# Lecture 06: CNNs II – Network Architectures

Xuming He

SIST, ShanghaiTech

Fall, 2019

# Summary of CNNs
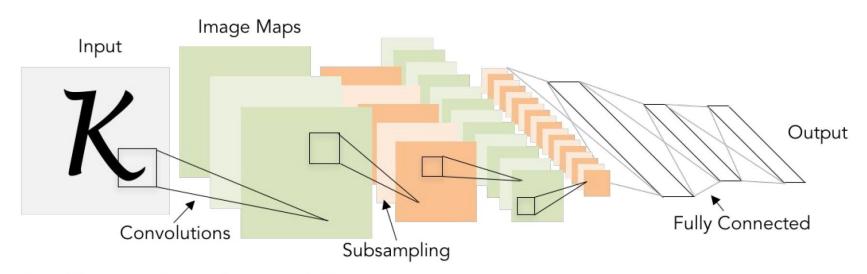
- **CNN properties** [Bronstein et al., 2018]
  - ☐ Convolutional (Translation invariance)
  - ☐ Scale Separation (Compositionality)
  - ☐ Filters localized in space (Deformation Stability)
  - ☐ $O(1)$ parameters per filter (independent of input image size n)
  - ☐ $O(n)$ complexity per layer (filtering done in the spatial domain)
  - ☐ $O(\log n)$ layers in classification tasks

# LeNet-5

- Handwritten digit recognition

[LeCun et al., 1998]



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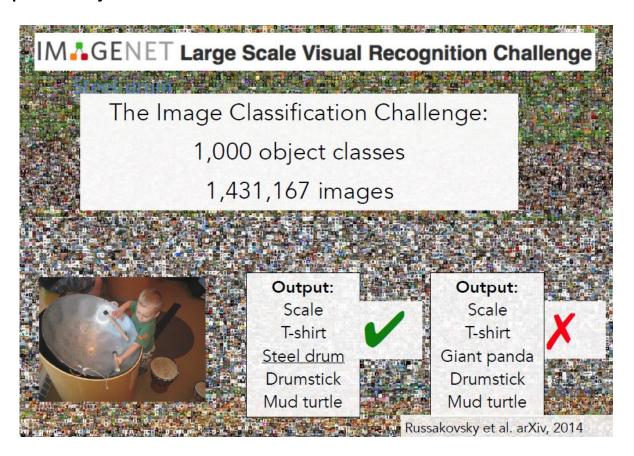i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

# Outline

- ## CNN architectures

  - ☐ Sequential structure: AlexNet/ZFNet/VGGNet

  - ☐ Parallel branches: GoogLeNet

  - ☐ Residual structure: ResNet/DenseNet

*Acknowledgement: Zemel et al's CSC411 and Feifei Li et al's cs231n notes*

# Background: Image/Object Classification

- **Problem Setup**
  - Input: Image
  - Output: Object class



The Image Classification Challenge:
1,000 object classes
1,431,167 images

| Output: | | Output: | |
|---|---|---|---|
| Scale | | Scale | |
| T-shirt | ✔ | T-shirt | ✘ |
| Steel drum | | Giant panda | |
| Drumstick | | Drumstick | |
| Mud turtle | | Mud turtle | |

Russakovsky et al. arXiv, 2014

# Sequential structure

- AlexNet/ZFNet
- VGGNet

# ImageNet (ILSVRC)

Xuming He – CS 280 Deep Learning

# AlexNet

- **Deeper network structure**
  - □ More convolution layers
  - □ Local contrast normalization
  - □ ReLu instead of Tanh
  - □ Dropout as regularization

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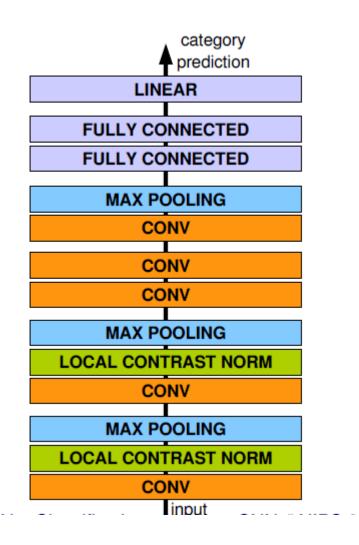CONV1
MAX POOL1
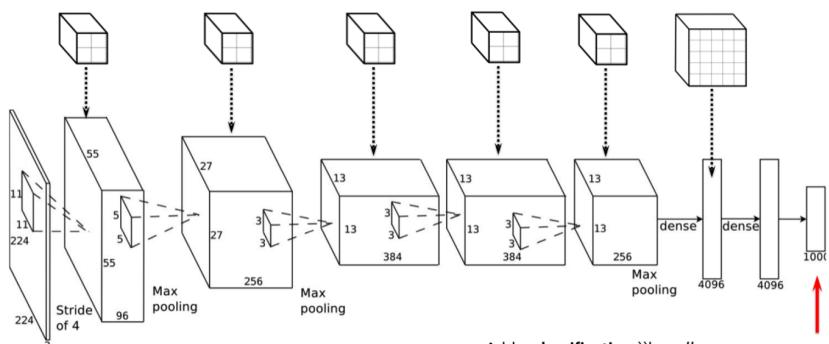NORM1
CONV2
MAX POOL2
NORM2
CONV3
CONV4
CONV5
Max POOL3
FC6
FC7
FC8

category prediction ↑

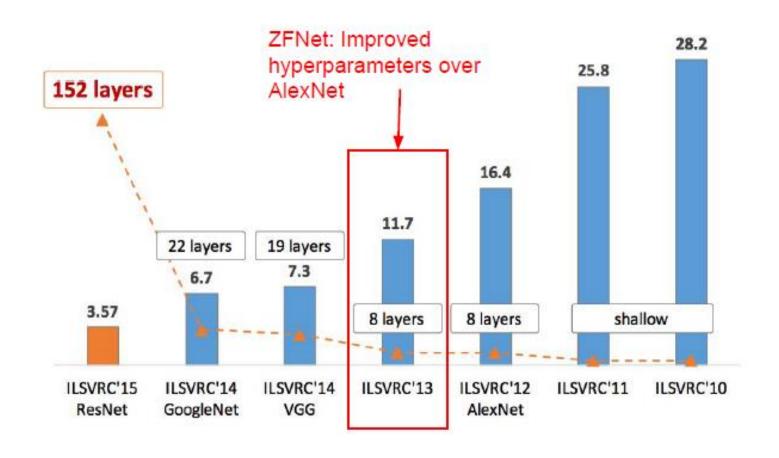| LINEAR |
| FULLY CONNECTED |
| FULLY CONNECTED |
| MAX POOLING |
| CONV |
| CONV |
| CONV |
| MAX POOLING |
| LOCAL CONTRAST NORM |
| CONV |
| MAX POOLING |
| LOCAL CONTRAST NORM |
| CONV |

input

# AlexNet



- 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
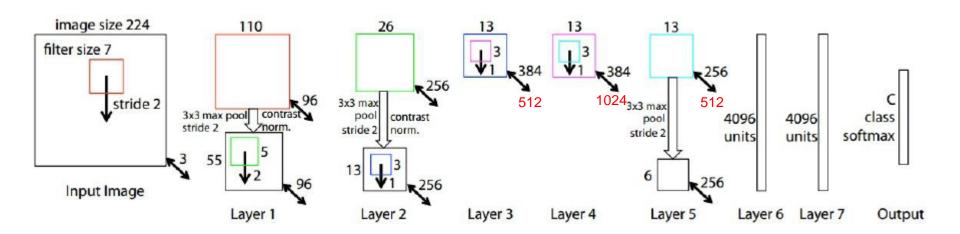- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Add a **classification** ``layer''.

For an input image, the value in a particular dimension of this vector tells you the probability of the corresponding object class.

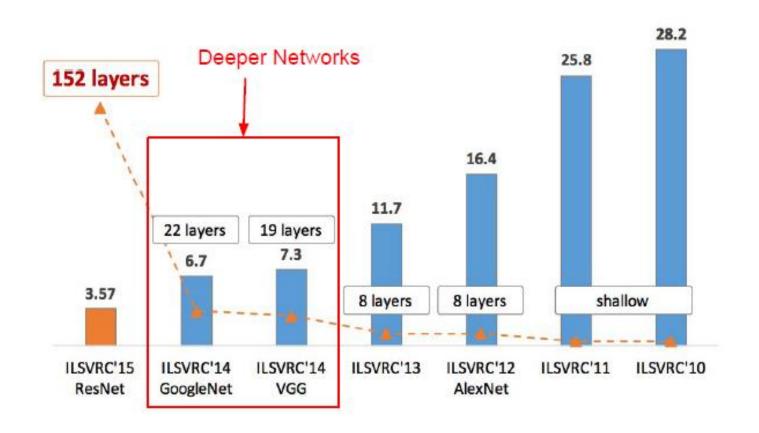# ImageNet (ILSVRC)



Xuming He – CS 280 Deep Learning

# ZFNet



AlexNet but:
CONV1: change from (11x11 stride 4) to (7x7 stride 2)
CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

ImageNet top 5 error: 16.4% -> 11.7%

# ImageNet (ILSVRC)

Xuming He – CS 280 Deep Learning

# VGGNet

## Case Study: VGGNet
*[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**

| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 384 |
| Pool |
| 3x3 conv, 384 |
| Pool |
| 5x5 conv, 256 |
| 11x11 conv, 96 |
| Input |

**VGG16**

| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

**VGG19**

| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

# VGGNet

## Case Study: VGGNet

*[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?

| AlexNet | VGG16 | VGG19 |
|---------|-------|-------|
| | | Softmax |
| | | FC 1000 |
| | Softmax | FC 4096 |
| | FC 1000 | FC 4096 |
| | FC 4096 | Pool |
| | FC 4096 | 3x3 conv, 512 |
| | Pool | 3x3 conv, 512 |
| | 3x3 conv, 512 | 3x3 conv, 512 |
| | 3x3 conv, 512 | 3x3 conv, 512 |
| | 3x3 conv, 512 | Pool |
| Softmax | Pool | 3x3 conv, 512 |
| FC 1000 | 3x3 conv, 512 | 3x3 conv, 512 |
| FC 4096 | 3x3 conv, 512 | 3x3 conv, 512 |
| FC 4096 | 3x3 conv, 512 | 3x3 conv, 512 |
| Pool | Pool | Pool |
| 3x3 conv, 256 | 3x3 conv, 256 | 3x3 conv, 256 |
| 3x3 conv, 384 | 3x3 conv, 256 | 3x3 conv, 256 |
| Pool | Pool | Pool |
| 3x3 conv, 384 | 3x3 conv, 128 | 3x3 conv, 128 |
| Pool | 3x3 conv, 128 | 3x3 conv, 128 |
| 5x5 conv, 256 | Pool | Pool |
| 11x11 conv, 96 | 3x3 conv, 64 | 3x3 conv, 64 |
| Input | 3x3 conv, 64 | 3x3 conv, 64 |
| | Input | Input |

# VGGNet

## Case Study: VGGNet

*[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: $3 * (3^2 C^2)$ vs. $7^2 C^2$ for C channels per layer

| AlexNet |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 384 |
| Pool |
| 3x3 conv, 384 |
| Pool |
| 5x5 conv, 256 |
| 11x11 conv, 96 |
| Input |

| VGG16 |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

| VGG19 |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

# VGGNet

■ ## Parameters

```
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
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
```

Note:

Most memory is in early CONV

Most params are in late FC

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

# VGGNet

- **Summary**

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
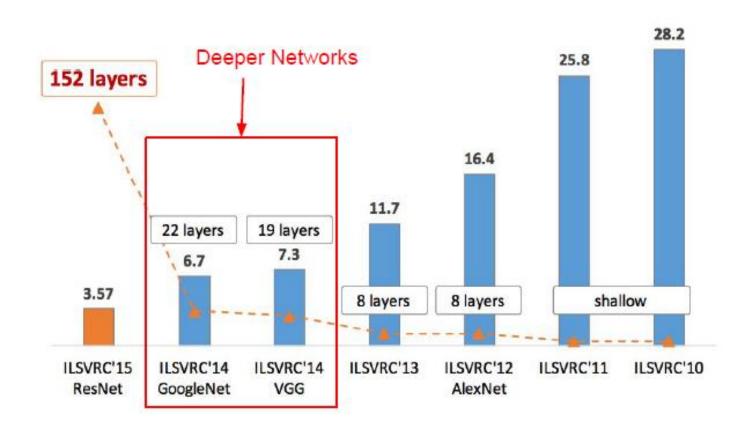- FC7 features generalize well to other tasks



VGG16          VGG19

# Parallel branches

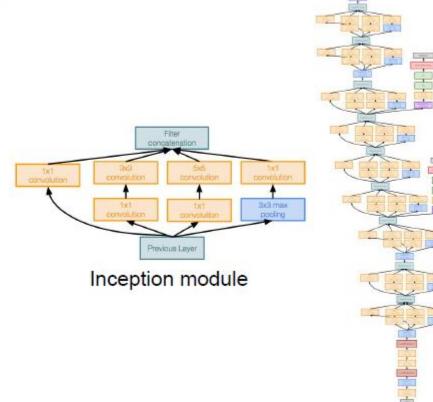- GooLeNet/Inception modules
- NiN (Network in Network)

# ImageNet (ILSVRC)

Xuming He – CS 280 Deep Learning

# GoogLeNet

## Case Study: GoogLeNet

*[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
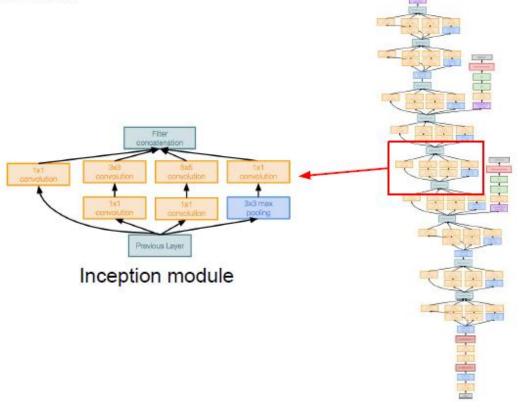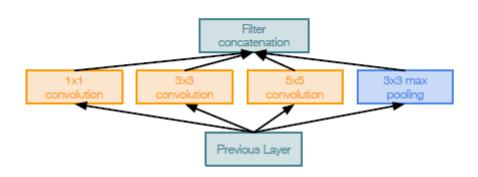- ILSVRC'14 classification winner (6.7% top 5 error)

Inception module

# GoogLeNet

## Case Study: GoogLeNet

*[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



Inception module

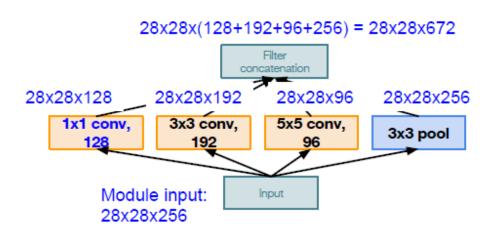Xuming He – CS 280 Deep Learning

# GoogLeNet

- **Inception Module**



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

# GoogLeNet

- ## Inception Module



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

Filter concatenation

28x28x128    28x28x192    28x28x96    28x28x256

1x1 conv, 128    3x3 conv, 192    5x5 conv, 96    3x3 pool

Module input: 28x28x256

Input

Naive Inception module

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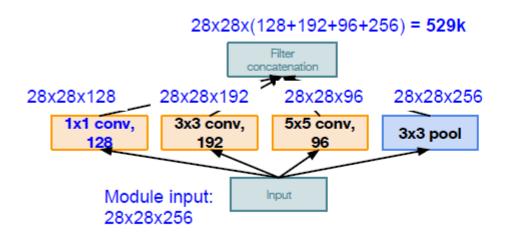[5x5 conv, 96]  28x28x96x5x5x256
**Total: 854M ops**

Very expensive compute

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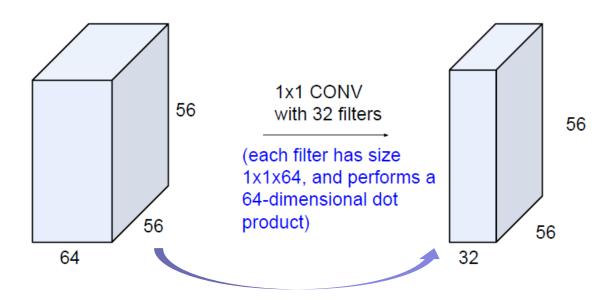
# GoogLeNet

- Inception Module

28x28x(128+192+96+256) = **529k**

Filter concatenation

28x28x128    28x28x192    28x28x96    28x28x256

1x1 conv, 128    3x3 conv, 192    5x5 conv, 96    3x3 pool

Module input: 28x28x256

Input

Naive Inception module

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

# GoogLeNet

- Bottleneck layer



1x1 CONV
with 32 filters

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

# GoogLeNet

- Inception Module



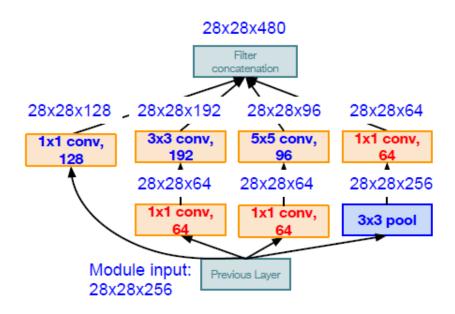1x1 conv "bottleneck" layers

Naive Inception module

Inception module with dimension reduction

# GoogLeNet

- Inception Module



28x28x480

Filter concatenation

28x28x128  28x28x192  28x28x96  28x28x64

1x1 conv, 128   3x3 conv, 192   5x5 conv, 96   1x1 conv, 64

28x28x64   28x28x64   28x28x256

1x1 conv, 64   1x1 conv, 64   3x3 pool

Module input: 28x28x256

Previous Layer

Inception module with dimension reduction

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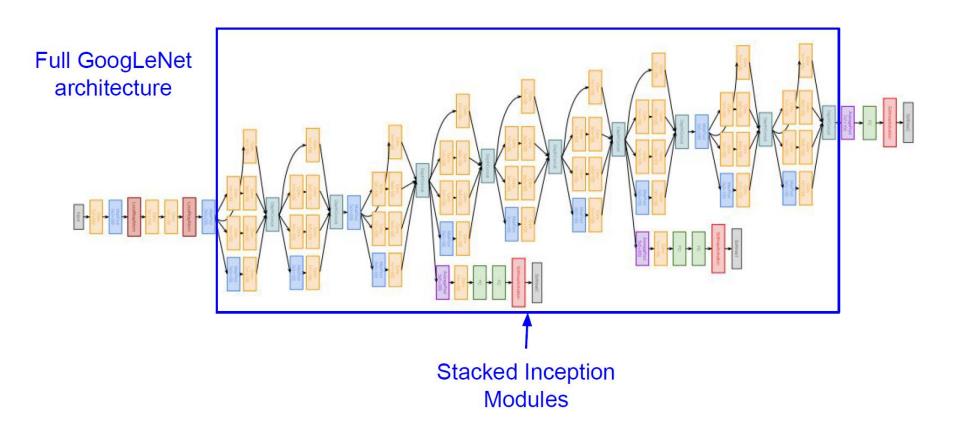[1x1 conv, 64]  28x28x64x1x1x256
**Total: 358M ops**

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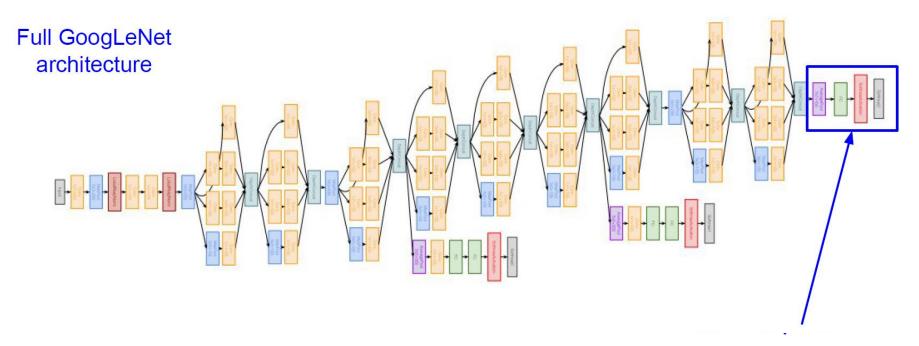Bottleneck can also reduce depth after pooling layer

# GoogLeNet

- **Overall network structure**



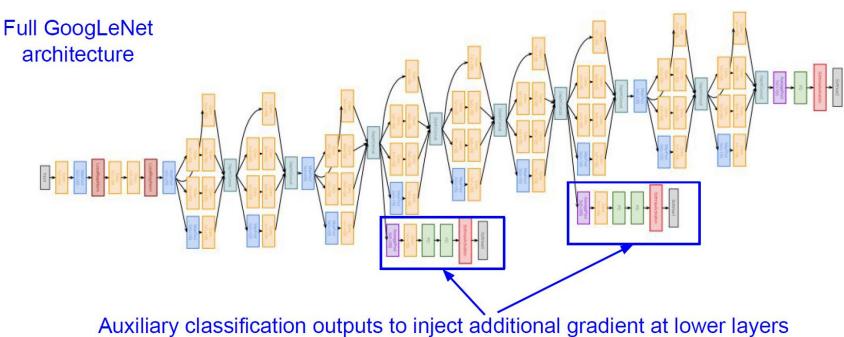Full GoogLeNet architecture

Stem Network:
Conv-Pool-
2x Conv-Pool

# GoogLeNet

- Overall network structure



Full GoogLeNet architecture

Stacked Inception Modules

# GoogLeNet

- Overall network structure

**Full GoogLeNet architecture**

**Classifier output
(removed expensive FC layers!)**
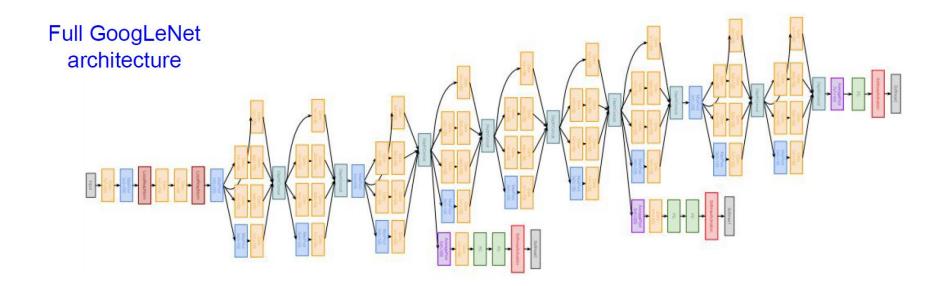
# GoogLeNet

- **Overall network structure**



Full GoogLeNet architecture

Auxiliary classification outputs to inject additional gradient at lower layers
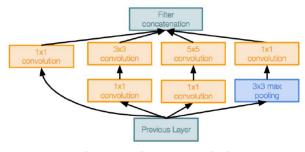(AvgPool-1x1Conv-FC-FC-Softmax)

# GoogLeNet

- Overall network structure

**Full GoogLeNet architecture**



22 total layers with weights (including each parallel layer in an Inception module)

# GoogLeNet

- **Summary**

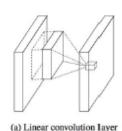Deeper networks, with computational efficiency

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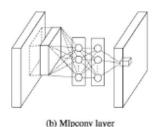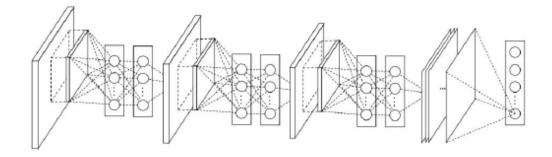- ILSVRC'14 classification winner (6.7% top 5 error)



Inception module

# Other: Network in Network (NiN)

- Mlpconv layer with "micronetwork" within each conv layer to compute more abstract features for local patches
- Micronetwork uses multilayer perceptron (FC, i.e. 1x1 conv layers)
- Precursor to GoogLeNet and ResNet "bottleneck" layers
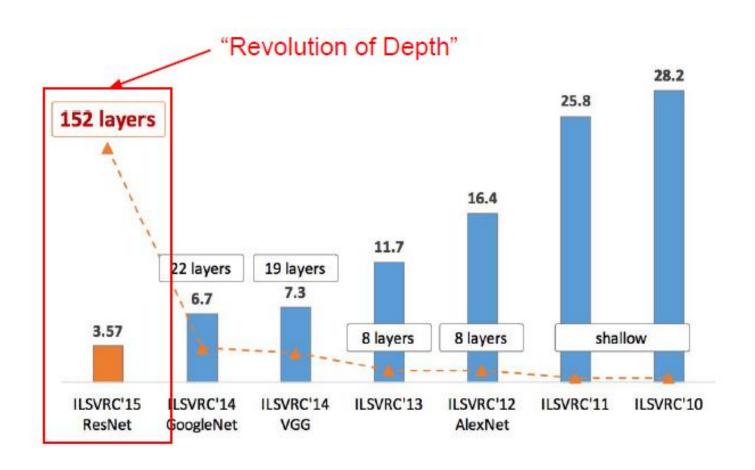- Philosophical inspiration for GoogLeNet



(a) Linear convolution layer

(b) Mlpconv layer

Xuming He – CS 280 Deep Learning

# Residual structure

- ResNet
- DenseNet

Xuming He – CS 280 Deep Learning

# ImageNet (ILSVRC)



"Revolution of Depth"

152 layers

22 layers — 6.7

19 layers — 7.3

11.7 — 8 layers

16.4 — 8 layers

25.8 — shallow

28.2

3.57

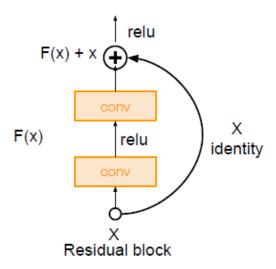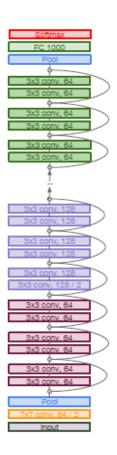| ILSVRC'15 ResNet | ILSVRC'14 GoogleNet | ILSVRC'14 VGG | ILSVRC'13 | ILSVRC'12 AlexNet | ILSVRC'11 | ILSVRC'10 |

Xuming He – CS 280 Deep Learning

# ResNet

## Case Study: ResNet

[He et al., 2015]

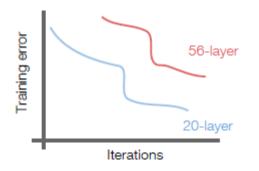Very deep networks using residual connections

- 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!



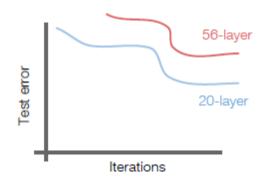Residual block

# ResNet

- What happens when stacking deeper plain conv layers?



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

# ResNet

- Hypothesis:
  - The problem is an optimization problem, deeper models are harder to optimize
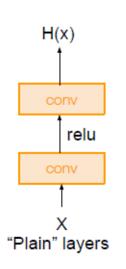
The deeper model should be able to perform at least as well as the shallower model.

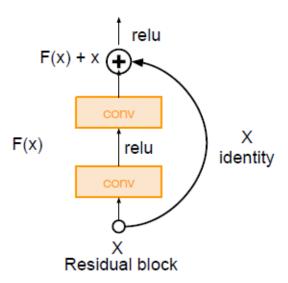A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping.

# ResNet

- Solution:
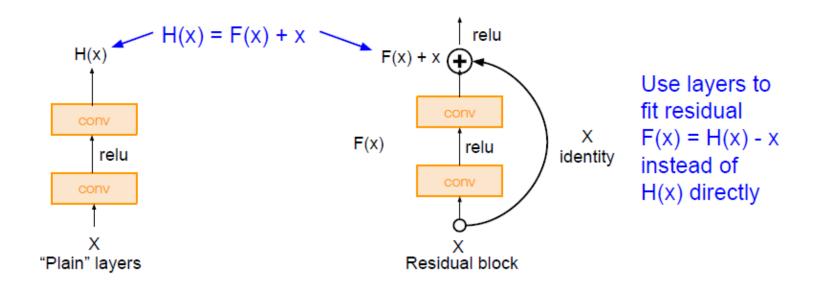  - Use network layers to fit a residual mapping



H(x)

conv

relu

conv

X
"Plain" layers

relu

F(x) + x ⊕

conv

F(x)          relu

conv

X
identity

X
Residual block

# ResNet

- **Solution:**
  - ☐ Use network layers to fit a residual mapping

# ResNet

## Case Study: ResNet

*[He et al., 2015]*

**Full ResNet architecture:**
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)

$F(x) + x$   relu

3x3 conv

$F(x)$   relu

3x3 conv

X

**Residual block**

X identity

Softmax
FC 1000
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, /2
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128, /2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64, / 2
Input

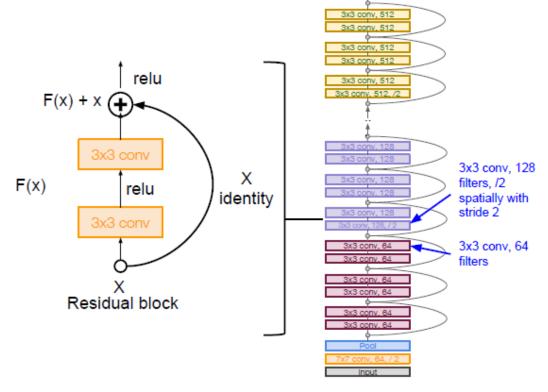3x3 conv, 128 filters, /2 spatially with stride 2

3x3 conv, 64 filters

# ResNet



## Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
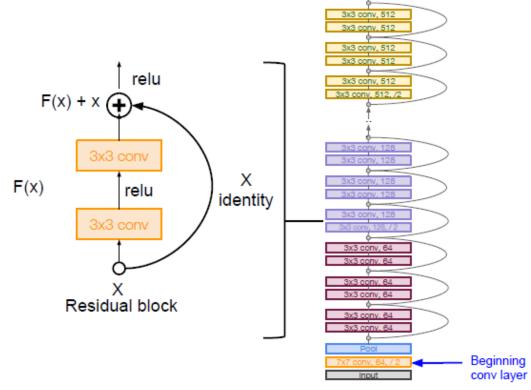- Additional conv layer at the beginning

relu

$F(x) + x$ ⊕

3x3 conv

$F(x)$  relu

3x3 conv

X

X identity

Residual block

Softmax
FC 1000
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, /2
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128, /2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64, / 2 ← Beginning conv layer
Input

# ResNet

## Case Study: ResNet

*[He et al., 2015]*

**Full ResNet architecture:**
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
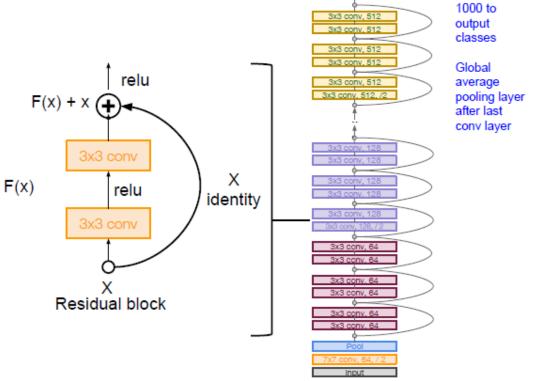- Additional conv layer at the beginning
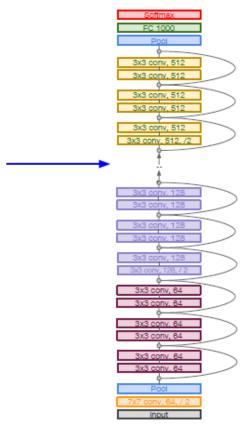- No FC layers at the end (only FC 1000 to output classes)

$F(x) + x$    relu    ⊕ ←   X identity

3x3 conv

$F(x)$    relu

3x3 conv

X

**Residual block**

No FC layers besides FC 1000 to output classes

Global average pooling layer after last conv layer

Softmax
FC 1000
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512, /2
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128, /2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64, /2
Input

# ResNet



## Case Study: ResNet

[He et al., 2015]

Total depths of 34, 50, 101, or 152 layers for ImageNet

# ResNet

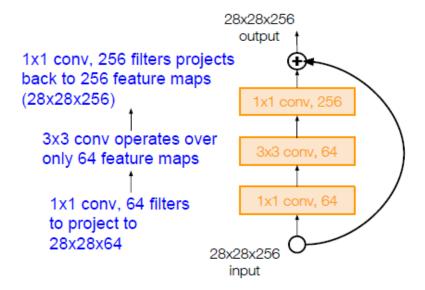## Case Study: ResNet

*[He et al., 2015]*

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

1x1 conv, 256 filters projects back to 256 feature maps (28x28x256)

3x3 conv operates over only 64 feature maps

1x1 conv, 64 filters to project to 28x28x64

28x28x256 output

1x1 conv, 256

3x3 conv, 64

1x1 conv, 64

28x28x256 input

# ResNet

- ## Training details

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

# ResNet

■ Results

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
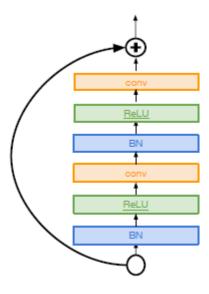
MSRA @ ILSVRC & COCO 2015 Competitions

- **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)

# Other: Identity Mappings in ResNet

- Improved ResNet block design from creators of ResNet
- Creates a more direct path for propagating information throughout network (moves activation to residual mapping pathway)
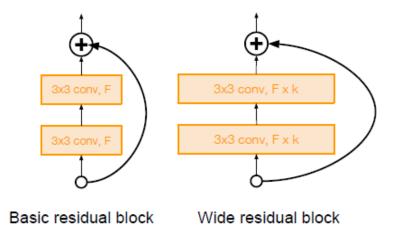- Gives better performance

# Other: Wide ResNets

## Wide Residual Networks

*[Zagoruyko et al. 2016]*

- Argues that residuals are the important factor, not depth
- User wider residual blocks (F x k filters instead of F filters in each layer)
- 50-layer wide ResNet outperforms 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)
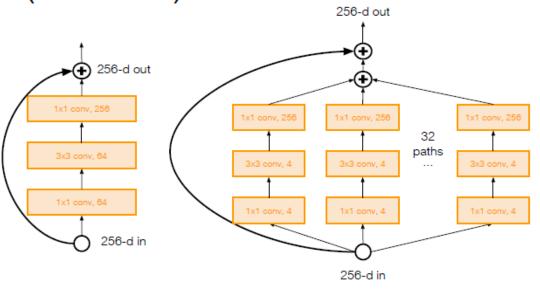


Basic residual block          Wide residual block

# Other: ResNeXt

## Aggregated Residual Transformations for Deep Neural Networks (ResNeXt)

*[Xie et al. 2016]*

- Also from creators of ResNet
- Increases width of residual block through multiple parallel pathways ("cardinality")
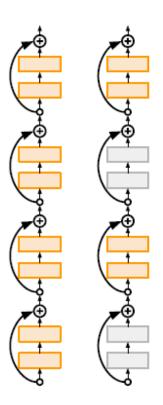- Parallel pathways similar in spirit to Inception module

# Other:ResNet with Stochastic Depth

## Deep Networks with Stochastic Depth

*[Huang et al. 2016]*

- Motivation: reduce vanishing gradients and training time through short networks during training
- Randomly drop a subset of layers during each training pass
- Bypass with identity function
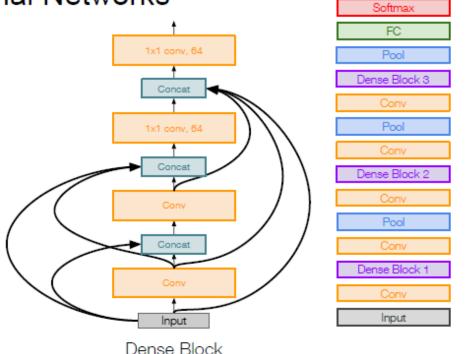- Use full deep network at test time

# DenseNet

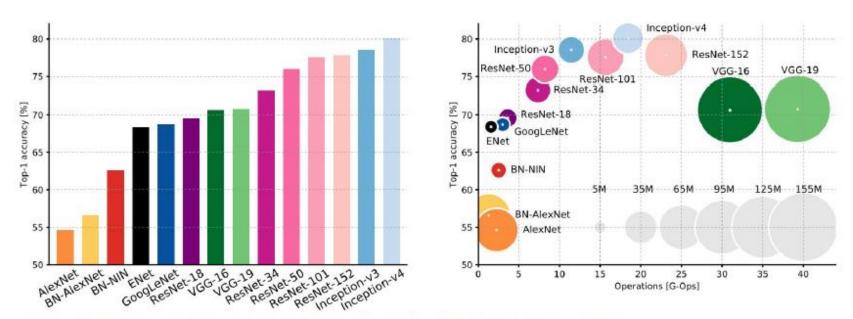## Densely Connected Convolutional Networks

*[Huang et al. 2017]*

- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



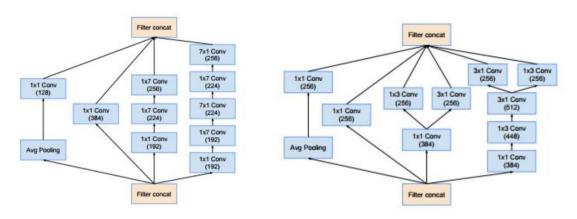Dense Block

# Model complexity



Comparing complexity...

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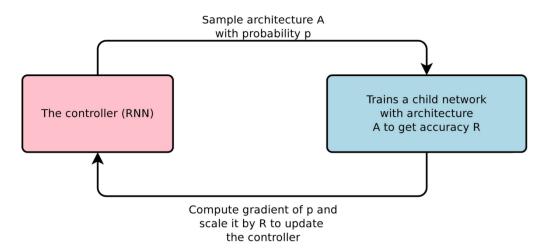
# Network Architecture

- **Problems with network architecture**
  - Designing NA is hard
  - Lots of human efforts go into tuning them
  - Not a lot of intuition into how to design them well
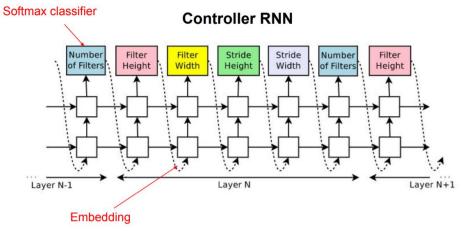  - Can we learn good architectures automatically?



Two layers from the famous Inception V4 computer vision model.
Szegedy et al, 2017

# Network Architecture

- ## Neural architecture search (Zoph and Le, ICLR 2016)



Sample architecture A with probability p

The controller (RNN) → Trains a child network with architecture A to get accuracy R

Compute gradient of p and scale it by R to update the controller

- □ CNN example:



Softmax classifier

**Controller RNN**

| Number of Filters | Filter Height | Filter Width | Stride Height | Stride Width | Number of Filters | Filter Height |

Layer N-1          Layer N          Layer N+1

Embedding

# Summary

- **CNNs for image/object classification**

    - VGG, GoogLeNet, ResNet all in wide use, available in model zoos
    - ResNet current best default
    - Trend towards extremely deep networks
    - Significant research centers around design of layer / skip connections and improving gradient flow
    - Even more recent trend towards examining necessity of depth vs. width and residual connections

- **Next time:**
    - Network learning algorithms