# DeepRob

**Lecture 7**
**Convolutional Neural Networks**
**University of Michigan and University of Minnesota**

# Project 1—Reminder

- Instructions and code available on the website
  - Here: deeprob.org/projects/project1/

- Uses Python, PyTorch and Google Colab

- Implement KNN, linear SVM, and linear softmax classifiers

- **Autograder is online and updated**

- **Due Thursday, January 26th 11:59 PM EST**

# Recap from Previous Lecture

Represent complex expressions as **computational graphs**

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**

$$s = Wx$$

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

$x$

$W$

$*$

**Hinge loss**

$+$

$R$

$R(W)$

L

**1. Forward pass**: Compute outputs

**2. Backward pass**: Compute gradients

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

Downstream Gradients

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

$x$

$y$

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

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

$f$

Local Gradients

$z$

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

Upstream Gradient

# Recap from Previous Lecture

$$f(x) = W_2 \max(0, W_1 x + b_1) + b_2$$



**Input:** 3072 → $x$ → $W_1$ → $h$ → $W_2$ → $s$

**Hidden Layer:** 100

**Output:** 10

**Problem**: So far our classifiers don't respect the spatial structure of images!

Stretch pixels into column



Input image (2, 2)

| 0.2 | -0.5 | 0.1 | 2.0 |
|-----|------|-----|-----|
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

W (3, 4)

| 56 |
|----|
| 231 |
| 24 |
| 2 |

(4,)

+

| 1.1 |
|-----|
| 3.2 |
| -1.2 |

b (3,)

=

| -96.8 |
|-------|
| 437.9 |
| 61.95 |

(3,)

# Recap from Previous Lecture



$$f(x) = W_2 \max(0, W_1 x + b_1) + b_2$$

**Input:** 3072

$x$  $W_1$  $h$  $W_2$  $s$

**Hidden Layer:** 100

**Output:** 10

**Problem**: So far our classifiers don't respect the spatial structure of images!

**Solution**: Define new computational nodes that operate on images!

DR
# Components of Fully-Connected Networks

**Fully-Connected Layers**

**Activation Functions**

# Components of Convolutional Neural Networks

## Fully-Connected Layers



## Activation Functions



## Convolution Layers



## Pooling Layers



## Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Components of Convolutional Neural Networks

**Fully-Connected Layers**

$$x \quad W_1 \quad h \quad W_2 \quad s$$

**Activation Functions**

**Convolution Layers**

**Pooling Layers**

**Normalization**

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Fully-Connected Layer

3x32x32 image ➝ stretch to 3072x1

**Input**

1 ▭ 3072

➝ $Wx$
10 x 3072
Weights

➝ 1 ▭

**Output**

10

# Fully-Connected Layer

3x32x32 image $\longrightarrow$ stretch to 3072x1



**Input**

1

3072

$Wx$

10 x 3072
Weights

**Output**

1

10

**1 number:**
The result of taking a dot product
between a row of W and the input

# Convolution Layer

3x32x32 image:  preserve spatial structure

3x5x5 filter

32   height

32   width

3   depth / channels

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer

3x32x32 image

Filters always extend the full depth of the input volume

3x5x5 filter

32   height

32   width

3   depth / channels

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer

3x32x32 image

3x5x5 filter

32

32

3

**1 number:**
The result of taking a dot product between the filter and a small 3x5x5 portion of the image (i.e. 3*5*5=75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer

3x32x32 image

1x28x28 activation map

3x5x5 filter

convolve (slide) over all spatial locations

32

32

3

28

28

1

# Convolution Layer

3x32x32 image                                    two 1x28x28 activation map



3x5x5 filter

**Consider repeating with a
second (green) filter**

convolve (slide) over all spatial locations

32

32

3

28

28

1  1

# Convolution Layer

3x32x32 image

six   1x28x28 activation map

**Consider 6 filters,
each 3x5x5**



Convolution
Layer

6x3x5x5
filters

32

32

3

Stack activations to get
a 6x28x28 output image

# Convolution Layer

3x32x32 image

six  1x28x28 activation map

Also 6-dim bias vector

Convolution Layer

6x3x5x5 filters

Stack activations to get a 6x28x28 output image

32

32

3

# Convolution Layer

**DR**

3x32x32 image

32

32

3

Also 6-dim bias vector

Convolution Layer

6x3x5x5 filters

28x28 grid, at each point a 6-dim vector

Stack activations to get a 6x28x28 output image

# Convolution Layer



2x3x32x32
batch of images

2x6x28x28
batch of outputs

Also 6-dim bias vector

Convolution Layer

6x3x5x5
filters

32
32
3

# Convolution Layer

N x $C_{in}$ x H  xW
batch of images

N x $C_{out}$ x H' x W'
batch of outputs

Also $C_{out}$-dim bias vector



Convolution
Layer

$C_{out}$ x $C_{in}$ x $K_h$ x $K_w$
filters

32

32

3

$C_{out}$

# Stacking Convolutions

input
N x 3 x 32 x 32

$W_1$: 6x3x5x5
$b_1$: 6

First hidden layer
N x 6 x 28 x 28

$W_2$: 10x6x3x3
$b_2$: 10

Second hidden layer
N x 10 x 26 x 26

$W_3$: 12x10x3x3
$b_3$: 12

# Stacking Convolutions

**Q:** What happens if we stack two convolution layers?



input
N x 3 x 32 x 32

First hidden layer
N x 6 x 28 x 28

Second hidden layer
N x 10 x 26 x 26

$W_1$: 6x3x5x5
$b_1$: 6

$W_2$: 10x6x3x3
$b_2$: 10

$W_3$: 12x10x3x3
$b_3$: 12

# Stacking Convolutions



**Q:** What happens if we stack two convolution layers?

(Recall $y=W_2W_1x$ is a linear classifier)

$W_1$: 6x3x5x5
$b_1$: 6

$W_2$: 10x6x3x3
$b_2$: 10

$W_3$: 12x10x3x3
$b_3$: 12

input
N x 3 x 32 x 32

First hidden layer
N x 6 x 28 x 28

Second hidden layer
N x 10 x 26 x 26

# Stacking Convolutions

**Q:** What happens if we stack two convolution layers?

(Recall $y=W_2W_1x$ is a linear classifier)

**A:** We get another convolution!



32
32
3

input
N x 3 x 32 x 32

$W_1$: 6x3x5x5
$b_1$: 6

28
28
6

First hidden layer
N x 6 x 28 x 28

$W_2$: 10x6x3x3
$b_2$: 10

26
26
10

Second hidden layer
N x 10 x 26 x 26

$W_3$: 12x10x3x3
$b_3$: 12

# Stacking Convolutions

**Q:** What happens if we stack two convolution layers?

(Recall $y=W_2W_1x$ is a linear classifier)

**A:** We get another convolution!



$W_1$: 6x3x5x5
$b_1$: 6

$W_2$: 10x6x3x3
$b_2$: 10

$W_3$: 12x10x3x3
$b_3$: 12

input
N x 3 x 32 x 32

First hidden layer
N x 6 x 28 x 28

Second hidden layer
N x 10 x 26 x 26

24

# What do convolutional filters learn?

input
N x 3 x 32 x 32

First hidden layer
N x 6 x 28 x 28

Second hidden layer
N x 10 x 26 x 26

$W_1$: 6x3x5x5
$b_1$: 6

$W_2$: 10x6x3x3
$b_2$: 10

$W_3$: 12x10x3x3
$b_3$: 12

# What do convolutional filters learn?



W$_1$: 6x3x5x5
b$_1$: 6

32
32
3
input
N x 3 x 32 x 32

28
28
6
First hidden layer
N x 6 x 28 x 28

Linear classifier: One template per class



master chef can
cracker box
sugar box
tomato soup can
mustard bottle

fish can
gelatin box
meat can
mug
large marker

# What do convolutional filters learn?

32

28

**Conv** → **ReLU** →

$W_1$: 6x3x5x5

$b_1$: 6

3

32

6

28

input
N x 3 x 32 x 32

First hidden layer
N x 6 x 28 x 28

MLP: Bank of whole-image templates

# What do convolutional filters learn?

32

28

Conv → ReLU →

$W_1$: 6x3x5x5

$b_1$: 6

3

input
N x 3 x 32 x 32

6

28

First hidden layer
N x 6 x 28 x 28

First-layer conv filters: local image templates
(often learns oriented edges, opposing colors)



AlexNet: 96 filters, each 3x11x11

# A closer look at spatial dimensions



input
N x 3 x 32 x 32

$W_1$: 6x3x5x5
$b_1$: 6

First hidden layer
N x 6 x 28 x 28

# A closer look at spatial dimensions

Input: 7x7
Filter: 3x3

7

7

# A closer look at spatial dimensions



Input: 7x7
Filter: 3x3

7

7

# A closer look at spatial dimensions



Input: 7x7
Filter: 3x3

# A closer look at spatial dimensions

Input: 7x7
Filter: 3x3

7

7

# A closer look at spatial dimensions

Input: 7x7
Filter: 3x3
Output: 5x5

7

7

# A closer look at spatial dimensions

Input: 7x7
Filter: 3x3
Output: 5x5

In general:
Input: W
Filter: K
Output: W − K + 1

Problem: Feature maps "shrink" with each layer!

7

7

# A closer look at spatial dimensions

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | | | | | | | | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Input: 7x7
Filter: 3x3
Output: 5x5

In general:
Input: W
Filter: K
Output: W – K + 1

Problem: Feature maps "shrink" with each layer!

Solution: **padding**
Add zeros around the input

# A closer look at spatial dimensions

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 |   |   |   |   |   |   |   | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Input: 7x7
Filter: 3x3
Output: 5x5

In general:
Input: W
Filter: K
Padding: P
Output: W − K + 1 + 2P

Very common:
Set P = (K − 1) / 2 to make output have same size as input!

# Receptive Fields

For convolution with kernel size K, each element in the output depends on a K x K **receptive field** in the input



Input                    Output

# Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



Input

Output

Be careful – "receptive field in the input" vs "receptive field in the previous layer"
Hopefully clear from context!

# Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



Input

Output

Problem: For large images we need many layers
for each output to "see" the whole image image

# Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



Input

Output

Problem: For large images we need many layers
for each output to "see" the whole image image

Solution: Downsample inside the network

# Strided Convolution

Input: 7x7
Filter: 3x3
Stride: 2

# Strided Convolution

Input: 7x7
Filter: 3x3
Stride: 2

# Strided Convolution

Input: 7x7
Filter: 3x3          Output: 3x3
Stride: 2

# Strided Convolution



Input: 7x7
Filter: 3x3          Output: 3x3
Stride: 2

In general:
Input: W
Filter: K
Padding: P
Stride: S
Output: $(W - K + 2P) / S + 1$

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

**Q:** What is the output volume size?

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

**Q:** What is the output volume size?

(32-5+2*2) / 1 + 1 = 32 spatially

So, 10 x 32 x 32 output

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32
**Q:** What is the number of learnable parameters?

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32
**Q:** What is the number of learnable parameters?

Parmeters per filter: (3*5*5) + 1 = 76

10 filters, so total is 10*76 = 760

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: 760

**Q:** What is the number of multiply-add operations?

# Convolution Example

Input volume: 3 x 32 x 32
10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32
Number of learnable parameters: 760
**Q:** What is the number of multiply-add operations?
10*32*32=10,240 outputs, each from inner product

of two 3x5x5 tensors, so total = 75 * 10,240 = **768,000**

# Example: 1x1 Convolution

1x1 Conv
with 32 filters

Each filter has size 1x1x64 and
performs a 64-dimensional dot product

56

56

64

56

56

32

# Example: 1x1 Convolution



1x1 Conv
with 32 filters

Each filter has size 1x1x64 and
performs a 64-dimensional dot product

Stacking 1x1 conv layers gives MLP
operating on each input position

Lin et al., "Network in Network", ICLR 2014

# Convolution Summary

**Input**: $C_{in}$ x H x W

**Hyperparameters**:
- **Kernel size**: $K_H$ x $K_W$
- **Number filters**: $C_{out}$
- **Padding**: P
- **Stride**: S

**Weight matrix**: $C_{out}$ x $C_{in}$ x $K_H$ x $K_W$
giving $C_{out}$ filters of size $C_{in}$ x $K_H$ x $K_W$

**Bias vector**: $C_{out}$

**Output size**: $C_{out}$ x H' x W' where:
- H' = (H − K + 2P) / S + 1
- W' = (W − K + 2P) / S + 1

# Convolution Summary

**Input**: $C_{in}$ x H x W

**Hyperparameters**:
- **Kernel size**: $K_H$ x $K_W$
- **Number filters**: $C_{out}$
- **Padding**: P
- **Stride**: S

**Weight matrix**: $C_{out}$ x $C_{in}$ x $K_H$ x $K_W$
giving $C_{out}$ filters of size $C_{in}$ x $K_H$ x $K_W$

**Bias vector**: $C_{out}$

**Output size**: $C_{out}$ x H' x W' where:
- H' = (H − K + 2P) / S + 1
- W' = (W − K + 2P) / S + 1

Common settings:
$K_H = K_W$  (Small square filters)
P = (K − 1) / 2  ("Same" padding)
$C_{in}$, $C_{out}$ = 32, 64, 128, 256 (powers of 2)
K = 3, P = 1, S = 1 (3x3 conv)
K = 5, P = 2, S = 1 (5x5 conv)
K = 1, P = 0, S = 1 (1x1 conv)
K = 3, P = 1, S = 2 (Downsample by 2)

# Other types of convolution

So far: 2D Convolution

Input: $C_{in}$ x H x W
Weights: $C_{out}$ x $C_{in}$ x K x K

H

W

$C_{in}$

# Other types of convolution

So far: 2D Convolution

Input: $C_{in}$ x H x W
Weights: $C_{out}$ x $C_{in}$ x K x K

H

W

$C_{in}$

1D Convolution

Input: $C_{in}$ x W
Weights: $C_{out}$ x $C_{in}$ x K

$C_{in}$

W

# Other types of convolution

So far: 2D Convolution

Input: $C_{in}$ x H x W
Weights: $C_{out}$ x $C_{in}$ x K x K

H

W

$C_{in}$

3D Convolution

Input: $C_{in}$ x H x W x D
Weights: $C_{out}$ x $C_{in}$ x K x K x K

$C_{in}$-dim vector
at each point
in the volume

H

D

W

# PyTorch Convolution Layer

## Conv2d

CLASS  `torch.nn.Conv2d(`*`in_channels, out_channels, kernel_size, stride=1, padding=0,`* *`dilation=1, groups=1, bias=True, padding_mode='zeros'`*`)`    [SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size $(N, C_{\text{in}}, H, W)$ and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) \star \text{input}(N_i, k)$$

# PyTorch Convolution Layer

Conv2d

| CLASS | torch.nn.Conv2d(*in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros'*) | [SOURCE] |

Conv1d

| CLASS | torch.nn.Conv1d(*in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros'*) | [SOURCE] 🔗 |

Conv3d

| CLASS | torch.nn.Conv3d(*in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros'*) | [SOURCE] |

**DR**

## Fully-Connected Layers



$x$  $W_1$  $x$  $W_2$  $h$  $S$  S  S

S

## Activation Functions



## Convolution Layers



## Pooling Layers



224x224x64

pool

112x112x64

224

downsampling

112

224

112

## Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Pooling Layers: Another way to downsample

# Pooling Layers: Another way to downsample



**Hyperparameters:**
Kernel size
Stride
Pooling function

# Max Pooling

Single depth slice



Max pooling with
2x2 kernel size
stride of 2

# Max Pooling

Single depth slice

| | | | |
|---|---|---|---|
| 1 | 1 | 2 | 4 |
| 5 | **6** | 7 | **8** |
| **3** | 2 | 1 | 0 |
| 1 | 2 | 3 | **4** |

x

y

Max pooling with 2x2 kernel size stride of 2

| | |
|---|---|
| 6 | 8 |
| 3 | 4 |

Introduces invariance to small spatial shifts

No learnable parameters!

# Pooling Summary

**Input**: C x H x W

**Hyperparameters**:
- Kernel size: K
- Stride: S
- Pooling function (max, avg)

**Output**: C x H' x W' where
- H' = (H − K) / S + 1
- W' = (W − K) / S + 1

**Learnable parameters**: None!

Common settings:
max, K = 2, S = 2
max, K = 3, S = 2 (AlexNet)

# Components of Convolutional Neural Networks

**Fully-Connected Layers**

**Activation Functions**

$x$  $W_1$  $h$  $W_2$  $s$

x  h  s

h  s

**Convolution Layers**

**Pooling Layers**

**Normalization**

224x224x64

pool

112x112x64

224

downsampling

112

224

112

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Convolutional Neural Networks

Classic architecture: [Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC

Example: LeNet-5



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|-------|-------------|-------------|
| Input | 1 x 28 x 28 | |

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |



Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |
| Linear (2450 -> 500) | 500 | 2450 x 500 |
| ReLU | 500 | |

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |
| Linear (2450 -> 500) | 500 | 2450 x 500 |
| ReLU | 500 | |
| Linear (500 -> 10) | 10 | 500 x 10 |



Input    Image Maps    Output

Convolutions    Subsampling    Fully Connected

Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |
| Linear (2450 -> 500) | 500 | 2450 x 500 |
| ReLU | 500 | |
| Linear (500 -> 10) | 10 | 500 x 10 |



As we progress through the network:

Spatial size decreases
(using pooling or striped convolution)

Number of channels increases
(total "volume" is preserved!)

Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Example: LeNet-5

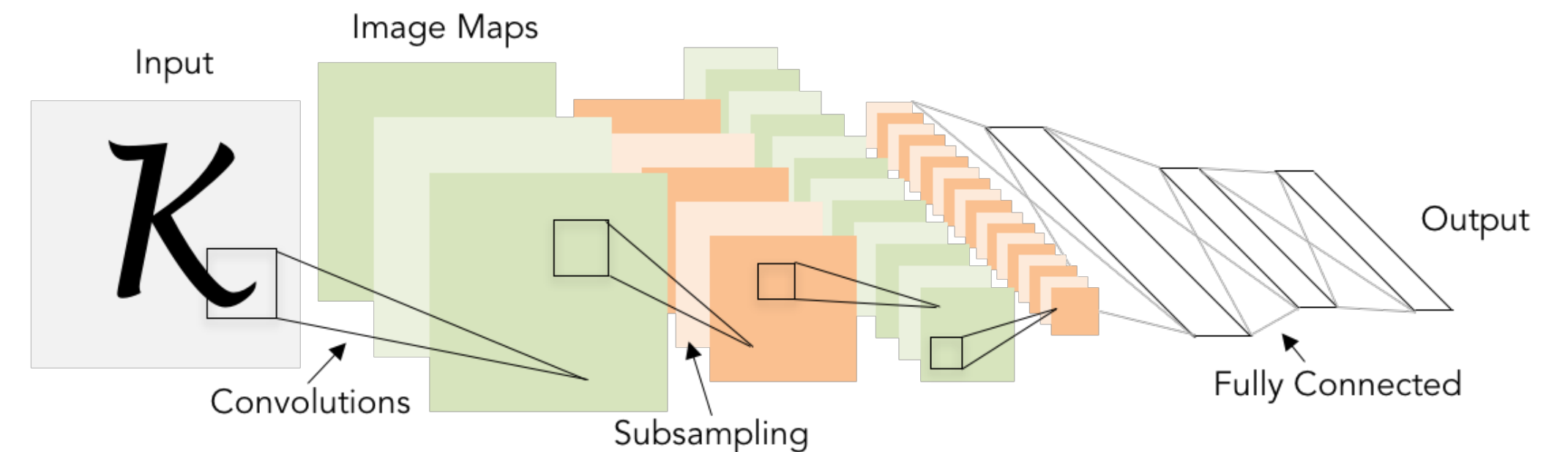| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |
| Linear (2450 -> 500) | 500 | 2450 x 500 |
| ReLU | 500 | |
| Linear (500 -> 10) | 10 | 500 x 10 |



As we progress through the network:

Spatial size **decreases**
(using pooling or striped convolution)

Number of channels **increases**
(total "volume" is preserved!)
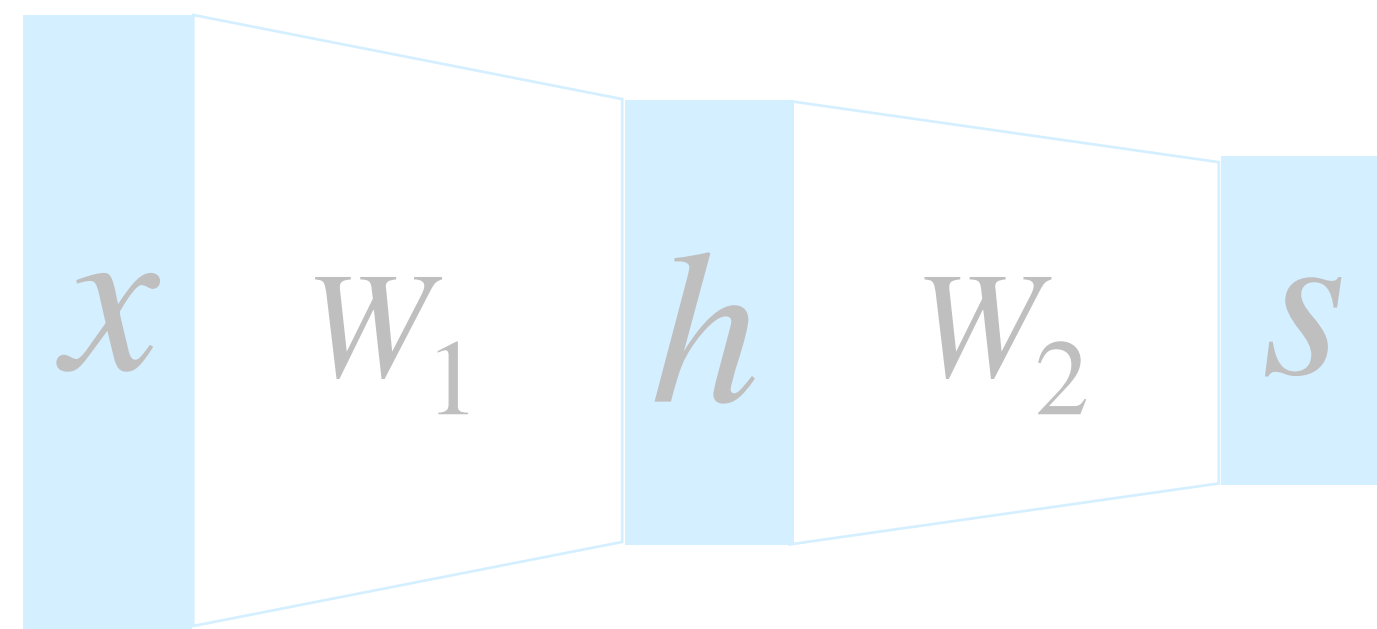
Some modern architectures
break this trend—stay tuned!

Lecun et al., "Gradient-based learning applied to document recognition", 1998

# Problem: Deep Networks very hard to train

# Components of Convolutional Neural Networks

## Fully-Connected Layers

$x$ $\quad W_1$ $\quad h$ $\quad x$ $\quad W_2$ $\quad h$ $\quad S$ $\quad s$ $\quad s$ $\quad s$

## Activation Functions

## Convolution Layers

## Pooling Layers

224x224x64

pool

112x112x64

224

downsampling

112

112

224

## Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Batch Normalization

Idea: "Normalize" the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization results

We can normalize a batch of activations using:

$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

# Batch Normalization

Idea: "Normalize" the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization results

We can normalize a batch of activations using:

$$\hat{x} = \frac{x - E[x]}{\sqrt{Var[x]}}$$

This is a **differentiable function**, so we can use it as an operator in our networks and backdrop through it!

Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

# Summary: Components of Convolutional Network
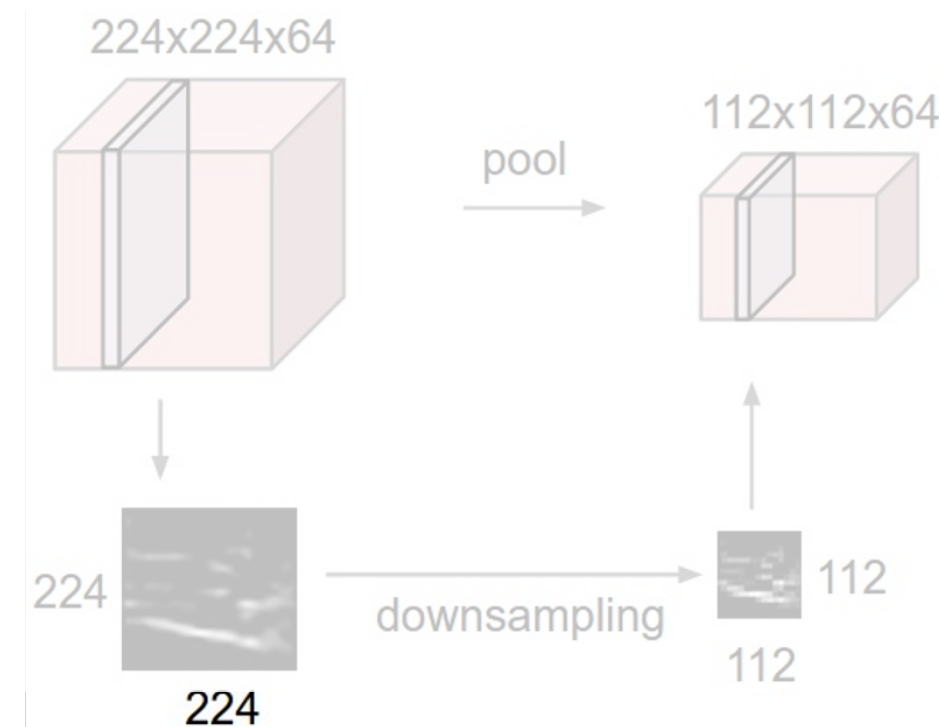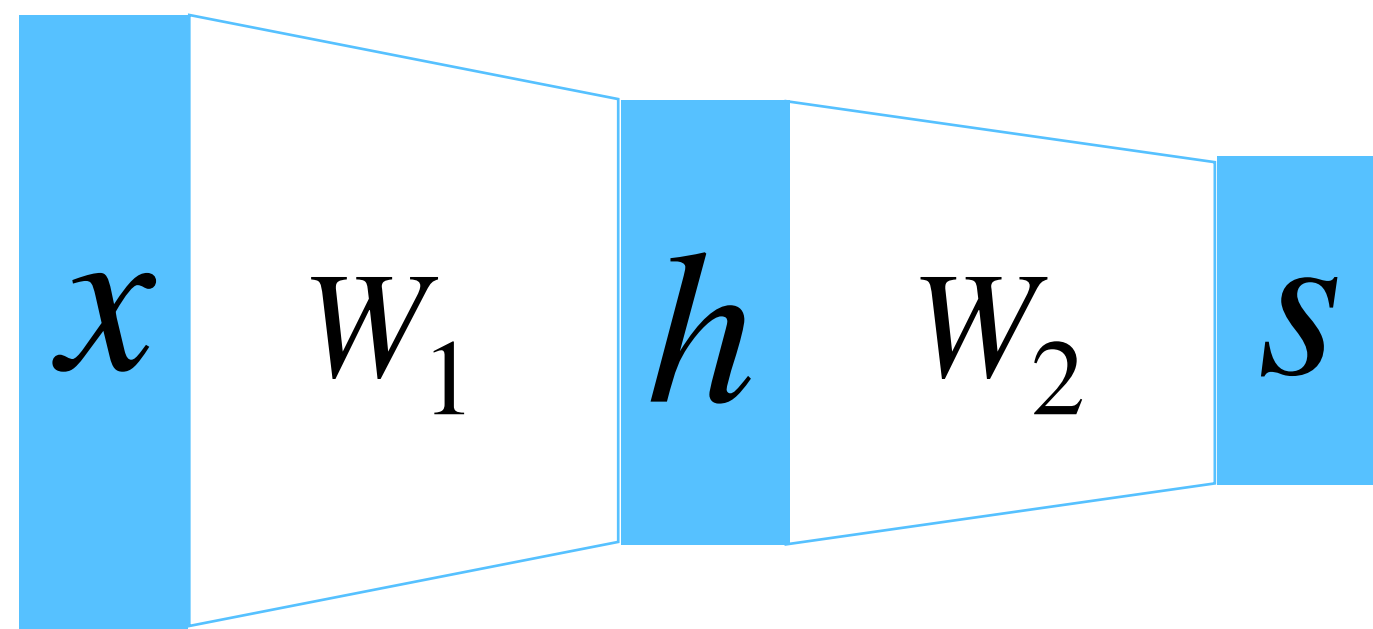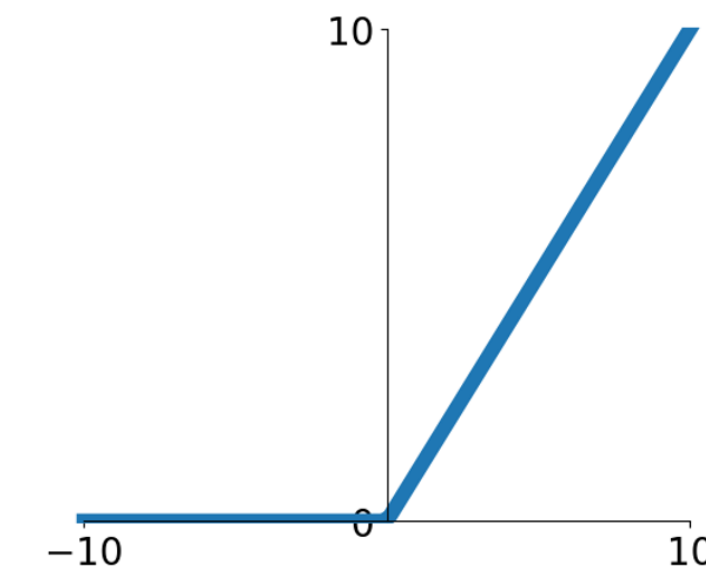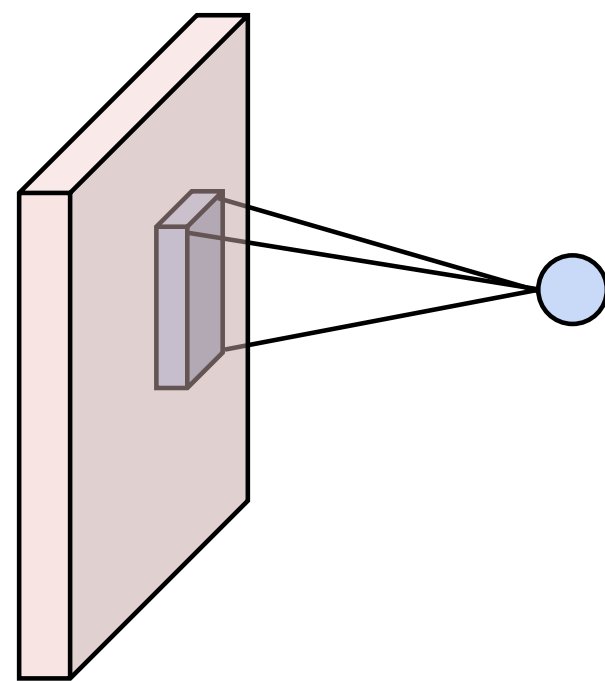
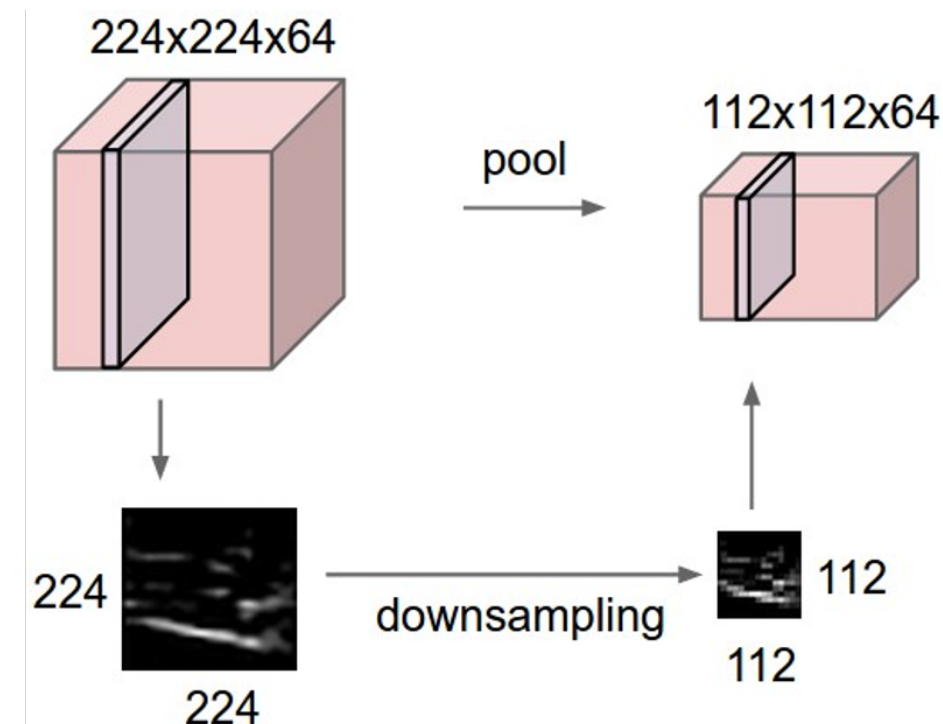**Fully-Connected Layers**

**Activation Functions**

**Convolution Layers**

**Pooling Layers**

**Normalization**

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

# Summary: Components of Convolutional Network

**Problem:** What is the right way to combine all these components?

# Next time: CNN Architectures

# Final Project Overview

- Research-oriented final project

  - Instead of a final exam!

- Objectives

  - Gain experience reading literature

  - Reproduce published results

  - Propose a new idea and test the results!

# Final Project Overview

- Research-oriented final project

  - Instead of a final exam!

Can be completed in teams of 1-3 people

- Objectives

  - Gain experience reading literature

  - Reproduce published results

  - Propose a new idea and test the results!

# Final Project Deliverables

1. A written paper review

2. In-class paper presentation

3. Reproduce published results

4. Extend results with new idea, technique or dataset

5. Document results in written report
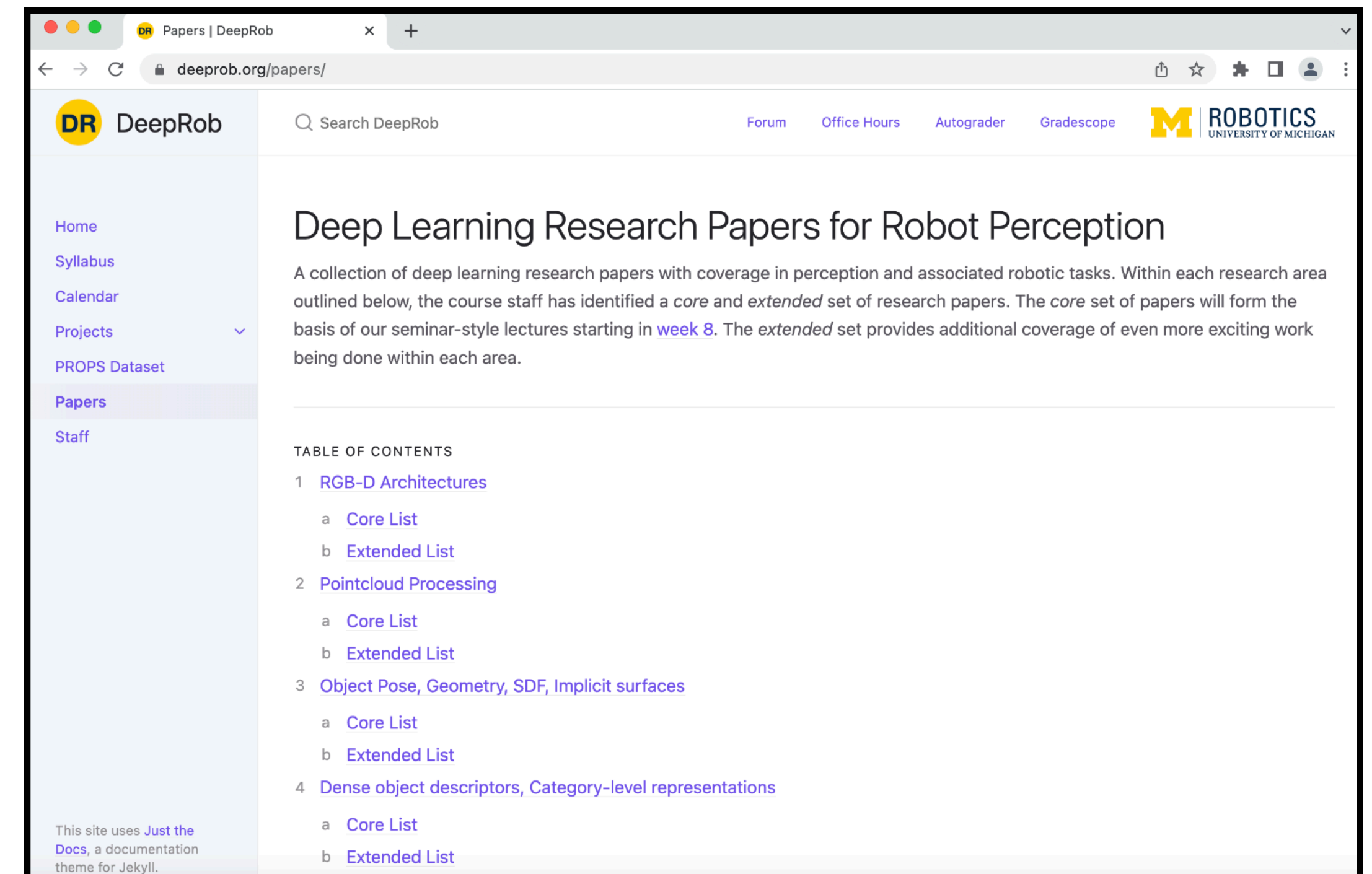
# (1) Paper Review and (2) Presentation

## Final project teams will be based on overlapping interest

Students will choose from the 'core' list
of papers on course website

Each team will be assigned one of the
'core' papers to review and present in-class

The 1-page paper review will be due **1-week before** the scheduled presentation

Presentation schedule will be based on
paper topic as shown in course calendar



More details on review and presentation
criteria in following lectures

# (3) Paper Reproduction and (4) Extension

Each team will choose a paper relating
to deep learning and robot perception

Doesn't have to be same paper you presented in class

Then reimplement and reproduce at least one of the
paper's published results **(not necessarily all the results)**

Then, each team will test one of their own ideas!

By extending the paper's model using new architecture or technique or dataset
**Your chance to experiment with deep learning and contribute to the field!**

More details on reproduction and extension
in following lectures

# (5) Project Report

- The final deliverable for your final project

- A 1-2 page paper

  - What problem within robot perception or manipulation?

  - What work has been done in this area?

  - What approach did you investigate?

  - What questions and directions exist for future work?

More details on report in following lectures

# Final Project Grading Overview

- Final Project:

  - Paper Review: 3%

  - Paper Presentation: 3%

  - Paper Reproduction: 6%

  - Algorithmic Extension: 6%

  - Written Report: 6%

More details on report in following lectures

# DeepRob

**Lecture 7**
**Convolutional Neural Networks**
**University of Michigan and University of Minnesota**