

Deep Gaussian Processes

- ▶ Deep probabilistic models;
- ▶ Composition of functions:

$$f(\mathbf{x}) = \left(h^{(N_h-1)} \left(\theta^{(N_h-1)} \right) \circ \dots \circ h^{(0)} \left(\theta^{(0)} \right) \right) (\mathbf{x});$$

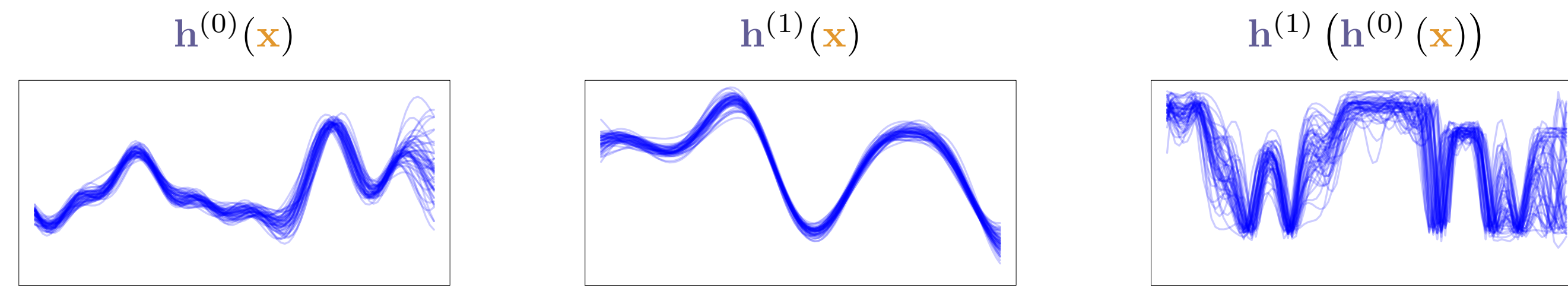


Fig. 1: Illustration of how stochastic processes may be composed.

- ▶ Inference requires calculating the marginal likelihood:

$$p(Y|X, \theta) = \int p(Y|F^{(N_h)}, \theta^{(N_h)}) \times p(F^{(N_h)}|F^{(N_h-1)}, \theta^{(N_h-1)}) \times \dots \times p(F^{(1)}|X, \theta^{(0)}) dF^{(N_h)} \dots dF^{(1)};$$

- ▶ **Extremely challenging!**

DGPs with Random Features

- ▶ GPs are single-layered Neural Nets with an infinite number of hidden units;
- ▶ Taking a weight-space view of a GP:

$$F = \Phi W;$$

- ▶ The priors over the weights are:

$$p(W_{.i}) = \mathcal{N}(\mathbf{0}, I);$$

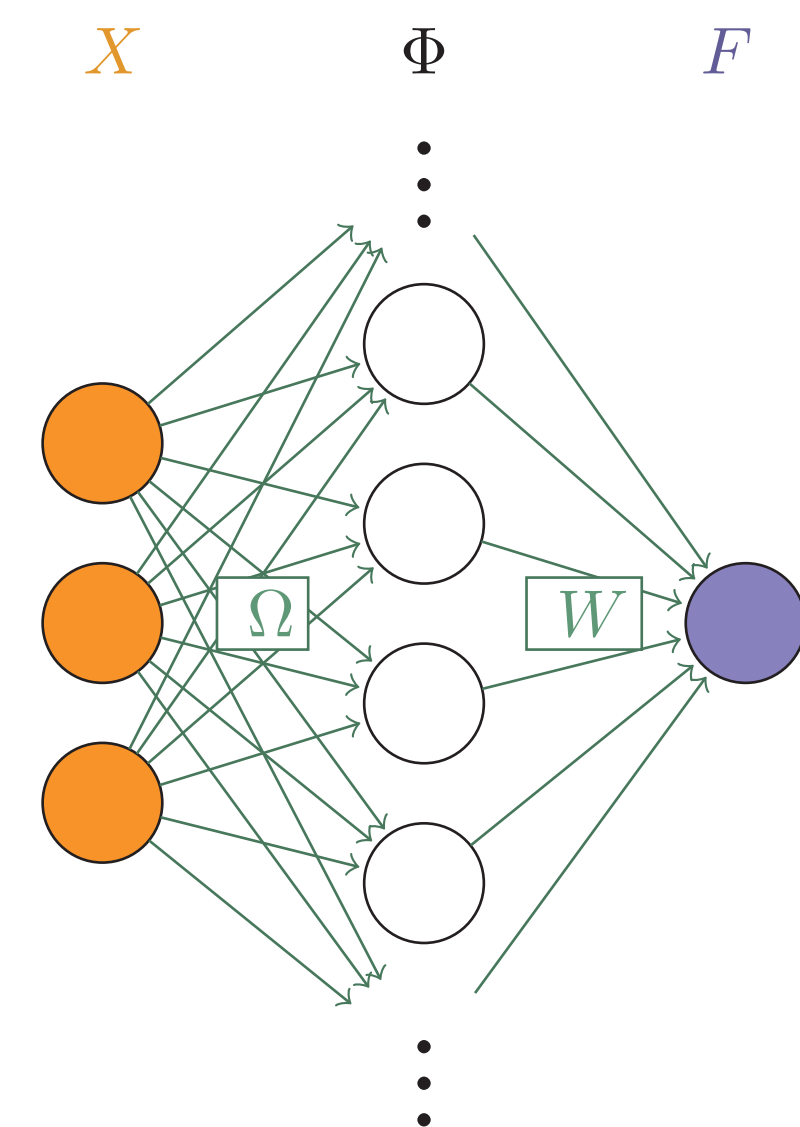


Fig. 2: Single-layered GP.

- ▶ The RBF kernel:

$$k_{\text{rbf}}(\mathbf{x}, \mathbf{x}') = \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}') \right] = \int p(\omega) \exp(\iota(\mathbf{x} - \mathbf{x}')^\top \omega) d\omega$$

can be approximated using trigonometric functions:

$$\Phi_{\text{rbf}} = \sqrt{\frac{\sigma^2}{N_{\text{RF}}}} [\cos(F\Omega), \sin(F\Omega)] \quad \text{with} \quad p(\Omega_{.j}|\theta) = \mathcal{N}(\mathbf{0}, \Lambda^{-1}),$$

allowing for scaling factors σ^2 and $\Lambda = \text{diag}(\iota_1^2, \dots, \iota_d^2)$ for the kernel and the features;

- ▶ Meanwhile, the first order Arc-Cosine kernel:

$$k_{\text{arc}}(\mathbf{x}, \mathbf{x}') = \frac{1}{\pi} \|\mathbf{x}\| \|\mathbf{x}'\| J(\alpha) = 2 \int \max(\mathbf{0}, \omega^\top \mathbf{x}) \max(\mathbf{0}, \omega^\top \mathbf{x}') \mathcal{N}(\omega|\mathbf{0}, I) d\omega \quad \text{where}$$

$$J(\alpha) = \sin \alpha + (\pi - \alpha) \cos \alpha \quad \text{and} \quad \alpha = \cos^{-1} \left(\frac{\mathbf{x}^\top \mathbf{x}'}{\|\mathbf{x}\| \|\mathbf{x}'\|} \right)$$

can be approximated using Rectified Linear Units (ReLU):

$$\Phi_{\text{arc}} = \sqrt{\frac{2\sigma^2}{N_{\text{RF}}}} \max(\mathbf{0}, F\Omega) \quad \text{with} \quad p(\Omega_{.j}|\theta) = \mathcal{N}(\mathbf{0}, \Lambda^{-1}).$$

Model Architecture

- ▶ DGPs with random features become DNNs with low-rank weight matrices!

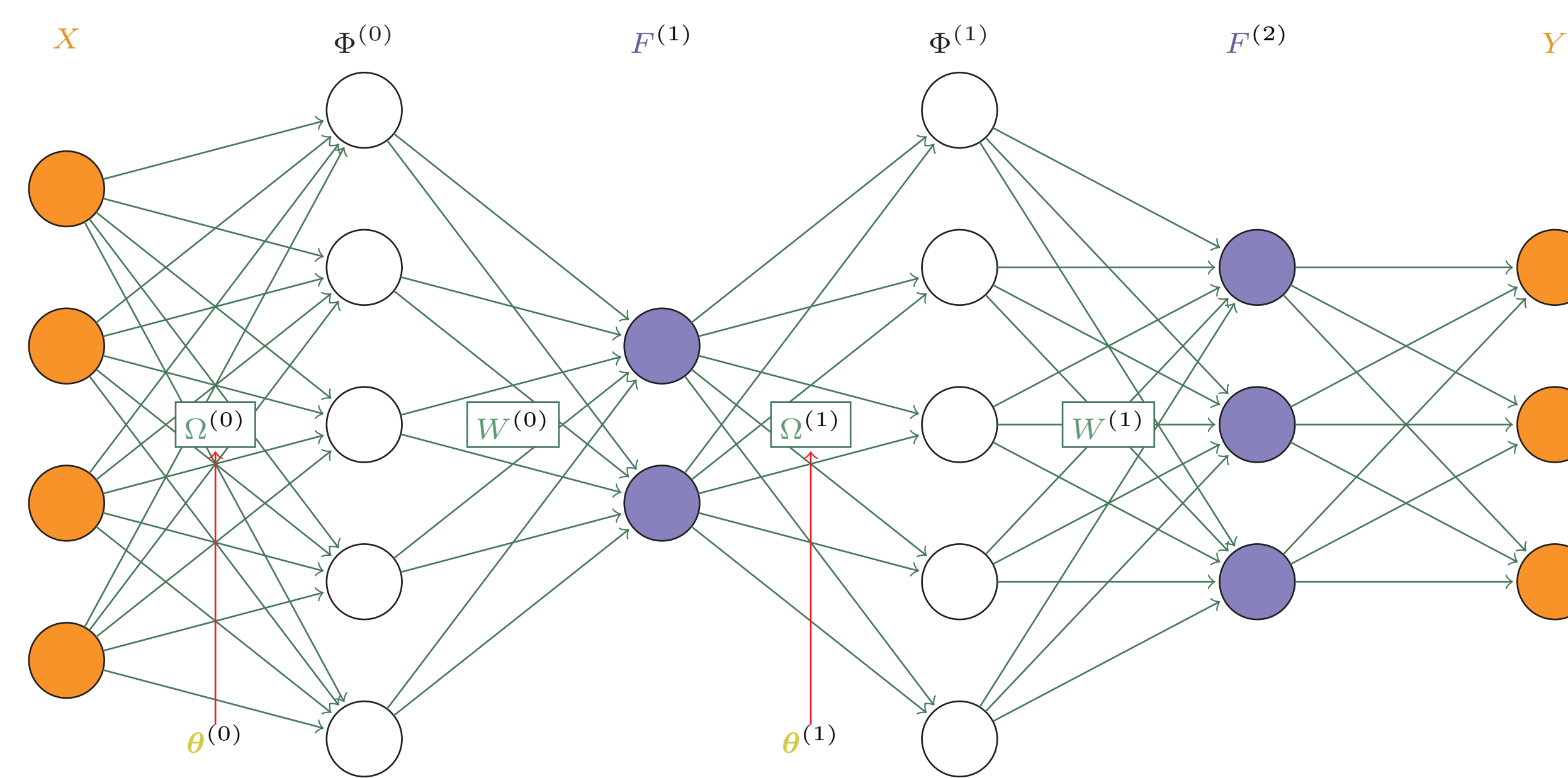


Fig. 3: Diagram of the proposed DGP model with random features.

Stochastic Variational Inference

- ▶ Define $\Psi = (\Omega^{(0)}, \dots, \Omega^{(L)}, W^{(0)}, \dots, W^{(L)});$

- ▶ Lower bound on marginal likelihood:

$$\log [p(Y|\theta)] \geq E_{q(\Psi)} (\log [p(Y|\Psi)]) - D_{\text{KL}} [q(\Psi) \| p(\Psi|\theta)],$$

where $q(\Psi)$ approximates $p(\Psi|Y, \theta);$

- ▶ Factorized approximate posterior:

$$q(\Psi) = \prod_{ijl} q(\Omega_{ij}^{(l)}) \prod_{ijl} q(W_{ij}^{(l)}),$$

with

$$q(W_{ij}^{(l)}) = \mathcal{N}(\mu_{ij}^{(l)}, (\sigma^2)_{ij}^{(l)}) \quad \text{and} \quad q(\Omega_{ij}^{(l)}) = \mathcal{N}(m_{ij}^{(l)}, (s^2)_{ij}^{(l)});$$

- ▶ Assuming factorized likelihood, we can use **mini-batch** stochastic gradient optimization:

$$\frac{n}{m} \sum_{k \in \mathcal{I}_m} E_{q(\Psi)} (\log [p(y_k|\Psi)]) - D_{\text{KL}} [q(\Psi) \| p(\Psi|\theta)];$$

- ▶ The expectation can be estimated using Monte Carlo:

$$E_{q(\Psi)} (\log [p(y_k|\Psi)]) \approx \frac{1}{N_{\text{MC}}} \sum_{r=1}^{N_{\text{MC}}} \log [p(y_k|\tilde{\Psi}_r)],$$

with $\tilde{\Psi}_r \sim q(\Psi);$

- ▶ Computational cost dominated by low-rank matrix multiplication - no inverses required;

- ▶ Various optimization strategies for Ω available.

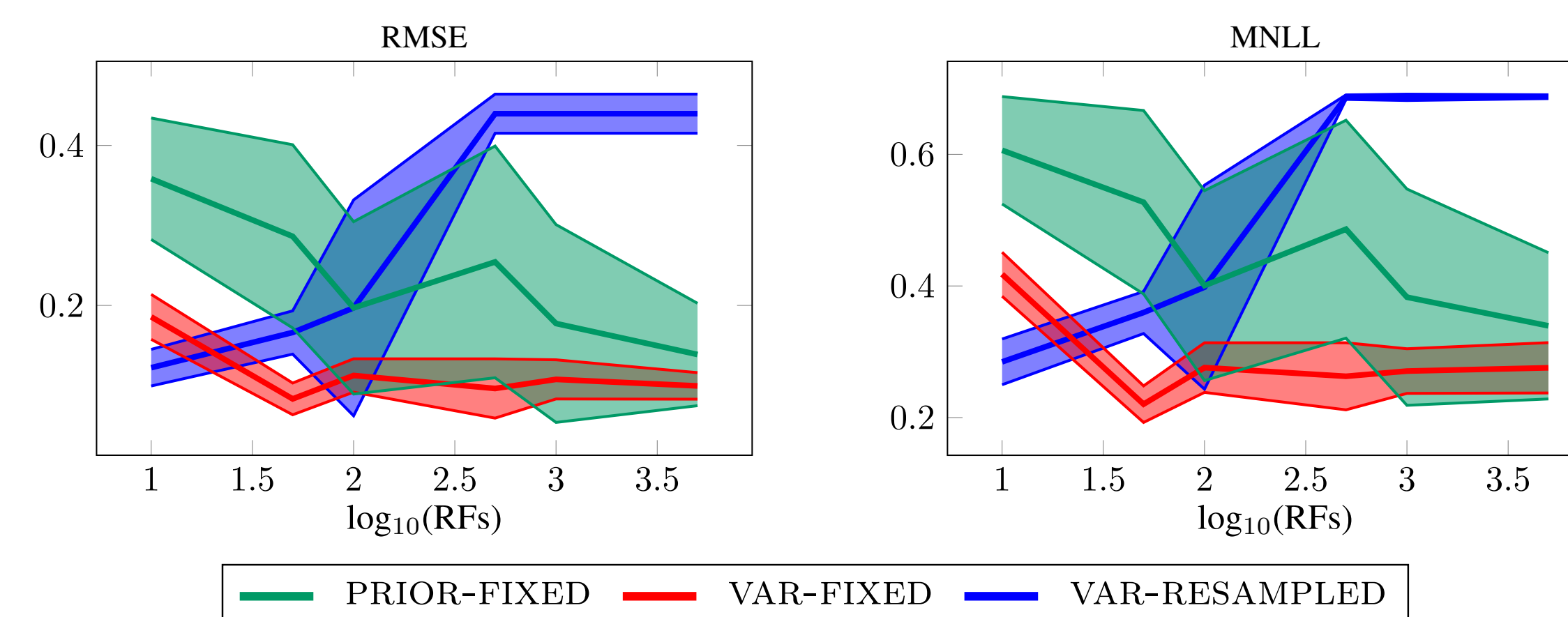


Fig. 4: Performance of different strategies for dealing with Ω as a function of the number of random features. These can be fixed (PRIOR-FIXED), or treated variationally (with fixed randomness VAR-FIXED or resampled at each iteration VAR-RESAMPLED).

Experimental Setup and Results

- ▶ Competing methods:

- DGP-RF with RBF Kernel (DGP-RBF);
- DGP-RF with first order Arc-Cosine Kernel (DGP-ARC);
- DGP with Expectation Propagation (DGP-EP);
- DNN with Dropout (DNN);
- Sparse GP with Variational Inference (VAR-GP).

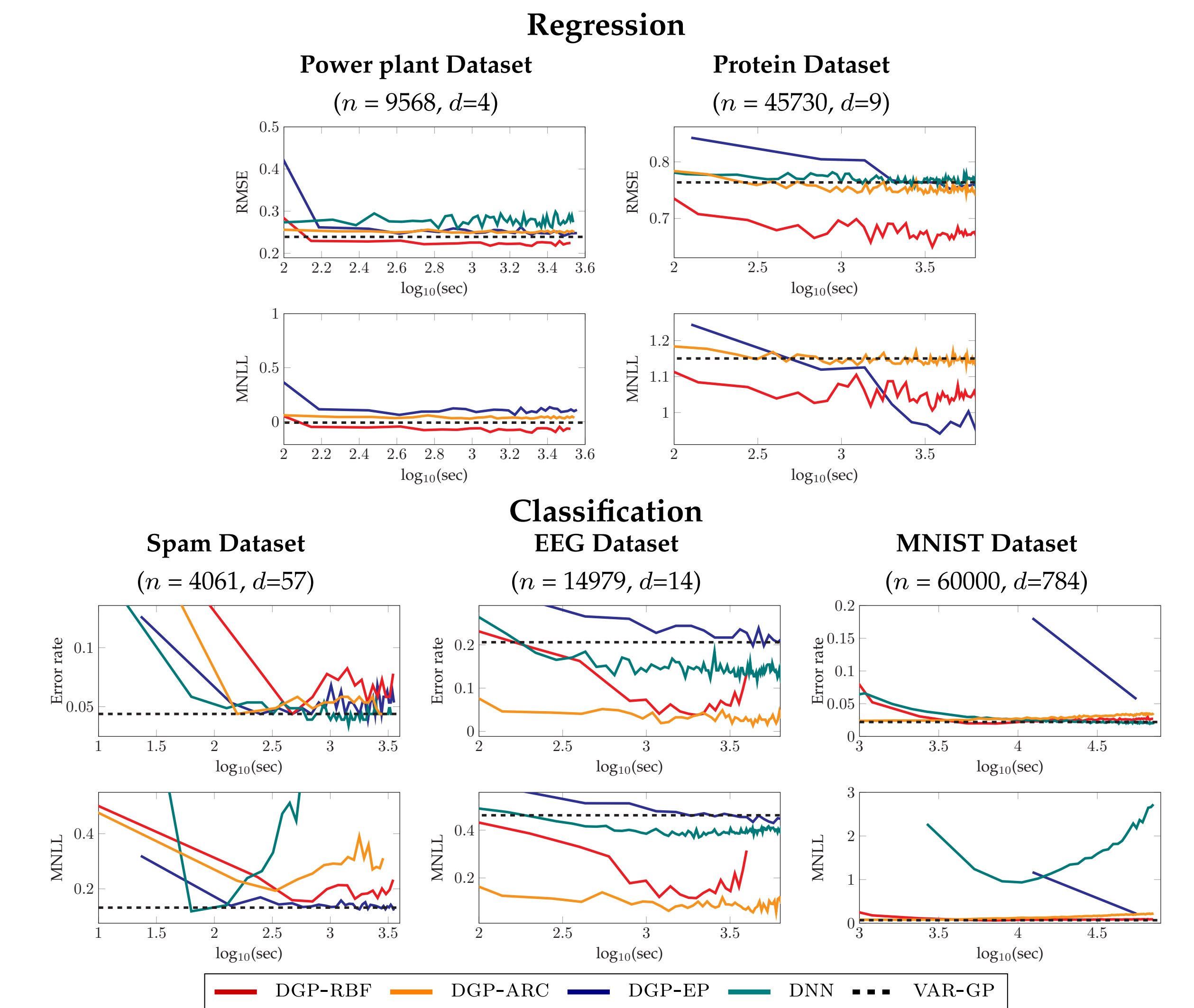


Fig. 5: Progression of RMSE and MNLL on test data over time for competing models.

- ▶ Easily extendable architecture with option to feed-forward inputs.

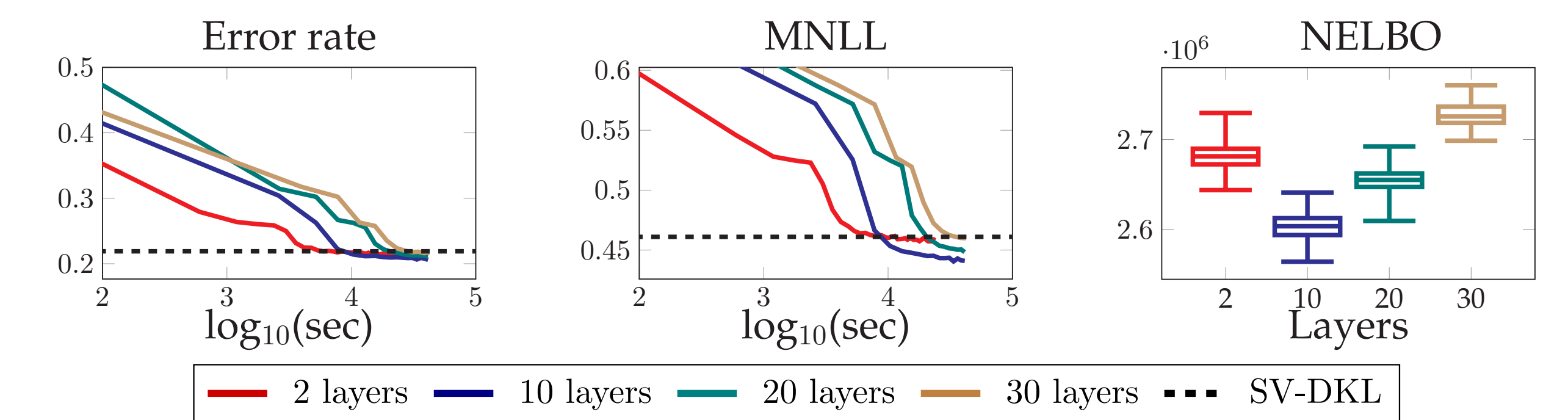


Fig. 6: Left and center - Performance of our model on the AIRLINE dataset as function of time for different depths. Right - The box plot of the negative evidence lower bound confirms this is a suitable objective for model selection.

Conclusions

- ▶ Our contributions:

- ✓ Complete specification and evaluation of DGPs based on random features;
- ✓ Scalable and practical DGP inference - no matrix inverses;
- ✓ Synchronous/asynchronous distributed implementation available;

- ▶ Ongoing work:

- Fastfood and other kernels;
- Convolutional GP layers for image processing.

References

- [1] A. Damianou and N. Lawrence. *Deep Gaussian Processes*, AISTATS 2013.
- [2] Y. Gal and Z. Ghahramani. *Dropout as a Bayesian Approximation*, ICML 2016.