# Introduction to Deep Learning, part 2

**- Rafał Cycoń, FORNAX**

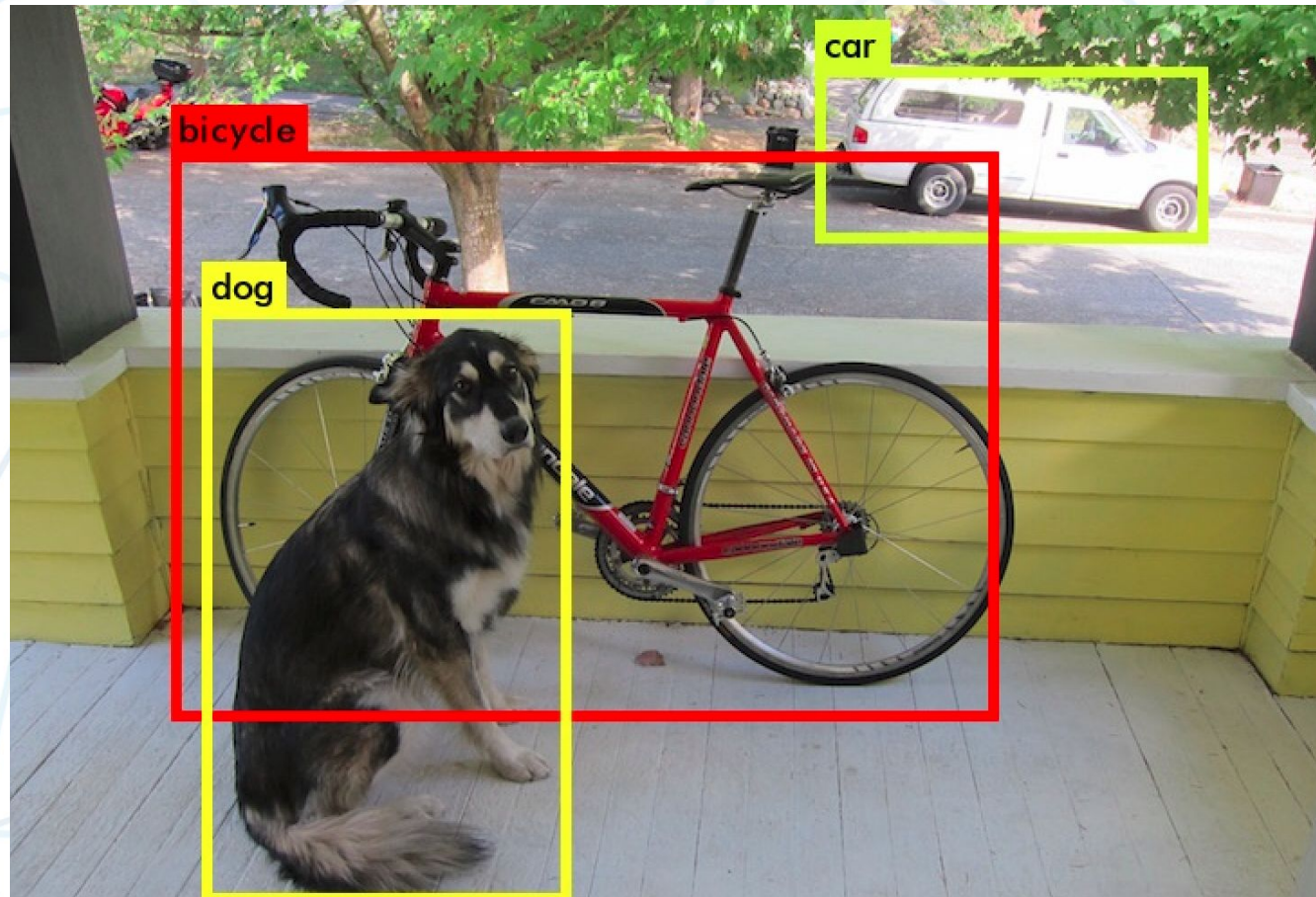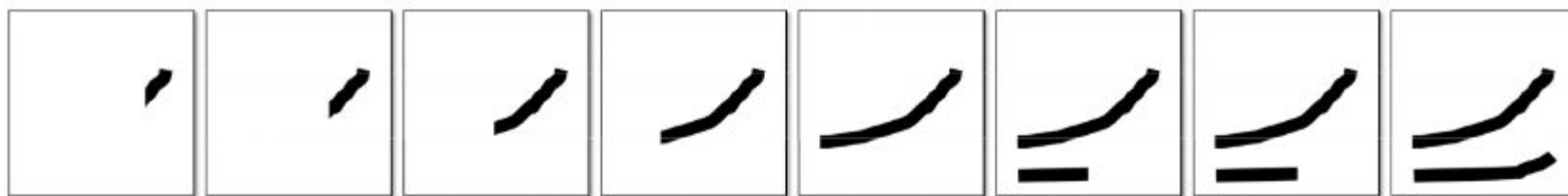# Image classification

# Localization

# Photorealistic style transfer

# Image generation



(a) User constraints $v_g$ at different update steps

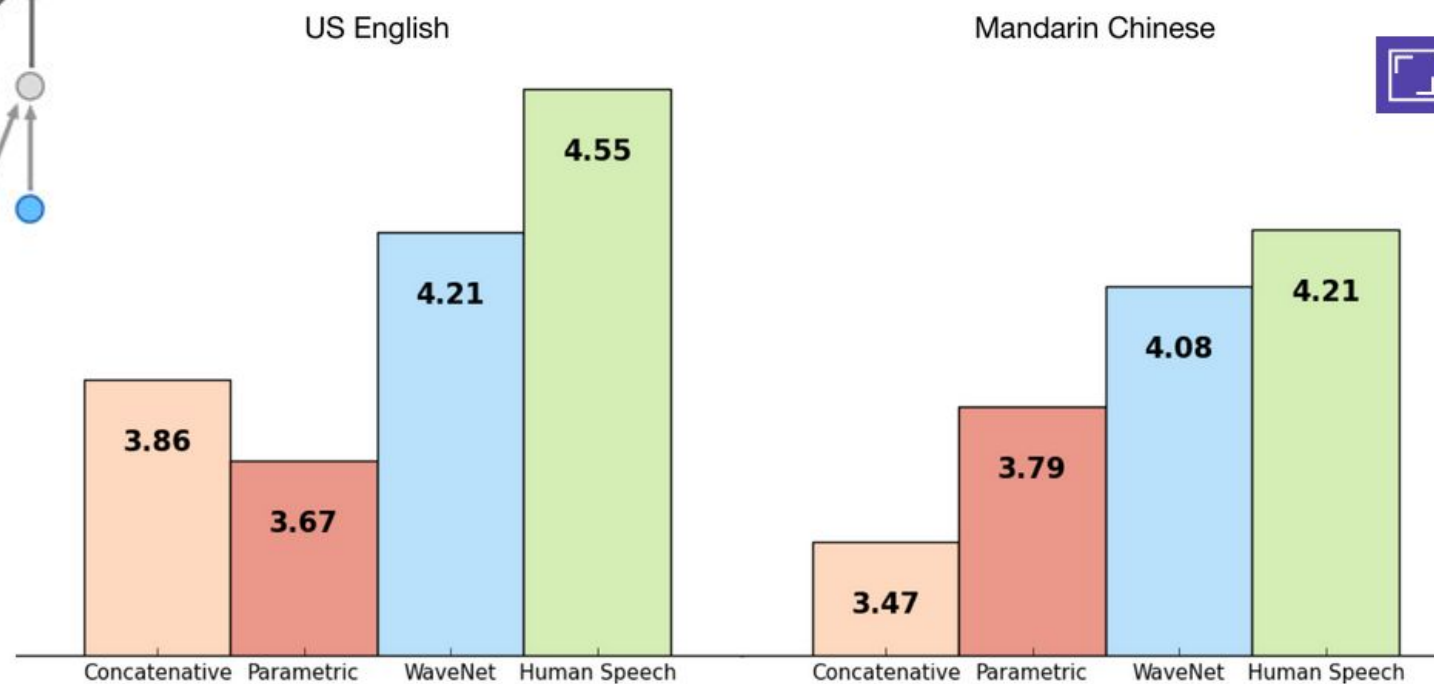$G(z_0)$     (b) Updated images according to user edits     $G(z_1)$

(c) Linear interpolation between $G(z_0)$ and $G(z_1)$

# Sound synthesis



Output

Hidden Layer

Hidden Layer

Hidden Layer

Input

US English

4.55

4.21

3.86

3.67

Concatenative  Parametric  WaveNet  Human Speech

Mandarin Chinese

4.21

4.08

3.79

3.47

Concatenative  Parametric  WaveNet  Human Speech
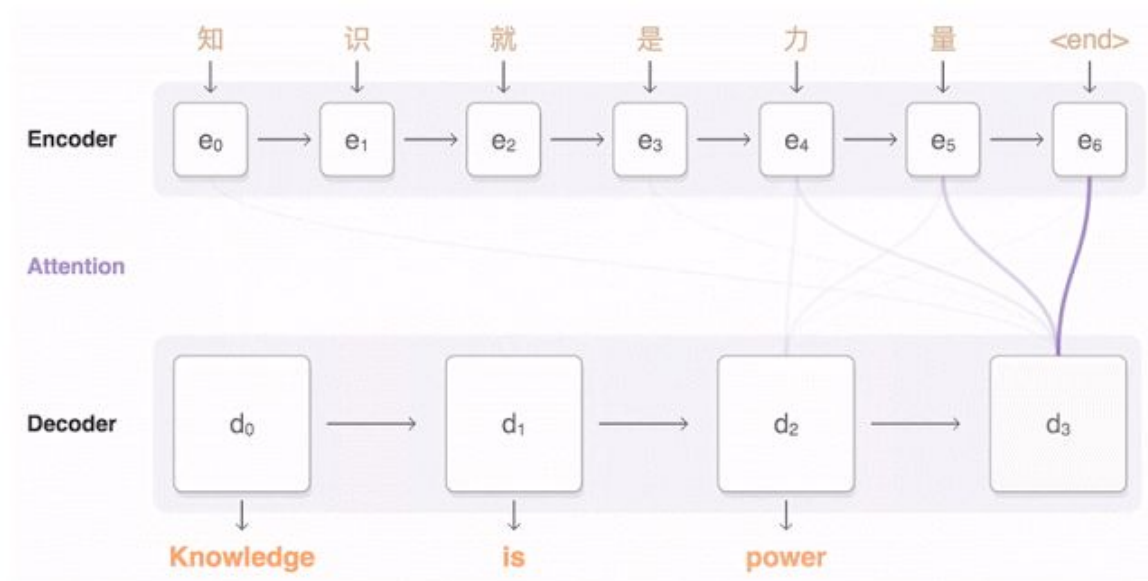
# Translation

# Text generation

PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
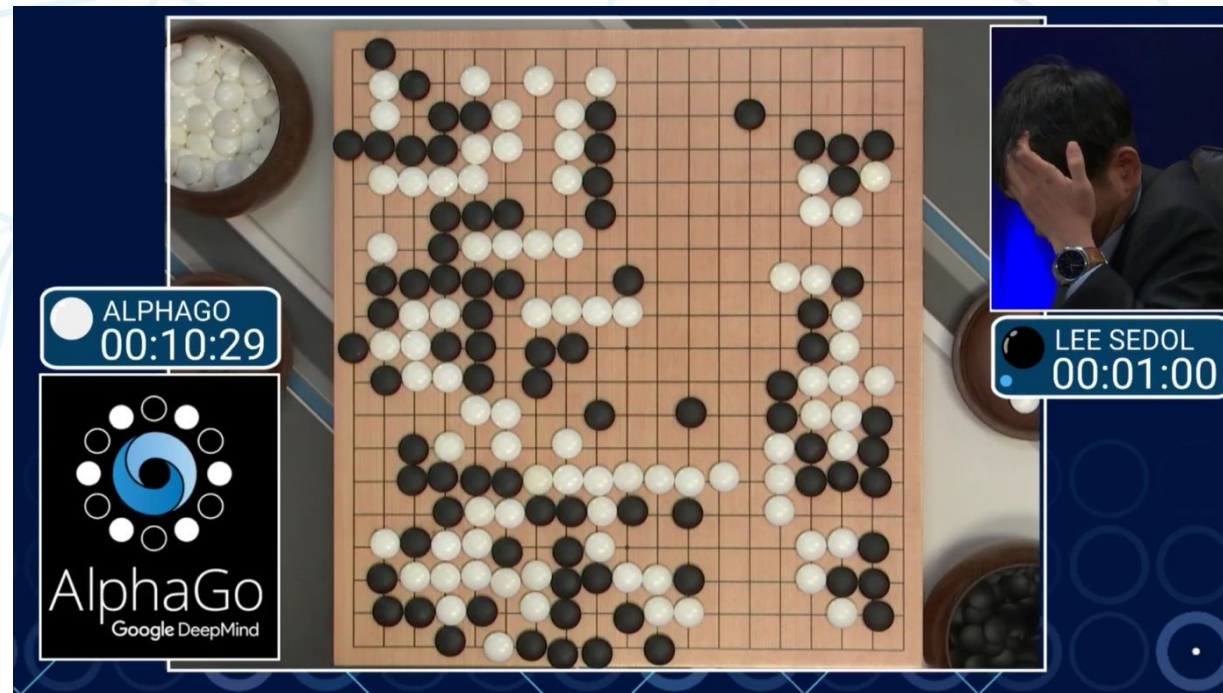Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.
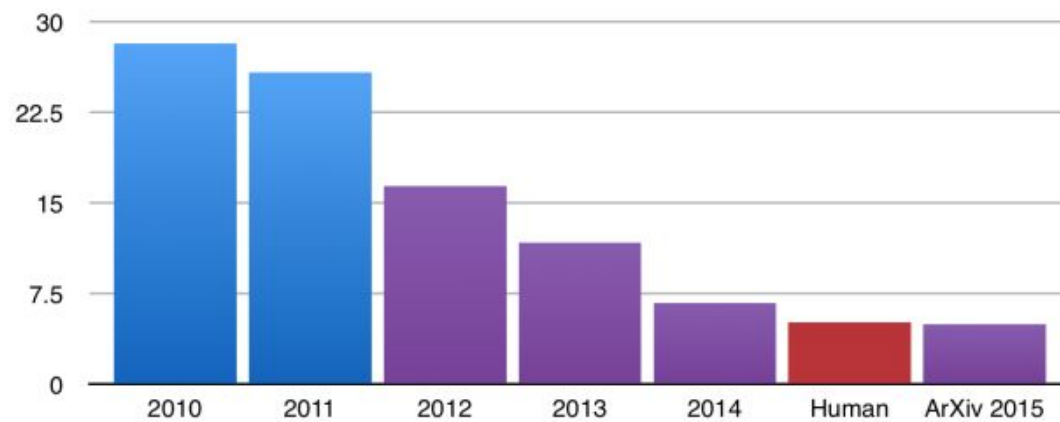
VIOLA:
I'll drink it.

# Reinforcement learning

# ImageNet

# Challenges in image processing



Viewpoint variation

Scale variation

Deformation
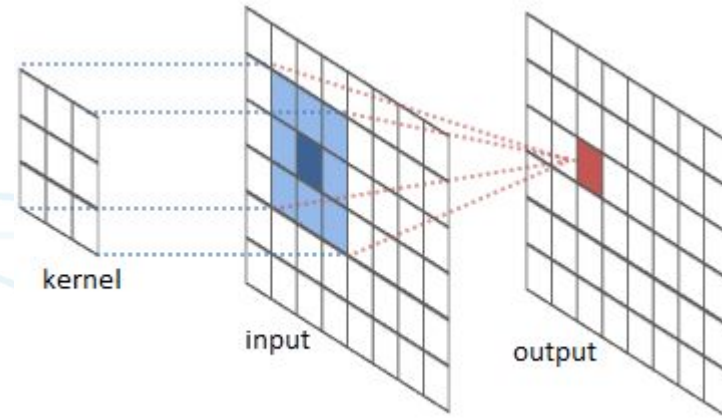
Occlusion

Illumination conditions

Background clutter

Intra-class variation

# Convolution

$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau)\, g(t - \tau)\, d\tau$$

$$= \int_{-\infty}^{\infty} f(t - \tau)\, g(\tau)\, d\tau.$$

$$(f * g)[n] \stackrel{\text{def}}{=} \sum_{m=-\infty}^{\infty} f[m]\, g[n - m]$$

$$= \sum_{m=-\infty}^{\infty} f[n - m]\, g[m].$$



kernel    input    output

# Convolution



Image

Convolved Feature

# Convolution

# Convolution



32x32x3 image
5x5x3 filter

convolve (slide) over all
spatial locations

activation map
feature map

$$w^T x + b$$

# Convolution



32x32x3 image
5x5x3 filter

32

32

3

activation maps

convolve (slide) over all spatial locations

28

28

1

$w^T x + b$

# Convolution Layer



32 × 32 × 3

6 filters, each of size 5x5x3

Convolution Layer

**activation maps**

28 × 28 × 6

# ConvNet

object models

object parts
(combination
of edges)

edges

pixels

WWW.FORNAX.CO

# Pooling

# Pooling



Single depth slice

max pool with 2x2 filters
and stride 2

Discards 75% of info!

WWW.FORNAX.CO

# Convolution



Input

# Nonlinearity



Input Feature Map — Black = negative; white = positive values
→ ReLU →
Rectified Feature Map — Only non-negative values

# Pooling



Pooling

Max

Sum

Only non-negative values

Rectified Feature Map

# AlexNet

# VGG-16

| \multicolumn{6}{c}{ConvNet Configuration} |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| \multicolumn{6}{c}{input ($224 \times 224$ RGB image)} |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| \multicolumn{6}{c}{maxpool} |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| \multicolumn{6}{c}{maxpool} |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| \multicolumn{6}{c}{maxpool} |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| \multicolumn{6}{c}{maxpool} |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| \multicolumn{6}{c}{maxpool} |
| \multicolumn{6}{c}{FC-4096} |
| \multicolumn{6}{c}{FC-4096} |
| \multicolumn{6}{c}{FC-1000} |
| \multicolumn{6}{c}{soft-max} |

# Smaller filters



Input     First Conv     Second Conv

# Smaller filters – what do we gain?

Consider an input with height H, width W, C channels.
Assume stride=1, padding such that H and W do not change on the output.

A single 7x7 filter:
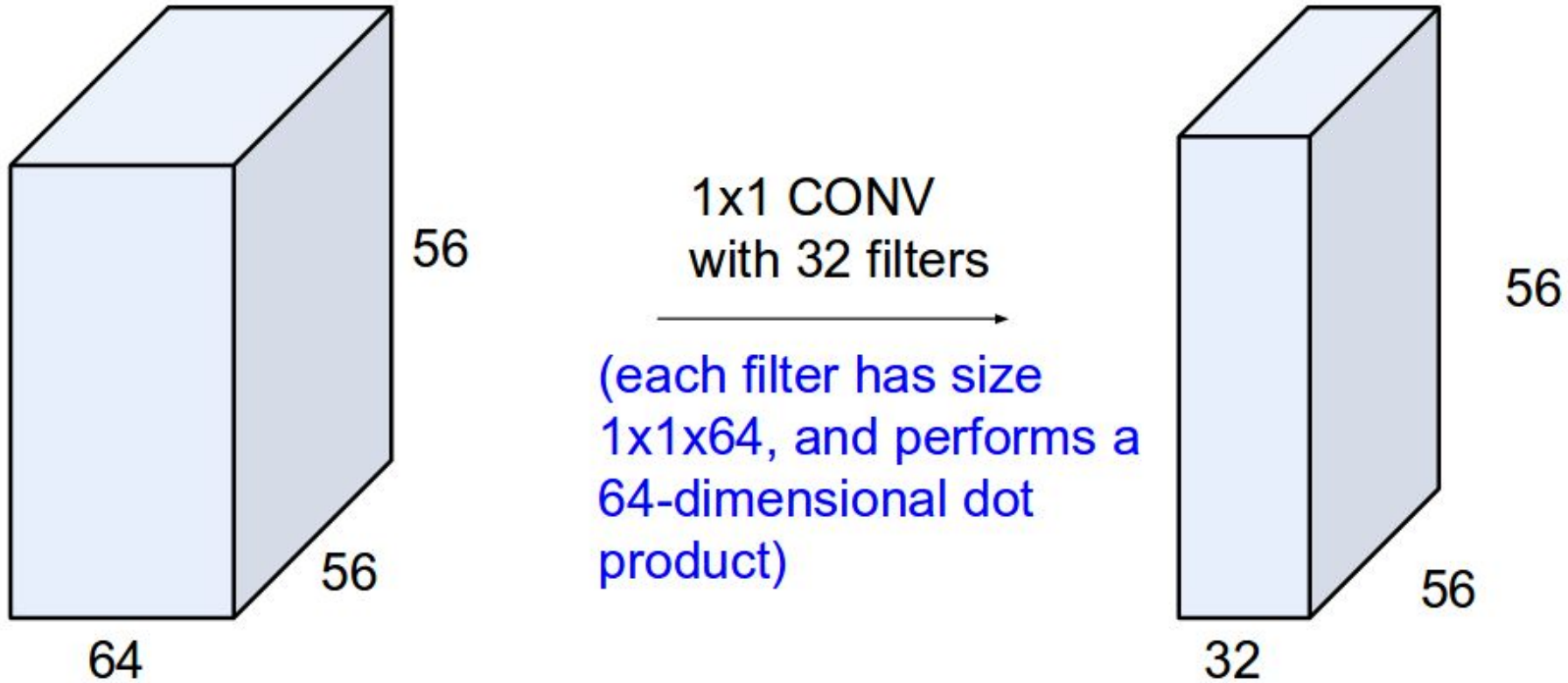- parameters: $C*(7*7*C) =$ **49**$C^2$  (and a bias)
- operations: $(H*W*C) * (7*7*C) =$ **49** $HWC^2$

A stack of three 3x3 filters (same receptive field):
- parameters: **3**$*C*(3*3*C) =$ **27**$C^2$  (and a bias)
- operations: $(H*W*C) *$ **3** $* (3*3*C) =$ **27** $HWC^2$

- Less parameters
- Less computations
- More nonlinearity

# 1x1 filters



1x1 CONV
with 32 filters

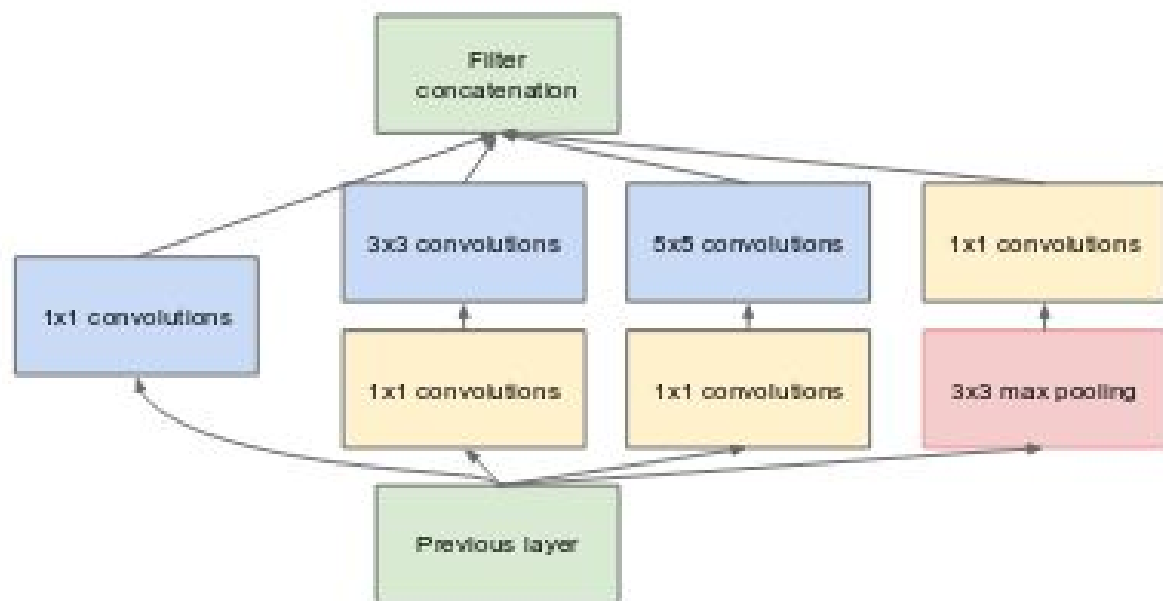(each filter has size 1x1x64, and performs a 64-dimensional dot product)

**Even less params, less computations, more nonlinearity**
**Can be used to increase/decrease number of channels (bottlenecks)**

WWW.FORNAX.CO

# GoogLeNet
https://arxiv.org/abs/1409.4842



(b) Inception module with dimensionality reduction

# Global pooling

**First appeared in "Network in Network"  ([https://arxiv.org/pdf/1312.4400.pdf](https://arxiv.org/pdf/1312.4400.pdf))**

**Facilitates the use of all conv networks (ConvNet can work with all input resolutions!)**
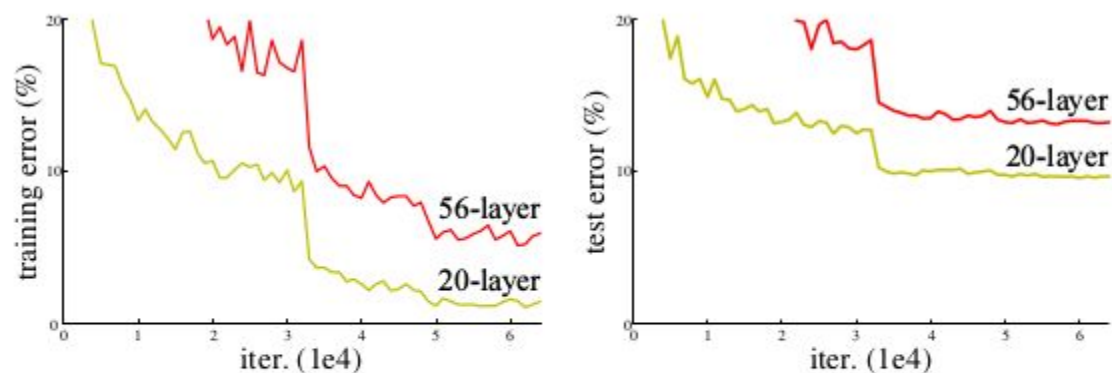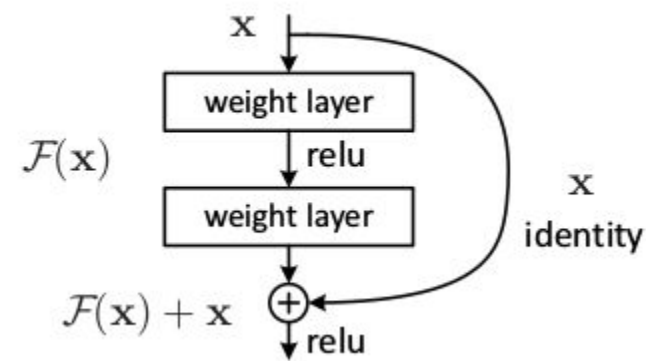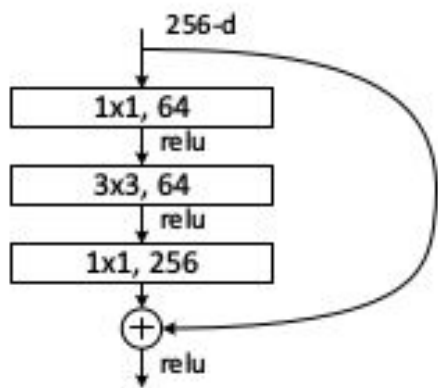
# ResNet

Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer "plain" networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.
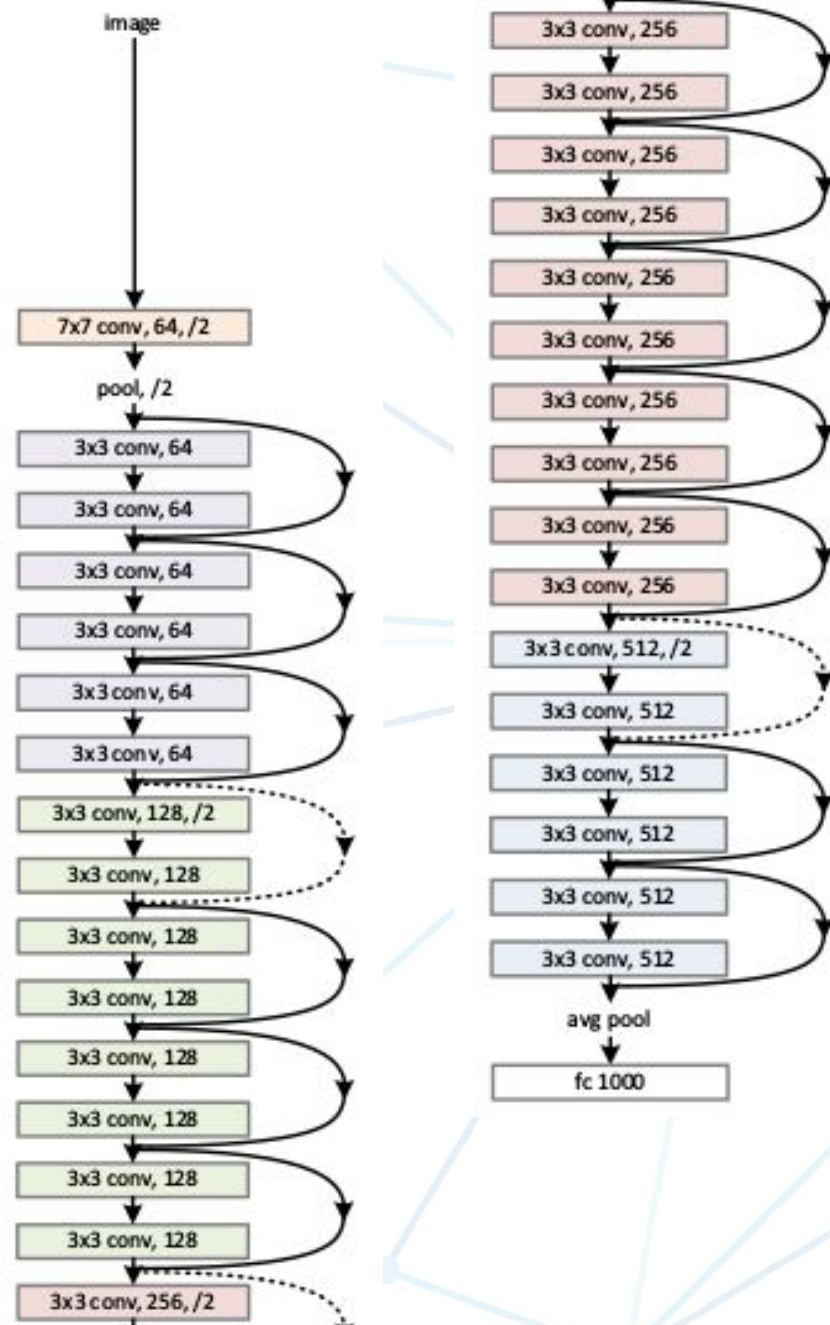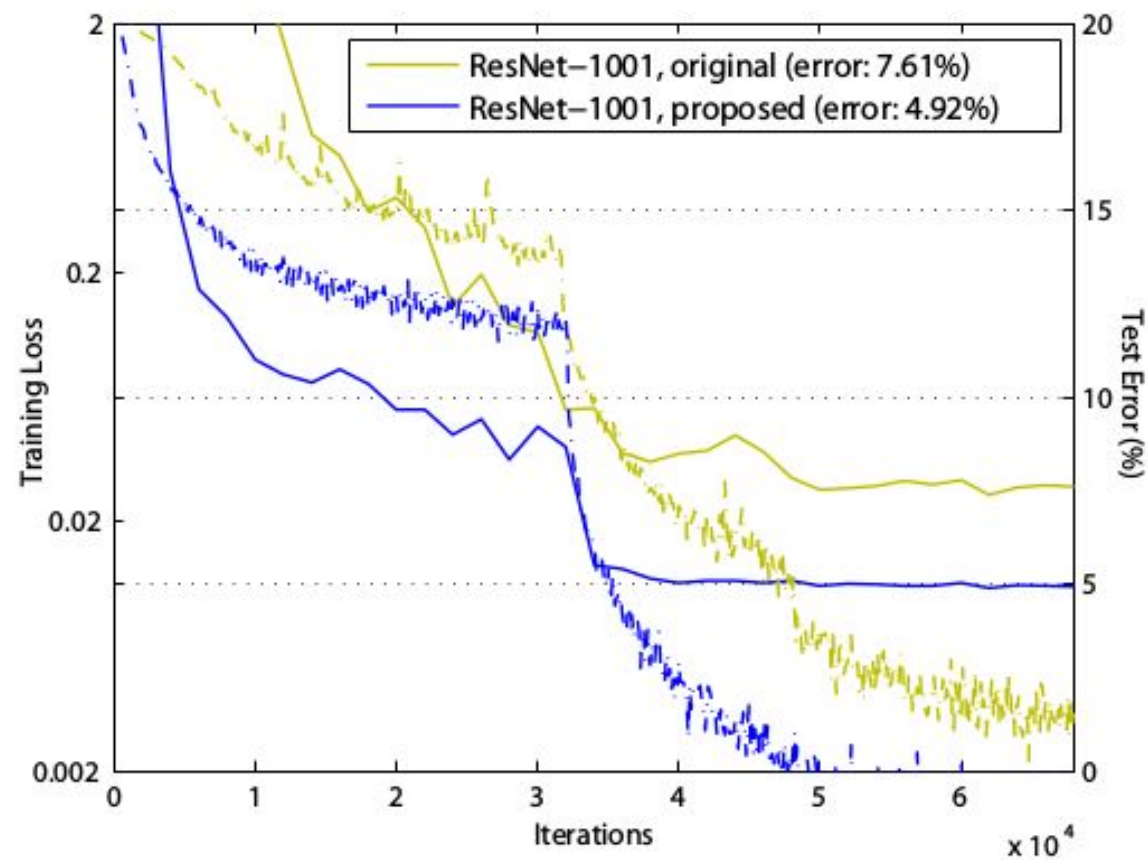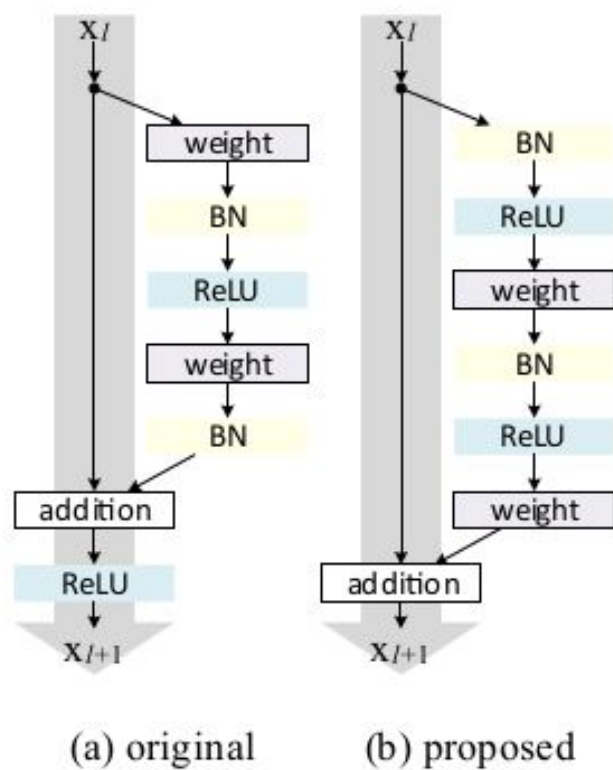
# ResNet

"Bottleneck" building block

# ResNet v2

(a) original   (b) proposed

# WideResNet

|  | depth-$k$ | # params | CIFAR-10 | CIFAR-100 |
|---|---|---|---|---|
| NIN [20] |  |  | 8.81 | 35.67 |
| DSN [19] |  |  | 8.22 | 34.57 |
| FitNet [24] |  |  | 8.39 | 35.04 |
| Highway [28] |  |  | 7.72 | 32.39 |
| ELU [5] |  |  | 6.55 | 24.28 |
| original-ResNet[11] | 110 | 1.7M | 6.43 | 25.16 |
|  | 1202 | 10.2M | 7.93 | 27.82 |
| stoc-depth[14] | 110 | 1.7M | 5.23 | 24.58 |
|  | 1202 | 10.2M | 4.91 | - |
| pre-act-ResNet[13] | 110 | 1.7M | 6.37 | - |
|  | 164 | 1.7M | 5.46 | 24.33 |
|  | 1001 | 10.2M | 4.92(4.64) | 22.71 |
| WRN (ours) | 40-4 | 8.9M | 4.53 | 21.18 |
|  | 16-8 | 11.0M | 4.27 | 20.43 |
|  | 28-10 | 36.5M | **4.00** | **19.25** |

| Model | top-1 err, % | top-5 err, % | #params | time/batch 16 |
|---|---|---|---|---|
| ResNet-50 | 24.01 | 7.02 | 25.6M | 49 |
| ResNet-101 | 22.44 | 6.21 | 44.5M | 82 |
| ResNet-152 | 22.16 | 6.16 | 60.2M | 115 |
| **WRN-50-2-bottleneck** | 21.9 | 6.03 | 68.9M | 93 |
| pre-ResNet-200 | 21.66 | 5.79 | 64.7M | 154 |

# Simple tips

**ConvNets need lots of data**
> **Use data augmentation or transfer learning**

**Lots of parameters = overfitting**
> **Regularize! L2, dropout, BatchNorm**

**Use residual blocks and bottleneck convs (e.g. 1x1 -> 3x3 -> 1x1)**

**Global pooling prior to softmax (or any other dense layer)**
> **Gives a possibility to process images of arbitrary resolution**

**We recommend using Adam as a state-of-the-art optimization algorithm**

# Transfer learning

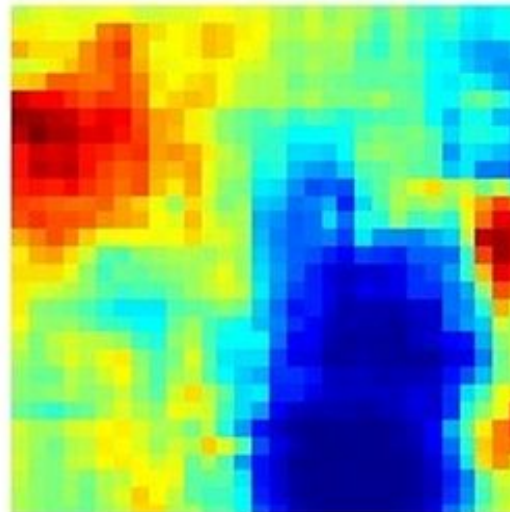Get a model trained on a similar dataset.

Got "lots" of data?
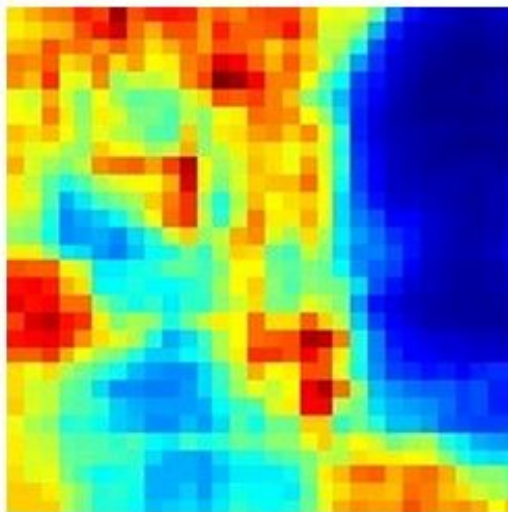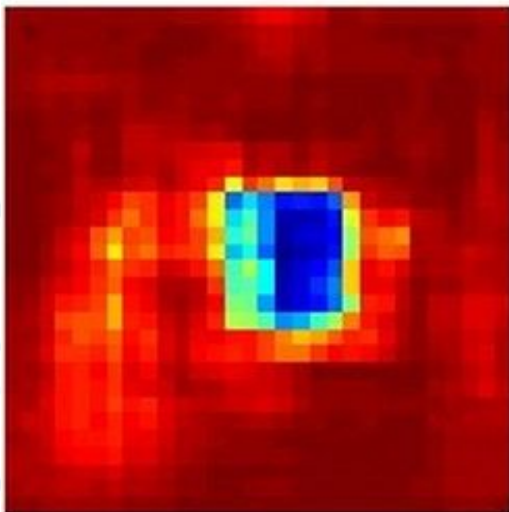**Yes!**
> Fine-tune last few layers, freeze the rest

**No…**
> Train a (simple) classifier on network output and/or on features (i.e. filter maps) from different layers

# Debugging



True Label: Pomeranian

True Label: Car Wheel

True Label: Afghan Hound

**Rafał Cycoń**
**Co-founder, CTO**

**+48 607 697 054**
rafal.cycon@fornax.co

**WWW.FORNAX.CO**