

具身智能-02

刘华平

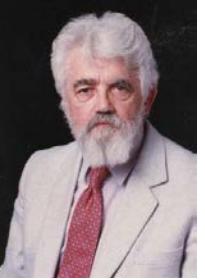
2024年2月26日

背景

➤ AI的历史



Marvin Lee Minsky
(1927-2016)



达特茅斯会议
(1956)



John Haugeland
(1945-2010)

GOF AI



Hopfield
(1933-)



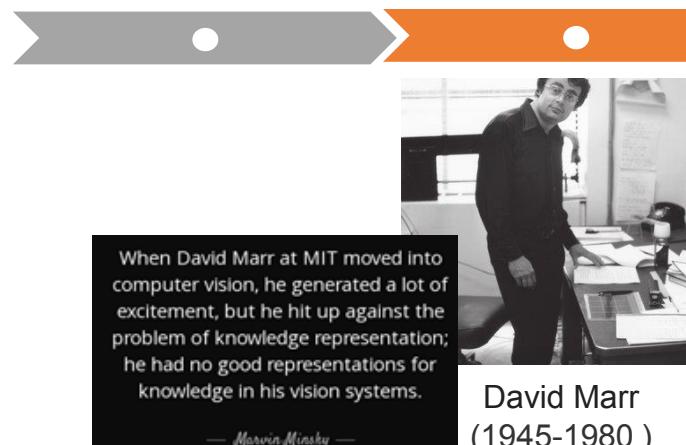
Rumelhart
(1942-2011)



Hilton
(1947-)

- 模式识别
- 数据挖掘
- 机器学习

神经网络->反向传播->深度学习



Pfeifer
(1947-) Moravec
(1948-) Brooks
(1954-) Cangelosi
(1967-)

- 机器人学
- 机构学
- 形态智能

• 智能是具身化和情境化的，智能需要一个身体

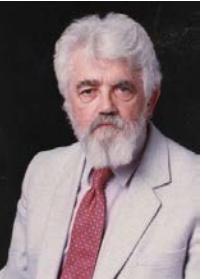
• 强化

背景

➤ 神经网络



Marvin Lee Minsky
(1927-2016)
John McCarthy
(1927-2011)
达特茅斯会议
(1956)



John Haugeland
(1945-2010)
GOF AI



Hopfield
(1933-)



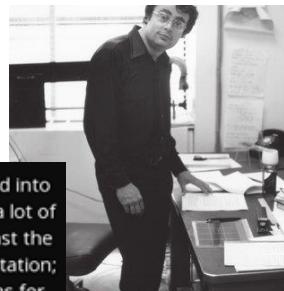
Rumelhart
(1942-2011)



Hilton
(1947-)

- 数据挖掘
- 机器学习
- 模式识别

神经网络->反向传播->深度学习



When David Marr at MIT moved into computer vision, he generated a lot of excitement, but he hit up against the problem of knowledge representation; he had no good representations for knowledge in his vision systems.
— Marvin Minsky —

David Marr
(1945-1980)

The Stochastic Neural Analog Reinforcement Calculator (SNARC) embodied a type of neural network that is designed to learn from experience and improve its performance through a process of trial and error (similar to reinforcement learning).

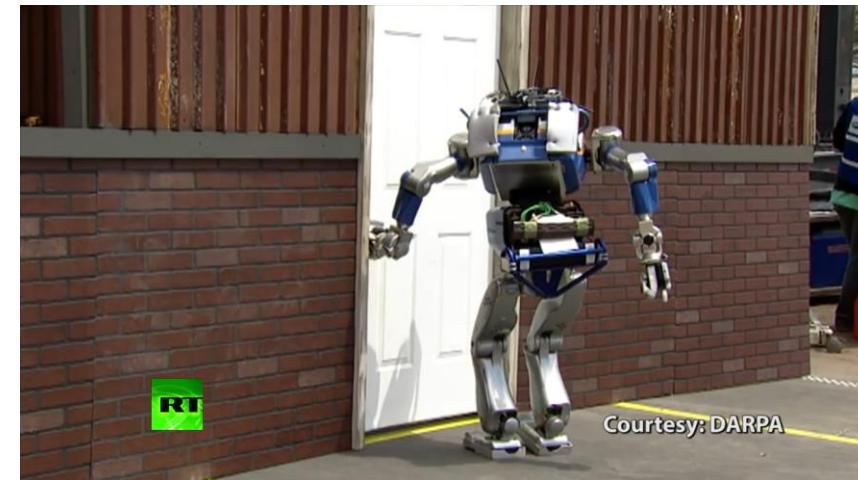
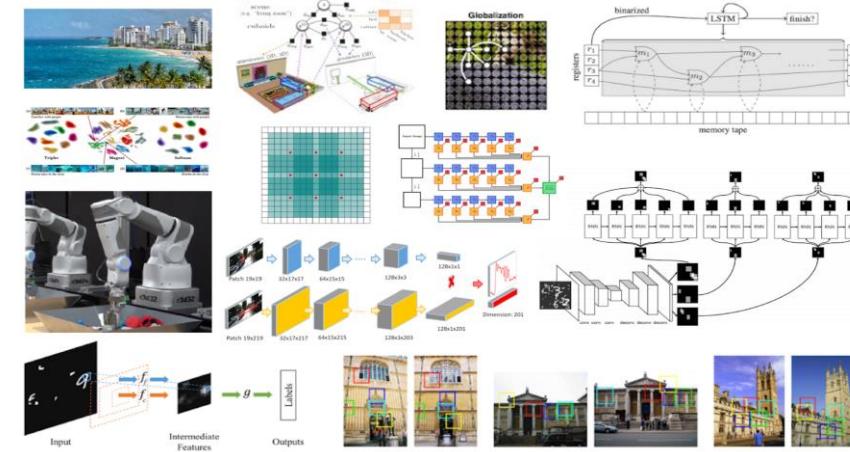
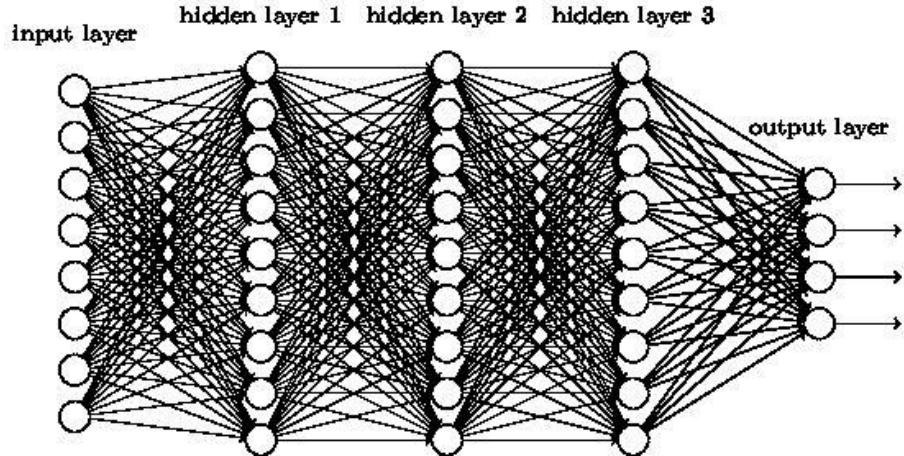


如果它现在不是，
那么总有一天会是



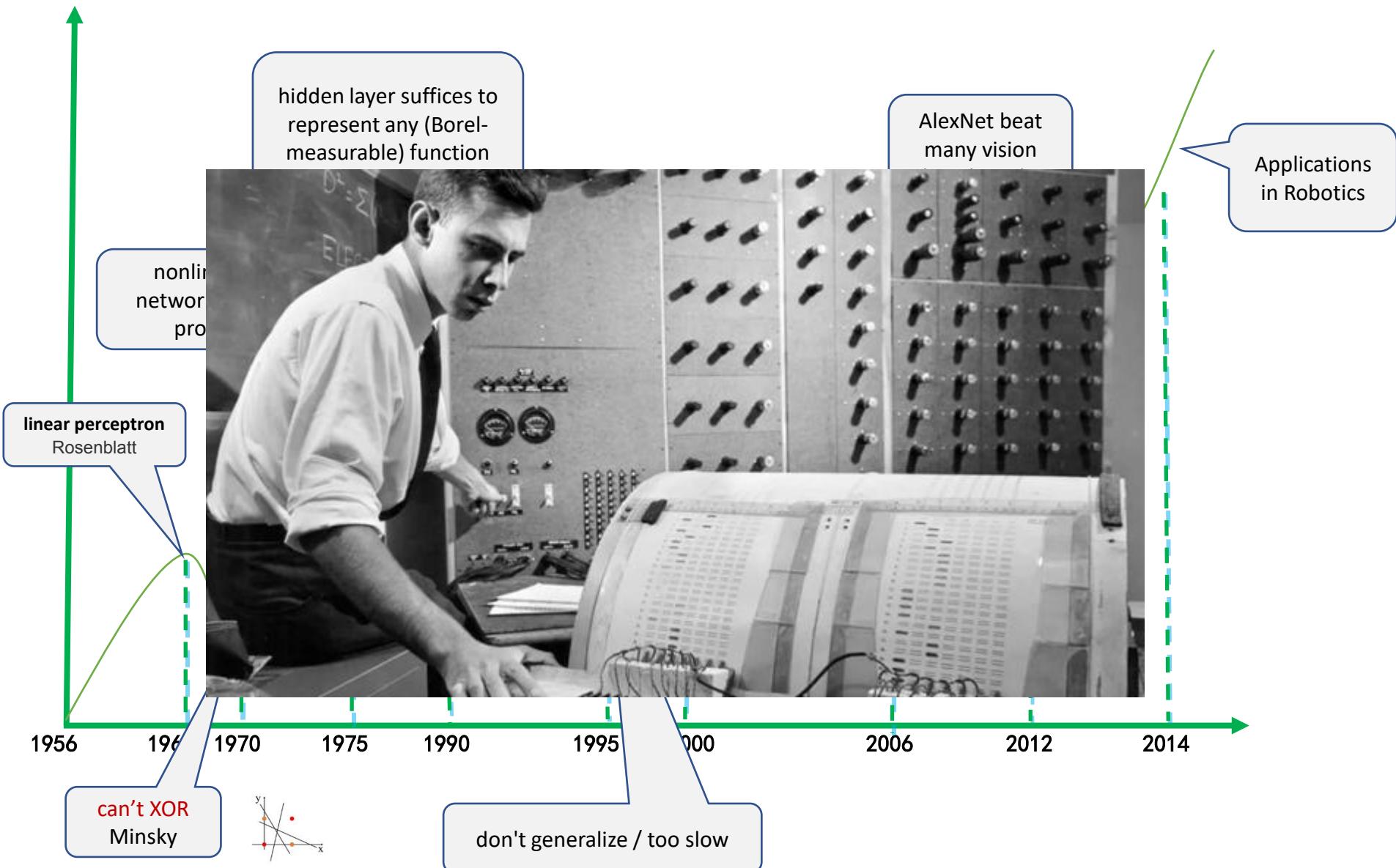
背景

➤ 神经网络

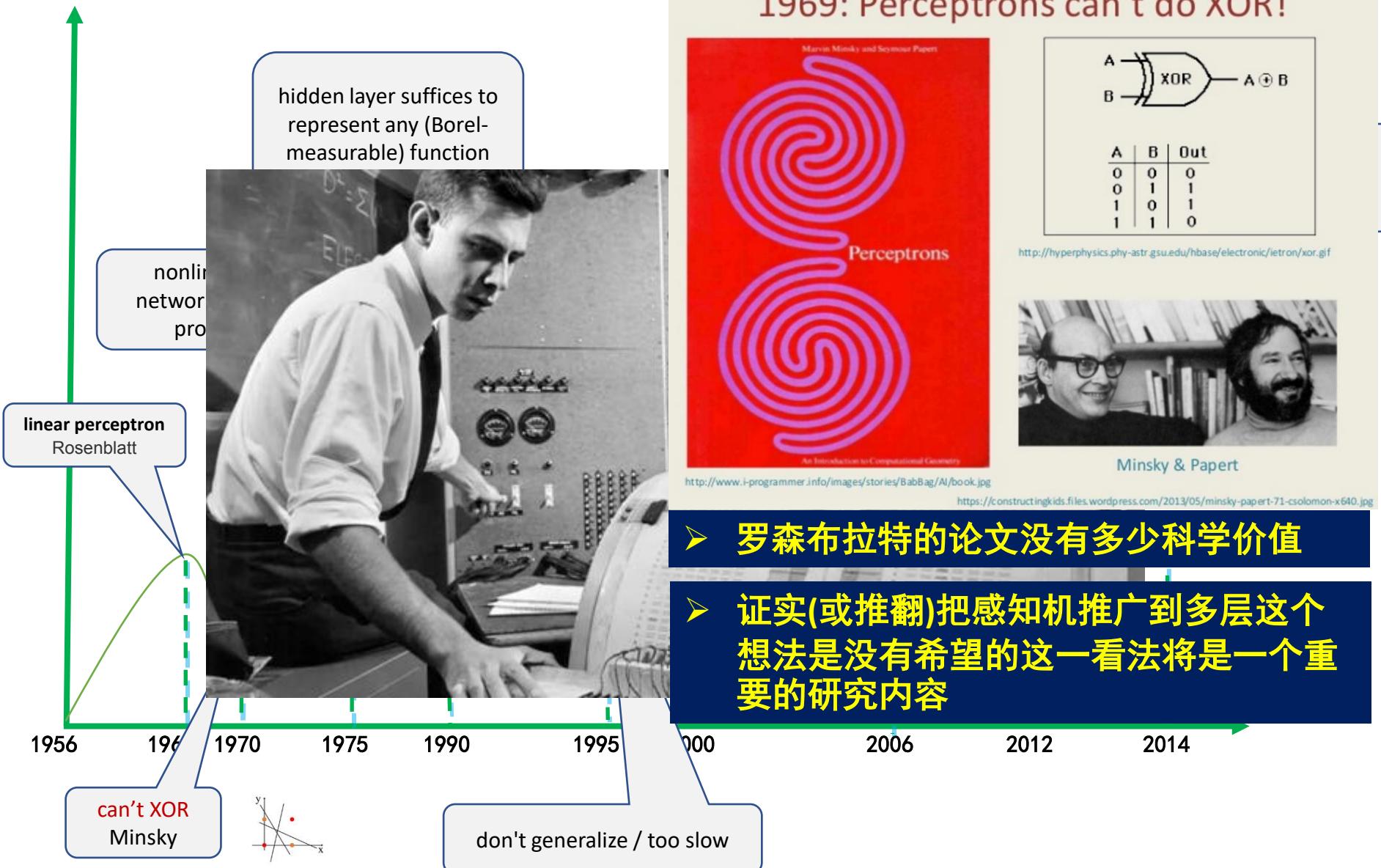


神经网络

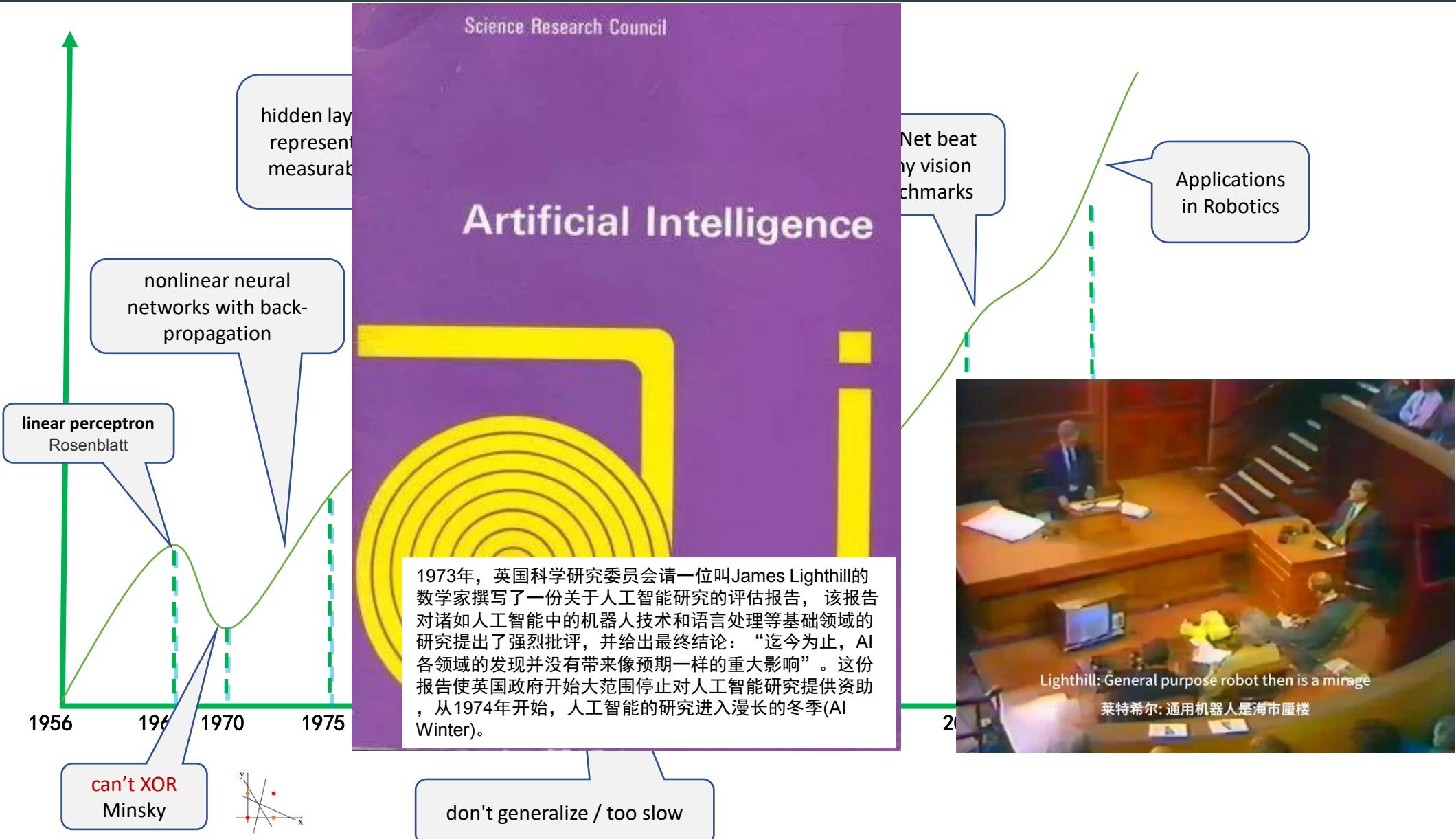
4



神经网络

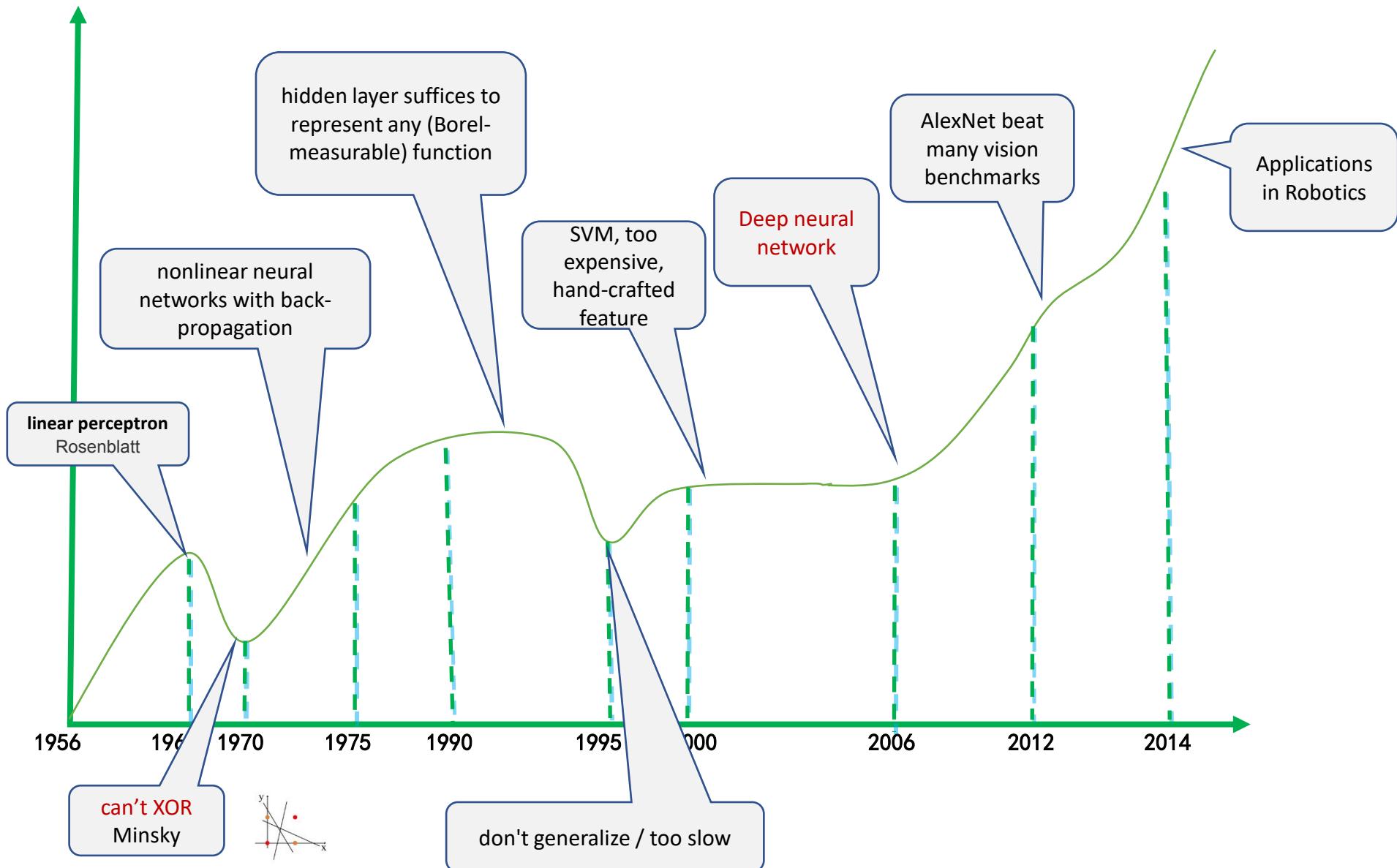


神经网络



神经网络

7



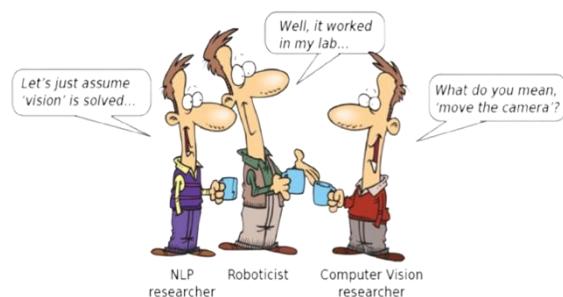
2010: Speech Recognition



2012: Computer Vision



2014: Machine Translation



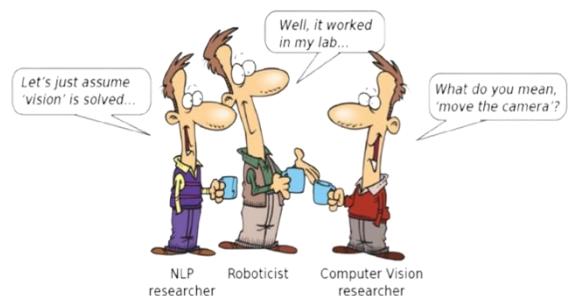
2010: Speech Recognition



2012: Computer Vision



2014: Machine Translation

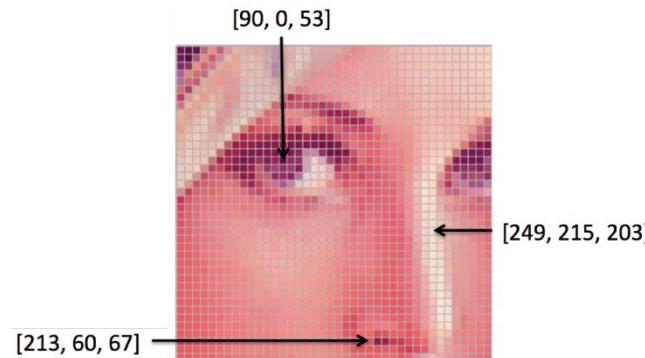


(1) 一般介绍

(2) 从多层感知机到卷积神经网络

(3) 应用

➤ 数字图像处理——图像滤波器



$$g[n, m] = \begin{cases} 255, & f[n, m] > 100 \\ 0, & \text{otherwise.} \end{cases}$$



23	25	30	35	30
25	30	35	37	40
45	40	37	43	45
38	40	43	42	46
35	40	42	45	47

				39

$$\text{Value} = (30 + 35 + 37 + 40 + 37 + 43 + 40 + 43 + 42) / 9 = 38.55 = 39$$

➤ 数字图像处理——图像滤波器



Original

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

= ?



Original

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 1 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

= ?



Original

$$\frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

= ?



original



+ a



sharpened

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix} = \begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$



original



-

smoothed (5x5)



=

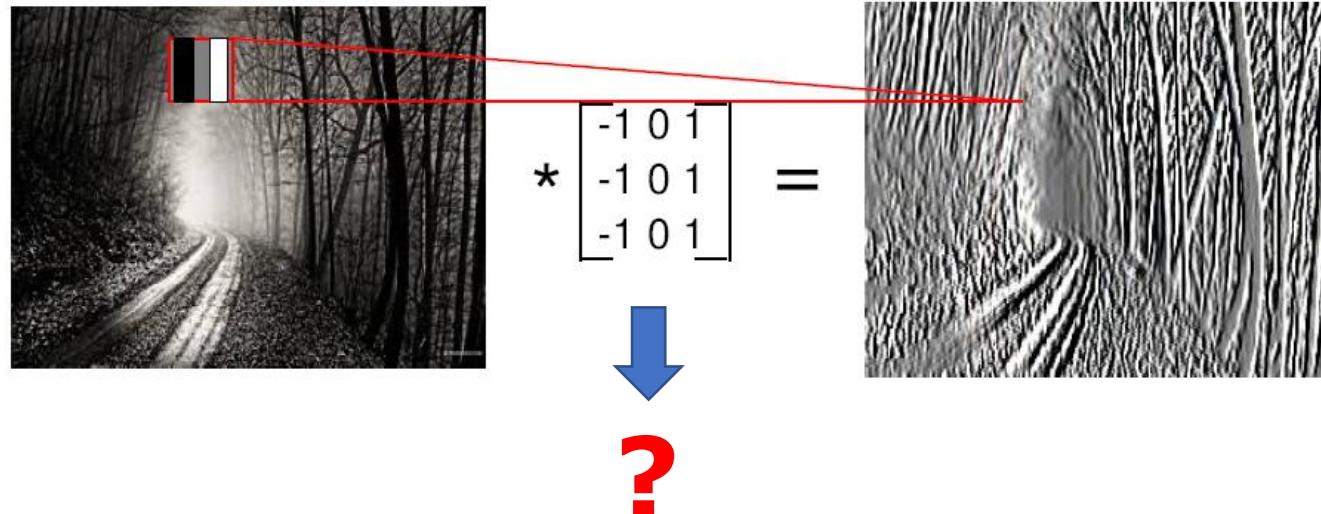
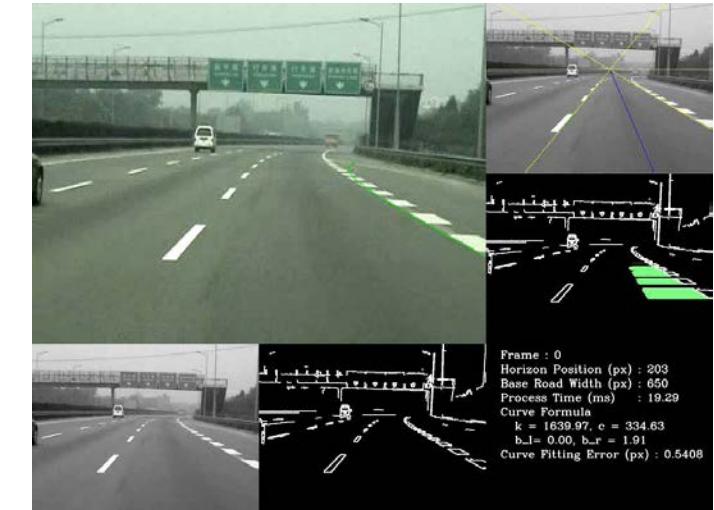
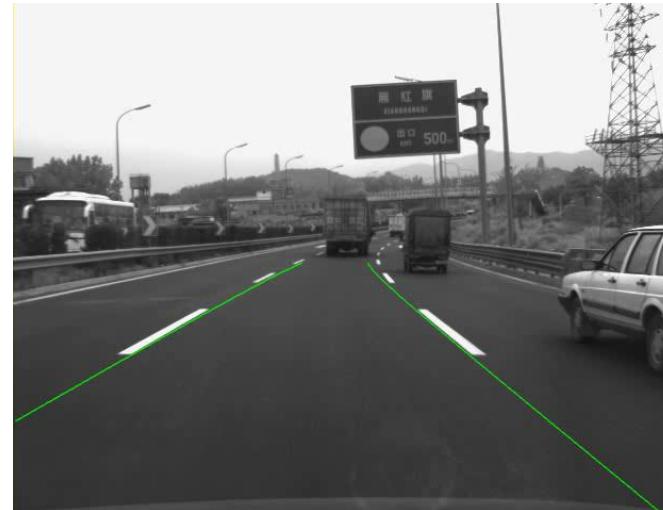
detail

$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix}$$

$$- \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

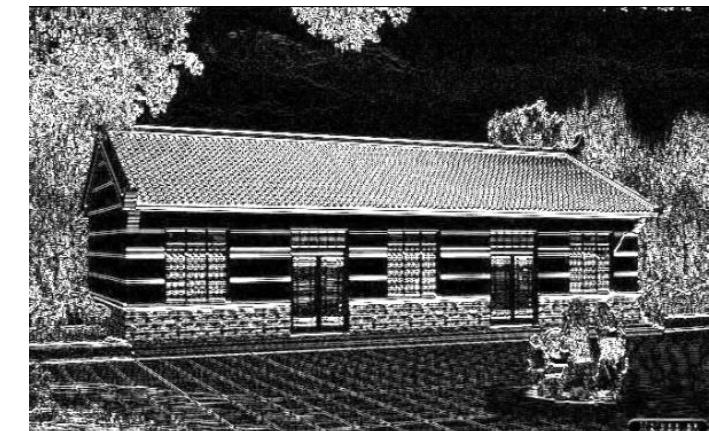
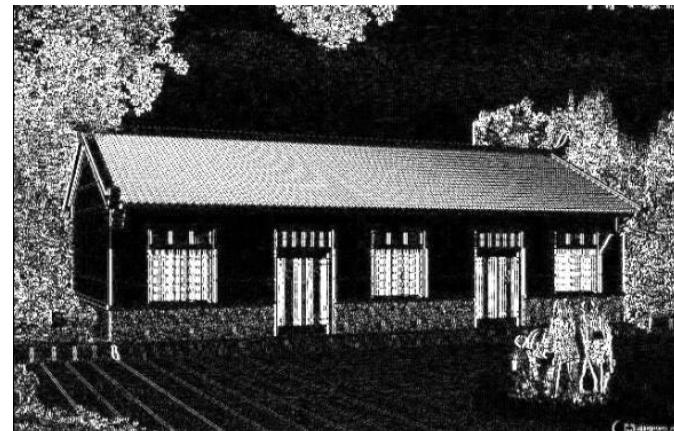


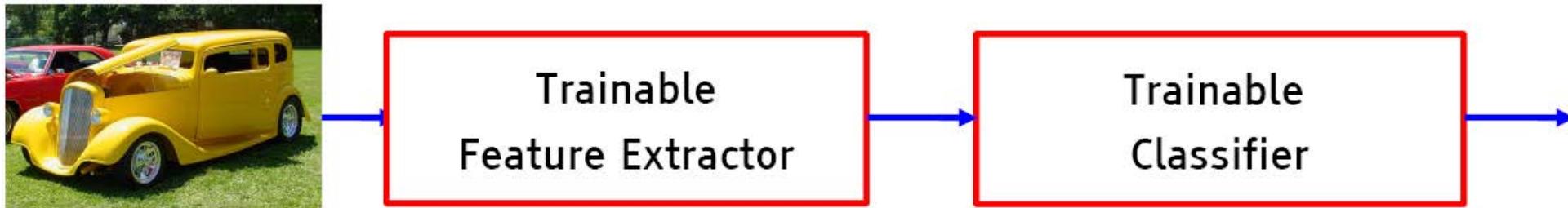
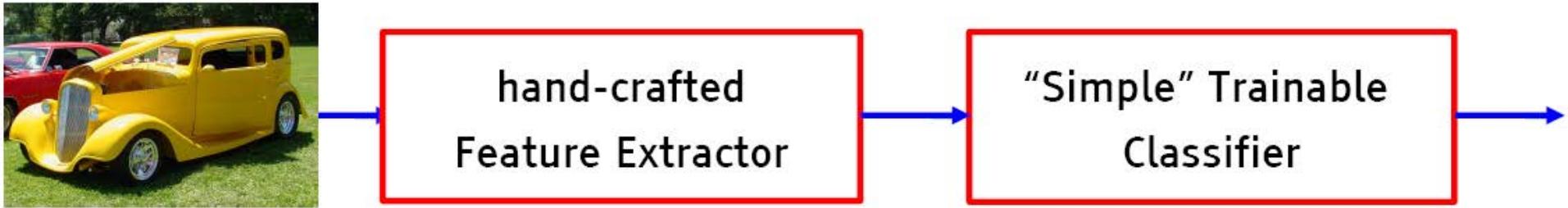
➤ 数字图像处理——图像滤波器



Sobel、Canny、Roberts、Prewitt,

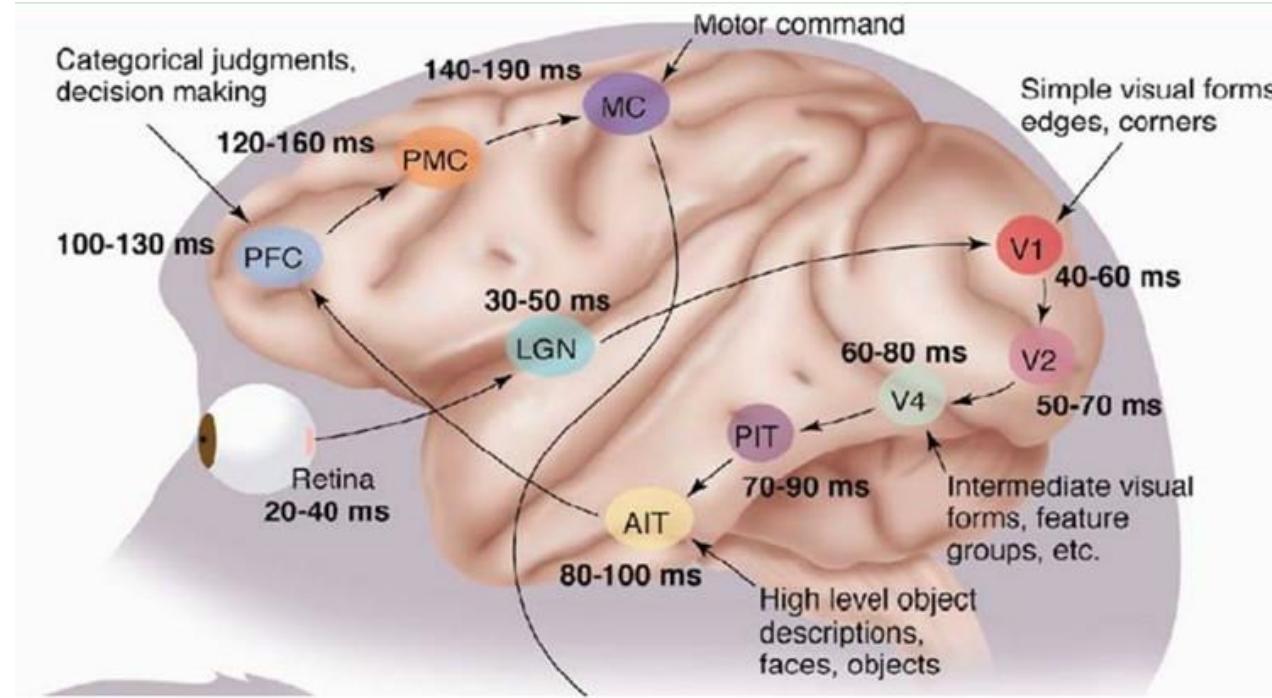
$$G_x = \begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix} * I \quad G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix} * I$$





Feature learning, representation learning

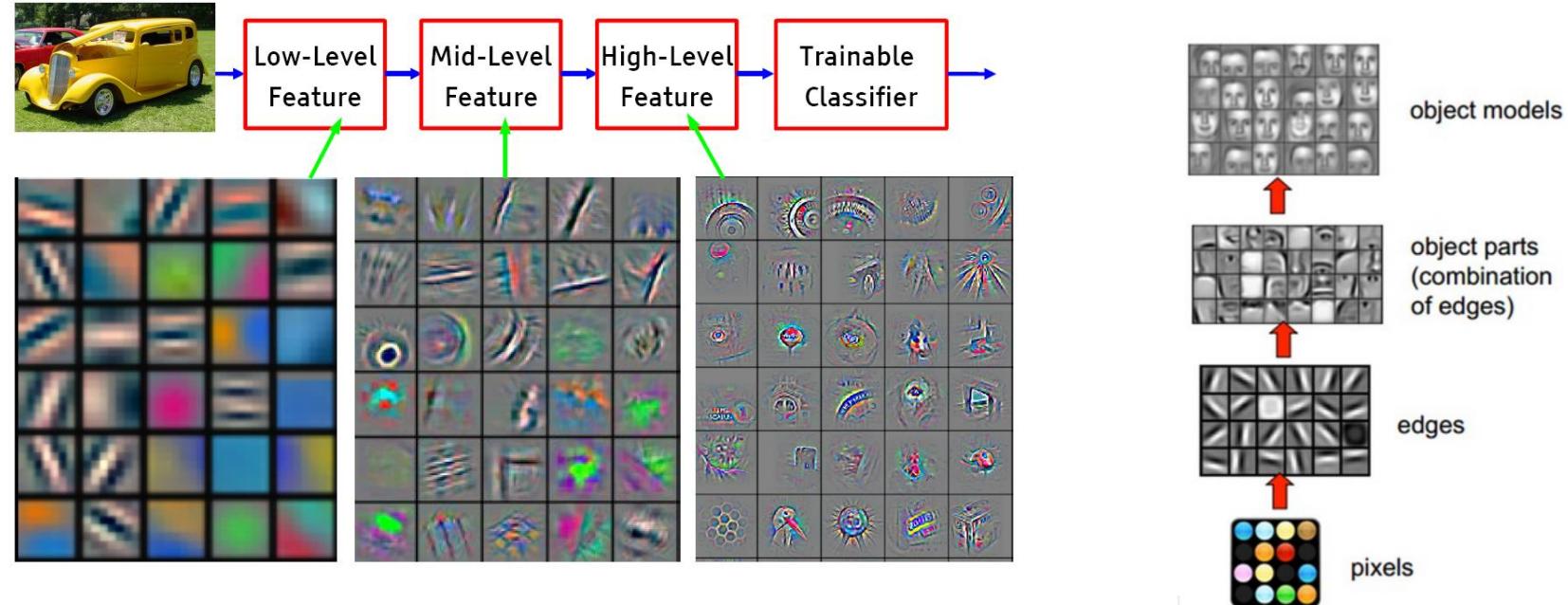
1981年，科学家们发现人的视觉系统的信息处理是分级的，人对物品的识别可能是一个不断迭代不断抽象的过程。



Retina - LGN - V1 - V2 - V4 - PIT - AIT

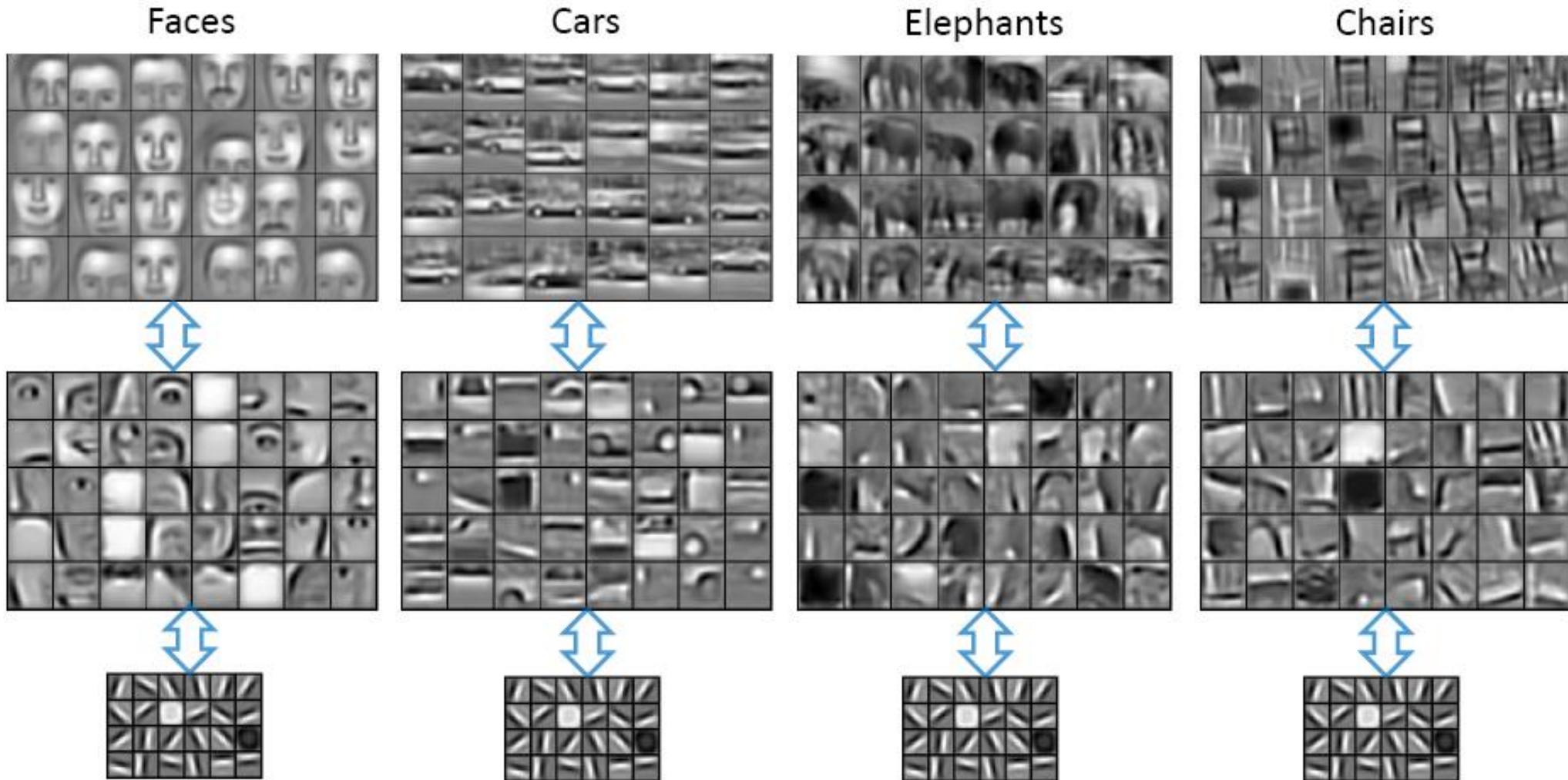
Lots of intermediate representations

Hierarchy of representations with increasing level of abstraction

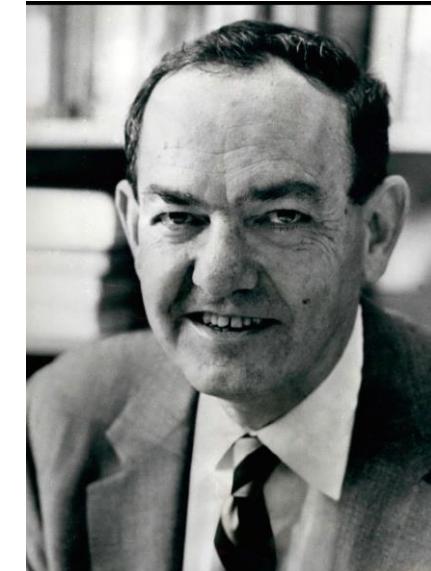


It's deep if it has more than one stage of non-linear feature transformation

Deep Learning addresses the problem of learning hierarchical representations with a single algorithm



“Solving a problem simply means **representing** it so as to make the solution transparent.”



Herbert A. Simon



ICLR 2020
8th International Conference
on Learning Representations

Addis Ababa, Ethiopia
April 26-30, 2020

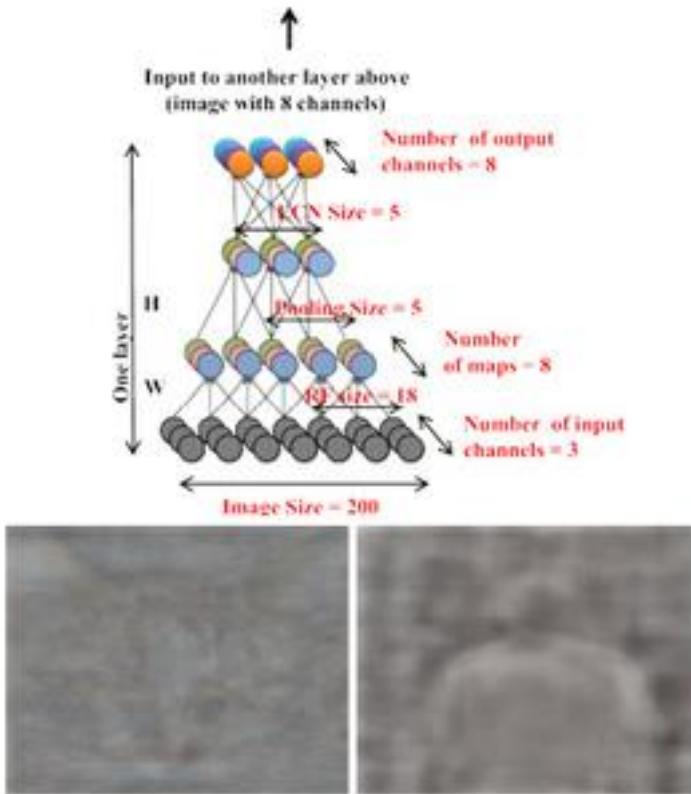


Figure 6. Visualization of the cat face neuron (left) and human body neuron (right).

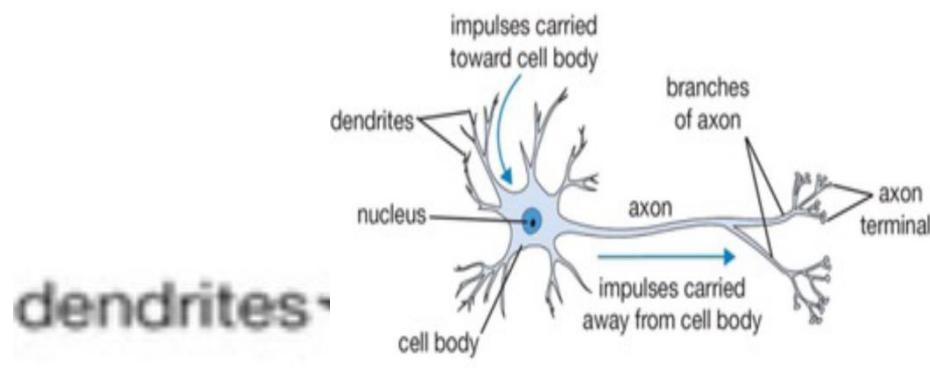
1000台机器，一共16000个Core训练了三天，参数的个数是1000000000个

(1) 一般介绍

(2) 从多层感知机到卷积神经网络

(3) 应用

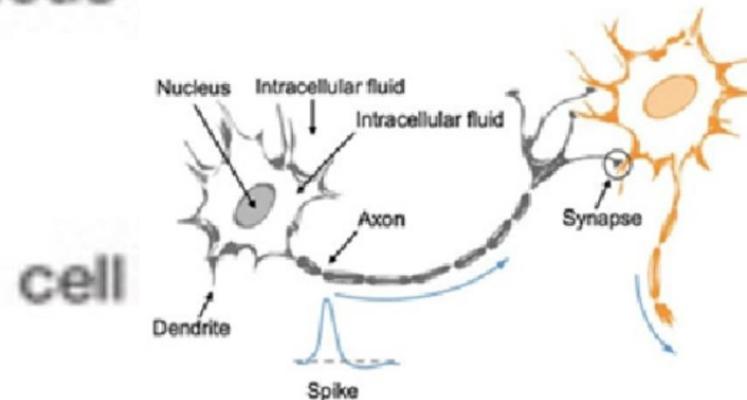
Biological Neuron versus Artificial Neural Network



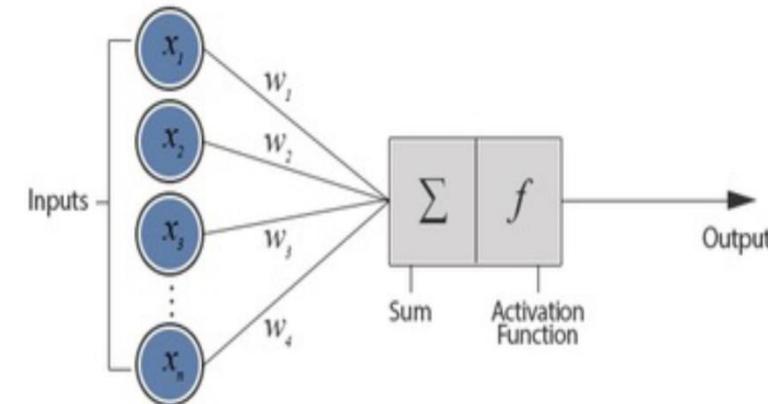
dendrites



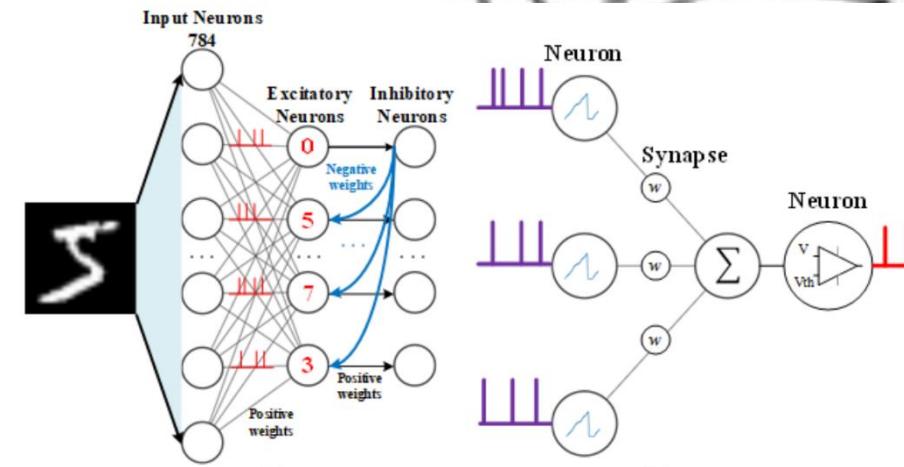
nucleus



cell

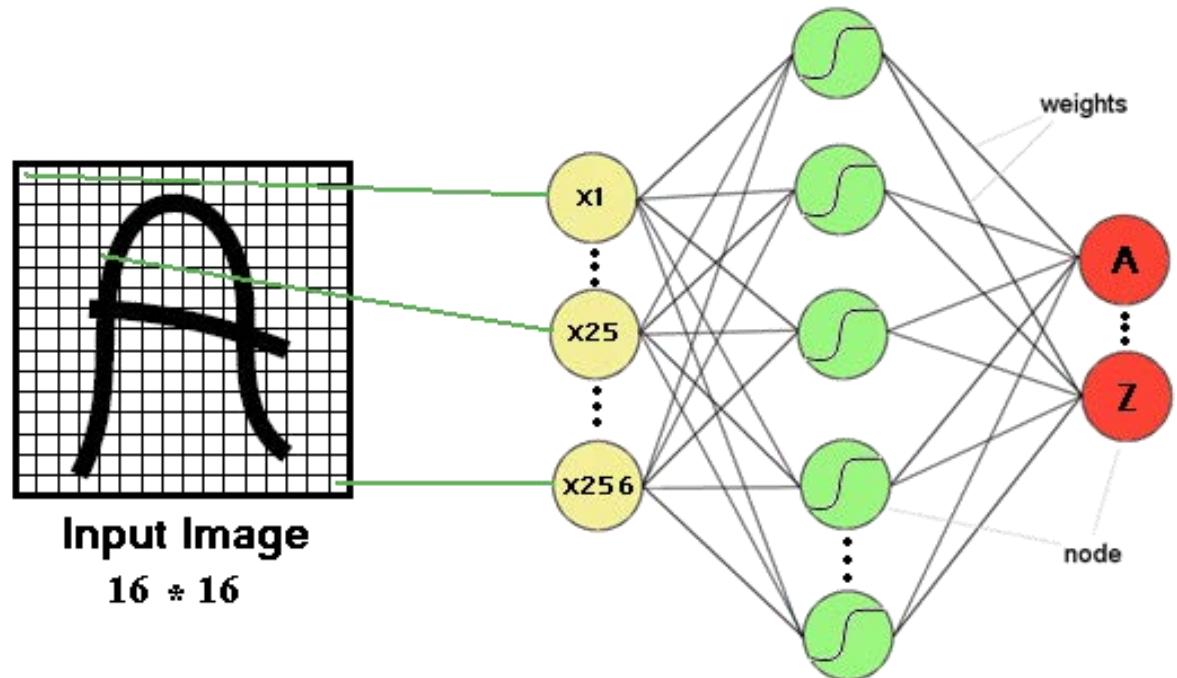


axon
minals

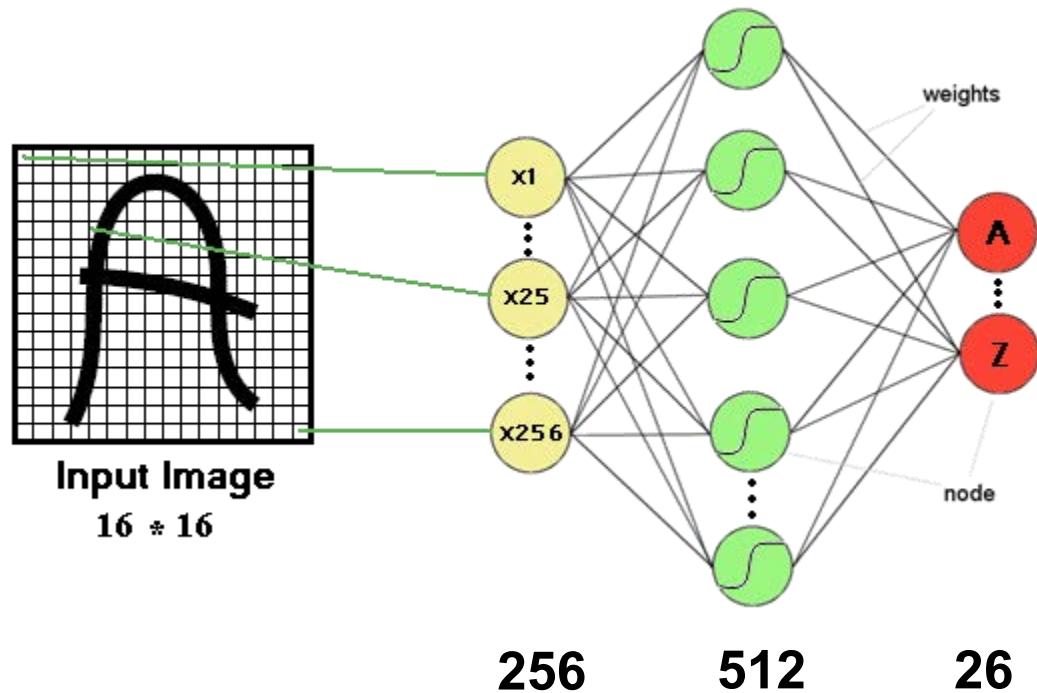


- 多个隐层
- Sigmoid激励函数

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



- $256 \times 512 + 512 \times 26$
- 大量的训练参数



Backpropagation (BP) algorithm

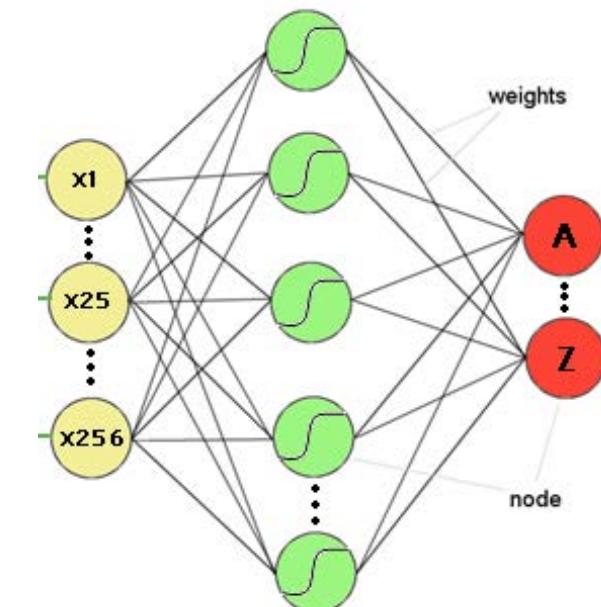
- 误差 $e_k(n) = T_k(n) - y_k(n)$

↓ ↓
target output

- 代价函数 $E = \frac{1}{2} \sum e_k^2(n)$

- 权值调整 $\Delta w_{kj}(l) = -\frac{\partial E}{\partial w_{kj}(l)}$

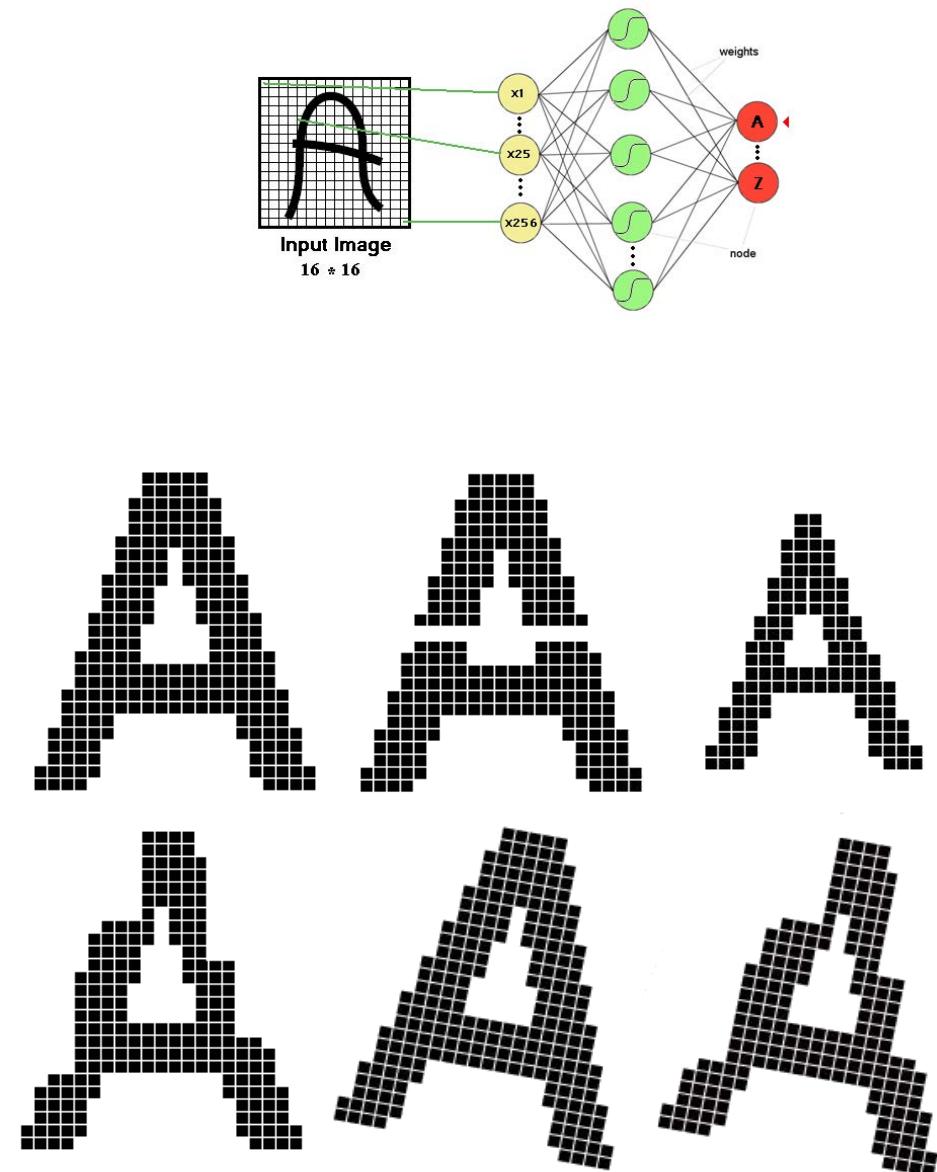
- BP algorithm was proposed by Paul Werbos in 1974 and reinvented in 1985 by Rumelhart et al.



- 对*shifting, scaling, rotation*等非常敏感
- 过拟合!

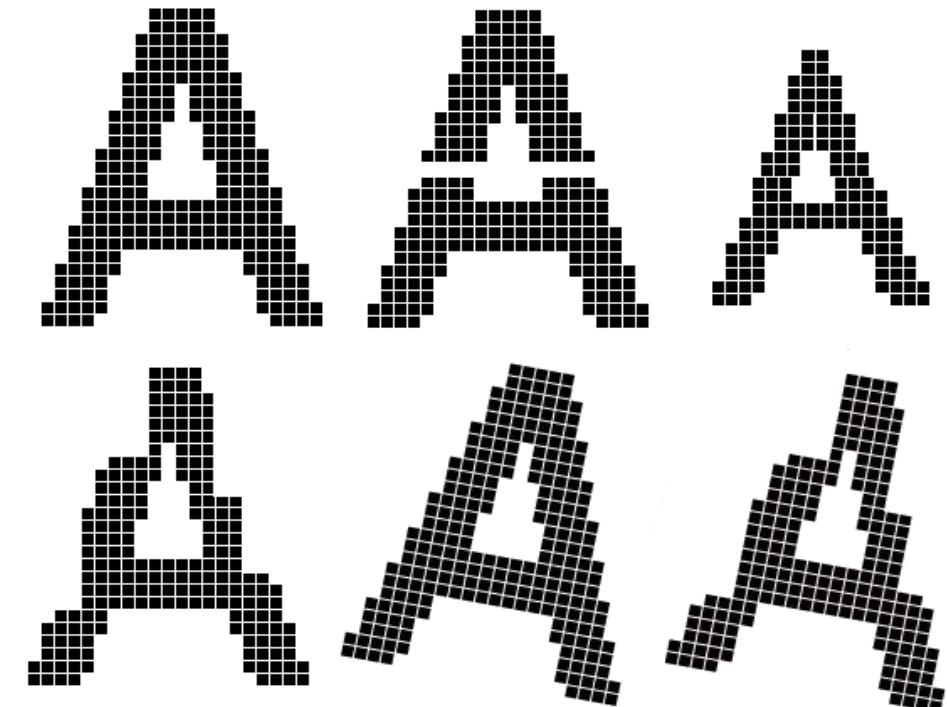


- ✓ 训练时间长
- ✓ 网络规模大
- ✓ 调整参数多



“好” 的特征

- 充分刻画结构
- 对*shifting, scaling, rotation*等不敏感
- 训练参数少



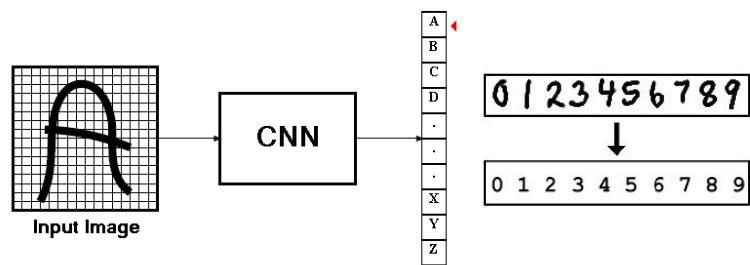
Convolutional Neural Network



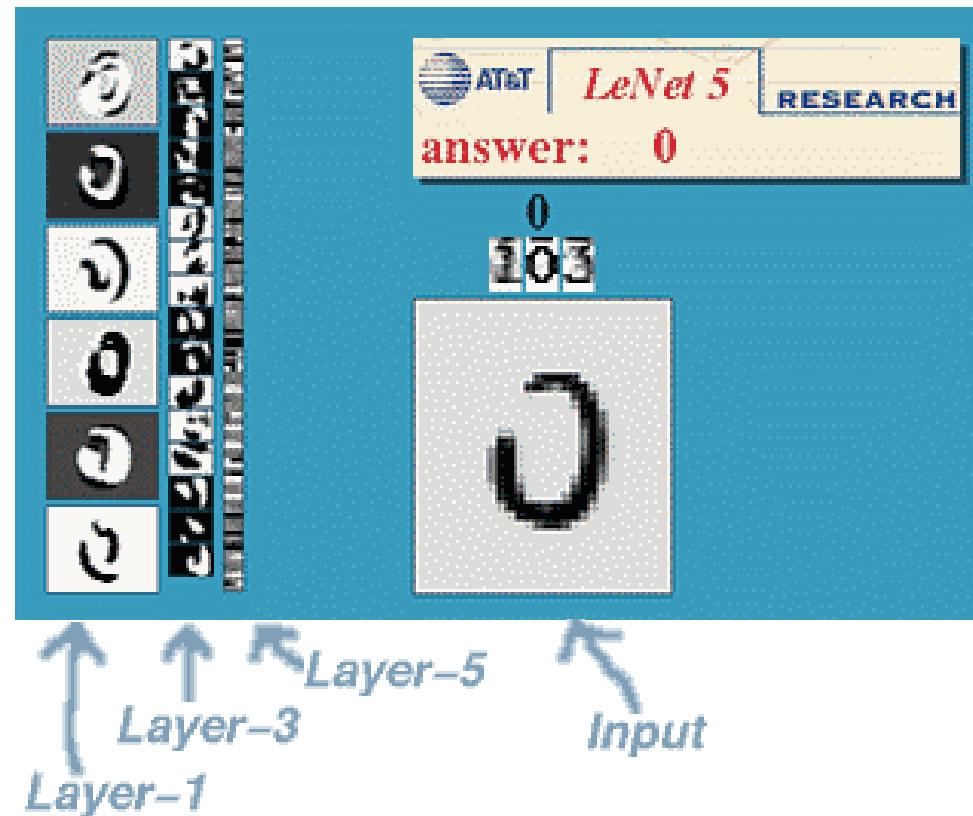
Yann LeCun



Yoshua Bengio

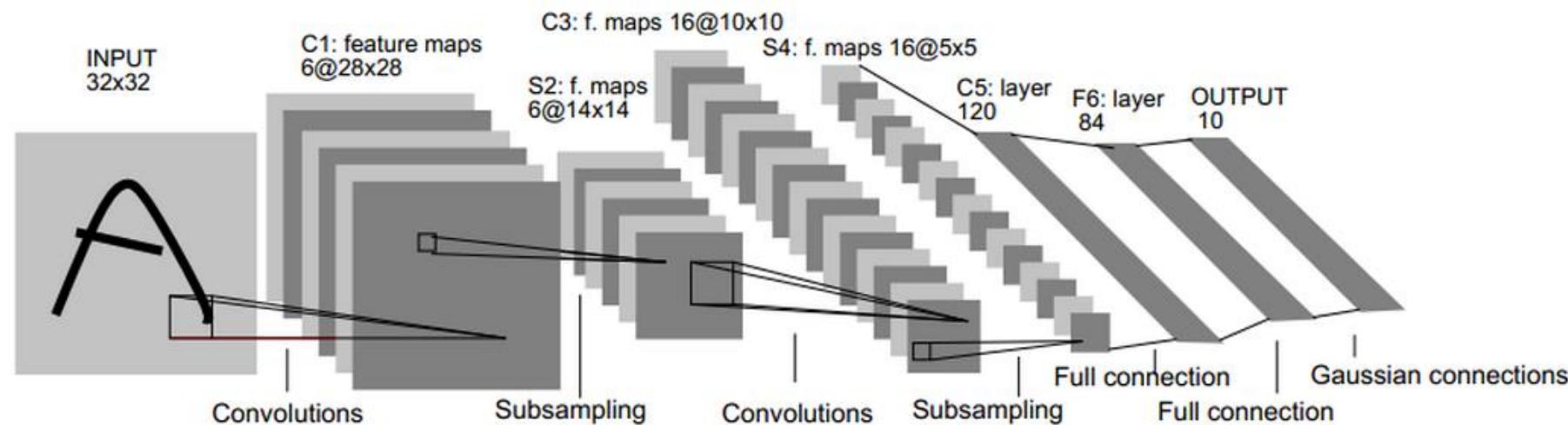


In 1995, **Yann LeCun** and **Yoshua Bengio** introduced the concept of convolutional neural networks.



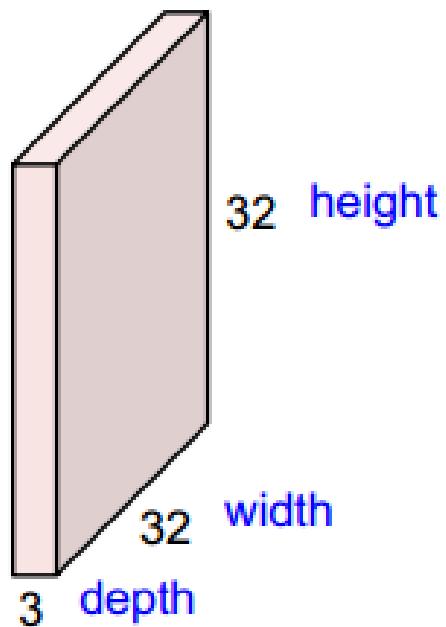
- C1,C3: Convolutional layer.
- 5×5 Convolution matrix.
- S2 , S4 : Subsampling layer.
- Subsampling by factor 2.
- F6 : Fully connected layer.

LeNet-5



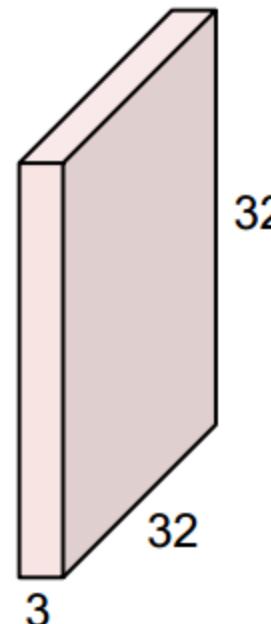
Convolution Layer

32x32x3 image



Convolution Layer

32x32x3 image



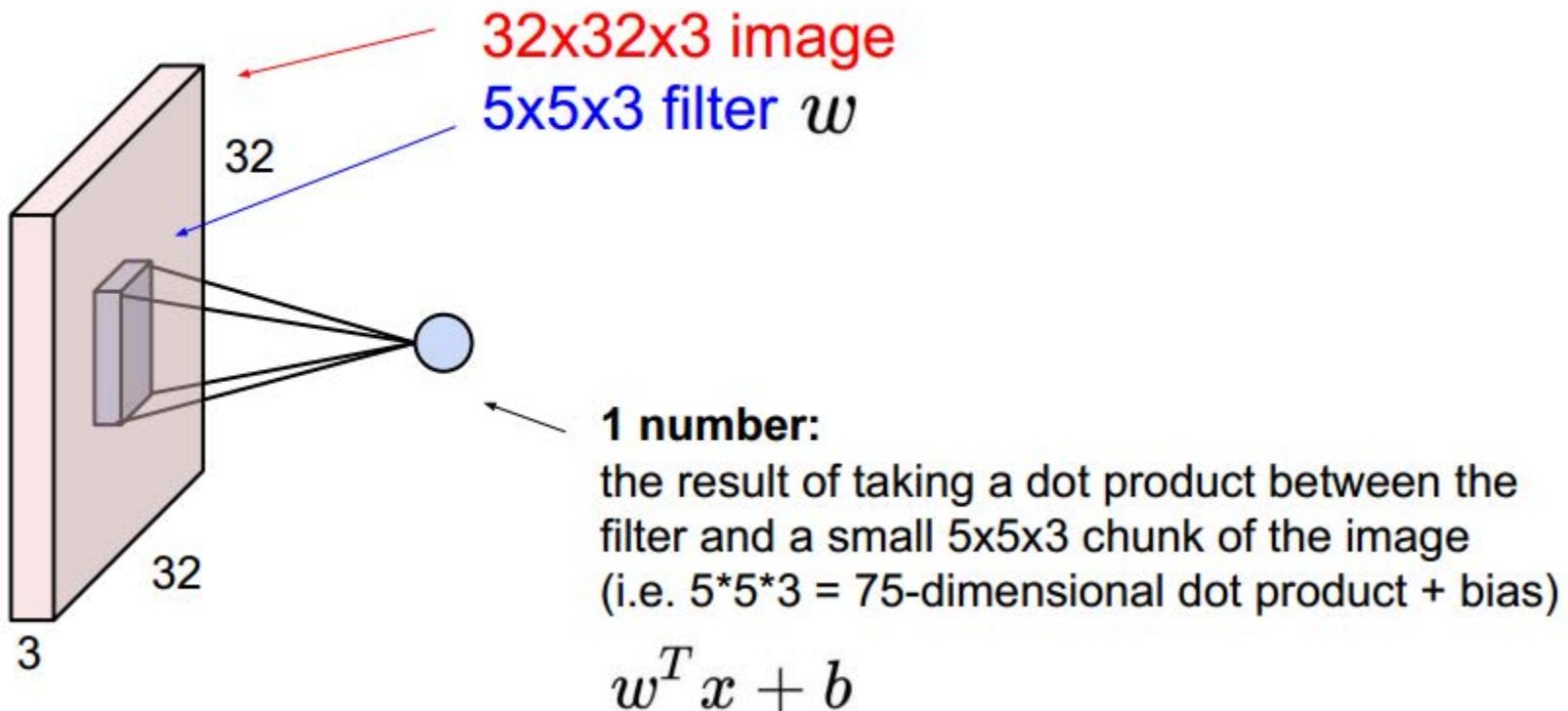
5x5x3 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

$$\ast \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} =$$

Convolution Layer

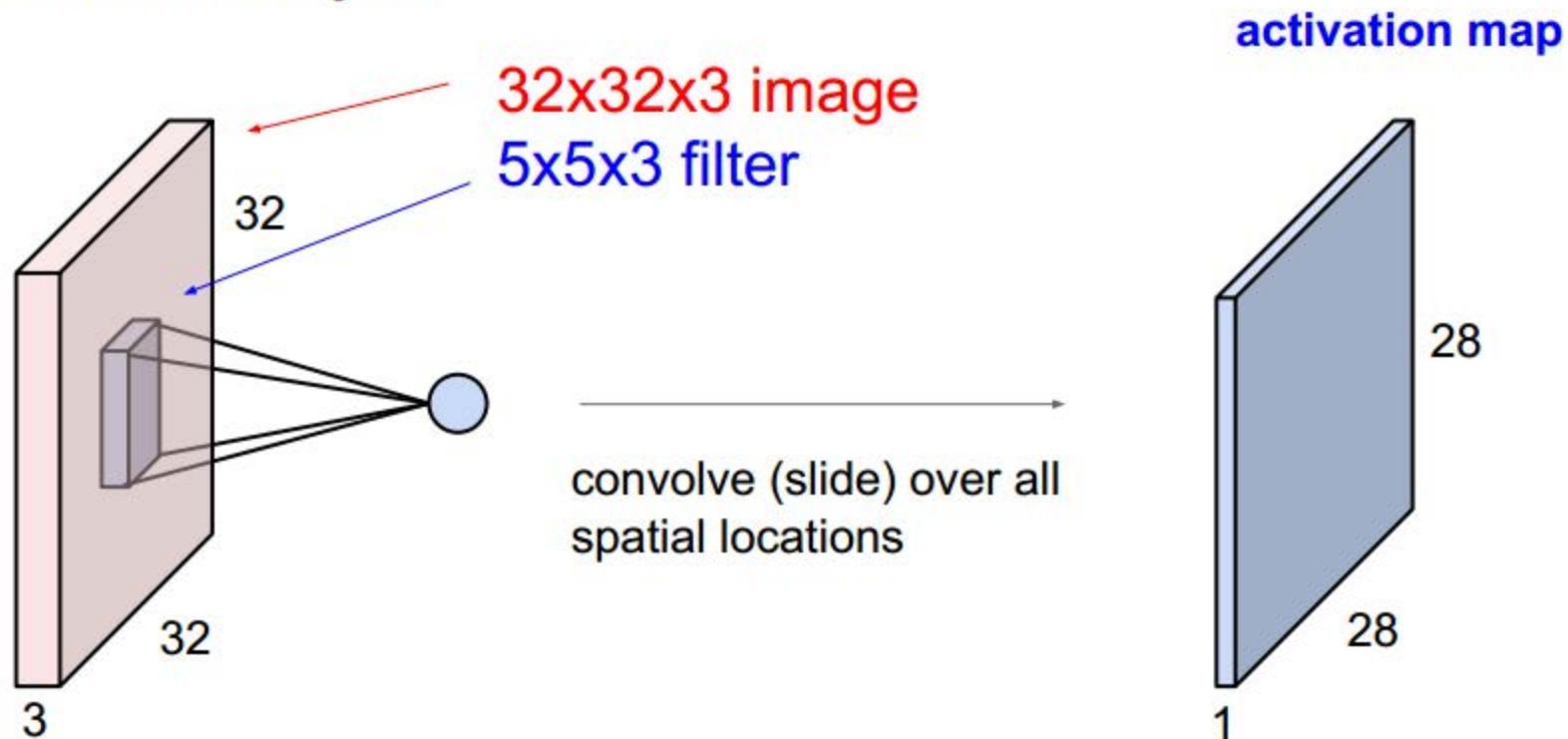


An illustration showing a convolution step. On the left is a grayscale image of a forest path. A small 3x3 input patch is highlighted with a red border. To its right is a 3x3 kernel matrix:

$$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

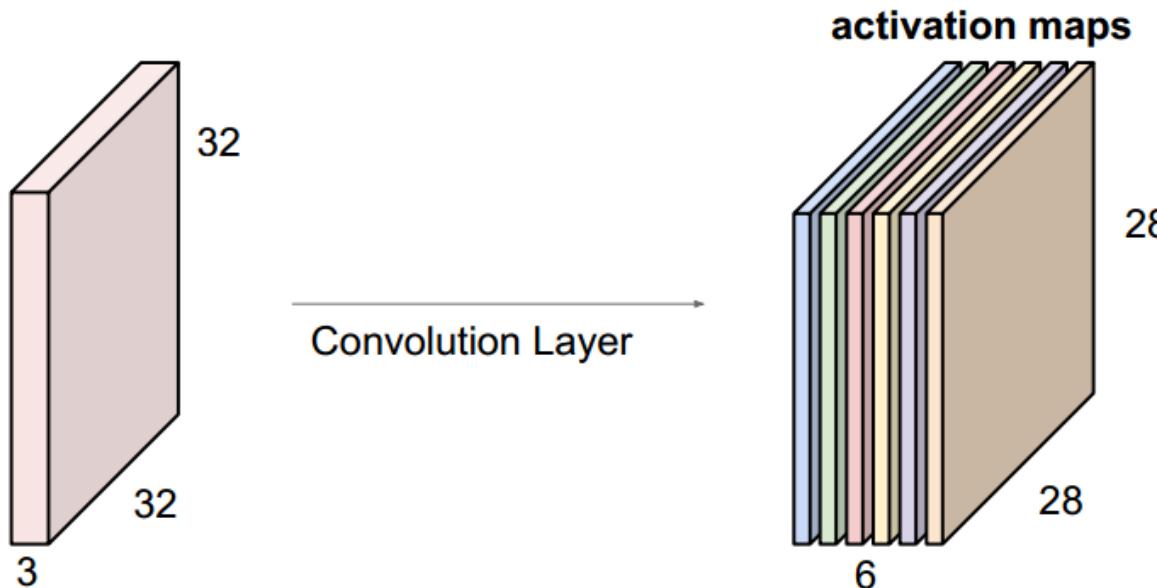
Below the kernel is an equals sign, followed by a blurred grayscale image showing the result of the convolution step.

Convolution Layer

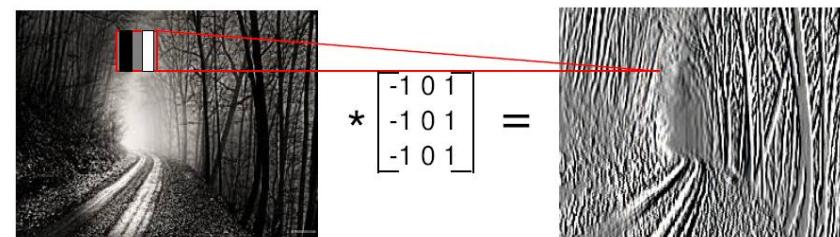


Convolution layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

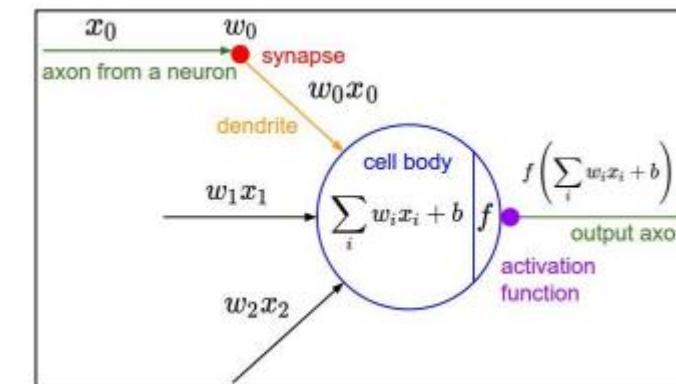
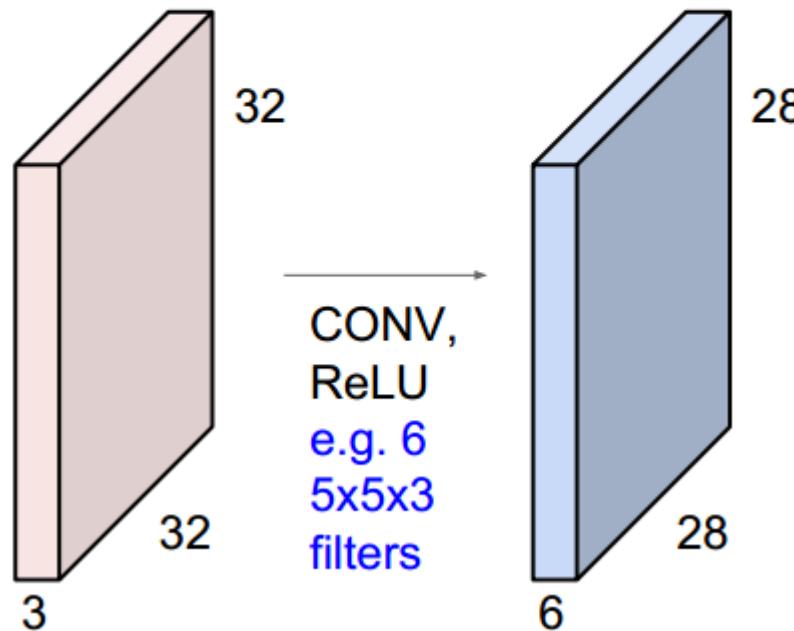


We stack these up to get a “new image” of size 28x28x6!

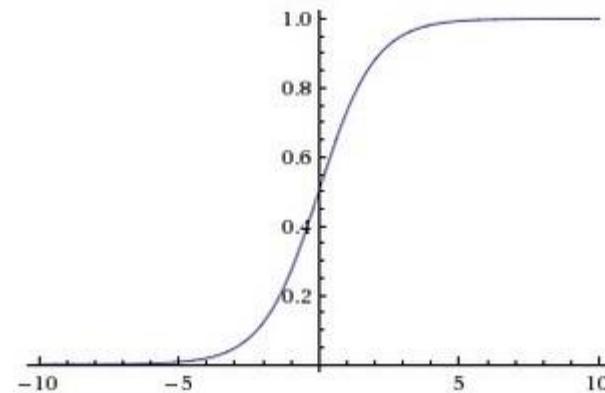


Convolution layer

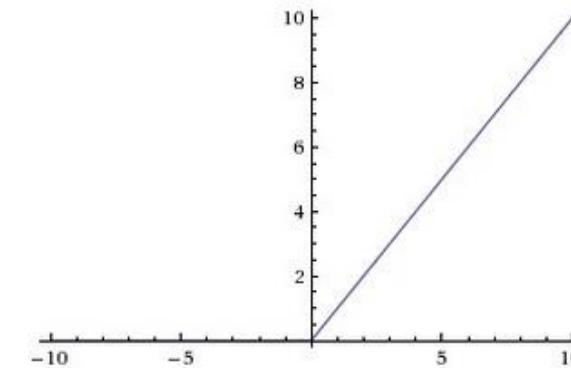
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



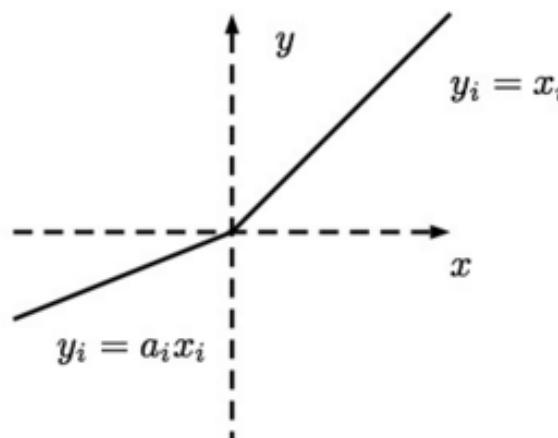
Activation function



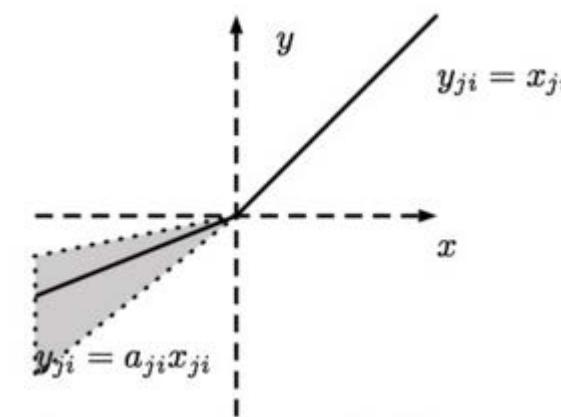
sigmoid



ReLU

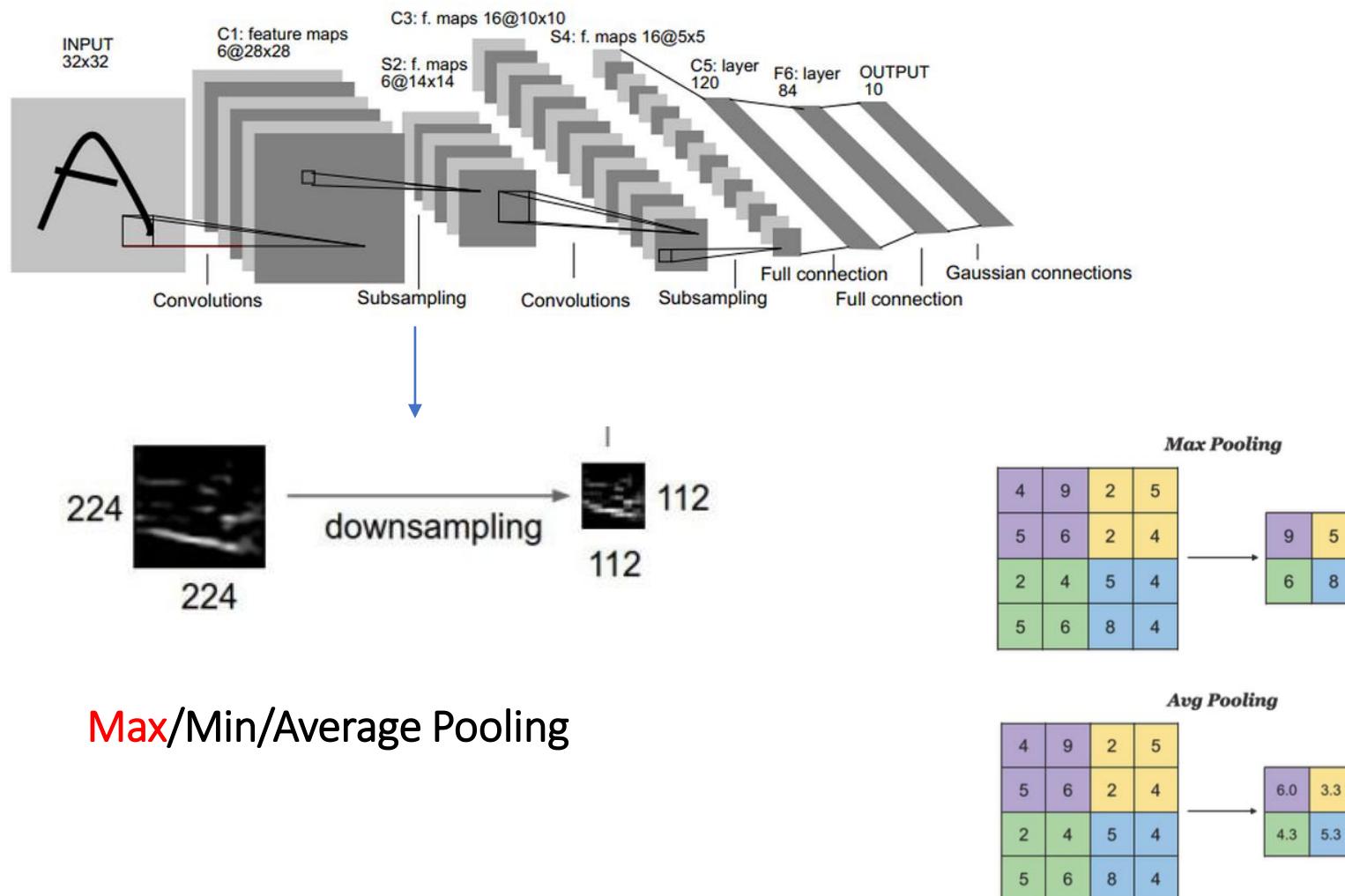


PReLU

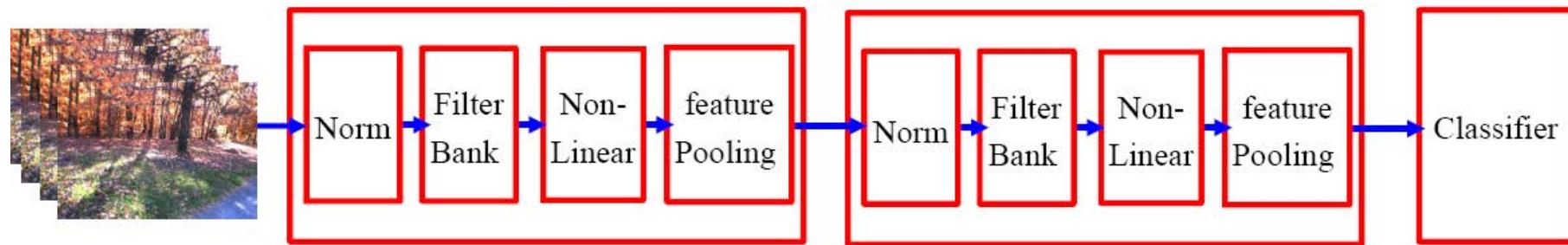


Random PReLU

Subsampling layer



➤ 总结：一般结构



Normalization: variations on whitening

Subtractive: average removal, high pass filtering

Divisive: local contrast normalization, variance normalization

Filter Bank: dimension expansion, projection on overcomplete basis

Non-Linearity: sparsification, saturation, lateral inhibition....

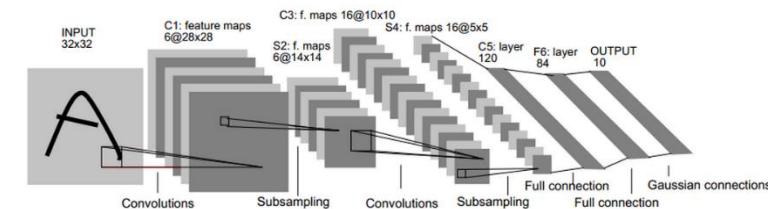
Rectification (ReLU), Component-wise shrinkage, tanh, winner-takes-all

Pooling: aggregation over space or feature type

➤ 总结：一般结构

传统的图像识别方法，需要预先使用例如SIFT/HOG等方法对图像进行特征提取，再使用SVM或稀疏编码等方法进行分类，面对复杂的问题往往难以提取合适的特征。

CNN算法能利用网络中间某一层的输出当做是数据的另一种表达，可以经过网络学习到适合于分类的特征。



卷积神经网络的优点

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

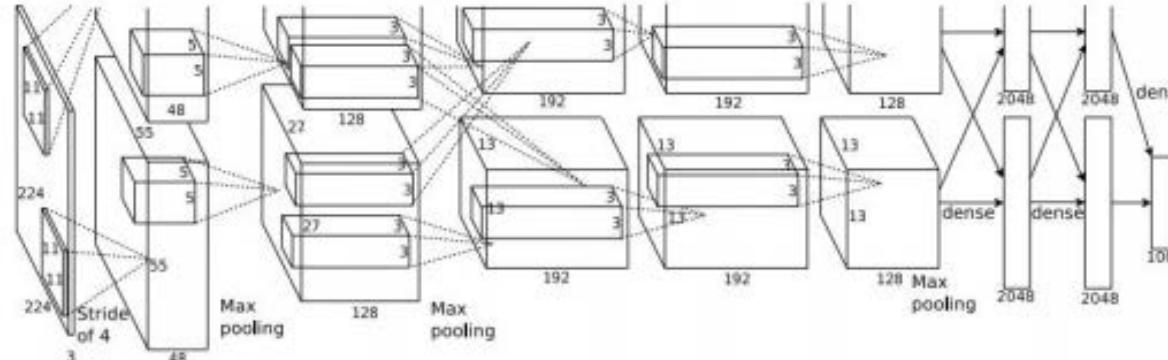
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

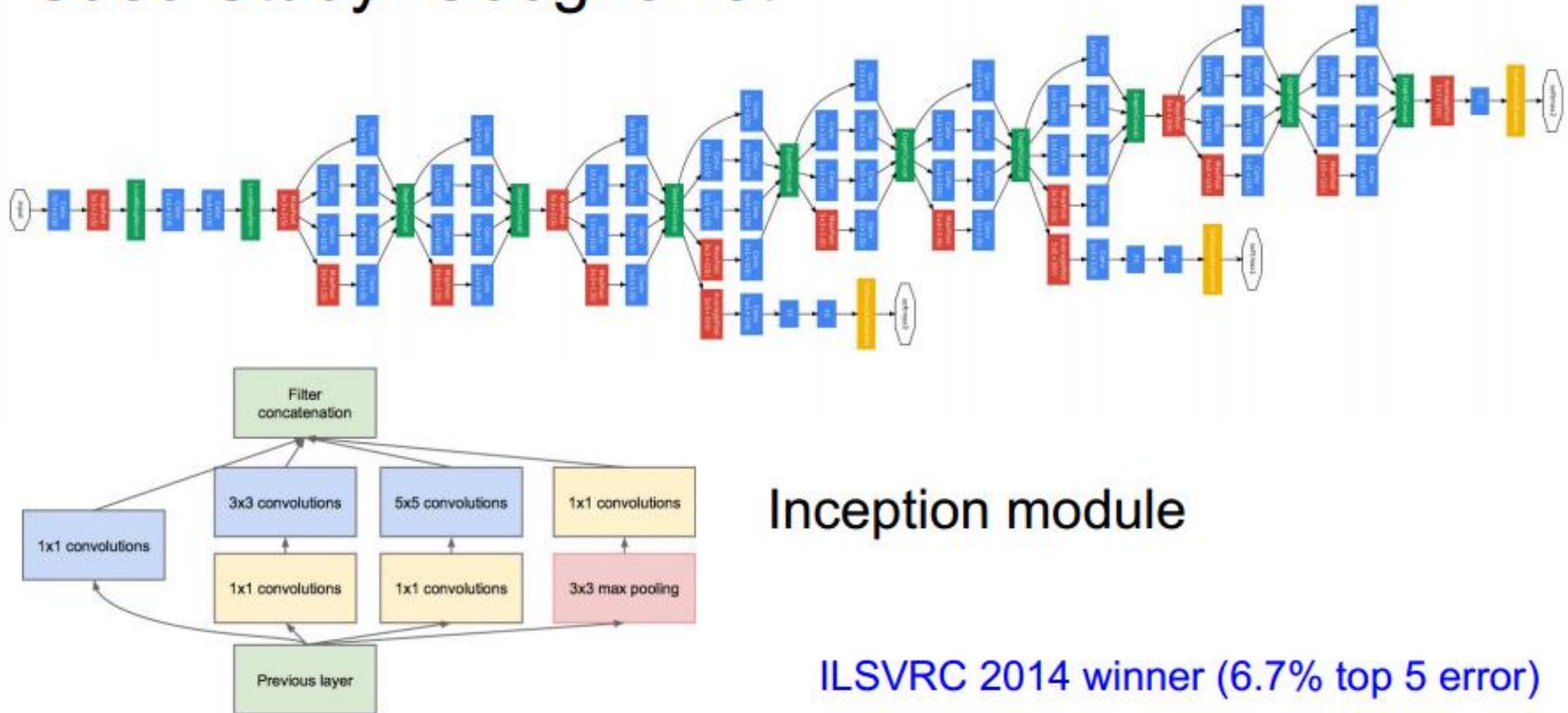
ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256	conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256
conv3-256	conv3-256	conv3-256	conv3-256 conv1-256	conv3-256 conv1-256	conv3-256 conv3-256
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	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
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

Case Study: GoogLeNet

[Szegedy et al., 2014]



Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)

Microsoft Research

MSRA @ ILSVRC & COCO 2015 Competitions

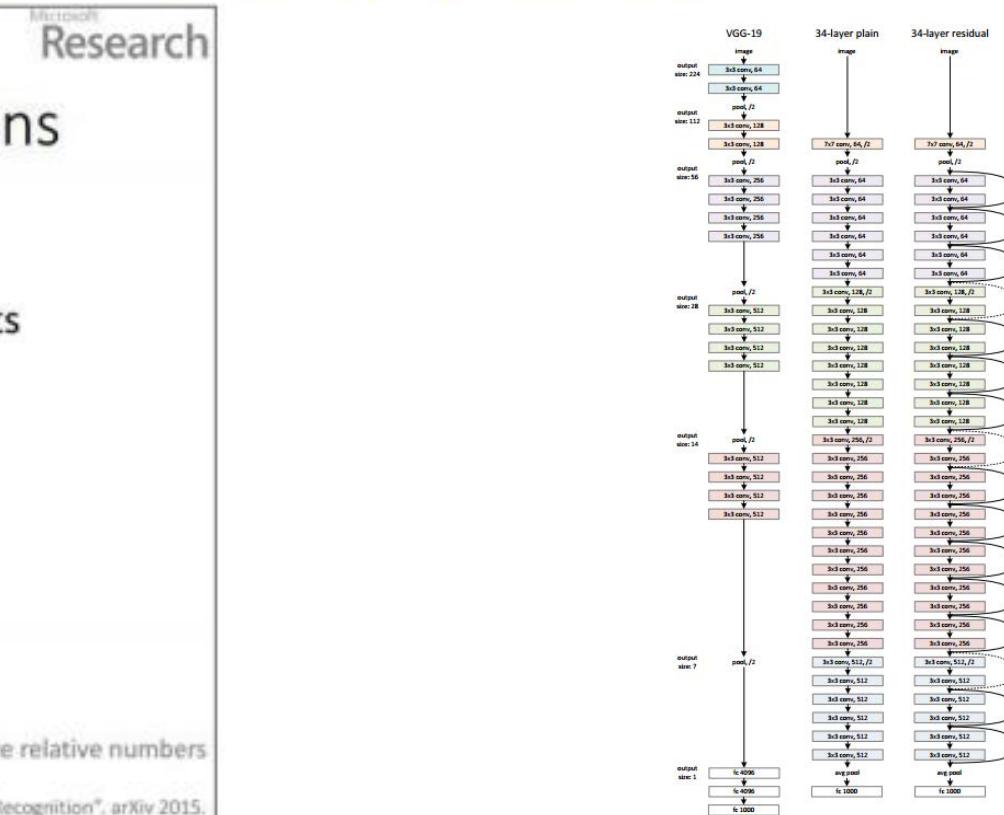
- **1st places** in all five main tracks
 - ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer** nets
 - ImageNet Detection: **16%** better than 2nd
 - ImageNet Localization: **27%** better than 2nd
 - COCO Detection: **11%** better than 2nd
 - COCO Segmentation: **12%** better than 2nd

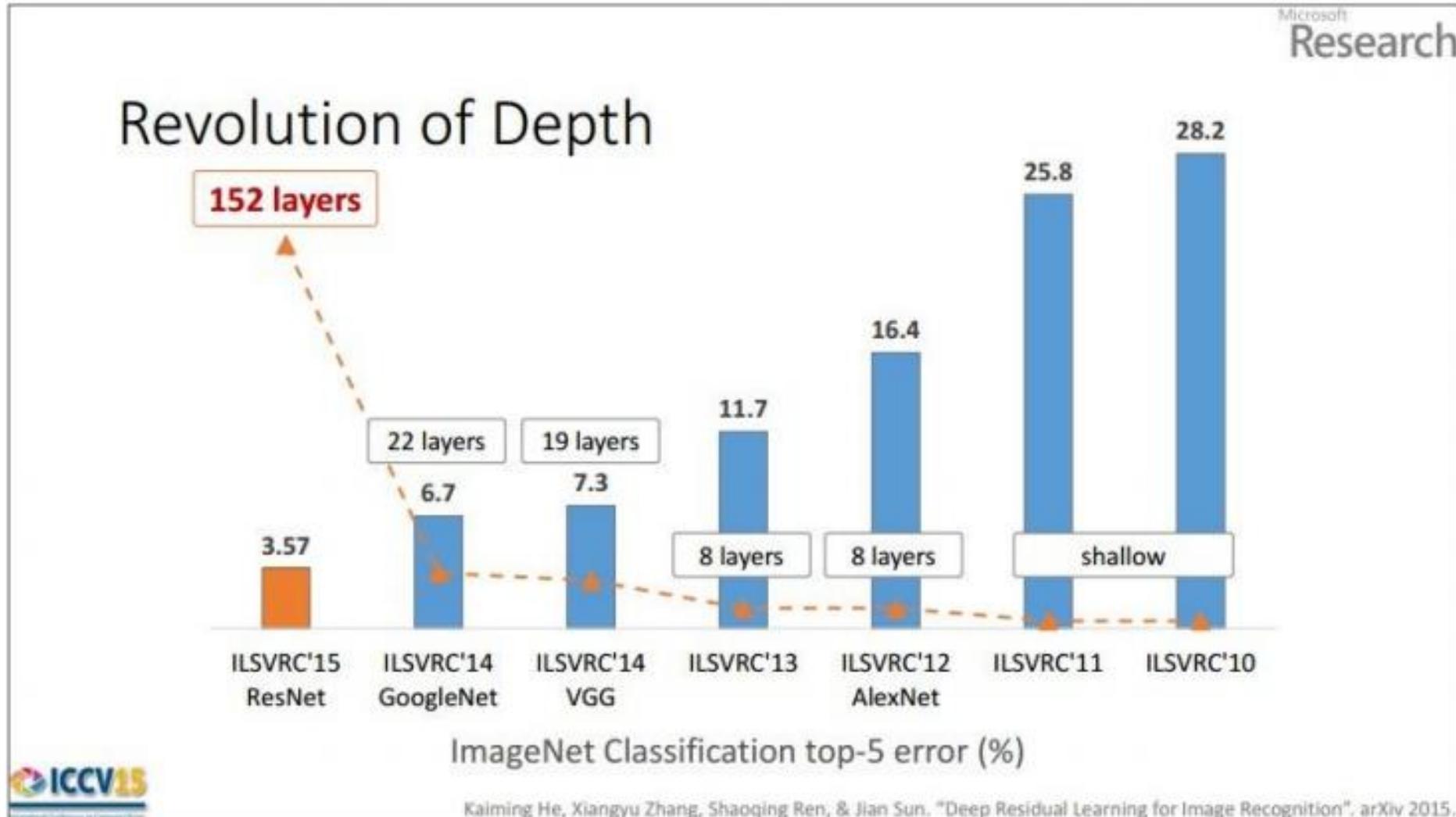
*improvements are relative numbers

ICCV15

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

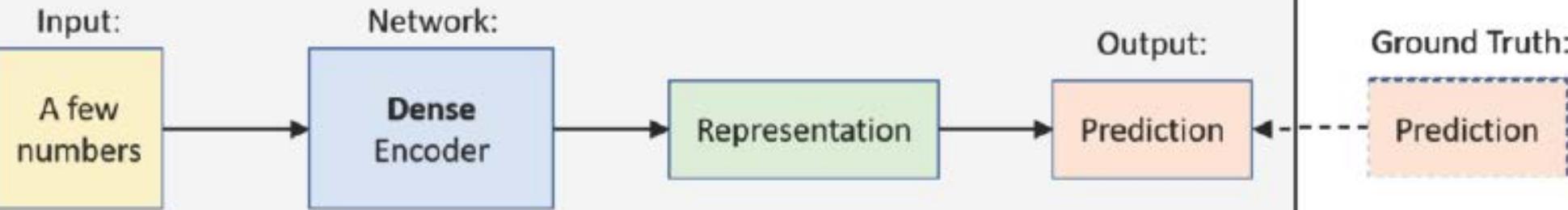
Slide from Kaiming He's recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>





(slide from Kaiming He's recent presentation)

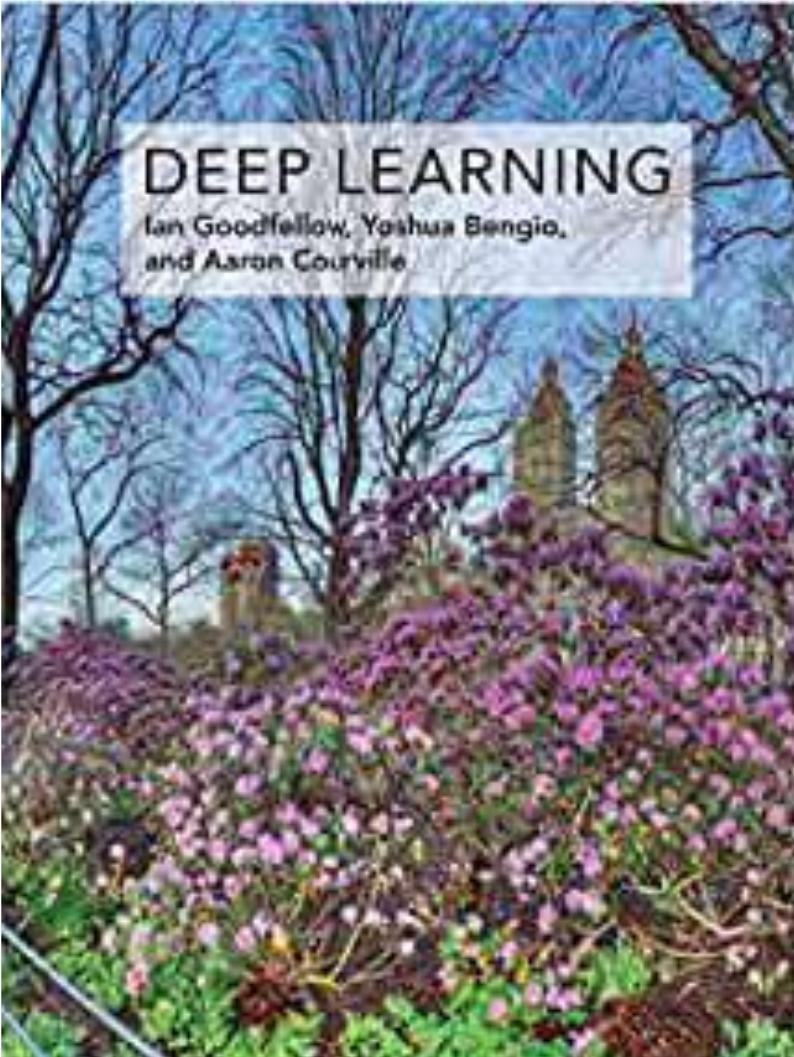
1. Feed Forward Neural Networks



2. Convolutional Neural Networks



$$\Delta w_{kj}(l) = -\frac{\partial E}{\partial w_{kj}(l)}$$



The International Journal of Robotics Research

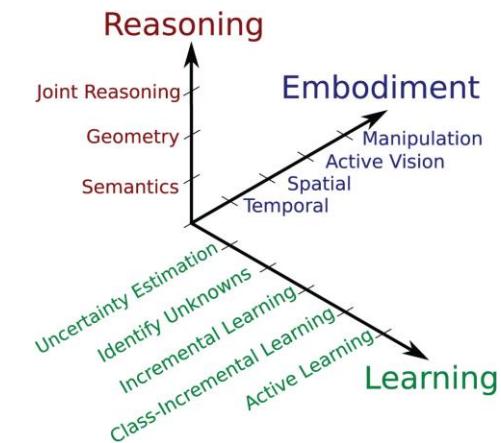
ijrr Impact Factor: 6.887 /

Available access | Research article | First published online April 27, 2018

The limits and potentials of deep learning for robotics

Niko Sünderhauf ✉, Oliver Brock, [...], and Peter Corke +8 View all authors and affiliations

Volume 37, Issue 4-5 | <https://doi.org/10.1177/0278364918770733>



(1) 一般介绍

(2) 从多层感知机到卷积神经网络

(3) 应用

➤ 目标检测



(x, y)



w, width, 宽度

→ CAT

h, height,
高度
→ (x, y, w, h)

Input: image



Neural Net

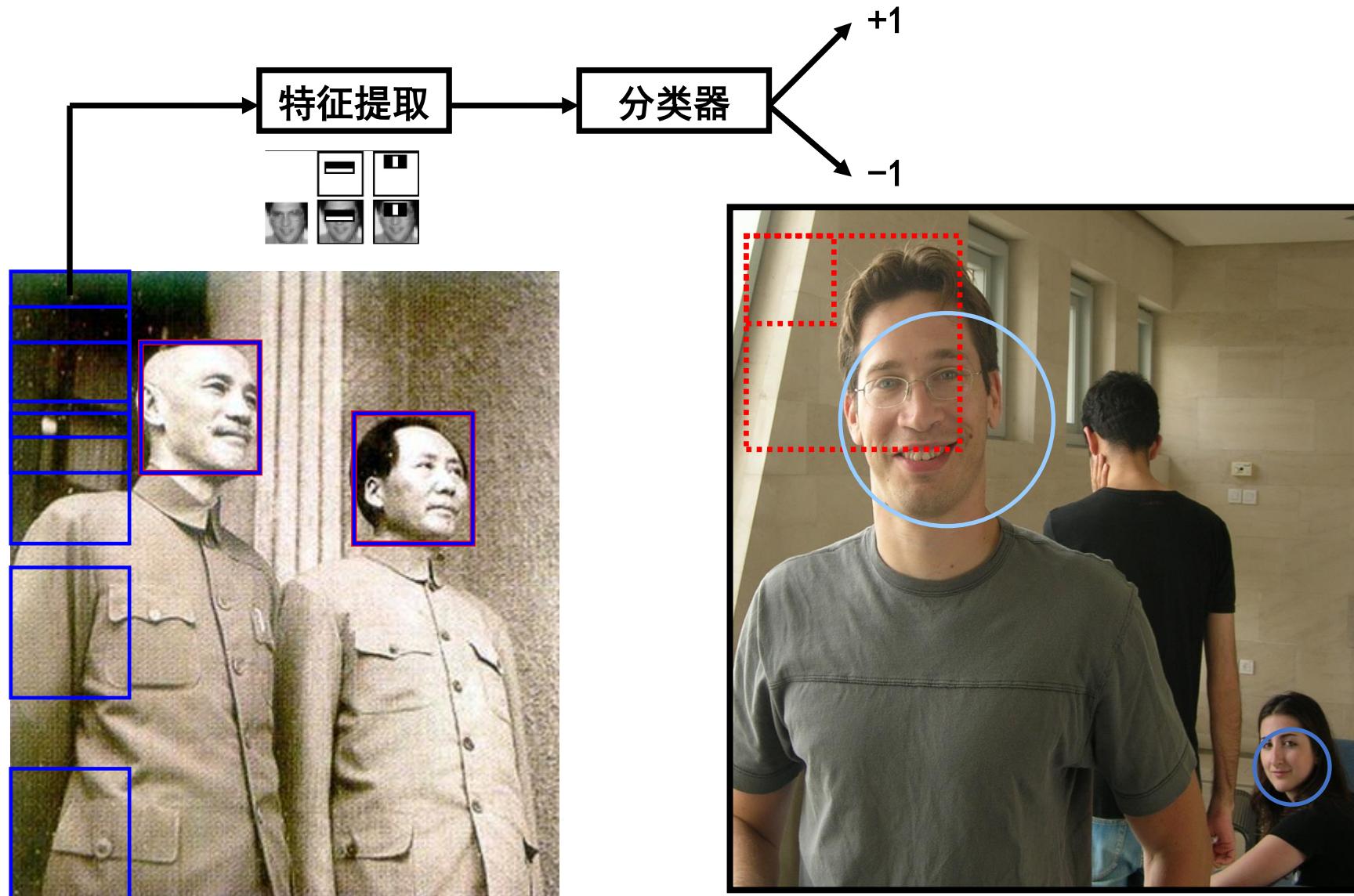
Output:
Box coordinates
(4 numbers)

Loss:
L2 distance
用欧氏距离作损失函数

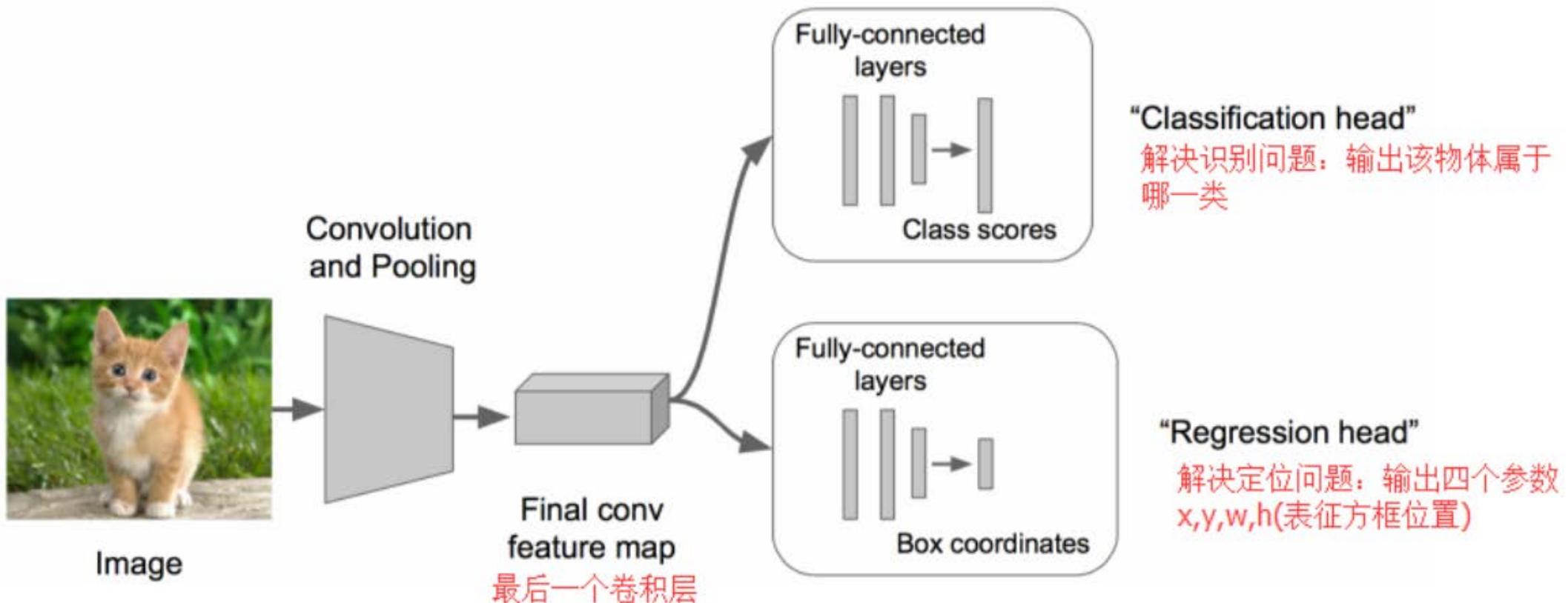
Correct output:
box coordinates
(4 numbers)

Only one object,
simpler than detection

➤ 目标检测：Adaboost



➤ 目标检测：深度学习的思路



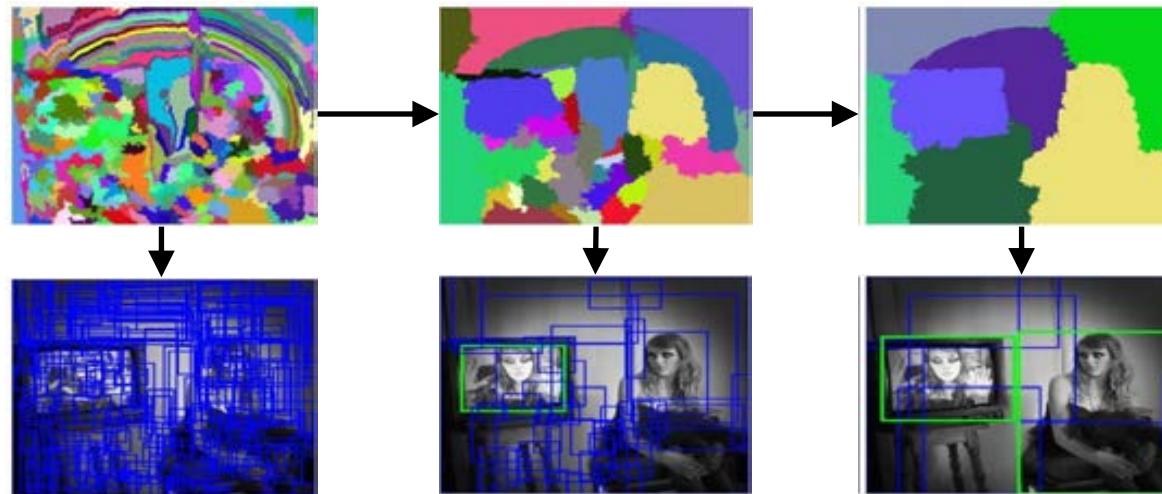
➤ 目标检测：深度学习的思路



Method	Approach	Outputs Segments	Outputs Score	Control #proposals	Time (sec.)	Repeatability	Recall Results	Detection Results
Bing [18]	Window scoring	✓	✓	✓	0.2	***	*	-
CPMC [19]	Grouping	✓	✓	✓	250	-	**	*
EdgeBoxes [20]	Window scoring	✓	✓	✓	0.3	**	***	***
Endres [21]	Grouping	✓	✓	✓	100	-	***	**
Geodesic [22]	Grouping	✓	✓	✓	1	*	***	**
MCG [23]	Grouping	✓	✓	✓	30	*	***	***
Objectness [24]	Window scoring	✓	✓	✓	3	-	*	-
Rahul [25]	Window scoring	✓	✓	✓	3	-	-	*
RandomizedPrim's [26]	Grouping	✓	✓	✓	1	*	*	**
Rantalaankila [27]	Grouping	✓	✓	✓	10	**	-	**
Rigor [28]	Grouping	✓	✓	✓	10	*	**	**
SelectiveSearch [29]	Grouping	✓	✓	✓	10	**	***	***
Gaussian SlidingWindow				✓	0	-	-	*
Superpixels				✓	0	***	-	-
Uniform				✓	0	-	-	-

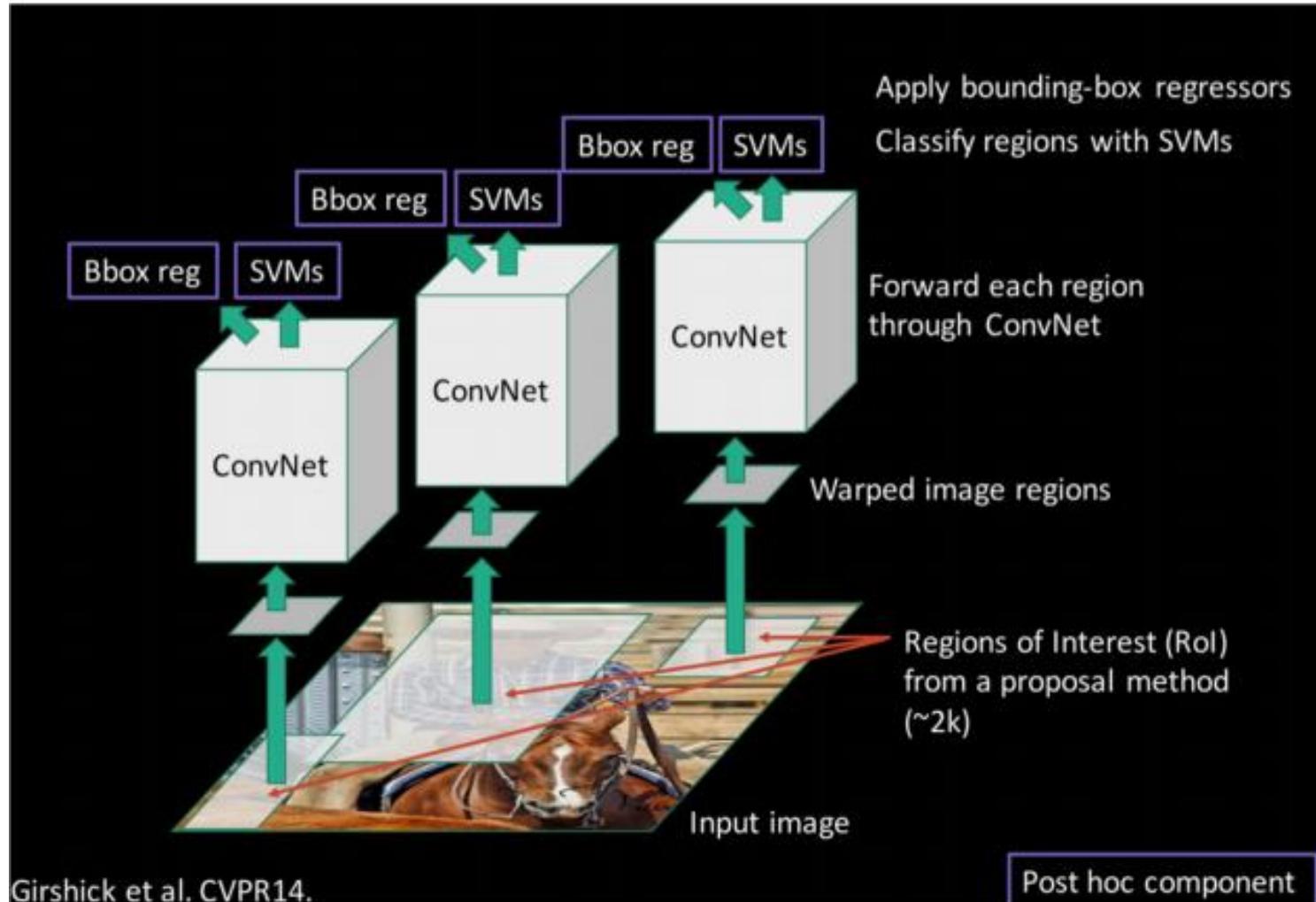
➤ 目标检测：深度学习的思路

Region proposals, Selective Search



Method	Approach	Outputs Segments	Outputs Score	Control #proposals	Time (sec.)	Repeatability	Recall Results	Detection Results
Bing [18]	Window scoring		✓	✓	0.2	***	*	.
CPMC [19]	Grouping	✓	✓	✓	250	-	**	*
EdgeBoxes [20]	Window scoring		✓	✓	0.3	**	***	***
Endres [21]	Grouping	✓	✓	✓	100	-	***	**
Geodesic [22]	Grouping	✓		✓	1	*	***	**
MCG [23]	Grouping	✓	✓	✓	30	*	***	***
Objectness [24]	Window scoring		✓	✓	3	.	*	.
Rahtu [25]	Window scoring		✓	✓	3	.	.	*
RandomizedPrim's [26]	Grouping	✓		✓	1	*	*	**
Rantalaikila [27]	Grouping	✓		✓	10	**	-	**
Rigor [28]	Grouping	✓		✓	10	*	**	**
SelectiveSearch [29]	Grouping	✓	✓	✓	10	**	***	***
Gaussian SlidingWindow				✓	0	.	.	*
Superpixels		✓		✓	0	***	.	.
Uniform				✓	1	*	.	.
					0	.	.	.

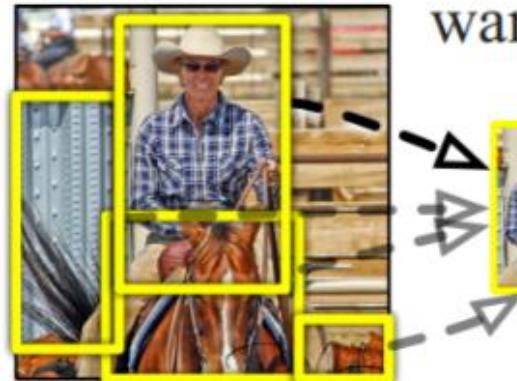
➤ 目标检测：R-CNN



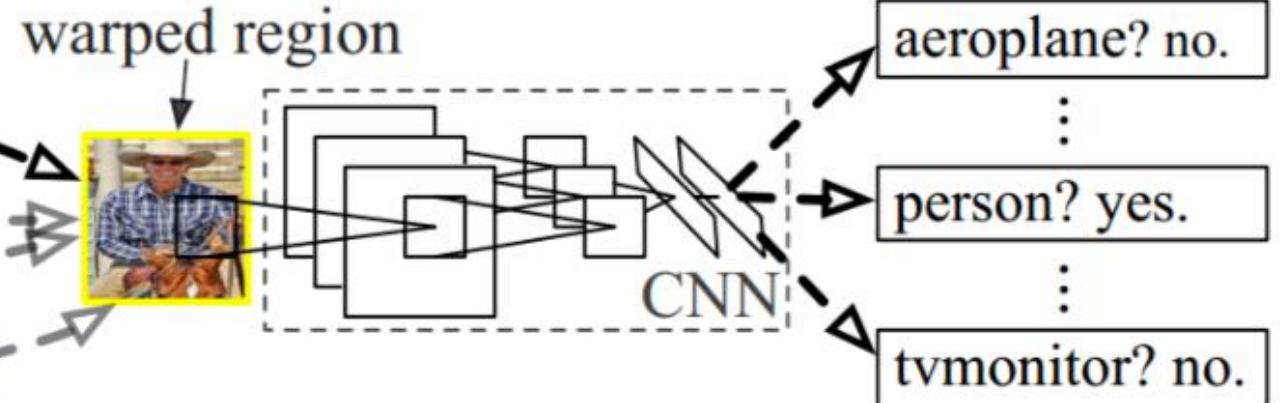
➤ 目标检测：R-CNN



1. Input
image



2. Extract region
proposals (~2k)

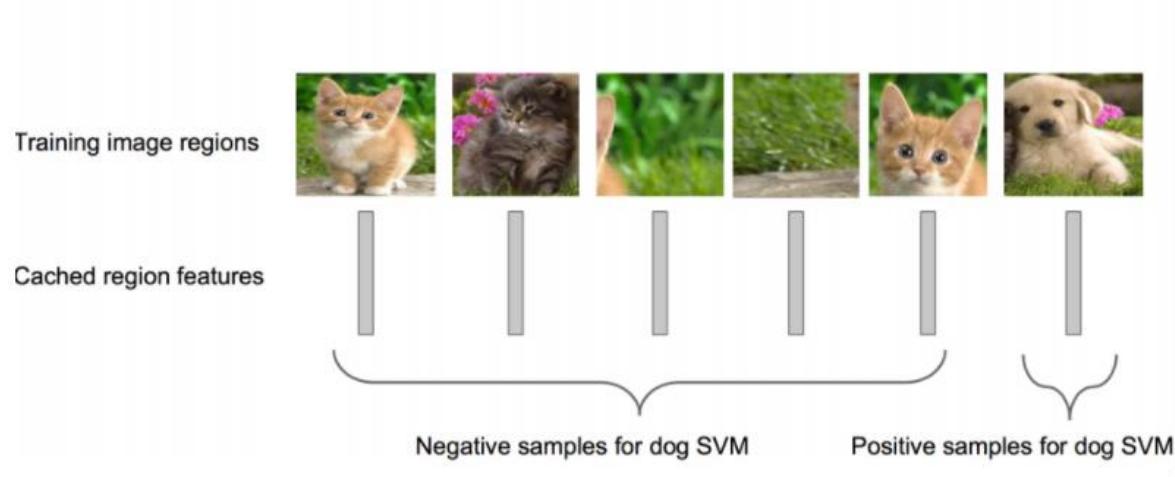


3. Compute
CNN features

4. Classify
regions

- R. Girshick, J. Donahue, T. Darrell, J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR 2014*

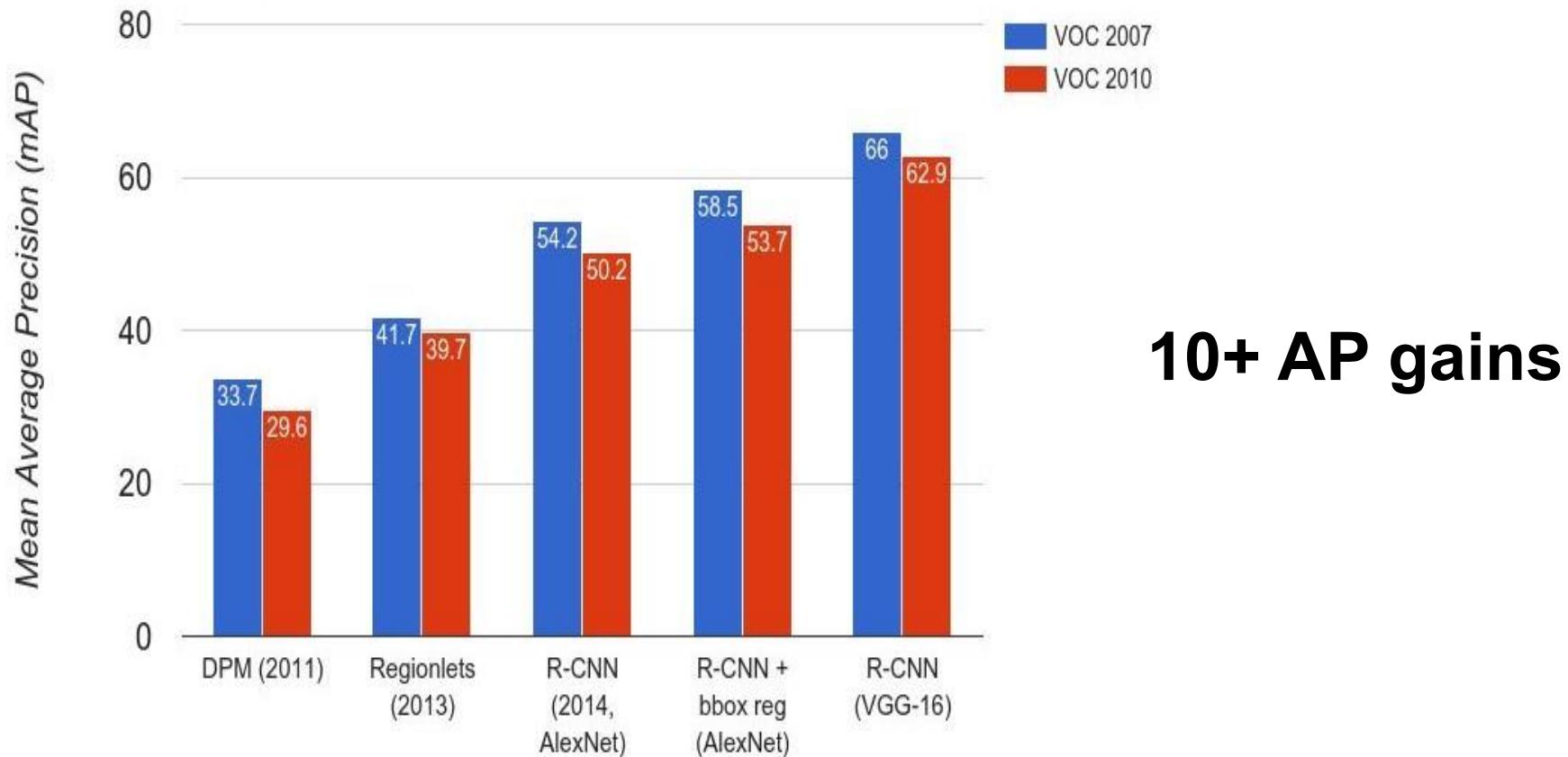
➤ 目标检测：R-CNN



步骤五：使用回归器精细修正候选框位置：对于每一个类，训练一个线性回归模型去判定这个框是否框得完美

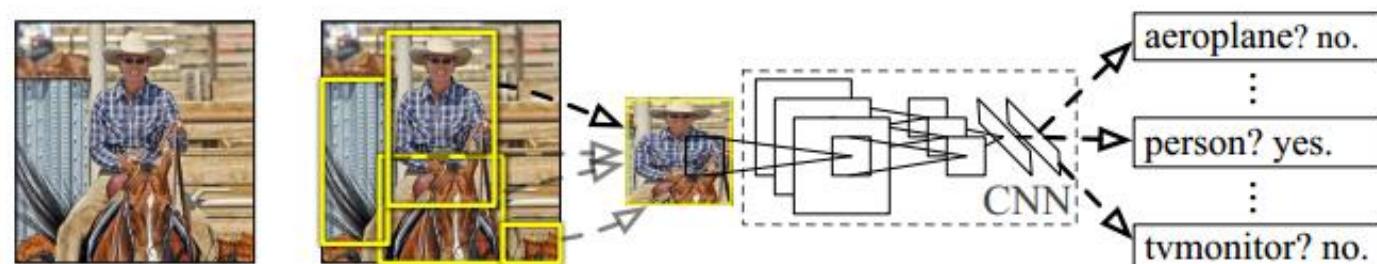
Training image regions			
Cached region features			
Regression targets (dx, dy, dw, dh) Normalized coordinates	(0, 0, 0, 0) Proposal is good 把猫完全框出来了， 很完美	(.25, 0, 0, 0) Proposal too far to left 框得偏左了	(0, 0, -0.125, 0) Proposal too wide 框画得太大了， 空白太多

➤ 目标检测：R-CNN

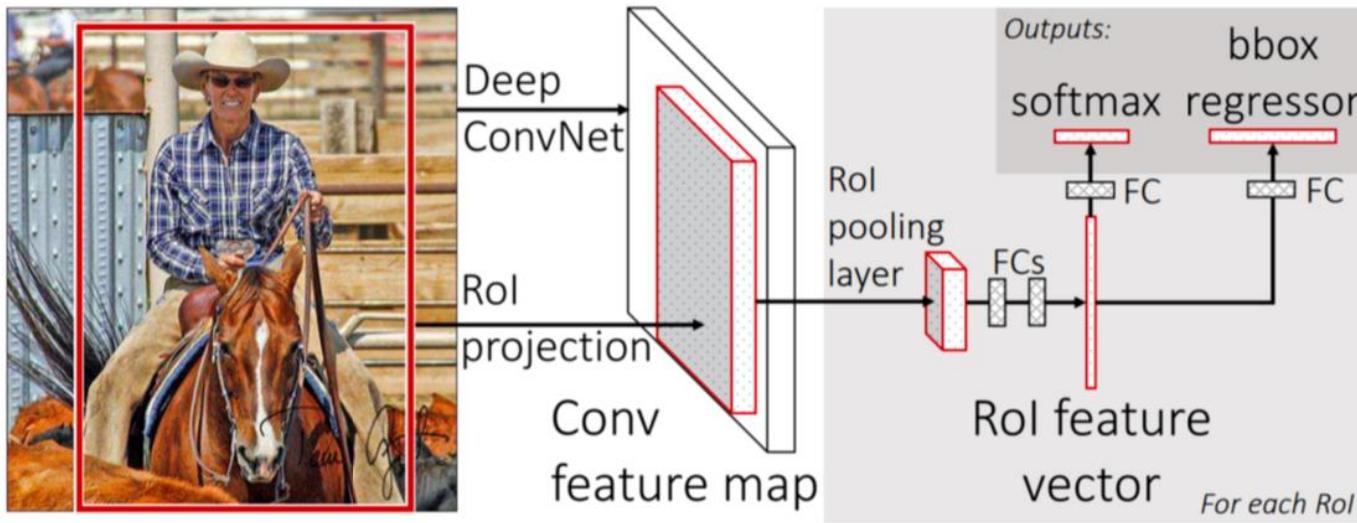


➤ 目标检测：R-CNN

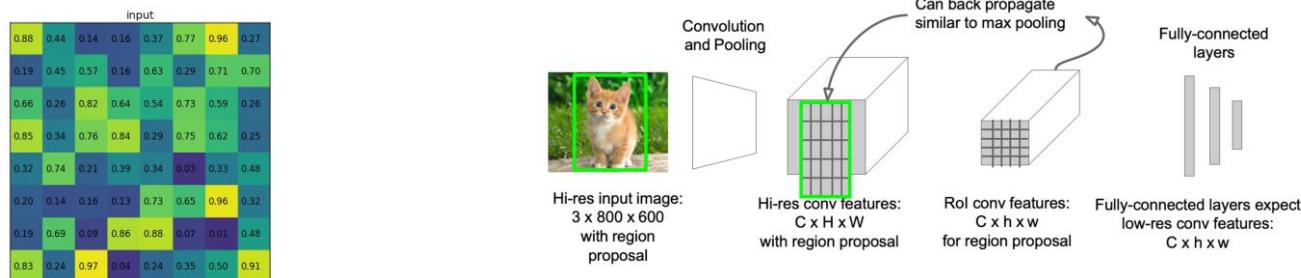
- Training and testing is **expensive in space and time**.
 - --Why: for each image, one need to extract CNN features for about 2k object proposals
 - and make classification. But the features are shared in space.
- Relays on the region proposal method (selective search)
- Training is a **multi-stage** pipeline



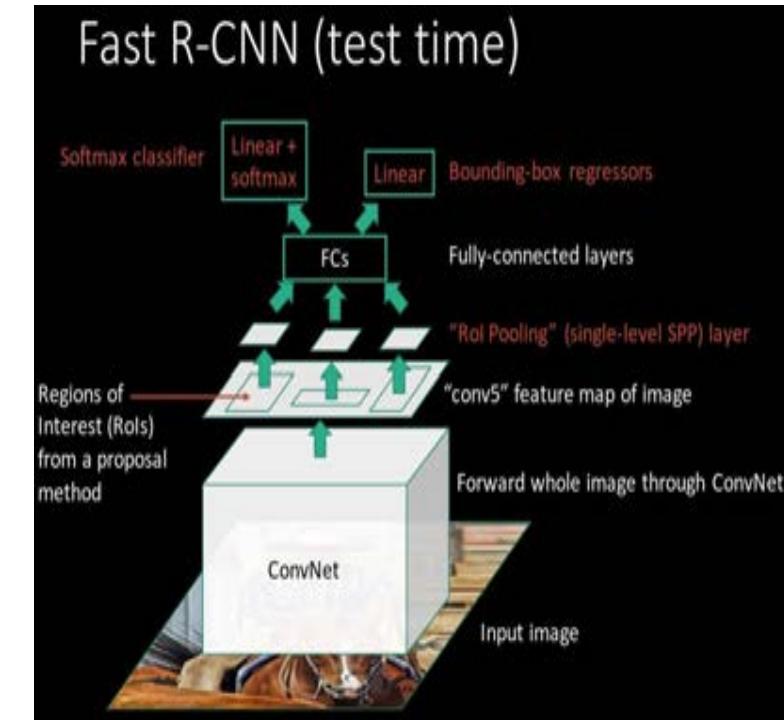
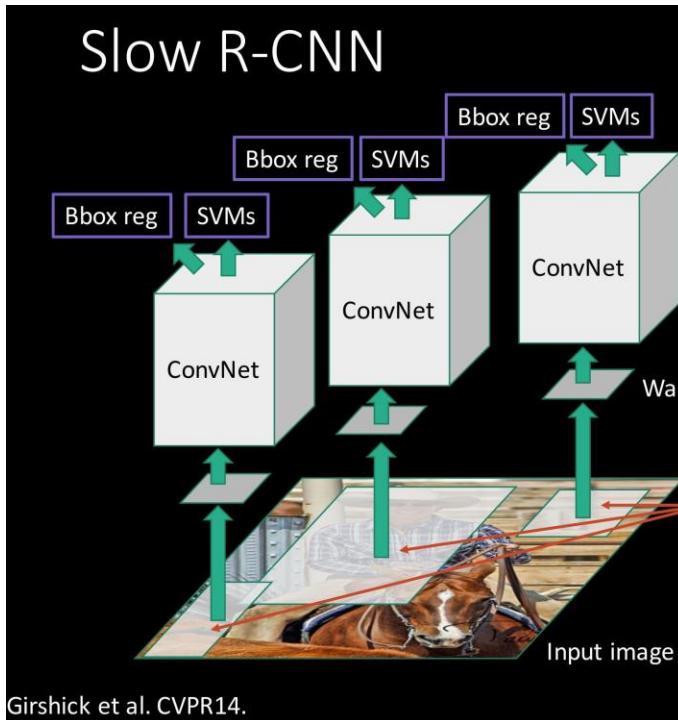
➤ 目标检测：Fast RCNN



Sliding windows on the feature maps



➤ 目标检测：Fast RCNN



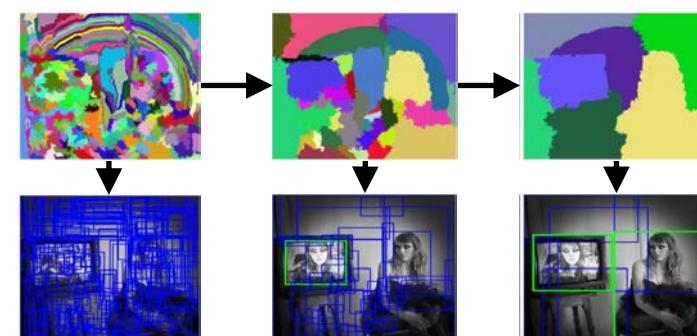
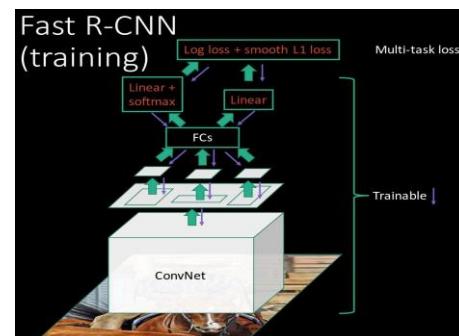
➤ 目标检测：Fast RCNN

		R-CNN	Fast R-CNN
Faster!	Training Time:	84 hours	9.5 hours
	(Speedup)	1x	8.8x
FASTER!	Test time per image	47 seconds	0.32 seconds
	(Speedup)	1x	146x
Better!	mAP (VOC 2007)	66.0	66.9

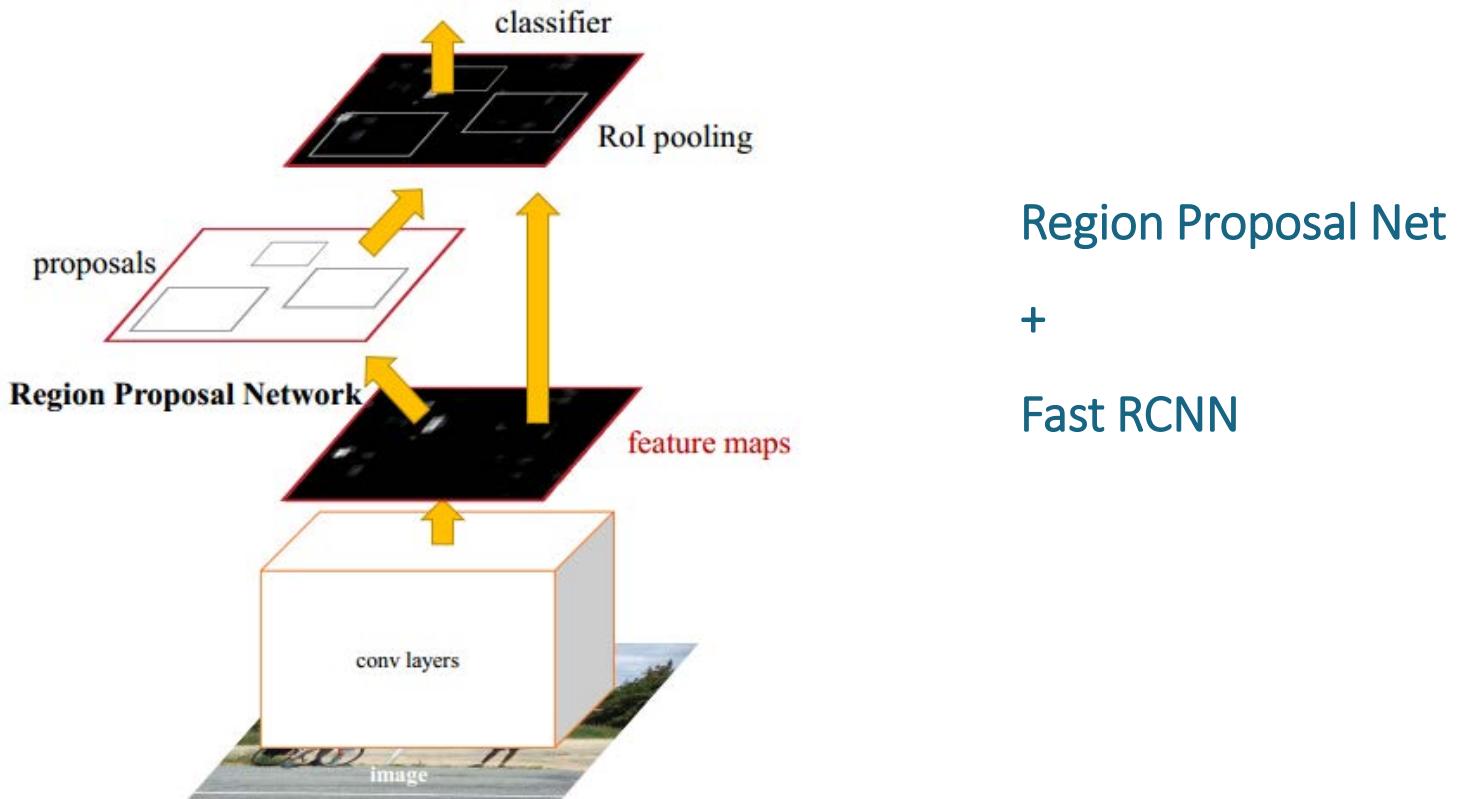
Using VGG-16 CNN on Pascal VOC 2007 dataset

➤ 目标检测：Fast RCNN

	R-CNN	Fast R-CNN
Test time per image	47 seconds	0.32 seconds
(Speedup)	1x	146x
Test time per image with Selective Search	50 seconds	2 seconds
(Speedup)	1x	25x

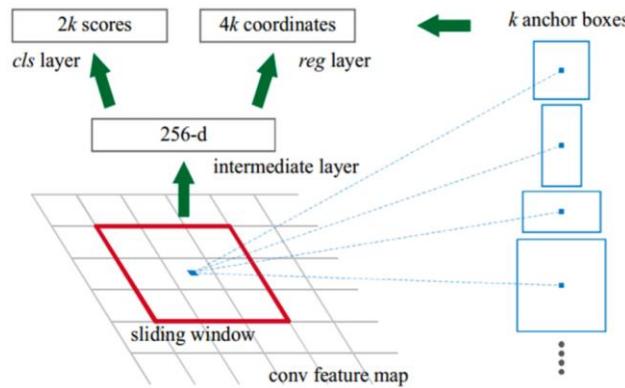


➤ 目标检测：Faster RCNN



➤ 目标检测：Faster RCNN

- Region proposal network

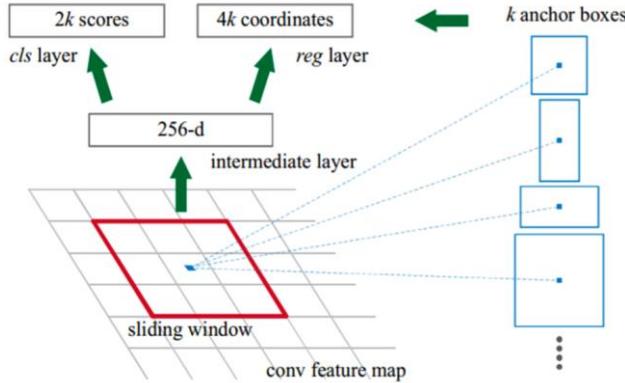


- Input an image of any size
- Generate **conv feature map**
- Map to **a lower-dimensional feature**
- Output **objectness score** and **bounding box**

The region proposal network is a FCN which outputs $K*(4+2)$ sized vectors.

➤ 目标检测：Faster RCNN

- Region proposal network



LOSS

- Positive: Among K anchors, one with highest IOU ($\text{IOU} \geq 0.7$)
- Negative: $\text{IOU} \leq 0.3$
- Others: Do not contribute

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*)$$

The mini-batch size(256)

The number of anchor locations

➤ 目标检测：Faster RCNN

	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image (with proposals)	50 seconds	2 seconds	0.2 seconds
(Speedup)	1x	25x	250x
mAP (VOC 2007)	66.0	66.9	73.2

➤ 目标检测：总结

- **RCNN**

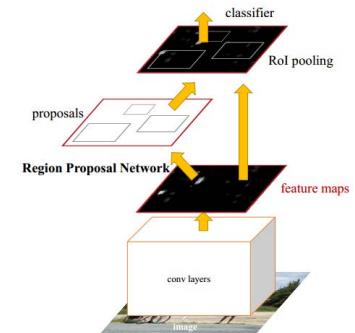
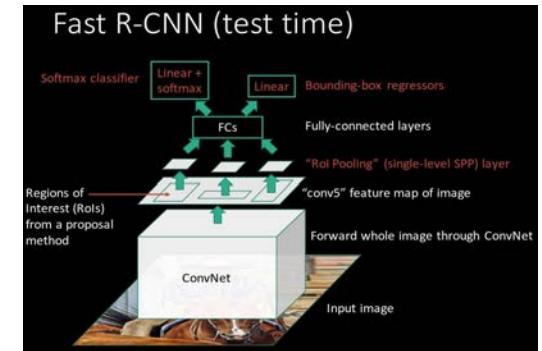
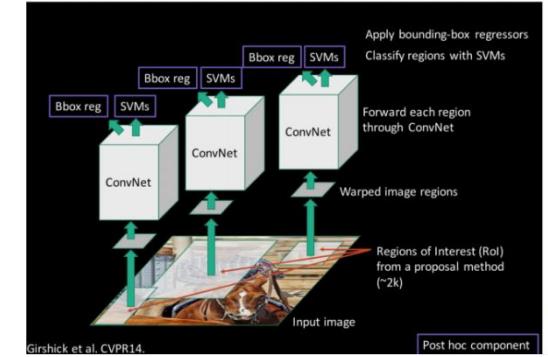
1. 在图像中确定约1000-2000个候选框 (使用选择性搜索)
2. 每个候选框内图像块缩放至相同大小，并输入到CNN内进行特征提取
3. 对候选框中提取出的特征，使用分类器判别是否属于一个特定类
4. 对于属于某一特征的候选框，用回归器进一步调整其位置

- **Fast RCNN**

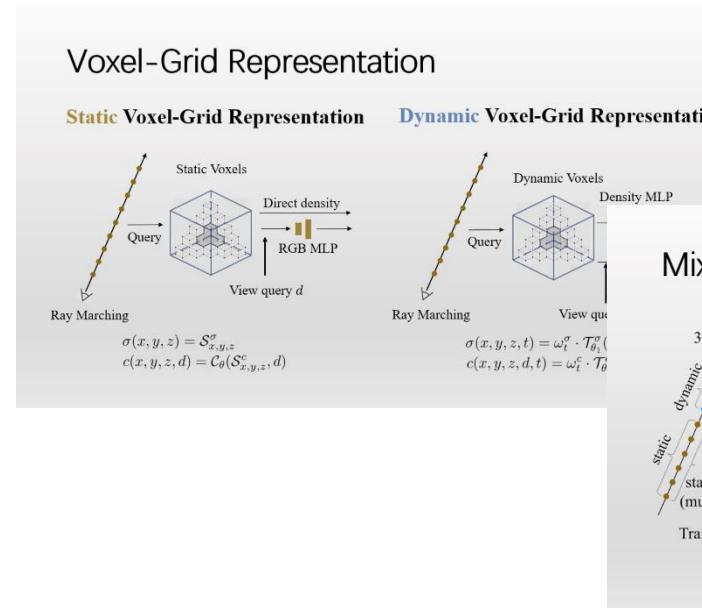
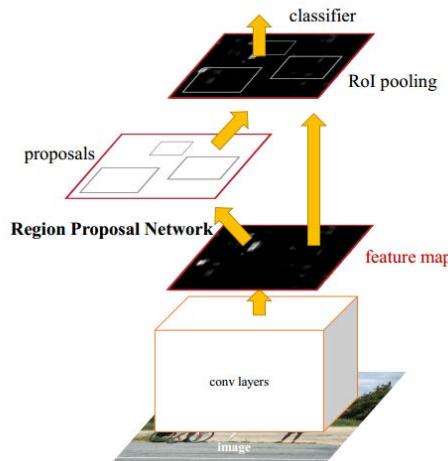
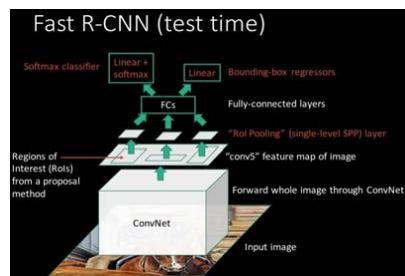
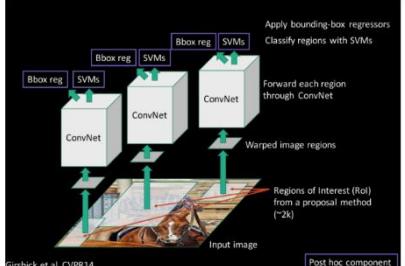
1. 在图像中确定约1000-2000个候选框 (使用选择性搜索)
2. 对整张图片输入CNN，得到feature map
3. 找到每个候选框在feature map上的映射patch，将此patch作为每个候选框的卷积特征输入到SPP layer和之后的层
4. 对候选框中提取出的特征，使用分类器判别是否属于一个特定类
5. 对于属于某一特征的候选框，用回归器进一步调整其位置

- **Faster RCNN**

1. 对整张图片输入CNN，得到feature map
2. 卷积特征输入到RPN，得到候选框的特征信息
3. 对候选框中提取出的特征，使用分类器判别是否属于一个特定类
4. 对于属于某一特征的候选框，用回归器进一步调整其位置

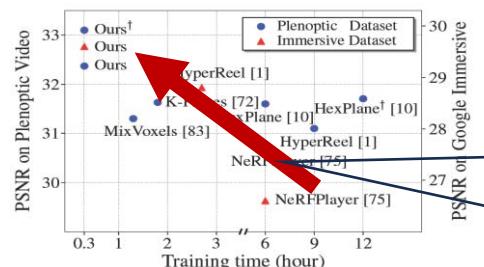


➤ 目标检测：扩展



- Accelerate Training by using lightweight static model to process static queries which is dominated in most scenes.
- Accelerate Rendering by using lightweight model to process static queries and only render once for static points.

动态场景的三维重构



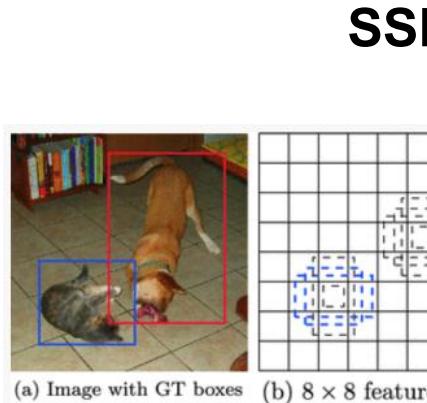
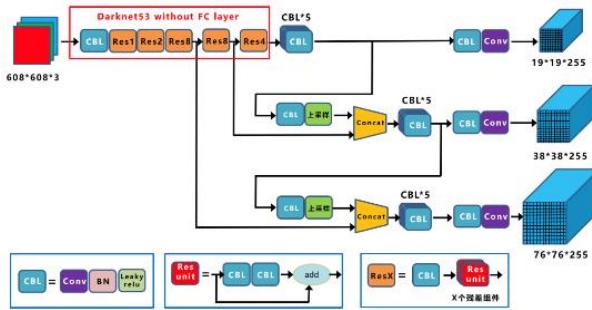
动态场景：渲染时间从12小时提升至20分钟，画质更高



- Mixed neural voxels for fast multi-view video synthesis, **ICCV**, 2023, **Oral**
- Masked space-time Hash encoding for efficient dynamic scene reconstruction, **NeurIPS**, 2023, **Spotlight**

➤ 目标检测：发展

Yolo



使用one-stage的方式，直接通过一系列CBL模块（Conv-Batch Normalization-Leaky Relu）提出特征金字塔，并在此基础上根据预先设置的anchor prior对各region进行分类和回归

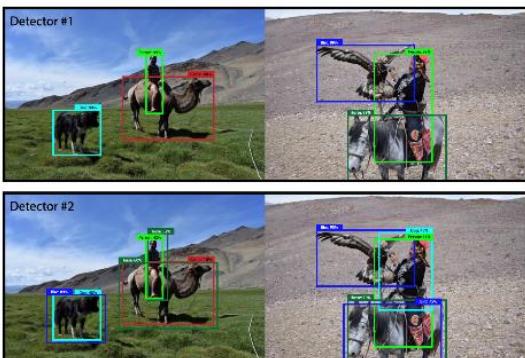
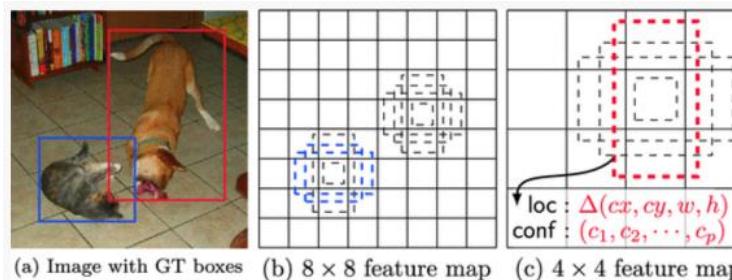


Figure 5. These two hypothetical detectors are perfect according to mAP over these two images. They are both perfect. Totally equal.

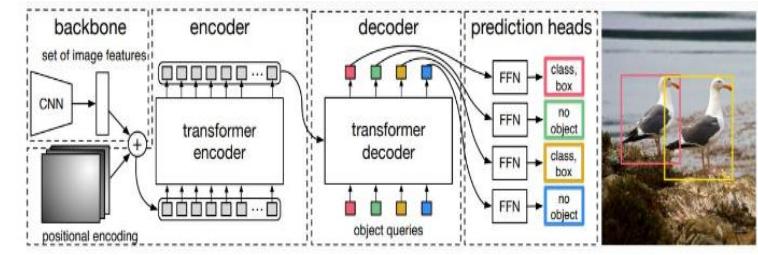
SSD



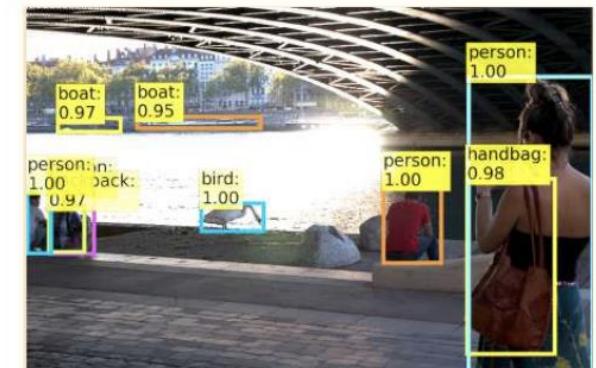
SSD是one-stage检测器，其速度要优于两阶段检测器，但是精度上与有些许损失。SSD较其他单阶段检测器相比其最大的亮点在于提出利用浅层特征检测小目标而深层特征检测大目标。



DeTR



将CNN的feature map拉平后视为一个序列，并通过transformer encoder进行编码。之后通过一些可学习的位置编码作为decoder的输入，将解码后的各vector送入到前向传播网络中，对其类别和包围盒进行预测和回归



➤ 目标检测：发展

YOLO-V5: <https://github.com/ultralytics/yolov5>, 最新的YOLO检测器地址，其安装和推理十分简便

▼ Inference

Inference with YOLOv5 and [PyTorch Hub . Models](#) download automatically from the latest YOLOv5 [release](#).

```
import torch

# Model
model = torch.hub.load('ultralytics/yolov5', 'yolov5s') # or yolov5m, yolov5l, yolov5x, custom

# Images
img = 'https://ultralytics.com/images/zidane.jpg' # or file, Path, PIL, OpenCV, numpy, list

# Inference
results = model(img)

# Results
results.print() # or .show(), .save(), .crop(), .pandas(), etc.
```

➤ 目标检测：发展

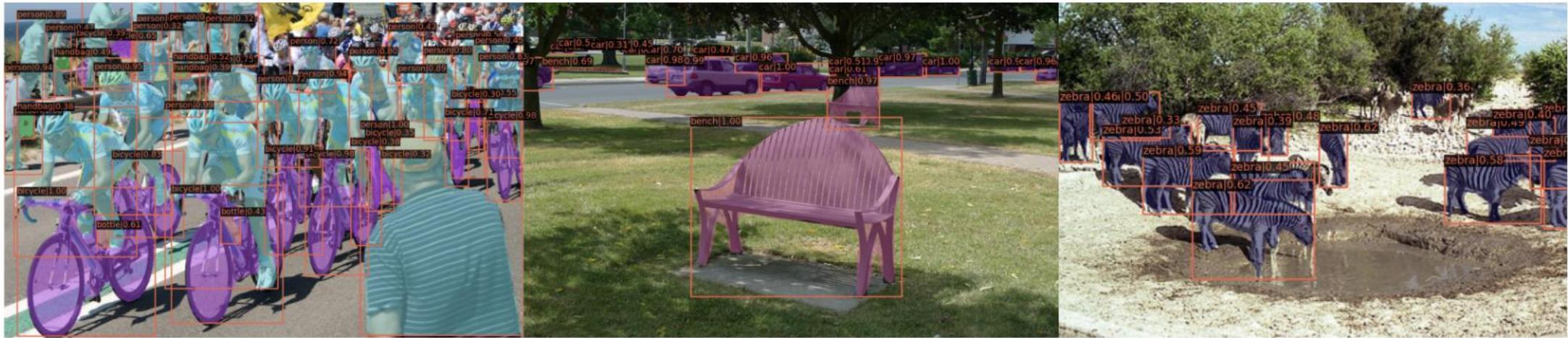
Faster RCNN: <https://github.com/jwyang/faster-rcnn.pytorch>, 通过pytorch实现的faster rcnn，需要先通过C++编译目标检测中常用的函数（如RoIPooling、NMS等），之后进行训练。Backbone网络包括VGG-16和ResNet-101可以使用。

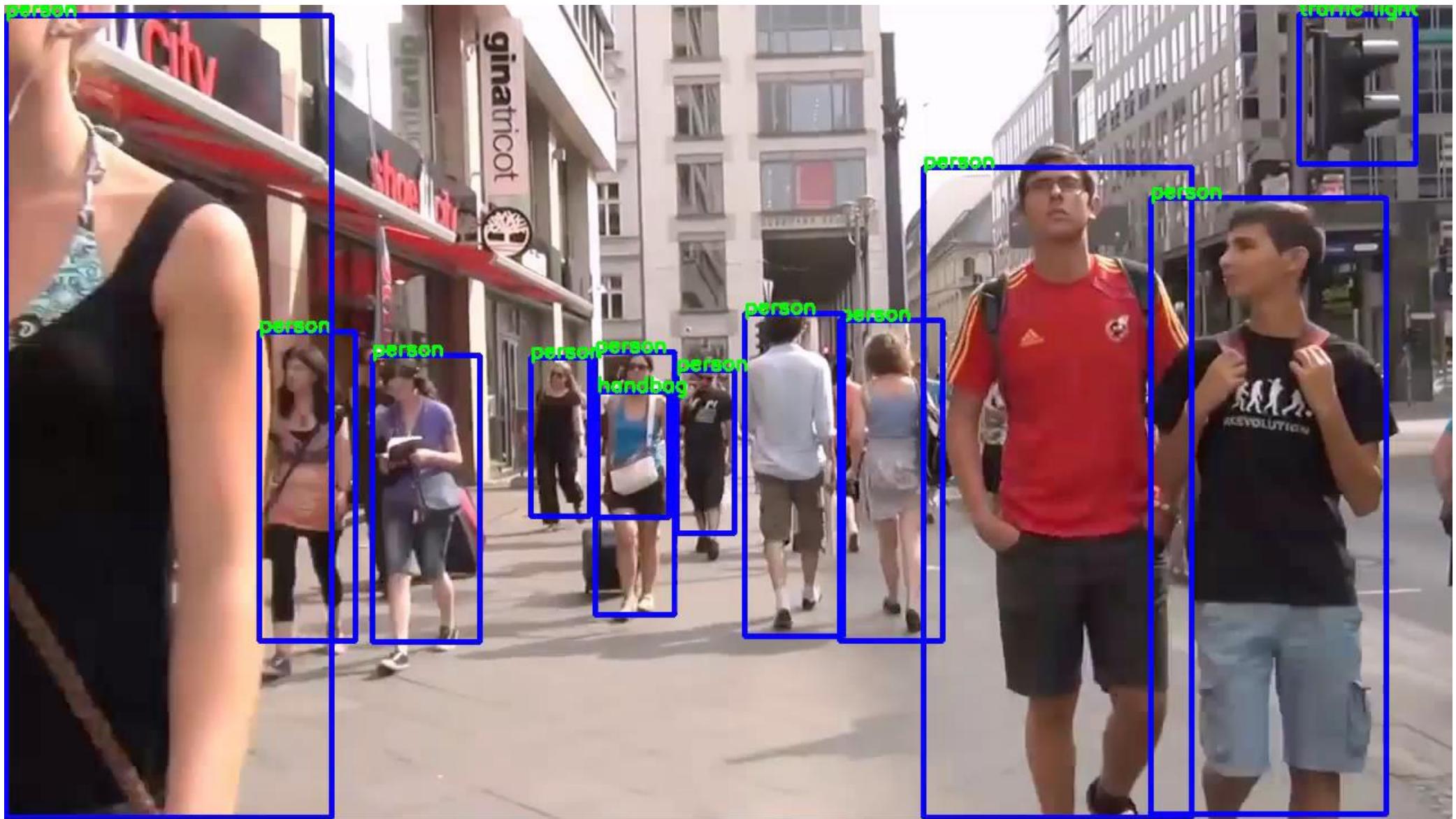


➤ 目标检测：发展

mmdetection: <https://github.com/open-mmlab/mmdetection>

mmdetection是通过pytorch实现的多个object detector的开源库，可通过制定config文件和模型参数来通过init_detector方法加载不同的detector，并通过inference_detector方法来对image进行目标检测。



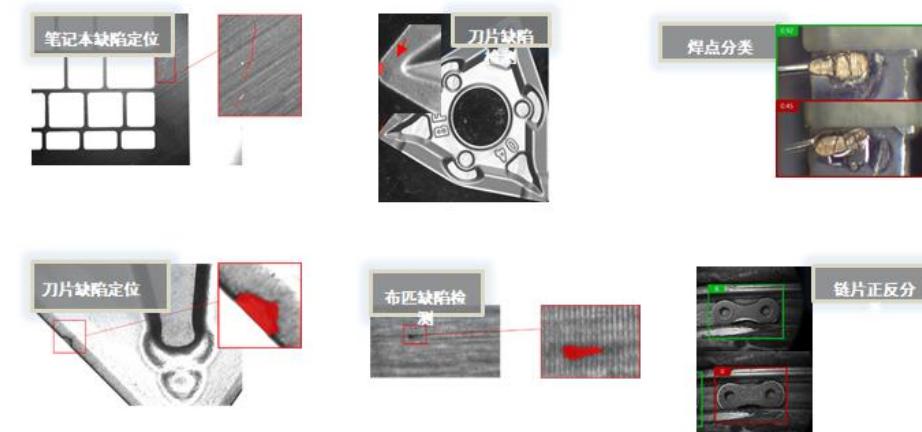
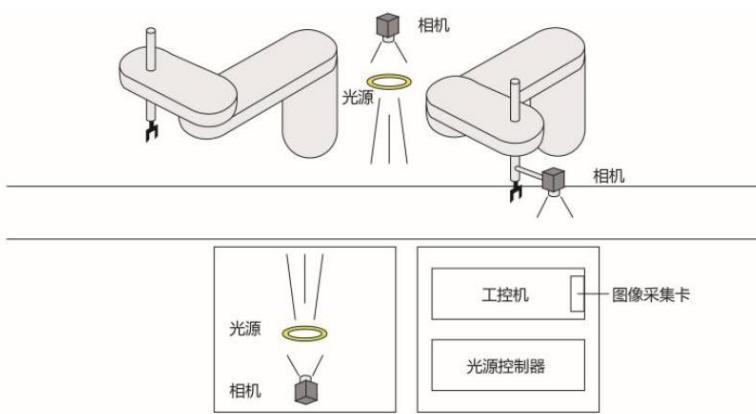
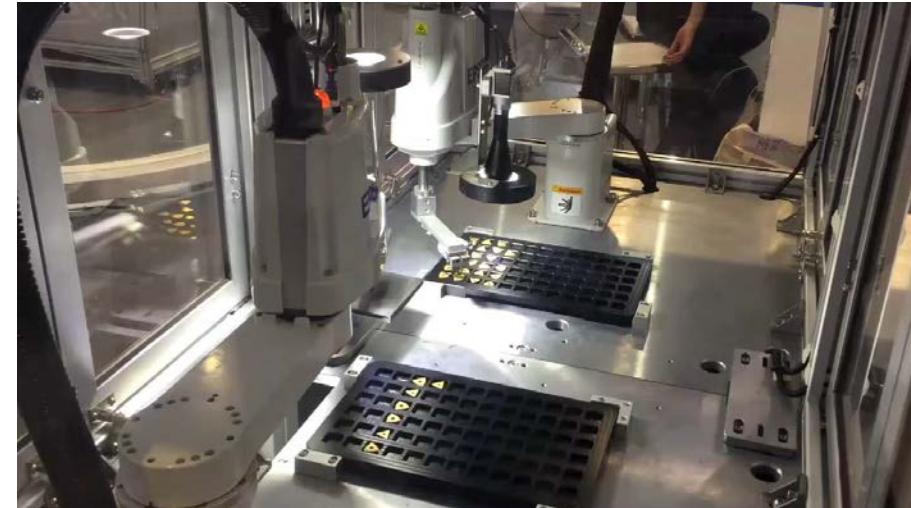


➤ 目标检测

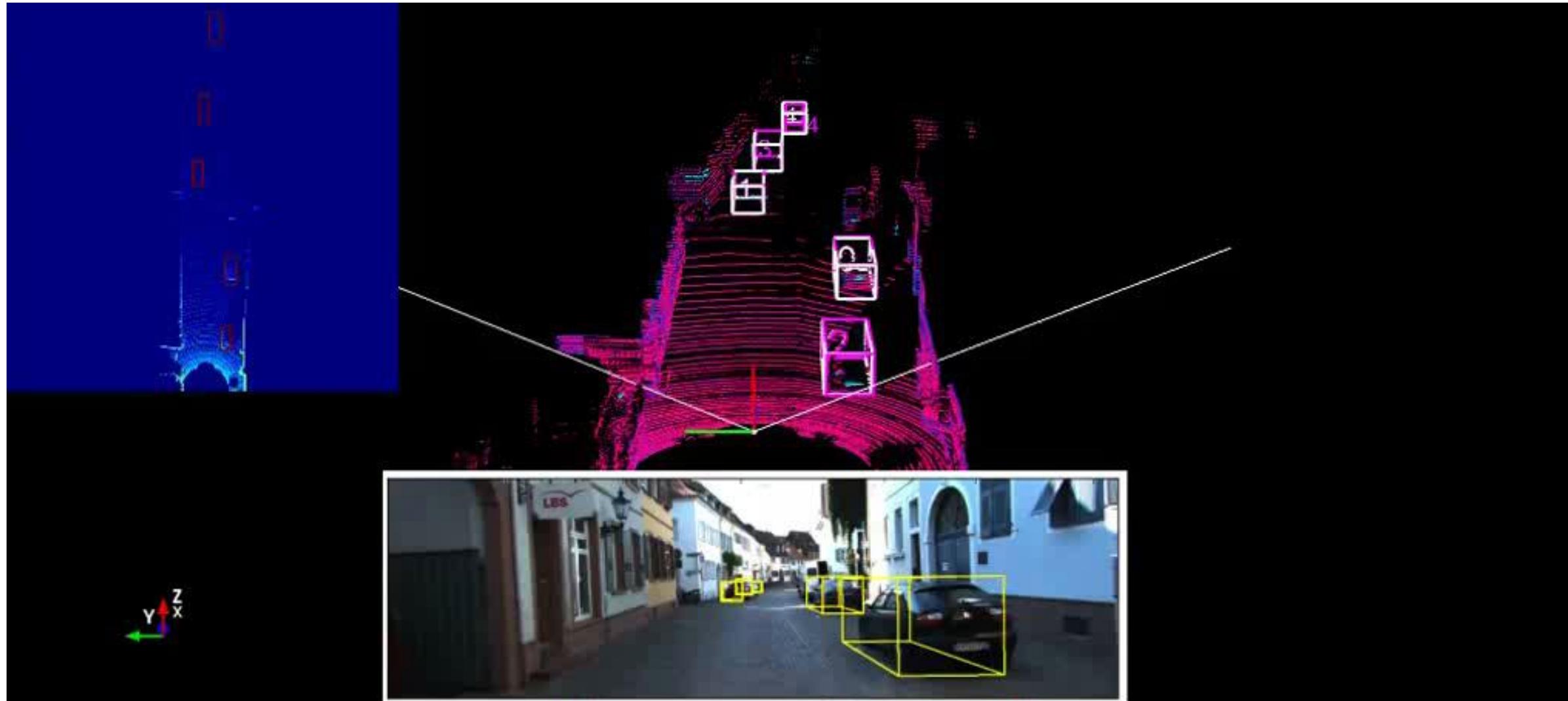


➤ 目标检测

具有持续自学习能力的智能视觉缺陷检测系统



➤ 目标检测：多模态



➤ 目标检测：多模态

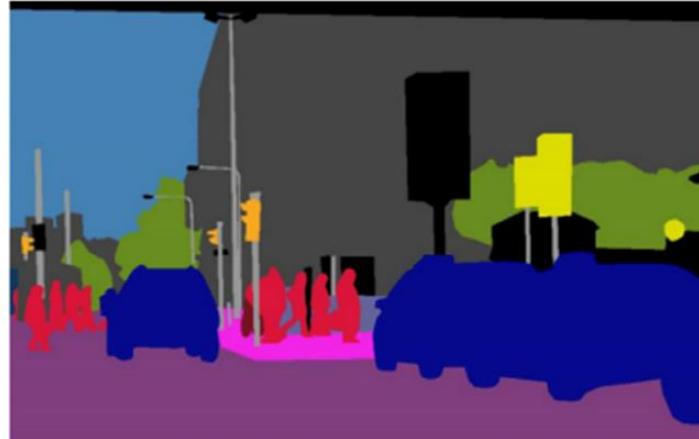
- 传统方法
- 机器学习



➤ 语义分割



(a) Image



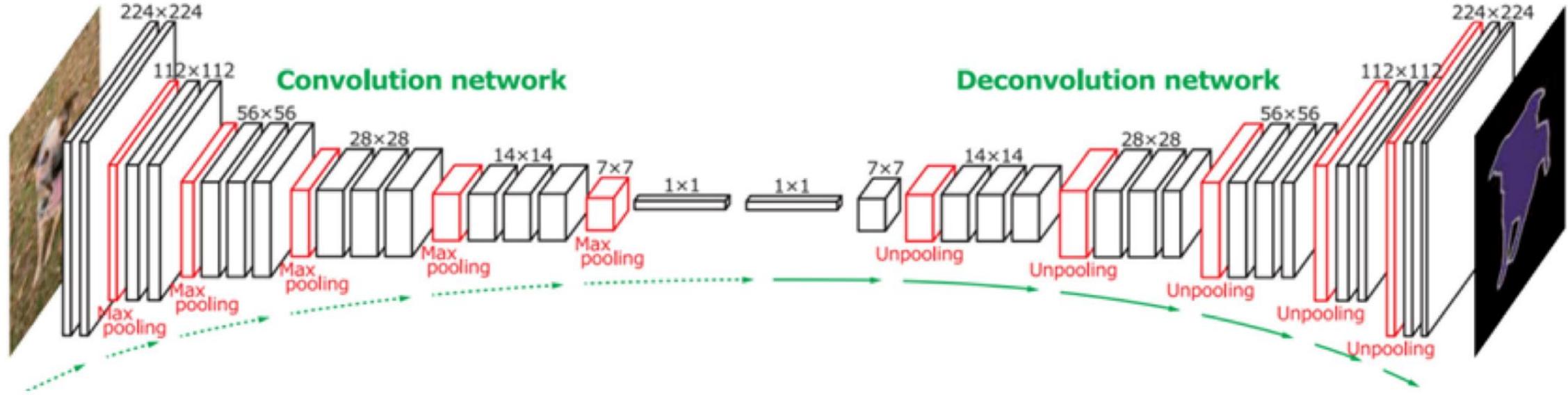
(b) Semantic segmentation



(c) Image classification

语义性是图像中很重要的性质，准确的语义信息可以提升许多下游任务的表现。例如机器人的操作（manipulation），导航（navigation），图像/视频的问答（Visual / Video QA），图像/视频的标注（Visual / Video captioning）等等。并且，最近用于自监督预训练的对比学习（contrastive learning）也是通过contrastive loss来帮助模型学习图像的语义性。因此，语义分割是计算机视觉中十分重要的应用，且语义分割的好坏可以一定程度上反应由网络抽取的视觉特征的好坏。

➤ 语义分割



经典的语义分割由一系列down sampling和up sampling组成。如上图所示，原始的224x224的图像先通过一系列卷积层（convolution layer）池化层（max pooling）得到了2048x1x1的向量。之后，通过一系列对称的反卷积操作将该向量上采样得到224x224的张量。并对其中每个像素点进行回归，来让网络学习每个pixel对应的semantic label。也可以在一系列下采样后直接在feature map上进行预测，并对预测的结果进行差值还原到与原始图像相同的尺寸。

➤ 语义分割

What is COCO?



COCO is a large-scale object detection, segmentation, and captioning dataset.
COCO has several features:

- ✓ Object segmentation
- ✓ Recognition in context
- ✓ Superpixel stuff segmentation
- ✓ 330K images (>200K labeled)
- ✓ 1.5 million object instances
- ✓ 80 object categories
- ✓ 91 stuff categories
- ✓ 5 captions per image
- ✓ 250,000 people with keypoints

<https://cocodataset.org/#home>

COCO dataset 是微软构建的一个数据集，其中包含了80类为物体，33万张图片，可以支持object detection、semantic segmentation等任务。

每个Image不但提供了各类物体的label和bounding box坐标，还额外给出了mask信息，即bounding box中哪些像素点是该物体。



2010: Speech Recognition



2012: Computer Vision



2014: Machine Translation

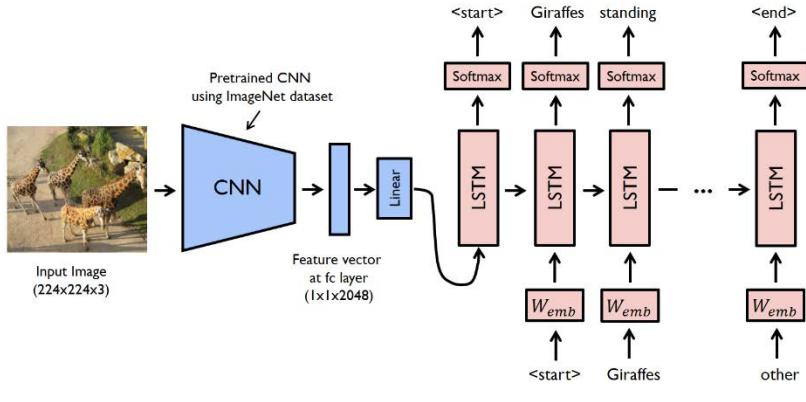


Source: 我 在 周日 看 了 一 本 书

Target: I read a book on Sunday



➤ 语言生成



前方道路有车辆行驶，请等待



谢 谢 !