# Convolutional Neural Networks

Samuel Cheng
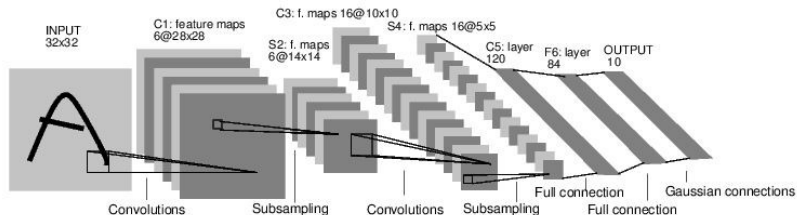
School of ECE
University of Oklahoma

Spring, 2018

# Table of Contents

# Convolutional Neural Networks



[LeNet-5, LeCun 1998]

Fei-Fei Li & Andrej Karpathy & Justin Johnson     Lecture 6 -65     25 Jan 2016
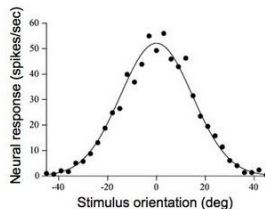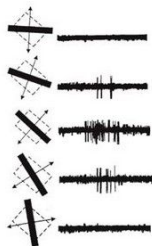
## CNN history

A bit of history:

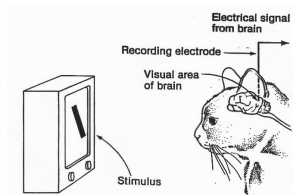**Hubel & Wiesel**,

1959
RECEPTIVE FIELDS OF SINGLE NEURONES IN
THE CAT'S STRIATE CORTEX

1962
RECEPTIVE FIELDS, BINOCULAR INTERACTION
AND FUNCTIONAL ARCHITECTURE IN
THE CAT'S VISUAL CORTEX

1968...



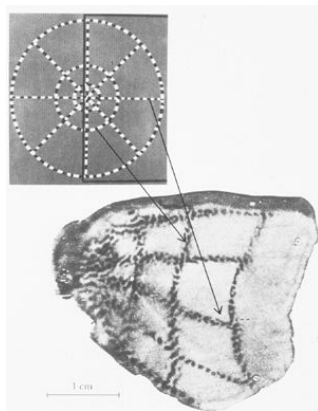Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 66    25 Jan 2016

# CNN history
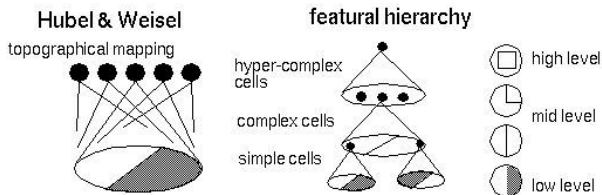
A bit of history

**Topographical mapping in the cortex:**
nearby cells in cortex represented
nearby regions in the visual field

# CNN history

## Hierarchical organization



LGB (lateral geniculate body) $\rightarrow$ simple cells $\rightarrow$ complex cells $\rightarrow$ lower order hypercomplex cells $\rightarrow$ higher order hypercomplex cells

Experiment video, explanation

# CNN history

A bit of history:

**Neurocognitron**
*[Fukushima 1980]*

"sandwich" architecture (SCSCSC…)
simple cells: modifiable parameters
complex cells: perform pooling



Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 70    25 Jan 2016

## CNN history

A bit of history:
**Gradient-based learning applied to document recognition**
*[LeCun, Bottou, Bengio, Haffner 1998]*



LeNet-5



Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 71    25 Jan 2016

S. Cheng  (OU-ECE)                  Convolutional Neural Networks                  Jan 2017        8 / 162

## CNN today

A bit of history:
**ImageNet Classification with Deep Convolutional Neural Networks**
*[Krizhevsky, Sutskever, Hinton, 2012]*

IM🅰️GENET



"AlexNet"

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 72    25 Jan 2016

# CNN today

## Fast-forward to today: ConvNets are everywhere

Classification

Retrieval



*[Krizhevsky 2012]*

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 73    25 Jan 2016

# CNN today

## Fast-forward to today: ConvNets are everywhere

Detection



Segmentation

*[Faster R-CNN: Ren, He, Girshick, Sun 2015]*

*[Farabet et al., 2012]*

Fei-Fei Li & Andrej Karpathy & Justin Johnson      Lecture 6 - 74      25 Jan 2016

## CNN today

### Fast-forward to today: ConvNets are everywhere



self-driving cars



NVIDIA Tegra X1

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 75    25 Jan 2016

## CNN today

### Fast-forward to today: ConvNets are everywhere



*[Taigman et al. 2014]*

*[Simonyan et al. 2014]*

*[Goodfellow 2014]*

Fei-Fei Li & Andrej Karpathy & Justin Johnson          Lecture 6 - 76          25 Jan 2016

## CNN today

### Fast-forward to today: ConvNets are everywhere



*[Toshev, Szegedy 2014]*



*[Mnih 2013]*

Fei-Fei Li & Andrej Karpathy & Justin Johnson     Lecture 6 - 77     25 Jan 2016

## CNN today

### Fast-forward to today: ConvNets are everywhere



[Ciresan et al. 2013]

[Sermanet et al. 2011]
[Ciresan et al.]

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 78    25 Jan 2016

S. Cheng (OU-ECE)    Convolutional Neural Networks    Jan 2017    15 / 162

# CNN today

## Fast-forward to today: ConvNets are everywhere



*[Turaga et al., 2010]*

*[Denil et al. 2014]*

Fei-Fei Li & Andrej Karpathy & Justin Johnson          Lecture 6 - 79          25 Jan 2016

## CNN today



*Whale recognition, Kaggle Challenge*



*Mnih and Hinton, 2010*

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 80    25 Jan 2016

# CNN today



Image Captioning

Describes without errors | Describes with minor errors | Somewhat related to the image | Unrelated to the image

A person riding a motorcycle on a dirt road.

Two dogs play in the grass.

A skateboarder does a trick on a ramp.

A dog is jumping to catch a frisbee.

A group of young people playing a game of frisbee.

Two hockey players are fighting over the puck.

A little girl in a pink hat is blowing bubbles.

A refrigerator filled with lots of food and drinks.

A herd of elephants walking across a dry grass field.

A close up of a cat laying on a couch.

A red motorcycle parked on the side of the road.

A yellow school bus parked in a parking lot.

[Vinyals et al., 2015]

Fei-Fei Li & Andrej Karpathy & Justin Johnson     Lecture 6 - 81     25 Jan 2016

# CNN today



reddit.com/r/deepdream

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 6 - 82    25 Jan 2016

# Motivation of CNN

- A same object under different viewpoints is very different in pixel domain
  - A slightly horizontally shifted image has change imperceivable to us but can confuse naive recognition system
- Ideally, we may want to have shift-invariant features
- In practice, if we have local feature suitable for a particular region, the same feature should work well with other region
  - Weight sharing across space $\rightarrow$ CNN

# Convolution Layer

32x32x3 image



32 height

32 width

3 depth

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 10    27 Jan 2016

# Convolution Layer

32x32x3 image



5x5x3 filter

**Convolve** the filter with the image
i.e. "slide over the image spatially,
computing dot products"

# Convolution Layer

Filters always extend the full depth of the input volume

32x32x3 image



32

32

3

5x5x3 filter



**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer



32x32x3 image

5x5x3 filter $w$

32

32

3

**1 number:**
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer



32x32x3 image
5x5x3 filter

activation map

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Convolution Layer

consider a second, green filter



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation maps**

28

28

1

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



**activation maps**

We stack these up to get a "new image" of size 28x28x6!

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



CONV,
ReLU
e.g. 6
5x5x3
filters

32
32
3

28
28
6

Fei-Fei Li & Andrej Karpathy & Justin Johnson      Lecture 7 - 17      27 Jan 2016

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



32

28

24

CONV, ReLU
e.g. 6
5x5x3
filters

CONV, ReLU
e.g. 10
5x5x**6**
filters

CONV, ReLU

....

32

28

24

3

6

10

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 18    27 Jan 2016

**Preview**



*[From recent Yann LeCun slides]*

Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]



Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 20    27 Jan 2016

one filter =>
one activation map

example 5x5 filters
(32 total)

Activations:

We call the layer convolutional
because it is related to convolution
of two signals:

$$f[x,y] * g[x,y] \quad = \quad \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1, y-n_2]$$

elementwise multiplication and sum of
a filter and the signal (image)

Fei-Fei Li & Andrej Karpathy & Justin Johnson          Lecture 7 - 21          27 Jan 2016

A closer look at spatial dimensions:



**activation map**

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**28**

**28**

1

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:



7

7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**

A closer look at spatial dimensions:

7



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

7

Fei-Fei Li & Andrej Karpathy & Justin Johnson        Lecture 7 - 29        27 Jan 2016

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2
=> 3x3 output!**

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

**doesn't fit!**
cannot apply 3x3 filter on
7x7 input with stride 3.

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 33    27 Jan 2016

Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

(recall:)
(N - F) / stride + 1

Fei-Fei Li & Andrej Karpathy & Justin Johnson        Lecture 7 - 35        27 Jan 2016

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**

Fei-Fei Li & Andrej Karpathy & Justin Johnson        Lecture 7 - 36        27 Jan 2016

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Output volume size:
(32+2*2-5)/1+1 = 32 spatially, so
**32x32x10**

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?



Fei-Fei Li & Andrej Karpathy & Justin Johnson      Lecture 7 - 41      27 Jan 2016

Examples time:

Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?
each filter has 5*5*3 + 1 = 76 params    (+1 for bias)
=> 76*10 = **760**

**Summary**. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters $K$,
  - their spatial extent $F$,
  - the stride $S$,
  - the amount of zero padding $P$.
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and $K$ biases.
- In the output volume, the $d$-th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the $d$-th filter over the input volume with a stride of $S$, and then offset by $d$-th bias.

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 43    27 Jan 2016

**Summary**. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters $K$,
  - their spatial extent $F$,
  - the stride $S$,
  - the amount of zero padding $P$.
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and $K$ biases.
- In the output volume, the $d$-th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the $d$-th filter over the input volume with a stride of $S$, and then offset by $d$-th bias.

Common settings:

K = (powers of 2, e.g. 32, 64, 128, 512)
- F = 3, S = 1, P = 1
- F = 5, S = 1, P = 2
- F = 5, S = 2, P = ? (whatever fits)
- F = 1, S = 1, P = 0

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 44    27 Jan 2016

(btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

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

The brain/neuron view of CONV Layer



32x32x3 image
5x5x3 filter

32

32

3

**1 number:**
the result of taking a dot product between
the filter and this part of the image
(i.e. 5*5*3 = 75-dimensional dot product)

The brain/neuron view of CONV Layer



32x32x3 image
5x5x3 filter

32

32

3

It's just a neuron with local connectivity...

**1 number:**
the result of taking a dot product between
the filter and this part of the image
(i.e. 5*5*3 = 75-dimensional dot product)

The brain/neuron view of CONV Layer



An activation map is a 28x28 sheet of neuron outputs:
1.  Each is connected to a small region in the input
2.  All of them share parameters

"5x5 filter" -> "5x5 receptive field for each neuron"

Fei-Fei Li & Andrej Karpathy & Justin Johnson      Lecture 7 - 51      27 Jan 2016

The brain/neuron view of CONV Layer



E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
(28x28x5)

There will be 5 different
neurons all looking at the same
region in the input volume



Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 52    27 Jan 2016

two more layers to go: POOL/FC



Fei-Fei Li & Andrej Karpathy & Justin Johnson     Lecture 7 - 53     27 Jan 2016

# Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:

# MAX POOLING

Single depth slice



max pool with 2x2 filters
and stride 2

Fei-Fei Li & Andrej Karpathy & Justin Johnson     Lecture 7 - 55     27 Jan 2016

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
    - their spatial extent $F$,
    - the stride $S$,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
    - $W_2 = (W_1 - F)/S + 1$
    - $H_2 = (H_1 - F)/S + 1$
    - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Fei-Fei Li & Andrej Karpathy & Justin Johnson        Lecture 7 - 56        27 Jan 2016

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
  - their spatial extent $F$,
  - the stride $S$,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

Common settings:

F = 2, S = 2
F = 3, S = 2

Fei-Fei Li & Andrej Karpathy & Justin Johnson       Lecture 7 - 57       27 Jan 2016

# Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks

# Demo

ConvNetJS cifar10 demo

# Case Study: 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-CONV-FC]

Fei-Fei Li & Andrej Karpathy & Justin Johnson    Lecture 7 - 60    27 Jan 2016

# AlexNet

## Case Study: AlexNet

*[Krizhevsky et al. 2012]*



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

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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 -  9          May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung            Lecture 9 - 10        May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 11      May 2, 2017

## AlexNet

# 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**

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 12        May 2, 2017

## AlexNet

# Case Study: AlexNet

*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung      Lecture 9 - 13      May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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?

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 14          May 2, 2017

## AlexNet

## Case Study: AlexNet

*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 15        May 2, 2017

## AlexNet

## Case Study: AlexNet

*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 16          May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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
[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)

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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 17          May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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
[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)

**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.

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 18          May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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
[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

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.

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

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 19        May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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
[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

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

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 20        May 2, 2017

# AlexNet

## Case Study: AlexNet

*[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
[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:
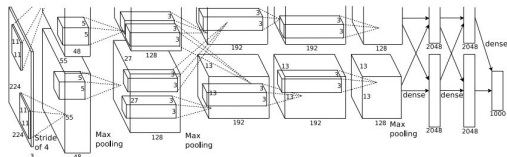Connections with all feature maps in
preceding layer, communication
across GPUs

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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 21        May 2, 2017

# AlexNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung      Lecture 9 - 22      May 2, 2017

# ZFNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 23    May 2, 2017

# ZFNet

## 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%

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 24          May 2, 2017

# VGGNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung         Lecture 9 - 25         May 2, 2017

# 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          VGG16          VGG19

# VGGNet

## Case Study: VGGNet
*[Simonyan and Zisserman, 2014]*

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



AlexNet        VGG16        VGG19

# 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

# 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

[7x7]



AlexNet          VGG16          VGG19

# 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

VGG16

VGG19

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 30          May 2, 2017

# VGGNet

(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



VGG16

# VGGNet

(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 (only forward! ~*2 for bwd)
TOTAL params: 138M parameters



VGG16

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 32          May 2, 2017

# VGGNet

(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
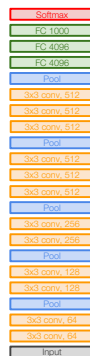CONV3-512: [14x14x512]   memory:  14*14*512=100K   params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]   memory:  14*14*512=100K   params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]         memory:  7*7*512=25K   params: 0
FC: [1x1x4096]  memory:  4096  params: 7*7*512*4096 = **102,760,448**
FC: [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 (only forward! ~*2 for bwd)
TOTAL params: 138M parameters

Note:

Most memory is in early CONV

Most params are in late FC

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 33          May 2, 2017

# VGGNet

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
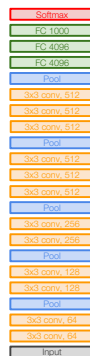POOL2: [7x7x512]  memory:  7*7*512=25K   params: 0
FC: [1x1x4096]  memory:  4096  params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]  memory:  4096  params: 4096*4096 = 16,777,216
FC: [1x1x1000]  memory:  1000 params: 4096*1000 = 4,096,000

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

VGG16

Common names

# VGGNet

## Case Study: VGGNet

*[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



AlexNet    VGG16    VGG19

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 35    May 2, 2017

# VGGNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 36        May 2, 2017

# 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

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 37        May 2, 2017

# 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

# GoogleNet

## 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

# GoogleNet

## 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]

Fei-Fei Li & Justin Johnson & Serena Yeung       Lecture 9 - 40       May 2, 2017

# GoogleNet

## Case Study: GoogLeNet
*[Szegedy et al., 2014]*

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

Example:



Naive Inception module

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 41        May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

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

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



Naive Inception module

Fei-Fei Li & Justin Johnson & Serena Yeung       Lecture 9 - 42       May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Example:

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

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



28x28x128

Module input:
28x28x256

Naive Inception module

# GoogleNet

## Case Study: GoogLeNet

*[Szegedy et al., 2014]*

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

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



Naive Inception module

# GoogleNet

# Case Study: GoogLeNet

*[Szegedy et al., 2014]*

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

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



Naive Inception module

# GoogleNet

## Case Study: GoogLeNet
*[Szegedy et al., 2014]*

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

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



Naive Inception module

# GoogleNet

## 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) = 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

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 47    May 2, 2017

# GoogleNet

## 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**

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 48    May 2, 2017

# GoogleNet

## Case Study: GoogLeNet
[Szegedy et al., 2014]

Example:

Q3: What is output size after filter concatenation?

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



28x28x128   28x28x192   28x28x96   28x28x256

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

Filter concatenation

Module input: 28x28x256

Input

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!

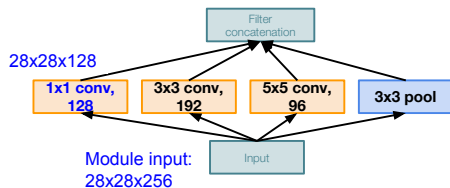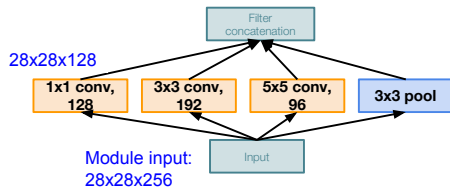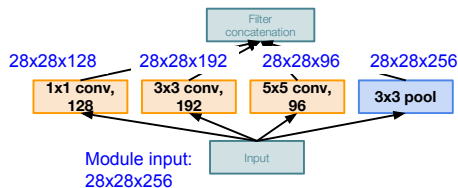Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 49        May 2, 2017
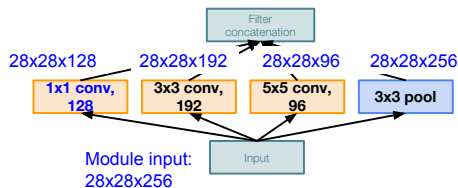
# GoogleNet

## Case Study: GoogLeNet

*[Szegedy et al., 2014]*

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

Example:

Q3: What is output size after filter concatenation?

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

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 |

Filter concatenation

Module input: 28x28x256

Input

Naive Inception module

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 50          May 2, 2017

# GoogleNet

## Reminder: 1x1 convolutions



1x1 CONV
with 32 filters
→

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

56

56

64

56

56

32

# GoogleNet

## Reminder: 1x1 convolutions



56

1x1 CONV
with 32 filters

56

preserves spatial
dimensions, reduces depth!

Projects depth to lower
dimension (combination of
feature maps)

64

56

32

56

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 52          May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]



Naive Inception module

Inception module with dimension reduction

# GoogleNet

## Case Study: GoogLeNet

*[Szegedy et al., 2014]*



Naive Inception module

1x1 conv "bottleneck" layers

Inception module with dimension reduction

# GoogleNet

## Case Study: GoogLeNet

*[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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 55          May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Stack Inception modules
with dimension reduction
on top of each other



Inception module

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 56          May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture



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

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 57        May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture



Stacked Inception
Modules

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 58          May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

[Szegedy et al., 2014]

Full GoogLeNet
architecture



Classifier output

# GoogleNet

## Case Study: GoogLeNet
*[Szegedy et al., 2014]*

**Full GoogLeNet architecture**



Classifier output
(removed expensive FC layers!)

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 60        May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Full GoogLeNet
architecture



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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 61          May 2, 2017

# GoogleNet

## Case Study: GoogLeNet

*[Szegedy et al., 2014]*

Full GoogLeNet
architecture



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

# GoogleNet

## Case Study: GoogLeNet

*[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

Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 63        May 2, 2017

# ResNet

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



Figure copyright Kaiming He, 2016. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 64    May 2, 2017

# 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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 65          May 2, 2017

## ResNet

## Case Study: ResNet

*[He et al., 2015]*

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

# ResNet

## 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]

# ResNet

## Case Study: ResNet

*[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!

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 68          May 2, 2017

## ResNet

## Case Study: ResNet
*[He et al., 2015]*

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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 69        May 2, 2017

# ResNet

## Case Study: ResNet
*[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.

Fei-Fei Li & Justin Johnson & Serena Yeung     Lecture 9 - 70     May 2, 2017

# ResNet

## 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

# ResNet

## Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:
- Stack residual blocks
- Every residual block has two 3x3 conv layers



Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 73        May 2, 2017

# 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



Residual block

Beginning conv layer

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 75          May 2, 2017

# 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)



relu

$F(x) + x$ ⊕

3x3 conv

$F(x)$    relu

3x3 conv

X
Residual block

X
identity

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 /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

Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 76    May 2, 2017

# ResNet

## Case Study: ResNet

*[He et al., 2015]*

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



Fei-Fei Li & Justin Johnson & Serena Yeung        Lecture 9 - 77        May 2, 2017

# ResNet

## Case Study: ResNet

*[He et al., 2015]*

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



28x28x256
output

1x1 conv, 256

3x3 conv, 64

1x1 conv, 64

28x28x256
input

# 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

Fei-Fei Li & Justin Johnson & Serena Yeung       Lecture 9 - 79       May 2, 2017

## ResNet

## Case Study: ResNet

*[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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 80          May 2, 2017

## ResNet

## 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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 81          May 2, 2017

## ResNet

## 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
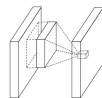
---

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
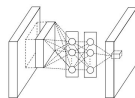  - COCO Segmentation: **12%** better than 2nd

---

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

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 82          May 2, 2017

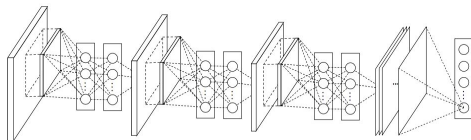# Other Architecture: Network in Network [Lin et al. 2014]

- Introduce MLPConv
  - Inspire inception modules in GoogleNet and residual blocks in ResNets
- Popularize 1x1 conv
- Popularize global average pooling in place of full connected layers



(a) Linear convolution layer    (b) Mlpconv layer

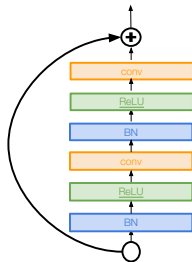# Other Architecture

Improving ResNets...

## Identity Mappings in Deep Residual Networks

*[He et al. 2016]*

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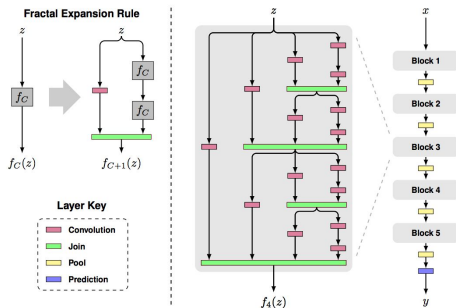- Gives better performance



Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 93          May 2, 2017

# More Skip Connection Tricks

Beyond ResNets...

## FractalNet: Ultra-Deep Neural Networks without Residuals

[Larsson et al. 2017]

- Argues that key is transitioning effectively from shallow to deep and residual representations are not necessary
- Fractal architecture with both shallow and deep paths to output
- Trained with dropping out sub-paths
- Full network at test time



Figures copyright Larsson et al., 2017. Reproduced with permission.

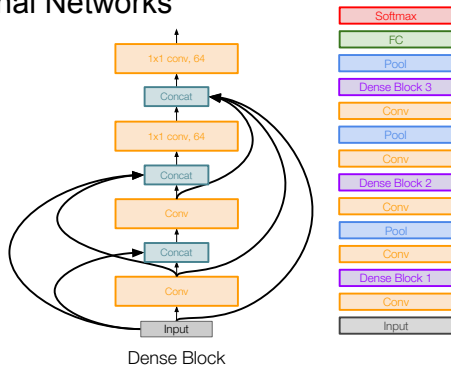Fei-Fei Li & Justin Johnson & Serena Yeung    Lecture 9 - 97    May 2, 2017

# More Skip Connection Tricks

Beyond ResNets...

## 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

# More Skip Connection Tricks

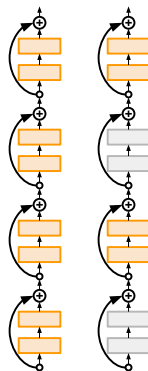Improving ResNets...

## 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



Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 - 96          May 2, 2017

# SqueezeNet Strategies

Strategy 1: Replace $3 \times 3$ by $1 \times 1$ filters

Strategy 2: Decrease # input channels of $3 \times 3$ filters

Strategy 3: Delay downsampling of the networks to increase the size of activation/feature map

# SqueezeNet (Con't)

Efficient networks...

SqueezeNet: AlexNet-level Accuracy With 50x Fewer Parameters and <0.5Mb Model Size

*[Iandola et al. 2017]*

- Fire modules consisting of a 'squeeze' layer with 1x1 filters feeding an 'expand' layer with 1x1 and 3x3 filters
- AlexNet level accuracy on ImageNet with 50x fewer parameters
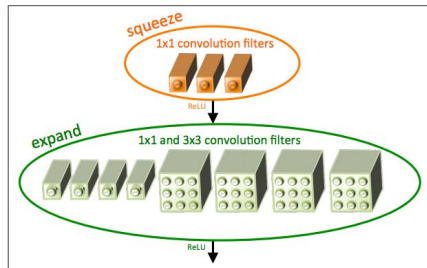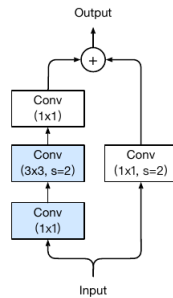- Can compress to 510x smaller than AlexNet (0.5Mb)



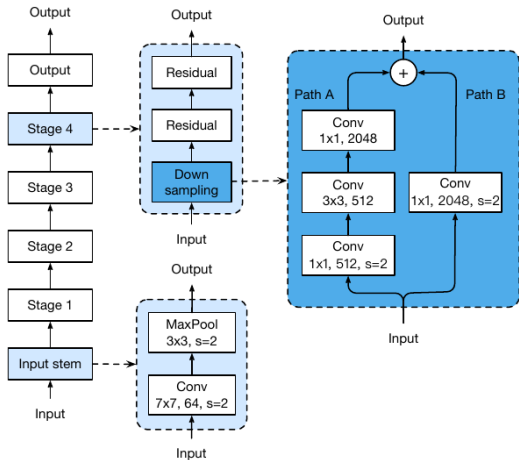Figure copyright Iandola, Han, Moskewicz, Ashraf, Dally, Keutzer, 2017. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung          Lecture 9 -          May 2, 2017

# Bag of tricks for ResNet



(a) ResNet-B     (b) ResNet-C     (c) ResNet-D

# Other Architecture

Improving 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

# ResNeXt [Xie et al. 2016]

- Also from creators of ResNet
- Increases width of residual block through parallel pathways ("cardinalities")
- Same spirit as Inception but simpler



Figure 1. **Left**: A block of ResNet [14]. **Right**: A block of ResNeXt with cardinality = 32, with roughly the same complexity. A layer is shown as (# in channels, filter size, # out channels).

# ResNeXt vs ResNet

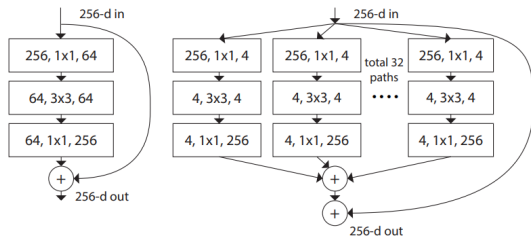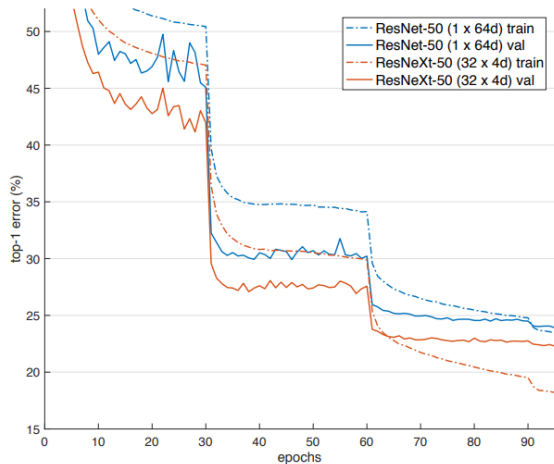| stage | output | ResNet-50 | ResNeXt-50 (32×4d) |
|---|---|---|---|
| conv1 | 112×112 | 7×7, 64, stride 2 | 7×7, 64, stride 2 |
| conv2 | 56×56 | 3×3 max pool, stride 2 $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix} \times 3$ | 3×3 max pool, stride 2 $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128,\ C{=}32 \\ 1\times1,\ 256 \end{bmatrix} \times 3$ |
| conv3 | 28×28 | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256,\ C{=}32 \\ 1\times1,\ 512 \end{bmatrix} \times 4$ |
| conv4 | 14×14 | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512,\ C{=}32 \\ 1\times1,\ 1024 \end{bmatrix} \times 6$ |
| conv5 | 7×7 | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1\times1,\ 1024 \\ 3\times3,\ 1024,\ C{=}32 \\ 1\times1,\ 2048 \end{bmatrix} \times 3$ |
| | 1×1 | global average pool 1000-d fc, softmax | global average pool 1000-d fc, softmax |
| # params. | | $25.5\times10^6$ | $25.0\times10^6$ |
| FLOPs | | $4.1\times10^9$ | $4.2\times10^9$ |

# Group filter/convolution



Image credit: Yani Ioannou

# Group filter interpretation of ResNeXt



*equivalent*

(a)    (b)    (c)

# Depthwise separable convolution (MobileNet)



Depth-wise convolution



Point-wise convolution
Image credit: Chi-Feng Wang

# Depthwise separable convolution (more channels)



Depth-wise convolution

Point-wise convolution
Image credit: Chi-Feng Wang

# Squeeze-and-excitation

# Squeeze-and-excitation

# Squeeze-and-excitation

# Inverted residual (MobileNetV2)

- We want to squeeze before conv in original ResNet to save compute and reduce # parameters
- It is not as much an issue with depthwise separable conv
  - It makes more sense to add skip connection to the more information-densed "squeezed" layer rather than a thicker layer

(a) Residual block          (b) Inverted residual block



Figure 3: The difference between residual block [8, 30] and inverted residual. Diagonally hatched layers do not use non-linearities. We use thickness of each block to indicate its relative number of channels. Note how classical residuals connects the layers with high number of channels, whereas the inverted residuals connect the bottlenecks. Best viewed in color.

# Fused-MBConv

- Use FusedMBConv in earlier layers
  - Not many channels yet
- Use MBConv in later layers
  - Many channels in later layers



Executes faster on Edge TPU

# Neural architecture search (works from Google brain)

- Neural architecture search with reinforcement learning
  - first paper in the area
- Learning transferable architectures for scalable image recognition
  - aka NASNet
  - Learn cells and duplicate them
- MnasNet: Platform-Aware Neural Architecture Search for Mobile
  - Include latency in its objective function
- EfficientNet
  - Introduce compound scaling: $d = \alpha^{\phi}, w = \beta^{\phi}, r = \gamma^{\phi}$
  - Optimize over flop rather than latency. Include memory usage in the objective function as well
- EfficientNetV2
  - More tweaks over EfficientNet. Include "fused-MBConv" as an option
  - Use progressive (curriculum) learning

# EfficientNet Result

| Model | Top-1 Acc. | Top-5 Acc. | #Params | Ratio-to-EfficientNet | #FLOPs | Ratio-to-EfficientNet |
|---|---|---|---|---|---|---|
| **EfficientNet-B0** | **77.1%** | **93.3%** | **5.3M** | **1x** | **0.39B** | **1x** |
| ResNet-50 (He et al., 2016) | 76.0% | 93.0% | 26M | 4.9x | 4.1B | 11x |
| DenseNet-169 (Huang et al., 2017) | 76.2% | 93.2% | 14M | 2.6x | 3.5B | 8.9x |
| **EfficientNet-B1** | **79.1%** | **94.4%** | **7.8M** | **1x** | **0.70B** | **1x** |
| ResNet-152 (He et al., 2016) | 77.8% | 93.8% | 60M | 7.6x | 11B | 16x |
| DenseNet-264 (Huang et al., 2017) | 77.9% | 93.9% | 34M | 4.3x | 6.0B | 8.6x |
| Inception-v3 (Szegedy et al., 2016) | 78.8% | 94.4% | 24M | 3.0x | 5.7B | 8.1x |
| Xception (Chollet, 2017) | 79.0% | 94.5% | 23M | 3.0x | 8.4B | 12x |
| **EfficientNet-B2** | **80.1%** | **94.9%** | **9.2M** | **1x** | **1.0B** | **1x** |
| Inception-v4 (Szegedy et al., 2017) | 80.0% | 95.0% | 48M | 5.2x | 13B | 13x |
| Inception-resnet-v2 (Szegedy et al., 2017) | 80.1% | 95.1% | 56M | 6.1x | 13B | 13x |
| **EfficientNet-B3** | **81.6%** | **95.7%** | **12M** | **1x** | **1.8B** | **1x** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 95.6% | 84M | 7.0x | 32B | 18x |
| PolyNet (Zhang et al., 2017) | 81.3% | 95.8% | 92M | 7.7x | 35B | 19x |
| **EfficientNet-B4** | **82.9%** | **96.4%** | **19M** | **1x** | **4.2B** | **1x** |
| SENet (Hu et al., 2018) | 82.7% | 96.2% | 146M | 7.7x | 42B | 10x |
| NASNet-A (Zoph et al., 2018) | 82.7% | 96.2% | 89M | 4.7x | 24B | 5.7x |
| AmoebaNet-A (Real et al., 2019) | 82.8% | 96.1% | 87M | 4.6x | 23B | 5.5x |
| PNASNet (Liu et al., 2018) | 82.9% | 96.2% | 86M | 4.5x | 23B | 6.0x |
| **EfficientNet-B5** | **83.6%** | **96.7%** | **30M** | **1x** | **9.9B** | **1x** |
| AmoebaNet-C (Cubuk et al., 2019) | 83.5% | 96.5% | 155M | 5.2x | 41B | 4.1x |
| **EfficientNet-B6** | **84.0%** | **96.8%** | **43M** | **1x** | **19B** | **1x** |
| **EfficientNet-B7** | **84.3%** | **97.0%** | **66M** | **1x** | **37B** | **1x** |
| GPipe (Huang et al., 2018) | 84.3% | 97.0% | 557M | 8.4x | - | - |

We omit ensemble and multi-crop models (Hu et al., 2018), or models pretrained on 3.5B Instagram images (Mahajan et al., 2018).

# EfficientNetV2 Result

|  | Model | Top-1 Acc. | Params | FLOPs | Infer-time(ms) | Train-time (hours) |
|---|---|---|---|---|---|---|
| | EfficientNet-B3 (Tan & Le, 2019a) | 81.5% | 12M | 1.9B | 19 | 10 |
| | EfficientNet-B4 (Tan & Le, 2019a) | 82.9% | 19M | 4.2B | 30 | 21 |
| | EfficientNet-B5 (Tan & Le, 2019a) | 83.7% | 30M | 10B | 60 | 43 |
| | EfficientNet-B6 (Tan & Le, 2019a) | 84.3% | 43M | 19B | 97 | 75 |
| | EfficientNet-B7 (Tan & Le, 2019a) | 84.7% | 66M | 38B | 170 | 139 |
| | RegNetY-8GF (Radosavovic et al., 2020) | 81.7% | 39M | 8B | 21 | - |
| | RegNetY-16GF (Radosavovic et al., 2020) | 82.9% | 84M | 16B | 32 | - |
| | ResNeSt-101 (Zhang et al., 2020) | 83.0% | 48M | 13B | 31 | - |
| | ResNeSt-200 (Zhang et al., 2020) | 83.9% | 70M | 36B | 76 | - |
| | ResNeSt-269 (Zhang et al., 2020) | 84.5% | 111M | 78B | 160 | - |
| ConvNets | TResNet-L (Ridnik et al., 2020) | 83.8% | 56M | - | 45 | - |
| & Hybrid | TResNet-XL (Ridnik et al., 2020) | 84.3% | 78M | - | 66 | - |
| | EfficientNet-X (Li et al., 2021) | 84.7% | 73M | 91B | - | - |
| | NFNet-F0 (Brock et al., 2021) | 83.6% | 72M | 12B | 30 | 8.9 |
| | NFNet-F1 (Brock et al., 2021) | 84.7% | 133M | 36B | 70 | 20 |
| | NFNet-F2 (Brock et al., 2021) | 85.1% | 194M | 63B | 124 | 36 |
| | NFNet-F3 (Brock et al., 2021) | 85.7% | 255M | 115B | 203 | 65 |
| | NFNet-F4 (Brock et al., 2021) | 85.9% | 316M | 215B | 309 | 126 |
| | LambdaResNet-420-hybrid (Bello, 2021) | 84.9% | 125M | - | - | 67 |
| | BotNet-T7-hybrid (Srinivas et al., 2021) | 84.7% | 75M | 46B | - | 95 |
| | BiT-M-R152x2 (21k) (Kolesnikov et al., 2020) | 85.2% | 236M | 135B | 500 | - |
| | ViT-B/32 (Dosovitskiy et al., 2021) | 73.4% | 88M | 13B | 13 | - |
| | ViT-B/16 (Dosovitskiy et al., 2021) | 74.9% | 87M | 56B | 68 | - |
| | DeiT-B (ViT+reg) (Touvron et al., 2021) | 81.8% | 86M | 18B | 19 | - |
| Vision | DeiT-B-384 (ViT+reg) (Touvron et al., 2021) | 83.1% | 86M | 56B | 68 | - |
| Transformers | T2T-ViT-19 (Yuan et al., 2021) | 81.4% | 39M | 8.4B | - | - |
| | T2T-ViT-24 (Yuan et al., 2021) | 82.2% | 64M | 13B | - | - |
| | ViT-B/16 (21k) (Dosovitskiy et al., 2021) | 84.6% | 87M | 56B | 68 | - |
| | ViT-L/16 (21k) (Dosovitskiy et al., 2021) | 85.3% | 304M | 192B | 195 | 172 |
| | **EfficientNetV2-S** | 83.9% | 22M | 8.8B | 24 | 7.1 |
| | **EfficientNetV2-M** | 85.1% | 54M | 24B | 57 | 13 |
| ConvNets | **EfficientNetV2-L** | 85.7% | 120M | 53B | 98 | 24 |
| (ours) | **EfficientNetV2-S (21k)** | 84.9% | 22M | 8.8B | 24 | 9.0 |
| | **EfficientNetV2-M (21k)** | 86.2% | 54M | 24B | 57 | 15 |
| | **EfficientNetV2-L (21k)** | 86.8% | 120M | 53B | 98 | 26 |
| | **EfficientNetV2-XL (21k)** | 87.3% | 208M | 94B | - | 45 |

We do not include models pretrained on non-public Instagram/JFT images, or models with extra distillation or ensemble.

# Summary of CNN tricks

- Parallel filters
  - Inception module
- Reduce number of parameters
  - Combination of small filters
  - Bottleneck layer using 1x1 conv
  - Group conv filter
  - Depth-wise separable filter (mobile net)
- Skip connections
  - Residual blocks
- Inverted residual blocks (combined with depth-wise separable filters)
  - Aka MBConv
  - Fused mbconv to match TPU architecture
- Channel importance scaling
  - Squeeze-and-excitation module
- Architecture search
  - RL learning with latency or #flops as regularization penalty