



# Graph Random Features for Scalable Gaussian Processes



Matthew Zhang<sup>1</sup>, Jihao Andreas Lin<sup>1,2</sup>, Krzysztof Choromanski<sup>3</sup>, Adrian Weller<sup>1,4</sup>, Richard E. Turner<sup>1,4</sup>, Isaac Reid<sup>1,3</sup>

## Efficient GPs on graphs

1. Under mild assumptions on the graph, the kernel estimate  $\hat{\mathbf{K}} = \Phi\Phi^\top$  is provably **sparse** and has a **bounded condition number**.

2. This lets us **estimate** quantities involving the kernel inverse like the posterior mean and variance

$$m_{|y}(x) = m(x) + k(x, x)[k(x, x) + \sigma_n^2 \mathbf{I}]^{-1}(y - m(x)),$$

$$k_{|y}(x, x') = k(x, x') - k(x, x)[k(x, x) + \sigma_n^2 \mathbf{I}]^{-1}k(x, x').$$

very efficiently using iterative linear system solvers, e.g. **conjugate gradient methods**.

5. Combining with **pathwise conditioning**

$$g_{|y}(\cdot) = g(\cdot) + \hat{\mathbf{K}}_{(\cdot)x}(\hat{\mathbf{K}}_{xx} + \sigma_n^2 \mathbf{I})^{-1}(y - (g(x) + \epsilon)).$$

we can efficiently approximate **inverse kernel-vector products**, unlocking Bayesian optimisation on massive graphs in  $\mathcal{O}(N^{3/2})$  time.

## Applications and experiments

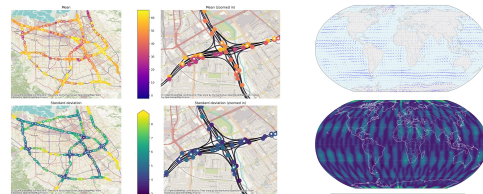


Fig 3. Regression tasks: traffic prediction and wind interpolation. GRFs for uncertainty-aware regression tasks real world graph datasets. Left: traffic speed prediction. Right: Wind velocity interpolation.

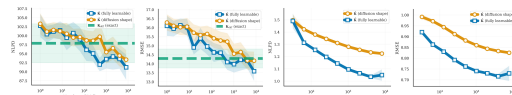


Fig 4. GRFs surpass the exact diffusion GPs with ~500 walkers per node. NLPD and RMSE versus the number of random walkers. From left to right: traffic NLPD, RMSE, wind NLPD, RMSE.

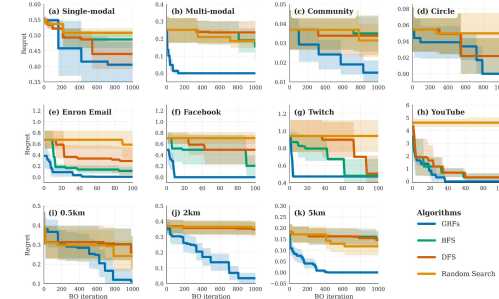


Fig 5. Bayesian optimisation. GRFs beat other scalable baselines for Bayesian optimisation tasks on graphs datasets with over 1M nodes.

You can efficiently **approximate** the graph node kernel by **sampling random walks**, which unlocks **scalable Bayesian optimisation** on graphs with **>1M nodes** on a **single computer chip**.



## Graph random features

1. Many graph node kernels are expressed as functions of a **weighted adjacency matrix**

$$\mathbf{K}_\alpha(\mathbf{W}) = \sum_{r=0}^{\infty} \alpha_r \mathbf{W}^r, \quad \alpha_r \in \mathbb{R} \quad \forall r \in (0, 1, \dots, \infty).$$

2. Powers of  $\mathbf{W}$  count random walks, so we can estimate this by importance sampling (a.k.a. ‘graph random features’, or **GRFs**)

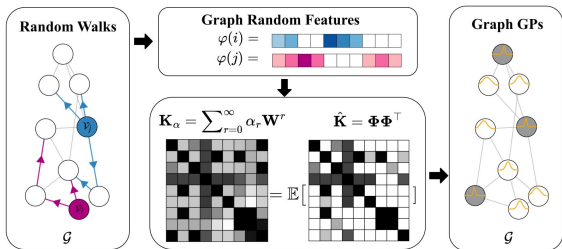


Fig 1. Schematic. GRFs for Bayesian Inference schematic

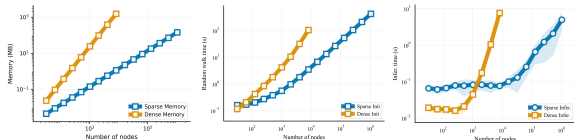


Fig 2. Scaling curves. GRF-GPs have sub-quadratic time scaling and linear memory scaling. From left to right: the scaling of memory, kernel computation (random walk sampling) time and inference time respectively.