

# Automatic construction and description of nonparametric models

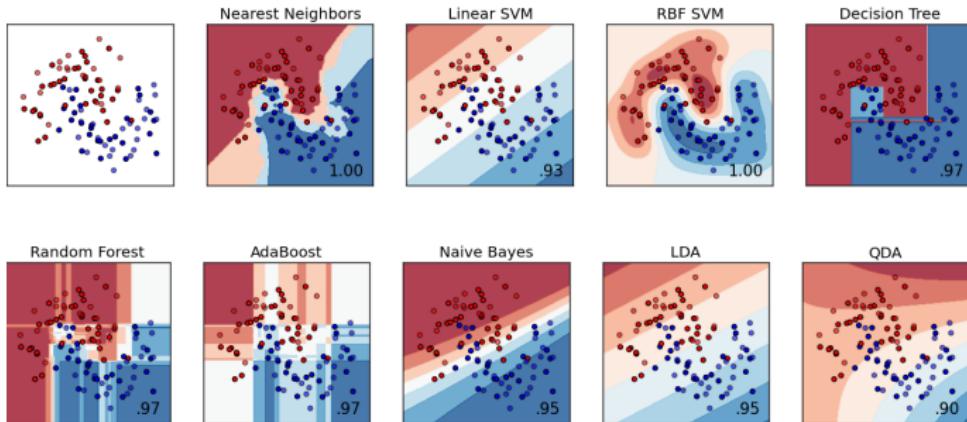


David Duvenaud, James Robert Lloyd, Roger Grosse,  
Joshua Tenenbaum, Zoubin Ghahramani

January 9, 2014

# TYPICAL STATISTICAL MODELLING

- ▶ models typically built by hand, or chosen from a fixed set

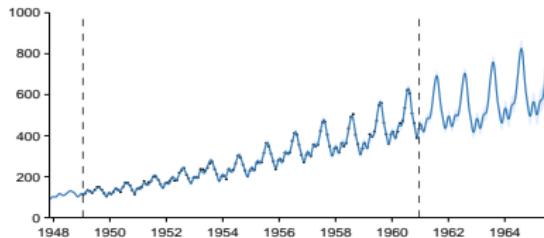


- ▶ Building by hand requires considerable expertise
- ▶ Just being nonparametric isn't good enough
  - ▶ Nonparametric does not mean assumption-free!
- ▶ Can silently fail
  - ▶ If none of the models tried fit the data well, how can you tell?

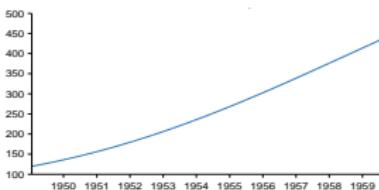
# CAN WE DO BETTER?

- ▶ Andrew Gelman asks: How can an AI do statistics?
- ▶ An artificial statistician would need:
  - ▶ a language that could describe arbitrarily complicated models
  - ▶ a method of searching over those models
  - ▶ a procedure to check model fit
- ▶ This talk: We construct such a language over regression models, a procedure to search over it, and a method to describe in natural language the properties of the resulting models
  - ▶ Working on automatic model-checking...

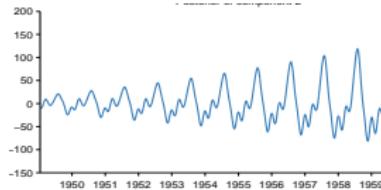
# EXAMPLE: AN AUTOMATIC ANALYSIS



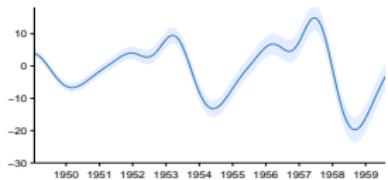
=



+



+



A very smooth, monotonically increasing function

An approximately periodic function with a period of 1.0 years and with approximately linearly increasing amplitude

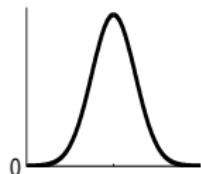
An exactly periodic function with a period of 4.3 years but with linearly increasing amplitude

# A LANGUAGE OF REGRESSION MODELS

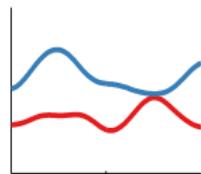
- ▶ We define a language of Gaussian process (GP) regression models by defining a language over kernel functions
- ▶ We start with a small set of base kernels and create a language with a generative grammar
  - ▶  $K \rightarrow K + K$
  - ▶  $K \rightarrow K \times K$
  - ▶  $K \rightarrow CP(K, K)$
  - ▶  $K \rightarrow \{\text{SE}, \text{LIN}, \text{PER}\}$
- ▶ The language is open-ended, but its structure makes natural-language description simple

# KERNELS DETERMINE STRUCTURE OF GPs

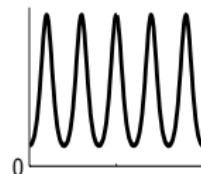
- ▶ Kernel determines almost all the properties of a GP prior
- ▶ Many different kinds, with very different properties:



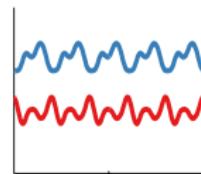
Squared-exp  
(SE)



local  
variation



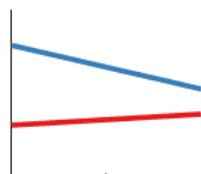
Periodic  
(PER)



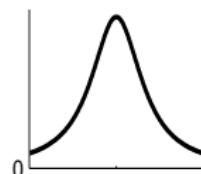
repeating  
structure



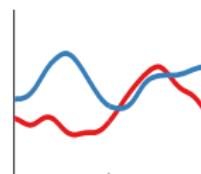
Linear (LIN)



linear  
functions



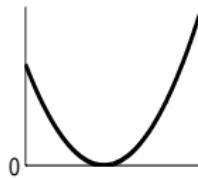
Rational-  
quadratic(RQ)



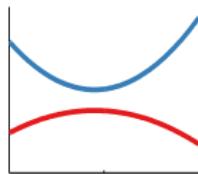
multi-scale  
variation

# KERNELS CAN BE COMPOSED

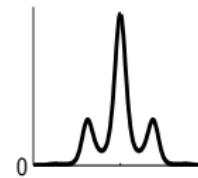
- ▶ Two main operations: addition, multiplication



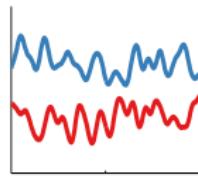
LIN × LIN



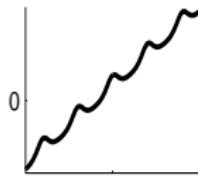
quadratic  
functions



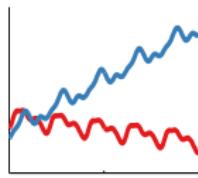
SE × PER



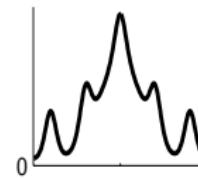
locally  
periodic



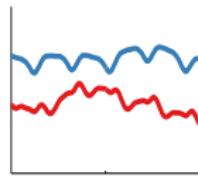
LIN + PER



periodic  
with trend



SE + PER

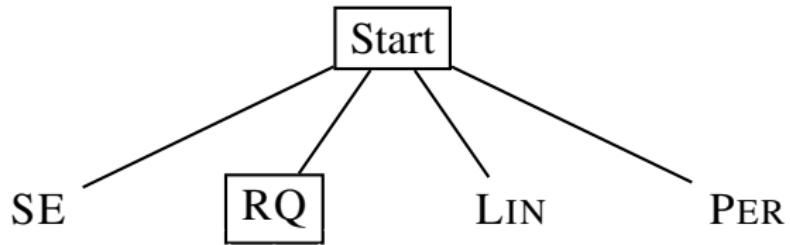
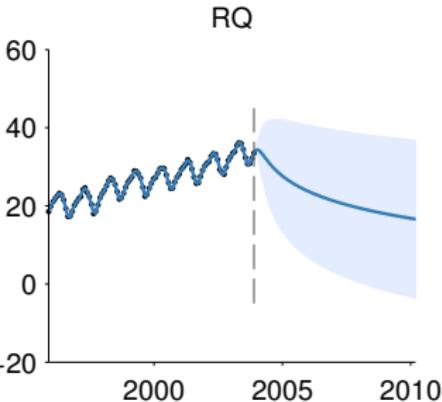


periodic  
with noise

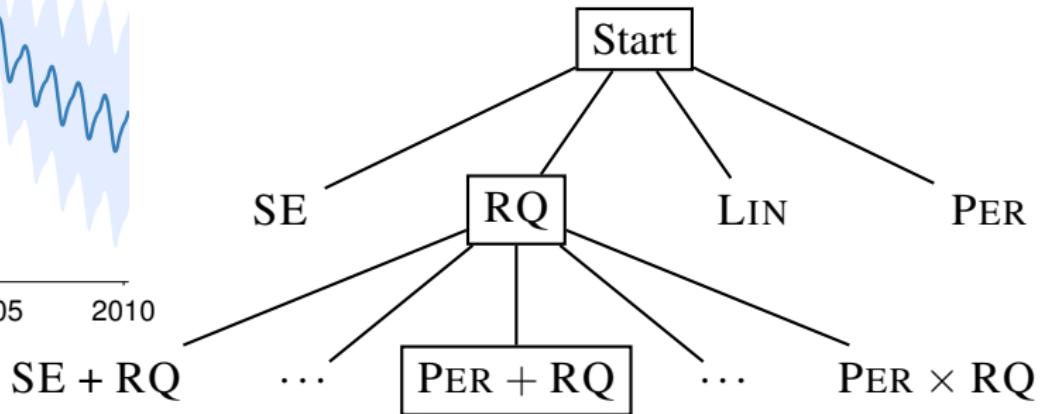
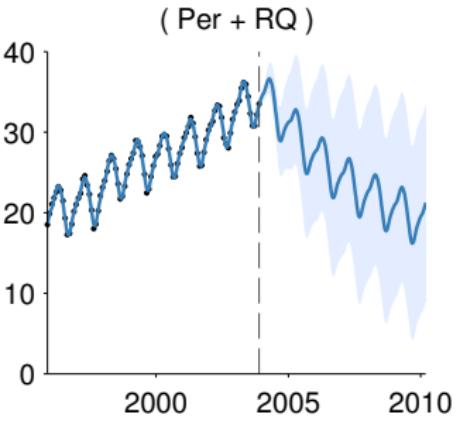
# SPECIAL CASES IN OUR LANGUAGE

| Regression motif     | Example kernel                 |
|----------------------|--------------------------------|
| Linear regression    | LIN                            |
| Quadratic regression | $\text{LIN} \times \text{LIN}$ |
| Fourier analysis     | $\sum \cos$                    |
| Spectral kernels     | $\sum \text{SE} \times \cos$   |
| Changepoints         | e.g. CP(PER, SE)               |
| Heteroscedasticity   | e.g. SE + LIN $\times$ SE      |

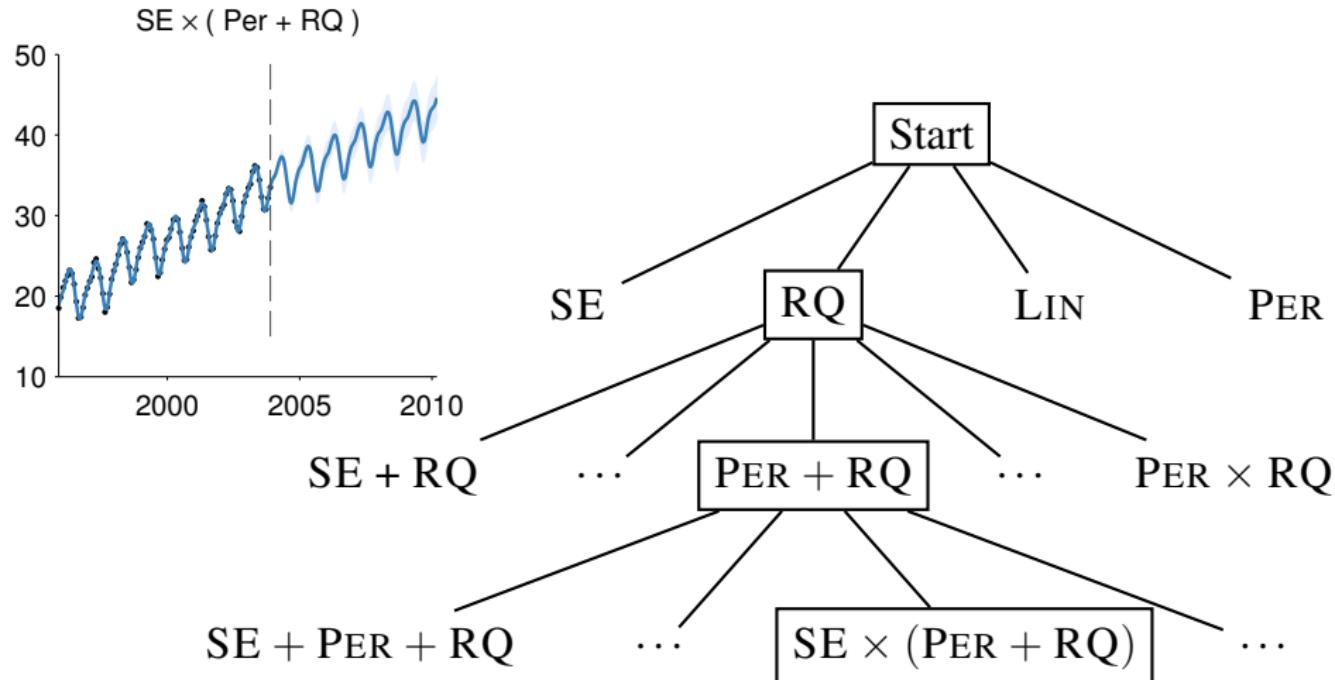
# COMPOSITIONAL STRUCTURE SEARCH



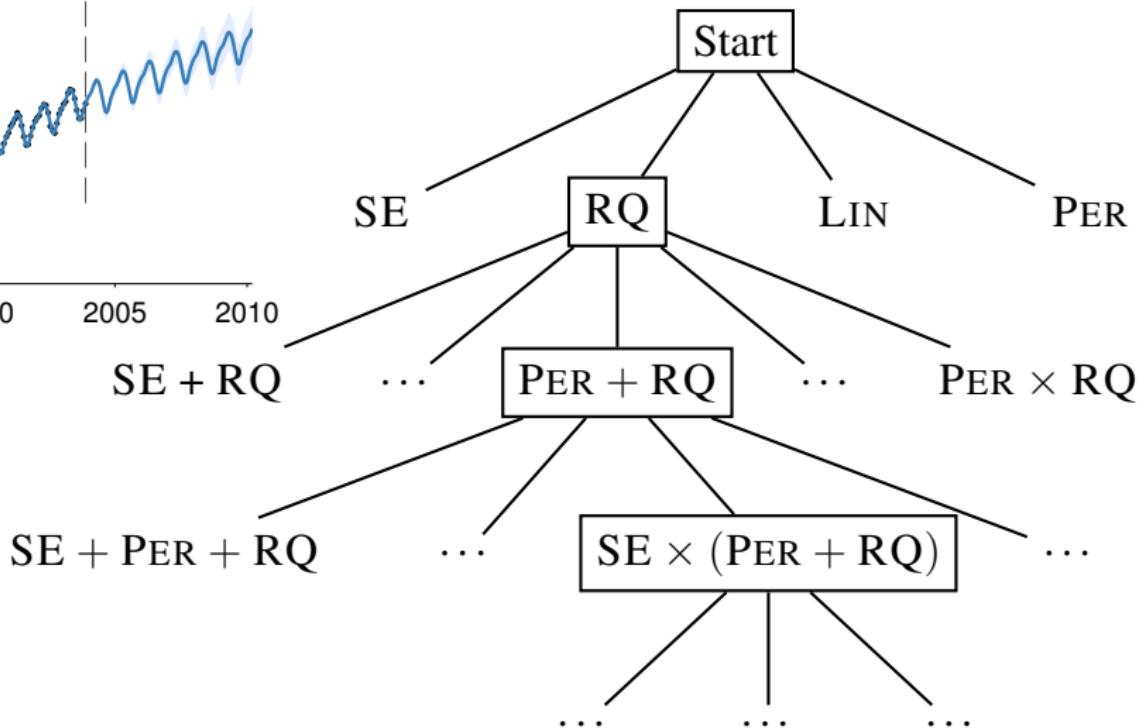
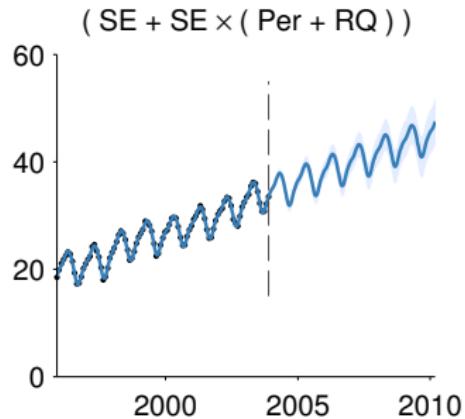
# COMPOSITIONAL STRUCTURE SEARCH



# COMPOSITIONAL STRUCTURE SEARCH



# COMPOSITIONAL STRUCTURE SEARCH



# DISTRIBUTIVITY HELPS INTERPRETABILITY

We can write all kernels as sums of products of base kernels:

$$\text{SE} \times (\text{RQ} + \text{LIN}) = (\text{SE} \times \text{RQ}) + (\text{SE} \times \text{LIN}).$$

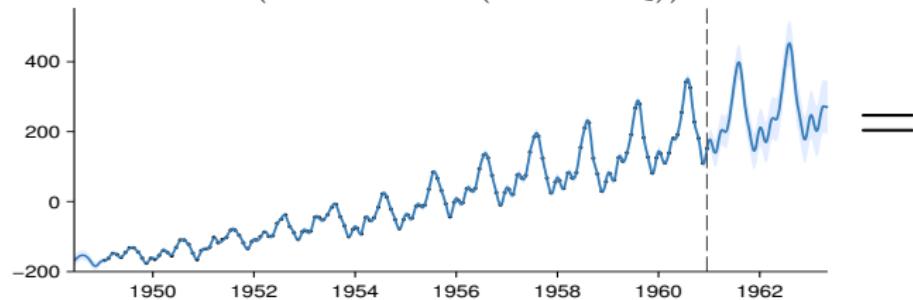
Sums of kernels are equivalent to sums of functions.

If  $f_1, f_2$  are independent, and  $f_1 \sim \mathcal{GP}(\mu_1, k_1), f_2 \sim \mathcal{GP}(\mu_2, k_2)$  then

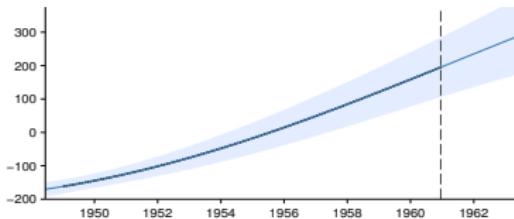
$$(f_1 + f_2) \sim \mathcal{GP}(\mu_1 + \mu_2, k_1 + k_2)$$

# EXAMPLE DECOMPOSITION: AIRLINE

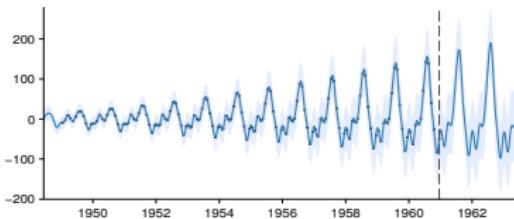
$$SE \times (LIN + LIN \times (PER + RQ))$$



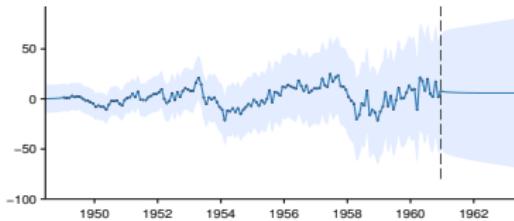
$$SE \times LIN$$



$$SE \times LIN \times PER$$

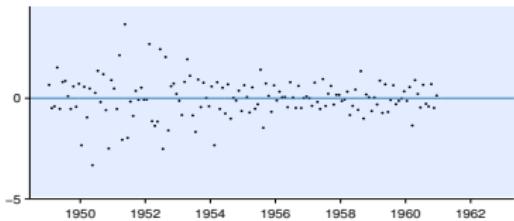


$$SE \times LIN \times RQ$$



+

$$\text{Residuals}$$



# DESCRIBING KERNELS

Products of same type of kernel collapse.

| Product of Kernels                                     | Becomes   |
|--|---|
| $\text{SE} \times \text{SE} \times \text{SE} \dots$    | SE  |
| $\text{LIN} \times \text{LIN} \times \text{LIN} \dots$ | A polynomial  |
| $\text{PER} \times \text{PER} \times \text{PER}$       | Same covariance as<br>product of periodic functions |

# DESCRIBING KERNELS

Each kernel acts as a modifier in a standard way: an “adjective”.

| Kernel                | Becomes                           |
|-----------------------|-----------------------------------|
| $K \times \text{SE}$  | 'locally' or 'approximately'      |
| $K \times \text{LIN}$ | 'with linearly growing amplitude' |
| $K \times \text{PER}$ | 'periodic'                        |
| $\text{CP}(K1, K2)$   | '... changing at $x$ to ...'      |

- ▶ Special cases for when they're on their own
- ▶ Extra adjectives depending on hyperparameters

# EXAMPLE KERNEL DESCRIPTIONS

| Product of Kernels           | Description  |
|------------------------------|--|
| PER                          | An exactly periodic function                                       |
| PER $\times$ SE              | An approximately periodic function                                 |
| PER $\times$ SE $\times$ LIN | An approximately periodic function with linearly varying amplitude |

# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

The raw data and full model posterior with extrapolations are shown in figure 1.

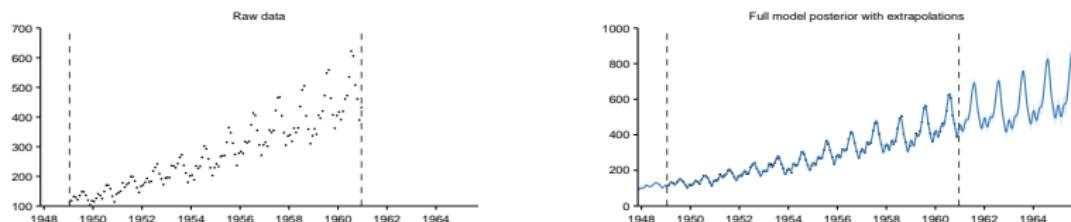


Figure 1: Raw data (left) and model posterior with extrapolation (right)

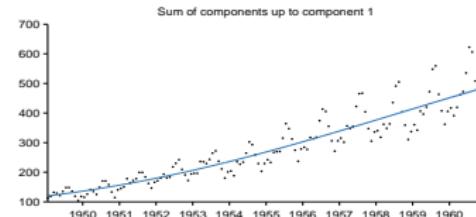
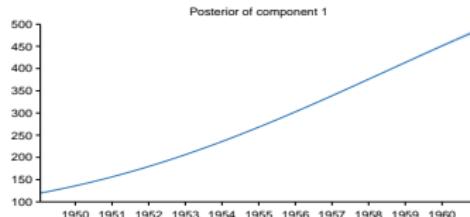
The structure search algorithm has identified four additive components in the data:

- A very smooth monotonically increasing function.
- An approximately periodic function with a period of 1.0 years and with approximately linearly increasing amplitude.
- An exactly periodic function with a period of 4.3 years but with linearly increasing amplitude.
- Uncorrelated noise with linearly increasing standard deviation.

# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.1 Component 1

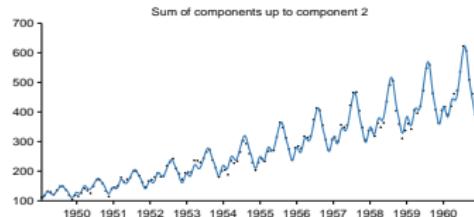
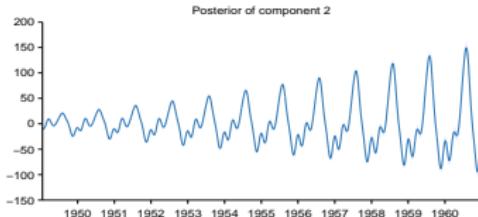
This component is a very smooth and monotonically increasing function.



# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.2 Component 2

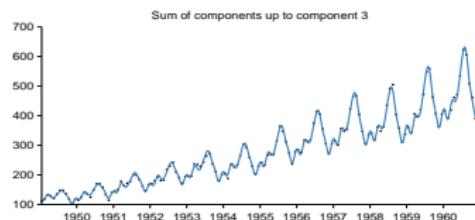
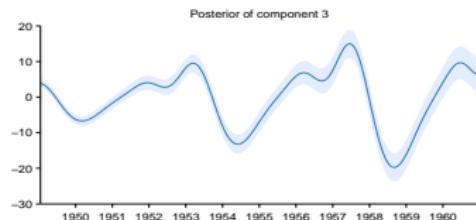
This component is approximately periodic with a period of 1.0 years and varying amplitude. Across periods the shape of this function varies very smoothly. The amplitude of the function increases approximately linearly. The shape of this function within each period has a typical lengthscale of 6.0 weeks.



# THIS ANALYSIS WAS AUTOMATICALLY GENERATED

## 2.3 Component 3

This component is exactly periodic with a period of 4.3 years but with varying amplitude. The amplitude of the function increases linearly. The shape of this function within each period has a typical lengthscale of 7.4 months.



# SUMMARY

---

- ▶ Constructed a language of regression models through kernel composition
- ▶ Searched over this language greedily
- ▶ Kernels sums and products modify prior in predictable ways, allowing automatic natural-language description of models
- ▶ Open questions:
  - ▶ Interpretability versus flexibility
  - ▶ Automatic Model-checking

# KERNEL LEARNING AS FEATURE LEARNING

- ▶ Kernels implicitly compute dot product between features of two items:  $k(x, x') = \phi(x)^\top \phi(x')$ .
- ▶ one-to-one mapping between kernels and feature maps
- ▶ feature learning with tractable marginal likelihood!
- ▶ Can compose feature maps:
- ▶ example: periodic kernel  $k_{per}(x, x') = \exp(-\sin^2(x - x'))$  is equivalent to  $k_{se}(\sin(x), \cos(x), \sin(x'), \cos(x'))$ .
- ▶ What can we do with feature compositions?

# DEEP KERNELS

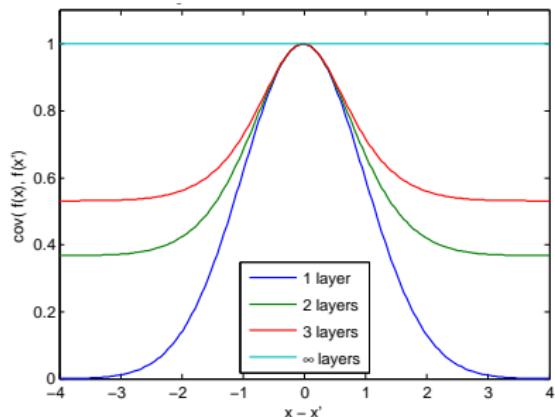
- ▶ (Cho, 2012) built kernels from multiple layers of feature mappings.
- ▶ given  $k_1(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^\top \Phi(\mathbf{x}')$ , we can consider  $k_2(\mathbf{x}, \mathbf{x}') = \Phi(\Phi(\mathbf{x}))^\top \Phi(\Phi(\mathbf{x}'))$ . For SE kernel:

$$\begin{aligned} k_2(\mathbf{x}, \mathbf{x}') &= \\ &= \exp\left(-\frac{1}{2}\|\Phi(\mathbf{x}) - \Phi(\mathbf{x}')\|_2^2\right) \\ &= \exp\left(-\frac{1}{2}[k_1(\mathbf{x}, \mathbf{x}) - 2k_1(\mathbf{x}, \mathbf{x}') + k_1(\mathbf{x}', \mathbf{x}')]\right) \\ &= \exp(k_1(\mathbf{x}, \mathbf{x}') - 1) \quad (\text{if } k_1(\mathbf{x}, \mathbf{x}) = 1) \end{aligned}$$

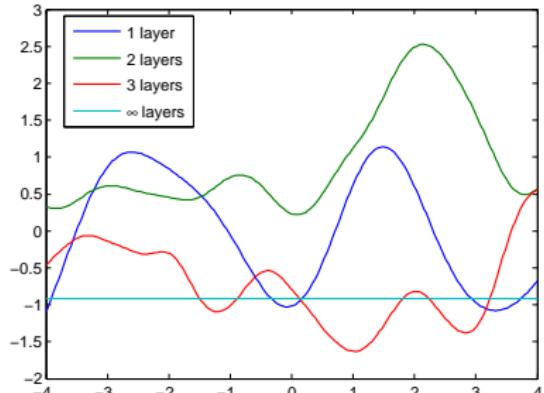
- ▶ Why not go deeper?

# INFINITELY DEEP KERNELS

- ▶ For SE kernel,  $k_{n+1}(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}') - 1)$ .
- ▶ What is the eventual limit?



Kernel



Draws from GP prior

- ▶  $k_\infty(\mathbf{x}, \mathbf{x}') = 1$  everywhere.

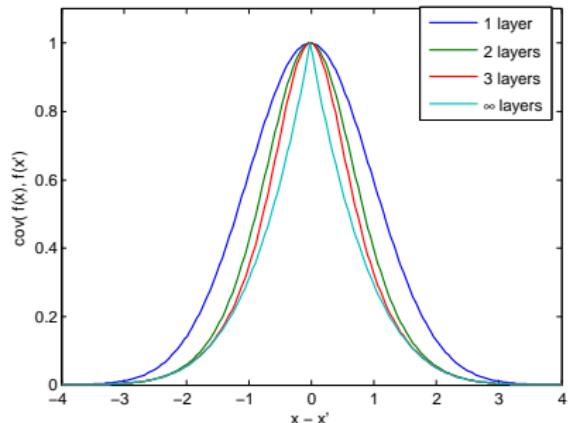
## A SIMPLE FIX...

- ▶ Following a suggestion from (Neal, 1995), we connect the inputs  $\mathbf{x}$  to each layer. (append  $\mathbf{x}$  to features at each layer:

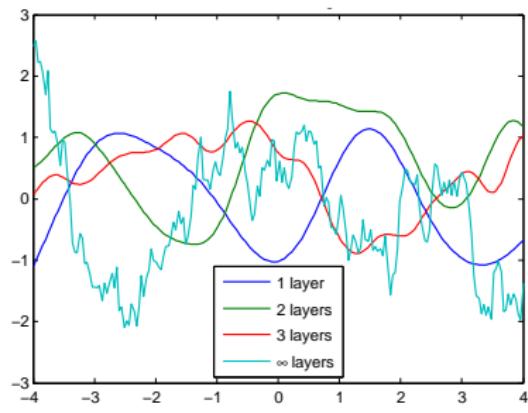
$$\begin{aligned} k_{n+1}(\mathbf{x}, \mathbf{x}') &= \\ &= \exp \left( -\frac{1}{2} \left\| \begin{bmatrix} \Phi_n(\mathbf{x}) \\ \mathbf{x} \end{bmatrix} - \begin{bmatrix} \Phi_n(\mathbf{x}') \\ \mathbf{x}' \end{bmatrix} \right\|_2^2 \right) \\ &= \exp \left( k_n(\mathbf{x}, \mathbf{x}') - 1 - \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \right) \quad (1) \end{aligned}$$

# INFINITELY DEEP KERNELS

- $k_{n+1}(\mathbf{x}, \mathbf{x}') = \exp(k_n(\mathbf{x}, \mathbf{x}') - 1 - \frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|_2^2)$ .
- What is the eventual limit?



Kernel



Draws from GP prior

- Looks like an OU process with skinny tails, samples are non-differentiable (fractal).

# WHAT WENT WRONG?

---

- ▶ Fixed feature mappings
- ▶ not capturing invariances
- ▶ not throwing away unnecessary information
- ▶ an illustration of why learning is usually necessary

# DEEP GAUSSIAN PROCESSES

- ▶ a prior over compositions of functions:

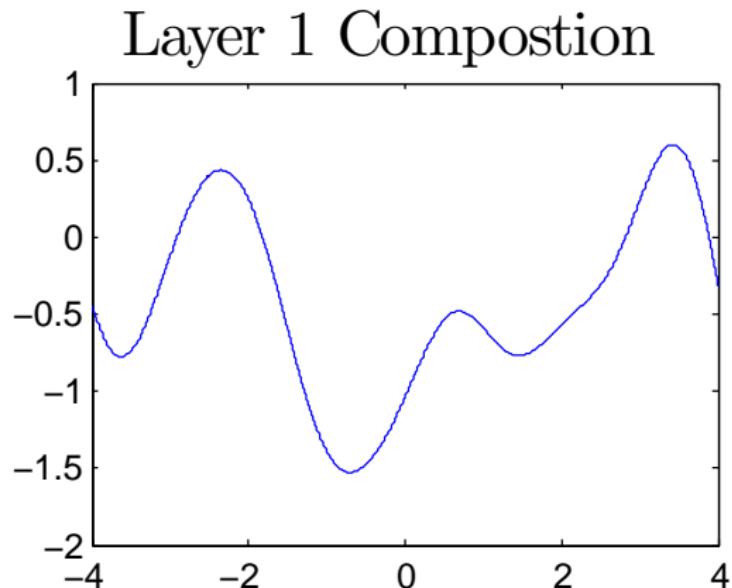
$$\mathbf{f}^{(1:L)}(\mathbf{x}) = \mathbf{f}^{(L)}\left(\mathbf{f}^{(L-1)}\left(\dots \mathbf{f}^{(2)}\left(\mathbf{f}^{(1)}(\mathbf{x})\right)\dots\right)\right) \quad (2)$$

where each  $\mathbf{f}_d^{(\ell)} \stackrel{\text{ind}}{\sim} \mathcal{GP}\left(0, k_d^\ell(\mathbf{x}, \mathbf{x}')\right)$ .

- ▶ Can be seen as a certain limit of infinitely-wide deep nets.
- ▶ inference is really hard.
- ▶ maybe we can learn something about deep models just from looking at prior draws?

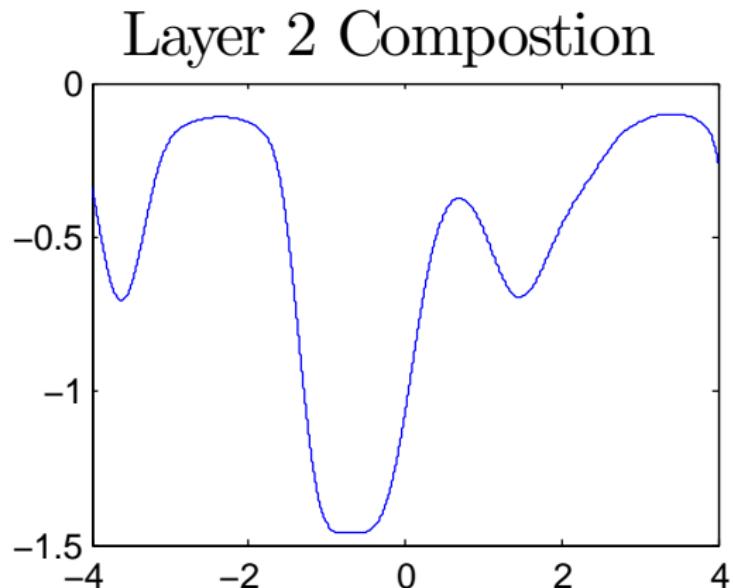
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



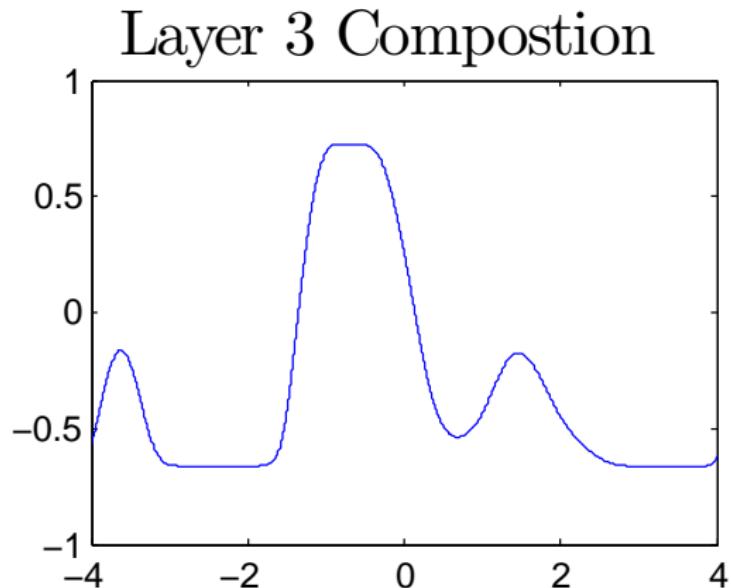
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



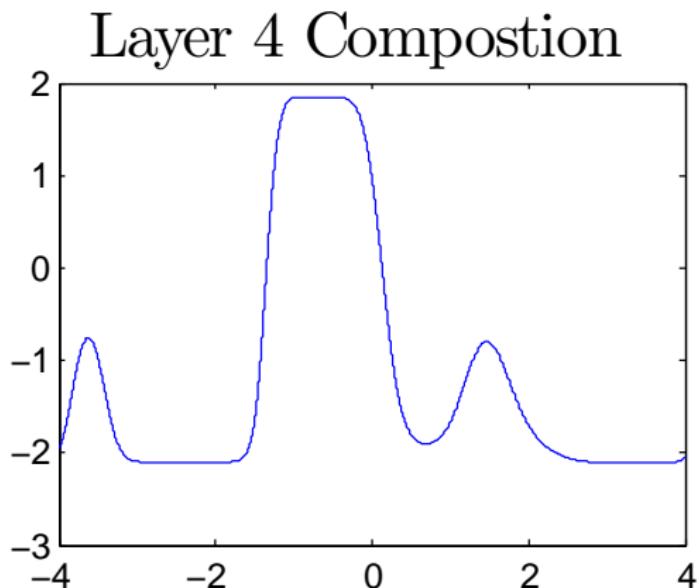
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



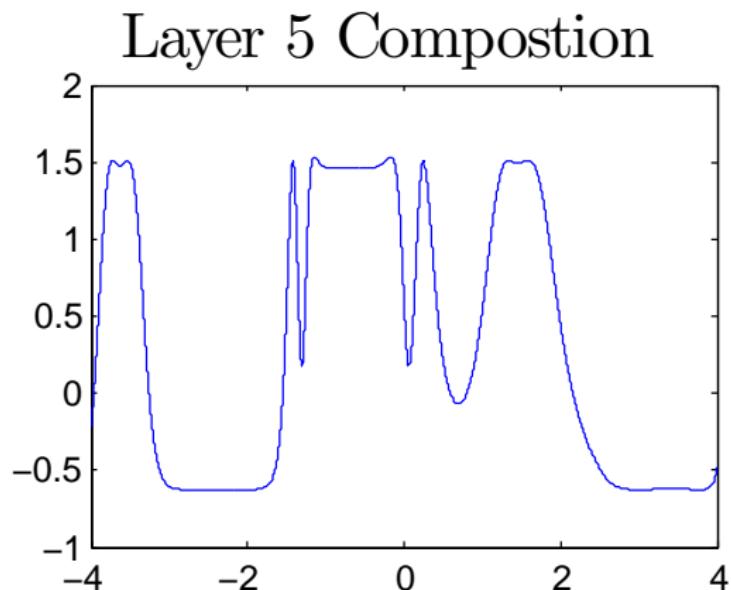
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



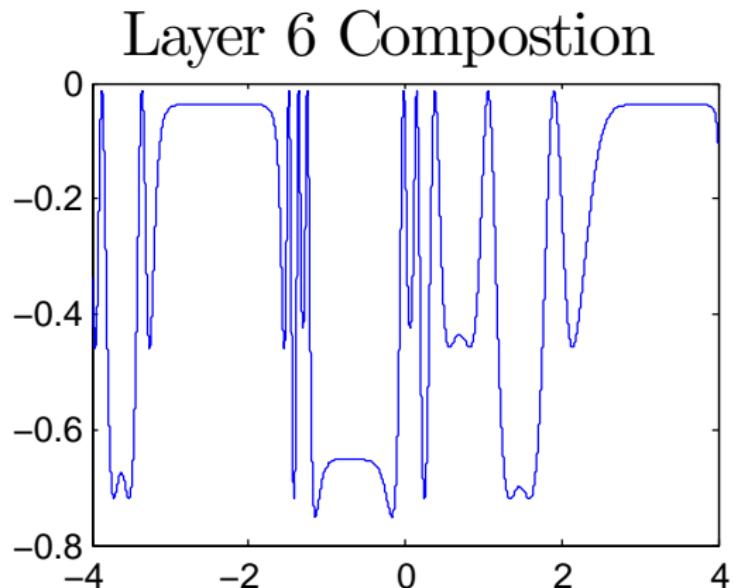
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



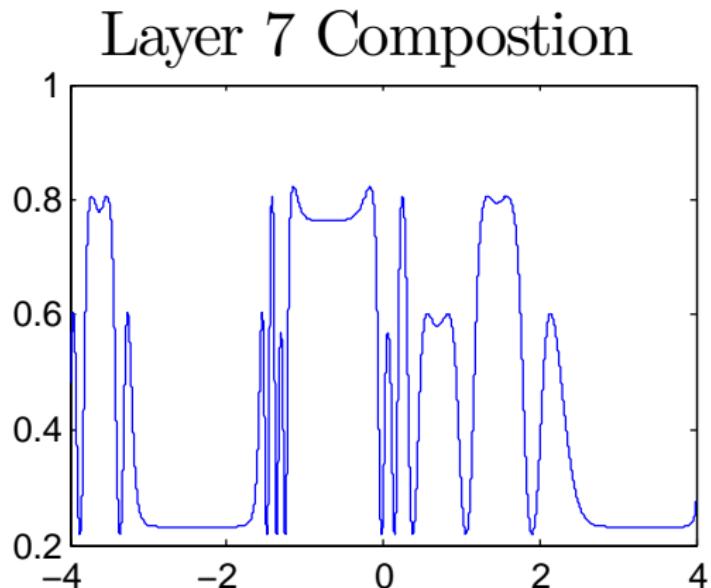
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



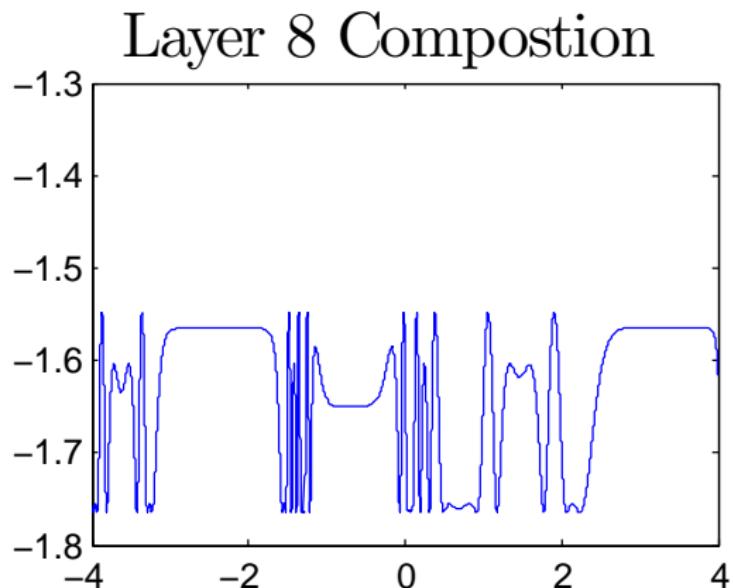
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



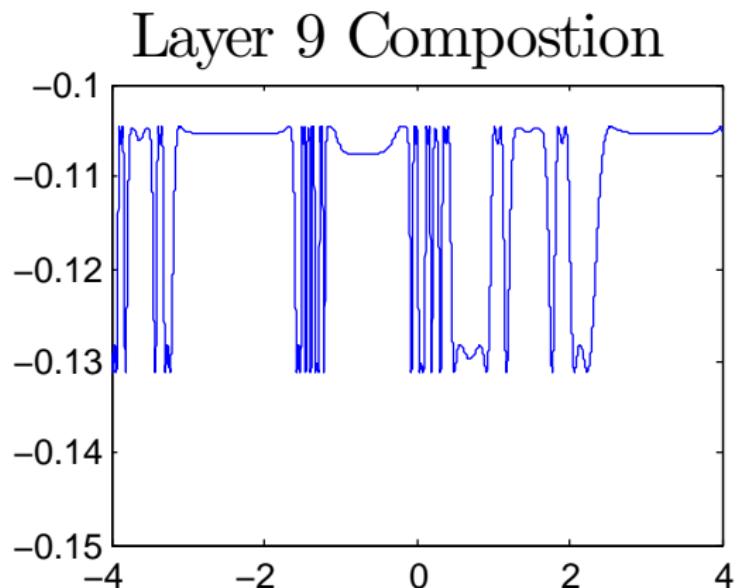
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



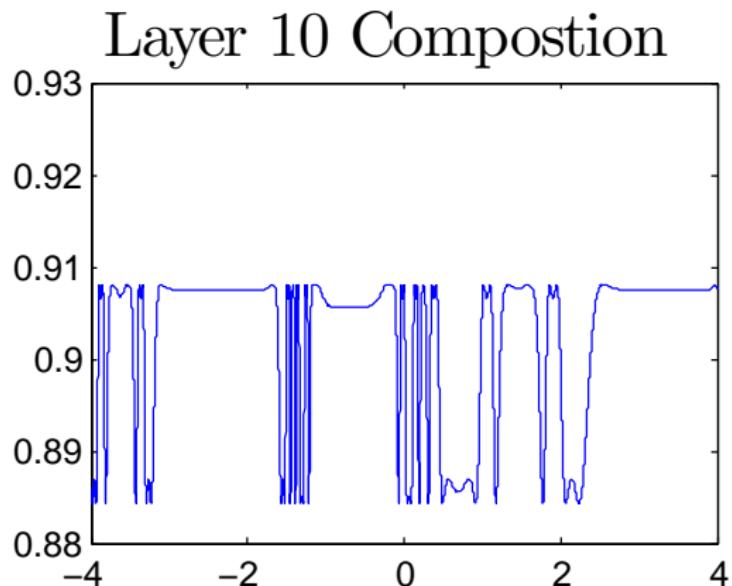
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



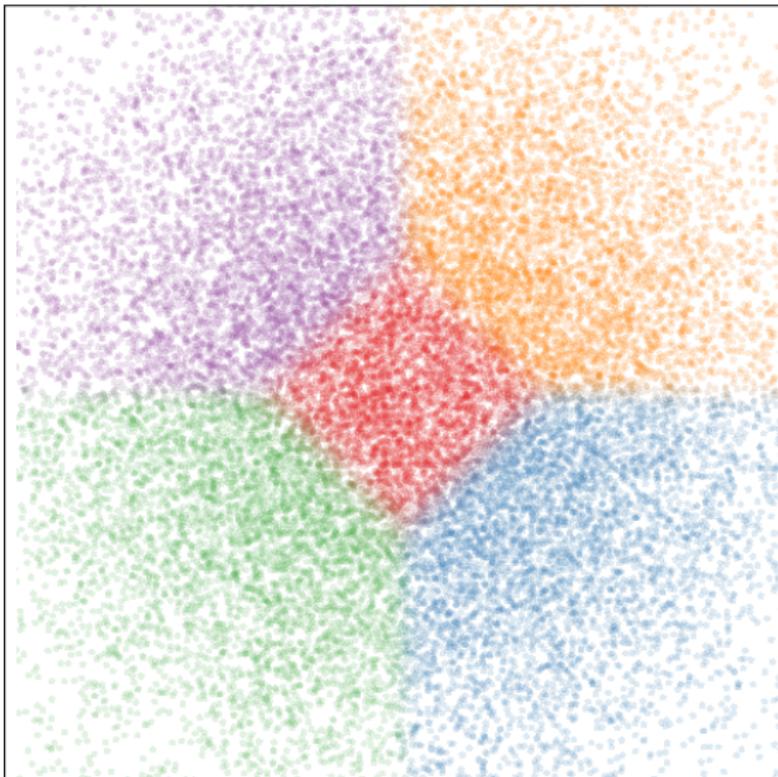
# DEEP GAUSSIAN PROCESSES

- ▶ Draws from one-dimensional deep GPs:



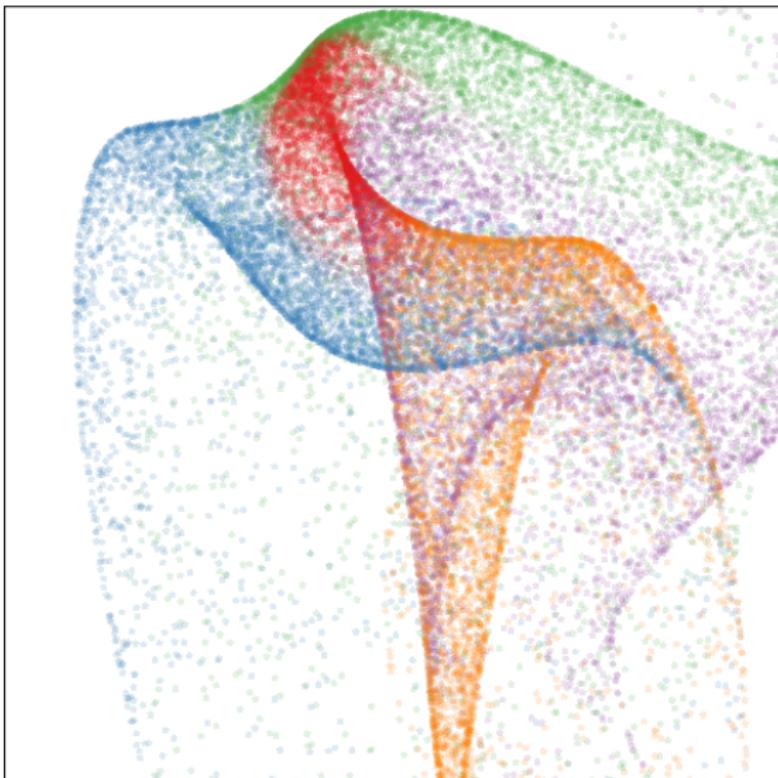
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



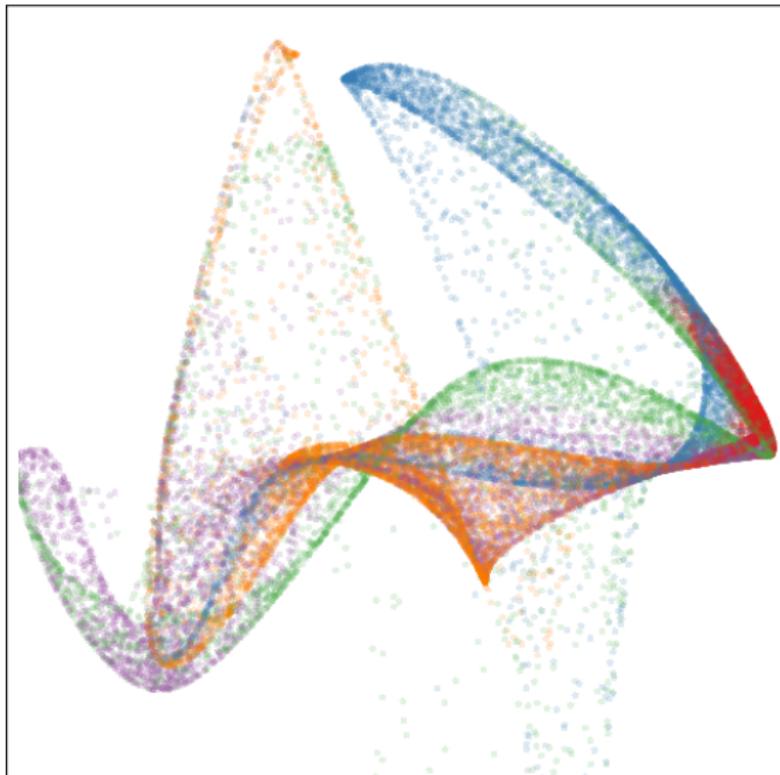
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



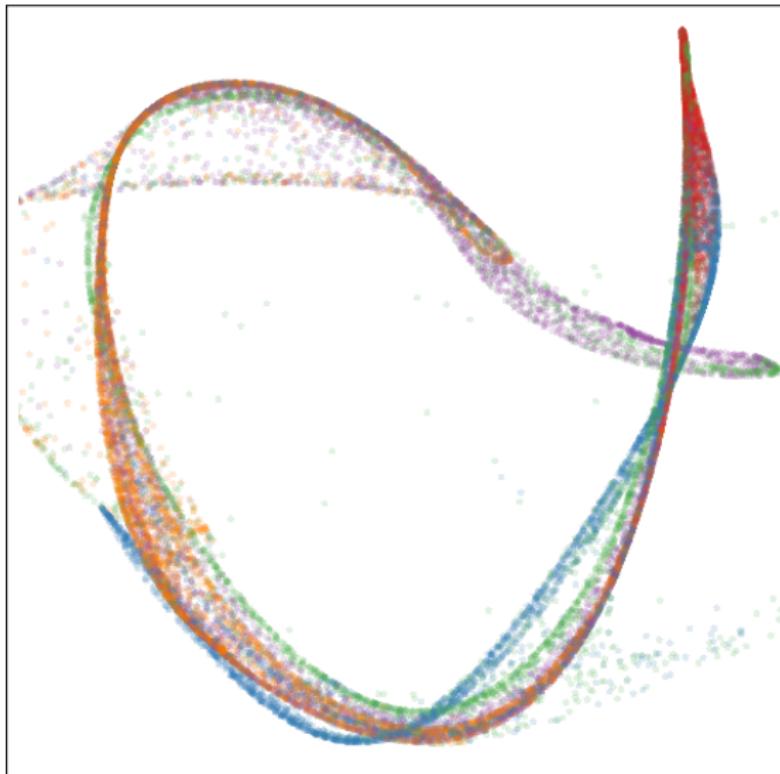
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



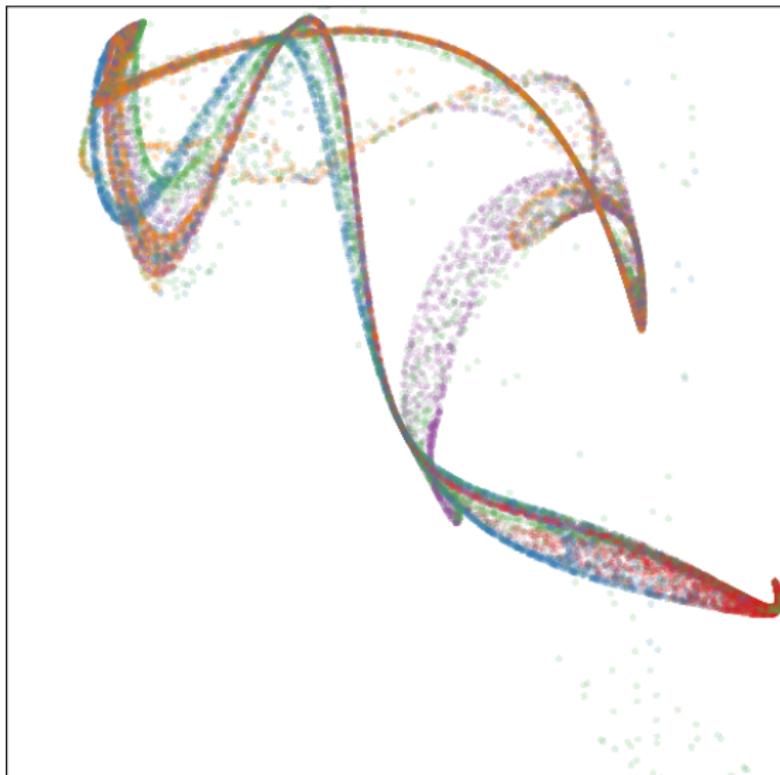
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



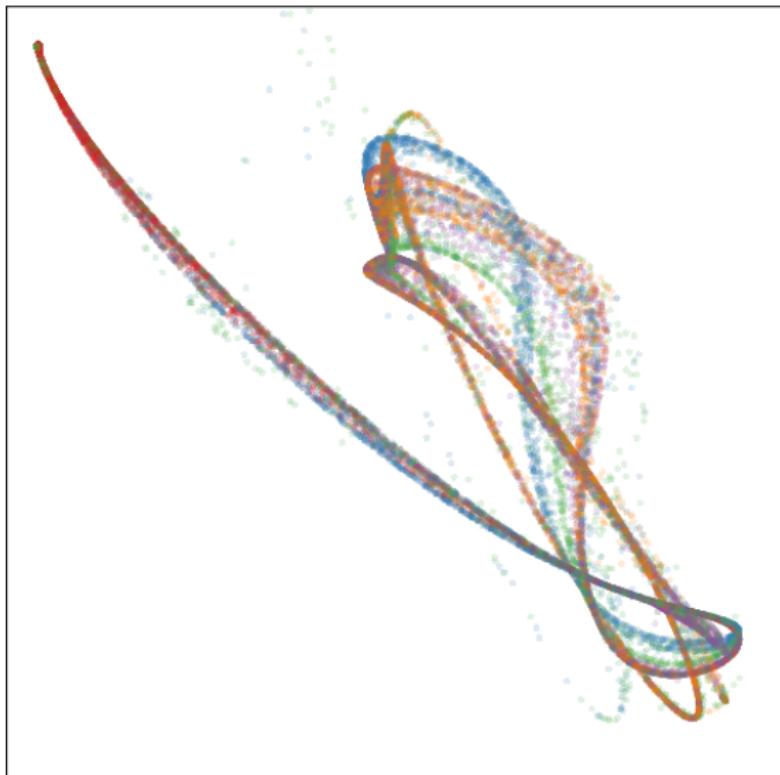
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



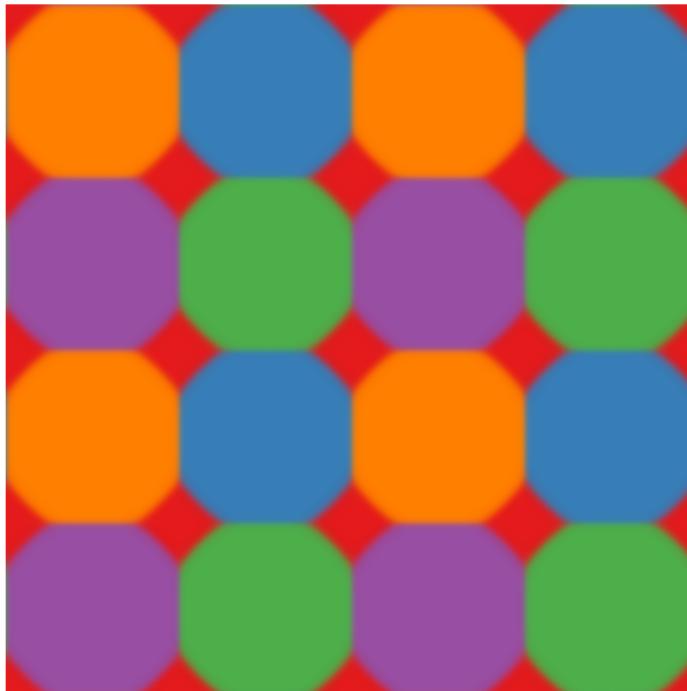
# DEEP GAUSSIAN PROCESSES

- ▶ 2D to 2D warpings of a Gaussian density:



# DEEP GAUSSIAN PROCESSES

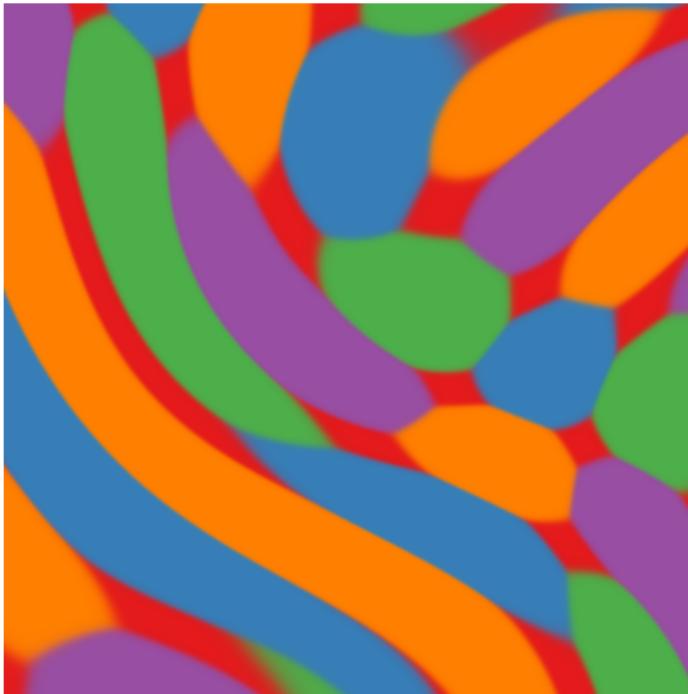
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



No warping

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



One Layers

# DEEP GAUSSIAN PROCESSES

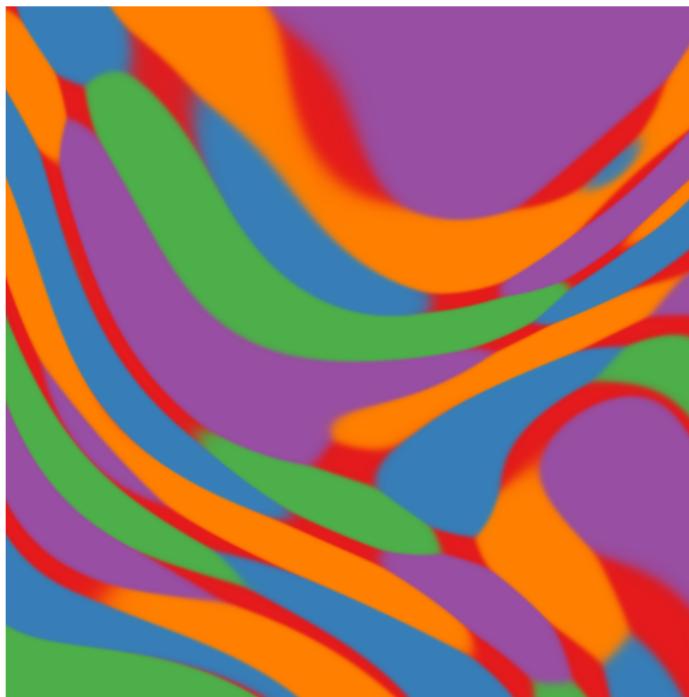
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Two Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Three Layers

# DEEP GAUSSIAN PROCESSES

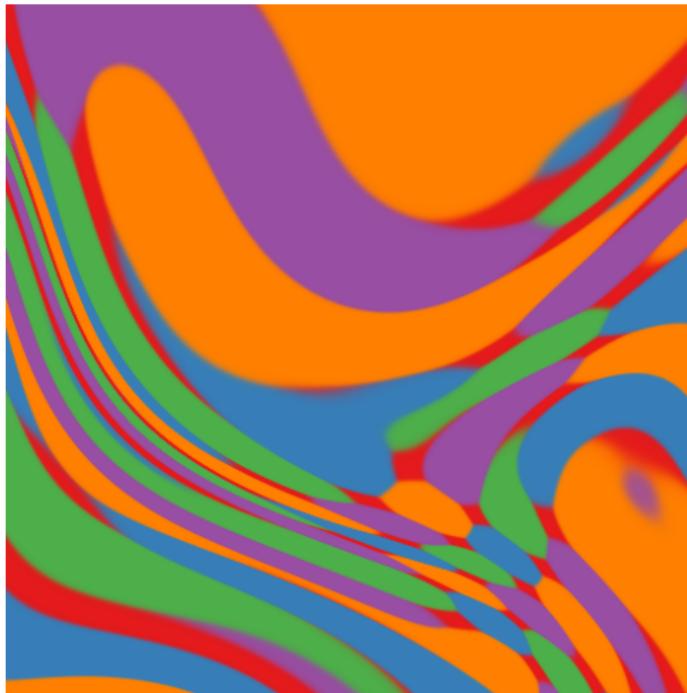
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Four Layers

# DEEP GAUSSIAN PROCESSES

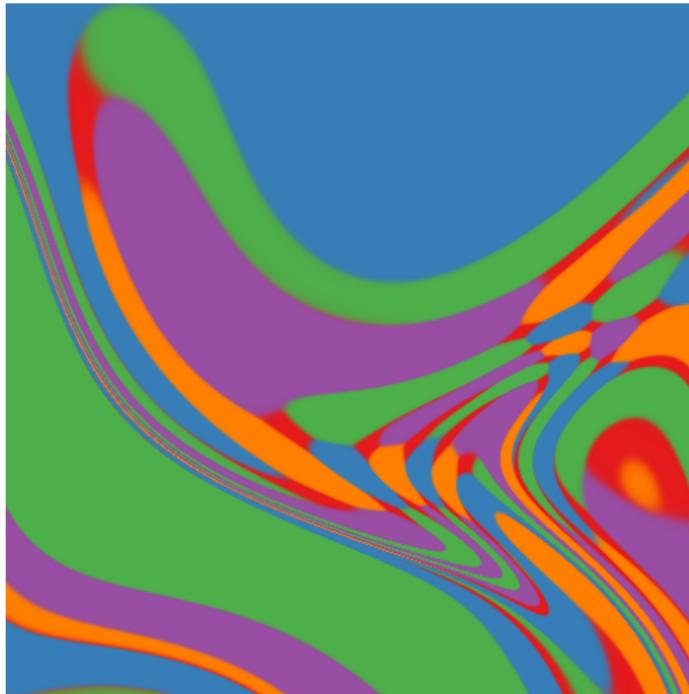
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Five Layers

# DEEP GAUSSIAN PROCESSES

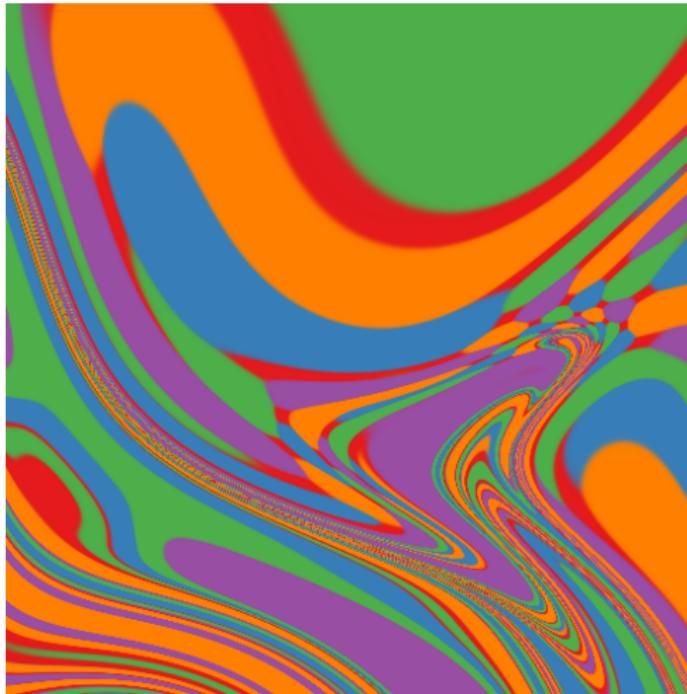
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Ten Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Twenty Layers

# DEEP GAUSSIAN PROCESSES

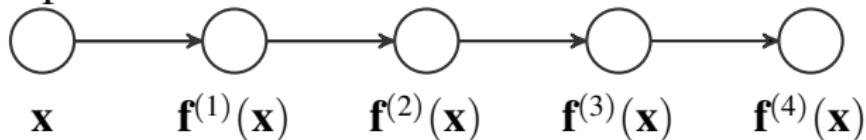
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



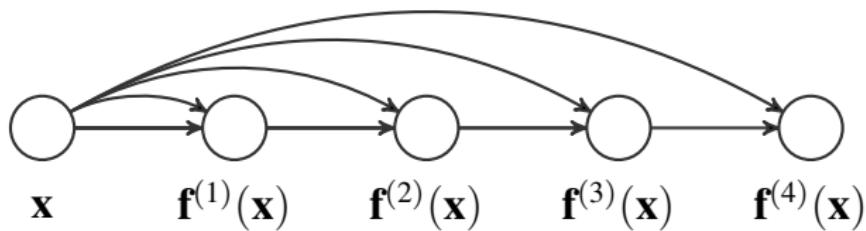
Forty Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Again following Radford's thesis, connect every layer to input:



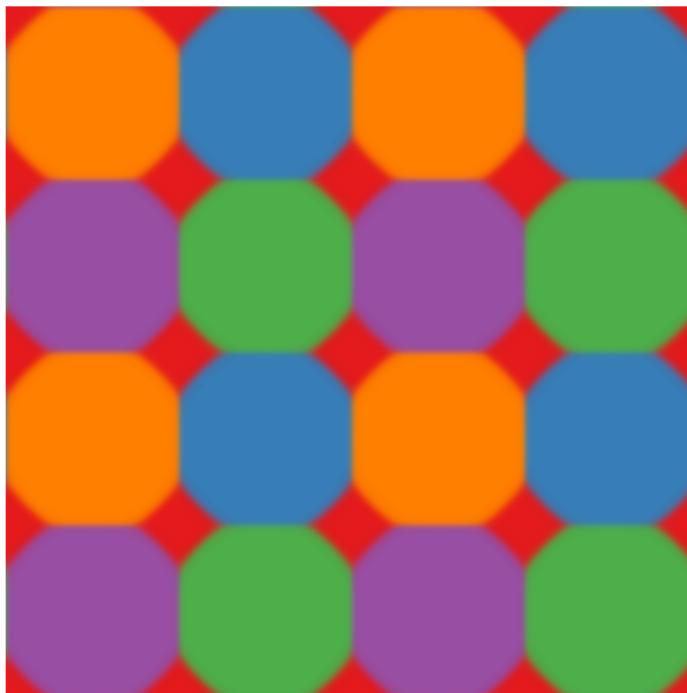
a) standard MLP architecture.



b) Input-connected architecture.

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



No warping

# DEEP GAUSSIAN PROCESSES

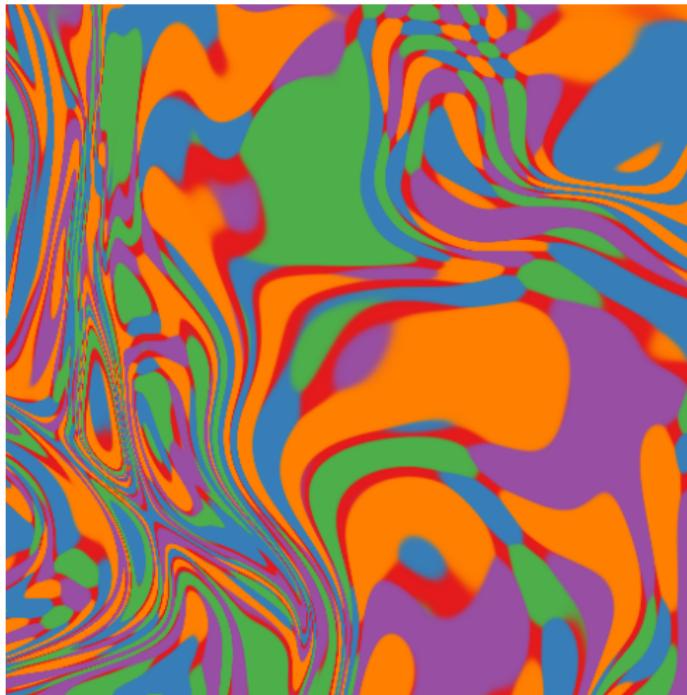
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Two Layers

# DEEP GAUSSIAN PROCESSES

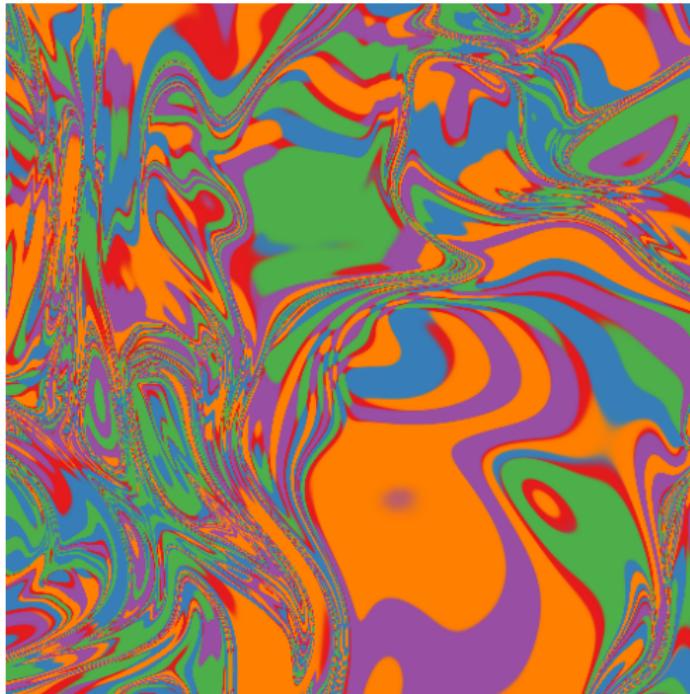
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Ten Layers

# DEEP GAUSSIAN PROCESSES

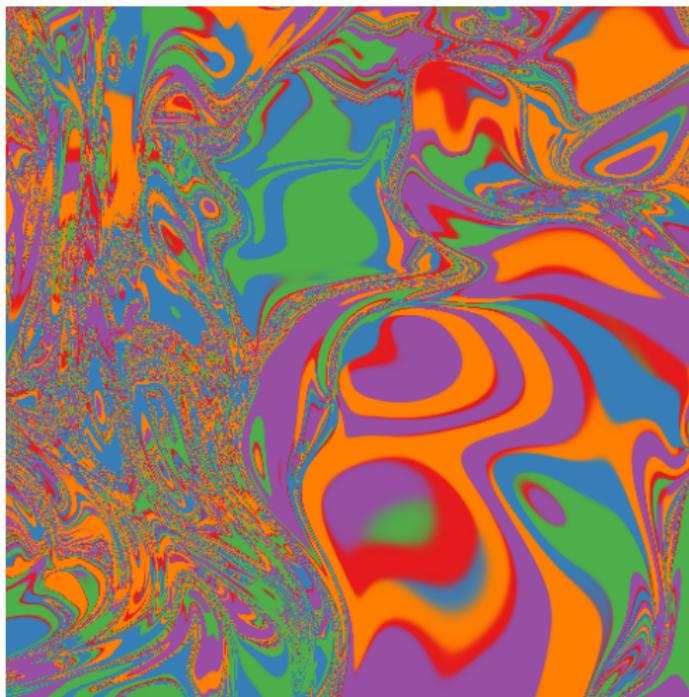
- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Twenty Layers

# DEEP GAUSSIAN PROCESSES

- ▶ Showing the  $x$  that gave rise to a particular  $y$
- ▶ (i.e. decision boundaries)



Forty Layers

# SUMMARY

- ▶ GPs let us compute marginal likelihood, which enables model search
- ▶ Kernel learning is an example of representation learning
- ▶ How to build priors over functions that put mass on the sorts of structures we'd like to be able to learn?
- ▶ This question might shed light on good architectures or initialization strategies in a data-independent way.
- ▶ Open questions:
  - ▶ When is discrete kernel search a good way to find representations?
  - ▶ What could we do if we could compute the marginal likelihood of a neural net?

# SUMMARY

- ▶ GPs let us compute marginal likelihood, which enables model search
- ▶ Kernel learning is an example of representation learning
- ▶ How to build priors over functions that put mass on the sorts of structures we'd like to be able to learn?
- ▶ This question might shed light on good architectures or initialization strategies in a data-independent way.
- ▶ Open questions:
  - ▶ When is discrete kernel search a good way to find representations?
  - ▶ What could we do if we could compute the marginal likelihood of a neural net?

Thanks!