

# An Improved Variational Approximate Posterior for the Deep Wishart Process

Sebastian W. Ober



Laurence Aitchison



Adam Yang



Edward Milsom



Ben Anson



shallow

deep

---

feature

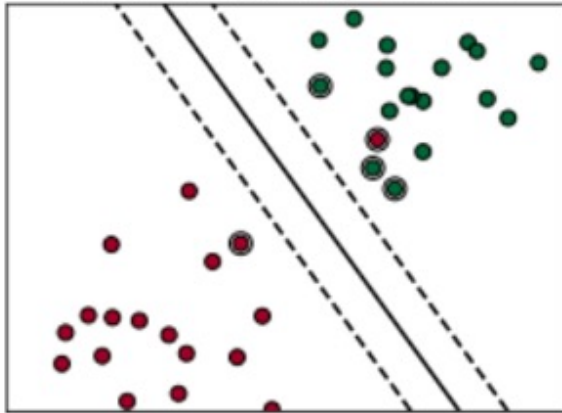
linear  
regression

kernel

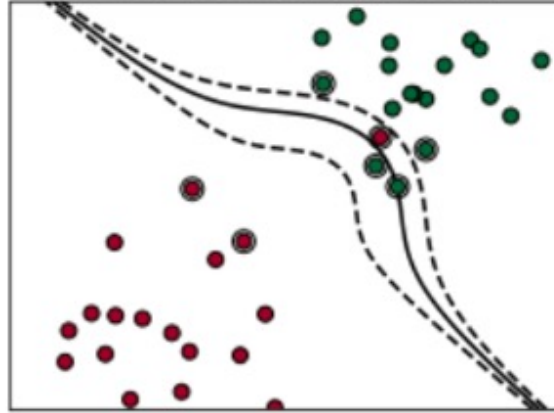
kernel ridge  
regression

For good performance, we need to choose a  
choosing a good feature extractor / kernel

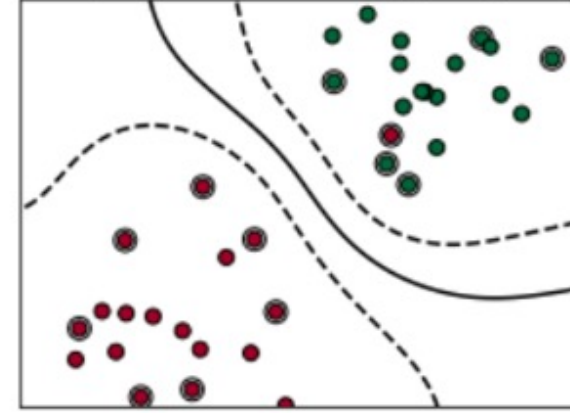
linear kernel



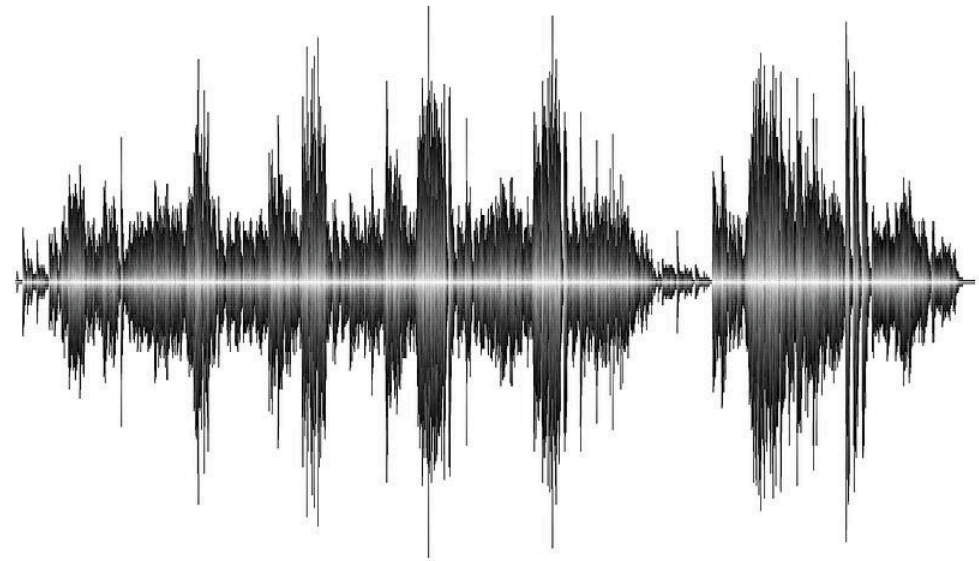
polynomial kernel



RBF kernel



Problem: can't choose a good feature extractor/kernel for complex data like images



shallow

---

feature

linear  
regression

kernel

kernel ridge  
regression

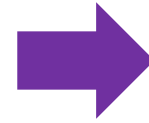
shallow

deep

---

feature

linear  
regression



neural net

Multiple layers

Flexibility at each layer

kernel

kernel ridge  
regression

shallow

deep

feature

linear  
regression



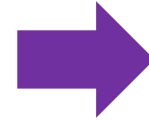
neural net

Multiple layers

Flexibility at each layer

kernel

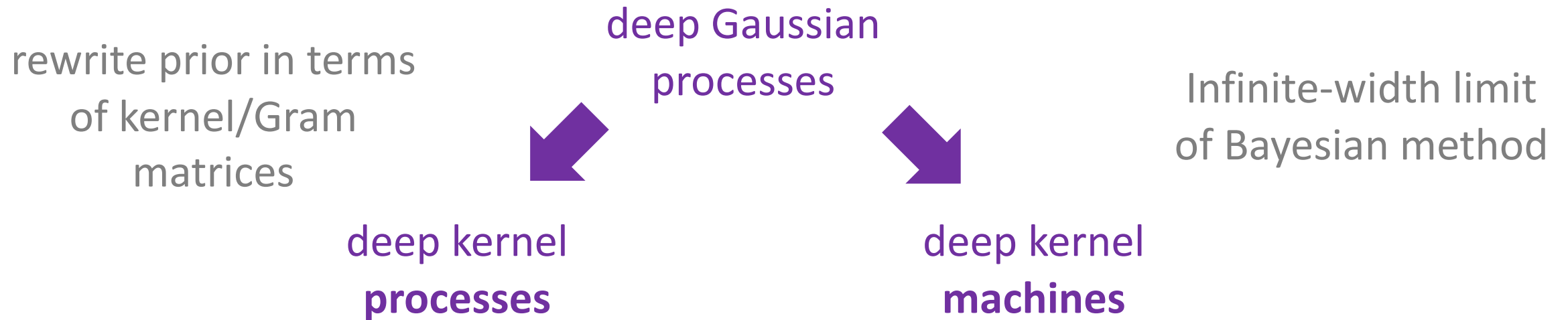
kernel ridge  
regression



**deep kernel  
methods**



# Summary



# Part 1



# In deep-kernel methods, we switch to working entirely with Gram matrices

$P$  = number of datapoints  
 $N_\ell$  = width of layer  $\ell$

## Gram matrices

$$\mathbf{G}_2 = \mathbf{F}_2 \mathbf{F}_2^T / N_2$$

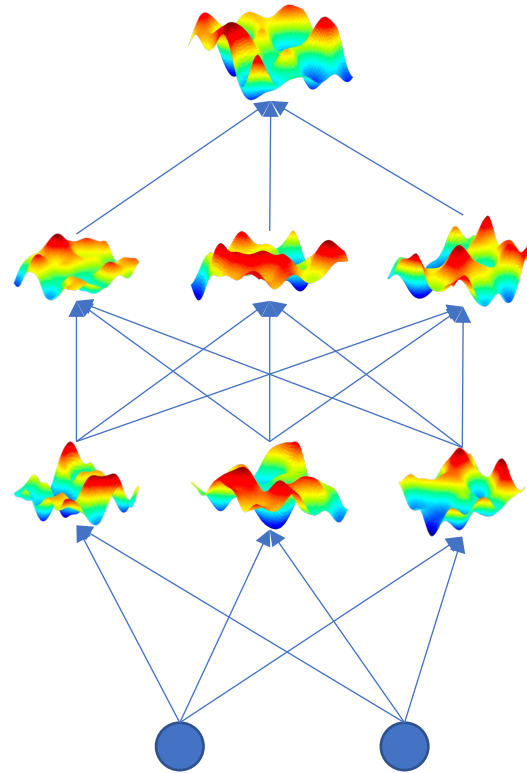
$P \times P$

$$\mathbf{G}_1 = \mathbf{F}_1 \mathbf{F}_1^T / N_1$$

$P \times P$

$$\mathbf{G}_0 = \mathbf{X} \mathbf{X}^T / N_X$$

$P \times P$



## DGP

outputs,  
 $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f(\mathbf{F}_2) + \sigma^2 \mathbf{I})$

$P \times P$

hidens,  
 $\mathbf{F}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f(\mathbf{F}_1))$

$P \times N_2$        $P \times P$

hidens,  
 $\mathbf{F}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f(\mathbf{X}))$

$P \times N_1$        $P \times P$

batch of  
input vectors,  $\mathbf{X}$

$P \times N_X$

# Trick 1: most kernels of interest can be computed from the Gram matrix

- True for e.g. arccos kernels used in infinite NNs (Cho and Saul 2009)
- Also true for standard GP kernels that only depend on distance between datapoints  $i$  and  $j$ , because we can recover distance from the Gram matrix, (Duvenaud et al. 2014)

$$\begin{aligned}R_{ij}(\mathbf{G}) &= \frac{1}{N} \sum_{\lambda=1}^N (F_{i\lambda} - F_{j\lambda})^2 \\ &= \frac{1}{N} \sum_{\lambda=1}^N ((F_{i\lambda})^2 - 2F_{i\lambda}F_{j\lambda} + (F_{j\lambda})^2) \\ &= G_{ii} - 2G_{ij} + G_{jj}\end{aligned}$$

- Overall:

$$\mathbf{K}_f(\mathbf{F}_\ell) = \mathbf{K}(\mathbf{G}_\ell)$$

## Trick 2: Gram matrices are Wishart distributed

To get next Gram matrix, we first sample a bunch of features,

$$\mathbf{F}_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{G}_{\ell-1}))$$

And then compute the Gram matrix

$$\mathbf{G}_\ell = \frac{1}{N_\ell} \mathbf{F}_\ell \mathbf{F}_\ell^T$$

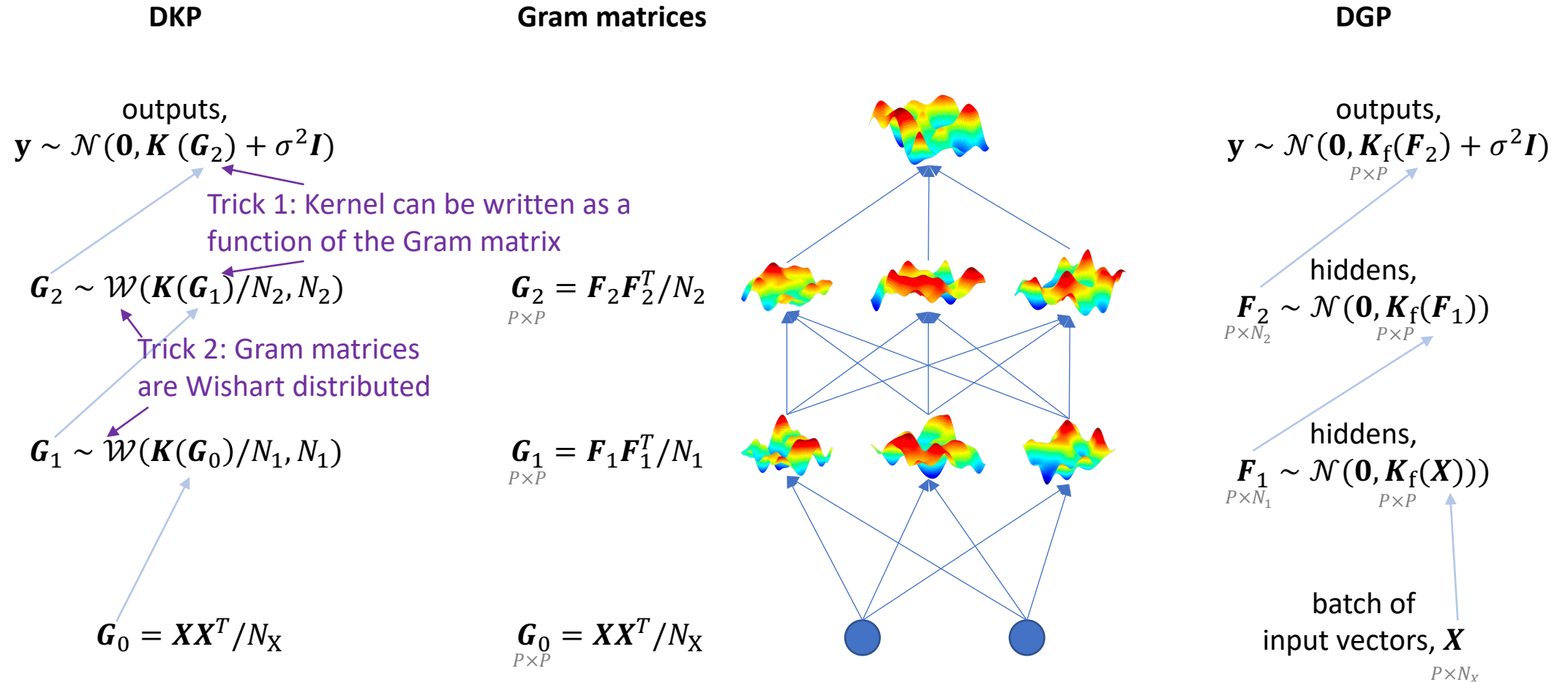
But this exactly matches the definition of the Wishart distribution!

$$\mathbf{G}_\ell \sim \mathcal{W}(\mathbf{K}(\mathbf{G}_{\ell-1})/N_\ell, N_\ell)$$

(see Wikipedia for pdf, moments etc.)

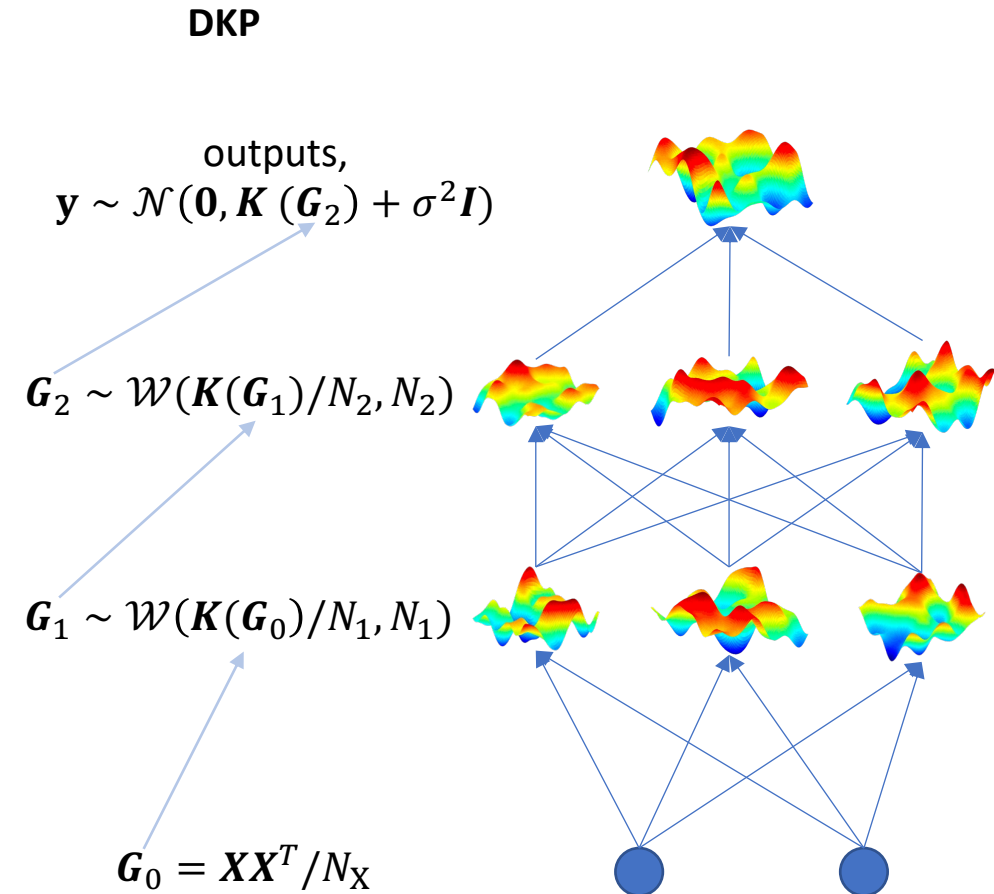
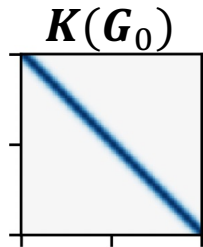
# In deep-kernel methods, we switch to working entirely with Gram matrices

$P$  = number of datapoints  
 $N_\ell$  = width of layer  $\ell$



# Sampling the prior in the kernelized DGP

- Next Gram matrix is “centered” on the kernel,  
 $E[\mathbf{G}_1 | \mathbf{G}_0] = \mathbf{K}(\mathbf{G}_0)$



# Developing practical methods + our results

We developed:

- Two processes: “deep Wishart process” and “deep *inverse* Wishart process”
- VI with priors + approximate posteriors over Gram matrices, *not* features.
- a bunch of approximate posteriors (e.g.  $Q_{GW}$  )

	Dataset	DGP	$Q_{GW}$	DWP $Q_{A-GW}$	$Q_{AB-GW}$
	BOSTON	$-2.43 \pm 0.04$	<b><math>-2.38 \pm 0.04</math></b>	$-2.39 \pm 0.05$	<b><math>-2.38 \pm 0.04</math></b>
	CONCRETE	$-3.13 \pm 0.02$	$-3.13 \pm 0.02$	<b><math>-3.07 \pm 0.02</math></b>	$-3.08 \pm 0.02$
	ENERGY	$-0.71 \pm 0.03$	$-0.71 \pm 0.03$	<b><math>-0.70 \pm 0.03</math></b>	<b><math>-0.70 \pm 0.03</math></b>
	KIN8NM	$1.38 \pm 0.00$	$1.40 \pm 0.01$	<b><math>1.41 \pm 0.01</math></b>	<b><math>1.41 \pm 0.01</math></b>
LL	NAVAL	$8.28 \pm 0.04$	$8.17 \pm 0.07$	<b><math>8.40 \pm 0.02</math></b>	$8.10 \pm 0.19$
	POWER	$-2.78 \pm 0.01$	$-2.77 \pm 0.01$	<b><math>-2.76 \pm 0.01</math></b>	<b><math>-2.76 \pm 0.01</math></b>
	PROTEIN	$-2.73 \pm 0.01$	$-2.72 \pm 0.01$	$-2.71 \pm 0.01$	<b><math>-2.70 \pm 0.00</math></b>
	WINE	$-0.96 \pm 0.01$	$-0.96 \pm 0.01$	$-0.96 \pm 0.01$	$-0.96 \pm 0.01$
	YACHT	$-0.73 \pm 0.07$	$-0.58 \pm 0.06$	$-0.22 \pm 0.09$	<b><math>-0.18 \pm 0.07</math></b>

[1] Aitchison, Yang and Ober. “Deep kernel processes” ICML (2021)

[2] Ober and Aitchison “An approximate posterior for the deep Wishart process” NeurIPS (2021)

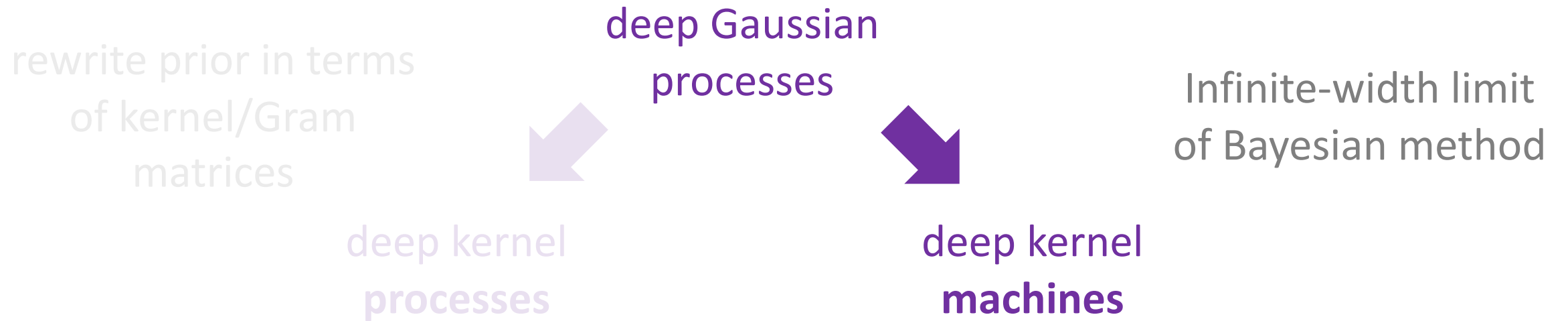
[3] Ober, Anson, Milsom and Aitchison “An improved approximate posterior for the deep Wishart process” UAI (2023)



# Part 1



# Part 2



# Why less(?) Bayesian deep kernel machines?

Less(?) Bayesian approach:

- simplifies implementation
- gives lower-variance updates that converge faster
- provides a cleaner link to NN / neuro-theory
- great preliminary results...

# We get a deep kernel machine by taking an infinite-width limit of a DGP

- True posterior over features becomes multivariate Gaussian [1]

$$P(\mathbf{F}_1, \dots, \mathbf{F}_L | \mathbf{X}, \mathbf{Y}) = \prod_{\ell=1}^L \sum_{\lambda=1}^{N_\ell} \mathcal{N}(\mathbf{f}_\lambda^\ell; \mathbf{0}, \mathbf{G}_\ell^*)$$

- We choose a family of approximate posteriors capturing the true posterior:

$$Q(\mathbf{F}_1, \dots, \mathbf{F}_L) = \prod_{\ell=1}^L \sum_{\lambda=1}^{N_\ell} \mathcal{N}(\mathbf{f}_\lambda^\ell; \mathbf{0}, \mathbf{G}_\ell)$$

- Gram matrices,  $\mathbf{G}_1, \dots, \mathbf{G}_L$ , are the same kind of thing as in deep kernel process!

$$\mathbf{G}_\ell = \frac{1}{N_\ell} \mathbf{F}_\ell \mathbf{F}_\ell^T$$

- But here, Gram matrices appear as parameters of approximate posterior
- So to find the Gram matrices, we optimize the ELBO!

# A DKM is an infinite-width limit of a DGP!

ELBO:

$$\mathcal{L}(\mathbf{G}_1, \dots, \mathbf{G}_L) = \underbrace{\log P(\mathbf{Y} | \mathbf{G}_L)}_{\text{likelihood}} - \sum_{\ell=1}^L v_{\ell} D_{\text{KL}}(\underbrace{\mathcal{N}(0, \mathbf{G}_{\ell})}_{\text{approx post}} \parallel \underbrace{\mathcal{N}(0, \mathbf{K}(\mathbf{G}_{\ell-1}))}_{\text{prior}})$$

- Optimizes intermediate layer Gram matrices,
- Encourages good performance
- Keeps approximate posterior covariance,  $\mathbf{G}_{\ell}$ , similar to prior covariance,  $\mathbf{K}(\mathbf{G}_{\ell-1})$

The objective only talks about  $P \times P$  Gram matrices,  $\mathbf{G}_{\ell}$ , and kernel matrices,  $\mathbf{K}(\mathbf{G}_{\ell-1})$ . So it is a deep kernel method!

# What is a deep kernel machine?

- A nonlinear function approximator
- With multiple layers
- Parameterised by *Gram matrices*, not features or weights
- Trained using the DKM objective:

$$\mathcal{L}(\mathbf{G}_1, \dots, \mathbf{G}_L) = \underbrace{\log P(\mathbf{Y} | \mathbf{G}_L)}_{\text{likelihood}} - \sum_{\ell=1}^L v_{\ell} D_{\text{KL}} \left( \underbrace{\mathcal{N}(0, \mathbf{G}_{\ell})}_{\text{approx post}} \parallel \underbrace{\mathcal{N}(0, \mathbf{K}(\mathbf{G}_{\ell-1}))}_{\text{prior}} \right)$$

- Optimizes intermediate layer Gram matrices,
- Encourages good performance
- Keeps approximate posterior covariance,  $\mathbf{G}_{\ell}$ , similar to prior covariance,  $\mathbf{K}(\mathbf{G}_{\ell-1})$

# Preliminary results for convolutional deep kernel machines

Table 3: Comparison of test accuracy against other kernel methods, including NNGP and NTK.

Paper	Method	CIFAR-10
This paper	DKM-DA-GAP (512 / 1024 / 2048)	91.98%
Novak et al (2018)	NNGP-GAP	77.43%
Arora et al (2019)	NNGP-GAP	83.75%
Lee et al (2020)	NNGP-GAP-DA	84.8%
Li et al (2019)	NNGP-LAP-flip	88.92%
Shankar et al (2020)	Myrtle10	89.80%
Adlam et al (2023)	Tuned Myrtle10 DA CG	91.2%

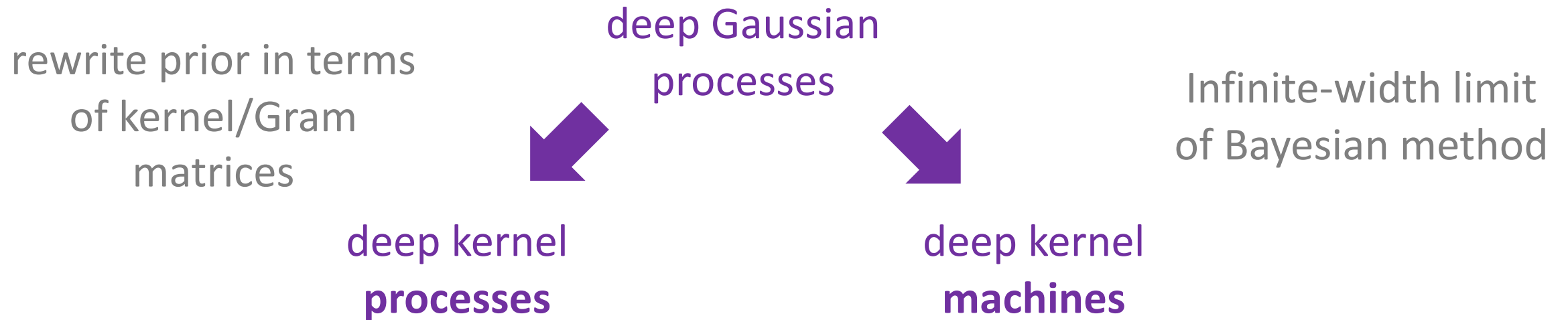


Edward Milsom

But how slow are DKMs? Surprisingly fast!

- We develop a novel inducing-point scheme
- Same FLOPs as CNN (computations ultimately look v. similar)
- Slower than a CNN (about 30 hours for largest model), as we're using float64 FLOPs
- Orders of magnitude faster than "full" kernel methods above.

# Summary





# Deep kernel landscape + our priorities

## Our priorities

- More architectures for DKMs (GNNs + transformers).
- Understanding Bayesian-ness of DKMs
- speed/scale-up:
  - memory efficiency
  - lower-precision
- user-friendly library (we can share preliminary work)

## Huge future opportunities:

	shallow	deep
feature	linear regression	neural net
kernel	kernel ridge regression	<b>deep kernel methods</b>

If you're interested, get in touch:  
[laurence.aitchison@bristol.ac.uk](mailto:laurence.aitchison@bristol.ac.uk)

[1] Aitchison, Yang and Ober. "Deep kernel processes" ICML (2021)

[2] Ober and Aitchison "An approximate posterior for the deep Wishart process" NeurIPS (2021)

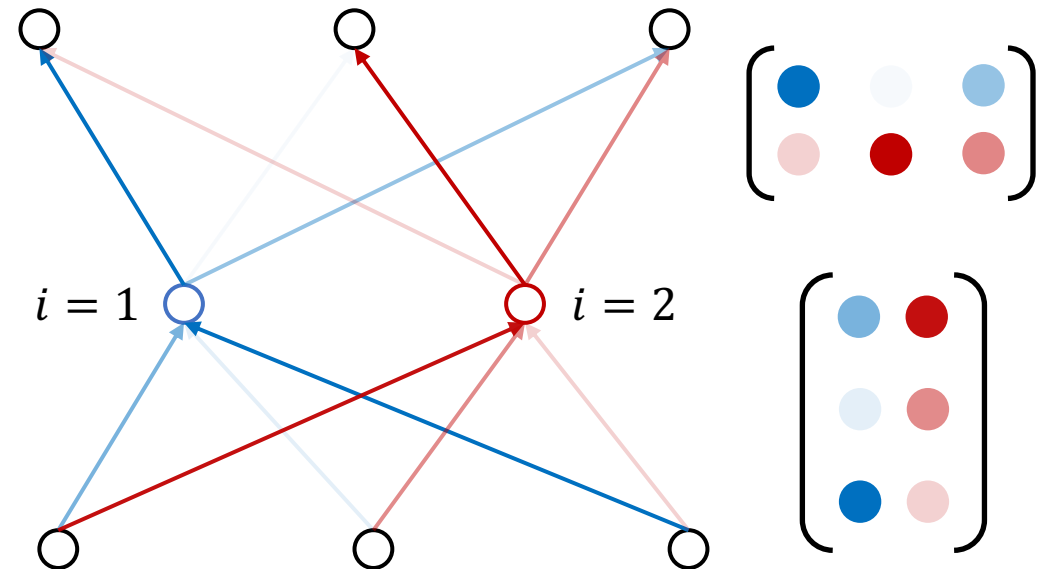
[3] Ober, Anson, Milsom and Aitchison "An improved approximate posterior for the deep Wishart process" UAI (2023)

[4] Yang, Robeyns, Milsom, Anson, Schoots, Aitchison "A theory of representation learning gives a deep generalisation of kernel methods" ICML (2023)

[5] Milsom, Anson, Aitchison "Convolutional deep kernel machines" (in prep)

Appendix slides

Deep kernel processes should work better because they have fewer local optima



# Deep kernel processes should work better because they have fewer local optima

- Implies loads of symmetric local optima...
- ...and local optima are bad if you have unimodal approximate posteriors.
- DKPs don't have these symmetries, so *far* fewer local optima!

