

# Visual scene real-time analysis for Intelligent Vehicles:

## Deep-Learning for visual scene analysis

Pr. Fabien MOUTARDE  
Center for Robotics,  
MINES ParisTech  
PSL Université Paris

[Fabien.Moutarde@mines-paristech.fr](mailto:Fabien.Moutarde@mines-paristech.fr)

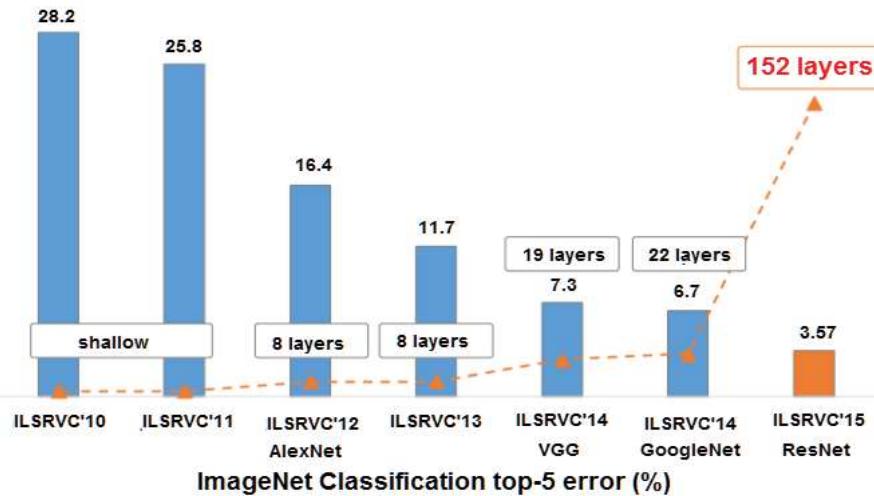
<http://people.mines-paristech.fr/fabien.moutarde>

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 1

## Outline

- **Recalls on Convolutional Neural Networks (CNN or ConvNets) and Deep-Learning**
- Transfer Learning
- Beyond Image Classification: DETECTION OF OBJECTS
- Instance segmentation with DeepLearning
- DL for Human pose inference and depth estimation
- Semantic segmentation with DeepLearning
- Interest and use of simulations / synthetic videos

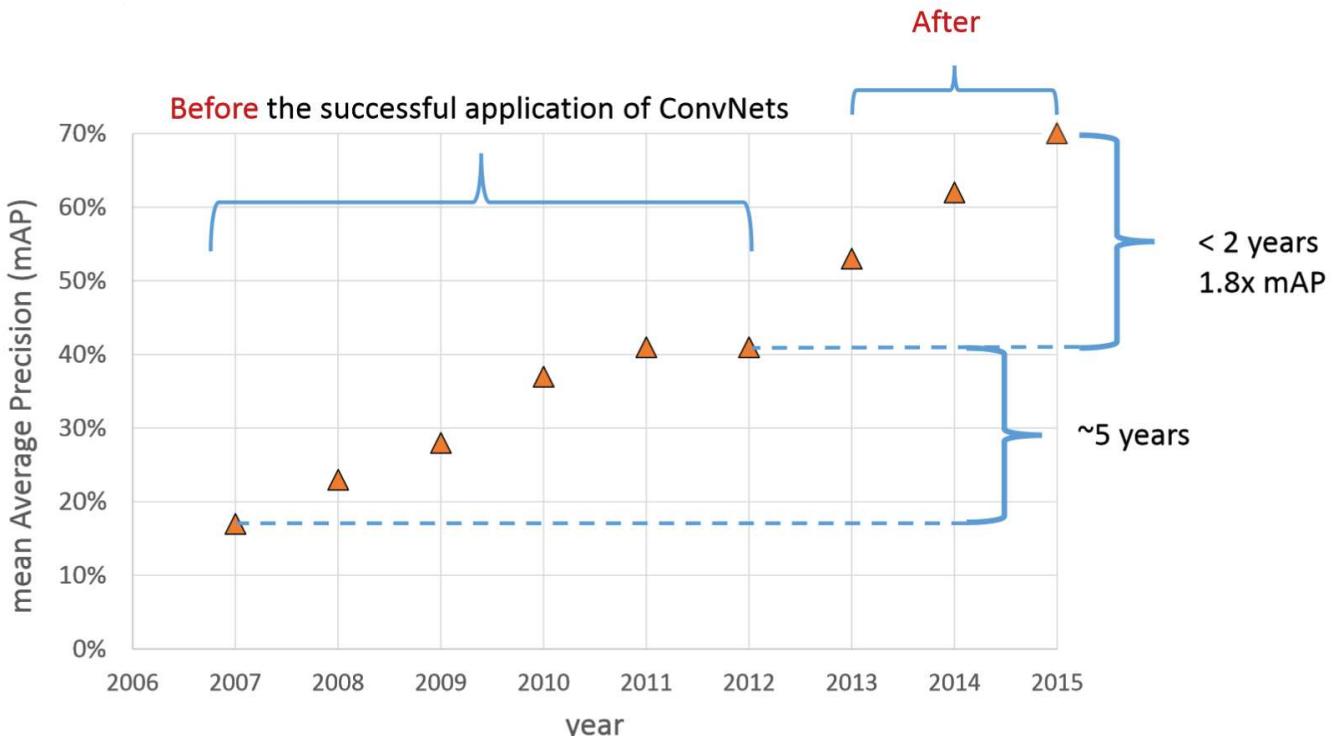
Since ~2012, Deep-Learning has brought very significant improvement over State-of-the-Art in Pattern Recognition and Image Semantic Analysis



- won many vision pattern recognition competitions (OCR, TSR, object categorization, facial expression,...)
- deployed in photo-tagging by Facebook, Google, Baidu,...

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 3

## Performance evolution of Pascal VOC object detection



Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 4

## ImageNet Large Scale Visual Recognition Challenge

- Public dataset and benchmark
- Worldwide research competition on image classification and visual objects detection & recognition/categorization

		PASCAL VOC 2012	ILSVRC 2013
Number of object classes		20	200
Training	Num images	5717	395909
	Num objects	13609	345854
Validation	Num images	5823	20121
	Num objects	13841	55502
Testing	Num images	10991	40152
	Num objects	---	---

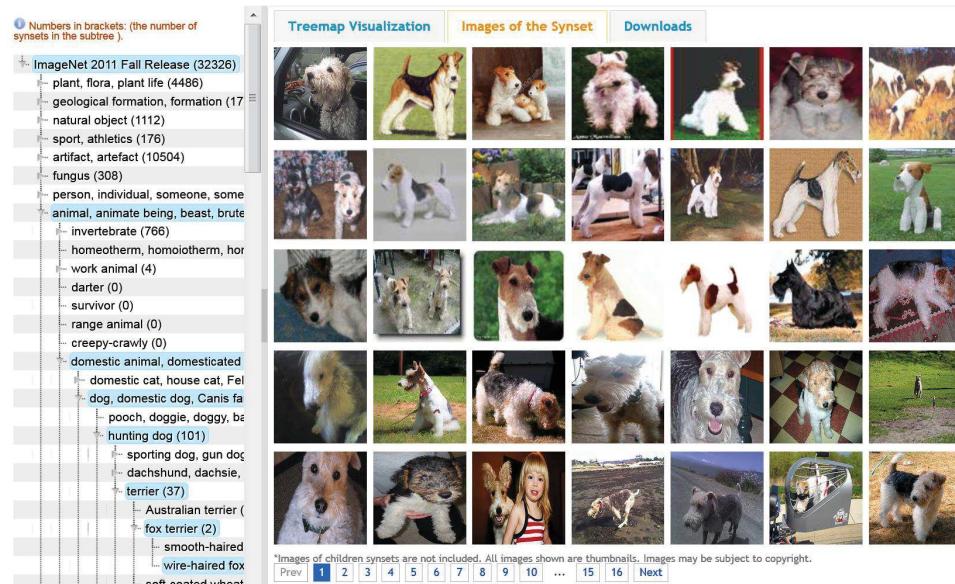
**Dramatic scale increase in image challenges in 2013**  
**Challenge won by Deep-Learning methods every year since 2012**

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 5

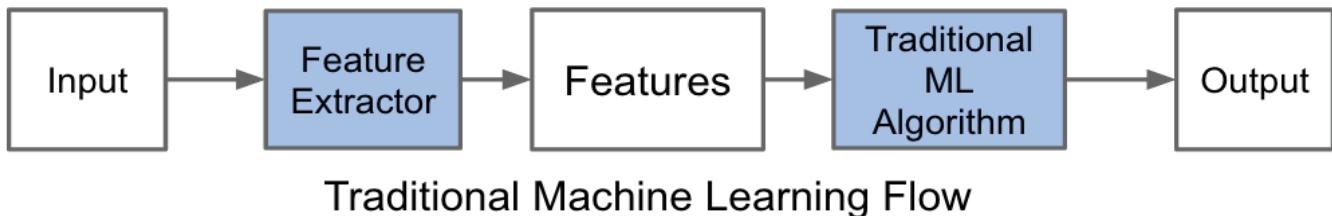
## ImageNet dataset

### Huge dataset of labelled images:

- 1000 classes of objects
- > 1 million labelled examples



# Importance of « features » in classical Machine-Learning

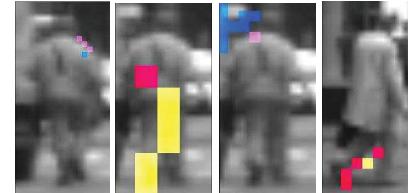


## Examples of *hand-crafted* features

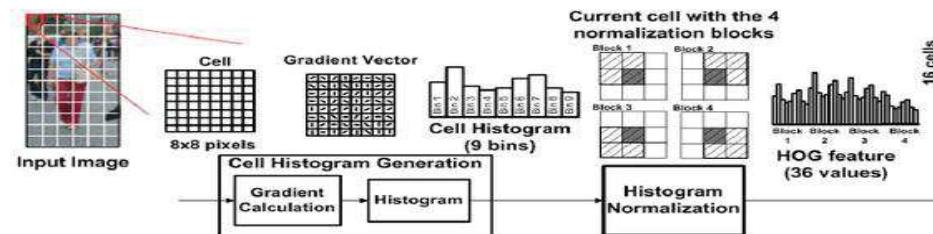
**Haar features**



**Control-points features**

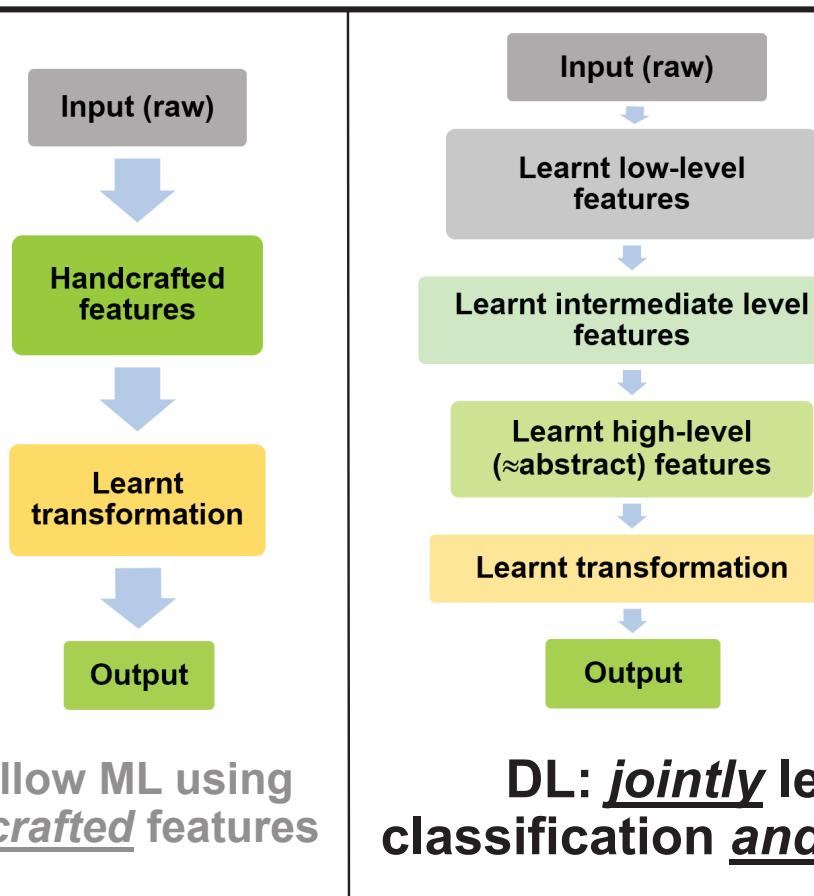


**HoG  
(Histogram  
of Gradients)**



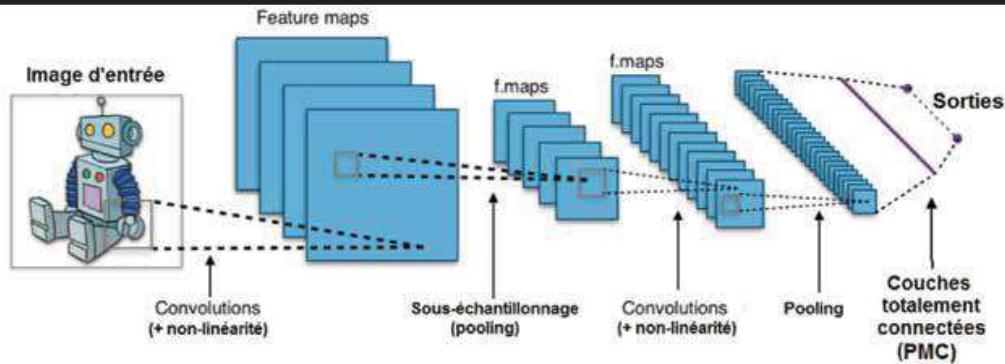
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 7

## Deep-Learning (DL) general principle

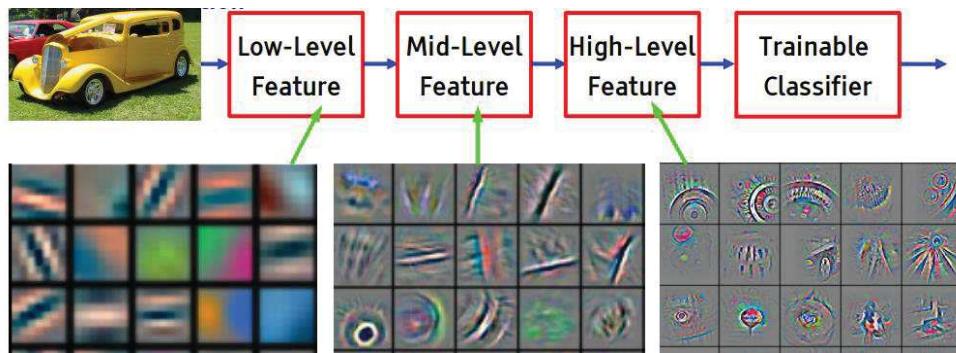


Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 8

# Convolutional Neural Networks (CNN, or ConvNet)

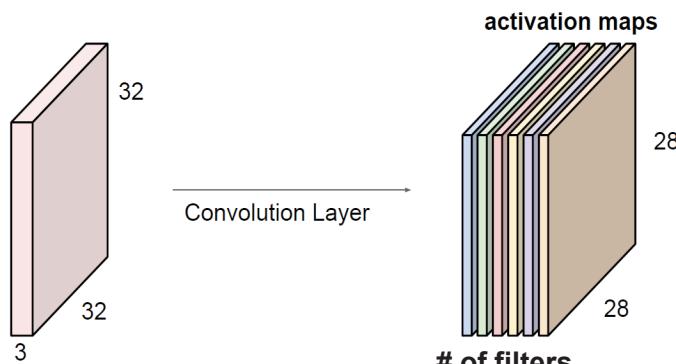


- For inputs with correlated dims (2D *image*, 1D signal,...)
- Succession of Convolutions and « pooling » layers



Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 9

## Convolution layers

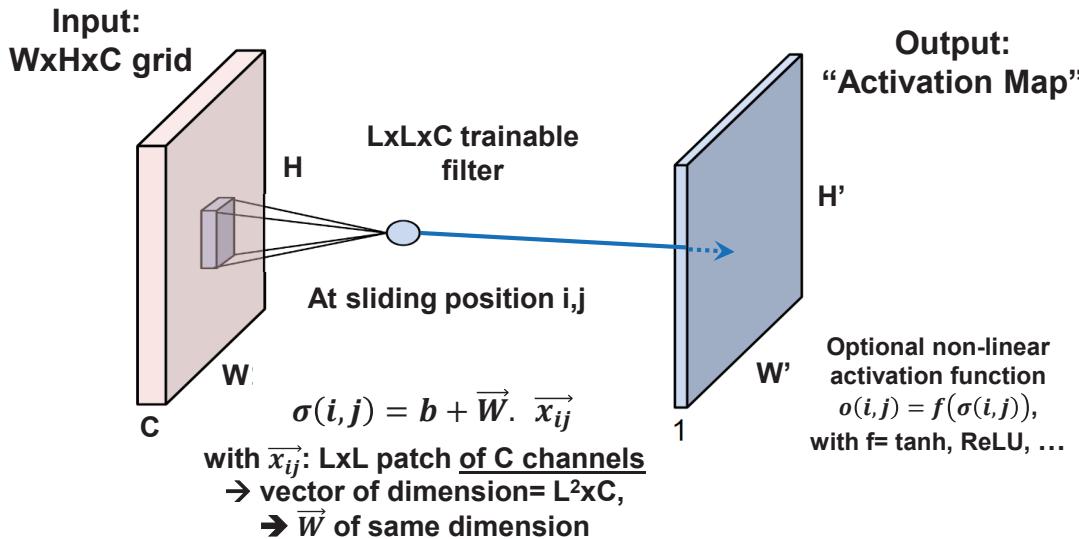


One “activation map”  
for each convolution filter



A Convolution layer applies several 3D filters to input image (or to input set of activation maps from previous layer)

# Convolution: sliding a 3D filter over a multi-channel image

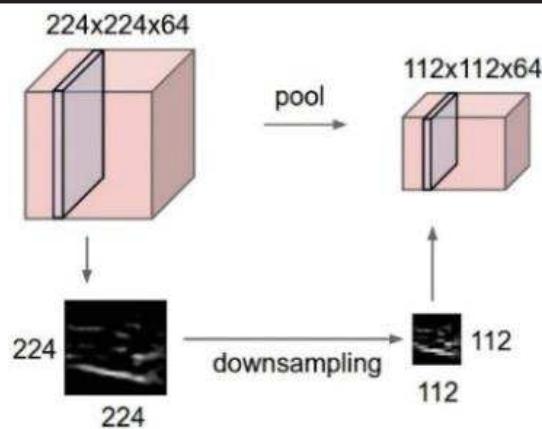


For each filter, a grid of  $W' \times H'$  neurons with shared weights  $\vec{W}$  (each neuron applies same filter at a different sliding position in input)

See *illustrative animation at:* <http://cs231n.github.io/convolutional-networks/>

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 11

## Pooling layers

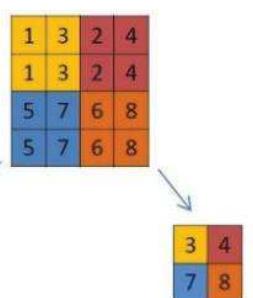


Pooling ≈ Downsampling

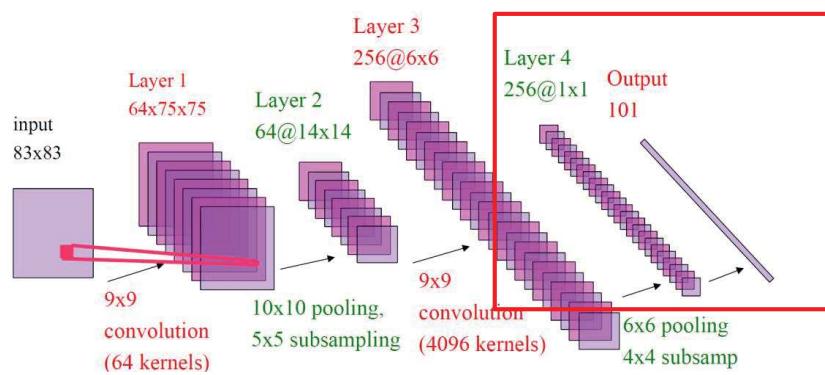
A pooling layer aggregates over space for:

- Dimension reduction
- Noise reduction
- Small translation and scaling invariance

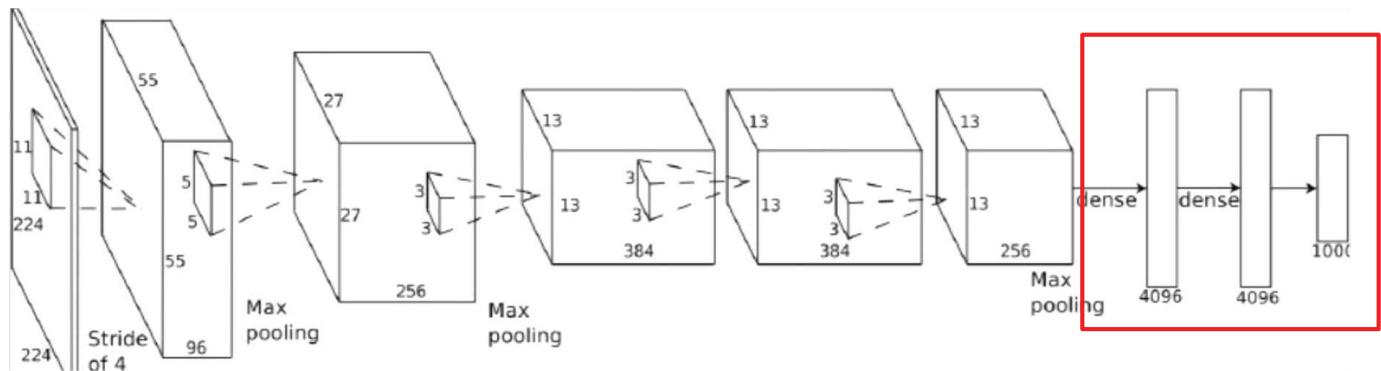
The pooling operation typically uses **average** or **max** on sets of  $2 \times 2$  (or  $p \times p$ ) pixels



# Final classification layer: often classical fully-connected MLP



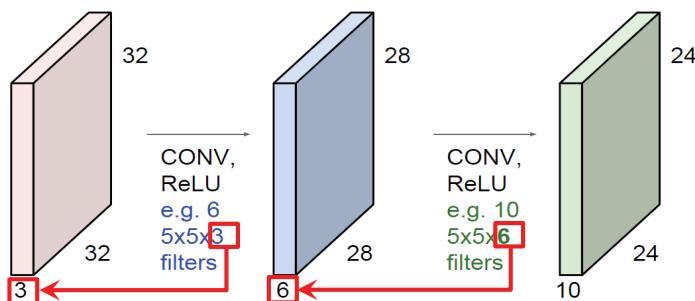
## AlexNet



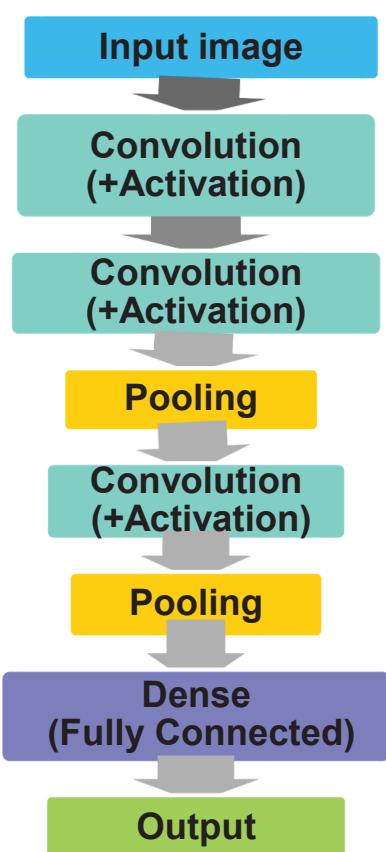
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 13

# Global architecture of a convNet

**Succession of Convolution (+ optional activation) layers and Pooling layers, which extract the hierarchy of features, followed by dense (fully connected) layer(s) for final classification**

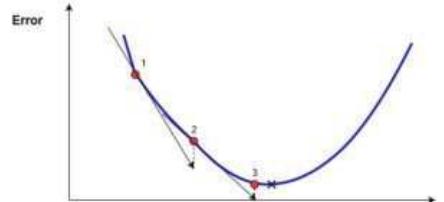


**NB: each convolution layer processes FULL DEPTH of previous set of activation maps**



All successive layers of a convNet forms a Deep neural network (with weigh-sharing inside each conv. Layer, and specific pooling layers).

Training = optimizing values of weights&biases  
Method used = gradient descent



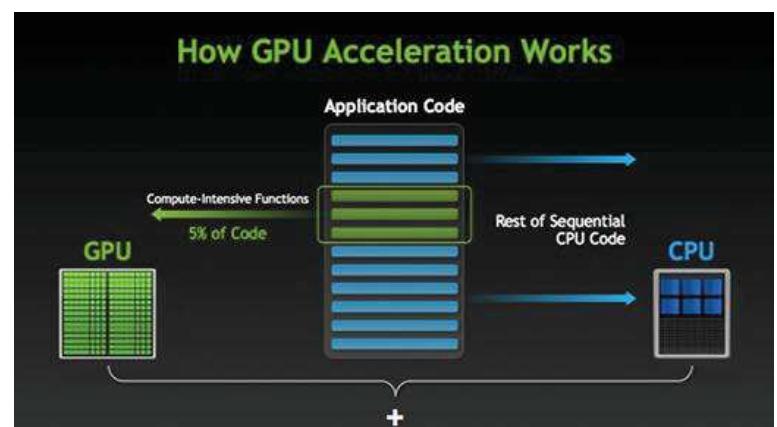
## → Stochastic Gradient Descent (SGD), using back-propagation:

- Input 1 (or a few) random training sample(s)
- Propagate
- Calculate error (loss)
- Back-propagate through all layers from end to input, to compute gradient
- Update convolution filter weights

Good convNets are very big (millions of parameters!)

Training generally performed on BIG datasets

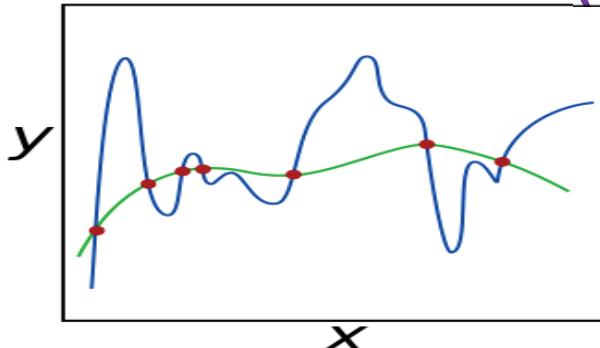
- Training can be extremely computer-intensive
- More manageable using **GPU (Graphical Processing Unit) acceleration** for ultra-parallel processing



- Importance of input normalization  
(zero mean, unit variance)
- Importance of weights initialization  
random but SMALL and prop. to  $1/\sqrt{\text{nbInputs}}$
- Decreasing (or adaptive) learning rate
- Importance of training set size  
ConvNets often have a LARGE number of free parameters  
→ train them with a sufficiently large training-set !
- Avoid overfitting by:
  - Use of L1 or L2 regularization (after some epochs)
  - Use « *Dropout* » regularization (esp. on large FC layers)

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 17

## Avoid overfitting using L1/L2 regularization « weight decay »)



Trying to fit too many free parameters with not enough information can lead to overfitting

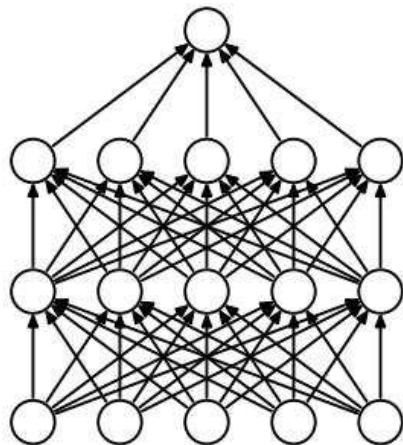
Regularization = penalizing too complex models

Often done by adding a special term to cost function

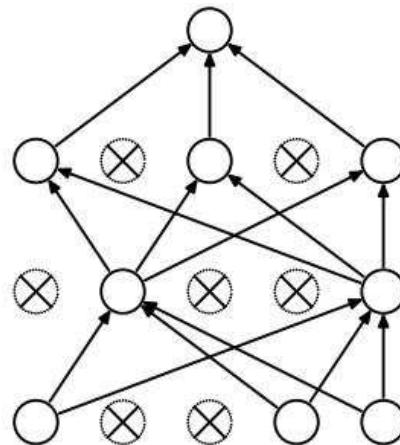
For neural network, the regularization term is just the L2- or L1- norm L2 of the vector of all weights:

$$K = \sum_m (\text{loss}(Y_m, D_m)) + \beta \sum_{ij} |W_{ij}|^p \quad \text{with } p=2 \text{ (L2) or } p=1 \text{ (L1)}$$

→ name “Weight decay”



(a) Standard Neural Net

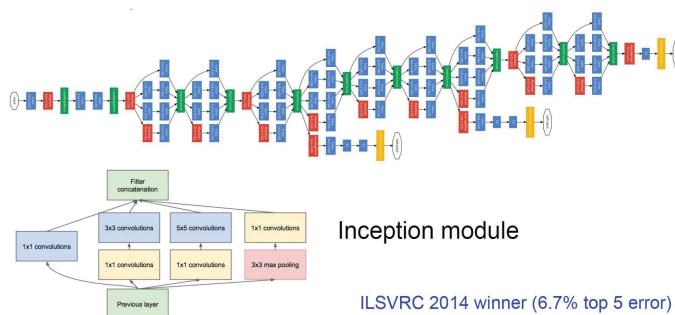


(b) After applying dropout.

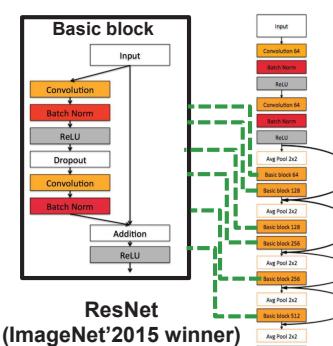
At each training stage, individual nodes can be temporarily "dropped out" of the net with probability  $p$  (usually  $\sim 0.5$ ), or re-installed with last values of weights

## ImageNet dataset and state-of-the-art convNets

- The most performant ConvNets have millions of trainable weights, so they must be trained on very large datasets
- Training on ImageNet, was an essential factor of their recognition performances



ILSVRC 2014 winner (6.7% top 5 error)

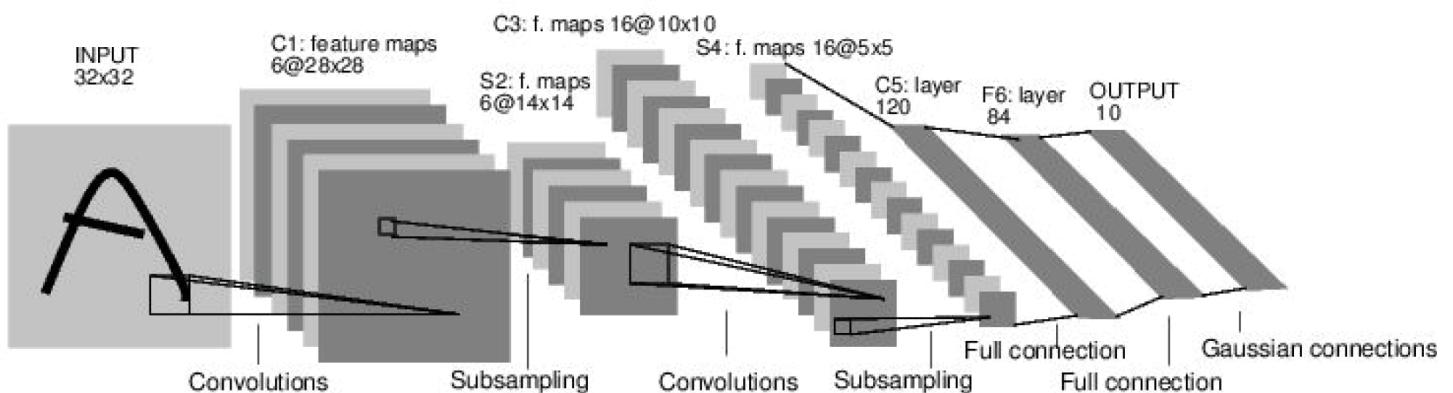


- Every year, a new ImageNet challenge winner, with better accuracy using new ConvNet architecture & algo
- Those general-purpose pre-trained image classifiers (AlexNet, GoogleNet\_Inception, ResNet, etc...) are publicly available

- **LeNet:** 1<sup>st</sup> successful applications of ConvNets, by Yann LeCun in 1990's. Used to read zip codes, digits, etc.
- **AlexNet:** Beginning of ConvNet “buzz”: largely outperformed competitors in ImageNet\_IISVRC2012 challenge. Developed by Alex Krizhevsky et al., architecture similar to LeNet (but deeper+larger, and some chained ConvLayers before Pooling). 60 M parameters !
- **GoogLeNet:** ILSVRC 2014 winner, developed by Google. Introduced an *Inception Module*, + AveragePooling instead of FullyConnected layer at output. Dramatic reduction of number of parameters (4M, compared to AlexNet with 60M).
- **VGGNet:** Runner-up in ILSVRC 2014. Very deep (16 CONV/FC layers) → 140M parameters !!
- **ResNet:** ILSVRC 2015, “Residual Network” introducing “skip” connections. Currently ~ SoA in convNet. Very long training but fast execution.

## LeNet, for digits/letters recognition [LeCun et al., 1998]

Input: 32x32 image



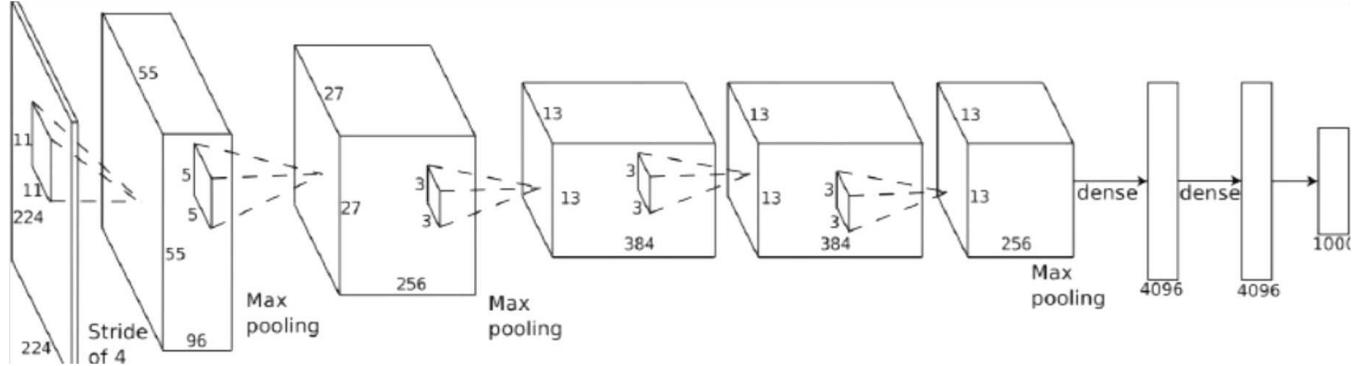
Conv filters were 5x5, applied at stride 1

Subsampling (Pooling) layers were 2x2 applied at stride 2  
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

# AlexNet, for image categorisation

[Krizhevsky et al. 2012]

**Input: 224x224x3 image**

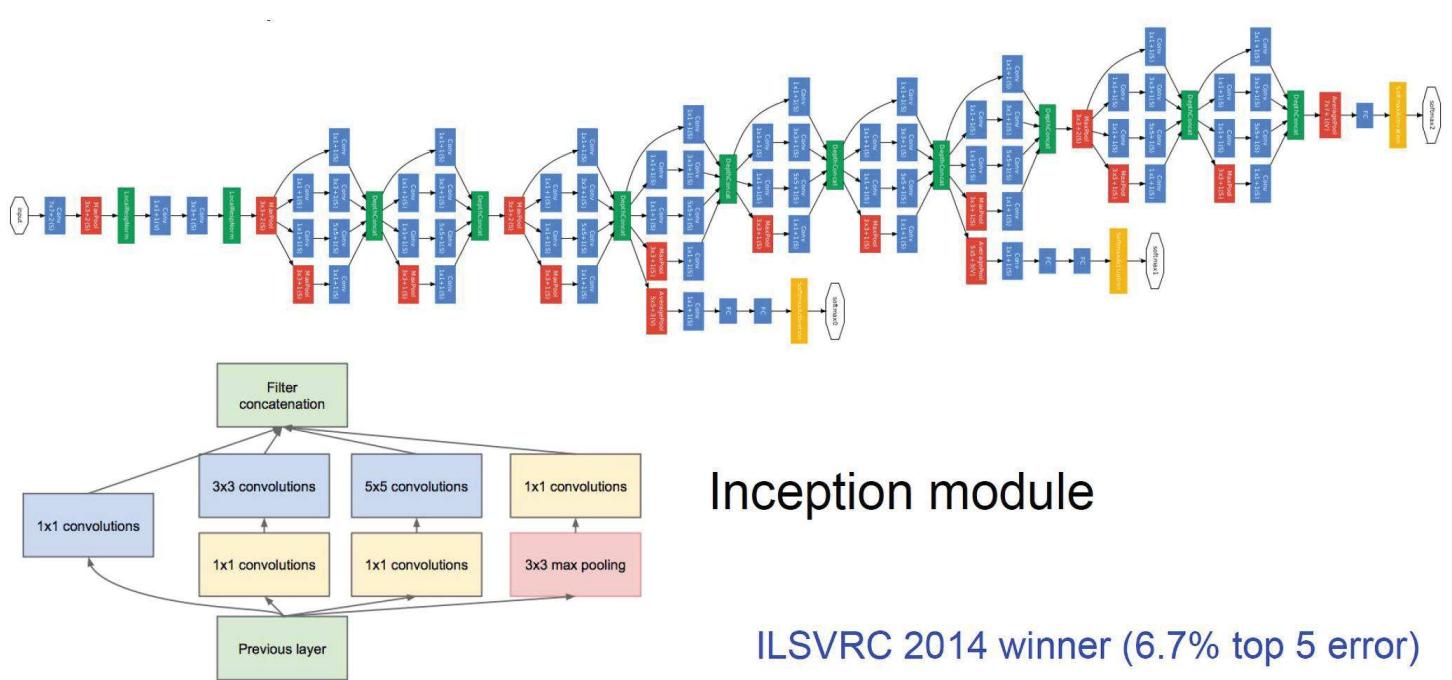


**60 million parameters !...**

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 23

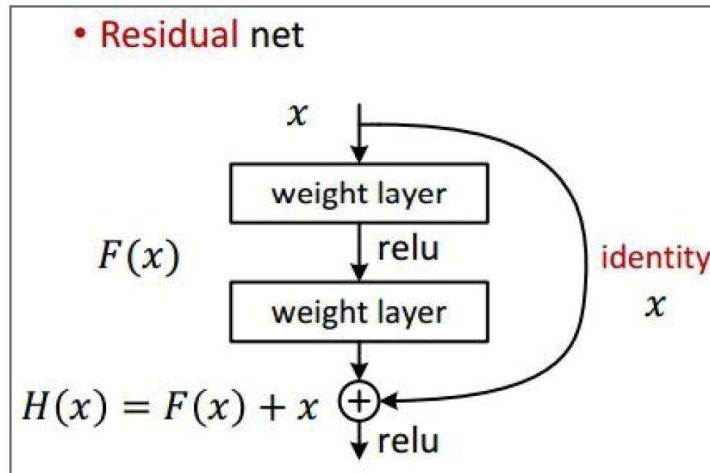
# GoogleNet

[Szegedy et al., 2014]



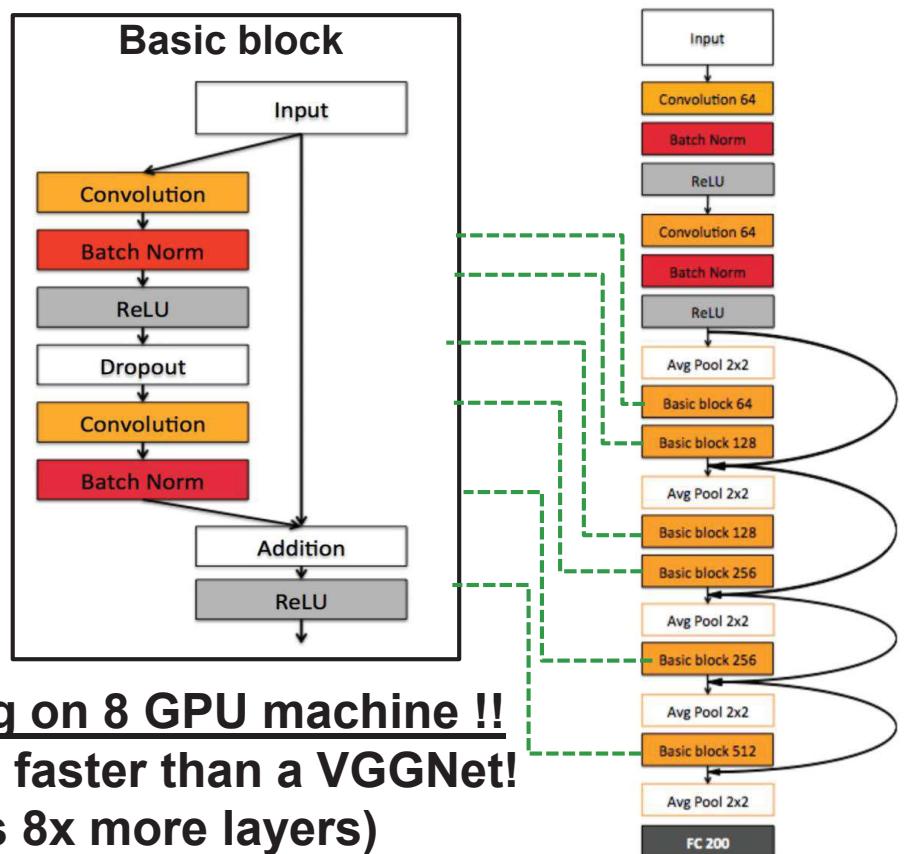
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 24

- ILSVRC 2015 large winner in 5 main tracks (3.6% top 5 error)
- 152 layers!!!
- But novelty = "skip" connections



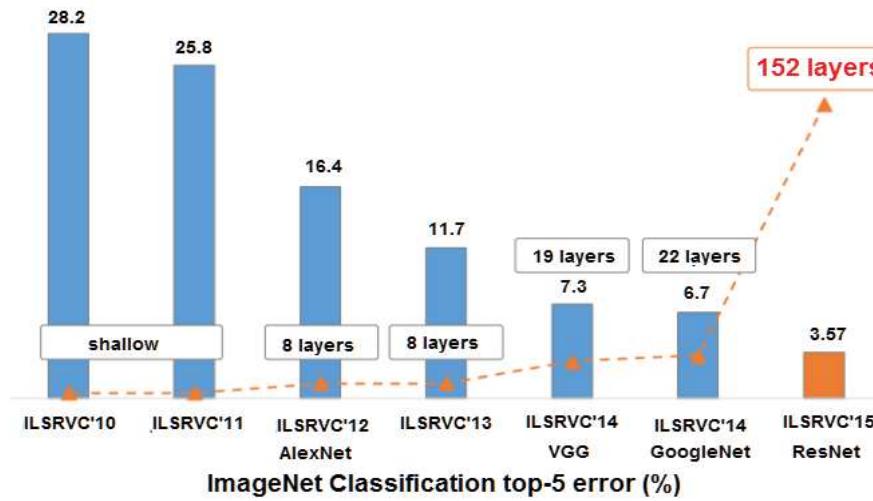
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 25

## ResNet global architecture



- 2-3 weeks of training on 8 GPU machine !!
- However, at runtime faster than a VGGNet! (even though it has 8x more layers)

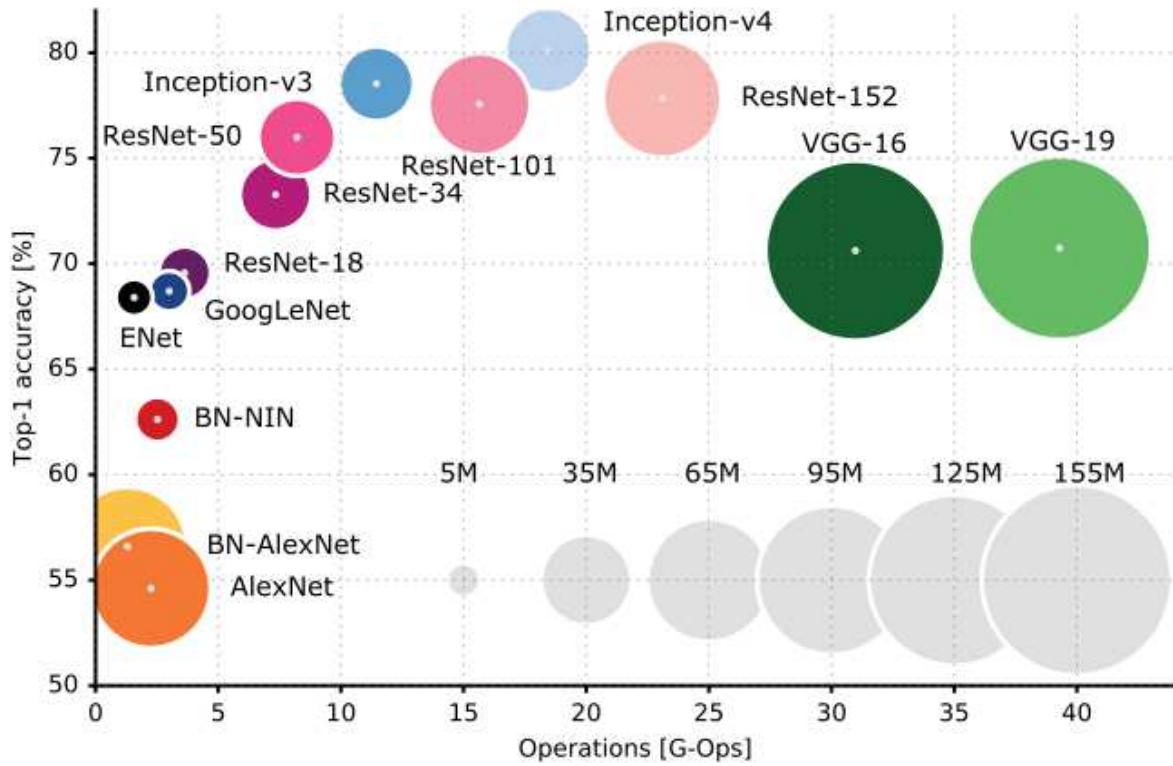
# Summary of recent ConvNet history



**But most important is the choice of  
ARCHITECTURAL STRUCTURE**

**+ More and more modified architectures  
for tasks other than just image classification**

## Performance comparisons of common pre-trained convNets



- **TensorFlow** <https://www.tensorflow.org>
- **KERAS** <https://keras.io>  
Python front-end APIs mapped either  
on Tensor-Flow or Theano back-end
- **PyTorch** <https://pytorch.org/>
- **Caffe** <http://caffe.berkeleyvision.org/>  
*C++ library, hooks from Python → notebooks*
- **Theano** <http://www.deeplearning.net/software/theano/>
- **Lasagne** <http://lasagne.readthedocs.io>  
*lightweight library to build+train neural nets in Theano*

All of them handle transparent use of GPU,  
and most of them are used in Python code/notebook

## Example of convNet code in Keras

```

model = Sequential()
# 1 set of (Convolution+Pooling) layers, with Dropout
model.add(Convolution2D(conv_depth_1, kernel_size, kernel_size,
                       border_mode='valid', input_shape=(depth, height, width)))
model.add( MaxPooling2D(pool_size=(pooling_size, pooling_size)) )
model.add(Activation('relu'))
model.add(Dropout(drop_prob))

# Now flatten to 1D, and apply 1 Fully_Connected layer
model.add(Flatten())
model.add(Dense(hidden_size1, init='lecun_uniform'))
model.add(Activation('sigmoid'))

# Finally add a Softmax output layer, with 1 neuron per class
model.add(Dense(num_classes, init='lecun_uniform'))
model.add(Activation('softmax'))

# Training "session"
sgd = SGD(lr=learning_rate, momentum=0.8) # Optimizer
model.compile(loss='categorical_crossentropy', optimizer=sgd)
model.fit(X_train, Y_train, batch_size=32, nb_epoch=2, verbose=1,
           validation_split=valid_proportion)

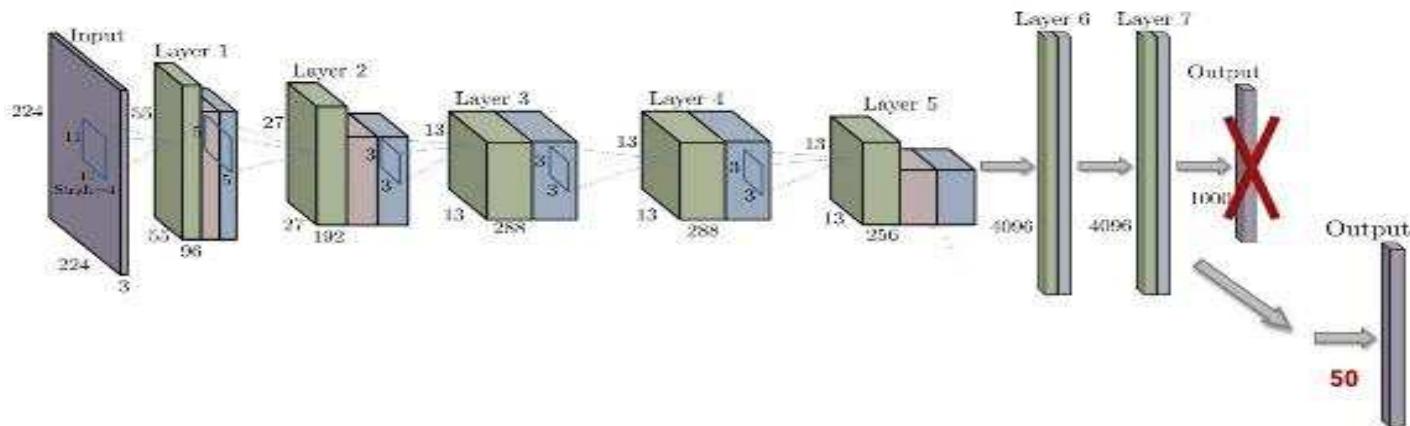
# Evaluate the trained model on the test set
model.evaluate(X_test, Y_test, verbose=1)

```

- **Recalls on Convolutional Neural Networks (CNN or ConvNets) and Deep-Learning**
- **Transfer Learning**
- **Beyond Image Classification: DETECTION OF OBJECTS**
- **Instance segmentation with DeepLearning**
- **DL for Human pose inference and depth estimation**
- **Semantic segmentation with DeepLearning**
- **Interest and use of simulations / synthetic videos**

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 31

## Generality of learnt representation + Transfer learning

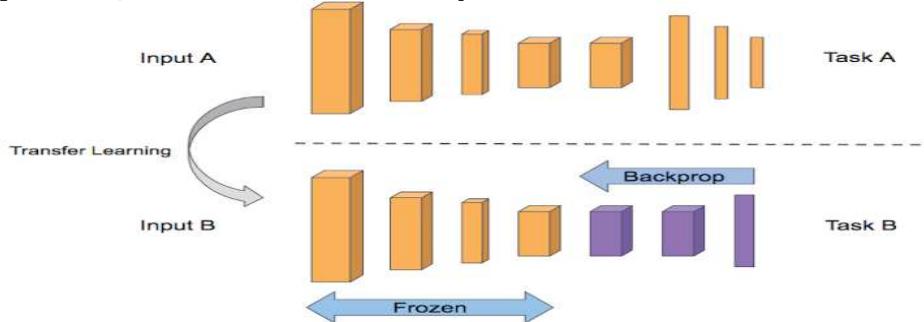


By removing last layer(s) (those for classification) of a convNet trained on ImageNet, one obtains a transformation of any input image into a semi-abstract representation, which can be used for learning SOMETHING ELSE (« transfer learning »):

- either by just using learnt representation as features
- or by creating new convNet output and perform learning of new output layers + fine-tuning of re-used layers

# Transfer learning

- SoA convNets winning ImageNet are image CLASSIFIERS for one object per image
  - Many object categories can be irrelevant (e.g. boat when onboard a car)
- For each particular application, models are usually obtained from state-of-the-art ConvNets pre-trained on ImageNet (winners of yearly challenge, eg: AlexNet, VGG, Inception, ResNet, etc...)

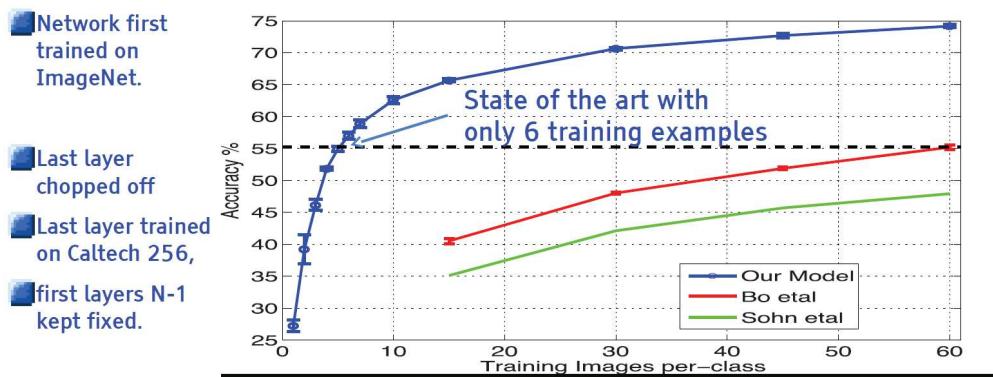


- Adaptation is performed by Transfer Learning, ie modification+training of last layers and/or fine-tuning of pre-trained weights of lower layers

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 33

# Transfer Learning and fine-tuning

- Using a CNN pre-trained on a large dataset, possible to adapt it to another task, using only a SMALL training set!



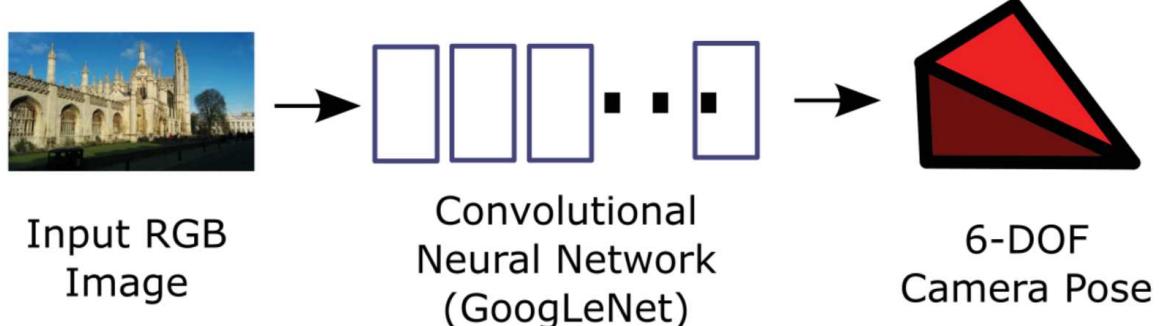
# Train	Acc % 15/class	Acc % 30/class	Acc % 45/class	Acc % 60/class
Sohn et al. [16]	35.1	42.1	45.7	47.9
Bo et al. [3]	$40.5 \pm 0.4$	$48.0 \pm 0.2$	$51.9 \pm 0.2$	$55.2 \pm 0.3$
Non-pretr.	$9.0 \pm 1.4$	$22.5 \pm 0.7$	$31.2 \pm 0.5$	$38.8 \pm 1.4$
ImageNet-pretr.	$65.7 \pm 0.2$	$70.6 \pm 0.2$	$72.7 \pm 0.4$	$74.2 \pm 0.3$

3: [Bo, Ren, Fox. CVPR, 2013] 16: [Sohn, Jung, Lee, Hero ICCV 2011]

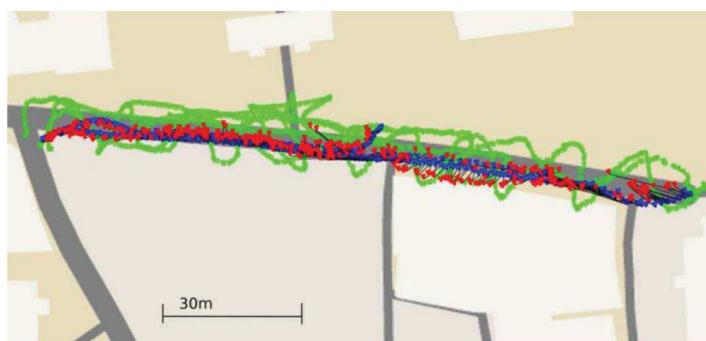
- Learning on simulated synthetic images + fine-tuning on real-world images
- Recognition/classification for OTHER categories or classes
- Training an objects detector (or a semantic segmenter)
  
- Precise localization (position+bearing) = PoseNet
- End-to-end driving (imitation Learning)
- 3D informations (depth map) from monovision!

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 35

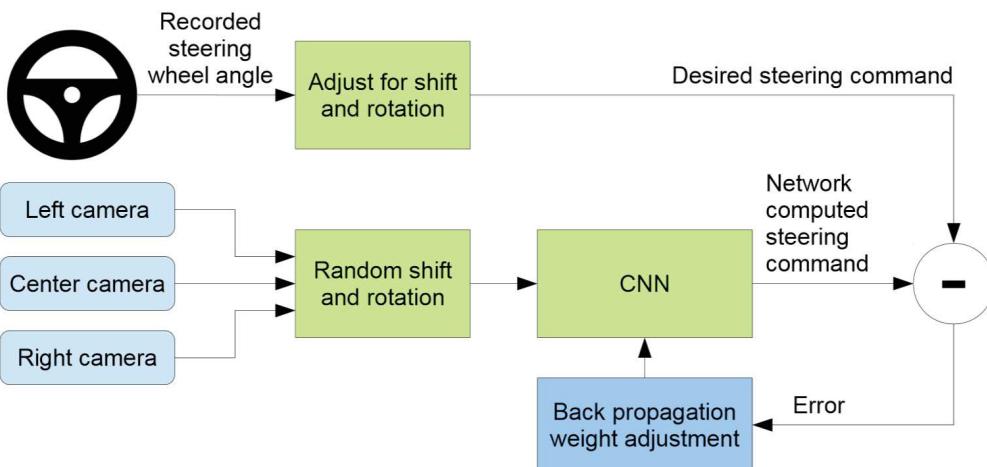
## Transfer-Learning for 6-DOF Camera Relocalization



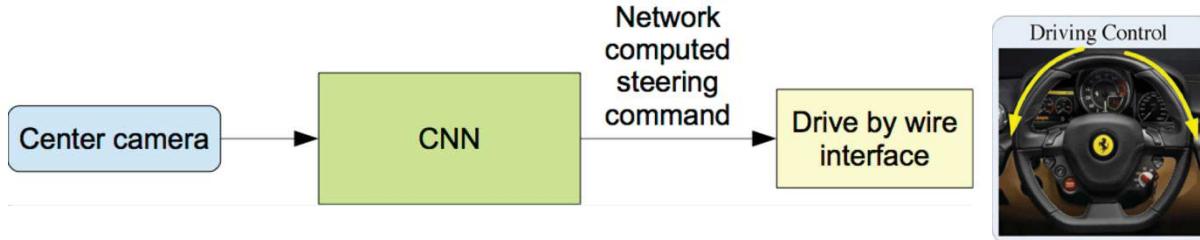
[A. Kendall, M. Grimes & R. Cipolla, "PoseNet: A Convolutional Network for Real-Time 6-DOF Camera Relocalization«, ICCV'2015, pp. 2938-2946]



King's College



- **End-to-end driving by «*imitation Learning*»**  
**(nVidia, Valeo)**



Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 37

## Transfer Learning code example in Keras

```

from keras.applications.inception_v3 import InceptionV3
from keras.preprocessing import image
from keras.models import Model
from keras.layers import Dense, GlobalAveragePooling2D
from keras import backend as K
# create the base pre-trained model
base_model = InceptionV3(weights='imagenet',
                           include_top=False)

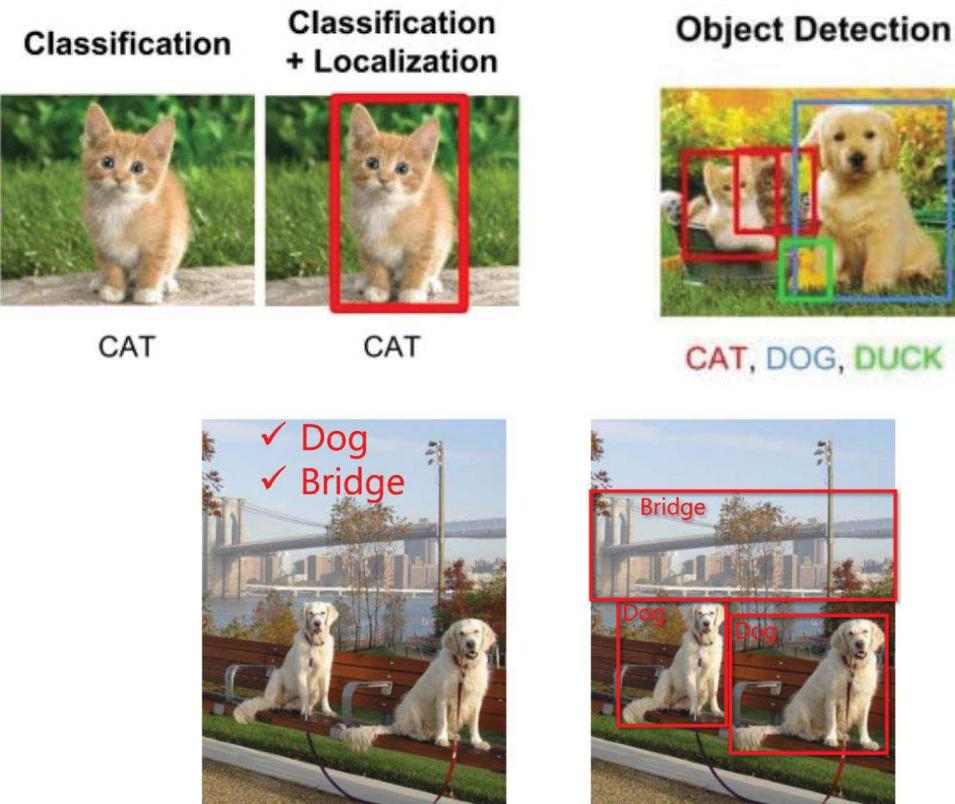
# add a global spatial average pooling layer
x = base_model.output
x = GlobalAveragePooling2D()(x)
# let's add a fully-connected layer
x = Dense(1024, activation='relu')(x)
# and a logistic layer -- let's say we have 200 classes
predictions = Dense(200, activation='softmax')(x)
# this is the model we will train
model = Model(input=base_model.input, output=predictions)
# first: train only the top layers (which were randomly initialized)
# i.e. freeze all convolutional InceptionV3 layers
for layer in base_model.layers:
    layer.trainable = False
# compile the model (should be done *after* setting layers to non-trainable)
model.compile(optimizer='rmsprop', loss='categorical_crossentropy')
# train the model on the new data for a few epochs
model.fit_generator(...)

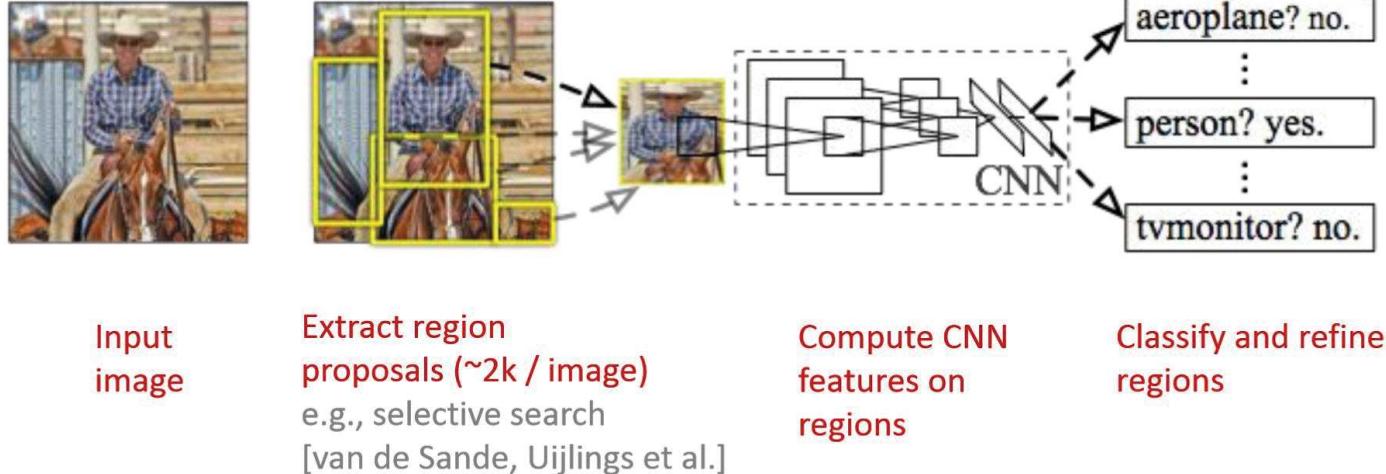
  
```

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 38

- **Recalls on Convolutional Neural Networks (CNN or ConvNets) and Deep-Learning**
- **Transfer Learning**
- **Beyond Image Classification: DETECTION OF OBJECTS**
- **Instance segmentation with DeepLearning**
- **DL for Human pose inference and depth estimation**
- **Semantic segmentation with DeepLearning**
- **Interest and use of simulations / synthetic videos**

## Classification vs Detection

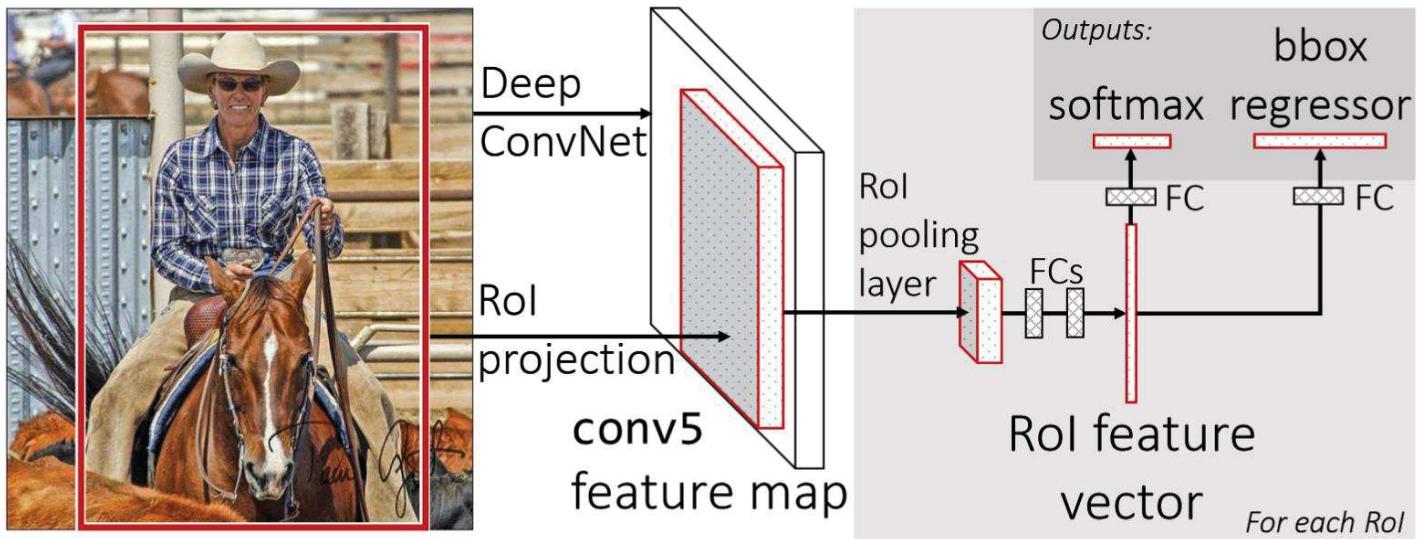




**Very slow + rather approximate bounding-boxes**

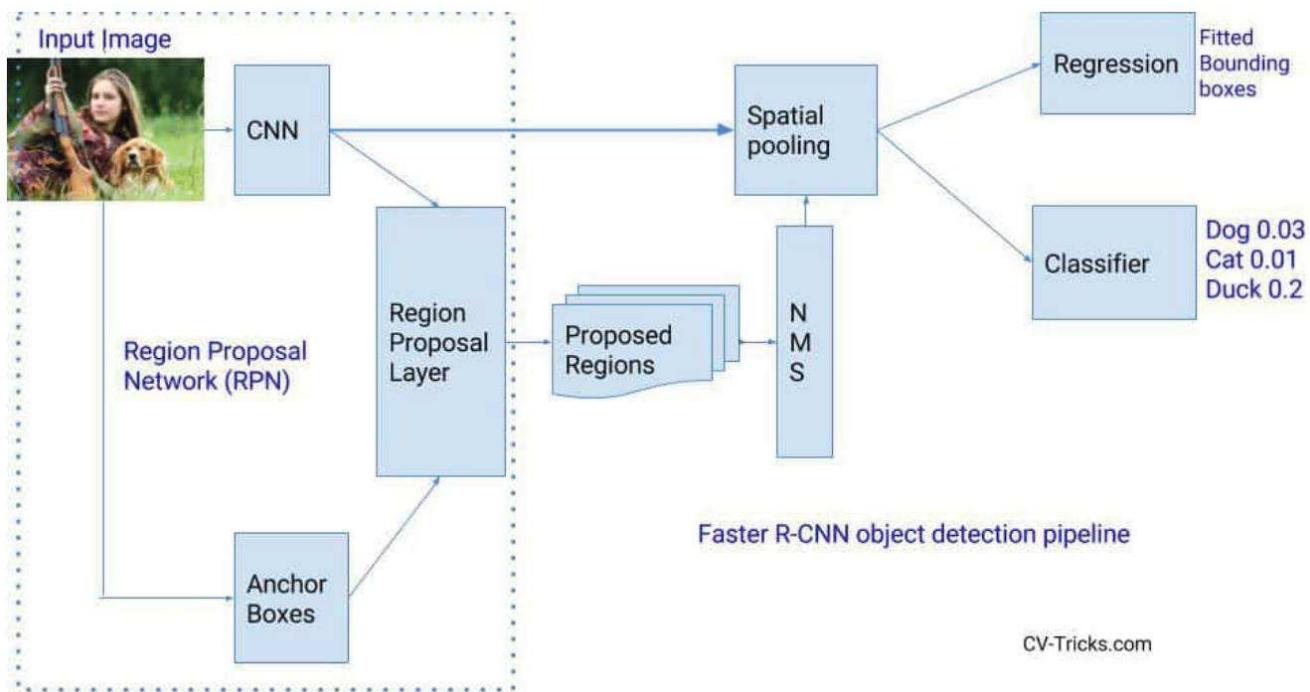
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 41

## Better: Fast R-CNN



**Learn a bounding-box regressor together with the class estimation (combination of 2 losses)**

# Faster\_R-CNN

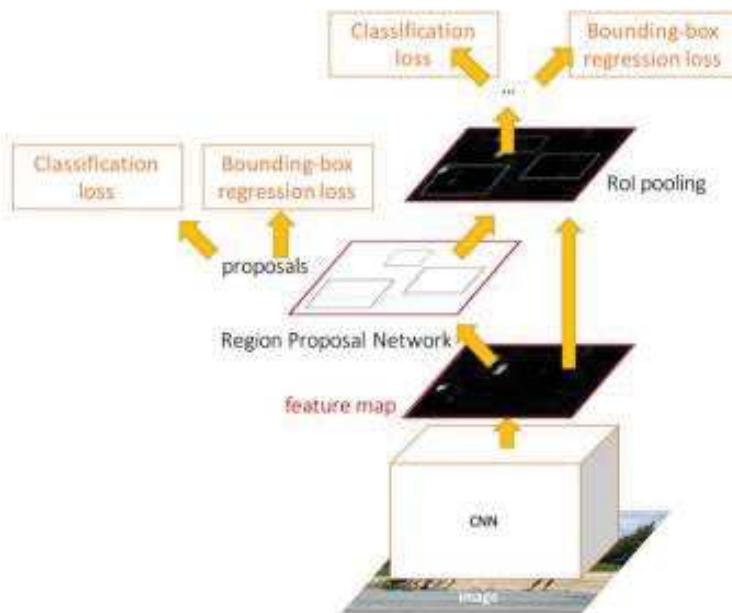


CV-Tricks.com

**Learn also a « Region Proposal Network »  
→ objects' bounding-boxes.**

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 43

# Training of Faster\_RCNN

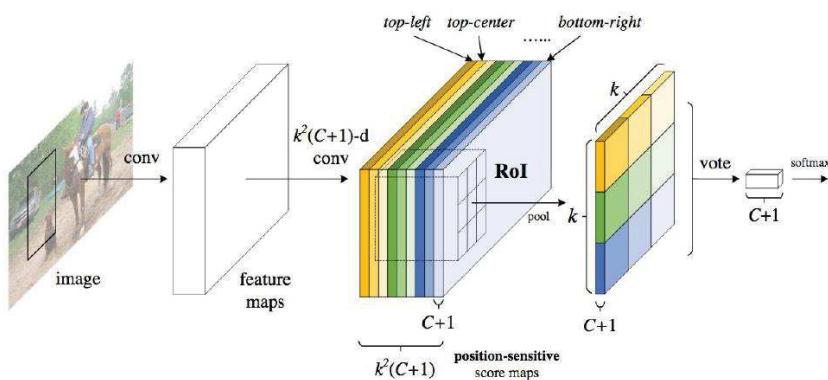


**Combining 4 losses!**

- **History**
  - **R-CNN**: Selective search → Cropped Image → CNN
  - **Fast R-CNN**: Selective search → Crop feature map of CNN
  - **Faster R-CNN**: CNN → Region-Proposal Network  
→ Crop feature map of CNN
- **Best performances, but longest run-time**
- **End-to-end, multi-task loss**

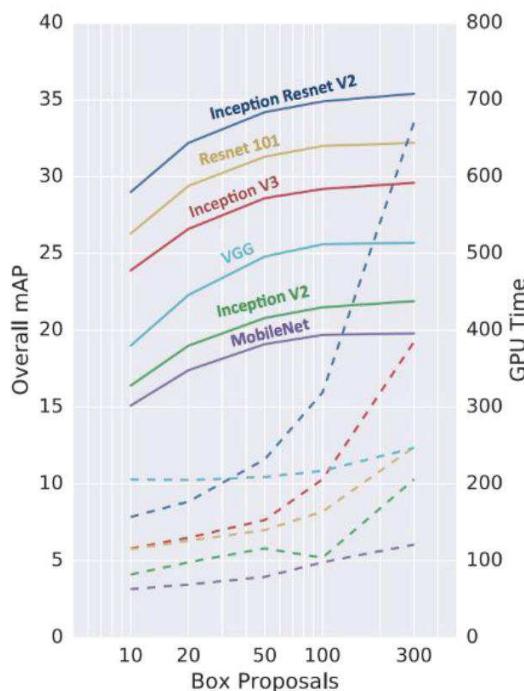
[<https://github.com/endernewton/tf-faster-rcnn>]

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 45

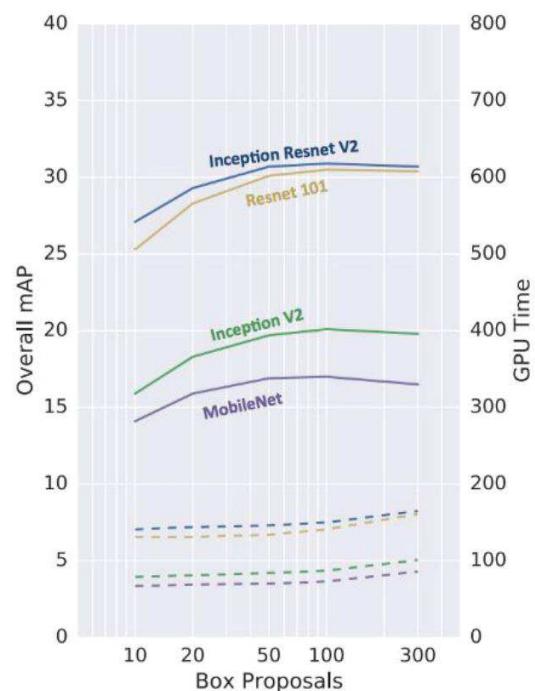


- **Addresses translation-variance in detection**
  - Position-sensitive ROI-pooling
- **Good balance between speed & performance**
  - 2.5 - 20x faster than Faster R-CNN

<https://github.com/daijifeng001/R-FCN>



(a) FRCNN



(b) RFCN

Image from: <https://arxiv.org/pdf/1611.10012.pdf>

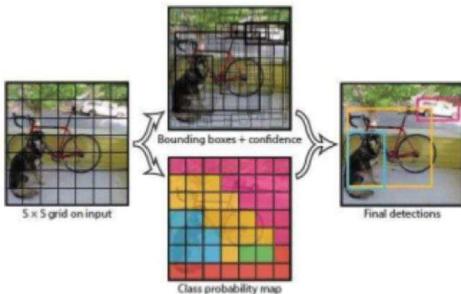
Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 47

## Example video of objects visual simultaneous detection and categorization with R-CNN

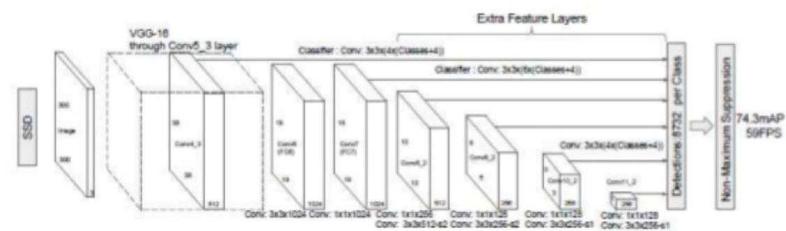


# Solve detection as a regression problem (“single-shot” detection)

## YOU ONLY LOOK ONCE(YOLO)



## SINGLE SHOT MULTIBOX DETECTOR(SSD)

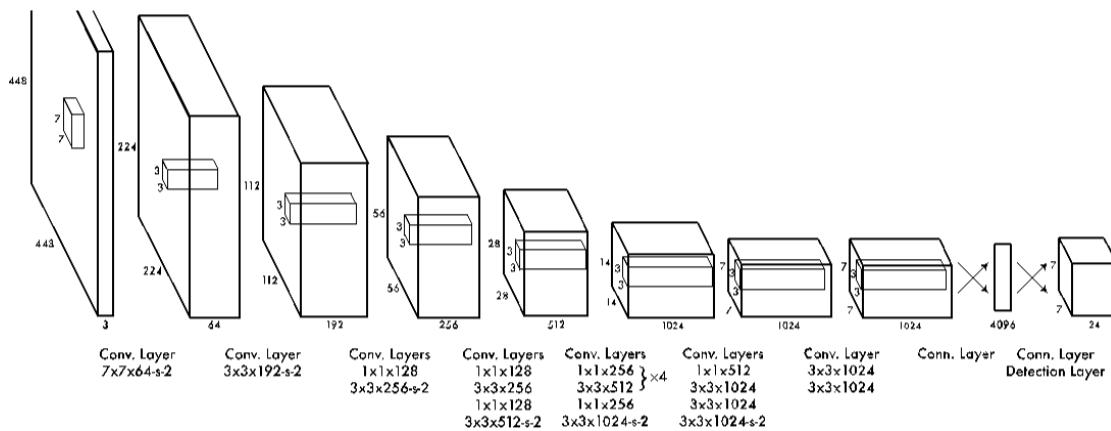


Images from: <https://www.slideshare.net/TaegyunJeon1/pr12-you-only-look-once-yolo-unified-realtime-object-detection>

Deep Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 49

YOLO

# You Only Look Once



- **Modified GoogleNet/Inception**
  - **Super fast (21~155 fps)**
  - **Finds objects in image grids *in parallel***
  - **Only slightly worse performance than Faster R-CNN**

## Unified Detection

- All BBox, All classes

1) Image  $\rightarrow S \times S$  grids

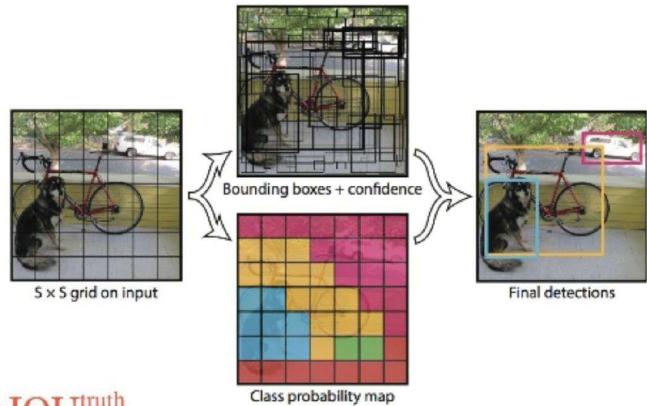
2) grid cell

$\rightarrow B$ : BBoxes and Confidence score



$\rightarrow C$ : class probabilities w.r.t #classes

$Pr(Class_i | Object)$



**Figure 2: The Model.** Our system models detection as a regression problem. It divides the image into an  $S \times S$  grid and for each grid cell predicts  $B$  bounding boxes, confidence for those boxes, and  $C$  class probabilities. These predictions are encoded as an  $S \times S \times (B * 5 + C)$  tensor.

Slide from: <https://www.slideshare.net/TaegyunJeon1/pr12-you-only-look-once-yolo-unified-realtime-object-detection>

Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 51

## YOLO limitations and v2 improvements

- Groups of small objects
- Unusual aspect ratios
- Localization error of bounding boxes

→ YOLOv2: Many improvements  
+ Custom architecture – Darknet  
(instead InceptionNet for YOLO)

YOLO (darknet) - <https://pjreddie.com/darknet/yolov1/> (C++)

YOLO v2 (darknet) - <https://pjreddie.com/darknet/yolov2/> (C++)

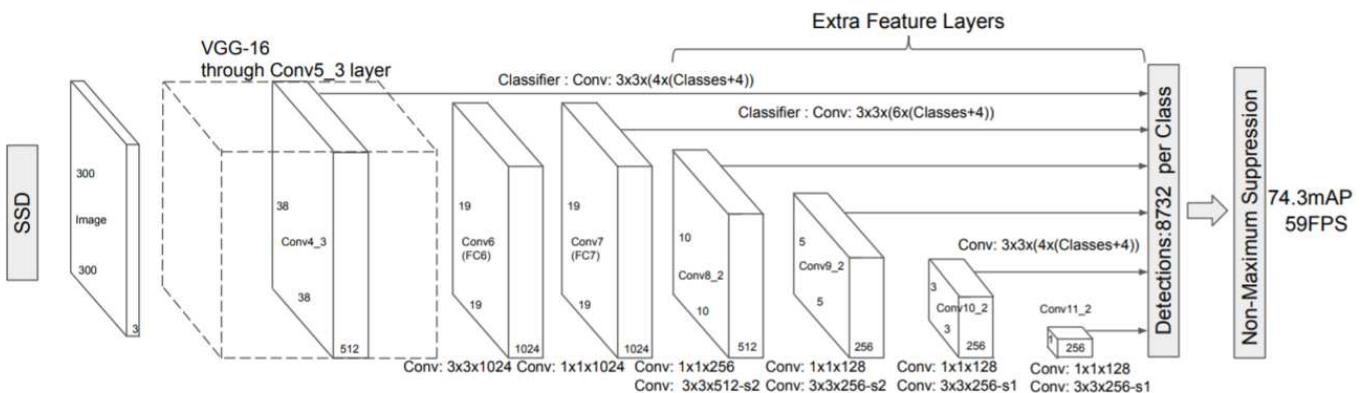
- Better and faster - 91 fps for 288 x 288

YOLO v3 (darknet) - <https://pjreddie.com/darknet/yolo/> (C++)

YOLO (caffe) - <https://github.com/xingwangsfu/caffe-yolo>

YOLO (tensorflow) - <https://github.com/thtrieu/darkflow>

## Architecture



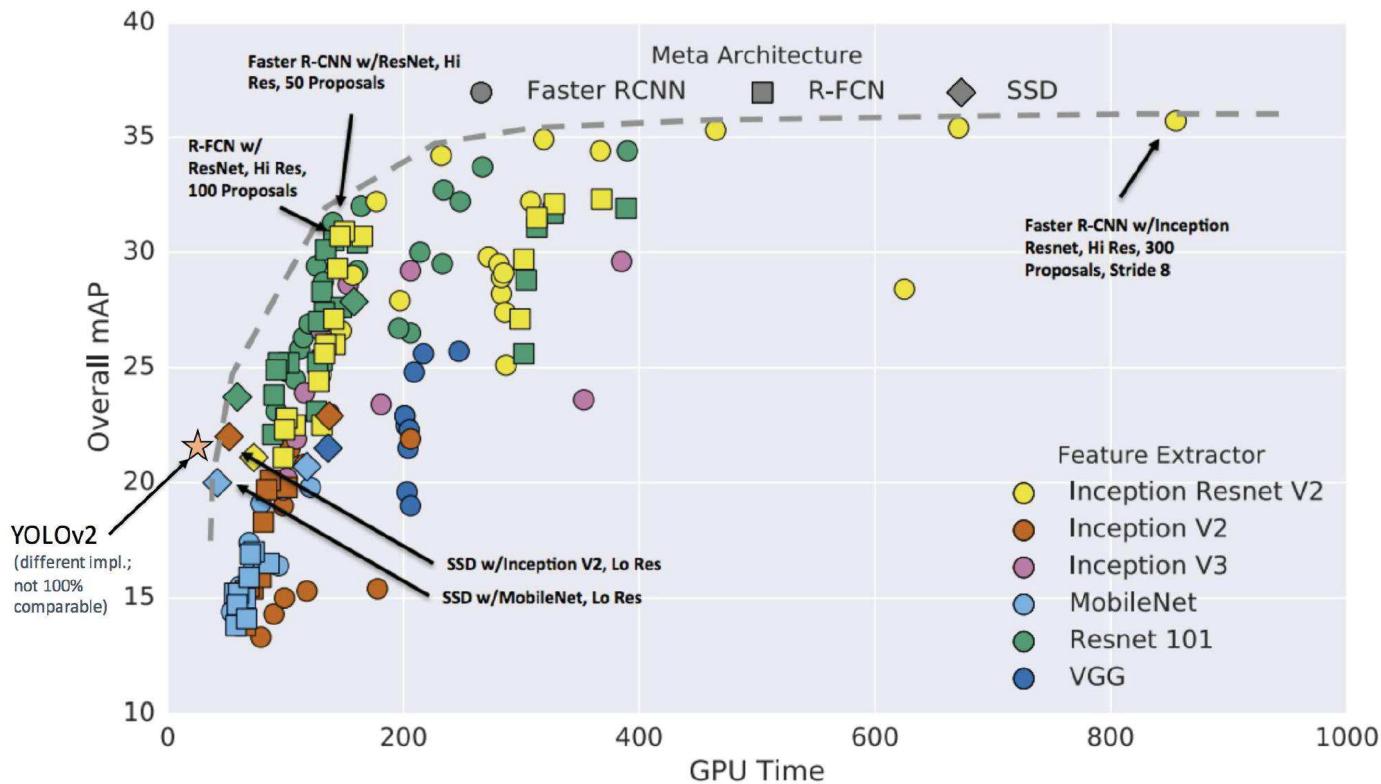
**Slower but more accurate than YOLO**  
**Faster but less accurate than Faster\_R-CNN**

SSD (caffe) - <https://github.com/weiliu89/caffe/tree/ssd>

SSD (tensorflow) - <https://github.com/balancap/SSD-Tensorflow>

SSD (pytorch) - <https://github.com/amdegroot/ssd.pytorch>

## Recent comparison of convNets for object detection



# Training sets for Visual objects detection

- Training a visual objects detector requires a training set containing images WITH BOUNDING-BOXES (or even mask) ANNOTATION
- Two main « reference » training sets of this type:
  - Pascal VOC (Visual Object Class)  
<http://host.robots.ox.ac.uk/pascal/VOC/>
  - Coco (Common Objects in Context)  
*[more classes + MASK annotations]*  
<http://cocodataset.org/>



Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 57

## VOC and COCO categories

### VOC categories

aeroplane  
bicycle  
bird  
boat  
bottle  
bus  
car  
cat  
chair  
cow  
diningtable  
dog  
horse  
motorbike  
person  
pottedplant  
sheep  
sofa  
train  
tvmonitor

### COCO categories

person	backpack	Apple	microwave
bicycle	umbrella	Sandwich	oven
car	handbag	Orange	toaster
motorbike	tie	broccoli	sink
aeroplane	suitcase	carrot	refrigerator
bus	frisbee	hot dog	book
train	skis	pizza	clock
truck	snowboard	donut	vase
boat	sports ball	cake	scissors
traffic light	Kite	chair	teddy bear
fire hydrant	baseball bat	Sofa	hair drier
stop sign	baseball glove	pottedplant	toothbrush
parking meter	skateboard	bed	
bench	surfboard	diningtable	
bird	tennis racket	toilet	
cat	Bottle	Tvmonitor	
Dog	wine glass	laptop	
horse	cup	mouse	
sheep	fork	remote	
Cow	knife	keyboard	
elephant	spoon	cell phone	
bear	bowl		
zebra	banana		
giraffe			

- If very fast inference is essential, better choose latest version of YOLO (or even MobileNet)
- If quality of detections (precision and recall) is more important, better choose Faster\_RCNN
- For a compromise, SSD can be considered
- Pre-trained ConvNet detectors are available for many pre-defined categories (those of VOC or COCO)

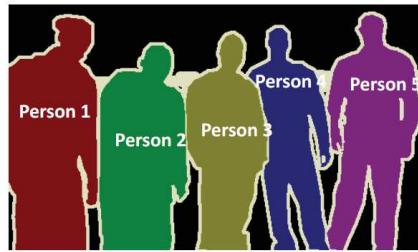
## Outline

- Recalls on Convolutional Neural Networks (CNN or ConvNets) and Deep-Learning
- Transfer Learning
- Beyond Image Classification: DETECTION OF OBJECTS
- Instance segmentation with DeepLearning
- DL for Human pose inference and depth estimation
- Semantic segmentation with DeepLearning
- Interest and use of simulations / synthetic videos

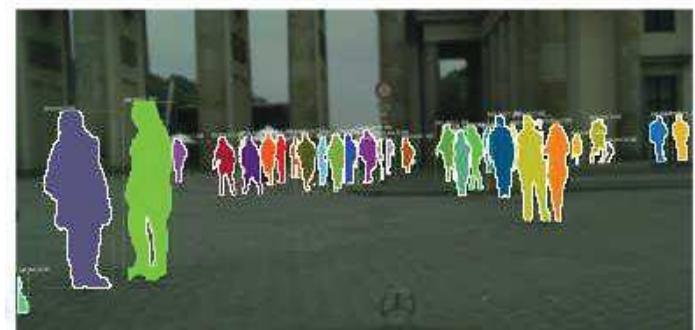
# Beyond bounding-boxes: getting detailed *contours* of objects of a given category



Object Detection



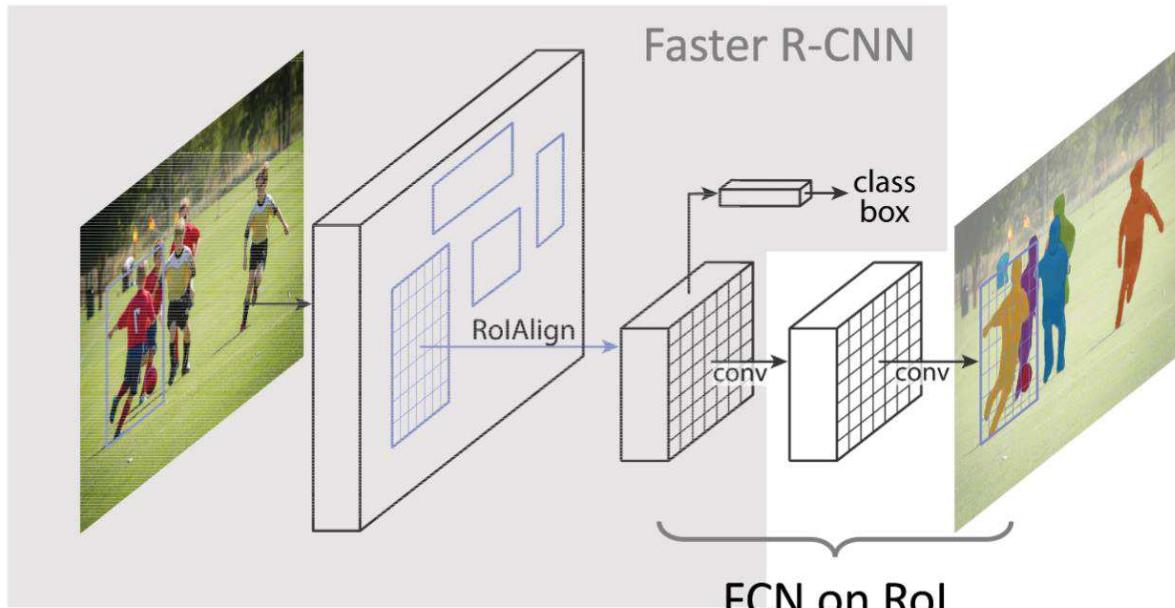
Instance Segmentation



Deep\_Learning for visual Scene Analysis (for IV), Pr. Fabien MOUTARDE, Center for Robotics, MINES ParisTech, PSL, Sept.2019 61

## Mask R-CNN principle

Mask R-CNN = Faster R-CNN with FCN on Rols



**Mask R-CNN architecture extract detailed contours and shape of objects instead of just bounding-boxes**

