

practical deep learning

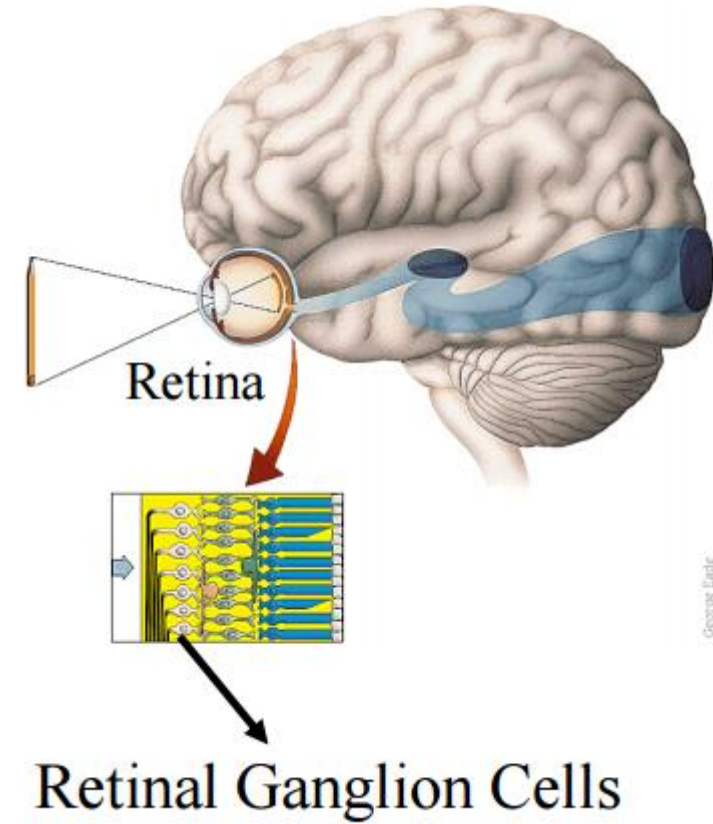
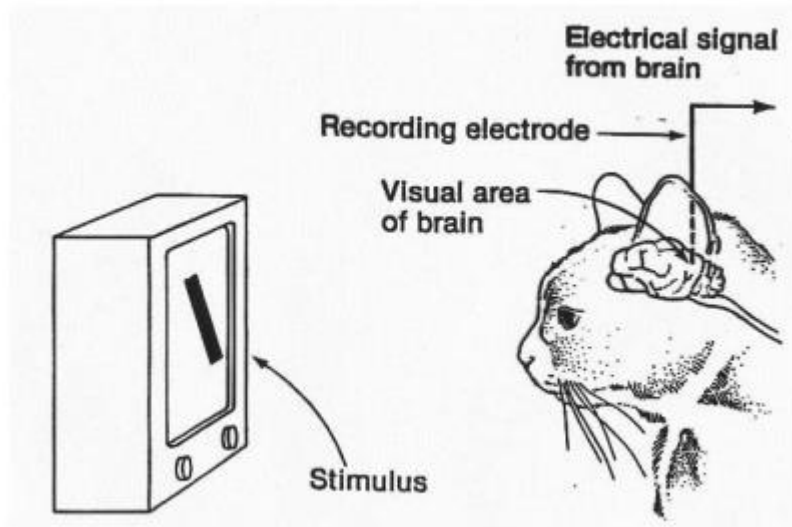
convolutional neural networks

Alex Honchar
University of Verona

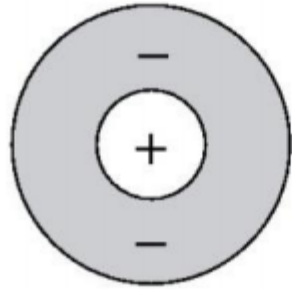
Day 3 goals

- You **understand** how convolutional neural network (CNN) works
- You **understand** modern approaches to CNNs
- You **can** train your own CNN for different computer vision problems
- You **can** use **Caffe** framework

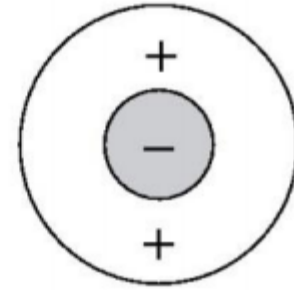
Mammals vision system



Mammals vision system

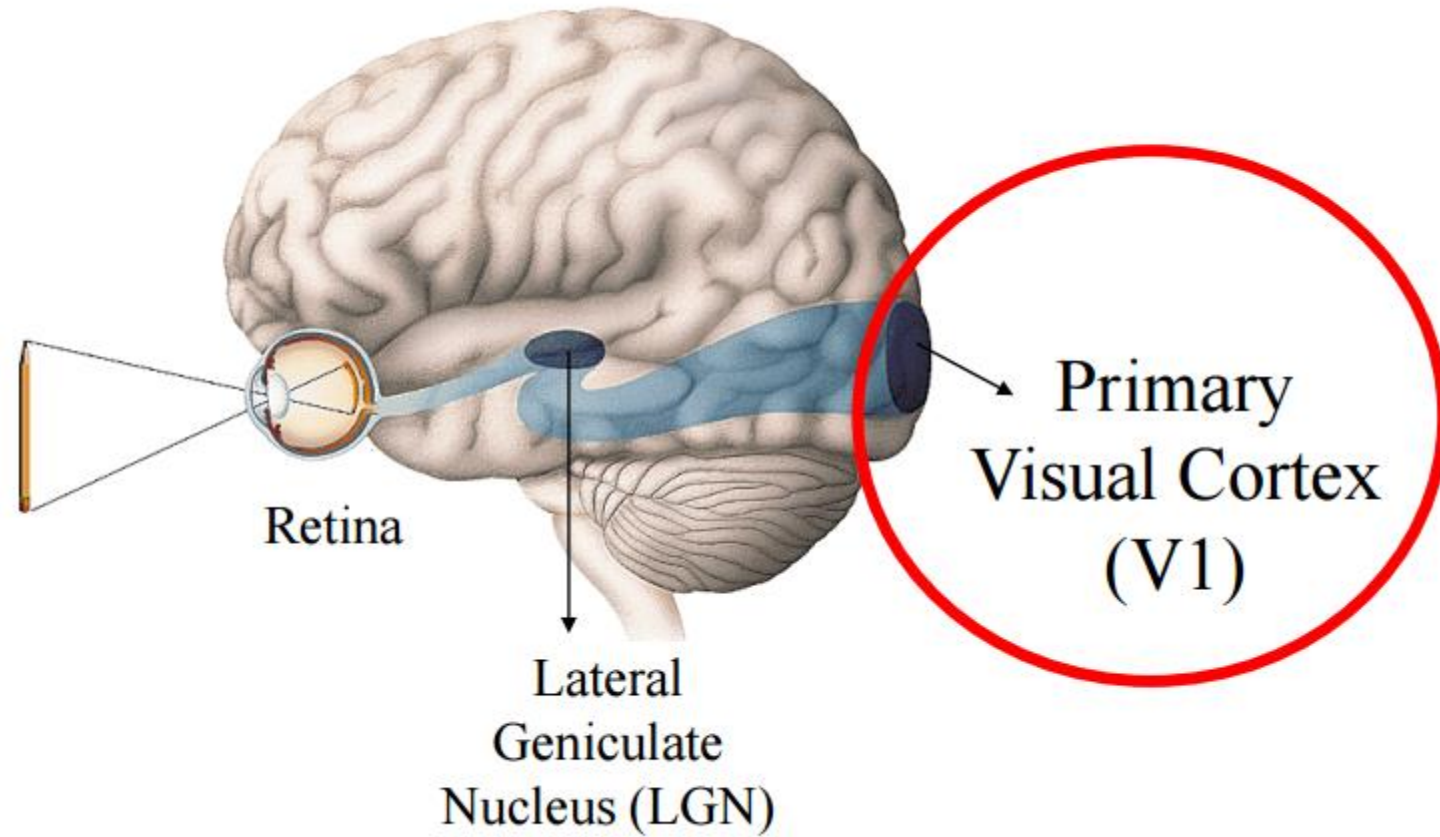


On-Center
Off-Surround
Receptive Field

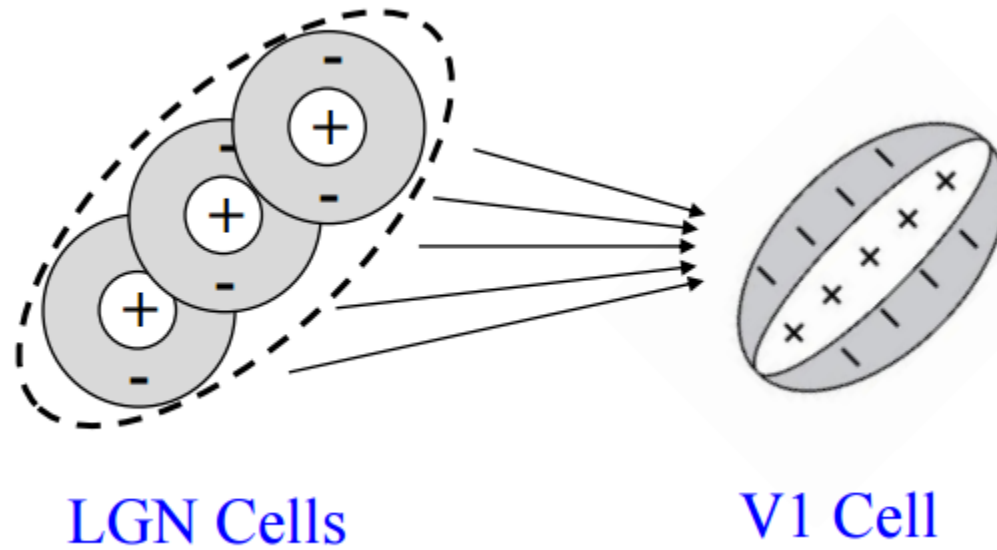


Off-Center
On-Surround
Receptive Field

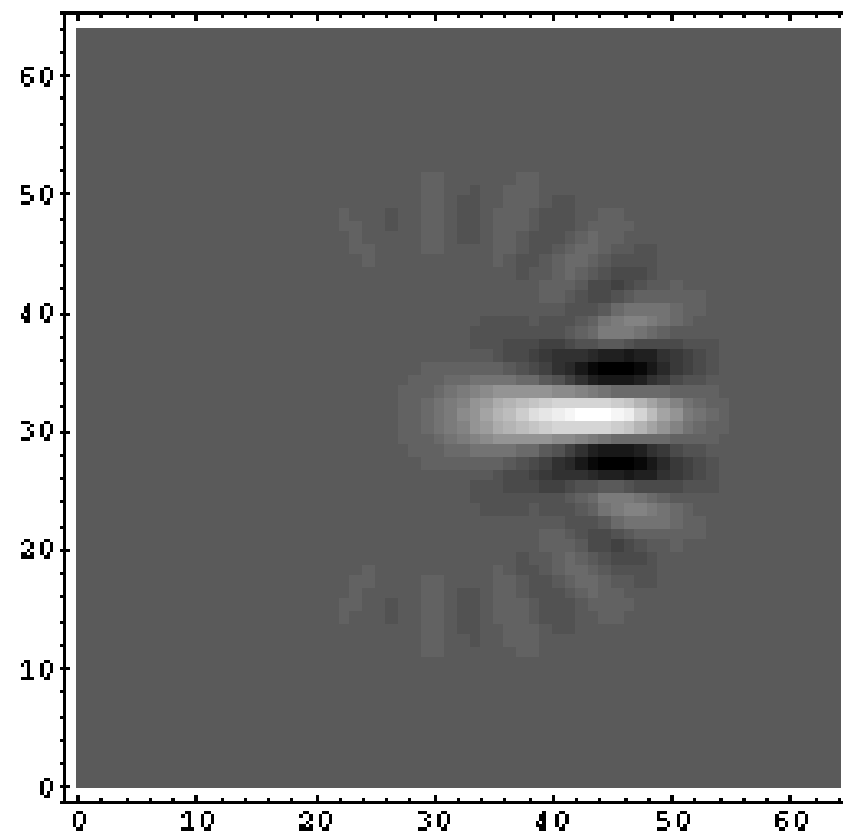
Mammals vision system



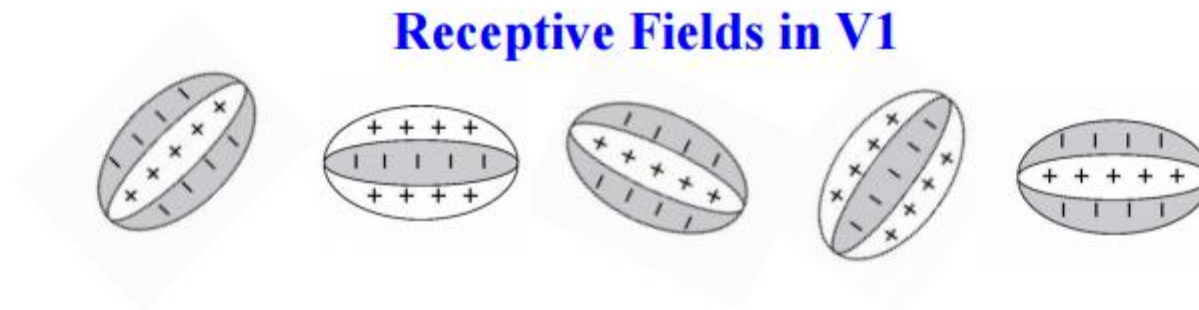
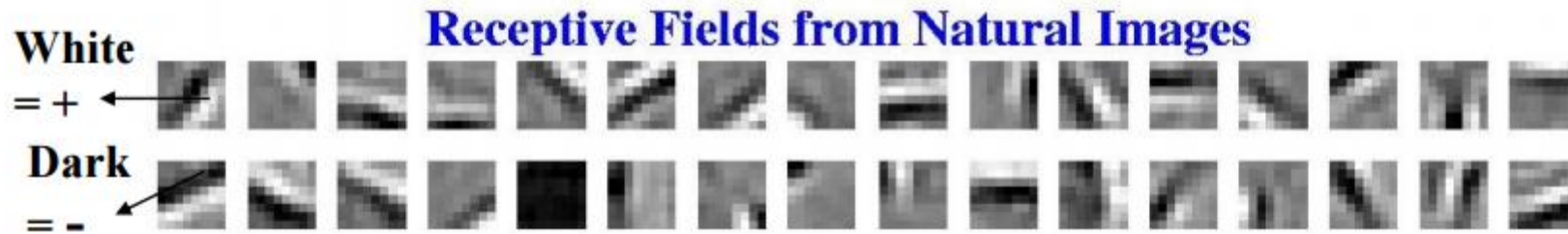
Mammals vision system



Mammals vision system

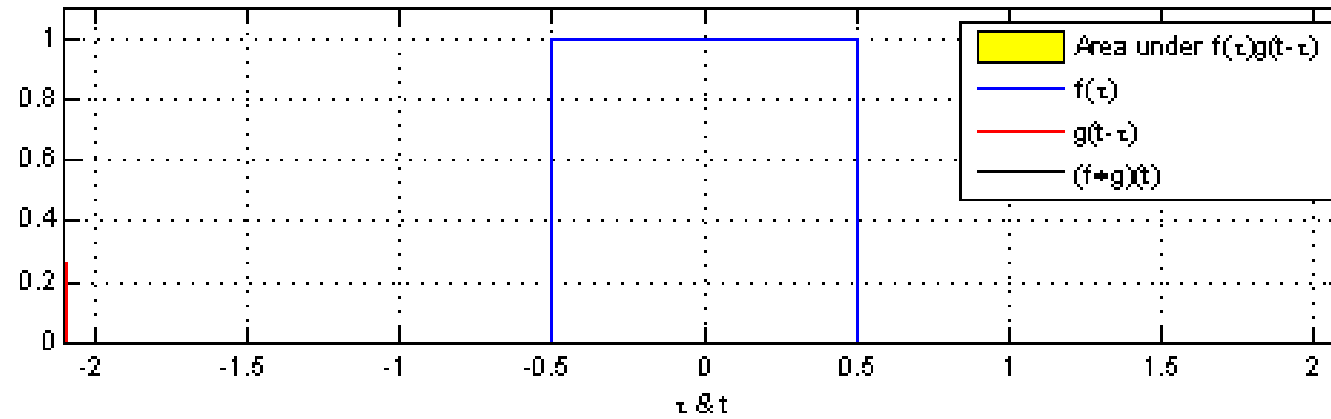


Mammals vision system



Convolution

$$\begin{aligned}(f * g)(t) &\stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau \\ &= \int_{-\infty}^{\infty} f(t - \tau) g(\tau) d\tau.\end{aligned}$$



Convolution as kernel



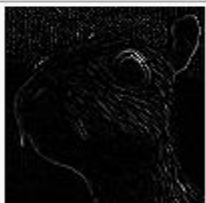

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

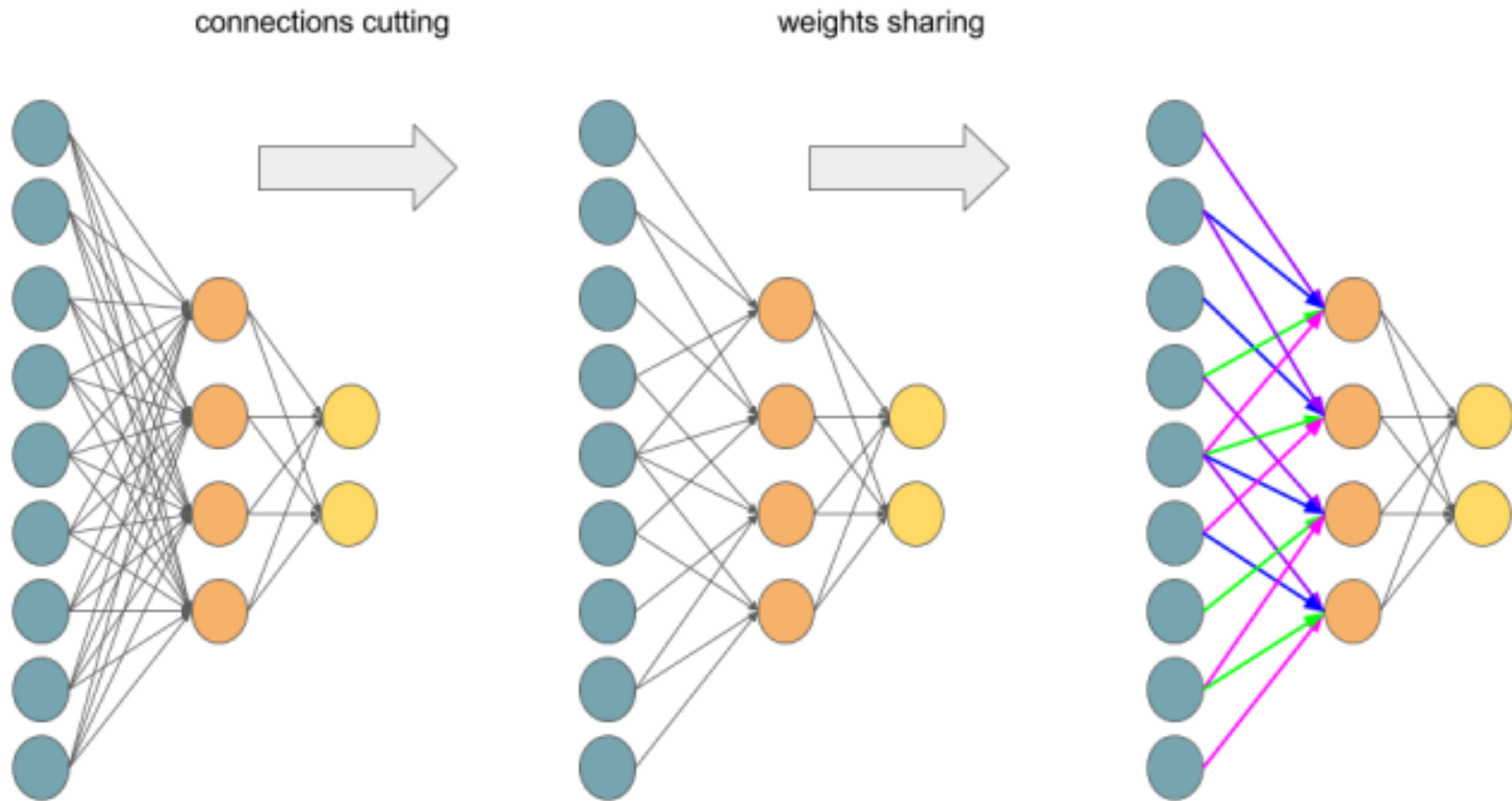
Convolved
Feature

Convolution as kernel

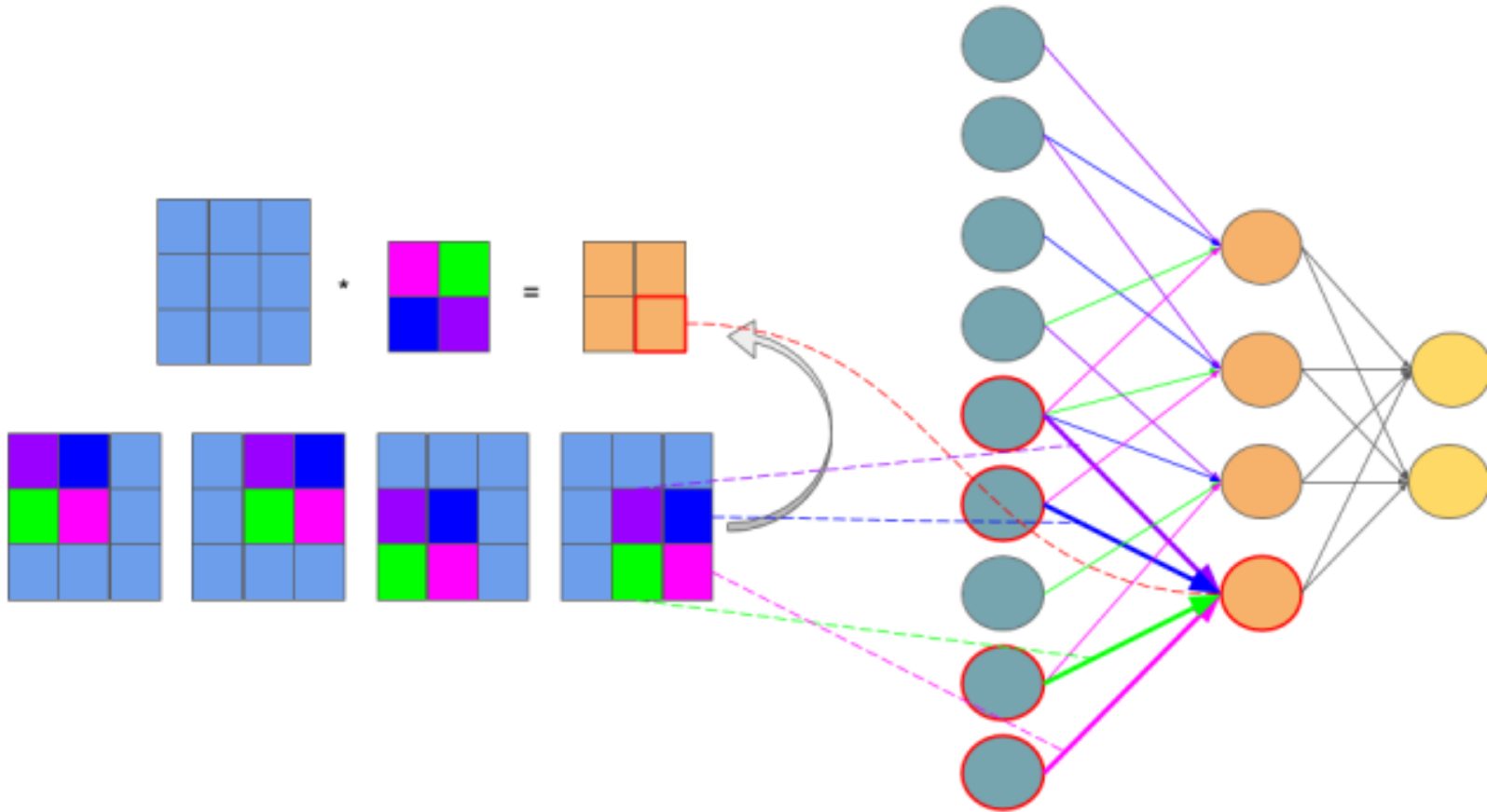
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	

Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	
5x5 Unsharp (with no image mask)	$\frac{-1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & -476 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$	

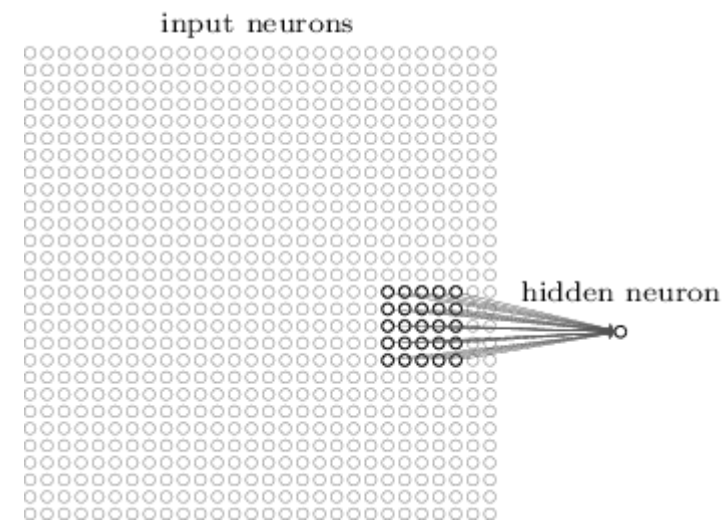
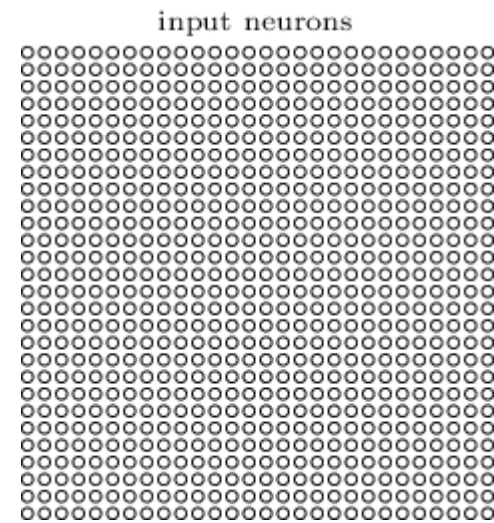
MLP \rightarrow CNN



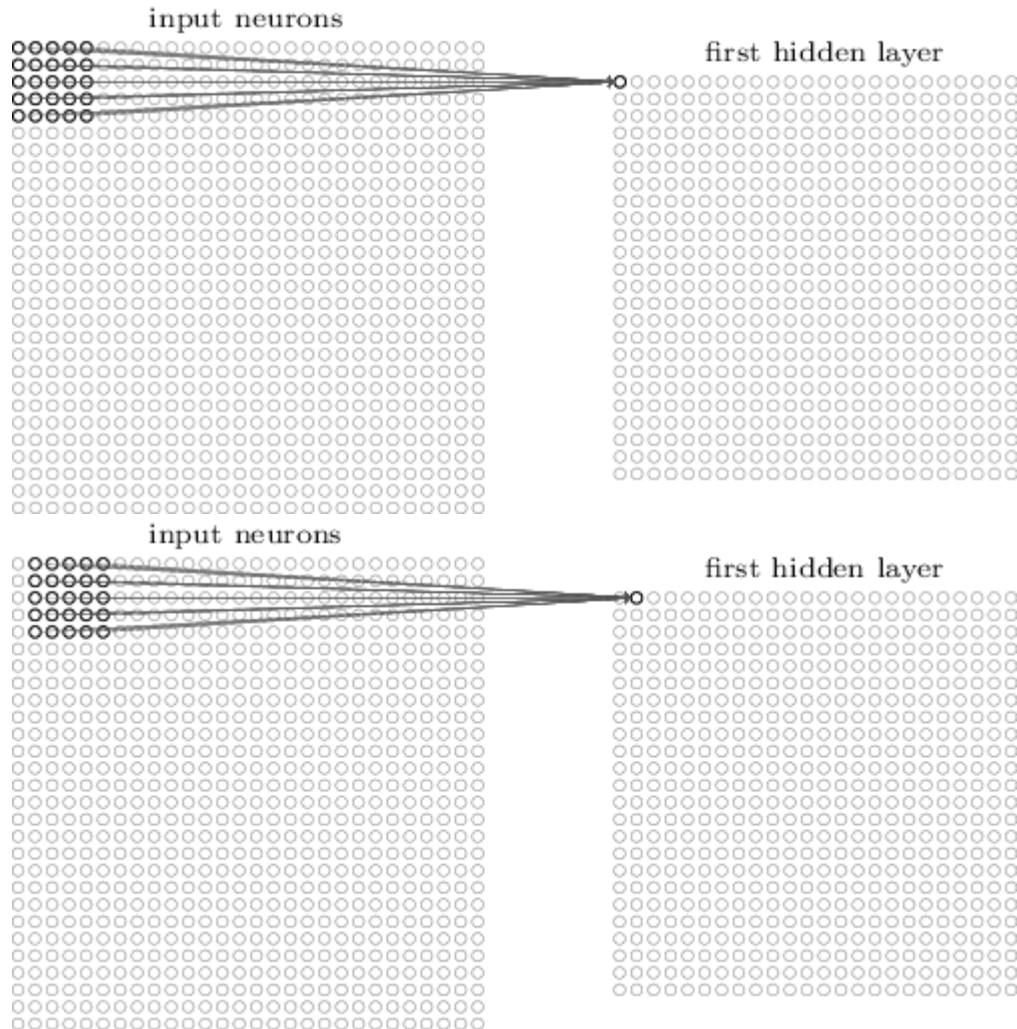
MLP -> CNN



CNN



CNN: convolutional layer



1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

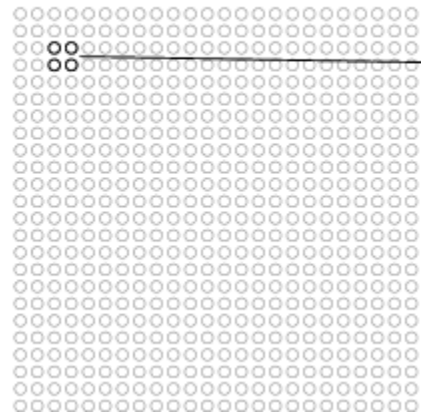
Image

4		

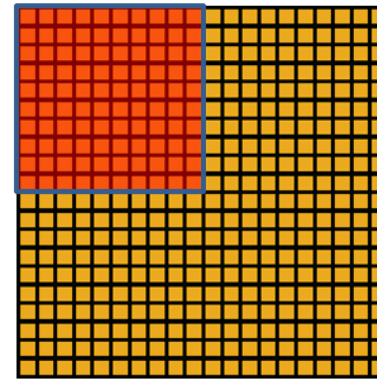
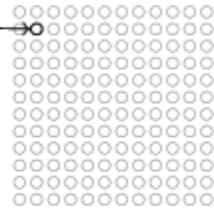
Convolved
Feature

CNN: pooling layer

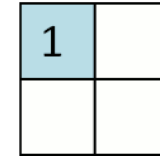
hidden neurons (output from feature map)



max-pooling units



Convolved
feature



Pooled
feature

CNN:building block

- Convolution (linear activations)
 - ReLU (detection stage)
- Pooling (modify output for next layer)

<https://www.youtube.com/watch?v=AgkfIQ4IGaM&t=141s>

<https://www.youtube.com/watch?v=GhqOMJIHD8A>

cnn backprobagation

<http://www.simon-hohberg.de/2014/10/10/conv-net.html>

CNN backpropagation

$$v(x, y) = \sum_{x_k=0}^{N_k-1} \sum_{y_k=0}^{N_k-1} i(x - x_k, y - y_k) k(x_k, y_k)$$

$$l(x, y) = \sum_{x_k=0}^{N_k-1} \sum_{y_k=0}^{N_k-1} i(x + x_k, y + y_k) k(x_k, y_k)$$

CNN backpropagation

$$i_{n,m}^l(x, y) = \sum o_n^{l-1}(x - x', y - y') \cdot w_{n,m}^l(x', y')$$

$$\begin{aligned} c_m^l(x, y) &= \sum_n i_{n,m}^l(x, y) + b_m^l \\ &= \sum_{n, x', y'} o_n^{l-1}(x - x', y - y') \cdot w_{n,m}^l(x', y') + b_m^l \end{aligned}$$

$$o_m^l(x, y) = a(c_m^l(x, y))$$

CNN weight update

$$\frac{\partial E}{\partial w_{n,m}^l(x, y)} = \sum_{x', y'} \underbrace{\frac{\partial E}{\partial c_m^l(x', y')}}_{\delta_m^l(x', y')} \cdot \frac{\partial c_m^l(x', y')}{\partial w_{n,m}^l(x, y)}$$

$$\begin{aligned} \frac{\partial c_m^l(x', y')}{\partial w_{n,m}^l(x, y)} &= \frac{\partial}{\partial w_{n,m}^l(x, y)} \left(\sum_{n, x'', y''} o_n^{l-1}(x' - x'', y' - y'') \cdot w_{n,m}^l(x'', y'') + b_m^l \right) \\ &= \frac{\partial}{\partial w_{n,m}^l(x, y)} \left(o_0^{l-1}(x' - 0, y' - 0) \cdot w_{0,m}^l(0, 0) + \dots + o_n^{l-1}(x' - x, y' - y) \cdot w_{n,m}^l(x, y) + \dots + b_m^l \right) \\ &= o_n^{l-1}(x' - x, y' - y) \end{aligned}$$

$$\frac{\partial E}{\partial w_{n,m}^l(x, y)} = \sum_{x', y'} o_n^{l-1}(x' - x, y' - y) \cdot \delta_m^l(x', y')$$

CNN weight update

$$\frac{\partial E}{\partial w_{n,m}^l(x, y)} = \text{rot}_{180} \left(\underbrace{\sum_{x', y'} o_n^{l-1}(x + x', y + y') \cdot \delta_m^l(x', y')}_{\text{Cross-Correlation}} \right)$$

CNN delta calculation

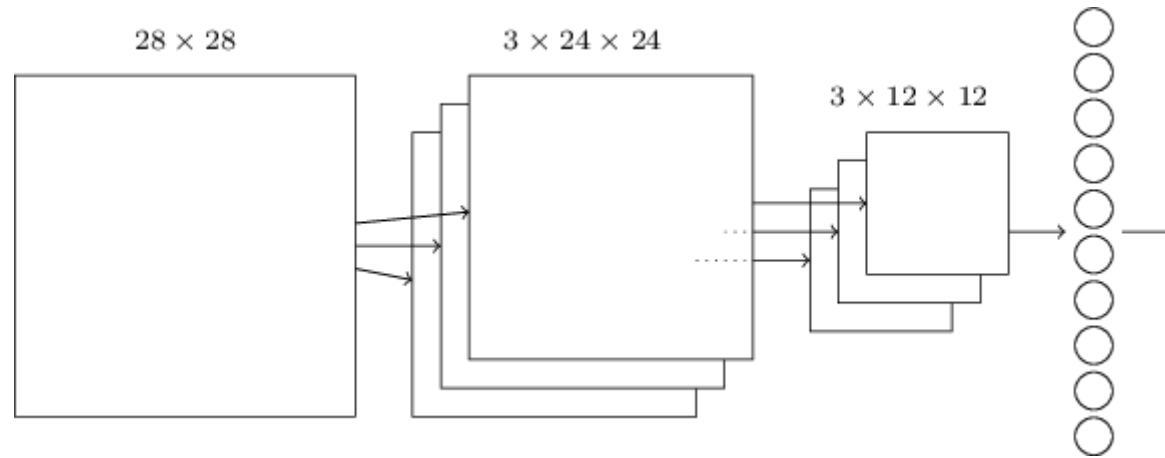
$$\delta_m^l(x, y) = \frac{\partial E}{\partial c_m^l(x, y)} = \sum_o \sum_{x', y'}^{N_w, N_w} \underbrace{\frac{\partial E}{\partial c_o^{l+1}(x + x', y + y')}}_{\delta_o^{l+1}(x+x', y+y')} \cdot \frac{\partial c_o^{l+1}(x + x', y + y')}{\partial c_m^l(x, y)}$$

$$\begin{aligned} \frac{\partial c_o^{l+1}(x + x', y + y')}{\partial c_m^l(x, y)} &= \frac{\partial}{\partial c_m^l(x, y)} \left(\sum_{m'} \sum_{x'', y''}^{N_w, N_w} o_{m'}^l(x + x' - x'', y + y' - y'') \cdot w_{m',o}^{l+1}(x'', y'') \right) \\ &= \frac{\partial}{\partial c_m^l(x, y)} \left(\sum_{m'} \sum_{x'', y''}^{N_w, N_w} a(c_{m'}^l(x + x' - x'', y + y' - y'')) \cdot w_{m',o}^{l+1}(x'', y'') \right) \\ &= \frac{\partial}{\partial c_m^l(x, y)} \left(w_{0,o}^{l+1}(0, 0) \cdot a(c_0^l(x' - 0, y' - 0)) + \dots + w_{m,o}^{l+1}(x'', y'') \cdot a(c_m^l(x, y)) + \dots \right) \\ &= w_{m,o}^{l+1}(x + x', y + y') \cdot \frac{\partial a(c_m^l(x, y))}{\partial c_m^l(x, y)} \end{aligned}$$

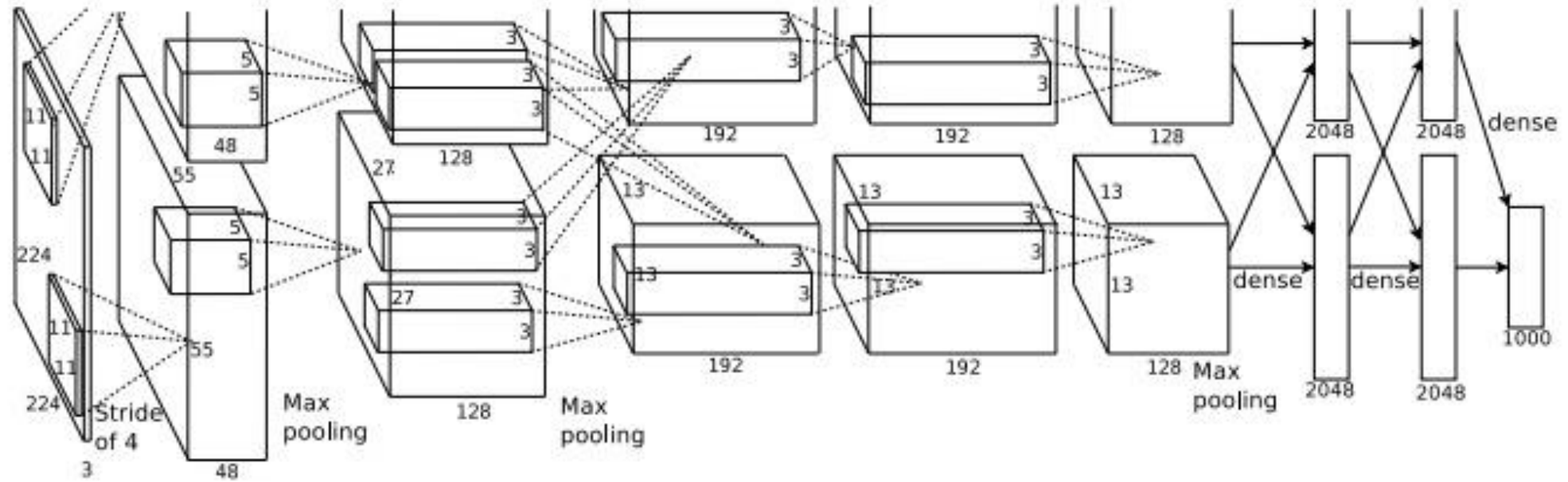
CNN delta calculation

$$\delta_m^l(x, y) = \underbrace{\sum_o \sum_{x', y'}^{N_w, N_w} \delta_o^{l+1}(x + x', y + y') \cdot w_{m,o}^{l+1}(x', y')}_{\text{Cross-Correlation, Backpropagated Error}} \cdot \frac{\partial a(c_m^l(x, y))}{\partial c_m^l(x, y)}$$

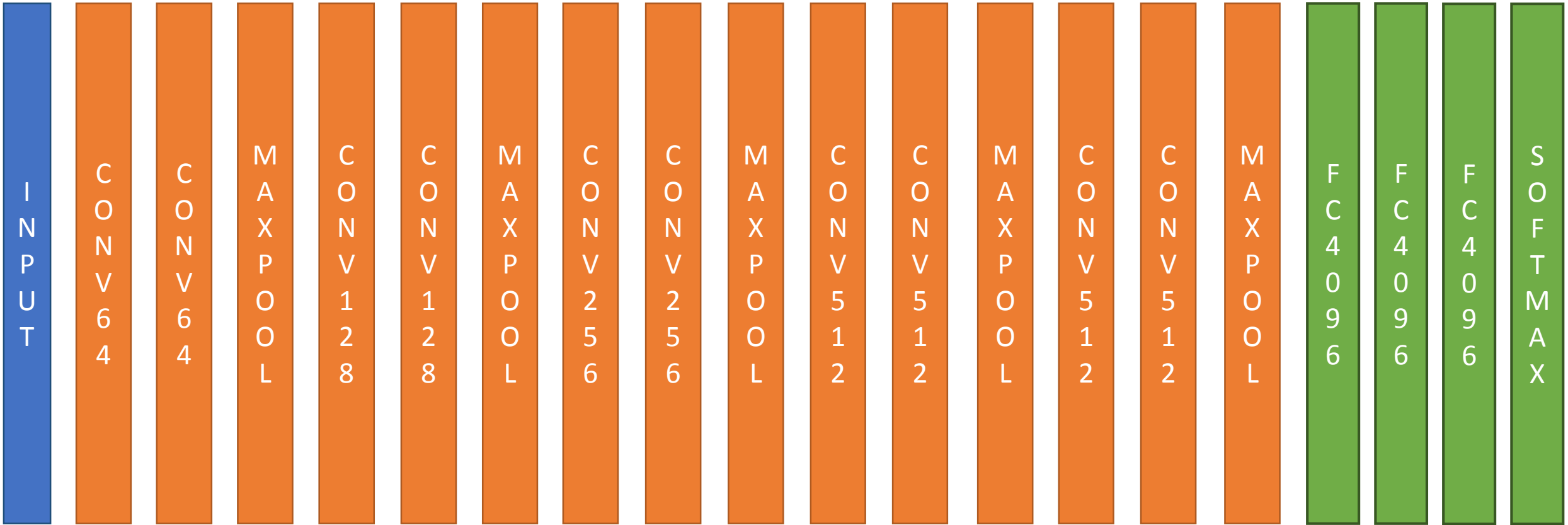
CNN: basic structure



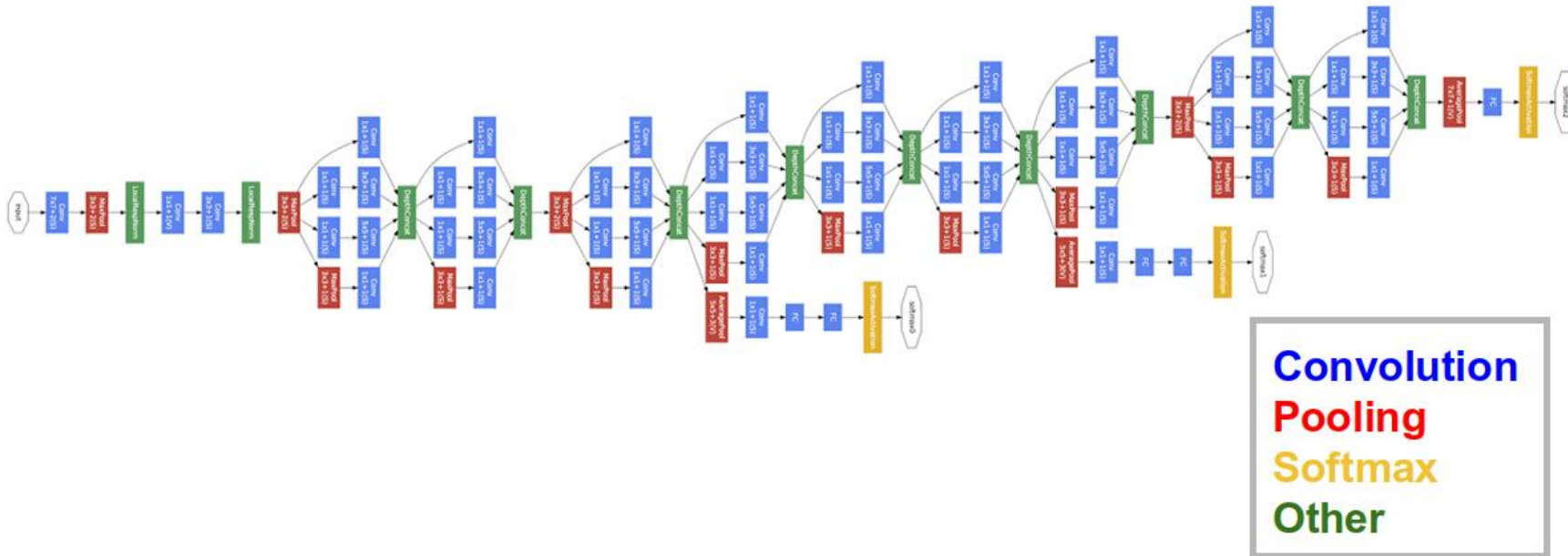
CNN: AlexNet 2012



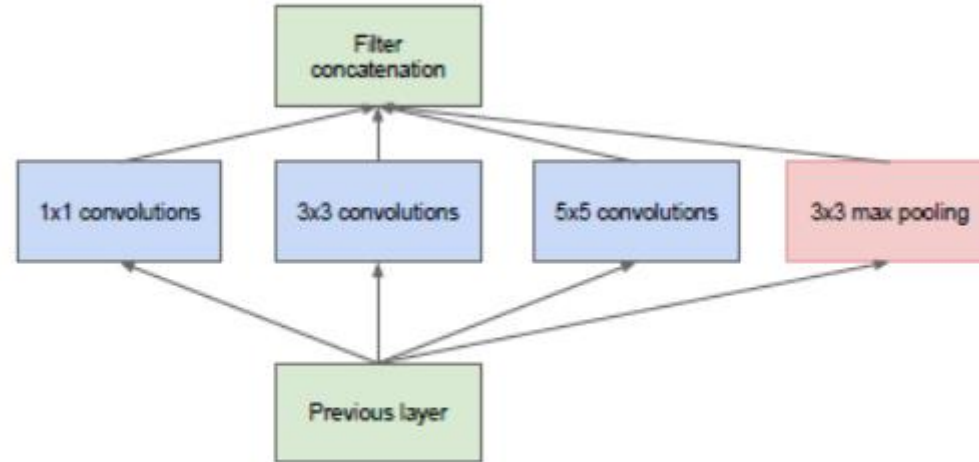
CNN: VGG Oxford 2014



CNN: GoogleNet 2014

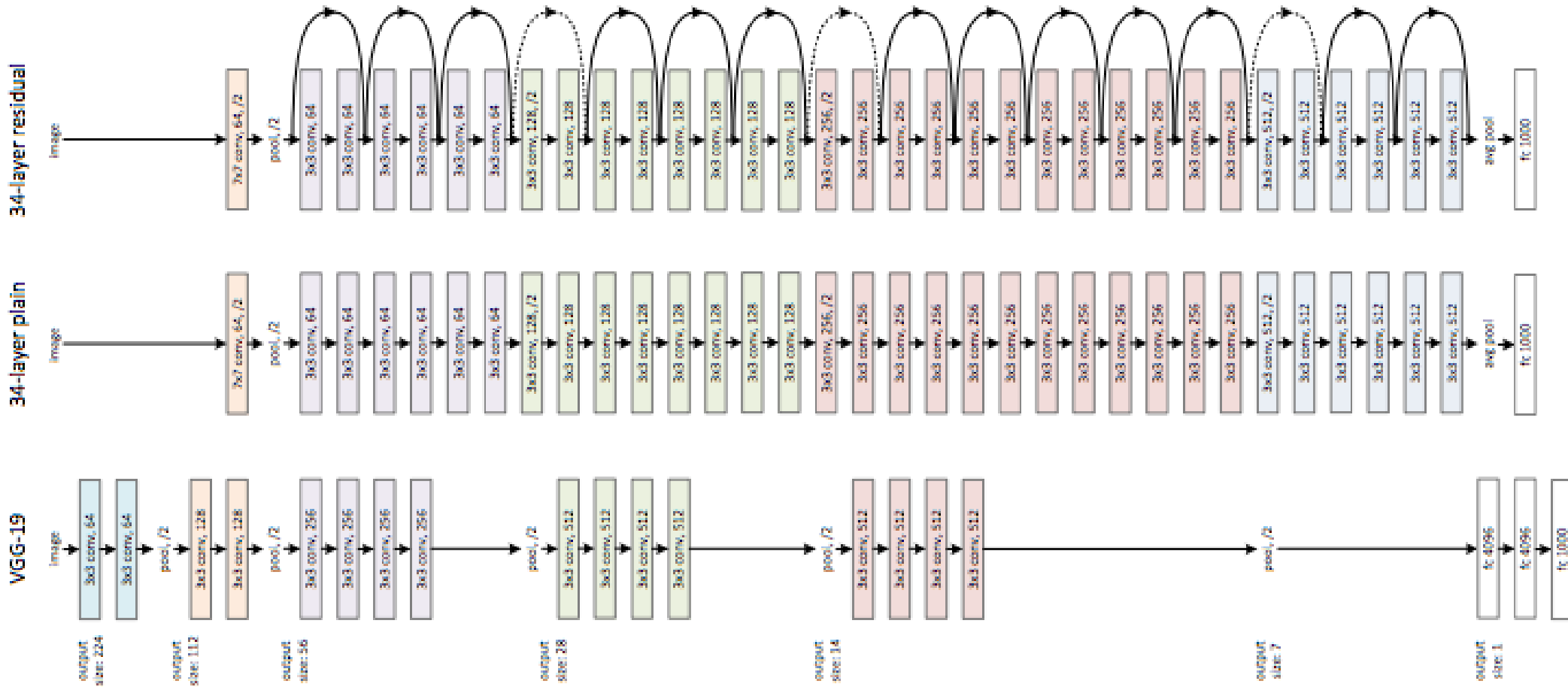


CNN: GoogleNet 2014



(a) Inception module, naïve version

CNN: ResNet 2015



CNN: ResNet 2015

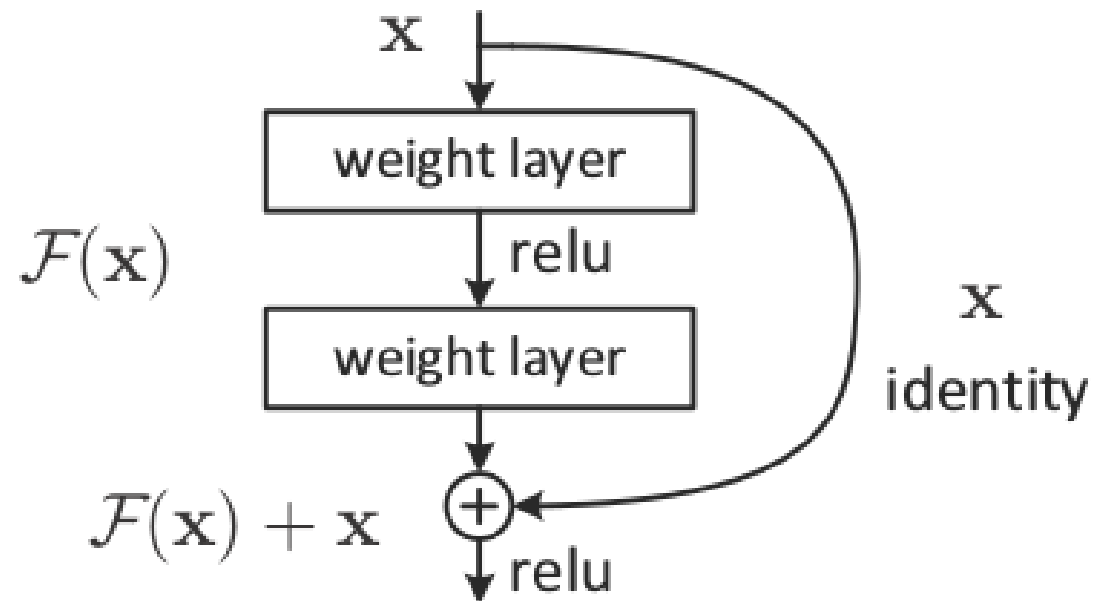
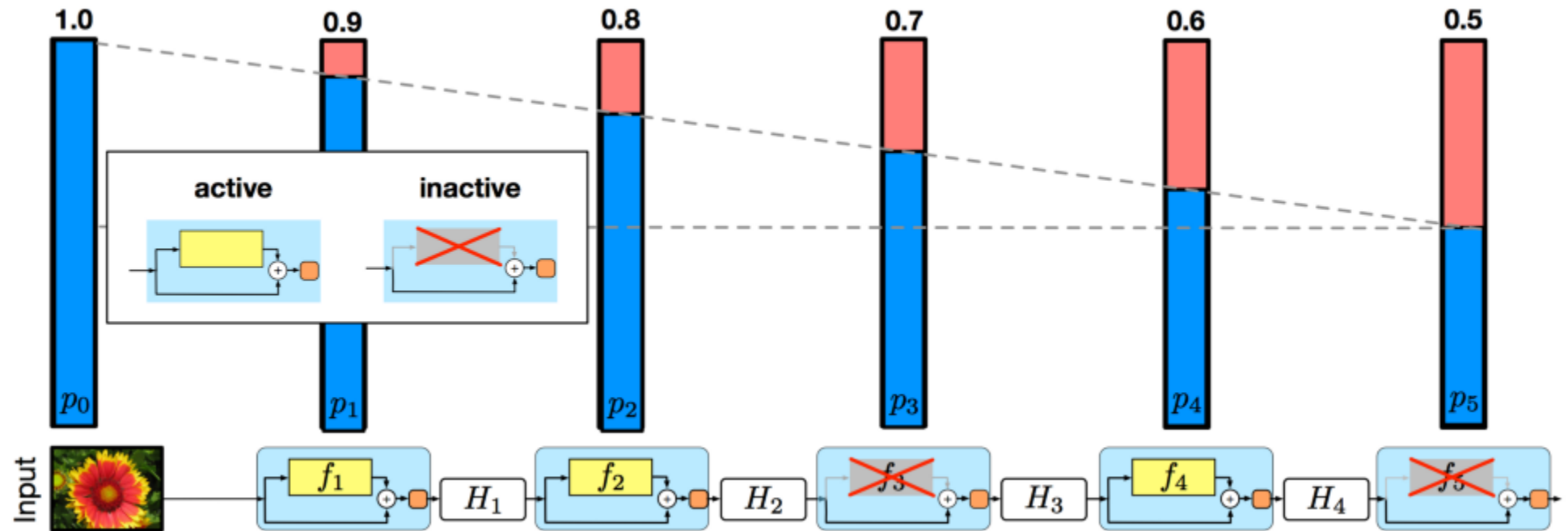
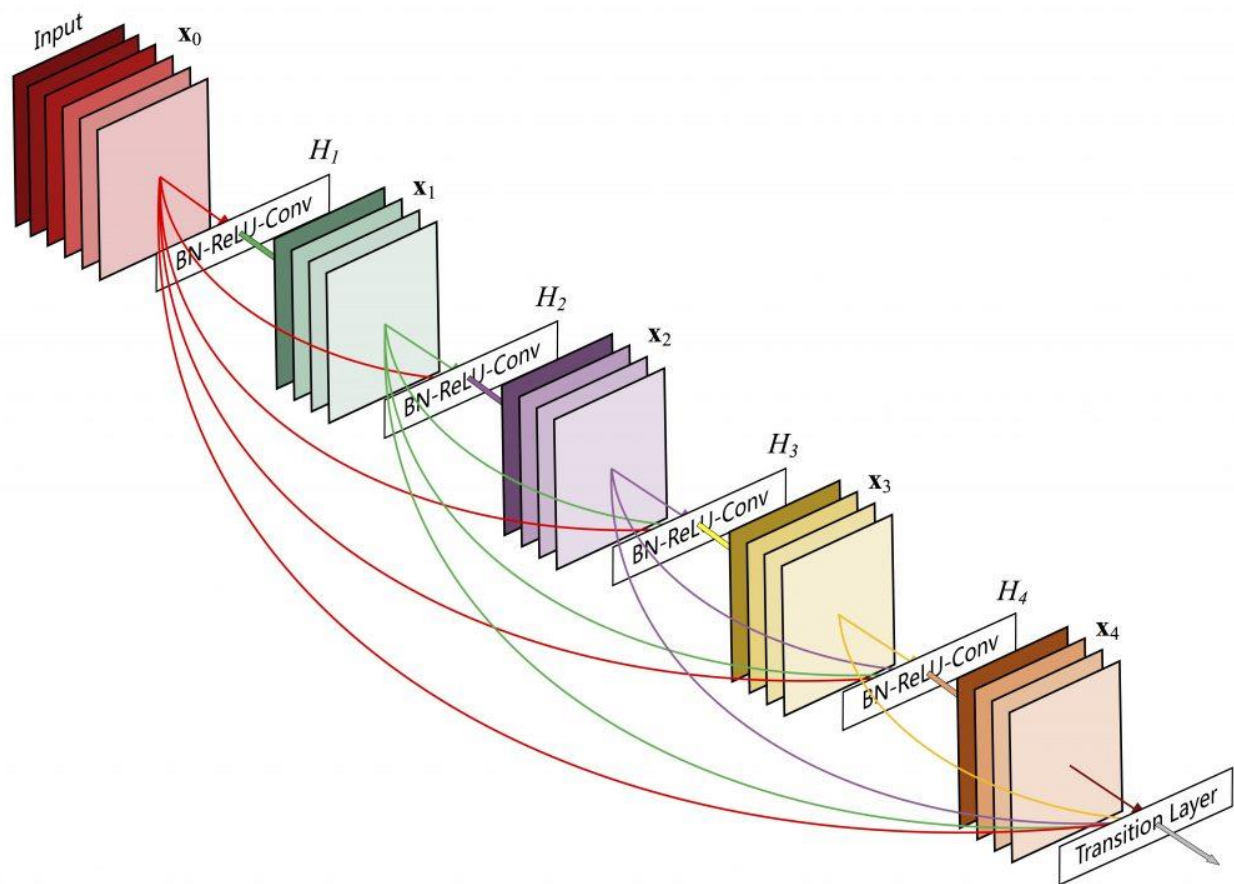


Figure 2. Residual learning: a building block.

CNN: Stochastic depth 2016



CNN: DenseNet 2016



CNN: Recent advances

- Batch normalization
- Different convolution types (dilated etc)
- Knowledge transfer
- Multi-task learning
- Convolution as FFT
- Several outputs
- Weights pruning / quantization

CNN: applications

- Classification (any object)
- Localization (any object)
- Segmentation (any object)
- Identification (face for example)
- Pose estimation
- Image captioning