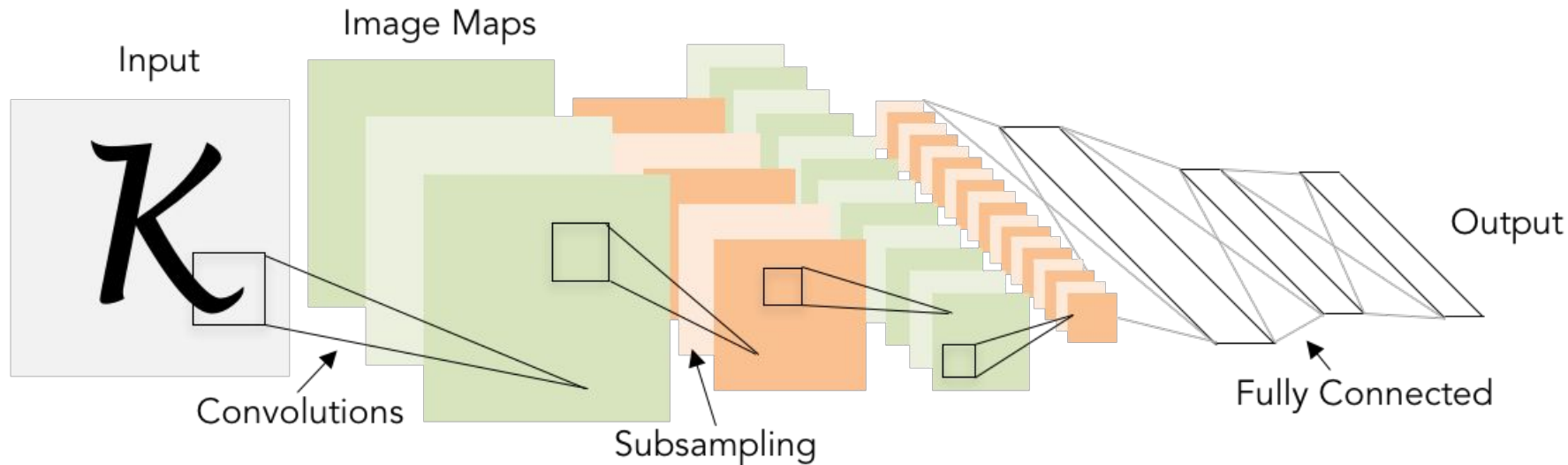# Lecture 9:
# CNN Architectures

# 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]

# Case Study: AlexNet

*[Krizhevsky et al. 2012]*



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

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

# Case Study: AlexNet

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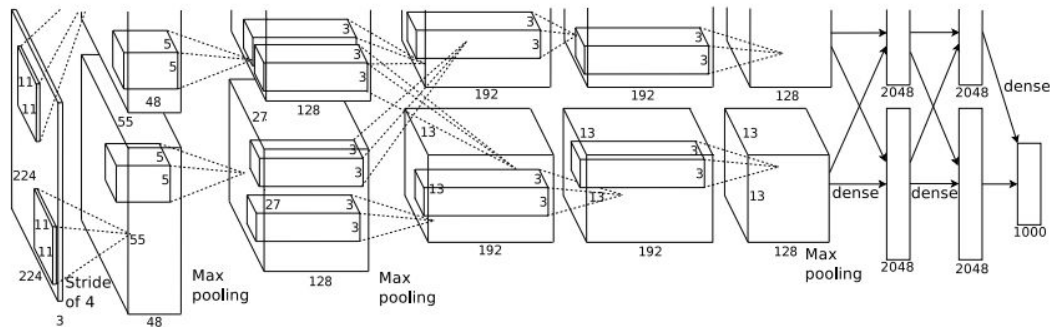
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ZFNet: Improved hyperparameters over AlexNet

152 layers

22 layers

19 layers

8 layers

8 layers

shallow

3.57
6.7
7.3
11.7
16.4
25.8
28.2

ILSVRC'15 ResNet | ILSVRC'14 GoogleNet | ILSVRC'14 VGG | ILSVRC'13 | ILSVRC'12 AlexNet | ILSVRC'11 | ILSVRC'10
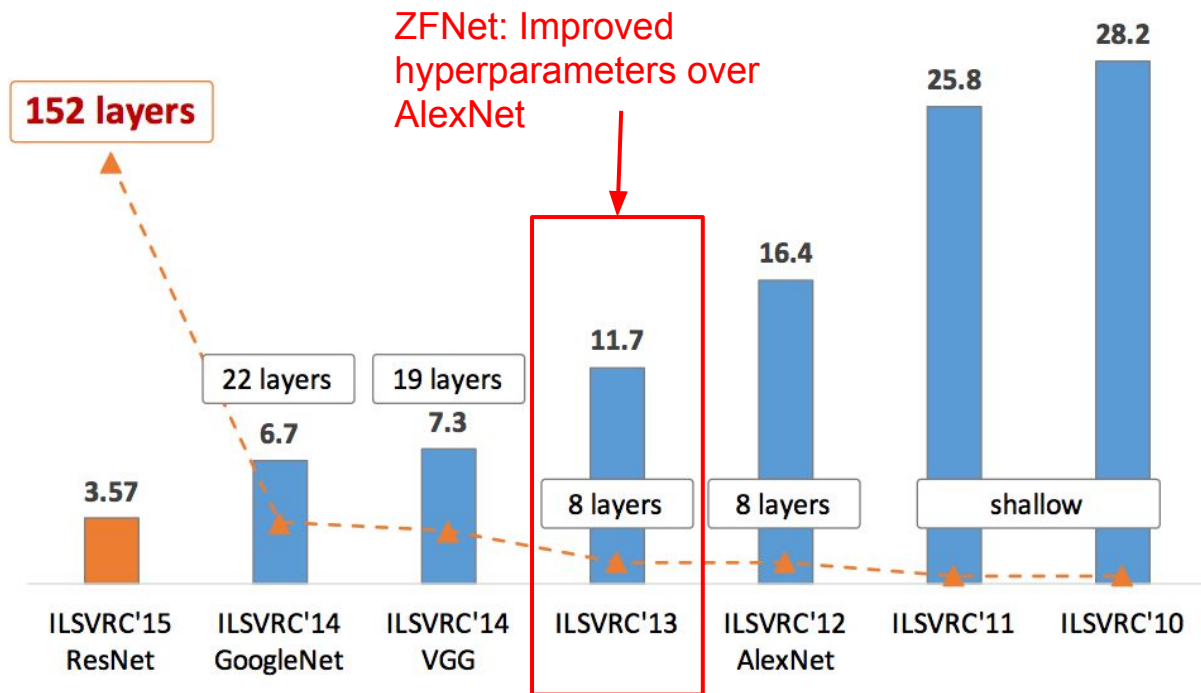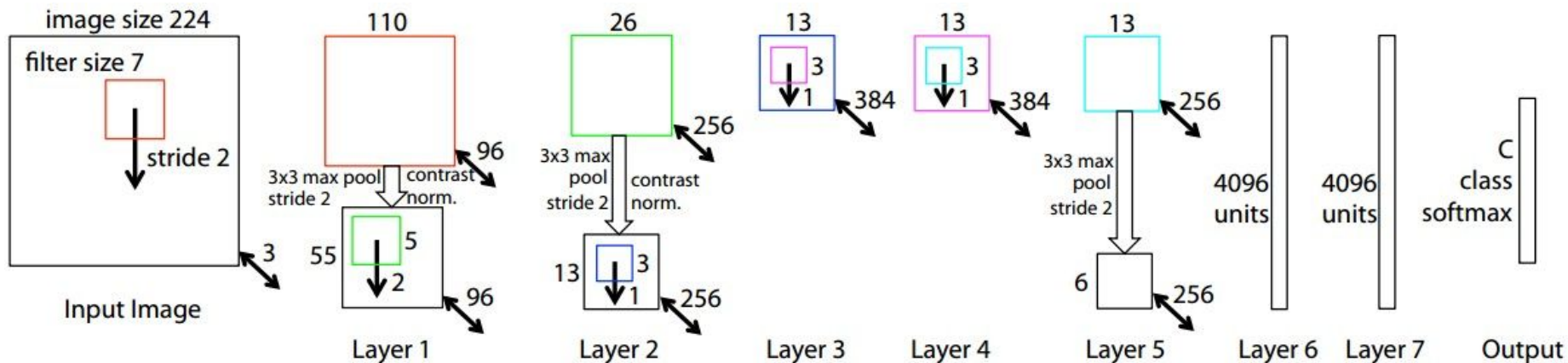
Figure copyright Kaiming He, 2016. Reproduced with permission.

# ZFNet

[Zeiler and Fergus, 2013]



TODO: remake figure

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%

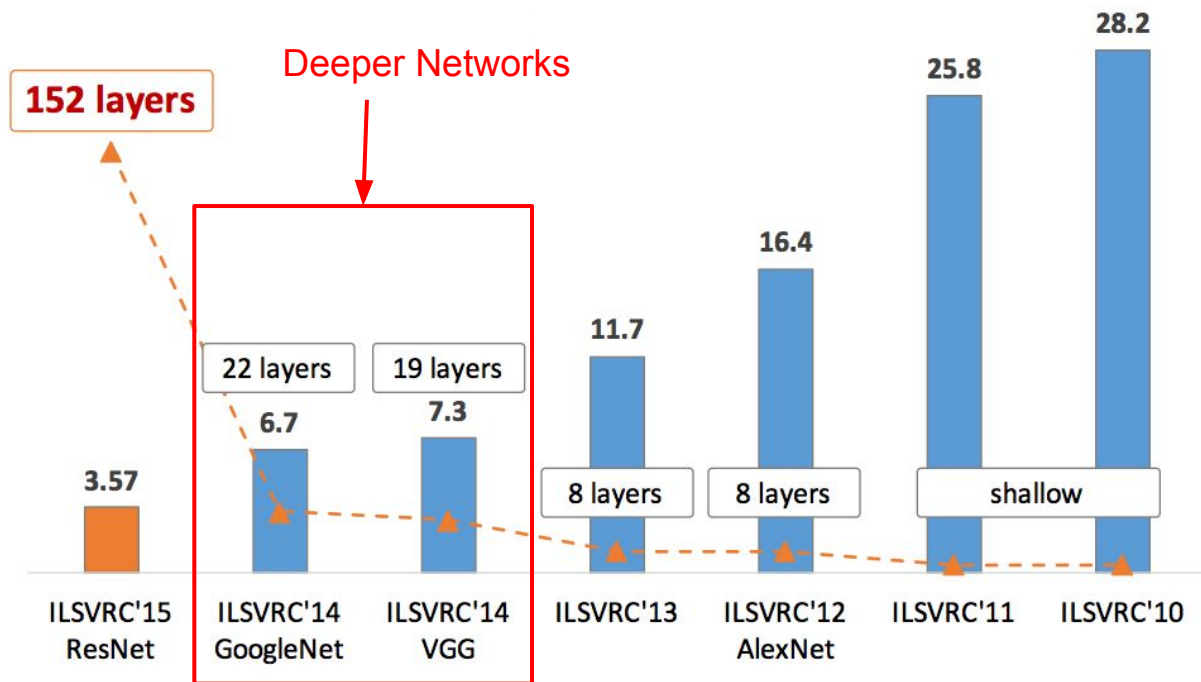# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

# Case Study: VGGNet

*[Simonyan and Zisserman, 2014]*
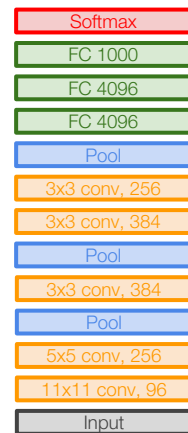
Small filters, Deeper networks

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

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

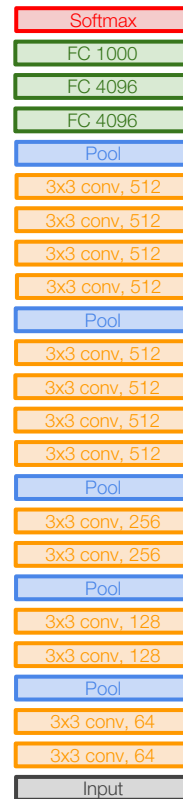11.7% top 5 error in ILSVRC'13
(ZFNet)
-> 7.3% top 5 error in ILSVRC'14



**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**

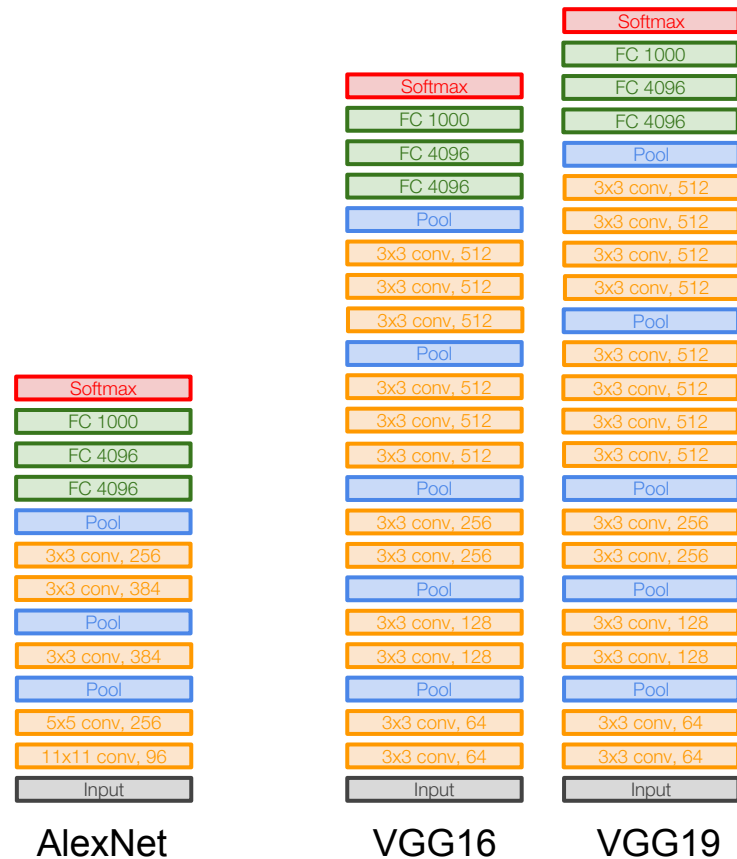| |
|---|
| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv,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

# 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

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

(not counting biases)

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

TOTAL params: 138M parameters

VGG16

Linear Neural Networks/

LeeSaeBom

# Dive into
# Deep Learning
## Linear Neural Networks

- *Network in Network*
- *Batch Normalixation*

Dive into Deep Learning

Interactive deep learning book with code, math, and discussions
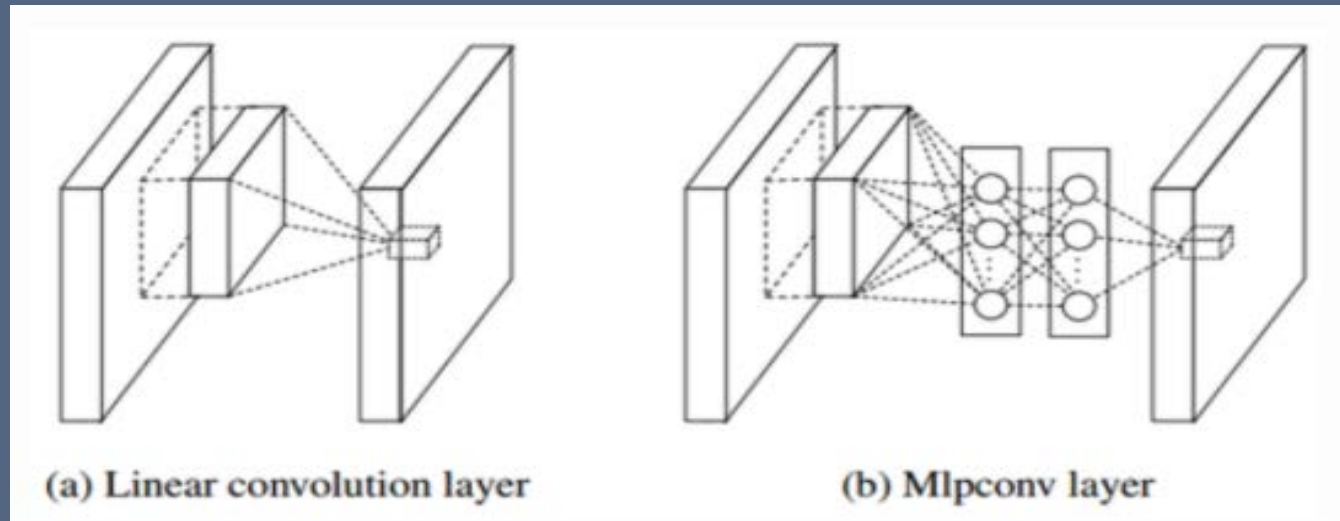
Implemented with **NumPy/MXNet**, **PyTorch**, and **TensorFlow**

Adopted at 140 universities from 35 countries

# Dive into
# Deep Learning

Network in Network

- **Network in Network**

- *CNN :* Filter을 이용하여 Stride만큼 이동하면서 CONV으로 Feature 추출함

- *NiN :* Filter대신에 MLP를 사용 = Mlpconv layer

*Network in Network*

*Batch Normalization*



(a) Linear convolution layer   (b) Mlpconv layer

- • *Network in Network*

- - *Why use the NiN?* 1x1 Conv를 통해 Feature map 개수를 줄일 수 있다
    = Parameter 수를 줄일 수 있다

- - NiN은 Mlpconv layer를 여러 개 쌓아 사용했으므로 네트워크 안에 네트워크가 있다는 개념을 NIN으로 불린다.

*Network in Network*

*Batch Normalization*



1x1 Convolution

# Dive into Deep Learning

Batch Normalization

- *Batch :* 신경망을 학습시킬 때, 한 번에 학습시키지 않고 조그만 단위로 분할해서 학습을 시키는데 이 때의 조그만 단위가 배치

- *Batch Normalization :* 배치별로 구분하고 각각의 출력값들의 정규화



(좌) Normalization 적용 전 / (우) Normalization 적용 후

- *Internal Covariance Shift*

- *Covariate Shift :* 이전 레이어의 파라미터 변화로 인하여 현재 레이어의 입력의 분포가 바뀌는 현상

- *Internal Covariance Shift* : 레이어를 통과할 때 마다 Covariate Shift가 일어나면서 입력의 분포가 약간씩 변하는 현상

# Dive into
# Deep Learning

Batch Normalization

*Batch Normalization*

- *Batch Normalization Algorithm*



**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.

- mini-batch의 평균과 분산을 구하고 입력 데이터를 평균이 0, 분산이 1로 되게 정규화를 진행함.
- scale(확대) and shift(이동)를 거쳐 학습 가능한 변수를 $\gamma, \beta$ 통해 실행
- $\gamma, \beta$ 는 역전파에 의해 학습된 변수

# Dive into
# Deep Learning

Batch Normalization

- 배치정규화 계산 그래프



- 신경망에서의 배치 정규화



without Batch Normalization

with Batch Normalization

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

# 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

# Case Study: GoogLeNet

*[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

# Case Study: GoogLeNet

*[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]

# Case Study: GoogLeNet

*[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]

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

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Example:

Q3:What is output size after filter concatenation?

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

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

Filter concatenation

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

28x28x128     28x28x192     28x28x96     28x28x256

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

Module input: 28x28x256

Input

Naive Inception module

# Reminder: 1x1 convolutions



56

1x1 CONV
with 32 filters

56

56

preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

64

56

32

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

1x1 conv "bottleneck"
layers



Naive Inception module

Inception module with dimension reduction

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*



**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

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

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

**Stack Inception modules with dimension reduction on top of each other**



Inception module

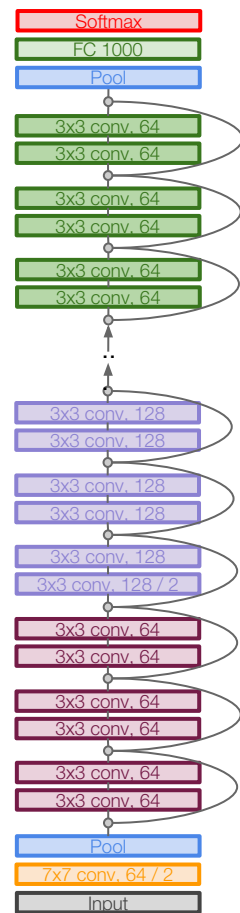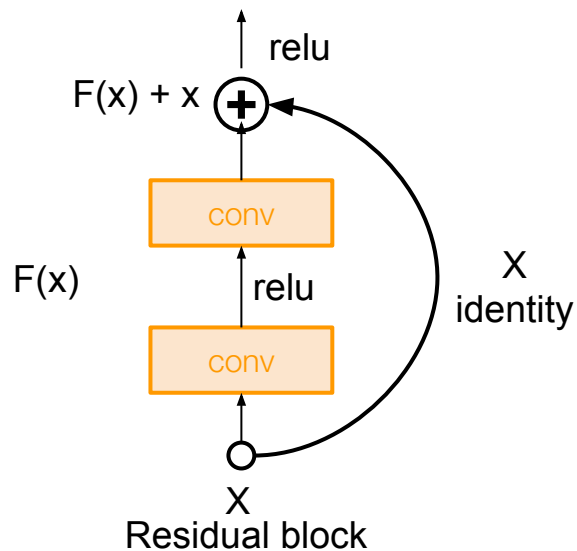# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



"Revolution of Depth"

152 layers

28.2

25.8

16.4

11.7

22 layers

19 layers

6.7

7.3

8 layers

8 layers

shallow

3.57

ILSVRC'15
ResNet

ILSVRC'14
GoogleNet

ILSVRC'14
VGG

ILSVRC'13

ILSVRC'12
AlexNet

ILSVRC'11

ILSVRC'10

Figure copyright Kaiming He, 2016. Reproduced with permission.

# 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!



F(x) + x

relu

conv

F(x)                relu        X
                                identity
conv

X
Residual block



Softmax
FC 1000
Pool
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
...
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

# Case Study: ResNet

*[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]

# Case Study: ResNet

*[He et al., 2015]*

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

# Case Study: ResNet

*[He et al., 2015]*

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



$H(x) = F(x) + x$

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

"Plain" layers

Residual block

# DENSELY CONNECTED CONVOLUTIONAL NETWORKS

Gao Huang*, Zhuang Liu*, Laurens van der Maaten, Kilian Q. Weinberger

Cornell University          Tsinghua University

Facebook AI Research

CVPR 2017

# CONVOLUTIONAL NETWORKS



LeNet

AlexNet

VGG

Inception

ResNet

# STANDARD CONNECTIVITY

# RESNET CONNECTIVITY

Identity mappings promote gradient propagation.
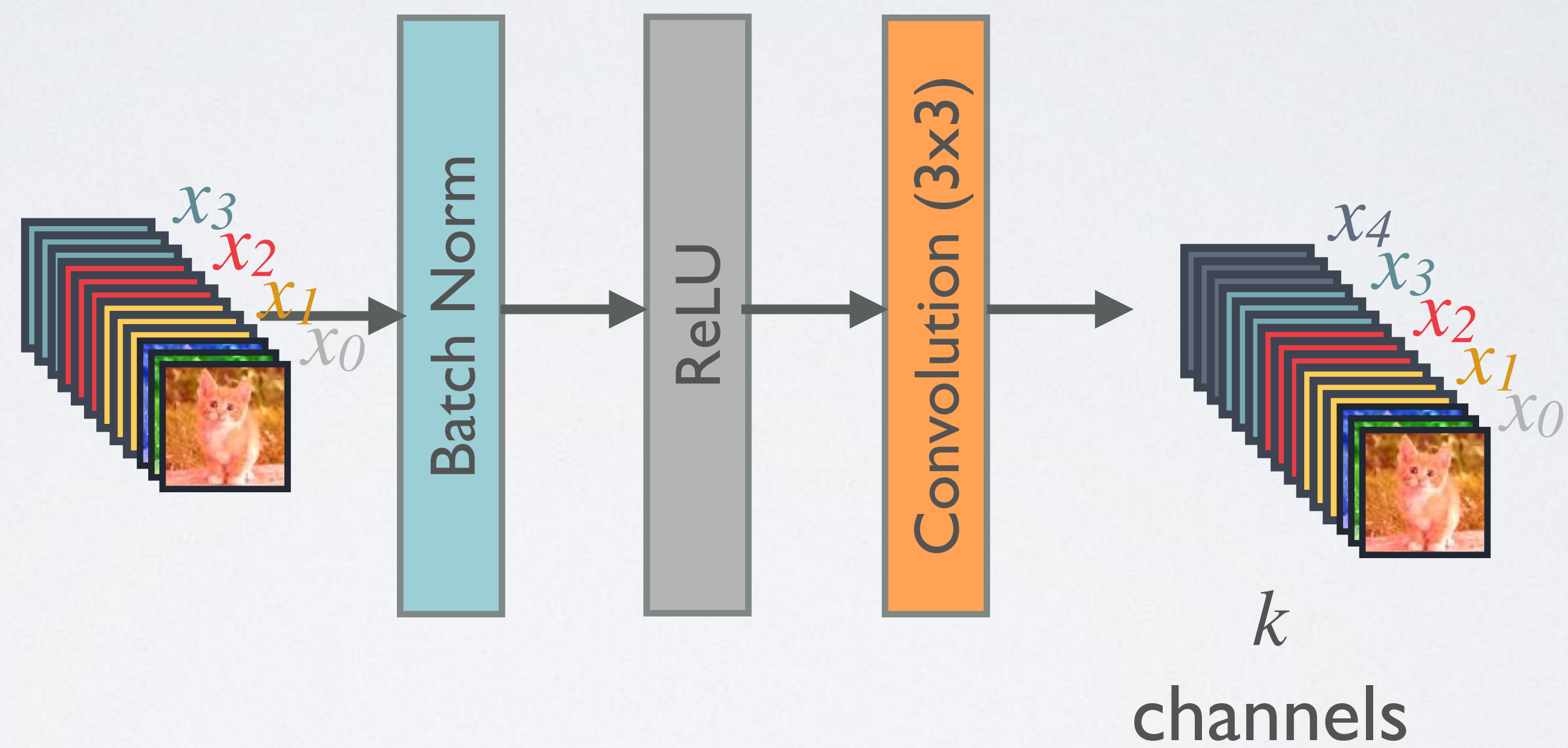


(+) : Element-wise addition

*Deep residual learning for image recognition: [He, Zhang, Ren, Sun] (CVPR 2015)*

# DENSE CONNECTIVITY



C : Channel-wise concatenation

# FORWARD PROPAGATION

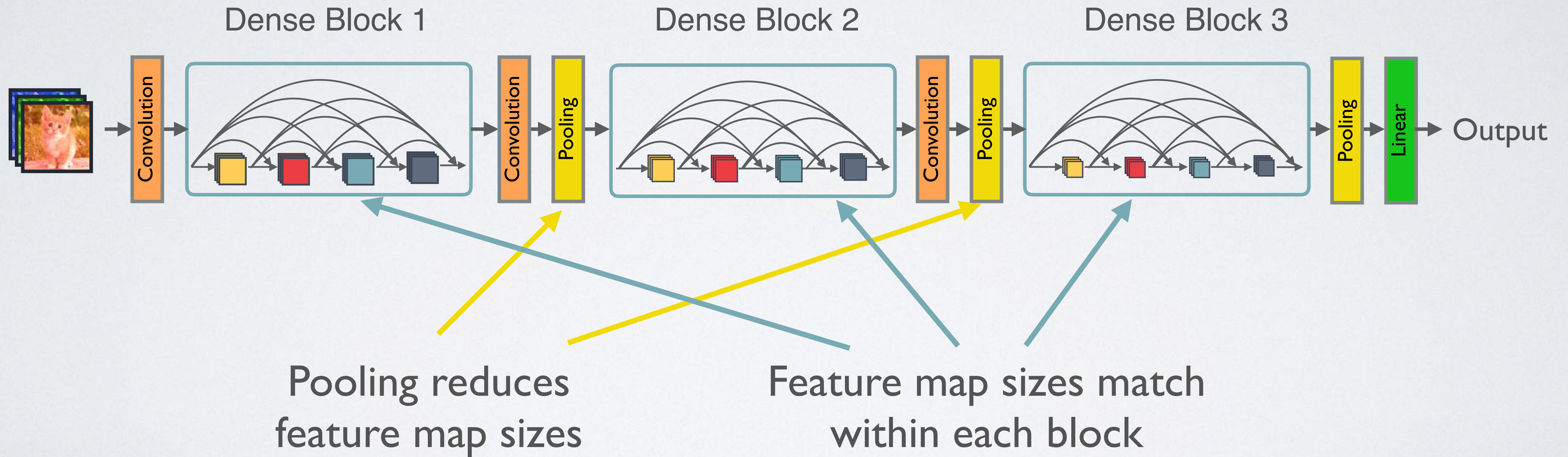# COMPOSITE LAYER IN DENSENET



$$x_5 = h_5([x_0, \ldots, x_4])$$

# DENSENET



Dense Block 1

Dense Block 2

Dense Block 3

Convolution

Convolution

Pooling

Convolution

Pooling

Pooling

Linear

Output

Pooling reduces feature map sizes

Feature map sizes match within each block

# ADVANTAGES OF DENSE CONNECTIVITY

# ADVANTAGE 1: STRONG GRADIENT FLOW



Error Signal

Implicit "deep supervision"

# ADVANTAGE 2: PARAMETER & COMPUTATIONAL EFFICIENCY

**Standard Connectivity:**

Classifier uses most complex (high level) features



$x$      $h_1(x)$      $h_2(x)$      $h_3(x)$      $h_4(x)$
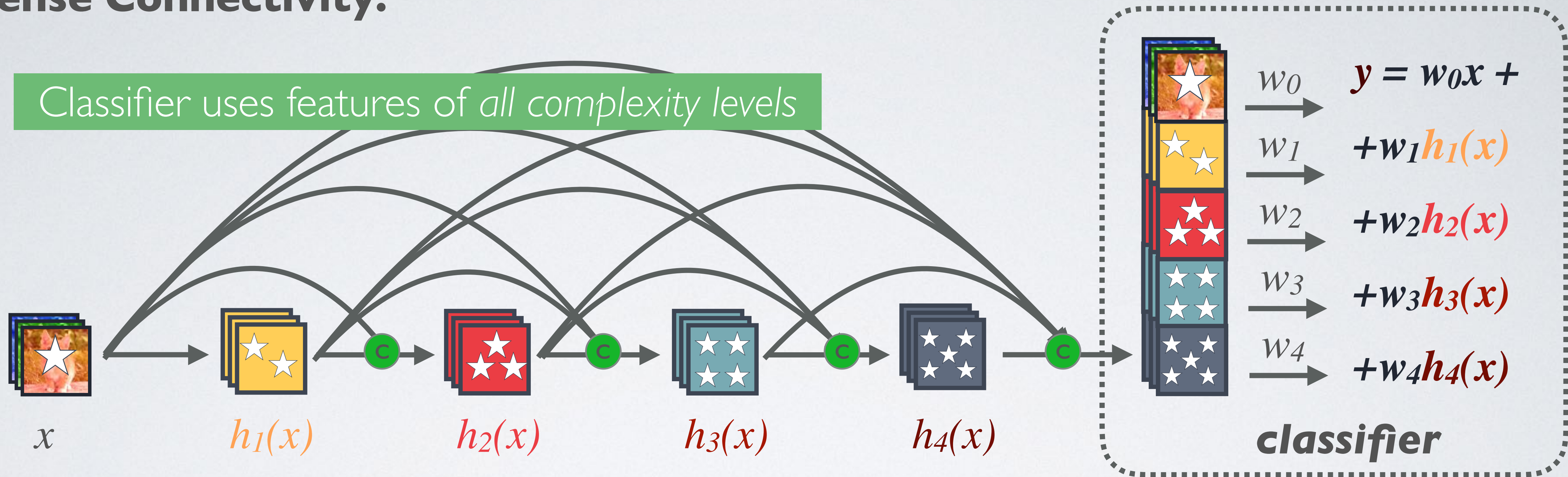
$w_4$

$$y = w_4 h_4(x)$$

*classifier*

★ Increasingly complex features

# ADVANTAGE 3: MAINTAINS LOW COMPLEXITY FEATURES

**Dense Connectivity:**

Classifier uses features of *all complexity levels*

$y = w_0x +$

$w_0$

$+w_1h_1(x)$

$w_1$

$+w_2h_2(x)$

$w_2$

$+w_3h_3(x)$

$w_3$

$+w_4h_4(x)$

$w_4$

*classifier*

$x$     $h_1(x)$     $h_2(x)$     $h_3(x)$     $h_4(x)$

Increasingly complex features

# RESULTS

# RESULTS ON **CIFAR-10**

# RESULTS ON **CIFAR-100**

# RESULTS ON **IMAGENET**



Top-1:  20.27%
Top-5:  5.17%

# MULTI-SCALE DENSENET (Preview)



Classifier 1          Classifier 2          Classifier 3          Classifier 4

cat: 0.2              cat: 0.4              cat: 0.6

0.2 ≱ threshold    0.4 ≱ threshold    0.6 > threshold

*Multi-Scale DenseNet: [Huang, Chen, Li, Wu, van der Maaten, Weinberger] (arXiv Preprint: 1703.09844)*

# MULTI-SCALE DENSENET (Preview)

Test Input

**Inference Speed:**
**~ 2.6x faster than ResNets**
**~ 1.3x faster than DenseNets**

Classifier 1    Classifier 2    Classifier 3    Classifier 4

"Easy" examples    "Hard" examples

*NEW* *Memory efficient Torch implementation:*
*https://github.com/liuzhuang13/DenseNet* *NEW*

## Other implementations:

Our Caffe Implementation
Our memory-efficient Caffe Implementation.
Our memory-efficient PyTorch Implementation.
PyTorch Implementation by Andreas Veit.
PyTorch Implementation by Brandon Amos.
MXNet Implementation by Nicatio.
MXNet Implementation (supports ImageNet) by Xiong Lin.
Tensorflow Implementation by Yixuan Li.
Tensorflow Implementation by Laurent Mazare.

Tensorflow Implementation (with BC structure) by Illarion Khlestov.
Lasagne Implementation by Jan Schlüter.
Keras Implementation by tdeboissiere.
Keras Implementation by Roberto de Moura Estevão Filho.
Keras Implementation (with BC structure) by Somshubra Majumdar.
Chainer Implementation by Toshinori Hanya.
Chainer Implementation by Yasunori Kudo.

# REFERENCES

- Kaiming He, et al. "Deep residual learning for image recognition" CVPR 2016
- Chen-Yu Lee, et al. "Deeply-supervised nets" *AISTATS* 2015
- Gao Huang, et al. "Deep networks with stochastic depth" *ECCV* 2016
- Gao Huang, et al. "Multi-Scale Dense Convolutional Networks for Efficient Prediction" *arXiv preprint arXiv:1703.09844* (2017)
- Geoff Pleiss, et al. "Memory-Efficient Implementation of DenseNets'', *arXiv preprint arXiv: 1707.06990* (2017)