

A middle-aged man with glasses and a goatee, wearing a light-colored suit jacket, a white shirt, and a dark tie with a circular pattern. He is speaking and looking slightly to the left of the camera. The background is solid black.

[www.grand-illusions.com](http://www.grand-illusions.com)



# Deep Learning Neural Net Basics

Computer Vision

James Hays

Many slides by Marc'Aurelio Ranzato

# Outline

- Neural Networks
- *Convolutional* Neural Networks
- Variants
  - Detection
  - Segmentation
  - Siamese Networks
- Visualization of Deep Networks

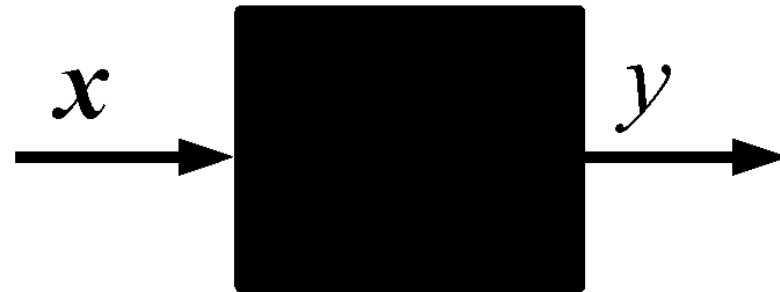
# Supervised Learning

$\{(\mathbf{x}^i, y^i), i=1 \dots P\}$  training dataset

$\mathbf{x}^i$  i-th input training example

$y^i$  i-th target label

$P$  number of training examples



Goal: predict the target label of unseen inputs.

# Supervised Learning: Examples

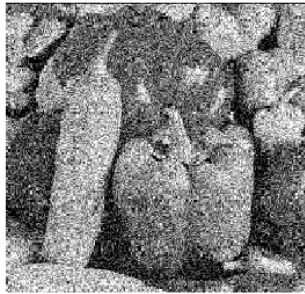
## Classification



“dog”

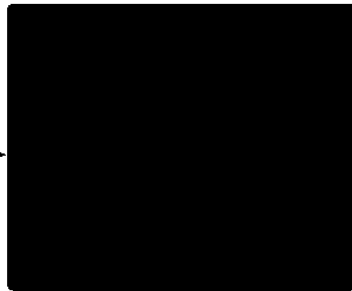
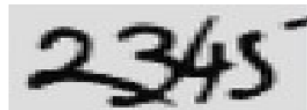
*classification*

## Denoising



*regression*

## OCR

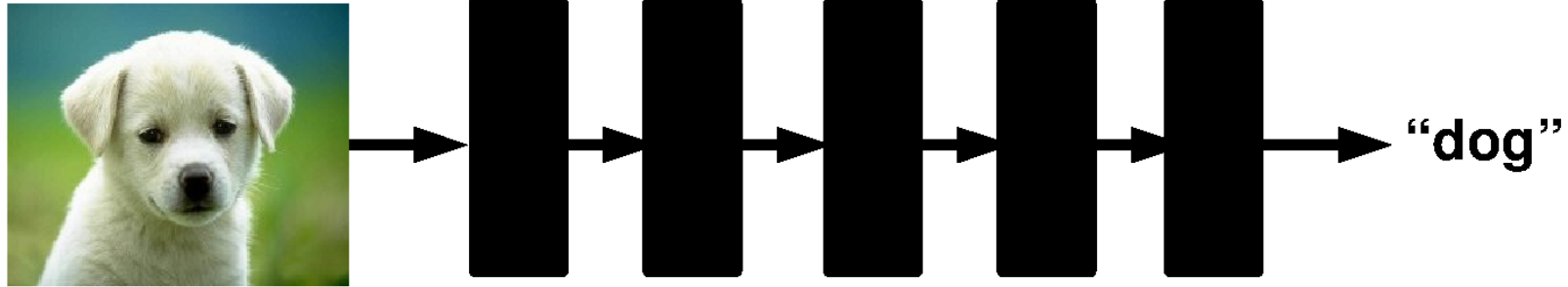


“2 3 4 5”

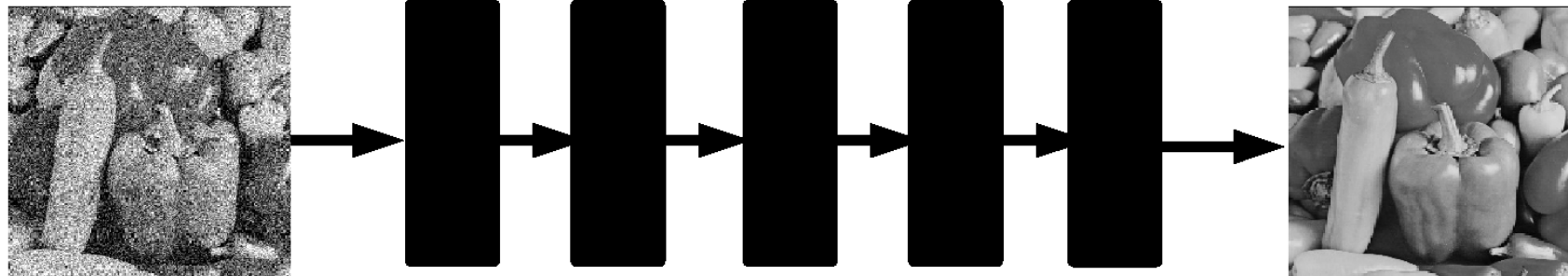
*structured prediction*

# Supervised Deep Learning

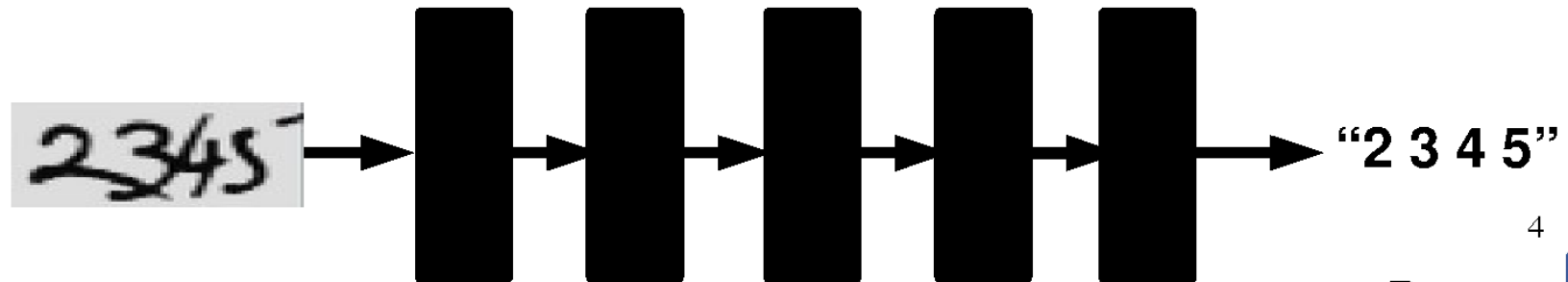
## Classification



## Denoising



## OCR



# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips



# Neural Networks

Assumptions (for the next few slides):

- The input image is vectorized (disregard the spatial layout of pixels)
- The target label is discrete (classification)

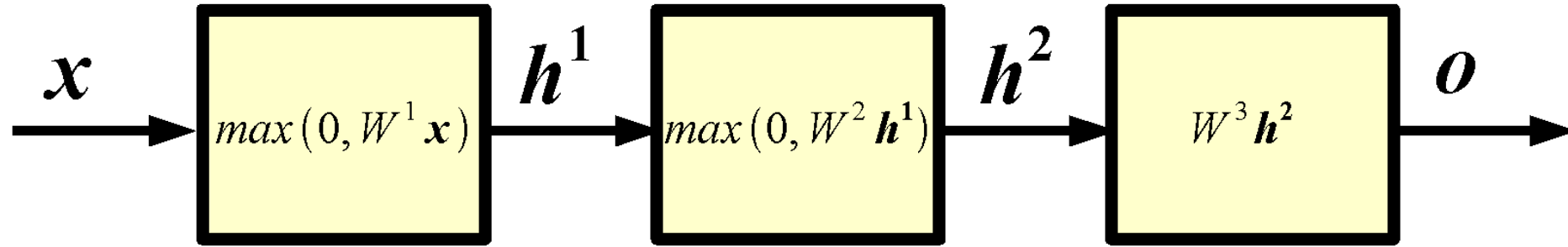
**Question:** what class of functions shall we consider to map the input into the output?

**Answer:** composition of simpler functions.

**Follow-up questions:** Why not a linear combination? What are the “simpler” functions? What is the interpretation?

**Answer:** later...

# Neural Networks: example



$x$  input

$h^1$  1-st layer hidden units

$h^2$  2-nd layer hidden units

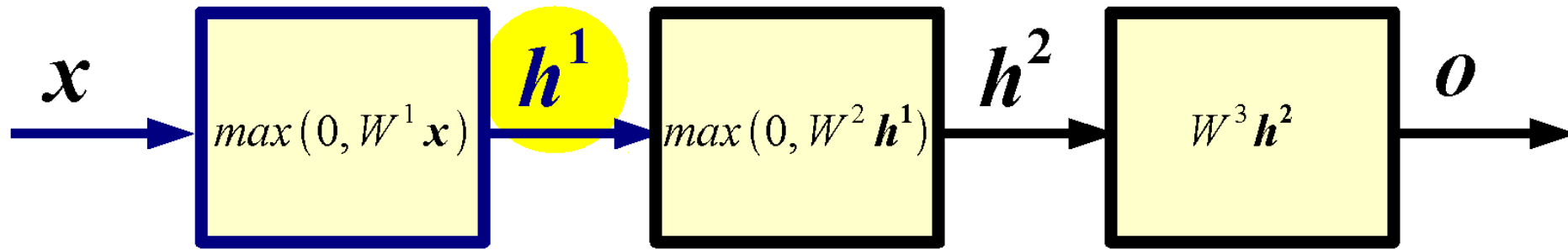
$o$  output

Example of a 2 hidden layer neural network (or 4 layer network, counting also input and output).

# Forward Propagation

**Def.:** Forward propagation is the process of computing the output of the network given its input.

# Forward Propagation



$$x \in R^D \quad W^1 \in R^{N_1 \times D} \quad b^1 \in R^{N_1} \quad h^1 \in R^{N_1}$$

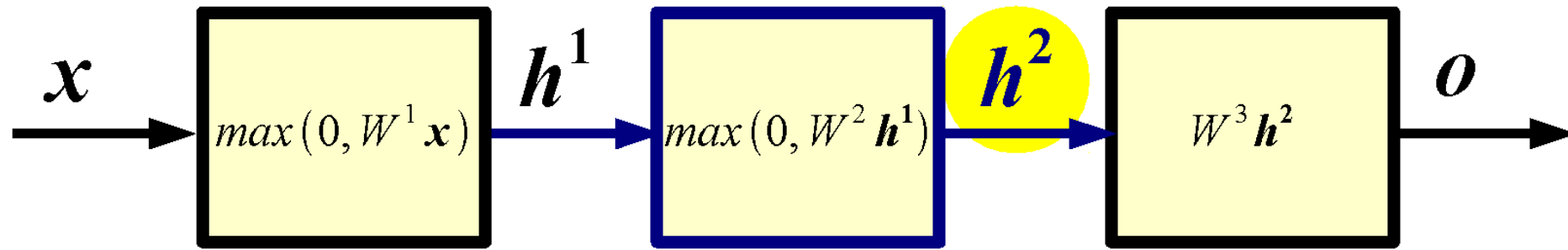
$$h^1 = \max(0, W^1 x + b^1)$$

$W^1$  1-st layer weight matrix or weights

$b^1$  1-st layer biases

The non-linearity  $u = \max(0, v)$  is called **ReLU** in the DL literature. Each output hidden unit takes as input all the units at the previous layer: each such layer is called “**fully connected**”.

# Forward Propagation



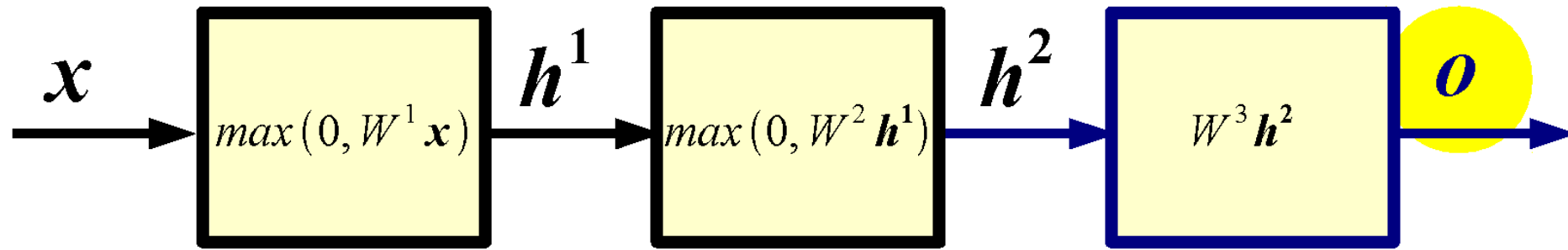
$$h^1 \in R^{N_1} \quad W^2 \in R^{N_2 \times N_1} \quad b^2 \in R^{N_2} \quad h^2 \in R^{N_2}$$

$$h^2 = \max(0, W^2 h^1 + b^2)$$

$W^2$  2-nd layer weight matrix or weights

$b^2$  2-nd layer biases

# Forward Propagation



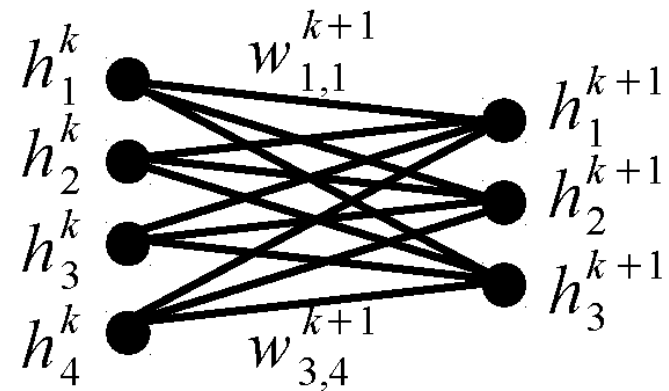
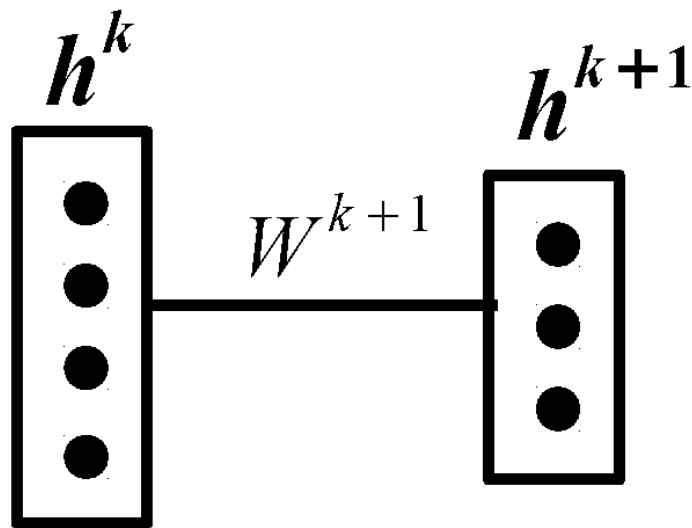
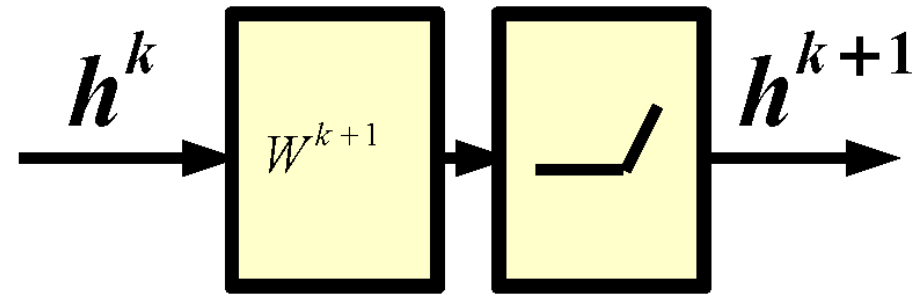
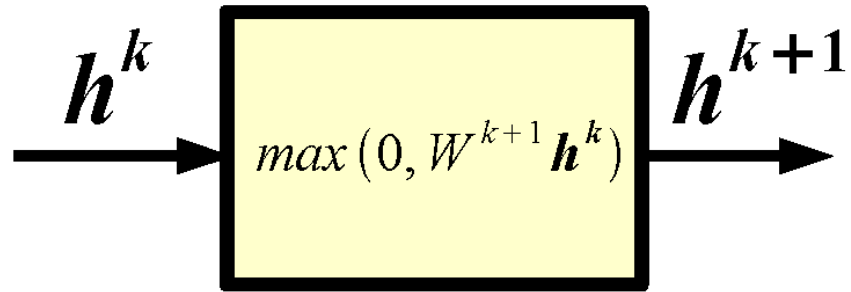
$$h^2 \in R^{N_2} \quad W^3 \in R^{N_3 \times N_2} \quad b^3 \in R^{N_3} \quad o \in R^{N_3}$$

$$o = \max(0, W^3 h^2 + b^3)$$

$W^3$  3-rd layer weight matrix or weights

$b^3$  3-rd layer biases

# Alternative Graphical Representation

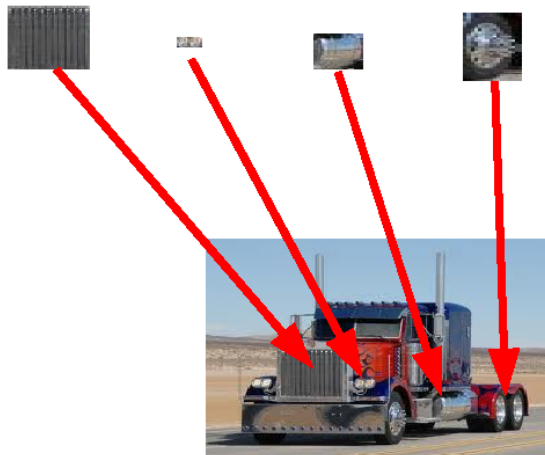


# Interpretation

**Question:** Why do we need many layers?

**Answer:** When input has hierarchical structure, the use of a hierarchical architecture is potentially more efficient because intermediate computations can be re-used. DL architectures are efficient also because they use **distributed representations** which are shared across classes.

[0 0 **1** 0 0 0 0 **1** 0 0 **1** **1** 0 0 **1** 0 ... ] truck feature



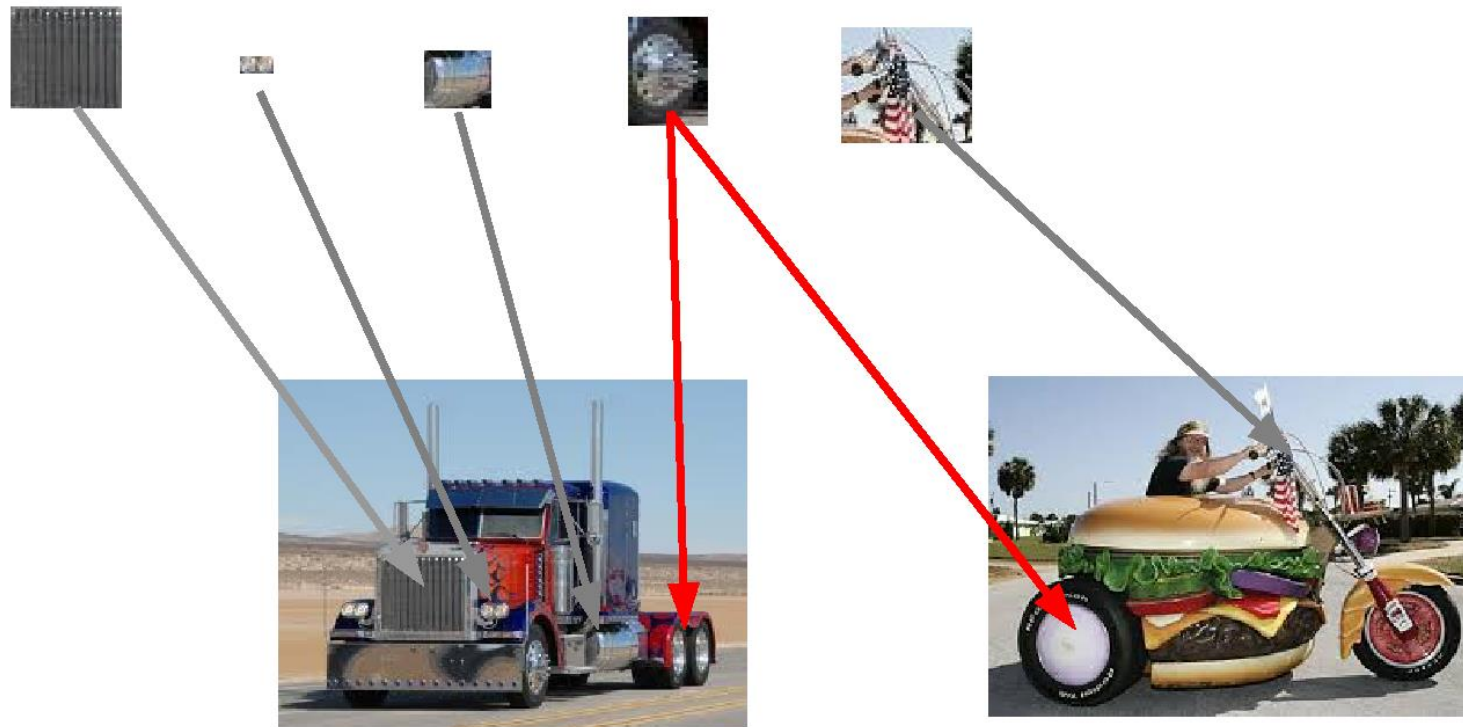
Exponentially more efficient than a 1-of-N representation (a la k-means)



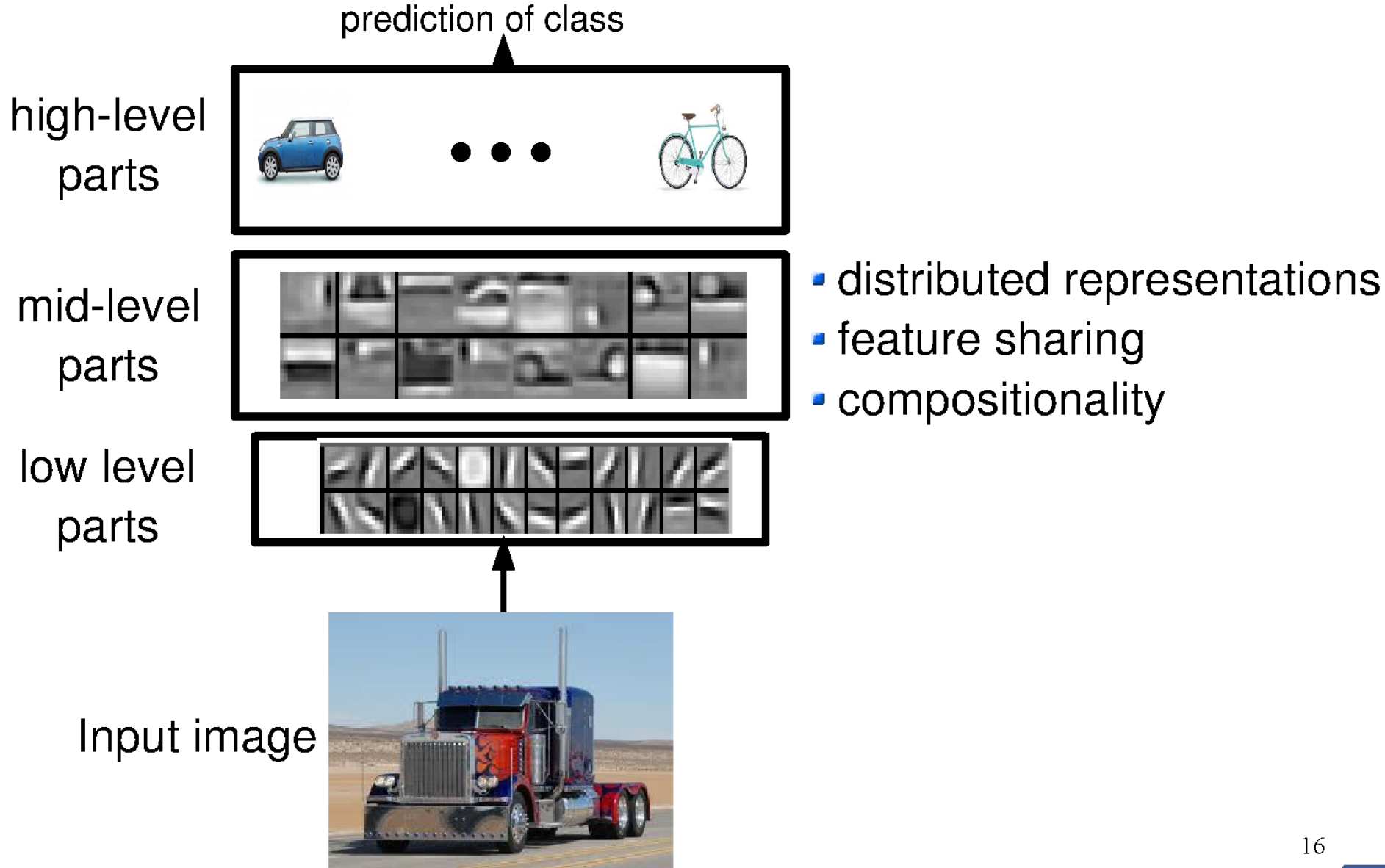
# Interpretation

[1 1 0 0 0 1 0 **1** 0 0 0 0 1 1 0 1... ] motorbike

[0 0 1 0 0 0 0 **1** 0 0 1 1 0 0 1 0 ... ] truck



# Interpretation



# Interpretation

**Question:** What does a hidden unit do?

**Answer:** It can be thought of as a classifier or feature detector.

**Question:** How many layers? How many hidden units?

**Answer:** Cross-validation or hyper-parameter search methods are the answer. In general, the wider and the deeper the network the more complicated the mapping.

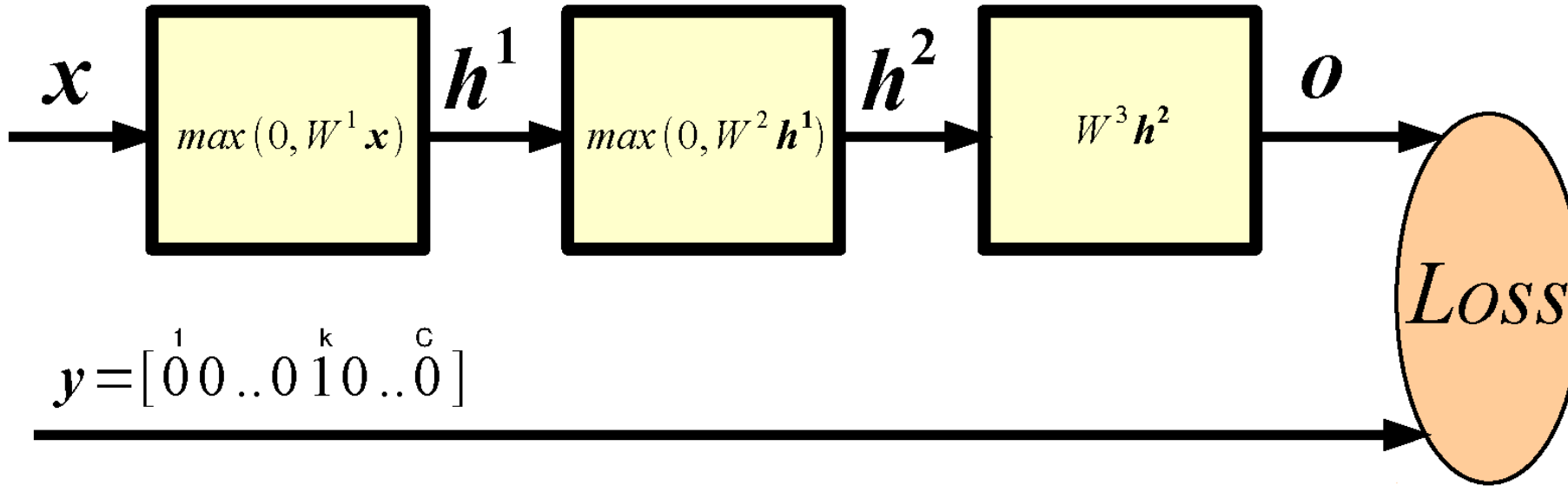
**Question:** How do I set the weight matrices?

**Answer:** Weight matrices and biases are learned.

First, we need to define a measure of quality of the current mapping.

Then, we need to define a procedure to adjust the parameters.

# How Good is a Network?



Probability of class  $k$  given input (softmax):

$$p(c_k = 1 | \mathbf{x}) = \frac{e^{o_k}}{\sum_{j=1}^C e^{o_j}}$$

(Per-sample) **Loss**; e.g., negative log-likelihood (good for classification of small number of classes):

$$L(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = - \sum_j y_j \log p(c_j | \mathbf{x})$$

# Training

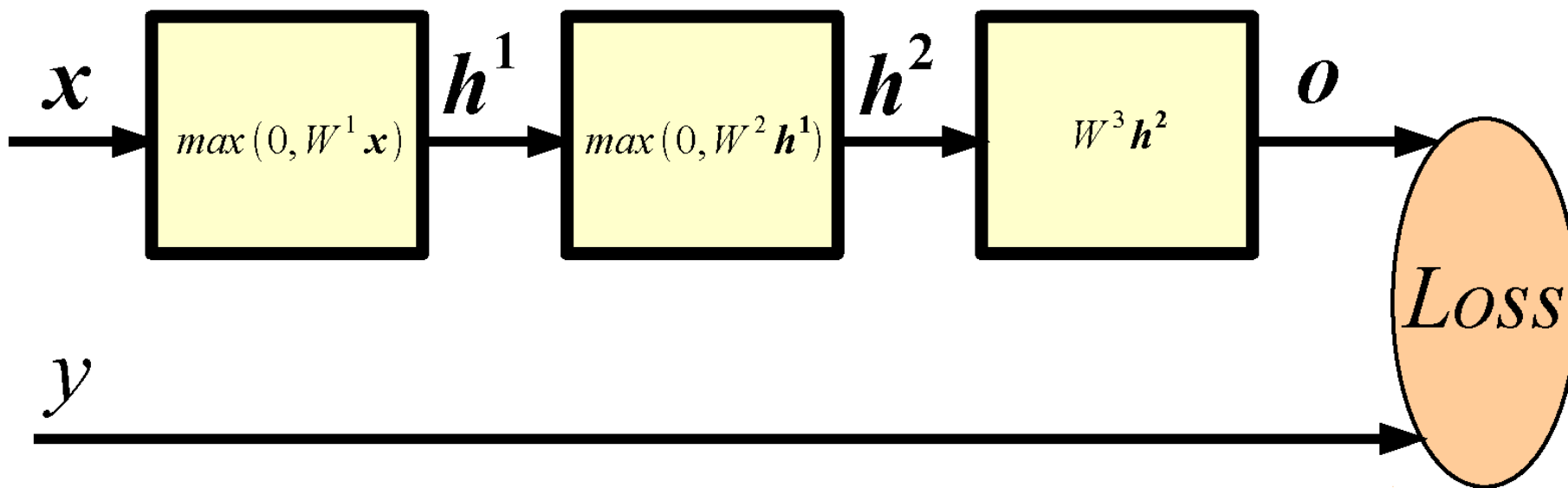
**Learning** consists of minimizing the loss (plus some regularization term) w.r.t. parameters over the whole training set.

$$\boldsymbol{\theta}^* = \mathit{arg\ min}_{\boldsymbol{\theta}} \sum_{n=1}^P L(\mathbf{x}^n, y^n; \boldsymbol{\theta})$$

**Question:** How to minimize a complicated function of the parameters?

**Answer:** Chain rule, a.k.a. **Backpropagation!** That is the procedure to compute gradients of the loss w.r.t. parameters in a multi-layer neural network.

# Key Idea: Wiggle To Decrease Loss



Let's say we want to decrease the loss by adjusting  $W_{i,j}^1$ .  
We could consider a very small  $\epsilon = 1e-6$  and compute:

$$L(\mathbf{x}, y; \boldsymbol{\theta})$$

$$L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon)$$

Then, update:

$$W_{i,j}^1 \leftarrow W_{i,j}^1 + \epsilon \operatorname{sgn}(L(\mathbf{x}, y; \boldsymbol{\theta}) - L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon))$$

# Derivative w.r.t. Input of Softmax

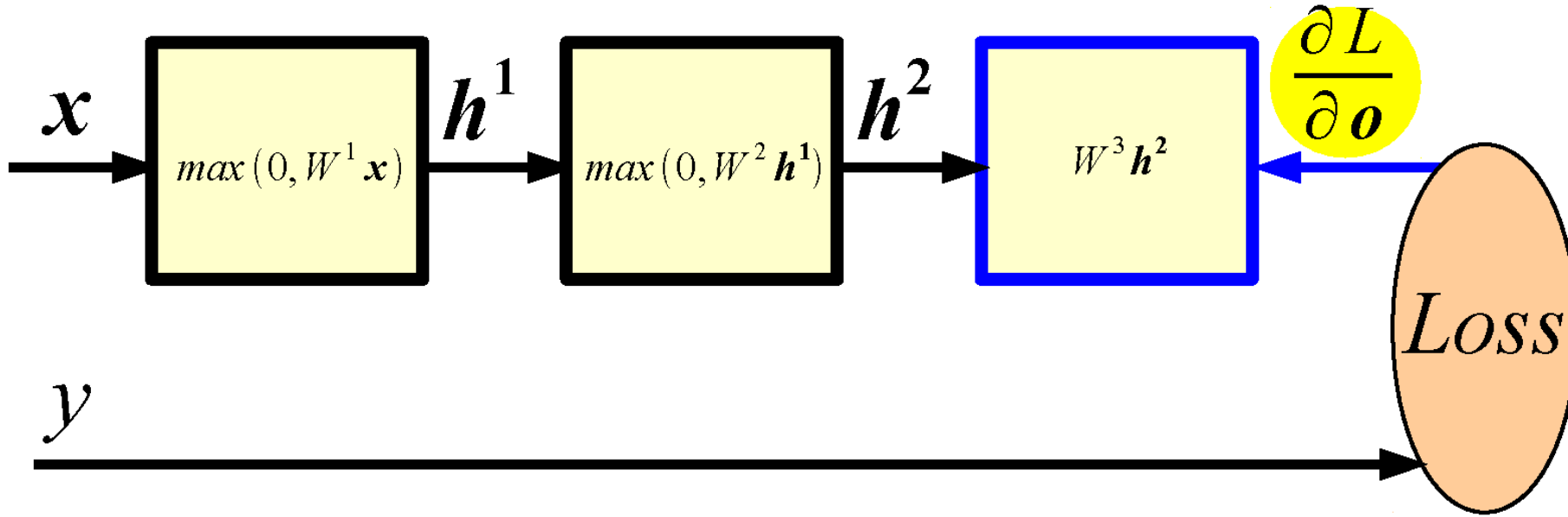
$$p(c_k = 1 | \mathbf{x}) = \frac{e^{o_k}}{\sum_j e^{o_j}}$$

$$L(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = - \sum_j y_j \log p(c_j | \mathbf{x}) \quad \mathbf{y} = [0 \overset{1}{0} \dots 0 \overset{k}{1} 0 \dots 0 \overset{c}{}]$$

By substituting the first formula in the second, and taking the derivative w.r.t.  $\boldsymbol{o}$  we get:

$$\frac{\partial L}{\partial \boldsymbol{o}} = p(c | \mathbf{x}) - \mathbf{y}$$

# Backward Propagation



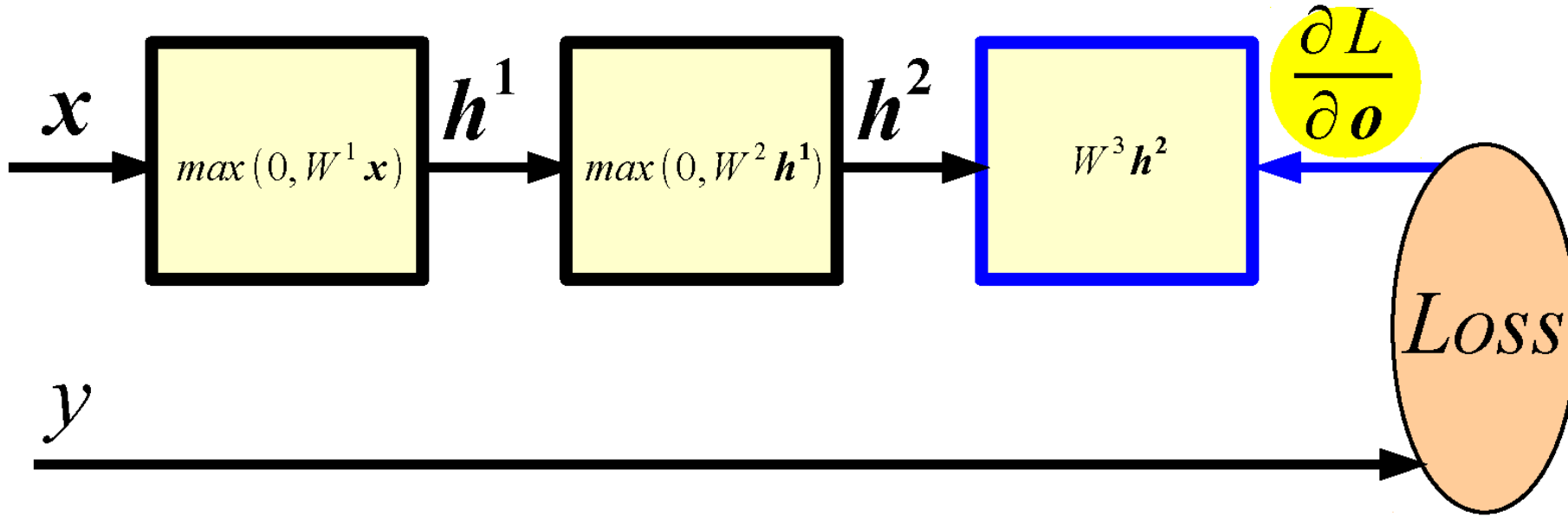
Given  $\frac{\partial L}{\partial o}$  and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial h^2}$$



# Backward Propagation



Given  $\frac{\partial L}{\partial o}$  and assuming we can easily compute the Jacobian of each module, we have:

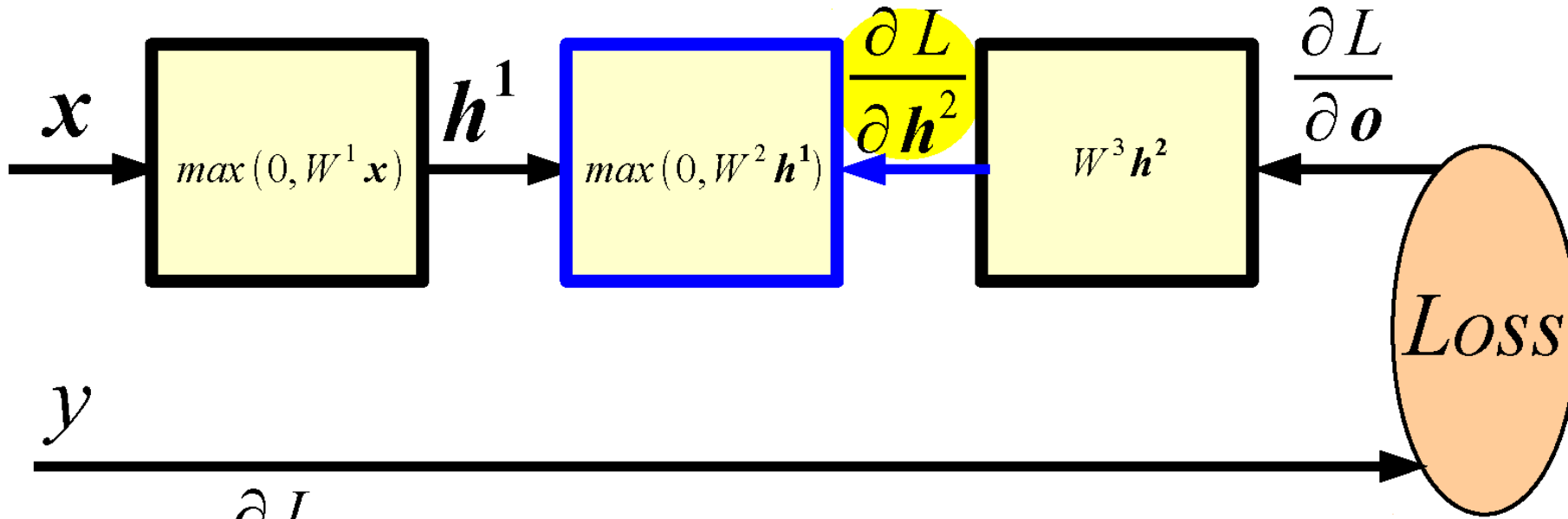
$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial W^3}$$

$$\frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial o} \frac{\partial o}{\partial h^2}$$

$$\frac{\partial L}{\partial W^3} = (p(c|x) - y) h^{2T}$$

$$\frac{\partial L}{\partial h^2} = W^{3T} (p(c|x) - y)$$

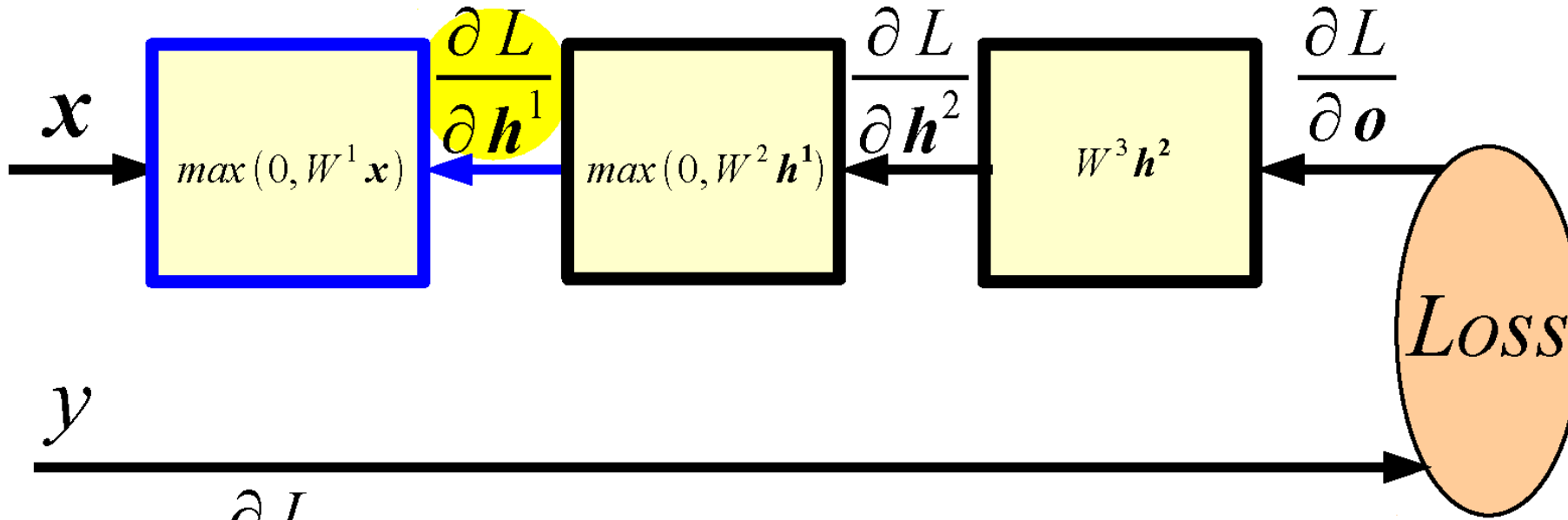
# Backward Propagation



Given  $\frac{\partial L}{\partial \mathbf{h}^2}$  we can compute now:

$$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial \mathbf{h}^2} \frac{\partial \mathbf{h}^2}{\partial W^2} \quad \frac{\partial L}{\partial \mathbf{h}^1} = \frac{\partial L}{\partial \mathbf{h}^2} \frac{\partial \mathbf{h}^2}{\partial \mathbf{h}^1}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial \mathbf{h}^1}$  we can compute now:

$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial \mathbf{h}^1} \frac{\partial \mathbf{h}^1}{\partial W^1}$$

# Backward Propagation

**Question:** Does BPROP work with ReLU layers only?

**Answer:** Nope, any a.e. differentiable transformation works.

**Question:** What's the computational cost of BPROP?

**Answer:** About twice FPROP (need to compute gradients w.r.t. input and parameters at every layer).

# Optimization

**Stochastic Gradient Descent (on mini-batches):**

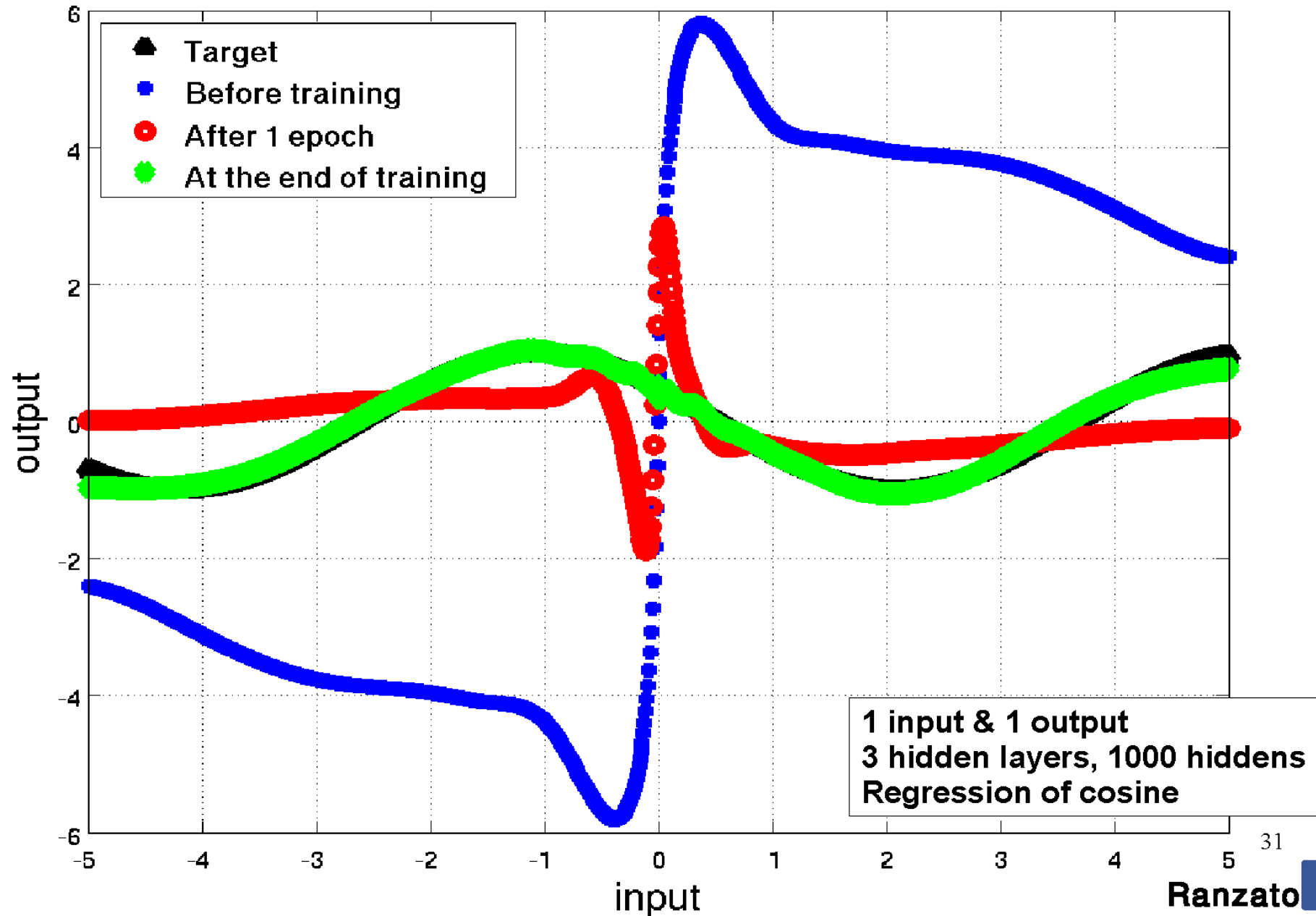
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \frac{\partial L}{\partial \boldsymbol{\theta}}, \eta \in (0, 1)$$

**Stochastic Gradient Descent with Momentum:**

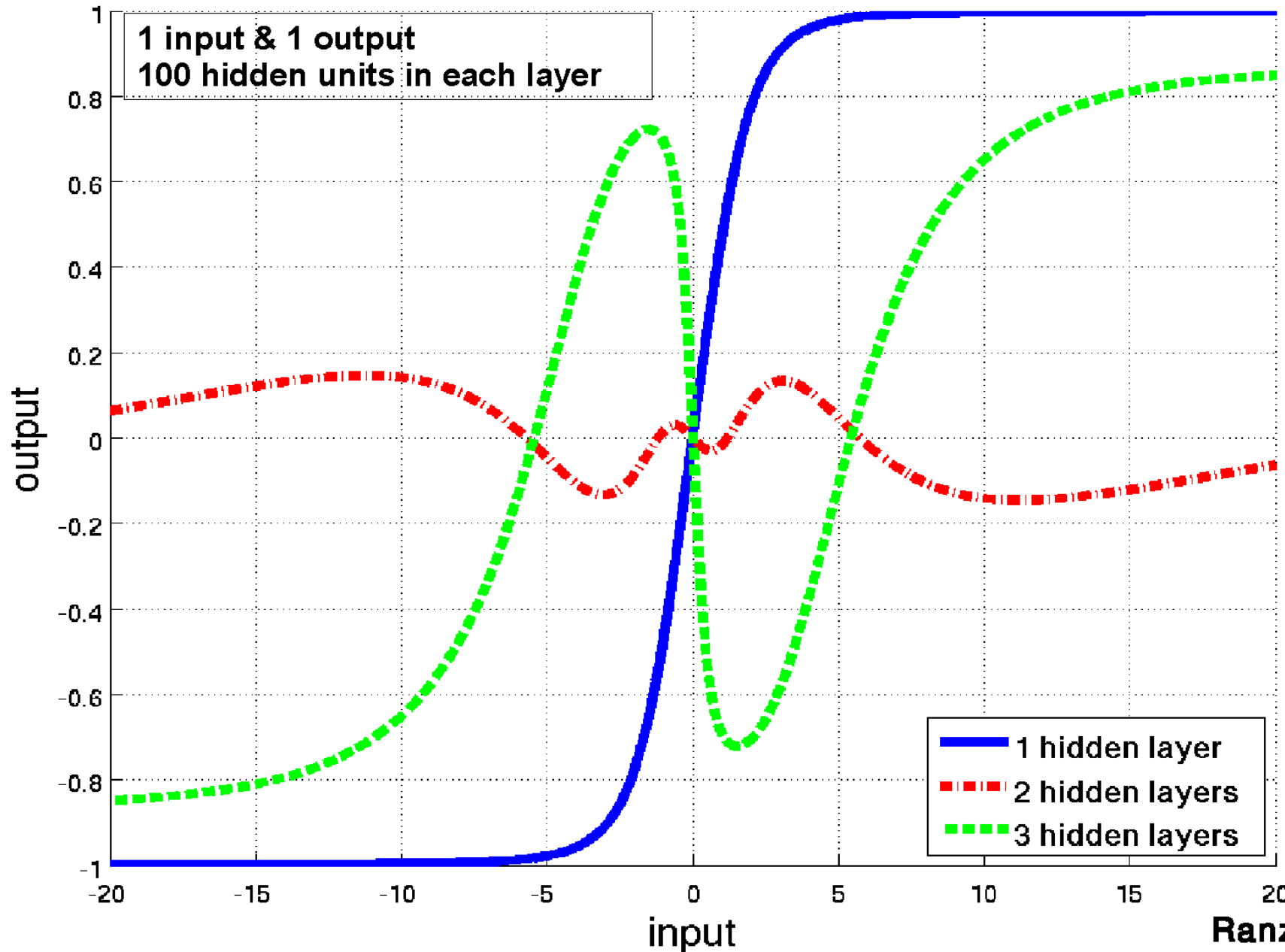
$$\begin{aligned}\boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \eta \boldsymbol{\Delta} \\ \boldsymbol{\Delta} &\leftarrow 0.9 \boldsymbol{\Delta} + \frac{\partial L}{\partial \boldsymbol{\theta}}\end{aligned}$$

**Note: there are many other variants...**

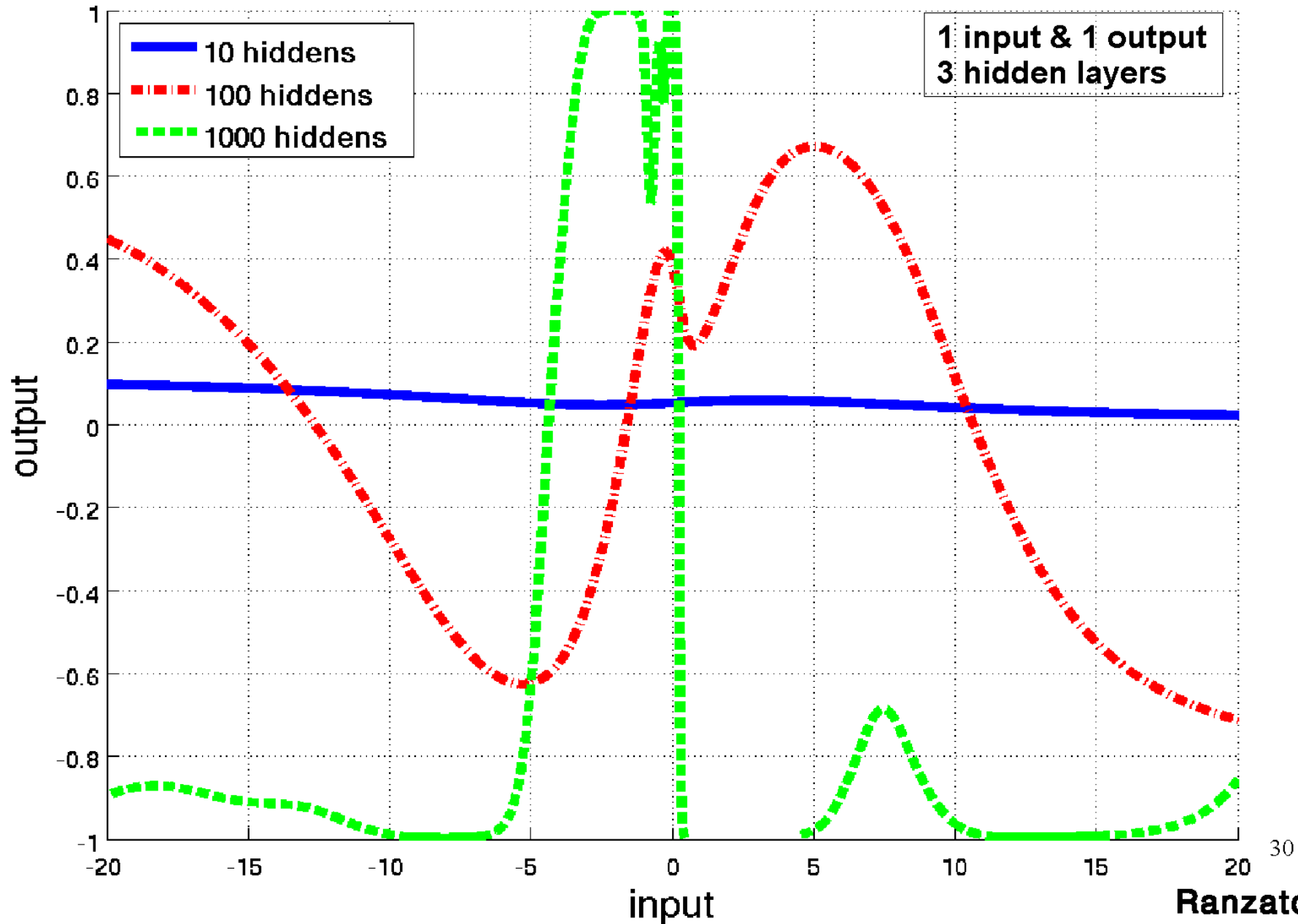
# Toy Example: Synthetic Data



# Toy Example: Synthetic Data



# Toy Example: Synthetic Data





# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration
Neural Network	1
Nearest Neighbor	10
Linear SVM	10
Non-linear SVM	5
Decision Tree or Random Forest	4

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration	Ease of interpretation
Neural Network	1	1
Nearest Neighbor	10	10
Linear SVM	10	9
Non-linear SVM	5	4
Decision Tree or Random Forest	4	4

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration	Ease of interpretation	Speed / memory when training
Neural Network	1	1	1
Nearest Neighbor	10	10	8
Linear SVM	10	9	10
Non-linear SVM	5	4	2
Decision Tree or Random Forest	4	4	4

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration	Ease of interpretation	Speed / memory when training	Speed / memory at test time
Neural Network	1	1	1	6
Nearest Neighbor	10	10	8	4
Linear SVM	10	9	10	10
Non-linear SVM	5	4	2	2
Decision Tree or Random Forest	4	4	4	8

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration	Ease of interpretation	Speed / memory when training	Speed / memory at test time	Accuracy w/ lots of data
Neural Network	1	1	1	6	10
Nearest Neighbor	10	10	8	4	7
Linear SVM	10	9	10	10	5
Non-linear SVM	5	4	2	2	8
Decision Tree or Random Forest	4	4	4	8	7

This all seems pretty complicated. Why are we using Neural Networks? James's rough assessment:

Learning method	Ease of configuration	Ease of interpretation	Speed / memory when training	Speed / memory at test time	Accuracy w/ lots of data
Neural Network	1	1	1	6	10
Nearest Neighbor	10	10	8	4	7
Linear SVM	10				
Non-linear SVM	5				
Decision Tree or Random Forest	4				

Representation design matters more for all of these

# Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips