

8. CONVOLUTIONAL NEURAL NETWORK ARCHITECTURES

Stephan Robert-Nicoud
HEIG-VD/HES-SO

Credit: Andres Perez-Urbe

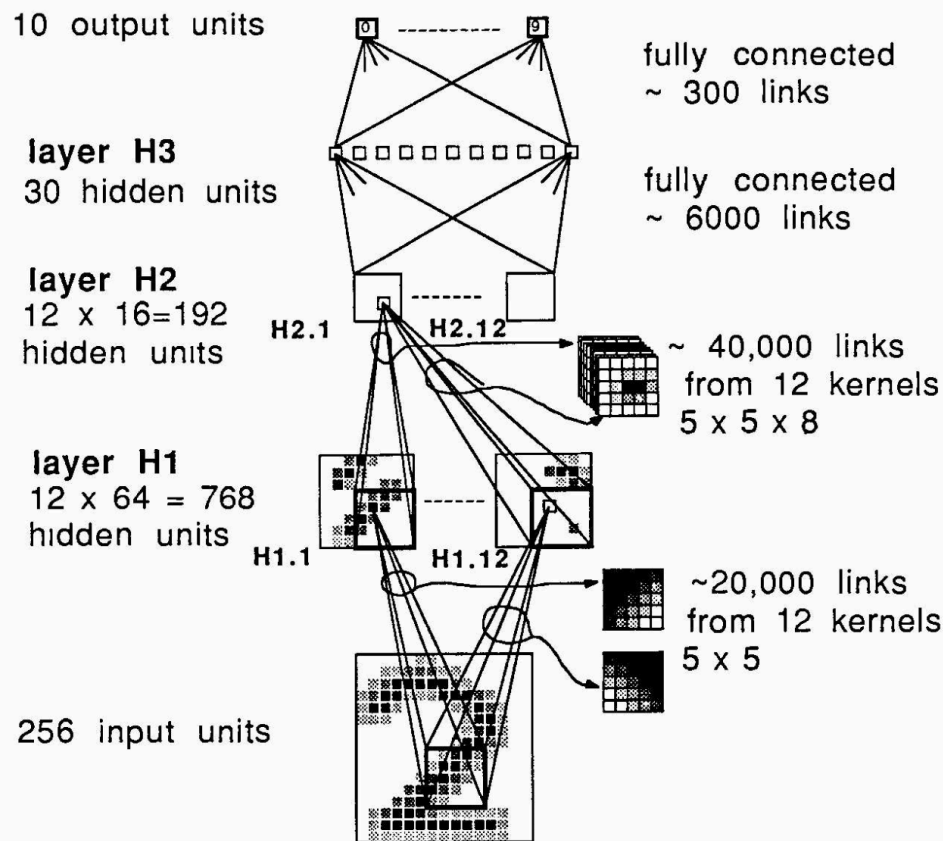


Objectives

- ❑ Understand the evolution of the architectures of Convolutional Neural Networks
- ❑ Understand the motivation and the characteristics of the main architectures of Convolutional Neural Networks
- ❑ Train a Convolutional Neural Network for an image recognition problem

LeNet-5 (1989)

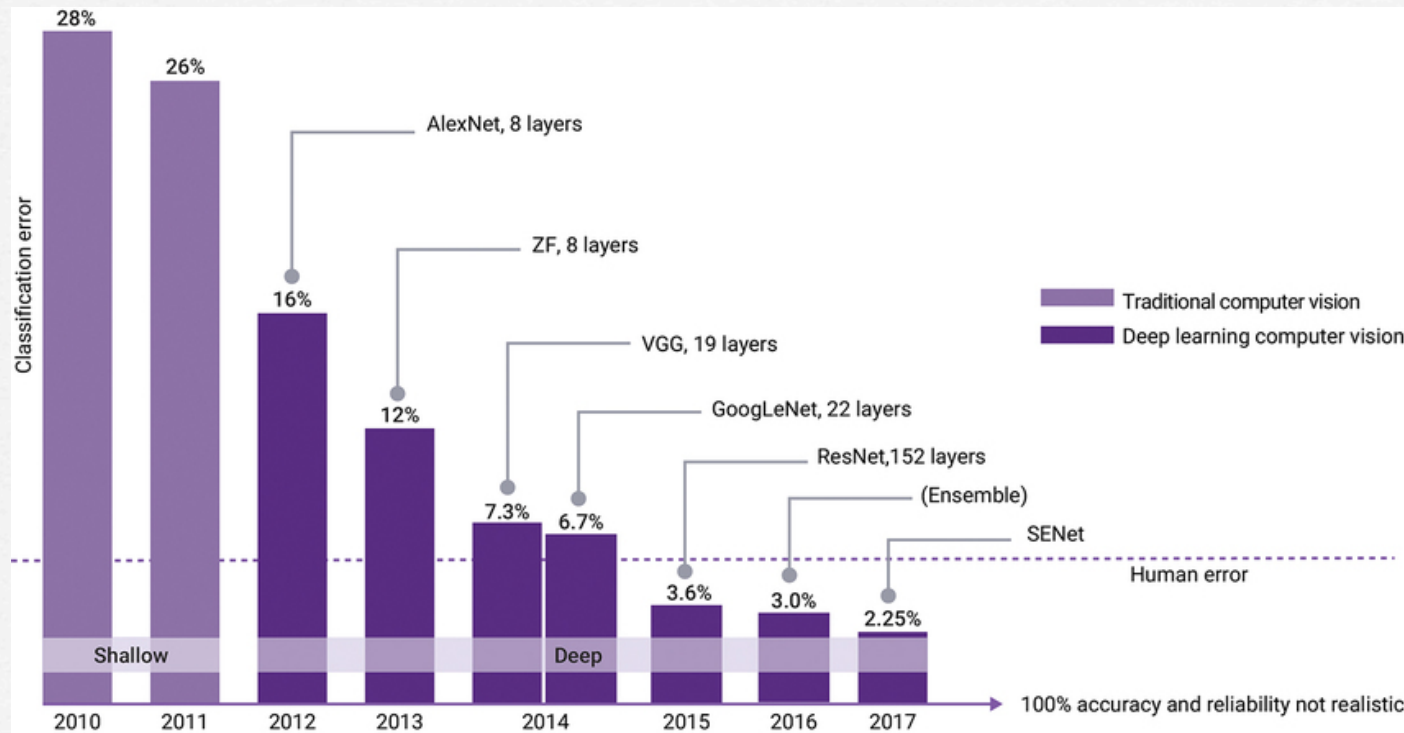
□ By Yann LeCun and colleagues (AT&T Bell Labs)



- "Backpropagation applied to Handwritten Zip code Recognition", Neural Computation 1
- "Large networks trained by Backpropagation can be applied to real image-recognition problems without complex pre-processing requiring detailed engineering"
- Trained on 9298 digits from US mail during 3 days on a Sun Sparc Station 1

The ILSVRC challenge

ImageNet Large Scale Visual Recognition Challenge

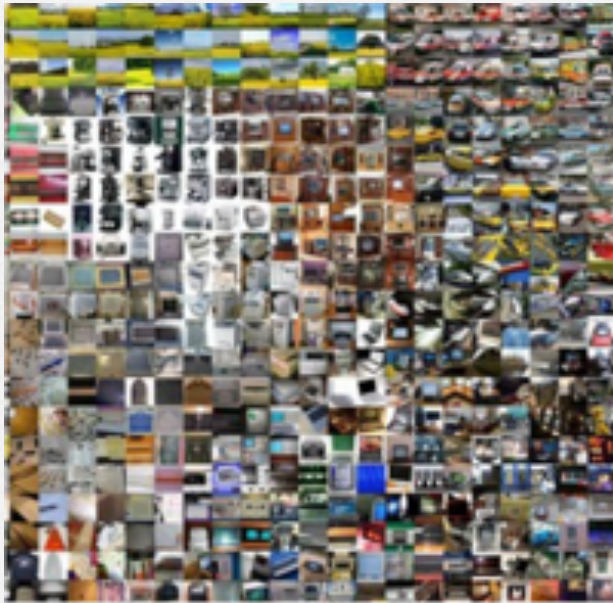


from synopsys.com

ImageNet: 15 million images (22k categories), human labeled

ILSVRC Challenge: ~1000 images of 1000 categories (top-1/top-5 error rates)

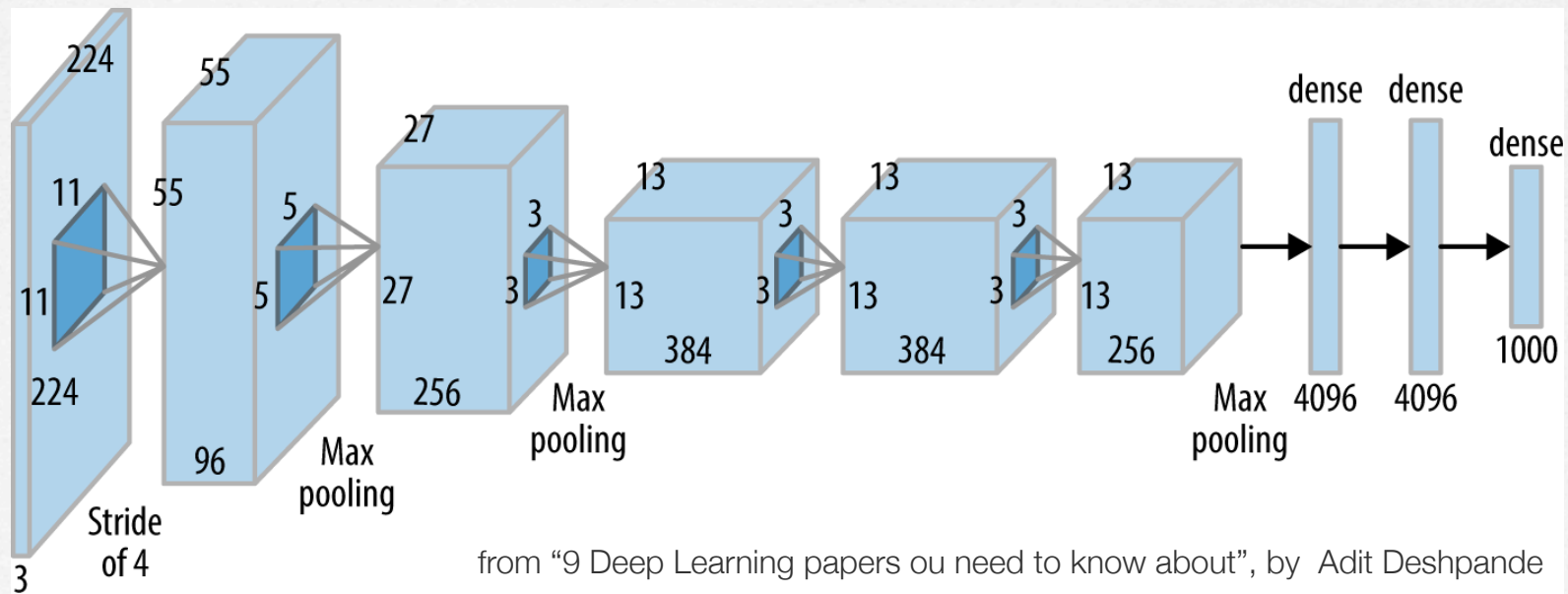
ImageNet



- 2007-2010 human labeling using Amazon Mechanical Turk
- 49K workers from 167 countries
- Human performance: ~5% (top-5 error) due to class unawareness, need for fine-grained recognition

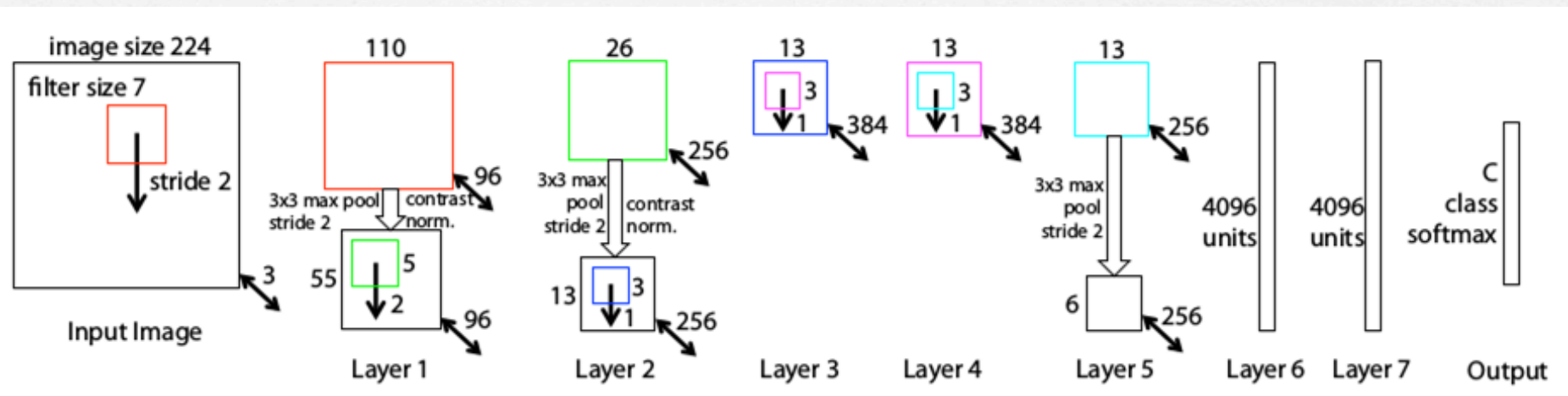
AlexNet (2012)

- By Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton
- 62-million parameters model trained on ImageNet data using batch stochastic gradient descent, with specific values for momentum and weight decay. Trained on two GTX 580 GPUs for 5 to 6 days.
- Used ReLus (instead of tanh), data augmentation and dropout.
- Won the 2012 ILSVRC (ImageNet Large-Scale Visual Recognition Challenge): error = 15.4% (second place 26.2%).



ZF Net (2013)

- ❑ by Matthew Zeiler and Rob Fergus from NYU
- ❑ Trained on only 1.3M images (ImageNet) using batch stochastic gradient descent to minimize cross-entropy.
- ❑ Used smaller kernel sizes (7x7) and increasing number of filters
- ❑ Trained on two GTX 580 GPUs for 12 days.
- ❑ Won the 2013 ILSVRC: error = 11.2%

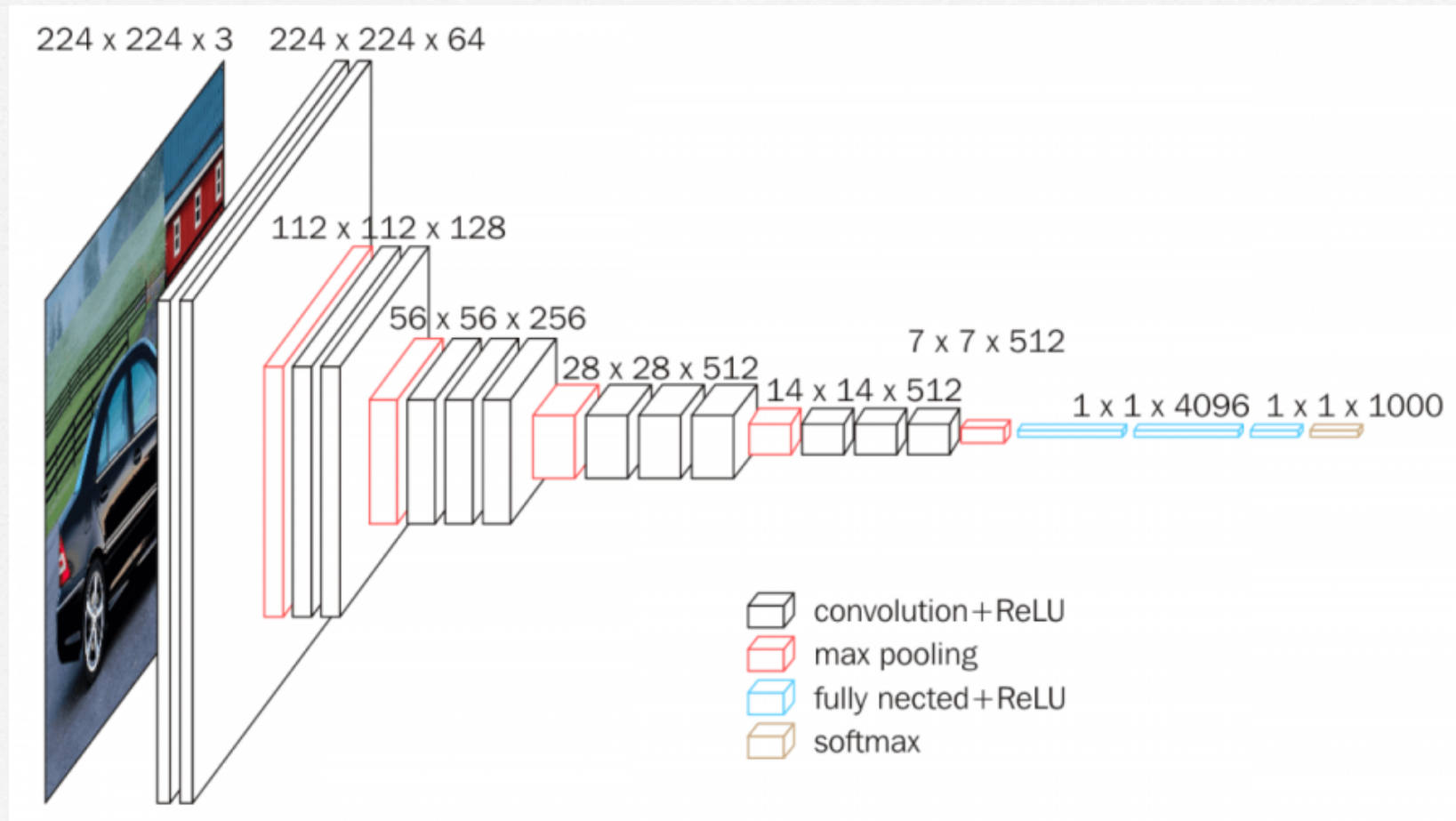


from “9 Deep Learning papers ou need to know about”, by Adit Deshpande
APE 2024

VGG Net (2014)

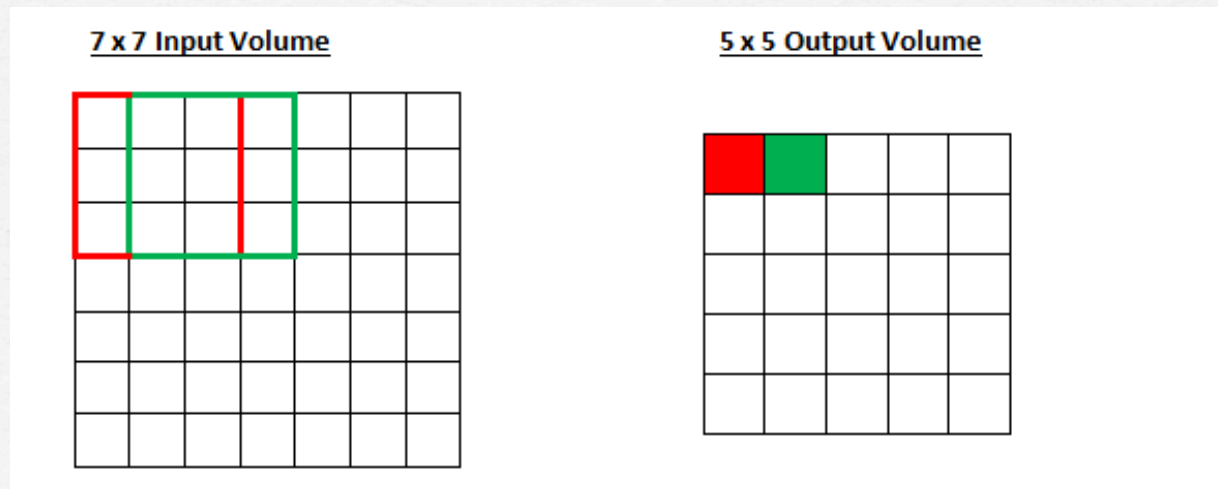
- ❑ by K. Simonyan and A. Zisserman from the Visual Geometry Group of the University of Oxford
- ❑ VGG-16 is a 16-layer CNN and VGG-19 is a 19-layer CNN that **strictly used 3x3 filters with stride and pad of 1, along with 2x2 maxpooling layers with stride 2**
- ❑ **Number of kernels: 64, 128, 256, 512, 512**
- ❑ Did not win ILSVRC (error 7.3%) but introduced a simple architecture
- ❑ Trained on 4 Nvidia Titan Black GPUs for two to three weeks
- ❑ VGG-16 has more than 138-million parameters and VGG-19 has almost 143-million parameters.

Ideas behind VGG Nets



Downsample input with stride

- ❑ Stride: the amount of movement between applications of the filter to the input image.
- ❑ Example: no padding and stride=1



The result is a downsampled image at the output, e.g., 5x5 instead of 7x7.

Ideas behind VGG Nets (2)

The Alexnet convolutional kernels of sizes 11x11, 5x5, and 3x3 can be replicated by making use of multiple 3x3 kernels as building blocks, but reducing the number of parameters to learn.

Input Feature Map and Receptive Field

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

Output for each receptive field

1	2	3
6	7	8
11	12	13

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

1	2	3
6	7	8
11	12	13

Output Feature Map of 1st conv layer

*	*	*
		*

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

1	2	3
6	7	8
11	12	13

Input Feature Map of 2nd conv layer

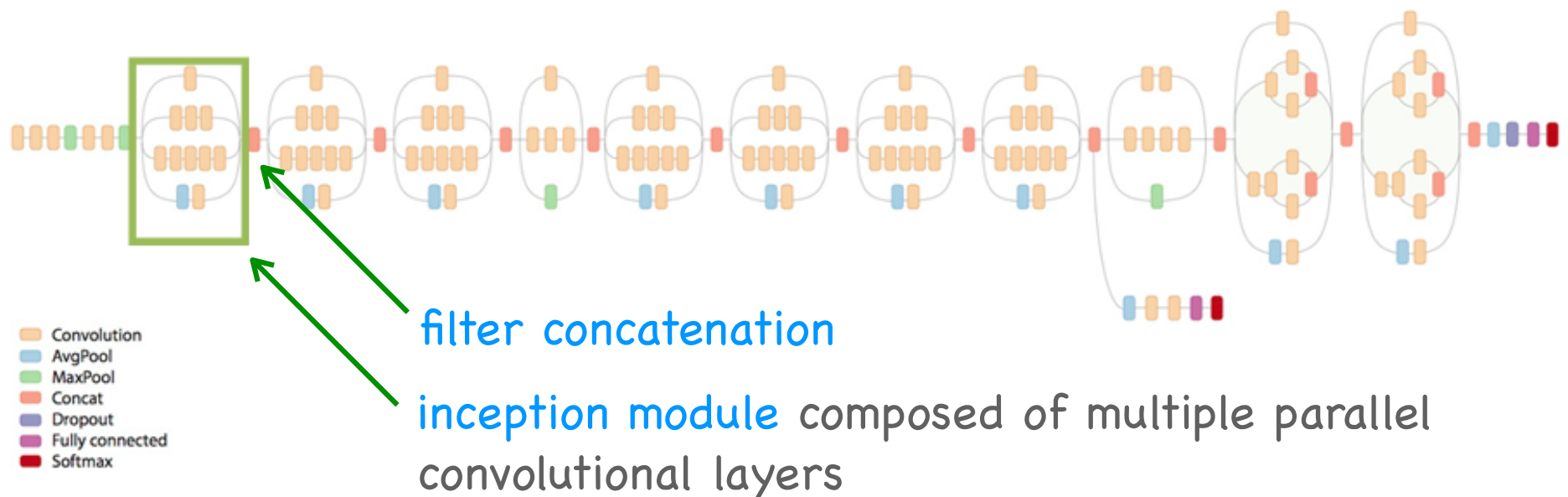
Output Feature Map of 2nd conv layer



- Suppose a 5x5x1 input. A 5x5 kernel (25 parameters) will produce a 1x1 output
- This can be computed by two layers of 3x3 kernels, stride=1 too (2 x 9 parameters)
- Similarly a 11x11 filter (121 parameters) can be computed by five 3x3 kernels, stride=1 (45 parameters)

GoogLeNet (2015)

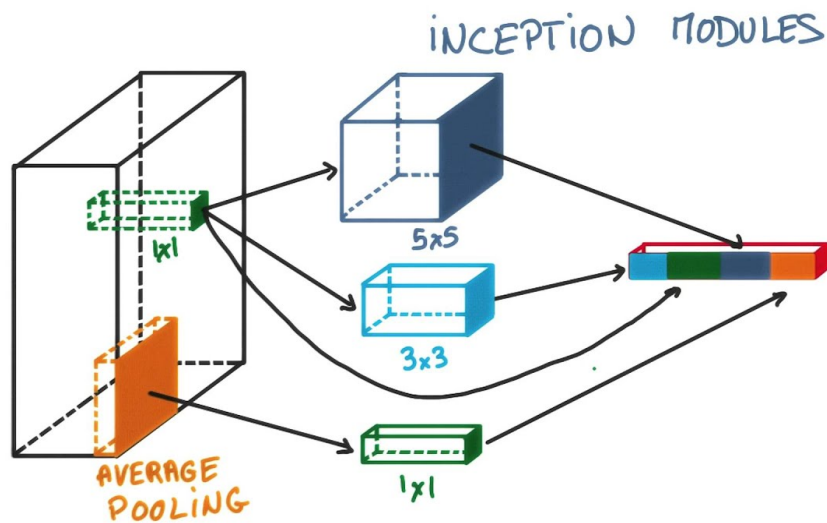
- ❑ Winner of ILSVRC'14 (error 6.7%)
- ❑ A new architecture with more than 100 layers in total, but not necessarily stacked up sequentially and **without fully-connected layers**. It uses 12x less weights than AlexNet



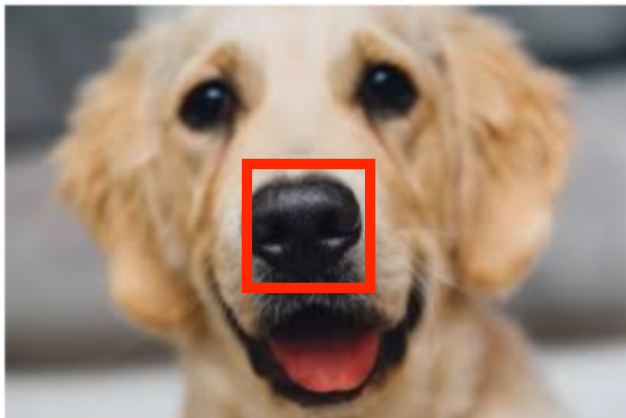
Green box shows parallel region of GoogLeNet

from "9 Deep Learning papers you need to know about", by Adit Deshpande
APE 2024

Ideas behind Inception Networks (1)

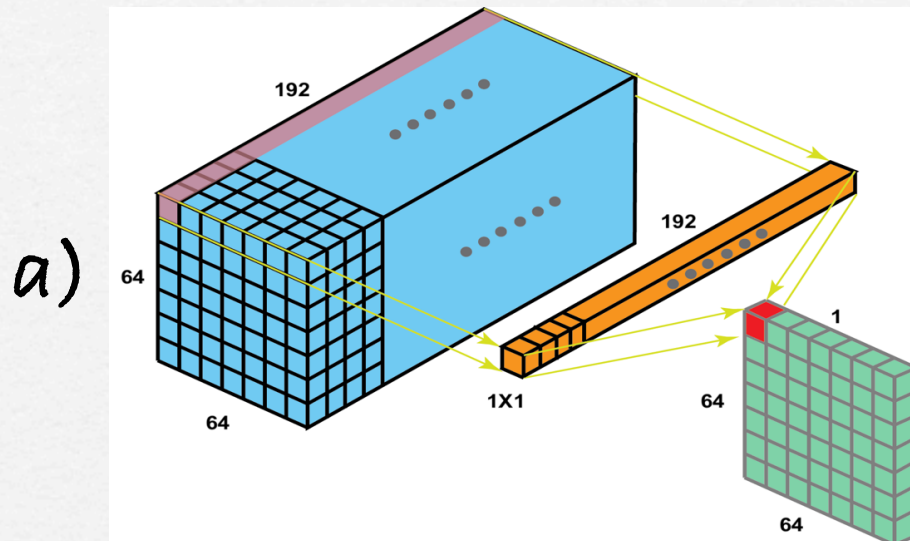


- Instead of choosing a single size (1×1 , 3×3 or 5×5) kernel for a layer, use them "all" and concatenate the outputs

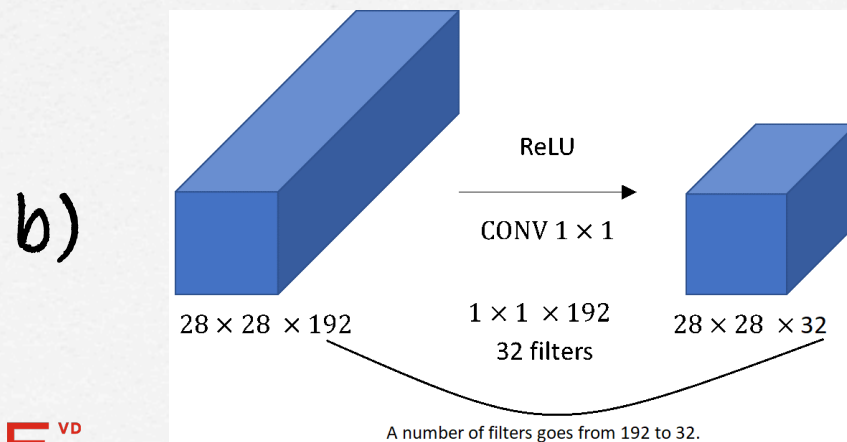


- different kernel sizes allows the identification of features at different scales

Ideas behind Inception Networks (2)

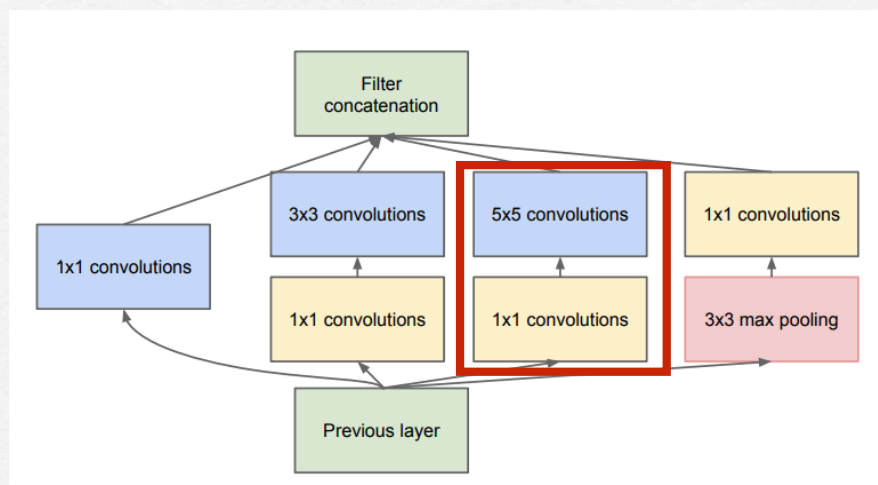


- Use **1x1 convolutions** (Fig a) to perform **feature pooling** or **dimensionality reduction**.



- Example: applying 32 1x1x192 filters to the 192 28x28 images in the figure b) produces 32 28x28 images

Ideas behind Inception Networks (3)



- Use **1x1 convolutions** before 3x3 or 5x5 convolutions to perform **feature pooling** and reduce the number of operations.

- Example: applying 16 1x1x192 filters to the incoming 192 28x28 images before applying 32 5x5 filters requires 12.4M multiplications instead of ~120M multiplications without the 1x1 convolutions.

Batch-normalization (1)

Sergey Ioffe and Christian Szegedy from Google (2015)

- ❑ **Problem:** when two inputs are in completely different scales, say a range of x_1 is $[1000-2000]$ and range of x_2 is $[0.1-0.5]$, using them as-is has implications on optimizing the loss function (e.g., by gradient descent). Intuitively, the model will tend to be more influenced by x_1 .
- ❑ To deal with this problem, we generally normalize inputs before using a neural network. **Normalization**, in general refers to squashing a diverse range of numbers to a fixed range.
- ❑ Ioffe and Szegedy did the same analysis within the network and noticed that the learning was enhanced when the inputs to a subsequent sub-network or layer are also normalized.

Batch-normalization (2)

A “learnable” normalization transform that is applied to every mini-batch during the learning process:

start by normalizing the inputs (per dimension) per batch (i.e., by setting the mean to zero and the standard deviation to 1), but then use the learnable parameters γ and β to scale and shift the input data. γ and β are learned over all the database and serve to restore the expressiveness of the network.

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

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

Batch-normalization (3)

- It has been shown that batch normalization accelerates training convergence (i.e., less learning steps are required), allows larger learning rate values, the training of models using saturating nonlinearities (e.g., sigmoids), and not needing dropout.
- It was originally proposed to be used before the non-linearity (e.g., relu), but it has been reported to work even better after it.

```
1 # example of batch normalization for an mlp
2 from keras.layers import Dense
3 from keras.layers import BatchNormalization
4 ...
5 model.add(Dense(32, activation='relu'))
6 model.add(BatchNormalization())
7 model.add(Dense(1))
8 ...
```

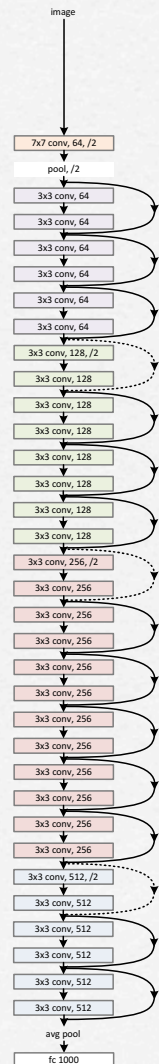
MLP model

```
1 # example of batch normalization for an cnn
2 from keras.layers import Dense
3 from keras.layers import Conv2D
4 from keras.layers import MaxPooling2D
5 from keras.layers import BatchNormalization
6 ...
7 model.add(Conv2D(32, (3,3), activation='relu'))
8 model.add(Conv2D(32, (3,3), activation='relu'))
9 model.add(BatchNormalization())
10 model.add(MaxPooling2D())
11 model.add(Dense(1))
12 ...
```

CNN model

ResNet (2015)

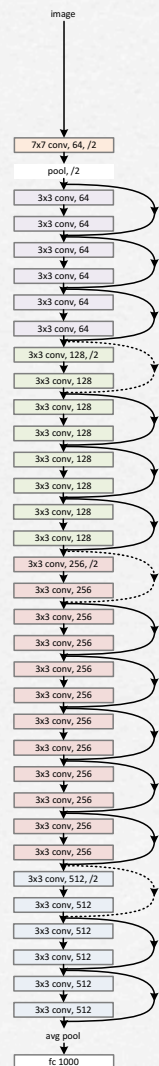
34-layer residual



- ❑ A degradation of performance was observed when going deeper (e.g., from 20 to >50 layers) on the ImageNet benchmark.
- ❑ A team from Microsoft won ILSVRC'15 (error: 3.6%) using a 152-layer CNN architecture exploiting « shortcut connections ».
- ❑ These so-called ResNet networks appear to have a better performance than plain networks with the same number of weights (ResNet152V2 has 60-million weights). Training took 3 weeks using 8 GPUs.

ResNet's residual blocks

34-layer residual



- The residual block assumes that it is easier to learn a mapping from $x \rightarrow F(x) + x$ than from $x \rightarrow$ to a new function $H(x)$
- A similar architecture called "Highway networks" was proposed by Schmidhuber et al. He claims to have successfully trained the first deep neural network with more than 100 layers.
- The team from Microsoft tested a 1202-layer network that performed less well.

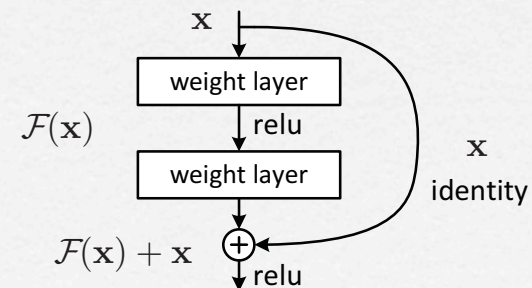
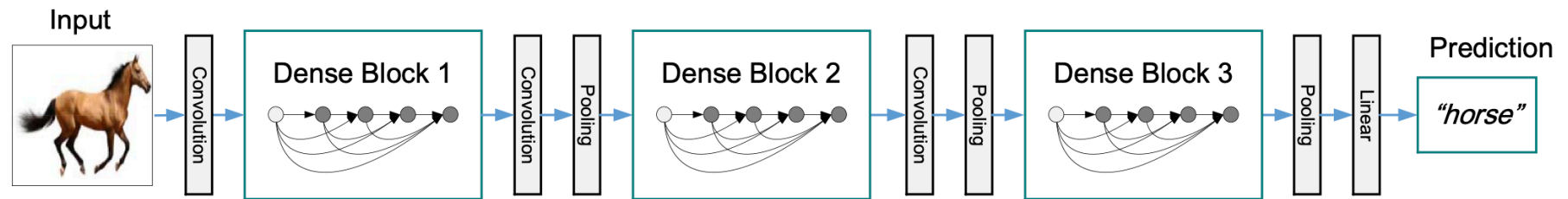


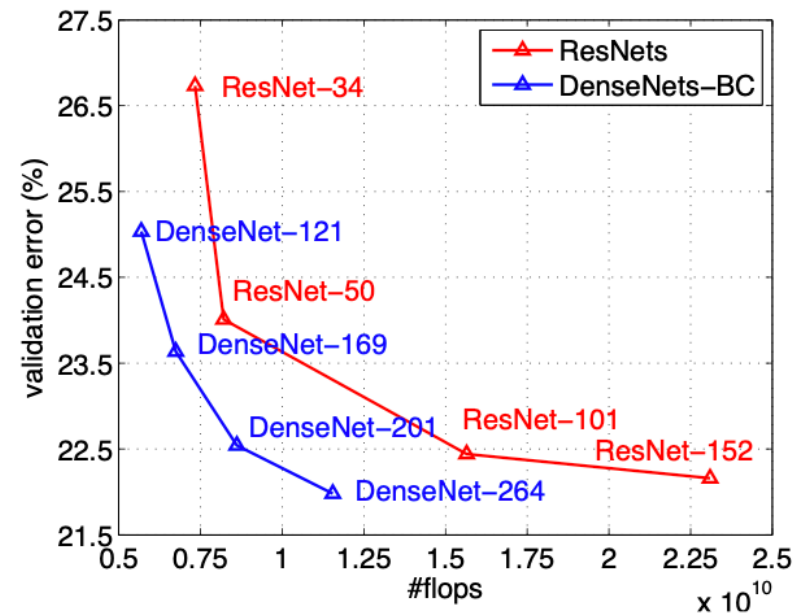
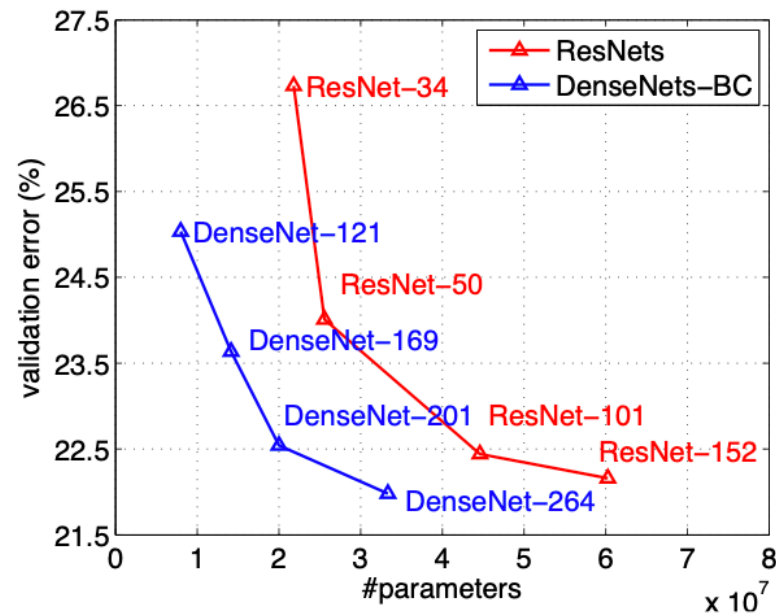
Figure 2. Residual learning: a building block.

DenseNet (2017)



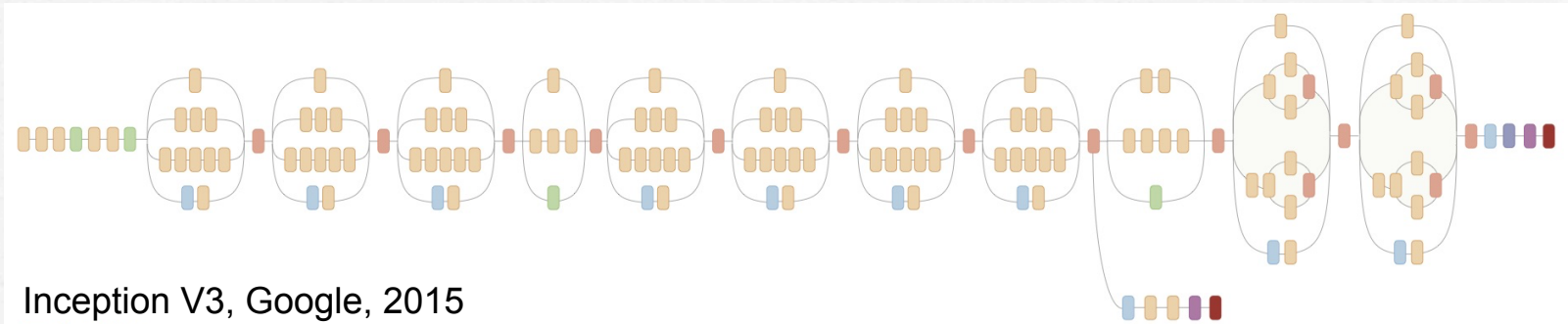
- ❑ Proposed by researchers from Cornell & Tsinghua universities and Facebook.
- ❑ For each layer, the feature-maps (i.e., the output of a convolution) of all preceding layers are used as inputs, and its own feature-maps are used as inputs into all subsequent layers.
- ❑ DenseNets have several advantages: they alleviate the vanishing-gradient problem, strengthen feature propagation, encourage feature reuse, and substantially reduce the number of parameters.

DenseNets vs ResNets

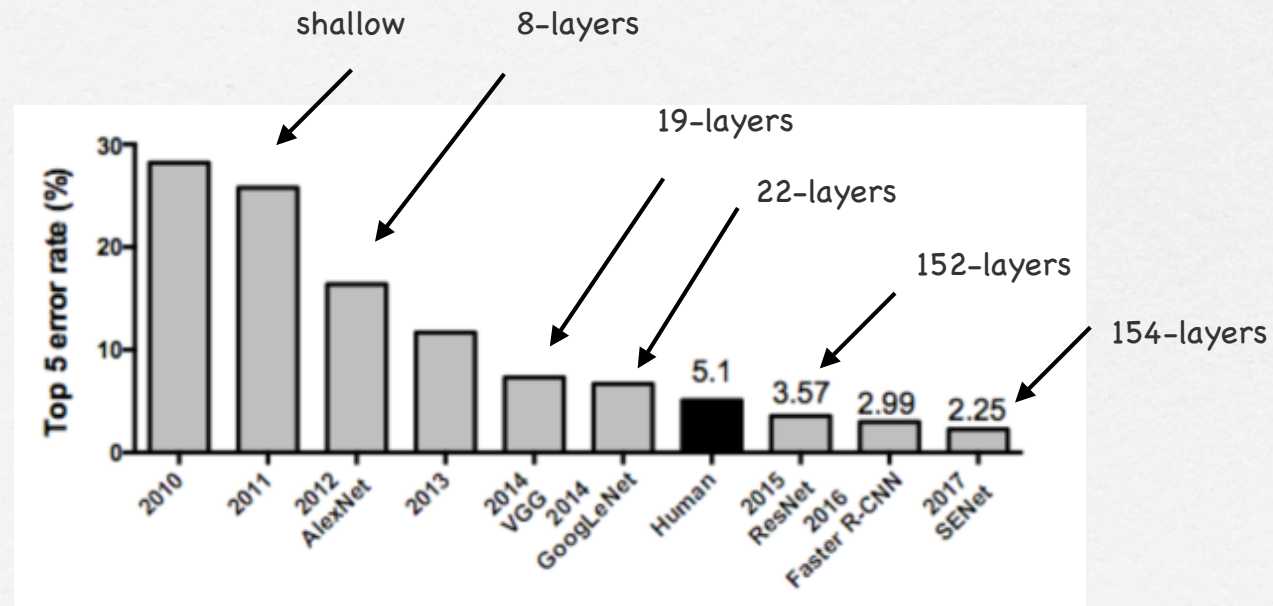


- DenseNets reach similar performance than ResNets with many less parameters (left) and require less computation (right).

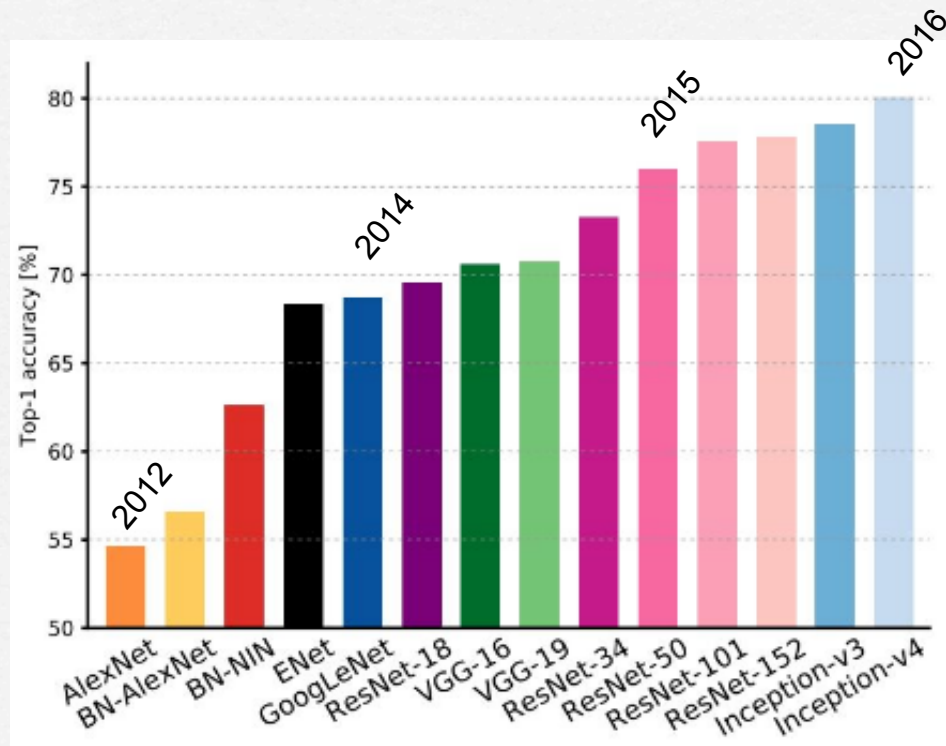
The alchemy of Machine Learning



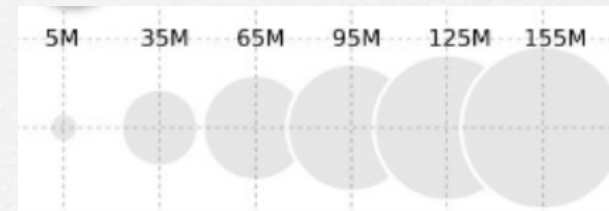
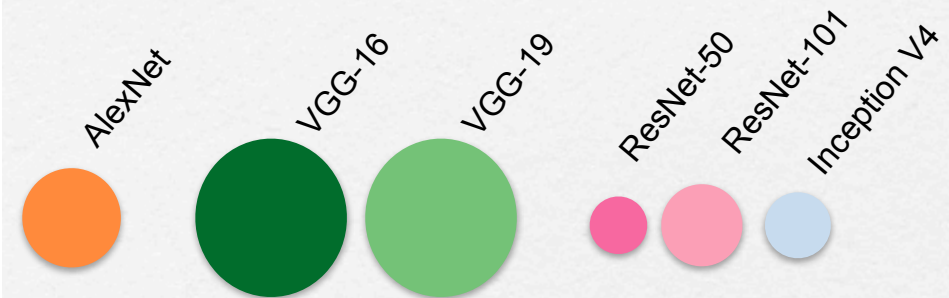
deep neural network performance on
ImageNet



Evolution of CNN architectures



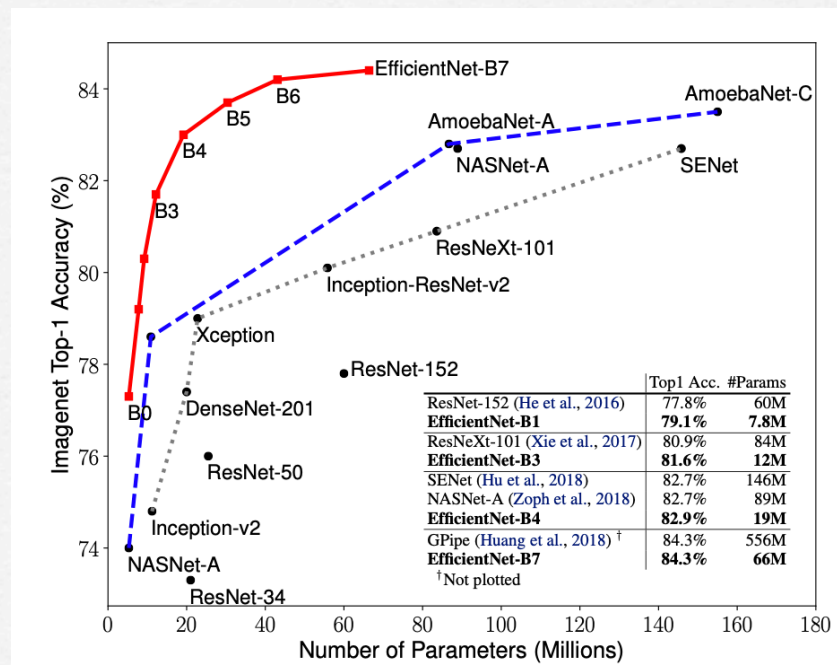
Canziani et al., 2017



Number of parameters

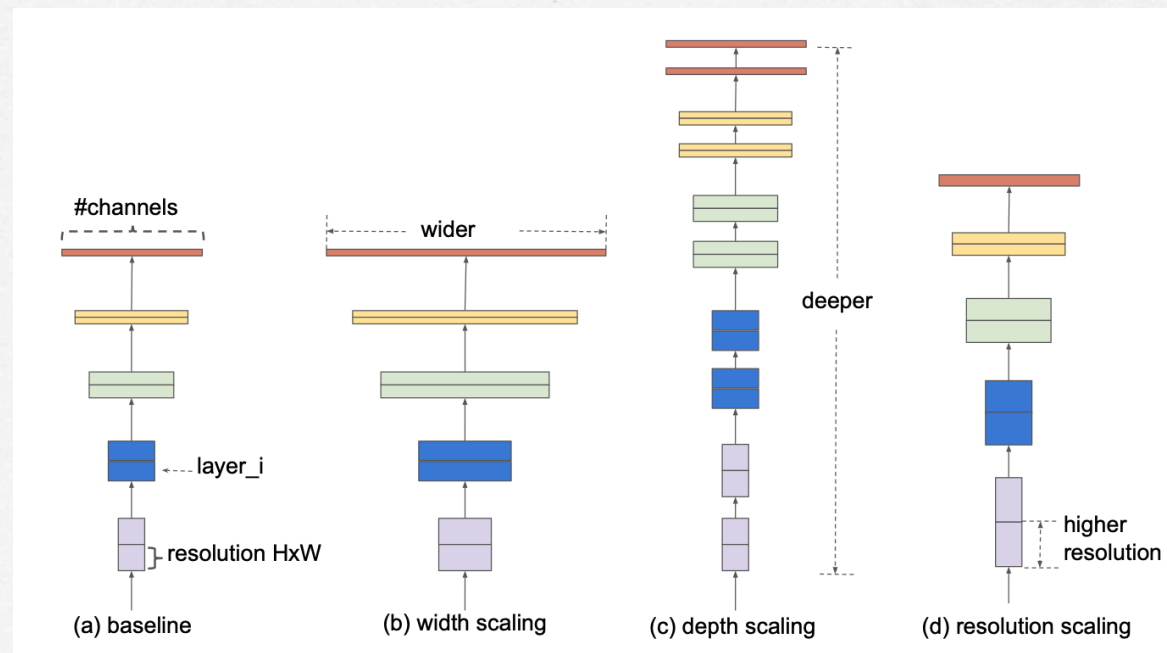
EfficientNets (2020)

- Mingxing Tan and Quoc Le from the Brain Team of Google Research systematically studied model scaling and identified that carefully balancing network depth, width, and resolution can lead to better performances.



EfficientNets scaling (1)

- Example: assume ResNet-18 as a reference architecture processing images of 224×224 resolution: $d=1.0$ (depth) $r=1.0$ (resolution)
- A new architecture with $d=2.0$ and $r=1.3$ will be a ResNet-36 processing 299×299 images. The width factor modifies the number of filters per layer.

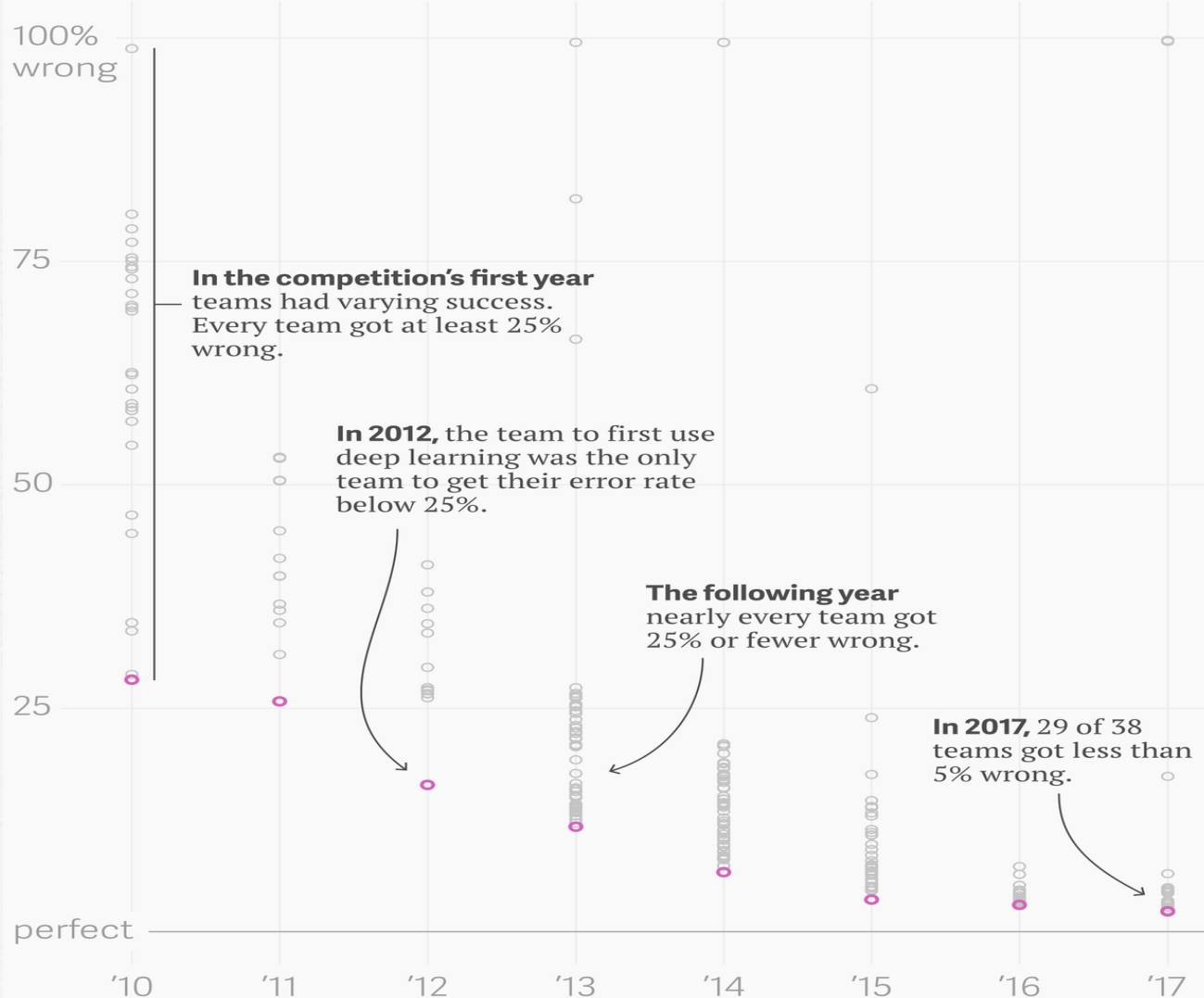


EfficientNets scaling (2)

- Tan et Le found that they could define a compound scaling factor Φ : $d = \alpha^\Phi$, $w = \beta^\Phi$, $r = \gamma^\Phi$, where $\alpha, \beta, \gamma \geq 1$
- Given that the computational load (FLOPS) of a deep network is proportional to d , w^2 and r^2 . Thus, using the constraint $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$, the FLOPS get proportional to 2^Φ .
- **EfficientNet-B0** performs better than ResNet-50 (~5x less parameters; 11x less FLOPS) and DenseNet-169 (~2.5x less parameters; 9x less FLOPS); **EfficientNet-B1** performs better than ResNet-152 (~7.5x less parameters; 16x less FLOPS), etc.

« Democratization » of Deep Learning

ImageNet Large Scale Visual Recognition Challenge results



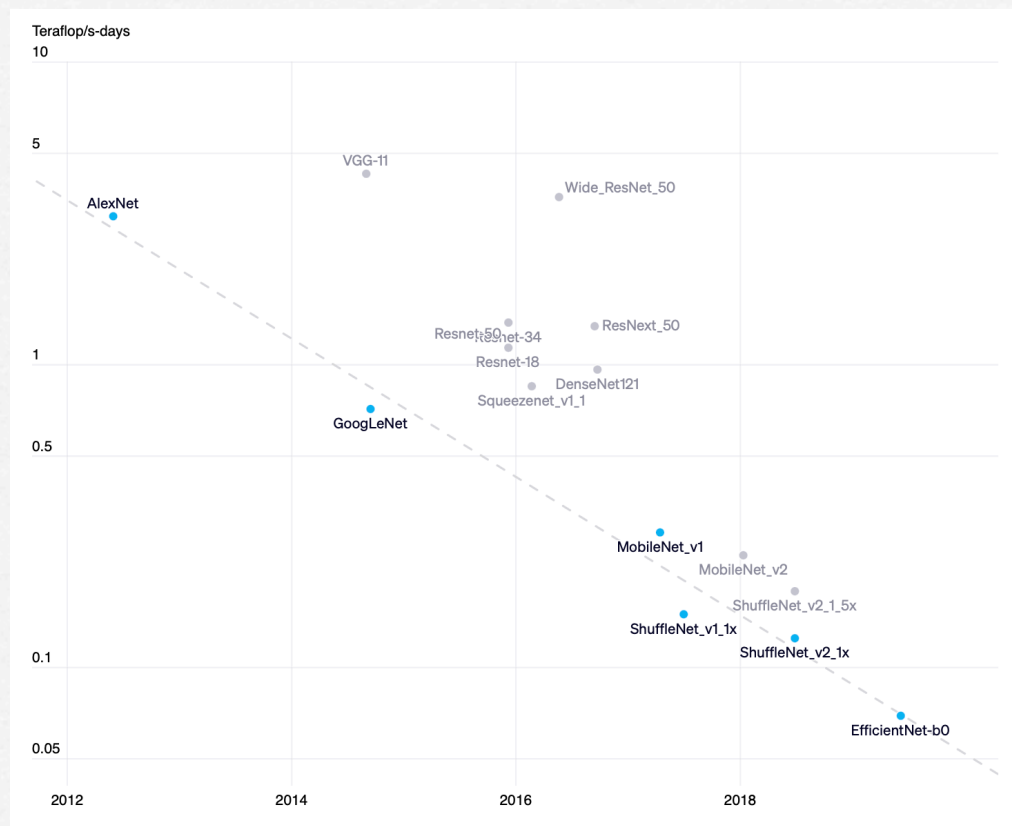
In the competition's first year teams had varying success. Every team got at least 25% wrong.

In 2012, the team to first use deep learning was the only team to get their error rate below 25%.

The following year nearly every team got 25% or fewer wrong.

In 2017, 29 of 38 teams got less than 5% wrong.

Training efficiency improvement



- Besides the improvement of hardware (Moore's law), the amount of compute needed to train a neural net has decreased by 2 every 16 months (thanks to better hyperparameter search, architectures and training procedures).