

Talk for ETH Zürich

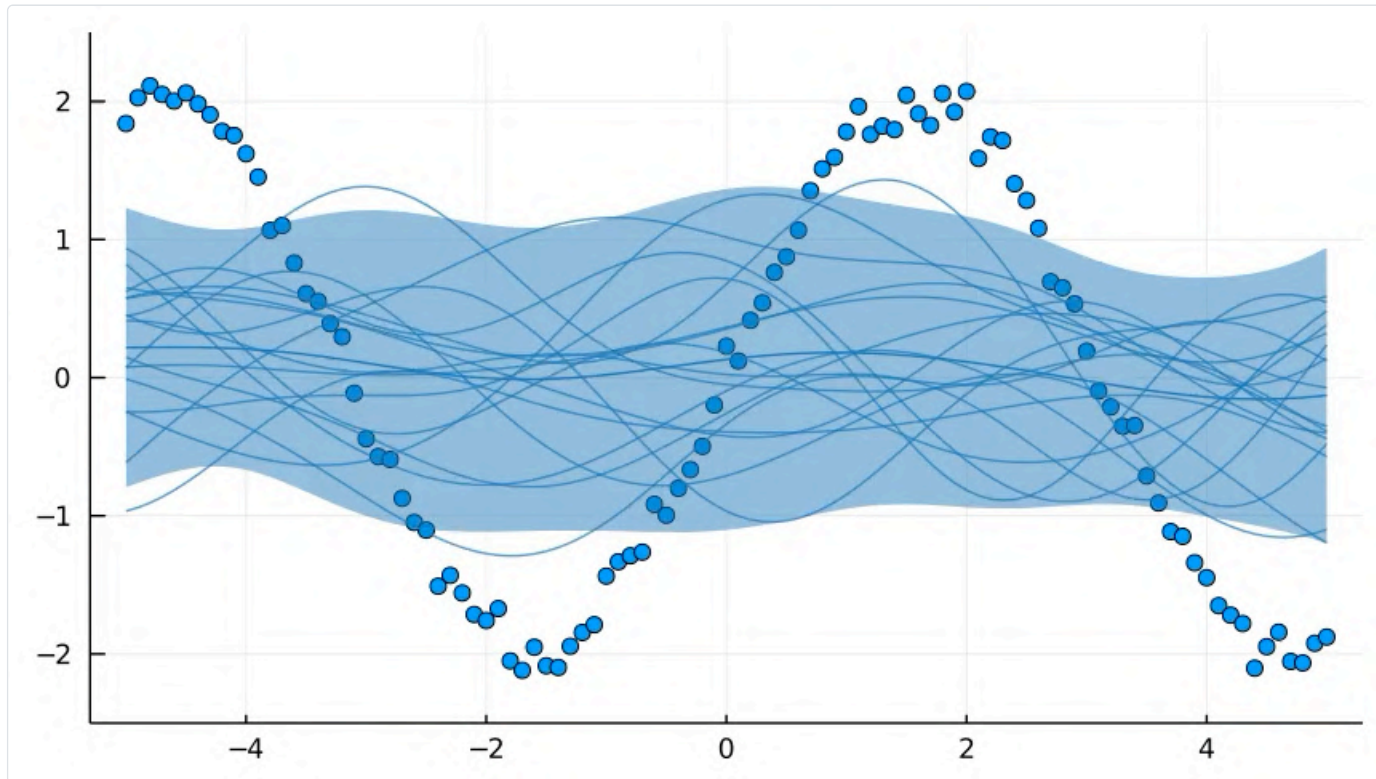
# Sampling from Gaussian Process Posteriors using Stochastic Gradient Descent



UNIVERSITY OF  
CAMBRIDGE

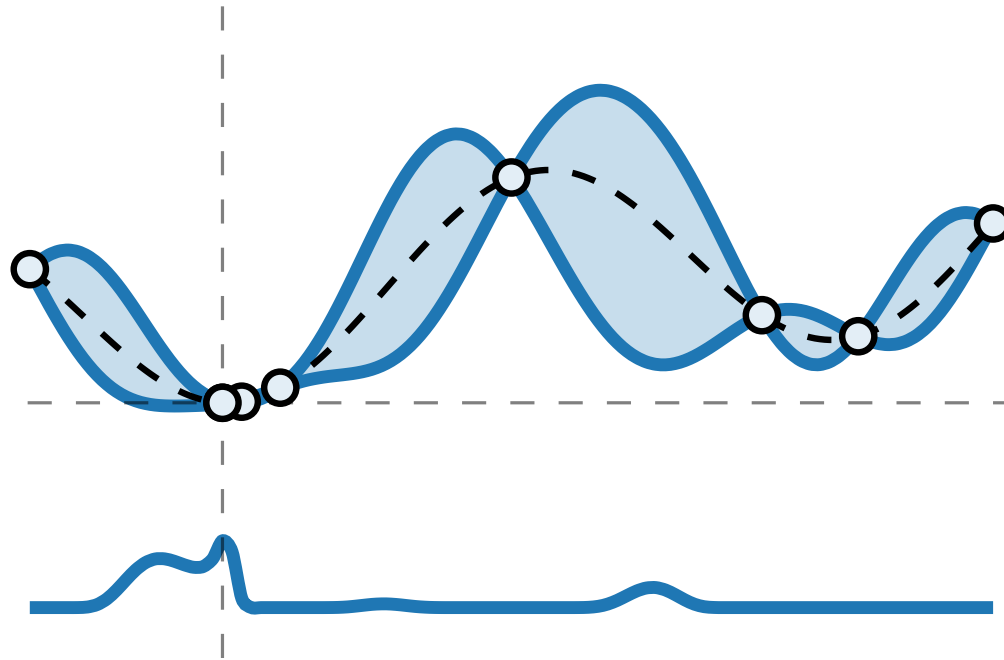
Alexander Terenin  
[HTTPS://AVT.IM/](https://AVT.IM/) ·  @AVT\_IM

# Gaussian Processes



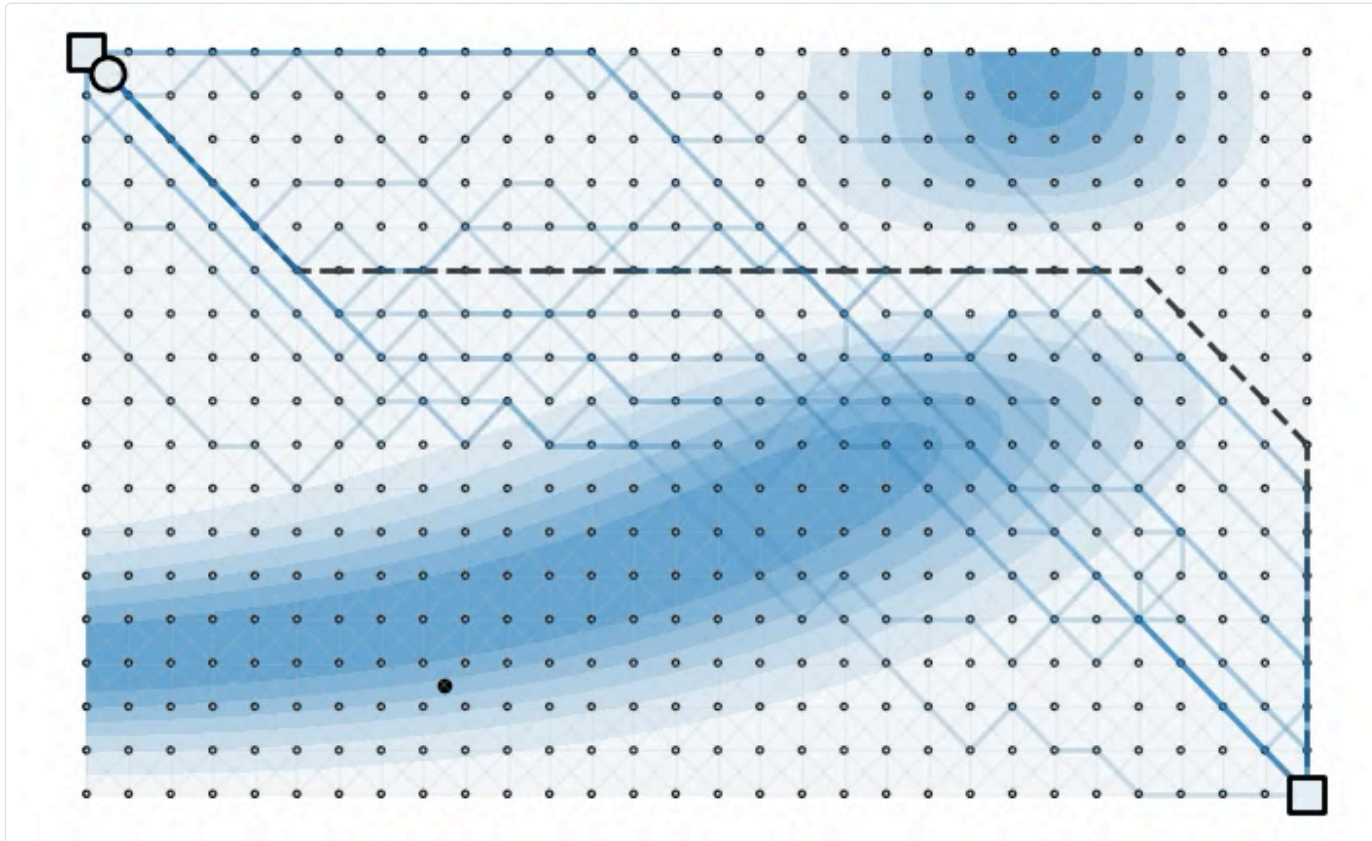
Probabilistic formulation provides uncertainty

# Bayesian Optimization

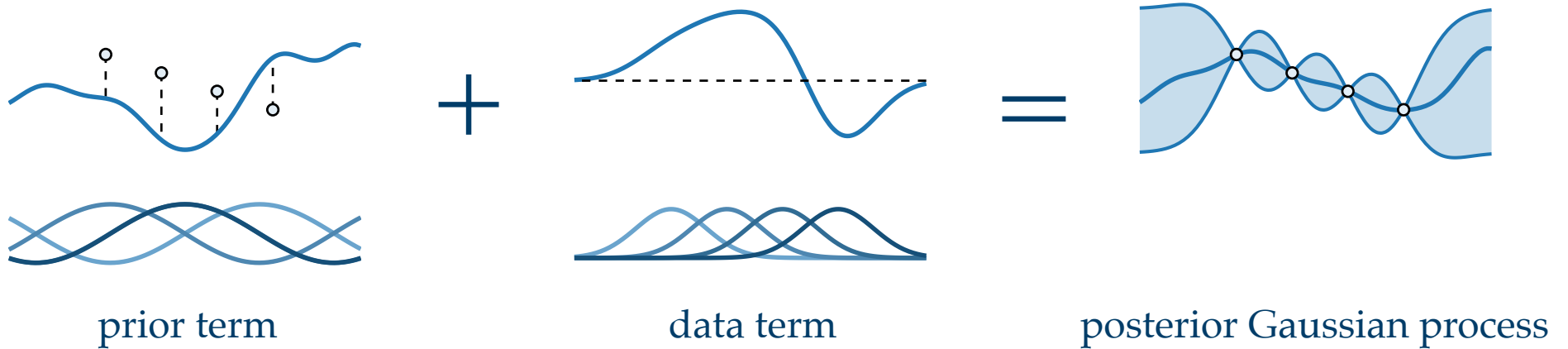


*Automatic explore-exploit tradeoff*

# From Bayesian Optimization to Bayesian Interactive Decision-making



# Pathwise Conditioning

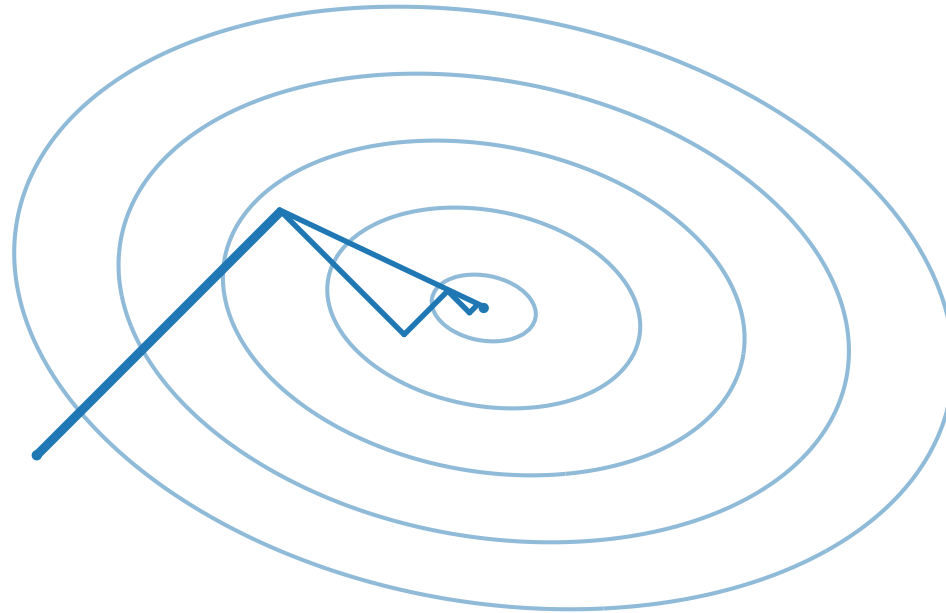


$$(f \mid \mathbf{y})(\cdot) = f(\cdot) + \sum_{i=1}^N v_i k(x_i, \cdot) \quad \mathbf{v} = \mathbf{K}_{xx}^{-1}(\mathbf{y} - f(\mathbf{x}))$$

$\mathbf{v}$ : representer weights

$k(x_i, \cdot)$ : canonical basis functions

# Conjugate Gradients



Refinement of gradient descent for solving linear systems  $\mathbf{A}^{-1}\mathbf{b}$

Convergence rate is much faster than gradient descent

Precise rate depends mainly on  $\text{cond}(\mathbf{A})$

# Numerical Stability

Condition number: quantifies difficulty of solving  $\mathbf{A}^{-1}\mathbf{b}$

$$\text{cond}(\mathbf{A}) = \lim_{\varepsilon \rightarrow 0} \sup_{\|\boldsymbol{\delta}\| \leq \varepsilon \|\mathbf{b}\|} \frac{\|\mathbf{A}^{-1}(\mathbf{b} + \boldsymbol{\delta}) - \mathbf{A}^{-1}\mathbf{b}\|_2}{\varepsilon \|\mathbf{A}^{-1}\mathbf{b}\|_2} = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})}$$

$\lambda_{\min}, \lambda_{\max}$ : eigenvalues

## Condition Numbers of Kernel Matrices

Are kernel matrices always well-conditioned? **No.**

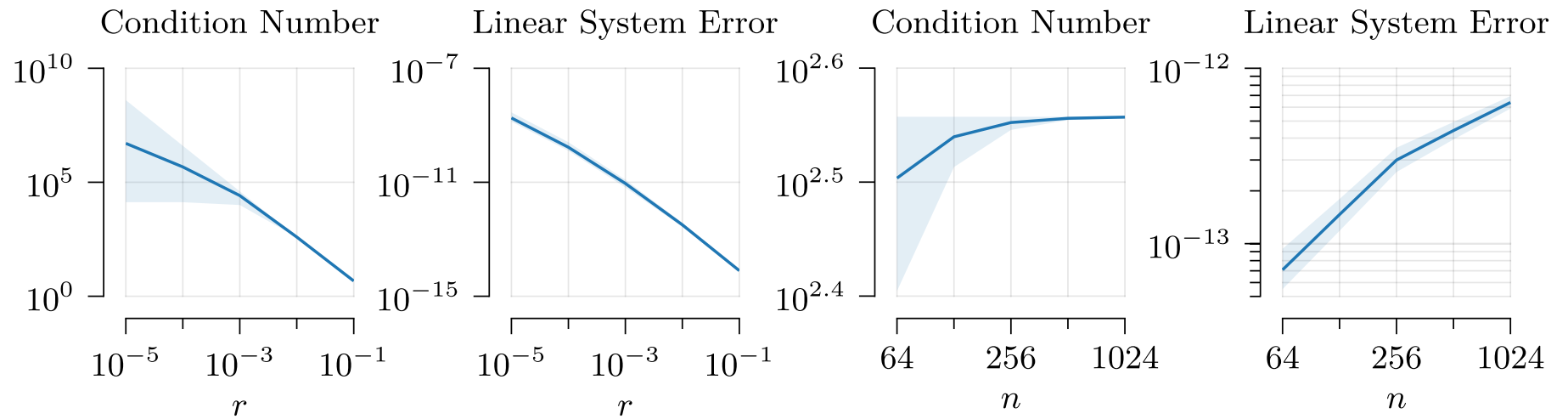
One-dimensional time series on grid: Kac–Murdock–Szegő matrix

$$\mathbf{K}_{xx} = \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & 1 & \rho & \dots & \rho^{n-2} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \rho^{n-1} & \rho^{n-2} & \rho^{n-3} & \dots & 1 \end{pmatrix}$$

for which  $\frac{1+2\rho+2\rho\varepsilon+\rho^2}{1-2\rho-2\rho\varepsilon+\rho^2} \leq \text{cond}(\mathbf{K}_{xx}) \leq \frac{(1+\rho)^2}{(1-\rho)^2}$ , where  $\varepsilon = \frac{\pi^2}{N+1}$ .

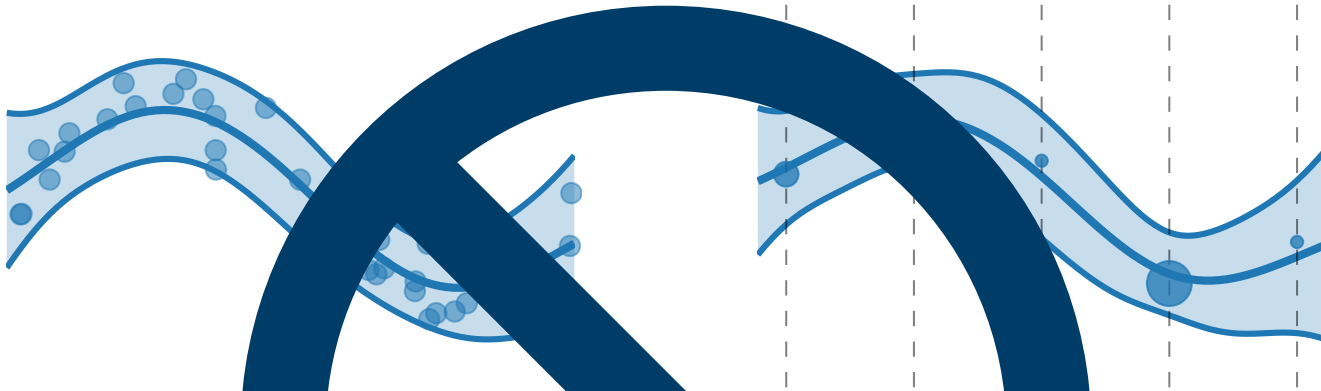


# Condition Numbers of Kernel Matrices



Problem: too much correlation  $\rightsquigarrow$  points too close by

# Minimum Separation



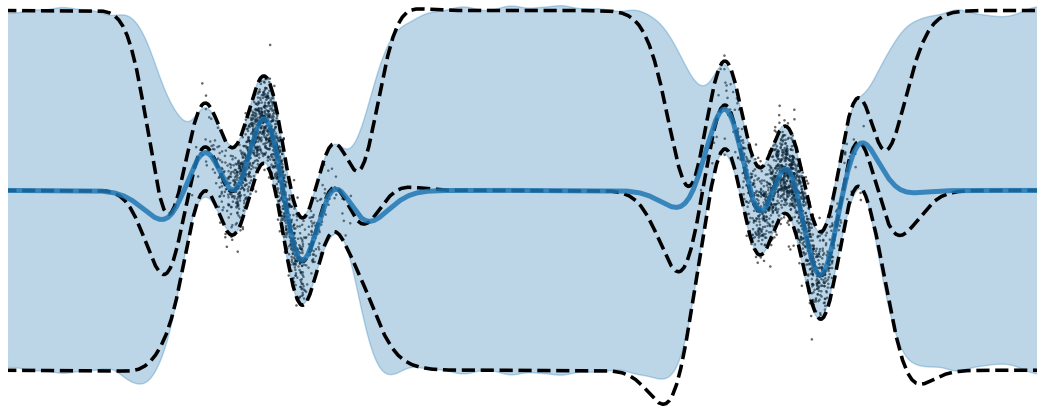
*Separation:* minimum distance between distinct  $z_i$  and  $z_j$

**Proposition.** Assuming spatial decay and stationarity,  $\kappa$  controls  $\text{cond}(\mathbf{K}_{zz})$  uniformly in  $M$ .

Idea: use this to select numerically stable inducing points

# Sampling from Gaussian Process Posteriors using Stochastic Gradient Descent

Jihao Andreas Lin,<sup>\*</sup> Javier Antorán,<sup>\*</sup> Shreyas Padhy,<sup>\*</sup>  
David Janz, José Miguel Hernández-Lobato, Alexander Terenin



<sup>\*</sup>equal contribution

Have you tried stochastic gradient descent?

Conventional wisdom in deep learning:

- SGD variants are empirically often the best optimization algorithms
- ADAM is extremely effective, even on non-convex problems
- Minibatch-based training critical part of scalability

Why not try it out for Gaussian process posterior sample paths?

# Gaussian Process Posteriors via Randomized Optimization Objectives

Split into posterior mean and uncertainty reduction terms

$$\begin{aligned}(f \mid \mathbf{y})(\cdot) &= f(\cdot) + \mathbf{K}_{(\cdot)\mathbf{x}}(\mathbf{K}_{\mathbf{x}\mathbf{x}} + \mathbf{\Sigma})^{-1}(\mathbf{y} - f(\mathbf{x}) - \boldsymbol{\varepsilon}) \\ &= f(\cdot) + \sum_{i=1}^N v_i^* k(x_i, \cdot) + \sum_{i=1}^N \alpha_i^* k(x_i, \cdot)\end{aligned}$$

where

$$\begin{aligned}\mathbf{v}^* &= \arg \min_{\mathbf{v} \in \mathbb{R}^N} \sum_{i=1}^N \frac{(y_i - \mathbf{K}_{x_i \mathbf{x}} \mathbf{v})^2}{\Sigma_{ii}} + \mathbf{v}^T \mathbf{K}_{\mathbf{x}\mathbf{x}} \mathbf{v} \\ \boldsymbol{\alpha}^* &= \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \sum_{i=1}^N \frac{(f(x_i) + \varepsilon_i - \mathbf{K}_{x_i \mathbf{x}} \boldsymbol{\alpha})^2}{\Sigma_{ii}} + \boldsymbol{\alpha}^T \mathbf{K}_{\mathbf{x}\mathbf{x}} \boldsymbol{\alpha}.\end{aligned}$$

# Gaussian Process Posteriors via Randomized Optimization Objectives

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathbb{R}^N} \sum_{i=1}^N \frac{(y_i - \mathbf{K}_{x_i x} \mathbf{v})^2}{\Sigma_{ii}} + \mathbf{v}^T \mathbf{K}_{xx} \mathbf{v}$$
$$\boldsymbol{\alpha}^* = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \sum_{i=1}^N \frac{(f(x_i) + \varepsilon_i - \mathbf{K}_{x_i x} \boldsymbol{\alpha})^2}{\Sigma_{ii}} + \boldsymbol{\alpha}^T \mathbf{K}_{xx} \boldsymbol{\alpha}.$$

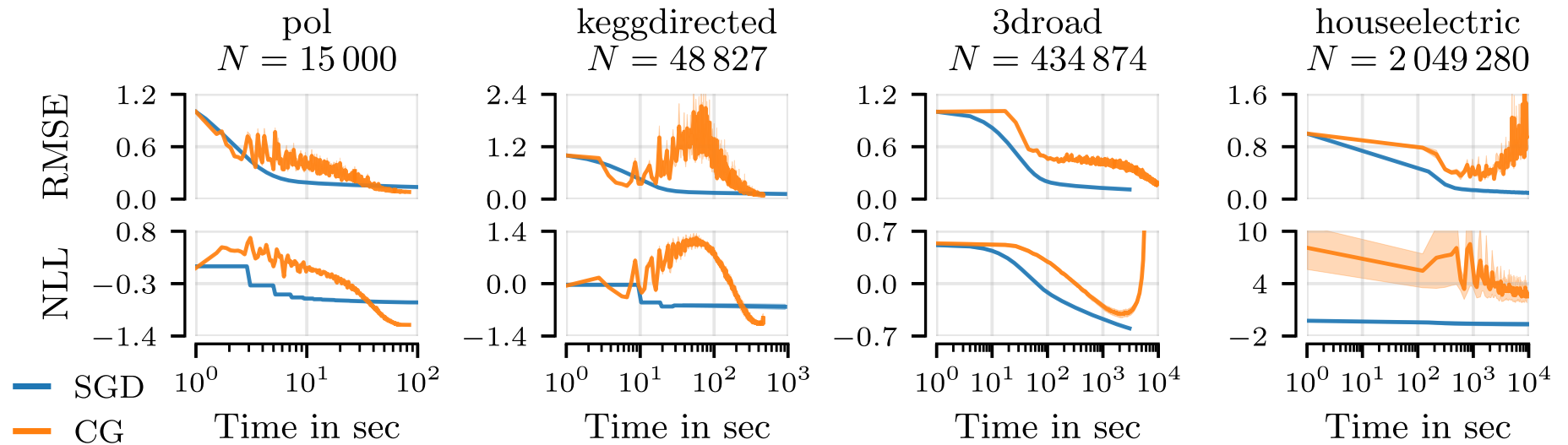
First term: apply mini-batch estimation

Second term: evaluate stochastically via random Fourier features

Variance reduction trick: shift  $\varepsilon_i$  into regularizer

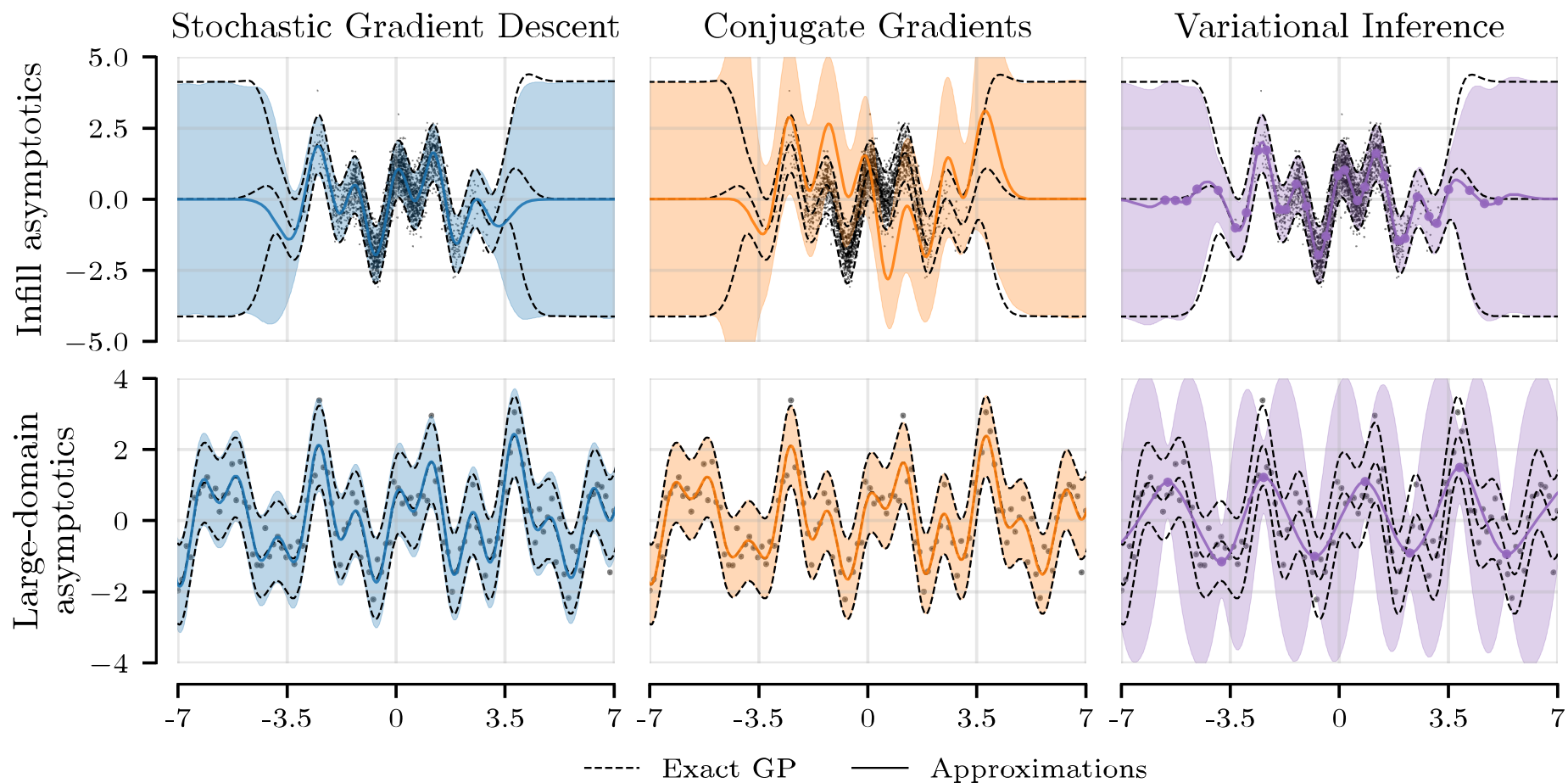
Use stochastic gradient descent with Polyak averaging

# Stochastic Gradient Descent



It works better than conjugate gradients on test data? *Wait, what?*

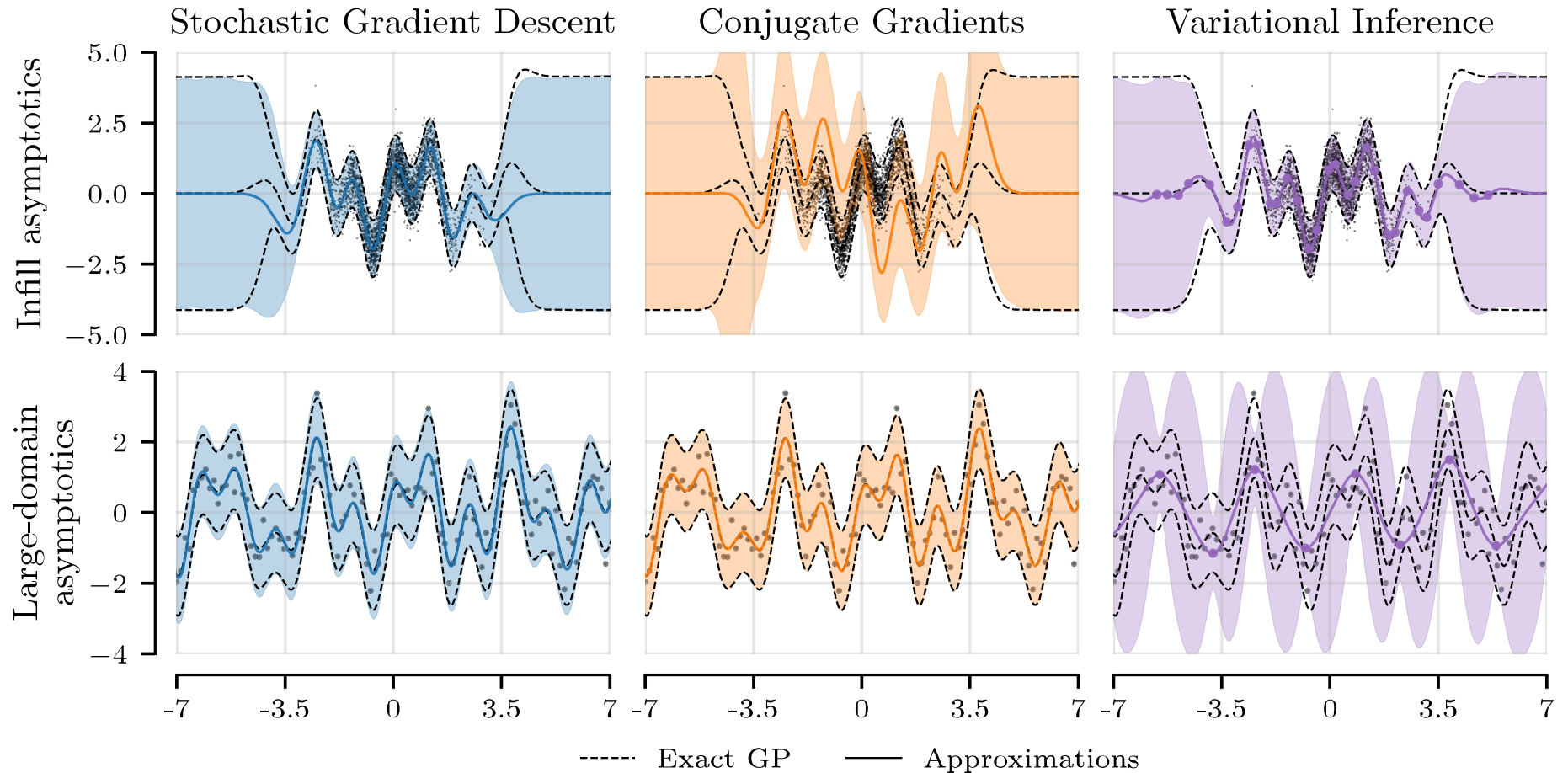
# What happens in one dimension?



Performance depends on data-generation asymptotics

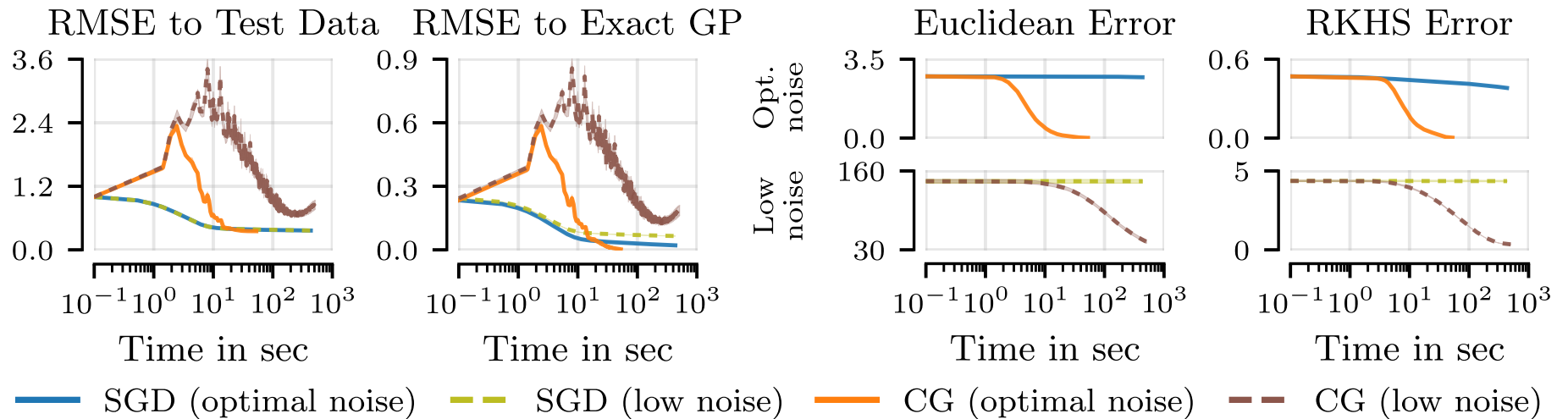


# What happens in one dimension?



SGD does not converge to the correct solution,  $\rightsquigarrow$  implicit bias?  
but still produces reasonable error bars

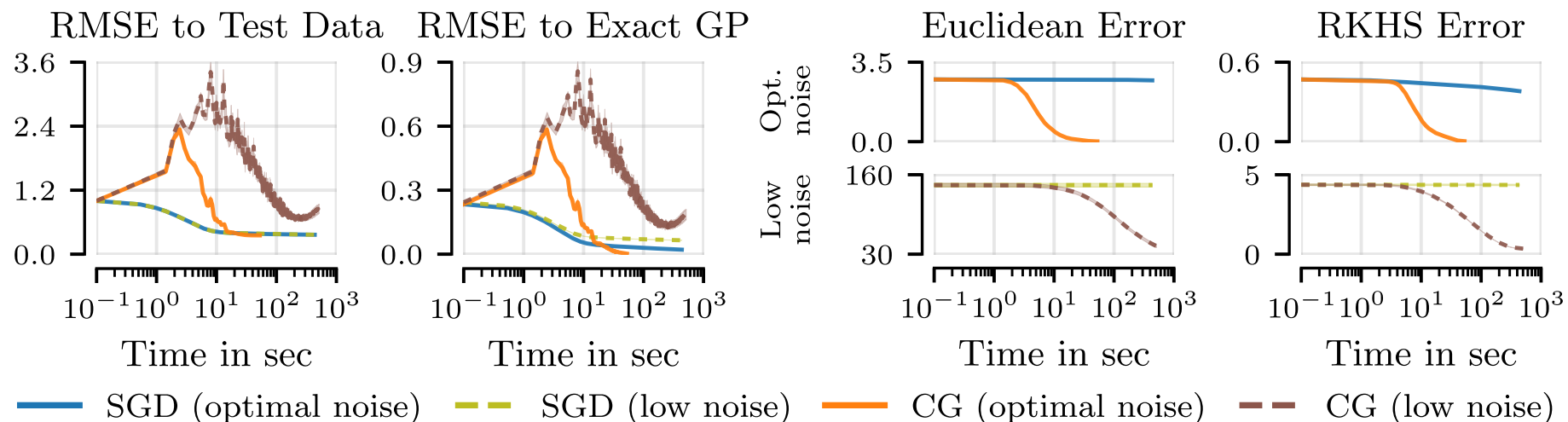
# Convergence: Euclidean and RKHS Norms



No convergence in representer weight space  
or in the reproducing kernel Hilbert space

Good test  
performance

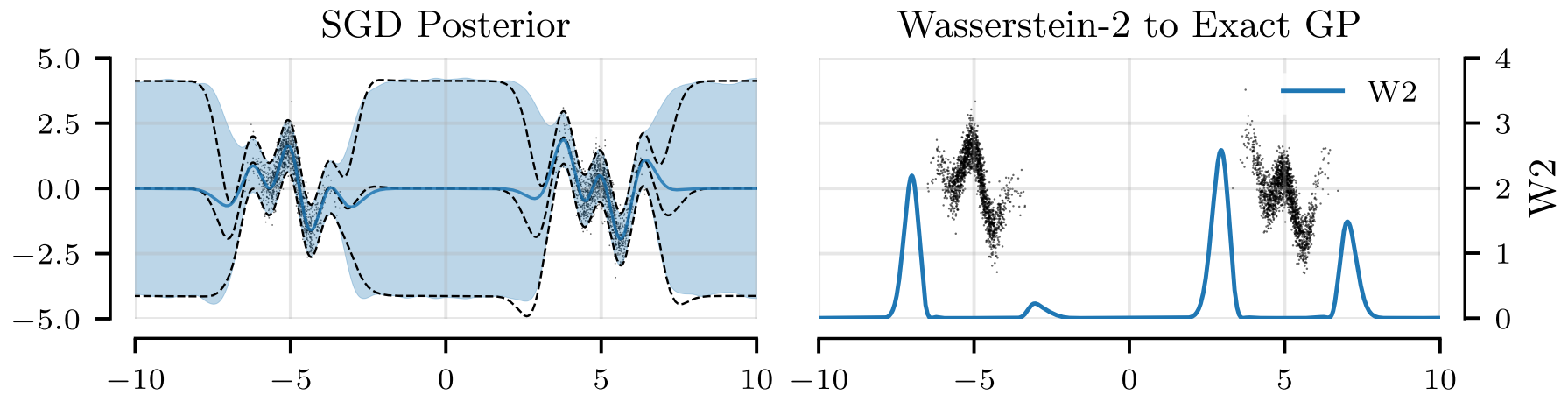
# Convergence: Euclidean and RKHS Norms



Performance not significantly affected by noise

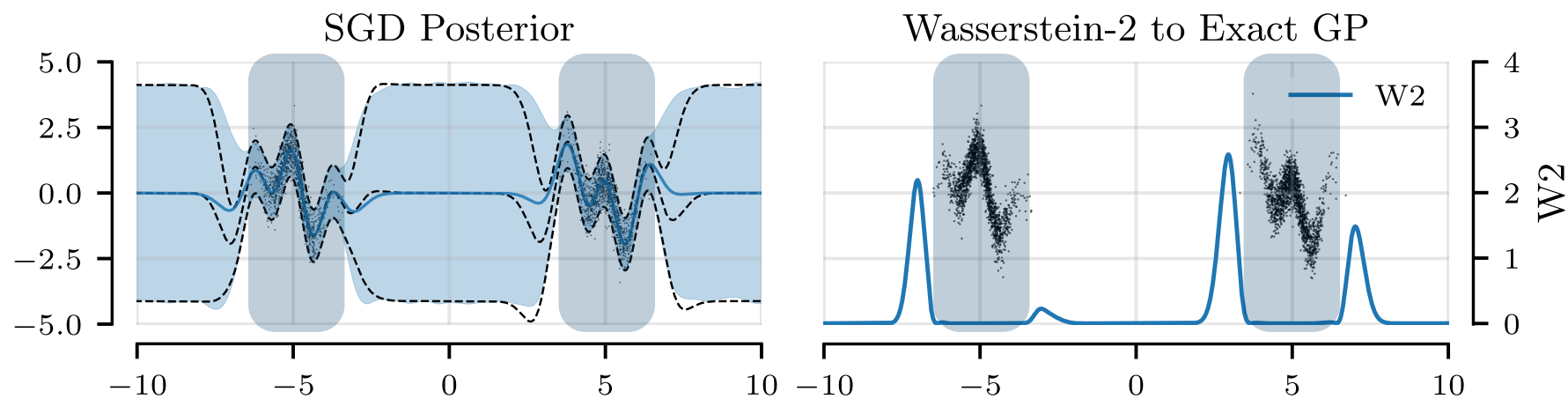
Unstable optimization problem  $\rightsquigarrow$  benign non-convergence  $\rightsquigarrow$  implicit bias

# Where does SGD's implicit bias affect predictions?



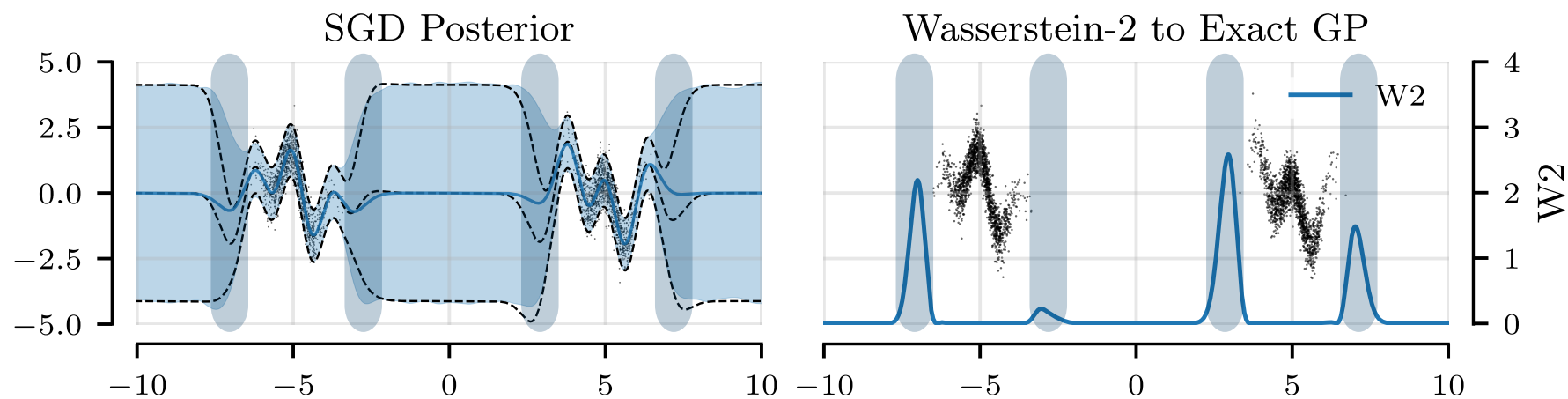
Error seems to concentrate away from data, but not too far away?

# Where does SGD's implicit bias affect predictions?



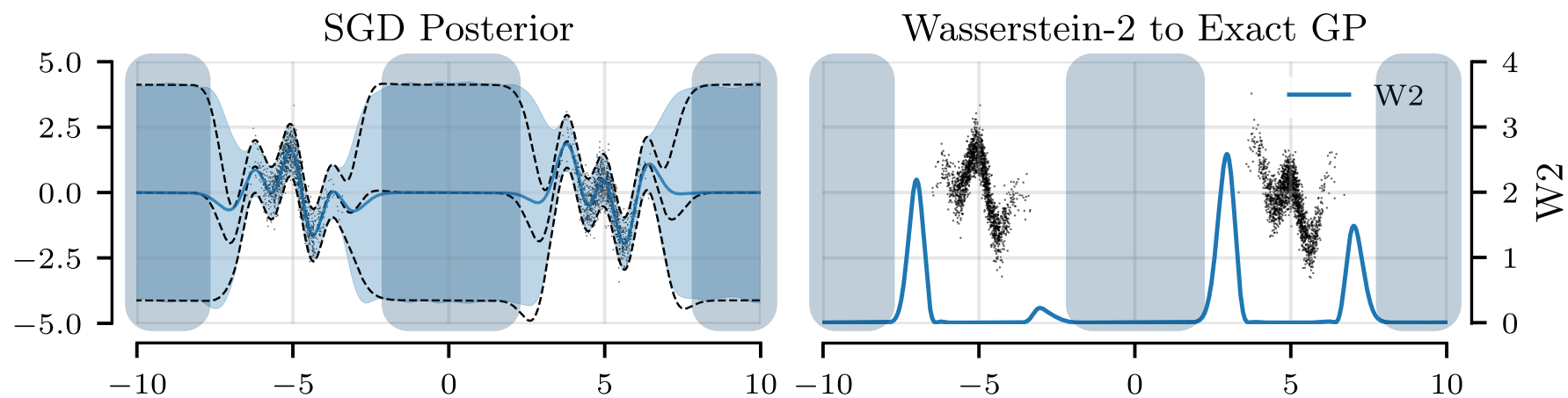
*Interpolation region*

# Where does SGD's implicit bias affect predictions?



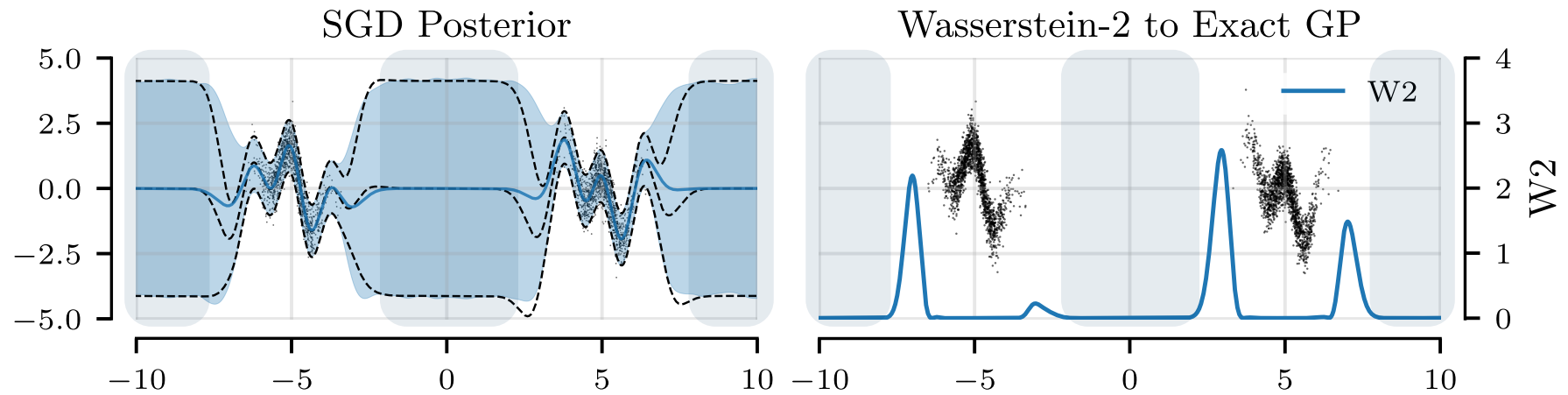
*Extrapolation region*

# Where does SGD's implicit bias affect predictions?



*Far-away region*

# The Far-away Region

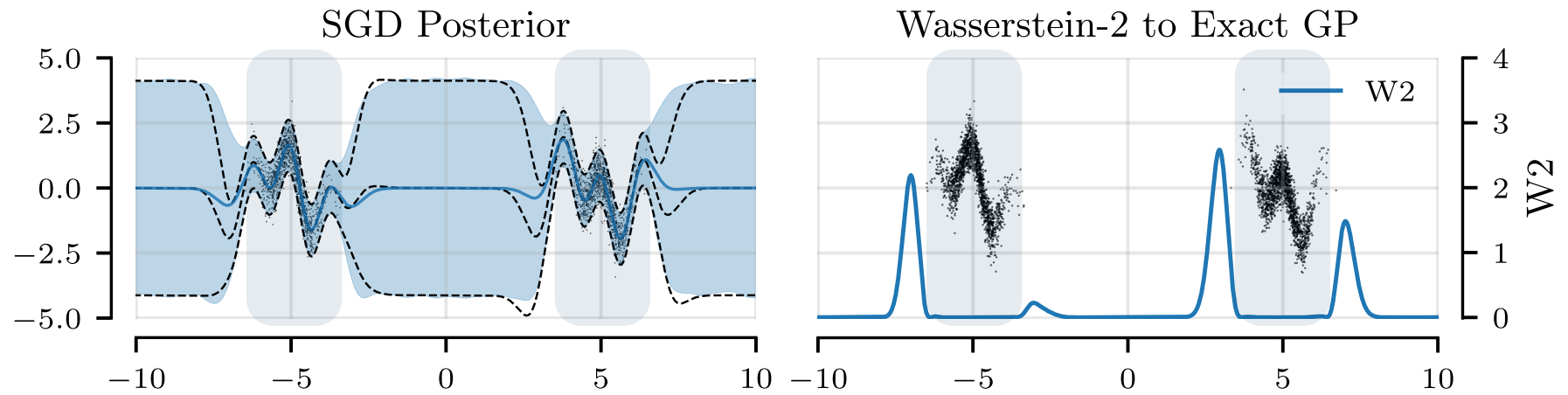


$$(f | \mathbf{y})(\cdot) = f(\cdot) + \sum_{i=1}^N v_i k(x_i, \cdot)$$

Kernel decays in space  $\rightsquigarrow$  predictions revert to prior



# The Interpolation Region



Low approximation error where data is dense

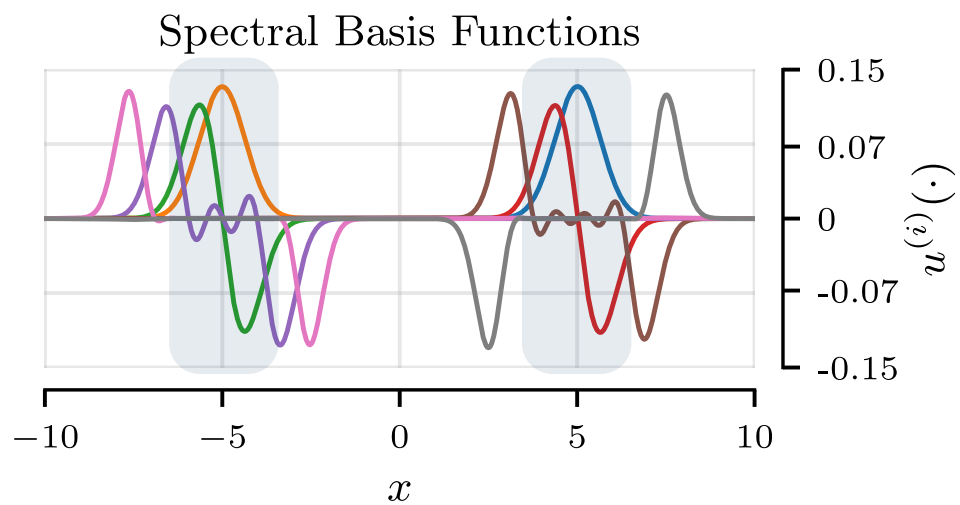
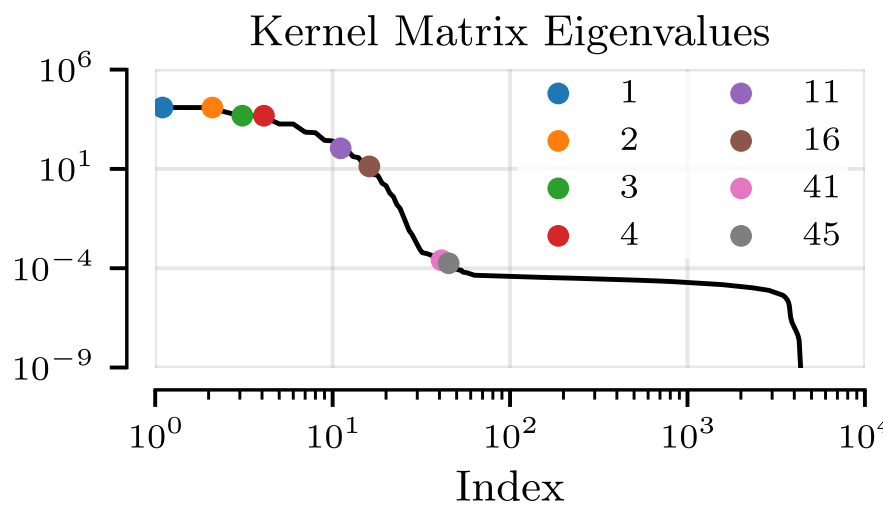
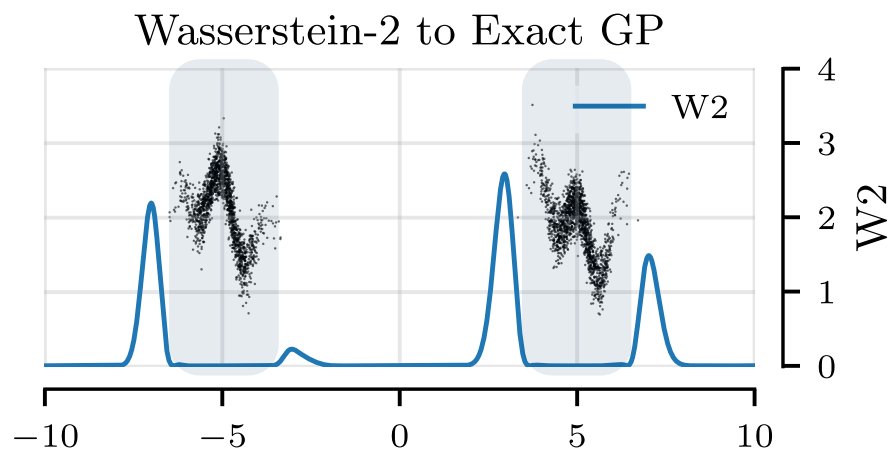
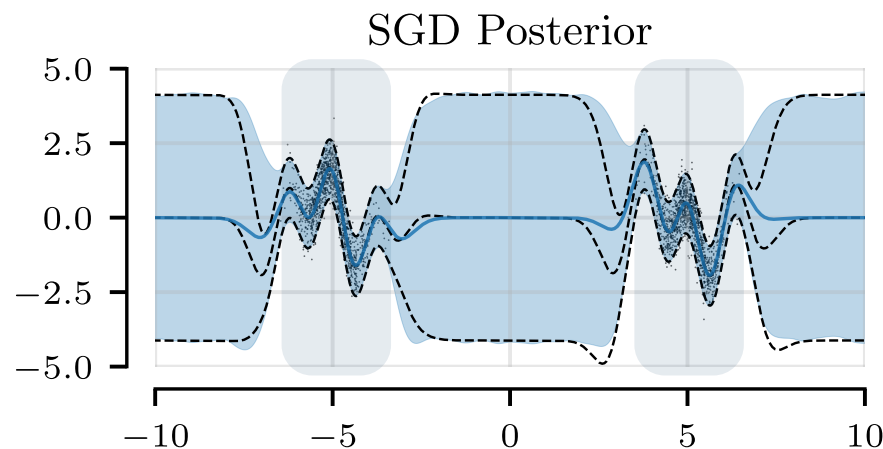
## The Interpolation Region

Idea: maybe SGD (a) converges fast on a subspace, and (b) obtains something *arbitrary but benign* on the rest of the space?

Let  $\mathbf{K}_{xx} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$  be the eigendecomposition of the kernel matrix. Define the *spectral basis functions*

$$u^{(i)}(\cdot) = \sum_{j=1}^N \frac{U_{ji}}{\sqrt{\lambda_i}} k(x_j, \cdot).$$

# The Interpolation Region



Large-eigenvalue spectral basis functions concentrate on data

# The Interpolation Region

*Proposition.* With probability  $1 - \delta$ , we have

$$\left\| \text{proj}_{u^{(i)}} \mu_{f|\mathbf{y}} - \text{proj}_{u^{(i)}} \mu_{\text{SGD}} \right\|_{H_k} \leq \frac{1}{\sqrt{\lambda_i t}} \left( \frac{\|\mathbf{y}\|_2}{\eta \sigma^2} + G \sqrt{2\eta \sigma^2 \log \frac{N}{\delta}} \right).$$

$\eta$ : learning  
rate

$\sigma^2$ : noise  
variance

$\lambda_i$ : kernel matrix  
eigenvalues

$G$ : gradient's sub-Gaussian  
coefficient

SGD converges fast with respect to top spectral basis functions

## The Interpolation Region

Where do the top spectral basis functions concentrate?

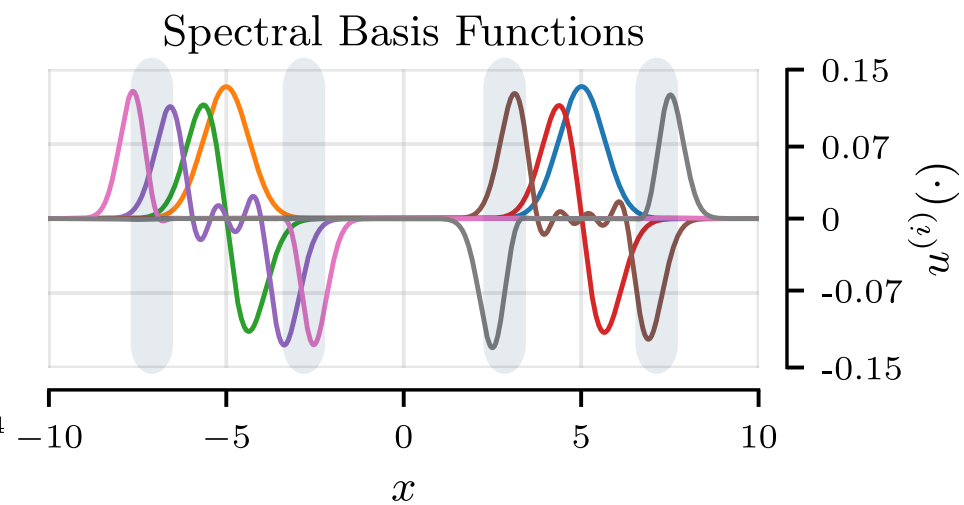
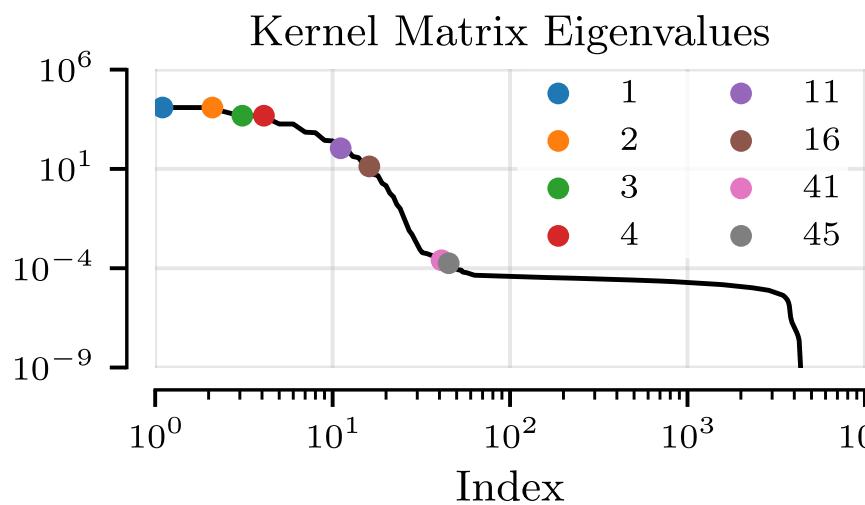
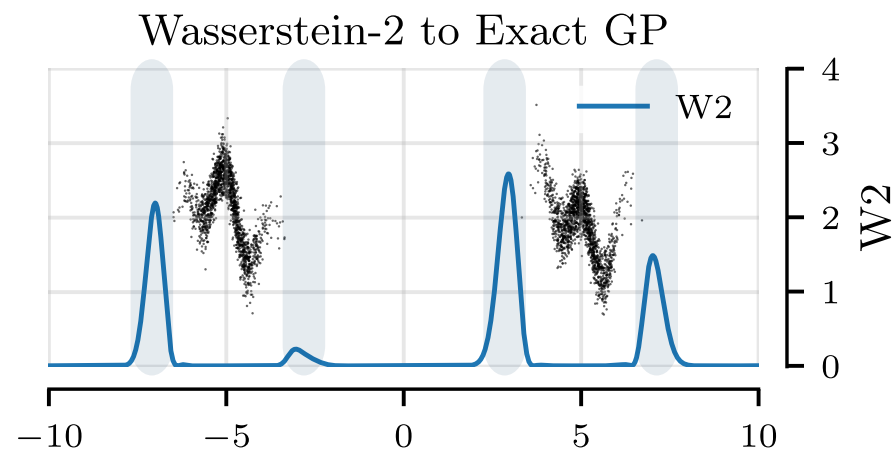
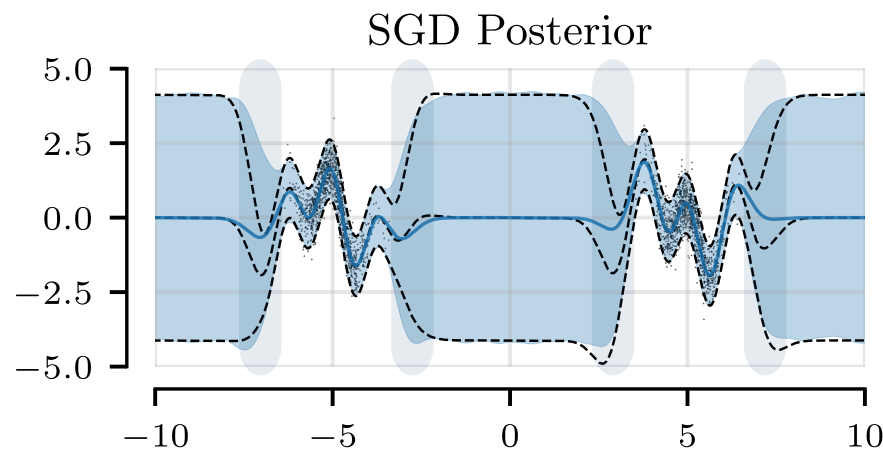
- Idea: lift Courant–Fischer eigenvector characterization to the RKHS

*Proposition.* The spectral basis functions can be written

$$u^{(i)}(\cdot) = \arg \max_{u \in H_k} \left\{ \sum_{i=1}^N u(x_i)^2 : \begin{array}{l} \|u\|_{H_k} = 1 \\ \langle u, u^{(j)} \rangle_{H_k} = 0, \forall j < i \end{array} \right\}.$$

Spectral basis functions concentrate on the data *as much as possible*, while remaining orthogonal to those with larger eigenvalues

# The Extrapolation Region



High error: where small-eigenvalue spectral basis functions concentrate

## SGD's implicit bias for Gaussian processes

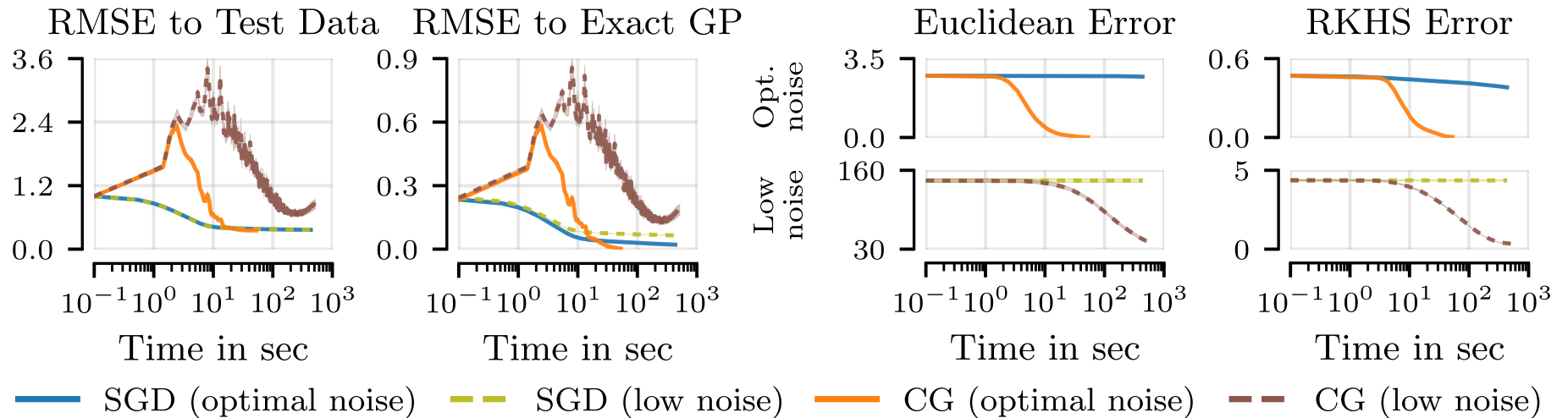
Implicit bias: SGD converges quickly near the data, and causes no harm far from the data

Error concentrates in regions (a) without much data, which also (b) aren't located too far from the data:

- Lack of data  $\rightsquigarrow$  predictions are mostly arbitrary
- Empirically: functions shrink to prior faster than exact posterior

Benign non-convergence  $\rightsquigarrow$  robustness to instability

# Performance



Conjugate gradients: non-monotonic test error, in spite of monotonic convergence  
SGD: almost always monotonic decrease in test error

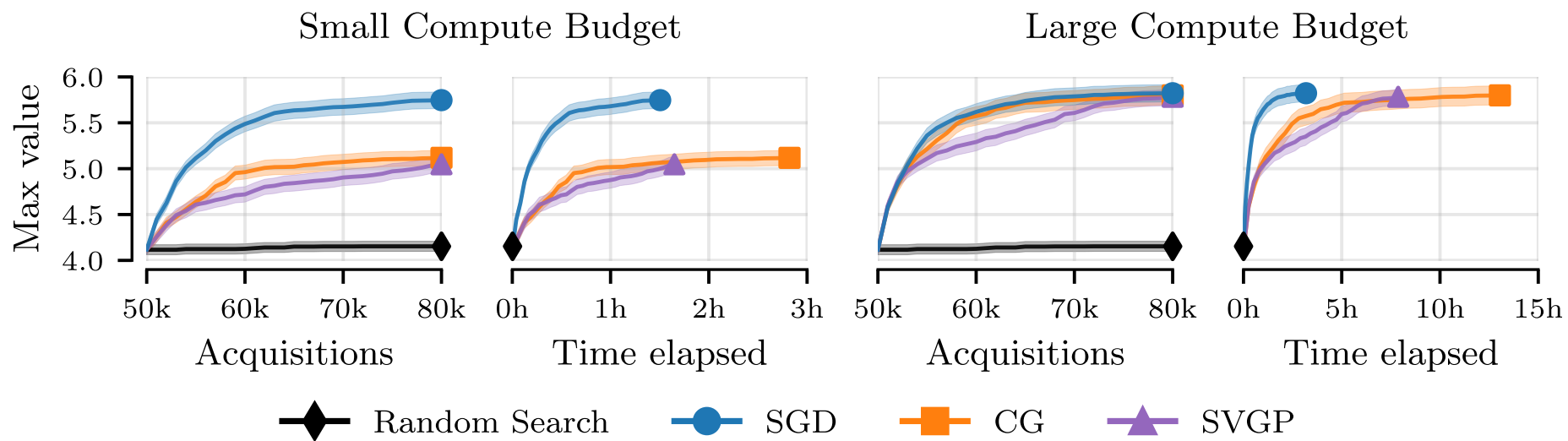


# Performance

Dataset		POL	ELEVATORS	BIKE	PROTEIN	KEGGDIR	3DROAD	SONG	BUZZ	HOUSEELEC
$N$		15000	16599	17379	45730	48827	434874	515345	583250	2049280
RMSE	SGD	$0.13 \pm 0.00$	$0.38 \pm 0.00$	$0.11 \pm 0.00$	<b><math>0.51 \pm 0.00</math></b>	$0.12 \pm 0.00$	<b><math>0.11 \pm 0.00</math></b>	<b><math>0.80 \pm 0.00</math></b>	$0.42 \pm 0.01$	<b><math>0.09 \pm 0.00</math></b>
	CG	<b><math>0.08 \pm 0.00</math></b>	<b><math>0.35 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.50 \pm 0.00</math></b>	<b><math>0.08 \pm 0.00</math></b>	$0.15 \pm 0.01$	$0.85 \pm 0.03$	$1.41 \pm 0.08$	$0.87 \pm 0.14$
	SVGP	$0.10 \pm 0.00$	$0.37 \pm 0.00$	$0.08 \pm 0.00$	$0.62 \pm 0.00$	$0.10 \pm 0.00$	$0.64 \pm 0.01$	$0.82 \pm 0.00$	<b><math>0.34 \pm 0.00</math></b>	$0.10 \pm 0.02$
RMSE <sup>†</sup>	SGD	<b><math>0.13 \pm 0.00</math></b>	<b><math>0.38 \pm 0.00</math></b>	$0.11 \pm 0.00$	<b><math>0.51 \pm 0.00</math></b>	<b><math>0.12 \pm 0.00</math></b>	<b><math>0.11 \pm 0.00</math></b>	<b><math>0.80 \pm 0.00</math></b>	<b><math>0.42 \pm 0.01</math></b>	<b><math>0.09 \pm 0.00</math></b>
	CG	$0.16 \pm 0.01$	$0.68 \pm 0.09$	<b><math>0.05 \pm 0.01</math></b>	$3.03 \pm 0.23$	$9.79 \pm 1.06$	$0.34 \pm 0.02$	$0.83 \pm 0.02$	$5.66 \pm 1.14$	$0.93 \pm 0.19$
	SVGP	—	—	—	—	—	—	—	—	—
Hours	SGD	$0.06 \pm 0.00$	$0.06 \pm 0.00$	$0.09 \pm 0.00$	$0.12 \pm 0.00$	$0.25 \pm 0.00$	$0.46 \pm 0.19$	$3.67 \pm 0.24$	$5.78 \pm 1.02$	$2.69 \pm 0.91$
	CG	$0.04 \pm 0.01$	<b><math>0.03 \pm 0.01</math></b>	$0.05 \pm 0.00$	$0.15 \pm 0.03$	$0.21 \pm 0.03$	$1.42 \pm 0.60$	$3.25 \pm 0.04$	$5.85 \pm 0.80$	$2.62 \pm 0.01$
	SVGP	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.05 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>	<b><math>0.04 \pm 0.00</math></b>
NLL	SGD	$-0.70 \pm 0.02$	$0.47 \pm 0.00$	$-0.48 \pm 0.08$	$0.64 \pm 0.01$	$-0.62 \pm 0.07$	<b><math>-0.60 \pm 0.00</math></b>	<b><math>1.21 \pm 0.00</math></b>	$0.83 \pm 0.07$	<b><math>-1.09 \pm 0.04</math></b>
	CG	<b><math>-1.17 \pm 0.01</math></b>	<b><math>0.38 \pm 0.00</math></b>	<b><math>-2.62 \pm 0.06</math></b>	<b><math>0.62 \pm 0.01</math></b>	<b><math>-0.92 \pm 0.10</math></b>	$16.27 \pm 0.45$	$1.36 \pm 0.07$	$2.38 \pm 0.08$	$2.07 \pm 0.58$
	SVGP	$-0.64 \pm 0.02$	$0.44 \pm 0.00$	$-1.47 \pm 0.02$	$0.94 \pm 0.01$	<b><math>-0.89 \pm 0.03</math></b>	$0.94 \pm 0.03$	$1.23 \pm 0.00$	<b><math>0.29 \pm 0.04</math></b>	$-0.94 \pm 0.13$

Strong predictive performance at sufficient scale

# Parallel Thompson Sampling



Uncertainty: strong decision-making performance

# Reflections

Numerical analysis conventional wisdom:

- Don't run gradient descent on quadratic objectives
  - It's slow, conjugate gradient works much better
- If CG is slow then your problem is unstable
  - Unstable problems are ill-posed, you should reformulate

This work: a very different way of looking at things

- Don't solve the linear system approximately if the solution isn't inherently needed
- Instead of a well-posed problem, a *well-posed subproblem* might be good enough
- Try SGD! It might work well, including for counterintuitive reasons
- A kernel matrix's eigenvectors carry information about data-density
  - To see this, adopt a function-analytic view given by the spectral basis functions

# Thank you!

[HTTPS://AVT.IM/](https://AVT.IM/) ·  @AVT\_IM

J. A. Lin,\* J. Antorán,\* S. Padhy,\* D. Janz, J. M. Hernández-Lobato, A. Terenin. Sampling from Gaussian Process Posteriors using Stochastic Gradient Descent. *arXiv:2306.11589*, 2023.

A. Terenin,\* D. R. Burt,\* A. Artemev, S. Flaxman, M. van der Wilk, C. E. Rasmussen, H. Ge. Numerically Stable Sparse Gaussian Processes via Minimum Separation using Cover Trees. *arXiv: 2210.07893*, 2022.

J. T. Wilson,\* V. Borovitskiy,\* P. Mostowsky,\* A. Terenin,\* M. P. Deisenroth. Efficiently Sampling Functions from Gaussian Process Posteriors. *International Conference on Machine Learning*, 2020. **Honorable Mention for Outstanding Paper Award.**

J. T. Wilson,\* V. Borovitskiy,\* P. Mostowsky,\* A. Terenin,\* M. P. Deisenroth. Pathwise Conditioning of Gaussian Process. *Journal of Machine Learning Research*, 2021.

\*Equal contribution



UNIVERSITY OF  
CAMBRIDGE