

# FAST SUMMATION OF RADIAL KERNELS VIA QMC SLICING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The fast computation of large kernel sums is a challenging task, which arises as a subproblem in any kernel method. We approach the problem by slicing, which relies on random projections to one-dimensional subspaces and fast Fourier summation. We prove bounds for the slicing error and propose a quasi-Monte Carlo (QMC) approach for selecting the projections based on spherical quadrature rules. Numerical examples demonstrate that our QMC-slicing approach significantly outperforms existing methods like (QMC-)random Fourier features, orthogonal Fourier features or non-QMC slicing on standard test datasets.

## 1 INTRODUCTION

We consider fast algorithms for computing the kernel sums

$$s_m = \sum_{n=1}^N w_n K(x_n, y_m), \quad \text{for all } m = 1, \dots, M, \quad (1)$$

where  $x_n, y_m \in \mathbb{R}^d$  and  $w_n \in \mathbb{R}$  for  $n = 1, \dots, N$ ,  $m = 1, \dots, M$  and  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a radial kernel, i.e.,  $K(x, y) = F(\|x - y\|)$  for the Euclidean norm  $\|\cdot\|$  and some  $F: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ . This summation problem appears in most kernel methods including kernel density estimation (Parzen, 1962; Rosenblatt, 1956), classification via support vector machines (Steinwart & Christmann, 2008), dimensionality reduction with kernelized principal component analysis (Schölkopf & Smola, 2002; Shawe-Taylor & Cristianini, 2004), distance measures on the space of probability measures like maximum mean discrepancies or the energy distance (Gretton et al., 2006; Székely, 2002), corresponding gradient flows (Arbel et al., 2019; Galashov et al., 2024; Hagemann et al., 2024; Hertrich et al., 2024; Kolouri et al., 2022), and methods for Bayesian inference like Stein variational gradient descent (Liu & Wang, 2016). Computing (1) exactly for all  $m = 1, \dots, M$  has complexity  $\mathcal{O}(MN)$ , which can be restricting if  $M$  and  $N$  are large.

In low dimensions, there is a rich literature on fast approximation algorithms, we include a (non-exhaustive) list in the “related work” section. One particular approach is the fast Fourier summation (Kunis et al., 2006; Potts et al., 2004), which approximates the kernel by a truncated Fourier series and reformulates (1) using the fast Fourier transform on non-equispaced data (Beylkin, 1995; Dutt & Rokhlin, 1993). We provide a short overview in Appendix I. This kind of methods usually provides a computational complexity of  $\mathcal{O}(M + N + N_{\text{ft}} \log(N_{\text{ft}}))$ , where  $N_{\text{ft}}$  is the number of relevant Fourier coefficients, and admits very fast error rates for  $N_{\text{ft}} \rightarrow \infty$  (even exponential if the kernel is sufficiently smooth). However, the number  $N_{\text{ft}}$  of relevant Fourier coefficients grows exponentially with the dimension  $d$ , such that this method is computationally intractable for dimensions larger than four.

As a remedy for higher dimensions, Rahimi & Recht (2007) proposed random Fourier features (RFF). They represent a positive definite kernel via Bochner’s theorem (Bochner, 1933) as the Fourier transform of a non-negative measure. Sampling randomly from this measure at  $D \in \mathbb{N}$  points (features) leads to an approximation algorithm with computational complexity  $\mathcal{O}(D(N+M))$  independent of the dimension  $d$ . However, the error decays only with rate  $\mathcal{O}(1/\sqrt{D})$ , which can be limiting if a high accuracy is required. Moreover, RFF are limited to positive definite kernels and do not apply to other kernels, like the negative distance kernel  $K(x, y) = -\|x - y\|$ , which has

054 applications, e.g., within the energy distance (Székely, 2002) that is used for defining a distance on  
 055 the space of probability measures.

056 A related approach is slicing (Hertrich, 2024), which represents the kernel sum (1) as an expectation  
 057 of one-dimensional kernel sums of the randomly projected data points with a different kernel. Dis-  
 058 cretizing the expectation by sampling at  $P$  random projections, the kernel sums (1) can be approx-  
 059 imated by  $P$  one-dimensional kernel sums, which can be computed efficiently, e.g., by fast Fourier  
 060 summation. Similarly as for RFF, this leads to a complexity of  $\mathcal{O}(P(N + M + N_{\text{ft}} \log N_{\text{ft}}))$ , where  
 061 the expected error can be bounded by  $\mathcal{O}(1/\sqrt{P})$ . For positive definite kernels, a close relation be-  
 062 tween RFF and slicing was established by Rux et al. (2024), see the short overview in Appendix G.  
 063 One advantage of slicing is the applicability to kernels that are not positive definite.

064 A way to improve the  $\mathcal{O}(1/\sqrt{P})$  rate is to replace the uniformly chosen directions with specific  
 065 sequences of points. This yields so-called quasi-Monte Carlo (QMC) algorithms on the sphere,  
 066 see (Brauchart et al., 2014). Note that there also exist QMC approaches for RFF (Avron et al.,  
 067 2016). However, they depend on the restrictive assumption that the measure from Bochner’s theorem  
 068 decouples over the dimension, which is true for the Gauss and  $L^1$ -Laplace<sup>1</sup> kernel, but false for most  
 069 other common kernels like the Laplace, Matérn or negative distance kernel.

070 **Contributions** Our contributions for fast kernel summation in  $\mathbb{R}^d$  via slicing are as follows:

- 071 - We derive bounds on the slicing error for various kernels in Theorem 1 including all positive  
 072 definite radial kernels. Furthermore, we exactly compute the variance for the negative  
 073 distance kernel, the thin plate spline, the Laplace kernel and the Gauss kernel.
- 074 - We exploit QMC sequences on the sphere in order to improve the error rate  $\mathcal{O}(1/\sqrt{P})$ . To  
 075 ensure the applicability of the QMC error bounds, we prove certain smoothness results for  
 076 the function which maps a direction  $\xi \in \mathbb{S}^{d-1}$  to the corresponding one-dimensional kernel  
 077 in Theorem 3. The improved error rates are outlined in Corollary 4.
- 078 - We conduct extensive numerical experiments on standard test datasets for several kernels  
 079 and different QMC sequences, and demonstrate that our QMC slicing approach with the  
 080 proposed distance QMC designs significantly outperforms the non-QMC slicing method as  
 081 well as (QMC-)RFF. While the advantage of QMC slicing is most significant in dimensions  
 082  $d \leq 100$ , it also performs better in higher dimensions.

083 **Outline** In Section 2, we first revisit the slicing approach in detail and present our improved error  
 084 bounds. Afterwards, in Section 3, we consider quadrature and QMC sequences on the sphere and  
 085 prove the applicability of the approaches for slicing. We present our numerical results in Section 4  
 086 and draw conclusions in Section 5. Additional proofs, plots and evaluations are contained in the  
 087 appendix. The code for the numerical examples is provided in the supplementary material.

## 088 RELATED WORK

089 **Low-Dimensional Kernel Summation** Fast summation algorithms have been extensively stud-  
 090 ied in the literature. They include fast kernel summations based on (non-)equispaced fast Fourier  
 091 transforms (Greengard et al., 2022; Kunis et al., 2006; Potts et al., 2004), fast multipole methods  
 092 (Beatson & Newsam, 1992; Greengard & Rokhlin, 1987; Lee & Gray, 2008; Yarvin & Rokhlin,  
 093 1999), tree-based methods (March et al., 2015a;b) or H- and mosaic-skeleton matrices (Hackbusch,  
 094 1999; Minden et al., 2017; Tyrtshnikov, 1996). For the Gauss kernel, the fast Gauss transform was  
 095 proposed by Greengard & Strain (1991) and improved by Yang et al. (2003; 2004). More general  
 096 fast kernel transforms were considered by Ryan et al. (2022).

097 **QMC and Quadrature on Spheres** QMC designs on spheres were studied by Brauchart et al.  
 098 (2014). Here, the quadrature points optimizing the worst-case error in certain Sobolev spaces are  
 099 given by spherical  $t$ -designs, which integrate all polynomials up to degree  $t$  on the sphere exactly  
 100 (Delsarte et al., 1977; Bannai & Bannai, 2009). The construction of spherical designs is highly  
 101 non-trivial and intractable in high dimensions. For  $\mathbb{S}^2$  and  $\mathbb{S}^3$ , several examples were computed  
 102

103 <sup>1</sup>In literature, there exist two versions of the Laplace kernel  $K(x, y) = \exp(-\alpha\|x - y\|_1)$  and  $K(x, y) =$   
 104  $\exp(-\alpha\|x - y\|)$ , which differ in the used norm. Since our analysis focuses on radial kernels, we will only  
 105 consider the second version in the rest of this paper.

numerically by Gräf & Potts (2011) and Womersley (2018). Gräf et al. (2012) related quadrature rules on the sphere with halftoning problems.

**Sliced Wasserstein Distance** The idea of slicing is also used in optimal transport. A sliced Wasserstein distance was proposed by Rabin et al. (2012). In contrast to the kernel summation problem, the sliced and non-sliced Wasserstein distance do not coincide and have different properties. QMC rules for the three-dimensional sliced Wasserstein distance were considered by Nguyen et al. (2024).

**Random Fourier Features** Random Fourier features (RFFs) were proposed by Rahimi & Recht (2007) and were further analyzed in several papers Bach (2017); Hashemi et al. (2023); Li et al. (2021). To improve the error rates, Avron et al. (2016) proposed a quasi-Monte Carlo approach for RFFs under the restrictive assumption that the measure from Bochner’s theorem decouples over the dimensions. This approach was refined by Huang et al. (2024) and Munkhoeva et al. (2018). Yu et al. (2016) proposed orthogonal random features. In a very recent preprint, Belhadji et al. (2024) derive an explicit quadrature rule in the Fourier space for efficient summations of the Gauss kernel.

## 2 SLICING OF RADIAL KERNELS

Let  $K: \mathbb{R}^d \times \mathbb{R}^d$  be a radial kernel of the form  $K(x, y) = F(\|x - y\|)$  for some basis function  $F: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ . Throughout this paper, we will assume that  $K$  has the form

$$K(x, y) = \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [k(\langle \xi, x \rangle, \langle \xi, y \rangle)]$$

for some one-dimensional radial kernel  $k: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  with basis function  $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$ , where we suppress the dependence of  $f$  and  $F$  onto the dimension  $d$ . By inserting the basis functions of the kernels, this corresponds to the slicing relation

$$F(\|x\|) = \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [f(|\langle \xi, x \rangle|)]. \quad (2)$$

A pair  $(F, f)$  of basis functions in  $L_{\text{loc}}^{\infty}(\mathbb{R}_{\geq 0})$  fulfills this relation if and only if  $F$  is the *generalized Riemann–Liouville fractional integral* transform given by

$$F(t) = \frac{2\Gamma(\frac{d}{2})}{\sqrt{\pi}\Gamma(\frac{d-1}{2})} \int_0^1 f(ts)(1-s^2)^{\frac{d-3}{2}} ds, \quad (3)$$

for  $2 \leq d \in \mathbb{N}$ , see (Hertrich, 2024, Prop 2 and Rubin, 2003). In order to find the one-dimensional basis function  $f$  for a given  $F$ , we have to invert the transform (3). This can be done explicitly if

- i)  $F$  is analytic on  $\mathbb{R}_{\geq 0}$ , i.e., there exists  $(a_n)_{n \in \mathbb{N}} \subset \mathbb{R}$  such that  $F(x) = \sum_{n=0}^{\infty} a_n x^n$  for all  $x \geq 0$ , or
- ii)  $F(\|\cdot\|): \mathbb{R}^d \rightarrow \mathbb{R}$  is continuous, bounded and positive definite, i.e., for all  $N \in \mathbb{N}$ , all pairwise distinct  $x_j \in \mathbb{R}^d$  and all  $a_j \in \mathbb{R}$  for  $j = 1, \dots, N$  it holds that  $\sum_{j,k=1}^N a_j a_k F(\|x_j - x_k\|) \geq 0$ ,

see (Hertrich, 2024, Thm 3 and Rux et al., 2024, Cor 4.11). We include a list of pairs  $(F, f)$  fulfilling (2) in Appendix A. In particular, it includes the basis functions of Gauss, Laplace and Matérn kernels. Moreover,  $f$  can be computed for other important choices that fulfill neither i) nor ii), e.g., the thin-plate spline and the generalized Riesz kernel.

### 2.1 FAST KERNEL SUMMATION VIA SLICING

In order to compute the kernel sums (1) efficiently, we approximate  $F(\|\cdot\|)$  by projections and the one-dimensional basis function  $f$ , i.e., we aim to find directions  $\xi_1, \dots, \xi_P \in \mathbb{S}^{d-1}$  such that

$$F(\|x\|) \approx \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \quad \text{for all } x \in \mathbb{R}^d. \quad (4)$$

Then, the kernel sums (1) read as

$$s_m = \sum_{n=1}^N w_n K(x_n, y_m) = \sum_{n=1}^N w_n F(\|x_n - y_m\|) \approx \frac{1}{P} \sum_{p=1}^P \sum_{n=1}^N w_n f(|\langle \xi_p, x_n - y_m \rangle|). \quad (5)$$

For computing the one-dimensional sums  $\sum_{n=1}^N w_n f(|\langle \xi_p, x_n - y_m \rangle|)$  for all  $m = 1, \dots, M$ , there exists algorithms with complexity  $\mathcal{O}(M + N)$  or  $\mathcal{O}((M + N) \log(M + N))$  in literature. These include fast summations based on non-equispaced Fourier transforms (Kunis et al., 2006; Potts et al., 2004), fast multipole methods (Greengard & Rokhlin, 1987) or sorting algorithms (Hertrich et al., 2024). In particular, we can approximate the vector  $s = (s_1, \dots, s_M)$  via (5) with a complexity of  $\mathcal{O}(P(M + N))$ .

## 2.2 ERROR BOUNDS FOR UNIFORMLY DISTRIBUTED SLICES

To bound the error of the slicing procedure from the previous subsection, we consider error estimates for the approximation in (4). To this end, we assume that the directions  $\xi_1, \dots, \xi_P$  are iid samples from the uniform distribution on the sphere. Then, we exactly compute the variance

$$\mathbb{V}_d[f](x) := \mathbb{E}_{\xi \sim \mathcal{U}_{S^{d-1}}} \left[ (f(|\langle \xi, x \rangle|) - F(\|x\|))^2 \right], \quad (6)$$

which bounds the mean squared error through the Bienaymé's identity as

$$\mathbb{E}_{\xi_1, \dots, \xi_P \sim \mathcal{U}_{S^{d-1}}} \left[ \left( \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) - F(\|x\|) \right)^2 \right] = \frac{\mathbb{V}_d[f](x)}{P}.$$

In particular, our results show relative error bounds of the negative distance kernel  $K(x, y) = -\|x - y\|$  and the thin plate spline (except around  $\|x\| = 1$ ), which do not depend on the dimension  $d$ . The proof is given in Appendix C.

**Theorem 1.** *Let  $F: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  and  $f: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  fulfill the slicing relation (2).*

i) *If  $F(\|\cdot\|)$  is continuous and positive definite on  $\mathbb{R}^d$ , then  $\mathbb{V}_d[f](x) \leq F(0)^2 - F(\|x\|)^2$ .*

ii) *For the generalized Riesz kernel  $F(\|x\|) = -\|x\|^r$  with  $r > 0$ , we have*

$$\mathbb{V}_d[f](x) = \left( \frac{\sqrt{\pi} \Gamma(r + \frac{1}{2}) \Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{r+1}{2})^2 \Gamma(\frac{d}{2}) \Gamma(r + \frac{d}{2})} - 1 \right) F(\|x\|)^2 < \frac{\sqrt{\pi} \Gamma(r + \frac{1}{2})}{\Gamma(\frac{r+1}{2})^2} F(\|x\|)^2. \quad (7)$$

iii) *For the thin plate spline  $F(\|x\|) = \|x\|^2 \log(\|x\|)$ , we have*

$$\mathbb{V}_d[f](x) = \left( 2 + \frac{c_1}{\log(\|x\|)} + \frac{c_2 + \mathcal{O}(d^{-1} \log d)}{\log(\|x\|)^2} \right) \left( 1 + \frac{2}{d} \right) F(\|x\|)^2,$$

with  $c_1 \approx 4.189$  and  $c_2 \approx 2.895$  given in (18).

iv) *For  $F(\|x\|) = \sum_{n=0}^{\infty} a_n \|x\|^n$  and  $d \geq 3$  odd, we have*

$$\mathbb{V}_d[f](x) = \sum_{n=0}^{\infty} \left( \sum_{k=0}^n \left( \prod_{i=1}^{(d-1)/2} \left( 1 + \frac{k(n-k)}{(2i-1)(n+2i-1)} \right) - 1 \right) a_k a_{n-k} \right) \|x\|^n.$$

In particular, for the Laplace kernel  $F(\|x\|) = \exp(-\alpha\|x\|)$ , we have

$$\mathbb{V}_3[f](x) = \frac{1}{4\alpha\|x\|} \left( 1 - (2(\alpha\|x\|)^2 + 2\alpha\|x\| + 1) F(\|x\|)^2 \right),$$

and for the Gauss kernel  $F(\|x\|) = \exp(-\|x\|^2/(2\sigma^2))$ , we have

$$\mathbb{V}_3[f](x) = \frac{\sigma^2}{2\|x\|^2} \left( 1 - \left( \frac{\|x\|^4}{2\sigma^4} + \frac{\|x\|^2}{\sigma^2} + 1 \right) F(\|x\|)^2 \right).$$

Some weaker error bounds for the Gauss and Riesz kernel were also shown in Hertrich (2024), see Appendix D.1 for a detailed comparison. In all cases, the variance is independent of the dimension  $d$ . The dependence on  $x$  differs between the kernels: For positive definite kernels, which are always bounded, we have an absolute error bound i) independent of  $\|x\|$ . For the Riesz kernel, we have a relative error bound in ii). For the thin plate spline, iii) behaves like a relative bound for  $\|x\| \rightarrow \infty$  and  $\|x\| \rightarrow 0$ , but as a constant around  $\|x\| = 1$ , which is a zero of  $F$ . For the Laplace and Gauss kernel, the dependence on  $\|x\|$  changes drastically between  $\|x\| \rightarrow \infty$  and  $\|x\| \rightarrow 0$ . In fact,  $\mathbb{V}_3[f](x)$  is monotonically increasing in  $\|x\|$  with global upper bound  $1/(4\alpha)$ , and converges quadratically in  $\|x\|$  to zero for  $\|x\| \rightarrow 0$ . For the case  $d > 3$ , we conjecture a similar behavior, see Appendix D.2 for the discussion.

### 3 QUASI-MONTE CARLO SLICING

For directions drawn independently from the uniform measure  $\mathcal{U}_{\mathbb{S}^{d-1}}$  on the sphere  $\mathbb{S}^{d-1}$ , our experiments from the numerical part suggest that the rate  $\mathcal{O}(1/\sqrt{P})$  from Theorem 1 is optimal. As a remedy, we employ quadrature and quasi-Monte Carlo designs on the sphere for improving these error rates. To this end, we first revisit the corresponding literature in Subsection 3.1. Afterwards, we apply these results for our slicing method in Subsection 3.2.

#### 3.1 QUASI-MONTE CARLO METHODS ON THE SPHERE

Let  $\xi^P = (\xi_1^P, \dots, \xi_P^P) \in (\mathbb{S}^{d-1})^P$  for  $P \in \mathbb{N}$ . In the following, we aim to construct  $\xi^P$  such that the worst case error in a certain Sobolev space is asymptotically optimal. The definition of the Sobolev space  $H^s(\mathbb{S}^{d-1})$  is given in Appendix B.

**Definition 2.** A sequence  $(\xi^P)_P$  with  $P \rightarrow \infty$  is called a sequence of QMC designs for  $H^s(\mathbb{S}^{d-1})$  if there exists some  $c(s, d) > 0$  independent of  $P$  such that the worst case error

$$\sup_{\substack{g \in H^s(\mathbb{S}^{d-1}) \\ \|g\|_{H^s(\mathbb{S}^{d-1})} \leq 1}} \left| \frac{1}{|\mathbb{S}^{d-1}|} \int_{\mathbb{S}^{d-1}} g(\xi) d\xi - \frac{1}{P} \sum_{p=1}^P g(\xi_p^P) \right| \leq \frac{c(s, d)}{P^{s/(d-1)}} \in \mathcal{O}(P^{-s/(d-1)}). \quad (8)$$

It was proven by Hesse (2006) that the error rate  $\mathcal{O}(P^{-s/(d-1)})$  is optimal, see also Brauchart et al. (2014). For  $s > \frac{d-1}{2}$ , the existence of sequences of QMC designs is ensured by so-called spherical designs. More precisely,  $\xi^P$  is called a spherical  $t$ -design if the quadrature at these points integrates all polynomials of degree  $t \in \mathbb{N}$  exactly, i.e., if it holds

$$\frac{1}{|\mathbb{S}^{d-1}|} \int_{\mathbb{S}^{d-1}} \psi(\xi) d\xi = \frac{1}{P} \sum_{p=1}^P \psi(\xi_p^P) \quad \text{for all polynomials } \psi \text{ of degree } \leq t.$$

It can be shown that for any  $t$  there exists a spherical  $t$ -design with  $P = \mathcal{O}(t^{d-1})$  points, see Bondarenko et al. (2013). By (Brauchart & Hesse, 2007, Cor 3.6), such a sequence of spherical  $t$ -designs is a sequence of QMC designs, see also (Brauchart et al., 2014, Sect 1) for a summary.

Unfortunately, the construction of spherical  $t$ -designs in arbitrary dimension is numerically intractable. Instead, many QMC methods rely on low-discrepancy point sets. It was shown in (Brauchart et al., 2014, Thm 14) that a sequence  $\xi^P$  that minimizes the sum of powers of Euclidean distances

$$\mathcal{E}(\xi^P) := - \sum_{p, q=1}^P \|\xi_p^P - \xi_q^P\|^{2s-d-1} \quad (9)$$

is a QMC design for  $H^s(\mathbb{S}^{d-1})$  for  $s \in (\frac{d-1}{2}, \frac{d+1}{2})$ . In the numerics, we will consider  $s = \frac{d}{2}$ , so that we get a QMC design for  $H^{d/2}(\mathbb{S}^{d-1})$ , which we call the distance QMC design. Note that, up to a constant, (9) coincides with the maximum mean discrepancy with the Riesz kernel  $K(x, y) = -\|x - y\|^{2s-d-1}$  between the probability measures  $\frac{1}{P} \sum_{i=1}^P \delta_{\xi_i^P}$  and  $\mathcal{U}_{\mathbb{S}^{d-1}}$ , which is also known as energy distance Székely (2002). Furthermore, one can easily transform a QMC sequence into an unbiased estimator in (4), see Appendix E.

#### 3.2 SMOOTHNESS OF ONE-DIMENSIONAL BASIS FUNCTIONS

In order to apply the above theorems for our approximation (2), we need to ensure that for any  $x \in \mathbb{R}^d$  the spherical function  $g_x : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$  given by

$$g_x(\xi) := f(|\langle \xi, x \rangle|) \quad (10)$$

is sufficiently smooth on  $\mathbb{S}^{d-1}$ . For some specific examples, this is verified in the next theorem, whose proof is given in Appendix F. For part ii), we explicitly compute the Sobolev norm of  $g_x$ . Note that, in the special case  $s = 0$ , this relates to the variance (6) by the formula  $\|g_x\|_{H^0(\mathbb{S}^{d-1})}^2 = |\mathbb{S}^{d-1}|(\nabla_d[f](x) + F(\|x\|)^2)$  for  $(f, F)$  fulfilling (2).

**Theorem 3.** Let  $x \in \mathbb{R}^d$  with  $x \neq 0$ . For the Gauss, Riesz and Matérn kernel, the following smoothness results hold true:

i) For  $F(t) = \exp(-\frac{t^2}{2\sigma^2})$ , we have  $g_x \in H^s(\mathbb{S}^{d-1})$  for all  $s \geq 0$ .

ii) For  $F(t) = t^r$  with  $t \geq 0$  and  $r > -1$ , we have  $g_x \in H^s(\mathbb{S}^{d-1})$  if and only if  $s < r + \frac{1}{2}$ .

iii) For  $F(t) = \frac{2^{1-\nu}}{\Gamma(\nu)} (\frac{\sqrt{2\nu}}{\beta} t)^\nu K_\nu(\frac{\sqrt{2\nu}}{\beta} t)$ ,  $t \geq 0$ , we have  $g_x \in H^s(\mathbb{S}^{d-1})$  if  $s < 2\nu + \frac{1}{2}$ .

Note that the theorem also includes the Laplace kernel, which is the Matérn kernel for  $\nu = \frac{1}{2}$ . Combining this theorem with the results from the previous subsection leads to improved error bounds for the Gauss and Matérn kernel in the following corollary. For the Riesz kernel, the last theorem can be seen as a negative result that the theory from the previous subsection is not applicable.

**Corollary 4.** Let  $d \in \mathbb{N}$  and  $s > \frac{d-1}{2}$ . Then there exists a constant  $c(s, d)$  and a sequence  $(\xi^P)_P$  with  $P \rightarrow \infty$  such that for the Gauss and Matérn kernel with basis functions  $F(t) = \exp(-\frac{t^2}{2\sigma^2})$  and  $F(t) = \frac{2^{1-\nu}}{\Gamma(\nu)} (\frac{\sqrt{2\nu}}{\beta} x)^\nu K_\nu(\frac{\sqrt{2\nu}}{\beta} t)$  with  $\nu > \frac{2s-1}{4}$ , respectively, it holds that

$$\sup_{x \in \mathbb{R}^d} \left| F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \right| \leq \frac{c(s, d)}{P^{\frac{s}{d-1}}}.$$

For  $s \in (\frac{d-1}{2}, \frac{d+1}{2})$ , such  $\xi^P$  are given by minimizers of (9).

In Appendix J, we derive the complexity of the slicing-based kernel summation with a QMC sequence and show in Proposition 5 that it is superior to the random Fourier feature approach for the Gauss kernel. The summation methods are described in Appendix I and Appendix G, respectively.

## 4 NUMERICAL EXAMPLES

In the following, we evaluate the kernel approximation with QMC slicing for several QMC sequences and compare our results with random Fourier feature-based (RFF, Rahimi & Recht, 2007) methods. We implement the comparison in Julia and Python, and provide the code in the supplementary material. In Subsection 4.1, we describe the used QMC sequences, RFF-methods and kernels. Afterwards, in Subsection 4.2, we numerically evaluate the approximation error in (4). Finally, we apply our approximation for fast kernel summations in Subsection 4.3. We include additional numerical examples in Appendix K.

### 4.1 QMC SEQUENCES AND KERNELS

**QMC Sequences** Beside standard slicing where the projections  $\xi^P$  are drawn iid from the uniform distribution on  $\mathbb{S}^{d-1}$ , we consider the following sequences. We would like to emphasize that for the first two of them it is not clear, whether they are QMC designs in the sense of Definition 2, even though they are sometimes called QMC sequences in the literature.

- **Sobol Sequence:** Two commonly used QMC sequences on  $[0, 1]^d$  are Sobol (Sobol', 1967) and Halton (Halton, 1960) sequences. They can be transformed to QMC sequences for the normal distribution by applying the inverse cumulative density function along each dimension of the sequence, which was used for deriving a QMC method for random Fourier features Avron et al. (2016). To obtain a potential QMC sequence on the sphere, Nguyen et al. (2024) proposed to project the QMC sequence for the multivariate normal distribution onto the sphere by the transformation  $\xi = \theta / \|\theta\|$ , see also Beltrán et al. (2023). It is not known whether this constitutes a QMC sequence in the sense of Definition 2. To generate the original Sobol sequence on  $[0, 1]^d$ , we use the implementation from SciPy (Virtanen et al., 2020) in Python and the `Sobol.jl` package in Julia. Our numerical experiments suggest that the Sobol sequence lead to slightly better results than the Halton sequence. Therefore, we omit the Halton sequence in our comparison.
- **Orthogonal:** Even though this is technically not a QMC sequence, we adapt the approach of orthogonal features (Yu et al., 2016) for slicing and generate directions  $\xi$  as follows: We

generate  $\lceil \frac{P}{d} \rceil$  orthogonal matrices from the uniform distribution on  $O(d)$  (this can be done by taking the Q-factor of the QR decomposition applied on a matrix with standard normally distributed entries). Together, these matrices have  $d \lceil \frac{P}{d} \rceil$  columns from which we choose  $\xi^P$  to be the first  $P$  of those.

- **Distance:** In Section 3.1, we considered the **distance QMC design**  $\xi^P$  for  $H^{\frac{d}{2}}(\mathbb{S}^{d-1})$ , which is a minimizer of  $\mathcal{E}(\xi^P) = -\sum_{p,q=1}^P \|\xi_p^P - \xi_q^P\|$ , see (9). In our application, we have the additional symmetry that  $f(|\langle x, \xi \rangle|) = f(|\langle x, -\xi \rangle|)$ . Therefore, we construct the **distance QMC designs**  $\xi^P$  by minimizing the functional  $\mathcal{E}_{\text{sym}}(\xi^P) := \mathcal{E}(\xi^P, -\xi^P)$ . We do this numerically with the Adam optimizer (Kingma & Ba, 2015) and the PyKeops package (Charlier et al., 2021), which takes from a couple of seconds (for  $d = 3$ ) up to one hour (for  $d = 200$  and  $P \approx 5000$ ) on an NVIDIA RTX 4090 GPU. In Appendix H, we show that if  $P \leq d$ , the orthogonal points from above minimize  $\mathcal{E}_{\text{sym}}$ , so this approach differs only if  $P > d$ .
- **Spherical Design:** For  $d = 3$ , several spherical  $t$ -designs on the  $\mathbb{S}^2$  were computed by Gräf & Potts (2011) up to  $t \leq 1000$  and  $P \leq 1002000$  and are available online<sup>2</sup>. Spherical  $t$ -designs for  $\mathbb{S}^3$  were computed by Womersley (2018) up to  $t \leq 31$  and  $P \leq 3642$ . Unfortunately, the computation in higher dimensions appears to be intractable such that we only use the spherical designs for  $d = 3$ .

**Compared Methods** We compare our results with the following methods:

- **Random Fourier Features** (RFF, Rahimi & Recht, 2007): see Appendix G for a description.
- **Orthogonal Random Features** (ORF, Yu et al., 2016): The directions of the RFF features are chosen in the same way as explained above for the orthogonal slicing.
- **QMC-Random Fourier Features** (Sobol RFF, Avron et al., 2016): For the Gauss kernel, we also compare with QMC random Fourier features, which are only applicable for kernels where the Fourier transform decouples as a product over the dimensions. For the kernels from Table 2, this is only true for the Gauss kernel. As a QMC sequence in  $[0, 1]^d$ , we choose the Sobol sequence.

**Kernels** We use the Gauss, Laplace, Matérn (with  $\nu = p + \frac{1}{2}$  for  $p \in \{1, 3\}$ ), the negative distance kernel (Riesz kernel with  $r = 1$ ) and the **thin plate spline kernel**, see Table 2 in the appendix for the pairs  $(f, F)$ . The parameters  $\sigma$ ,  $\alpha$  and  $\beta$  are chosen by the median rule (see, e.g., Garreau et al., 2017 for an overview and history). That is, we choose  $\sigma = \beta = \frac{1}{\alpha} = \gamma m$ , where  $m$  is the median of all considered input norms  $\|x\|$  of the basis functions  $F$  and  $\gamma$  is some scaling factor which we set to  $\gamma \in \{\frac{1}{2}, 1, 2\}$ .

## 4.2 NUMERICAL EVALUATION OF THE SLICING ERROR

We examine the approximation error in (4) numerically. To this end, we draw a sample  $x$  from  $\mathcal{N}(0, 0.1I)$  and evaluate the absolute error  $|F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|)|$ . We average this error over 50 realizations of  $\xi^P$  (whenever  $\xi^P$  is random) and 1000 samples of  $x$ . The average results for the scale factor  $\gamma = 1$  and dimension  $d \in \{3, 10, 50\}$  with the Gauss, Laplace and Matérn kernel are given in Figure 1. We observe that all methods despite the spherical designs for the Gauss kernel converge with rate  $\mathcal{O}(P^{-r})$  for some  $r > 0$ . To estimate the rate  $r$  numerically, we fit a regression line in the loglog plot. The resulting rates  $r$  are given in Table 1. Further plots and tables are given in Appendix K.1 considering the negative distance kernel, higher dimensions and smaller/larger kernel widths.

Overall, the **distance QMC designs** perform best in most examples, except when the (provably optimal) **spherical designs** are applicable, which are only computable in  $d \leq 4$  and outperform the **distance QMC designs** for smooth kernels as they reach machine precision already for  $P \approx 250$ . In accordance with Corollary 4, the **benefits of QMC slicing are better for smooth kernels and in low dimensions**. But also for  $d = 50$ , a significant advantage of QMC slicing is visible. In particular,

<sup>2</sup><https://www-user.tu-chemnitz.de/~potts/workgroup/graef/quadrature/index.php.en>

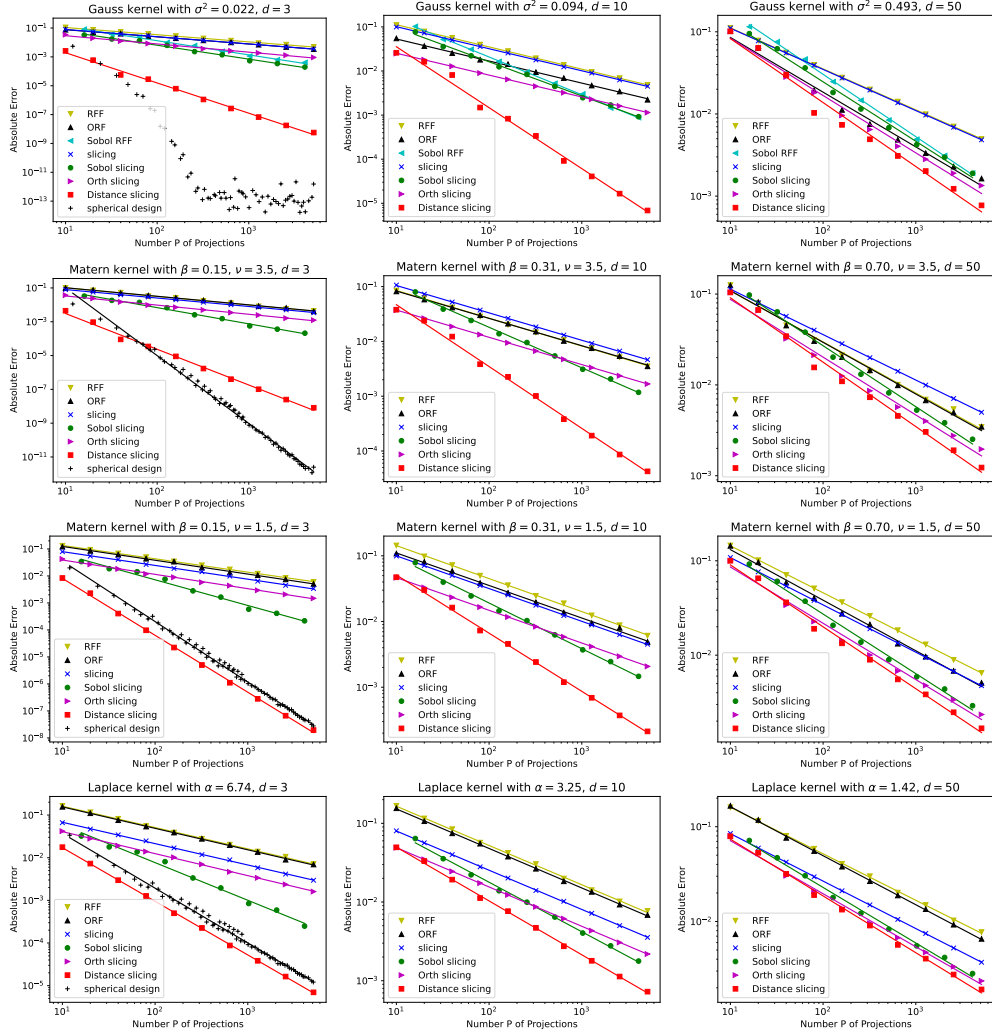


Figure 1: Loglog plot of the approximation error  $|F(\|x\|) - \frac{1}{P} \sum_p f(|\langle \xi_p, x \rangle|)|$  for approximating the function  $F$  by slicing (4) versus the number  $P$  of projections (or the number  $D = P$  of features for RFF and ORF) for different kernels and dimensions (left  $d = 3$ , middle  $d = 10$ , right  $d = 50$ ). The results are averaged over 50 realizations of  $\xi^P$  and 1000 realizations of  $x$ . The kernel parameters are set by the median rule with scaling factor  $\gamma = 1$ . We fit a regression line in the loglog plot for each method to estimate the convergence rate  $r$ , see also Table 1.

Table 1: Estimated convergence rates for the different methods. We estimate the rate  $r$  by fitting a regression line in the loglog plot. Then, we obtain the estimated convergence rate  $P^{-r}$  for some  $r > 0$ . Consequently, larger values of  $r$  correspond to a faster convergence. The resulting values of  $r$  are given in the below tables, the best values are highlighted in bold. The kernel parameters are the same as in Figure 1 (median rule with scaling factor  $\gamma = 1$ ).

Gauss kernel with median rule and scaling $\gamma = 1$							
Dimension	RFF-based			Slicing-based			
	RFF	Sobol	ORF	Slicing	Sobol	Orth	Distance
$d = 3$	0.50	0.98	0.50	0.50	0.96	0.57	<b>2.10</b>
$d = 10$	0.50	0.86	0.50	0.50	0.78	0.50	<b>1.38</b>
$d = 50$	0.50	0.76	0.67	0.50	0.72	0.70	<b>0.78</b>

Matérn kernel with $\nu = 3 + \frac{1}{2}$ and median rule with scaling $\gamma = 1$							
Dimension	RFF-based			Slicing-based			
	RFF	Sobol	ORF	Slicing	Sobol	Orth	Distance
$d = 3$	0.51	0.51	0.50	0.50	0.96	0.54	2.11
$d = 10$	0.51	0.50	0.50	0.50	0.74	0.50	<b>1.13</b>
$d = 50$	0.56	0.57	0.50	0.50	0.67	0.64	<b>0.71</b>

Matérn kernel with $\nu = 1 + \frac{1}{2}$ and median rule with scaling $\gamma = 1$							
Dimension	RFF-based			Slicing-based			
	RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design
$d = 3$	0.50	0.51	0.51	0.95	0.53	2.11	<b>2.24</b>
$d = 10$	0.50	0.50	0.50	0.70	0.50	<b>0.89</b>	-
$d = 50$	0.50	0.54	0.50	0.63	0.60	<b>0.66</b>	-

Laplace kernel with median rule and scaling $\gamma = 1$							
Dimension	RFF-based			Slicing-based			
	RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design
$d = 3$	0.50	0.50	0.50	0.88	0.52	1.26	<b>1.28</b>
$d = 10$	0.50	0.50	0.50	0.63	0.50	<b>0.68</b>	-
$d = 50$	0.49	0.52	0.50	0.59	0.56	<b>0.60</b>	-



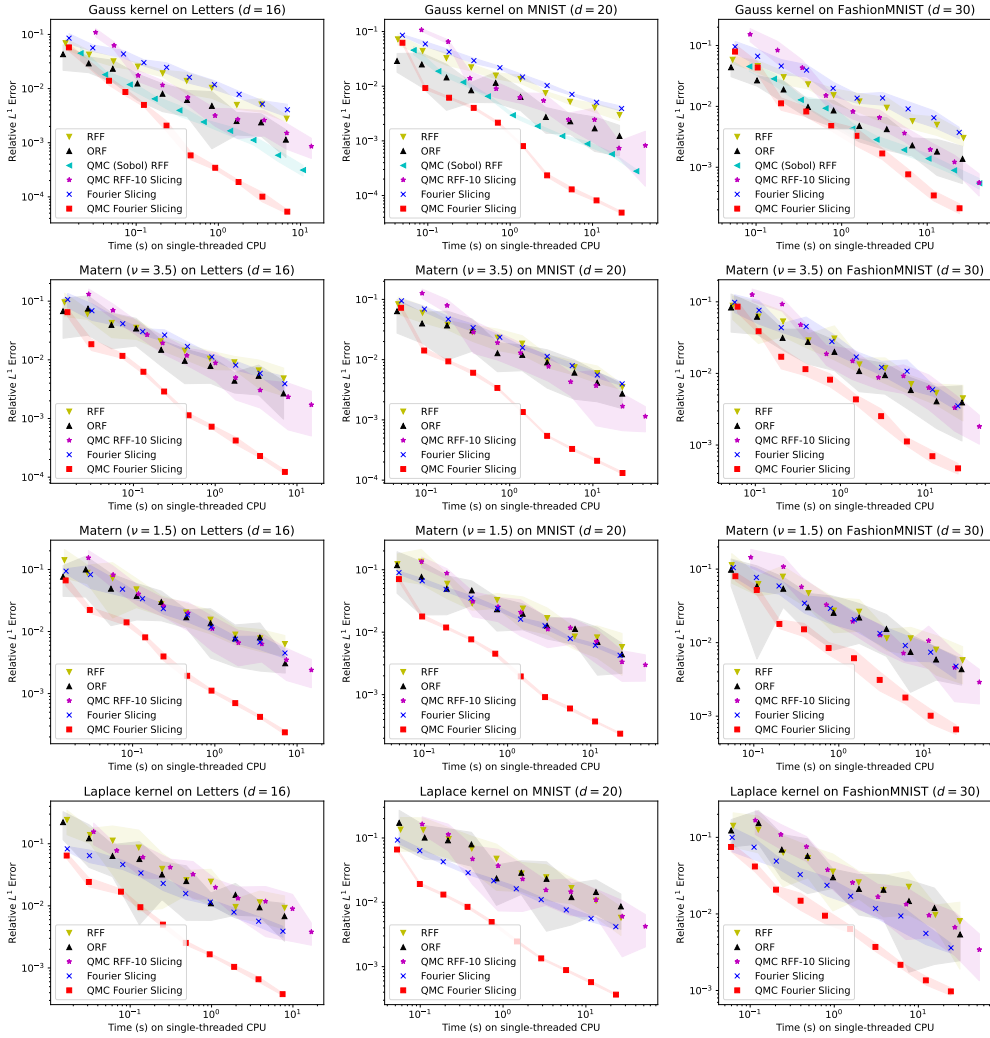


Figure 2: Loglog plot of the relative  $L^1$  approximation error versus computation time for computing the kernel summations (1) with different kernels and methods. We use the Letters dataset ( $M = N = 20000$  points), MNIST (reduced to dimension  $d = 20$  via PCA,  $M = N = 60000$  points) and FashionMNIST (reduced to dimension  $d = 30$  via PCA,  $M = N = 60000$  points). We run each method 10 times. The shaded area indicates the standard deviation of the error. For Fourier slicing, we use  $P = 5 \cdot 2^k$  slices for  $k = 1, \dots, 10$ . In order to obtain similar computation times, we use  $5 \cdot 2^{k-1}$  slices for RFF-10 slicing and  $D = 2P$  features for RFF and ORF.

we often observe much faster convergence rates than the proven worst case error rates  $r = \frac{d}{2(d-1)}$  on  $H^s(\mathbb{S}^{d-1})$  with the distance QMC designs, see Section 3. Furthermore, the slight advantage of the distance QMC designs versus the spherical designs for the Laplace kernel is because the former ones are chosen specifically for symmetric functions.

### 4.3 FAST KERNEL SUMMATION

Finally, we test our kernel approximation for computing the whole kernel sums from (1). For computing the one-dimensional kernel sums  $\sum_{n=1}^N w_n f(|\langle \xi_p, x_n - y_m \rangle|)$ , we use the following methods combined with either random or QMC points on the sphere:

- **(QMC) Sorting-Slicing:** For the negative distance kernel, we use the sorting algorithm from Hertrich et al. (2024), see also (Hertrich, 2024, Sec 3.2).

- **(QMC) Fourier-Slicing:** For the Gauss, Matérn and Laplace kernel, we use fast Fourier summation based on the non-equispaced fast Fourier transform (NFFT) for the one-dimensional kernel summation. A general overview on NFFTs and fast Fourier summation can be found in the text book (Plonka et al., 2023, Sec 7.5). Similar as in (Hertrich, 2024, Sec 2.3), we do not evaluate the one-dimensional basis functions, but directly compute the Fourier transforms. We revisit the background on the one-dimensional fast Fourier summation and specify the used parameters in Appendix I.
- **(QMC) RFF- $k$  Slicing:** For positive definite kernels (Gauss, Laplace, Matérn), we use one-dimensional random Fourier features with  $k$  features for the basis function  $f$ . For iid or orthogonal slices and  $k = 1$ , this approach is related to RFF and ORF, as outlined in Appendix G.

For  $\xi^P$ , we use the distance QMC designs, since we have seen in the previous subsection that it performs best among the QMC rules. Moreover, we use the randomization from Appendix E for the QMC design to obtain an unbiased estimator. We evaluate the kernel sums on Letters dataset ( $d = 16$ , Slate, 1991), MNIST (reduced to  $d = 20$  dimensions via PCA, LeCun et al., 1998) and FashionMNIST (reduced to  $d = 30$  dimensions via PCA, Xiao et al., 2017), where  $(x_1, \dots, x_N)$  and  $(y_1, \dots, y_M)$  constitute the whole dataset and the weights  $(w_1, \dots, w_N)$  are set to 1. In particular, we have  $M = N = 20000$  for the Letters dataset and  $M = N = 60000$  for MNIST and FashionMNIST. Then, we approximate the vector  $s = (s_1, \dots, s_M)$  from (1) and report the absolute error  $\|s - s_{\text{true}}\|_1$ . We choose the kernel parameters by the median rule with scale factor  $\gamma = 1$  based on 1000 example pairs  $(x, y)$ . We benchmark the computation times on a single thread of an AMD Ryzen Threadripper 7960X CPU and compare our results with RFF, ORF and (non-QMC) slicing. Since the QMC designs  $\xi^P$  depend neither on the dataset nor on the kernel, we consider its construction not as a part of the computation time. For the Gauss kernel, we also compare with QMC (Sobol) RFF (Avron et al., 2016), which is not applicable for the other kernels. We use  $P = 5 \cdot 2^k$  slices for  $k = 1, \dots, 10$  in the slicing method. In order to obtain similar computation times, we use  $5 \cdot 2^{k-1}$  slices for RFF-10 slicing and  $D = 2P$  features for RFF and ORF.

We visualize the approximation error (including standard deviations) in Figure 2 for the Gauss, Laplace and Matérn kernel. [The results for the negative distance kernel, for the thin plate spline kernel, for higher dimensional datasets \(including the full MNIST and FashionMNIST with  \$d = 784\$ \) and a GPU comparison are included in Appendix K.2. In addition, we apply the fast kernel summation to computing MMD gradient flows in Appendix K.3.](#) We can see that QMC Fourier-slicing has a significantly smaller error than the other methods. Moreover, it has the smallest standard deviation of the error.

## 5 CONCLUSIONS

**Summary and Outlook** We proposed a slicing approach to compute large kernel sums in  $\mathcal{O}(P(N + M + N_{\text{ft}} \log N_{\text{ft}}))$  instead of the naïve  $\mathcal{O}(NM)$  arithmetic operations. In the case of iid directions, we proved error bounds with rate  $\mathcal{O}(1/\sqrt{P})$ . To improve this rate, we proposed a QMC approach based on spherical quadrature rules. We demonstrated by numerical methods that our QMC slicing approach outperforms existing methods, where the advantage is most significant for dimensions  $d \leq 100$ . In the future, we want to improve our theoretical analysis on QMC slicing in order to match the convergence rate from the numerical section. One possible approach for that could be to study worst case errors for symmetric functions on the sphere, since the mappings  $\xi \rightarrow g_x(\xi)$  from Section 3.2 are always symmetric, [see Appendix K.4 for details](#). From a practical side, we want to apply the slicing approach in some actual applications.

**Limitations** In the numerical part, we observe significantly better error rates than we can prove theoretically, [see Appendix K.4 for a discussion](#). Moreover, the computation of the QMC directions can be very costly and depends strongly on the chosen method. For the spherical designs, it is even intractable for high dimensions. Finally, the advantage of QMC slicing becomes smaller for higher dimensions, which is a well-known effect for most QMC methods.

## REFERENCES

- 540  
541  
542 Fabian Altekrüger, Johannes Hertrich, and Gabriele Steidl. Neural Wasserstein gradient flows for  
543 discrepancies with Riesz kernels. In *International Conference on Machine Learning*, pp. 664–  
544 690. PMLR, 2023.
- 545 Michael Arbel, Anna Korba, Adil Salim, and Arthur Gretton. Maximum mean discrepancy gradient  
546 flow. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 547  
548 Kendall Atkinson. *An Introduction to Numerical Analysis*. Wiley, 1991. ISBN 978-0-471-62489-9.
- 549  
550 Kendall Atkinson and Weimin Han. *Spherical Harmonics and Approximations on the Unit Sphere:  
551 An Introduction*. Springer, Heidelberg, 2012. doi: 10.1007/978-3-642-25983-8.
- 552 Haim Avron, Vikas Sindhwani, Jiyan Yang, and Michael W. Mahoney. Quasi-Monte Carlo feature  
553 maps for shift-invariant kernels. *Journal of Machine Learning Research*, 17(120):1–38, 2016.
- 554  
555 Francis Bach. On the equivalence between kernel quadrature rules and random feature expansions.  
556 *The Journal of Machine Learning Research*, 18(1):714–751, 2017.
- 557 Eiichi Bannai and Etsuko Bannai. A survey on spherical designs and algebraic combinatorics on  
558 spheres. *European Journal of Combinatorics*, 30(6):1392–1425, 2009. doi: 10.1016/j.ejc.2008.  
559 11.007.
- 560  
561 Richard K Beatson and Garry N Newsam. Fast evaluation of radial basis functions: I. *Computers &  
562 Mathematics with Applications*, 24(12):7–19, 1992.
- 563 Ayoub Belhadji, Qianyu Julie Zhu, and Youssef Marzouk. On the design of scalable, high-precision  
564 spherical-radial Fourier features. *arXiv preprint arXiv:2408.13231*, 2024.
- 565  
566 Carlos Beltrán, Damir Ferizović, and Pedro R López-Gómez. Measure-preserving mappings from  
567 the unit cube to some symmetric spaces. *arXiv preprint 2303.00405*, 2023.
- 568  
569 Gregory Beylkin. On the fast Fourier transform of functions with singularities. *Applied and Com-  
570 putational Harmonic Analysis*, 2(4):363–381, 1995.
- 571  
572 Salomon Bochner. Monotone Funktionen, Stieltjessche Integrale und harmonische Analyse. *Math-  
573 ematische Annalen*, 108(1):378–410, 1933.
- 574  
575 Andriy Bondarenko, Danylo Radchenko, and Maryna Viazovska. Optimal asymptotic bounds for  
576 spherical designs. *Annals of Mathematics*, 178(2):443–452, 2013. doi: 10.4007/annals.2013.178.  
2.2.
- 577  
578 Johann Brauchart and Kerstin Hesse. Numerical integration over spheres of arbitrary dimension.  
579 *Constructive Approximation*, 25(1):41–71, 2007.
- 580  
581 Johann Brauchart, Edward Saff, Ian Sloan, and Robert Womersley. QMC designs: optimal order  
582 quasi Monte Carlo integration schemes on the sphere. *Mathematics of Computation*, 83(290):  
2821–2851, 2014.
- 583  
584 Benjamin Charlier, Jean Feydy, Joan Alexis Glaunès, François-David Collin, and Ghislain Durif.  
585 Kernel operations on the GPU, with autodiff, without memory overflows. *Journal of Machine  
586 Learning Research*, 22(74):1–6, 2021.
- 587  
588 Zonghao Chen, Aratrika Mustafi, Pierre Glaser, Anna Korba, Arthur Gretton, and Bharath K  
589 Sriperumbudur. (de)-regularized maximum mean discrepancy gradient flow. *arXiv preprint  
arXiv:2409.14980*, 2024.
- 590  
591 Feng Dai and Yuan Xu. *Approximation Theory and Harmonic Analysis on Spheres and Balls*.  
592 Springer, New York, 2013. doi: 10.1007/978-1-4614-6660-4.
- 593  
Philippe Delsarte, J. M. Goethals, and Johan Jacob Seidel. Spherical codes and designs. *Geometriae  
Dedicata*, 6(3):363–388, 1977. ISSN 1572-9168. doi: 10.1007/bf03187604.

- 594 Alok Dutt and Vladimir Rokhlin. Fast Fourier transforms for nonequispaced data. *SIAM Journal on*  
595 *Scientific Computing*, 14(6):1368–1393, 1993.
- 596
- 597 Alexandre Galashov, Valentin de Bortoli, and Arthur Gretton. Deep MMD gradient flow without  
598 adversarial training. *arXiv preprint arXiv:2405.06780*, 2024.
- 599 Damien Garreau, Wittawat Jitkrittum, and Motonobu Kanagawa. Large sample analysis of the  
600 median heuristic. *arXiv preprint arXiv:1707.07269*, 2017.
- 601
- 602 Izrail S Gradshteyn and Iosif M Ryzhik. *Table of Integrals, Series, and Products*. Academic Press  
603 New York, seventh edition, 2007.
- 604 Manuel Gräf and Daniel Potts. On the computation of spherical designs by a new optimization  
605 approach based on fast spherical Fourier transforms. *Numerische Mathematik*, 119(4):699–724,  
606 2011.
- 607 Manuel Gräf, Daniel Potts, and Gabriele Steidl. Quadrature rules, discrepancies and their relations to  
608 half-toning on the torus and the sphere. *SIAM Journal on Scientific Computing*, 34(5):2760–2791,  
609 2012.
- 610
- 611 Leslie Greengard and Vladimir Rokhlin. A fast algorithm for particle simulations. *Journal of*  
612 *Computational Physics*, 73(2):325–348, 1987.
- 613
- 614 Leslie Greengard and John Strain. The fast Gauss transform. *SIAM Journal on Scientific and*  
615 *Statistical Computing*, 12(1):79–94, 1991.
- 616 Philip Greengard, Manas Rachh, and Alex Barnett. Equispaced Fourier representations for efficient  
617 Gaussian process regression from a billion data points. *arXiv preprint arXiv:2210.10210*, 2022.
- 618
- 619 Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. A kernel  
620 method for the two-sample-problem. *Advances in Neural Information Processing Systems*, 19,  
621 2006.
- 622 Wolfgang Hackbusch. A sparse matrix arithmetic based on H-Matrices. Part I: Introduction to H-  
623 Matrices. *Computing*, 62(2):89–108, 1999.
- 624 Paul Hagemann, Johannes Hertrich, Fabian Altekrüger, Robert Beinert, Jannis Chemseddine, and  
625 Gabriele Steidl. Posterior sampling based on gradient flows of the MMD with negative distance  
626 kernel. In *International Conference on Learning Representations*, 2024.
- 627
- 628 John H Halton. On the efficiency of certain quasi-random sequences of points in evaluating multi-  
629 dimensional integrals. *Numerische Mathematik*, 2:84–90, 1960.
- 630 Abolfazl Hashemi, Hayden Schaeffer, Robert Shi, Ufuk Topcu, Giang Tran, and Rachel Ward. Gen-  
631 eralization bounds for sparse random feature expansions. *Applied and Computational Harmonic*  
632 *Analysis*, 62:310–330, 2023.
- 633
- 634 Johannes Hertrich. Fast kernel summation in high dimensions via slicing and Fourier transforms.  
635 *arXiv preprint arXiv:2401.08260*, 2024.
- 636 Johannes Hertrich, Christian Wald, Fabian Altekrüger, and Paul Hagemann. Generative sliced MMD  
637 flows with Riesz kernels. In *International Conference on Learning Representations*, 2024.
- 638
- 639 Kerstin Hesse. A lower bound for the worst-case cubature error on spheres of arbitrary dimension.  
640 *Numerische Mathematik*, 103:413–433, 2006.
- 641 Zhen Huang, Jiajin Sun, and Yian Huang. Quasi-Monte Carlo features for kernel approximation. In  
642 *International Conference on Machine Learning*, 2024.
- 643
- 644 Jens Keiner, Stefan Kunis, and Daniel Potts. Using NFFT3 - a software library for various noneq-  
645 uispaced fast Fourier transforms. *ACM Transactions on Mathematical Software*, 36:Article 19,  
646 1–30, 2009.
- 647 Diederik P Kingma and Jimmy Ba. Adam: a method for stochastic optimization. In *Proceedings of*  
*the ICLR '15*, 2015.

- 648 Tobias Knopp, Marija Boberg, and Mirco Grosser. NFFT.jl: Generic and fast julia implementation  
649 of the nonequidistant fast Fourier transform. *SIAM Journal on Scientific Computing*, 45(3):C179–  
650 C205, 2023.
- 651 Soheil Kolouri, Kimia Nadjahi, Shahin Shahrapour, and Umut Şimşekli. Generalized sliced prob-  
652 ability metrics. In *IEEE International Conference on Acoustics, Speech and Signal Processing*,  
653 pp. 4513–4517, 2022.
- 654 Stefan Kunis, Daniel Potts, and Gabriele Steidl. Fast Gauss transforms with complex parameters  
655 using NFFTs. *Journal of Numerical Mathematics*, 14(4):295, 2006.
- 656 Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
657 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 658 Dongryeol Lee and Alexander Gray. Fast high-dimensional kernel summations using the Monte  
659 Carlo multipole method. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou (eds.), *Advances  
660 in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008.
- 661 Zhu Li, Jean-Francois Ton, Dino Oglic, and Dino Sejdinovic. Towards a unified analysis of random  
662 Fourier features. *The Journal of Machine Learning Research*, 22(1):4887–4937, 2021.
- 663 Jen Ning Lim, Juan Kuntz, Samuel Power, and Adam Michael Johansen. Momentum particle max-  
664 imum likelihood. In *Proceedings of the 41st International Conference on Machine Learning*,  
665 volume 235 of *Proceedings of Machine Learning Research*, pp. 29816–29871. PMLR, 2024.
- 666 Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose Bayesian inference  
667 algorithm. *Advances in Neural Information Processing Systems*, 29, 2016.
- 668 William B March, Bo Xiao, and George Biros. ASKIT: Approximate skeletonization kernel-  
669 independent treecode in high dimensions. *SIAM Journal on Scientific Computing*, 37(2):A1089–  
670 A1110, 2015a.
- 671 William B March, Bo Xiao, D Yu Chenhan, and George Biros. An algebraic parallel treecode in ar-  
672 bitrary dimensions. In *2015 IEEE International Parallel and Distributed Processing Symposium*,  
673 pp. 571–580. IEEE, 2015b.
- 674 Victor Minden, Anil Damle, Kenneth L Ho, and Lexing Ying. Fast spatial Gaussian process maxi-  
675 mum likelihood estimation via skeletonization factorizations. *Multiscale Modeling & Simulation*,  
676 15(4):1584–1611, 2017.
- 677 Marina Munkhoeva, Yermek Kapushev, Evgeny Burnaev, and Ivan Oseledets. Quadrature-based  
678 features for kernel approximation. *Advances in Neural Information Processing Systems*, 31, 2018.
- 679 Khai Nguyen, Nicola Barileto, and Nhat Ho. Quasi-Monte Carlo for 3d sliced Wasserstein. In  
680 *International Conference on Learning Representations*, 2024.
- 681 NIST. NIST Digital Library of Mathematical Functions. <https://dlmf.nist.gov/>, Release  
682 1.2.1 of 2024-06-15, 2024.
- 683 Emanuel Parzen. On estimation of a probability density function and mode. *The Annals of Mathe-  
684 matical Statistics*, 33(3):1065–1076, 1962.
- 685 Gerlind Plonka, Daniel Potts, Gabriele Steidl, and Manfred Tasche. *Numerical Fourier Analysis*.  
686 Springer, 2 edition, 2023.
- 687 Daniel Potts and Gabriele Steidl. Fast summation at nonequispaced knots by NFFT. *SIAM Journal  
688 on Scientific Computing*, 24(6):2013–2037, 2003.
- 689 Daniel Potts, Gabriele Steidl, and Arthur Nieslony. Fast convolution with radial kernels at nonequi-  
690 spaced knots. *Numerische Mathematik*, 98:329–351, 2004.
- 691 Michael Quellmalz. The Funk–Radon transform for hyperplane sections through a common point.  
692 *Analysis and Mathematical Physics*, 10(38), 2020. doi: 10.1007/s13324-020-00383-2.

- 702 Julien Rabin, Gabriel Peyré, Julie Delon, and Marc Bernet. Wasserstein barycenter and its applica-  
703 tion to texture mixing. In *International Conference on Scale Space and Variational Methods in*  
704 *Computer Vision*, pp. 435–446. Springer, 2012.
- 705 David L. Ragozin. Rotation invariant measure algebras on Euclidean space. *Indiana University*  
706 *Mathematics Journal*, 23(12):1139–54, 1974.
- 707 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in*  
708 *Neural Information Processing Systems*, 20, 2007.
- 709 Murray Rosenblatt. Remarks on some nonparametric estimates of a density function. *The Annals of*  
710 *Mathematical Statistics*, pp. 832–837, 1956.
- 711 Boris Rubin. Notes on Radon transforms in integral geometry. *Fractional Calculus and Applied*  
712 *Analysis*, 6(1):25–72, 2003.
- 713 Boris Rubin. The  $\lambda$ -cosine transforms with odd kernel and the hemispherical transform. *Fractional*  
714 *Calculus and Applied Analysis*, 17(3):765–806, 2014. doi: 10.2478/s13540-014-0198-9.
- 715 Nicolaj Rux, Michael Quellmalz, and Gabriele Steidl. Slicing of radial functions: a dimension walk  
716 in the Fourier space. *arXiv preprint arXiv:2408.11612*, 2024.
- 717 John P Ryan, Sebastian E Ament, Carla P Gomes, and Anil Damle. The fast kernel transform. In  
718 *International Conference on Artificial Intelligence and Statistics*, pp. 11669–11690. PMLR, 2022.
- 719 Bernhard Schölkopf and Alexander J Smola. *Learning with Kernels: Support Vector Machines,*  
720 *Regularization, Optimization, and Beyond*. MIT Press, 2002.
- 721 John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge Univer-  
722 sity Press, 2004.
- 723 Yu-hsuan Shih, Garrett Wright, Joakim Andén, Johannes Blaschke, and Alex H Barnett. cuFIN-  
724 UFFT: a load-balanced GPU library for general-purpose nonuniform FFTs. In *IEEE Inter-*  
725 *national Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2021. doi:  
726 10.1109/IPDPSW52791.2021.00105.
- 727 David Slate. Letter Recognition. UCI Machine Learning Repository, 1991. DOI:  
728 <https://doi.org/10.24432/C5ZP40>.
- 729 Il’ya Meerovich Sobol’. On the distribution of points in a cube and the approximate evaluation of  
730 integrals. *Zhurnal Vychislitel’noi Matematiki i Matematicheskoi Fiziki*, 7(4):784–802, 1967.
- 731 Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer Science & Business  
732 Media, 2008.
- 733 Danica J Sutherland and Jeff Schneider. On the error of random Fourier features. In *Proceedings of*  
734 *the Thirty-First Conference on Uncertainty in Artificial Intelligence*, UAI’15, pp. 862–871. AUAI  
735 Press, 2015.
- 736 Gabor Székely. E-statistics: The energy of statistical samples. *Technical Report*, Bowling Green  
737 University, 2002.
- 738 Eugene Tyrtshnikov. Mosaic-skeleton approximations. *Calcolo*, 33:47–57, 1996.
- 739 Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau,  
740 Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. Scipy 1.0: funda-  
741 mental algorithms for scientific computing in Python. *Nature methods*, 17(3):261–272, 2020.
- 742 Holger Wendland. *Scattered Data Approximation*. Cambridge University Press, 2004. doi: 10.1017/  
743 CBO9780511617539.
- 744 Christopher KI Williams and Carl Edward Rasmussen. *Gaussian Processes for Machine Learning*,  
745 volume 2. MIT Press, 2006.

756 Robert S Womersley. Efficient spherical designs with good geometric properties. In J. Dick, F. Kuo,  
757 and H. Wozniakowski (eds.), *Contemporary Computational Mathematics - A Celebration of the*  
758 *80th Birthday of Ian Sloan*, pp. 1243–1285. Springer, 2018. doi: 10.1007/978-3-319-72456-0\_57.  
759

760 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for bench-  
761 marking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

762 Changjiang Yang, Ramani Duraiswami, Nail A Gumerov, and Larry S Davis. Improved fast Gauss  
763 transform and efficient kernel density estimation. In *IEEE International Conference on Computer*  
764 *Vision*, pp. 664–671 vol.1, 2003.

765 Changjiang Yang, Ramani Duraiswami, and Larry S Davis. Efficient kernel machines using the  
766 improved fast Gauss transform. *Advances in Neural Information Processing Systems*, 17, 2004.  
767

768 Norman Yarvin and Vladimir Rokhlin. An improved fast multipole algorithm for potential fields on  
769 the line. *SIAM Journal on Numerical Analysis*, 36(2):629–666, 1999.

770 Felix Xinnan X Yu, Ananda Theertha Suresh, Krzysztof M Choromanski, Daniel N Holtmann-Rice,  
771 and Sanjiv Kumar. Orthogonal random features. *Advances in Neural Information Processing*  
772 *Systems*, 29, 2016.  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A BASIS FUNCTION PAIRS $(F, f)$

We include a list of pairs of basis functions  $(F, f)$  that fulfill the slicing formula and the corresponding Fourier transforms in Table 2. The pairs are taken from (Hertrich, 2024, Table 1). We use the convention that the Fourier transform of a function  $g \in L_1(\mathbb{R}^d)$  is defined by

$$\mathcal{F}_d[g](\omega) = \int_{\mathbb{R}^d} g(x) e^{-2\pi i \langle \omega, x \rangle} dx. \quad (11)$$

The basis functions from the table involve a couple of special functions, which are defined as follows:

- Gamma function:  $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$ ,  $\text{Re}(z) > 0$ ,
- Modified Bessel function of first kind:  $I_\alpha(x) = \sum_{m=0}^\infty \frac{1}{m! \Gamma(m+\alpha+1)} \left(\frac{x}{2}\right)^{2m+\alpha}$ ,
- Modified Bessel function of second kind:  $K_\alpha(x) = \frac{\pi}{2} \frac{I_{-\alpha}(x) - I_\alpha(x)}{\sin(\alpha\pi)}$ ,
- Modified Struve function:  $\mathbf{L}_\alpha(x) = \sum_{m=0}^\infty \frac{1}{\Gamma(m+\frac{3}{2})\Gamma(m+\alpha+\frac{3}{2})} \left(\frac{x}{2}\right)^{2m+\alpha+1}$ .

The formula for the Fourier transform of the Matérn kernel (and thus for the Laplace kernel with  $\nu = 1/2$ ) can be found in (Williams & Rasmussen, 2006, (4.15)). Consequently, we recover the well-known results that the Fourier transforms of the Gauss, Laplace and Matérn kernels are the Fourier transforms of the Gauss, Cauchy and Student- $t$  (with  $2\nu$  degrees of freedom) distribution. To compute the Fourier transforms  $\mathcal{F}_1^{-1}[f(|\cdot|)]$ , we apply (Rux et al., 2024, Prop 3.1), see also (Hertrich, 2024, Lem 6) for the Gauss kernel. For the Riesz and thin plate spline kernel, the Fourier transform does not exist in a classical sense, but only in a distributional one, see Rux et al. (2024) and (Wendland, 2004, Sect 8.3) for details. For the sliced Laplace  $f(x) = \exp(-\alpha x)$ , we have by (Hertrich, 2024, (3)) and (Gradshteyn & Ryzhik, 2007, 3.387.5)

$$F(t) = \int_0^1 \exp(-\alpha st) (1-s^2)^{\frac{d-3}{2}} ds = \frac{\sqrt{\pi} 2^{\frac{d-4}{2}} \Gamma(\frac{d-1}{2})}{(\alpha t)^{\frac{d-2}{2}}} \left( I_{\frac{d-2}{2}}(-\alpha t) + \mathbf{L}_{\frac{d-2}{2}}(-\alpha t) \right). \quad (12)$$

Table 2: Basis functions  $F$  for different kernels  $K(x, y) = F(\|x - y\|)$  and corresponding basis functions  $f$  from  $k(x, y) = f(\|x - y\|)$ . We added the inverse Fourier transforms  $\mathcal{F}_d^{-1}[F(\|\cdot\|)]$  and  $\mathcal{F}_1^{-1}[f(|\cdot|)]$  to the table.

Kernel	$F(x)$	$\mathcal{F}_d^{-1}[F(\ \cdot\ )](\ \omega\ )$	$f(x)$	$\mathcal{F}_1^{-1}[f( \omega )]$
Gauss	$\exp(-\frac{x^2}{2\sigma^2})$	$(2\pi\sigma^2)^{d/2} \exp(-2\pi^2\sigma^2\omega^2)$	${}_1F_1(\frac{d}{2}; \frac{1}{2}; -\frac{x^2}{2\sigma^2})$	$\frac{\pi\sigma \exp(-2\pi^2\sigma^2\omega^2) (2\pi^2\sigma^2\omega^2)^{(d-1)/2}}{2\Gamma(\frac{d}{2})}$
Laplace	$\exp(-\alpha x)$	$\frac{\Gamma(\frac{d+1}{2}) 2^d \pi^{\frac{d-1}{2}}}{\alpha^d} (1 + \frac{4\pi^2\omega^2}{\alpha^2})^{-\frac{d+1}{2}}$	$\sum_{n=0}^\infty \frac{(-1)^n \alpha^n \sqrt{\pi} \Gamma(\frac{d+n}{2})}{n! \Gamma(\frac{d}{2}) \Gamma(\frac{d+n}{2})} x^n$	$\frac{\Gamma(\frac{d+1}{2}) 2^d \pi^{\frac{d-1}{2}}  \omega ^{d-1}}{\Gamma(\frac{d}{2}) \alpha^d} (1 + \frac{4\pi^2\omega^2}{\alpha^2})^{-\frac{d+1}{2}}$
Sliced Laplace	See (12)	$\frac{\Gamma(\frac{d}{2})}{\pi^{d/2}  \omega ^{d-1}} \frac{2\alpha}{\alpha^2 + 4\pi^2\omega^2}$	$\exp(-\alpha x)$	$\frac{2\alpha}{\alpha^2 + 4\pi^2\omega^2}$
Matérn	$\frac{2^{1-\nu}}{\Gamma(\nu)} (\frac{\sqrt{2\nu} x}{\beta})^\nu K_\nu(\frac{\sqrt{2\nu} x}{\beta})$	$\frac{\Gamma(\frac{2\nu+d}{2}) 2^{\frac{d}{2}} \pi^{\frac{d}{2}} \beta^d}{\Gamma(\nu) \nu^{\frac{d}{2}}} (1 + 2\pi^2 \beta^2 \omega^2)^{-\frac{2\nu+d}{2}}$	(Hertrich, 2024, Appendix C)	$\frac{\Gamma(\frac{2\nu+d}{2}) 2^{\frac{d}{2}} \pi^{\frac{d}{2}} \beta^d  \omega ^{d-1}}{\Gamma(\frac{d}{2}) \Gamma(\nu) \nu^{\frac{d}{2}}} (1 + \frac{2\pi^2 \beta^2 \omega^2}{\nu})^{-\frac{2\nu+d}{2}}$
Riesz for $r \in (0, 2)$	$-x^r$	not a function	$-\frac{\sqrt{\pi} \Gamma(\frac{d+r}{2})}{\Gamma(\frac{d}{2}) \Gamma(\frac{d+r}{2})} x^r$	not a function
Thin Plate Spline	$x^2 \log(x)$	not a function	$dx^2 \log(x) + C_d x^2$ , with $C_d = \frac{d}{2}(H_{\frac{d}{2}} - 2 + \log(4))$	not a function

## B SPHERICAL SOBOLEV SPACES

We briefly introduce spherical harmonics and Sobolev spaces following Atkinson & Han (2012). We denote by  $|\mathbb{S}^{d-1}| = \frac{2\pi^{d/2}}{\Gamma(d/2)}$  the volume of the unit sphere  $\mathbb{S}^{d-1}$ . Let  $Y_n^k$ ,  $k = 1, \dots, N_{n,d}$  be an  $L_2$ -orthonormal basis of spherical harmonics of degree  $n \in \mathbb{N}_0$ , i.e., harmonic polynomials in  $d$  variables that are homogeneous of degree  $n$ . Here, the dimension of the space of spherical harmonics of degree  $n$  is

$$N_{n,d} = \frac{(2n+d-2)(n+d-3)!}{n!(d-2)!} \simeq \frac{2}{(d-2)!} n^{d-2} \quad \text{for } n \rightarrow \infty. \quad (13)$$

The spherical harmonics  $Y_n^k$ ,  $n \in \mathbb{N}_0$ ,  $k = 1, \dots, N_{n,d}$  form an orthonormal basis of  $L_2(\mathbb{S}^{d-1})$ .



The spherical Sobolev space  $H^s(\mathbb{S}^{d-1})$  for  $s \in \mathbb{R}$  can be defined as the completion of  $C^\infty(\mathbb{S}^{d-1})$  with respect to the norm

$$\|g\|_{H^s(\mathbb{S}^{d-1})}^2 = \sum_{n=0}^{\infty} \sum_{k=1}^{N_{n,d}} \left(n + \frac{d-2}{2}\right)^{2s} |\langle g, Y_n^k \rangle_{L_2(\mathbb{S}^{d-1})}|^2, \quad (14)$$

where  $\langle g, Y_n^k \rangle_{L_2(\mathbb{S}^{d-1})} = \int_{\mathbb{S}^{d-1}} g(\xi) Y_n^k(\xi) d\xi$ . Note that the factor  $(n + \frac{d-2}{2})^{2d}$  can be replaced by another one with the same asymptotic behavior with respect to  $n$ , yielding an equivalent norm. For  $s = 0$ , we can identify  $H^s(\mathbb{S}^{d-1})$  with  $L_2(\mathbb{S}^{d-1})$ . The Sobolev spaces are nested in the sense that  $H^s(\mathbb{S}^{d-1}) \subset H^t(\mathbb{S}^{d-1})$  whenever  $s > t$ . If  $s > \frac{d-1}{2}$ , each function in the Sobolev space  $H^s(\mathbb{S}^{d-1})$  is continuous (more specifically, it has a continuous representative). If  $s$  is an integer, then  $H^s(\mathbb{S}^{d-1})$  consists of all functions whose (distributional) derivatives up to order  $s$  are square integrable, cf. (Dai & Xu, 2013, Sect 4.5 and 4.8). An alternative characterization of the Sobolev norm uses the Laplace–Beltrami operator  $\Delta_*$ , which consists of the spherical part of the Laplace, and is given by

$$\|g\|_{H^s(\mathbb{S}^{d-1})} = \left\| \left(-\Delta_* + \left(\frac{d-2}{2}\right)^2\right)^{s/2} g \right\|_{L_2(\mathbb{S}^{d-1})}.$$

If  $s$  is an even integer, the operator applied to  $g$  is a usual differentiable operator, otherwise it is a pseudodifferential operator.

## C PROOF OF THEOREM 1

**i):** Since  $F$  is continuous and positive definite, so is  $f$  by (Rux et al., 2024, Corollary 4.11). Further, because  $\mathbb{E}[f(|\langle \xi, x \rangle|)] = F(\|x\|)$  and due to the fact that positive definite functions are maximal in the origin, we deduce

$$\mathbb{V}_d[f] = \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [f(|\langle \xi, x \rangle|)^2] - F(\|x\|)^2 \leq f(0)^2 - F(\|x\|)^2.$$

**ii):** We write  $\xi = (\xi_1, \dots, \xi_d)^T$ , where  $\xi_d \in [-1, 1]$  denotes the  $d$ -th component of  $\xi$ . We assume w.l.o.g. that  $x = \lambda e_d$  with  $\lambda \neq 0$  and recall that we can decompose the unnormalized surface measure on  $\mathbb{S}^{d-1}$  as

$$d\mathbb{S}^{d-1}(\xi) = d\mathbb{S}^{d-2}(\eta)(1-t^2)^{\frac{d-3}{2}} dt,$$

where  $\xi = \sqrt{1-t^2}\eta + te_d$  and  $\eta \in \mathbb{S}^{d-2} \times \{0\}$ , see (Atkinson & Han, 2012, (1.16)). Consequently,

$$\begin{aligned} \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [|\langle \xi, x \rangle|^{2r}] &= \frac{\|x\|^{2r}}{|\mathbb{S}^{d-1}|} \int_{\mathbb{S}^{d-1}} |\xi_d|^{2r} d\mathbb{S}^{d-1}(\xi) \\ &= \frac{\|x\|^{2r}}{|\mathbb{S}^{d-1}|} \int_{\mathbb{S}^{d-2}} \int_{-1}^1 |t|^{2r} (1-t^2)^{\frac{d-3}{2}} d\mathbb{S}^{d-2}(\eta) dt = \frac{\|x\|^{2r}}{|\mathbb{S}^{d-1}|} |\mathbb{S}^{d-2}| \int_{-1}^1 (t^2)^r (1-t^2)^{\frac{d-3}{2}} dt. \end{aligned}$$

Further, with  $B(\cdot, \cdot)$  denoting the Beta function and the substitution  $u = t^2$ , we have

$$\int_{-1}^1 t^{2r} (1-t^2)^{\frac{d-3}{2}} dt = \int_0^1 u^{r-\frac{1}{2}} (1-u)^{\frac{d-1}{2}-1} du = B\left(r + \frac{1}{2}, \frac{d-1}{2}\right) = \frac{\Gamma\left(r + \frac{1}{2}\right) \Gamma\left(\frac{d-1}{2}\right)}{\Gamma\left(r + \frac{d}{2}\right)}, \quad (15)$$

and

$$\frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} = \frac{2\pi^{\frac{d-1}{2}} \Gamma\left(\frac{d}{2}\right)}{\Gamma\left(\frac{d-1}{2}\right) 2\pi^{\frac{d}{2}}} = \frac{\Gamma\left(\frac{d}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{d-1}{2}\right)}.$$

Combining both expressions yields

$$\mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [|\langle \xi, x \rangle|^{2r}] = \|x\|^{2r} \frac{\Gamma\left(r + \frac{1}{2}\right) \Gamma\left(\frac{d}{2}\right)}{\Gamma\left(r + \frac{d}{2}\right) \sqrt{\pi}}, \quad (16)$$

and hence

$$\frac{\pi \Gamma\left(\frac{d+r}{2}\right)^2}{\Gamma\left(\frac{d}{2}\right)^2 \Gamma\left(\frac{r+1}{2}\right)^2} \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [|\langle \xi, x \rangle|^{2r}] = \frac{\|x\|^{2r} \sqrt{\pi} \Gamma\left(r + \frac{1}{2}\right)}{\Gamma\left(\frac{r+1}{2}\right)^2} \frac{\Gamma\left(\frac{d+r}{2}\right)^2}{\Gamma\left(\frac{d}{2}\right) \Gamma\left(r + \frac{d}{2}\right)}.$$

The asymptotic expansion  $\Gamma(z+a)/\Gamma(z) \sim z^a$  for  $z \rightarrow \infty$ , see (NIST, 5.11.12), implies that

$$\lim_{d \rightarrow \infty} \frac{\Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{d}{2})\Gamma(r+\frac{d}{2})} = 1.$$

On the other hand,  $\frac{\Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{d}{2})\Gamma(r+\frac{d}{2})}$  is increasing with respect to  $d$  because the identity  $\Gamma(z+1) = z\Gamma(z)$  yields that

$$\frac{\Gamma(\frac{d+2+r}{2})^2}{\Gamma(\frac{d+2}{2})\Gamma(r+\frac{d+2}{2})} \bigg/ \frac{\Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{d}{2})\Gamma(r+\frac{d}{2})} = \frac{(\frac{d+r}{2})^2}{(\frac{d}{2})(r+\frac{d}{2})} = \frac{(d+r)^2}{2dr+d^2} \geq 1.$$

Hence, we see that

$$\frac{\Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{d}{2})\Gamma(r+\frac{d}{2})} \leq 1,$$

and thus finally

$$\begin{aligned} \mathbb{V}_d[f] &\leq \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [f(|\langle \xi, x \rangle|)^2] = \frac{\pi \Gamma(\frac{d+r}{2})^2}{\Gamma(\frac{d}{2})^2 \Gamma(\frac{r+1}{2})^2} \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [|\langle \xi, x \rangle|^{2r}] \\ &\leq \frac{\sqrt{\pi} \Gamma(r+\frac{1}{2})}{\Gamma(\frac{r+1}{2})^2} \|x\|^{2r}, \end{aligned}$$

which proves (7).

**iii):** We move on to the thin plate spline kernel with  $f(t) = d|t|^2 \log(|t|) + C_d|t|^2$ . As above, w.l.o.g. we assume that  $x = se_d$  with  $s \geq 0$ . Then

$$\begin{aligned} &\mathbb{E} [f(|\langle \xi, x \rangle|)^2] \\ &= \frac{|\mathbb{S}^{d-2}|}{|\mathbb{S}^{d-1}|} 2 \int_0^1 (ds^2 \xi_d^2 \log(s\xi_d) - C_d s^2 \xi_d^2)^2 (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \\ &= \frac{\Gamma(\frac{d}{2})}{\sqrt{\pi} \Gamma(\frac{d-1}{2})} \frac{d^2}{2} \int_0^1 s^4 \xi_d^4 \left( 2 \log(\xi_d) + 2 \log(s) + H_{\frac{d}{2}} - 2 + \log(4) \right)^2 (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \\ &= \frac{d^2 \Gamma(\frac{d}{2})}{2\sqrt{\pi} \Gamma(\frac{d-1}{2})} s^4 \left( 4 \int_0^1 \xi_d^4 \log^2(\xi_d) (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \right. \\ &\quad \left. + 4(2 \log(s) + H_{\frac{d}{2}} - 2 + \log(4)) \int_0^1 \xi_d^4 \log(\xi_d) (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \right. \\ &\quad \left. + \left( 2 \log(s) + H_{\frac{d}{2}} - 2 + \log(4) \right)^2 \int_0^1 \xi_d^4 (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \right). \end{aligned} \tag{17}$$

We analyze the terms separately. Denote by  $\psi^{(n)}$  the  $n$ -th derivative of the digamma function and by  $\gamma \approx 0.57$  the Euler–Mascheroni constant. In the following, we use the asymptotics

$$H_x = \log(x) + \gamma + \mathcal{O}(x^{-1}), \quad \psi^{(0)}(x) = \log(x) + \mathcal{O}(x^{-1}), \quad \psi^{(1)}(x) = \mathcal{O}(x^{-1}).$$

By (Gradshteyn & Ryzhik, 2007, 4.261.21), we have

$$\begin{aligned} &\int_0^1 \xi_d^4 \log^2(\xi_d) (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d \\ &= \frac{3\sqrt{\pi} \Gamma(\frac{d-1}{2})}{32\Gamma(\frac{d+4}{2})} \left( \psi^{(1)}\left(\frac{5}{2}\right) - \psi^{(1)}\left(\frac{d+4}{2}\right) + \left( \psi^{(0)}\left(\frac{5}{2}\right) - \psi^{(0)}\left(\frac{d+4}{2}\right) \right)^2 \right), \\ &= \frac{3\sqrt{\pi} \Gamma(\frac{d-1}{2})}{32\Gamma(\frac{d+4}{2})} \left( \log(d)^2 - 2 \left( \psi^{(0)}\left(\frac{5}{2}\right) + \log(2) + \mathcal{O}\left(\frac{1}{d}\right) \right) \log(d) \right. \\ &\quad \left. + \psi^{(1)}\left(\frac{5}{2}\right) + \left( \psi^{(0)}\left(\frac{5}{2}\right) + \log(2) \right)^2 + \mathcal{O}(d^{-1}) \right), \end{aligned}$$

972 and

$$973 \int_0^1 \xi_d^4 \log(\xi_d) (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d = -\frac{\sqrt{\pi}\Gamma(\frac{d-1}{2})}{16\Gamma(\frac{d+4}{2})} \left(-8 + 3H_{1+\frac{d}{2}} + \log(64)\right)$$

$$974 = -\frac{\sqrt{\pi}\Gamma(\frac{d-1}{2})}{16\Gamma(\frac{d+4}{2})} (3\log(d) + 3\gamma - 8 + \log(8) + \mathcal{O}(d^{-1})).$$

975 Recall from (15) that

$$976 \int_0^1 \xi_d^4 (1 - \xi_d^2)^{\frac{d-3}{2}} d\xi_d = \frac{3\sqrt{\pi}\Gamma(\frac{d-1}{2})}{8\Gamma(\frac{d+4}{2})}.$$

977 Furthermore, using that  $\log(4) = 2\log(2)$ , we obtain

$$978 2\log(s) + H_{\frac{d}{2}} - 2 + \log(4) = 2\log(s) + \log(d) + \gamma + \log(2) - 2 + \mathcal{O}(d^{-1}),$$

$$979 \left(2\log(s) + H_{\frac{d}{2}} - 2 + \log(4)\right)^2 = (2\log(s) + \log(d) + \gamma + \log(2) - 2 + \mathcal{O}(d^{-1}))^2$$

$$980 = 4\log(s)^2 + \log(d)^2 + 4\log(s)\log(d)$$

$$981 + 2(\gamma + \log(2) - 2)\log(d) + 4\log(s)(\gamma + \log(2) - 2)$$

$$982 + (\gamma + \log(2) - 2)^2 + \mathcal{O}(d^{-1}).$$

983 Plugging this into (17) yields

$$984 4 \frac{\Gamma(\frac{d+4}{2})}{\sqrt{\pi}\Gamma(\frac{d-1}{2})} \frac{\sqrt{\pi}\Gamma(\frac{d-1}{2})}{\Gamma(\frac{d}{2}) d^2 s^4} \mathbb{E}[f(|\langle \xi, x \rangle|)^2] = \frac{1+2/d}{s^4} \mathbb{E}[f(|\langle \xi, x \rangle|)^2]$$

$$985 = \frac{3}{4} \left( \log(d)^2 - 2 \left( \psi^{(0)}(\frac{5}{2}) + \log(2) + \mathcal{O}(\frac{1}{d}) \right) \log(d) + \psi^{(1)}(\frac{5}{2}) + \left( \psi^{(0)}(\frac{5}{2}) + \log(2) \right)^2 + \mathcal{O}(\frac{1}{d}) \right)$$

$$986 - \frac{1}{8} (2\log(s) + \log(d) + \gamma + \log(2) - 2 + \mathcal{O}(\frac{1}{d})) (3\log(d) + 3\gamma - 8 + \log(8) + \mathcal{O}(\frac{1}{d}))$$

$$987 + \frac{3}{4} (4\log(s)^2 + \log(d)^2 + 4\log(s)\log(d) + 2(\gamma + \log(2) - 2)\log(d)$$

$$988 + 4\log(s)(\gamma + \log(2) - 2) + (\gamma + \log(2) - 2)^2 + \mathcal{O}(d^{-1}))$$

$$989 = \log(d)^2 \left( \frac{3}{4} - \frac{3}{2} + \frac{3}{4} \right)$$

$$990 + \log(d) \left( -\frac{3}{2} \left( \psi^{(0)}(\frac{5}{2}) + \log(2) + \mathcal{O}(\frac{1}{d}) \right) - \frac{1}{2} (3\gamma - 8 + \log(8)) \right.$$

$$991 \left. - \frac{3}{2} ((2\log(s) + \gamma + \log(2) - 2) + \frac{3}{4} (4\log(s) + 2(\gamma + \log(2) - 2))) \right)$$

$$992 + 3\log(s)^2 - \log(s) (3\gamma - 8 + \log(8))$$

$$993 + \frac{3}{4} \left( \psi^{(1)}(\frac{5}{2}) + \left( \psi^{(0)}(\frac{5}{2}) + \log(2) \right)^2 \right) + 2(\gamma + \log(2) - 2)^2 + \mathcal{O}(d^{-1})$$

$$994 = \log(d) \left( -\frac{3}{2} \psi^{(0)}(\frac{5}{2}) - \frac{3}{2} \gamma - 3\log(2) + \mathcal{O}(\frac{1}{d}) \right) + 3\log(s)^2 - \log(s) (3\gamma + \log(8) - 8)$$

$$995 + \frac{3}{4} \left( \psi^{(1)}(\frac{5}{2}) + \left( \psi^{(0)}(\frac{5}{2}) + \log(2) \right)^2 \right) + 2(\gamma + \log(2) - 2)^2 + \mathcal{O}(d^{-1}).$$

996 With the identity  $\psi^{(0)}(\frac{5}{2}) = -2\log(2) - \gamma + \frac{8}{3}$  and

$$997 c_1 := -3\gamma - \log(8) + 8 \approx 4.189$$

$$998 c_2 := \frac{3}{4} (\psi^{(1)}(\frac{5}{2}) + (\psi^{(0)}(\frac{5}{2}) + \log(2))^2) + 2(\gamma + \log(2) - 2)^2 \approx 2.895, \quad (18)$$

999 we obtain

$$1000 \frac{1+2/d}{s^4} \mathbb{E}[f(|\langle \xi, x \rangle|)^2] = 3\log(s)^2 + c_1 \log(s) + c_2 + \mathcal{O}(d^{-1} \log d).$$

1026 **iv):** By (Hertrich, 2024, Thm 3), the transformed function has the form

$$1027 f(s) = \sum_{n=0}^{\infty} b_n x^n \quad \text{with} \quad b_n = \frac{\sqrt{\pi} \Gamma\left(\frac{n+d}{2}\right)}{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{n+1}{2}\right)} a_n.$$

1028 Moreover,

$$1029 f(s)^2 = \sum_{n=0}^{\infty} c_n s^n := \sum_{n=0}^{\infty} \left( \sum_{k=0}^n b_k b_{n-k} \right) s^n$$

1030 and, by (16),

$$1031 \mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}} [f(|\langle \xi, x \rangle|)^2] = \sum_{n=0}^{\infty} \frac{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n+d}{2}\right)} c_n \|x\|^n.$$

1032 Using that  $d$  is odd and applying the identities  $\Gamma(z+1) = \Gamma(z)z$  and  $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$ , we have

$$\begin{aligned} 1033 \tilde{c}_{k,n,d} &:= \frac{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi} \Gamma\left(\frac{n+d}{2}\right)} \frac{\sqrt{\pi} \Gamma\left(\frac{k+d}{2}\right)}{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{k+1}{2}\right)} \frac{\sqrt{\pi} \Gamma\left(\frac{n-k+d}{2}\right)}{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{n-k+1}{2}\right)} \\ 1034 &= \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n+1+d-1}{2}\right)} \frac{\Gamma\left(\frac{k+1+d-1}{2}\right)}{\Gamma\left(\frac{k+1}{2}\right)} \frac{\sqrt{\pi} \Gamma\left(\frac{n-k+1+d-1}{2}\right)}{\Gamma\left(\frac{d}{2}\right) \Gamma\left(\frac{n-k+1}{2}\right)} \\ 1035 &= \prod_{j=1}^{\frac{d-1}{2}} \frac{(k+2j-1)(n-k+2j-1)}{(n+2j-1)(2j-1)}. \end{aligned}$$

1036 Hence, we obtain

$$1037 \mathbb{V}_d[f](x) = \sum_{n=0}^{\infty} \left( \sum_{k=0}^n (\tilde{c}_{k,n,d} - 1) a_k a_{n-k} \right) \|x\|^n.$$

1038 We apply the identity  $4xy = (x+y)^2 - (x-y)^2$  to the numerator (with  $x = k+2j-1$ ,  $y = n-k+2j-1$ ) and denominator (with  $x = n+2j-1$ ,  $y = 2j-1$ ), and obtain

$$\begin{aligned} 1039 \tilde{c}_{k,n,d} &= \prod_{j=1}^{\frac{d-1}{2}} \frac{(k+2j-1)(n-k+2j-1)}{(n+2j-1)(2j-1)} = \prod_{j=1}^{\frac{d-1}{2}} \frac{(n+2(2j-1))^2 - (n-2k)^2}{(n+2(2j-1))^2 - n^2} \\ 1040 &= \prod_{j=1}^{\frac{d-1}{2}} \left( 1 + \frac{n^2 - (n-2k)^2}{(n+2(2j-1))^2 - n^2} \right) = \prod_{j=1}^{\frac{d-1}{2}} \left( 1 + \frac{k(n-k)}{(2j-1)(n+2j-1)} \right). \end{aligned}$$

1041 For  $F$  the Laplace kernel, assuming w.l.o.g. that  $\alpha = 1$ , we have  $a_n := (-1)^n/n!$ . Consequently,

$$\begin{aligned} 1042 \mathbb{V}_3[f](x) &= \sum_{n=0}^{\infty} \sum_{k=0}^n \frac{k(n-k)}{n+1} \frac{1}{k!(n-k)!} (-\|x\|)^n = \sum_{n=0}^{\infty} \frac{(-\|x\|)^n}{(n+1)!} \sum_{k=0}^n \binom{n}{k} (kn - k^2) \\ 1043 &= \sum_{n=0}^{\infty} \frac{(-\|x\|)^n}{(n+1)!} (n^2 2^{n-1} - n(n+1)2^{n-2}) = \frac{1}{4} \sum_{n=0}^{\infty} \frac{(-2\|x\|)^n}{(n+1)!} n(n-1) \\ 1044 &= \frac{1}{4} \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n+1)!} (n+1-2)n = \frac{1}{4} \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n-1)!} - \frac{1}{2} \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n+1)!} n. \quad (19) \end{aligned}$$

1045 For the first term, we obtain

$$1046 \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n-1)!} = -2\|x\| \sum_{n=1}^{\infty} \frac{(-2\|x\|)^n}{n!} = -2\|x\| (e^{-2\|x\|} - 1), \quad (20)$$

1047 and for the second

$$\begin{aligned} 1048 \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n+1)!} (n+1-1) &= \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{n!} - \sum_{n=2}^{\infty} \frac{(-2\|x\|)^n}{(n+1)!} \\ 1049 &= e^{-2\|x\|} - 1 + 2\|x\| + \frac{1}{2\|x\|} (e^{-2\|x\|} - 1 + 2\|x\| - 2\|x\|^2). \quad (21) \end{aligned}$$

1080 Plugging (20) and (21) into (19) finally yields

$$\begin{aligned}
1081 \mathbb{V}_3[f](x) &= \frac{-\|x\|}{2} (e^{-2\|x\|} - 1) - \frac{1}{2} \left( e^{-2\|x\|} - 1 + 2\|x\| + \frac{e^{-2\|x\|}}{2\|x\|} - \frac{1}{2\|x\|} + 1 - \|x\| \right) \\
1082 &= \frac{e^{-2\|x\|}}{4\|x\|} \left( e^{2\|x\|} - (2\|x\|^2 + 2\|x\| + 1) \right) \\
1083 &= \frac{1}{4\|x\|} \left( e^{2\|x\|} - (2\|x\|^2 + 2\|x\| + 1) \right) F(\|x\|)^2.
\end{aligned}$$

1084 For the Gauss kernel, we can follow exactly the same lines, after considering w.l.o.g.  $\sigma^2 = 1/2$  and  
1085 replacing  $\|x\|$  with  $\|x\|^2$ .

## 1093 D DISCUSSION OF THEOREM 1

### 1095 D.1 COMPARISON TO THE ERROR BOUNDS FROM HERTRICH (2024)

1096 In the following, we give a short summary how Theorem 1 improves the error bounds from Hertrich  
1097 (2024).

- 1098 - For positive definite kernels and Riesz kernels, the bound from Theorem 1 is dimension-  
1099 independent. In contrast, Hertrich (2024) only proves a dimension-independent absolute  
1100 error bound for the Gauss kernel (and conjectures that this is also true for the Laplace and  
1101 Matérn kernel). For the Riesz kernel, the error bound in Hertrich (2024) depends on the  
1102 dimension by  $\mathcal{O}(\sqrt{d})$ .
- 1103 - Theorem 1 bounds the error of the thin plate spline kernel. For this kernel Hertrich (2024)  
1104 does not provide a bound.
- 1105 - For kernels with analytic basis functions and for the Riesz kernel, we exactly compute  
1106 the variance  $\mathbb{V}_d[f]$ . In particular, Theorem 1 provides an *exact* calculation for the mean  
1107 square error, i.e., no tighter estimation is possible. E.g., for the Gauss and Laplace kernel,  
1108 this yields an improvement of the absolute error bounds given in Hertrich (2024), as we  
1109 observe that this bounds actually decay to zero for  $\|x\| \rightarrow \infty$ .

### 1113 D.2 THE VARIANCE OF GAUSS AND LAPLACE KERNEL FOR $d > 3$

1114 For  $d > 3$ ,  $d$  odd, we evaluated  $\mathbb{V}_d(f)$  for  $d = 5, \dots, 15$ , and make the following conjecture. Let  
1115  $T_n(f) = \sum_{k=0}^n \frac{f^{(k)}(0)}{k!} x^k$  denote the Taylor expansion of  $f$  in zero up to order  $n$ . We conjecture  
1116 that, for  $d \geq 3$ ,  $d$  odd, and the Laplace kernel  $F(x) = e^{-\alpha x}$ , it holds that

$$1117 \mathbb{V}_d[f](x) = \frac{e^{-2\alpha\|x\|}}{c_d(\alpha\|x\|)^{d-2}} \left( \sum_{i=0}^{\frac{d-3}{2}} (\alpha\|x\|)^{2i} (-1)^{\frac{d-5}{2}-i} c_{i,d} T_{d-2i}(e^{2\alpha\|x\|}) + \sum_{i=0}^{\frac{d-5}{2}} b_{i,d} (-1)^i (\alpha\|x\|)^{d+i} \right),$$

1118 where

$$\begin{aligned}
1119 c_d &= (12 + 4(d-5)) c_{d-2} \quad \text{for } d \geq 5 \quad \text{with } c_3 = 4, \\
1120 c_{0,d} &= (d-2)^2 (d-4) c_{0,d-2} \quad \text{for } d \geq 5 \quad \text{with } c_{0,3} = 1, \\
1121 b_{\frac{d-5}{2},d} &= 2b_{\frac{d-5}{2},d-2} \quad \text{for } d \geq 7 \quad \text{with } b_{0,5} = 4.
\end{aligned}$$

1122 For the remaining coefficients, we have not found a general simple rule. Our numerical computations  
1123 indicate that they are positive and that  $\mathbb{V}_d[f](x)\|x\|$  is monotonically increasing in  $\|x\|$ , with upper  
1124 bound  $s_d/\alpha$  and quadratic convergence to 0 for  $x \rightarrow 0$ . The upper bound  $s_d$  increases slowly in  $d$ ,  
1125 and our simulations hint that it might be bounded by a constant independently of  $d$ , see Figure 3.  
1126 For the Gauss kernel, the variance has the same form, except that  $\alpha^{-1}$  is replaced with  $\sigma/\sqrt{2}$  and  
1127  $\|x\|$  with  $\|x\|^2$  at all occurrences.

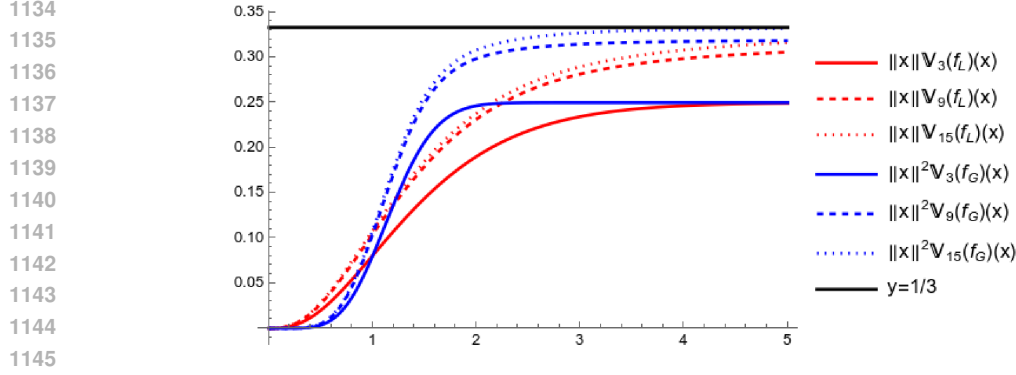


Figure 3: The scaled variance of the Laplace kernel  $f_L$  with  $\alpha = 1$  (in red) and the Gauss kernel  $f_G$  with  $\sigma = 1/\sqrt{2}$  (in blue) for dimension  $d = 3$  (solid lines),  $d = 9$  (dashed lines) and  $d = 15$  (dotted lines). The variance is multiplied times  $\|x\|$  (Laplace kernel) and  $\|x\|^2$  (Gauss kernel). We observe that the scaled variance increases monotonically in all cases, seemingly bounded from above by a constant (black solid line).

## E RANDOMIZATION OF A QMC SEQUENCE

Having a QMC sequence  $(\xi^P)_P$ , see Section 3.1, we can easily obtain an unbiased estimator in (4) by considering the QMC sequence  $(A\xi^P)_P$  with a uniformly chosen orthogonal matrix  $A \sim \mathcal{U}_{O(d)}$  because, according to (Ragozin, 1974, (2.3)), we have

$$\mathbb{E}_{A \sim \mathcal{U}_{O(d)}}[\psi(A\xi_p^P)] = \int_{O(d)} \psi(A\xi_p^P) d\mathcal{U}_{O(d)}(A) = \int_{\mathbb{S}^{d-1}} \psi(\eta) d\mathcal{U}_{\mathbb{S}^{d-1}}(\eta),$$

where  $\mathcal{U}_{O(d)}$  is the uniform distribution on the set  $O(d)$  of orthogonal  $d \times d$  matrices. Furthermore, note that  $\frac{1}{P} \sum_{i=1}^P f(\xi_i^P)$  converges to  $\int_{\mathbb{S}^{d-1}} f(\xi) d\mathcal{U}_{\mathbb{S}^{d-1}}(\xi)$  for  $P \rightarrow \infty$  and all continuous functions  $f$  because  $H^s(\mathbb{S}^{d-1})$  is dense in  $C(\mathbb{S}^{d-1})$  and the weights are all one, cf. (Atkinson, 1991, Sect 5.2).

## F PROOF OF THEOREM 3

**Gauss kernel:** A sufficient criterion for the function  $g_x$  being in the Sobolev space  $H^s(\mathbb{S}^{d-1})$  is that  $f(|\cdot|)$  is  $s$  times continuously differentiable. For  $F(t) = \exp(-\frac{t^2}{2\sigma^2})$ , the one-dimensional basis function is given by the confluent hypergeometric function  $f(t) = {}_1F_1(\frac{d}{2}; \frac{1}{2}; -\frac{t^2}{2\sigma^2})$ . We know by (Hertrich, 2024, Lem 6) that  $f = \mathcal{F}_1(g)$  with  $g(\omega) = \frac{d\pi\sigma \exp(-2\pi^2\sigma^2\omega^2)(2\pi^2\sigma^2\omega^2)^{(d-1)/2}}{\sqrt{2}\Gamma(\frac{d+2}{2})}$ , where  $\mathcal{F}_1$  denotes the one-dimensional Fourier transform. Then,  $g_k(\omega) := \omega^k g(\omega)$  is absolutely integrable for any  $k \in \mathbb{N}$  such that the differentiation/multiplication formula for the Fourier transform yields that the  $k$ -th derivative of  $f$  exists and is given by the continuous function  $f^{(k)} = (2\pi i)^k \mathcal{F}_1(g_k)$ .

**Riesz kernel:** With fixed  $x \in \mathbb{R}^d \setminus \{0\}$ , we have  $g_x(\xi) = \frac{\sqrt{\pi}\Gamma(\frac{d+r}{2})}{\Gamma(\frac{d}{2})\Gamma(\frac{r+1}{2})} |\langle \xi, x \rangle|^r$ . We set  $\eta := \frac{x}{\|x\|} \in \mathbb{S}^{d-1}$ , then  $x = \|x\|\eta$  and accordingly  $g_x(\xi) = \|x\|^r g_\eta(\xi)$ . The Sobolev norm (14) reads

$$\|g_x\|_{H^s(\mathbb{S}^{d-1})}^2 = \sum_{n=0}^{\infty} \sum_{k=1}^{N_{n,d}} \left(n + \frac{d-2}{2}\right)^{2s} \left\| \|x\|^r \int_{\mathbb{S}^{d-1}} g_\eta(\xi) \overline{Y_n^k(\xi)} d\xi \right\|^2.$$

The last integral can be computed with the help of the so-called  $\lambda$ -cosine transform. By (Rubin, 2014, (3.5) and (3.8)), the  $\lambda$ -cosine transform of a function  $h \in L^1(\mathbb{S}^{d-1})$  for  $\lambda > -1$  is defined by

$$\mathcal{C}^\lambda[h](\omega) = \frac{1}{|\mathbb{S}^{d-1}|} \frac{\sqrt{\pi}\Gamma(-\frac{\lambda}{2})}{\Gamma(\frac{d}{2})\Gamma(\frac{\lambda+1}{2})} \int_{\mathbb{S}^{d-1}} h(\theta) |\langle \omega, \theta \rangle|^\lambda d\theta,$$

for  $\omega \in \mathbb{S}^{d-1}$  and satisfies

$$\mathcal{C}^\lambda[Y_n^k](\omega) = Y_n^k(\omega) \begin{cases} (-1)^{\frac{n}{2}} \frac{\Gamma(\frac{n-\lambda}{2})}{\Gamma(\frac{n+d+\lambda}{2})}, & n \text{ even,} \\ 0, & n \text{ odd.} \end{cases}$$

Hence, we obtain that

$$\int_{\mathbb{S}^{d-1}} g_\eta(\xi) \overline{Y_n^k(\xi)} d\xi = \frac{|\mathbb{S}^{d-1}| \Gamma(\frac{d+r}{2})}{\Gamma(-\frac{r}{2})} \mathcal{C}^r[Y_n^k](\eta) = (-1)^{\frac{n}{2}} \frac{|\mathbb{S}^{d-1}| \Gamma(\frac{n-r}{2}) \Gamma(\frac{d+r}{2})}{\Gamma(\frac{n+d+r}{2}) \Gamma(-\frac{r}{2})} \overline{Y_n^k(\eta)}$$

if  $n$  is even and 0 if  $n$  is odd. The addition formula for spherical harmonics,

$$\sum_{k=1}^{N_{n,d}} |Y_n^k(\eta)|^2 = \frac{N_{n,d}}{|\mathbb{S}^{d-1}|},$$

see Atkinson & Han (2012), and the substitution  $n = 2m$  yield

$$\|g_x\|_{H^s(\mathbb{S}^{d-1})}^2 = \sum_{m=0}^{\infty} \left(2m + \frac{d-2}{2}\right)^{2s} \|x\|^{2r} \left(\frac{\Gamma(m - \frac{r}{2}) \Gamma(\frac{d+r}{2})}{\Gamma(m + \frac{d+r}{2}) \Gamma(-\frac{r}{2})}\right)^2 |\mathbb{S}^{d-1}| N_{2m,d}.$$

The asymptotic relations  $N_{n,d} \sim n^{d-2}$  from (13) and  $\Gamma(z+a)/\Gamma(z) \sim z^a$  for  $z \rightarrow \infty$ , see (NIST, 5.11.12), yield that the  $m$ -th summand behaves asymptotically like a multiple of  $(2m)^{2s-2r-2}$ , therefore the series converges if and only if  $s < r + \frac{1}{2}$ .

**Matérn kernel:** The one-dimensional basis function of the Matérn kernel is (Hertrich, 2024, Appx C.1)

$$f(t) = \underbrace{{}_1F_2\left(\frac{d}{2}; \frac{1}{2}, 1 - \nu; \frac{\nu t^2}{2\beta^2}\right)}_{=: f_1(t)} - \underbrace{|t|^{2\nu}}_{=: f_2(t)} \underbrace{\frac{\Gamma(1-\nu)\Gamma(\nu + \frac{d}{2}(2\nu)^\nu)}{\Gamma(\frac{d}{2})\Gamma(2\nu+1)\beta^{2\nu}} {}_1F_2\left(\nu + \frac{d}{2}; \nu + \frac{1}{2}, \nu + 1; \frac{\nu t^2}{2\beta^2}\right)}_{=: f_3(t)}.$$

The above hypergeometric functions  $f_1$  and  $f_3$ , whose parameters do not contain any nonpositive integers, are entire functions in  $\mathbb{R}$ . Hence, the corresponding spherical functions  $(g_1)_x$  and  $(g_3)_x$  from (10) are in  $C^\infty(\mathbb{S}^{d-1})$ , because they can be extended to smooth functions defined on a neighborhood of  $\mathbb{S}^{d-1}$ . Furthermore, since  $f_2$  is a multiple of the one-dimensional basis function of the Riesz kernel, we know from above that  $(g_2)_x \in H^s(\mathbb{S}^{d-1})$  for  $s < 2\nu + \frac{1}{2}$ . Hence, by (Quellmalz, 2020, Thm 5.2), we see that the product  $(g_2)_x(g_3)_x$  is also in  $H^s(\mathbb{S}^{d-1})$ . Note that the referenced theorem is only formulated for in integer  $s$ , but this can be easily extended using an interpolation argument (Quellmalz, 2020, Sect 5.3), because the multiplication with a smooth function constitutes a continuous operator both  $H^{\lfloor s \rfloor}(\mathbb{S}^{d-1}) \rightarrow H^{\lfloor s \rfloor}(\mathbb{S}^{d-1})$  and  $H^{\lfloor s \rfloor+1}(\mathbb{S}^{d-1}) \rightarrow H^{\lfloor s \rfloor+1}(\mathbb{S}^{d-1})$ .

## G BACKGROUND ON RFF AND RELATION WITH SLICING

We denote by  $\mathcal{M}_+(\mathbb{R}^d)$  the space of finite positive Borel measures on  $\mathbb{R}^d$ . Such measures can be identified with linear functional on the space  $C_0(\mathbb{R}^d)$  of continuous functions that vanish at infinity. The Fourier transform of measures is a linear operator defined by

$$\mathcal{F}_d: \mathcal{M}_+(\mathbb{R}^d) \rightarrow C_0(\mathbb{R}^d), \quad \mathcal{F}_d[\mu](x) = \int_{\mathbb{R}^d} e^{-2\pi i \langle x, v \rangle} d\mu(v),$$

cf. (Plonka et al., 2023, Sect 4.4). By Bochner's theorem, the Fourier transform is bijective from  $\mathcal{M}_+(\mathbb{R}^d)$  to the set of positive definite functions on  $\mathbb{R}^d$ , see ii) in Section 2 for the definition. If  $\mu$  is a probability measure, i.e.,  $\mu(\mathbb{R}^d) = 1$ , we have  $\mathcal{F}_d[\mu](0) = 1$ .

In the following, let  $F \circ \|\cdot\|$  be a positive definite function on  $\mathbb{R}^d$  with  $F(0) = 1$ .

**RFF:** Random Fourier features (RFF), see Rahimi & Recht (2007), use that, by Bochner's theorem,  $F(\|\cdot\|)$  is the Fourier transform of a probability measure  $\mu \in \mathcal{M}_+(\mathbb{R}^d)$ , i.e.,

$$F(\|x\|) = \mathcal{F}_d[\mu](x) = \mathbb{E}_{v \sim \mu} [\exp(2\pi i \langle x, v \rangle)] \approx \frac{1}{D} \sum_{p=1}^D \exp(2\pi i \langle x, v_p \rangle), \quad (22)$$

where we sample  $D \in \mathbb{N}$  instances  $v_p$  iid from  $\mu$ . Using that  $F$  is real-valued, we can replace the exponential by the cosine in (22). By (Altekrüger et al., 2023, Prop C.2), the radial measure  $\mu$  can be decomposed as follows. We define  $\iota: \mathbb{R}^d \rightarrow [0, \infty)$ ,  $\iota(x) = \|x\|$  and its pushforward measure  $\tilde{\mu} := \iota_*\mu = \mu \circ \iota^{-1} \in \mathcal{M}_+([0, \infty))$ , then we have

$$\mu = T(\tilde{\mu} \otimes \mathcal{U}_{\mathbb{S}^{d-1}}), \quad \text{where } T(r, \xi) = r\xi.$$

Hence, sampling from  $\mu$  can be realized by  $v_p = r_p \xi_p$ , where  $\xi_p \sim \mathcal{U}_{\mathbb{S}^{d-1}}$  and  $r_p \sim \tilde{\mu}$ . Then (22) becomes

$$F(\|x\|) \approx \frac{1}{D} \sum_{p=1}^D \cos(2\pi r_p \langle x, \xi_p \rangle). \quad (23)$$

**RFF Summation:** For approximating the kernel sum (1), we insert the RFF (22) and obtain

$$s_m \approx \sum_{n=1}^N w_n \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m - x_n, v_p \rangle} = \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m, v_p \rangle} \sum_{n=1}^N w_n e^{-2\pi i \langle x_n, v_p \rangle}, \quad (24)$$

where  $v_1, \dots, v_D$  are iid samples of  $\mu$ . Since the inner sum over  $n$  is independent of  $m$  and therefore has to be evaluated only  $D$  times, the total computational complexity of computing (24) for all  $m = 1, \dots, M$  is  $\mathcal{O}(D(N + M))$ .

**Slicing:** The slicing approach uses the approximation (4), i.e.,

$$F(\|x\|) \approx \frac{1}{P} \sum_{p=1}^P f(|\langle x, \xi_p \rangle|),$$

where  $\xi_p \sim \mathcal{U}_{\mathbb{S}^{d-1}}$ . By (Rux et al., 2024, Cor 4.11), the function  $f(|\cdot|)$  is positive definite and hence possesses a Fourier transform, which is a probability measure on  $\mathbb{R}$  because  $f(0) = F(0) = 1$ . Applying RFF with  $Q$  points to the one-dimensional function  $f(|\cdot|)$ , we obtain

$$F(\|x\|) \approx \frac{1}{PQ} \sum_{p=1}^P \sum_{q=1}^Q \cos(2\pi r_{p,q} \langle x, \xi_p \rangle), \quad (25)$$

where  $r_{p,q} \sim \mathcal{F}_1^{-1}[f(|\cdot|)]$ . According to (Rux et al., 2024, Cor 4.11), we have

$$\mathcal{F}_1^{-1}[f(|\cdot|)] = \mathcal{A}_1^* \iota_* \mu = \mathcal{A}_1^* \tilde{\mu},$$

where  $\mathcal{A}_1^*: \mathcal{M}_+([0, \infty)) \rightarrow \mathcal{M}_+(\mathbb{R})$  is the symmetrization operator that extends a measure on  $[0, \infty)$  to an even measure on the real line and is defined for any  $\nu \in \mathcal{M}_+([0, \infty))$  and  $g \in C_0(\mathbb{R})$  by  $\langle \mathcal{A}_1^* \nu, g \rangle = \langle \nu, g + g(-\cdot) \rangle$ . Since the right-hand side of (25) is independent of the sign of  $r_{p,q}$ , it stays the same when we sample  $r_{p,q}$  from  $\tilde{\mu}$  instead of  $\mathcal{A}_1^* \tilde{\mu}$ . Therefore, if  $Q = 1$ , we see that the right-hand side of (25) is the same as the right-hand side of (23). In particular, RFF can be viewed as a special case of slicing. The respective slicing summation is given by (5) combined with the one-dimensional summation in Appendix I below.

## H RELATION BETWEEN DISTANCE QMC DESIGNS AND ORTHOGONAL POINTS

Consider the QMC design  $\xi^P$  that is a minimizer (9) for  $s = \frac{d}{2}$ . To improve the error, Womersley (2018) suggested symmetric QMC designs, meaning that for every point  $\xi$  it contains also the antipodal point  $-\xi$ . However, since the function (10) we want to integrate is symmetric, i.e.,  $g_x(\xi) = g_x(-\xi)$ , we can discard one of the antipodal points  $\xi$  and  $-\xi$  and get the same result. Therefore, we minimize the functional

$$\mathcal{E}_{\text{sym}}(\xi^P) = \mathcal{E}((\xi^P, -\xi^P)) = -2 \sum_{p,q=1}^P (\|\xi_p^P - \xi_q^P\| + \|\xi_p^P + \xi_q^P\|).$$



Our numerical trials indicate that indeed the minimizers of  $\mathcal{E}_{\text{sym}}$  yield a smaller integration error than the minimizers of  $\mathcal{E}$ . Using  $\|\xi_p\| = 1$ , we see that

$$\begin{aligned} -(\|\xi_p^P - \xi_q^P\| + \|\xi_p^P + \xi_q^P\|)^2 &= -\left(\sqrt{2 - 2\langle \xi_p^P, \xi_q^P \rangle} + \sqrt{2 + 2\langle \xi_p^P, \xi_q^P \rangle}\right)^2 \\ &= -4 - 2\sqrt{4 - 4\langle \xi_p^P, \xi_q^P \rangle}, \end{aligned}$$

attains its minimum if and only if  $\langle \xi_p^P, \xi_q^P \rangle = 0$ . Hence, if  $P \leq d$  and  $\xi^P$  is an orthonormal system in  $\mathbb{R}^d$ , then  $\xi^P$  a minimizer of  $\mathcal{E}_{\text{sym}}$ . However, this argumentation does not work if  $P > d$ , as we can only choose  $d$  orthogonal vectors.

## I BACKGROUND ON ONE-DIMENSIONAL FAST FOURIER SUMMATION

In this section, we review literature about one-dimensional fast Fourier summation used in Section 4.3 and specify the parameters used in our numerical examples. Fast Fourier summations were proposed in Kunis et al. (2006); Potts et al. (2004) based on the non-equispaced fast Fourier transform (Beylkin, 1995; Dutt & Rokhlin, 1993), which is implemented in several libraries (Knopp et al., 2023; Keiner et al., 2009; Shih et al., 2021). In our numerical examples, we use the Julia library Knopp et al. (2023). Here, we follow a similar workflow as in Hertrich (2024).

Let  $x_1, \dots, x_N \in \mathbb{R}^d$ ,  $y_1, \dots, y_M \in \mathbb{R}^d$ ,  $w_1, \dots, w_N \in \mathbb{R}$  and  $k(x, y) = f(|x - y|) = g(x - y)$ . We want to compute for  $\xi \in \mathbb{S}^{d-1}$ ,  $x_{n,\xi} := \langle x_n, \xi \rangle$  and  $y_{m,\xi} := \langle y_m, \xi \rangle$  the one-dimensional kernel sums

$$t_m = \sum_{n=1}^N w_n k(x_{n,\xi}, y_{m,\xi}) = \sum_{n=1}^N w_n g(x_{n,\xi} - y_{m,\xi}).$$

**Step 1: Rescaling** For Step 2 and 3, we will need two properties. First, since we will use discrete (fast) Fourier transforms, we require that  $x_{n,\xi}, y_{m,\xi}, x_{n,\xi} - y_{m,\xi} \in [-\frac{1}{2}, \frac{1}{2}]$ . For the important example of positive definite kernels, which decay to zero, we often can derive explicit formulas for the Fourier transform of  $g$  via Bochner’s integral and Rux et al. (2024), see Table 2 for some examples. In order to use this explicit formula for the fast Fourier summation, we will additionally require that  $g(x) \approx 0$  for  $|x| > \frac{1}{2}$ .

Both properties can be achieved by rescaling the problem. More precisely, let  $T < 0.5$  and  $c = \max_{n=1, \dots, N} \|x_n\| + \max_{m=1, \dots, M} \|y_m\|$ . For the case of decaying positive definite kernels, assume that  $g(x) \approx 0$  for  $|x| > g_{\max}$ . Then, it holds that

$$t_m = \sum_{n=1}^N w_n g(x_{n,\xi} - y_{m,\xi}) = \sum_{n=1}^N w_n \tilde{g}(\tilde{x}_{n,\xi} - \tilde{y}_{m,\xi}),$$

with

$$\tilde{g}(x) := g\left(\frac{x}{\tau}\right), \quad \tilde{x}_{n,\xi} := \tau x_{n,\xi} = \langle \tau x_n, \xi \rangle, \quad \tilde{y}_{m,\xi} := \tau y_{m,\xi} = \langle \tau y_m, \xi \rangle,$$

where the constant  $\tau := \min\left\{\frac{T}{c}, \frac{1}{2g_{\max}}\right\}$  does not depend on  $\xi$ . Then, by definition, the two properties from above are fulfilled. For the rest of the section, we will denote the rescaled points  $\tilde{x}_{n,\xi}, \tilde{y}_{m,\xi}$  and the rescaled kernel function  $\tilde{g}$  again by  $x_{n,\xi}, y_{m,\xi}$  and  $g$ .

In the case of the Gauss, Laplace or Matérn kernels, the rescaled kernel  $\tilde{k}(x, y) = \tilde{g}(x - y)$  is again a Gauss, Laplace or Matérn kernel with the altered parameter  $\tilde{\sigma} = \tau\sigma$ ,  $\tilde{\alpha} = \alpha/\tau$  or  $\tilde{\beta} = \tau\beta$ . In our numerics, we set  $g_{\max} = 5m$  with  $m = \sigma = \beta = \frac{1}{\alpha}$  for the Gauss, Laplace and Matérn kernel. Moreover, we set the threshold  $T$  to 0.3 for the Gauss kernel, to 0.2 for the Matérn kernel and to 0.1 for the Laplace kernel.

**Step 2: Computation of the Fourier Coefficients of the Kernel** In the next step, we expand  $g$  into its Fourier series on  $[-\frac{1}{2}, \frac{1}{2}]$  and truncate it by

$$g(x) = \sum_{k \in \mathbb{Z}} c_k(g) e^{2\pi i k x} \approx \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{2\pi i k x}, \quad c_k(g) = \int_{-\frac{1}{2}}^{\frac{1}{2}} g(x) e^{-2\pi i k x} dx$$

with some even  $N_{\text{ft}} \in \mathbb{N}$ . To compute the Fourier coefficients  $c_k(g)$ , we employ that  $g(x) \approx 0$  for  $|x| > \frac{1}{2}$  such that Poisson’s summation formula (see, e.g., Plonka et al., 2023, Thm 2.28) implies

$$c_k(g) \approx c_k \left( \sum_{l \in \mathbb{Z}} g(\cdot + l) \right) = \mathcal{F}_1[g](k),$$

where  $\mathcal{F}_1[g](\omega)$  is the Fourier transform (11).

In our experiments, we choose  $N_{\text{ft}} = 128$  for the Gauss,  $N_{\text{ft}} = 512$  for the Matérn and  $N_{\text{ft}} = 1024$  for the Laplace kernel. Note that the coefficients  $c_k(g)$  do not depend on the input points and need to be computed only once for different choices of  $\xi$ . The function  $\mathcal{F}_1[g]$  is analytically given for the Gauss, Laplace and Matérn kernel in Table 2.

**Step 3: Fast Fourier Summation** Finally, we use this expansion to compute the kernel sums

$$t_m \approx \sum_{n=1}^N \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} w_n c_k(g) e^{2\pi i k (y_{m,\xi} - x_{n,\xi})} = \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{2\pi i k y_{m,\xi}} \underbrace{\sum_{n=1}^N w_n e^{-2\pi i k x_{n,\xi}}}_{=: \hat{w}_k}. \quad (26)$$

The computation of the second sum  $\hat{w}_k$  is the adjoint discrete Fourier transform of the vector  $w = (w_1, \dots, w_N)$  at the non-equispaced knots  $(-x_{1,\xi}, \dots, -x_{N,\xi})$ . Afterward, the computation of the vector  $t = (t_1, \dots, t_M)$  is the Fourier transform of the vector  $(c_k(g) \hat{w}_k)_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1}$  at the non-equispaced knots  $(-y_{1,\xi}, \dots, -y_{M,\xi})$ . These Fourier transforms can be computed by the NFFT in time complexity  $\mathcal{O}(N + N_{\text{ft}} \log N_{\text{ft}})$  and  $\mathcal{O}(M + N_{\text{ft}} \log N_{\text{ft}})$  leading to an overall complexity of  $\mathcal{O}(N + M + N_{\text{ft}} \log N_{\text{ft}})$ .

## J COMPUTATIONAL COMPLEXITY

We consider the computational complexity of the random Fourier features (RFF) and the Fourier slicing summation. We denote the RFF summation (24) with  $D$  features by

$$\tilde{s}_m^{\text{RFF}} := \sum_{n=1}^N w_n \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m - x_n, v_p \rangle} = \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m, v_p \rangle} \sum_{n=1}^N w_n e^{-2\pi i \langle x_n, v_p \rangle}. \quad (27)$$

Computing  $\tilde{s}_m^{\text{RFF}}$  for all  $m = 1, \dots, M$  has a complexity of  $\mathcal{O}(D(N + M))$ .

We consider the Fourier slicing as described in Appendix I, with the slight modification that instead of rescaling the problem, we take a  $2T$ -periodic Fourier series. To this end, we choose

$$R \geq \max_{n=1, \dots, N} \|x_n\| + \max_{m=1, \dots, M} \|y_m\|. \quad (28)$$

Then, in the one-dimensional sum in (5), it holds that  $\langle x_n - y_m, \xi \rangle \leq R$  for all  $\xi \in \mathbb{S}^{d-1}$ . Therefore, we can replace  $f(|\cdot|)$  in (5) by any function  $g: \mathbb{R} \rightarrow \mathbb{R}$  that satisfies  $g(t) = f(|t|)$  for all  $|t| \leq R$  without changing the sum. In particular, we choose  $g$  as a sufficiently smooth  $2T$ -periodic function for some  $T > R$ , in order to achieve a convergent Fourier series of  $g$ . With this, we insert the one-dimensional summation (26) with a  $2T$ -periodic function  $g$  into the sliced kernel sum (5) and obtain

$$\tilde{s}_m^{\text{FS}} := \frac{1}{P} \sum_{p=1}^P \sum_{n=1}^N \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} w_n c_k(g) e^{\pi i k \langle y_m - x_n, \xi_p \rangle / T} \quad (29)$$

with the Fourier coefficients  $c_k(g) = (2T)^{-1} \int_{-T}^T g(t) \exp(-\pi i k t / T) dt$ . As already noted in Appendix I, we interchange the sums in order to achieve a fast summation

$$\tilde{s}_m^{\text{FS}} = \frac{1}{P} \sum_{p=1}^P \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{\pi i k \langle y_m, \xi_p \rangle / T} \sum_{n=1}^N w_n e^{-\pi i k \langle x_n, \xi_p \rangle / T}, \quad m = 1, \dots, M.$$

Utilizing the NFFT for the sums over  $n$  and  $k$ , this has a complexity of

$$\mathcal{O}(P(N + M + N_{\text{ft}} \log N_{\text{ft}}))$$

with  $P$  the number of slices and  $N_{\text{ft}}$  the number of Fourier coefficients.

The asymptotic complexities with respect to  $M, N$  of the two methods depend on different parameters, most notably the number  $D$  of random features for RFF and the number  $P$  of points/slices in the QMC sequence for Fourier slicing, which need to be chosen depending on the desired accuracy. The following proposition compares their complexity to achieve a relative error

$$\frac{1}{\|w\|_1} \mathbb{E} [\|s - \tilde{s}^{\text{RFF}}\|_\infty] \in \mathcal{O}(\varepsilon) \quad \text{and} \quad \frac{1}{\|w\|_1} \|s - \tilde{s}^{\text{FS}}\|_\infty \in \mathcal{O}(\varepsilon) \quad \text{for } \varepsilon \downarrow 0,$$

where  $\|w\|_1 := \sum_{n=1}^N |w_n|$  and  $\|s\|_\infty := \max_{m=1, \dots, M} |s_m|$ , and, in case of RFF, the expectation is taken with respect to the random  $v_p \sim \mu$ .

**Proposition 5.** *Let  $F(t) = \exp(-\frac{t^2}{2\sigma^2})$  be the Gauss kernel with  $\sigma > 0$  and sliced kernel  $f$ . Assume there exists  $R > 0$  satisfying (28) independently of  $N$  and  $M$ .*

- *The computation of the RFF sum  $\tilde{s}^{\text{RFF}} = (\tilde{s}_m^{\text{RFF}})_{m=1}^M$ , see (27), with the parameter  $D \sim \varepsilon^2 |\log \varepsilon|$  for  $\varepsilon \downarrow 0$  achieves the relative error  $\mathbb{E}[\|s - \tilde{s}\|_\infty] / \|w\|_1 \in \mathcal{O}(\varepsilon)$  and the numerical complexity*

$$\mathcal{O}((N + M) \varepsilon^{-2} |\log \varepsilon|).$$

- *For  $q \in \mathbb{N}$  arbitrary, let  $g \in C^q(\mathbb{R}^d)$  be a  $2T$ -periodic function with  $T > R$  that satisfies  $g(t) = f(|t|)$  for  $|t| \leq R$ . Furthermore, assume that the slicing error has rate  $r > 0$ , i.e.,*

$$\sup_{x \in \mathbb{R}^d} \left| F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \right| \in \mathcal{O}(P^{-r}).$$

*Then the computation of the Fourier slicing sum  $\tilde{s}^{\text{FS}} = (\tilde{s}_m^{\text{FS}})_{m=1}^M$ , see (29), with the parameters  $P \sim \varepsilon^{-1/r}$  and  $N_{\text{ft}} \sim \varepsilon^{-1/q} |\log \varepsilon|^{-1}$  for  $\varepsilon \downarrow 0$  achieves the relative error  $\|s - \tilde{s}^{\text{FS}}\|_\infty / \|w\|_1 \in \mathcal{O}(\varepsilon)$  and the numerical complexity*

$$\mathcal{O}(\varepsilon^{-1/r} (\varepsilon^{-1/q} + N + M)).$$

The assumption that the slicing error is bounded with rate  $r$  is fulfilled with  $r = \frac{d}{2(d-1)} > \frac{1}{2}$  for the distance QMC designs, which minimize (9) with  $s = d/2$ , as we proved in Corollary 4. However, our numerical results suggest that  $r$  might be even larger, cf. Table 1. Because  $q \in \mathbb{N}$  can be chosen arbitrarily and  $r > 1/2$ , the asymptotic complexity of the slicing summation is lower than for RFF. Furthermore, we note that the second part of the proposition holds for any kernel function  $F$  for which  $t \mapsto f(|t|)$  is in  $C^q(\mathbb{R})$ .

*Proof. RFF:* By (Sutherland & Schneider, 2015, Prop 3), we have for the Gauss kernel

$$\begin{aligned} \mathbb{E}_{v_1, \dots, v_D \sim \mu} \left[ \sup_{x \in B_T} \left| F(\|\cdot\|) - \frac{1}{D} \sum_{p=1}^D e^{-2\pi i \langle \cdot, v_p \rangle} \right| \right] \\ \leq \frac{24\sqrt{dT}(e^{-1/2} + \sqrt{d} + \sqrt{2\log(D)})}{\sigma\sqrt{D}} \in \mathcal{O}\left(\sqrt{\frac{\log(D)}{D}}\right), \end{aligned}$$

where  $B_T$  the ball in  $\mathbb{R}^d$  of radius  $T > 0$ . Hence, the error of the RFF summation is bounded by

$$\begin{aligned} \mathbb{E}[\|s - \tilde{s}^{\text{RFF}}\|_\infty] &\leq \sum_{n=1}^N |w_n| \mathbb{E} \left[ \max_{m=1, \dots, M} \left| F(\|x_n - y_m\|) - \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m - x_n, v_p \rangle} \right| \right] \\ &\leq \|w\|_1 \mathbb{E} \left[ \max_{n=1, \dots, N} \left| F(\|x_n - y_m\|) - \frac{1}{D} \sum_{p=1}^D e^{2\pi i \langle y_m - x_n, v_p \rangle} \right| \right] \in \mathcal{O}\left(\sqrt{\frac{\log(D)}{D}}\right). \end{aligned}$$

For the desired accuracy  $\varepsilon$ , we choose  $D \sim -\varepsilon^{-2} \log(\varepsilon)$ . Then  $\frac{\log(D)}{D} = \varepsilon^2 \frac{\log(-\log(\varepsilon)) - 2\log(\varepsilon)}{-\log(\varepsilon)} \in \mathcal{O}(\varepsilon^2)$ , so that we obtain a relative error  $\mathcal{O}(\varepsilon)$  with a complexity of

$$\mathcal{O}(D(N + M)) = \mathcal{O}((N + M) \varepsilon^{-2} |\log(\varepsilon)|).$$

**Slicing:** For the Gauss kernel, the sliced kernel  $t \mapsto f(|t|)$  is an analytic function, see Table 2. Therefore,  $g$  can be constructed via two-point Taylor approximation, see Potts & Steidl (2003). We estimate the error

$$\|s - \tilde{s}^{\text{FS}}\|_{\infty} \leq \|w\|_1 \max_{\substack{n=1,\dots,N \\ m=1,\dots,M}} \left| F(\|x_n - y_m\|) - \frac{1}{P} \sum_{p=1}^P \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{\pi i k \langle y_m - x_n, \xi_p \rangle / T} \right|.$$

Since,  $R > \|x_n - y_m\|$  for all  $n$  and  $m$  by (28), we have

$$\begin{aligned} \frac{\|s - \tilde{s}^{\text{FS}}\|_{\infty}}{\|w\|_1} &\leq \sup_{\|x\| \leq R} \left| F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \right| \\ &\quad + \sup_{\|x\| \leq R} \left| \frac{1}{P} \sum_{p=1}^P \left( f(|\langle \xi_p, x \rangle|) - \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{\pi i k \langle x, \xi_p \rangle / T} \right) \right|. \end{aligned}$$

Since  $g(t) = f(|t|)$  for all  $|t| \leq R$ , we have

$$\frac{\|s - \tilde{s}^{\text{FS}}\|_{\infty}}{\|w\|_1} \leq \sup_{x \in \mathbb{R}^d} \left| F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \right| + \sup_{|t| \leq R} \left| g(t) - \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{\pi i k t / T} \right|.$$

For the first term, we use our assumption that

$$\sup_{x \in \mathbb{R}^d} \left| F(\|x\|) - \frac{1}{P} \sum_{p=1}^P f(|\langle \xi_p, x \rangle|) \right| \in \mathcal{O}(P^{-r})$$

with some  $r > 0$ . For the second term, Bernstein's theorem for the convergence of Fourier series implies that there exists  $C_g$  such that

$$\sup_{t \in \mathbb{R}} \left| g(t) - \sum_{k=-N_{\text{ft}}/2}^{N_{\text{ft}}/2-1} c_k(g) e^{\pi i k t / T} \right| \leq C_g N_{\text{ft}}^{-q} \log(N_{\text{ft}}),$$

cf. (Plonka et al., 2023, Thm 1.39). Combining these estimates, we obtain

$$\frac{\|s - \tilde{s}^{\text{FS}}\|_{\infty}}{\|w\|_1} \in \mathcal{O}(P^{-r} + N_{\text{ft}}^{-q} \log(N_{\text{ft}})).$$

Choosing  $P \sim \varepsilon^{-1/r}$  and  $N_{\text{ft}} \sim \varepsilon^{-1/q} / \log(\varepsilon^{-1})$  for  $\varepsilon \downarrow 0$ , we see that

$$P^{-r} + N_{\text{ft}}^{-q} \log(N_{\text{ft}}) \sim \varepsilon + \frac{\varepsilon(q^{-1} \log(\varepsilon^{-1}) - \log(\log(\varepsilon^{-1})))}{\log(\varepsilon^{-1})^q} \in \mathcal{O}(\varepsilon).$$

Hence, the complexity of the Fourier slicing to achieve relative error  $\varepsilon$  is

$$\mathcal{O}(P(N + M + N_{\text{ft}} \log N_{\text{ft}})) = \mathcal{O}(\varepsilon^{-1/r} (\varepsilon^{-1/q} + N + M)). \quad \square$$

## K ADDITIONAL NUMERICAL RESULTS

### K.1 ADDITIONAL PLOTS AND TABLES FOR SECTION 4.2

In the following, we redo the experiments from Figure 1 and Table 1 and vary some parameters. More precisely, we redo it for the negative distance kernel, choose other length scale parameters of the kernel and perform it in higher dimensions ( $d = 200$ ).

**Negative Distance Kernel** We do the same experiment as in Figure 1 and Table 1 with the negative distance kernel. The results are given in Figure 4 and Table 3. We can see that the advantage of QMC slicing is not as large as for smooth kernels, which is expected considering the theoretical results from Section 3. In particular, the spherical function (10) is not in  $H^{d/2}(\mathbb{S}^{d-1})$  if  $d \geq 3/2$ , so the

assumptions of the bound (8) are not fulfilled. Nevertheless, QMC slicing is still significantly more accurate than non-QMC slicing. Note that RFF based methods are not available for the negative distance kernel since it is not positive definite and therefore Bochner’s theorem does not apply.

**Other Length Scales of the Kernels** We redo the experiment from Figure 1 and Table 1 with scale factors  $s = \frac{1}{2}$  and  $s = 2$  of the kernel parameter. The results are given in Figure 5 and Table 4 for  $s = \frac{1}{2}$  and in Figure 6 and Table 5 for  $s = 2$ . We observe that the advantage of QMC is more significant of for larger scale factors, which is expected since the function  $\xi \mapsto f(|\langle \xi, x \rangle|)$  is more regular for larger  $s$  than for smaller  $s$ .

**Higher Dimensions** Finally, we do the same experiment as in Figure 1 and Table 1 for the higher dimension  $d = 200$ . Here, we use the negative distance kernel, the Matérn kernel with  $\nu = 3 + \frac{1}{2}$  and the Gauss kernel, where the parameters are chosen by the median rule with scale factor  $\gamma = 1$ . The results are given in Figure 7 and Table 6. The advantage of QMC is less pronounced in such high dimensions, but still visible.

## K.2 ADDITIONAL RESULTS FOR SECTION 4.3

**Negative Distance Kernel** We redo the experiment from Section 4.3 for the negative distance kernel and the thin-plate spline kernel. The results are given in Figure 8. We can see that QMC Fourier slicing outperforms standard slicing clearly in all cases. Note that RFF based methods are not available for these kernels since they are not positive definite and Bochner’s theorem does not apply.

**Higher Dimensions** We run the same experiment as in Section 4.3 on the MNIST and FashionMNIST dataset without dimension reduction and therefore  $d = 784$ . The results are given in Figure 9. We can see that the advantage of QMC slicing is smaller than for the lower-dimensional examples but still clearly visible for some kernels. In accordance with the considerations of Appendix H, the advantage comes in when  $P > d$ .

**GPU Comparison** We want to demonstrate the advantage of our method in a GPU-comparison with a large number of data points. As a test dataset we concatenate the MNIST and FashionMNIST in all eight orientations arising rotating and mirroring the images and reduce the dimension via PCA to  $d = 30$ . The arising dataset has  $N = M = 960000$  entries. Then, we compare RFF, ORF, QMC (Sobol) RFF, Slicing and QMC Slicing, where the QMC directions for slicing are chosen by minimizing the distance functional, see Section 4.1. This experiment is implemented in Python using PyTorch and we use brute-force kernel summation by the PyKeOps library (Charlier et al., 2021) as a baseline. The results are given in Figure 10. Even though the RFF-based methods parallelize a bit better on the GPU than the fast Fourier summations, the conclusions are mainly the same as for the CPU experiments. We can clearly see the advantage of QMC slicing over the comparisons. Particularly, for non-smooth kernels, slicing-based methods work much better than RFF-based methods.

## K.3 MMD GRADIENT FLOWS

Finally, we use our fast summation method in a specific application. Here, we consider gradient flows of the maximum mean discrepancy (MMD), which have been considered in several papers for generative modeling and other applications, see, e.g., Arbel et al. (2019); Chen et al. (2024); Galashov et al. (2024); Hertrich et al. (2024); Lim et al. (2024).

**Background** For a dataset  $\mathbf{y} = (y_1, \dots, y_M)$  of target particles, we consider the discrete MMD functional  $F_{\mathbf{y}}: (\mathbb{R}^d)^N \rightarrow \mathbb{R}$  defined by

$$G_{\mathbf{y}}(\mathbf{x}) = \text{MMD}_K(\mathbf{x}, \mathbf{y}) = \frac{1}{2N^2} \sum_{i,j=1}^N K(x_i, x_j) - \frac{1}{MN} \sum_{i,j=1}^{N,M} K(x_i, y_j) + \frac{1}{2M^2} \sum_{i,j=1}^M K(y_i, y_j).$$

We note that (under suitable assumptions on the kernel), the MMD is a metric on the space of probability distributions on  $\mathbb{R}^d$  such that  $F_{\mathbf{y}}(\mathbf{x})$  is always non-negative and zero if and only if the particles  $\mathbf{x}$  and  $\mathbf{y}$  coincide. Here, we interpret  $\mathbf{x}$  as the discrete probability measure  $\frac{1}{N} \sum_{n=1}^N \delta_{x_n}$ .

Now, we minimize  $G_{\mathbf{y}}$  by simulating the gradient flow

$$\dot{\mathbf{x}} = -\nabla G_{\mathbf{y}}(\mathbf{x}), \quad \mathbf{x}(0) = \mathbf{x}^{(0)},$$

starting with some random initial samples  $\mathbf{x}^{(0)} \in (\mathbb{R}^d)^N$  using the explicit Euler discretization

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau \nabla G_{\mathbf{y}}(\mathbf{x}^{(k)}), \quad (30)$$

where the gradient  $\nabla G_{\mathbf{y}}$  is computed via backpropagation. Throughout the flow, the particles  $\mathbf{x}$  will then move towards the target distribution  $\mathbf{y}$ .

**Experimental Setup** We choose the negative distance kernel  $K(x, y) = -\|x - y\|$  and consider the CIFAR10 dataset ( $M = 50000$  and  $d = 3072$ ) as target points  $\mathbf{y}$ . For  $\mathbf{x}$ , we choose  $N = M = 50000$  samples. Then, we run the (discretized) gradient flow from (30), with initial particles  $\mathbf{x}^{(0)}$  drawn iid from a standard normal distribution. For speeding up the convergence, we follow Hertrich et al. (2024); Lim et al. (2024) and add a momentum parameter  $m = 0.9$ . That is, we modify the equation (30) to

$$\begin{aligned} \mathbf{v}^{(k+1)} &= \nabla G_{\mathbf{y}}(\mathbf{x}^{(k)}) + m\mathbf{v}^{(k)} \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \tau\mathbf{v}^{(k+1)} \end{aligned}$$

with initial value  $\mathbf{v}^{(0)} = 0$ . We run the MMD flow for 50000 steps with step size  $\tau = 1$ , where we compute the function  $G_{\mathbf{y}}$  by (Monte Carlo) Slicing and by QMC Slicing with  $P = 1000$  projections.

**Remark 6.** *We would like to point out that running the gradient flow with exact computation of  $G_{\mathbf{y}}$  is computationally intractable. Using (QMC-)Slicing, one iteration (30) takes between 0.2 and 0.3 seconds on an NVIDIA RTX 4090 GPU. On the other hand, the exact gradient evaluation via PyKeOps takes about one hour. Considering that we are running the flow for 50000 steps, this underlines the need of (QMC-)Slicing.*

**Results** We plot the objective value of  $G_{\mathbf{y}}(\mathbf{x})$  versus the computation time in Figure 11. We observe that the smaller error in the gradient evaluation by QMC Slicing significantly improves the convergence behavior.

#### K.4 ON THE GAP BETWEEN THEORETICAL GUARANTEES AND NUMERICAL RESULTS

In our numerical part, we observe significantly better error rates than we can prove theoretically. One possible explanation is the following. Our theoretical guarantees for QMC Slicing are based on worst-case errors in Sobolev spaces on the sphere and consequently rely on the smoothness of the functions  $g_x$  from (10), which is not satisfied for some kernels in Theorem 3. However, these results are only worst-case error rates that do not account for the specific properties of  $g_x$ . First, the function  $g_x(\xi)$  depends only on  $\langle x, \xi \rangle$ , thus having a lower effective dimension. Furthermore, by construction  $g_x(\xi) = g_x(-\xi)$  such that for all  $x \in \mathbb{R}^d$  it holds that  $\mathbb{E}_{\xi \sim \mathcal{U}_{\mathbb{S}^{d-1}}}[g_x(\xi)] = \mathbb{E}_{\xi \sim \mathcal{U}_{\{\zeta \in \mathbb{S}^{d-1}, \langle \zeta, x \rangle > 0\}}}[g_x(\xi)]$ . Moreover,  $g_x$  is infinitely often differentiable for on the hemisphere  $\{\zeta \in \mathbb{S}^{d-1} : \langle \zeta, x \rangle > 0\}$  for all considered kernels of Appendix A such that tighter error bounds could apply. Consequently, exploring QMC designs on the hemisphere could be an interesting direction for further improving our theoretical analysis.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

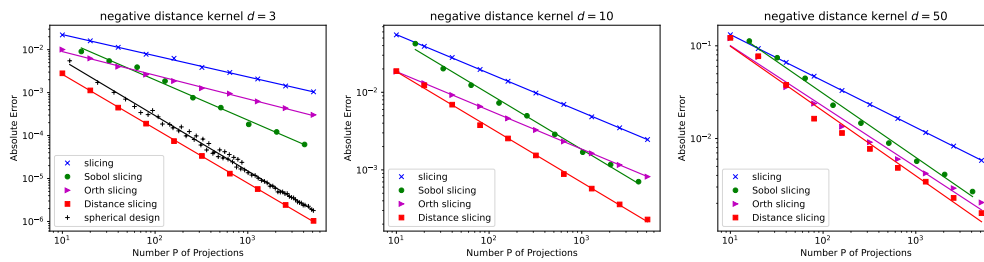


Figure 4: Loglog plot of the approximation error in (4) versus the number  $P$  of projections for the negative distance kernel. The results are averaged over 50 realizations of  $\xi^P$  and 1000 realizations of  $x$ . We fit a regression line in the loglog plot for each method to estimate the convergence rate, see also Table 3.

Table 3: Estimated convergence rates for the different methods and the negative distance kernel. We estimate the rate  $r$  by fitting a regression line in the loglog plot. Then, we obtain the estimated convergence rate  $\mathcal{O}(P^{-r})$  for some  $r > 0$ . Consequently, larger values of  $r$  correspond to a faster convergence. The resulting values of  $r$  are given in the below tables, the best values are highlighted in bold. See Figure 4 for a visualization.

Negative distance kernel					
Slicing-based					
Dimension	Slicing	Sobol	Orth	Distance	spherical design
$d = 3$	0.49	0.94	0.55	1.27	<b>1.29</b>
$d = 10$	0.50	0.72	0.50	<b>0.71</b>	-
$d = 50$	0.50	0.69	0.65	<b>0.70</b>	-

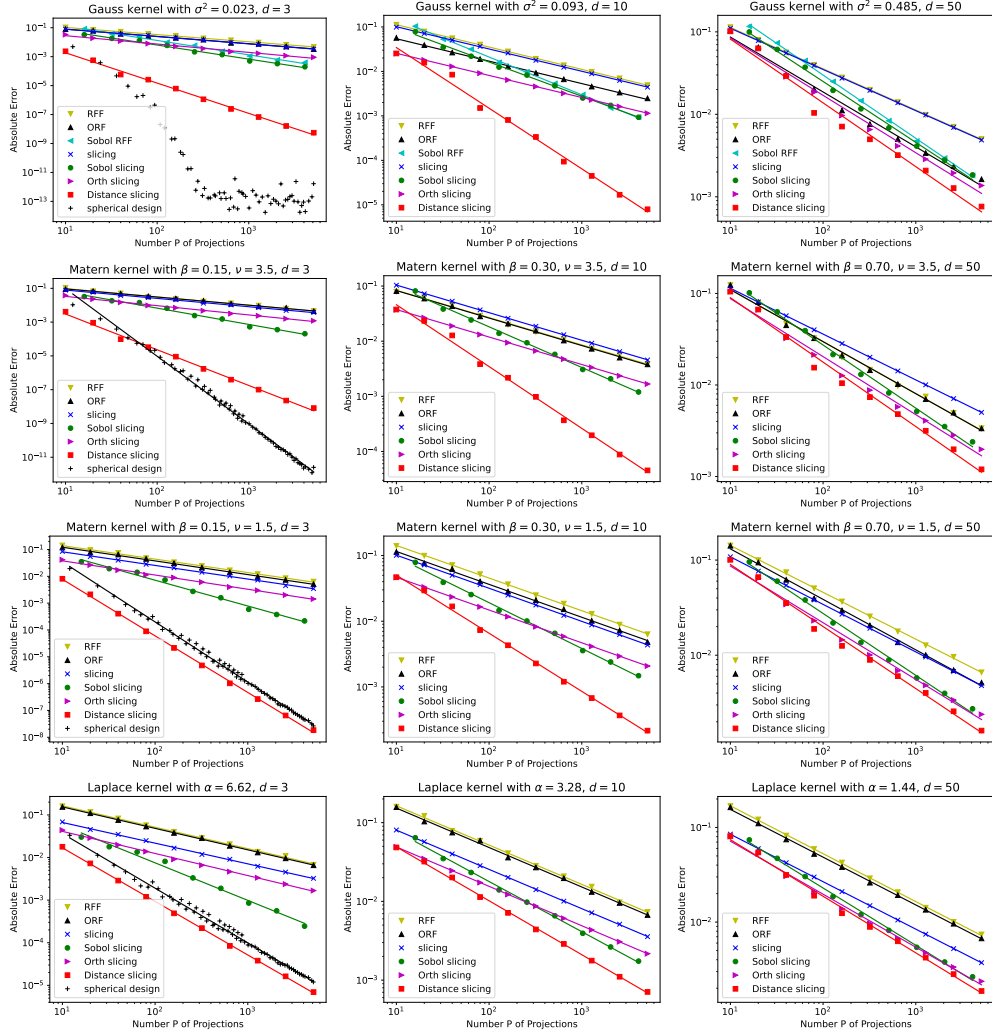


Figure 5: Loglog plot of the approximation error in (4) versus the number  $P$  of projections for different kernels and dimensions (left  $d = 3$ , middle  $d = 10$ , right  $d = 50$ ). The results are averaged over 50 realizations of  $\xi^P$  and 1000 realizations of  $x$ . The kernel parameters are set by the median rule with scale factor  $\gamma = \frac{1}{2}$ . We fit a regression line in the loglog plot for each method to estimate the convergence rate, see also Table 4.

Table 4: Estimated convergence rates for the different methods. We estimate the rate  $r$  by fitting a regression line in the loglog plot. Then, we obtain the estimated convergence rate  $\mathcal{O}(P^{-r})$  for some  $r > 0$ . Consequently, larger values of  $r$  correspond to a faster convergence. The resulting values of  $r$  are given in the below tables, the best values are highlighted in bolt. The kernel parameters are the same as in Figure 5 (scale factor  $\gamma = \frac{1}{2}$ ).

Gauss kernel with kernel scaling $\gamma = \frac{1}{2}$								Matérn kernel with $\nu = 3 + \frac{1}{2}$ and kernel scaling $\gamma = \frac{1}{2}$							
Dimension	RFF-based			Slicing-based				Dimension	RFF-based		Slicing-based				spherical design
	RFF	Sobol	ORF	Slicing	Sobol	Orth	Distance		RFF	ORF	Slicing	Sobol	Orth	Distance	
$d = 3$	0.50	0.99	0.50	0.51	0.97	0.58	<b>2.09</b>	$d = 3$	0.49	0.48	0.49	0.96	0.55	2.11	<b>4.01</b>
$d = 10$	0.50	0.85	0.50	0.50	0.77	0.50	<b>1.36</b>	$d = 10$	0.49	0.50	0.50	0.74	0.50	<b>1.12</b>	-
$d = 50$	0.50	<b>0.77</b>	0.67	0.50	0.74	0.70	<b>0.77</b>	$d = 50$	0.57	0.57	0.50	0.69	0.63	<b>0.70</b>	-
Matérn kernel with $\nu = 1 + \frac{1}{2}$ and kernel scaling $\gamma = \frac{1}{2}$								Laplace kernel with kernel scaling $\gamma = \frac{1}{2}$							
Dimension	RFF-based			Slicing-based				Dimension	RFF-based		Slicing-based				spherical design
	RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design		RFF	ORF	Slicing	Sobol	Orth	Distance	
$d = 3$	0.50	0.51	0.51	0.96	0.53	2.11	<b>2.24</b>	$d = 3$	0.50	0.50	0.49	0.88	0.52	1.26	<b>1.28</b>
$d = 10$	0.50	0.50	0.50	0.70	0.50	<b>0.88</b>	-	$d = 10$	0.50	0.50	0.50	0.64	0.50	<b>0.69</b>	-
$d = 50$	0.49	0.53	0.50	<b>0.65</b>	0.60	<b>0.65</b>	-	$d = 50$	0.51	0.50	0.50	<b>0.61</b>	0.56	0.60	-



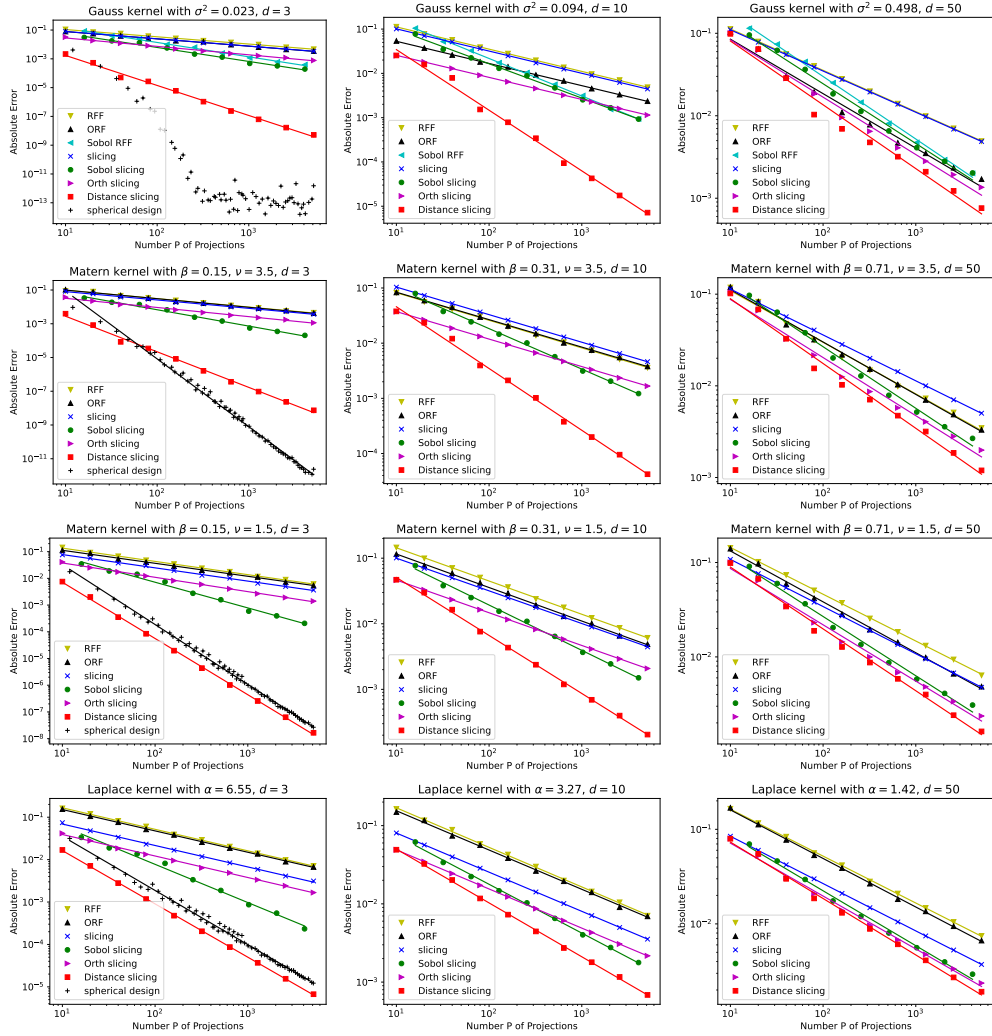


Figure 6: Loglog plot of the approximation error in (4) versus the number  $P$  of projections for different kernels and dimensions (left  $d = 3$ , middle  $d = 10$ , right  $d = 50$ ). The results are averaged over 50 realizations of  $\xi^P$  and 1000 realizations of  $x$ . The kernel parameters are set by the median rule with scale factor  $\gamma = 2$ . We fit a regression line in the loglog plot for each method to estimate the convergence rate, see also Table 5.

Table 5: Estimated convergence rates for the different methods. We estimate the rate  $r$  by fitting a regression line in the loglog plot. Then, we obtain the estimated convergence rate  $\mathcal{O}(P^{-r})$  for some  $r > 0$ . Consequently, larger values of  $r$  correspond to a faster convergence. The resulting values of  $r$  are given in the below tables, the best values are highlighted in bolt. The kernel parameters are the same as in Figure 6 (scale factor  $\gamma = 2$ ).

Gauss kernel with kernel scaling $\gamma = 2$							Matérn kernel with $\nu = 3 + \frac{1}{4}$ and kernel scaling $\gamma = 2$								
Dimension	RFF-based			Slicing-based				Dimension	RFF-based			Slicing-based			
	RFF	Sobol	ORF	Slicing	Sobol	Orth	Distance		RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design
$d = 3$	0.50	1.00	0.51	0.51	0.97	0.59	<b>2.08</b>	$d = 3$	0.50	0.51	0.50	0.97	0.54	2.11	<b>4.02</b>
$d = 10$	0.50	0.85	0.50	0.50	0.77	0.50	<b>1.37</b>	$d = 10$	0.50	0.50	0.50	0.73	0.50	<b>1.12</b>	-
$d = 50$	0.50	0.76	0.66	0.50	0.72	0.70	<b>0.77</b>	$d = 50$	0.56	0.57	0.50	0.67	0.63	<b>0.71</b>	-
Matérn kernel with $\nu = 1 + \frac{1}{4}$ and kernel scaling $\gamma = 2$							Laplace kernel with kernel scaling $\gamma = 2$								
Dimension	RFF-based			Slicing-based				Dimension	RFF-based			Slicing-based			
	RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design		RFF	ORF	Slicing	Sobol	Orth	Distance	spherical design
$d = 3$	0.50	0.49	0.50	0.96	0.54	2.10	<b>2.24</b>	$d = 3$	0.51	0.50	0.51	0.90	0.51	1.26	<b>1.28</b>
$d = 10$	0.51	0.51	0.50	0.69	0.50	<b>0.88</b>	-	$d = 10$	0.51	0.50	0.50	0.62	0.50	<b>0.69</b>	-
$d = 50$	0.50	0.54	0.50	0.63	0.60	<b>0.65</b>	-	$d = 50$	0.50	0.51	0.50	0.58	0.56	<b>0.60</b>	-

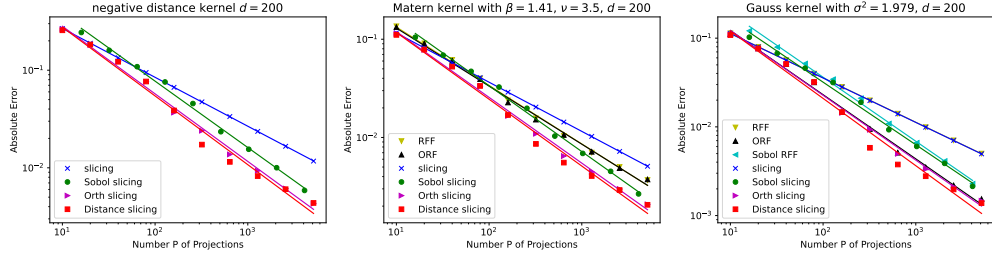


Figure 7: Loglog plot of the approximation error in (4) versus the number  $P$  of projections for  $d = 200$ . The results are averaged over 50 realizations of  $\xi^P$  and 1000 realizations of  $x$ . The kernel parameters are chosen by the median rule with scale factor  $\gamma = 1$ . We fit a regression line in the loglog plot for each method to estimate the convergence rate, see also Table 6.

Table 6: Estimated convergence rates for the different methods and kernels for  $d = 200$ . We estimate the rate  $r$  by fitting a regression line in the loglog plot. Then, we obtain the estimated convergence rate  $\mathcal{O}(P^{-r})$  for some  $r > 0$ . Consequently, larger values of  $r$  correspond to a faster convergence. The resulting values of  $r$  are given in the below tables, the best values are highlighted in bolt. The kernel parameters are the same as in Figure 7 (scale factor  $\gamma = 1$ ).

Dimension $d = 200$							
Kernel	RFF-based			Slicing-based			
	RFF	Sobol	ORF	Slicing	Sobol	Orth	Distance
Negative Distance	-	-	-	0.50	0.68	0.69	<b>0.70</b>
Matérn, $\nu = 3 + \frac{1}{2}$	0.60	-	0.59	0.50	0.67	0.67	<b>0.68</b>
Gauss	0.50	0.72	0.72	0.50	0.71	0.73	<b>0.76</b>

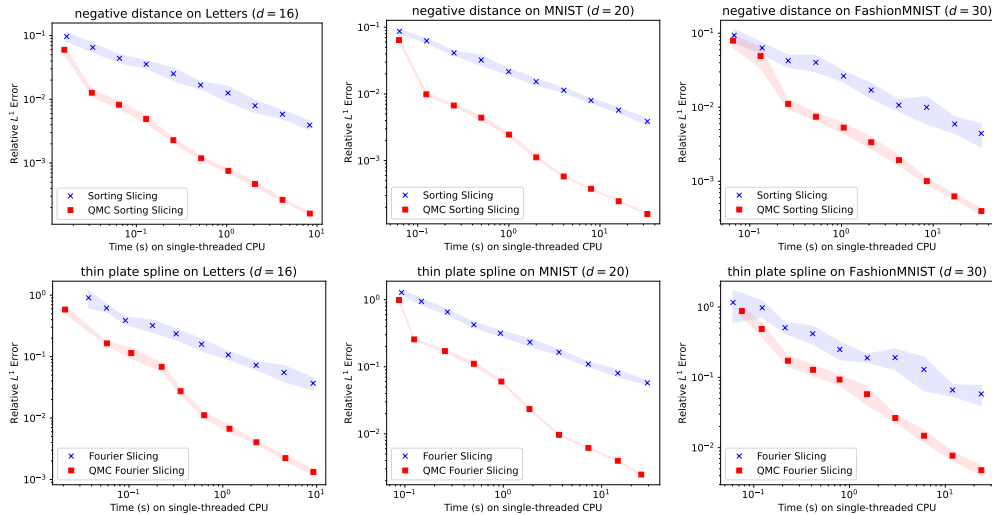


Figure 8: Loglog plot of the relative  $L^1$  approximation error versus computation time for computing the kernel summations (1) with the negative distance and thin plate spline kernel and different methods and datasets. MNIST and FashionMNIST are reduced to dimension  $d = 20$  and  $d = 30$  via PCA. We run each method 10 times. The shaded area indicates the standard deviation of the error. For the slicing method, we use  $P = 10 \cdot 2^k$  slices for  $k = 0, \dots, 9$ .

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

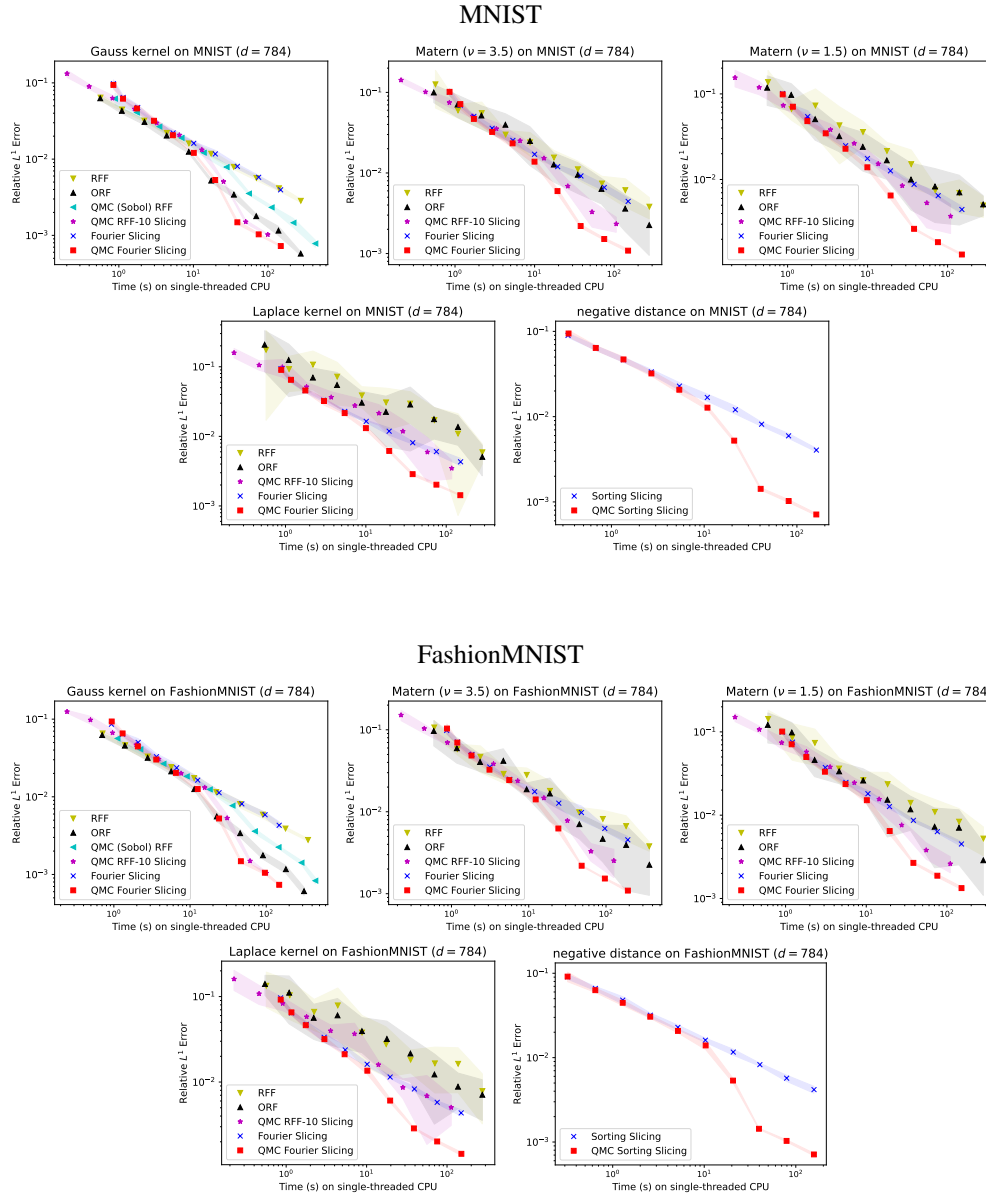


Figure 9: Loglog plot of the relative  $L^1$  approximation error versus computation time for computing the kernel summations (1) with different kernels and methods on the MNIST and FashionMNIST dataset without dimension reduction. We run each method 10 times. The shaded area indicates the standard deviation of the error. For the slicing method, we use  $P = 10 \cdot 2^k$  slices for  $k = 0, \dots, 9$ . In order to obtain similar computation times, we run RFF and ORF with  $D = 2P$  features.

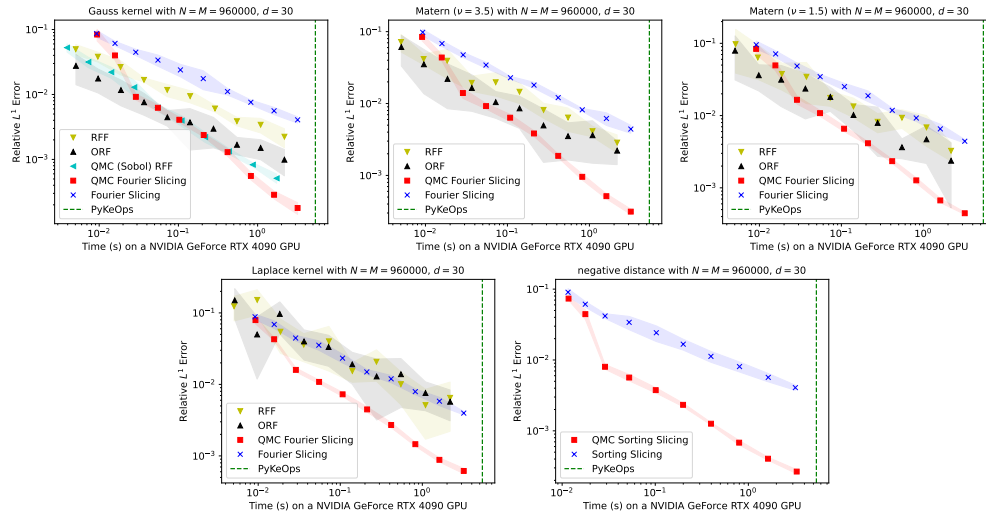


Figure 10: Loglog plot of the relative  $L^1$  approximation error versus GPU computation time for computing the kernel summations (1) with different kernels on a large dataset ( $M = N = 960000$ ). We run each method 10 times. The shaded area indicates the standard deviation of the error. For the slicing method, we use  $P = 10 \cdot 2^k$  slices for  $k = 0, \dots, 9$ . In order to obtain similar computation times, we run RFF and ORF with  $D = 4P$  features. Since PyKeOps computes the exact kernel sum, the computation time is indicated by a vertical line.

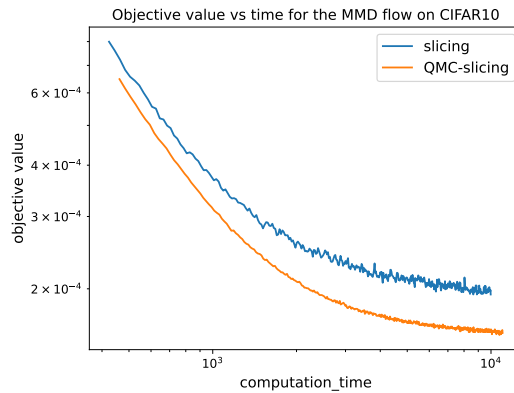


Figure 11: Plot of the objective value  $F_{\mathbf{y}}(x)$  for the MMD gradient flow on the CIFAR10 dataset on an NVIDIA GeForce RTX 4090 GPU. Note that computing this flow with exact kernel summations is computationally intractable. The plot omits the first 2000 steps of the gradient flow to improve the scaling of the plot.