

DS-GA 3001.008 Modelling time series data

L8a. Advanced topics in GPs.

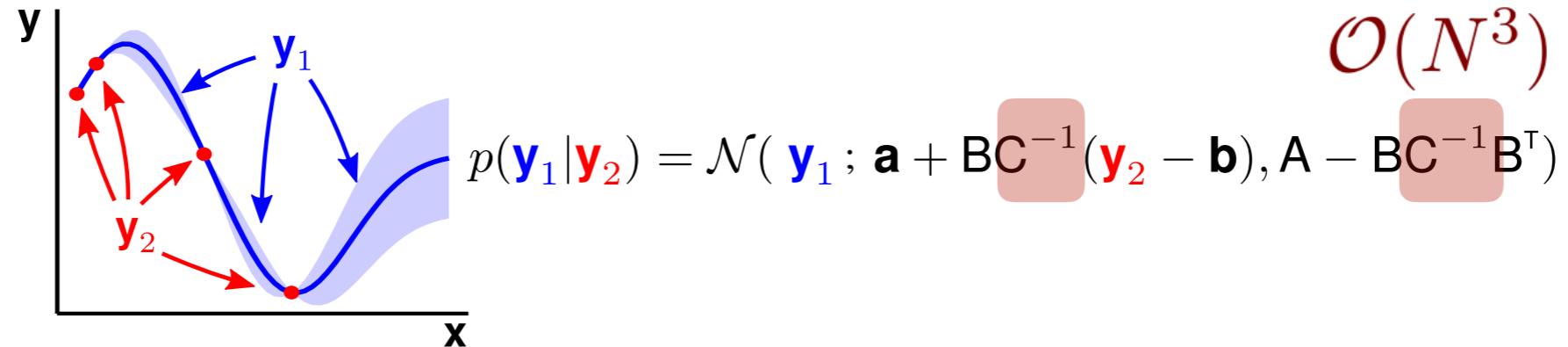
Instructor: Cristina Savin
NYU, CNS & CDS

Def: A **Gaussian Process** is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

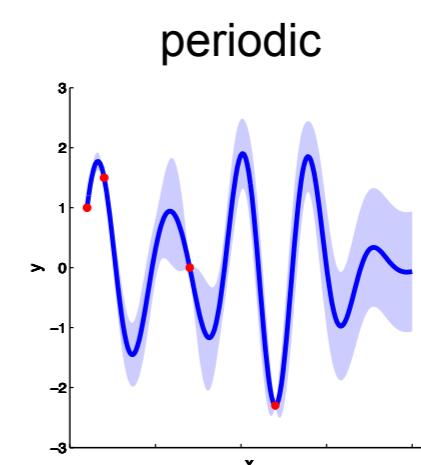
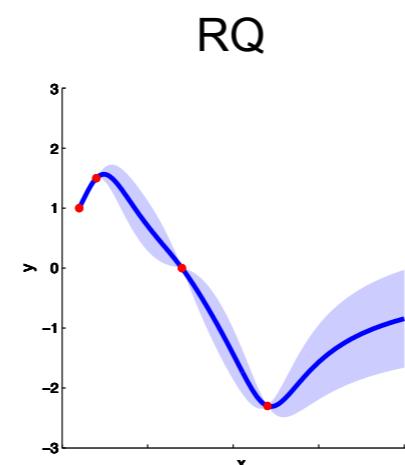
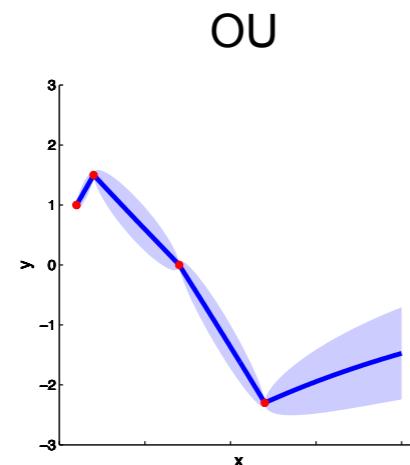
It is essentially a infinite dimensional generalization of a multivariate gaussian;

We call this a '**non-parametric**' model because the number of parameters (the covariance matrix) increases with the number of data points, N.

Inference



Kernels:



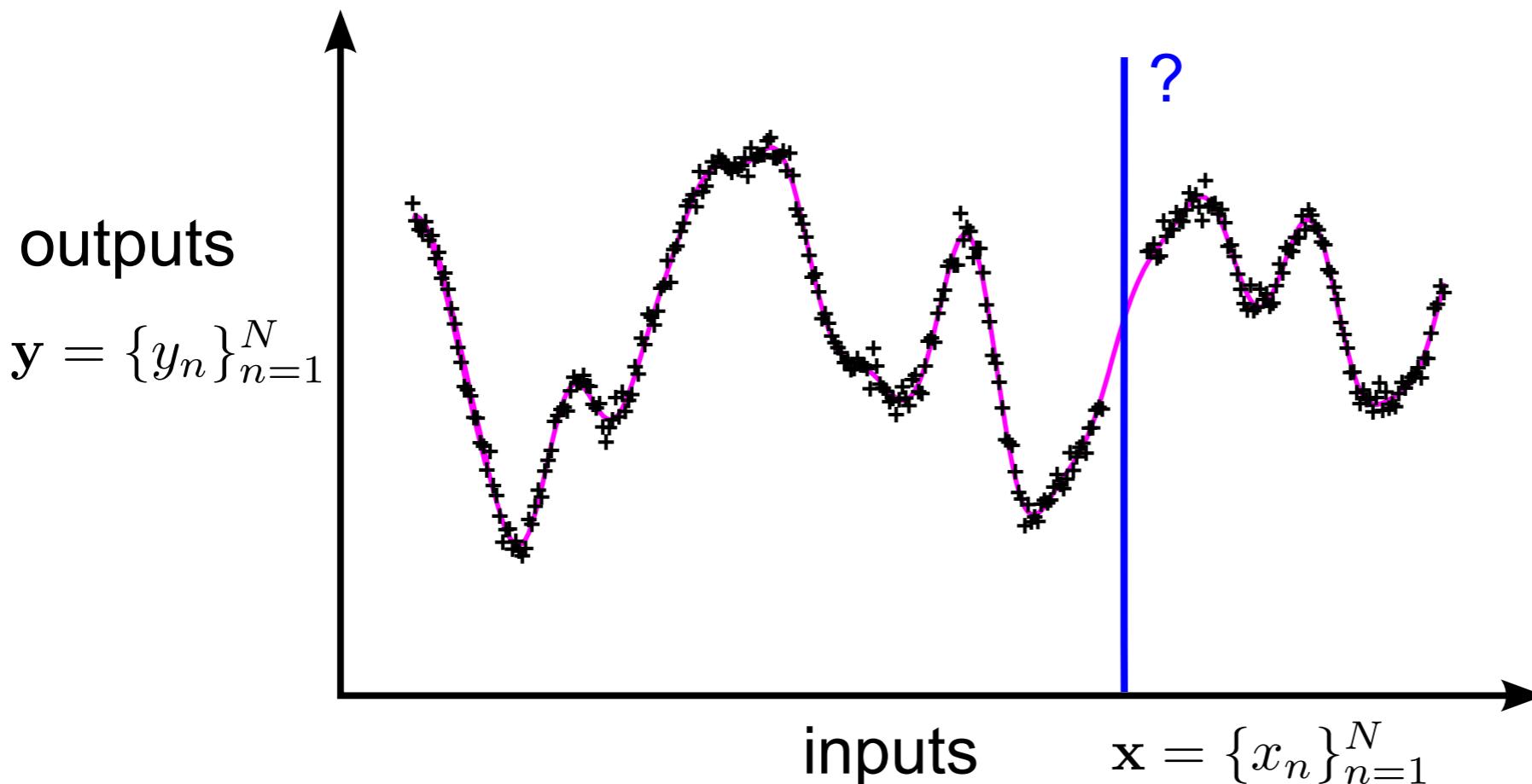
GP credits: Rich Turner

Efficient inference for large N

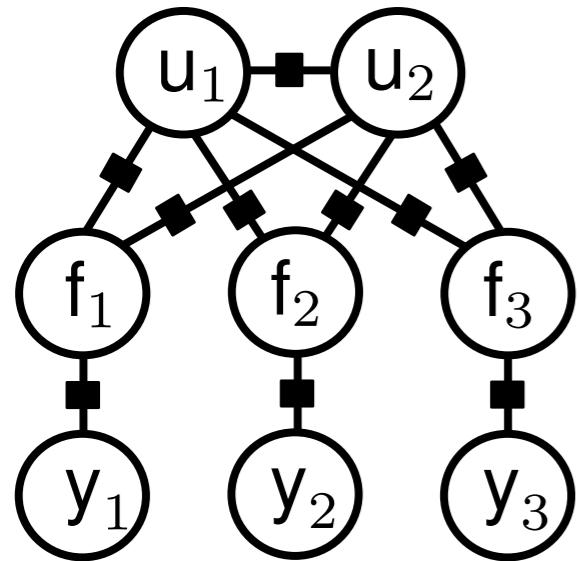
$$p(f|\theta) = \mathcal{GP}(f; 0, K_\theta)$$

$$p(y_n|f, x_n, \theta)$$

IDEA: maybe we don't **really** need **all** the data



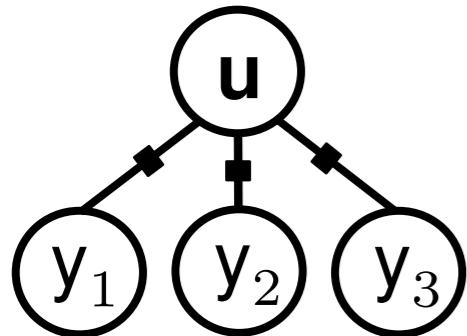
FITC: fully independent training conditional approximation



Introduce $M < N$ inducing points (pseudo-data) \mathbf{u}

These \mathbf{u} variables summarize (as much as possible) the dependencies in \mathbf{f} , while keeping costs low

$$p(\mathbf{f}_t | \mathbf{u}) = \mathcal{N}(\mathbf{f}_t; \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{u}, \mathbf{K}_{\mathbf{f}_t \mathbf{f}_t} - \mathbf{K}_{\mathbf{f}_t \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{f}_t})$$



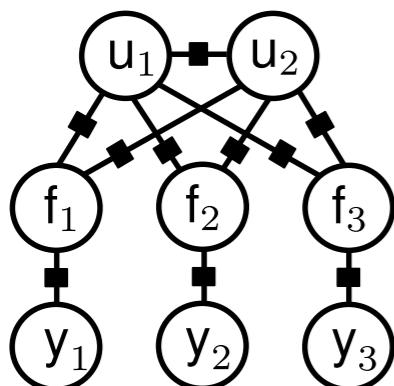
$$p(\mathbf{y}_t | \theta) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{K}_{\mathbf{f} \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{u}} \mathbf{K}_{\mathbf{u} \mathbf{f}}^{-1} \mathbf{K}_{\mathbf{u} \mathbf{f}} + \mathbf{D} + \sigma_y^2 \mathbf{I})$$

A Brief History of Gaussian Process Approximations

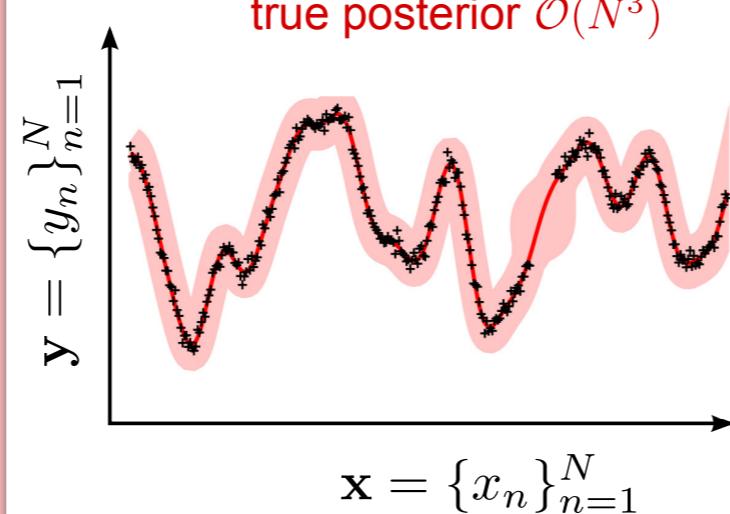
approximate generative model
exact inference

$$\text{div}[p(\mathbf{f}, \mathbf{y}) || q(\mathbf{f}, \mathbf{y})]$$

A Unifying View of Sparse
Approximate Gaussian
Process Regression
Quinonero-Candela &
Rasmussen, 2005
(FITC, PITC, DTC)



methods employing
pseudo-data



exact generative model
approximate inference

$$\text{div}[p(\mathbf{f}|\mathbf{y}) || q(\mathbf{f})]$$

A Unifying Framework for
Sparse Gaussian Process
Approximation using
Power Expectation
Propagation
Bui, Yan and Turner, 2016
(VFE, EP, FITC, PITC ...)

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

EP: Csato and Opper 2002 / Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

VFE: Titsias "Variational Learning of Inducing Variables in Sparse Gaussian Processes"

DTC / PP: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

**Other tricks: building extra structure in the covariance matrix
so you don't have to do the inversion the usual way.**

Kronecker tricks

$$\mathbf{K} = \bigotimes_i \mathbf{K}_i$$

Kronecker product

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{11}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}$$

Matrix operations decompose nicely

$$\text{Transpose : } \mathbf{K}^T = \bigotimes_i \mathbf{K}_i^T$$

$$\text{Inverse : } \mathbf{K}^{-1} = \bigotimes_i \mathbf{K}_i^{-1}$$

$$\text{Eigen - decomposition : } \mathbf{K} = \bigotimes_i \mathbf{Q}_i \cdot \bigotimes_i \boldsymbol{\Lambda}_i \cdot \bigotimes_i \mathbf{Q}_i^T$$

$$\text{Determinants : } \det(\mathbf{K}) = \prod_i \det(\mathbf{K}_i)^{D_i}$$

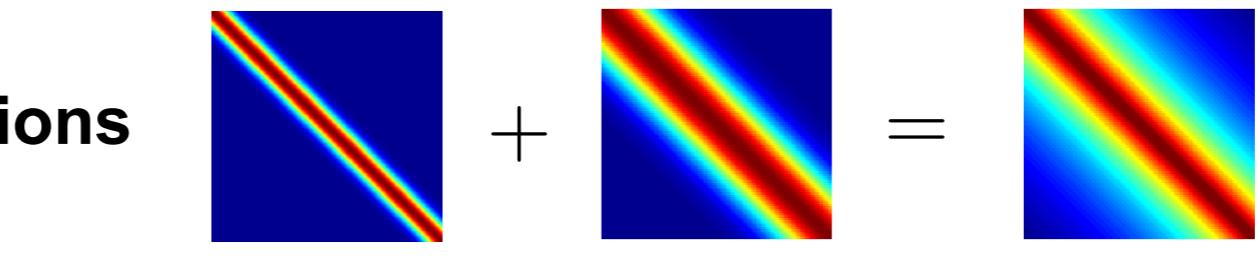
$$\text{Trace : } \text{tr}(\mathbf{K}) = \prod_i \text{tr}(\mathbf{K}_i)$$

$$\text{Vec : } \text{vec}(\mathbf{AXB}) = \mathbf{A} \otimes \mathbf{B} \text{ vec}(\mathbf{X})$$

How to choose the kernel?

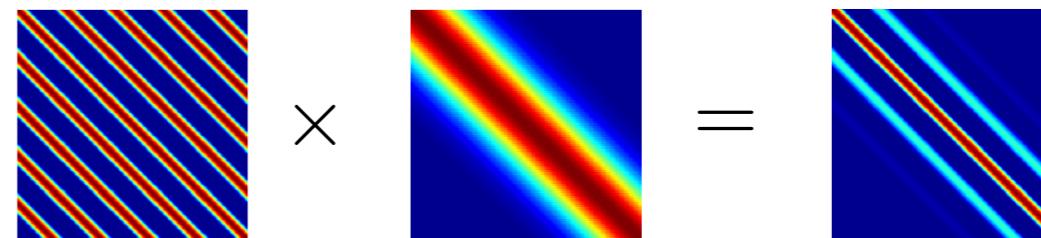
A1: mix and match

(positive) linear combinations
of covariance functions



e.g. scale mixture of SE = rational quadratic

multiplication of covariance
functions



e.g. periodic SE $\cos(\omega \Delta x) \exp(-\frac{1}{2l^2} \Delta x^2)$

derivative of GP = GP

$$\frac{d}{dx} y(x) = \frac{d}{dx} \sum_{k=1}^{\infty} \gamma_k g_k(x) = \sum_{k=1}^{\infty} \gamma_k \frac{d}{dx} g_k(x)$$

new basis: $g'_k(x) = \frac{d}{dx} g_k(x)$

new covariance: $K'(x, x') = \frac{d}{dx} \frac{d}{dx'} K(x, x')$

integral of GP = GP

$$\int dx y(x) = \sum_{k=1}^{\infty} \gamma_k \int dx g_k(x)$$

new basis: $g'_k(x) = \int dx g_k(x)$

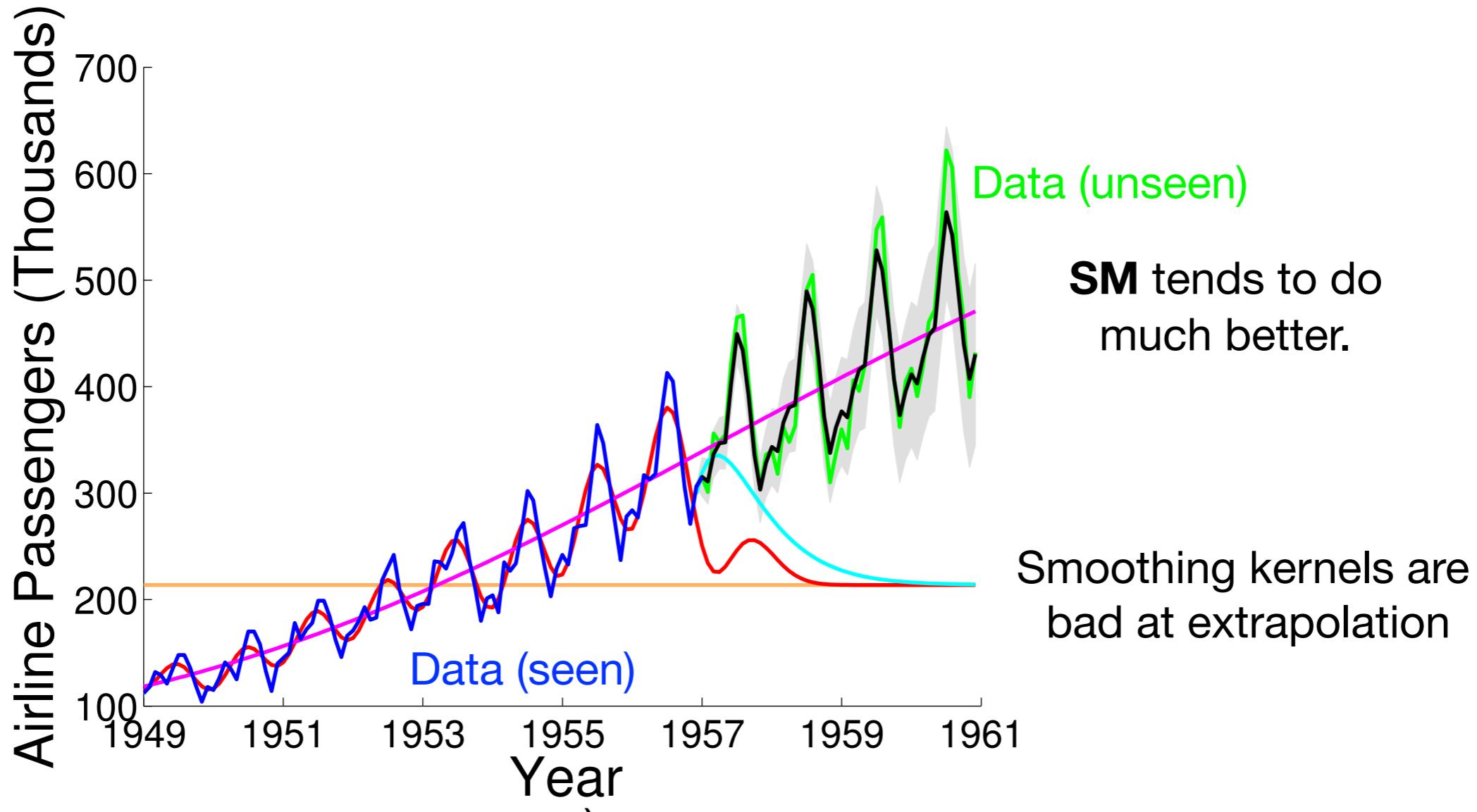
new covariance: $K'(x, x') = \int \int dx dx' K(x, x')$

How to choose the kernel?

A2: don't bother, learn that from data too (**SM**)

$$k(\tau) = \sum_{q=1}^Q w_q \prod_{p=1}^P \exp\{-2\pi^2 \tau_p^2 v_q^{(p)}\} \cos(2\pi \tau_p \mu_q^{(p)}).$$

Mixture of
Q components
Gaussian-like
components
(in Fourier domain)



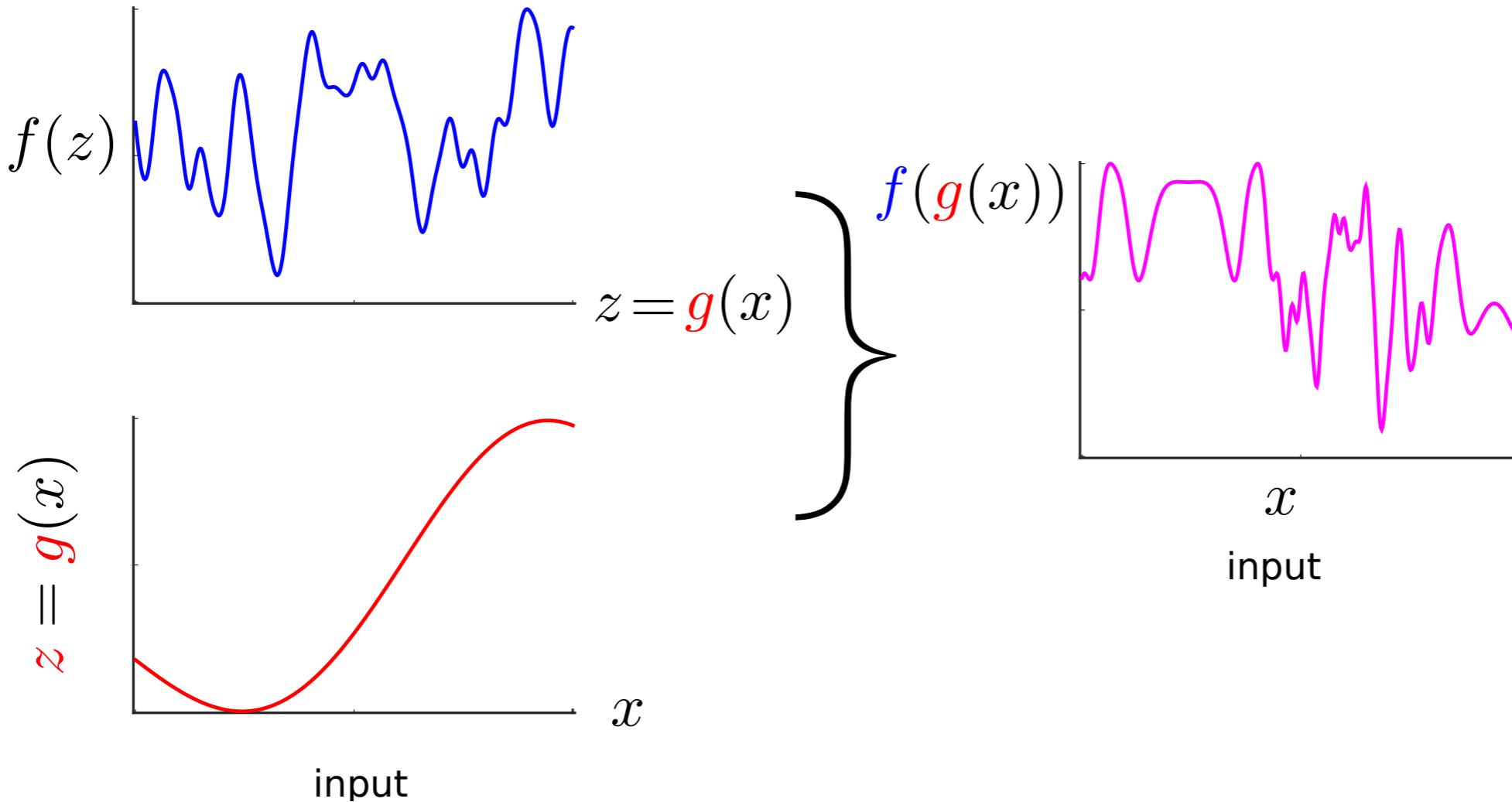
How to choose the kernel?

A3: go deep!

$$y(x) = \textcolor{blue}{f}(\textcolor{red}{g}(x)) + \sigma_y \epsilon$$

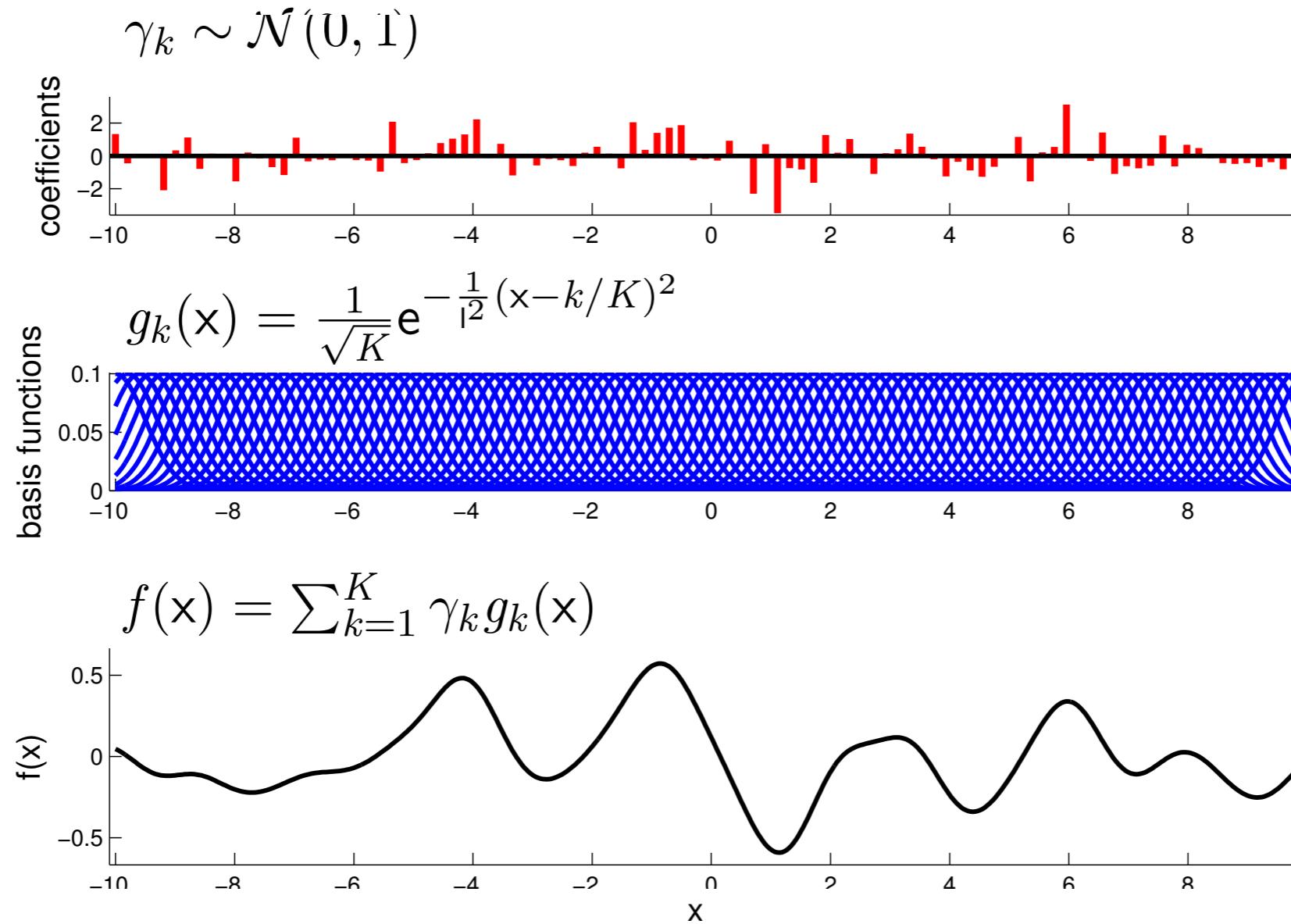
$$f(x) = \mathcal{GP}(0, K_f(x, x'))$$

$$g(x) = \mathcal{GP}(0, K_g(x, x'))$$

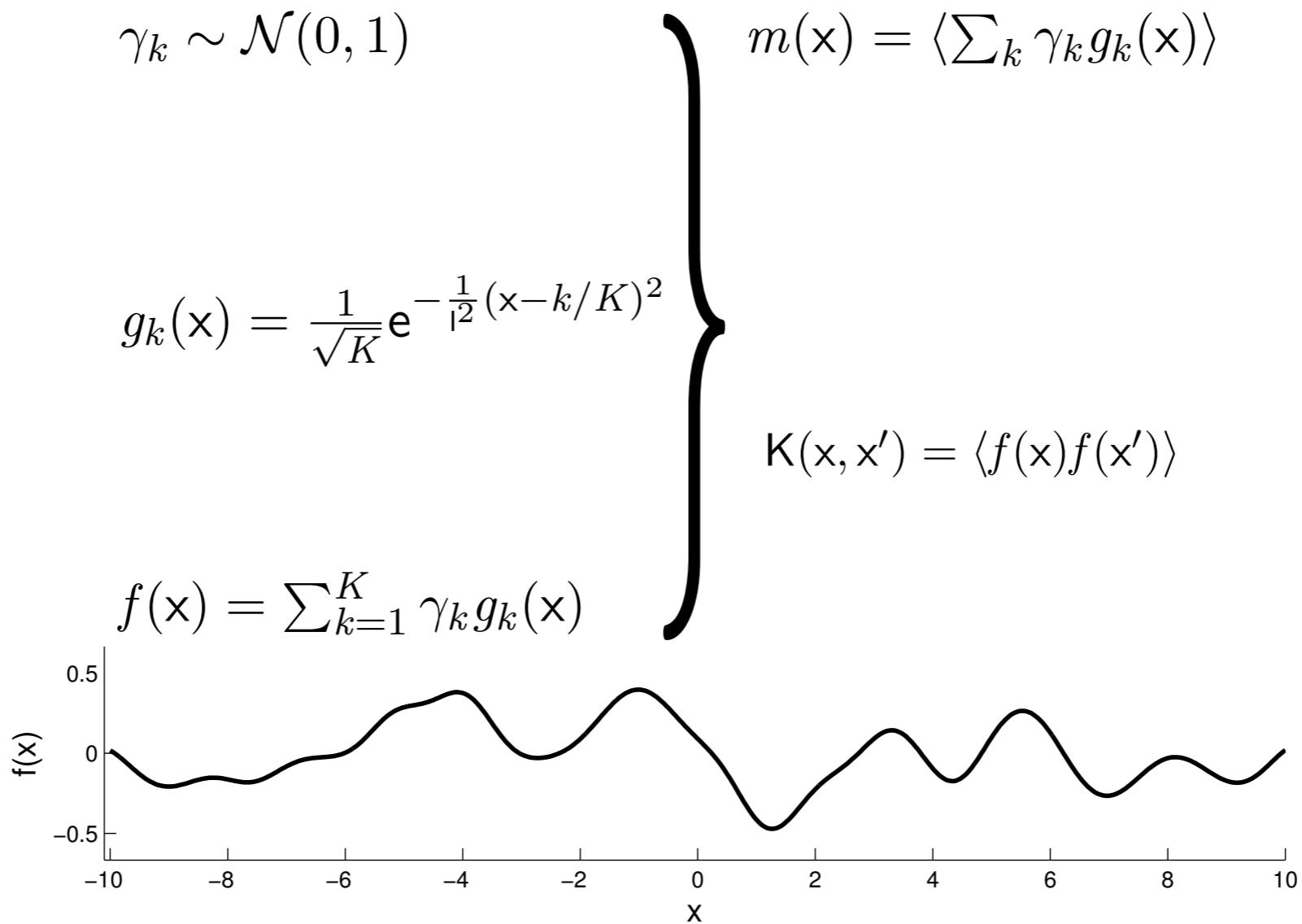


What is a kernel anyways?

We can formally link **GPs** with usual kernel methods (**SVM and co**) and **NNs**

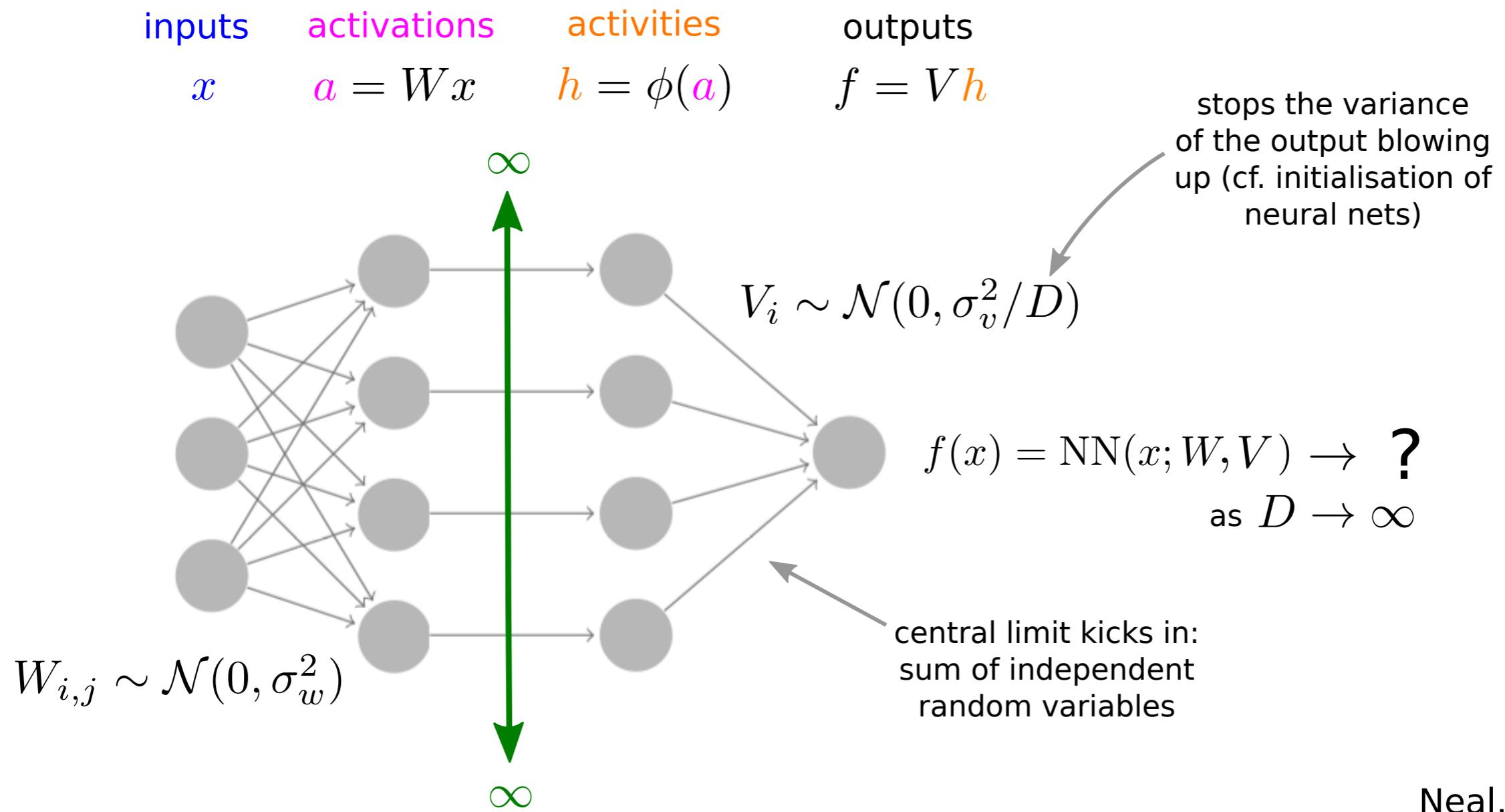


Basis function view of Gaussian processes



GP = infinite basis

Infinitely wide neural nets as GPs



Summary

GPs are a simple and elegant way of defining distributions over functions

Building block for more complex models (audio textures in last lecture)

Proper treatment of uncertainty (optimal control, active sensing), Occam's razor

Deep mathematical connections to a variety of popular ML models



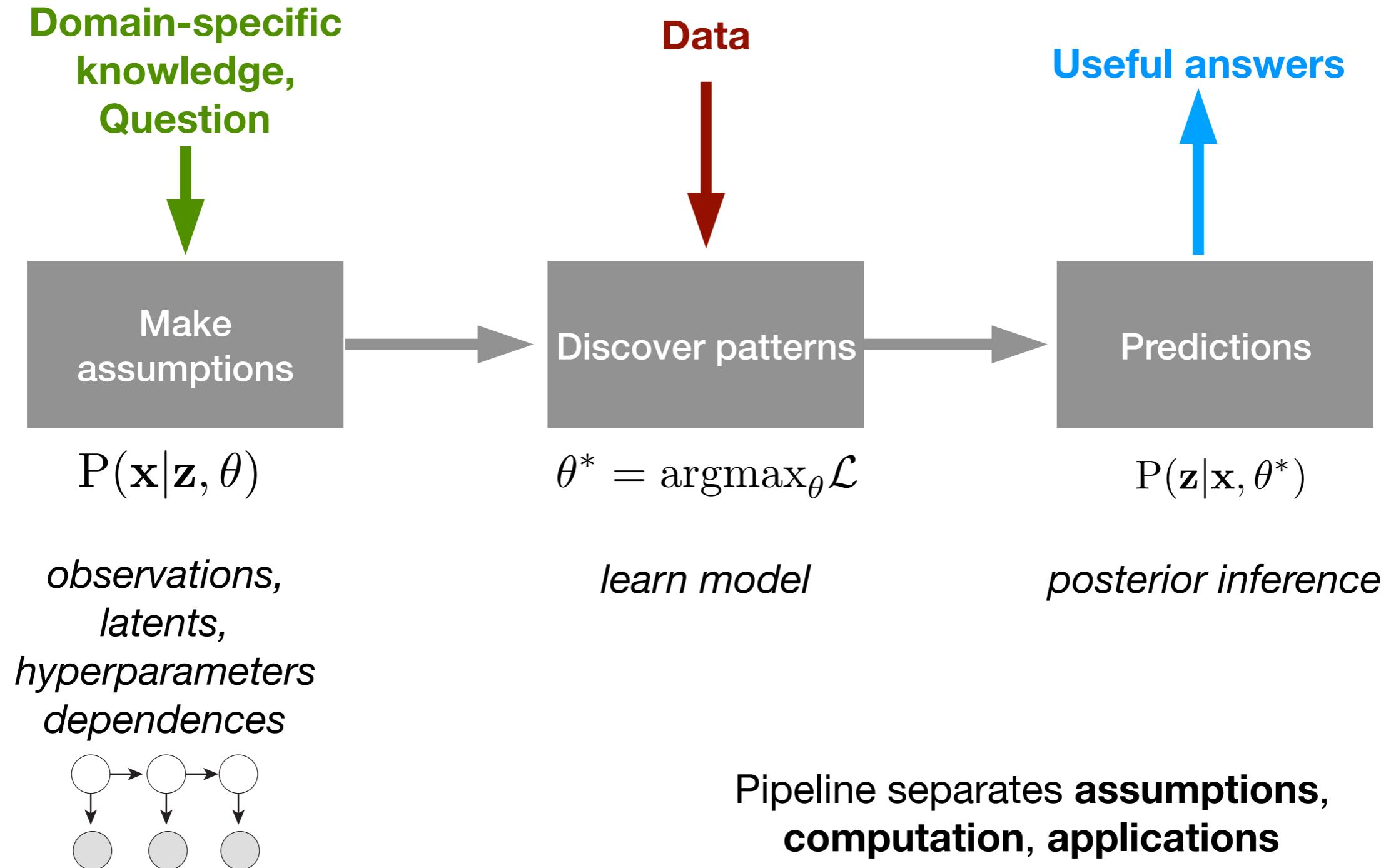
DS-GA 3001.008 Modelling time series data

L8b. Intro to RNNs

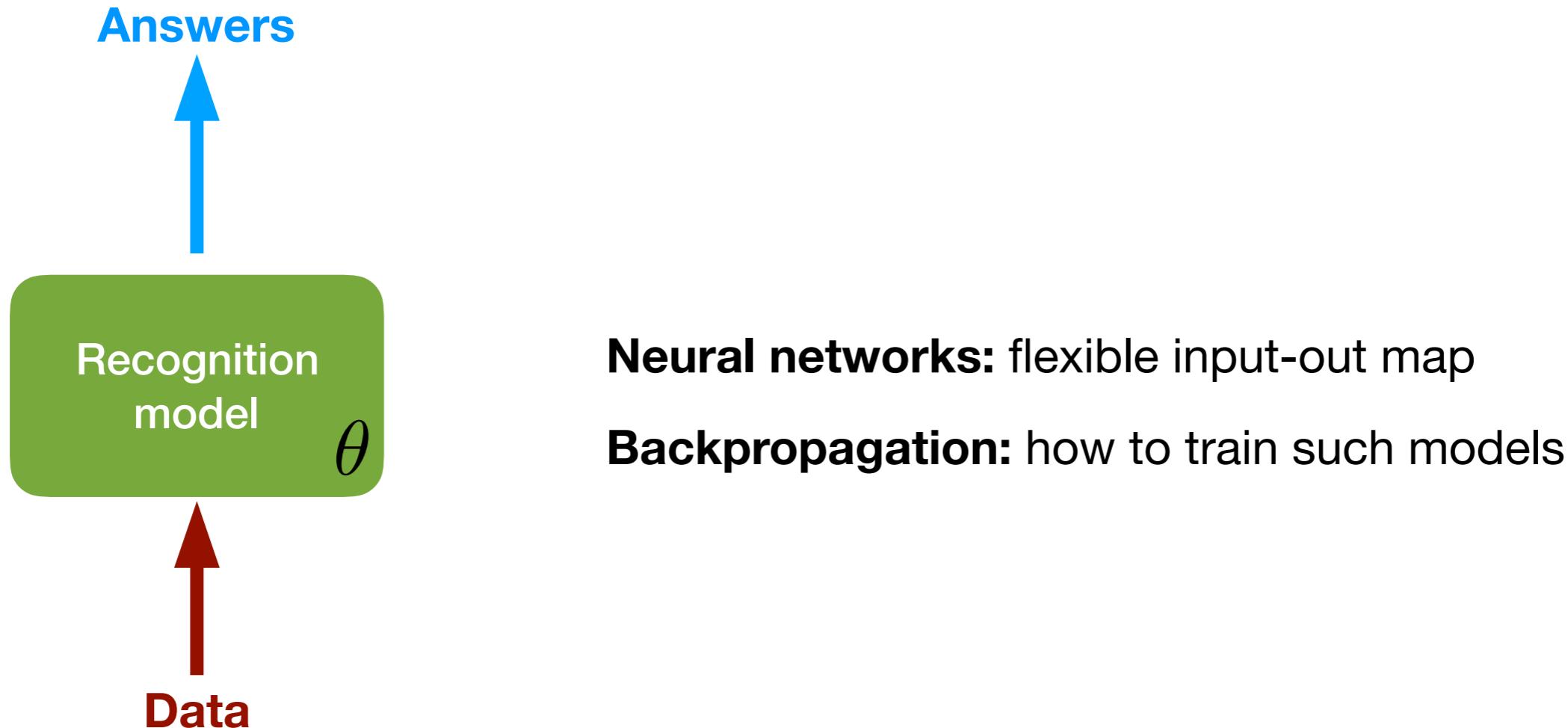
Instructor: Cristina Savin

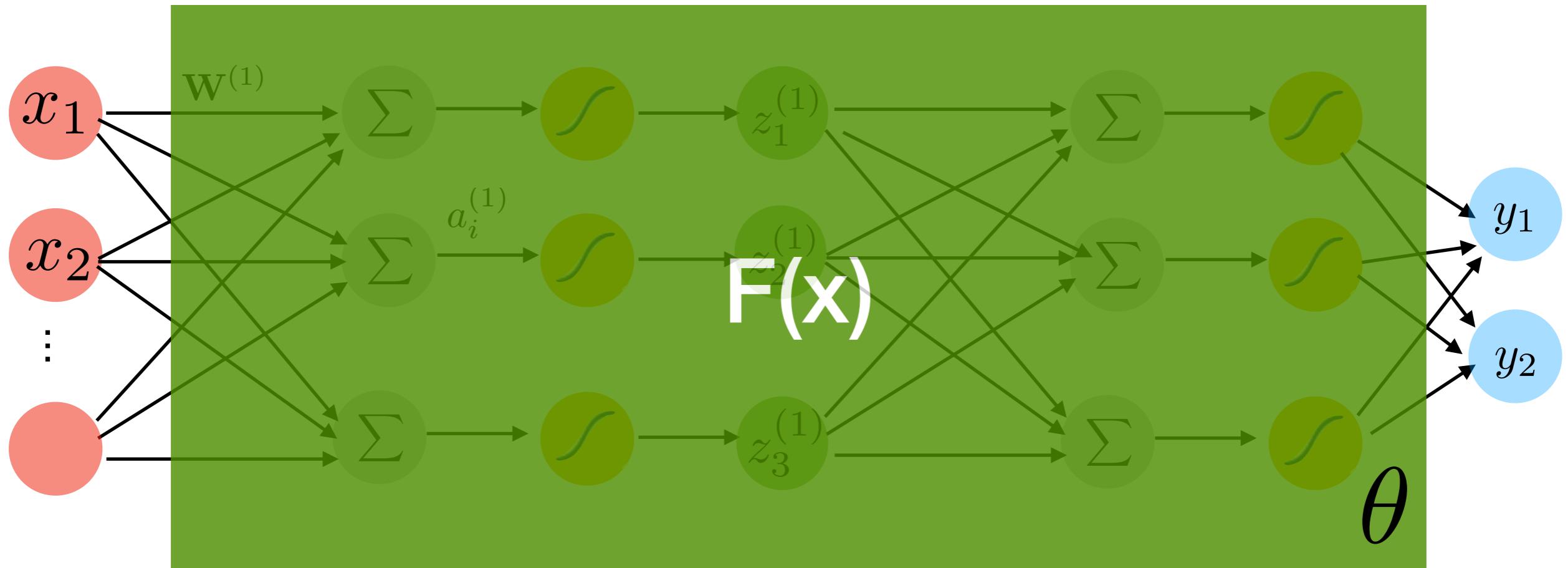
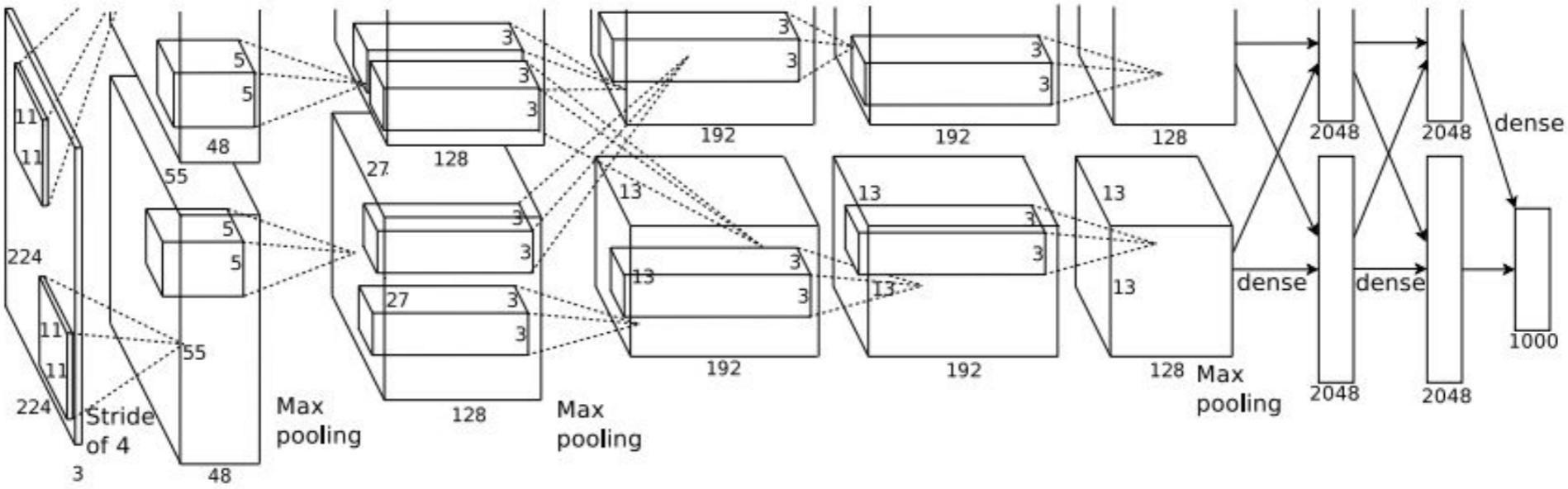
NYU, CNS & CDS

The probabilistic ‘pipeline’



The alternative: just build a recognition model from the get go

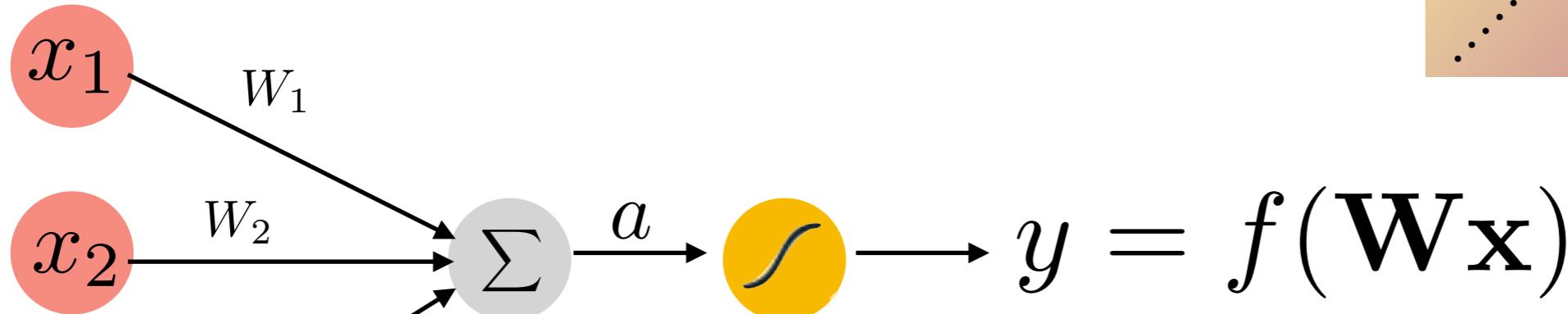




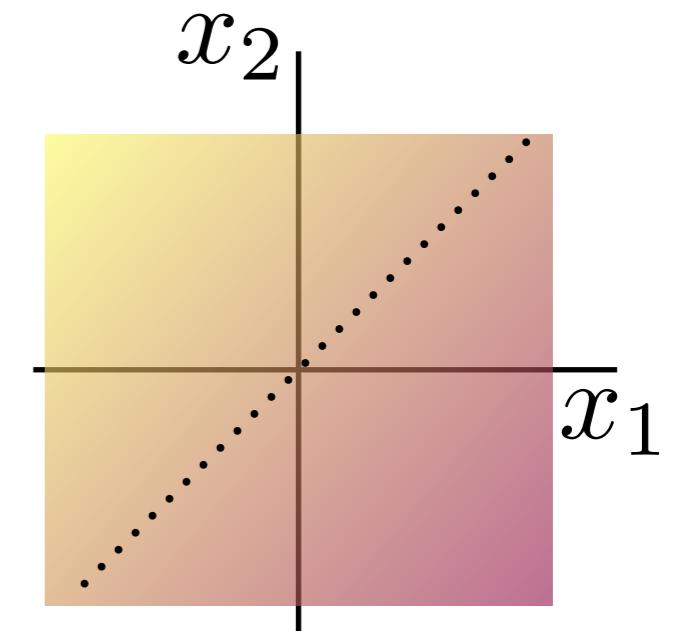
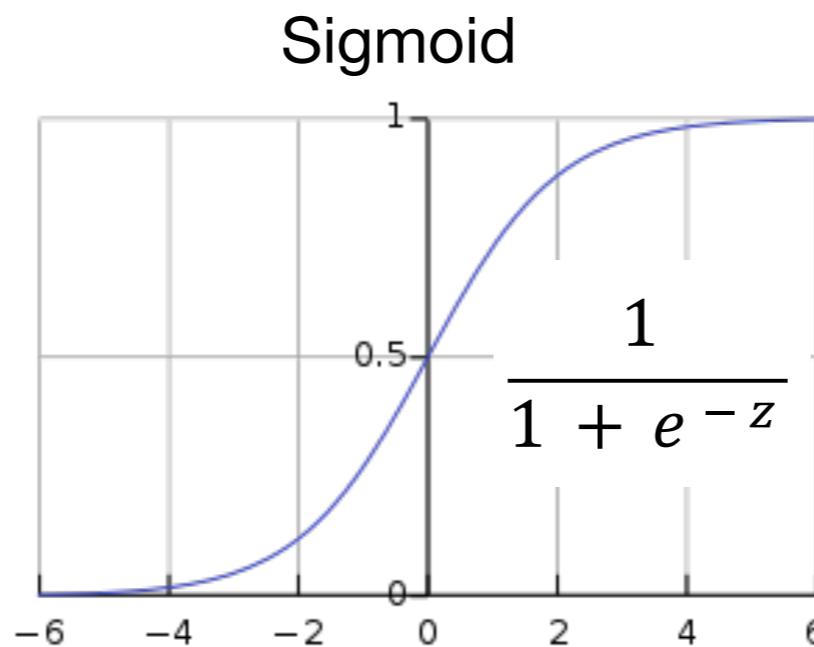
A cascade of linear-nonlinear steps-> a differentiable map

The basic building block: perceptron

inputs



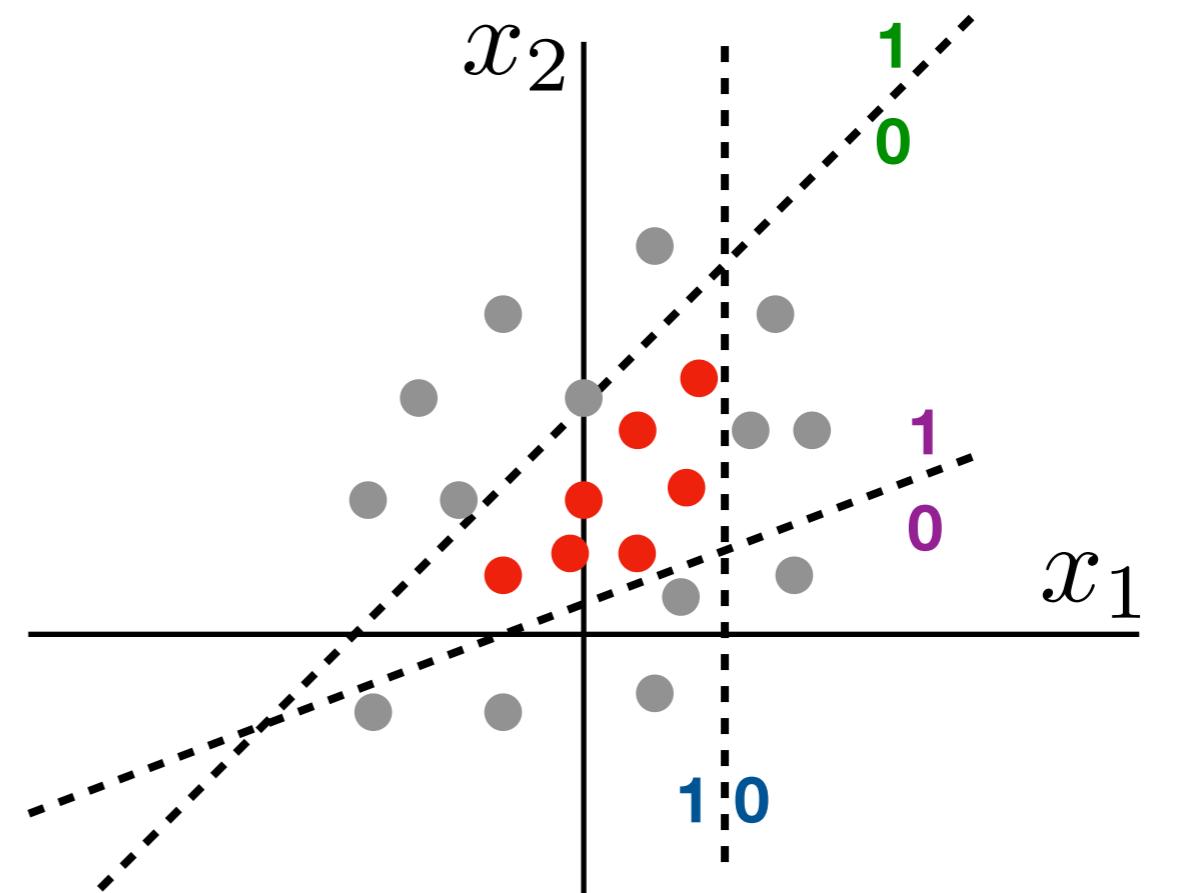
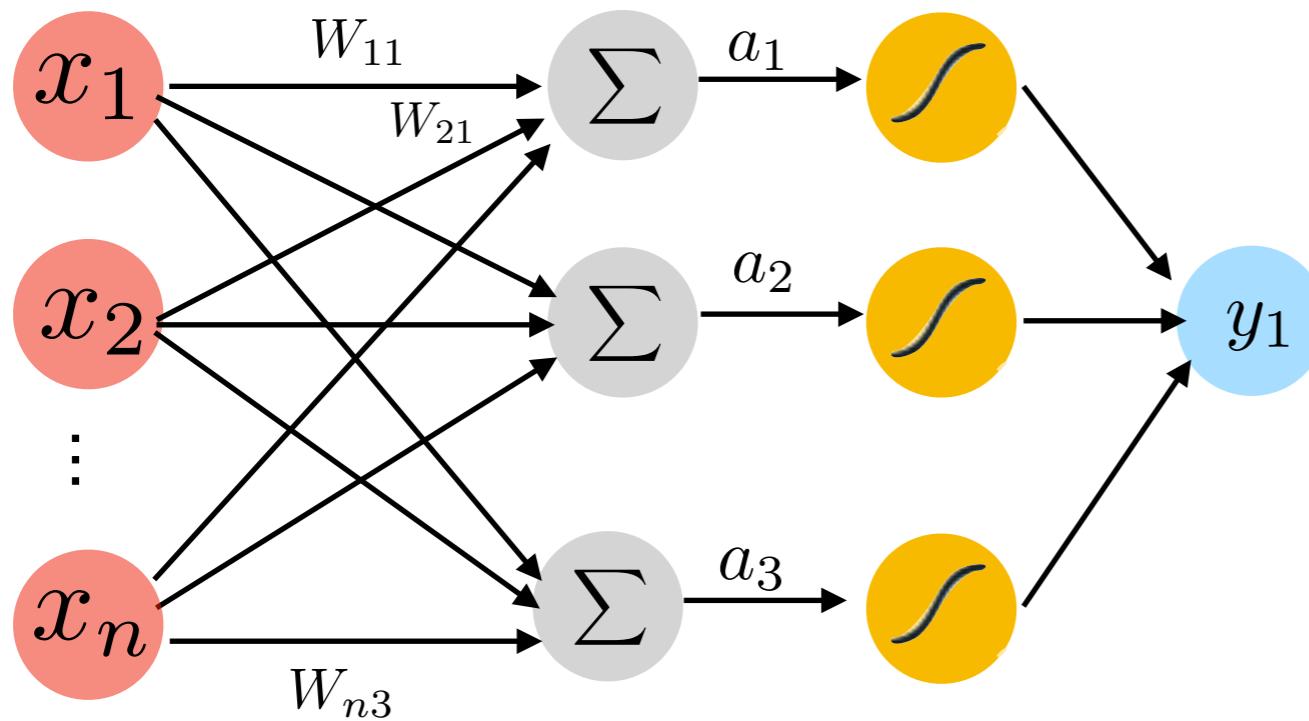
nonlinear
activation



tanh:
$$\frac{e^z - e^{-z}}{e^z + e^{-z}}$$

ReLU:
$$\max(0, z)$$

One layer



**Efficient partitioning
of input space**

Universal approximation theorem (Hornik, '91)

“A single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units.”

This works for any standard nonlinearities f .

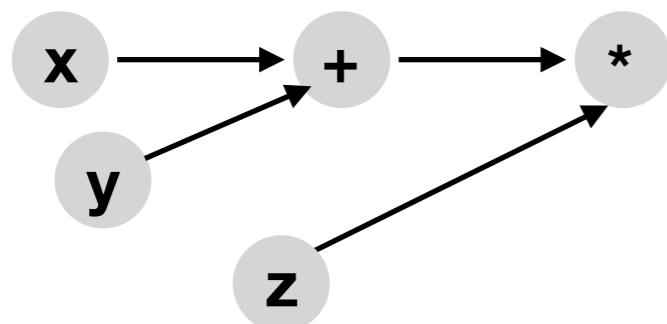
Caveat: It does not necessarily mean that there exists a learning algorithm that finds the correct solution

How do we compute the gradient?

$$\frac{\partial J(\theta)}{\partial \theta}$$

Simple trick: repeatedly apply the chain rule

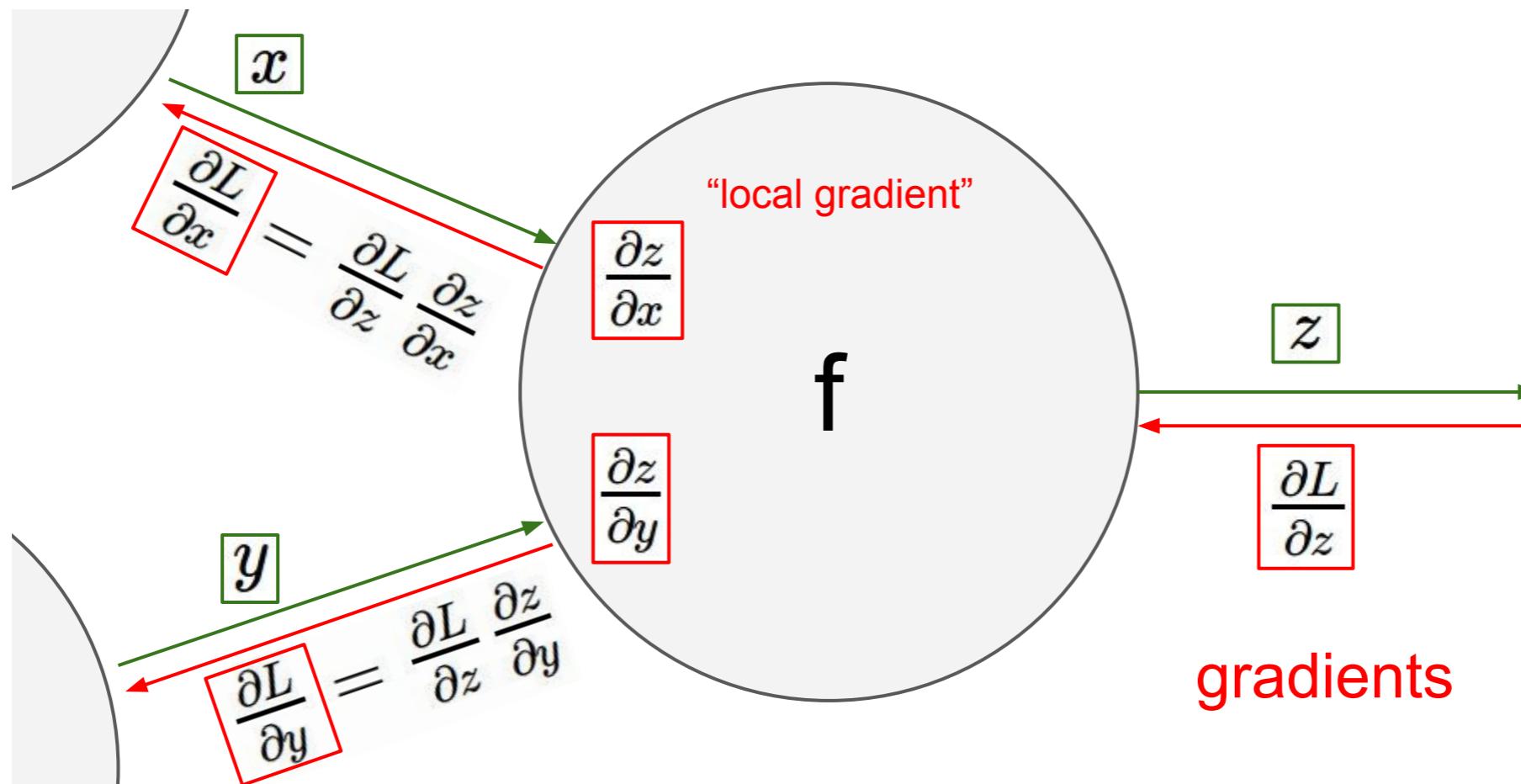
Let's take some simple examples (on the board)



$$f(x, y, z) = (x + y)z$$

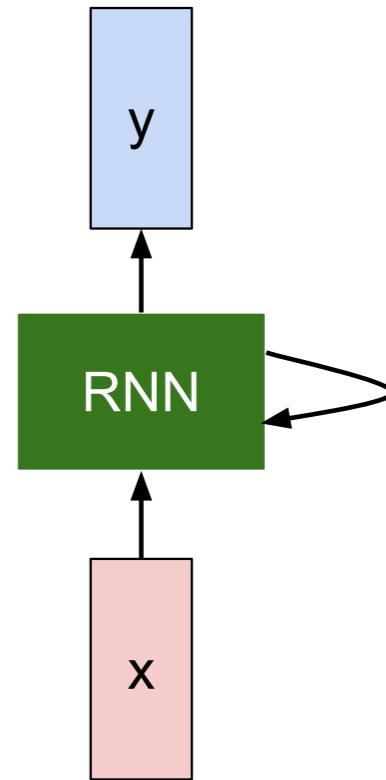
$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

Backpropagation



Vanilla RNNs math

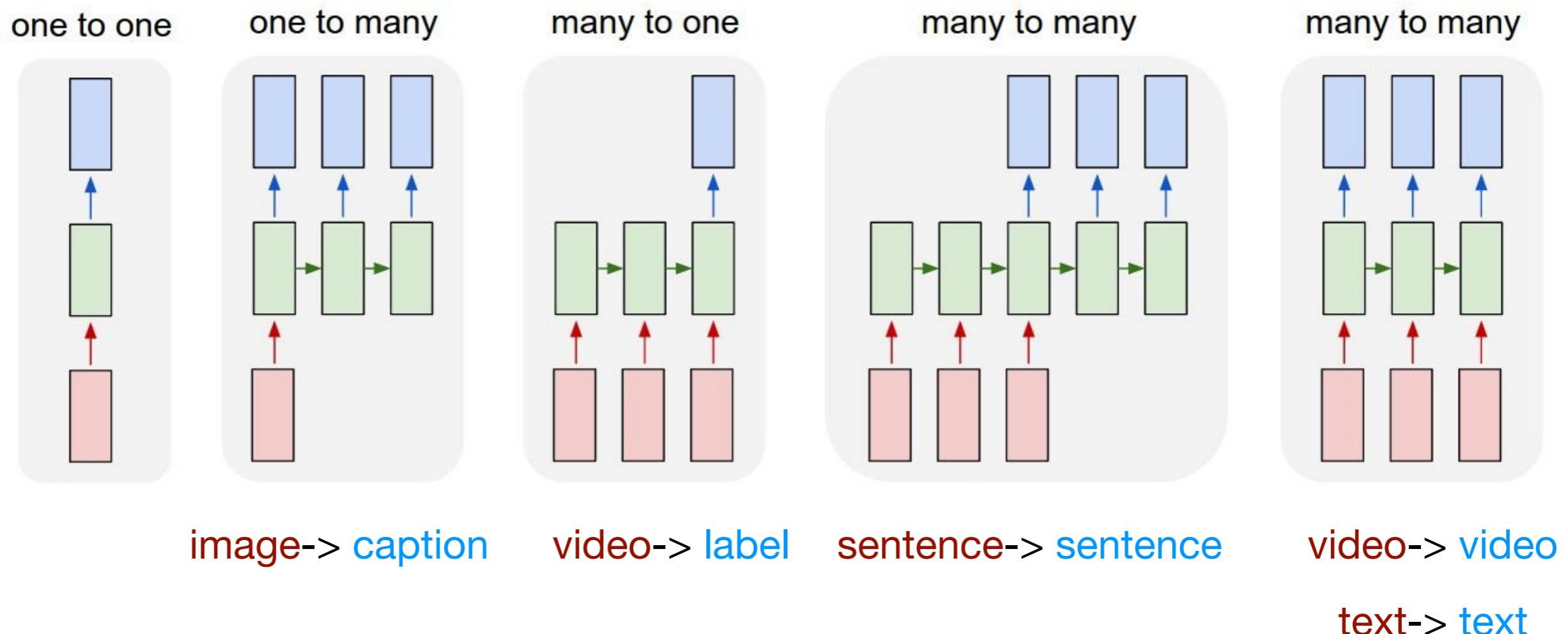
We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:



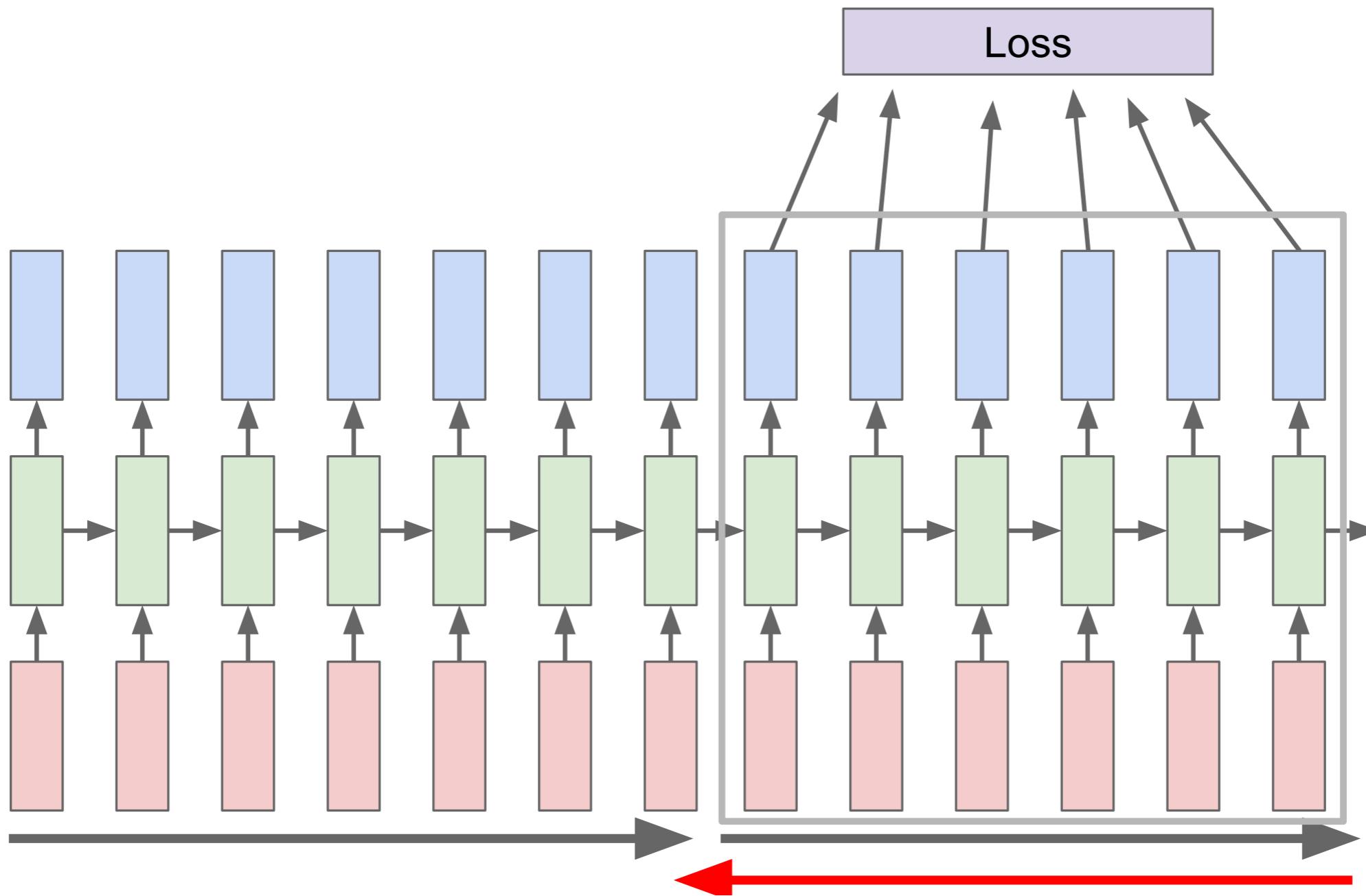
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy} h_t$$

RNN architectures for time series

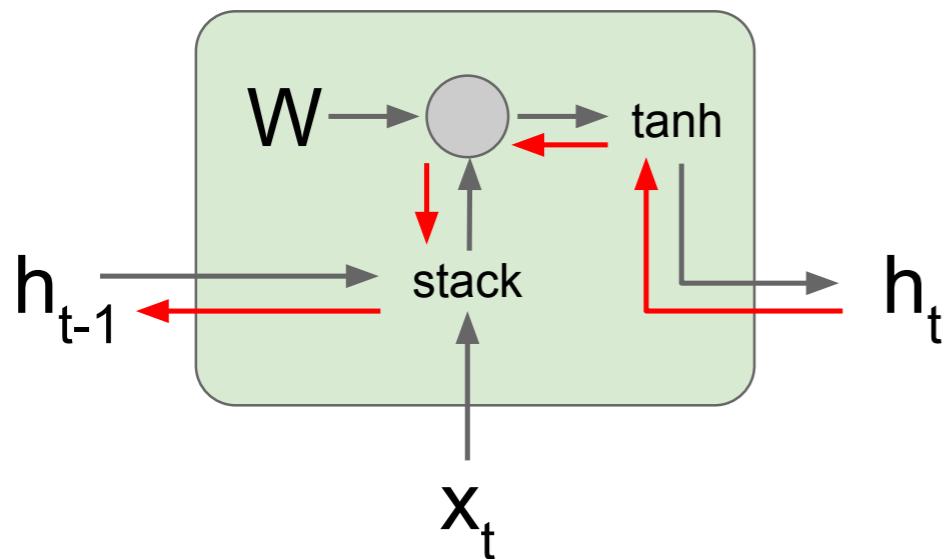


Truncated backprop through time



Life is harder in the RNN world

Backpropagation from h_t to h_{t-1} multiplies by W
(actually W_{hh}^T)



$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

K steps back: W^k

Exponential decay or explosion of gradients

If gradients are exploding: clip to a max value

Fancy regularization: unitary matrices

Smarter architectures!

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

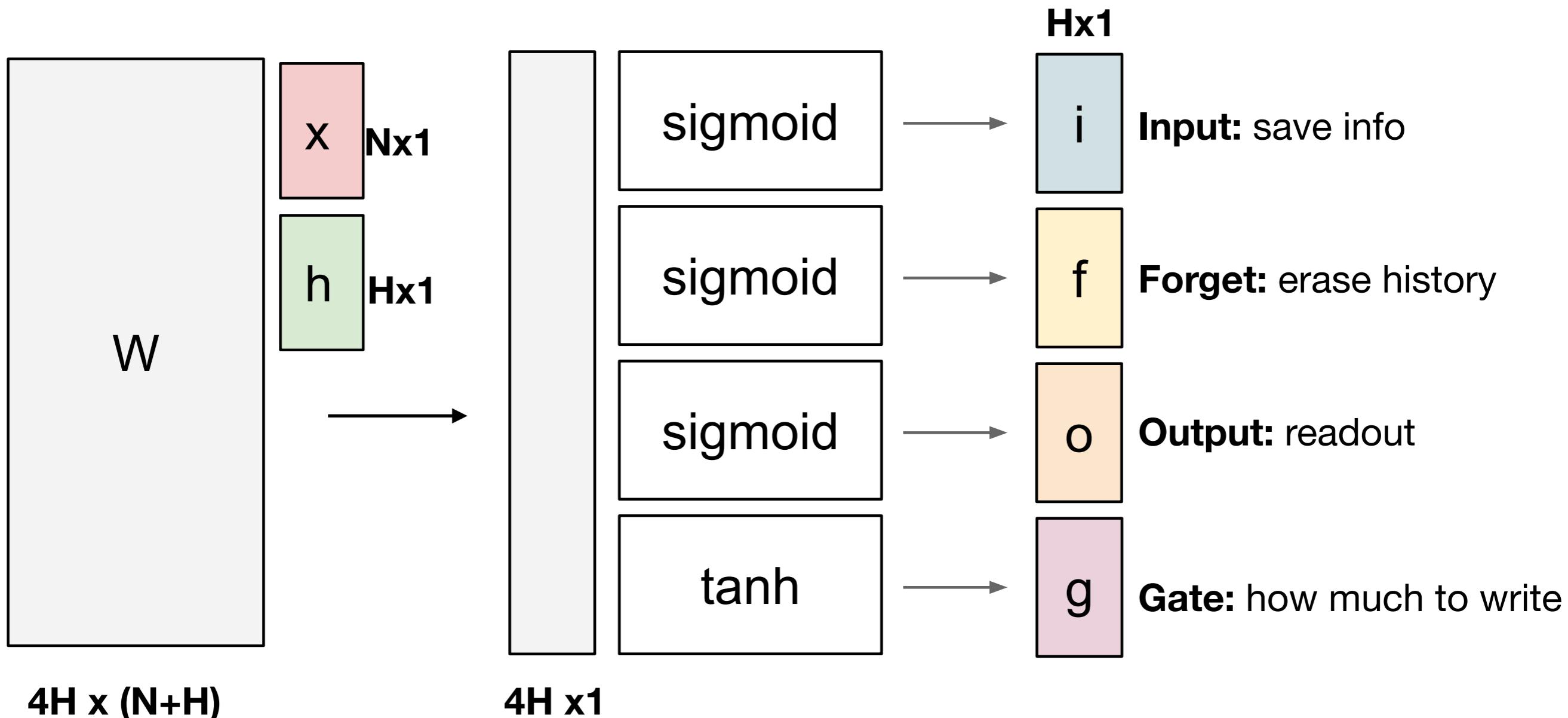
$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Hochschreiter & Schmidhuber, 1997

⊕ Element-wise multiplication

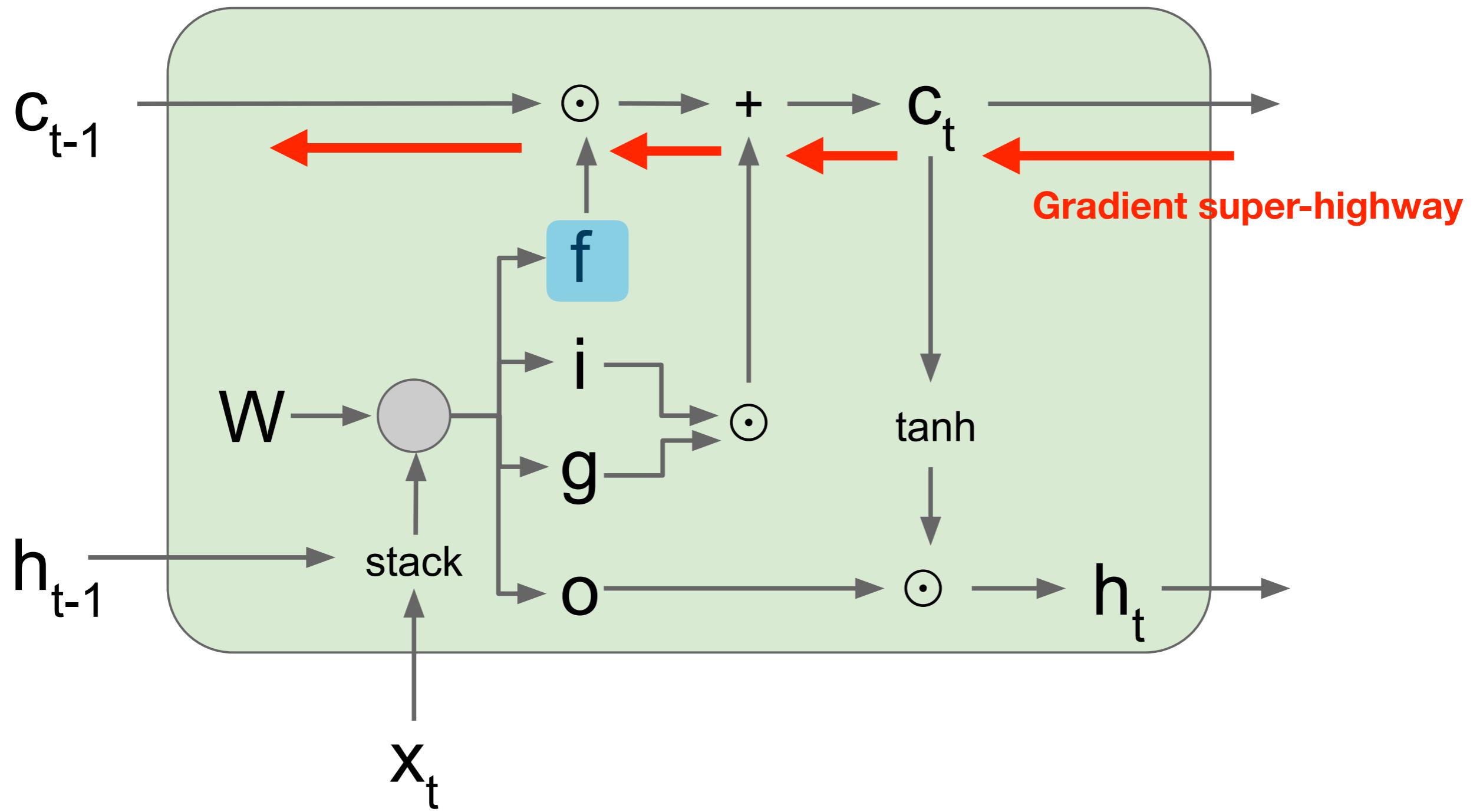
LSTM details



$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

**Gradient for c multiplied (element wise) with f ,
which varies over time and bounded
so much more numerically stable**



Other architectures

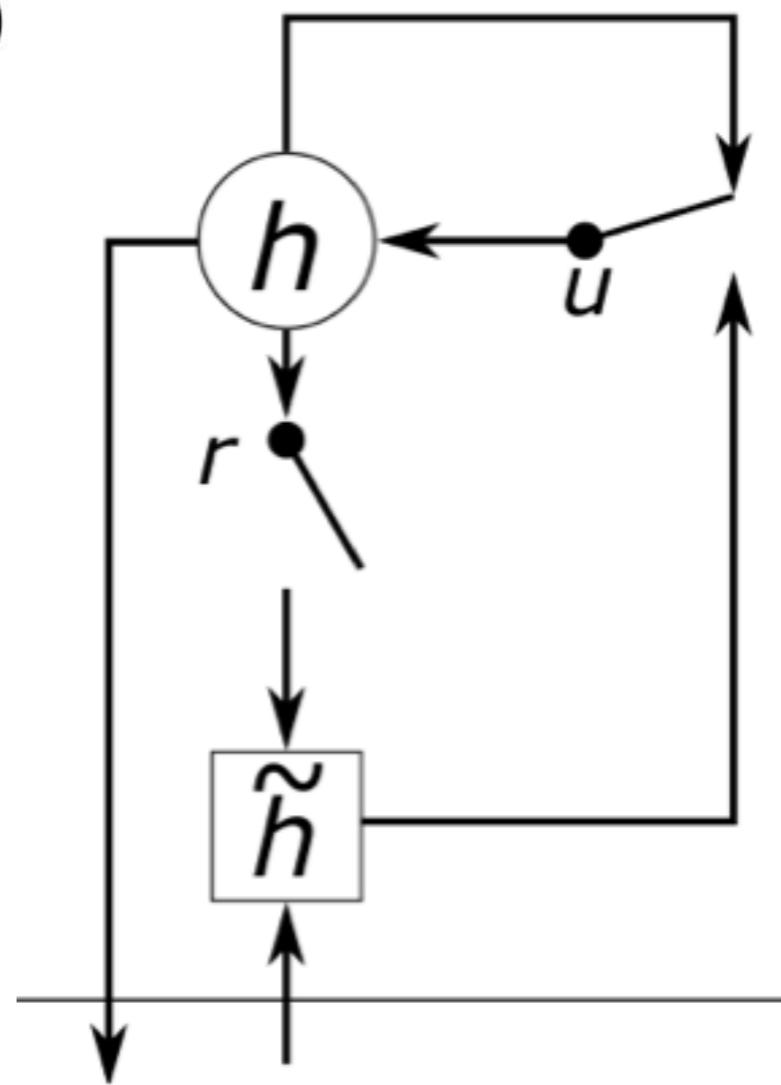
GRU (Cho et al, 2013)

$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t$$



Summary

RNNs provide a **flexible** way for modeling temporal dependencies at scale.

Neural networks (deep learning) enable us to handle

- High-dimensional observations (often 100s-1000s)
- High-dimensional latent variables (often 100s-1000s)
- Large-scale data (often millions or more)
- Neural networks (deep learning) enable us to capture
 - Nonlinear dynamics
 - Long-term dependencies
- Neural networks (deep learning) enable us to work with
 - Intractable probabilistic models
- Neural networks enable us to transfer knowledge across problems
 - Parameter sharing across multiple datasets
 - One neural network with a task indicator input

Summary

Data hungry, and not necessarily easy to train

LSTMs and GRUs used in practice

Good at capturing long-term dependencies but not high frequency structure

Ongoing research: better architectures, theoretical understanding.

Greff et al. *LSTM: A Search Space Odyssey*, 2015.

Goodfellow, Courville, Bengio. *Deep Learning*, 2016.

Goldberg. *Neural Network Methods for Natural Language Processing*, 2017.

K Cho lecture notes: https://github.com/nyu-dl/NLP_DL_Lecture_Note