

Network structures: CNN, ResNet, DenseNet

Shikui Tu

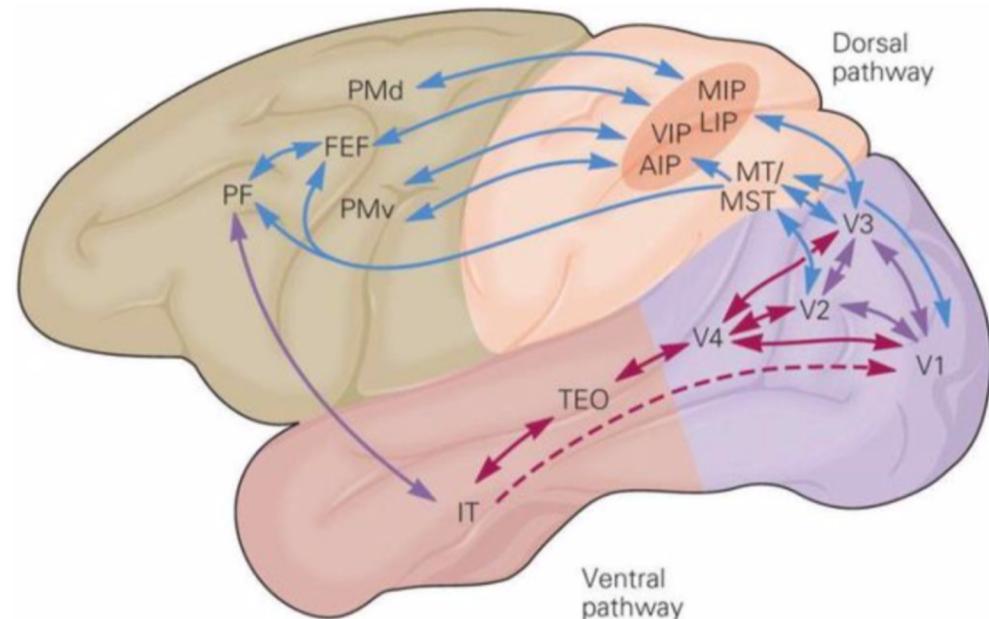
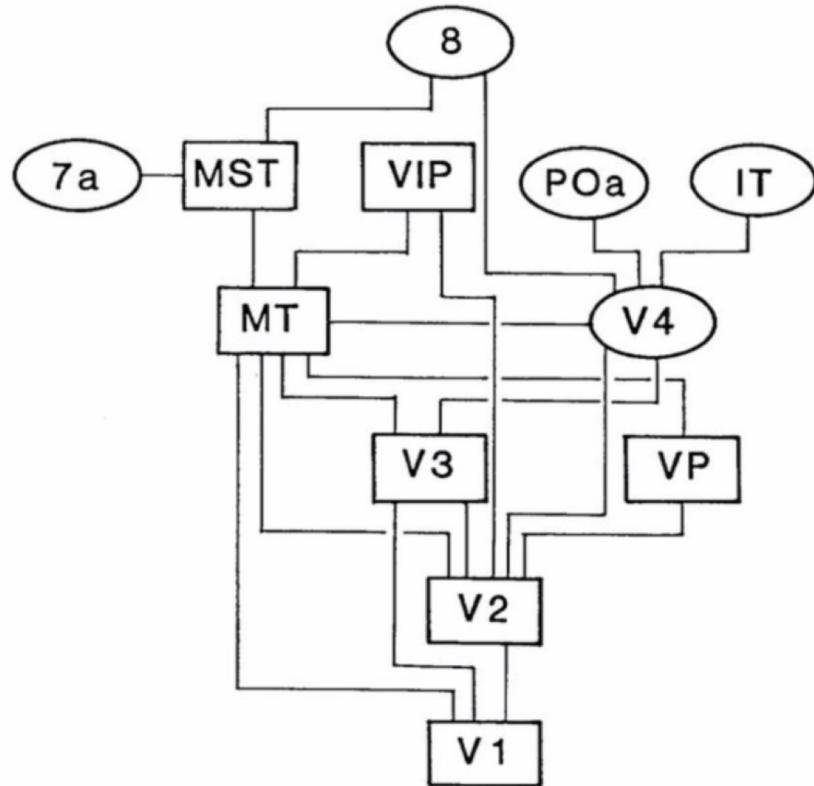
Department of Computer Science and Engineering,
Shanghai Jiao Tong University

2019-05-21

Outline

- Background introduction
- Convolutional neural networks (CNN)
 - Overview
 - Network details: convolution, pooling, activation, loss functions
- ResNet
- DenseNet
- RNN

Visual cortex



Starting from V1 primary visual cortex, visual signal is transmitted upwards, becoming more complicated and abstract.

视觉信息处理过程

- 生物视觉信息处理的大致过程
 - -0ms: 光刺激发生
 - -10ms: 视网膜光电转换和编码
 - -35ms: 视网膜“模数”转换和神经脉冲
 - -45ms: 外侧膝状体(LGN)中继转发
 - -55ms: V1 方向选择细胞响应
 - -85ms: V4 大范围整合
 - -100ms: 人脸选择细胞响应
 - -160-220ms: 对象分类识别

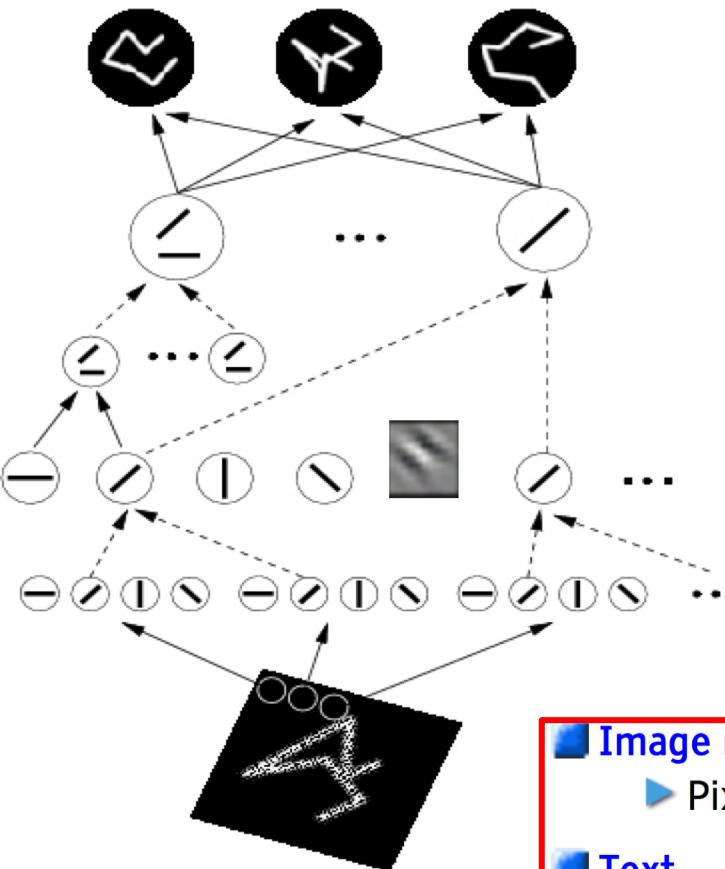
“只因在人群中多看了你一眼”

Maunsell and Gibson 1992; Raiguel et al. 1989; Nowak et al. 1995;
Schmolesky et al. 1998; Thorpe, Fize& Marlot 1996

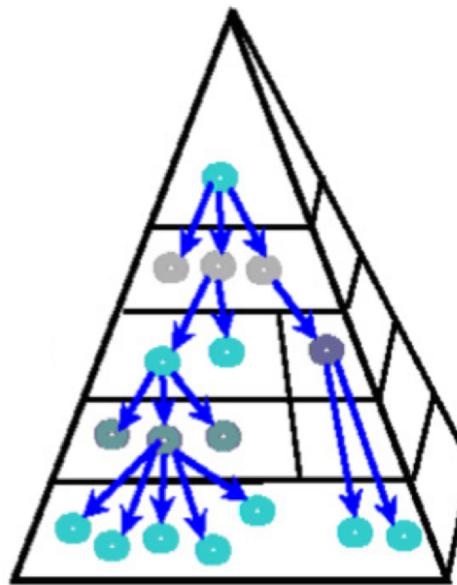
Feature Detection Theory 特征检测理论 (1960)

Hubel

Wiesel



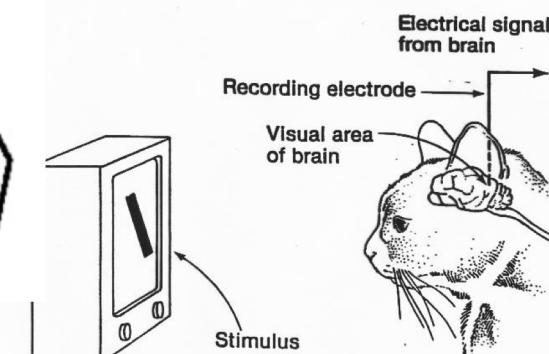
具备广泛的普遍性！



David H. Hubel

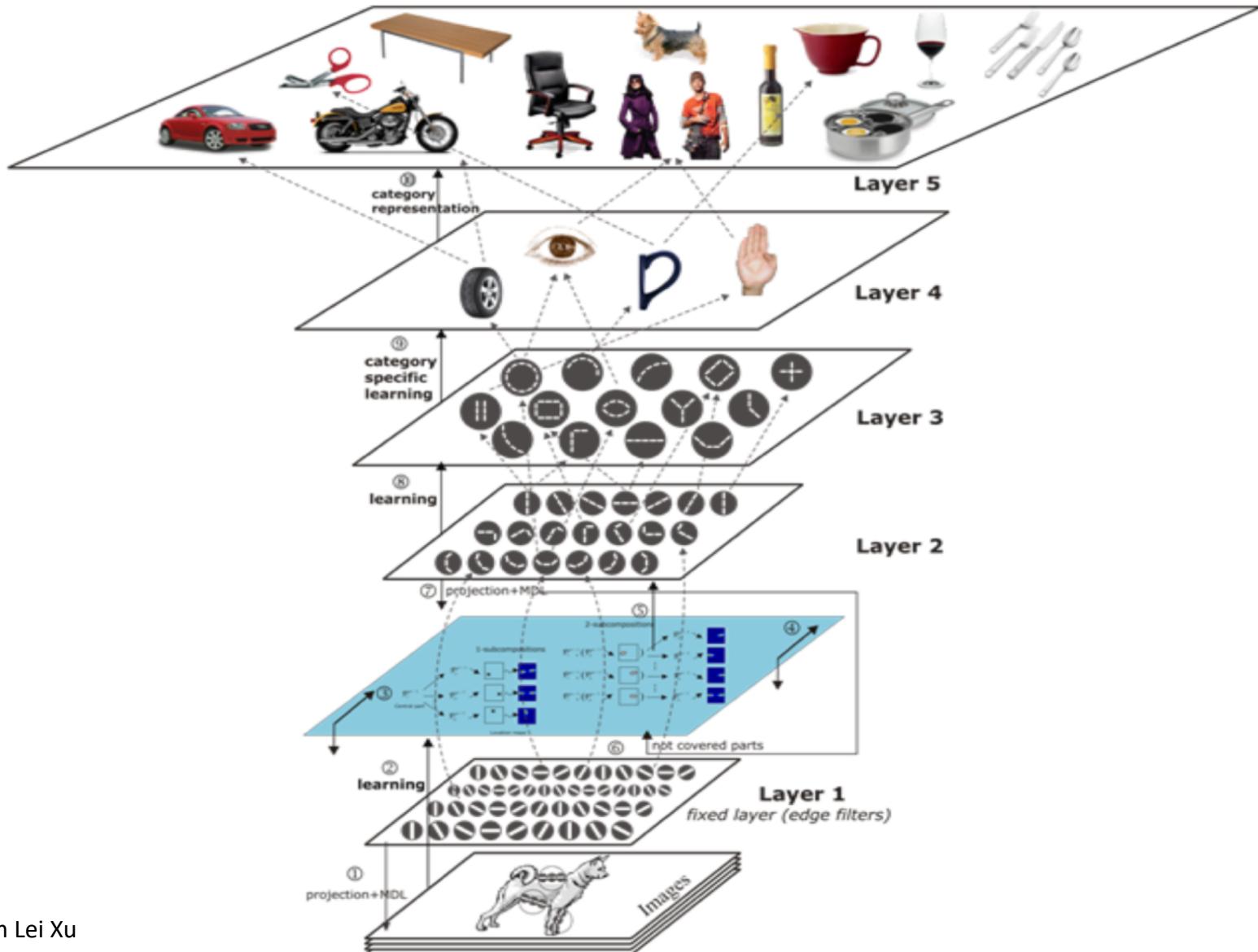


Torsten N. Wiesel

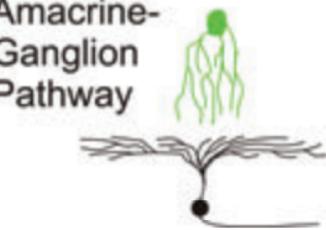
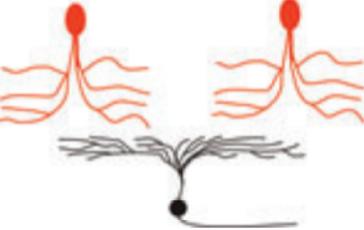
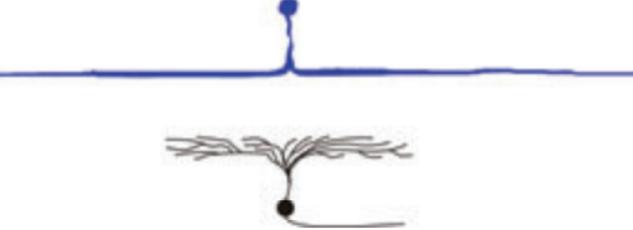


- Image recognition**
 - ▶ Pixel → edge → texton → motif → part → object
- Text**
 - ▶ Character → word → word group → clause → sentence → story
- Speech**
 - ▶ Sample → spectral band → sound → ... → phone → phoneme → word →

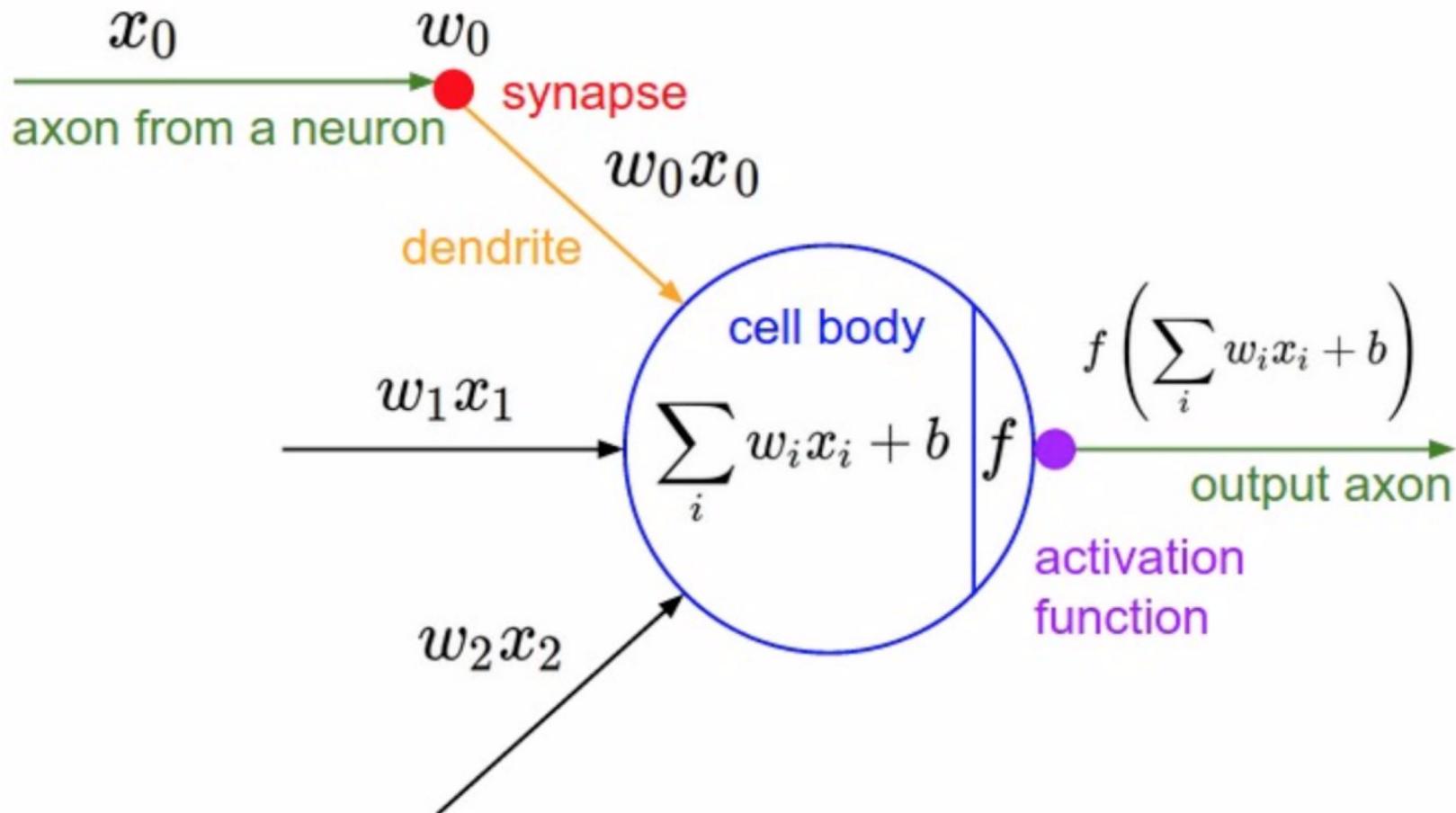
Hierarchical representations



narrow field, medium field and wide field amacrine cells

	A Narrow Field	B Medium Field	C Wide Field
Morphology			
Transmitter	Glycine	GABA	GABA
Process Spread	\leftrightarrow $<100 \text{ um}$ 	\leftrightarrow $\sim200 \text{ um}$ 	\leftrightarrow $>1000 \text{ um}$ 
Spatial Influence	$<100 \text{ um}$	$\sim200 \text{ um}$	$>1000 \text{ um}$
Latency	$<160 \text{ msec}$	$>160 \text{ msec}$	$<100 \text{ msec}$
Amacrine-Ganglion Pathway			

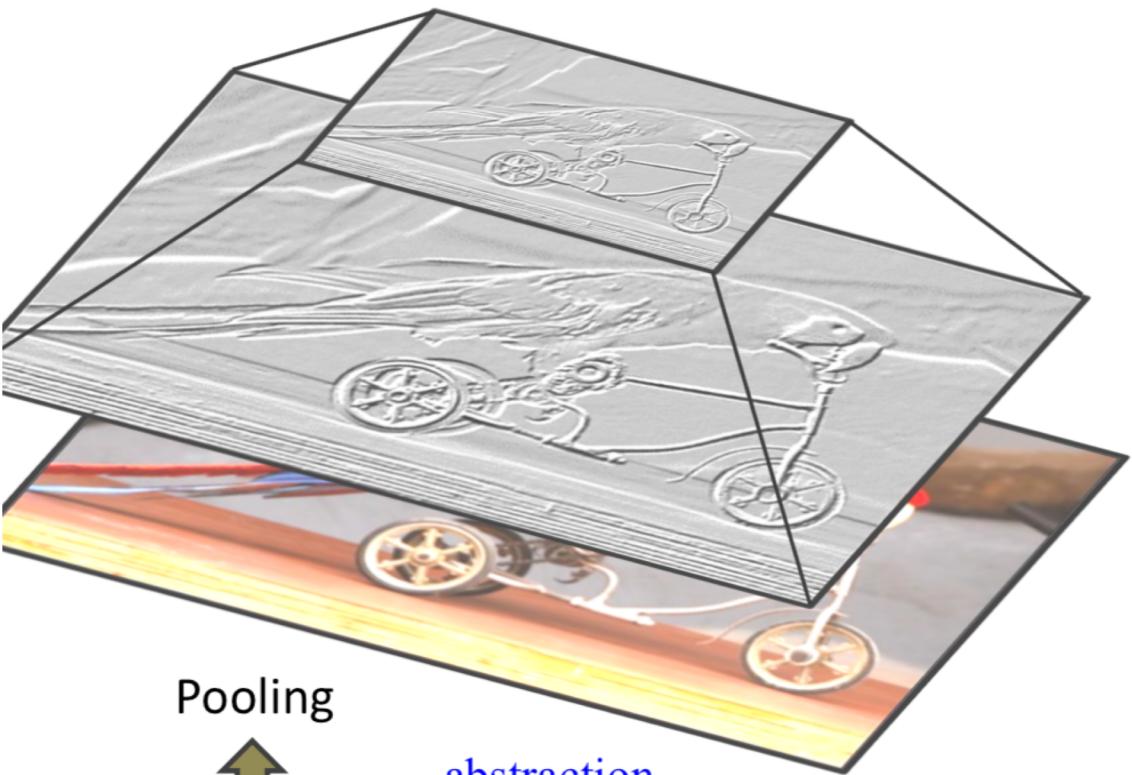
The neuron model



Outline

- Background introduction
- **Convolutional neural networks (CNN)**
 - Overview
 - Network details: convolution, pooling, activation, loss functions
- ResNet
- DenseNet
- RNN

A Basic Module of the Convolutional Neural Nets



Convolution

summary



Image



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)

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



Basic functions: constant, successor, projection

$$f(x_1, \dots, x_k) = n, \quad f(x) = x + 1, \quad P_i(x_1, \dots, x_k) = x_i$$

Composition operator

$$h \circ (g_1, \dots, g_m) \stackrel{\text{def}}{=} f(x_1, \dots, x_k) = h(g_1(x_1, \dots, x_k), \dots, g_m(x_1, \dots, x_k)).$$

Primitive recursion operator $\rho(g, h) \stackrel{\text{def}}{=} f(y, x_1, \dots, x_k)$ where

$$f(0, x_1, \dots, x_k) = g(x_1, \dots, x_k)$$

$$f(y + 1, x_1, \dots, x_k) = h(y, f(y, x_1, \dots, x_k), x_1, \dots, x_k)$$

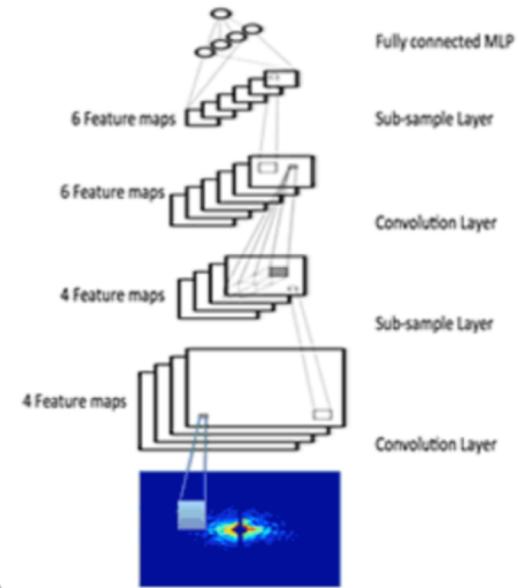
Minimization operator

Given a $(k+1)$ -ary total function $f(y, x_1, \dots, x_k)$

$$\mu(f)(x_1, \dots, x_k) = z \stackrel{\text{def}}{\iff} f(z, x_1, \dots, x_k) = 0 \quad \text{and}$$

$$f(i, x_1, \dots, x_k) > 0 \quad \text{for } i = 0, \dots, z - 1.$$

$$div(x, y) = \begin{cases} \text{integer portion of } x/y \text{ if } y \neq 0 \\ \text{undefined if } y = 0 \end{cases}$$



$$f(x, 0) = g(x)$$

$$f(x, y + 1) = h(x, y, f(x, y))$$

$$h(x, y, x) = z + x^y x$$

$$f(\bar{x}, 3) = h(\bar{x}, 2, f(\bar{x}, 2))$$

$$f(\bar{x}, 2) = h(\bar{x}, 1, f(\bar{x}, 1))$$

$$f(\bar{x}, 1) = h(\bar{x}, 0, f(\bar{x}, 0))$$

$$f(\bar{x}, 0) = g(\bar{x})$$

Early CNN: LeNet-5

LeNet-5, a pioneering 7-level convolutional network by LeCun et al. in 1998, that classifies digits, was applied by several banks to recognize hand-written numbers on checks digitized in 32x32 pixel images.

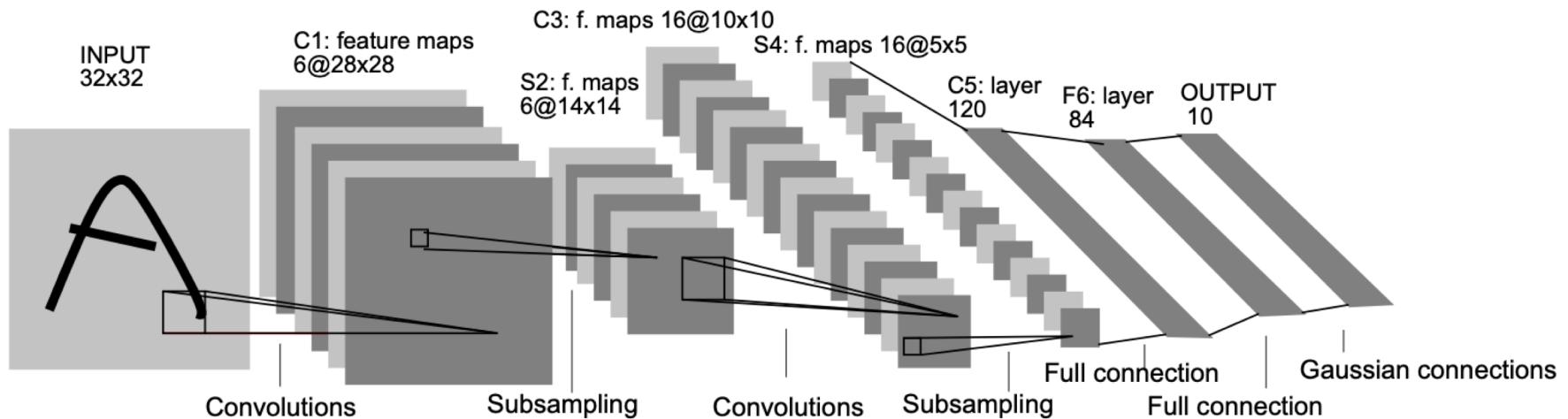
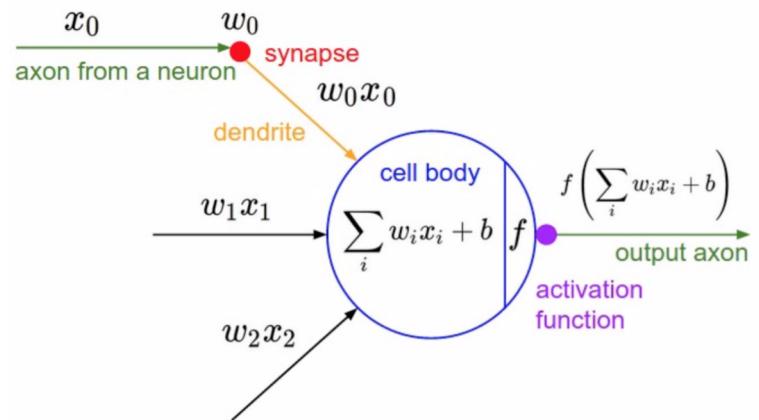


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Affine transformations

- A vector is received as input and is multiplied with a matrix to produce an output (to which a bias vector is usually added before passing the result through a non-linearity).
- Applicable to any type of input
 - flattened into a vector
 - Image, sound clip, etc.



Affine transformations are not good for image-like structured data

- Images, sound clips and many other similar kinds of data have an intrinsic structure.
- They are stored as **multi-dimensional** arrays.
- They feature one or more axes for which **ordering matters** (e.g., width and height axes for an image, time axis for a sound clip).
- One axis, called the channel axis, is used to access **different views** of the data (e.g., the red, green and blue channels of a color image, or the left and right channels of a stereo audio track).

Discrete convolution

Kernel (3x3)

0	1	2
2	2	0
0	1	2

Input feature map

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

Output feature map

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

$$\begin{aligned} & 0 \times 3 + 1 \times 3 + 2 \times 2 \\ & + 2 \times 0 + 2 \times 0 + 0 \times 1 \\ & + 0 \times 3 + 1 \times 1 + 2 \times 2 \\ & = 12.0 \end{aligned}$$

Discrete convolution

Kernel (3x3)

0	1	2
2	2	0
0	1	2

Input

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

output

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 ₀	2 ₁	1 ₂	0
0	0 ₂	1 ₂	3 ₀	1
3	1 ₀	2 ₁	2 ₂	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2 ₀	1 ₁	0 ₂
0	0	1 ₂	3 ₂	1 ₀
3	1	2 ₀	2 ₁	3 ₂
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 ₀	0 ₁	1 ₂	3	1
3 ₂	1 ₂	2 ₀	2	3
2 ₀	0 ₁	0 ₂	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0 ₀	1 ₁	3 ₂	1
3	1 ₂	2 ₂	2 ₀	3
2	0 ₀	0 ₁	2 ₂	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1 ₀	3 ₁	1 ₂
3	1	2 ₂	2 ₂	3 ₀
2	0	0 ₀	2 ₁	2 ₂
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 ₀	1 ₁	2 ₂	2	3
2 ₂	0 ₂	0 ₀	2	2
2 ₀	0 ₁	0 ₂	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 ₁	1 ₀	2 ₁	2 ₂	3
2	0 ₂	0 ₂	2 ₀	2
2	0 ₀	0 ₁	0 ₂	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3	1	2 ₀	2 ₁	3 ₂
2	0	0 ₂	2 ₂	2 ₀
2	0	0 ₀	0 ₁	1 ₂

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

N-D convolution (N=2,3,...)

axis $j = 1, 2, \dots, N$.

$n \equiv$ number of output feature maps,

$m \equiv$ number of input feature maps,

$k_j \equiv$ kernel size along axis j .

- i_j : input size along axis j ,
- k_j : kernel size along axis j ,
- s_j : stride (distance between two consecutive positions of the kernel) along axis j ,
- p_j : zero padding (number of zeros concatenated at the beginning and at the end of an axis) along axis j .

Example with zero paddings

for $N = 2$, $i_1 = i_2 = 5$, $k_1 = k_2 = 3$, $s_1 = s_2 = 2$, and $p_1 = p_2 = 1$.

0 ₀	0 ₁	0 ₂	0	0	0	0	0
0 ₂	3 ₂	3 ₀	2	1	0	0	0
0 ₀	0 ₁	0 ₂	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	1	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0 ₀	0 ₁	0 ₂	0	0	0
0	3	3 ₂	2 ₂	1 ₀	0	0	0
0	0	0 ₀	1 ₁	3 ₂	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	1

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0 ₀	0 ₁	0 ₂
0	3	3	2	1 ₂	0 ₂	0 ₀	0
0	0	0	1	3 ₀	1 ₁	0 ₂	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	1 ₂	0 ₂	0	2	2	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3 ₂	1	0	0
0	3	1 ₂	2 ₂	2 ₀	3	0	0
0	2	0 ₀	0 ₁	2 ₂	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3 ₀	1 ₁	0 ₂	0
0	3	1	2	2 ₂	3 ₂	0 ₀	0
0	2	0	0	2 ₀	2 ₁	0 ₂	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	1 ₂	0 ₂	0	2	2	0
0	2	2 ₂	0 ₀	0	0	1	0
0	0	1 ₂	0 ₂	0	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0 ₀	0 ₁	2 ₂	2	0	0
0	2	0 ₂	0 ₀	0 ₀	1	0	0
0	0	0 ₀	0 ₁	0 ₂	0	0	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

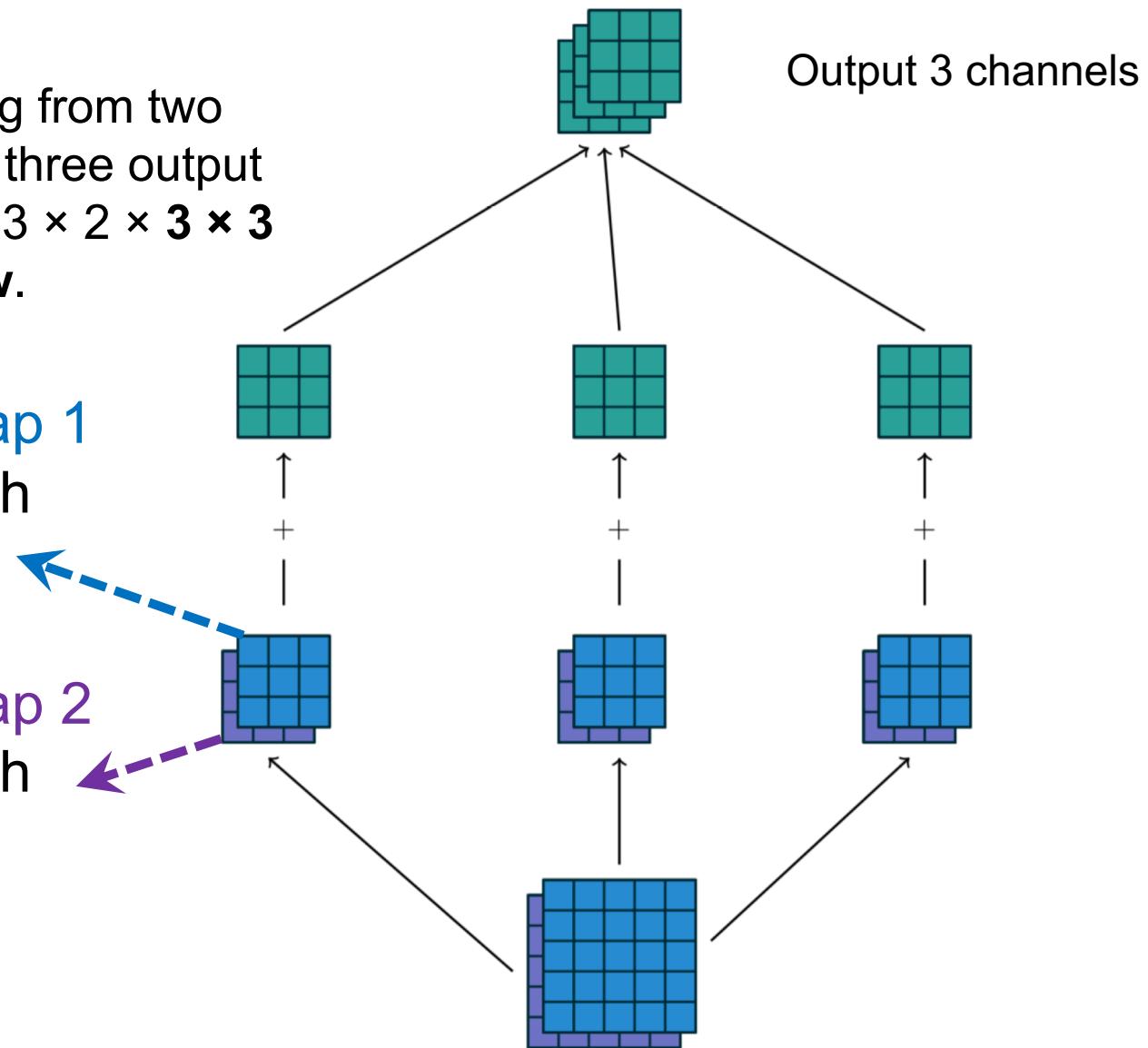
0	0	0	0	0	0	0	0
0	3	3	2	1	0	0	0
0	0	0	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2 ₀	2 ₁	0 ₂	0
0	2	0	0	0 ₂	1 ₂	0 ₀	0
0	0	0	0	0 ₀	0 ₁	0 ₂	0

6.0	14.0	17.0
14.0	12.0	12.0
8.0	10.0	17.0

Example with three output feature maps

A convolution mapping from two input feature maps to three output feature maps using a $3 \times 2 \times 3 \times 3$ collection of kernels w .

- input **feature map 1** is convolved with kernel $w_{1,1}$.
- input **feature map 2** is convolved with kernel $w_{1,2}$.



Pooling

- Pooling operations reduce the size of feature maps by using some function to summarize subregions, such as taking the **average** or the **maximum** value.

Average pooling

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

1.7	1.7	1.7
1.0	1.2	1.8
1.1	0.8	1.3

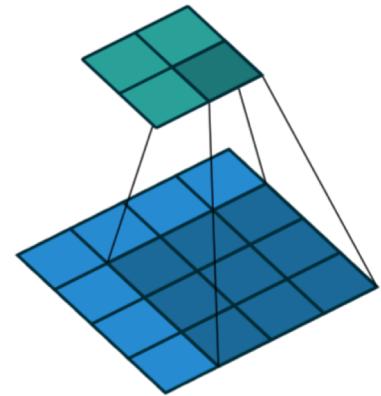
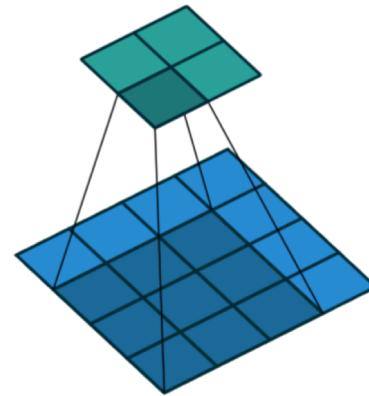
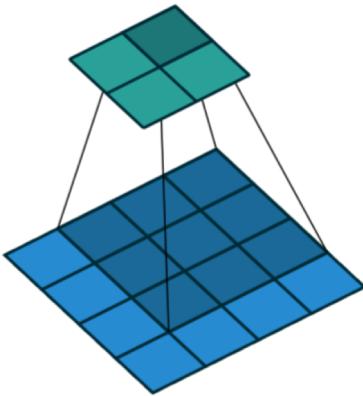
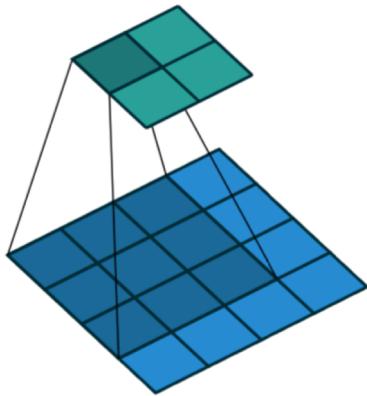
Maximum pooling

3	3	2	1	0
0	0	1	3	1
3	1	2	2	3
2	0	0	2	2
2	0	0	0	1

3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

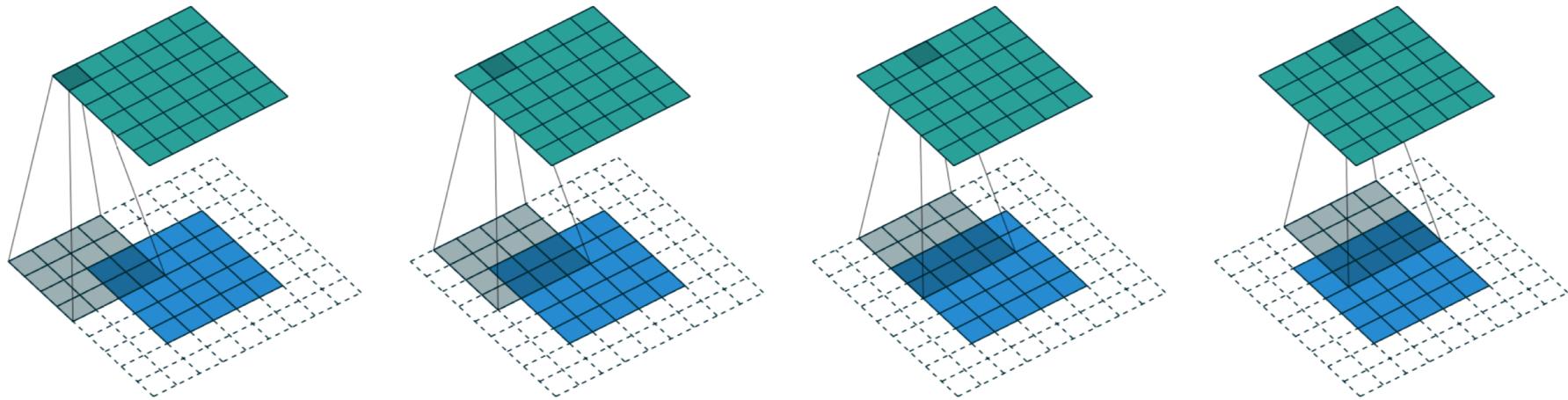
Example -- No padding, no strides

Convolving a 3×3 kernel over a 4×4 input using unit strides
(i.e., $i = 4$, $k = 3$, $s = 1$ and $p = 0$).

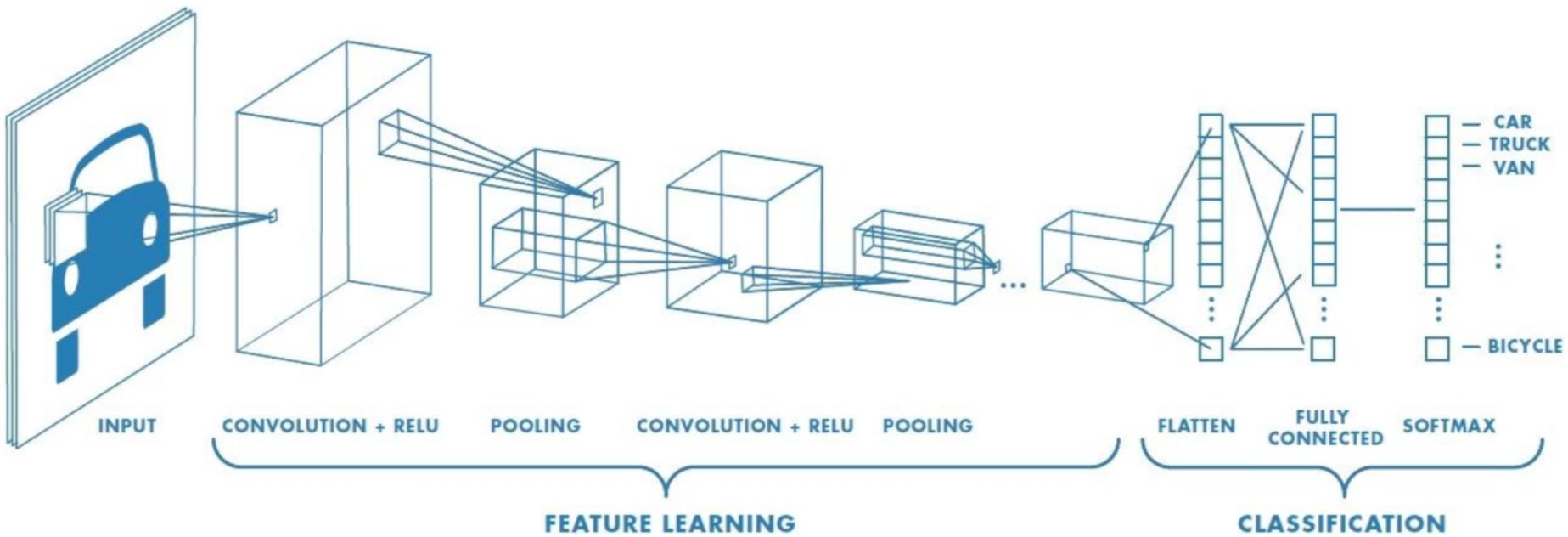


Example -- Arbitrary padding, no strides

Convolving a 4×4 kernel over a 5×5 input padded with a 2×2 border of zeros using unit strides (i.e., $i = 5$, $k = 4$, $s = 1$ and $p = 2$).



Sample architecture of CNN

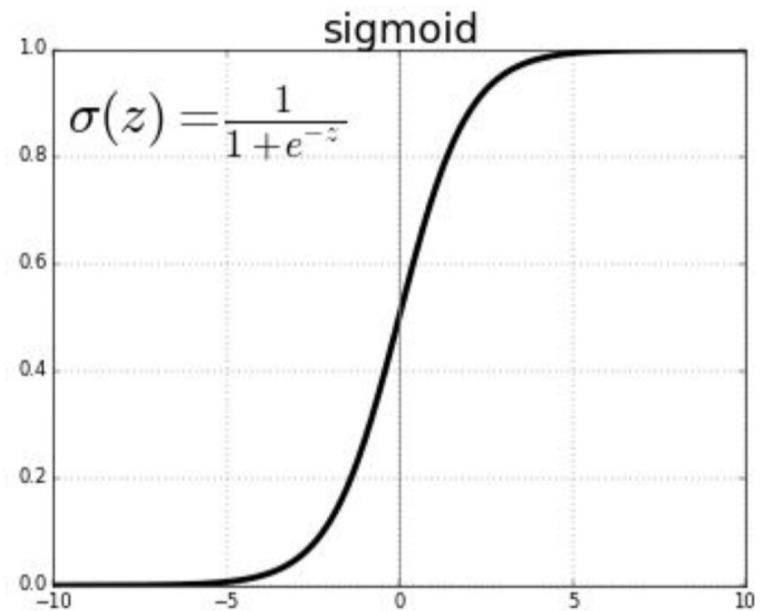
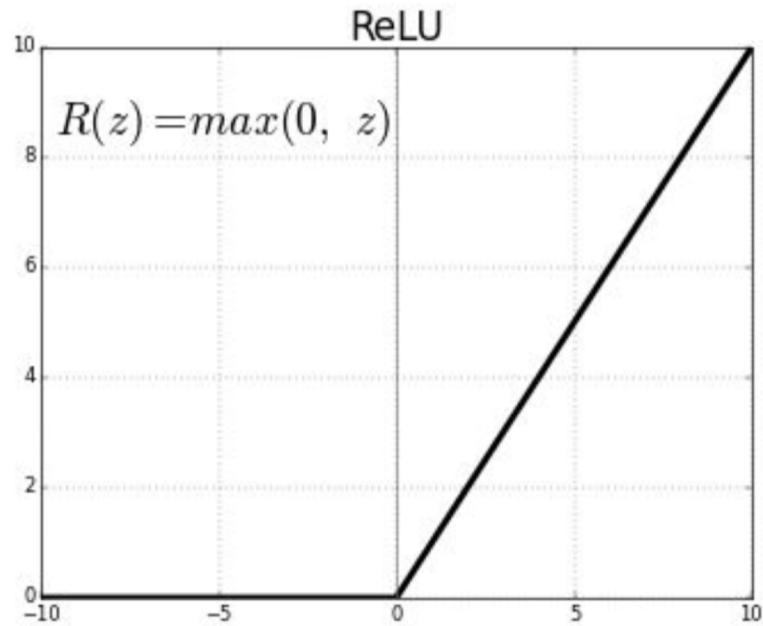


Activation layer

- Used to increase non-linearity of the network without affecting receptive fields of conv layers
- Prefer **ReLU**, results in faster training
- **LeakyReLU** addresses the vanishing gradient problem

Other types:

- Leaky ReLU,
- Randomized Leaky ReLU,
- Parameterized ReLU
- Exponential Linear Units (ELU),
- Scaled Exponential Linear Units Tanh,
- hardtanh, softtanh, softsign, softmax, softplus...



Softmax

- A special kind of activation layer, usually at the end of FC layer outputs
- Can be viewed as a fancy **normalizer** (a.k.a. Normalized exponential function)
- Produce a **discrete probability distribution** vector
- Very convenient when combined with cross-entropy loss

$$P(y = j \mid \mathbf{x}) = \frac{e^{\mathbf{x}^\top \mathbf{w}_j}}{\sum_{k=1}^K e^{\mathbf{x}^\top \mathbf{w}_k}}$$

- Sample vector \mathbf{x}
- Weight vector \mathbf{w}

Loss layer

- L1, L2 loss
- **Cross-Entropy** loss
(works well for classification, e.g., image classification)
- **Hinge** Loss
- **Huber** Loss, more resilient to outliers with smooth gradient
- **Minimum Squared Error** (works well for regression task, e.g., Behavioral Cloning)

$$J = - \sum_{i=1}^N y_i \log(h_\theta(x_i)) + (1 - y_i) \log(1 - h_\theta(x_i))$$

$$p(y_i|x_i) = [h_\theta(x_i)]^{(y_i)}[1 - h_\theta(x_i)]^{(1-y_i)}$$

$$p(y_i = 1|x_i) = h_\theta(x_i) \quad p(y_i = 0|x_i) = 1 - h_\theta(x_i)$$

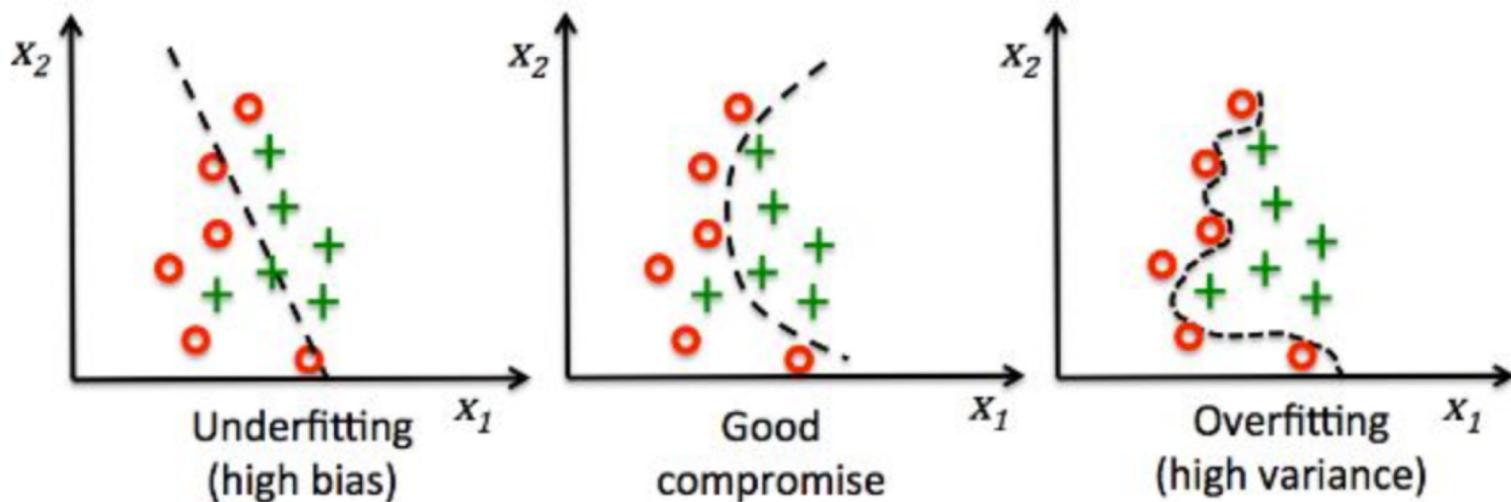
$$\sum_i \max(0, 1 - y_i * h_\theta(x_i))$$

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

$$J = \frac{1}{N} \sum_{i=1}^N (y_i - h_\theta(x_i))^2$$

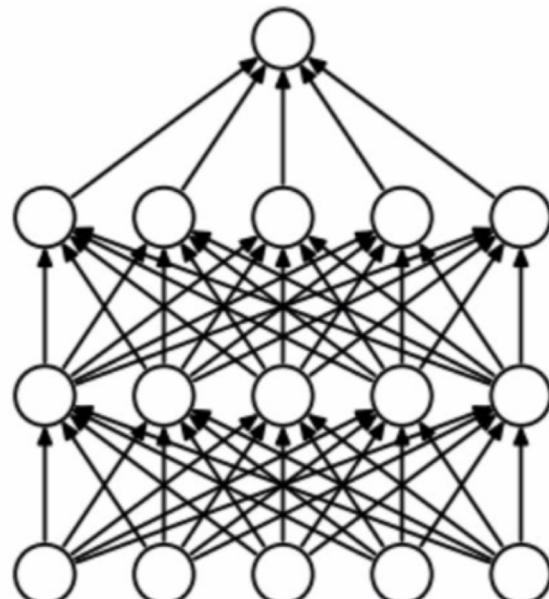
Regularization

- L1 / L2
- Dropout
- Batch norm
- Gradient clipping
- Max norm constraint

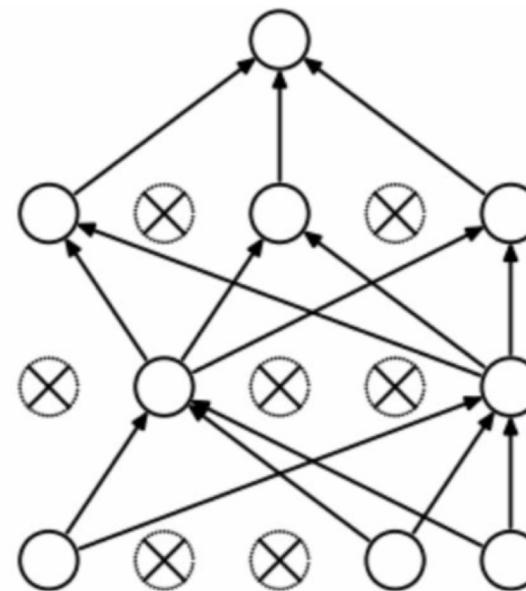


Dropout

- During training, randomly ignore activations by probability p
- During testing, use all activations but scale them by p
- Effectively prevent overfitting by reducing correlation between neurons



(a) Standard Neural Net



(b) After applying dropout.

Batch Normalization

- Makes networks robust to bad initialization of weights
- Usually inserted right before activation layers
- Reduce covariance shift by normalizing and scaling inputs
- The scale and shift parameters are trainable to avoid losing stability of the network

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Outline

- Background introduction
- Convolutional neural networks (CNN)
 - Overview
 - Network details: convolution, pooling, activation, loss functions
- ResNet
- DenseNet
- RNN

ResNet

- Residual Network, by Kaiming He (2015)
- Heavy usage of “**skip connections**” which are similar to RNN Gated Recurrent Units (GRU)
- Commonly used as visual feature extractor in all kinds of learning tasks, ResNet50, ResNet101, ResNet152
- 3.57% Top-5 accuracy, **beats human**

[Deep Residual Learning for Image Recognition](#)

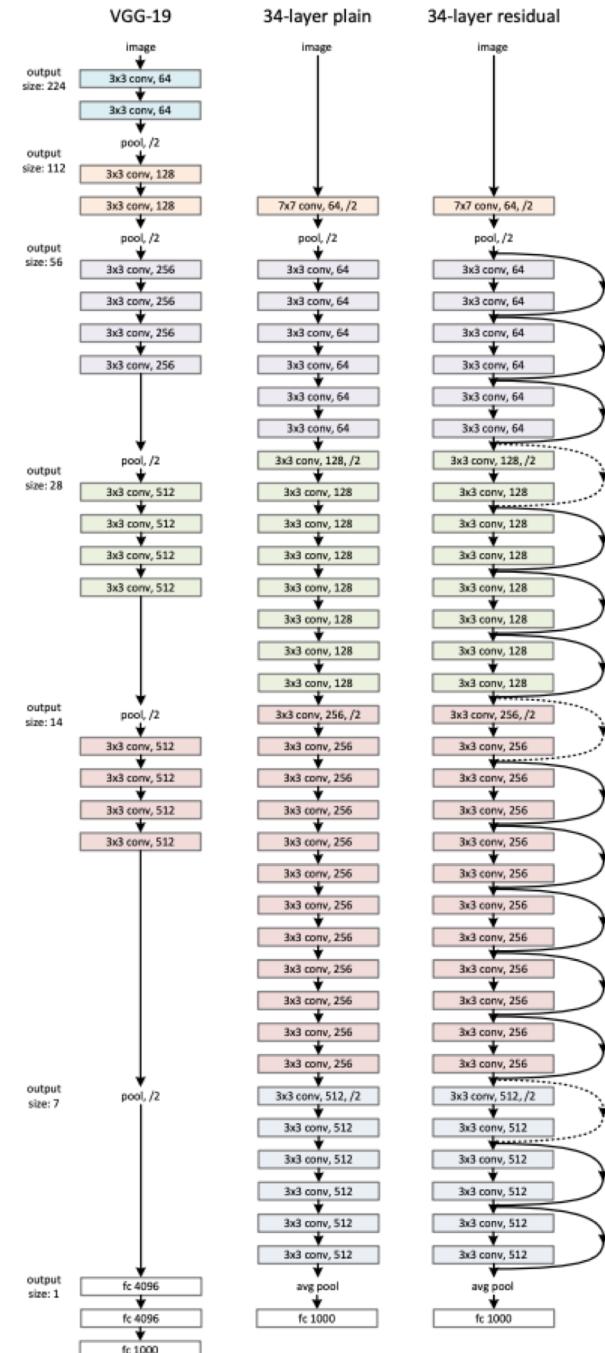
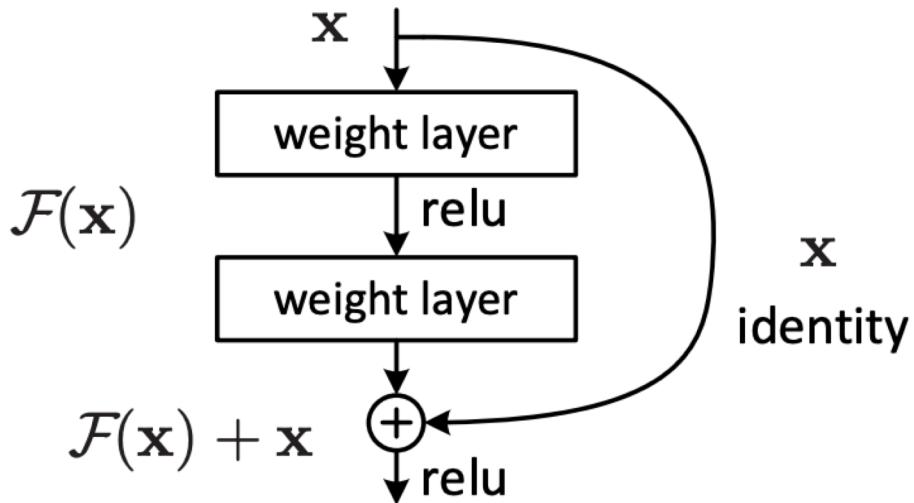
<https://arxiv.org> > cs ▾

by K He - 2015 - Cited by 21903 - Related articles

Dec 10, 2015 - Computer Science > Computer Vision and Pattern Recognition ... We explicitly reformulate the layers as learning residual functions with ...

ResNet architecture

Residual learning: a building block



DenseNet

Densely Connected Convolutional Networks

<https://arxiv.org/> > cs

by G Huang - 2016 - Cited by 3591 - Related articles

In this paper, we embrace this observation and introduce the Dense Convolutional Network (DenseNet), which connects each layer to every other layer in a ...

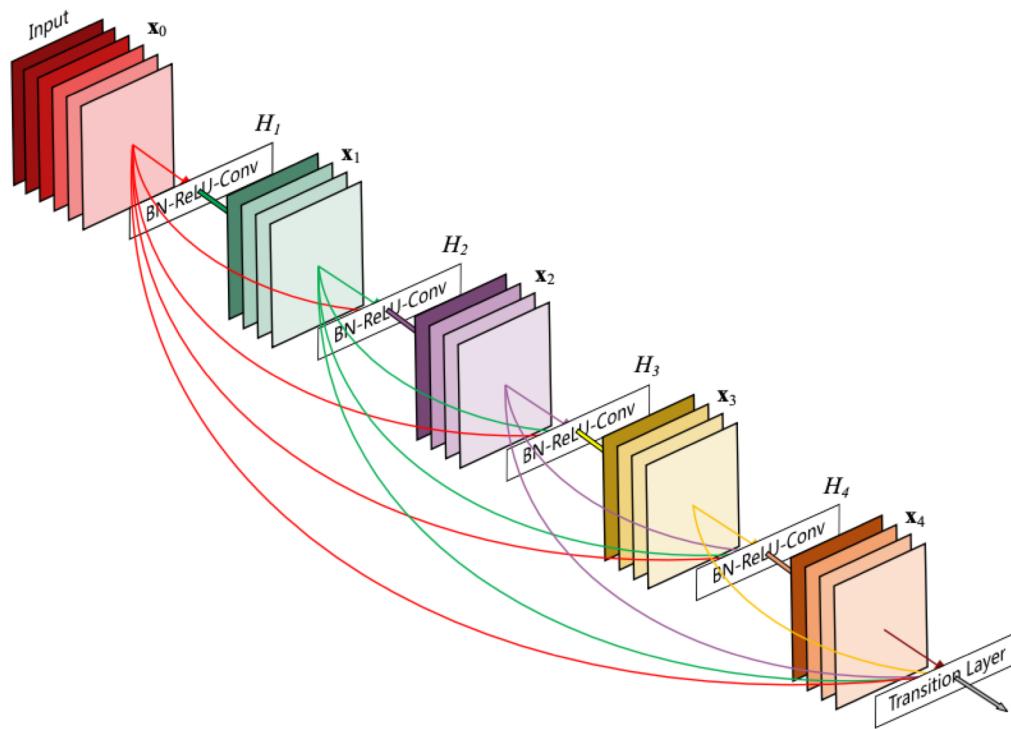
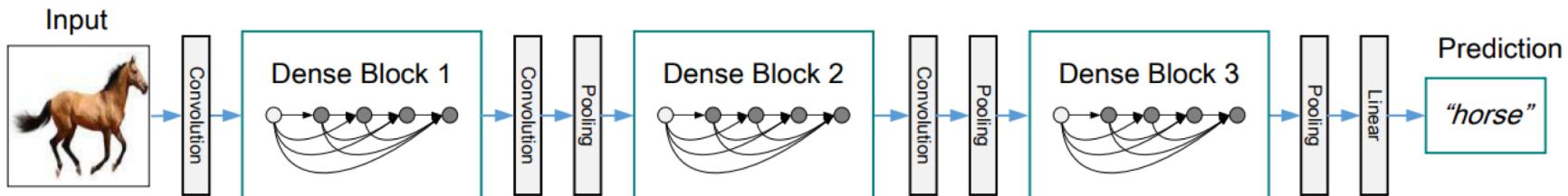
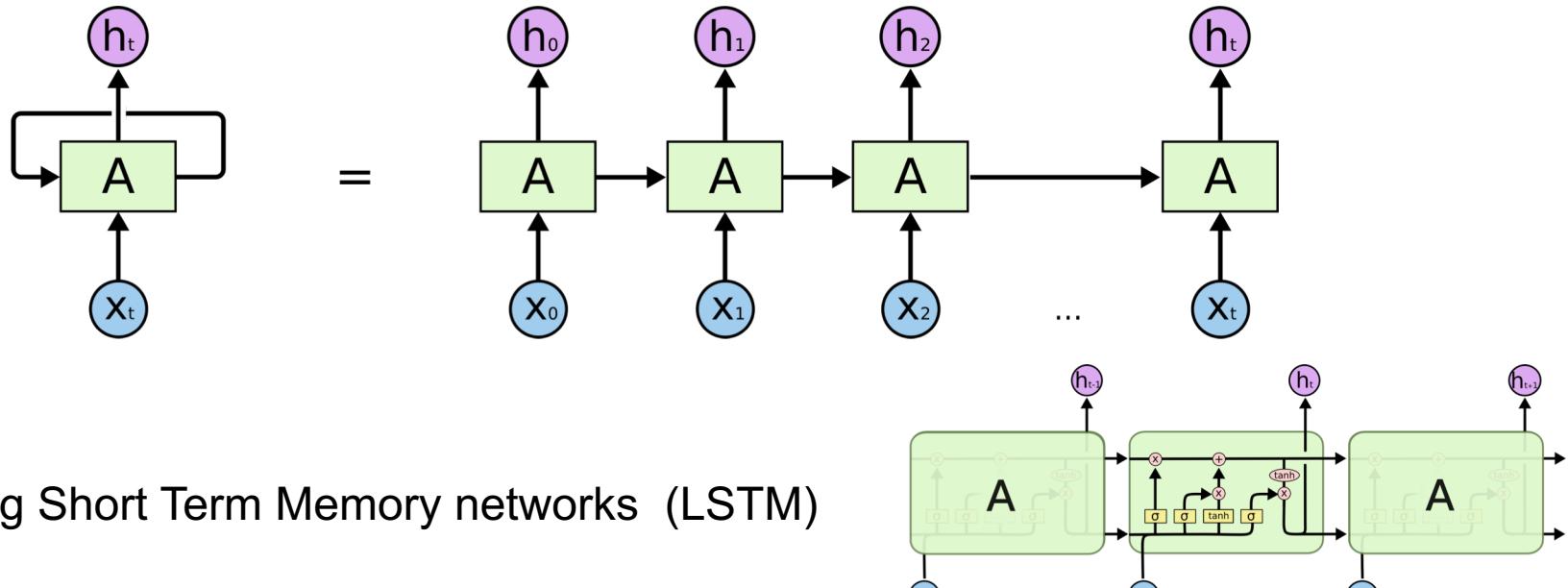


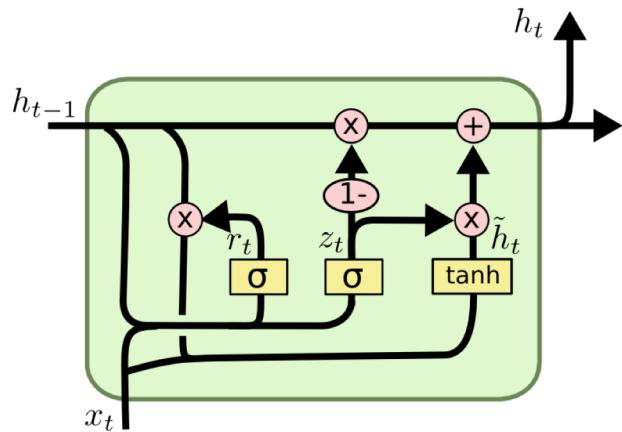
Figure 1: A 5-layer dense block with a growth rate of $k = 4$.
Each layer takes all preceding feature-maps as input.



Recurrent neural networks



Long Short Term Memory networks (LSTM)

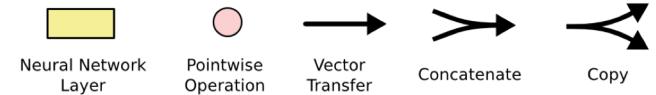


$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Thank you!