

### THE UNIVERSITY OF BRITISH COLUMBIA

# Topics in AI (CPSC 532S): **Multimodal Learning with Vision, Language and Sound**

Lecture 5: Convolutional Neural Networks (Part 2)



# Logistics:

# Assignment 2 is out

Groups + idea for project by January 31st

Paper list to be posted soon-ish

# Last Time: Convolutional Neural Networks



VGG-16 Network

# Improving Single Model

# Regularization

- L2, L1
- Dropout / Inverted Dropout
- Data augmentation

L2 Regularization: Learn a more (dense) distributed representation

# $R(\mathbf{W}) = ||\mathbf{W}|$

L1 Regularization: Learn a sparse representation

 $R(\mathbf{W}) = ||\mathbf{W}|$ 



Dropout

$$||_{2} = \sum_{i} \sum_{j} \mathbf{W}_{i,j}^{2}$$
  
n (few non-zero wight elements)

$$\|_1 = \sum_i \sum_j |\mathbf{W}_{i,j}|$$







Horizontal flips

Random crops & scales

Color Jitter

# **Horizontal flips**

### Random crops & scales





Color Jitter

Horizontal flips

# **Training:** sample random crops and scales e.g., ResNet:

- 1. Pick random L in range [256, 480]
- 2. Resize training image, short size = L
- 3. Sample random 224x224 patch

# **Testing:** average a fix set of crops e.g., ResNet:

Resize image to 5 scales (224, 256, 384, 480, 640) 2. For each image use 10 224x224 crops: 4 corners + center, + flips

### **Random crops & scales**

### Color Jitter



Horizontal flips

### Random perturbations in contrast and brightness



### Random crops & scales

### **Color Jitter**



# **Regularization:** Stochastic Depth

Effectively "dropout" but for layers

**some layer** (for each batch)



Huang et al., ECCV 2016]



**Common "Wisdom":** You need a lot of data to train a CNN



# This strategy is PERVASIVE.





**Solution: Transfer learning** — taking a model trained on the task that has lots of data and adopting it to the task that may not

### Train on **ImageNet**

| FC-1000  |
|----------|
| FC-4096  |
| FC-4096  |
|          |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-256 |
| Conv-256 |
| MaxPool  |
| Conv-128 |
| Conv-128 |
| MaxPool  |
| Conv-64  |
| Conv-64  |
| Image    |

- Why on **ImageNet**?
  - Convenience, lots of data
  - We know how to train these well



[Yosinski et al., NIPS 2014] Donahue et al., ICML 2014 Razavian et al., CVPR Workshop 2014

However, for some tasks we would need to start with something else (e.g., videos for optical flow)

### Train on **ImageNet**

| FC-1000  |
|----------|
| FC-4096  |
| FC-4096  |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-256 |
| Conv-256 |
| MaxPool  |
| Conv-128 |
| Conv-128 |
| MaxPool  |
| Conv-64  |
| Conv-64  |
| Image    |

# Re-initialize and train



[Yosinski et al., NIPS 2014] Donahue et al., ICML 2014 Razavian et al., CVPR Workshop 2014

**Small dataset** with C classes

# Lower levels of the CNN are at task independent anyways

Freeze these layers

### Train on **ImageNet**

### **Small dataset** with C classes

| FC-1000  |
|----------|
| FC-4096  |
| FC-4096  |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-512 |
| Conv-512 |
| MaxPool  |
| Conv-256 |
| Conv-256 |
| MaxPool  |
| Conv-128 |
| Conv-128 |
| MaxPool  |
| Conv-64  |
| Conv-64  |
| Image    |

# Re-initialize and train



[Yosinski et al., NIPS 2014] [Donahue et al., ICML 2014] [Razavian et al., CVPR Workshop 2014]

### Larger dataset

\* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

Freeze these layers



<sup>[</sup>Yosinski et al., NIPS 2014]

# Model **Ensemble**

**Training:** Train multiple independent models **Test:** Average their results

### ~ 2% improved performance in practice

**Alternative:** Multiple snapshots of the single model during training!

- **Improvement:** Instead of using the actual parameter vector, keep a moving average of the parameter vector and use that at test time (Polyak averaging)

# CPU vs. GPU (Why do we need Azure?)



Data from <a href="https://github.com/jcjohnson/cnn-benchmarks">https://github.com/jcjohnson/cnn-benchmarks</a>

# Frameworks: Super quick overview

# 1. Easily **build computational graphs**

# 2. Easily **compute gradients** in computational graphs

3. Run it all efficiently on a GPU (weap cuDNN, cuBLAS, etc.)



# Frameworks: Super quick overview

### Core DNN Frameworks

Caffe (UC Berkeley) Caffe 2 (Facebook)

(Baidu)

Torch (NYU/Facebook)

**PyTorch** (Facebook)

Theano (U Montreal) **TensorFlow** (Google)

Puddle

CNTK (Microsoft)

MXNet (Amazon)

### Wrapper Libraries

**Keras** TFLearn TensorLayer tf.layers **TF-Slim** tf.contrib.learn Pretty Tensor

# Frameworks: PyTorch vs. TensorFlow

Dynamic vs. Static computational graphs

# Frameworks: PyTorch vs. TensorFlow

Dynamic vs. **Static** computational graphs

### With static graphs, framework can optimize the graph for you before it runs!



# Frameworks: PyTorch vs. TensorFlow

**Dynamic** vs. Static computational graphs

Graph building and execution is intertwined. Graph can be different for every sample.

# The cat ate a big rat

# **PyTorch:** Three levels of abstraction

### **Tensor:** Imperative ndarray, but runs on GPU

# Variable: Node in a computational graph; stores data and gradients

### Module: A neural network layer; may store state or learnable weights



### Categorization



Horse Multi-class: Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person



### Categorization

### Detection





Multi-**class:** Horse Church Toothbrush Person IM GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





### Categorization

### Detection





Multi-**class:** Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





### Segmentation

Horse Person



### Categorization

### Detection





Multi-class: Horse Church Toothbrush Person **IM** GENET

Multi-label: Horse

Church Toothbrush Person

Horse (x, y, w, h) Horse (x, y, w, h) Person (x, y, w, h) Person (x, y, w, h)





### Segmentation

Horse Person



### Instance Segmentation

Horse1 Horse<sub>2</sub> Person1 Person2



# Object Classification



# Problem: For each image predict which category it belongs to out of a fixed set





# Object Classification





|  | Category | Predictio |
|--|----------|-----------|
|  | Dog      | No        |
|  | Cat      | No        |
|  | Couch    | No        |
|  | Flowers  | No        |
|  | Leopard  | Yes       |
|  |          |           |
|  |          |           |

**Problem:** For each image predict which category it belongs to out of a fixed set







# Object Classification







**Problem:** For each image predict which category it belongs to out of a fixed set

 $\mathbf{x}^t$ 

# **CNN Architectures:** LeNet-5



# Architecture: $CONV \longrightarrow POOL \longrightarrow CONV \longrightarrow POOL \longrightarrow FC \longrightarrow FC$ Conv filters: 5x5, Stride: 1 **Pooling:** 2x2, Stride: 2

[LeCun et al., 1998]

# ImageNet Dataset

Over **14 million** (high resolution) web **images** Roughly labeled with **22K synset** categories Labeled on Amazon Mechanical Turk (AMT)



### **Popular Synsets**

### Animal

fish bird mammal invertebrate

### Plant

tree flower vegetable

### Activity

sport

### Material

### fabric

### Instrumentation

utensil appliance tool musical instrument

### Scene

room geological formation

### Food

beverage



# ImageNet Competition (ILSVRC)

Annual competition of image classification at scale Focuses on a subset of **1K synset** categories **Scoring:** need to predict true label within top K (K=5)





# AlexNet

Architecture: CONV1 MAX POOL1 NORM1 CONV2 MAX POOL2 NORM2 CONV3 CONV4 CONV5 Max POOL3 FC6 FC7 FC8

**Output:** 55 x 55 x 96



[Krizhevsky et al., 2012]

### **Input:** 227 x 227 x 3 images

# **CONV1:** 96 11 x 11 filters applied at stride 4

Parameters: 35K

### MAX POOL1: 96 11 x 11 filters applied at stride 4 **Output:** 27 x 27 x 96 **Parameters:** 0

# AlexNet

Full (simplified) AlexNet architecture: [227x227x3] INPUT [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0 [27x27x96] MAX POOL1: 3x3 filters at stride 2 [27x27x96] NORM1: Normalization layer [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2 [13x13x256] MAX POOL2: 3x3 filters at stride 2 [13x13x256] NORM2: Normalization layer [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1 [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1 [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1 [6x6x256] MAX POOL3: 3x3 filters at stride 2 [4096] FC6: 4096 neurons [4096] FC7: 4096 neurons [1000] FC8: 1000 neurons (class scores)



### Krizhevsky et al., 2012

### **Details / Comments**

- First use of ReLU
- Used contrast normalization layers
- Heavy data augmentation
- Dropout of 0.5
- Batch size of 128
- SGD Momentum (0.9)
- Learning rate (1e-2) reduced by 10 manually when validation accuracy plateaus
- L2 weight decay
- -7 CNN ensample: 18.2% -> 15.4%



# **ILSVRC** winner 2012



# **ZF** Net



### **AlexNet with small modifications:**

- CONV1 (11 x 11 stride 4) to (7 x 7 stride 2)
- CONV3 # of filters 384 -> 512
- CONV4 # of filters 384 -> 1024
- CONV5 # of filters 256 -> 512

[Zeiler and Fergus, 2013]

# **ILSVRC** winner 2012



# VGG Net

# Trend:

- -smaller filters (3 x 3)
- -deeper network (16 or 19 vs. 8 in Alex

# Why?

- receptive field of a 3 layer ConvNet with filter size = is the same as 1 layer ConvNet with filter 7x7 (at stride

- deeper = more non-linearity (so richer filters)
- fewer parameters

### [Simonyan and Zisserman, 2014]

|            |                |               | EC 1000       |
|------------|----------------|---------------|---------------|
|            |                | Softmax       | FC 1000       |
|            |                | FC 1000       | FC 4096       |
|            |                | FC 4096       | Pool          |
|            |                | FC 4096       | 3x3 conv, 512 |
|            |                | Pool          | 3x3 conv, 512 |
| xNet)      |                | 3x3 conv, 512 | 3x3 conv, 512 |
|            |                | 3x3 conv, 512 | 3x3 conv, 512 |
|            |                | 3x3 conv, 512 | Pool          |
|            |                | Pool          | 3x3 conv, 512 |
|            | Softmax        | 3x3 conv, 512 | 3x3 conv, 512 |
|            | FC 1000        | 3x3 conv, 512 | 3x3 conv, 512 |
|            | FC 4096        | 3x3 conv, 512 | 3x3 conv, 512 |
|            | FC 4096        | Pool          | Pool          |
| = 3x3      | Pool           | 3x3 conv, 256 | 3x3 conv, 256 |
| 1)         | 3x3 conv, 256  | 3x3 conv, 256 | 3x3 conv, 256 |
| I <i>)</i> | 3x3 conv, 384  | Pool          | Pool          |
|            | Pool           | 3x3 conv, 128 | 3x3 conv, 128 |
|            | 3x3 conv, 384  | 3x3 conv, 128 | 3x3 conv, 128 |
|            | Pool           | Pool          | Pool          |
|            | 5x5 conv, 256  | 3x3 conv, 64  | 3x3 conv, 64  |
|            | 11x11 conv, 96 | 3x3 conv, 64  | 3x3 conv, 64  |
|            | Input          | Input         | Input         |
|            | AlexNet        | VGG16         | VGG19         |





# VGG Net

INPUT: [224x224x3] memory: 224\*224\*3=150K params: 0 CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params: (3\*3\*3)\*64 = 1,728 CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params: (3\*3\*64)\*64 = 36,864 POOL2: [112x112x64] memory: 112\*112\*64=800K params: 0 CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params: (3\*3\*64)\*128 = 73,728 CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params: (3\*3\*128)\*128 = 147,456 POOL2: [56x56x128] memory: 56\*56\*128=400K params: 0 CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*128)\*256 = 294,912 CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*256)\*256 = 589,824 CONV3-256: [56x56x256] memory: 56\*56\*256=800K params: (3\*3\*256)\*256 = 589,824 POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0 CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*256)\*512 = 1,179,648 CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512 = 2,359,296 CONV3-512: [28x28x512] memory: 28\*28\*512=400K params: (3\*3\*512)\*512 = 2,359,296 POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0 CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296 CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296 CONV3-512: [14x14x512] memory: 14\*14\*512=100K params: (3\*3\*512)\*512 = 2,359,296 POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0 FC: [1x1x4096] memory: 4096 params: 7\*7\*512\*4096 = 102,760,448 FC: [1x1x4096] memory: 4096 params: 4096\*4096 = 16,777,216 FC: [1x1x1000] memory: 1000 params: 4096\*1000 = 4,096,000

TOTAL memory: 24M \* 4 bytes ~= 96MB / image (only forward! ~\*2 for bwd) TOTAL params: 138M parameters

### [Simonyan and Zisserman, 2014]

```
(not counting biases)
```

\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford



VGG16



# **ILSVRC** winner 2012



### even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

[Szegedy et al., 2014]



these modules

Szegedy et al., 2014]

### Idea: design good local topology ("network within network") and then stack

![](_page_44_Figure_4.jpeg)

these modules

# Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

**Concatenate** all filter outputs together at output depth-wise Szegedy et al., 2014

### Idea: design good local topology ("network within network") and then stack

### What's the problem?

![](_page_45_Figure_9.jpeg)

### Naive Inception module

these modules

# Apply parallel filter operations on the input from previous layer

 Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)

- Pooling operation (3x3)

**Concatenate** all filter outputs together at output depth-wise

Szegedy et al., 2014

### **Idea:** design good local topology ("network within network") and then stack

![](_page_46_Figure_8.jpeg)

# Convolutional Layer: 1x1 convolutions

### 56 x 56 x 64 **image**

![](_page_47_Figure_2.jpeg)

### Idea: design good local topology ("network within network") and then stack these modules

![](_page_48_Figure_2.jpeg)

### Naive Inception module

Szegedy et al., 2014]

1x1 "bottleneck" layers

![](_page_48_Figure_7.jpeg)

Inception module with dimension reduction

### saves approximately 60% of computations

### even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

[Szegedy et al., 2014]

![](_page_49_Figure_9.jpeg)

### even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

### [Szegedy et al., 2014]

![](_page_50_Figure_9.jpeg)

### even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

[Szegedy et al., 2014]

![](_page_51_Figure_9.jpeg)

### even deeper network with computational efficiency

- -22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
- (12x less than AlexNet!)
- Better performance (@6.7 top 5 error)

### [Szegedy et al., 2014]

![](_page_52_Figure_9.jpeg)

# **ILSVRC** winner 2012

![](_page_53_Figure_1.jpeg)