

Prof Xiaowei Huang

<https://cgi.csc.liv.ac.uk/~xiaowei/>

(Attendance Code: **908791**)

Lecture 19 -- Convolutional Neural Networks



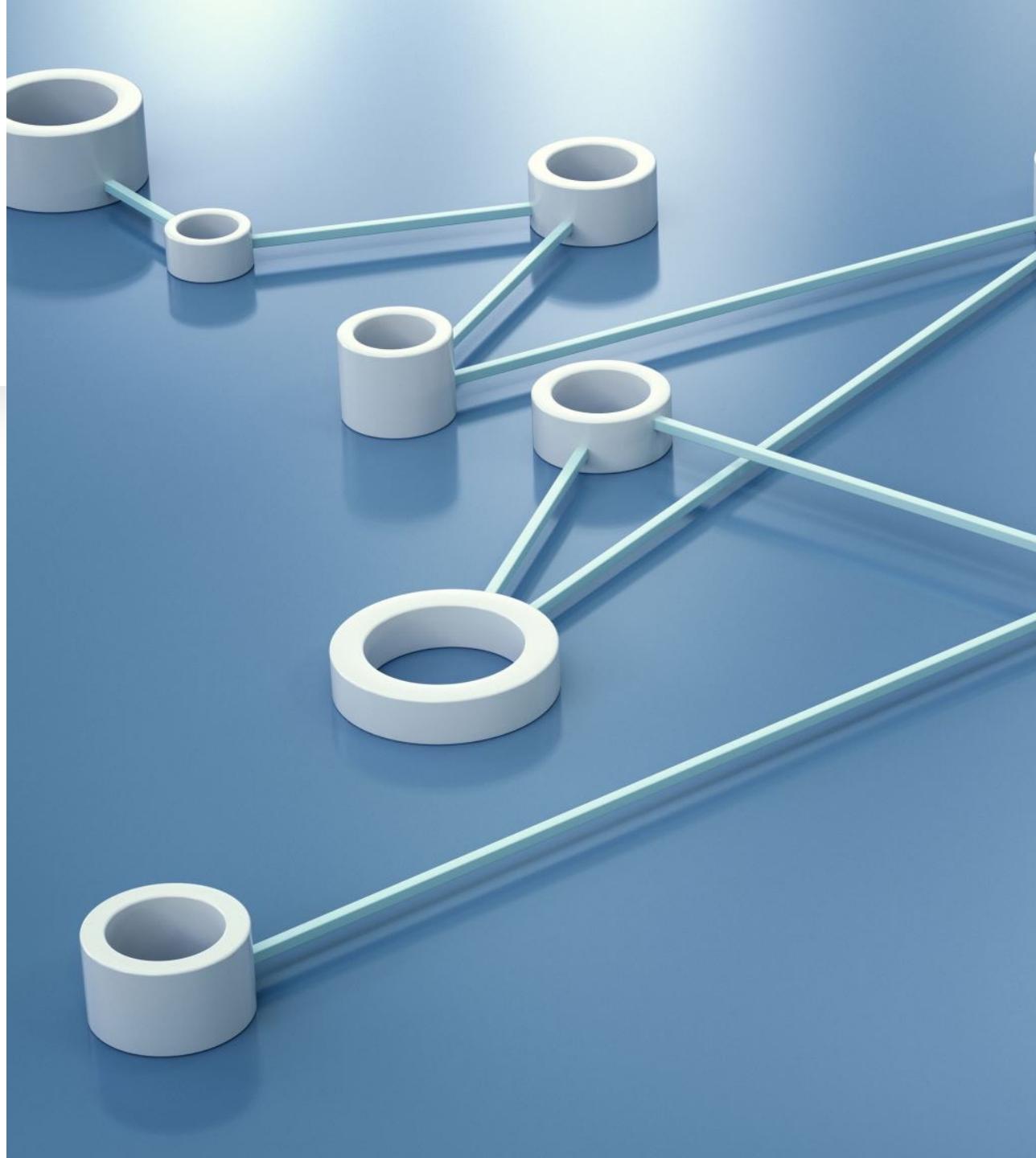


Up to now,

- Traditional Machine Learning Algorithms
- Adversarial Attack and Defence
- Deep learning
 - Introduction to Deep Learning
 - Functional view and features
 - Backward and forward computation (including backpropogation and chain rule)

Topics

- Convolutional neural networks (CNN)
 - Fully-connected
 - Convolutional Layer
 - Advantage of Convolutional Layer
 - Zero-padding Layer
 - ReLU Layer
 - Pooling Layer
 - Residual Block
 - Dropout
 - BatchNorm Layer
 - Softmax Layer
- Preprocessing data
- Example
 - LeNet



```
model = Sequential()

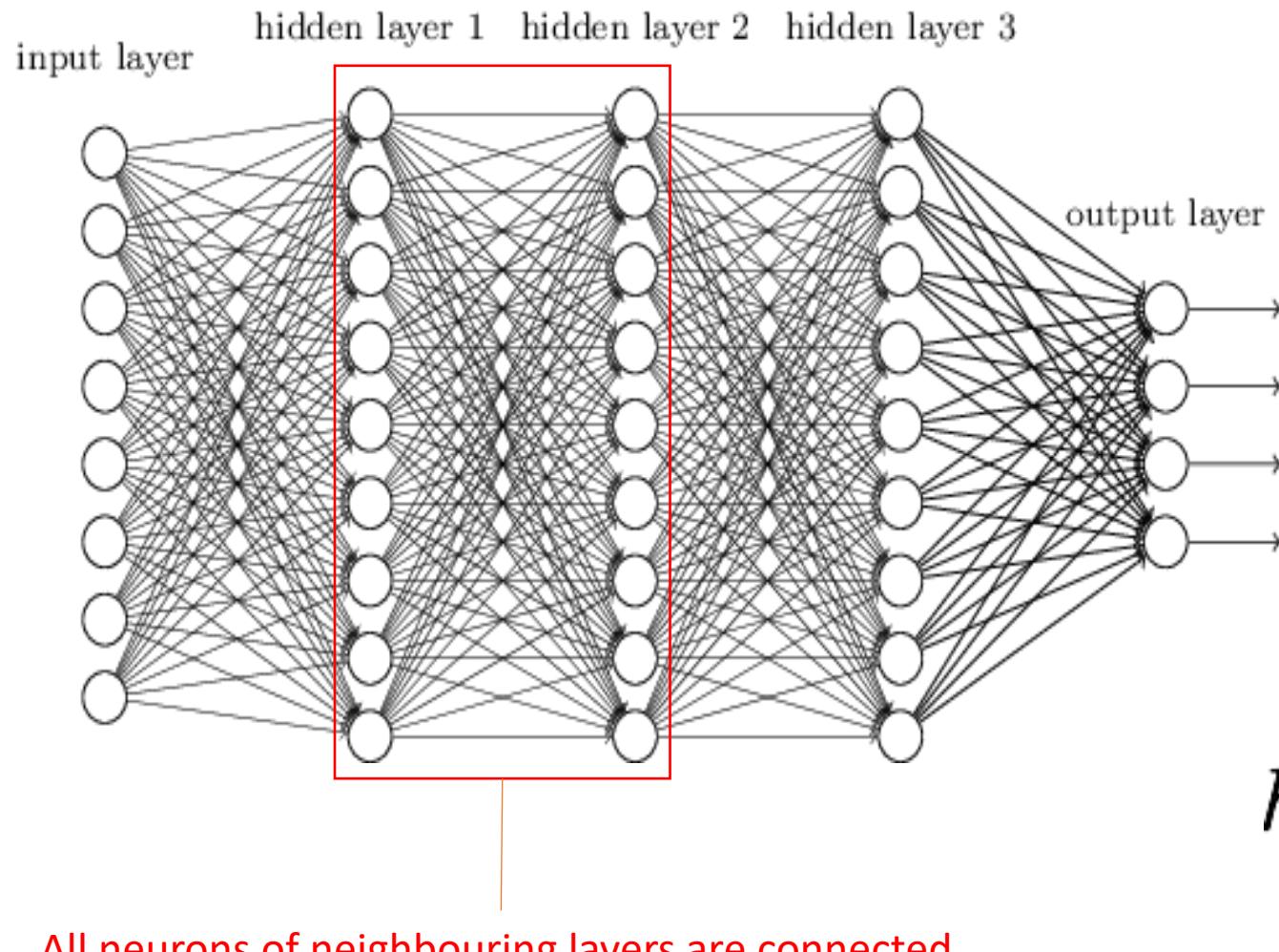
model.add(Convolution2D(nb_filters, nb_conv, nb_conv,
                      border_mode='valid',
                      input_shape=(1, img_rows, img_cols)))          Convolutional layer
model.add(Activation('relu'))                         ReLU layer
model.add(Convolution2D(nb_filters, nb_conv, nb_conv))          Convolutional layer
model.add(Activation('relu'))                         ReLU layer
model.add(MaxPooling2D(pool_size=(nb_pool, nb_pool)))      Maxpooling layer
model.add(Dropout(0.25))                            Dropout layer: for
                                                       regularisation

model.add(Flatten())                                Flatten layer: from
                                                       convolutional to fully-
                                                       connected
model.add(Dense(128))
model.add(Activation('relu'))
model.add(Dropout(0.5))
model.add(Dense(nb_classes))
model.add(Activation('softmax'))                     Fully-connected layer

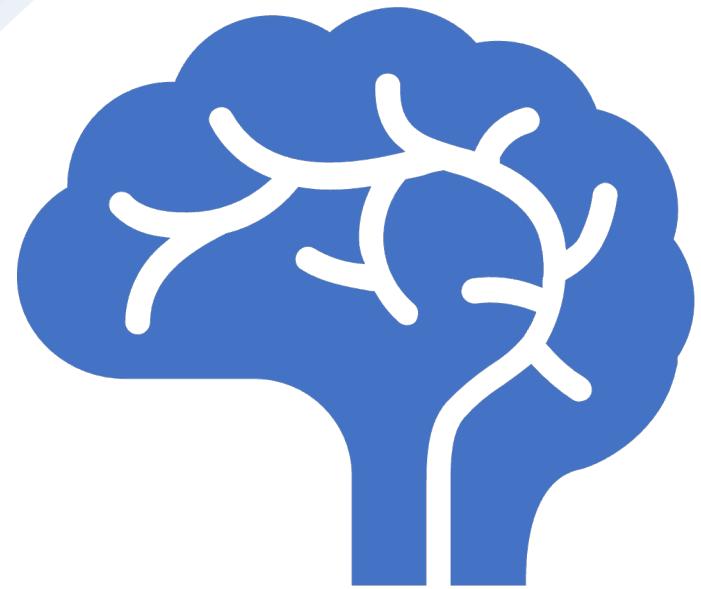
model.compile(loss='categorical_crossentropy',
              optimizer='adadelta',
              metrics=['accuracy'])
```



Fully-connected Layer



$$h = \sigma(W^T x + b)$$



Convolutional Layer

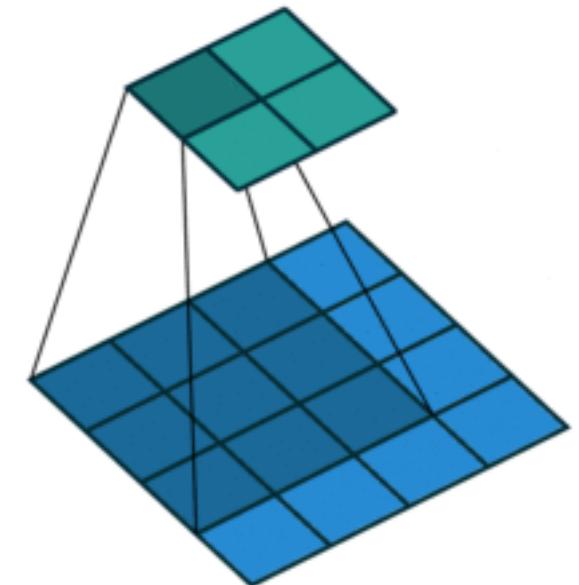
Convolutional neural networks

- Strong empirical application performance
- Convolutional networks: neural networks that **use convolution in place of general matrix multiplication** in at least one of their layers

$$h = \sigma(W^T x + b)$$

for a **specific kind of weight matrix W**

Will be explained later



Convolutional layer illustrated

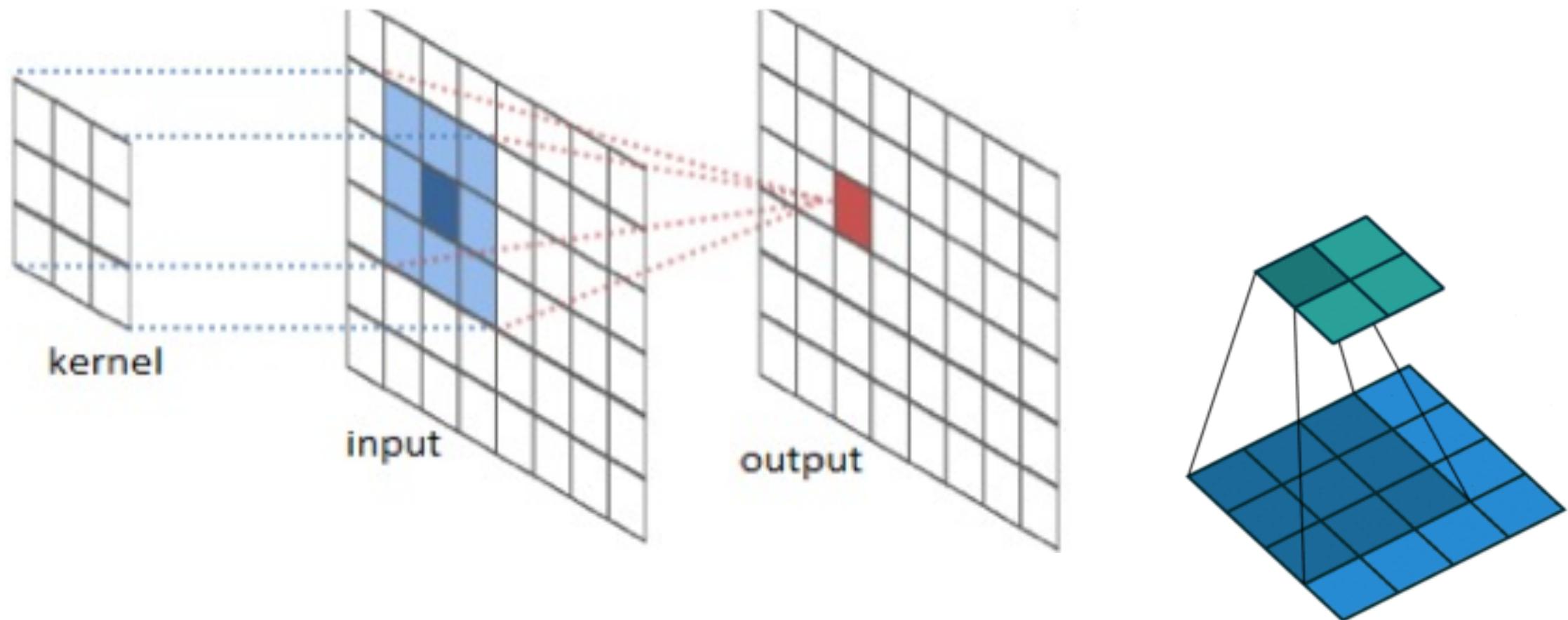


Illustration 1

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$



Illustration 1

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$

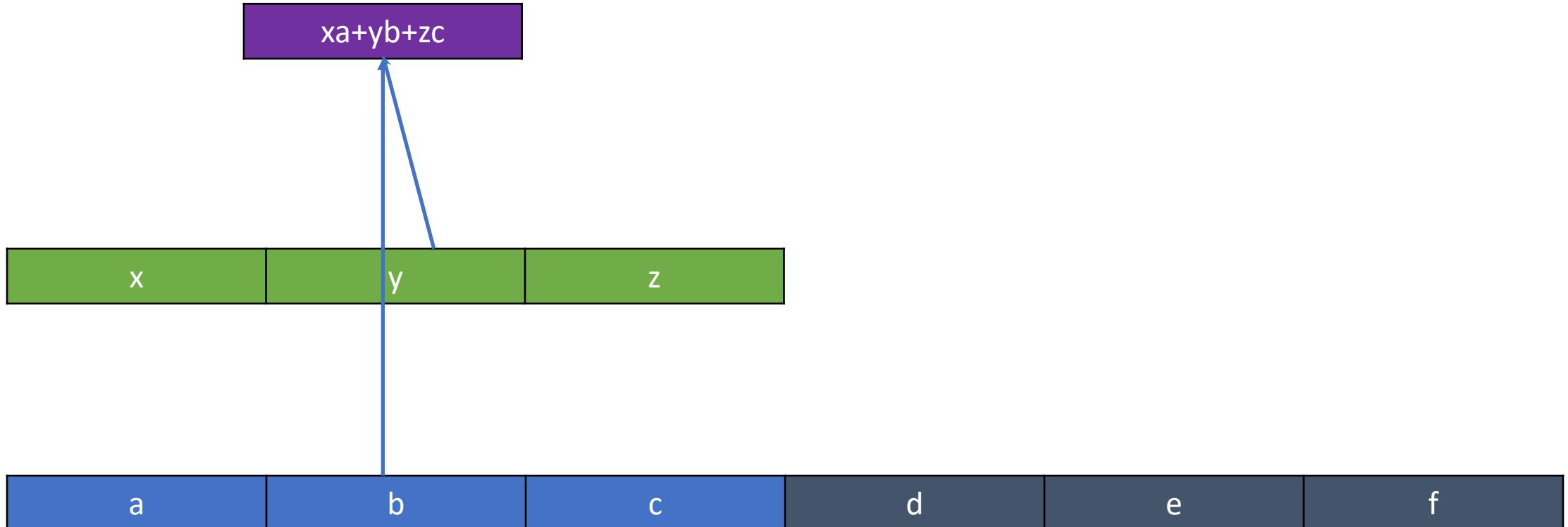


Illustration 1

$$w = [z, y, x]$$
$$u = [a, b, c, d, e, f]$$

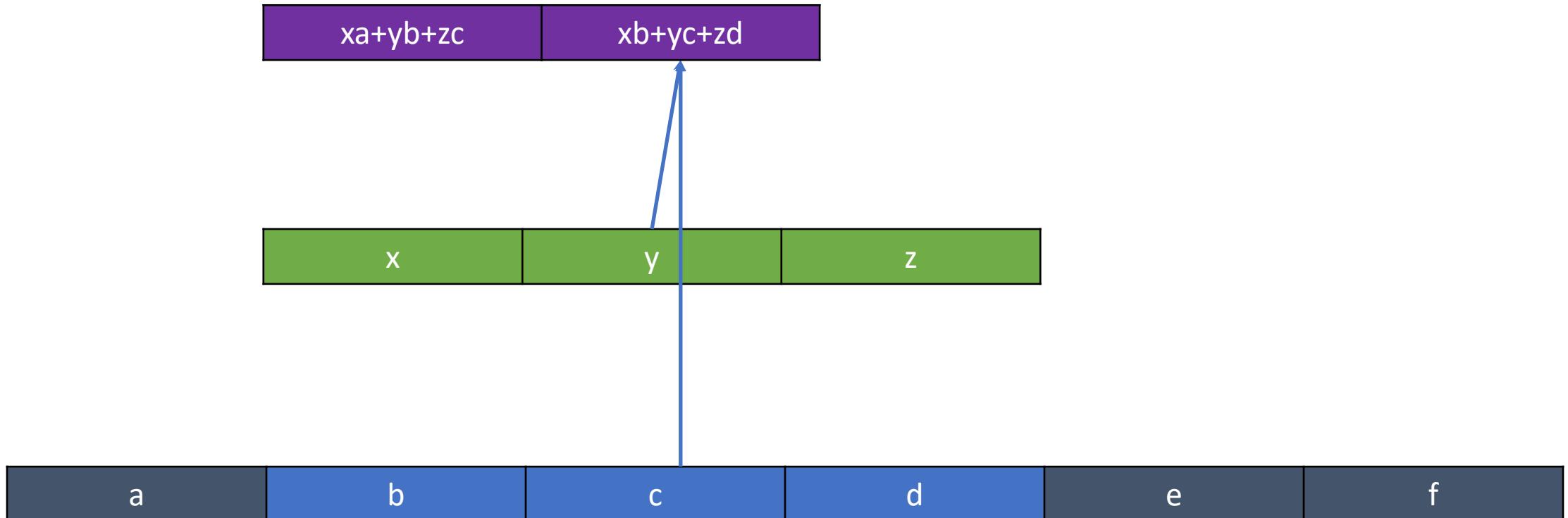


Illustration 1

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$

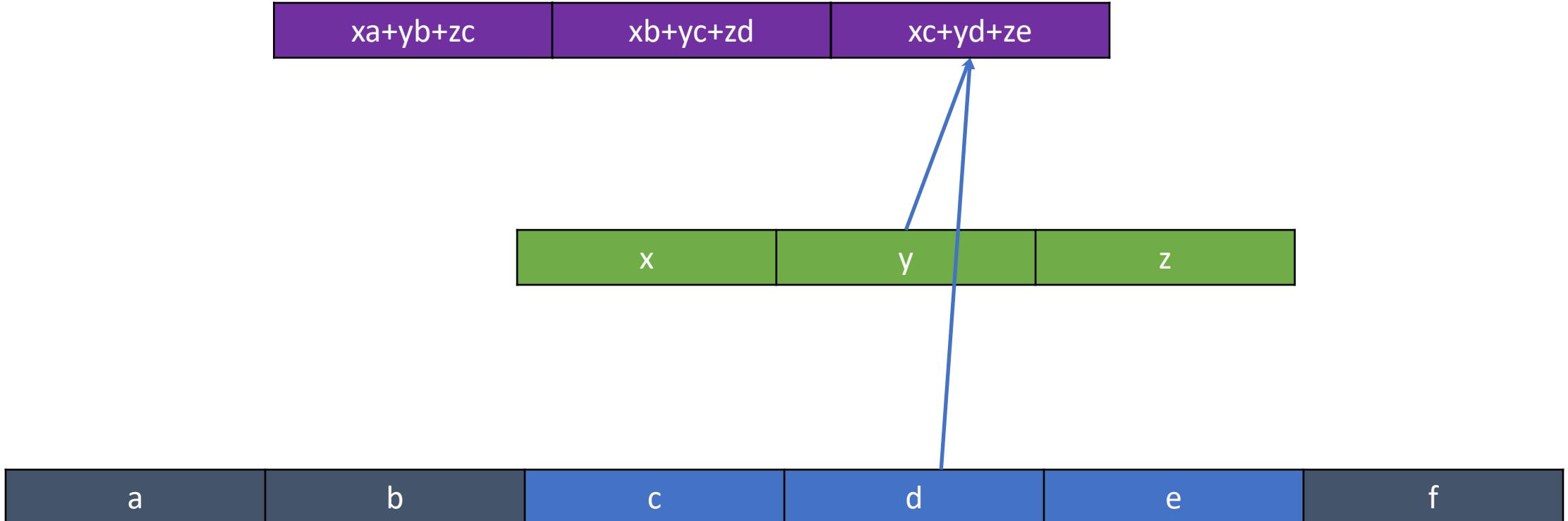


Illustration 1

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$

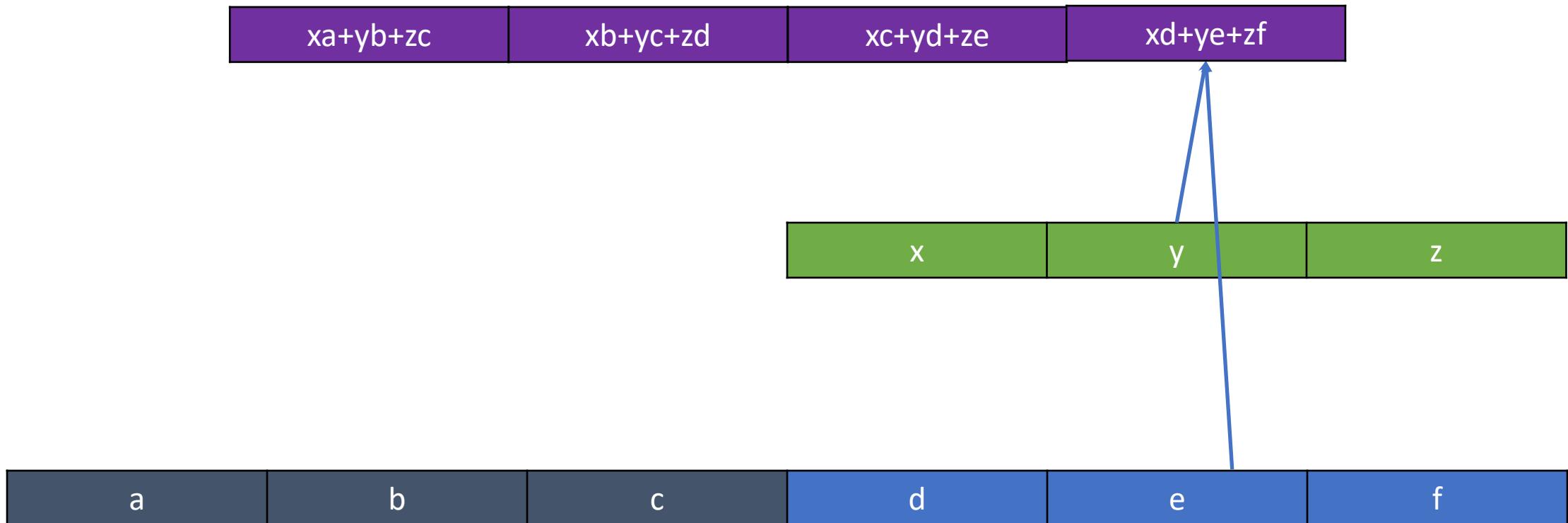


Illustration 1: boundary case

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$

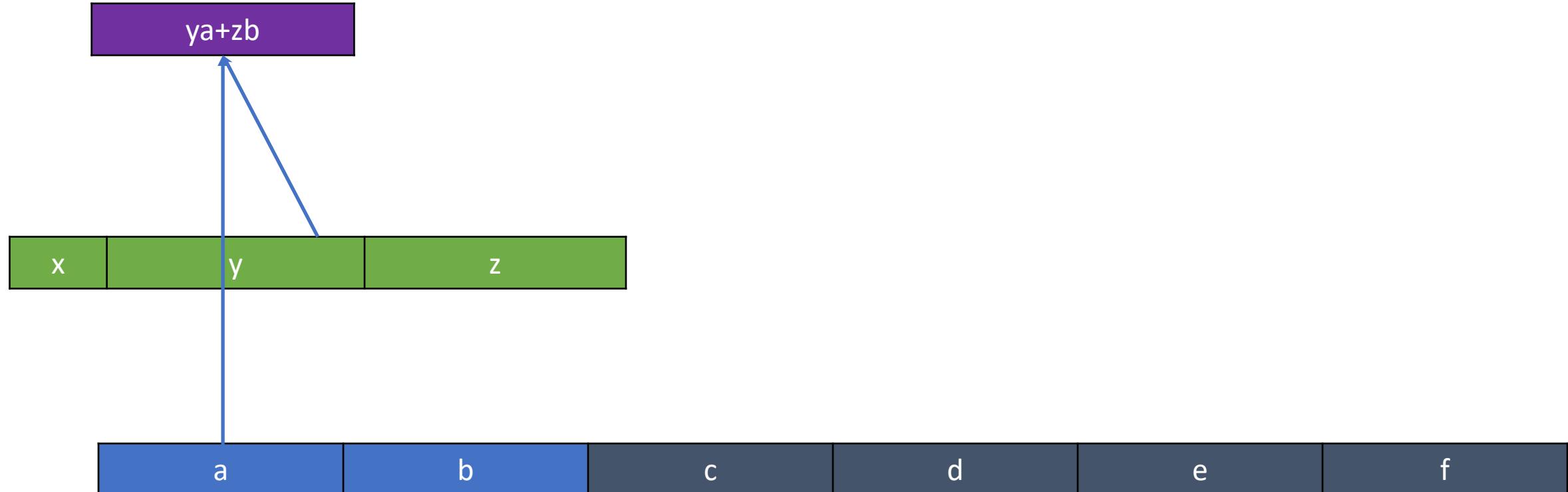


Illustration 1: boundary case

$$w = [z, y, x]$$

$$u = [a, b, c, d, e, f]$$

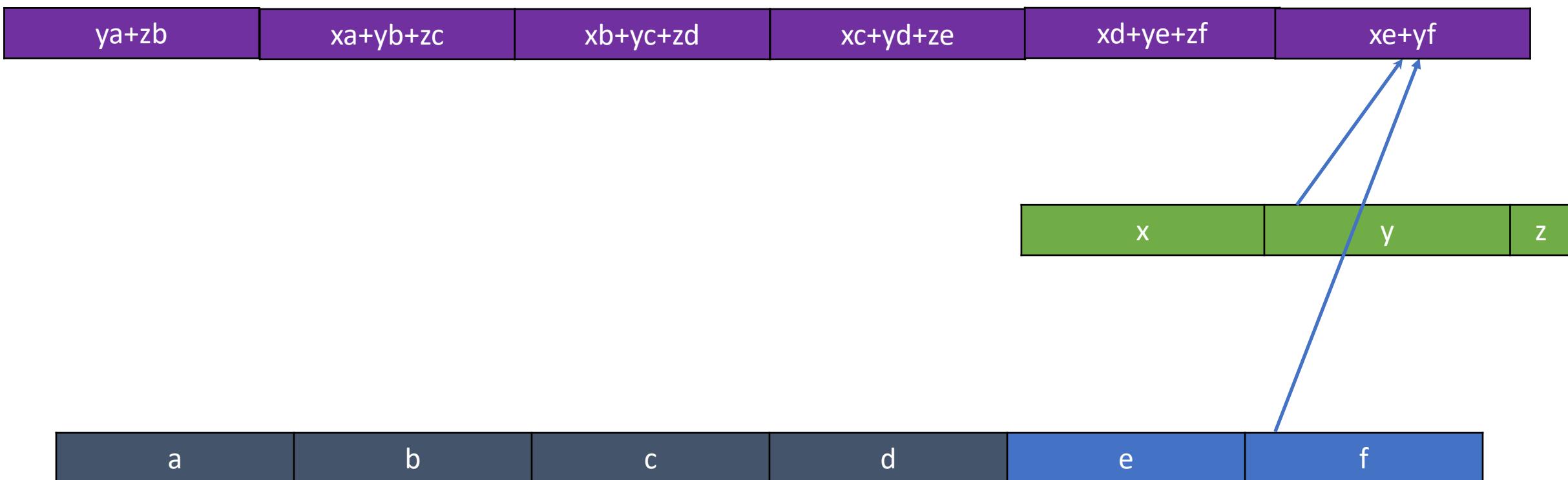
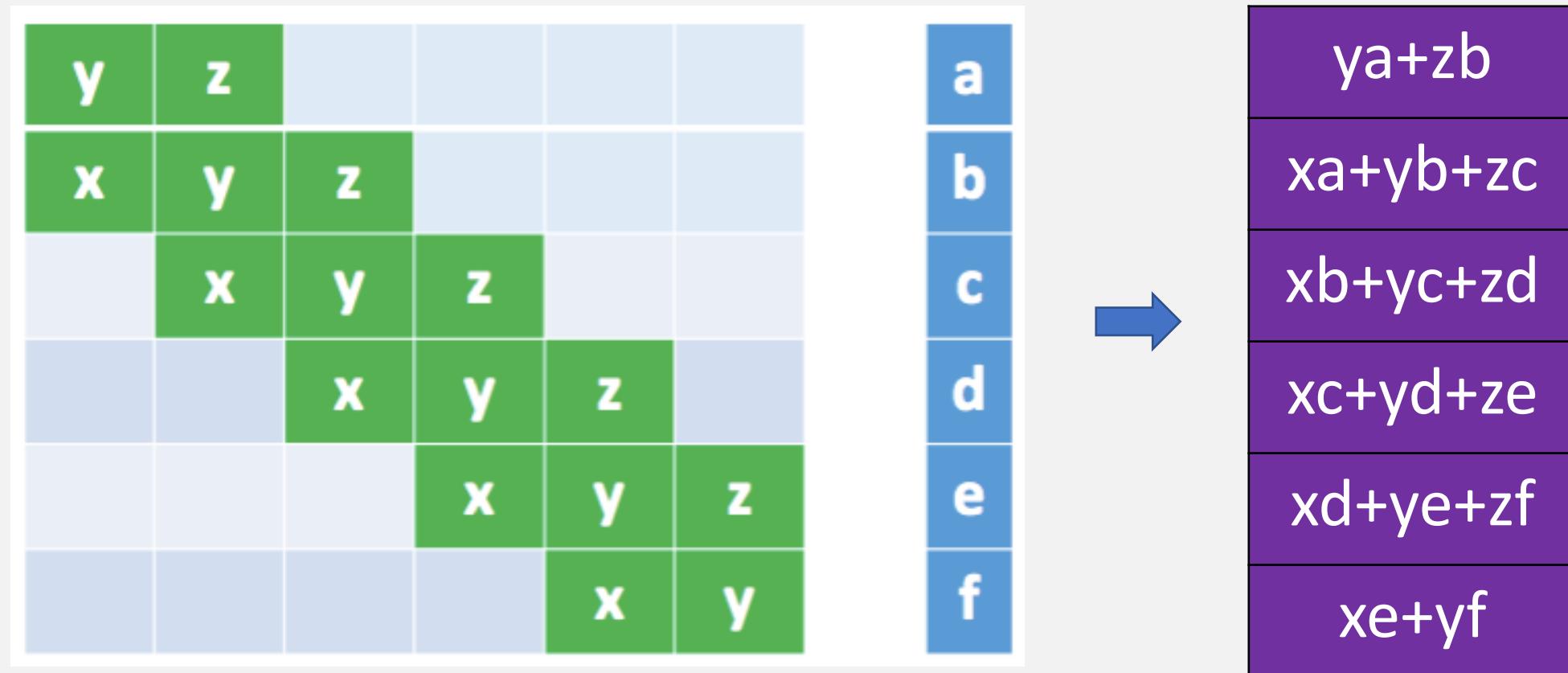


Illustration 1: as matrix multiplication



It also explains why convolutional layer is a special type of fully connected layer
– weights of non-green cells are known to be 0 before training.

Illustration 2: two dimensional case

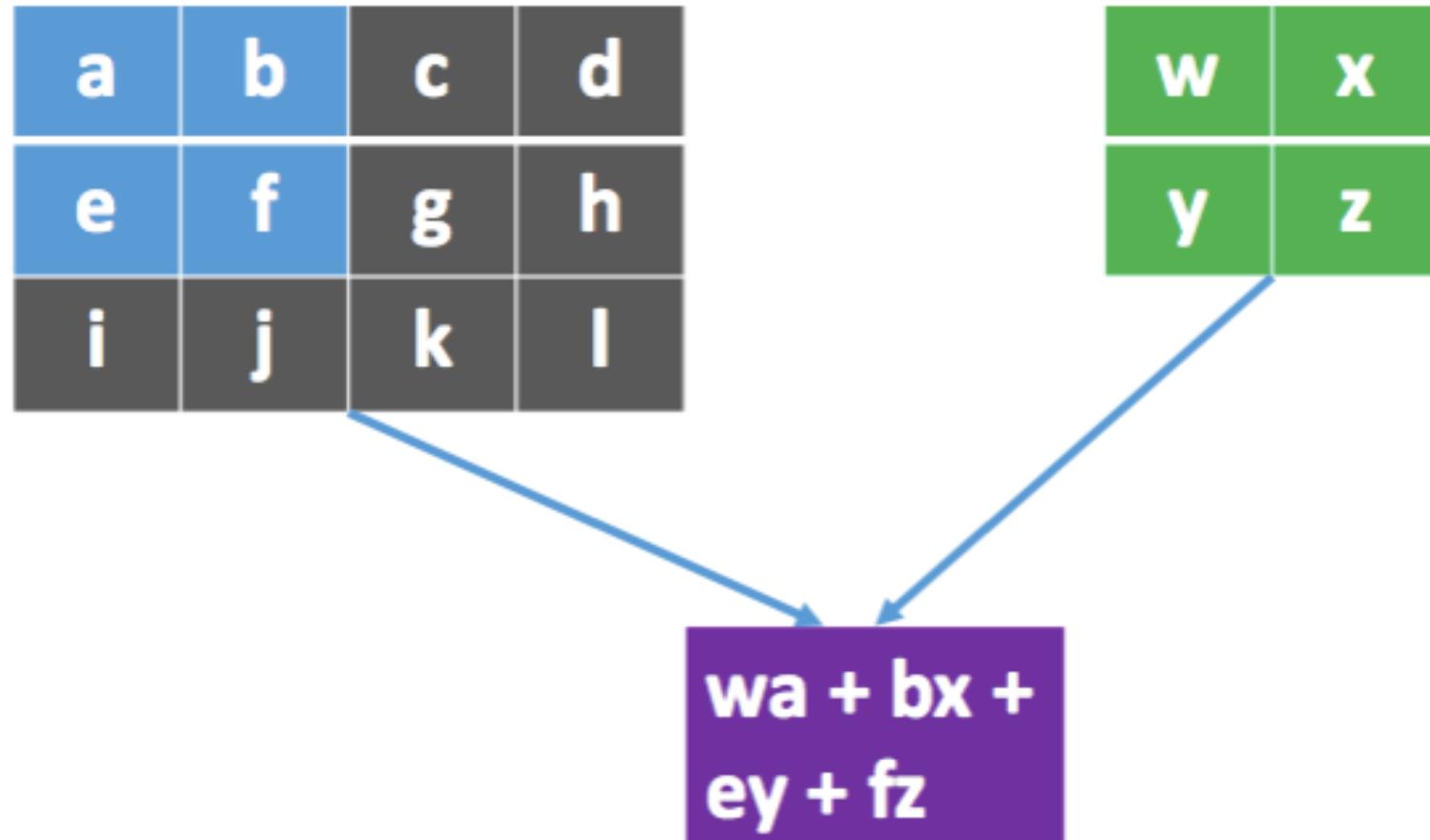


Illustration 2

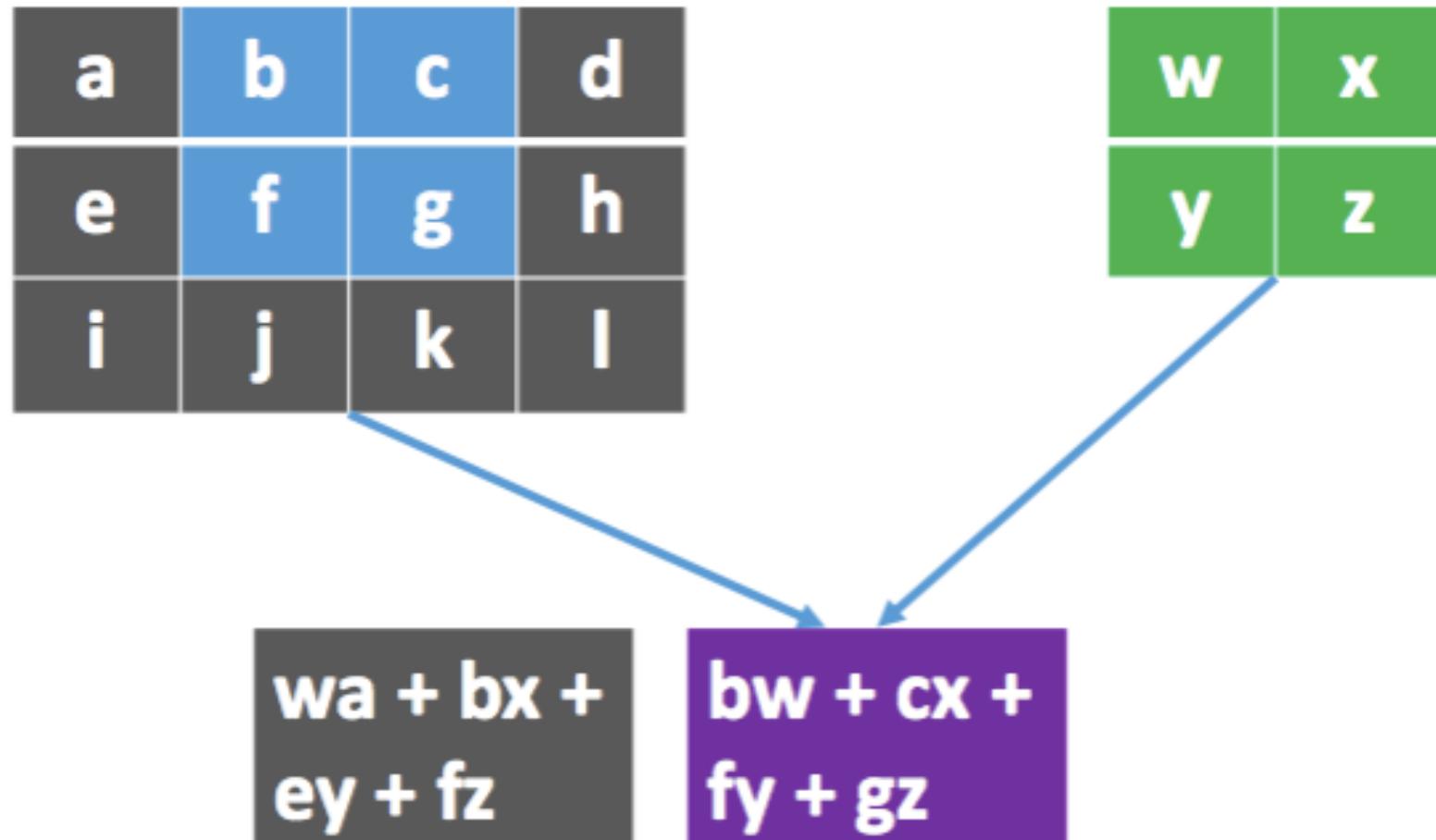
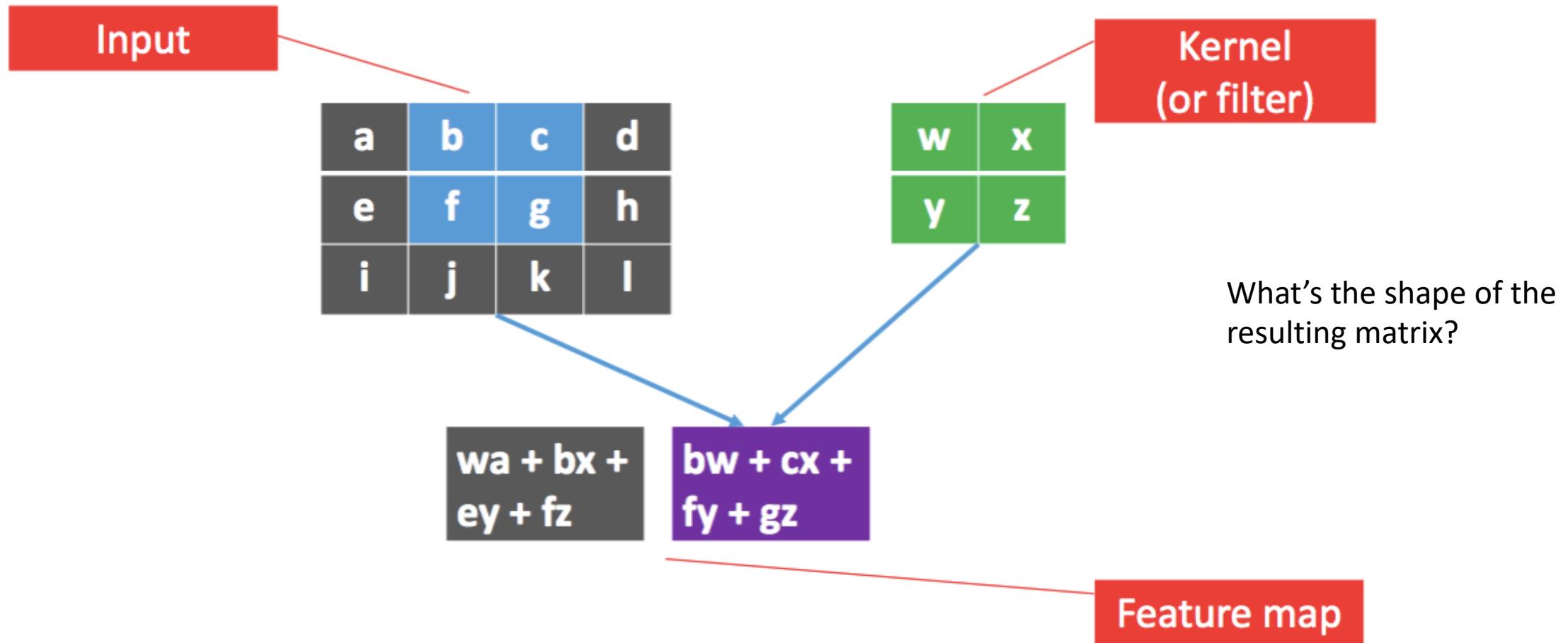
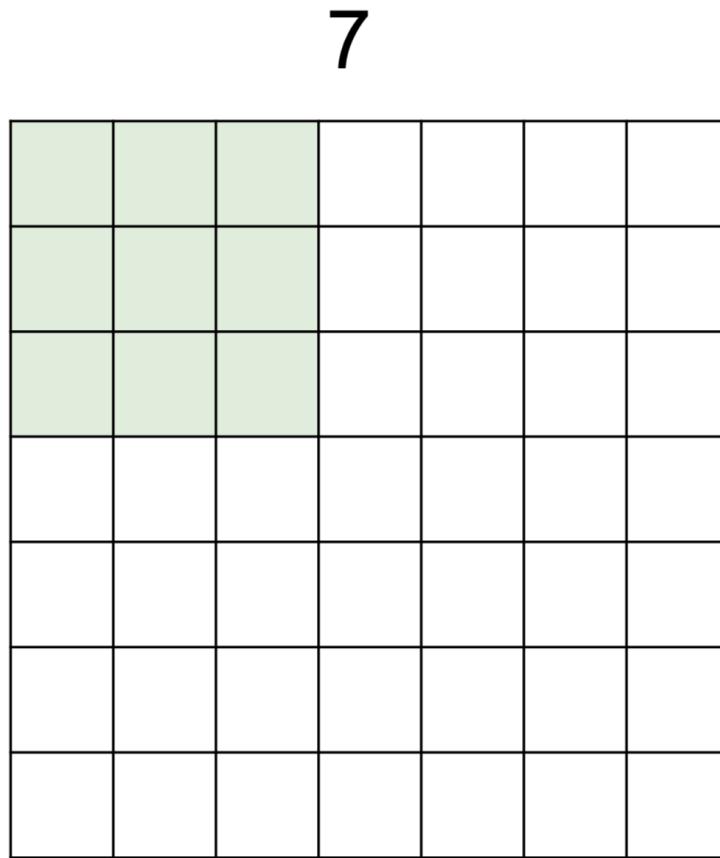


Illustration 2

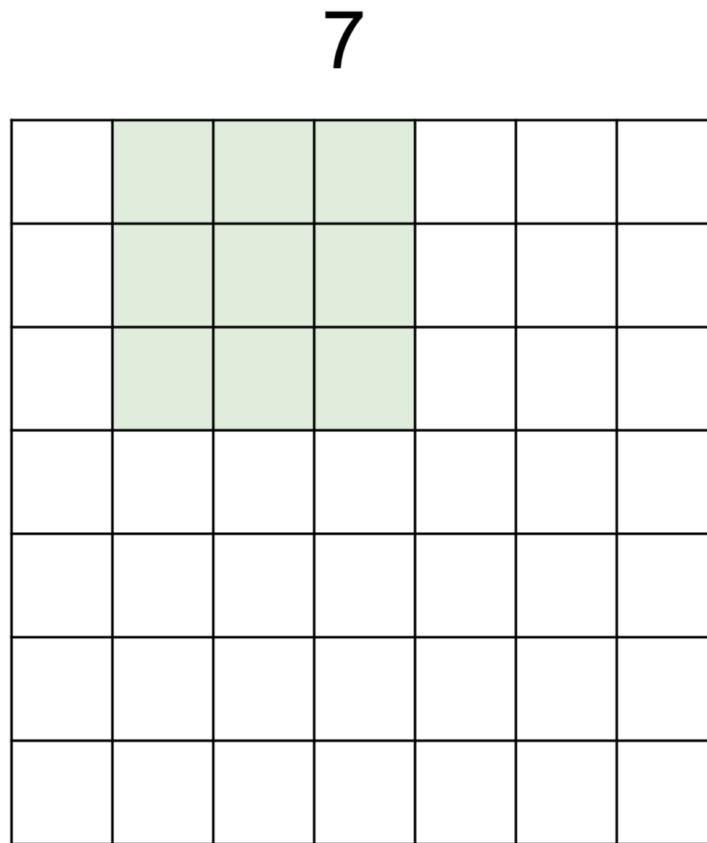


A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

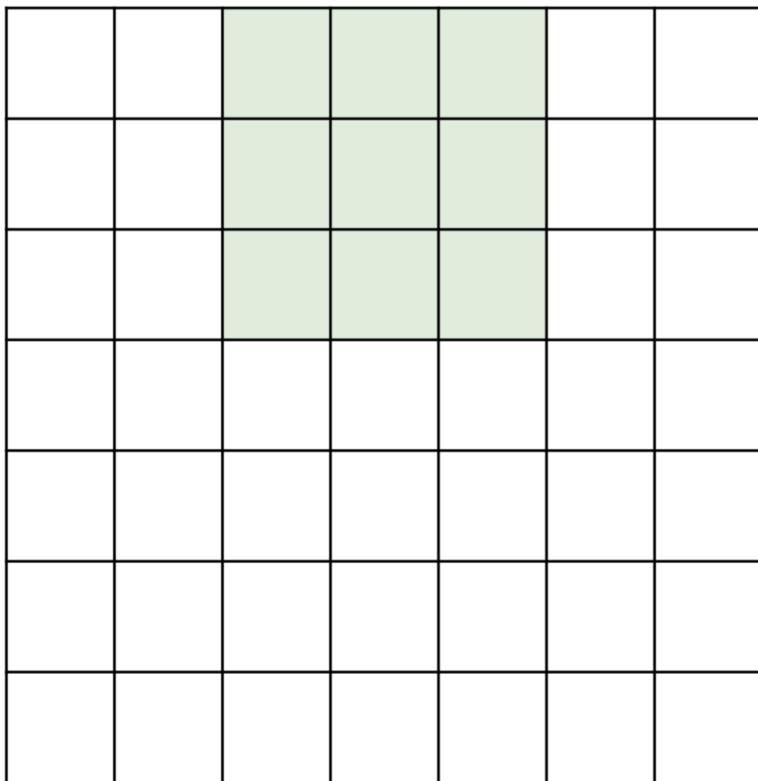


7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

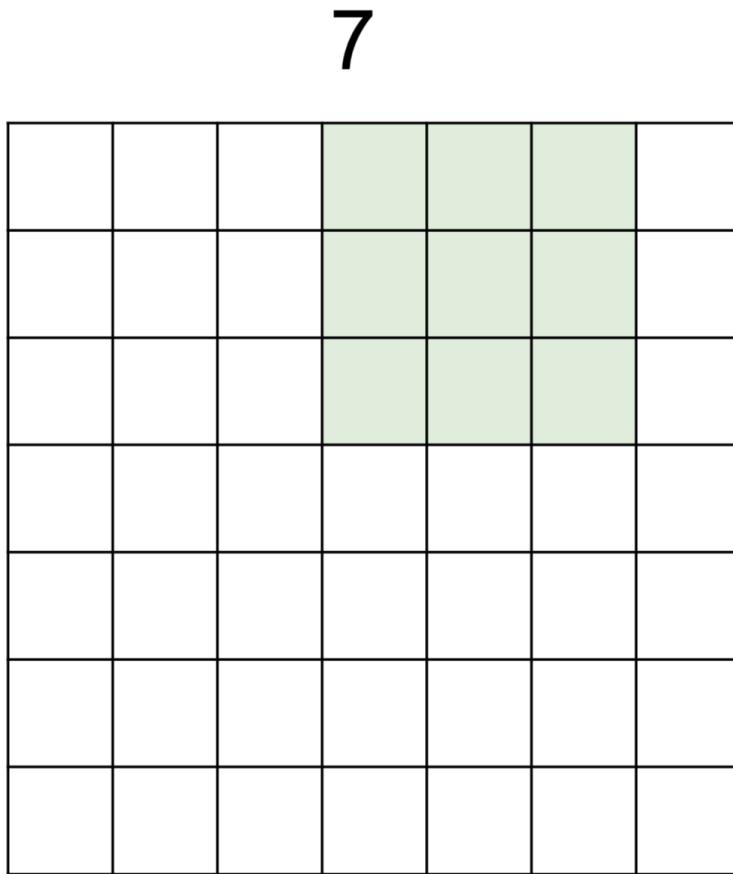
7



7x7 input (spatially)
assume 3x3 filter

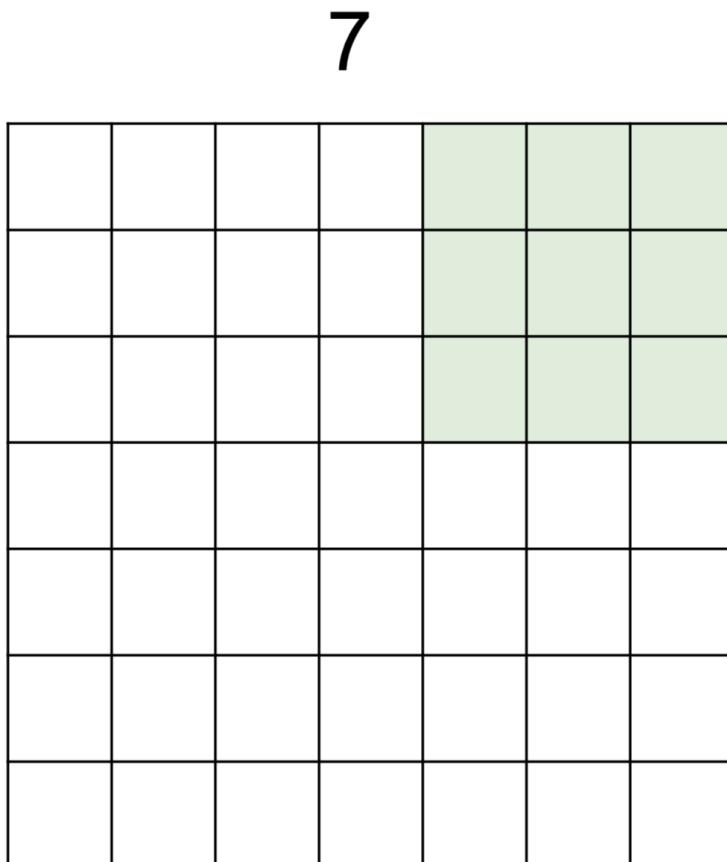
7

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

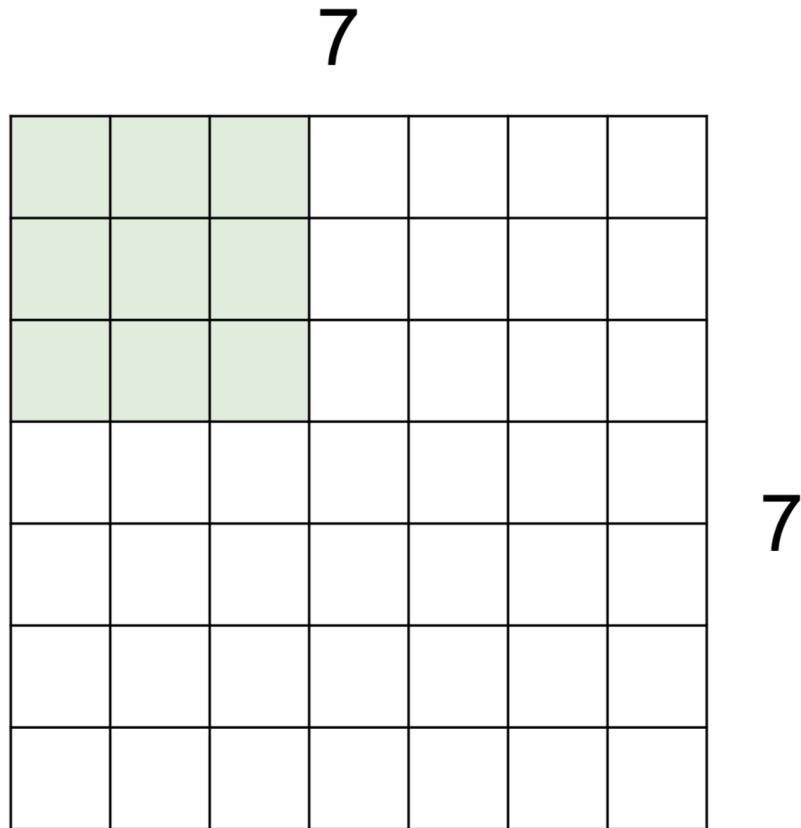
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

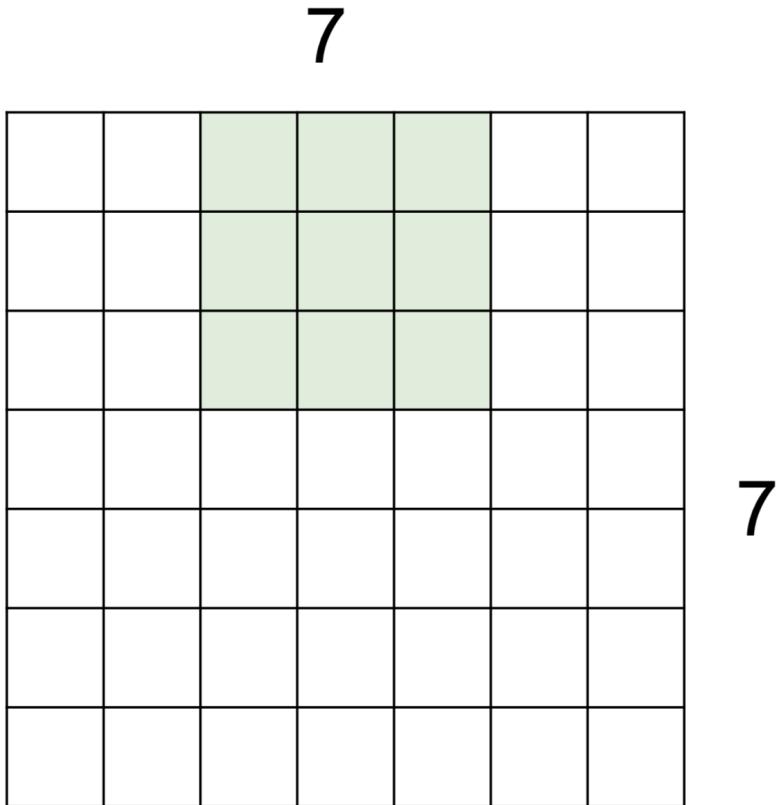
=> 5x5 output

A closer look at spatial dimensions:



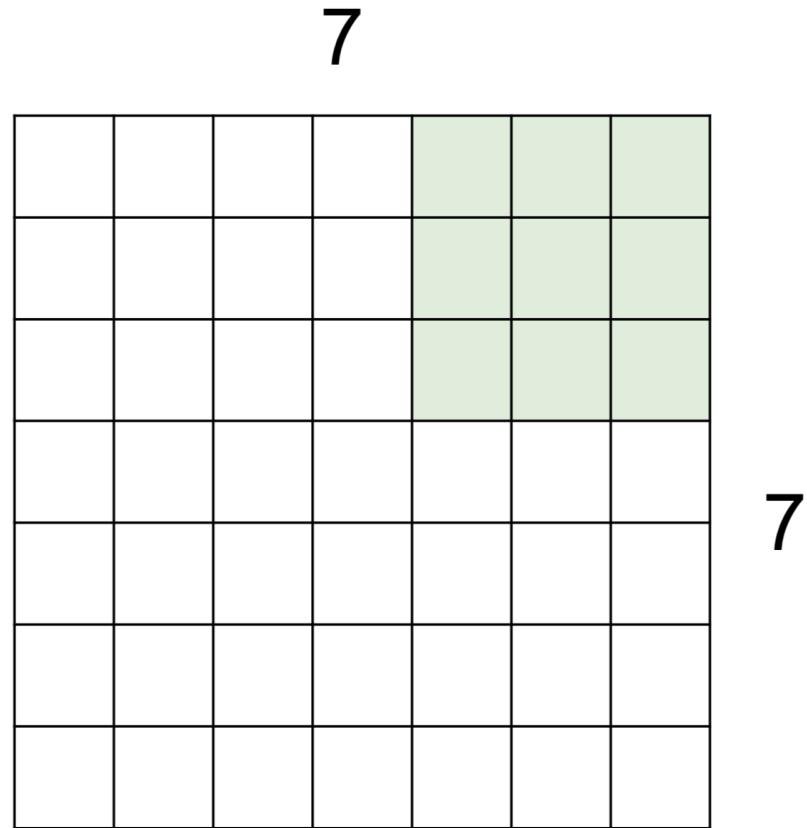
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



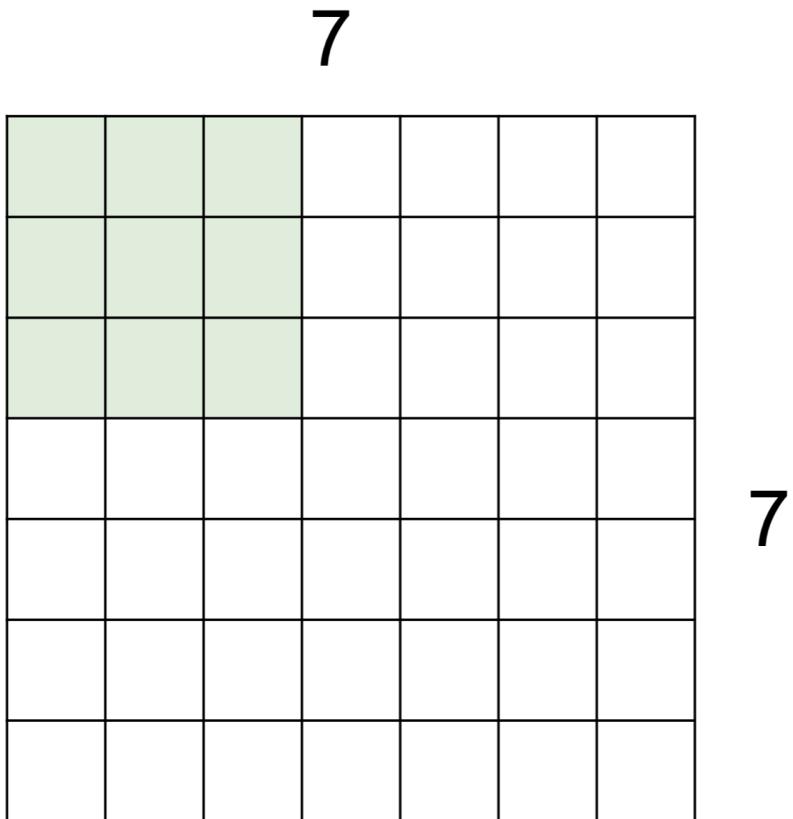
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



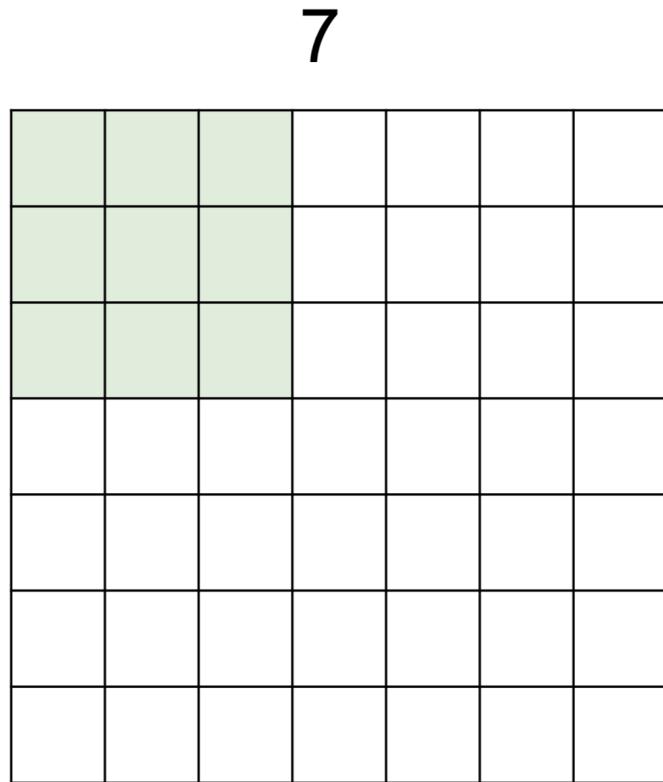
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

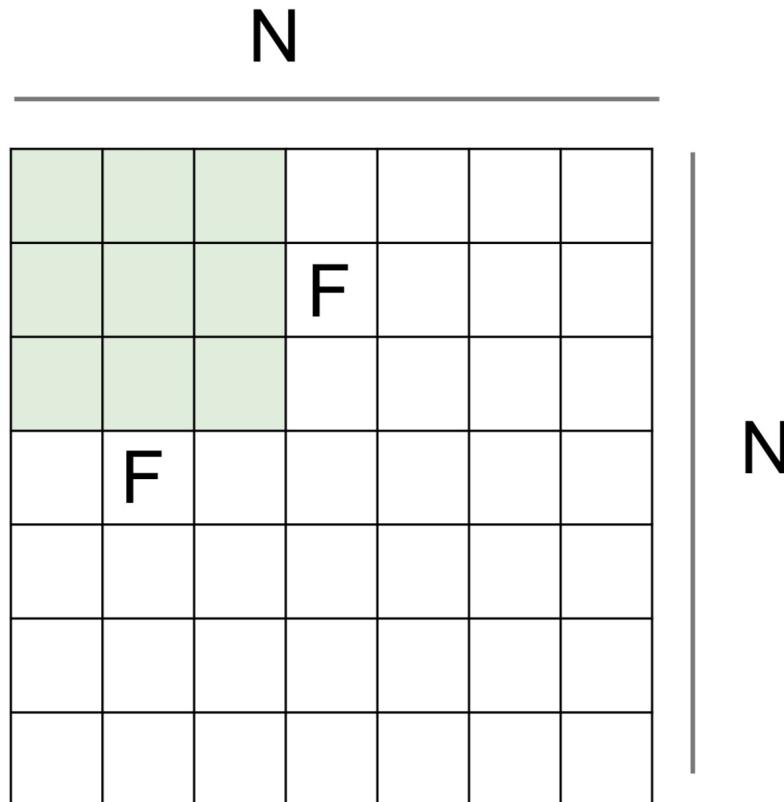
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

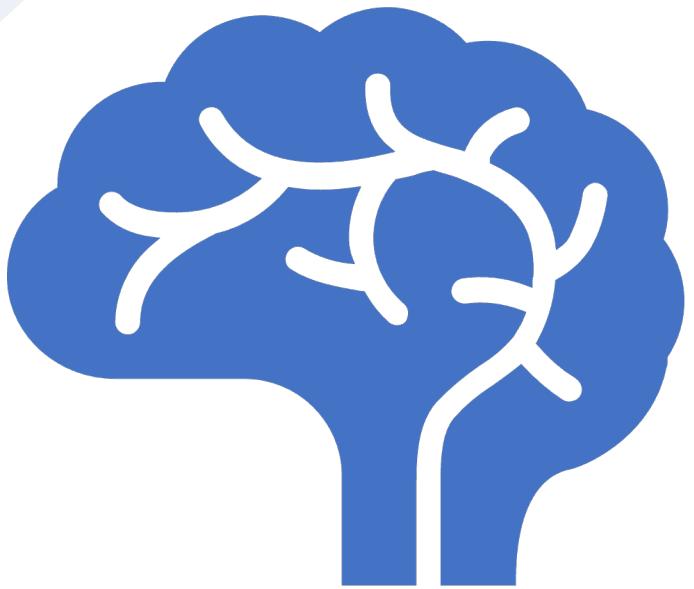
doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

A closer look at spatial dimensions:



Output size:
 $(N - F) / \text{stride} + 1$

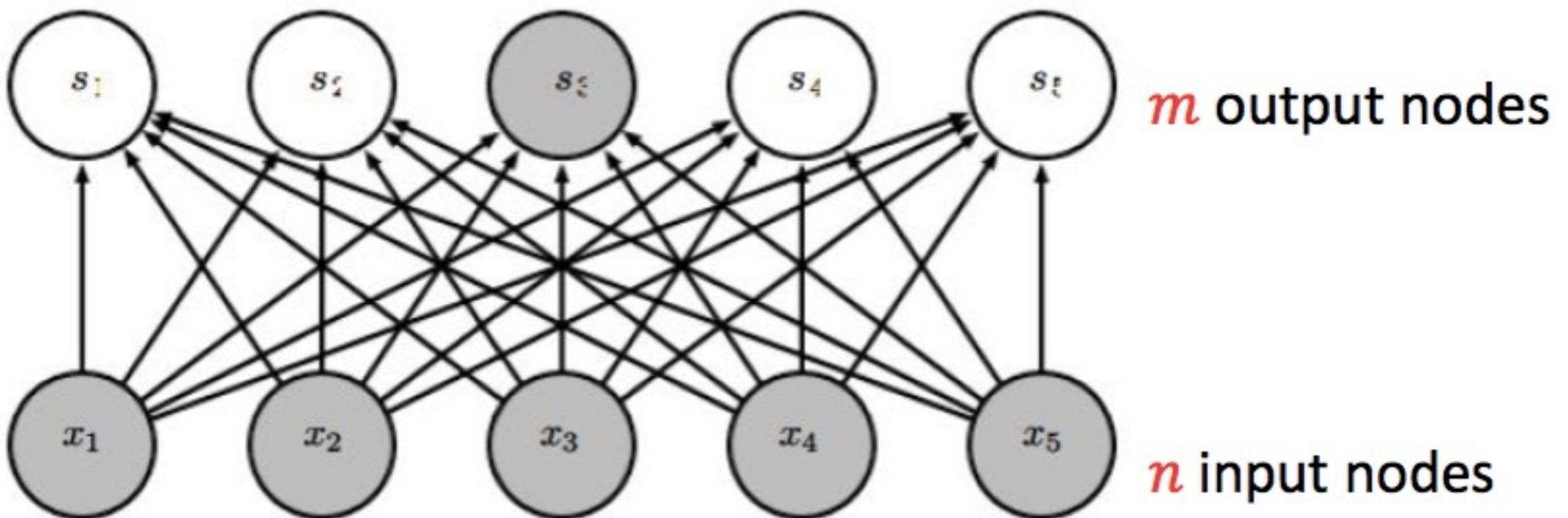
e.g. $N = 7$, $F = 3$:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 :\backslash$



Advantage of Convolutional Layer

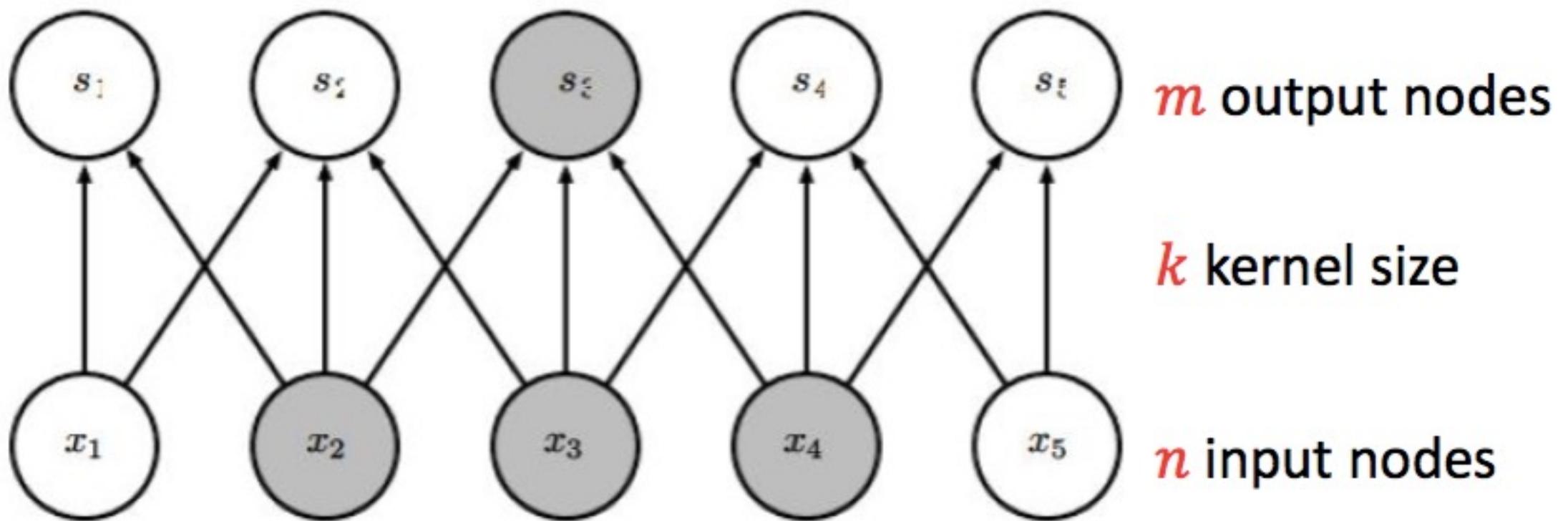
Advantage: sparse interaction

Fully connected layer, $m \times n$ edges



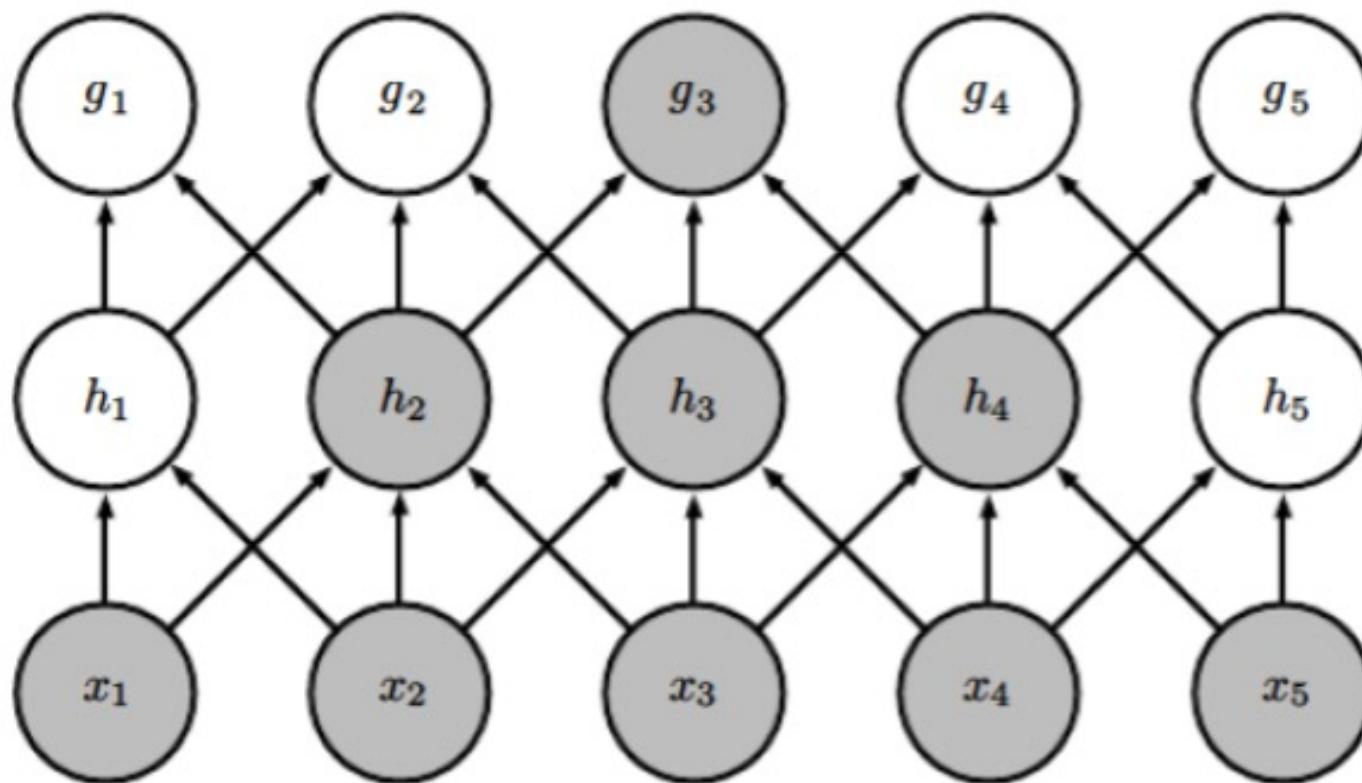
Advantage: sparse interaction

Convolutional layer, $\leq m \times k$ edges

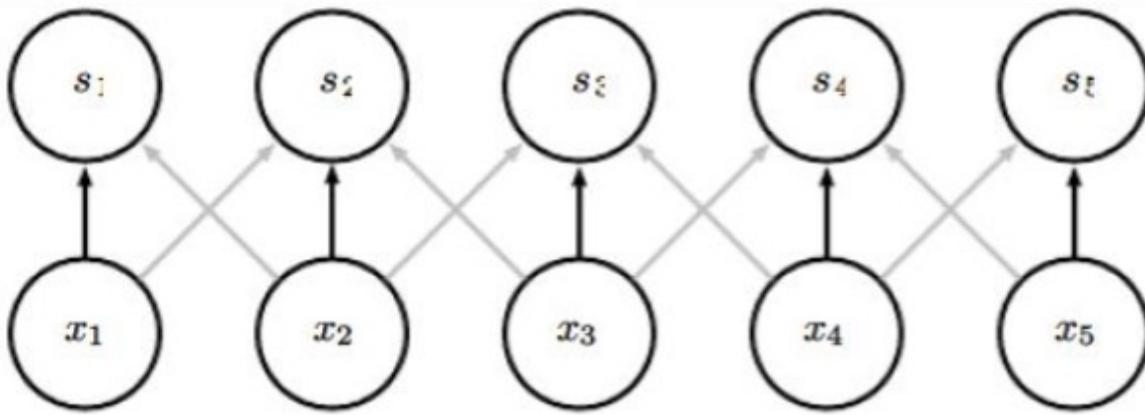


Advantage: sparse interaction

Multiple convolutional layers: larger receptive field



Advantage: parameter sharing/weight tying



The same kernel
are used repeatedly.
E.g., the black edge
is the same weight
in the kernel.

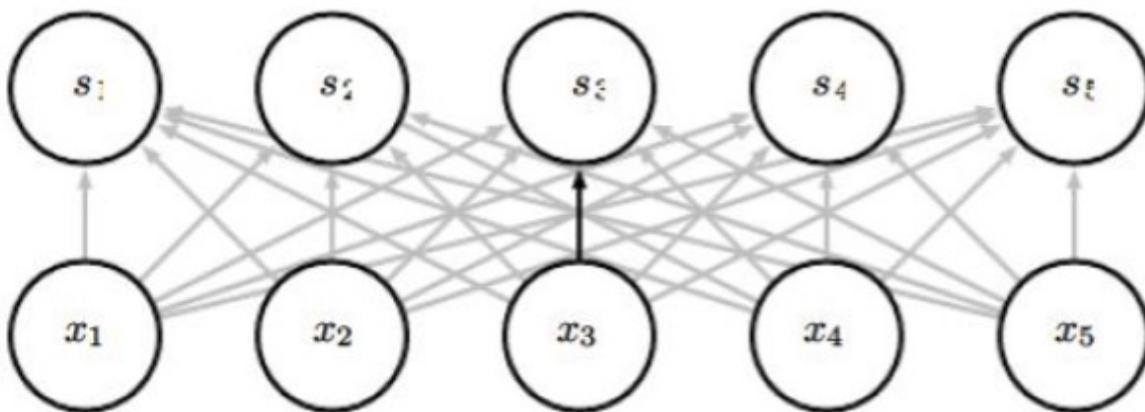
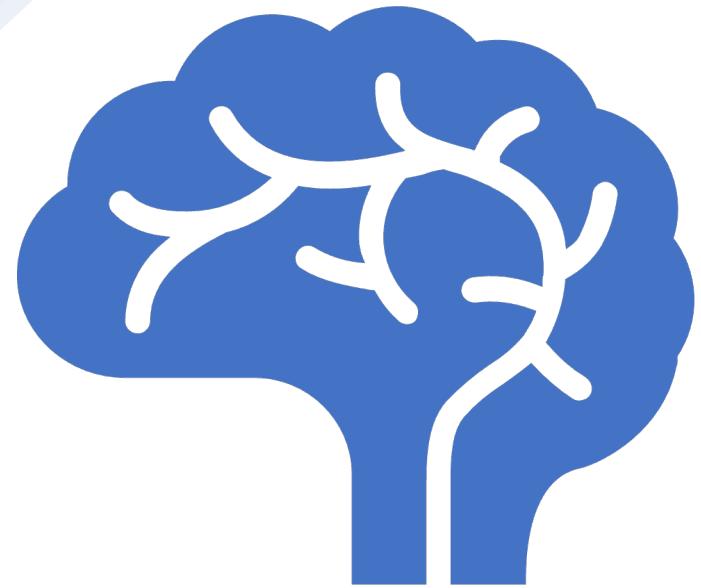


Figure from *Deep Learning*,
by Goodfellow, Bengio,
and Courville

Advantage: equivariant representations

- Equivariant: transforming the input = transforming the output
- Example: input is an image, transformation is shifting
- $\text{Convolution}(\text{shift}(\text{input})) = \text{shift}(\text{Convolution}(\text{input}))$
- Useful when care only about the **existence** of a pattern, rather than the **location**



Zero-Padding Layer

Zero-Padding

filter

w	x
y	z

Input

a	b	c	d
e	f	g	h
i	j	k	l



0	0	0	0	0	0
0	a	b	c	d	0
0	e	f	g	h	0
0	i	j	k	l	0
0	0	0	0	0	0

Convolutional layer changes the input dimensions,
But sometimes we like the input and output of a
convolutional layer to be of the same shape.

What's the shape of the
resulting matrix?



ReLU Layer

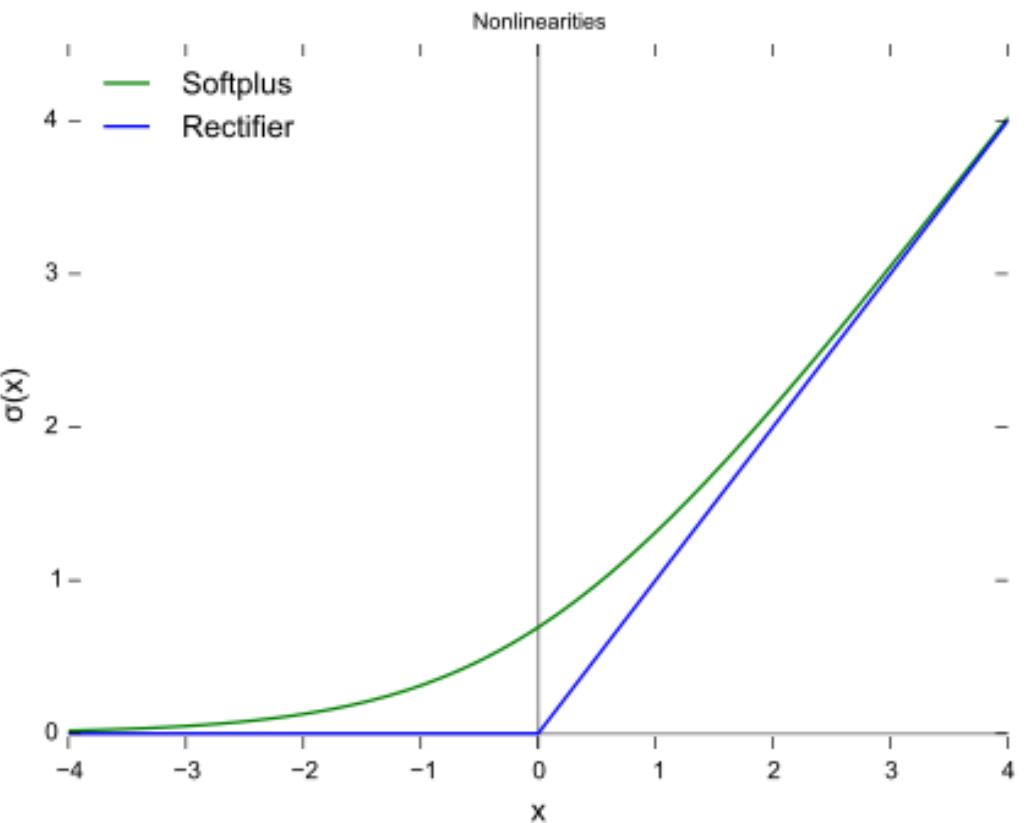
ReLU (rectified linear unit)

- **rectifier** is an activation function defined as the positive part of its argument

$$f(x) = \max(0, x)$$

- A smooth approximation to the rectifier is the analytic function

$$f(x) = \log(1+e^x)$$





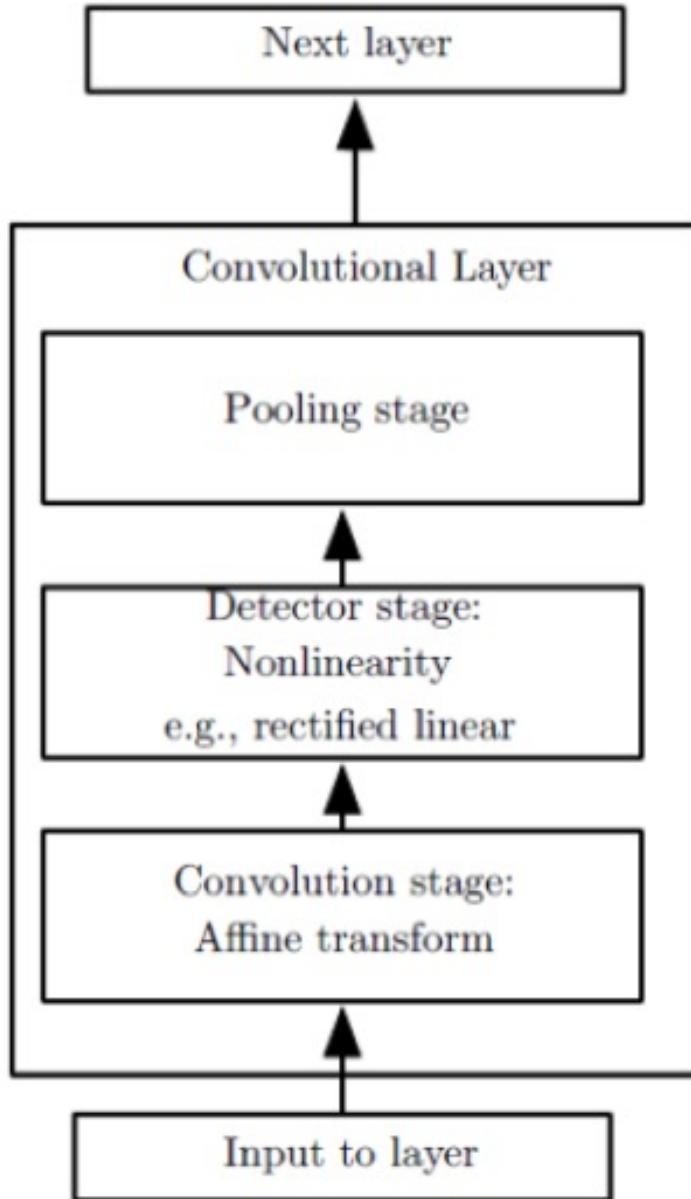
Pooling Layer

Pooling

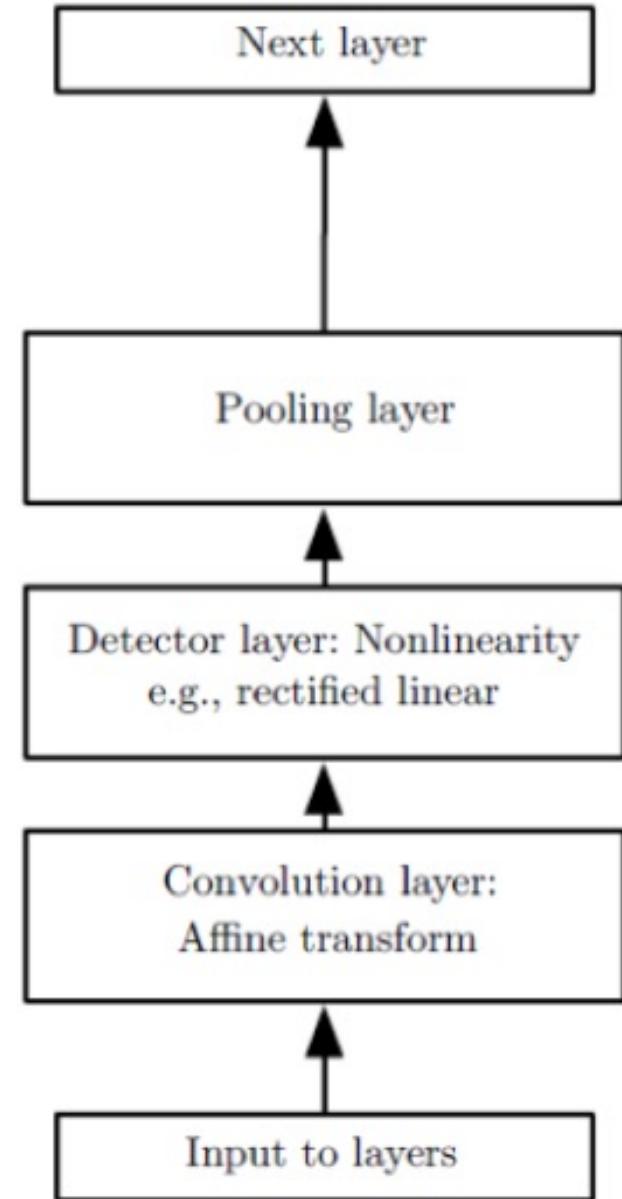
- Pooling layer is frequently used in convolutional neural networks with the purpose to progressively reduce the spatial size of the representation to reduce the amount of features and the computational complexity of the network.

Terminology

Complex layer terminology

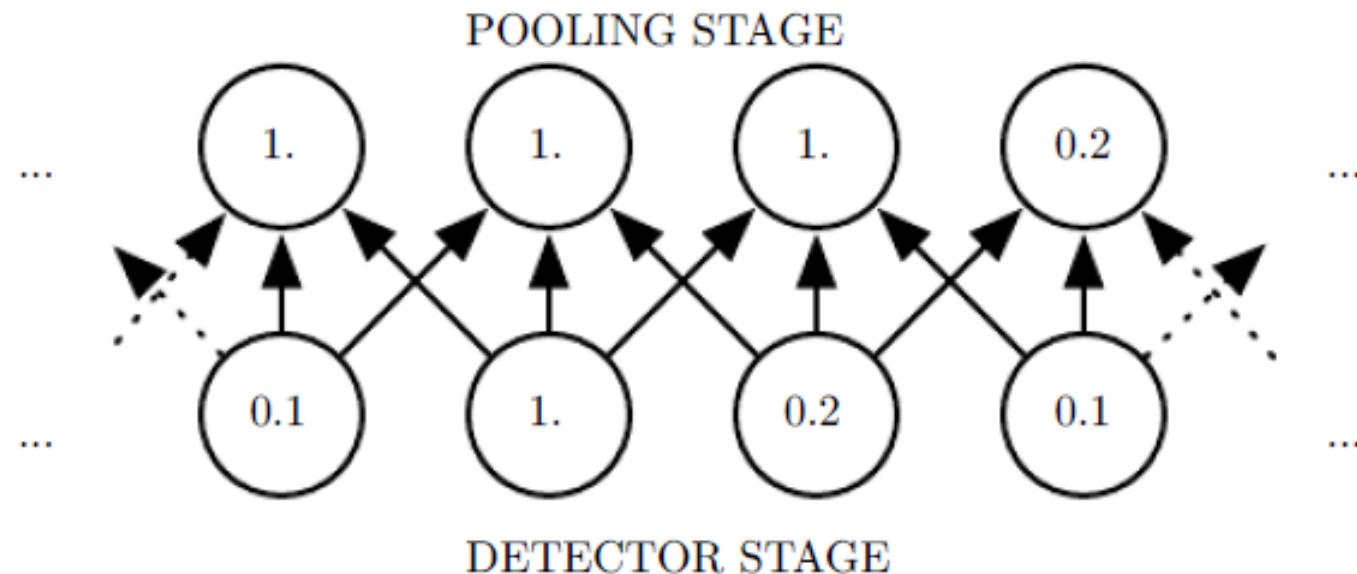


Simple layer terminology

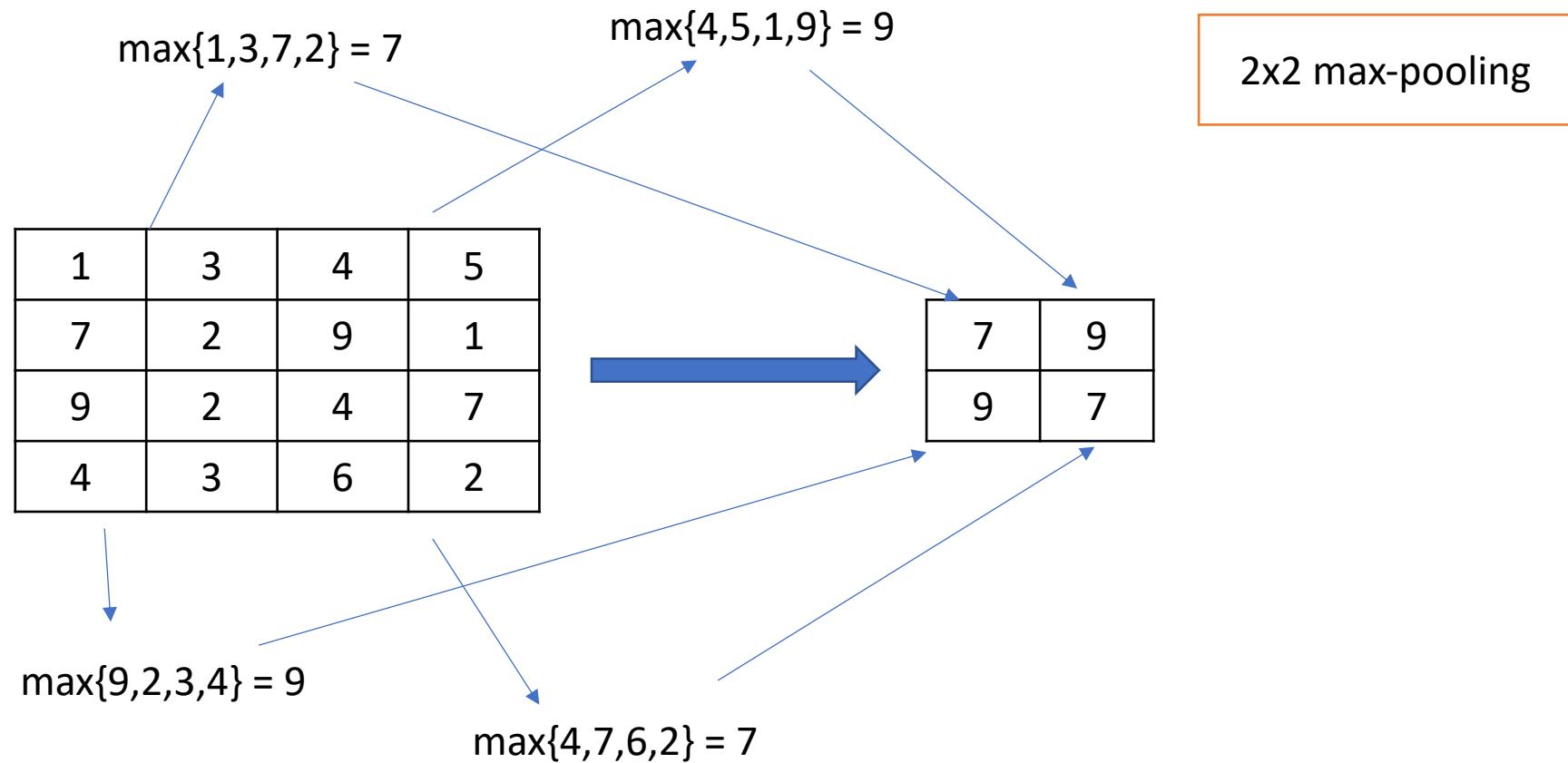


Pooling

- Summarizing the input (i.e., output the max of the input)

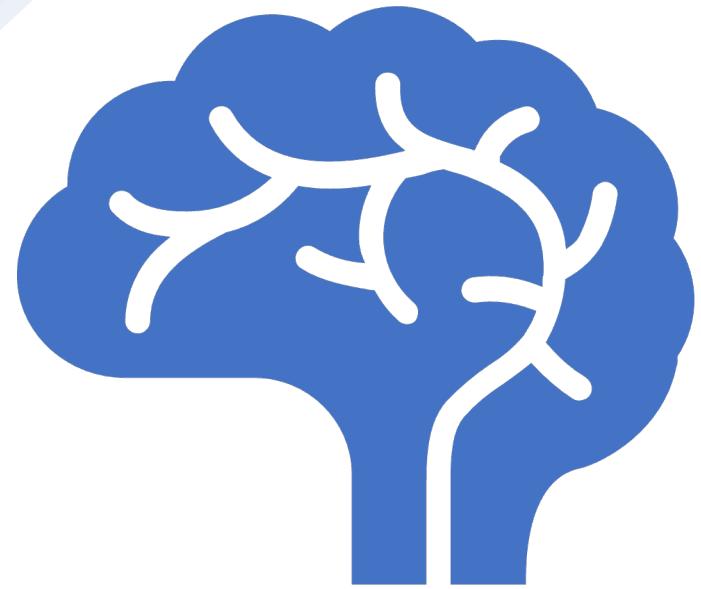


Example: Max-pooling



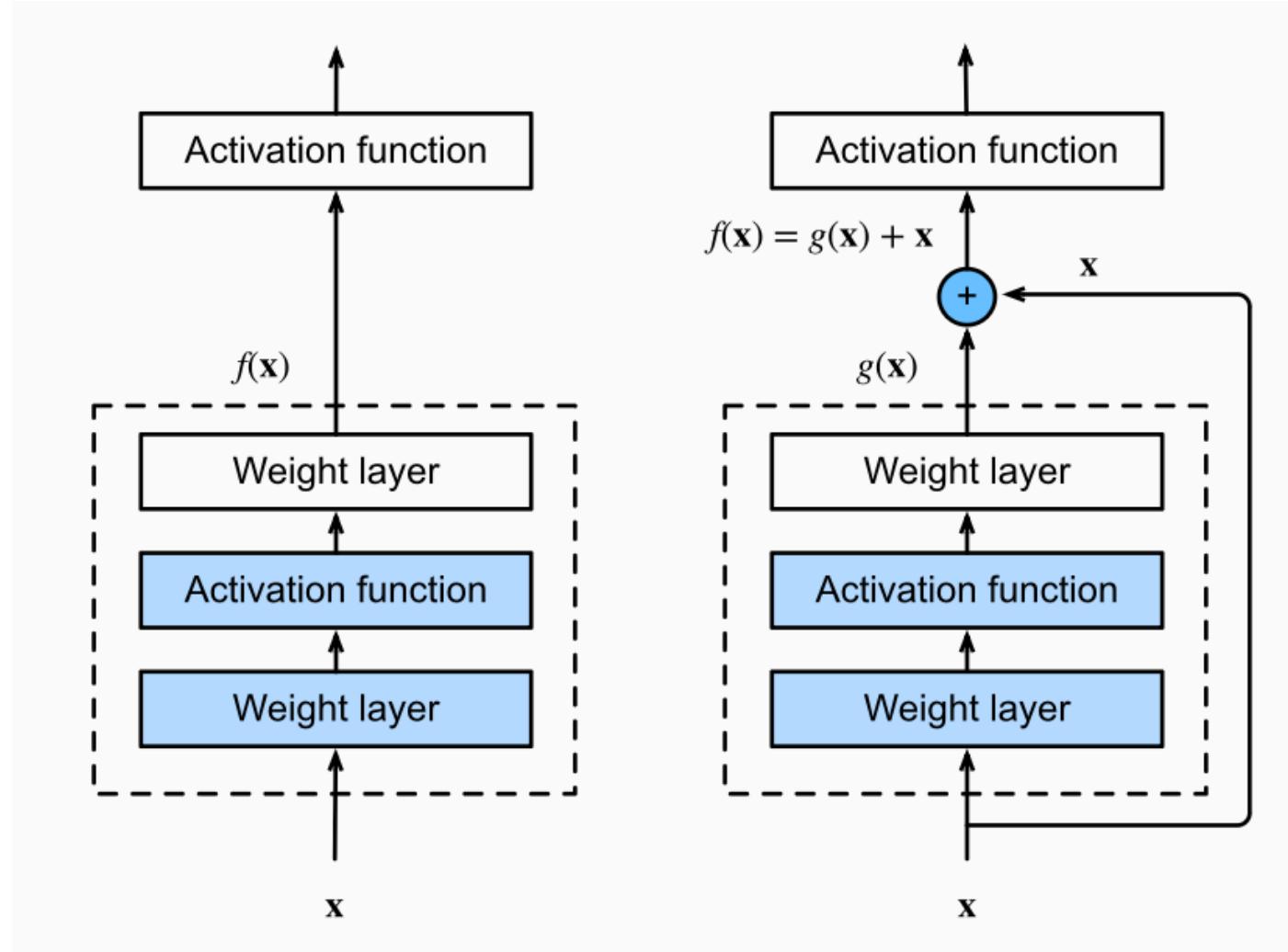
Motivation from neuroscience

- David Hubel and Torsten Wiesel studied early visual system in human brain (V1 or primary visual cortex), and won Nobel prize for this
- V1 properties
 - 2D spatial arrangement
 - Simple cells: inspire convolution layers
 - Complex cells: inspire pooling layers



Residual Block

Residual Block



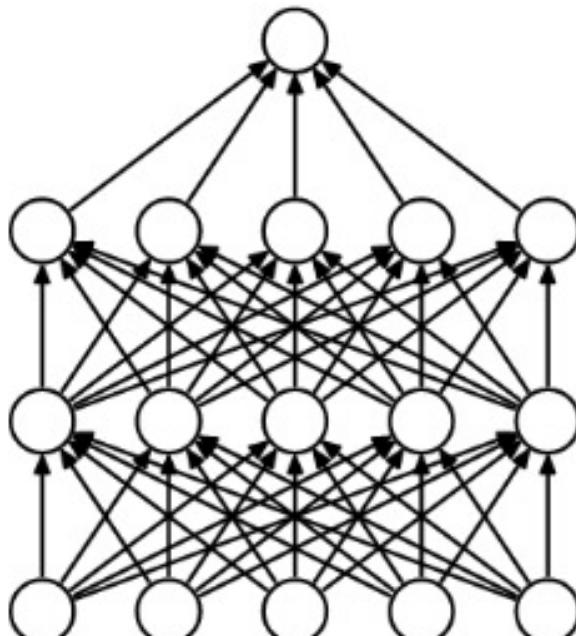
regular block

residual block

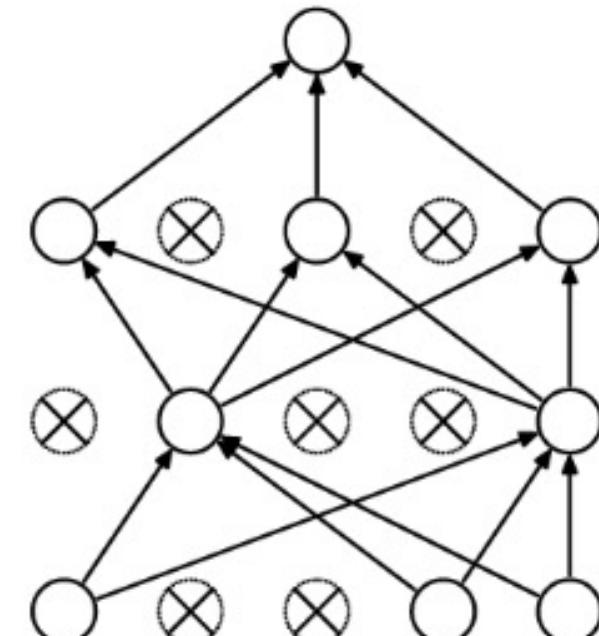


Dropout

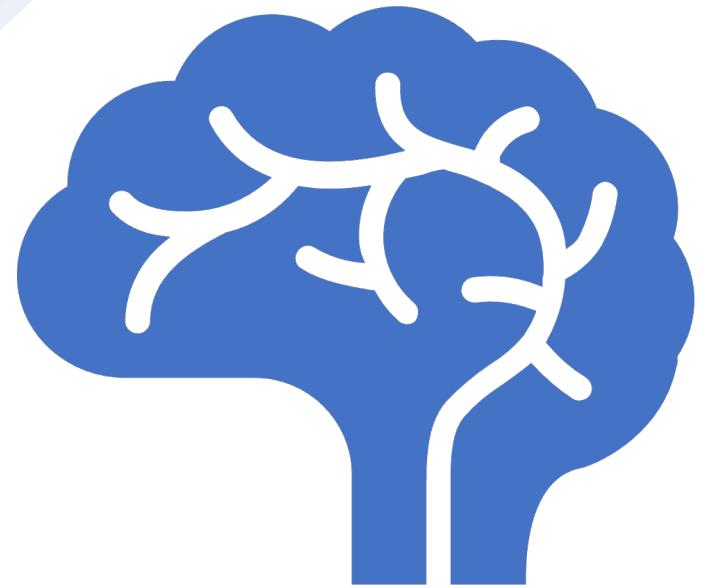
Dropout



(a) Standard Neural Net

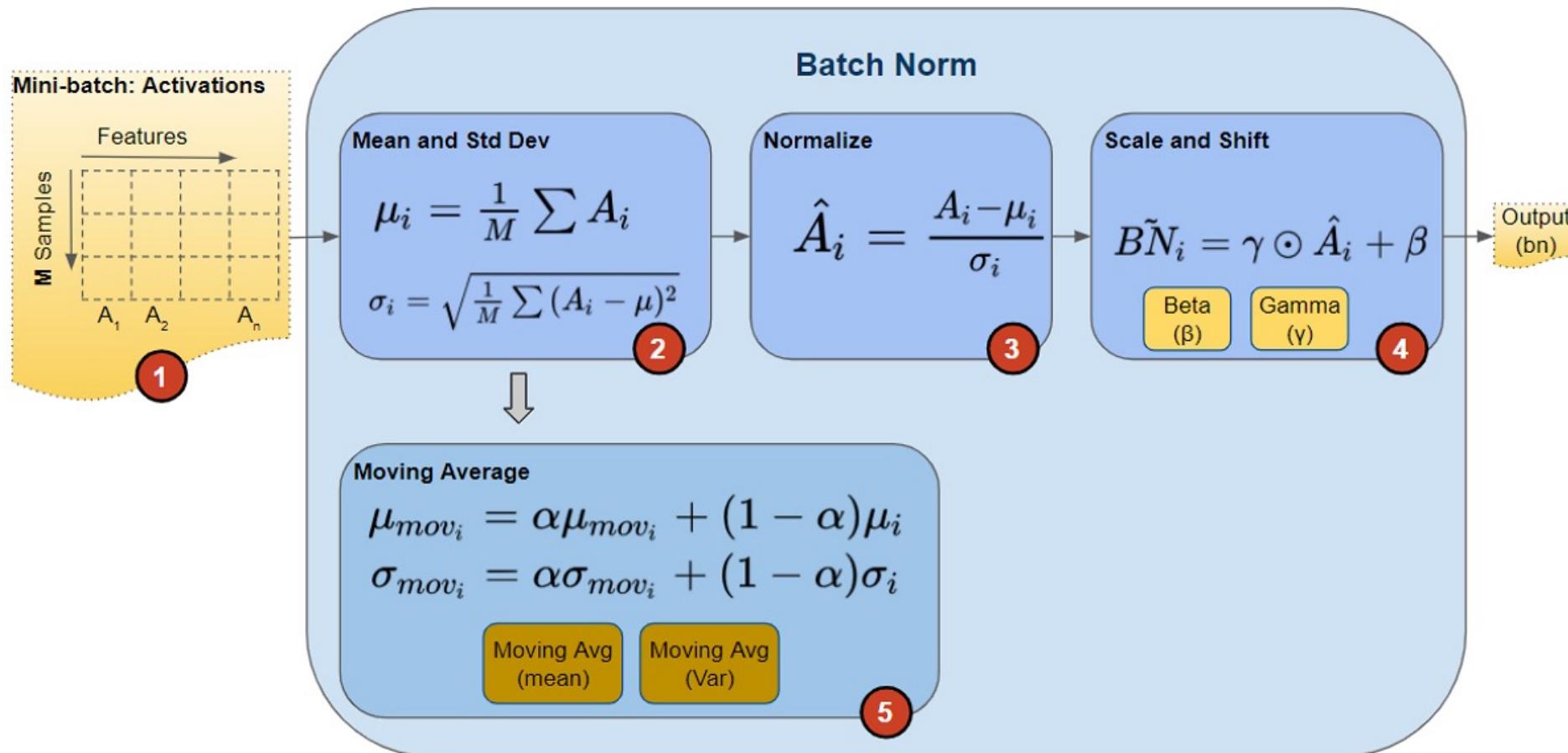


(b) After applying dropout.

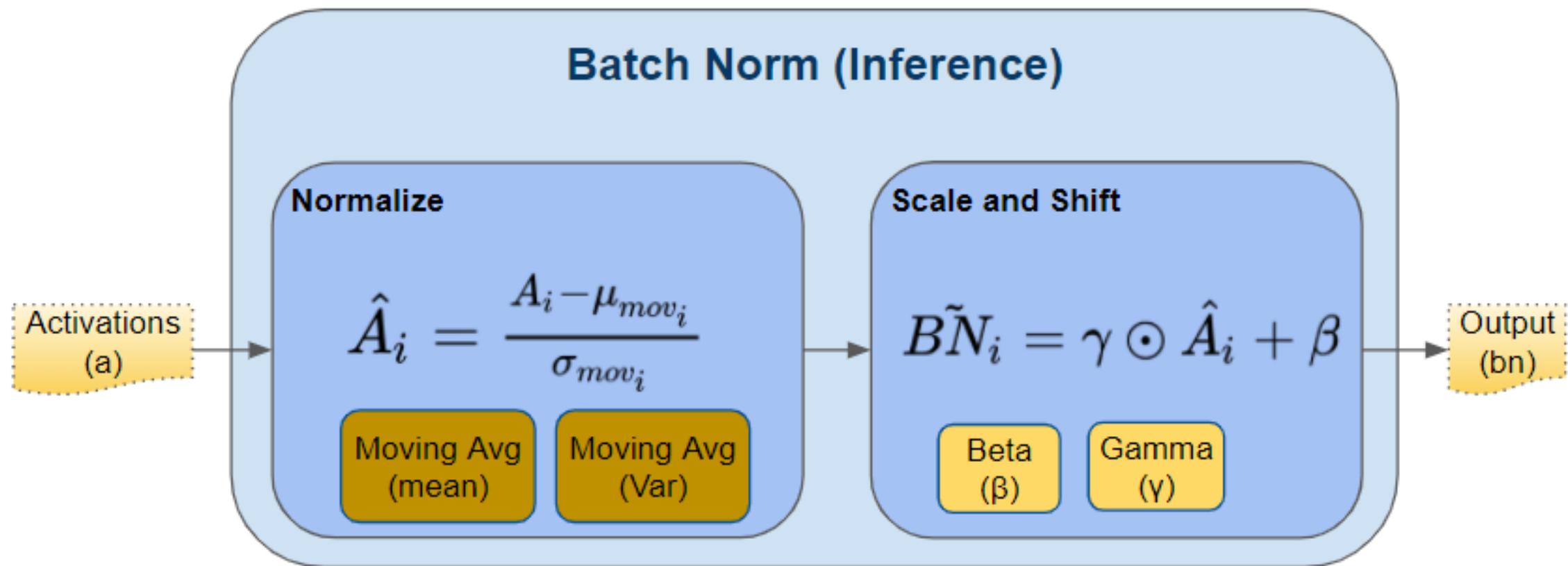


BatchNorm Layer

BatchNorm -- training



BatchNorm -- inference



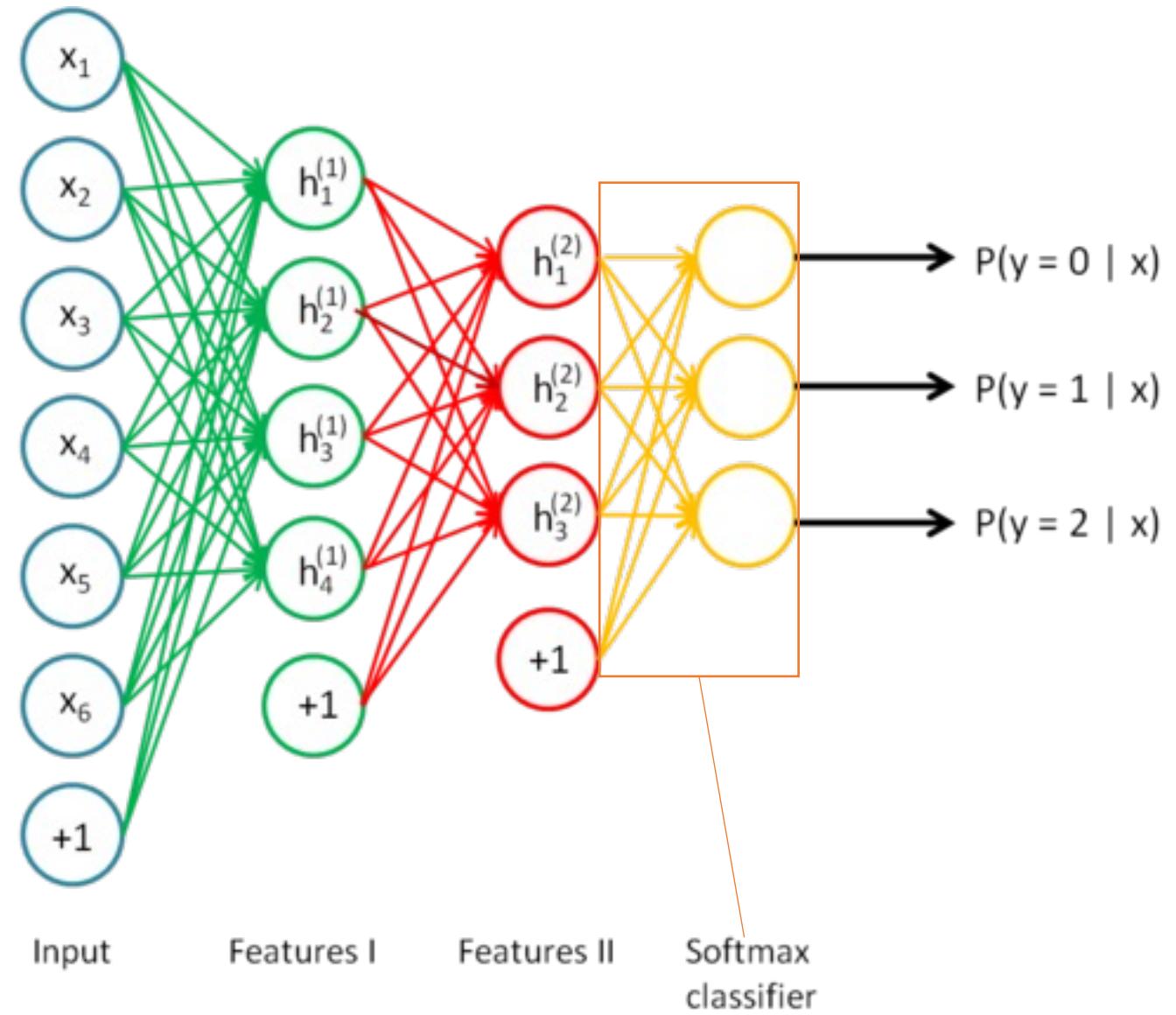


Softmax Layer

Softmax

- Recall that [logistic regression](#) produces a decimal between 0 and 1.0. For example, a logistic regression output of 0.8 from an email classifier suggests an 80% chance of an email being spam and a 20% chance of it being not spam. Clearly, the sum of the probabilities of an email being either spam or not spam is 1.0.
- **Softmax** extends this idea into a multi-class world. That is, Softmax assigns decimal probabilities to each class in a multi-class problem. Those decimal probabilities must add up to 1.0. This additional constraint helps training converge more quickly than it otherwise would.

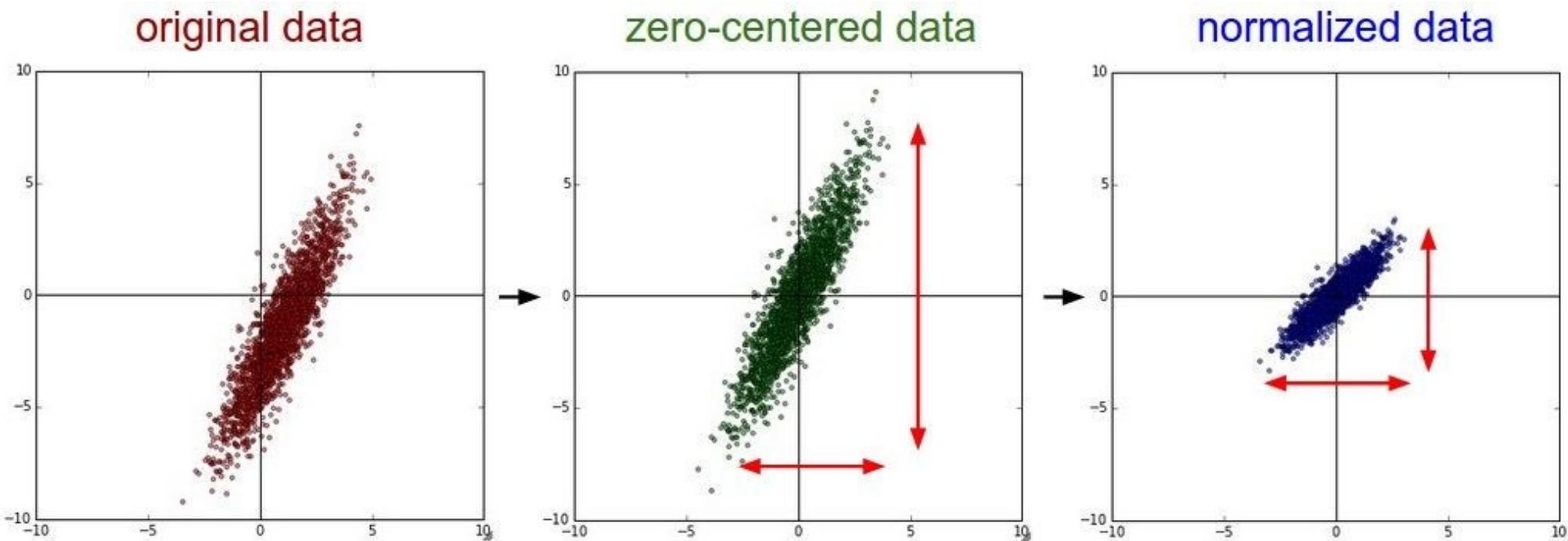
Softmax



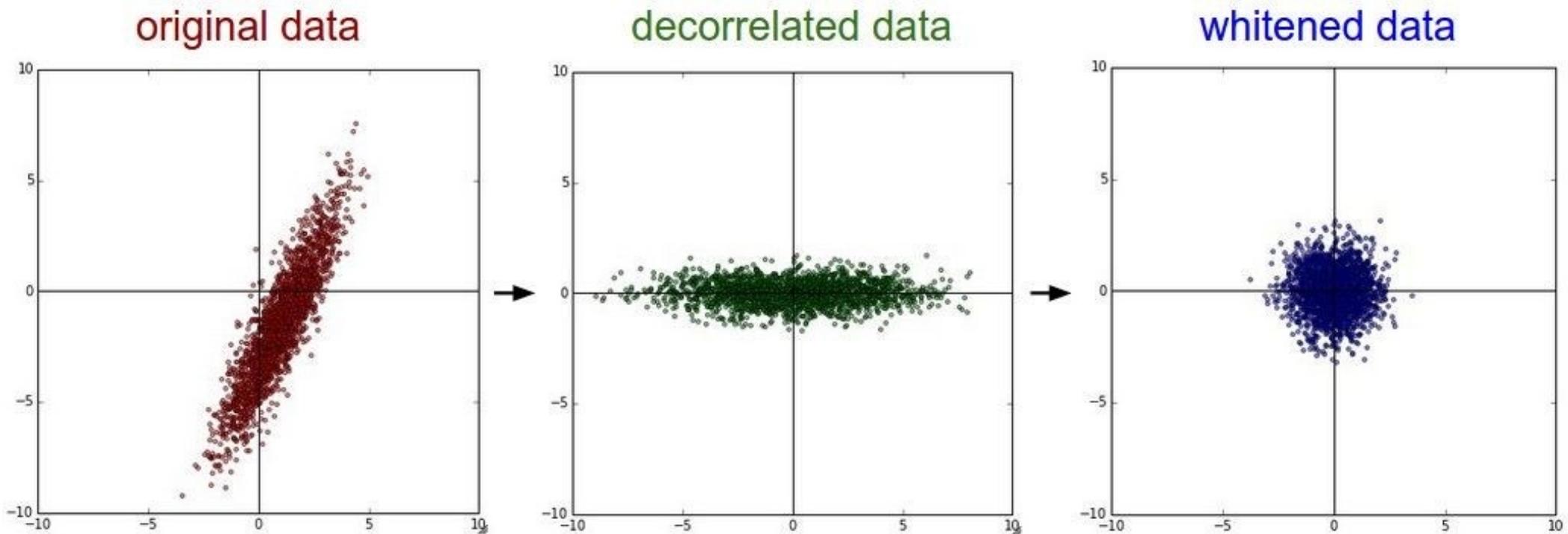


Preprocessing data

Preprocessing data



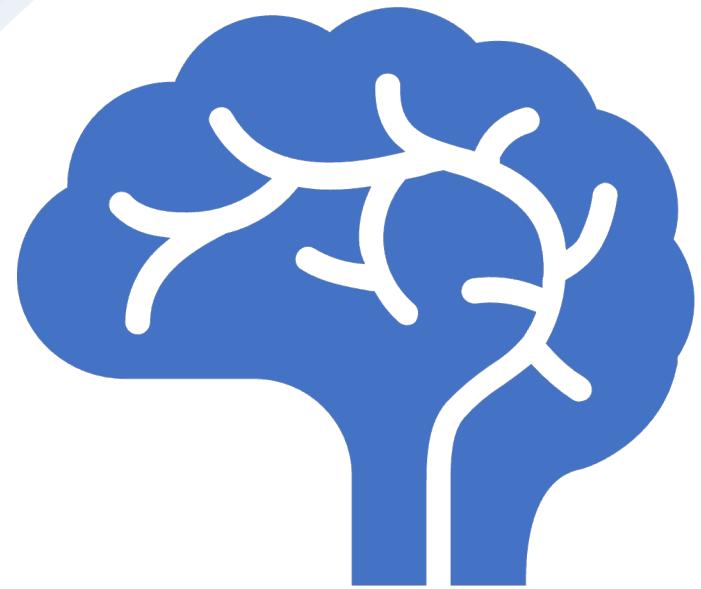
Preprocessing data



Preprocessing data

e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet) (mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet) (mean along each channel = 3 numbers)

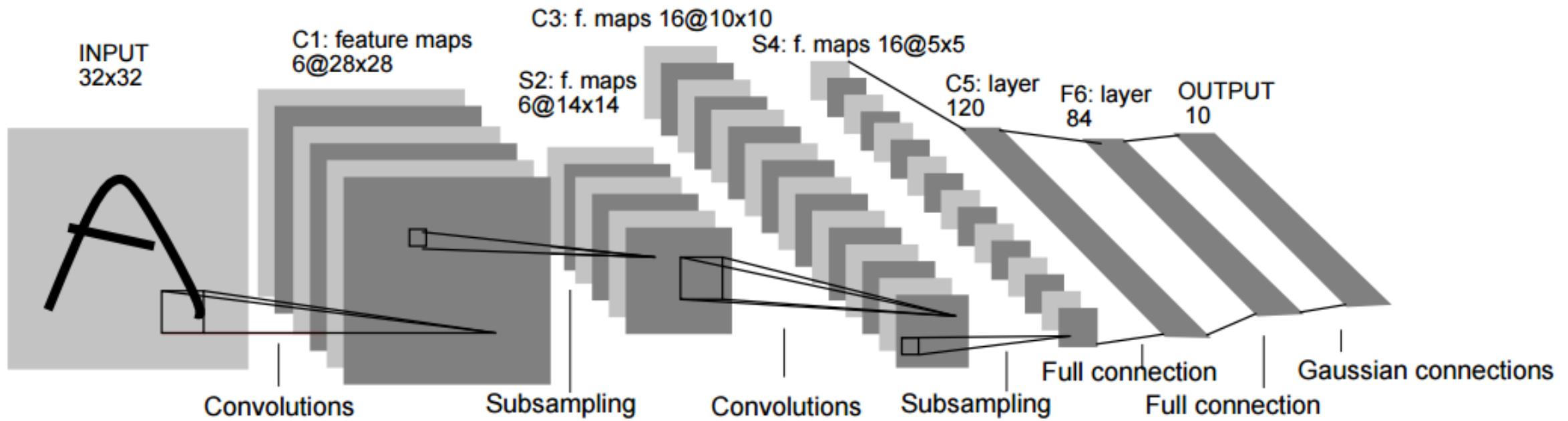


Example: LeNet

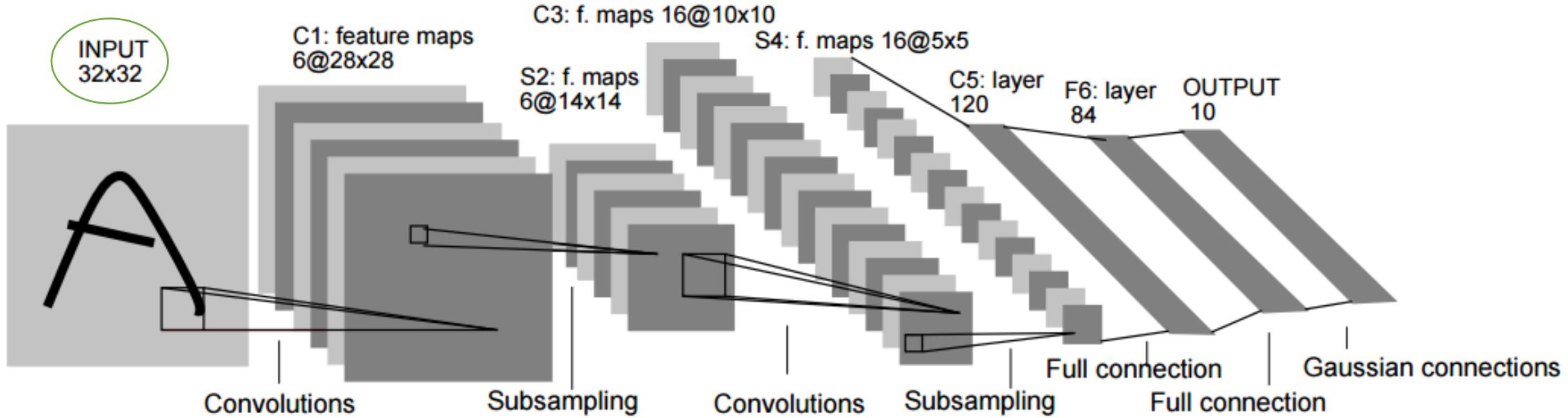
LeNet-5

- Proposed in “*Gradient-based learning applied to document recognition*”,
*by Yann LeCun, Leon Bottou, Yoshua Bengio and Patrick Haffner, in
Proceedings of the IEEE, 1998*
- Apply **convolution** on 2D images (MNIST) and use **backpropagation**
- Structure: 2 convolutional layers (with pooling) + 3 fully connected layers
 - Input size: 32x32x1
 - Convolution kernel size: 5x5
 - Pooling: 2x2

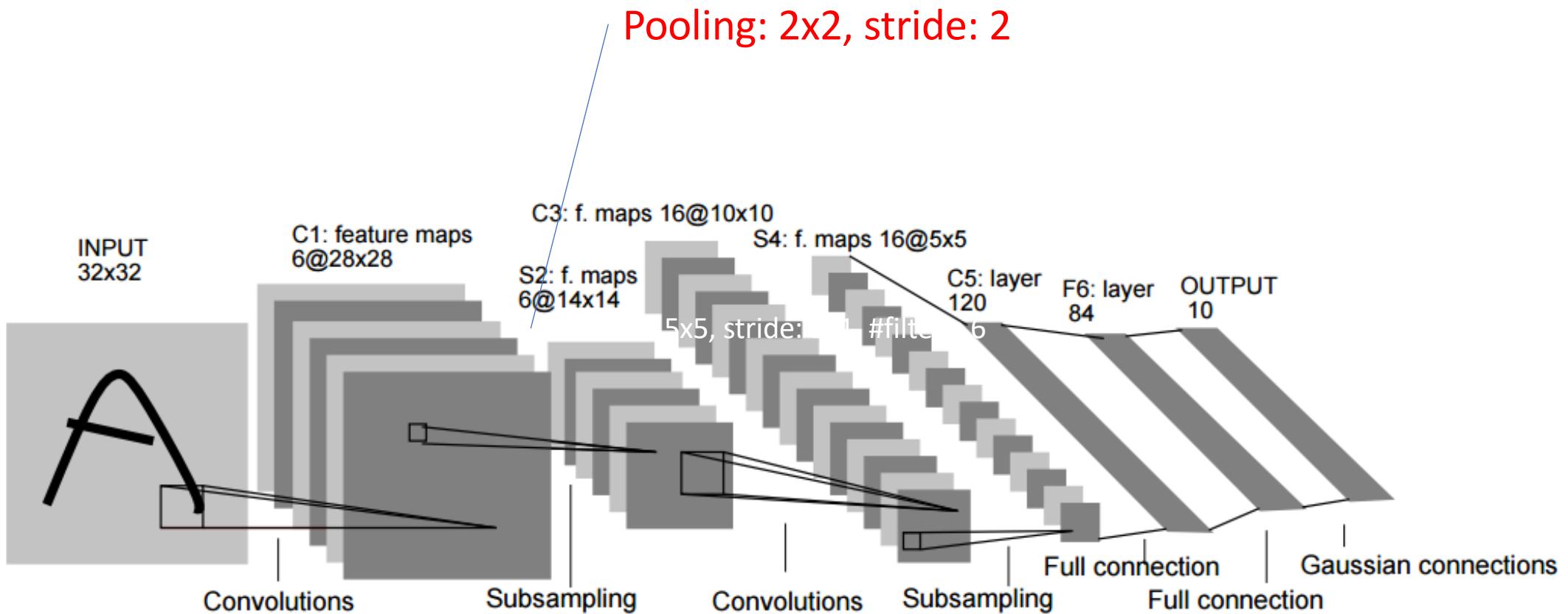
LeNet-5



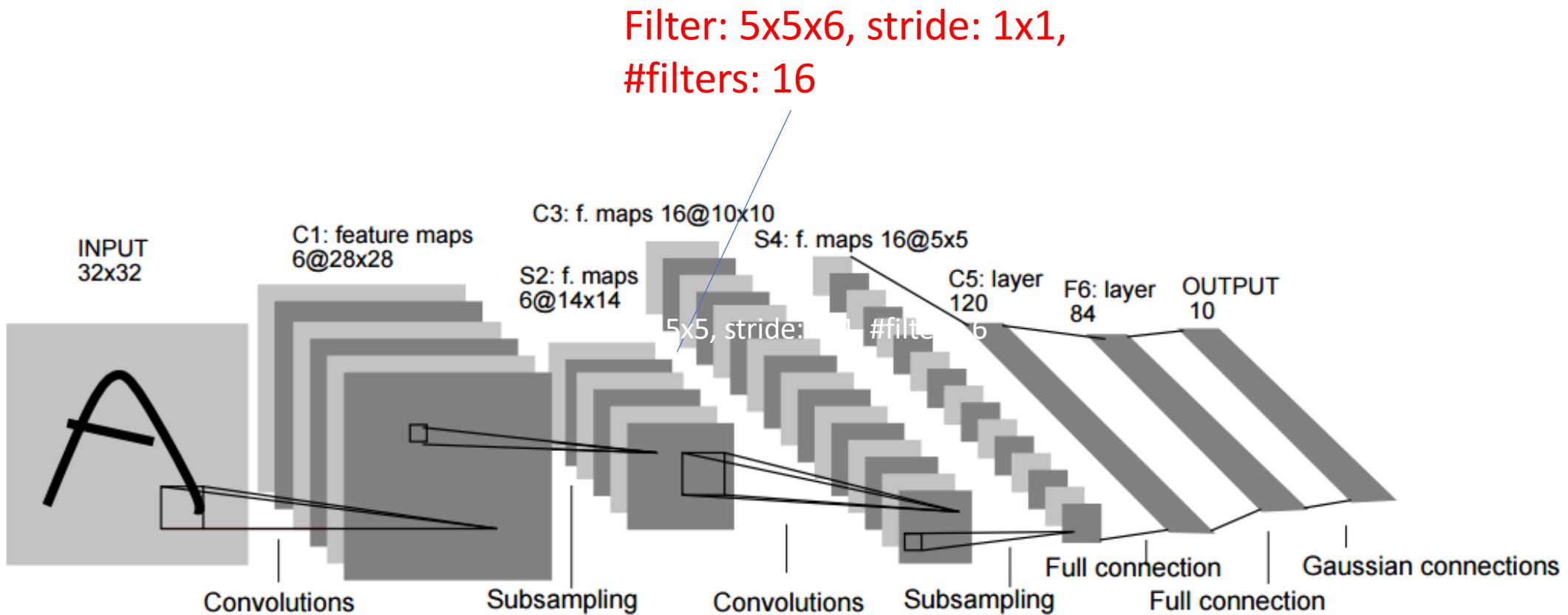
LeNet-5



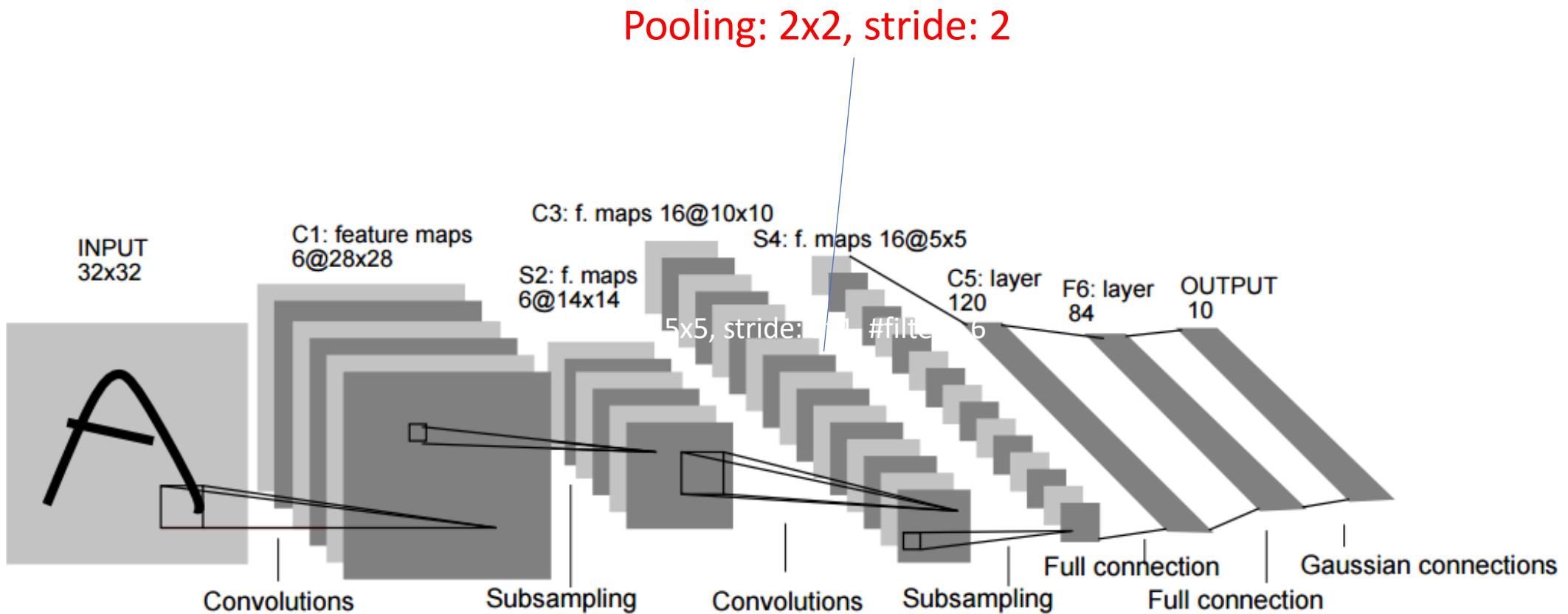
LeNet-5



LeNet-5

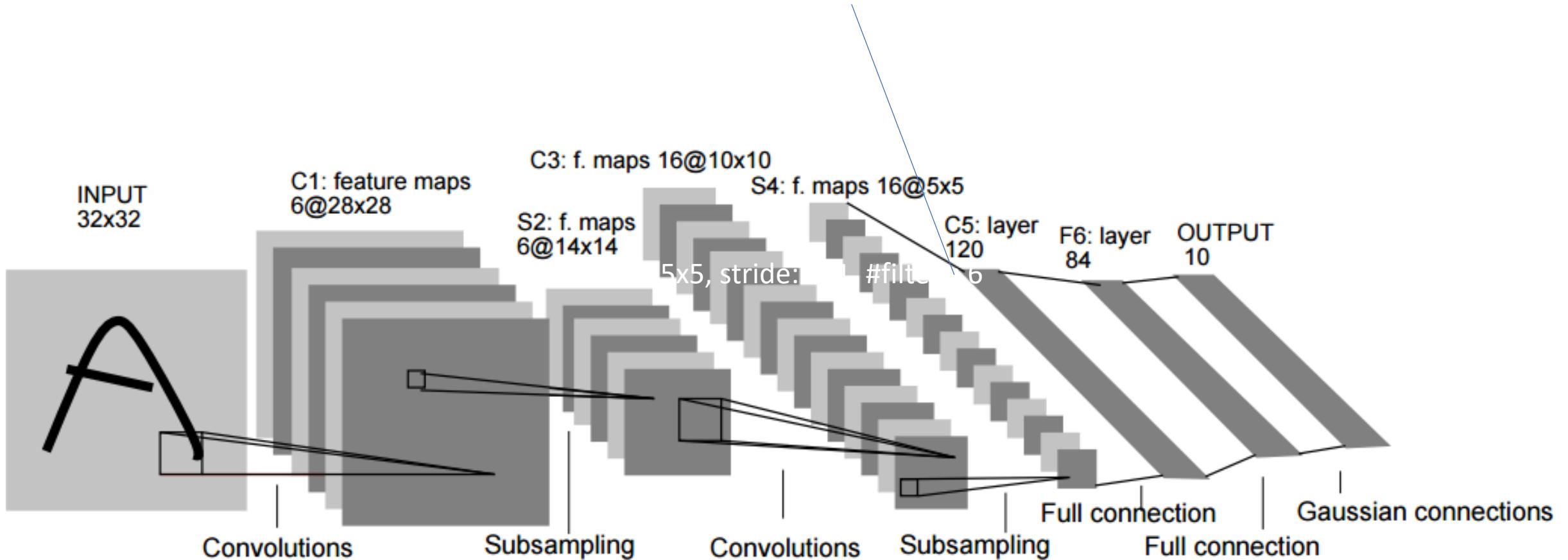


LeNet-5



LeNet-5

Weight matrix: 400x120



LeNet-5

