# High Level Computer Vision

# Backpropagation & Convolutional Neural Networks
# @ April 24, 2019

**Bernt Schiele & Mario Fritz**

**www.mpi-inf.mpg.de/hlcv/**

**Max Planck Institute for Informatics & Saarland University,**
**Saarland Informatics Campus Saarbrücken**

# Overview Today's Lecture

- Backpropagation - Gradient Descent

  ▸ illustrated using computational graphs

  ▸ chain rule - upstream and local gradients

  ▸ modularization simple

- What is Deep Learning

  ▸ intuition why deep learning can help

  ▸ integrated learning of features and classifier

- Convolutional Neural Networks (CNNs)

  ▸ one of the (few) highly successful NNs

# Where we are

$$s = f(x; W) = Wx$$  scores function

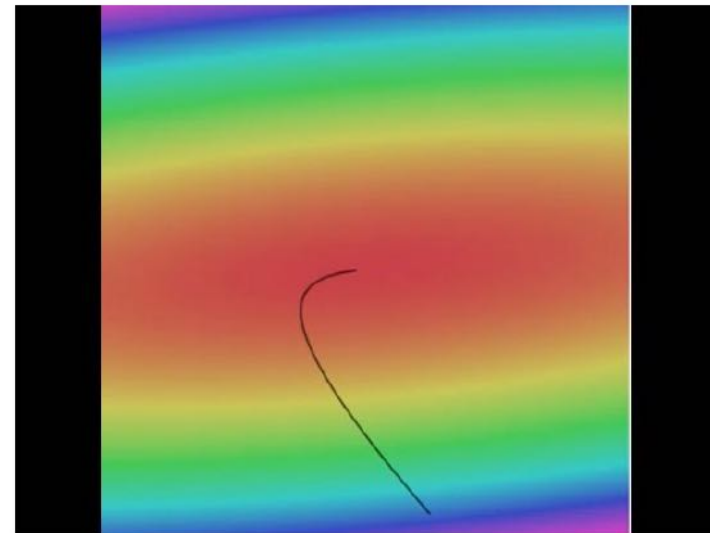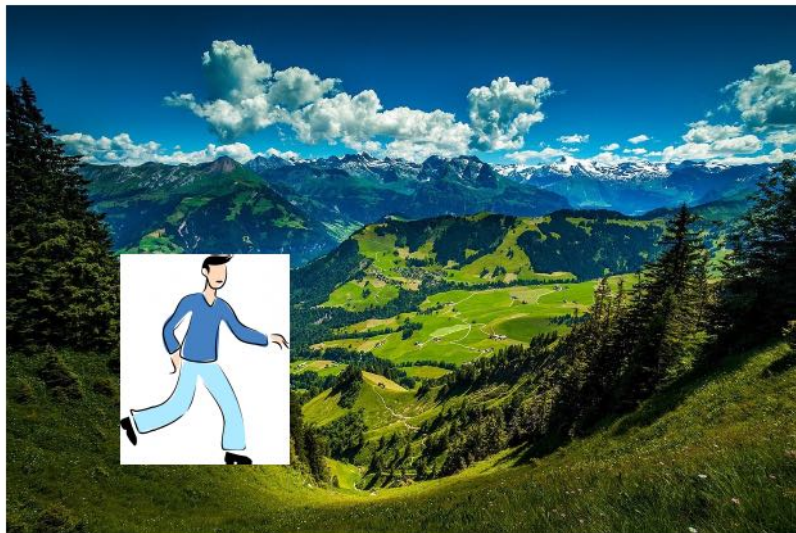$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$  SVM loss

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + \sum_k W_k^2$$  data loss + regularization

want  $\boxed{\nabla_W L}$

# Optimization





```
# Vanilla Gradient Descent

while True:
  weights_grad = evaluate_gradient(loss_fun, data, weights)
  weights += - step_size * weights_grad # perform parameter update
```

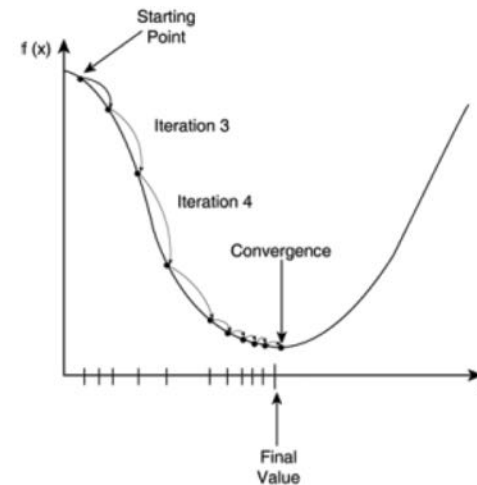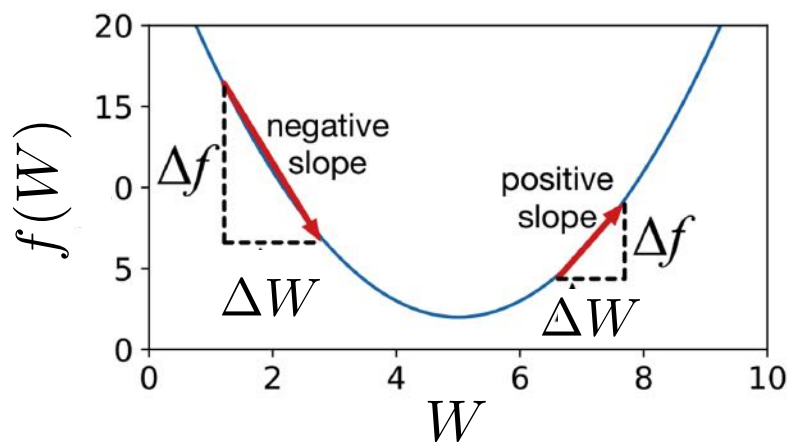slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Gradient Descent

**General gradient descent:** Start with initial point $W_0$

$$\text{Sequence:} \quad W_{t+1} = W_t + \alpha_t d_t$$

**Steepest Descent:**

$d_t = -\nabla f(W_t)$ (we move in the opposite direction of the gradient).



```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# Gradient Descent - Variants…

- Assume Loss to be:

$$L(W) = \frac{1}{n} \sum_{i=1}^{n} L_i(W)$$

  - ▸ with n the number of training samples
  - ▸ $L_i$ the loss for training sample $x_i$

- **Stochastic Gradient Descent**:

  - ▸ randomly choose one training sample $x_i$
  - ▸ update weights based on loss $L_i(W)$

- **Mini-batch training**:

  - ▸ process a subset of training samples $M \subset \{1, \ldots, n\}$
  - ▸ update weights based on $L_M(W) = \frac{1}{|M|} \sum_{i \in M} L_i(W)$

- **Batch training**:

  - ▸ process all training samples
  - ▸ update weights based on $L(W) = \frac{1}{n} \sum_{i=1}^{n} L_i(W)$

# Computational Graphs



$$f = Wx \qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$R(W)$$

# Neural Network Example…

Convolutional network
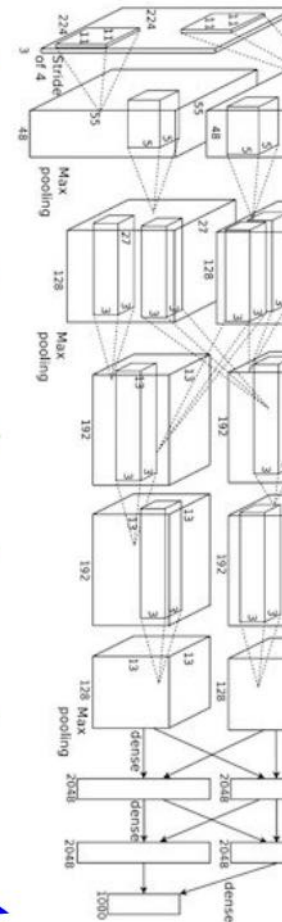(AlexNet)

input image

weights

loss



Figure copyright Alex Krizhevsky, Ilya Sutskever, and
Geoffrey Hinton, 2012. Reproduced with permission.

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Neural Network Example…
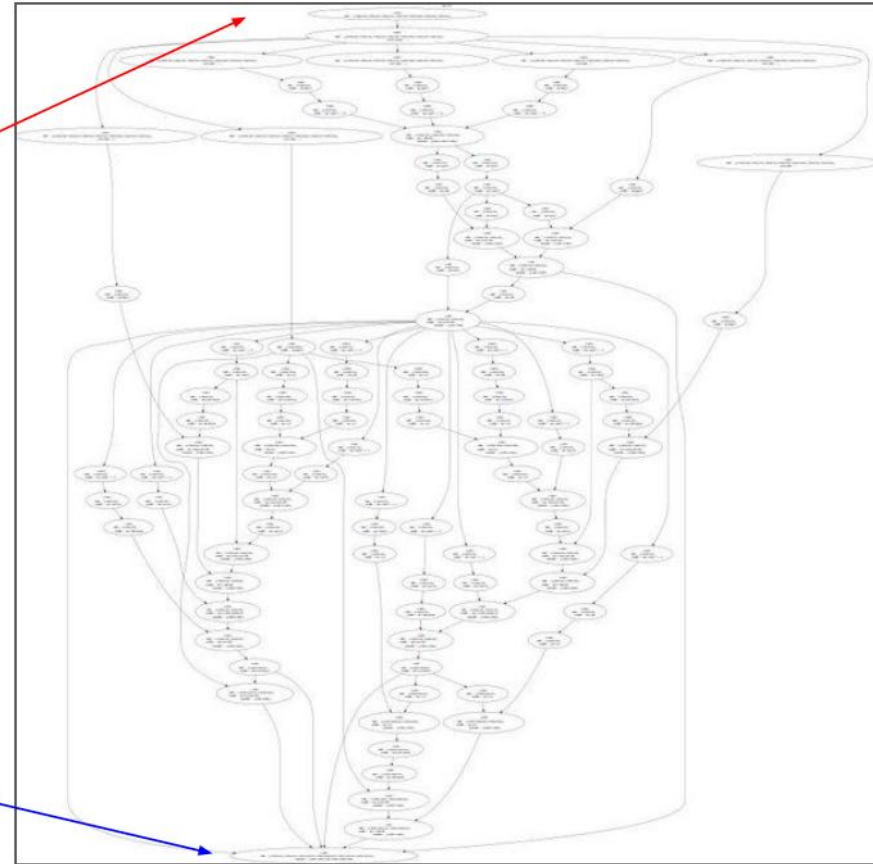


Neural Turing Machine

input image

loss

Figure reproduced with permission from a Twitter post by Andrej Karpathy.

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4



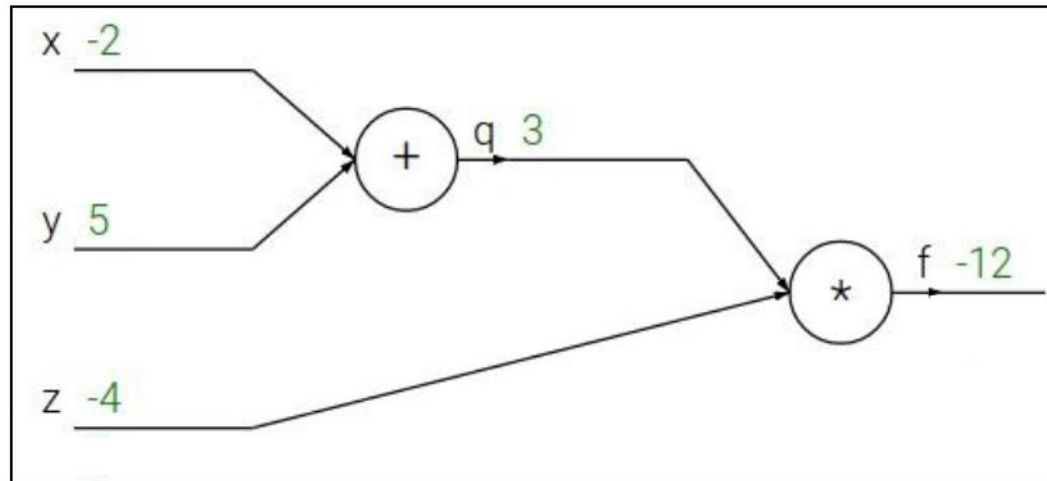$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial f}$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
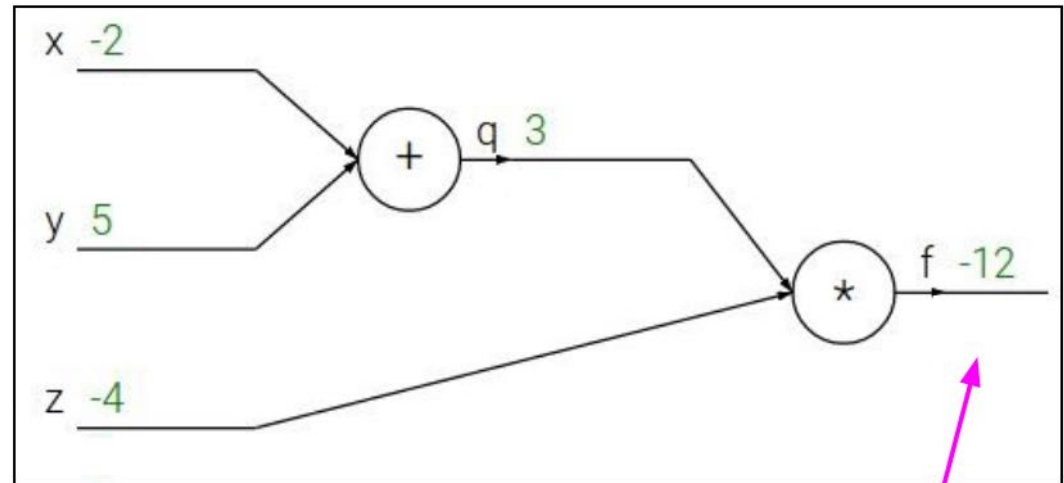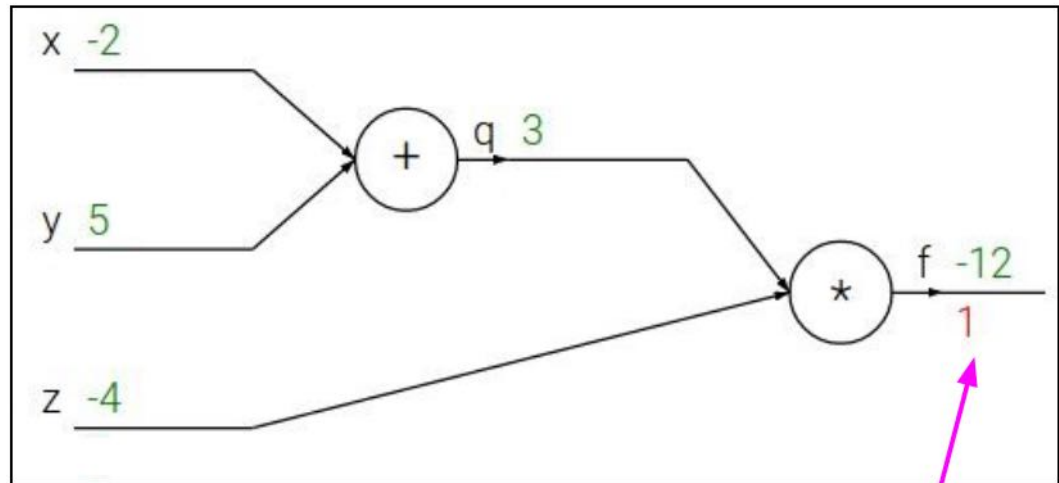
slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$\dfrac{\partial f}{\partial f}$

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

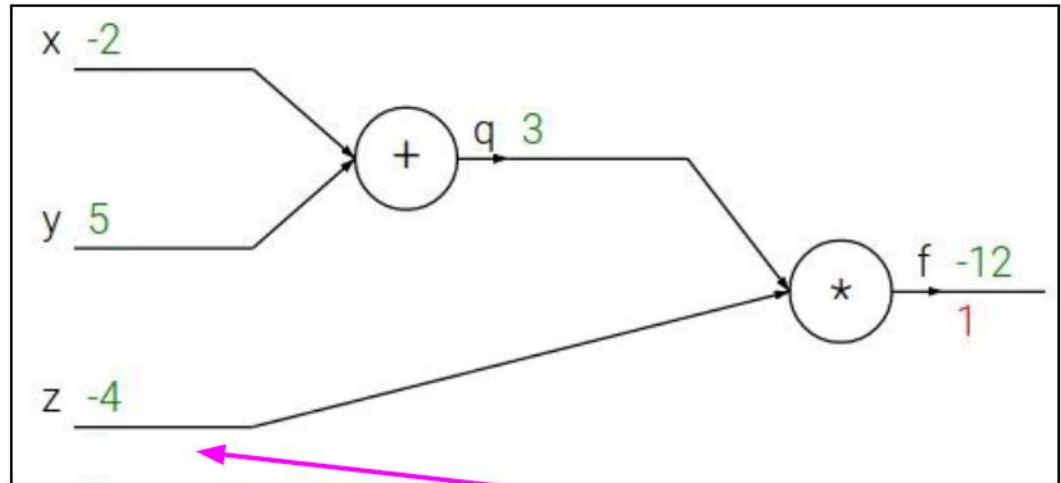Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

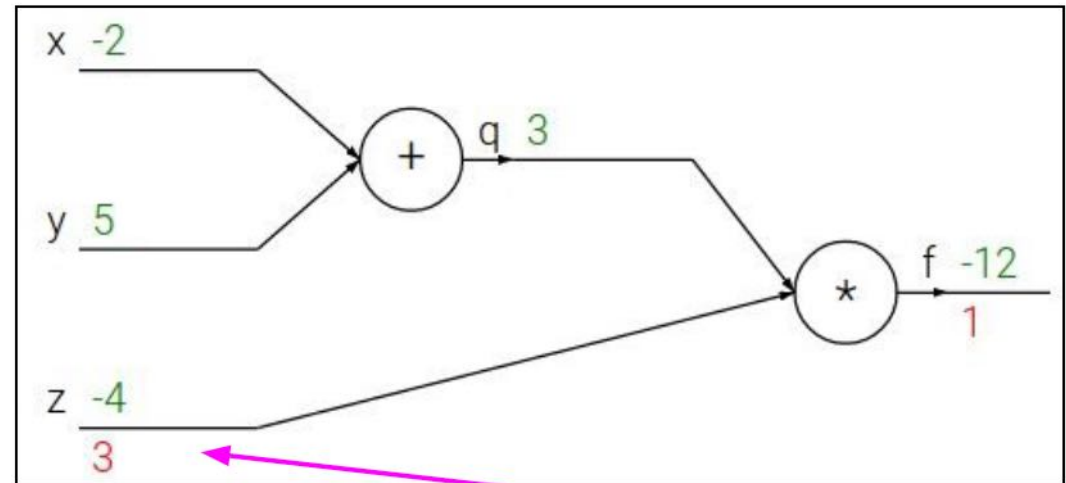# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial q}$$

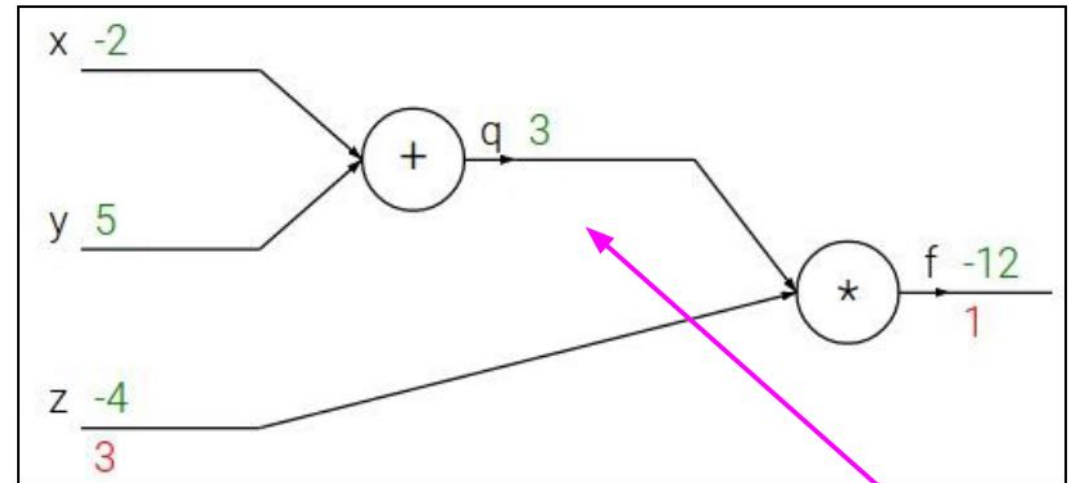# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial q}$$

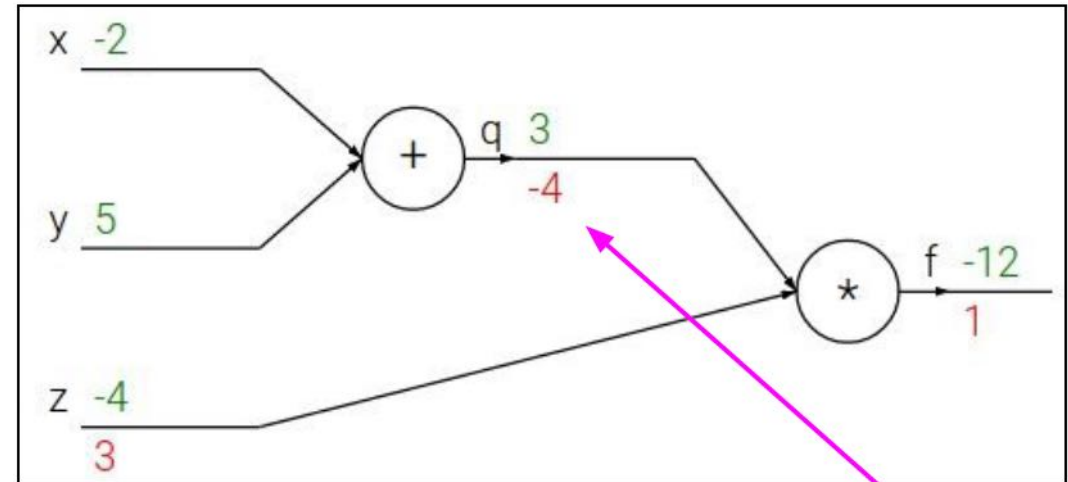# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$\dfrac{\partial f}{\partial y}$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream gradient       Local gradient

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

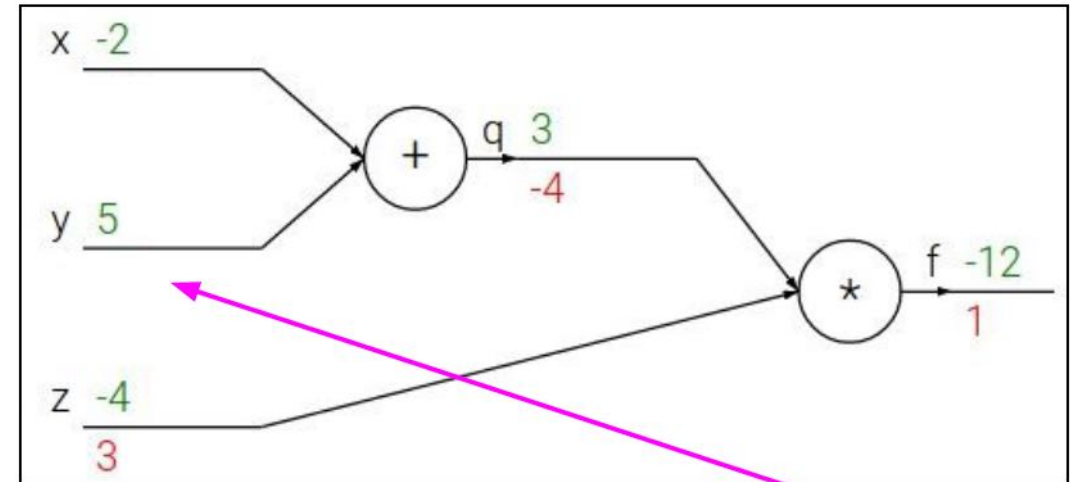# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial y}$$

Chain rule:

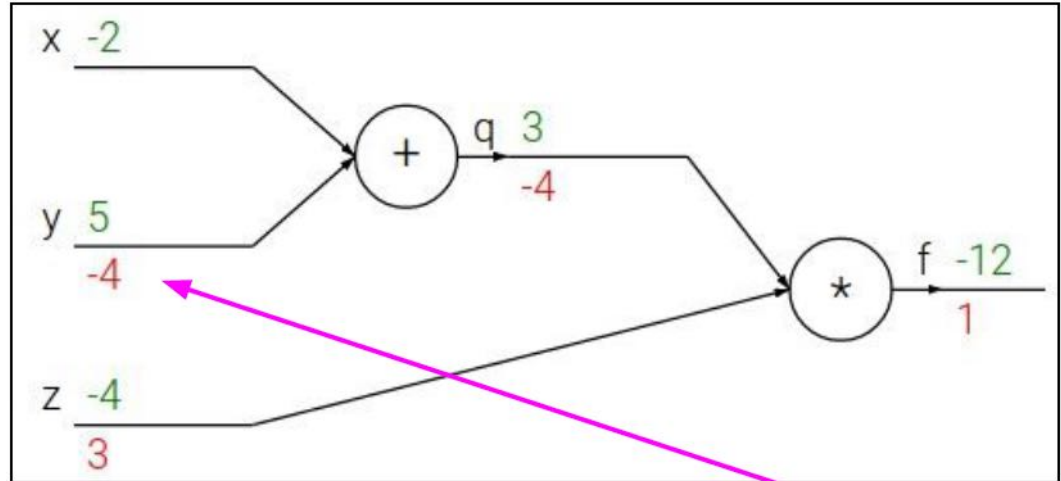$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream    Local
gradient    gradient

# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial x}$$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream gradient    Local gradient

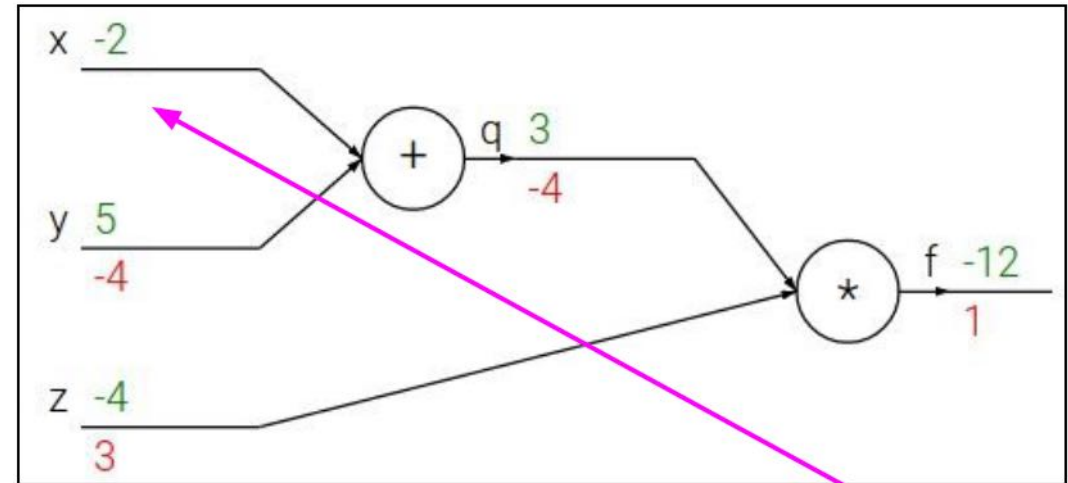# Backpropagation - A Simple Example

Backpropagation: a simple example

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

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

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



$$\frac{\partial f}{\partial x}$$

Chain rule:

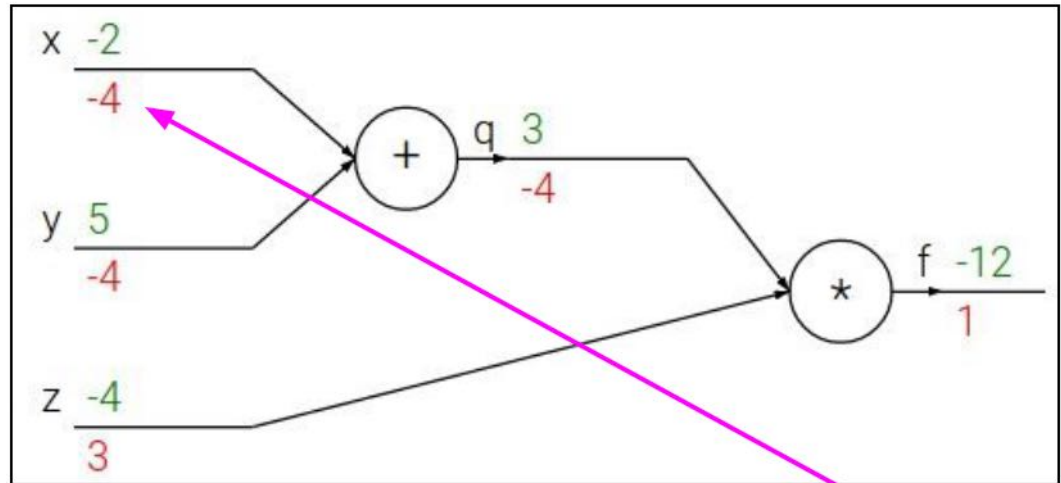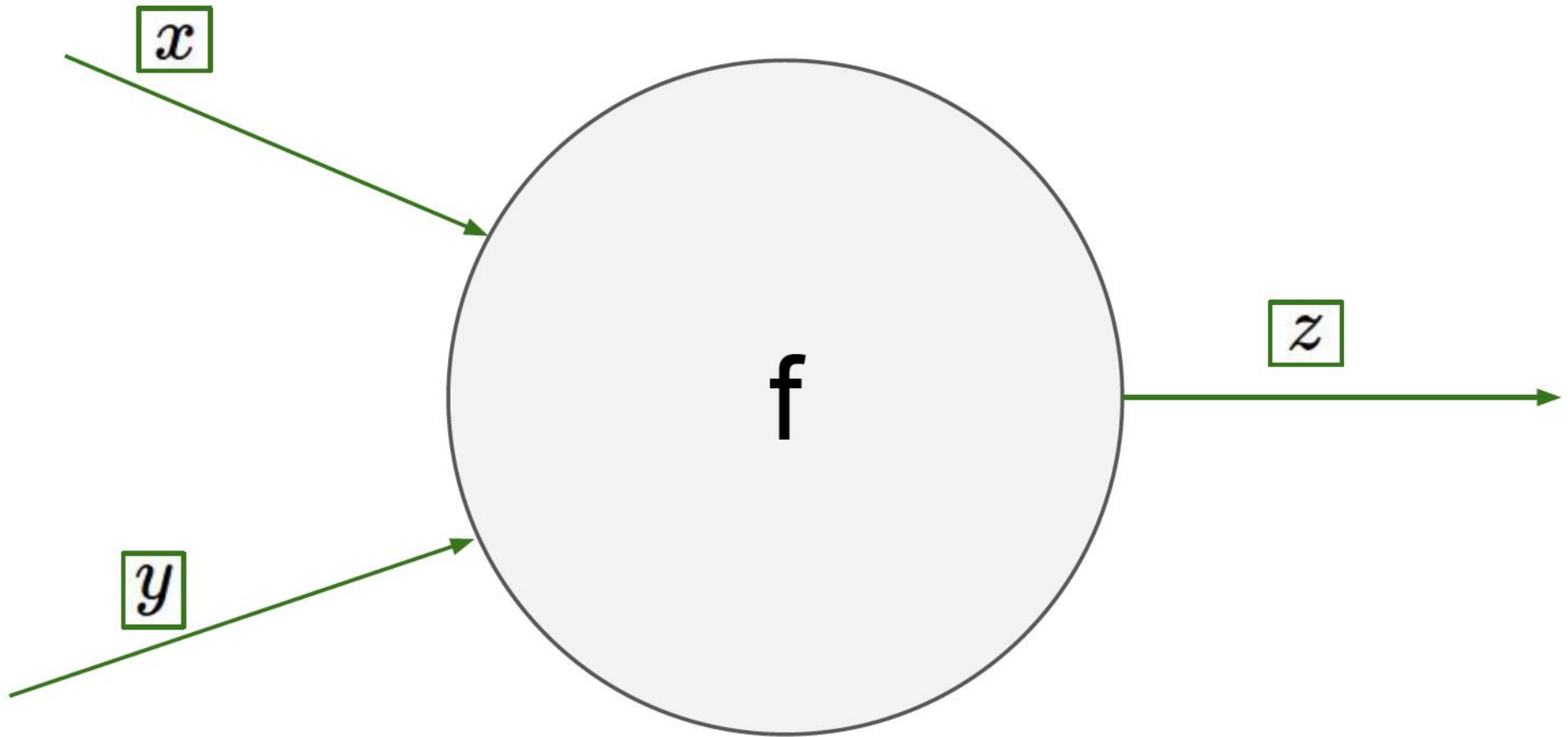$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream gradient    Local gradient

# Backpropagation - Local View

# Backpropagation - Local View

# Backpropagation - Local View



"local gradient"

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

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

$x$

$y$

$f$

$z$

$$\frac{\partial L}{\partial z}$$

"upstream gradient"

gradients

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Backpropagation - Local View



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

"local gradient"

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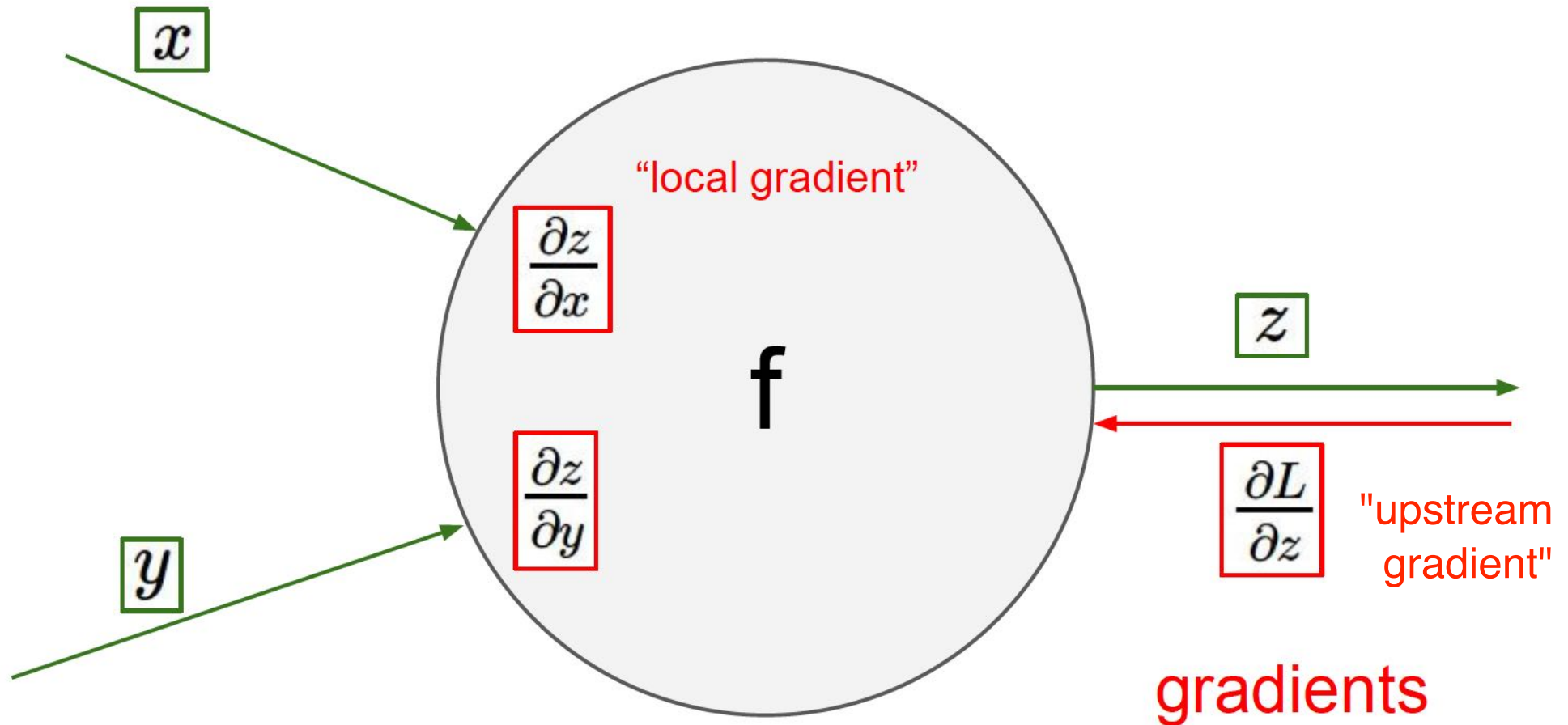
$$\frac{\partial z}{\partial y}$$

f

x

y

z

$$\frac{\partial L}{\partial z}$$

"upstream gradient"

gradients

# Backpropagation - Local View



"local gradient"

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

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

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

f

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

$x$

$y$

$z$

$$\frac{\partial L}{\partial z}$$

"upstream gradient"

gradients

# Backpropagation - Local View



"local gradient"

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

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

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

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

f

$x$

$y$

$z$

$$\frac{\partial L}{\partial z}$$

"upstream gradient"

gradients

# Backpropagation

Another example:  $f(w,x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

# Backpropagation

Another example:
$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

# Backpropagation

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example:  $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient    Local gradient

$$(1.00)\left(\frac{-1}{1.37^2}\right) = -0.53$$

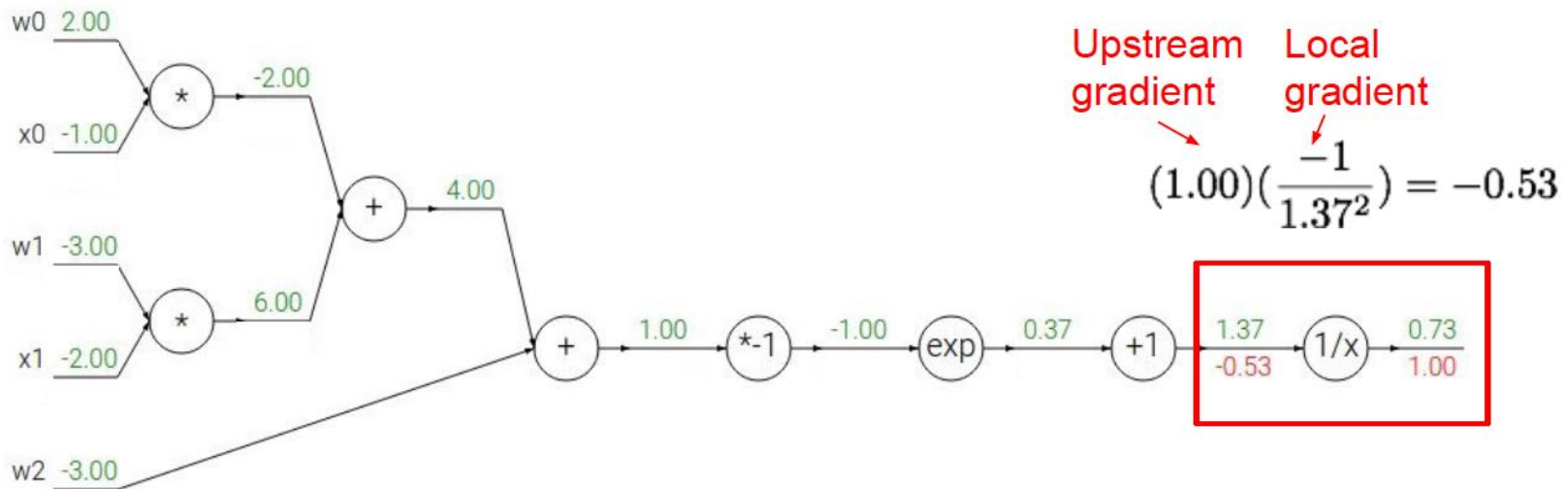| $f(x) = e^x$ | $\rightarrow$ | $\frac{df}{dx} = e^x$ | $f(x) = \frac{1}{x}$ | $\rightarrow$ | $\frac{df}{dx} = -1/x^2$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\frac{df}{dx} = a$ | $f_c(x) = c + x$ | $\rightarrow$ | $\frac{df}{dx} = 1$ |

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Backpropagation

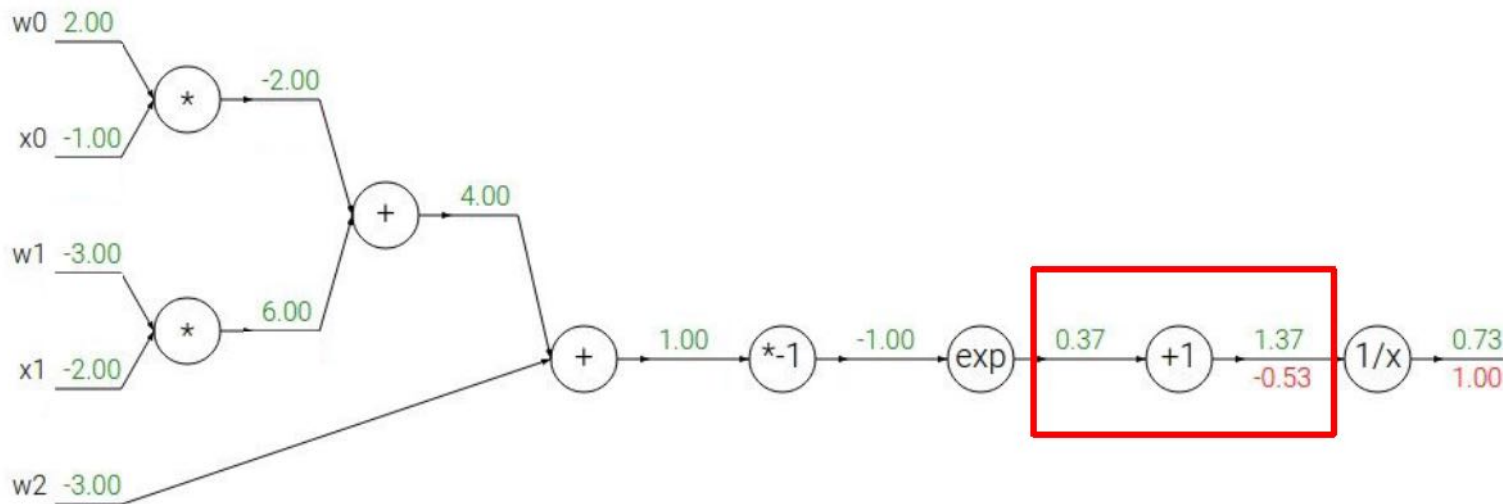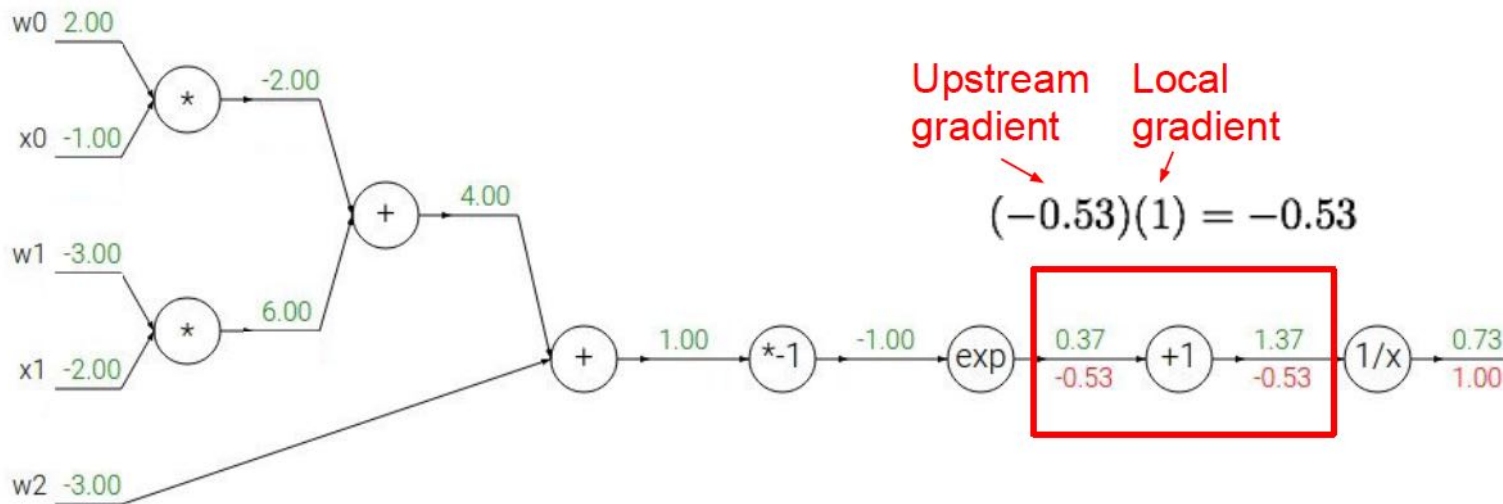Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x \qquad \Bigg| \qquad f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a \qquad \Bigg| \qquad f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

# Backpropagation

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Upstream gradient   Local gradient

$$(-0.53)(1) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example:   $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



Upstream gradient     Local gradient

$$(-0.53)(e^{-1}) = -0.20$$

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ |

| | | |
|---|---|---|
| $f(x) = \dfrac{1}{x}$ | $\rightarrow$ | $\dfrac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\dfrac{df}{dx} = 1$ |

# Backpropagation

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



Upstream gradient, Local gradient

$(-0.20)(-1) = 0.20$

$f(x) = e^x \qquad \rightarrow \qquad \dfrac{df}{dx} = e^x$

$f_a(x) = ax \qquad \rightarrow \qquad \dfrac{df}{dx} = a$

$f(x) = \dfrac{1}{x} \qquad \rightarrow \qquad \dfrac{df}{dx} = -1/x^2$

$f_c(x) = c + x \qquad \rightarrow \qquad \dfrac{df}{dx} = 1$

# Backpropagation

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \quad \bigg| \quad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \quad \bigg| \quad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Backpropagation

Another example:
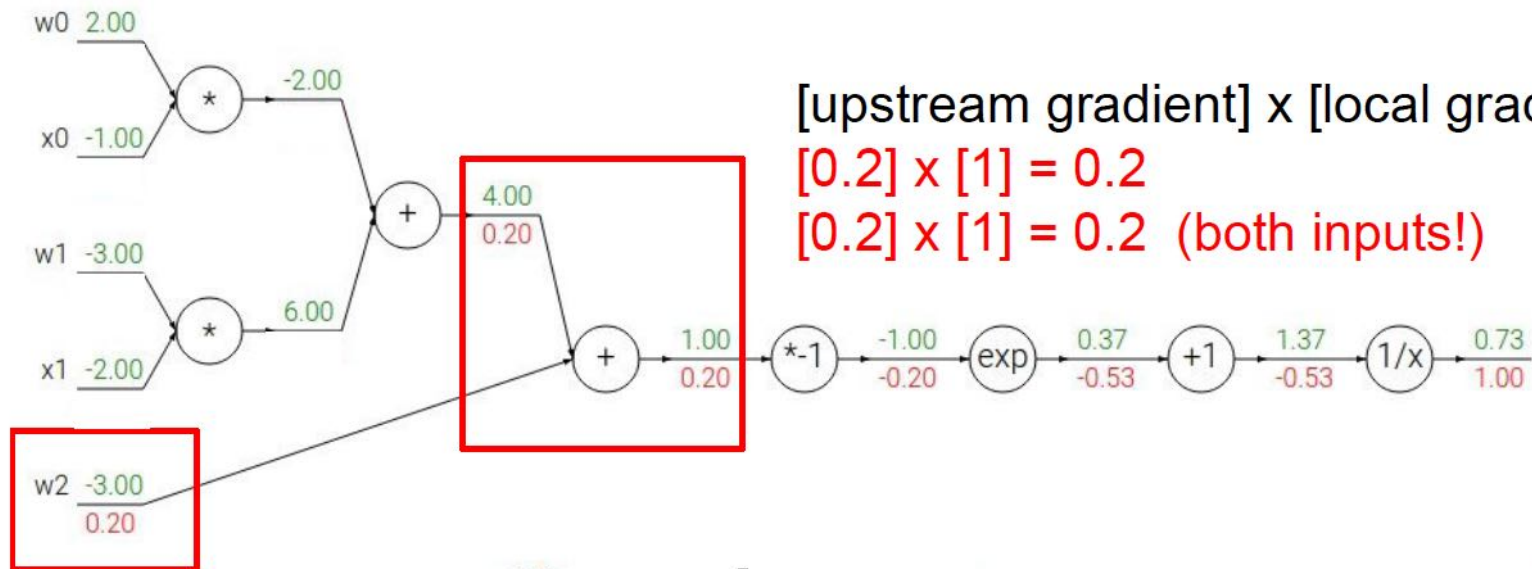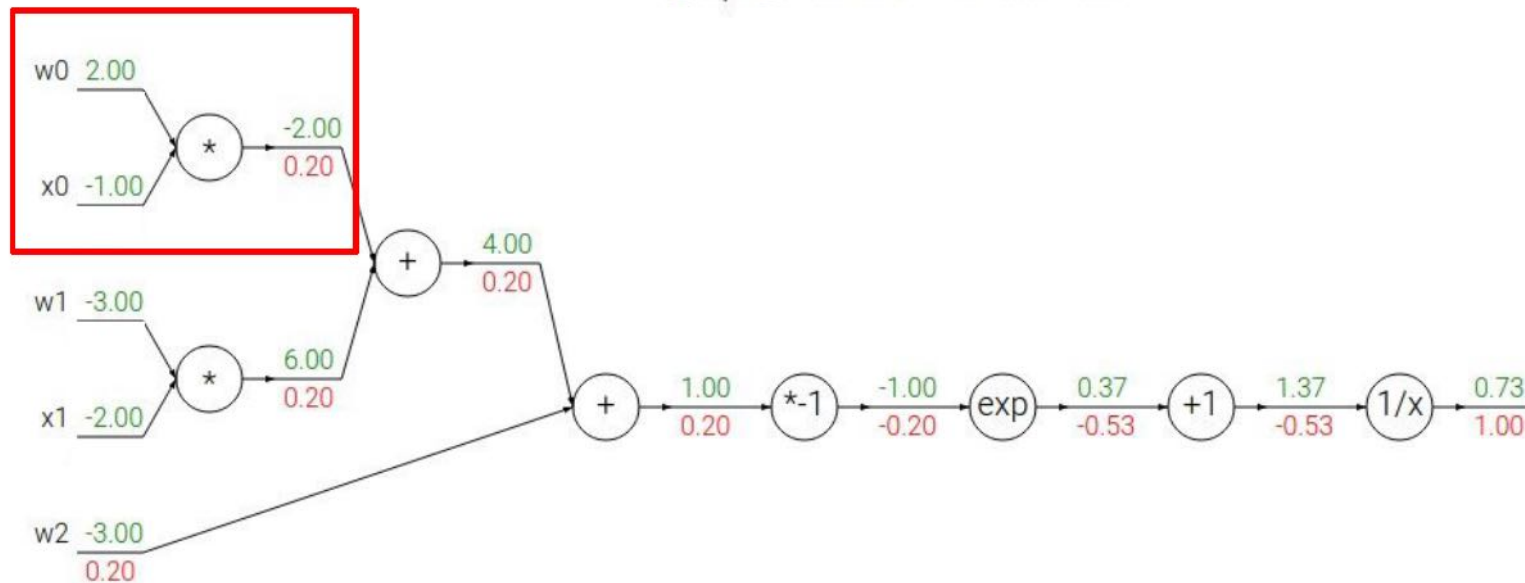$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
[0.2] x [1] = 0.2
[0.2] x [1] = 0.2  (both inputs!)

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Backpropagation

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[upstream gradient] x [local gradient]
x0: [0.2] x [2] = 0.4
w0: [0.2] x [-1] = -0.2

$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

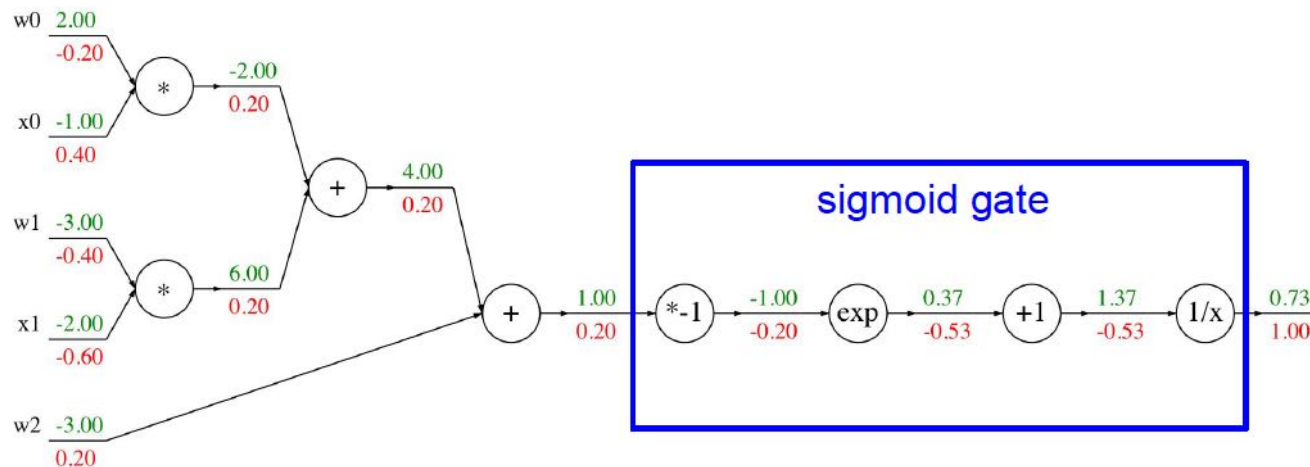$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

# Backpropagation

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

$$\boxed{\sigma(x) = \frac{1}{1 + e^{-x}}} \qquad \frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right)\left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\,\sigma(x)$$
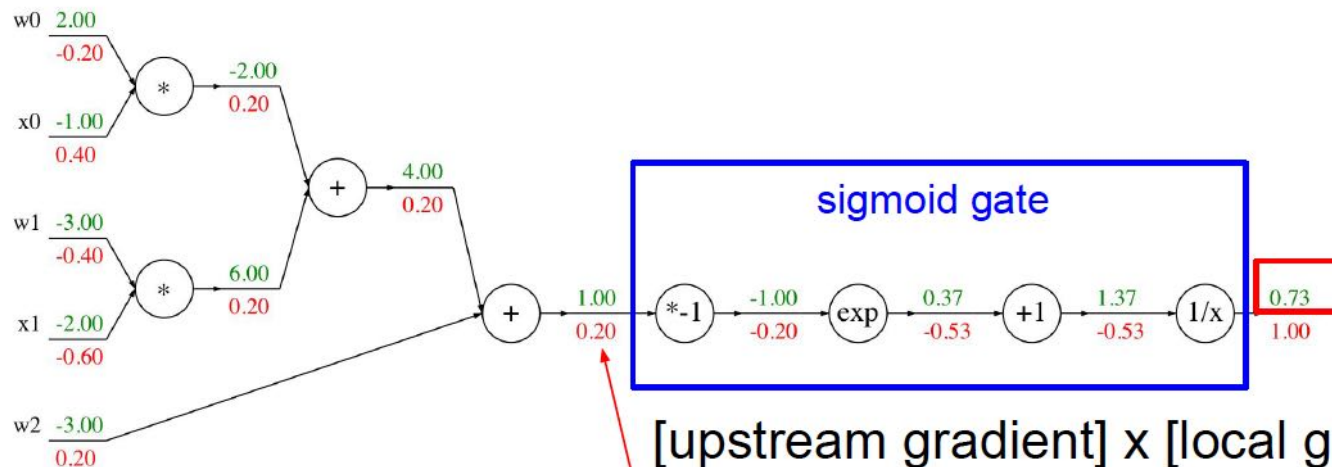
sigmoid function



sigmoid gate

# Backpropagation

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

$$\boxed{\sigma(x) = \frac{1}{1 + e^{-x}}} \qquad \frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right)\left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\sigma(x)$$
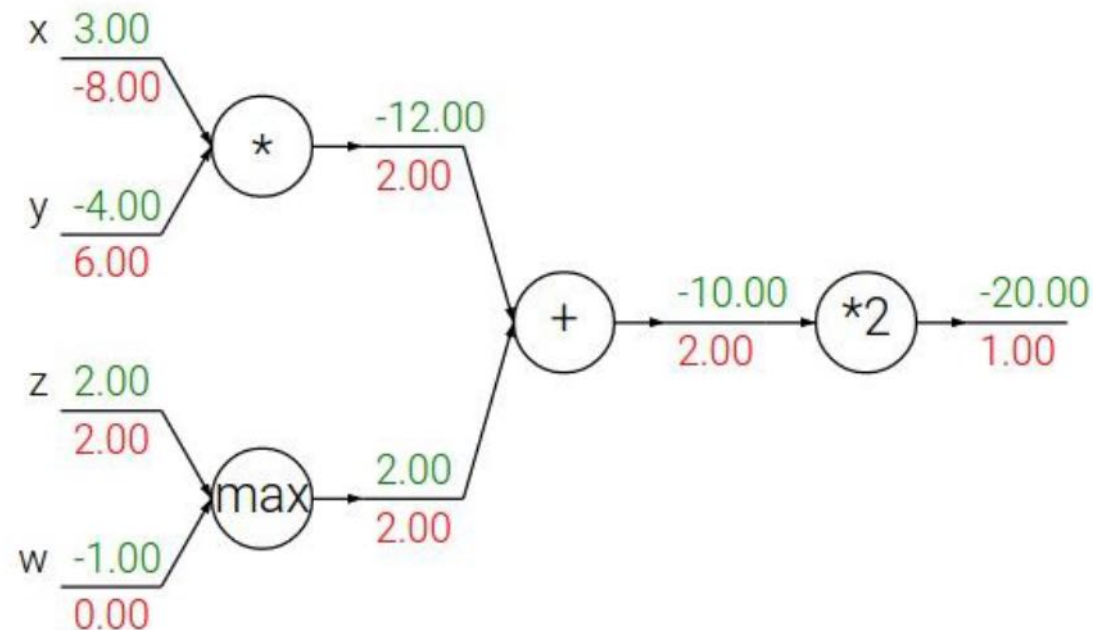
sigmoid function



sigmoid gate

[upstream gradient] x [local gradient]
[1.00] x [(1 - 0.73) (0.73)]= 0.2

# Backpropagation

## Patterns in backward flow

**add** gate: gradient distributor

# Backpropagation

## Patterns in backward flow

**add** gate: gradient distributor

Q: What is a **max** gate?

# Backpropagation

## Patterns in backward flow

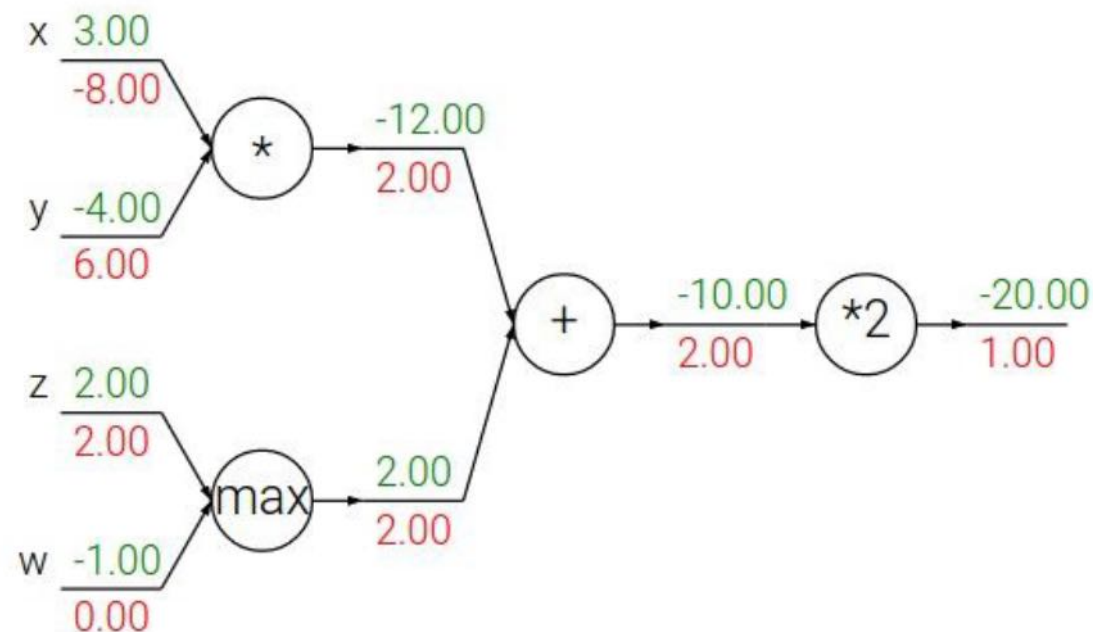**add** gate: gradient distributor

**max** gate: gradient router

# Backpropagation

## Patterns in backward flow

**add** gate: gradient distributor

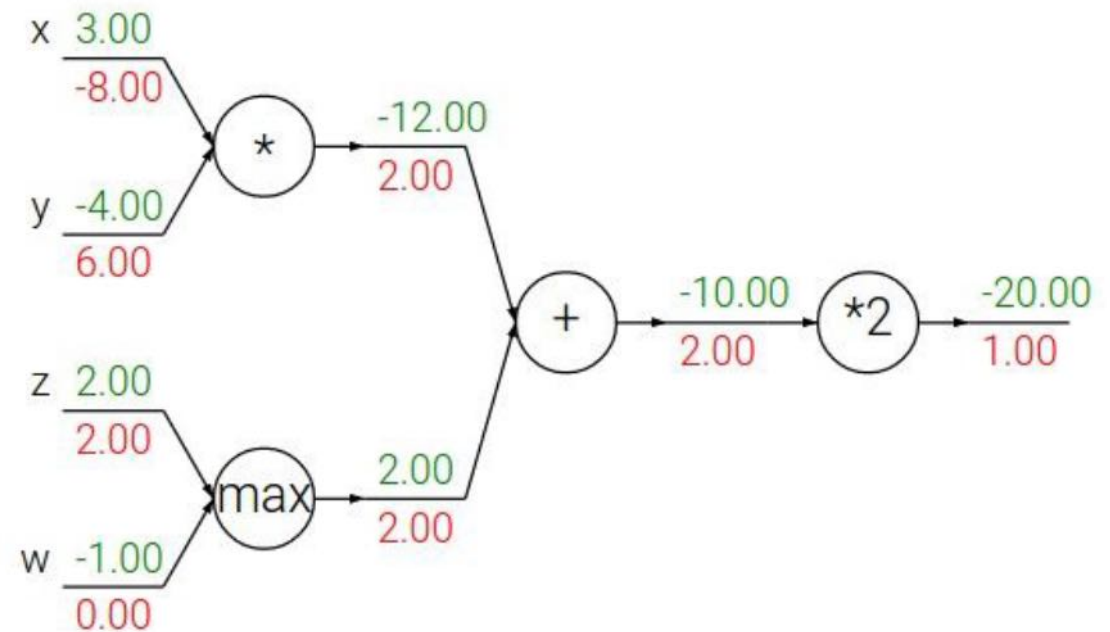**max** gate: gradient router

Q: What is a **mul** gate?

# Backpropagation

## Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

**mul** gate: gradient switcher

# Backpropagation

Gradients add at branches

# Backpropagation

Gradients for vectorized code

(x,y,z are now vectors)

This is now the **Jacobian matrix** (derivative of each element of z w.r.t. each element of x)



$x$

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

"local gradient"

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

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

f

$y$

$z$

$\frac{\partial L}{\partial z}$

gradients

# Backpropagation

## Modularized implementation: forward / backward API



Graph (or Net) object  *(rough pseudo code)*

```python
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Backpropagation

## Modularized implementation: forward / backward API
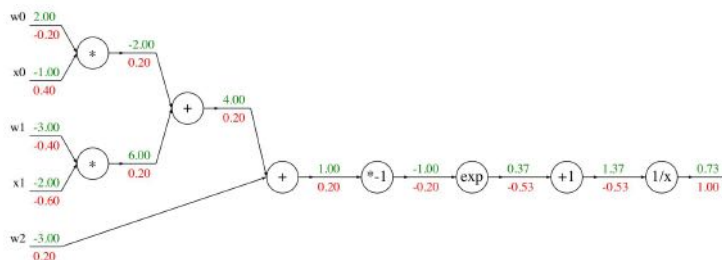


```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

X

Z

*

y

(x,y,z are scalars)

Local gradient          Upstream gradient variable

# Gradient & Training

- Only two things need to be implemented to define new layer:
  - ▶ forward pass
  - ▶ error back propagation

- Watch overfitting
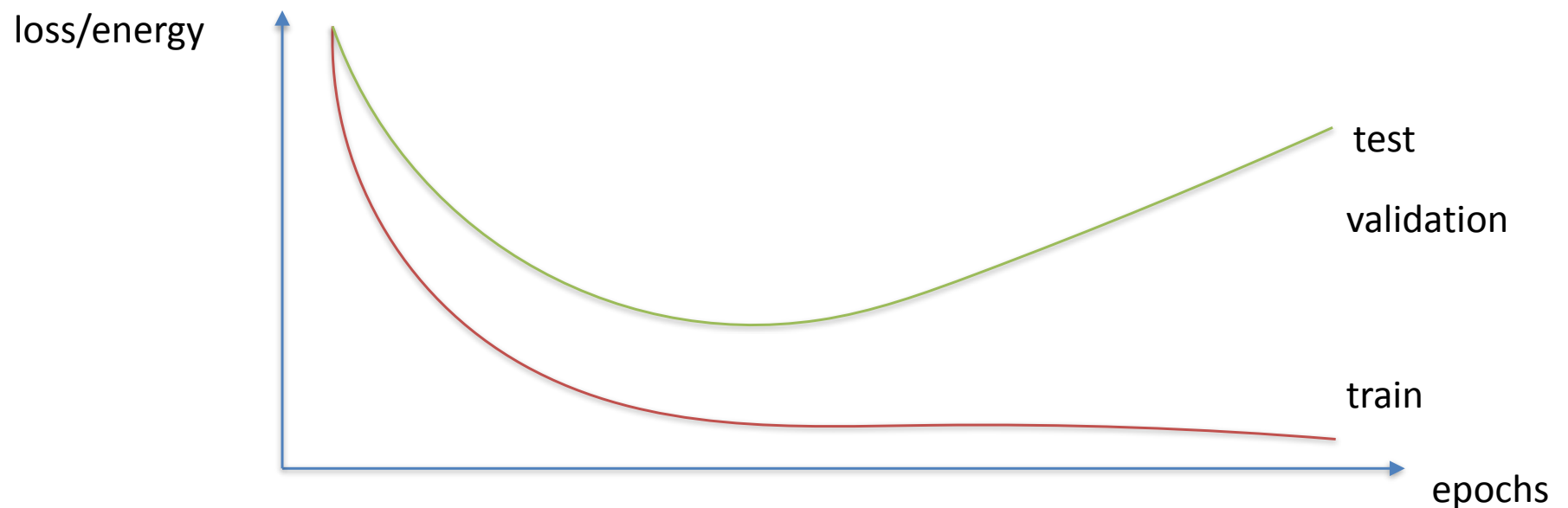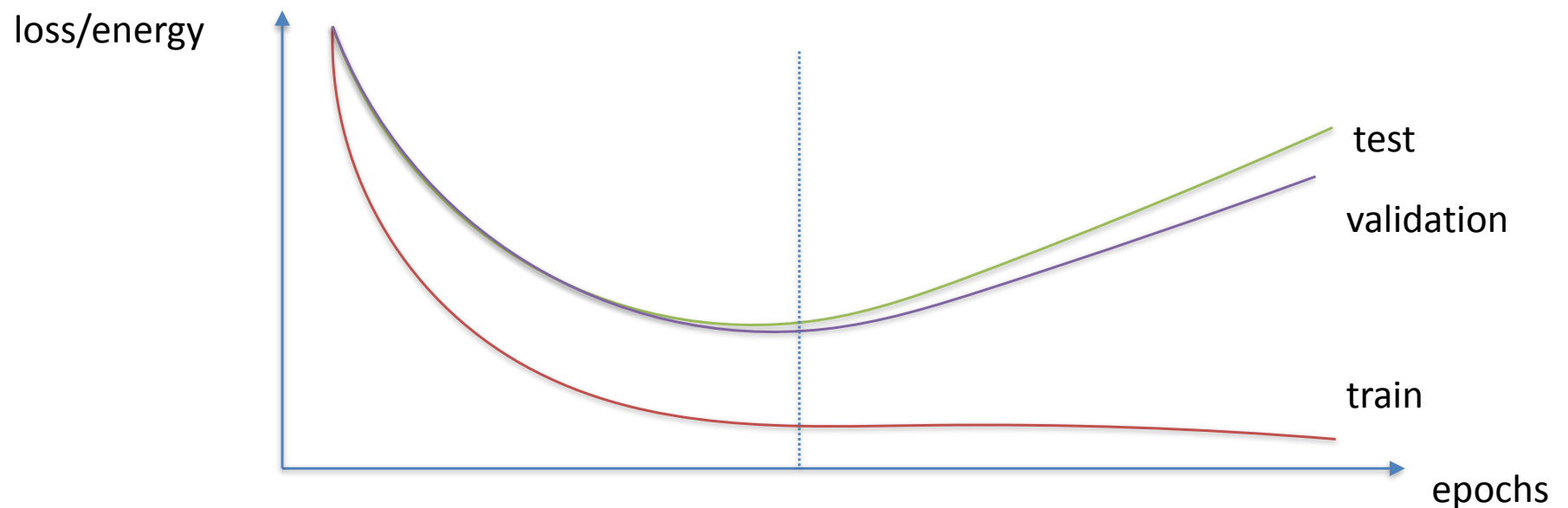
loss/energy

test

validation

train

epochs

# Gradient & Training

- Only two things need to be implemented to define new layer:

  ▶ forward pass

  ▶ error back propagation

- Watch overfitting

# Gradient & Training

- Only two things need to be implemented to define new layer:

  ▶ forward pass

  ▶ error back propagation

- Watch overfitting

# Summary so far...

- neural nets will be very large: impractical to write down gradient formula by hand for all parameters
- **backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates
- implementations maintain a graph structure, where the nodes implement the **forward**() / **backward**() API
- **forward**: compute result of an operation and save any intermediates needed for gradient computation in memory
- **backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs

# Overview Today's Lecture

- Backpropagation - Gradient Descent

  ‣ illustrated using computational graphs

  ‣ chain rule - upstream and local gradients

  ‣ modularization simple

- **What is Deep Learning**

  ‣ intuition why deep learning can help

  ‣ integrated learning of features and classifier

- Convolutional Neural Networks (CNNs)

  ‣ one of the (few) highly successful NNs

# Simple Neural Networks

(**Before**) Linear score function: $f = Wx$

# Simple Neural Networks
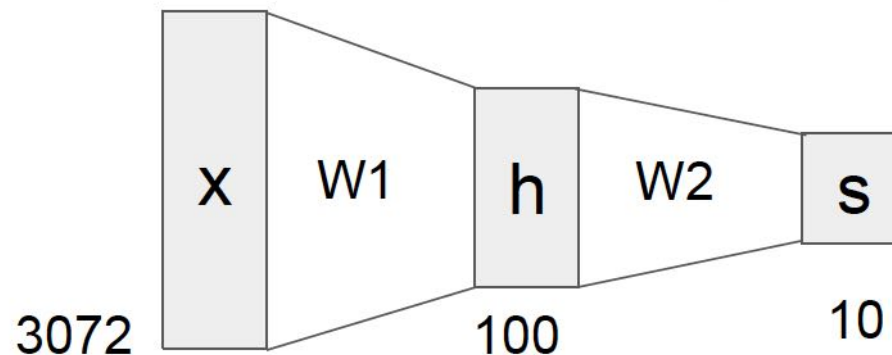
(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

# Simple Neural Networks

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

# Simple Neural Networks

(**Before**) Linear score function: $f = Wx$

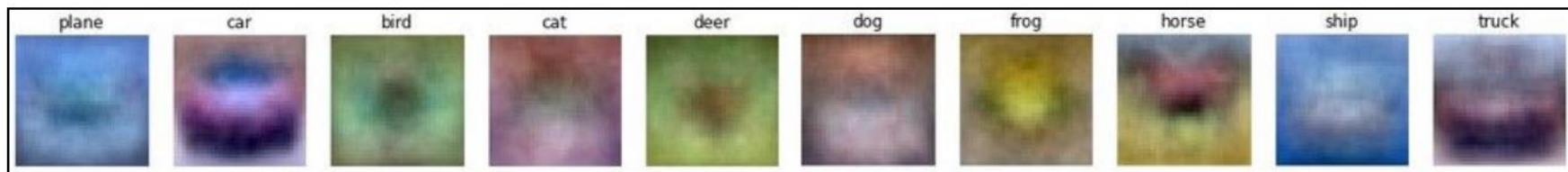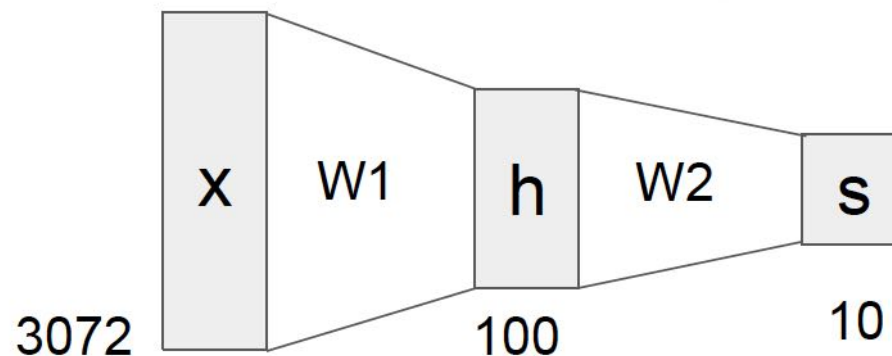(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

# Simple Neural Networks

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $\quad f = W_2 \max(0, W_1 x)$
or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

# Some Inspiration from the (Human) Brain



Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Some Inspiration from the (Human) Brain



Impulses carried toward cell body

dendrite

presynaptic terminal

cell body

axon

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

$x_0$

$w_0$

axon from a neuron

synapse

$w_0 x_0$

dendrite

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$

$f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Some Inspiration from the (Human) Brain

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

$x_0$     $w_0$

axon from a neuron     synapse

$w_0 x_0$

dendrite

cell body

$f\left(\sum_i w_i x_i + b\right)$

$w_1 x_1$

$\sum_i w_i x_i + b$     $f$

output axon

activation function

$w_2 x_2$

sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Some Inspiration from the (Human) Brain

Be very careful with your brain analogies!

**Biological Neurons:**
- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

[Dendritic Computation. London and Hausser]

# Activation Functions

## Activation functions

**Sigmoid**

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



**tanh**

$\tanh(x)$



**ReLU**

$\max(0, x)$



**Leaky ReLU**

$\max(0.1x, x)$



**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

# Neural Networks - Classic Architectures

- (Multi-Layer) Perceptron



input layer
hidden layer
output layer

"2-layer Neural Net", or
"1-hidden-layer Neural Net"

**"Fully-connected" layers**

input layer
hidden layer 1    hidden layer 2
output layer

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

# Overview of Neural Network Architectures



SUPERVISED

Recurrent Neural Net

Boosting

Convolutional Neural Net

Neural Net

Perceptron

SVM

DEEP

SHALLOW

Deep (sparse/denoising) Autoencoder AutoencoderNeural Net

Sparse Coding

SP

GMM

Deep Belief Net

Restricted BM

BayesNP

UNSUPERVISED

Slide: M. Ranzato

# Deep Learning Ingredients



- Deep Learning is based on
  - ▶ Availability of large datasets
  - ▶ Massive parallel compute power
  - ▶ Advances in machine learning over the years

- Strong improvements due to
  - ▶ Internet (availability of large-scale data)
  - ▶ GPUs (availability of parallel compute power)
  - ▶ Deep / hierarchical models with end-to-end learning

# Image Features vs. Deep Learning



Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012. Figure copyright Krizhevsky, Sutskever, and Hinton, 2012. Reproduced with permission.

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Traditional Approach



$x$ → feature extraction (hand crafted) → Classification → **Car** $y$

- Feature extraction
  - ▸ often hand crafted and fixed
  - ▸ might be too general (not task-specific enough)
  - ▸ might be too specific (does not generalize to other tasks)
- How to achieve best classification performance
  - ▸ more complex classifier (e.g. multi-feature, non-linear)?
  - ▸ how specialized for the task?

# Features Alone Can Explain (Almost) 10 Years of Progress

[Benenson,Omran,Hosang,Schiele@ECCV workshop'14]



Legend:
- 94.73% VJ (I)
- 89.06% VJLike (C)
- 68.46% HOG (I)
- 63.59% HOGLike–L1 (C)
- 51.46% HOGLike–L2 (C)
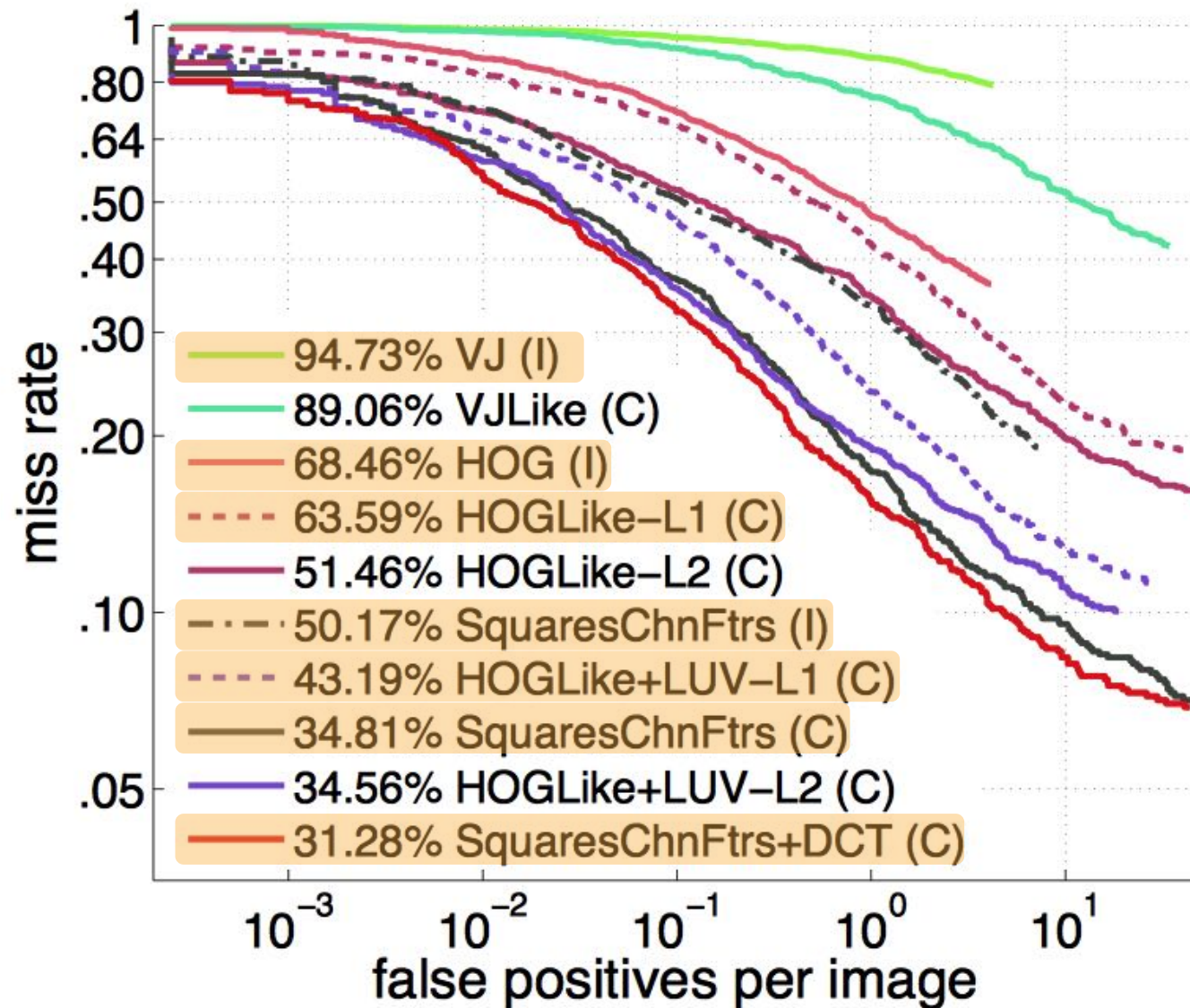- 50.17% SquaresChnFtrs (I)
- 43.19% HOGLike+LUV–L1 (C)
- 34.81% SquaresChnFtrs (C)
- 34.56% HOGLike+LUV–L2 (C)
- 31.28% SquaresChnFtrs+DCT (C)

Axes: miss rate (y) vs false positives per image (x)

# Motivation

- Features are key to recent progress in recognition
- Multitude of hand-designed features currently in use
  - SIFT, HOG, LBP, MSER, Color-SIFT………..
- Where next? Better classifiers? Or keep building more features?



Felzenszwalb, Girshick,
McAllester and Ramanan, PAMI 2007



- Low level features: SIFT and its variants, LBP, HOG.
- Dense sampling and interest point detector;
- Represented as Bags of Words;

Yan & Huang
(Winner of PASCAL 2010 classification competition)

slide credit: Rob Fergus

# Hand-Crafted Features

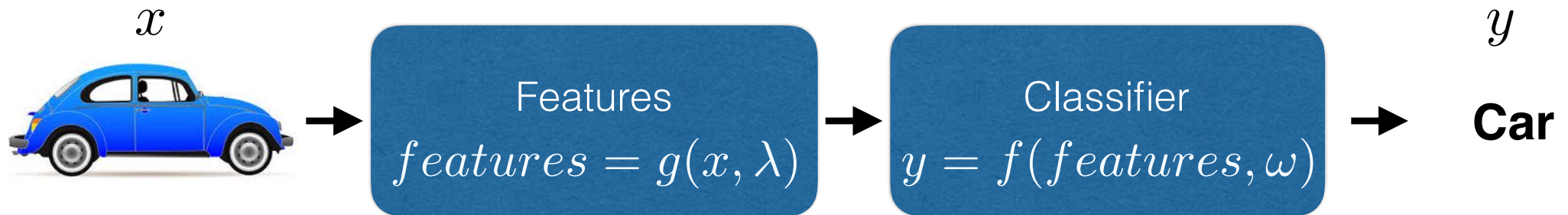- LP-β  Multiple Kernel Learning (MKL)
  - Gehler and Nowozin, On Feature Combination for Multiclass Object Classification, ICCV'09

- 39 different kernels
  - PHOG, SIFT, V1S+, Region Cov. Etc.

- MKL only gets few % gain over averaging features

→    Features are doing the work



Caltech–256 (39 kernels)

Legend:
- best feature
- product
- average
- MKL
- LP–β
- LP–B
- Griffin, Holub and Perona (TR06)
- Pinto, Cox and DiCarlo (PLOS08)

Axes: accuracy (y) vs #training examples (x)

# Deep Learning: Trainable features

$x$



Features
$$features = g(x, \lambda)$$

Classifier
$$y = f(features, \omega)$$

$y$

**Car**

- Parameterized feature extraction

- Features should be

  ▸ efficient to compute

  ▸ efficient to train (differentiable)

# Deep Learning: Joint Training of all Parameters

$x$



Features

$$features = g(x, \lambda)$$

Classifier

$$y = f(features, \omega)$$

$y$

**Car**

**"End-to-End" System**

- Parameterized feature extraction

- Features should be

  ▸ efficient to compute

  ▸ efficient to train (differentiable)

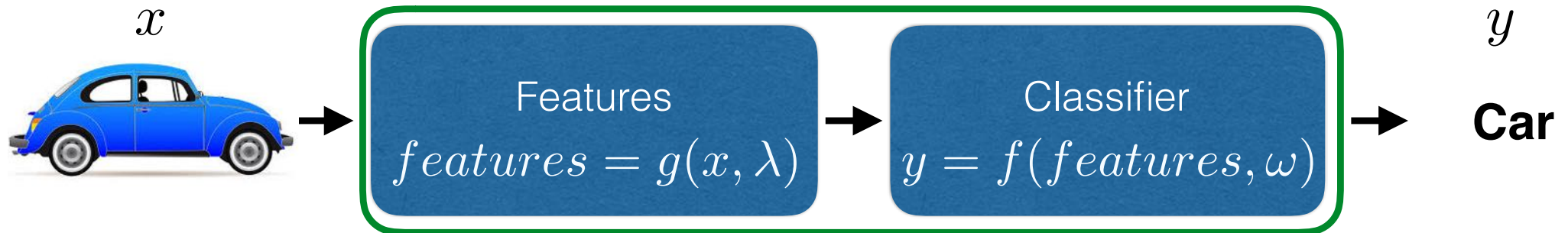- **Joint** training of **feature** extraction and **classification**

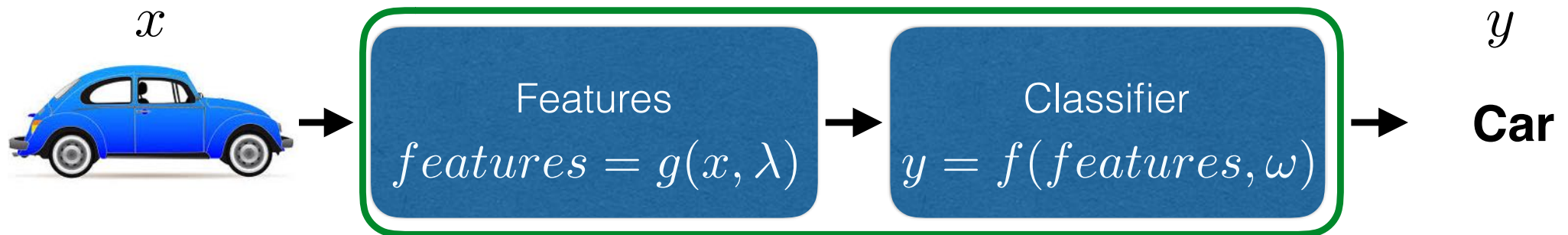- Feature extraction and classification merge into one pipeline

# Deep Learning: Joint Training of all Parameters

$x$



Features

$$features = g(x, \lambda)$$

Classifier

$$y = f(features, \omega)$$

$y$

**Car**

**"End-to-End" System**

- All parts are adaptive
- No differentiation between feature extraction and classification
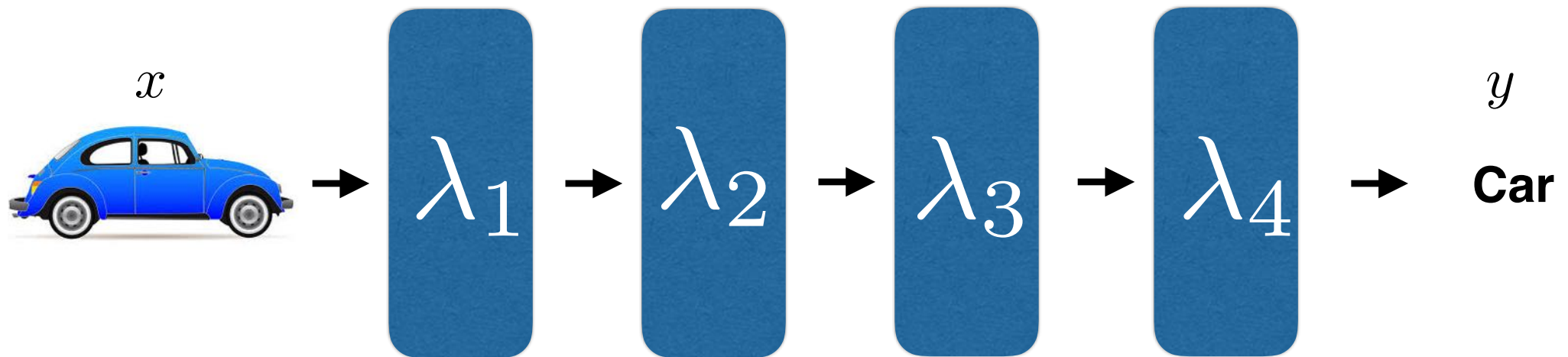- Non linear transformation from input to desired output

# Deep Learning: Complex Functions by Composition

$x$



$y$

**Car**

Features

$$features = g(x, \lambda)$$

Classifier

$$y = f(features, \omega)$$

**"End-to-End" System**

- How can we build such systems?

- What is the parameterization (hypothesis)?

- Composition of simple building blocks can lead to complex systems (e.g. neurons - brain)

# Deep Learning: Complex Functions by Composition



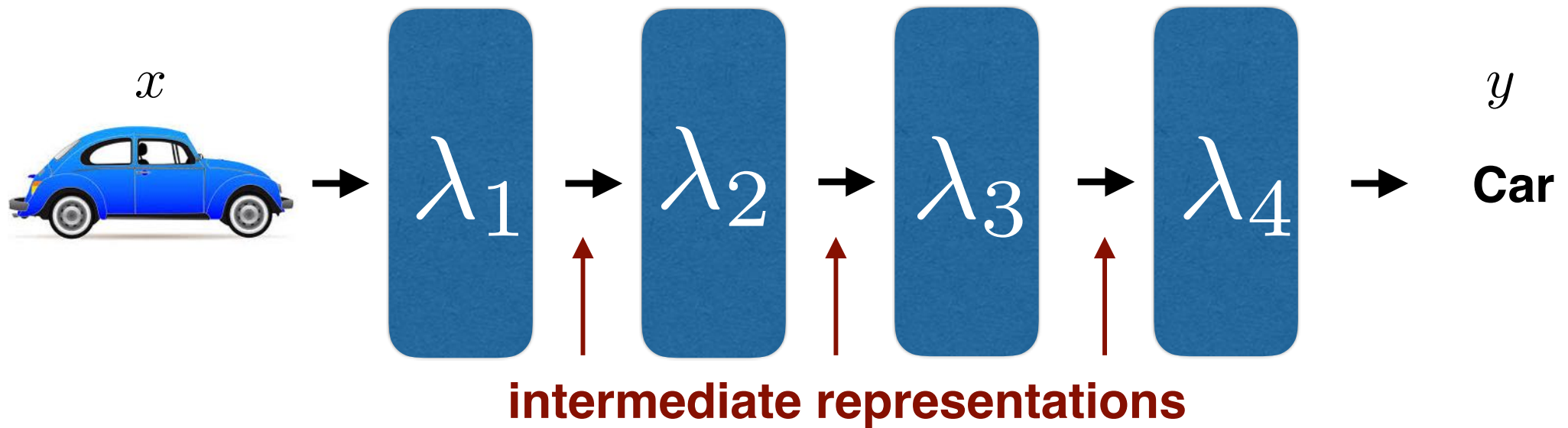$x$ → $\lambda_1$ → $\lambda_2$ → $\lambda_3$ → $\lambda_4$ → $y$ Car

- How can we build such systems?

- What is the parameterization (hypothesis)?

- Composition of simple building blocks can lead to complex systems (e.g. neurons - brain) each block has trainable parameters

- Each block has trainable parameters $\lambda_i$

# Deep Learning: Complex Functions by Composition



$x$ → $\lambda_1$ → $\lambda_2$ → $\lambda_3$ → $\lambda_4$ → $y$ **Car**
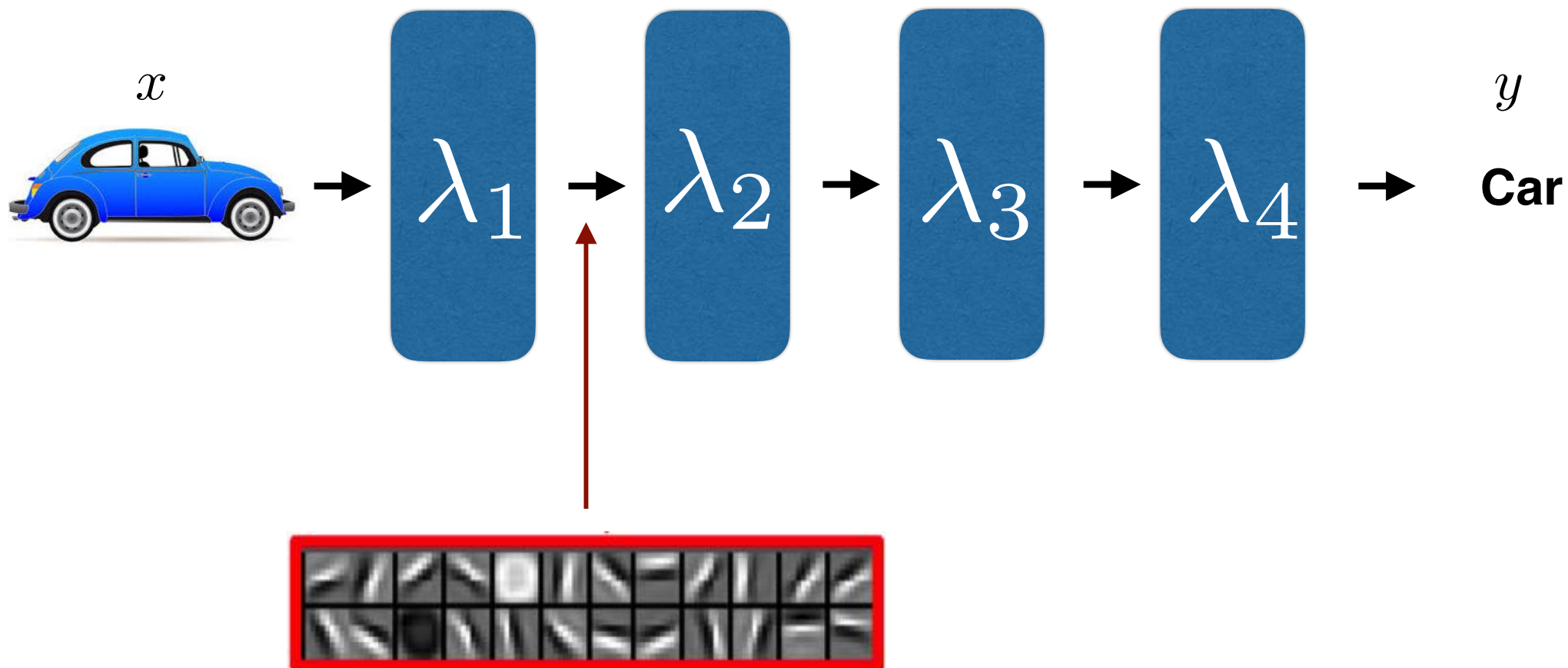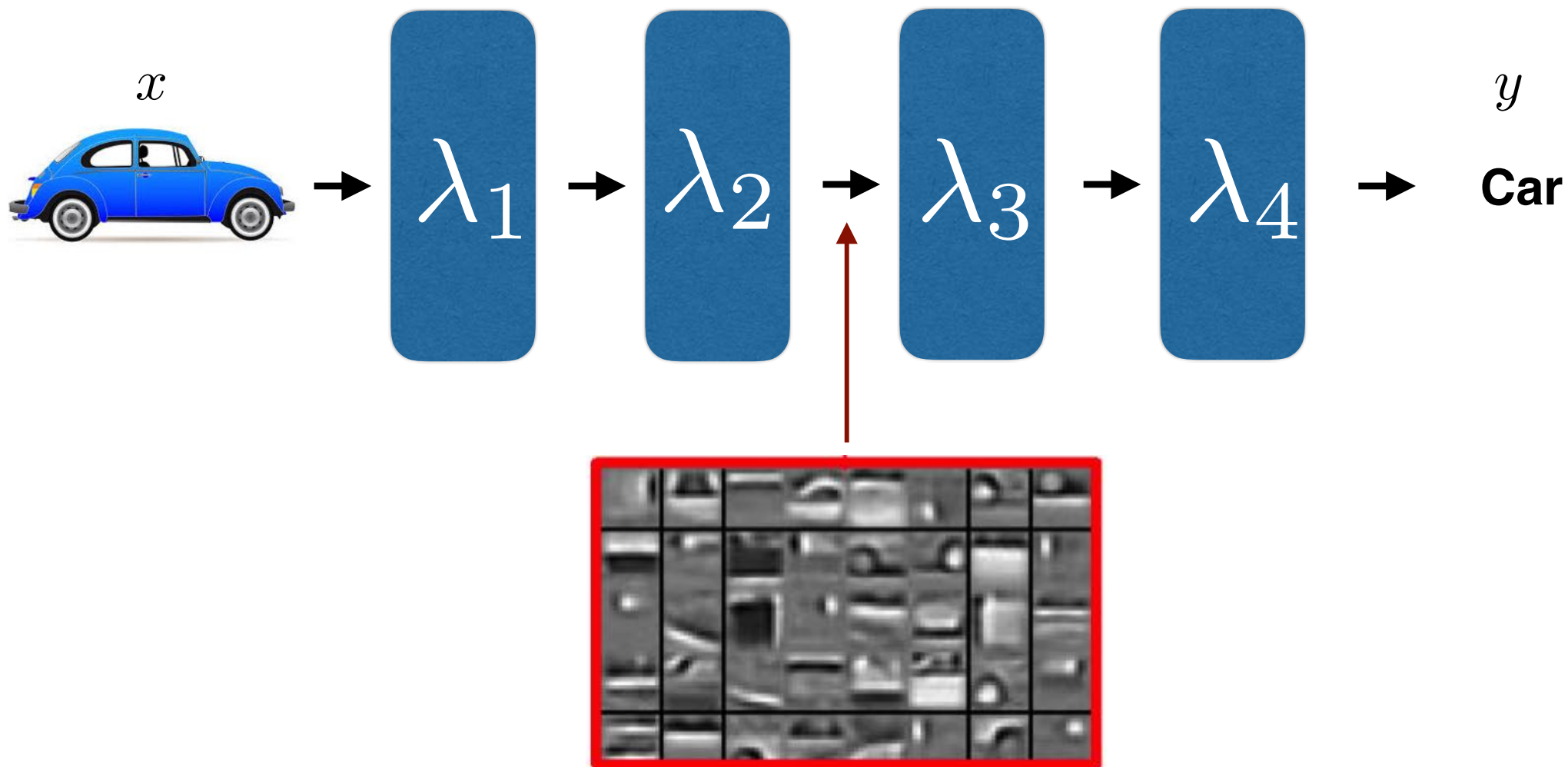
**intermediate representations**

- How can we build such systems?

- What is the parameterization (hypothesis)?

- Composition of simple building blocks can lead to complex systems (e.g. neurons - brain) each block has trainable parameters

- Each block has trainable parameters $\lambda_i$

# Deep Learning: Complex Functions by Composition



$$x \rightarrow \lambda_1 \rightarrow \lambda_2 \rightarrow \lambda_3 \rightarrow \lambda_4 \rightarrow \text{Car} \quad (y)$$
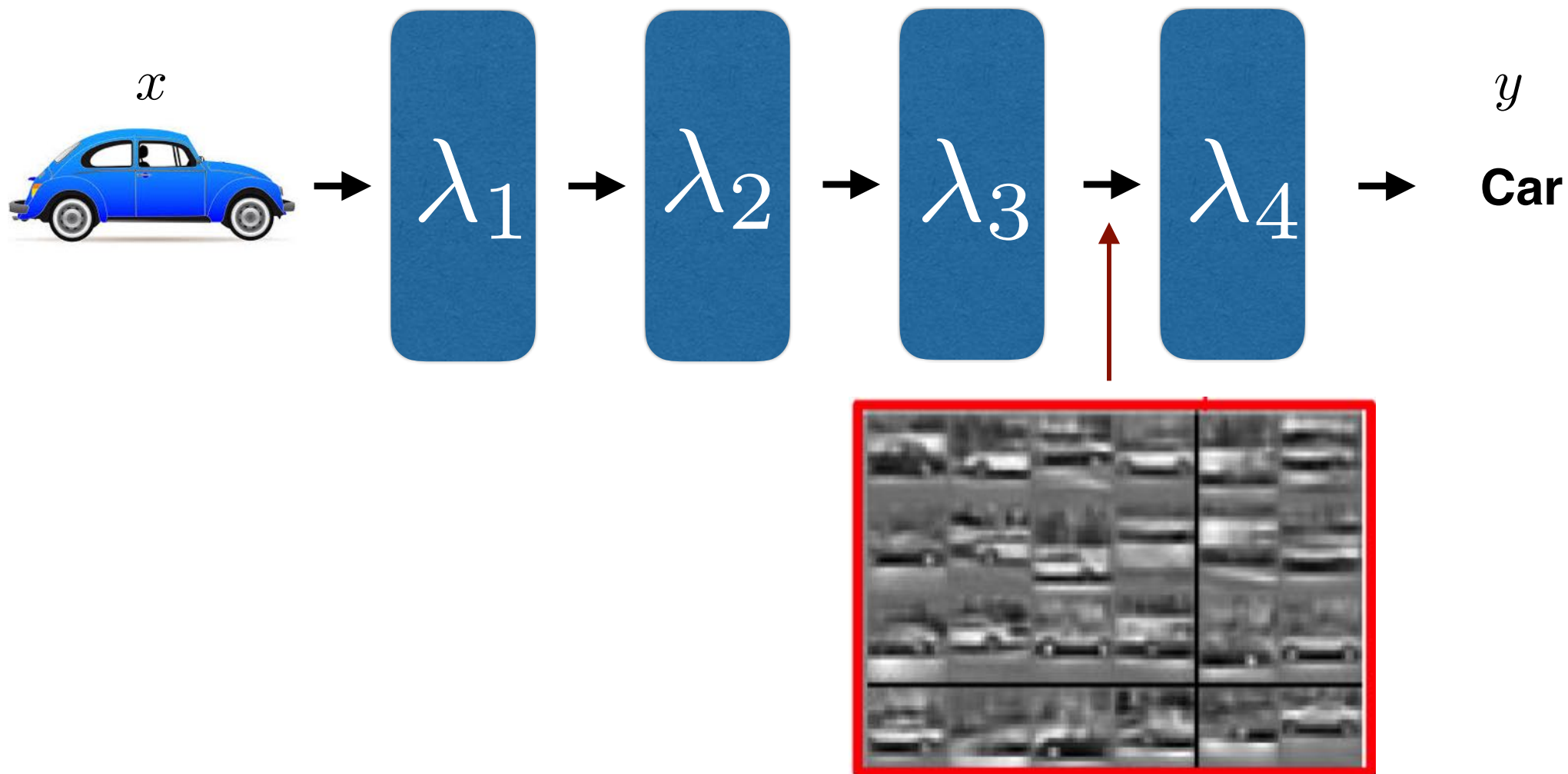
Lee et al. "Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations"

# Deep Learning: Complex Functions by Composition



$x$       $\lambda_1 \rightarrow \lambda_2 \rightarrow \lambda_3 \rightarrow \lambda_4 \rightarrow$    $y$   **Car**
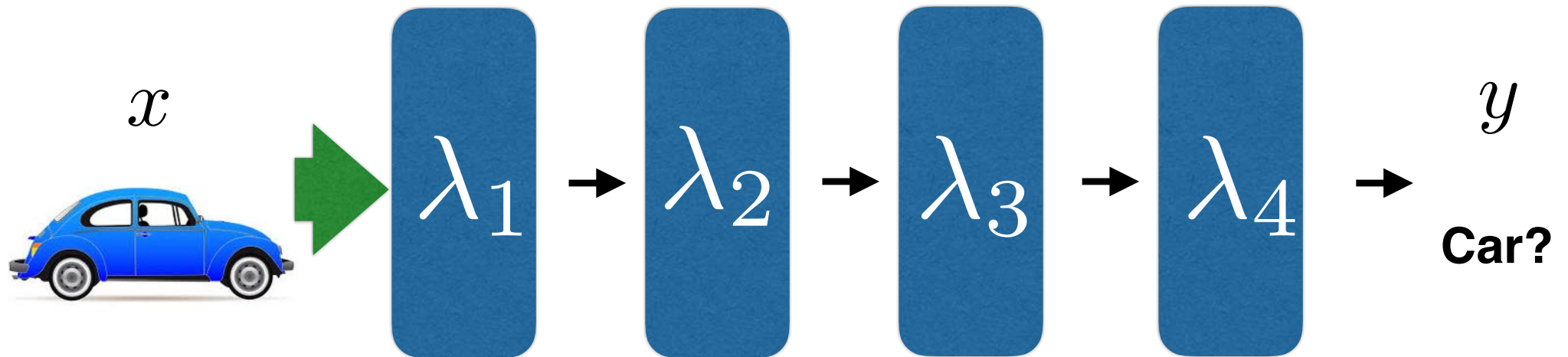
Lee et al. "Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations"

# Deep Learning: Complex Functions by Composition



$$x \rightarrow \lambda_1 \rightarrow \lambda_2 \rightarrow \lambda_3 \rightarrow \lambda_4 \rightarrow \text{Car} \quad (y)$$
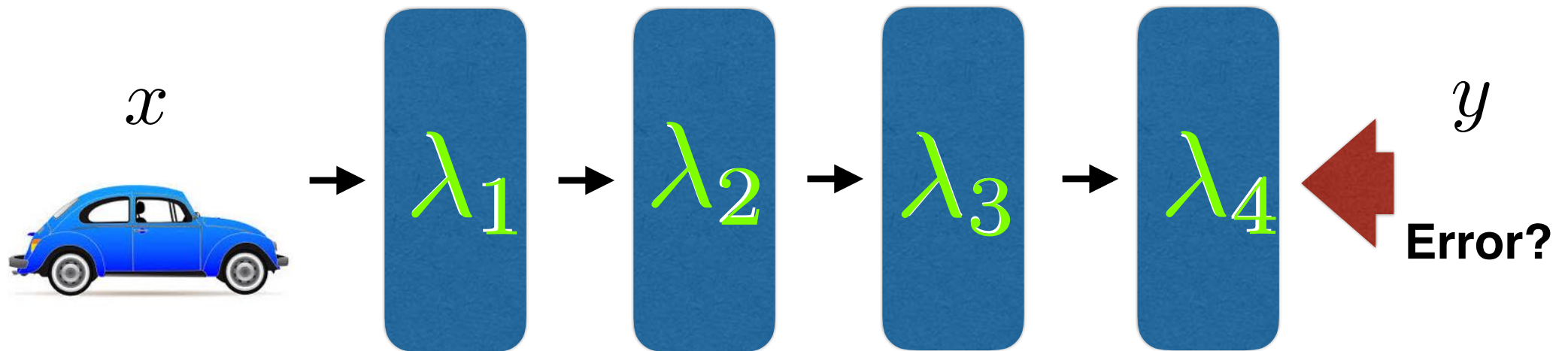
Lee et al. "Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations"

# Training: Overview



- Setting
  - generate output y for input x (forward pass)

# Training: Overview



- Setting
  - generate output y for input x (forward pass)
  - when there is an error, propagate error backwards to update weights (error back propagation)
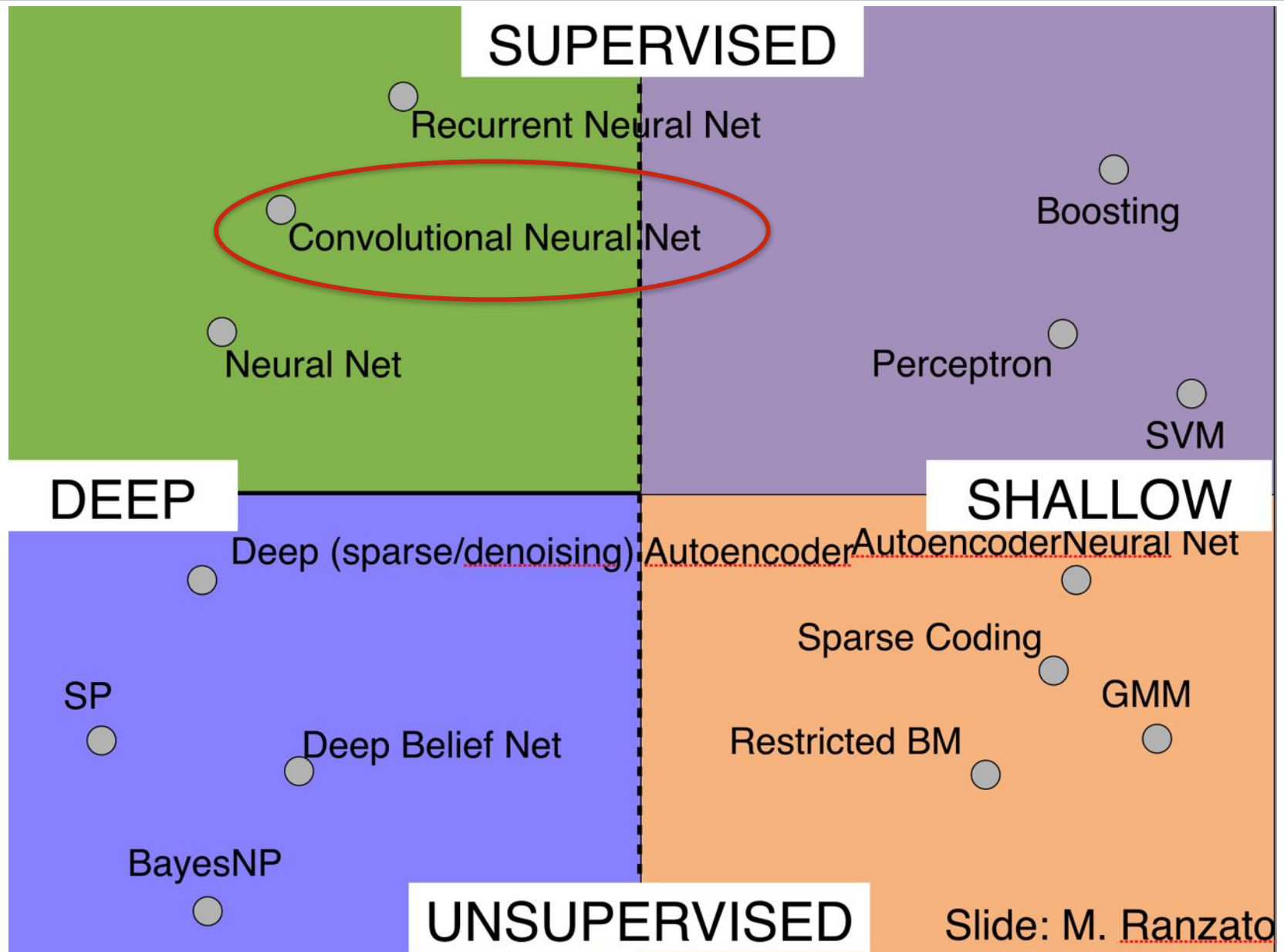
# Summary of Main Ideas in Deep Learning

1. **Learning of feature extraction** (across many layers)

2. **Efficient** and **trainable** systems by **differentiable** building blocks

3. Composition of deep architectures via **non-linear modules**

4. "**End-to-End**" training: no differentiation between feature extraction and classification

# Overview Today's Lecture

- Backpropagation - Gradient Descent

  ‣ illustrated using computational graphs

  ‣ chain rule - upstream and local gradients

  ‣ modularization simple

- What is Deep Learning

  ‣ intuition why deep learning can help

  ‣ integrated learning of features and classifier

- **Convolutional Neural Networks (CNNs)**
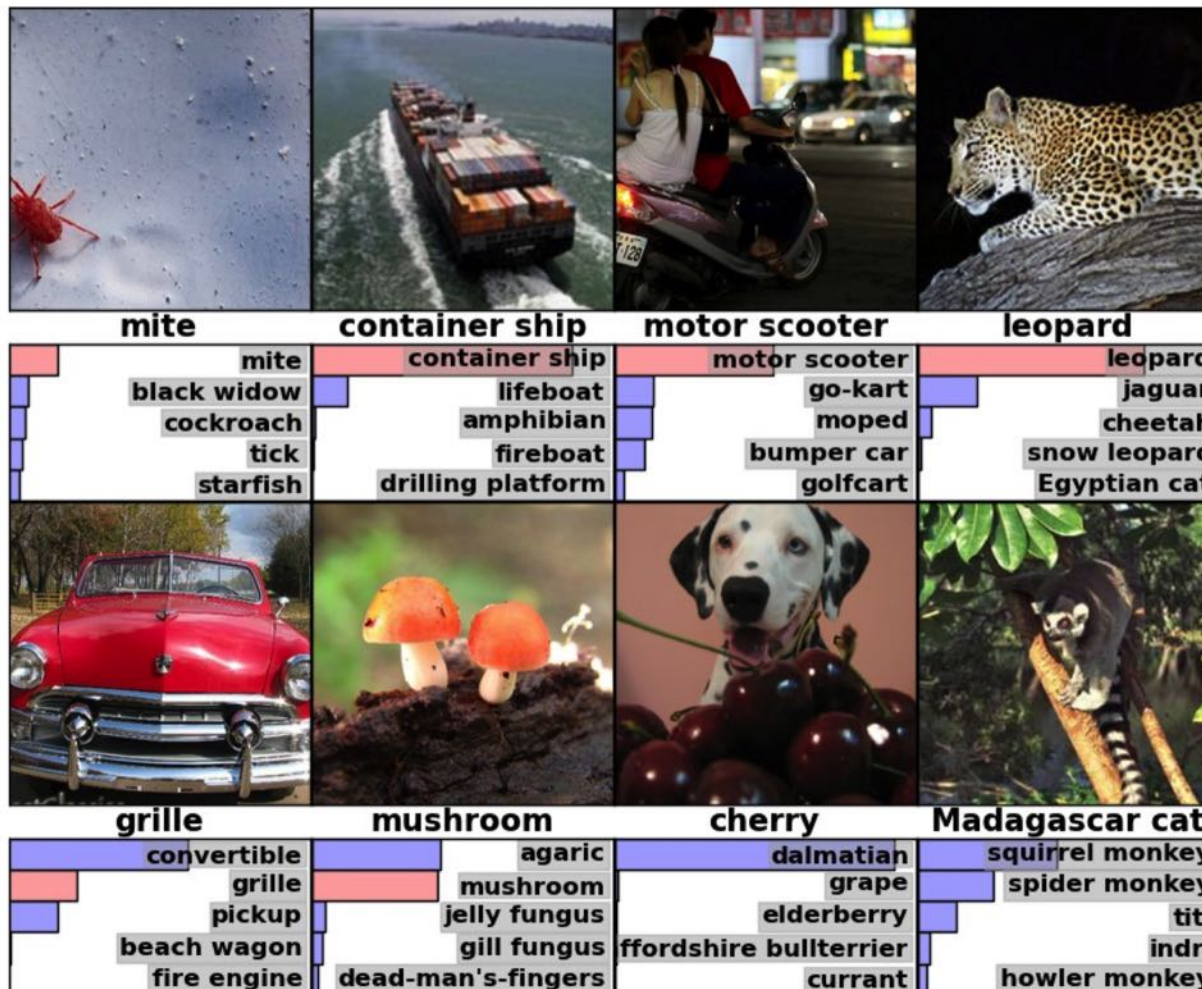
  ‣ one of the (few) highly successful NNs
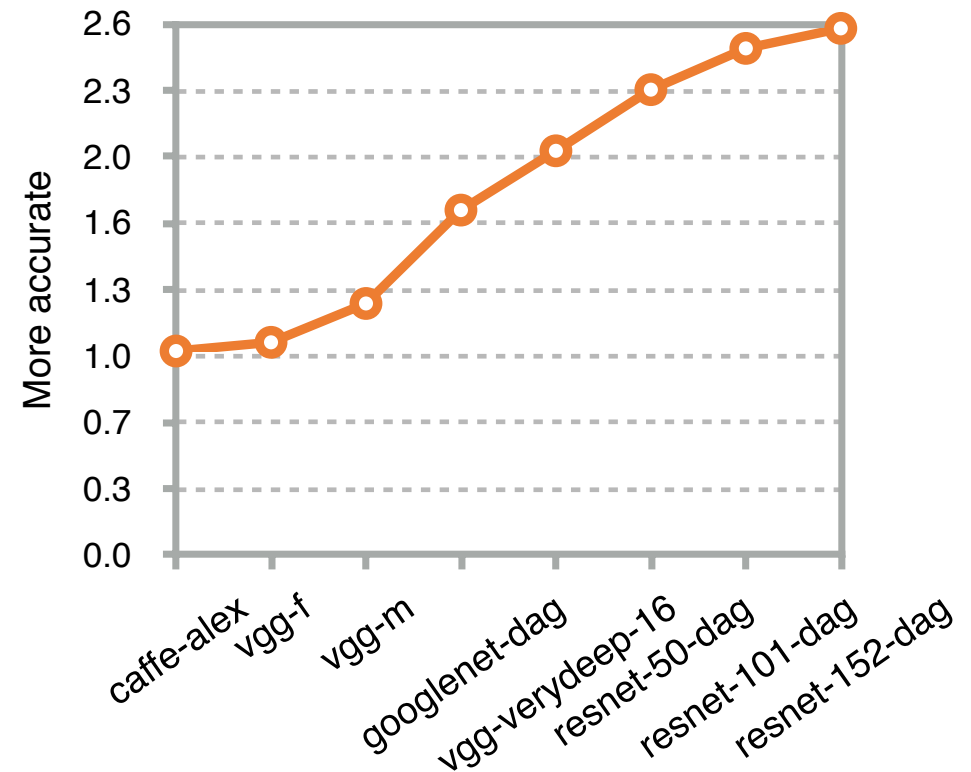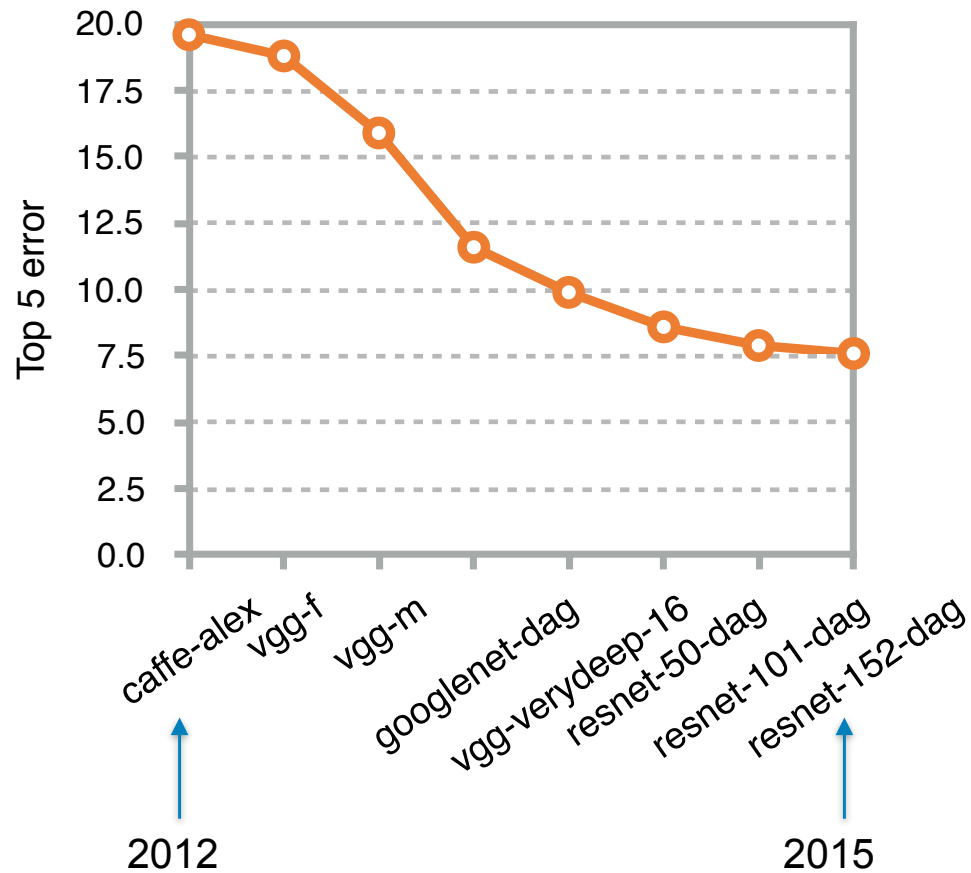
# Overview of Deep Learning



Slide: M. Ranzato

# Computer Vision: Image Classification

1000 classes; 1ms for an image;  92.5% recognition rate

http://demo.caffe.berkeleyvision.org

# Architectures get deeper



~ 2.6x improvement in 3 years

# Architectures get deeper

GoogLeNet (2014) ──────────────────────────── ResNet 50 (2015)

VGG-VD-16 (2014) ──────────────────────────── ResNet 152 (2015)

VGG-M (2013) ────────────

AlexNet (2012) ────────────

16 convolutional layers ──────────▶

50 convolutional layers ──────────▶

152 convolutional layers ──────────▶

Krizhevsky, I. Sutskever, and G. E. Hinton. *ImageNet classification with deep convolutional neural networks*. In Proc. NIPS, 2012.

C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. *Going deeper with convolutions*. In Proc. CVPR, 2015.

K. Simonyan and A. Zisserman. *Very deep convolutional networks for large-scale image recognition*. In Proc. ICLR, 2015.

K. He, X. Zhang, S. Ren, and J. Sun. *Deep residual learning for image recognition*. In Proc. CVPR, 2016.
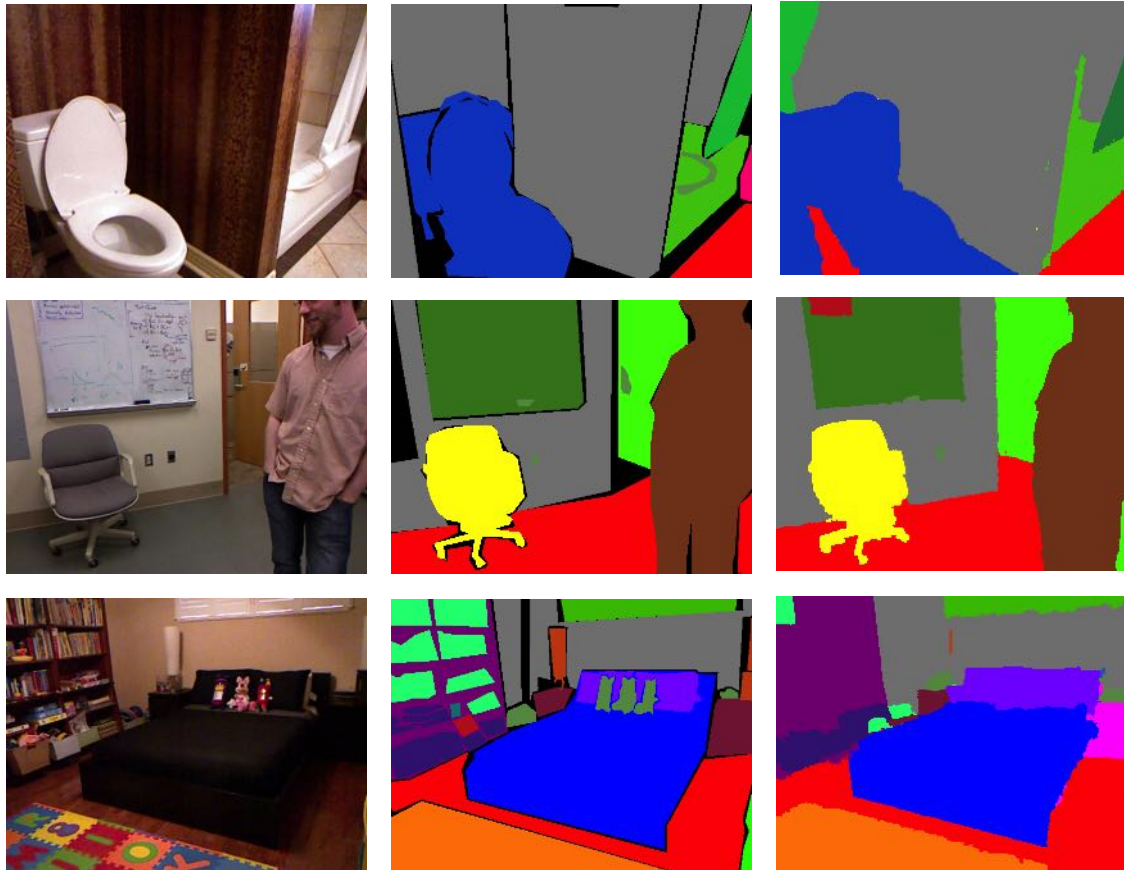
# Computer Vision: Semantic Segmentation



Badrinarayanan et al. 2015
SegNet: A Deep Convolutional Encoder-Decoder Architecture for Robust Semantic Pixel-Wise Labelling

# Computer Vision: Semantic Segmentation



Input          Groundtruth          Output

Yang He; Wei-Chen Chiu; Margret Keuper; Mario Fritz
**STD2P: RGBD Semantic Segmentation Using Spatio-Temporal Data-Driven Pooling**
IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017
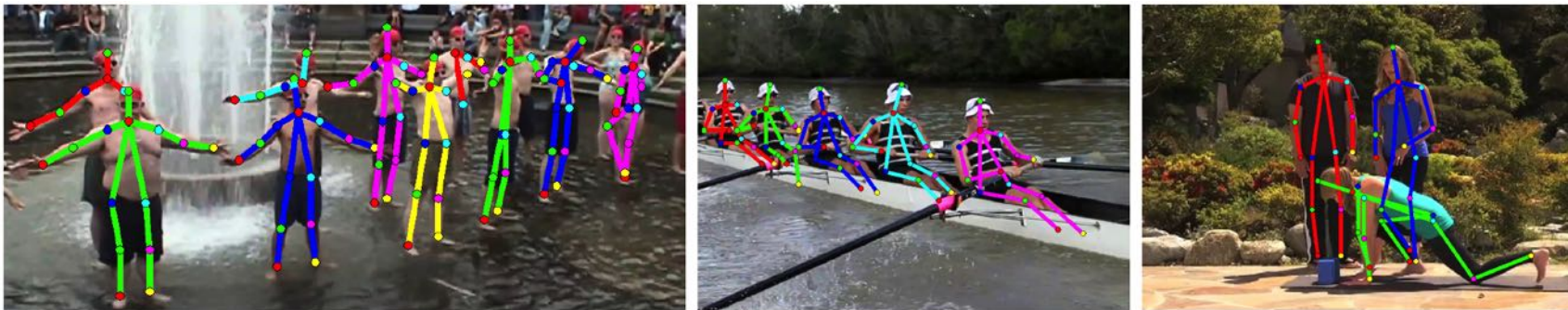
# Deep(er)Cut: Joint Subset Partition & Labeling for Multi Person Pose Estimation

**DeepCut: Joint Subset Partition and Labeling for Multi-Person Pose Estimation**
L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. Gehler, and B. Schiele, CVPR'16
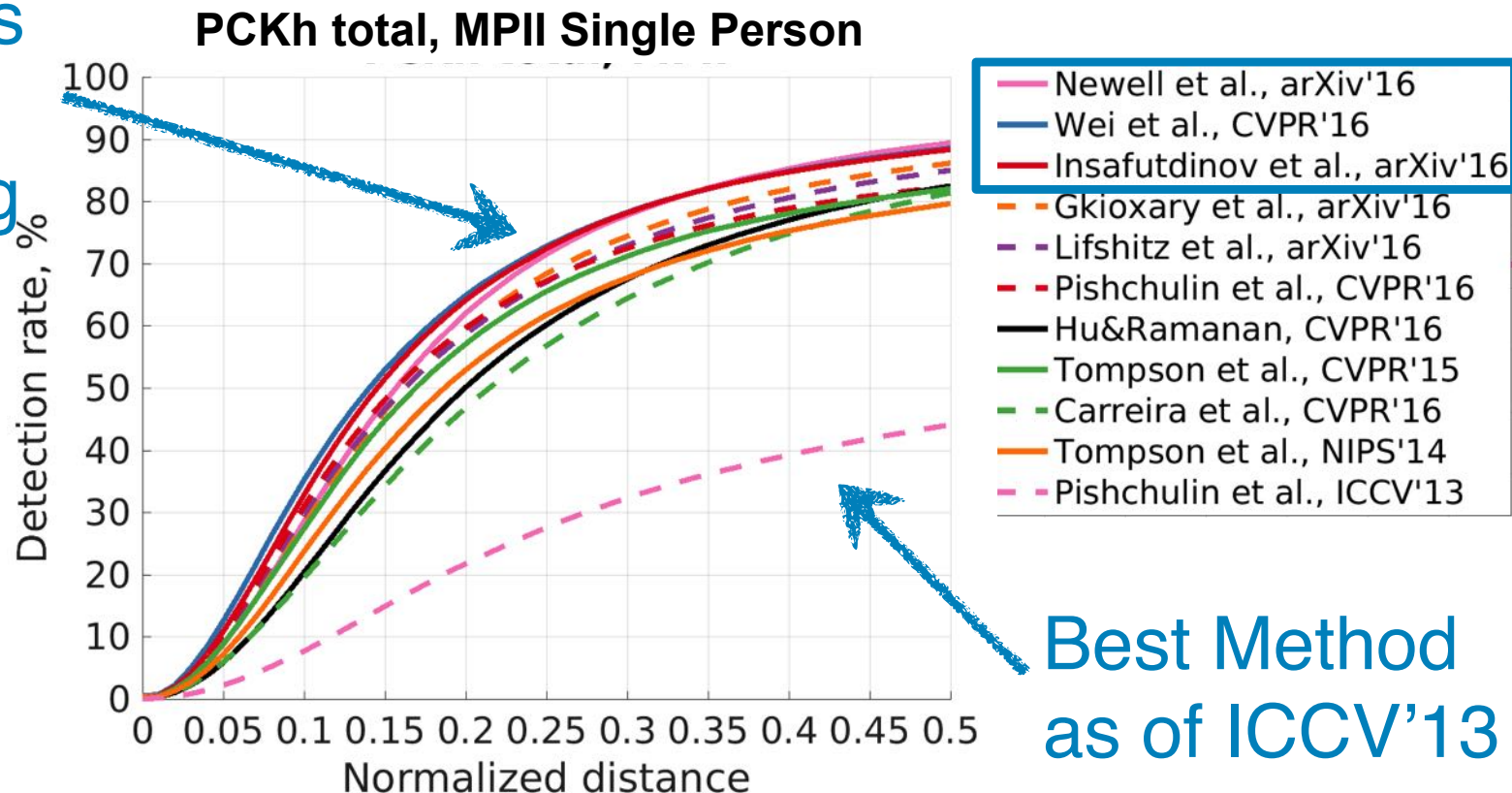
**DeeperCut: A Deeper, Stronger, and Faster Multi-Person Pose Estimation Model**
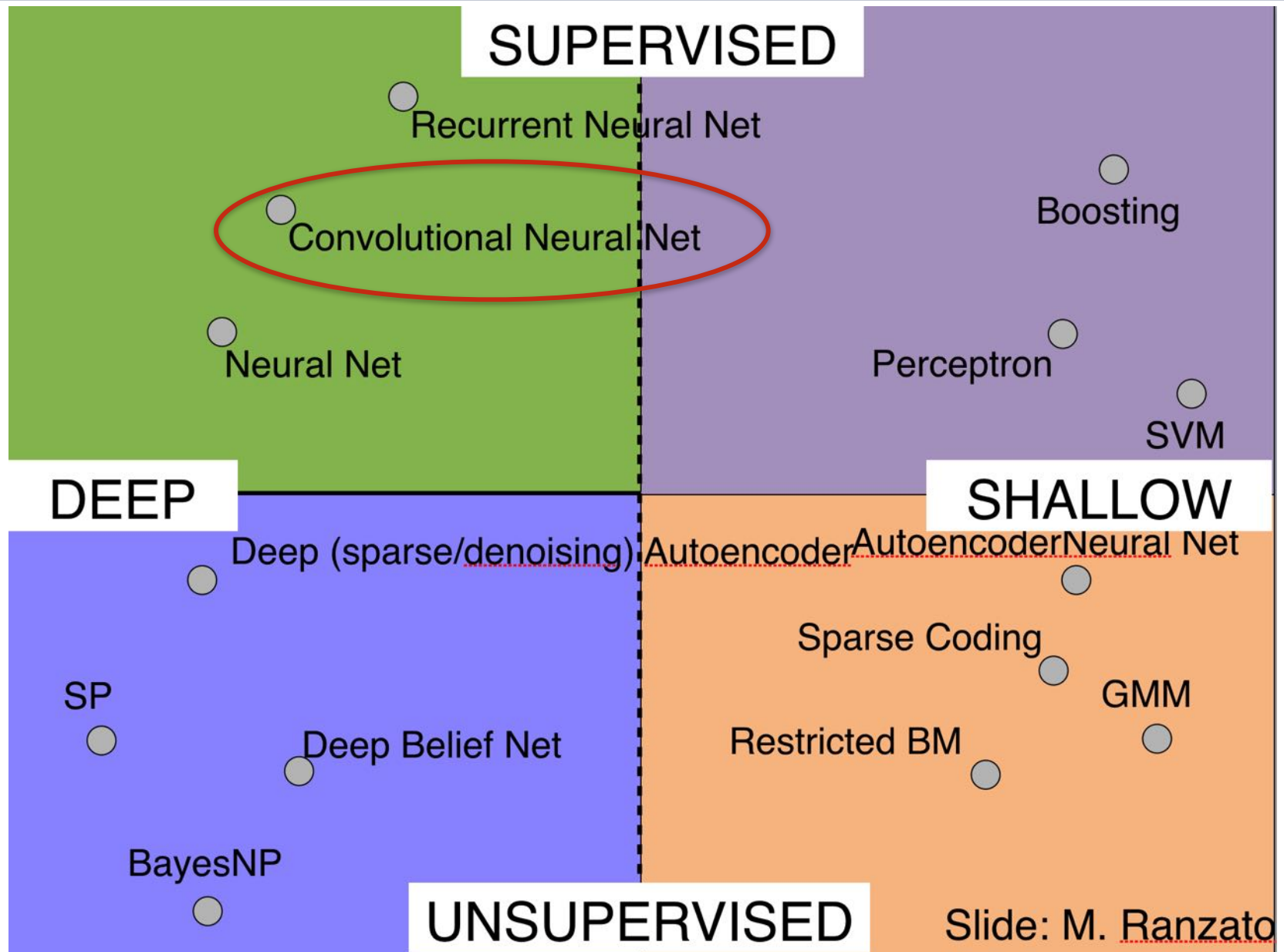E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, ECCV'16

# Analysis - overall performance

**Best Methods now: deep learning "takes" over**



PCKh total, MPII Single Person

Legend:
- Newell et al., arXiv'16
- Wei et al., CVPR'16
- Insafutdinov et al., arXiv'16
- Gkioxary et al., arXiv'16
- Lifshitz et al., arXiv'16
- Pishchulin et al., CVPR'16
- Hu&Ramanan, CVPR'16
- Tompson et al., CVPR'15
- Carreira et al., CVPR'16
- Tompson et al., NIPS'14
- Pishchulin et al., ICCV'13

Y-axis: Detection rate, %
X-axis: Normalized distance

**Best Method as of ICCV'13**

✓ since CVPR'14, dataset has become **de-facto standard benchmark**

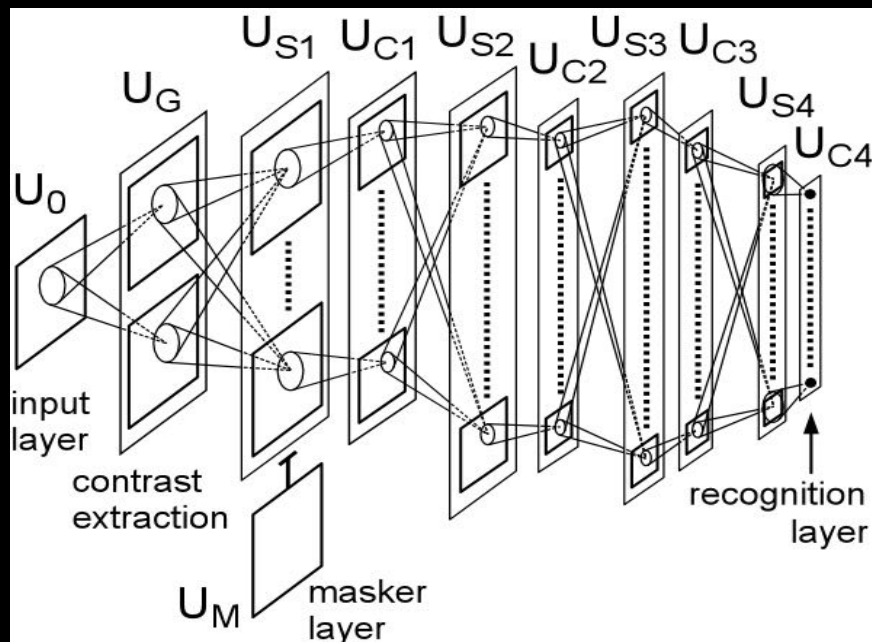✓ **large training set** facilitated development of **deep learning methods**

# Overview of Deep Learning



SUPERVISED

Recurrent Neural Net

Convolutional Neural Net

Boosting

Neural Net

Perceptron

SVM

DEEP

SHALLOW

Deep (sparse/denoising) Autoencoder   AutoencoderNeural Net

Sparse Coding

SP

GMM

Deep Belief Net

Restricted BM

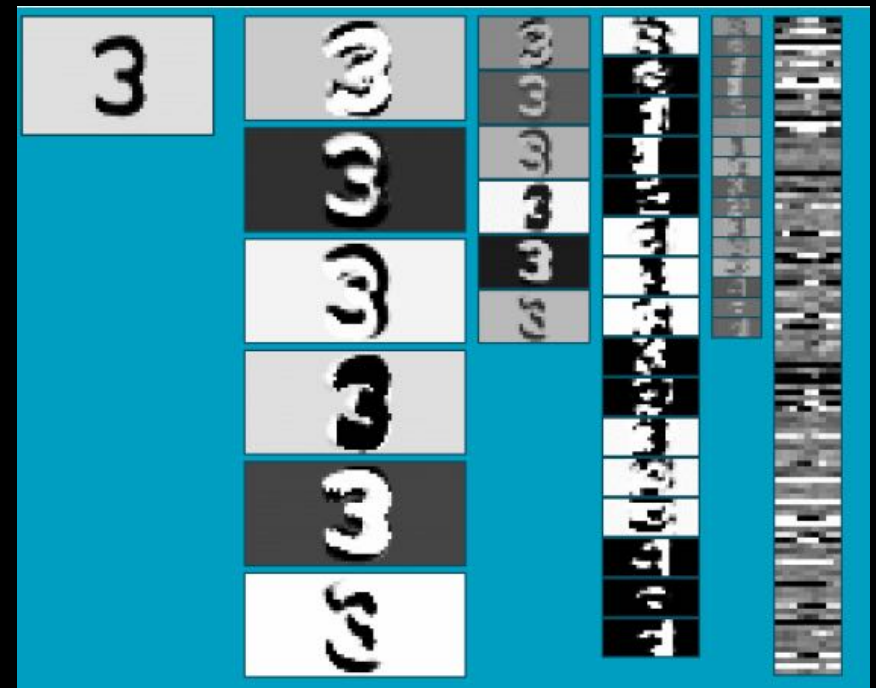BayesNP

UNSUPERVISED

Slide: M. Ranzato

# Multistage Hubel&Wiesel Architecture

- [Hubel & Wiesel 1962]
  - simple cells detect local features
  - complex cells "pool" the outputs of simple cells within a retinotopic neighborhood.





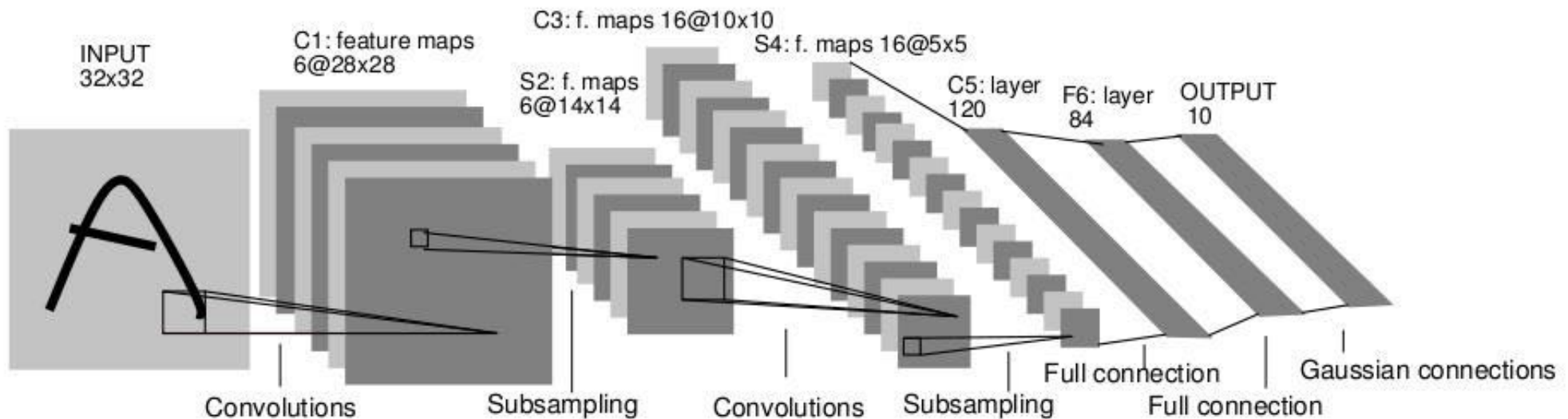Cognitron / Neocognitron
[Fukushima 1971-1982]

- Also HMAX [Poggio 2002-2006]



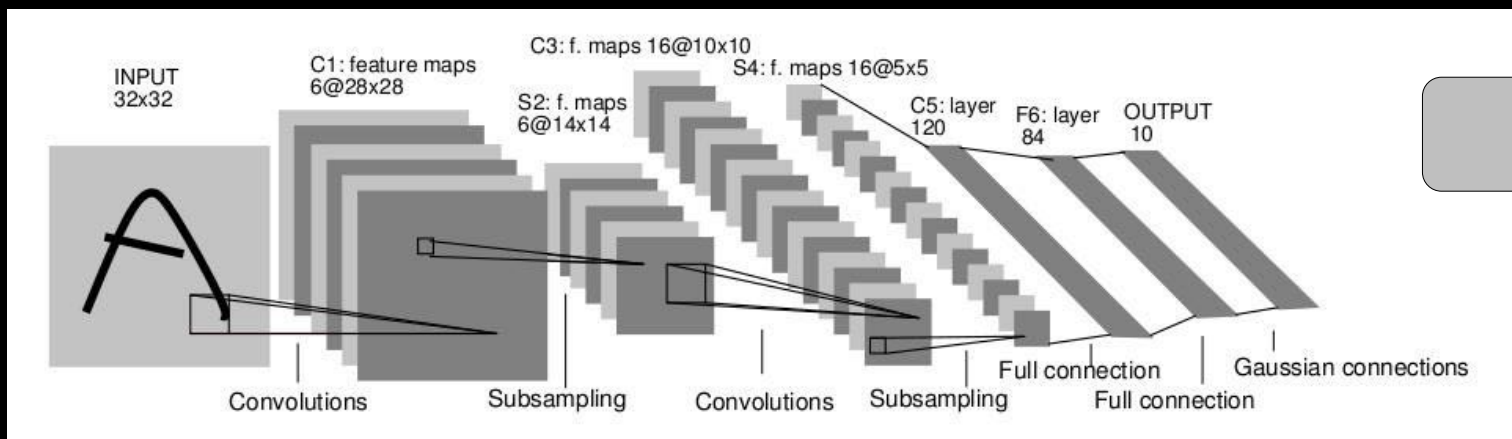Convolutional Networks
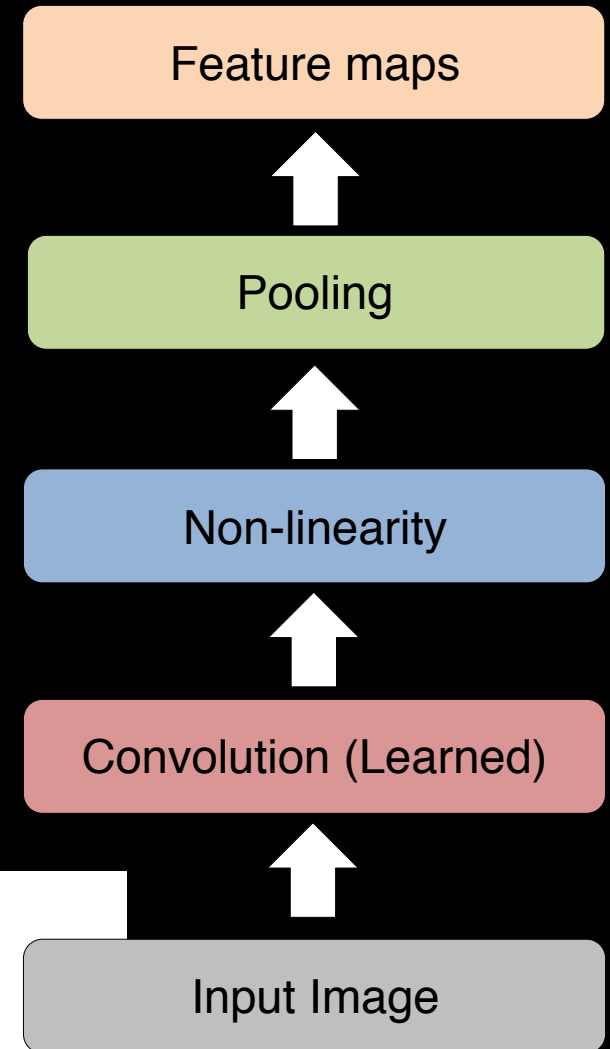[LeCun 1988-present]

# Convolutional Neural Networks

- LeCun et al. 1989

- Neural network with specialized connectivity structure

# Characteristics of Convnets

- Feed-forward:
  - Convolve input
  - Non-linearity (rectified linear)
  - Pooling (local max) / (=subsampling)
- Supervised
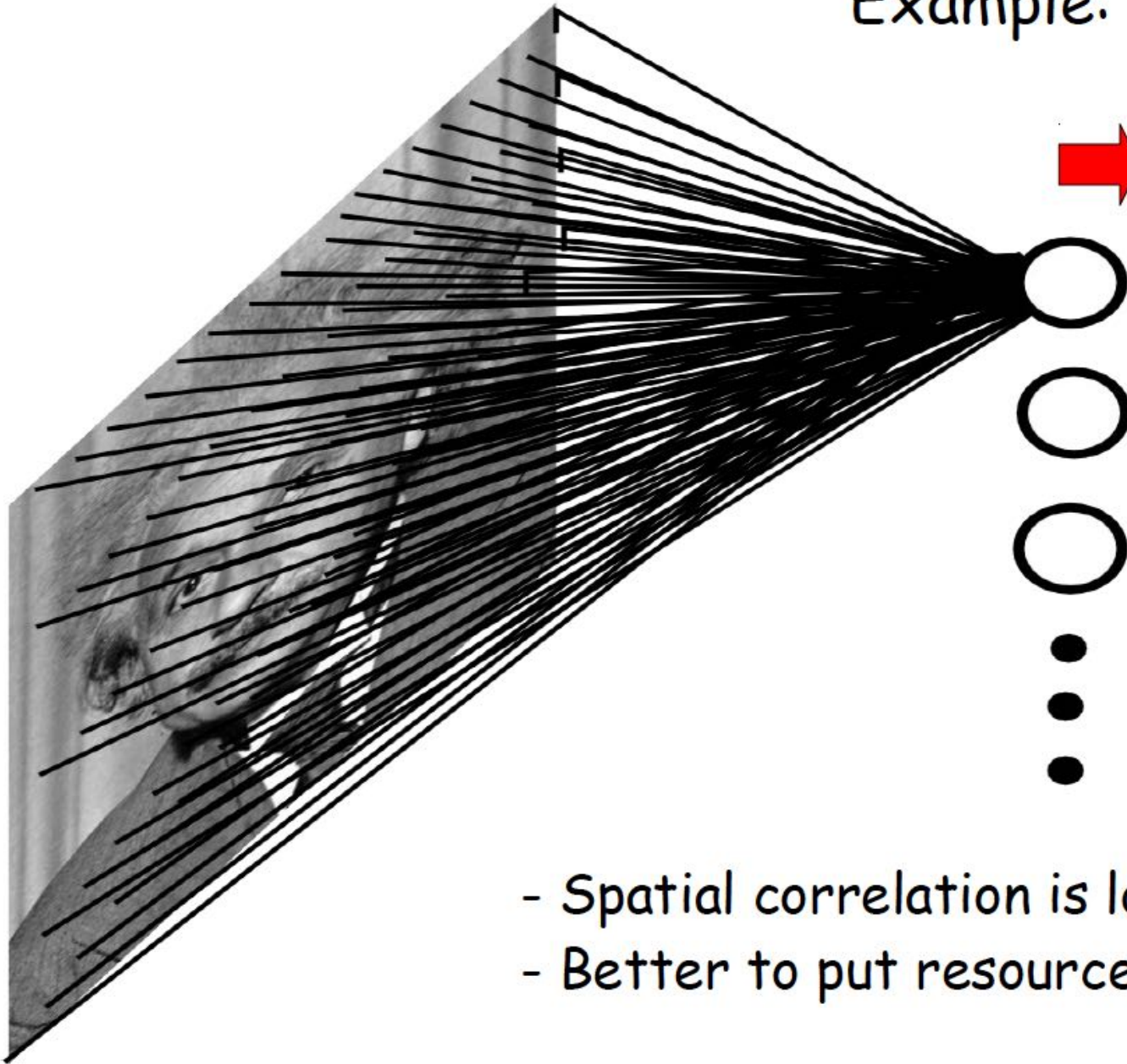- Train convolutional filters by back-propagating classification error



**Feature maps**

↑

**Pooling**

↑

**Non-linearity**

↑

**Convolution (Learned)**

↑

**Input Image**



[LeCun et al. 1989]
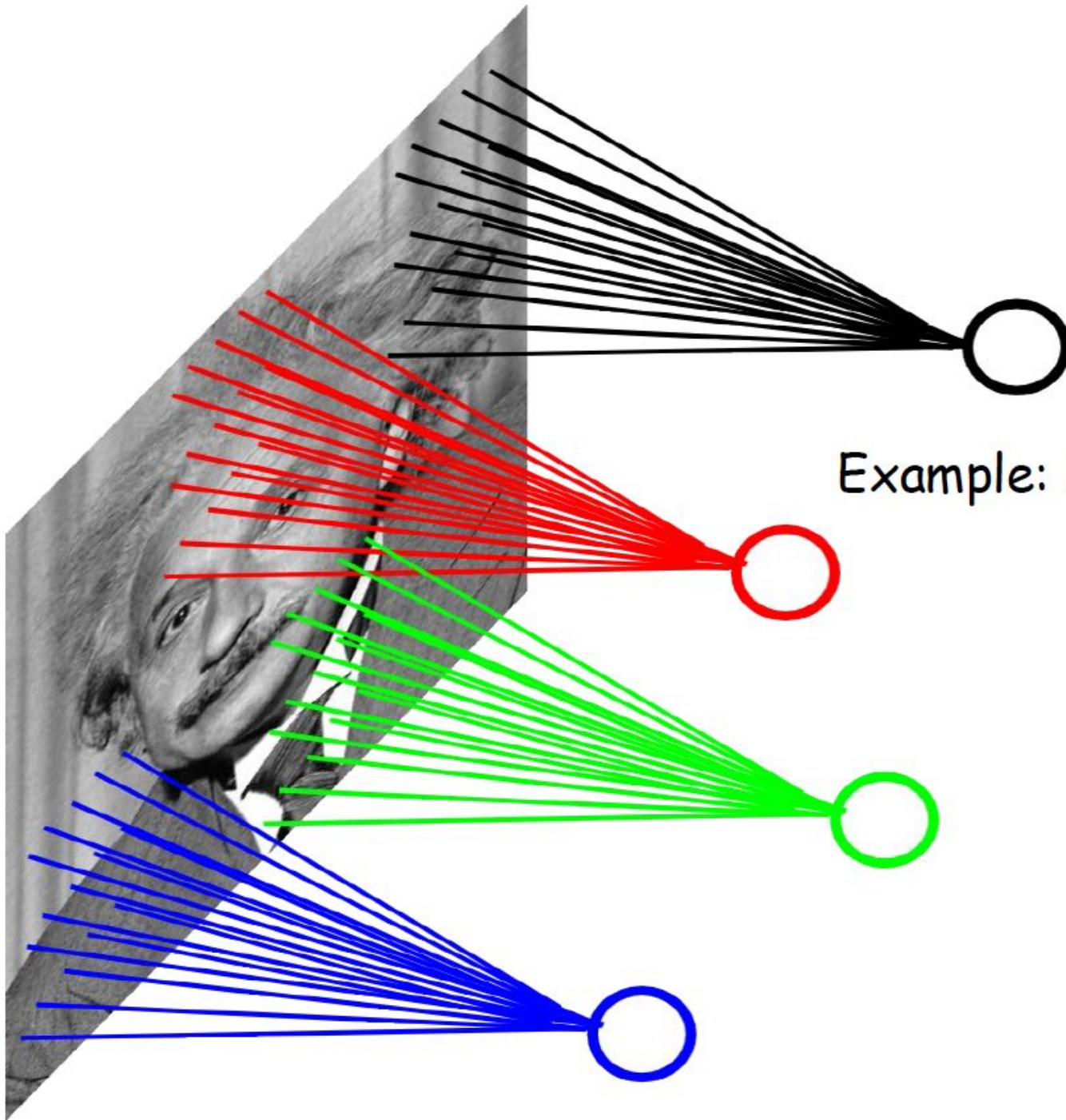
# FULLY CONNECTED NEURAL NET

Example: 1000x1000 image
1M hidden units
➡ 10^12 parameters!!!



- Spatial correlation is local
- Better to put resources elsewhere!

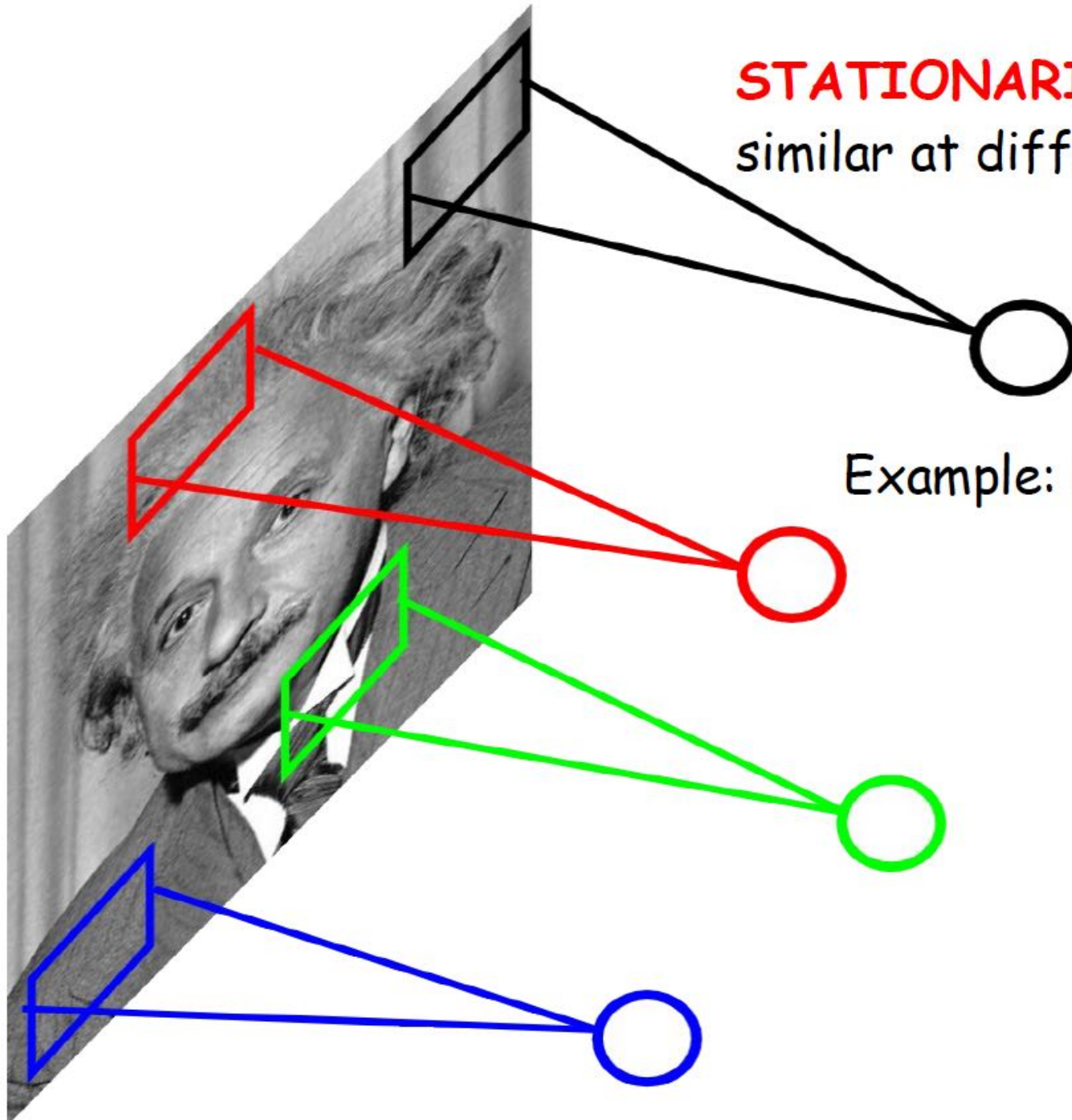Ranzato

# LOCALLY CONNECTED NEURAL NET

Example: 1000x1000 image
   1M hidden units
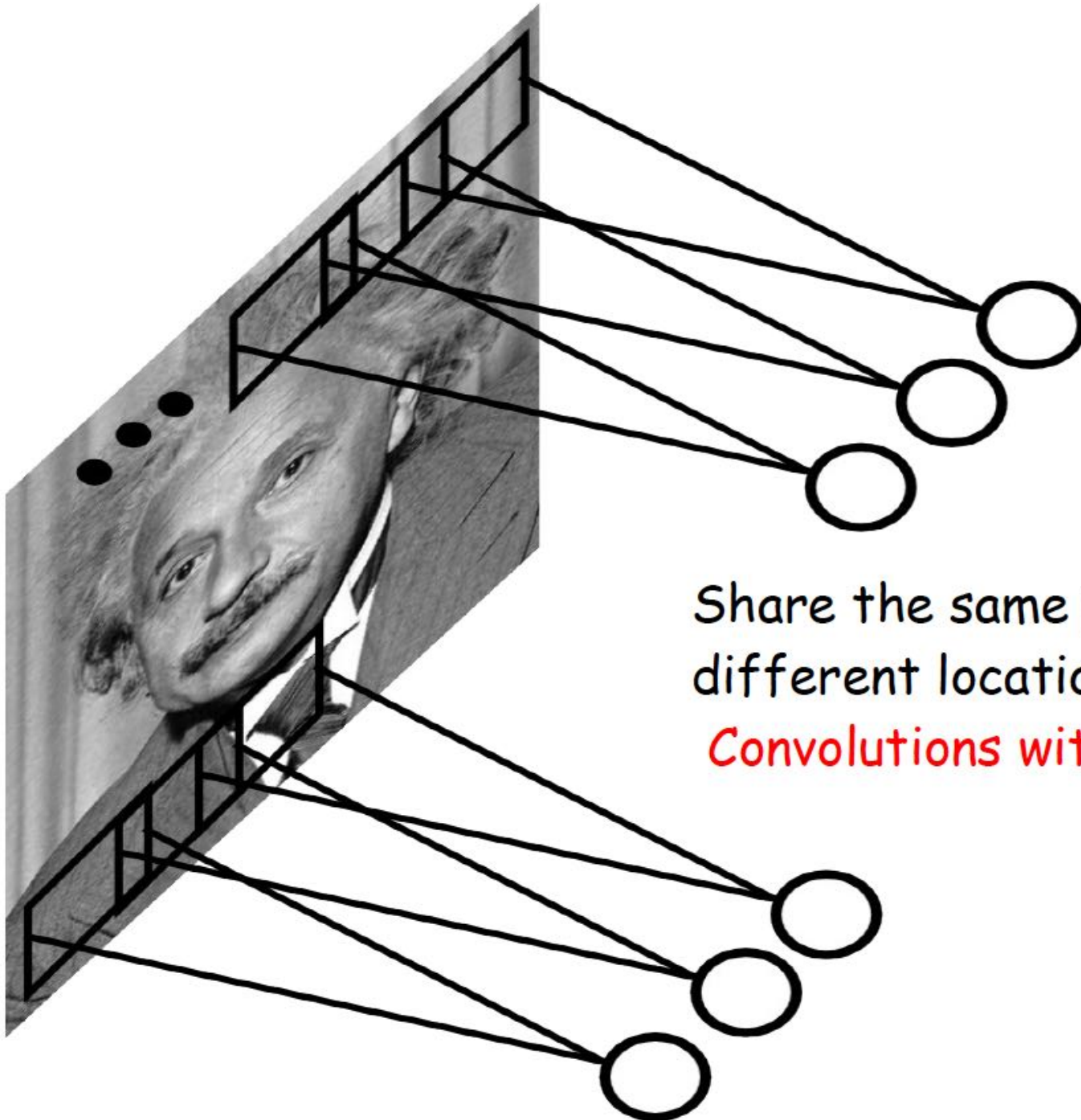   Filter size: 10x10
   100M parameters

Ranzato

# LOCALLY CONNECTED NEURAL NET



**STATIONARITY?** Statistics is similar at different locations

Example: 1000x1000 image
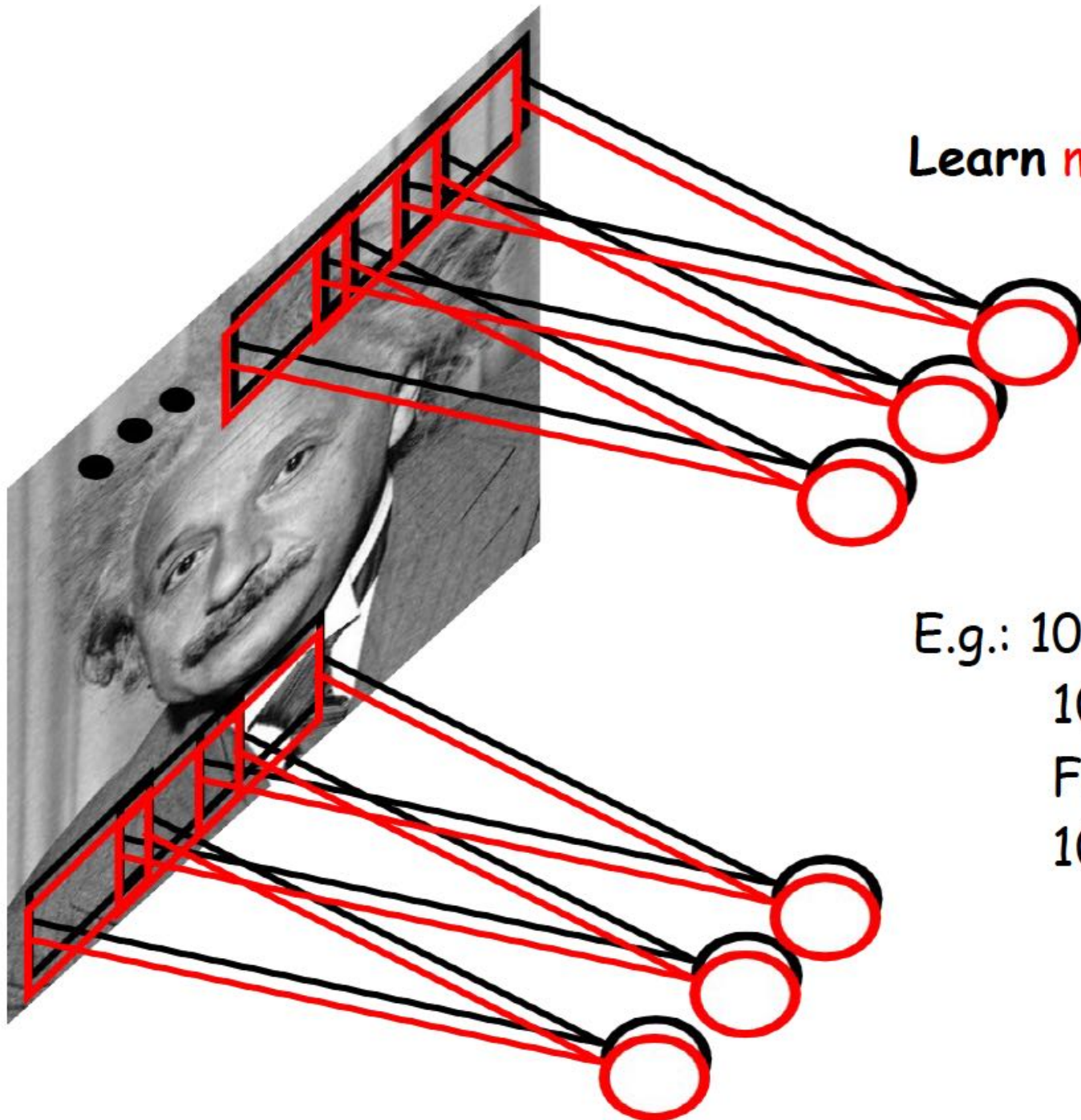1M hidden units
Filter size: 10x10
100M parameters

62

Ranzato

# CONVOLUTIONAL NET



Share the same parameters across different locations:
Convolutions with learned kernels

63

Ranzato

# CONVOLUTIONAL NET



**Learn** multiple filters.

E.g.: 1000x1000 image
    100 Filters
    Filter size: 10x10
    10K parameters

64

Ranzato

# NEURAL NETS FOR VISION

A standard neural net applied to images:

- scales quadratically with the size of the input
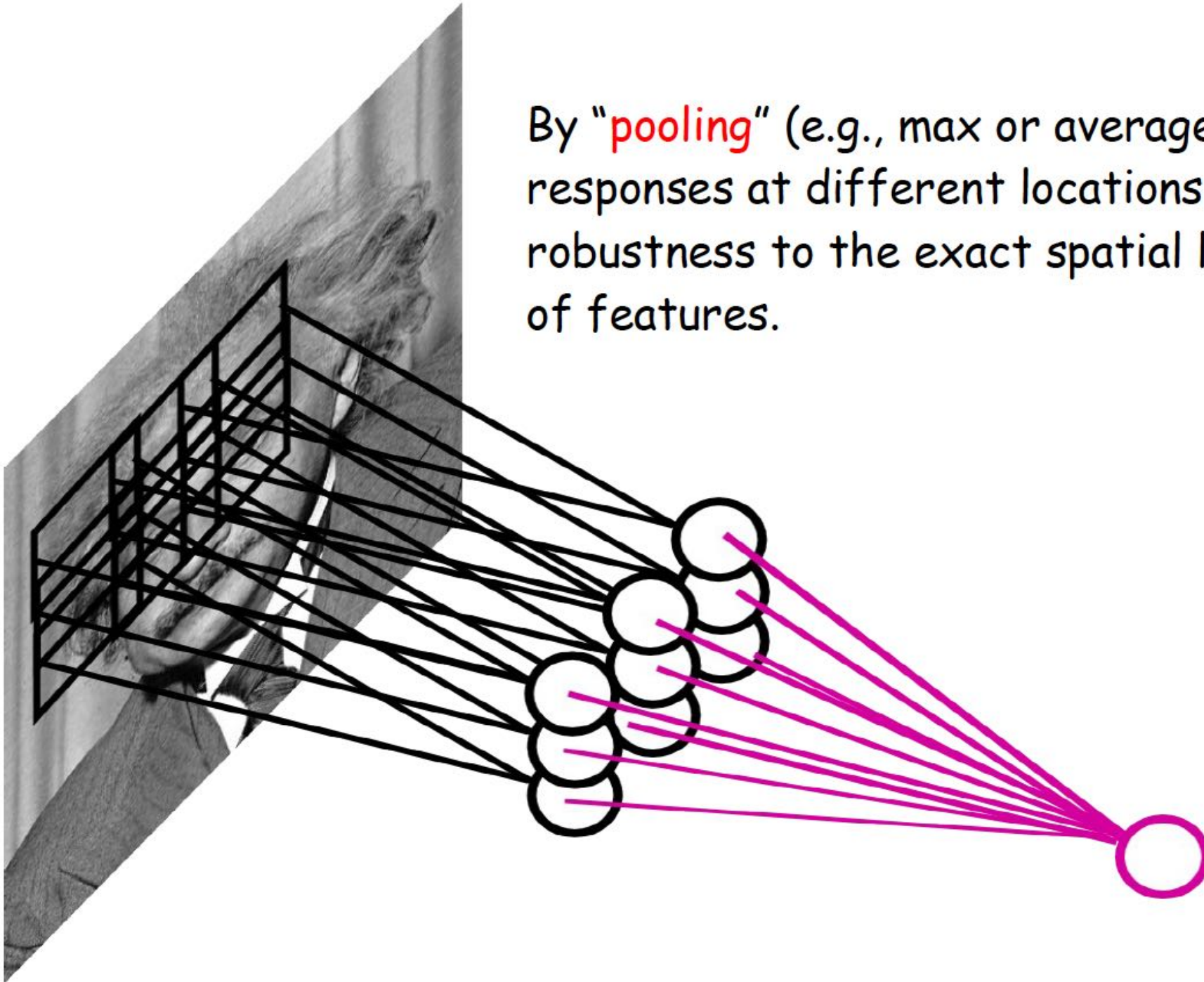
- does not leverage stationarity

Solution:

- connect each hidden unit to a small patch of the input

- share the weight across hidden units

This is called: **convolutional network.**

*LeCun et al. "Gradient-based learning applied to document recognition" IEEE 1998*
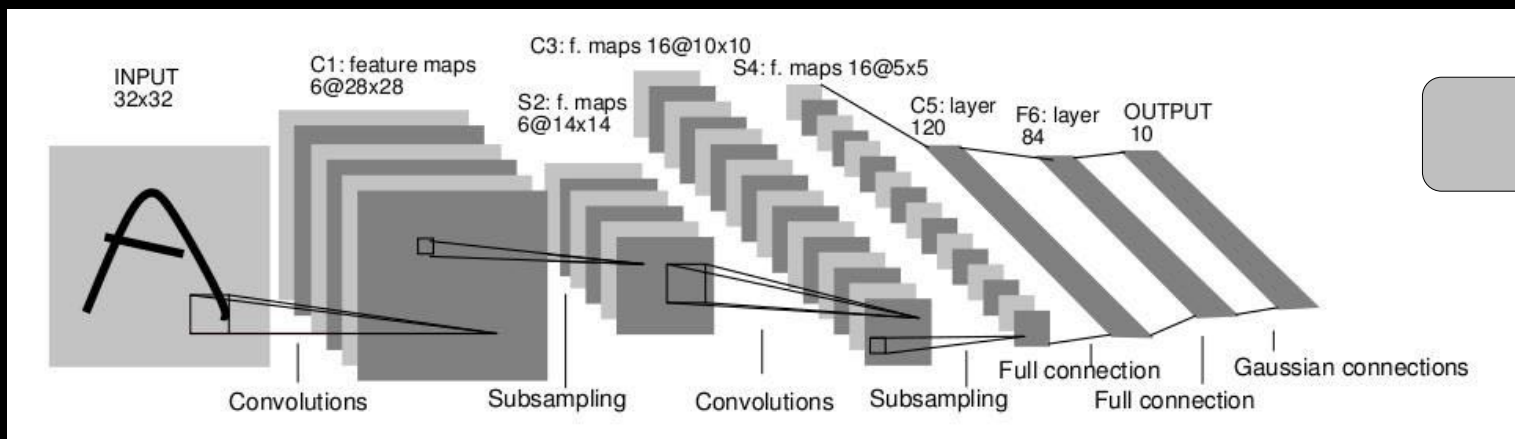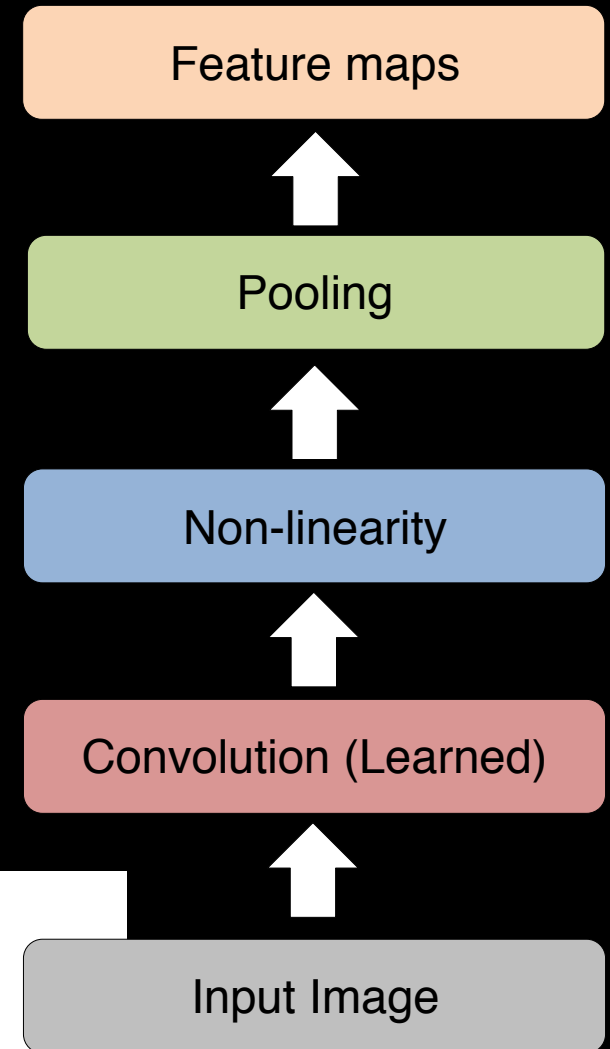
Ranzato

# CONVOLUTIONAL NET

By "pooling" (e.g., max or average) filter responses at different locations we gain robustness to the exact spatial location of features.

Ranzato

# Characteristics of Convnets

- Feed-forward:
  - Convolve input
  - Non-linearity (rectified linear)
  - Pooling (local max) / (=subsampling)
- Supervised
- Train convolutional filters by back-propagating classification error



Feature maps

↑

Pooling

↑

Non-linearity

↑

Convolution (Learned)

↑

Input Image

[LeCun et al. 1989]
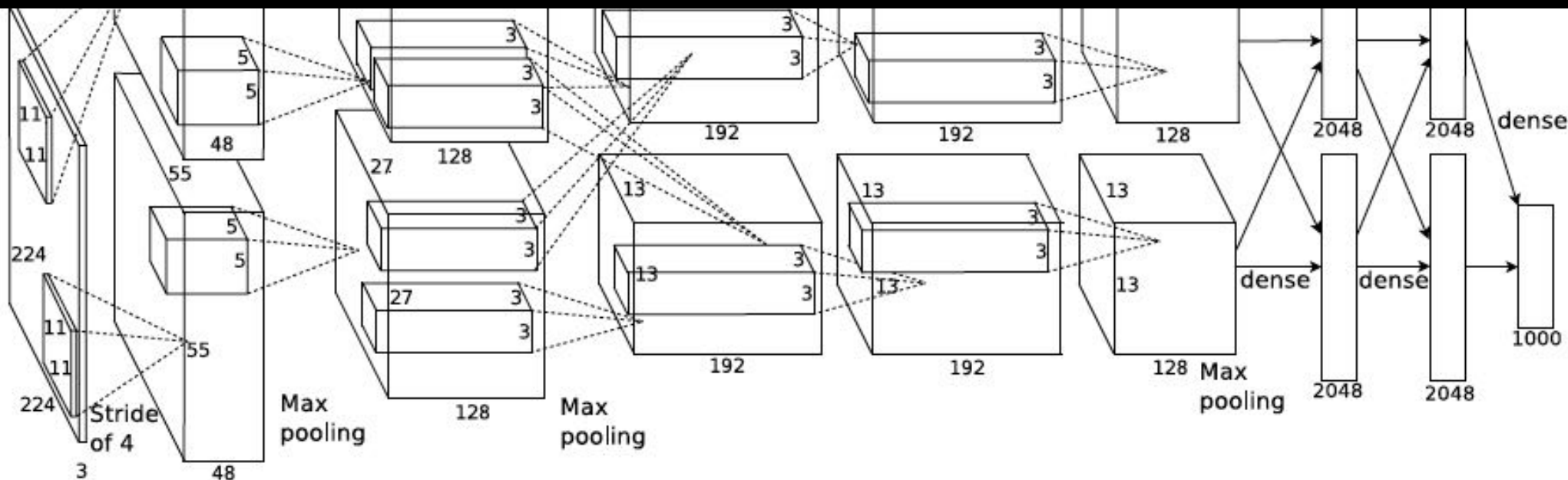
# Krizhevsky et al. [NIPS2012]

- Same model as LeCun'98 but:
  - Bigger model (8 layers)
    - More data ($10^6$ vs $10^3$ images)
    - GPU implementation (50x speedup over CPU)
    - Better regularization (DropOut)

**ImageNet Classification with Deep Convolutional Neural Networks**

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

- 7 hidden layers, 650,000 neurons, 60,000,000 parameters
- Trained on 2 GPUs for a week

**IMAGENET Large Scale Visual Recognition Challenge**

The Image Classification Challenge:
1,000 object classes
1,431,167 images

Output:
Scale
T-shirt
<u>Steel drum</u>
Drumstick
Mud turtle
✔

Output:
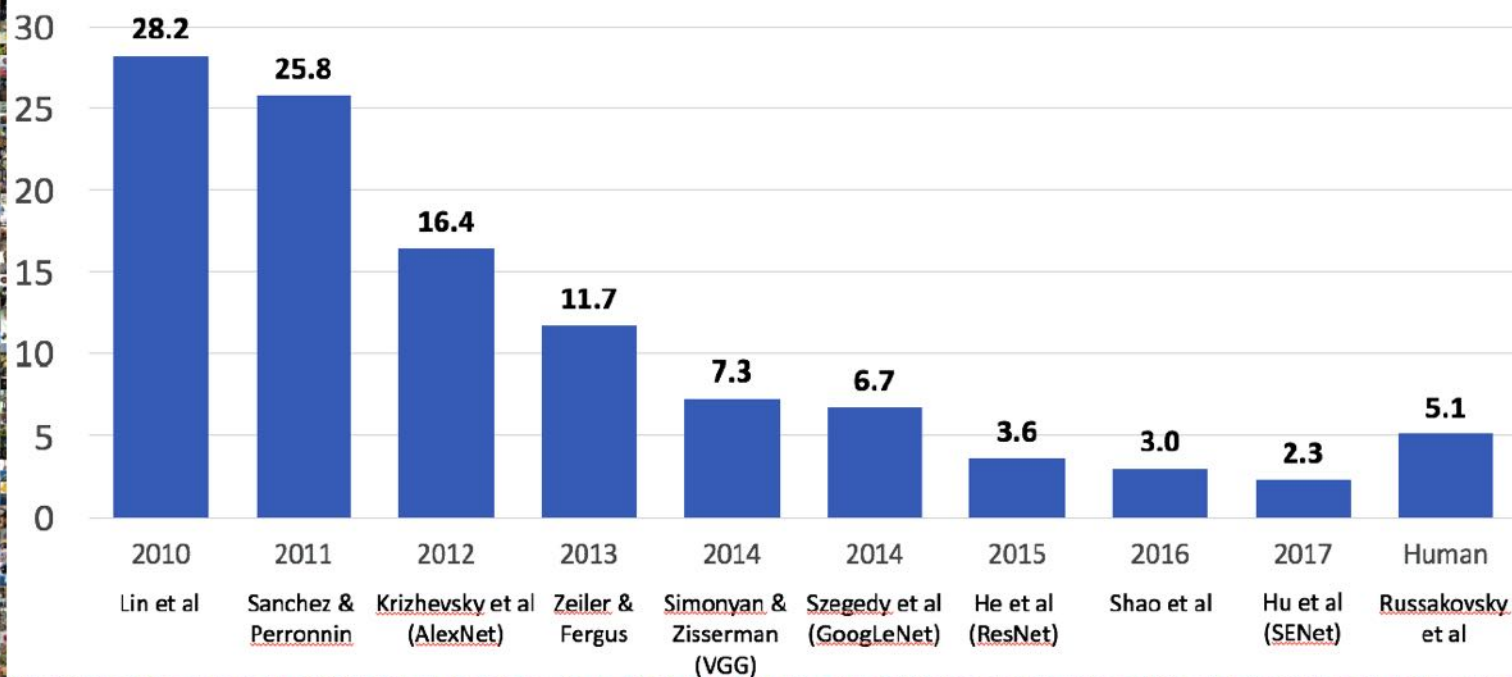Scale
T-shirt
Giant panda
Drumstick
Mud turtle
✗

Russakovsky et al. IJCV 2015

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

The Image Classification Challenge:
1,000 object classes
1,431,167 images

Russakovsky et al. IJCV 2015