

(Convolutional) deep kernel machines

Laurence Aitchison
Assistant Prof.
University of Bristol



Sebastian Ober



Adam Yang



Edward Milsom



Ben Anson



Neural
Networks



Kernel Methods



Edward
Milsom

shallow

deep

feature

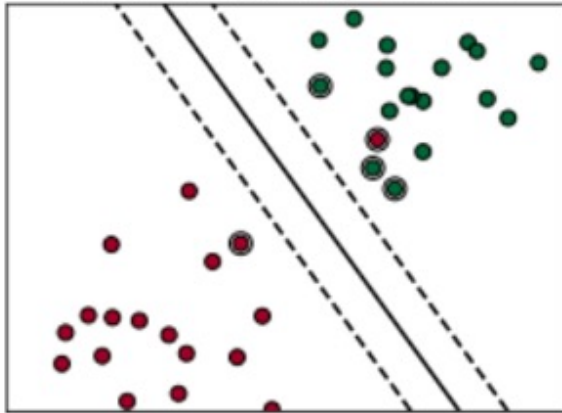
linear
regression

kernel

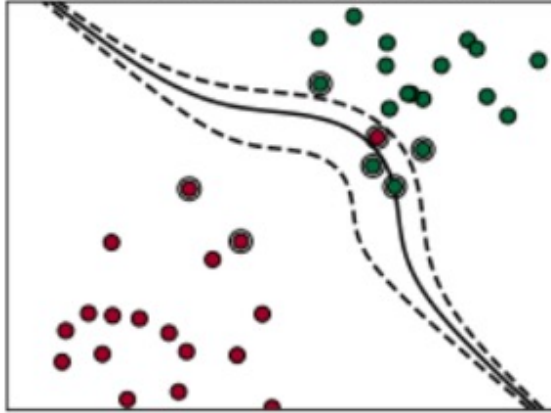
kernel ridge
regression

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

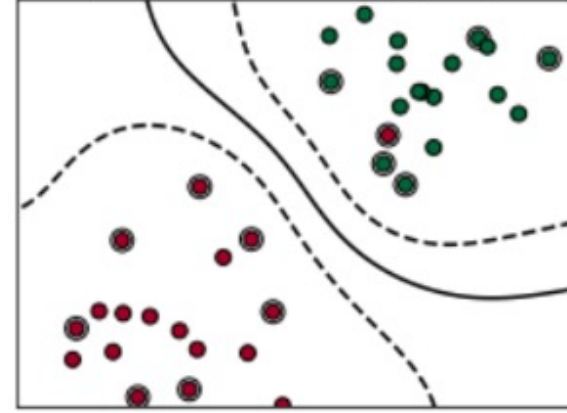
linear kernel



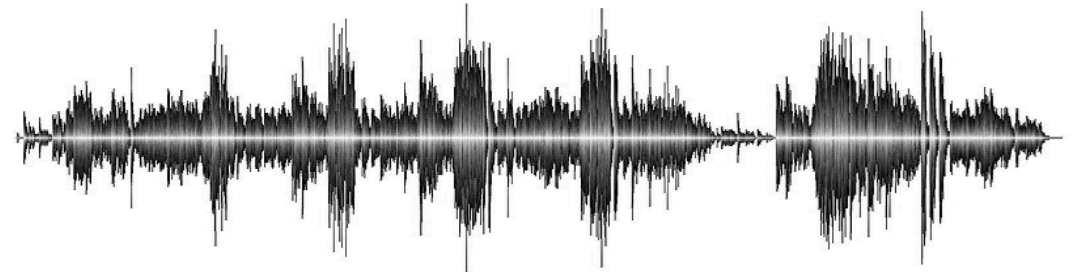
polynomial kernel



RBF kernel



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



need to learn a representation
(i.e. feature extractor / kernel!)

shallow

feature

linear
regression

kernel

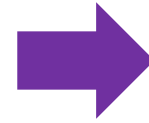
kernel ridge
regression

shallow

deep

feature

linear
regression



neural net

Multiple layers

Flexibility at each layer

kernel

kernel ridge
regression

shallow

deep

feature

linear
regression



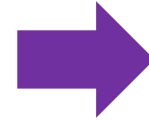
neural net

Multiple layers

Flexibility at each layer

kernel

kernel ridge
regression



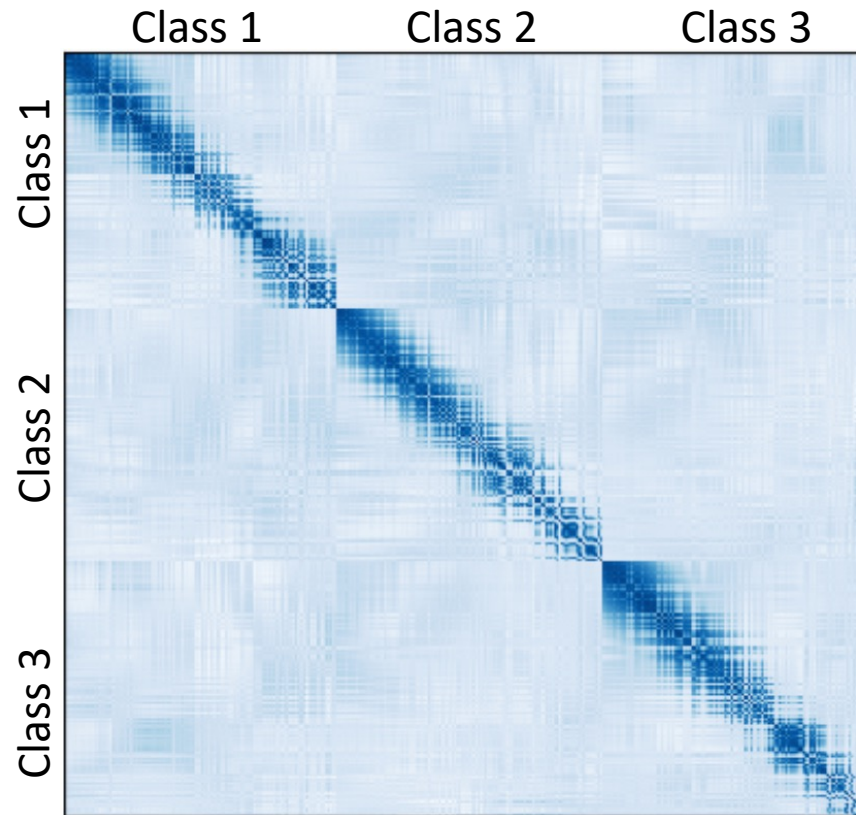
**deep kernel
methods**

no representation
learning

has representation
learning

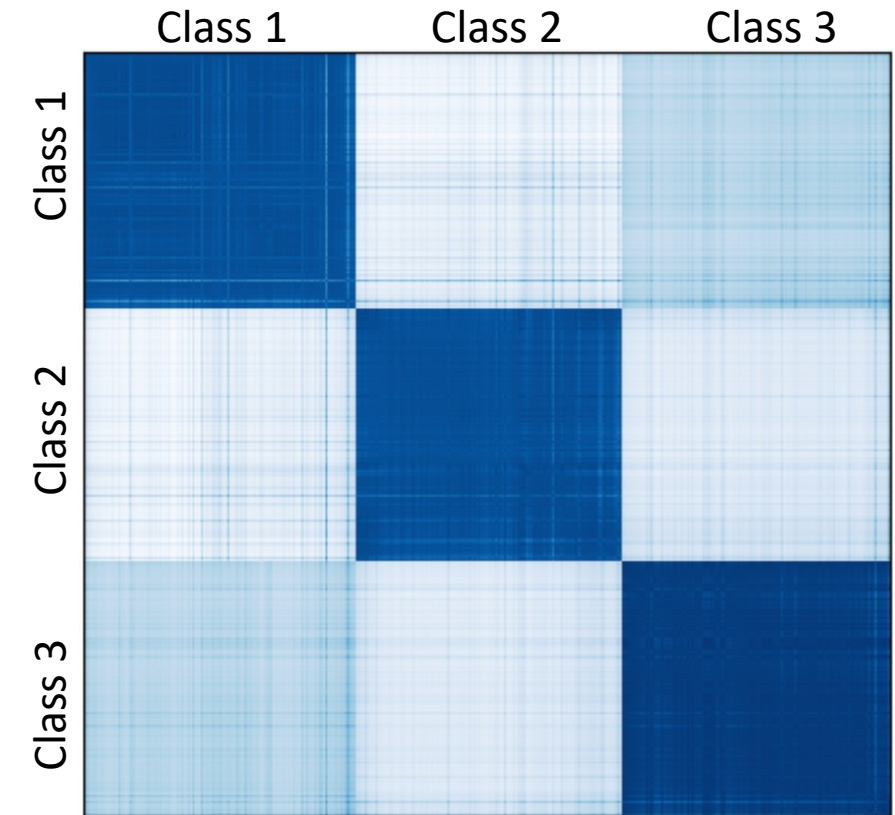
What representation learning looks like

- dot product of neural activations, for all pairs of training examples.
 - intuition: basically the cosine similarity
- This matrix = the kernel!



before/no representation learning:

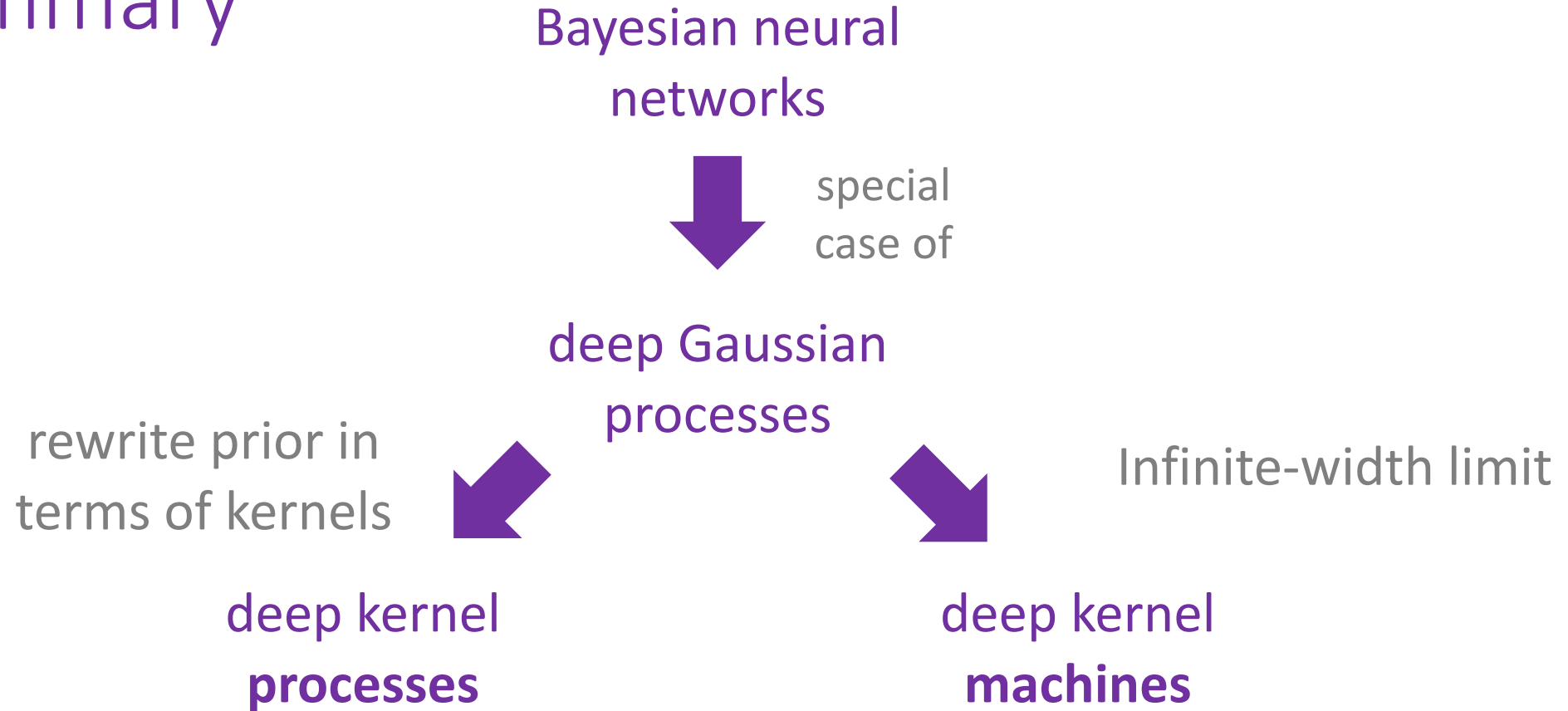
- Hand-picked kernel/features
- Finite randomly initialized NN
- Infinite randomly initialized network (NNGP/NTK)



learned representation learning:

- Trained NN
- Trained deep kernel method.

Summary



Part 0

Bayesian neural
networks



special
case of

deep Gaussian
processes

rewrite prior in
terms of kernels



deep kernel
processes



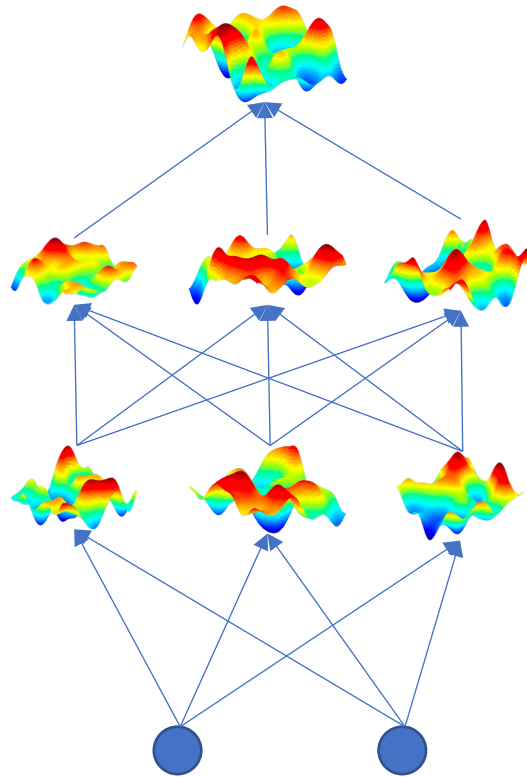
Infinite-width limit

deep kernel
machines

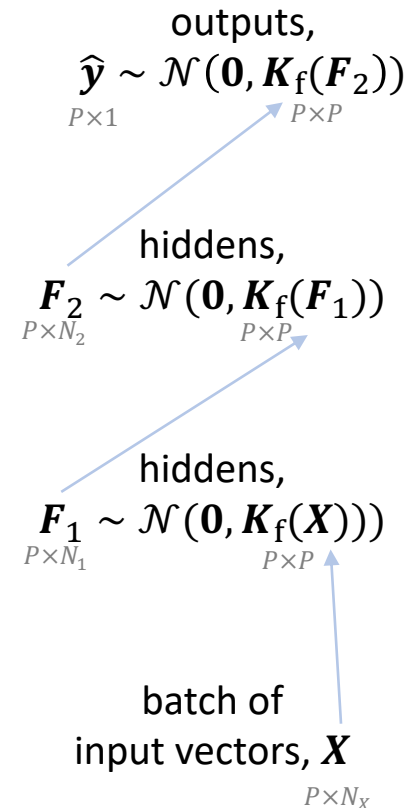
Bayesian neural networks (BNNs) are a special case of deep Gaussian processes (DGPs)

BNNs are a special case of DGPs, with a particular choice of kernel:

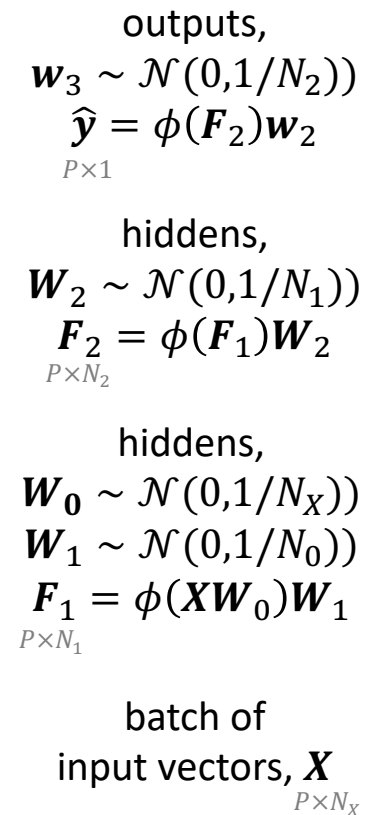
$$\mathbf{K}_f(\mathbf{F}) = \underbrace{\phi(\mathbf{F})}_{P \times N} \underbrace{(\phi(\mathbf{F}))^T}_{N \times P}$$



DGP



BNN



Part 0

Bayesian neural
networks



special
case of

deep Gaussian
processes

rewrite prior in
terms of kernels



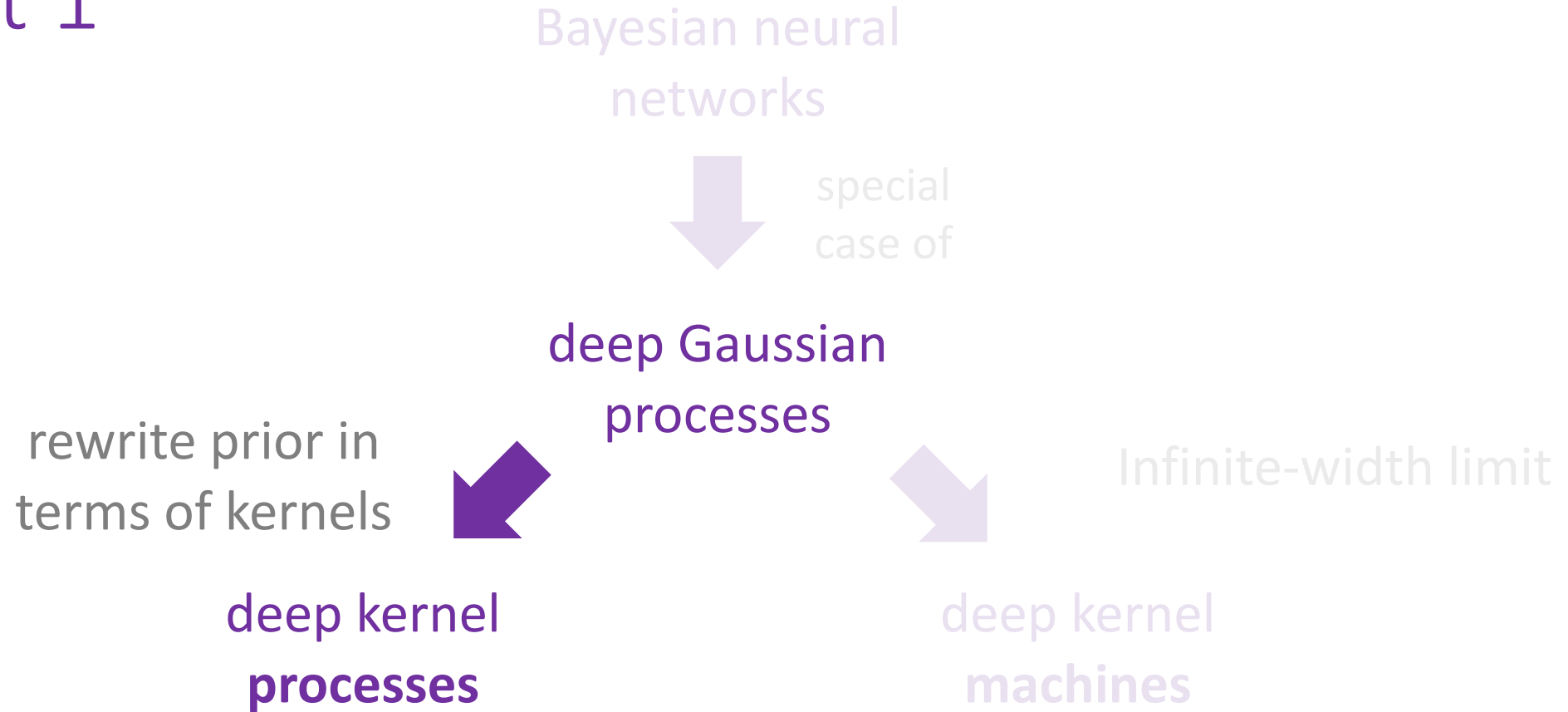
deep kernel
processes



Infinite-width limit

deep kernel
machines

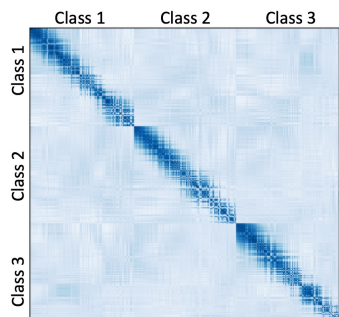
Part 1



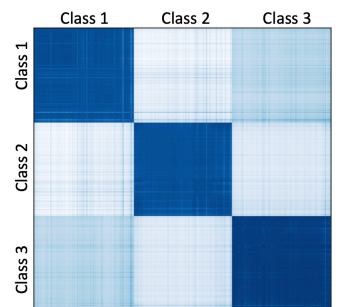
P = number of datapoints
 N_ℓ = width of layer ℓ

Deep kernel processes are “reparametrized” deep Gaussian processes

Gram matrices

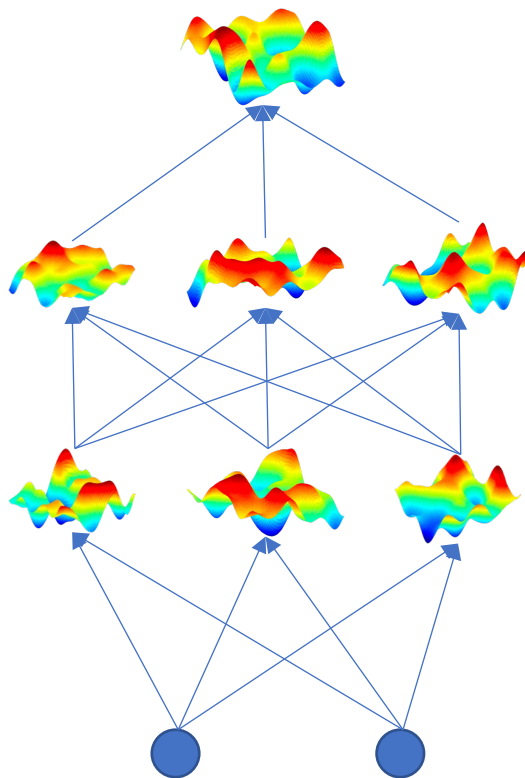


$$\mathbf{G}_2 = \mathbf{F}_2 \mathbf{F}_2^T / N_2$$

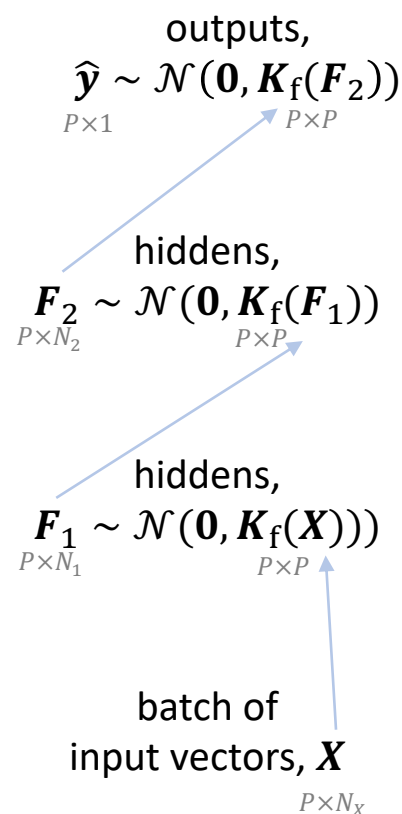


$$\mathbf{G}_1 = \mathbf{F}_1 \mathbf{F}_1^T / N_1$$

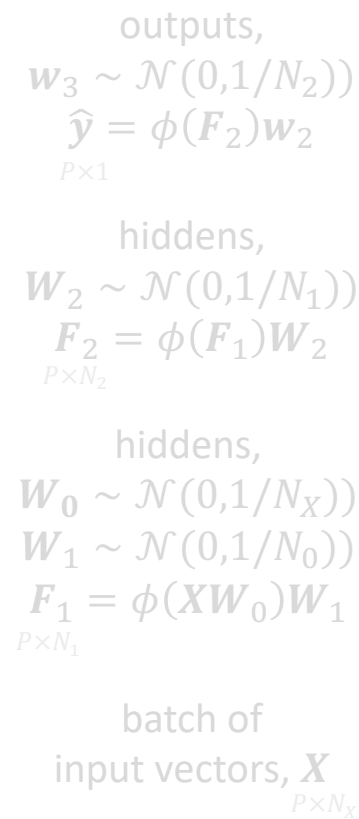
$$\mathbf{G}_0 = \mathbf{X} \mathbf{X}^T / N_X$$



DGP



BNN



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

$$\mathbf{K}_f(\mathbf{F}_\ell) = \mathbf{K}(\mathbf{G}_\ell)$$
$$\mathbf{G}_\ell = \mathbf{F}_\ell \mathbf{F}_\ell^T / N_\ell$$

- Holds for kernels corresponding to infinite-width BNNs!
- Doesn't quite hold for $\mathbf{K}_f(\mathbf{F}) = \phi(\mathbf{F})(\phi(\mathbf{F}))^T$, which corresponds to finite BNNs.
- 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)

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

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

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

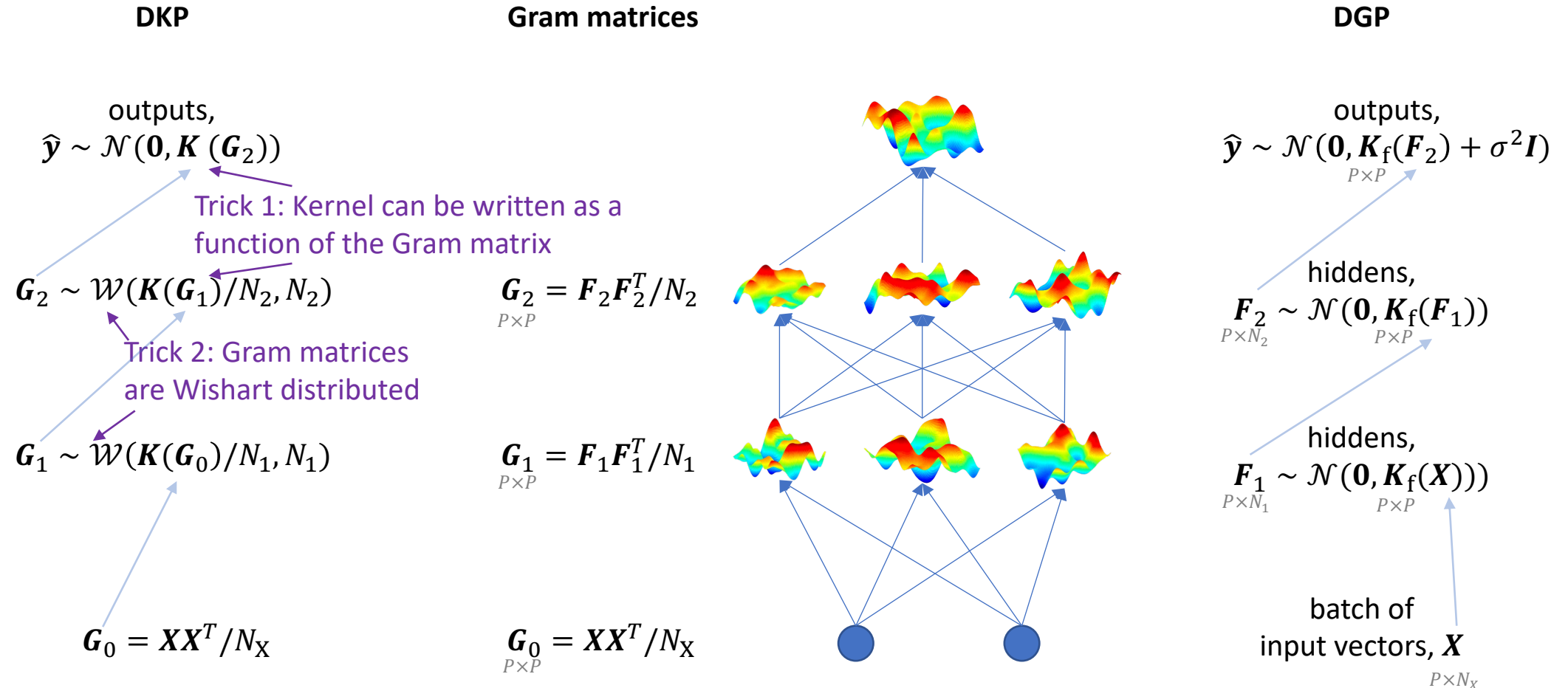
But this exactly matches the definition of the Wishart distribution!

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

(e.g. see Wikipedia for pdf, moments etc.)

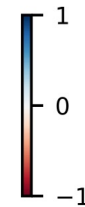
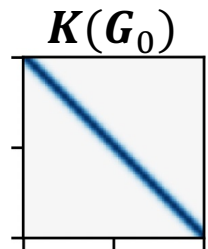
In deep-kernel methods, we switch to working entirely with Gram matrices

P = number of datapoints
 N_ℓ = width of layer ℓ

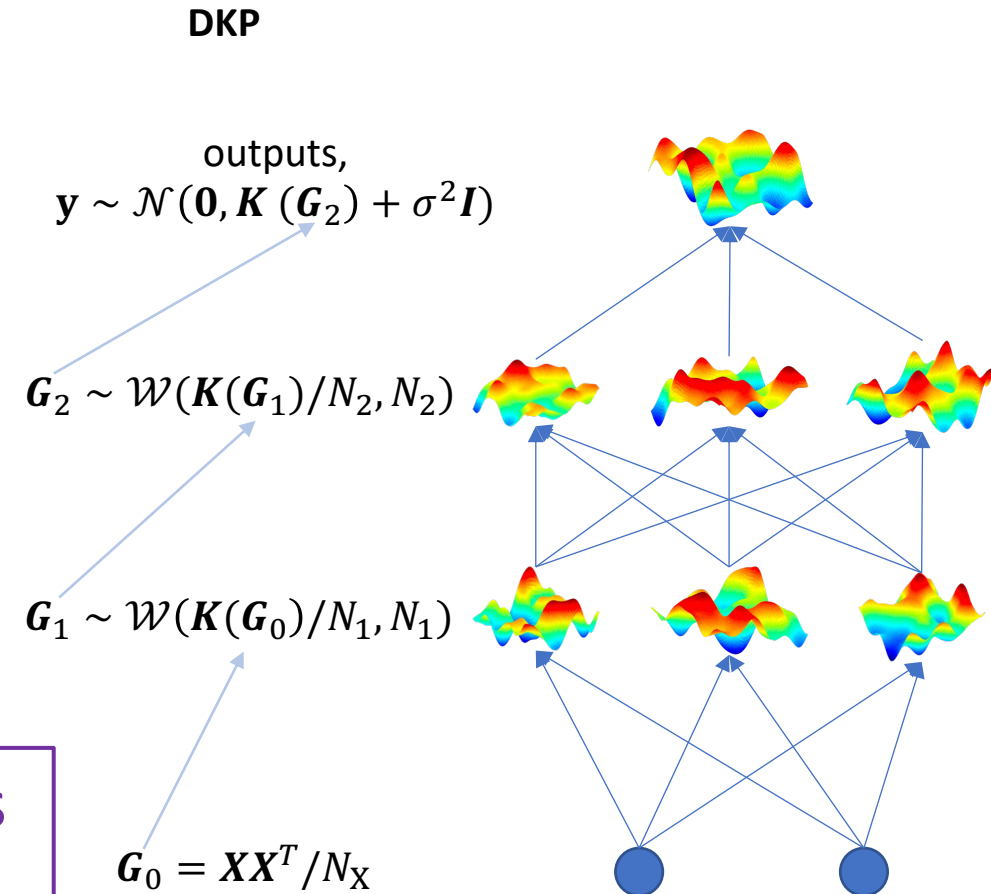


Sampling the prior in the kernelized DGP

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



intuitive, understandable prior over functions
in a deep, nonlinear function approximator!



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})

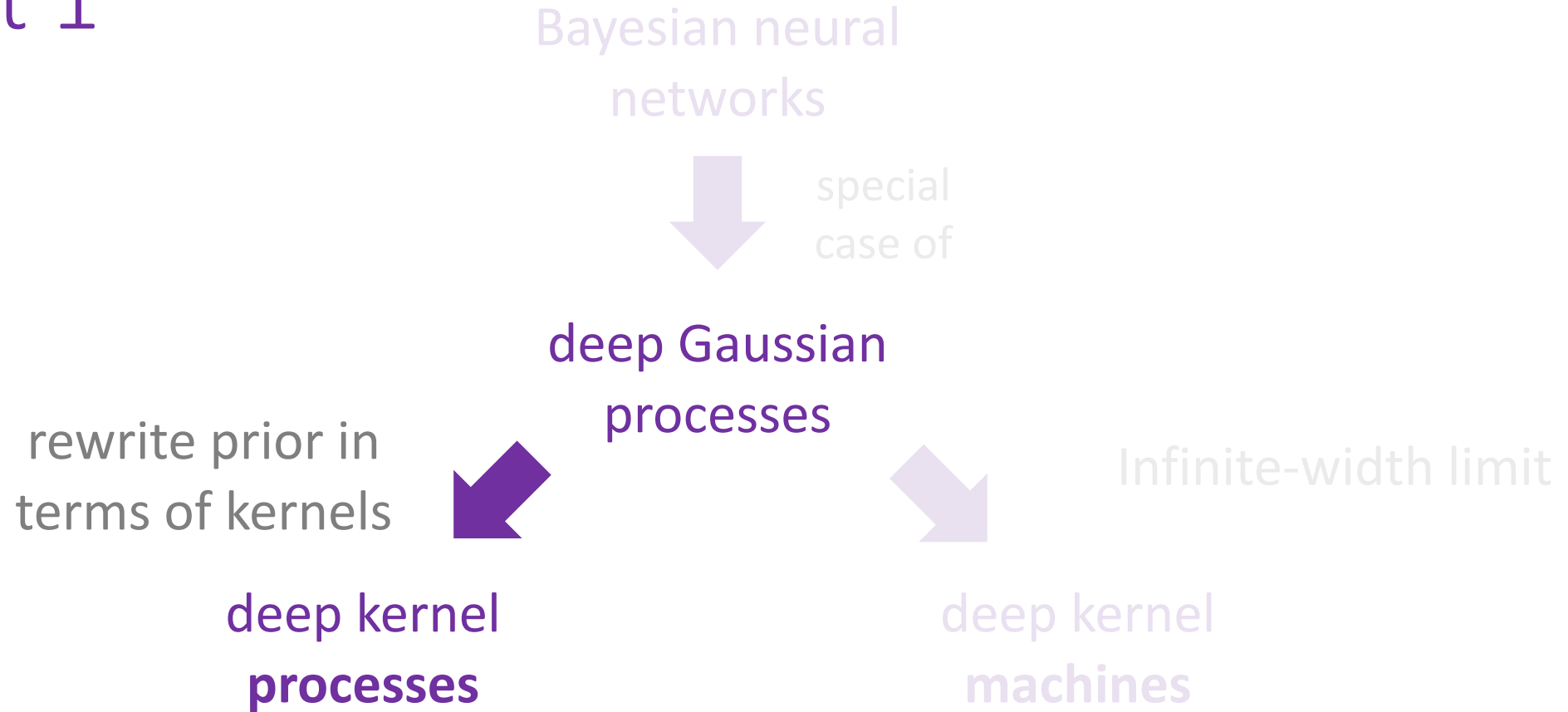
	Dataset	DGP	Q_{GW}	DWP Q_{A-GW}	Q_{AB-GW}
	BOSTON	-2.43 ± 0.04	-2.38 ± 0.04	-2.39 ± 0.05	-2.38 ± 0.04
	CONCRETE	-3.13 ± 0.02	-3.13 ± 0.02	-3.07 ± 0.02	-3.08 ± 0.02
	ENERGY	-0.71 ± 0.03	-0.71 ± 0.03	-0.70 ± 0.03	-0.70 ± 0.03
	KIN8NM	1.38 ± 0.00	1.40 ± 0.01	1.41 ± 0.01	1.41 ± 0.01
LL	NAVAL	8.28 ± 0.04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

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

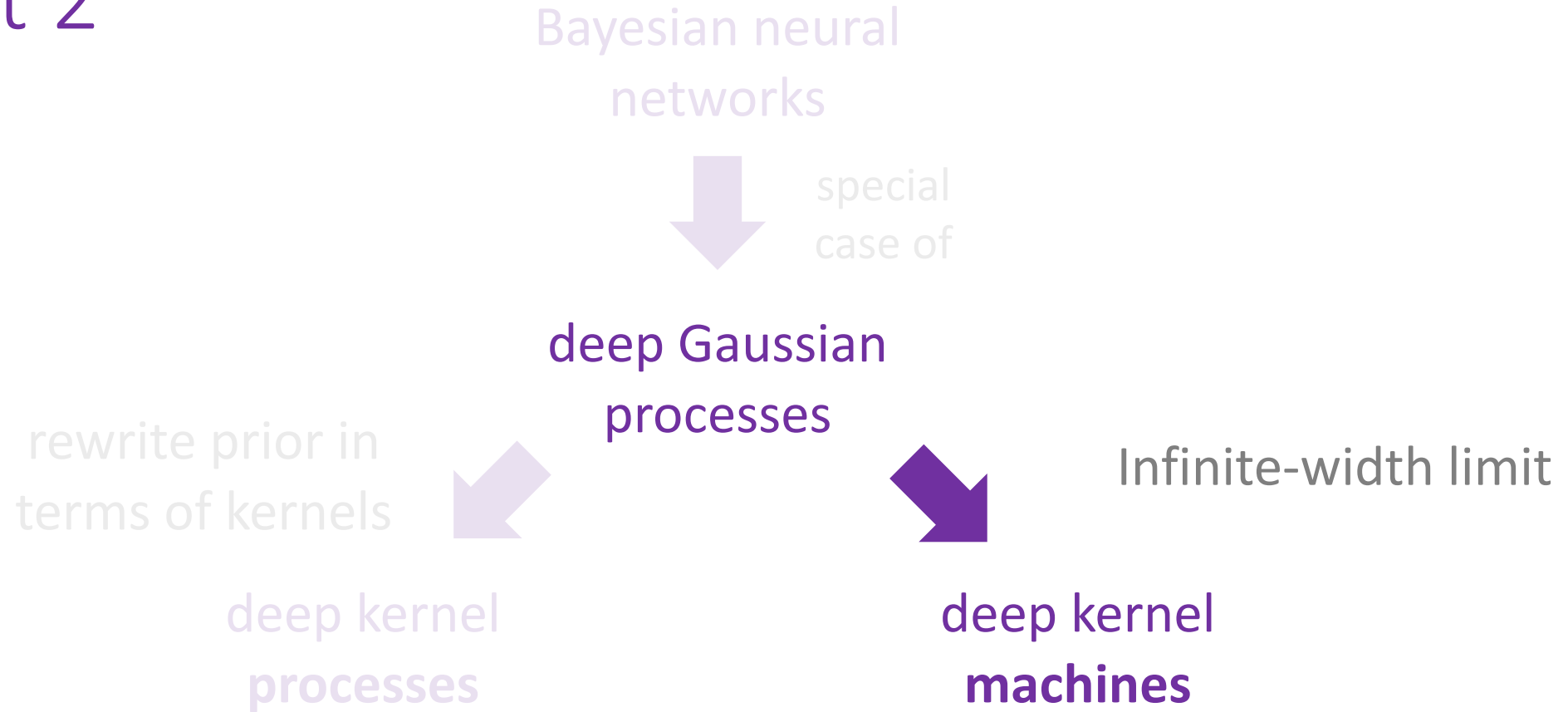
[2] Ober and Aitchison “An approximate posterior for the deep Wishart process” NeurIPS (2021)

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

Part 1



Part 2



Why take an infinite width limit?

- Connect to theory of NNs (which is also infinite-width).
- Get rid of stochasticity (which is a pain in practice).
- Develop effective, practical deep kernel methods.

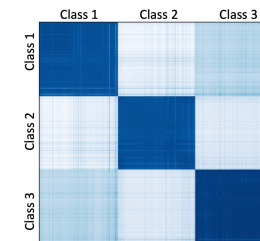
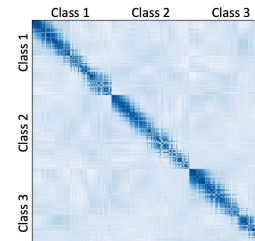
Taking the infinite-width limit of a DKP/DGP...

- VI in an infinite-width DKP/DGP.
- Need to be careful with the limit to make sure we keep representation learning

ELBO:

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

- Optimizes intermediate layer Gram matrices,
- Encourages good performance
- Keeps learned Gram matrix, \mathbf{G}_{ℓ} , similar to NNGP Gram matrix, $\mathbf{K}(\mathbf{G}_{\ell-1})$

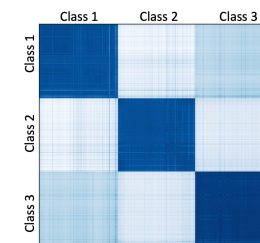
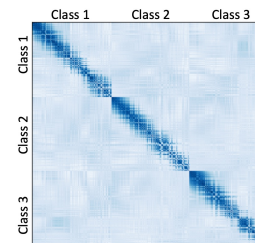


What is a deep kernel machine?

- 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 learned Gram matrix, \mathbf{G}_{ℓ} , similar to NNGP Gram matrix, $\mathbf{K}(\mathbf{G}_{\ell-1})$



Convolutional deep kernel machines are “kernel SOTA”

Table 4: Test accuracy of various kernel methods on CIFAR-10.

Paper	Method	CIFAR-10
This paper	DKM	93.22%
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%

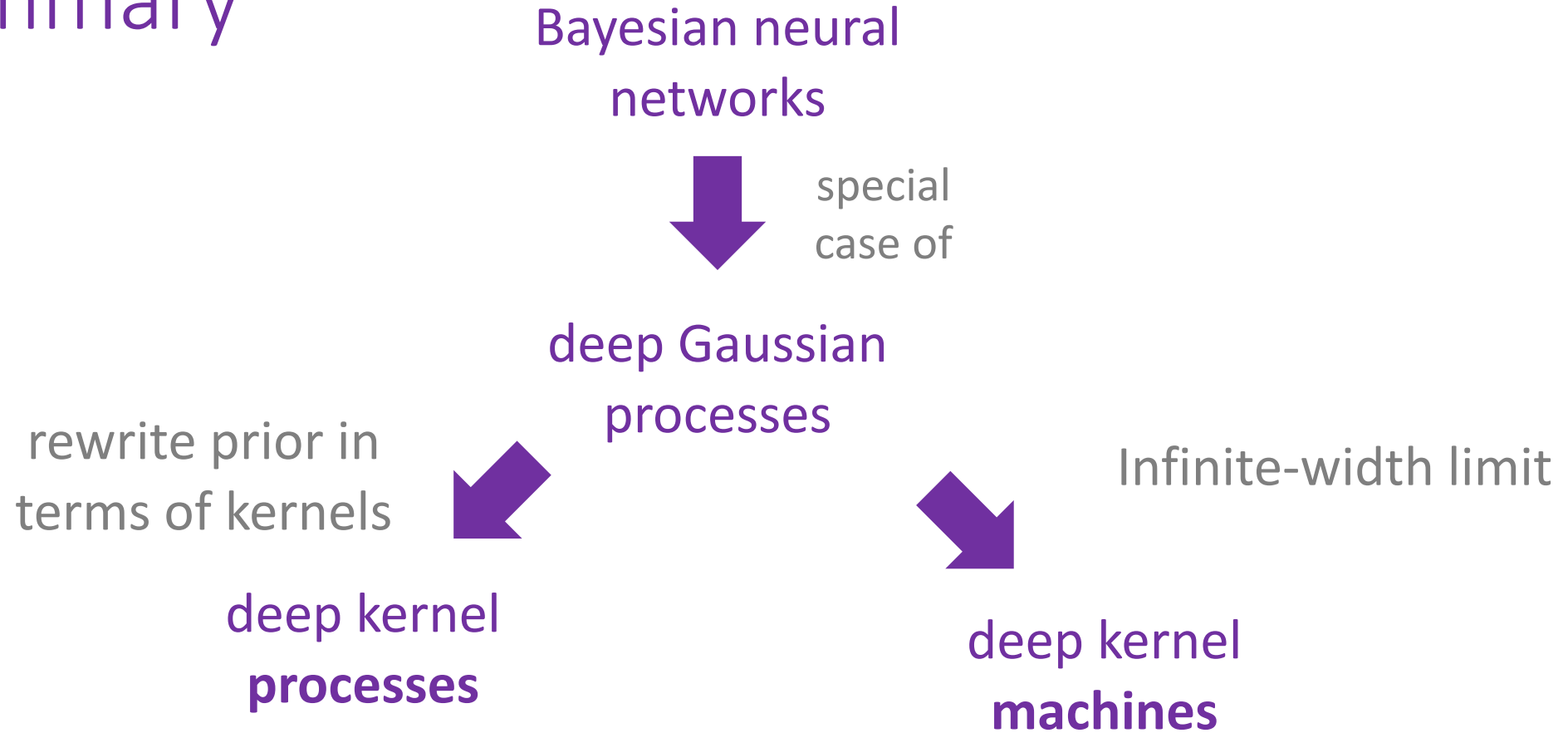


Edward Milsom

But how slow are DKMs? Surprisingly quick!

- We develop a novel inducing-point scheme...
- *...which was roughly same cost as training a standard CNN.*
- So the DKM was orders of magnitude faster than the kernel methods in the table...
- ...but still slow compared to CNNs, (3 days/1200 epochs on 1 GPU for this result)
- But *lots* of structure in the parameter space: we’re working on a natural gradient method that should converge *much* faster.

Summary



Deep kernel landscape + our priorities

Our priorities

- Library: user friendly, "drop in" replacement for NNs.
- Deep kernel transformers
- Speed up (natural gradients)

Huge future opportunities:

	shallow	deep
feature	linear regression	neural net
kernel	kernel ridge regression	deep kernel methods

If you're interested, get in touch:
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" arXiv

Appendix slides

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:

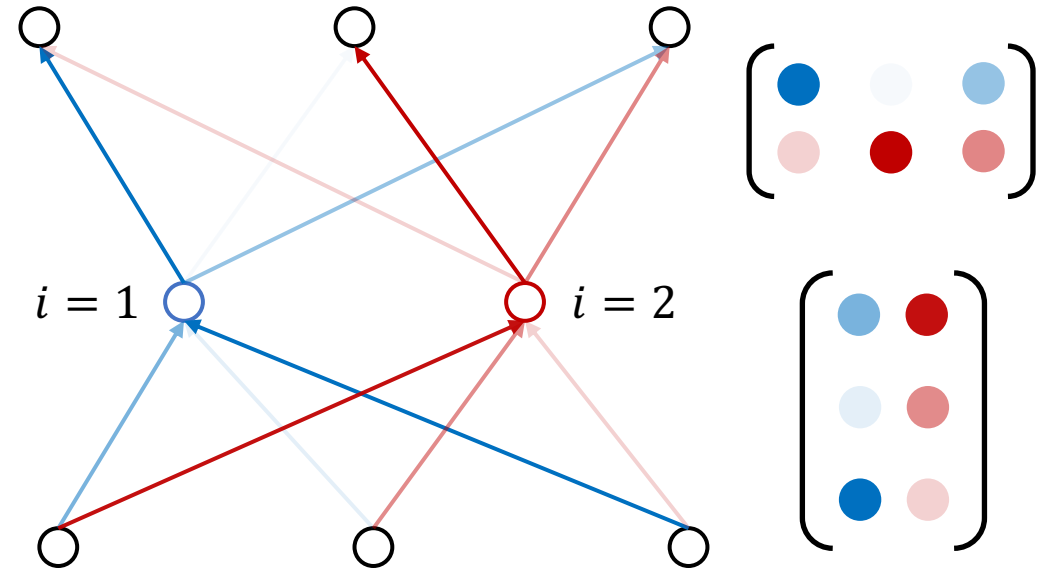
$$Q(\mathbf{F}_1, \dots, \mathbf{F}_L) = \prod_{\ell=1}^L \sum_{\lambda=1}^{N_\ell} \mathcal{N}(\mathbf{f}_\lambda^\ell; \mathbf{0}, \mathbf{G}_\ell)$$

- 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!

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!

