

CS 291A: Deep Learning for NLP

Convolutional Neural Networks

William Wang
UCSB Computer Science
wiliam@cs.ucsb.edu

Slides adapted from V. Chen and M. Chang.

Midterm Report

- Due today to reader Ke Ni at ke00@ucsb.edu.
- By now you should have obtained some preliminary results for your project.

Talk: Adaptive Adversarial Learning for a Diverse Visual World

- **Judy Hoffman (UC Berkeley)**
- **Date:** Thursday, February 22, 2018 - 3:30pm
- **Location:** HFH 1132
- **Speaker:**
- Judy Hoffman



Convolutional Neural Networks

- We need a course to talk about this topic
 - <http://cs231n.stanford.edu/syllabus.html>
- However, we only have one lecture.

Outline

- CNN(Convolutional Neural Networks) Introduction
- Evolution of CNN
- Visualizing the Features
- Sentiment Analysis with CNN

Outline

- **CNN(Convolutional Neural Networks) Introduction**
- Evolution of CNN
- Visualizing the Features
- Sentiment Analysis by CNN

Image Recognition



mite

container ship

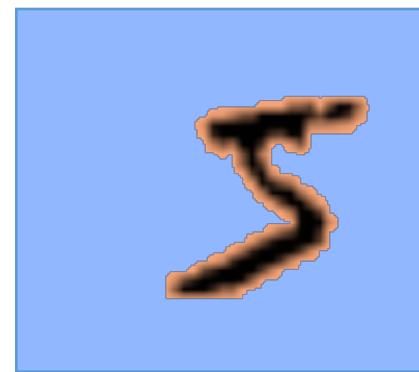
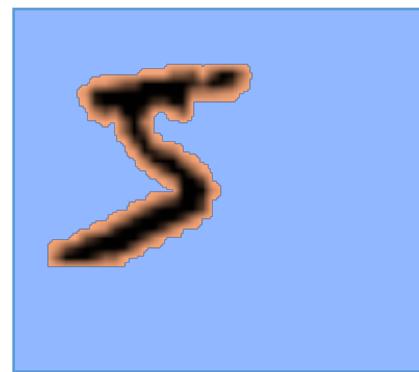
motor scooter

leopard

mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat

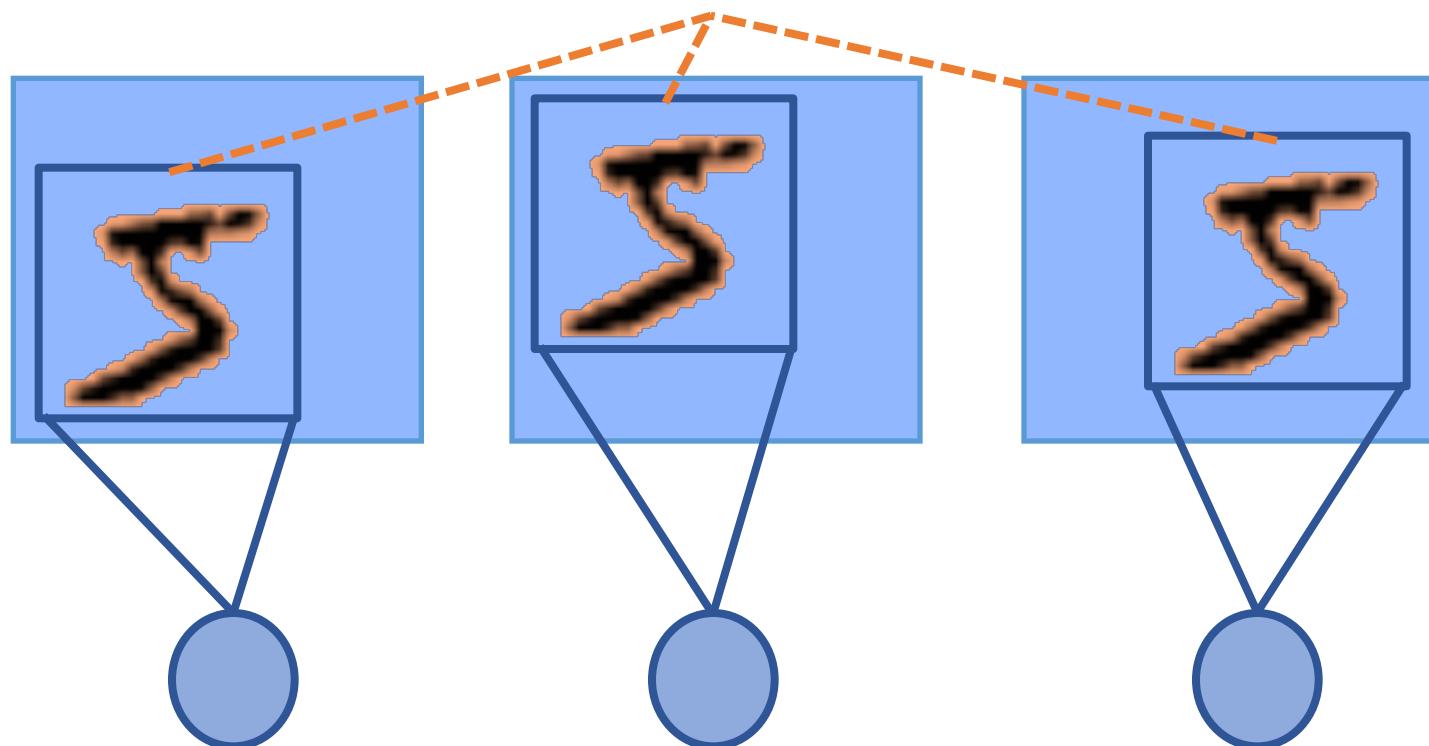
<http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>

Image Recognition



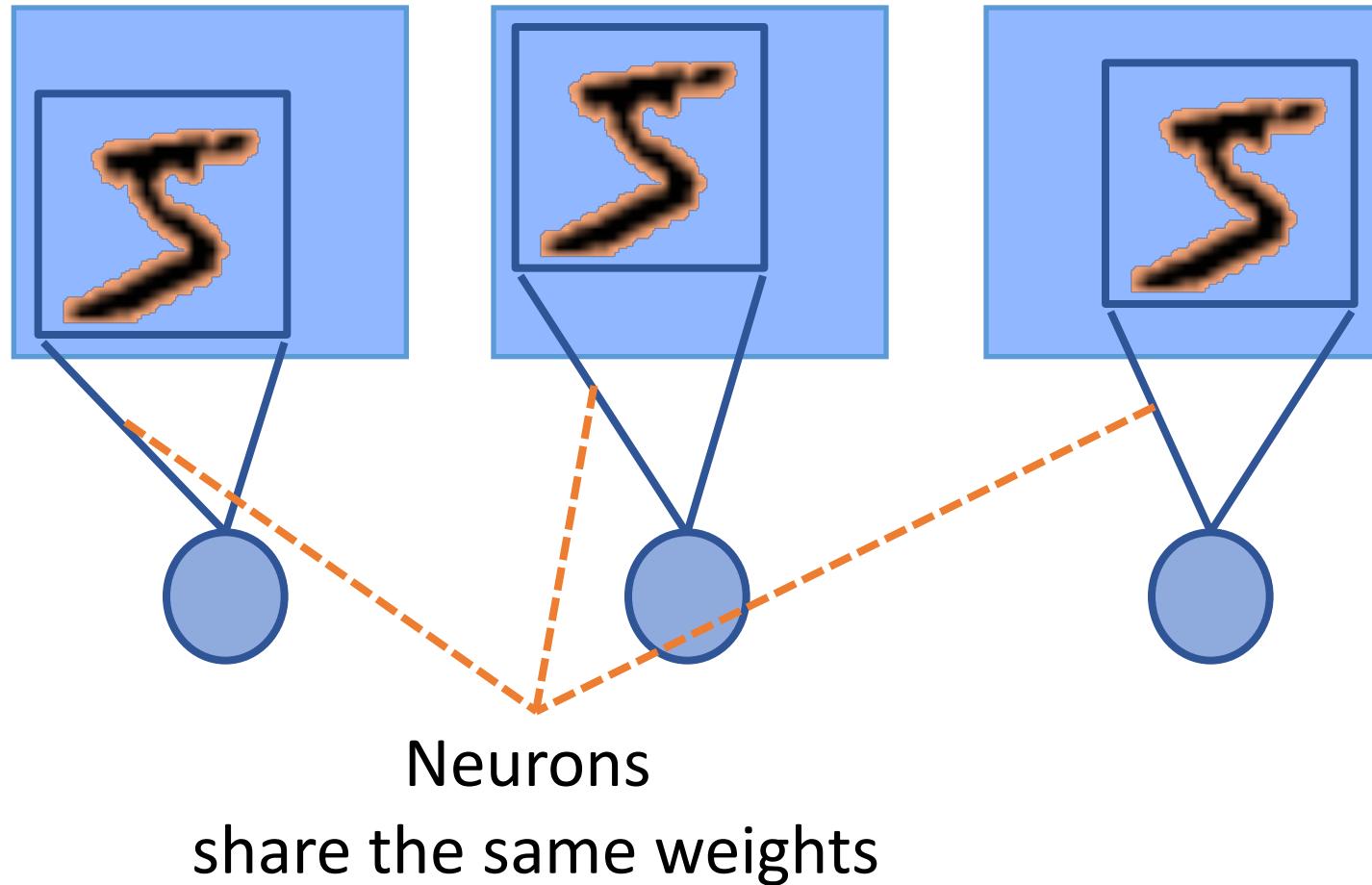
Local Connectivity

Neurons connect to a small region



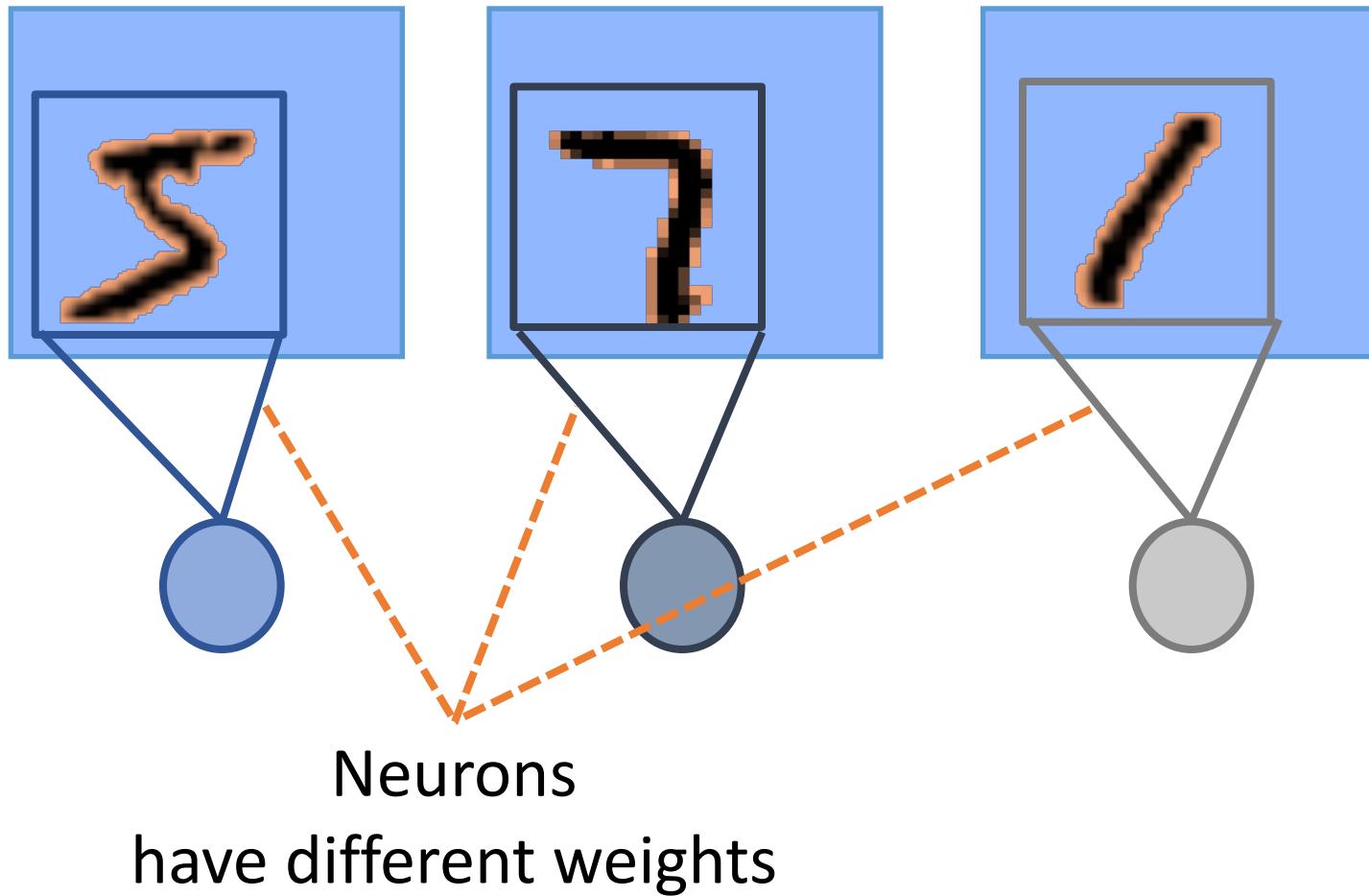
Parameter Sharing

- The same feature in different positions

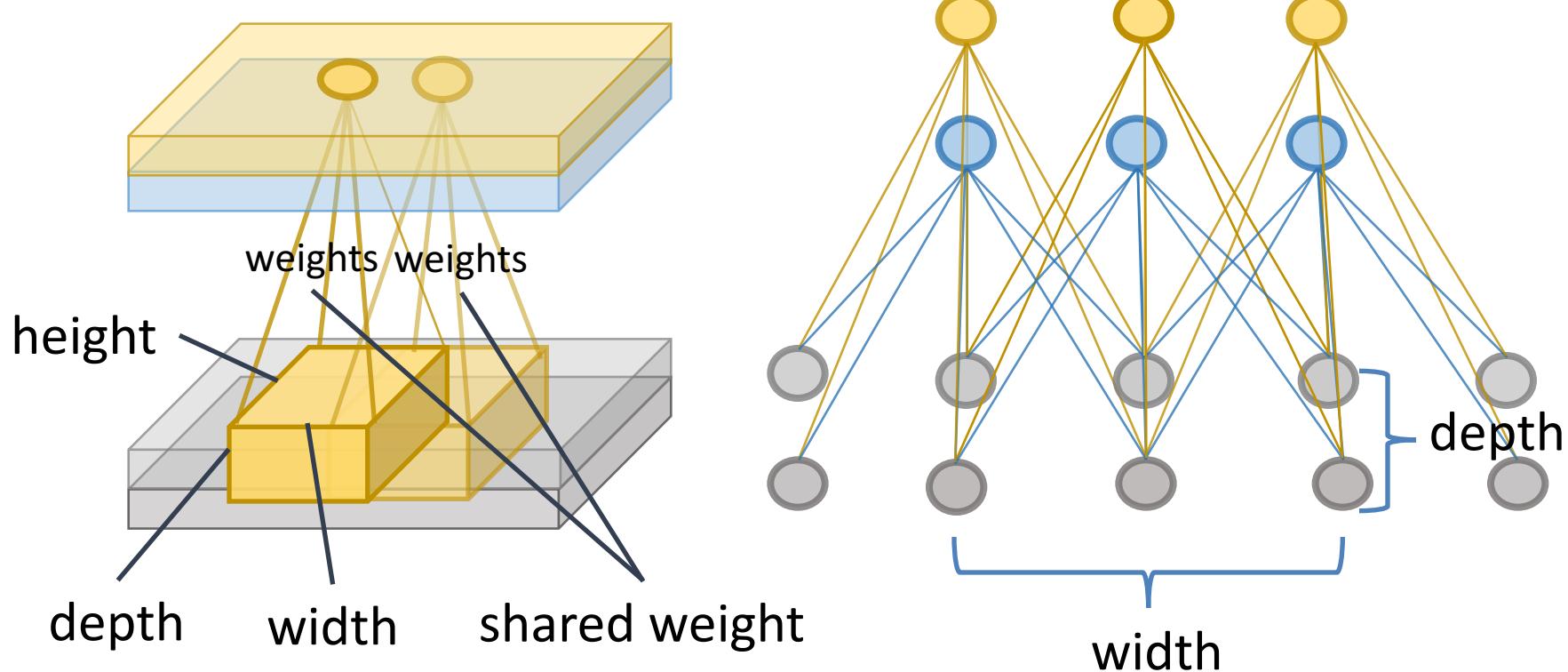


Parameter Sharing

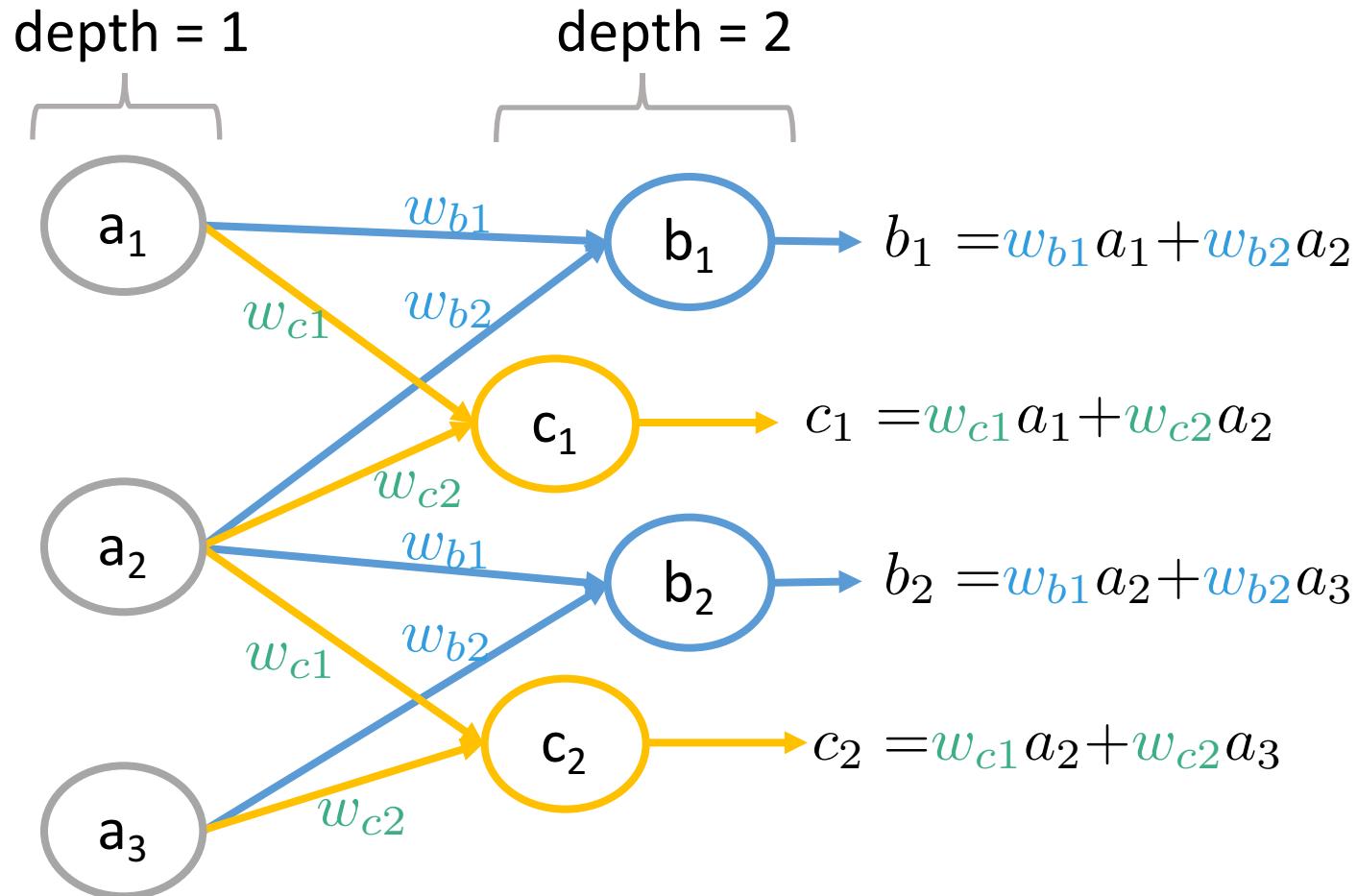
- Different features in the same position



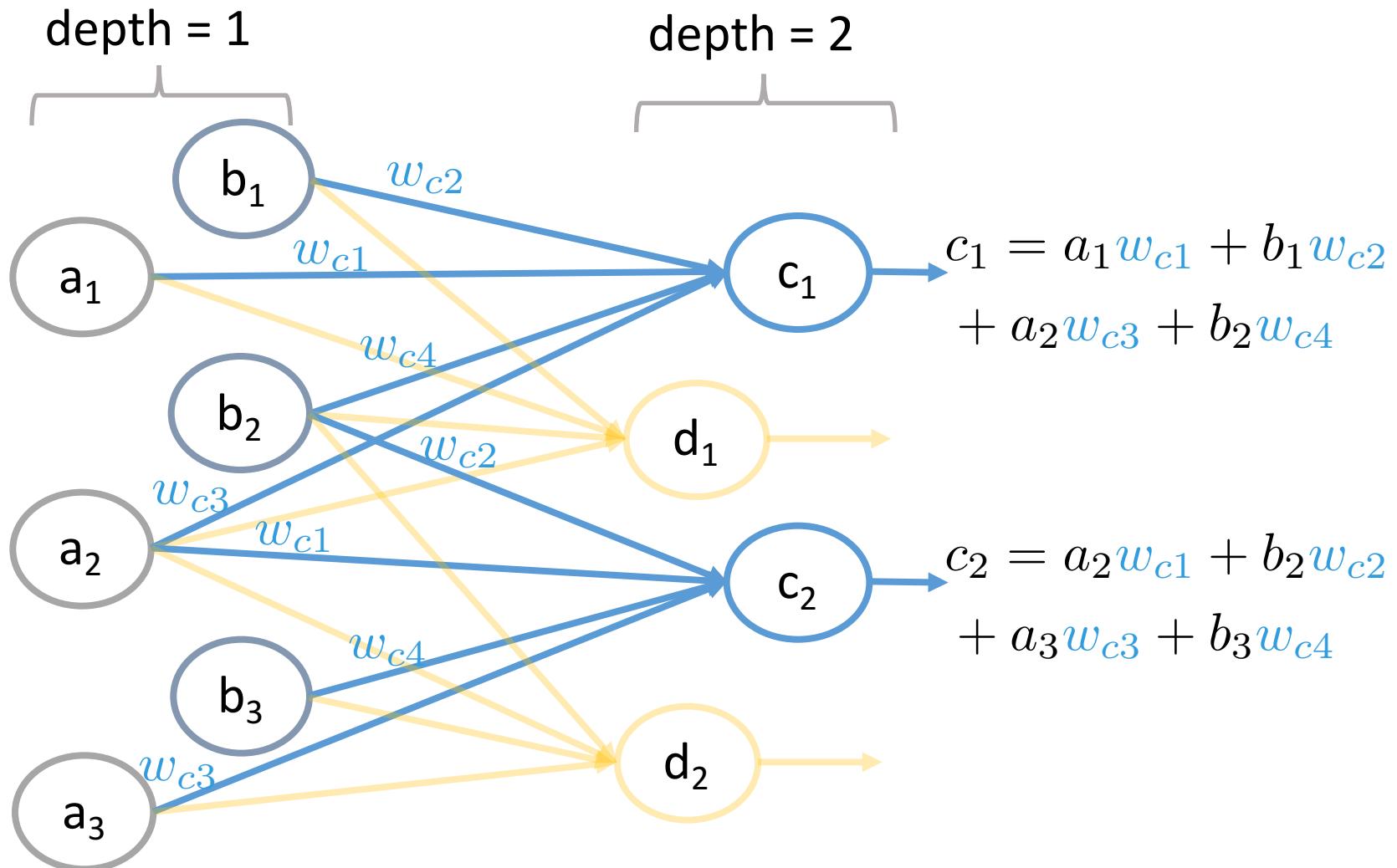
Convolutional Layers: model local contexts.
Sliding window in an image / sentence.



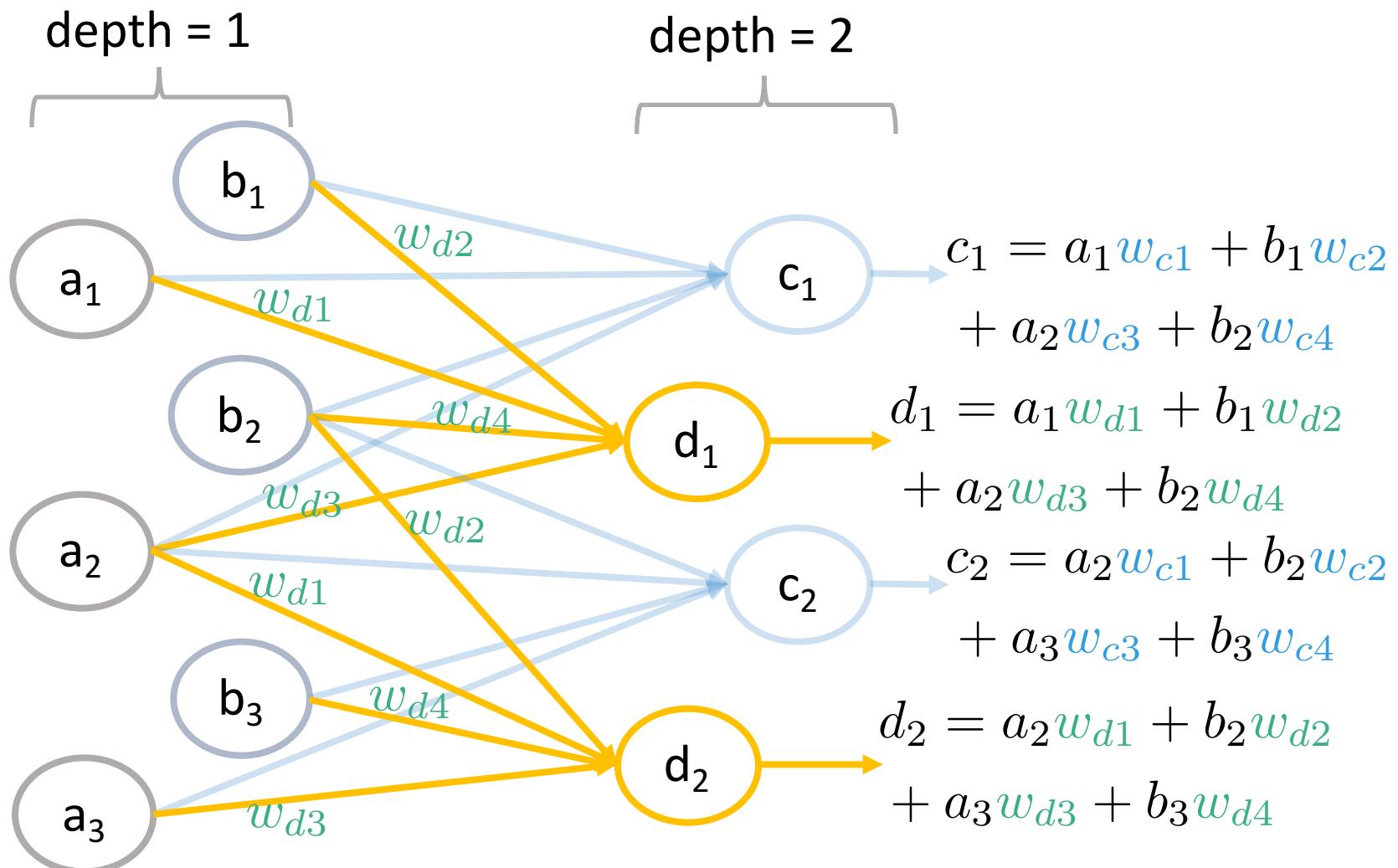
Convolutional Layers: Window Size of 2.



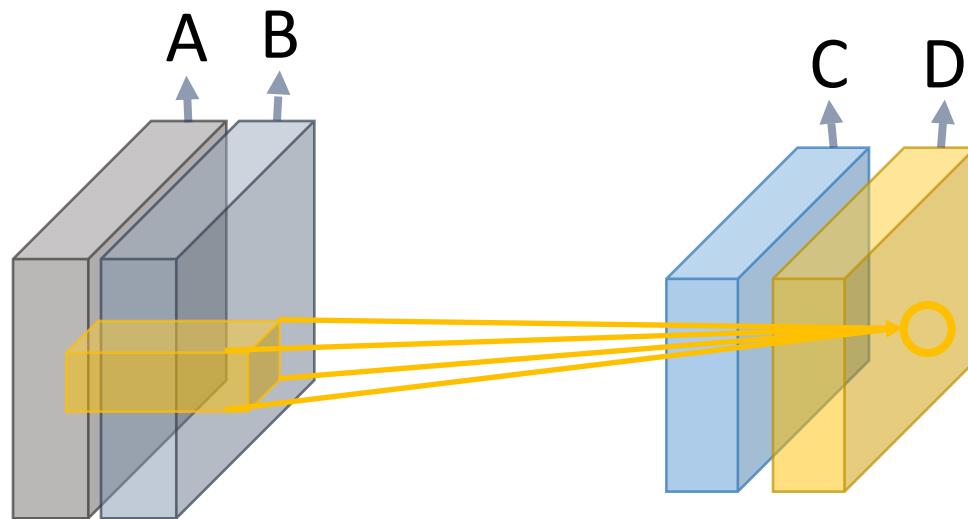
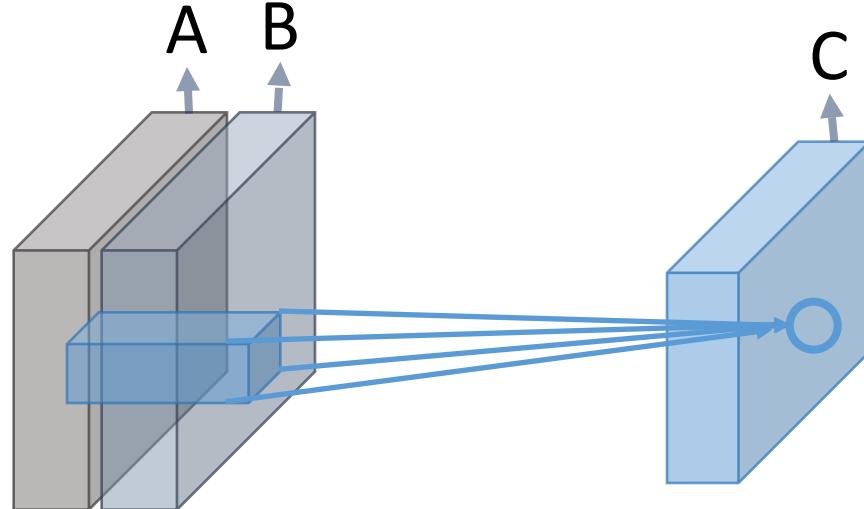
Convolutional Layers: Window Size of 4.



Convolutional Layers: Window Size of 4.

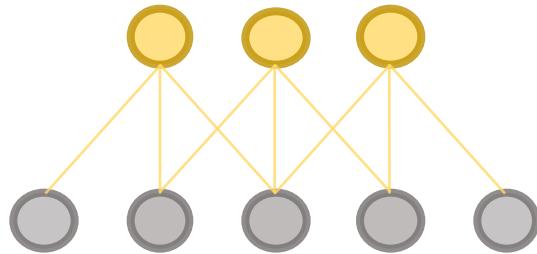


Convolutional Layers



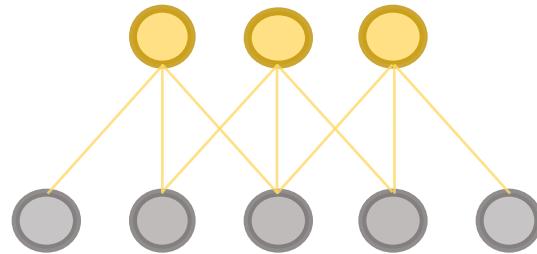
Hyper-parameters of CNN

- Stride

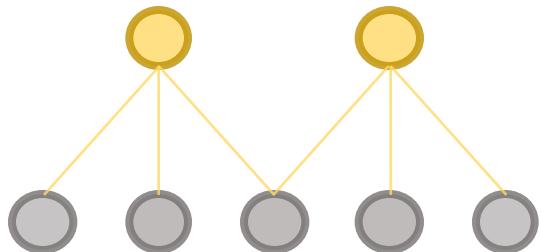


Stride = 1

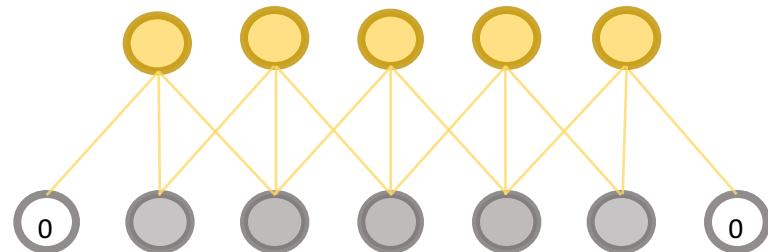
- Padding



Padding = 0



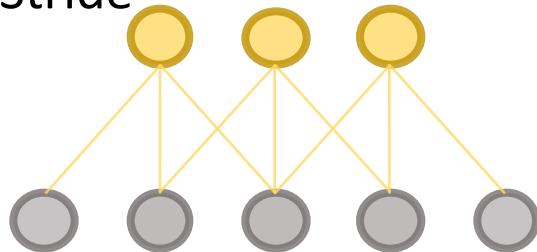
Stride = 2



Padding = 1

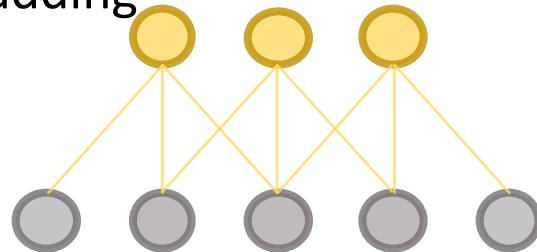
Piazza Poll: Why do we need padding?

- Stride

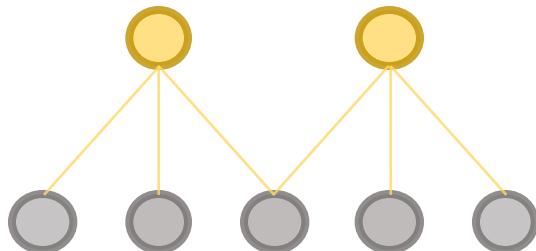


Stride = 1

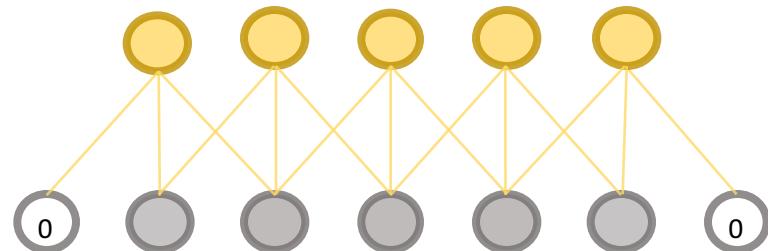
- Padding



Padding = 0



Stride = 2



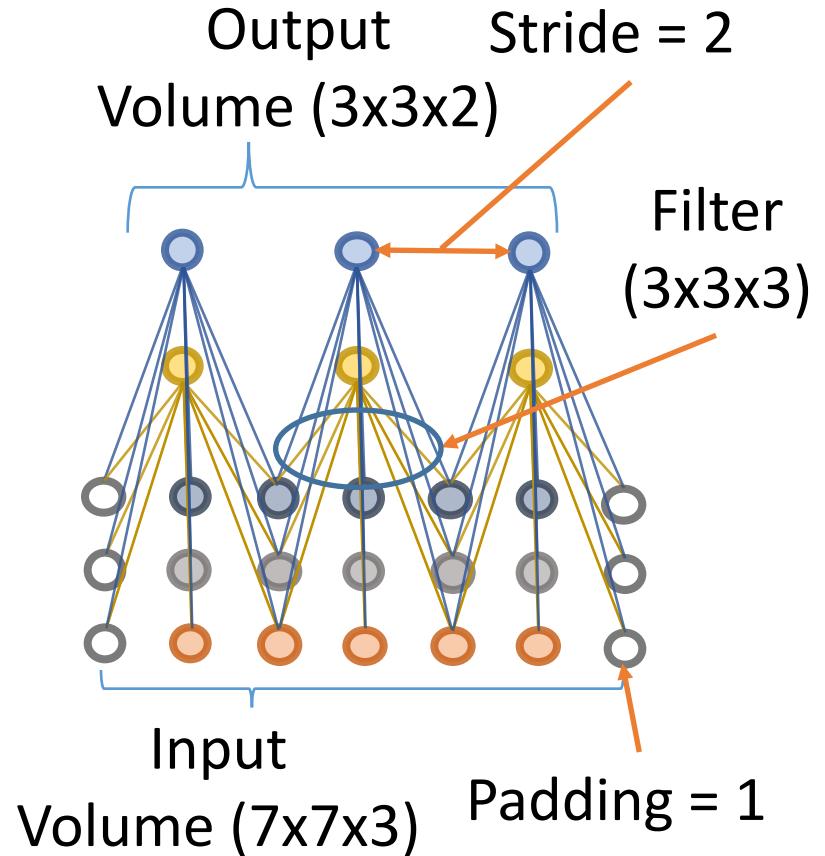
Padding = 1

CNN Example

Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)
$x[:, :, 0]$	$w0[:, :, 0]$
0 0 0 0 0 0 0 0 0 2 0 0 1 0 0 0 0 0 1 0 2 1 0 0 0 0 1 1 0 2 0 0 0 1 0 2 0 1 0 0 0 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0	-1 -1 0 -1 1 1 1 -1 1 0 1 1 -1 1 1 -1 -1 -1 0 1 1
$x[:, :, 1]$	$w0[:, :, 1]$
0 0 0 0 0 0 0 0 0 0 0 1 2 0 0 0 0 2 1 1 1 1 0 0 0 0 0 0 2 1 0 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	0 1 1 0 0 -1 0 1 1
$x[:, :, 2]$	$w0[:, :, 2]$
0 0 0 0 0 0 0 0 0 1 1 2 0 2 0 0 0 1 1 1 2 1 0 0 0 0 2 0 0 0 2 0 0 1 1 0 2 0 0 0 0 0 1 2 2 0 0 0 0 0 0 0 0 0 0 0	1

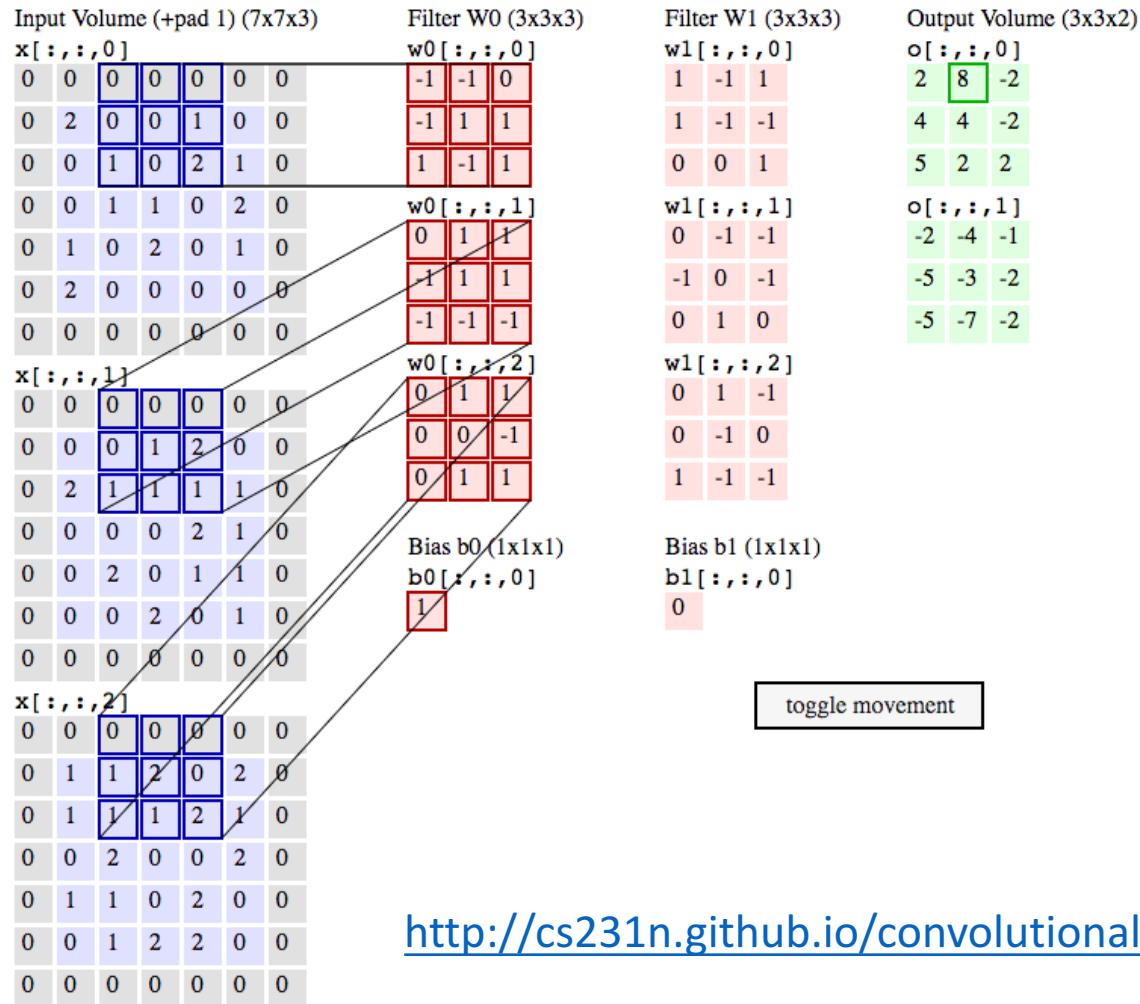
Filter W1 (3x3x3)	Output Volume (3x3x2)
$w1[:, :, 0]$	$o[:, :, 0]$
1 -1 1 1 -1 -1 0 0 1	2 8 -2 4 4 -2 5 2 2
$w1[:, :, 1]$	$o[:, :, 1]$
0 -1 -1 -1 0 -1 0 1 0	-2 -4 -1 -5 -3 -2 -5 -7 -2
$w1[:, :, 2]$	
0 1 -1 0 -1 0 1 -1 -1	
Bias b0 (1x1x1) $b0[:, :, 0]$	0

toggle movement

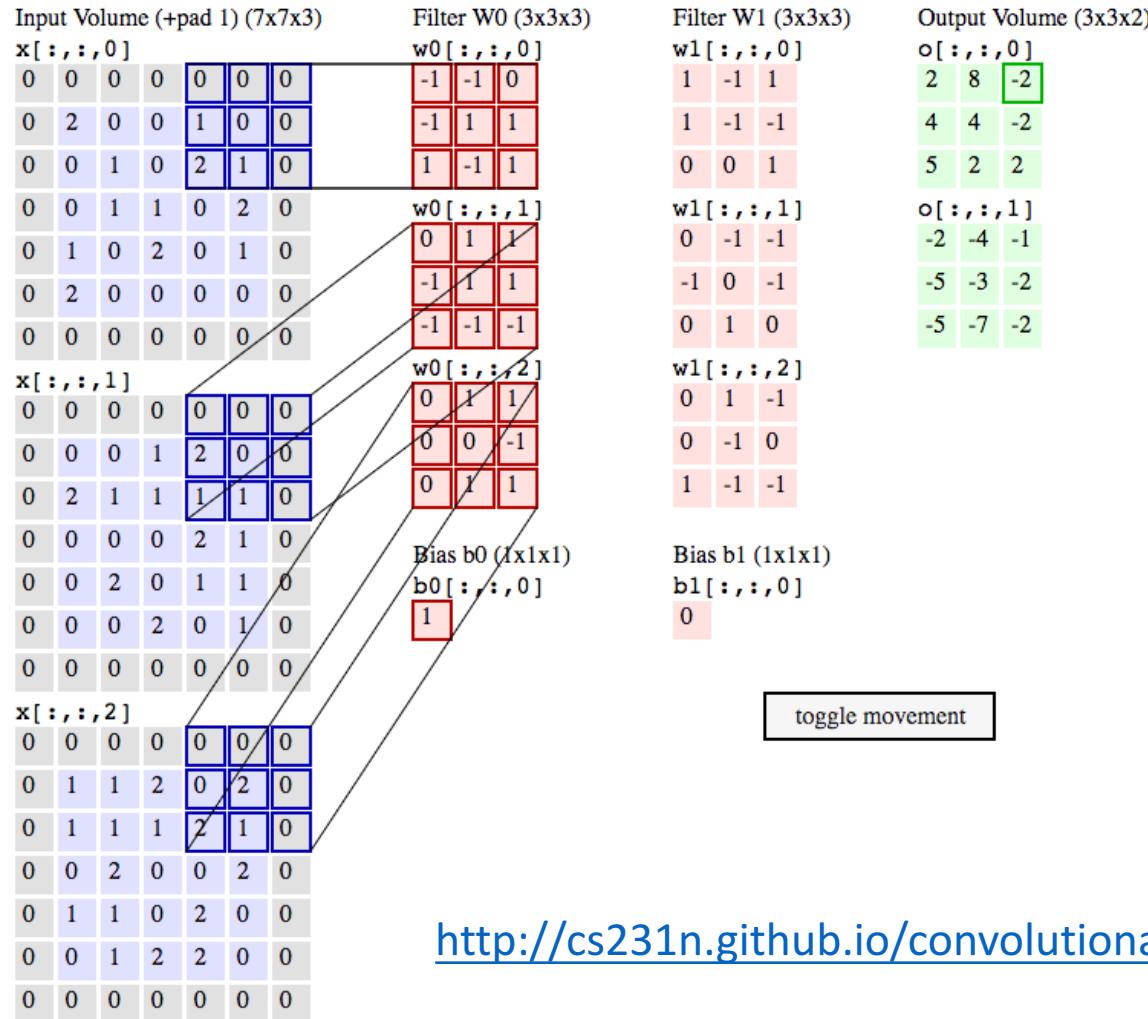


<http://cs231n.github.io/convolutional-networks/>

Convolutional Layers

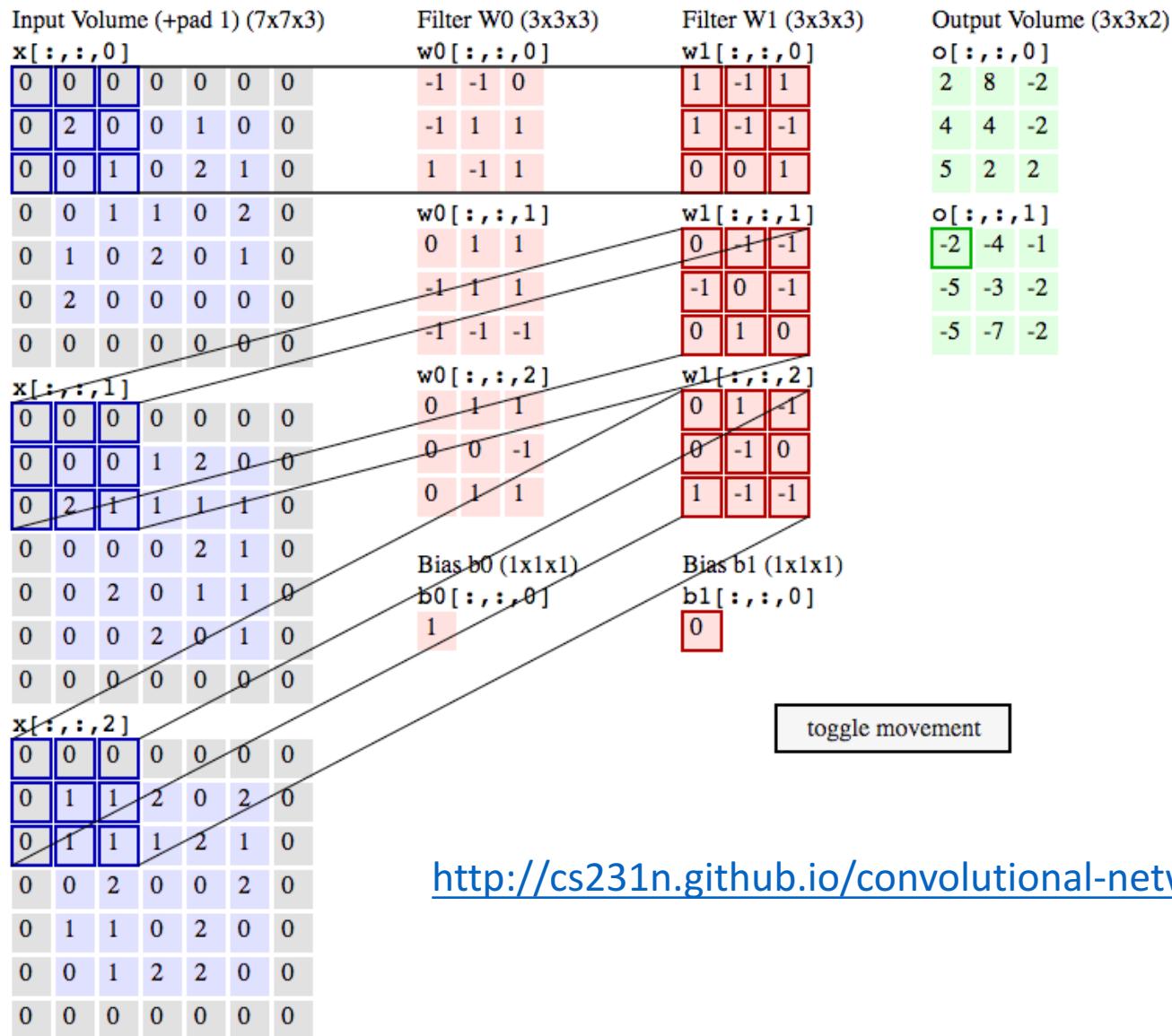


Convolutional Layers



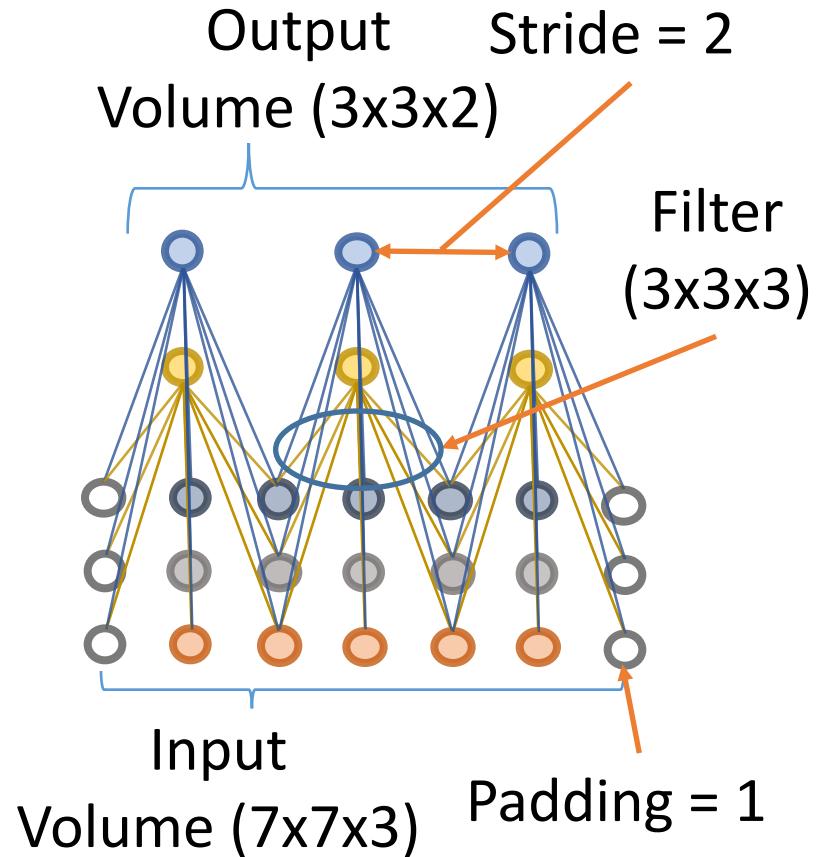
<http://cs231n.github.io/convolutional-networks/>

Convolutional Layers



Quick question: how to change one parameter and generate an output of 5x5x2?

Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)	Filter W1 (3x3x3)	Output Volume (3x3x2)
$x[:, :, 0]$	$w0[:, :, 0]$	$w1[:, :, 0]$	$o[:, :, 0]$
0 0 0 0 0 0 0 0 2 0 0 1 0 0 0 0 1 0 2 1 0 0 0 1 1 0 2 0 0 1 0 2 0 1 0 0 2 0 0 0 0 0 0 0 0 0 0 0 0	-1 -1 0 -1 1 1 1 -1 1 0 1 1 -1 1 1 -1 -1 -1 0 1 1 0 0 0 0 2 1 1 1 0 0 0 0 0 2 1 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	1 -1 1 1 -1 -1 0 0 1 0 -1 -1 -1 0 -1 0 1 0 0 1 -1 0 -1 0 1 -1 -1 1	2 8 -2 4 4 -2 5 2 2 -2 -4 -1 -5 -3 -2 -5 -7 -2 0
$x[:, :, 1]$	$w0[:, :, 1]$	$w1[:, :, 1]$	$o[:, :, 1]$
0 0 0 0 0 0 0 0 0 0 1 2 0 0 0 2 1 1 1 1 0 0 0 0 0 2 1 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	0 1 1 0 0 -1 0 1 1 0 1 1 0 0 -1 1 -1 -1 0 1 1 0 0 0 0 2 1 1 1 0 0 0 0 0 2 1 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	0 1 -1 0 -1 0 1 -1 -1 0 1 0 0 1 -1 0 -1 0 1 -1 -1 0 1 0 0 0 0 0 2 1 1 1 0 0 0 0 0 2 1 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	-2 -4 -1 -5 -3 -2 -5 -7 -2 0
$x[:, :, 2]$	$w0[:, :, 2]$	$w1[:, :, 2]$	
0 0 0 0 0 0 0 0 1 1 2 0 2 0 0 1 1 1 2 1 0 0 0 2 0 0 2 0 0 1 1 0 2 0 0 0 0 1 2 2 0 0 0 0 0 0 0 0 0	0 1 1 0 0 -1 1 -1 -1 1	0 1 -1 0 -1 0 1 -1 -1 0 1 0 0 1 -1 0 -1 0 1 -1 -1 0 1 0 0 0 0 0 2 1 1 1 0 0 0 0 0 2 1 0 0 0 2 0 1 1 0 0 0 0 2 0 1 0 0 0 0 0 0 0 0	



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

Pooling Layer

1	3	2	4
5	7	6	8
0	0	3	3
5	5	0	0

Maximum
Pooling



7	8
5	3

Average
Pooling

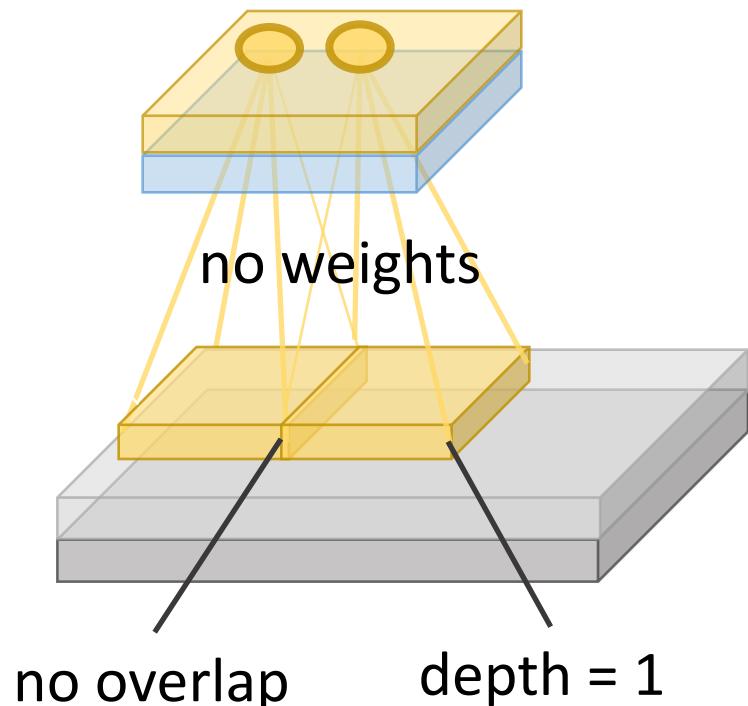


4	5
5	3

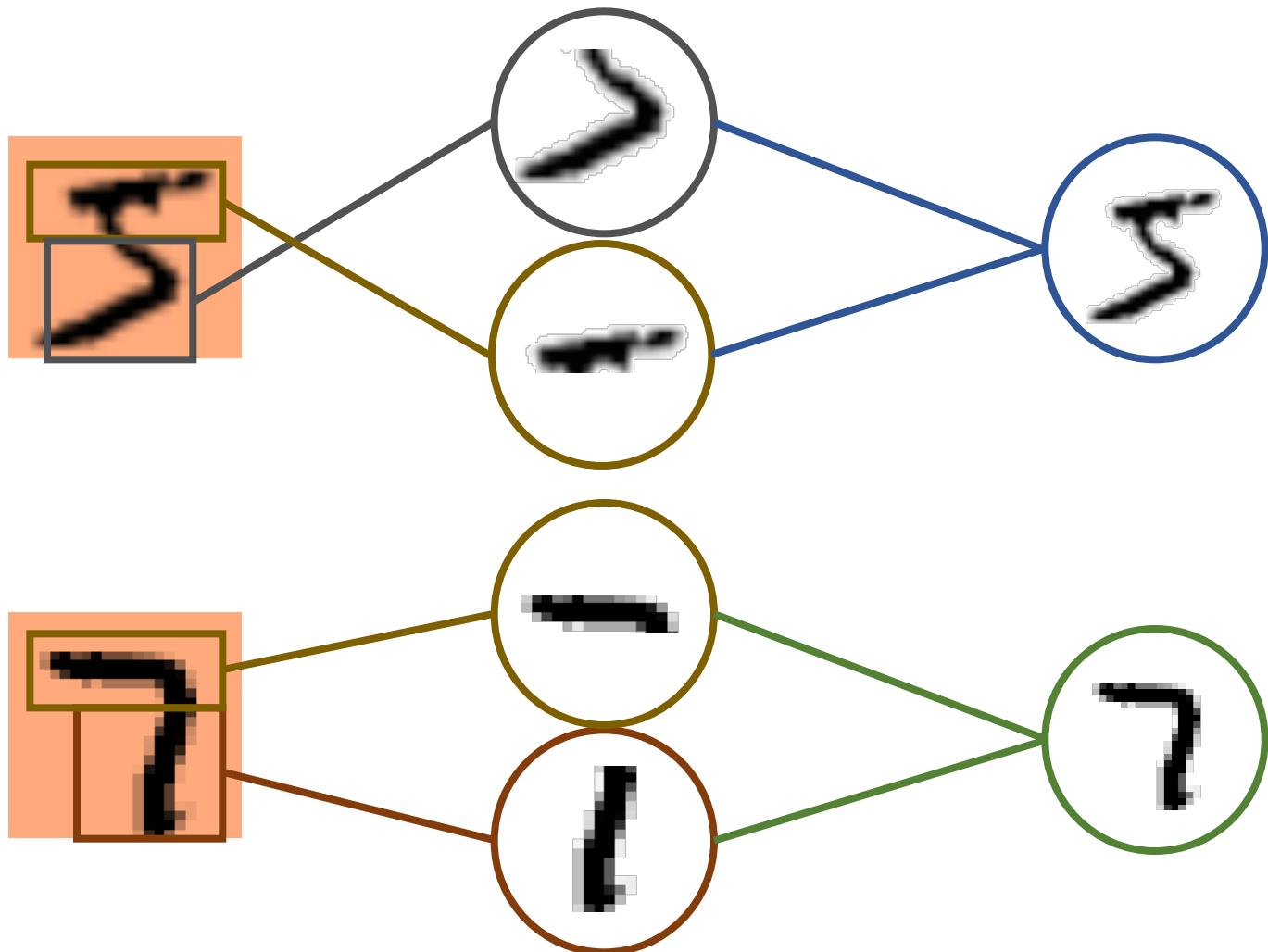
$$\text{Max}(1,3,5,7) = 7$$

$$\text{Avg}(1,3,5,7) = 4$$

$$\text{Max}(0,0,5,5) = 5$$

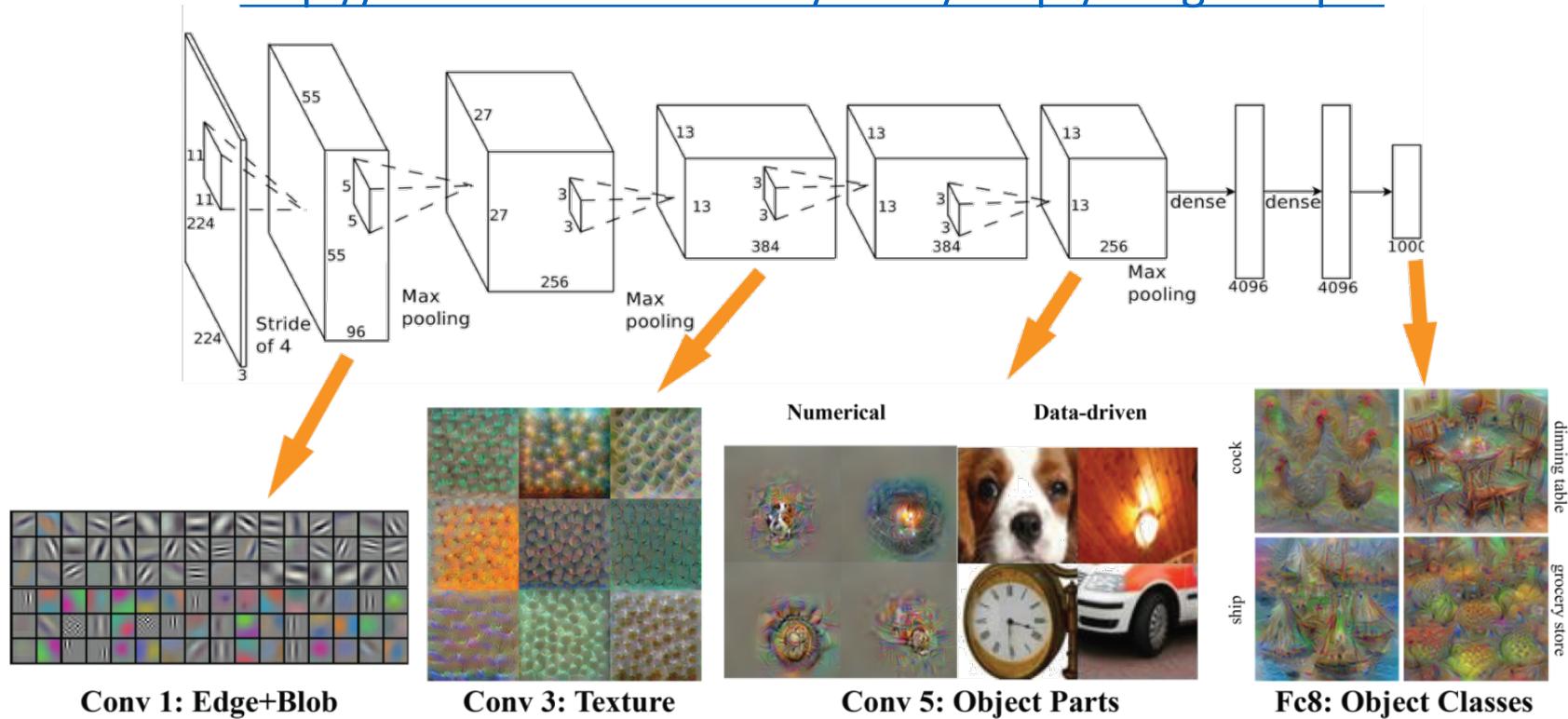


Why “Deep” Learning?



AlexNet

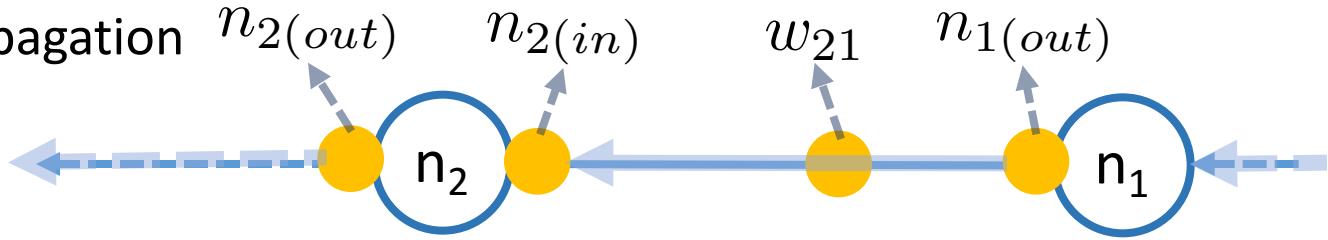
- Alexnet <http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>



http://vision03.csail.mit.edu/cnn_art/data/single_layer.png

Training

- Forward Propagation

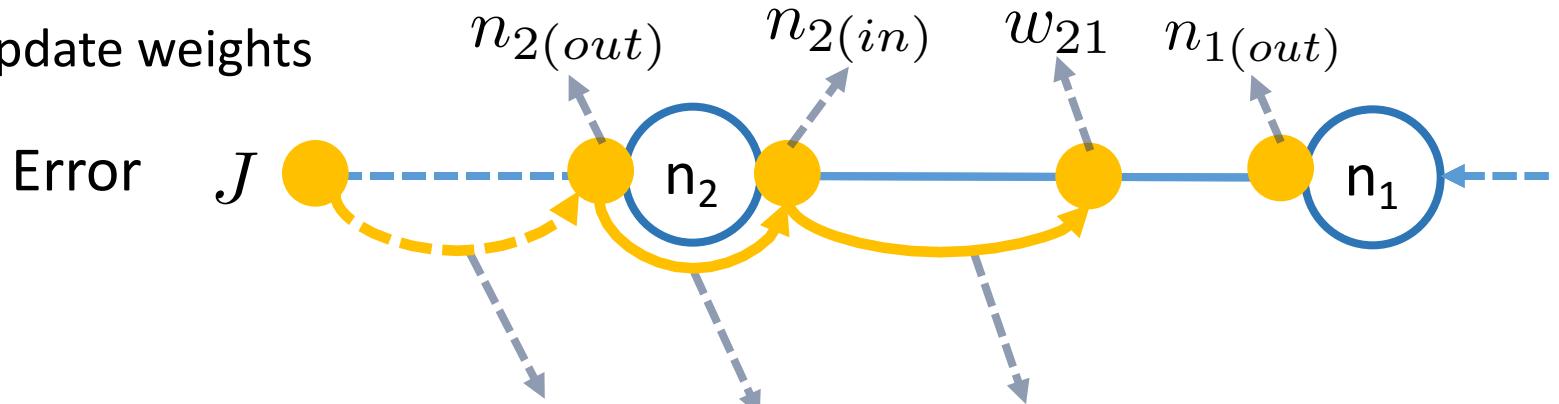


$$n_2(in) = w_{21} n_1(out)$$

$$n_2(out) = g(n_2(in)), \quad g \text{ is activation function}$$

Training

- Update weights



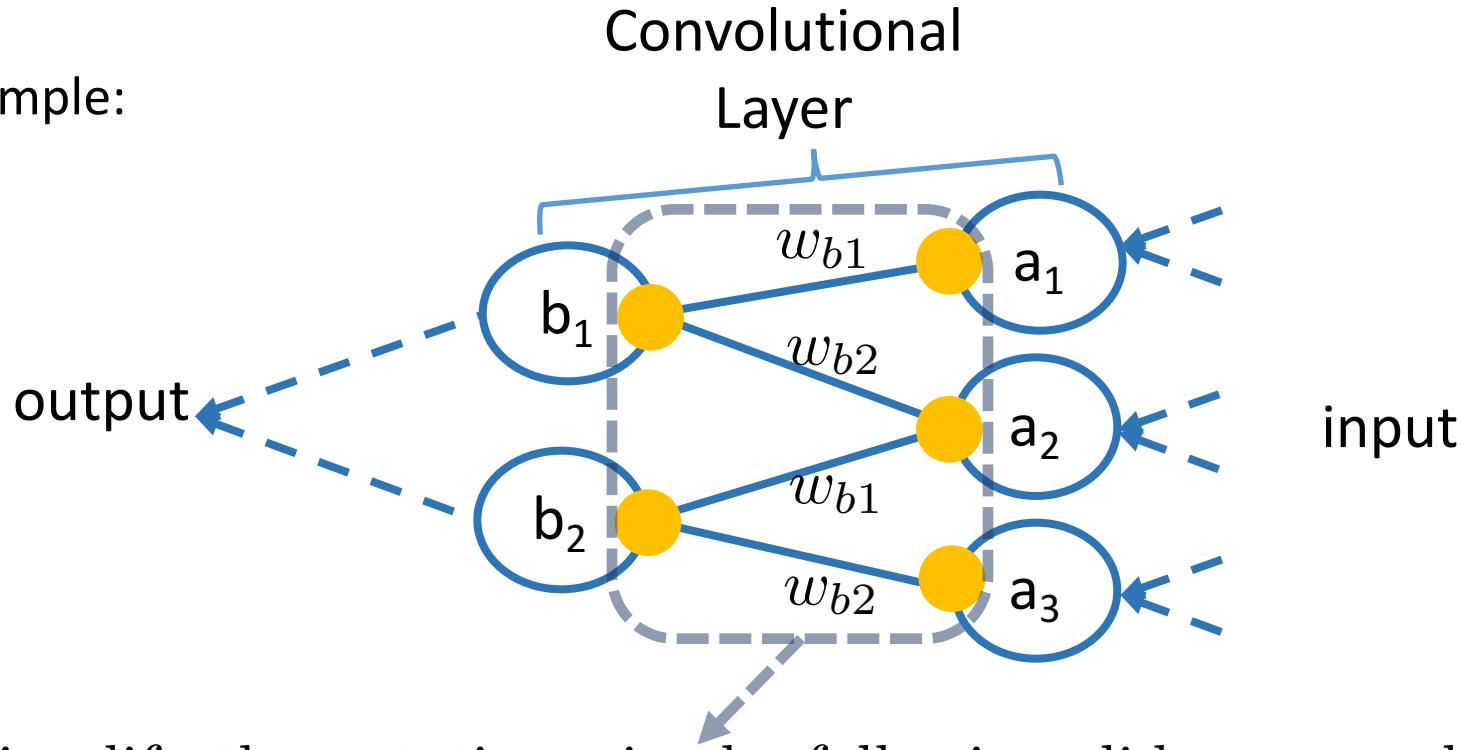
$$\frac{\partial J}{\partial w_{21}} = \frac{\partial J}{\partial n_{2(out)}} \frac{\partial n_{2(out)}}{\partial n_{2(in)}} \frac{\partial n_{2(in)}}{\partial w_{21}}$$

$$w_{21} \leftarrow w_{21} - \eta \frac{\partial J}{\partial w_{21}}$$

$$\Rightarrow w_{21} \leftarrow w_{21} - \eta \frac{\partial J}{\partial n_{2(out)}} \frac{\partial n_{2(out)}}{\partial n_{2(in)}} \frac{\partial n_{2(in)}}{\partial w_{21}}$$

Training Convolutional Layers

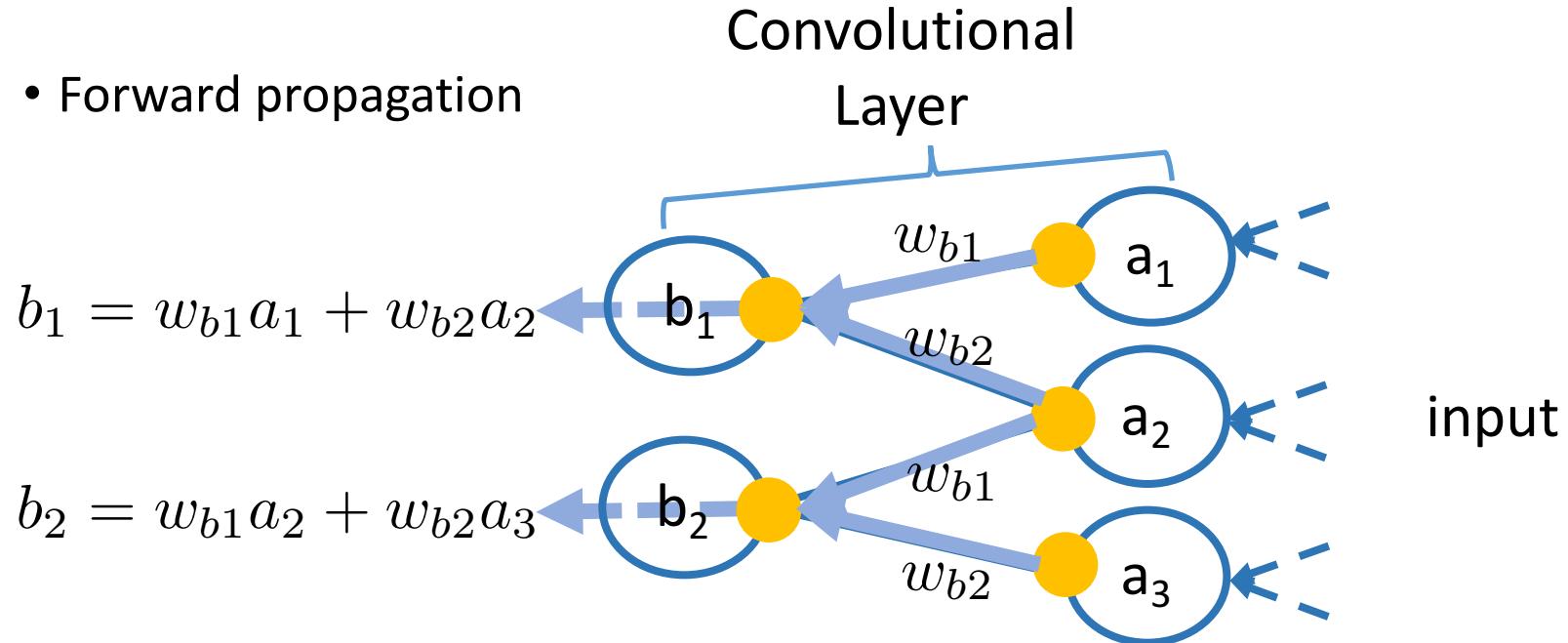
- example:



To simplify the notations, in the following slides, we make:
 b_1 means $b_{1(in)}$, a_1 means $a_{1(out)}$, and so on.

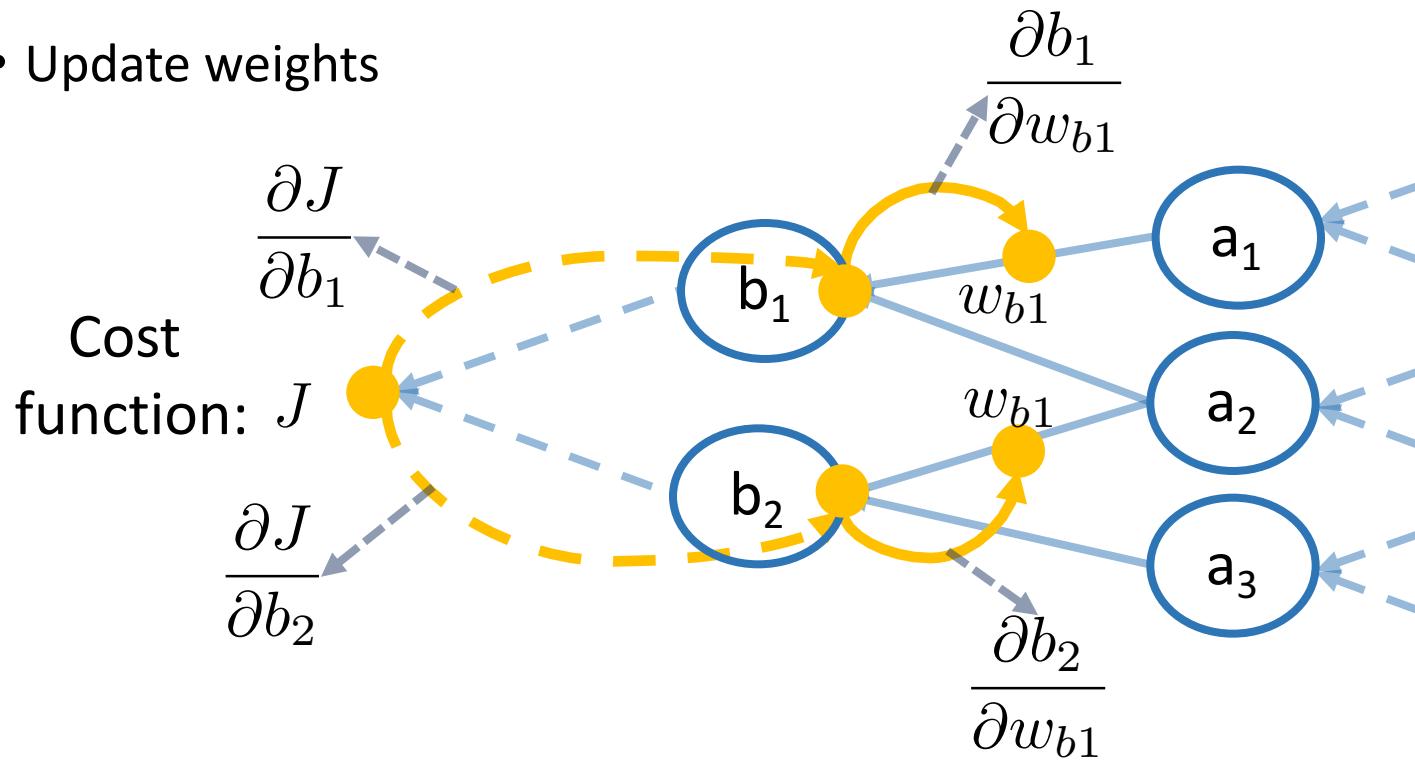
Training Convolutional Layers

- Forward propagation



Training Convolutional Layers

- Update weights

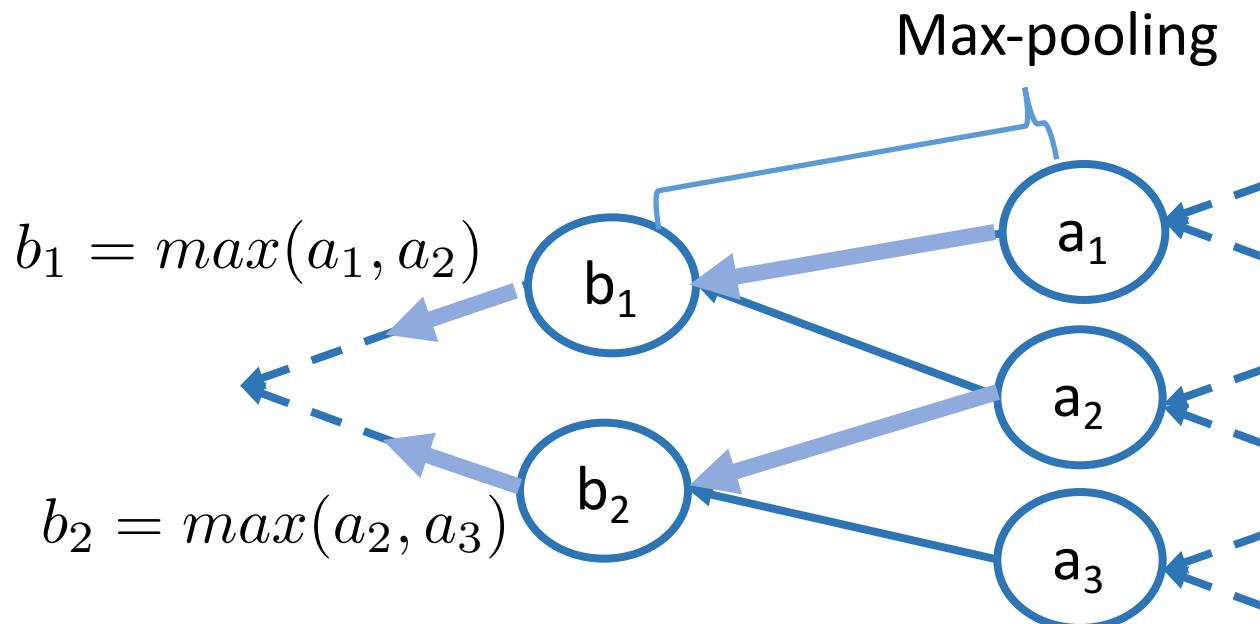


$$w_{b1} \leftarrow w_{b1} - \eta \left(\frac{\partial J}{\partial b_1} \frac{\partial b_1}{\partial w_{b1}} + \frac{\partial J}{\partial b_2} \frac{\partial b_2}{\partial w_{b1}} \right)$$

Piazza Poll:

Max-Pooling Layers during Training

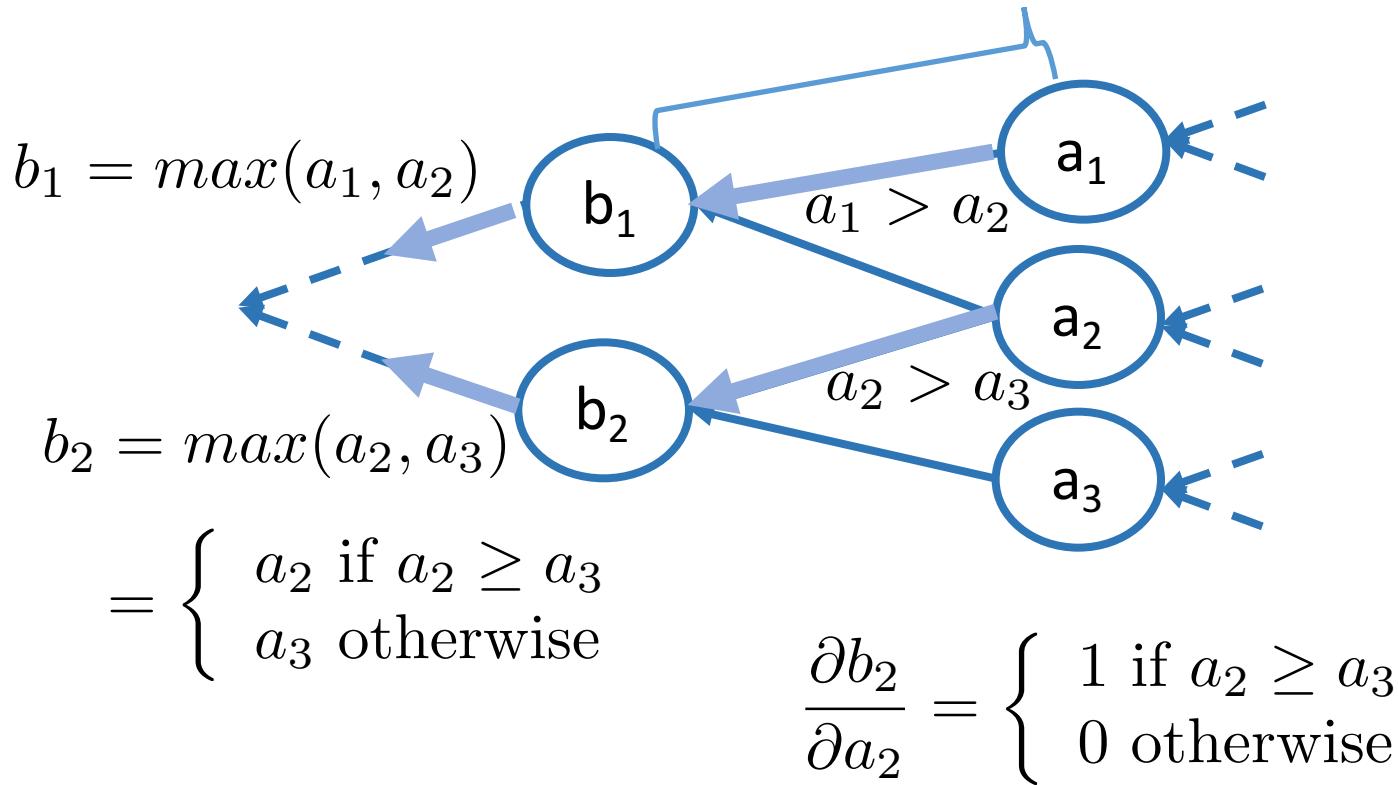
- How to update weights for max-pooling layers?



Max-Pooling Layers during Training

- Pooling layers have no weights
- No need to update weights

Max-pooling

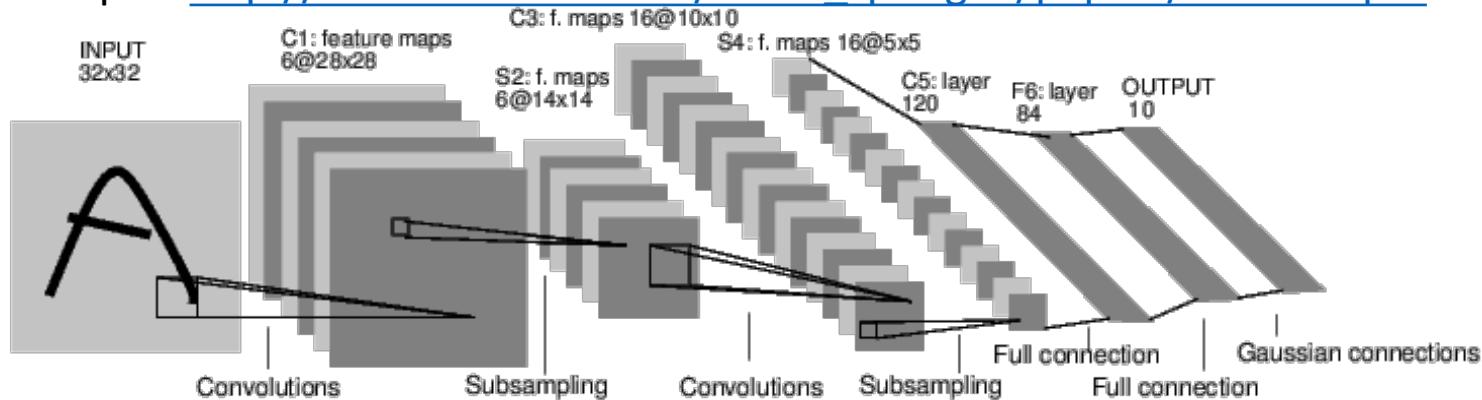


Outline

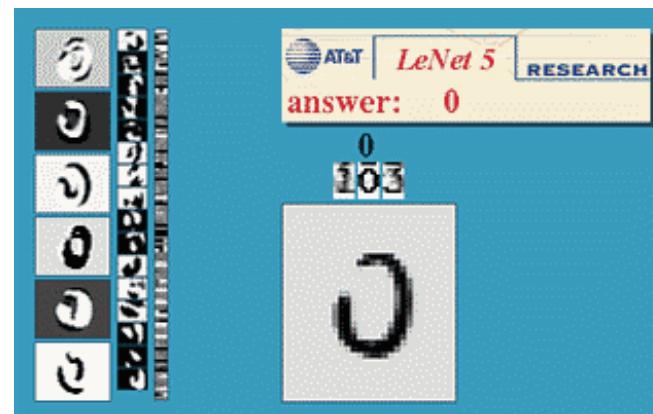
- CNN(Convolutional Neural Networks) Introduction
- **Evolution of CNN**
- Visualizing the Features
- CNN as Artist
- Sentiment Analysis by CNN

LeNet

- Paper: http://vision.stanford.edu/cs598_spring07/papers/Lecun98.pdf



Yann LeCun



<http://yann.lecun.com/exdb/lenet/>

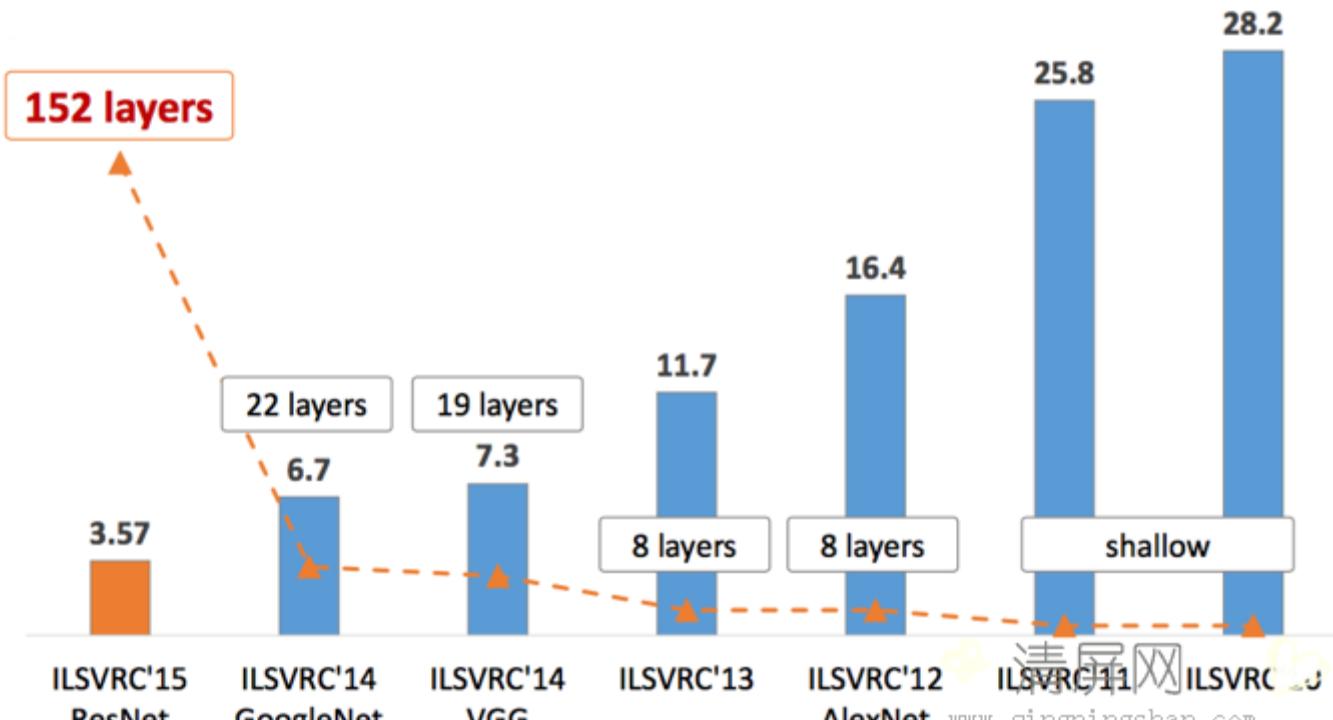
ImageNet Challenge

- ImageNet Large Scale Visual Recognition Challenge
 - <http://image-net.org/challenges/LSVRC/>
- Dataset :
 - 1000 categories
 - Training: 1,200,000
 - Validation: 50,000
 - Testing: 100,000



http://vision.stanford.edu/Datasets/collage_s.png

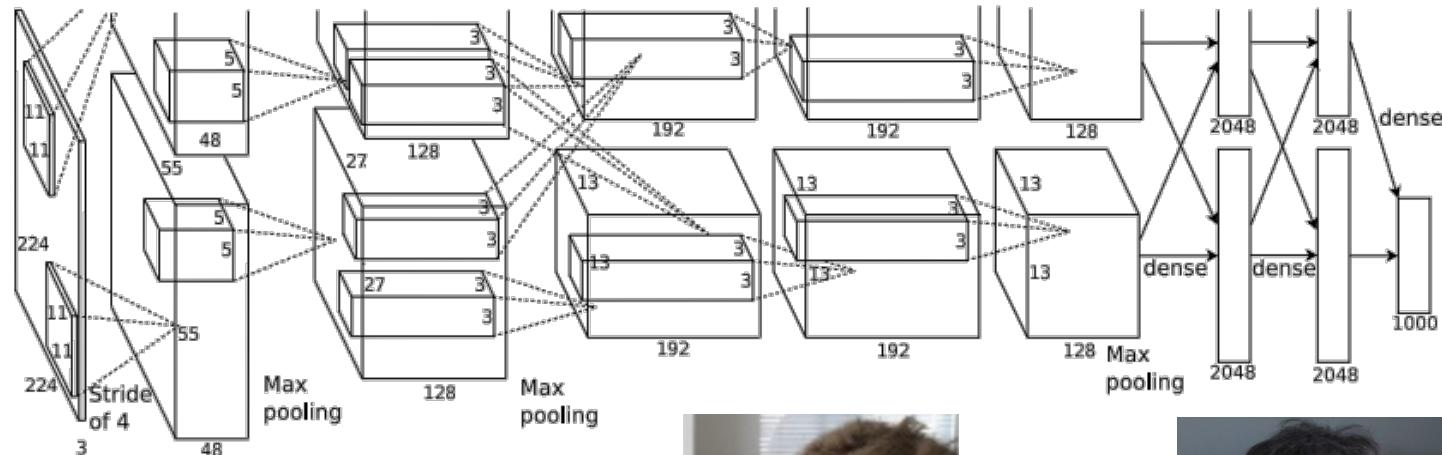
ImageNet Challenge



<http://www.qingpingshan.com/uploads/allimg/160818/1J22QI5-0.png>

AlexNet (2012)

- Paper: <http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf>
- The resurgence of Deep Learning



Alex Krizhevsky Geoffrey Hinton

VGGNet (2014)

- Paper: <https://arxiv.org/abs/1409.1556>

More layers & smaller filters (3x3) is better
More non-linearity, fewer parameters

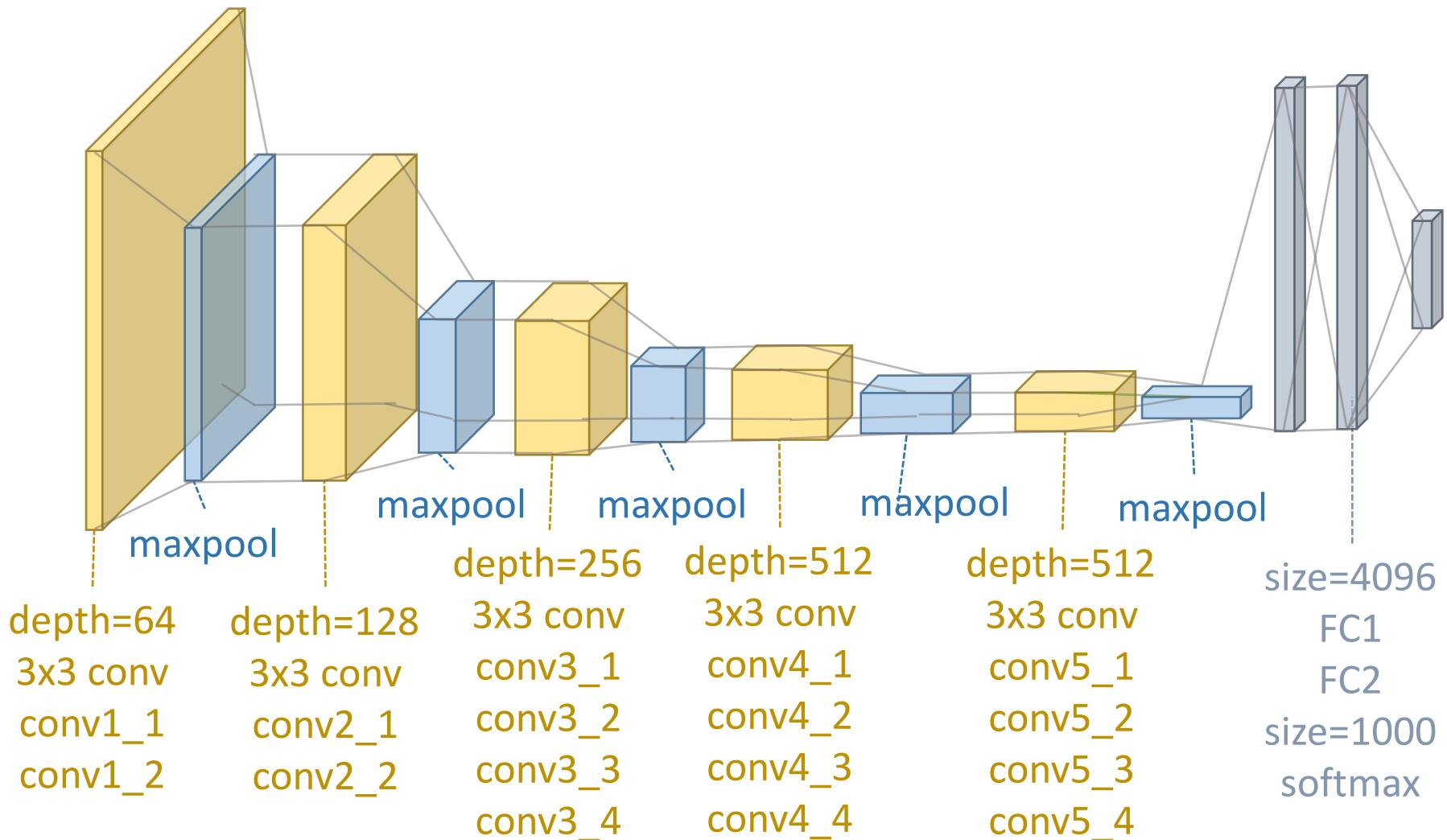
ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

D: VGG16

E: VGG19

All filters are 3x3

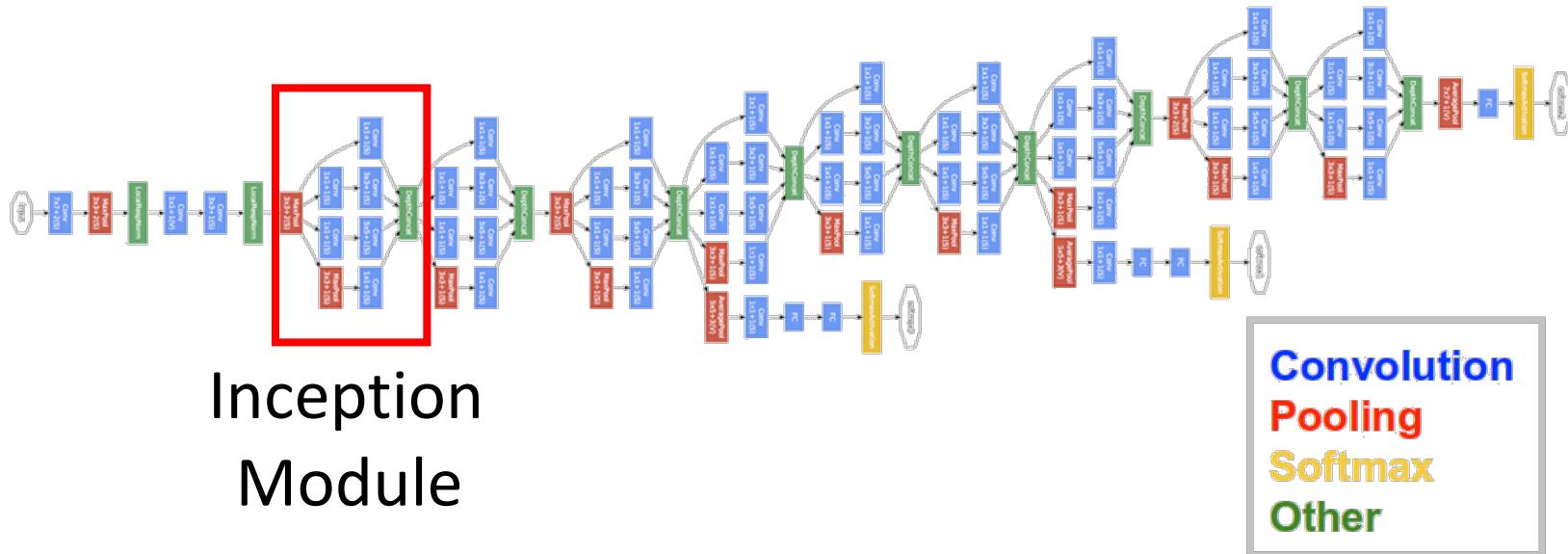
VGG 19



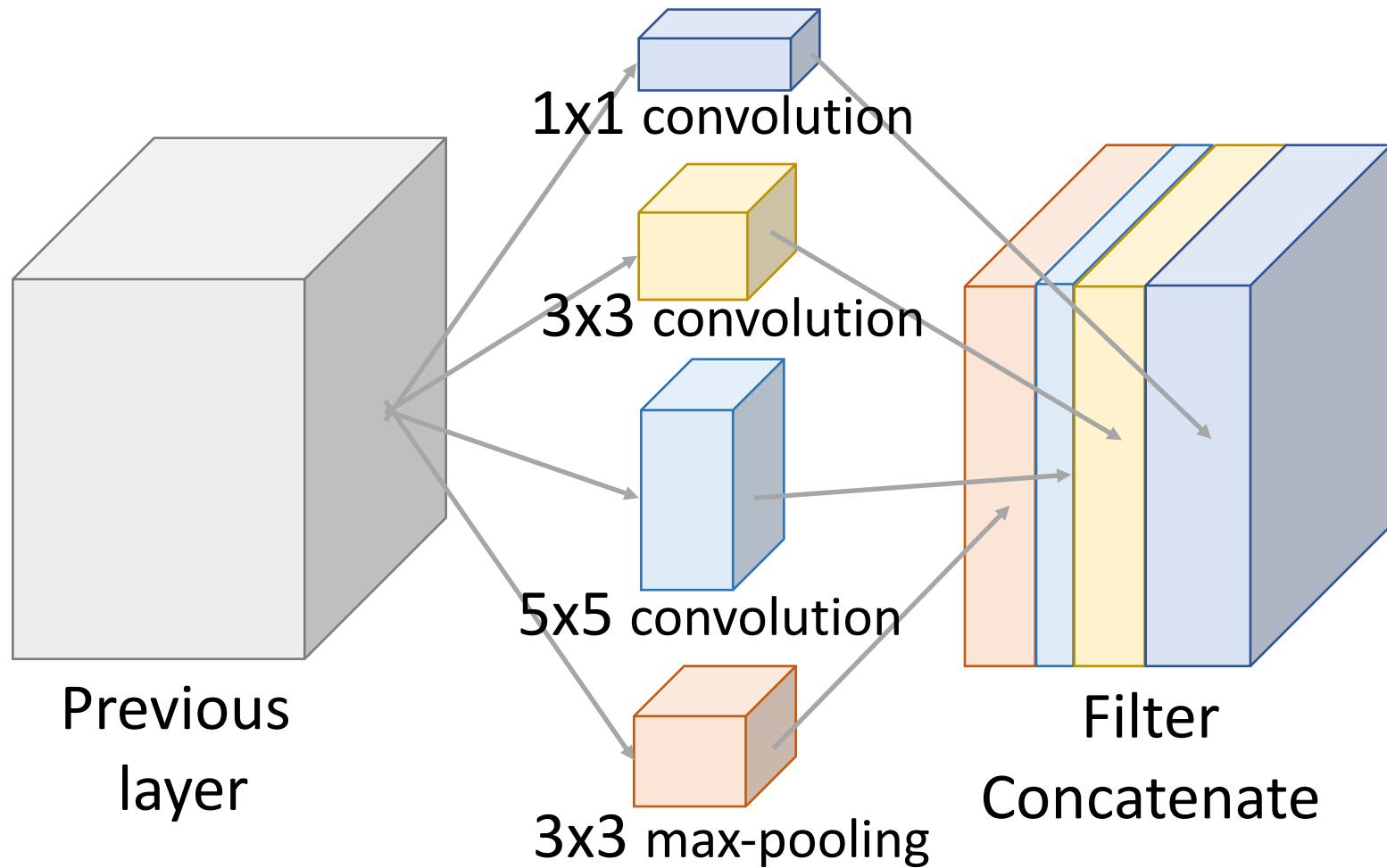
GoogLeNet (2014)

- Paper: <http://www.cs.unc.edu/~wliu/papers/GoogLeNet.pdf>

22 layers deep network

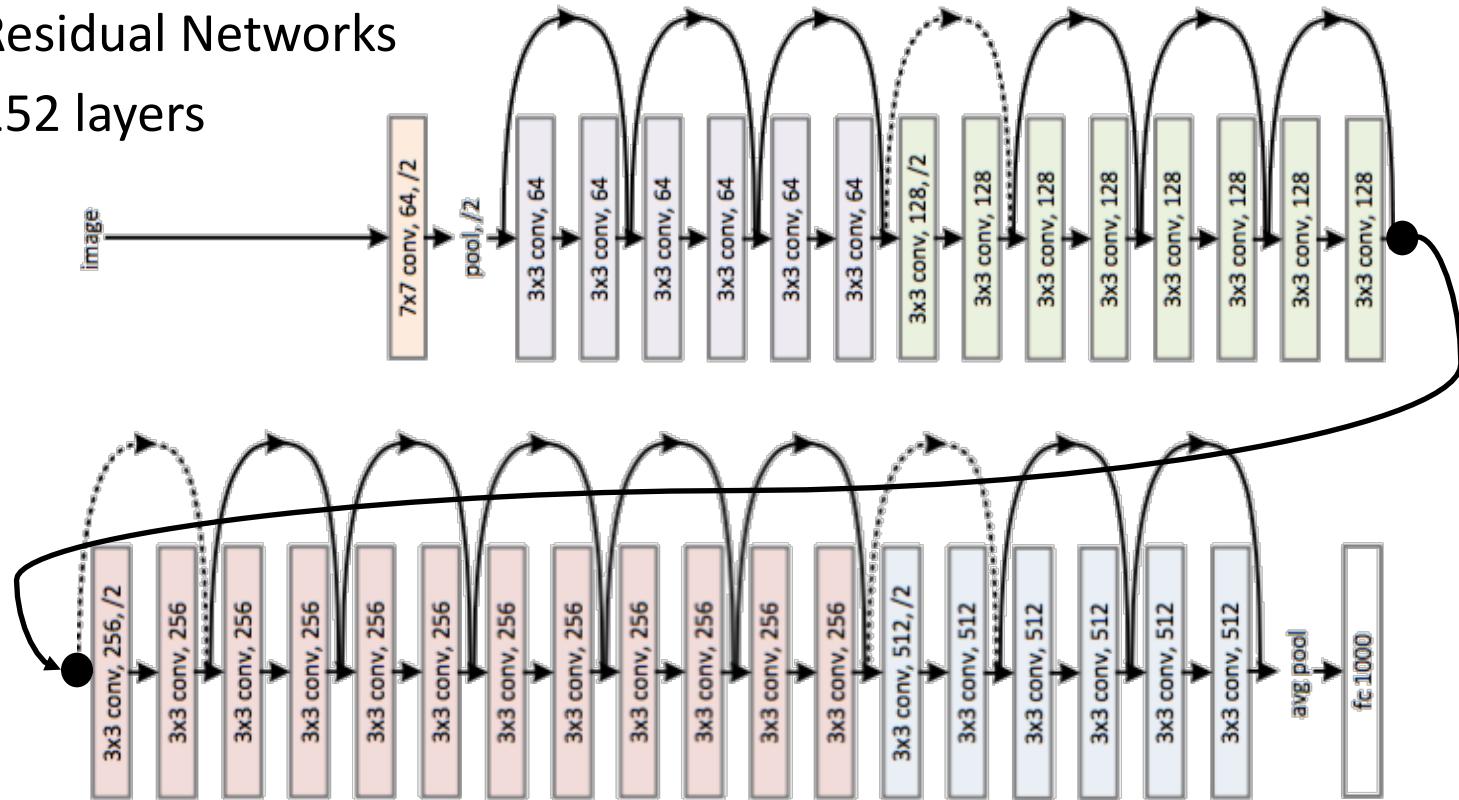


Inception Module



ResNet (2015)

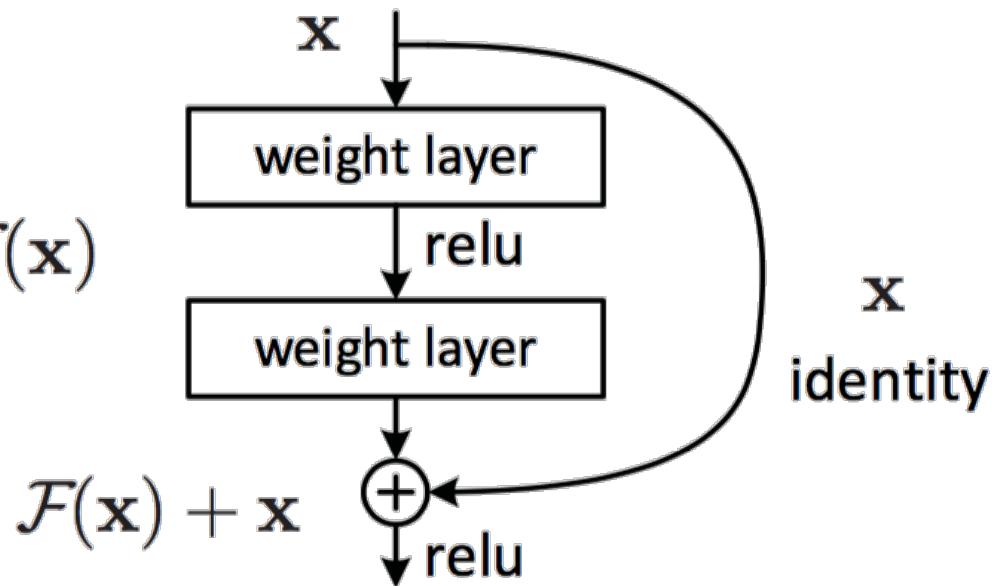
- Paper: <https://arxiv.org/abs/1512.03385>
- Residual Networks
- 152 layers



ResNet

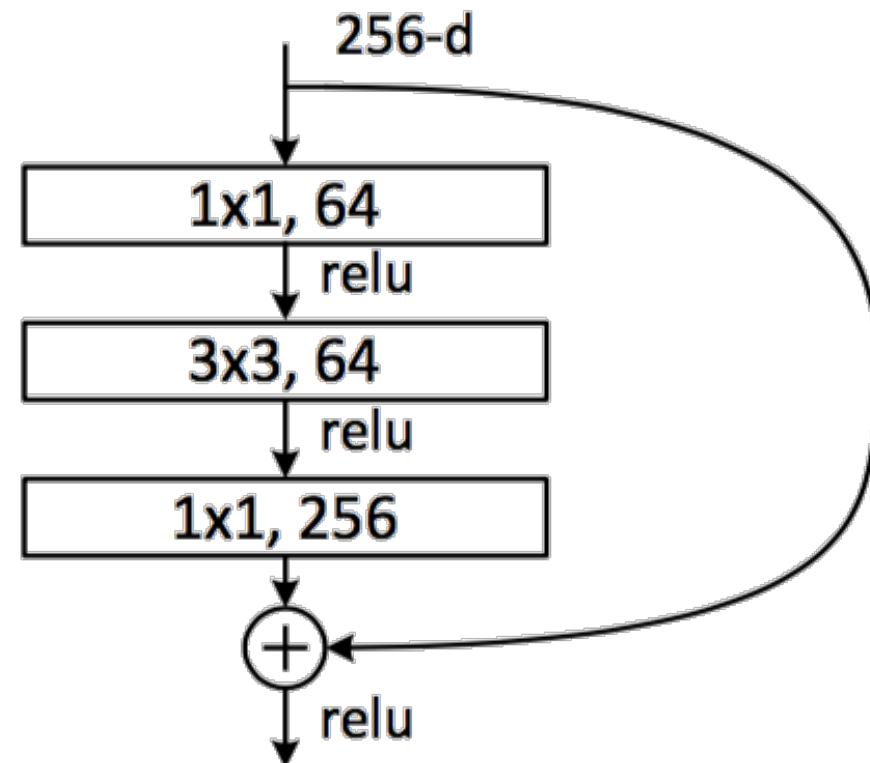
- Residual learning: a building block

Residual
function $\mathcal{F}(x)$



Residual Learning with Dimension Reduction

- using 1x1 filters

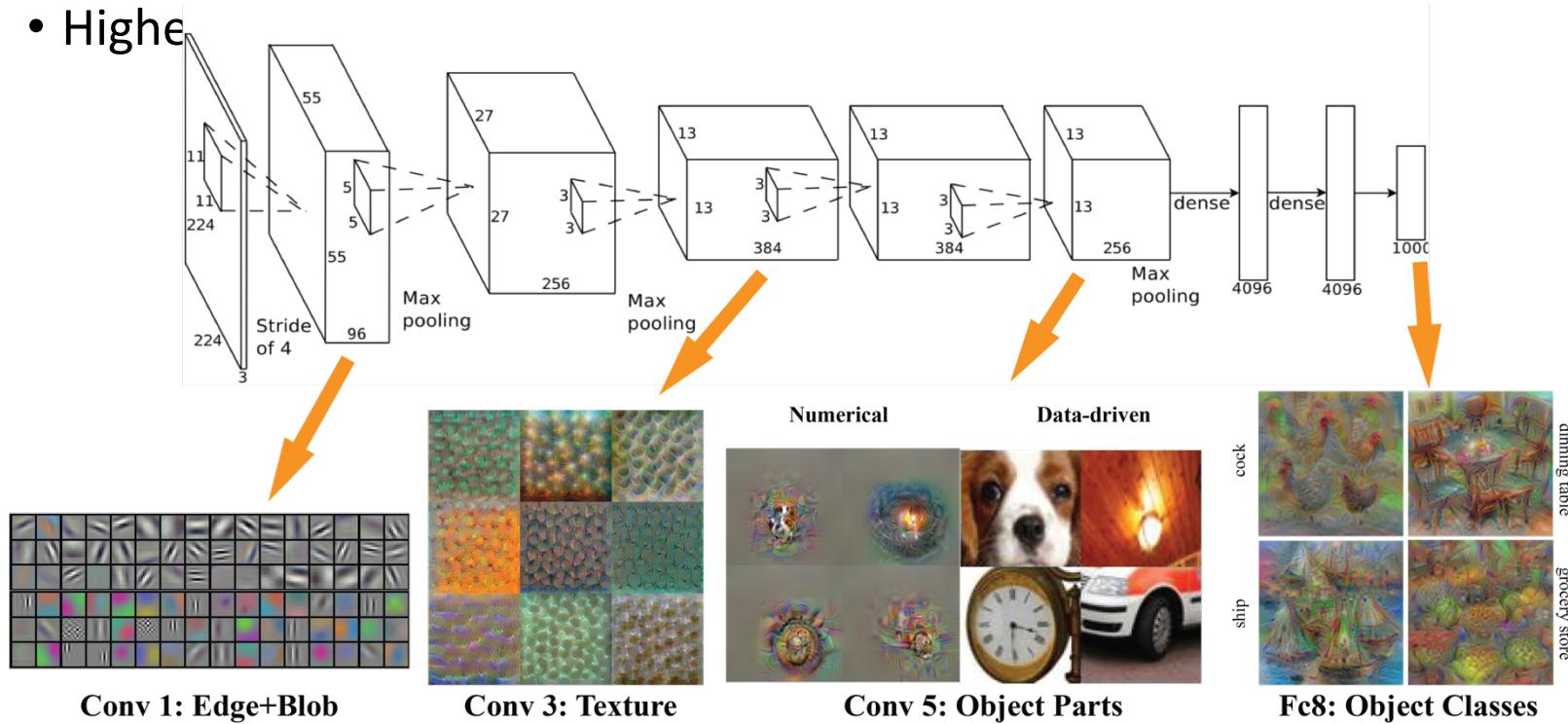


Pretrained Model Download

- <http://www.vlfeat.org/matconvnet/pretrained/>
 - Alexnet:
 - <http://www.vlfeat.org/matconvnet/models/imagenet-matconvnet-alex.mat>
 - VGG19:
 - <http://www.vlfeat.org/matconvnet/models/imagenet-vgg-verydeep-19.mat>
 - GoogLeNet:
 - <http://www.vlfeat.org/matconvnet/models/imagenet-googlenet-dag.mat>
 - ResNet
 - <http://www.vlfeat.org/matconvnet/models/imagenet-resnet-152-dag.mat>

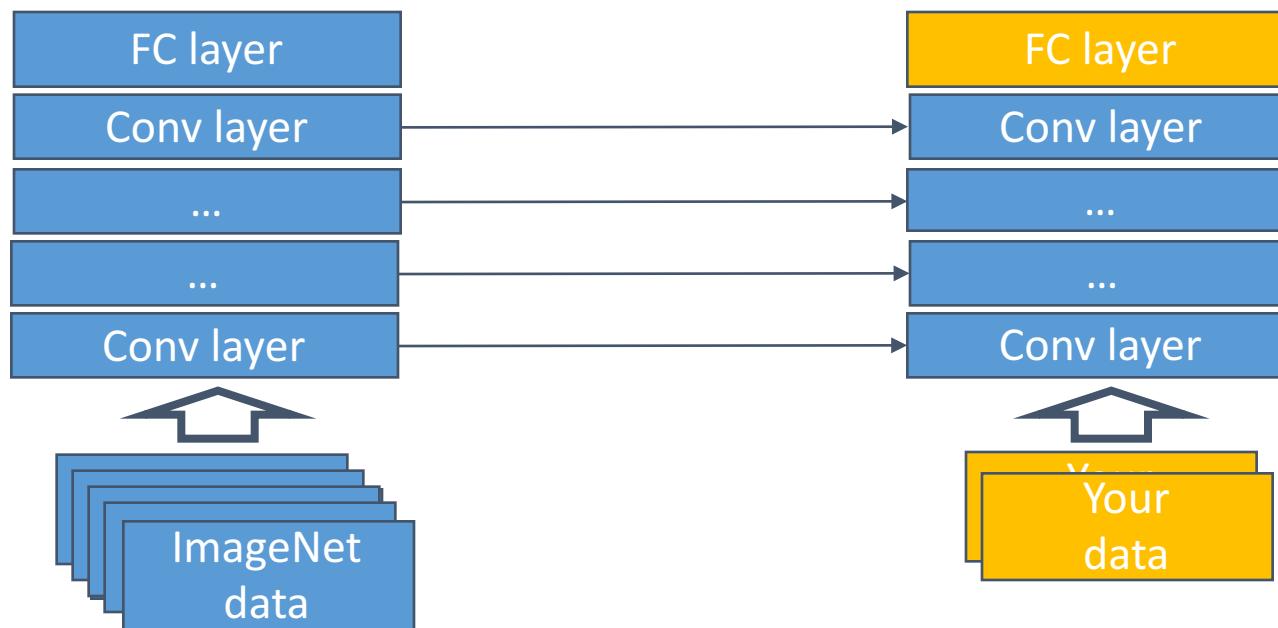
Using Pretrained Model

- Lower layers : edge, blob, texture (more general)
- Higher



Transfer Learning

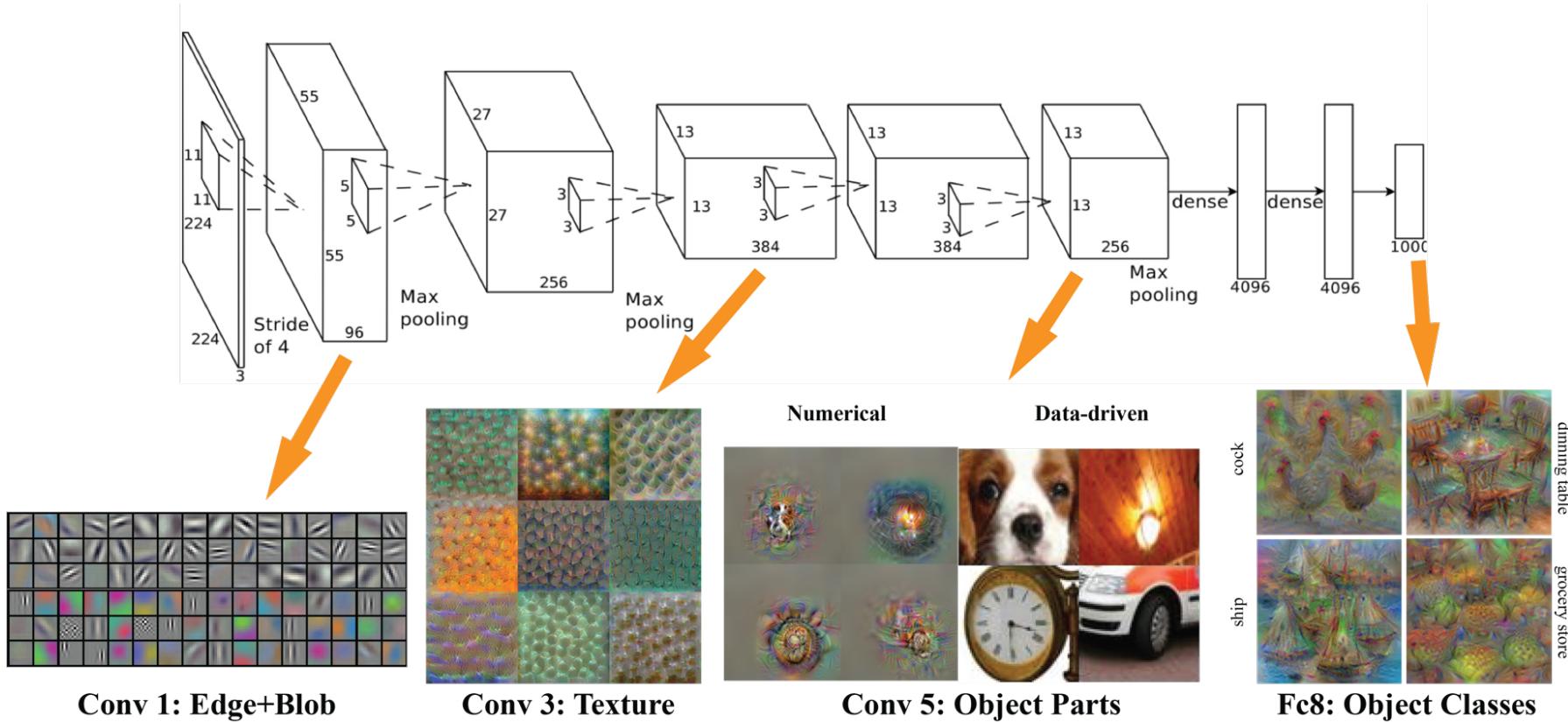
- The Pretrained Model is trained on ImageNet dataset
- If your data is similar to the ImageNet data
 - Fix all CNN Layers
 - Train FC layer



Outline

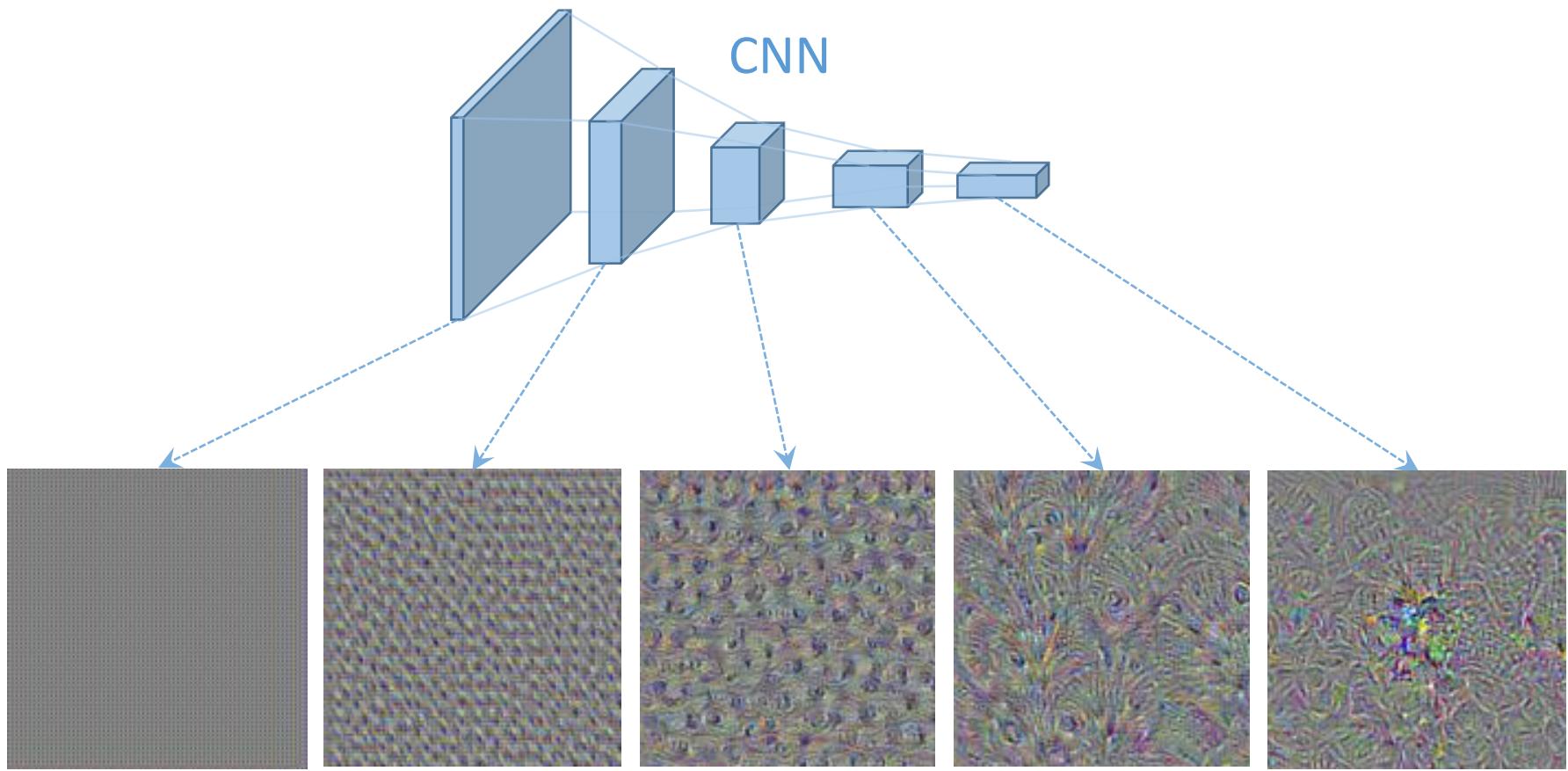
- CNN(Convolutional Neural Networks) Introduction
- Evolution of CNN
- **Visualizing the Features**
- Sentiment Analysis by CNN

Visualizing CNN



http://vision03.csail.mit.edu/cnn_art/data/single_layer.png

Different Layers of Visualization



Multiscale Image Generation

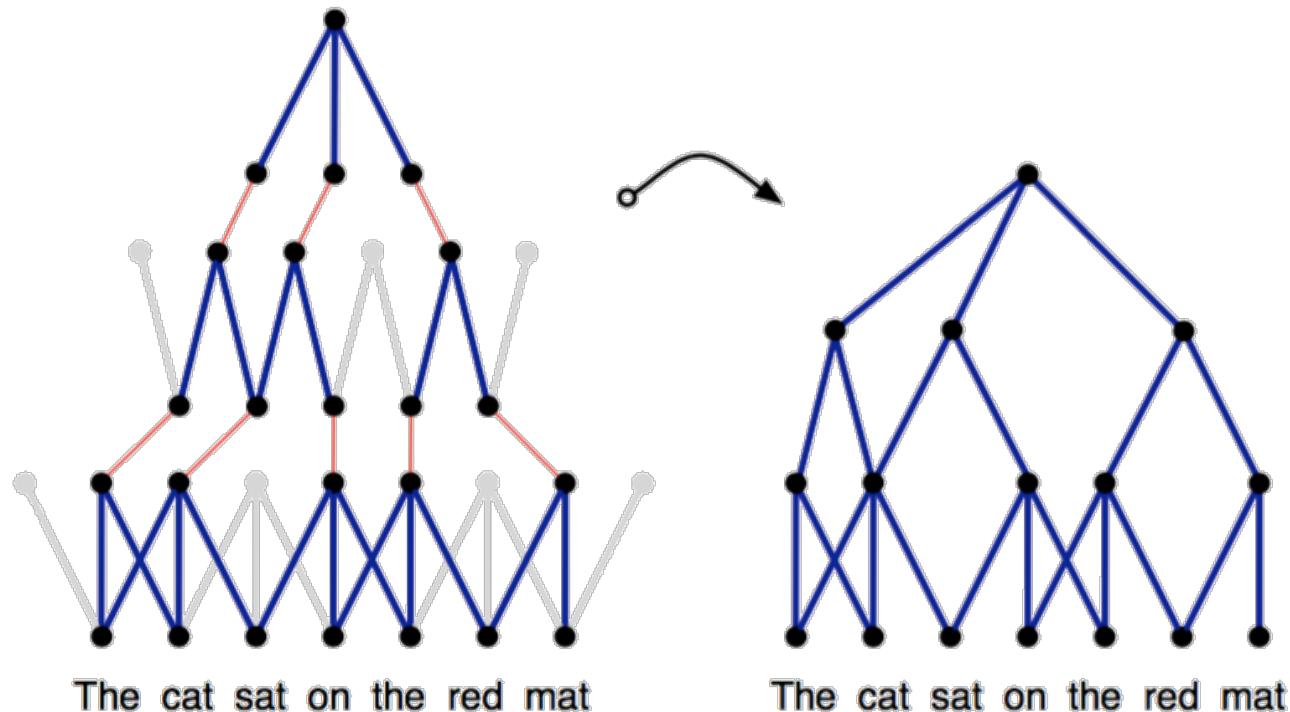


Outline

- CNN(Convolutional Neural Networks) Introduction
- Evolution of CNN
- Visualizing the Features
- Sentiment Analysis by CNN

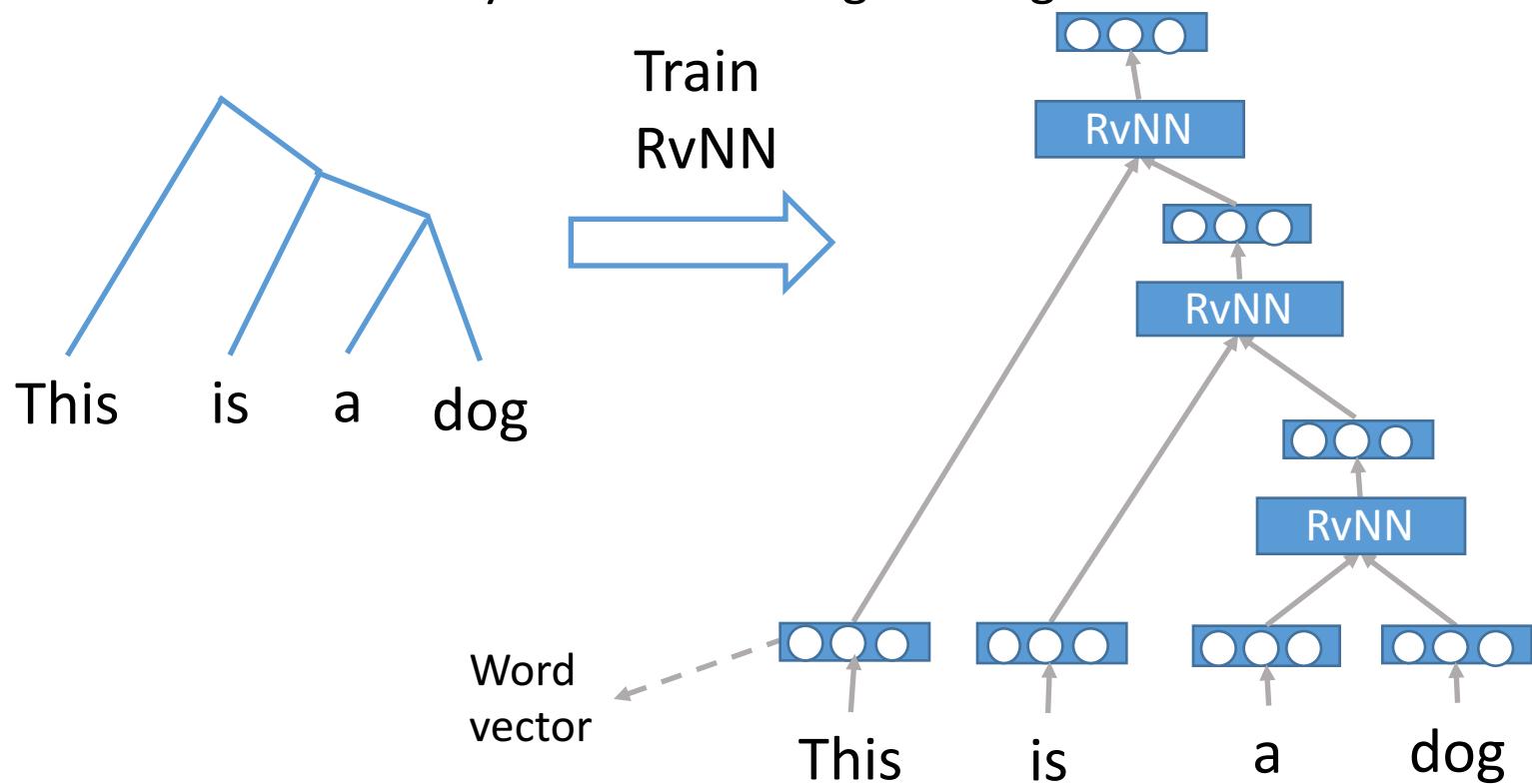
A Convolutional Neural Network for Modelling Sentences

- Paper: <https://arxiv.org/abs/1404.2188>
- Source code: <https://github.com/FredericGodin/DynamicCNN>



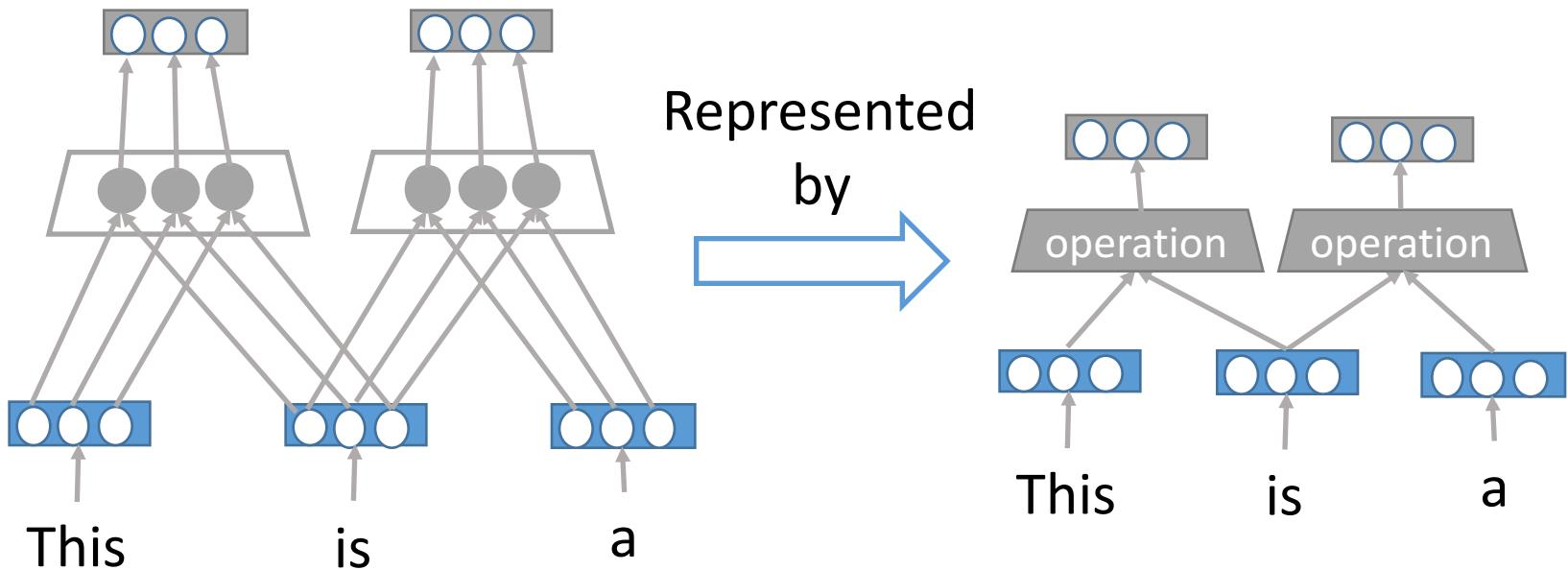
Drawbacks of Recursive Neural Networks(RvNN)

- Need human-labeled syntax tree during training



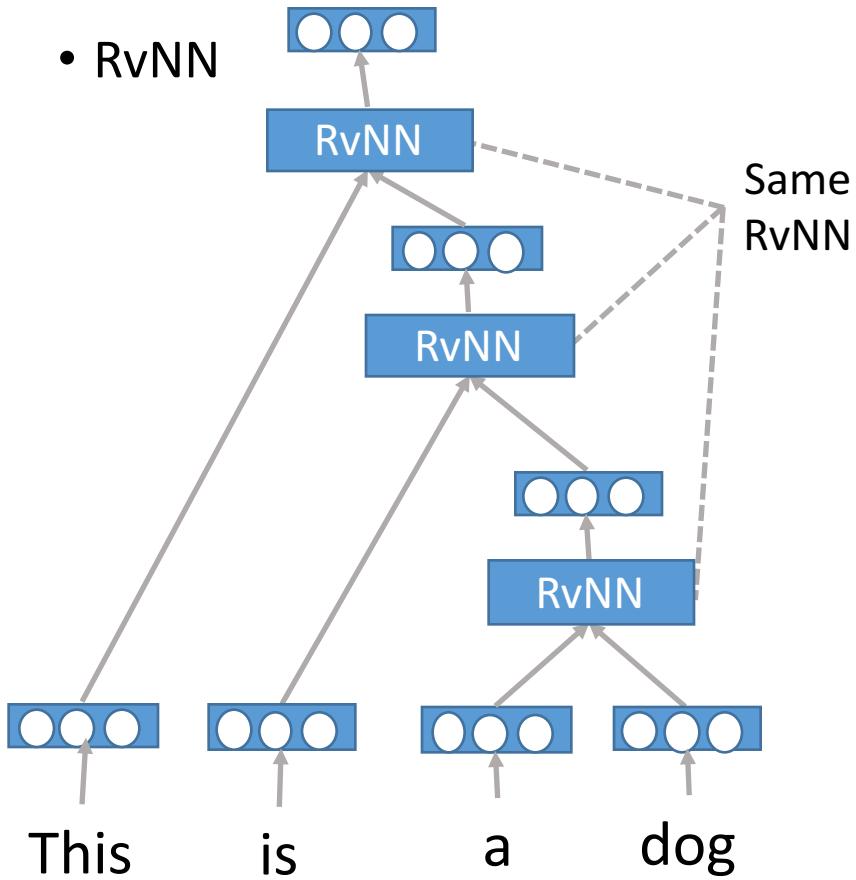
Element-wise 1D operations on word vectors

- 1D Convolution or 1D Pooling

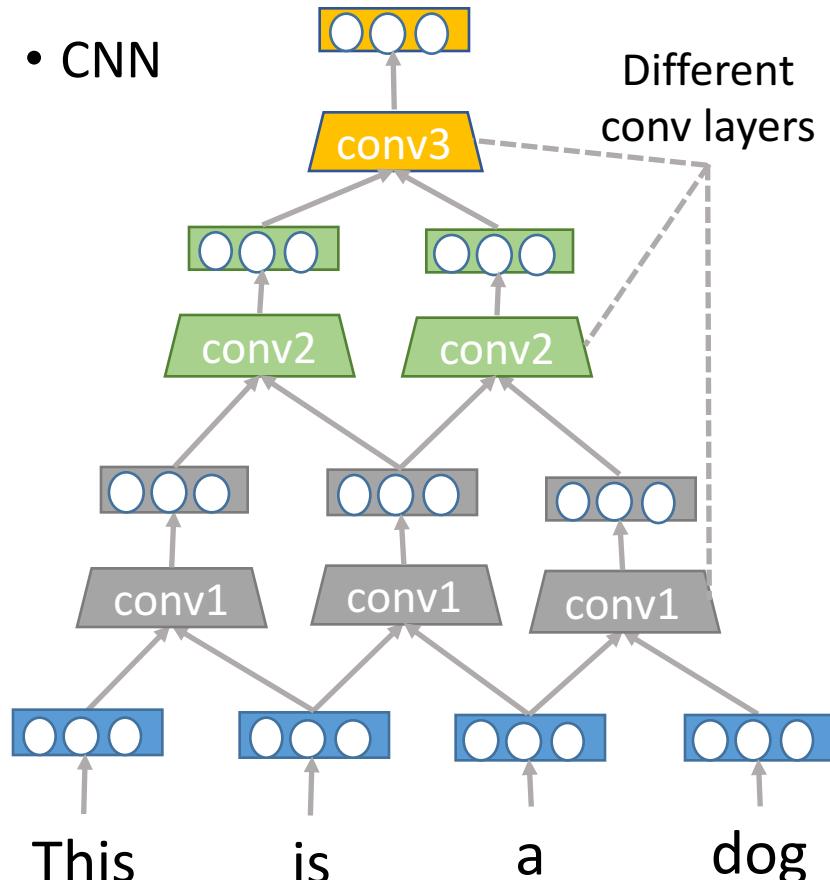


From RvNN to CNN

- RvNN

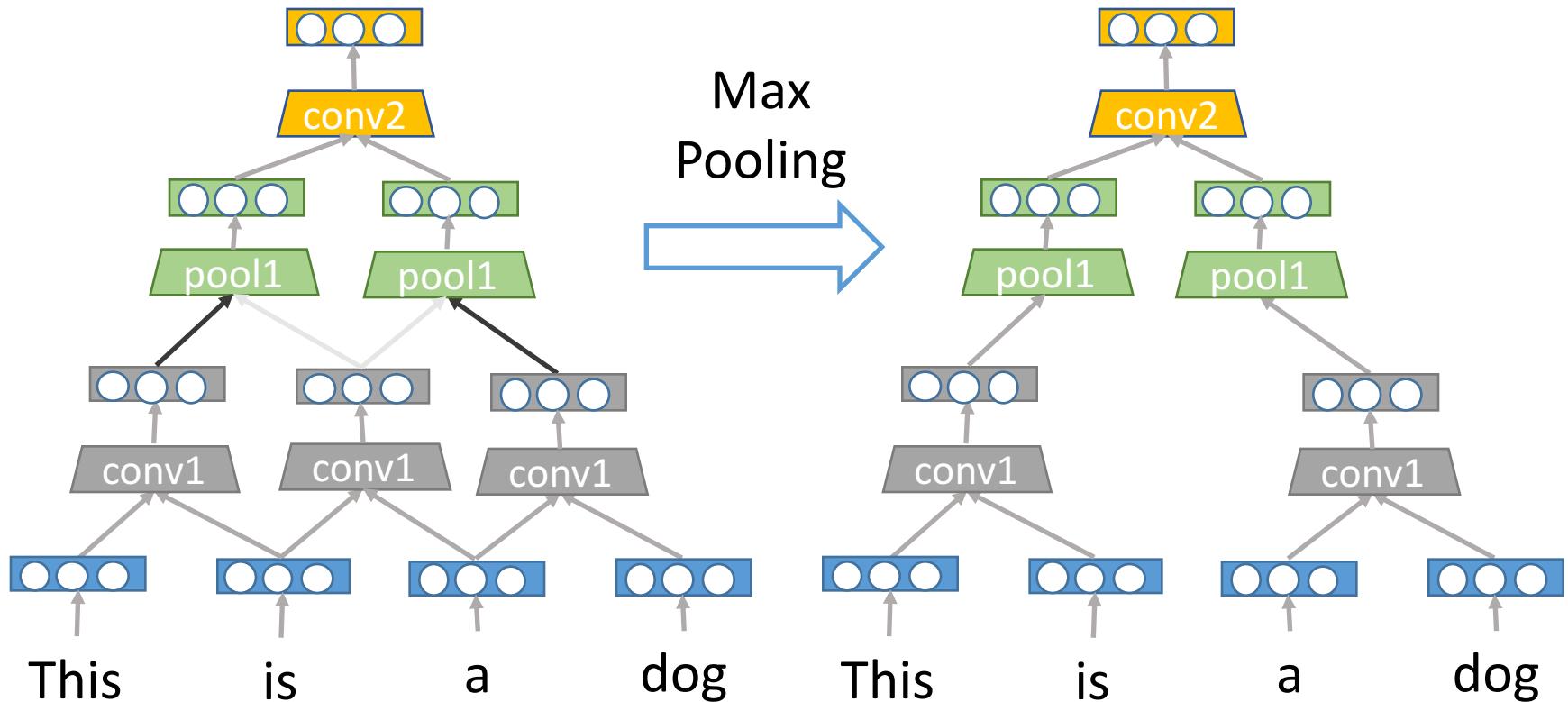


- CNN



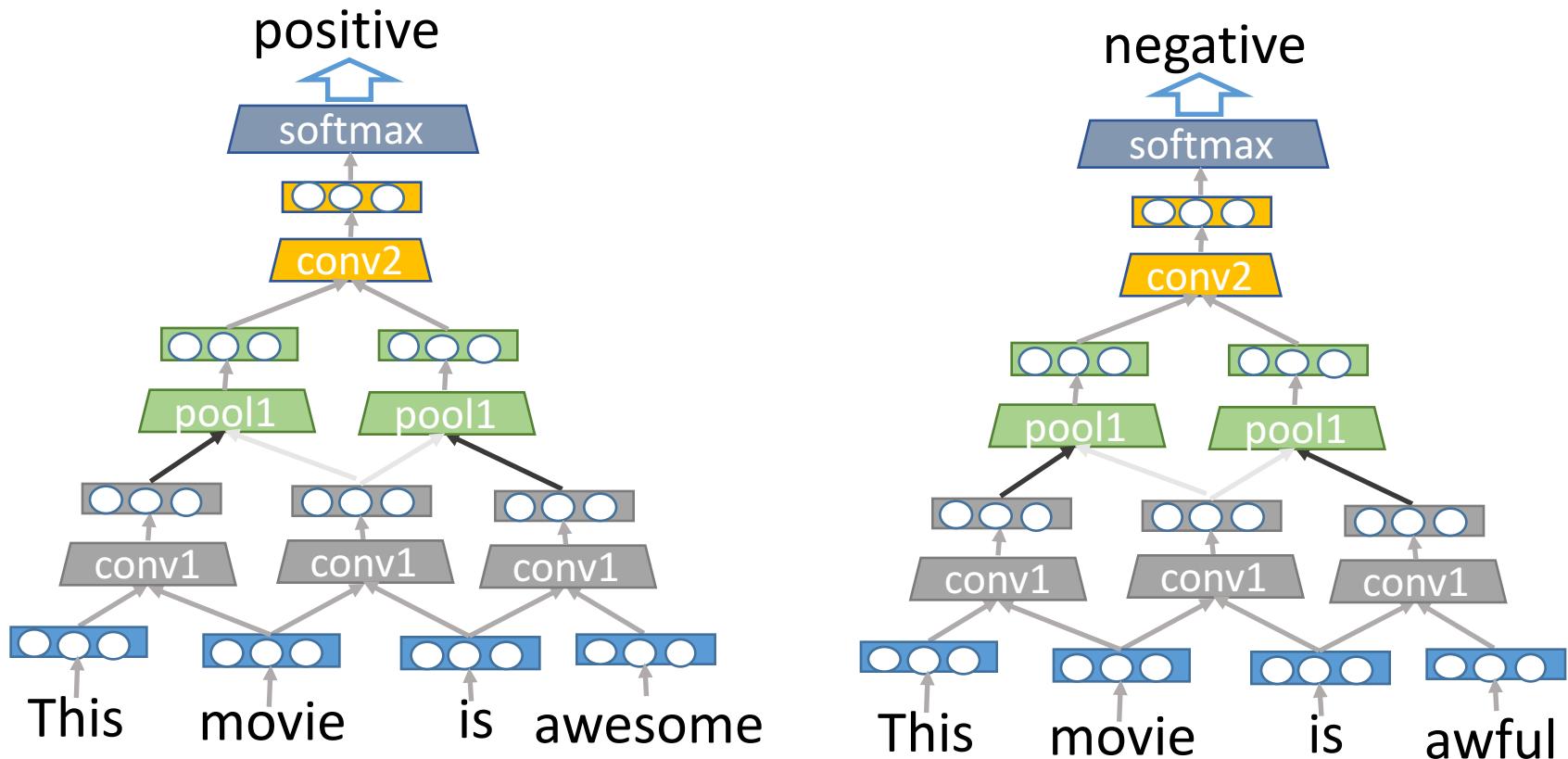
CNN with Max-Pooling Layers

- Similar to syntax tree
- But human-labeled syntax tree is not needed



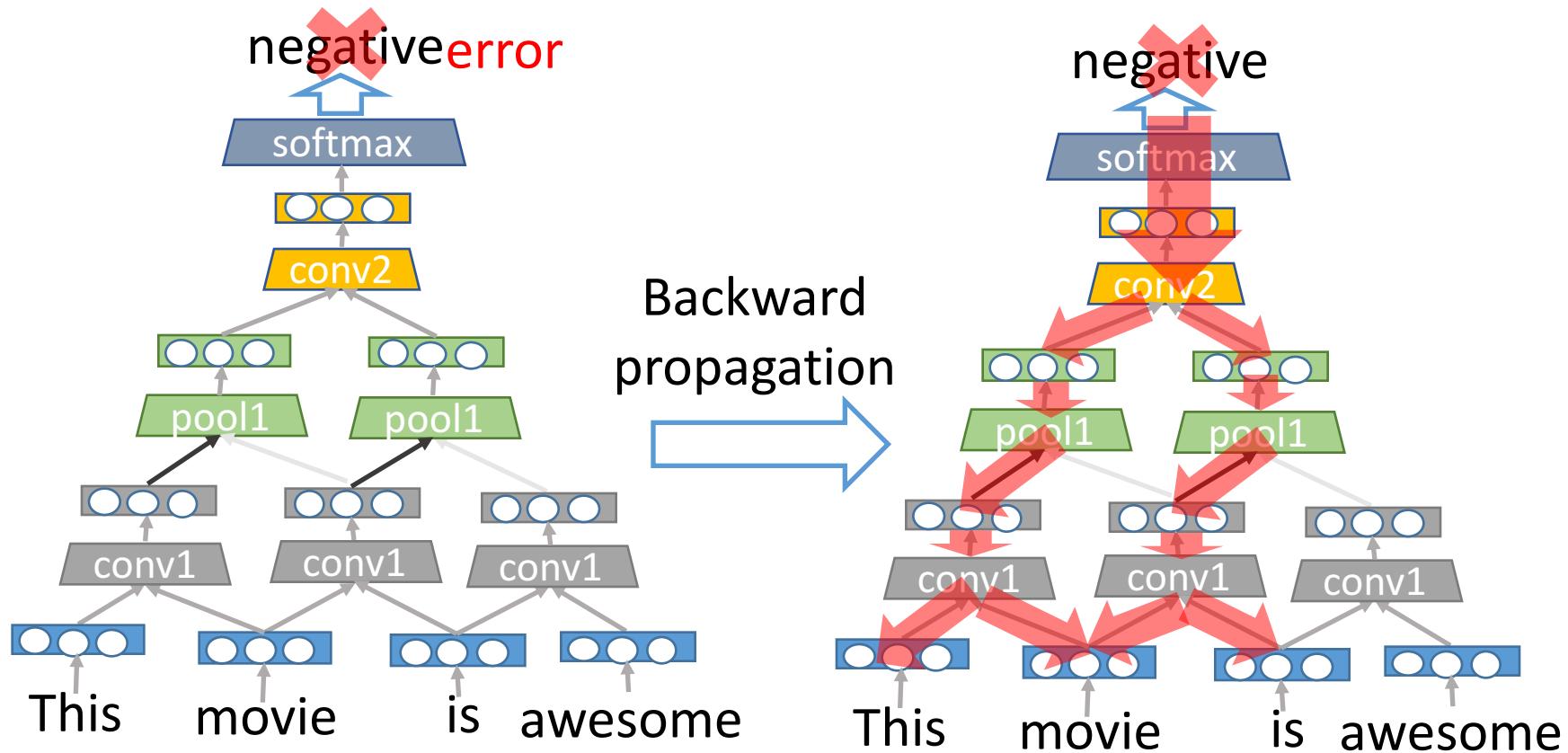
Sentiment Analysis by CNN

- Use softmax layer to classify the sentiments



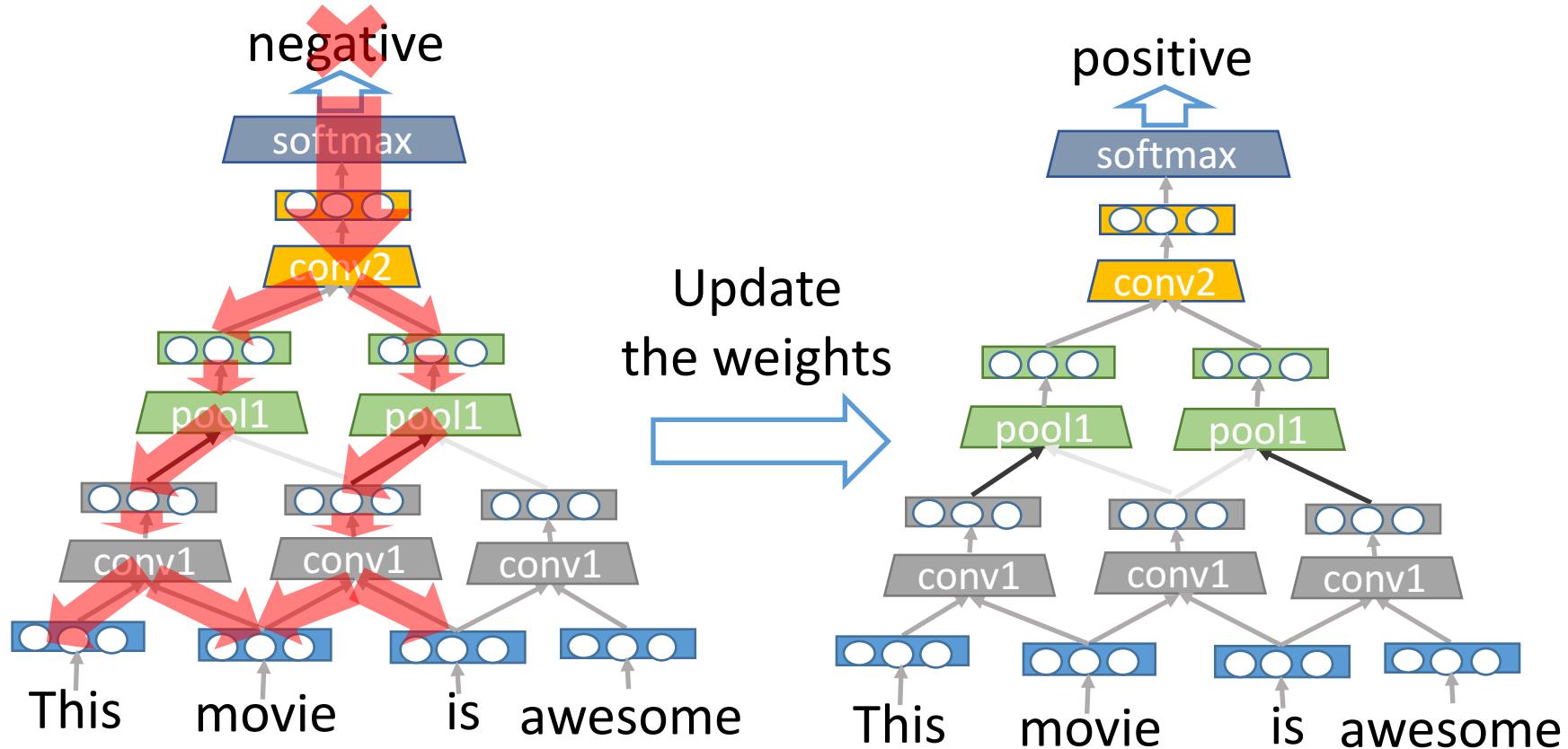
Sentiment Analysis by CNN

- Build the “correct syntax tree” by training



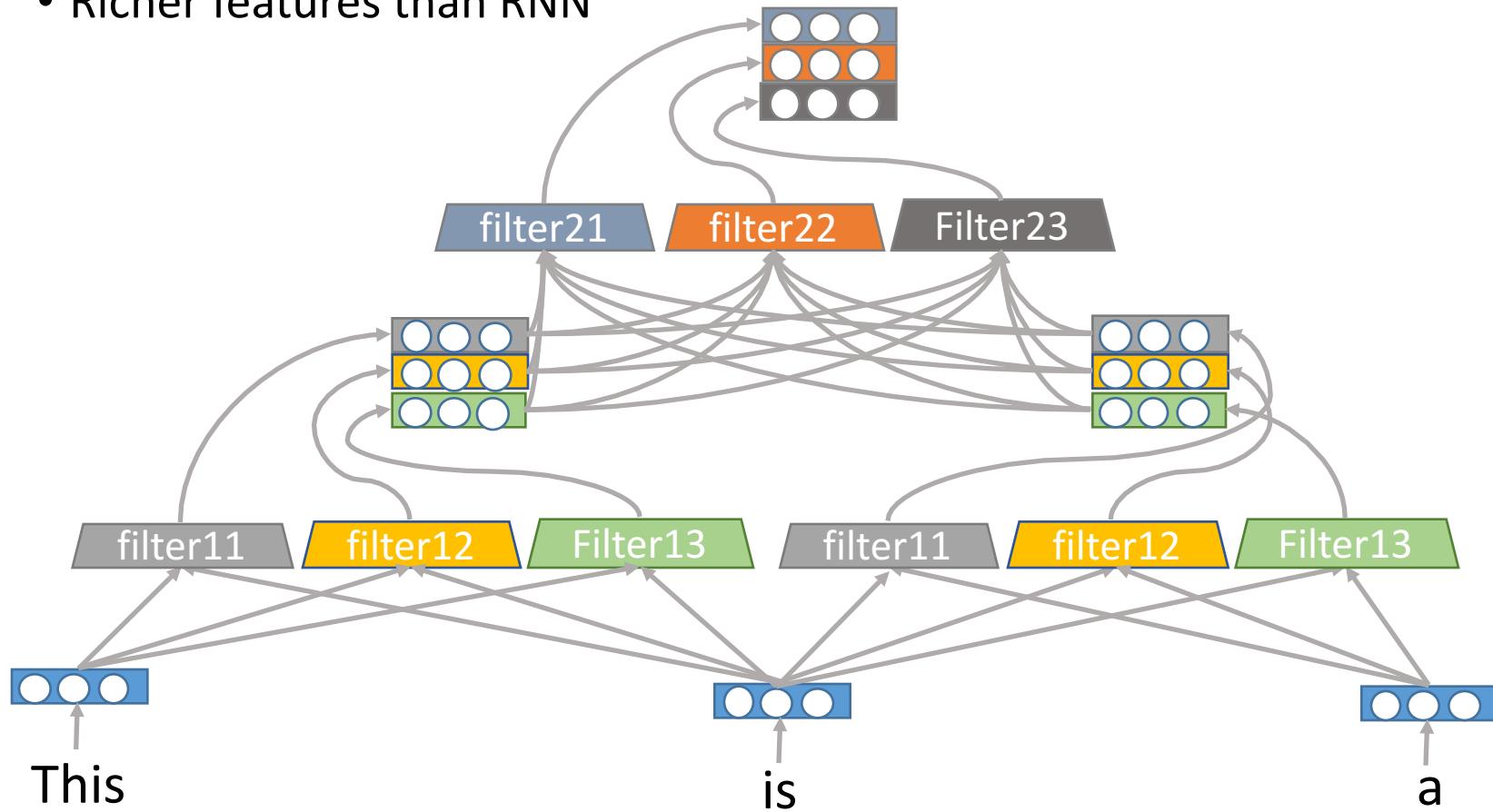
Sentiment Analysis by CNN

- Build the “correct syntax tree” by training



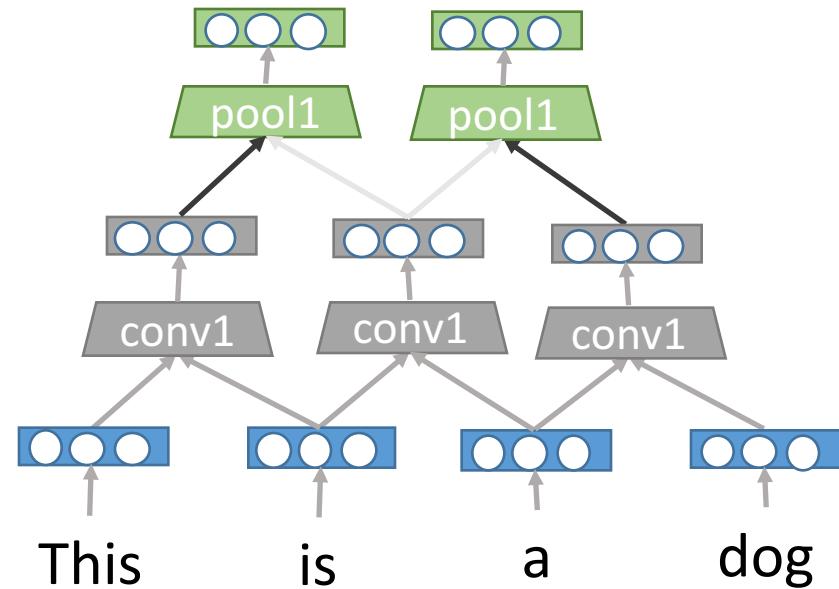
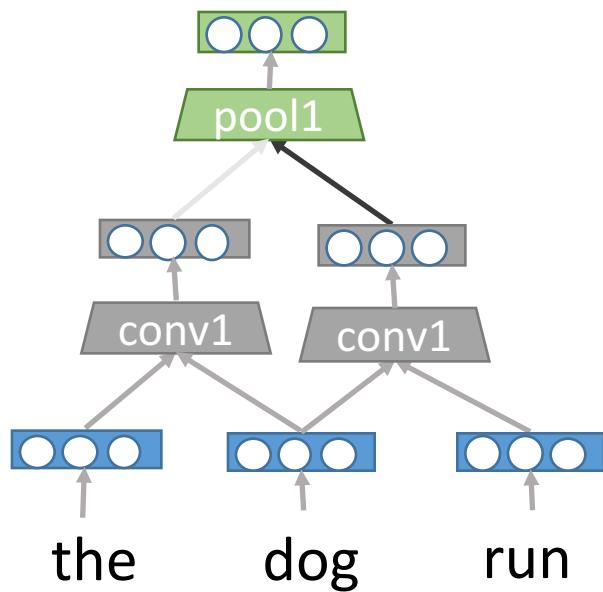
Multiple Filters

- Richer features than RNN



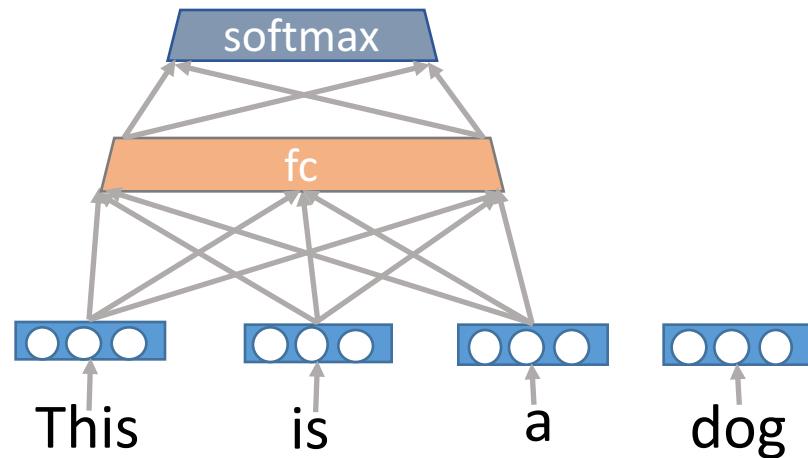
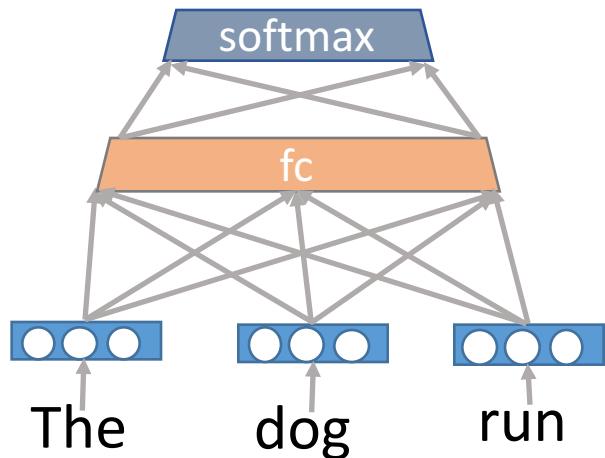
Various Input Size

- Convolutional layers and pooling layers
 - can handle input with various size



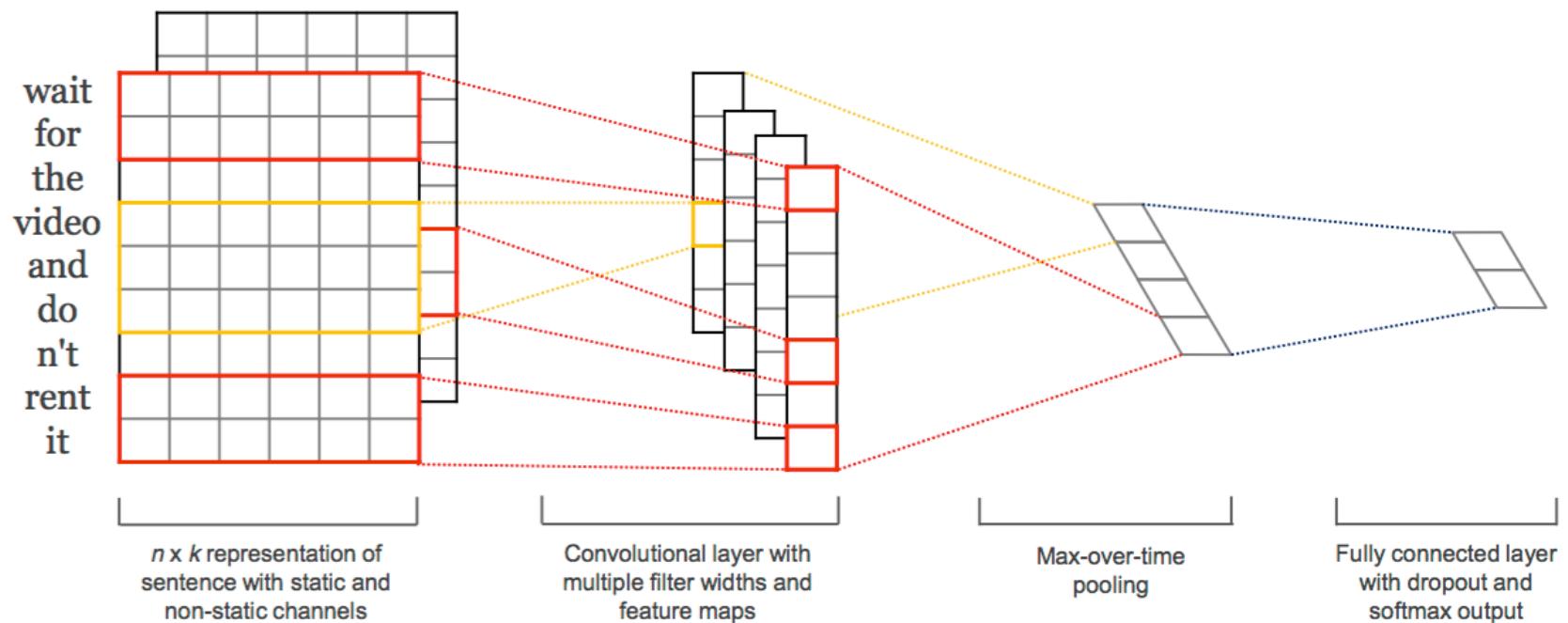
Various Input Size

- Fully-connected layer and softmax layer
 - need fixed-size input



Convolutional Neural Networks for Sentence Classification

- Paper: <http://www.aclweb.org/anthology/D14-1181>
- Sourcee code: https://github.com/yoonkim/CNN_sentence



Static & Non-Static Channel

- Pretrained by word2vec
- Static: fix the values during training
- Non-Static: update the values during training

