

# MLFA-Conv-Net

15-Mar-2023 at 9:52 PM

# Convolutional Neural Networks

## Motivation:

One strategy to better fit the data is to tune the network architecture to better represent the constraints and traits of the problem  
⇒ Inductive bias

## Example:

### Vision:

- ⇒ Nearby pixels are correlated
- ⇒ There are local features like edges and corners
- ⇒ An object (e.g. handwritten character) can be represented as a combination of such features.

### Speech:

- ⇒ Locality in terms of time
- ⇒ inputs that are temporally close can be grouped
- ⇒ Combining these primitives, utterances like speech phoneme can be defined.

\*\* Hidden units can be designed not to consider connections from all inputs.

⇒ Hidden unit that define a window over the input space

⇒ Connected to a local subset of inputs

⇒ Decrease in the number of parameters.

This idea of local aggregation can be applied in the successive layers

⇒ Each layer is connected to a small number of local units below

⇒ Combining features from the previous layer to construct more complex features.

Example?

pixels → edges → arcs, corners → circle, rectangle  
↓  
digit

Concern with number of parameters:

⇒ units checking for different parameters in all local patches.

⇒ Too many parameters

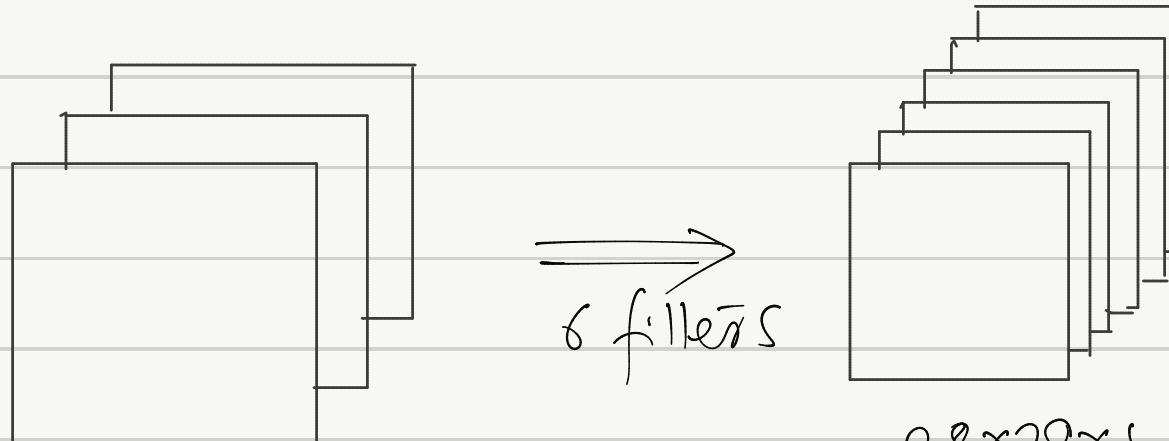
⇒ Reduction in no of parameters: Weight sharing

When looking for features like oriented edges, they may be present in different parts of the input space.

## Design

⇒ image data is represented as 3D matrices ⇒ TENSOR  
 width × height × no of channels (RGB)

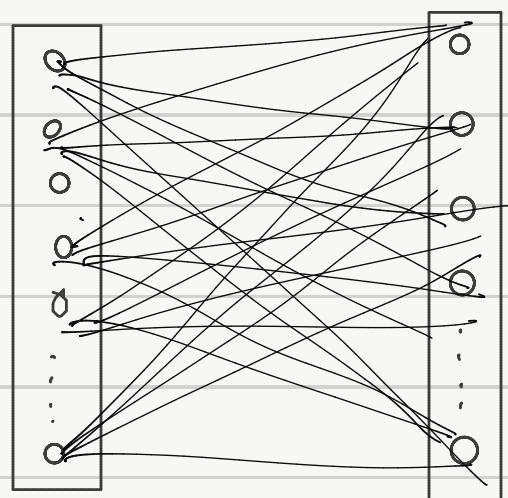
## Filters



$32 \times 32 \times 3$

## Parameter sharing

If we had used a fully connected NN



$32 \times 32 \times 3$   
 $= 3072$

$28 \times 28 \times 6$   
 $= 4704$

$3072 \times 4704 > 14 \text{ million}$   
 parameter  
 only for single  
 layer

Observations: 14 million parameters are not independent from each other

⇒ Each of the 6 filters represents a fixed transformation

⇒ The hidden units for a given filter performs

same kind of computation in every part  
of the image

⇒ The weights can be kept same  
across the input space

### Sparse Connection

Let us consider a vertical edge detector

A vertical edge may reside in a tiny part of  
the image and hence detection will depend  
on a small fraction of the input pixels

⇒ ignoring contributions from the other  
pixels

⇒ Sparse Connection

### Convolution Operation:

Image matrix :  $A \in \mathbb{R}^{w \times h}$

$f < w, f < h$

Filter or Kernel :  $F \in \mathbb{R}^{f \times f}$

→ learned by  $i^{\text{th}}$

$A * F \in \mathbb{R}^{w-f+1, h-f+1}$

neural net

$$(A * F)_{ij} = \sum_{u=1}^f \sum_{v=1}^f A_{i+u-1, j+v-1} \cdot f_{uv}$$

$$1 \leq i \leq w-f+1$$

$$1 \leq j \leq h-f+1$$

1	0	-1		4	5	6
1	2	-1		4	5	6
1	2	0	-1	1	3	5
1	6	0	-1	3	2	1
1	6	2		5	3	4
4	1	5	2	6	3	
2	1	6	5	4	3	

6x6

$$\begin{aligned}
 & 1 \times 1 + 0 \times 2 + (-1) \times 3 + 1 \times 2 + \\
 & 0 \times 4 + (-1) \times 6 + 1 \times 6 + 0 \times 5 + \\
 & -1 \times 4 = -4
 \end{aligned}$$

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} *
 \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = 
 \begin{array}{|c|c|c|} \hline -4 & 3 & 3 \\ \hline -3 & 6 & 4 \\ \hline 0 & 2 & 0 \\ \hline -6 & -4 & 0 \\ \hline 2 & & \\ \hline \end{array}$$

$$\begin{array}{l} 6-3+1 \times 6-3+1 \\ 4 \times 2 \end{array}$$

Filter vertical Edge Detector:

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} *
 \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = 
 \begin{array}{|c|c|c|} \hline 0 & 30 & 30 \\ \hline \end{array}$$

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

vertical edge indicator



1	1	1
1	1	1
1	1	1

Kernel that

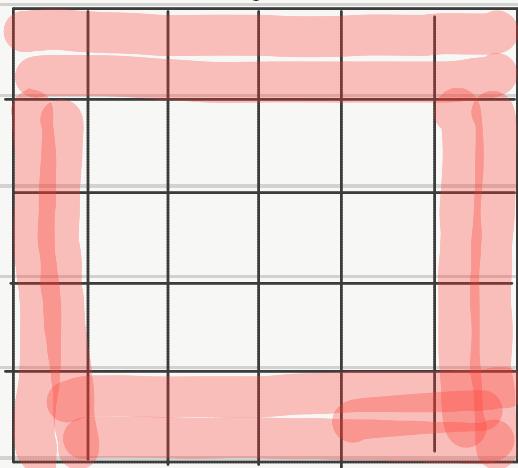
sharpens an  
image

0	-1	0
-1	5	-1
0	-1	0

Kernel that

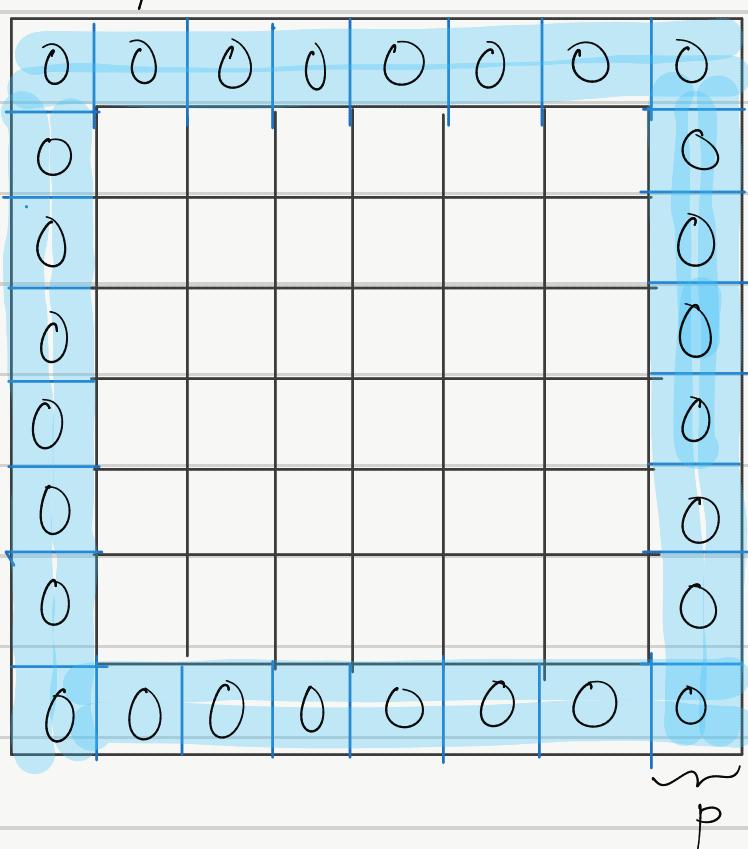
sharpens an  
image

## Padding:



- \* The convolution operator misses in first and last column or row  
 $\Rightarrow$  if there are important info in these rows and columns they are missed out.

Padding ( $p$ )



- \* Size of the output is smaller than the input matrix

$\Rightarrow$  shrinking

$\Rightarrow$  after few conv layers in image shrinks significantly

Input :  $w \times h$  || Padding :  $p$

Input after padding :  $w+2p \times h+2p$

Filles :  $f \times f$

Output :  $(w+2p-f+1) \times (h+2p-f+1)$

## Choosing padding size:

$\Rightarrow$  "valid" convolution :  $p=0$

$\Rightarrow$  "same" convolution

$\Rightarrow$   $p$  is chosen such that input and output size is same.

Input size :  $W \times h$

$$\text{Output size} = (W + 2p - f + 1) \times (h + 2p - f + 1)$$

$$W + 2p - f + 1 = W$$

$$p = \frac{f-1}{2}$$

Typically  $f$  is chosen to be odd so that  $\frac{f-1}{2}$  becomes an integer division.

Strides in convolution :

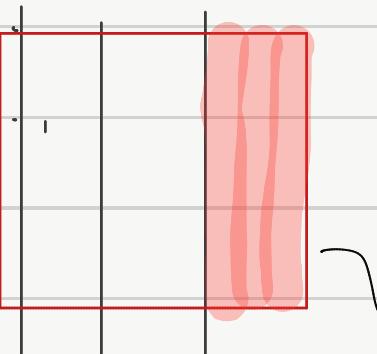
$\Rightarrow$  Shifting filter by  $s$  columns or rows

Input :  $W \times h$ , Padding :  $p$ . Filter :  $f \times f$   
Stride :  $s$

$$\text{Output size} : \left(\left\lfloor \frac{W+2p-f}{s} \right\rfloor + 1\right) \times \left(\left\lfloor \frac{h+2p-f}{s} \right\rfloor + 1\right)$$

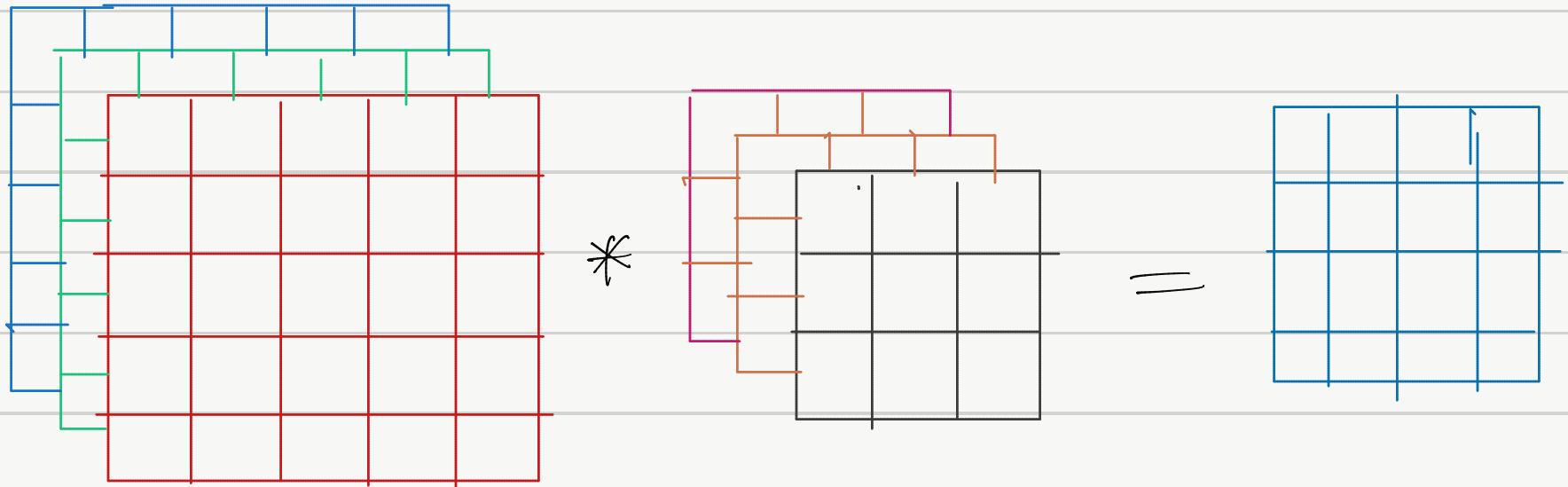
Floor function:

$\Rightarrow$  if in filter matrix is hanging at an edge, discard the computation of the hanging part



ignored.

## Convolution on Tensors:

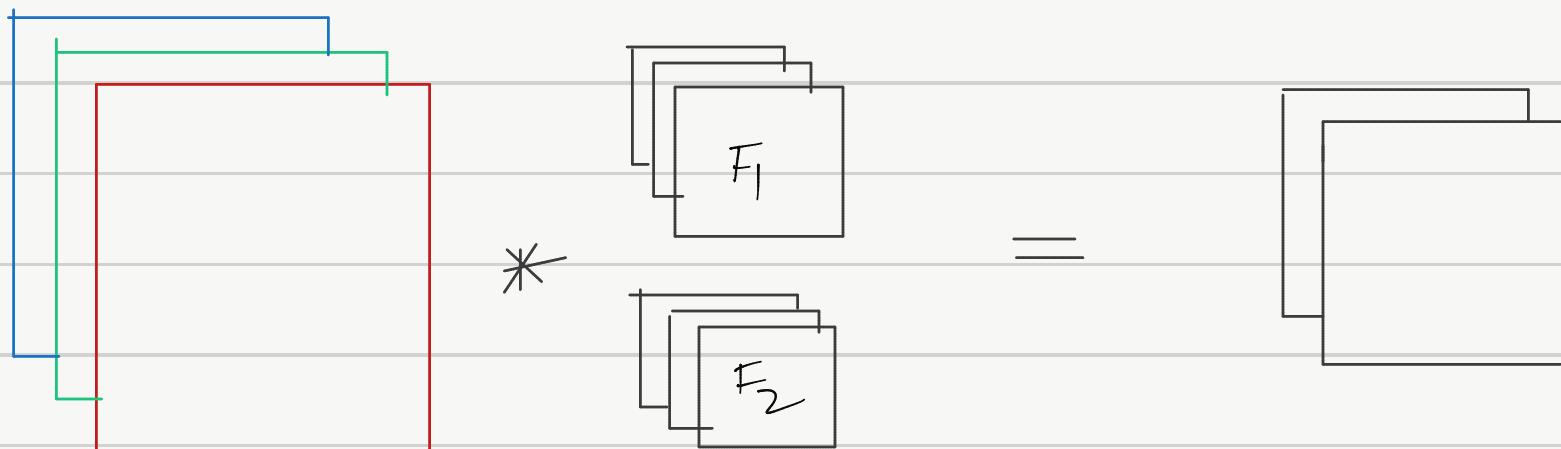


Input:  $A \in \mathbb{R}^{w \times h \times c}$ , Filter:  $F \in \mathbb{R}^{f \times f \times c}$

$$A * F \in \mathbb{R}^{(w-f+1) \times (h-f+1)}$$

$$(A * F)_{ij} = \sum_{u=1}^f \sum_{v=1}^f \sum_{k=1}^c A_{i+u-1, j+v-1, k} \cdot F_{u, v, k}$$

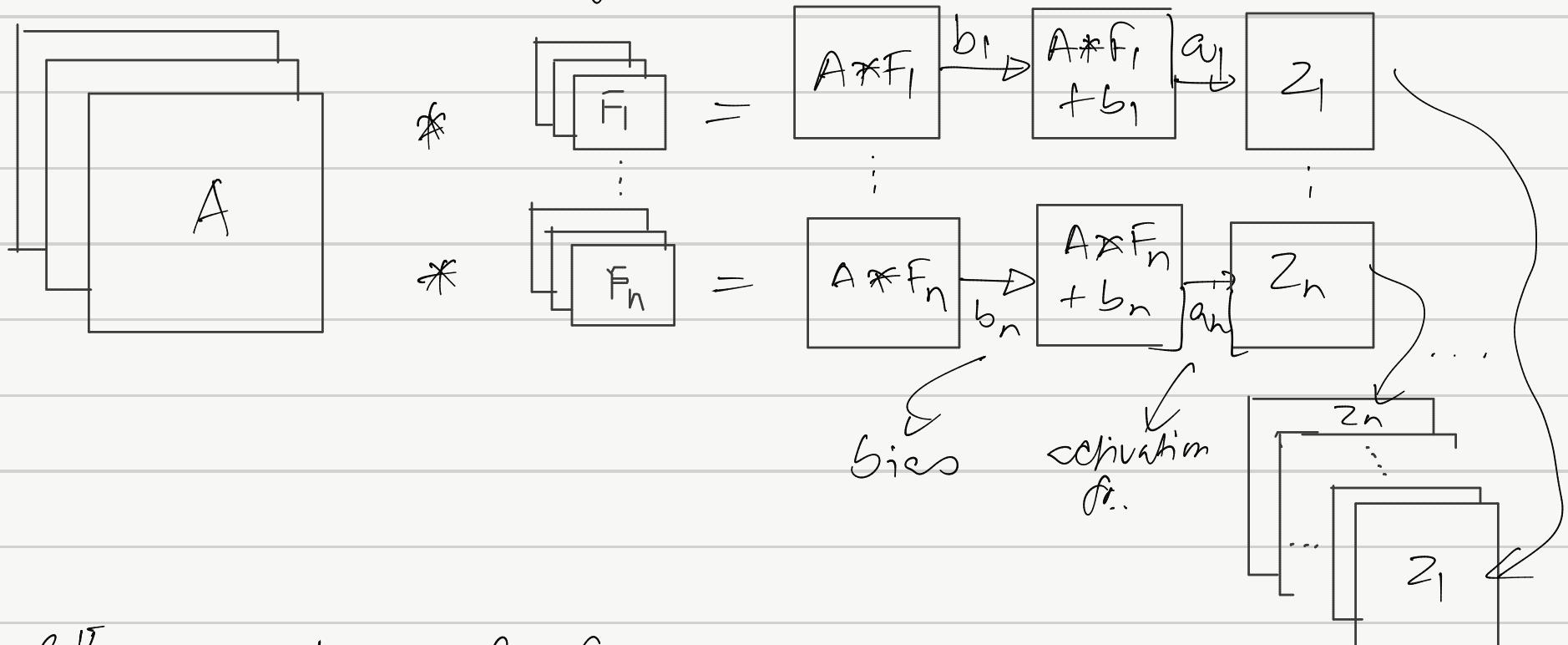
## Multiple filters:



Input:  $A \in \mathbb{R}^{w \times h \times c}$        $F_i \in \mathbb{R}^{f \times f \times c}$        $F = [F_1, \dots, F_{n_f}]$

$$A * F \in \mathbb{R}^{(w-f+1) \times (h-f+1) \times n_f}$$

# Convolutional Layers:



## $\ell^{th}$ convolutional Layer

- Filter Size :  $f^{(\ell)}$
- No of filters :  $n_f^{(\ell)}$
- padding size :  $p^{(\ell)}$
- stride length :  $s^{(\ell)}$

Input :  $w^{(\ell-1)} \times h^{(\ell-1)} \times c^{(\ell-1)}$

Output :  $w^{(\ell)} \times h^{(\ell)} \times c^{(\ell)}$

$$w^{(\ell)} = \left\lceil \frac{w^{(\ell-1)} + 2p^{(\ell)} - f^{(\ell)}}{s^{(\ell)}} \right\rceil + 1$$

$h^{(\ell)} \rightarrow$  similar

$$c^{(\ell)} = n_f^{(\ell)}$$

## Parameters:

- For each filter :  $f^{(\ell)}, f^{(\ell)}, c^{(\ell)}$
- For each filter : 1 bias
- Total :  $n_f^{(\ell)} [f^{(\ell)}, f^{(\ell)}, c^{(\ell-1)} + 1]$

## Pooling Layers:

- Similar to applying filter in convolution layer
- MAX, AVG pooling
- used to shrink the input size

Filter size:  $f$

Stride length:  $s$

Type of pooling: MAX or AVG

Padding:  $p$  [typically 0]

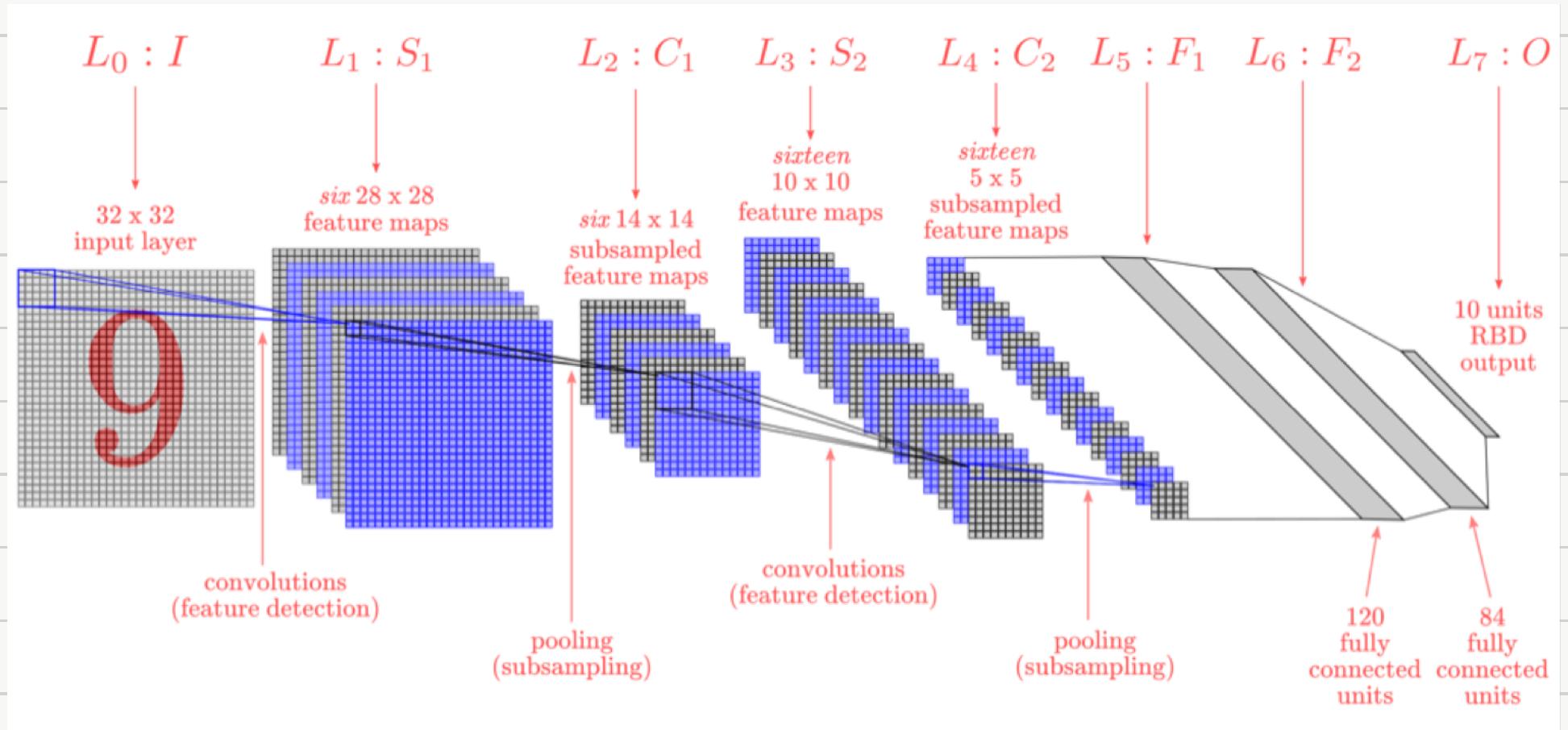
input:  $N \times h \times c$

output:  $(\lfloor \frac{N-f}{s} \rfloor + 1) \times (\lfloor \frac{h-f}{s} \rfloor + 1) \times c$

Pooling layer has fixed computation  
⇒ no parameters to learn

## Architecture of CNN:

LeNet - 5 (Lecun et al 1998)



Parameters ~ 60K

## Training a CNN:

→ Similar to training a fully connected network

