := \{x : |x x_{\text{goal}}| \le \tau\}$ with $x_{\text{goal}} = +1$. On first entry into \mathcal{G} , the episode terminates (absorbing goal).
- Reward: One-shot sparse success with step penalty:

$$r(x, a, x') = \mathbb{1}\{x' \in \mathcal{G}\} - \lambda_{\text{step}}.$$

• Initial state: μ_1 is a point mass near the left well, $X_1 \approx -1$ (small Gaussian jitter).

Absorption ensures V_1^* is O(1) and aligned with the simulated environment.

Default parameters (reproduced). Unless stated otherwise, the experiments in the main text use:

$$H = 40$$
, $\Delta t = 0.10$, $u_0 = 1.0$, $\sigma = 0.07$, $D = \frac{\sigma^2}{2\Delta t}$, $\tau = 0.10$, $\lambda_{\text{step}} = 0.01$.

We run K=100 episodes and average over three seeds. For numerical kernels and projection:

state kernel length $\ell=0.6$, state-action kernel length k=0.35, $\rho=3\times10^{-4}$, B=1.2

The projection grid uses m=81 anchor states; the DP benchmark grid uses M=121 points. We cap each stage buffer to at most 120 tuples to bound kernel linear–algebra cost. We warm-start with 5 random episodes and then plan every 3 episodes (plan-every-K schedule).

Kernels and surrogates. The state RKHS $(\mathcal{H}_{\ell}, \ell)$ uses the RBF kernel $\ell(x, x') = \exp(-\frac{(x - x')^2}{2\ell^2})$. For state–action surrogates we use the product kernel

$$\kappa((x,a),(x',a')) = \ell_k(x,x') \mathbb{1}\{a=a'\}, \qquad \ell_k(x,x') = \exp(-\frac{(x-x')^2}{2k^2}).$$

These choices are standard (Schölkopf & Smola, 2002) and make the Gram matrices PSD.

Projection (QCQP, always). At each stage h, KOVI-Proj projects the optimistic targets $v_h \in \mathbb{R}^m$ onto the feasible RKHS ball with range constraints:

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{m} \|L\alpha - v_h\|_2^2 \quad \text{s.t.} \quad \alpha^\top L\alpha \le B^2, \qquad 0 \le (L\alpha)_j \le H - h + 1 \quad (j = 1, \dots, m) \quad (7)$$

where $L_{ij} = \ell(s_i, s_j)$ for the anchor grid $\{s_j\}_{j=1}^m$. By a constrained representer theorem, the optimizer lies in span $\{\ell(\cdot, s_j)\}$ (Kimeldorf & Wahba, 1971; Schölkopf & Smola, 2002). We solve equation 7 via cvxpy using either MOSEK or SCS; to avoid numerical PSD certification issues on L, we symmetrize $L \leftarrow \frac{1}{2}(L + L^\top)$, add a 10^{-10} ridge, and wrap it with psd_wrap in the quadratic constraint. Projection is performed always; there is no ridge fallback.

Optimistic targets and uncertainty. Given a stage dataset $\mathcal{D}_h = \{(z_i = (x_i, a_i), x_i')\}$ with $y_i := V_{h+1}(x_i')$, we form the KRR mean $\widehat{f}_h(z) = k(z, Z)^{\top}(K + \rho I)^{-1}y$ and its variance via a Cholesky factor of $K + \rho I$,

$$\sigma_h^2(z) = k(z, z) - k(z, Z)^{\top} (K + \rho I)^{-1} k(Z, z)$$

The optimistic action–value uses $\widetilde{Q}_h(x,a) = \widehat{f}_h(x,a) + \beta_h \, \sigma_h(x,a) + r(x,\cdot)$ with a logarithmic scale $\beta_h = \beta(h, |\mathcal{D}_h|)$ (details and schedules in the main text).

Vectorized DP benchmark V^* . We compute V^* on discretization $\{s_j\}_{j=1}^M$ without Monte Carlo by using row–stochastic Gaussian weight matrices. Let $b_j := b(s_j)$ and $\mu_j(a) = s_j + \Delta t \, (b_j + a)$. For each action a, define

$$W_a(i,j) \propto \exp\left(-\frac{(s_j - \mu_i(a))^2}{2\sigma^2}\right), \qquad \sum_{j=1}^M W_a(i,j) = 1,$$

and an absorbing mask $goal(j) = \mathbb{1}\{|s_j - x_{goal}| \le \tau\}$. With $r_j = goal(j) - \lambda_{step}$ and $V_{H+1} \equiv 0$, we recurse

$$V_h(i) = \max_{a \in \mathcal{A}} \sum_{j=1}^{M} W_a(i,j) \Big(r_j + \mathbb{1} \{ \text{goal}(j) = 0 \} V_{h+1}(j) \Big), \qquad h = H, H - 1, \dots, 1$$

This enforces zero continuation from goal bins (absorption) and avoids high-variance MC estimation. Lookup $V_h(x)$ is done by nearest-neighbor interpolation on $\{s_i\}$.

Preprocessing and amortization. We do a warm-start for each learner with 5 random episodes to populate $\{\mathcal{D}_h\}$ before applying optimism; thereafter we perform a full backward planning pass every 3 episodes (*plan-every-K*) and cap per–stage replay by 120 pairs to control kernel linear-algebra cost. These engineering choices do not affect the statement of the algorithms and keep QCQP solves tractable.

B REMARKS ON OTHER WORKS

Kernel function approximation for RL. Kernel methods have long served as nonparametric function approximators in reinforcement learning, bridging linear models and certain infinite-width neural networks. A modern line of work instantiates *optimistic least-squares value iteration* (LSVI) with kernels, coupling kernel ridge regression (KRR) backups with exploration bonuses (Yang et al., 2020). Analytically, these approaches often invoke a union bound over a *data-dependent, evolving* class of optimistic value proxies, bringing in a covering-number penalty that may scale as $\Omega(\sqrt{T})$ for common kernels. This term can spoil no-regret guarantees in long horizons and large time budgets, and it is one of the central obstacles our work circumvents by replacing the union bound with a uniform, CME-based confidence statement that holds simultaneously for all value proxies inside a fixed state-RKHS ball.

Optimistic closure via conditional mean embeddings. A complementary kernel-RL line replaces the evolving-cover argument with a structural assumption: *optimistic closure*, i.e., every optimistic value proxy produced by the algorithm lies in a common, fixed state-RKHS ball. Chowdhury and Oliveira (Chowdhury & Oliveira, 2023) operationalize this idea using *conditional mean embeddings* (CMEs) to map one-step lookahead into a linear functional on the state RKHS. This recovers clean, GP/KRR-style uncertainty quantification, but at the cost of a strong structural premise on the optimizer's iterates. In contrast, our analysis similarly leverages CMEs, yet *dispenses with optimistic closure*: we enforce the bounded-norm property algorithmically by an explicit RKHS projection of the optimistic proxy each step, and then prove a *uniform* confidence bound that applies to *all* functions in the ball *without* any data-dependent covering.

Vector-valued RKHS and CMEs. Our development relies on classical results on vector-valued RKHSs and conditional mean embeddings. The CME view represents the Bellman image as an inner product $[P_hV](z) = \langle \mu_h(z), V \rangle_{\mathcal{H}_\ell}$ with an \mathcal{H}_ℓ -valued map μ_h ; this viewpoint is extensively surveyed by Muandet et al. (Muandet et al., 2017b). The required functional-analytic foundations for vector-valued RKHSs with operator-valued kernels such as K(z,z')=k(z,z')I—can be found in Carmeli, De Vito, Toigo, and co-authors (Carmeli et al., 2010). Building on these tools, we show that (i) scalar KRR predictions with labels V(s') can be written as an inner product with a *vector-valued* KRR estimator of the CME, and (ii) a single Hilbert-space self-normalized concentration argument yields uniform confidence for the whole state-ball $\{V: \|V\|_{\mathcal{H}_\ell} \leq B\}$, removing the covering-number penalty.

 Kernel bandits, information gain, and elliptical potentials. Our regret analysis adopts the standard information-gain and elliptical-potential machinery developed for kernelized bandits and GP regression. In particular, Chowdhury and Gopalan (Chowdhury & Gopalan, 2017) provide clean, modular bounds in terms of the (regularized) information gain $\gamma(n,\rho)$, which we adapt to the multi-step RL setting by summing per-step potentials (with a block-diagonal argument across steps). The combination of CME-based linearization and information-gain control yields the $\tilde{\mathcal{O}}(\sqrt{T\gamma(HT,\rho)})$ -type scaling in our main result, while avoiding data-dependent covers.

Positioning within kernel RL. Putting these threads together, our contribution can be viewed as a third route to kernel-RL optimism: (i) unlike covering-number analyses for kernelized LSVI (Yang et al., 2020), we avoid data-dependent covers; (ii) unlike *optimistic closure* (Chowdhury & Oliveira, 2023), we do not assume a priori that all optimistic proxies already lie in a fixed state-RKHS ball; instead, (iii) we *enforce* the bounded-norm property by projection and prove a *uniform* CME-based confidence bound that holds for all functions in the ball simultaneously. This uniformization is central to obtaining sublinear regret without the $\Omega(\sqrt{T})$ covering penalty.

Context in open problems. The broader agenda of obtaining sharp or order-optimal regret guarantees for kernel-based RL has been highlighted as an open challenge (Vakili, 2024). Our analysis: via vector-valued RKHS concentration for CMEs and a projection step that replaces optimistic closure addresses a prominent bottleneck identified in that discussion: removing the covering-number dependence while retaining principled uncertainty quantification in kernelized optimistic value iteration.

On horizon dependence and refinements. As in kernelized optimistic LSVI, our H^2 scaling arises from a standard telescoping decomposition and coarse coupling of stepwise estimation errors. While we expect refined Bellman-error coupling or variance-aware decompositions to reduce this to $H^{3/2}$ or even H, the present focus is on eliminating the covering-number obstruction under a natural CME boundedness condition closing a gap emphasized in prior work (Yang et al., 2020; Chowdhury & Oliveira, 2023; Vakili, 2024).

C PROOF OF THEOREM 1

Proposition C.1 (Scalar KRR = inner product with a vector-valued KRR). Fix a step h and data $\{(z_i, s_i')\}_{i=1}^n$. Let ℓ be a kernel on S with RKHS $(\mathcal{H}_{\ell}, \langle \cdot, \cdot \rangle_{\mathcal{H}_{\ell}})$ and feature map $\phi : S \to \mathcal{H}_{\ell}$ so that $\ell(s, s') = \langle \phi(s), \phi(s') \rangle_{\mathcal{H}_{\ell}}$. Let k be a kernel on $\mathcal{Z} = S \times \mathcal{A}$ with Gram matrix $K_n = [k(z_i, z_j)]_{i,j=1}^n$ and, for $z \in \mathcal{Z}$, define $k_n(z) = [k(z, z_1), \dots, k(z, z_n)]^\top$. For a ridge parameter $\rho > 0$, define

$$\widehat{\mu}_n(z) := \sum_{i=1}^n \alpha_i(z) \, \phi(s_i') \in \mathcal{H}_\ell, \qquad \alpha(z) := (K_n + \rho I)^{-1} k_n(z).$$

Then for every $V \in \mathcal{H}_{\ell}$ and $z \in \mathcal{Z}$,

$$\widehat{f}_{h,n}^{V}(z) = \langle \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}},$$

where $\widehat{f}_{h,n}^V$ is the scalar KRR predictor trained on labels $y_i^{(V)} := V(s_i') = \langle \phi(s_i'), V \rangle_{\mathcal{H}_\ell}$, i.e. $\widehat{f}_{h,n}^V(z) = k_n(z)^\top (K_n + \rho I)^{-1} \boldsymbol{y}^{(V)}$ with $\boldsymbol{y}^{(V)} = (y_1^{(V)}, \dots, y_n^{(V)})^\top$.

Proof. We give a self-contained argument in two steps.

Step 1: Scalar KRR with inner-product labels. Fix $V \in \mathcal{H}_{\ell}$. Consider the *scalar* KRR problem on the input space \mathcal{Z} with kernel k and training labels

$$y_i^{(V)} := V(s_i') = \langle \phi(s_i'), V \rangle_{\mathcal{H}_{\epsilon}}, \qquad i = 1, \dots, n.$$

It is standard that the KRR predictor at a test point $z \in \mathcal{Z}$ is

$$\widehat{f}_{h,n}^{V}(z) = k_n(z)^{\top} (K_n + \rho I)^{-1} \boldsymbol{y}^{(V)}.$$
(8)

Step 2: Vector-valued KRR and the CME estimator. Define the *vector-valued* RKHS on \mathcal{Z} with operator-valued kernel $K(z,z'):=k(z,z')\,I_{\mathcal{H}_\ell}$; this space can be identified with the tensor-product RKHS $\mathcal{H}_k\otimes\mathcal{H}_\ell$. Consider the vector-valued KRR problem that regresses the \mathcal{H}_ℓ -valued observations $\phi_i:=\phi(s_i')\in\mathcal{H}_\ell$ on the inputs z_i :

$$\widehat{\mu}_n \in \arg\min_{g \in \mathcal{H}_k \otimes \mathcal{H}_\ell} \left\{ \sum_{i=1}^n \|\phi_i - g(z_i)\|_{\mathcal{H}_\ell}^2 + \rho \|g\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}^2 \right\}.$$
 (9)

By the (vector-valued) representer theorem, the minimizer has the finite form

$$\widehat{\mu}_n(\cdot) = \sum_{i=1}^n K(\cdot, z_i) c_i = \sum_{i=1}^n k(\cdot, z_i) c_i, \qquad c_i \in \mathcal{H}_{\ell}$$

Let $C = [c_1, \dots, c_n]$ be the column tuple and note that $g(z_j) = \sum_{i=1}^n k(z_j, z_i)c_i$. The normal equations for equation 9 read

$$(K_n + \rho I) C^{\top} = \Phi^{\top}, \quad \text{where } \Phi : \mathbb{R}^n \to \mathcal{H}_{\ell}, \ \Phi e_i = \phi_i,$$

so that $C^{\top} = (K_n + \rho I)^{-1} \Phi^{\top}$. Therefore, for any $z \in \mathcal{Z}$,

$$\widehat{\mu}_n(z) = \sum_{i=1}^n k(z, z_i) c_i = \sum_{i=1}^n \alpha_i(z) \phi_i = \sum_{i=1}^n \alpha_i(z) \phi(s_i'), \qquad \alpha(z) := (K_n + \rho I)^{-1} k_n(z),$$
(10)

which matches the stated definition.

Equality of predictons. Combining equation 8 and equation 10, and recalling $y_i^{(V)} = \langle \phi(s_i'), V \rangle_{\mathcal{H}_s}$, we compute

$$\widehat{f}_{h,n}^{V}(z) = k_n(z)^{\top} (K_n + \rho I)^{-1} \mathbf{y}^{(V)} = \sum_{i=1}^{n} \alpha_i(z) y_i^{(V)} = \sum_{i=1}^{n} \alpha_i(z) \langle \phi(s_i'), V \rangle_{\mathcal{H}_{\ell}} = \langle \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}}.$$

This holds for every $V \in \mathcal{H}_{\ell}$ and every $z \in \mathcal{Z}$, as claimed.

D Proof of Theorem 5.1

Proof. Step 1: A uniform "good" event. Apply Theorem 3.3 with a union bound over all steps $h \in [H]$, episodes $t \in [T]$, and query points z (the latter handled by the supremum in Theorem 3.3). Using the per-step confidence radius in equation 3 with the $\log(2HT/\delta)$ factor, there exists an event

$$\mathcal{G}$$
 with $\Pr(\mathcal{G}) \geq 1 - \delta$

such that, *simultaneously* for all h, t and all $z \in \mathcal{Z}$,

$$[P_h V_{h+1,t}](z) \le \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \beta_{h,t} \, \sigma_{h,t}(z) \tag{11}$$

where $\beta_{h,t} = B\left(\sqrt{\rho}\,U + \frac{\sigma}{\sqrt{\rho}}\sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{2HT}{\delta}}\right)$ and $\sigma_{h,t}$ is as in equation 2. See proof in H.1. The projection step (Section 4) guarantees $\|V_{h,t}\|_{\mathcal{H}_{\ell}} \leq B$, ensuring applicability of Theorem 3.3 to the *actual* proxies the algorithm uses.

Remark D.1 (Why the projection step matters?). Theorem 3.3 provides a high-probability confidence bound that holds uniformly for all value functions V whose RKHS norm is bounded by B, i.e., for all $V \in \{V : ||V||_{\mathcal{H}_{\ell}} \leq B\}$. The optimistic proxy $\widetilde{V}_{h,t}$ produced by the backup (§4) need not lie in this ball a priori. The projection step maps $\widetilde{V}_{h,t}$ to

$$V_{h,t} \in \arg\min_{\|V\|_{\mathcal{H}_s} \leq B} \|V - \widetilde{V}_{h,t}\|_{L^2(\nu)}$$
 (with range clipping),

thereby guaranteeing $||V_{h,t}||_{\mathcal{H}_{\ell}} \leq B$. Consequently, every value proxy the algorithm actually uses satisfies the assumptions of Theorem 3.3, and the uniform confidence bound applies directly to the algorithm's updates without any additional covering or closure assumptions.

 Step 2: Optimism up to agnostic error. Fix (h,t) and z=(s,a). By definition of Q_h^* and by boundedness of the value range,

$$Q_h^*(z) = r_h(z) + [P_h V_{h+1}^*](z) \le r_h(z) + [P_h V_{h+1,t}](z) + ||V_{h+1}^* - V_{h+1,t}||_{\infty}.$$

By the definition of the "worst case" agnostic approximation level $\varepsilon_B := \max_h \sup_{\|V\|_{\mathcal{H}_\ell} \leq B} \|V_h^* - V\|_{\infty}$ and since $\|V_{h+1,t}\|_{\mathcal{H}_\ell} \leq B$, we have $\|V_{h+1}^* - V_{h+1,t}\|_{\infty} \leq \varepsilon_B$. See proof in D.2. Combining with equation 11 and the definition equation 3 of $\widetilde{Q}_{h,t}$ gives

$$Q_h^*(z) \le \widetilde{Q}_{h,t}(z) + \varepsilon_B$$
 for all h, t, z on the event \mathcal{G} . (12)

See proof in Remark D.3. Maximizing over a further yields $V_h^*(s) \leq \widetilde{V}_{h,t}(s) + \varepsilon_B$.

Step 3: Telescoping regret decomposition. Let $z_{h,t} = (s_{h,t}, a_{h,t})$ be the state-action chosen by KOVI-Proj at step h of episode t. From equation 12 and the greedy action choice $a_{h,t} \in \arg\max_a \widetilde{Q}_{h,t}(s_{h,t},a)$,

$$V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t}) \leq \sum_{h=1}^H \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t}) \right) + H \, \varepsilon_B.$$

See Remark D.3. Summing over episodes and using equation 3 then gives

$$R(T) \leq \sum_{t=1}^{T} \sum_{h=1}^{H} \beta_{h,t} \, \sigma_{h,t}(z_{h,t}) + HT \, \varepsilon_B \quad \text{on } \mathcal{G}.$$
 (13)

Step 4: Elliptical-potential bound across steps. For each fixed step h, let $n_{h,T}$ be the number of transitions observed at step h up to episode T. Denote by $\sigma_{h,\tau-1}(z_{h,\tau})$ the posterior standard deviation just before the τ -th observation at step h (this is exactly $\sigma_{h,t}(z_{h,t})$ when the τ -th observation occurs in episode t). The standard GP/RKHS potential argument applied to the (adaptively chosen) design at step t0 yields

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \; \leq \; 2 \, \gamma(n_{h,T},\rho), \qquad \sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}(z_{h,\tau}) \; \leq \; \sqrt{2 \, n_{h,T} \, \gamma(n_{h,T},\rho)}.$$

See detailed proof in Lemma D.5. Summing over h and using Cauchy-Schwarz,

$$\sum_{t=1}^{T} \sum_{h=1}^{H} \sigma_{h,t}(z_{h,t}) \leq \sum_{h=1}^{H} \sqrt{2 \, n_{h,T} \, \gamma(n_{h,T}, \rho)} \leq \sqrt{2 \left(\sum_{h} n_{h,T} \right) \left(\sum_{h} \gamma(n_{h,T}, \rho) \right)} = \sqrt{2 \, HT \, \Gamma_{T}}.$$

See Remark D.9 for last equality. Let K_h be the Gram matrix of the design at step h and $K_{\text{blk}} := \text{diag}(K_1, \dots, K_H)$. Then

$$\Gamma_T = \frac{1}{2} \sum_{h=1}^{H} \log \det(I + \rho^{-1} K_h) = \frac{1}{2} \log \det(I + \rho^{-1} K_{\text{blk}}) \le \frac{1}{2} \log \det(I + \rho^{-1} K_{\text{all}})$$
 $\le \gamma(HT, \rho)$

where $K_{\rm all}$ is full Gram matrix over the concateneted HT design points and the last inequality uses that adding nonnegative off-diagonal blocks (cross-step similarities) increases the determinant. Therefore,

$$\sum_{t=1}^{T} \sum_{h=1}^{H} \sigma_{h,t}(z_{h,t}) \leq \sqrt{2 HT \gamma(HT, \rho)}.$$
 (14)

Step 5: Putting it together. From equation 13 and equation 14, and using that (see proof in Remark D.8)

$$\beta_{h,t} \leq \tilde{\mathcal{O}}\left(B\left(\sqrt{\rho}\,U + \frac{\sigma}{\sqrt{\rho}}\sqrt{\gamma(HT,\rho)}\right)\right)$$
 uniformly over $h,t,$

we obtain

$$R(T) \; \leq \; \tilde{\mathcal{O}}\!\!\left(H^2\,B\!\left(\sqrt{\rho}\,U + \frac{\sigma}{\sqrt{\rho}}\sqrt{\gamma(HT,\rho)}\right)\sqrt{T\,\gamma(HT,\rho)}\right) \; + \; HT\,\varepsilon_B,$$

which is the claimed bound. This completes the proof.

 Lemma D.2 (Agnostic approximation bound for projected proxies). For each $h \in [H]$, define the worst-case (supremum) approximation error of the RKHS ball

$$\varepsilon_B(h) := \sup_{\|V\|_{\mathcal{H}_{\ell}} \leq B} \|V_h^* - V\|_{\infty}, \qquad \varepsilon_B := \max_{j \in [H]} \varepsilon_B(j).$$

If the algorithm's projection guarantees $||V_{h,t}||_{\mathcal{H}_{\ell}} \leq B$ for all h, t, then for every $h \in [H]$ and $t \in [T]$,

$$\|V_h^* - V_{h,t}\|_{\infty} \le \varepsilon_B(h) \le \varepsilon_B.$$

In particular, with $h\mapsto h+1$ we get $\|V_{h+1}^*-V_{h+1,t}\|_\infty\leq \varepsilon_B$

Proof. Fix $h \in [H]$ and $t \in [T]$. By projecton, $||V_{h,t}||_{\mathcal{H}_{\ell}} \leq B$, so $V_{h,t}$ belongs to the admissible set in the definition of $\varepsilon_B(h)$. Since $\varepsilon_B(h)$ is a *supremum* over that set, it dominates the error at the particular choice $V_{h,t}$:

$$\left\| V_h^* - V_{h,t} \right\|_{\infty} \le \sup_{\|V\|_{\mathcal{H}_s} \le B} \left\| V_h^* - V \right\|_{\infty} = \varepsilon_B(h).$$

Finally, by definition $\varepsilon_B(h) \leq \max_{j \in [H]} \varepsilon_B(j) = \varepsilon_B$, which yields the second inequality. The special case $h \mapsto h + 1$ is immediate.

Remark D.3 (Telescoping bound from optimism up to ε_B). From equation 12, for every step h and episode t and every z = (s, a),

$$Q_h^*(z) \leq \widetilde{Q}_{h,t}(z) + \varepsilon_B.$$

Evaluating at the algorithm's visited pair $z_{h,t} = (s_{h,t}, a_{h,t})$ and using the Bellman identities

$$V_h^*(s_{h,t}) = r_h(z_{h,t}) + [P_h V_{h+1}^*](z_{h,t}), \qquad V_h^{\pi_t}(s_{h,t}) = r_h(z_{h,t}) + [P_h V_{h+1}^{\pi_t}](z_{h,t}),$$

we obtain the one-step inequality

$$\begin{split} V_h^*(s_{h,t}) - V_h^{\pi_t}(s_{h,t}) &= [P_h V_{h+1}^*](z_{h,t}) - [P_h V_{h+1}^{\pi_t}](z_{h,t}) \\ &\leq \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t})\right) - [P_h V_{h+1,t}](z_{h,t}) \\ &+ \underbrace{\left([P_h V_{h+1}^*] - [P_h V_{h+1}^{\pi_t}]\right)(z_{h,t})}_{= \mathbb{E}\left[V_{h+1}^*(s_{h+1,t}) - V_{h+1}^{\pi_t}(s_{h+1,t}) \mid s_{h,t}, a_{h,t}\right]} + \varepsilon_B. \end{split}$$

Taking conditional expectation on the episode's history up to step h (which leaves the displayed conditional expectation unchanged), and summing this inequality over h = 1, ..., H makes the middle terms telescope (See Remark D.4):

$$\sum_{h=1}^{H} \mathbb{E} \big[V_{h+1}^*(s_{h+1,t}) - V_{h+1}^{\pi_t}(s_{h+1,t}) \, \big| \, \textit{history up to } h \big] = \mathbb{E} \big[V_{H+1}^*(s_{H+1,t}) - V_{H+1}^{\pi_t}(s_{H+1,t}) \big] = 0.$$

since $V_{H+1}^* \equiv V_{H+1}^{\pi_t} \equiv 0$. Therefore,

$$V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t}) \leq \sum_{h=1}^H \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t}) \right) + H \varepsilon_B,$$

which is the claimed bound.

Remark D.4. How the middle terms telescope. Let \mathcal{F}_h be the history (sigma-field) up to step h in episode t, and define

$$\Delta_h := V_h^*(s_{h,t}) - V_h^{\pi_t}(s_{h,t}), \qquad h = 1, \dots, H, \quad \text{with} \quad \Delta_{H+1} = 0$$

From equation 12 we derived, for each h.

$$\Delta_h \leq \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t})\right) + \mathbb{E}[\Delta_{h+1} \mid \mathcal{F}_h] + \varepsilon_B. \tag{15}$$

Rearrange equation 15 to isolate the martingale increment

$$\Delta_h - \mathbb{E}[\Delta_{h+1} \mid \mathcal{F}_h] \leq (\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t})) + \varepsilon_B.$$

Summing this inequality over h = 1, ..., H and using linearity gives

$$\sum_{h=1}^{H} \left(\Delta_h - \mathbb{E}[\Delta_{h+1} \mid \mathcal{F}_h] \right) \leq \sum_{h=1}^{H} \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t}) \right) + H \varepsilon_B.$$

The left-hand side telescopes by the tower property:

$$\sum_{h=1}^{H} \left(\Delta_h - \mathbb{E}[\Delta_{h+1} \mid \mathcal{F}_h] \right) = \Delta_1 - \mathbb{E}[\Delta_{H+1} \mid \mathcal{F}_H] = \Delta_1 - 0 = V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t}),$$

because $\Delta_{H+1} = V_{H+1}^*(s_{H+1,t}) - V_{H+1}^{\pi_t}(s_{H+1,t}) \equiv 0$. Thus we obtain

$$V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t}) \leq \sum_{h=1}^H \left(\widetilde{Q}_{h,t}(z_{h,t}) - r_h(z_{h,t}) - [P_h V_{h+1,t}](z_{h,t}) \right) + H \, \varepsilon_B$$

Concrete cancellation for H=3 (**illustration**). Writing equation 15 for h=1,2,3 and subtracting the conditional expectations:

$$\begin{array}{ll} \Delta_1 - \mathbb{E}[\Delta_2 \mid \mathcal{F}_1] & \leq \ bonus_1 + \varepsilon_B, \\ \Delta_2 - \mathbb{E}[\Delta_3 \mid \mathcal{F}_2] & \leq \ bonus_2 + \varepsilon_B, \\ \Delta_3 - \mathbb{E}[\Delta_4 \mid \mathcal{F}_3] & \leq \ bonus_3 + \varepsilon_B & (\Delta_4 \equiv 0). \end{array}$$

Summing yields

$$\left(\Delta_{1} - \mathbb{E}[\Delta_{2} \mid \mathcal{F}_{1}]\right) + \left(\Delta_{2} - \mathbb{E}[\Delta_{3} \mid \mathcal{F}_{2}]\right) + \left(\Delta_{3} - \mathbb{E}[\Delta_{4} \mid \mathcal{F}_{3}]\right) \leq bonus_{1} + bonus_{2} + bonus_{3} + 3\varepsilon_{B}$$

The middle terms cancel pairwise by the tower property: $-\mathbb{E}[\Delta_2 \mid \mathcal{F}_1] + \Delta_2$ and $-\mathbb{E}[\Delta_3 \mid \mathcal{F}_2] + \Delta_3$ vanish after taking expectations step by step, and $\mathbb{E}[\Delta_4 \mid \mathcal{F}_3] = 0$. What remains is exactly Δ_1 on the left, i.e., $V_1^*(s_{1,t}) - V_1^{\pi_t}(s_{1,t})$, which proves the claim.

Lemma D.5 (Elliptical potential / information-gain bound at a fixed step). Fix a step h and let $\{z_{h,T}\}_{\tau=1}^{n_{h,T}}$ be the (adaptively chosen) design points collected at this step up to time T. Let

$$\sigma_{h,\tau-1}^2(z) := k(z,z) - k_{h,\tau-1}(z)^{\top} (K_{h,\tau-1} + \rho I)^{-1} k_{h,\tau-1}(z),$$

where $K_{h,\tau-1} = [k(z_{h,i}, z_{h,j})]_{i,j=1}^{\tau-1}$ and $k_{h,\tau-1}(z) = [k(z, z_{h,1}), \dots, k(z, z_{h,\tau-1})]^{\top}$. Then, for any $\rho > 0$,

$$\sum_{\tau=1}^{n_{h,T}} \log \left(1 + \frac{\sigma_{h,\tau-1}^2(z_{h,\tau})}{\rho} \right) = \frac{1}{2} \log \det \left(I + \rho^{-1} K_{h,n_{h,T}} \right) =: \gamma(n_{h,T}, \rho), \tag{16}$$

and consequently

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}(z_{h,\tau}) \leq \sqrt{n_{h,T} \sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau})}$$
 (by Cauchy-Schwarz). (17)

Moreover, under the common normalization $k(z, z) \leq 1$ *and* $\rho = 1$ *,*

$$\sum_{l=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \leq 2\gamma(n_{h,T},1), \qquad \sum_{l=1}^{n_{h,T}} \sigma_{h,\tau-1}(z_{h,\tau}) \leq \sqrt{2n_{h,T}\gamma(n_{h,T},1)}. \tag{18}$$

Proof. We prove in following three steps.

Step 1: Determinant telescoping (matrix determinant lemma). Let $A_{\tau-1} := K_{h,\tau-1} + \rho I$ (with $A_0 = \rho I$). Consider augmenting $A_{\tau-1}$ by the new point $z_{h,\tau}$, i.e., the block matrix

$$A_{\tau} = \begin{bmatrix} K_{h,\tau-1} + \rho I & k_{h,\tau-1}(z_{h,\tau}) \\ k_{h,\tau-1}(z_{h,\tau})^{\top} & k(z_{h,\tau}, z_{h,\tau}) + \rho \end{bmatrix}.$$

 By the Schur complement (or the matrix determinant lemma),

$$\det(A_{\tau}) = \det(A_{\tau-1}) \left(\rho + k(z_{h,\tau}, z_{h,\tau}) - k_{h,\tau-1}(z_{h,\tau})^{\top} A_{\tau-1}^{-1} k_{h,\tau-1}(z_{h,\tau}) \right)$$
$$= \det(A_{\tau-1}) \left(\rho + \sigma_{h,\tau-1}^2(z_{h,\tau}) \right).$$

Divide both sides by ρ^{τ} and take logs. Telescoping over $\tau = 1, \dots, n_{h,T}$ gives

$$\log \det (I + \rho^{-1} K_{h,n_{h,T}}) = \sum_{\tau=1}^{n_{h,T}} \log \left(1 + \frac{\sigma_{h,\tau-1}^2(z_{h,\tau})}{\rho} \right)$$

which is equation 16 after multiplying by 1/2 to match the definition $\gamma(n,\rho)=\frac{1}{2}\log\det(I+\rho^{-1}K)$.

Step 2: From equation 16 to bounds on sums. The second display equation 17 is a direct application of Cauchy-Schwarz: $\sum a_{\tau} \leq \sqrt{(\sum 1)(\sum a_{\tau}^2)}$.

To control $\sum \sigma^2$ in terms of γ , one can use standard scalar inequalities relating $\log(1+x)$ and x. A common (and sharp) form in the GP literature (see, e.g., Srinivas et al., 2010, or Chowdhury & Gopalan, 2017) is

$$\sum_{\tau=1}^{n_{h,T}} \min \left\{ 1, \ \frac{\sigma_{h,\tau-1}^2(z_{h,\tau})}{\rho} \right\} \le 2 \sum_{\tau=1}^{n_{h,T}} \log \left(1 + \frac{\sigma_{h,\tau-1}^2(z_{h,\tau})}{\rho} \right) \ = \ 4 \, \gamma(n_{h,T},\rho).$$

In particular, under the normalization $k(z, z) \le 1$ and $\rho = 1$, we have $0 \le \sigma_{h, \tau - 1}^2(z_{h, \tau}) \le 1$ so that $\min\{1, \sigma^2\} = \sigma^2$. See proof in D.6. Thus

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \leq 2 \log \det \left(I + K_{h,n_{h,T}}\right) = 2 \cdot 2 \gamma(n_{h,T}, 1) = 4 \gamma(n_{h,T}, 1).$$

See detailed proof in Remark D.7. A slightly refined inequality (using, for $x \in [0,1]$, that $\log(1+x) \ge x - x^2/2$ together with $\sum \sigma^4 \le \sum \sigma^2$) improves the constant and yields

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \leq 2 \gamma(n_{h,T}, 1),$$

as stated in equation 18. Finally, combining with equation 17 gives

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}(z_{h,\tau}) \le \sqrt{2 \, n_{h,T} \, \gamma(n_{h,T}, 1)}$$

Remark on constants. All bounds above hold up to universel constants that can be made explicit; the versions in equation 18 are the ones commonly used in GP-UCB analyses (with $k(z,z) \le 1$, $\rho = 1$). For general $\rho > 0$, one obtains $\sum \sigma^2 \lesssim \rho \gamma(n,\rho)$ and hence $\sum \sigma \lesssim \sqrt{\rho n \gamma(n,\rho)}$.

Remark D.6. Why $\min\{1, \sigma^2\} = \sigma^2$ when $k(z, z) \le 1$ and $\rho = 1$. Recall the posterior deviation at time $\tau - 1$:

$$\sigma_{h,\tau-1}^2(z) = k(z,z) - k_{h,\tau-1}(z)^{\top} (K_{h,\tau-1} + I)^{-1} k_{h,\tau-1}(z).$$

Two facts imply $0 \le \sigma_{h,\tau-1}^2(z) \le 1$:

1. Nonnegativity. The block matrix $\binom{K_{h,\tau-1}+I}{k_{h,\tau-1}(z)^{\top}}$ is positive semidefinite, so its Schur complement is nonnegative:

$$k(z,z) - k_{h,\tau-1}(z)^{\top} (K_{h,\tau-1} + I)^{-1} k_{h,\tau-1}(z) \ge 0$$

2. Upper bound by k(z,z). Since the subtracted term is nonnegative, $\sigma_{h,\tau-1}^2(z) \leq k(z,z) \leq 1$ under the normalization $k(z,z) \leq 1$

Therefore, pointwise for every querid $z_{h,\tau}$,

$$0 \leq \sigma_{h,\tau-1}^2(z_{h,\tau}) \leq 1,$$

and hence $\min\{1, \sigma_{h,\tau-1}^2(z_{h,\tau})\} = \sigma_{h,\tau-1}^2(z_{h,\tau}).$

Remark D.7. From $\sum \min\{1, \sigma^2\}$ to $\sum \sigma^2$ and γ . Under the normalization $k(z, z) \leq 1$ and $\rho = 1$ we have $0 \leq \sigma_{h, \tau - 1}^2(z_{h, \tau}) \leq 1$, hence $\min\{1, \sigma_{h, \tau - 1}^2(z_{h, \tau})\} = \sigma_{h, \tau - 1}^2(z_{h, \tau})$. A standard scalar inequalty used in GP/KRR analyses (see, e.g., GP-UCB) states that

$$\sum_{\tau=1}^{n_{h,T}} \min \left\{ 1, \ \sigma_{h,\tau-1}^2(z_{h,\tau}) \right\} \le 2 \sum_{\tau=1}^{n_{h,T}} \log \left(1 + \sigma_{h,\tau-1}^2(z_{h,\tau}) \right)$$

Therefore,

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \leq 2 \sum_{\tau=1}^{n_{h,T}} \log \Big(1 + \sigma_{h,\tau-1}^2(z_{h,\tau})\Big).$$

Using the determinant telescoping identity $\sum_{\tau=1}^{n_{h,T}} \log(1+\sigma_{h,\tau-1}^2(z_{h,\tau})) = \log \det(I+K_{h,n_{h,T}})$ (at $\rho=1$), we obtain

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \le 2 \log \det (I + K_{h,n_{h,T}}).$$

Finally, by definition $\gamma(n_{h,T},1) = \frac{1}{2} \log \det(I + K_{h,n_{h,T}})$, so

$$2 \log \det(I + K_{h,n_{h,T}}) = 2 \cdot 2 \gamma(n_{h,T}, 1) = 4 \gamma(n_{h,T}, 1)$$

Hence

$$\sum_{\tau=1}^{n_{h,T}} \sigma_{h,\tau-1}^2(z_{h,\tau}) \le 4 \gamma(n_{h,T}, 1)$$

where two factors of "2" come from (i) the scalr inequality linking $\min\{1, \sigma^2\}$ to $\log(1 + \sigma^2)$ and (ii) the definition $\gamma = \frac{1}{2} \log \det(\cdot)$.

Remark D.8 (Uniform bound on $\beta_{h,t}$). Recall

$$\beta_{h,t} = B\left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{2HT}{\delta}}\right),$$

where $n_{h,t-1} = |\mathcal{D}_{h,t-1}|$ is the number of step-h samples before episode t and $\gamma(\cdot,\rho)$ is the (regularized) information gain. Since $n_{h,t-1} \leq \sum_{h'=1}^{H} n_{h',t-1} \leq HT$ and $\gamma(n,\rho)$ is nondecreasing in n,

$$\gamma(n_{h,t-1},\rho) \leq \gamma(HT,\rho)$$
 for all h,t .

Therefore,

$$\beta_{h,t} \leq B\left(\sqrt{\rho}U + \frac{\sigma}{\sqrt{\rho}}\sqrt{2\gamma(HT,\rho) + 2\log\frac{2HT}{\delta}}\right) \leq \tilde{\mathcal{O}}\left(B\left(\sqrt{\rho}U + \frac{\sigma}{\sqrt{\rho}}\sqrt{\gamma(HT,\rho)}\right)\right),$$

uniformly over h, t, where $\tilde{\mathcal{O}}(\cdot)$ hides polylogarithmic factors in $(H, T, 1/\delta)$ and absolute constants. The last step uses the elementary inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and absorbs the $\sqrt{\log(2HT/\delta)}$ term into the $\tilde{\mathcal{O}}(\cdot)$ notation.

Remark D.9. Why $\sqrt{2\left(\sum_{h}n_{h,T}\right)\left(\sum_{h}\gamma(n_{h,T},\rho)\right)} = \sqrt{2HT\Gamma_{T}}$. By definition we set

$$\Gamma_T := \sum_{h=1}^H \gamma(n_{h,T}, \rho).$$

Also, over T episodes and H steps per episode, total number of design points across all steps is

$$\sum_{h=1}^{H} n_{h,T} = HT$$

Substituting these two identities into $\sqrt{2\left(\sum_{h}n_{h,T}\right)\left(\sum_{h}\gamma(n_{h,T},\rho)\right)}$ gives

$$\sqrt{2\left(\,\sum_{h}n_{h,T}\right)\left(\,\sum_{h}\gamma(n_{h,T},\rho)\right)}\;=\;\sqrt{\,2\,(HT)\,\Gamma_{T}}\;=\;\sqrt{2\,HT\,\Gamma_{T}}\,.$$

E PROOF OF LEMMA

Lemma E.1 (Vector-valued kernel ridge concentration). Suppose Assumption 2.1 holds, $k(z,z) \leq \kappa_k^2$, and $\ell(s,s) \leq \kappa_\ell^2$. Let $\rho > 0$ and define $\sigma_{h,n}(\cdot)$ by equation 2. Then for any $\delta \in (0,1)$, with probability at least $1 - \delta$, simultaneously for all $z \in \mathcal{Z}$,

$$\|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}} \le \left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right) \sigma_{h,n}(z).$$

Proof. Recall the vector-valued KRR estimator $\widehat{\mu}_n : \mathcal{Z} \to \mathcal{H}_{\ell}$ defined by

$$\widehat{\mu}_n(z) = \sum_{i=1}^n \alpha_i(z) \, \phi(s_i'), \qquad \alpha(z) = (K_n + \rho I)^{-1} k_n(z),$$

where $K_n = [k(z_i, z_j)]_{i,j=1}^n$, $k_n(z) = [k(z, z_1), \dots, k(z, z_n)]^\top$, and $\phi(s')$ is the canonical feature map of ℓ . Let $\Phi : \mathbb{R}^n \to \mathcal{H}_\ell$ be the linear map $\Phi b = \sum_{i=1}^n b_i \, \phi(s_i')$, so $\widehat{\mu}_n(z) = \Phi^\top (K_n + \rho I)^{-1} k_n(z)$. By the data model (Section 2),

$$\phi(s_i') = \mu(z_i) + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i \mid \mathcal{F}_{i-1}] = 0, \quad \|\varepsilon_i\|_{\mathcal{H}_\ell} \le \kappa_\ell, \quad \sigma$$
-sub-Gaussian in \mathcal{H}_ℓ ,

where $\{\mathcal{F}_i\}$ is the natural filtration.

Error decomposition. Let $\mu \in \mathcal{H}_k \otimes \mathcal{H}_\ell$ denote the (unknown) CME map $z \mapsto \mu(z)$. Write $\Phi = \underbrace{M}_{\text{signal}} + \underbrace{E}_{\text{noise}}$, where $M\mathbf{b} = \sum_i b_i \mu(z_i)$ and $E\mathbf{b} = \sum_i b_i \varepsilon_i$. Then, for any $z \in \mathcal{Z}$,

$$\mu(z) - \widehat{\mu}_n(z) = \underbrace{\mu(z) - M^{\top}(K_n + \rho I)^{-1}k_n(z)}_{\text{bias}} - \underbrace{E^{\top}(K_n + \rho I)^{-1}k_n(z)}_{\text{noise}}$$
(19)

We next bound the two terms separately and then combine via the triangle inequality.

Bias term. Let $\mathcal{H}_{k,I}$ be the vector-valued RKHS over \mathcal{Z} with operator-valued kernel $K(z,z')=k(z,z')I_{\mathcal{H}_\ell}$ and norm $\|\cdot\|_{\mathcal{H}_k\otimes\mathcal{H}_\ell}$. Denote by $\Pi_{n,\rho}$ the ρ -regularized orthogonal projector onto the finite-dimensional subspace span $\{K(\cdot,z_i)u:i\in[n],u\in\mathcal{H}_\ell\}\subset\mathcal{H}_{k,I}$. It is standard (vector-valued representer theorem and Tikhonov interpolation inequality, see Lemma G.2 in Appendix) that

$$\|\mu(z) - M^{\top}(K_n + \rho I)^{-1}k_n(z)\|_{\mathcal{H}_{\ell}} = \|\mu(z) - \Pi_{n,\rho}\mu(z)\|_{\mathcal{H}_{\ell}} \le \sqrt{\rho} \|\mu\|_{\mathcal{H}_k \otimes \mathcal{H}_{\ell}} \sigma_{h,n}(z).$$
(20)

By Assumption 2.1 we have $\|\mu\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell} \leq U$, hence the bias is bounded by $\sqrt{\rho} U \sigma_{h,n}(z)$.

Noise term (Hilbert-space self-normalized bound). Consider the random element $N(z) := E^{\top}(K_n + \rho I)^{-1}k_n(z) = \sum_{i=1}^n \alpha_i(z)\,\varepsilon_i \in \mathcal{H}_{\ell}$ with $\alpha(z) = (K_n + \rho I)^{-1}k_n(z)$. We will show that, with probability at least $1 - \delta$, simultaneously for all $z \in \mathcal{Z}$,

$$\|\mathsf{N}(z)\|_{\mathcal{H}_{\ell}} \leq \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}} \ \sigma_{h,n}(z) \tag{21}$$

Derivation. For any fixed z, write $\mathsf{N}(z) = \sum_{i=1}^n \alpha_i(z) \varepsilon_i$. Let $\langle \cdot, \cdot \rangle$ denote the inner product in \mathcal{H}_ℓ and let $\mathbb{S} := \{u \in \mathcal{H}_\ell : \|u\|_{\mathcal{H}_\ell} = 1\}$. By duality,

$$\|\mathsf{N}(z)\|_{\mathcal{H}_{\ell}} = \sup_{u \in \mathbb{S}} \sum_{i=1}^{n} \alpha_{i}(z) \langle \varepsilon_{i}, u \rangle.$$

Define, for each $u \in \mathbb{S}$, the scalar martingale difference sequence $\xi_i^{(u)} := \langle \varepsilon_i, u \rangle$, which is conditionally σ -sub-Gaussian (by assumption) and satisfies $|\xi_i^{(u)}| \le \kappa_\ell$ a.s. Let $\boldsymbol{\xi}^{(u)} := (\xi_1^{(u)}, \dots, \xi_n^{(u)})^\top$. Then

$$\sum_{i=1}^{n} \alpha_i(z) \, \xi_i^{(u)} = k_n(z)^{\top} (K_n + \rho I)^{-1} \boldsymbol{\xi}^{(u)}.$$

We invoke the standard *kernel self-normalized concentration* for adaptively chosen designs (the proof appears below as: for any $\delta \in (0,1)$, with probability at least $1-\delta$,

$$\left| k_n(z)^\top (K_n + \rho I)^{-1} \boldsymbol{\xi}^{(u)} \right| \le \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}} \ \sigma_{h,n}(z), \tag{22}$$

simultaneously for all $z \in \mathcal{Z}$ and fixed $u \in \mathbb{S}$. The inequality equation 22 is proved below. Since the right-hand side does not depend on u, taking the supremum over $u \in \mathbb{S}$ yields equation 21.

Proof of equation 22. Fix $u \in \mathbb{S}$. Let $A_n := K_n + \rho I$ and note that $\gamma(n,\rho) = \frac{1}{2} \log \det(I + \rho^{-1}K_n) = \frac{1}{2} \log \det(A_n) - \frac{n}{2} \log \rho$. For any $\lambda > 0$, by the sub-Gaussian mgf bound and the fact that the design may be adaptive but A_n is \mathcal{F}_n -measurable, one can show (see Abbasi-Yadkori et al. (2011); Chowdhury & Gopalan (2017)) the *mixture* supermartingale

$$\mathcal{M} := \exp\left(\frac{1}{2\sigma^2} \boldsymbol{\xi}^{(u)\top} A_n^{-1} \boldsymbol{\xi}^{(u)}\right) \left(\frac{\rho^{n/2}}{\det(A_n)^{1/2}}\right)$$

satisfies $\mathbb{E}[\mathcal{M}] \leq 1$ (this is the standard Laplace method; see, e.g., the scalar KRR analyses for kernelized bandits). By Markov's inequality, see proof Lemma E.2,

$$\Pr\left(\boldsymbol{\xi}^{(u)\top} A_n^{-1} \boldsymbol{\xi}^{(u)} \geq 2\sigma^2 \left(\gamma(n,\rho) + \log \frac{1}{\delta}\right)\right) \leq \delta$$

On this event, for any z,

$$\left|k_n(z)^{\top} A_n^{-1} \boldsymbol{\xi}^{(u)}\right| \leq \|A_n^{-1/2} k_n(z)\|_2 \cdot \|A_n^{-1/2} \boldsymbol{\xi}^{(u)}\|_2 \leq \sqrt{2} \, \sigma \, \sqrt{\gamma(n,\rho) + \log \frac{1}{\delta}} \, \|A_n^{-1/2} k_n(z)\|_2.$$

Finally, using the identity

$$\sigma_{h,n}^2(z) \ = \ k(z,z) - k_n(z)^\top A_n^{-1} k_n(z) \ = \ k(z,z) - \|A_n^{-1/2} k_n(z)\|_2^2$$

and the inequality $||A_n^{-1/2}k_n(z)||_2 \le \rho^{-1/2}\sqrt{k(z,z)-k_n(z)^\top A_n^{-1}k_n(z)}$ (which follows from $A_n \succeq \rho I$), we obtain

$$||A_n^{-1/2}k_n(z)||_2 \le \frac{1}{\sqrt{\rho}} \, \sigma_{h,n}(z)$$

Combining the last two displays gives equation 22, completing the proof of the scalar self-normalized bound.

Combine bias and noise. From equation 19, equation 20, and equation 21, with probability at least $1 - \delta$,

$$\|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}} \le \sqrt{\rho} U \, \sigma_{h,n}(z) + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log \frac{1}{\delta}} \, \sigma_{h,n}(z)$$

simultaneously for all $z \in \mathcal{Z}$, as claimed.

Lemma E.2 (Self-normalized tail bound by Markov). Let $(\mathcal{F}_t)_{t=0}^n$ be a filtration and let $\boldsymbol{\xi}^{(u)} = (\xi_1, \dots, \xi_n)^{\top}$ be an \mathcal{F}_t -adapted martingale difference sequence that is conditionally σ -sub-Gaussian: $\mathbb{E}[\exp\{\lambda \xi_t\} \mid \mathcal{F}_{t-1}] \leq \exp(\frac{\sigma^2 \lambda^2}{2})$ for all $\lambda \in \mathbb{R}$ and $t=1,\dots,n$. Let $A_n \in \mathbb{R}^{n \times n}$ be \mathcal{F}_n -measurble, symmetric positive definite (e.g., $A_n = K_n + \rho I$ with ridge $\rho > 0$ and a design-dependent Gram matrix $K_n \succeq 0$). Define the (design-dependent) information term (Scarlett & Bogunovic (2018); Srinivas et al. (2010))

$$\gamma(n,\rho) := \frac{1}{2} \log \frac{\det(A_n)}{\rho^n} = \frac{1}{2} \log \det \left(\mathbf{I} + \rho^{-1} K_n \right).$$

Then for every $\delta \in (0,1)$,

$$\mathbb{P}\!\!\left(\boldsymbol{\xi}^{(u)\top}A_n^{-1}\boldsymbol{\xi}^{(u)} \;\geq\; 2\sigma^2\!\left(\gamma(n,\rho) + \log\tfrac{1}{\delta}\right)\right) \;\leq\; \delta$$

Proof. Consider the *mixture/Laplace* supermartingale (proved, e.g., in Abbasi-Yadkori et al. (2011); Chowdhury & Gopalan (2017))

$$\mathcal{M} \,:=\, \exp\Bigl(\frac{1}{2\sigma^2}\, \boldsymbol{\xi}^{(u)\top} A_n^{-1} \boldsymbol{\xi}^{(u)}\Bigr) \, \Bigl(\frac{\rho^{n/2}}{\det(A_n)^{1/2}}\Bigr), \qquad \text{which satisfies} \quad \mathbb{E}[\mathcal{M}] \leq 1.$$

 Fix $\delta \in (0,1)$. By the definition of $\gamma(n,\rho)$, $\exp{\{\gamma(n,\rho)\}} = \det(A_n)^{1/2}/\rho^{n/2}$. Therefore, the event $\boldsymbol{\xi}^{(u)\top}A_n^{-1}\boldsymbol{\xi}^{(u)} > 2\sigma^2(\gamma(n,\rho) + \log\frac{1}{5})$

is equivalent to

$$\exp\left(\frac{1}{2\sigma^2}\boldsymbol{\xi}^{(u)\top}A_n^{-1}\boldsymbol{\xi}^{(u)}\right) \geq \exp\{\gamma(n,\rho)\} \cdot \frac{1}{\delta} = \frac{\det(A_n)^{1/2}}{\rho^{n/2}} \cdot \frac{1}{\delta}$$

$$\iff \mathcal{M} \geq \frac{1}{\delta}.$$

Hence,

$$\mathbb{P}\!\!\left(\boldsymbol{\xi}^{(u)\top}A_n^{-1}\boldsymbol{\xi}^{(u)} \;\geq\; 2\sigma^2\!\left(\gamma(n,\rho) + \log\tfrac{1}{\delta}\right)\right) \;=\; \mathbb{P}(\mathcal{M} \geq \delta^{-1}) \;\leq\; \delta\,\mathbb{E}[\mathcal{M}] \;\leq\; \delta$$

where we used Markov's inequality in the first inequality and $\mathbb{E}[\mathcal{M}] \leq 1$ in the second. This proves the claim.

Remark E.3 (Interpretation). When $A_n = K_n + \rho I$ with $\rho > 0$, the quantity $\gamma(n, \rho) = \frac{1}{2} \log \det(I + \rho^{-1}K_n)$ coincides with the standard information gain in kernel bandits/GP regression; the lemma is the usual self-normalized tail bound obtained directly from the mixture supermartingale via Markov

F UNIFORM CI

Theorem F.1 (Uniform CI for all $||V||_{\mathcal{H}_{\ell}} \leq B$). Under conditions of Lemma 3.2, for any B > 0 and $\delta \in (0,1)$, with probability at least $1 - \delta$, for all $V \in \mathcal{H}_{\ell}$ with $||V||_{\mathcal{H}_{\ell}} \leq B$ and all $z \in \mathcal{Z}$,

$$\left| [P_h V](z) - \widehat{f}_{h,n}^V(z) \right| \leq \beta_{n,\delta} \, \sigma_{h,n}(z), \qquad \beta_{n,\delta} \ := \ B\left(\sqrt{\rho} \, U \ + \ \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right).$$

Proof. We proceed in three steps and keep the step index h implicit to lighten notation. Throughout, recall the following definitions:

(i) (Bellman image as a CME inner product) Under Assumption 2.1, for every $V \in \mathcal{H}_{\ell}$ and $z \in \mathcal{Z}$,

$$[P_h V](z) = \langle \mu(z), V \rangle_{\mathcal{H}_{\ell}}, \tag{23}$$

where $\mu: \mathcal{Z} \to \mathcal{H}_{\ell}$ is the conditional mean embedding (CME) with $\|\mu\|_{\mathcal{H}_k \otimes \mathcal{H}_{\ell}} \leq U$.

(ii) (Scalar and vector KRR) Given data $\{(z_i, s_i')\}_{i=1}^n$, define the scalar KRR predictor for labels $y_i^{(V)} := V(s_i')$

$$\widehat{f}_{h,n}^{V}(z) = k_n(z)^{\top} (K_n + \rho I)^{-1} \boldsymbol{y}^{(V)}, \qquad \sigma_{h,n}^{2}(z) = k(z,z) - k_n(z)^{\top} (K_n + \rho I)^{-1} k_n(z), \tag{24}$$

and the vector-valued KRR CME estimator

$$\widehat{\mu}_n(z) := \sum_{i=1}^n \alpha_i(z) \, \phi(s_i'), \qquad \alpha(z) := (K_n + \rho I)^{-1} k_n(z). \tag{25}$$

(iii) (Scalar-vector identity) By Proposition 3.1,

$$\widehat{f}_{h,n}^{V}(z) = \langle \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}} \quad \text{for all } V \in \mathcal{H}_{\ell}, z \in \mathcal{Z}.$$
 (26)

Step 1: Reduce scalar error to a vector errer via inner products. Combining equation 23 and equation 26, for any $V \in \mathcal{H}_{\ell}$ and $z \in \mathcal{Z}$,

$$[P_h V](z) - \widehat{f}_{h,n}^V(z) = \langle \mu(z), V \rangle_{\mathcal{H}_{\ell}} - \langle \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}} = \langle \mu(z) - \widehat{\mu}_n(z), V \rangle_{\mathcal{H}_{\ell}}$$
(27)

Step 2: Apply Cauchy-Schwarz + take a supremum over RKHS ball. By Cauchy-Schwarz in \mathcal{H}_{ℓ} .

 $|[P_h V](z) - \hat{f}_{h,n}^V(z)| \le ||\mu(z) - \hat{\mu}_n(z)||_{\mathcal{H}_{\ell}} \cdot ||V||_{\mathcal{H}_{\ell}}.$ (28)

Hence, uniformly over all V in the RKHS ball $\{V : ||V||_{\mathcal{H}_{\ell}} \leq B\}$,

$$\sup_{\|V\|_{\mathcal{H}_{\ell}} \le B} \left| [P_h V](z) - \widehat{f}_{h,n}^V(z) \right| \le B \|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}}$$
(29)

Note that the right-hand side depends on the data and on z, but *not* on V; this is the key to obtaining a *uniform* statement over the entire ball.

Step 3: Invoke vector-valued KRR concentration (Lemma 3.2). Lemma 3.2 asserts that, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\|\mu(z) - \widehat{\mu}_n(z)\|_{\mathcal{H}_{\ell}} \le \left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right) \sigma_{h,n}(z) \qquad \text{simultaneously for all } z \in \mathcal{Z}.$$
(30)

Multiplying both sides of equation 30 by B and plugging into equation 29 gives, on the same high-probability event,

$$\sup_{\|V\|_{\mathcal{H}_{\ell}} \leq B} \left| [P_h V](z) - \widehat{f}_{h,n}^V(z) \right| \leq B \left(\sqrt{\rho} \, U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2 \gamma(n,\rho) + 2 \log \frac{1}{\delta}} \, \right) \sigma_{h,n}(z) \qquad \text{for all } z \in \mathcal{Z}.$$

Since the left-hand side is an upper bound on *each* particular V with $\|V\|_{\mathcal{H}_{\ell}} \leq B$, we conclude that, with probability at least $1 - \delta$, *simultaneously for all* V *with* $\|V\|_{\mathcal{H}_{\ell}} \leq B$ *and all* $z \in \mathcal{Z}$,

$$|[P_h V](z) - \widehat{f}_{h,n}^V(z)| \leq \beta_{n,\delta} \, \sigma_{h,n}(z),$$

with

$$\beta_{n,\delta} := B\left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n,\rho) + 2\log\frac{1}{\delta}}\right).$$

This is exactly the claimed bound.

G ADDITIONAL RESULTS

Definition G.1 (ρ -regularized orthogonal projector (Tikhonov projector)). Let \mathcal{H} be a Hilbert space and $\mathcal{S} \subset \mathcal{H}$ a finite-dimensional subspace with basis $\{s_1, \ldots, s_m\}$. For $\rho > 0$, the ρ -regularized orthogonal projector (or Tikhonov projector) $\Pi_{\mathcal{S},\rho}: \mathcal{H} \to \mathcal{S}$ maps any $f \in \mathcal{H}$ to the unique element $g \in \mathcal{S}$ that solves the ridge-regularized least-squares problem

$$g = \arg\min_{h \in \mathcal{S}} \|f - h\|_{\mathcal{H}}^2 + \rho \|h\|_{\mathcal{H}}^2$$

Equivalently, if $S: \mathbb{R}^m \to \mathcal{H}$ denots the synthesis operator $S\mathbf{c} = \sum_{j=1}^m c_j s_j$ and $G = S^*S$ is the Gram matrix of $\{s_i\}$ in \mathcal{H} , then

$$\Pi_{\mathcal{S},\rho} f = S (G + \rho I)^{-1} S^* f$$

which reduces to standard orthogonal projector as $\rho \downarrow 0$ (provided G is invertible).

Lemma G.2 (Bias inequality using Tikhonov interpolaton). Let $K(z,z') = k(z,z') I_{\mathcal{H}_{\ell}}$ be the operator-valued kernel on \mathcal{Z} with scalar kernel k and output space \mathcal{H}_{ℓ} , and let $\mathcal{H}_{k,I}$ denote the associated vector-valued RKHS (isometric to $\mathcal{H}_k \otimes \mathcal{H}_{\ell}$). Given training inputs $z_{1:n}$, define the finite-dimensional subspace

$$S_n := \operatorname{span} \{ K(\cdot, z_i)u : i = 1, \dots, n, u \in \mathcal{H}_{\ell} \} \subset \mathcal{H}_{k,I}$$

and let $\Pi_{n,\rho}: \mathcal{H}_{k,I} \to \mathcal{S}_n$ be the ρ -regularized orthogonal projector (Tikhonv projector) onto \mathcal{S}_n . Let $\mu \in \mathcal{H}_{k,I}$ be the (vector-valued) target and $M^{\top}: \mathbb{R}^n \to \mathcal{H}_{\ell}$ be the linear operator $M^{\top} \mathbf{b} = \sum_{i=1}^n b_i \, \mu(z_i)$. Then, for every $z \in \mathcal{Z}$,

$$\|\mu(z) - M^{\top}(K_n + \rho I)^{-1}k_n(z)\|_{\mathcal{H}_{\ell}} = \|\mu(z) - \Pi_{n,\rho}\mu(z)\|_{\mathcal{H}_{\ell}} \le \sqrt{\rho} \|\mu\|_{\mathcal{H}_{k,I}} \, \sigma_{h,n}(z), \quad (31)$$

where $K_n = [k(z_i, z_j)]_{i,j}$, $k_n(z) = [k(z, z_1), \dots, k(z, z_n)]^\top$, and $\sigma_{h,n}^2(z) = k(z, z) - k_n(z)^\top (K_n + \rho I)^{-1} k_n(z)$.

Proof. We first recall that in the vector-valued RKHS with kernel K=kI, the evaluation functional at z is represented by $K(\cdot,z)=k(\cdot,z)I_{\mathcal{H}_\ell}$, and the *regularized* orthogonal projection $\Pi_{n,\rho}$ onto \mathcal{S}_n satisfies the normal equations (see Lemma G.3)

$$\Pi_{n,\rho}\mu(\cdot) = \sum_{i=1}^{n} K(\cdot, z_i) c_i^{\star}, \quad \text{with} \quad (K_n + \rho I) C^{\star \top} = M^{\top},$$

where $C^* = [c_1^*, \dots, c_n^*]$ and $M^\top : \mathbb{R}^n \to \mathcal{H}_\ell$ maps $e_i \mapsto \mu(z_i)$. Evaluating at z and using $K(\cdot, z_i) = k(\cdot, z_i)I$, we obtain

$$\Pi_{n,\rho}\mu(z) = \sum_{i=1}^{n} k(z, z_i) c_i^{\star} = \left(k_n(z)^{\top} (K_n + \rho I)^{-1}\right) M^{\top} = M^{\top} (K_n + \rho I)^{-1} k_n(z)$$

which proves the first equalty in equation 31.

For the inequality, we use the standard Tikhonov interpolation error bound in RKHSs (vector-valued case with kernel K=kI). Let $g^*=\Pi_{n,\rho}\mu$. Then, for any z,

$$\|\mu(z) - g^{\star}(z)\|_{\mathcal{H}_{\ell}} \leq \|\mu - g^{\star}\|_{\mathcal{H}_{k,I}} \|K(\cdot,z)\|_{\mathcal{H}_{k,I}} \leq \sqrt{\rho} \|\mu\|_{\mathcal{H}_{k,I}} \|(K_n + \rho I)^{-1/2} k_n(z)\|_2,$$

where the last step uses the optimality of g^* for the Tikhonov problem and the standard interpolation inequality (see, e.g., Steinwart & Christmann, 2008; Carmeli et al., 2010, see Lemma G.8 for details). Finally,

$$\|(K_n + \rho I)^{-1/2} k_n(z)\|_2^2 = k_n(z)^{\top} (K_n + \rho I)^{-1} k_n(z) = k(z, z) - \sigma_{h,n}^2(z),$$

and since $K_n + \rho I \succeq \rho I$,

$$\|(K_n + \rho I)^{-1/2} k_n(z)\|_2 \le \frac{1}{\sqrt{\rho}} \sigma_{h,n}(z).$$

Combining the last two displays yields $\|\mu(z) - g^*(z)\|_{\mathcal{H}_{\ell}} \leq \sqrt{\rho} \|\mu\|_{\mathcal{H}_{k,I}} \sigma_{h,n}(z)$, which is equation 31.

Lemma G.3 (Normal equations for the Tikhonov projector onto S_n). Let $K(z, z') = k(z, z') I_{\mathcal{H}_{\ell}}$ be the operator-valued kernel on Z with scalar kernel k and output space \mathcal{H}_{ℓ} , and let $\mathcal{H}_{k,I}$ be the associated vector-valued RKHS. Given inputs $z_{1:n}$, define

$$S_n := \operatorname{span} \{ K(\cdot, z_i)u : i = 1, \dots, n, u \in \mathcal{H}_{\ell} \} \subset \mathcal{H}_{k,I}$$

For $\rho > 0$, the ρ -regularized orthogonal projection $\Pi_{n,\rho}: \mathcal{H}_{k,I} \to \mathcal{S}_n$ of any $g \in \mathcal{H}_{k,I}$ is the (unique) minimizer of

$$\min_{h \in \mathcal{S}_n} \|g - h\|_{\mathcal{H}_{k,I}}^2 + \rho \|h\|_{\mathcal{H}_{k,I}}^2.$$

In particular, for $g = \mu$ and $h(\cdot) = \sum_{i=1}^{n} K(\cdot, z_i) c_i$ with coefficients $c_i \in \mathcal{H}_{\ell}$, optimal coefficients c_i^* satisfy the normal equations

$$(K_n + \rho I) C^{\star \top} = M^{\top} \tag{32}$$

where $K_n = [k(z_i, z_j)]_{i,j=1}^n$, $C^* = [c_1^*, \dots, c_n^*]$, and $M^\top : \mathbb{R}^n \to \mathcal{H}_\ell$ is defined by $M^\top e_i = \mu(z_i)$ Consequently,

$$\Pi_{n,\rho}\mu(\cdot) = \sum_{i=1}^{n} K(\cdot, z_i) c_i^{\star}.$$

Proof. Write $h(\cdot) = \sum_{i=1}^n K(\cdot, z_i) c_i$ with $c_i \in \mathcal{H}_\ell$, and define the synthesis operator $S: \mathcal{H}^n_\ell \to \mathcal{H}_{k,I}$ by $S(c_1, \ldots, c_n) = \sum_{i=1}^n K(\cdot, z_i) c_i$. The objective is

$$J(c_1, \dots, c_n) = \| \mu - SC \|_{\mathcal{H}_{b,I}}^2 + \rho \| SC \|_{\mathcal{H}_{b,I}}^2, \qquad C = (c_1, \dots, c_n) \in \mathcal{H}_{\ell}^n.$$

The RKHS inner product with kernl K = kI implies $S^*S = K_n \otimes I_{\mathcal{H}_\ell}$ and $S^*\mu = (\mu(z_1), \dots, \mu(z_n))$, i.e., $M^\top : \mathbb{R}^n \to \mathcal{H}_\ell$ maps $e_i \mapsto \mu(z_i)$ (see Lemma G.4). Expanding and taking the Fréchet derivative with respect to C yields the normal equations (see Lemma G.6)

$$(S^*S + \rho I) C^* = S^*\mu,$$

or equivalently,

$$((K_n \otimes I_{\mathcal{H}_{\ell}}) + \rho I) C^* = M$$

where we regard C^\star as a vector in \mathcal{H}^n_ℓ and $M=(\mu(z_1),\ldots,\mu(z_n))$. Grouping by coordinates in \mathcal{H}_ℓ gives equation 32: $(K_n+\rho I)\,C^{\star\top}=M^\top$ Substituting C^\star back into $h=SC^\star$ shows that the minimizer is $\Pi_{n,\rho}\mu(\cdot)=\sum_{i=1}^n K(\cdot,z_i)\,c_i^\star$. Uniqueness follows from strict convexity of J for $\rho>0$.

Lemma G.4 (Adjoint identities for synthesis operator). Let $K(z, z') = k(z, z') I_{\mathcal{H}_{\ell}}$ be the operator-valued kernel on \mathcal{Z} with scalar kernel k and output Hilbert space \mathcal{H}_{ℓ} , and let $\mathcal{H}_{k,I}$ be the associated vector-valued RKHS. Fix inputs $z_{1:n}$ and define the synthesis operator

$$S: \mathcal{H}_{\ell}^n \longrightarrow \mathcal{H}_{k,I}, \qquad S(c_1,\ldots,c_n) := \sum_{i=1}^n K(\cdot,z_i) c_i = \sum_{i=1}^n k(\cdot,z_i) c_i$$

Then its adjoint $S^*: \mathcal{H}_{k,I} \to \mathcal{H}^n_{\ell}$ satisfies

$$S^*S = K_n \otimes I_{\mathcal{H}_{\ell}}, \qquad S^*\mu = (\mu(z_1), \dots, \mu(z_n)),$$

where $K_n = [k(z_i, z_j)]_{i,j=1}^n$, $\mu : \mathcal{Z} \to \mathcal{H}_{\ell}$ is any \mathcal{H}_{ℓ} -valued function, and \otimes denotes the Kronecker product (acting as the identity on \mathcal{H}_{ℓ}).

Proof. We characterize S^* using the defining relation $\langle SC, g \rangle_{\mathcal{H}_{k,I}} = \langle C, S^*g \rangle_{\mathcal{H}^n_\ell}$ for all $C = (c_1, \ldots, c_n) \in \mathcal{H}^n_\ell$ and $g \in \mathcal{H}_{k,I}$. First, by the reproducing property in the vector-valued RKHS with kernel K = kI (see Lemma G.5),

$$\langle K(\cdot, z_i) c_i, g \rangle_{\mathcal{H}_{k,I}} = \langle c_i, g(z_i) \rangle_{\mathcal{H}_{\ell}}$$

Summing over *i*,

$$\langle SC, g \rangle_{\mathcal{H}_{k,I}} = \sum_{i=1}^{n} \langle c_i, g(z_i) \rangle_{\mathcal{H}_{\ell}} = \langle C, (g(z_1), \dots, g(z_n)) \rangle_{\mathcal{H}_{\ell}^n}.$$

Hence $S^*g = (g(z_1), \dots, g(z_n)) \in \mathcal{H}_{\ell}^n$.

Now take $g = SC' = \sum_{j=1}^n K(\cdot, z_j) c'_j$ with $C' = (c'_1, \dots, c'_n) \in \mathcal{H}^n_\ell$. Then

$$S^*SC' = (SC')(z_1), \dots, (SC')(z_n) = \left(\sum_{j=1}^n K(z_1, z_j)c'_j, \dots, \sum_{j=1}^n K(z_n, z_j)c'_j\right)$$
(33)

$$= \left(\sum_{j=1}^{n} k(z_1, z_j)c'_j, \dots, \sum_{j=1}^{n} k(z_n, z_j)c'_j\right).$$
(34)

This is exactly $(K_n \otimes I_{\mathcal{H}_{\ell}}) C'$, proving $S^*S = K_n \otimes I_{\mathcal{H}_{\ell}}$.

Finally, for any $\mu: \mathcal{Z} \to \mathcal{H}_{\ell}$, $S^*\mu = (\mu(z_1), \dots, \mu(z_n))$, by the first identity with $g = \mu$. Writing $M^{\top}: \mathbb{R}^n \to \mathcal{H}_{\ell}$ for the linear map $M^{\top}e_i = \mu(z_i)$, this is the same as the stacked vector of evaluations.

Lemma G.5 (Vector-valued reproducing property for K = k I). Let $K(z, z') = k(z, z') I_{\mathcal{H}_{\ell}}$ be the operator-valued kernel on \mathcal{Z} , where k is a scalar positive-definite kernel and $I_{\mathcal{H}_{\ell}}$ is the identity on the Hilbert space \mathcal{H}_{ℓ} . Let $\mathcal{H}_{k,I}$ be the associated vector-valued RKHS of \mathcal{H}_{ℓ} -valued functions on \mathcal{Z} . Then for every $z \in \mathcal{Z}$, $c \in \mathcal{H}_{\ell}$, and $g \in \mathcal{H}_{k,I}$,

$$\langle K(\cdot, z) c, g \rangle_{\mathcal{H}_{k,I}} = \langle c, g(z) \rangle_{\mathcal{H}_{\ell}}$$

Proof. By definition of a vector-valued RKHS with kernel K, the evaluation at z is a bounded linear functional from $\mathcal{H}_{k,I}$ to \mathcal{H}_{ℓ} , represented by $K(\cdot,z)$ in the sense that for all $g \in \mathcal{H}_{k,I}$,

$$g(z) = \langle g, K(\cdot, z) \rangle_{\mathcal{H}_{k,I}},$$

where the right-hand side is an element of \mathcal{H}_{ℓ} obtained by the Riesz representation (here, the inner product in $\mathcal{H}_{k,I}$ takes values in \mathcal{H}_{ℓ} when pairing with $K(\cdot,z)$). Concretely, for any $c\in\mathcal{H}_{\ell}$, taking inner products with c in \mathcal{H}_{ℓ} yields

$$\langle c, g(z) \rangle_{\mathcal{H}_{\ell}} = \langle c, \langle g, K(\cdot, z) \rangle_{\mathcal{H}_{k,I}} \rangle_{\mathcal{H}_{\ell}} = \langle g, K(\cdot, z) c \rangle_{\mathcal{H}_{k,I}}$$

where last equality uses bilinearity and the fact that $K(\cdot,z)$ acts on c via $I_{\mathcal{H}_{\ell}}$. Symmetry of the inner product gives the displayed identity. (For a formal construction, see the standard vector-valued RKHS references; e.g., Carmeli, De Vito, and Toigo, 2010.)

Lemma G.6 (Normal equations for ridge in coefficient space). Let $K(z, z') = k(z, z') I_{\mathcal{H}_{\ell}}$ be the operator-valued kernel on \mathcal{Z} with scalar kernel k and output Hilbert space \mathcal{H}_{ℓ} , and let $\mathcal{H}_{k,I}$ be the associated vector-valued RKHS. Fix inputs $z_{1:n}$ and define the synthesis operator

$$S: \mathcal{H}_{\ell}^n \longrightarrow \mathcal{H}_{k,I}, \qquad S(c_1,\ldots,c_n) := \sum_{i=1}^n K(\cdot,z_i) c_i.$$

Equip \mathcal{H}_{ℓ}^n with the product inner product $\langle C, D \rangle_{\mathcal{H}_{\ell}^n} = \sum_{i=1}^n \langle c_i, d_i \rangle_{\mathcal{H}_{\ell}}$ for $C = (c_1, \dots, c_n)$, $D = (d_1, \dots, d_n)$. For a target $\mu \in \mathcal{H}_{k,I}$ and ridge parameter $\rho > 0$, consider the Tikhonov objective in coefficient space

$$J(C) := \| \mu - SC \|_{\mathcal{H}_{k,I}}^2 + \rho \| C \|_{\mathcal{H}_{\ell}^n}^2, \qquad C \in \mathcal{H}_{\ell}^n$$

Then J is strictly convex and also Fréchet differentiable, and its unique minimizer C^* would satisfy the normal equations

$$(S^*S + \rho I) C^* = S^*\mu, \tag{35}$$

where $S^*: \mathcal{H}_{k,I} \to \mathcal{H}^n_{\ell}$ is the adjoint of S. Moreover, using the identities $S^*S = K_n \otimes I_{\mathcal{H}_{\ell}}$ and $S^*\mu = (\mu(z_1), \dots, \mu(z_n)) =: M$ (cf. Lemma G.4), equation 35 is equivalent to

$$((K_n \otimes I_{\mathcal{H}_{\ell}}) + \rho I) C^{\star} = M, \tag{36}$$

with $K_n = [k(z_i, z_j)]_{i,j=1}^n$

Proof. Fréchet derivative. For any direction $D \in \mathcal{H}_{\ell}^n$ and $\varepsilon \in \mathbb{R}$,

$$J(C + \varepsilon D) = \|\mu - SC - \varepsilon SD\|_{\mathcal{H}_{k, I}}^{2} + \rho \|C + \varepsilon D\|_{\mathcal{H}_{n}}^{2}$$

Differentiate at $\varepsilon = 0$ (Gâteaux/Fréchet derivative), we have

$$\frac{\mathrm{d}}{\mathrm{d}\varepsilon}J(C+\varepsilon D)\Big|_{\varepsilon=0} = -2\langle \mu - SC, SD\rangle_{\mathcal{H}_{k,I}} + 2\rho\langle C, D\rangle_{\mathcal{H}_{\ell}^{n}}$$

$$= -2\langle S^{*}(\mu - SC), D\rangle_{\mathcal{H}_{\ell}^{n}} + 2\rho\langle C, D\rangle_{\mathcal{H}_{\ell}^{n}}$$

$$= 2\langle (-S^{*}\mu + S^{*}SC + \rho C), D\rangle_{\mathcal{H}_{\ell}^{n}}$$

where we have used definition of adjoint S^* defind as $\langle SC, g \rangle_{\mathcal{H}_{k,I}} = \langle C, S^*g \rangle_{\mathcal{H}_{\ell}^n}$. The gradient of J at C is therefore $\nabla J(C) = 2((S^*S + \rho I)C - S^*\mu)$.

Optimality and normal equations. Since J is strictly convex (sum of a convex quadratic and a strongly convex quadratic), it has a unique minimizer C^* characterized by $\nabla J(C^*) = 0$, i.e.

$$(S^*S + \rho I) C^* = S^* \mu$$

which is equation 35.

Equivalence to Gram form. By Lemma G.4, $S^*S = K_n \otimes I_{\mathcal{H}_{\ell}}$ and $S^*\mu = (\mu(z_1), \dots, \mu(z_n)) =: M$. Substituting these into equation 35 yields equation 36.

Remark G.7 (Form of Tikhonov projection). Let $S:\mathcal{H}^n_\ell\to\mathcal{H}_{k,I}$ be the synthesis operator $S(c_1,\ldots,c_n)=\sum_{i=1}^n K(\cdot,z_i)\,c_i$, and let $C^\star\in\mathcal{H}^n_\ell$ be the unique solution of the normal equations

$$(S^*S + \rho I) C^* = S^* \mu$$
 (equivalently, $((K_n \otimes I_{\mathcal{H}_\ell}) + \rho I) C^* = M$)

By definition of the ρ -regularized orthogonal projection $\Pi_{n,\rho}$ onto the span $S_n = \operatorname{span}\{K(\cdot,z_i)u: i \in [n], u \in \mathcal{H}_\ell\}$, the minimizer of $\min_{h \in S_n} \|\mu - h\|_{\mathcal{H}_{k,I}}^2 + \rho \|h\|_{\mathcal{H}_{k,I}}^2$ is $h^* = SC^*$. Therefore,

$$\Pi_{n,\rho}\mu(\cdot) = h^{\star}(\cdot) = \sum_{i=1}^{n} K(\cdot, z_i) c_i^{\star}$$

In words: Tikhonov projector onto S_n retains the finite-span form with coefficients given by the ridge normel equations.

Lemma G.8 (Tikhonov interpolation bound in the vector-valued RKHS). Let $K(z,z') = k(z,z') I_{\mathcal{H}_{\ell}}$ be an operator-valued kernel on \mathcal{Z} with scalar kernel k and output Hilbert space \mathcal{H}_{ℓ} , and let $\mathcal{H}_{k,I}$ be the associated vector-valued RKHS. Fix training inputs $z_{1:n}$ and ridge $\rho > 0$. For a target $\mu \in \mathcal{H}_{k,I}$, let

$$g^* := \Pi_{n,\rho}\mu \in \operatorname{span}\{K(\cdot, z_i)u : i \in [n], u \in \mathcal{H}_\ell\}$$

be the ρ -regularized orthogonel projection of μ onto the finite span (the Tikhonov projector). Then, for every $z \in \mathcal{Z}$,

$$\|\mu(z) - g^{\star}(z)\|_{\mathcal{H}_{\ell}} \le \sqrt{\rho} \|\mu\|_{\mathcal{H}_{k,I}} \|(K_n + \rho I)^{-1/2} k_n(z)\|_2$$
(37)

where $K_n = [k(z_i, z_j)]_{i,j=1}^n$ and $k_n(z) = [k(z, z_1), \dots, k(z, z_n)]^{\top}$

Proof. Step 1: A residual representer. Define the linear *evaluation* functional at z by $E_z: \mathcal{H}_{k,I} \to \mathcal{H}_{\ell}, E_z(h) = h(z)$. Let $S: \mathcal{H}^n_{\ell} \to \mathcal{H}_{k,I}$ be the synthesis operator $S(c_1, \ldots, c_n) = \sum_{i=1}^n K(\cdot, z_i)c_i$, and $S^*: \mathcal{H}_{k,I} \to \mathcal{H}^n_{\ell}$ its adjoint (Lemma G.4 gives $S^*h = (h(z_1), \ldots, h(z_n))$ and $S^*S = K_n \otimes I_{\mathcal{H}_{\ell}}$). Let $\alpha(z) := (K_n + \rho I)^{-1}k_n(z)$ and define the *residual representer*

$$r_z(\cdot) := K(\cdot, z) - S\alpha(z) \in \mathcal{H}_{k,I}$$
(38)

By vector-valued reproducing prop. (Lemma G.5), for any $h \in \mathcal{H}_{k,I}$,

$$\langle h, r_z \rangle_{\mathcal{H}_{k,I}} = \langle h, K(\cdot, z) \rangle_{\mathcal{H}_{k,I}} - \langle h, S\alpha(z) \rangle_{\mathcal{H}_{k,I}} = \langle h(z), \cdot \rangle_{\mathcal{H}_{\ell}} - \langle S^*h, \alpha(z) \rangle_{\mathcal{H}_{\ell}^n}$$
$$= h(z) - \sum_{i=1}^n \alpha_i(z) h(z_i).$$

In particular, for $h = \mu$ and $h = g^*$, we obtain

$$\mu(z) - g^{\star}(z) = \left\langle \mu - g^{\star}, r_z \right\rangle_{\mathcal{H}_{k,I}}.$$
 (39)

Step 2: Tikhonov orthogonality and swapping the residual. By optimality of $g^{\star} = \Pi_{n,\rho}\mu$ for the Tikhonov problem $\min_{h \in \text{span}} \| \mu - h \|_{\mathcal{H}_{k,I}}^2 + \rho \|h\|_{\mathcal{H}_{k,I}}^2$, the Fréchet first-order condition reads

$$\langle \mu - g^{\star}, SC \rangle_{\mathcal{H}_{\bullet}, \Gamma} + \rho \langle g^{\star}, SC \rangle_{\mathcal{H}_{\bullet}, \Gamma} = 0 \quad \forall C \in \mathcal{H}_{\ell}^{n}.$$

Equivalently, with S^* , $S^*(\mu - g^*) = -\rho S^*g^*$, and therefore, for every z,

$$\langle \mu - g^{\star}, S \alpha(z) \rangle_{\mathcal{H}_{k,I}} = \langle S^{*}(\mu - g^{\star}), \alpha(z) \rangle_{\mathcal{H}_{\ell}^{n}} = -\rho \langle S^{*}g^{\star}, \alpha(z) \rangle_{\mathcal{H}_{\ell}^{n}}$$

Thus equation 39 can be rewrittn as

$$\mu(z) - g^{\star}(z) = \left\langle \mu - g^{\star}, K(\cdot, z) \right\rangle_{\mathcal{H}_{k, I}} + \rho \left\langle S^{*} g^{\star}, \alpha(z) \right\rangle_{\mathcal{H}_{k}^{n}}.$$

Using $g^* = SC^*$ and the normal equations $(S^*S + \rho I)C^* = S^*\mu$ (Lemma G.6), one checks that the second term equals $\rho \langle C^*, \alpha(z) \rangle_{\mathcal{H}^n_\ell} = \langle S^*\mu - S^*SC^*, \alpha(z) \rangle = \langle \mu - g^*, S\alpha(z) \rangle_{\mathcal{H}_{k,I}}$. Therefore

$$\mu(z) - g^{\star}(z) = \langle \mu - g^{\star}, K(\cdot, z) - S \alpha(z) \rangle_{\mathcal{H}_{k,I}} = \langle \mu - g^{\star}, r_z \rangle_{\mathcal{H}_{k,I}}$$

This recovers equation 39 and shows r_z as the Riesz representer of the linear functional $h \mapsto h(z) - \sum_i \alpha_i(z)h(z_i)$

Step 3: Bounding residual via the powar function. By using Cauchy-Schwarz,

$$\|\mu(z) - g^{\star}(z)\|_{\mathcal{H}_{\ell}} \leq \|\mu - g^{\star}\|_{\mathcal{H}_{k,I}} \|r_z\|_{\mathcal{H}_{k,I}}$$

A standard computatin (the "power function" calculation; see, e.g., Steinwart & Christmann, 2008, or Carmeli et al., 2010) gives

$$||r_z||_{\mathcal{H}_{k,I}}^2 = \langle r_z, r_z \rangle_{\mathcal{H}_{k,I}} = \rho ||(K_n + \rho I)^{-1/2} k_n(z)||_2^2$$

whence

$$||r_z||_{\mathcal{H}_{k,I}} = \sqrt{\rho} ||(K_n + \rho I)^{-1/2} k_n(z)||_2$$
 (40)

Finally, Tikhonov optimality inequalty $\|\mu - g^*\|_{\mathcal{H}_{k,I}}^2 + \rho \|g^*\|_{\mathcal{H}_{k,I}}^2 \leq \|\mu\|_{\mathcal{H}_{k,I}}^2$ implies $\|\mu - g^*\|_{\mathcal{H}_{k,I}} \leq \|\mu\|_{\mathcal{H}_{k,I}}$. Combining with equation 40 yields equation 37.

Remark G.9 (On the power functin identity). The equality $||r_z||_{\mathcal{H}_{k,I}}^2 = \rho ||(K_n + \rho I)^{-1/2} k_n(z)||_2^2$ follows from expanding $r_z = K(\cdot,z) - S(K_n + \rho I)^{-1} k_n(z)$ in the RKHS inner product, using $S^*S = K_n \otimes I_{\mathcal{H}_\ell}$ and $S^*K(\cdot,z) = k_n(z)$ (Lemma G.4), and the matrix identity $(K_n + \rho I)^{-1} K_n(K_n + \rho I)^{-1} = (K_n + \rho I)^{-1} - \rho(K_n + \rho I)^{-2}$.

H Additional Results for Theorem 5.1

Lemma H.1 (Global "good event" via a union bound). Fix $\delta \in (0,1)$. For each step $h \in [H]$ and episode $t \in [T]$, let $n_{h,t-1} = |\mathcal{D}_{h,t-1}|$ be the number of transitions collected at step h before episode t, and define the per-step confidence radius (as in equation 3)

$$\beta_{h,t} := B\left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{2HT}{\delta}}\right)$$

Assume algorithm's projection guarantees $||V_{h+1,t}||_{\mathcal{H}_{\ell}} \leq B$ for all h,t. Then there exists an event \mathcal{G} with

$$\Pr(\mathcal{G}) \geq 1 - \delta$$

such that, simultaneously for all $h \in [H]$, $t \in [T]$, and all $z \in \mathcal{Z}$, Eq equation 11 copied below

$$[P_h V_{h+1,t}](z) \le \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \beta_{h,t} \sigma_{h,t}(z). \tag{41}$$

Proof (with union bound). **Step 1: A per-**(h,t) **confidence event.** Fix a particular pair (h,t). Apply the *uniform* confidence theorem (Theorem 3.3) at step h using the dataset $\mathcal{D}_{h,t-1}$ and failure probability

$$\delta_{h,t} := \frac{\delta}{HT}$$

Because the algorithm projects onto the RKHS ball, we have $||V_{h+1,t}||_{\mathcal{H}_{\ell}} \leq B$. Therefore, Theorem 3.3 (with δ replaced by $\delta_{h,t}$ and n replaced by $n_{h,t-1}$) gives a high-probability event $\mathcal{G}_{h,t}$ (depending on the random data collected up to episode t) on which, simultaneously for all $z \in \mathcal{Z}$,

$$\left| \left[P_h V_{h+1,t} \right](z) - \widehat{f}_{h,t}^{V_{h+1,t}}(z) \right| \leq B \left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n_{h,t-1},\rho) + 2\log \frac{HT}{\delta}} \right) \sigma_{h,t}(z).$$

Since the left-hand side is an absolute deviation, it implies the desired one-sided inequality

$$[P_h V_{h+1,t}](z) \leq \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \underbrace{B\left(\sqrt{\rho} U + \frac{\sigma}{\sqrt{\rho}} \sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{HT}{\delta}}\right)}_{\beta_{h,t}^{(\min)}} \sigma_{h,t}(z), \quad \forall z \in \mathcal{Z},$$

with probability at least $1 - \delta_{h,t}$ (i.e., $\Pr(\mathcal{G}_{h,t}) \ge 1 - \delta/(HT)$).

Step 2: Uniformity across all (h, t) by a union bound. There are at most HT such pairs (h, t). The union bound³ yields

$$\Pr\left(\bigcap_{h=1}^{H}\bigcap_{t=1}^{T}\mathcal{G}_{h,t}\right) \geq 1 - \sum_{h,t}\Pr(\mathcal{G}_{h,t}^{c}) \geq 1 - HT \cdot \frac{\delta}{HT} = 1 - \delta.$$

³If events E_1, \ldots, E_m each fail with probability at most ϵ , then $\Pr(\bigcap_i E_i) \geq 1 - m\epsilon$.

Let $\mathcal{G} := \bigcap_{h,t} \mathcal{G}_{h,t}$; then $\Pr(\mathcal{G}) \ge 1 - \delta$ and, on \mathcal{G} , the one-sided bound above holds for *every* pair (h,t) and *every* z.

Step 3: Using the slightly larger radius in equation 3. In the algorithm we instantiate the per-step radius with the slightly larger log factor,

$$\beta_{h,t} = B\left(\sqrt{\rho}U + \frac{\sigma}{\sqrt{\rho}}\sqrt{2\gamma(n_{h,t-1},\rho) + 2\log\frac{2HT}{\delta}}\right) \geq \beta_{h,t}^{(\min)},$$

since $\log\left(\frac{2HT}{\delta}\right) \geq \log\left(\frac{HT}{\delta}\right)$. Using a *larger* (more conservative) radius can only make the inequality easier to satisfy. Therefore, on the same event \mathcal{G} ,

$$[P_h V_{h+1,t}](z) \leq \widehat{f}_{h,t}^{V_{h+1,t}}(z) + \beta_{h,t} \, \sigma_{h,t}(z) \qquad \text{for all } h,t \text{ and all } z \in \mathcal{Z},$$

which is exactly equation 11.

I ADDITIONAL RESULTS

Lemma I.1 (Finite-dimensional reduction of the RKHS projection). Let $(\mathcal{H}_{\ell}, \langle \cdot, \cdot \rangle_{\mathcal{H}_{\ell}})$ be an RKHS with the reproducing kernel $\ell: \mathcal{S} \times \mathcal{S} \to \mathbb{R}$. We fix atoms $\bar{s}_1, \dots, \bar{s}_{m_h} \in \mathcal{S}$ and we write Gram matrix $L_h \in \mathbb{R}^{m_h \times m_h}$ as $(L_h)_{ij} = \ell(\bar{s}_i, \bar{s}_j)$. For a target vector $v_{h,t} \in \mathbb{R}^{m_h}$, we consider the (empirical) projection problem over the feasible class

$$\mathcal{F} := \{ V \in \mathcal{H}_{\ell} : \|V\|_{\mathcal{H}_{\ell}} \le B, \ 0 \le V(\bar{s}_j) \le U \ \forall j \in [m_h] \}, \qquad U := H - h + 1.$$

That is,

$$\min_{V \in \mathcal{F}} \frac{1}{m_h} \sum_{j=1}^{m_h} \left(V(\bar{s}_j) - v_{h,t}(j) \right)^2 \tag{42}$$

Then there exists an optimal solution of the form $V^*(\cdot) = \sum_{j=1}^{m_h} \alpha_j \, \ell(\cdot, \bar{s}_j)$ and, by parameterizing by $\alpha \in \mathbb{R}^{m_h}$, equation 42 is equivalent to the convex quadratic program

$$\min_{\alpha \in \mathbb{R}^{m_h}} \frac{1}{m_h} \left\| L_h \alpha - v_{h,t} \right\|_2^2 \quad \text{s.t.} \quad \alpha^\top L_h \alpha \leq B^2, \qquad 0 \leq (L_h \alpha)_j \leq U \ \forall j \in [m_h]. \tag{43}$$

Moreover, equation 5 is a convex program: its objective has PSD Hessian $\frac{2}{m_h}L_h^{\top}L_h$, the quadratic constraint uses the PSD matrix $L_h \succeq 0$, and the box constraints are linear.

Proof. Let $\mathcal{H}_S := \operatorname{span}\{\ell(\cdot, \bar{s}_j): j \in [m_h]\} \subseteq \mathcal{H}_\ell$ and let $P_S : \mathcal{H}_\ell \to \mathcal{H}_S$ denote orthogonel projection (in RKHS inner product). For any $V \in \mathcal{H}_\ell$, write the orthogonal decomposition $V = P_S V + (I - P_S)V =: V_S + V_\perp$ with $V_S \in \mathcal{H}_S$ and $V_\perp \in \mathcal{H}_S^\perp$.

(i) Loss depends only on V_S . By the reproducing property, for every j,

$$V_{\perp}(\bar{s}_j) \; = \; \langle V_{\perp}, \, \ell(\cdot, \bar{s}_j) \rangle_{\mathcal{H}_{\ell}} \; = \; 0 \quad \text{since} \; \; \ell(\cdot, \bar{s}_j) \in \mathcal{H}_S \; \perp \; V_{\perp}$$

Hence we have $V(\bar{s}_j) = V_S(\bar{s}_j)$ for all j, so the empirical loss in equation 42 equals $\frac{1}{m_h} \sum_j (V_S(\bar{s}_j) - v_{h,t}(j))^2$, independent of V_\perp

- (ii) Feasibility is preserved (and improved) by dropping V_{\perp} . The box constraints $0 \leq V(\bar{s}_j) \leq U$ involve only the evaluations at \bar{s}_j and thus are unchanged when replacing V by V_S (by (i)). For the norm constraint, $\|V\|_{\mathcal{H}_{\ell}}^2 = \|V_S\|_{\mathcal{H}_{\ell}}^2 + \|V_{\perp}\|_{\mathcal{H}_{\ell}}^2 \geq \|V_S\|_{\mathcal{H}_{\ell}}^2$, so $\|V\| \leq B$ implies $\|V_S\| \leq B$.
- (iii) Reduction to \mathcal{H}_S . Given any feasible V, the function V_S is also feasible and achieves the same objective value; therefore an optimal solution exists in \mathcal{H}_S
- (iv) Parameterization by coefficients. Every $V \in \mathcal{H}_S$ can be written as $V(\cdot) = \sum_{j=1}^{m_h} \alpha_j \, \ell(\cdot, \bar{s}_j)$ for some $\alpha \in \mathbb{R}^{m_h}$. The vector of evaluations at the atoms is then

$$(V(\bar{s}_1), \dots, V(\bar{s}_{m_h}))^{\top} = L_h \alpha, \qquad (L_h)_{ij} = \ell(\bar{s}_i, \bar{s}_j)$$

The RKHS norm satisfies $\|V\|_{\mathcal{H}_{\ell}}^2 = \sum_{i,j} \alpha_i \alpha_j \, \ell(\bar{s}_i, \bar{s}_j) = \alpha^\top L_h \alpha$ (standard RKHS identity). Substituting these relations into equation 42 and the constraints yields equation 5. (v) Convexity. Since L_h is a (symmetric) Gram metrix, $L_h \succeq 0$. The objective $\frac{1}{m_h} \|L_h \alpha - v_{h,t}\|_2^2$ is convex with Hessian $\frac{2}{m_h}L_h^{\top}L_h \succeq 0$. The quadratic constraint $\alpha^{\top}L_h\alpha \leq B^2$ defines a convex set because the quadratic form is convex for $L_h \succeq 0$. The bounds $0 \leq (L_h \alpha)_j \leq U$ are linear inequalities in α Thus equation 5 is a convex quadratic program. Remark I.2 (Representer viewpoint). Argument above is a constrained version of the representer theorem: because both the objective and the constraints depend on V only through its evaluations at $\{ar{s}_i\}$ and its RKHS norm, the optimizer lies in the span of kernel sections at these points Kimeldorf & Wahba (1971); Schölkopf & Smola (2002); Schölkopf et al. (2001).