# Neural Networks and Deep Learning

## Architectures for Image Classification

Slides credit: Pau Rodriguez, Ariadna Romero

# What´s in a scene



| Type of environment | Indoor |
|---|---|
| Scene categories | martial_arts_gym (0.204) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy |
| Informative región for predicting the top category |  |

*Try it yourself:* *http://places2.csail.mit.edu/demo.html*

# What´s in a scene





| Type of environment | Indoor | Outdoor |
|---|---|---|
| Scene categories | martial_arts_gym (0.204) | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy | natural light, man-made, open area, no horizon, sunny, glass, vegetation, vertical components, leaves |
| Informative región for predicting the top category |  |  |

*Try it yourself: http://places2.csail.mit.edu/demo.html*

# What´s in a scene



Binary classification

| Type of environment | Indoor | Outdoor |
|---|---|---|
| Scene categories | martial_arts_gym (0.204) | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy | natural light, man-made, open area, no horizon, sunny, glass, vegetation, vertical components, leaves |
| Informative región for predicting the top category |  |  |

*Try it yourself: http://places2.csail.mit.edu/demo.html*

# What´s in a scene



Multi-class classification

| Type of environment | Indoor | Outdoor |
|---|---|---|
| Scene categories | martial_arts_gym (0.204) | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy | natural light, man-made, open area, no horizon, sunny, glass, vegetation, vertical components, leaves |
| Informative región for predicting the top category | | |

# What´s in a scene

Multi-label classification

| Type of environment | Indoor | Outdoor |
|---|---|---|
| Scene categories | martial_arts_gym (0.204) | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy | natural light, man-made, open area, no horizon, sunny, glass, vegetation, vertical components, leaves |
| Informative región for predicting the top category | | |

*Try it yourself:*

# What´s in a scene





Understanding / visualising CNNs

| Type of environment | Indoor | Outdoor |
|---|---|---|
| Scene categories | martial_arts_gym (0.204) | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| Scene attributes | no horizon, enclosed area, man-made, cloth, natural light, competing, exercise, sports, glossy | natural light, man-made, open area, no horizon, sunny, glass, vegetation, vertical components, leaves |
| Informative región for predicting the top category |  |  |

*Try it yourself: http://places2.csail.mit.edu/demo.html*

Image Classification

# PROBLEM METRICS AND DATASETS

# Evaluation Metrics



$$\text{Accuracy} = \frac{\# \; correct \; predictions}{\# \; total \; predictions}$$

$$\text{Misclassification Error} = 1 - Accuracy$$

| Scene categories | **apartment_building/outdoor** (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
| --- | --- |

# Evaluation Metrics



$$\text{Accuracy} = \frac{\#\ correct\ predictions}{\#\ total\ predictions}$$
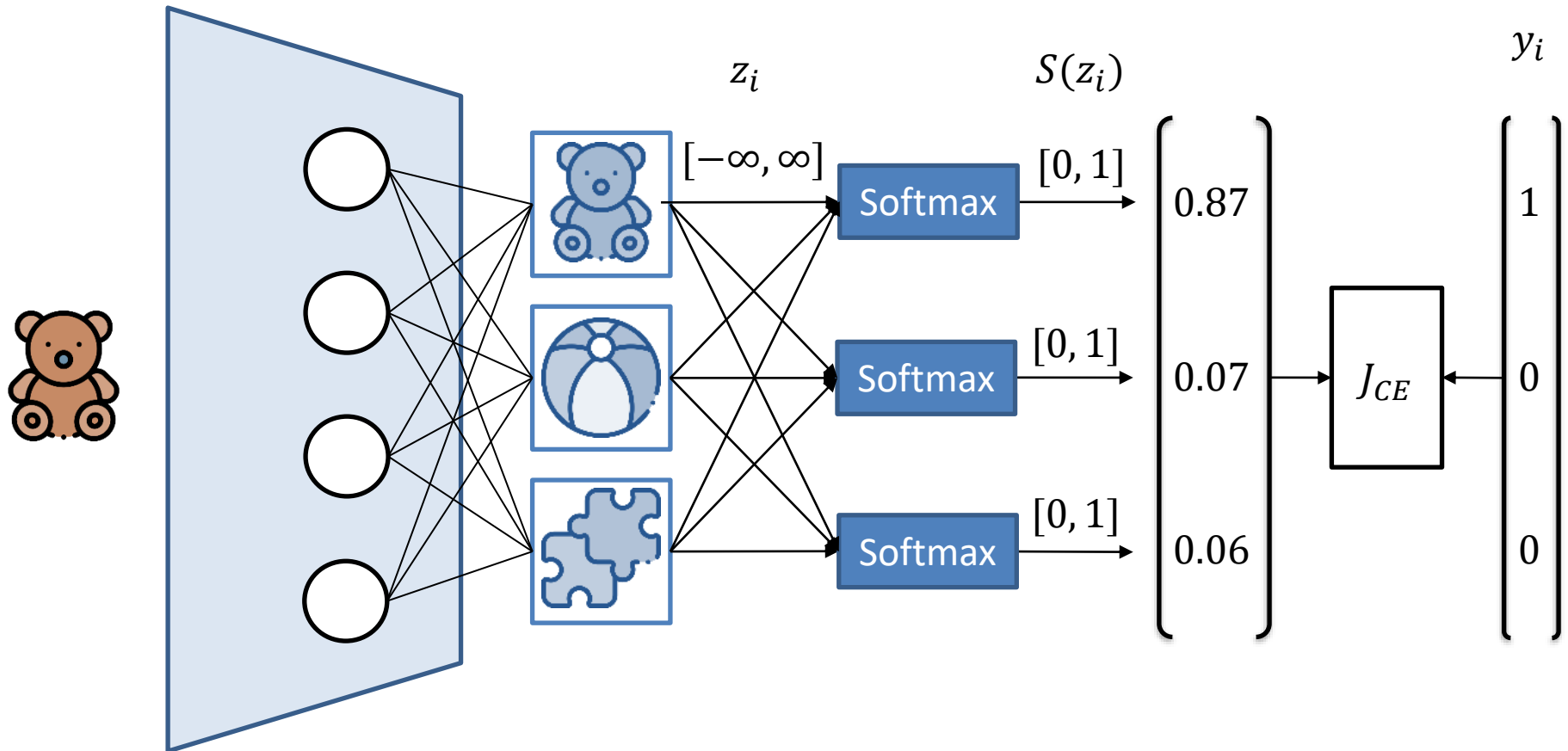
$$\text{Misclassification Error} = 1 - Accuracy$$

| Scene categories | apartment_building/outdoor (0.195), office_building (0.140), hospital (0.131), parking_garage/outdoor (0.126) |
|---|---|

$$\text{Top}-\text{k Accuracy} = \frac{\#\ correct\ prediction\ in\ the\ top\ k}{\#\ total\ predictions}$$

$$\text{Top}-\text{k Classification Error} = 1 - \text{Top}-\text{k Accuracy}$$
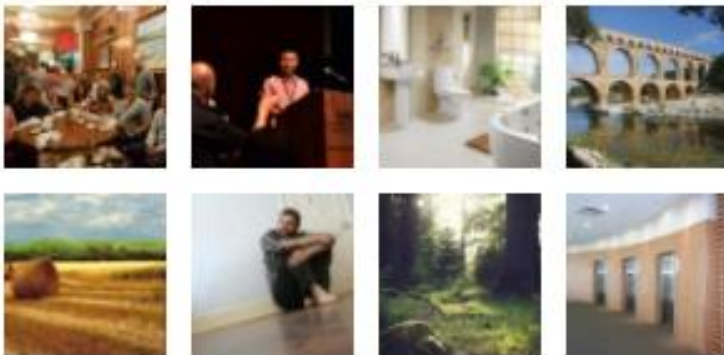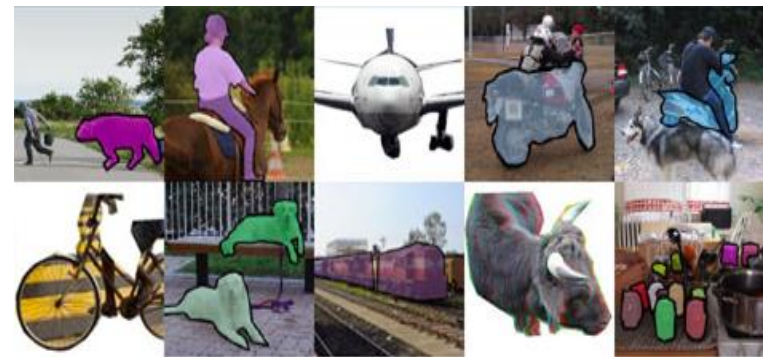
# Multi-class classification

# Datasets



- Scene photographs from image search engines, non-uniform distribution of images per category.
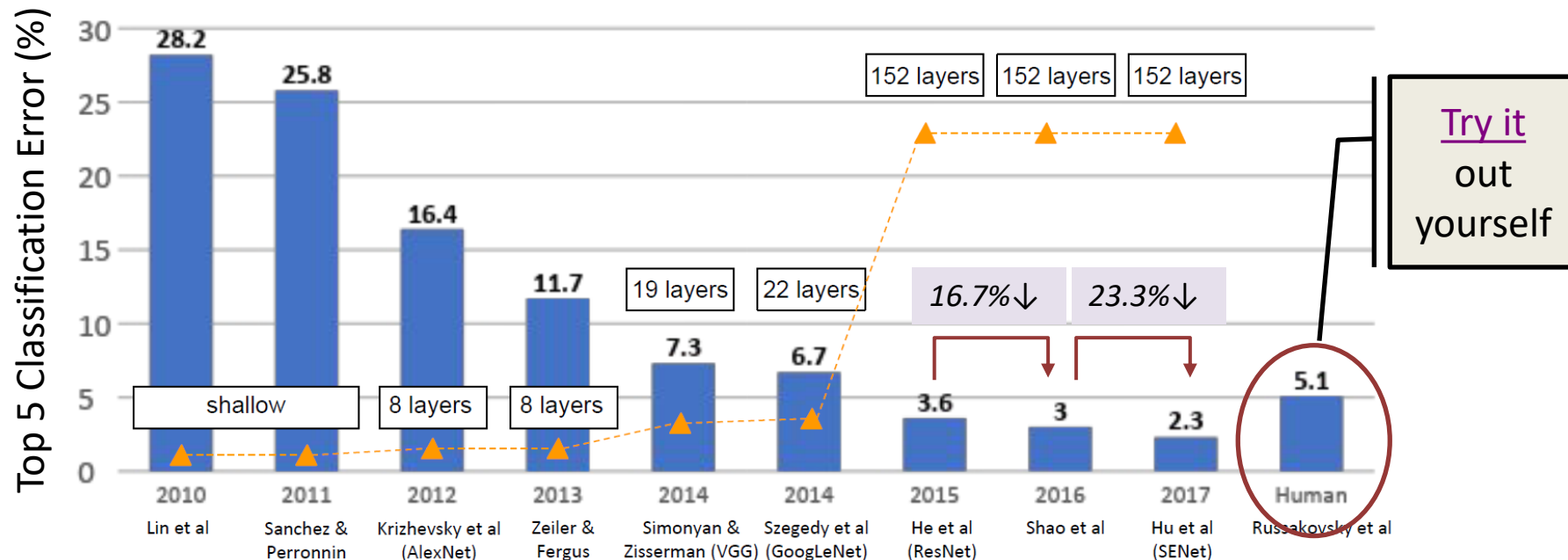
- 8M train, 36K val, 328K test, 1000 classes.



- Photos of 91 object types in the context of the broader question of scene understanding that would be easily recognizable by a 4 year old.

- 2.5 million labeled instances in 328k images.



from ILSVRC 2016, Places Dataset



from the COCO Dataset

# ImageNet Large Scale Visual Recognition Challenge

- Standard benchmark for image classification
- Standardised evaluation (top-k classification error)
- Evolution over time

ILSVRC: *https://image-net.org/challenges/LSVRC/2016/#scene*
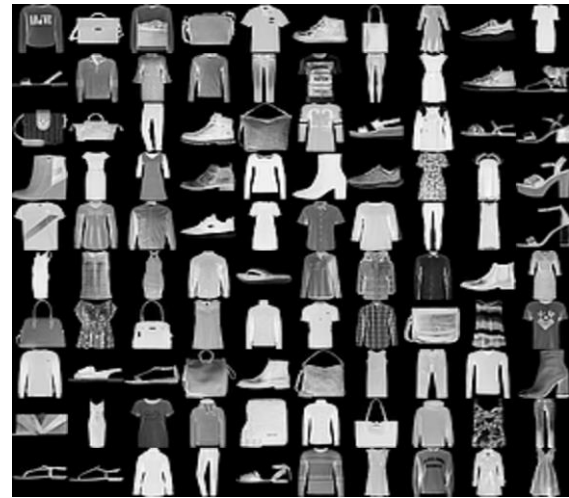
# Prototype Datasets

**MNIST**



**Fashion-Mnist**



- 60000 images of 10 handwritten digits of 28x28 pixels

- Accuracy: 99.9 %

- Usage: test if your model runs correctly

- 50000 images of 10 different clothes from Zalando of 28x28 pixels

- Accuracy: 96 %

- Usage: test if your model runs correctly
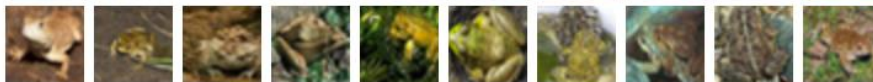
# Prototype Datasets

## CIFAR



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

- ✦ 60000 images of 32x32 color pixels

- ✦ Organized in 10 or 100 classes (CIFAR10 and CIFAR100)

# BASIC CNN MODELS

# Typical CNN architecture

# Different CNN Diagrams



Pooling (max, avg, sum) and strided convolution

Fully Connected

Flattening

Convolution and non-linearity

Image
Conv-64
Conv-64
maxpool
Conv-128
Conv-128
maxpool
Conv-256
Conv-256
maxpool
Conv-512
Conv-512
maxpool
FC-4096
FC-4096
FC-1000
Softmax

# CNN architecture

Input

Output class probabilities

$$S(z)_i = \frac{e^{z_i}}{\sum_j^C e^{z_j}} \qquad J_{CCE} = -\sum_i^C y_c \log(S(z)_i)$$

$z_i$

$S(z_i)$

Softmax

$J_{CE}(S(z_i), y_i)$

Cross Entropy Loss

# LeNet (1998)

Input

Output
class
probabilities

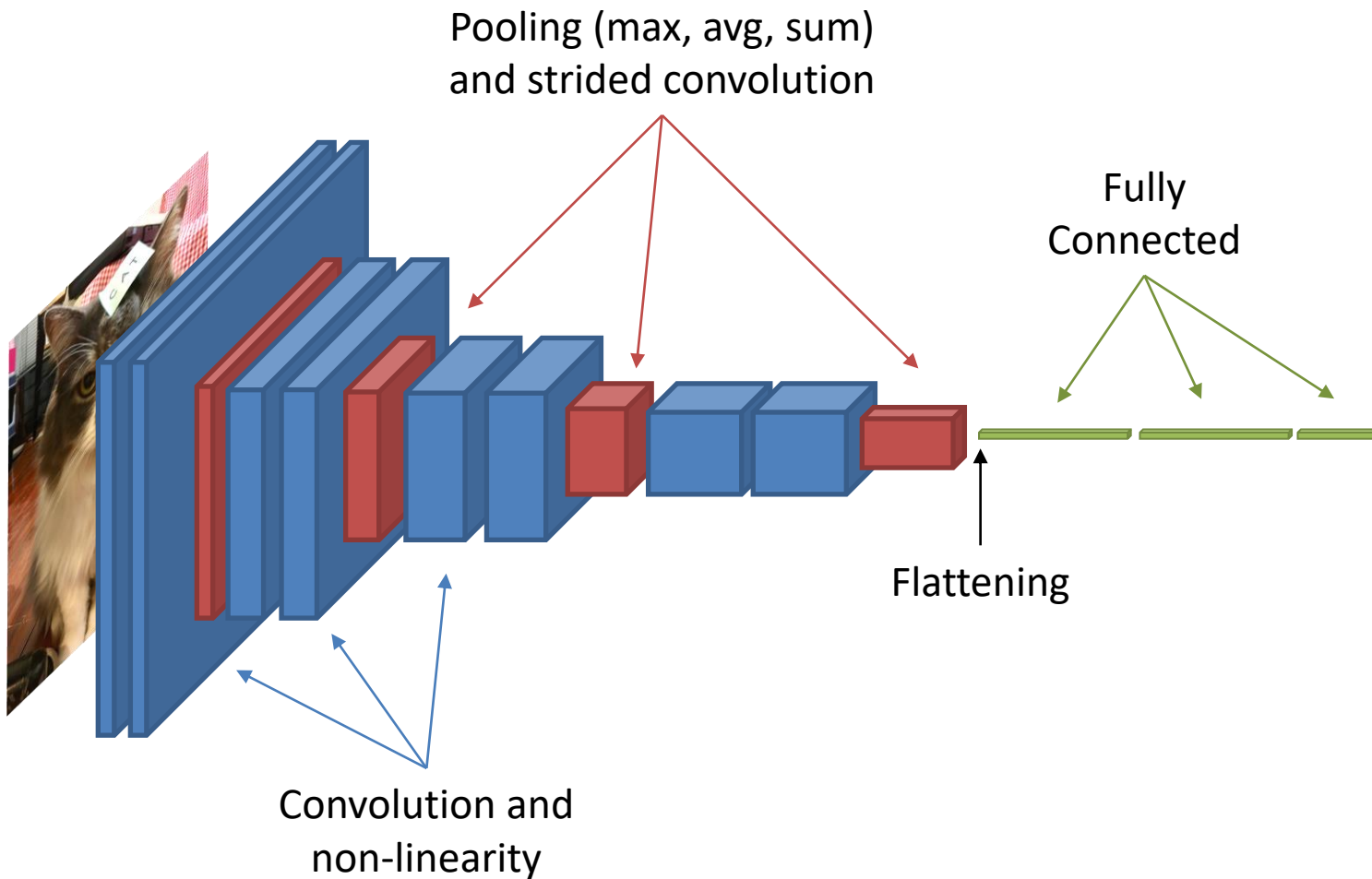| Conv. 1 | Sub. 1 | Conv. 2 | Sub. 2 | Conv. 2 | FC 1 | Output |
|---------|--------|---------|--------|---------|------|--------|
| (5x5, 6) | (2x2) | (5x5, 16) | (2x2) | (5x5, 120) | (84) | (10) |

Nonlinearity used: $tanh$

*LeCun, Y.; Boser, B.; Denker, J. S.; Henderson, D.; Howard, R. E.; Hubbard, W. & Jackel, L. D. (1989).*
*Backpropagation applied to handwritten zip code recognition. Neural Computation*

# Image Classification Evolution

training of deep nets
unsuccessful (except
for CNNs)

AI winter

Large Datasets &
Computational Power

DL breakthrough

LeNet

1998      2012

AlexNet, VGG, GoogLeNet

# 2012 - 2014

# AlexNet (2012)

8 convolutional (11x11, 5x5, 3x3, 3x3, 3x3) and 3 fully-connected layers.
Took ~6 days on two GTX 580 3GB GPUs (own implementation)



Key concepts introduced to mainstream CNN methodology:
- **ReLU non-linearity** [see here for details]
- **Dropout** [see here for details]
- **Data Augmentation**

*Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25, 1097-1105.*

# VGG (2014)

"Our main contribution is a thorough evaluation of **networks of increasing depth** using an architecture with **very small (3 × 3) convolution** filters […] significant improvement can be achieved by **pushing the depth to 16–19** weight layers."
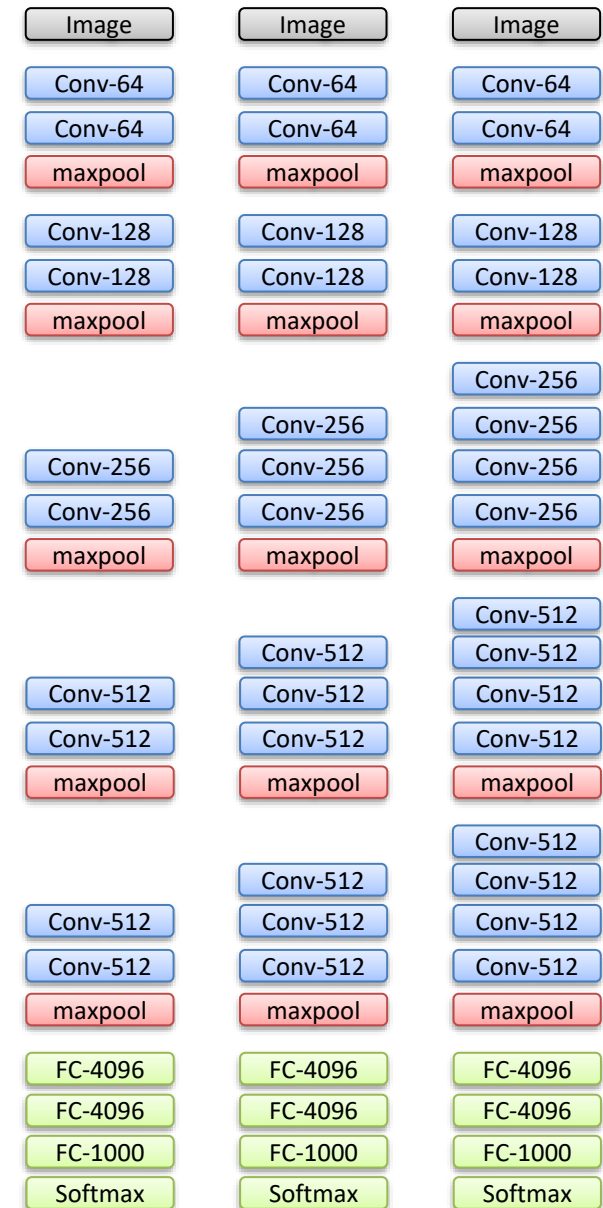
- "Very" deep network (up to 19 layers)
- **3x3 filters**, stride 1, padding 1
- **Stacked** convolutions
- 2x2 non-overlapping max-pooling
- # features increases as we go deeper
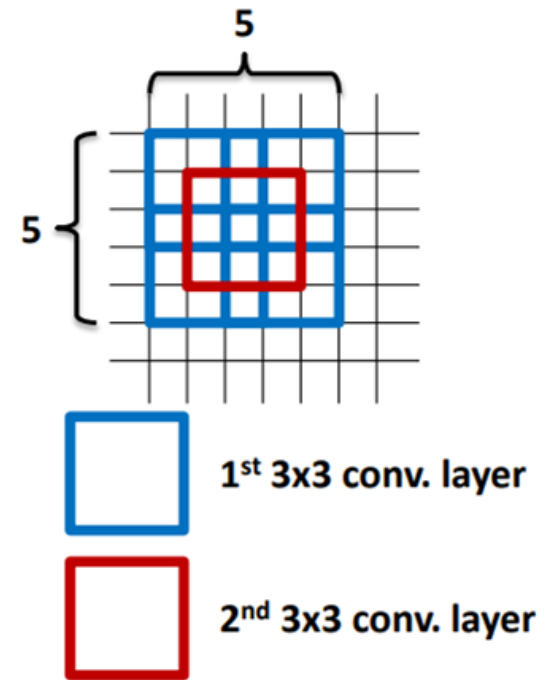- ReLU non-linearity
- Data Augmentation
- Dropout

*Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556*

| Image | Image | Image |
|---|---|---|
| Conv-64 | Conv-64 | Conv-64 |
| Conv-64 | Conv-64 | Conv-64 |
| maxpool | maxpool | maxpool |
| Conv-128 | Conv-128 | Conv-128 |
| Conv-128 | Conv-128 | Conv-128 |
| maxpool | maxpool | maxpool |
| | | Conv-256 |
| | Conv-256 | Conv-256 |
| Conv-256 | Conv-256 | Conv-256 |
| Conv-256 | Conv-256 | Conv-256 |
| maxpool | maxpool | maxpool |
| | | Conv-512 |
| | Conv-512 | Conv-512 |
| Conv-512 | Conv-512 | Conv-512 |
| Conv-512 | Conv-512 | Conv-512 |
| maxpool | maxpool | maxpool |
| | | Conv-512 |
| | Conv-512 | Conv-512 |
| Conv-512 | Conv-512 | Conv-512 |
| Conv-512 | Conv-512 | Conv-512 |
| maxpool | maxpool | maxpool |
| FC-4096 | FC-4096 | FC-4096 |
| FC-4096 | FC-4096 | FC-4096 |
| FC-1000 | FC-1000 | FC-1000 |
| Softmax | Softmax | Softmax |

24

# VGG (2014)

Why 3x3 layers?

What is the receptive field if we stack two $3 \times 3$ convolutions?



1st 3x3 conv. layer

2nd 3x3 conv. layer

# VGG (2014)

Why 3x3 layers?

How many parameters are in a layer if we use (a) $3 \times 3$ filters? How many if we use (b) $5 \times 5$ filters?

| Conv-256 |
|:---:|
| Conv-256 |
| Conv-256 |
| Conv-256 |
| maxpool |

(a) $3 \times 3 \times C \times C + C \cong 590\text{k}$ (for $C = 256$)

(b) $5 \times 5 \times C \times C + C \cong 1.6\text{M}$ (for $C = 256$)
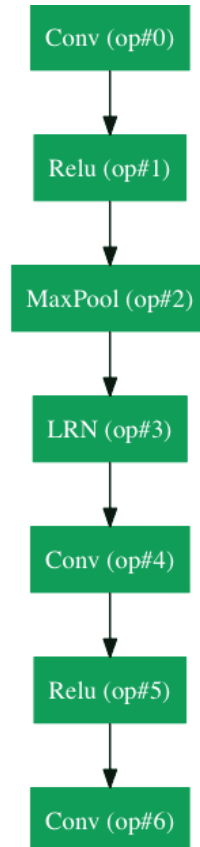
# VGG (2014)

Why 3x3 layers?

- Stacked conv. layers have a large receptive field
    - two 3x3 layers: *5x5 receptive field*
    - three 3x3 layers: *7x7 receptive field*
- More non-linearity
- Less parameters to learn (compared to using large kernels)
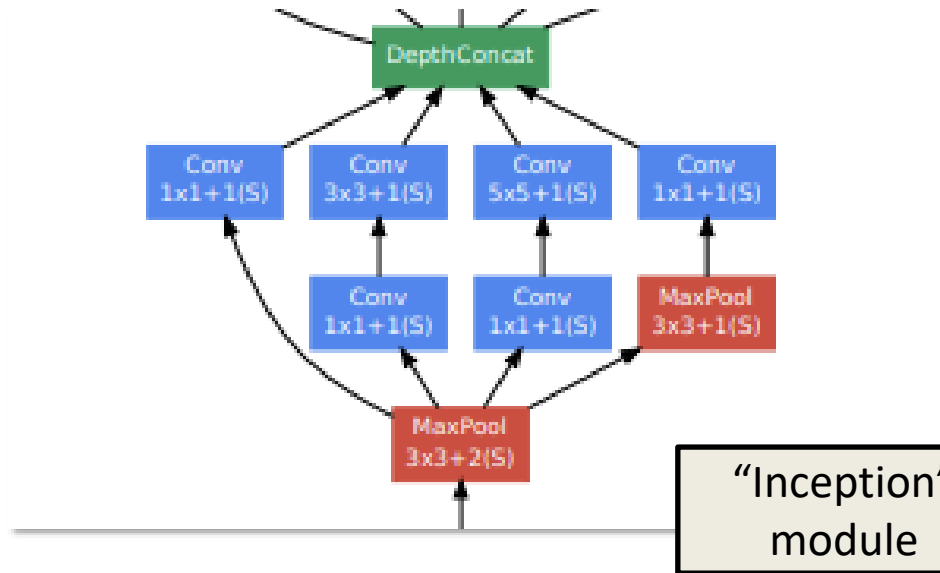    - ~140M per net

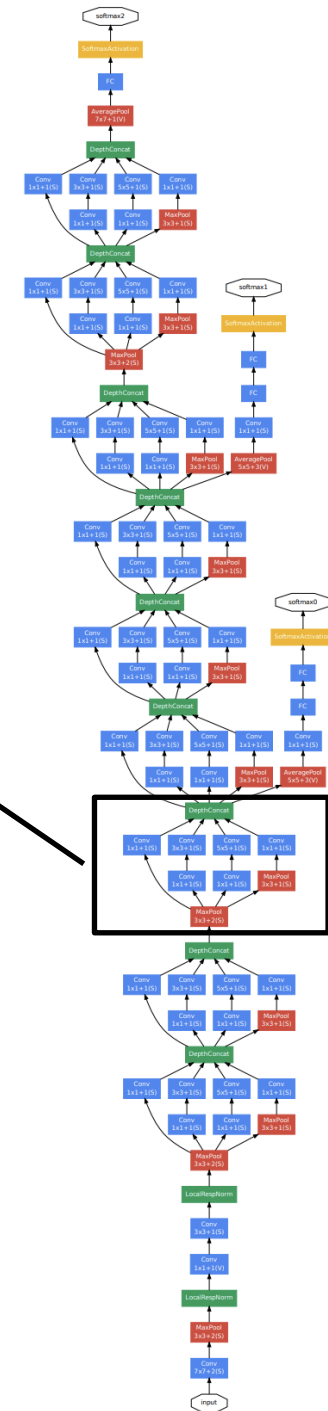# GoogLeNet (2015)
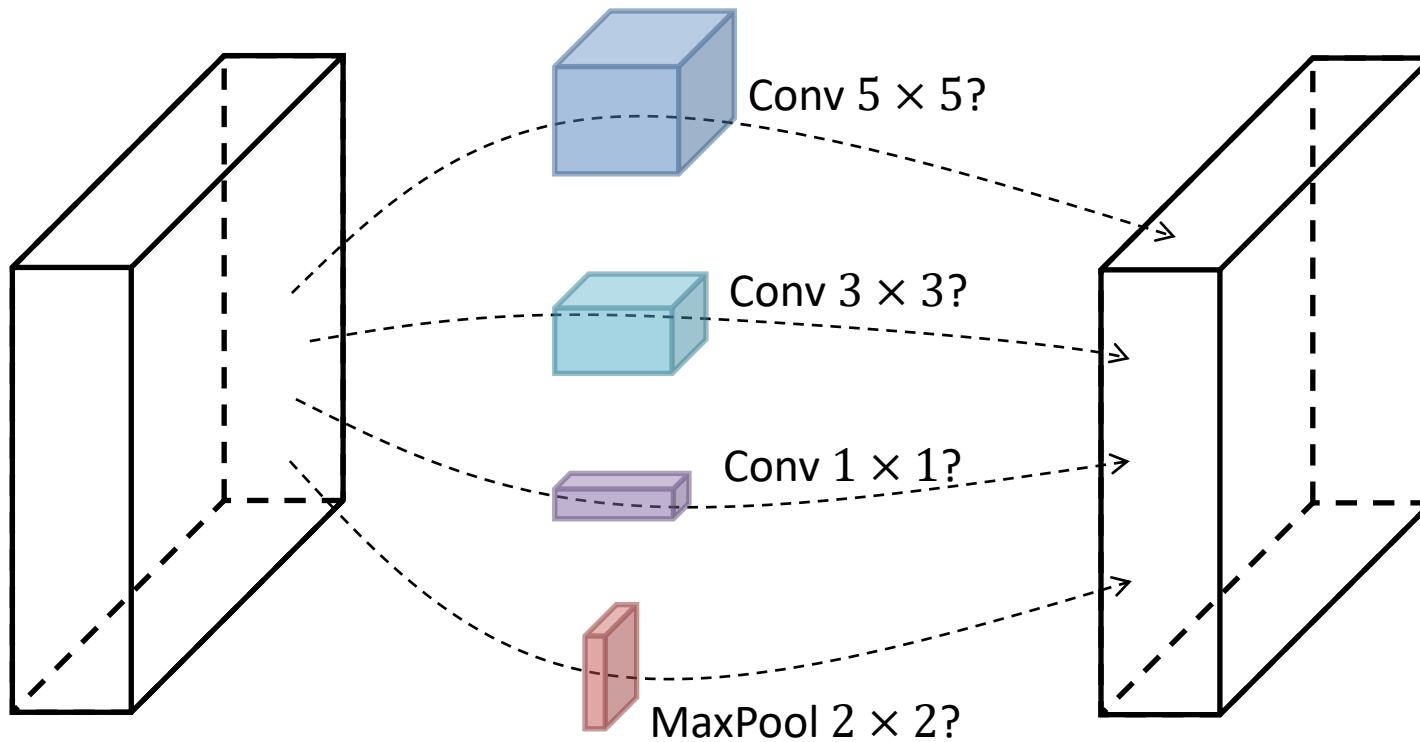
# GoogLeNet (2015)

# GoogLeNet (2015)



"Inception" module

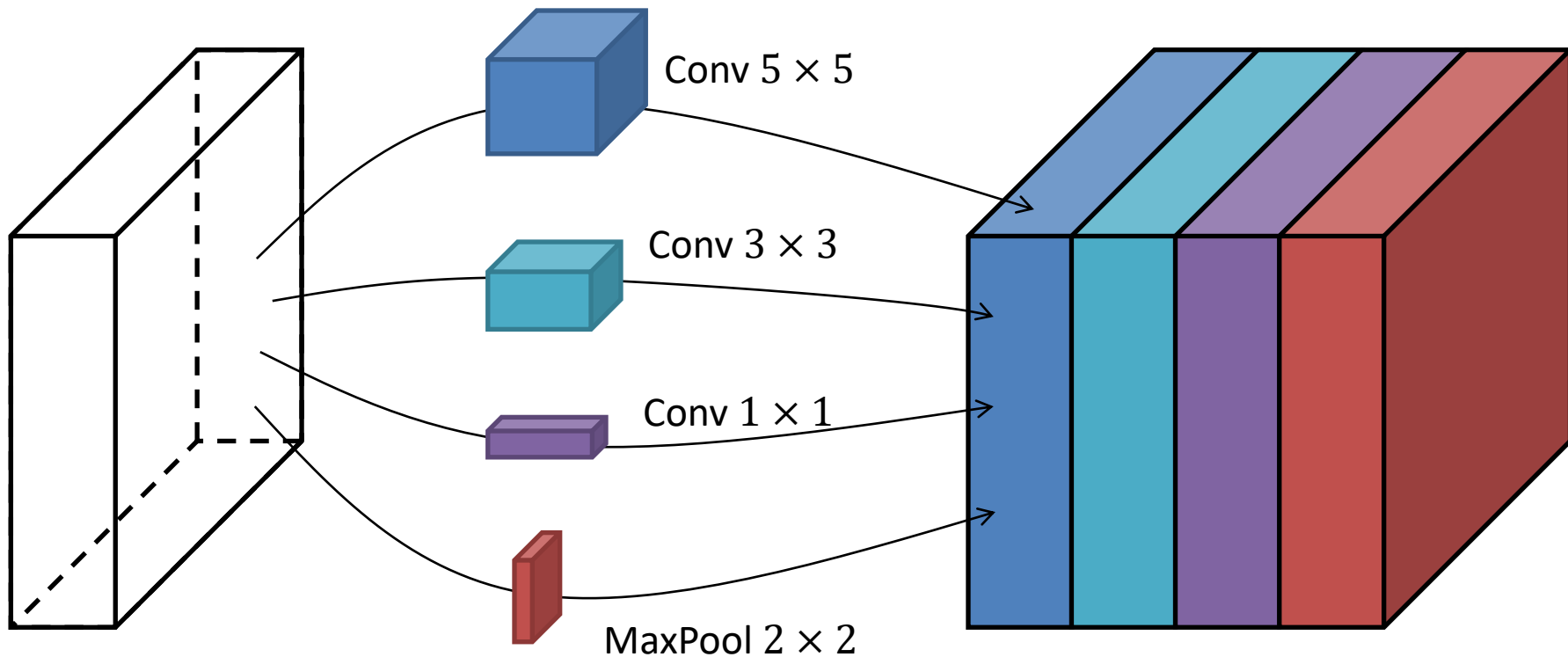■ Convolution / FC

■ Max / Avg Pooling

■ Softmax

■ Concatenation

*Szegedy et al, "Going deeper with convolutions", 2014*

# GoogLeNet (2015) – Inception Module

Conv 5 × 5?

Conv 3 × 3?

Conv 1 × 1?

MaxPool 2 × 2?

Which operation should we choose?

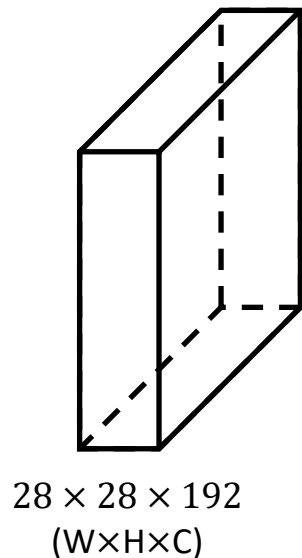# GoogLeNet (2015) – Inception Module



Conv 5 × 5

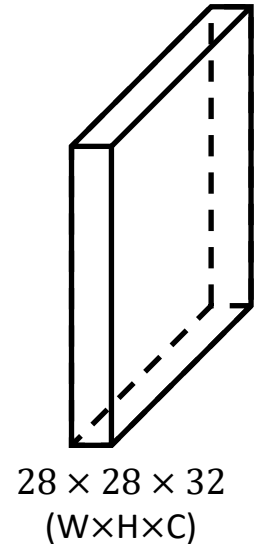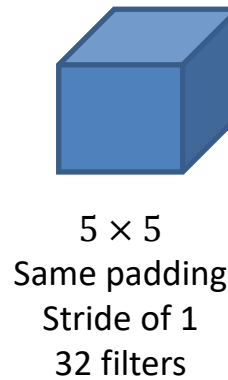Conv 3 × 3

Conv 1 × 1

MaxPool 2 × 2

Let's do them all at once! Multi-scale

But the size explodes…

# Reducing the Computational Cost

Number of multiplications:
$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = \mathbf{120,422,400}$$

$5 \times 5$
Same padding
Stride of 1
32 filters

$28 \times 28 \times 192$
(W×H×C)

$28 \times 28 \times 32$
(W×H×C)

$$28 \times 28 \times 16 \times 1 \times 1 \times 192 = \mathbf{2,408,448}$$

$$28 \times 28 \times 32 \times 5 \times 5 \times 16 = \mathbf{10,035,200}$$

$1 \times 1$
Same padding
Stride of 1
16 filters

$5 \times 5$
Same padding
Stride of 1
32 filters

$28 \times 28 \times 192$
(W×H×C)

$28 \times 28 \times 16$
(W×H×C)

$28 \times 28 \times 32$
(W×H×C)

# GoogLeNet (2015) – Inception Module



Reduce the number of channels, using $1 \times 1$ convolutions

# GoogLeNet (2015)

"…By adding **auxiliary classifiers** connected to these intermediate layers, we would expect to **encourage discrimination** in the lower stages in the classifier, **increase the gradient signal** that gets propagated back, and provide **additional regularization**…"

Auxiliary supervision

- 🟦 Convolution / FC
- 🟥 Max / Avg Pooling
- 🟨 Softmax
- 🟩 Concatenation

*Szegedy et al, "Going deeper with convolutions", 2014*

37

# GoogLeNet (2015)

22-layer deep CNN, but reduced the number of parameters **from 60 million (AlexNet) to 4 million (12x)**

Low use of fully connected layers, they use average pooling instead

■ Convolution / FC

■ Max / Avg Pooling

■ Softmax

■ Concatenation



*Szegedy et al, "Going deeper with convolutions", 2014*

# ImageNet Classification

# Wrap Up

**time consuming, high memory demanding**

**training of deep nets unsuccessful (except for CNNs)**

*Large Datasets & Computational Power*

*Dropout & ReLU*

*LeNet*

*AlexNet*

*VGG*

*GoogLeNet*

1998    2012    2015

DSNs, FitNets, Highway Nets and ResNets

# 2015 - 2016

# Deeply Supervised Networks

Adding intermediate supervision (discrimitative loss)



*Lee, Chen-Yu, et al. "Deeply-supervised nets." Artificial intelligence and statistics. PMLR, 2015*

# FitNets



Teacher
# parameters = 32

Student
# parameters = 24

Knowledge transfer: trade width for depth to reduce the number of parameters

Key intuition: bigger is not better, just easier to optimise! Over-parameterisation helps optimisation.

*A. Romero et al. "FitNets: Hints for Thin Deep Nets" ICLR 2015*

# Highway Networks



$x$: Input

$H(x)$: Layer transform
(a block of operations)

$T(x)$ : Transform gate
(how much we transform the input)

$1 - T(x)$: Carry gate
(how much we preserve the input)

$$output = H(x)T(x) + x(1 - T(x))$$

Key idea: provide an alternative, direct path for the gradient to flow back

*Srivastava, Rupesh Kumar, Klaus Greff, and Jürgen Schmidhuber. "Highway networks." ICML 2017*

# Residual connections



Traditional feedforward

Feedforward with
**residual** connection

One way to avoid vanishing gradient is to provide alternative routes for your gradients to flow through, a.k.a. residual connections

# Residual Networks

Residual skip-connection



$x$

$F(x)$

$\oplus$

$F(x) + x$

$$F(x) = \begin{pmatrix} \text{BN} \\ + \\ \text{ReLU} \\ + \\ \text{Conv.} \end{pmatrix} * n$$

- Compute a small change to the input to get a slightly altered representation
- Gradient flows easily to the bottom layers of the network
- **Ultra deep** architecture (152 layers)
- Adopts Batch Normalisation

*S. Ioffe, C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. ICML 2015*
*Kaiming, et al. "Deep residual learning for image recognition." CVPR. 2016*

# Residual Networks

# Residual Blocks



Basic Block

Bottleneck block

# Residual Networks



Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

# Residual Networks



(a) without skip connections

(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

*Li, Hao, et al. "Visualizing the loss landscape of neural nets." arXiv preprint arXiv:1712.09913 (2017)*

Stochastic depth, Dense Nets

# 2015 - 2016

# Stochastic Depth

*Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016*

# Stochastic Depth

*Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016*

# Stochastic Depth

*Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016*

# Stochastic Depth

*Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016*

# Stochastic Depth



training data

**short** during **training**

test data

**deep** during **testing**

Key aspects:
- Implicit **ensemble** of $2^L$ models
- Improved **gradient flow**
- 25% speedup during training
- Lower error

*Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016*

# DenseNet

Connect every layer to every other layer of the same filter size (dense blocks).



1x1 Conv. + pool

BN
+
ReLU
+
Conv.
+
Dropout

*Huang, Gao, et al. "Densely connected convolutional networks." CVPR 2017*

# DenseNet

Example Dense Block (4 layers, growth rate 4):



DenseNet can have very narrow layers (**small growth rate**). One explanation for this is that each layer has access to all the preceding feature-maps in its block. Each layer then adds a little bit of information ($k$ feature-maps) to the global state of the network



*Huang, Gao, et al. "Densely connected convolutional networks." CVPR 2017*

# ResNeXt

ResNeXt introduces the concept of "cardinality". Apart from depth (how many layers) and width (how many filters).



Key concept: create a lot of paths with a very small number of channels using $1 \times 1$ convolutions

*Xie, et al. "Aggregated Residual Transformations for Deep Neural Networks." CVPR 2017*

# Squeeze-&-Excitation Networks



Squeeze        Excitation        Output:

Channel-wise

**Squeeze**: average pooling - captures channel-wise global info

**Excitation**: non-linear channel interaction + gating (sigmoid)

**Output**: scales activations, leveraging global information at all levels

*Hu, et al. "Squeeze-and-Excitation Network." CVPR 2018*

Top-1 accuracy [%] vs Operations [G-FLOPs]

NASNet-A-Large

SE-ResNeXt-101(32x4d)

SE-ResNeXt-50(32x4d)

SENet-154

Inception-ResNet-v2
Inception-v4

Xception

DualPathNet-131

SE-ResNet-101

DualPathNet-98

SE-ResNet-152

ResNeXt-101(64x4d)

SE-ResNet-50

ResNeXt-101(32x4d)

ResNet-152

DenseNet-201

Inception-v3

DenseNet-161

ResNet-101

FB-ResNet-152

DualPathNet-68

ResNet-50

Caffe-ResNet-101

VGG-19_BN

DenseNet-169

DenseNet-121

VGG-16_BN

NASNet-A-Mobile

ResNet-34

BN-Inception

VGG-13_BN

MobileNet-v2

VGG-11_BN

VGG-19

ResNet-18

VGG-16

MobileNet-v1

VGG-13

ShuffleNet

VGG-11

GoogLeNet

1M  5M  10M  50M  75M  100M  150M

SqueezeNet-v1.1

SqueezeNet-v1.0

AlexNet

66

# ImageNet Classification

# Wrap Up



training of deep nets unsuccessful (except for CNNs)

time consuming, high memory demanding

addressing optimization issues

way of training
architecture design

Dropout & ReLU

Batch Norm

Large Datasets & Computational Power

LeNet

AlexNet

VGG

GoogLeNet

Deeply Supervised Networks

FitNets

Highway Networks

Residual Networks

Stochastic Depth

DenseNets & ResNeXt

SE Networks

1998    2012    2015    2017    2018

Neural Architecture Search

# 2018 - 2019

# Neural Architecture Search (NAS)



*Zoph, B., & Le, Q. V. (2016). Neural architecture search with reinforcement learning. ICLR 2017*

# Neural Architecture Search (NAS)



Zoph, B., & Le, Q. V. (2016). Neural architecture search with reinforcement learning. ICLR 2017

# CIFAR-10 Results



**800** GPUs!
**28** Days!

Misclassification error (%)

CIFAR-10

- DSN
- FitNets
- Highway
- ResNets
- Stoch. Depth
- NAS
- ResNeXt
- DenseNets

8.22  8.39  7.60  6.41  5.25  3.65  3.58  3.46

# DARTS: Differentiable Architecture Search

Key concept: Start with all operations connected. Then slowly prune connections that are not being used. Completely remove them at the end for test.



**1**        **2**        **3**        **4**

Relaxed ────────────────────────────────▶ Constrained

*Chu, X.,et al. "Fair darts: Eliminating unfair advantages in differentiable architecture search". ECCV 2020*

# CIFAR-10 Results

# Wrap Up



training of deep nets unsuccessful (except for CNNs)

time consuming, high memory demanding

addressing optimization issues

understanding new architectures

Automatic Architecture Search

way of training
architecture design

Dropout & ReLU

Large Datasets & Computational Power
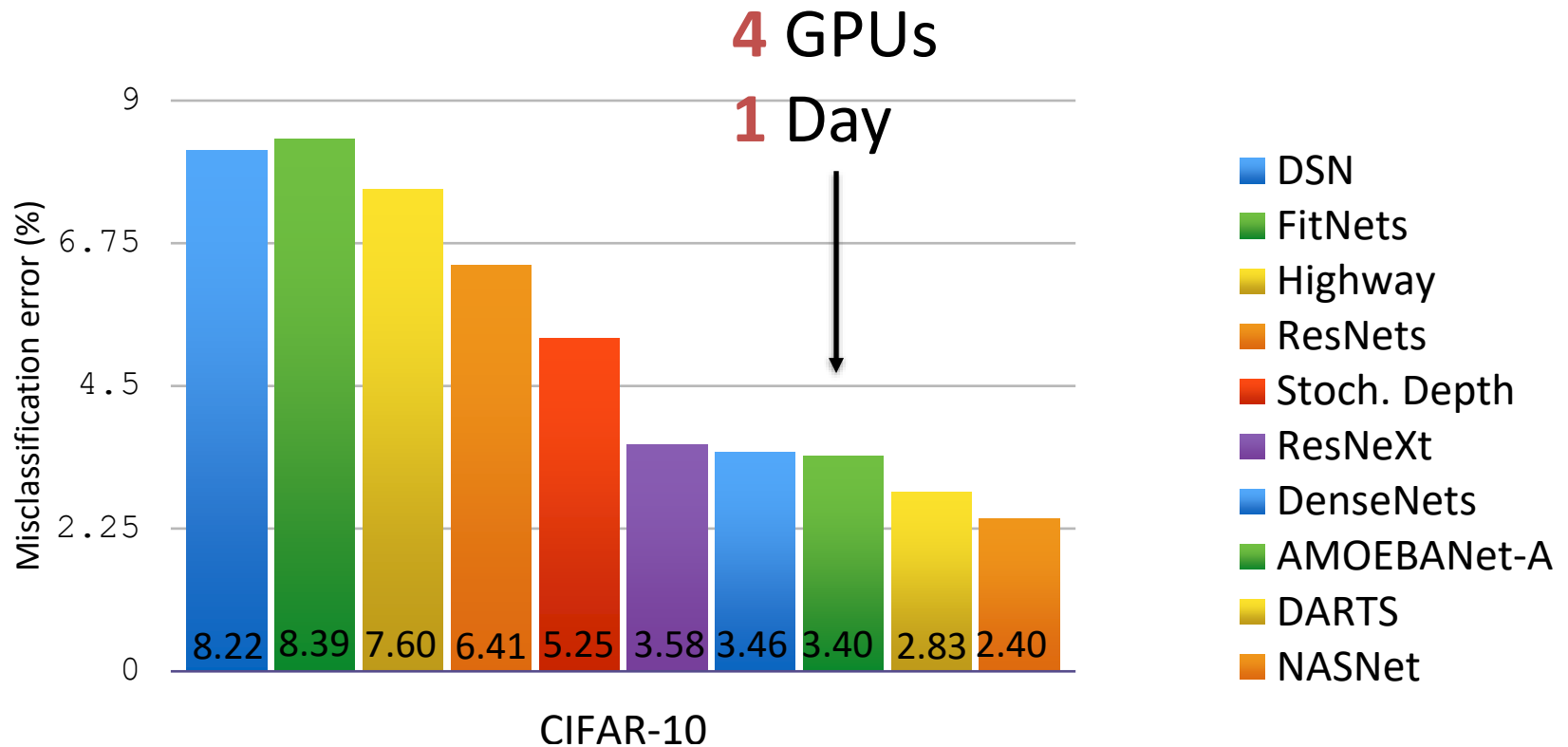
Batch Norm

???

Block-wise Architecture Search

LeNet
AlexNet
VGG
GoogLeNet
Deeply Supervised Networks
FitNets
Highway Networks
Residual Networks
Stochastic Depth
DenseNets & ResNeXt
Unraveled view
Unrolled iterative estimation view
SE Networks
Neural Architecture Search
Neural Block Architecture Search
Neural Evolution
DARTS
EfficientNet

1998        2012        2015        2017    2018        2019

# Wrap Up

*Slide credit: Pau Rodriguez and Adriana Romero. MCV M3 Architectures for image classification.*
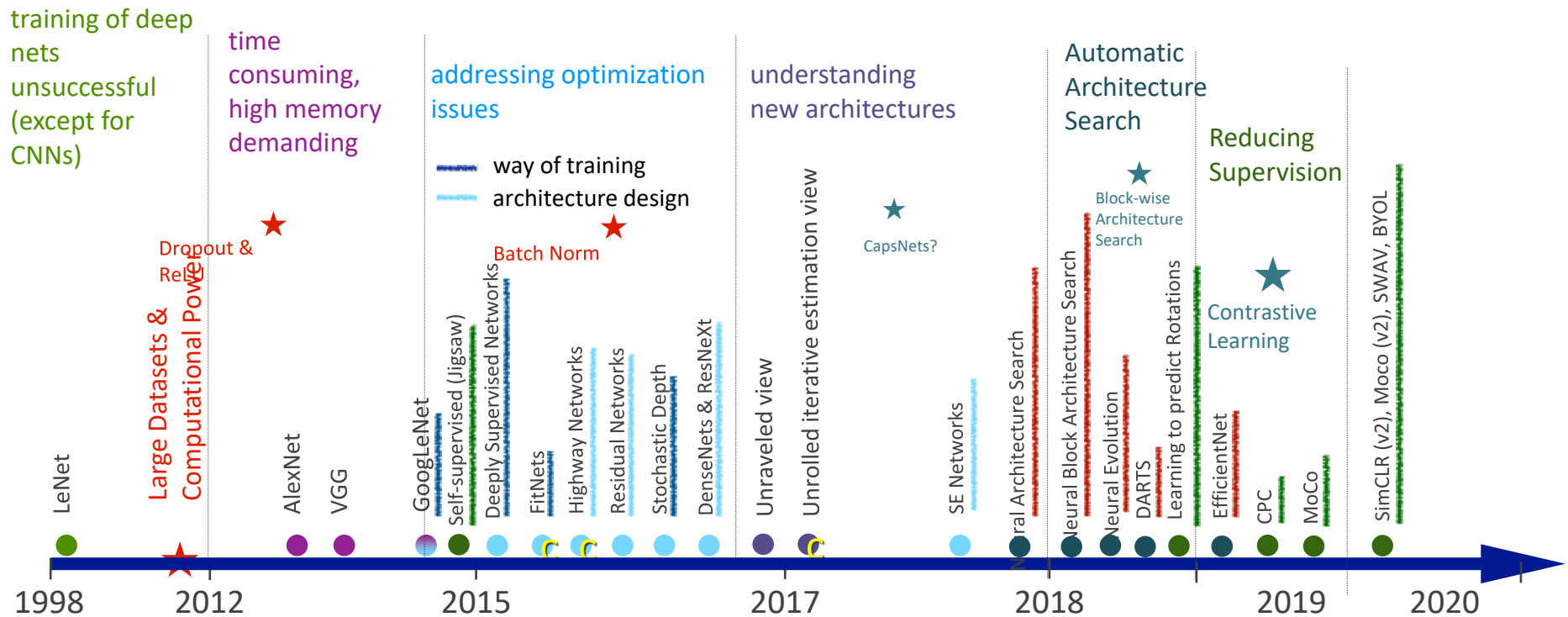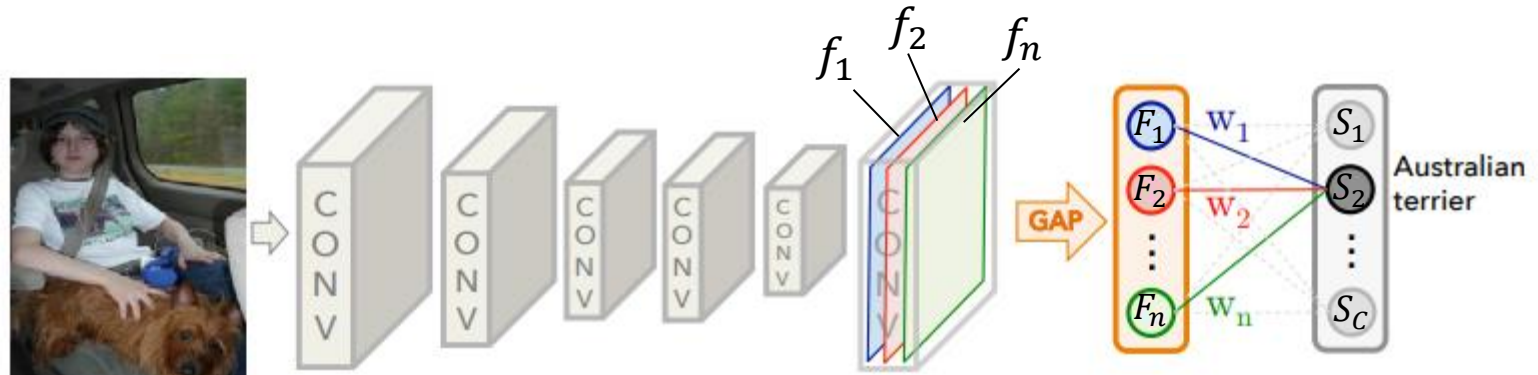
76

# GETTING INSIGHTS

# Class Activation Mapping



$$F_k = \sum_{x,y} f_k(x,y)$$

$$S_c = \sum_{k} w_k^c \sum_{x,y} f_k(x,y)$$

*B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. CVPR'16*

# Class Activation Mapping



$$M_c(x, y) = \sum_k w_k^c f_k(x, y)$$

B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning Deep Features for Discriminative Localization. CVPR'16

# Class Activation Mapping

Seen in previous classes

# RECAP OF IMPORTANT IDEAS

# Dropout



Simulate bagging
by dropping each
unit in the network
with probability 1/2

# Dropout in PyTorch

`torch.nn.Dropout(p=0.5, inplace=False)`

https://pytorch.org/docs/stable/generated/torch.nn.Dropout.html



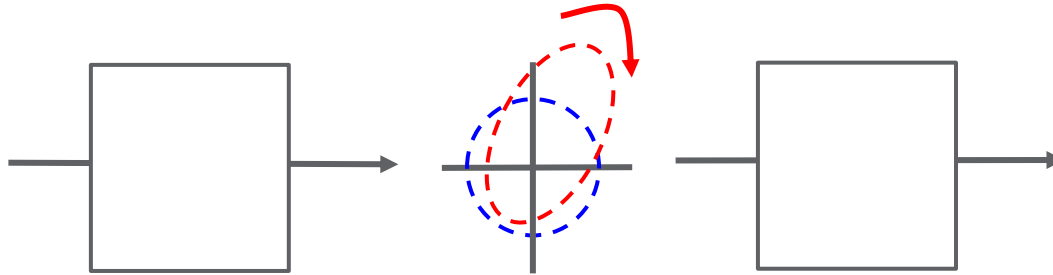https://losslandscape.com/          https://youtu.be/2PqTW_p1fls
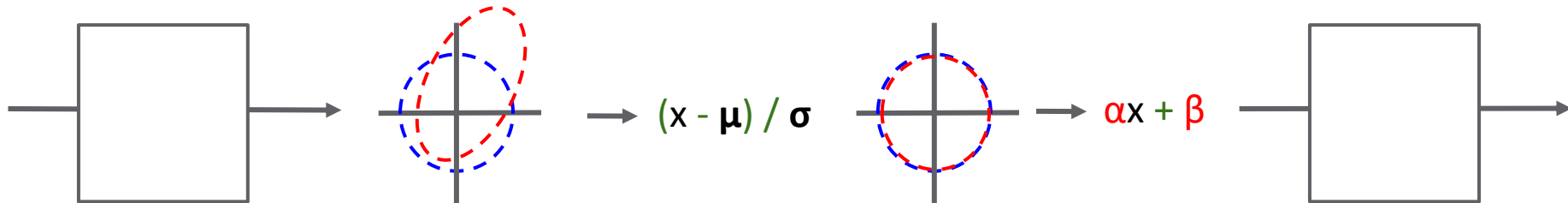
# Batch Normalization

Before:

After:

# Batch Normalization in PyTorch

```
torch.nn.BatchNorm1d(num_features, eps=1e-05, momentum=0.1,
                          affine=True, track_running_stats=True)

torch.nn.BatchNorm2d(num_features, eps=1e-05, momentum=0.1,
                          affine=True, track_running_stats=True)
```

https://pytorch.org/docs/stable/nn.html#normalization-layers

# Data Augmentation

Best way to make a machine learning model generalize better is to train it on more data -> **Create new data!**



Image source: Building powerful image classification models using very little data

# Data Augmentation in PyTorch

```python
From torchvision import transforms

rgb_mean = (0.4914, 0.4822, 0.4465)
rgb_std = (0.2023, 0.1994, 0.2010)

transform_train = transforms.Compose([
    transforms.RandomCrop(32, padding=4),
    transforms.RandomHorizontalFlip(),
    transforms.ToTensor(),
    transforms.Normalize(rgb_mean, rgb_std),
])
```

https://pytorch.org/vision/stable/transforms.html

# Transfer Learning (supervised pre-training)

- Initialize the weights of your model using the optimal weights learnt on a similar task!

- Fine-tuning can provide a reasonably good model even when we have very few training data.
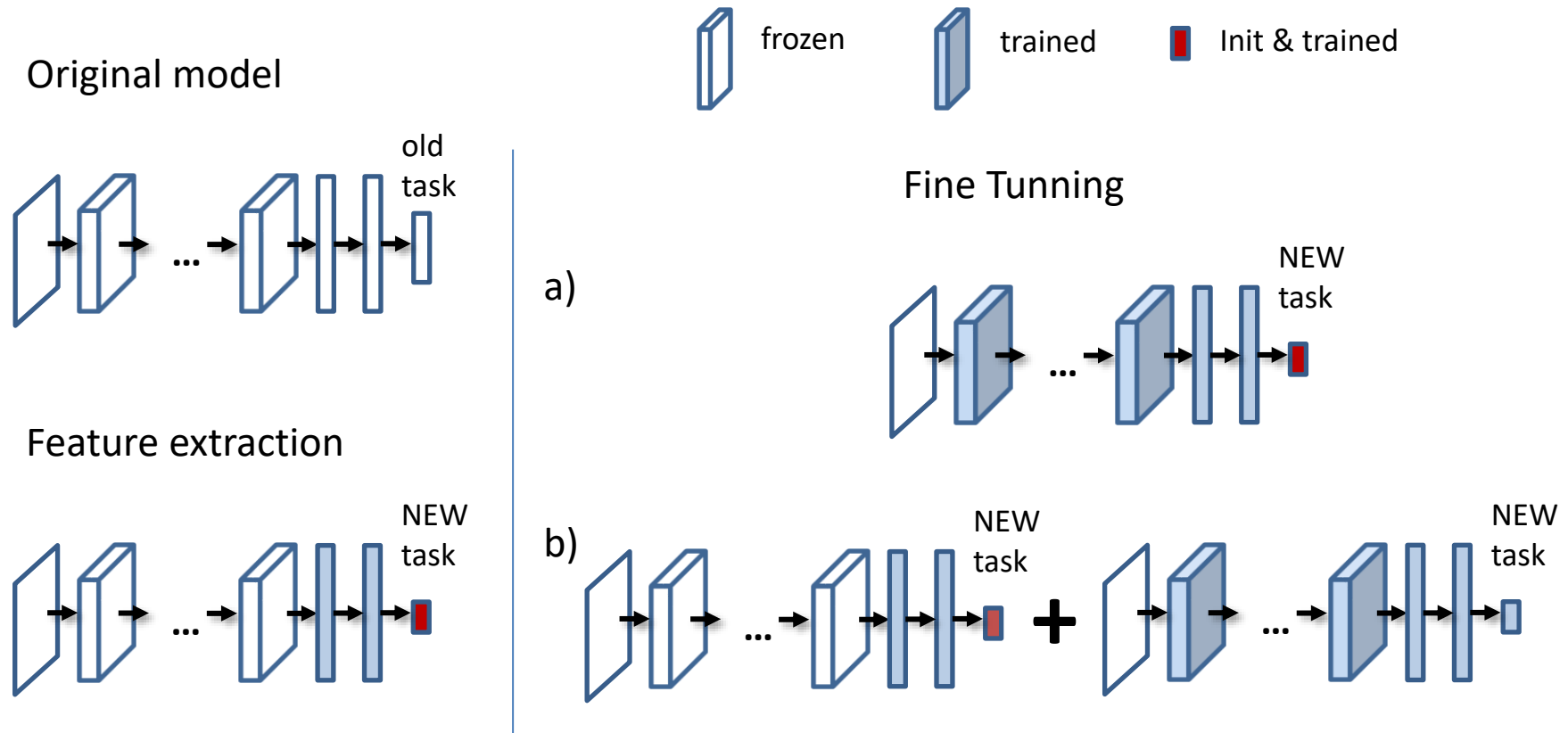
If you know how to recognize...                    You will be able to recognize...



Transfer Learning

Image source: Yannis Ghazouani,  2016.

# Transfer Learning (supervised pre-training)

Original model

frozen     trained     Init & trained



old task

Fine Tunning

a)

NEW task

Feature extraction

NEW task

b)

NEW task

+

NEW task

Sharif Razavian, Ali, et al. "CNN features off-the-shelf: an astounding baseline for recognition." Proceedings of the IEEE CVPRW. 2014.

# Transfer Learning

In practice, very few people train a CNN from scratch (with random initialization). Don't have sufficient data and/or resources.

**It is common to use pretrained CNNs** and use them either as an initialization or a fixed feature extractor for the new task.

**When and how to transfer learning?**

- New dataset is small and similar to original dataset: feature ext.
- New dataset is large and similar to the original dataset: feat ext+fine tunn
- New dataset is small but very different from the original dataset: feat ext
- New dataset is large and very different from the original dataset: pre trained

**Practical advice**

- Constraints from pretrained models.
- Learning rates.

https://cs231n.github.io/transfer-learning/