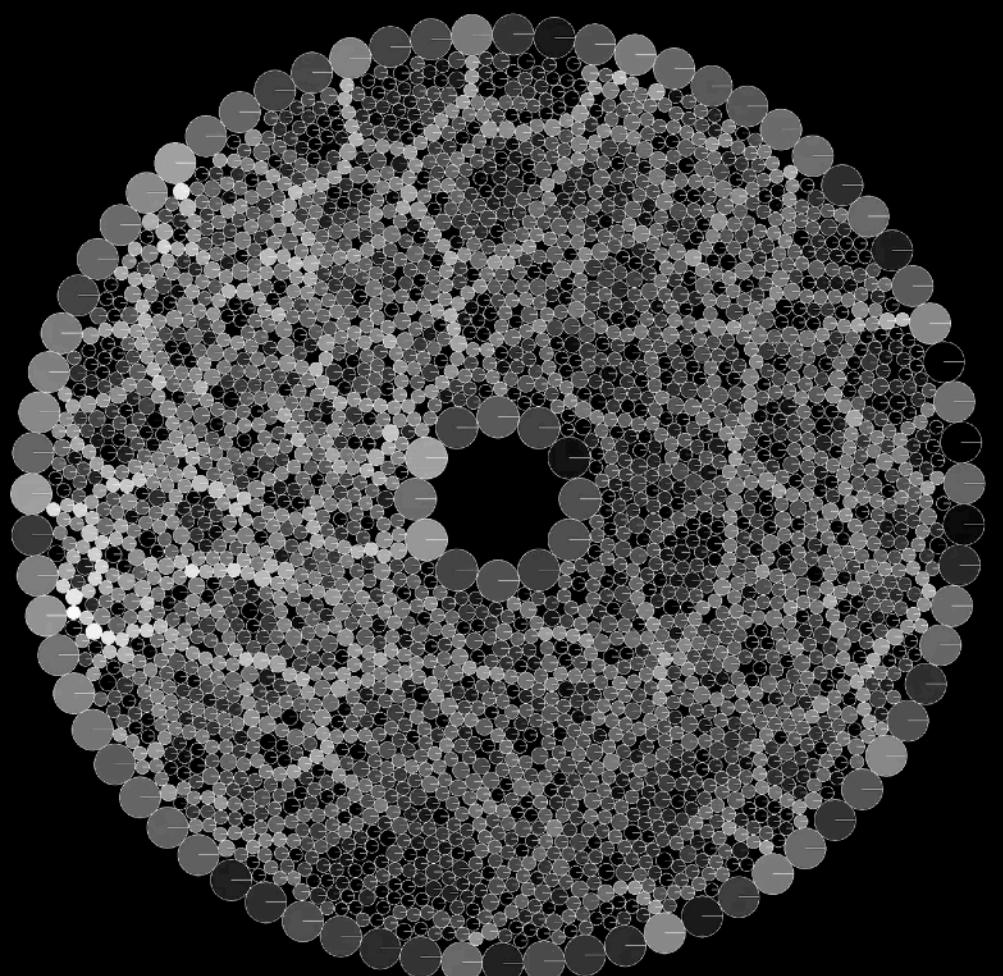


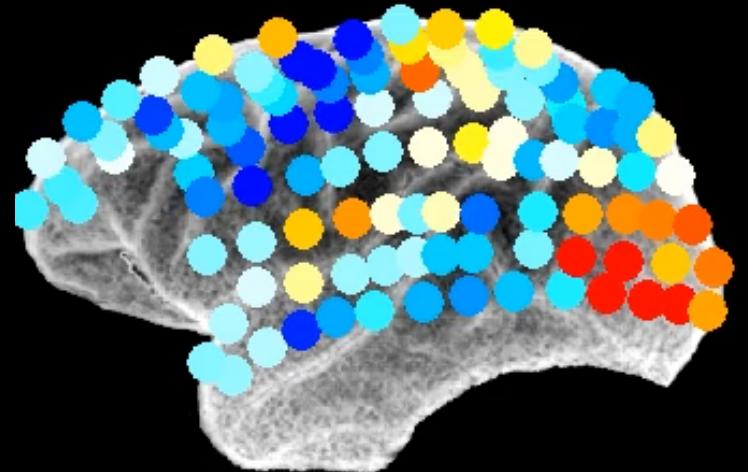
A Brief Intro to Deep Learning

Joseph Bakarji
Olivia Thomas
Jan William

Complex systems are **high dimensional**



Granular Material (DEM Simulation)



Raut, Ryan V., et al. Science advances (2021)



lasefist
<https://youtu.be/uI-KrtV0PJA>

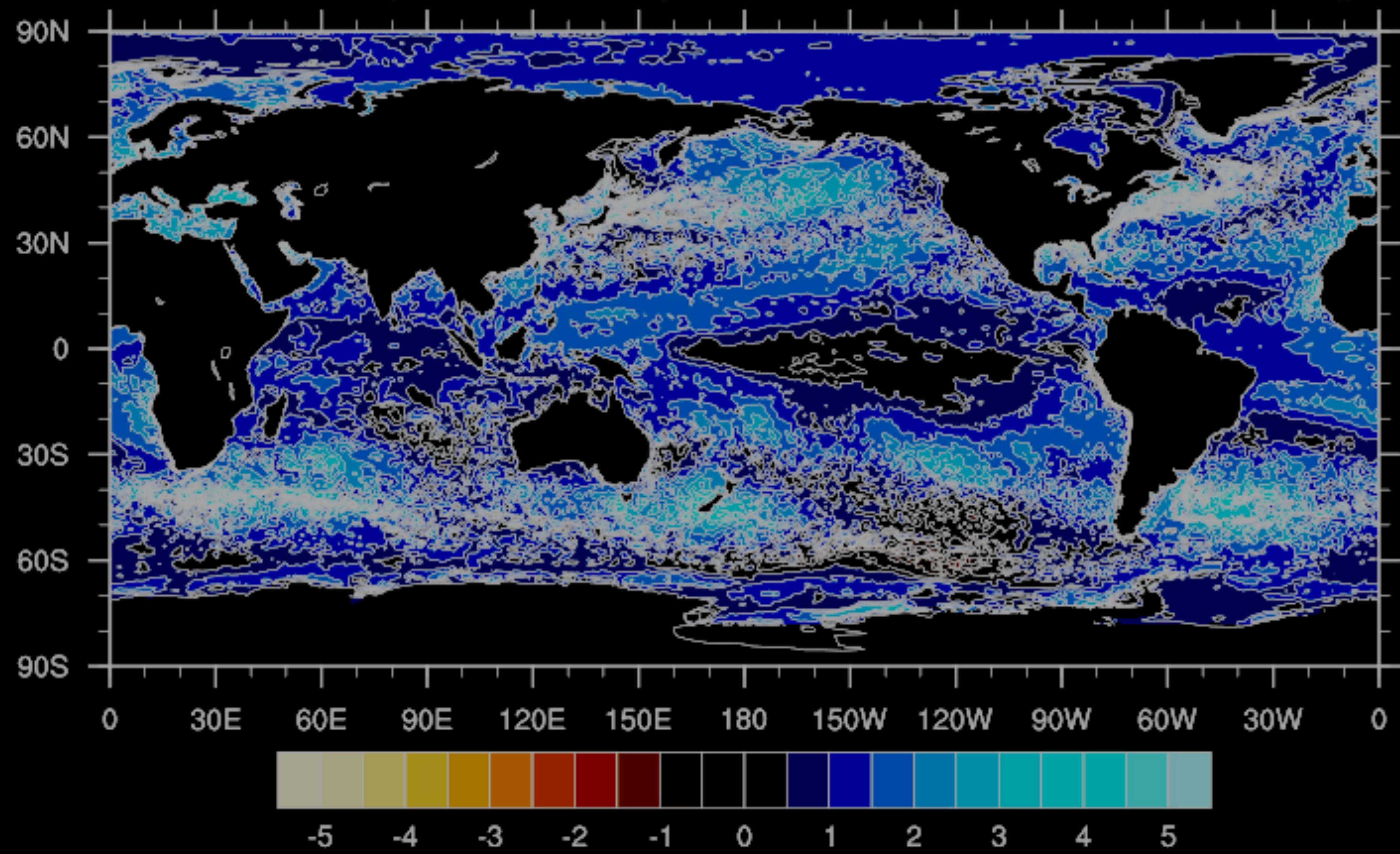


National Geographic
https://youtu.be/RtUQ_pz5wlo

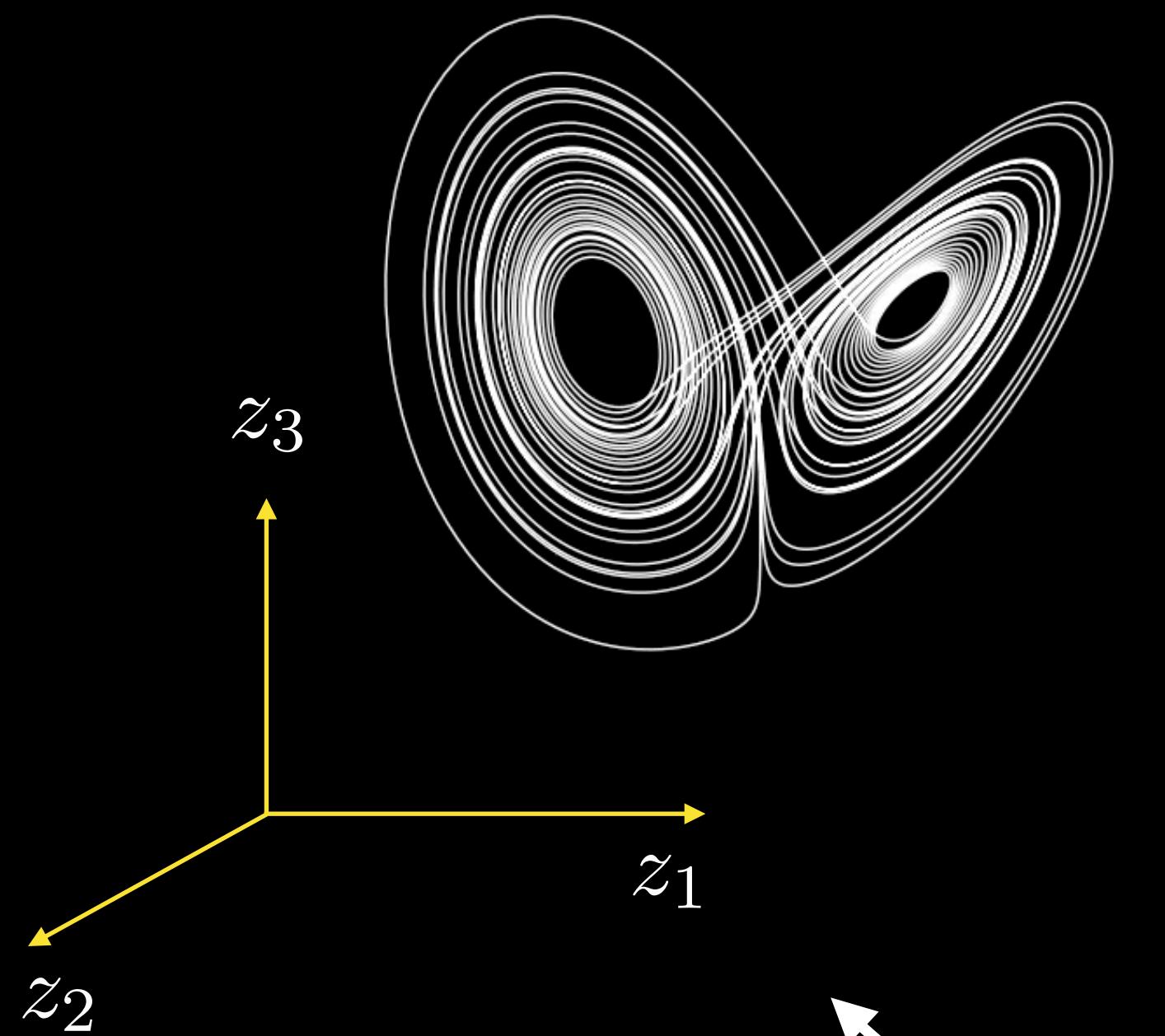
We only have access to **partial measurements**

1991-2020 LTM

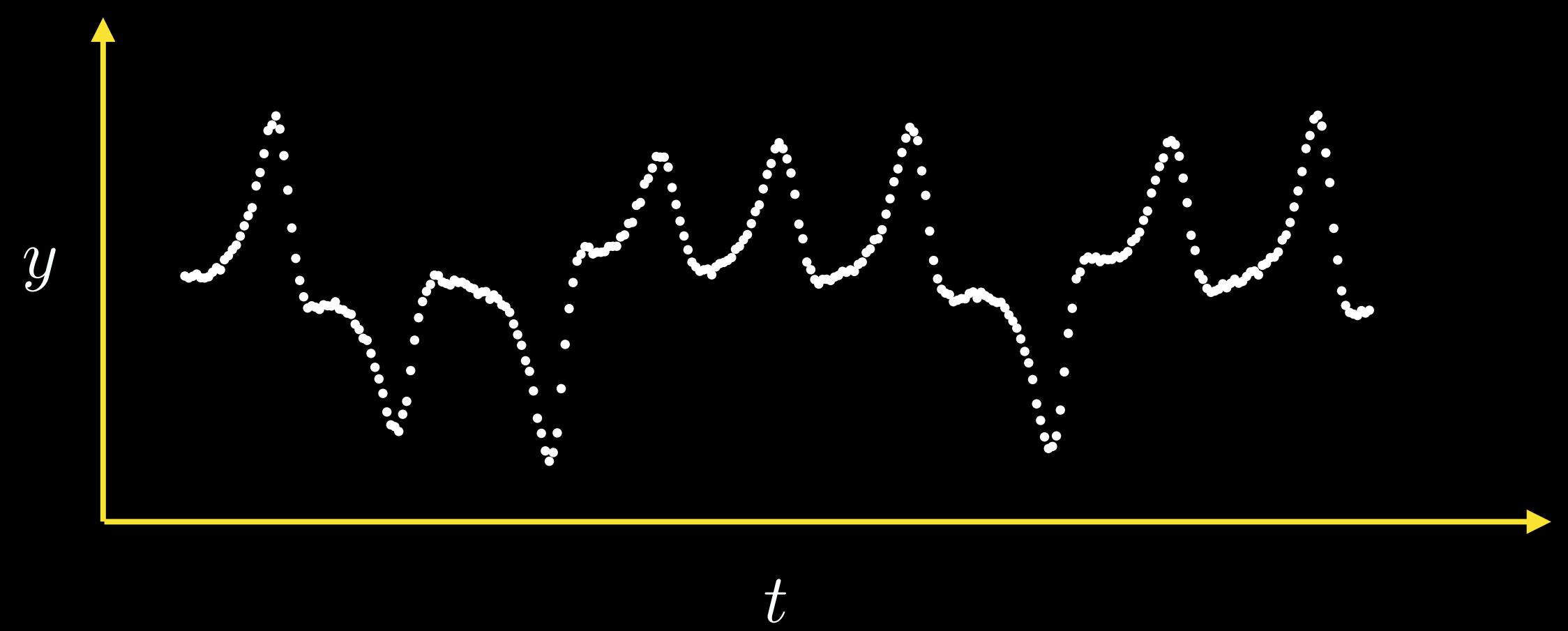
Sea Surface Temperature Anomaly



High Dimensional System



Partial Measurements



Measure



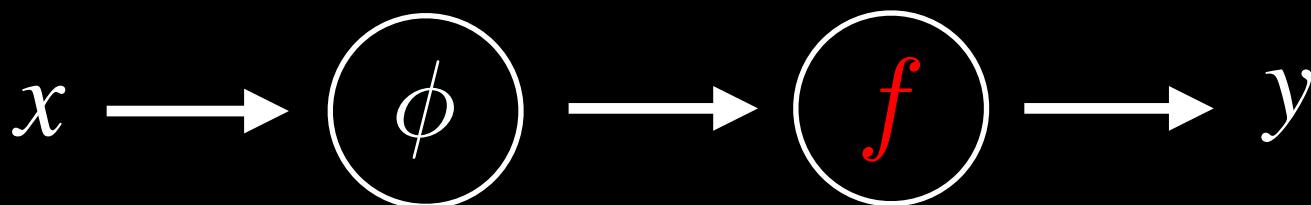
Reconstruct

?



Linear predictors

(x_1, y_1)
(x_2, y_2)
(x_3, y_3)
(x_4, y_4)
\vdots
(x_n, y_n)



feature vector

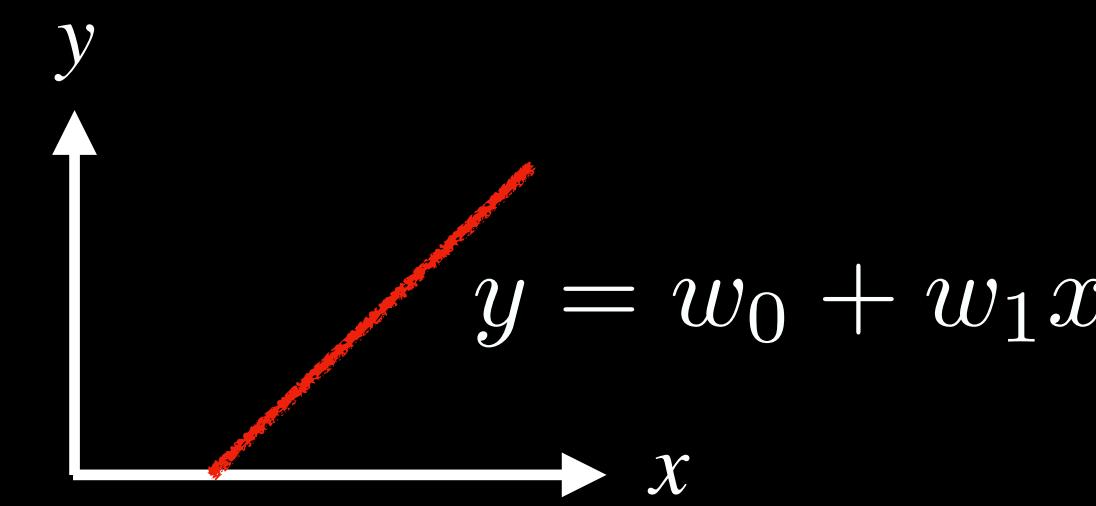
$$\phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_d(x)]$$

linear predictor

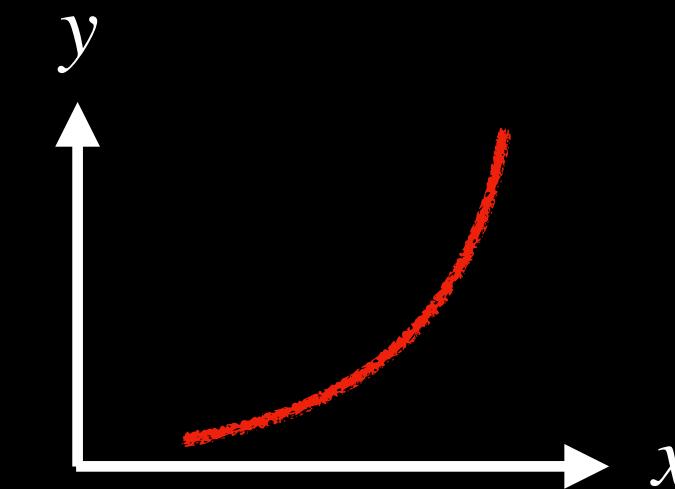
$$\begin{aligned}f_{\mathbf{w}} &= \mathbf{w} \cdot \phi(x) \\&= w_1\phi_1(x) + w_2\phi_2(x) + \cdots + w_d\phi_d(x)\end{aligned}$$

What is linear?

Adding a bias: $\phi(x) = [1, x]$

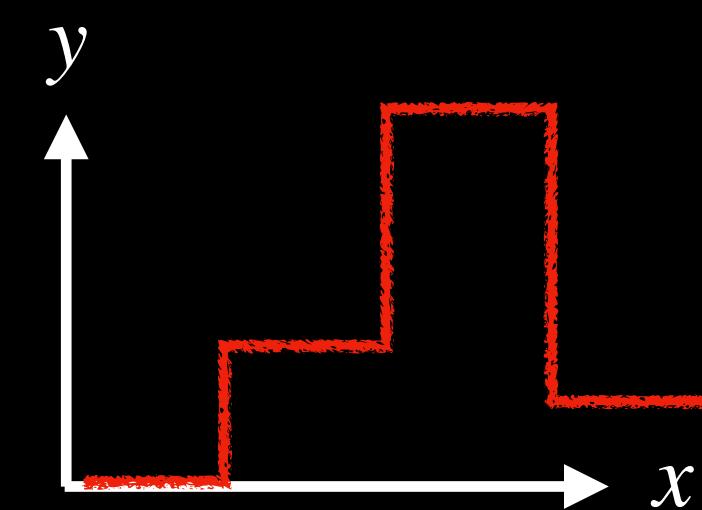


A polynomial predictor: $\phi(x) = [1, x, x^2, x^3]$



Sines and cosines: $\phi(x) = [1, x, \sin(3x)]$

Indicator functions: $\phi(x) = [\mathbf{1}[0 < x \leq 1], \mathbf{1}[1 < x \leq 2], \mathbf{1}[2 < x \leq 3]]$



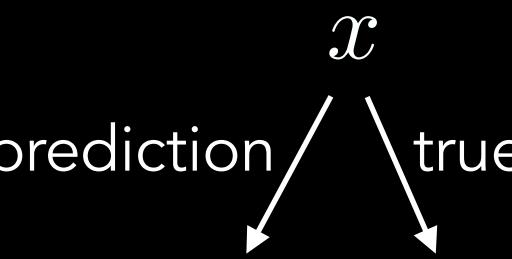
Derivatives: $\phi(x, t) = \left[x^2, \frac{\Delta x}{\Delta t} \right]$

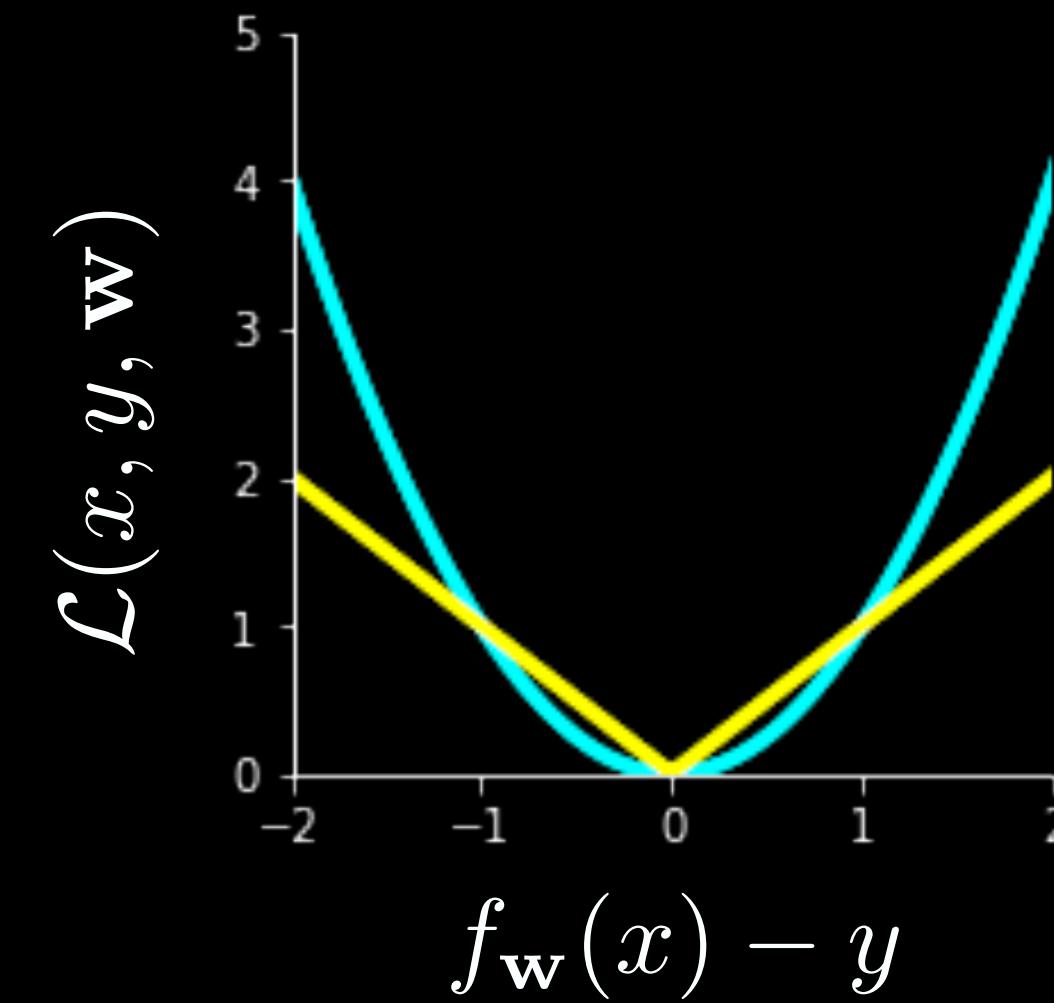
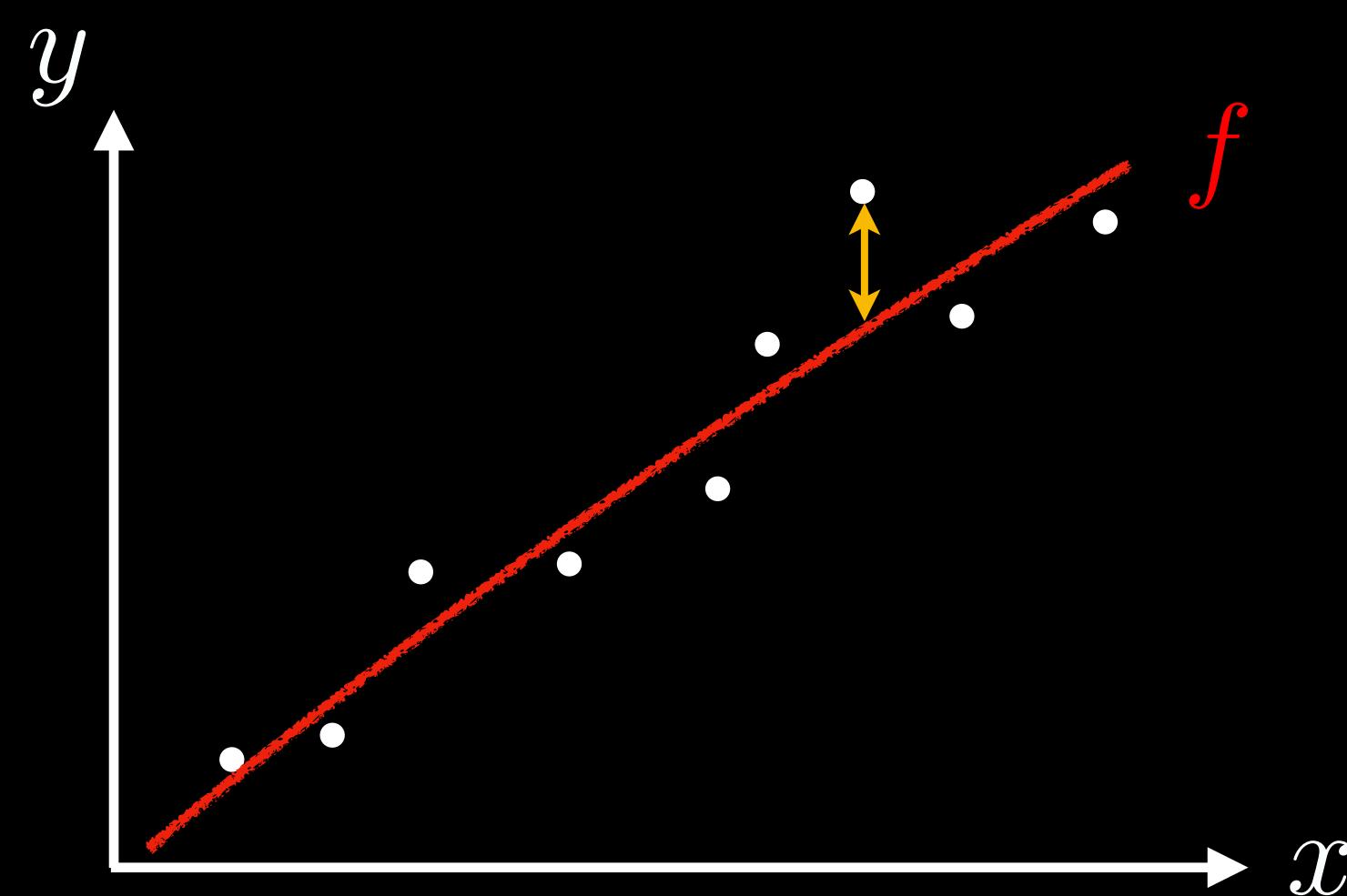
Linear in $\phi(x)$ and w but not in x

The loss function

A **loss function** quantifies how *bad* the predictor is

$$\mathcal{L}(x, y, \mathbf{w}) = \text{distance}(f_{\mathbf{w}}(x), y)$$





$$\mathcal{L}_{\text{sq}}(x, y, \mathbf{w}) = (f_{\mathbf{w}}(x) - y)^2$$

$$\mathcal{L}_{\text{abs}}(x, y, \mathbf{w}) = |f_{\mathbf{w}}(x) - y|$$

Minimize the loss

The **training loss** is the average loss over the data set

$$\mathcal{L}_{\text{train}}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \mathcal{L}(x, y, \mathbf{w})$$

objective

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\mathbf{w})$$

optimal predictor

$$f_{\hat{\mathbf{w}}}(x) = \hat{\mathbf{w}} \cdot \phi(x)$$

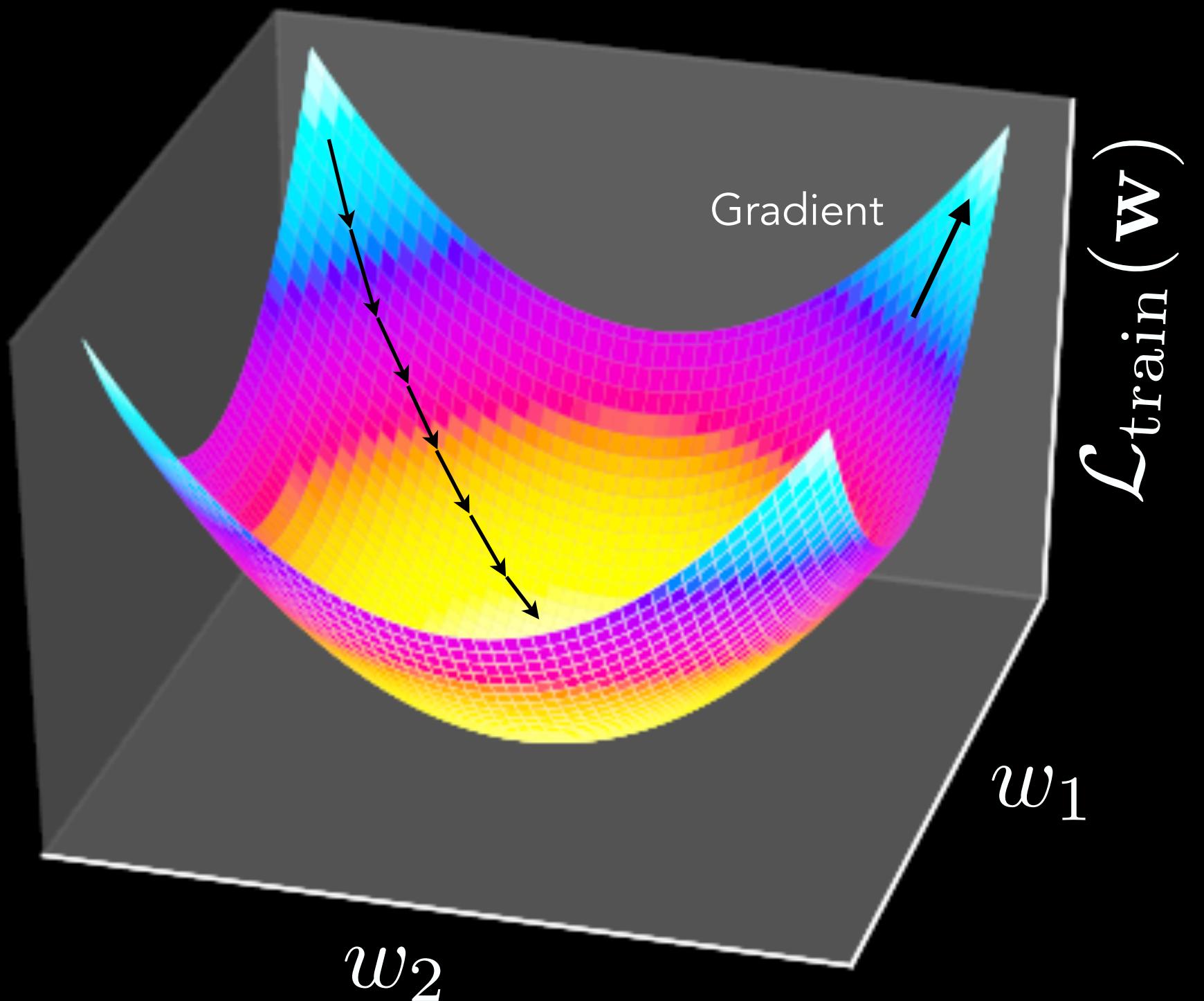
Gradient Descent

Training loss gradient

$$\nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}(\mathbf{w}) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \underbrace{2(\mathbf{w} \cdot \phi(x) - y)\phi(x)}_{\text{Prediction} - \text{True value}}$$

Gradient descent algorithm

```
epochs  
initialize  $\mathbf{w} = [0, \dots, 0]$ ;  
for  $t = 1, \dots, T$  do  
|  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}$   
end  
step size
```



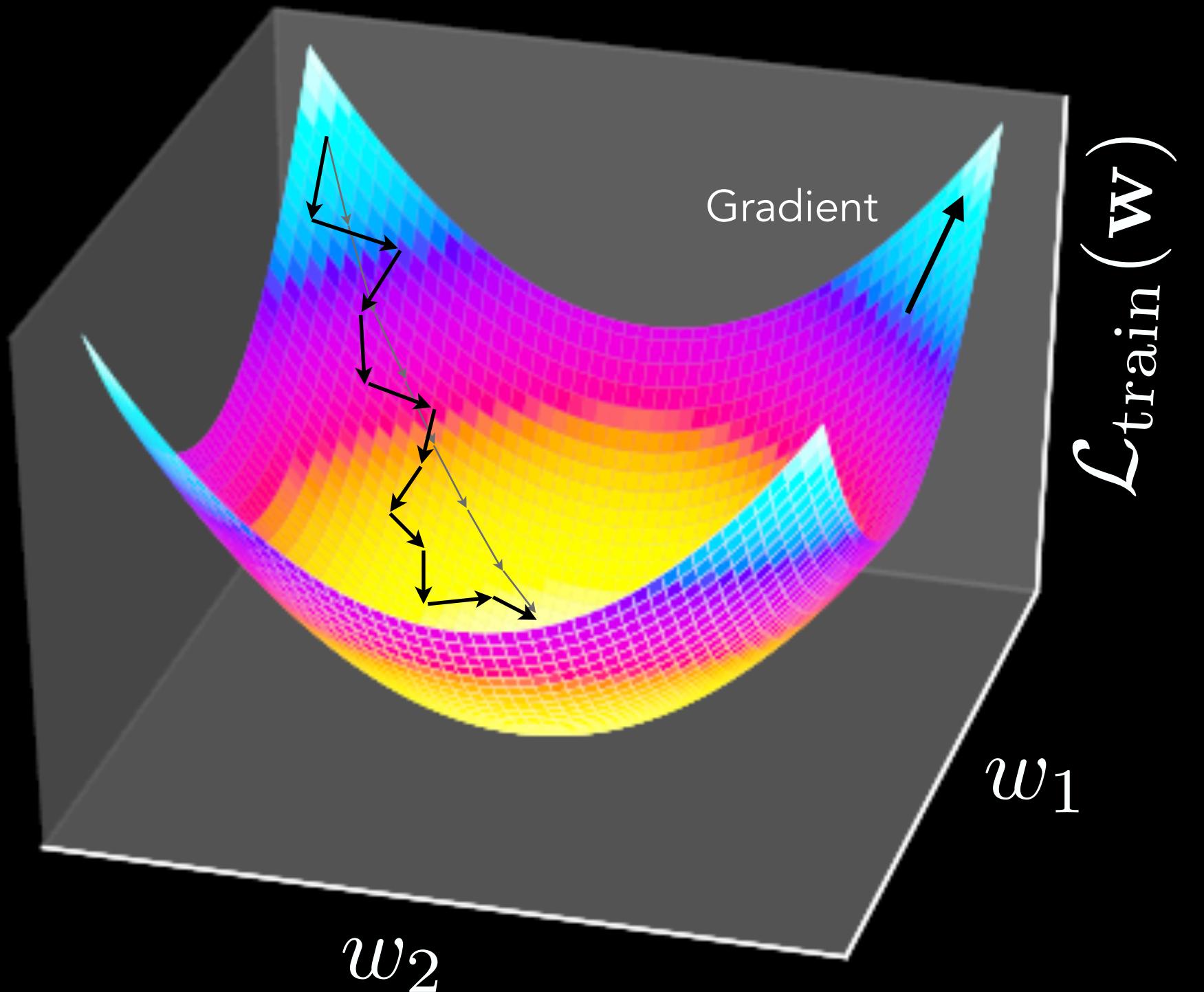
Stochastic Gradient Descent

Gradient descent is too expensive

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}_{\text{train}}$$

Stochastic gradient descent

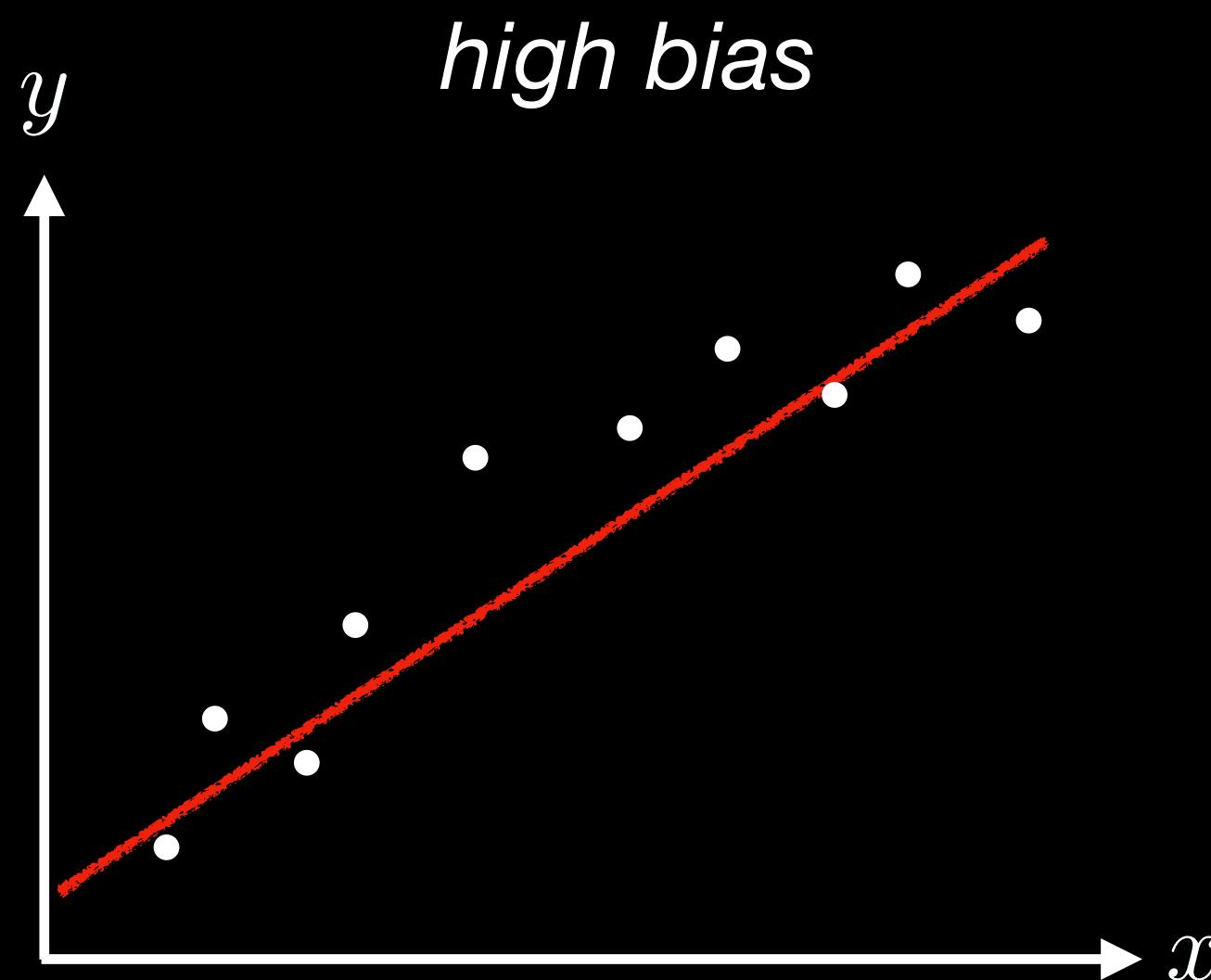
```
initialize  $\mathbf{w} = [0, \dots, 0]$ ;
for  $t = 1, \dots, T$  do
    for  $(x, y) \in \mathcal{D}_{\text{train}}$  do
         $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}(x, y, \mathbf{w})$ 
    end
end
```



Variance vs. Bias

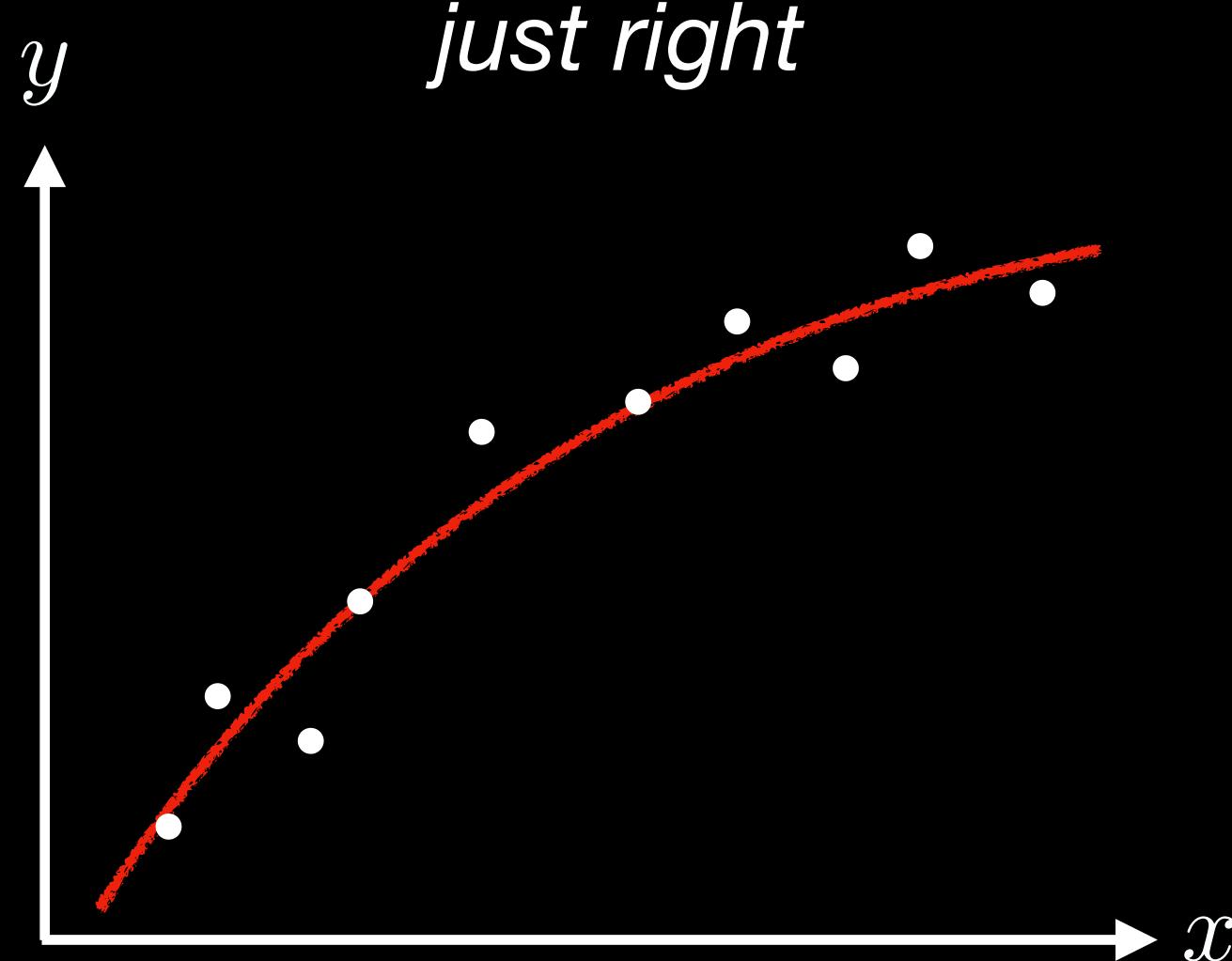
How do you choose $\phi(\cdot)$?

under-fitting



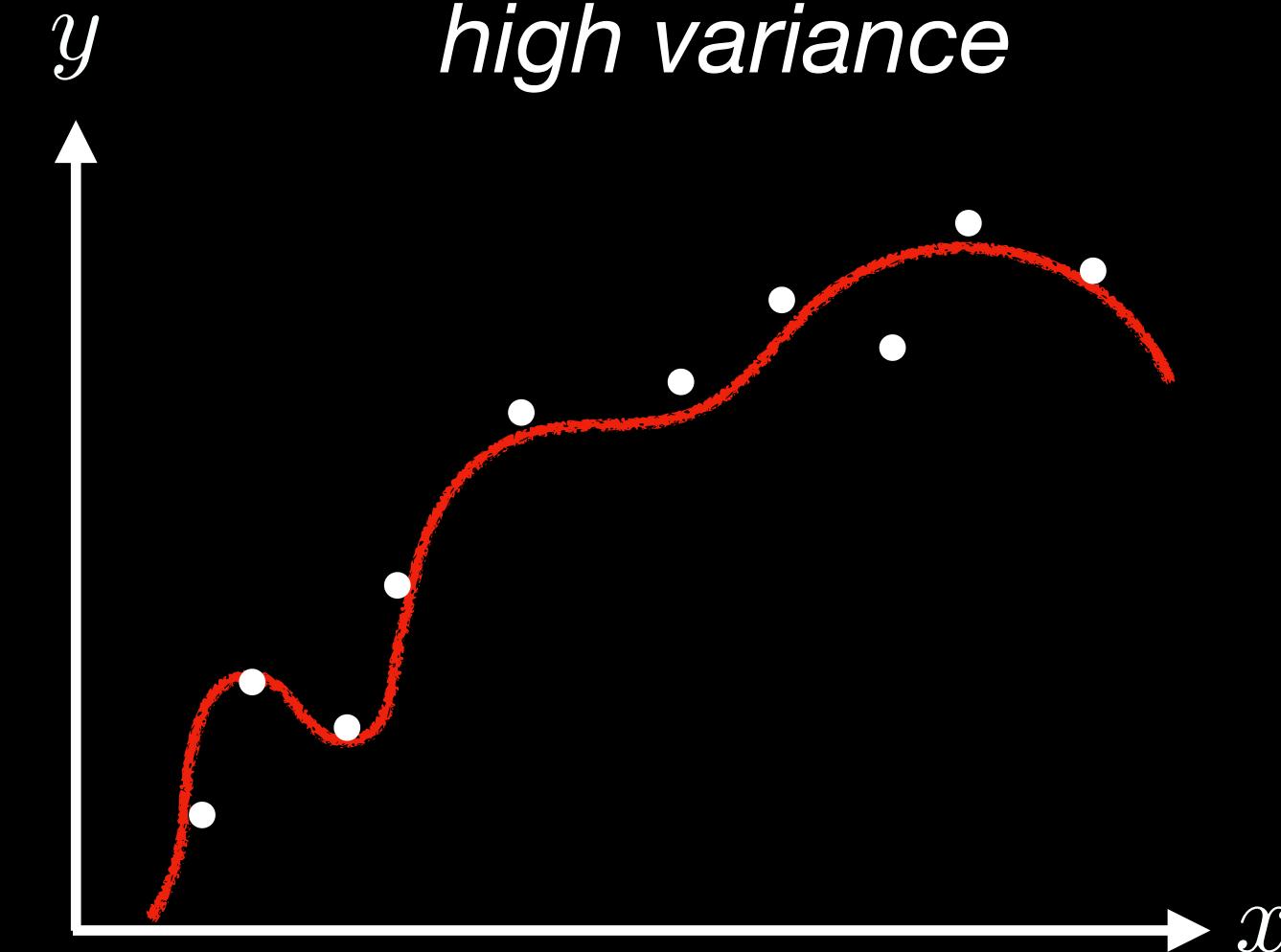
$$\phi(x) = [1, x]$$

just right



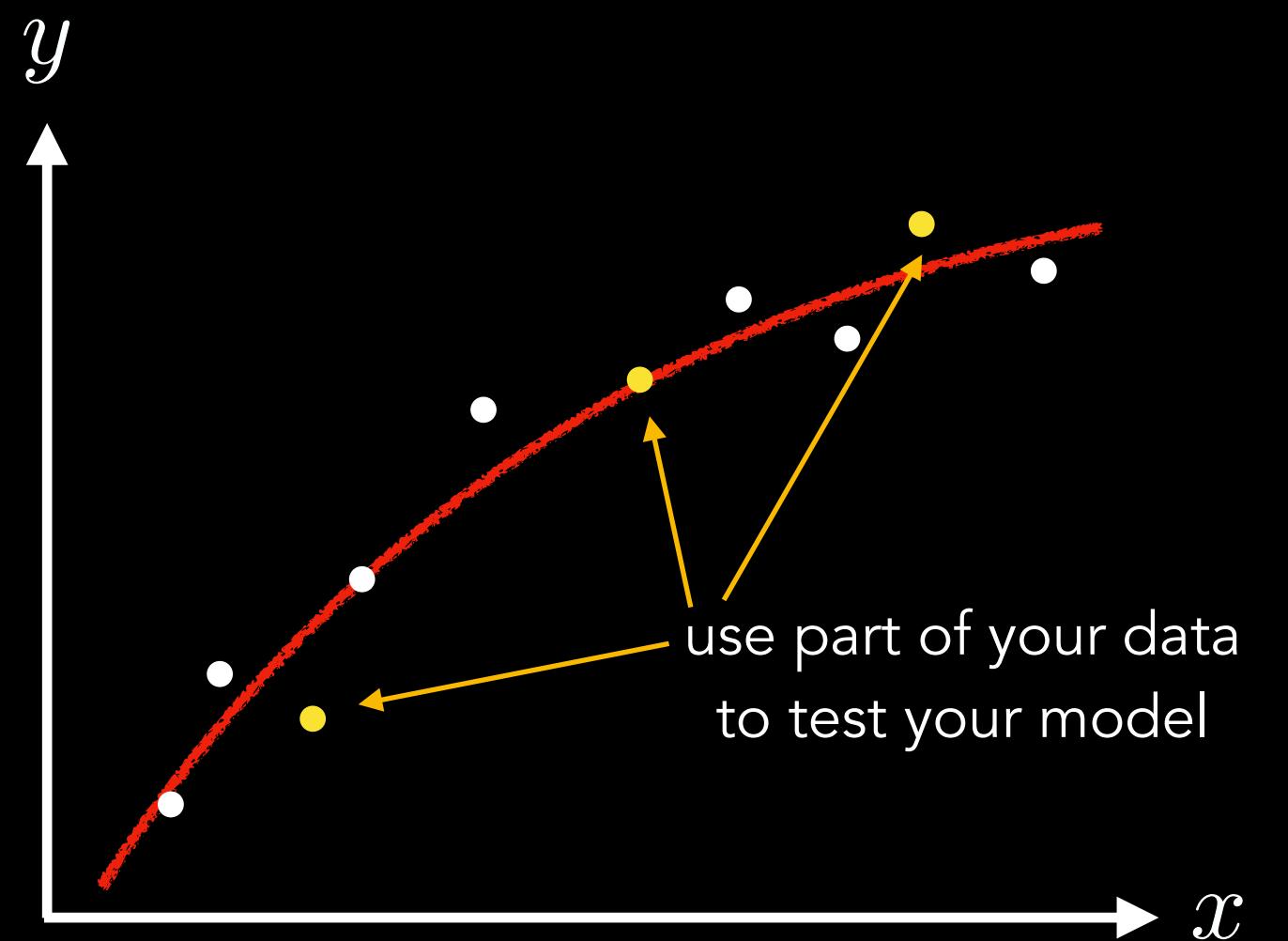
$$\phi(x) = [1, x, x^2]$$

over-fitting
high variance



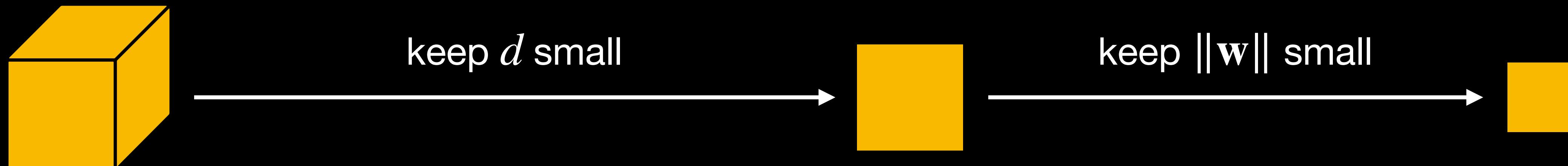
$$\phi(x) = [1, x, x^2, \dots, x^n]$$

Use a test set



Designing the hypothesis class

$$\mathbf{w} \in \mathbb{R}^d$$



Remove features if they don't help

Use automatic feature selection

$$\mathcal{L} = \mathcal{L}_{\text{sq}} + \lambda \|\mathbf{w}\|_0$$

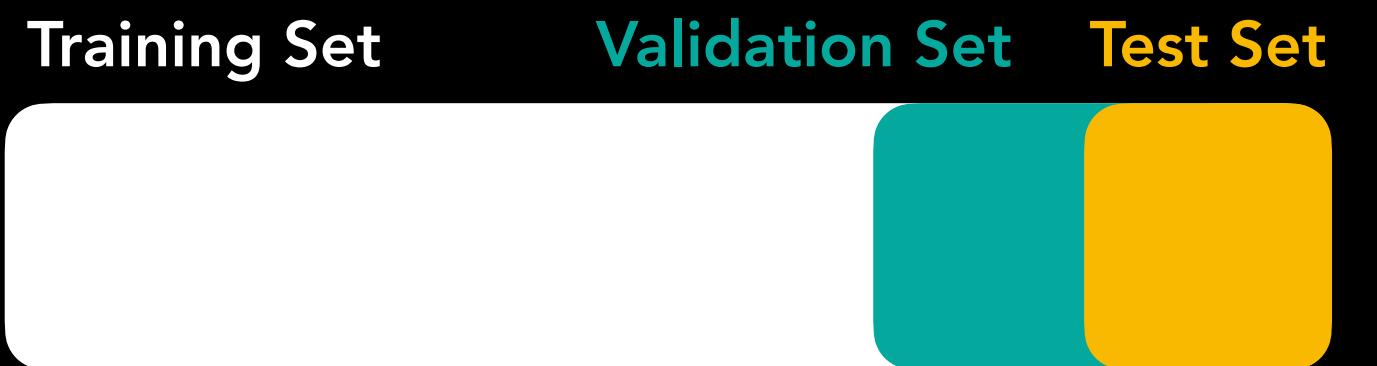
number of nonzero
elements

Use L_2 regularization
 $\mathcal{L} = \mathcal{L}_{\text{sq}} + \lambda \|\mathbf{w}\|^2$

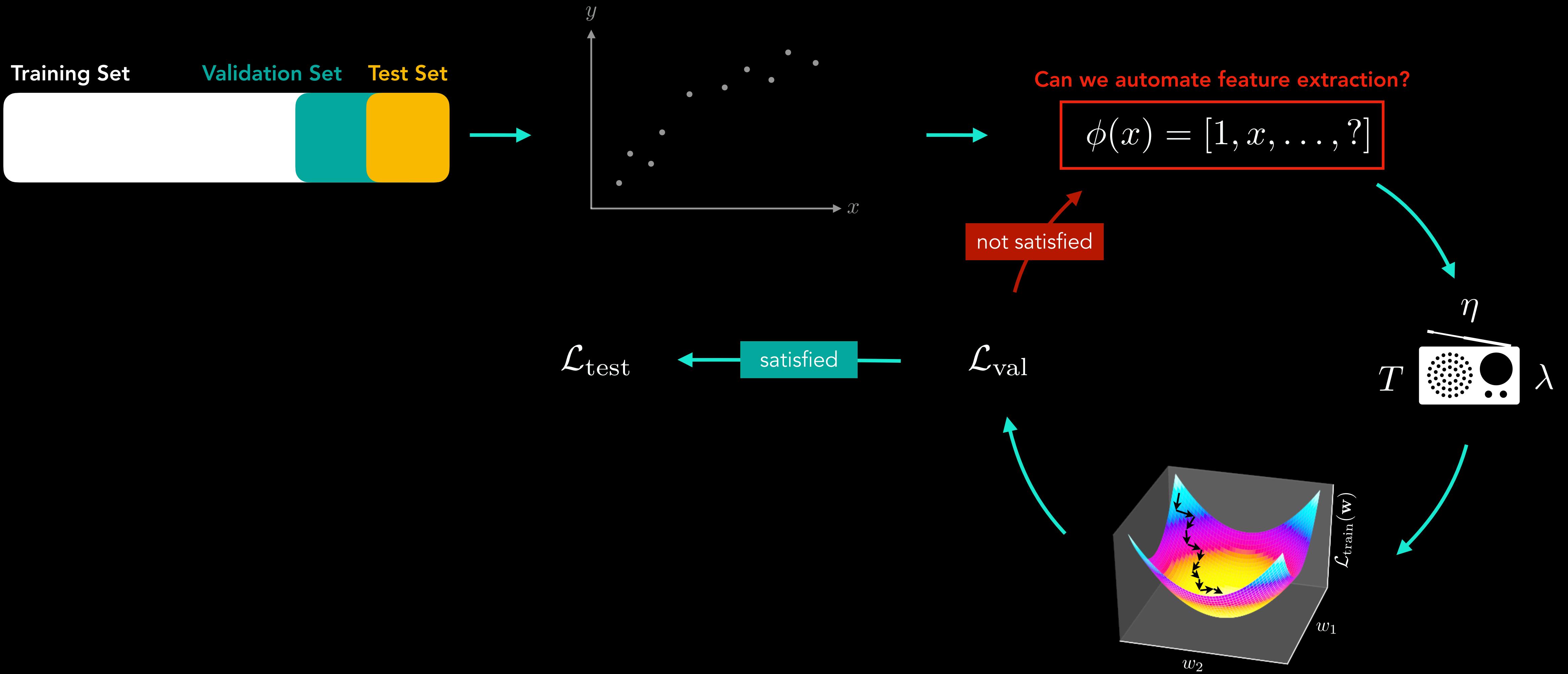
Hyperparameters

How do you choose the

- Regularization parameter
- Number of iterations
- Steps size
- ...



The ML workflow

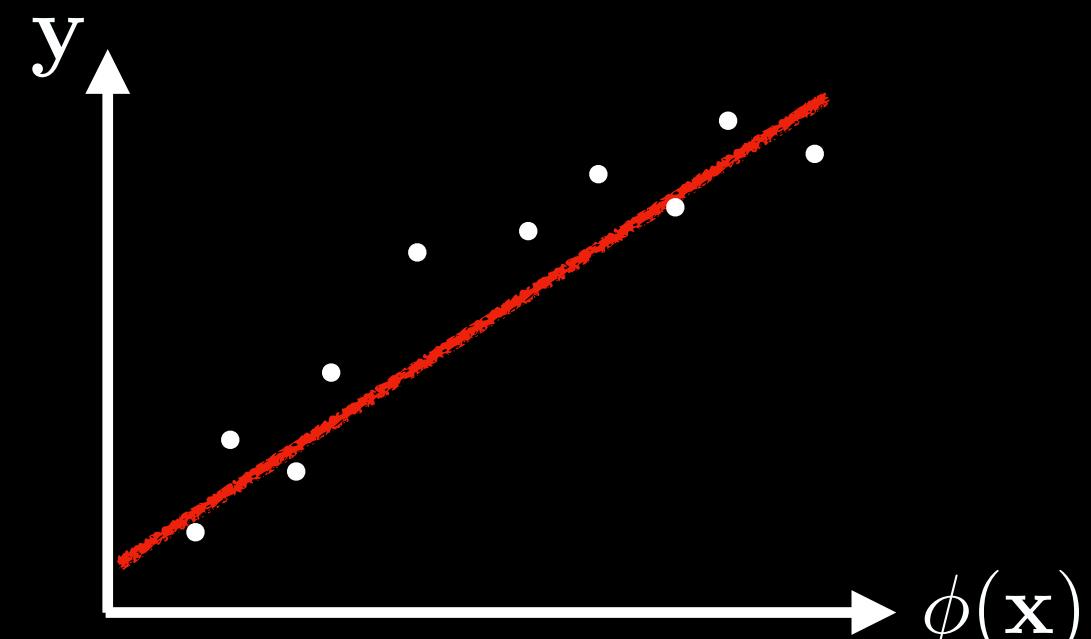


Linear Predictor

$$f_{\mathbf{w}}(x) = \mathbf{w} \cdot \phi(x)$$
$$\phi(x) = [1, x]$$
$$\phi(x) = [1, x, x^2, x^3]$$
$$\phi(x) = [1, x, \sin(3x)]$$

Higher dimensions

$$\hat{\mathbf{y}} = \mathbf{W} \phi(\mathbf{x})$$
$$\begin{matrix} | & = & | \\ m \times 1 & & m \times n & n \times 1 \end{matrix}$$



Linear Predictor

Matrix multiplication

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{x}$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

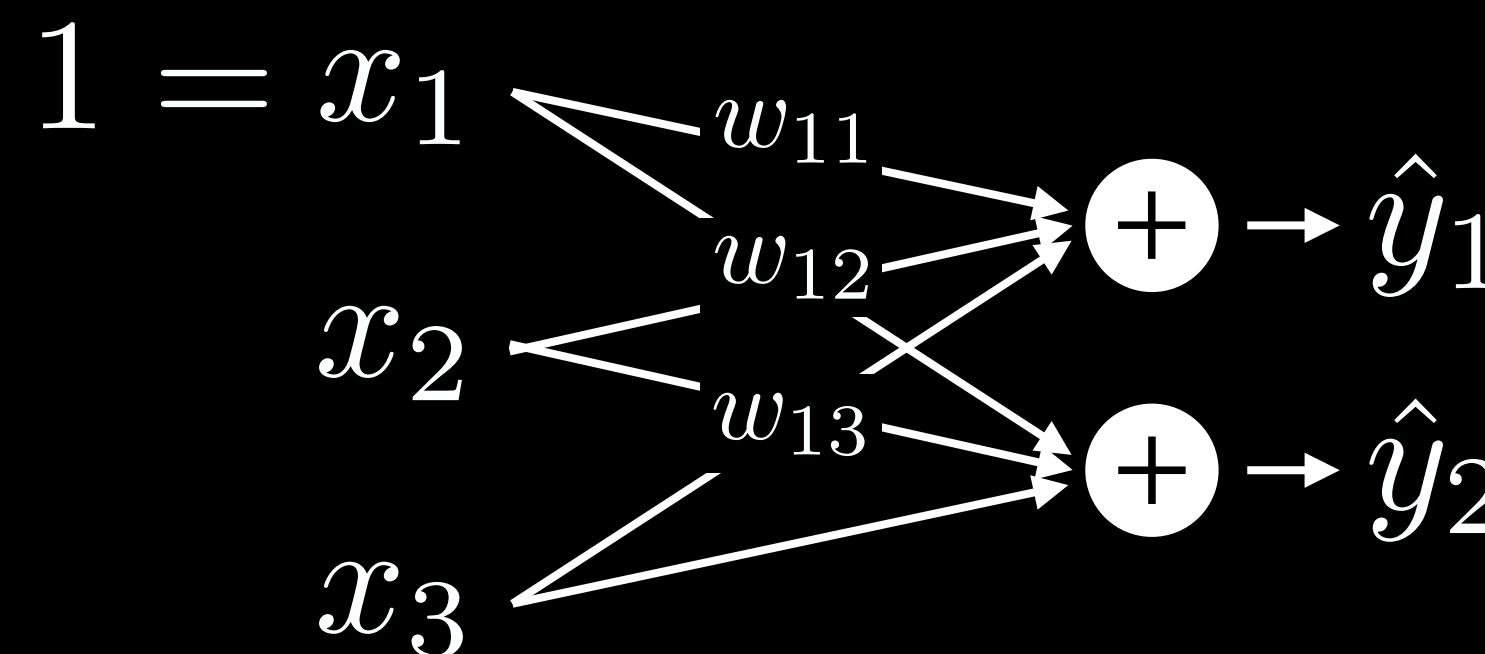
Index notation

$$\hat{y}_i = \sum_{j=1}^n w_{ij} x_j$$

Define features such that
first component accounts for bias

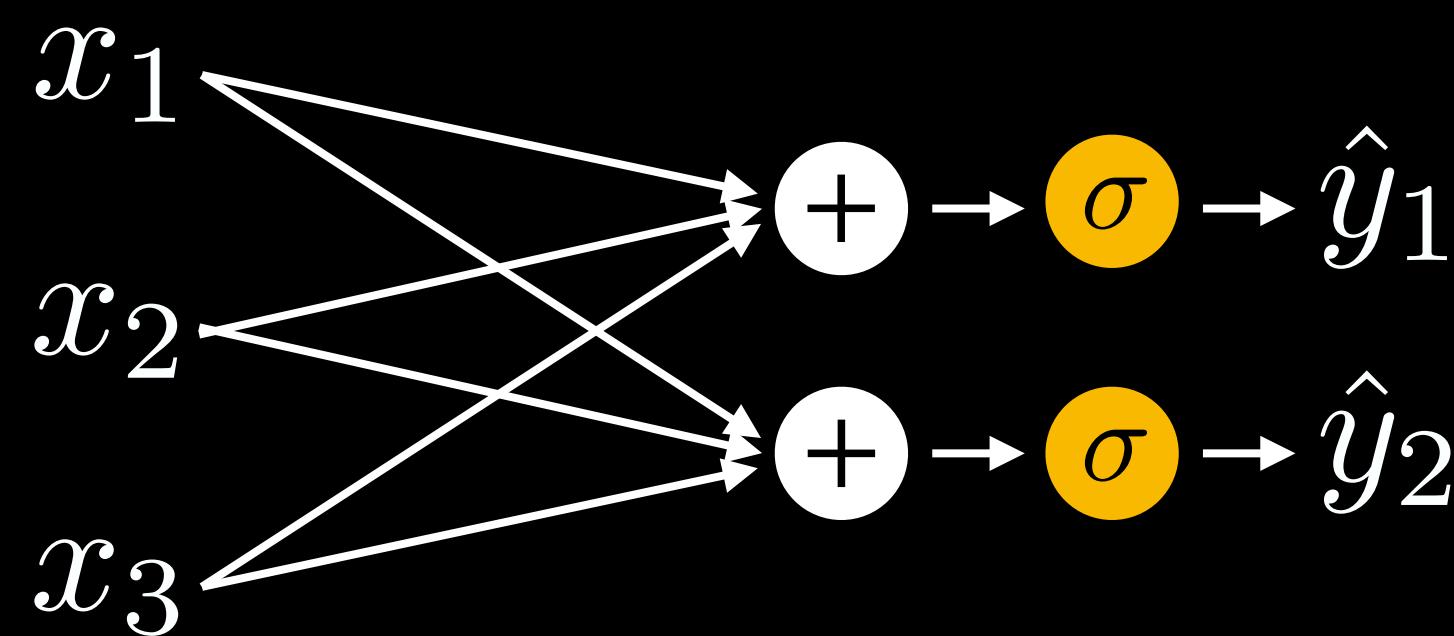


Network representation



Nonlinear Predictor

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{x})$$

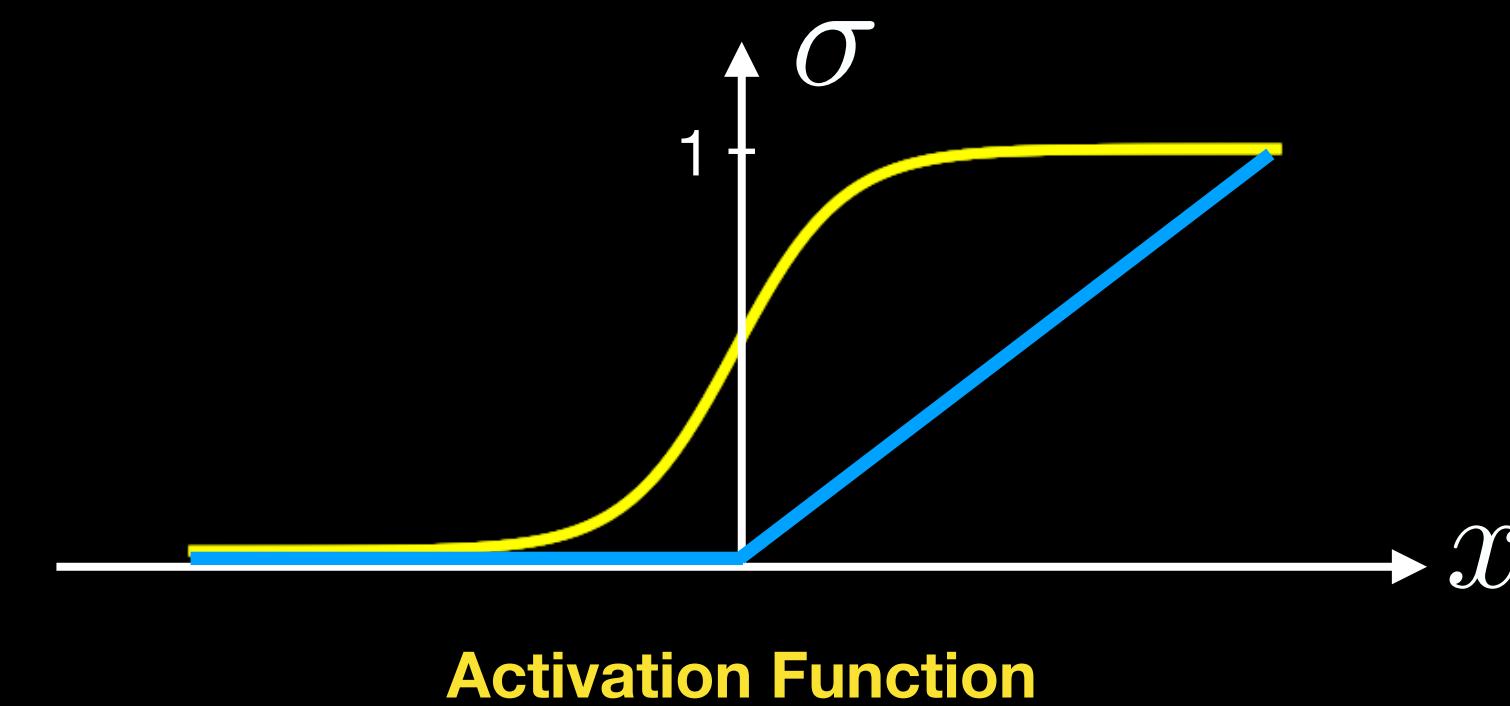
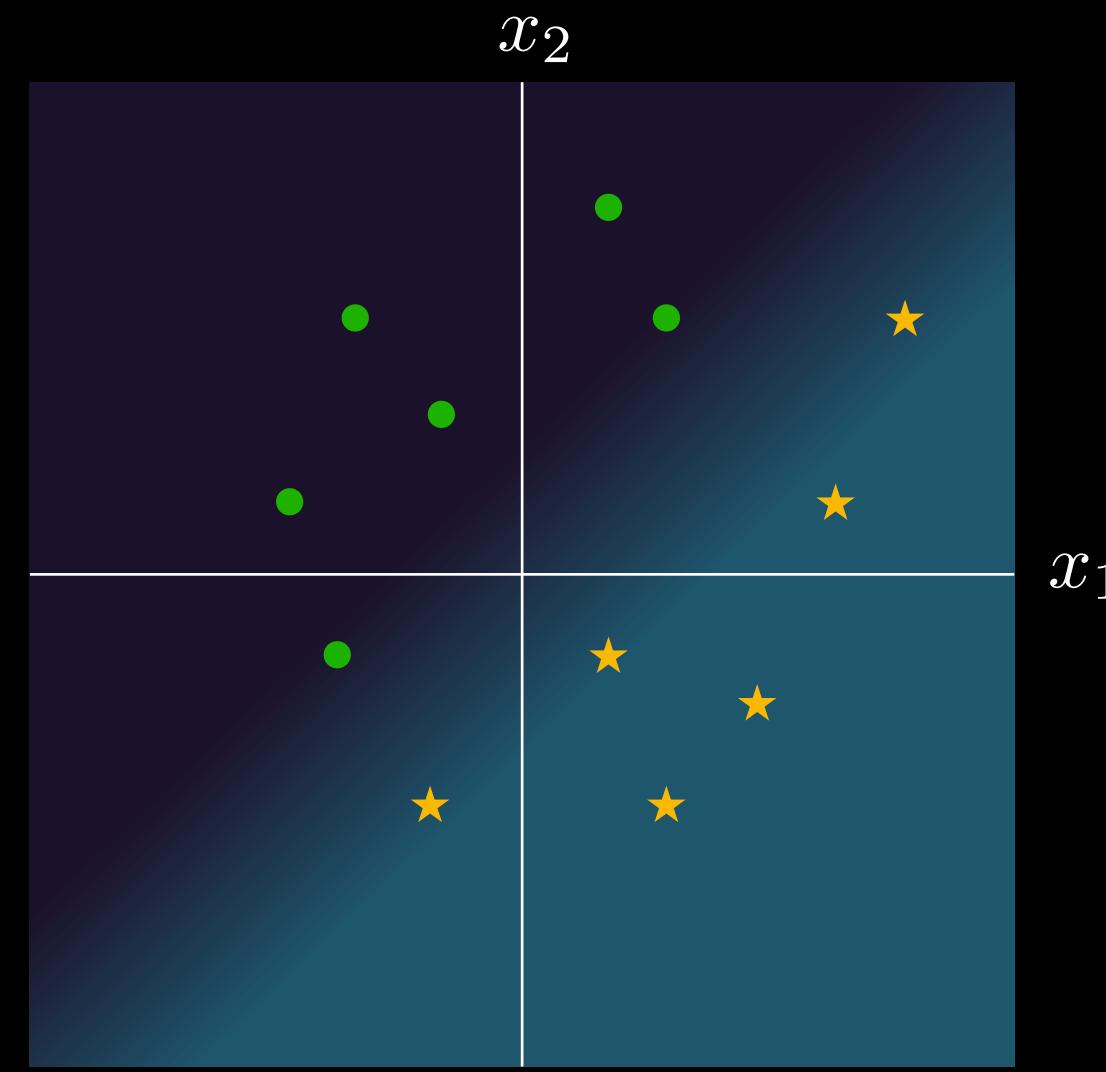


Logistic function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

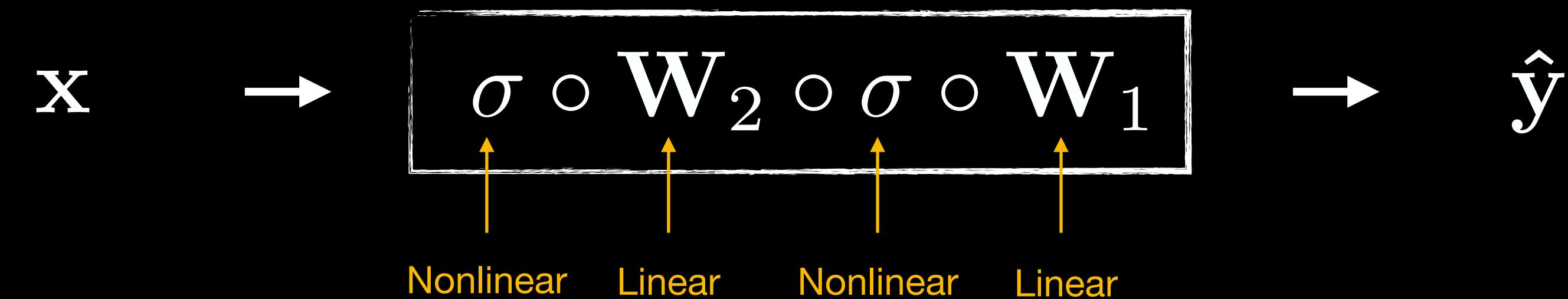
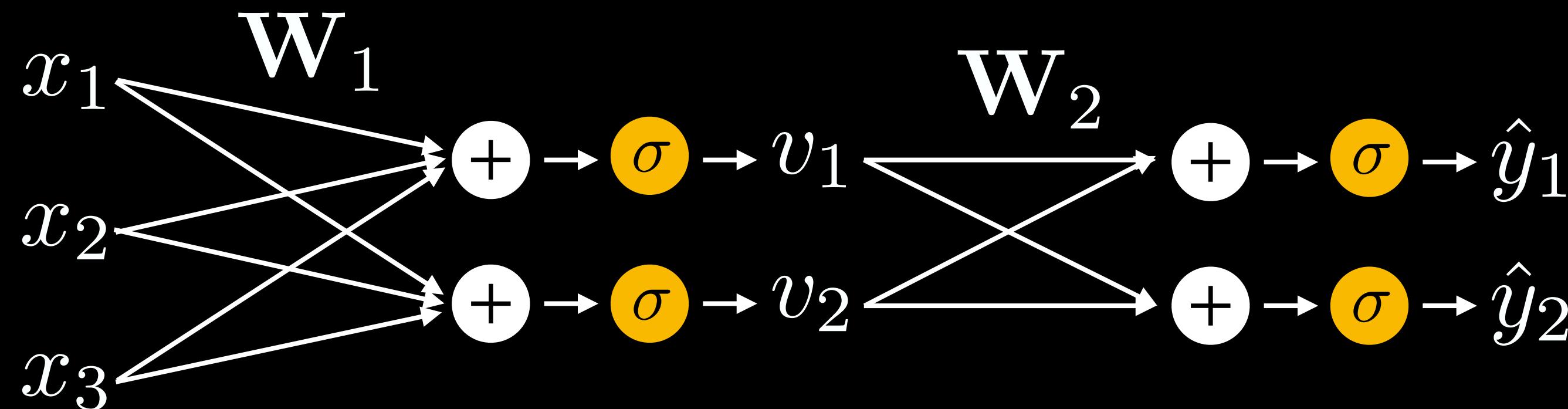
ReLU:

$$\sigma(x) = xH(x)$$



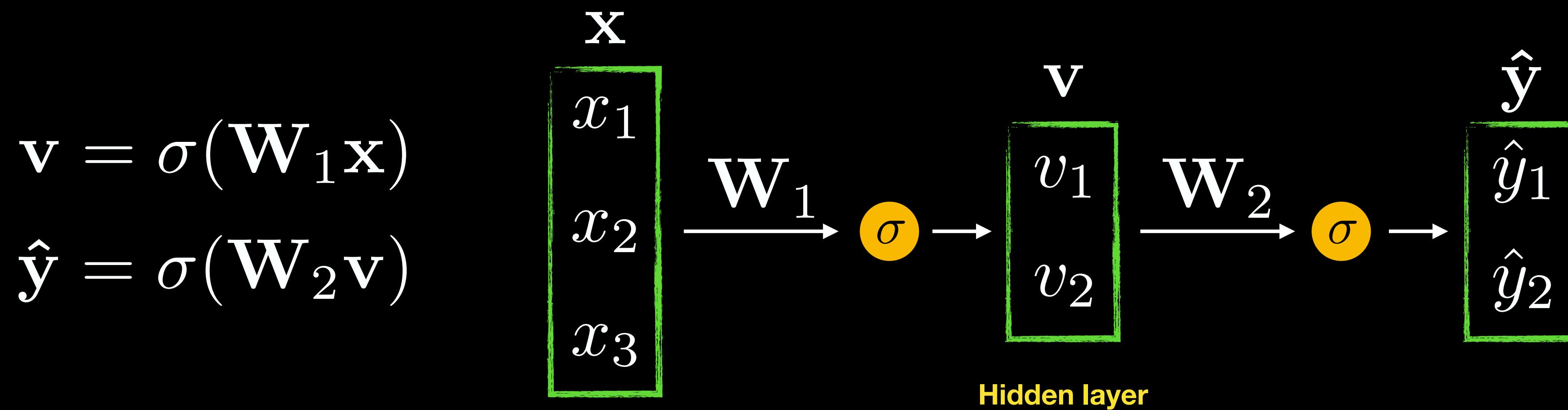
Neural network

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$



Neural network

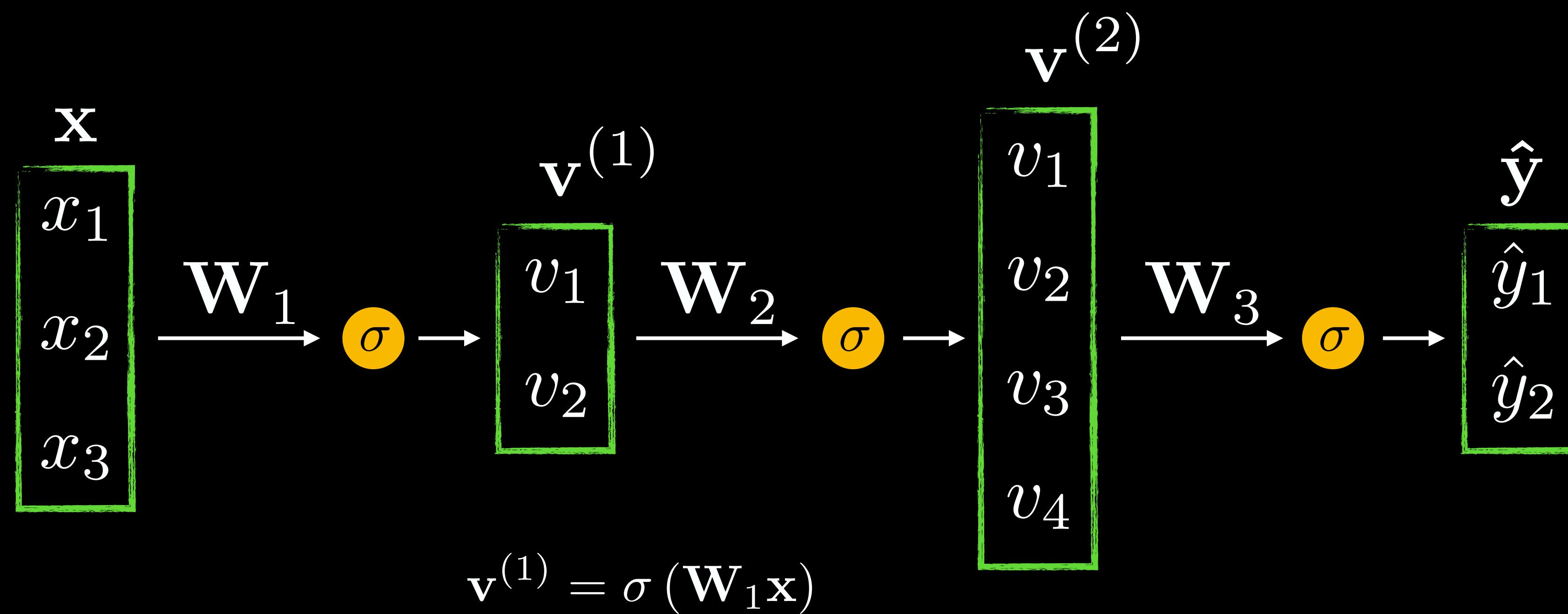
$$\hat{\mathbf{y}} = \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$



Can be interpreted
as a learned $\phi(\mathbf{x})$

Deep network

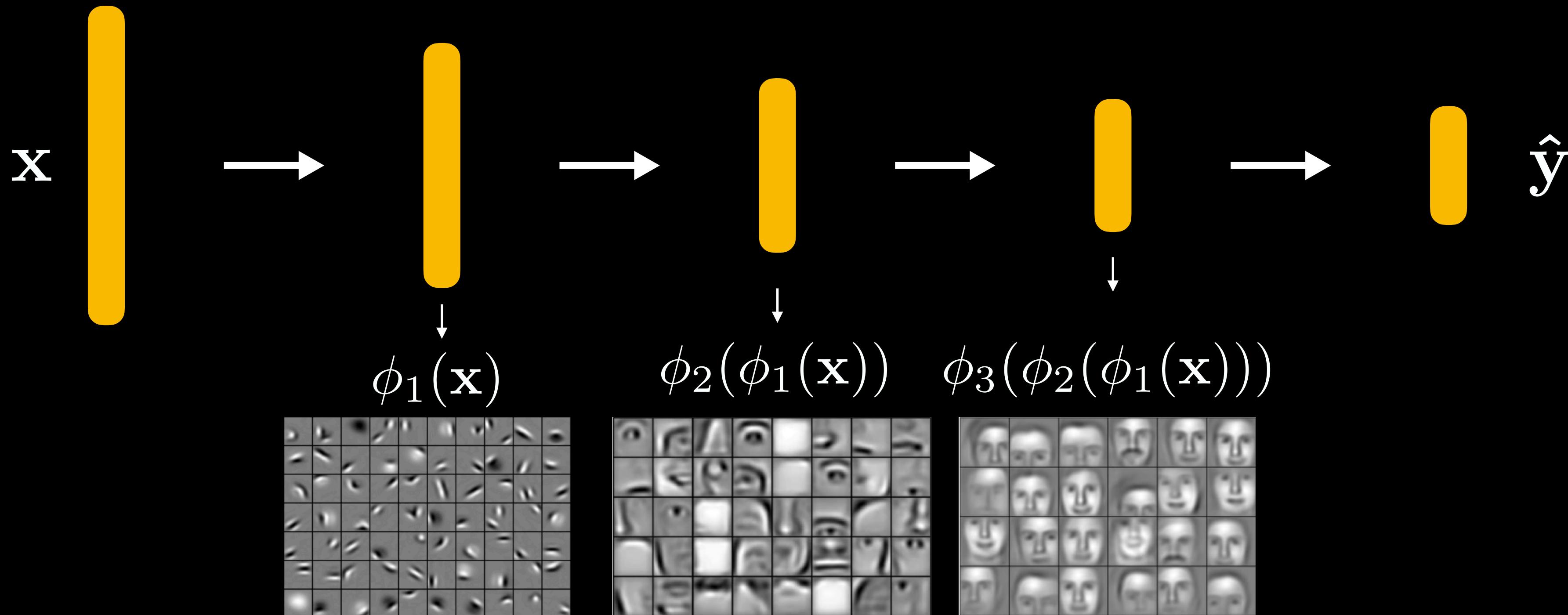
$$\hat{\mathbf{y}} = \sigma(\mathbf{W}_3 \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})))$$



$$\mathbf{v}^{(2)} = \sigma(\mathbf{W}_2 \mathbf{v}^{(1)})$$

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}_3 \mathbf{v}^{(2)})$$

Why deep learning?



Feature learning!

Loss function

$$f_{\mathbf{W}_1 \mathbf{W}_2}(\mathbf{x}) = \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \mathbf{W}_1, \mathbf{W}_2) = \|f_{\mathbf{W}_1 \mathbf{W}_2}(\mathbf{x}) - \mathbf{y}\|^2$$

Stochastic gradient descent update

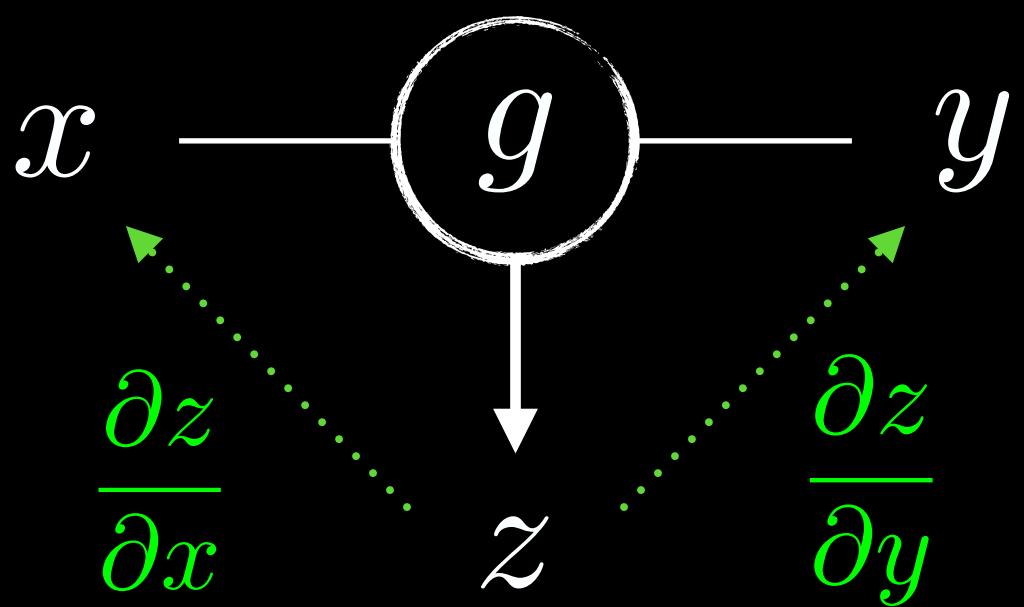
$$\mathbf{W}_1 \leftarrow \mathbf{W}_1 - \eta \nabla_{\mathbf{W}_1} \mathcal{L}(\mathbf{x}, \mathbf{y}, \mathbf{W}_1, \mathbf{W}_2)$$

$$\mathbf{W}_2 \leftarrow \mathbf{W}_2 - \eta \nabla_{\mathbf{W}_2} \mathcal{L}(\mathbf{x}, \mathbf{y}, \mathbf{W}_1, \mathbf{W}_2)$$

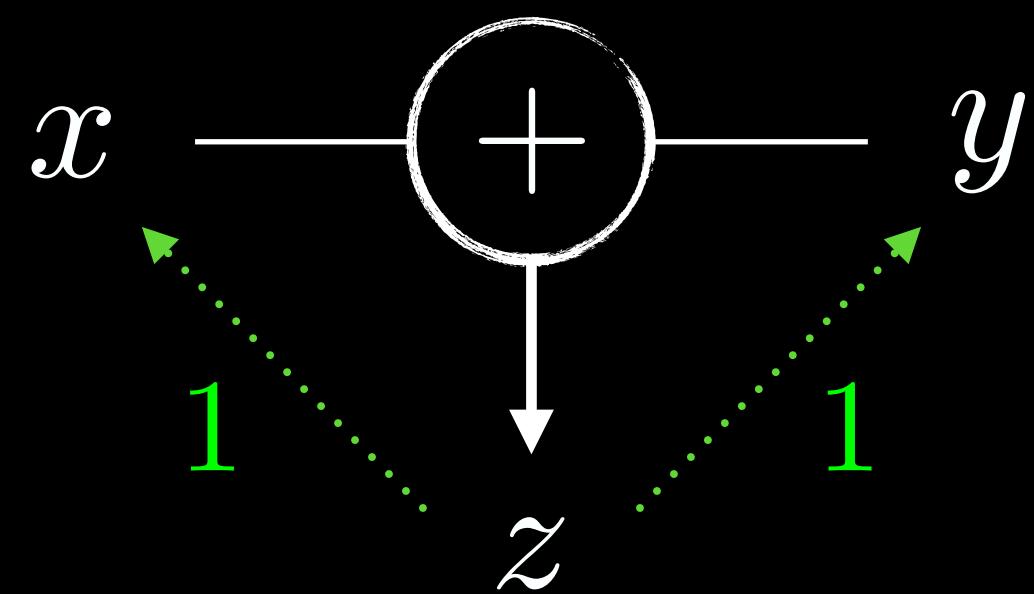
How do we calculate the gradients?

Computation graphs

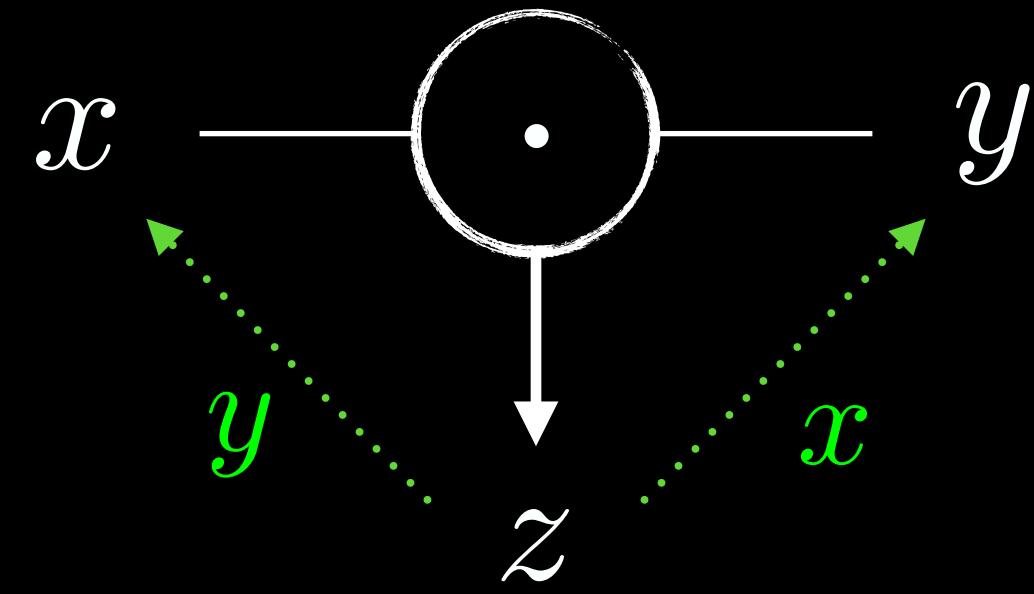
$$z = g(x, y)$$



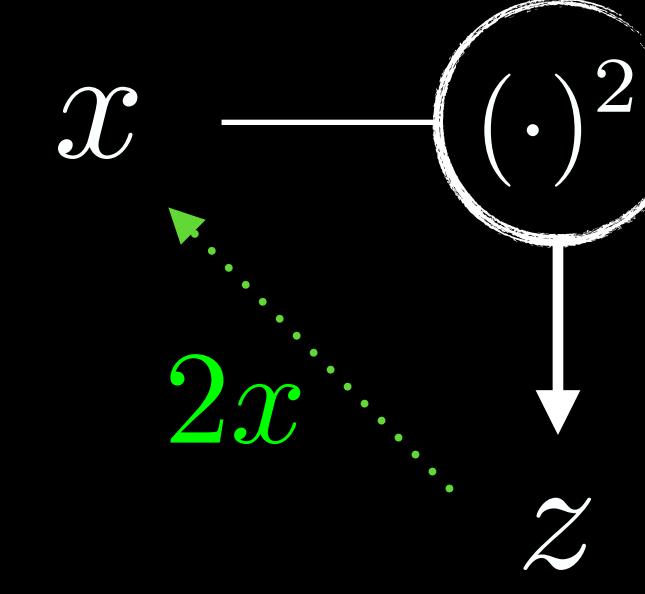
$$z = x + y$$



$$z = xy$$



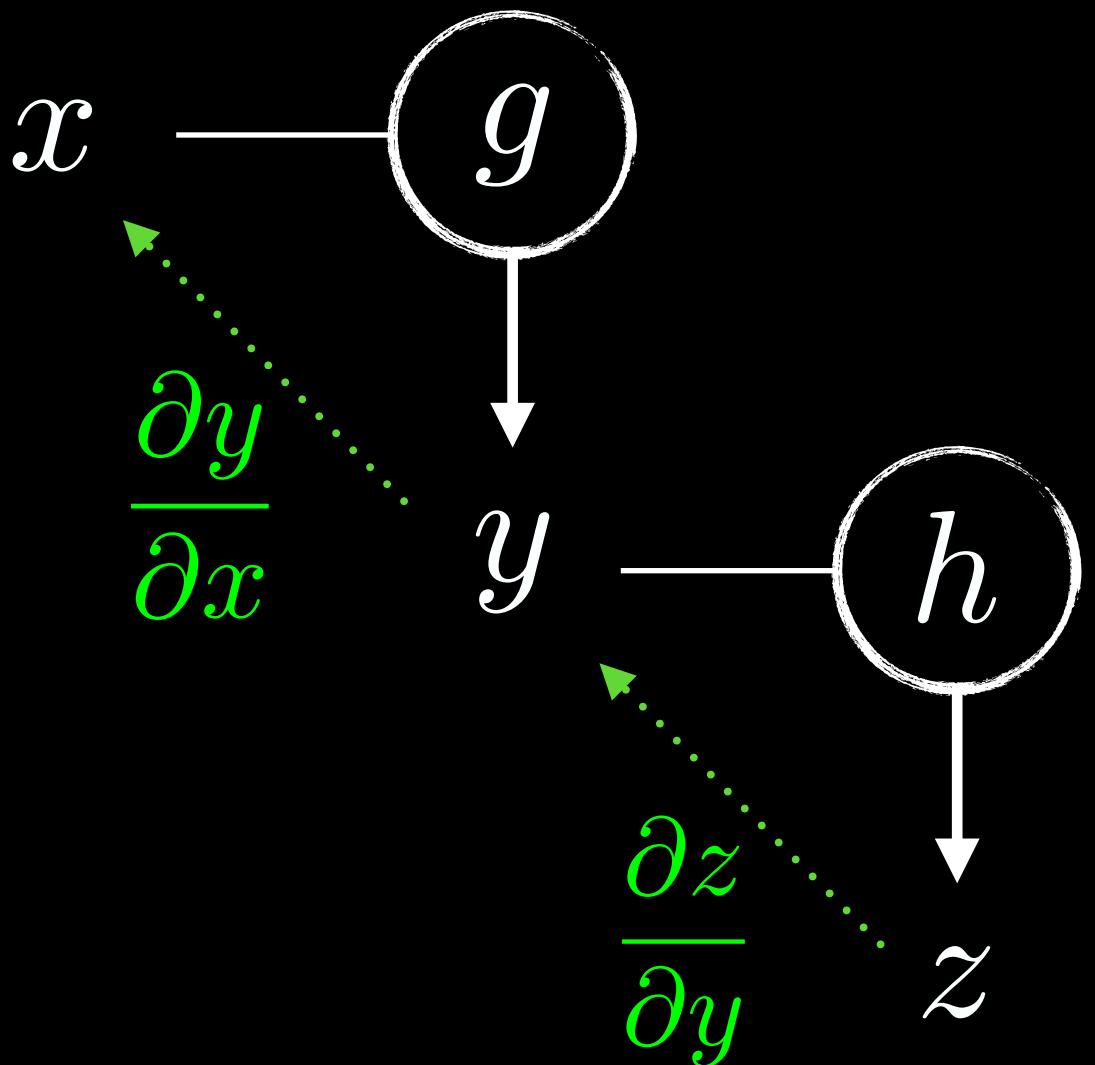
$$z = x^2$$



Chain rule

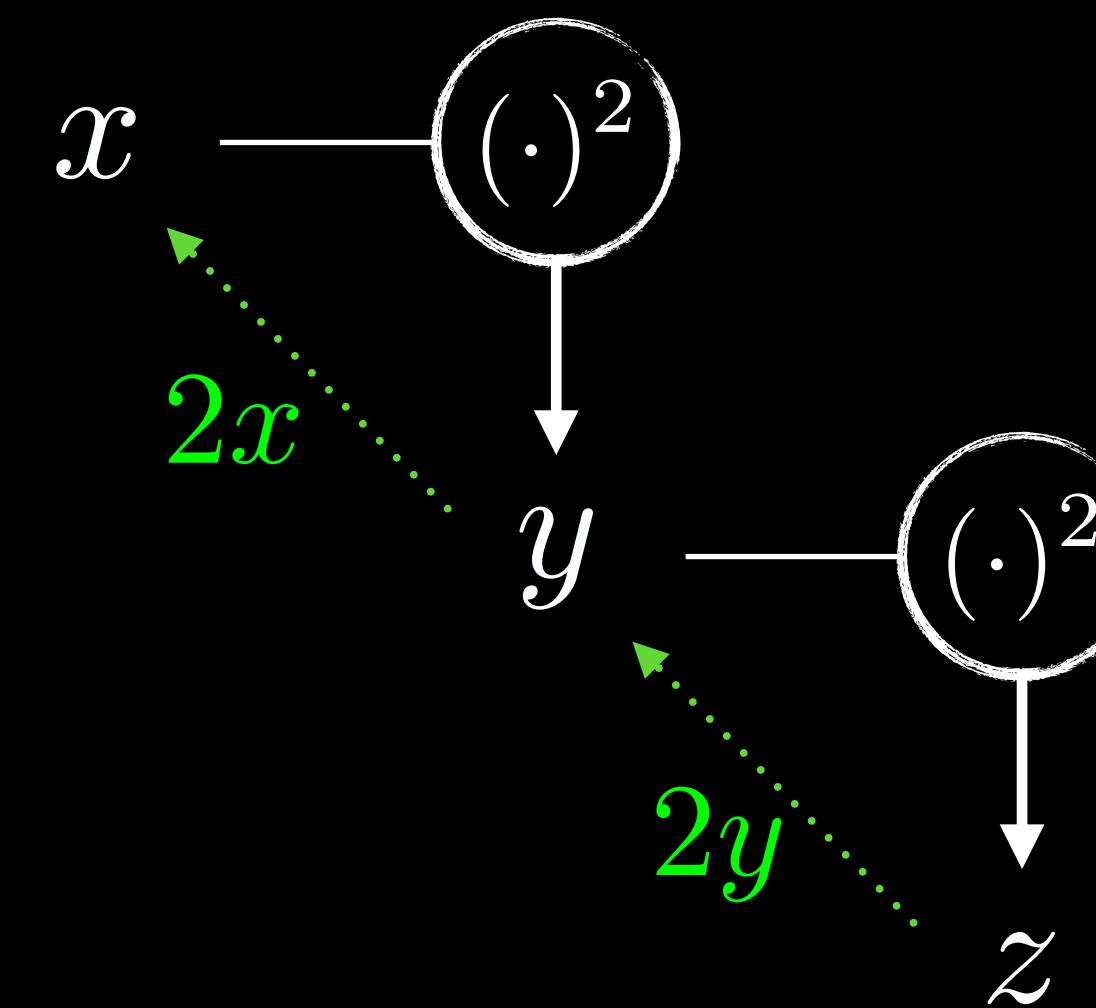
$$y = g(x)$$

$$z = h(y)$$



$$y = x^2$$

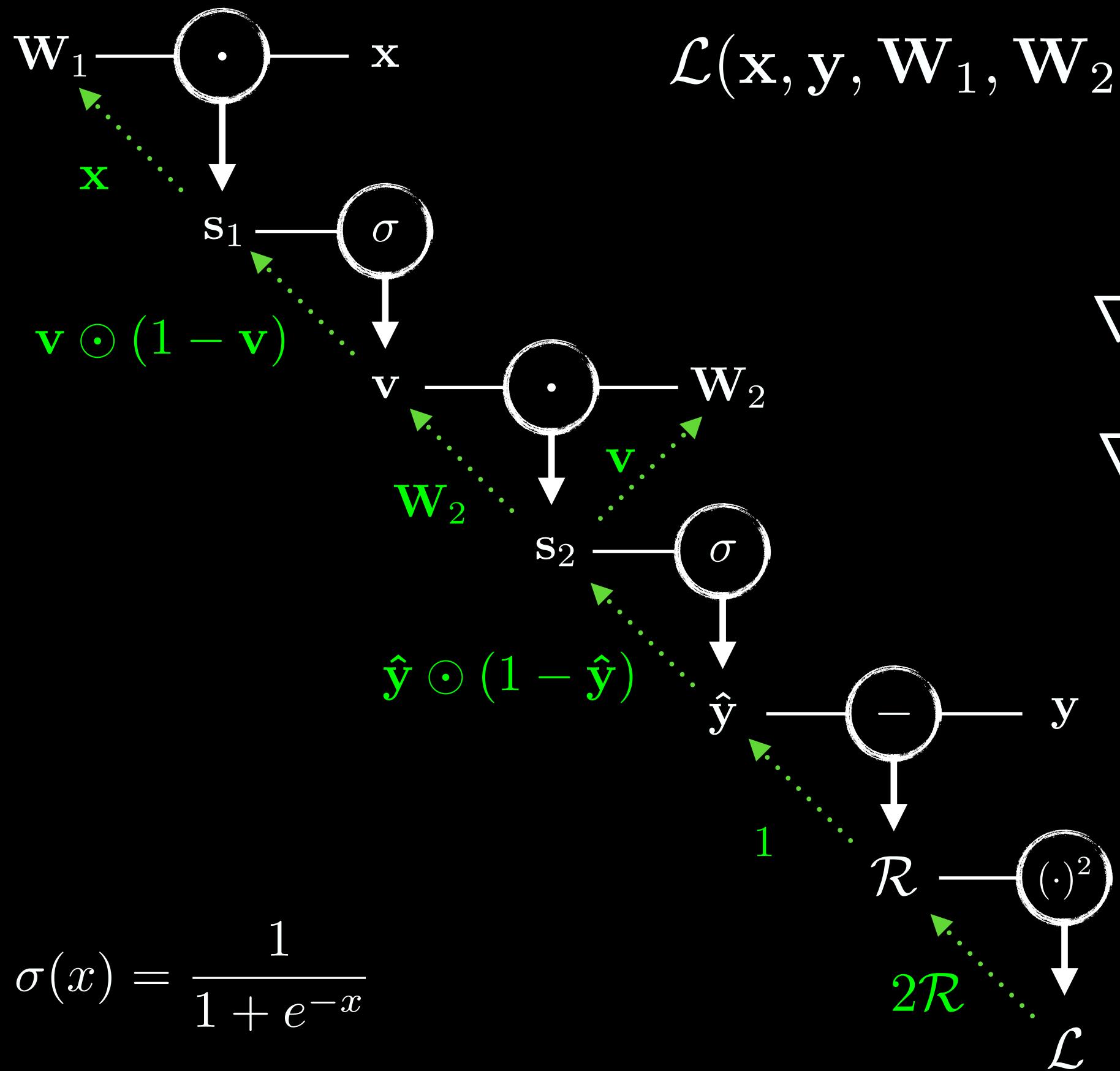
$$z = y^2$$



$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$$\frac{\partial z}{\partial x} = 4xy = 4x^3$$

Backpropagation



$$\mathcal{L}(x, y, W_1, W_2) = \|\sigma(W_2 \sigma(W_1 x)) - y\|^2$$

$$\nabla_{W_1} \mathcal{L} = 2W_2^\top \mathcal{R} \odot \hat{y} \odot (1 - \hat{y}) \odot v \odot (1 - v)x^\top$$

$$\nabla_{W_2} \mathcal{L} = 2\mathcal{R} \odot \hat{y} \odot (1 - \hat{y})v^\top$$

assuming $\sigma(x) = \frac{1}{1 + e^{-x}}$

Optimization

Training loss

$$\mathcal{L}_{\text{train}}(\mathbf{W}_1, \mathbf{W}_2) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}_{\text{train}}} \mathcal{L}(\mathbf{x}, \mathbf{y}, \mathbf{W}_1, \mathbf{W}_2)$$

Objective

$$\hat{\mathbf{W}}_1, \hat{\mathbf{W}}_2 = \operatorname{argmin}_{\mathbf{W}_1, \mathbf{W}_2} \mathcal{L}_{\text{train}}(\mathbf{W}_1, \mathbf{W}_2)$$

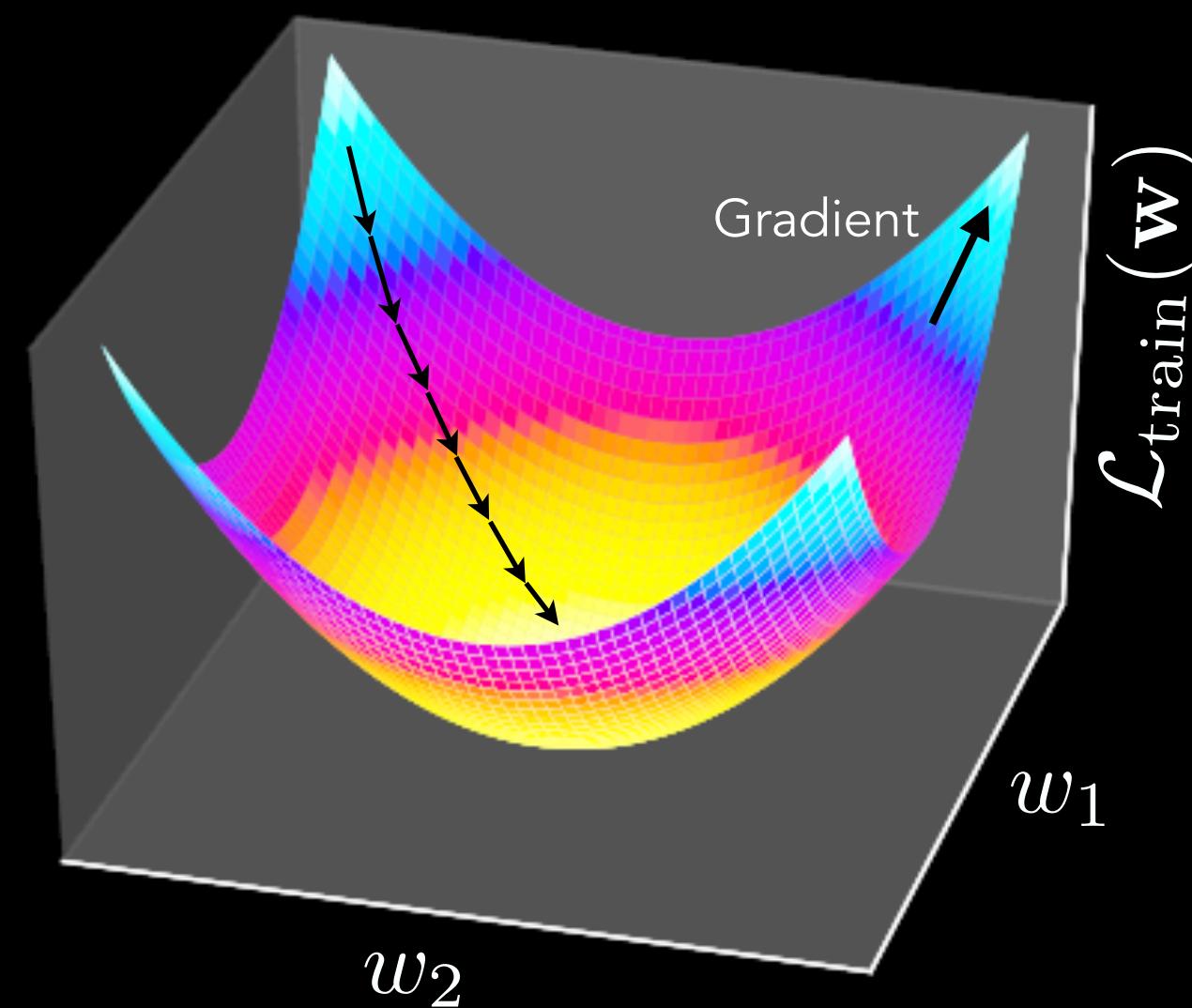
Optimal predictor

$$f_{\hat{\mathbf{W}}_1 \hat{\mathbf{W}}_2}(\mathbf{x}) = \sigma \left(\hat{\mathbf{W}}_2 \sigma \left(\hat{\mathbf{W}}_1 \mathbf{x} \right) \right)$$

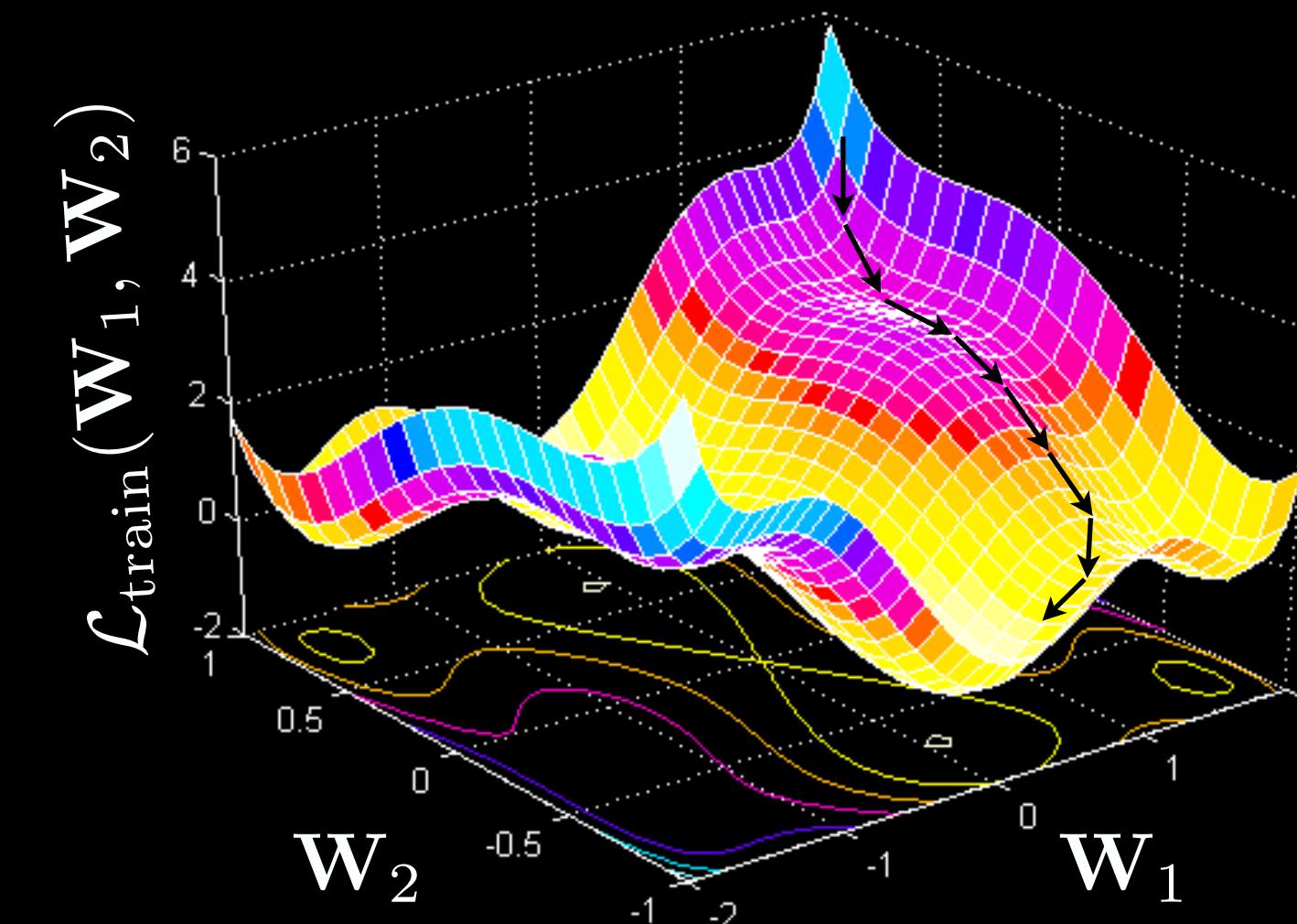
Non-convexity

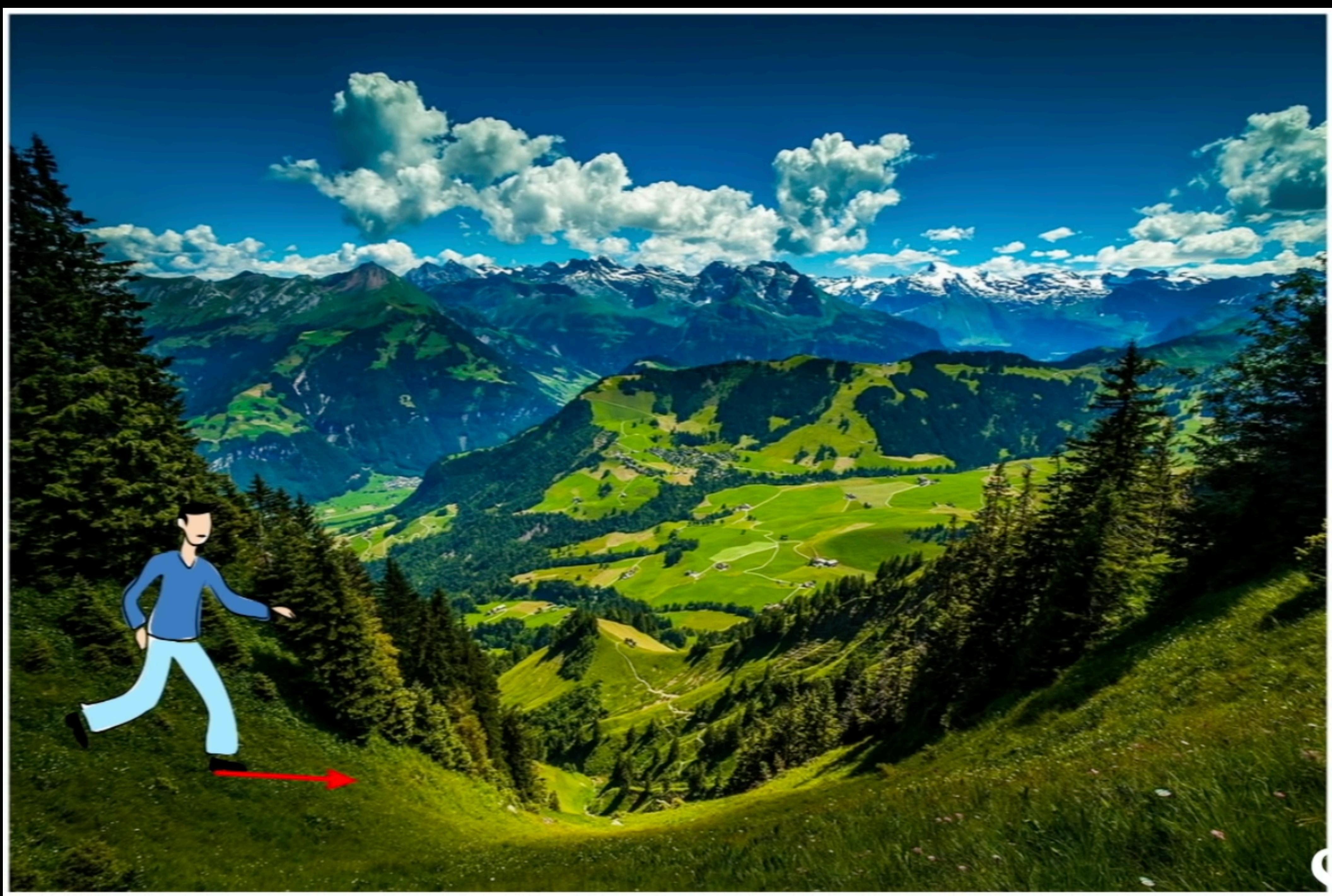
$$\hat{\mathbf{W}}_1, \hat{\mathbf{W}}_2 = \underset{\mathbf{W}_1, \mathbf{W}_2}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\mathbf{W}_1, \mathbf{W}_2)$$

Linear predictor loss



Neural network loss

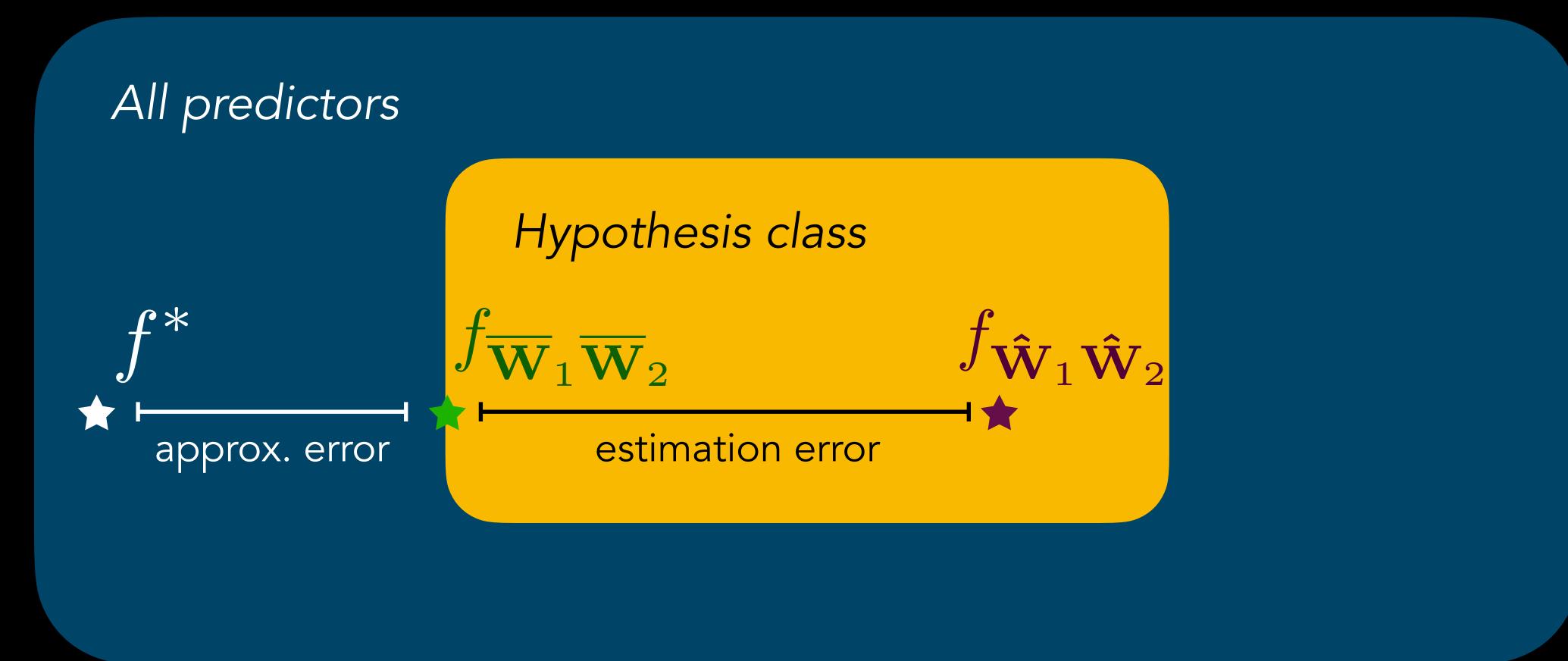




Hypothesis Class

$$f_{\mathbf{W}_1 \mathbf{W}_2}(\mathbf{x}) = \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$

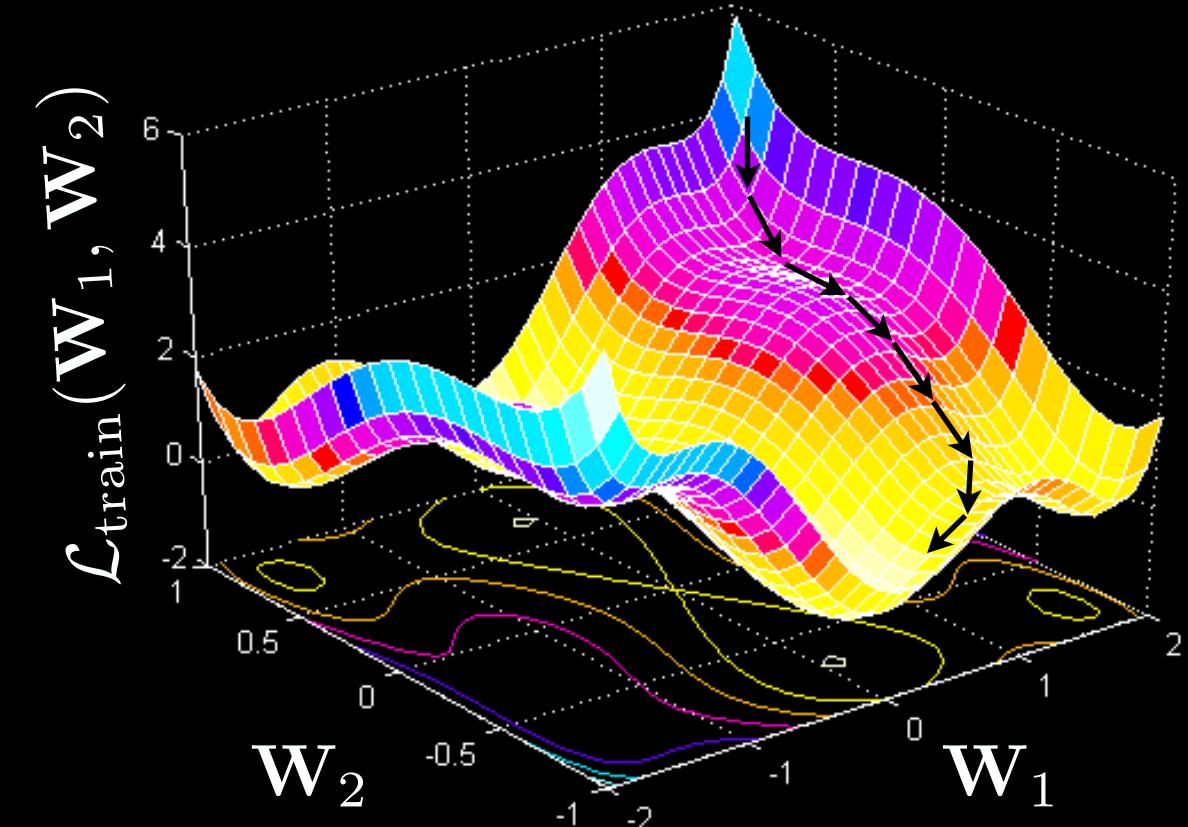
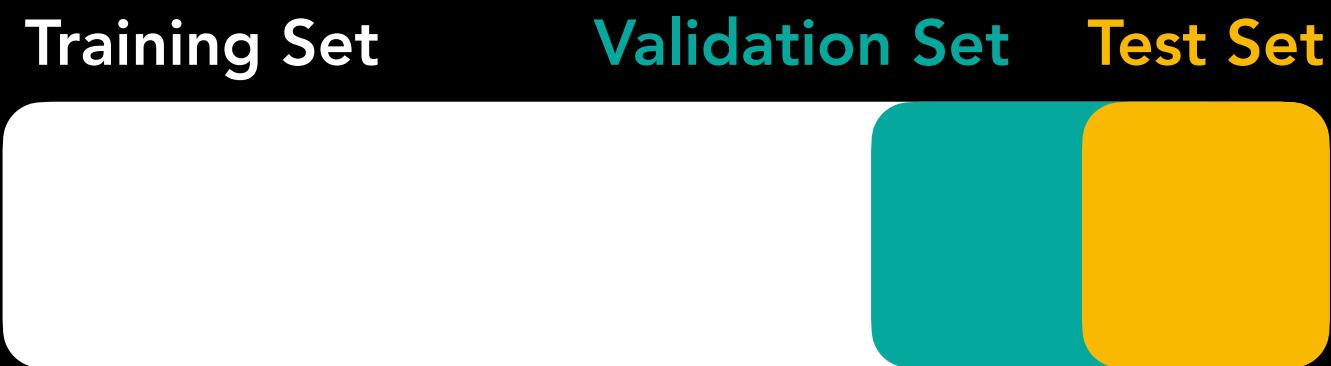
$$\mathcal{F} = \{ f_{\mathbf{W}_1 \mathbf{W}_2}(\mathbf{x}) \mid \mathbf{W}_1 \in \mathbb{R}^{k \times n}, \mathbf{W}_2 \in \mathbb{R}^{m \times k} \}$$



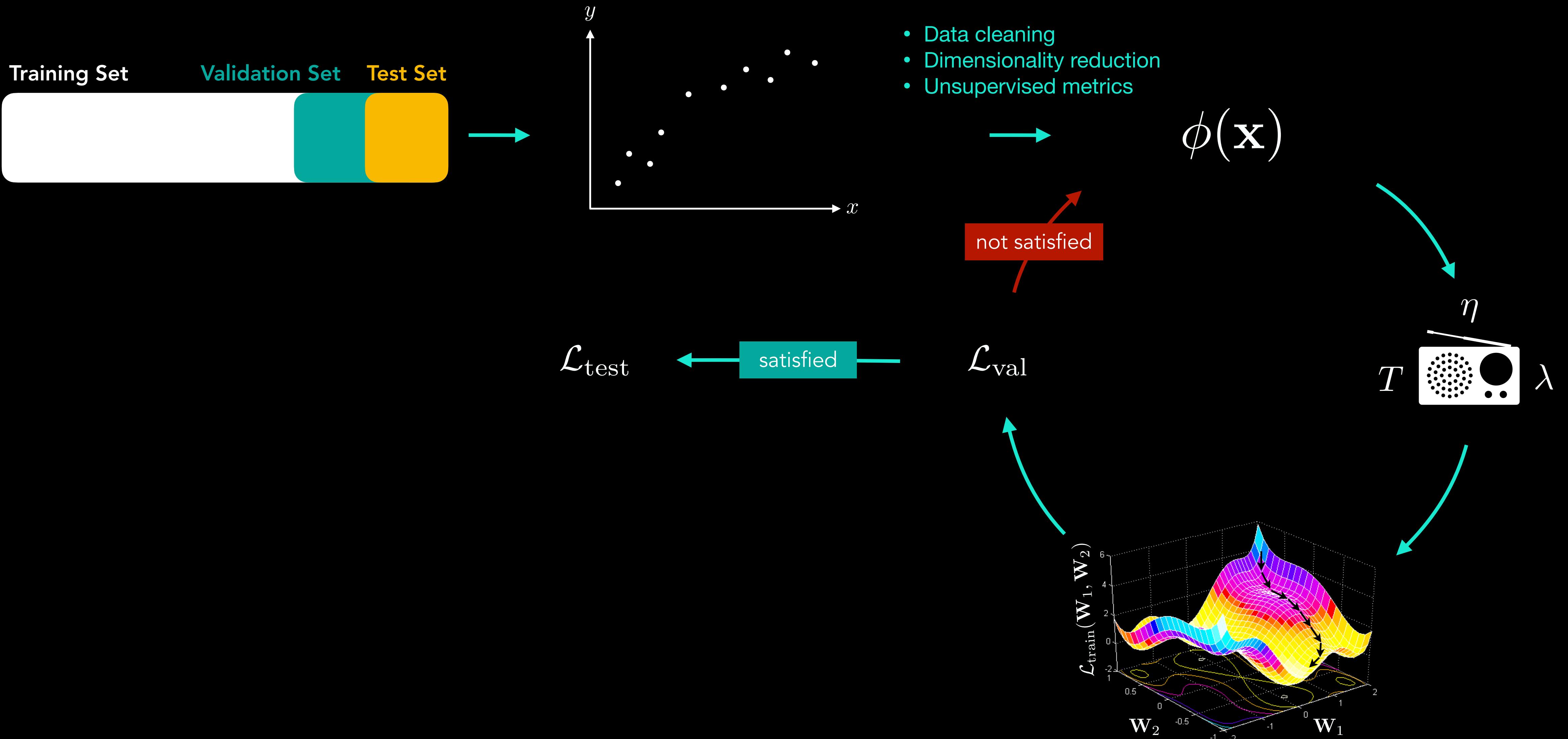
Hyperparameters

How do you train a deep network?

- Use many hidden layers for abstraction
- Use adaptive time steps
- Use hyper-parameter optimization



The ML workflow

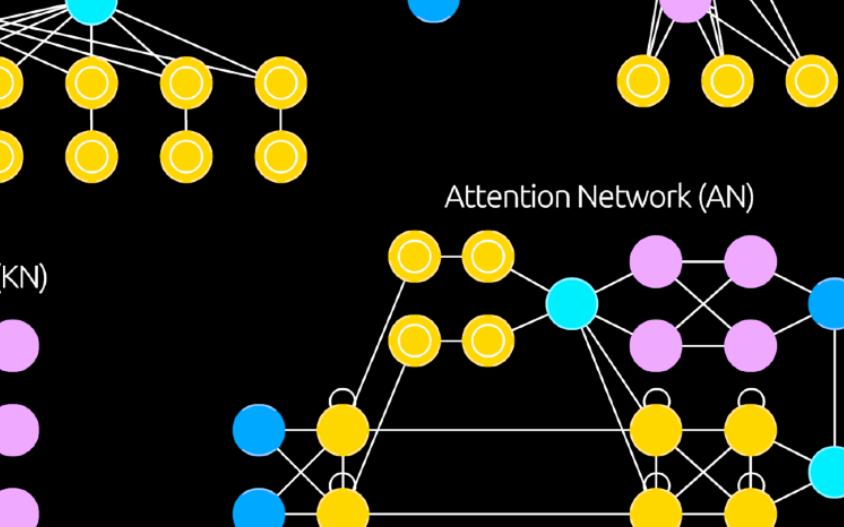
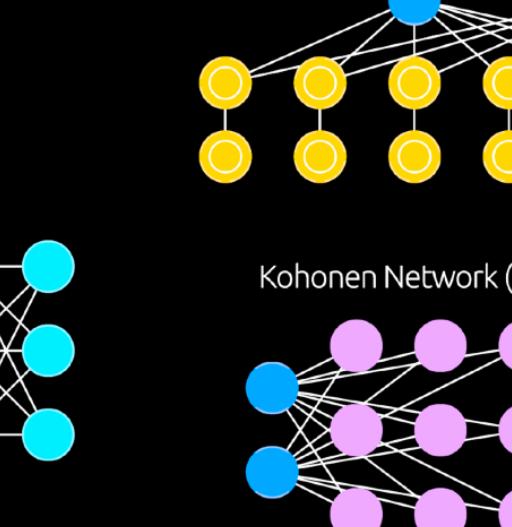
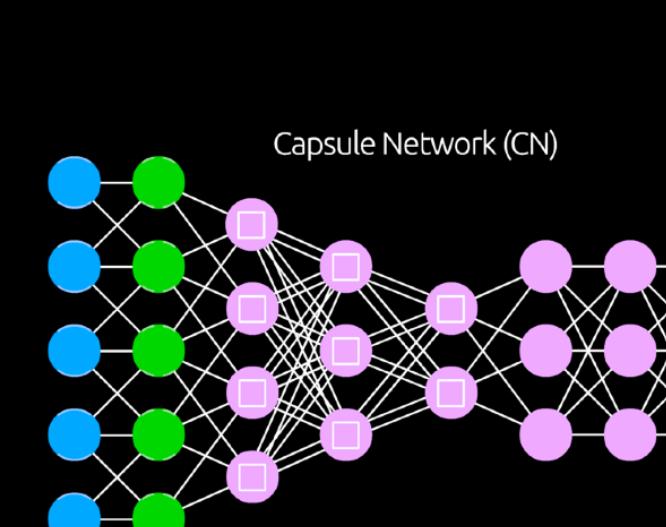
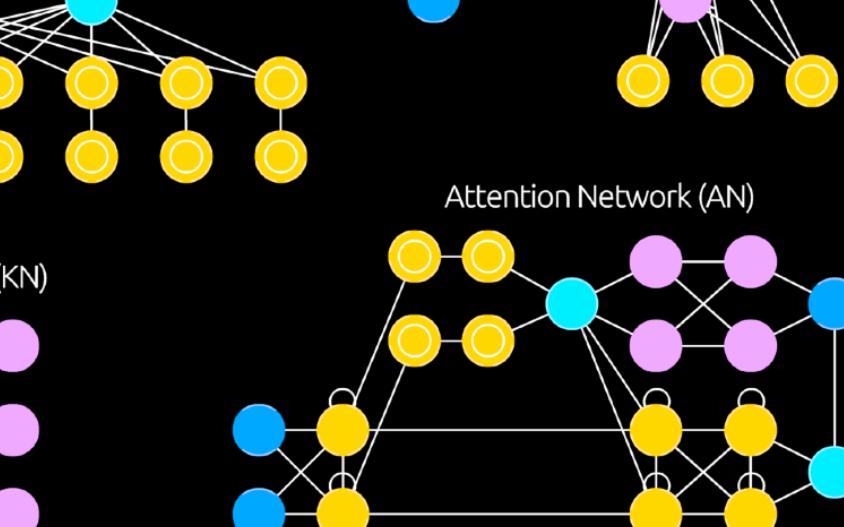
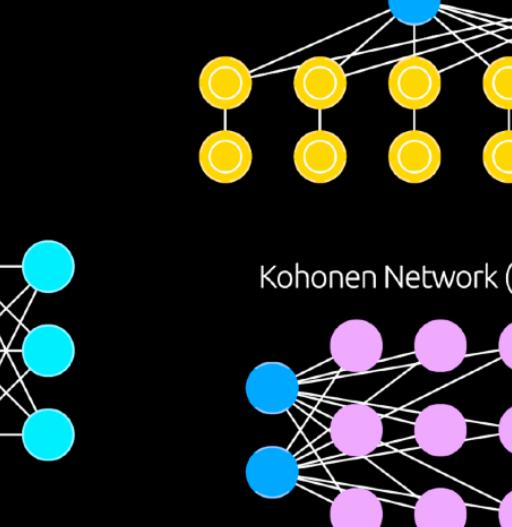
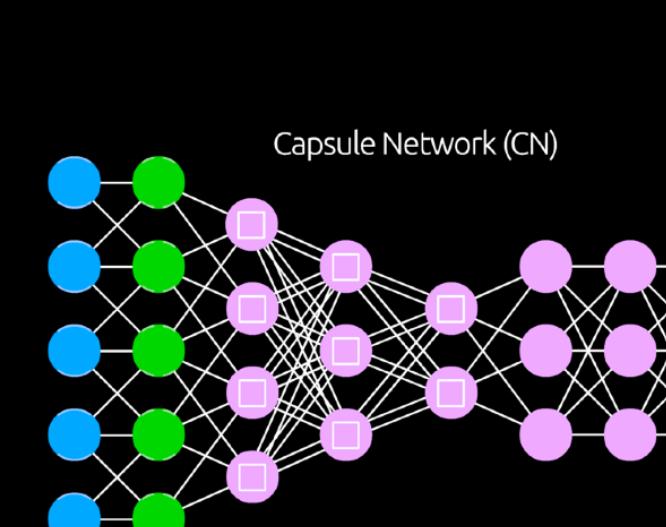
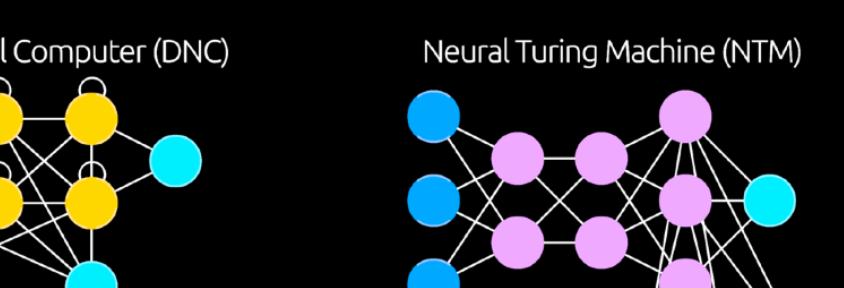
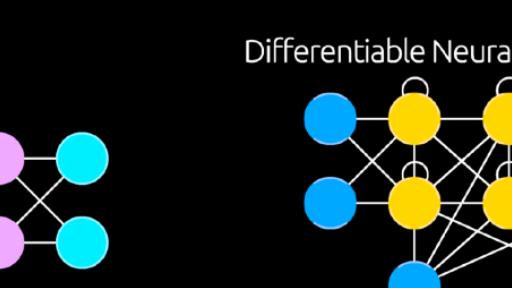
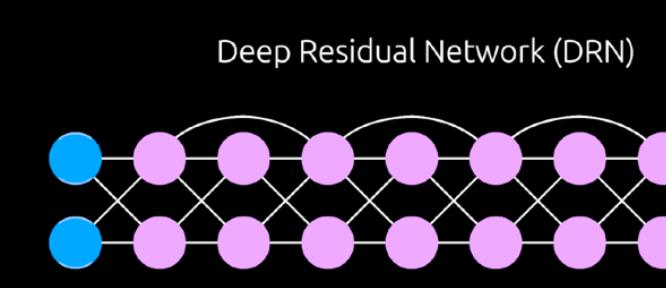
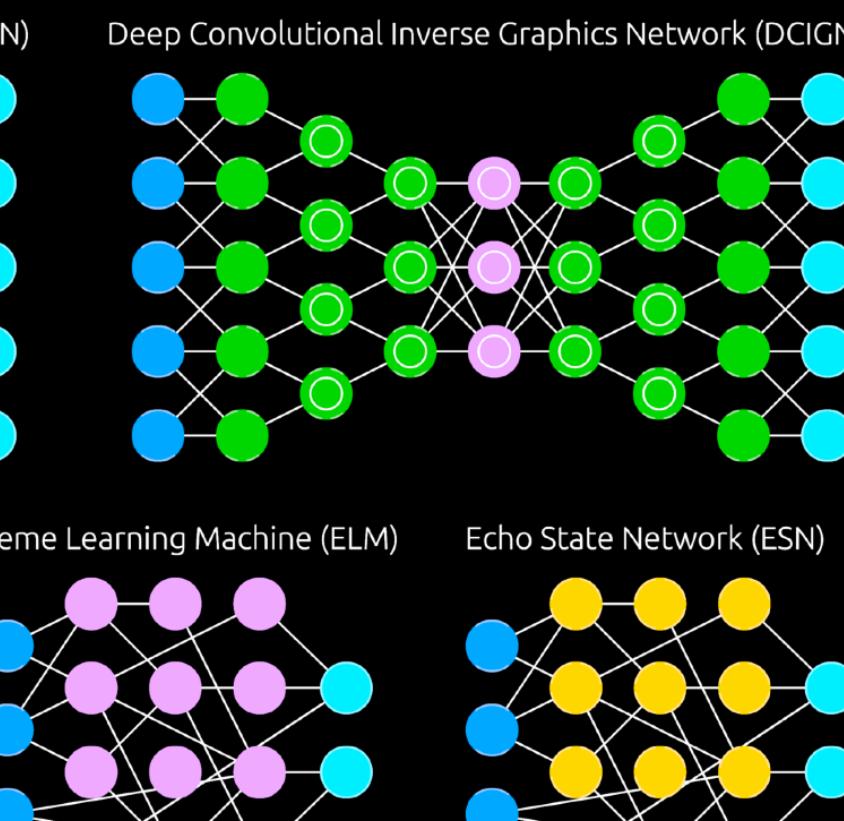
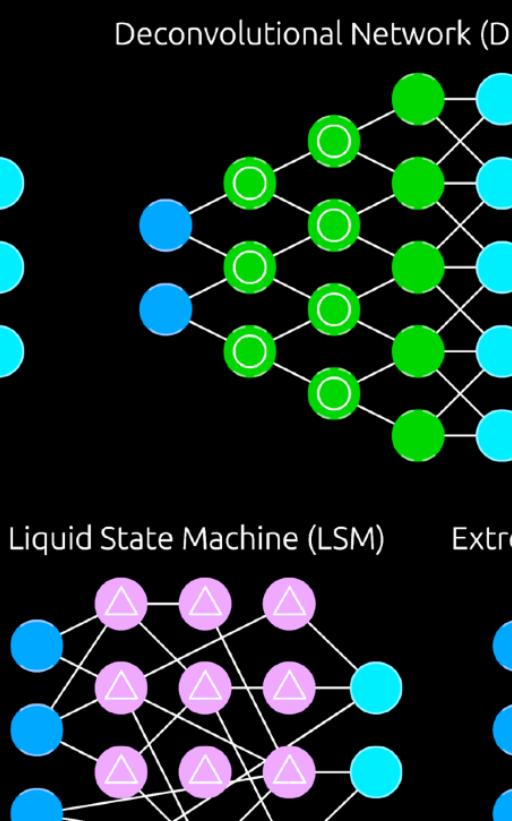
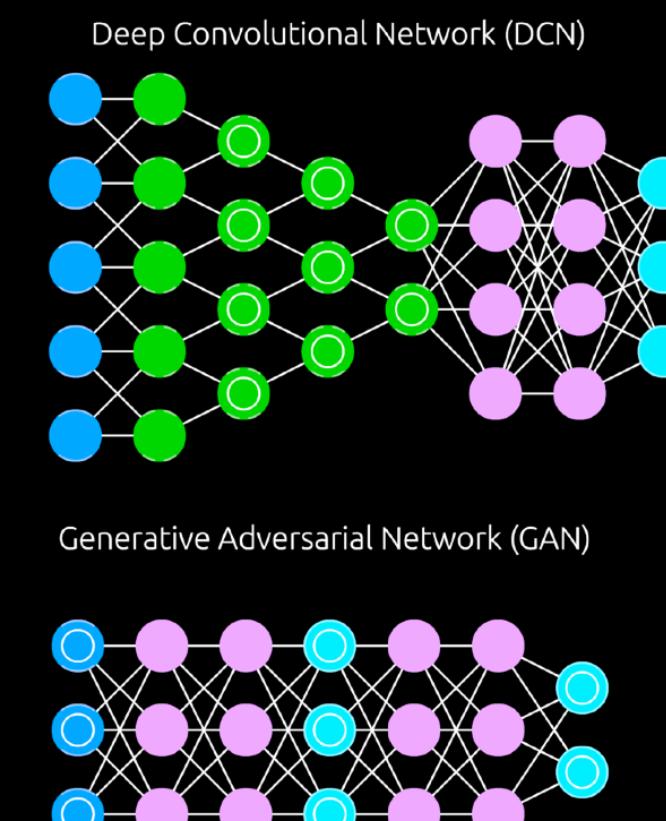
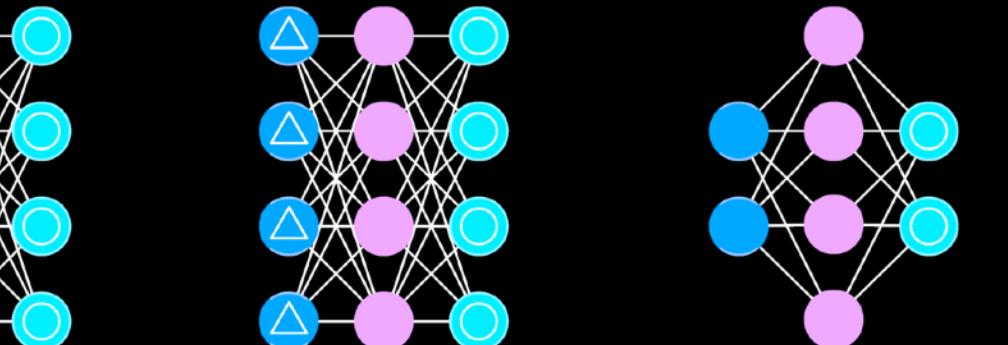
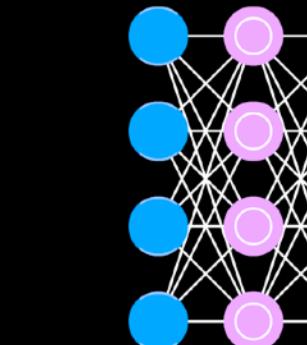
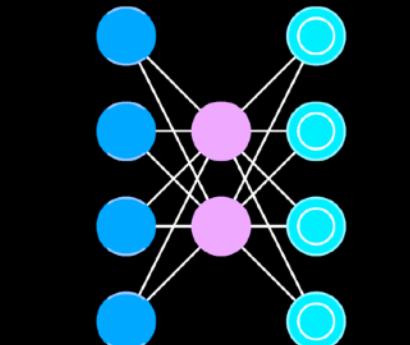
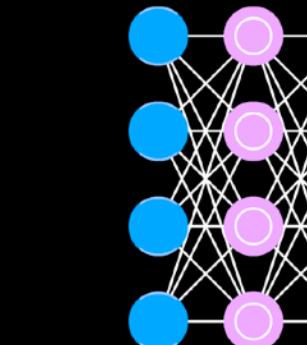
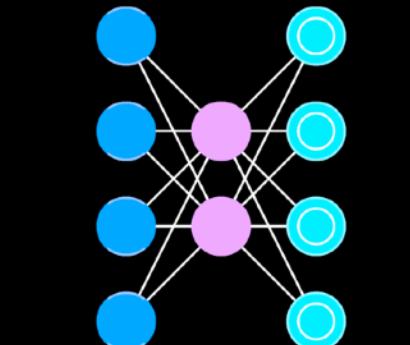
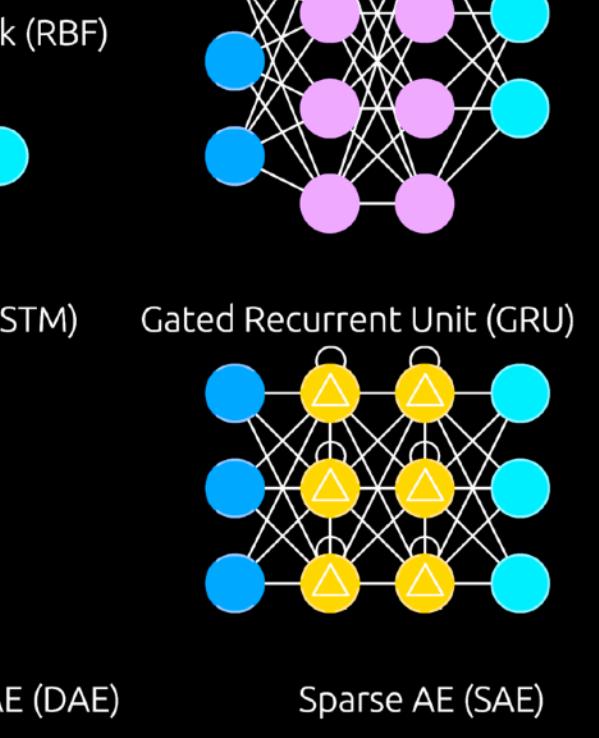
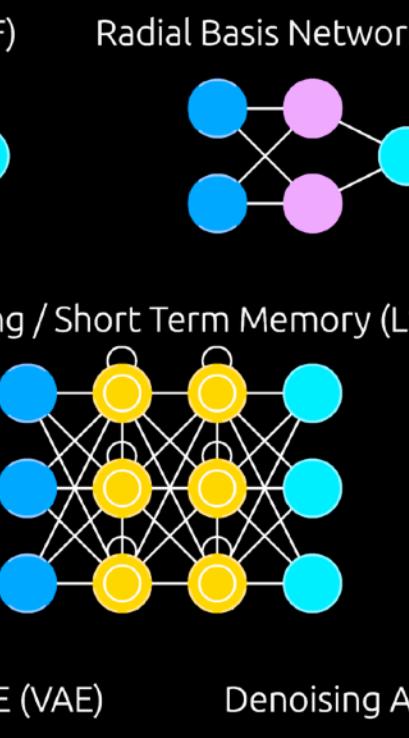
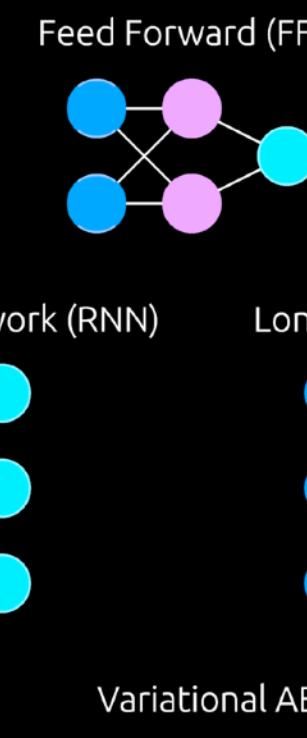
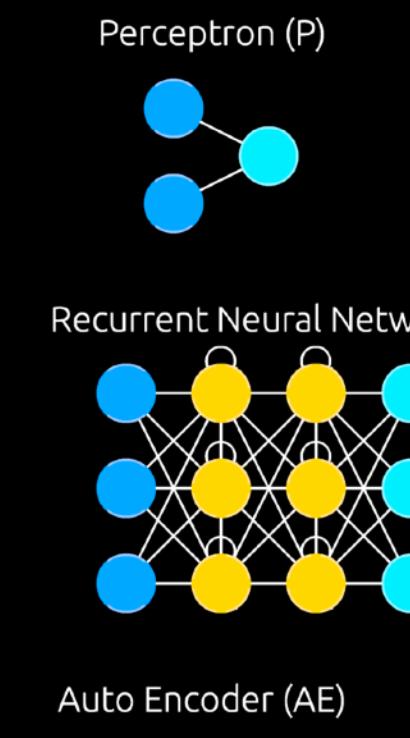
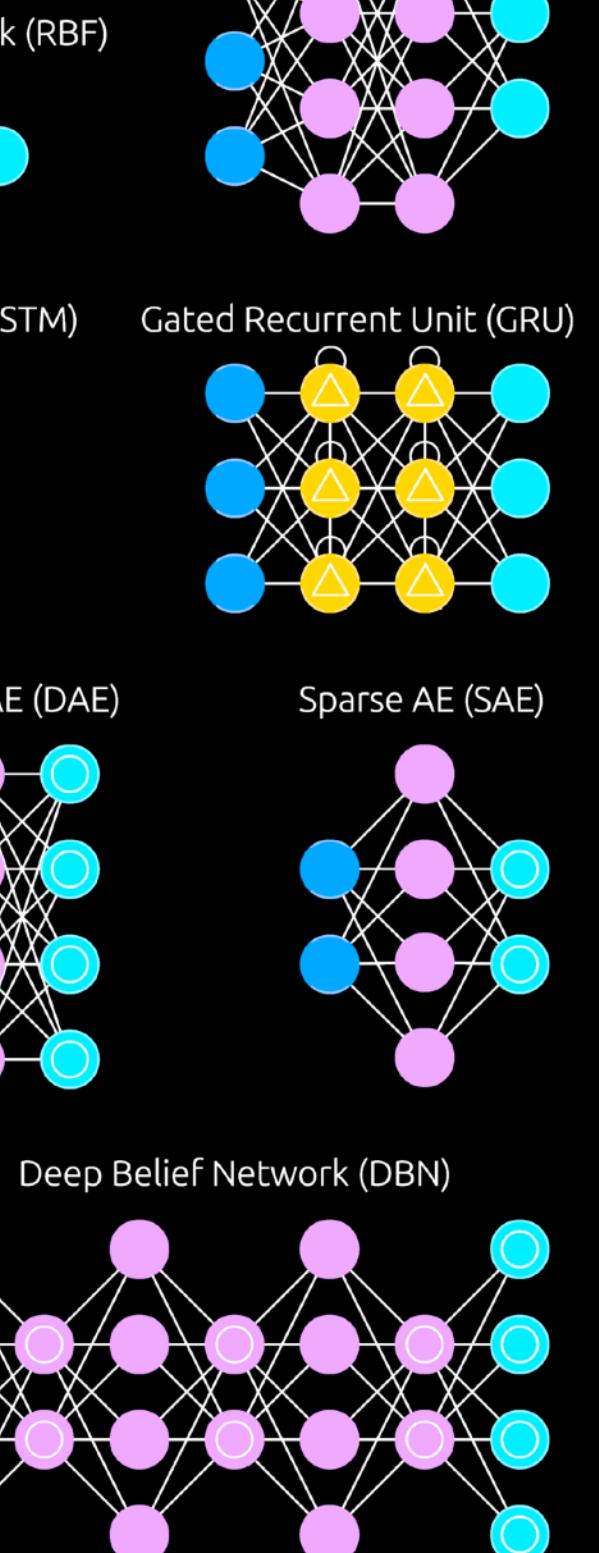
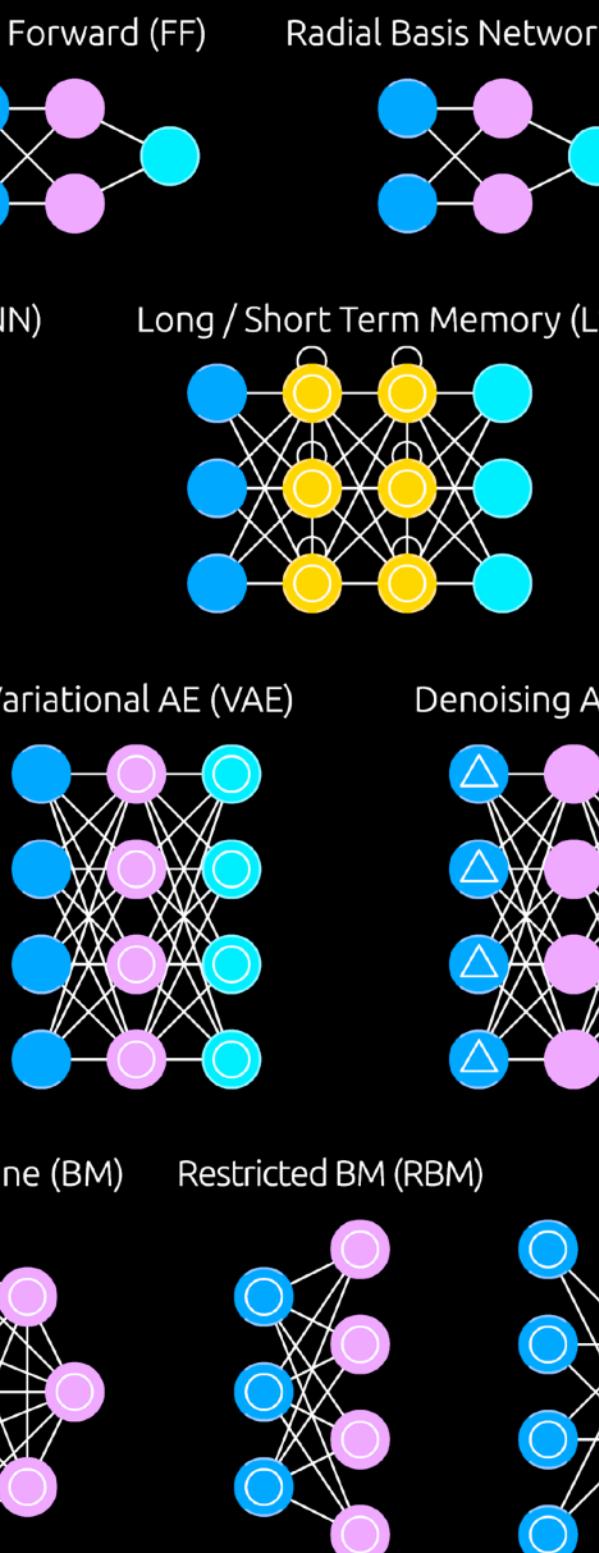
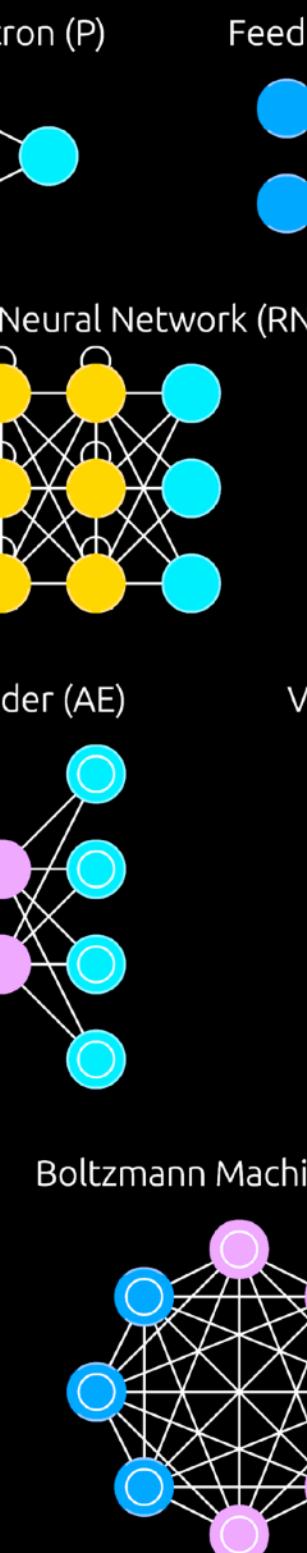
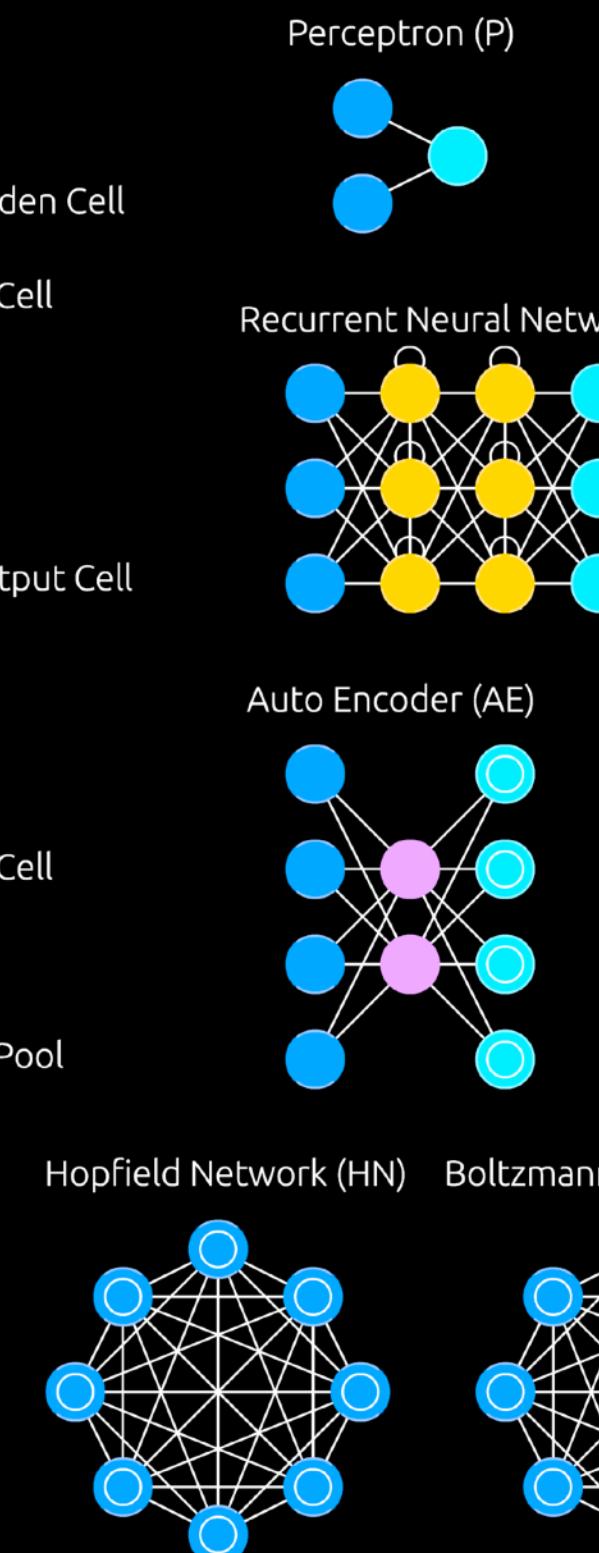
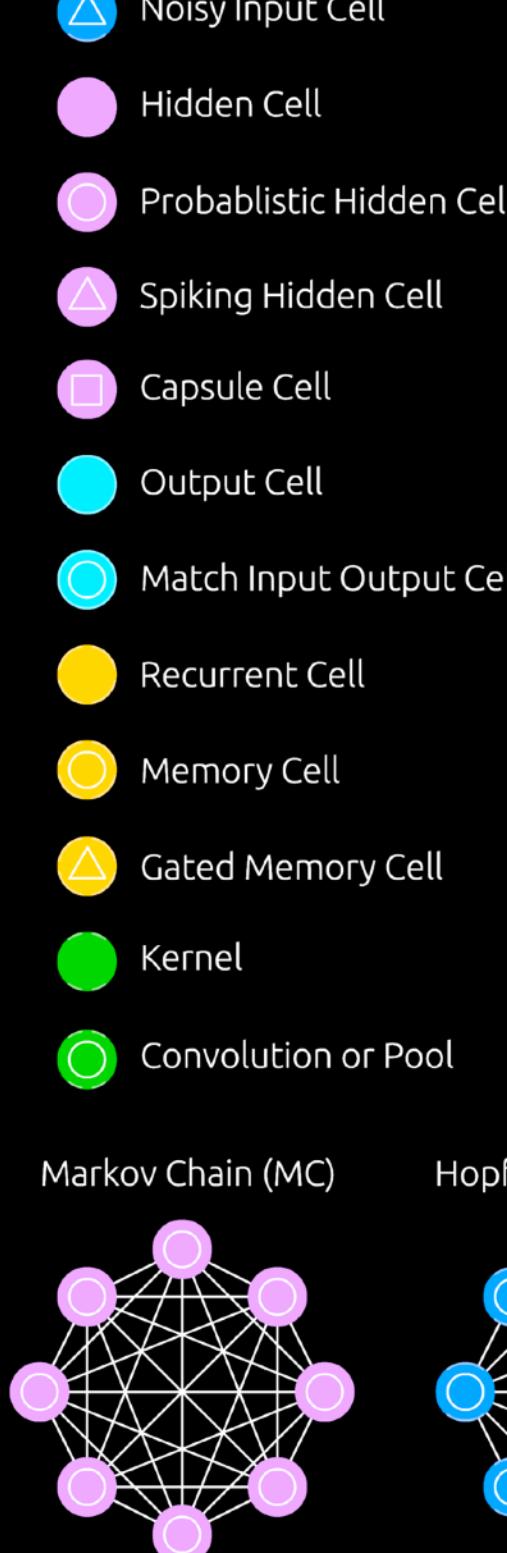


Neural network zoo

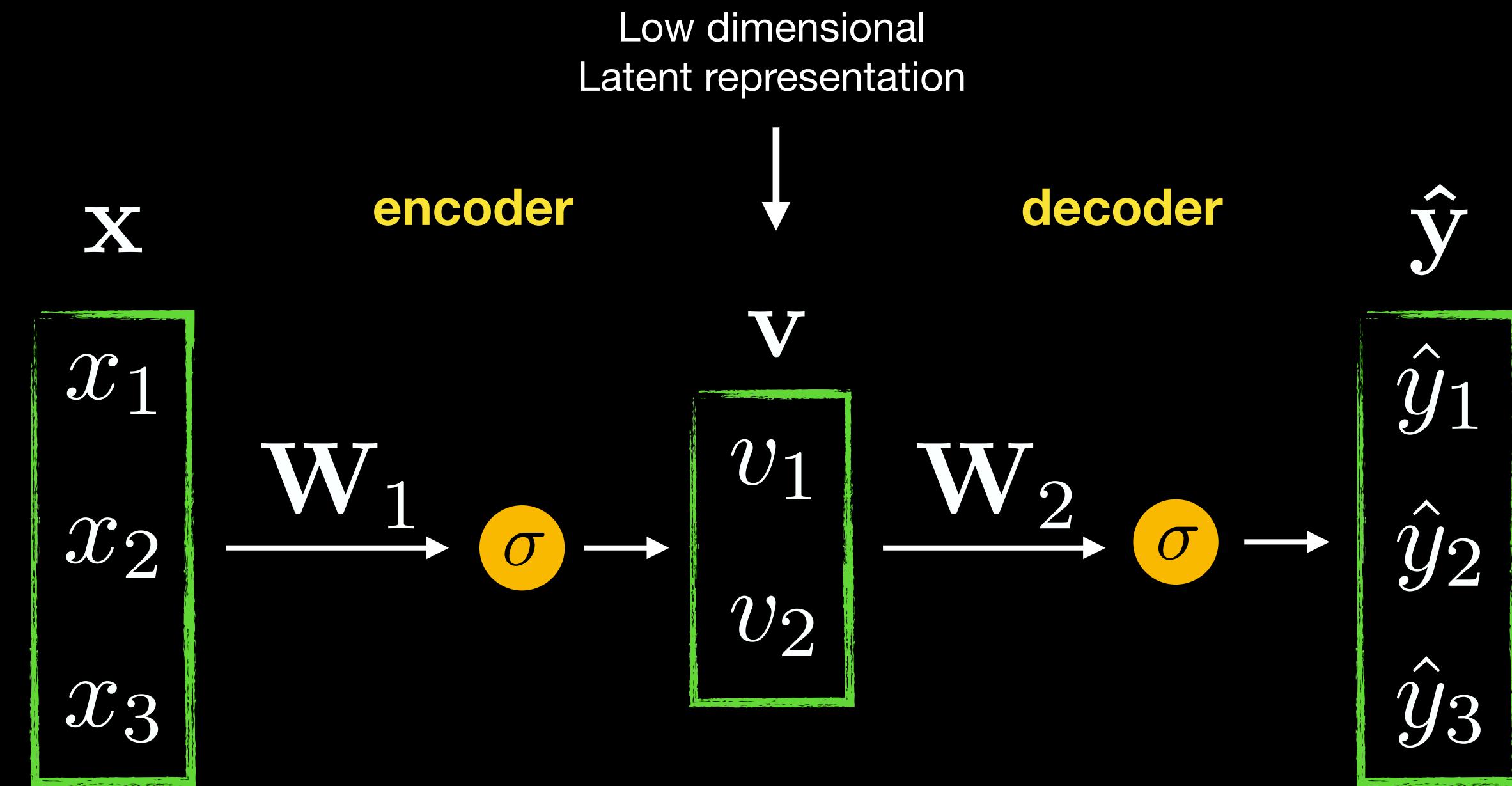
A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

- Input Cell
- Backfed Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Gated Memory Cell
- Kernel
- Convolution or Pool



Auto-encoders

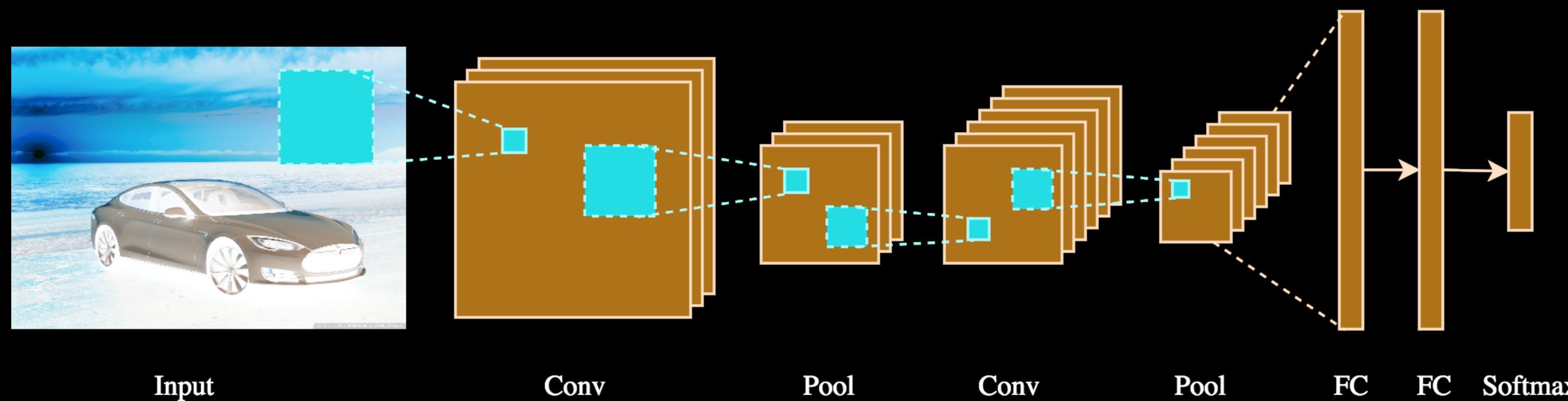
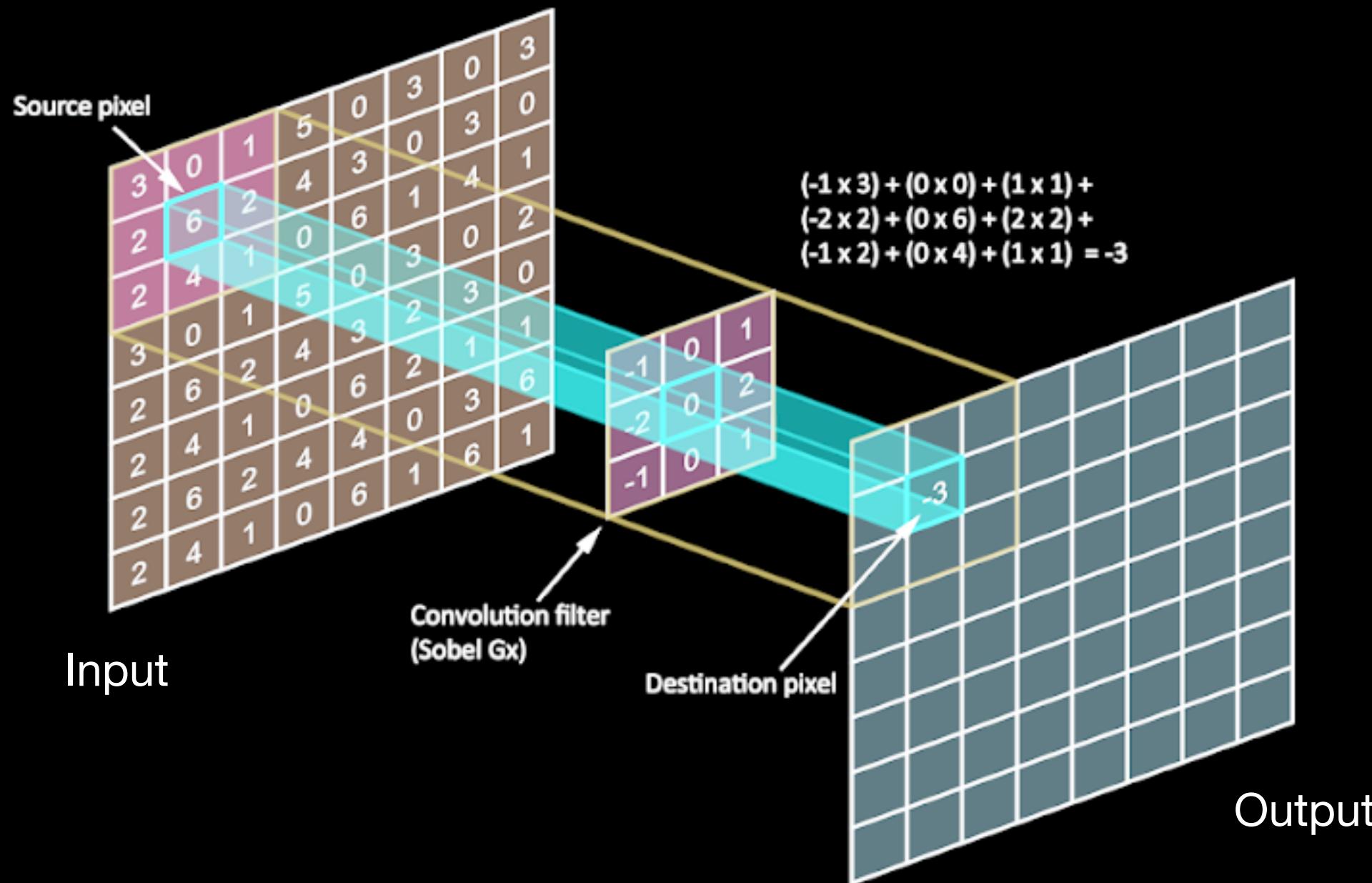


$$\mathcal{L} = \|f_{\mathbf{W}_1 \mathbf{W}_2}(\mathbf{x}) - \mathbf{x}\|^2$$

if $\sigma = I$, network performs SVD decomposition

$$\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^*$$

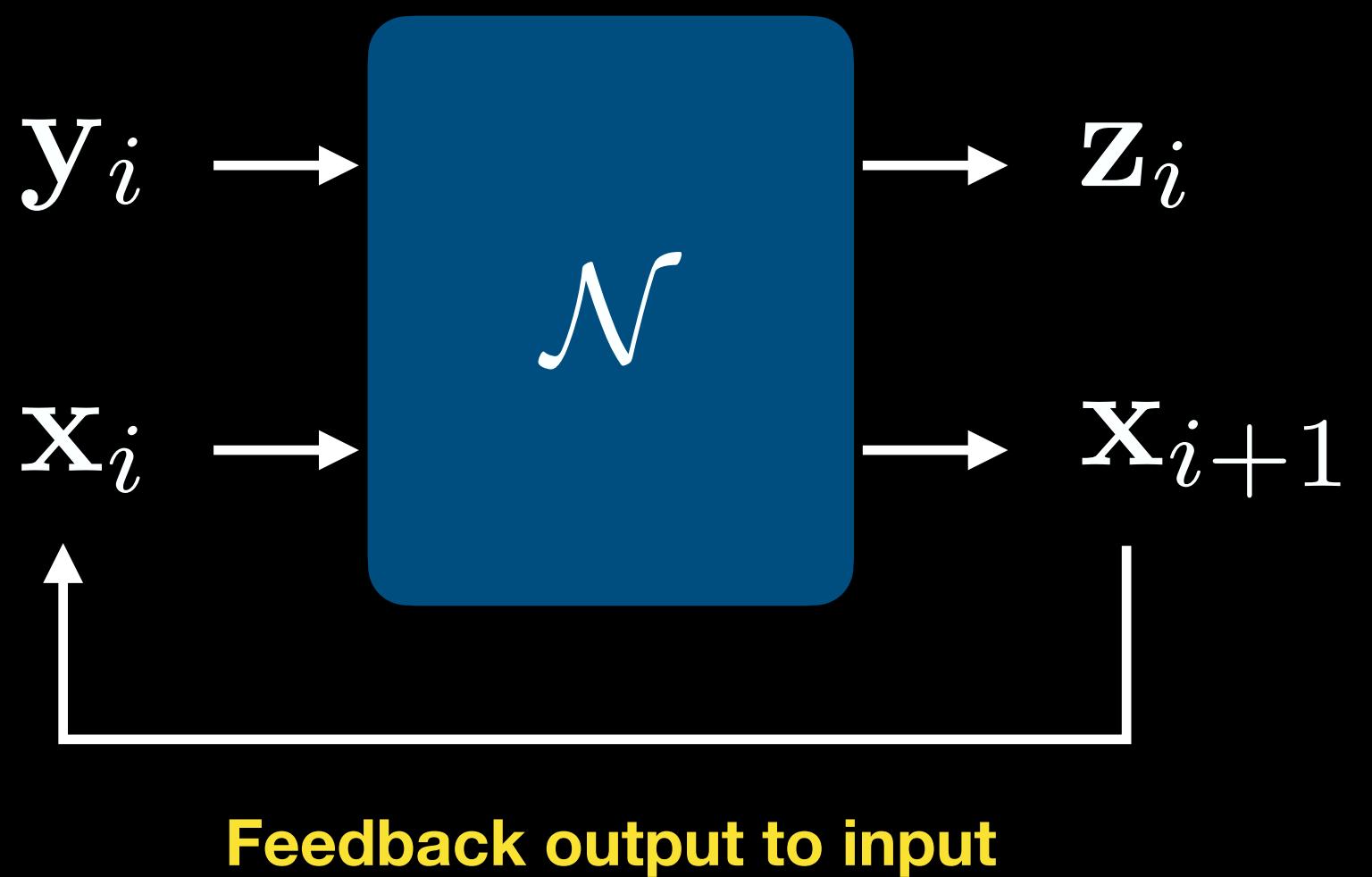
Convolutional neural networks



Recurrent neural networks

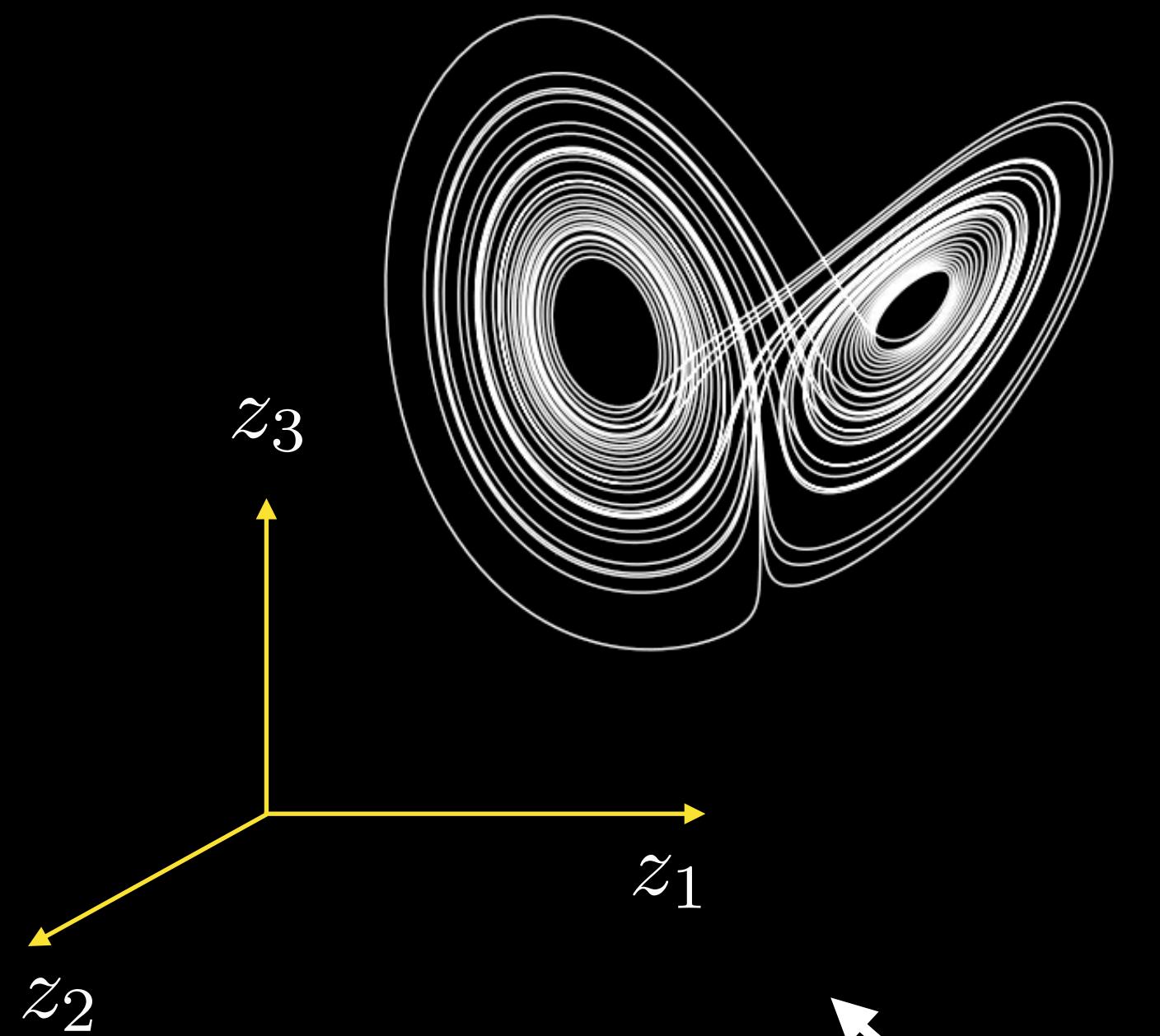
$$\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n$$

$$[\mathbf{x}_{i+1}, \mathbf{z}_i] = \mathcal{N}(\mathbf{x}_i, \mathbf{y}_i)$$

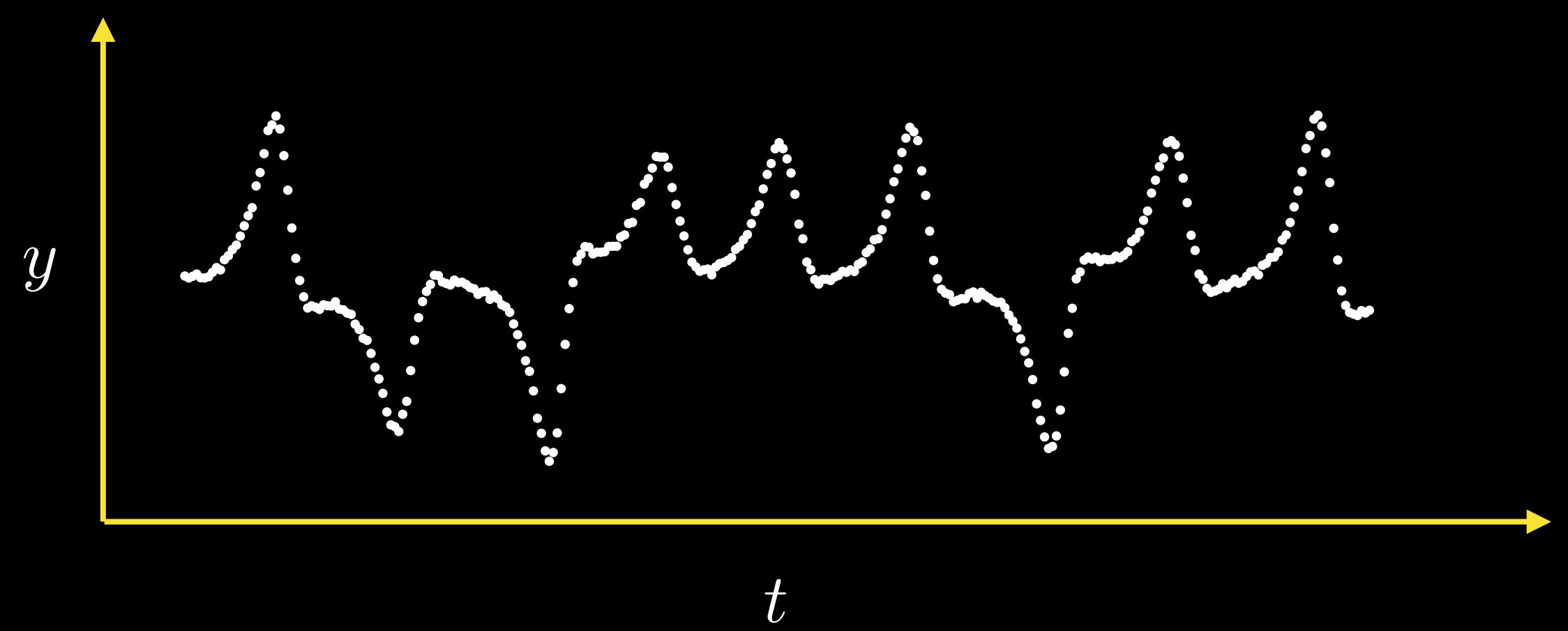


Numerical solvers are recurrence relations!

High Dimensional System



Partial Measurements



Measure



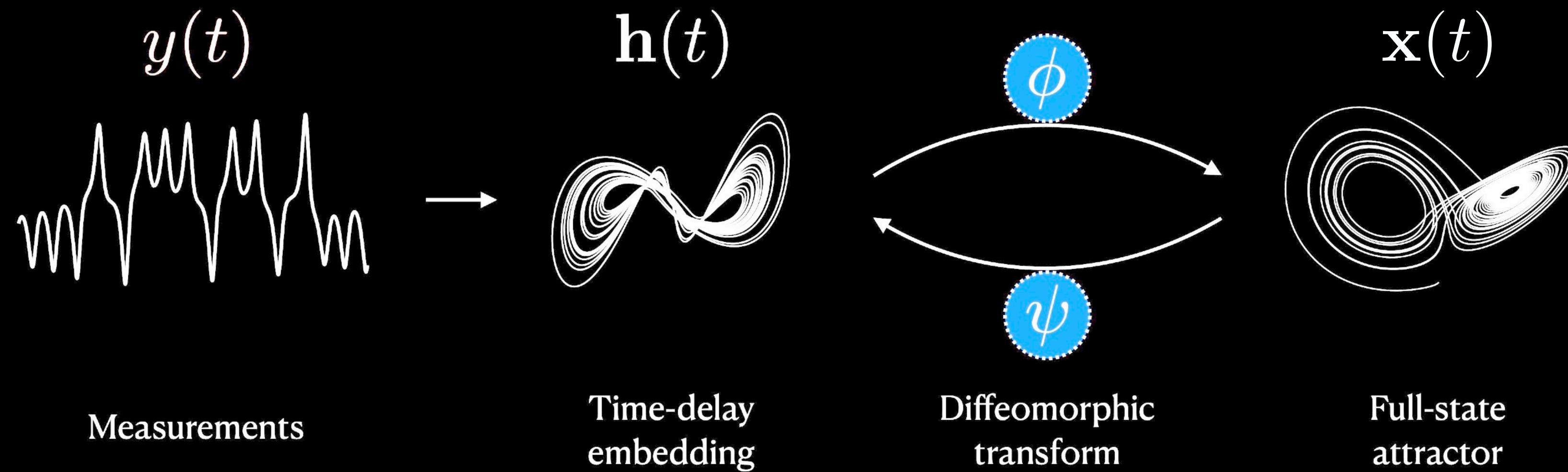
Reconstruct

?

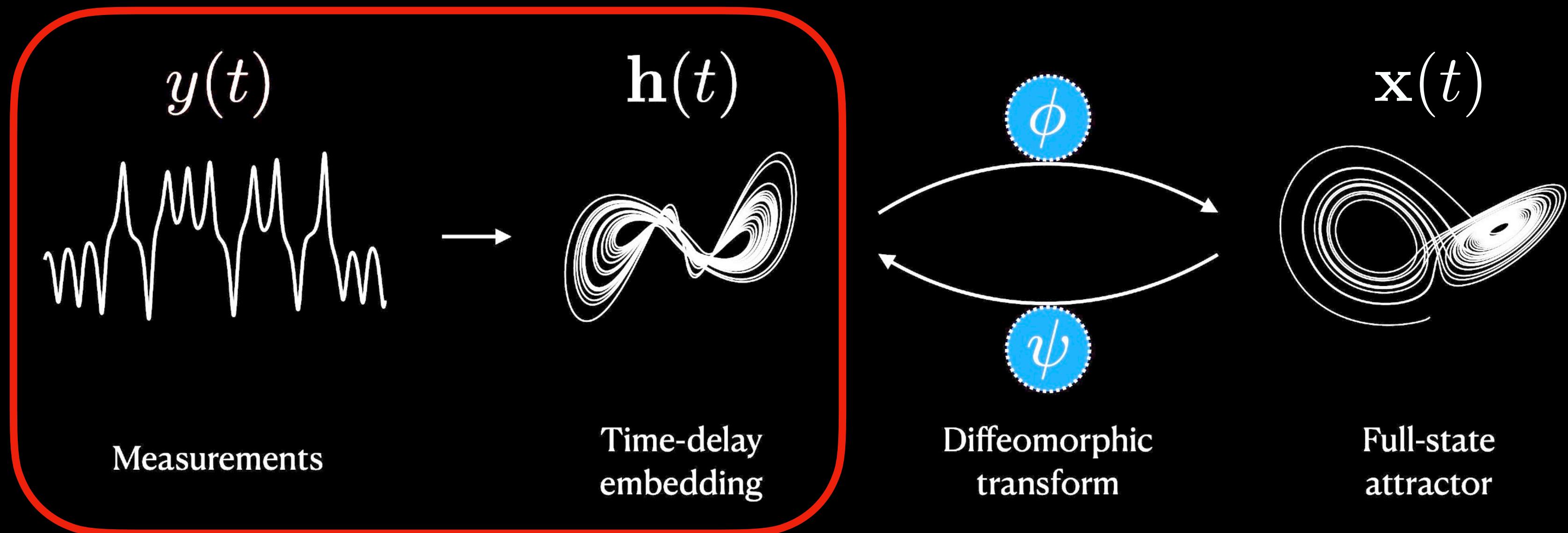


Takens' Theorem

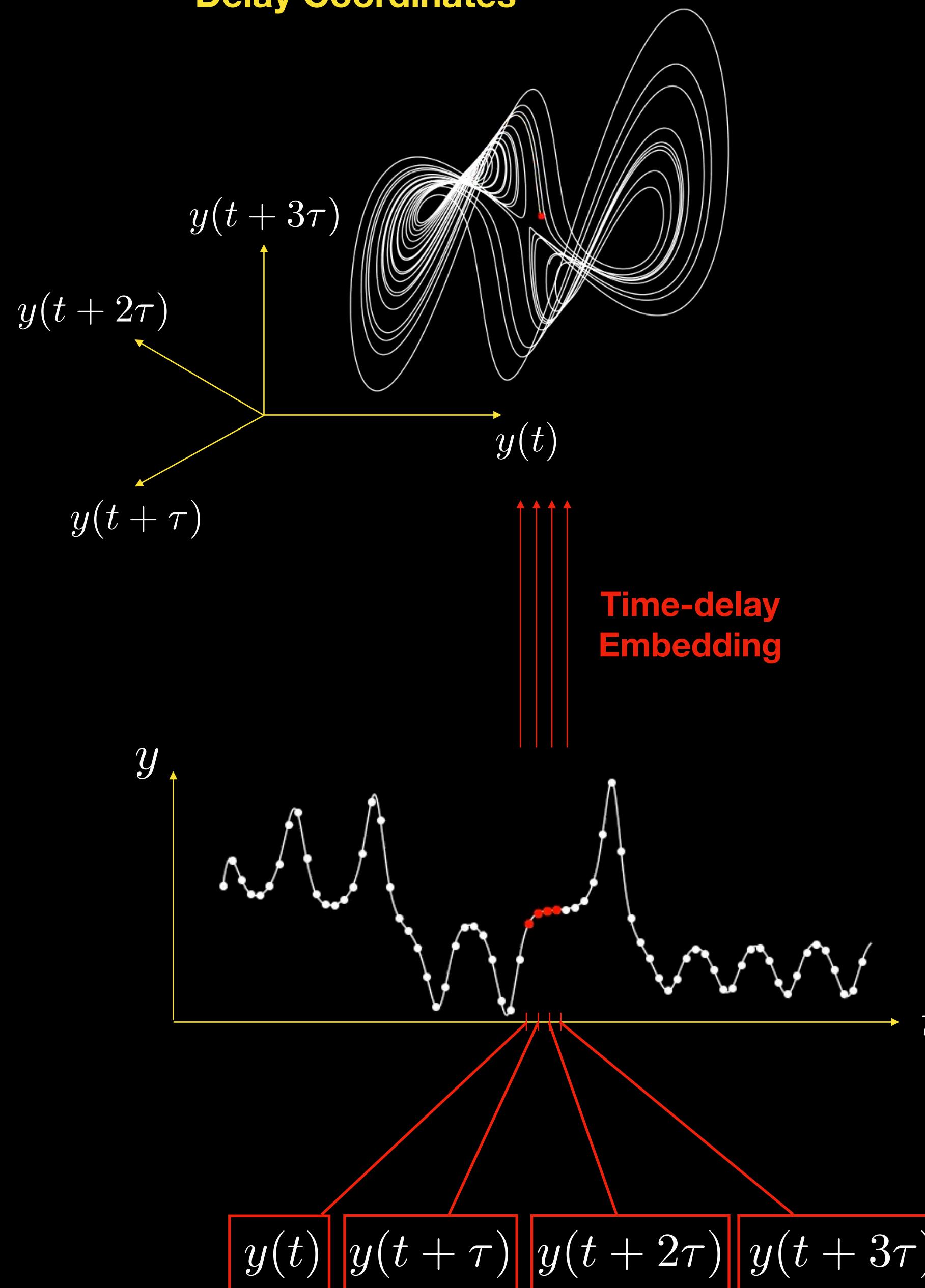
Detecting strange attractors in turbulence (1981)



Time-delay embedding



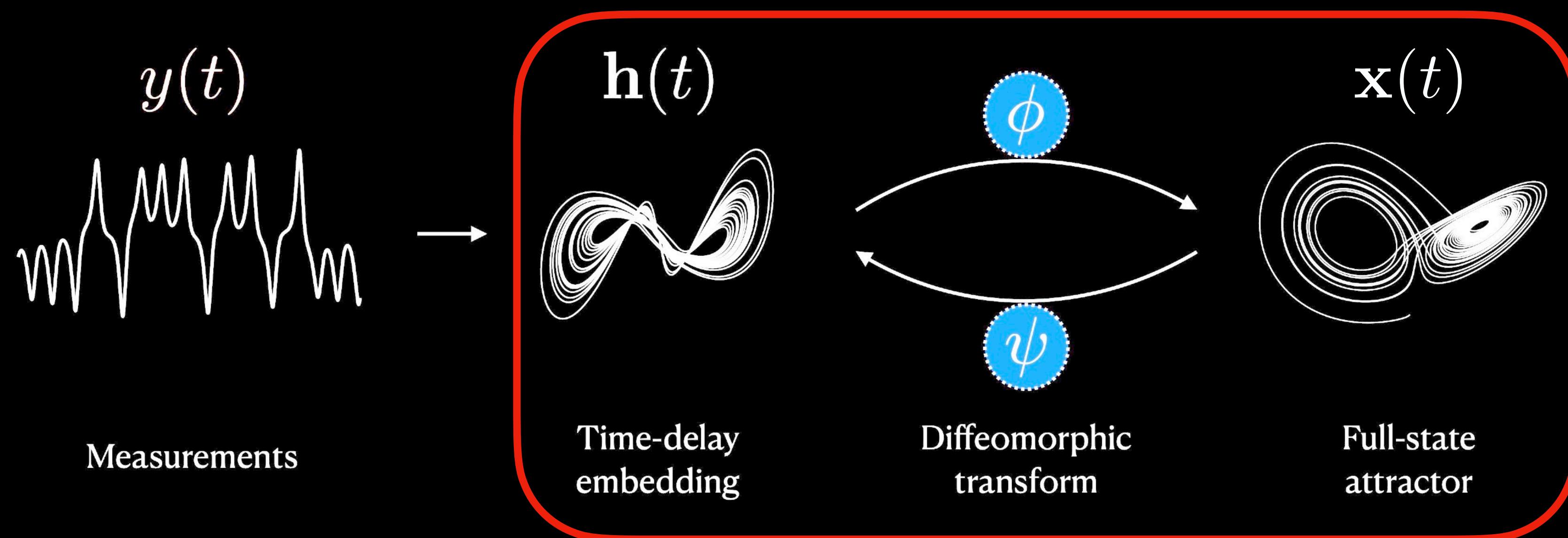
Delay Coordinates



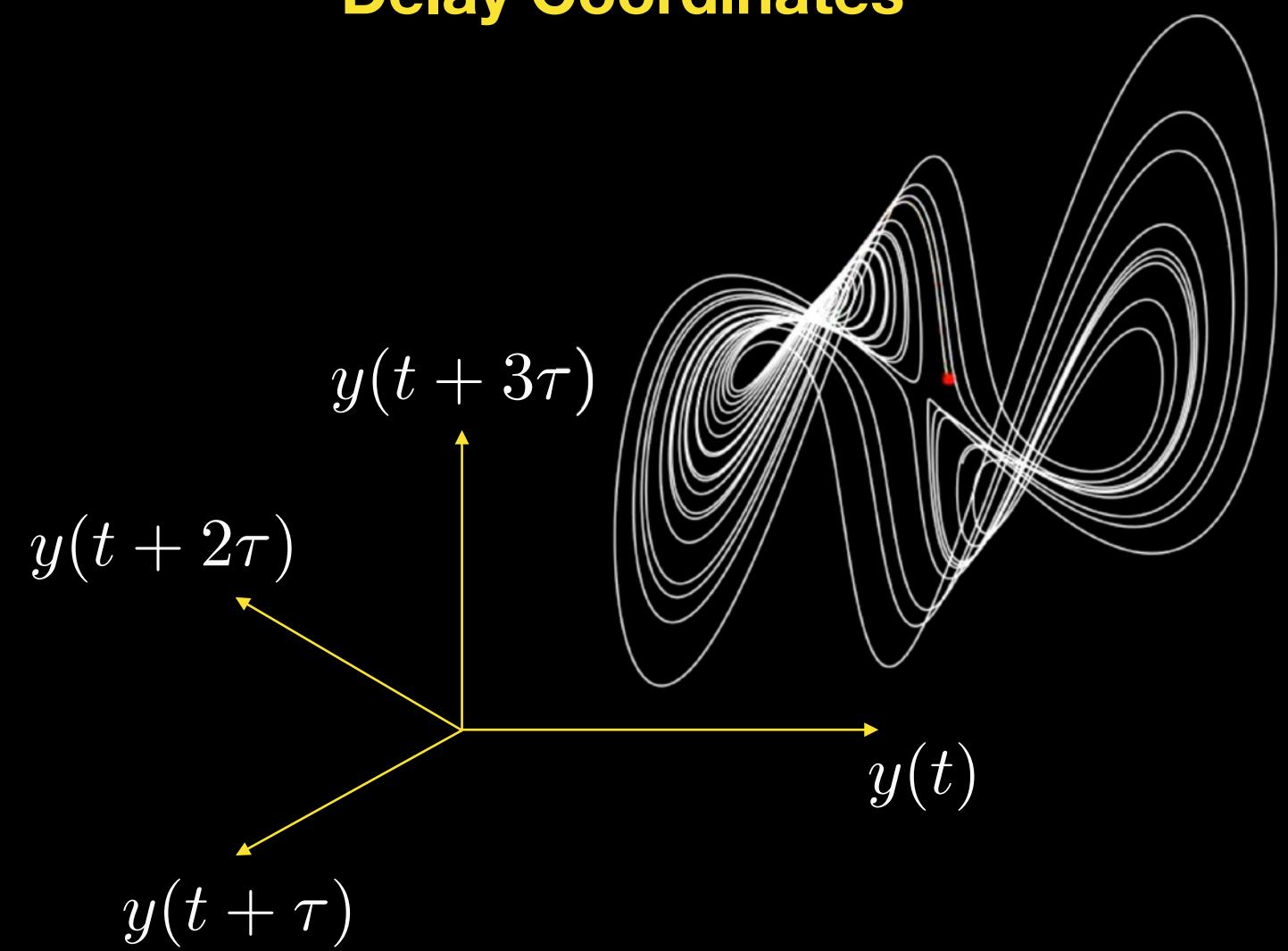
Hankel Matrix

$$\mathbf{H} = \begin{bmatrix} y(t_1) & y(t_2) & \dots & y(t_q) \\ y(t_2) & y(t_3) & \dots & y(t_{q+1}) \\ \vdots & \vdots & \ddots & \vdots \\ y(t_n) & y(t_{n+1}) & \dots & y(t_{n+q+1}) \end{bmatrix} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_q].$$

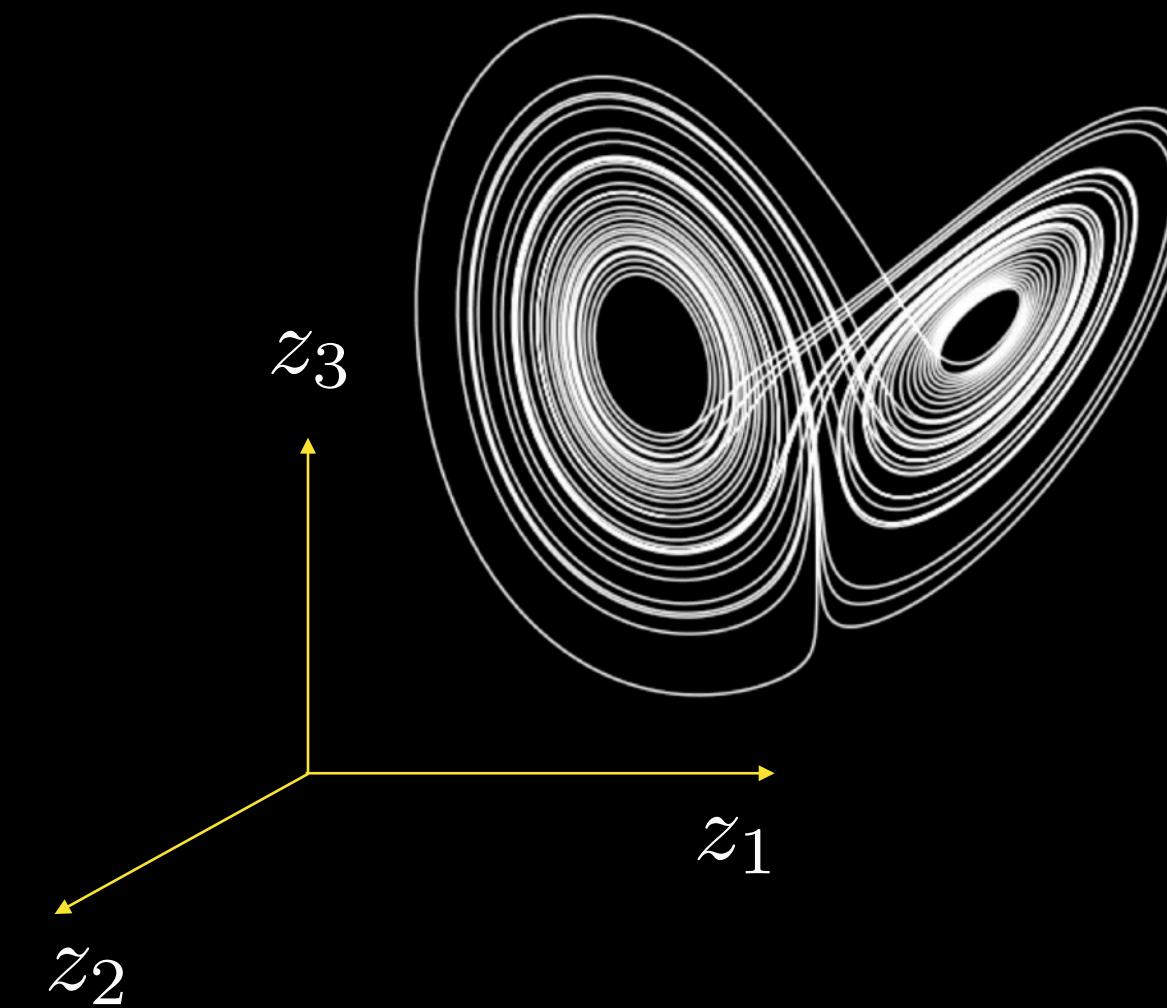
Takens' Theorem



Delay Coordinates



Full State Coordinates



Diffeomorphic



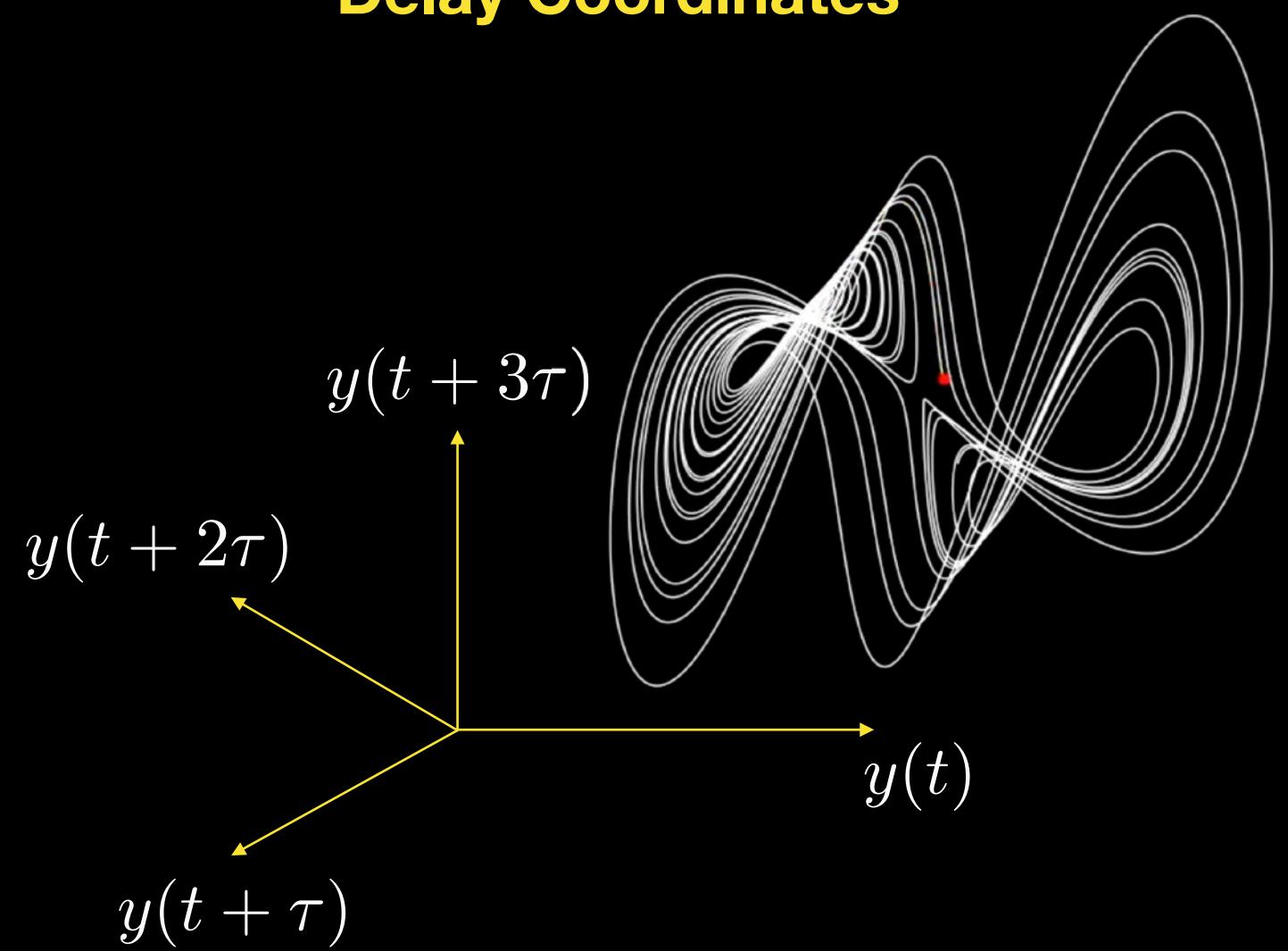
Takens' Theorem

If the dynamical system and the observed quantity are generic,
then the **delay-coordinate map** from a d -dimensional smooth
compact manifold \mathcal{M} to \mathbb{R}^{2d+1} is a diffeomorphism on \mathcal{M} .

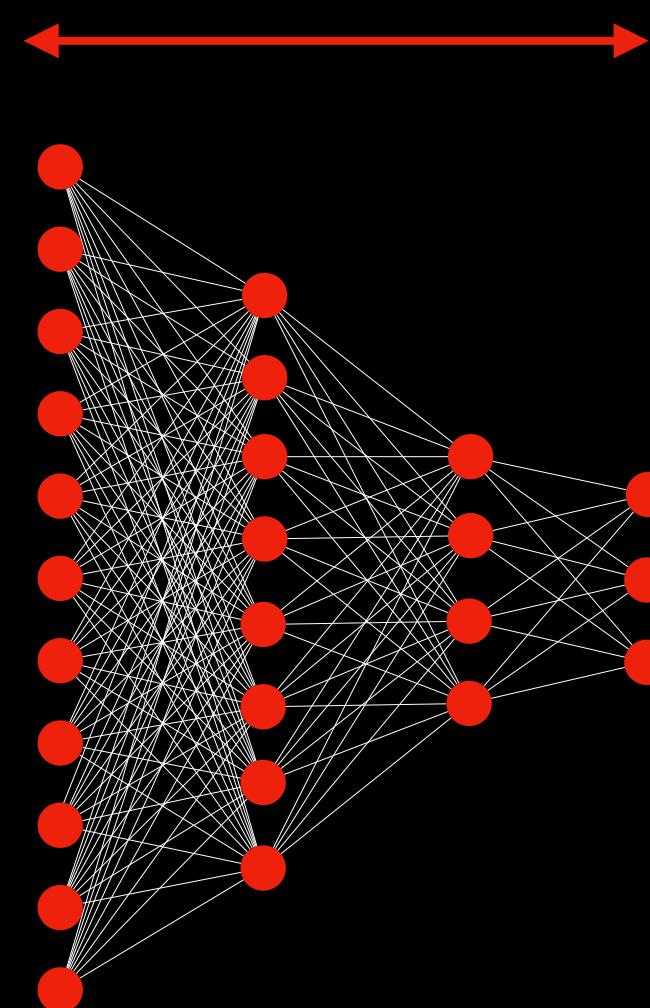
F. Takens

Detecting strange attractors in turbulence (1981)

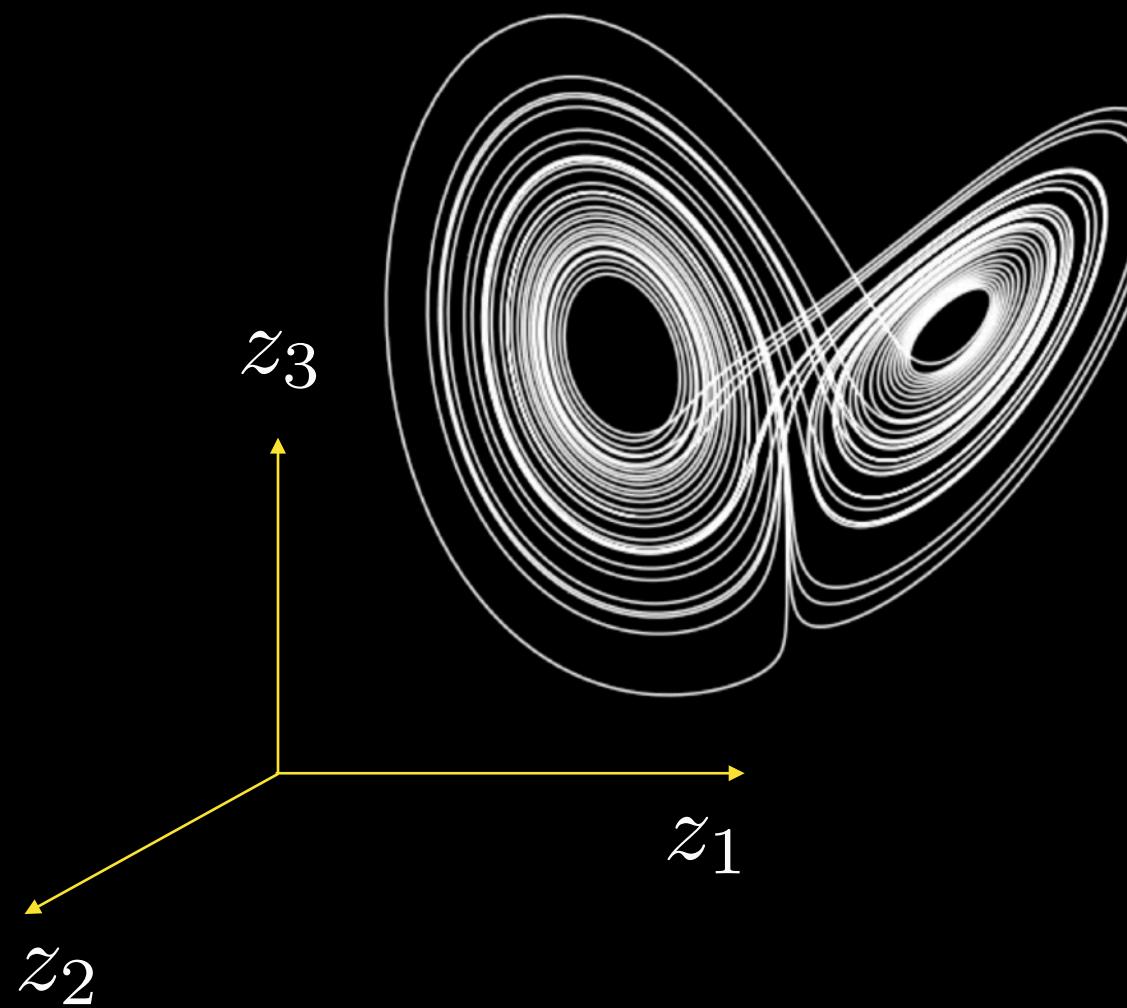
Delay Coordinates



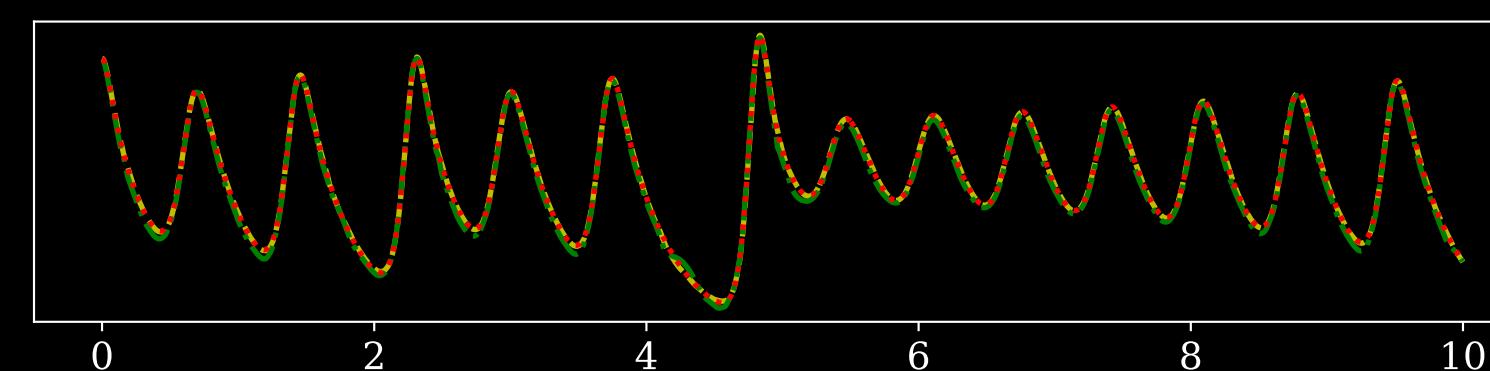
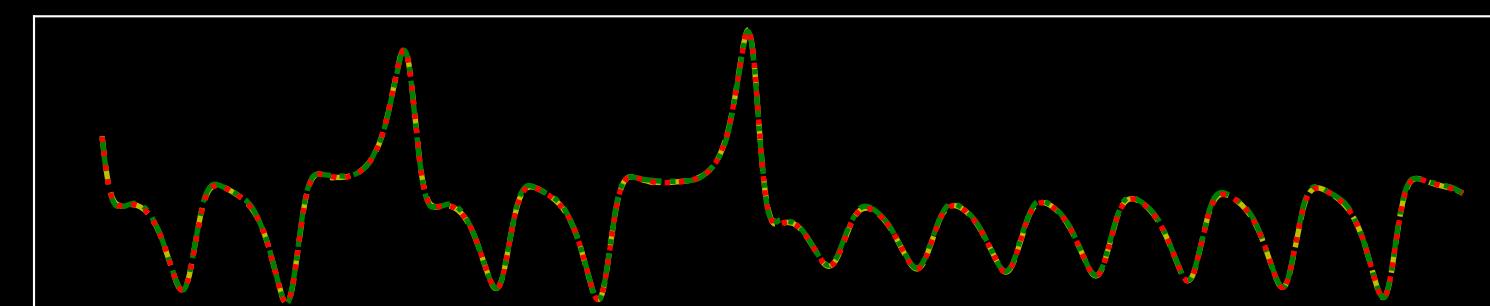
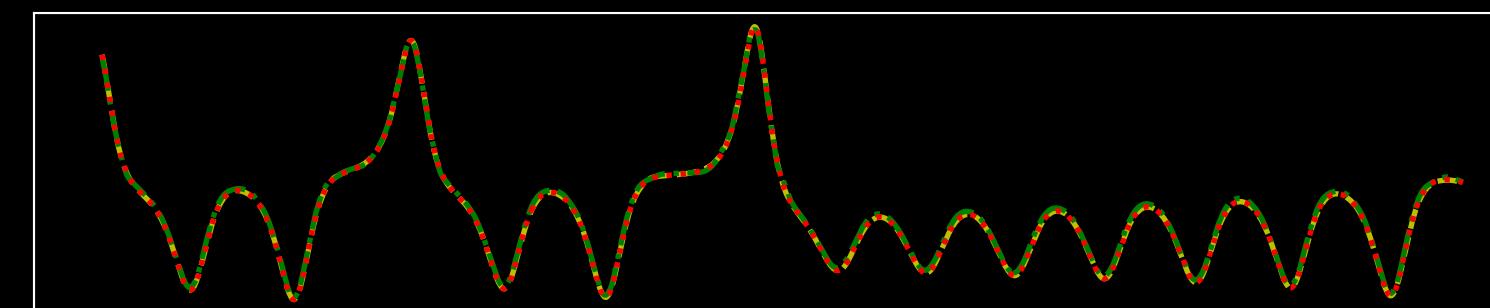
Diffeomorphic



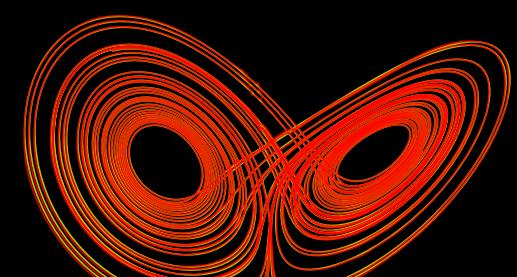
Full State Coordinates



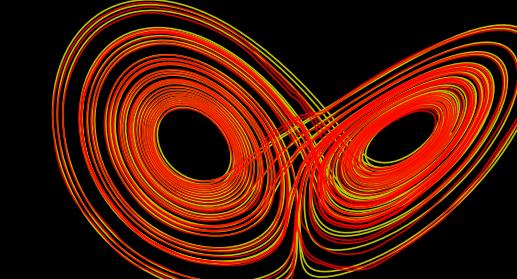
— Supervised Network
— Supervised SVD + Network
... True



— True
— Prediction



Network



SVD + Network