

Convolutional Neural Networks

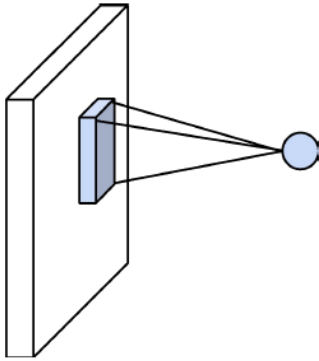
Longbin Jin



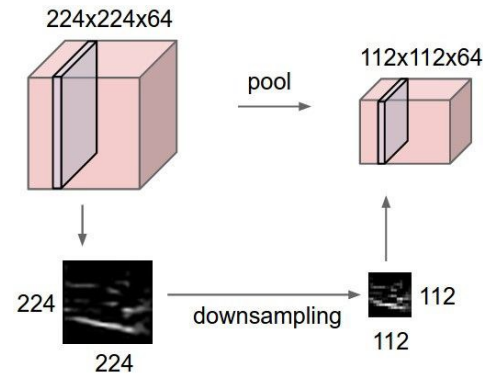
Artificial Intelligence
& Computer Vision
L a b o r a t o r y

Components of CNNs

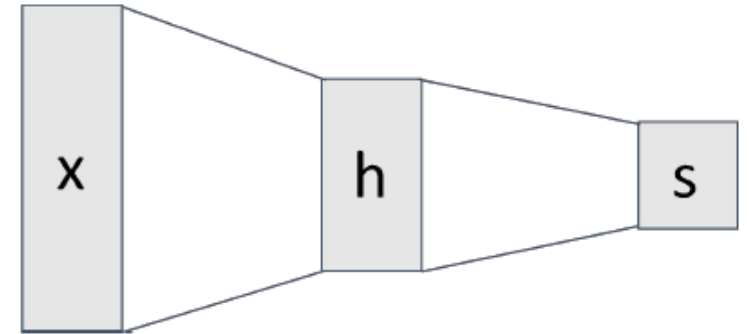
Convolution Layers



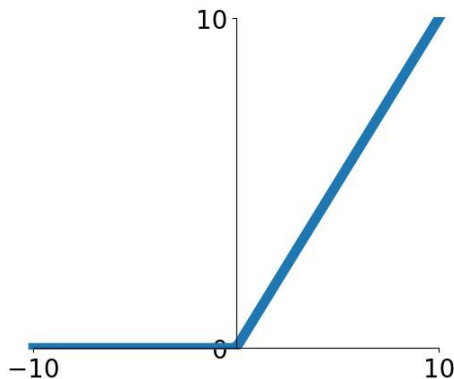
Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

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

Question: How should we put them together?

CNN Architectures

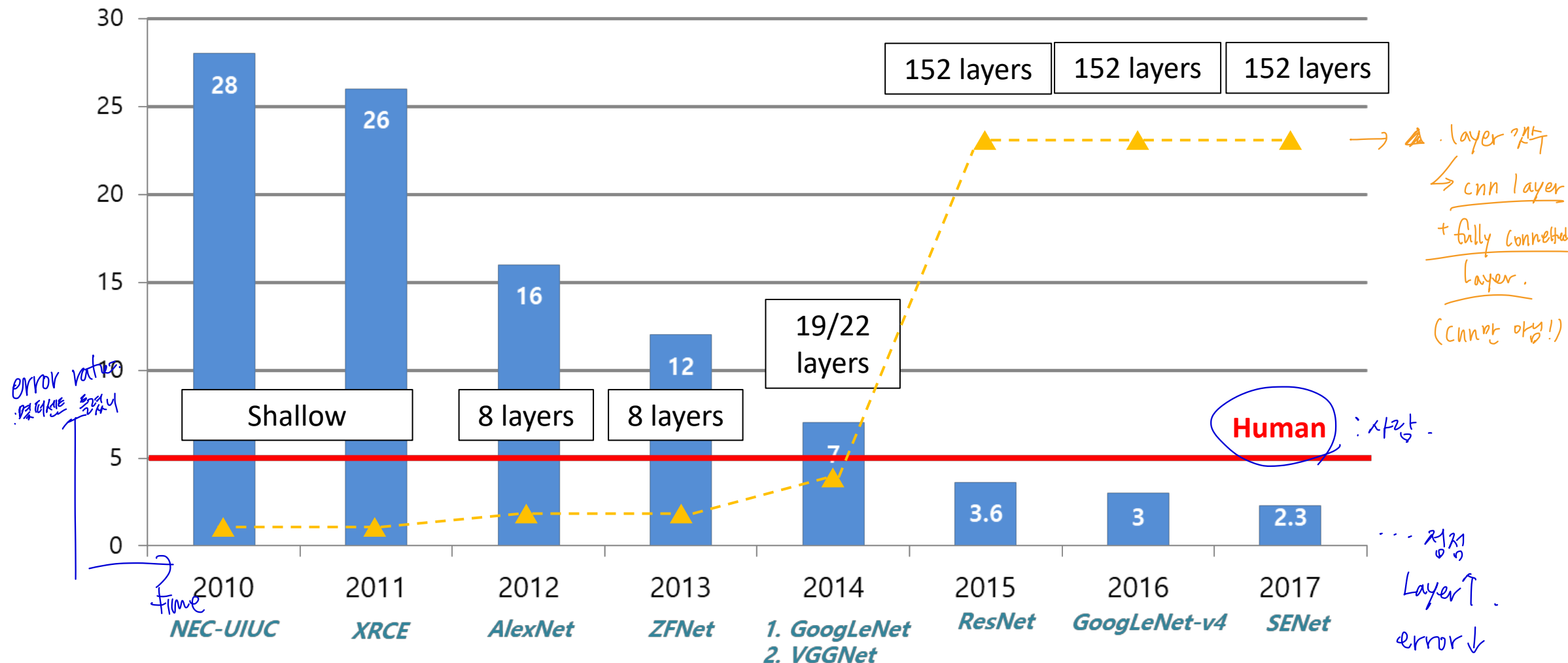
- AlexNet
- ZFNet
- VGG
- GoogleNet
- ResNet
- SENet
- MobileNets

ImageNet *(dataset)*

- ImageNet-1k
 - 1000 classes
 - 1000 images/class
- ImageNet-21k
 - 21841 classes
 - 14197122 images



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



AlexNet

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

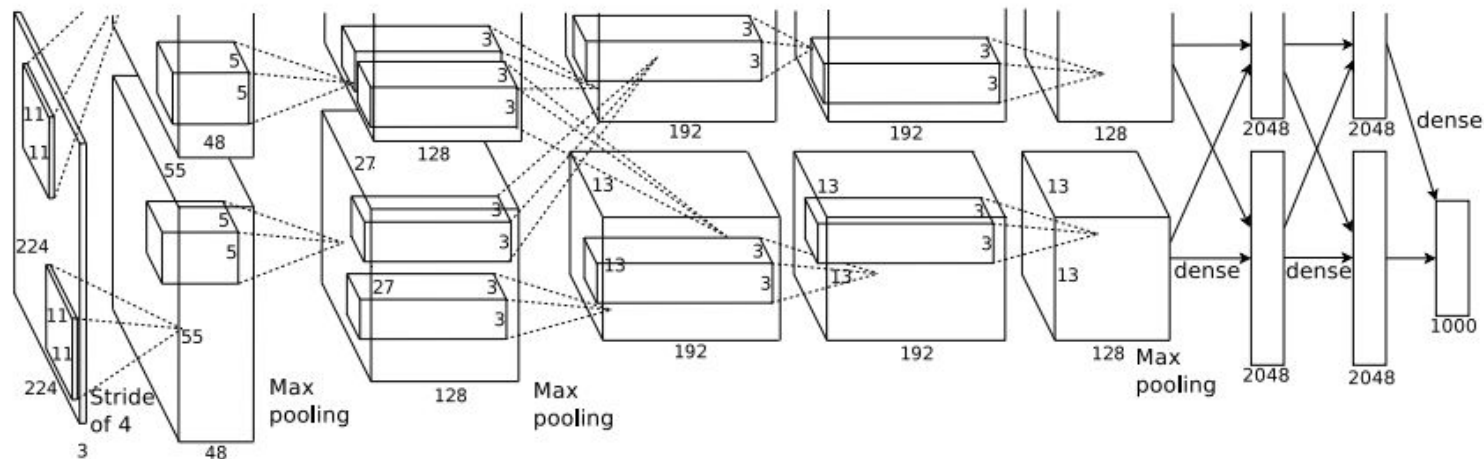
Max POOL3

FC6

FC7

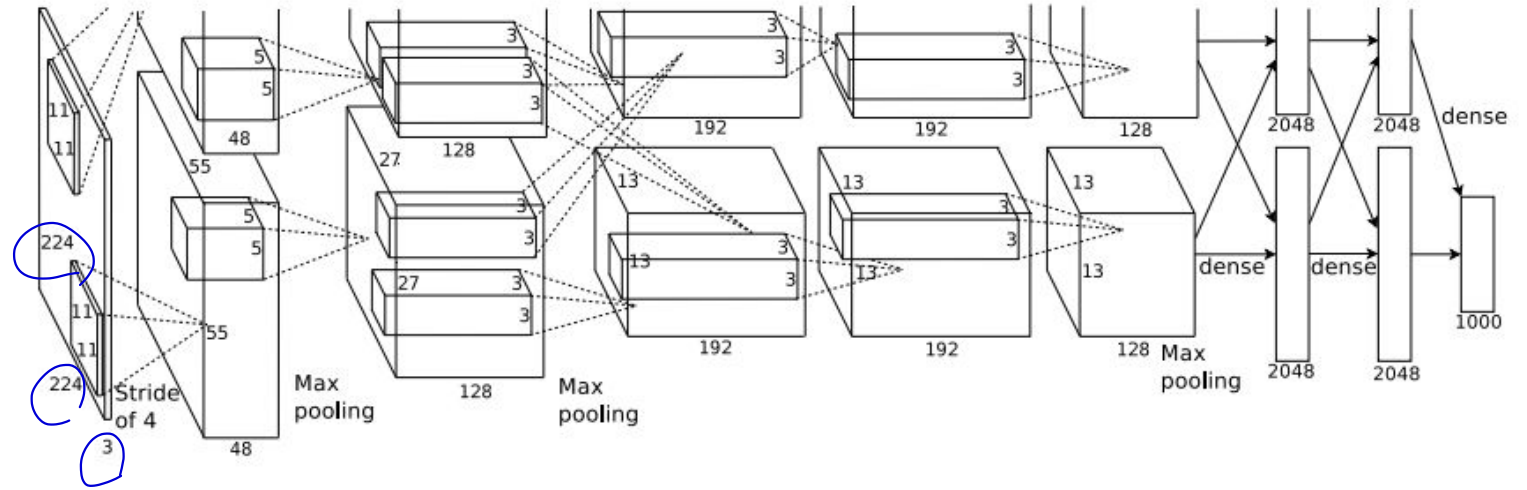
FC8

8x11 layer
3x layer AlexNet
belu 8x11.



AlexNet

<314444 Layer>



Input: 227x227x3 images → 원래 Image는 224x224x3 인데,
정확히 모르겠다? 227x227x3.

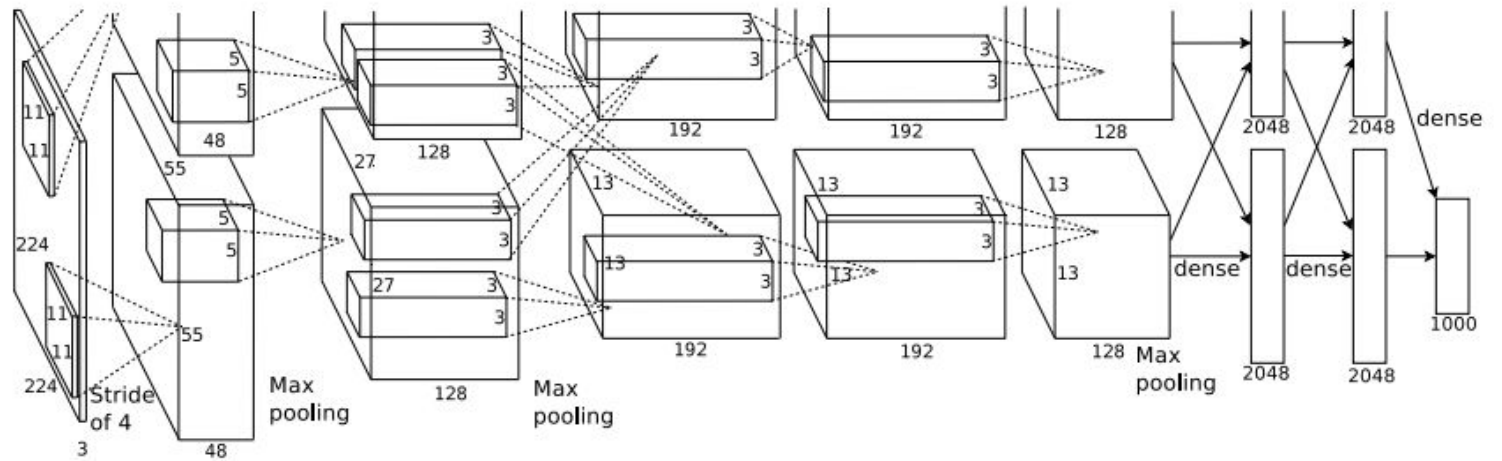
$$W' = (W - F + 2P) / S + 1$$

First layer (CONV1): 96 11x11 filters applied at stride 4

padding = 0.

Q: what is the output volume size?

AlexNet

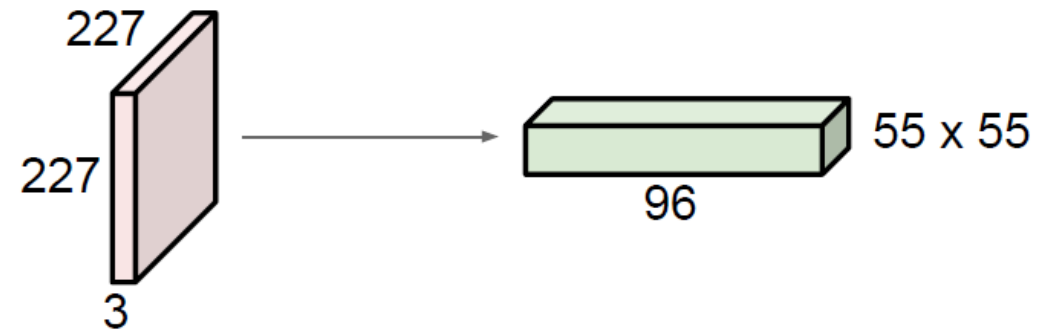


Input: 227x227x3 images

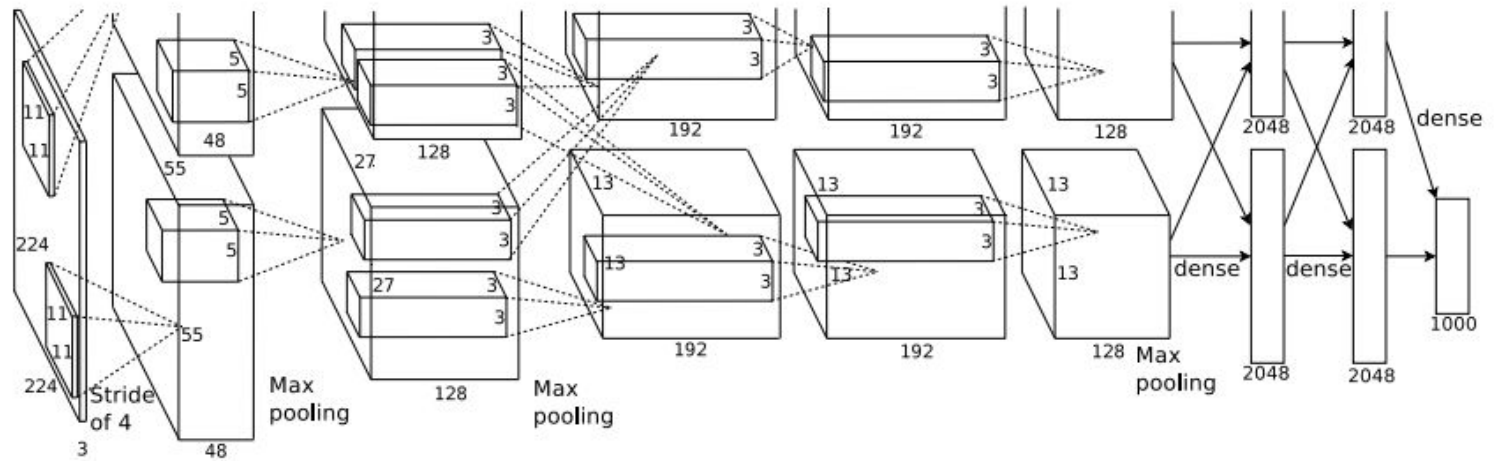
$$W' = (W - F + 2P) / S + 1$$

First layer (CONV1): 96 11x11 filters applied at stride 4

Output volume 55x55x96



AlexNet



Input: 227x227x3 images

$\therefore 11 \times 11 \times 3$

채널 차원 \times \Rightarrow 입력 채널 차원과 같게 설정하면 됨.

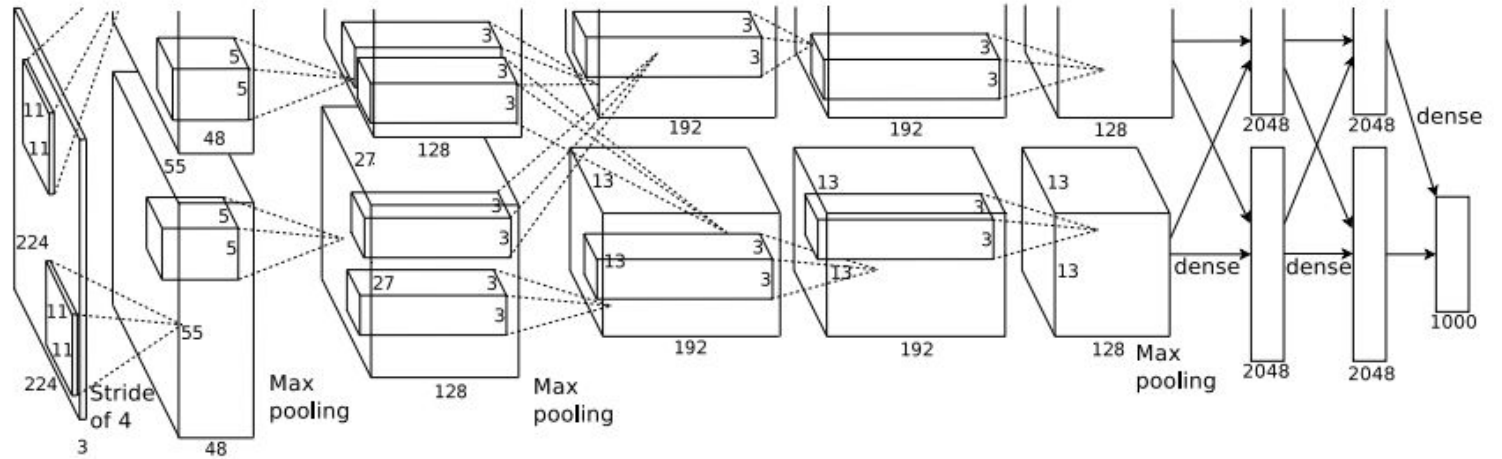
$$W' = (W - F + 2P) / S + 1$$

First layer (CONV1): 96 11×11 filters applied at stride 4

Output volume **55x55x96**

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

AlexNet



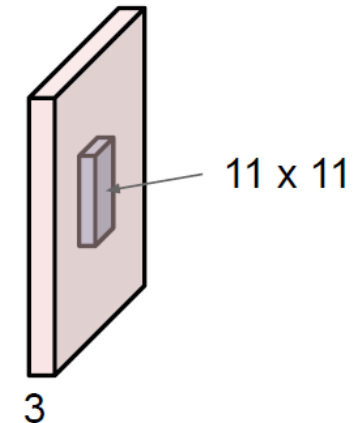
Input: 227x227x3 images

First layer (CONV1): 96 11x11 filters applied at stride 4

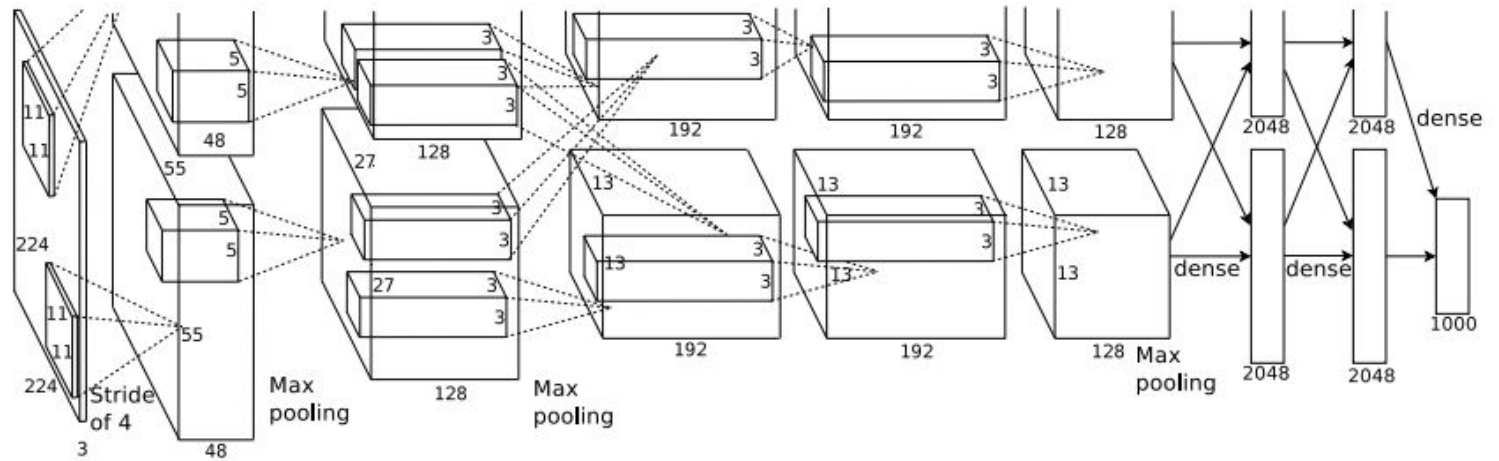
Output volume **55x55x96**

Parameters: $(11*11*3 + 1)*96 = 35K$

$$W' = (W - F + 2P) / S + 1$$



AlexNet



Input: 227x227x3 images

After CONV1: 55x55x96

output

input

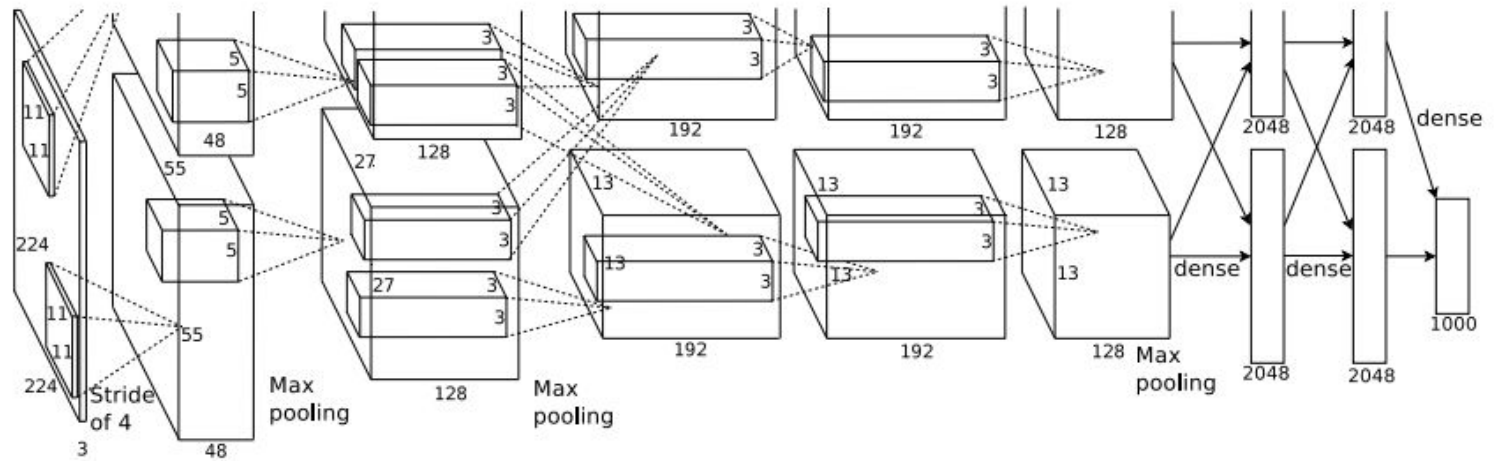
Second layer (POOL1): 3x3 filters applied at stride 2

Q: what is the output volume size?

conv filter output

$$W' = (W - F + 2P) / S + 1$$

AlexNet



Input: 227x227x3 images

After CONV1: 55x55x96

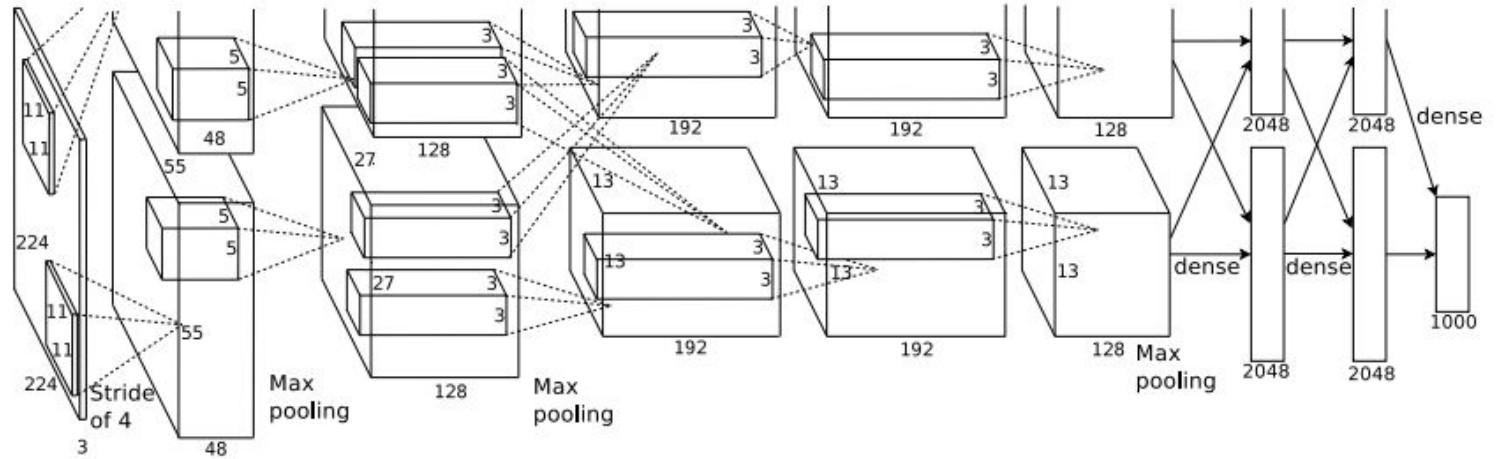
$$W' = (W - F + 2P) / S + 1$$

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

Output volume: 27x27x96 → Pooling은 채널 차원을 절반으로 X

Q: what is the number of parameters in this layer? ⇒ Pooling layer에는 parameter 존재 X!

AlexNet



Input: 227x227x3 images

After CONV1: 55x55x96

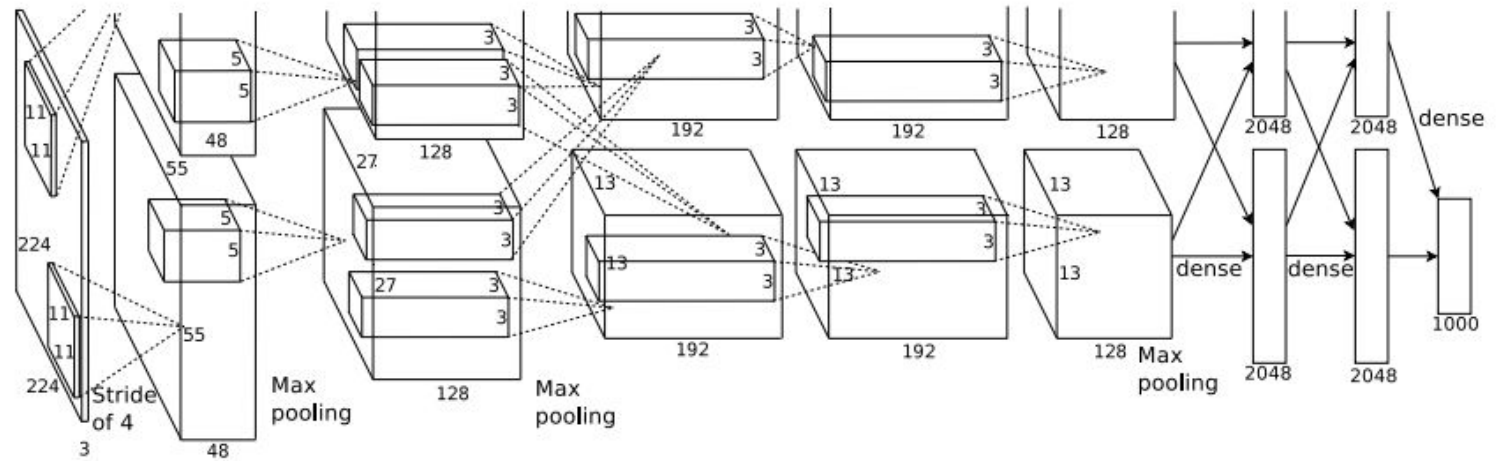
Second layer (POOL1): 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0!

$$W' = (W - F + 2P) / S + 1$$

AlexNet



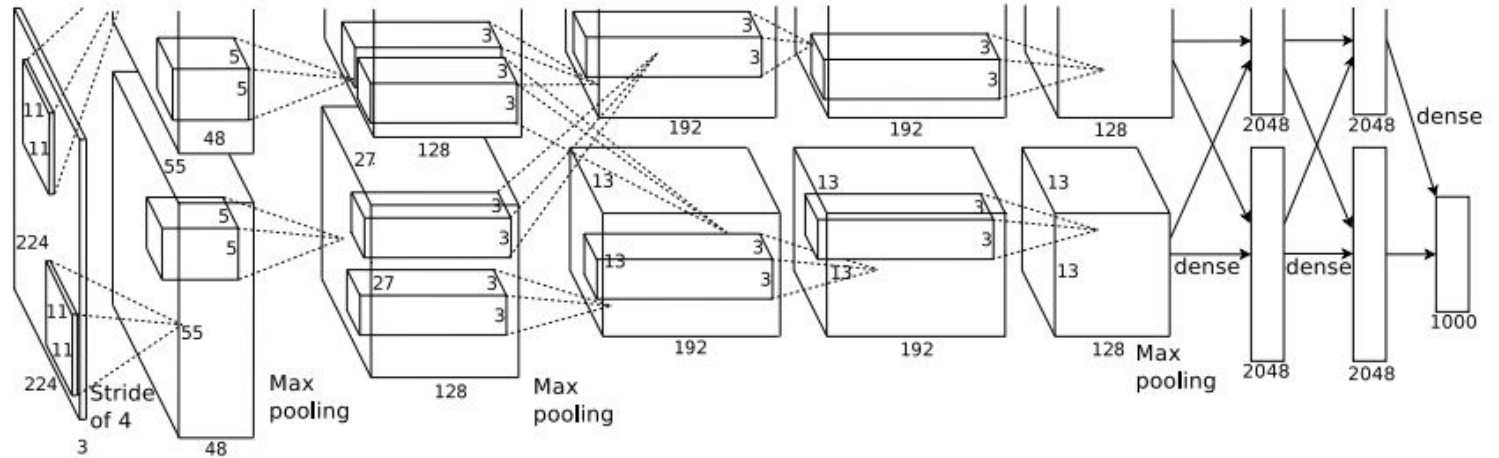
Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

...

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)

→ 인접한 픽셀 사이에 중복되는 특징이 있어서,
큰 영역에서 추출함.

filter size ↓↓.

→ 2중 특징이 아니어서 점점 작아
영역을 추출함.

* stride 1, pad 1 이면

입력 size == 출력 size.

6x6x256을 넣을 때,
9천 몇천개를 4천 번 계산하고,
다시 1000 클래스로 넣음.
(3차원 아님)

중복값, parameter 값
계산할 수 있어야함!

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

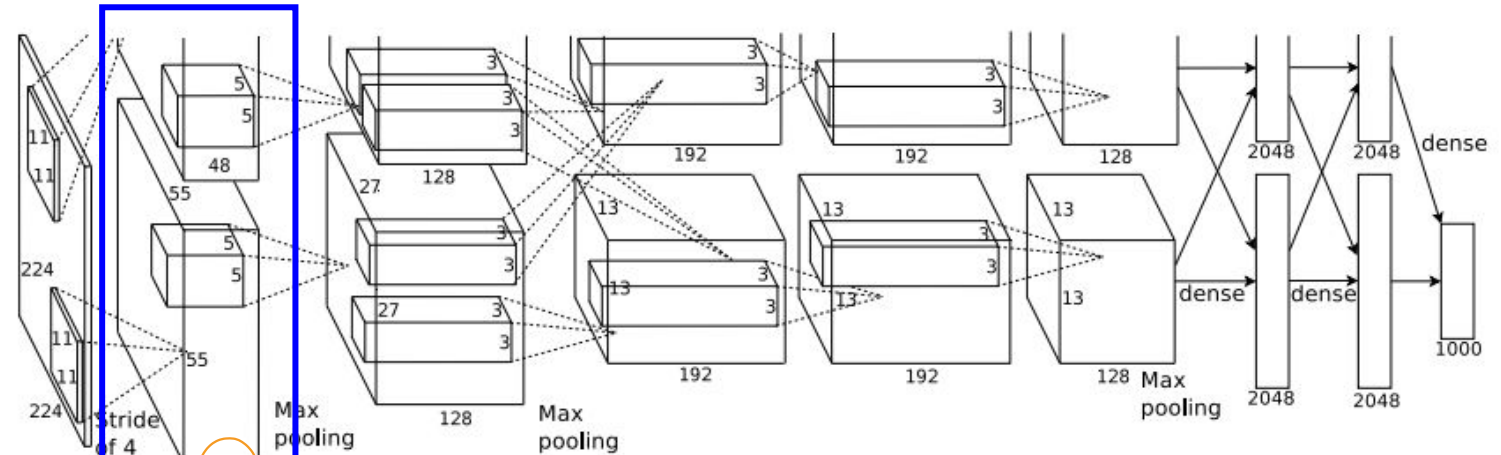
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



[55x55x48] x 2 \Rightarrow 이 지점에서는 GPU 용량이 적어서 나눠서 계산함.

(Group Conv 라는 생각 형태로 무방함)

Historical note: Trained on GTX 580 GPU with only 3 GB of memory. Network spread across 2 GPUs, half the neurons (feature maps) on each GPU.

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

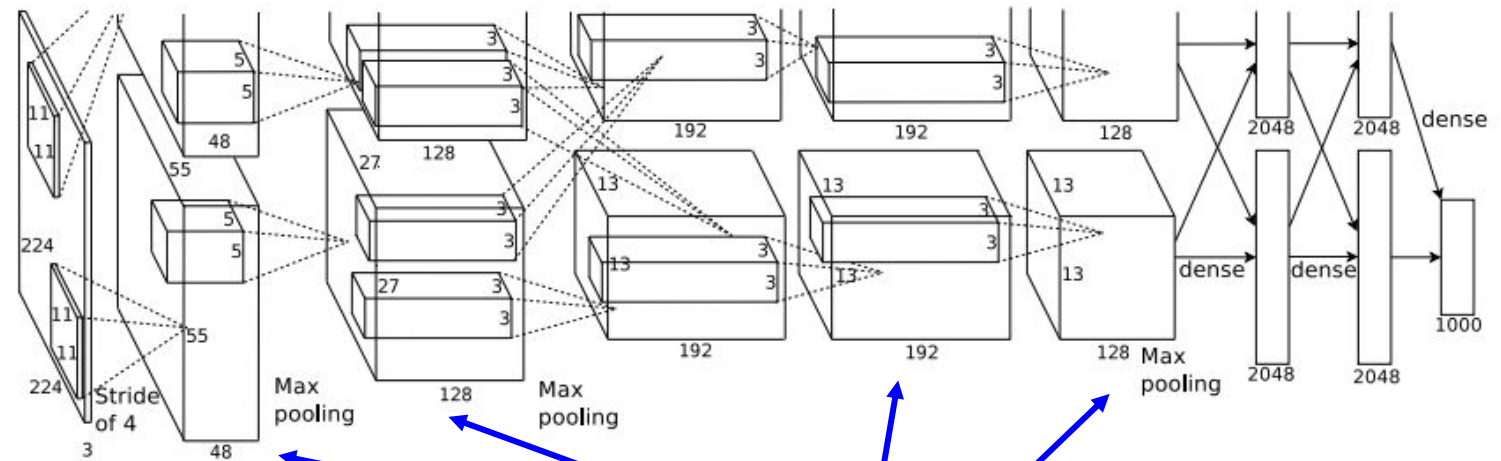
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

[1000] **FC8**: 1000 neurons (class scores)



CONV1, CONV2, CONV4, CONV5:
Connections only with feature maps
on same GPU

group conv?

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

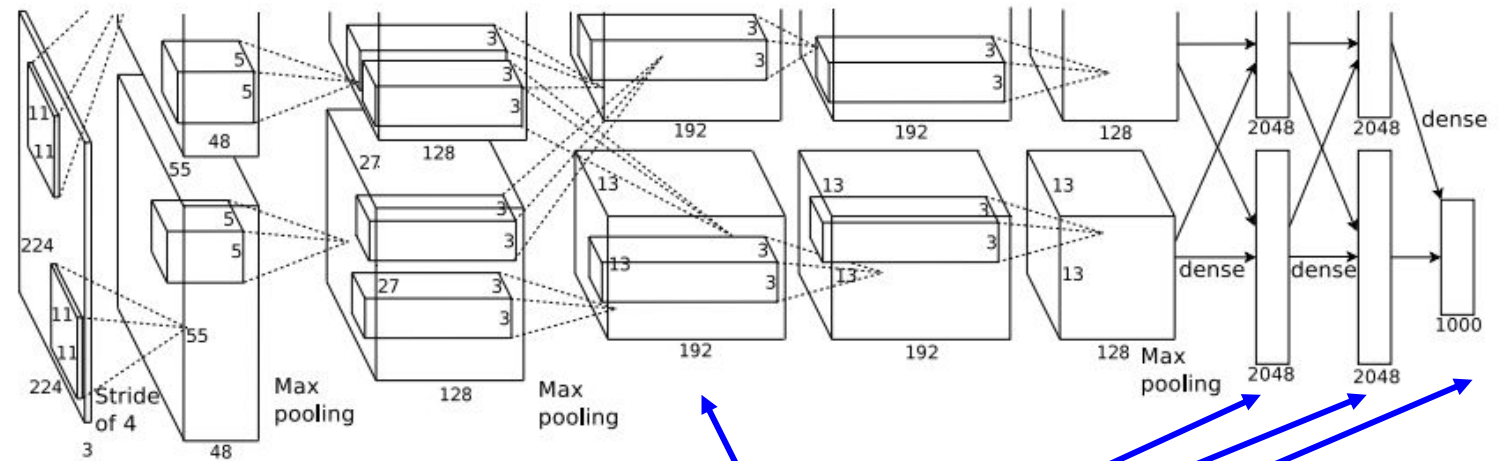
[13x13x256] **CONV5**: 256 3x3 filters at stride 1, pad 1

[6x6x256] **MAX POOL3**: 3x3 filters at stride 2

[4096] **FC6**: 4096 neurons

[4096] **FC7**: 4096 neurons

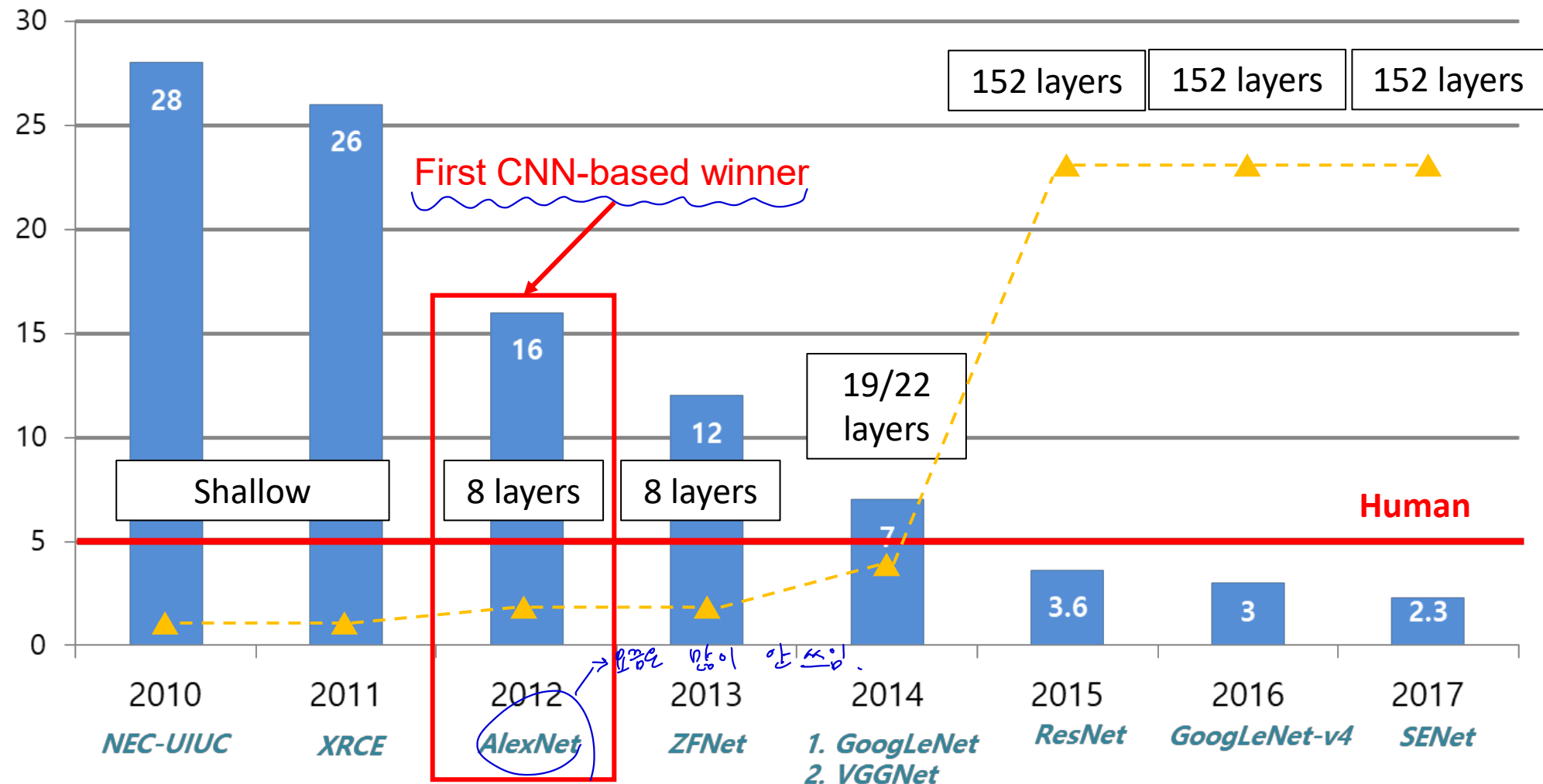
[1000] **FC8**: 1000 neurons (class scores)



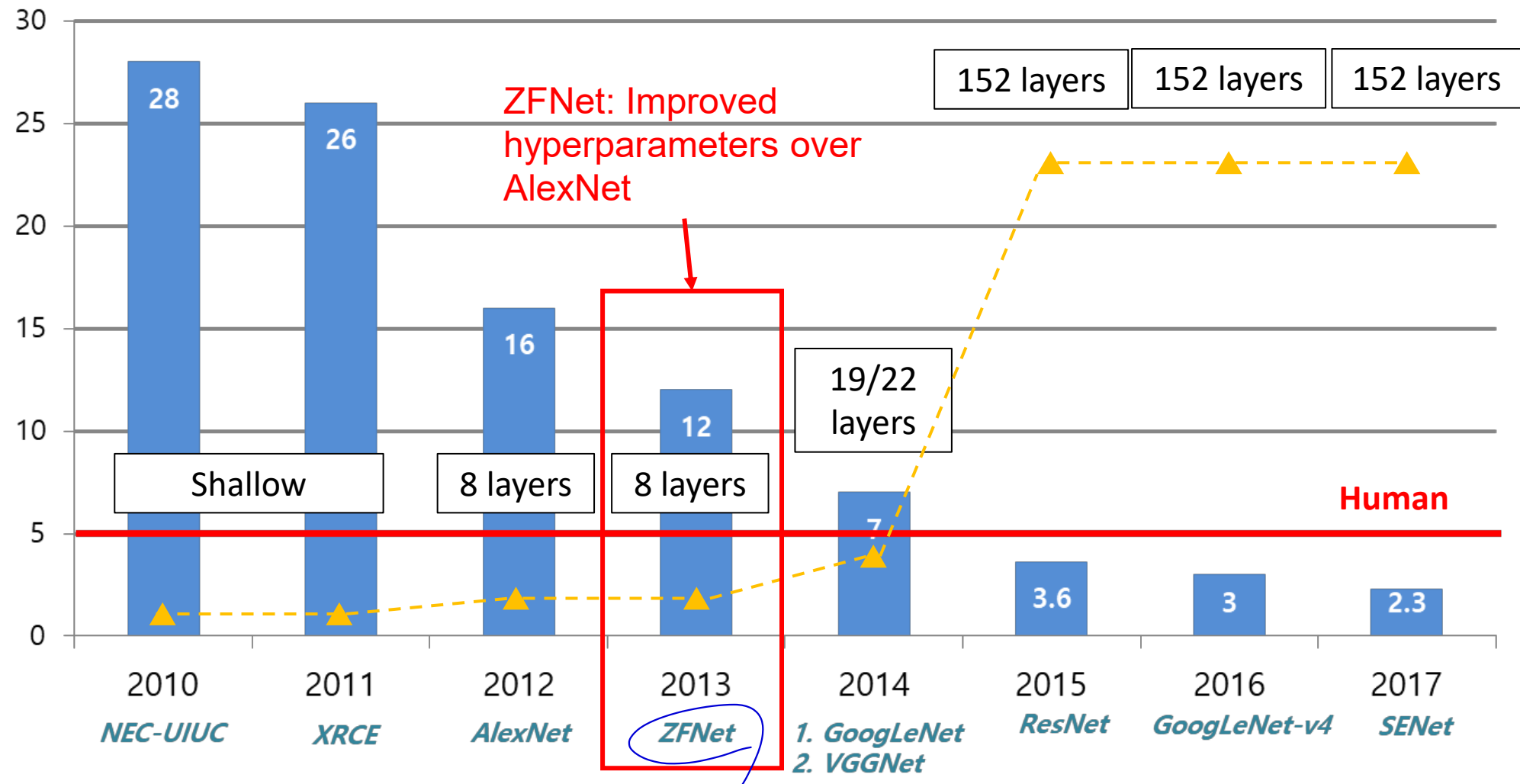
CONV3, FC6, FC7, FC8:
Connections with all feature maps in
preceding layer, communication
across GPUs

group x

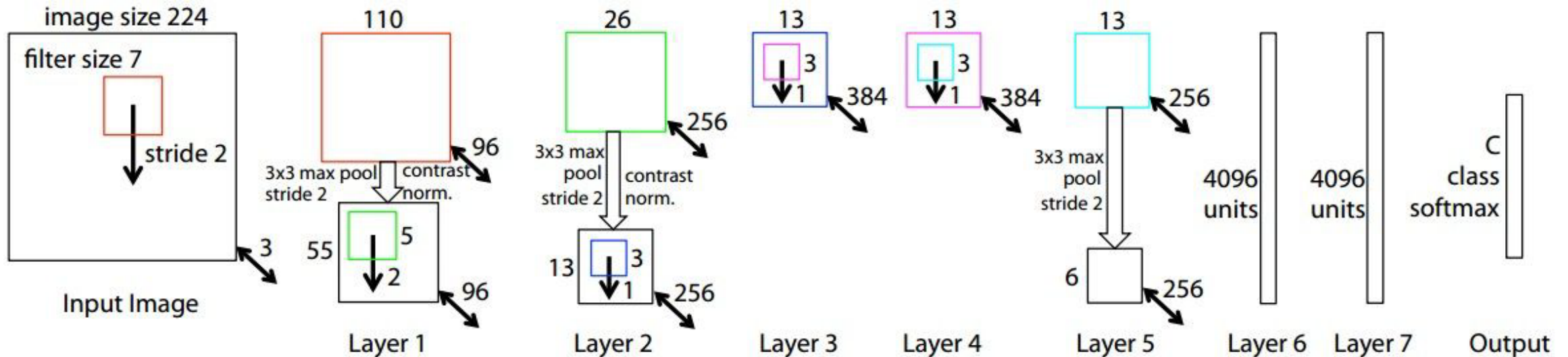
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ZFNet



AlexNet but:

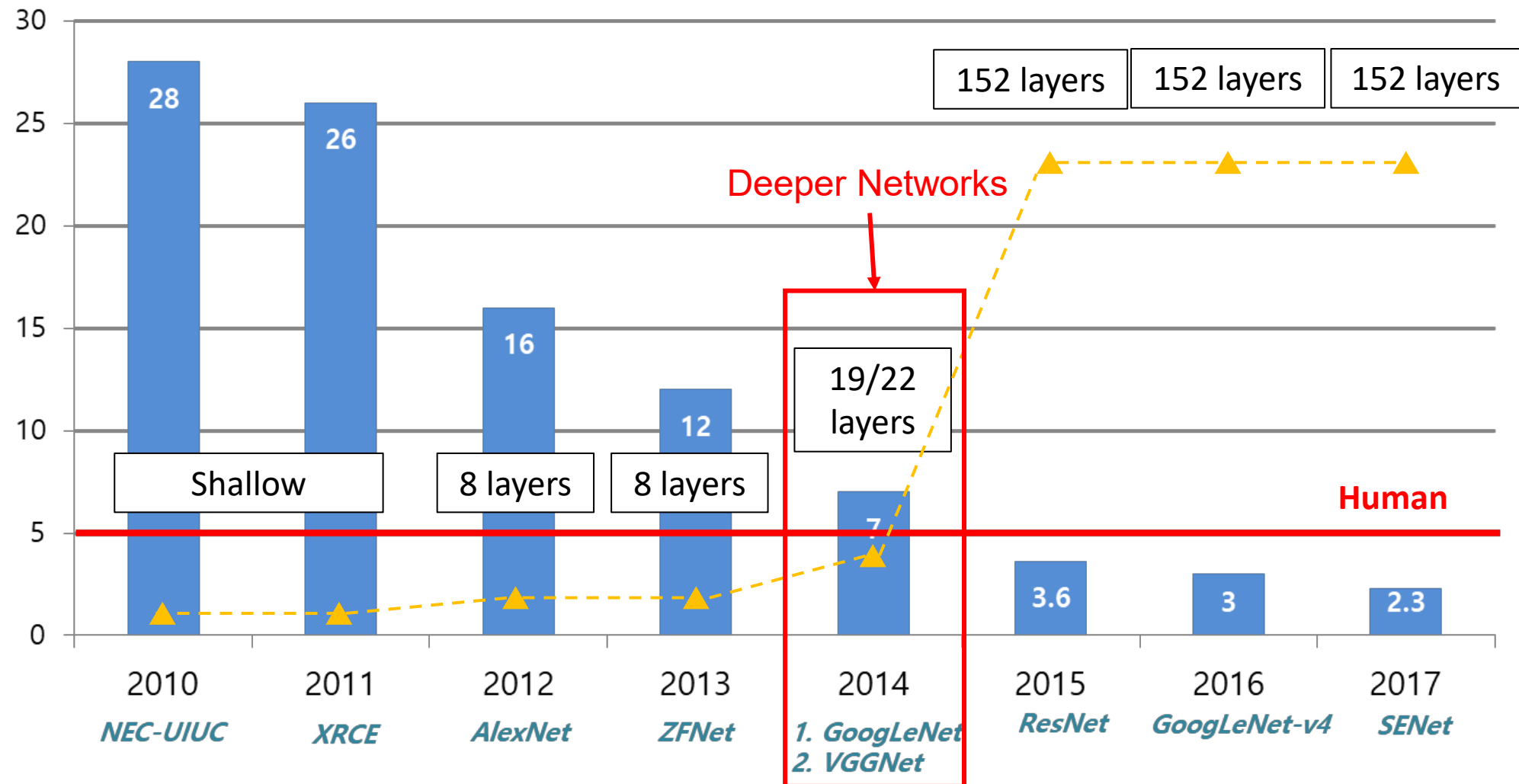
CONV1: change from (11x11 stride 4) to (7x7 stride 2)

CONV2: change from (5x5 stride 1) to (5x5 stride 2)

ImageNet top 5 error: 16.4% -> 11.7%

→ AlexNet
ZFNet.
→ 두 개만 다른 점
이거 2배씩.

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



VGGNet

→ AlexNet 사용했던 11x11 같은
큰 필터 사이즈 사용 X.

Small filters, Deeper networks

8 layers (AlexNet) ↘ 2x

-> 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

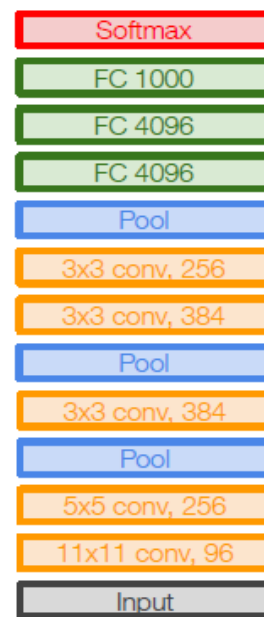


VGGNet

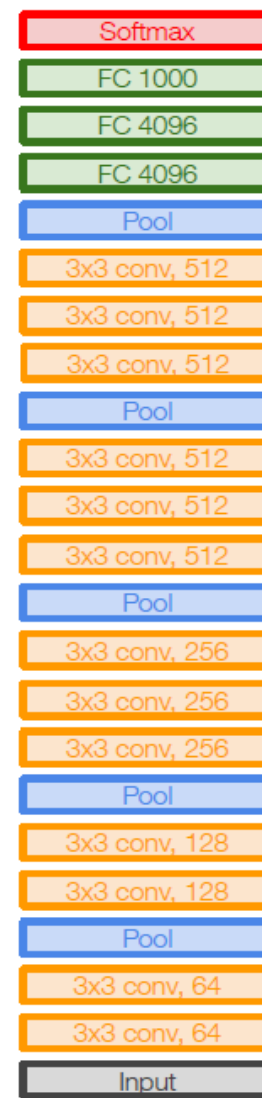
Q: Why use smaller filters? (3x3 conv)

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

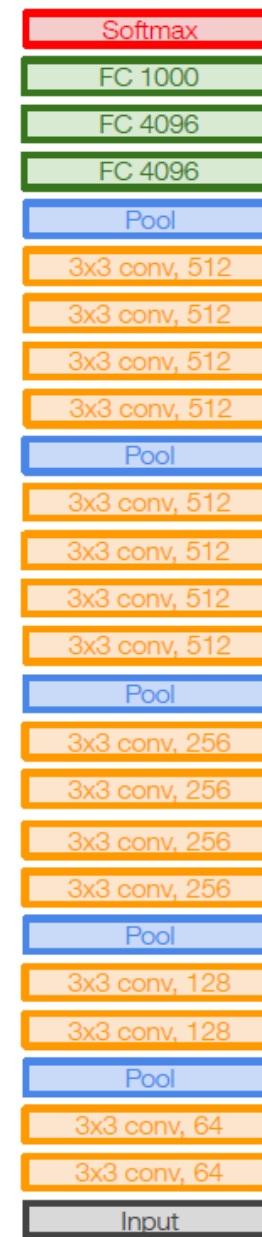
↓
7x7 conv



AlexNet



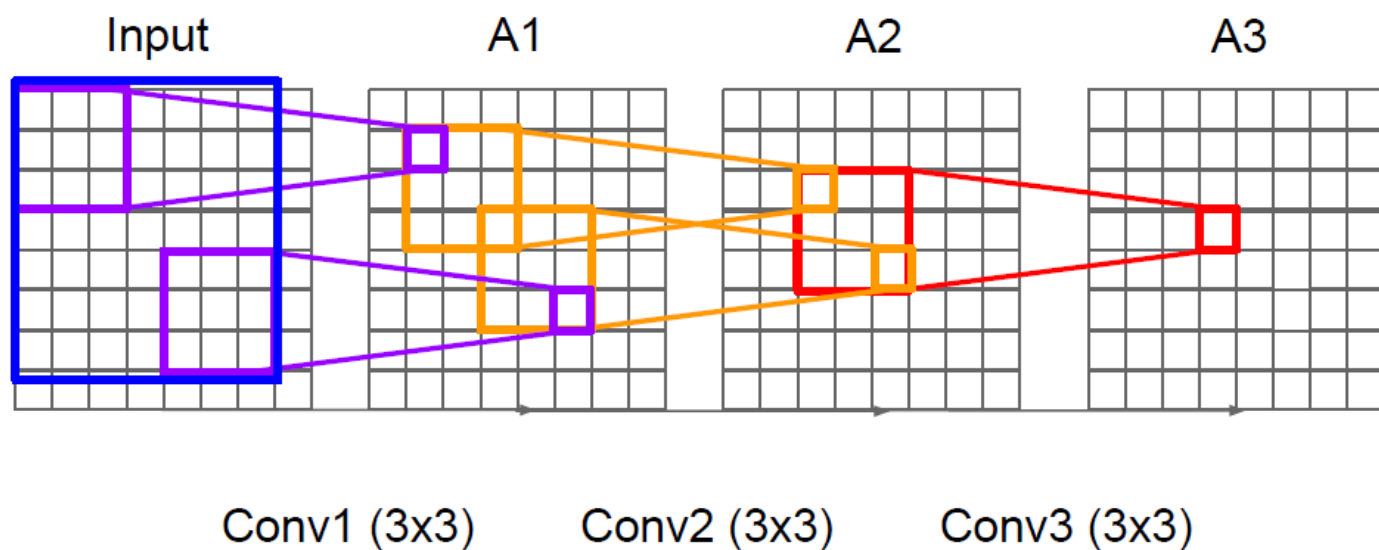
VGG16



VGG19

VGGNet

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



VGGNet

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

Layer를 사용하면
성능이 좋아지는 이유.



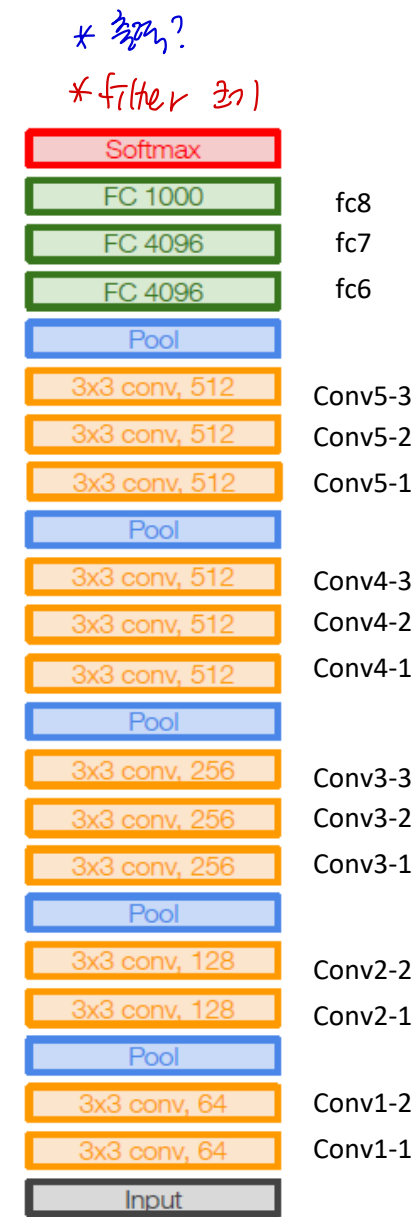
중간중간 relu (non-linear) 가 필요해서,
1x1 한번 사용한다
3x3을 7번 사용한다
2번 필요해서,
그게 param은 2배 정도



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



VGG16

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: $24M * 4 \text{ bytes} \sim 96MB / \text{image}$ (only forward! ~ 2 for bwd)

TOTAL params: 138M parameters

Note:

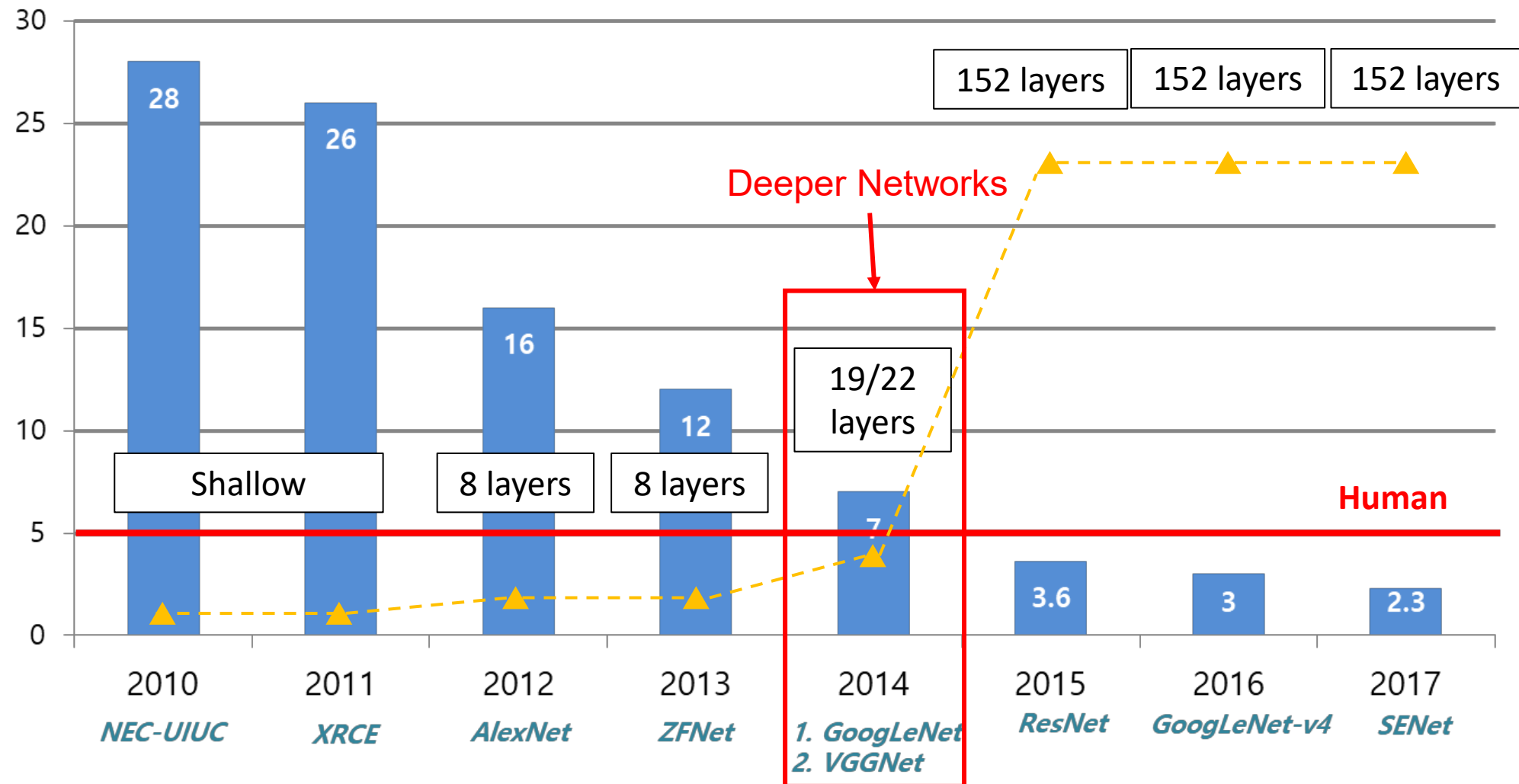
Most memory is in early CONV

여기서 제일 크다.
점점 작아짐.

Most params are in late FC

conv param 다 합쳐도,
fully connected param을 못이김.
so, 잔정이다. 이해됨.

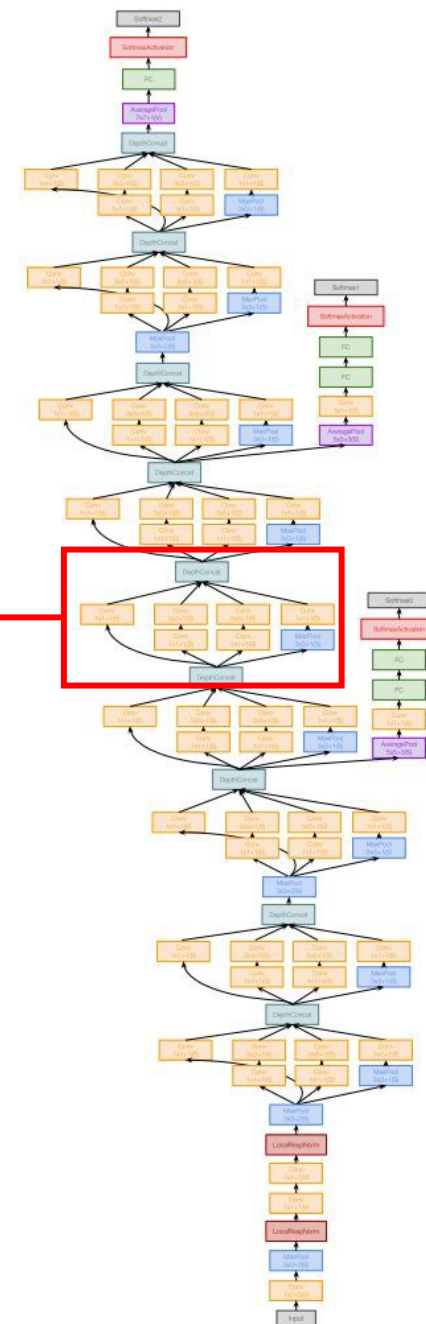
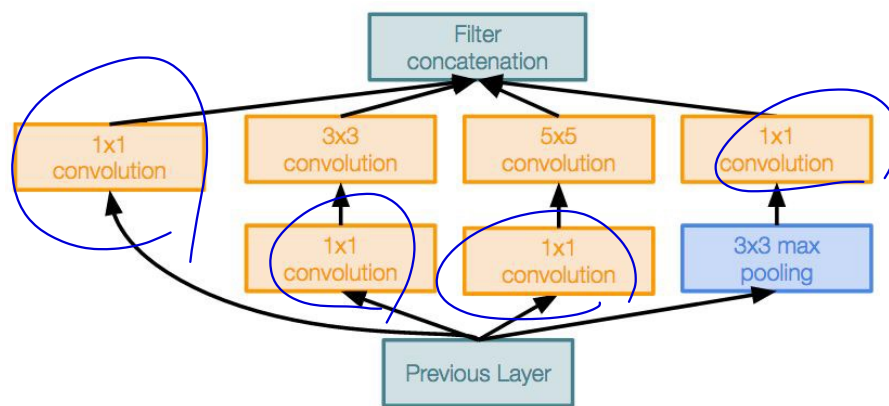
ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



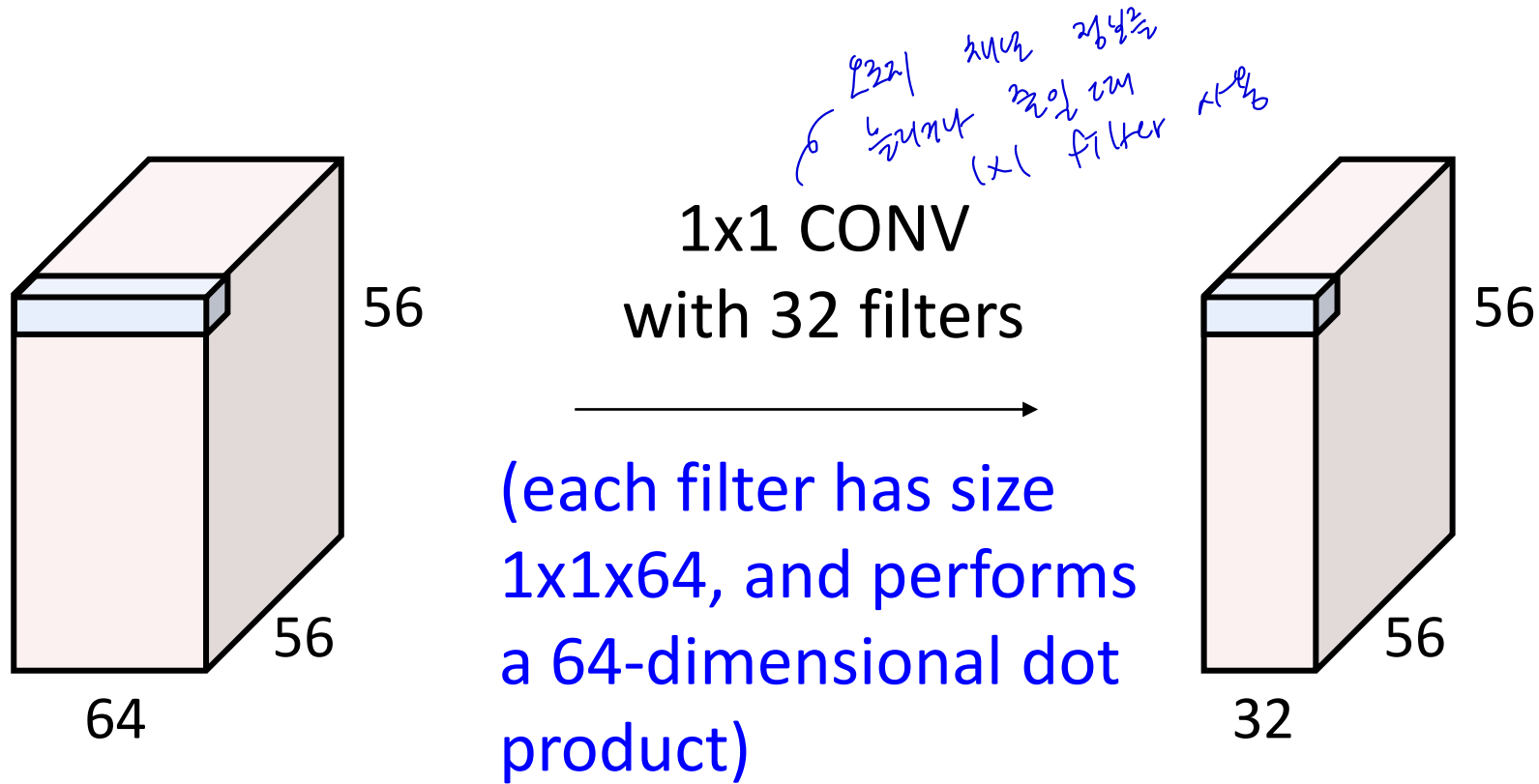
GoogLeNet

이것 미친 수준.

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

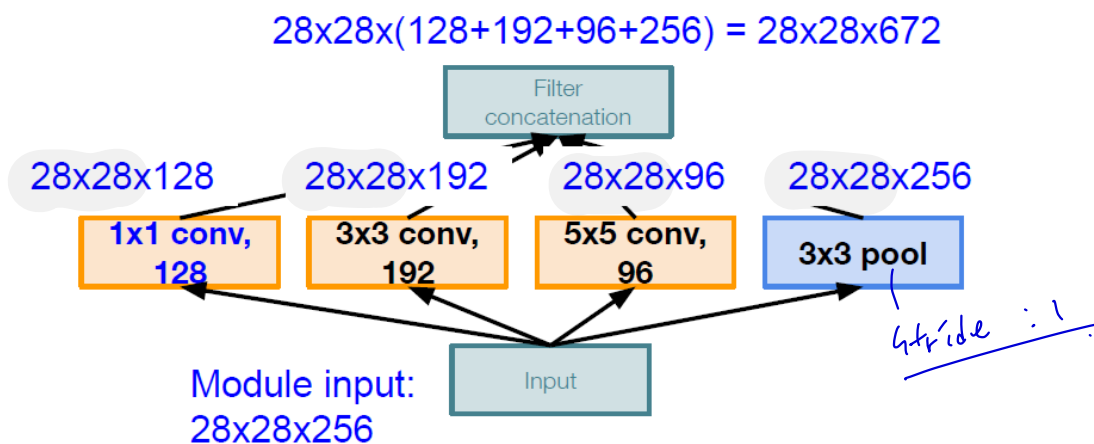


1x1 convolution layers make perfect sense

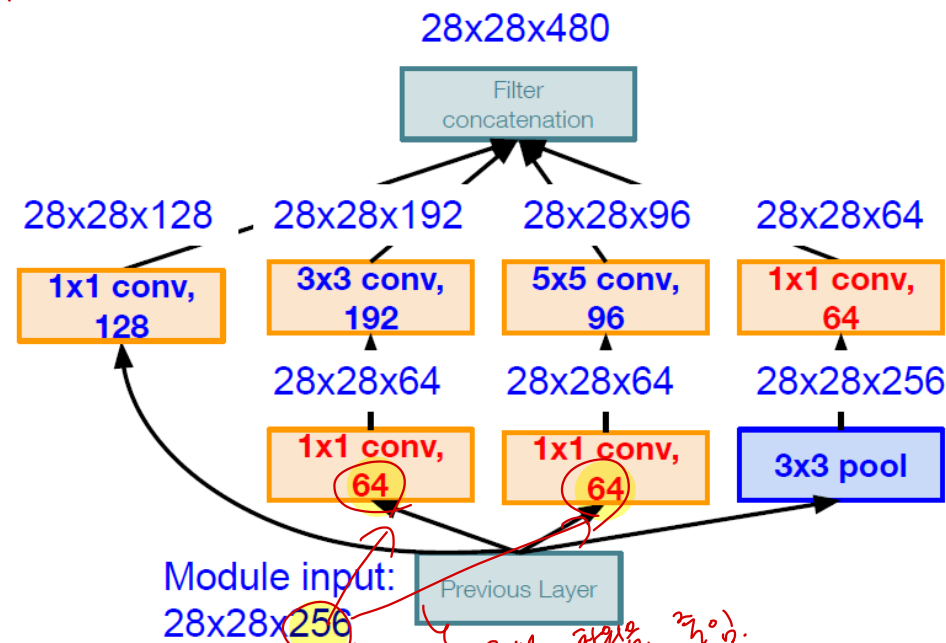


GoogLeNet

☆ 곁에 연산이
필요해서.

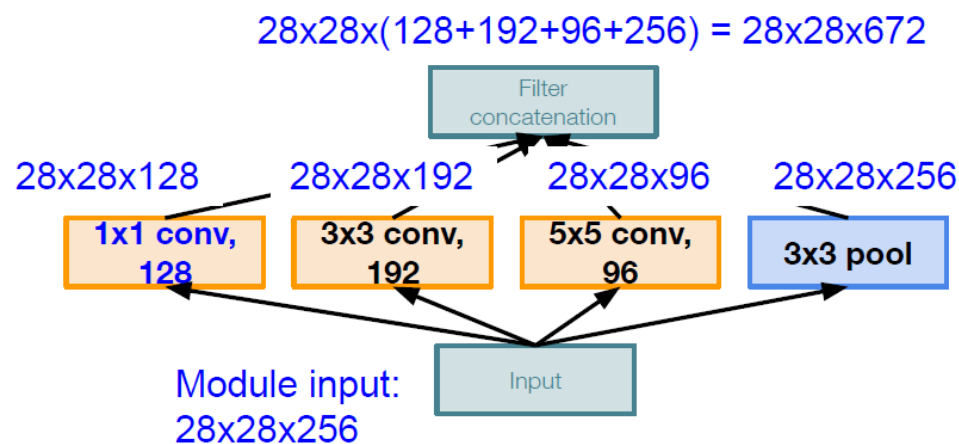


Conv Ops:
[1x1 conv, 128]
[3x3 conv, 192]
[5x5 conv, 96]
Total:



Conv Ops:
[1x1 conv, 64]
[1x1 conv, 64]
[1x1 conv, 128]
[3x3 conv, 192]
[5x5 conv, 96]
[1x1 conv, 64]
Total:

GoogLeNet



Conv Ops:

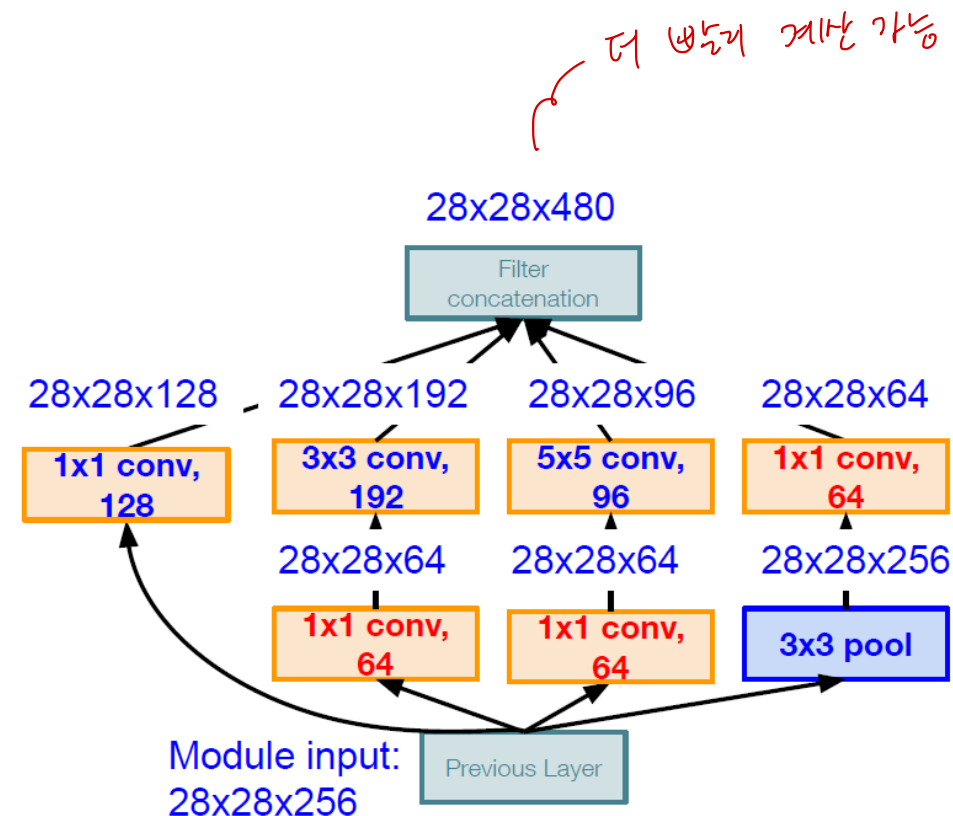
[1×1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$

[3×3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$

[5×5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$

Total: 854M ops

Very expensive compute



Conv Ops:

[1×1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$

[1×1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$

[1×1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$

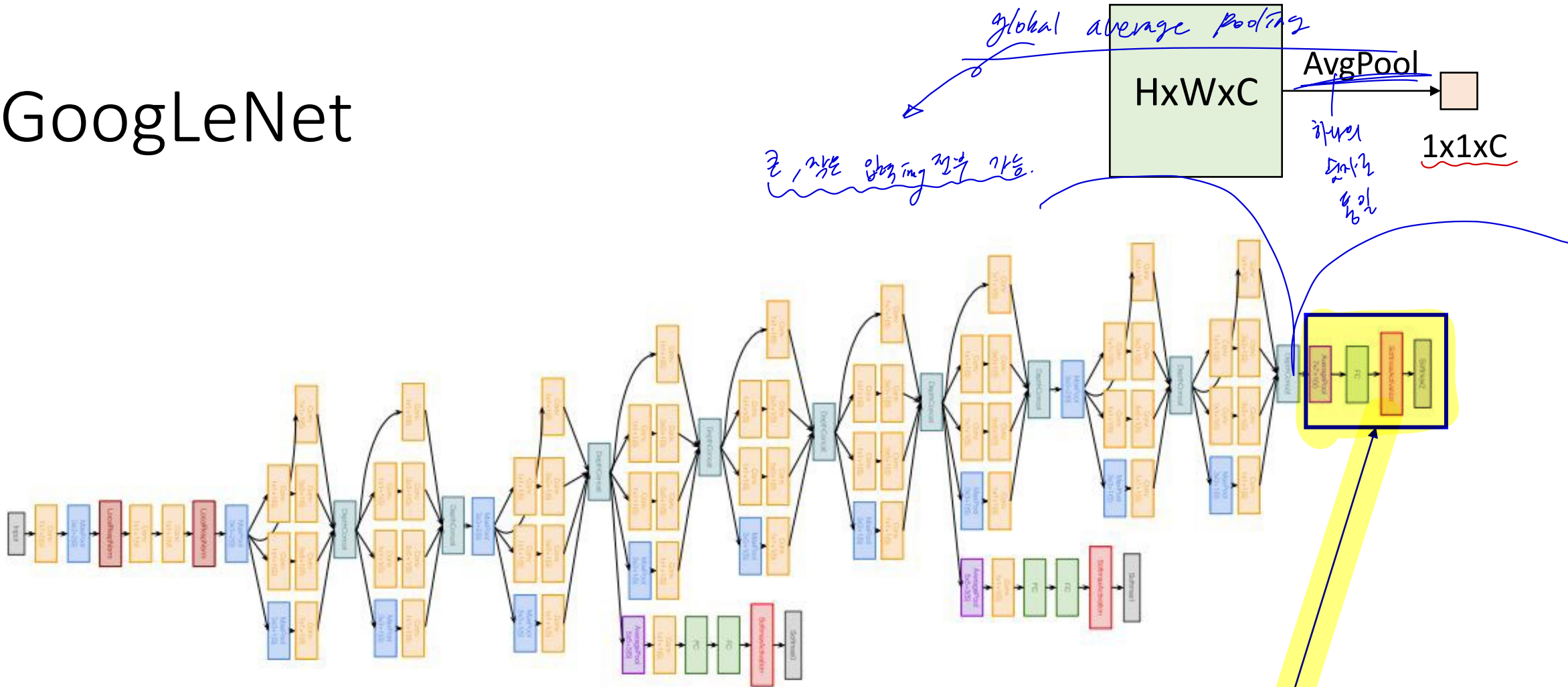
[3×3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 64$

[5×5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 64$

[1×1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$

Total: 358M ops

GoogLeNet

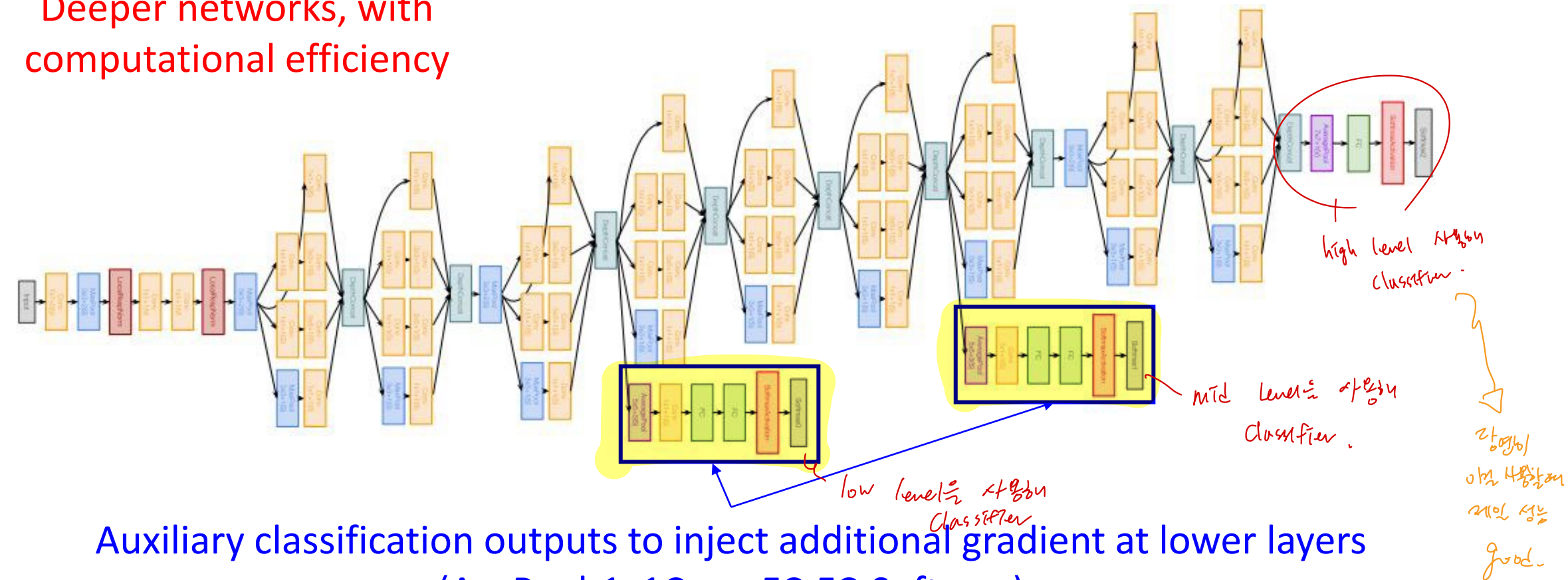


Note: after the last convolutional layer, a global average pooling layer is used that spatially averages across each feature map, before final FC layer. No longer multiple expensive FC layers!

Classifier output

GoogLeNet

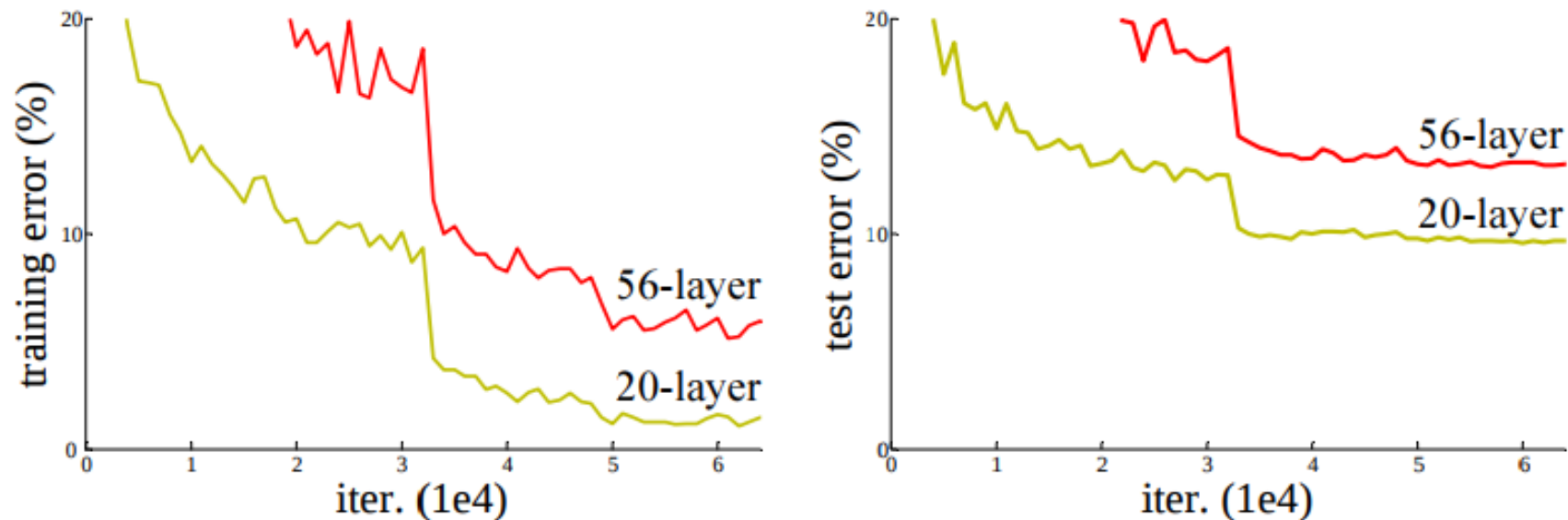
Deeper networks, with computational efficiency



앞단의 ∇ 는 gradient를 미분해서, gradient를 곱해서 나온 식이므로 문제 0.
그러나 앞단에서는 gradient를 따로 빼서 항하게 됨.

ResNet

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



56-layer model performs worse on both test and training error

-> The deeper model performs worse, but it's not caused by **overfitting**!

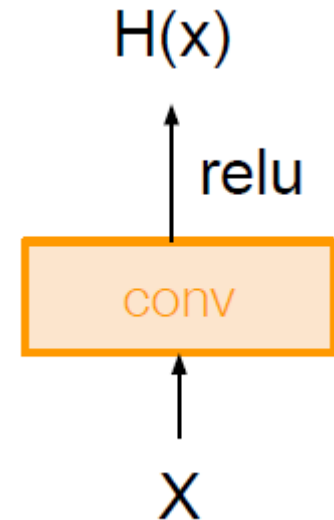
ResNet

Fact: Deep models have more representation power (more parameters) than shallower models.

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

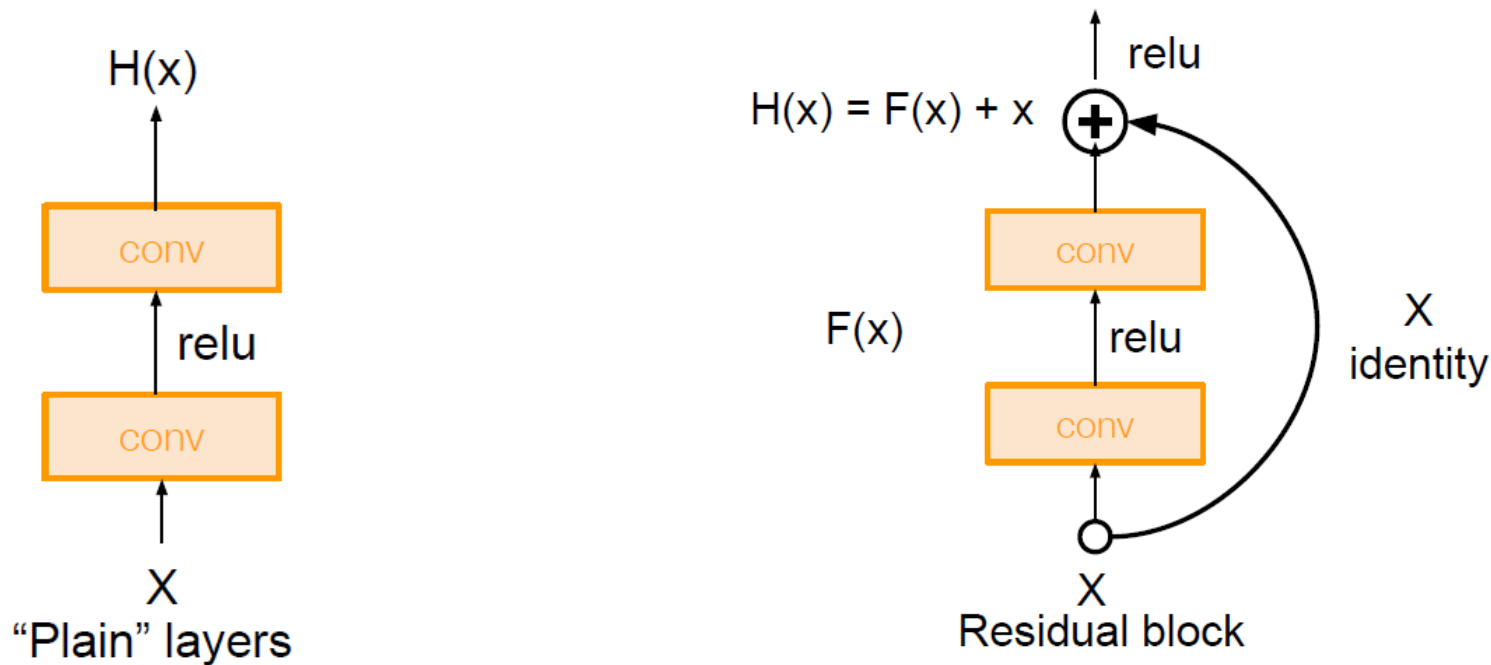
What should the deeper model learn to be at least as good as the shallower model?

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



ResNet

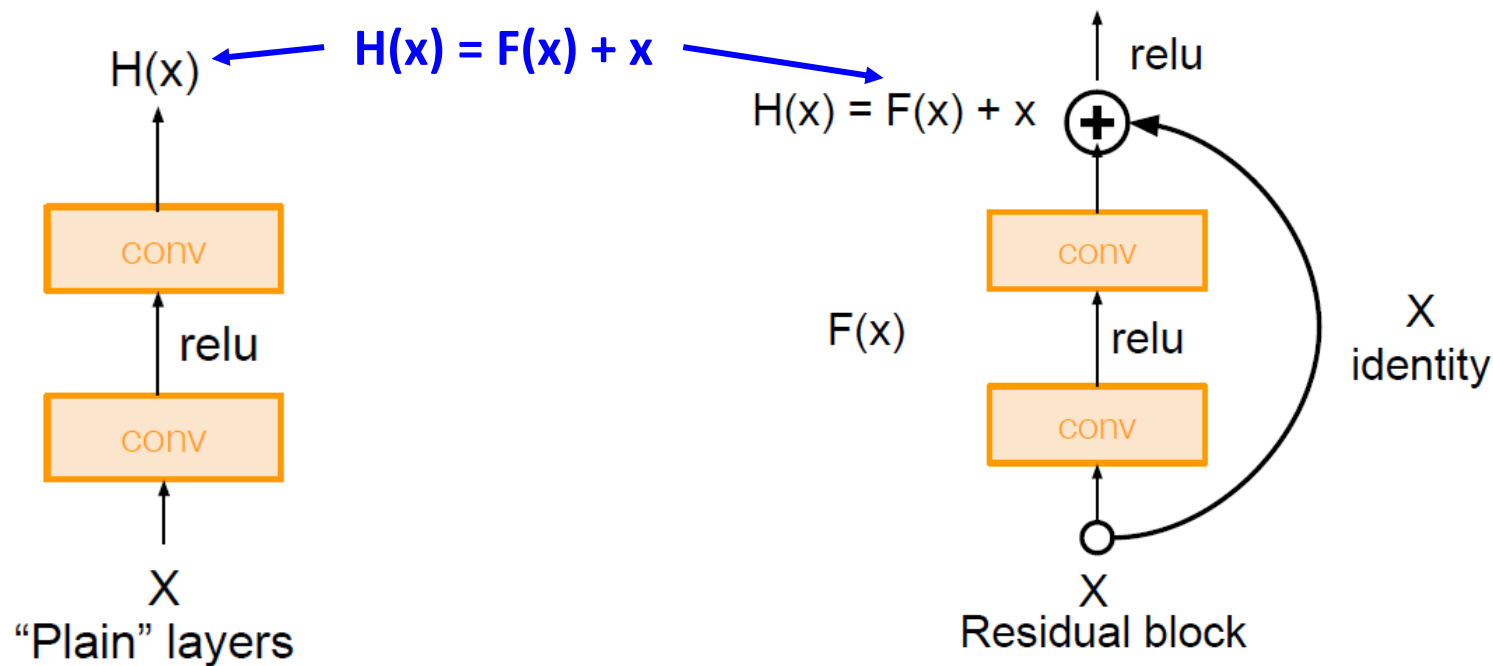
Solution: Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



Identity mapping:
 $H(x) = x$ if $F(x) = 0$

ResNet

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



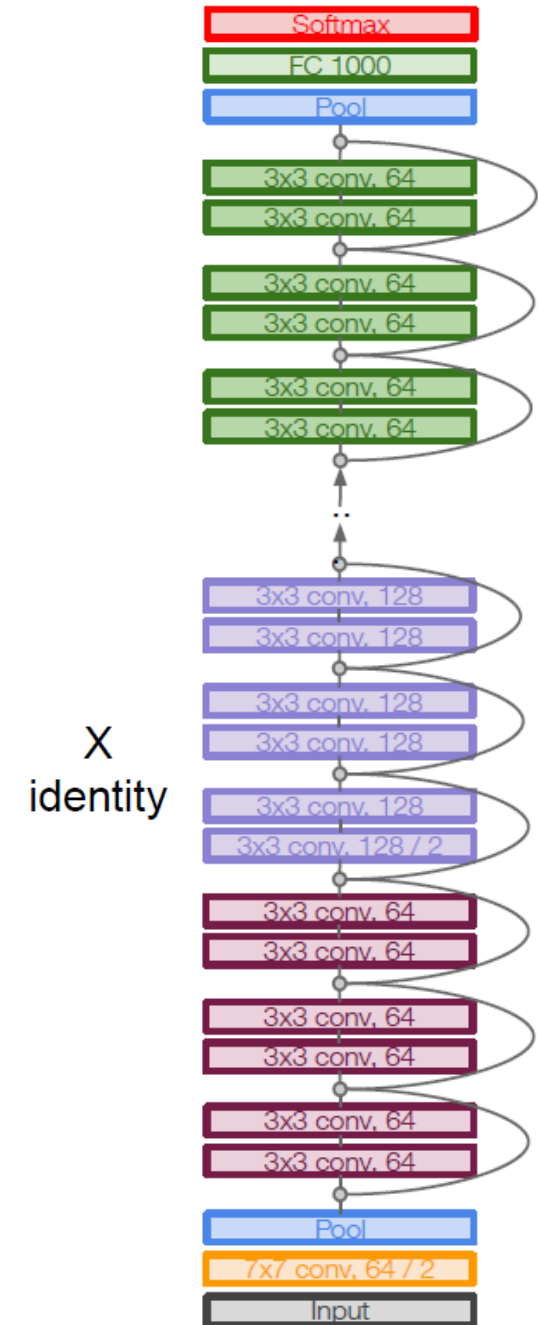
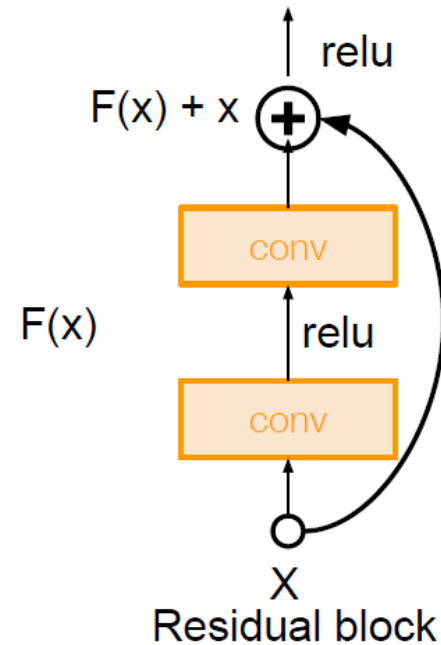
Identity mapping:
 $H(x) = x$ if $F(x) = 0$

Use layers to
fit **residual**
 $F(x) = H(x) - x$
instead of
 $H(x)$ directly

ResNet

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!

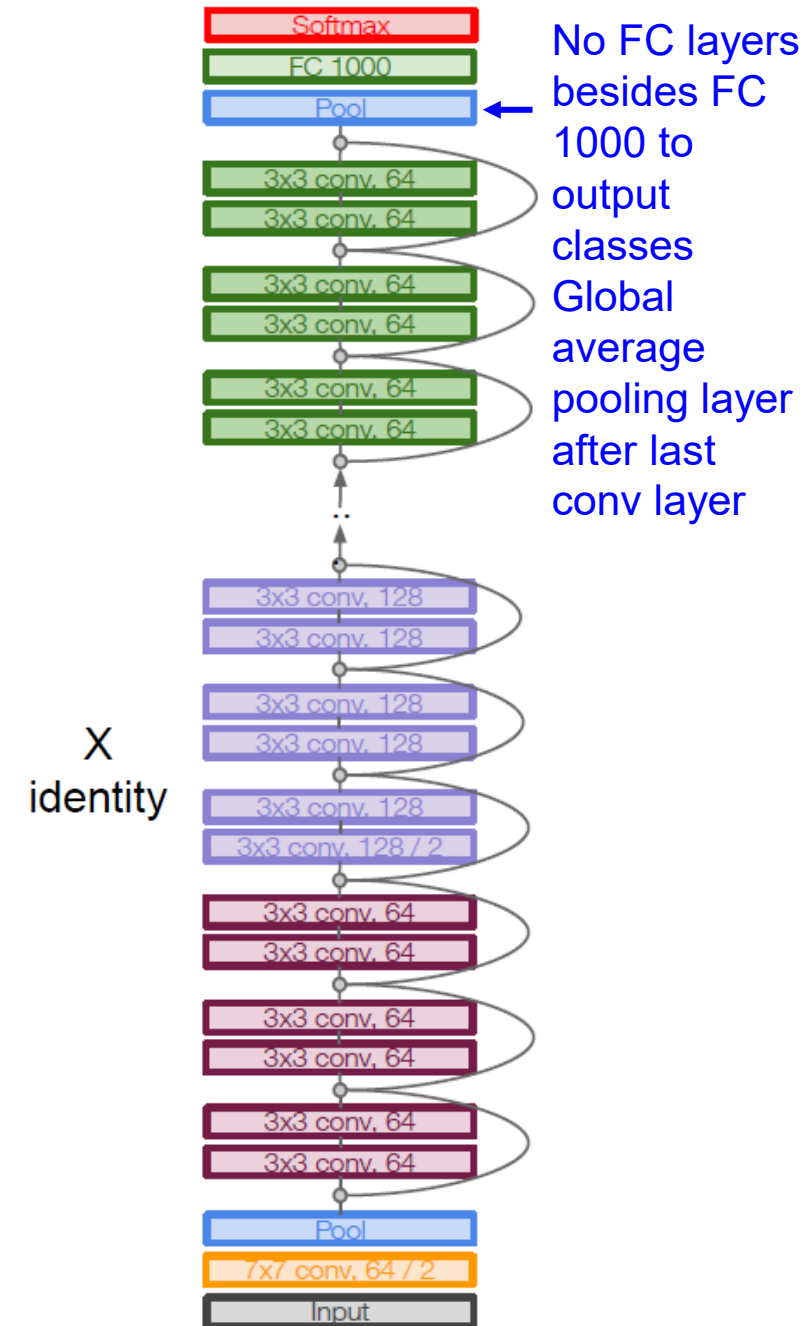
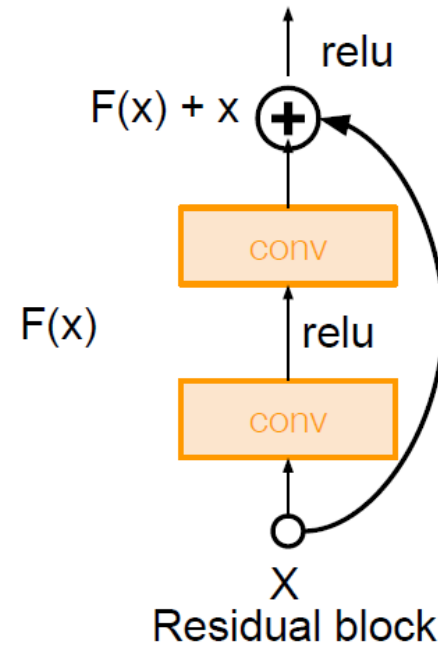


ResNet

Full ResNet architecture:

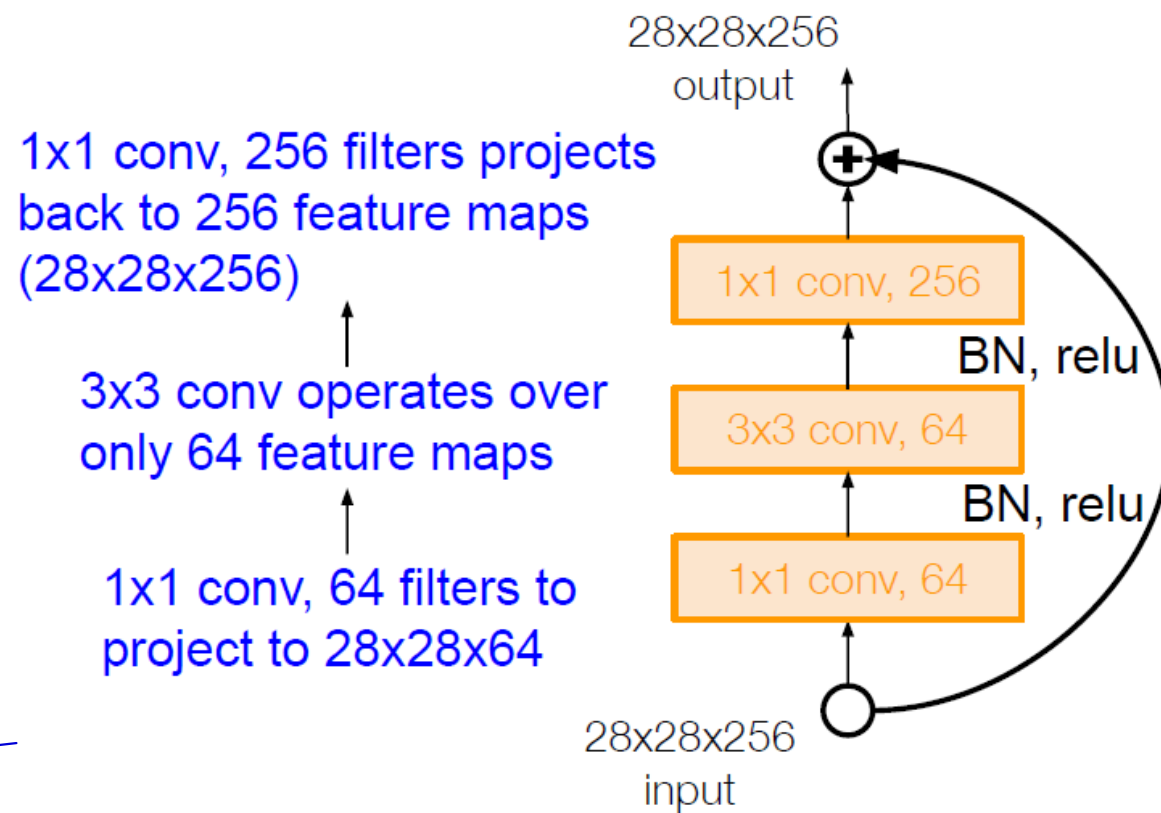
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning (stem)
- No FC layers at the end (only FC 1000 to output classes)
- (In theory, you can train a ResNet with input image of variable sizes)

Conv layer : 2×3



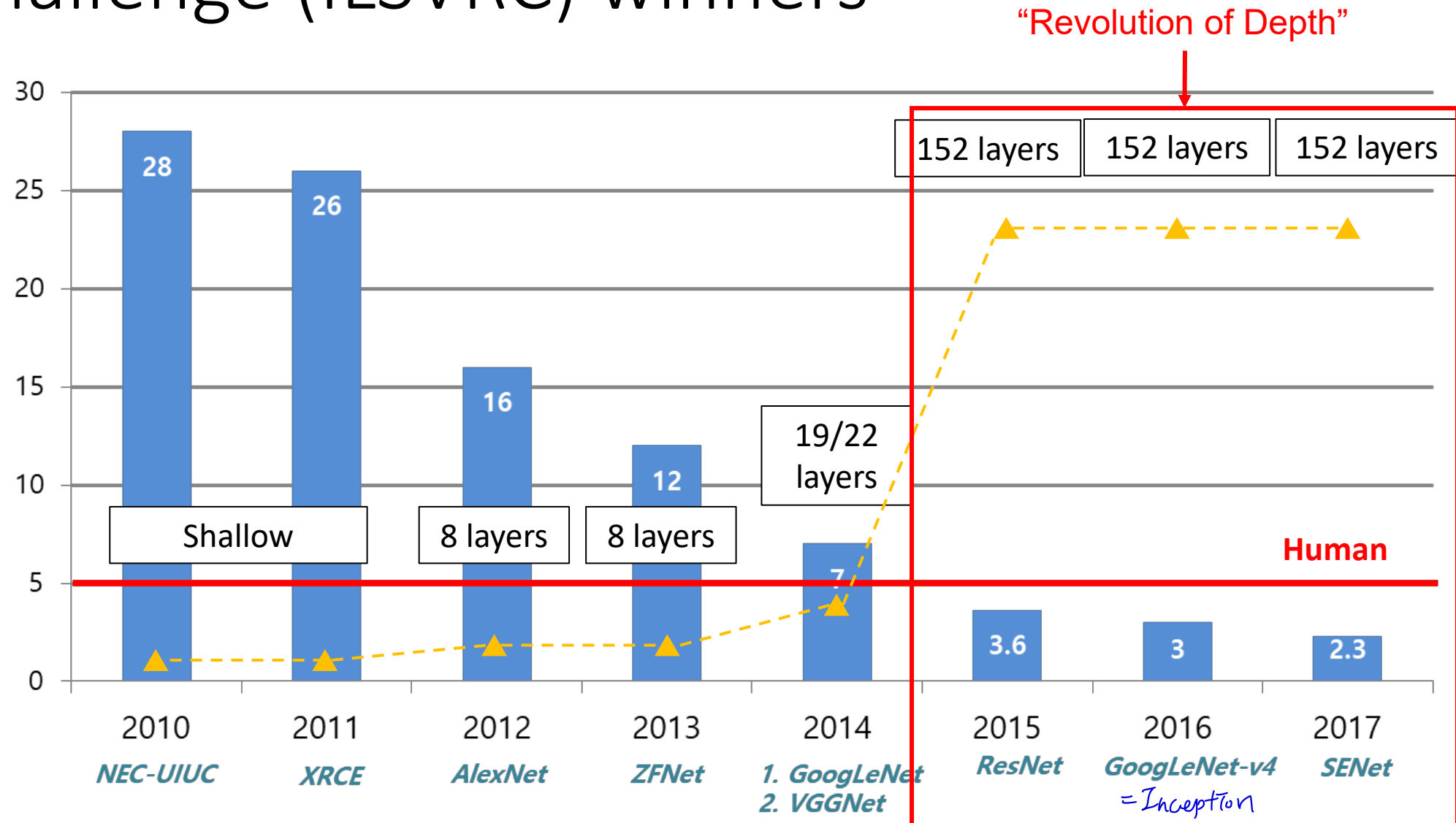
ResNet

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



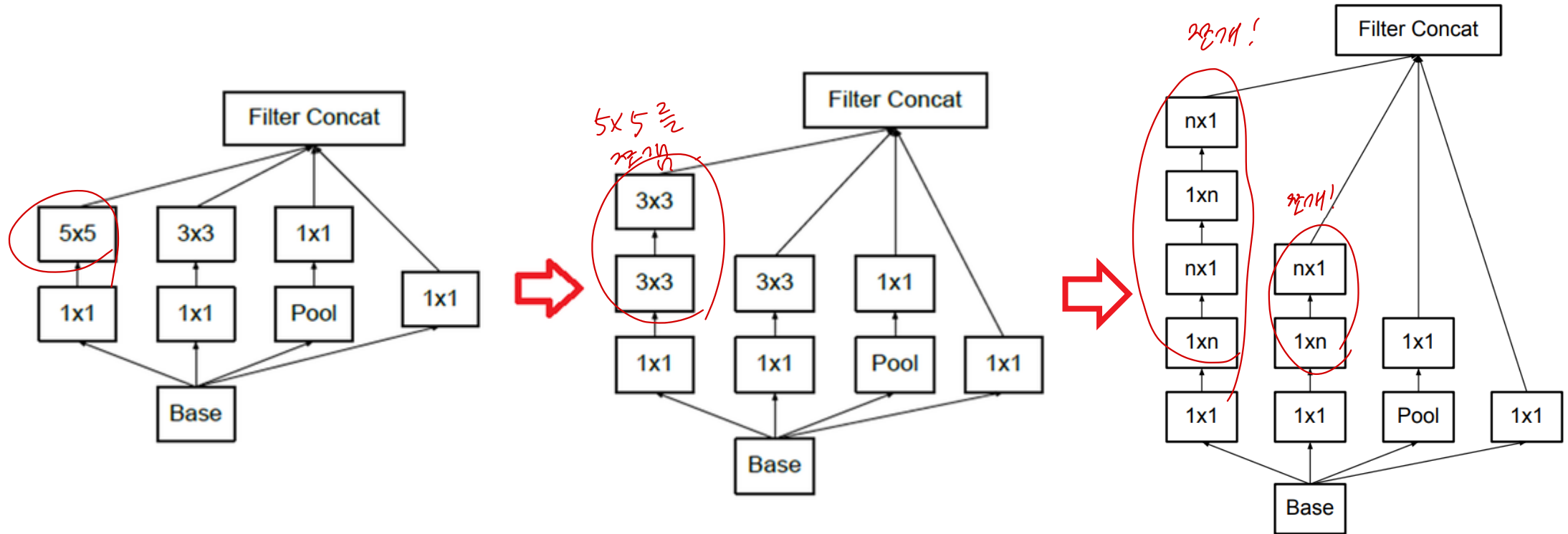
28x28x256
28x28x64 28x28x256?

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



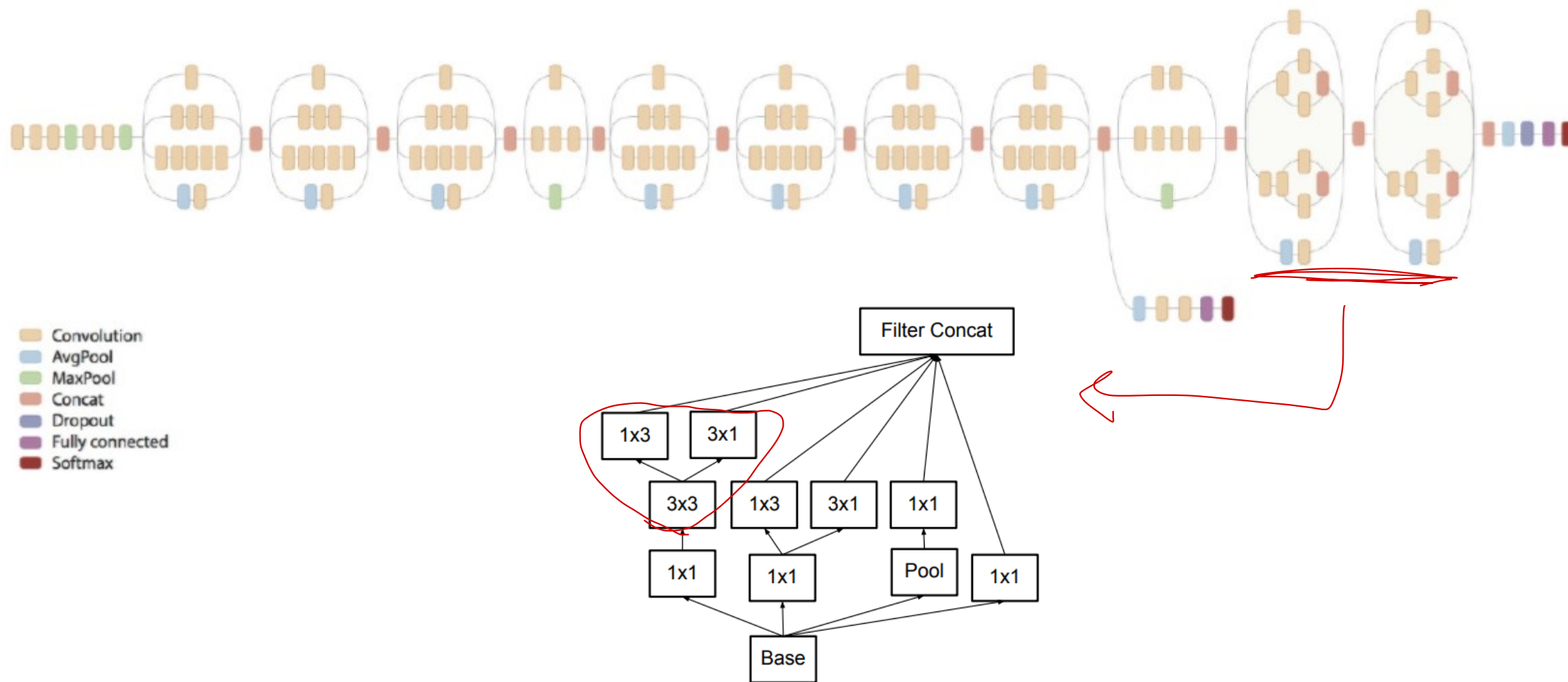
google net '합친 ver.

Inception v2



Reduce the model size and computational cost

Inception v2



Inception v3

흔히 학습 parameter 학습량? 만 큼임.

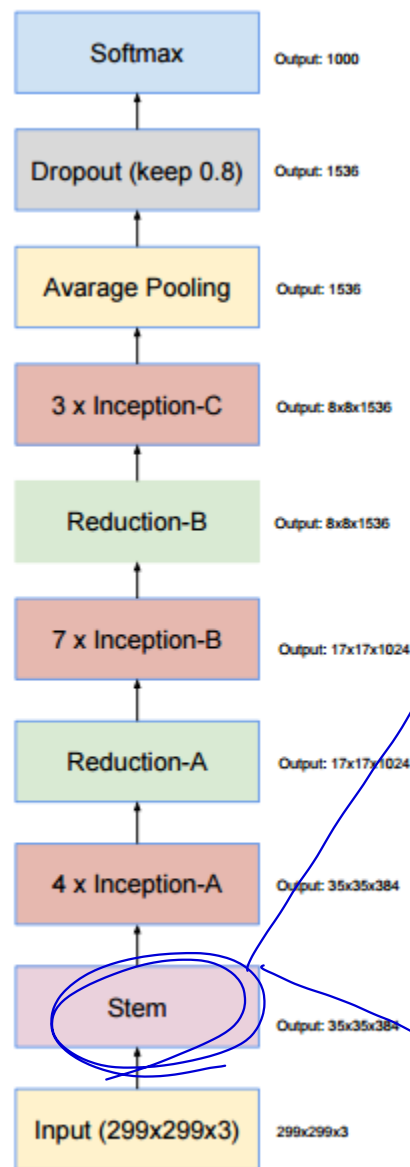
- Optimizer
 - RMSProp

- Label Smoothing → 허버 커리언 레이블 00001 이런 식이 아니라,
예측한 만큼? , 0.8%는 리라라, 0.2% 라라라 를 씀. 이런식으로

- Add BatchNorm after last Fully Connected layer

Inception v4

Inception 모듈 사용
 (A, B, C) 반복.
 특징을 만들어서
 사용.



이 부분을 zoom

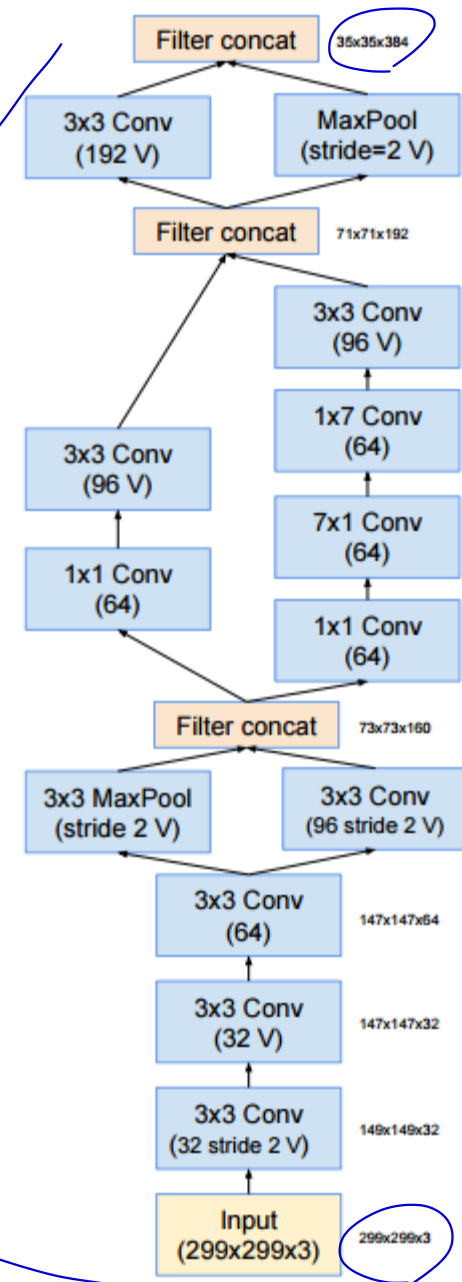


Figure 9. The overall schema of the Inception-v4 network. For the detailed modules, please refer to Figures 3, 4, 5, 6, 7 and 8 for the detailed structure of the various components.

Figure 3. The schema for stem of the pure Inception-v4 and Inception-ResNet-v2 networks. This is the input part of those networks. Cf. Figures 9 and 15

Inception v4

Inception-A

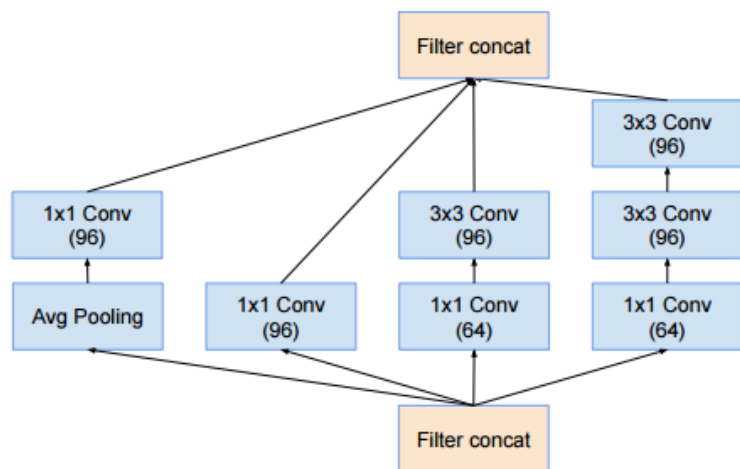


Figure 4. The schema for 35×35 grid modules of the pure Inception-v4 network. This is the Inception-A block of Figure 9.

low level

Inception-B

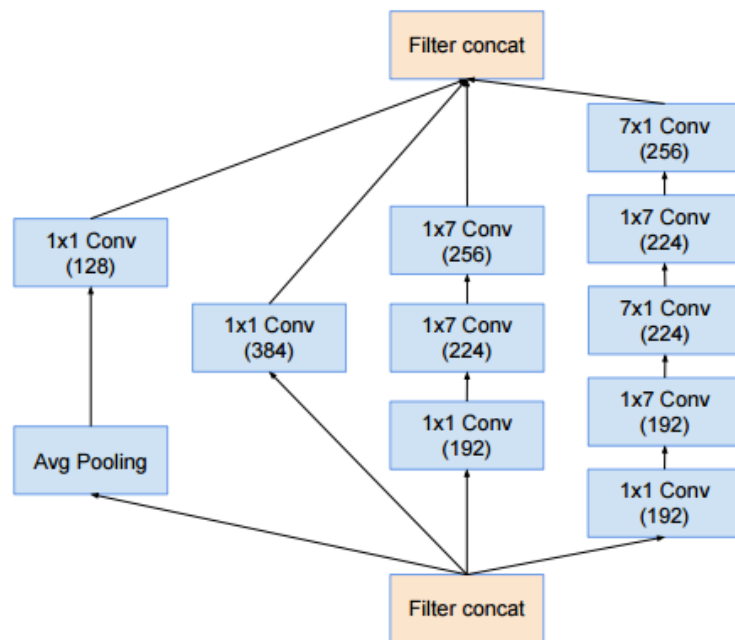


Figure 5. The schema for 17×17 grid modules of the pure Inception-v4 network. This is the Inception-B block of Figure 9.

mid level

Inception-C

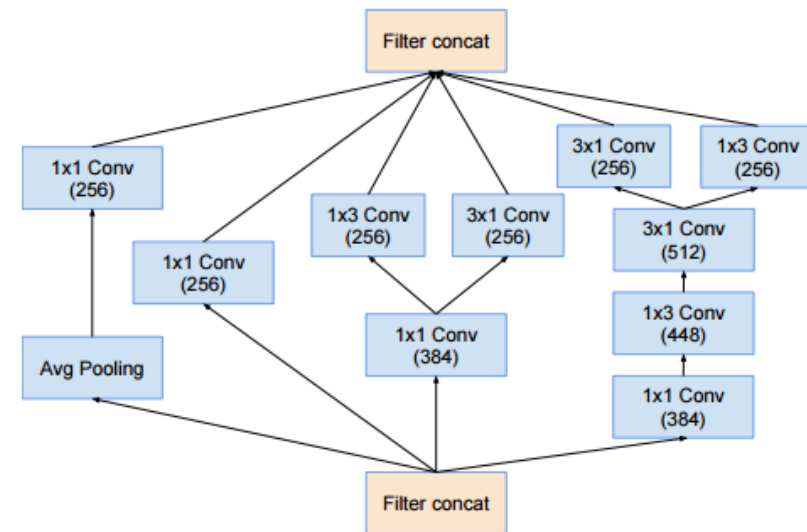


Figure 6. The schema for 8×8 grid modules of the pure Inception-v4 network. This is the Inception-C block of Figure 9.

high level

Inception v4

Pooling

Reduction-A

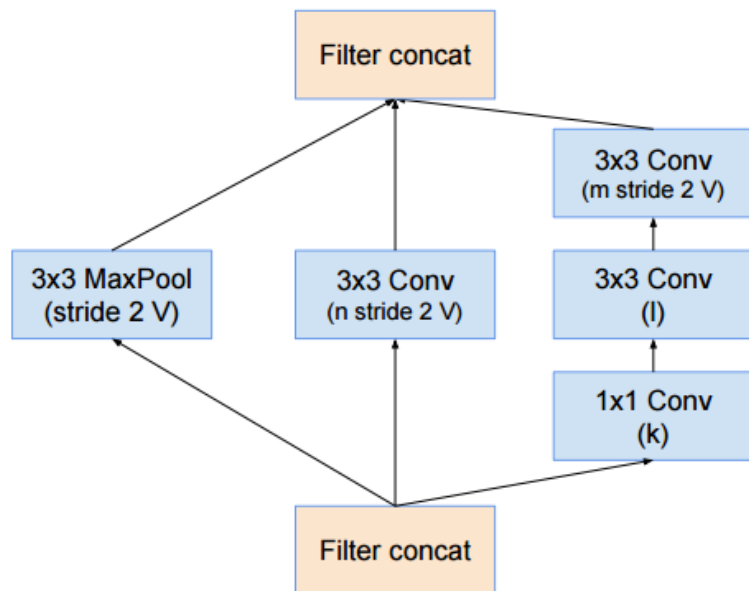


Figure 7. The schema for 35×35 to 17×17 reduction module. Different variants of this blocks (with various number of filters) are used in Figure 9, and 15 in each of the new Inception(-v4, -ResNet-v1, -ResNet-v2) variants presented in this paper. The k, l, m, n numbers represent filter bank sizes which can be looked up in Table 1.

Reduction-B

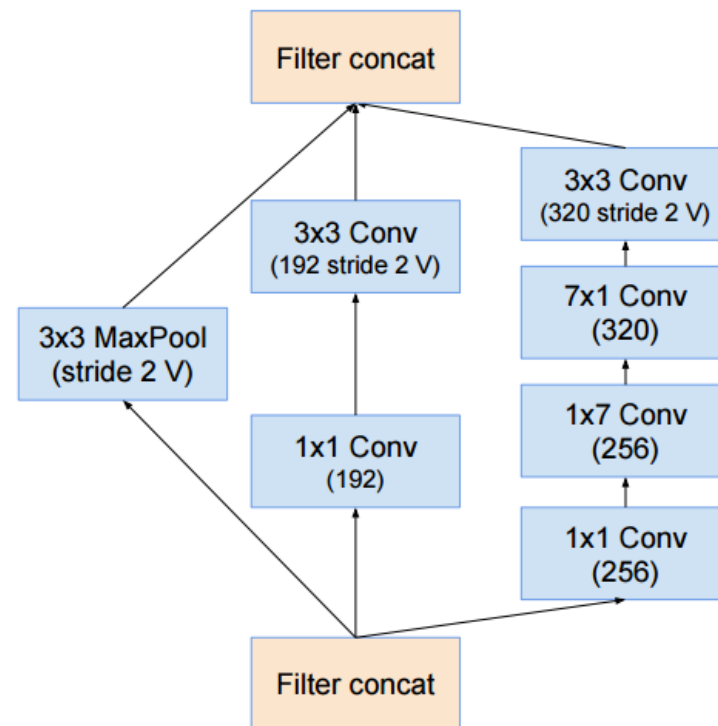
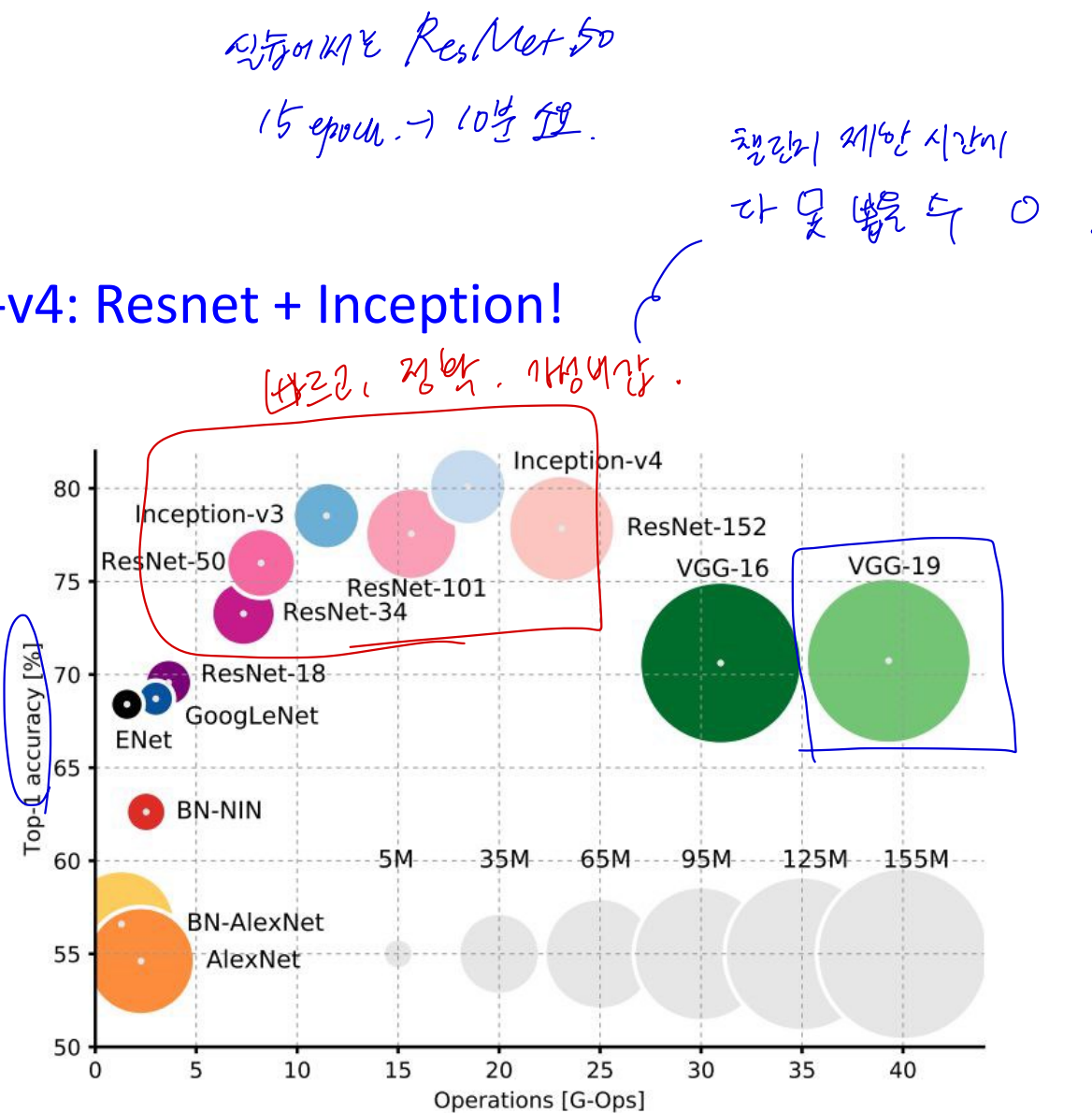
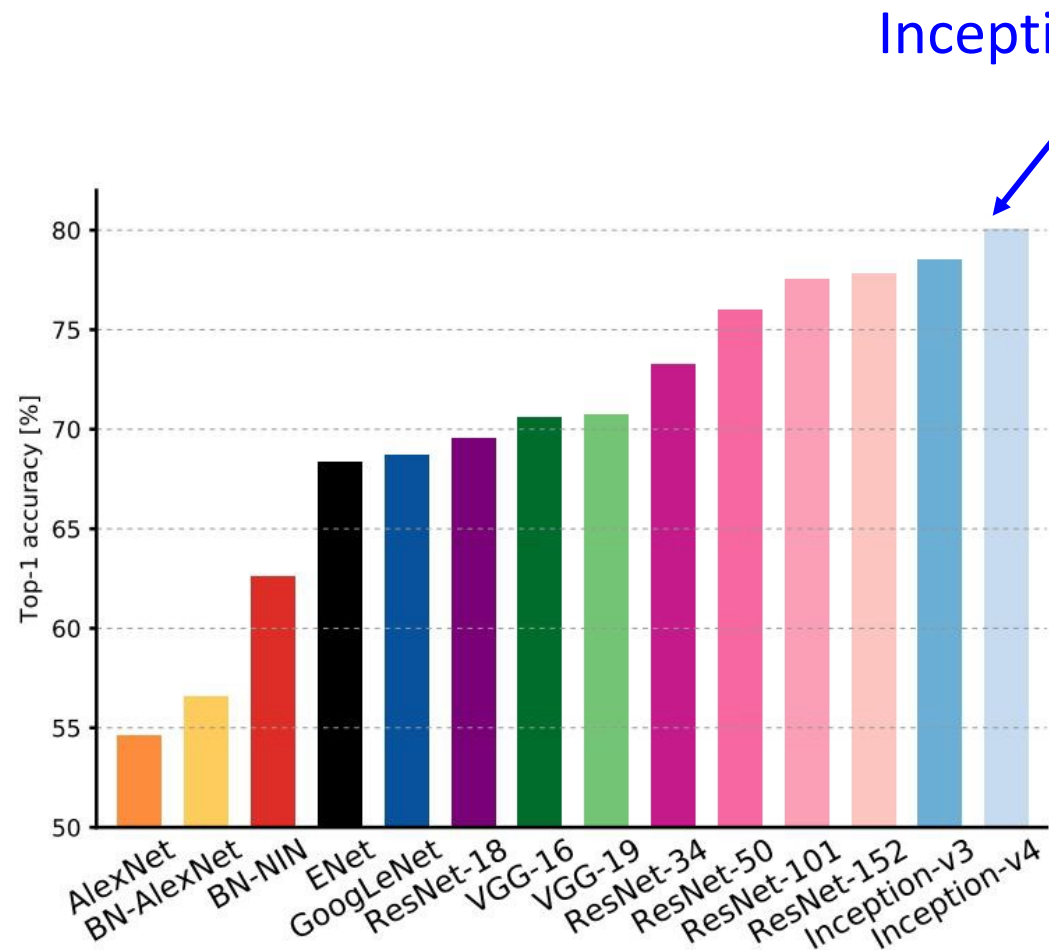
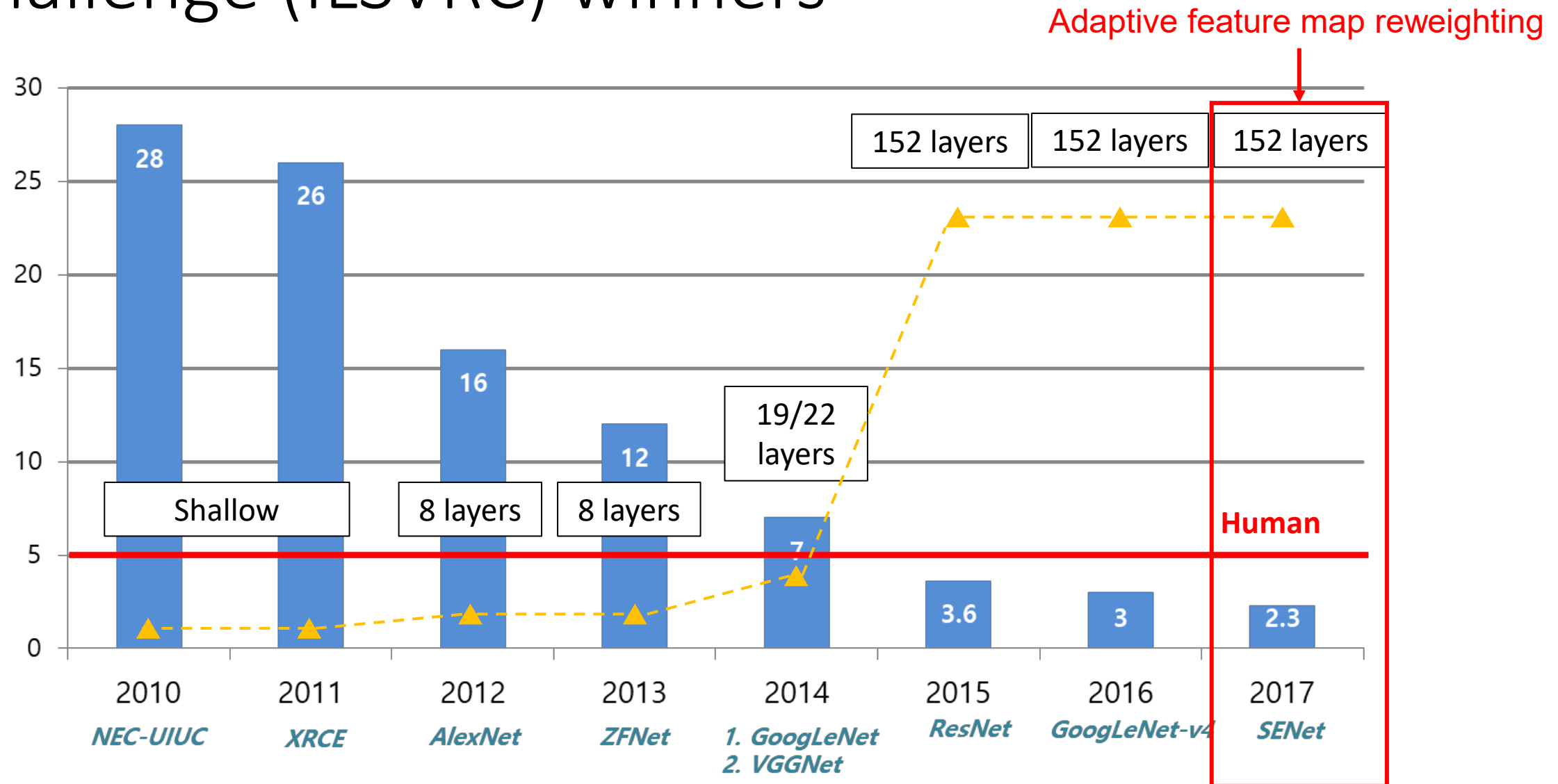


Figure 8. The schema for 17×17 to 8×8 grid-reduction module. This is the reduction module used by the pure Inception-v4 network in Figure 9.

Comparing complexity

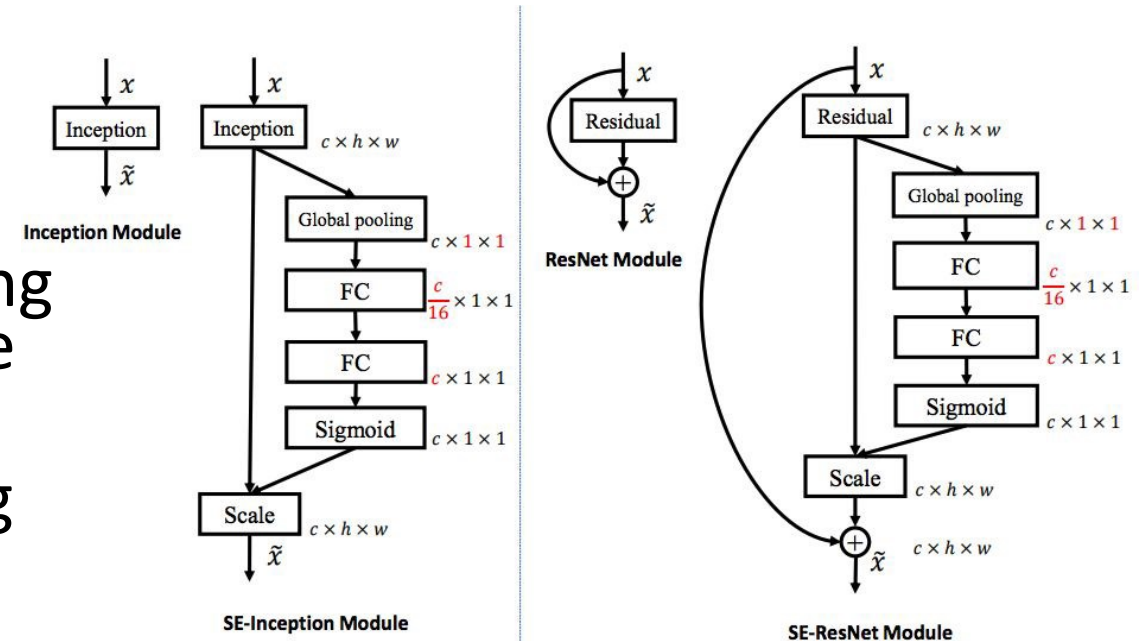


ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

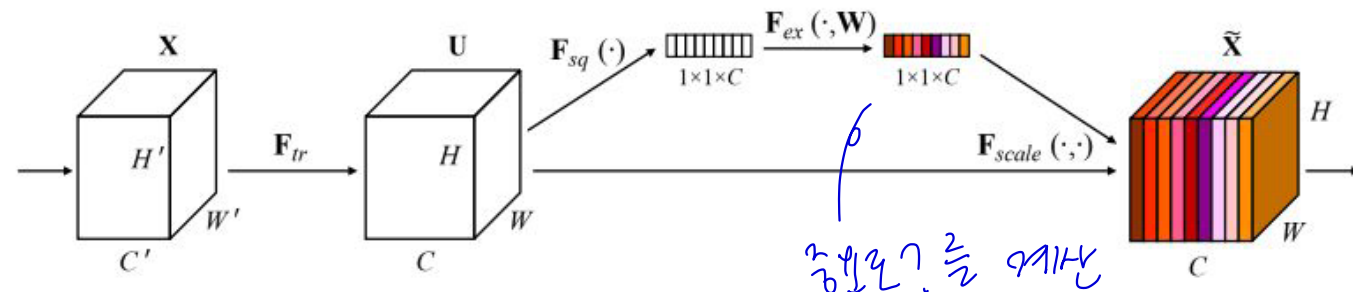


Squeeze-and-Excitation Networks (SENet)

- Add a “feature recalibration” module that learns to adaptively reweight feature maps
- Global information (global avg. pooling layer) + 2 FC layers used to determine feature map weights
- ILSVRC'17 classification winner (using ResNeXt-152 as a base architecture)



공간에 정보를 제한?



Efficient networks

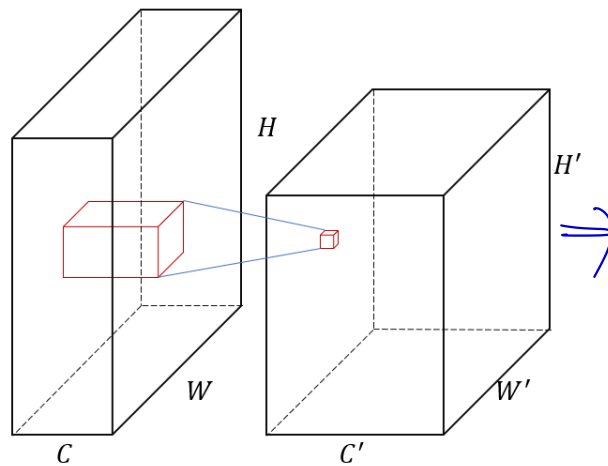
batch sized
증가.

차이점 다시 보기

• MobileNets

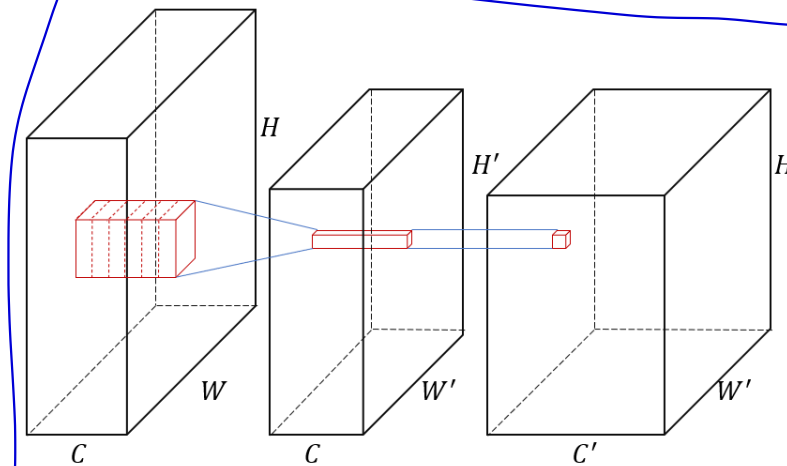
group conv를 사용.

- Depthwise separable convolutions replace standard convolutions by factorizing them into a depthwise convolution and a 1x1 convolution
- Much more efficient, with little loss in accuracy



$k \times k$ Convolution

#Parameters: $C \times k \times k \times C'$



$k \times k$ Depthwise Convolution \rightarrow 1x1 Pointwise Convolution

#Parameters: $C \times (k \times k + C')$

\rightarrow 1x1으로
서로의 group을 연결

\rightarrow C개의
2중 Conv 진행,

Main takeaways

AlexNet showed that you can use CNNs to train Computer Vision models.

ZFNet, **VGG** shows that bigger networks work better

GoogLeNet is one of the first to focus on efficiency using 1x1 bottleneck convolutions and global avg pool instead of FC layers

ResNet showed us how to train extremely deep networks

- Limited only by GPU & memory!
- Showed diminishing returns as networks got bigger

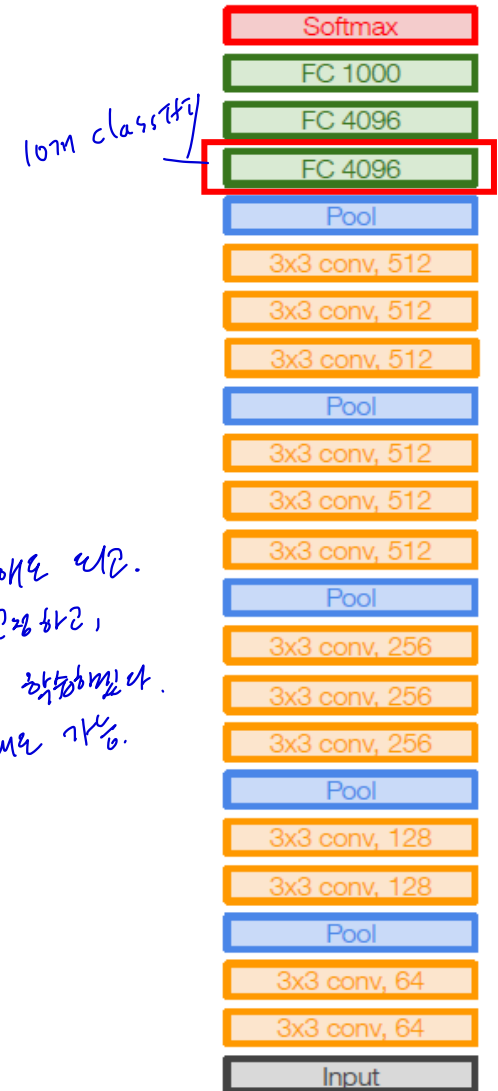
After ResNet: CNNs were better than the human metric and focus shifted to Efficient networks:

- Lots of tiny networks aimed at mobile devices: **MobileNet**

Transfer Learning with CNNs

Query(Test) image

L2 Nearest neighbors in feature space



다 학습해도 되고.
많은 경우 하고,
뒤만 학습하면 된다.
이런 것.

VGG16