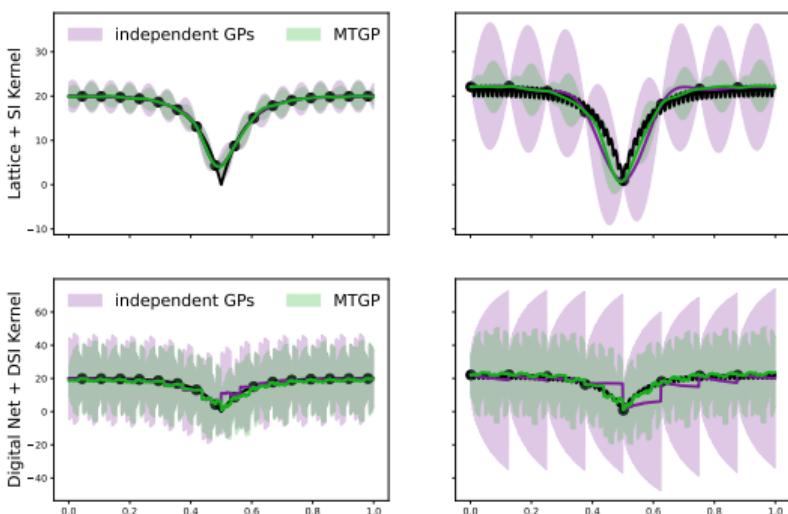


# Quasi-Monte Carlo and Fast Multi-Task Gaussian Process Regression

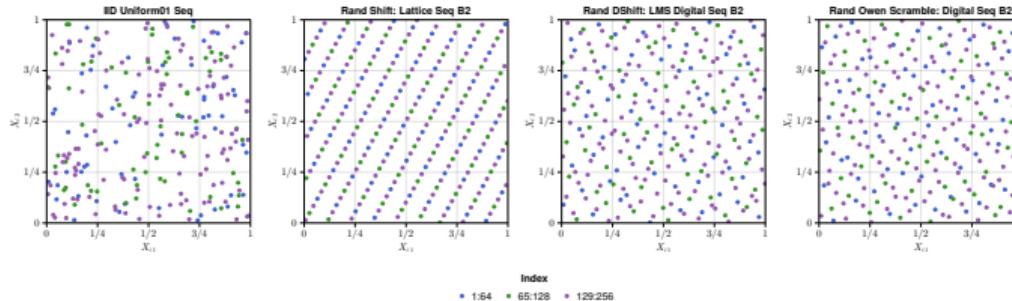
Aleksei G Sorokin<sup>12</sup>, Fred J Hickernell<sup>1</sup>, Pieterjan M Robbe<sup>2</sup>  
IIT, Department of Applied Math<sup>1</sup>. Sandia National Lab<sup>2</sup>.



## Quasi-Monte Carlo Methods

$$\mu = \mathbb{E}[g(\mathbf{T})] = \mathbb{E}[f(\mathbf{X})] \quad \approx \quad \frac{1}{N} \sum_{i=0}^{N-1} f(\mathbf{x}_i) = \hat{\mu}, \quad \mathbf{X} \sim \mathcal{U}[0,1]^d$$

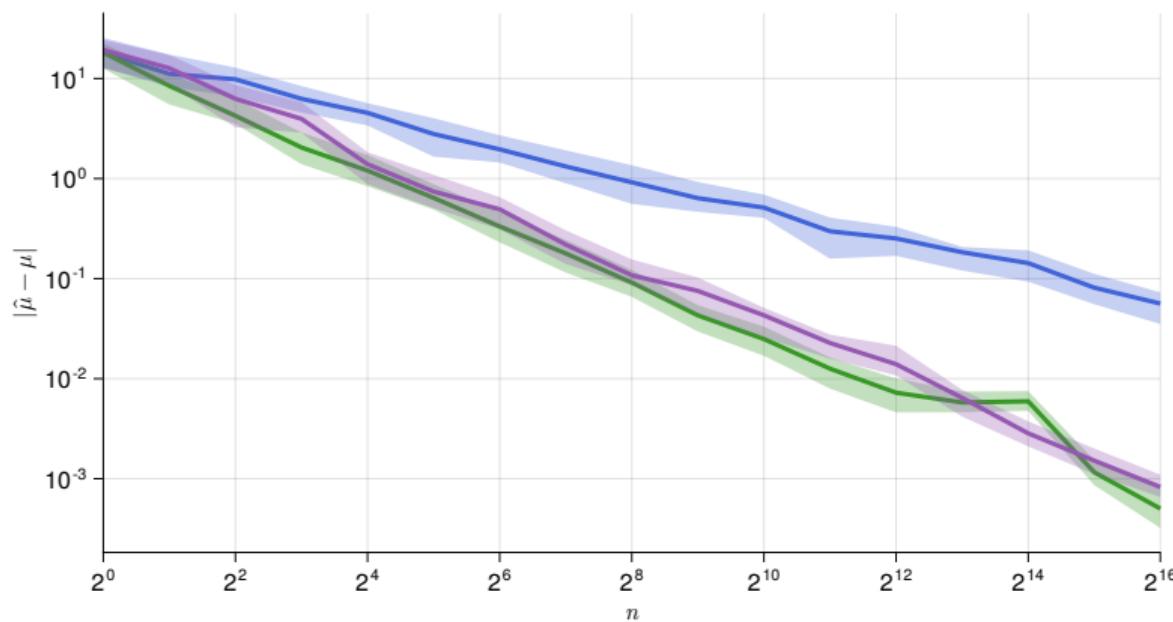
- True mean  $\mu$  is a function of true integrand  $g$  and true measure  $T$
  - Transformation  $T \sim \psi(X)$  makes  $f = g \circ \psi$ ,  
e.g.,  $T \sim \mathcal{N}(\mathbf{m}, \Sigma = \mathbf{A}\mathbf{A}^T)$  then  $\psi(X) = \mathbf{A}\Phi^{-1}(X) + \mathbf{m}$ , standard normal CDF  $\Phi$
  - Sampling nodes  $x_0, x_1, \dots \in [0, 1]^d$  chosen to be
    - IID  $\rightarrow$  Monte Carlo (MC)  $\rightarrow$  error  $\mathcal{O}(N^{-1/2})$
    - Low discrepancy (LD)  $\rightarrow$  Quasi-Monte Carlo (QMC)  $\rightarrow$  error  $\approx \mathcal{O}(N^{-1})$  or better  
[Dick and Pillichshammer, 2010, Dick et al., 2013]



# Monte Carlo vs Quasi-Monte Carlo

MC error like  $\mathcal{O}(N^{-1/2})$ . QMC error like  $\mathcal{O}(N^{-1+\delta})$  for  $\delta > 0$  arbitrarily small

— IID Uniform01 Seq — Rand Shift: Lattice Seq B2 — Rand DShift: LMS Digital Seq B2



$d = 7$  dimensional function from [Keister, 1996]

## Monte Carlo Stopping Criteria

How to choose  $N$  so that  $|\mu - \hat{\mu}| < \varepsilon$  for some user specified error tolerance  $\varepsilon > 0$ ?

e.g., want to approximate the value of a financial option to within 1 penny

Approximations should hold with sufficiently high probability or with guarantees

Central Limit Theorem (CLT) Stopping Criteria for IID-Monte Carlo

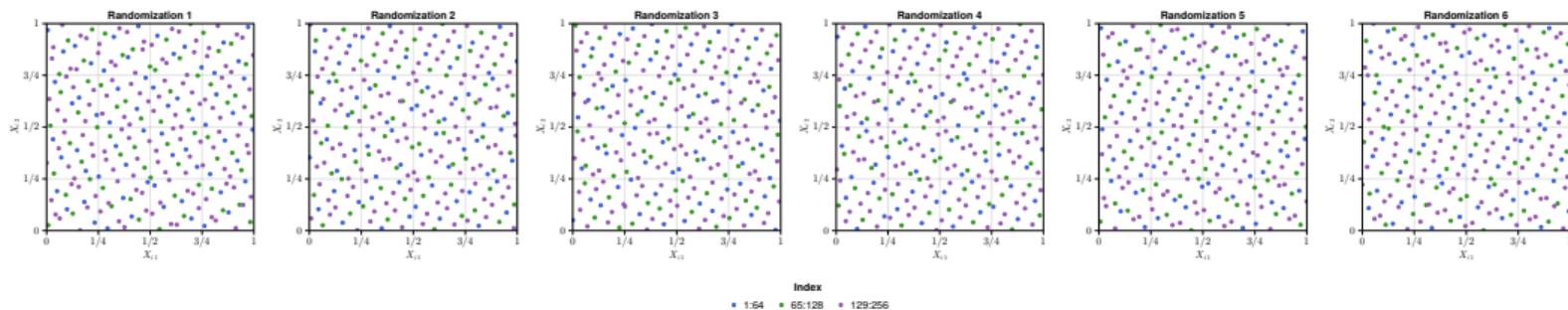
$$N \geq [2.58\sigma/\varepsilon]^2 \quad \rightarrow \quad |\mu - \hat{\mu}| < 2.58\sigma/\sqrt{N} \leq \varepsilon \text{ with probability } 99\%.$$

- Only a heuristic algorithm as CLT is asymptotic in  $N$
  - May use initial sample to approximate  $\sigma$
  - [Hickernell et al., 2013] gives a guaranteed version of this two-step method for finite  $N$  using Berry-Esseen inequalities and assuming a bounded Kurtosis

# Quasi-Monte Carlo Stopping Criteria

See [Owen, 2024] for a recent review of error estimation for QMC

1. Student- $T$  intervals for  $R$  randomizations of a LD point set [L'Ecuyer et al., 2023]
2. Quickly track decay of coefficients for functions in cones [Hickernell et al., 2017]
  - Fourier coefficients using LD lattices [Jiménez Rugama and Hickernell, 2014]
  - Walsh coefficients using LD digital nets [Hickernell and Jiménez Rugama, 2014]
3. Fast Bayesian cubature for functions in cones [Rathinavel, 2019]
  - Gram matrices diagonalizable by FFT [Rathinavel and Hickernell, 2019]
  - Gram matrices diagonalizable by FWHT<sup>1</sup> [Rathinavel and Hickernell, 2022]



<sup>1</sup>Fast Walsh Hadamard Transform (FWHT) [Fino and Algazi, 1976]

# Gaussian Process Regression (GPR)

$$f \sim \text{GP}(0, K)$$

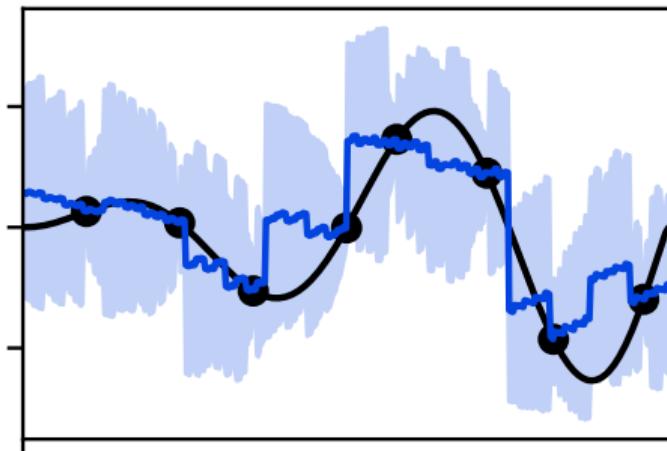
Posterior mean and covariance

$$\mathbb{E}[f(\mathbf{x})|\mathbf{X}, \mathbf{f}] = \mathbf{K}_X(\mathbf{x})\mathbf{K}^{-1}\mathbf{f}$$

$$\text{Cov}[f(\mathbf{x}), f(\mathbf{x}')|\mathbf{X}, \mathbf{f}] = K(\mathbf{x}, \mathbf{x}') - \mathbf{K}_X(\mathbf{x})^T \mathbf{K}^{-1} \mathbf{K}_X(\mathbf{x}')$$

- SPD  $K : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$
- Sampling locations  $\mathbf{X} = \{\mathbf{x}_i\}_{i=0}^{N-1}$
- Sample values  $\mathbf{f} = \{f(\mathbf{x}_i)\}_{i=0}^{N-1}$
- Kernel vector  $\mathbf{K}_X(\mathbf{x}) = \{K(\mathbf{x}, \mathbf{x}_i)\}_{i=0}^{N-1}$
- Gram matrix  $\mathbf{K} = \{K(\mathbf{x}_i, \mathbf{x}_{i'})\}_{i,i'=0}^{N-1}$

Standard GPR cost is  $\mathcal{O}(N^3)$



## Bayesian Cubature

$$f \sim \text{GP}(0, K)$$

$\mu = \mathbb{E}_X[f(\mathbf{X})]$ ,  $\mathbf{X} \sim \mathcal{U}[0,1]^d$  is Gaussian with posterior cubature mean and variance

$$\hat{\mu} := \mathbb{E}_f[\mu | \mathbf{X}, \mathbf{f}] = \tilde{\mathbf{K}}_{\mathbf{X}}^T \mathbf{K}^{-1} \mathbf{f}$$

$$\hat{\sigma}^2 := \text{Var}_f[\mu | \mathbf{X}, \mathbf{f}] = \tilde{K} - \tilde{\mathbf{K}}_{\mathbf{X}}^T \mathbf{K}^{-1} \tilde{\mathbf{K}}_{\mathbf{X}}$$

$\tilde{\mathbf{K}}_{\mathbf{X}} = \mathbb{E}_{\mathbf{X}}[\mathbf{K}_{\mathbf{X}}(\mathbf{X})]$  and  $\tilde{K} = \mathbb{E}_{(\mathbf{X}, \mathbf{X}')}[K(\mathbf{X}, \mathbf{X}')]$  for independent  $\mathbf{X}, \mathbf{X}' \sim \mathcal{U}[0,1]^d$

$$\therefore |\mu - \hat{\mu}| < 2.58\hat{\sigma} \text{ with probability 99\%}$$

If  $\tilde{\mathbf{K}}_{\mathbf{X}}^T \mathbf{K}^{-1} = \mathbf{1}/N$  then  $\hat{\mu}$  is the sample average as in (Q)MC

# Fast GPs Pairing LD Points with Special Kernels

## 1. LD Lattices + Shift Invariant (SI) Kernels

- Give circulant Gram matrices  $K = \{K(\mathbf{x}_i, \mathbf{x}_{i'})\}_{i,i'=0}^{N-1}$
- ∴ Eigendecomp  $K = V \Lambda \bar{V}$  where  $\bar{V}$  is the DFT matrix → FFT in  $\mathcal{O}(N \log N)$
- [Rathinavel and Hickernell, 2019]

## 2. LD Digital Nets + Digitally Shift Invariant (DSI) Kernels

- Give Recursive Symmetric Block Toeplitz (RSBT) Gram matrices K
- ∴ Eigendecomp  $K = V \Lambda \bar{V}$  where  $\bar{V}$  is the Hadamard matrix → FWHT in  $\mathcal{O}(N \log N)$
- [Rathinavel and Hickernell, 2022]

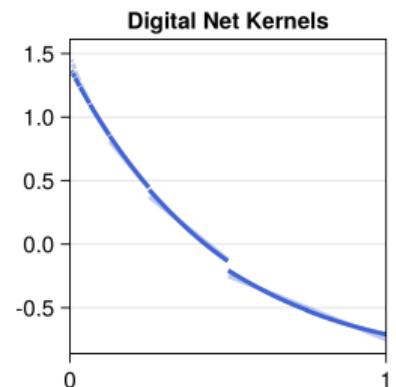
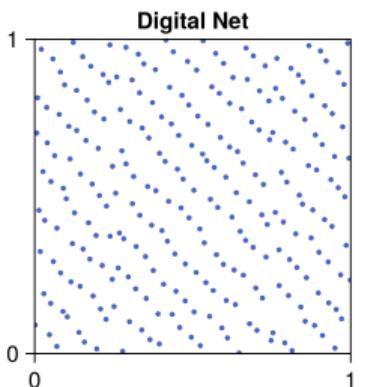
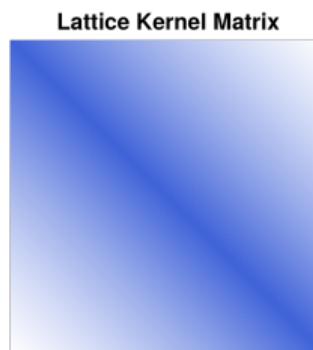
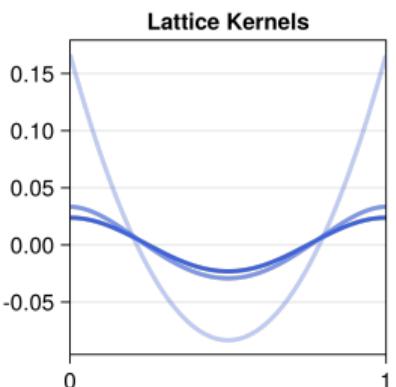
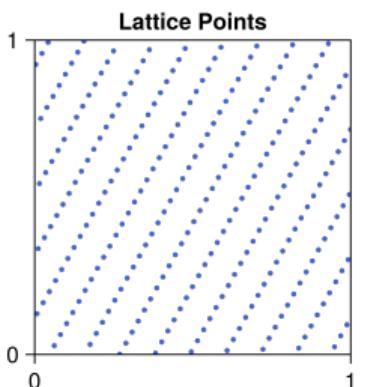
Let  $v_1 = \mathbf{1}/\sqrt{N}$  and  $k_1$  be the first columns of  $\bar{V}$  and K respectively:

$$\lambda := \Lambda \mathbf{1} = \sqrt{N} \Lambda \bar{v}_1 = \sqrt{N} \bar{V} V \Lambda \bar{v}_1 = \sqrt{N} \bar{V} k_1$$

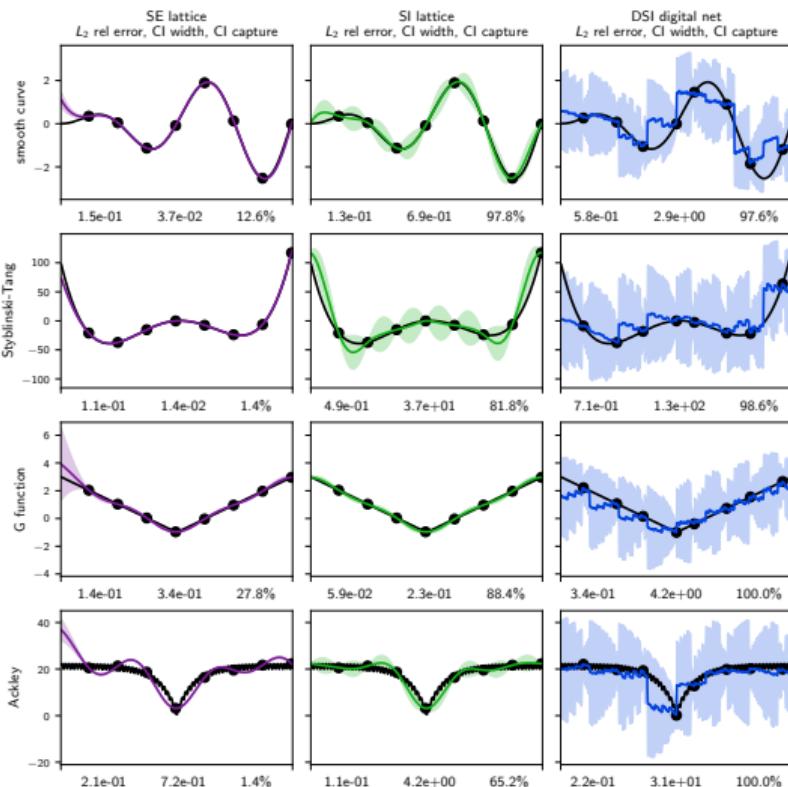
- $Ka$ ,  $K^{-1}a$ , and  $|K|$  can all be computed in  $\mathcal{O}(N \log N)$  computations
- Only requires evaluating and storing the first column of K

Originally developed in the context of fast Bayesian Cubature [Rathinavel, 2019]

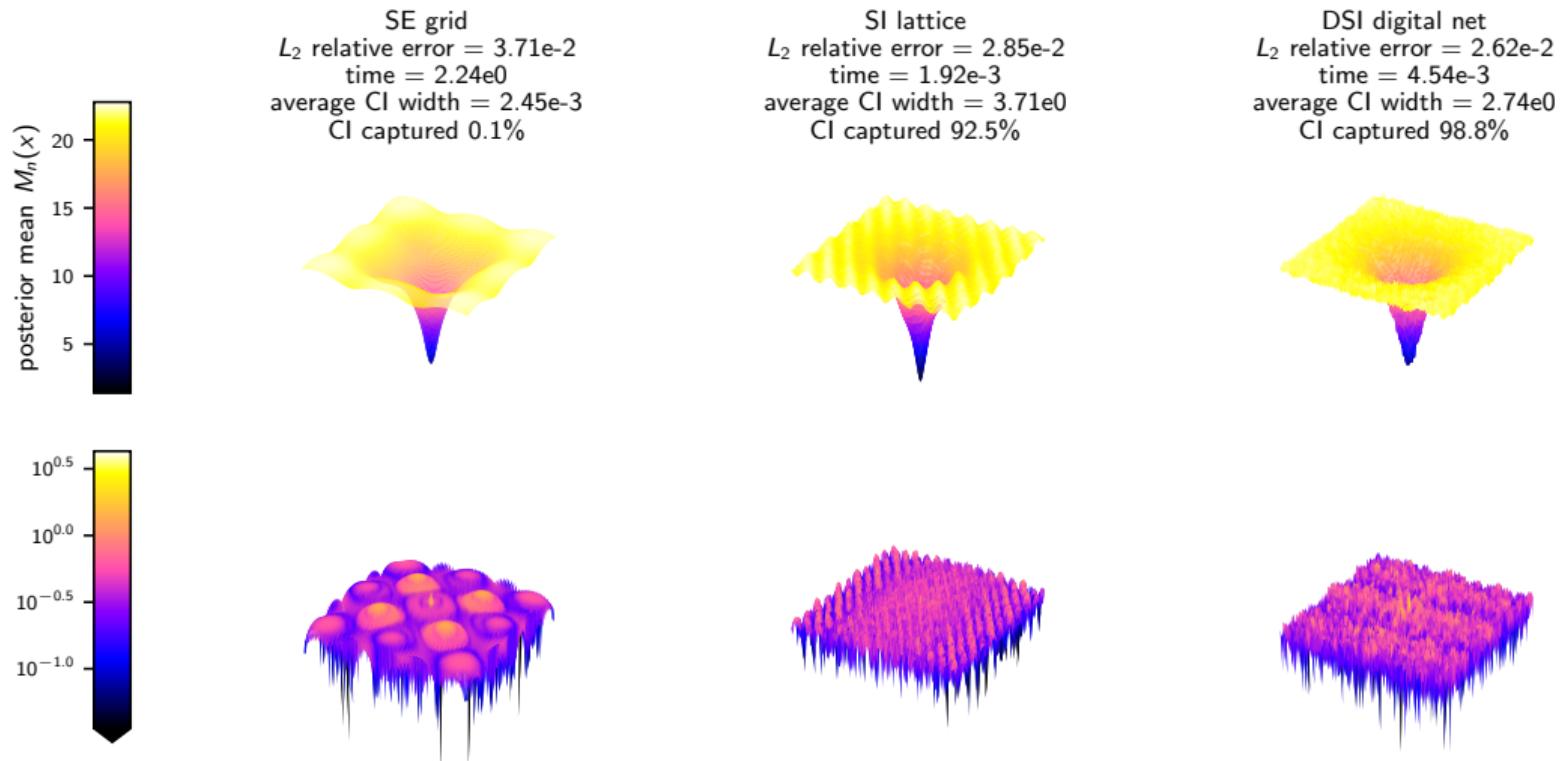
Lattices + SI  $K$  = Circulant  $K$  and Digital Nets + DSII  $K$  = RSBT  $K$



# Examples in One Dimension [Surjanovic and Bingham]



# Ackley Function in Two Dimensions [Surjanovic and Bingham]

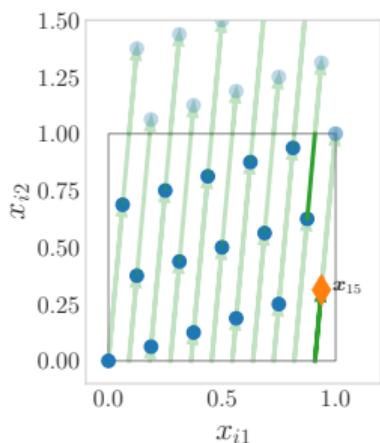
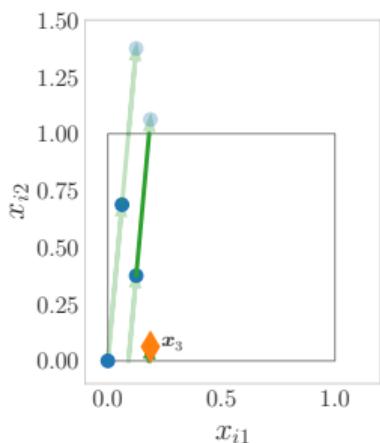
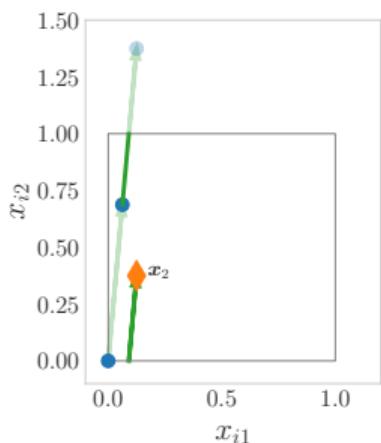
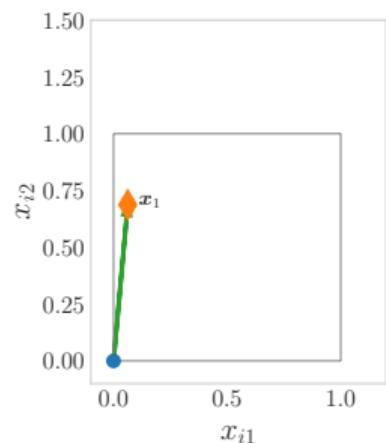


## Unrandomized Lattice

Generating vector  $\mathbf{g} \in \{1, \dots, N-1\}^d$  gives the **unrandomized lattice**

$$\mathbf{x}_i^{\text{lat}} = i\mathbf{g}/N \pmod{1}, \quad i = 0, \dots, N-1$$

Lattice  $\mathbf{x}_i^{\text{lat}} = i(1, 11)/16 \pmod{1}, \quad i = 0, \dots, 15$

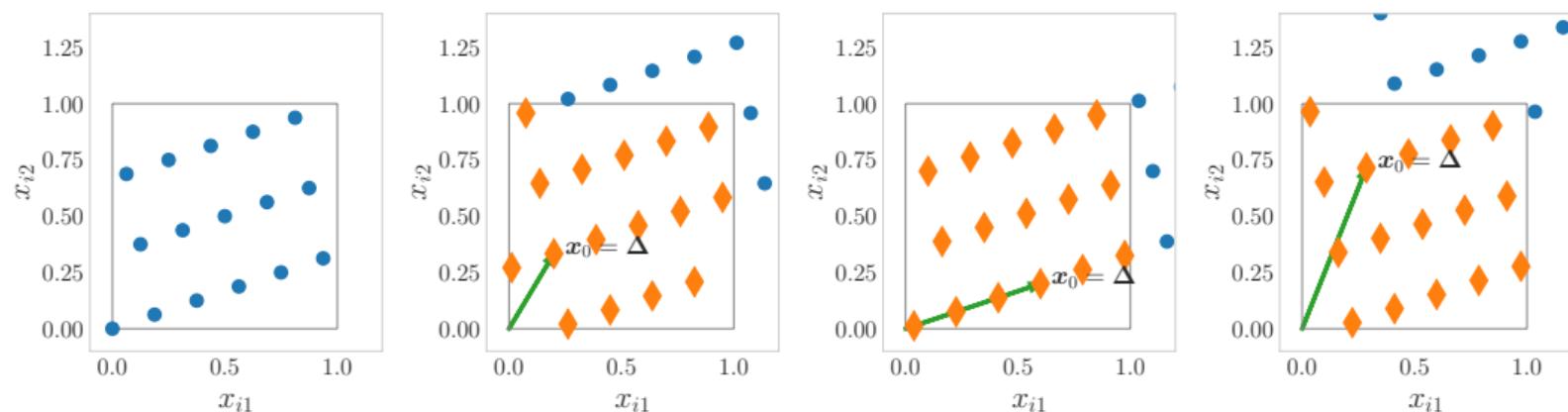


## Shifted Lattice

Generating vector  $\mathbf{g} \in \{1, \dots, N-1\}^d$  and shift  $\Delta \in [0, 1]^d$  gives the **shifted lattice**

$$\mathbf{x}_i = (\mathbf{x}_{i\text{lat}} + \Delta) \bmod 1 = (i\mathbf{g}/N + \Delta) \bmod 1, \quad i = 0, \dots, N-1$$

$$\text{Lattice } \mathbf{x}_i = i(1, 11)/16 + \Delta \pmod{1}, \quad i = 0, \dots, 15$$



## Shifted Lattices + SI Kernels = Circulant Matrices

**Shift invariant (SI)** kernels  $K$  may be written in terms of some  $\hat{K} : [0, 1]^d \rightarrow \mathbb{R}$  s.t.

$$K(\mathbf{x}, \mathbf{x}') = \hat{K}(\mathbf{x} \ominus \mathbf{x}'), \quad \mathbf{x} \ominus \mathbf{x}' = (\mathbf{x} - \mathbf{x}') \mod 1$$

For  $\{\mathbf{x}_i\}_{i=0}^{N-1}$  a **shifted lattice**,

$$\begin{aligned}\mathbf{x}_i \ominus \mathbf{x}_{i'} &= \left( \frac{i\mathbf{g}}{N} \ominus \Delta \right) \ominus \left( \frac{i'\mathbf{g}}{N} \ominus \Delta \right) \\ &= \left\{ \frac{(i - i') \mod N}{N} \mathbf{g} \right\} = \mathbf{x}_0 \ominus \mathbf{x}_{(i' - i) \mod N} \\ \therefore K_{i,i'} &= K_{0,[(i' - i) \mod N]}\end{aligned}$$

so  $K$  is circulant

## Examples of Shift Invariant (SI) Kernels

For  $f : [0, 1] \rightarrow \mathbb{R}$ , if  $f^{(\alpha)}$  has an absolutely convergent Fourier series for some  $\alpha \in \mathbb{N}_0$  and  $f^{(\beta)}$  periodic for all  $\beta \in \{0, \dots, \alpha - 1\}$  then

$$f(x) = \sum_{k \in \mathbb{Z}} \widehat{f}(k) e^{2\pi i k x} \quad \longrightarrow \quad f^{(\alpha)}(x) = \sum_{k \in \mathbb{Z}} \widehat{f^{(\alpha)}}(k) e^{2\pi i k x}, \quad \widehat{f^{(\alpha)}}(k) = (2\pi i k)^\alpha \widehat{f}(k)$$

For  $B_i$  the  $i^{\text{th}}$  Bernoulli polynomial,

$$\mathring{K}_\alpha(x, x') = \sum_{k \in \mathbb{Z} \setminus \{0\}} \frac{e^{2\pi i(x-x')}}{k^{2\alpha}} = \frac{(-1)^{\alpha+1} (2\pi)^{2\alpha}}{(2\alpha)!} B_{2\alpha}(x \ominus x') =: \mathring{K}_\alpha(x \ominus x')$$

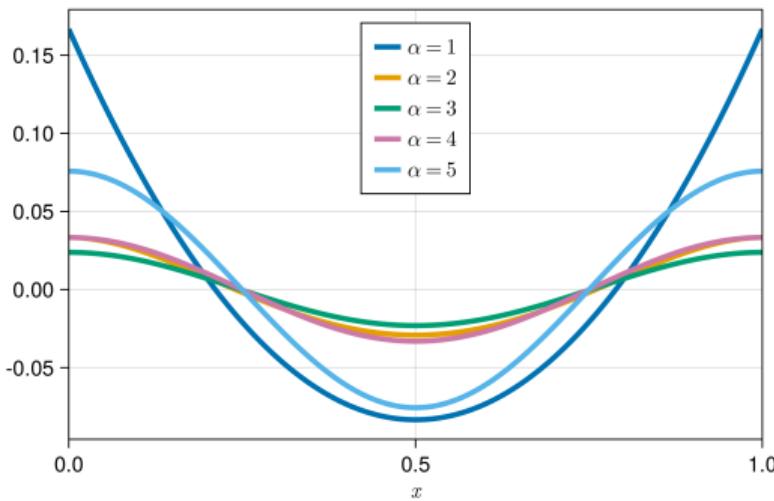
is the kernel of the Sobolev RKHS  $\mathring{H}^\alpha$  with inner product

$$\langle f, g \rangle_\alpha = (-1)^\alpha (2\pi)^{2\alpha} \int_0^1 f^{(\alpha)}(x) g^{(\alpha)}(x) dx$$

Let  $\mathring{H}^\alpha$  be the RKHS with SI kernel

$$\mathring{K}_\alpha(\mathbf{x}, \mathbf{x}') = \gamma \prod_{j=1}^d [1 + \eta_j \mathring{K}_{\alpha_j}(x_j \ominus x'_j)]$$

$$\mathring{K}_\alpha(x) = \frac{(-1)^{\alpha+1} (2\pi)^{2\alpha}}{(2\alpha)!} B_{2\alpha}(x)$$



If  $f \in \mathring{H}^{\alpha,1}$  then QMC error  $\left| \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} - \frac{1}{N} \sum_{i=0}^{N-1} f(\mathbf{x}_i) \right|$  is  $\mathcal{O}(n^{-\alpha+\delta})$  for certain lattices  $\{\mathbf{x}_i\}_{i=0}^{N-1}$  and certain weights  $\boldsymbol{\eta}$  (with  $\delta > 0$  arbitrarily small)

## Digitwise Operations

Prime base  $b \geq 2$  expansion of  $x \in [0, 1)$  is

$$x = .x_1 x_2 x_3 \cdots_b = \sum_{\ell \in \mathbb{N}} x_\ell b^{-\ell}, \quad \text{e.g.} \quad .375 = .011_2,$$

with digitwise subtraction (digitwise exclusive or in base  $b = 2$ )

$$x \ominus y := \sum_{\ell \in \mathbb{N}} ((x_\ell - y_\ell) \bmod b) b^{-\ell}, \quad \text{e.g.} \quad .375 \ominus .625 = .011_2 \ominus .101_2 = .110_2 = .75.$$

Similarly for  $k \in \mathbb{N}_0$

$$k = \cdots k_2 k_1 k_0.0_b = \sum_{\ell \in \mathbb{N}_0} k_\ell b^\ell, \quad \text{e.g.} \quad 5 = 101_2,$$

$$k \ominus h := \sum_{\ell \in \mathbb{N}_0} ((k_\ell + h_\ell) \bmod b) b^\ell, \quad \text{e.g.} \quad 5 \ominus 6 = 101_2 \ominus 110_2 = 011_2 = 3.$$

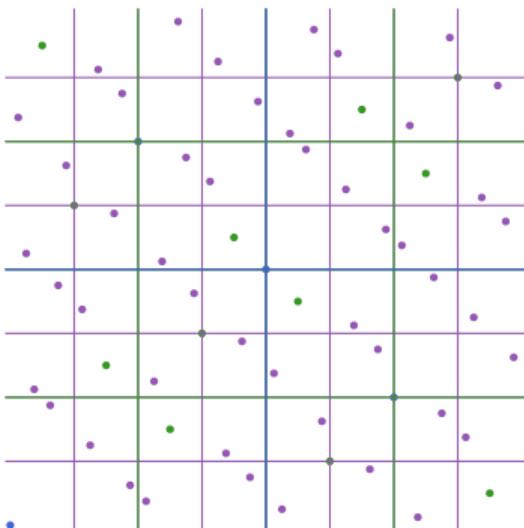
Apply elementwise to vectors

$$\mathbf{x} \ominus \mathbf{x}' = (x_1 \ominus x'_1, \dots, x_d \ominus x'_d), \quad \mathbf{k} \ominus \mathbf{k}' = (k_1 \ominus k'_1, \dots, k_d \ominus k'_d)$$

## Digital Nets

$N = b^P$  and fixed  $\mathbf{g}_{b^0}, \mathbf{g}_{b^1}, \dots, \mathbf{g}_{b^{P-1}} \in [0,1)^d$  give the **unrandomized digital net**

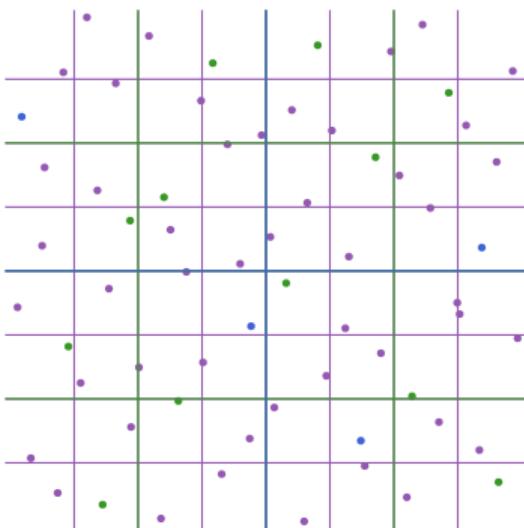
$$\mathbf{x}_i^{\text{dig}} = \bigoplus_{\ell=0}^{P-1} \mathbf{i}_\ell \mathbf{g}_{b^\ell}, \quad i = \sum_{\ell=0}^{P-1} \mathbf{i}_\ell \in \{0, \dots, b^P - 1\}$$



## Digitally Shifted Digital Nets (and other Randomizations)

$\mathbf{g}_{b^0}, \mathbf{g}_{b^1}, \dots, \mathbf{g}_{b^{P-1}} \in [0,1]^d$  and digital shift (DS)  $\Delta \in [0,1]^d$  gives the **DS digital net**

$$\mathbf{x}_i = \mathbf{x}_i^{\text{dig}} \ominus \Delta = \bigoplus_{\ell=0}^{P-1} i_\ell \mathbf{g}_{2^\ell} \ominus \Delta, \quad i = \sum_{\ell=0}^{P-1} i_\ell \in \{0, \dots, b^P - 1\}$$



Linear Matrix Scrambling (LMS) and Owen Scrambling [Owen, 2003] also available

# Digitally Shifted Digital Nets + DSI Kernels = RSBT Matrices

**Digitally Shift Invariant (DSI)** kernels  $K$  may be written as

$$K(\mathbf{x}, \mathbf{x}') = \hat{K}(\mathbf{x} \ominus \mathbf{x}') \quad \text{for some } \hat{K} : [0, 1]^d \rightarrow \mathbb{R}$$

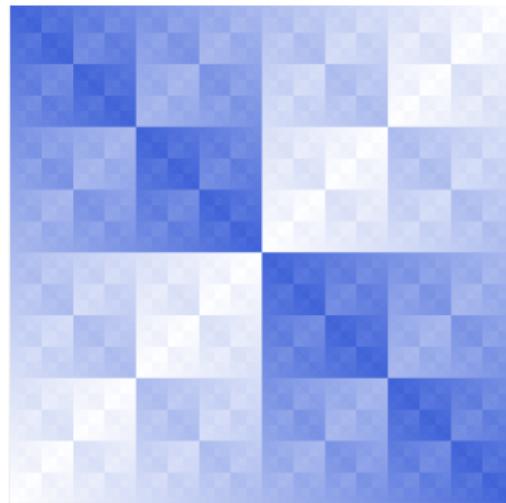
For  $\{\mathbf{x}_i\}_{i=0}^{2^P-1}$  a  $b = 2$  **DS digital net** and  $0 \leq p' < P$  and  $i, i' \in \{0, \dots, 2^{p'}-1\}$

- $K_{i+2^{p'}, i'+2^{p'}} = K_{i, i'}$  since

$$\begin{aligned} & \mathbf{x}_{i+2^{p'}} \ominus \mathbf{x}_{i'+2^{p'}} \\ &= (\mathbf{x}_i \ominus \mathbf{g}_{2^{p'}}) \ominus (\mathbf{x}_{i'} \ominus \mathbf{g}_{2^{p'}}) \\ &= \mathbf{x}_i \ominus \mathbf{x}_{i'} \end{aligned}$$

- $K_{i+2^{p'}, i'} = K_{i, i'+2^{p'}}$  since

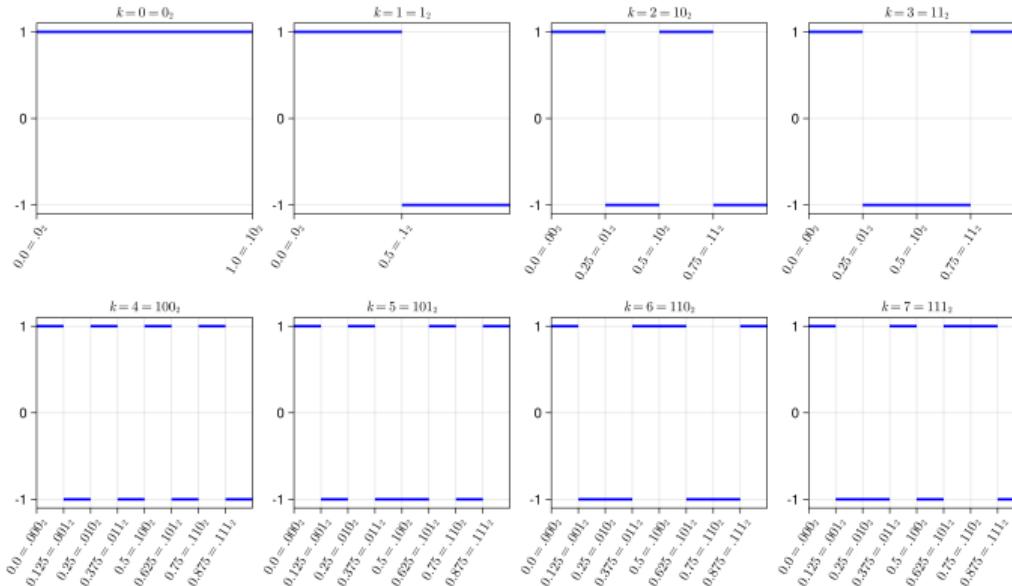
$$\begin{aligned} & \mathbf{x}_{i+2^{p'}} \ominus \mathbf{x}_{i'} \\ &= \mathbf{x}_i \ominus \mathbf{g}_{2^{p'}} \ominus \mathbf{x}_{i'} \\ &= \mathbf{x}_i \ominus \mathbf{x}_{i'+2^{p'}} \end{aligned}$$



so  $K$  is Recursive Symmetric Block Toeplitz (RSBT)

## $b = 2$ Walsh Functions $\text{wal}_k(x)$ and Weight Functions $\mu_\alpha(k)$

$\text{wal}_k(x) = (-1)^{\sum_{\ell \in \mathbb{N}_0} k_\ell x_{\ell+1}}$  with binary expansions  $x = \sum_{\ell \in \mathbb{N}} x_\ell b^{-\ell}$  and  $k = \sum_{\ell \in \mathbb{N}_0} k_\ell b^\ell$



$\mu_\alpha(k)$  sums the  $\alpha$  largest indices of non-zero digits in the base  $b$  expansion of  $k$ , e.g., for  $b = 2$ ,  $k = 13 = 1101_2$  has nonzero digit indices  $(4, 3, 1)$  so

$$\mu_1(13) = 4, \quad \mu_2(13) = 4 + 3, \quad \mu_3(13) = \mu_4(13) = \dots = 4 + 3 + 1$$

## Decay of Walsh Coefficients and QMC

For  $\alpha \geq 2$ , let the  $d = 1$  Sobolev RKHS  $H^\alpha$  have kernel  $K_\alpha$  and inner product<sup>2</sup>

$$\langle f, g \rangle_\alpha = \sum_{\beta=1}^{\alpha-1} \int_0^1 f^{(\beta)}(x) dx \int_0^1 g^{(\beta)}(x) dx + \int_0^1 f^{(\alpha)}(x) g^{(\alpha)}(x) dx$$

[Dick, 2008, 2009] show that  $\exists C_{f,\alpha} > 0$  s.t. if  $f \in H^\alpha$  then

$$|\hat{f}(k)| = \left| \int_0^1 f(x) \overline{\text{walk}(x)} dx \right| < C_{f,\alpha} b^{-\mu_\alpha(k)}$$

Let  $H^\alpha$  be the RKHS with kernel  $K_\alpha(\mathbf{x}, \mathbf{x}') = \gamma \prod_{j=1}^d [1 + \eta_j K_{\alpha_j}(x_j, x'_j)]$

If  $f \in H^{\alpha \mathbf{1}}$  then QMC error  $\left| \int_{[0,1]^d} f(\mathbf{x}) d\mathbf{x} - \frac{1}{N} \sum_{i=0}^{N-1} f(\mathbf{x}_i) \right|$  is  $\mathcal{O}(n^{-\alpha+\delta})$  for certain (higher order) digital nets  $\{\mathbf{x}_i\}_{i=0}^{N-1}$  and certain weights  $\boldsymbol{\eta}$  (with  $\delta > 0$  arbitrarily small)

---

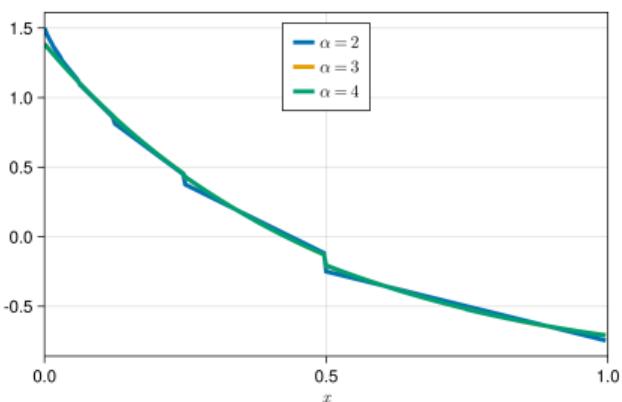
<sup>2</sup>For the  $\alpha = 1$  case see [Dick and Pillichshammer, 2005]

## Examples of Digitally Shift Invariant (DSI) Kernels

- \*  $K_\alpha$ , the kernel of Sobolev RKHS  $H^\alpha$  is **not DSI**, however  $H^\alpha \subset \tilde{H}^\alpha$  where  $\tilde{H}^\alpha$  is an RKHS with DSI kernel

$$\tilde{K}_\alpha(x, y) = \sum_{k \in \mathbb{N}} \frac{\text{wal}_k(x \ominus y)}{b^{\mu_\alpha(k)}} =: \tilde{K}_\alpha(x \ominus y)$$

Below  $\tilde{K}_\alpha(x)$  with  $b = 2$  is shown. Discontinuities at  $\{2^{-a} : a \in \mathbb{N}\}$  among others.



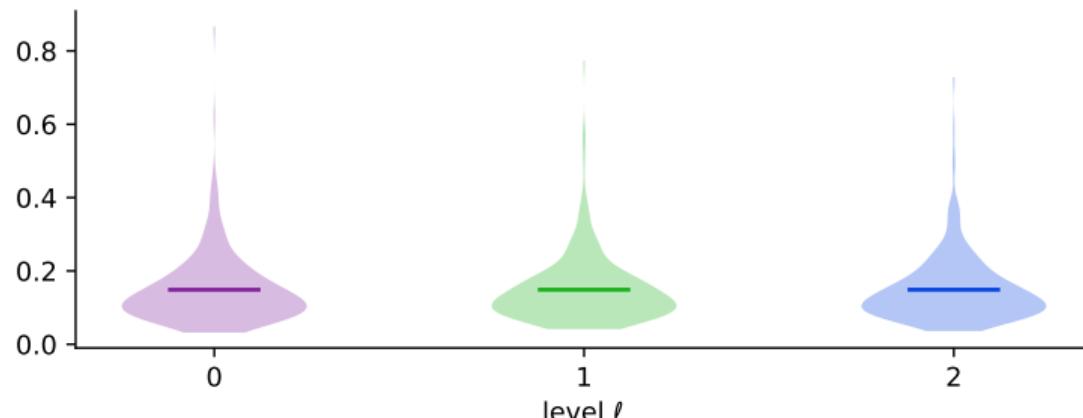
# Multilevel (Multi-Task) Modeling

Given a multilevel simulation

$$f : \{1, \dots, L\} \times [0, 1]^d \rightarrow \mathbb{R},$$

we want to model  $f(L, \cdot)$ , the true (maximum-fidelity) simulation.

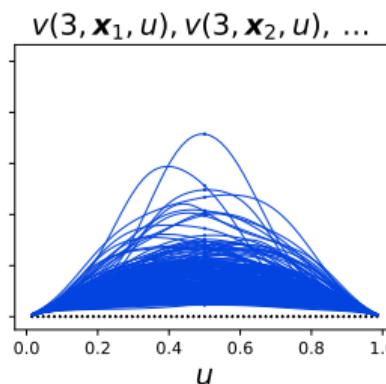
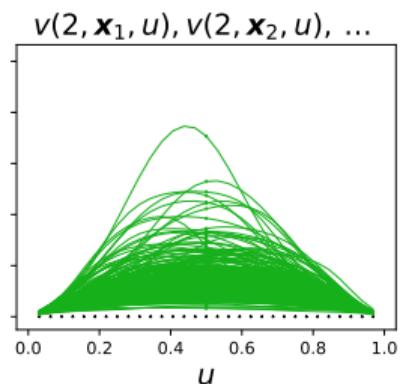
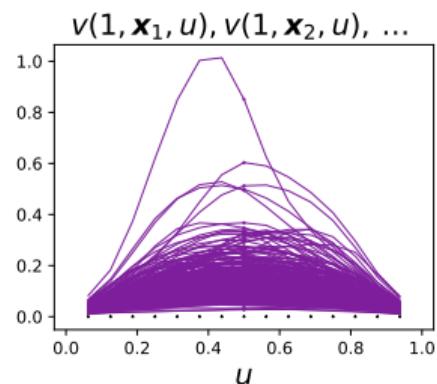
- $f(\ell, \mathbf{x})$  simulates at level  $\ell \in \{1, \dots, L\}$  and with parameters  $\mathbf{x} \in [0, 1]^d$
- Higher levels are typically more expensive to evaluate
- $f(1, \cdot), f(2, \cdot), \dots, f(L, \cdot)$  are typically highly correlated



Numerical Solutions of PDEs with Random Coefficients

$$f(\ell, \mathbf{x}) = \mathcal{F}(v(\ell, \mathbf{x}, \cdot))$$

- $v : \{1, \dots, L\} \times [0, 1]^d \times \Omega$  is the numerical solution to the PDE
  - $\boldsymbol{x}$  represent random coefficients, e.g. coefficients in a Karhunen-Loéve expansion
  - $\ell$  controls the fidelity of the numerical solver e.g. the mesh width is  $2^{-\ell}$
  - $\mathcal{F}$  is a (possibly non-linear) functional of the PDE solution, e.g.,
    - $\mathcal{F}(v(\ell, \boldsymbol{x}, \cdot)) = \mathbb{E}[v(\ell, \boldsymbol{x}, \mathbf{U})]$  where  $\mathbf{U} \sim \mathcal{U}(\Omega)$ , or
    - $\mathcal{F}(v(\ell, \boldsymbol{x}, \cdot)) = v(\ell, \boldsymbol{x}, u)$  for some  $u \in \Omega$ , e.g.,  $u = 1/2$  shown below



# Multi-Task Gaussian Processes (MTGPs)

$$f \sim \text{GP}(0, K)$$

$N = N_1 + \dots + N_L$  sampling locations  $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_L$  where  $\mathcal{D}_\ell = \{(\ell, \mathbf{x}_{\ell i})\}_{i=1}^{N_\ell}$ .  
Posterior mean and covariance

$$\mathbb{E}[f(\ell, \mathbf{x})] = \mathbf{K}^T(\ell, \mathbf{x}) \mathbf{K}^{-1} \mathbf{f}$$

$$\mathbb{C}[f(\ell, \mathbf{x}), f(\ell', \mathbf{x}')] = K((\ell, \mathbf{x}), (\ell', \mathbf{x}')) - \mathbf{K}^T(\ell', \mathbf{x}') \mathbf{K}^{-1} \mathbf{K}(\ell', \mathbf{x}')$$

- $\mathbf{K}(\ell, \mathbf{x}) = K(\mathcal{D}, (\ell, \mathbf{x}))$  and  $\mathbf{f} = f(\mathcal{D})$  are length  $N$  vectors
- $\mathbf{K} = K(\mathcal{D}, \mathcal{D}^T)$  is the  $N \times N$  Gram matrix

Kernel  $K$  depends on hyperparameters  $\theta$  e.g. global scale, lengthscale, etc.

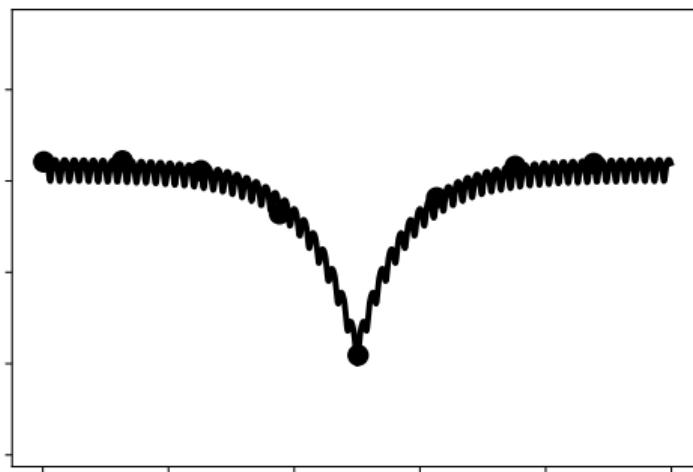
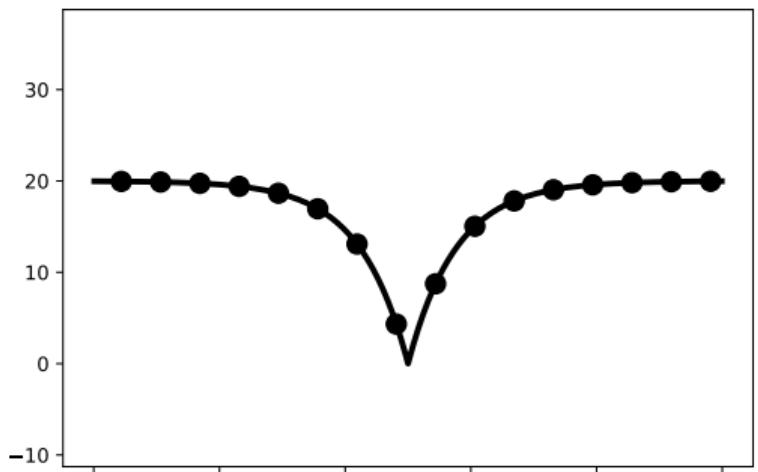
Hyperparameters  $\theta$  often chosen to minimize negative marginal log likelihood (NMLL)

$$\text{NMLL} \propto \mathbf{f}^T \mathbf{K}^{-1} \mathbf{f} + \log|\mathbf{K}| + \log(2\pi)N$$

$\therefore$  MTGP fitting requires computing  $\mathbf{K}^{-1} \mathbf{f}$  and  $\log|\mathbf{K}| \implies$  standard cost  $\mathcal{O}(N^3)$

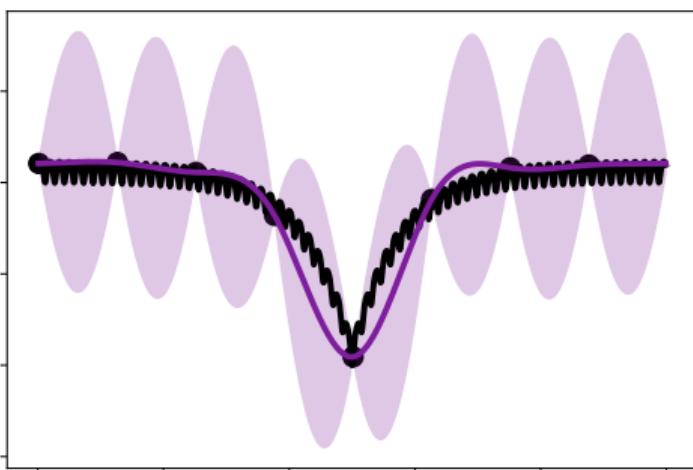
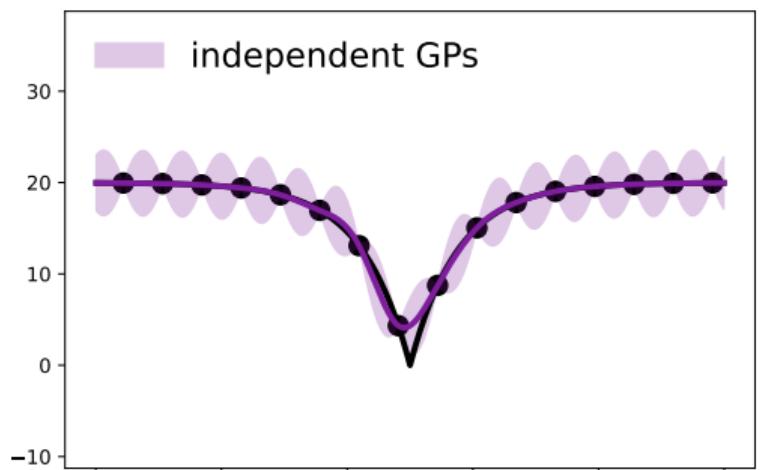
# Independent GPs vs a Multi-Task GP Example

Low Fidelity Left, High Fidelity Right



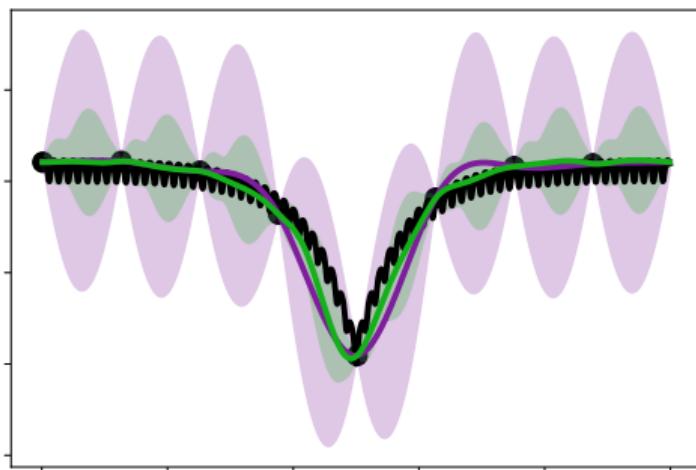
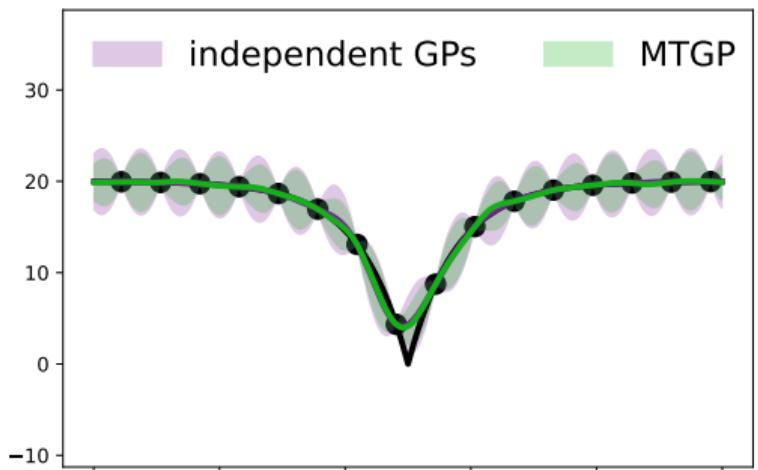
# Independent GPs vs a Multi-Task GP Example

Low Fidelity Left, High Fidelity Right



# Independent GPs vs a Multi-Task GP Example

Low Fidelity Left, High Fidelity Right



# Multilevel Monte Carlo (MLMC) with MTGPs

Quantity of interest

$$\tilde{f}_\ell = \int f(\ell, \mathbf{x})$$

In our Bayesian setting, the posterior cubature mean and covariance are

$$\mathbb{E}\left[\tilde{f}_\ell\right] = \widetilde{\mathbf{K}}_\ell^T \mathbf{K}^{-1} \mathbf{f}$$

$$\mathbb{C}\left[\tilde{f}_\ell, \tilde{f}_{\ell'}\right] = \widetilde{K}_{\ell\ell'} - \widetilde{\mathbf{K}}_\ell^T \mathbf{K}^{-1} \widetilde{\mathbf{K}}_{\ell'}$$

- Integrals understood to be over  $[0, 1]^d$  with respect to  $\mathbf{x}$
  - $\widetilde{\mathbf{K}}_\ell = \int \mathbf{K}(\ell, \mathbf{x})$  a length  $N$  vector
  - Scalar  $\widetilde{K}_{\ell\ell'} = \iint K((\ell, \mathbf{x}), (\ell', \mathbf{x}')) d\mathbf{x} d\mathbf{x}'$
- . $\therefore$  Quantity of interest

$$\tilde{f}_L \sim \mathcal{N}\left(\widetilde{\mathbf{K}}_L^T \mathbf{K}^{-1} \mathbf{f}, \widetilde{K}_{LL} - \widetilde{\mathbf{K}}_L^T \mathbf{K}^{-1} \widetilde{\mathbf{K}}_L\right)$$

# Product Kernels for Multi-Task GPs

Common to assume

$$K((\ell, \mathbf{x}), (\ell', \mathbf{x}')) = R(\ell, \ell') Q(\mathbf{x}, \mathbf{x}')$$

- $R : \{1, \dots, L\} \times \{1, \dots, L\} \rightarrow \mathbb{R}$  an SPD kernel over levels e.g.

$$\mathbf{R} = \{R(\ell, \ell')\}_{\ell, \ell'=1}^L = \mathbf{B}\mathbf{B}^T + \text{diag}(\boldsymbol{\nu}), \quad \boldsymbol{\nu} \in \mathbb{R}_+^L, \quad \mathbf{B} \in \mathbb{R}^{L \times r}, \quad \text{rank } r \leq L$$

- $Q : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$  an SPD kernel over parameters

$$\mathbf{K} = \begin{pmatrix} \mathbf{K}_{11} & \cdots & \mathbf{K}_{1L} \\ \vdots & \ddots & \vdots \\ \overline{\mathbf{K}_{1L}} & \cdots & \mathbf{K}_{LL} \end{pmatrix} = \begin{pmatrix} R_{11} \mathbf{Q}_{11} & \cdots & R_{1L} \mathbf{Q}_{1L} \\ \vdots & \ddots & \vdots \\ \overline{R_{1L} \mathbf{Q}_{1L}} & \cdots & R_{LL} \mathbf{Q}_{LL} \end{pmatrix}$$

- $R_{\ell\ell'} = R(\ell, \ell')$  is a scalar
- $\mathbf{Q}_{\ell\ell'} = Q(\mathcal{D}_\ell, \mathcal{D}_{\ell'}^T)$  is an  $N_\ell \times N_{\ell'}$  Gram matrix

## Fast MTGPs

$$\mathbf{K} = \begin{pmatrix} R_{11}\mathbf{Q}_{11} & \cdots & R_{1L}\mathbf{Q}_{1L} \\ \vdots & \ddots & \vdots \\ \overline{R_{1L}\mathbf{Q}_{1L}} & \cdots & R_{LL}\mathbf{Q}_{LL} \end{pmatrix}$$

**Idea:** Force “nice” structure in  $\mathbf{Q}_{\ell\ell'}$  through special pairings of  $X_\ell = \{\mathbf{x}_{\ell i}\}_{i=1}^{N_\ell}$  and  $Q$

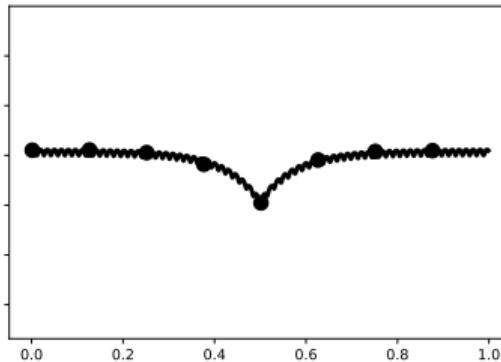
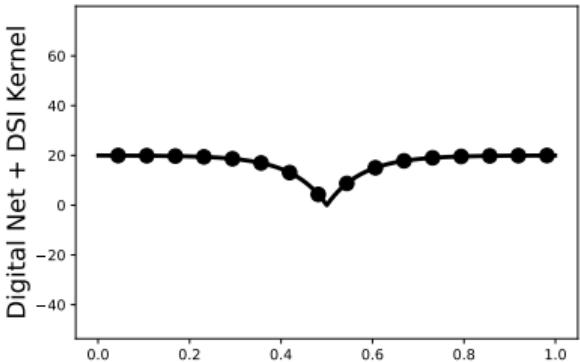
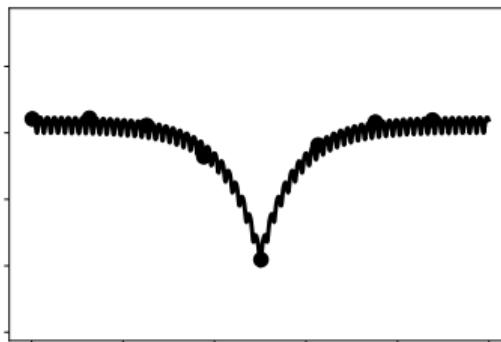
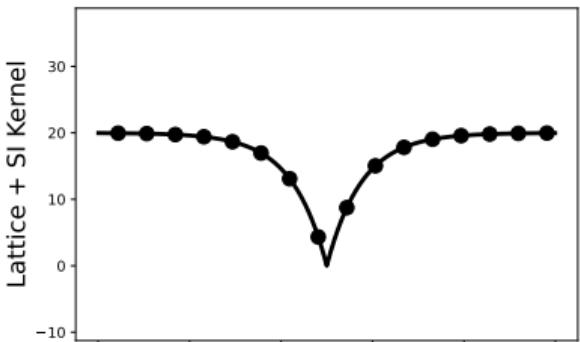
1.  $X_\ell$  a lattice and  $Q$  a shift invariant (SI) kernel  $\implies \mathbf{Q}_{\ell\ell'}$  block circulant
2.  $X_\ell$  a (base 2) digital net and  $Q$  a digitally SI (DSI) kernel  $\implies \mathbf{Q}_{\ell\ell'}$  block RSBT

### Technicalities

- Lattices and digital nets require sample sizes  $N_\ell = 2^{m_\ell}$
- Lattices  $X_1, \dots, X_L$ : same generating vector, possibly different random shifts
- Circulant matrices diagonalizable by FFT
- Digital nets  $X_1, \dots, X_L$ : same generating matrices, possibly different digital shifts
- RSBT matrices diagonalizable by FWHT

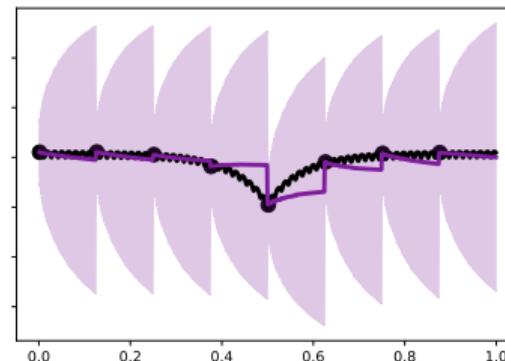
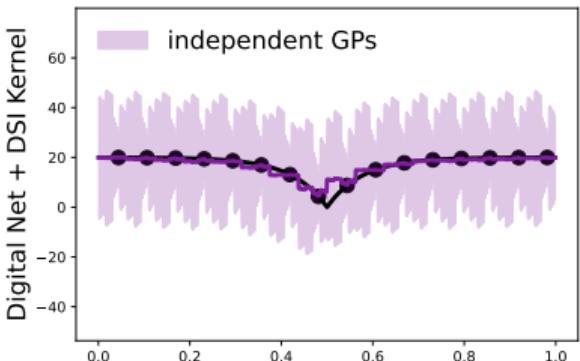
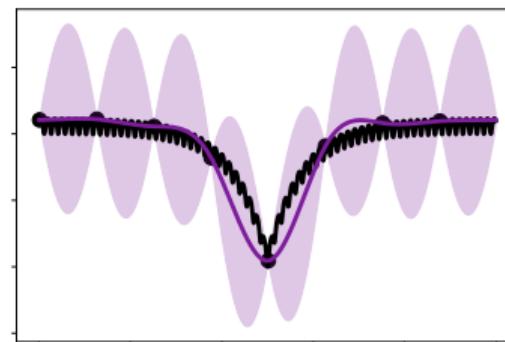
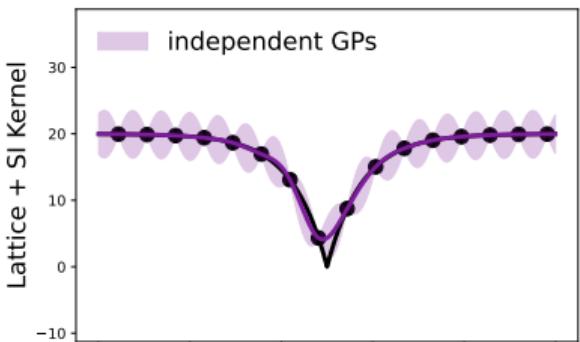
# Independent Fast GPs vs Fast MTGPs Example

Low Fidelity Left, High Fidelity Right



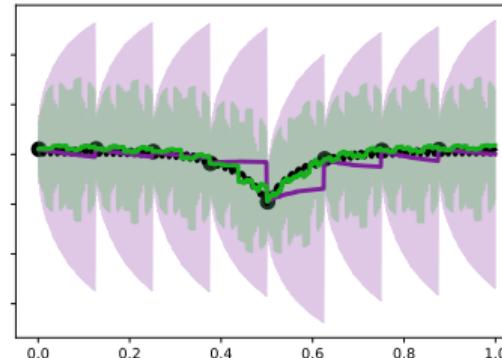
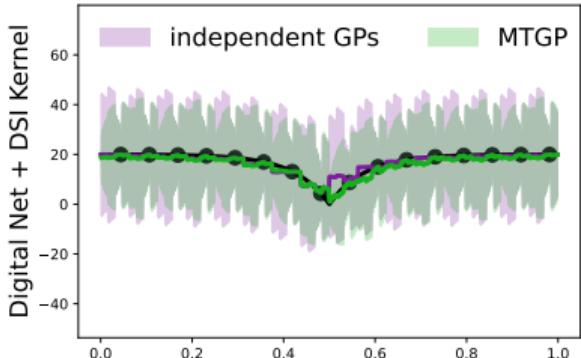
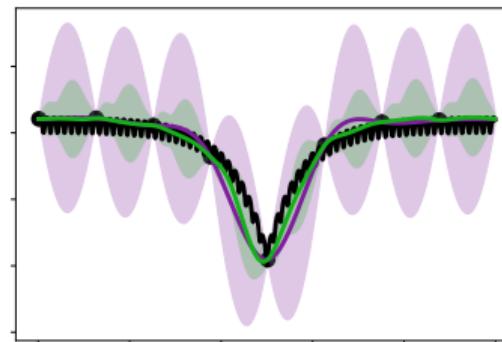
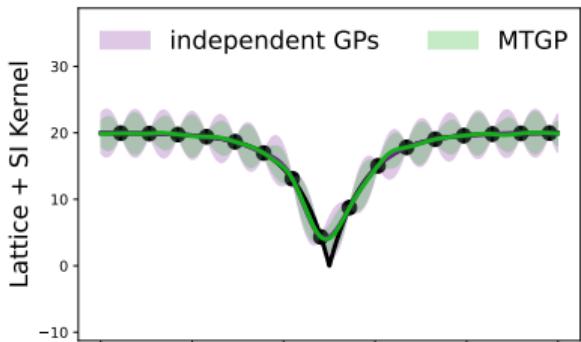
# Independent Fast GPs vs Fast MTGPs Example

Low Fidelity Left, High Fidelity Right



# Independent Fast GPs vs Fast MTGPs Example

Low Fidelity Left, High Fidelity Right



## Fast MTGPs Continued

$$\mathbf{K}_{\ell\ell'} = \mathbf{R}_{\ell\ell'} \mathbf{Q}_{\ell\ell'} = \mathbf{V}_{m_\ell} \boldsymbol{\Sigma}_{\ell\ell'} \overline{\mathbf{V}_{m_{\ell'}}}$$

- $\overline{\mathbf{V}_m}$  a  $2^m \times 2^m$  *fast transform matrix*
  1. Lattice  $X_\ell$  with SI  $Q$  makes  $\overline{\mathbf{V}_{m_\ell}}$  the Fast Fourier Transform
  2. Digital Net  $X_\ell$  with DSI  $Q$  makes  $\overline{\mathbf{V}_{m_\ell}}$  the Fast Walsh Hadamard Transform
- $\mathbf{V}_m \mathbf{a}$  and  $\overline{\mathbf{V}_m} \mathbf{a}$  both cost only  $\mathcal{O}(m2^m)$  to compute
- The first column of  $\overline{\mathbf{V}_m}$  is  $\mathbf{1}_m / \sqrt{2^m}$
- $\boldsymbol{\Sigma}_{\ell\ell'}$  a diagonal block matrix characterized by

$$\sigma_{\ell\ell'} = \boldsymbol{\Sigma}_{\ell\ell'} \mathbf{1}_{m_{\ell'}} = \sqrt{2^{m_{\ell'}}} \overline{\mathbf{V}_{m_\ell}} \mathbf{k}_{\ell\ell',1}$$

where  $\mathbf{k}_{\ell\ell',1}$  is the first column of  $\mathbf{K}_{\ell\ell'}$  and we assume  $m_\ell \geq m_{\ell'}$

## Fast MTGPs NMLL

$$\mathbf{K} = \begin{pmatrix} \mathbf{V}_{m_1} & & \\ & \ddots & \\ & & \mathbf{V}_{m_L} \end{pmatrix} \begin{pmatrix} \Sigma_{11} & \cdots & \Sigma_{1L} \\ \vdots & \ddots & \vdots \\ \overline{\Sigma_{1L}} & \cdots & \Sigma_{LL} \end{pmatrix} \begin{pmatrix} \sqrt{\mathbf{V}_{m_1}} & & \\ & \ddots & \\ & & \sqrt{\mathbf{V}_{m_L}} \end{pmatrix} =: \mathbf{V} \Sigma \mathbf{\bar{V}}$$

$$\text{NMLL} \propto \hat{\mathbf{f}}^T \Sigma^{-1} \hat{\mathbf{f}} + \log|\Sigma| + \log(2\pi)N, \quad \hat{\mathbf{f}} = \mathbf{\bar{V}} \mathbf{f}$$

∴ Fast MTGP fitting requires computing  $\Sigma^{-1} \hat{\mathbf{f}}$  and  $\log|\Sigma|$

For example, if  $N_1 = 8$ ,  $N_2 = 4$ , and  $N_3 = 2$  then  $\Sigma$  has the following structure

$$\Sigma = \left[ \begin{array}{c|c|c|c} \cdot & & & \cdot \\ \hline & \cdot & & \cdot \\ \hline & & \cdot & \cdot \\ \hline & & & \cdot \\ \hline \cdot & & & \cdot \\ \hline & \cdot & & \cdot \\ \hline & & \cdot & \cdot \\ \hline & & & \cdot \end{array} \right]$$

## Fast MTGPs Storage and Costs

- Assume evaluating  $f(\ell, \mathbf{x})$  costs  $C_\ell$  for any  $\mathbf{x} \in [0, 1]^d$
- Assume evaluating  $K((\ell, \mathbf{x}), (\ell', \mathbf{x}'))$  costs  $\mathcal{O}(d)$
- Assume  $m_1 \geq \dots \geq m_L$  i.e. less samples on higher levels (or reorder levels)

$$\text{NMLL} \propto \hat{\mathbf{f}}^T \Sigma^{-1} \hat{\mathbf{f}} + \log|\Sigma| + \log(2\pi)N$$

- Evaluate  $\mathbf{f} = f(\mathcal{D})$  at cost  $\mathcal{O}\left(\sum_{\ell=1}^L C_\ell 2^{m_\ell}\right)$
- Evaluate  $\hat{\mathbf{f}} = \bar{V}\mathbf{f}$  at cost  $\mathcal{O}\left(\sum_{\ell=1}^L m_\ell 2^{m_\ell}\right)$
- Evaluate only first columns of  $K_{\ell\ell'}$  at total cost  $\mathcal{O}\left(d \sum_{\ell=1}^L (L - \ell + 1) 2^{m_\ell}\right)$
- Evaluate  $\Sigma$  at cost  $\mathcal{O}\left(\sum_{\ell=1}^L (L - \ell + 1) m_\ell 2^{m_\ell}\right)$
- Store  $\Sigma$  in  $\mathcal{O}\left(\sum_{\ell=1}^L (L - \ell + 1) 2^{m_\ell}\right)$  (possibly complex) floats
- Evaluate  $\Sigma^{-1} \hat{\mathbf{f}}$  and  $\log|\Sigma|$  at cost  $\mathcal{O}\left(\sum_{\ell'=1}^L 2^{-m_{\ell'}} \left[ \sum_{\ell=1}^{\ell'-1} 2^{m_\ell} \right]^2\right)$

▶ algorithm

# Python Software

QMCPy [[github.com/QMCSwift/QMCSwift](https://github.com/QMCSwift/QMCSwift)]

- Low discrepancy sequences including lattices and digital nets
- Measures with automatic transforms to integrals over  $[0, 1]^d$
- Adaptive stopping criteria for IID Monte Carlo and QMC

FastGaussianProcesses [[github.com/alegresor/FastGaussianProcesses](https://github.com/alegresor/FastGaussianProcesses)]

- Fast GPs
- Fast MTGPs
- Uses PyTorch for efficient optimizing and GPU compatibility
- Efficient MTGP updates when increasing sample sizes via caching

MLQMCPy [[github.com/PieterjanRobbe/mlqmcpy](https://github.com/PieterjanRobbe/mlqmcpy)]

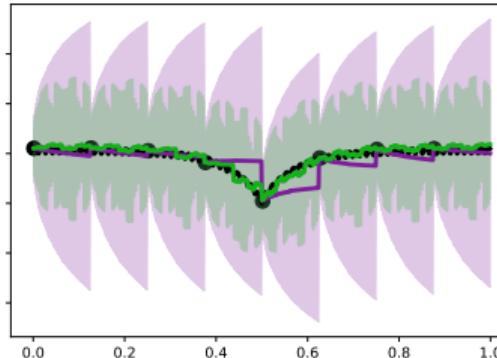
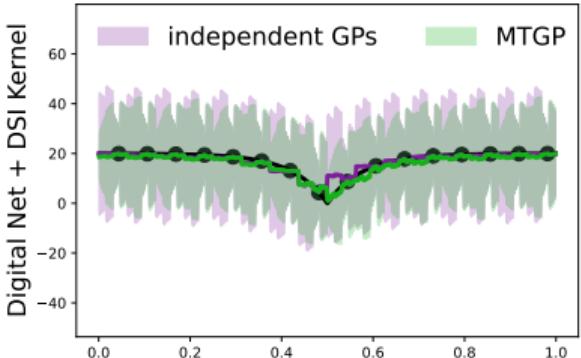
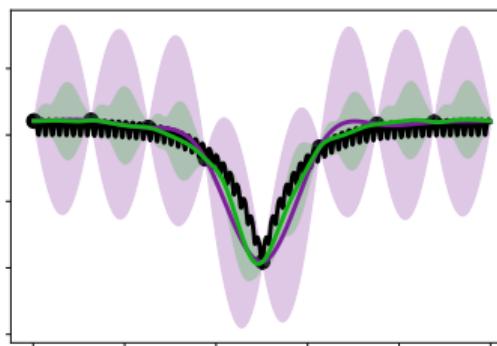
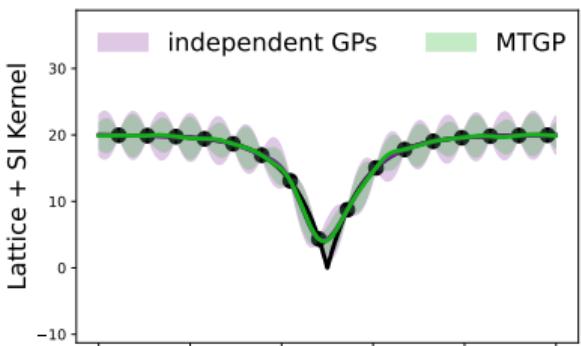
- Multilevel IID Monte Carlo algorithms
- Multilevel QMC algorithms

QMC  
ooooGPs  
ooFast GPs  
ooooLattices + SI  $K$   
ooooooDigital Nets + DSI  $K$   
ooooooooMulti-Task GPs  
ooooooFast MTGPs  
ooooooPython Software  
o●

References

Appendix  
oooooo

# Thank you for listening!



## References I

- Josef Dick. Walsh spaces containing smooth functions and quasi-monte carlo rules of arbitrary high order. *SIAM Journal on Numerical Analysis*, 46(3):1519–1553, 2008.
- Josef Dick. The decay of the walsh coefficients of smooth functions. *Bulletin of the Australian Mathematical Society*, 80(3):430–453, 2009.
- Josef Dick and Friedrich Pillichshammer. Multivariate integration in weighted hilbert spaces based on walsh functions and weighted sobolev spaces. *Journal of Complexity*, 21(2):149–195, 2005.
- Josef Dick and Friedrich Pillichshammer. *Digital nets and sequences: discrepancy theory and quasi-Monte Carlo integration*. Cambridge University Press, 2010.
- Josef Dick, Frances Y Kuo, and Ian H Sloan. High-dimensional integration: the quasi-monte carlo way. *Acta Numerica*, 22:133–288, 2013.
- Nathan Jacob Fine. On the walsh functions. *Transactions of the American Mathematical Society*, 65(3):372–414, 1949.

## References II

Fino and Algazi. Unified matrix treatment of the fast walsh-hadamard transform. *IEEE Transactions on Computers*, 100(11):1142–1146, 1976.

Fred J. Hickernell and Lluís Antoni Jiménez Rugama. Reliable adaptive cubature using digital sequences, 2014.

Fred J Hickernell, Lan Jiang, Yuewei Liu, and Art B Owen. Guaranteed conservative fixed width confidence intervals via monte carlo sampling. In *Monte Carlo and Quasi-Monte Carlo Methods 2012*, pages 105–128. Springer, 2013.

Fred J. Hickernell, Lluís Antoni Jiménez Rugama, and Da Li. Adaptive quasi-monte carlo methods for cubature, 2017.

Lluís Antoni Jiménez Rugama and Fred J. Hickernell. Adaptive multidimensional integration based on rank-1 lattices, 2014.

B. D. Keister. Multidimensional quadrature algorithms. *Computers in Physics*, 10: 119–122, 1996. doi: 10.1063/1.168565.

## References III

Gerhard Larcher, Harald Niederreiter, and Wolfgang Ch Schmid. Digital nets and sequences constructed over finite rings and their application to quasi-monte carlo integration. *Monatshefte für Mathematik*, 121:231–253, 1996.

Pierre L'Ecuyer, Marvin K Nakayama, Art B Owen, and Bruno Tuffin. Confidence intervals for randomized quasi-monte carlo estimators. In *2023 Winter Simulation Conference (WSC)*, pages 445–456. IEEE, 2023.

Art B Owen. Variance and discrepancy with alternative scramblings. *ACM Transactions of Modeling and Computer Simulation*, 13(4), 2003.

Art B Owen. Error estimation for quasi-monte carlo. *arXiv preprint arXiv:2501.00150*, 2024.

Jagadeeswaran Rathinavel. *Fast automatic Bayesian cubature using matching kernels and designs*. Illinois Institute of Technology, 2019.

## References IV

Jagadeeswaran Rathinavel and Fred J. Hickernell. Fast automatic bayesian cubature using lattice sampling. *Statistics and Computing*, 29(6):1215–1229, Sep 2019. ISSN 1573-1375. doi: 10.1007/s11222-019-09895-9. URL <http://dx.doi.org/10.1007/s11222-019-09895-9>.

Jagadeeswaran Rathinavel and Fred J Hickernell. Fast automatic bayesian cubature using sobol' sampling. In *Advances in Modeling and Simulation: Festschrift for Pierre L'Ecuyer*, pages 301–318. Springer, 2022.

S. Surjanovic and D. Bingham. Virtual library of simulation experiments: Test functions and datasets. Retrieved April 9, 2025, from <http://www.sfu.ca/~ssurjano>.

Joseph L Walsh. A closed set of normal orthogonal functions. *American Journal of Mathematics*, 45(1):5–24, 1923.

## Walsh Functions

Introduced for base  $b = 2$  in [Walsh, 1923] with important results in [Fine, 1949]. Generalized to finite abelian group with a bijection in [Larcher et al., 1996].

For  $k \in \mathbb{N}_0$  with  $\mathbf{k} = (k_0, k_1, \dots)$  and  $x \in [0, 1)$  with  $\mathbf{x} = (x_1, x_2, \dots)$ ,

$$\text{wal}_k(x) = e^{2\pi i/b \sum_{\ell=0}^{\infty} k_{\ell} x_{\ell+1}} = e^{2\pi i/b \mathbf{k} \cdot \mathbf{x}}$$

e.g. for  $b = 2$ ,  $\text{wal}_6(.75) = (-1)^{(0,1,1).(1,1,0)} = -1$ .

For any fixed  $b$ ,  $\{\text{wal}_k : k \in \mathbb{N}_0\}$  is a complete orthonormal system in  $\mathcal{L}_2([0, 1))$ . Notice similarity to complex exponential basis  $\{e^{2\pi i k x} : k \in \mathbb{Z}\}$  for Fourier series.

## Walsh Function Properties

For any  $x, y \in [0, 1)$  and  $k, h \in \mathbb{N}_0$  and  $f \in \mathcal{L}_2([0, 1))$

1.  $\text{wal}_k(x)\text{wal}_h(x) = \text{wal}_{k\ominus h}(x)$  and  $\text{wal}_k(x)\text{wal}_k(y) = \text{wal}_k(x\ominus y)$

2.

$$\int_0^1 \text{wal}_k(x)dx = \begin{cases} 1, & k = 0 \\ 0, & k > 0 \end{cases}$$

3.

$$\int_0^1 f(\sigma)d\sigma = \int_0^1 f(x\ominus\sigma)d\sigma$$

4.

$$\sum_{k=0}^{b^a-1} \text{wal}_k(x) = \begin{cases} b^a, & a < \mathcal{I}(x) - 1 \\ 0, & \text{otherwise} \end{cases}$$

where  $\mathcal{I}(x) = -[\log_b(x)]$  is the index first non-zero digit in the base  $b$  expansion  
e.g. with  $b = 2$  then  $\mathcal{I}(.375) = \mathcal{I}(.011_2) = 2$ .

## Weight Function

Write  $k \in \mathbb{N}$  as

$$k = \sum_{\ell=1}^{\#k} k_{a_\ell} b^{a_\ell}$$

with  $a_1 > \dots > a_{\#k} \geq 0$  and  $k_{a_\ell} \in \{1, \dots, b-1\}$ .

Weight function for  $\alpha \in \mathbb{N}_0$  has

$$\mu_\alpha(k) = \sum_{\ell=1}^{\min(\alpha, \#k)} (a_\ell + 1)$$

with  $\mu_0(k) = \mu_\alpha(0) = 0$ .  $\mu$  sums indices of non-zero digits. For example, with  $b = 2$

$$k = 13 = 1101_2 \quad \text{has} \quad (a_1, a_2, a_3) = (3, 2, 0)$$

$$\mu_1(k) = (3+1), \quad \mu_2(k) = (3+1) + (2+1), \quad \mu_3(k) = (3+1) + (2+1) + (0+1) = \mu_4(k) = \dots$$

## Walsh Series of Smooth Functions

For  $\alpha \geq 2$  the Sobolev RKHS  $H^\alpha$  with inner product

$$\langle f, g \rangle_\alpha = \sum_{\beta=1}^{\alpha-1} \int_0^1 f^{(\beta)}(x) dx \int_0^1 g^{(\beta)}(x) dx + \int_0^1 f^{(\alpha)}(x) g^{(\alpha)}(x) dx$$

has kernel

$$K_\alpha(x, x') = \sum_{\beta=1}^{\alpha-1} \frac{B_\beta(x) B_\beta(x')}{(\beta!)^2} + \overbrace{(-1)^{\alpha+1} \frac{B_{2\alpha}(\{x-x'\})}{(2\alpha)!}}^{\hat{K}_\alpha((x-x') \bmod 1)/(2\pi)^{2\alpha}}.$$

[Dick, 2008, 2009] show that if  $f \in H^\alpha$  then for  $\hat{f}(k) = \int_0^1 f(x) \overline{\text{wal}_k(x)} dx$  we have

$$\sup_{k \in \mathbb{N}_0} |\hat{f}(k)| b^{\mu_\alpha(k)} < \infty \quad \text{i.e.} \quad \exists C_{f,\alpha} > 0 \quad \text{s.t.} \quad |\hat{f}(k)| \leq \frac{C_{f,\alpha}}{b^{\mu_\alpha(k)}}.$$

For the  $\alpha = 1$  case see [Dick and Pillichshammer, 2005].

► Return to slide on Fast MTGPs Storage and Costs

---

## Algorithm 1 Inverse and Determinant of $\Sigma$

---

**Require:**  $\Sigma$  diagonal block matrix with  $m_1 \geq \dots \geq m_L$ .

$A \leftarrow \Sigma_{11}^{-1}$  diagonal

$\rho \leftarrow |\Sigma_{11}|$

$\ell \leftarrow 2$

**while**  $\ell \leq L$  **do**

$D \leftarrow \Sigma_{\ell\ell}$  diagonal

$C \leftarrow \Sigma_{1:\ell-1,\ell}$  diagonal blocks

$S_\ell \leftarrow D - \bar{C}A^{-1}C$  diagonal Schur complement

$A \leftarrow \begin{pmatrix} A^{-1} + A^{-1}CS_\ell^{-1}\bar{C}A^{-1} & -A^{-1}CS_\ell^{-1} \\ -S_\ell^{-1}\bar{C}A^{-1} & S_\ell^{-1} \end{pmatrix}$

$\rho \leftarrow \rho|S_\ell|$

$\ell \leftarrow \ell + 1$

**end while**

**return**  $A = \Sigma^{-1}$  and  $\rho = |\Sigma|$

---