

Gaussian Processes: The Story So Far

Gaussian processes (GPs) are renowned for their exceptional data efficiency, reliable uncertainty estimation, flexibility, and built-in mechanisms to mitigate against overfitting. However, they are often unfavorably compared to deep learning approaches due to limited scalability and their inability to capture hierarchies of abstract representations.

Sparse variational GPs (svGPs) [5] address scalability by introducing auxiliary inducing variables $\mathbf{u} \triangleq f(\mathbf{Z}) \in \mathbb{R}^M$ at pseudo-inputs $\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_M]^\top$. Approximate the posterior $p(\mathbf{f}, \mathbf{u} | \mathbf{y})$ with $q^*(\mathbf{f}, \mathbf{u}) = \arg \min_q \text{KL}[q(\mathbf{f}, \mathbf{u}) \| p(\mathbf{f}, \mathbf{u} | \mathbf{y})]$ where $q(\mathbf{f}, \mathbf{u}) \triangleq p(\mathbf{f} | \mathbf{u})q(\mathbf{u})$ and $q(\mathbf{u}) \triangleq \mathcal{N}(\mathbf{m}_{\mathbf{u}}, \mathbf{C}_{\mathbf{u}})$.

Leads to predictive density:

$$q(f(\mathbf{x})) = \mathcal{GP}(\mathbf{k}_{\mathbf{u}}^\top(\mathbf{x})\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{m}_{\mathbf{u}}, k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{\mathbf{u}}^\top(\mathbf{x})\mathbf{K}_{\mathbf{uu}}^{-1}(\mathbf{K}_{\mathbf{uu}} - \mathbf{C}_{\mathbf{u}})\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{k}_{\mathbf{u}}(\mathbf{x}')) \quad (1)$$

where $[\mathbf{K}_{\mathbf{uu}}]_{mm'} \triangleq \text{Cov}(u_m, u_{m'})$.

- Reduces cost from $\mathcal{O}(N^3)$ to $\mathcal{O}(M^3)$ (assuming $M \ll N$)
- Unlocks greater flexibility in model specification

Basis functions are effectively $\mathbf{k}_{\mathbf{u}}: \mathcal{X} \rightarrow \mathbb{R}^M$ where each element $[\mathbf{k}_{\mathbf{u}}(\mathbf{x})]_m \triangleq \text{Cov}(f(\mathbf{x}), u_m)$

- Standard Inducing Points are values of f evaluated at pseudo-inputs

$$u_m \triangleq f(\mathbf{z}_m) \Rightarrow [\mathbf{K}_{\mathbf{uu}}]_{mm'} = k(\mathbf{z}_m, \mathbf{z}_{m'}) \text{ and } [\mathbf{k}_{\mathbf{u}}(\mathbf{x})]_m = k(\mathbf{z}_m, \mathbf{x})$$

Mapping is static: solely determined by fixed kernel k and local influence of \mathbf{z}_m .

- Inter-domain Inducing Features [3] are a generalisation involving scalar projection of f onto some ϕ_m in the reproducing kernel Hilbert space (RKHS) \mathcal{H} of k ,

$$u_m \triangleq \langle f, \phi_m \rangle_{\mathcal{H}} \Rightarrow [\mathbf{K}_{\mathbf{uu}}]_{mm'} = \langle \phi_m, \phi_{m'} \rangle_{\mathcal{H}} \text{ and } [\mathbf{k}_{\mathbf{u}}(\mathbf{x})]_m = \phi_m(\mathbf{x})$$

Mapping is adaptive: can result in sparser representations that lead to greater scalability.

Spherical Inducing Features

Spherical harmonics, an extension of the Fourier basis to multiple dimensions, can be used to form ϕ_m [1]. Orthogonality leads to diagonal covariance:

$$[\mathbf{K}_{\mathbf{uu}}]_{mm'} = \lambda_m^{-1} \delta_{mm'}$$

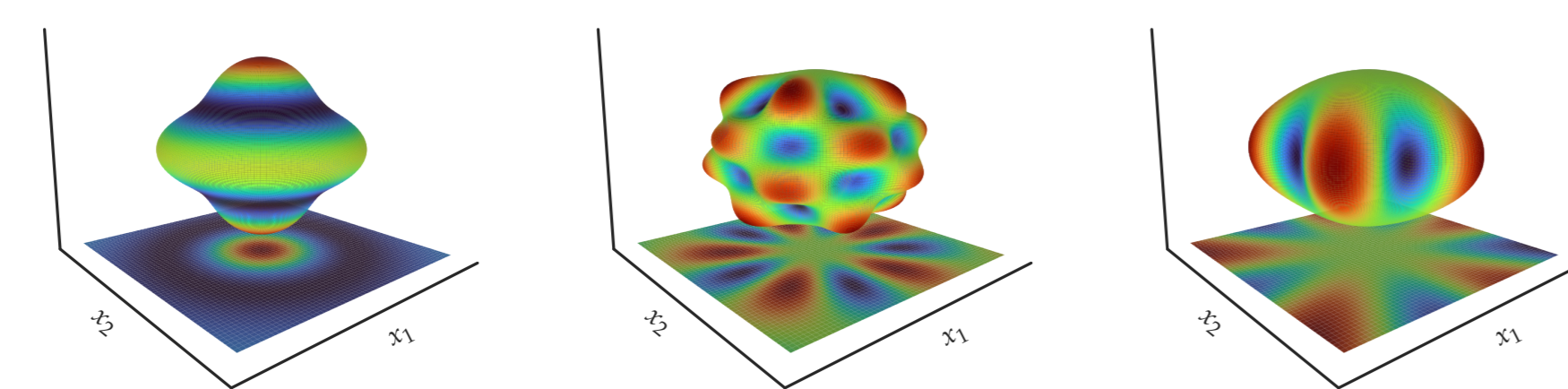


Figure 1: A few example spherical harmonics in 3D.

- Reduces cost from $\mathcal{O}(M^3)$ to $\mathcal{O}(M)$ (!)

Spherical neural network (NN) activations. To make $\mathbf{k}_{\mathbf{u}}(\mathbf{x})$ resemble a hidden layer in a feedforward NN [2], define ϕ_m as the m th hidden unit with nonlinear activation σ ,

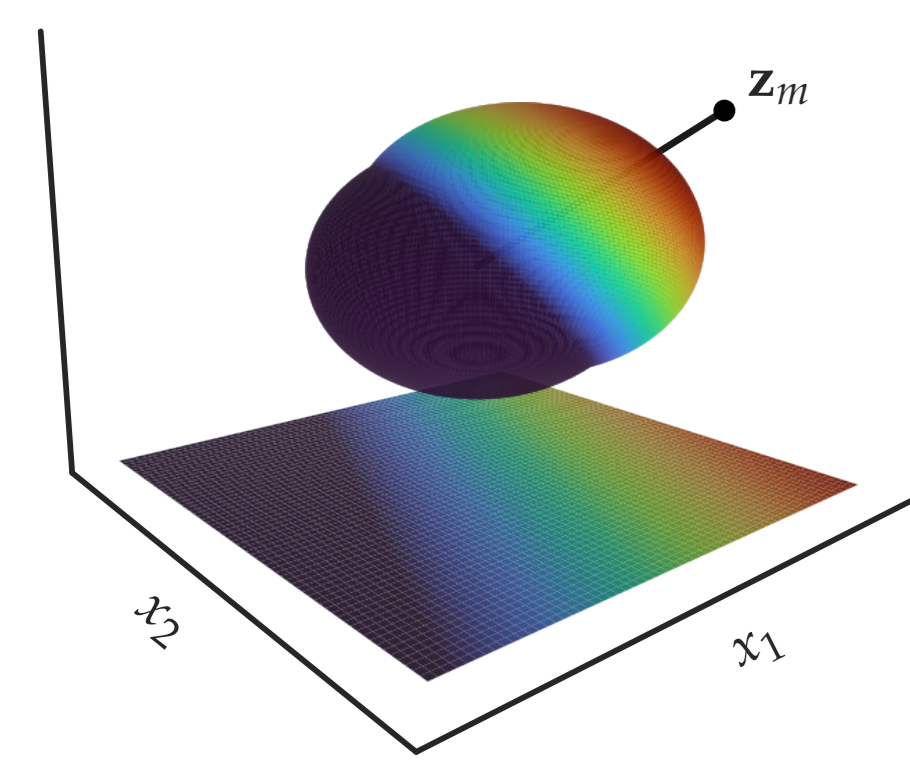
$$\phi_m(\mathbf{x}) \triangleq \|\mathbf{z}_m\| \|\mathbf{x}\| \cdot \sigma\left(\frac{\mathbf{z}_m^\top \mathbf{x}}{\|\mathbf{z}_m\| \|\mathbf{x}\|}\right)$$

- Predictive mean becomes a single-layer feedforward NN:

$$\mathbf{k}_{\mathbf{u}}(\mathbf{x})^\top \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{m}_{\mathbf{u}} = \sum_{m=1}^M \beta_m \phi_m(\mathbf{x}), \quad \beta \triangleq \mathbf{K}_{\mathbf{uu}}^{-1} \mathbf{m}_{\mathbf{u}} \in \mathbb{R}^M$$

- When stacked to form a deep GP (DGP), the propagation of predictive means emulates forward pass through a deep NN (DNN)
- Obtain predictive uncertainty in DNNs for free as a byproduct (!)

Figure 2: A RELU-activated hidden unit on the sphere in 3D.



Orthogonally-Decoupled Gaussian Processes

- Decouple GP as sum of two independent GPs [4]:

$$f(\mathbf{x}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}')) \Leftrightarrow f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x}),$$

where

$$g(\mathbf{x}) \sim \mathcal{GP}(0, \mathbf{k}_{\mathbf{u}}^\top(\mathbf{x})\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{k}_{\mathbf{u}}(\mathbf{x}')),$$

$$h(\mathbf{x}) \sim \mathcal{GP}(0, s(\mathbf{x}, \mathbf{x}'))$$

for $s(\mathbf{x}, \mathbf{x}') \triangleq k(\mathbf{x}, \mathbf{x}') - \mathbf{k}_{\mathbf{u}}^\top(\mathbf{x})\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{k}_{\mathbf{u}}(\mathbf{x}')$.

- Introduce orthogonal inducing variables $\mathbf{v} \triangleq f(\mathbf{W}), \mathbf{v}' \triangleq h(\mathbf{W}) \in \mathbb{R}^K$ at pseudo-inputs $\mathbf{W} \triangleq [\mathbf{w}_1 \dots \mathbf{w}_K]^\top$.

- Approximate posterior $q(\mathbf{v}') \triangleq \mathcal{N}(\mathbf{m}_{\mathbf{v}'}, \mathbf{C}_{\mathbf{v}'})$.

Leads to predictive density:

$$q(\mathbf{f}_*) = \mathcal{N}(\mathbf{K}_{*u}\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{m}_{\mathbf{u}} + \underbrace{\mathbf{S}_{*v}\mathbf{S}_{\mathbf{v}\mathbf{v}}^{-1}\mathbf{m}_{\mathbf{v}'}}_{\text{orthogonal bases}}, \mathbf{K}_{**} - \mathbf{K}_{*u}\mathbf{K}_{\mathbf{uu}}^{-1}(\mathbf{K}_{\mathbf{uu}} - \mathbf{C}_{\mathbf{u}})\mathbf{K}_{\mathbf{uu}}^{-1}\mathbf{K}_{u*} - \underbrace{\mathbf{S}_{*v}\mathbf{S}_{\mathbf{v}\mathbf{v}}^{-1}(\mathbf{S}_{\mathbf{v}\mathbf{v}} - \mathbf{C}_{\mathbf{v}'})\mathbf{S}_{\mathbf{v}\mathbf{v}}^{-1}\mathbf{S}_{v*}}_{\text{orthogonal bases}}) \quad (2)$$

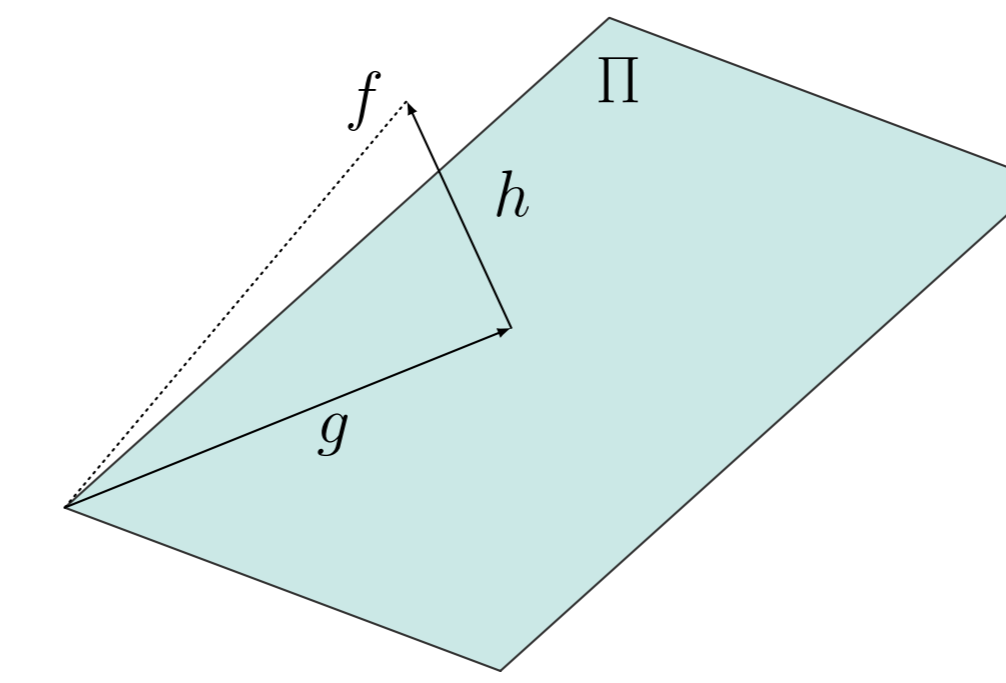


Figure 3: Hyperplane $\Pi \triangleq \{\alpha^\top \mathbf{k}_{\mathbf{u}}(\cdot); \alpha \in \mathbb{R}^M\}$

Technical Issues with Spherical NN Activation Features

In practice, several widely-used kernels and activation functions are incompatible:

- Spectra mismatch.** For the Matérn kernel, discrepancies in its Fourier coefficients (nonzero) with those of the activation features (zero) lead to overestimation of the predictive variance.
- RKHS inner product.** For the RELU activation features, its (squared) Fourier coefficients decay at the same rate as those of numerous kernels, resulting in an indeterminate RKHS inner product.

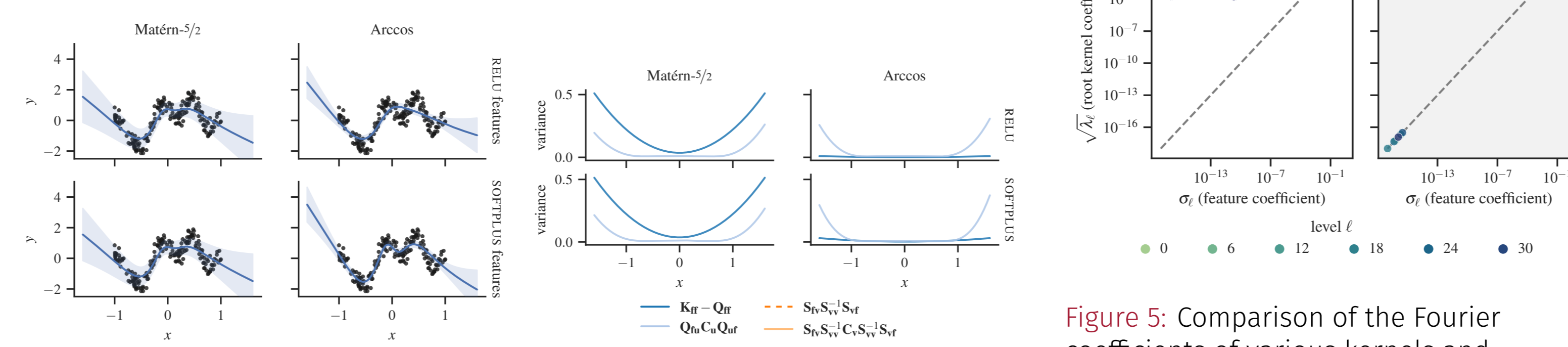


Figure 4: Posterior of svGPs with various kernels and activations; at $L = 8$ levels.

Our Solution

Extend the orthogonally-decoupled GP framework with inter-domain inducing features: let

$$u_m \triangleq \langle f, \phi_m \rangle_{\mathcal{H}}, \text{ and } v_k \triangleq \langle f, \psi_k \rangle_{\mathcal{H}}$$

for some choices of $\phi_m, \psi_k \in \mathcal{H}$.

To obtain $\mathbf{S}_{\mathbf{v}\mathbf{v}}, \mathbf{S}_{\mathbf{v}\mathbf{f}}$ in eq. 2, need to compute prior covariances $\mathbf{K}_{\mathbf{v}\mathbf{f}}, \mathbf{K}_{\mathbf{v}\mathbf{u}}, \mathbf{K}_{\mathbf{v}\mathbf{v}}$

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \\ \mathbf{v} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{f}\mathbf{f}} & \mathbf{K}_{\mathbf{f}\mathbf{u}}^\top & \mathbf{K}_{\mathbf{f}\mathbf{v}}^\top \\ \mathbf{K}_{\mathbf{u}\mathbf{f}} & \mathbf{K}_{\mathbf{u}\mathbf{u}} & \mathbf{K}_{\mathbf{u}\mathbf{v}}^\top \\ \mathbf{K}_{\mathbf{v}\mathbf{f}} & \mathbf{K}_{\mathbf{v}\mathbf{u}} & \mathbf{K}_{\mathbf{v}\mathbf{v}} \end{bmatrix}\right)$$

In this work:

- ϕ_m : m th unit of the spherical activation layer
- $\psi_k(\mathbf{x}) \triangleq k(\mathbf{w}_k, \mathbf{x})$

Leads to covariances:

$$[\mathbf{K}_{\mathbf{v}\mathbf{f}}]_{kn} \triangleq \text{Cov}(v_k, f(\mathbf{x}_n)) = k(\mathbf{w}_k, \mathbf{x}_n),$$

$$[\mathbf{K}_{\mathbf{v}\mathbf{u}}]_{km} \triangleq \text{Cov}(v_k, u_m) = \phi_m(\mathbf{w}_k),$$

$$[\mathbf{K}_{\mathbf{v}\mathbf{v}}]_{kk'} \triangleq \text{Cov}(v_k, v_{k'}) = k(\mathbf{w}_k, \mathbf{w}_{k'}).$$

Cross-covariance $\mathbf{K}_{\mathbf{v}\mathbf{u}}$ consists of forward-pass of pseudo-inputs \mathbf{w}_k through neurons ϕ_m

Experimental Results

Regression on Synthetic 1D Dataset

Incorporating a small handful of $K = 8$ orthogonal inducing variables costs roughly the same as doubling the truncation level L but leads to substantial improvements.

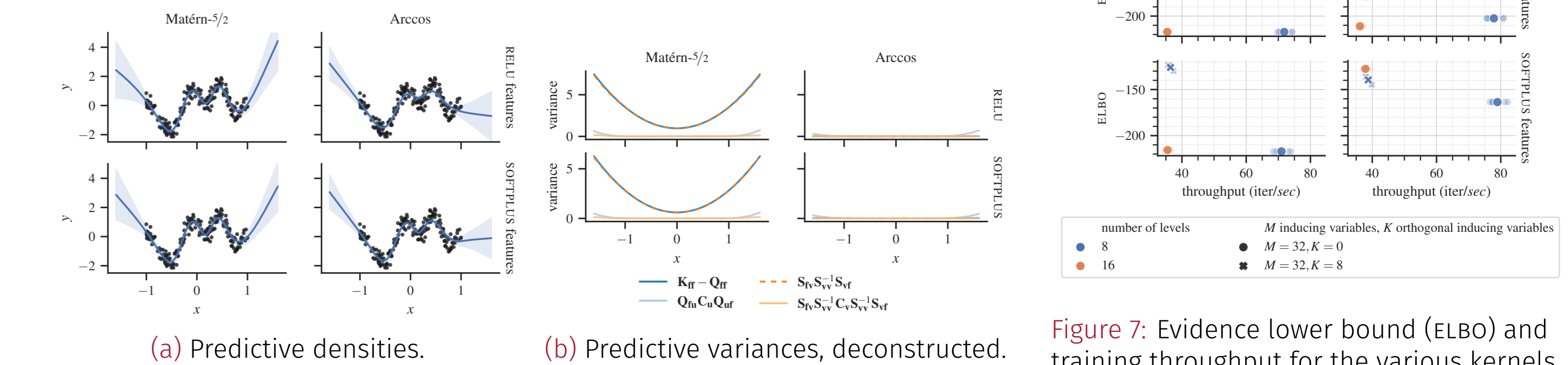


Figure 6: Posterior of svGPs with various kernels and activations, $K = 8$ orthogonal bases; at $L = 8$ levels. New term $\mathbf{S}_{*v}\mathbf{S}_{\mathbf{v}\mathbf{v}}^{-1}\mathbf{S}_{v*}$ offsets errors from the original basis.

Regression on UCI Repository Datasets

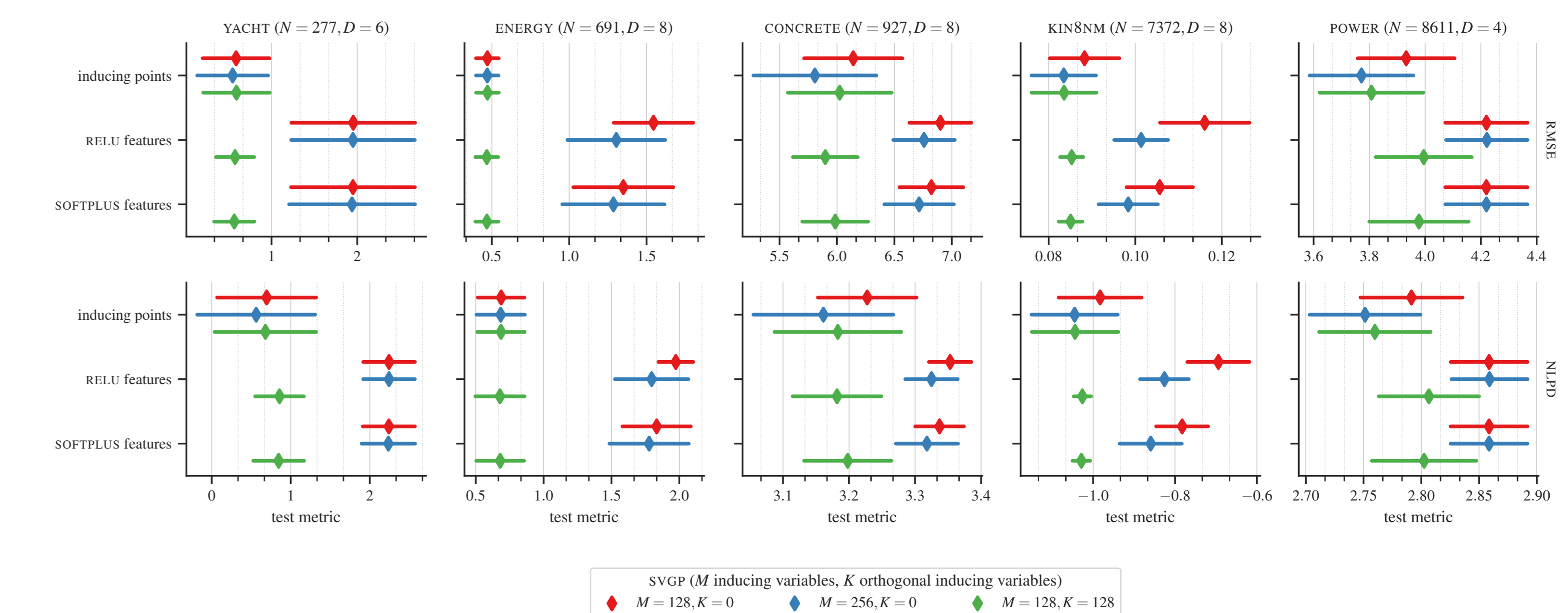


Figure 8: Test metrics, RMSE and NLPD, on the UCI regression datasets using the Arccos kernel with various activation features. Along the rows labeled "inducing points", the red and blue markers (\diamond, \blacklozenge) represent the original svGP model [5], while the green markers (\blacklozenge) represent solveGP [4]. Along the remaining rows, the red and blue markers ($\blacklozenge, \blacklozenge$) represent the activated svGP [2], while the green markers (\blacklozenge) represent our proposed approach.

References

- Vincent Dutordoir, Nicolas Durrande, and James Hensman. Sparse Gaussian processes with spherical harmonic features. In *International Conference on Machine Learning*, pages 2793–2802. PMLR, 2020.
- Vincent Dutordoir, James Hensman, Mark van der Wilk, Carl Henrik Ek, Zoubin Ghahramani, and Nicolas Durrande. Deep neural networks as point estimates for deep Gaussian processes. *Advances in Neural Information Processing Systems*, 34, 2021.
- Miguel Lázaro-Gredilla and Anibal Figueiras-Vidal. Inter-domain Gaussian processes for sparse inference using inducing features. *Advances in Neural Information Processing Systems*, 22, 2009.
- Jiaxin Shi, Michalis Titsias, and Andriy Mnih. Sparse orthogonal variational inference for Gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, pages 1932–1942. PMLR, 2020.
- Michalis Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *Artificial Intelligence and Statistics*, pages 567–574. PMLR, 2009.