



**POLITECNICO**  
MILANO 1863



# Hardware Architectures for Embedded and Edge AI

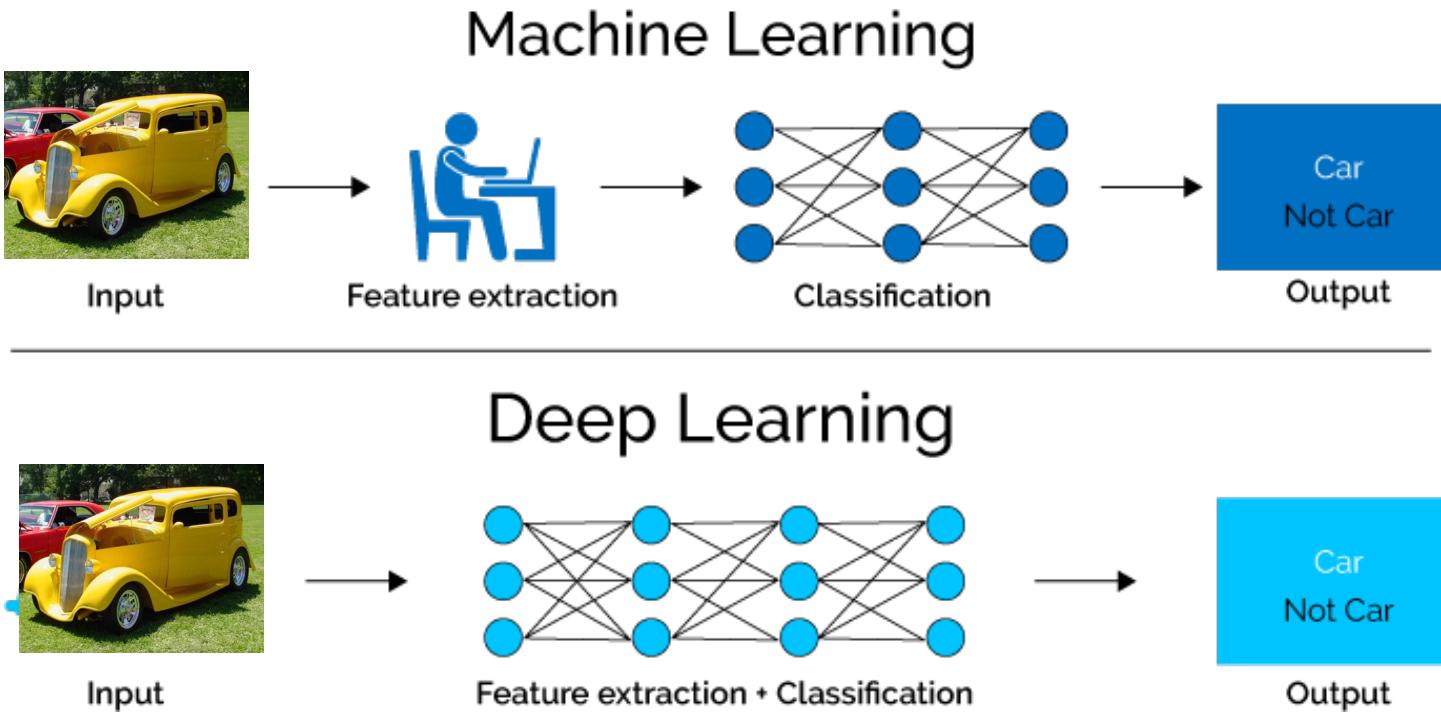
*Prof Manuel Roveri – manuel.roveri@polimi.it*

*Lecture 5 – Deep Learning for Embedded and Edge AI*

# The basics of deep learning



# Deep representation of knowledge



# How does TinyML work?

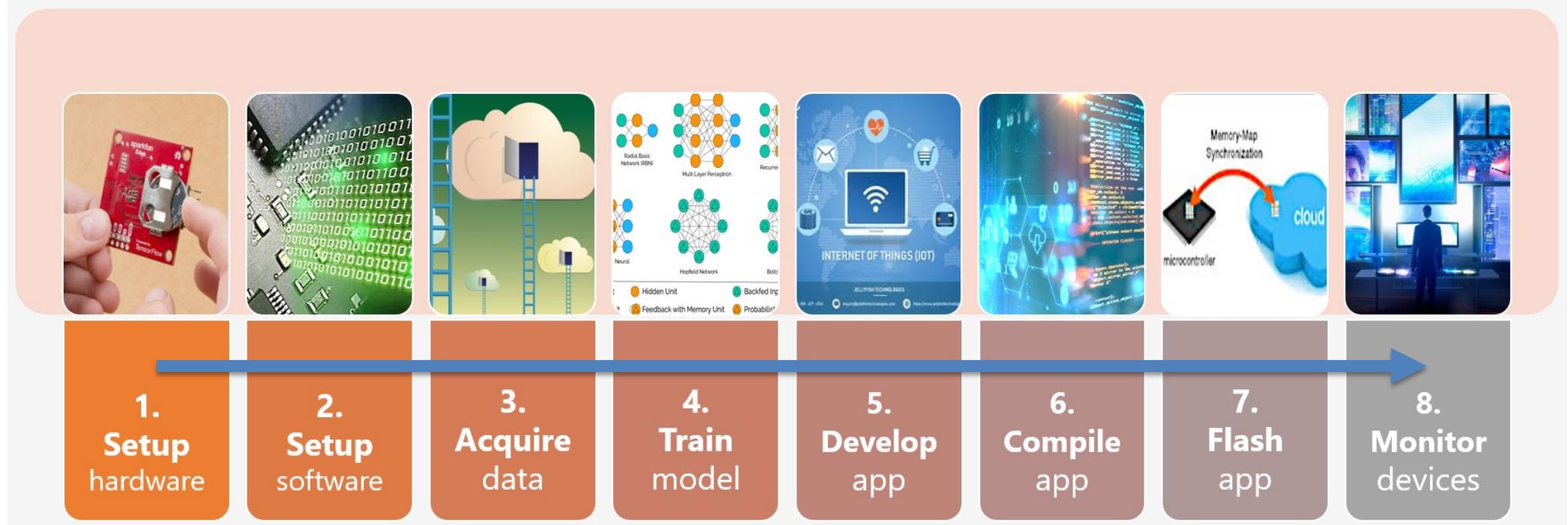
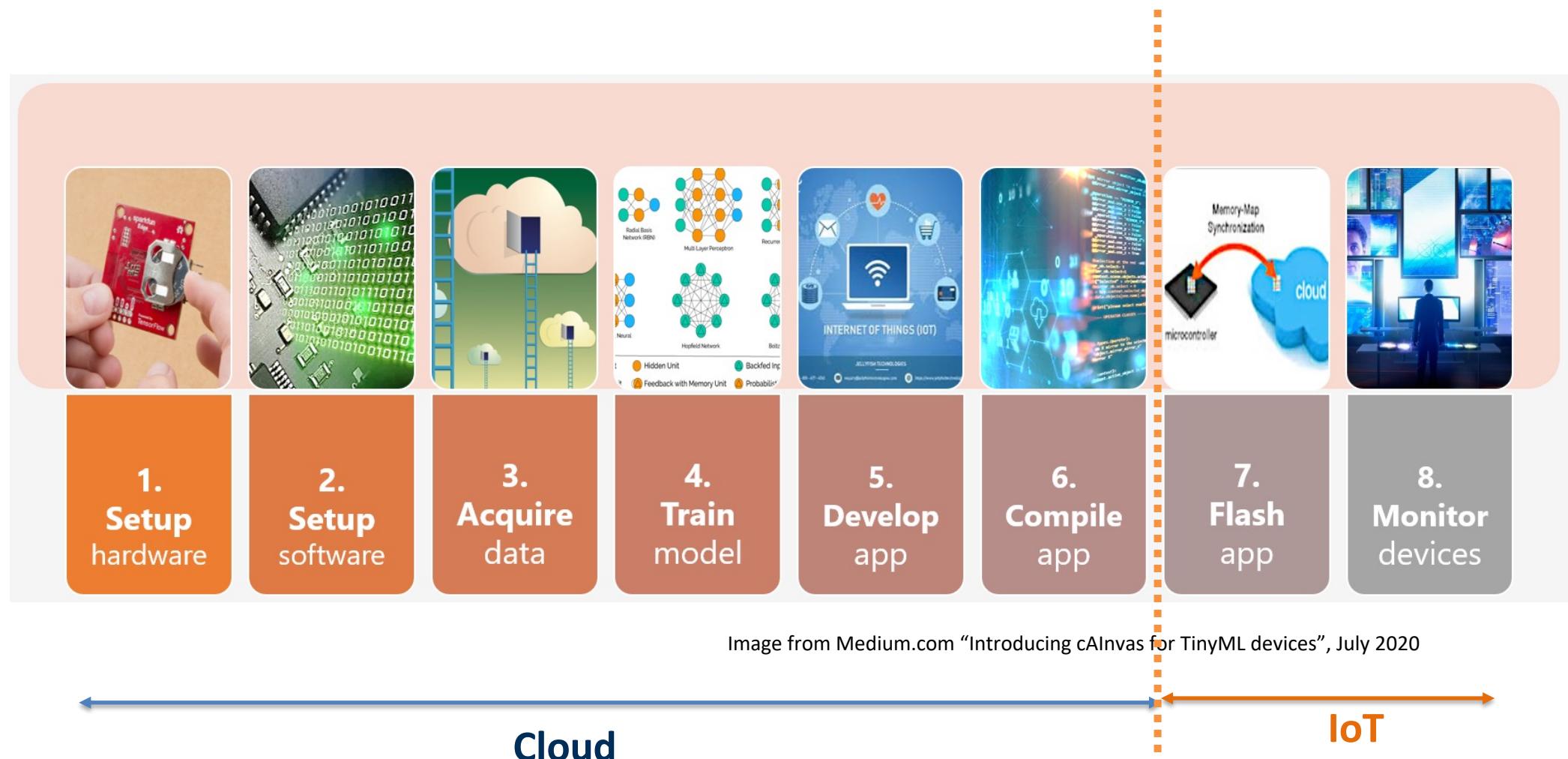


Image from Medium.com “Introducing cAlvas for TinyML devices”, July 2020

# How does TinyML work?



# Open points?

What about  
Deep Learning Models ?



A provocative question ...

Is Tiny Deep Learning the new Deep  
Learning?



# Deep Learning in 2020s



# The 70s hardware...



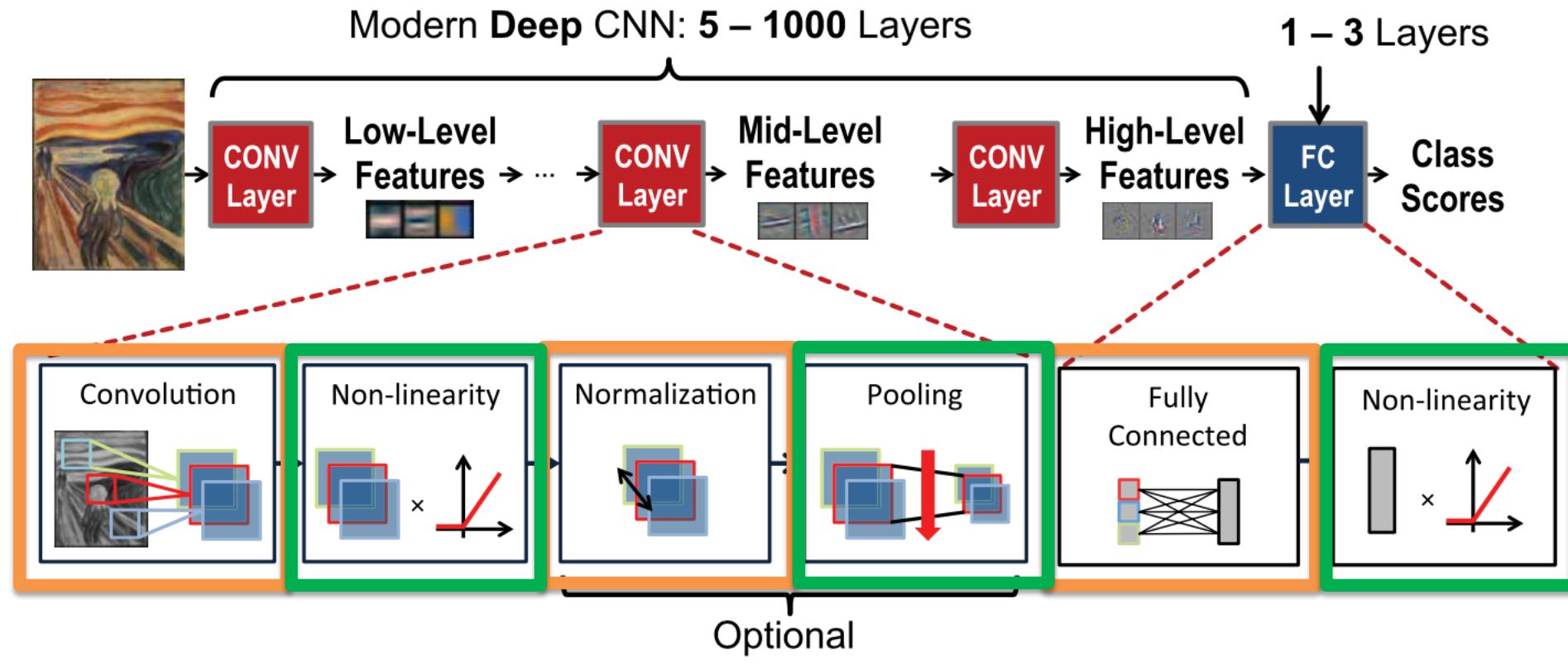
# What happened in the 2000s?



Which are the sources of complexity (memory footprint and computational demand) in deep learning models?



# Let's focus on Convolutional Neural Networks



Processing layer with trainable parameters

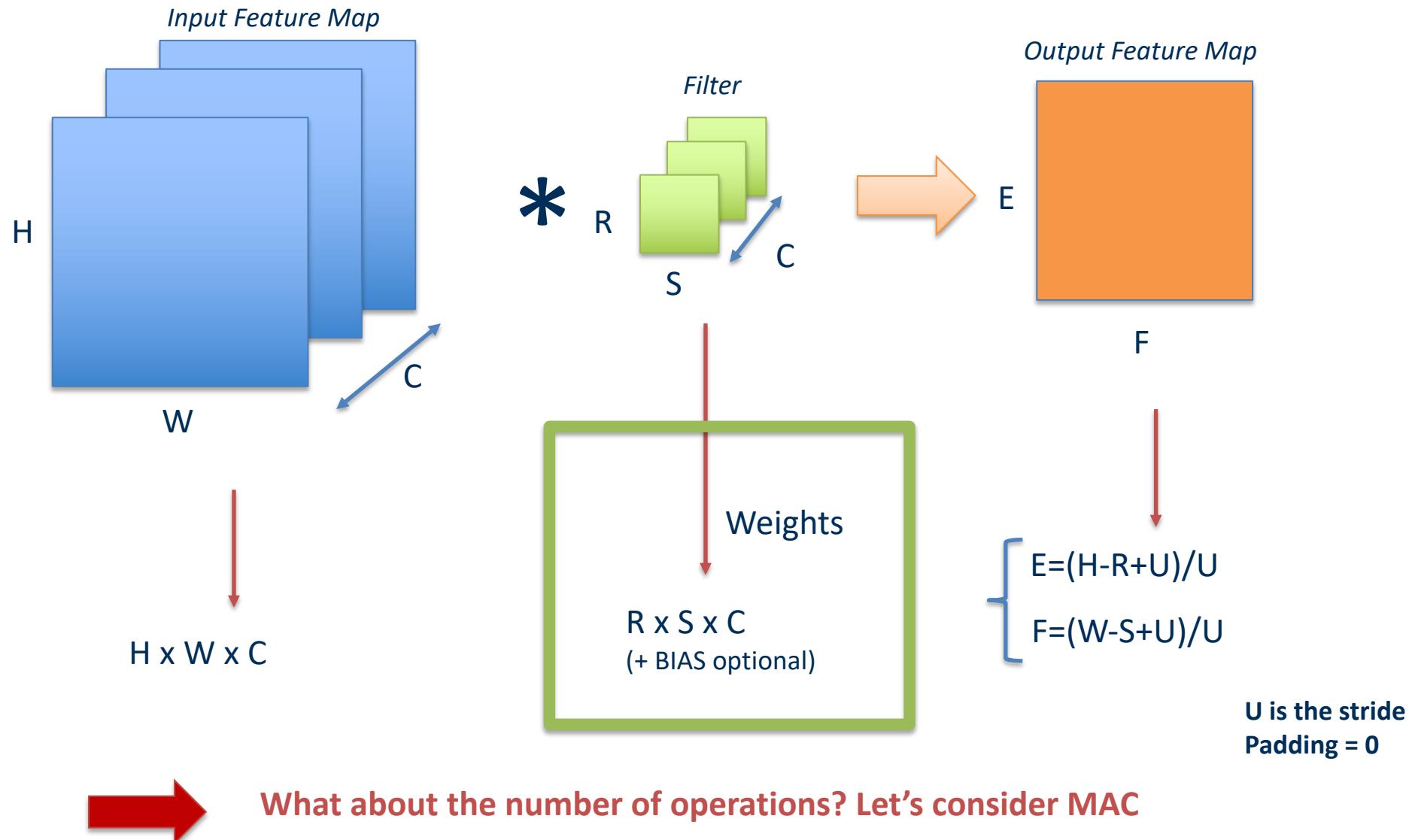
CONV and FC

Processing layer without trainable parameters

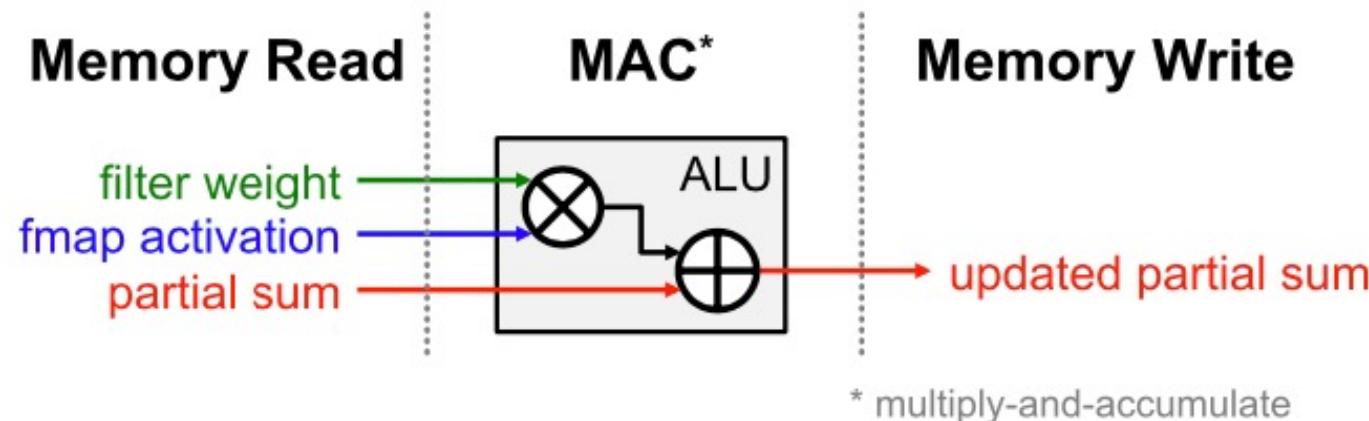
Picture taken from [1]



# Convolutional layers (single filter): weights and FMs



# What is a Multiply And Accumulate (MAC)?



Let's consider a  $[3 \times 3]$  filter W:

$$o = W[0] * I[0] + W[1] * I[1] + \dots + W[8] * I[8]$$

MAC

Total = 9 MAC

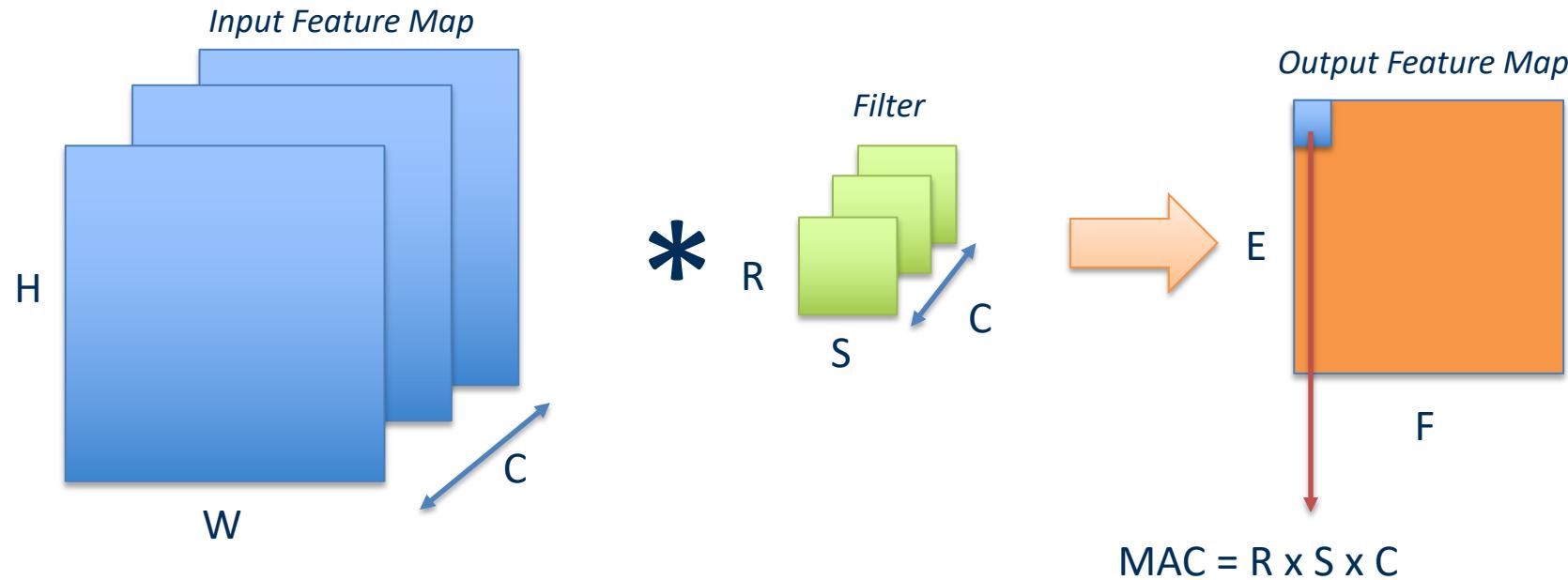


We have  $9+(9-1)$  Floating point OPerations (FLOPs) in total

(1MAC  $\approx$  2FLOPs)

Picture taken from [1]

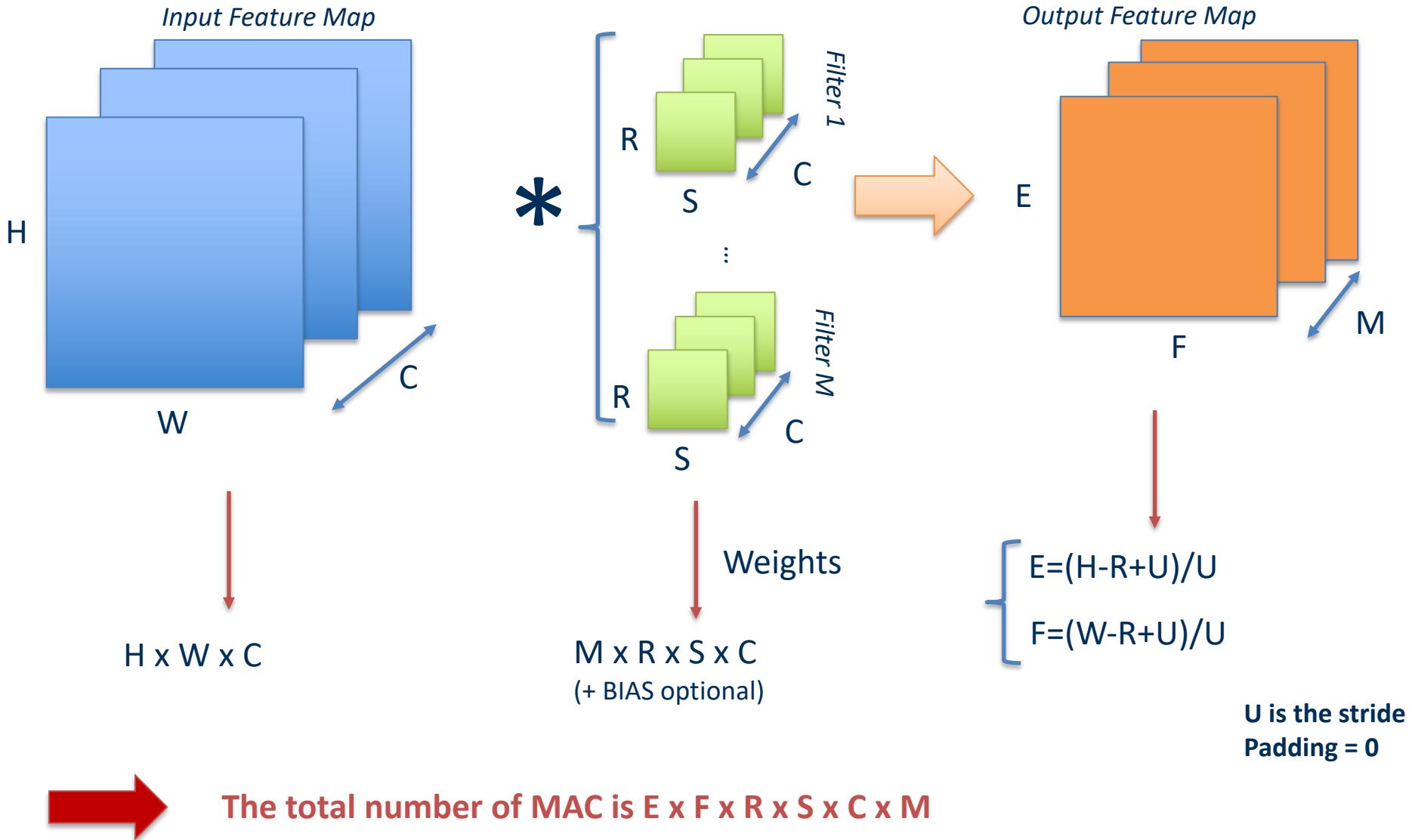
# Convolutional layers (single filter): MAC



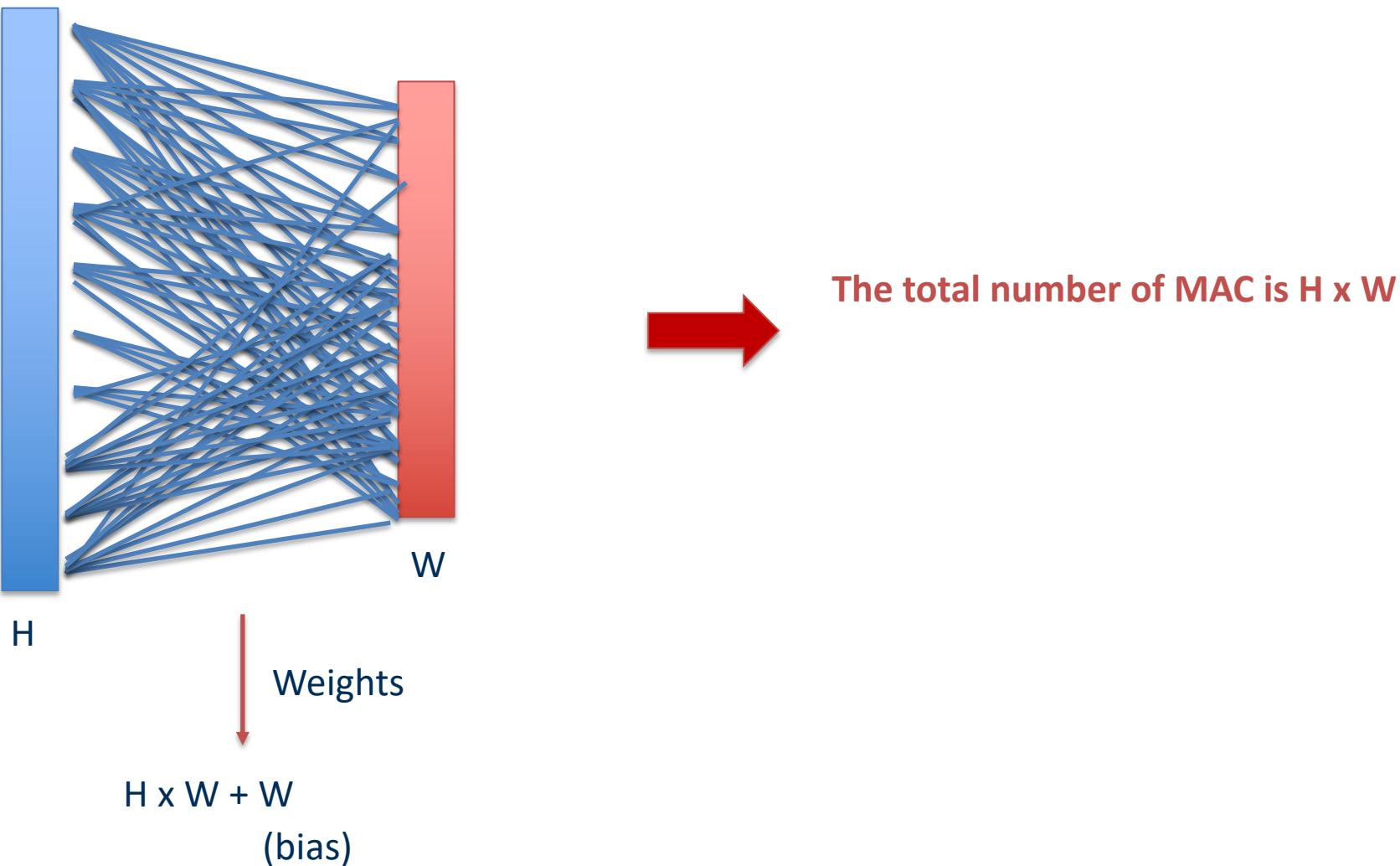
Total number of MAC operations:  $E \times F \times R \times S \times C$

We are here neglecting the bias

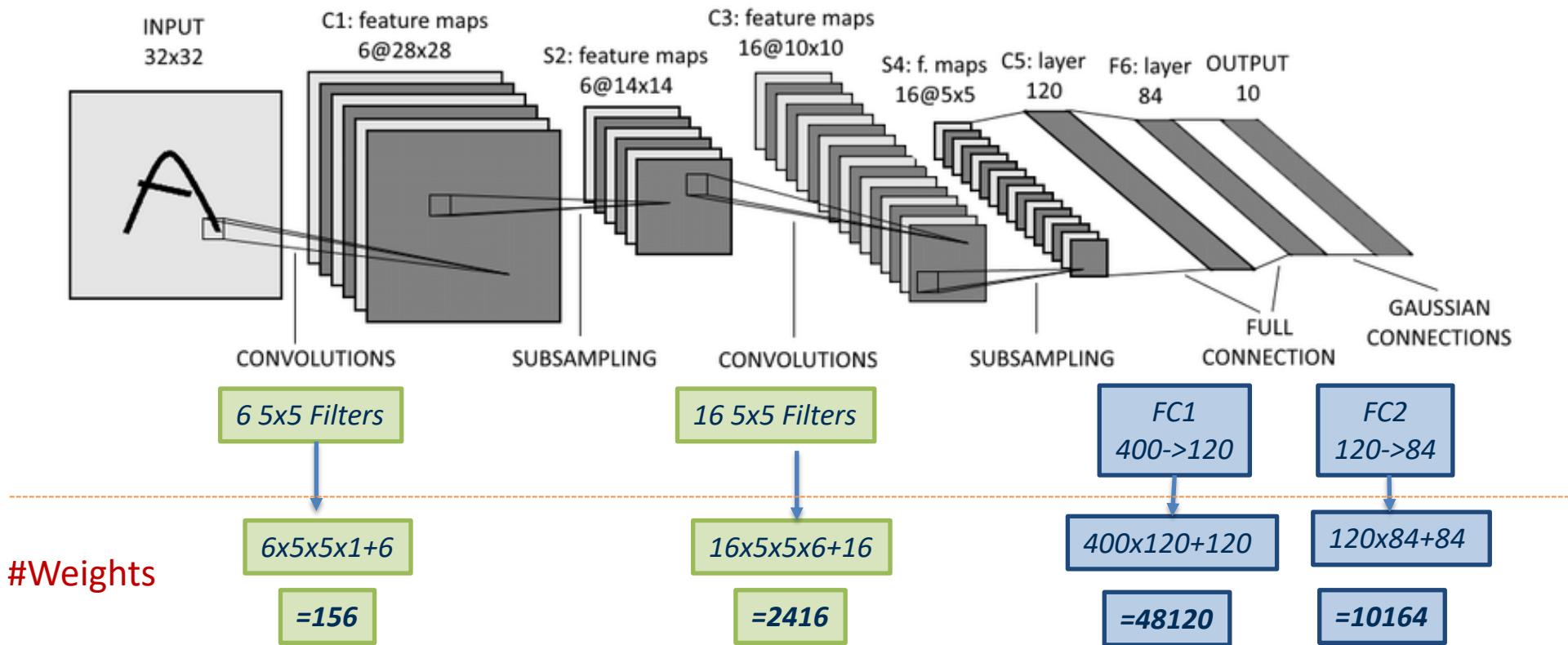
## Convolutional layers (multiple filters): weights, FPs, MAC



# Fully Connected Layer: weights and MAC



# Let's consider a simple CNN: the Le-Net (1989)



Picture taken from [2]

# Le-Net: the Tensorflow implementation and summary

```
model = models.Sequential()
model.add(layers.Conv2D(6, 5, activation='tanh',
input_shape=x_train.shape[1:]))
model.add(layers.AveragePooling2D(2))
model.add(layers.Activation('sigmoid'))
model.add(layers.Conv2D(16, 5, activation='tanh'))
model.add(layers.AveragePooling2D(2))
model.add(layers.Activation('sigmoid'))
model.add(layers.Conv2D(120, 5, activation='tanh'))
model.add(layers.Flatten())
model.add(layers.Dense(84, activation='tanh'))
model.add(layers.Dense(10, activation='softmax'))
model.summary()
```

Model: "sequential"

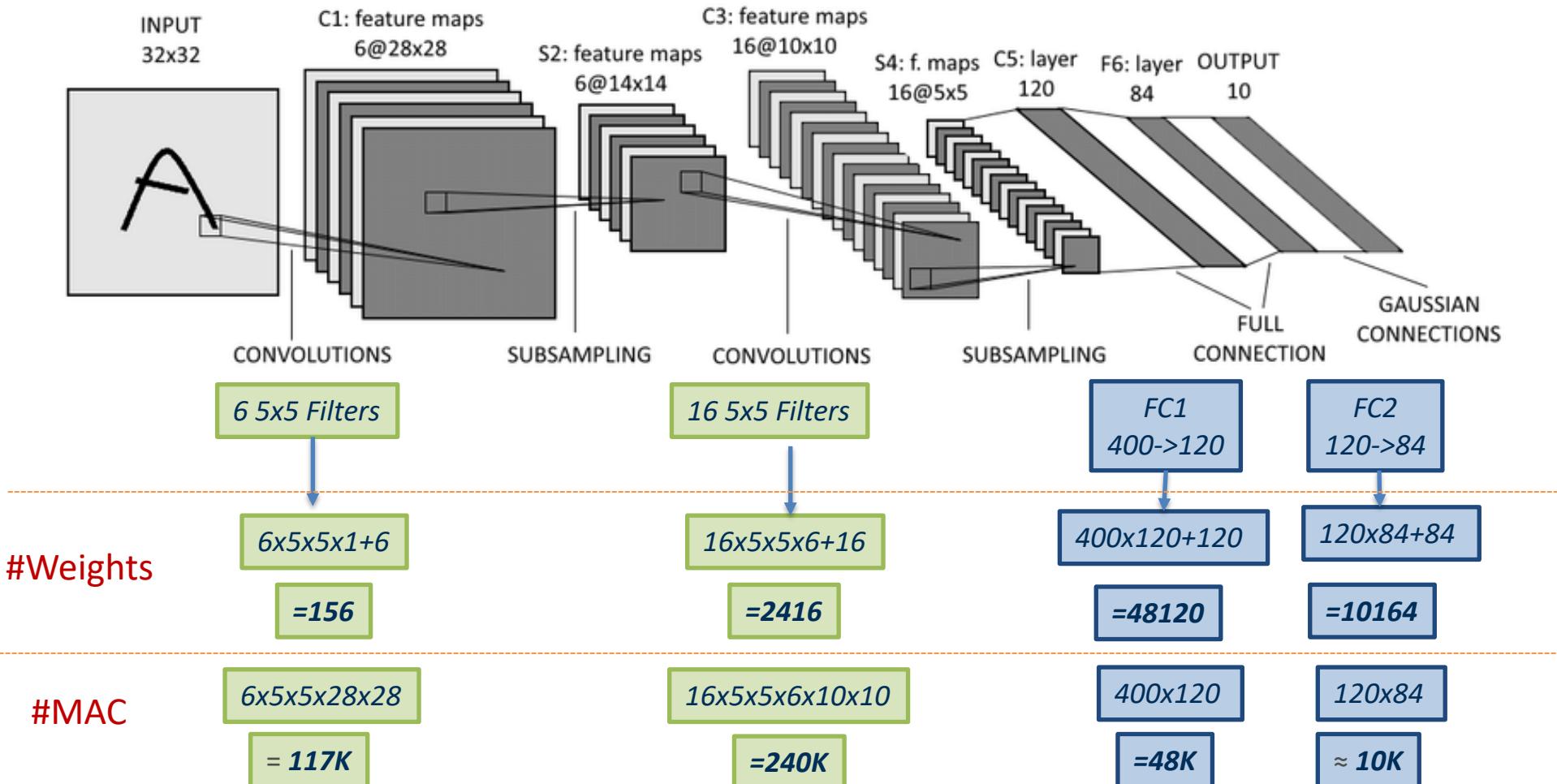
Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 28, 28, 6)	156
average_pooling2d (AveragePo	(None, 14, 14, 6)	0
activation (Activation)	(None, 14, 14, 6)	0
conv2d_1 (Conv2D)	(None, 10, 10, 16)	2416
average_pooling2d_1 (Average	(None, 5, 5, 16)	0
activation_1 (Activation)	(None, 5, 5, 16)	0
conv2d_2 (Conv2D)	(None, 1, 1, 120)	48120
flatten (Flatten)	(None, 120)	0
dense (Dense)	(None, 84)	10164
dense_1 (Dense)	(None, 10)	850

Total params: 61,706  
Trainable params: 61,706  
Non-trainable params: 0

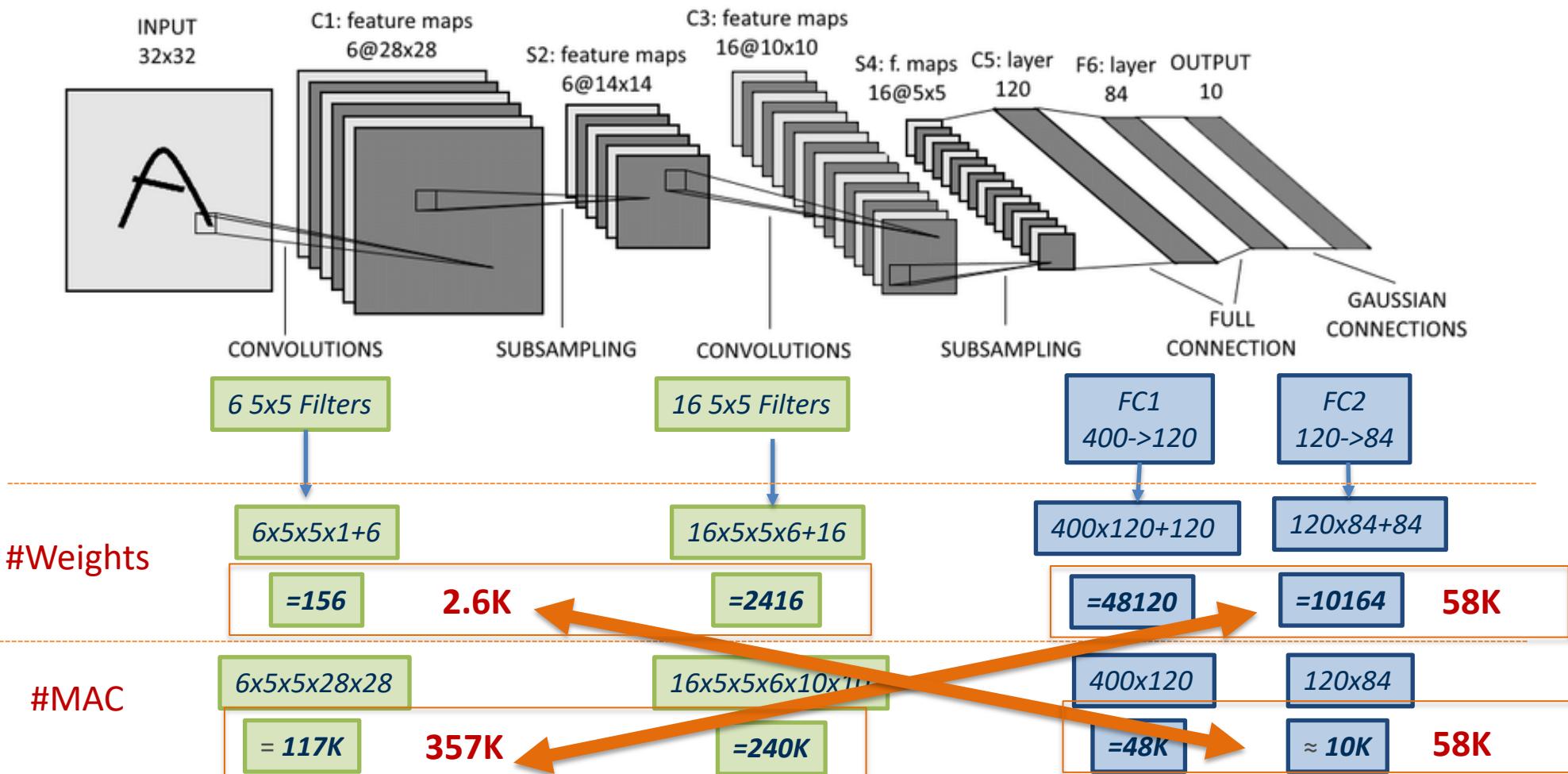
Picture taken from [2]



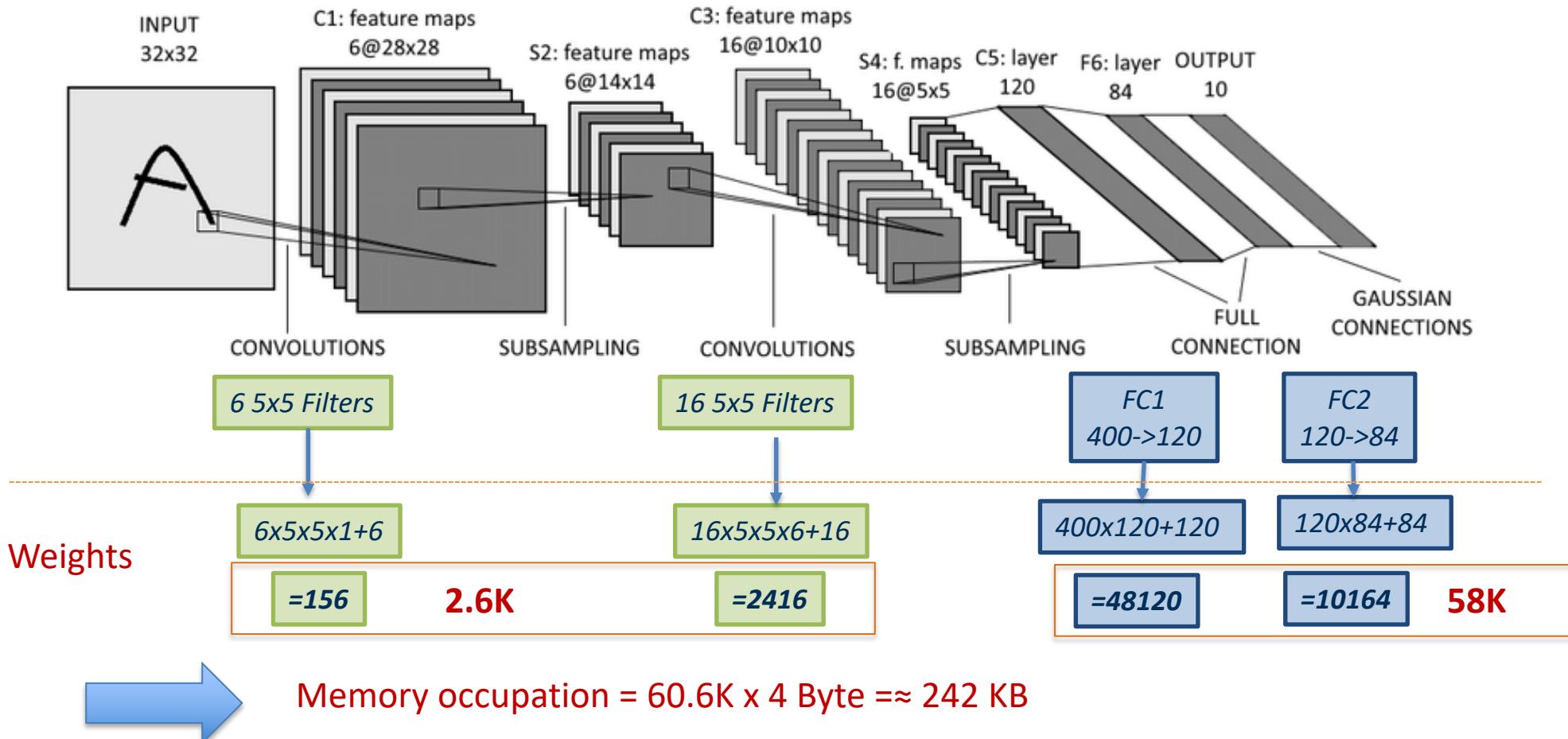
# Le-Net: weights and MAC



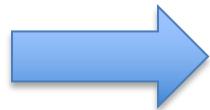
# Le-Net: weights and MAC



# Le-Net: memory occupation of the weights



# Can we port our LeNet on a IoT unit with 256 KB of RAM?

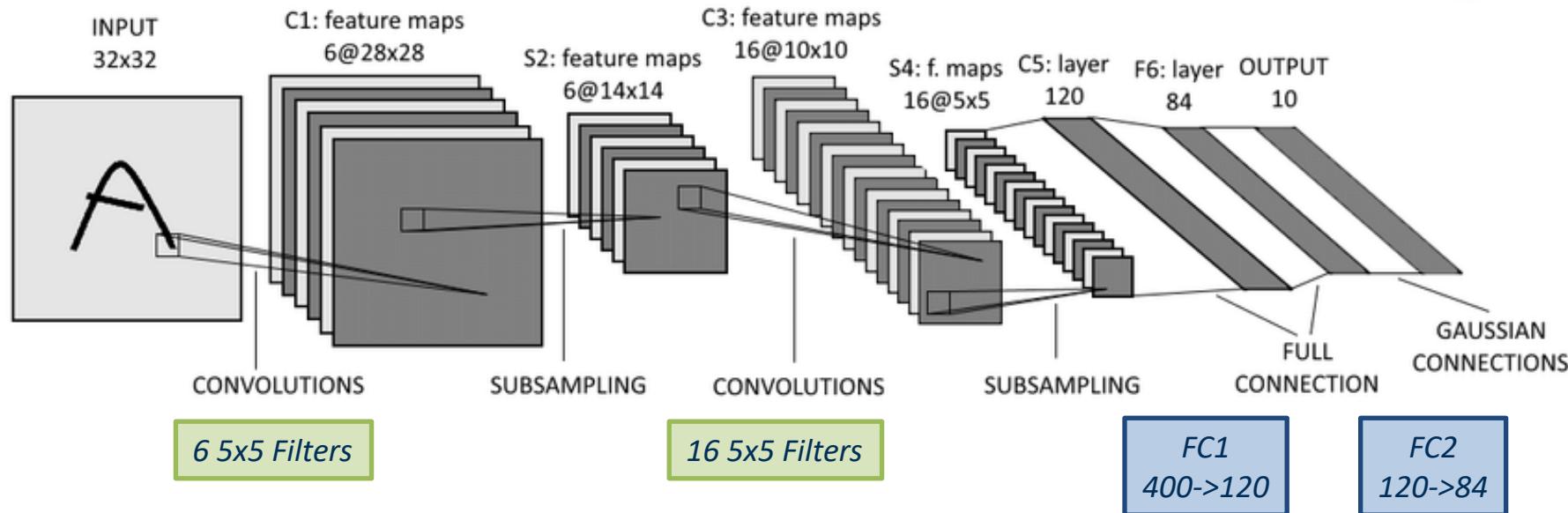


What are we missing?



# What are we missing?

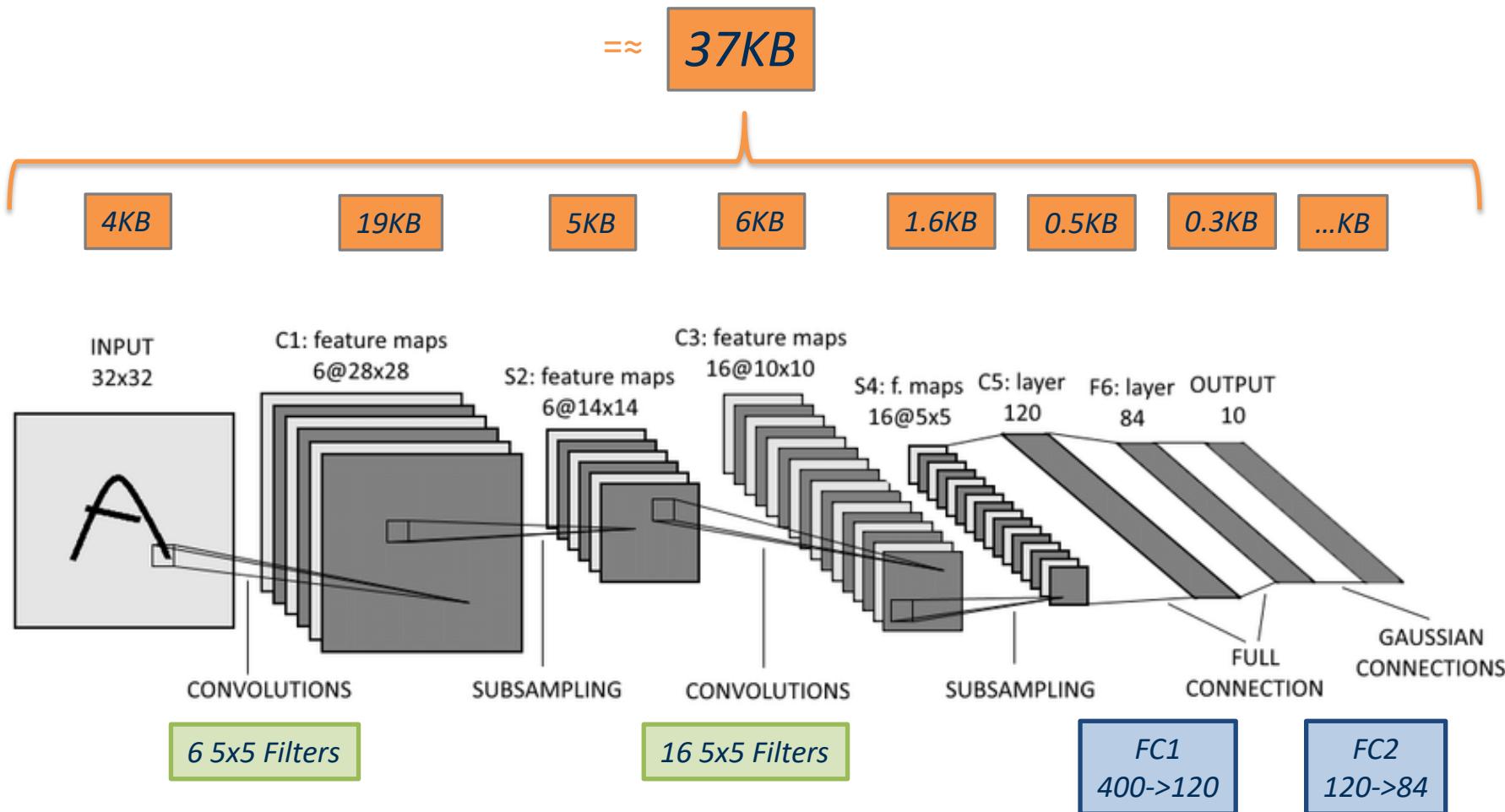
What about the feature maps ???



Weight Memory occupation =  $60.6K \times 4 \text{ Byte} = 242 \text{ KB}$

Picture taken from [2]

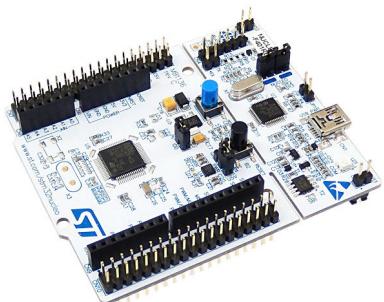
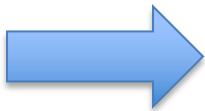
# Le-Net: memory occupation of the feature maps



**Weight Memory occupation =  $60.6K \times 4 \text{ Byte} = 240 \text{ KB}$**

Picture taken from [2]

# Can we port our LeNet on a IoT unit with 384 KB of RAM?



Product lines	$F_{CPU}$ (MHz)	Flash (Kbytes)	RAM (KB)
STM32F469 <sup>2</sup>	180	512 K to 2056 K	384



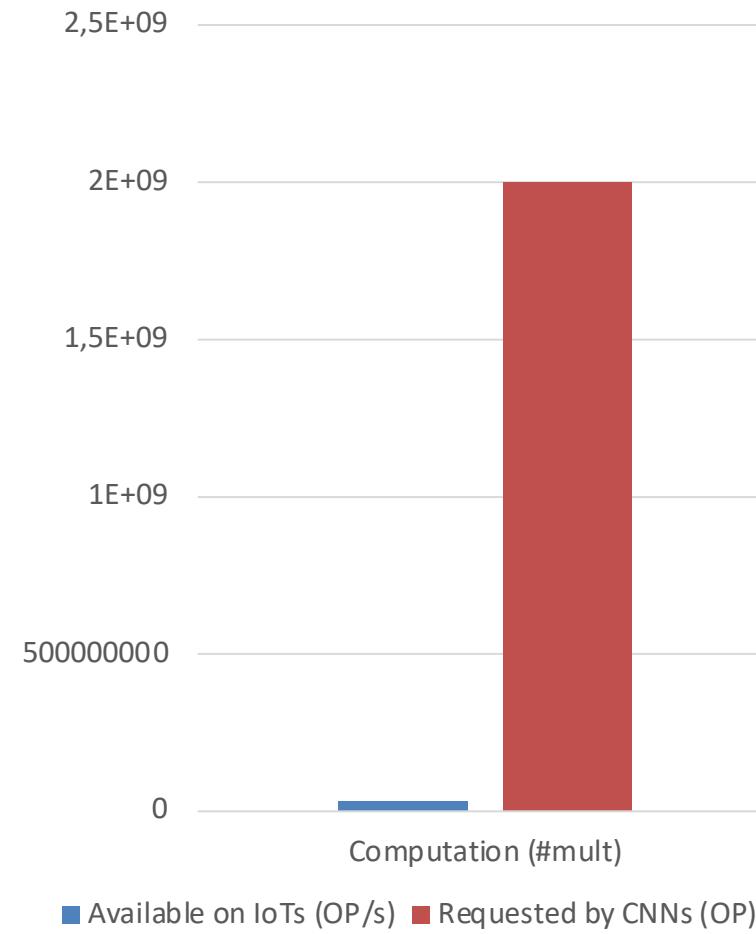
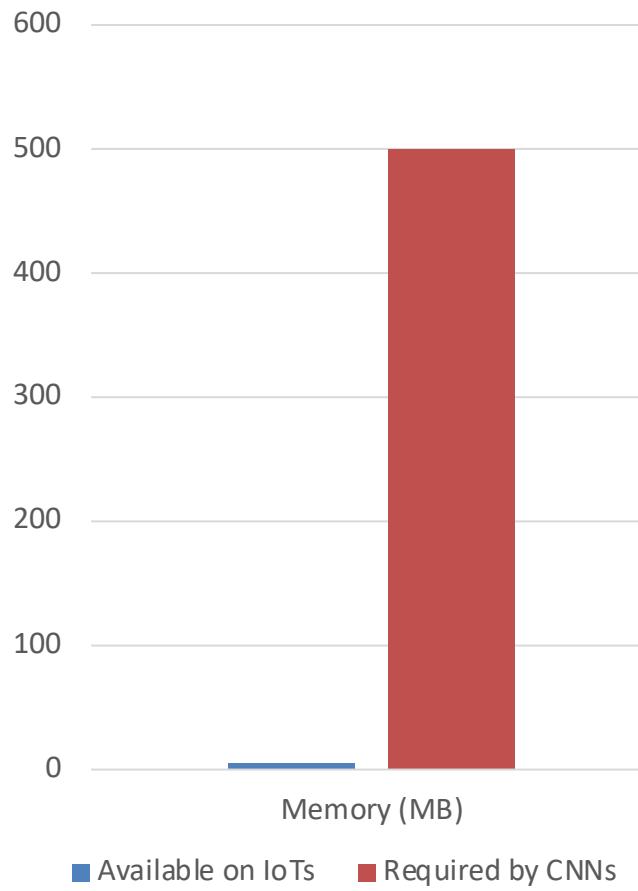
# CNNs memory and computational demand

Metrics	AlexNet	Overfeat fast	VGG 16	GoogLeNet v1	ResNet 50
<b>Top-5 error<sup>†</sup></b>	16.4	14.2	7.4	6.7	5.3
<b>Top-5 error (single crop)<sup>†</sup></b>	19.8	17.0	8.8	10.7	7.0
<b>Input Size</b>	227×227	231×231	224×224	224×224	224×224
<b># of CONV Layers</b>	5	5	13	57	53
<b>Depth in # of CONV Layers</b>	5	5	13	21	49
<b>Filter Sizes</b>	3,5,11	3,5,11	3	1,3,5,7	1,3,7
<b># of Channels</b>	3-256	3-1024	3-512	3-832	3-2048
<b># of Filters</b>	96-384	96-1024	64-512	16-384	64-2048
<b>Stride</b>	1,4	1,4	1	1,2	1,2
<b>Weights</b>	2.3M	16M	14.7M	6.0M	23.5M
<b>MACs</b>	666M	2.67G	15.3G	1.43G	3.86G
<b># of FC Layers</b>	3	3	3	1	1
<b>Filter Sizes</b>	1,6	1,6,12	1,7	1	1
<b># of Channels</b>	256-4096	1024-4096	512-4096	1024	2048
<b># of Filters</b>	1000-4096	1000-4096	1000-4096	1000	1000
<b>Weights</b>	58.6M	130M	124M	1M	2M
<b>MACs</b>	58.6M	130M	124M	1M	2M
<b>Total Weights</b>	61M	146M	138M	7M	25.5M
<b>Total MACs</b>	724M	2.8G	15.5G	1.43G	3.9G
<b>Pretrained Model Website</b>	[57, 58]	n/a	[57-59]	[57-59]	[57-59]

<sup>†</sup>Accuracy is Measured Based on Top-5 Error on ImageNet [14].

Picture taken from [1]

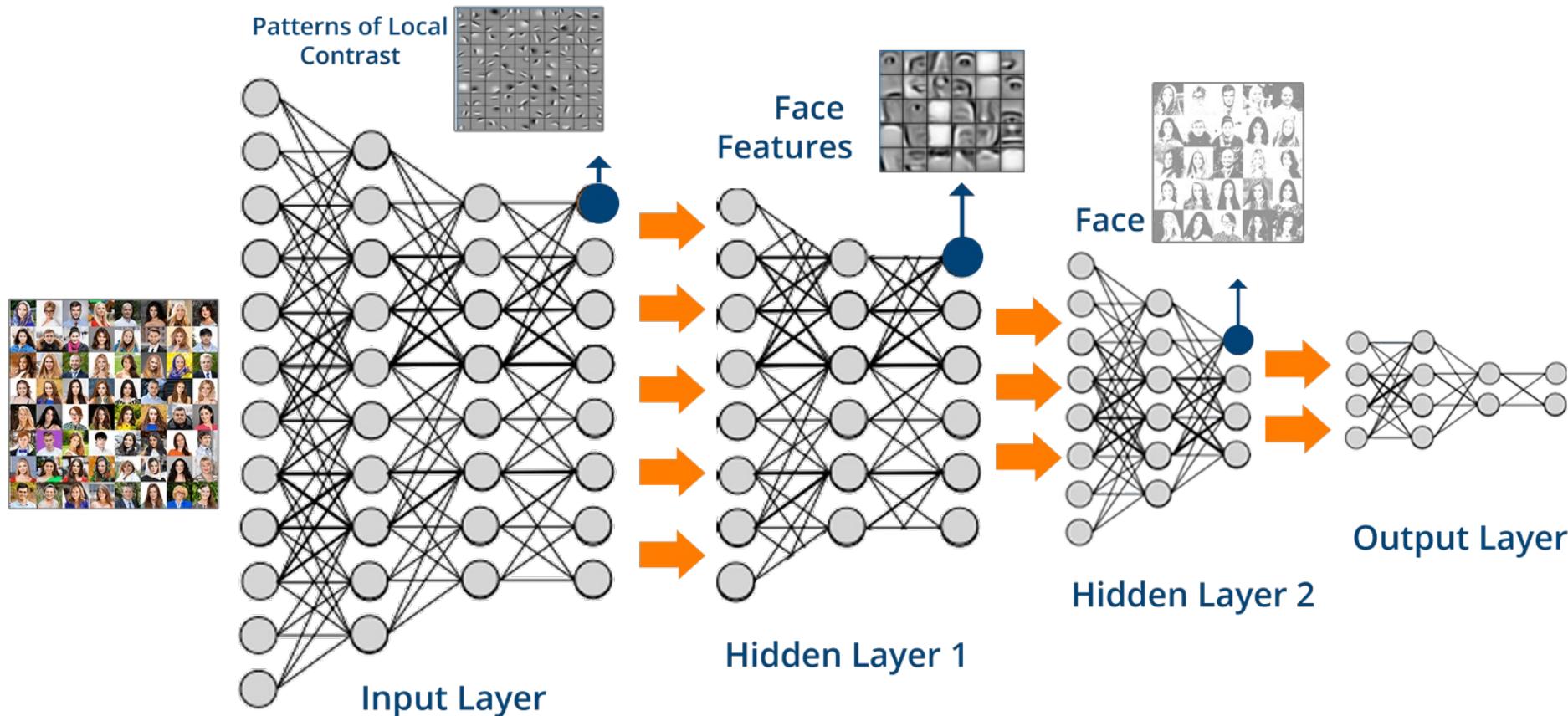
# What about available memory and computation in IoT units?





How can we match these extremes?

# Approximate Deep Learning for IoT: Tiny Deep Learning



# Approximate Deep Learning for IoT: Tiny Deep Learning

