

Advanced Architectures Object Detection & Segmentation Bright & Dark Side of Context

Prof. Dr. Mario Fritz CISPA Helmholtz Center for Information Security

Research



Computer Vision

- Scene Understanding
- Segmentation
- QA and Captioning







What sport is this man enjoying?

Ours: a skier is headed down a steep slope

Security & Privacy

- Trustworthy AI
- Understanding & Controlling Privacy of Data & Models
- Adversarial Machine Learning/Attacks
- AI/ML for Cybersecurity
- Detecting Deep Fakes



Machine Learning

- Deep Learning
- Domain Adaptation
- Generative Adversarial Networks, Variational Autoencoders
- Uncertainty
- Explainability



https://fritz.cispa.saarland

Review: LeNet-5

[LeCun et al., 1998]



Conv filters were 5x5, applied at stride 1 Subsampling (Pooling) layers were 2x2 applied at stride 2 i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

[Krizhevsky et al. 2012]

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



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

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4 => Q: what is the output volume size? Hint: (227-11)/4+1 = 55

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

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4 => Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?

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

[Krizhevsky et al. 2012]



Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4 => Output volume **[55x55x96]**

Parameters: (11*11*3)*96 = **35K**

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

[Krizhevsky et al. 2012]



Input: 227x227x3 images After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2

Q: what is the output volume size? Hint: (55-3)/2+1 = 27

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

[Krizhevsky et al. 2012]



Input: 227x227x3 images After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2 Output volume: 27x27x96

Q: what is the number of parameters in this layer?

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

[Krizhevsky et al. 2012]



Input: 227x227x3 images After CONV1: 55x55x96

Second layer (POOL1): 3x3 filters applied at stride 2 Output volume: 27x27x96 Parameters: 0!

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

[Krizhevsky et al. 2012]

. . .



Input: 227x227x3 images After CONV1: 55x55x96 After POOL1: 27x27x96

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

[Krizhevsky et al. 2012]

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 [27x27x26] 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] MAX POOL3: 3x3 filters at stride 1, pad 1 [13x13x256] MAX POOL3: 3x3 filters at stride 2 [4096] FC6: 4096 neurons [4096] FC7: 4096 neurons [1000] FC8: 1000 neurons (class scores)



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

[Krizhevsky et al. 2012]

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 [27x27x26] 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] MAX POOL3: 3x3 filters at stride 1, pad 1 [13x13x256] MAX POOL3: 3x3 filters at stride 2 [4096] FC6: 4096 neurons [4096] FC7: 4096 neurons [1000] FC8: 1000 neurons (class scores)

dense 192 192 128 48 128 dense dense 1000 128 Max 192 192 2048 2048 pooling Max 128 pooling pooling

Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10
- manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

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







[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

8 layers (AlexNet) -> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1 and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13 (ZFNet)

-> 7.3% top 5 error in ILSVRC'14





AlexNet

VGG19

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

VGG16

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)



[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same effective receptive field as one 7x7 conv layer

Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



Softmax

FC 1000 FC 4096

FC 4096

Pool

Pool

Pool

Input

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same **effective receptive field** as one 7x7 conv layer

[7x7]



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

Softmax FC 1000

FC 4096 FC 4096

Pool

Pool

Pool

1x11 conv

Input

AlexNet

[Simonyan and Zisserman, 2014]

Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same effective receptive field as one 7x7 conv layer

But deeper, more non-linearities

And fewer parameters: $3 * (3^2C^2)$ vs. 7²C² for C channels per layer



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

Softmax

FC 1000 FC 4096

FC 4096

Pool

Pool

Pool

Input

(not counting biases) 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

Softmax FC 1000 FC 4096 FC 4096 Pool Pool Pool Pool x3 conv, 12 Pool k3 conv, Input

VGG16

(not counting biases) 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 (for a forward pass) TOTAL params: 138M parameters

Softmax FC 1000 FC 4096 FC 4096 Pool 3x3 conv, 51 Pool Pool Pool x3 conv, 12 Pool k3 conv, Input

VGG16

INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases) CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728	Note:
CONV3-64: [224x224x64] memory: 224*224*64=3.2M arams: (3*3*64)*64 = 36,864	
POOL2: [112x112x64] memory: 112*112*64=800K params: 0	Most memory is in
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	Most params are
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296	in late FC
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	

(not counting biases) memory: 224*224*3=150K params: 0 INPUT: [224x224x3] 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]

Details:

- ILSVRC'14 2nd in classification, 1st in localization
- Similar training procedure as Krizhevsky 2012
- No Local Response Normalisation (LRN)
- Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
- Use ensembles for best results
- FC7 features generalize well to other tasks



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

fc7

fc6



[Szegedy et al., 2014]

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
 12x less than AlexNet
- ILSVRC'14 classification winner (6.7% top 5 error)





[Szegedy et al., 2014]

"Inception module": design a good local network topology (network within a network) and then stack these modules on top of each other



[Szegedy et al., 2014]



Naive Inception module

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 depth-wise

[Szegedy et al., 2014]



Naive Inception module

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 depth-wise

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Q: What is the problem with this? [Hint: Computational complexity]

Example:



Naive Inception module

[Szegedy et al., 2014]

Example:

Q1: What is the output size of the 1x1 conv, with 128 filters?



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Example:

Q1: What is the output size of the 1x1 conv, with 128 filters?



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Example:

Q2: What are the output sizes of all different filter operations?



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Example:

Q2: What are the output sizes of all different filter operations?



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]
[Szegedy et al., 2014]

Example:

Q3:What is output size after filter concatenation?



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Example: Q3:What is output size after filter concatenation?

28x28x(128+192+96+256) = 28x28x672



Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

[Szegedy et al., 2014]

Example:

Q3:What is output size after filter concatenation?

28x28x(128+192+96+256) = 28x28x672



Q: What is the problem with this? [Hint: Computational complexity]

Conv Ops:

[1x1 conv, 128] 28x28x128x1x1x256 [3x3 conv, 192] 28x28x192x3x3x256 [5x5 conv, 96] 28x28x96x5x5x256 **Total: 854M ops**

Naive Inception module

[Szegedy et al., 2014]

Example:

Q3:What is output size after filter concatenation?





Naive Inception module

Q: What is the problem with this? [Hint: Computational complexity]

Conv Ops:

[1x1 conv, 128] 28x28x128x1x1x256 [3x3 conv, 192] 28x28x192x3x3x256 [5x5 conv, 96] 28x28x96x5x5x256 **Total: 854M ops**

Very expensive compute

Pooling layer also preserves feature depth, which means total depth after concatenation can only grow at every layer!

[Szegedy et al., 2014]

Example:

Q3:What is output size after filter concatenation?



Q: What is the problem with this? [Hint: Computational complexity]

Solution: "bottleneck" layers that use 1x1 convolutions to reduce feature depth

Naive Inception module

Reminder: 1x1 convolutions



Reminder: 1x1 convolutions



[Szegedy et al., 2014]



Inception module with dimension reduction

[Szegedy et al., 2014]

1x1 conv "bottleneck" layers



Inception module with dimension reduction

[Szegedy et al., 2014]



Inception module with dimension reduction

Using same parallel layers as naive example, and adding "1x1 conv, 64 filter" bottlenecks:

Conv Ops:

[1x1 conv, 64] 28x28x64x1x1x256 [1x1 conv, 64] 28x28x64x1x1x256 [1x1 conv, 128] 28x28x128x1x1x256 [3x3 conv, 192] 28x28x192x3x3x64 [5x5 conv, 96] 28x28x96x5x5x64 [1x1 conv, 64] 28x28x64x1x1x256 **Total: 358M ops**

Compared to 854M ops for naive version Bottleneck can also reduce depth after pooling layer



[Szegedy et al., 2014]





[Szegedy et al., 2014]



[Szegedy et al., 2014]



[Szegedy et al., 2014]



[Szegedy et al., 2014]



22 total layers with weights

(parallel layers count as 1 layer => 2 layers per Inception module. Don't count auxiliary output layers)

[Szegedy et al., 2014]

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- No FC layers
- 12x less params than AlexNet
- ILSVRC'14 classification winner (6.7% top 5 error)



Inception module





ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners "Revolution of Depth"

[He et al., 2015]

Very deep networks using residual connections

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!





[He et al., 2015]

What happens when we continue stacking deeper layers on a "plain" convolutional neural network?

[He et al., 2015]

What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



Q: What's strange about these training and test curves? [Hint: look at the order of the curves]

[He et al., 2015]

What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



56-layer model performs worse on both training and test error -> The deeper model performs worse, but it's not caused by overfitting!

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

[He et al., 2015]

Hypothesis: the problem is an *optimization* problem, deeper models are harder to optimize

The deeper model should be able to perform at least as well as the shallower model.

A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



[He et al., 2015]

Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Softmax Case Study: ResNet FC 1000 [He et al., 2015] 3x3 conv, 512 3x3 conv, 512 3x3 conv, 512 Full ResNet architecture: 3x3 conv. 512 relu 3x3 conv, 512 Stack residual blocks 3x3 conv, 512, /2 F(x) + xEvery residual block has two 3x3 conv layers 3x3 conv Х 3x<u>3 conv.</u> F(x) relu identity 3x3 conv 3x3 conv, 128 3x3 conv, 64 3x3 conv, 64 3x3 conv, 64 Х 3x3 conv. 64 **Residual block** 3x3 conv. 64 3x3 conv, 64

_

-

Softmax Case Study: ResNet FC 1000 [He et al., 2015] 3x3 conv, 512 3x3 conv. 512 3x3 conv, 512 Full ResNet architecture: 3x3 conv. 512 relu 3x3 conv, 512 Stack residual blocks 3x3 conv. 512. /2 F(x) + x (+ Every residual block has two 3x3 conv layers 3x3 conv Periodically, double # of -3x3 conv, 128 Х filters and downsample F(x)filters, /2 relu identity spatially with spatially using stride 2 stride 2 3x3 conv 3x3 conv (/2 in each dimension) 3x3 conv, 64 3x3 conv, 64 3x3 conv. 64 filters 3x3 conv, 64 Х 3x3 conv. 64 **Residual block** 3x3 conv. 64 3x3 conv. 64

Case Study: ResNet [He et al., 2015] Full ResNet architecture: relu Stack residual blocks F(x) + x (+ Every residual block has two 3x3 conv layers 3x3 conv Periodically, double # of -Х filters and downsample F(x)relu

spatially using stride 2 (/2 in each dimension)

 Additional conv layer at the beginning



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

Softmax

Softmax Case Study: ResNet **No FC layers** FC 1000 besides FC 1000 to [He et al., 2015] 3x3 conv, 512 output 3x3 conv. 512 classes 3x3 conv, 512 Full ResNet architecture: 3x3 conv. 512 Global relu average 3x3 conv, 512 Stack residual blocks 3x3 conv. 512. pooling layer F(x) + xEvery residual block has after last conv layer two 3x3 conv layers 3x3 conv Periodically, double # of Х filters and downsample 3x3 conv. F(x) relu identity spatially using stride 2 3x3 conv 3x3 conv, 128 (/2 in each dimension) 3x3 conv. 64 Additional conv layer at 3x3 conv. 64 the beginning 3x3 conv, 64 3x3 conv. 64 No FC layers at the end **Residual block** 3x3 conv. 64 (only FC 1000 to output 3x3 conv, 64 classes)



[He et al., 2015]

For deeper networks (ResNet-50+), use "bottleneck" layer to improve efficiency (similar to GoogLeNet)



[He et al., 2015]

For deeper networks (ResNet-50+), use "bottleneck" layer to improve efficiency (similar to GoogLeNet)



[He et al., 2015]

Training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier 2/ initialization from He et al.
- SGD + Momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

[He et al., 2015]

Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lowing training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

1st places in all five main tracks

- ImageNet Classification: "Ultra-deep" (quote Yann) 152-layer nets
- ImageNet Detection: 16% better than 2nd
- ImageNet Localization: 27% better than 2nd
- COCO Detection: 11% better than 2nd
- COCO Segmentation: 12% better than 2nd
Case Study: ResNet

[He et al., 2015]

Experimental Results

- Able to train very deep networks without degrading (152 layers on ImageNet, 1202 on Cifar)
- Deeper networks now achieve lowing training error as expected
- Swept 1st place in all ILSVRC and COCO 2015 competitions

MSRA @ ILSVRC & COCO 2015 Competitions

1st places in all five main tracks

- ImageNet Classification: "Ultra-deep" (quote Yann) 152-layer nets
- ImageNet Detection: 16% better than 2nd
- ImageNet Localization: 27% better than 2nd
- COCO Detection: 11% better than 2nd
- COCO Segmentation: 12% better than 2nd

ILSVRC 2015 classification winner (3.6% top 5 error) -- better than "human performance"! (Russakovsky 2014)

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners







An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.



An Analysis of Deep Neural Network Models for Practical Applications, 2017.



Comparing complexity...

An Analysis of Deep Neural Network Models for Practical Applications, 2017.

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

VGG: Highest



GoogLeNet:

most efficient

Comparing complexity...

An Analysis of Deep Neural Network Models for Practical Applications, 2017.



AlexNet:

Smaller compute, still memory

Comparing complexity...

An Analysis of Deep Neural Network Models for Practical Applications, 2017.



Comparing complexity...

ResNet: Moderate efficiency depending on model, highest accuracy

An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Forward pass time and power consumption



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.

Other Computer Vision Tasks

Semantic Segmentation

Classification + Localization

Object Detection

Instance Segmentation



Other Computer Vision Tasks

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

2D Object Detection



DOG, DOG, CAT

Object categories + 2D bounding boxes

3D Object Detection



Car

Object categories + 3D bounding boxes

This image is CC0 public domain

Semantic Segmentation

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

2D Object Detection



DOG, DOG, CAT

Object categories + 2D bounding boxes 3D Object Detection



Car

Object categories + 3D bounding boxes

This image is CC0 public domain

Semantic Segmentation



This image is CC0 public domain

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window



Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network upsampling: "Unpooling"



In-Network upsampling: "Max Unpooling"



Convolution Transposed

Learnable Upsampling: Transpose Convolution



Input Filter ax a x ay b y az + bx z by bz

Learnable Upsampling: 1D Example

Output contains copies of the filter weighted by the input, summing at where at overlaps in the output

Need to crop one pixel from output to make output exactly 2x input

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

Output



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

2D Object Detection

Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

2D Object Detection



DOG, DOG, CAT

Object categories + 2D bounding boxes

3D Object Detection



Car

Object categories + 3D bounding boxes

This image is CC0 public domain

Classification + Localization









Object Detection as Regression?





CAT: (x, y, w, h)



DOG: (x, y, w, h) DOG: (x, y, w, h) CAT: (x, y, w, h)

DUCK: (x, y, w, h) DUCK: (x, y, w, h)

. . .

Object Detection as Regression?





Each image needs a different number of outputs!

CAT: (x, y, w, h) 4 numbers

DOG: (x, y, w, h) DOG: (x, y, w, h) **16 numbers** CAT: (x, y, w, h)



DUCK: (x, y, w, h) Many DUCK: (x, y, w, h) numbers!

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? NO Cat? NO Background? YES

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? YES Cat? NO Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? YES Cat? NO Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background





Dog? NO Cat? YES Background? NO

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Dog? NO Cat? YES Background? NO

Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!
Region Proposals / Selective Search

- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU



Alexe et al, "Measuring the objectness of image windows", TPAMI 2012 Uijlings et al, "Selective Search for Object Recognition", IJCV 2013 Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014 Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV 2014



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.



Regions of Interest (RoI) from a proposal method (~2k)

> Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.



Regions of Interest (Rol) from a proposal method (~2k)

> Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.







Linear Regression for bounding box offsets

R-CNN: Problems

- Ad hoc training objectives
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
- Training is slow (84h), takes a lot of disk space
- Inference (detection) is slow
 - 47s / image with VGG16 [Simonyan & Zisserman. ICLR15]
 - Fixed by SPP-net [He et al. ECCV14]



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. Slide copyright Ross Girshick, 2015; source. Reproduced with permission.



Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.







Project proposal proposal into 7x7 onto features **Fully-connected** grid, max-pool within each cell layers CNN Hi-res input image: Rol conv features: Hi-res conv features: Fully-connected layers expect 3 x 640 x 480 512 x 7 x 7 low-res conv features: 512 x 20 x 15; with region for region proposal 512 x 7 x 7 proposal **Projected region** proposal is e.g. 512 x 18 x 8 Girshick, "Fast R-CNN", ICCV 2015. (varies per proposal)

Divide projected

Fast R-CNN: Rol Pooling









Figure copyright Ross Girshick, 2015; <u>source</u>. Reproduced with permission.

R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

R-CNN vs SPP vs Fast R-CNN



Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR 2014. He et al, "Spatial pyramid pooling in deep convolutional networks for visual recognition", ECCV 2014 Girshick, "Fast R-CNN", ICCV 2015

Faster R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

Jointly train with 4 losses:

- 1. RPN classify object / not object
- 2. RPN regress box coordinates
- 3. Final classification score (object classes)
- 4. Final box coordinates

Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015 Figure copyright 2015, Ross Girshick; reproduced with permission



Fast<u>er</u> R-CNN:

Make CNN do proposals!



Detection without Proposals: YOLO / SSD



Input image 3 x H x W

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016



Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3 Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
 - (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)

Output: 7 x 7 x (5 * B + C)

Detection without Proposals: YOLO / SSD

Go from input image to tensor of scores with one big convolutional network!



Input image 3 x H x W

Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016



Divide image into grid 7 x 7

Image a set of **base boxes** centered at each grid cell Here B = 3 Within each grid cell:

- Regress from each of the B base boxes to a final box with 5 numbers:
 - (dx, dy, dh, dw, confidence)
- Predict scores for each of C classes (including background as a class)

Output: 7 x 7 x (5 * B + C)

Object Detection: Impact of Deep Learning



Reproduced with permission.





[Antonio Torralba]



"<u>Context Models and Out-of-context Objects</u>" Myung Jin Choi, Antonio Torralba, and Alan S. Willsky Pattern Recognition Letters, 2012.



Rakshith Shetty; Bernt Schiele; Mario Fritz; Not Using the Car to See the Sidewalk: Quantifying and Controlling the Effects of Context in Classification and Segmentation ;IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

Mario Fritz - Elements of Data Science and Artificial Intelligence

Object without **Context**

Context without **Object**



Rakshith Shetty; Bernt Schiele; Mario Fritz; Not Using the Car to See the Sidewalk: Quantifying and Controlling the Effects of Context in Classification and Segmentation ;IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

Mario Fritz - Elements of Data Science and Artificial Intelligence



Rakshith Shetty; Bernt Schiele; Mario Fritz; Not Using the Car to See the Sidewalk: Quantifying and Controlling the Effects of Context in Classification and Segmentation ;IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

Attacks on Object Detectors

- Attacks against ML Models / Object Detectors
- Attacker can craft small change in the image that makes detector
 - Fail
 - Behaves to the liking of adversary



https://bair.berkeley.edu/blog/2017/12/30/yolo-attack/



Advanced Architectures Object Detection & Segmentation Bright & Dark Side of Context

Prof. Dr. Mario Fritz CISPA Helmholtz Center for Information Security