

The Inductive Bias of Quantum Kernels

Jonas M. Kübler* Simon Buchholz* Bernhard Schölkopf
 Max Planck Institute for Intelligent Systems
 Tübingen, Germany
 {jmkuebler, sbuchholz, bs}@tue.mpg.de

June 8, 2021

Abstract

It has been hypothesized that quantum computers may lend themselves well to applications in machine learning. In the present work, we analyze function classes defined via *quantum kernels*. Quantum computers offer the possibility to efficiently compute inner products of exponentially large density operators that are classically hard to compute. However, having an exponentially large feature space renders the problem of generalization hard. Furthermore, being able to evaluate inner products in high dimensional spaces efficiently by itself does not guarantee a quantum advantage, as already classically tractable kernels can correspond to high- or infinite-dimensional reproducing kernel Hilbert spaces (RKHS).

We analyze the spectral properties of quantum kernels and find that we can expect an advantage if their RKHS is low dimensional and contains functions that are hard to compute classically. If the target function is known to lie in this class, this implies a quantum advantage, as the quantum computer can encode this *inductive bias*, whereas there is no classically efficient way to constrain the function class in the same way. However, we show that finding suitable quantum kernels is not easy because the kernel evaluation might require exponentially many measurements.

In conclusion, our message is a somewhat sobering one: we conjecture that quantum machine learning models can offer speed-ups only if we manage to encode knowledge about the problem at hand into quantum circuits, while encoding the same bias into a classical model would be hard. These situations may plausibly occur when learning on data generated by a quantum process, however, they appear to be harder to come by for classical datasets.

1 Introduction

In recent years, much attention has been dedicated to studies of how small and noisy quantum devices [1] could be used for near term applications to showcase the power of quantum computers. Besides fundamental demonstrations [2], potential applications that have been discussed are in quantum chemistry [3], discrete optimization [4] and machine learning (ML) [5–12].

Initiated by the seminal HHL algorithm [13], early work in quantum machine learning (QML) was focused on speeding up linear algebra subroutines, commonly used in ML, offering the perspective of a runtime logarithmic in the problem size [14–17]. However, most of these works have an inverse polynomial scaling of the runtime in the error and it was shown rigorously

* JMK and SB contributed equally and are ordered randomly.

by Ciliberto et al. [18] that due to the quantum mechanical measurement process a runtime complexity $O(\sqrt{n})$ is necessary for convergence rate $1/\sqrt{n}$.

Rather than speeding up linear algebra subroutines, we focus on more recent suggestions that use a quantum device to define and implement the function class and do the optimization on a classical computer. There are two ways to that: the first are so-called *Quantum Neural Networks* (QNN) or parametrized quantum circuits [5–7] which can be trained via gradient based optimization [5, 19–23]. The second approach is to use a predefined way of encoding the data in the quantum system and defining a *quantum kernel* as the inner product of two quantum states [7–11]. These two approaches essentially provide a parametric and a non-parametric path to quantum machine learning, which are closely related to each other [11]. Since the optimization of QNNs is non-convex and suffers from so-called Barren Plateaus [24], we here focus on quantum kernels, which allow for convex problems and thus lend themselves more readily to theoretical analysis.

The central idea of using a QML model is that it enables to do computations that are exponentially hard classically. However, also in classical ML, kernel methods allow us to implicitly work with high- or infinite dimensional function spaces [25, 26]. Thus, purely studying the expressivity of QML models [27] is not sufficient to understand when we can expect speed-ups. Only recently first steps were taken into this direction [10, 12, 28]. Assuming classical hardness of computing discrete logarithms, Liu et al. [10] proposed a task based on the computation of the discrete logarithm where the quantum computer, equipped with the right feature mapping, can learn the target function with exponentially less data than any classical (efficient) algorithm. Similarly, Huang et al. [12] analyzed generalization bounds and realized that the expressivity of quantum models can hinder generalization. They proposed a heuristic to optimize the labels of a dataset such that it can be learned well by a quantum computer but not a classical machine.

In this work, we relate the discussion of quantum advantages to the classical concept of *inductive bias*. The *no free lunch* theorem informally states that no learning algorithm can outperform other algorithms on all problems. This implies that an algorithm that performs well on one type of problem necessarily performs poorly on other problems. A standard inductive bias in ML is to prefer functions that are continuous. An algorithm with that bias, however, will then struggle to learn functions that are discontinuous. For a QML model to have an edge over classical ML models, we could thus ensure that it is equipped with an inductive bias that cannot be encoded (efficiently) with a classical machine. If a given dataset fits this inductive bias, the QML model will outperform any classical algorithm. For kernel methods, the qualitative concept of inductive bias can be formalized by analyzing the spectrum of the kernel and relating it to the target function [25, 29–33].

Our main contribution is the analysis of the inductive bias of quantum machine learning models based on the spectral properties of quantum kernels. First, we show that quantum kernel methods will fail to generalize as soon as the data embedding into the quantum Hilbert space is too expressive (Theorem 1). Then we note that projecting the quantum kernel appropriately allows to construct inductive biases that are hard to create classically (Figure 1). However, our Theorem 2 also implies that estimating the biased kernel requires exponential measurements, a phenomenon reminiscent of the Barren plateaus observed in quantum neural networks. Finally we show experiments supporting our main claims.

While our work gives guidance to find a quantum advantage in ML, this yields no recipe for obtaining a quantum advantage on a classical dataset. We conjecture that unless we have a clear idea how the data generating process can be described with a quantum computer, we cannot expect an advantage by using a quantum model in place of a classical machine learning model.

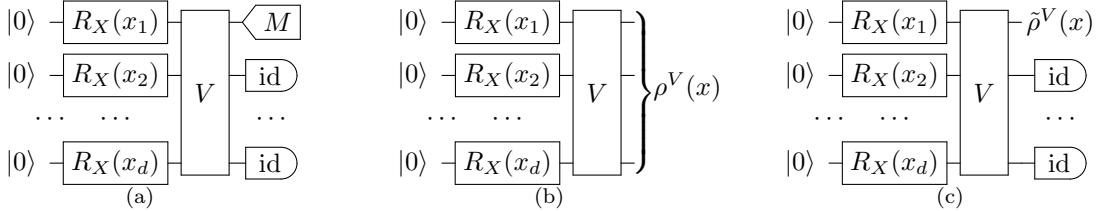


Figure 1: **Quantum advantage via inductive bias:** (a) Data generating quantum circuit $f(x) = \text{Tr} [\rho^V(x)(M \otimes \text{id})] = \text{Tr} [\tilde{\rho}^V(x)M]$. (b) The full quantum kernel $k(x, x') = \text{Tr} [\rho^V(x)\rho^V(x')]$ is too general and cannot learn f efficiently. (c) The biased quantum kernel $q(x, x') = \text{Tr} [\tilde{\rho}^V(x)\tilde{\rho}^V(x')]$ meaningfully constrains the function space and allows to learn f with little data.

2 Supervised learning

We briefly introduce the setting and notation for supervised learning as a preparation for our analysis of quantum mechanical methods in this context. The goal of supervised learning is the estimation of a functional mechanism based on data generated from this mechanism. For concreteness we focus on the regression setting where we assume data is generated according to $Y = f^*(X) + \varepsilon$ where ε denotes zero-mean noise. We focus on $X \in \mathcal{X} \subset \mathbb{R}^d$, $Y \in \mathbb{R}$. We denote the joint probability law of (X, Y) by \mathcal{D} and we are given n i.i.d. observations D_n from \mathcal{D} . We will refer to the marginal law of X as μ , define the L_μ^2 inner product $\langle f, g \rangle = \int f(x)g(x) \mu(dx)$ and denote the corresponding norm by $\|\cdot\|$. The least square risk and the empirical risk of some hypothesis $h : \mathcal{X} \rightarrow \mathbb{R}$ is defined by $R(h) = \mathbb{E}_{\mathcal{D}} [(h(X) - Y)^2]$ and $R_n(h) = \mathbb{E}_{D_n} [(h(X) - Y)^2]$.

In supervised machine learning, one typically considers a hypothesis space H of functions $h : \mathcal{X} \rightarrow \mathbb{R}$ and tries to infer $\text{argmin}_{h \in H} R(h)$ (assuming for simplicity that the minimizer exists). Typically this is done by (regularized) empirical risk minimization $\text{argmin}_{h \in H} R_n(h) + \lambda \Omega(h)$, where $\lambda > 0$ and Ω determine the regularization. The risk of h can then be decomposed in generalization and training error $R(h) = (R(h) - R_n(h)) + R_n(h)$.

Kernel ridge regression. We will focus on solving the regression problem over a reproducing kernel Hilbert space (RKHS) [25, 26]. An RKHS \mathcal{F} associated with a positive definite kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is the space of functions such that for all $x \in \mathcal{X}$ and $h \in \mathcal{F}$ the *reproducing* property $h(x) = \langle h, k(x, \cdot) \rangle_{\mathcal{F}}$ holds. Kernel ridge regression regularizes the RKHS norm, i.e., $\Omega(h) = \|h\|_{\mathcal{F}}^2$. With observations $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ we can compute the kernel matrix $K(X, X)_{ij} = k(x^{(i)}, x^{(j)})$ and the Representer Theorem [34] ensures that the empirical risk minimizer of kernel ridge regression is of the form $\hat{f}_n^\lambda(\cdot) = \sum_{i=1}^n \alpha_i k(x^{(i)}, \cdot)$, with $\alpha = (K(X, X) + \lambda \text{id})^{-1}y$. The goal of our work is to study when a (quantum) kernel is suitable for learning a particular problem. The central object to study this is the integral operator.

Spectral properties and inductive bias. For kernel k and marginal law μ , the integral operator K , is defined as $(Kf)(x) = \int k(x, x')f(x')\mu(dx')$. Mercer's Theorem ensures that there exist a spectral decomposition of K with (possibly infinitely many) eigenvalues γ_i (ordered non-increasingly) and corresponding eigenfunctions ϕ_i , which are orthonormal in L_μ^2 , i.e., $\langle \phi_i, \phi_j \rangle = \delta_{i,j}$. We will assume that $\text{Tr}[K] = \sum_i \gamma_i = 1$ which we can ensure by rescaling the kernel. We can then write $k(x, x') = \sum_i \gamma_i \phi_i(x)\phi_i(x')$. While the functions ϕ form a basis of \mathcal{F} they might not

completely span L^2_μ . In this case we simply complete the basis and implicitly take $\gamma = 0$ for the added functions. Then we can decompose functions in L^2_μ as

$$f(x) = \sum_i a_i \phi_i(x). \quad (1)$$

We have $\|f\|^2 = \sum_i a_i^2$ and $\|f\|_{\mathcal{F}}^2 = \sum_i \frac{a_i^2}{\gamma_i}$ (if $f \in \mathcal{F}$). Kernel ridge regression penalizes the RKHS norm of functions. The components corresponding to zero eigenvalues are infinitely penalized and cannot be learned since they are not in the RKHS. For large regularization λ the solution \hat{f}_n^λ is heavily *biased* towards learning only the coefficients of the principal components (corresponding to the largest eigenvalues) and keeps the other coefficients small (at the risk of *underfitting*). Decreasing the regularization allows ridge regression to also fit the other components, however, at the potential risk of overfitting to the noise in the empirical data. Finding good choices of λ thus balances this *bias-variance* tradeoff.

We are less concerned with the choice of λ , but rather with the spectral properties of a kernel that allow for a quantum advantage. Similar to the above considerations, a target function f can easily be learned if it is well *aligned* with the principal components of a kernel. In the easiest case, the kernel only has a single non-zero eigenvalue and is just $k(x, x') = f(x)f(x')$. Such a construction is arguably the simplest path to a quantum advantage in ML.

Example 1 (Trivial Quantum Advantage). Let f be a scalar function that is easily computable on a quantum device but requires exponential resources to approximate classically. Generate data as $Y = f(X) + \epsilon$. The kernel $k(x, x') = f(x)f(x')$ then has an exponential advantage for learning f from data.

To go beyond this trivial case, we introduce two qualitative measures to judge the quality of a kernel for learning the function f . The *kernel target alignment* of Cristianini et al. [30] is

$$A(k, f) = \frac{\langle k, f \otimes f \rangle}{\langle k, k \rangle^{1/2} \langle f \otimes f, f \otimes f \rangle^{1/2}} = \frac{\sum_i \gamma_i a_i^2}{(\sum_i \gamma_i^2)^{1/2} \sum_i a_i^2} \quad (2)$$

and measures how well the kernel fits f . If $A = 1$, learning reduces to estimating a single real parameter, whereas for $A = 0$, learning is infeasible. We note that the kernel target alignment also weighs the contributions of f depending on the corresponding eigenvalue, i.e., the alignment is better if large $|a_i|$ correspond to large γ_i . The kernel target alignment was used extensively to optimize kernel functions [31] and recently also used to optimize quantum kernels [35].

In a similar spirit, the *task-model alignment* of Canatar et al. [32] measures how much of the signal of f is captured in the first i principal components: $C(i) = \sum_{j \leq i} a_j^2 (\sum_j a_j^2)^{-1}$. The slower $C(i)$ approaches 1, the harder it is to learn as the target function is more spread over the eigenfunctions.

3 Quantum computation in machine learning

In this section we introduce hypothesis spaces containing functions whose output is given by the result of a quantum computation. For a general introduction to concepts of quantum computation we refer to the book of Nielsen and Chuang [36].

We will consider quantum systems comprising $d \in \mathbb{N}$ qubits. Discussing such systems and their algebraic properties does not require in-depth knowledge of quantum mechanics. A *pure state* of a single qubit is described by vector $(\alpha, \beta)^\top \in \mathbb{C}^2$ s.t. $|\alpha|^2 + |\beta|^2 = 1$ and we write

$|\psi\rangle = \alpha|0\rangle + \beta|1\rangle$, where $\{|0\rangle, |1\rangle\}$ forms the computational basis. A d qubit pure state lives in the tensor product of the single qubit state spaces, i.e., it is described by a normalized vector in \mathbb{C}^{2^d} . A *mixed state* of a d -qubit system can be described by a density operator $\rho \in \mathbb{C}^{2^d \times 2^d}$, i.e., a positive definite matrix ($\rho \geq 0$) with unit trace ($\text{Tr}[\rho] = 1$). For a pure state $|\psi\rangle$ the corresponding density operator is $\rho = |\psi\rangle\langle\psi|$ (here, $\langle\psi|$ is the complex conjugate transpose of $|\psi\rangle$). A general density operator can be thought of as a classical probabilistic mixture of pure states. We can extract information from ρ by estimating (through repeated measurements) the expectation of a suitable *observable*, i.e., a Hermitian operator $M = M^\dagger$ (where the adjoint $(\cdot)^\dagger$ is the complex conjugate of the transpose), as

$$\text{Tr}[\rho M]. \quad (3)$$

Put simply, the potential advantage of a quantum computer arises from its state space being exponentially large in the number of qubits d , thus computing general expressions like (3) on a classical computer is exponentially hard. However, besides the huge obstacles in building quantum devices with high fidelity, the fact that the outcome of the quantum computation (3) has to be estimated from measurements often prohibits to easily harness this power. We will discuss this in the context of quantum kernels in Section 4.

We consider parameter dependent quantum states $\rho(x) = U(x)\rho_0U^\dagger(x)$ that are generated by evolving the initial state ρ_0 with the data dependent unitary transformation $U(x)$ [7, 11]. Most often we will without loss of generality assume that the initial state is $\rho_0 = (|0\rangle\langle 0|)^{\otimes d}$. We then define quantum machine learning models via observables M of the data dependent state

$$f_M(x) = \text{Tr}[U(x)\rho_0U^\dagger(x)M] = \text{Tr}[\rho(x)M]. \quad (4)$$

In the following we introduce the two most common function classes suggested for quantum machine learning. We note that there also exist proposals that do not fit into the form of Eq. (4) [27, 35, 37].

Quantum neural networks. A "quantum neural network" (QNN) is defined via a *variational quantum circuit* (VQC) [6, 38, 39]. Here the observable M_θ is parametrized by $p \in \mathbb{N}$ classical parameters $\theta \in \Theta \subseteq \mathbb{R}^p$. This defines a parametric function class $\mathcal{F}_\Theta = \{f_{M_\theta} | \theta \in \Theta\}$. The most common ansatz is to consider $M_\theta = U(\Theta)MU^\dagger(\Theta)$ where $U(\Theta) = \prod_i U(\theta_i)$ is the composition of unitary evolutions each acting on few qubits. For this and other common models of the parametric circuit it is possible to analytically compute gradients and specific optimizers for quantum circuits based on gradient descent have been developed [5, 19–23]. Nevertheless, the optimization is usually a non-convex problem and suffers from additional difficulties due to oftentimes exponentially (in d) vanishing gradients [24]. This hinders a theoretical analysis. Note that the non-convexity does not arise from the fact that the QNN can learn non-linear functions, but rather because the observable M_θ depends non-linearly on the parameters. In fact, the QNN functions are linear in the *fixed* feature mapping $\rho(x)$. Therefore the analogy to classical neural networks is somewhat incomplete.

Quantum kernels. The class of functions we consider are RKHS functions where the kernel is expressed by a quantum computation. The key observation is that (4) is linear in $\rho(x)$. Instead of optimizing over the parametric function class \mathcal{F}_Θ , we can define the nonparametric class of functions $\mathcal{F} = \{f_M | f_M(\cdot) = \text{Tr}[\rho(\cdot)M], M = M^\dagger\}$.¹ To endow this function class with the

¹ \mathcal{F} is defined for a fixed feature mapping $x \mapsto \rho(x)$. Although M is finite dimensional and thus \mathcal{F} can be seen as a parametric function class, we will be interested in the case where M is exponentially large in d and we can only access functions from \mathcal{F} implicitly. Therefore we refer to it as nonparametric class of functions.

structure of an RKHS, observe that the expression $\text{Tr}[\rho_1 \rho_2]$ defines a scalar product on density matrices. We then define kernels via the inner product of data-dependent density matrices:

Definition 1 (Quantum Kernel [7, 8, 11]). Let $\rho : x \mapsto \rho(x)$ be a fixed feature mapping from \mathcal{X} to density matrices. Then the corresponding *quantum kernel* is $k(x, x') = \text{Tr}[\rho(x)\rho(x')]$.

The Representer Theorem [34] reduces the empirical risk minimization over the exponentially large function class \mathcal{F} to an optimization problem with a set of parameters whose dimensionality corresponds to the training set size. Since the ridge regression objective is convex (and so are many other common objective functions in ML), this can be solved efficiently with a classical computer.

In the described setting, the quantum computer is only used to estimate the kernel. For pure state encodings, this is done by inverting the data encoding transformation (taking its conjugate transpose) and measuring the probability that the resulting state equals the initial state ρ_0 . To see this we use the cyclic property of the trace $k(x, x') = \text{Tr}[\rho(x)\rho(x')] = \text{Tr}[U(x)\rho_0 U^\dagger(x) U(x')\rho_0 U^\dagger(x')] = \text{Tr}[(U^\dagger(x')U(x)\rho_0 U^\dagger(x)U(x'))\rho_0]$. If $\rho_0 = (|0\rangle\langle 0|)^{\otimes d}$, then $k(x, x')$ corresponds to the probability of observing every qubit in the '0' state after the initial state was evolved with $U^\dagger(x')U(x)$. To evaluate the kernel, we thus need to estimate this probability from a finite number of measurements. For our theoretical analysis we work with the exact value of the kernel and in our experiments we also simulate the full quantum state. We discuss the difficulties related to measurements in Sec. 4.

4 The inductive bias of simple quantum kernels

We now study the inductive bias for simple quantum kernels and their learning performance. We first give a high level discussion of a general hurdle for quantum machine learning models to surpass classical methods and then analyze two specific kernel approaches in more detail.

Continuity in classical machine learning. Arguably the most important bias in nonparametric regression are continuity assumptions on the regression function. This becomes particularly apparent in, e.g., nearest neighbour regression or random regression forests [40] where the regression function is a weighted average of close points. Here we want to briefly discuss two potential barriers that might prevent major improvements of quantum machine learning methods compared to classical methods. First there is a long list of results concerning the minimax optimality of kernel methods for regression problems [41–43]. In particular these results show that asymptotically the convergence of kernel ridge regression of, e.g., Sobolev functions reaches the statistical limits which also apply to any quantum method. A second barrier concerns the approximation of quantum methods by classical methods. Suppose we are given a quantum kernel $k(x, x')$ on \mathbb{R}^d that is continuous (so that functions in the corresponding RKHS are continuous) and translation invariant. In this case it is known that the learning performance of k can be well approximated using random Fourier features [44] where ε^{-1} features (up to logarithmic terms) are necessary to get the same risk as kernel ridge regression up to $O(\varepsilon + n^{-1/2})$ [45]. Moreover random Fourier feature regression can be efficiently implemented (in polynomial time) on a classical computer. Therefore efficient classical learning algorithms with similar performance as the quantum model exist in this case. It might, however, be difficult to construct the classical model for a given quantum circuit. For non-translation invariant kernels (which are more relevant in quantum machine learning) random feature methods are less studied but there are some extensions which might be applicable in similar settings [46].

A simple quantum kernel. We now restrict our attention to rather simple kernels to facilitate a theoretical analysis. As indicated above we consider data in $\mathcal{X} \subset \mathbb{R}^d$ and we assume that the distribution μ of the data factorizes over the coordinates (i.e. μ can be written as $\mu = \bigotimes \mu_i$). This data is embedded in a d -qubit quantum circuit. Let us emphasize here that the RKHS based on a quantum state of d -qubits is at most 4^d dimensional, i.e., finite dimensional and in the infinite data limit $n \rightarrow \infty$ standard convergence guarantees from parametric statistics apply. Here we consider growing dimension $d \rightarrow \infty$, and sample size polynomial in the dimension $n = n(d) \in \text{Poly}(d)$. In particular the sample size $n \ll 4^d$ will be much smaller than the dimension of the feature space and bounds from the parametric literature do not apply.

Here we consider embeddings where each coordinate is embedded into a single qubit using the map φ_i followed by an arbitrary unitary transformation V , so that we can express the embedding in the quantum Hilbert space as $|\psi(x)\rangle = V \bigotimes |\varphi_i(x_i)\rangle$ with corresponding density matrix (feature map) $\rho(x) = |\psi(x)\rangle \langle \psi(x)|$. Note that the corresponding kernel $k(x, x') = \text{Tr}[\rho(x)\rho(x')]$ is independent of V and factorizes $k(x, x') = \text{Tr}[\bigotimes \rho_i(x_i) \bigotimes \rho_i(x'_i)] = \prod \text{Tr}[\rho_i(x_i)\rho_i(x'_i)]$ where $\rho_i(x_i) = |\varphi_i(x_i)\rangle \langle \varphi_i(x_i)|$. The product structure of the kernel allows us to characterize the RKHS generated by k based on the one dimensional case. The embedding of a single variable can be parametrized by complex valued functions $a(x), b(x)$ as

$$|\varphi_i(x)\rangle = a(x)|0\rangle + b(x)|1\rangle. \quad (5)$$

One important object characterizing this embedding turns out to be the mean density matrix of this embedding given by $\rho_{\mu_i} = \int \rho_i(y) \mu_i(dy) = \int |\varphi_i(y)\rangle \langle \varphi_i(y)| \mu_i(dy)$. This can be identified with the kernel mean embedding of the distribution [47]. Note that for factorizing base measure μ the factorization $\rho_\mu = \bigotimes \rho_{\mu_i}$ holds. Let us give a concrete example to clarify the setting, see Fig. 1(b).

Example 2. [11, Example III.1.] We consider the cosine kernel where $a(x) = \cos(x/2)$, $b(x) = i \sin(x/2)$. This embedding can be realized using a single quantum $R_X(x) = \exp(-i\frac{x}{2}\sigma_x)$ gate such that $|\psi(x)\rangle = R_X(x)|0\rangle = \cos(x/2)|0\rangle + i \sin(x/2)|1\rangle$. In this case the kernel is given by

$$k(x, x') = |\langle 0 | R_X^\dagger(x) R_X(x) | 0 \rangle|^2 = |\cos(\frac{x}{2}) \cos(\frac{x'}{2}) + \sin(\frac{x}{2}) \sin(\frac{x'}{2})|^2 = \cos(\frac{x-x'}{2})^2. \quad (6)$$

As a reference measure μ we consider the uniform measure on $[-\pi, \pi]$. Then the mean density matrix is the completely mixed state $\rho_\mu = \frac{1}{2} \text{id}$. For \mathbb{R}^d valued data whose coordinates are encoded independently the kernel is given by $k(x, x') = \prod \cos^2((x_i - x'_i)/2)$ and $\rho_\mu = 2^{-d} \text{id}_{2^d \times 2^d}$. We emphasize that due to the kernel trick this kernel can be evaluated classically in runtime $O(d)$.

Quantum RKHS. We now characterize the RKHS and the eigenvalues of the integral operator for quantum kernels. The RKHS consists of all function $f \in \mathcal{F}$ that can be written as $f(x) = \text{Tr}[\rho(x)M]$ where $M \in \mathbb{C}^{2^d \times 2^d}$ is a Hermitian operator. Using this characterization of the finite dimensional RKHS we can rewrite the infinite dimensional eigenvalue problem of the integral operator as a finite dimensional problem. The action of the corresponding integral operator on f can be written as

$$\begin{aligned} (Kf)(x) &= \int f(y)k(y, x) \mu(dy) = \int \text{Tr}[M\rho(y)] \text{Tr}[\rho(y)\rho(x)] \mu(dy) \\ &= \int \text{Tr}[(M \otimes \rho(x))(\rho(y) \otimes \rho(y))] \mu(dy) = \text{Tr} \left[(M \otimes \rho(x)) \int \rho(y) \otimes \rho(y) \mu(dy) \right] \end{aligned} \quad (7)$$

We denote the operator $O_\mu = \int \rho(y) \otimes \rho(y) \mu(dy)$ for which $\text{Tr}[O_\mu] = 1$ holds. Then we can write

$$(Kf)(x) = \text{Tr}[O_\mu(M \otimes \rho(x))] = \text{Tr}[O_\mu(M \otimes \text{id})(\text{id} \otimes \rho(x))] = \text{Tr}[\text{Tr}_1[O_\mu(M \otimes \text{id})]\rho(x)].$$

The eigenvalues of K can now be identified with the eigenvalues of the linear map T_μ mapping $M \rightarrow \text{Tr}_1[O_\mu(M \otimes \text{id})]$. As shown in the appendix there is an eigendecomposition such that $T_\mu(M) = \sum \lambda_i A_i \text{Tr}[A_i M]$ where A_i are orthonormal Hermitian matrices (for details, a proof and an example we refer to Appendix C). The eigenfunctions of K are given by $f_i(x) = \text{Tr}[\rho(x)A_i]$.

We now state a bound that controls the largest eigenvalue of the integral operator K in terms of the eigenvalues of the mean density matrix ρ_μ (Proof in Appendix C.2).

Lemma 1. *The largest eigenvalue γ_{\max} of K satisfies the bound $\gamma_{\max} \leq \sqrt{\text{Tr}[\rho_\mu^2]}$.*

The lemma above shows that the squared eigenvalues of K are bounded by $\text{Tr}[\rho_\mu^2]$, an expression known as the *purity* [36] of the density matrix ρ_μ , which measures the diversity of the data embedding. On the other hand the eigenvalues of K are closely related to the learning guarantees of kernel ridge regression. In particular, standard generalization bounds for kernel ridge regression [48] become vacuous when γ_{\max} is exponentially smaller than the training sample size (if $\text{Tr}[K] = 1$ which holds for pure state embeddings). The next result shows that this is not just a matter of bounds.

Theorem 1. *Suppose the purity of the embeddings ρ_{μ_i} satisfies $\text{Tr}[\rho_{\mu_i}^2] \leq \delta < 1$ and training sample size satisfies $n \leq d^k$. Then there exists $d_0 = d_0(\delta, k, \varepsilon)$ such that for all $d \geq d_0$ no function can be learned using kernel ridge regression with the d -qubit kernel $k(x, x') = \text{Tr}[\rho(x)\rho(x')]$ in the sense that for any f , with probability at least $1 - \varepsilon$ for all λ*

$$R(\hat{f}_n^\lambda) \geq (1 - \varepsilon)\|f\|^2. \quad (8)$$

Theorem 1 (Proof in Appendix D) implies that generalization is only possible when the mean embedding of most coordinates is close to a pure state, i.e. the embedding $x \rightarrow |\varphi_i(x)\rangle$ is almost constant. To make learning from data feasible we cannot use the full expressive power of the quantum Hilbert space but instead only very restricted embeddings allow to learn from data. This generalizes an observation already made in [12]. Since also classical methods allow to handle high-dimensional and infinite dimensional RKHS the same problem occurs for classical kernels where one solution is to adapt the bandwidth of the kernel to control the expressivity of the RKHS. In principle this is also possible in the quantum context, e.g., for the cosine embedding.

Biased kernels. We have discussed that without any inductive bias, the introduced quantum kernel cannot learn any function for large d . One suggestion to reduce the expressive power of the kernel is the use of *projected* kernels [12]. They are defined using reduced density matrices given by $\tilde{\rho}_m^V(x) = \text{Tr}_{m+1\dots d}[\rho^V(x)]$ where $\text{Tr}_{m+1\dots d}[\cdot]$ denotes the partial trace over qubits $m + 1$ to d (definition in Appendix A). Then they consider the usual quantum kernel for this embedding $q_m^V(x, x') = \text{Tr}[\tilde{\rho}_m^V(x)\tilde{\rho}_m^V(x')]$. Physically, this corresponds to just measuring the first m qubits and the functions f in the RKHS can be written in terms of a hermitian operator M acting on m qubits so that $f(x) = \text{Tr}[\rho^V(x)(M \otimes \text{id})] = \text{Tr}[\tilde{\rho}_m^V(x)M]$. If V is sufficiently complex it is assumed that f is hard to compute classically [49].

Indeed above procedure reduces the generalization gap. But this comes at the price of an increased *approximation error* if the remaining RKHS cannot fully express the target function f^* anymore, i.e., the learned function *underfits*. Without any reason to believe that the target function

is well represented via the projected kernel, we cannot hope for a performance improvement by simply reducing the size of the RKHS in an arbitrary way. However, if we *know* something about the data generating process than this can lead to a meaningful inductive bias. For the projected kernel this could be that we *know* that the target function can be expressed as $f^*(x) = \text{Tr} [\tilde{\rho}_m^V(x) M^*]$, see Fig. 1. In this case using q_m^V improves the generalization error without increasing the approximation error. To emphasize this, we will henceforth refer to q_m^V as *biased kernel*.

We now investigate the RKHS for reduced density matrices where V is a Haar-distributed random unitary matrix (proof in Appendix E).

Theorem 2. *Suppose V is distributed according to the Haar measure on the group of unitary matrices. Fix m . Then the following two statements hold:*

- a) *The reduced density operator satisfies with high probability $\tilde{\rho}_m^V = 2^{-m}\text{id} + O(2^{-d/2})$ and the projected kernel satisfies with high probability $q_m^V(x, x') = 2^{-m} + O(2^{-d/2})$ as $d \rightarrow \infty$.*
- b) *Let $T_{\mu, m}^V$ denote the linear integral operator for the kernel q_m^V as defined above. Then the averaged operator $\mathbb{E}_V [T_{\mu, m}^V]$ has one eigenvalue $2^{-m} + O(2^{-2d})$ whose eigenfunction is constant (up to higher order terms of order $O(2^{-2d})$) and $2^m - 1$ eigenvalues $2^{-m-d} + O(2^{-2d})$.*

The second result for the averaged kernel is not in itself meaningful as this operator has no direct meaning. However, we expect a similar result to hold with high probability for fixed V , but a proof would require the evaluation of 8-th order polynomials over the unitary group. Nevertheless the result indicates that for generic rotations V , learning with the biased kernel q_m^V is possible as soon as the training sample size satisfies $n \geq \dim(\mathcal{F})$ (which is bounded by $\dim(\mathcal{F}) \leq 4^m$).

Let us now focus on the case $m = 1$, that is the biased kernel is solely defined via the first qubit. Assuming that Theorem 2b) also holds for fixed V we can assume that the biased kernel has the form

$$q(x, x') \equiv q_1^V(x, x') = \gamma_0 \phi_0(x) \phi_0(x') + \sum_{i=1}^3 \gamma_i \phi_i(x) \phi_i(x'), \quad (9)$$

where $\gamma_0 = 1/2 + O(2^{-2d})$ and $\phi_0(x) = 1$ is the constant function up to terms of order $O(2^{-2d})$. For $i = 1, 2, 3$ we have $\gamma_i = O(2^{-d-1}) = O(2^{-d})$ (Fig. 2) and ϕ_i is a function that conjectured to be exponentially hard in d to compute classically [49]. It is thus impossible to include a bias towards those three eigenfunctions classically. On the other hand we can include a strong bias towards the constant eigenfunction also classically. The straightforward way to do this is to center the data in the RKHS (see Appendix B for details).

Barren plateaus. Another conclusion from Theorem 2a) is that the fluctuations of the reduced density matrix around its mean are exponentially vanishing in the number of qubits. In practice the value of the kernel would be determined by measurements and exponentially many measurements are necessary to obtain exponential accuracy of the kernel function. Therefore the theorem suggests that it is not possible in practice to learn anything beyond the constant function from generic biased kernels for (modestly) large values of d . This observation is closely related to the fact that for many quantum neural networks architectures, the gradient of the parameters with respect to the loss is exponentially vanishing with the system size d , a phenomenon known as *Barren plateaus* [24, 50].

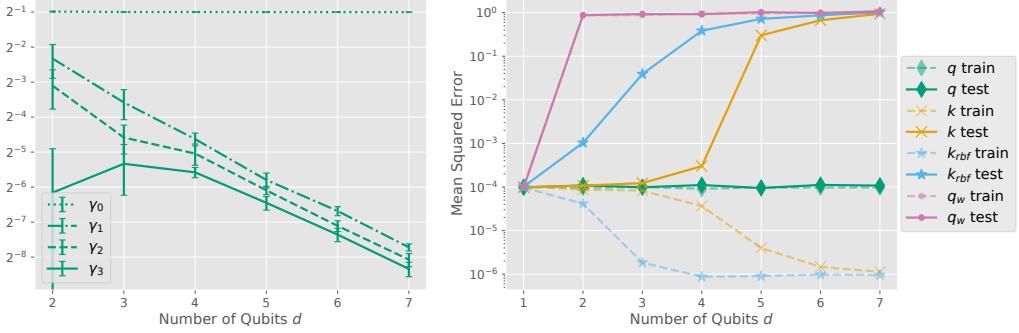


Figure 2: **Left:** Spectral behavior of biased kernel q , see Theorem 2b) and Equation (9) **Right:** The biased kernel q , equipped with prior knowledge, easily learns the function for arbitrary number of qubits and achieves optimal mean squared error (MSE). Models that are ignorant to the structure of f^* fail to learn the function.

5 Experiments

Since for small d we can simulate the biased kernel efficiently, we illustrate our theoretical findings in the following experiments. Our implementation, building on standard open source packages [51, 52], is available in the supplementary information. We consider the case described above where we *know* that the data was generated by measuring an observable on the first qubit, i.e., $f^*(x) = \text{Tr} [\tilde{\rho}_1^V(x) M^*]$, but we do not know M^* , see Fig. 1. We use the full kernel k and the biased kernel q for the case $m = 1$. To show the effect of selecting the wrong bias, we also include the behavior of a biased kernel defined only on the second qubit, denoted as q_w . As a classical reference we also include the performance of a radial basis function kernel $k_{\text{rbf}}(x, x') = \exp(-\|x - x'\|^2/2)$. For the experiments we fix a single qubit observable $M^* = \sigma_z$ and perform the experiment for varying number d of qubits. First we draw a random unitary V . The dataset is then generated by drawing $N = 200$ realizations $\{x^{(i)}\}_{i=1}^N$ from the d dimensional uniform distribution on $[0, 2\pi]^d$. We then define the labels as $y^{(i)} = c f^*(x^{(i)}) + \epsilon^{(i)}$, where $f^*(x) = \text{Tr} [\tilde{\rho}^V(x) \sigma_z]$, $\epsilon^{(i)}$ is Gaussian noise with $\text{Var}[\epsilon] = 10^{-4}$, and c is chosen such that $\text{Var}[f(X)] = 1$. Keeping the variances fixed ensures that we can interpret the behavior for varying d .

We first verify our findings from Theorem 2b) and Equation (9) by estimating the spectrum of q . Fig. 2 (left) shows that Theorem 2b) also holds for individual V with high probability. We then use 2/3 of the data for training kernel ridge regression (we fit the mean separately) with preset regularization, and use 1/3 to estimate the test error. We average the results over ten random seeds (random V , $x^{(i)}, \epsilon^{(i)}$) and results are reported in the right panel of Fig. 2. This showcases that as the number of qubits increases, it is impossible to learn f^* without the appropriate spectral bias. k and k_{rbf} have to little bias and overfit, whereas q_w has the wrong bias and severely underfits. The performance of q_w underlines that randomly biasing the kernel does not significantly improve the performance over the full kernel k . In the appendix we show that this is not due to a bad choice of regularization, by reporting cherry-picked results over a range of regularizations.

To further illustrate the spectral properties, we empirically estimate the kernel target alignment [30] and the task-model alignment [32] that we introduced in Sec. 2. By using the centered kernel matrix (see App. B) and centering the data we can ignore the first eigenvalue in (9) corresponding

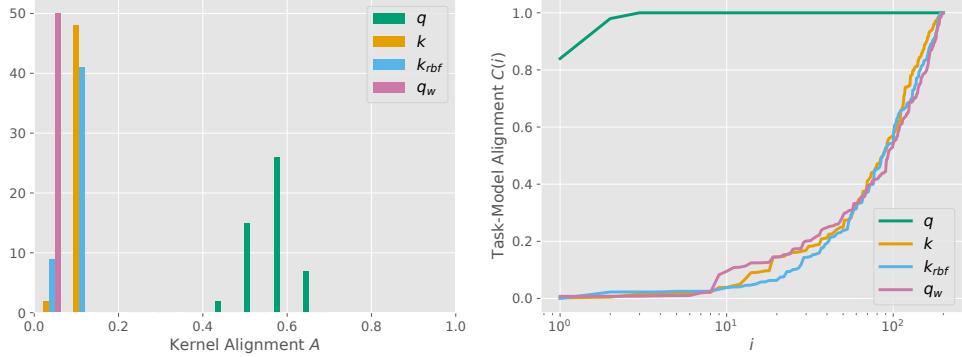


Figure 3: Kernel target alignment (left) and task model alignment (right) for $d = 7$.

the constant function. In Figure 3 (left) we show the empirical (centered) kernel target alignment for 50 random seeds. The biased kernel is the only one well aligned with the task. The right panel of Fig. 3 shows the task model alignment. This shows that f^* can be completely expressed with the first four components of the biased kernel, while the other kernels essentially need the entire spectrum (we use a sample size of 200, hence the empirical kernel matrix is only 200 dimensional) and thus are unable to learn. Note that the kernel q_w is four dimensional, and so higher contributions correspond to functions outside its RKHS that it actually cannot even learn at all.

6 Discussion

We provided an analysis of the reproducing kernel Hilbert space (RKHS) and the inductive bias of quantum kernel methods. Rather than the dimensionality of the RKHS, its spectral properties determine whether learning is feasible. Working with exponentially large RKHS comes with the risk of having a correspondingly small inductive bias. This situation indeed occurs for naive quantum encodings, and hinders learning unless datasets are of exponential size. To enable learning, we necessarily need to consider models with a stronger inductive bias. Encoding a bias towards continuous functions is likely not a promising path for a quantum advantage, as this is where classical machine learning models excel.

Our results suggest that we can only achieve a quantum advantage if we *know* something about the data generating process and cannot efficiently encode this classically, yet are able use this information to bias the quantum model. We indeed observe an exponential advantage in the case where we know that the data comes from a single qubit observable and constrain the RKHS accordingly. However, we find that evaluating the kernel requires exponentially many measurements, an issue related to Barren Plateaus encountered in quantum neural networks.

With fully error-corrected quantum computers it becomes feasible to define kernels with a strong bias that do not require exponentially many measurements. An example of this kind was recently presented by Liu et al. [10]: Here one knows that the target function is extremely simple after computing the discrete logarithm. A quantum computer is able to encode this inductive bias by using an efficient algorithm for computing the discrete logarithm.

However, even for fully coherent quantum computers it is unclear how we can reasonably

encode a strong inductive bias about a classical dataset (e.g., images of cancer cells, climate-data, etc.). The situation might be better when working with *quantum data*, i.e., data that is collected via observations of systems at a quantum mechanical scale. To summarize, we conclude that there is no indication that quantum machine learning will substantially improve supervised learning on classical datasets.

References

- [1] John Preskill. Quantum computing in the NISQ era and beyond. *Quantum*, 2:79, 2018.
- [2] Frank Arute, Kunal Arya, Ryan Babbush, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019.
- [3] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J Love, Alán Aspuru-Guzik, and Jeremy L O’brien. A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, 5:4213, 2014.
- [4] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. A quantum approximate optimization algorithm. *arXiv:1411.4028*, 2014.
- [5] Kosuke Mitarai, Makoto Negoro, Masahiro Kitagawa, and Keisuke Fujii. Quantum circuit learning. *Phys. Rev. A*, 98:032309, 2018.
- [6] Edward Farhi and Hartmut Neven. Classification with quantum neural networks on near term processors. *arXiv:1802.06002*, 2018.
- [7] Vojtěch Havlíček, Antonio D Cárocoles, Kristan Temme, Aram W Harrow, Abhinav Kandala, Jerry M Chow, and Jay M Gambetta. Supervised learning with quantum-enhanced feature spaces. *Nature*, 567(7747):209, 2019.
- [8] Maria Schuld and Nathan Killoran. Quantum machine learning in feature Hilbert spaces. *Phys. Rev. Lett.*, 122:040504, 2019.
- [9] Jonas M. Kübler, Krikamol Muandet, and Bernhard Schölkopf. Quantum mean embedding of probability distributions. *Phys. Rev. Research*, 1:033159, 2019.
- [10] Yunchao Liu, Srinivasan Arunachalam, and Kristan Temme. A rigorous and robust quantum speed-up in supervised machine learning. *arXiv:2010.02174*, 2020.
- [11] Maria Schuld. Quantum machine learning models are kernel methods. *arXiv:2101.11020*, 2021.
- [12] Hsin-Yuan Huang, Michael Broughton, Masoud Mohseni, Ryan Babbush, Sergio Boixo, Hartmut Neven, and Jarrod R. McClean. Power of data in quantum machine learning. *Nature Communications*, 12(1):2631, 2021.
- [13] Aram W Harrow, Avinatan Hassidim, and Seth Lloyd. Quantum algorithm for linear systems of equations. *Phys. Rev. Lett.*, 103(15):150502, 2009.
- [14] Patrick Rebentrost, Masoud Mohseni, and Seth Lloyd. Quantum support vector machine for big data classification. *Phys. Rev. Lett.*, 113(13), 2014.

- [15] Nathan Wiebe, Daniel Braun, and Seth Lloyd. Quantum algorithm for data fitting. *Phys. Rev. Lett.*, 109:050505, 2012.
- [16] Seth Lloyd, Masoud Mohseni, and Patrick Rebentrost. Quantum algorithms for supervised and unsupervised machine learning. *arXiv:1307.0411*, 2013.
- [17] Iordanis Kerenidis and Anupam Prakash. Quantum gradient descent for linear systems and least squares. *Phys. Rev. A*, 101:022316, 2020.
- [18] Carlo Ciliberto, Andrea Rocchetto, Alessandro Rudi, and Leonard Wossnig. Statistical limits of supervised quantum learning. *Physical Review A*, 102(4), 2020.
- [19] Gian Giacomo Guerreschi and Mikhail Smelyanskiy. Practical optimization for hybrid quantum-classical algorithms. *arXiv:1701.01450*, 2017.
- [20] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran. Evaluating analytic gradients on quantum hardware. *Physical Review A*, 99(3):032331, 2019.
- [21] James Stokes, Josh Izaac, Nathan Killoran, and Giuseppe Carleo. Quantum natural gradient. *Quantum*, 4:269, 2020.
- [22] Ryan Sweke, Frederik Wilde, Johannes Jakob Meyer, Maria Schuld, Paul K Fährmann, Barthélémy Meynard-Piganeau, and Jens Eisert. Stochastic gradient descent for hybrid quantum-classical optimization. *Quantum*, 4:314, 2020.
- [23] Jonas M. Kübler, Andrew Arrasmith, Lukasz Cincio, and Patrick J. Coles. An Adaptive Optimizer for Measurement-Frugal Variational Algorithms. *Quantum*, 4:263, 2020.
- [24] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven. Barren plateaus in quantum neural network training landscapes. *Nature Communications*, 9:4812, 2018.
- [25] B. Schölkopf and A. J. Smola. *Learning with Kernels*. MIT Press, Cambridge, MA, USA, 2002.
- [26] John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [27] Maria Schuld, Ryan Sweke, and Johannes Jakob Meyer. Effect of data encoding on the expressive power of variational quantum-machine-learning models. *Phys. Rev. A*, 103:032430, 2021.
- [28] Hsim-Yuan Huang, Richard Kueng, and John Preskill. Information-theoretic bounds on quantum advantage in machine learning. *Phys. Rev. Lett.*, 126:190505, 2021.
- [29] R. C. Williamson, A. J. Smola, and B. Schölkopf. Generalization performance of regularization networks and support vector machines via entropy numbers of compact operators. *IEEE Transactions on Information Theory*, 47(6):2516–2532, 2001.
- [30] Nello Cristianini, Jaz Kandola, Andre Elisseeff, and John Shawe-Taylor. On kernel target alignment. In Dawn E. Holmes and Lakhmi C. Jain, editors, *Innovations in Machine Learning: Theory and Applications*, pages 205–256. Springer Berlin Heidelberg, 2006.
- [31] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13(28):795–828, 2012.

- [32] Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature Communications*, 12(1):2914, 2021.
- [33] Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clement Hongler, and Franck Gabriel. Kernel alignment risk estimator: Risk prediction from training data. In *NeurIPS*, 2020.
- [34] B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *COLT*, 2001.
- [35] Thomas Hubregtsen, David Wierichs, Elies Gil-Fuster, Peter-Jan H. S. Derkx, Paul K. Faehrmann, and Johannes Jakob Meyer. Training quantum embedding kernels on near-term quantum computers. *arXiv:2105.02276*, 2021.
- [36] Michael A. Nielsen and Isaac L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, 2010.
- [37] Adrián Pérez-Salinas, Alba Cervera-Lierta, Elies Gil-Fuster, and José I. Latorre. Data re-uploading for a universal quantum classifier. *Quantum*, 4:226, 2020.
- [38] M. Cerezo, Andrew Arrasmith, Ryan Babbush, et al. Variational quantum algorithms. *arXiv:2012.09265*, 2020.
- [39] Kishor Bharti, Alba Cervera-Lierta, Thi Ha Kyaw, et al. Noisy intermediate-scale quantum (nisq) algorithms. *arXiv:2101.08448*, 2021.
- [40] Leo Breiman. Random forests. *Mach. Learn.*, 45(1):5–32, 2001.
- [41] A. Caponnetto and E. De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- [42] Ingo Steinwart, Don R. Hush, and Clint Scovel. Optimal rates for regularized least squares regression. In *COLT*, 2009.
- [43] Simon Fischer and Ingo Steinwart. Sobolev norm learning rates for regularized least-squares algorithms. *Journal of Machine Learning Research*, 21(205):1–38, 2020.
- [44] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NeurIPS*, 2007.
- [45] Zhu Li, Jean-Francois Ton, Dino Oglic, and Dino Sejdinovic. Towards a unified analysis of random fourier features. In *ICML*, 2019.
- [46] Jean-Francois Ton, Seth Flaxman, Dino Sejdinovic, and Samir Bhatt. Spatial mapping with gaussian processes and nonstationary fourier features. *Spatial Statistics*, 28:59–78, 2018.
- [47] Krikamol Muandet, Kenji Fukumizu, Bharath K. Sriperumbudur, and Bernhard Schölkopf. Kernel mean embedding of distributions: A review and beyond. *Found. Trends Mach. Learn.*, 10(1-2):1–141, 2017.
- [48] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive computation and machine learning, MIT Press, 2012.

- [49] Aram W. Harrow and Ashley Montanaro. Quantum computational supremacy. *Nature*, 549(7671):203–209, 2017.
- [50] M Cerezo, Akira Sone, Tyler Volkoff, Lukasz Cincio, and Patrick J Coles. Cost-function-dependent barren plateaus in shallow quantum neural networks. *arXiv:2001.00550*, 2020.
- [51] F. Pedregosa, G. Varoquaux, Gramfort, et al. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [52] Ville Bergholm, Josh Izaac, Maria Schuld, Christian Gogolin, and Nathan Killoran. Pennylane: Automatic differentiation of hybrid quantum-classical computations. *arXiv:1811.04968*, 2018.
- [53] Vern I. Paulsen and Mrinal Raghupathi. *An Introduction to the Theory of Reproducing Kernel Hilbert Spaces*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2016.
- [54] Z. Puchała and J.A. Miszczak. Symbolic integration with respect to the haar measure on the unitary groups. *Bulletin of the Polish Academy of Sciences: Technical Sciences*, 65(No 1):21–27, 2017.
- [55] Christoph Dankert, Richard Cleve, Joseph Emerson, and Etera Livine. Exact and approximate unitary 2-designs and their application to fidelity estimation. *Phys. Rev. A*, 80:012304, 2009.

The Inductive Bias of Quantum Kernels

Supplementary Material

A Partial trace in quantum mechanics

Here we provide the definition of the partial trace used for the biased quantum kernels. For details we refer to [36]. The state space of the union of two quantum systems with state space \mathcal{H}_1 and \mathcal{H}_2 is given by the tensor product $\mathcal{H}_1 \otimes \mathcal{H}_2$. A general mixed state is described by a density matrix ρ_{12} which is hermitian positive linear operator on $\mathcal{H}_1 \otimes \mathcal{H}_2$ with $\text{Tr}[\rho_{12}] = 1$. The state ρ_1 on the subsystem \mathcal{H}_1 is obtained by the partial trace operation $\rho_1 = \text{Tr}_2[\rho_{12}]$. The partial trace can be defined as the linear map $\text{Tr}_2 : \mathcal{L}(\mathcal{H}_1 \otimes \mathcal{H}_2) \rightarrow \mathcal{L}(\mathcal{H}_1)$ that satisfies

$$\text{Tr}_2[S \otimes T] = \text{Tr}[T]S \quad (10)$$

for all $S \in \mathcal{L}(\mathcal{H}_1)$, $T \in \mathcal{L}(\mathcal{H}_2)$. It can be shown that this map exists and is unique. Picking a basis on \mathcal{H}_1 and \mathcal{H}_2 we consider the tensor product basis on $\mathcal{H}_1 \otimes \mathcal{H}_2$. In coordinates given by this basis we can write

$$(\rho_1)_{i_1 j_1} = \text{Tr}_2[\rho_{12}]_{i_1 j_1} = \sum_{k=1}^{\dim(\mathcal{H}_2)} (\rho_{12})_{i_1 k, j_1 k}. \quad (11)$$

B General results about RKHS

In this section we briefly discuss basic results on centering in RKHS and the RKHS of tensor product kernels.

B.1 Centering in the RKHS

As shown in Section 4, the constant function plays a special role for typical biased kernels as the corresponding eigenvalue is much larger than the remaining eigenvalues. Clearly, it is also possible classically to treat the constant function separately. To do so, it is natural to center the data by subtracting the mean $\bar{y} = n^{-1} \sum_{i=1}^n y_i$ and to consider the *centered* kernel. This corresponds to putting no penalty on the constant function which is also common in linear models where no penalty is put on the intercept. Formally, for a kernel k , the centered kernel is defined as

$$k_c(x, x') = k(x, x') - \mathbb{E}_X[k(X, x')] - \mathbb{E}_{X'}[k(x, X')] + \mathbb{E}_{X, X'}[k(X, X')]. \quad (12)$$

In analogy we can center the kernel matrix as $K_c(X, X) = (\text{id} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)K(X, X)(\text{id} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)$, where $\mathbf{1}$ is a vector of all ones.

Let k be a kernel with Mercer decomposition

$$k(x, x') = \sum \gamma_i \phi_i(x) \phi_i(x'), \quad (13)$$

and define $a_i = \int \phi_i(x) \mu(dx)$. Then the centered kernel can be written as

$$k_c(x, x') = \sum \gamma_i (\phi_i(x) - a_i)(\phi_i(x') - a_i). \quad (14)$$

To make things explicit let us focus on the biased kernel of Equation (9). Ignoring terms of order $\mathcal{O}(2^{-2d})$, the constant function is an eigenfunction of the kernel. In such a case centering corresponds to setting the corresponding eigenvalue γ_0 to zero, while the other terms in the spectral decomposition are invariant under centering (by orthogonality we have $a_i = 0$ for $i \neq 0$). The centered biased kernel of Eq. (9) is thus

$$q_c(x, x') = \sum_{i=1}^3 \gamma_i \phi_i(x) \phi_i(x'). \quad (15)$$

By Theorem 2 we expect that all the eigenvalues of the centered biased kernel are similarly large. Further we know that the centered part of the target function can completely be expressed in terms of the eigenfunctions of the centered biased kernel $f^*(x) - \bar{f}^* = \sum_{i=1}^3 a_i \phi_i(x)$, where $\bar{f}^* = \mathbb{E}[f^*(X)]$. Let us assume that all the three eigenvalues are completely equal. Then we can compute the kernel target alignment of Eq. (2)

$$A(q_c, f^* - \bar{f}^*) = \frac{\sum_{i=1}^3 \gamma a_i^2}{(\sum_{i=1}^3 \gamma^2)^{1/2} \sum_{i=1}^3 a_i^2} = \frac{\gamma \sum_{i=1}^3 a_i^2}{\sqrt{3}\gamma \sum_{i=1}^3 a_i^2} = \frac{1}{\sqrt{3}} \approx 0.58. \quad (16)$$

We emphasize that this expectation is in good accordance with the results of our experiments reported in Fig. 3. Further, note that computing the kernel target alignment after centering is quite common in the kernel literature and is used to optimize the kernel function [31].

B.2 Tensor product of kernels

In this section we describe the construction of product kernels on product spaces. More details can be found in any textbook on RKHS [25]. Let (X_1, k_1) and (X_2, k_2) be two spaces with positive definite kernels with RKHS \mathcal{F}_1 and \mathcal{F}_2 . Then the function

$$k((x_1, x_2), (y_1, y_2)) = k_1(x_1, y_1) k_2(x_2, y_2) \quad (17)$$

defines a positive definite kernel on $X_1 \times X_2$ and the RKHS is given by $\{f_1(x_1)f_2(x_2) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$. Moreover, if we are given a product measure $\mu = \mu_1 \otimes \mu_2$ on $X_1 \times X_2$ then the integral operator for the kernel k factorizes, i.e., for functions $f(x_1, x_2) = f_1(x_1)f_2(x_2)$

$$\begin{aligned} (Kf)(x_1, x_2) &= \int f(y) k(y, x) \mu(dy) \\ &= \int f_1(y_1) k_1(y_1, x_1) \mu_1(dx_1) \int f_2(y_2) k_2(y_2, x_2) \mu_2(dx_2) \\ &= (K_1 f_1)(x_1) (K_2 f_2)(x_2). \end{aligned} \quad (18)$$

Therefore the eigenvalue problems of the integral operators decouple and the eigenvalues of K are given by $\{\gamma^1 \gamma^2 : \gamma^1 \in E_1, \gamma^2 \in E_2\}$ where E_i denotes the eigenvalues of K_i .

It can be derived from the results above that the RKHS of the product kernel $k(x, x') = k_1(x, x') k_2(x, x')$ on X is given by $\{f_1(x)f_2(x) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}$ where \mathcal{F}_i denotes the RKHS of (X, k_i) . There is no simple relation for the integral operators.

C More details on quantum kernels for classical data

In this section we analyze in more detail the properties of quantum kernel methods for classical data.

C.1 Description of the RKHS

To understand the quantum kernel better we give a description of the RKHS for the quantum kernels. We consider the one-qubit embedding $x \rightarrow a(x)|0\rangle + b(x)|1\rangle$. The RKHS $\tilde{\mathcal{F}}$ corresponding to the (non-physical) kernel $\tilde{k}(x, y) = \langle \varphi(x), \varphi(y) \rangle$ is then generated by $a(x), b(x)$. Moreover, the RKHS corresponding to the physical kernel $k(x, x') = \text{Tr}[\rho(x)\rho(x')] = |\langle \varphi(x), \varphi(x') \rangle|^2 = |\tilde{k}(x, x')|^2 = \tilde{k}(x, x')\tilde{k}(x, x')$ is the vector space \mathcal{F} generated by $\{f \cdot \bar{g} : f, g \in \tilde{\mathcal{F}}\}$ [53] (to obtain the real valued RKHS which is more relevant in the learning theoretic setting we consider the real and imaginary part). This result can also be obtained by looking at the feature map $x \rightarrow \rho(x)$ of the physical kernel directly. When we consider data from \mathbb{R}^d where all dimensions are encoded independently in a single qubit the resulting RKHS has the tensor product structure $\mathcal{F} = \bigotimes \mathcal{F}_i$ where \mathcal{F}_i are the RKHS for the single coordinate embeddings.

C.2 Proof of Lemma 1

Here we analyze the integral operators in a bit more detail and prove Lemma 1. For the proof of Lemma 1 we need to briefly look again at the simpler non-physical kernel $\tilde{k}(x, y) = \langle \varphi(x), \varphi(y) \rangle$ and its integral operator. Suppose data has distribution μ on \mathbb{R} . We consider the integral operator \tilde{K} acting on $f(x) = \langle \omega, \varphi(x) \rangle$ defined by

$$\tilde{K}f(x) = \int f(y)\tilde{k}(y, x)\mu(dy) = \int \langle \omega, \varphi(y) \rangle \langle \varphi(y), \varphi(x) \rangle \mu(dy) = \langle \omega, \rho_\mu \varphi(x) \rangle \quad (19)$$

where $\rho_\mu = \int |\varphi(y)\rangle\langle\varphi(y)|\mu(dy)$ denotes the mean density matrix associated with the measure μ . We observe that the eigenvalues γ_i of \tilde{K} agree with the eigenvalues of the density matrix ρ_μ . In particular we conclude

$$\|\tilde{K}\|_{HS}^2 = \sum \gamma_i^2 = \|\rho_\mu\|_{HS}^2 = \text{Tr}[\rho_\mu^2]. \quad (20)$$

This observation corresponds to the fact that for the linear kernel the eigenvalues of the integral operator agree with the eigenvalues of the covariance matrix.

Now we can give the simple proof of Lemma 1. For convenience we restate the lemma.

Lemma 2 (Lemma 1 in the main part). *The largest eigenvalue γ_{max} of K satisfies the bound $\gamma_{max} \leq \sqrt{\text{Tr}[\rho_\mu^2]}$.*

Proof. We observe, denoting the constant function with value 1 by $\mathbf{1}$,

$$\int \mathbf{1}(x)k(x, y)\mathbf{1}(y)\mu(dx)\mu(dy) = \int |\tilde{k}(x, y)|^2\mu(dx)\mu(dy) = \|\tilde{K}\|_{HS}^2 = \|\rho_\mu\|_{HS}^2 = \text{Tr}[\rho_\mu^2] \quad (21)$$

where we used (20) in the last two steps. Suppose that f is a normalized eigenfunction for the eigenvalue γ_{max} . From the Mercer decomposition we obtain

$$1 = K(x, x) \geq \gamma_{max} f(x)^2. \quad (22)$$

Hence f is bounded by $\sqrt{\gamma_{max}}^{-1}$ and we conclude that

$$\begin{aligned} \gamma_{max} &= \int f(x)(Kf)(x)\mu(dx) = \int f(x)k(x, y)f(y)\mu(dx)\mu(dy) \\ &\leq \gamma_{max}^{-1} \int \mathbf{1}(x)k(x, y)\mathbf{1}(y)\mu(dx)\mu(dy) = \gamma_{max}^{-1} \text{Tr}[\rho_\mu^2] \end{aligned}$$

where we used that k is pointwise positive. This ends the proof. \square

Let us look at this result in our main setting where each coordinate of d -dimensional data is embedded in a single qubit. If the measure μ on \mathbb{R}^d factorizes as $\mu = \bigotimes \mu_i$. The integral operator factorizes over the d coordinates and the eigenvalues of the integral operator are given by $\{\prod_{j=1}^d \gamma_{i_j}, \gamma_{i_j} \in E_j\}$ with E_j denoting the eigenvalues of the one-dimensional integral operators. In particular the largest eigenvalue will be exponentially small (in d) as soon as $\max(E_j) \leq \delta < 1$ for a fixed δ which holds if the individual embeddings satisfy $\text{Tr}[\rho_{\mu_i}^2] < \delta$. Note that $\text{Tr}[\rho_{\mu_i}^2] = 1$ if and only if the embedding is constant.

C.3 Spectral decomposition of the integral operator

As shown in the main text the integral operator K applied to $f(x) = \text{Tr}[\rho(x)M]$ can be written as

$$(Kf)(x) = \text{Tr}[O_\mu(M \otimes \rho(x))] = \text{Tr}[O_\mu(M \otimes \text{id})(\text{id} \otimes \rho(x))] = \text{Tr}[\text{Tr}_1[O_\mu(M \otimes \text{id})]\rho(x)] \quad (23)$$

where $O_\mu = \int \rho(y) \otimes \rho(y) \mu(dy)$. Note that this reformulation makes the isomorphism of $\mathcal{L}(\mathcal{H}, \mathcal{H}) \otimes \mathcal{L}(\mathcal{H}, \mathcal{H}) \simeq \mathcal{L}(\mathcal{H} \otimes \mathcal{H}, \mathcal{H} \otimes \mathcal{H}) \simeq \mathcal{L}(\mathcal{L}(\mathcal{H}, \mathcal{H}), \mathcal{L}(\mathcal{H}, \mathcal{H}))$ explicit. The spectrum of K thus agrees with the eigenvalues of the linear map T acting on matrices by

$$T(M) = \text{Tr}_2[O_\mu(\text{id} \otimes M)]. \quad (24)$$

We claim that there is an eigendecomposition

$$T(M) = \sum \gamma_i A_i \text{Tr}[A_i M] \quad (25)$$

where A_i are orthonormal hermitian matrices. Moreover, the eigenfunctions of K are $f_i(x) = \text{Tr}[\rho(x)A_i]$. This result follows from standard results in linear algebra, we give all details in the next subsection.

C.4 Spectral decomposition of linear maps preserving hermitian matrices

We consider the space of matrices $\mathbb{C}^{n \times n}$ equipped with the usual scalar product $\langle A, B \rangle = \text{Tr}[A^\dagger B]$ which agrees with the standard scalar product on \mathbb{C}^{n^2} after vectorisation. We will need them following fact: For hermitian matrices A, B the scalar product $\langle A, B \rangle \in \mathbb{R}$ is real.

Lemma 3. *Let $T : \mathbb{C}^{n \times n} \rightarrow \mathbb{C}^{n \times n}$ be a linear and hermitian map that maps hermitian matrices to hermitian matrices. Then there is a eigendecomposition (γ_i, H_i) with real eigenvalues γ_i and orthonormal hermitian matrices H_i such that*

$$T(A) = \sum_i \gamma_i H_i \text{Tr}[H_i^\dagger A]. \quad (26)$$

Proof. Hermitian matrices can be diagonalized with real values γ_i so we can write

$$T(A) = \sum_i \gamma_i X_i \text{Tr}[X_i^\dagger A] \quad (27)$$

where X_i form an orthonormal eigenbasis. It remains to show that we can find such a decomposition where the X_i are hermitian. We decompose $X_i = \tilde{H}_i + i\tilde{S}_i$ where \tilde{H}_i and \tilde{S}_i are hermitian. Then we observe

$$\gamma_i(\tilde{H}_i + i\tilde{S}_i) = \gamma_i X_i = T(X_i) = T(\tilde{H}_i) + iT(\tilde{S}_i). \quad (28)$$

Using the invariances of T on hermitian matrices we conclude that \tilde{S}_i and \tilde{H}_i are again eigenvectors with eigenvalue γ_i . Now we can iteratively replace X_i by either \tilde{S}_i or \tilde{H}_i so that the set of vectors remains a basis. Finally we orthonormalize the resulting basis of all eigenspaces using the Gram-Schmidt procedure. Since scalar products of hermitian matrices are real we obtain an orthonormal eigenbasis H_i consisting of hermitian matrices. \square

We now apply this to the integral operator for the quantum embedding. Recall that the linear map T acting on matrices was defined by

$$T(M) = \text{Tr}_2 [O_\mu(\text{id} \otimes M)]. \quad (29)$$

Clearly, T is linear. To show that T is hermitian we observe that

$$\begin{aligned} \langle M, T(M) \rangle &= \text{Tr} [M^\dagger \text{Tr}_2 [O_\mu(\text{id} \otimes M)]] = \int \text{Tr} [M^\dagger \text{Tr}_2 [\rho(y) \otimes \rho(y)(\text{id} \otimes M)]] \mu(dy) \\ &= \int \text{Tr} [M^\dagger \rho(y)] \text{Tr} [\rho(y)M] \mu(dy) \in \mathbb{R}. \end{aligned} \quad (30)$$

Similarly we see that T preserves hermitian matrices, indeed, if $M = M^\dagger$

$$T(M) = \int \text{Tr}_2 [\rho(y) \otimes \rho(y)(\text{id} \otimes M)] \mu(dy) = \int \rho(y) \text{Tr} [\rho(y)M] \mu(dy). \quad (31)$$

which is hermitian because $\rho(y)$ is hermitian and the scalar product of hermitian matrices is real. Using Lemma 3 above we conclude that we can write $T(M) = \sum_i \gamma_i A_i \text{Tr}[A_i M]$ where γ_i are the eigenvalues of T which agree with the eigenvalues of the corresponding integral operator and the eigenfunctions are given by $x \rightarrow \text{Tr}[\rho(x)A_i]$.

C.5 A complete example

To illustrate the analysis above we consider the setting from Example 2 where $x \rightarrow \cos(x/2)|0\rangle + i \sin(x/2)|1\rangle$. Then $\tilde{\mathcal{F}} = \langle \sin(x), \cos(x) \rangle$ and $\mathcal{F} = \langle \sin^2(x), \cos^2(x), \sin(x) \cos(x) \rangle$. Note that the RKHS has dimension 4 when the relative phase between $a(x)$ and $b(x)$ is not constant (then ab and $a\bar{b}$ are not linearly dependent). The feature map of the physical kernel for our example is

$$\rho(y) = \begin{pmatrix} \cos^2(\frac{y}{2}) & -i \cos(\frac{y}{2}) \sin(\frac{y}{2}) \\ i \cos(\frac{y}{2}) \sin(\frac{y}{2}) & \sin^2(\frac{y}{2}) \end{pmatrix}. \quad (32)$$

For the analysis of the integral operator we need the matrix elements of the linear map T . We observe that in index notation using the Einstein summation convention and denoting complex conjugation without transposition by $*$

$$T(M)_{ij} = \int \rho(y)_{ij} \rho(y)_{kl} M_{lk} \mu(dy) = \int \rho(y)_{ij} \rho^*(y)_{lk} M_{lk} \mu(dy). \quad (33)$$

Using vectorisation we obtain

$$\text{Vec}(T(M)) = \int \text{Vec}(\rho(y)) \text{Vec}(\rho(y))^\top \mu(dy) \text{Vec}(M) = A_\mu \text{Vec}M. \quad (34)$$

In our example we obtain

$$\begin{aligned}
A_\mu &= \frac{1}{\pi} \int_0^\pi \begin{pmatrix} \cos^2(\frac{y}{2}) \\ -i \cos(\frac{y}{2}) \sin(\frac{y}{2}) \\ i \cos(\frac{y}{2}) \sin(\frac{y}{2}) \\ \sin^2(\frac{y}{2}) \end{pmatrix} (\cos^2(\frac{y}{2}) \quad i \cos(\frac{y}{2}) \sin(\frac{y}{2}) \quad -i \cos(\frac{y}{2}) \sin(\frac{y}{2}) \quad \sin^2(\frac{y}{2})) dy \\
&= \frac{1}{\pi} \int_0^\pi \begin{pmatrix} \cos^4(\frac{y}{2}) & i \cos^3(\frac{y}{2}) \sin(\frac{y}{2}) & -i \cos^3(\frac{y}{2}) \sin(\frac{y}{2}) & \cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) \\ -i \cos^3(\frac{y}{2}) \sin(\frac{y}{2}) & \cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) & -\cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) & -i \cos(\frac{y}{2}) \sin^3(\frac{y}{2}) \\ i \cos^3(\frac{y}{2}) \sin(\frac{y}{2}) & -\cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) & \cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) & i \cos(\frac{y}{2}) \sin^3(\frac{y}{2}) \\ \cos^2(\frac{y}{2}) \sin^2(\frac{y}{2}) & i \cos(\frac{y}{2}) \sin^3(\frac{y}{2}) & -i \cos(\frac{y}{2}) \sin^3(\frac{y}{2}) & \sin^4(\frac{y}{2}) \end{pmatrix} dy \\
&= \frac{1}{8} \begin{pmatrix} 3 & 0 & 0 & 1 \\ 0 & 1 & -1 & 0 \\ 0 & -1 & 1 & 0 \\ 1 & 0 & 0 & 3 \end{pmatrix}
\end{aligned} \tag{35}$$

We obtain the eigenvalues $\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0$ the eigenvectors are, in matrix notation,

$$H_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad H_2 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \quad H_3 = \begin{pmatrix} 0 & i \\ -i & 0 \end{pmatrix}, \quad H_4 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}. \tag{36}$$

The corresponding eigenfunctions f_i of the integral operator are given by $x \rightarrow \text{Tr}[\rho(x)H_i]$, i.e.,

$$\begin{aligned}
f_1(x) &= 1, & f_2(x) &= \cos^2(\frac{x}{2}) - \sin^2(\frac{x}{2}) = \cos(x), \\
f_3(x) &= 2 \cos(\frac{x}{2}) \sin(\frac{x}{2}) = \sin(x), & f_4(x) &= 0.
\end{aligned} \tag{37}$$

We can also parametrize the functions in the RKHS by $a \cos(x + b) + c$ with $a, b, c \in \mathbb{R}$.

Let us also look at the generalization to the vector valued case with d -qubits. Then the RKHS is given by all functions of the form

$$x \rightarrow \prod_{i=1}^d (a_i \cos(x_i + b_i) + c_i). \tag{38}$$

The eigenfunctions of the integral operator are given by

$$\prod_{i=1}^d \sin^{\alpha_i}(x_i) \cos^{\beta_i}(x_i) \tag{39}$$

where α_i, β_i are non-negative integers satisfying $\alpha_i + \beta_i \leq 1$. The corresponding eigenvalue is $2^{-d-\sum(\alpha_i+\beta_i)}$. The degeneracy of the eigenvalue 2^{-d-l} can be calculated to $2^l \binom{d}{l}$.

D Proof of Theorem 1

In this section we prove Theorem 1 which will follow easily from the result below. We remark that the following theorem is by no means sharp but a detailed analysis when learning is not possible is of limited interest. Note that again typical lower bounds for the learning performance are focused on the case $n \rightarrow \infty$ [41].

Theorem 3. Consider a measure space (X, μ) such that $\mu(X) = 1$ with a kernel k satisfying $k(x, x) = 1$ for all $x \in X$. Denote by γ_{\max} the largest eigenvalue of the corresponding integral operator. Suppose we have n training points $\mathcal{D}_n = \{(x_i, y_i), 1 \leq i \leq n\}$ with $(x_i, y_i) \in X \times \mathbb{R}$ where x_i are i.i.d. draws from μ and $y_i = f(x_i)$ for some square integrable function f . Then, for any $\varepsilon > 0$ with probability at least $1 - \varepsilon - \gamma_{\max} n^4$

$$\|f - \hat{f}_n^\lambda\|_2 \geq \left(1 - \sqrt{\frac{2\gamma_{\max} n^2}{\varepsilon}}\right) \|f\|_2 \quad (40)$$

for all λ where \hat{f}_n^λ denotes the kernel ridge regression estimator for training data (x_i, y_i) .

Proof. Denote the eigenvalues of the integral operator by γ_i with $\gamma_1 = \gamma_{\max}$. Standard results for integral operators imply

$$\sum_i \gamma_i = \int k(x, x) \mu(dx) = 1 \quad (41)$$

$$\sum_i \gamma_i^2 = \int k(x, y)^2 \mu(dx) \mu(dy) = \|k\|_2^2. \quad (42)$$

We conclude that

$$\|k\|_2^2 = \sum_i \gamma_i^2 \leq \gamma_{\max} \sum_i \gamma_i = \gamma_{\max}. \quad (43)$$

This implies $\mathbb{P}_{\mu \otimes \mu}(|k(x, y)| \geq \varepsilon) \leq \frac{\gamma_{\max}}{\varepsilon^2}$. Let $A_n = \{|k(x_i, x_j)| \leq \frac{1}{2n} \text{ for all } i \neq j\}$. Using the union bound we conclude that

$$\mathbb{P}_{\mathcal{D}_n}(A_n) \geq 1 - n^2 \mathbb{P}_{\mu \otimes \mu}\left(|k(x, y)| \geq \frac{1}{2n}\right) \geq 1 - 4n^4 \gamma_{\max}. \quad (44)$$

Conditioned on A_n we can bound the eigenvalues of the kernel matrix $K(X, X)_{i,j} = k(x_i, x_j)$ using Gershgorin circles by $1 - n \frac{1}{2n} = \frac{1}{2}$ and thus

$$K(X, X)^{-1} \leq 2 \cdot \text{id}_n. \quad (45)$$

Let us denote the Mercer decomposition of k by

$$k(x, y) = \sum_i \gamma_i f_i(x) f_i(y) \quad (46)$$

where f_i are the orthonormal eigenfunctions. Then we can bound

$$\begin{aligned} |k(x, \cdot)|_2^2 &= \int k(x, y)^2 \mu(dy) = \int \sum_i \gamma_i f_i(x) f_i(y) \sum_j \gamma_j f_j(x) f_j(y) \mu(dy) \\ &= \sum_{i,j} \delta_{ij} \gamma_i \gamma_j f_i(x) f_j(x) \leq \gamma_{\max} \sum_i \gamma_i f_i(x)^2 = \gamma_{\max}. \end{aligned} \quad (47)$$

The kernel ridge regression function \hat{f}_n^λ can be written as

$$\hat{f}_n^\lambda = \sum_i \alpha_i k(x_i, \cdot) \quad (48)$$

where the vector α is given by $\alpha = (K(X, X) + \lambda)^{-1}y$. Using (45) we conclude that conditioned on A_n we have $\|\alpha\|^2 \leq 2\|y\|^2$. Moreover, for any $\varepsilon > 0$ with probability $1 - \varepsilon$ we have $\|y\|^2 \leq \frac{n}{\varepsilon}\|f\|_2^2$. The L^2 norm of f_n^λ satisfies now with probability $1 - \varepsilon - \gamma_{max}n^4$ the bound

$$\begin{aligned} \|f_n^\lambda\|_2 &\leq \sum_i |\alpha_i| \|k(x_i, \cdot)\|_2 \leq \sqrt{\gamma_{max}} \sum_i |\alpha_i| \leq \sqrt{\gamma_{max}n} \sqrt{\sum_i \alpha_i^2} \\ &\leq \sqrt{\gamma_{max}n} \sqrt{\frac{2n\|f\|_2^2}{\varepsilon}} \leq \sqrt{\frac{2\gamma_{max}n^2}{\varepsilon}} \|f\|. \end{aligned} \quad (49)$$

We conclude that with probability $1 - \varepsilon - \gamma_{max}n^4$

$$\|f - f_n^\lambda\|_2 \geq \|f\|_2 - \|f_n^\lambda\|_2 \geq \|f\|_2 \left(1 - \sqrt{\frac{2\gamma_{max}n^2}{\varepsilon}}\right). \quad (50)$$

□

The proof of Theorem 1 is now a simple consequence of the result above.

Proof of Theorem 1. We first note that, using the assumption $\mu = \bigotimes \mu_i$

$$\rho_\mu = \bigotimes \rho_{\mu_i} \quad (51)$$

and thus

$$\text{Tr} [\rho_\mu] = \prod \text{Tr} [\rho_{\mu_i}] \leq \delta^d. \quad (52)$$

Lemma 1 implies that the largest eigenvalue of the integral operator is bounded by $\gamma_{max} \leq \delta^{d/2}$. The claim now follows from Theorem 3 as soon as $\delta^{d/2} < \varepsilon n^{-4}/2 \leq \varepsilon d^{-4k}/2$ (ensuring $1 - \varepsilon/2 - \gamma_{max}n^4 \geq 1 - \varepsilon$) and $\delta^{d/2} < \varepsilon^2 n^{-2}/4 \leq \varepsilon^2 d^{-2k}/4$ (ensuring $\sqrt{2\gamma_{max}n^2(\varepsilon/2)^{-1}} \leq \varepsilon$). □

E Proof of Theorem 2

We introduce some theory and notation necessary for the proof. We investigate the behavior of reduced density matrices when V is distributed according to the Haar-measure on the group of unitary matrices. The first even moments of the Haar measure on $U(2^d)$ are given by (see e.g., [54])

$$\begin{aligned} \int V_{ij} V_{i'j'}^* \mu(dV) &= \frac{\delta_{ii'} \delta_{jj'}}{2^d} \\ \int V_{i_1 j_1} V_{i_2 j_2} V_{i'_1 j'_1}^* V_{i'_2 j'_2}^* \mu(dV) &= \frac{1}{2^{2d} - 1} \left(\delta_{i_1 i'_1} \delta_{j_1 j'_1} \delta_{i_2 i'_2} \delta_{j_2 j'_2} + \delta_{i_1 i'_2} \delta_{j_1 j'_2} \delta_{i_2 i'_1} \delta_{j_2 j'_1} \right. \\ &\quad \left. - \frac{1}{2^d} (\delta_{i_1 i'_1} \delta_{j_1 j'_2} \delta_{i_2 i'_2} \delta_{j_2 j'_1} + \delta_{i_1 i'_2} \delta_{j_1 j'_1} \delta_{i_2 i'_1} \delta_{j_2 j'_2}) \right). \end{aligned} \quad (53)$$

Let us remark that while random circuits that output Haar-distributed unitaries require an exponential (in d) number of gates our arguments actually only require unitary t -designs which are point distributions that match the first t moments of the Haar measure. In particular a 2-design is a measure with finite support on unitary matrices satisfying (53) (and odd moments

of lower order vanish). Those can be implemented using polynomially many gates. For details and further information we refer to the literature [55].

Recall the definition of the projected quantum kernel

$$\tilde{\rho}_m^V(x) = \text{Tr}_{m+1\dots d} [\rho^V(x)]. \quad (54)$$

To denote the partial trace in index notation we split the index $i \in \{1, \dots, 2^d\}$ in $(\alpha, \tilde{\alpha})$ where $\alpha \in \{1, \dots, 2^m\}$ denotes the index corresponding to the first m qubits and $\tilde{\alpha} \in \{1, \dots, 2^{d-m}\}$ denotes the index corresponding to the remaining $d - m$ qubits. In particular, summing 1 over $\tilde{\alpha}$ results in 2^{d-m} . We are now ready to prove Theorem 2.

Proof of Theorem 2. We start to prove the asymptotic expression for the reduced density matrix which is a standard result. We can write using Einstein summation convention

$$\mathbb{E}_V [\tilde{\rho}_m^V(x)_{\alpha\beta}] = \mathbb{E}_V [V_{\alpha\tilde{\alpha},j} \rho(x)_{jj'} V_{\beta\tilde{\alpha},j'}^*] = \frac{2^{d-m}}{2^d} \delta_{\alpha\beta} \delta_{jj'} \rho(x)_{jj'} = 2^{-m} \delta_{\alpha\beta} \text{Tr} [\rho(x)]. \quad (55)$$

To show the concentration around the expectation value we need to calculate the variance of this expression. We calculate the second moment of the reduced density matrix

$$\begin{aligned} \mathbb{E}_V [\tilde{\rho}_m^V(x)_{\alpha\beta} (\tilde{\rho}_m^V(y)_{\gamma\delta})^*] &= \mathbb{E}_V \left[V_{\alpha\tilde{\alpha},j_1} \rho(x)_{j_1 j'_1} V_{\beta\tilde{\alpha},j'_1}^* V_{\gamma\tilde{\gamma},j'_2}^* (\rho(y)_{j'_2 j_2})^* V_{\delta\tilde{\gamma},j_2} \right] \\ &= A_1 + A_2 + A_3 + A_4. \end{aligned} \quad (56)$$

The four terms can be evaluated to (assuming that $\text{Tr} [\rho(x)] = 1$ for all x)

$$\begin{aligned} A_1 &= \frac{1}{2^{2d}-1} \delta_{\alpha\beta} 2^{d-m} \text{Tr} [\rho(x)] \delta_{\gamma\delta} 2^{d-m} \text{Tr} [\rho(y)] = \frac{2^{2d}}{2^{2d}-1} 2^{-2m} \delta_{\alpha\beta} \delta_{\gamma\delta} \\ &= 2^{-2m} \delta_{\alpha\beta} \delta_{\gamma\delta} + \frac{1}{2^{2d}(2^{2d}-1)} 2^{-2m} \delta_{\alpha\beta} \delta_{\gamma\delta} \\ A_2 &= \frac{1}{2^{2d}-1} \delta_{\alpha\gamma} \delta_{\tilde{\alpha}\tilde{\gamma}} \delta_{j_1 j'_2} \delta_{\beta\delta} \delta_{\tilde{\alpha}\tilde{\gamma}} \delta_{j_2 j'_1} \rho(x)_{j_1 j'_1} \rho^*(y)_{j'_2 j_2} \\ &= \frac{1}{2^{2d}-1} \delta_{\tilde{\alpha}\tilde{\alpha}} \rho(x)_{j_1 j'_1} \rho^*(y)_{j_1 j'_1} \delta_{\alpha\gamma} \delta_{\beta\delta} = \frac{2^{d-m}}{2^{2d}-1} \text{Tr} [\rho(x) \rho^\dagger(y)] \delta_{\alpha\gamma} \delta_{\beta\delta} \end{aligned} \quad (57)$$

$$\begin{aligned} A_3 &= -\frac{1}{2^d(2^{2d}-1)} \delta_{\alpha\beta} \delta_{\tilde{\alpha}\tilde{\alpha}} \delta_{j_1 j'_2} \delta_{\gamma\delta} \delta_{\tilde{\gamma}\tilde{\gamma}} \delta_{j_2 j'_1} \rho(x)_{j_1 j'_1} \rho^*(y)_{j'_2 j_2} \\ &= -\frac{2^{2d-2m}}{2^d(2^{2d}-1)} \rho(x)_{j_1 j'_1} \rho^*(y)_{j_1 j'_1} \delta_{\alpha\beta} \delta_{\gamma\delta} = -\frac{2^{d-2m}}{(2^{2d}-1)} \text{Tr} [\rho(x) \rho^\dagger(y)] \delta_{\alpha\beta} \delta_{\gamma\delta} \\ A_4 &= -\frac{1}{2^d(2^{2d}-1)} \delta_{\alpha\gamma} \delta_{\tilde{\alpha}\tilde{\gamma}} \delta_{j_1 j'_1} \delta_{\beta\delta} \delta_{\tilde{\alpha}\tilde{\gamma}} \delta_{j_2 j'_2} \rho(x)_{j_1 j'_1} \rho^*(y)_{j'_2 j_2} \\ &= -\frac{1}{2^d(2^{2d}-1)} \delta_{\tilde{\alpha}\tilde{\alpha}} \rho(x)_{j_1 j_1} \rho^*(y)_{j_2 j_2} \delta_{\alpha\gamma} \delta_{\beta\delta} = -\frac{2^{-m}}{2^{2d}-1} \delta_{\alpha\gamma} \delta_{\beta\delta} \end{aligned}$$

Collecting all terms we can bound the variance of the entries of $\tilde{\rho}(x)$ by

$$\begin{aligned} \mathbb{E}_V [\tilde{\rho}_m^V(x)_{\alpha\beta} (\tilde{\rho}_m^V(x)_{\alpha\beta})^*] - \mathbb{E}_V [\tilde{\rho}_m^V(x)_{\alpha\beta}] \mathbb{E}_V [(\tilde{\rho}_m^V(x)_{\alpha\beta})^*] \\ = 2^{-2m} \delta_{\alpha\beta} - (2^{-m})^2 \delta_{\alpha\beta} + O(2^{-d}) = O(2^{-d}). \end{aligned} \quad (58)$$

This shows that $\tilde{\rho}^V(x)$ is close to 2^{-m}id with high probability for large d .

We now turn to the evaluation of the averaged operator $\mathbb{E}_V [O_\mu]$ and the corresponding operator

$$T(M) = \text{Tr}_2 [\mathbb{E}_V [O_\mu](\text{id} \otimes M)] \quad (59)$$

whose matrix elements we denote by $T_{\alpha\beta,\gamma\delta}$. We assume that $\rho(x)$ is pure for all x , i.e., $\text{Tr} [\rho(x)^2] = 1$. We have seen in (33) that the matrix elements of this operator are given by

$$\mathbb{E}_V \left[\int \tilde{\rho}_m^V(y) \otimes (\tilde{\rho}_m^V(y))^* \mu(dy) \right] = \int \mathbb{E}_V [\tilde{\rho}_m^V(y) \otimes (\tilde{\rho}_m^V(y))^*] \mu(dy). \quad (60)$$

From (56) and (57) we obtain for the matrix elements

$$\mathbb{E}_V [\tilde{\rho}^V(y)_{\alpha\beta} (\tilde{\rho}^V(y)_{\gamma\delta})^*] = 2^{-2m} \delta_{\alpha\beta} \delta_{\gamma\delta} + \frac{2^{-m}}{2^d} \delta_{\alpha\gamma} \delta_{\beta\delta} - \frac{2^{-2m}}{2^d} \delta_{\alpha\beta} \delta_{\gamma\delta} + O(2^{-2d}). \quad (61)$$

Since this is independent of x we can write the matrix elements of T as

$$T_{\alpha\beta,\gamma\delta} = 2^{-2m} (1 - 2^{-d}) \delta_{\alpha\beta} \delta_{\gamma\delta} + \frac{2^{-m}}{2^d} \delta_{\alpha\gamma} \delta_{\beta\delta} + O(2^{-2d}) \quad (62)$$

From here we conclude that T is the sum of a multiple of the identity and a rank one perturbation (plus higher order terms):

$$T(M) = \frac{2^{-m}}{2^d} M + 2^{-2m} (1 - 2^{-d}) \text{id}_{2^m \times 2^m} \text{Tr} [\text{id}_{2^m \times 2^m} M] + O(2^{-2d}). \quad (63)$$

In particular the eigenvalues neglecting the perturbation are

$$\gamma_1 = 2^{-2m} (1 - 2^{-d}) \text{Tr} [\text{id}_{2^m \times 2^m} \text{id}_{2^m \times 2^m}] + 2^{-m-d} = 2^{-m} (1 - 2^{-d}) + 2^{-m-d} = 2^{-m} \quad (64)$$

with eigenvector $M_1 = \text{id}_{2^m \times 2^m}$ and $\gamma_2 = \dots = \gamma_{2^m \times 2^m} = 2^{-m-d}$ with traceless eigenvectors, i.e., $\text{Tr} [\text{id}_{2^m \times 2^m} M_i] = 0$ for $i \neq 1$. Standard bounds show that the higher order terms change the eigenvalues only by a term of order $O(2^{-2d})$. Finally, we observe that the function mapping $x \rightarrow \text{Tr} [\tilde{\rho}_m^V(x) M_1]$ is a constant function for any V . Indeed,

$$\text{Tr} [\tilde{\rho}_m^V(x) M_1] = \text{Tr} [\text{Tr}_{m+1\dots d} [V \rho(x) V^\dagger]] = \text{Tr} [V \rho(x) V^\dagger] = \text{Tr} [\rho(x)] = 1. \quad (65)$$

□

F More on experiments

For details on the implementation we refer to the provided code. We emphasize that our experiments simulate the full quantum state and thus work with the true values of the quantum

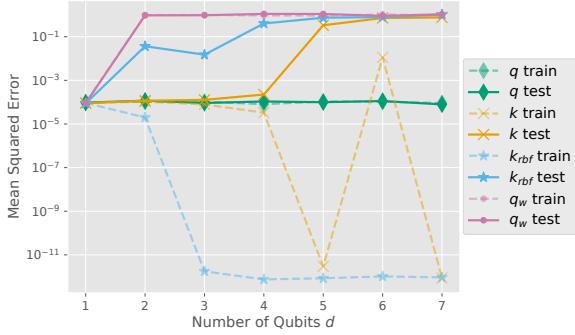


Figure 4: Similar as in Fig. 2. However, for the full quantum kernel k and the rbf kernel, we compute train and test loss over multiple choices of the regularization parameter. For each number of qubits, we only report the loss of the method that achieved smallest test loss. Note that, although this is invalid to asses the power of the full and rbf kernel, it shows, that the poor performance is not due to the choice of regularization. Since we cherry-pick on the test loss, it can happen that an underfitting regularization has the best test loss, which explains the outlier on k at $d = 6$.

kernel. This is an idealized setting and neglects the effect of finite measurements. Please see our discussion on Barren Plateaus in the main paper.

To reduce computations and speed-up the simulation, we compute the full quantum kernel $k(x, x') = \prod \cos^2((x_i - x'_i)/2)$ directly without simulating a quantum circuit. For the biased kernels we recommend (and implement it that way) to completely simulate $\rho_1^V(x_i)$ for all $i = 1, \dots, n$ and store the reduced density matrices (2×2 hermitian matrices). On a real quantum device this would correspond to doing quantum state tomography [36]. The benefit of this is that we only need to simulate the quantum circuit n times and can then directly compute the biased kernels via matrix products and tracing. If we chose to compute each entry of the kernel matrix individually we would have to simulate the circuit n^2 times.

Random generation of V . In order to generate random unitary matrices V we use the PennyLane function RANDOMLAYERS [52]. For d qubits we use d^2 layers of single qubit rotations and 2-qubit entangling gates. For more details and the used seeds please refer to the provided implementation.

Choice of regularization. For the biased kernels q, q_w regularization does not matter much, since they have only a four-dimensional RKHS and we consider sample sizes much larger than that. The RKHS simply does not have enough capacity to overfit to random noise. We therefore set the regularization $\lambda = 0$ for the biased kernels. On the other hand for the higher dimensional kernels k, k_{rbf} , the regularization strongly influences their performance. For the experiment in the main paper we set $\lambda = 10^{-3}$ for the latter methods. Note that in a real application one should use cross-validation or other model selection techniques to find good hyperparameters, which we omitted for simplicity. To exclude that the bad performance of k and k_{rbf} stems from a bad choice of regularization, we include experiments where we fit kernel ridge regression for 15 values of λ on a logarithmic grid from 10^{-6} to 10^4 . We then cherry-pick only the solution that performs best and report it in Figure 4. Note that such an approach is of course not legit to asses the actual

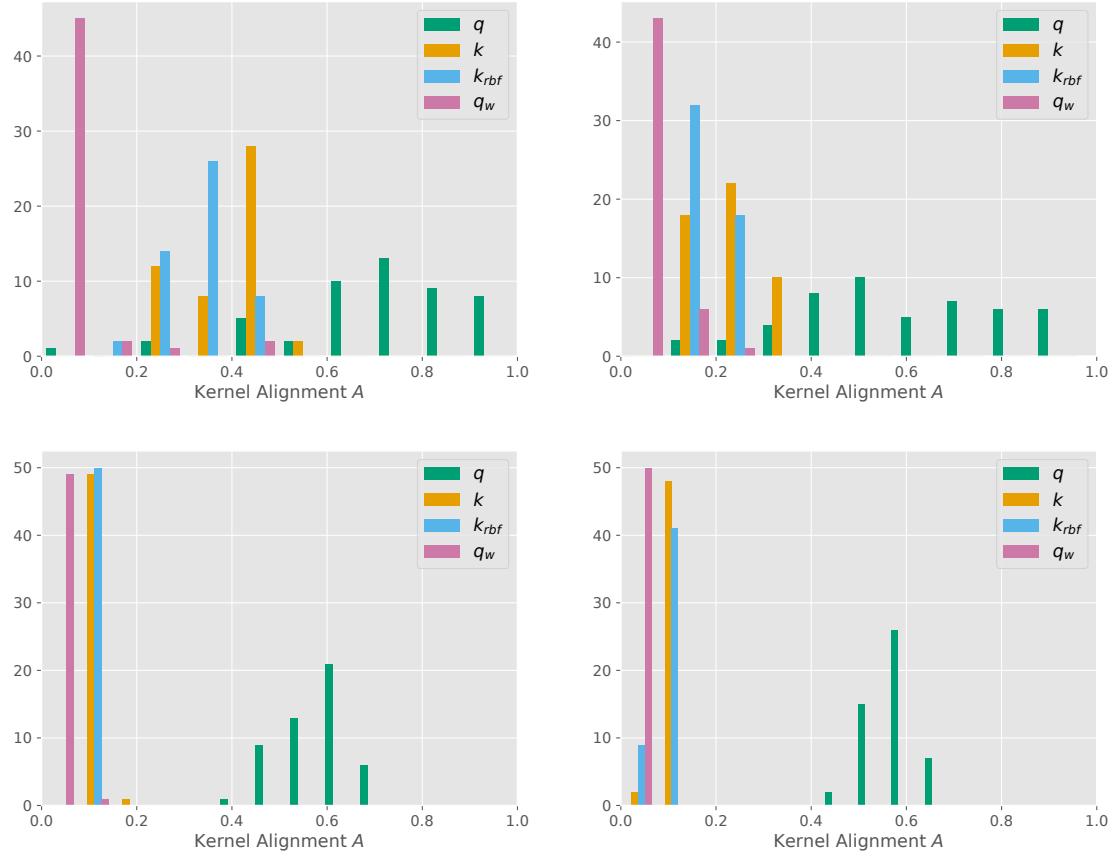


Figure 5: Kernel Target Alignment for $d = 1, 3, 5, 7$.

performance. However, it serves to bound the performance for the optimal choice of regularization. Our observations show that the behavior does not significantly change and we conclude that the performance difference indeed comes from the spectral bias as predicted by our theory.

Additional experiments. To show how the kernel target alignment changes as we increase the number of qubits d , we include further histograms in Figure 5. The estimated kernel alignment correlates with the learning performance reported in Figure 2.