

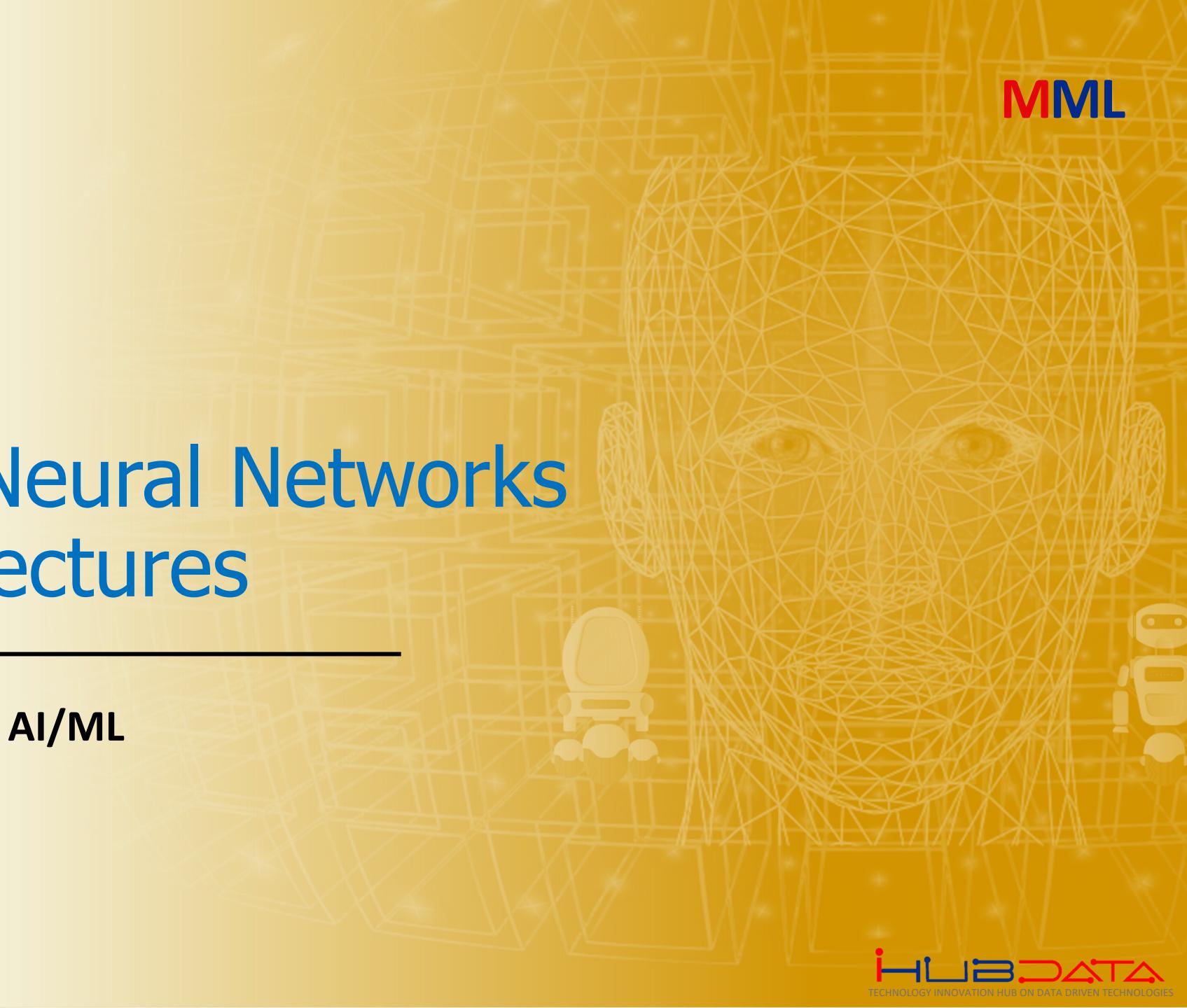


MML

Convolutional Neural Networks Architectures

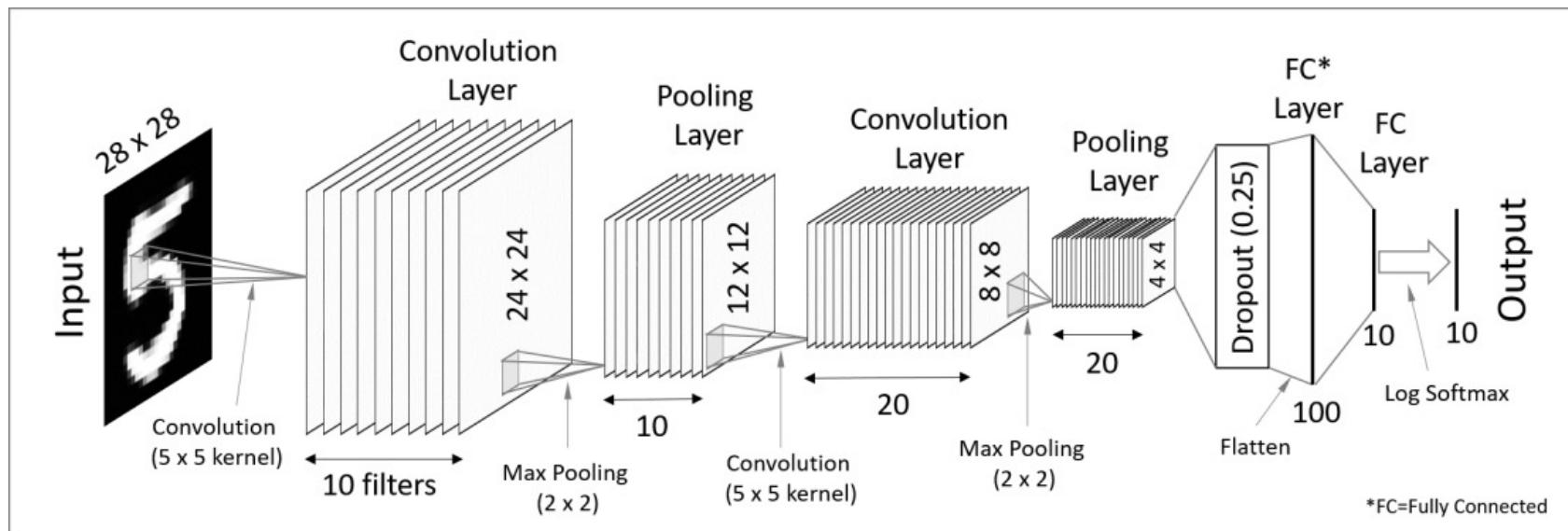
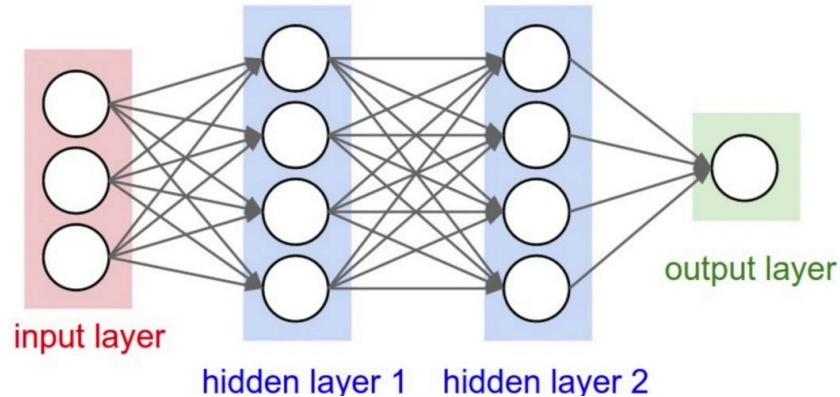
Executive Program on AI/ML

IIIT Hyderabad



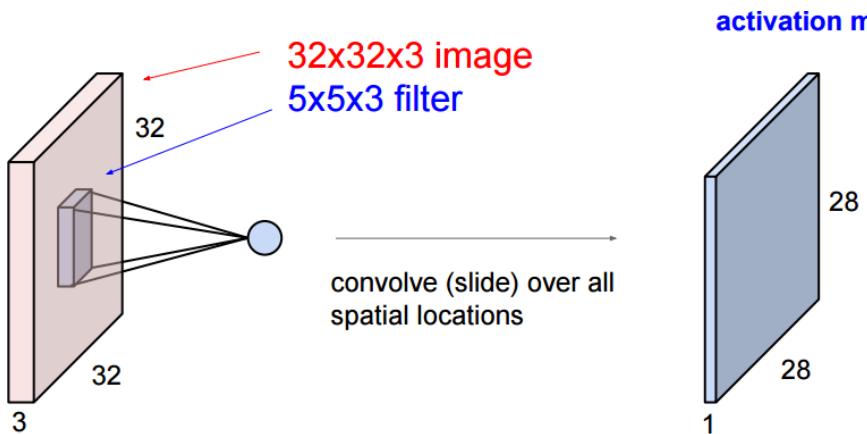


Convolutional Neural Network



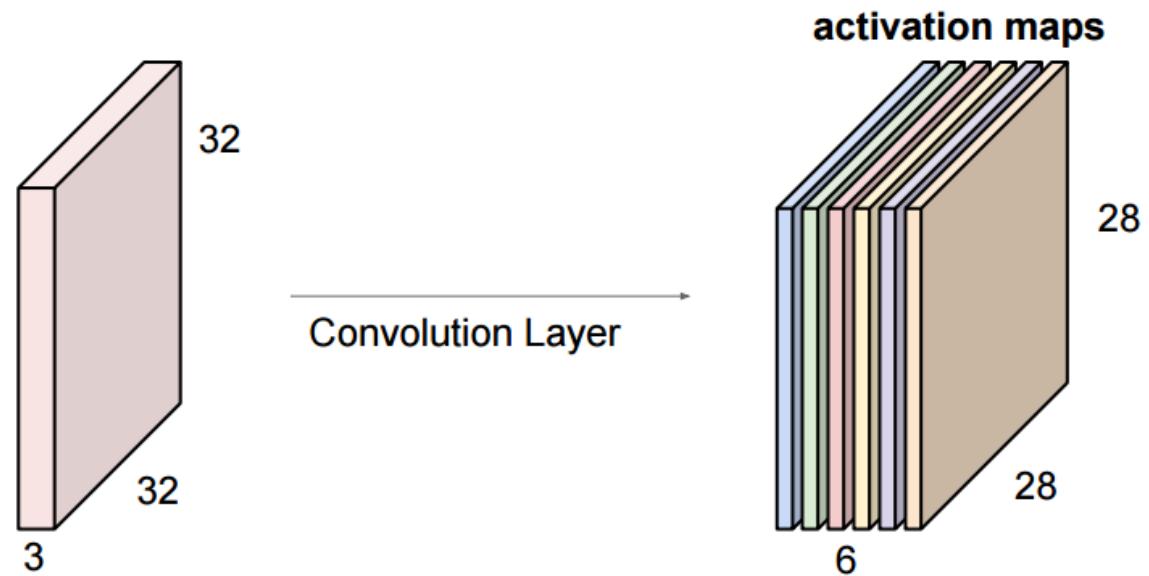
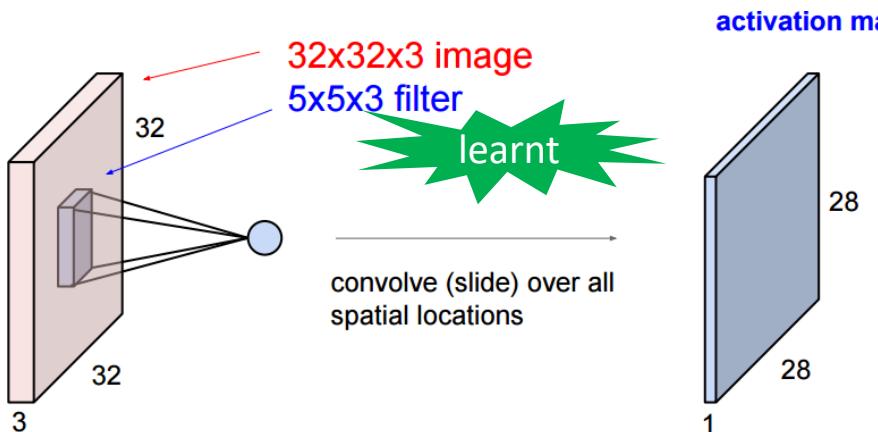


Convolutional Neural Network



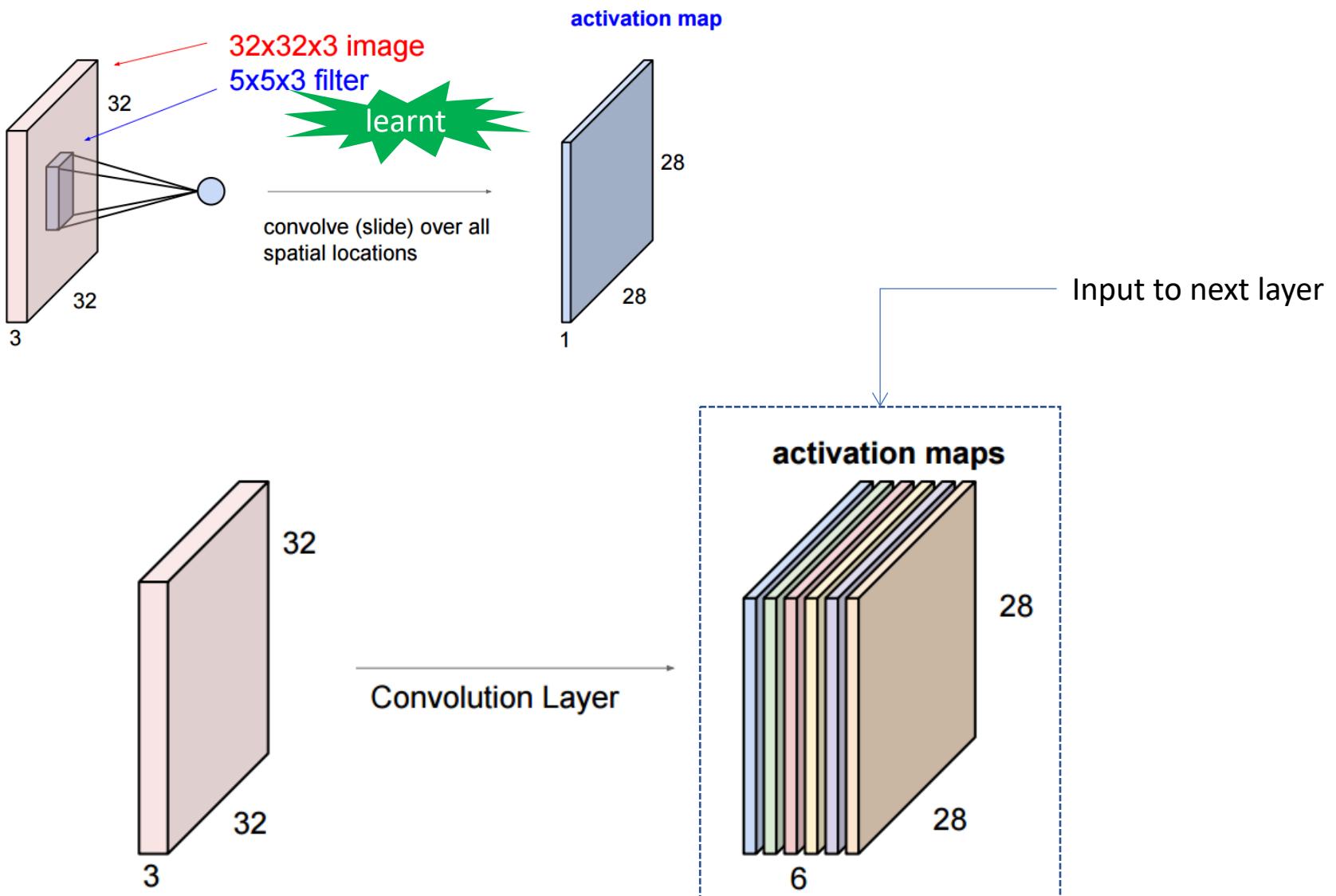


Convolution Layer



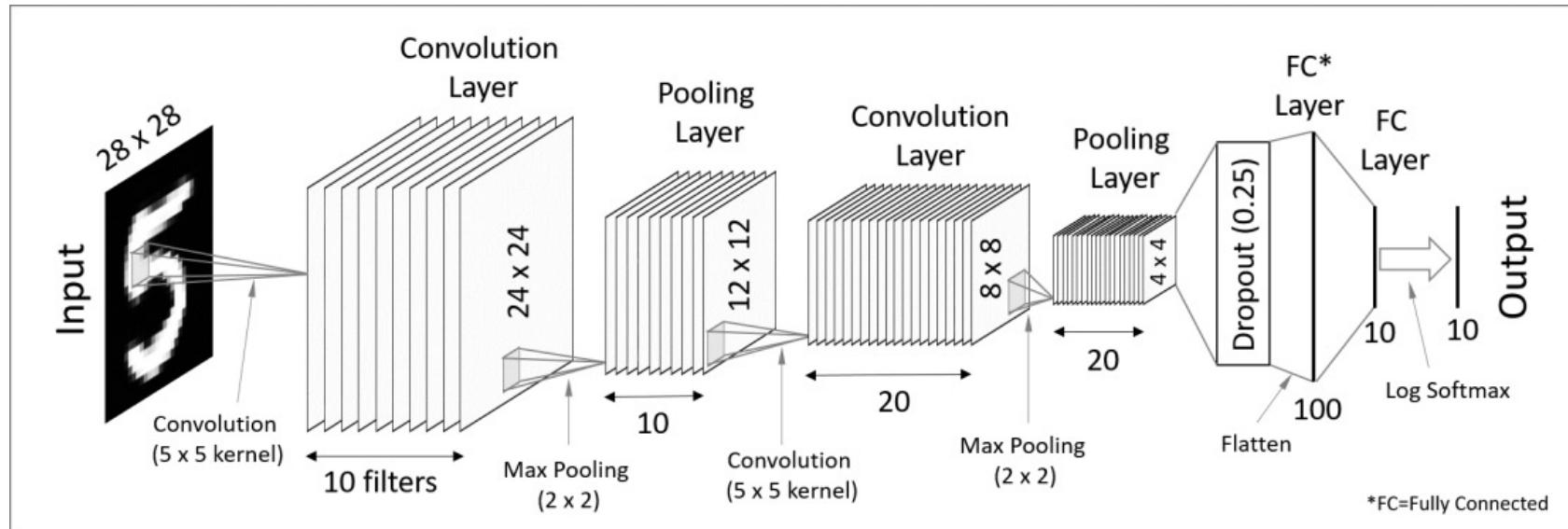


Convolution Layer



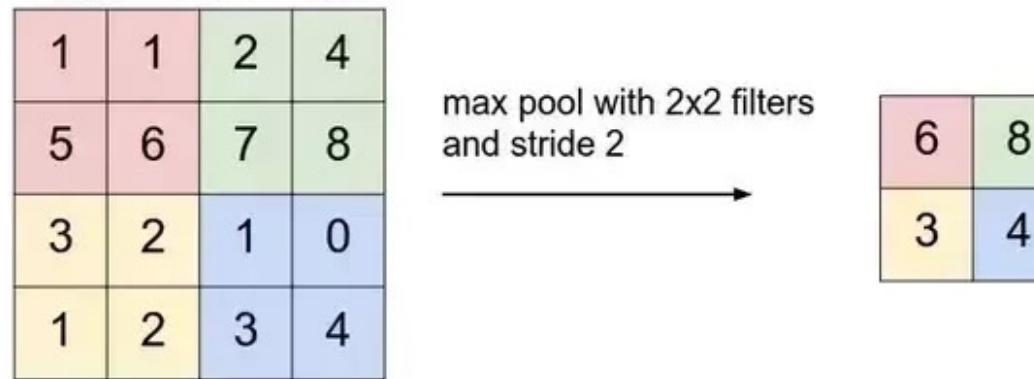


Convolutional Neural Network

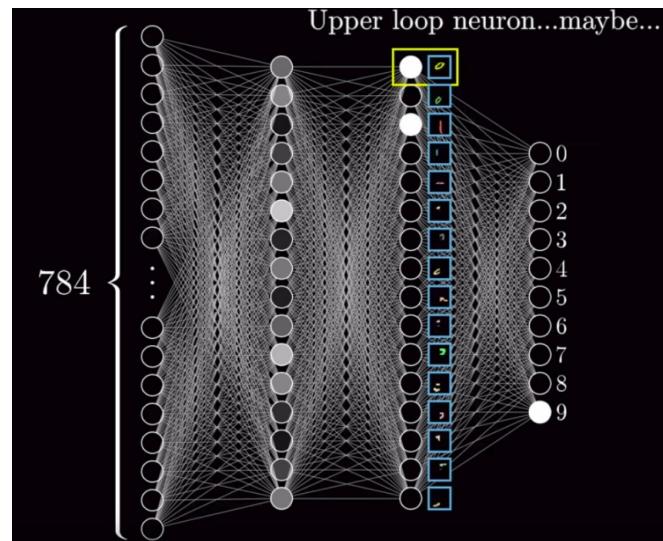




Pooling Layer



Fully-connected Layer





MML

Questions?





MML

ImageNet Challenge and AlexNet

The net that triggered the deep learning revolution

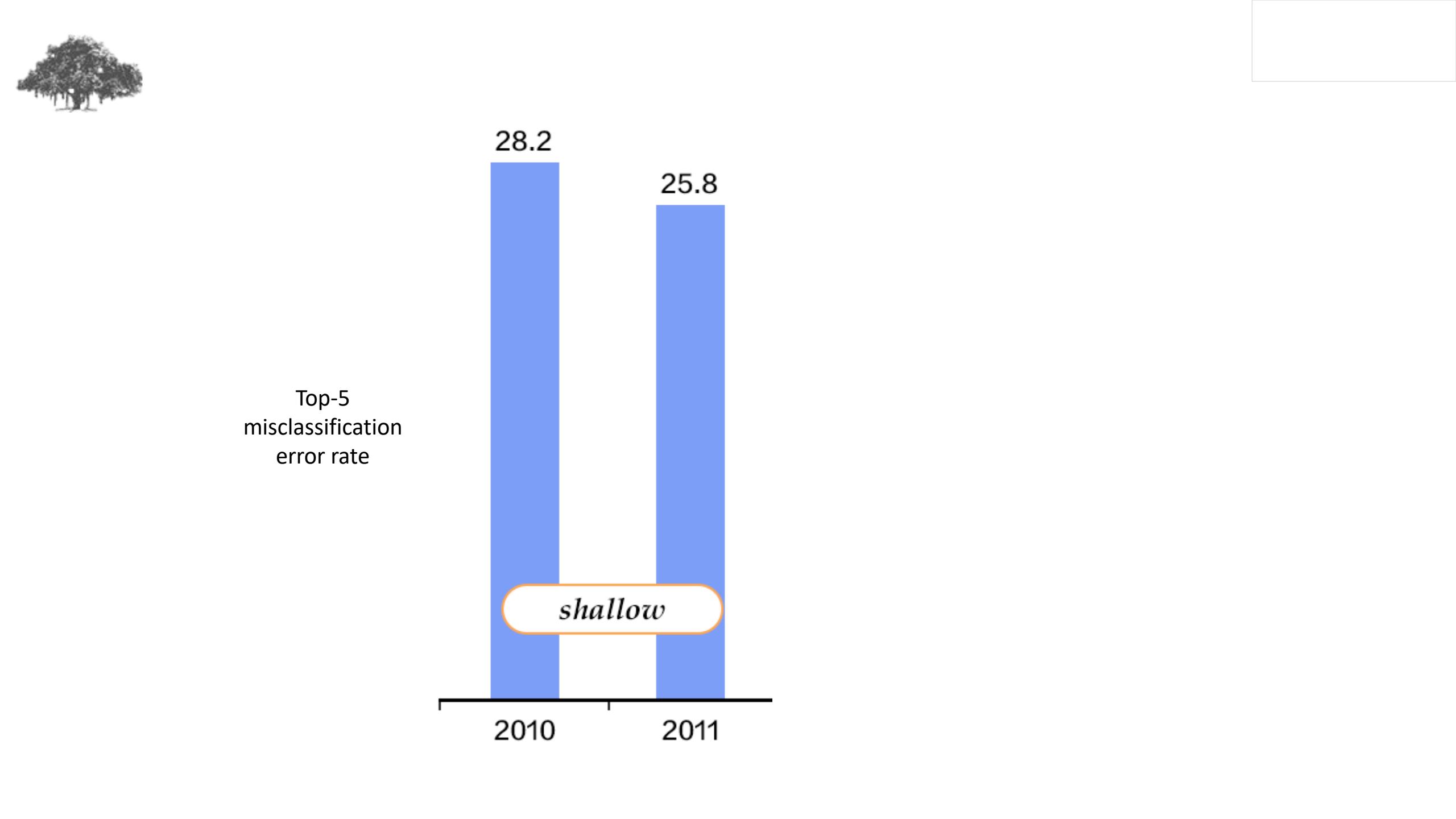


ImageNet Challenge

IMAGENET

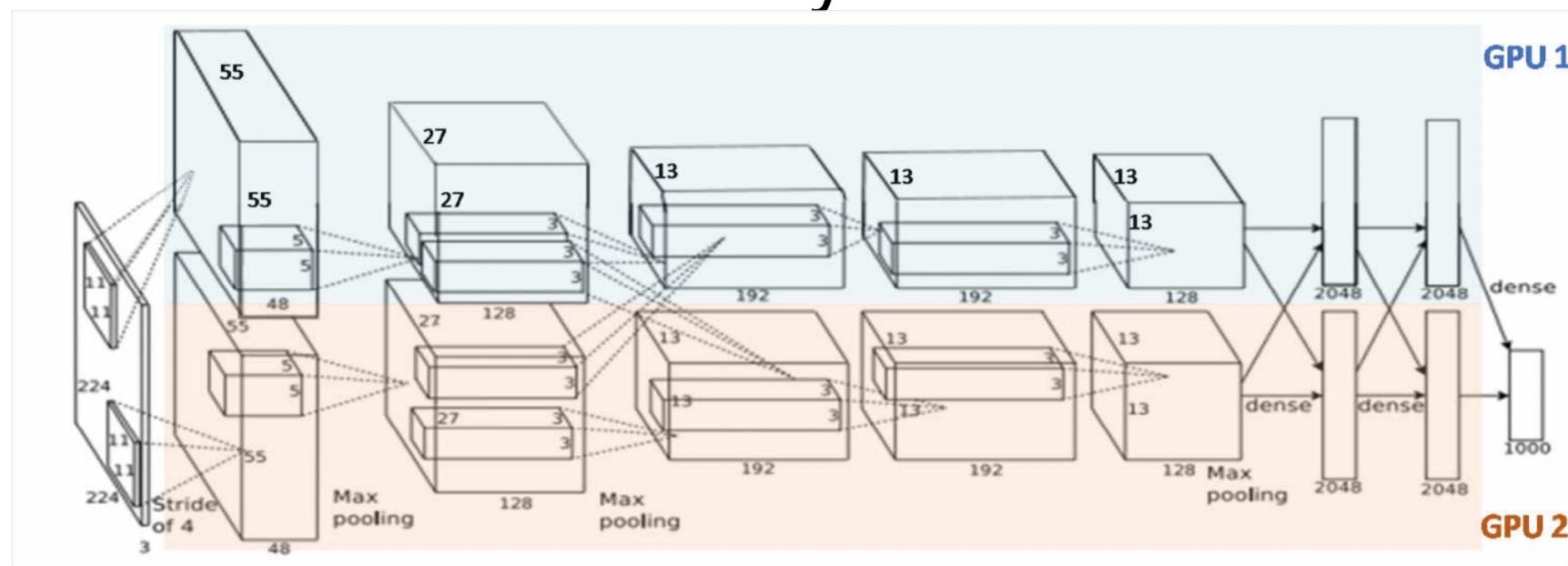
- 1,000 object classes (categories).
- Images:
 - 1.2 M train
 - 100k test.







Case Study: AlexNet

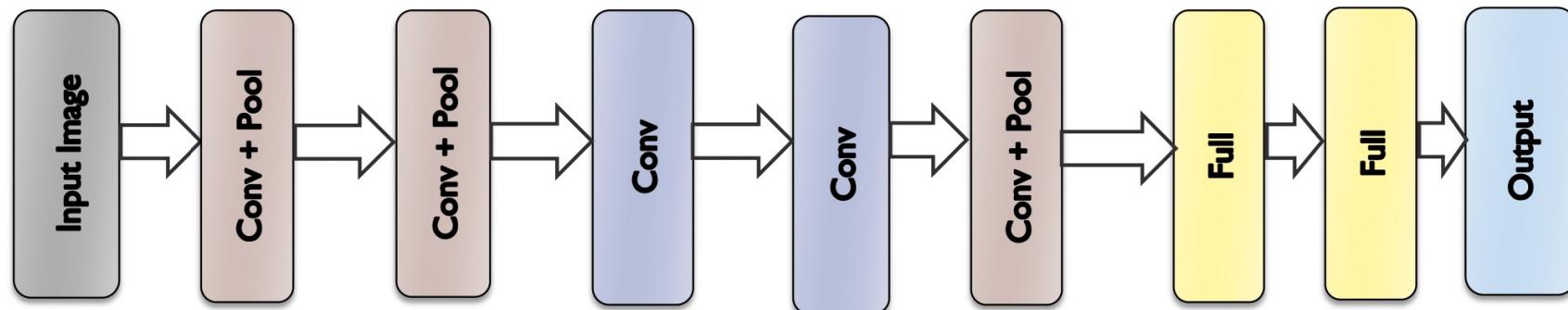
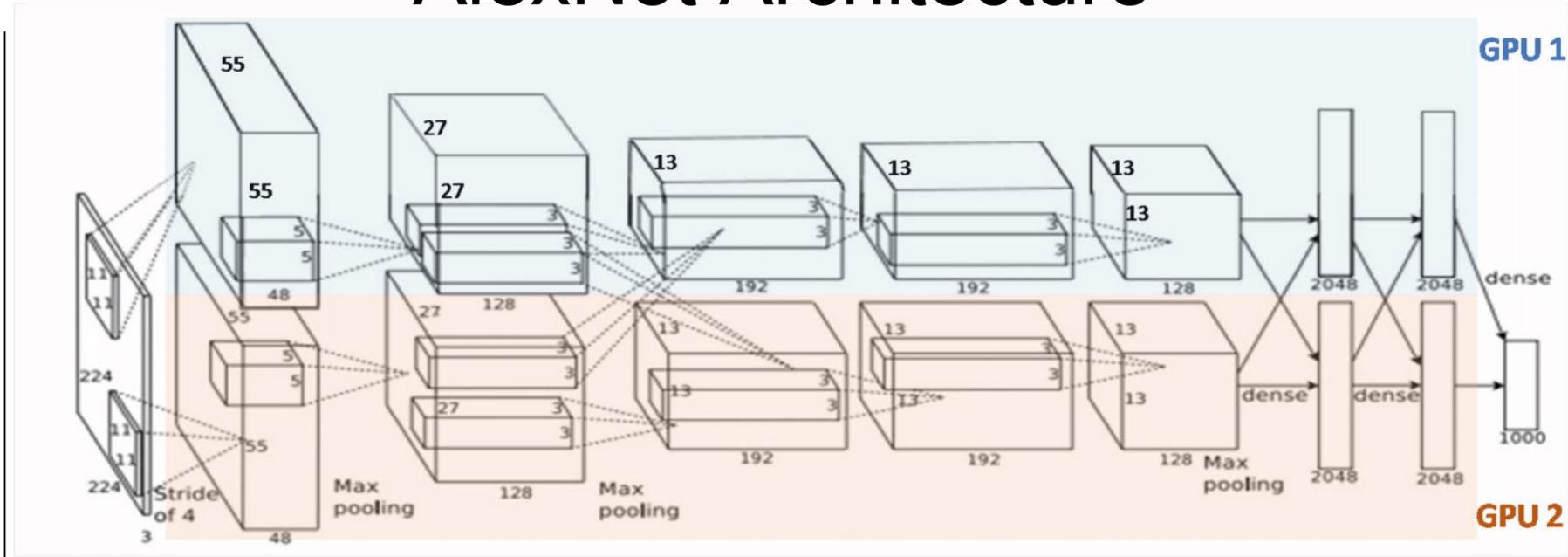


- Winner of ImageNet LSVRC-2010.
- Trained over 1.2M images using SGD with regularization.
- Deep architecture (60M parameters.)
- Optimized GPU implementation (cuda-convnet)

Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *NIPS* 2012.

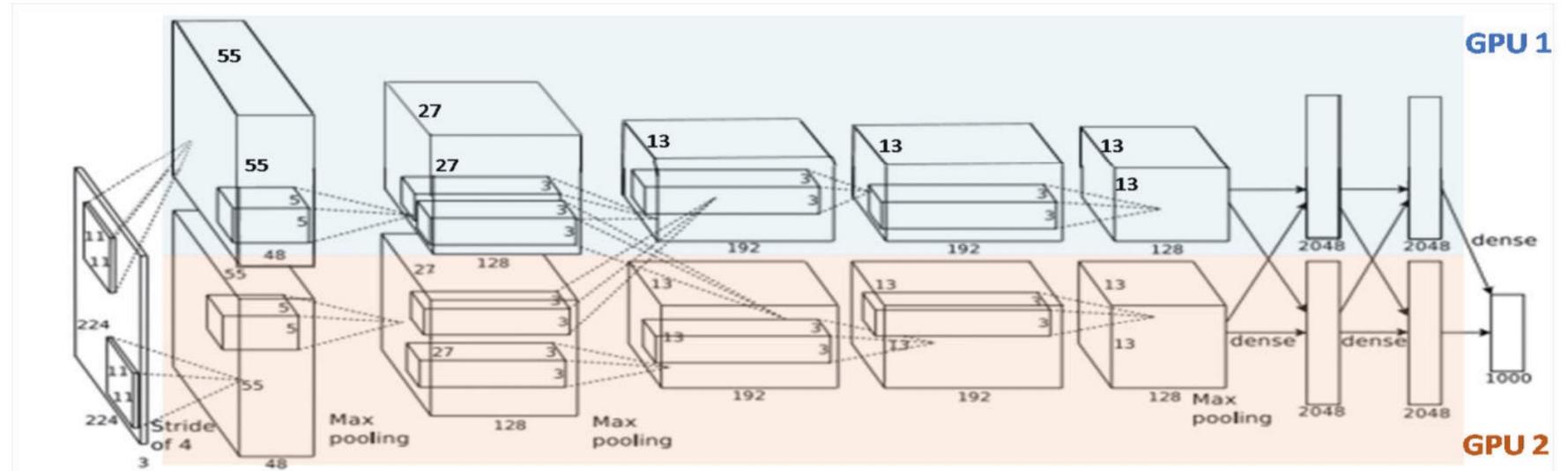


AlexNet Architecture





AlexNet Architecture

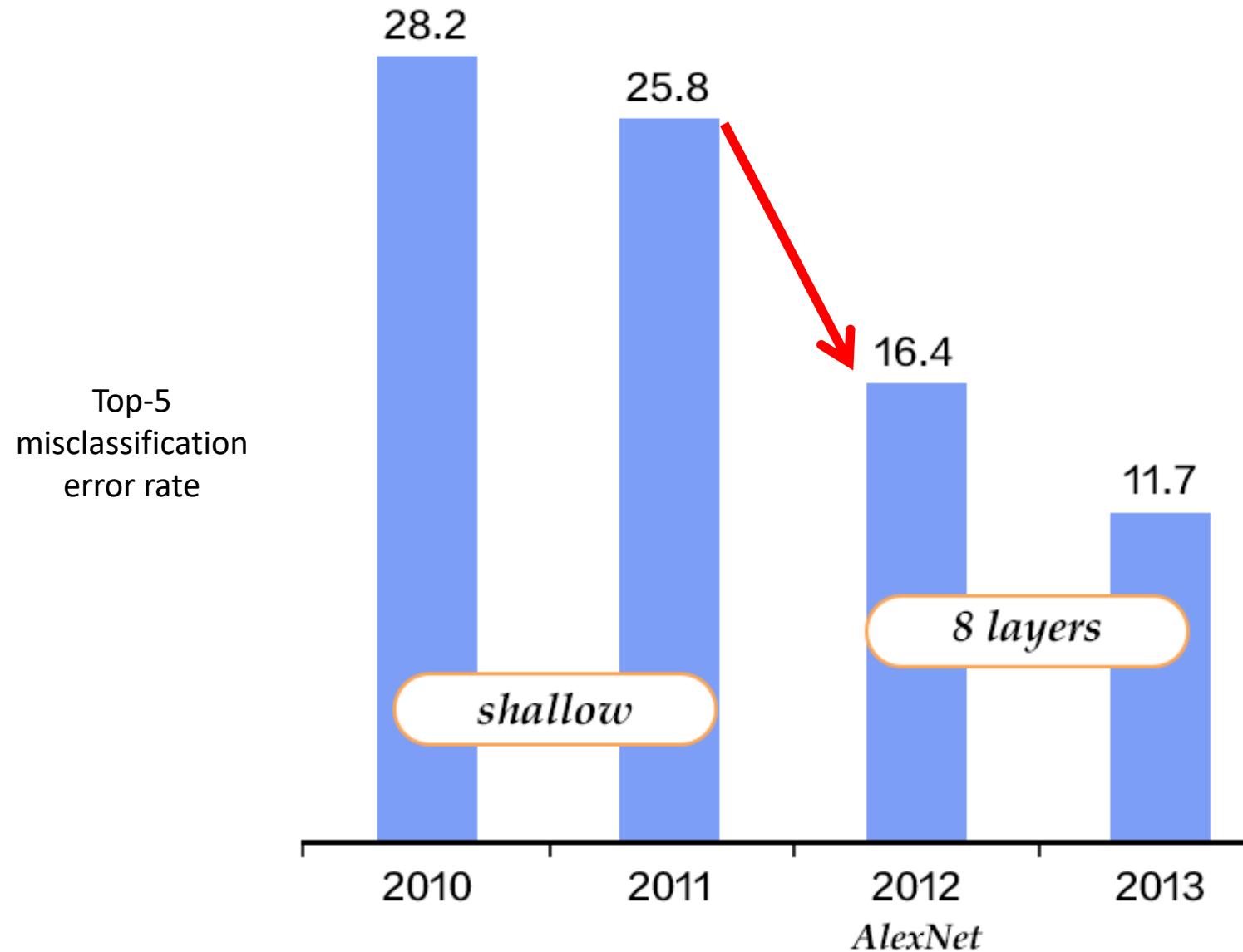


About 57 M parameters are in
the fully connected layers

Total

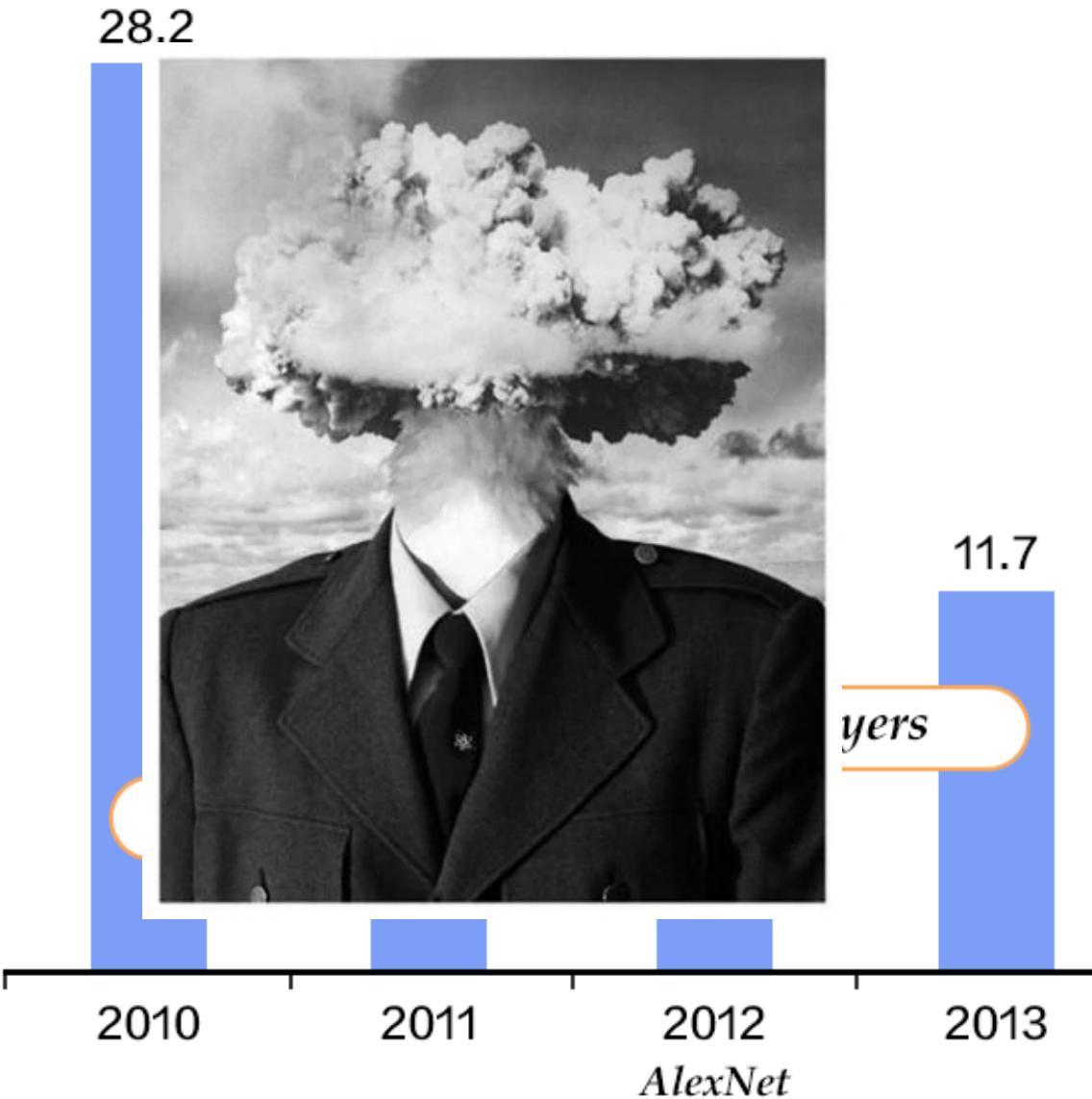
| Parameters : | $[(11 \times 11 \times 3) + 1] \times 96 = 35 \text{ K}$ | $[5 \times 5 \times 48] \times 256 = 307 \text{ K}$ | $[3 \times 3 \times 256] \times 384 = 884 \text{ K}$ | 663 K | 442 K | 37 M | 16 M | 4 M | 60 M |
|--------------|--|---|--|------------------------------------|------------------------------------|------|------|------|--------|
| Neurons : | 253,440 | $27 \times 27 \times 256 = 186,624$ | $13 \times 13 \times 384 = 64,896$ | $13 \times 13 \times 384 = 64,896$ | $13 \times 13 \times 256 = 43,264$ | 4096 | 4096 | 1000 | 0.63 M |

- Convolutional layers cumulatively contain about 90-95% of computation, only about 5% of the parameters
- Fully-connected layers contain about 95% of parameters.



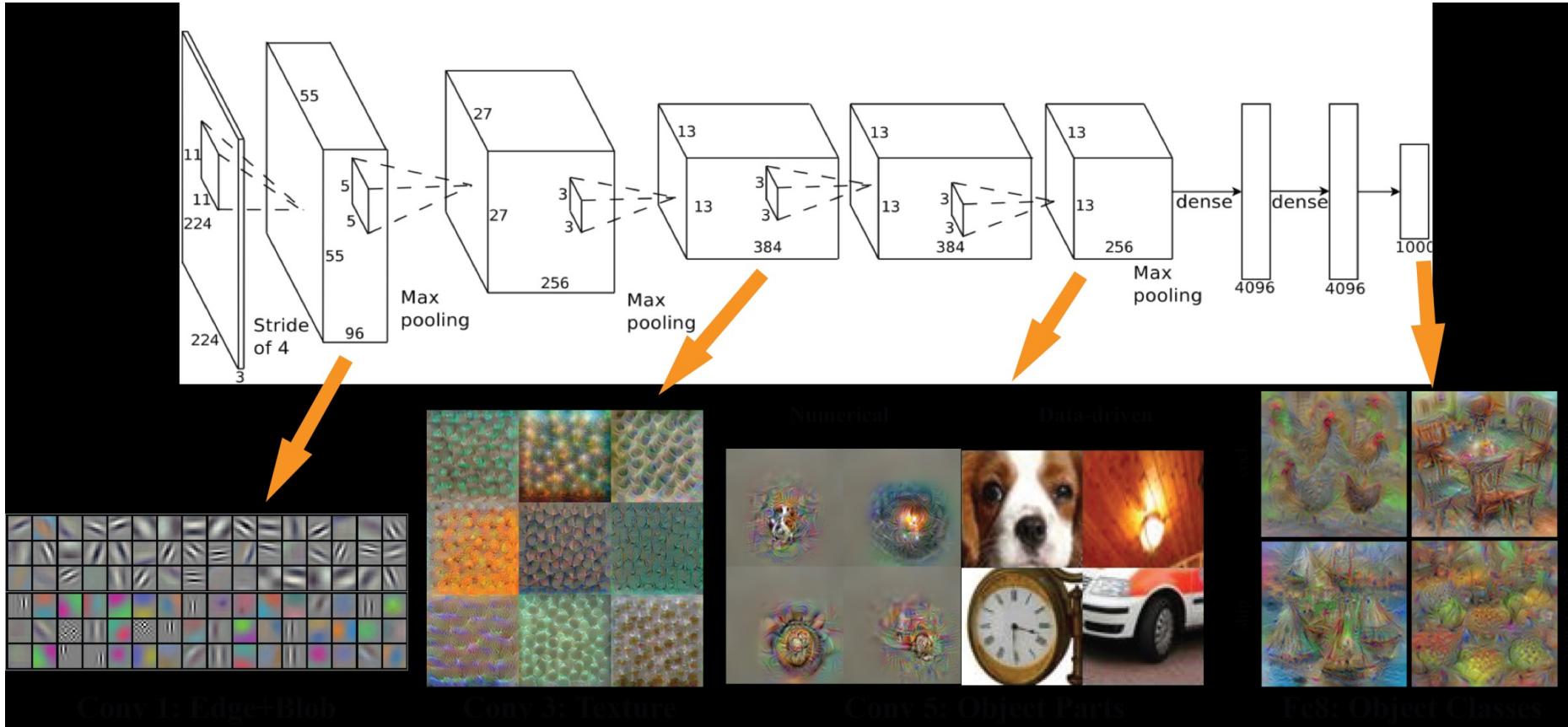


Top-5
misclassification
error rate





Filters learnt by AlexNet





MML

Questions?





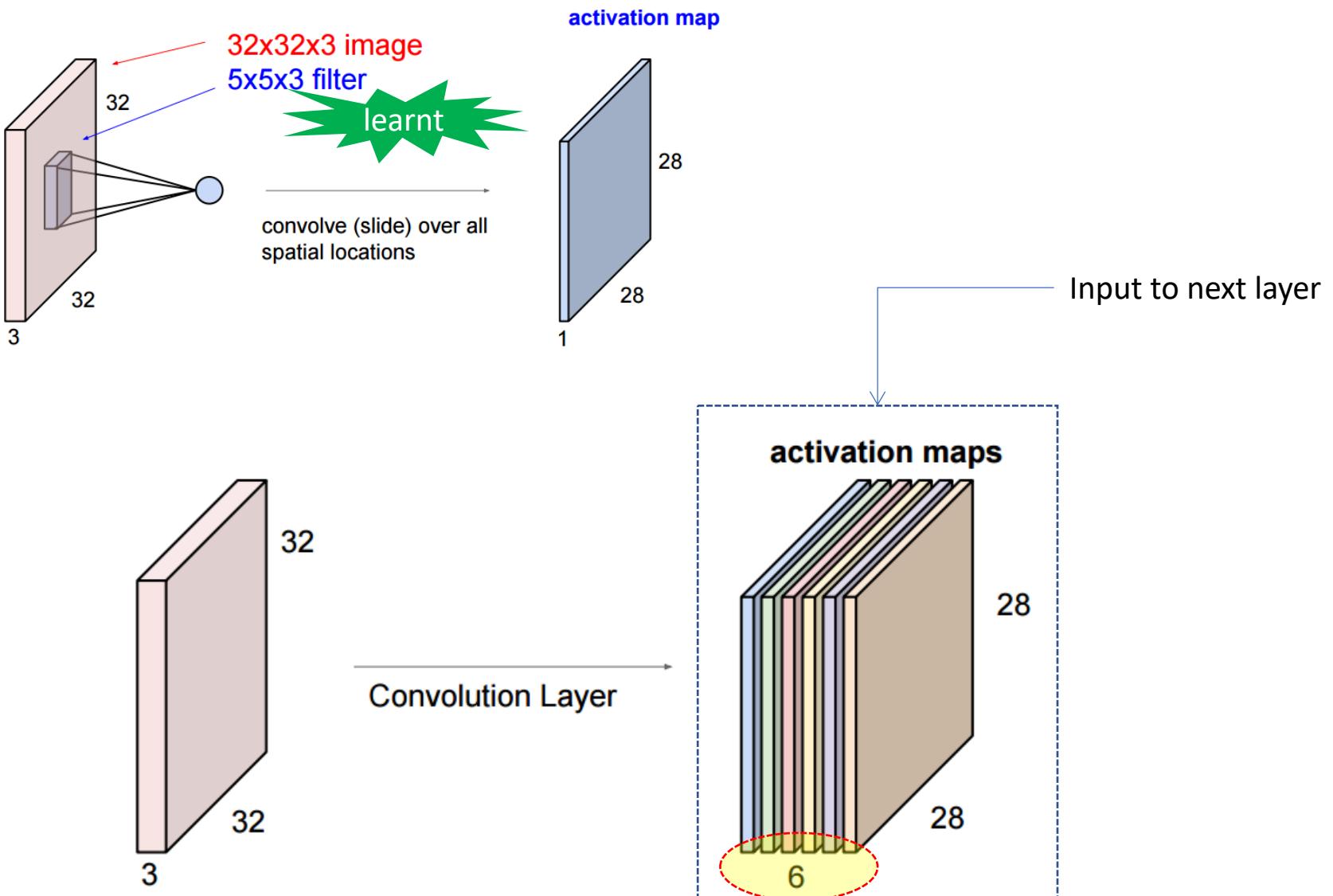
MML

Popular CNN architectures

Improvements and Extensions

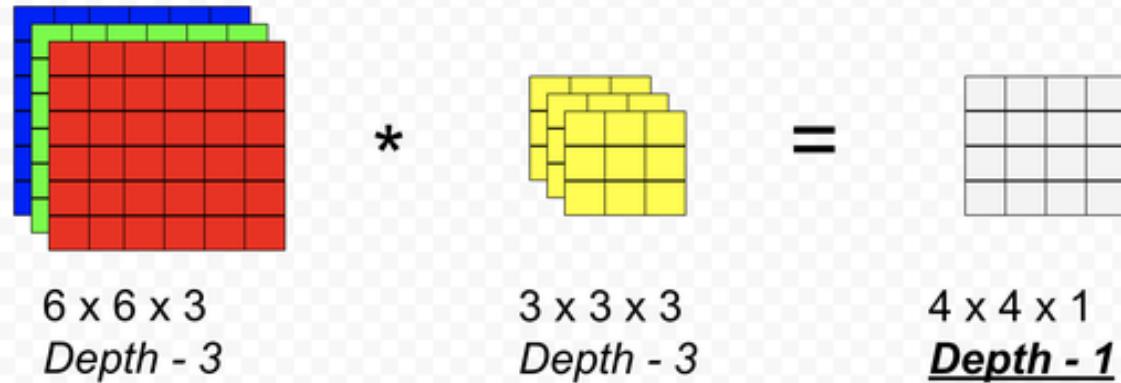


Convolution Layer





Convolution Layer





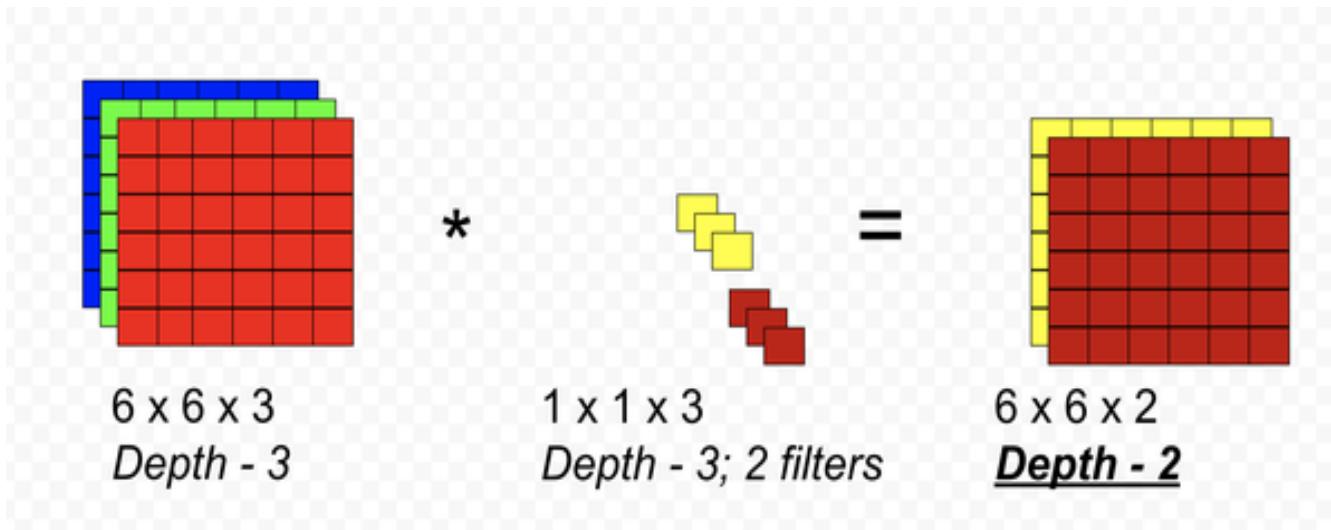
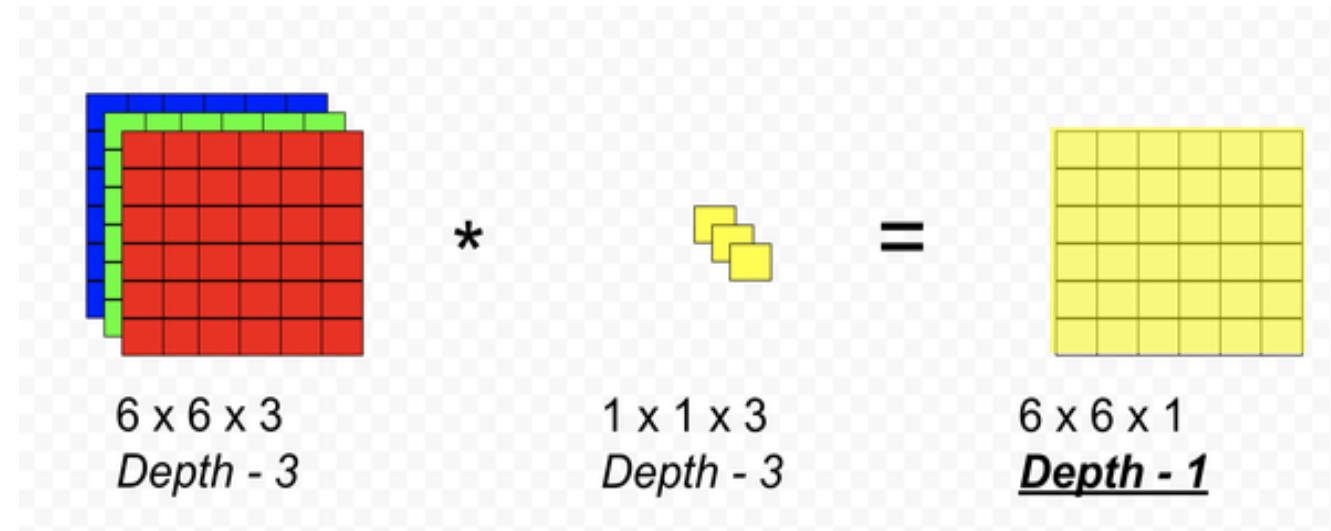
1x1 (Bottleneck) Convolution Layer

$$\begin{array}{ccc} \text{Diagram of a } 6 \times 6 \times 3 \text{ tensor with depth 3, colored blue, green, and red.} & * & \text{Diagram of a } 3 \times 3 \times 3 \text{ tensor with depth 3, colored yellow.} \\ 6 \times 6 \times 3 \\ Depth - 3 & & 3 \times 3 \times 3 \\ Depth - 3 \end{array} = \begin{array}{c} \text{Diagram of a } 4 \times 4 \times 1 \text{ tensor with depth 1, colored grey.} \\ 4 \times 4 \times 1 \\ \underline{\text{Depth - 1}} \end{array}$$

$$\begin{array}{ccc} \text{Diagram of a } 6 \times 6 \times 3 \text{ tensor with depth 3, colored blue, green, and red.} & * & \text{Diagram of a } 1 \times 1 \times 3 \text{ tensor with depth 3, colored yellow.} \\ 6 \times 6 \times 3 \\ Depth - 3 & & 1 \times 1 \times 3 \\ Depth - 3 \end{array} = \begin{array}{c} \text{Diagram of a } 6 \times 6 \times 1 \text{ tensor with depth 1, colored grey.} \\ 6 \times 6 \times 1 \\ \underline{\text{Depth - 1}} \end{array}$$



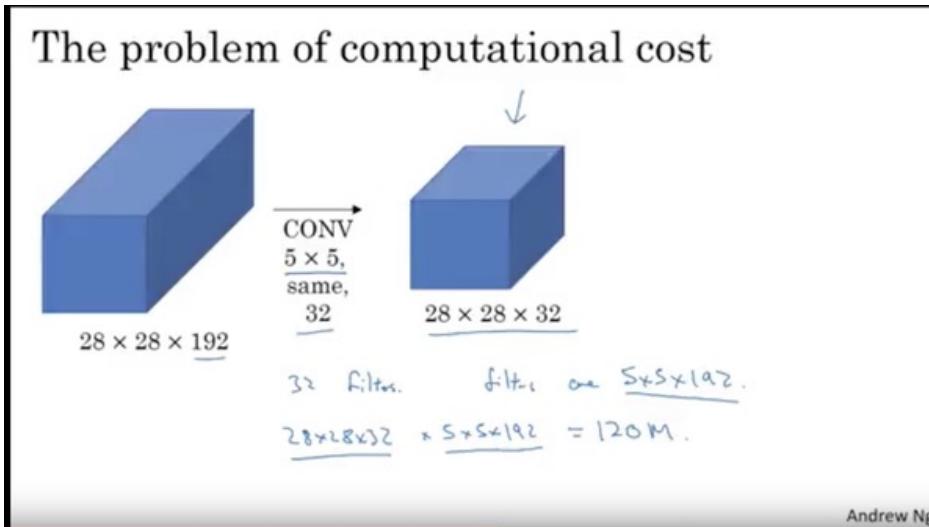
1x1 (Bottleneck) Convolution Layer





1x1 (Bottleneck) Convolution Layer

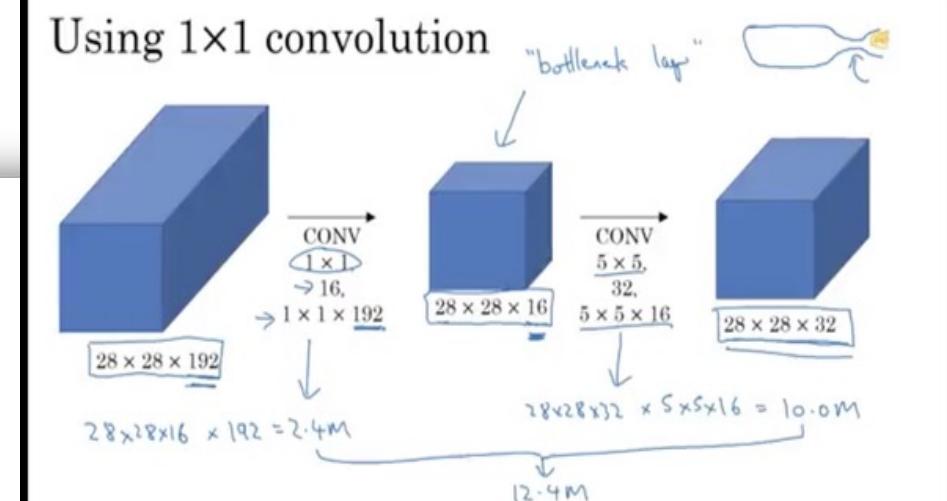
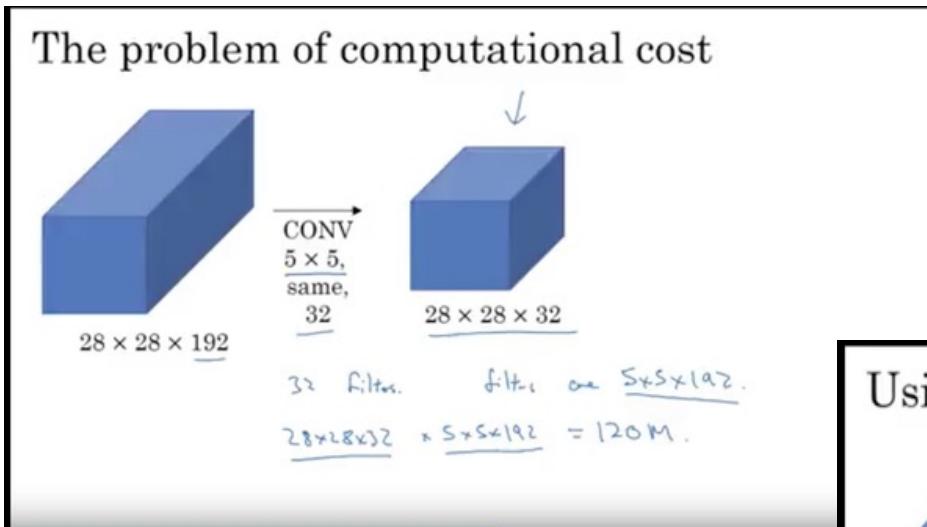
- Enable dimensionality reduction (computational savings)
- Add more non-linearity to representation space





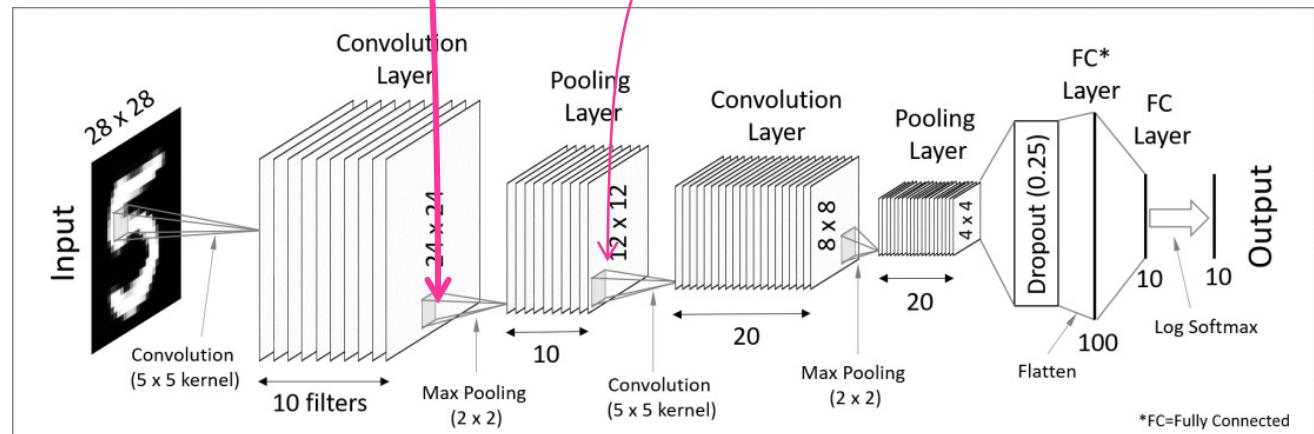
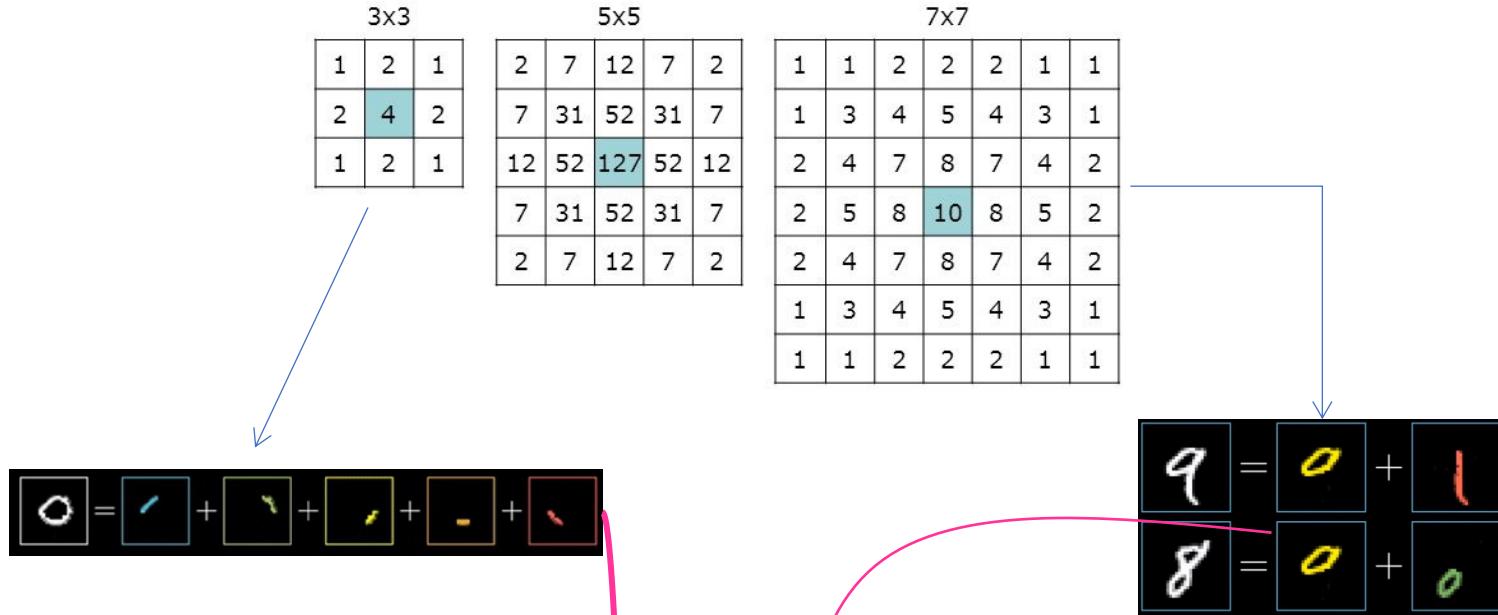
1x1 (Bottleneck) Convolution Layer

- Enable dimensionality reduction (computational savings)
- Add more non-linearity to representation space



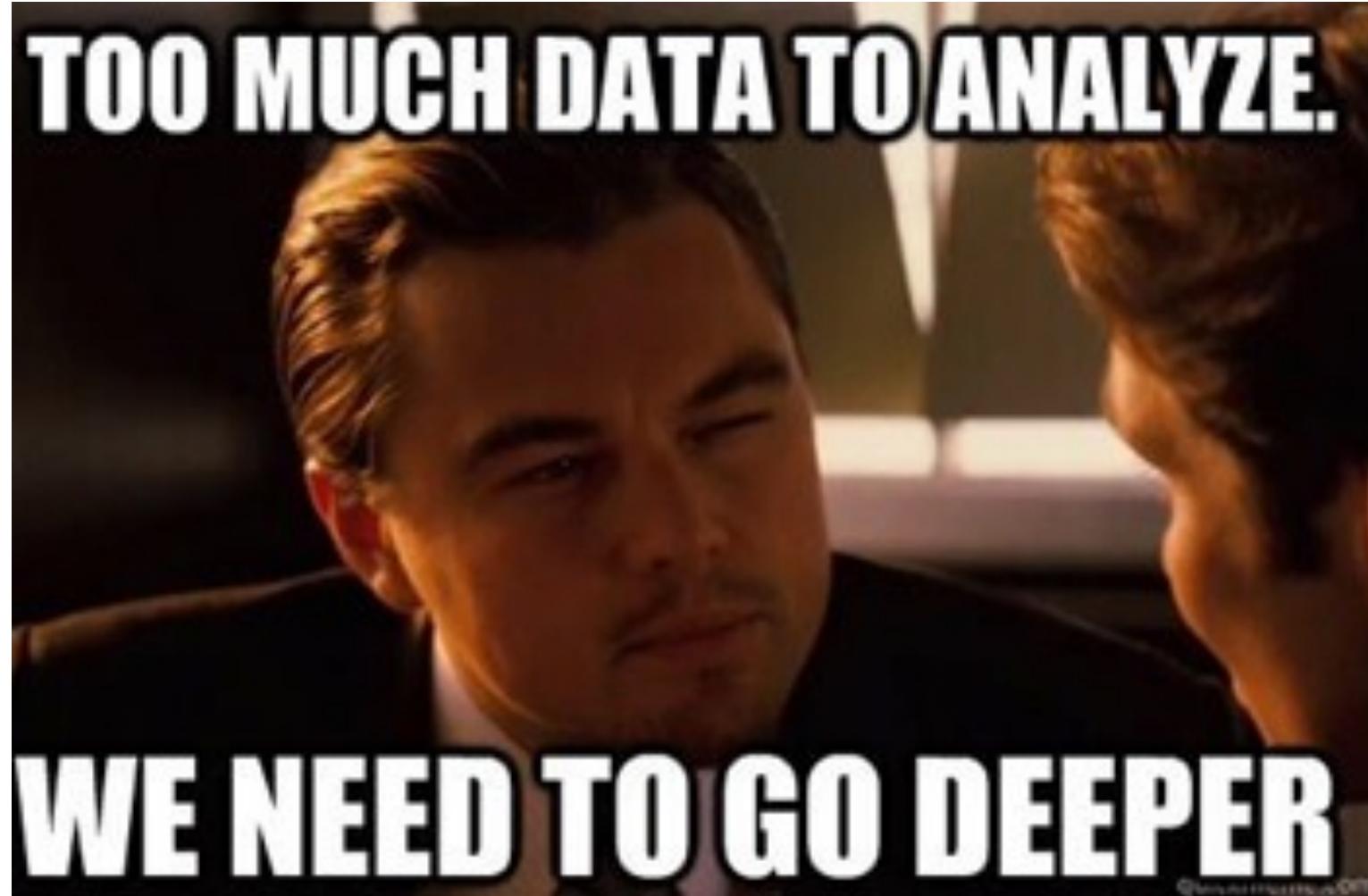


Committing to a single filter size at each layer seems somewhat rigid !



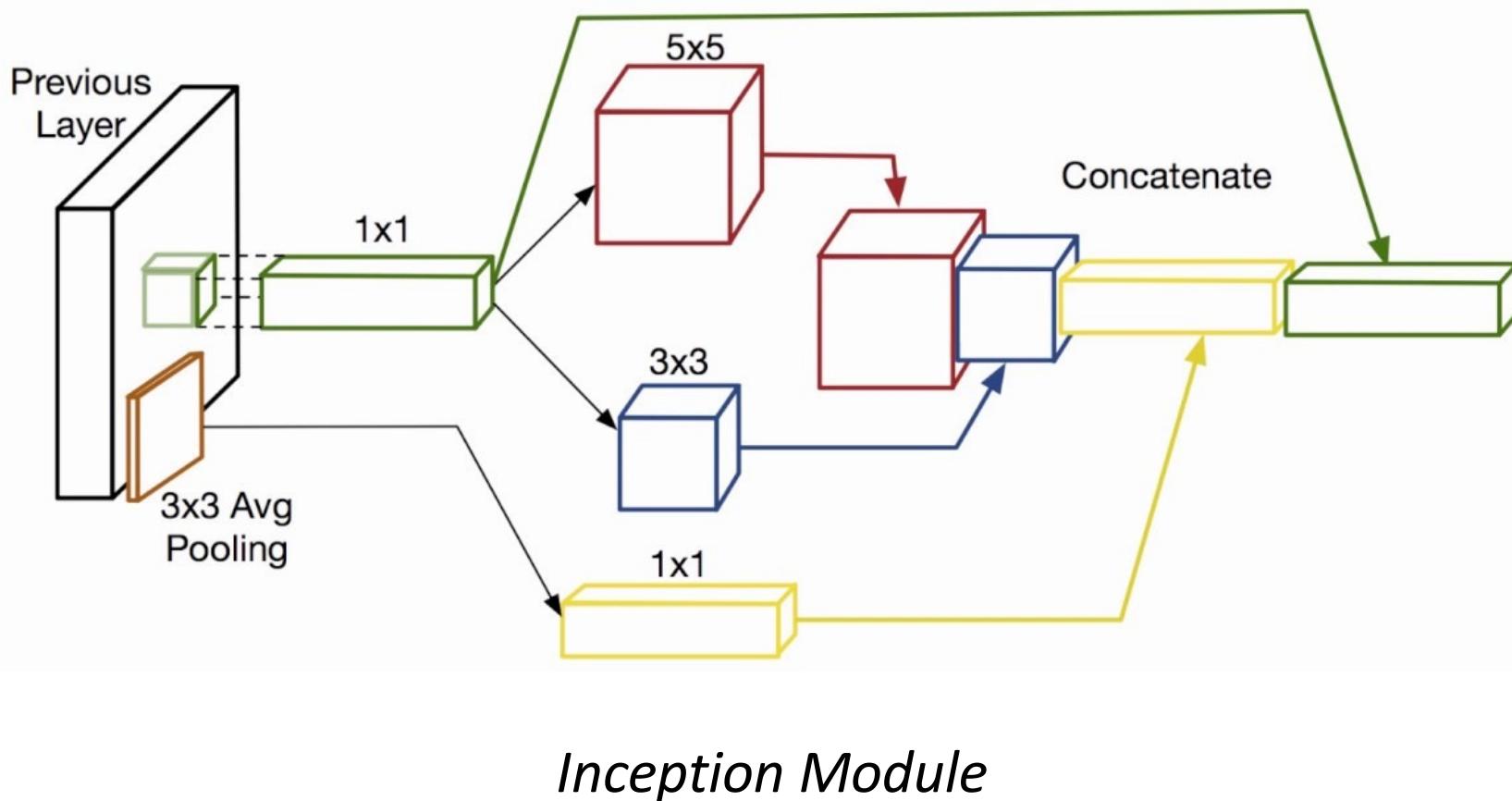


Why not have multiple levels (scales)
of filtering ?





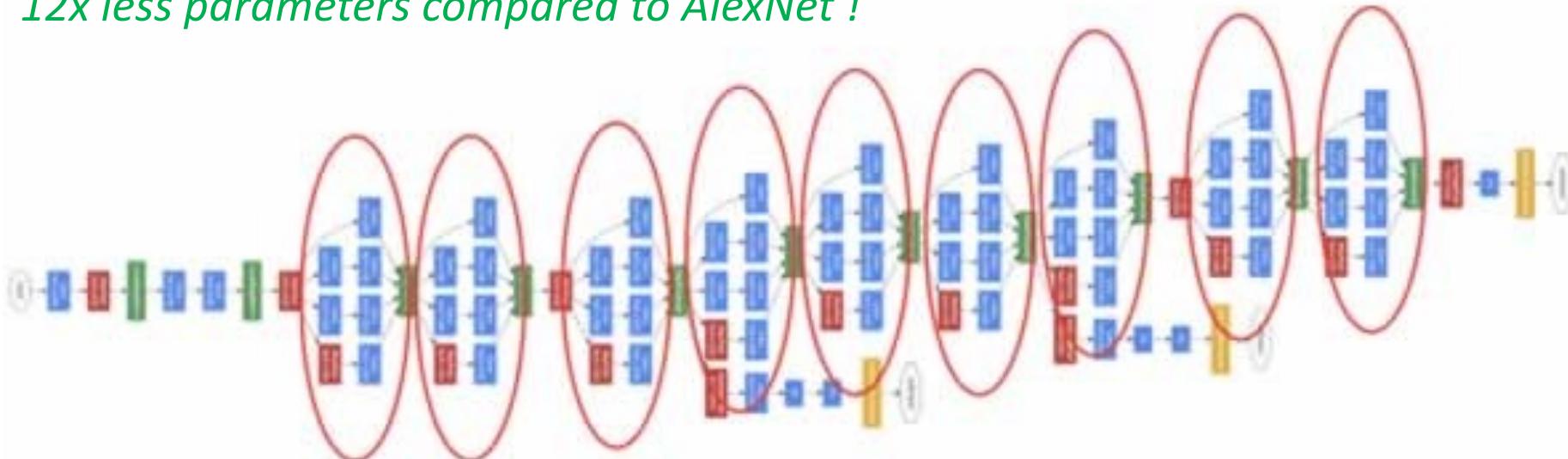
Why not have multiple levels (scales) of filtering ?





GoogLeNet

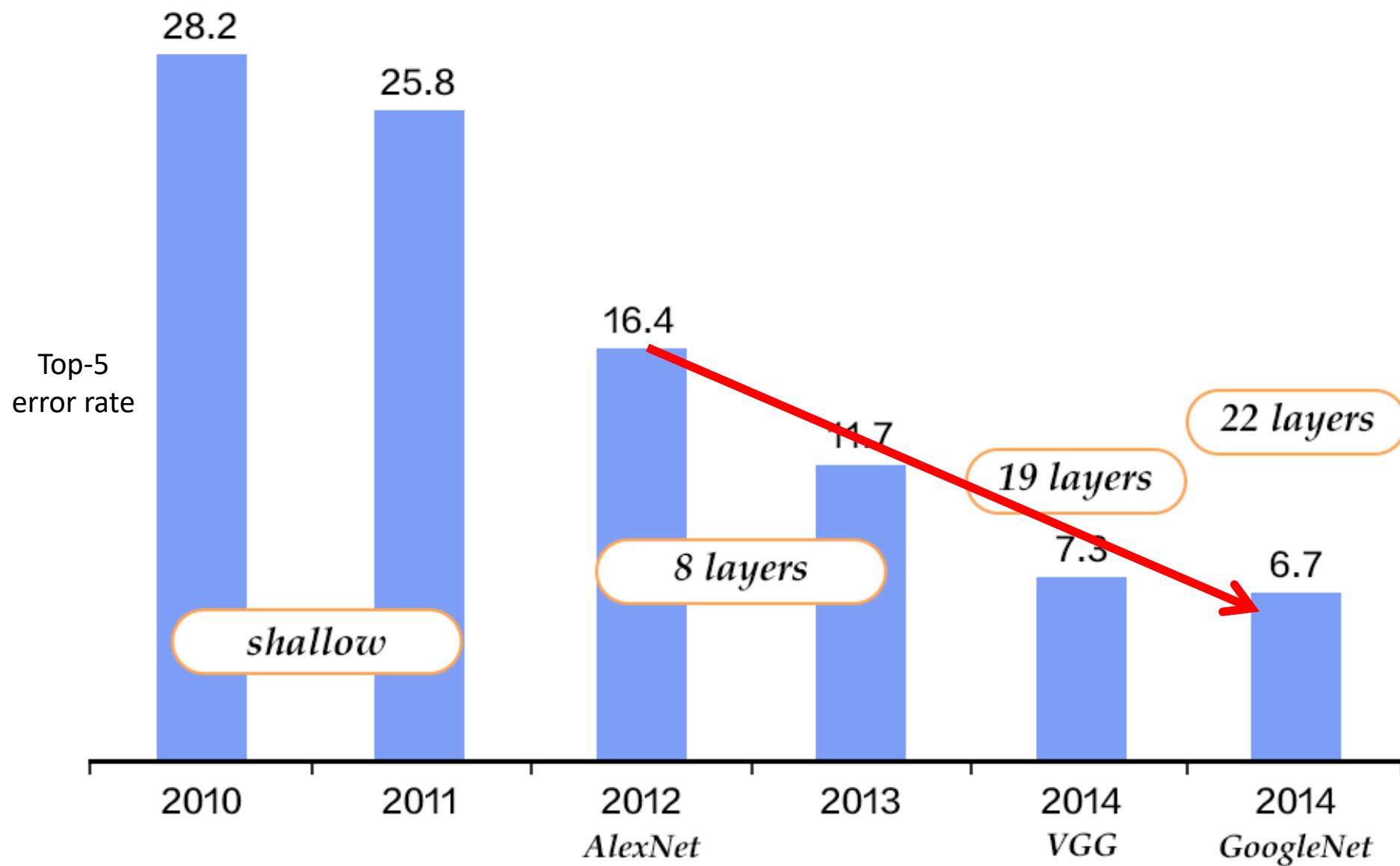
12x less parameters compared to AlexNet !



9 Inception modules

Network in a network in a network...

Convolution
Pooling
Softmax
Other





How deep can the network get ?

- 50 layers ? 100 ? 1000 ?

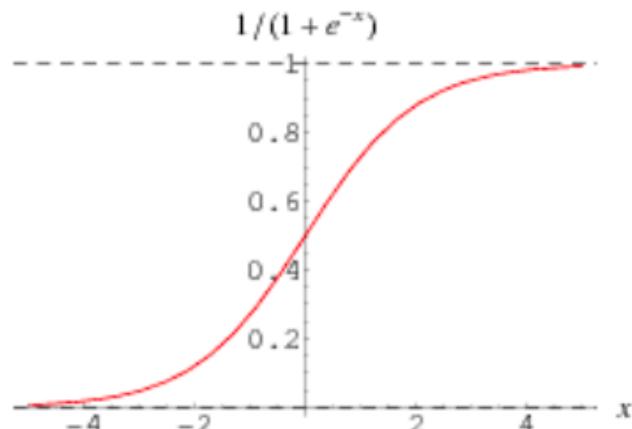


How deep can the network get ?

- 50 layers ? 100 ? 1000 ?
- Issue: Vanishing gradients

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

Sigmoid



Plot of Sigmoid function

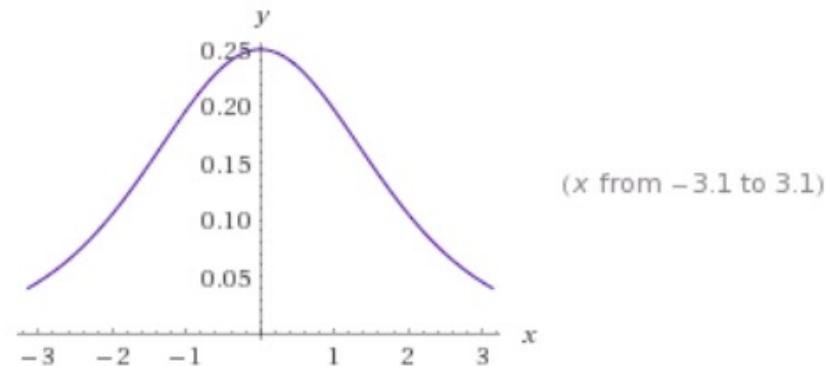


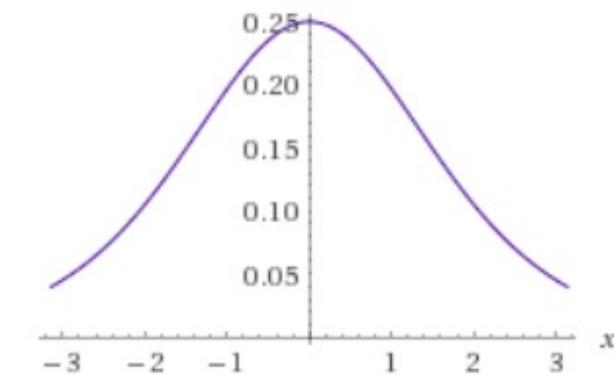
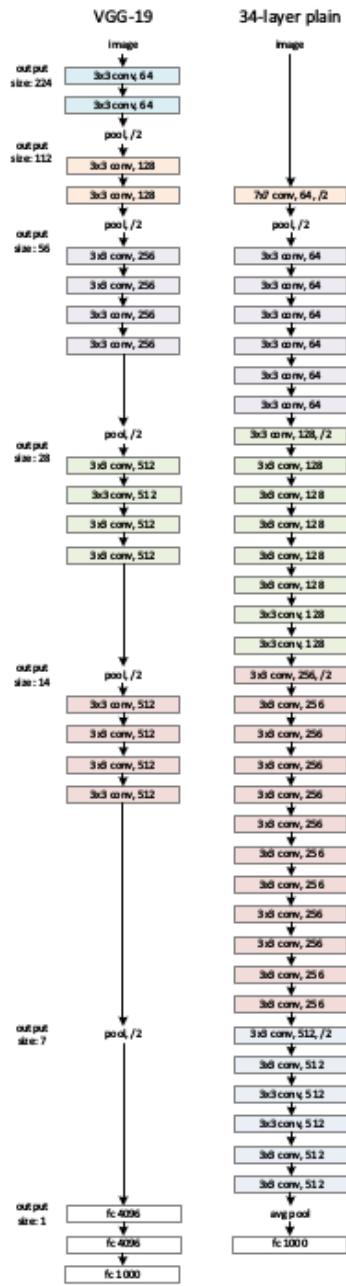
Vanishing gradients

$$\sigma'(x) = \frac{1}{(1 + e^{-x})} \left(1 - \frac{1}{(1 + e^{-x})}\right)$$

$$\sigma'(x) = \sigma(1 - \sigma)$$

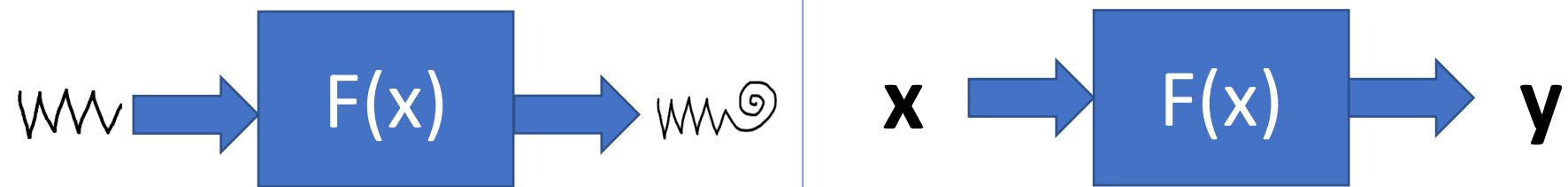
Derivative of sigmoid can be represented as





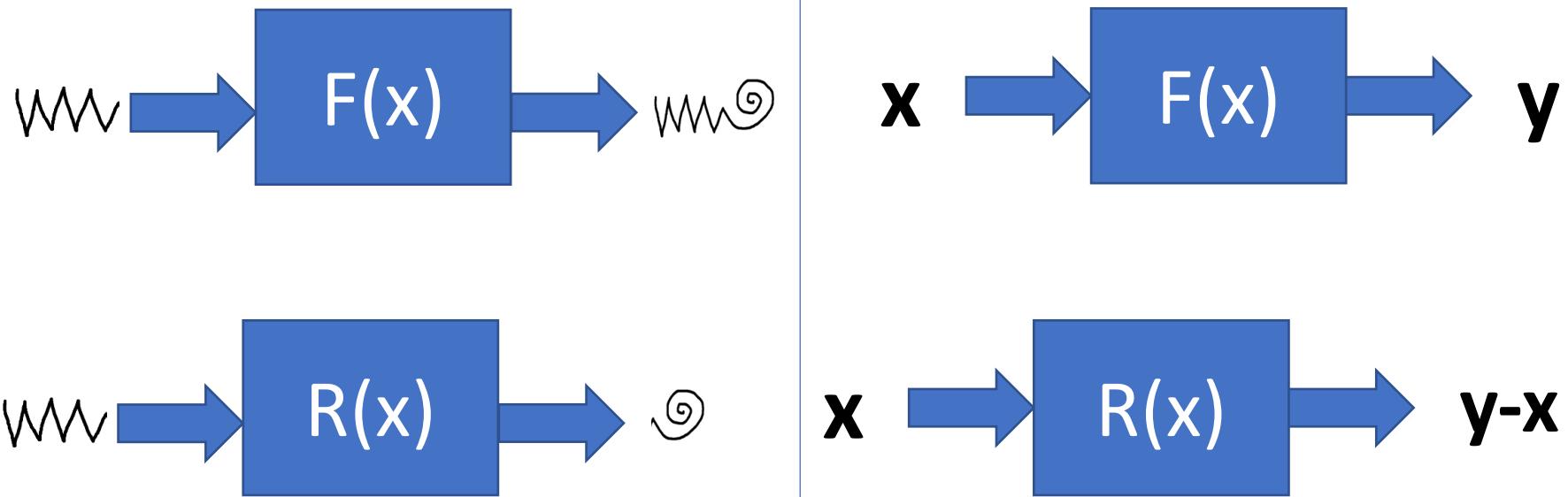


Residual Networks – A crayon analogy



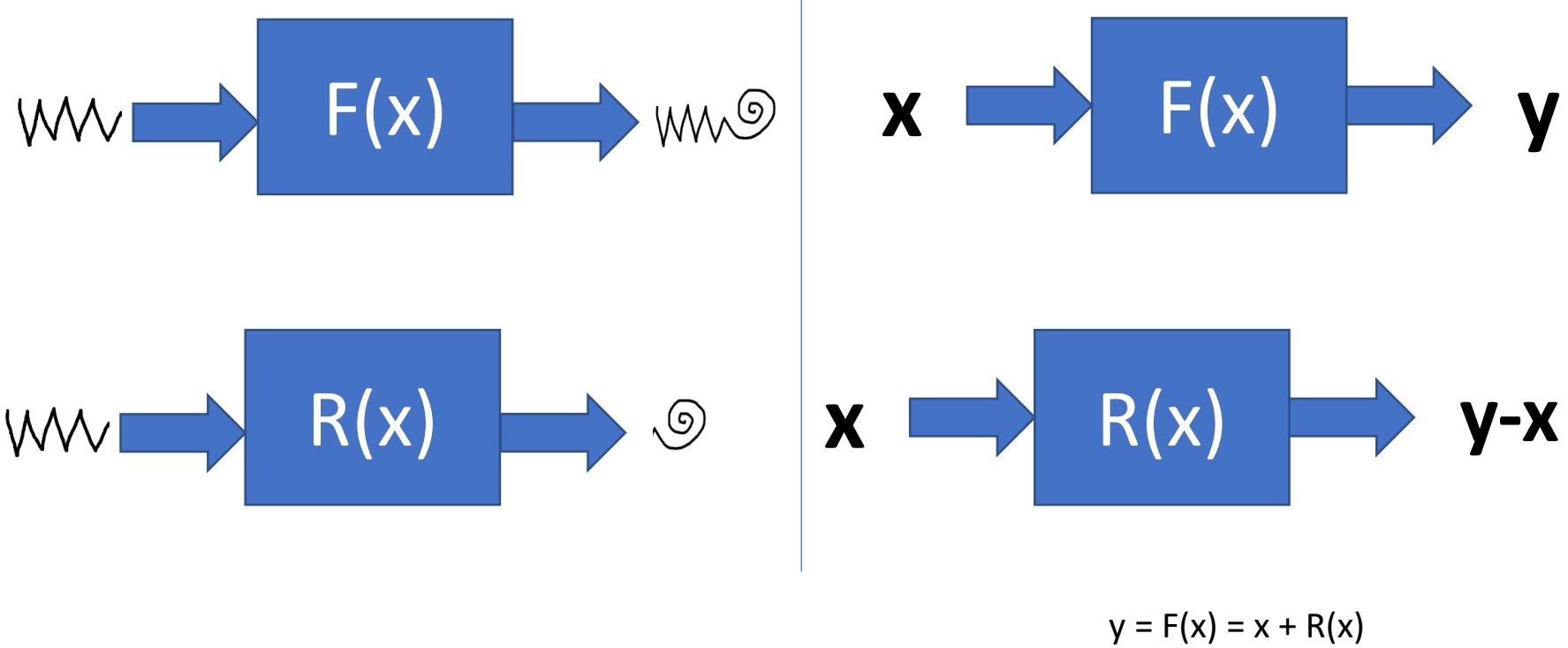


Residual Networks – A crayon analogy



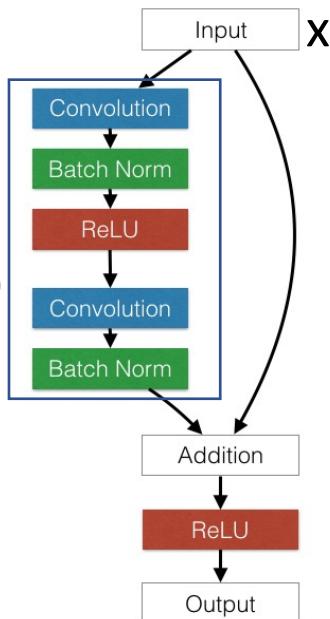
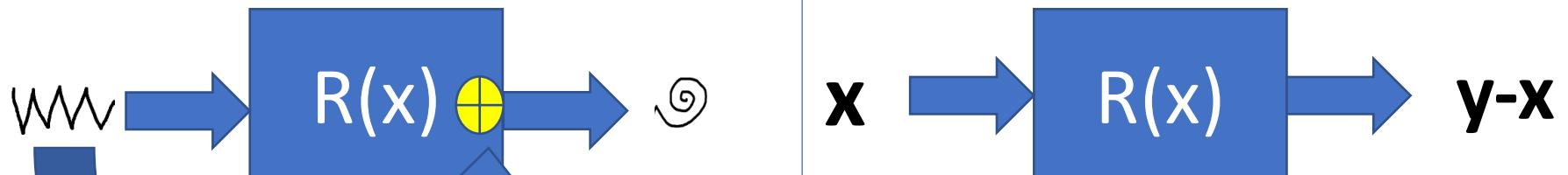
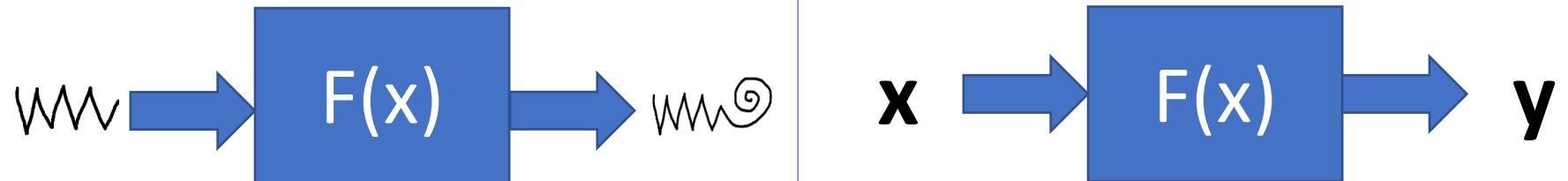


Residual Networks – A crayon analogy





Residual Networks – A crayon analogy

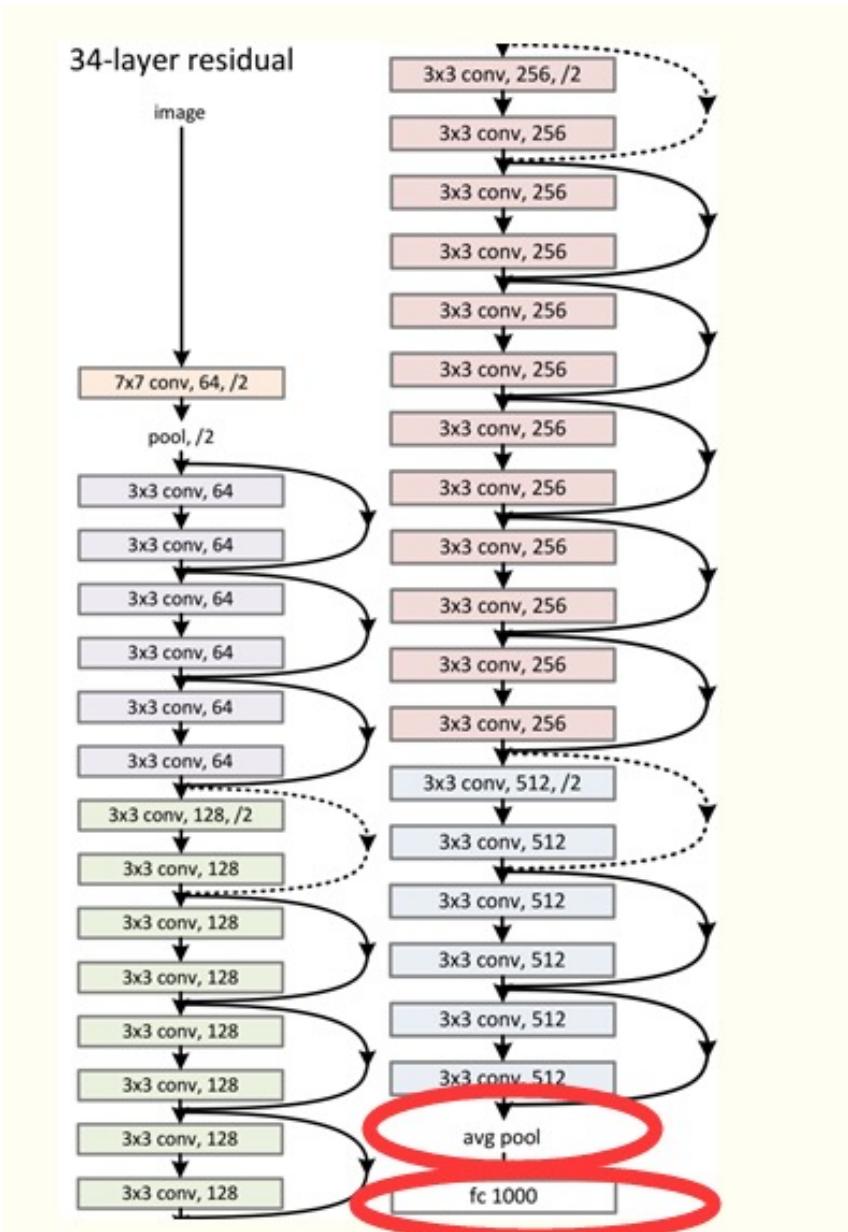


Architecture of a residual block

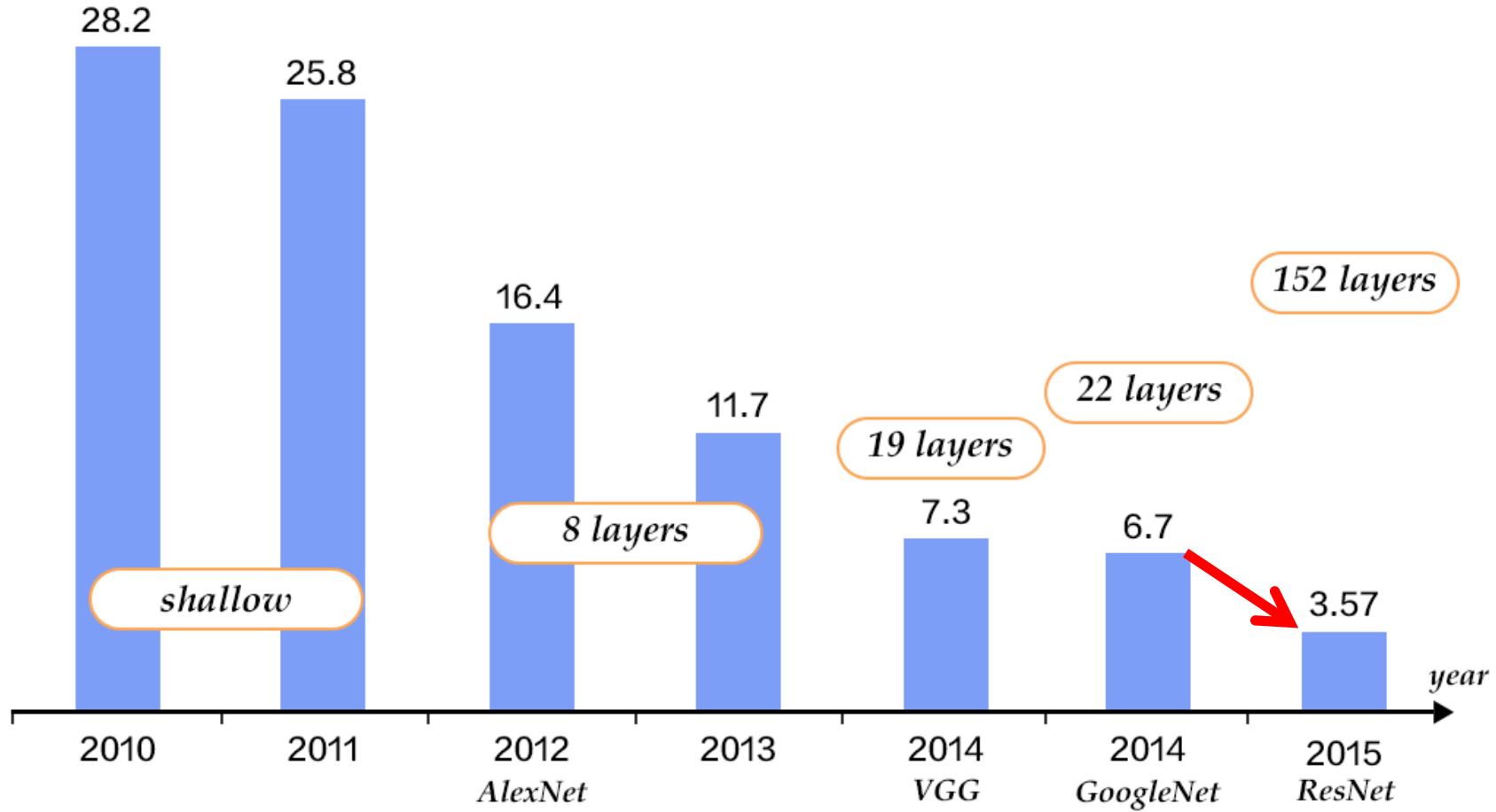
$$F(x) = x + R(x)$$



Residual Networks



Gradients can be supplied to shallower layers directly if needed (via skips)





MML

Questions?





MML

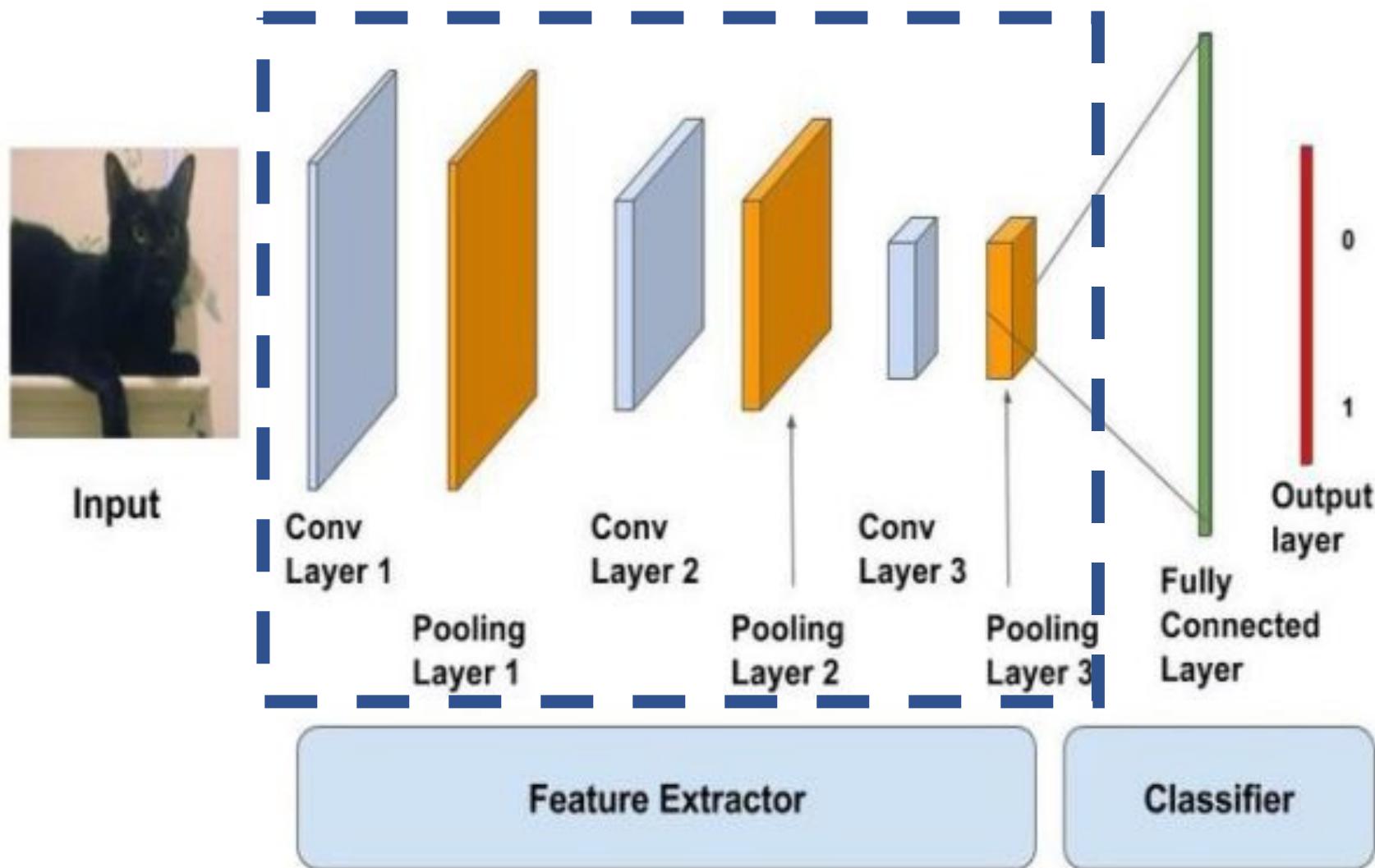
Landscape of CNNs



Applications and Architectures

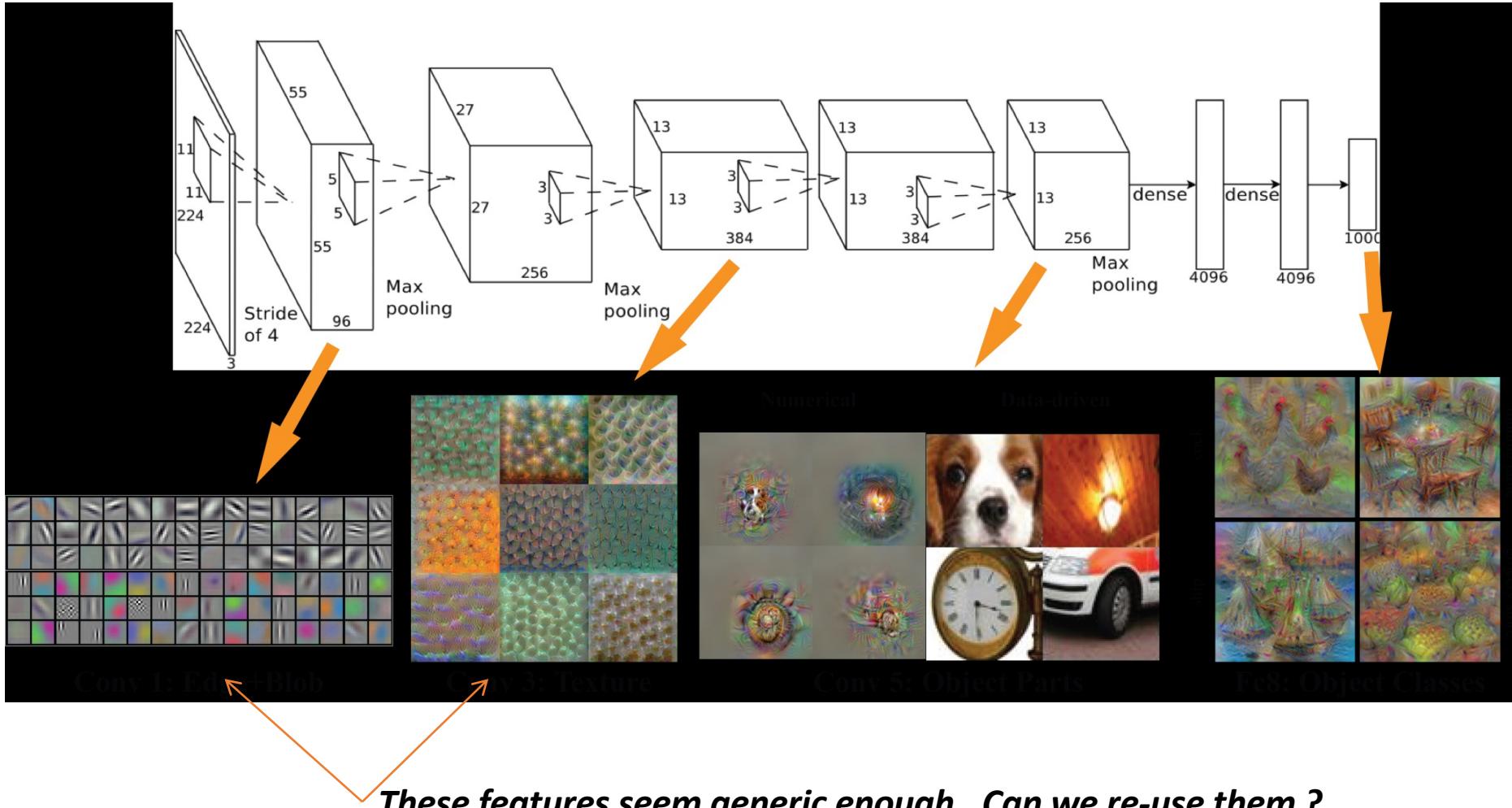


Transfer Learning





Filters learnt by AlexNet



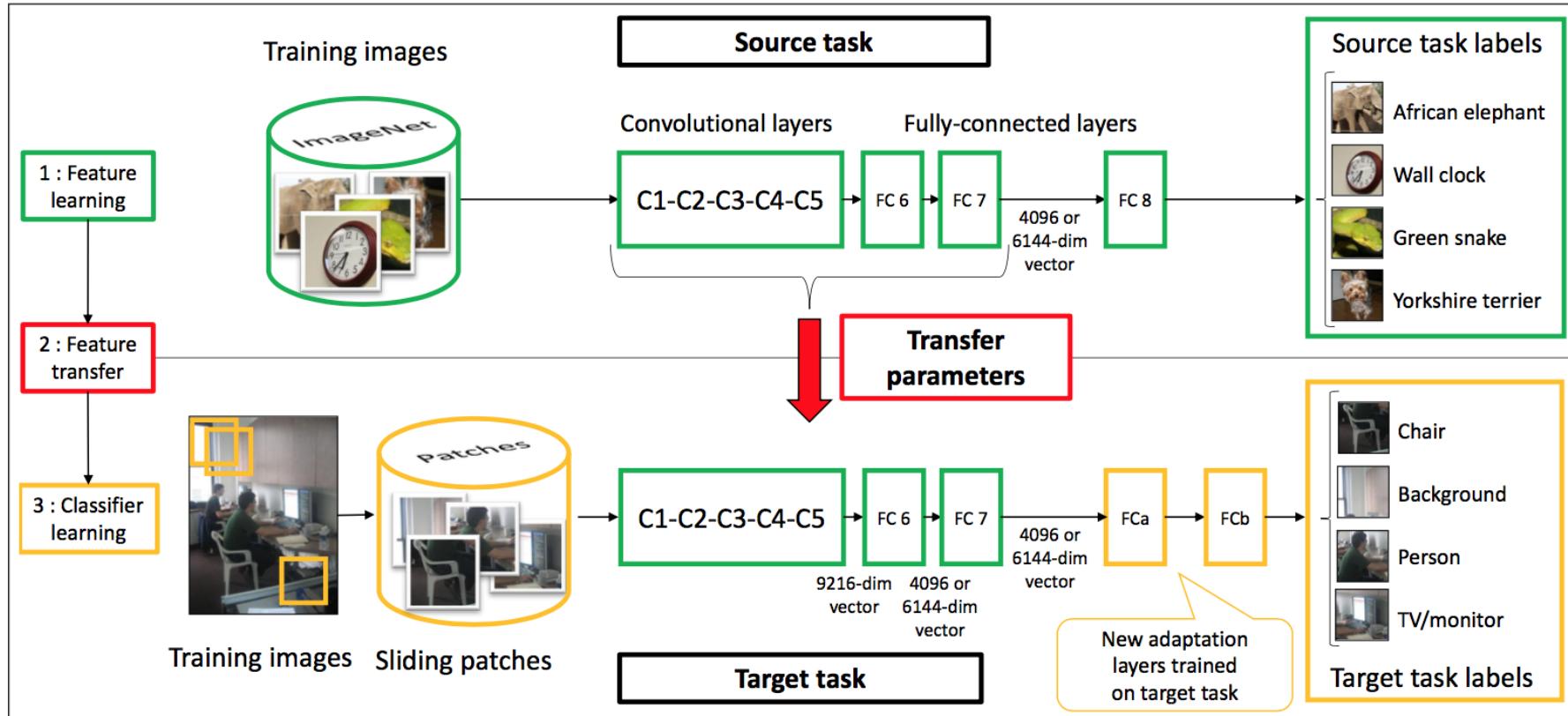


Transfer Learning



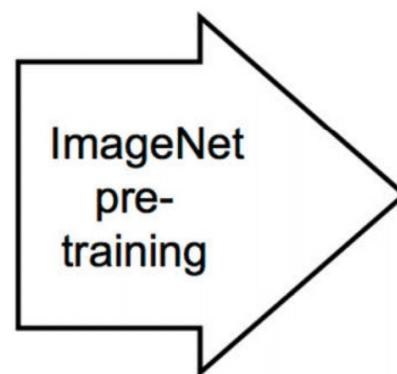


Transfer Learning





Examples



Pre-trained features are quite general

Object classification (PASCAL VOC)

- ▶ [Chatfield et al. 2014, Razavian et al. 2014, Zeiler et al. 2014]

Object detection (PASCAL VOC)

- ▶ R-CNN [Girshick et al. 2014]

Fine-grained classification (UCSD birds)

- ▶ Part-R-CNN [Zhang et al. 2014]

MIT 67 scene classification

- ▶ [Razavin et al. 2014]

...



Semantic segmentation



Input



Segmentation [9]



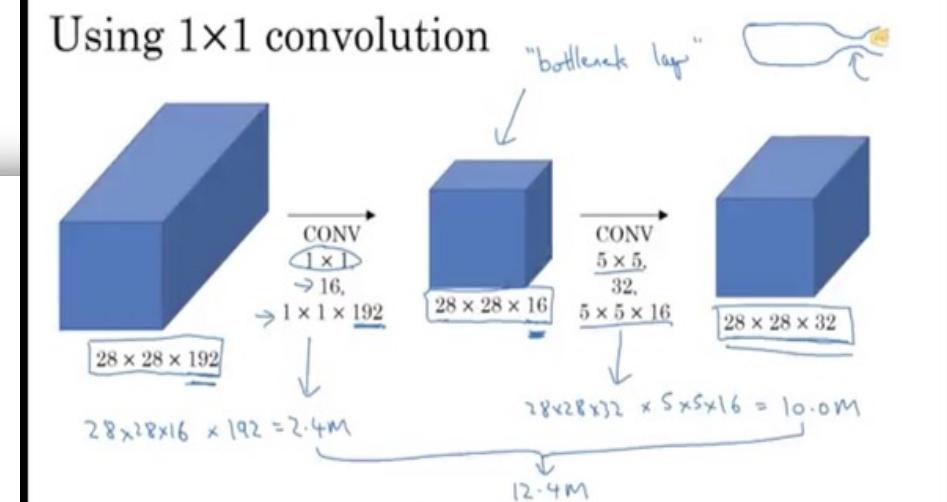
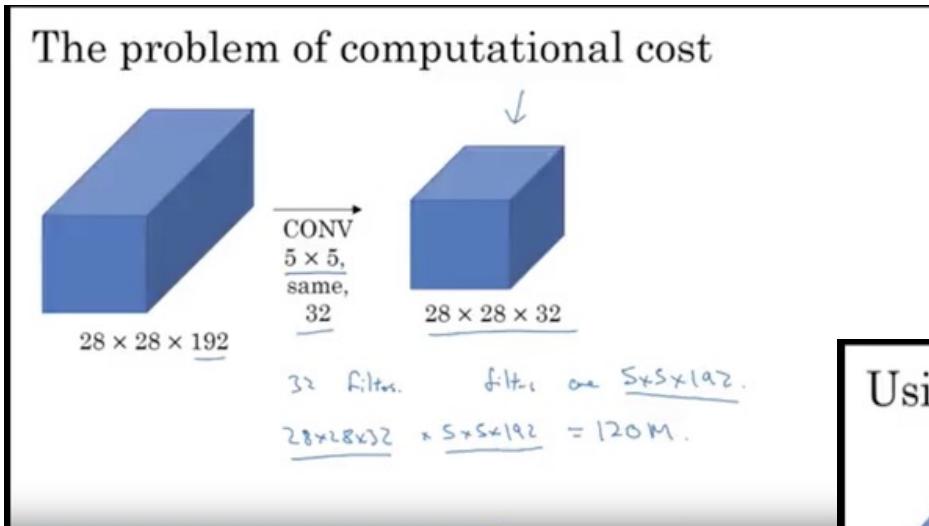
Semantic segmentation





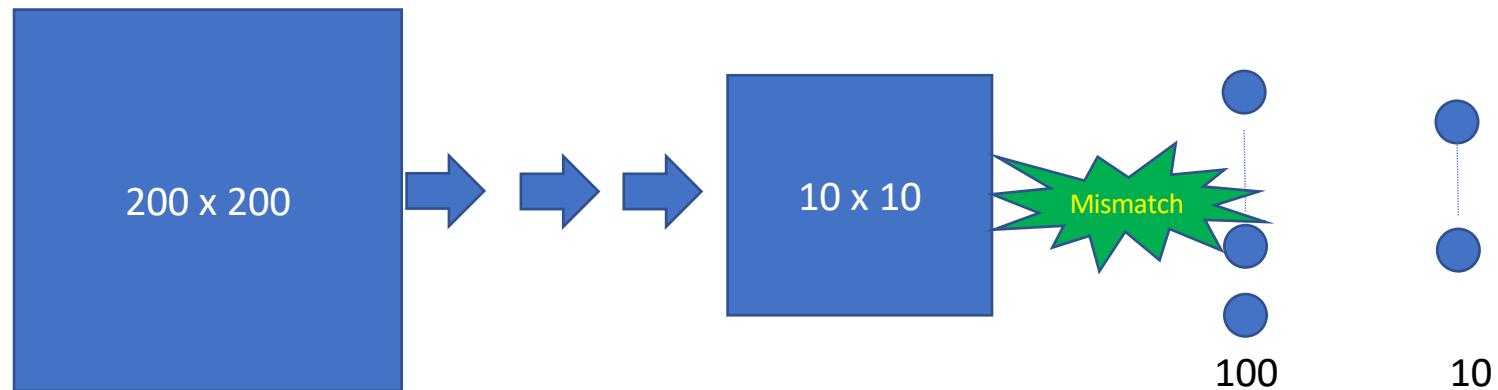
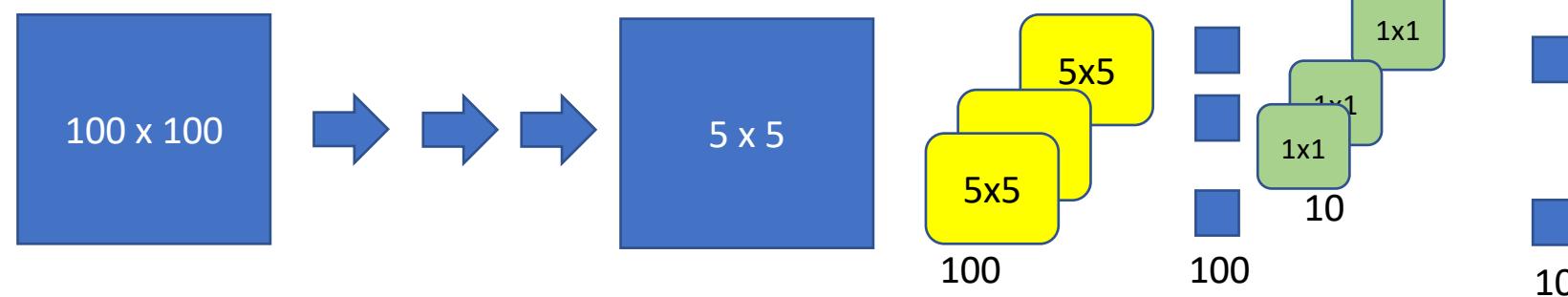
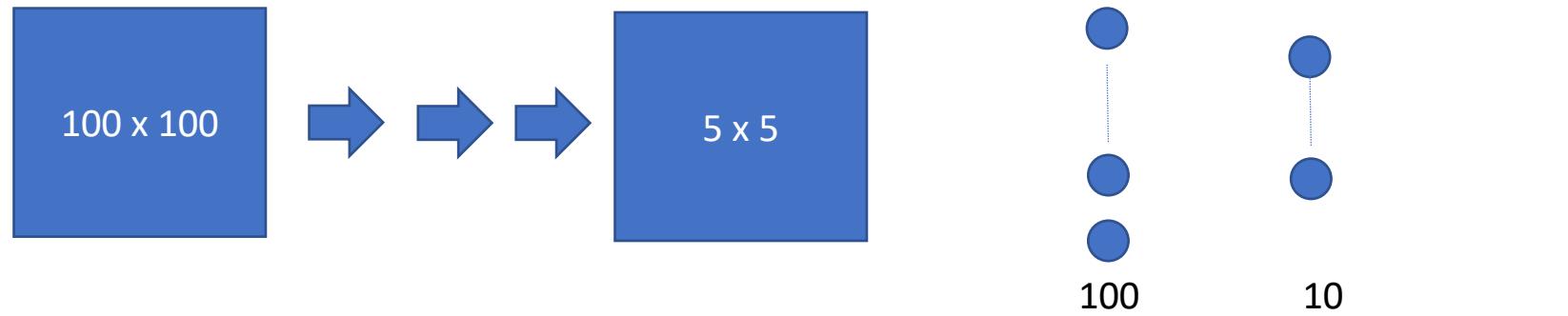
1x1 (Bottleneck) Convolution Layer

- Enable dimensionality reduction (computational savings)
- Add more non-linearity to representation space



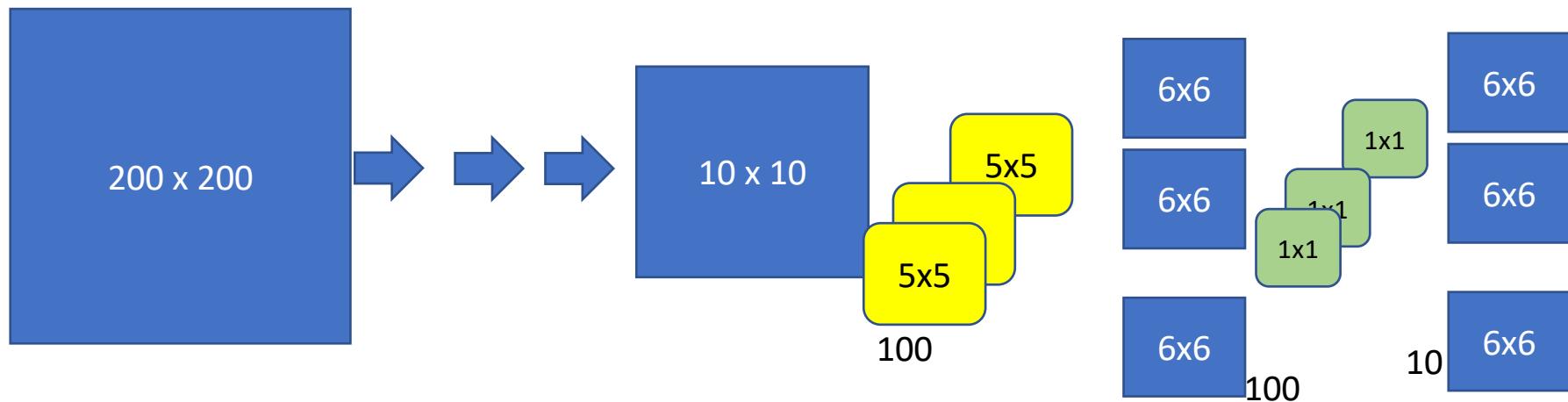
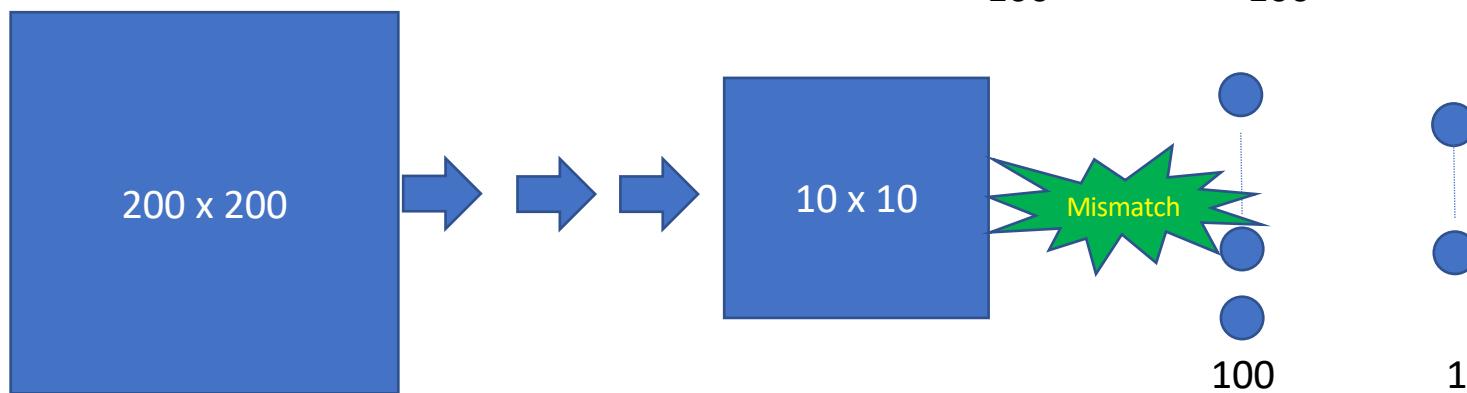
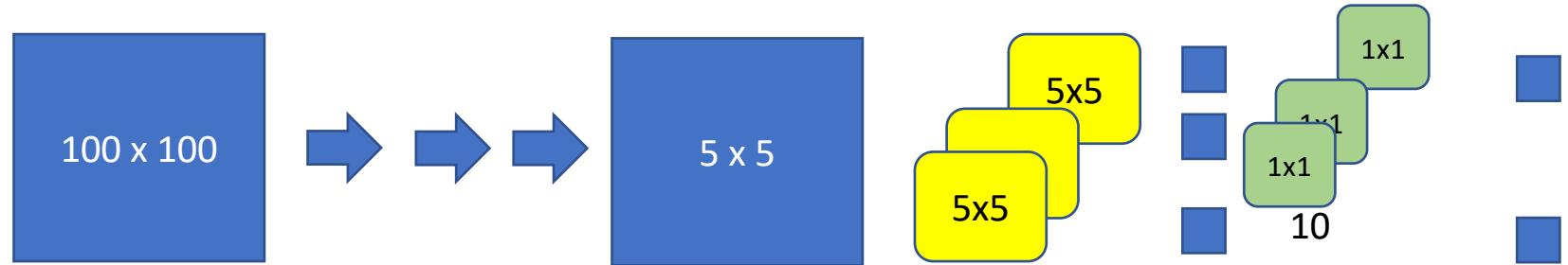


Converting FC Layers to Conv Layers



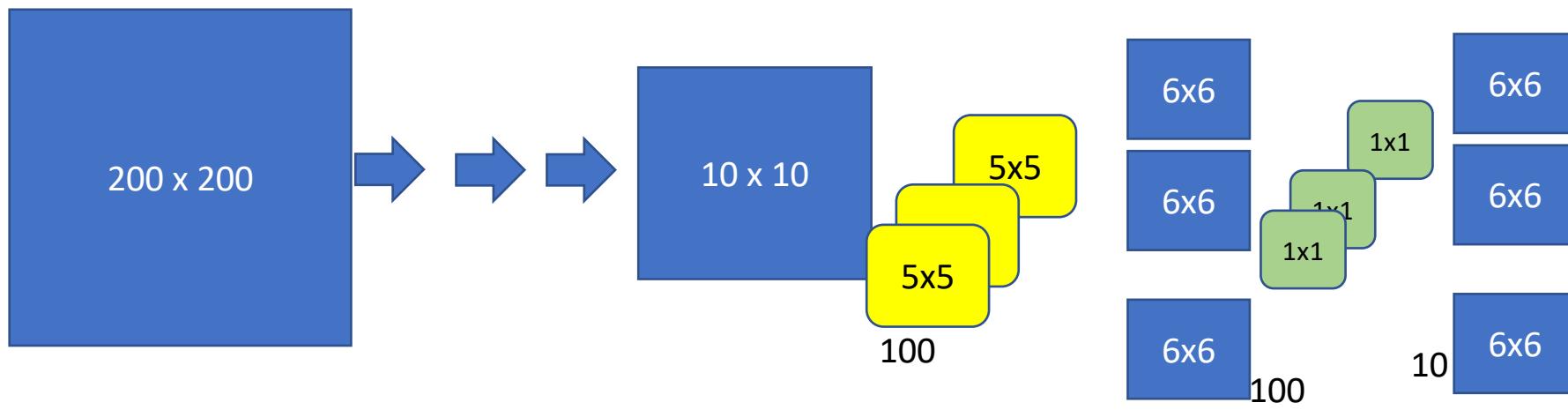
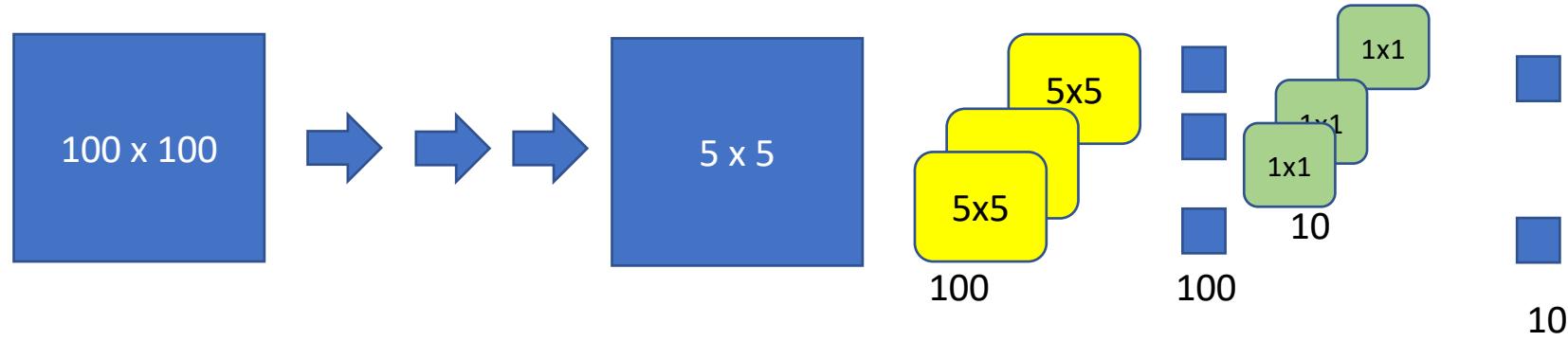


Converting FC Layers to Conv Layers





Converting FC Layers to Conv Layers

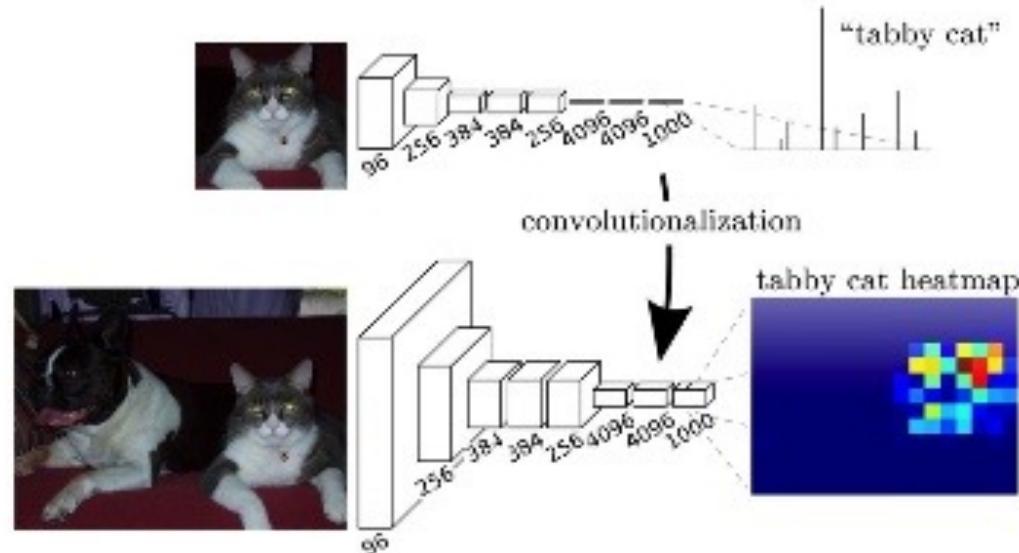


Scores for local image patches



Semantic segmentation

Reinterpret classification as coarse prediction

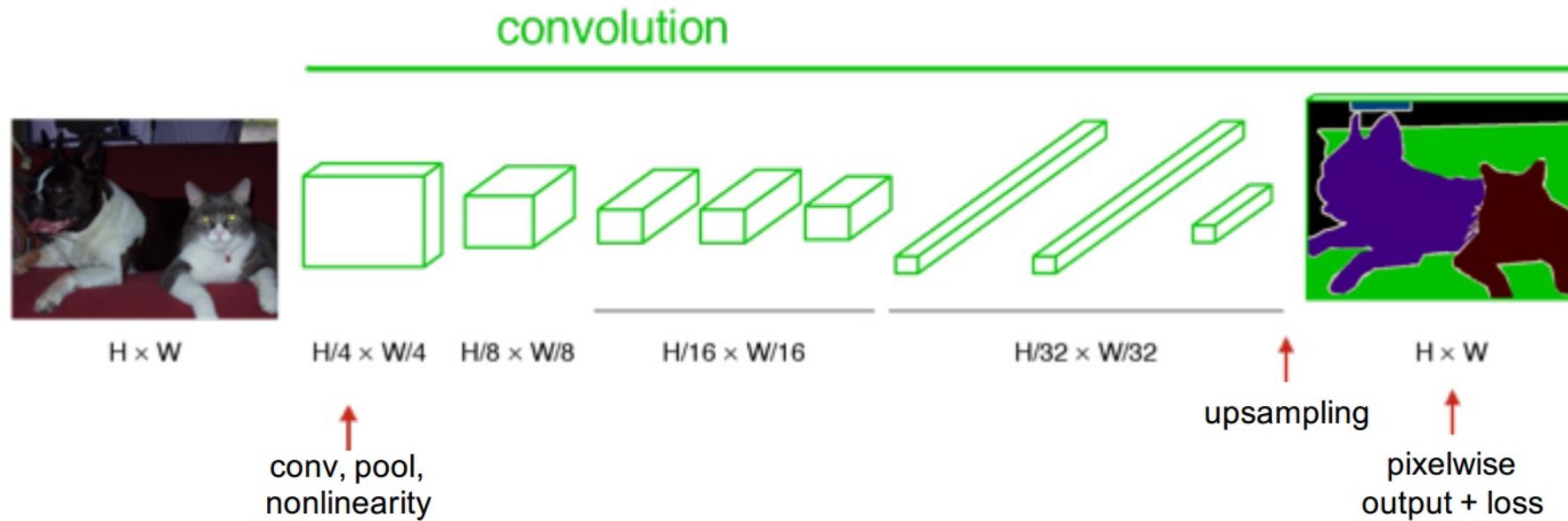


Shelhamer, Long, Darrell, [Fully Convolutional Networks for Semantic Segmentation](#), 2014-2016



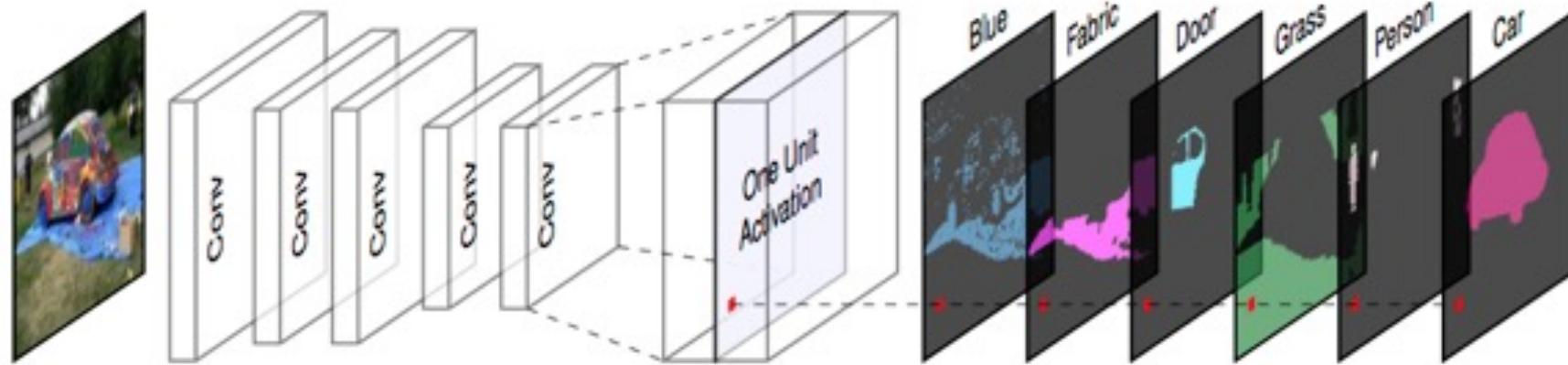
Semantic segmentation

Reinterpret classification as coarse prediction





Semantic segmentation



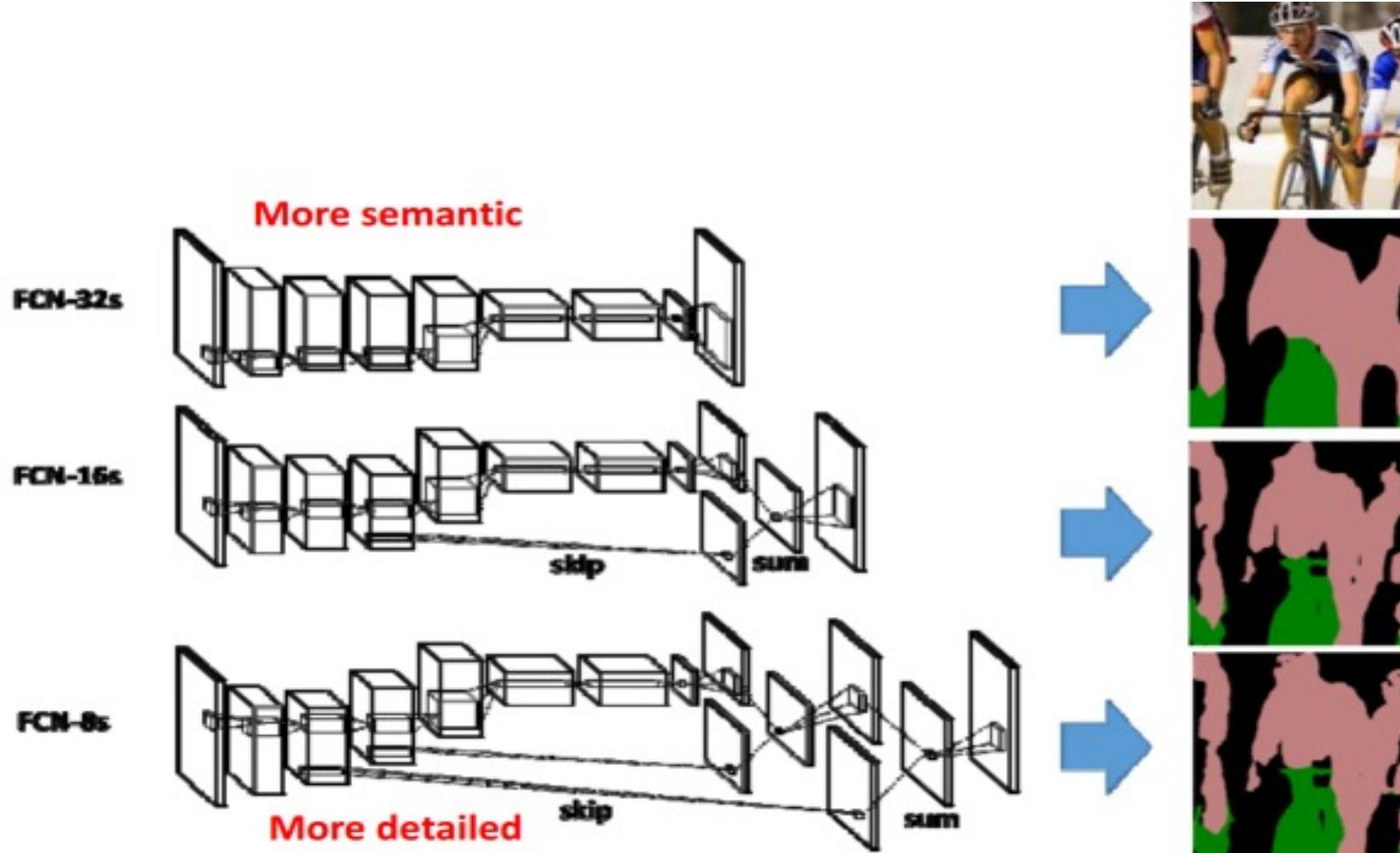
<http://netdissect.csail.mit.edu/thumb/slides-thumbnail.png>

https://leonardoaraujosantos.gitbooks.io/artificial-intelligence/content/image_segmentation.html



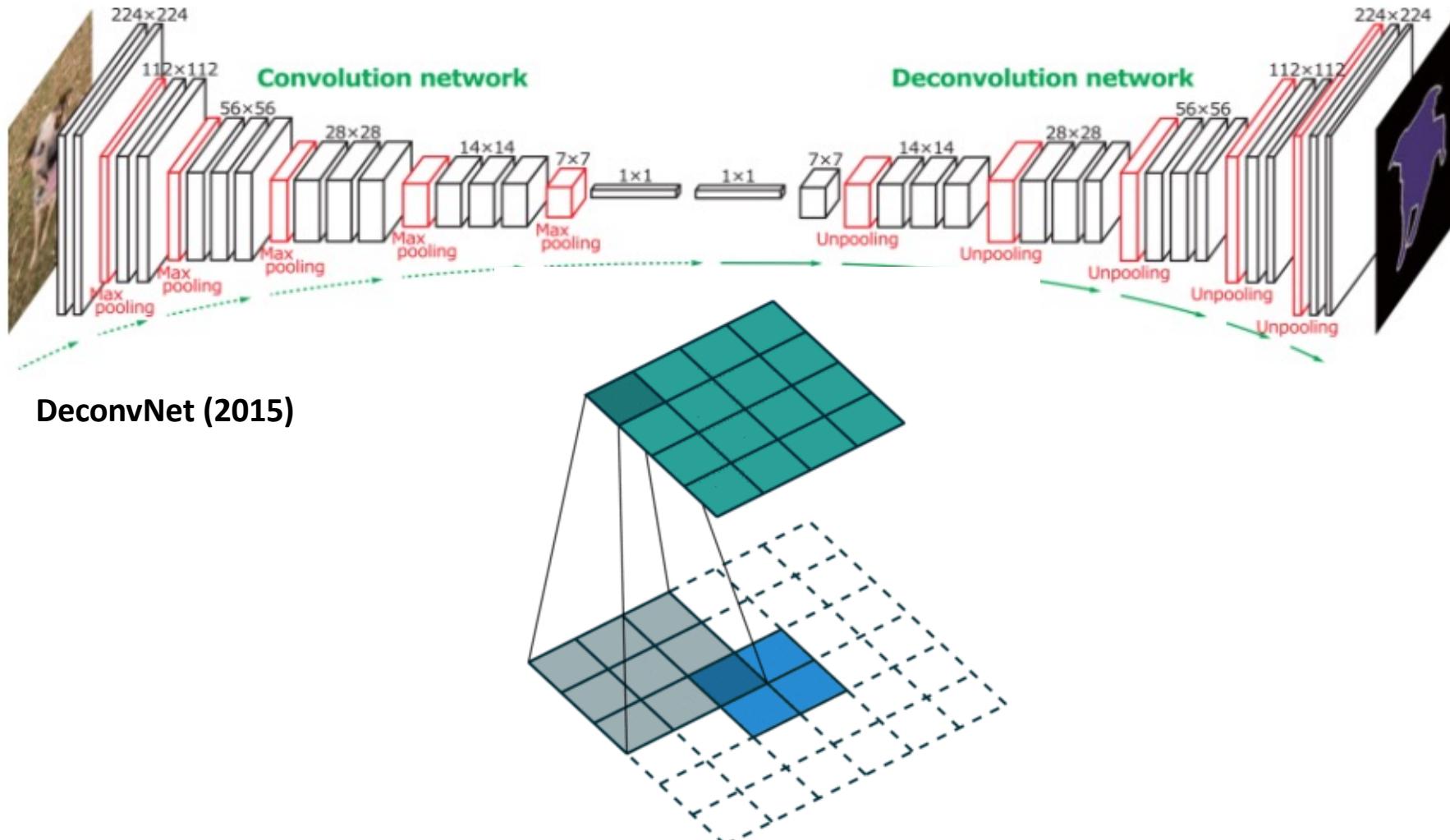
Semantic segmentation

Skip connections



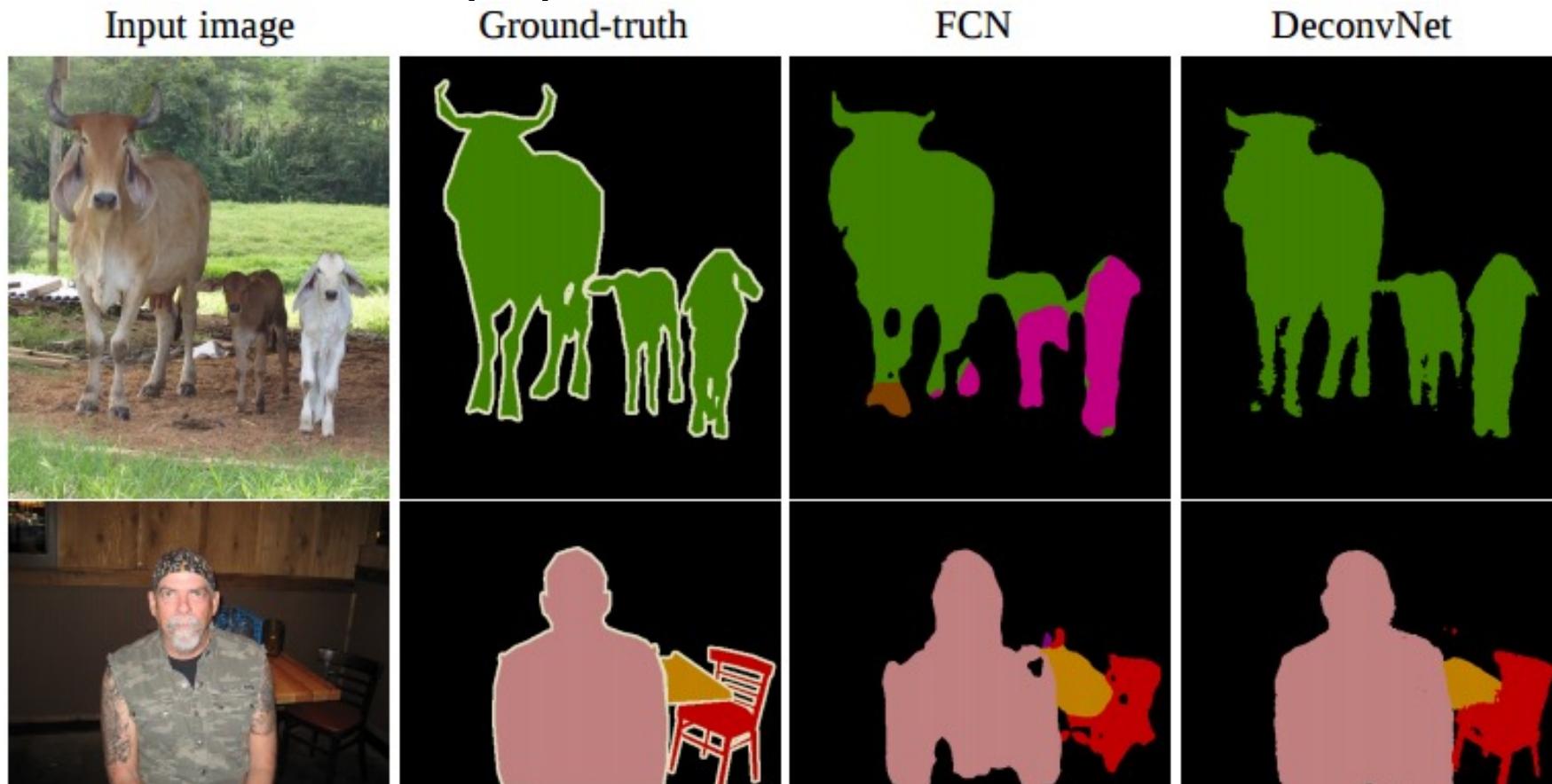


(Extreme) Semantic





(Extreme) Semantic



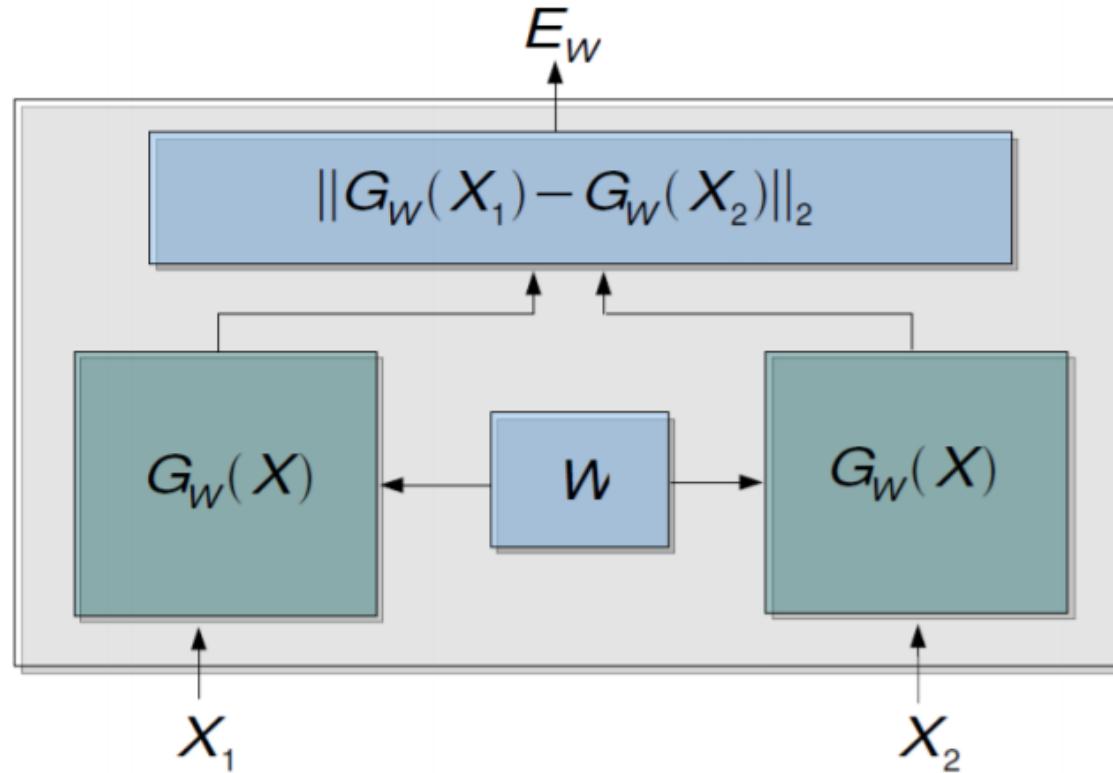


Siamese Network





Siamese Network

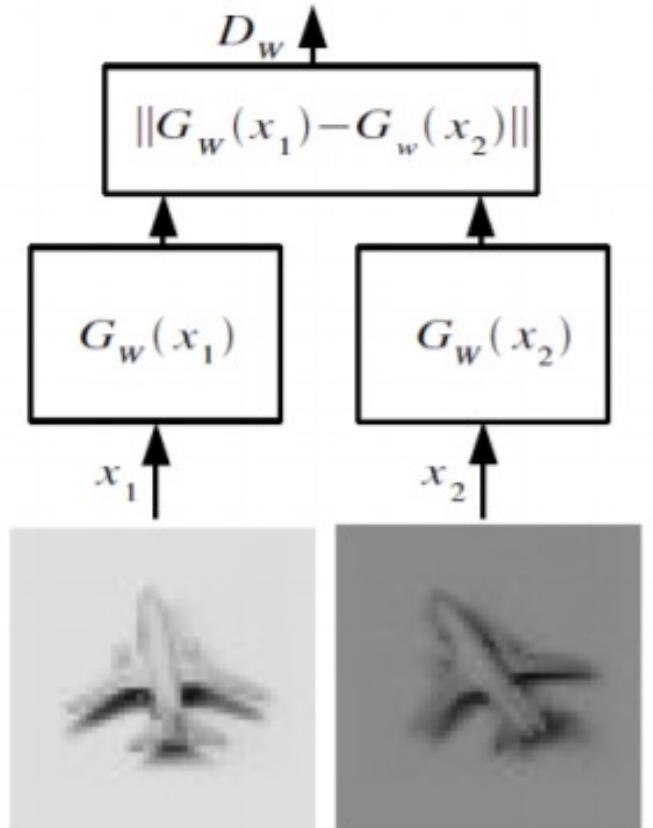




Siamese Network

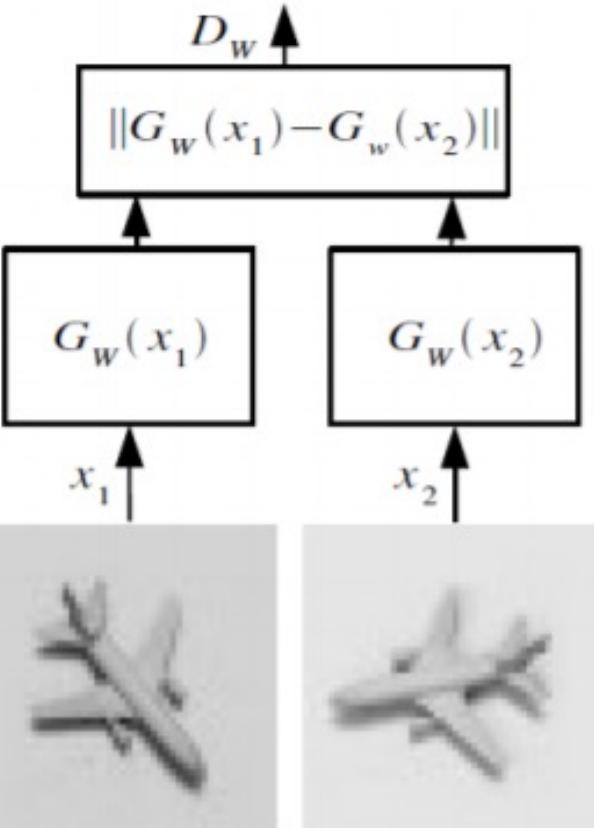
$$(1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Make this small

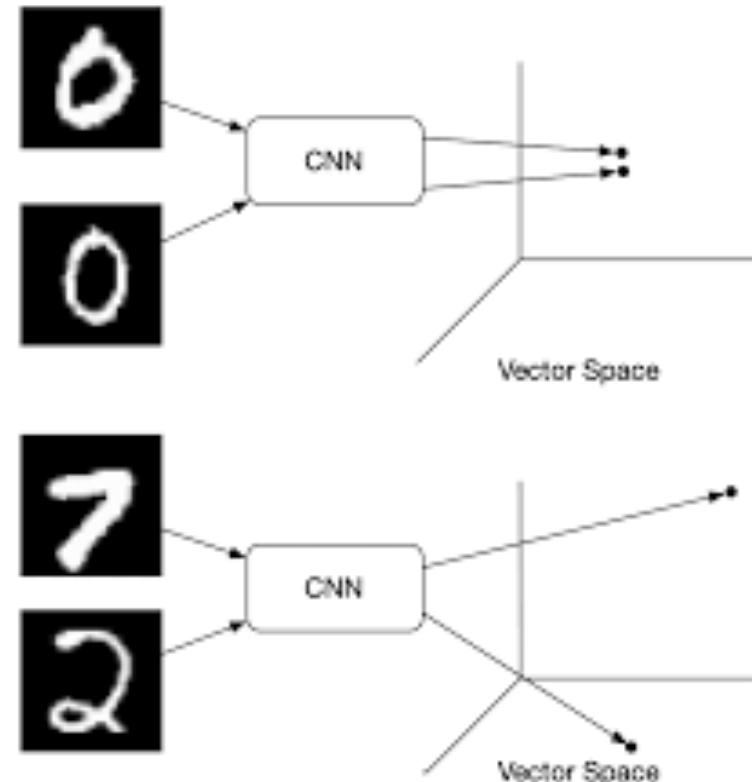


Similar Images ($Y=0$)

Make this large



Dissimilar Images ($Y=1$)

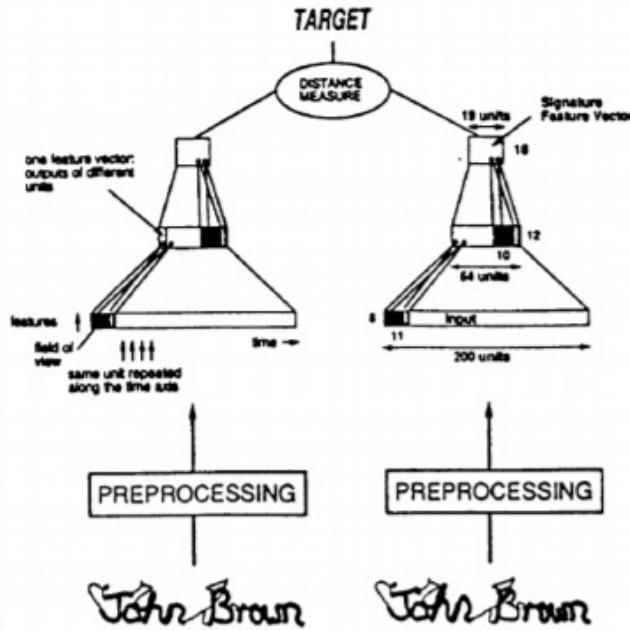




Siamese Network

Application in Signature Verification

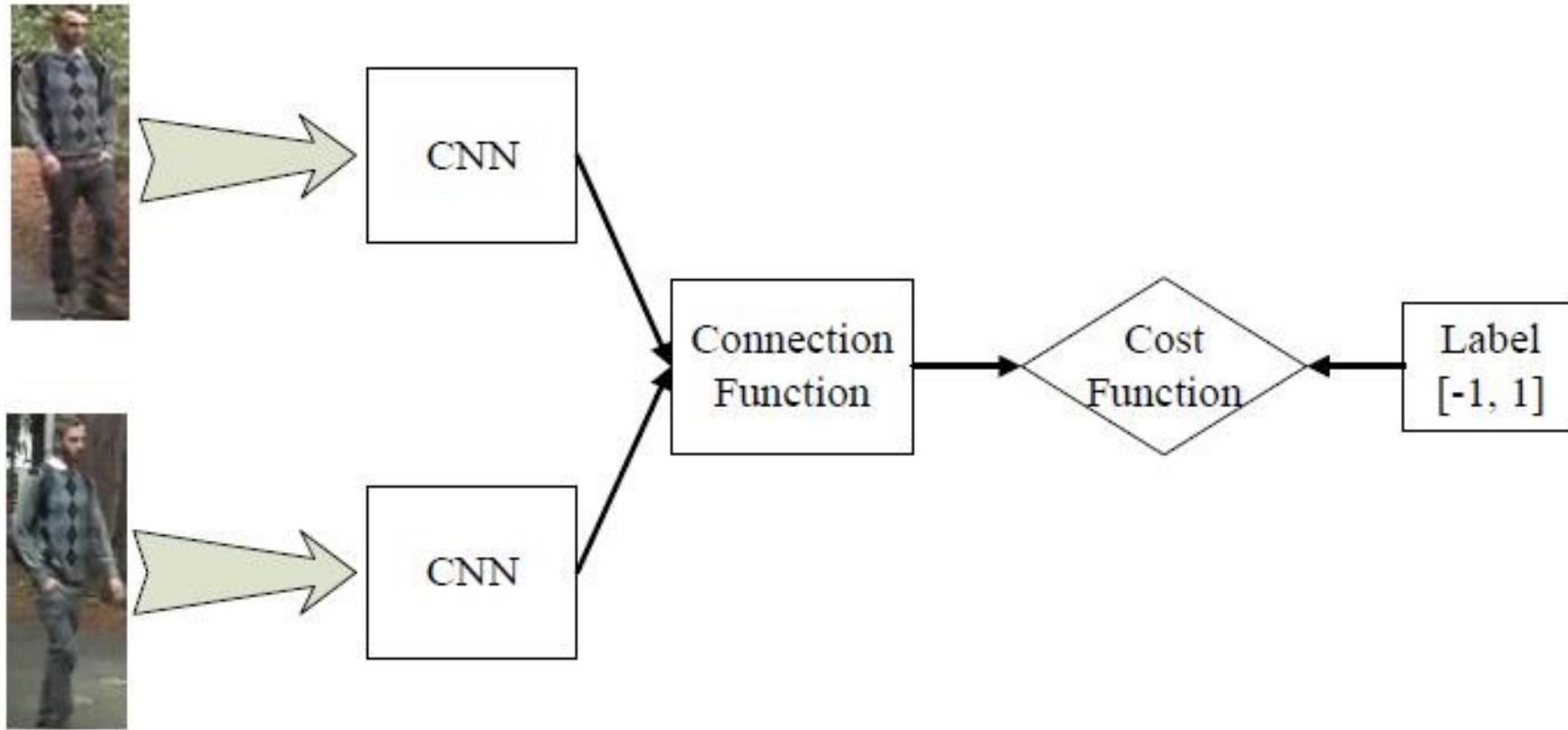
- The input is 8(feature) x 200(time) units.
- The cosine distance was used, (1 for genuine pairs, -1 for forgery pairs)



Bromley J, Guyon I, Lecun Y, et al. Signature Verification using a "Siamese" Time Delay Neural Network, NIPS Proc. 1994



Siamese Network (Person re-id)



http://www.fubin.org/research/Person_ReID/Person_ReID.html



MML

Questions?





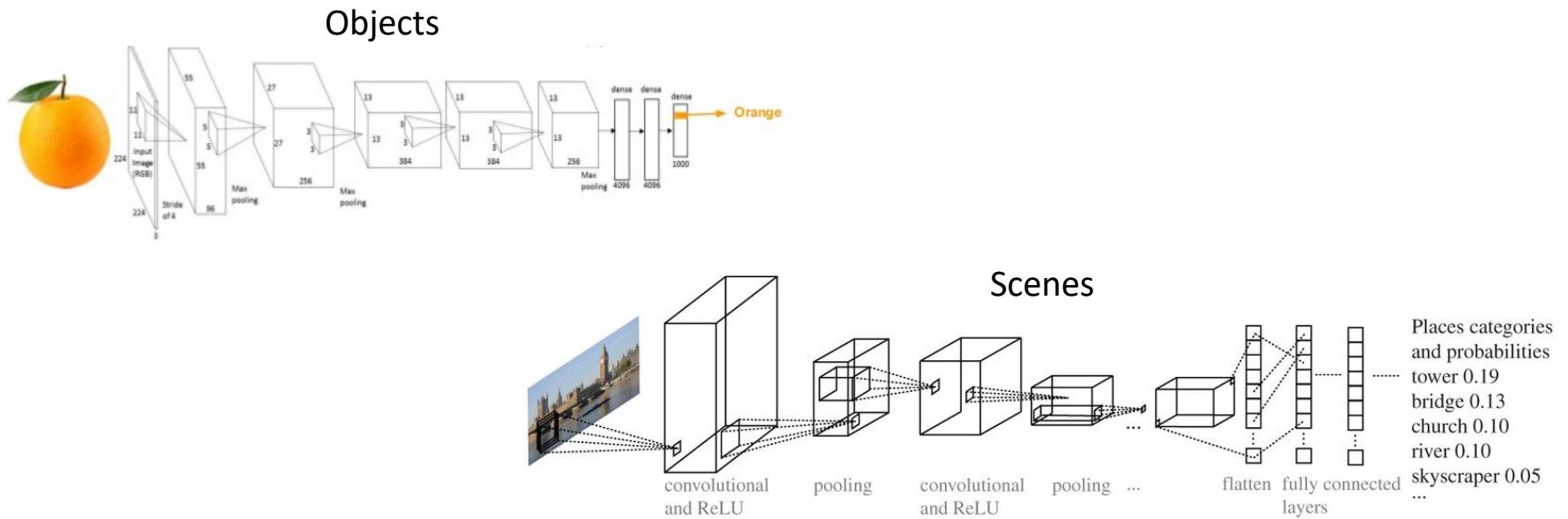
MML

Landscape of CNNs

Applications



Image → Label



DeepFace Architecture

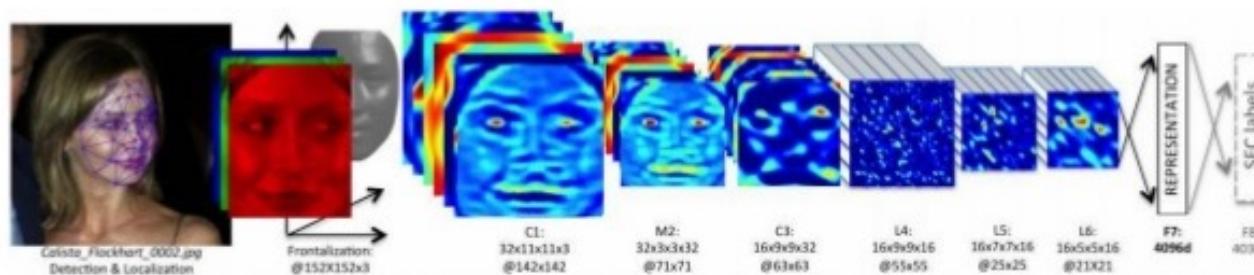
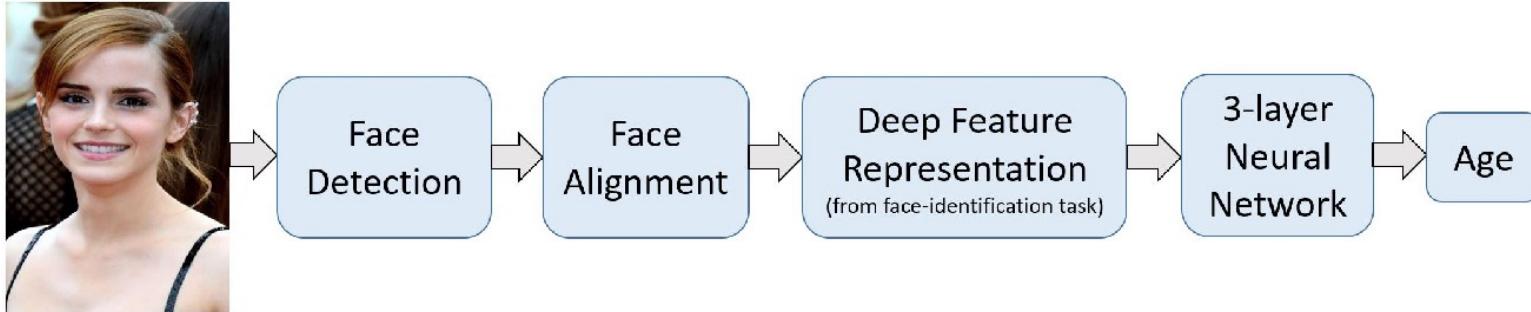
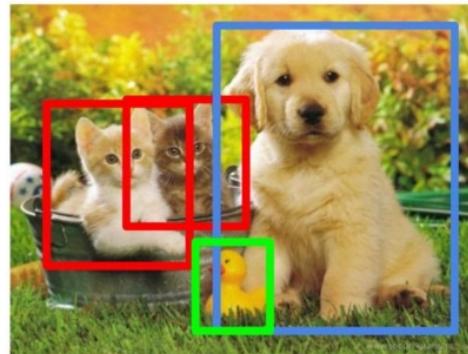




Image → Number



Object Detection

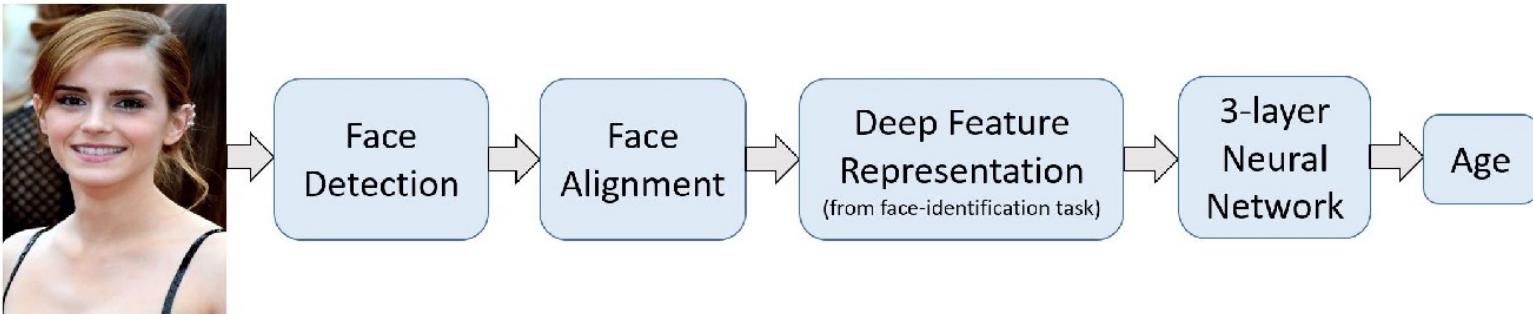


CAT, DOG, DUCK



Image → Number

Age Estimation



Object detection

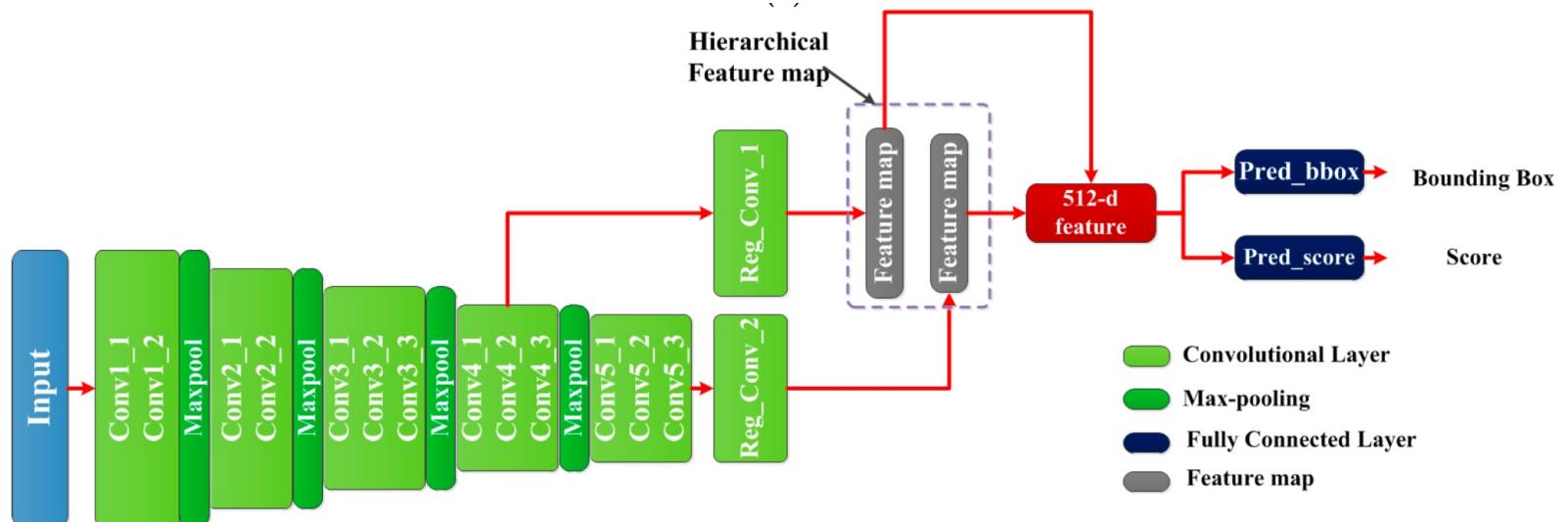
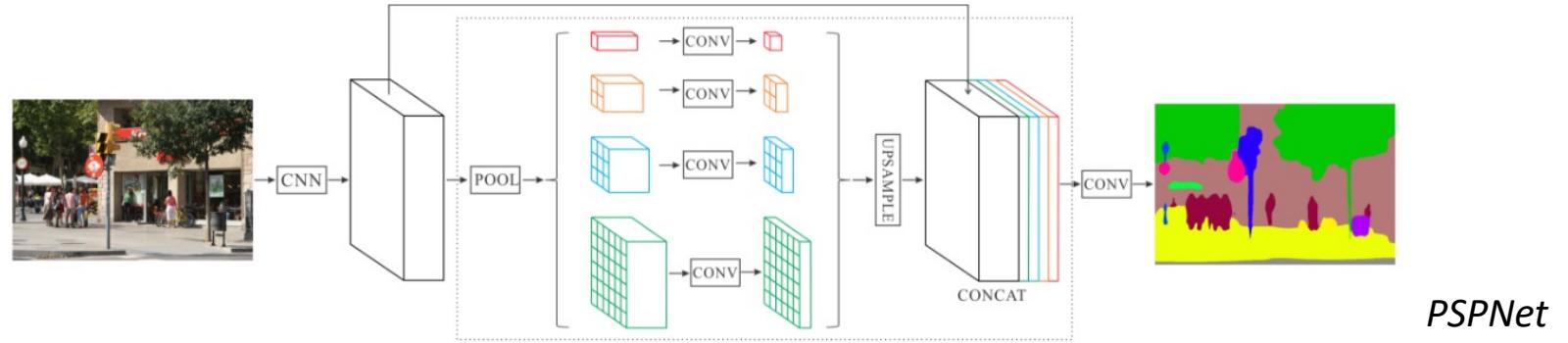




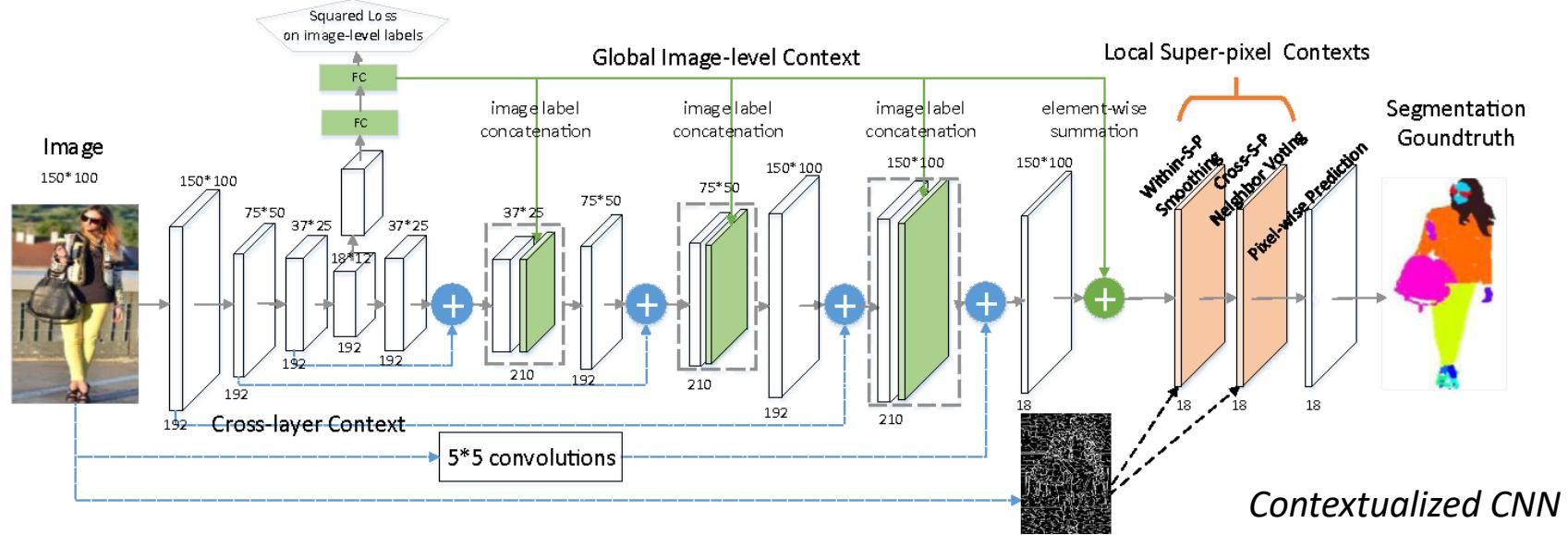
Image → Label Image

Scene Parsing



PSPNet

Object Parsing

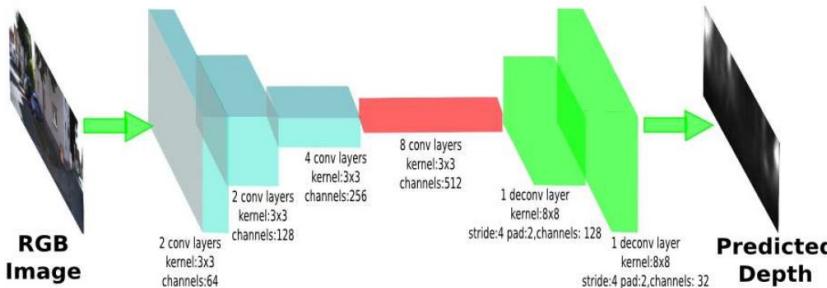


Contextualized CNN

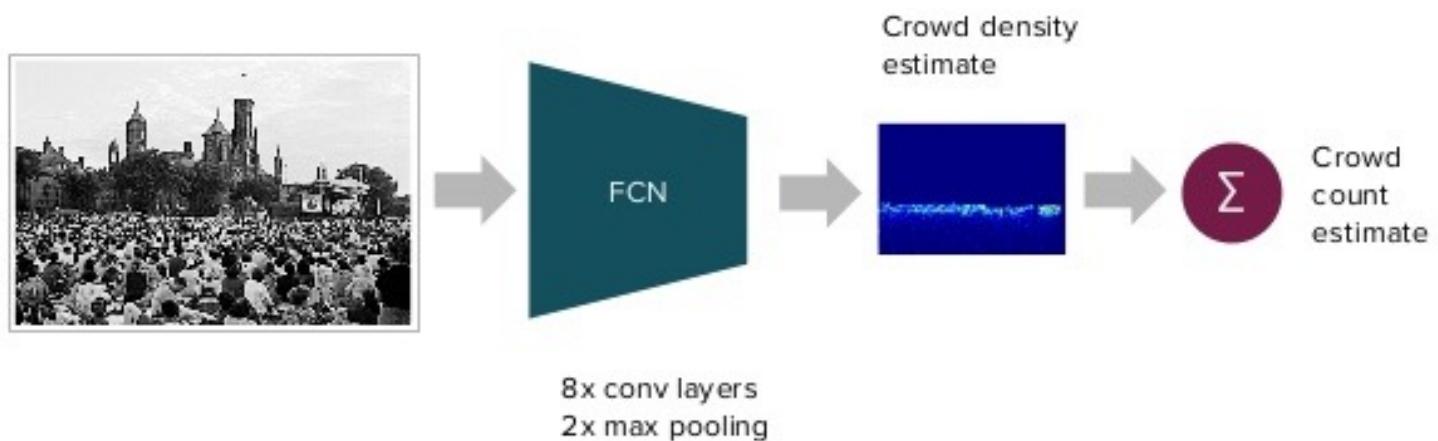


Image → Image

Depth Estimation



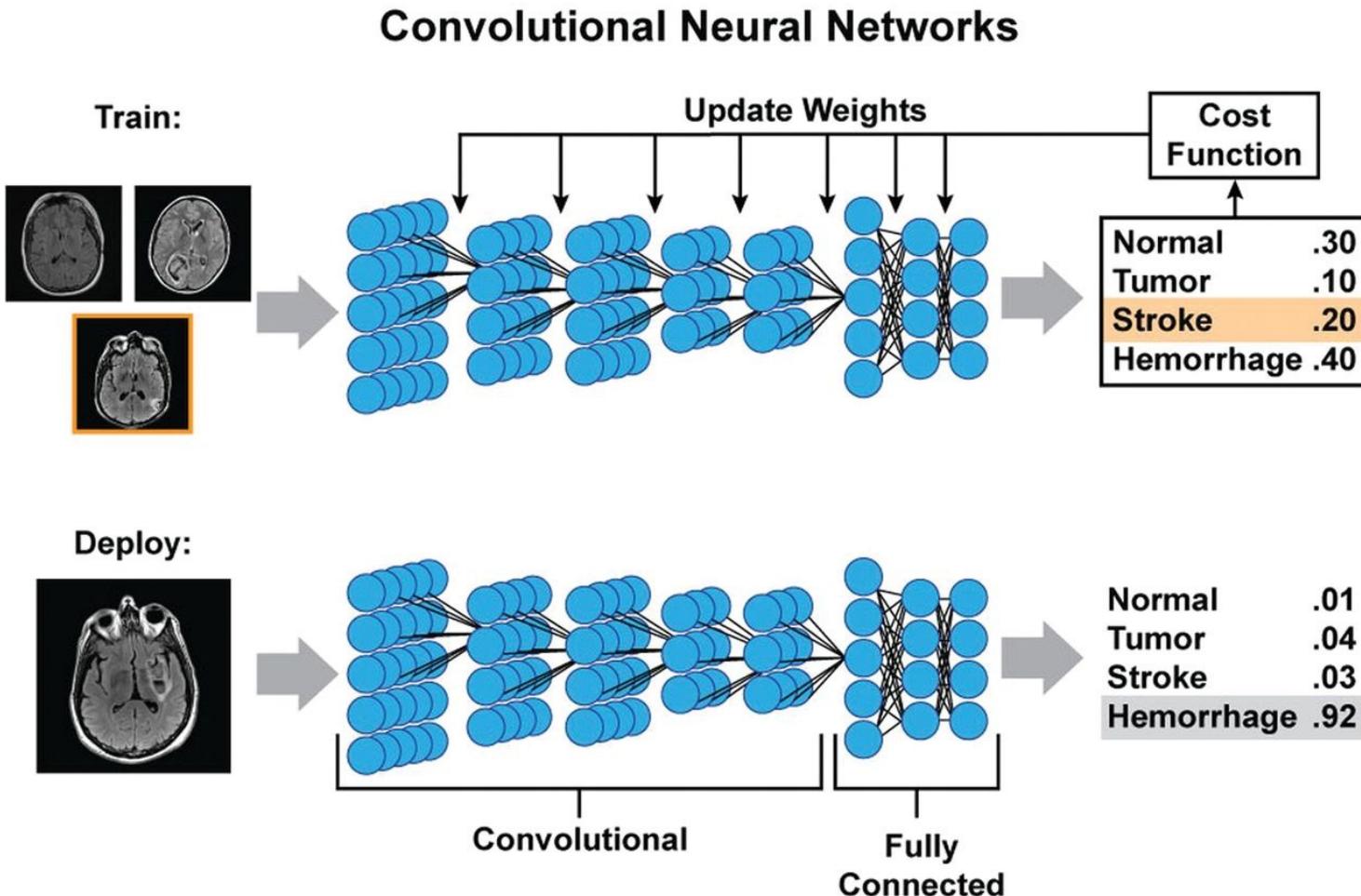
Crowd Counting





Multi-channel input

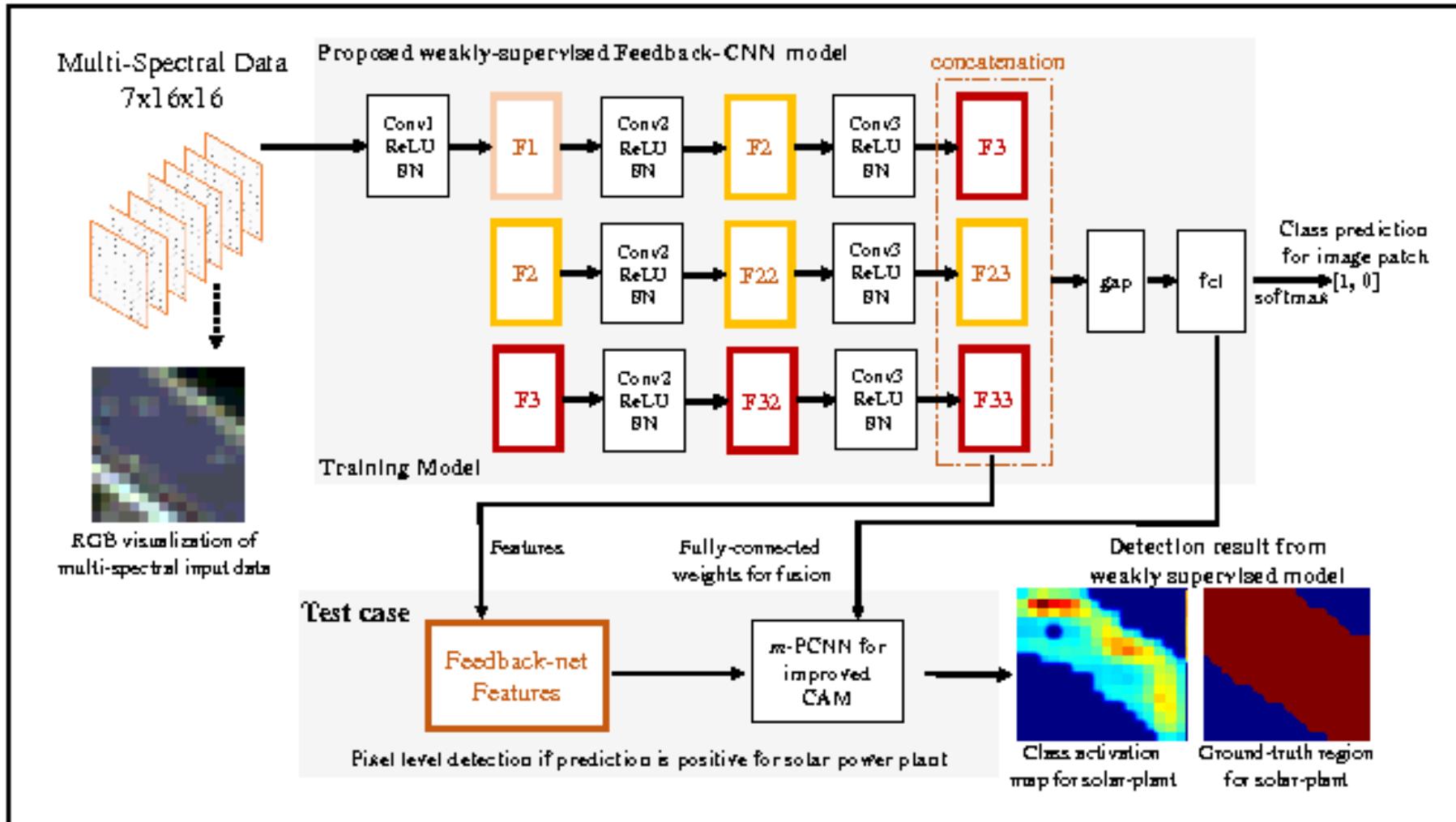
NeuroRadiology (fMRI)





Multi-channel input

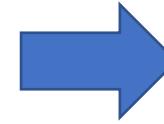
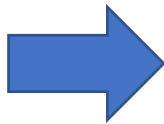
Solar-power plant detection from multi-spectral data





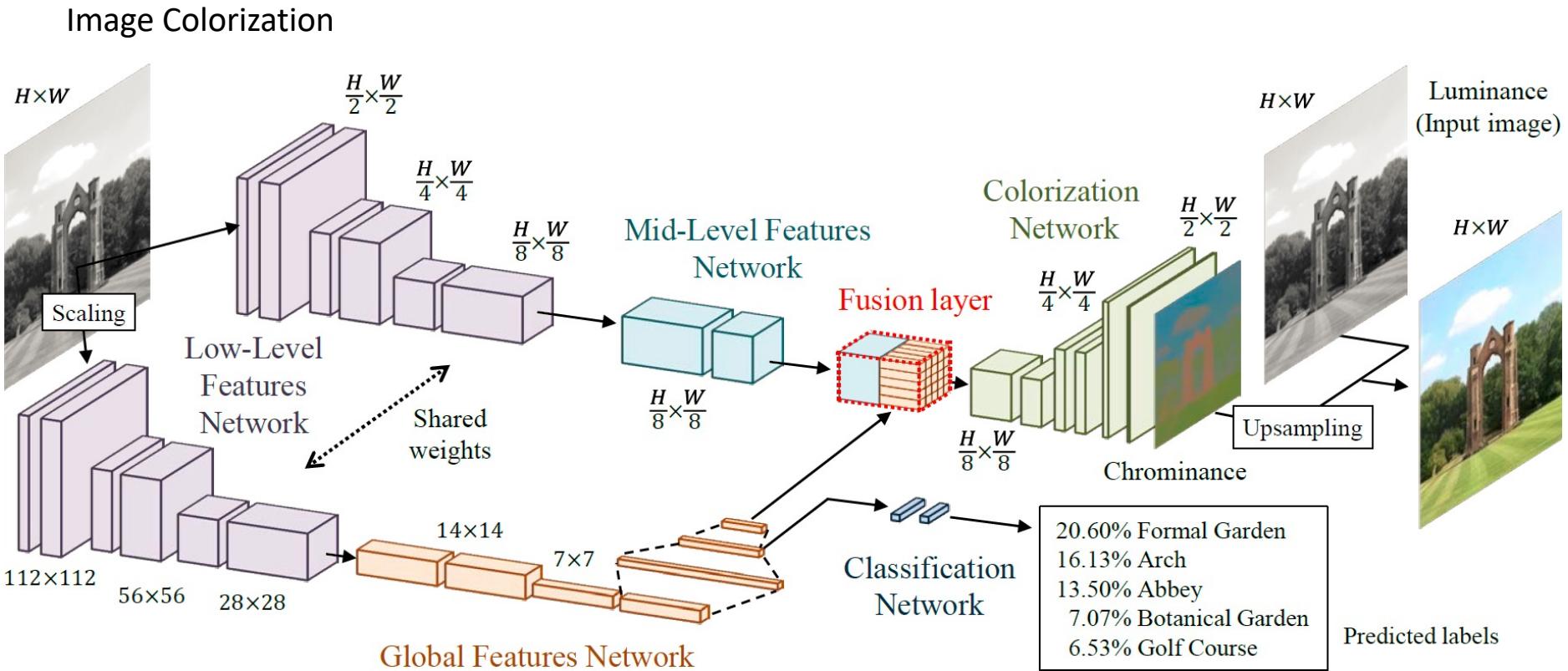
Multi-branch input

Image Colorization





Multi-branch input

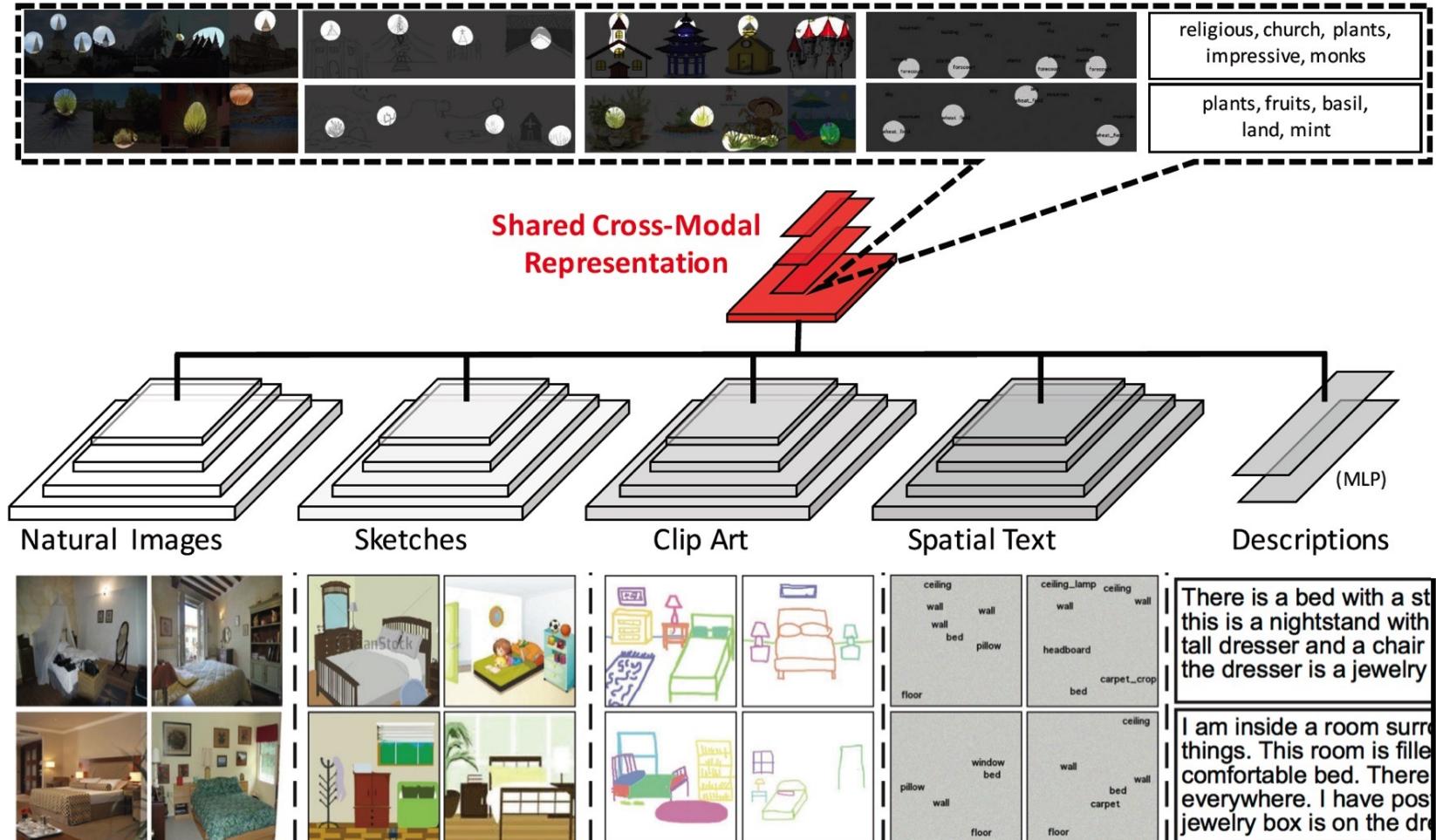


<http://hi.cs.waseda.ac.jp/~iizuka/projects/colorization/en/>



Multi-branch input

Cross-modal Scene Classification





Multi-branch input

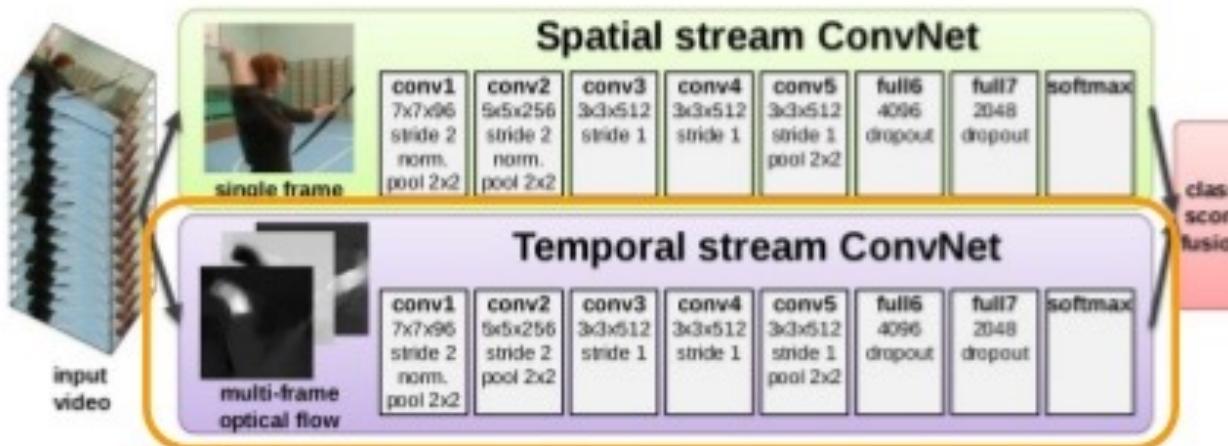
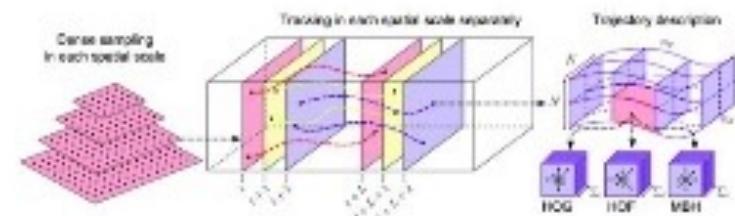
Activity Recognition

Recognition: Two stream

Two CNNs in parallel:

- One for RGB images
- One for Optical flow (hand-crafted features)

Fusion after the softmax layer

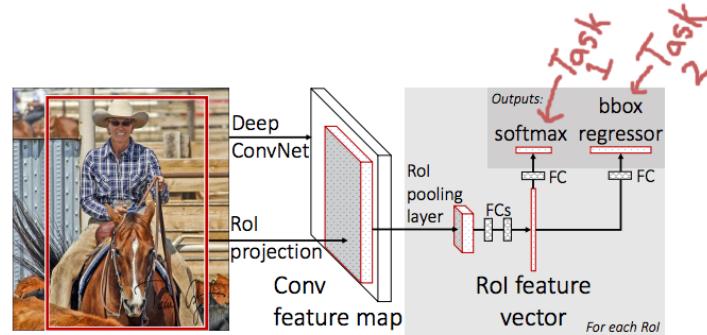


Simonyan, Karen, and Andrew Zisserman. ["Two-stream convolutional networks for action recognition in videos."](#) NIPS 2014.



Multi-branch output

Object Detection, Classification



Scene Parsing

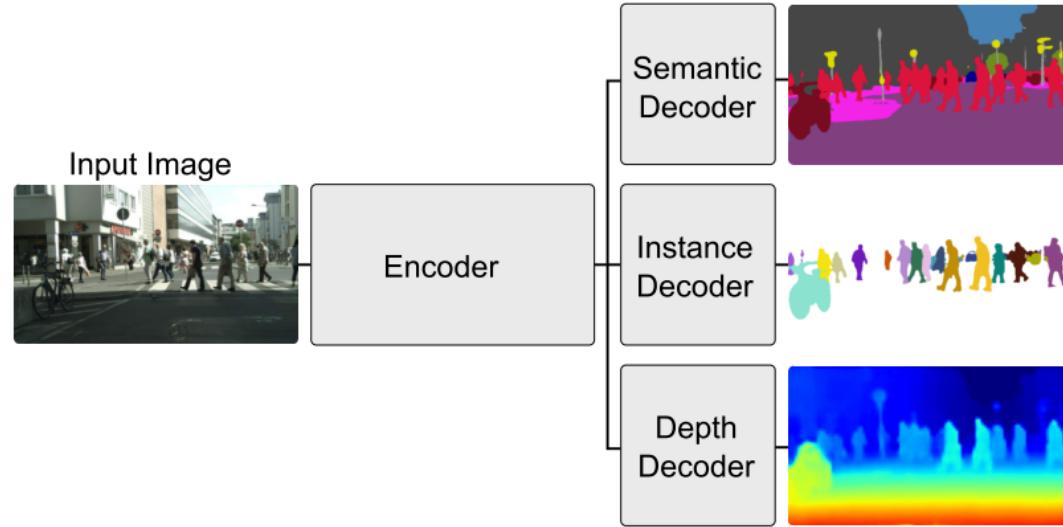
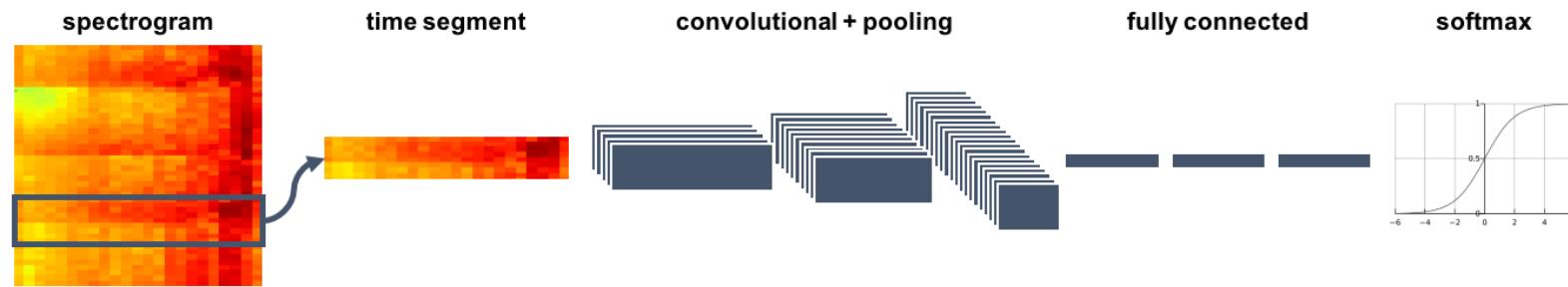




Image CNNs for non-image data

Audio Beat Detection





MML

Questions?

