

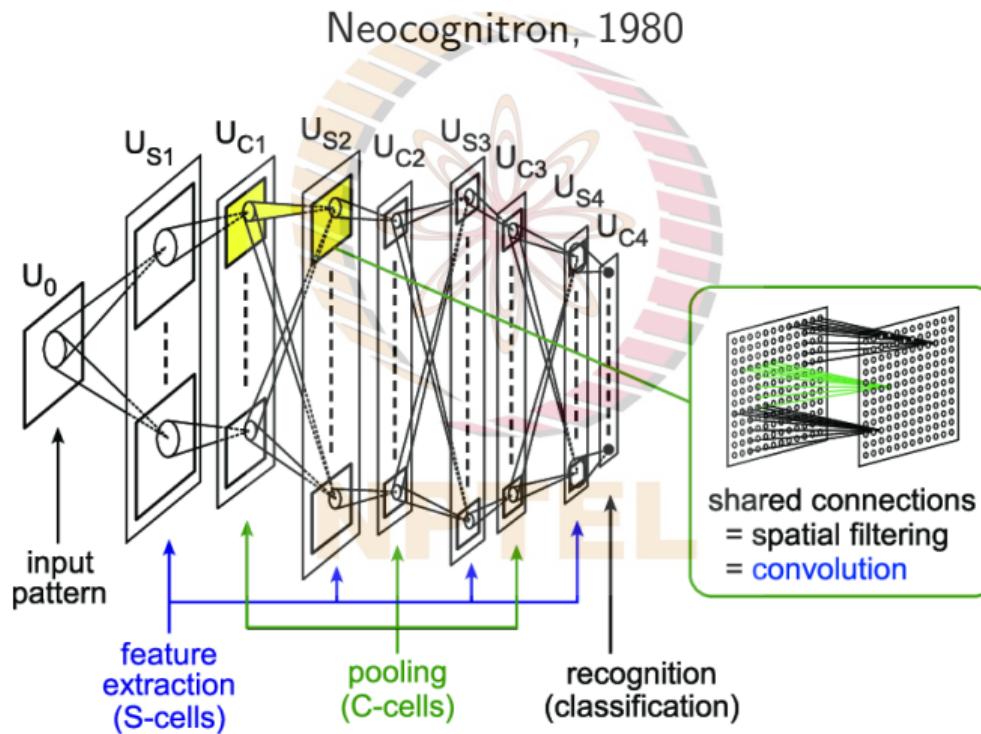
CNN Architectures for Image Classification: AlexNet, VGG

Vineeth N Balasubramanian

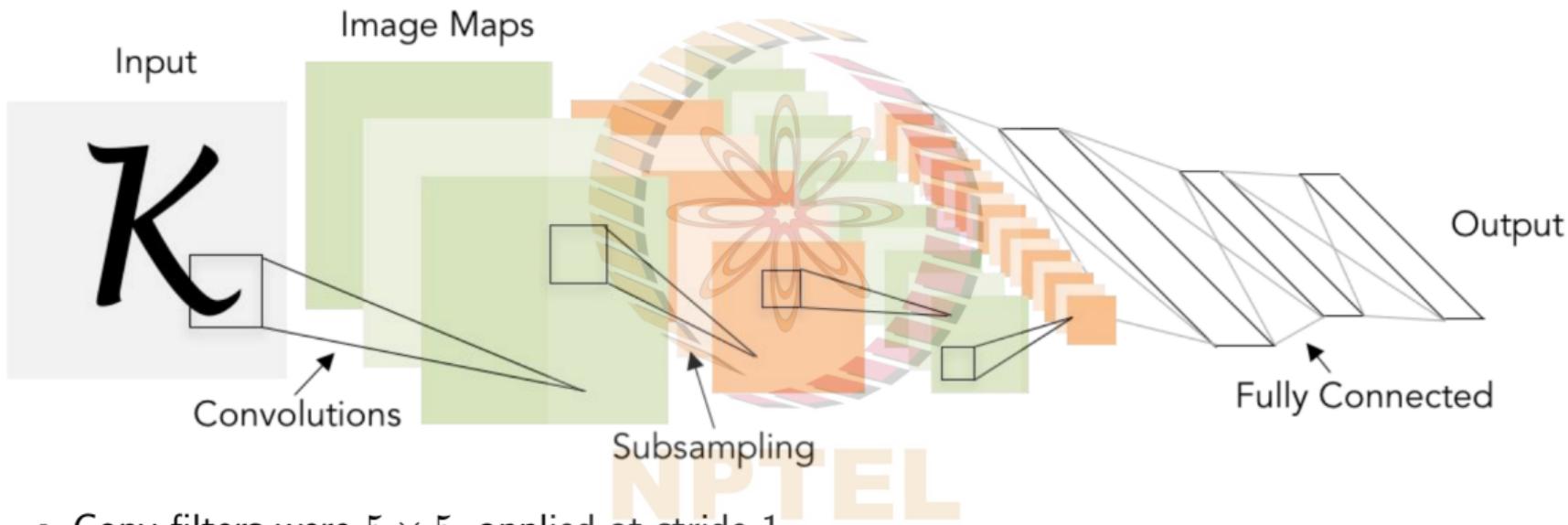
Department of Computer Science and Engineering
Indian Institute of Technology, Hyderabad



History of CNNs



LeNet-5 (1989-1998)



- Conv filters were 5×5 , applied at stride 1
- Subsampling (Pooling) layers were 2×2 applied at stride 2
- **Overall Architecture:** [CONV-POOL-CONV-POOL-FC-FC]

Credit: Fei-Fei Li, Justin Johnson and Serena Yeung, CS231n course, Stanford, Spring 2019

ImageNet Classification Challenge

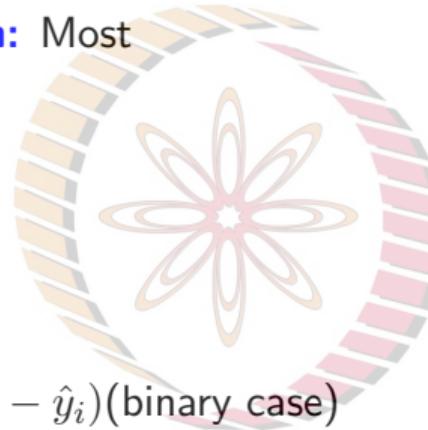
- Image database organized according to WordNet hierarchy (currently only nouns)
- Currently, over five hundred images per node
- Started the ImageNet LSVRC in 2010, for benchmarking of methods for image classification
- Performance measure in Top-1 error and Top-5 error
- <http://www.image-net.org/>



Loss Functions: Beyond Mean Square Error

- **Cross-Entropy Loss Function:** Most popular for classification
- Given by:

$$\begin{aligned} L &= -\frac{1}{C} \sum_{i=1}^C y_i \log \hat{y}_i \\ &= -y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \text{(binary case)} \end{aligned}$$

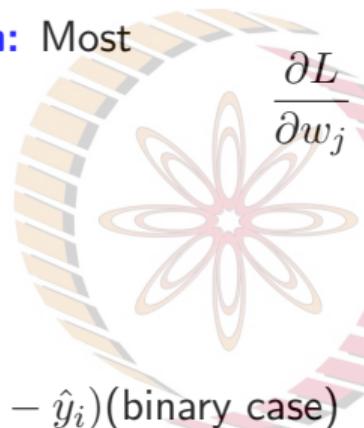


- When activation function is sigmoid $(\sigma(x) = \frac{1}{1+e^{-x}})$, derivative of cross-entropy loss function, $\frac{\partial L}{\partial w_j}$, w.r.t. a weight in last layer, w_j , is:

Loss Functions: Beyond Mean Square Error

- **Cross-Entropy Loss Function:** Most popular for classification
- Given by:

$$L = -\frac{1}{C} \sum_{i=1}^C y_i \log \hat{y}_i$$
$$= -y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \text{(binary case)}$$


$$\begin{aligned}\frac{\partial L}{\partial w_j} &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{1-y}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{1-y}{1-\sigma(z)} \right) \sigma'(z) x_j \\ &= \frac{1}{n} \sum_x \frac{\sigma'(z) x_j}{\sigma(z)(1-\sigma(z))} (\sigma(z) - y) \\ &= \frac{1}{n} \sum_x x_j (\sigma(z) - y)\end{aligned}$$

- When activation function is sigmoid $(\sigma(x) = \frac{1}{1+e^{-x}})$, derivative of cross-entropy loss function, $\frac{\partial L}{\partial w_j}$, w.r.t. a weight in last layer, w_j , is:

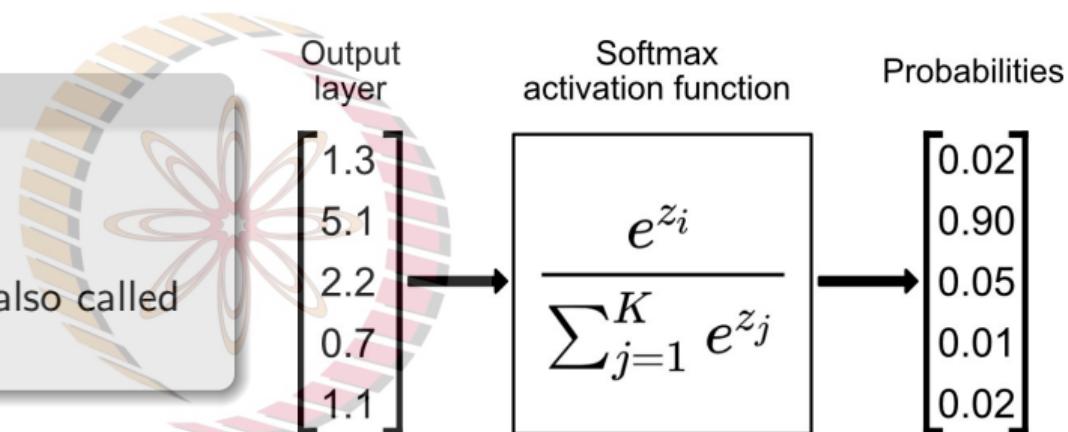
Note the last term in the final expression, very similar to gradient of MSE loss function

Activation Function in Output Layer

Softmax activation function

$$a_j = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

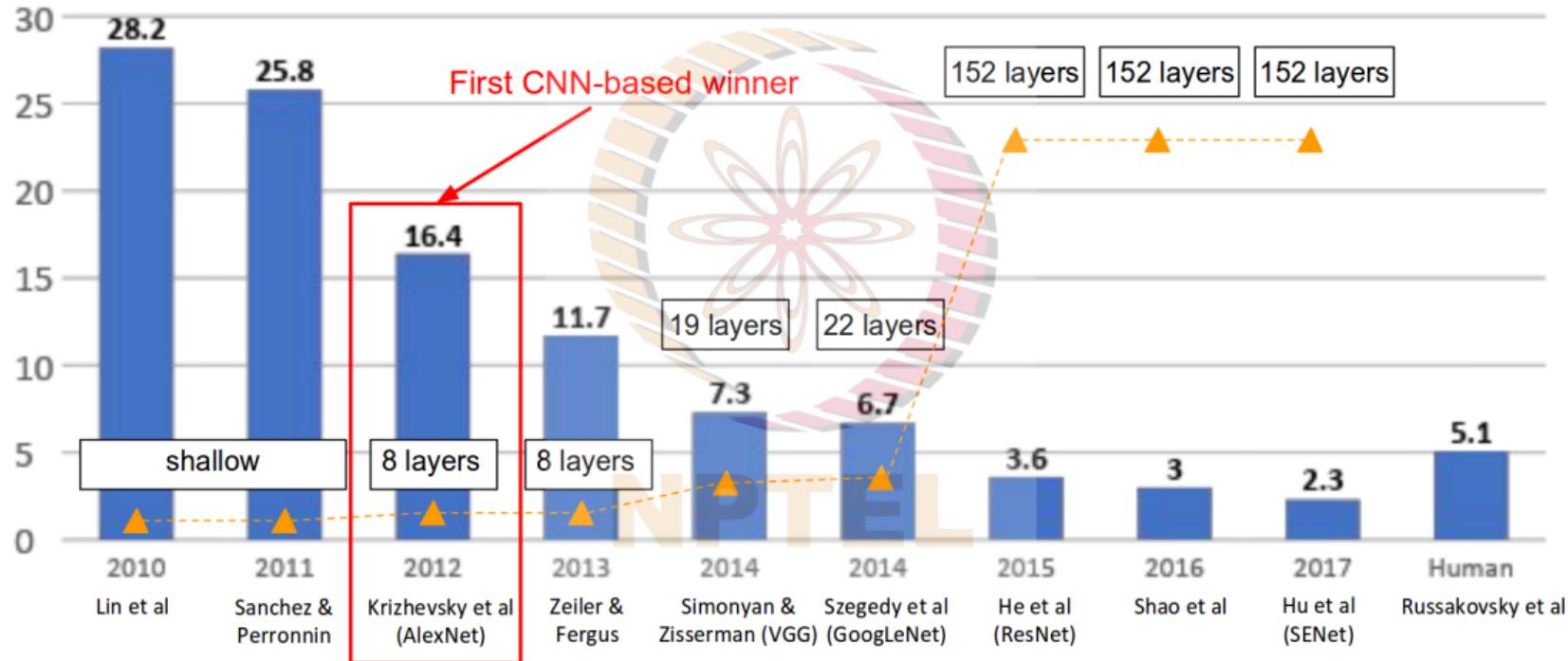
Helps convert output layer values (also called **logits**) to probability scores



NPTEL

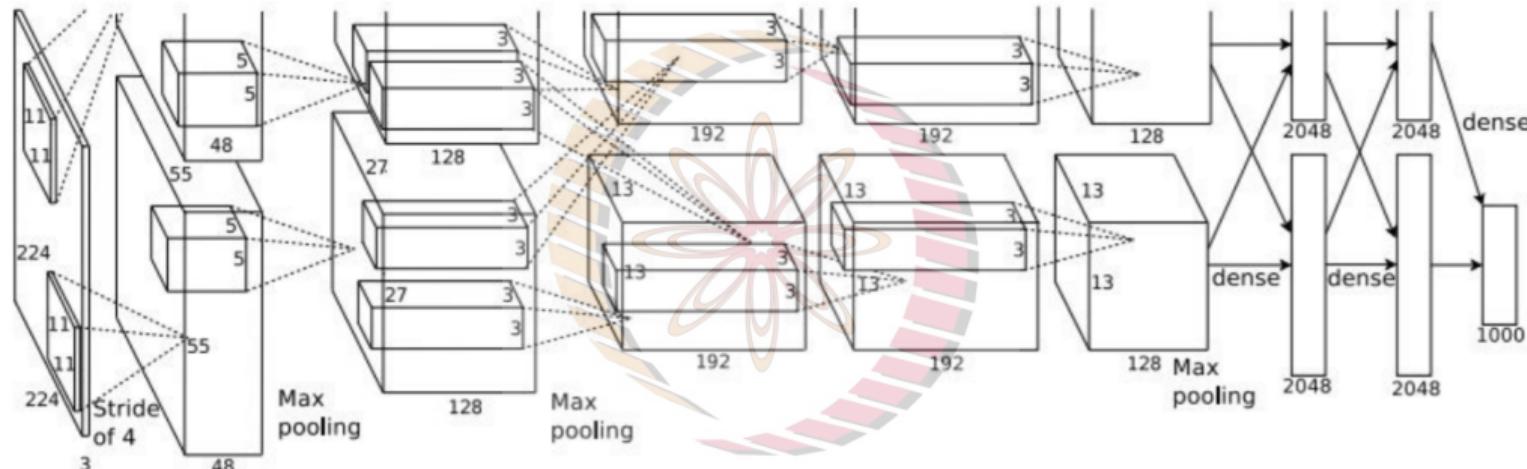
Credit: Dario Redicic, TowardsDataScience blog

Winners of ImageNet Classification Challenge



Credit: Fei-Fei Li, Justin Johnson and Serena Yeung, CS231n course, Stanford, Spring 2019

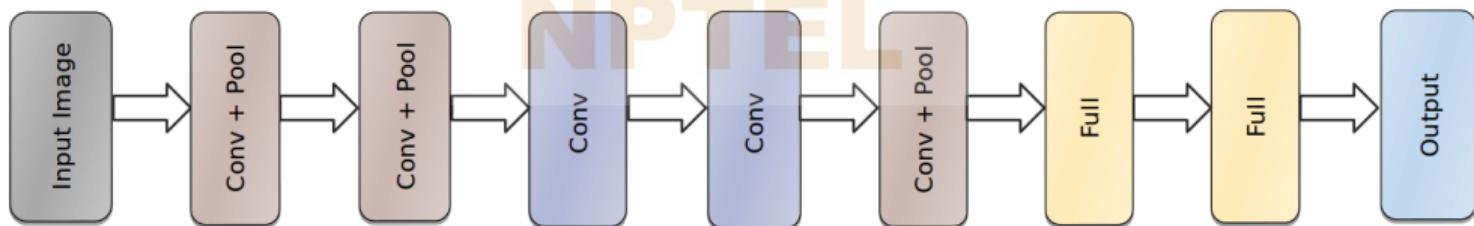
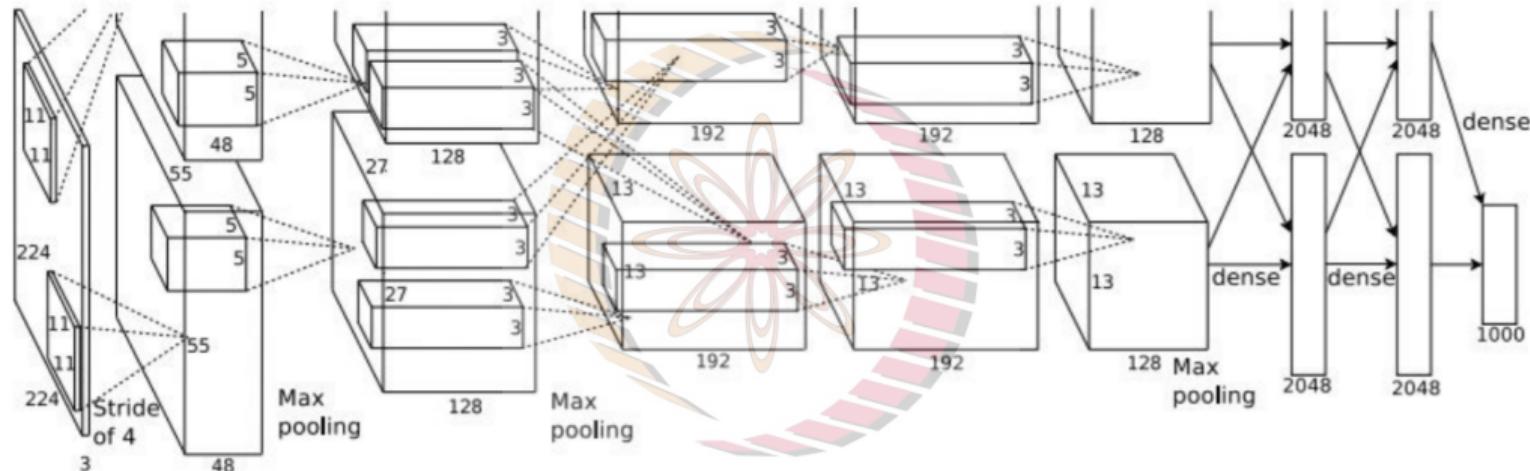
AlexNet¹



- Winner of ImageNet LSVRC-2012
- Overall architecture design similar to LeNet; but deeper with conv layers stacked on top of each other
- Trained over 1.2M images using SGD with regularization

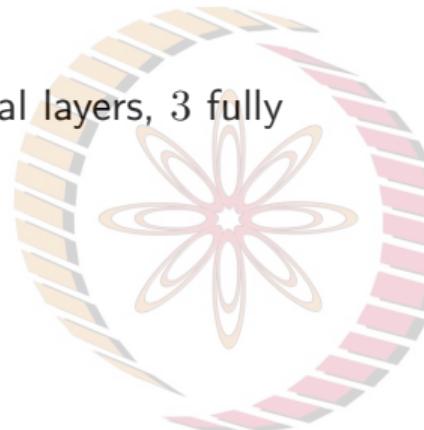
¹Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." NIPS 2012.

AlexNet

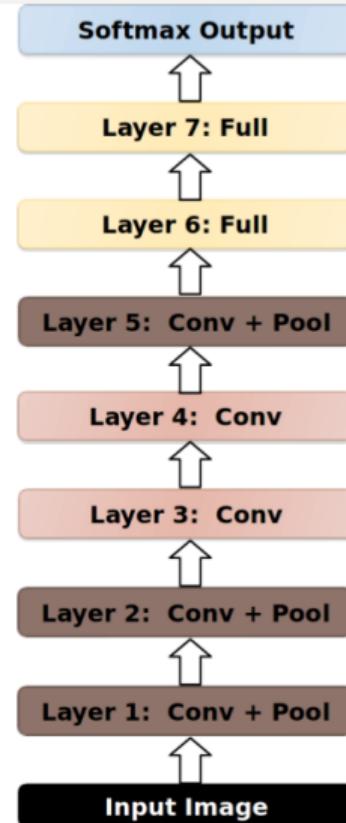


AlexNet

- 8 layers in total (5 convolutional layers, 3 fully connected layers)
- Trained on ImageNet Dataset

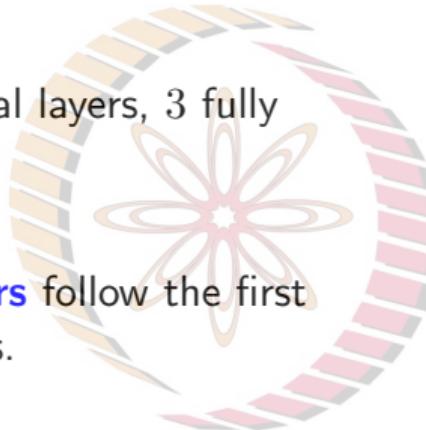


NPTEL

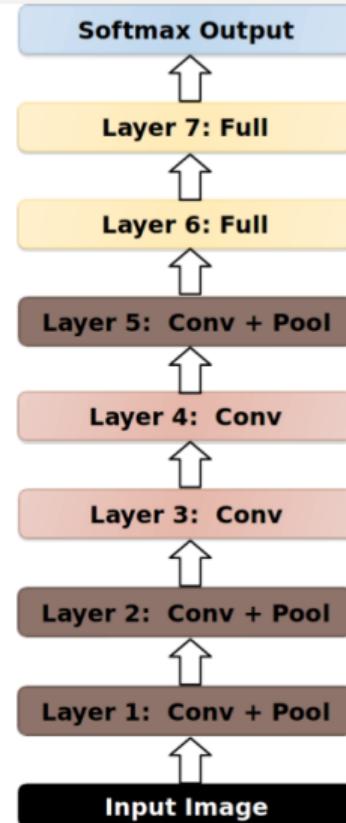


AlexNet

- 8 layers in total (5 convolutional layers, 3 fully connected layers)
- Trained on ImageNet Dataset
- **Response normalization layers** follow the first and second convolutional layers.

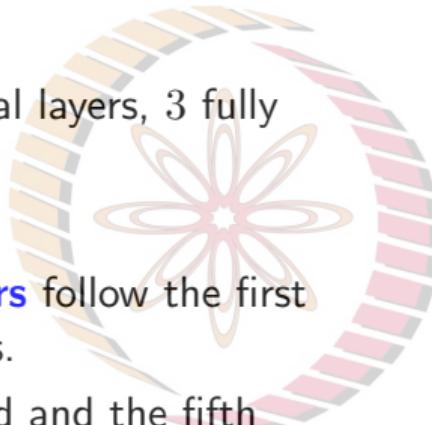


NPTEL

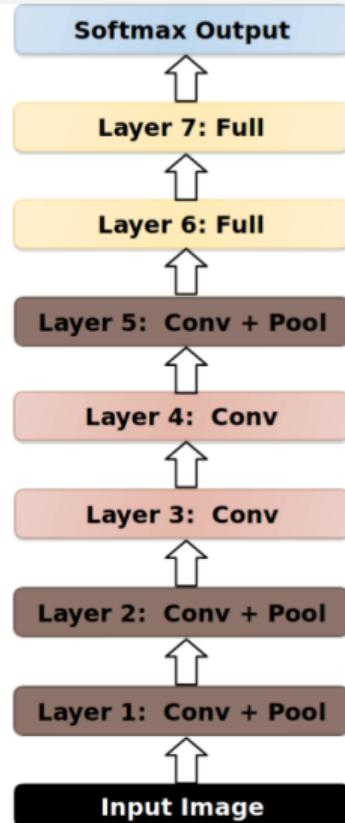


AlexNet

- 8 layers in total (5 convolutional layers, 3 fully connected layers)
- Trained on ImageNet Dataset
- **Response normalization layers** follow the first and second convolutional layers.
- Max-pooling follow first, second and the fifth convolutional layers

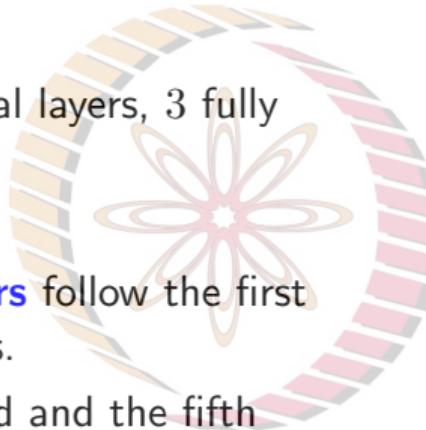


NPTEL

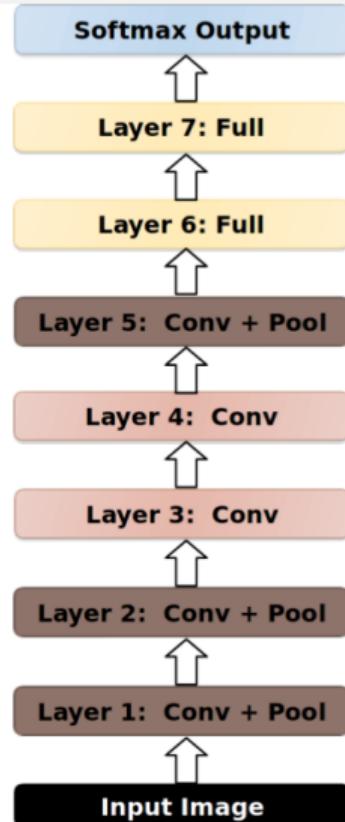


AlexNet

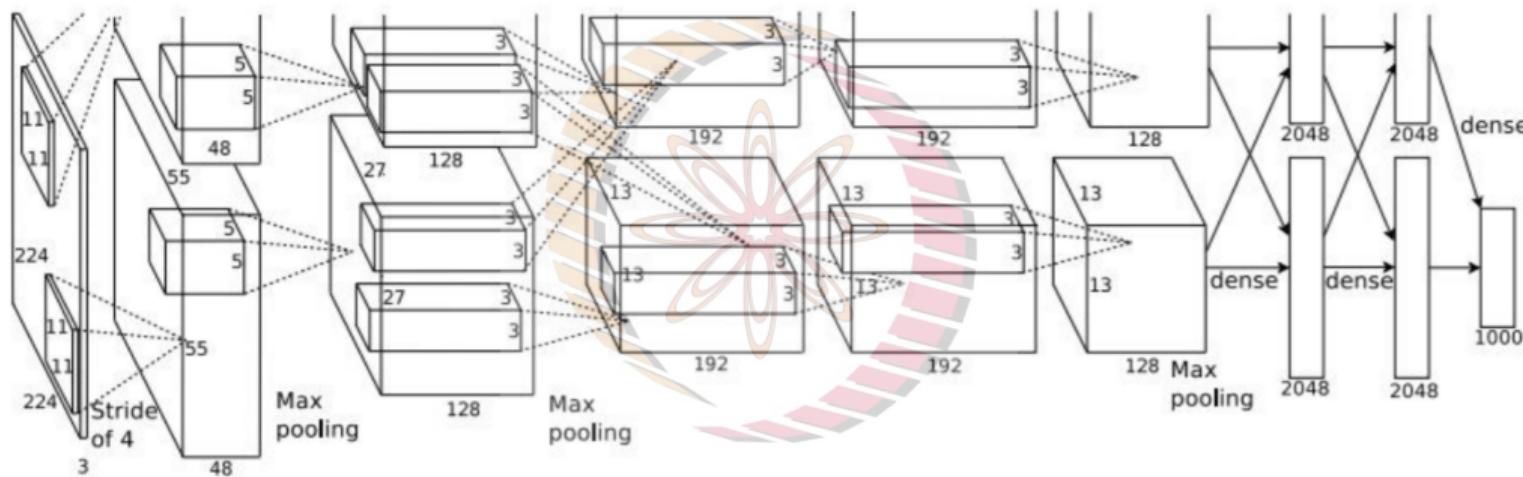
- 8 layers in total (5 convolutional layers, 3 fully connected layers)
- Trained on ImageNet Dataset
- **Response normalization layers** follow the first and second convolutional layers.
- Max-pooling follow first, second and the fifth convolutional layers
- The ReLU non-linearity is applied to the output of every layer



NPTEL



AlexNet



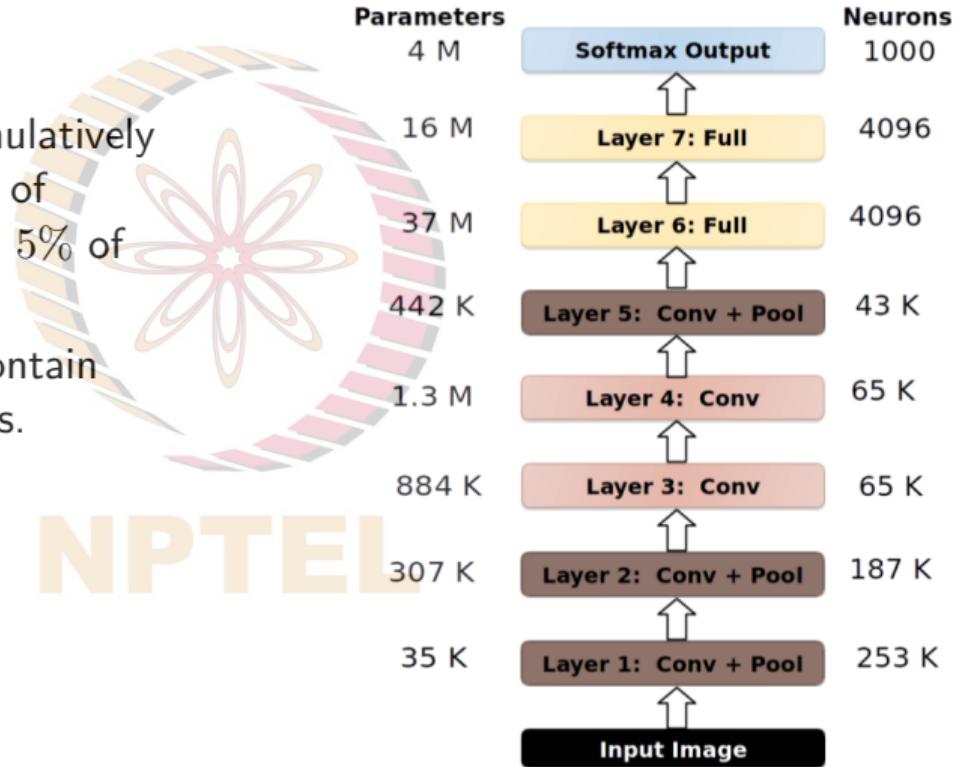
About 57 M parameters are in the fully connected layers

Total

Parameters :	$[(11 \times 11 \times 3) + 1] \times 96 = 35 \text{ K}$	$[5 \times 5 \times 48] \times 256 = 307 \text{ K}$	$[3 \times 3 \times 256] \times 384 = 884 \text{ K}$	663 K	442 K	37 M	16 M	4 M	60 M
Neurons :	253,440	$27 \times 27 \times 256 = 186,624$	$13 \times 13 \times 384 = 64,896$	$13 \times 13 \times 384 = 64,896$	$13 \times 13 \times 256 = 43,264$	4096	4096	1000	0.63 M

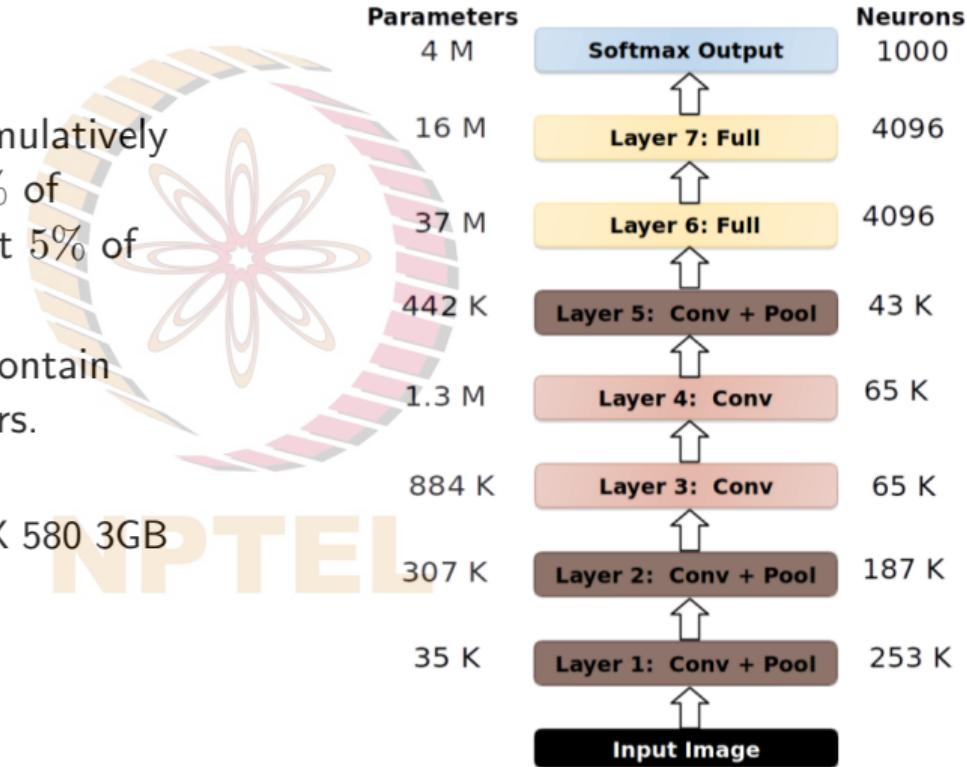
AlexNet

- Convolutional layers cumulatively contain about 90 – 95% of computation, only about 5% of the parameters
- Fully-connected layers contain about 95% of parameters.

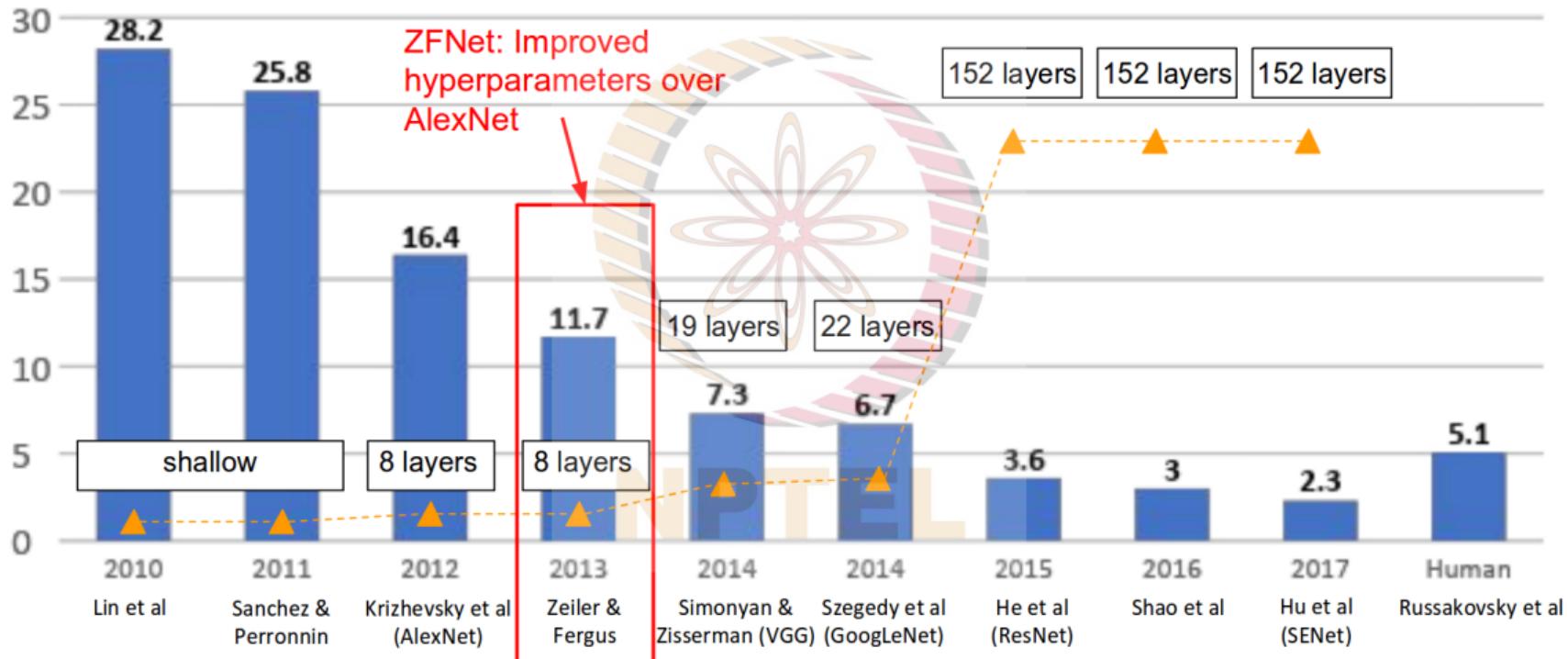


AlexNet

- Convolutional layers cumulatively contain about 90 – 95% of computation, only about 5% of the parameters
- Fully-connected layers contain about 95% of parameters.
- Trained with SGD
 - on **two** NVIDIA GTX 580 3GB GPUs
 - for about a week

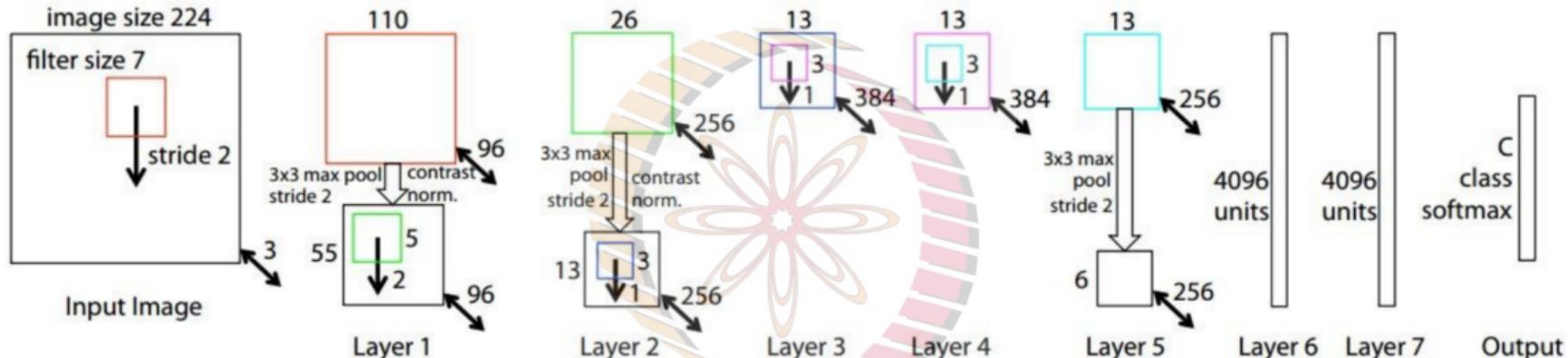


Winners of ImageNet Classification Challenge



Credit: Fei-Fei Li, Justin Johnson and Serena Yeung, CS231n course, Stanford, Spring 2019

ZFNet (2013)²

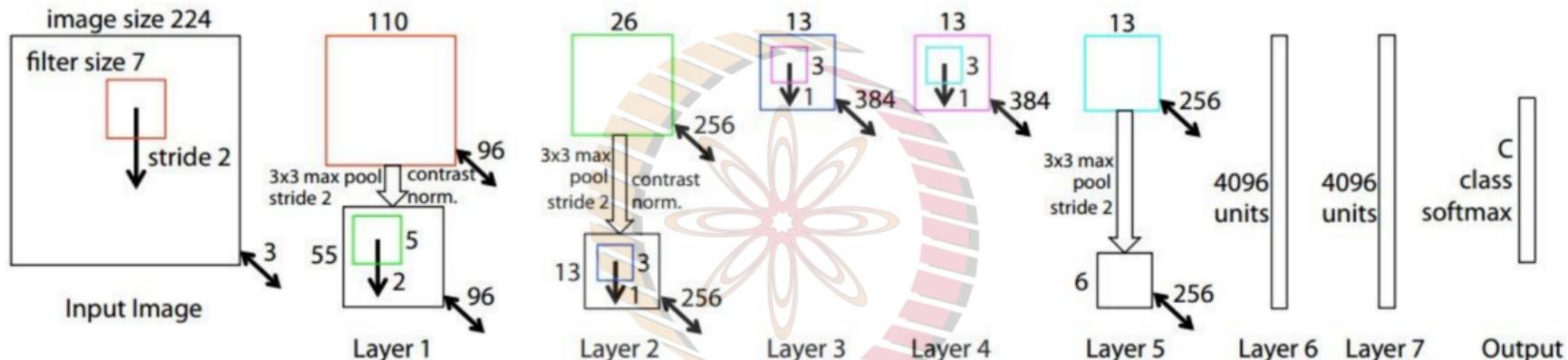


- Similar to AlexNet but:

NPTEL

²Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

ZFNet (2013)²

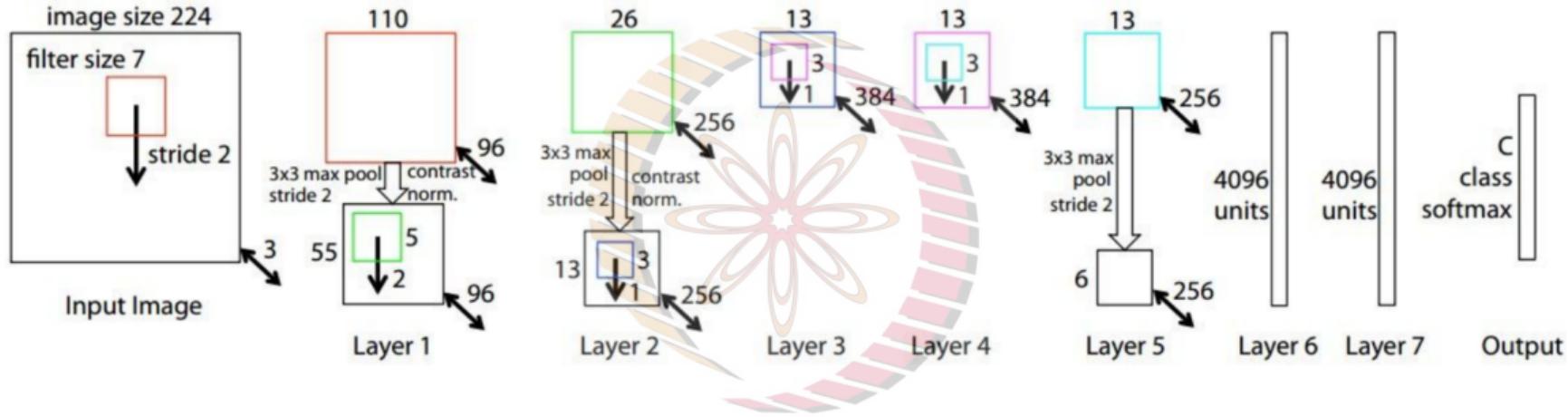


- Similar to AlexNet but:

NPTEL

²Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

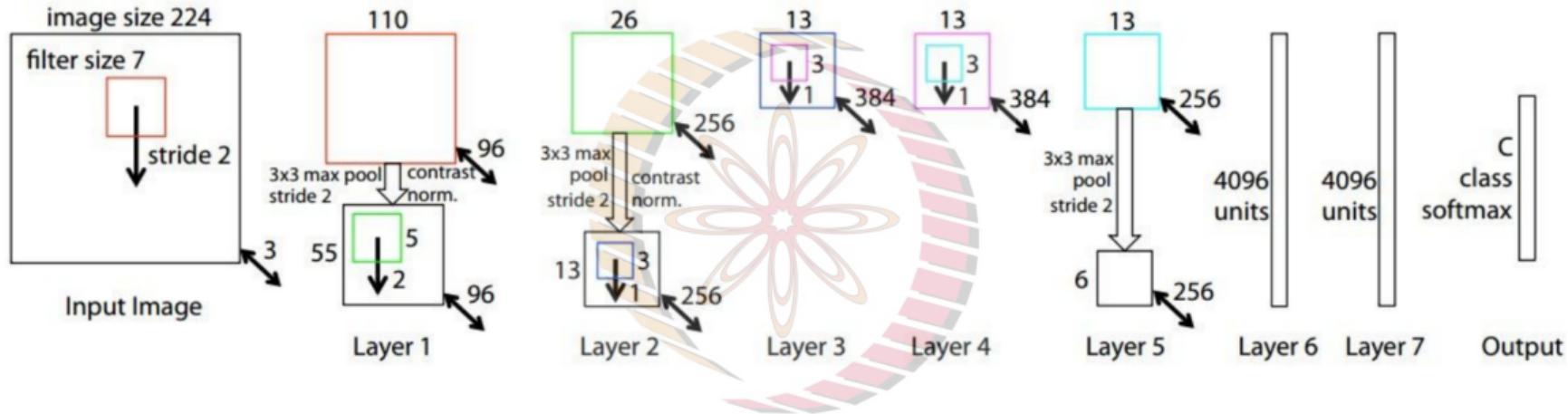
ZFNet (2013)²



- Similar to AlexNet but:
 - CONV1: change from $(11 \times 11 \text{ stride } 4)$ to $(7 \times 7 \text{ stride } 2)$

²Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

ZFNet (2013)²

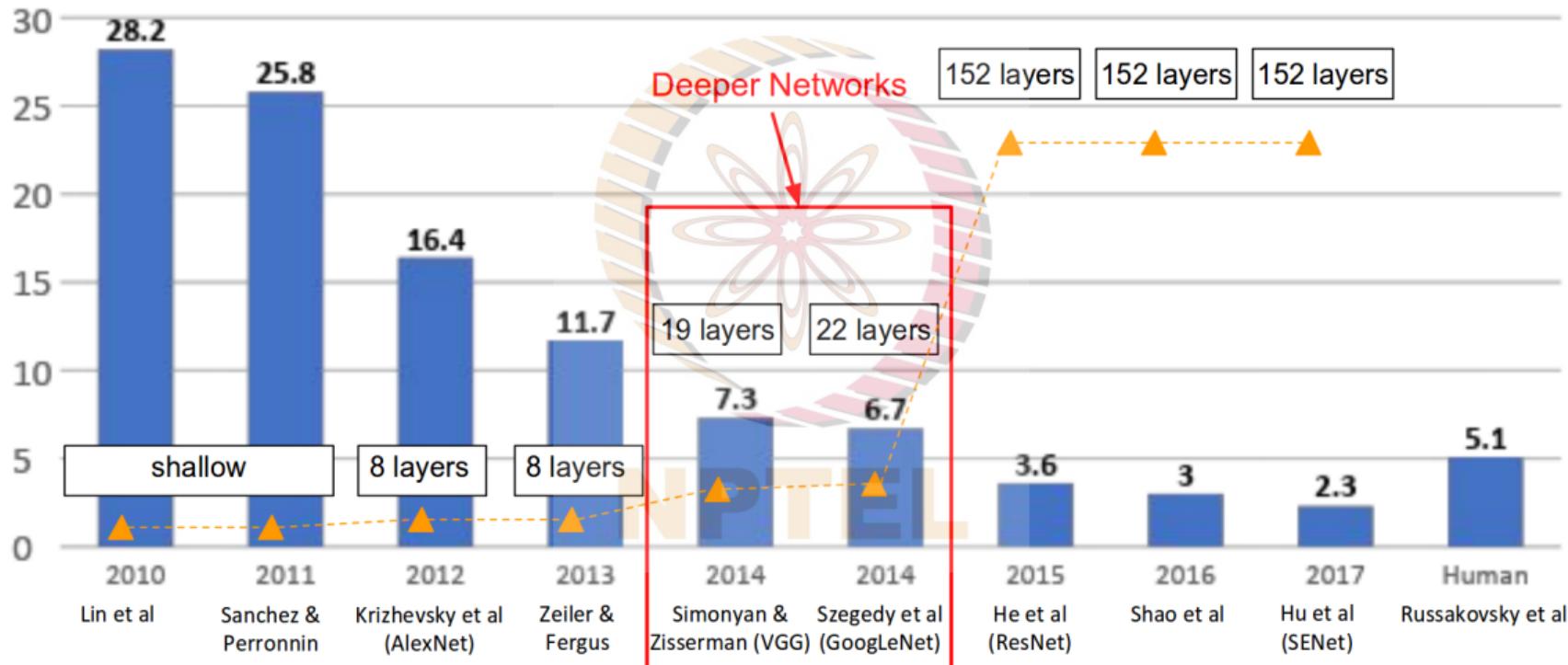


- Similar to AlexNet but:
 - CONV1: change from $(11 \times 11 \text{ stride } 4)$ to $(7 \times 7 \text{ stride } 2)$
 - CONV3,4,5: instead of 384, 384, 256 filters, use 512, 1024, 512
- ImageNet top-5 error: $16.4\% \rightarrow 11.7\%$

Credit: Fei-Fei Li, Justin Johnson and Serena Yeung, CS231n course, Stanford, Spring 2019

²Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Winners of ImageNet Classification Challenge



Credit: Fei-Fei Li, Justin Johnson and Serena Yeung, CS231n course, Stanford, Spring 2019

VGGNet

image

conv-64

maxpool

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

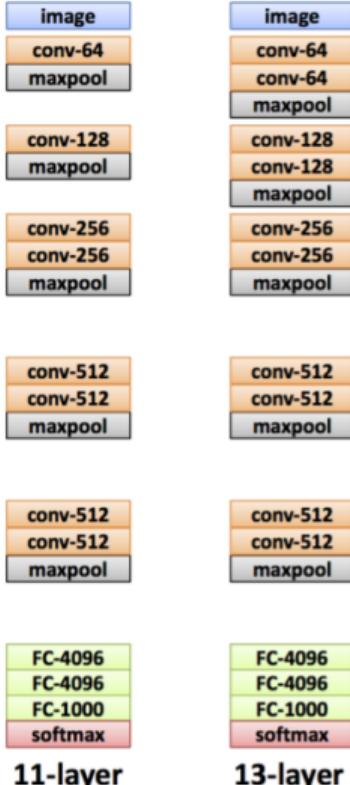
softmax

11-layer



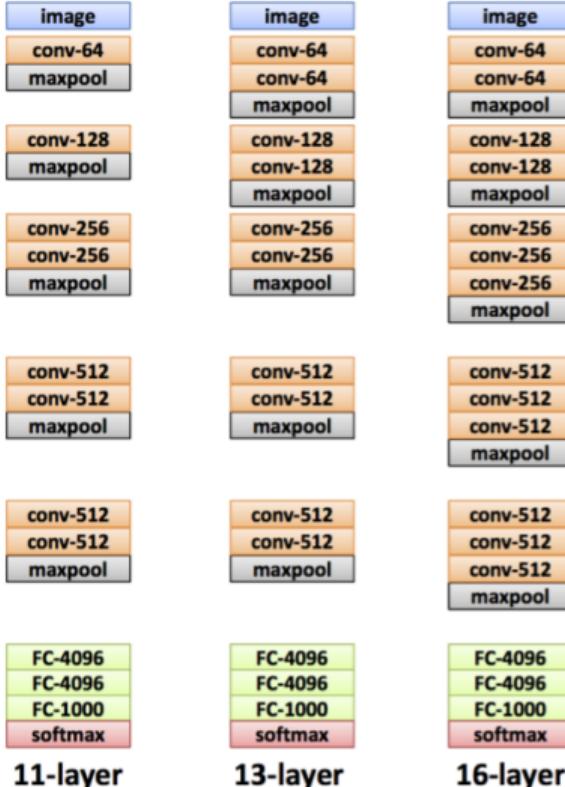
NPTEL

VGGNet



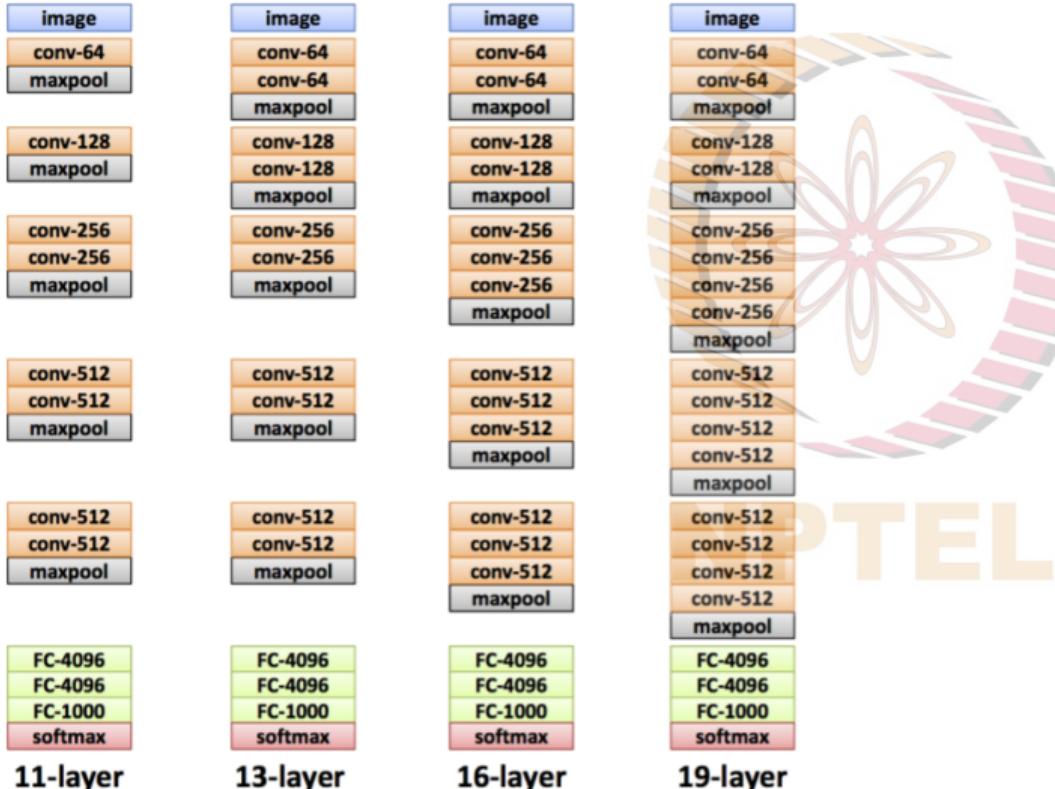
NPTEL

VGGNet

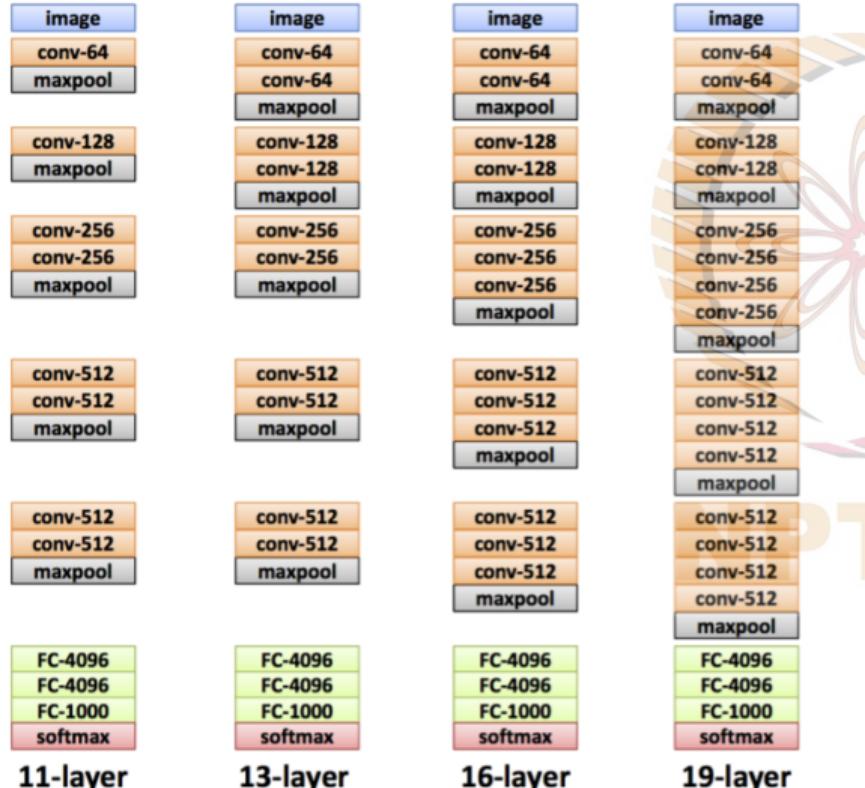


NPTEL

VGGNet

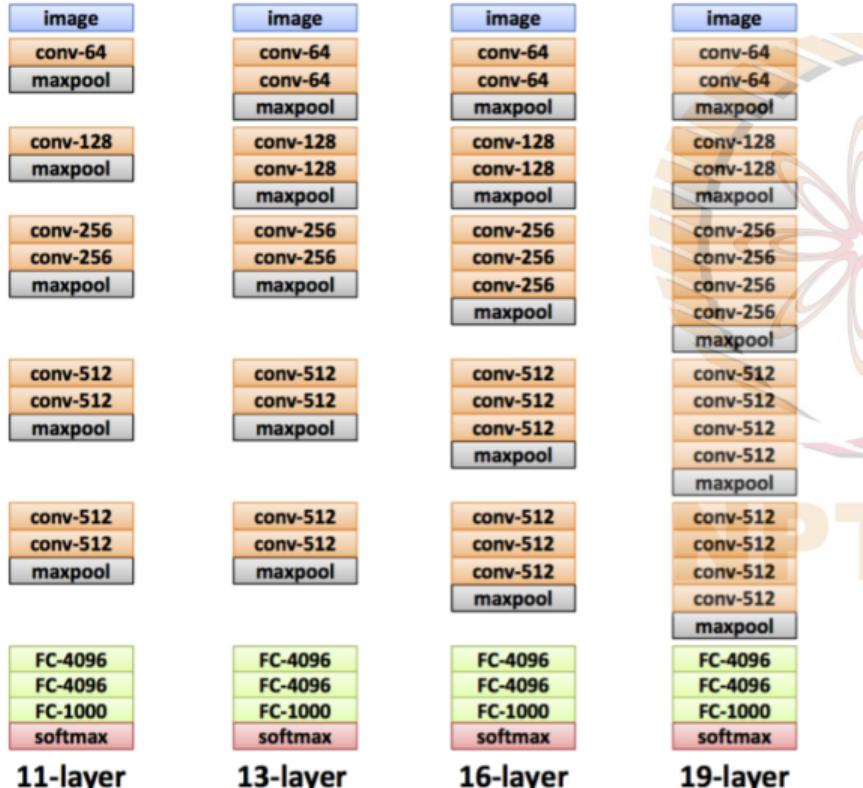


VGGNet



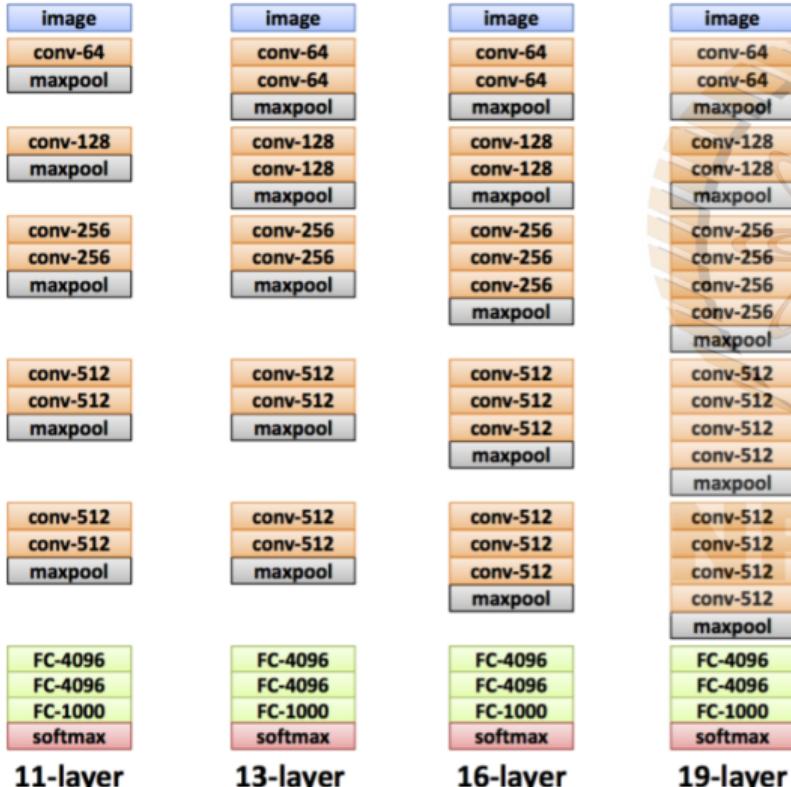
Runner-up in ILSVRC 2014

VGGNet



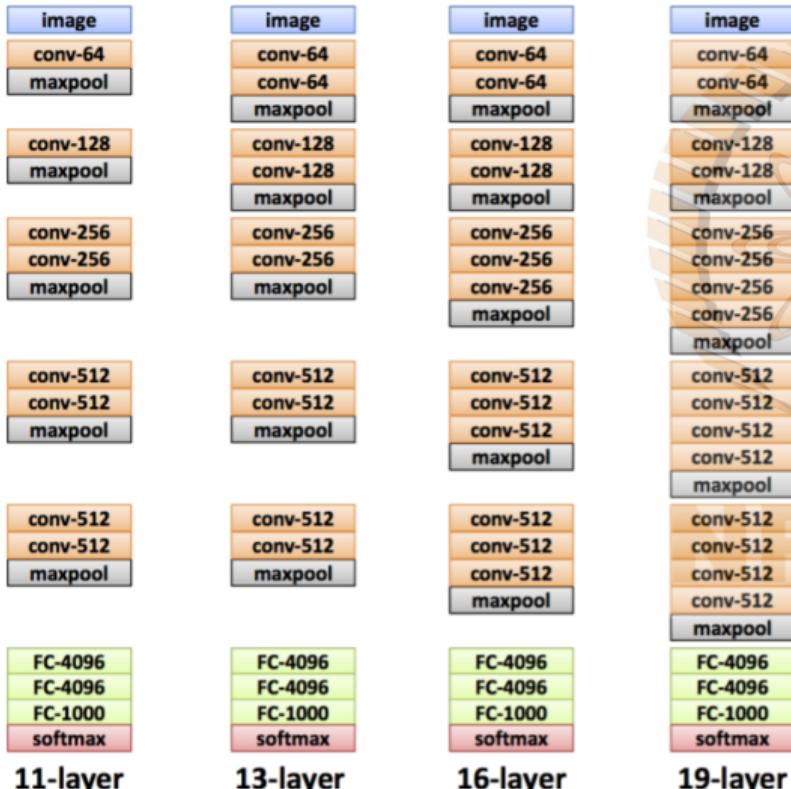
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities

VGGNet



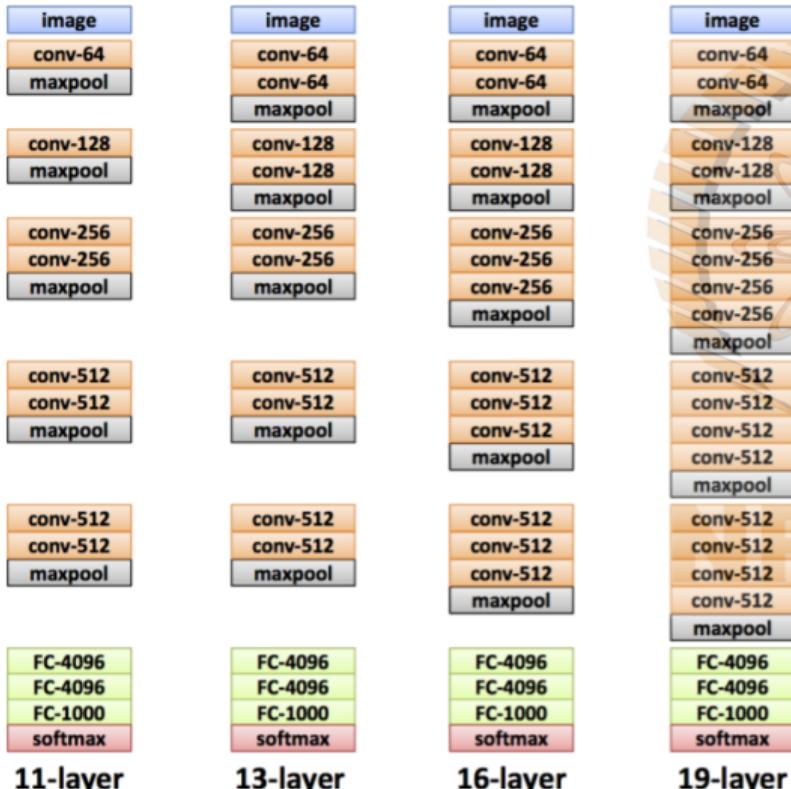
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance

VGGNet



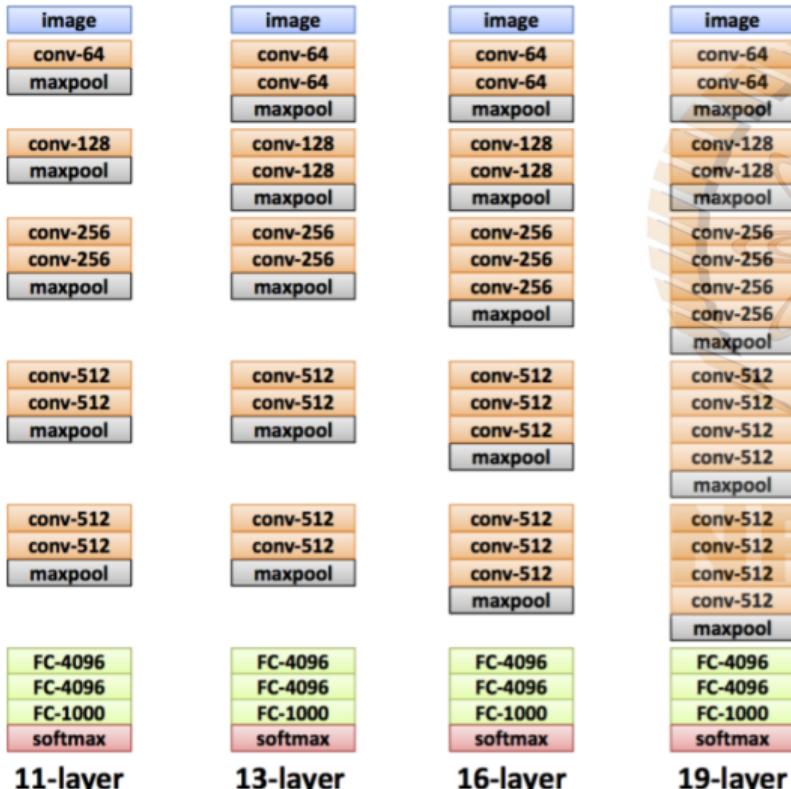
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:

VGGNet



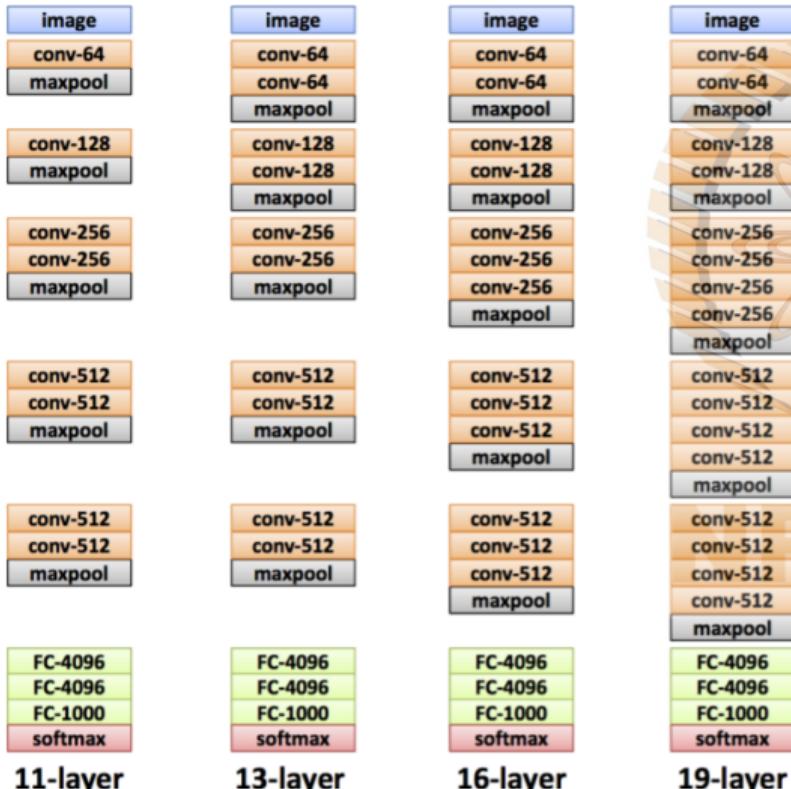
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 × 3 CONV stride 1 pad 1

VGGNet



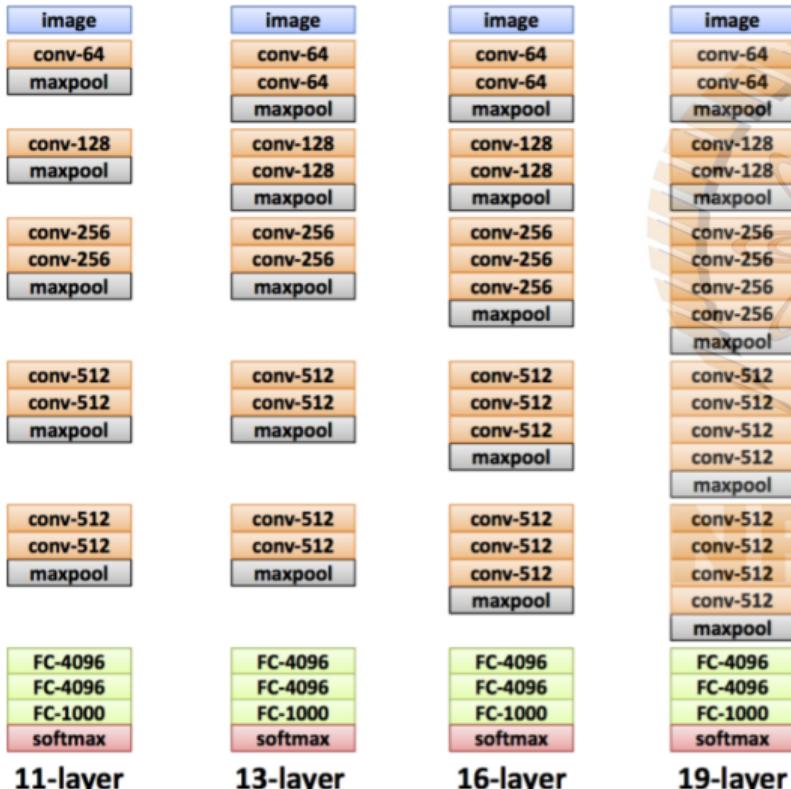
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 x 3 CONV stride 1 pad 1
 - 2 x 2 MAX POOL stride 2

VGGNet



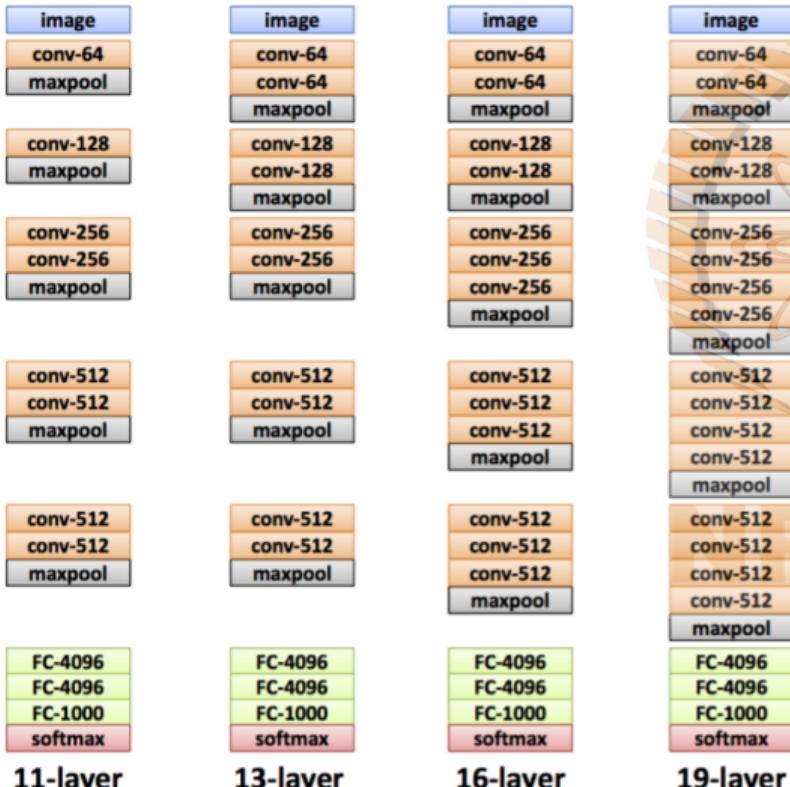
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 x 3 CONV stride 1 pad 1
 - 2 x 2 MAX POOL stride 2
- **Smaller receptive fields:**

VGGNet



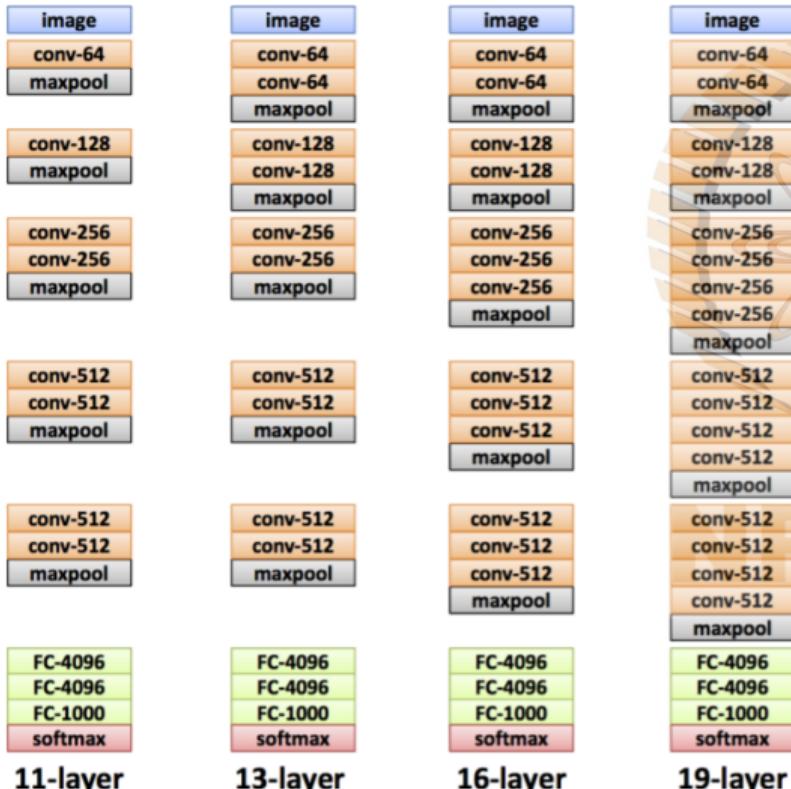
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 x 3 CONV stride 1 pad 1
 - 2 x 2 MAX POOL stride 2
- **Smaller receptive fields:**
 - less parameters; faster

VGGNet



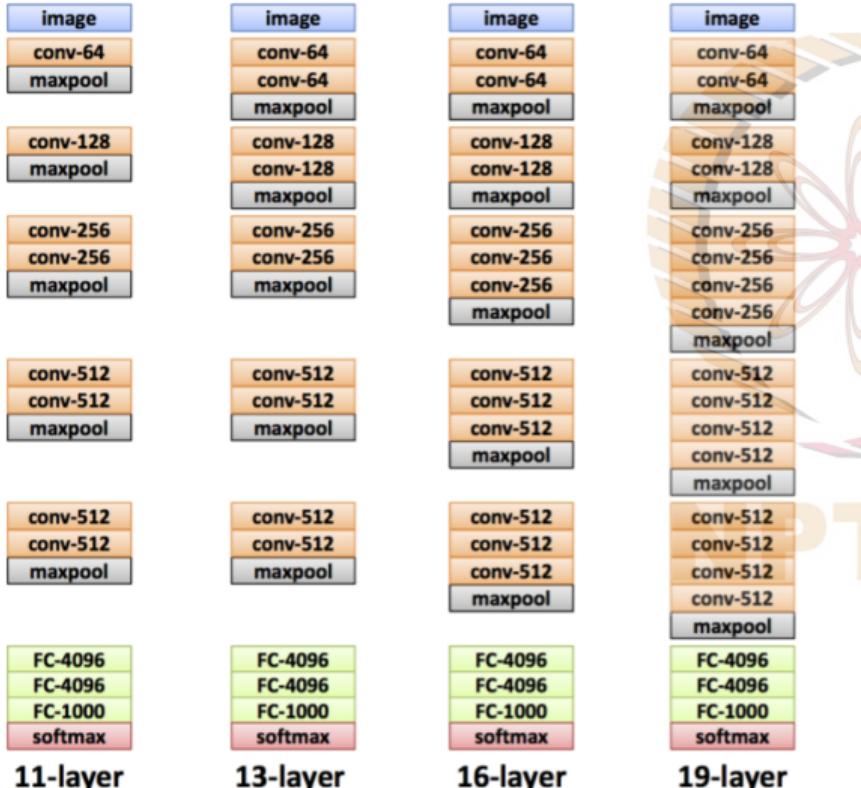
- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 x 3 CONV stride 1 pad 1
 - 2 x 2 MAX POOL stride 2
- **Smaller receptive fields:**
 - less parameters; faster
 - two 3 x 3 conv has same receptive field as a single 5 x 5 conv; three 3 x 3 conv has same receptive field as a single 7 x 7 conv

VGGNet



- Runner-up in ILSVRC 2014
- More layers lead to more nonlinearities
- **Key contribution:** Depth of the network is a critical component for good performance
- **Homogeneous Architecture:** From beginning to end:
 - 3 x 3 CONV stride 1 pad 1
 - 2 x 2 MAX POOL stride 2
- **Smaller receptive fields:**
 - less parameters; faster
 - two 3 x 3 conv has same receptive field as a single 5 x 5 conv; three 3 x 3 conv has same receptive field as a single 7 x 7 conv
- Fewer parameters: $3 \times 3^2 C^2$ (vs) $7^2 C^2$

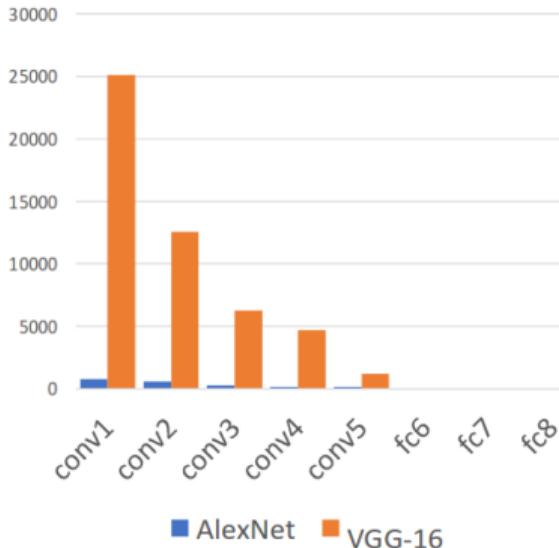
VGGNet



- VGG19 only slightly better than VGG16
- Used ensembles of networks for best results

VGGNet

AlexNet vs VGG-16
(Memory, KB)

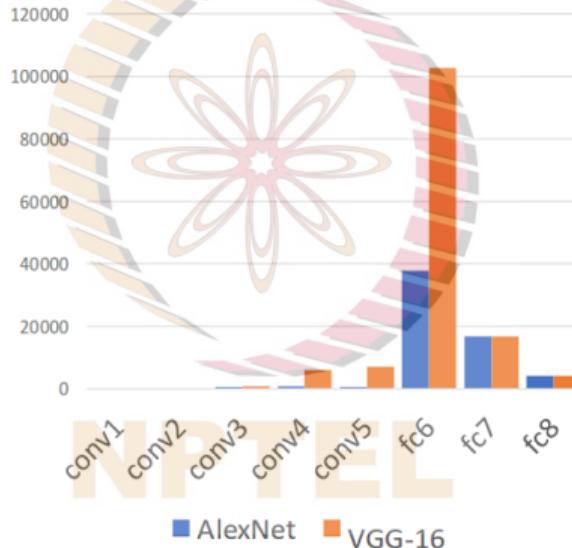


AlexNet total: 1.9 MB

VGG-16 total: 48.6 MB (25x)

Credit: Justin Johnson, Univ of Michigan

AlexNet vs VGG-16
(Params, M)

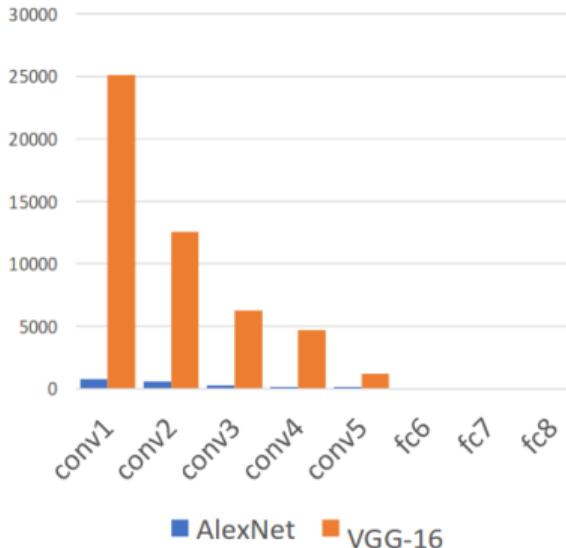


AlexNet total: 61M

VGG-16 total: 138M (2.3x)

VGGNet

AlexNet vs VGG-16
(Memory, KB)

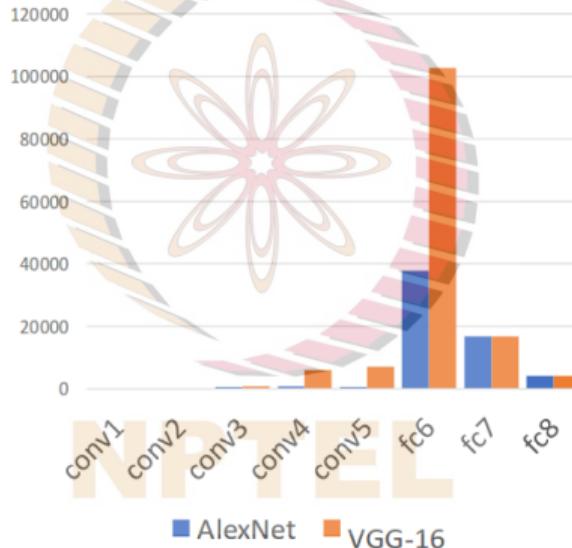


AlexNet total: 1.9 MB

VGG-16 total: 48.6 MB (25x)

Credit: Justin Johnson, Univ of Michigan

AlexNet vs VGG-16
(Params, M)



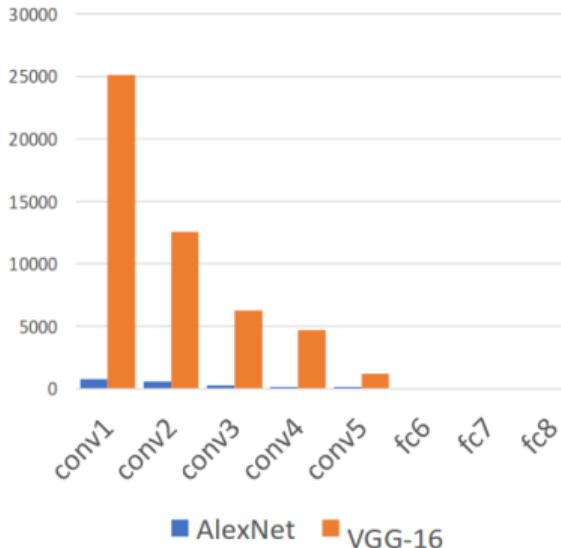
AlexNet total: 61M

VGG-16 total: 138M (2.3x)

- Uses a lot more memory and parameters

VGGNet

AlexNet vs VGG-16
(Memory, KB)

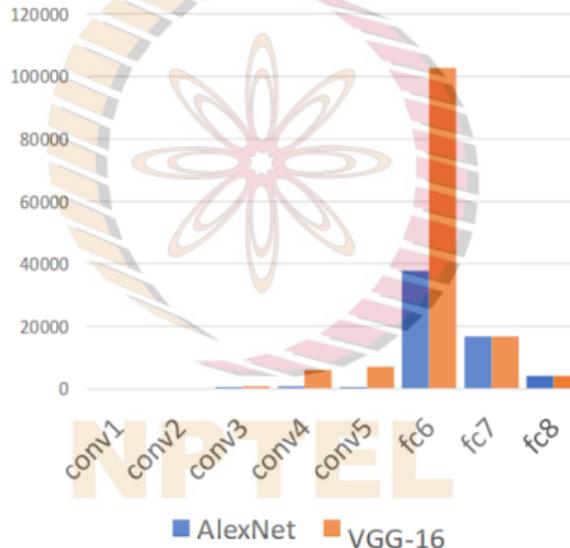


AlexNet total: 1.9 MB

VGG-16 total: 48.6 MB (25x)

Credit: Justin Johnson, Univ of Michigan

AlexNet vs VGG-16
(Params, M)



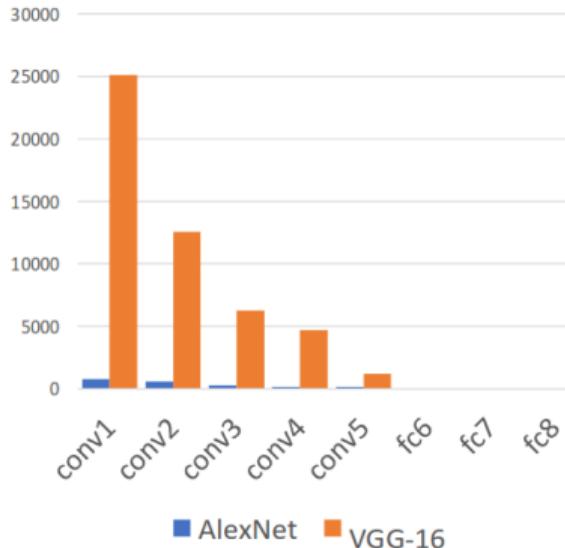
AlexNet total: 61M

VGG-16 total: 138M (2.3x)

- Uses a lot more memory and parameters
- Most of these parameters are in the first fully connected layer

VGGNet

AlexNet vs VGG-16
(Memory, KB)

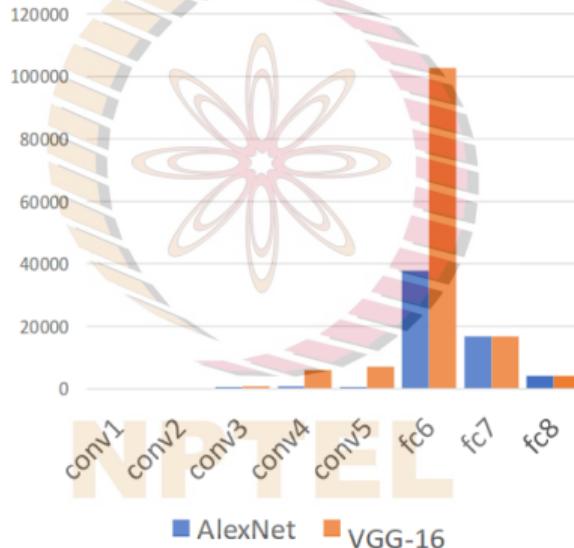


AlexNet total: 1.9 MB

VGG-16 total: 48.6 MB (25x)

Credit: Justin Johnson, Univ of Michigan

AlexNet vs VGG-16
(Params, M)



AlexNet total: 61M

VGG-16 total: 138M (2.3x)

- Uses a lot more memory and parameters
- Most of these parameters are in the first fully connected layer
- Most of the memory is used in early CONV layer

Homework

Readings

- Tutorial: [Illustrated: 10 CNN Architectures](#)
 - Read the first 3: LeNet, AlexNet, VGG
- (Optional) For more details, skim through the following papers:
 - [ImageNet Classification with Deep Convolutional Neural Networks](#)
 - [Very Deep Convolutional Networks for Large-Scale Image Recognition](#)



References

- 
- [1] Yann LeCun et al. "Gradient-based learning applied to document recognition". In: 1998.
 - [2] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. "ImageNet Classification with Deep Convolutional Neural Networks". In: *NIPS*. 2012.
 - [3] Karen Simonyan and Andrew Zisserman. "Very Deep Convolutional Networks for Large-Scale Image Recognition". In: *CoRR* abs/1409.1556 (2015).
 - [4] Christian Szegedy et al. "Going deeper with convolutions". In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 1–9.
 - [5] Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.
 - [6] Johnson, Justin, EECS 498-007 / 598-005 - Deep Learning for Computer Vision (Fall 2019). URL: <https://web.eecs.umich.edu/~justincj/teaching/eecs498/> (visited on 06/29/2020).
 - [7] Li, Fei-Fei; Johnson, Justin; Serena, Yeung; CS 231n - Convolutional Neural Networks for Visual Recognition (Spring 2019). URL: <http://cs231n.stanford.edu/2019/> (visited on 06/29/2020).