



# 10-601 Introduction to Machine Learning

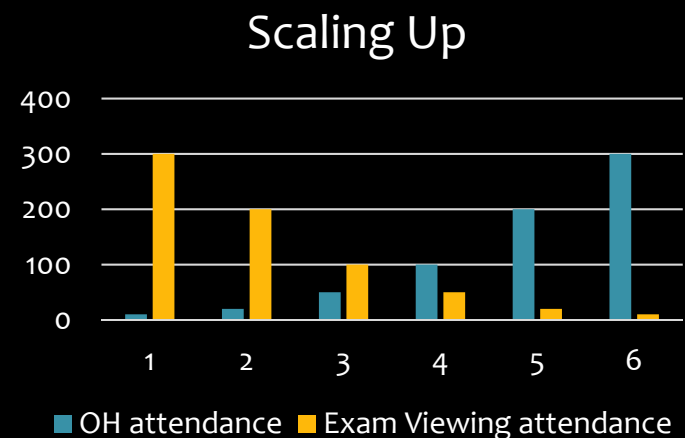
Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

# Backpropagation + Deep Learning

Matt Gormley  
Lecture 14  
Feb. 28, 2020

# Reminders

- **Homework 4: Logistic Regression**
  - Out: Wed, Feb. 19
  - Due: Fri, Feb. 28 at 11:59pm
- **Homework 5: Neural Networks**
  - Out: Fri, Feb. 28
  - Due: Wed, Mar. 18 at 11:59pm
- **Today's In-Class Poll**
  - <http://poll.mlcourse.org>
- **Exam Viewing**



# Q&A

**Q:** Do I need to know Matrix Calculus to derive the backprop algorithms used in this class?

**A:** No. We've carefully constructed our assignments so that you do **not** need to know Matrix Calculus.

That said, it's kind of handy.

# Matrix Calculus

Let  $y, x \in \mathbb{R}$  be scalars,  
 $\mathbf{y} \in \mathbb{R}^M$  and  $\mathbf{x} \in \mathbb{R}^P$   
 be vectors, and  
 $\mathbf{Y} \in \mathbb{R}^{M \times N}$  and  $\mathbf{X} \in \mathbb{R}^{P \times Q}$   
 be matrices

		Numerator		
Types of Derivatives		scalar	vector	matrix
Denominator	scalar	$\frac{\partial y}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{Y}}{\partial x}$
	gradient vector	$\frac{\partial y}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$
	matrix	$\frac{\partial y}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$



# Matrix Calculus

Types of Derivatives	<b>scalar</b>
<b>scalar</b>	$\frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x} \right]$
<b>vector</b>	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$
<b>matrix</b>	$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \frac{\partial y}{\partial X_{12}} & \cdots & \frac{\partial y}{\partial X_{1Q}} \\ \frac{\partial y}{\partial X_{21}} & \frac{\partial y}{\partial X_{22}} & \cdots & \frac{\partial y}{\partial X_{2Q}} \\ \vdots & & & \vdots \\ \frac{\partial y}{\partial X_{P1}} & \frac{\partial y}{\partial X_{P2}} & \cdots & \frac{\partial y}{\partial X_{PQ}} \end{bmatrix}$

# Matrix Calculus

Types of Derivatives	scalar	vector
scalar	$\frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x} \right]$	$\frac{\partial \mathbf{y}}{\partial x} = \left[ \frac{\partial y_1}{\partial x} \quad \frac{\partial y_2}{\partial x} \quad \dots \quad \frac{\partial y_N}{\partial x} \right]$
vector	$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \dots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$

# Matrix Calculus

## Common Vector Derivatives

Let  $\frac{\partial f(\vec{x})}{\partial \vec{x}} = \nabla_x f(\vec{x})$  be the vector derivative of  $f$ ,  $B \in \mathbb{R}^{m \times m}$   
 $x \in \mathbb{R}^m$

### Scalar Derivative

$$f(x) \rightarrow \frac{df}{dx}$$

.....

$$bx \rightarrow b$$

$$xb \rightarrow b$$

$$x^2 \rightarrow 2x$$

$$bx^2 \rightarrow 2bx$$

### Vector Derivative

$$f(x) \rightarrow \frac{df}{dx}$$

.....

~~$$xB \rightarrow B$$~~

$$x^T b \rightarrow b$$

$$x^T x \rightarrow 2x$$

$$x^T B x \rightarrow 2Bx$$

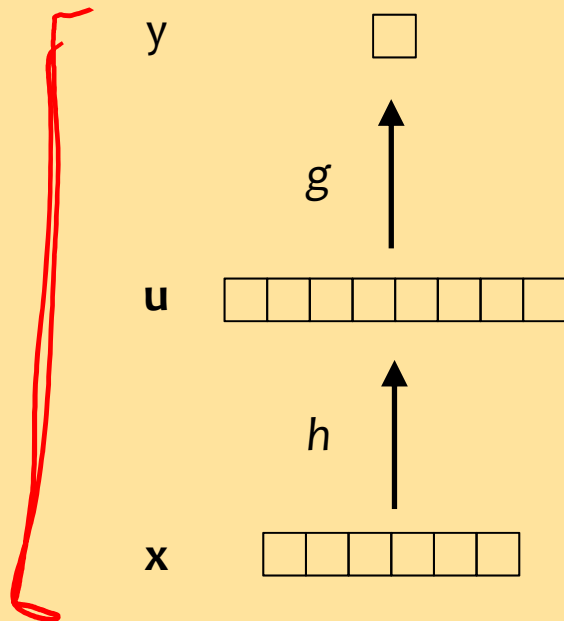
↑  $B$  symmetric

Not pol!

# Matrix Calculus

## Question:

Suppose  $y = g(\mathbf{u})$  and  $\mathbf{u} = h(\mathbf{x})$



Which of the following is the correct definition of the chain rule?

Recall:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_N}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_P} & \frac{\partial y_2}{\partial x_P} & \dots & \frac{\partial y_N}{\partial x_P} \end{bmatrix}$$

## Answer:

$$\frac{\partial y}{\partial \mathbf{x}} = \dots$$

- A.  $\frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$
- B.  $\frac{\partial y}{\partial \mathbf{u}}^T \frac{\partial \mathbf{u}}{\partial \mathbf{x}}$
- C.  $\frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}^T}{\partial \mathbf{x}}$
- D.  $\frac{\partial y}{\partial \mathbf{u}}^T \frac{\partial \mathbf{u}^T}{\partial \mathbf{x}}$
- E.  $\left(\frac{\partial y}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{x}}\right)^T$
- F. None of the above

Algorithm

# **BACKPROPAGATION**

Training

# Backpropagation

## Chalkboard

- Example: Backpropagation for Chain Rule #1

### Differentiation Quiz #1:

Suppose  $x = 2$  and  $z = 3$ , what are  $dy/dx$  and  $dy/dz$  for the function below? **Round your answer to the nearest integer.**

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Training

# Backpropagation

## *Chalkboard*

- SGD for Neural Network
- Example: Backpropagation for Neural Network

# Training

# Backpropagation

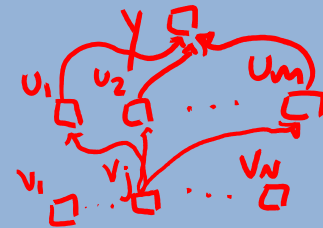
## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

### Forward Computation

1. Write an **algorithm** for evaluating the function  $y = f(x)$ . The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the “**computation graph**”)
2. Visit each node in **topological order**.  
For variable  $u_i$  with inputs  $v_1, \dots, v_N$ 
  - a. Compute  $u_i = g_i(v_1, \dots, v_N)$
  - b. Store the result at the node

### Backward Computation (Version A)

1. **Initialize**  $dy/dy = 1$ .
2. Visit each node  $v_j$  in **reverse topological order**.  
Let  $u_1, \dots, u_M$  denote all the nodes with  $v_j$  as an input  
Assuming that  $y = h(\mathbf{u}) = h(u_1, \dots, u_M)$   
and  $\mathbf{u} = g(\mathbf{v})$  or equivalently  $u_i = g_i(v_1, \dots, v_j, \dots, v_N)$  for all  $i$ 
  - a. We already know  $dy/du_i$  for all  $i$
  - b. Compute  $dy/dv_j$  as below (Choice of algorithm ensures computing  $(du_i/dv_j)$  is easy)



Return partial derivatives  $dy/du_i$  for all variables



## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

### Forward Computation

1. Write an **algorithm** for evaluating the function  $y = f(x)$ . The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the “**computation graph**”)
2. Visit each node in **topological order**.  
For variable  $u_i$  with inputs  $v_1, \dots, v_N$ 
  - a. Compute  $u_i = g_i(v_1, \dots, v_N)$
  - b. Store the result at the node

### Backward Computation (Version B)

1. **Initialize** all partial derivatives  $dy/du_i$  to 0 and  $dy/dy = 1$ .
2. Visit each node in **reverse topological order**.  
For variable  $u_i = g_i(v_1, \dots, v_N)$ 
  - a. We already know  $dy/du_i$
  - b. Increment  $dy/dv_j$  by  $(dy/du_i)(du_i/dv_j)$   
(Choice of algorithm ensures computing  $(du_i/dv_j)$  is easy)



**Return** partial derivatives  $dy/du_i$  for all variables

*Why is the backpropagation algorithm efficient?*

1. Reuses **computation from the forward pass** in the backward pass
2. Reuses **partial derivatives** throughout the backward pass (*but only if the algorithm reuses shared computation in the forward pass*)

(Key idea: partial derivatives in the backward pass should be thought of as variables stored for reuse)

Example: 1-Hidden Layer Neural Network

---

## Algorithm 1 Stochastic Gradient Descent (SGD)

---

```

1: procedure SGD(Training data  $\mathcal{D}$ , test data  $\mathcal{D}_t$ )
2:   Initialize parameters  $\alpha, \beta$  randomly, zero
3:   for  $e \in \{1, 2, \dots, E\}$  do ← epochs
4:     for  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$  do ←
5:       Compute neural network layers:
6:        $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{b}, \mathbf{z}, \hat{\mathbf{y}}, J) = \text{NNFORWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta)$ 
7:       Compute gradients via backprop:
8:        $\left. \begin{array}{l} \mathbf{g}_\alpha = \nabla_\alpha J \\ \mathbf{g}_\beta = \nabla_\beta J \end{array} \right\} = \text{NNBACKWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta, \mathbf{o})$ 
9:       Update parameters:
10:       $\alpha \leftarrow \alpha - \gamma \mathbf{g}_\alpha$ 
11:       $\beta \leftarrow \beta - \gamma \mathbf{g}_\beta$ 
12:      Evaluate training mean cross-entropy  $J_{\mathcal{D}}(\alpha, \beta)$ 
13:      Evaluate test validation mean cross-entropy  $J_{\mathcal{D}_t}(\alpha, \beta)$ 
14:      return parameters  $\alpha, \beta$ 

```

**Simple Example:** The goal is to compute  $J = \cos(\sin(x^2) + 3x^2)$  on the forward pass and the derivative  $\frac{dJ}{dx}$  on the backward pass.

Forward

$$J = \cos(u)$$

$$u = u_1 + u_2$$

$$u_1 = \sin(t)$$

$$u_2 = 3t$$

$$t = x^2$$

# Training

# Backpropagation

**Simple Example:** The goal is to compute  $J = \cos(\sin(x^2) + 3x^2)$  on the forward pass and the derivative  $\frac{dJ}{dx}$  on the backward pass.

Forward	Backward
$J = \cos(u)$	$\frac{dJ}{du} += -\sin(u)$
$u = u_1 + u_2$	$\frac{dJ}{du_1} += \frac{dJ}{du} \frac{du}{du_1}, \quad \frac{du}{du_1} = 1$ $\frac{dJ}{du_2} += \frac{dJ}{du} \frac{du}{du_2}, \quad \frac{du}{du_2} = 1$
$u_1 = \sin(t)$	$\frac{dJ}{dt} += \frac{dJ}{du_1} \frac{du_1}{dt}, \quad \frac{du_1}{dt} = \cos(t)$
$u_2 = 3t$	$\frac{dJ}{dt} += \frac{dJ}{du_2} \frac{du_2}{dt}, \quad \frac{du_2}{dt} = 3$
$t = x^2$	$\frac{dJ}{dx} += \frac{dJ}{dt} \frac{dt}{dx}, \quad \frac{dt}{dx} = 2x$

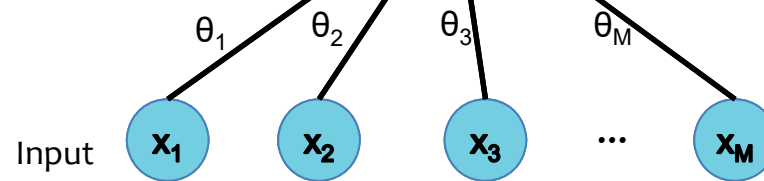
# Training

# Backpropagation

Output

y

Case 1:  
Logistic  
Regression



Forward

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$y = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{j=0}^D \theta_j x_j$$

Backward

$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

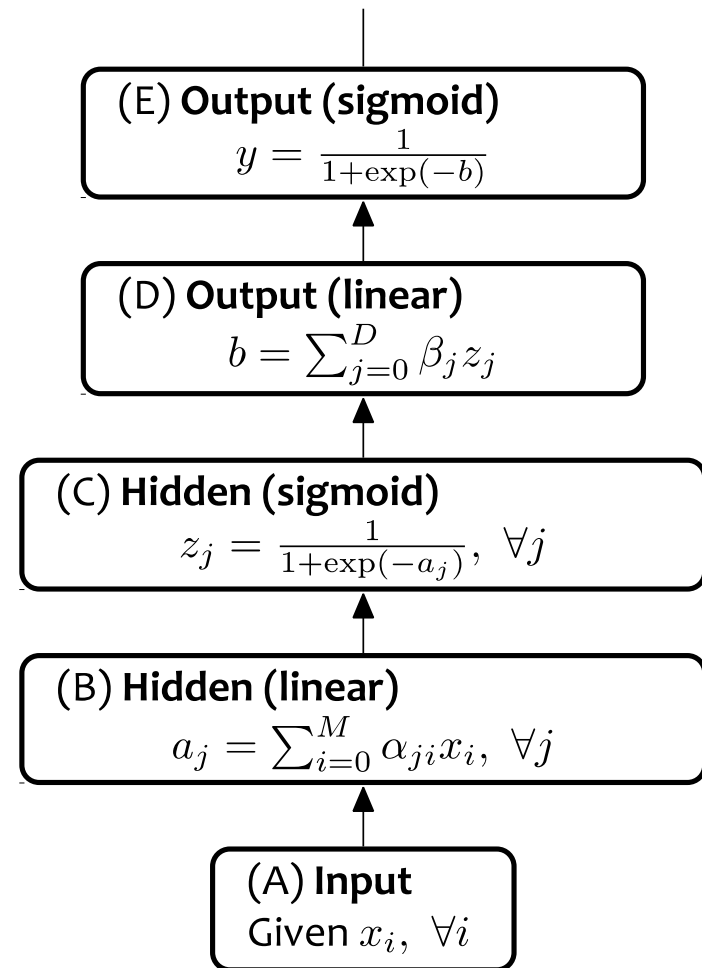
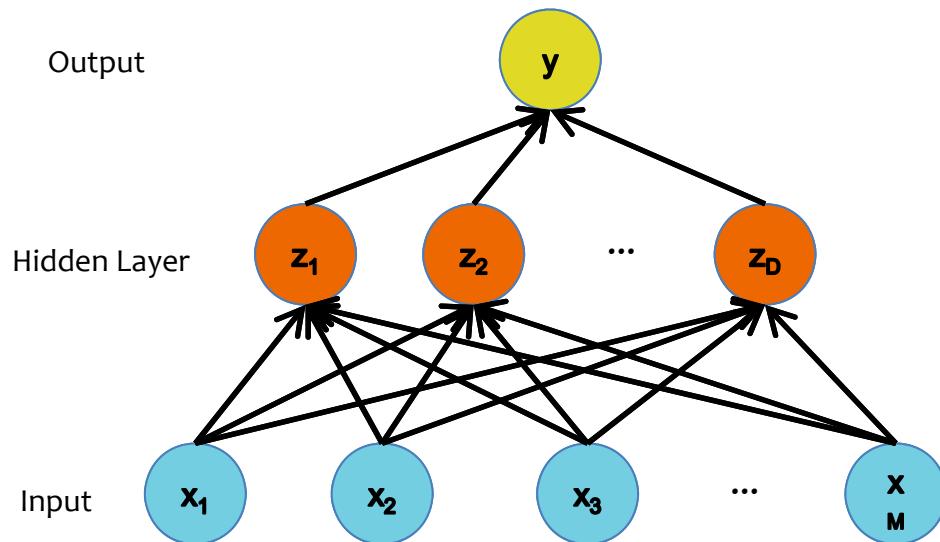
$$\frac{dJ}{da} = \frac{dJ}{dy} \frac{dy}{da}, \quad \frac{dy}{da} = \frac{\exp(-a)}{(\exp(-a) + 1)^2}$$

$$\frac{dJ}{d\theta_j} = \frac{dJ}{da} \frac{da}{d\theta_j}, \quad \frac{da}{d\theta_j} = x_j$$

$$\frac{dJ}{dx_j} = \frac{dJ}{da} \frac{da}{dx_j}, \quad \frac{da}{dx_j} = \theta_j$$

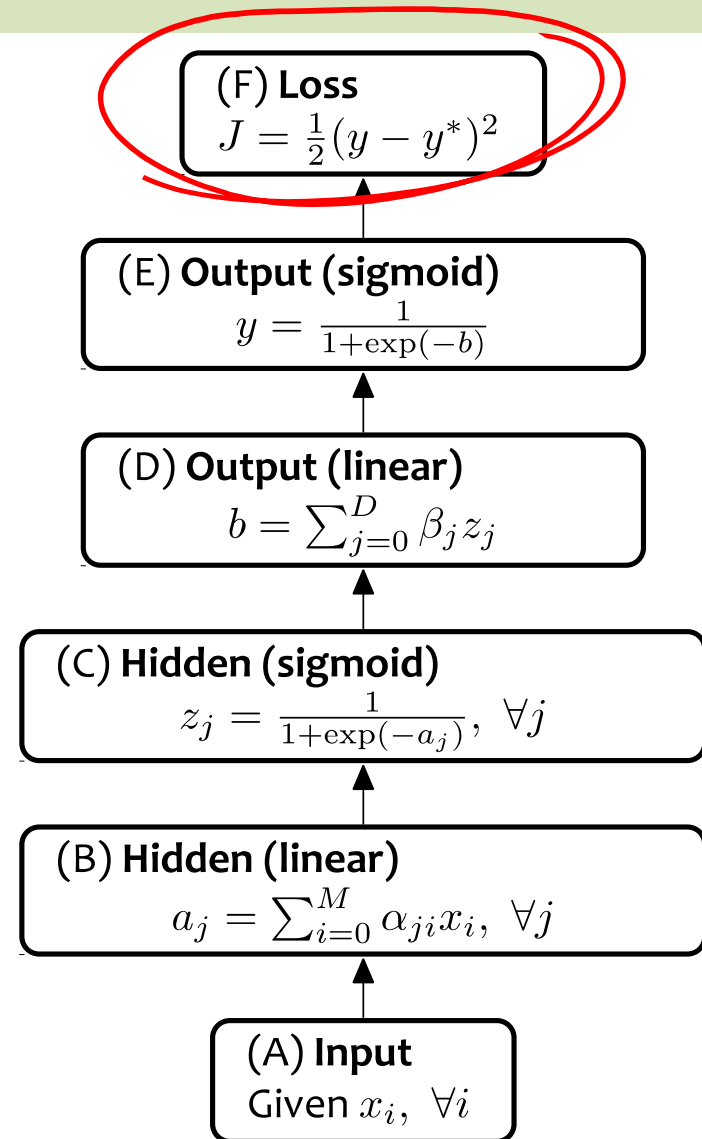
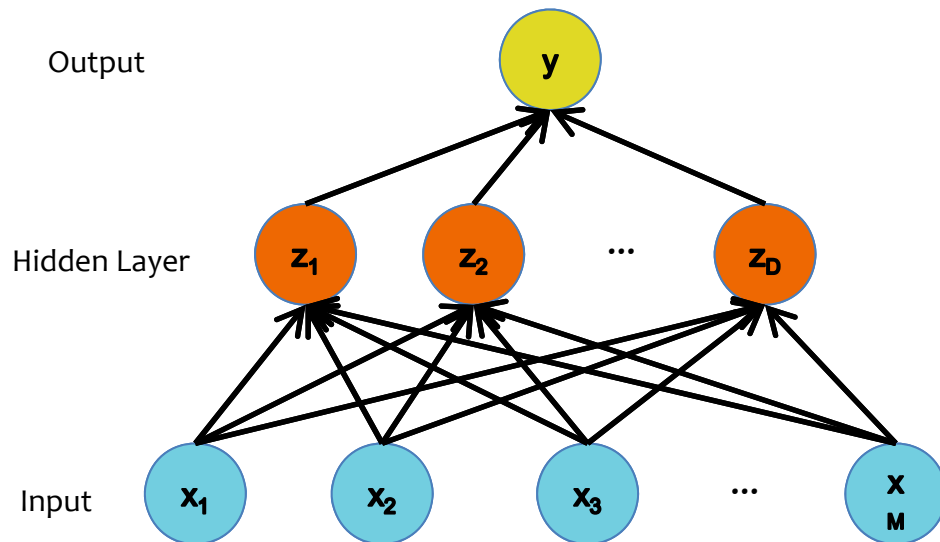
# Training

# Backpropagation



# Training

# Backpropagation

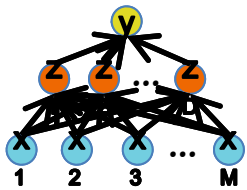




# Training

# Backpropagation

## Case 2: Neural Network



Forward

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$y = \frac{1}{1 + \exp(-b)}$$

$$b = \sum_{j=0}^D \beta_j z_j$$

$$z_j = \frac{1}{1 + \exp(-a_j)}$$

$$a_j = \sum_{i=0}^M \alpha_{ji} x_i$$

Backward

$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

$$\frac{dJ}{db} = \frac{dJ}{dy} \frac{dy}{db}, \quad \frac{dy}{db} = \frac{\exp(-b)}{(\exp(-b) + 1)^2}$$

$$\frac{dJ}{d\beta_j} = \frac{dJ}{db} \frac{db}{d\beta_j}, \quad \frac{db}{d\beta_j} = z_j$$

$$\frac{dJ}{dz_j} = \frac{dJ}{db} \frac{db}{dz_j}, \quad \frac{db}{dz_j} = \beta_j$$

$$\frac{dJ}{da_j} = \frac{dJ}{dz_j} \frac{dz_j}{da_j}, \quad \frac{dz_j}{da_j} = \frac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$$

$$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j} \frac{da_j}{d\alpha_{ji}}, \quad \frac{da_j}{d\alpha_{ji}} = x_i$$

$$\frac{dJ}{dx_i} = \sum_{j=0}^D \frac{dJ}{da_j} \frac{da_j}{dx_i}, \quad \frac{da_j}{dx_i} = \alpha_{ji}$$

# Training

# Backpropagation

Case 2:

Forward

Backward

Loss

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

Sigmoid

$$y = \frac{1}{1 + \exp(-b)}$$

$$\frac{dJ}{db} = \frac{dJ}{dy} \frac{dy}{db}, \quad \frac{dy}{db} = \frac{\exp(-b)}{(\exp(-b) + 1)^2}$$

Linear

$$b = \sum_{j=0}^D \beta_j z_j$$

$$\frac{dJ}{d\beta_j} = \frac{dJ}{db} \frac{db}{d\beta_j}, \quad \frac{db}{d\beta_j} = z_j$$

$$\frac{dJ}{dz_j} = \frac{dJ}{db} \frac{db}{dz_j}, \quad \frac{db}{dz_j} = \beta_j$$

Sigmoid

$$z_j = \frac{1}{1 + \exp(-a_j)}$$

$$\frac{dJ}{da_j} = \frac{dJ}{dz_j} \frac{dz_j}{da_j}, \quad \frac{dz_j}{da_j} = \frac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$$

Linear

$$a_j = \sum_{i=0}^M \alpha_{ji} x_i$$

$$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j} \frac{da_j}{d\alpha_{ji}}, \quad \frac{da_j}{d\alpha_{ji}} = x_i$$

$$\frac{dJ}{dx_i} = \sum_{j=0}^D \frac{dJ}{da_j} \frac{da_j}{dx_i}, \quad \frac{da_j}{dx_i} = \alpha_{ji}$$

# Derivative of a Sigmoid

First suppose that

$$s = \frac{1}{1 + \exp(-b)} \quad (1)$$

To obtain the simplified form of the derivative of a sigmoid.

$$\frac{ds}{db} = \frac{\exp(-b)}{(\exp(-b) + 1)^2} \quad (2)$$

$$= \frac{\exp(-b) + 1 - 1}{(\exp(-b) + 1 + 1 - 1)^2} \quad (3)$$

$$= \frac{\exp(-b) + 1 - 1}{(\exp(-b) + 1)^2} \quad (4)$$

$$= \frac{\exp(-b) + 1}{(\exp(-b) + 1)^2} - \frac{1}{(\exp(-b) + 1)^2} \quad (5)$$

$$= \frac{1}{(\exp(-b) + 1)} - \frac{1}{(\exp(-b) + 1)^2} \quad (6)$$

$$= \frac{1}{(\exp(-b) + 1)} - \left( \frac{1}{(\exp(-b) + 1)} \frac{1}{(\exp(-b) + 1)} \right) \quad (7)$$

$$= \frac{1}{(\exp(-b) + 1)} \left( 1 - \frac{1}{(\exp(-b) + 1)} \right) \quad (8)$$

$$= s(1 - s) \quad (9)$$

# Training

# Backpropagation

Case 2:

Forward

Backward

Loss

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$$

Sigmoid

$$y = \frac{1}{1 + \exp(-b)}$$

$$\frac{dJ}{db} = \frac{dJ}{dy} \frac{dy}{db}, \quad \frac{dy}{db} = \frac{\exp(-b)}{(\exp(-b) + 1)^2}$$

Linear

$$b = \sum_{j=0}^D \beta_j z_j$$

$$\frac{dJ}{d\beta_j} = \frac{dJ}{db} \frac{db}{d\beta_j}, \quad \frac{db}{d\beta_j} = z_j$$

Sigmoid

$$z_j = \frac{1}{1 + \exp(-a_j)}$$

$$\frac{dJ}{da_j} = \frac{dJ}{dz_j} \frac{dz_j}{da_j}, \quad \frac{dz_j}{da_j} = \frac{\exp(-a_j)}{(\exp(-a_j) + 1)^2}$$

Linear

$$a_j = \sum_{i=0}^M \alpha_{ji} x_i$$

$$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j} \frac{da_j}{d\alpha_{ji}}, \quad \frac{da_j}{d\alpha_{ji}} = x_i$$

$$\frac{dJ}{dx_i} = \sum_{j=0}^D \frac{dJ}{da_j} \frac{da_j}{dx_i}, \quad \frac{da_j}{dx_i} = \alpha_{ji}$$

# Training

# Backpropagation

Case 2:	Forward	Backward
Loss	$J = y^* \log y + (1 - y^*) \log(1 - y)$	$\frac{dJ}{dy} = \frac{y^*}{y} + \frac{(1 - y^*)}{y - 1}$
Sigmoid	$y = \frac{1}{1 + \exp(-b)}$	$\frac{dJ}{db} = \frac{dJ}{dy} \frac{dy}{db}, \frac{dy}{db} = y(1 - y)$
Linear	$b = \sum_{j=0}^D \beta_j z_j$	$\frac{dJ}{d\beta_j} = \frac{dJ}{db} \frac{db}{d\beta_j}, \frac{db}{d\beta_j} = z_j$
Sigmoid	$z_j = \frac{1}{1 + \exp(-a_j)}$	$\frac{dJ}{da_j} = \frac{dJ}{dz_j} \frac{dz_j}{da_j}, \frac{dz_j}{da_j} = z_j(1 - z_j)$
Linear	$a_j = \sum_{i=0}^M \alpha_{ji} x_i$	$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j} \frac{da_j}{d\alpha_{ji}}, \frac{da_j}{d\alpha_{ji}} = x_i$ $\frac{dJ}{dx_i} = \sum_{j=0}^D \frac{dJ}{da_j} \frac{da_j}{dx_i}, \frac{da_j}{dx_i} = \alpha_{ji}$

# Training

# Backpropagation

## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

### Forward Computation

1. Write an **algorithm** for evaluating the function  $y = f(\mathbf{x})$ . The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the “**computation graph**”)
2. Visit each node in **topological order**.  
For variable  $u_i$  with inputs  $v_1, \dots, v_N$ 
  - a. Compute  $u_i = g_i(v_1, \dots, v_N)$
  - b. Store the result at the node

### Backward Computation (Version A)

1. **Initialize**  $dy/dy = 1$ .
2. Visit each node  $v_j$  in **reverse topological order**.  
Let  $u_1, \dots, u_M$  denote all the nodes with  $v_j$  as an input  
Assuming that  $y = h(\mathbf{u}) = h(u_1, \dots, u_M)$   
and  $\mathbf{u} = g(\mathbf{v})$  or equivalently  $u_i = g_i(v_1, \dots, v_j, \dots, v_N)$  for all  $i$ 
  - a. We already know  $dy/du_i$  for all  $i$
  - b. Compute  $dy/dv_j$  as below (Choice of algorithm ensures computing  $(du_i/dv_j)$  is easy)

$$\frac{dy}{dv_j} = \sum_{i=1}^M \frac{dy}{du_i} \frac{du_i}{dv_j}$$

**Return** partial derivatives  $dy/du_i$  for all variables

# Training

# Backpropagation

## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

### Forward Computation

1. Write an **algorithm** for evaluating the function  $y = f(\mathbf{x})$ . The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the “**computation graph**”)
2. Visit each node in **topological order**.  
For variable  $u_i$  with inputs  $v_1, \dots, v_N$ 
  - a. Compute  $u_i = g_i(v_1, \dots, v_N)$
  - b. Store the result at the node

### Backward Computation (Version B)

1. **Initialize** all partial derivatives  $dy/du_i$  to 0 and  $dy/dy = 1$ .
2. Visit each node in **reverse topological order**.  
For variable  $u_i = g_i(v_1, \dots, v_N)$ 
  - a. We already know  $dy/du_i$
  - b. Increment  $dy/dv_j$  by  $(dy/du_i)(du_i/dv_j)$   
(Choice of algorithm ensures computing  $(du_i/dv_j)$  is easy)

**Return** partial derivatives  $dy/du_i$  for all variables

Example: 1-Hidden Layer Neural Network

---

**Algorithm 1** Stochastic Gradient Descent (SGD)

---

```
1: procedure SGD(Training data  $\mathcal{D}$ , test data  $\mathcal{D}_t$ )
2:   Initialize parameters  $\alpha, \beta$ 
3:   for  $e \in \{1, 2, \dots, E\}$  do
4:     for  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$  do
5:       Compute neural network layers:
6:        $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{b}, \mathbf{z}, \hat{\mathbf{y}}, J) = \text{NNFORWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta)$ 
7:       Compute gradients via backprop:
8:        $\left. \begin{array}{l} \mathbf{g}_\alpha = \nabla_\alpha J \\ \mathbf{g}_\beta = \nabla_\beta J \end{array} \right\} = \text{NNBACKWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta, \mathbf{o})$ 
9:       Update parameters:
10:       $\alpha \leftarrow \alpha - \gamma \mathbf{g}_\alpha$ 
11:       $\beta \leftarrow \beta - \gamma \mathbf{g}_\beta$ 
12:      Evaluate training mean cross-entropy  $J_{\mathcal{D}}(\alpha, \beta)$ 
13:      Evaluate test mean cross-entropy  $J_{\mathcal{D}_t}(\alpha, \beta)$ 
14:   return parameters  $\alpha, \beta$ 
```

---



## Background

# Gradients

1. Given training data

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of the

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

opposite the gradient)


$$\theta^{(t)} \leftarrow \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$



# **OTHER APPROACHES TO DIFFERENTIATION**

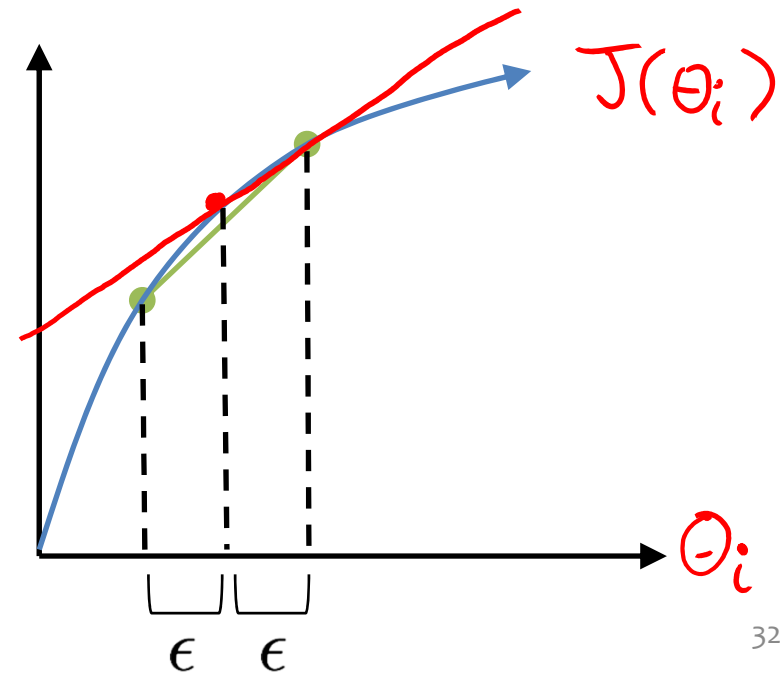
The *centered* finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \mathbf{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \mathbf{d}_i))}{2\epsilon} \quad (1)$$

where  $\mathbf{d}_i$  is a 1-hot vector consisting of all zeros except for the  $i$ th entry of  $\mathbf{d}_i$ , which has value 1.

### Notes:

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon



## Differentiation Quiz

Speed Quiz:  
2 minute time limit.

## Differentiation Quiz #1:

Suppose  $x = 2$  and  $z = 3$ , what are  $dy/dx$  and  $dy/dz$  for the function below? **Round your answer to the nearest integer.**

$$y = \cos(xz) + \frac{\log(xz)}{z} + \frac{\log(x)}{\exp(\sin(z))}$$

calamity

**Answer:** Answers below are in the form  $[dy/dx, dy/dz]$

~~A. [4.2, -7.2]~~

B. [7.2, -4.2]

C. [1.10, 2.7]

D. [2.7, 1.10]

E. [8.1, 12.2]

F. [12.2, 8.1]

G. [-0.97, -1.28]

H. [-1.28, -0.97]

## Differentiation Quiz #2:

A neural network with 2 hidden layers can be written as:

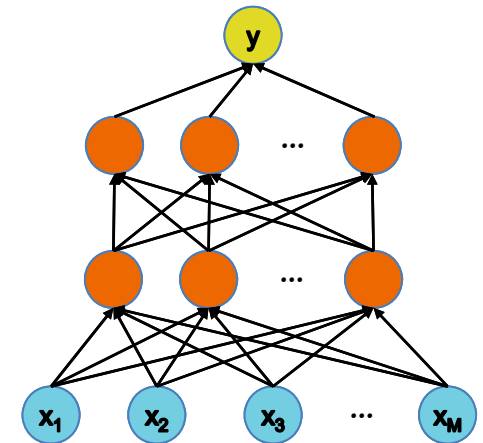
$$y = \sigma(\boldsymbol{\beta}^T \sigma((\boldsymbol{\alpha}^{(2)})^T \sigma((\boldsymbol{\alpha}^{(1)})^T \mathbf{x})))$$

where  $y \in \mathbb{R}$ ,  $\mathbf{x} \in \mathbb{R}^{D^{(0)}}$ ,  $\boldsymbol{\beta} \in \mathbb{R}^{D^{(2)}}$  and  $\boldsymbol{\alpha}^{(i)}$  is a  $D^{(i)} \times D^{(i-1)}$  matrix. Nonlinear functions are applied elementwise:

$$\sigma(\mathbf{a}) = [\sigma(a_1), \dots, \sigma(a_K)]^T$$

Let  $\sigma$  be sigmoid:  $\sigma(a) = \frac{1}{1 + \exp(-a)}$

What is  $\frac{\partial y}{\partial \beta_j}$  and  $\frac{\partial y}{\partial \alpha_j^{(i)}}$  for all  $i, j$ .



# Summary

## 1. **Neural Networks...**

- provide a way of learning features
- are highly nonlinear prediction functions
- (can be) a highly parallel network of logistic regression classifiers
- discover useful hidden representations of the input

## 2. **Backpropagation...**

- provides an efficient way to compute gradients
- is a special case of reverse-mode automatic differentiation

# Backprop Objectives

*You should be able to...*

- Construct a computation graph for a function as specified by an algorithm
- Carry out the backpropagation on an arbitrary computation graph
- Construct a computation graph for a neural network, identifying all the given and intermediate quantities that are relevant
- Instantiate the backpropagation algorithm for a neural network
- Instantiate an optimization method (e.g. SGD) and a regularizer (e.g. L2) when the parameters of a model are comprised of several matrices corresponding to different layers of a neural network
- Apply the empirical risk minimization framework to learn a neural network
- Use the finite difference method to evaluate the gradient of a function
- Identify when the gradient of a function can be computed at all and when it can be computed efficiently

# DEEP LEARNING

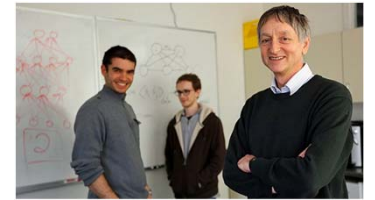


# Deep Learning Outline

- **Background: Computer Vision**
  - Image Classification
  - ILSVRC 2010 - 2016
  - Traditional Feature Extraction Methods
  - Convolution as Feature Extraction
- **Convolutional Neural Networks (CNNs)**
  - Learning Feature Abstractions
  - Common CNN Layers:
    - Convolutional Layer
    - Max-Pooling Layer
    - Fully-connected Layer (w/tensor input)
    - Softmax Layer
    - ReLU Layer
  - Background: Subgradient
  - Architecture: LeNet
  - Architecture: AlexNet
- **Training a CNN**
  - SGD for CNNs
  - Backpropagation for CNNs

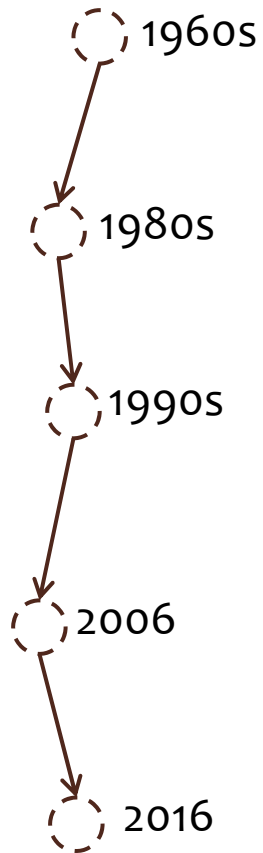
# Why is everyone talking about Deep Learning?

- Because a lot of money is invested in it...
  - DeepMind: Acquired by Google for **\$400 million**
  - DNNResearch: **Three person startup** (including Geoff Hinton) acquired by Google for unknown price tag
  - Enlitic, Ersatz, MetaMind, Nervana, Skylab: Deep Learning startups commanding **millions of VC dollars**
- Because it made the **front page** of the New York Times



**The New York Times**

# Why is everyone talking about Deep Learning?



## Deep learning:

- Has won numerous pattern recognition competitions
- Does so with minimal feature engineering

### This wasn't always the case!

Since 1980s: Form of models hasn't changed much, but lots of new tricks...

- More hidden units
- Better (online) optimization
- New nonlinear functions (ReLUs)
- Faster computers (CPUs and GPUs)

# **BACKGROUND: COMPUTER VISION**

# Example: Image Classification

- ImageNet LSVRC-2011 contest:
  - **Dataset:** 1.2 million labeled images, 1000 classes
  - **Task:** Given a new image, label it with the correct class
  - **Multiclass** classification problem
- Examples from <http://image-net.org/>

## Bird

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126 pictures

92.85% Popularity Percentile

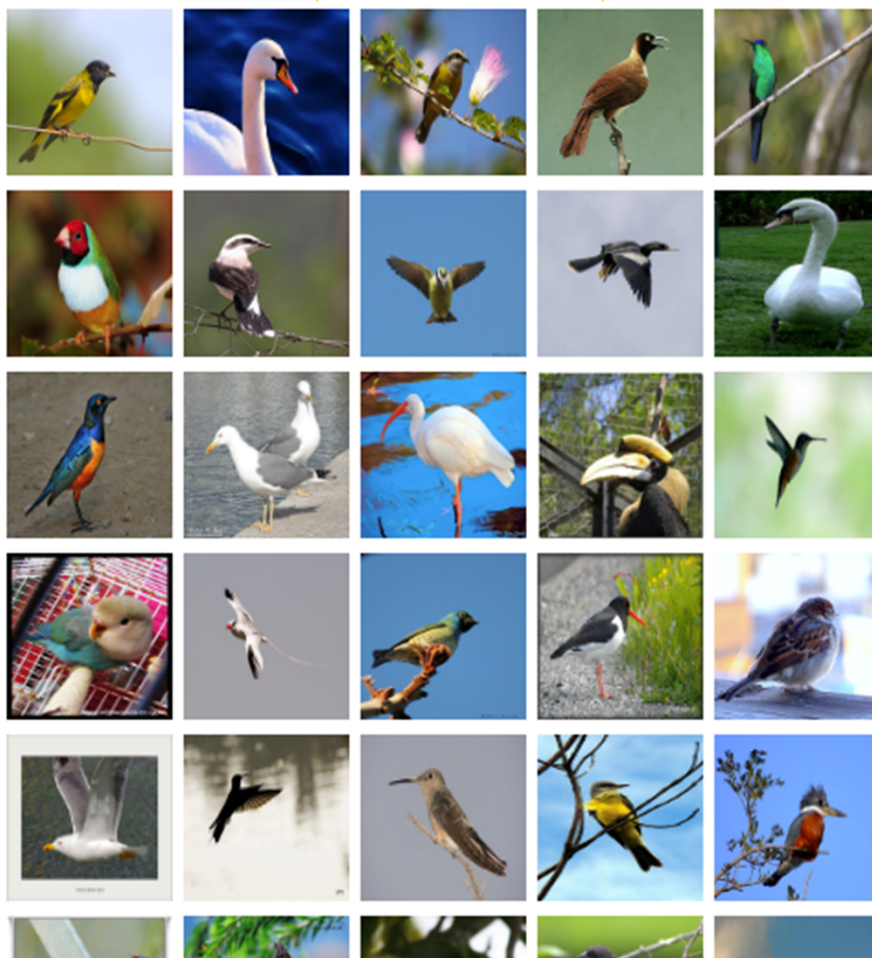


- marine animal, marine creature, sea animal, sea creature (1)
- scavenger (1)
- biped (0)
- predator, predatory animal (1)
- larva (49)
- acrodont (0)
- feeder (0)
- stunt (0)
- chordate (3087)**
  - tunicate, urochordate, urochord (6)
  - cephalochordate (1)
  - vertebrate, craniate (3077)**
    - mammal, mammalian (1169)
    - bird (871)**
      - dickeybird, dickey-bird, dickybird, dicky-bird (0)
      - cock (1)
      - hen (0)
      - nester (0)
      - night bird (1)
      - bird of passage (0)
      - protoavis (0)
      - archaeopteryx, archeopteryx, Archaeopteryx lithographi (0)
      - Sinornis (0)
      - Ibero-mesornis (0)
      - archaeornis (0)
      - ratite, ratite bird, flightless bird (10)
      - carinate, carinate bird, flying bird (0)
      - passerine, passeriform bird (279)
      - nonpasserine bird (0)
      - bird of prey, raptor, raptorial bird (80)
      - gallinaceous bird, gallinacean (114)

### Treemap Visualization

### Images of the Synset

### Downloads





# German iris, *Iris kochii*

Iris of northern Italy having deep blue-purple flowers; similar to but smaller than *Iris germanica*

469 pictures

49.6% Popularity Percentile



- halophyte (0)
- succulent (39)
- cultivar (0)
- cultivated plant (0)
- weed (54)
- evergreen, evergreen plant (0)
- deciduous plant (0)
- vine (272)
- creeper (0)
- woody plant, ligneous plant (1868)
- geophyte (0)
- desert plant, xerophyte, xerophytic plant, xerophile, xerophilic mesophyte, mesophytic plant (0)
- aquatic plant, water plant, hydrophyte, hydrophytic plant (11)
- tuberous plant (0)
- bulbous plant (179)**
  - iridaceous plant (27)**
    - iris, flag, fleur-de-lis, sword lily (19)**
      - bearded iris (4)**
        - Florentine iris, orris, *Iris germanica florentina*, Iris
        - German iris, *Iris germanica* (0)
        - German iris, *Iris kochii* (0)**
        - Dalmatian iris, *Iris pallida* (0)
      - beardless iris (4)
      - bulbous iris (0)
      - dwarf iris, *Iris cristata* (0)
      - stinking iris, gladdon, gladdon iris, stinking gladwyn,
      - Persian iris, *Iris persica* (0)
      - yellow iris, yellow flag, yellow water flag, *Iris pseudo*
      - dwarf iris, vernal iris, *Iris verna* (0)
      - blue flag, *Iris versicolor* (0)

Treemap Visualization

**Images of the Synset**

Downloads



# Court, courtyard

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

165 pictures

92.61% Popularity Percentile



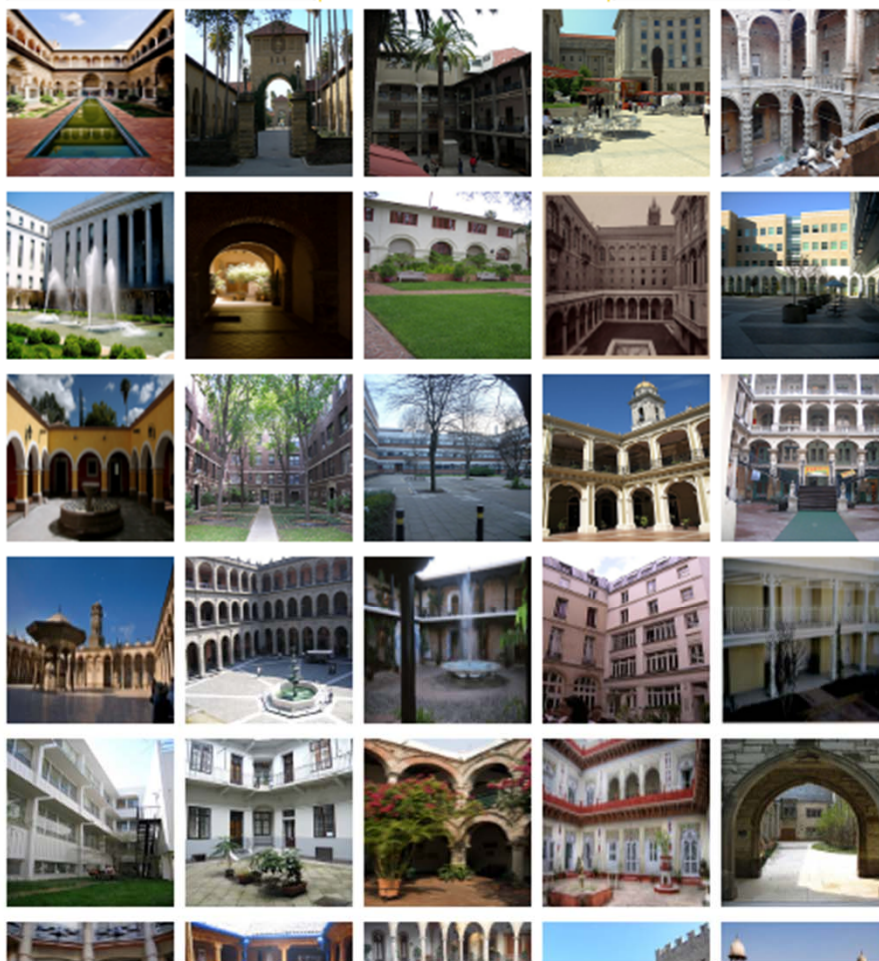
Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32326)
  - plant, flora, plant life (4486)
  - geological formation, formation (175)
  - natural object (1112)
  - sport, athletics (176)
  - artifact, artefact (10504)
    - instrumentality, instrumentation (5494)
    - structure, construction (1405)
      - airdock, hangar, repair shed (0)
      - altar (1)
      - arcade, colonnade (1)
      - arch (31)
      - area (344)
        - aisle (0)
        - auditorium (1)
        - baggage claim (0)
        - box (1)
        - breakfast area, breakfast nook (0)
        - bullpen (0)
        - chancel, sanctuary, bema (0)
        - choir (0)
        - corner, nook (2)
        - court, courtyard (6)
          - atrium (0)
          - bailey (0)
          - cloister (0)
          - food court (0)
          - forecourt (0)
          - narvis (0)

## Treemap Visualization

## Images of the Synset

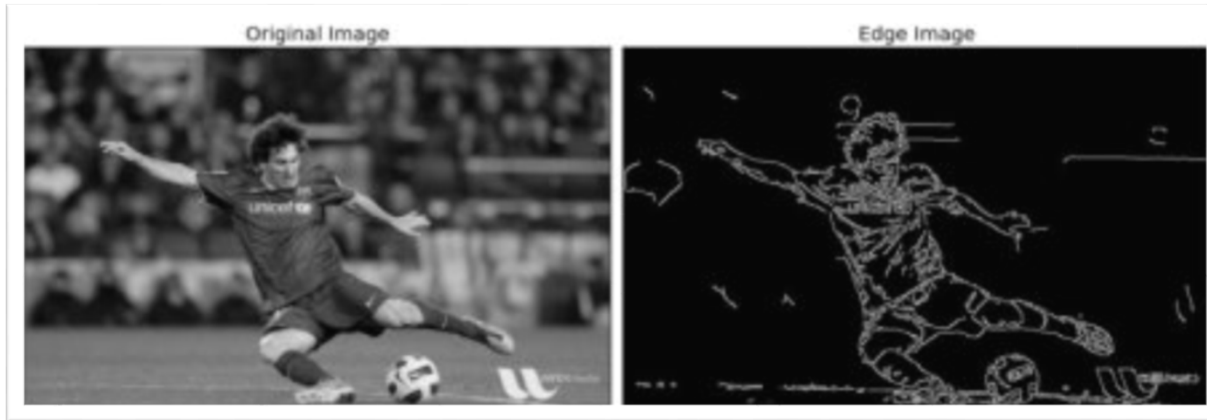
## Downloads



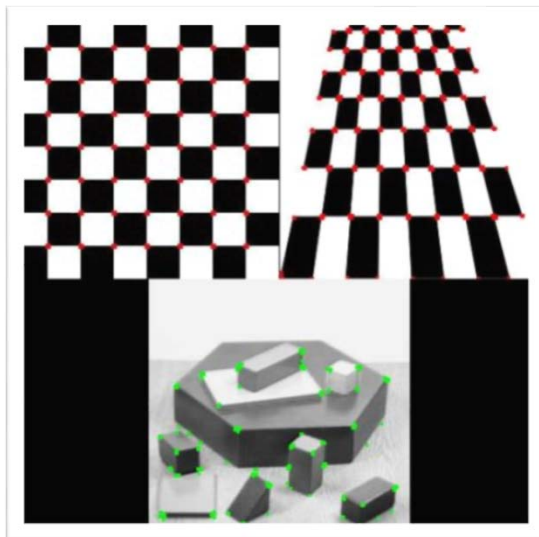


# Feature Engineering for CV

## Edge detection (Canny)



## Corner Detection (Harris)



## Scale Invariant Feature Transform (SIFT)



Figure 3: Model images of planar objects are shown in the top row. Recognition results below show model outlines and match keys used for matching.

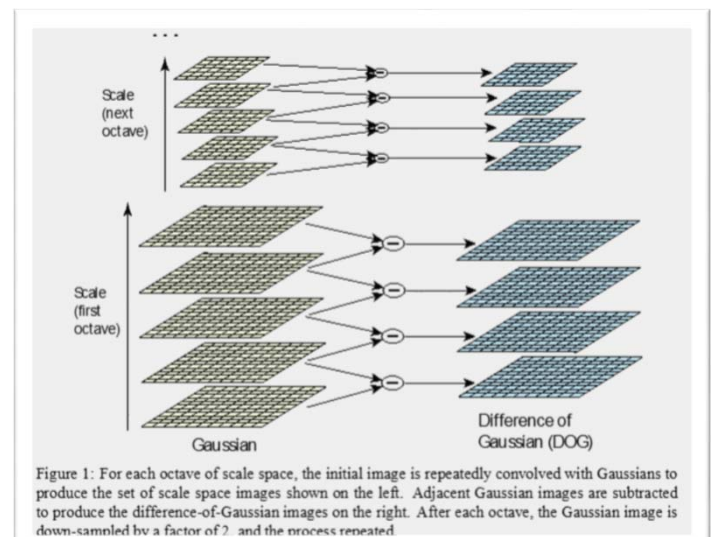


Figure 1: For each octave of scale space, the initial image is repeatedly convolved with Gaussians to produce the set of scale space images shown on the left. Adjacent Gaussian images are subtracted to produce the difference-of-Gaussian images on the right. After each octave, the Gaussian image is down-sampled by a factor of 2, and the process repeated.

# Example: Image Classification

## CNN for Image Classification

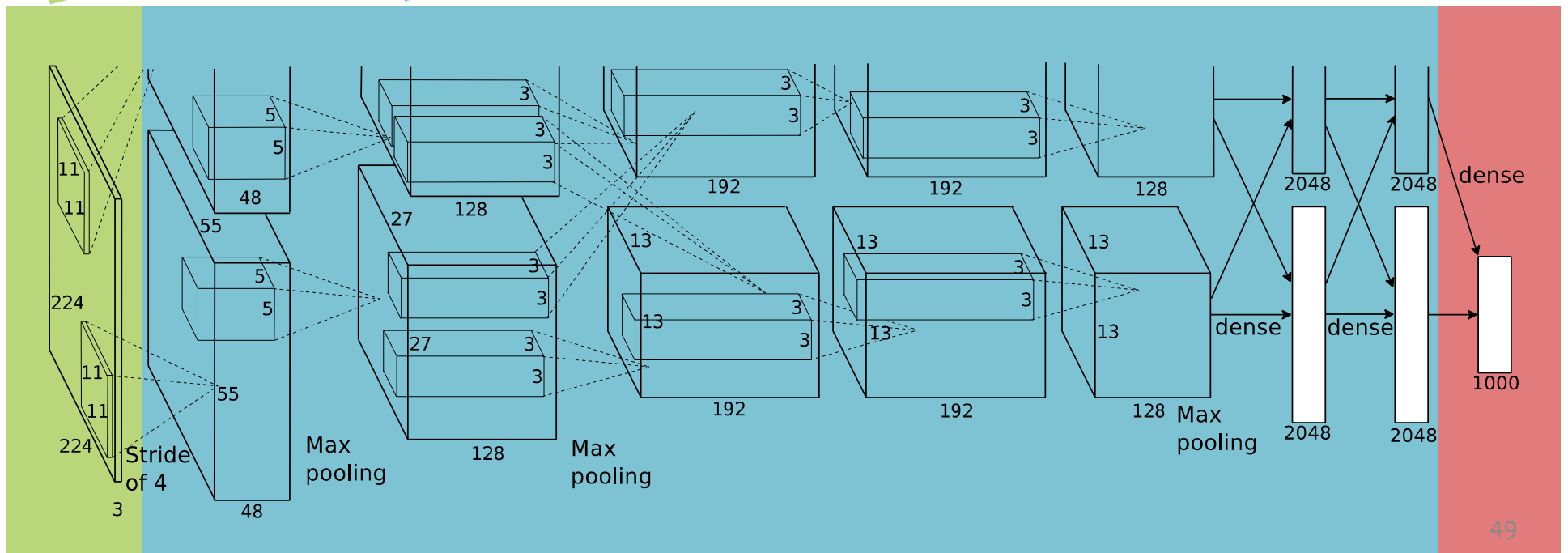
(Krizhevsky, Sutskever & Hinton, 2012)

15.3% error on ImageNet LSVRC-2012 contest

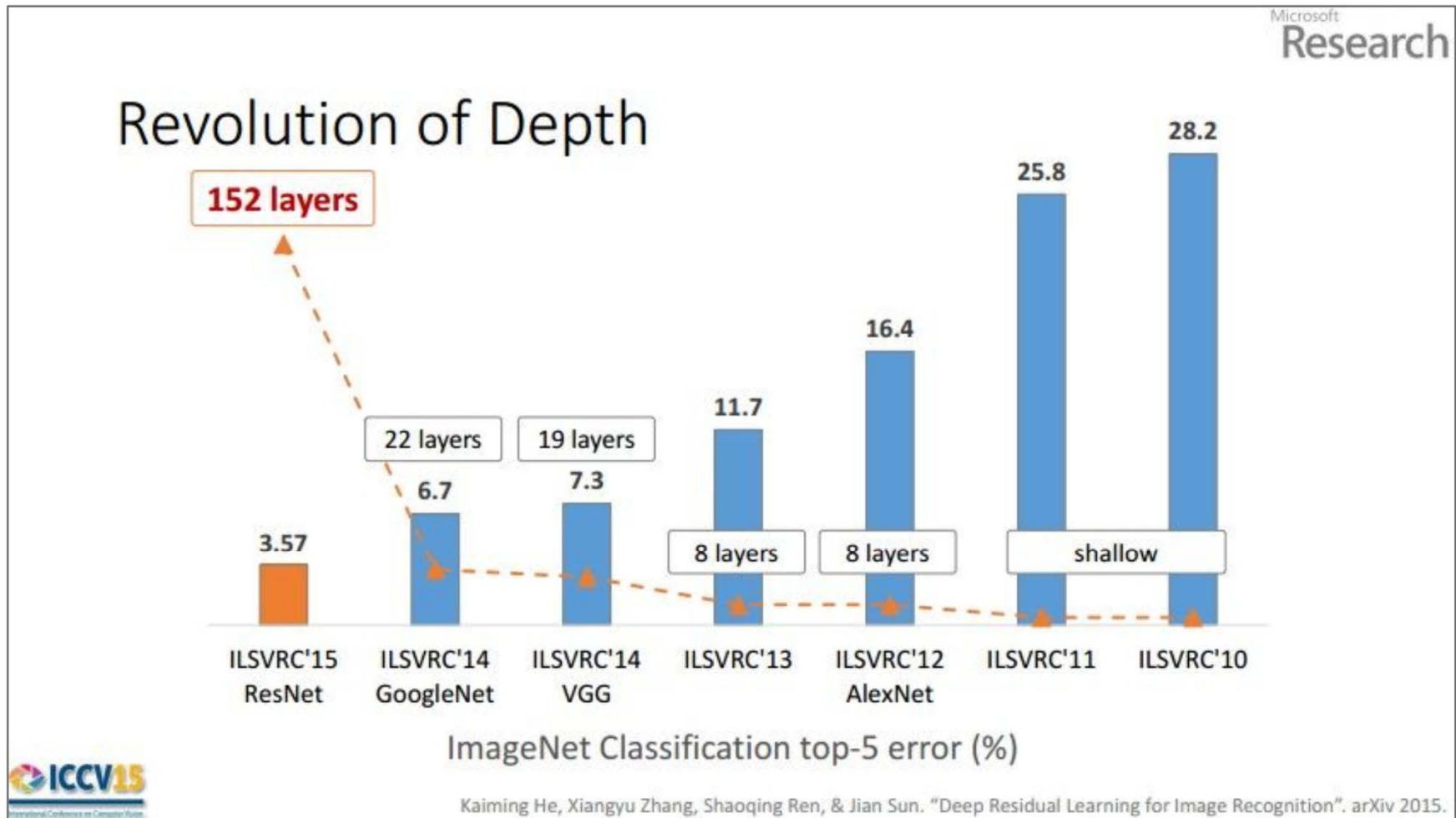
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



# CNNs for Image Recognition

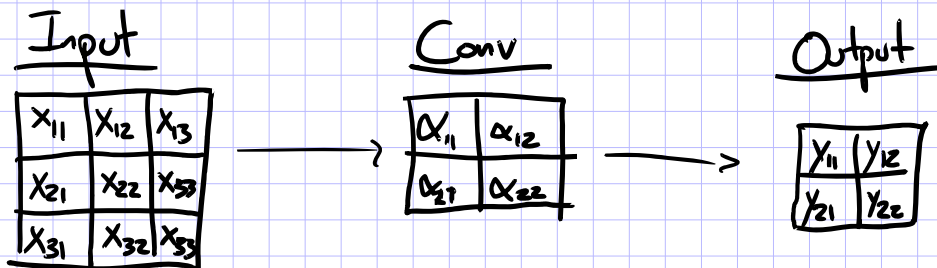


# CONVOLUTION

# What's a convolution?

- Basic idea:
  - Pick a 3x3 matrix F of weights
  - Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation
- Key point:
  - Different convolutions extract different types of low-level “features” from an image
  - All that we need to vary to generate these different features is the weights of F

Ex: 1 input channel, 1 output channel



$$\begin{aligned}y_{11} &= \alpha_{11} x_{11} + \alpha_{12} x_{12} + \alpha_{21} x_{21} + \alpha_{22} x_{22} + \alpha_0 \\y_{12} &= \alpha_{11} x_{12} + \alpha_{12} x_{13} + \alpha_{21} x_{22} + \alpha_{22} x_{23} + \alpha_0 \\y_{21} &= \alpha_{11} x_{21} + \alpha_{12} x_{22} + \alpha_{21} x_{31} + \alpha_{22} x_{32} + \alpha_0 \\y_{22} &= \alpha_{11} x_{22} + \alpha_{12} x_{23} + \alpha_{21} x_{32} + \alpha_{22} x_{33} + \alpha_0\end{aligned}$$

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	1
0	1	0

Convolved Image

1	1	1	1	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	0	0	0	0

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	1
0	1	0

Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

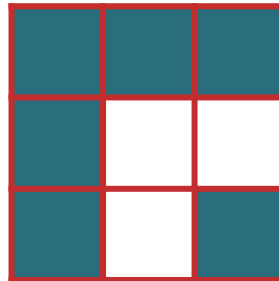
# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0



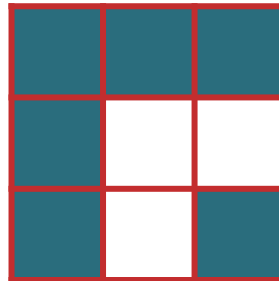
# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

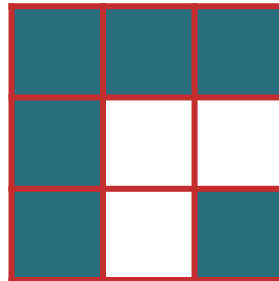
# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

			0	0	0	0
	1	1	1	1	1	0
	1		0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution


Convolved Image

3				

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	0	1	1	1	1	0
0	0	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	1
0	1	0

Convolved Image

3	2	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	1	1	1	0	0
0	1	1	1	1	1	0
0	1	1	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

1	1	1
1	0	0
1	0	1

Convolved Image

3	2	2	0	0
0	1	1	0	1
0	1	0	1	1
0	0	1	1	1
0	1	1	1	1

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	1	1	1	0
0	1	1	1	1	1	0
0	1	0	1	1	1	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

1	1	1
1	0	0
1	0	1

Convolved Image

3	2	2	3	
	1	1		1
	1		1	1
		1	1	1
	1	1	1	1

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	1	1	1
0	1	1	1	1	1	0
0	1	0	0	1	0	1
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

1	1	1
1	0	0
1	0	1

Convolved Image

3	2	2	3	1
0	1	1	0	1
0	1	0	1	1
0	0	1	1	1
0	1	1	1	1

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	0	0	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	1	0
0	1	0

Convolved Image

3	2	2	3	1
2	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0



# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	0	0	0	1	1	0
0	0	0	0	1	0	0
0	0	0	0	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution

0	0	0
0	0	0
0	0	0

Convolved Image

3	2	2	3	1
2	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

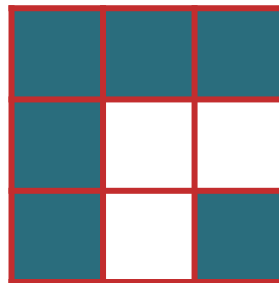
# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Convolution



Convolved Image

3	2	2	3	1
2	0	2	1	0
2	2	1	0	0
3	1	0	0	0
1	0	0	0	0

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Identity  
Convolution

0	0	0
0	1	0
0	0	0

Convolved Image

1	1	1	1	1
1	0	0	1	0
1	0	1	0	0
1	1	0	0	0
1	0	0	0	0

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0

Blurring  
Convolution

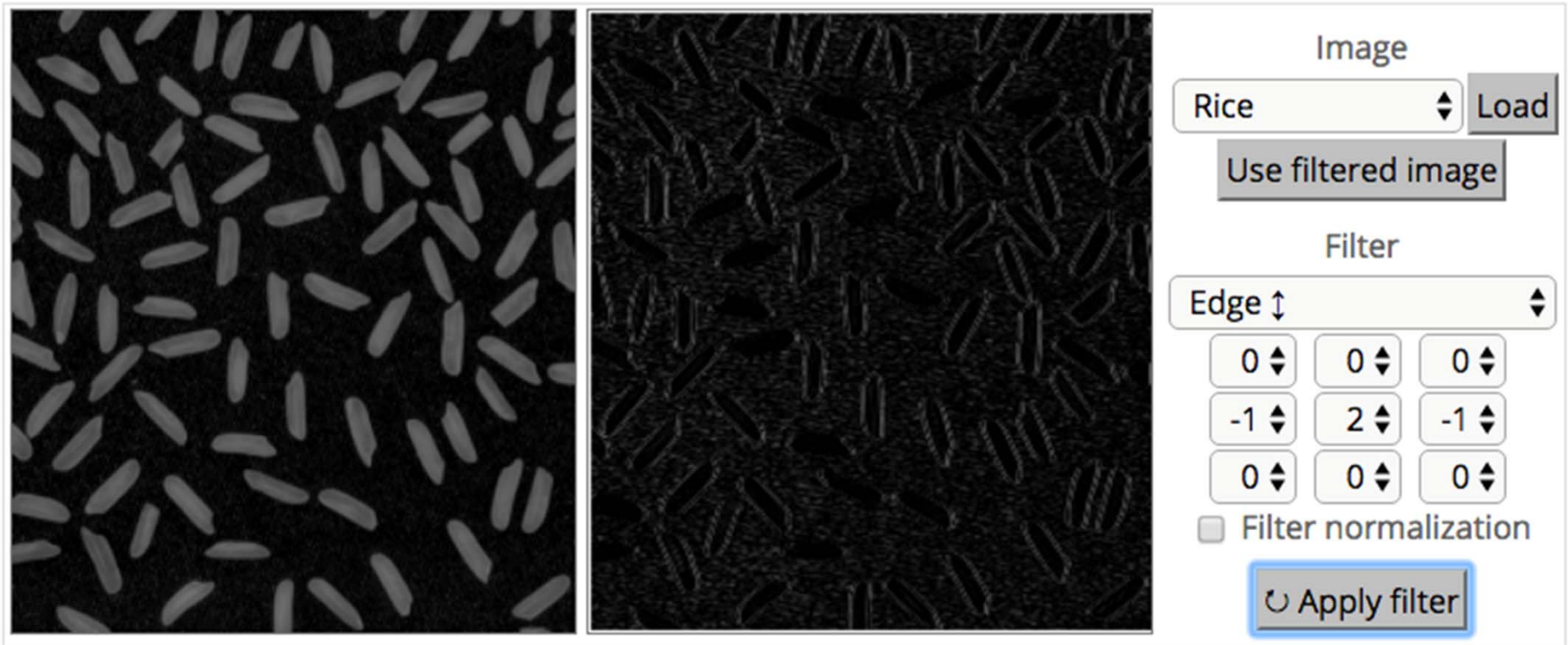
.1	.1	.1
.1	.2	.1
.1	.1	.1

Convolved Image

.4	.5	.5	.5	.4
.4	.2	.3	.6	.3
.5	.4	.4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0

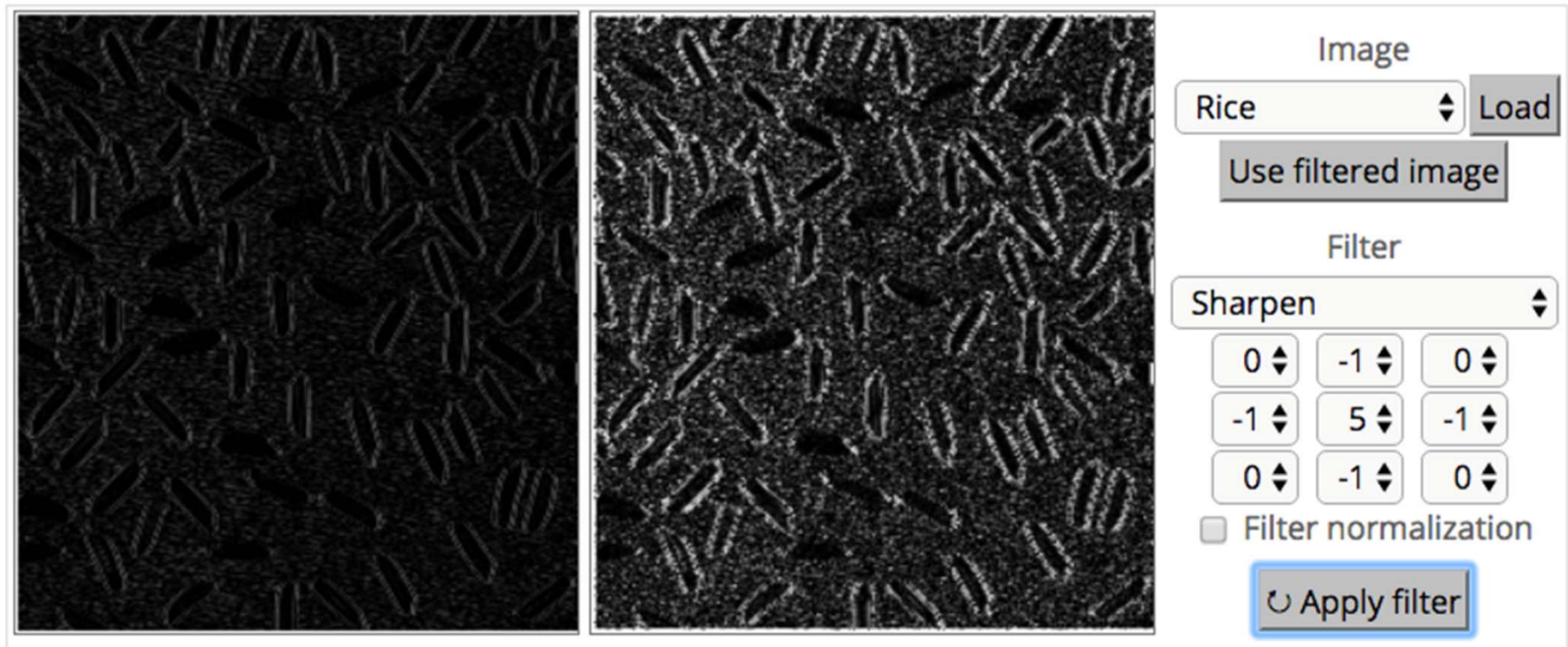
# What's a convolution?

<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>



# What's a convolution?

<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>



# What's a convolution?

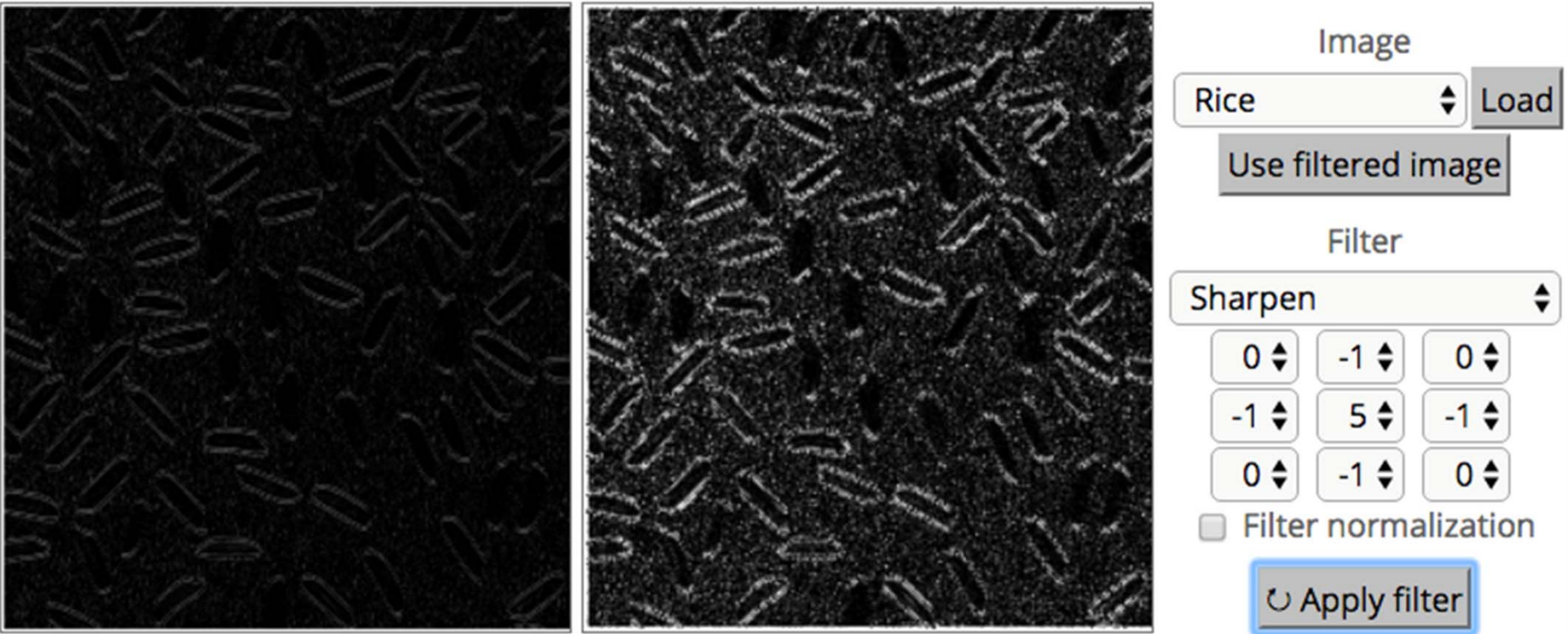
<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>





# What's a convolution?

<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>



The image shows a side-by-side comparison of a grayscale image of rice grains. The left image is the original, which is somewhat blurry. The right image is the result of applying a sharpening filter, showing significantly more detail and contrast. To the right of the images is a control panel for the online demo.

**Image**

Rice

**Filter**

Sharpen

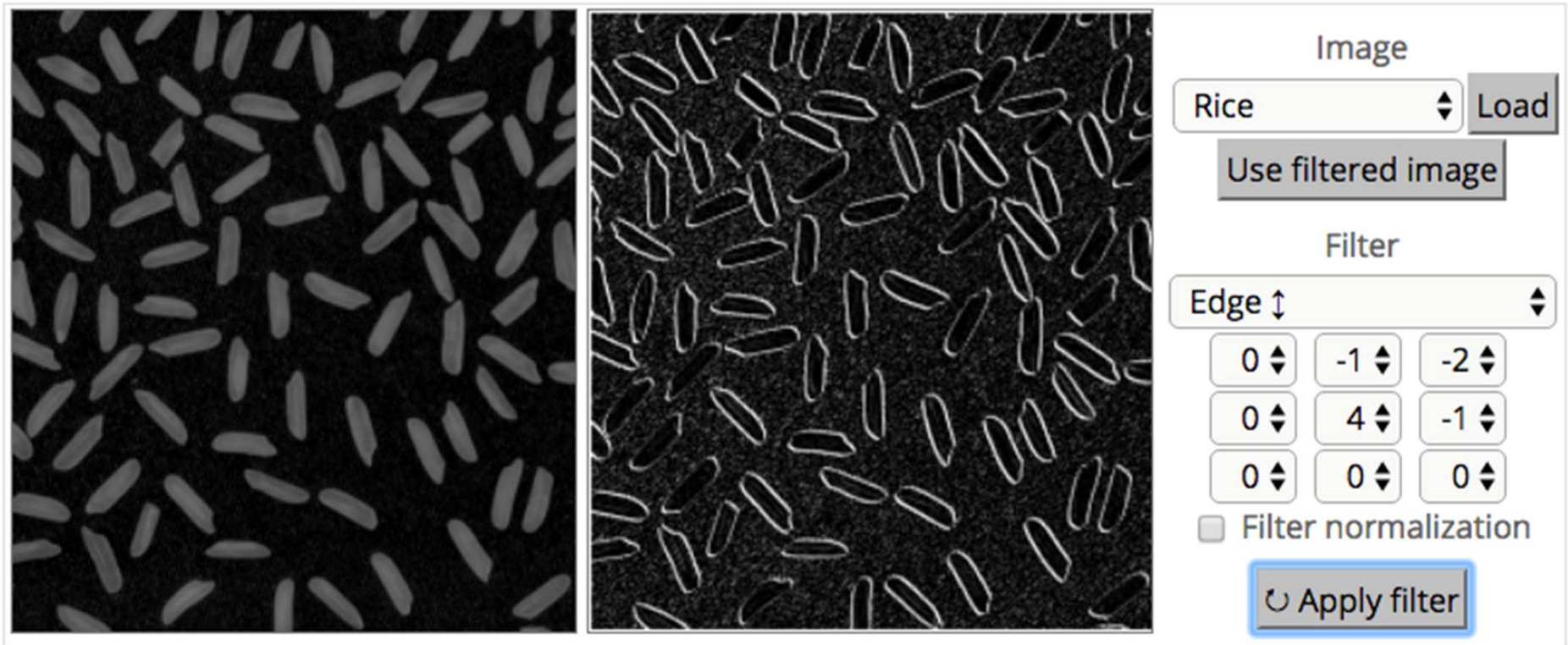
0 <input type="button" value="↓"/>	-1 <input type="button" value="↓"/>	0 <input type="button" value="↓"/>
-1 <input type="button" value="↓"/>	5 <input type="button" value="↓"/>	-1 <input type="button" value="↓"/>
0 <input type="button" value="↓"/>	-1 <input type="button" value="↓"/>	0 <input type="button" value="↓"/>

Filter normalization



# What's a convolution?

<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>



# What's a convolution?

<http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo>

Image

Rice Load

Use filtered image

Filter

Edge

-2	-1	0
-1	4	0
0	0	0

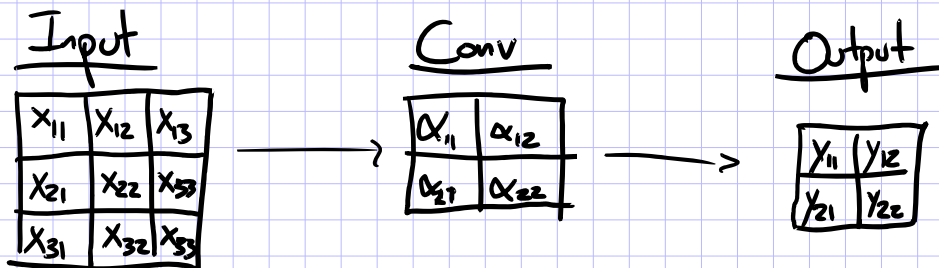
Filter normalization

Apply filter

# What's a convolution?

- Basic idea:
  - Pick a 3x3 matrix F of weights
  - Slide this over an image and compute the “inner product” (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation
- Key point:
  - Different convolutions extract different types of low-level “features” from an image
  - All that we need to vary to generate these different features is the weights of F

Ex: 1 input channel, 1 output channel



$$\begin{aligned}y_{11} &= \alpha_{11}x_{11} + \alpha_{12}x_{12} + \alpha_{21}x_{21} + \alpha_{22}x_{22} + \alpha_0 \\y_{12} &= \alpha_{11}x_{12} + \alpha_{12}x_{13} + \alpha_{21}x_{22} + \alpha_{22}x_{23} + \alpha_0 \\y_{21} &= \alpha_{11}x_{21} + \alpha_{12}x_{22} + \alpha_{21}x_{31} + \alpha_{22}x_{32} + \alpha_0 \\y_{22} &= \alpha_{11}x_{22} + \alpha_{12}x_{23} + \alpha_{21}x_{32} + \alpha_{22}x_{33} + \alpha_0\end{aligned}$$

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image


# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3		

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3		



# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1		

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1	0	

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

1	1
1	1

Convolved Image

3	3	1
3	1	0
1	0	0

# CONVOLUTIONAL NEURAL NETS

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

## Background

# A Recipe for Machine Learning

1. Convolutional Neural Networks (CNNs) provide another form of **decision function**
  - Let's see what they look like...

2. Choose each of these.

- Decision function

$$\hat{y} = f_{\theta}(\mathbf{x}_i)$$

- Loss function

$$\ell(\hat{y}, \mathbf{y}_i) \in \mathbb{R}$$

4. Train with SGD:

(Take small steps opposite the gradient)

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$



# Convolutional Neural Network (CNN)

- Typical layers include:
  - Convolutional layer
  - Max-pooling layer
  - Fully-connected (Linear) layer
  - ReLU layer (or some other nonlinear activation function)
  - Softmax
- These can be arranged into arbitrarily deep topologies

## Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

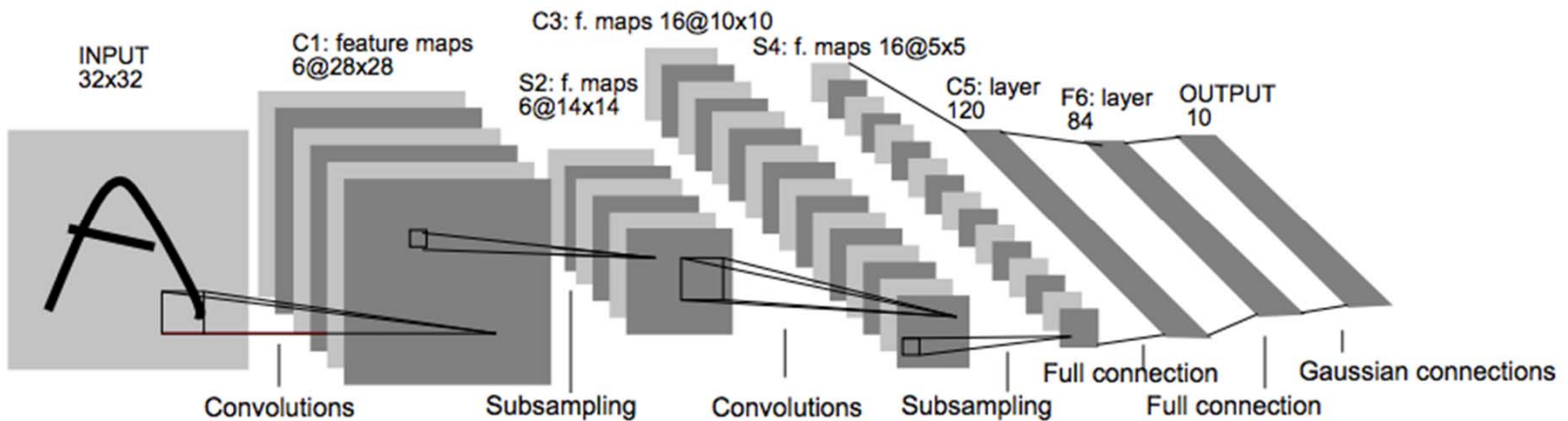


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

# Convolutional Layer

**CNN key idea:**  
Treat convolution matrix as parameters and learn them!

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0



Learned Convolution

$\theta_{11}$	$\theta_{12}$	$\theta_{13}$
$\theta_{21}$	$\theta_{22}$	$\theta_{23}$
$\theta_{31}$	$\theta_{32}$	$\theta_{33}$

Convolved Image

.4	.5	.5	.5	.4
.4	.2	.3	.6	.3
.5	.4	.4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0

# Downsampling by Averaging

- Downsampling by averaging **used to be** a common approach
- This is a special case of convolution where the weights are fixed to a uniform distribution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

$1/4$	$1/4$
$1/4$	$1/4$

Convolved Image

$3/4$	$3/4$	$1/4$
$3/4$	$1/4$	0
$1/4$	0	0

# Max-Pooling

- Max-pooling is another (common) form of downsampling
- Instead of averaging, we take the max value within the same range as the equivalently-sized convolution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Max-pooling

$x_{i,j}$	$x_{i,j+1}$
$x_{i+1,j}$	$x_{i+1,j+1}$

Max-Pooled Image

1	1	1
1	1	0
1	0	0

$$y_{ij} = \max(x_{ij}, x_{i,j+1}, x_{i+1,j}, x_{i+1,j+1})$$

# TRAINING CNNs

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

## Background

# A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of the following:

– Decision function

$$\hat{\mathbf{y}} = f_{\theta}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

- Q: Now that we have the CNN as a decision function, how do we compute the gradient?
- A: Backpropagation of course!

opposite the gradient)


$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

# SGD for CNNs

## SGD for CNNs

Ex: Architecture: Given  $\vec{x}, y^*$

$$J = \ell(y, y^*)$$

$$y = \text{softmax}(z^{(5)})$$

$$z^{(5)} = \text{linear}(z^{(4)}, W)$$

$$z^{(4)} = \text{relu}(z^{(3)})$$

$$z^{(3)} = \text{conv}(z^{(2)}, \beta)$$

$$z^{(2)} = \text{max-pool}(z^{(1)})$$

$$z^{(1)} = \text{conv}(\vec{x}, \alpha)$$

Parameters

$$\vec{\theta} = [\alpha, \beta, W]$$

SGD:

① Init  $\vec{\theta} \leftarrow$

② While not converged:

Sample  $i \in \{1, \dots, N\}$

Forward:  $y = h_{\theta}(\vec{x}^{(i)})$ ,  $J_i(\theta) = \ell(y, y^*)$

Backward:  $\nabla_{\vec{\theta}} J_i(\theta) = \dots$

Update:  $\vec{\theta} \leftarrow \vec{\theta} - \lambda \nabla_{\vec{\theta}} J_i(\theta)$



# LAYERS OF A CNN

# Common CNN Layers

## *Whiteboard*

- ReLU Layer
- Background: Subgradient
- Fully-connected Layer (w/tensor input)
- Softmax Layer
- Convolutional Layer
- Max-Pooling Layer

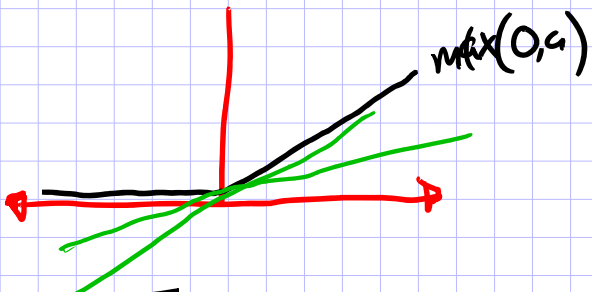
# ReLU Layer

ReLU Layer Input:  $\vec{x} \in \mathbb{R}^k$  Output:  $\vec{y} \in \mathbb{R}^k$

Forward:

$$\vec{y} = \sigma(\vec{x}) \leftarrow \text{element-wise}$$

$$\sigma(a) = \max(0, a)$$



Backward:

$$\frac{dJ}{dx_i} = \frac{dJ}{dy_i} \frac{dy_i}{dx_i}$$

where  $\frac{dy_i}{dx_i} = \begin{cases} 1 & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$

subderivative

# Softmax Layer

Softmax Layer

Input:  $\vec{x} \in \mathbb{R}^k$  Output:  $\vec{y} \in \mathbb{R}^k$

Forward:

$$y_i = \frac{\exp(x_i)}{\sum_{k=1}^k \exp(x_k)}$$

Backward:

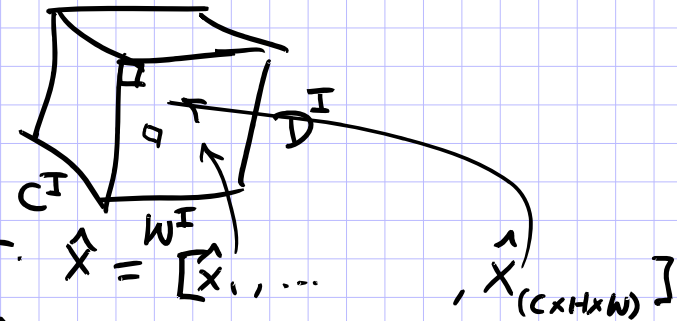
$$\frac{dJ}{dx_j} = \sum_{i=1}^k \frac{dJ}{dy_i} \frac{dy_i}{dx_j}$$

$$\text{where } \frac{dy_i}{dx_j} = \begin{cases} y_i(1-y_i) & \text{if } i=j \\ -y_i y_j & \text{otherwise} \end{cases}$$

# Fully-Connected Layer

Fully Connected Layer (w/ tensor input)

- Suppose input is a 3D Tensor:  $X =$



- stretch out into a long vector.  $\hat{x} = [x_1, \dots$

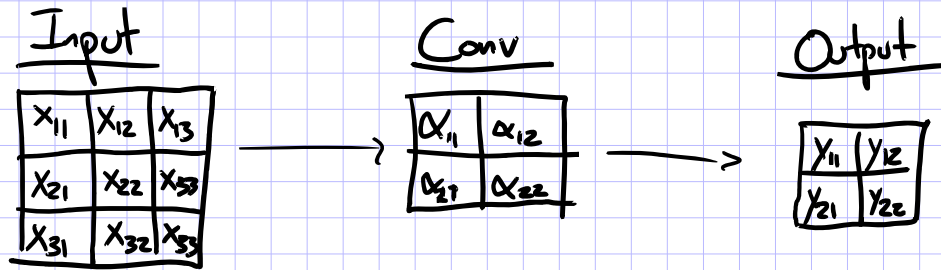
- then standard linear layer:

$$y = \alpha^T \hat{x} + \alpha_0 \quad \text{where } \alpha \in \mathbb{R}^{A \times B}$$

$|\hat{x}| = A, |y| = B$

# Convolutional Layer

Ex: 1 input channel, 1 output channel



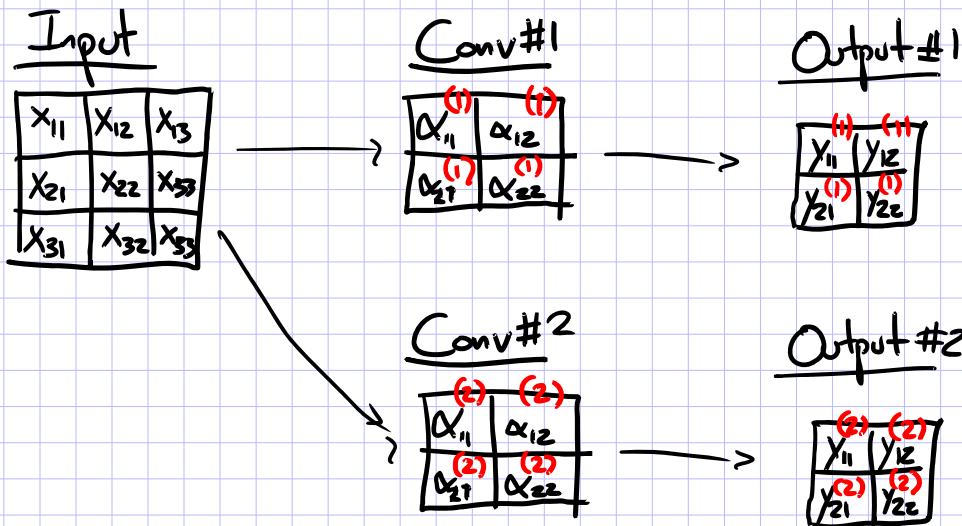
$$y_{11} = \alpha_{11}x_{11} + \alpha_{12}x_{12} + \alpha_{21}x_{21} + \alpha_{22}x_{22} + \alpha_0$$

$$y_{12} = \alpha_{11}x_{12} + \alpha_{12}x_{13} + \alpha_{21}x_{22} + \alpha_{22}x_{23} + \alpha_0$$

$$y_{21} = \alpha_{11}x_{21} + \alpha_{12}x_{22} + \alpha_{21}x_{31} + \alpha_{22}x_{32} + \alpha_0$$

$$y_{22} = \alpha_{11}x_{22} + \alpha_{12}x_{23} + \alpha_{21}x_{32} + \alpha_{22}x_{33} + \alpha_0$$

Ex: 1 input channel, 2 output channels



$$y_{11}^{(1)} = \alpha_{11}^{(1)}x_{11} + \alpha_{12}^{(1)}x_{12} + \alpha_{21}^{(1)}x_{21} + \alpha_{22}^{(1)}x_{22} + \alpha_0^{(1)}$$

$$y_{12}^{(1)} = \dots$$

$$y_{21}^{(1)} = \dots$$

$$y_{22}^{(1)} = \alpha_{11}^{(1)}x_{22} + \alpha_{12}^{(1)}x_{23} + \alpha_{21}^{(1)}x_{32} + \alpha_{22}^{(1)}x_{33} + \alpha_0^{(1)}$$

$$y_{11}^{(2)} = \alpha_{11}^{(2)}x_{11} + \alpha_{12}^{(2)}x_{12} + \alpha_{21}^{(2)}x_{21} + \alpha_{22}^{(2)}x_{22} + \alpha_0^{(2)}$$

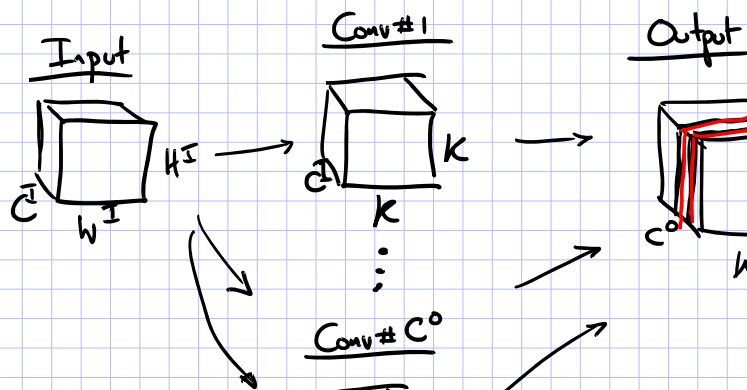
$$y_{12}^{(2)} = \dots$$

$$y_{21}^{(2)} = \dots$$

$$y_{22}^{(2)} = \alpha_{11}^{(2)}x_{22} + \alpha_{12}^{(2)}x_{23} + \alpha_{21}^{(2)}x_{32} + \alpha_{22}^{(2)}x_{33} + \alpha_0^{(2)}$$

# Convolutional Layer

Ex:  $C^I$  input channels,  $C^O$  output channels

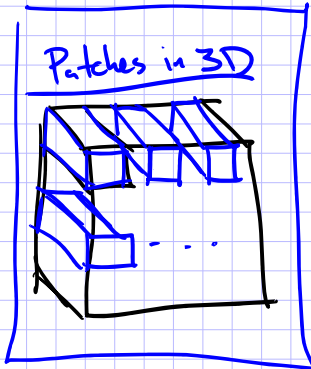


*j-th slice is output from j-th convolution matrix*

$$H^O = \lfloor (H^I + 2p - K) / s + 1 \rfloor$$

$$W^O = \lfloor (W^I + 2p - K) / s + 1 \rfloor$$

where  $p$  = # pixels of padding on input  
 $k$  = size of conv. matrix  
 $s$  = stride length



Forward:

$$y_{ij}^{(k)} = \alpha_0^{(k)} + \sum_{c=1}^{C^I} \sum_{q=1}^K \sum_{r=1}^K \alpha_{qr}^{(c)} x_{mn}^{(c)} \quad \text{where } m = s(i-1) + q, n = s(j-1) + r$$

Backward:

$$\frac{dJ}{d\alpha_0^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_0^{(k)}}$$

$$\frac{dJ}{d\alpha_{qr}^{(c)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_{qr}^{(c)}}$$

$$\frac{dJ}{dx_{mn}^{(c)}} = \sum_i \sum_j \sum_k \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(c)}}$$

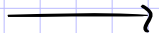
*just some calculus*

# Max-Pooling Layer

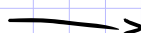
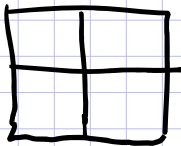
Ex: 1 input channel, 1 output channel, stride of 1

Input

$x_{11}$	$x_{12}$	$x_{13}$
$x_{21}$	$x_{22}$	$x_{23}$
$x_{31}$	$x_{32}$	$x_{33}$



Pool Size



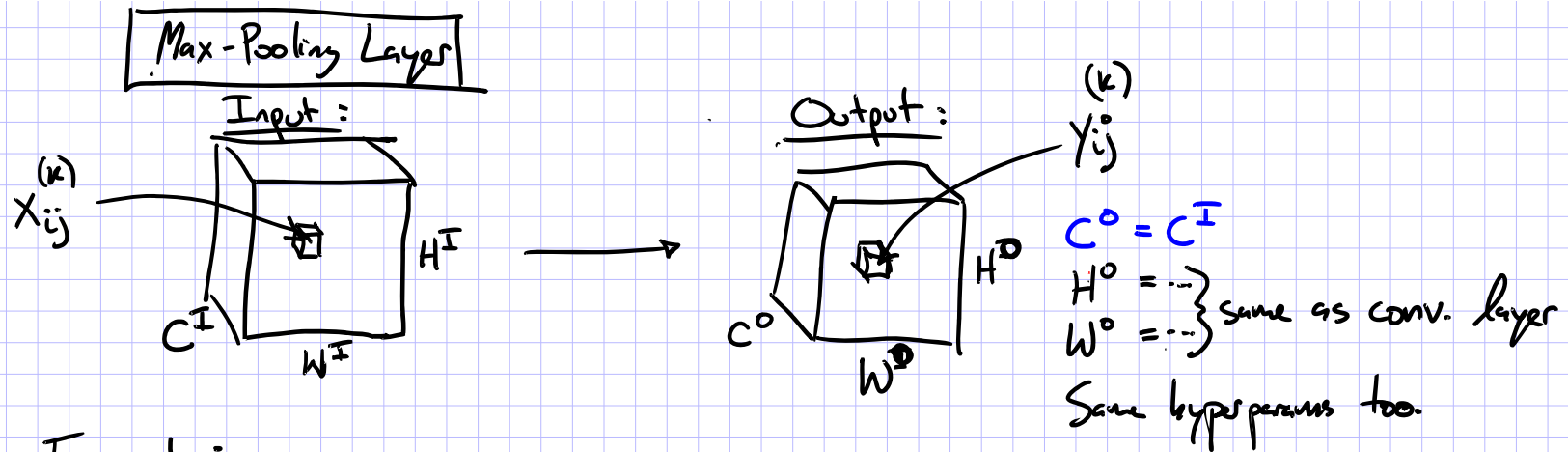
Output

$y_{11}$	$y_{12}$
$y_{21}$	$y_{22}$

$$y_{11} = \max(x_{11}, x_{12}, x_{21}, x_{22})$$
$$y_{12} = \max(x_{12}, x_{13}, x_{22}, x_{23})$$
$$y_{21} = \max(x_{21}, x_{22}, x_{31}, x_{32})$$
$$y_{22} = \max(x_{22}, x_{23}, x_{32}, x_{33})$$



# Max-Pooling Layer



Forward:

$$Y_{ij}^{(k)} = \max_{\substack{q \in \{1, \dots, k\} \\ r \in \{1, \dots, k\}}} X_{mn}^{(k)} \text{ where } \begin{cases} m = s(i-1) + q \\ n = s(j-1) + r \end{cases}$$

Backward:

$$\frac{dJ}{dx_{mn}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(k)}}$$

## Subderivatives

- +  $\max()$  is not differentiable, but subdifferentiable.
- + There are a set of derivatives and we can just choose one for SGD.

$$y = \max(a, b)$$

$$\Rightarrow \frac{dJ}{da} = \frac{dJ}{dy} \frac{dy}{da} \text{ where } \frac{dy}{da} = \begin{cases} 1 & \text{if } a > b \\ 0 & \text{otherwise} \end{cases}$$

# Convolutional Neural Network (CNN)

- Typical layers include:
  - Convolutional layer
  - Max-pooling layer
  - Fully-connected (Linear) layer
  - ReLU layer (or some other nonlinear activation function)
  - Softmax
- These can be arranged into arbitrarily deep topologies

## Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

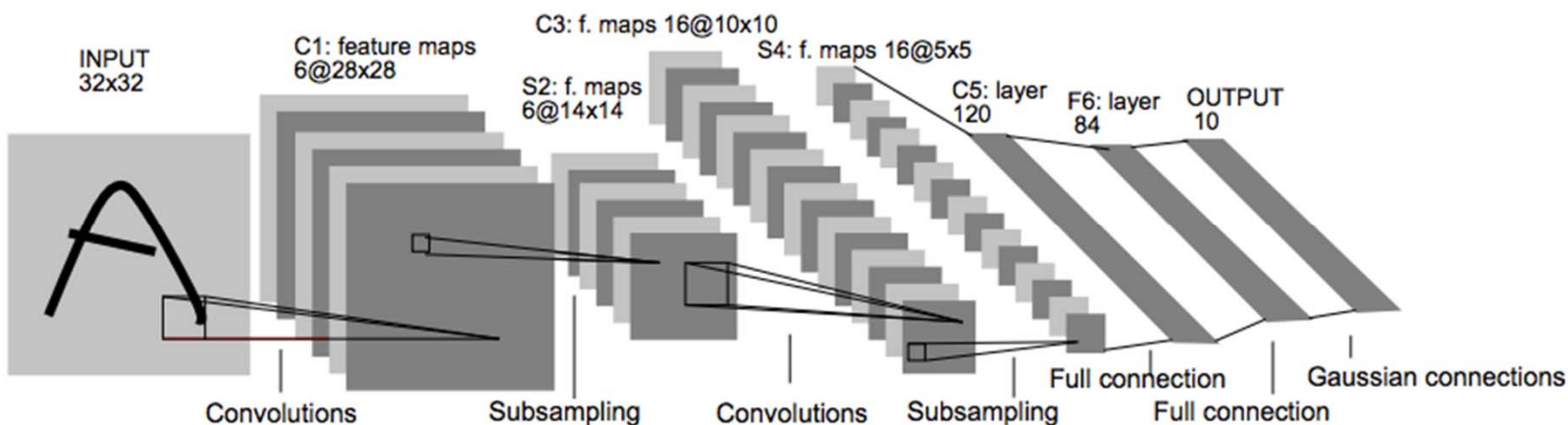


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

# Architecture #2: AlexNet

## CNN for Image Classification

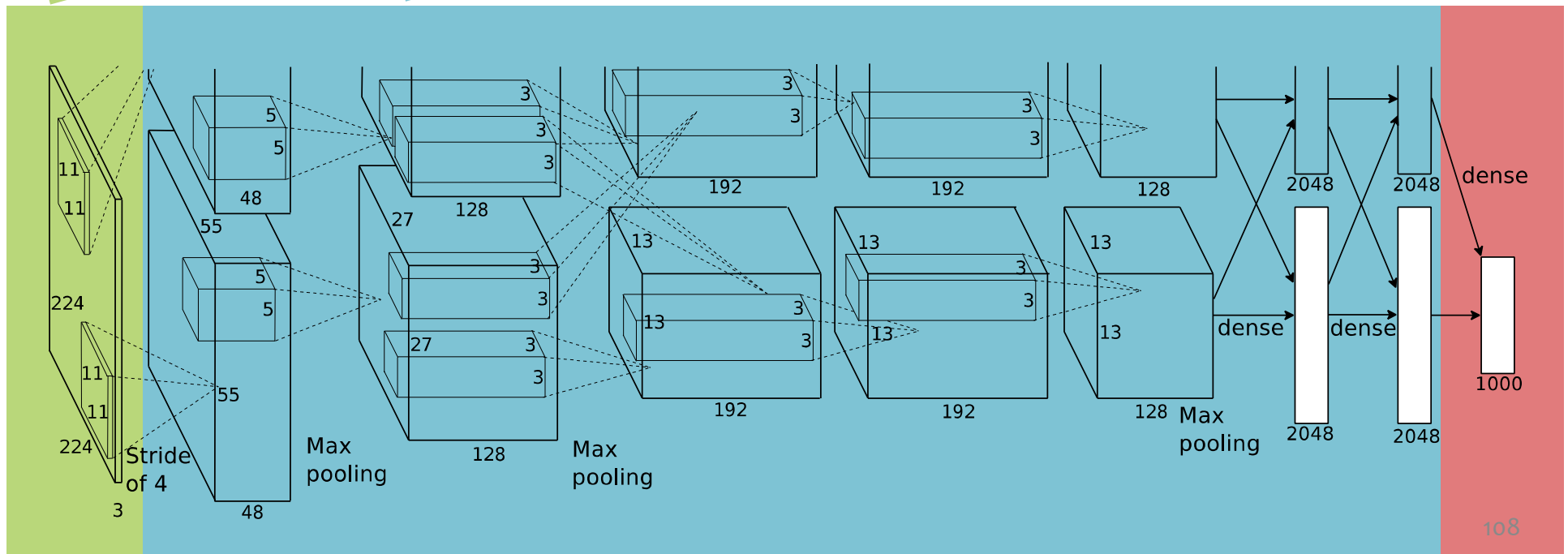
(Krizhevsky, Sutskever & Hinton, 2012)

15.3% error on ImageNet LSVRC-2012 contest

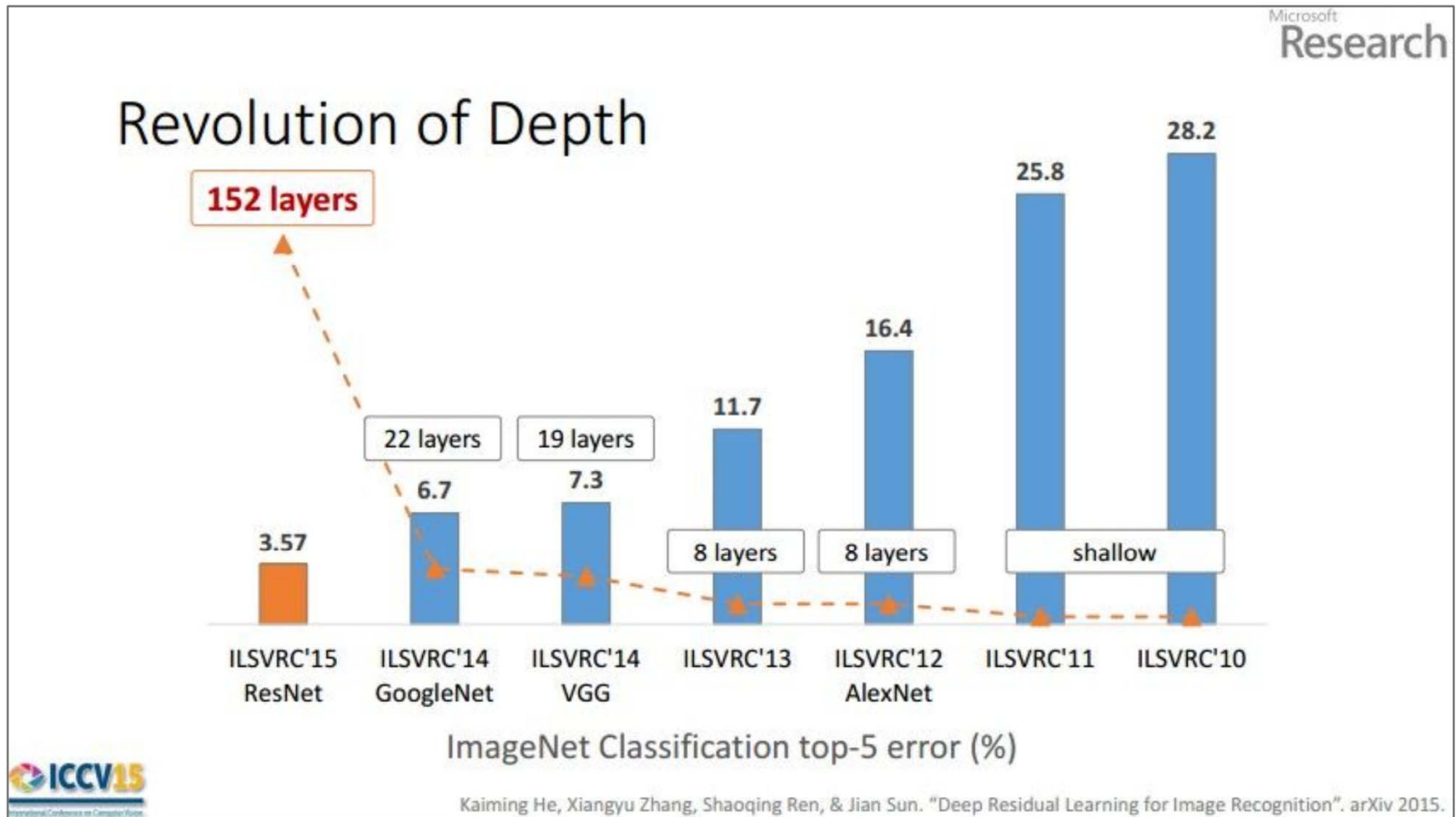
Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax



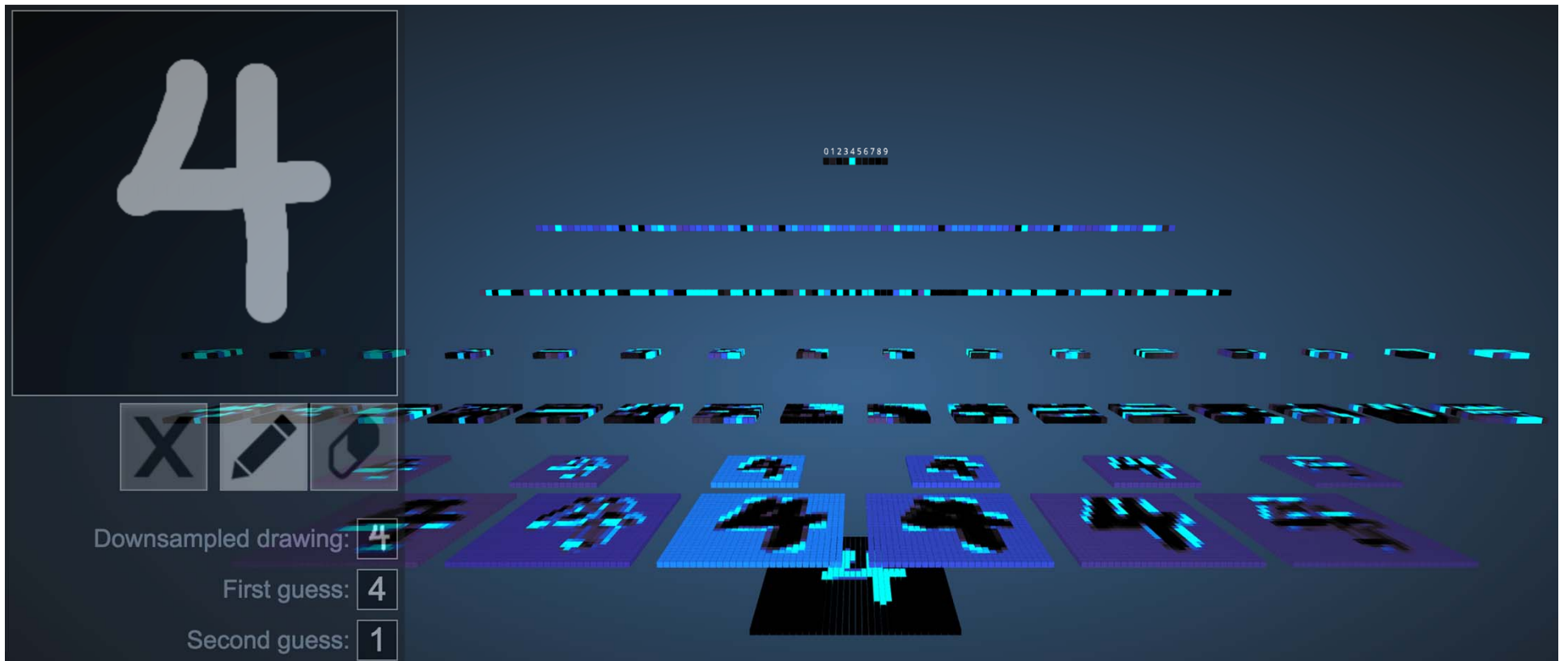
# CNNs for Image Recognition



# CNN VISUALIZATIONS

# 3D Visualization of CNN

<http://scs.ryerson.ca/~aharley/vis/conv/>



# Convolution of a Color Image

- Color images consist of 3 floats per pixel for RGB (red, green blue) color values
- Convolution must also be 3-dimensional

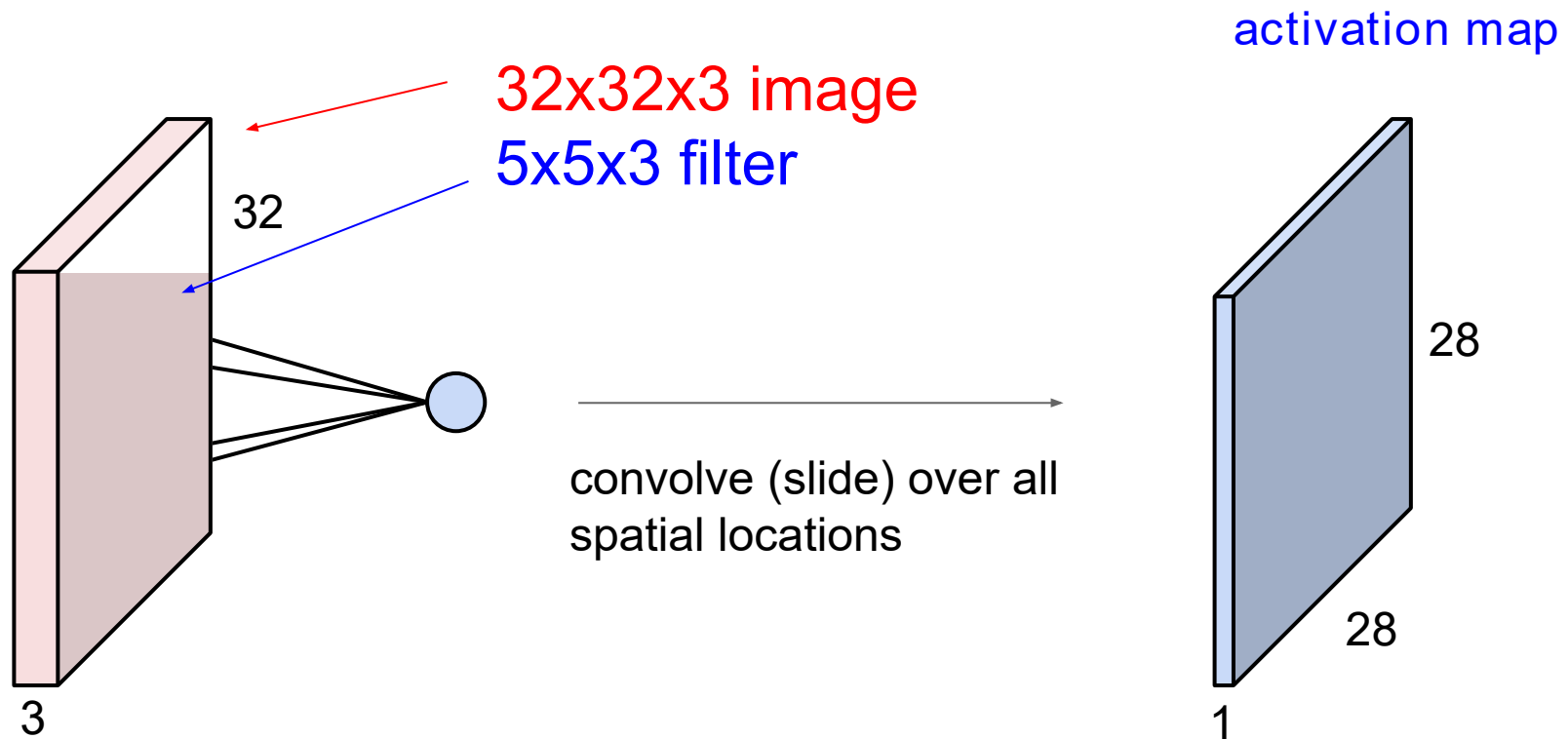


Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)

# Animation of 3D Convolution

<http://cs231n.github.io/convolutional-networks/>

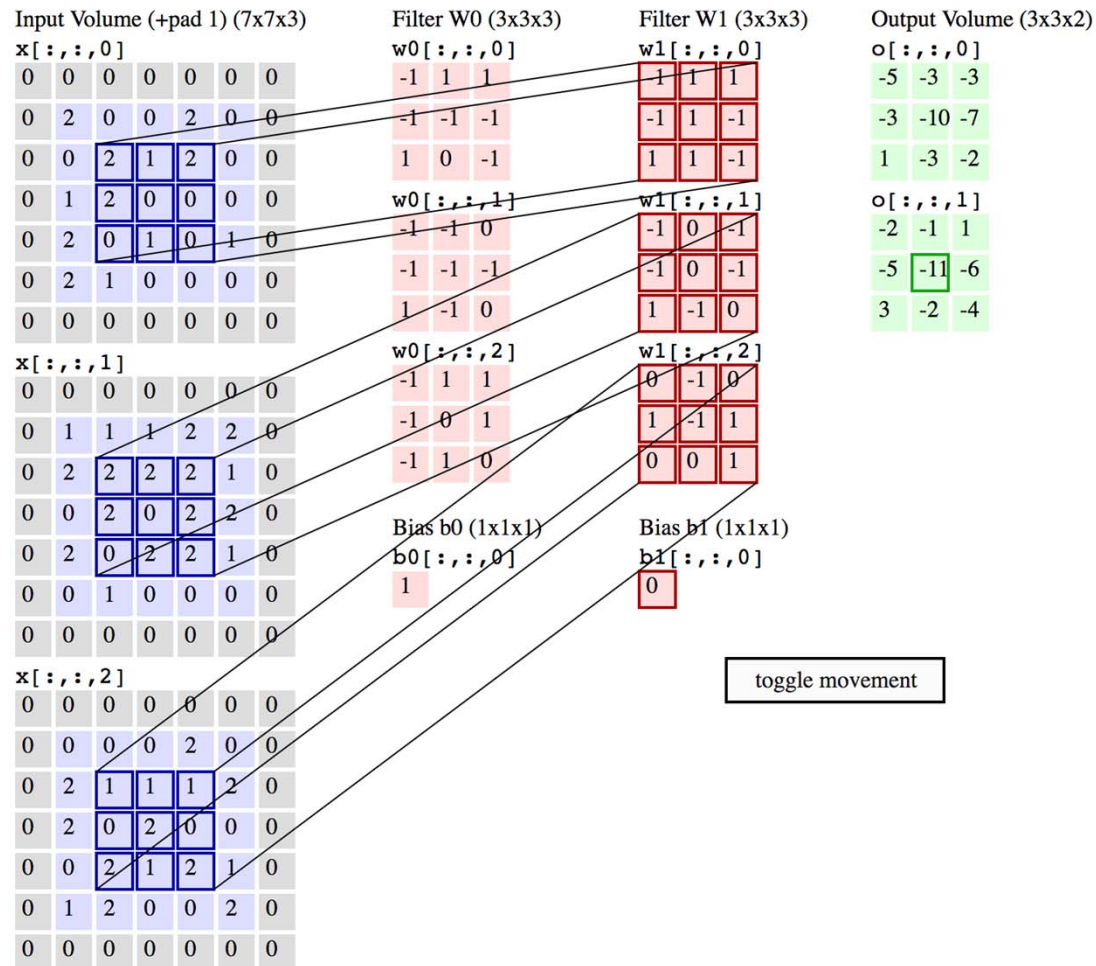


Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)



# MNIST Digit Recognition with CNNs (in your browser)

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

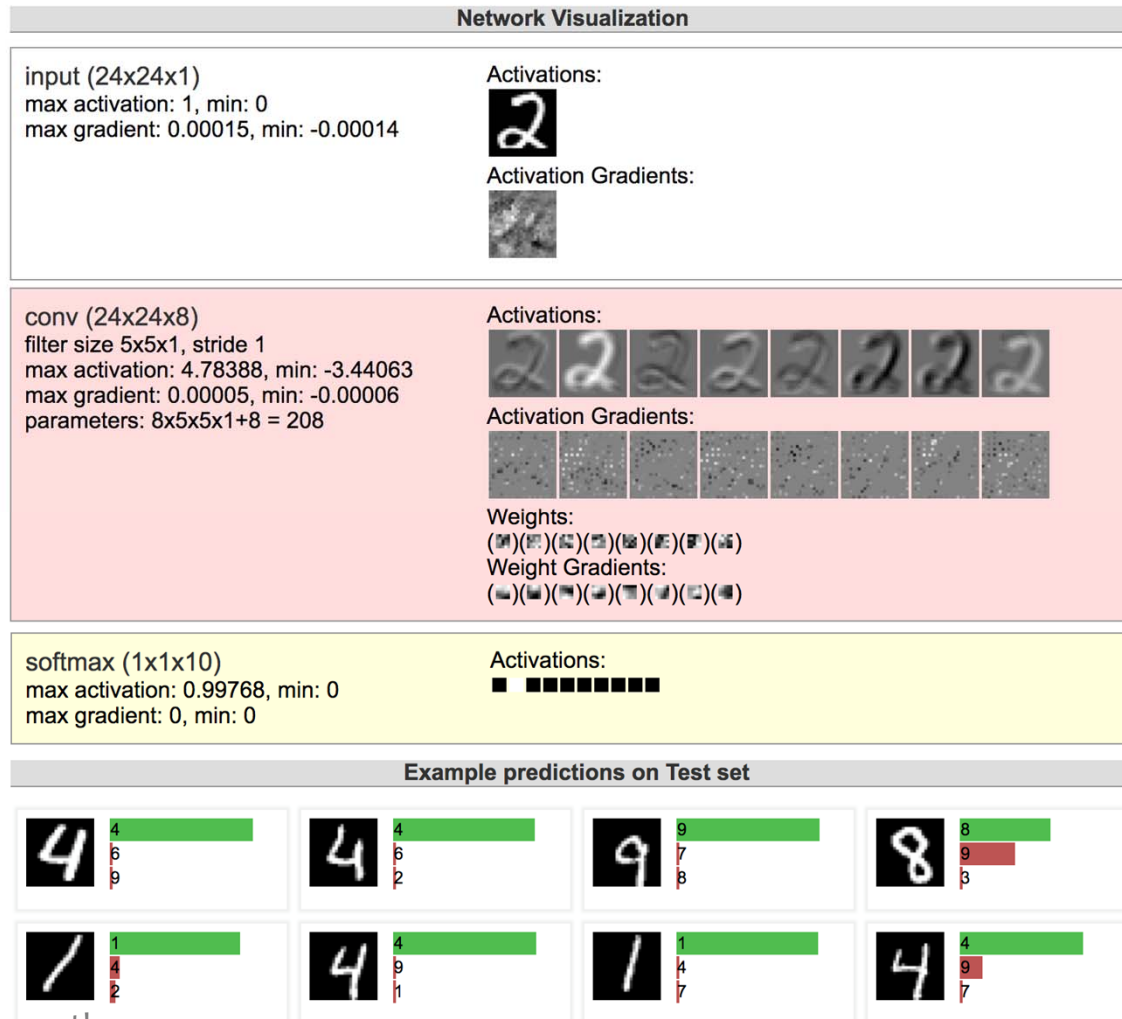


Figure from Andrej Karpathy

# CNN Summary

## CNNs

- Are used for all aspects of **computer vision**, and have won numerous pattern recognition competitions
- Able learn **interpretable features** at different levels of abstraction
- Typically, consist of **convolution** layers, **pooling** layers, **nonlinearities**, and **fully connected** layers

## Other Resources:

- Readings on course website
- Andrej Karpathy, CS231n Notes  
<http://cs231n.github.io/convolutional-networks/>

# **RECURRENT NEURAL NETWORKS**

# Dataset for Supervised Part-of-Speech (POS) Tagging

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

Sample 1:						} $y^{(1)}$
						} $x^{(1)}$
Sample 2:						} $y^{(2)}$
						} $x^{(2)}$
Sample 3:						} $y^{(3)}$
						} $x^{(3)}$
Sample 4:						} $y^{(4)}$
						} $x^{(4)}$

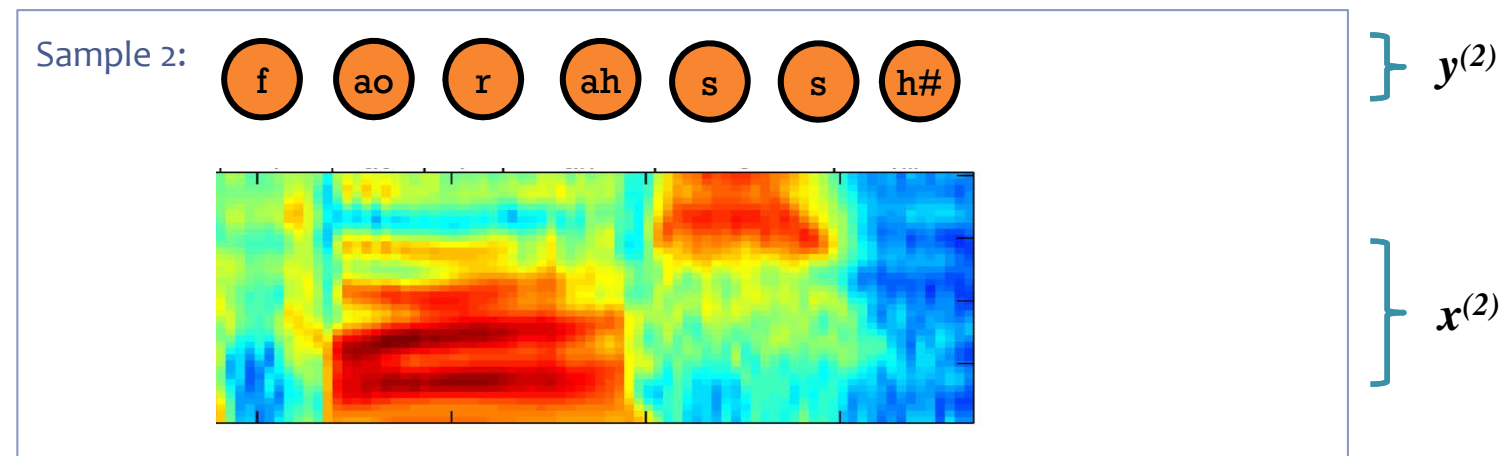
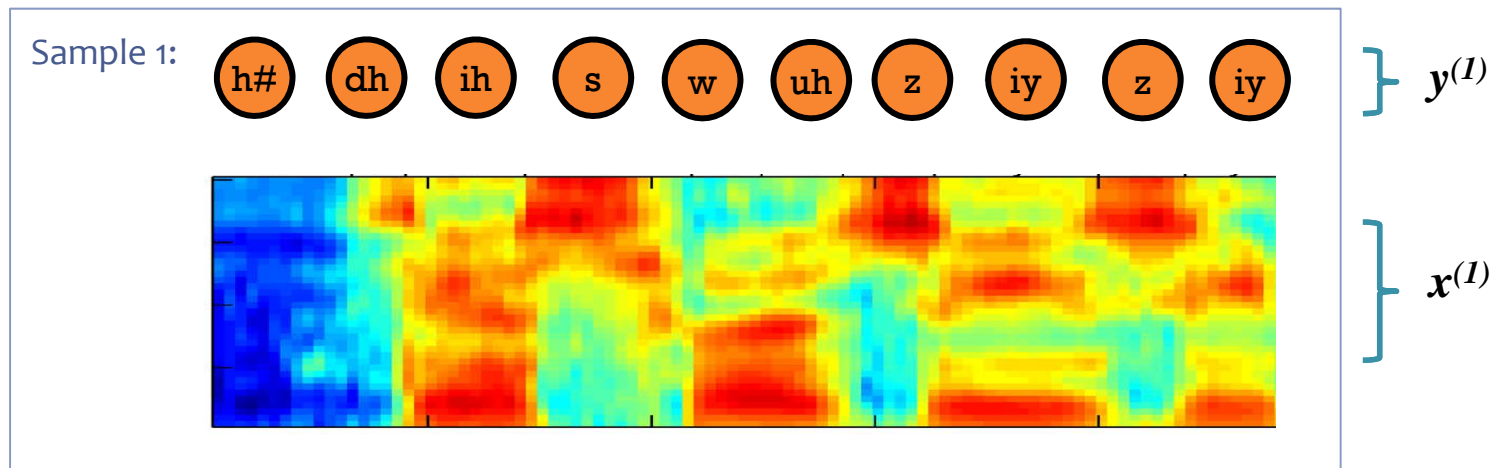
# Dataset for Supervised Handwriting Recognition

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$



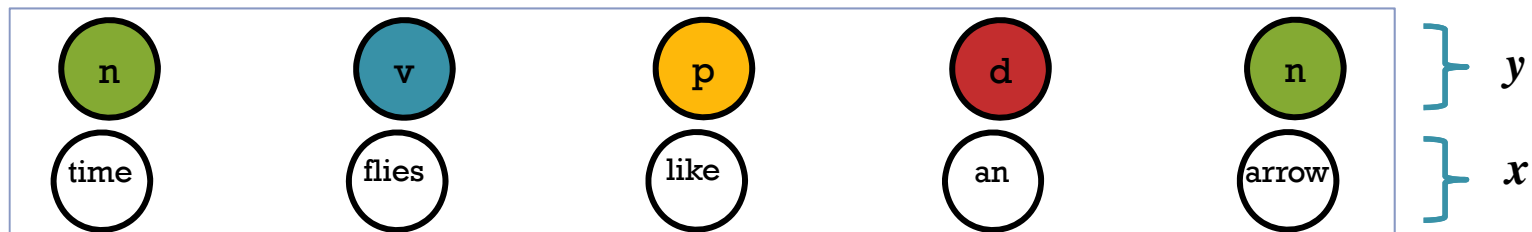
# Dataset for Supervised Phoneme (Speech) Recognition

Data:  $\mathcal{D} = \{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$



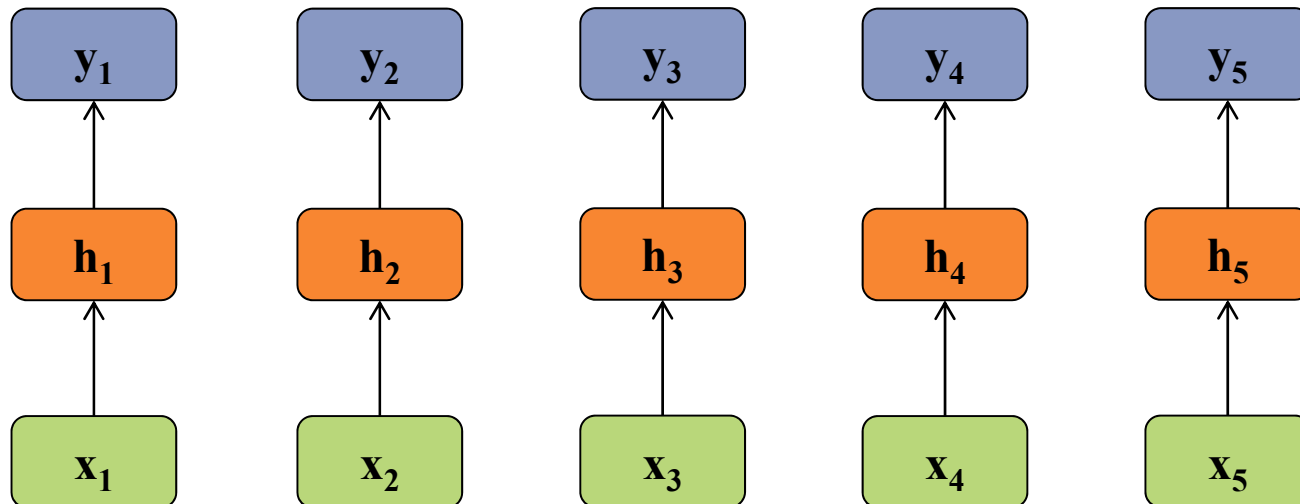
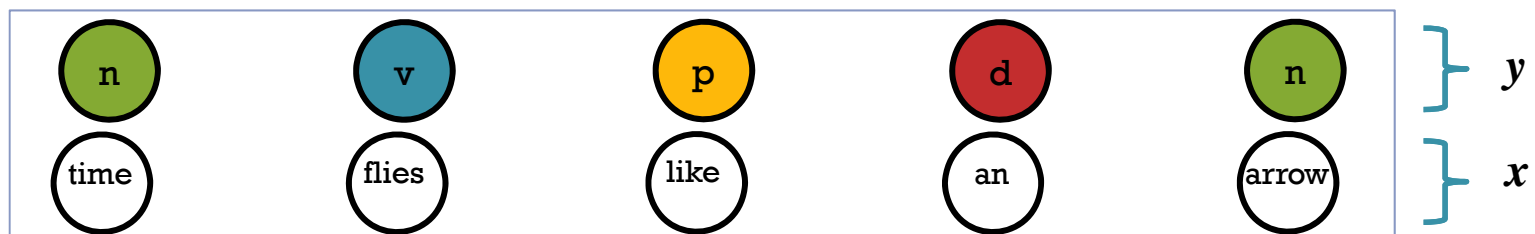
# Time Series Data

**Question 1:** How could we apply the neural networks we've seen so far (which expect **fixed size input/output**) to a prediction task with **variable length input/output**?



# Time Series Data

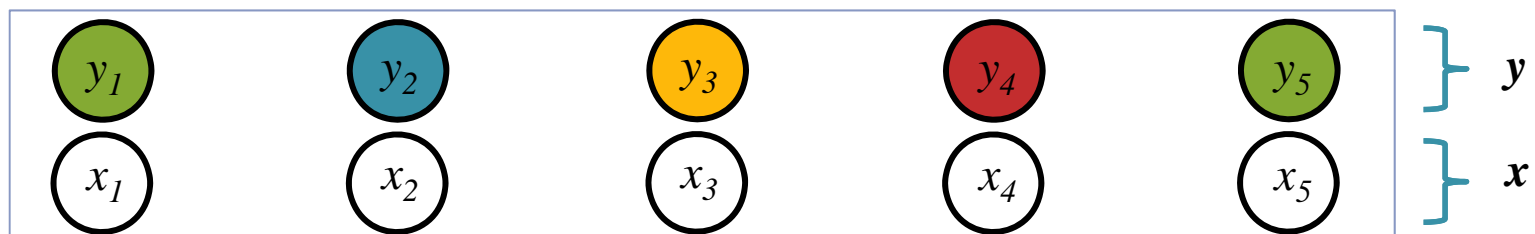
**Question 1:** How could we apply the neural networks we've seen so far (which expect **fixed size input/output**) to a prediction task with **variable length input/output**?





# Time Series Data

**Question 2:** How could we incorporate context (e.g. words to the left/right, or tags to the left/right) into our solution?



## Multiple Choice:

Working left-to-right, use features of...

	$x_{i-1}$	$x_i$	$x_{i+1}$	$y_{i-1}$	$y_i$	$y_{i+1}$
A	✓					
B				✓		
C	✓			✓		
D	✓			✓	✓	✓
E	✓	✓		✓	✓	✓
F	✓	✓	✓	✓		
G	✓	✓	✓	✓	✓	
H	✓	✓	✓	✓	✓	✓

# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

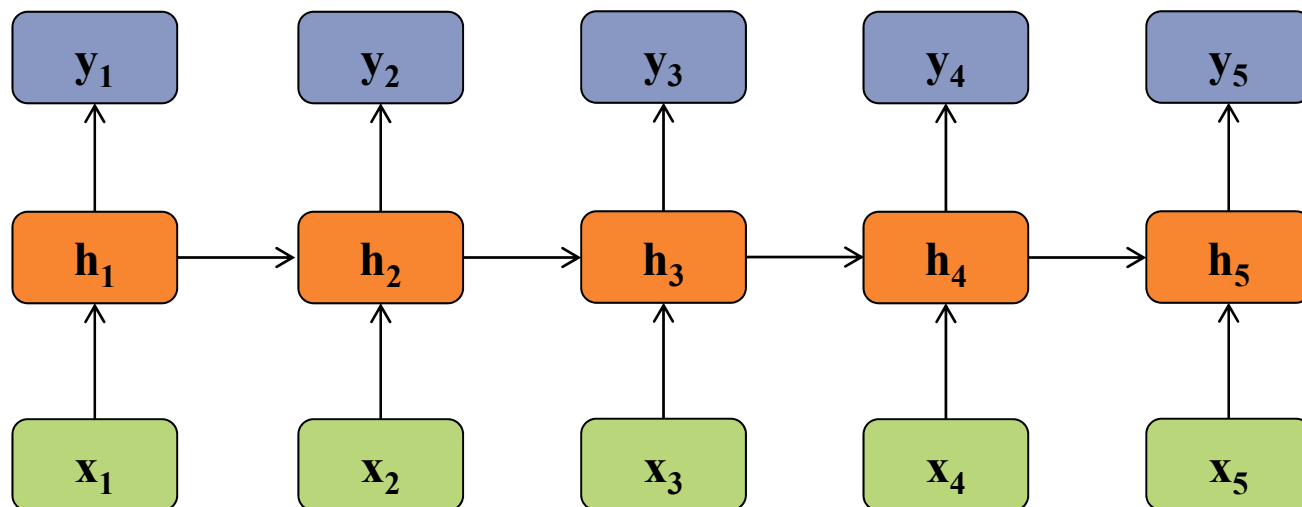
outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity:  $\mathcal{H}$

Definition of the RNN:

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

$$y_t = W_{hy}h_t + b_y$$



# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

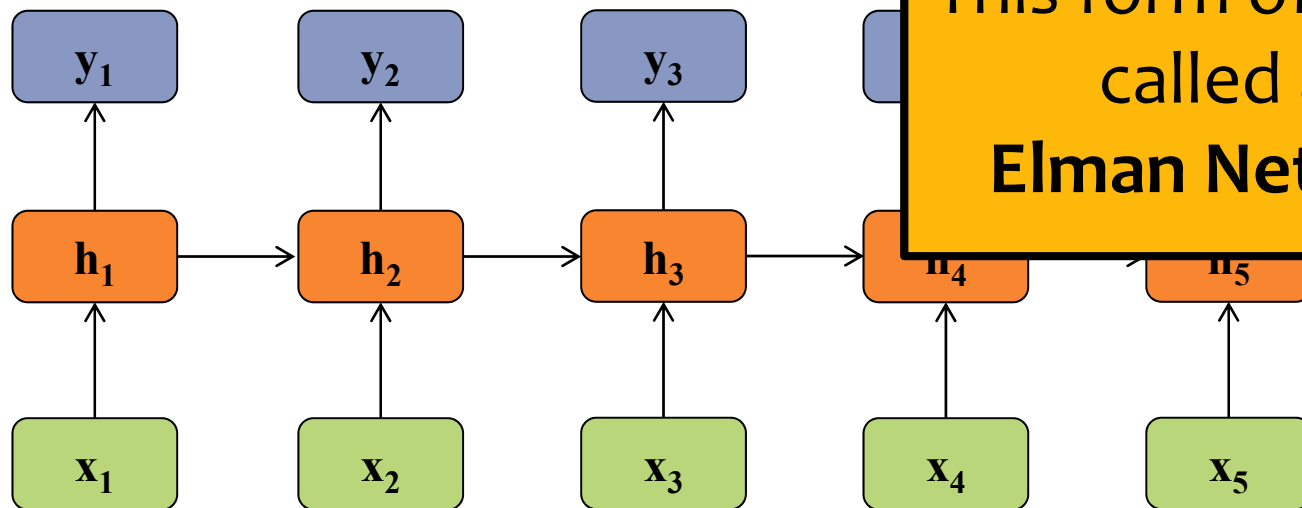
outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity:  $\mathcal{H}$

Definition of the RNN:

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

$$y_t = W_{hy}h_t + b_y$$



# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

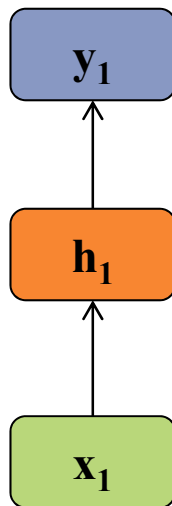
outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity:  $\mathcal{H}$

Definition of the RNN:

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

$$y_t = W_{hy}h_t + b_y$$



- If  $T=1$ , then we have a standard feed-forward neural net with one hidden layer
- All of the deep nets from last lecture required **fixed size inputs/outputs**

# Background A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

(take small steps  
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

# Background A Recipe for Machine Learning

1.
  - Recurrent Neural Networks (RNNs) provide another form of **decision function**
  - An RNN is just another differential function

2. CHOOSE EACH OF THESE:

- Decision function

$$\hat{y} = f_{\theta}(x_i)$$

4. Train with SGD:  
(take small steps opposite the gradient)

- We'll just need a method of computing the gradient efficiently
- Let's use Backpropagation Through Time...

$$-\eta_t \nabla \ell(f_{\theta}(x_i), y_i)$$

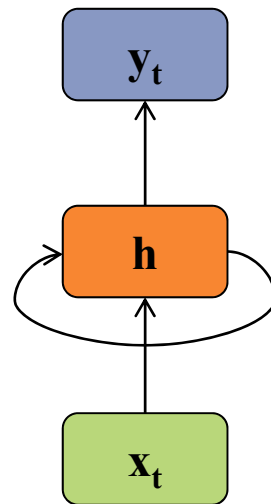
# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$   
hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$   
outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$   
nonlinearity:  $\mathcal{H}$

Definition of the RNN:

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

$$y_t = W_{hy}h_t + b_y$$



# Recurrent Neural Networks (RNNs)

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\mathbf{h} = (h_1, h_2, \dots, h_T), h_i \in \mathcal{R}^J$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

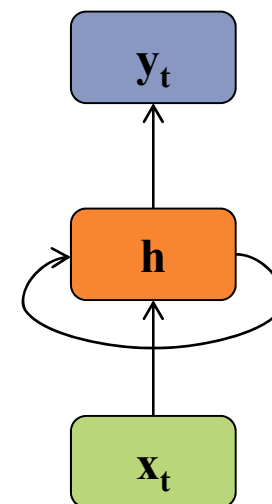
nonlinearity:  $\mathcal{H}$

Definition of the RNN:

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

$$y_t = W_{hy}h_t + b_y$$

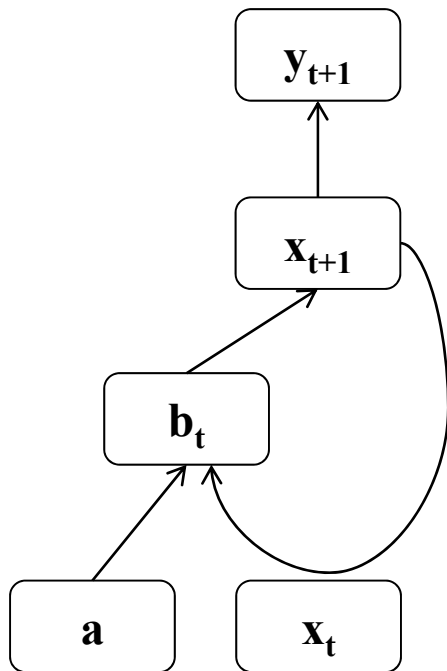
- By unrolling the RNN through time, we can **share parameters** and accommodate **arbitrary length** input/output pairs
- Applications: **time-series data** such as sentences, speech, stock-market, signal data, etc.





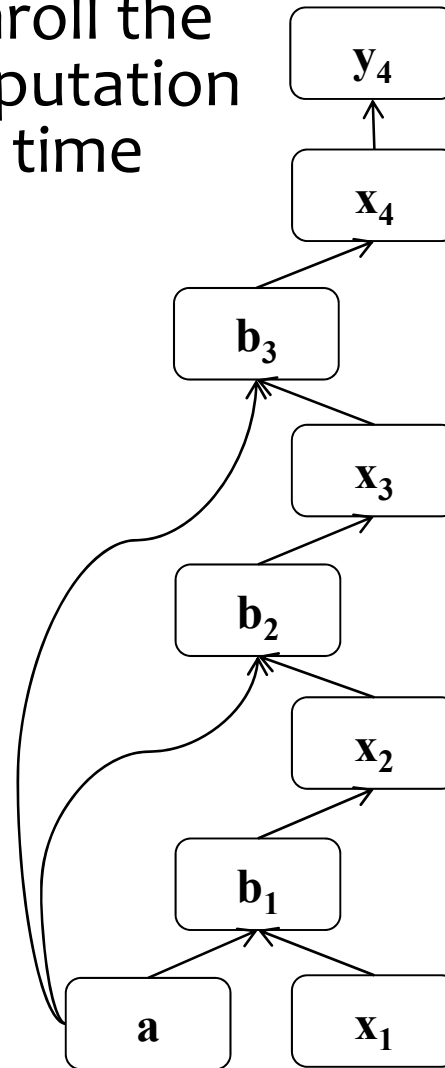
# Background: Backprop through time

## Recurrent neural network:



## BPTT:

1. Unroll the computation over time



2. Run backprop through the resulting feed-forward network

(Robinson & Fallside, 1987)  
(Werbos, 1988)  
(Mozer, 1995)



# Bidirectional RNN

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\vec{\mathbf{h}}$  and  $\overleftarrow{\mathbf{h}}$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

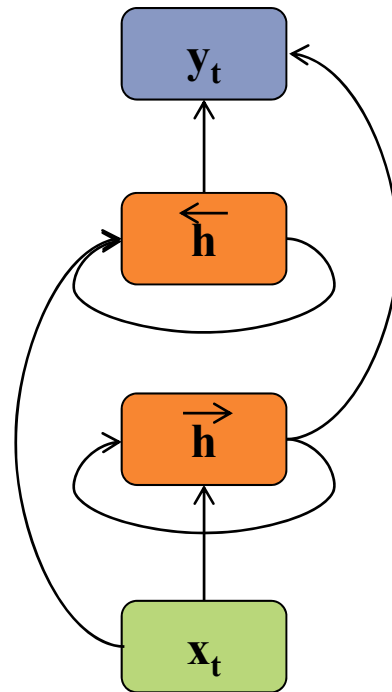
nonlinearity:  $\mathcal{H}$

Recursive Definition:

$$\vec{h}_t = \mathcal{H} \left( W_{x\vec{h}} x_t + W_{\vec{h}\vec{h}} \vec{h}_{t-1} + b_{\vec{h}} \right)$$

$$\overleftarrow{h}_t = \mathcal{H} \left( W_{x\overleftarrow{h}} x_t + W_{\overleftarrow{h}\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}} \right)$$

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_y$$



# Bidirectional RNN

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\vec{\mathbf{h}}$  and  $\overleftarrow{\mathbf{h}}$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

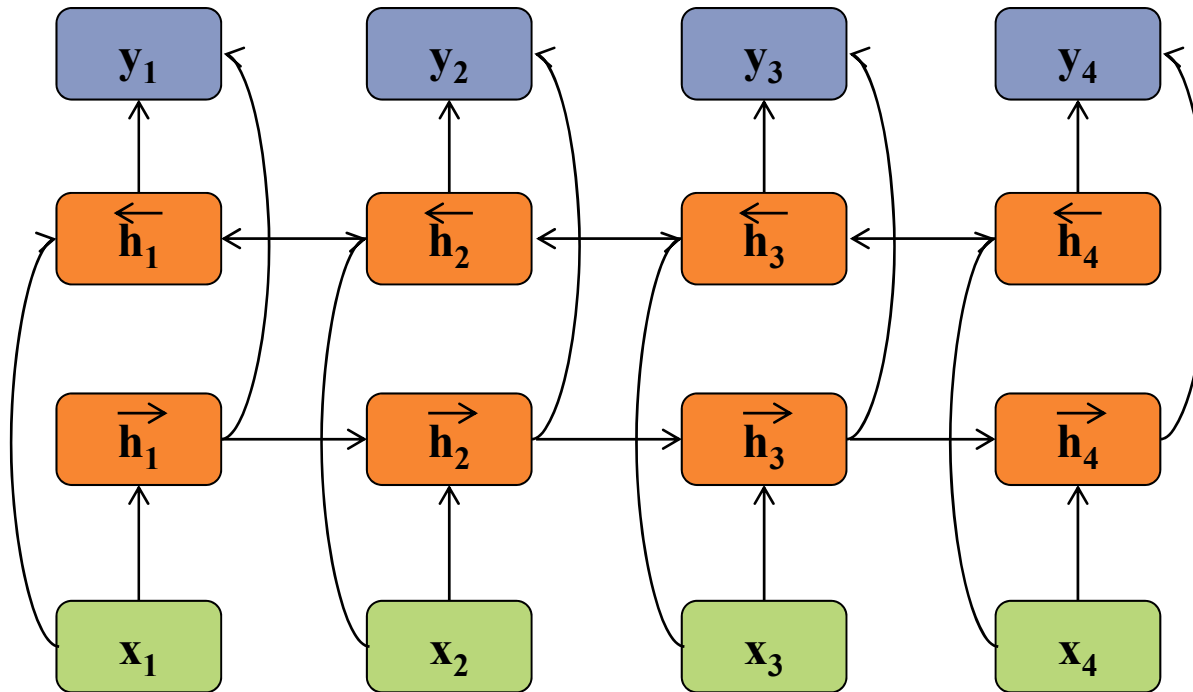
nonlinearity:  $\mathcal{H}$

Recursive Definition:

$$\vec{h}_t = \mathcal{H} \left( W_{x\vec{h}} x_t + W_{\vec{h}\vec{h}} \vec{h}_{t-1} + b_{\vec{h}} \right)$$

$$\overleftarrow{h}_t = \mathcal{H} \left( W_{x\overleftarrow{h}} x_t + W_{\overleftarrow{h}\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}} \right)$$

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_y$$



# Bidirectional RNN

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\vec{\mathbf{h}}$  and  $\overleftarrow{\mathbf{h}}$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

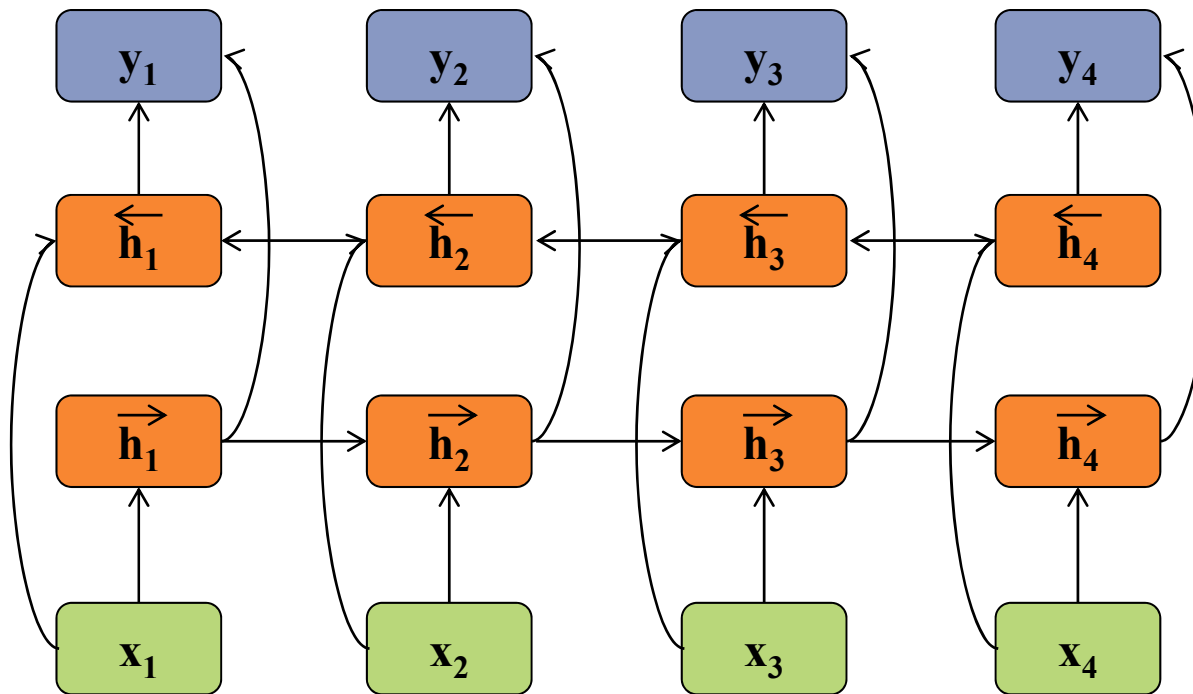
nonlinearity:  $\mathcal{H}$

Recursive Definition:

$$\vec{h}_t = \mathcal{H} \left( W_{x\vec{h}} x_t + W_{\vec{h}\vec{h}} \vec{h}_{t-1} + b_{\vec{h}} \right)$$

$$\overleftarrow{h}_t = \mathcal{H} \left( W_{x\overleftarrow{h}} x_t + W_{\overleftarrow{h}\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}} \right)$$

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_y$$



# Bidirectional RNN

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

hidden units:  $\vec{\mathbf{h}}$  and  $\overleftarrow{\mathbf{h}}$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

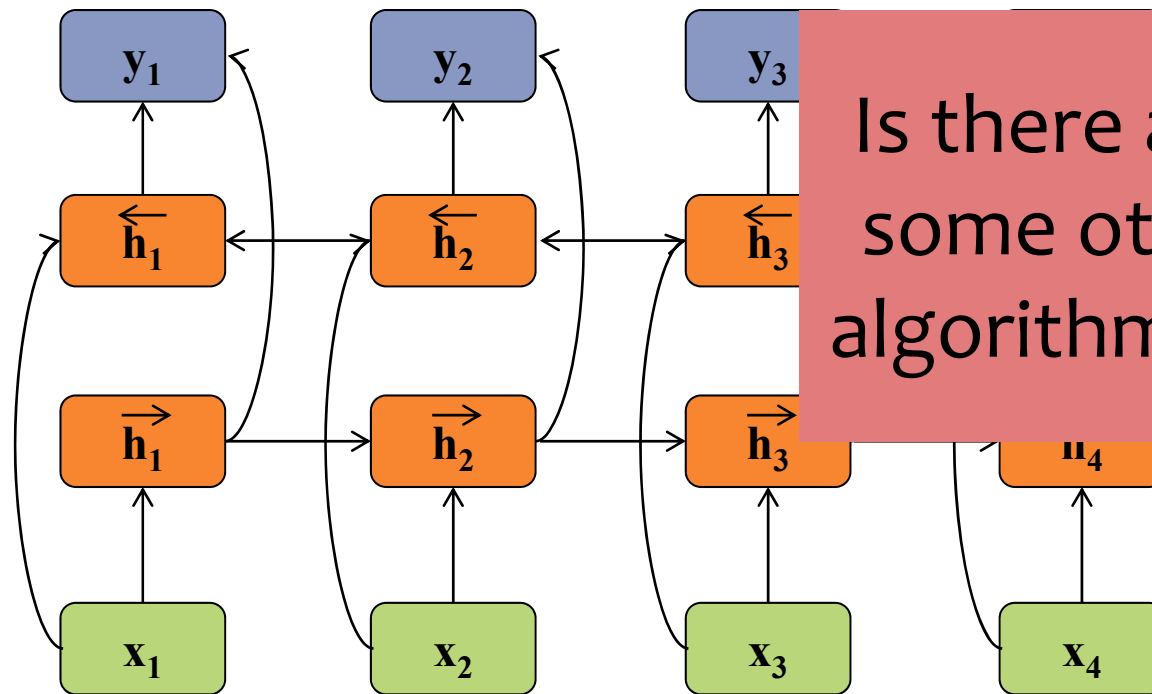
nonlinearity:  $\mathcal{H}$

Recursive Definition:

$$\vec{h}_t = \mathcal{H} \left( W_{x\vec{h}} x_t + W_{\vec{h}\vec{h}} \vec{h}_{t-1} + b_{\vec{h}} \right)$$

$$\overleftarrow{h}_t = \mathcal{H} \left( W_{x\overleftarrow{h}} x_t + W_{\overleftarrow{h}\overleftarrow{h}} \overleftarrow{h}_{t+1} + b_{\overleftarrow{h}} \right)$$

$$y_t = W_{\vec{h}y} \vec{h}_t + W_{\overleftarrow{h}y} \overleftarrow{h}_t + b_y$$



Is there an analogy to some other recursive algorithm(s) we know?

# Deep RNNs

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity:  $\mathcal{H}$

Recursive Definition:

$$h_t^n = \mathcal{H}(W_{h^{n-1}h^n} h_t^{n-1} + W_{h^n h^n} h_{t-1}^n + b_h^n)$$

$$y_t = W_{h^N y} h_t^N + b_y$$

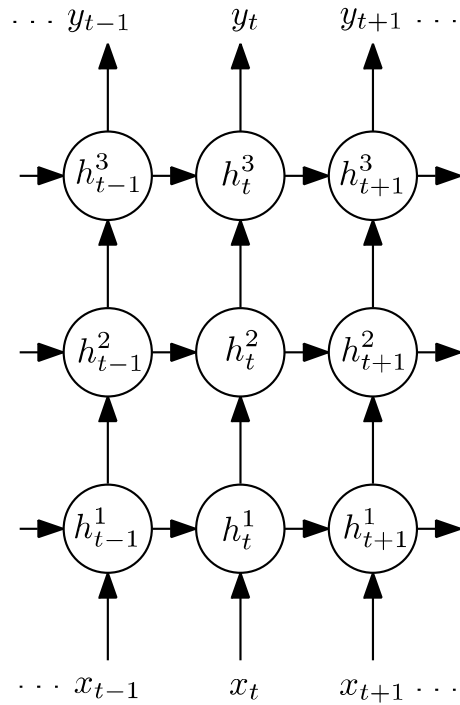


Figure from (Graves et al., 2013)

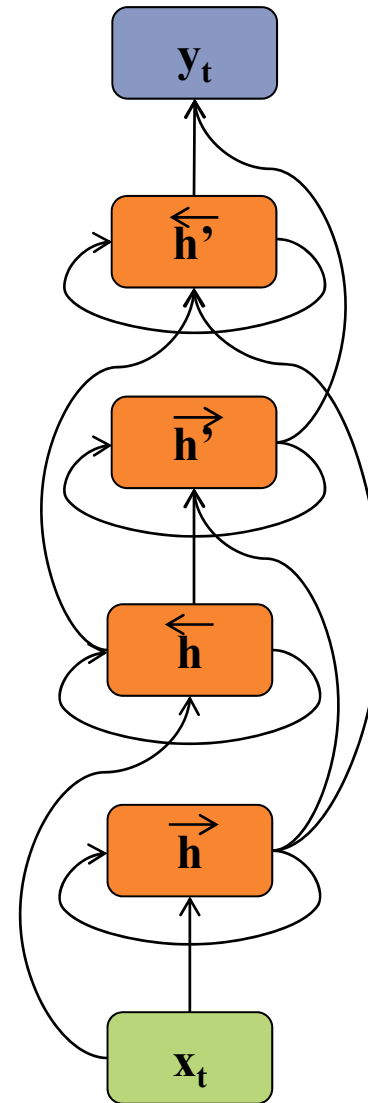
# Deep Bidirectional RNNs

inputs:  $\mathbf{x} = (x_1, x_2, \dots, x_T), x_i \in \mathcal{R}^I$

outputs:  $\mathbf{y} = (y_1, y_2, \dots, y_T), y_i \in \mathcal{R}^K$

nonlinearity:  $\mathcal{H}$

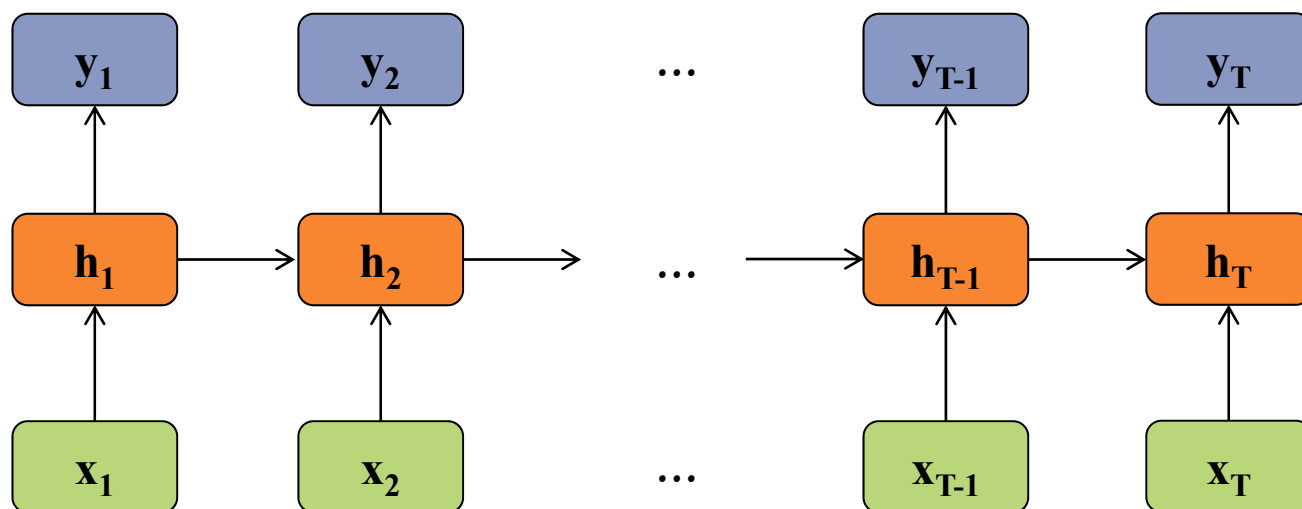
- Notice that the upper level hidden units have input from **two previous layers** (i.e. wider input)
- Likewise for the output layer
- What analogy can we draw to DNNs, DBNs, DBMs?



# Long Short-Term Memory (LSTM)

Motivation:

- Standard RNNs have trouble learning long distance dependencies
- LSTMs combat this issue





# Long Short-Term Memory (LSTM)

Motivation:

- Vanishing gradient problem for Standard RNNs
- Figure shows sensitivity (darker = more sensitive) to the input at time  $t=1$

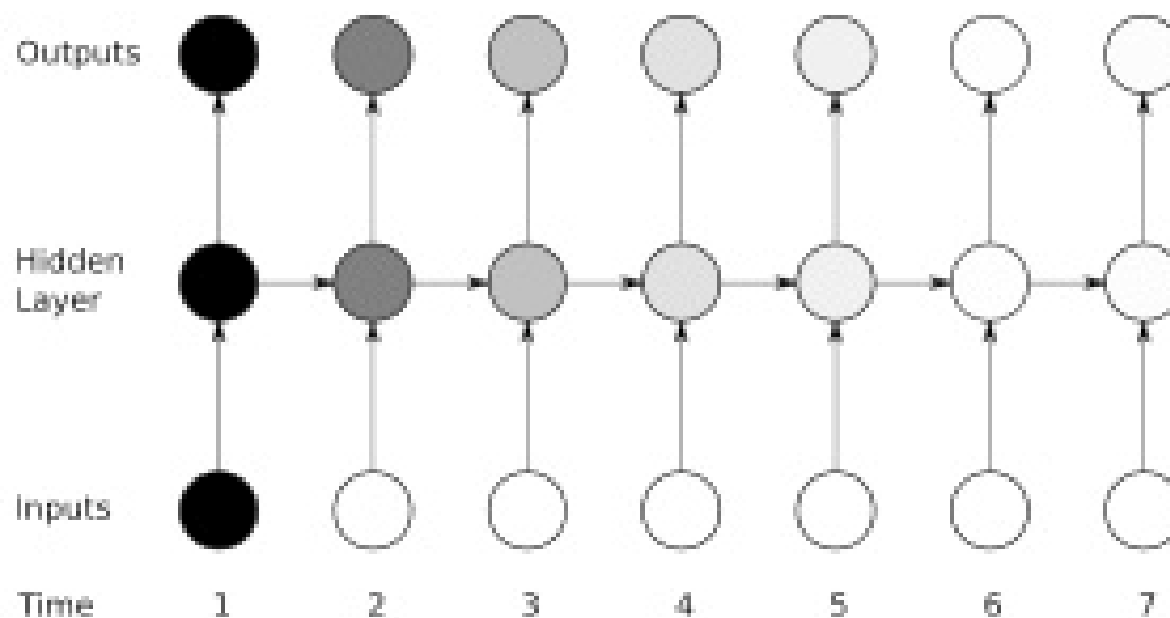


Figure from (Graves, 2012)

# Long Short-Term Memory (LSTM)

Motivation:

- LSTM units have a rich internal structure
- The various “gates” determine the propagation of information and can choose to “remember” or “forget” information

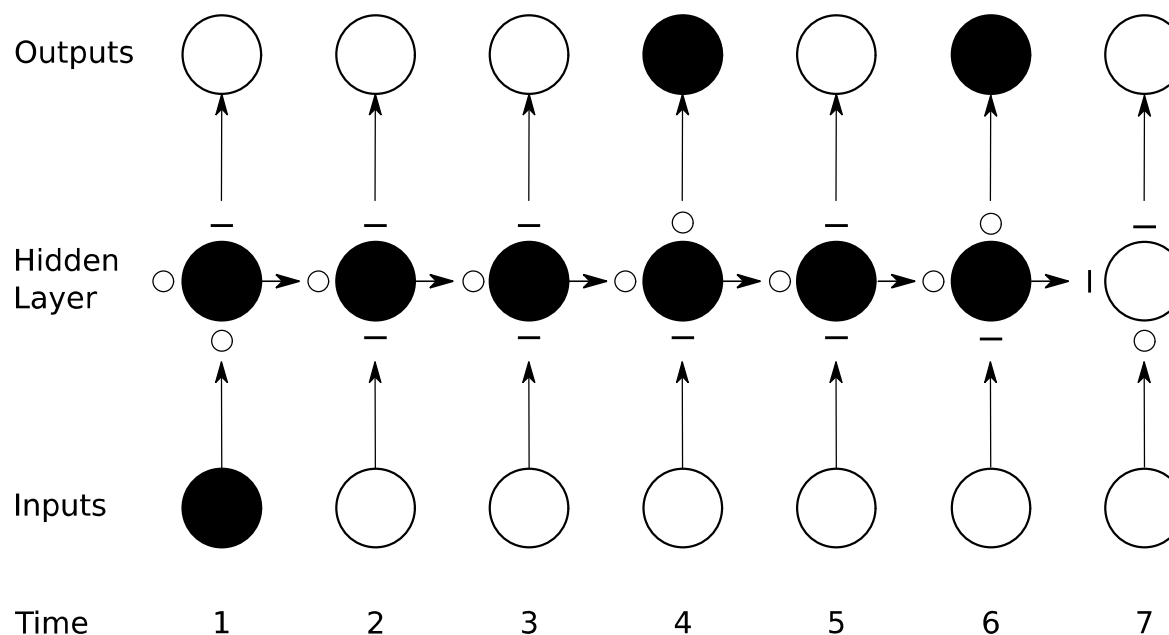
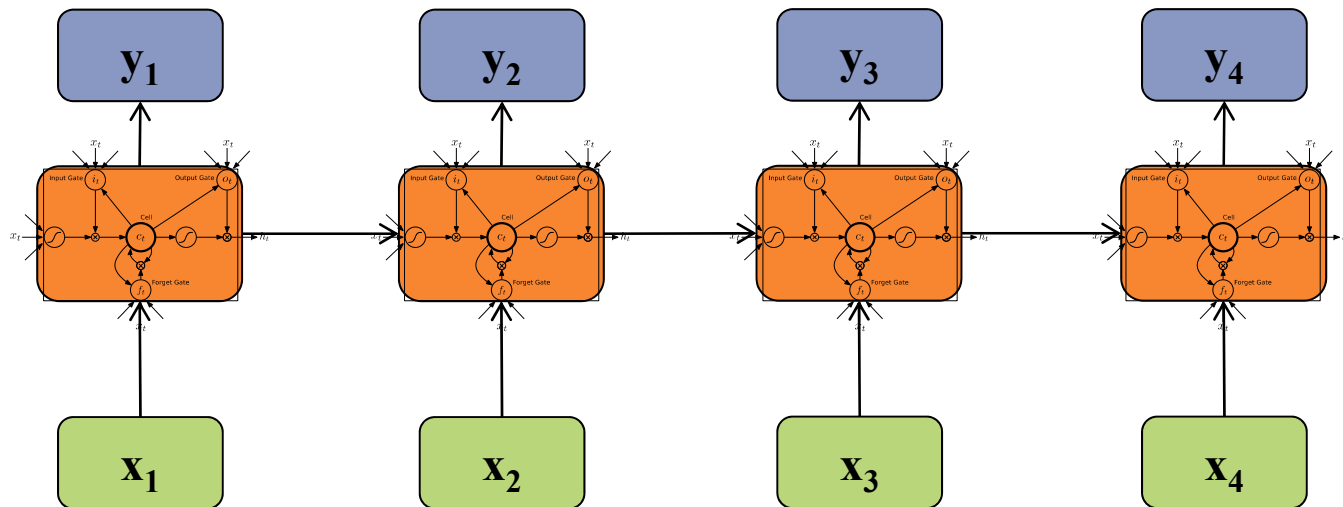


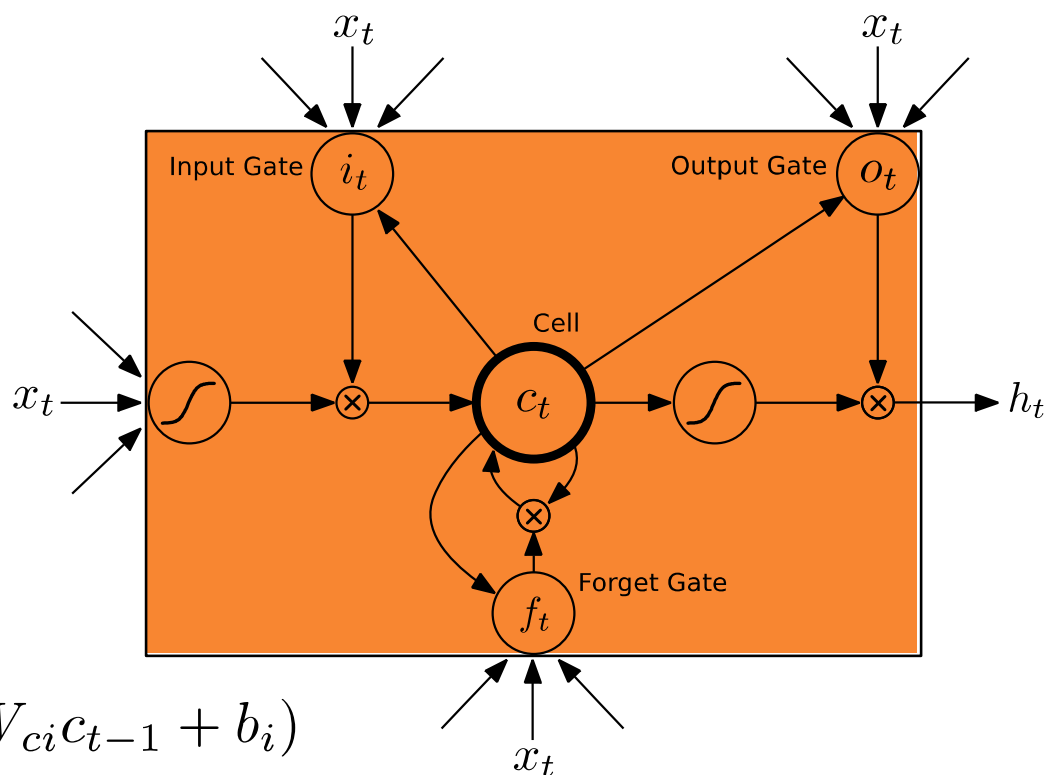
Figure from (Graves, 2012)

# Long Short-Term Memory (LSTM)



# Long Short-Term Memory (LSTM)

- **Input gate:** masks out the standard RNN inputs
- **Forget gate:** masks out the previous cell
- **Cell:** stores the input/forget mixture
- **Output gate:** masks out the values of the next hidden



$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$$

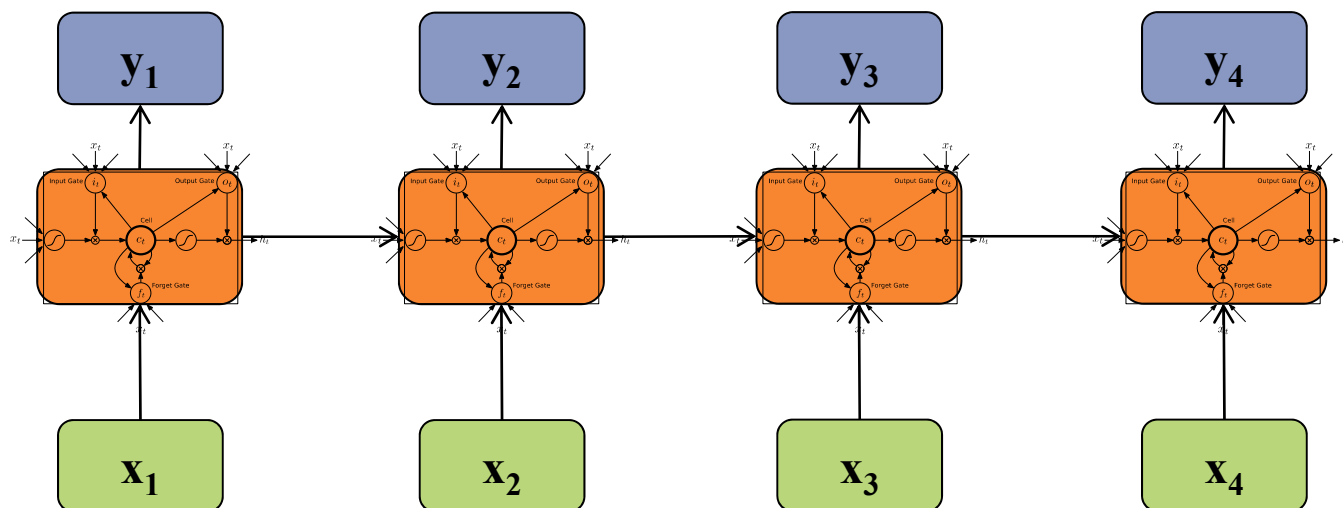
$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o)$$

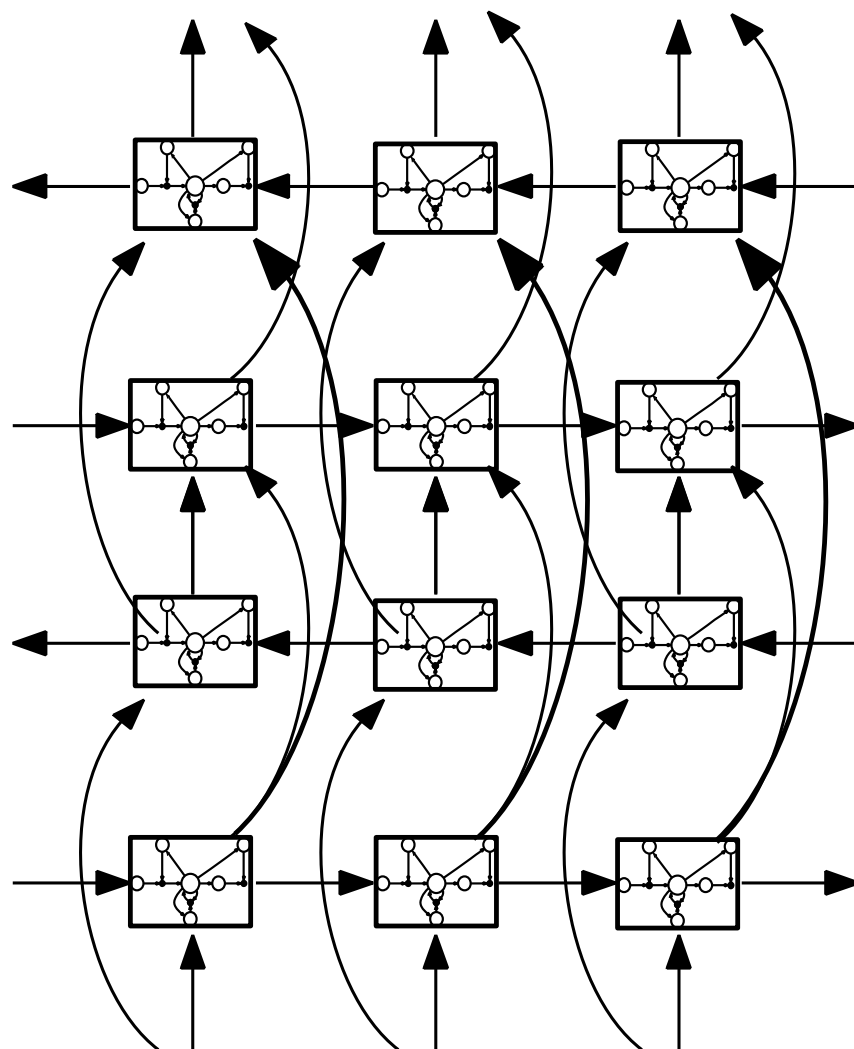
$$h_t = o_t \tanh(c_t)$$

Figure from (Graves et al., 2013)

# Long Short-Term Memory (LSTM)

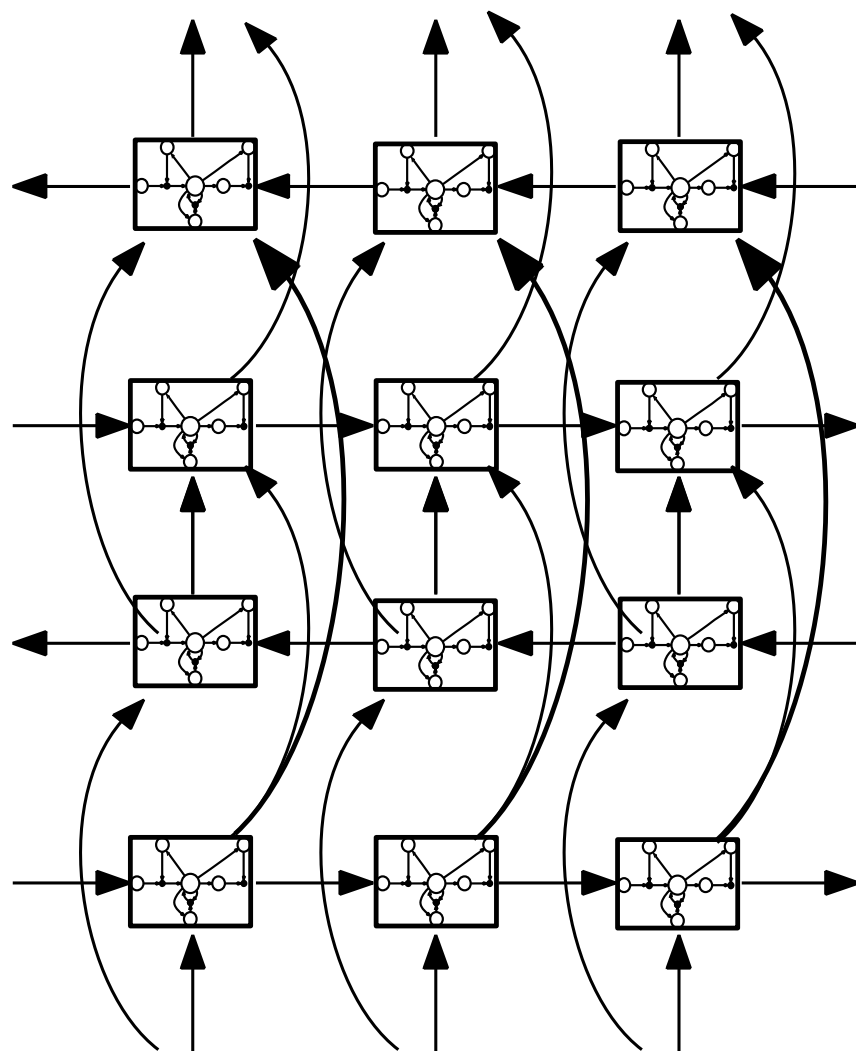


# Deep Bidirectional LSTM (DBLSTM)



- Figure: input/output layers not shown
- **Same general topology** as a Deep Bidirectional RNN, but with **LSTM units** in the hidden layers
- No additional **representational power** over DBRNN, but **easier to learn** in practice

# Deep Bidirectional LSTM (DBLSTM)



How important is this particular architecture?

Jozefowicz et al. (2015) **evaluated 10,000 different LSTM-like architectures** and found several variants that worked just as well on several tasks.

# RNN Training Tricks

- Deep Learning models tend to consist largely of **matrix multiplications**
- Training tricks:
  - **mini-batching with masking**

	Metric	DyC++	DyPy	Chainer	DyC++ Seq	Theano	TF
RNNLM (MB=1)	words/sec	190	190	114	494	189	298
RNNLM (MB=4)	words/sec	830	825	295	1510	567	473
RNNLM (MB=16)	words/sec	1820	1880	794	2400	1100	606
RNNLM (MB=64)	words/sec	2440	2470	1340	2820	1260	636

- **sorting into buckets of similar-length sequences**, so that mini-batches have same length sentences
- **truncated BPTT**, when sequences are too long, divide sequences into chunks and use the final vector of the previous chunk as the initial vector for the next chunk (but don't backprop from next chunk to previous chunk)



# RNN Summary

- **RNNs**
  - Applicable to tasks such as **sequence labeling**, speech recognition, machine translation, etc.
  - Able to **learn context features** for time series data
  - Vanishing gradients are still a problem – but **LSTM units** can help
- **Other Resources**
  - Christopher Olah’s blog post on LSTMs  
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>