# Lecture #26 Convolutional Neural Networks

## Gian Antonio Susto

Reference Material (used for this presentation):
- MIT *Introduction to Deep Learning* http://introtodeeplearning.com
- L. Fridman MIT Deep Learning https://deeplearning.mit.edu/
- I. Goodfellow, Y. Bengio, A. Courville. *Deep learning*. MIT press, 2016.
- ImageNET http://www.image-net.org/
- C. Szegedy et al. *Going Deeper with Convolutions* CVPR2015
- F. Li, J. Johnson, S. Yeung. *Stanford CS231n Convolutional Neural Networks for Visual Recognition* *http://cs231n.stanford.edu/* + *http://cs231n.github.io/convolutional-networks/*
- Y. LeCun http://yann.lecun.com/exdb/lenet/
- https://www.kaggle.com/zalando-research/fashionmnist
- MIT https://groups.csail.mit.edu/vision/TinyImages/
- ImageNET http://www.image-net.org/
- https://github.com/GoogleCloudPlatform/tensorflow-without-a-phd
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. *Imagenet classification with deep convolutional neural networks*. Advances in neural information processing systems. 2012.
- C. Szegedy et al. *Going Deeper with Convolutions* CVPR2015
- Lee, Honglak, et al. "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations." Proceedings of the 26th annual international conference on machine learning. ACM, 2009
- Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014
- Lee, Honglak, et al. "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations." Proceedings of the 26th annual international conference on machine learning. ACM, 2009
- Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014
- The complete history of Lenna http://www.ee.cityu.edu.hk/~lmpo/lenna/Lenna97.html
- Style transfer https://deepart.io/
- https://paperswithcode.com/sota/image-classification-on-imagenet

# From lecture #02:

Different tasks (objectives), different models, different data type... and different stages of development!

# With DL architectures we are able to handle all sort of data



A mostly complete chart of
## Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

**Legend:**
- Input Cell
- Backfed Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Gated Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P)    Feed Forward (FF)    Radial Basis Network (RBF)    Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)    Long / Short Term Memory (LSTM)    Gated Recurrent Unit (GRU)

Auto Encoder (AE)    Variational AE (VAE)    Denoising AE (DAE)    Sparse AE (SAE)

Markov Chain (MC)    Hopfield Network (HN)    Boltzmann Machine (BM)    Restricted BM (RBM)    Deep Belief Network (DBN)

Deep Convolutional Network (DCN)    Deconvolutional Network (DN)    Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)    Liquid State Machine (LSM)    Extreme Learning Machine (ELM)    Echo State Network (ESN)

Deep Residual Network (DRN)    Differentiable Neural Computer (DNC)    Neural Turing Machine (NTM)

Capsule Network (CN)    Kohonen Network (KN)    Attention Network (AN)

# Today: Convolutional Neural Networks (CNNs) to deal with images/videos

Today: Convolutional Neural Networks (CNNs) to deal with images/videos

We will not see in this course Recurrent Neural Networks (RNNs) to deal with sequence learning tasks

Outline:

- CNN building blocks

- Tasks in Computer Vision

- Datasets for Computer Vision & Historical CNNs

- What a CNN sees?

# Human Vision



A system trained on 540 million years of data

# Computer Vision



Started in the 60s aiming at automatically doing what human vision can do

CV is undoubtly THE field that mostly benefit from Deep Learning

# DL has been disruptive in computer vision!



A dataset with more than 14 Million annotated images of 20K categories

# Example Task #1: Classification Chihuahua vs Muffin

# Example Task #2: Semantic Segmentation



**Input:** 3 x H x W

**High-res:** D₁ x H/2 x W/2

**Med-res:** D₂ x H/4 x W/4

**Low-res:** D₃ x H/4 x W/4

**Med-res:** D₂ x H/4 x W/4

**High-res:** D₁ x H/2 x W/2

**Predictions:** H x W

All convolutional layers with downsampling and upsampling operations

# Example Task #2: Semantic Segmentation

# Example Task #3: Object Detection



1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

1. Find regions that we think have objects

2. Classify such objects

# Example Task #4: Image Captioning

Combination of CNN and sequence learning

# Example Task #4: Image Captioning

Combination of CNN and sequence learning





-MIT *Introduction to Deep Learning* http://introtodeeplearning.com

# Example Task #5: Style Transfer



Content Representations

Convolutional Neural Network

# Example Task #5: Style Transfer



-https://deepart.io/ (not working anymore)

# Example Task #5: Style Transfer

# Example Task #5: Style Transfer

# Example Task #5: Style Transfer

Some style transfer works
- Gatys, L. A., Ecker, A. S., & Bethge, M. (2016). Image style transfer using convolutional neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 2414-2423).
- Johnson, J., Alahi, A., & Fei-Fei, L. (2016, October). Perceptual losses for real-time style transfer and super-resolution. In European conference on computer vision (pp. 694-711). Springer, Cham.
- Luan, F., Paris, S., Shechtman, E., & Bala, K. (2017). Deep photo style transfer. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 4990-4998).
- i, X., Liu, S., Kautz, J., & Yang, M. H. (2019). Learning linear transformations for fast image and video style transfer. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 3809-3817)

# Example Task #6: Image Generation

# Example Task #7: Visual Anomaly Detection

# Images are numbers



What the computer sees

L. Fridman MIT Deep Learning https://deeplearning.mit.edu/

JPG 260 X 194

260 X 194 X 3

8,11,0, 55,13,25,19
15,241,2,155,13,35,65
14,211,0,255,23,45,11
05,255,1,255,10,17,23
77,167,9,112,56,16,90
45,245,0,145,22,55,48

# Images are numbers



L. Fridman MIT Deep Learning https://deeplearning.mit.edu/

Any (simple) ideas on how to feed this into a model?

# Using FFNN is not a good idea...

- 2D images converted to arrays: **spatial information is lost!**

- Each hidden unit is connected to all units in the previous layers: **lots of parameters**!

# Using FFNN is not a good

- 2D images converted to array

Does the value in this specific pixel really matters by itself?



What the computer sees

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|157|153|174|168|150|152|129|151|172|161|155|156|
|155|182|163|74|75|62|33|17|110|210|180|154|
|180|180|50|14|34|6|10|33|48|106|159|181|
|206|109|5|124|131|111|120|204|166|15|56|180|
|194|68|137|251|237|239|239|228|227|87|71|201|
|172|105|207|233|233|214|220|239|228|98|74|206|
|188|88|179|209|185|215|211|158|139|75|20|169|
|189|97|165|84|10|168|134|11|31|62|22|148|
|199|168|191|193|158|227|178|143|182|106|36|190|
|205|174|155|252|236|231|149|178|228|43|95|234|
|190|216|116|149|236|187|86|150|79|38|218|241|
|190|224|147|108|227|210|127|102|36|101|255|224|
|190|214|173|66|103|143|96|50|2|109|249|215|
|187|196|235|75|1|81|47|0|6|217|255|211|
|183|202|237|145|0|0|12|108|200|138|243|236|
|195|206|123|207|177|121|123|200|175|13|96|218|

Input layer $L_1$   Hidden layer $L_2$   Output layer $L_3$

| 1 | 2 |
|---|---|
| 3 | 4 |

→

| 1 |
|---|
| 2 |
| 3 |
| 4 |

→

$x_1 \rightarrow$
$x_2 \rightarrow$
$x_3 \rightarrow$
$x_4 \rightarrow$

$\rightarrow f(x)$

# In Computer Vision, feature engineering is hard

# As seen before: no feature engineering is needed with DL



L. Fridman MIT Deep Learning https://deeplearning.mit.edu/

I. Goodfellow, Y. Bengio, A. Courville. Deep learning. MIT press, 2016.

# Features extracted by Deep Neural Networks



Low Level Features — Lines & Edges

Mid Level Features — Eyes & Nose & Ears

High Level Features — Facial Structure

These features were obtained in an automatic fashion!

Lee, Honglak, et al. "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations." *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009.

# Idea: using the spatial structure (patches)

- Connect patch in input layer to a single neuron in subsequent layer.

- Use a sliding window to define connections.

# Idea: 2D convolutions



| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

We would like to extract something useful out of this patch! Why not doing a 'linear combination' as we did with tabular data?

# Idea: 2D convolutions



| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |
| 7 | 8 | 9 |

| 1 | 2 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -2 | -1 |

The operation is the 2D convolution: we multiply element-wise 2 matrices (the input and the kernel – weights), then we sum and derive some 'features'

Input Image — Applying padding of 1 on 3x3 — Padded Image

$$y[0,0] = \sum_j \sum_i x[i,j] \cdot h[0-i, 0-j]$$

$$= \quad x[-1,-1] \cdot h[1,1] + x[0,-1] \cdot h[0,1] + x[1,-1] \cdot h[-1,1]$$

$$+ \, x[-1,0] \cdot h[1,0] \quad + x[0,0] \cdot h[0,0] \quad + x[1,0] \cdot h[-1,0]$$

$$+ \, x[-1,1] \cdot h[1,-1] + x[0,1] \cdot h[0,-1] + x[1,1] \cdot h[-1,-1]$$

$$= \quad 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1$$

$$+ \, 0 \cdot 0 + 1 \cdot 0 + 2 \cdot 0$$

$$+ \, 0 \cdot (-1) + 4 \cdot (-2) + 5 \cdot (-1)$$

$$= -13$$

From https://www.songho.ca/dsp/convolution/convolution2d_example.html

$$y[1,0] = \sum_j \sum_i x[i,j] \cdot h[1-i, 0-j]$$

$$= \quad x[0,-1] \cdot h[1,1] + x[1,-1] \cdot h[0,1] + x[2,-1] \cdot h[-1,1]$$

$$+ x[0,0] \cdot h[1,0] \quad + x[1,0] \cdot h[0,0] \quad + x[2,0] \cdot h[-1,0]$$

$$+ x[0,1] \cdot h[1,-1] + x[1,1] \cdot h[0,-1] + x[2,1] \cdot h[-1,-1]$$

$$= \quad 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1$$

$$+ 1 \cdot 0 + 2 \cdot 0 + 3 \cdot 0$$

$$+ 4 \cdot (-1) + 5 \cdot (-2) + 6 \cdot (-1)$$

$$= -20$$

From https://www.songho.ca/dsp/convolution/convolution2d_example.html

$$y[2,0] = \sum_j \sum_i x[i,j] \cdot h[2-i, 0-j]$$

$$= \quad x[1,-1] \cdot h[1,1] + x[2,-1] \cdot h[0,1] + x[3,-1] \cdot h[-1,1]$$

$$+ x[1,0] \cdot h[1,0] \quad + x[2,0] \cdot h[0,0] \quad + x[3,0] \cdot h[-1,0]$$

$$+ x[1,1] \cdot h[1,-1] + x[2,1] \cdot h[0,-1] + x[3,1] \cdot h[-1,-1]$$

$$= \quad 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1$$

$$+ 2 \cdot 0 + 3 \cdot 0 + 0 \cdot 0$$

$$+ 5 \cdot (-1) + 6 \cdot (-2) + 0 \cdot (-1)$$

$$= -17$$

From https://www.songho.ca/dsp/convolution/convolution2d_example.html

We then 'stride' vertically!

$$y[0,1] = \sum_j \sum_i x[i,j] \cdot h[0-i, 1-j]$$

$$= \quad x[-1,0] \cdot h[1,1] \quad + x[0,0] \cdot h[0,1] \quad + x[1,0] \cdot h[-1,1]$$
$$+ x[-1,1] \cdot h[1,0] \quad + x[0,1] \cdot h[0,0] \quad + x[1,1] \cdot h[-1,0]$$
$$+ x[-1,2] \cdot h[1,-1] + x[0,2] \cdot h[0,-1] + x[1,2] \cdot h[-1,-1]$$

$$= \quad 0 \cdot 1 + 1 \cdot 2 + 2 \cdot 1$$
$$+ 0 \cdot 0 + 4 \cdot 0 + 5 \cdot 0$$
$$+ 0 \cdot (-1) + 7 \cdot (-2) + 8 \cdot (-1)$$

$$= -18$$

After completing the stride:

| Input | | |
|---|---|---|
| 1 | 2 | 3 |
| 4 | 5 | 6 |
| 7 | 8 | 9 |

Kernel

| | m | | |
| n | -1 | 0 | 1 |
| -1 | -1 | -2 | -1 |
| 0 | 0 | 0 | 0 |
| 1 | 1 | 2 | 1 |

| Output | | |
|---|---|---|
| -13 | -20 | -17 |
| -18 | -24 | -18 |
| 13 | 20 | 17 |

From https://www.songho.ca/dsp/convolution/convolution2d_example.html

# Idea: 2D convolutions

# Idea: 2D convolutions



The filter can be of different size! And we can have multiple filters!

# Idea: 2D convolutions

# Idea: 2D convolutions



Padding = 1
Stride = 2

# Idea: 2D convolutions



With padding we give 'right' importance to pixels in the border, while with stride we can tune the inner dimensions of the network

Padding = 1
Stride = 2

# Feature maps induced by convolutions: some examples



Original

Sharpen

Edge Detect

"Strong" Edge Detect

-MIT *Introduction to Deep Learning* http://introtodeeplearning.com

# Feature maps induced by convolutions: some examples



Original

Sharpen

| 0 | -1 | 0 |
|---|----|---|
| -1 | 5 | -1 |
| 0 | -1 | 0 |

Edge Detect

| 0 | 1 | 0 |
|---|---|---|
| 1 | -4 | 1 |
| 0 | 1 | 0 |

"Strong" Edge Detect

| -1 | -2 | -1 |
|----|----|----|
| 0 | 0 | 0 |
| 1 | 2 | 1 |

| ? | ? | ? |
|---|---|---|
| ? | ? | ? |
| ? | ? | ? |

# Recap



32x32x3 image

5x5x3 filter

32

32

3

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

F. Li, J. Johnson, S. Yeung. *Stanford CS231n Convolutional Neural Networks for Visual Recognition* http://cs231n.stanford.edu/

# Recap



32x32x3 image
5x5x3 filter

convolve (slide) over all spatial locations

activation map

stride = 1, no padding

# Recap



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

**activation map**
**activation maps**

28

28

1

stride = 1,
no padding

F. Li, J. Johnson, S. Yeung. *Stanford CS231n Convolutional Neural Networks for Visual Recognition* http://cs231n.stanford.edu/

# Recap



activation maps

32
32
3

Convolution Layer

28
28
6

F. Li, J. Johnson, S. Yeung. *Stanford CS231n Convolutional Neural Networks for Visual Recognition* http://cs231n.stanford.edu/

# Receptive field



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

F. Li, J. Johnson, S. Yeung. *Stanford CS231n Convolutional Neural Networks for Visual Recognition* http://cs231n.stanford.edu/

# Pros of 2D convolutions

- Use convolution filters (weights) to extract **local features**

- Use multiple filters to extract **different features**

- Filter parameters are **shared spatially** across the whole input



**Layer Dimensions:**

$$h \ x \ w \ x \ d$$

where h and w are spatial dimensions
d (depth) = number of filters

**Stride:**
Filter step size

**Receptive Field:**
Locations in input image that
a node is path connected to

# Deep Learning: we stack different layers!

# Make extracted features rich, robust and low-dimensional

- Introduce non-linearity: apply non-linear activation to convolution output

- Downsample and preserve spatial invariance: pooling



max pool with 2x2 filters and stride 2

1) Reduced dimensionality
2) Spatial invariance

Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$

# Dimensionality reduction (optional)

- While to reduce width and height we can exploit pooling (or bigger strides), to reduce depth we can use 1d convolutions

## 1d convolutions

Why 1x1 convolutions?

- **Efficiency**: reduces the depth (number of channels). Width and height are unchanged. To reduce the horizontal dimensions, you would use pooling (or increase the stride of the conv).

- The 1x1 conv computes a weighted sum of input channels (or features). This allows it to "**select**" certain combinations of features that are useful downstream.

10

Weights:

1x1x10

1

# Now put it all together to perform <mark>classification</mark>

- We apply back-propagation to train model (weights for convolution and dense layers)

- Cross entropy loss

$$J(\boldsymbol{\theta}) = \sum_i y^{(i)} \log(\hat{\boldsymbol{y}}^{(i)})$$

# Now put it all together to perform classification

- Extract high-level features from last conv and pool layer

- Add fully connected layers to learn complex transformation of the features

- Output expressed as probability distribution over classes



INPUT    CONVOLUTION + RELU    POOLING    CONVOLUTION + RELU    POOLING      FLATTEN    FULLY CONNECTED    SOFTMAX

— CAR
— TRUCK
— VAN
— BICYCLE

We'll see softmax on next lecture

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

-MIT *Introduction to Deep Learning* http://introtodeeplearning.com

# Now put it all together to perform <mark>classification</mark>

- Convolution to learn local features from input image

- Introduce non-linear activation

- Reduce dimensionality and preserve local translation invariance with pooling

# One feature extraction tool, many tasks!



**FEATURE LEARNING**

You can use feature learning part
to perform a variety of tasks!

An overview of famous architectures and datasets

# LeNet (LeNet-5)

Once upon a time… LeNet [LeCun et al., 1998]

- 5x5 conv filters with stride 1

- 2x2 pooling with stride 2



CONV-POOL-CONV-POOL-FC-FC architecture

# MNIST (1998)

- Handwritten digits recognition

- 10 classes (of course!)

- 60k training images (28x28), 10k test images

- Grayscale images

- MNIST: Modified National Institute of Standards and Technology database

- Original NIST: training dataset was taken from American Census Bureau employees, while the testing dataset was taken from American high school students

- Available on Tensorflow

- LeNet on MNIST was the first real NN application!

- MNIST is not much of a challenge anymore…

# Fashion MNIST (2017)

- Clothing classification

- 10 classes

- 60k training images (28x28), 10k test images

- Grayscale images

- Available on Tensorflow

- More challenging then MNIST

# Cifar10 (2010)

- 10 classes

- 60k images (32x32)

- Balanced dataset: 6000 images each

- Colour images

- It is a labeled subset of the 80 million tiny image dataset

- Available on Tensorflow

- Widely used for research nowadays

# Cifar100 (2010)

- 100 classes

- 60k images (32x32)

- Balanced dataset: 600 images each

- Colour images

- It is a labeled subset of the 80 million tiny image dataset

- Available on Tensorflow

- Widely used for research nowadays (less than Cifar10)

# Imagenet (2009)

- 21841 classes

- 14M images with different dimensions and resolutions (many apply resize to 256x256)

- Unbalanced dataset

- Colour images

- Lead developer Fei-Fei Li



"Elongated crescent-shaped yellow fruit with soft sweet flesh"



1409 pictures of bananas.

ImageNet Large Scale Visual Recognition Challenges

Classification task: produce a list of object categories present in image. 1000 categories.

"Top 5 error": rate at which the model does not output correct label in top 5 predictions

# An overview on the most famous architectures

Imagenet – visual recognition challenge with 1000 classes.

Winners:



First CNN-based winner

# AlexNet

AlexNet (2012) renewed interest in CNN.

It uses:

- RELU (first use)

- Layer normalization

- Dropout

- MaxPooling

- Momentum

- Data augmentation in training



AlexNet

Full (simplified) AlexNet architecture:
[227x227x3] INPUT
[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer
[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer
[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2
[4096] FC6: 4096 neurons
[4096] FC7: 4096 neurons
[1000] FC8: 1000 neurons (class scores)

# AlexNet

AlexNet (2012) renewed interest in CNN.

It uses:

-      RELU (first use)

-      Layer normalization

-      Dropout

-      MaxPooling

-      Momentum

-      **Data augmentation in training**



**Data Augmentation:**

a. No augmentation (= 1 image)

224x224

b. Flip augmentation (= 2 images)

224x224

c. Crop+Flip augmentation (= 10 images)

224x224

+ flips

From: http://cs231n.stanford.edu - Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.

# AlexNet

AlexNet (2012) renewed interest in CNN.

It uses:

- RELU (first use)

- Layer normalization

- Dropout

- MaxPooling

- Momentum

- **Data augmentation in training**



62.3 million parameters: the original paper showed that the depth of the model was essential for its high performance, which was computationally expensive, but made feasible due to the utilization of (GPUs) during training.

# An overview on the most famous architectures

Imagenet – visual recognition challenge with 1000 classes.

Winners:



First really «deep» network

# VGG

VGG (2014):

- many small convolutional filters stacked together before pooling

- very deep (at the time) with 16/19 layers



From: http://cs231n.stanford.edu

# VGG



Wait. Why using smaller conv filters?

AlexNet          VGG16          VGG19

# VGG

Consider a stack of three 3x3 conv layers.

Which is the receptive field of this hidden unit?



Layer 3

Layer 2

Layer 1

**AlexNet**

**VGG16**

**VGG19**

# VGG

Consider a stack of three 3x3 conv layers.

Which is the receptive field of this hidden unit?



Layer 3

Layer 2

Layer 1

**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 |

From: http://cs231n.stanford.edu

# VGG

Consider a stack of three 3x3 conv layers.

Which is the receptive field of this hidden unit?

Same as one layer with 7x7 conv filters.

Layer 1

Still, 138M parameters for VGG16!

But three layers mean more non-linearities, i.e. more complex features...

... and fewer parameters!

C channels (e.g. C=3 for RGB images):

- 7x7 conv has $7^2 \times C = 147$ parameters

- 3 layers of 3x3 conv have $3 \times (3^2 \times C) = 81$

AlexNet

| | |
|---|---|
| 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 |

# An overview on the most famous architectures

Imagenet – visual recognition challenge with 1000 classes.

Winners:

# GoogLeNet

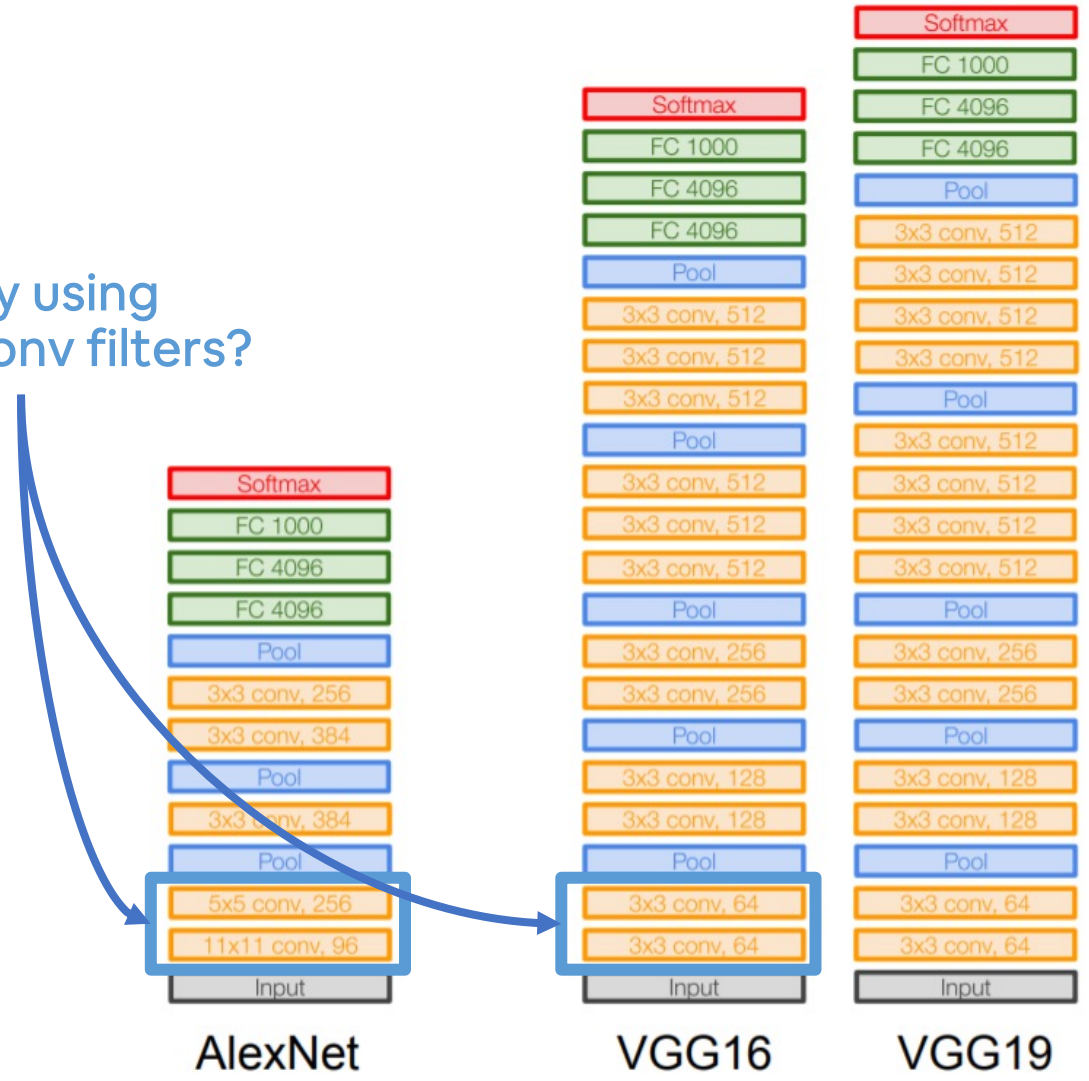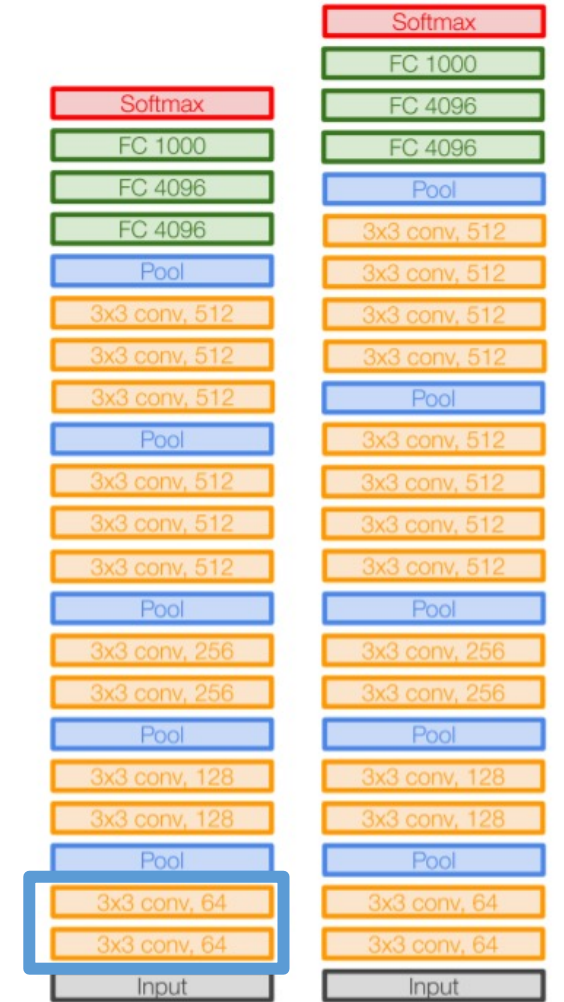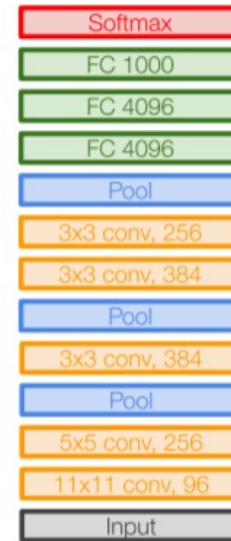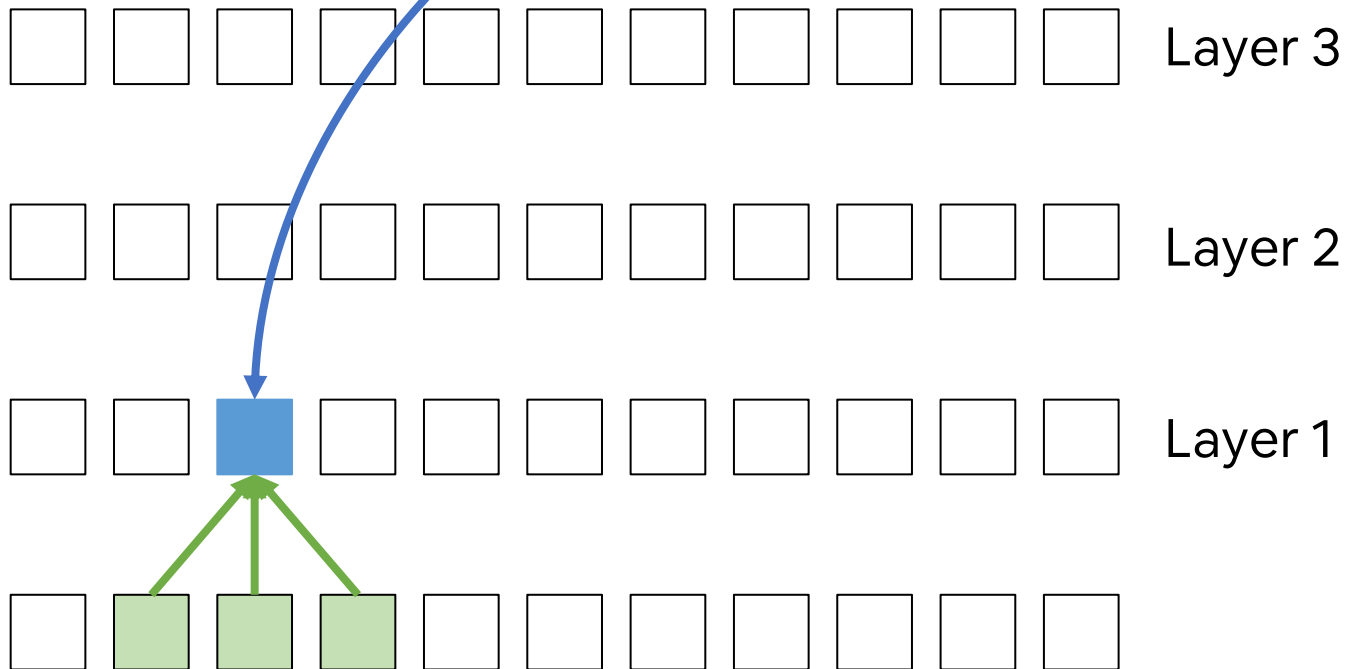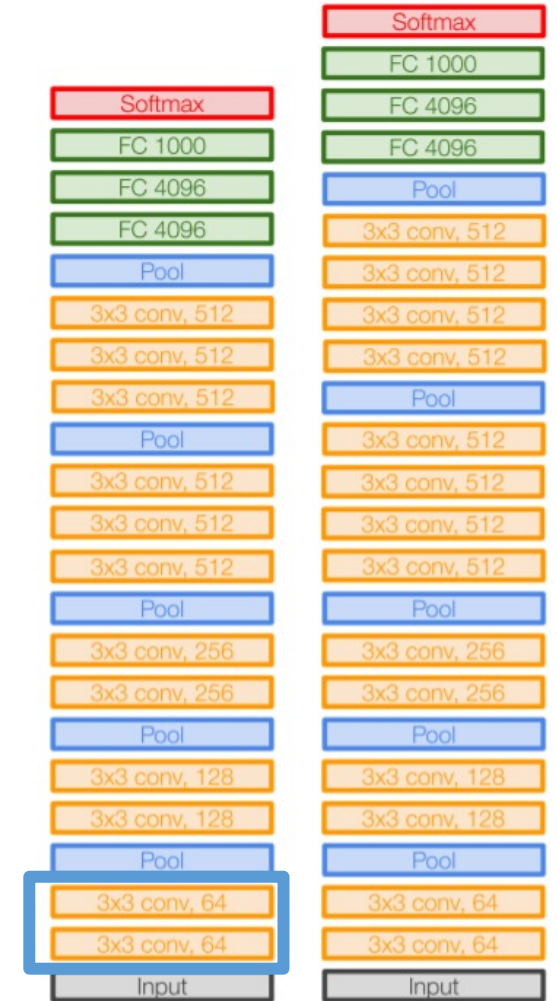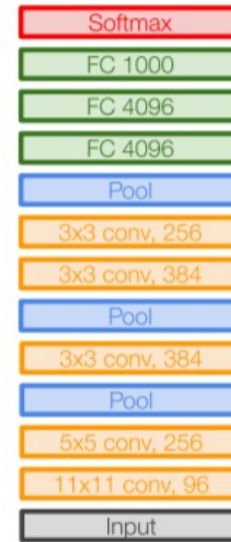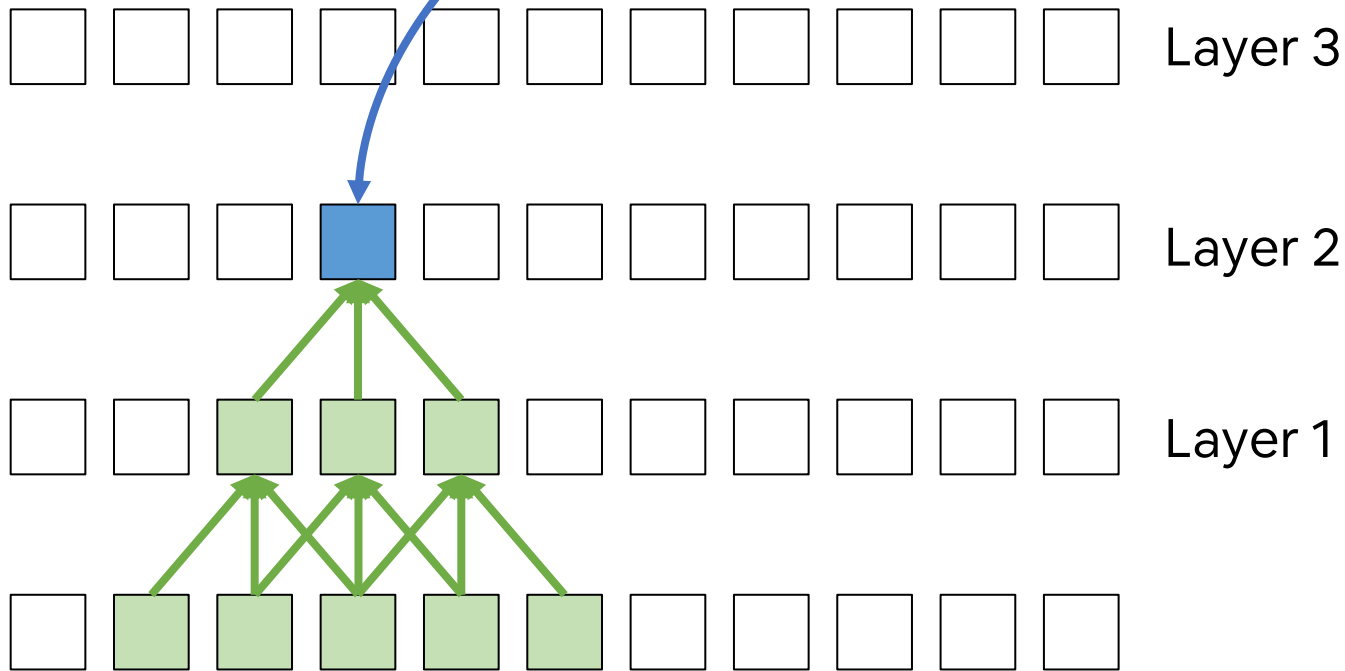Inception + GoogLeNet (2015)– introduces
parallel conv blocks (inception)



Inception module

# GoogLeNet

Inception + GoogLeNet (2015)– introduces
parallel conv blocks (inception)

Cleverly uses 1x1 convolutions

Preserves spatial dimensions, but reduces depth! Feature maps (depth) are projected to lower dimension



1x1 CONV with 32 filters

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



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

Filter concatenation

28x28x128    28x28x192    28x28x96    28x28x256

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

Module input: 28x28x256

Input

**Naive Inception module**



28x28x480

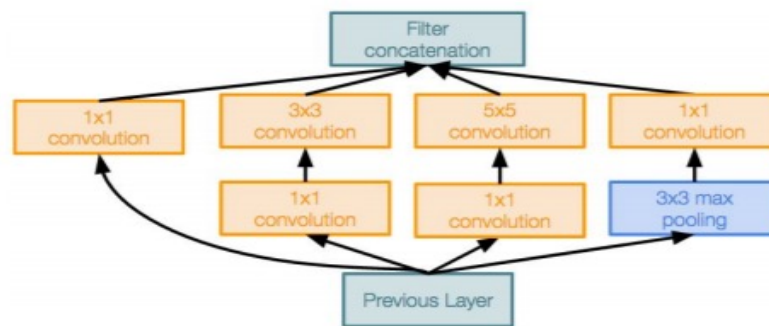Filter concatenation

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**

# GoogLeNet

Inception + GoogLeNet (2015)– introduces
parallel c...

Cleverly

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

Preserves spatial dimensions, but reduces depth! Feature maps (depth) are projected to lower dimension

**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**



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

Filter concatenation

28x28x128     28x28x192     28x28x96     28x28x256

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

Module input: 28x28x256     Input

Naive Inception module



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

# GoogLeNet



**Full GoogLeNet architecture**

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

# An overview on the most famous architectures

Imagenet – visual recognition challenge with 1000 classes.

Winners:

# ResNet

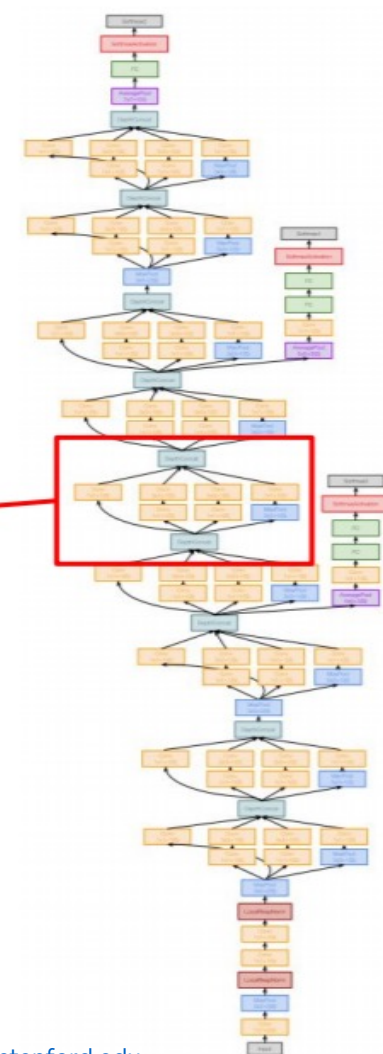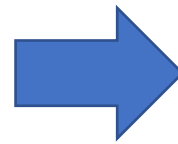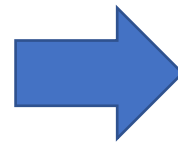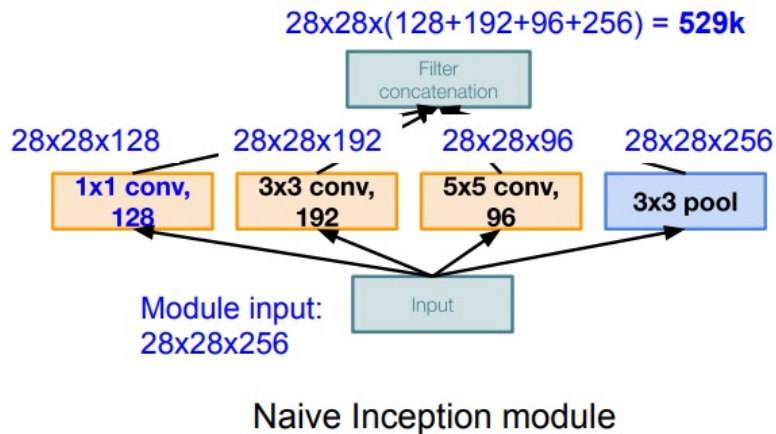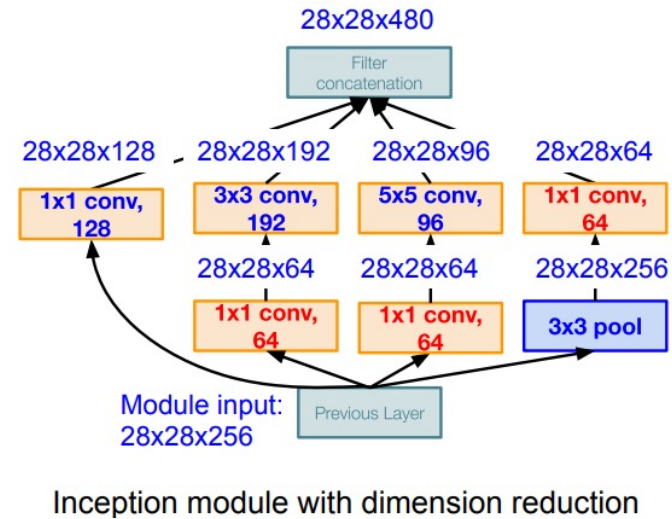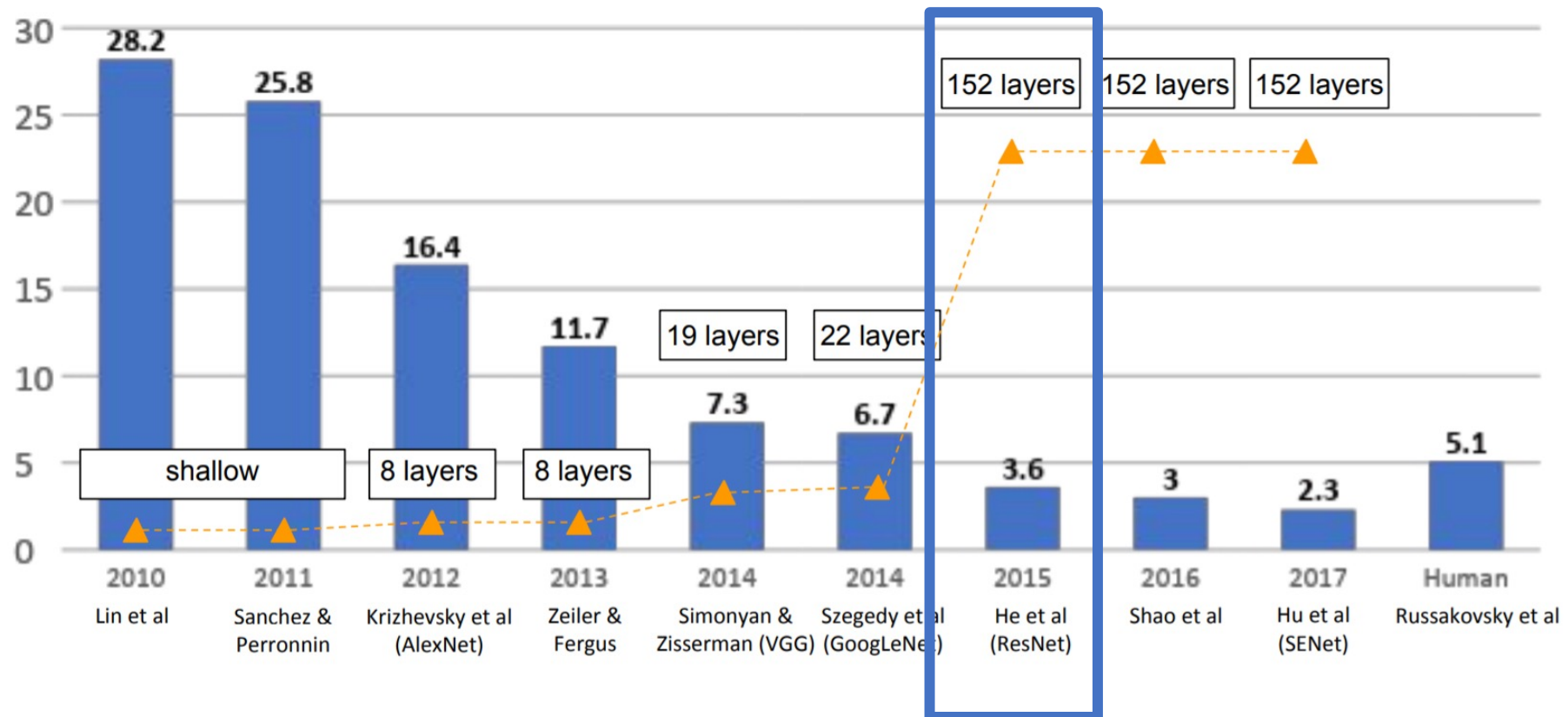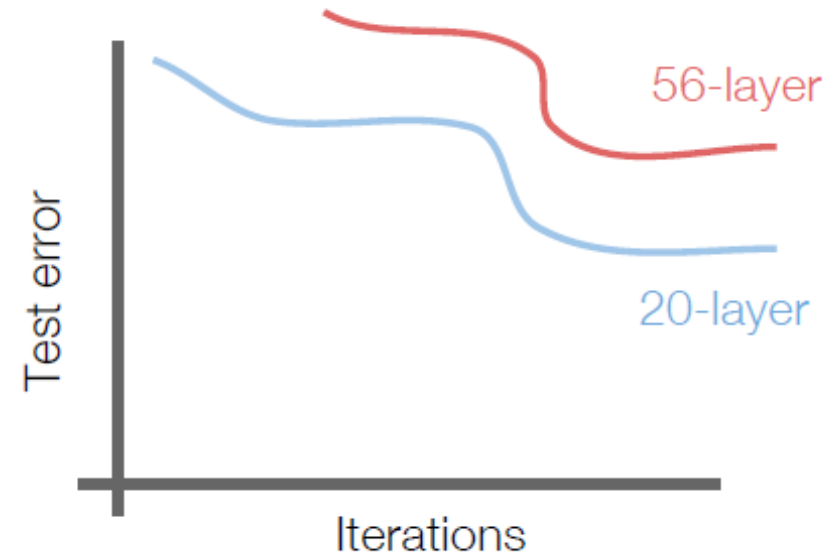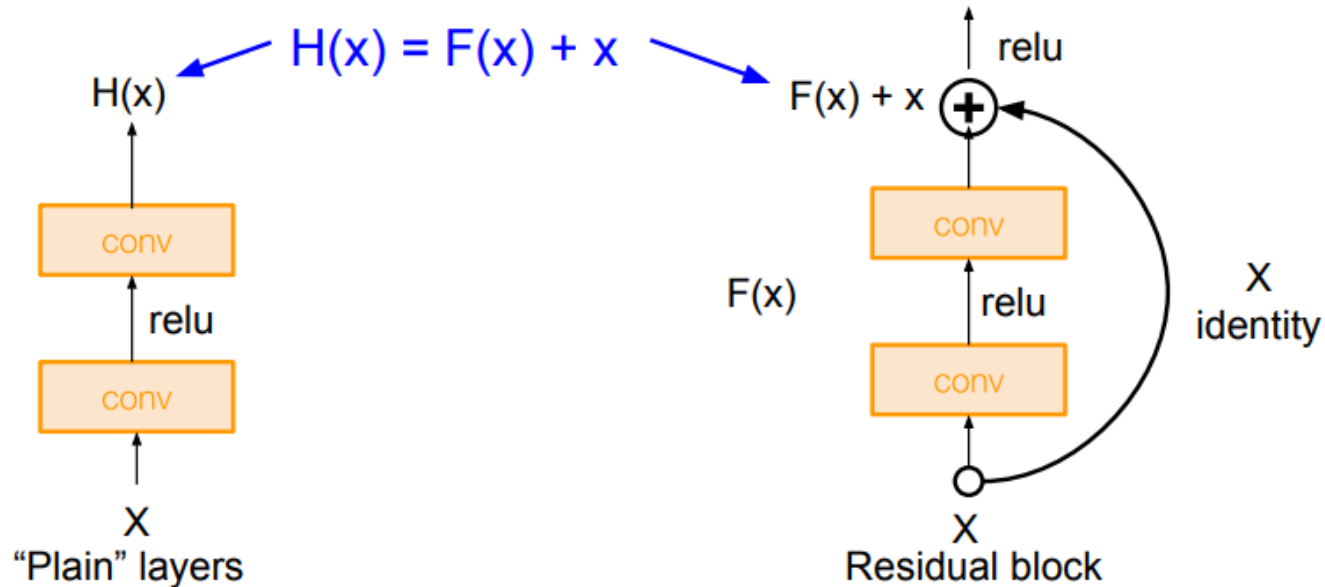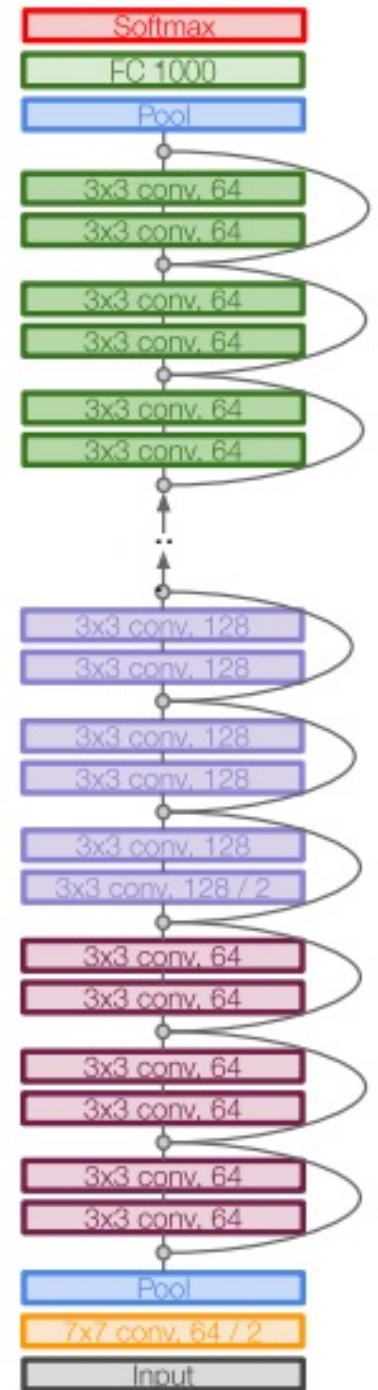Why not stacking more and more layers?

# ResNet

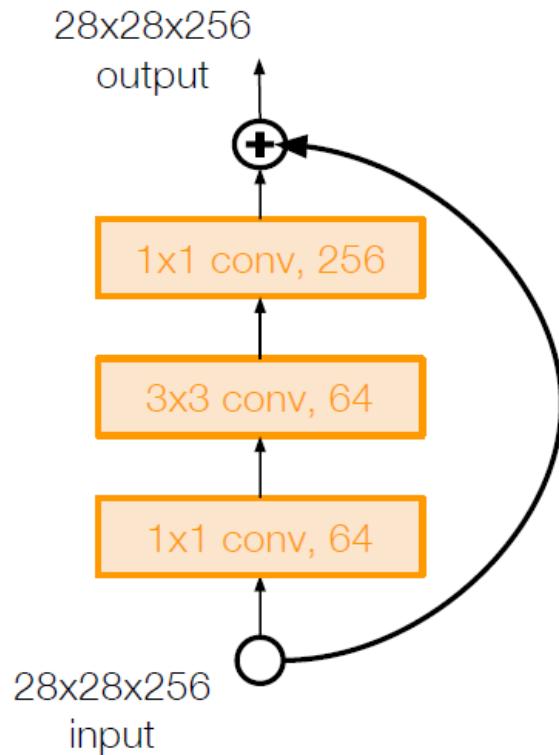ResNet (2015) – very deep model (152 layer) with shortcut connections



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

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

"Plain" layers

Residual block

# ResNet

28x28x256
output



For ResNet with more than 50 layers

**ResNet training:**

- 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

# State of the art is always on the move...

# Is not all about accuracy...
# EfficientNet

Tan, M., & Le, Q. (2019, May). Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning* (pp. 6105-6114). PMLR.
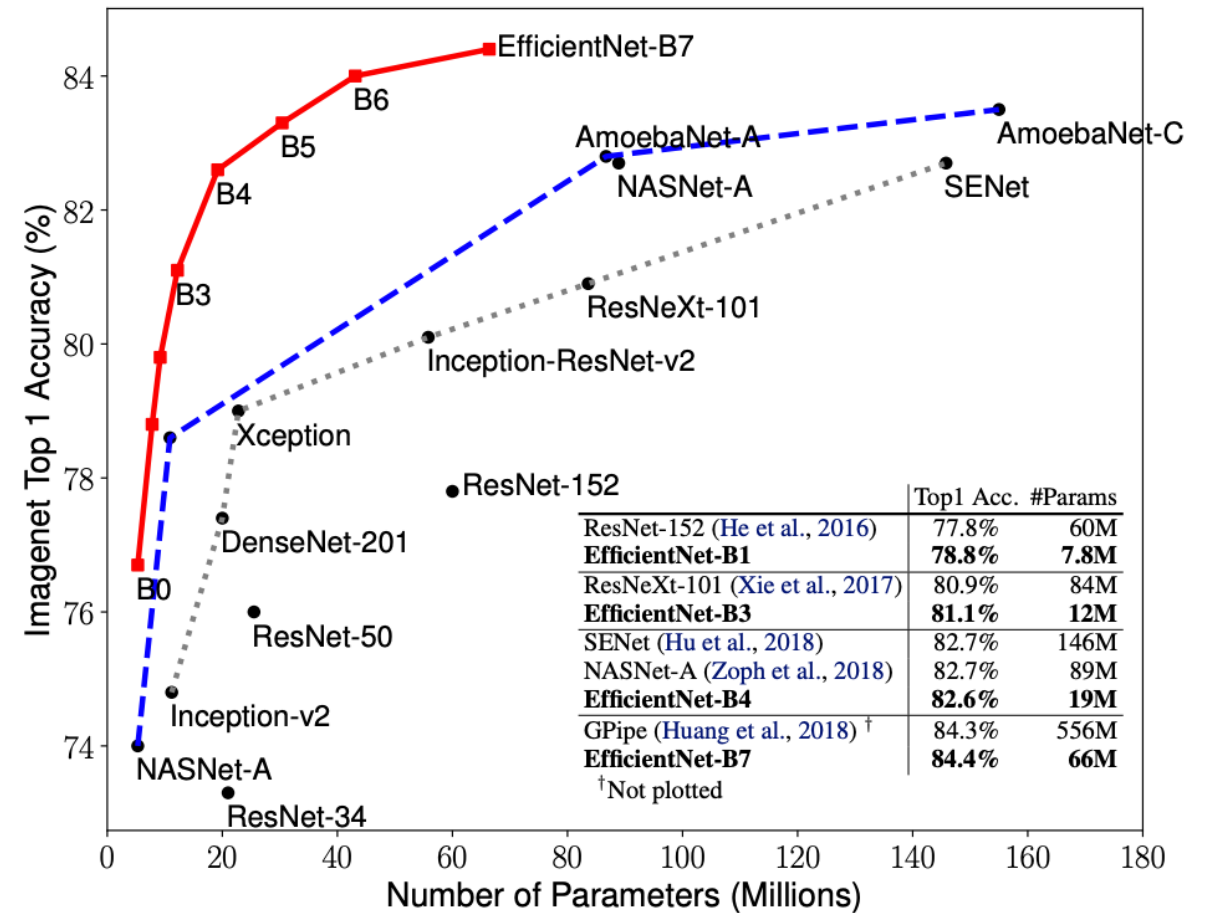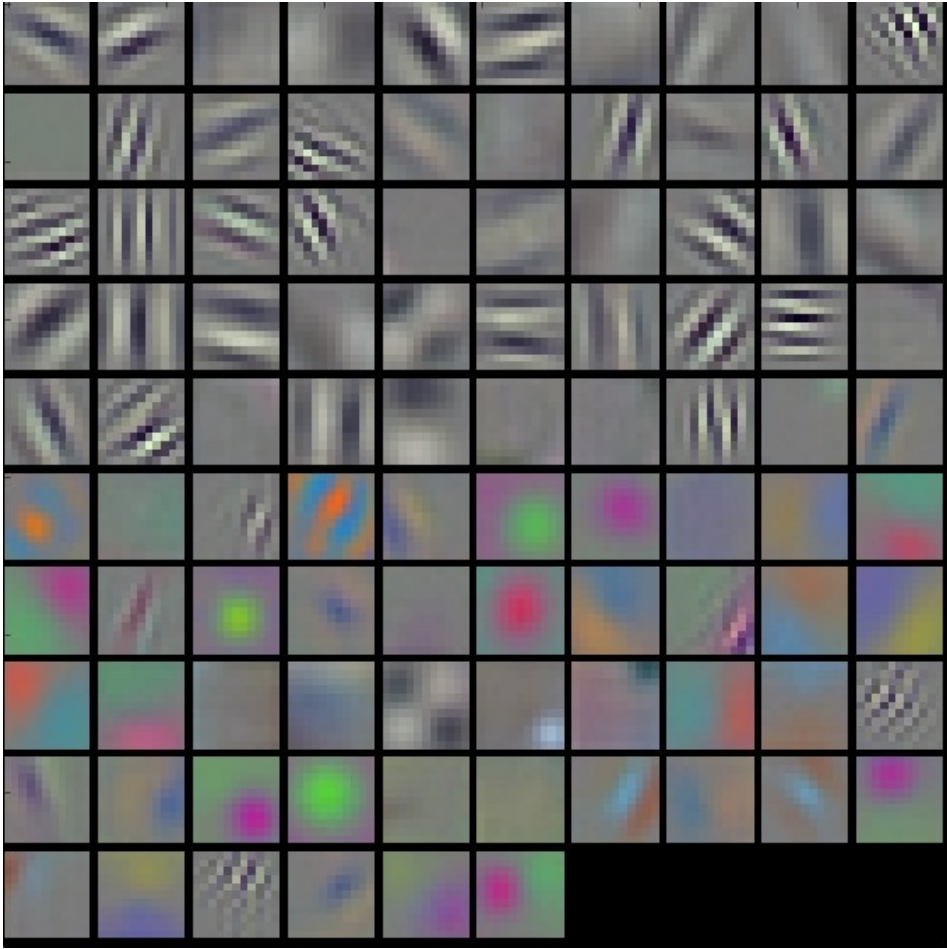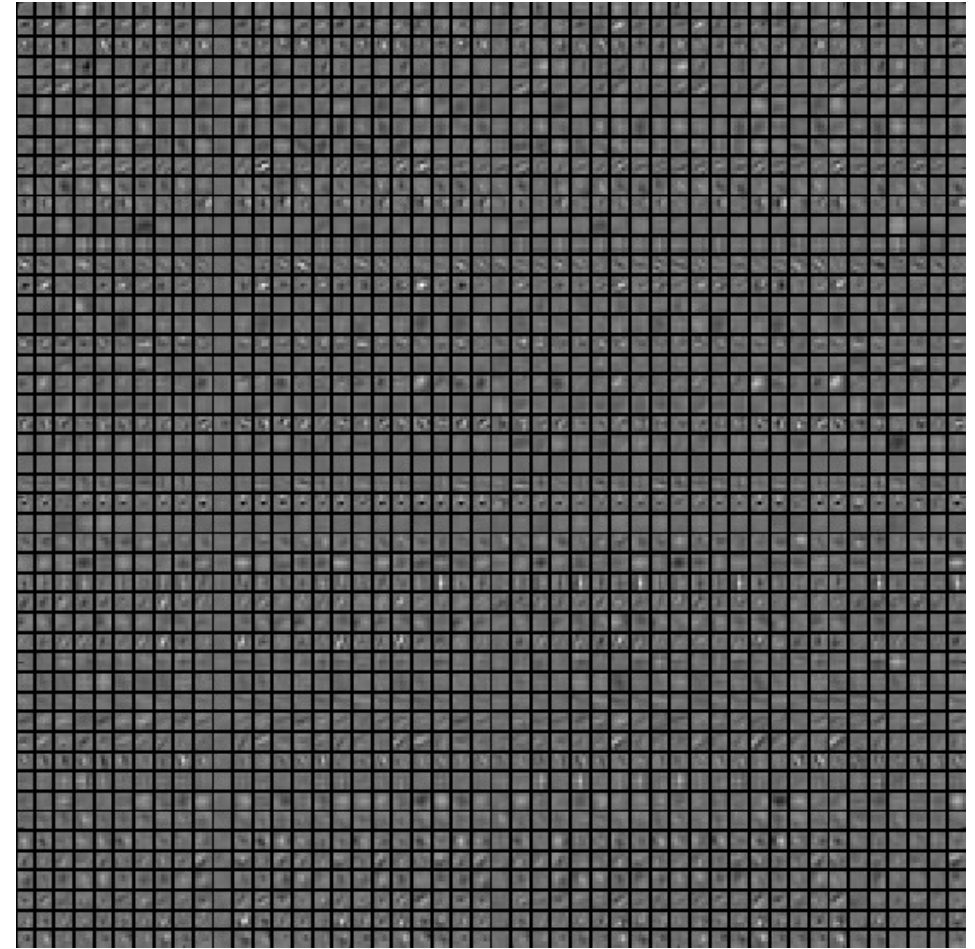


| | Top1 Acc. | #Params |
|---|---|---|
| ResNet-152 (He et al., 2016) | 77.8% | 60M |
| **EfficientNet-B1** | **78.8%** | **7.8M** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 84M |
| **EfficientNet-B3** | **81.1%** | **12M** |
| SENet (Hu et al., 2018) | 82.7% | 146M |
| NASNet-A (Zoph et al., 2018) | 82.7% | 89M |
| **EfficientNet-B4** | **82.6%** | **19M** |
| GPipe (Huang et al., 2018) † | 84.3% | 556M |
| **EfficientNet-B7** | **84.4%** | **66M** |
| †Not plotted | | |

*Figure 1.* **Model Size vs. ImageNet Accuracy.** All numbers are for single-crop, single-model. Our EfficientNets significantly outperform other ConvNets. In particular, EfficientNet-B7 achieves new state-of-the-art 84.4% top-1 accuracy but being 8.4x smaller and 6.1x faster than GPipe. EfficientNet-B1 is 7.6x smaller and 5.7x faster than ResNet-152. Details are in Table 2 and 4.

# What do CCN see?

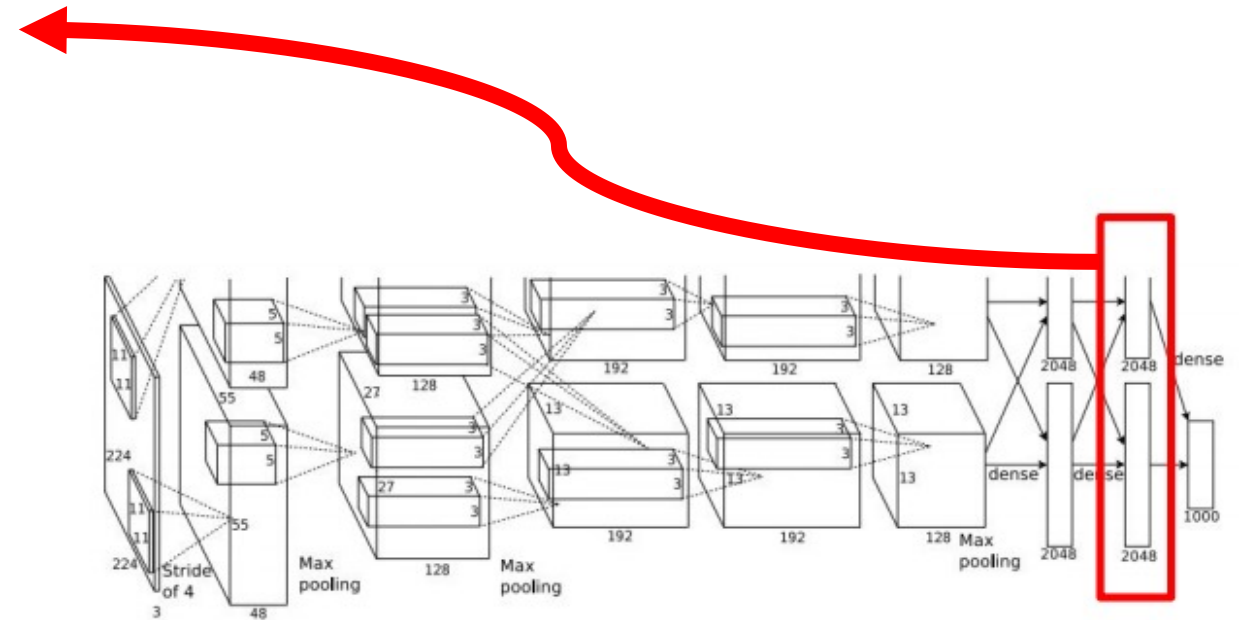# What do CNN see? Visualize the layers



AlexNet 1CONV layer

AlexNet 2CONV layer

# What do CNN see? Embedding space for features

We can **consider k-nearest neighbors in embedding space for last FC layer:**



Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.

# What do CNN see? Embedding space for features

We can **plot final FC embedding layer by means of dimensionality reduction**, e.g. tSNE (more powerful than PCA) or UMAP

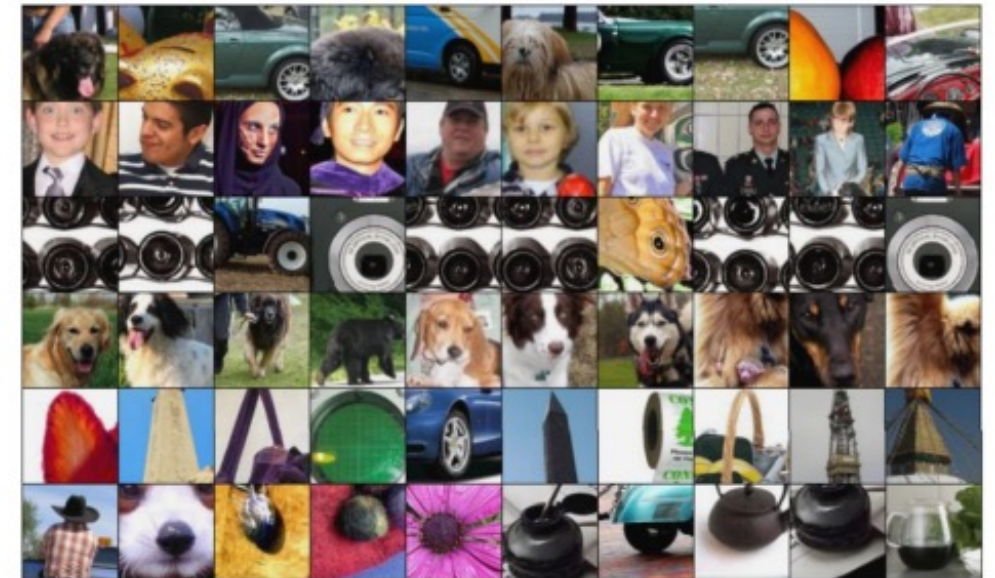Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008

From: http://cs231n.stanford.edu

# What do CNN see? Maximally activating neuros

We can **compute maximally activating patches.**

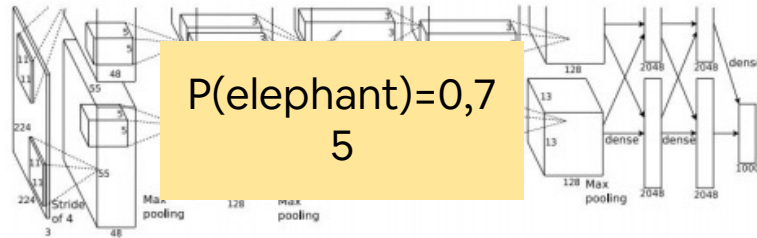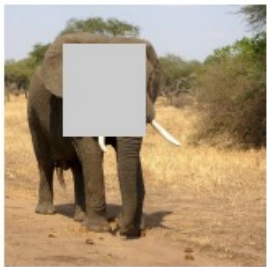Run many images through the network, record values of chosen channel (e.g. channel 17/128 in conv5).
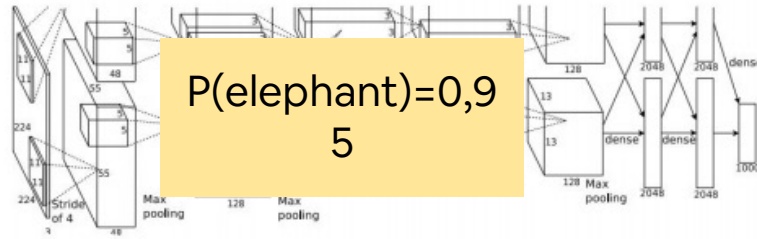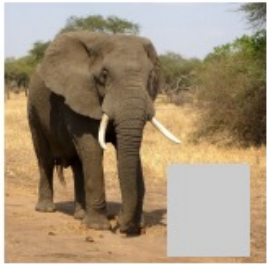
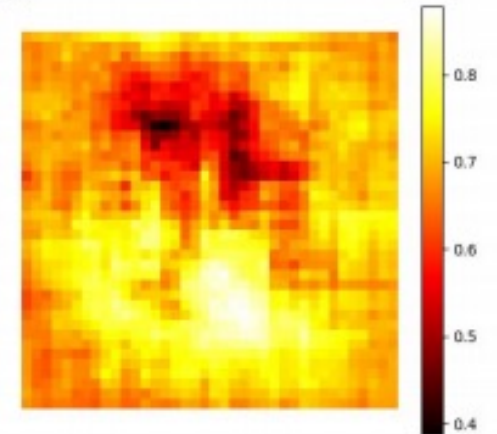Visualize image patches that correspond to maximal activations.

Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015;



From: http://cs231n.stanford.edu

# What do CNN see? Most relevant pixels

**Saliency maps**, e.g. by occlusion:



P(elephant)=0,95

P(elephant)=0,75

African elephant, Loxodonta africana

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

# Thank you!

## Gian Antonio Susto