



浙江大学爱丁堡大学联合学院

ZJU-UoE Institute

## Lecture 13 - CNN architectures

---

Nicola Romanò - [nicola.romano@ed.ac.uk](mailto:nicola.romano@ed.ac.uk)

# Introduction

Today we are going to discuss a few classic papers using CNN for image analysis.

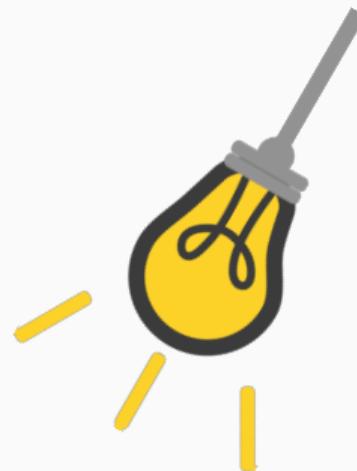
We will analyse the following architectures:

- LeNET-5
- AlexNet
- VGG
- GoogLeNet
- ResNet

The idea is to get some **intuition** about these architectures and how they work.

## Learning objectives

- Describe commonly used patterns in CNN architectures
- Describe and explain the advantages of different CNN architectures



## LeNET-5

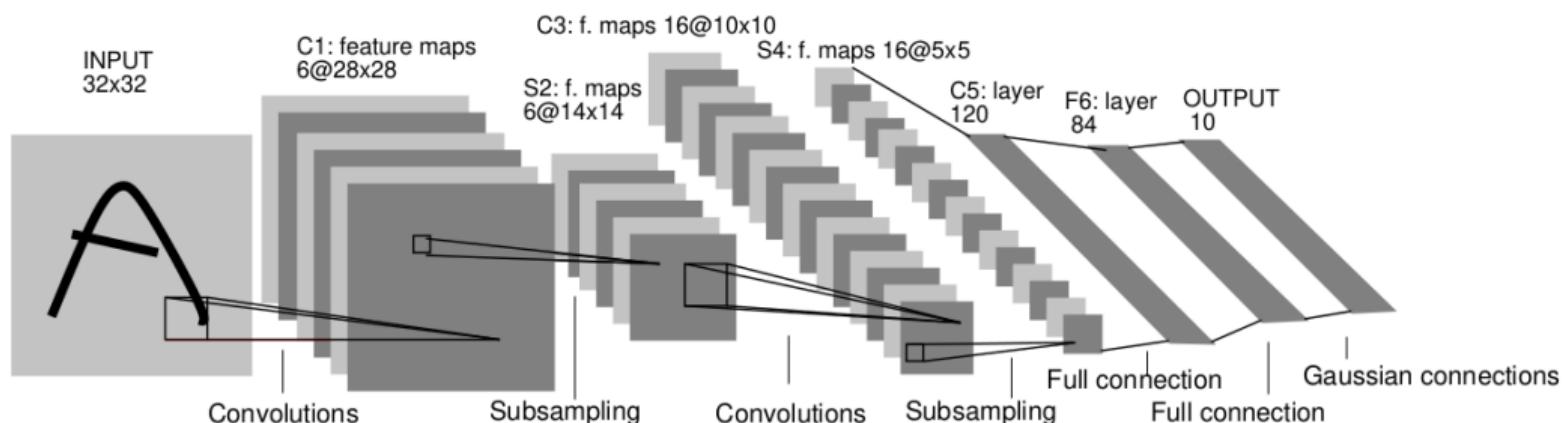
---

## LeNET-5

- "Gradient Based Learning Applied to Document Recognition", Yann LeCun et al. 1998
- A seminal paper describing the use of CNN in image analysis
- Simple architecture with convolutional layers, average pooling and fully-connected layers
- Task: recognition of handwritten digits to be used for processing of bank cheques

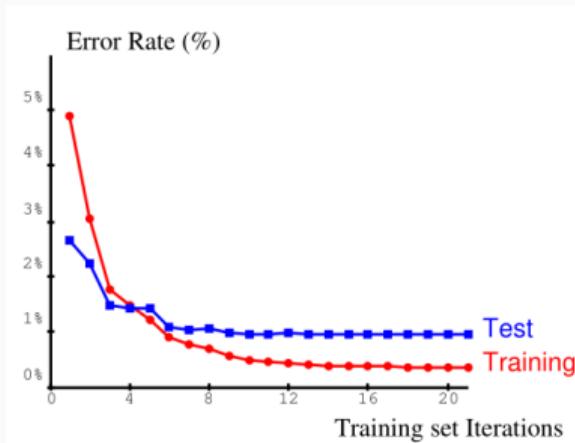
# Gradient-Based Learning Applied to Document Recognition

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner



# LeNet-5 training

- Trained on MNIST dataset (60000 handwritten digits)
- Error rate: <1%



## LeNet-5 take home points

- A simple architecture with convolutional layers, average pooling and fully-connected layers
- Introduced the  $[\text{Conv} + \text{Pool}]_n + \text{FC}$  pattern
- This is mostly interesting from a historical perspective, not really used nowadays.

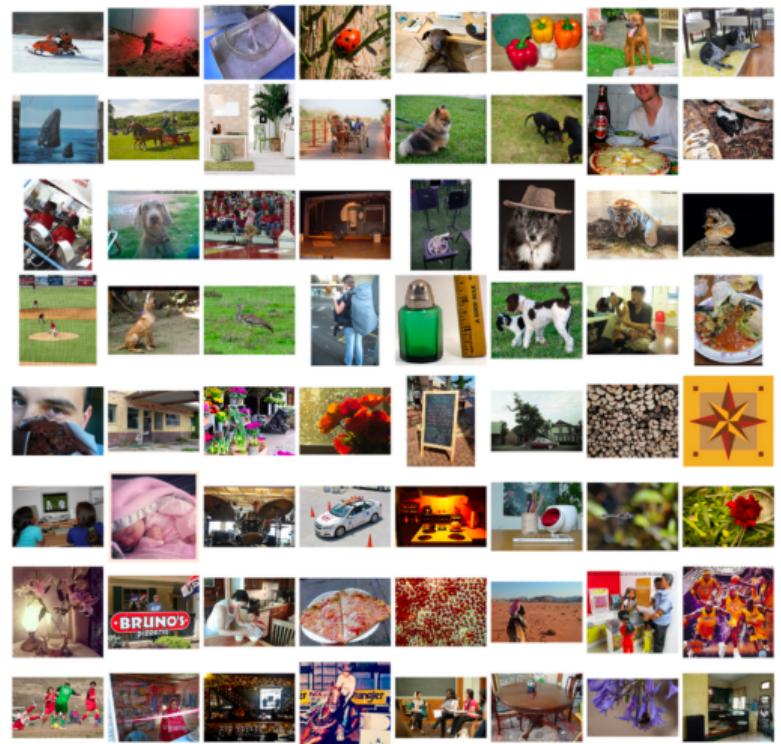
## AlexNet

---

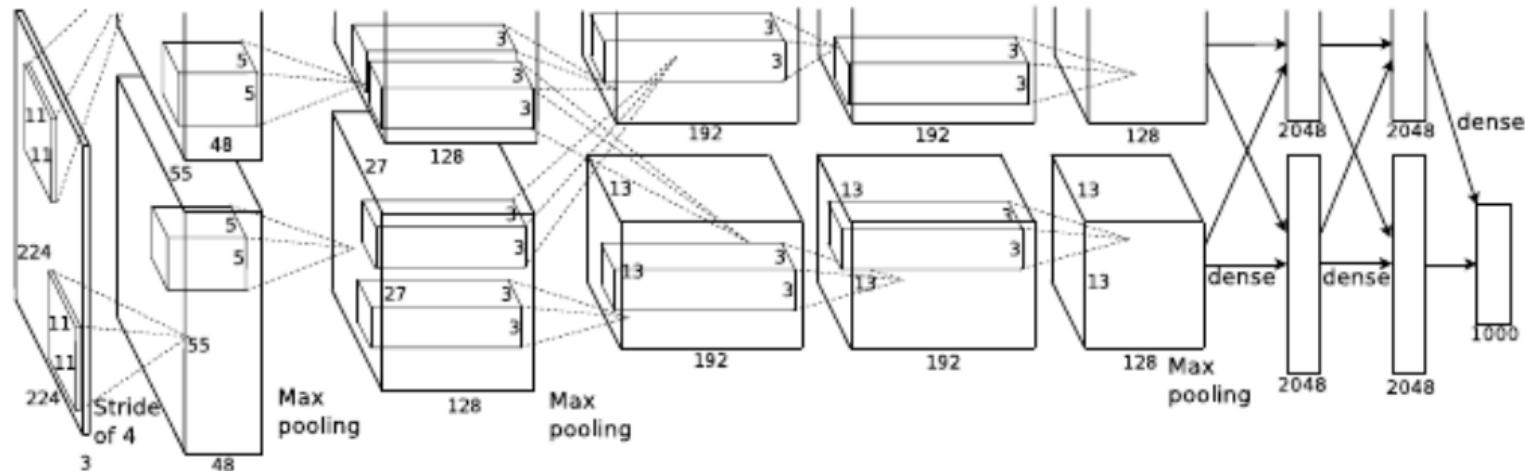
- "ImageNet Classification with Deep Convolutional Neural Networks", Alex Krizhevsky et al. 2012
- Widely considered as one of the most influential papers that boosted research in CNN for image analysis
- Similar architecture to LeNet-5, but with more convolutional layers
- Much bigger network (LeNet-5 60k parameters, AlexNet 60M parameters)
- Winner of the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2012.

# The ImageNet Large Scale Visual Recognition Challenge

- ImageNet is a database of images of various objects, used for training and testing deep neural networks.
- Introduced in Deng et al., 2009 - ImageNet: A large-scale hierarchical image database
- It contains >14 million images of various objects, labelled with >20000 classes.
- The ILSVRC is a competition to define new algorithms for image classification.
- ILSVRC uses a subset of ImageNet, containing 1000 classes and 1.3M training images, 50k validation images and 100k test images.



# ImageNet Classification with Deep Convolutional Neural Networks

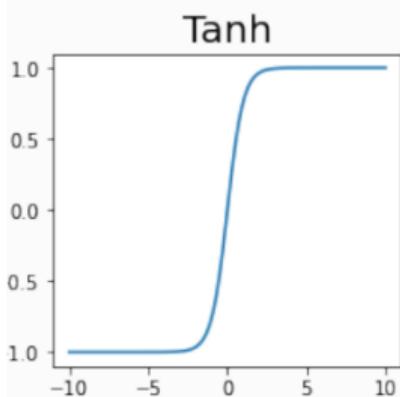


## Improvements in AlexNet

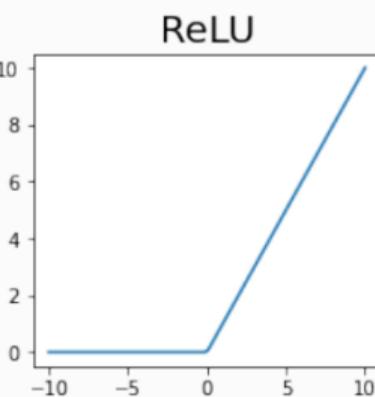
AlexNet introduced a number of improvements:

- ReLU activation function (instead of sigmoid/tanh)
- Dropout
- Softmax activation function in the last layer
- Data augmentation
- Local response normalization
- Training on multiple GPUs

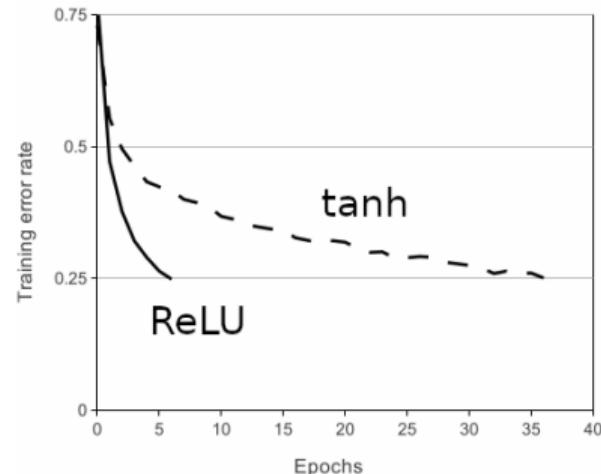
## ReLU activation function



$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



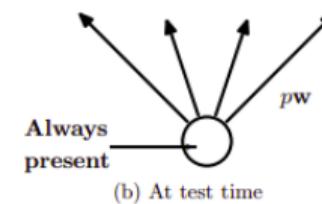
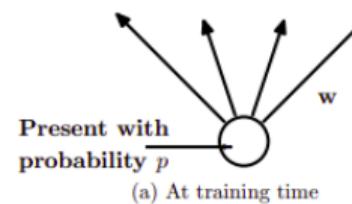
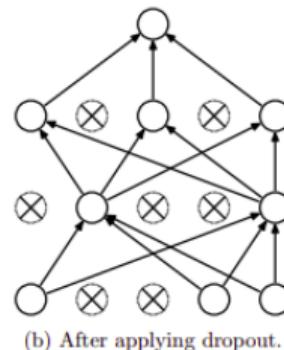
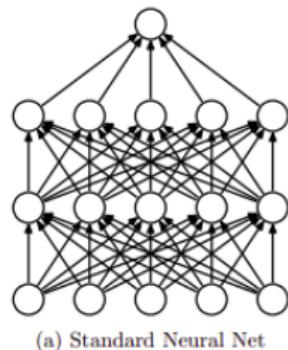
$$\text{ReLU}(x) = \max(0, x)$$



Introduction of the ReLU activation function (instead of sigmoid/tanh) allowed to train much faster.

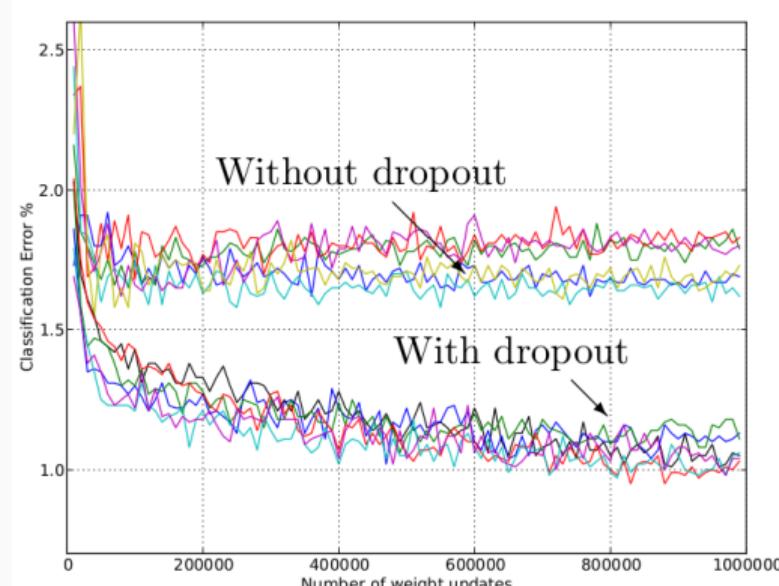
# Dropout

- A type of "regularization" technique, used to prevent overfitting
- A random subset of the weights is set to zero at each training step.
- Originally introduced in "Dropout: A Simple Way to Prevent Neural Networks from Overfitting", Srivastava et al. 2014



# Dropout

- A type of "regularization" technique, used to prevent overfitting
- A random subset of the weights is set to zero at each training step.
- Originally introduced in "Dropout: A Simple Way to Prevent Neural Networks from Overfitting", Srivastava et al. 2014



## The softmax activation function

- Generalization of the logistic function to multiple classes

## The softmax activation function

- Generalization of the logistic function to multiple classes
- Common choice for classification problems in ANNs.

## The softmax activation function

- Generalization of the logistic function to multiple classes
- Common choice for classification problems in ANNs.
- Defined as

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

## The softmax activation function

- Generalization of the logistic function to multiple classes
- Common choice for classification problems in ANNs.
- Defined as

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

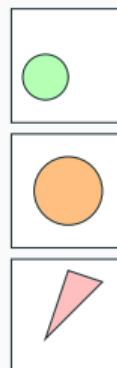
- It is used in the last layer of a CNN to compute the probability of each class.

## The softmax activation function

- Generalization of the logistic function to multiple classes
- Common choice for classification problems in ANNs.
- Defined as

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

- It is used in the last layer of a CNN to compute the probability of each class.



Feedforward  
propagation

Circle	Square	Triangle
8	4	1
9	3	2
5	4	10

Softmax

Circle	Square	Triangle
0.98	0.018	0.002
0.99	0.002	0.008
0.006	0.002	0.992

## AlexNet take home points

- Similar architecture to LeNet-5, but with more convolutional layers
- **ReLU activation functions** - faster computation, more efficient training
- **Dropout** to prevent overfitting
- Training on multiple GPUs

**VGG**

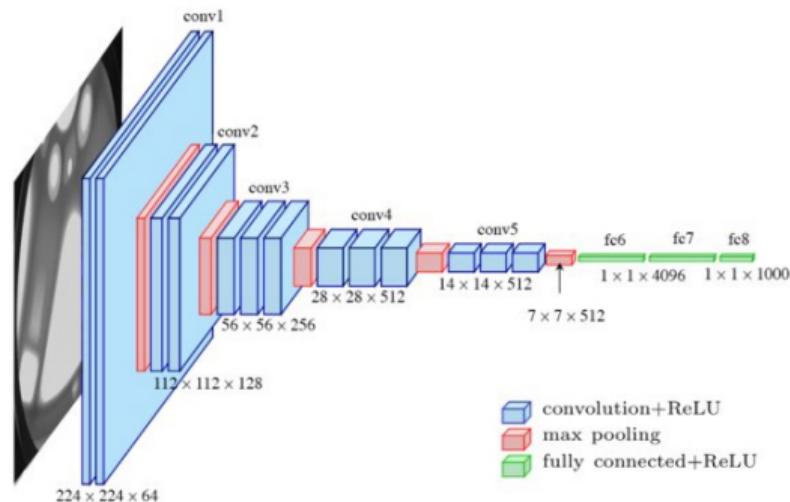
---

- "Very Deep Convolutional Networks for Large-Scale Image Recognition", Karen Simonyan and Andrew Zisserman, 2015
- Very popular architecture for image analysis
- Very deep network, with 16 layers (VGG-16) or 19 layers (VGG-19). 130M parameters
- Winner of ILSVRC in 2015.
- VGG-19 is slightly better, but more computationally expensive (in practice VGG-16 more common).

## VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION

Karen Simonyan\* & Andrew Zisserman<sup>†</sup>

Visual Geometry Group, Department of Engineering Science, University of Oxford  
`{karen, az}@robots.ox.ac.uk`



## VGG take home points

- Very deep network, 130M parameters
- Uses small convolutions ( $3 \times 3$ ) with stride 1
- All layers have same configuration (simplified hyperparameter choice)
- $1 \times 1$  convolutions to increase non-linearity

## **GoogLeNet**

---

- "Going Deeper with Convolutions", Szegedy et al. 2014
- Moves away from the structure we've seen so far
- Introduces "Inception" modules
- 12x less parameters than AlexNet but much more accurate!
- Newer versions (Inception v3, v4) have more powerful architectures



## Going deeper with convolutions

Christian Szegedy

Google Inc.

Wei Liu

University of North Carolina, Chapel Hill

Yangqing Jia

Google Inc.

Pierre Sermanet

Google Inc.

Scott Reed

University of Michigan

Dragonir Anguelov

Google Inc.

Dumitru Erhan

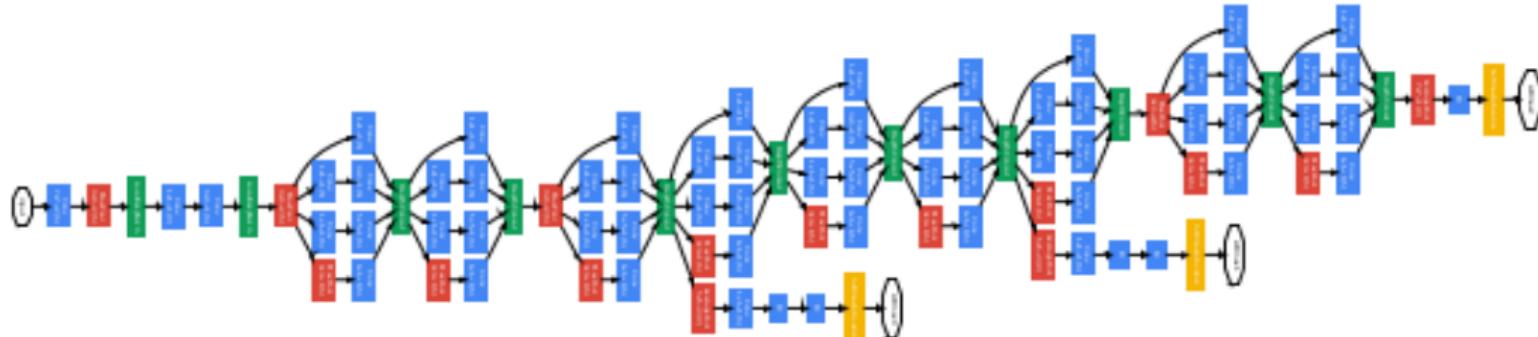
Google Inc.

Vincent Vanhoucke

Google Inc.

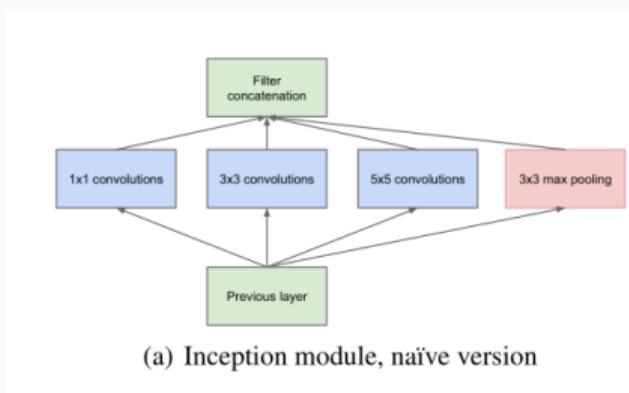
Andrew Rabinovich

Google Inc.



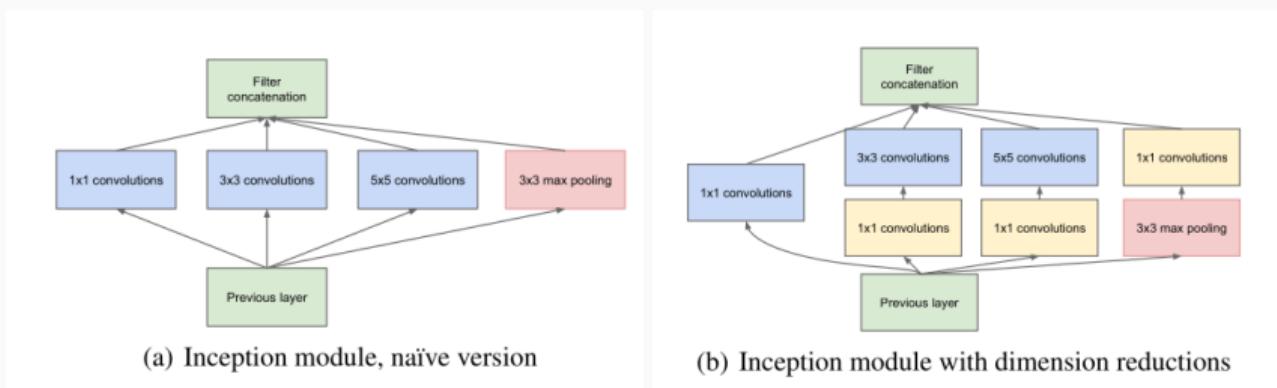
## Inception module

- The inception module is a building block of the GoogLeNet architecture
- It is a combination of convolutions with different kernel sizes
- It is a way to increase the non-linearity of the network
- It is a way to reduce the number of parameters



# Inception module

- The inception module is a building block of the GoogLeNet architecture
- It is a combination of convolutions with different kernel sizes
- It is a way to increase the non-linearity of the network
- It is a way to reduce the number of parameters



## GoogLeNet take home points

- 22 layers
- Heavily relies on  $1 \times 1$  convolutions
- Inception modules allow multi-scale feature extraction
- Drops FC layers
- Extra "side" classifications to improve gradient optimization in earlier layers

## **ResNet**

---

- He 2015 - Deep Residual Learning for Image Recognition
- Tackles the problem of degraded performance in larger networks
- Introduces *skip connections* between layers
- Up to 1000+ layers!

# ResNet architecture

## Deep Residual Learning for Image Recognition

Kaiming He

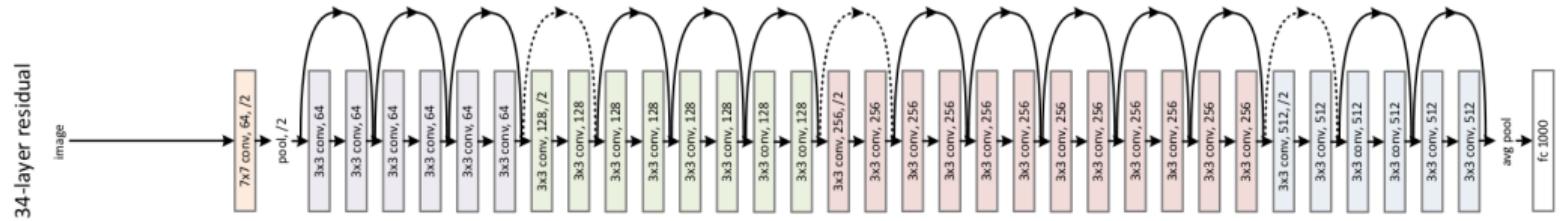
Xiangyu Zhang

Shaoqing Ren

Jian Sun

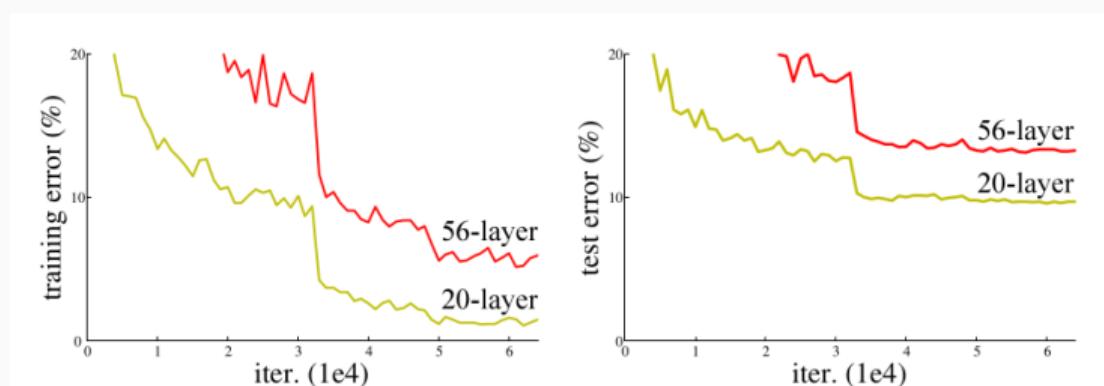
Microsoft Research

{kahe, v-xiangz, v-shren, jiansun}@microsoft.com



## The problem with deep networks

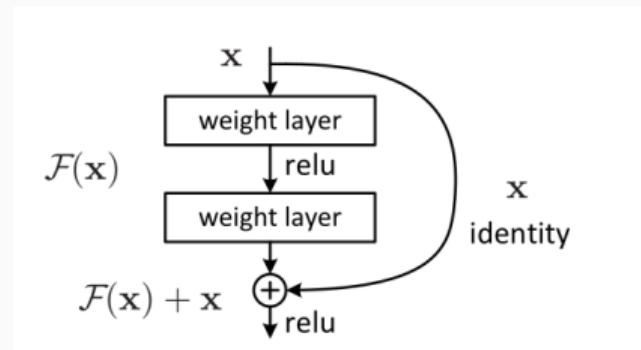
As we add more layers, the network becomes more difficult to train



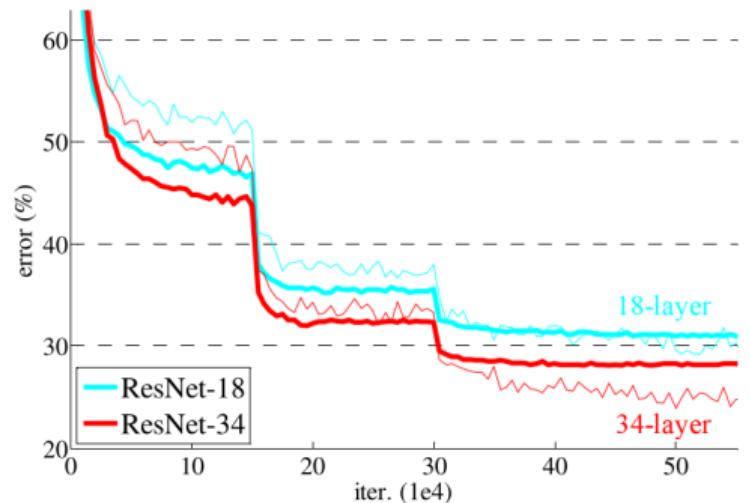
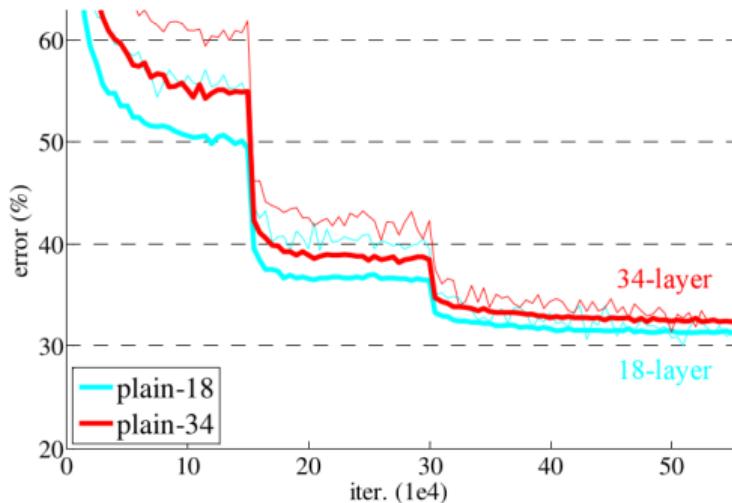
This is surprising because we expect a deeper network to be at least as good as a shallower network. Indeed the extra layers could "simply" learn the identity function  $I(x) = x$ .

## Skip connections

By adding a skip connection, the network can learn the residual function  $F(x) = H(x) - x$ . The intuition is that the residual function is easier to learn than the identity function as you could simply zero out some of the layers.



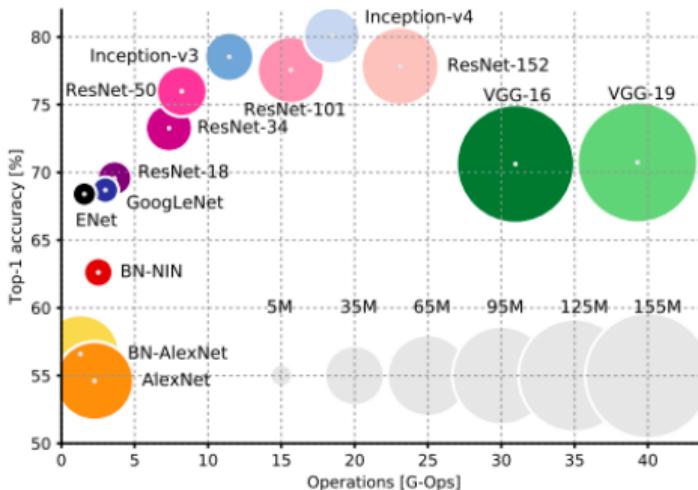
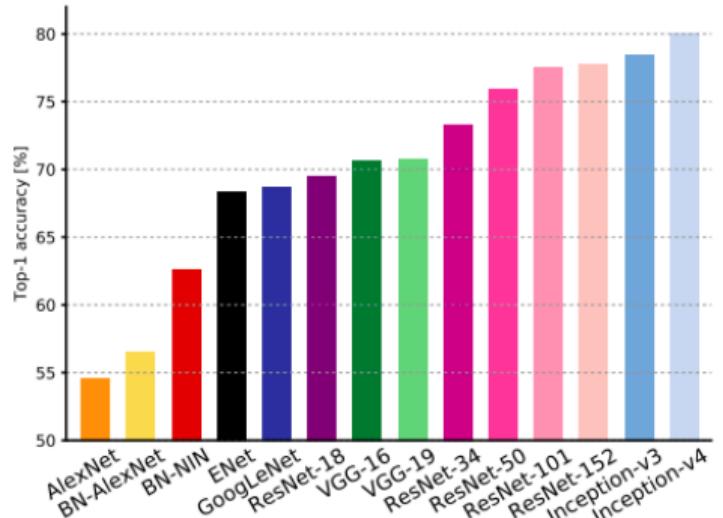
## Skip connections improve performance



## ResNet take home points

- Very deep network (up to 1000+ layers)
- Uses *skip connections* between layers
- Uses *bottleneck* blocks (similar to GoogLeNet)

# Comparison of CNN architectures



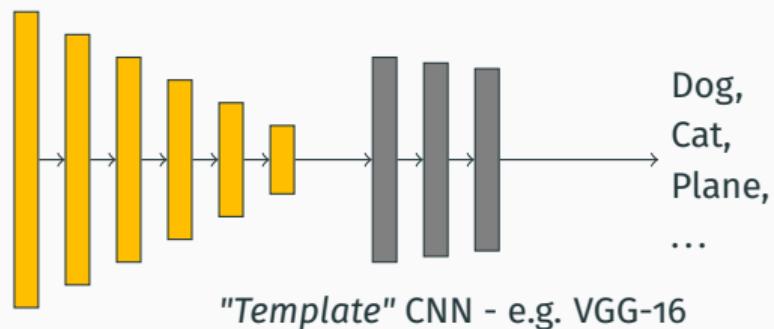
## Transfer learning

- Training a CNN is expensive in terms of time and computer memory.
- Luckily, you can make use of pre-trained models to speed up the training process!
- This process is called **transfer learning**.
- For example, if you wanted to classify images, rather than start from scratch you could begin with VGG-16, or with GoogLeNet and modify them for your needs! (i.e. do not reinvent the wheel!)

# Transfer learning



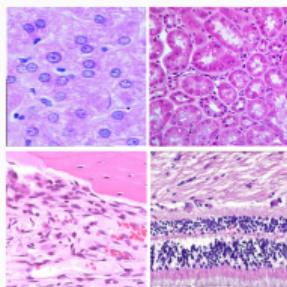
Generic images



# Transfer learning



Generic images



Task-specific images

Conv layers

FC layers

