

# ENM 540: Data-driven modeling and probabilistic scientific computing

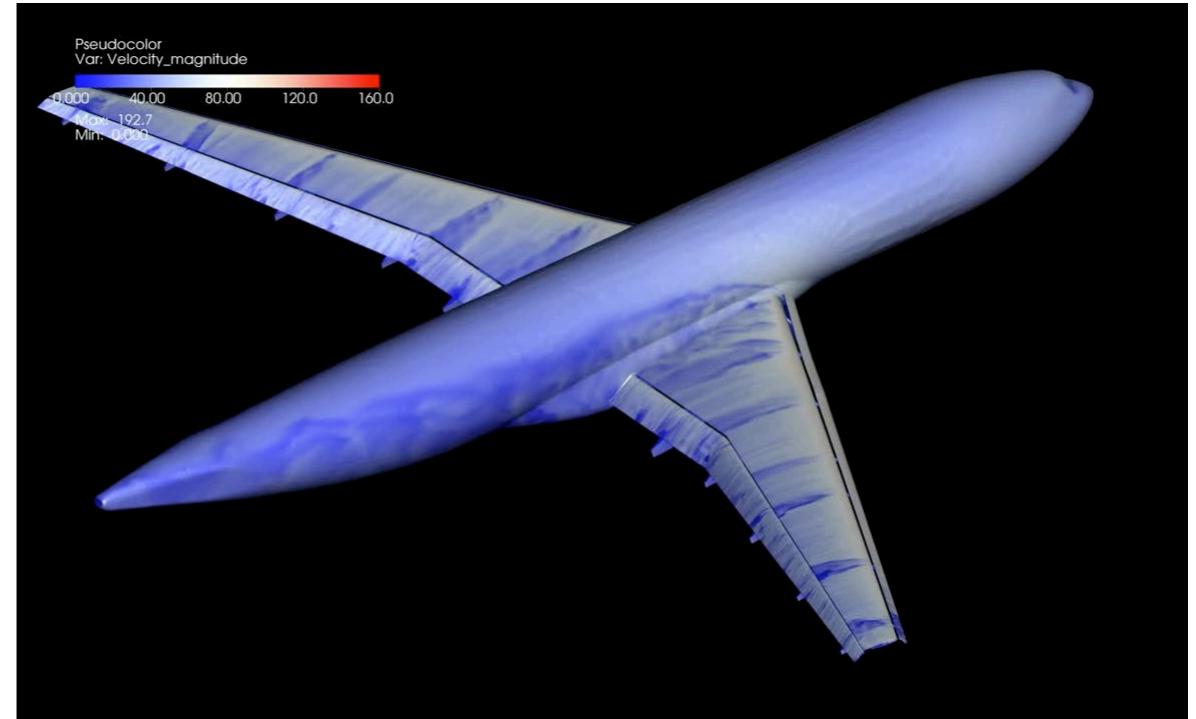
## *Summary*

Paris Perdikaris  
April 19, 2018

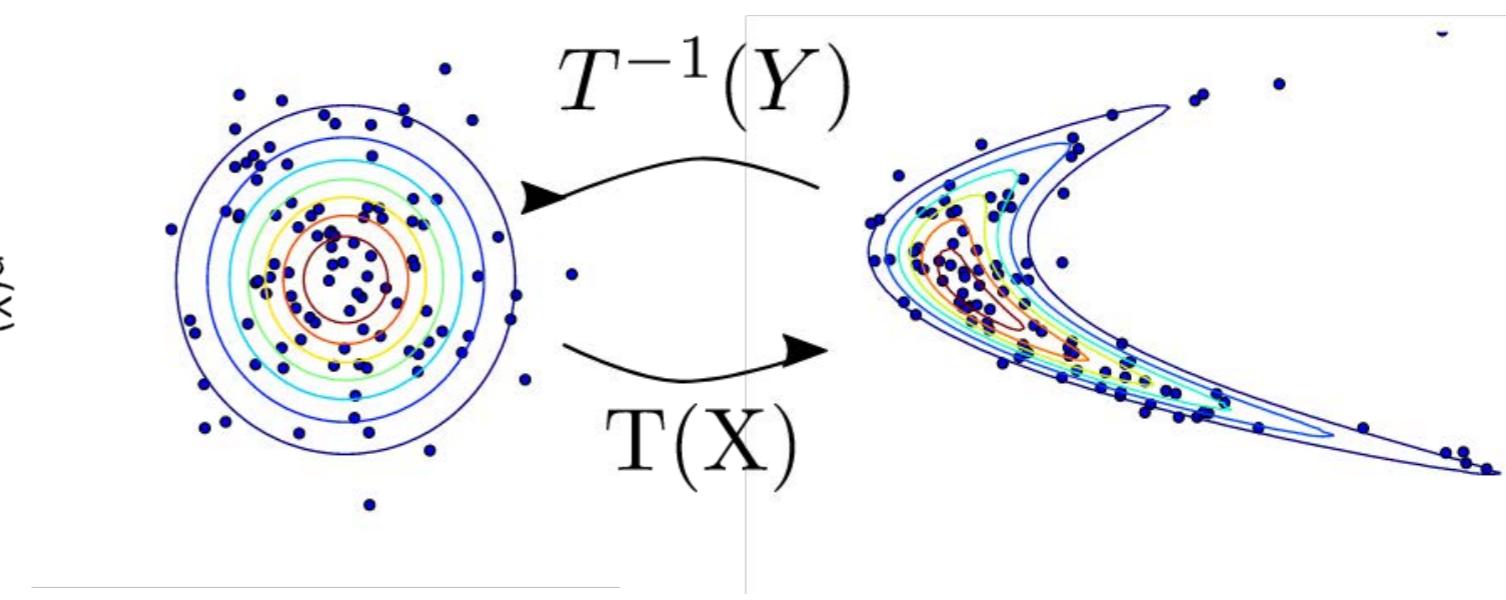
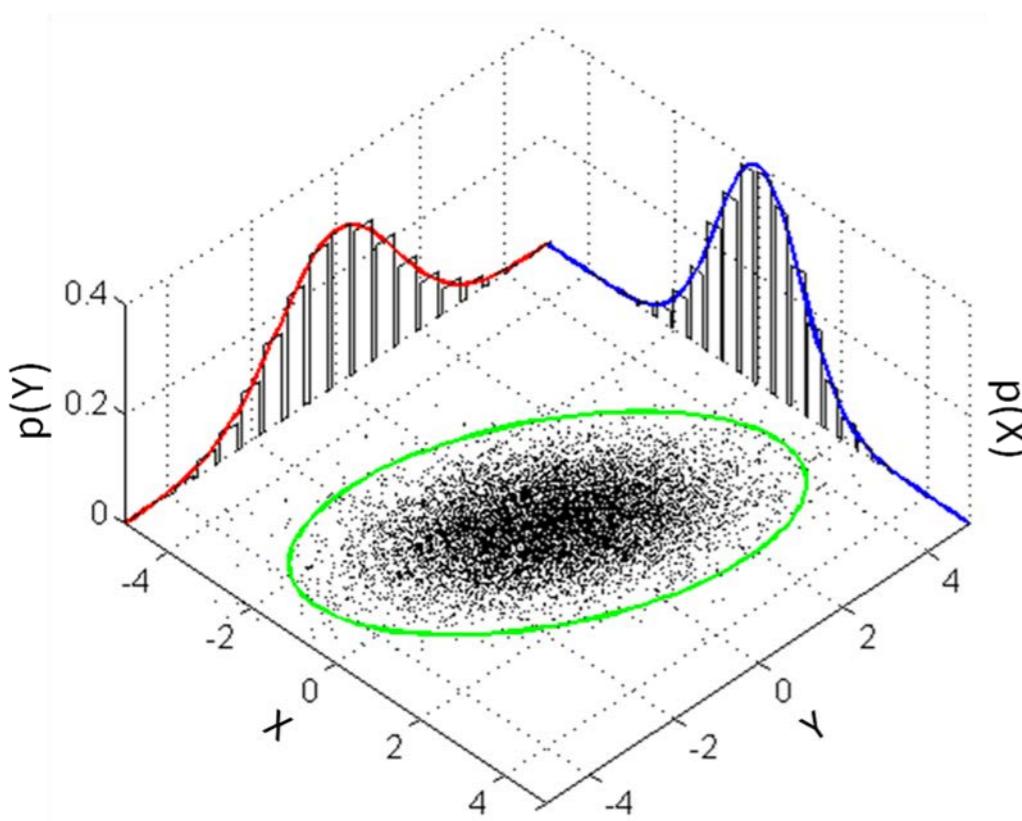
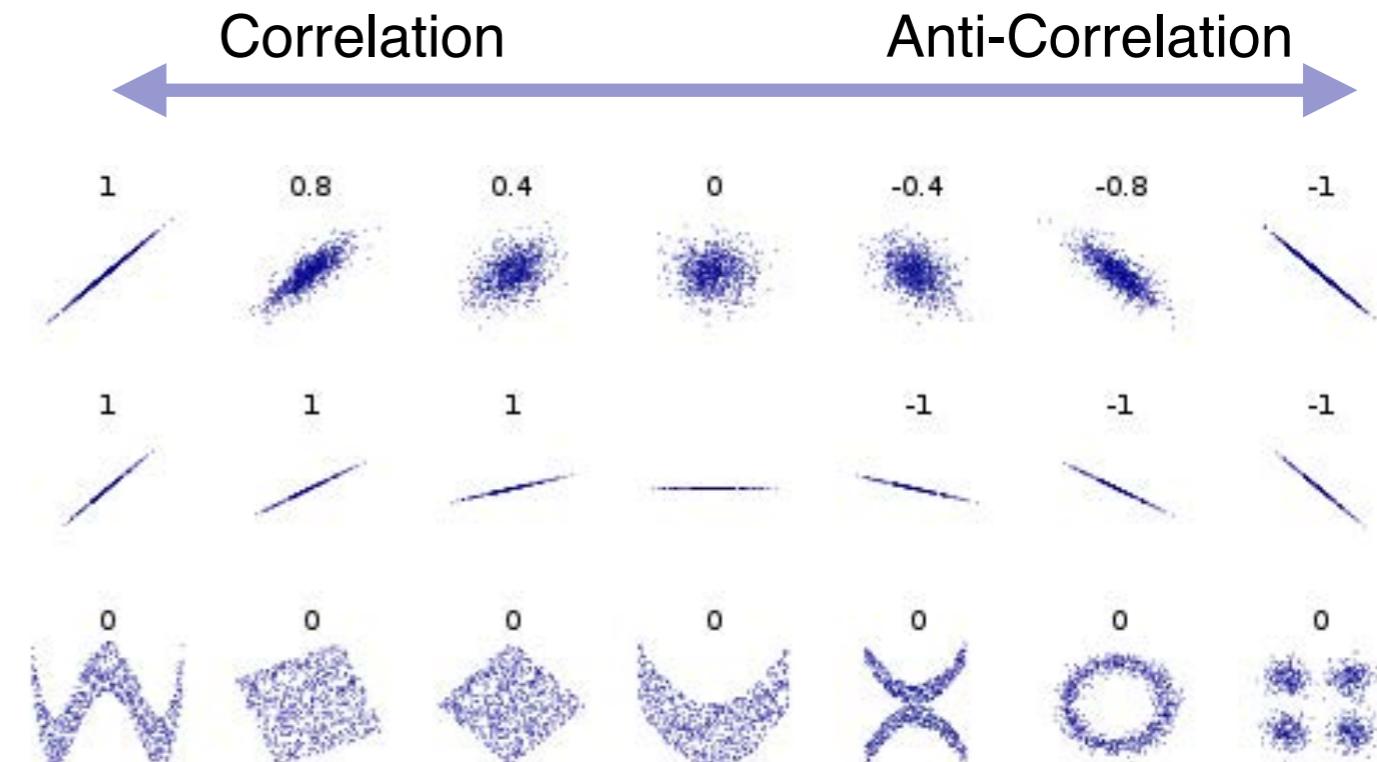
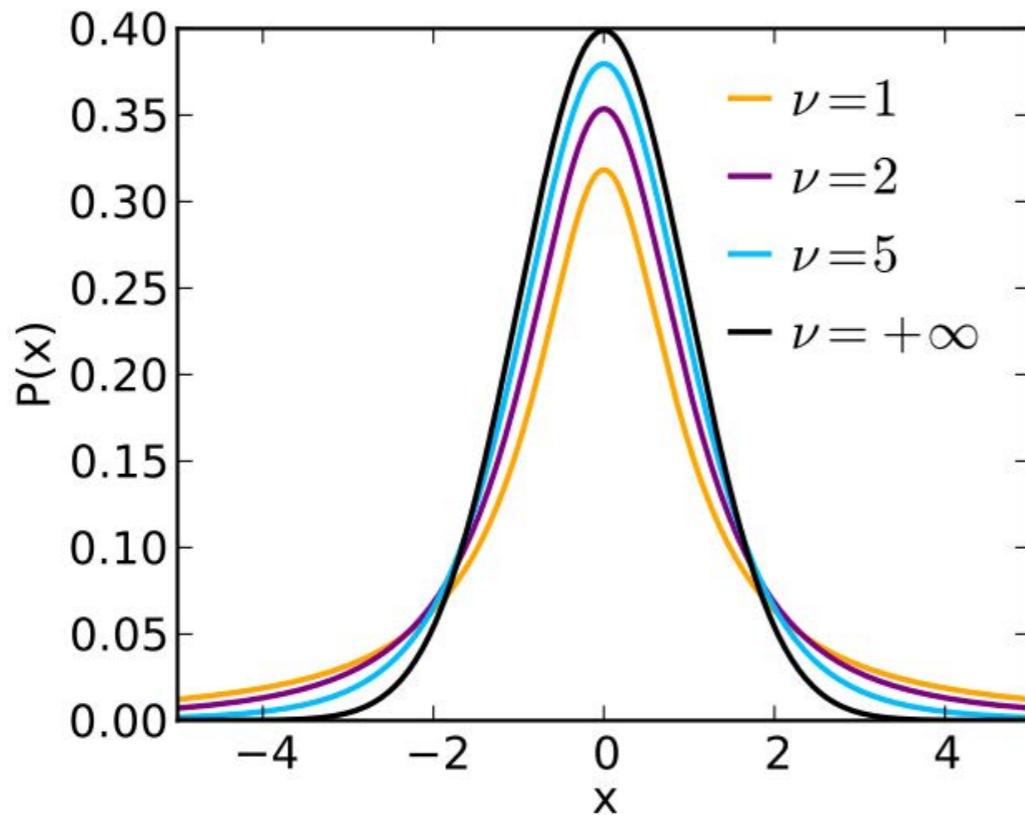


*Machine  
learning*

# Data-driven modeling

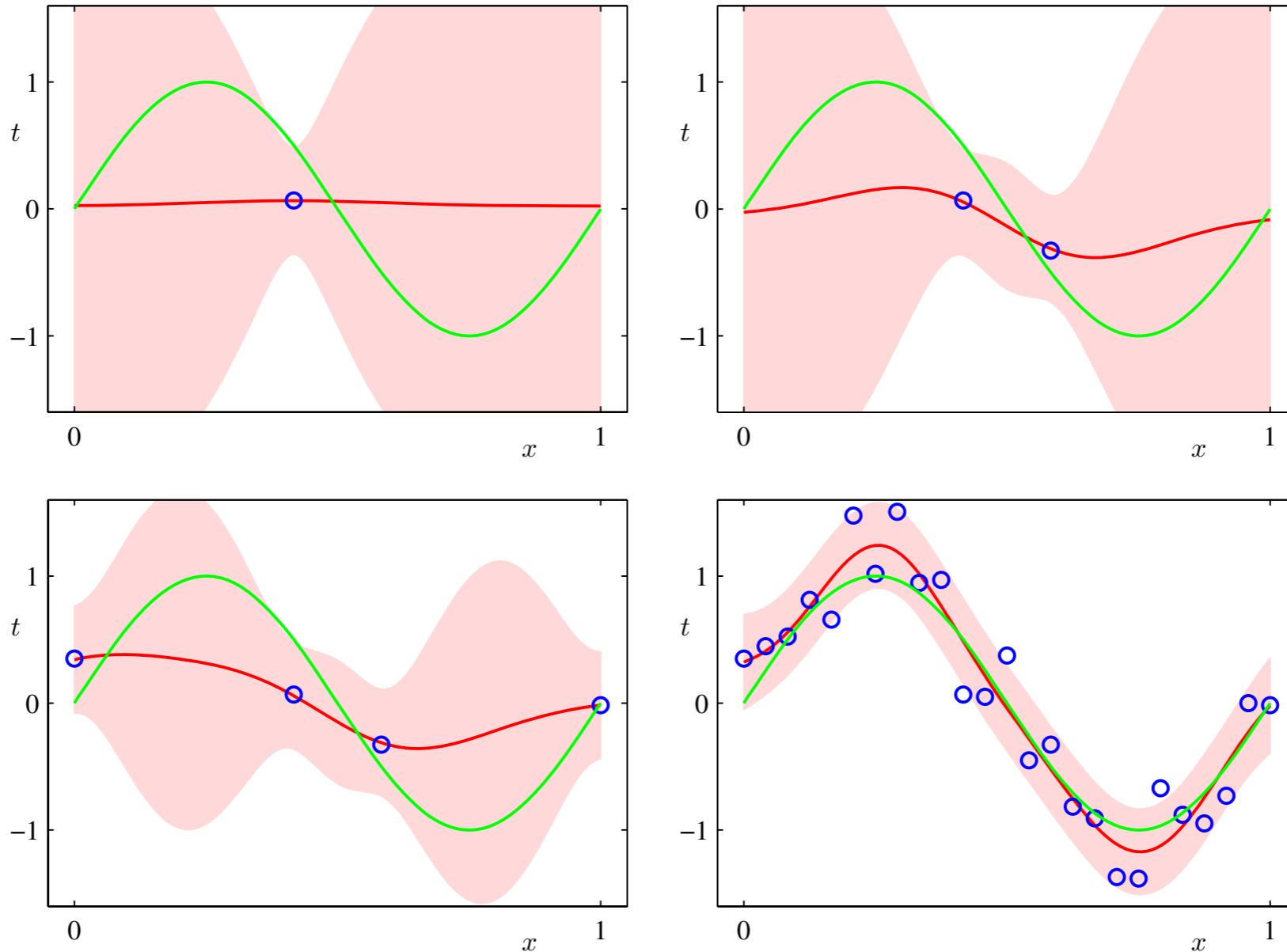


# A probabilistic perspective



# A probabilistic perspective

Classical models can be rigorously formulated from a probabilistic standpoint!



**Figure 3.8** Examples of the predictive distribution (3.58) for a model consisting of 9 Gaussian basis functions of the form (3.4) using the synthetic sinusoidal data set of Section 1.1. See the text for a detailed discussion.

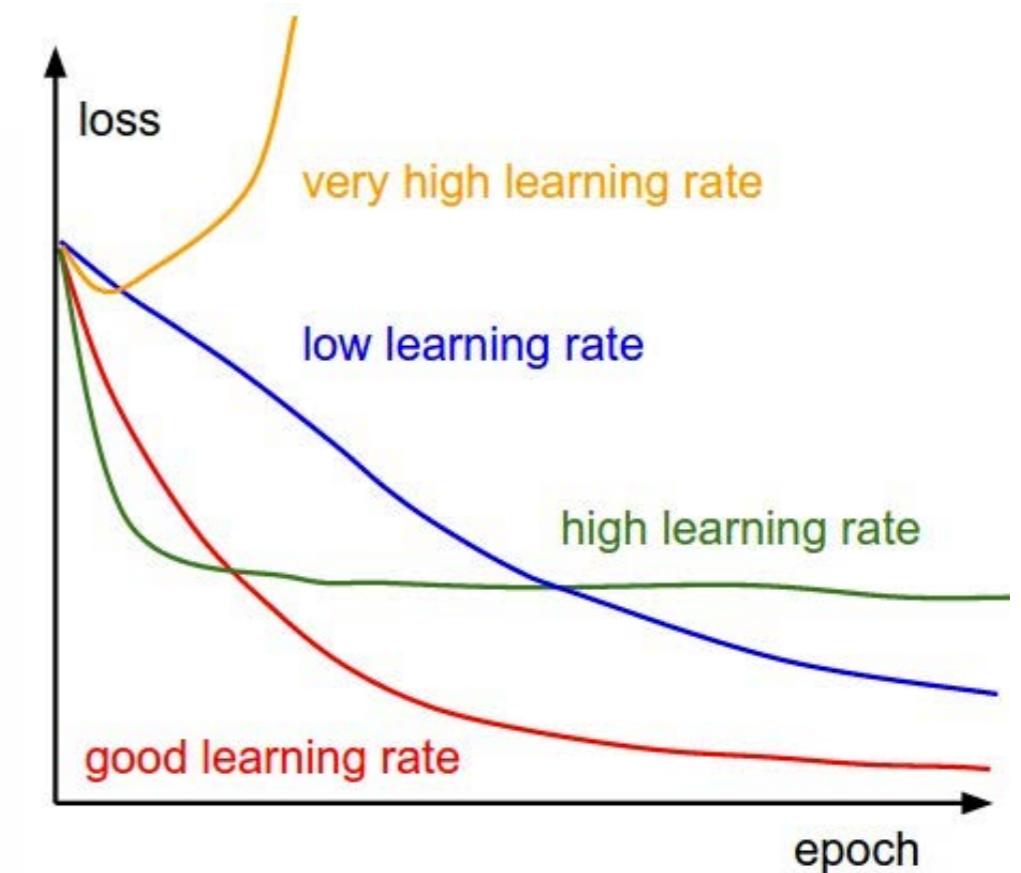
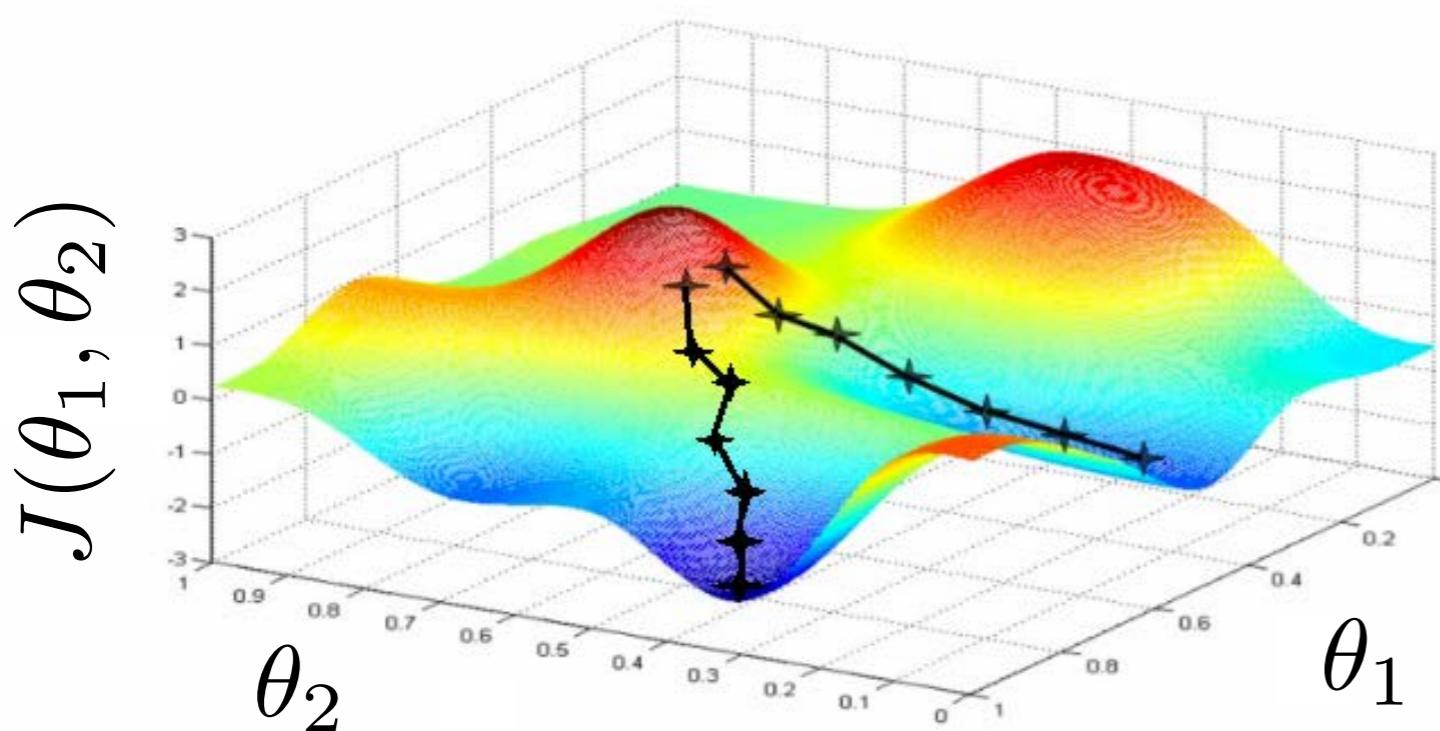
$$\mathbf{y} = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$$

# Optimization

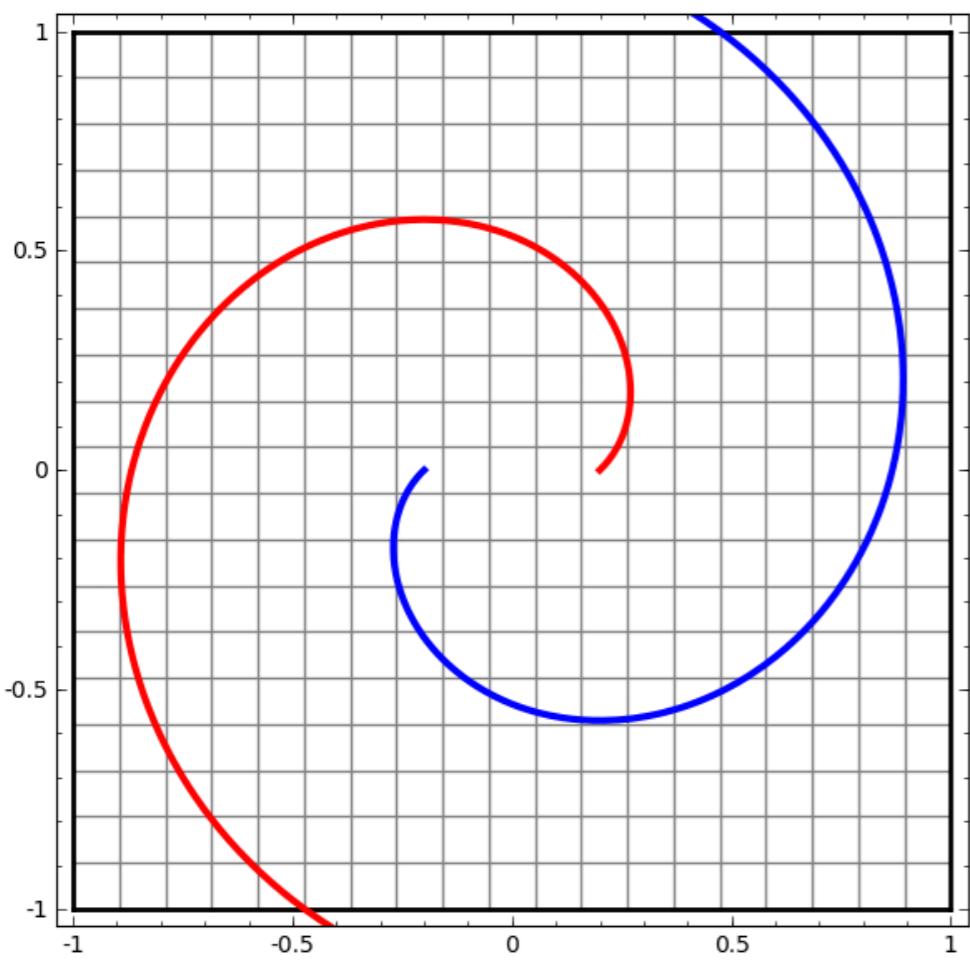
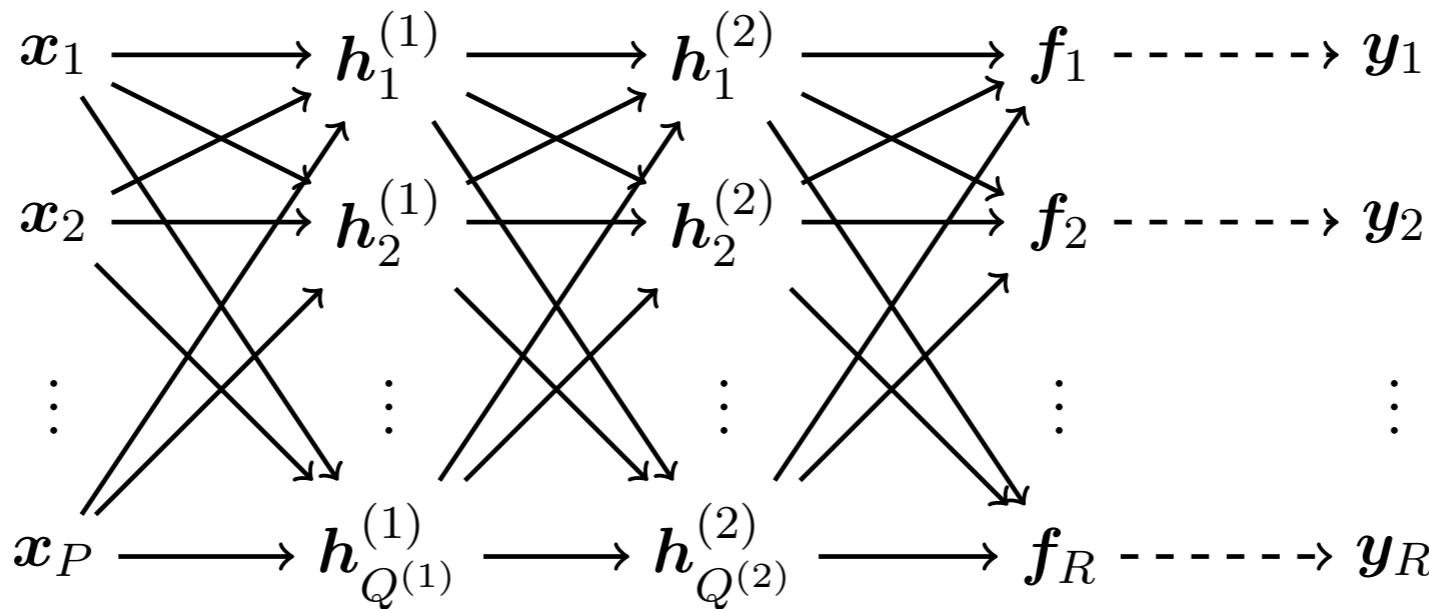
Modern statistics and ML pose great challenges on optimization and define new exciting directions for research!

$$\theta^* = \arg \min_{\theta \in \Theta} J(\theta)$$

$$\theta_{n+1} = \theta_n - \eta \nabla_{\theta} J(\theta)$$

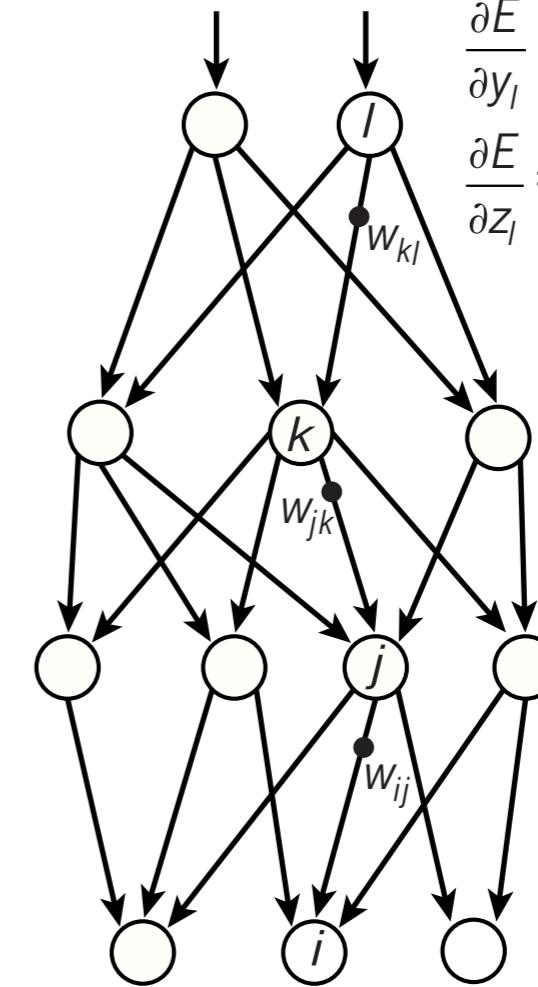


# Neural networks



$$\frac{\partial E}{\partial y_k} = \sum_{l \in \text{out}} w_{kl} \frac{\partial E}{\partial z_l}$$

$$\frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k}$$



$$\begin{aligned}\frac{\partial E}{\partial y_I} &= y_I - t_I \\ \frac{\partial E}{\partial z_I} &= \frac{\partial E}{\partial y_I} \frac{\partial y_I}{\partial z_I}\end{aligned}$$

$$\begin{aligned}\frac{\partial E}{\partial y_j} &= \sum_{k \in H2} w_{jk} \frac{\partial E}{\partial z_k} \\ \frac{\partial E}{\partial z_j} &= \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_j}\end{aligned}$$

# Automatic differentiation

## The chain rule, forward and reverse accumulation [\[edit\]](#)

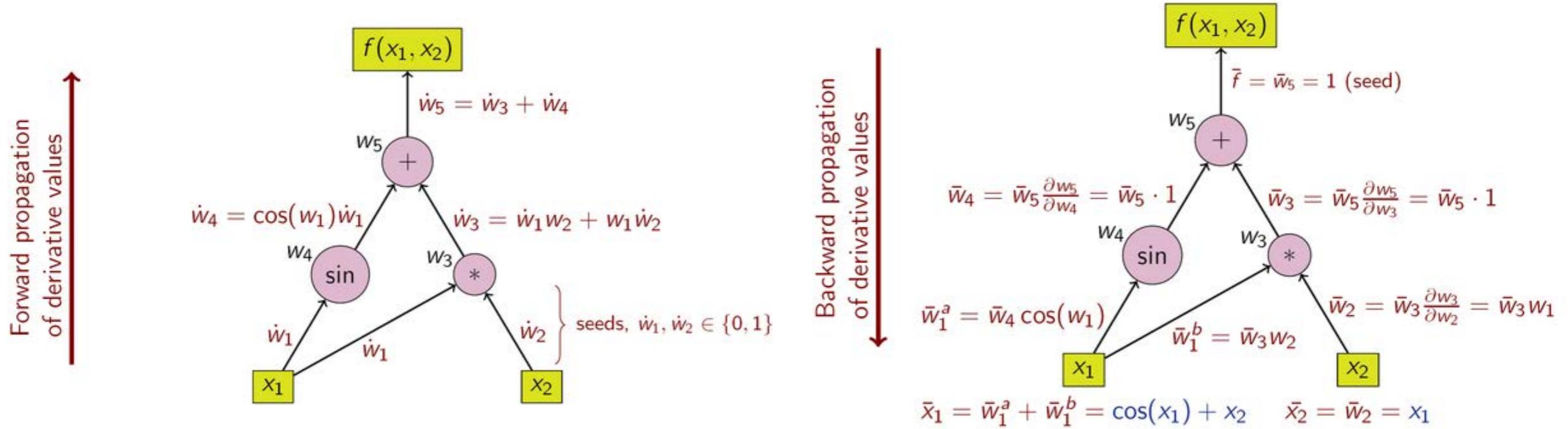
Fundamental to AD is the decomposition of differentials provided by the [chain rule](#). For the simple composition  $y = f(g(h(x))) = f(g(h(w_0))) = f(g(w_1)) = f(w_2) = w_3$  the chain rule gives

$$\frac{dy}{dx} = \frac{dy}{dw_2} \frac{dw_2}{dw_1} \frac{dw_1}{dx}$$

Usually, two distinct modes of AD are presented, **forward accumulation** (or **forward mode**) and **reverse accumulation** (or **reverse mode**). Forward accumulation specifies that one traverses the chain rule from inside to outside (that is, first compute  $dw_1/dx$  and then  $dw_2/dx$  and at last  $dy/dx$ ), while reverse accumulation has the traversal from outside to inside (first compute  $dy/dw_2$  and then  $dy/dw_1$  and at last  $dy/dx$ ). More succinctly,

1. **forward accumulation** computes the recursive relation:  $\frac{dw_i}{dx} = \frac{dw_i}{dw_{i-1}} \frac{dw_{i-1}}{dx}$  with  $w_3 = y$ , and,
2. **reverse accumulation** computes the recursive relation:  $\frac{dy}{dw_i} = \frac{dy}{dw_{i+1}} \frac{dw_{i+1}}{dw_i}$  with  $w_0 = x$ .

**Example:**  $z = f(x_1, x_2) = x_1 x_2 + \sin x_1$



# Practical tips

- Data normalization
- Stochastic optimization
- Network initialization
- Batch normalization
- Dropout

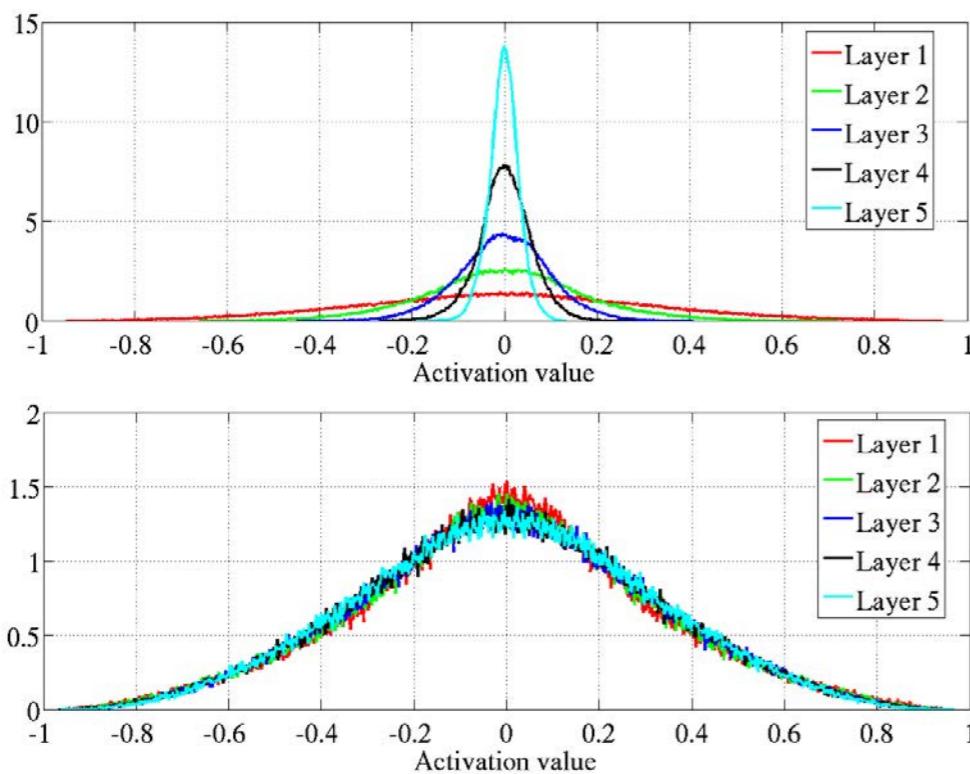
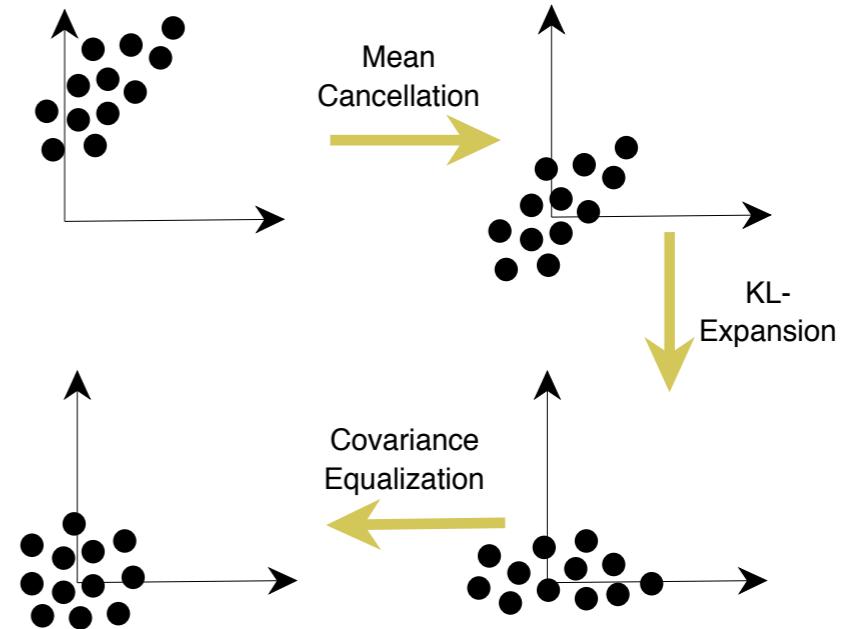


Figure 6: Activation values normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized initialization (bottom). Top: 0-peak increases for higher layers.

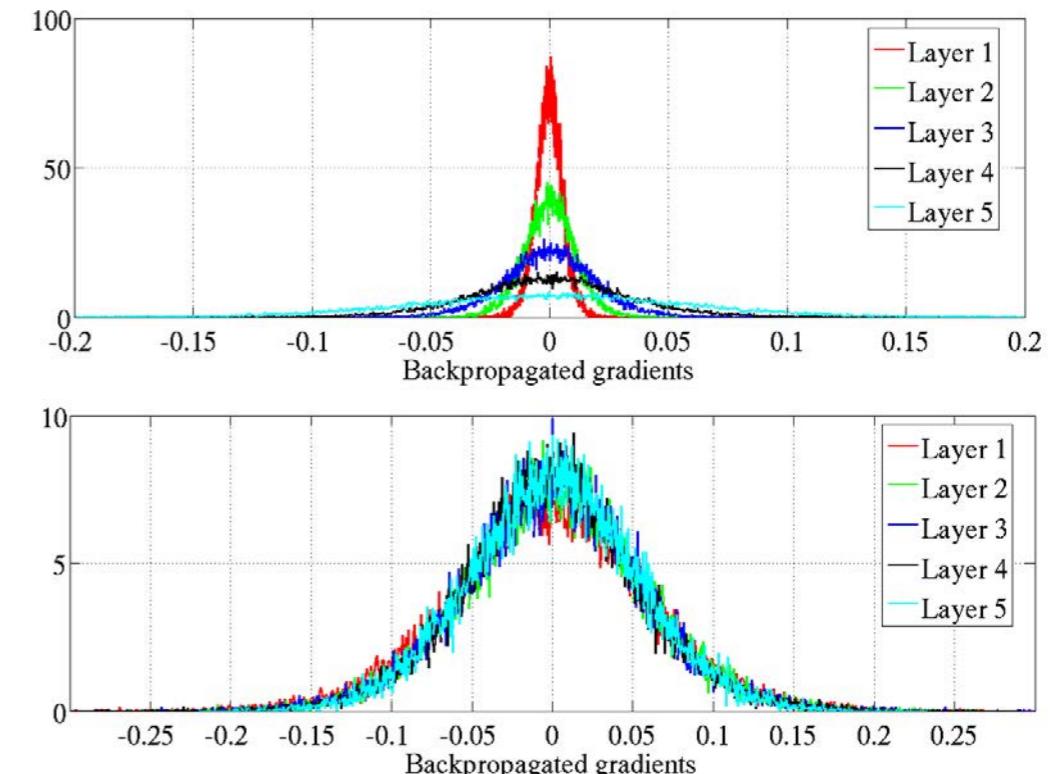


Figure 7: Back-propagated gradients normalized histograms with hyperbolic tangent activation, with standard (top) vs normalized (bottom) initialization. Top: 0-peak decreases for higher layers.

# Applications in scientific computing

**Example:** Constraining the outputs of NNs to satisfy differential equations

$$u_t + uu_x - (0.01/\pi)u_{xx} = 0, \quad x \in [-1, 1], \quad t \in [0, 1], \quad (3)$$
$$u(0, x) = -\sin(\pi x),$$
$$u(t, -1) = u(t, 1) = 0.$$

Let us define  $f(t, x)$  to be given by

$$f := u_t + uu_x - (0.01/\pi)u_{xx},$$

---

```
def u(t, x):
    u = neural_net(tf.concat([t,x],1), weights, biases)
    return u
```

---

Correspondingly, the *physics informed neural network*  $f(t, x)$  takes the form

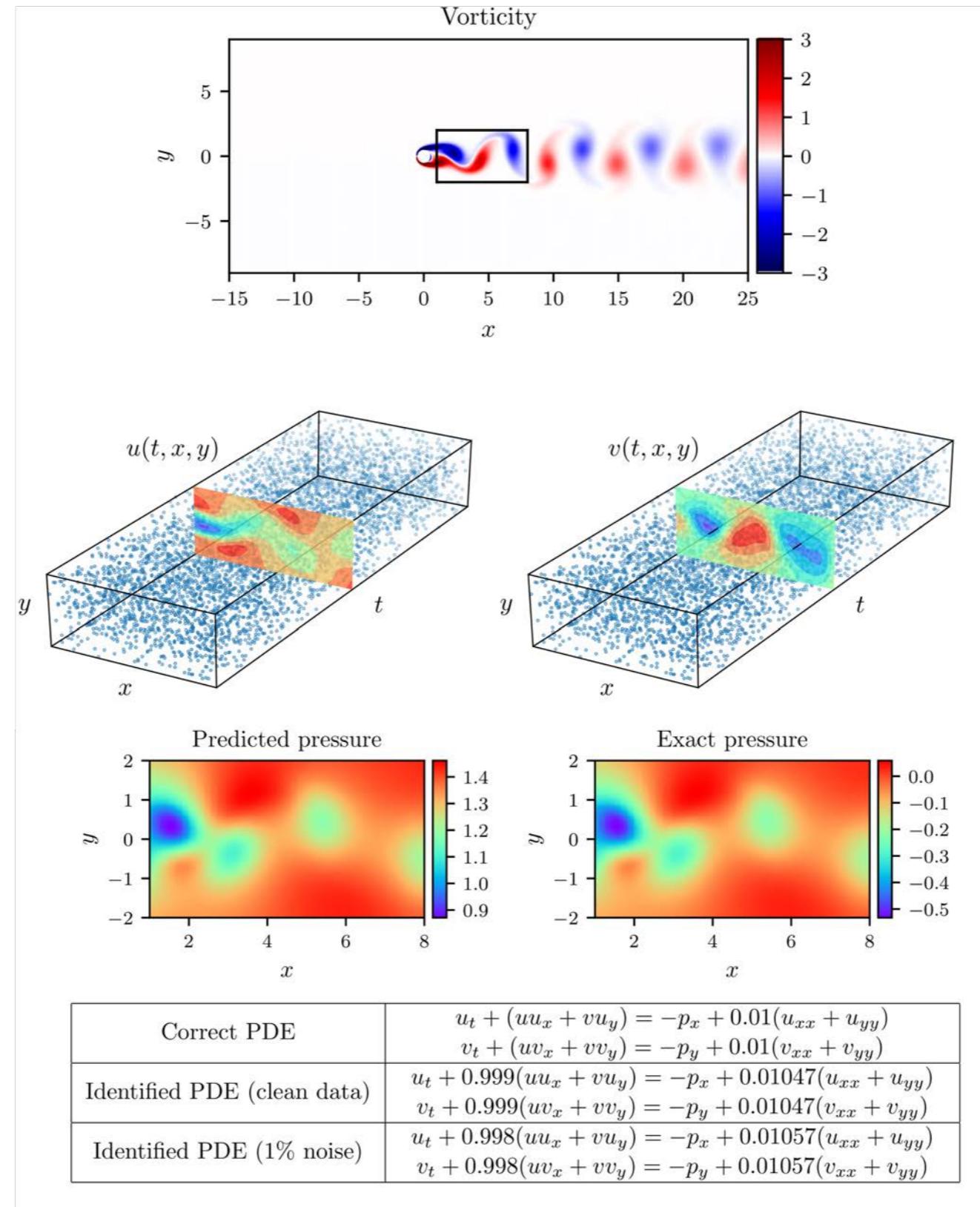
---

```
def f(t, x):
    u = u(t, x)
    u_t = tf.gradients(u, t)[0]
    u_x = tf.gradients(u, x)[0]
    u_xx = tf.gradients(u_x, x)[0]
    f = u_t + u*u_x - (0.01/tf.pi)*u_xx
    return f
```

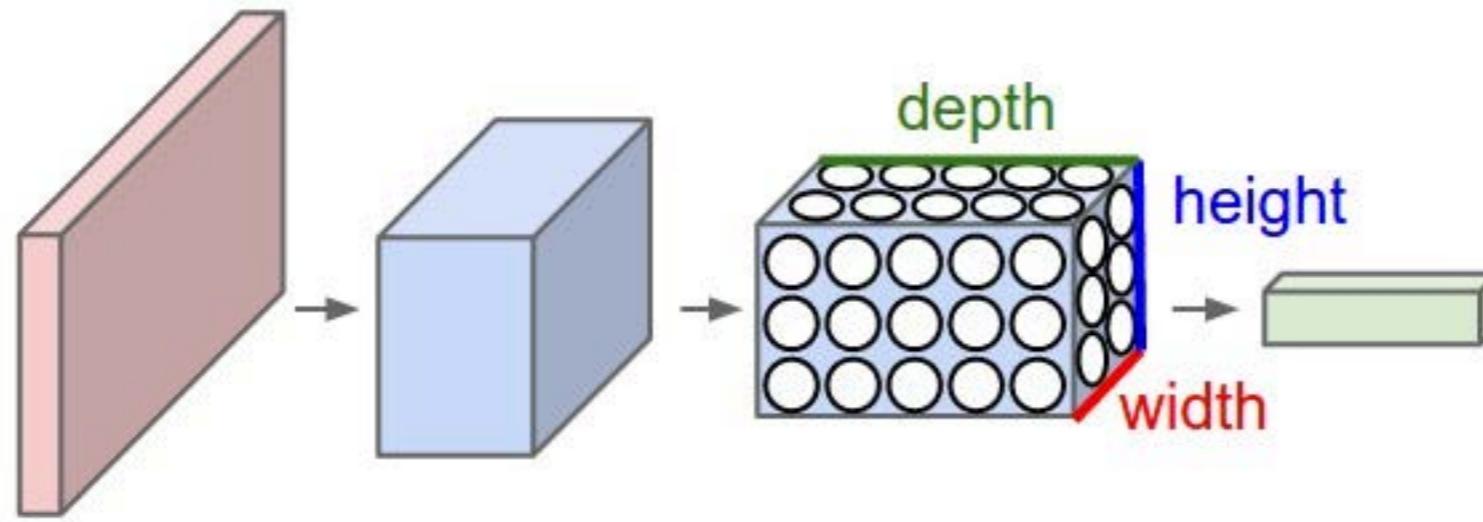
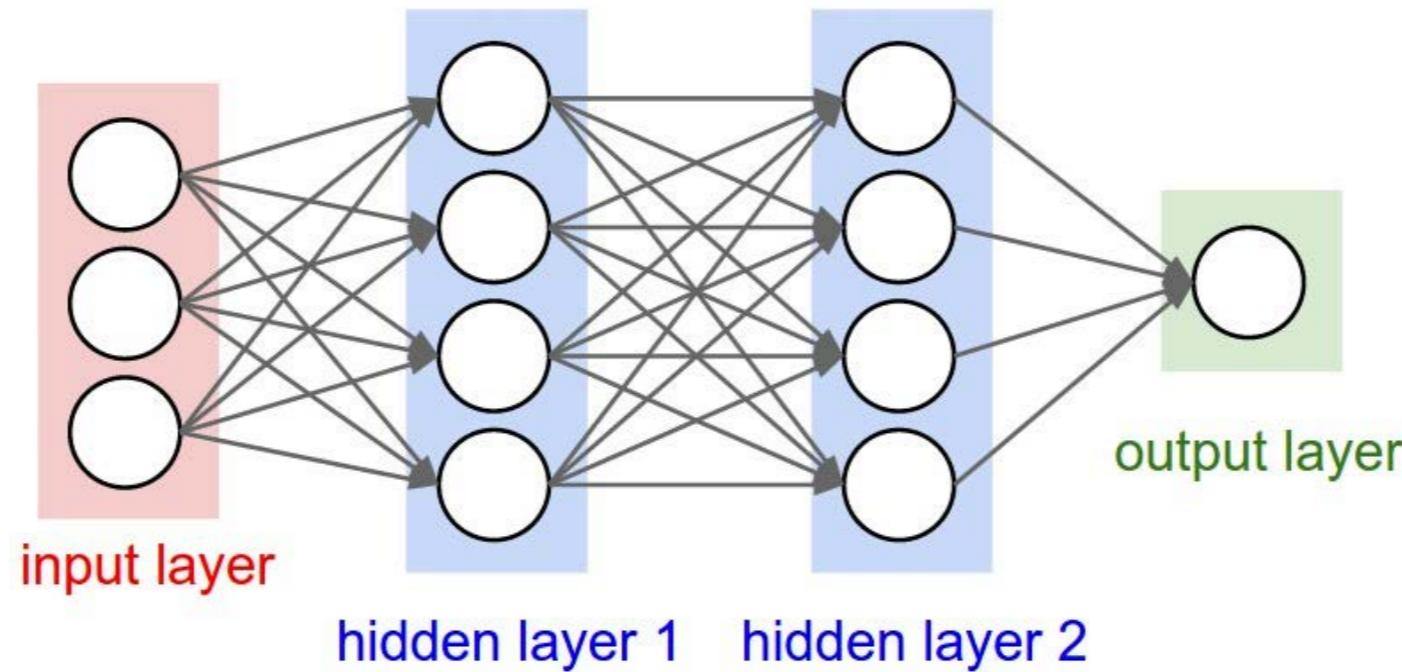
---

# Applications in scientific computing

**Example:** Discovering new physical models and calibrating parameters from data



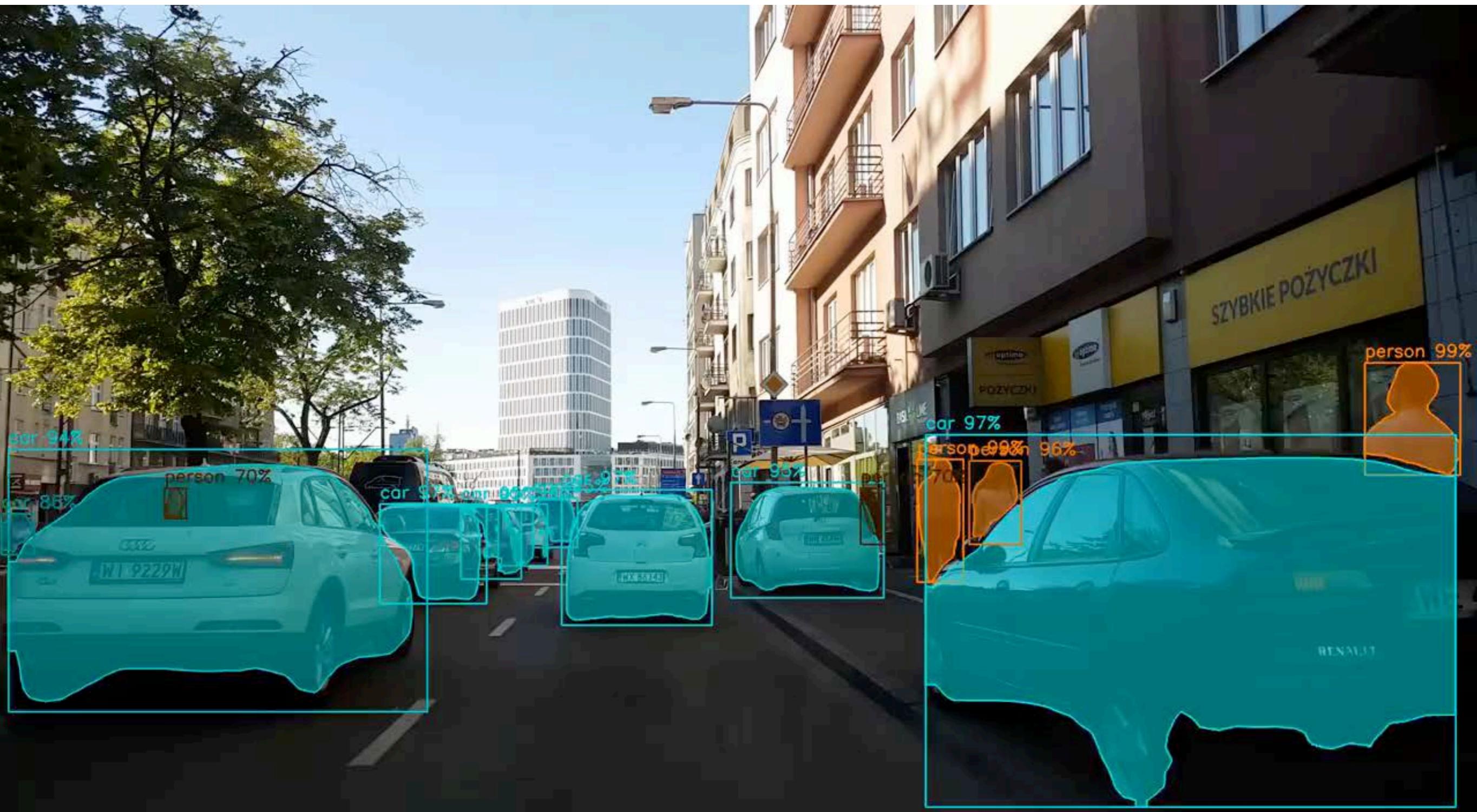
# Convolutional neural networks



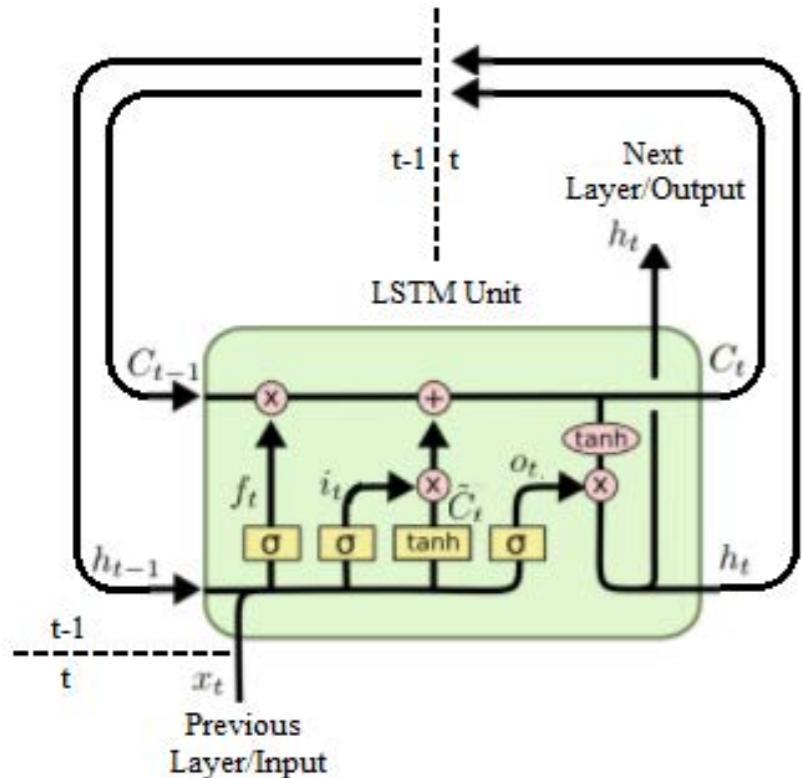
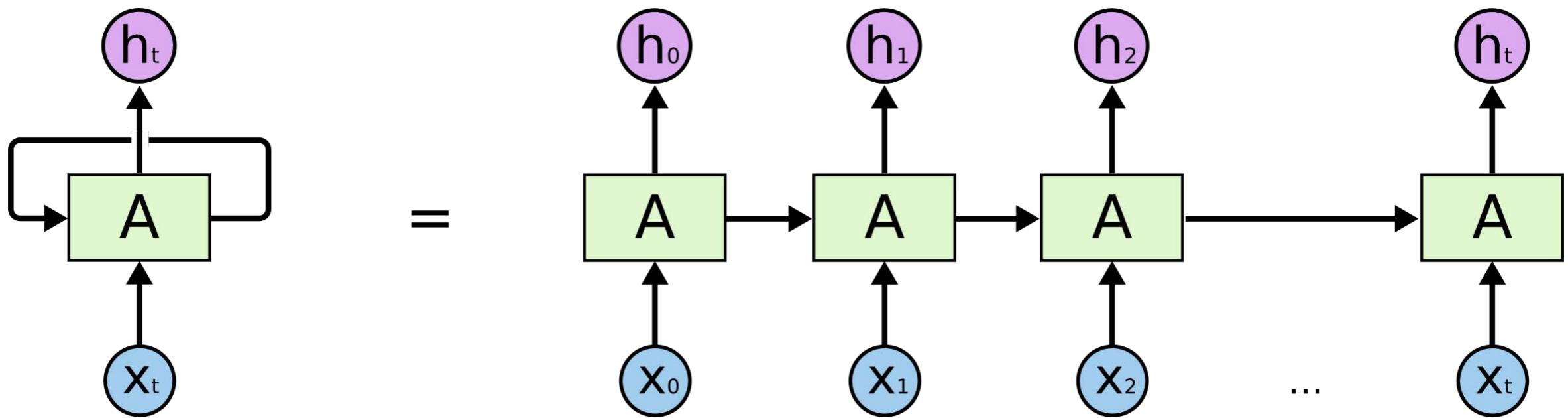
Top: A regular 3-layer Neural Network.

Bottom: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

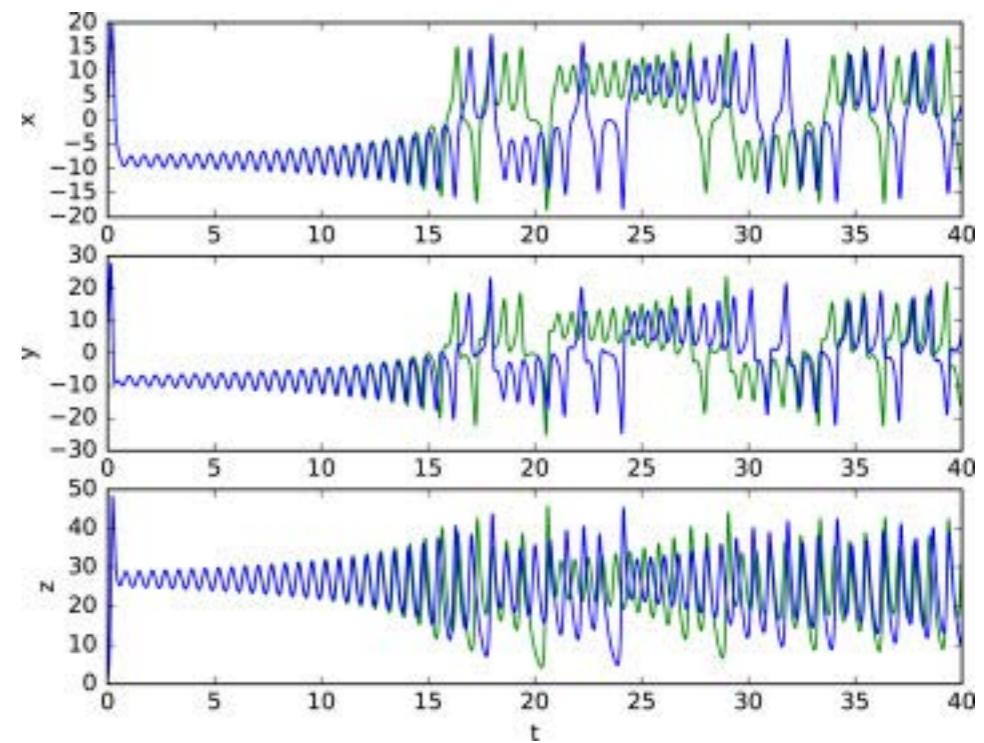
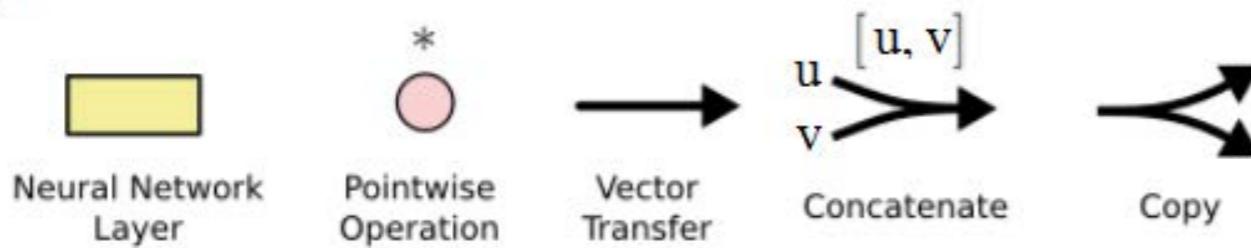
# Convolutional neural networks



# Recurrent neural networks



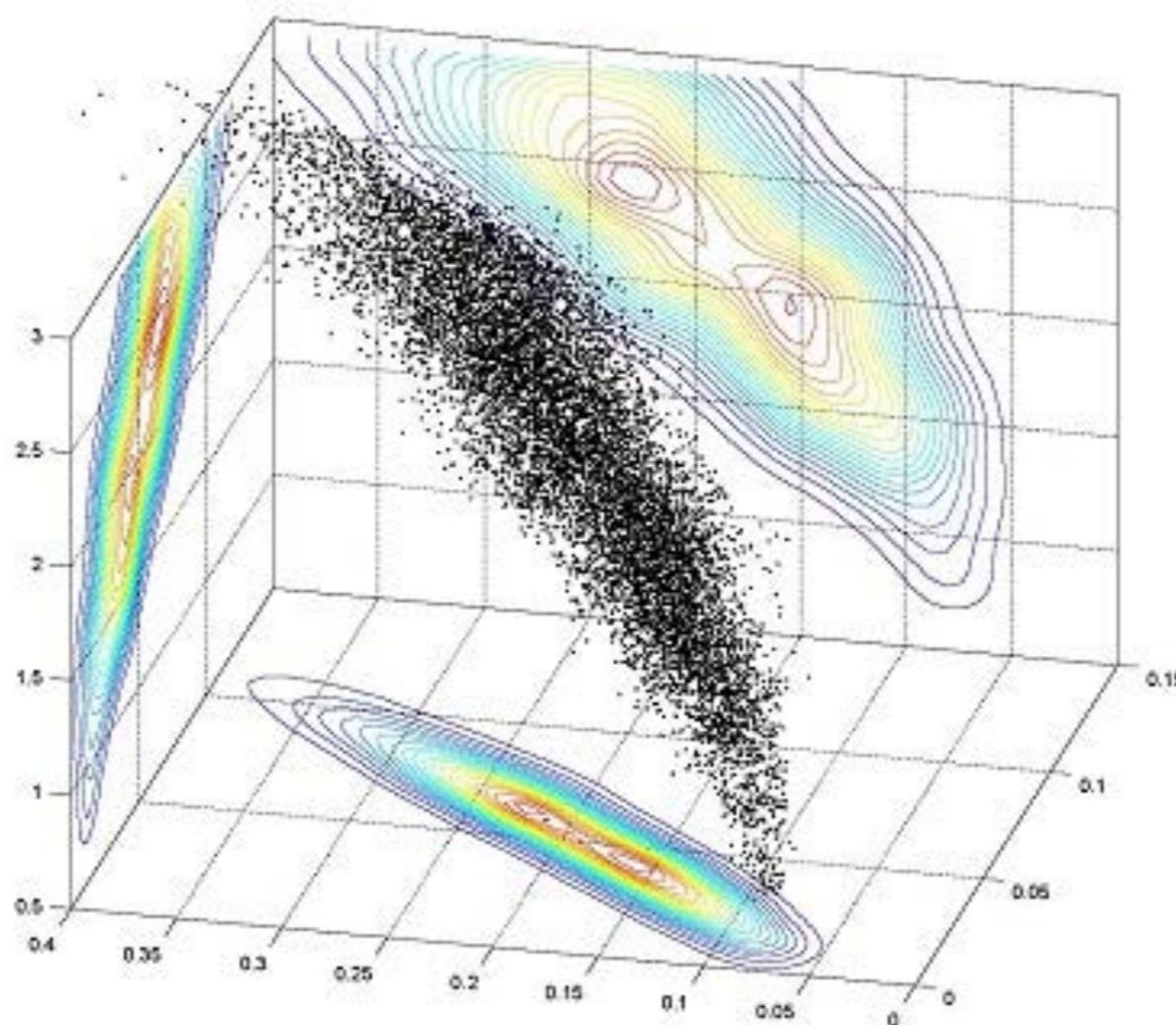
$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \\
 \tilde{C}_t &= \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \\
 C_t &= f_t * C_{t-1} + i_t * \tilde{C}_t \\
 o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\
 h_t &= o_t * \tanh(C_t)
 \end{aligned}$$



## Approximate inference: Sampling methods

$$\mathbb{E}_{x \sim p(x)}[f(x)] = \int f(x)p(x)dx \approx \frac{1}{n} \sum_{i=1}^n f(x_i),$$

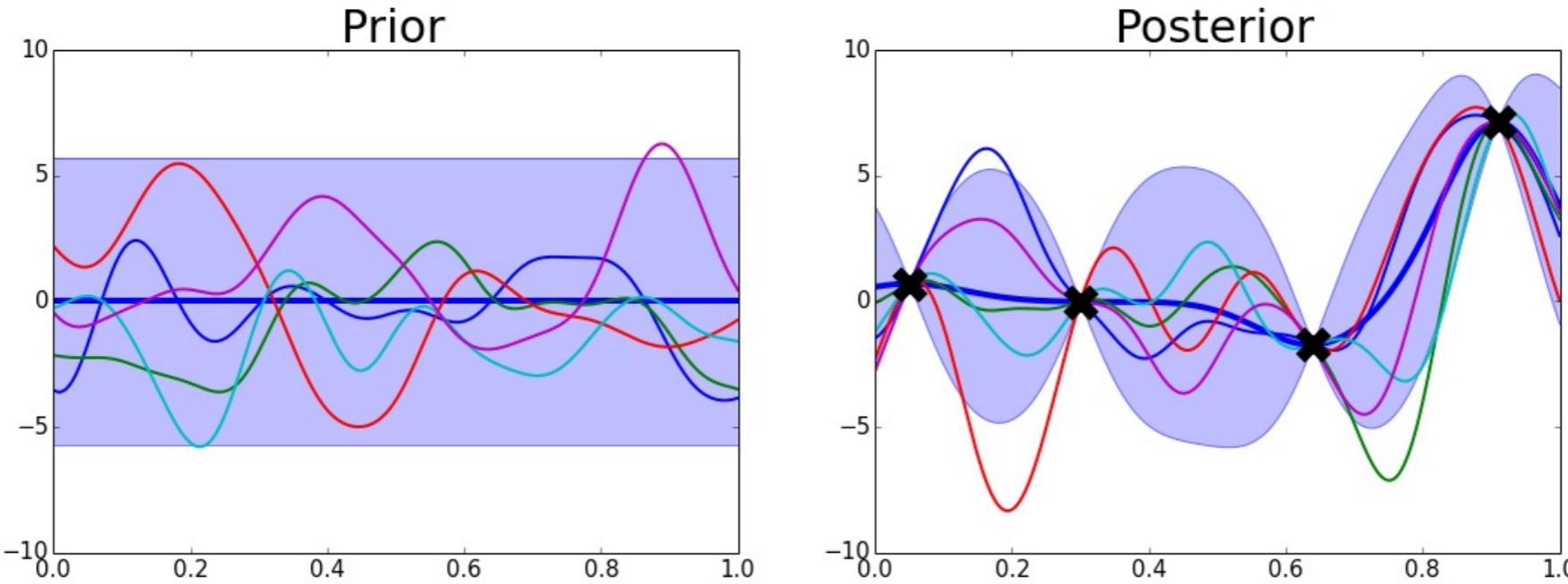
where  $x_i$  are drawn iid from  $p(x)$



# Gaussian process regression

$$y = f(\mathbf{x}) + \epsilon$$

$$f \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}))$$



**Training via maximizing the marginal likelihood**

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\mathbf{K} + \sigma_\epsilon^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log 2\pi$$

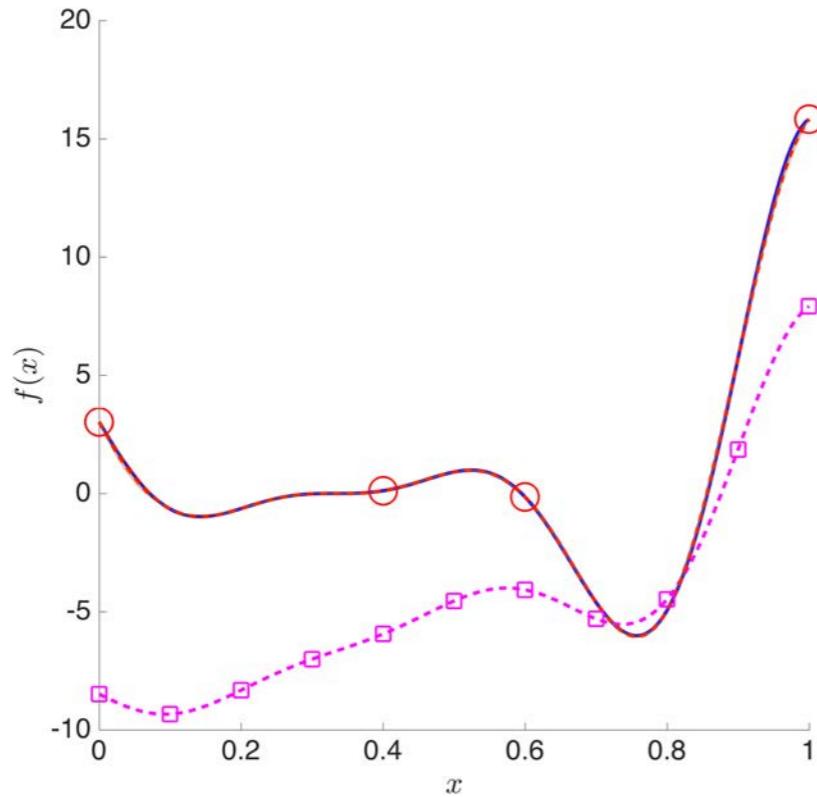
**Prediction via conditioning on available data**

$$p(f_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$\mu_*(\mathbf{x}_*) = \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\sigma_*^2(\mathbf{x}_*) = \mathbf{k}_{**} - \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_{N*},$$

# Multi-fidelity regression



Training:

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_L \\ \mathbf{y}_H \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k_L(\mathbf{x}_L, \mathbf{x}'_L; \theta_L) + \sigma_{\epsilon_L}^2 \mathbf{I} & \rho k_L(\mathbf{x}_L, \mathbf{x}'_H; \theta_L) \\ \rho k_L(\mathbf{x}_H, \mathbf{x}'_L; \theta_L) & \rho^2 k_L(\mathbf{x}_H, \mathbf{x}'_H; \theta_L) + k_H(\mathbf{x}_H, \mathbf{x}'_H; \theta_H) + \sigma_{\epsilon_H}^2 \mathbf{I} \end{bmatrix} \right)$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_L \\ \mathbf{x}_H \end{bmatrix} \quad -\log p(\mathbf{y} | \mathbf{X}, \theta_L, \theta_H, \rho, \sigma_{\epsilon_L}^2, \sigma_{\epsilon_H}^2) = \frac{1}{2} \log |\mathbf{K}| + \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{N_L + N_H}{2} \log 2\pi$$

Prediction:

$$p(f(\mathbf{x}^*) | \mathbf{y}, \mathbf{X}, \mathbf{x}^*) \sim \mathcal{N}(f(\mathbf{x}^*) | \mu(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$\mu(\mathbf{x}^*) = \mathbf{k}(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} \mathbf{y}$$

$$\sigma(\mathbf{x}^*) = \mathbf{k}(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}^*)$$

M.C Kennedy, and A. O'Hagan. *Predicting the output from a complex computer code when fast approximations are available*, 2000.

Demo code: <https://github.com/PredictiveIntelligenceLab/GPTutorial>

Multi-fidelity observations:

$$\mathbf{y}_L = f_L(\mathbf{x}_L) + \boldsymbol{\epsilon}_L$$

$$\mathbf{y}_H = f_H(\mathbf{x}_H) + \boldsymbol{\epsilon}_H$$

Probabilistic model:

$$f_H(\mathbf{x}) = \rho f_L(\mathbf{x}) + \delta(\mathbf{x})$$

$$f_L(\mathbf{x}) \sim \mathcal{GP}(0, k_L(\mathbf{x}, \mathbf{x}'; \theta_L))$$

$$\delta(\mathbf{x}) \sim \mathcal{GP}(0, k_H(\mathbf{x}, \mathbf{x}'; \theta_H))$$

$$\boldsymbol{\epsilon}_L \sim \mathcal{N}(0, \sigma_{\epsilon_L}^2 \mathbf{I})$$

$$\boldsymbol{\epsilon}_H \sim \mathcal{N}(0, \sigma_{\epsilon_H}^2 \mathbf{I})$$

# Bayesian optimization

**Goal:** Estimate the global minimum of a function:  $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$  (potentially intractable)

**Setup:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  using a GP surrogate:  $y = f(\mathbf{x}) + \epsilon$ ,  $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

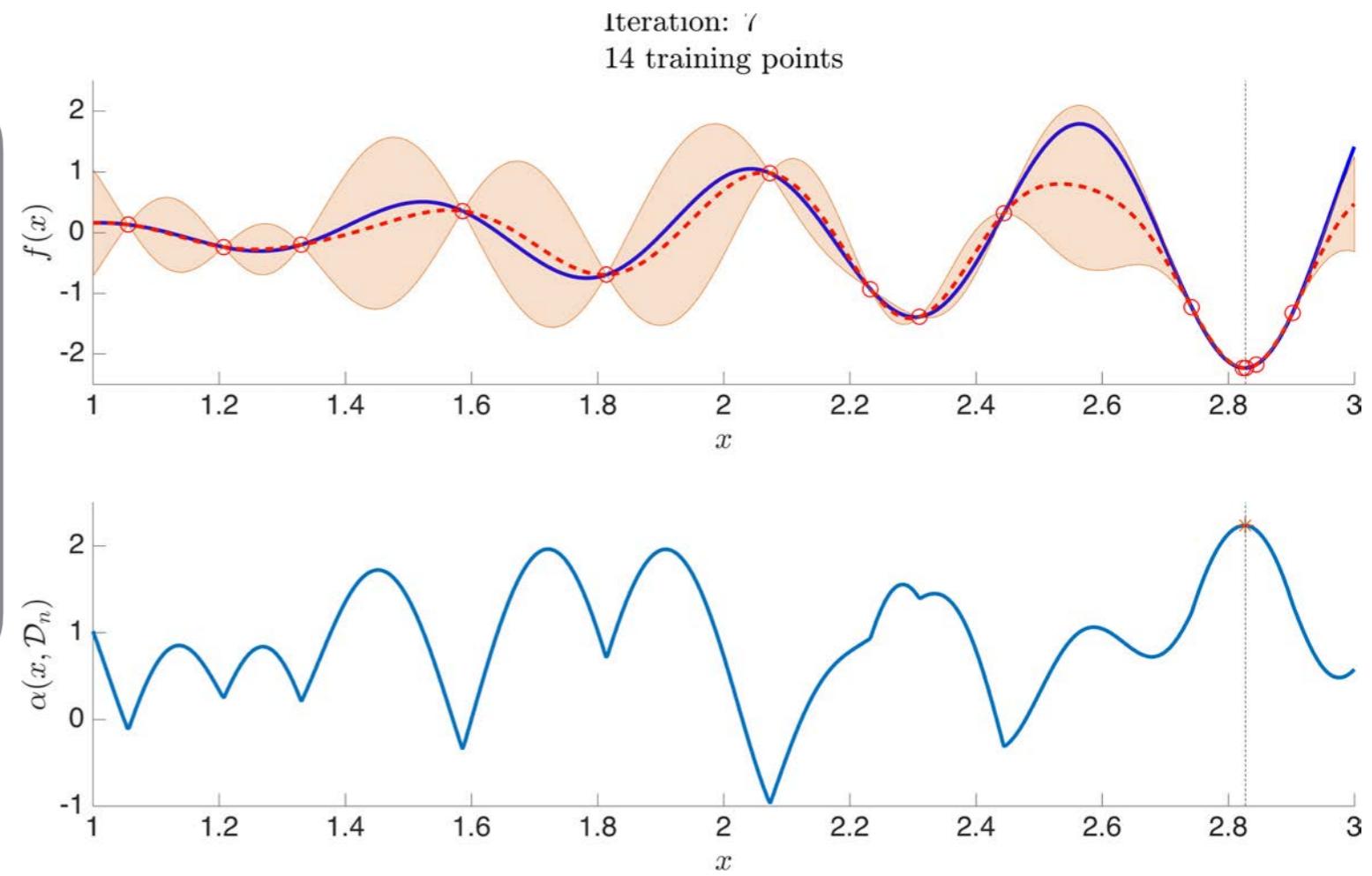
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

## Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower super-quintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

# A new paradigm in scientific computing

$$\mathcal{L}_x u(x) = f(x)$$

$$u(x) \sim \mathcal{GP}(0, g(x, x'; \theta)) \longrightarrow f(x) \sim \mathcal{GP}(0, k(x, x'; \theta))$$

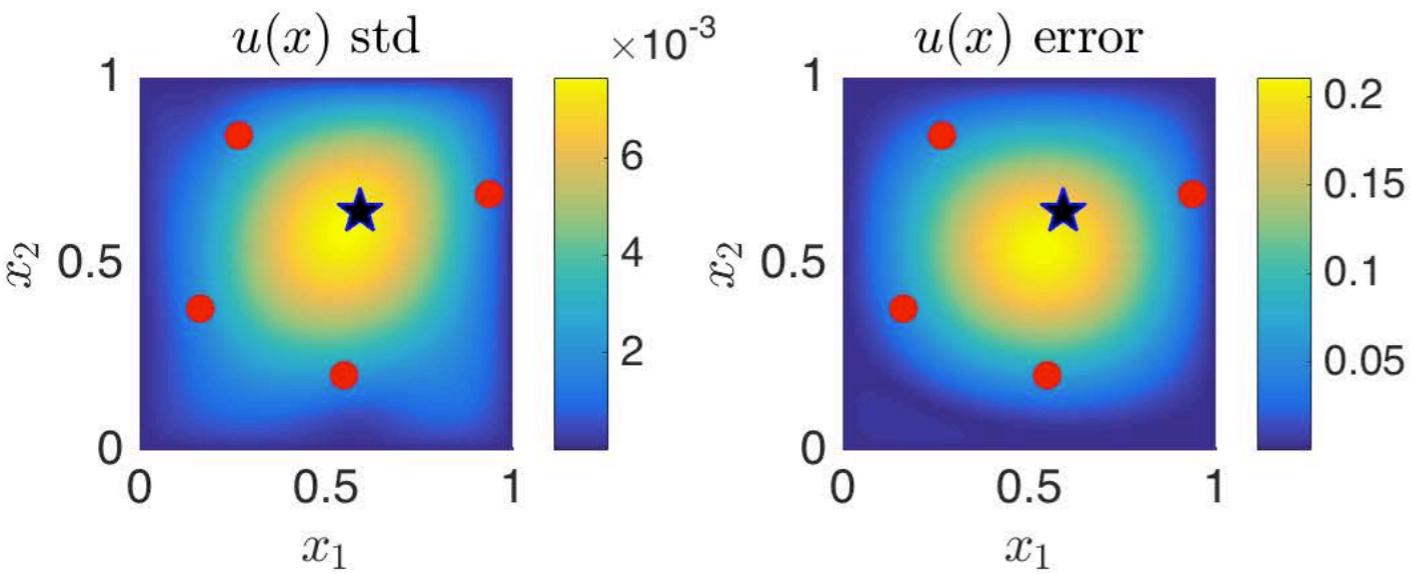
$$k(x, x'; \theta) = \mathcal{L}_x \mathcal{L}_{x'} g(x, x'; \theta)$$

$$-\log p(\mathbf{y}|\phi, \theta, \sigma_{n_u}^2, \sigma_{n_f}^2) = \frac{1}{2} \log |\mathbf{K}| + \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} + \frac{N}{2} \log 2\pi,$$

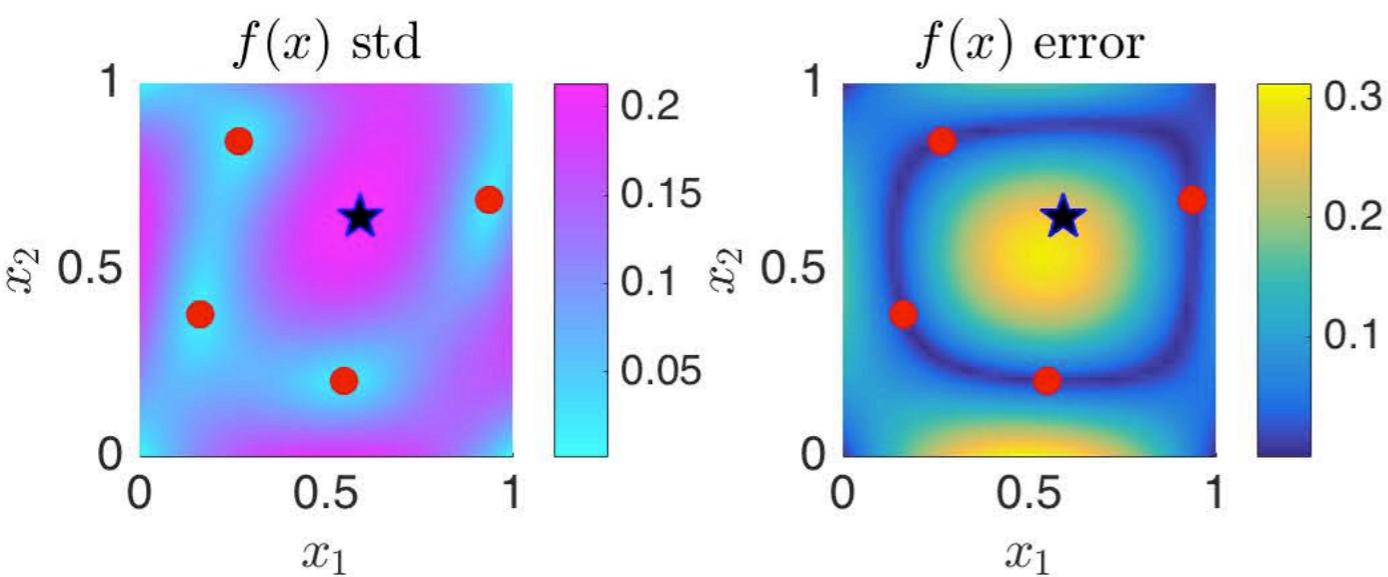
where  $\mathbf{y} = \begin{bmatrix} \mathbf{y}_u \\ \mathbf{y}_f \end{bmatrix}$ ,  $p(\mathbf{y}|\phi, \theta, \sigma_{n_u}^2, \sigma_{n_f}^2) = \mathcal{N}(\mathbf{0}, \mathbf{K})$ , and  $\mathbf{K}$  is given by

$$\mathbf{K} = \begin{bmatrix} k_{uu}(\mathbf{X}_u, \mathbf{X}_u; \theta) + \sigma_{n_u}^2 \mathbf{I}_{n_u} & k_{uf}(\mathbf{X}_u, \mathbf{X}_f; \theta, \phi) \\ k_{fu}(\mathbf{X}_f, \mathbf{X}_u; \theta, \phi) & k_{ff}(\mathbf{X}_f, \mathbf{X}_f; \theta, \phi) + \sigma_{n_f}^2 \mathbf{I}_{n_f} \end{bmatrix}.$$

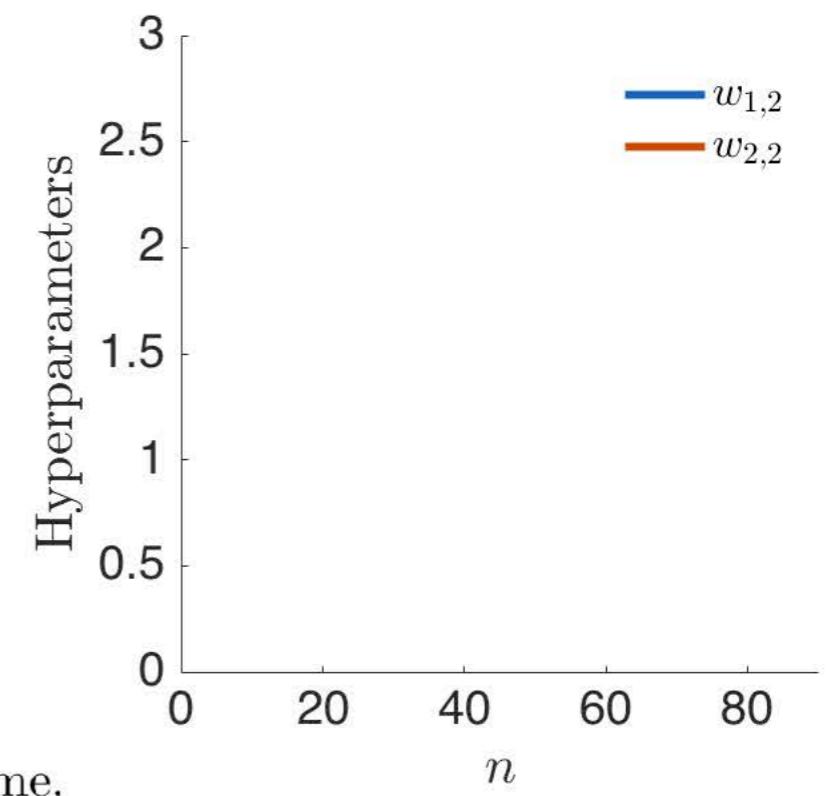
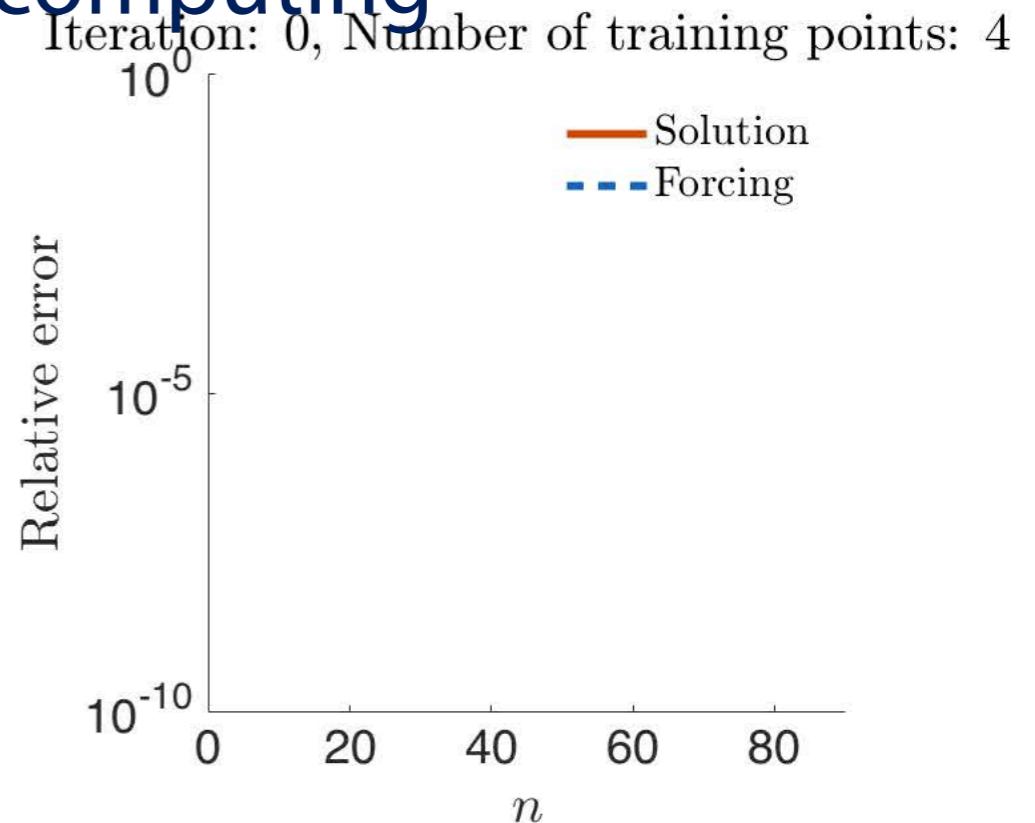
# A new paradigm in scientific computing



$$\frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2} = f(x_1, x_2)$$



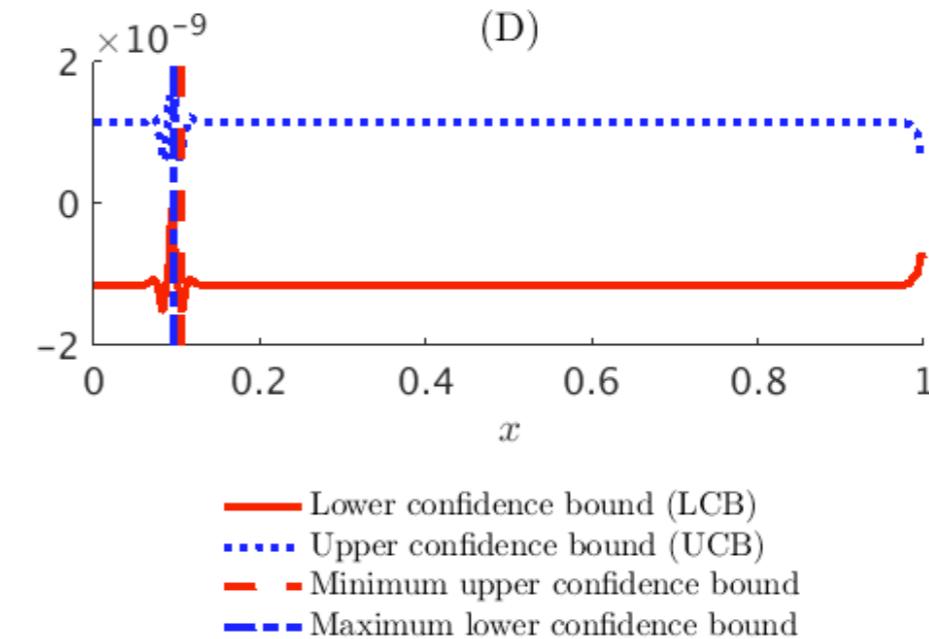
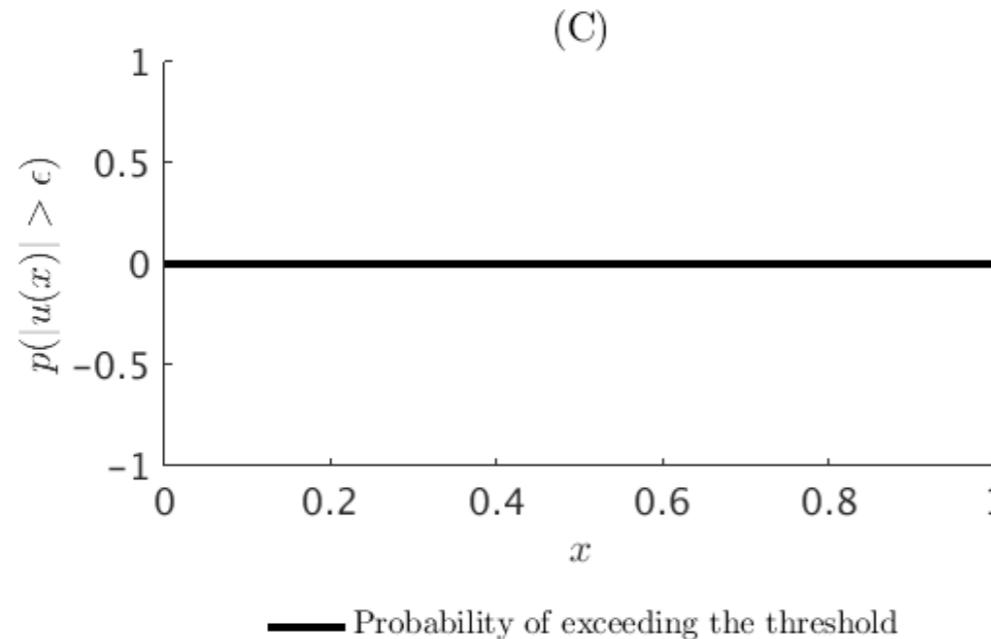
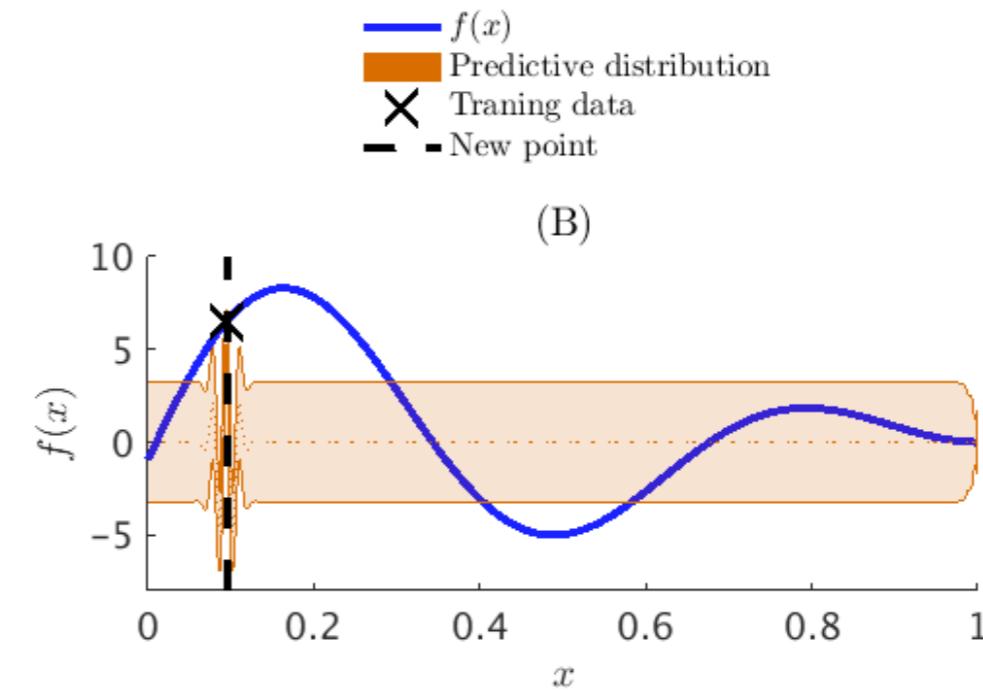
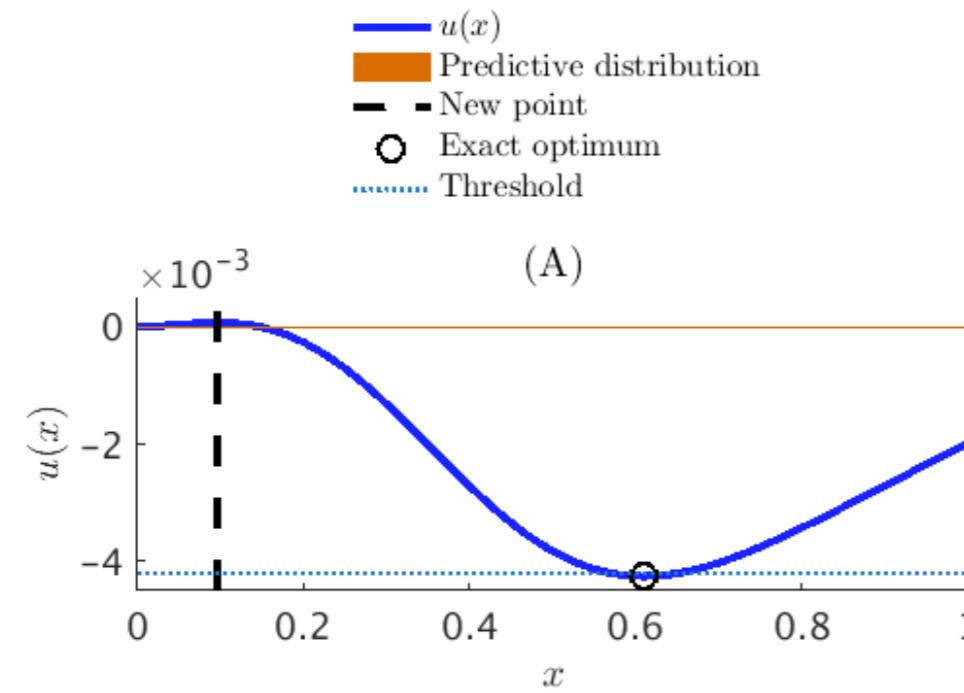
- denotes the training data actively collected by the scheme.
- ★ denotes the next sampling point suggested by the active learning scheme.



# A new paradigm in scientific computing

$$\begin{aligned}\mathcal{L}_x u(x) &:= \frac{d^4}{dx^4} u(x) = f(x) \\ u(0) &= u'(0) = 0 \\ u''(1) &= 0 \\ u'''(1) &= f(1)\end{aligned}$$

- ➊ Given noisy observations of the loading  $f(x)$ , solve for the displacement  $u(x)$ .
- ➋ Find the maximum displacement  $|u(x)|$ .
- ➌ Given the threshold  $\epsilon$ , find the probability of failure.



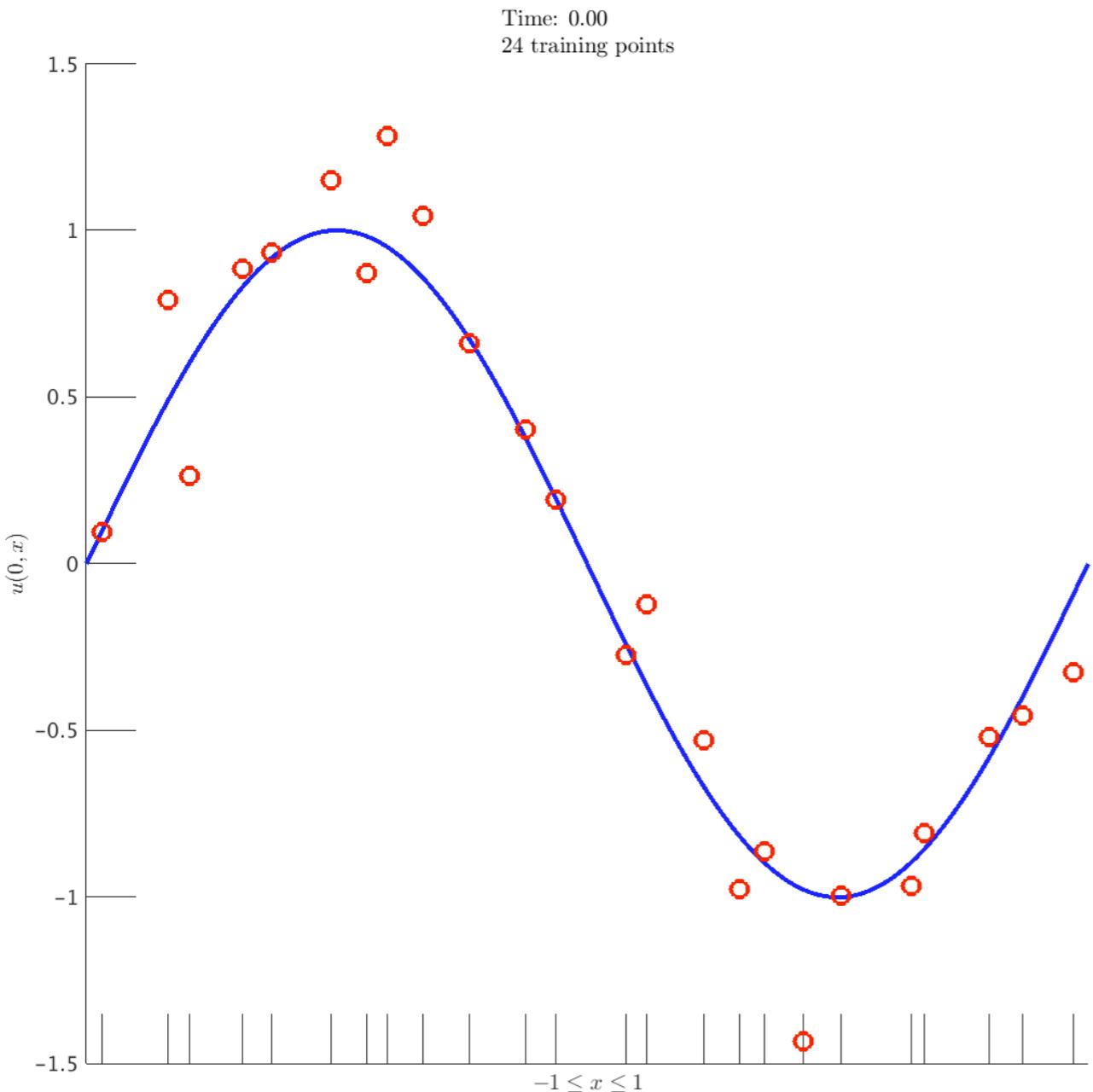
# A new paradigm in scientific computing

Example: 1D viscous Burgers  $\rightarrow$  The equation, along with the choice of a time-stepping scheme define a GP prior!

$$\begin{aligned} \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} &= \nu \frac{\partial^2 u}{\partial x^2}, \quad x \in [-1, 1], \quad t \in [0, 1], \\ u(0, x) &= u^0(x) \\ u(t, 0) = u(t, 1) &= 0. \quad \nu = \pi/100 \end{aligned}$$

$$+ \quad \mathcal{L}_x u^{n+1}(x) := u^{n+1}(x) + \Delta t \ u^n(x) \frac{du^{n+1}(x)}{dx} - \Delta t \ \nu \frac{d^2 u^{n+1}(x)}{dx^2} = u^n(x)$$

e.g., Backward Euler time-stepping

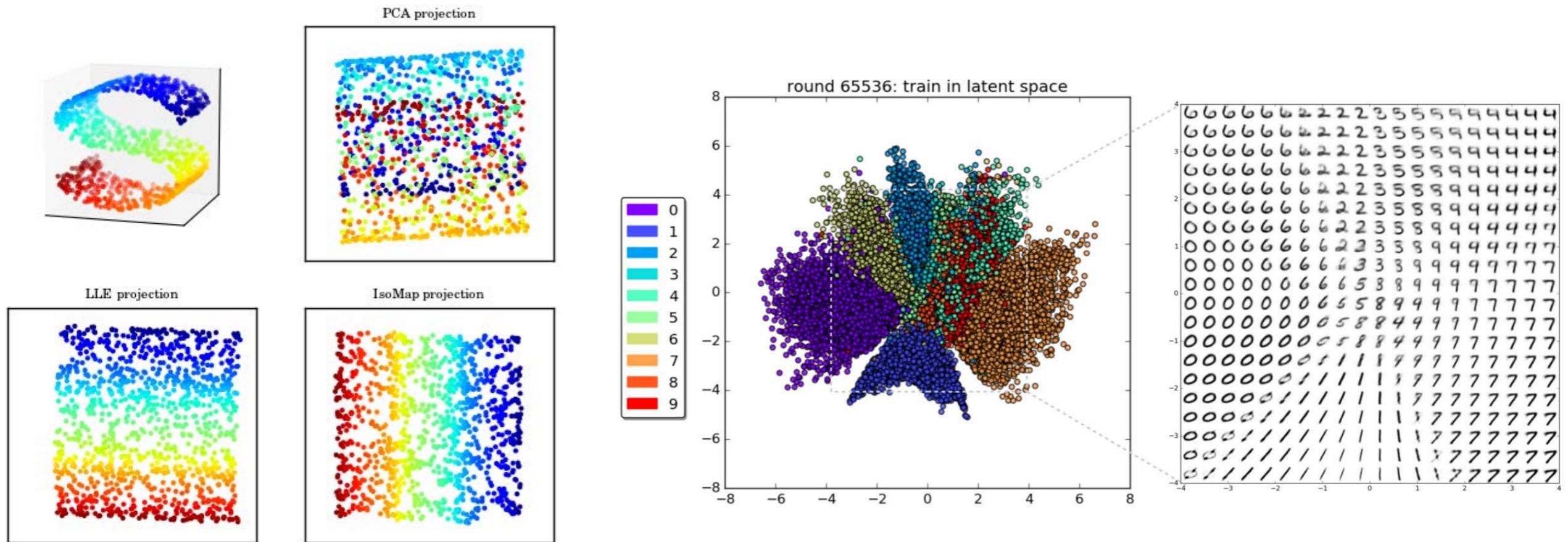


## Remarks:

- Despite starting from a stationary prior:  
 $u^{n+1,n+1}(x) \sim \mathcal{GP}(0, k^{n+1,n+1}(x, x'; \theta))$   
the structure of the kernel:  
 $k^{n,n}(x, x'; \theta) = \mathcal{L}_x \mathcal{L}_{x'} k^{n+1,n+1}(x, x'; \theta)$   
can generate discontinuous solutions!
- This is a general approach applicable to any nonlinear equation and any time-stepping scheme.
- We only need to derive the kernel  $k^{n,n}(x, x'; \theta)$
- Under this setup, fully implicit time-stepping schemes have the same complexity as their explicit counterparts. Hence, one can obtain highly accurate and stable schemes at no extra cost.

# Unsupervised learning

- Density estimation
- Learning to draw samples from a distribution
- Learning to denoise samples from a distribution
- Find a low-dimensional manifold that the data lies near
- Cluster the data into groups of related examples



## A classic unsupervised task:

Find the “best” representation of the data.

\*By “best” we can mean different things, but generally speaking we are looking for a representation that preserves as much information about  $x$  as possible while obeying some penalty or constraint aimed at keeping the representation simpler or more accessible than  $x$  itself.

## PCA and probabilistic PCA

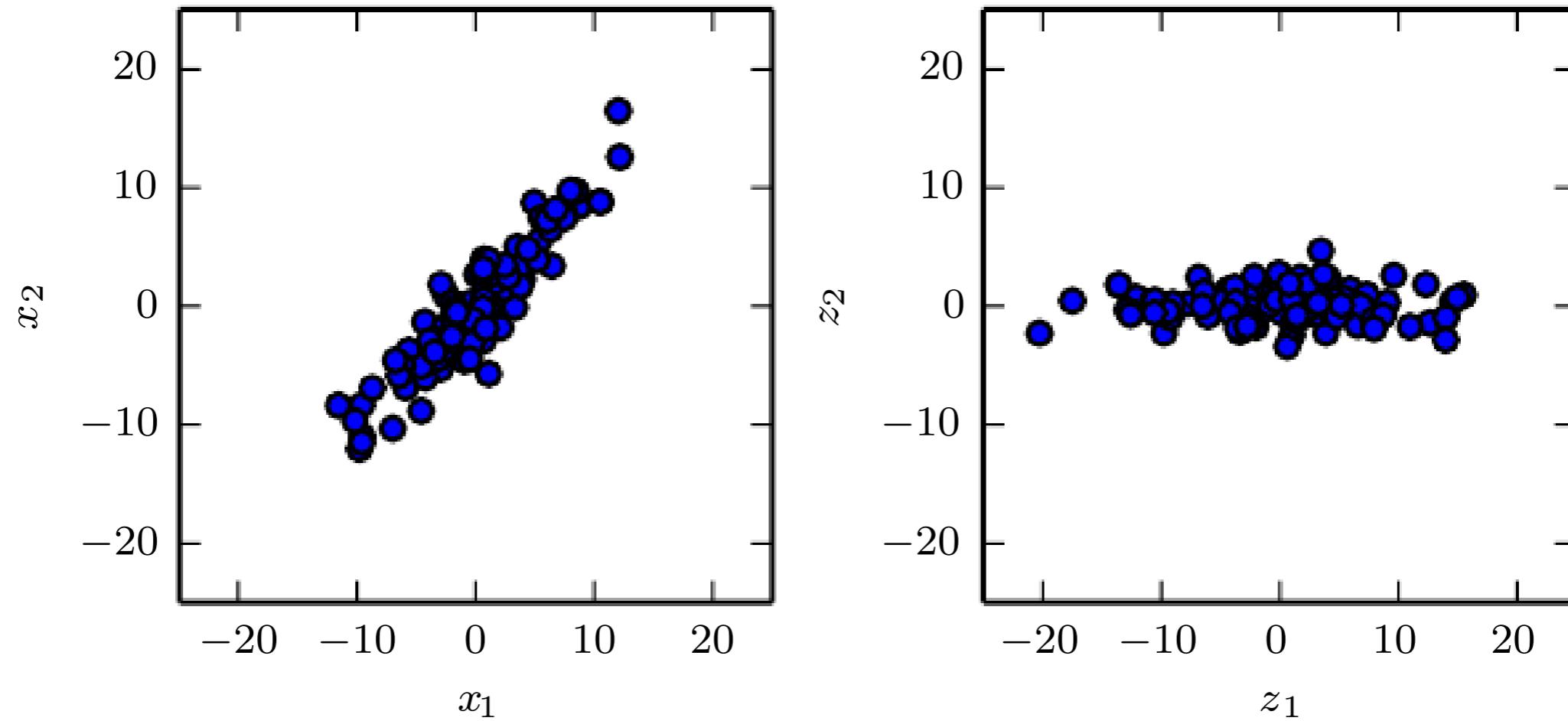


Figure 5.8: PCA learns a linear projection that aligns the direction of greatest variance with the axes of the new space. (*Left*) The original data consists of samples of  $\mathbf{x}$ . In this space, the variance might occur along directions that are not axis-aligned. (*Right*) The transformed data  $\mathbf{z} = \mathbf{x}^\top \mathbf{W}$  now varies most along the axis  $z_1$ . The direction of second most variance is now along  $z_2$ .

# Tricks of the trade

*Jensen's inequality, importance sampling and variational bounds:*

$$\begin{aligned}-\log p(x) &= -\log \int p(x, y) dy \\&= -\log \int q(y|x) \frac{p(y, x)}{q(y|x)} dy \\&\leq -\int q(y|x) \log \frac{p(y, x)}{q(y|x)} dy\end{aligned}$$

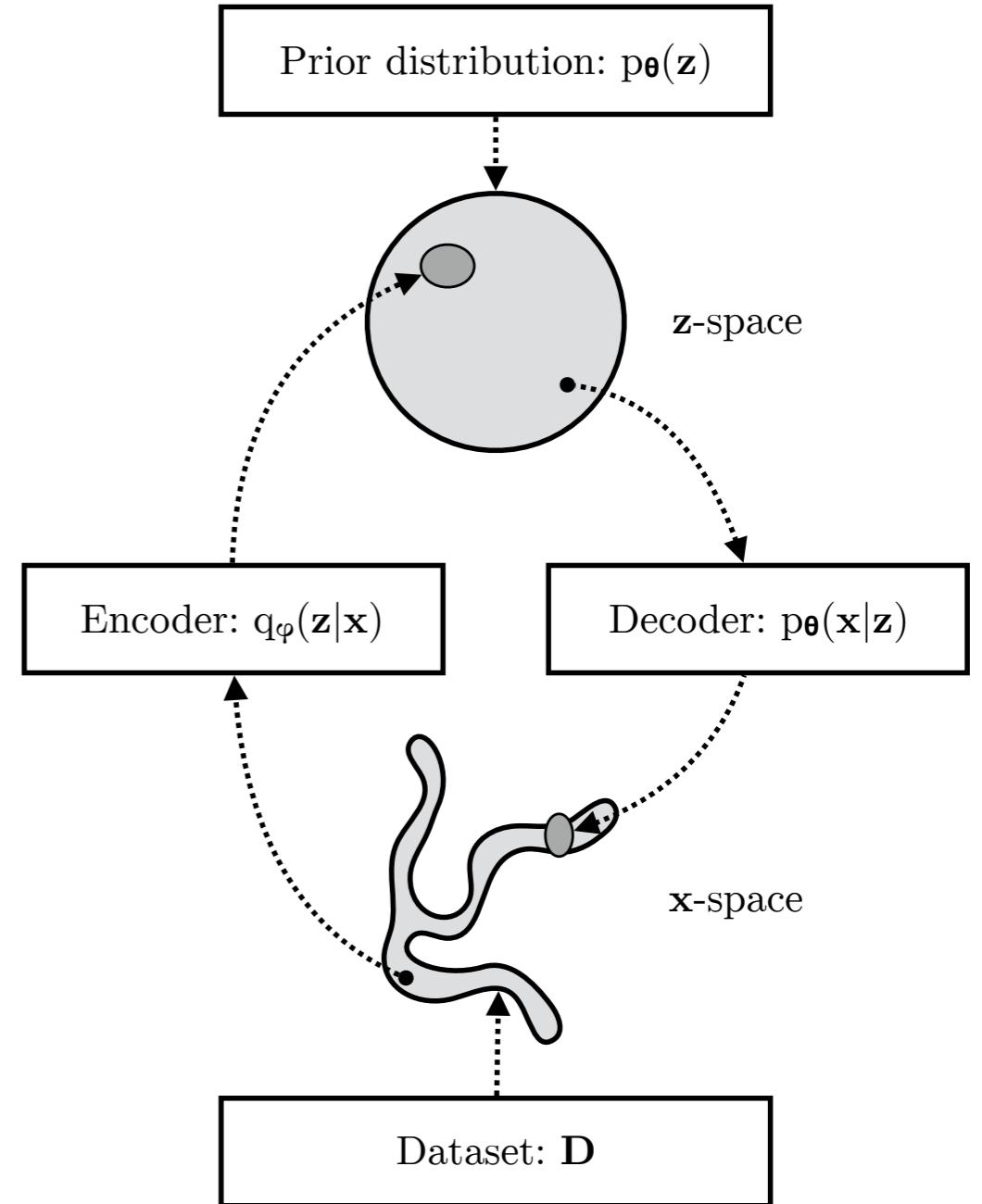
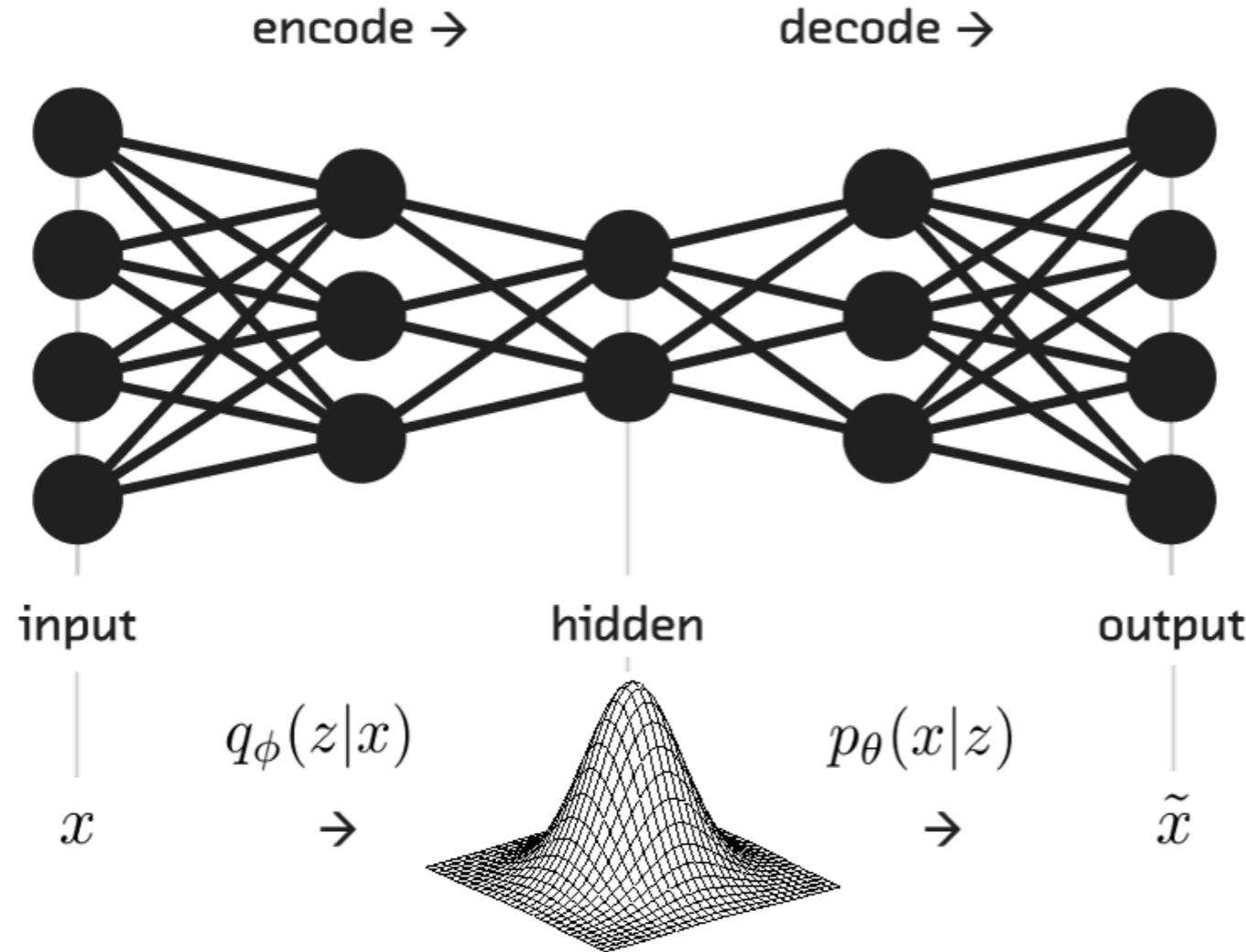
*Unbiased gradients through the re-parametrization trick*

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p(z; \theta)}[f(z)] &= \nabla_{\theta} \int p(z; \theta) f(z) dz \\&= \nabla_{\theta} \int p(\epsilon) f(z) d\epsilon = \nabla_{\theta} \int p(\epsilon) f(g(\epsilon, \theta)) d\epsilon \\&= \nabla_{\theta} \mathbb{E}_{p(\epsilon)}[f(g(\epsilon, \theta))] = \mathbb{E}_{p(\epsilon)}[\nabla_{\theta} f(g(\epsilon, \theta))]\end{aligned}$$

*Density ratio estimation from samples using binary classifiers*

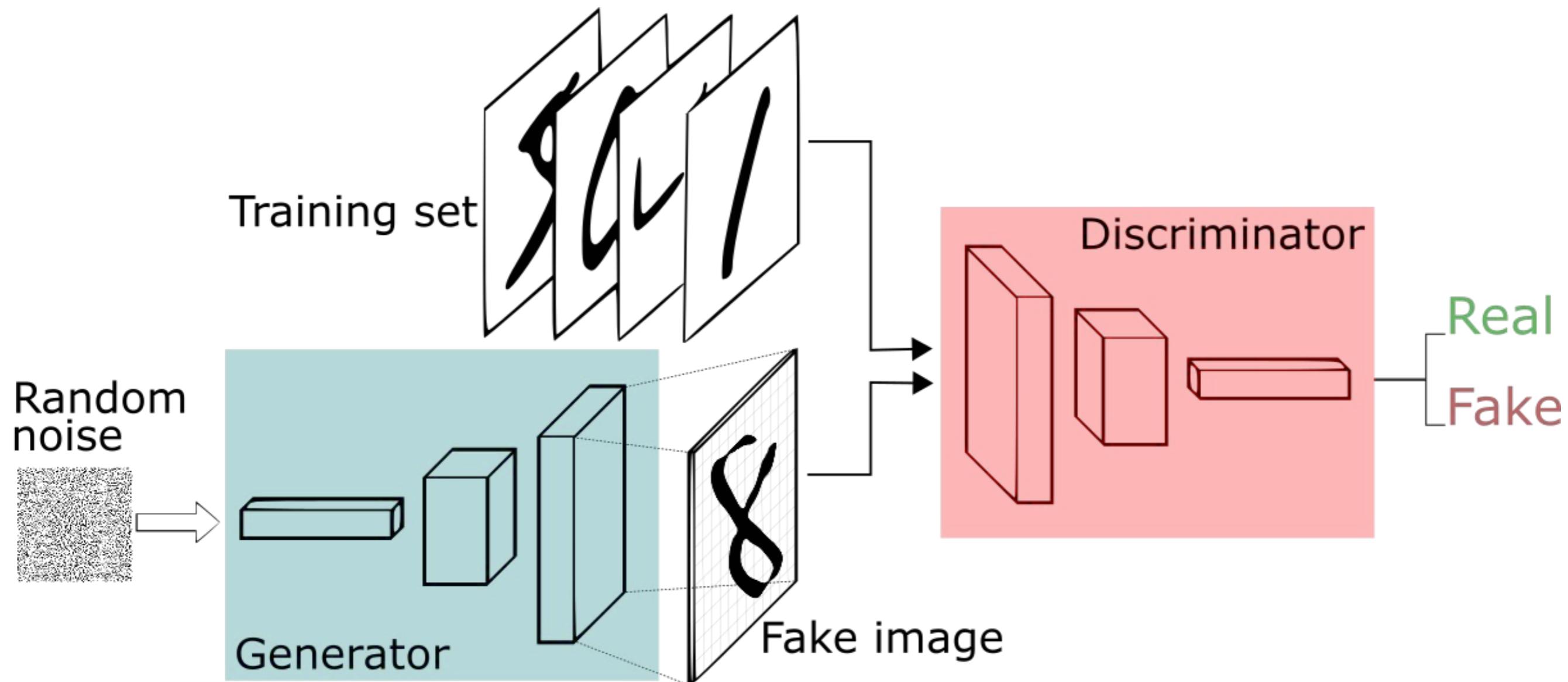
$$\begin{aligned}r(x) &= \frac{\rho(x)}{q(x)} = \frac{p(x|y=+1)}{p(x|y=-1)} \\&= \frac{p(y=+1|x)p(x)}{p(y=+1)} \Bigg/ \frac{p(y=-1|x)p(x)}{p(y=-1)} \\&= \frac{p(y=+1|x)}{p(y=-1|x)} = \frac{p(y=+1|x)}{1-p(y=+1|x)} = \frac{S(x)}{1-S(x)}\end{aligned}$$

# Variational auto-encoders



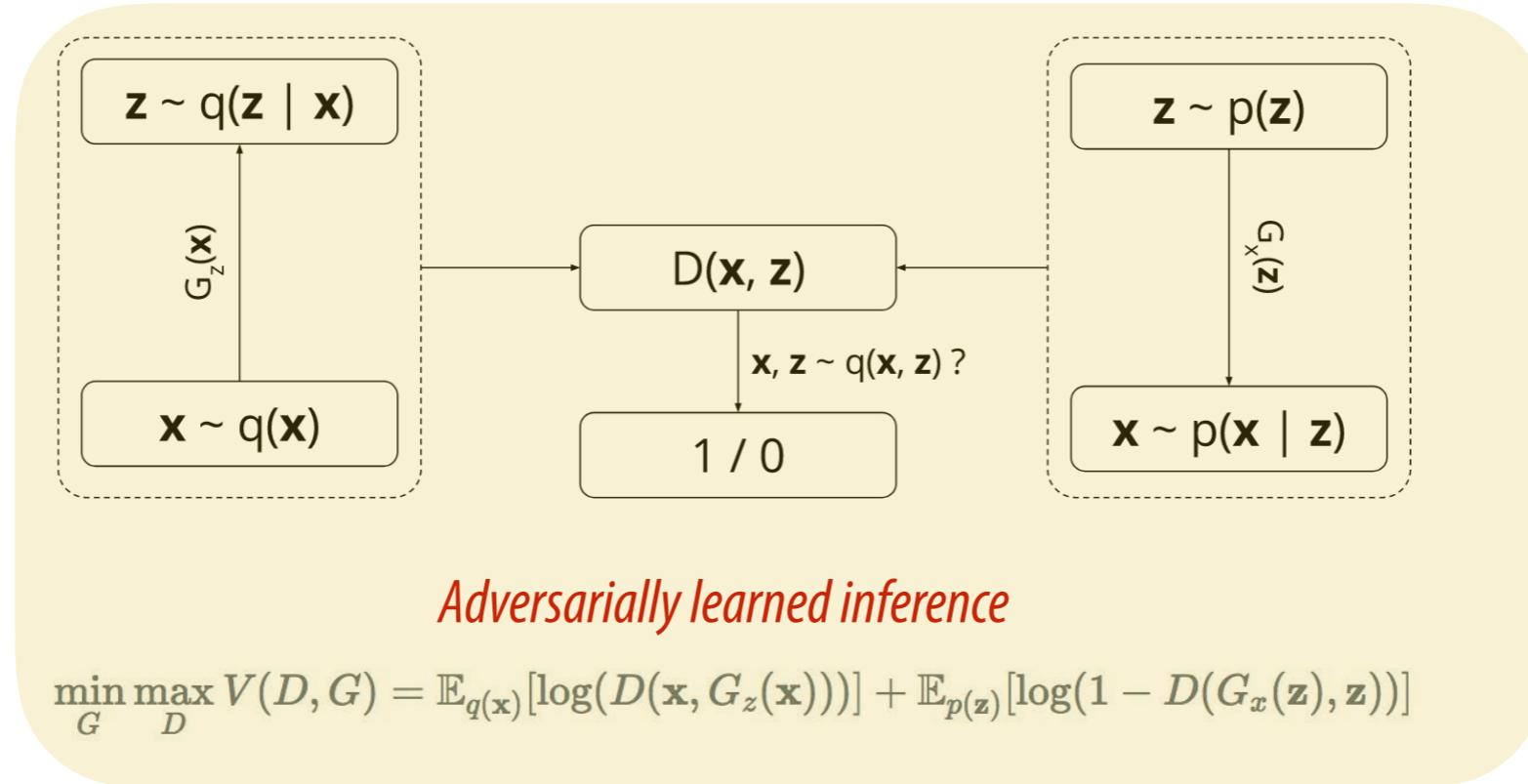
A VAE learns stochastic mappings between the observed  $x$ -space, whose empirical distribution  $q_D(x)$  is typically complicated, and a latent  $z$ -space, whose distribution can be relatively simple (such as spherical, as in this figure). The generative model learns a joint distribution  $p_\theta(x,z) = p_\theta(x|z)p_\theta(z)$ , factorized into a prior distribution over latent space,  $p_\theta(z)$ , and a stochastic decoder  $p_\theta(x|z)$ . The stochastic encoder  $q_\phi(z|x)$ , also called inference model, approximates the true but intractable posterior  $p_\theta(z|x)$  of the generative model.

# Generative adversarial networks

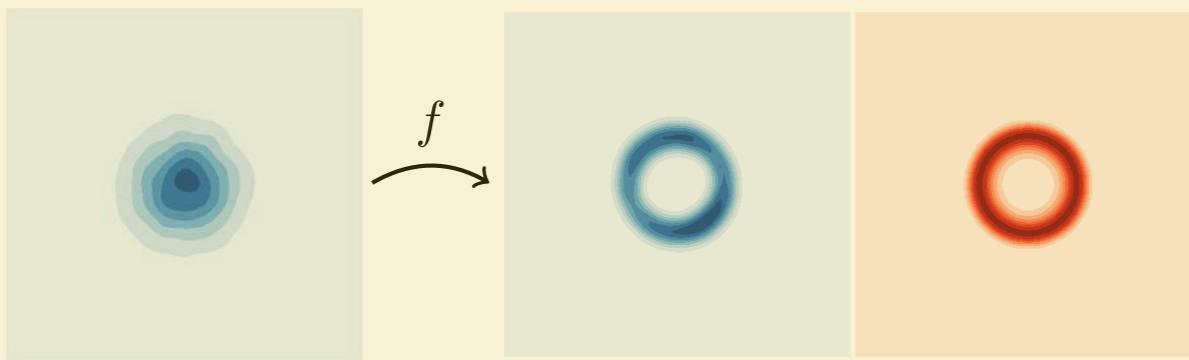


$$\min_G \max_D V(D, G) = \mathbb{E}_{q(\mathbf{x})} [\log(D(\mathbf{x}))] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

# Convergence of variational and adversarial learning



Instead of using approximating distributions of a given pre-defined form (e.g. Gaussian) we can implicitly parametrize them using deep neural networks.



*Adversarial Variational Bayes*

