

# Thin and Deep Gaussian processes

Scan me to see the poster, blogpost and more!



Daniel Augusto de Souza, Alexander Nikitin, S. T. John, Magnus Ross, Mauricio A Álvarez, Marc Peter Deisenroth, João P. P. Gomes, Diego Mesquita, César L. C. Mattos

## Why (deep) Gaussian processes?

- Gaussian processes (GPs) are widely used in machine learning for their simple uncertainty quantification;
- GP's kernel function → Modelling and uncertainty quality;
- Common stationary kernels, like square exponential and Matérn are unsuitable for non-stationary data;
  - Many geospatial processes, such as sea surface height, bathymetry.
- Popular research direction: how to make non-stationary kernels from common stationary ones.

## From stationary to non-stationary

- Stationarity:  $k(\mathbf{a}, \mathbf{b}) = k(\mathbf{a} - \mathbf{b}, 0)$ ;
  - Kernel is effectively a one-argument function, every slice looks the same.
- Isotropic stationary:  $k(\mathbf{a}, \mathbf{b}) = \pi_k((\mathbf{a} - \mathbf{b})^T \Delta^{-1} (\mathbf{a} - \mathbf{b}))$ ;
  - Kernel is a function of the distance, weighted by the lengthscale matrix  $\Delta$ .
- E.g., squared exponential kernel  $\pi_k(d^2) = \exp[-\frac{1}{2}d^2]$

## Compositional kernels

- Non-linear warping function:  $\tau(\cdot) \rightarrow k(\tau(\mathbf{a}), \tau(\mathbf{b}))$ ;
  - If  $\tau(x) = \ell^{-1} \cdot x$ , then  $\ell$  gets absorbed in the kernel's lengthscales
- If  $\tau(\mathbf{x})$  is a neural network → Deep Kernel Learning;
- If  $\tau(\mathbf{x})$  is a GP → Compositional Deep Gaussian process;

## Limitations

- The more complex  $\tau(\cdot)$  is, the harder the model is to interpret;
- For DKL, parameters are not Bayesian, thus the model overfits.
- For DGP,  $\tau(\cdot)$  cannot have zero prior mean → model collapse with increasing depth;

## Lengthscale mixture kernels

- Positive semi-definite function field  $\Delta(\cdot)$ 

$$k_{\Delta}(\mathbf{a}, \mathbf{b}) = \frac{\sqrt{|\Delta(\mathbf{a})|} \sqrt{|\Delta(\mathbf{b})|}}{\sqrt{|\Delta(\mathbf{a}) + \Delta(\mathbf{b})|}} \pi_k\left((\mathbf{a} - \mathbf{b})^T \left[\frac{\Delta(\mathbf{a}) + \Delta(\mathbf{b})}{2}\right]^{-1} (\mathbf{a} - \mathbf{b})\right)$$
- If  $\sigma(\Delta(\cdot))$  is a GP w/ linking function  $\sigma(\cdot)$  → Deeply Non-stationary GP

## Limitations

- Does not encode latent spaces; • Kernel scale affected by the pre-factor, giving rise to unwanted correlations;

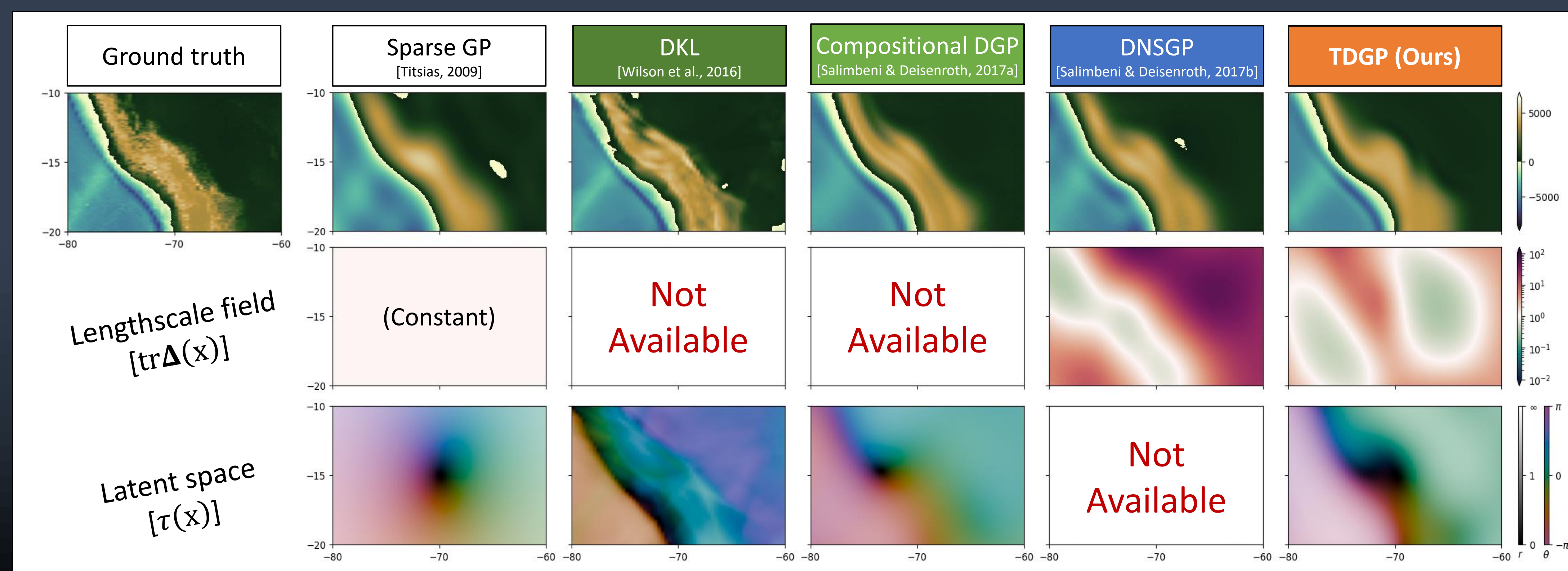
## Our hybrid proposal

- We propose a compositional kernel:  $\tau(\mathbf{x}) = \mathbf{W}(\mathbf{x}) \cdot \mathbf{x}$ .
  - For deeper layers,  $\tau^{(2)}(\tau^{(1)}(\mathbf{x})) = \mathbf{W}^{(2)}(\tau(\mathbf{x})) \cdot \mathbf{W}^{(1)}(\mathbf{x}) \cdot \mathbf{x}$ ;
- At the neighborhood of  $\mathbf{x}$ , we have  $\Delta(\mathbf{x}) = [\mathbf{W}^T(\mathbf{x})\mathbf{W}(\mathbf{x})]^{-1}$ ;
- Therefore, we learn latent spaces and lengthscale fields;
- If  $\mathbf{W}(\cdot)$  is a GP, then we obtain our deep model.  $\mathbf{W}(\cdot)$  can be zero mean and with learned variances, we optimize for dim. reduction

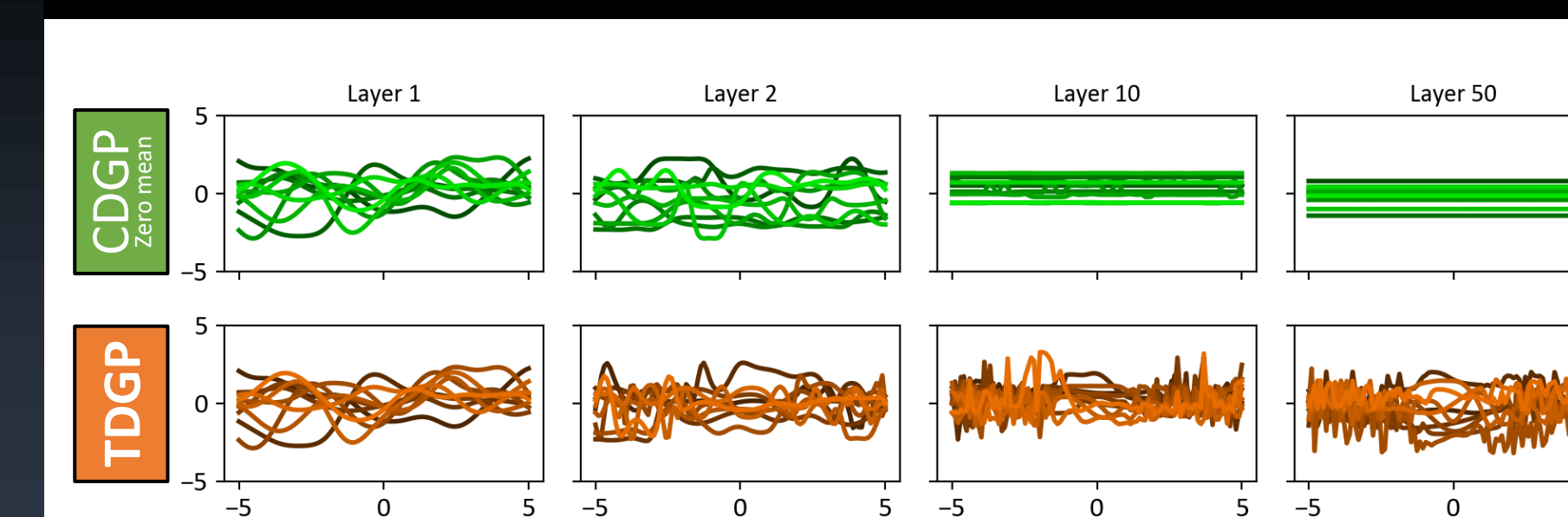
## Limitations

- Adds  $d \times q$  more GPs, making training harder; • Without adding a bias dimension to input, neighborhood of 0 is unaffected;

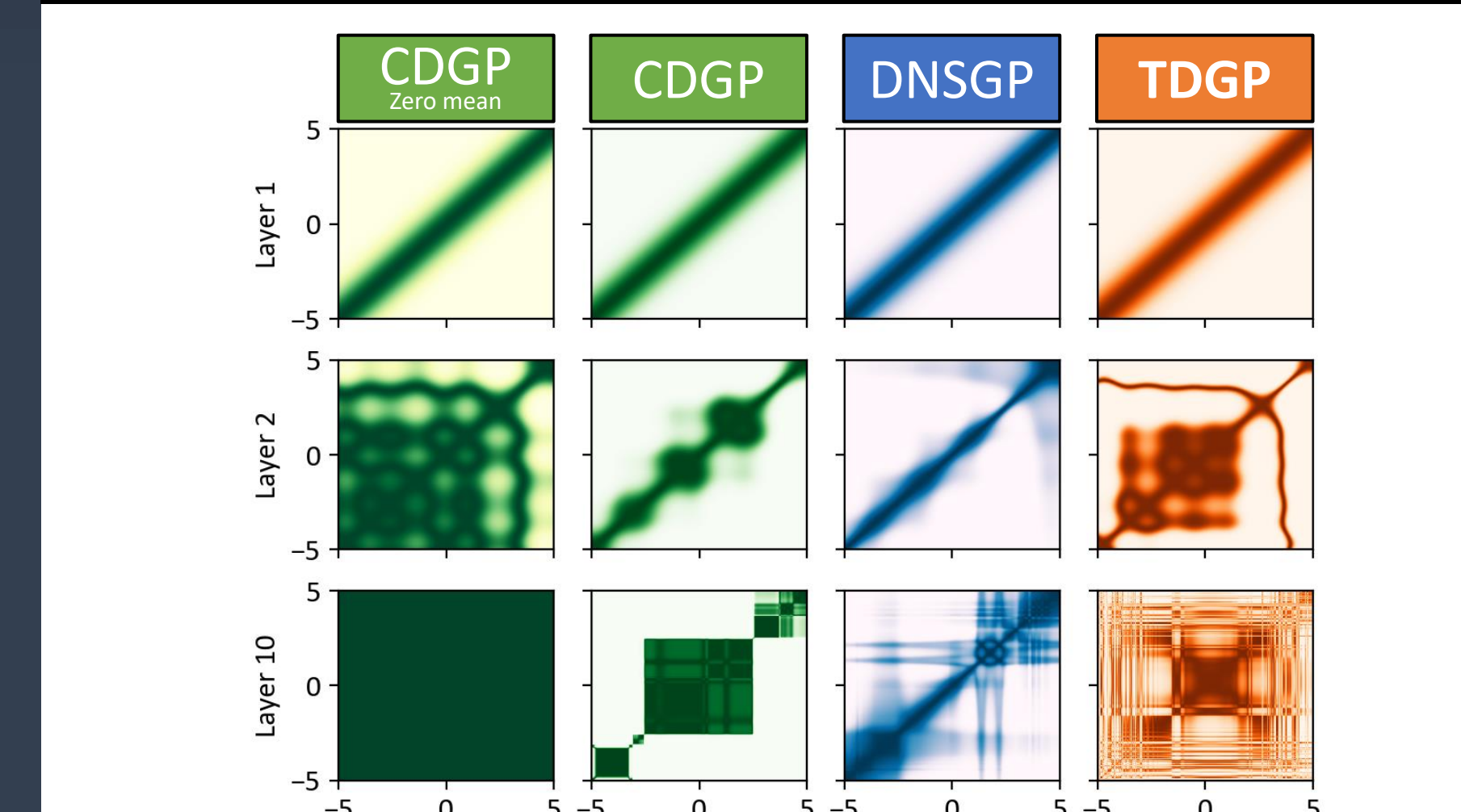
By unifying different approaches to deep Gaussian processes, we build probabilistic models that are more interpretable whilst learning lower-dimensional latent representations for complex, non-stationary data



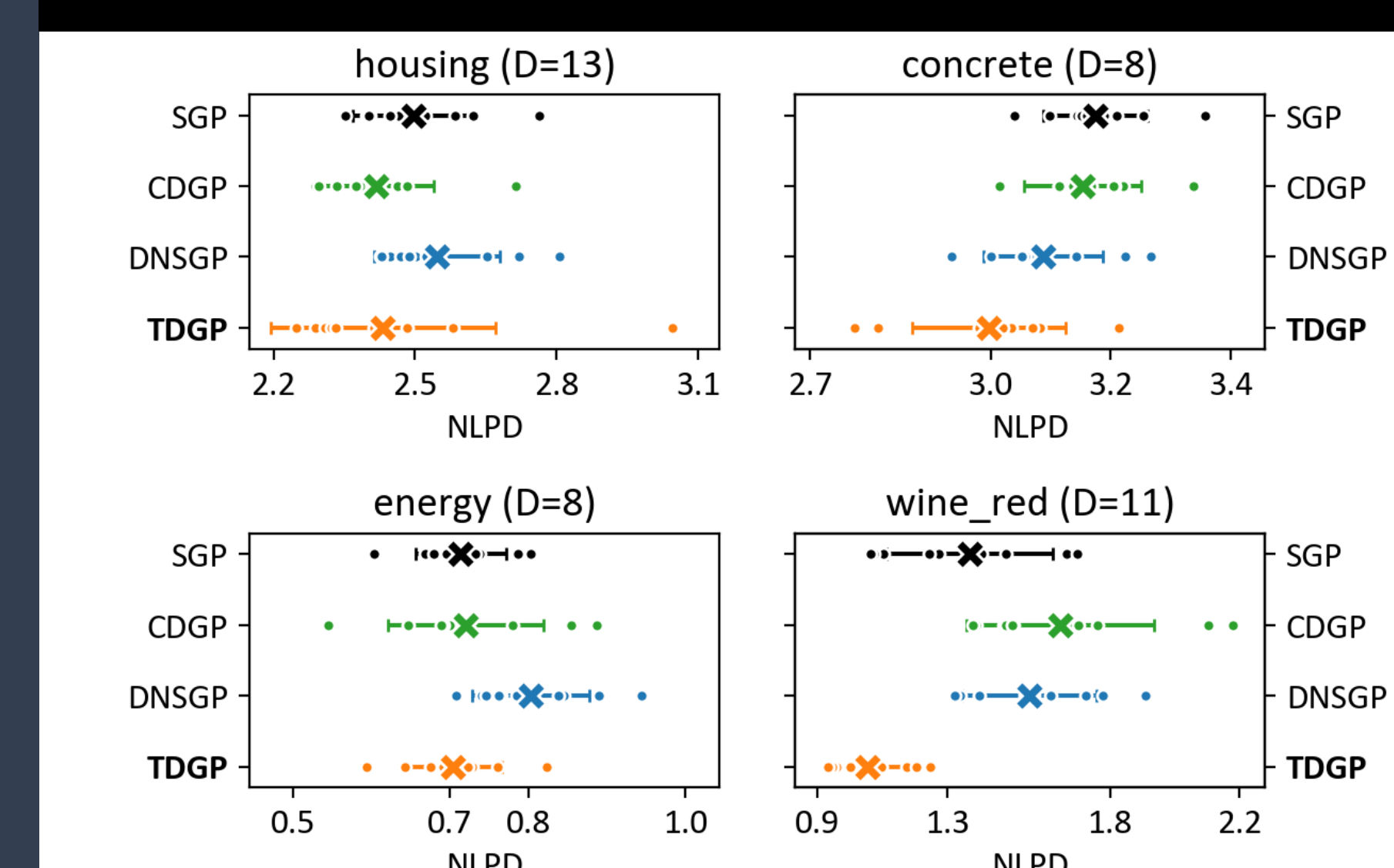
## Very deep GP priors



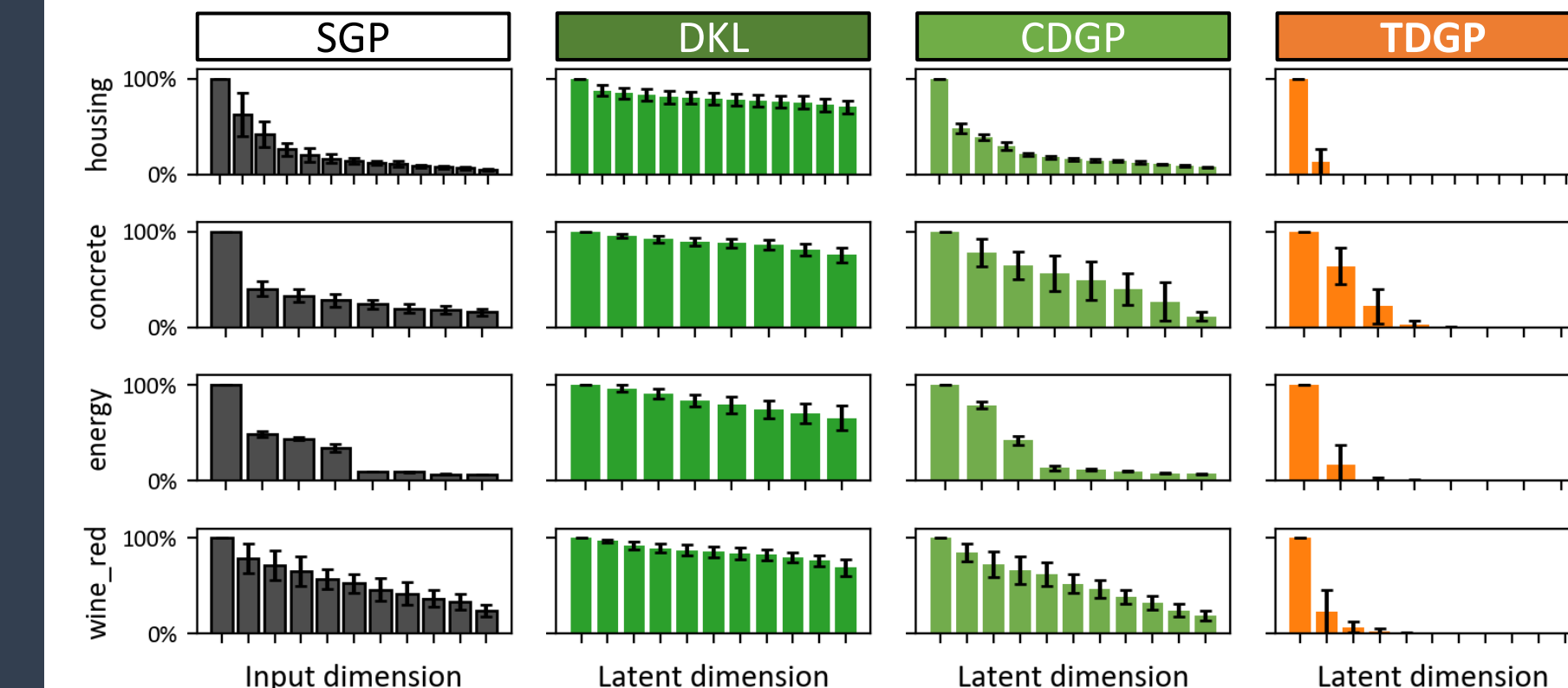
## Prior kernel matrices



## UCI datasets



## Dimensionality reduction



## Bathymetry on Central Andes

	NLPD	MRAE
Sparse GP	-0.13 ± 0.09	1.19 ± 0.63
Deep Kernel Learning	3.85 ± 0.92	<b>0.59 ± 0.31</b>
Compositional DGP	-0.44 ± 0.12	0.83 ± 0.56
Deeply Nonstationary GP	-0.31 ± 0.12	1.12 ± 0.75
<b>TDGP (Ours)</b>	<b>-0.53 ± 0.10</b>	0.66 ± 0.43

## Learned correlation

