

Convolutional Neural Networks

Neural Networks Design And Application

LeNet-5 in 1999

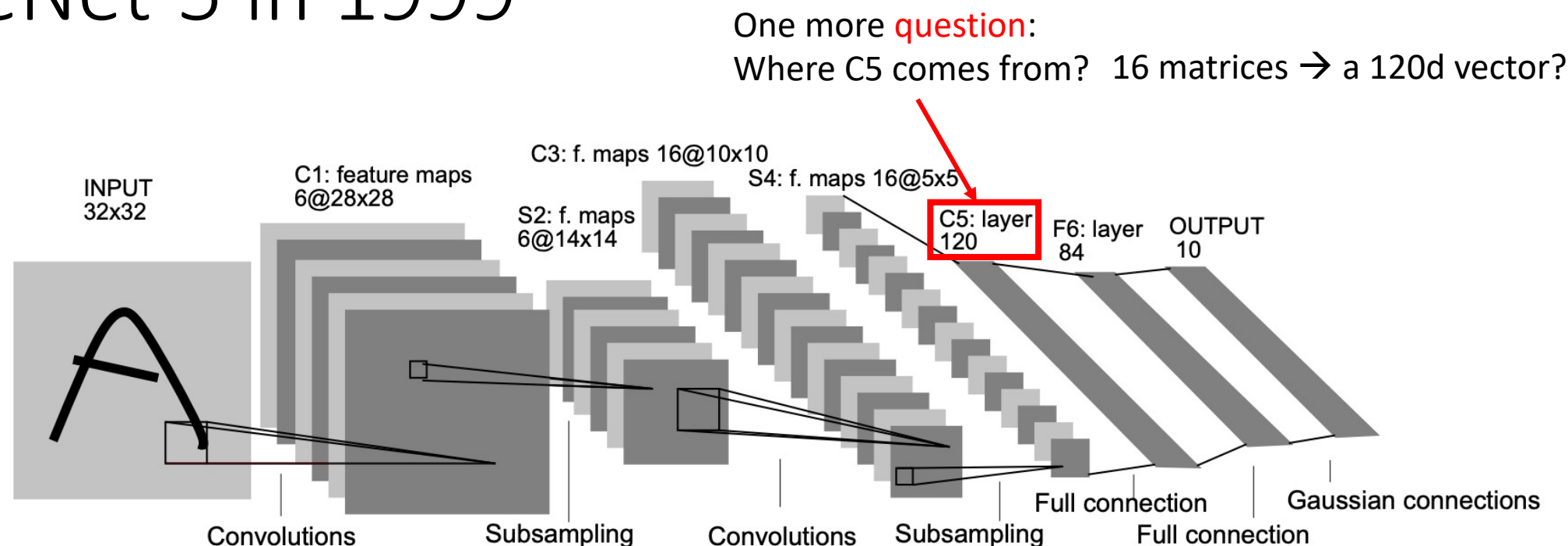
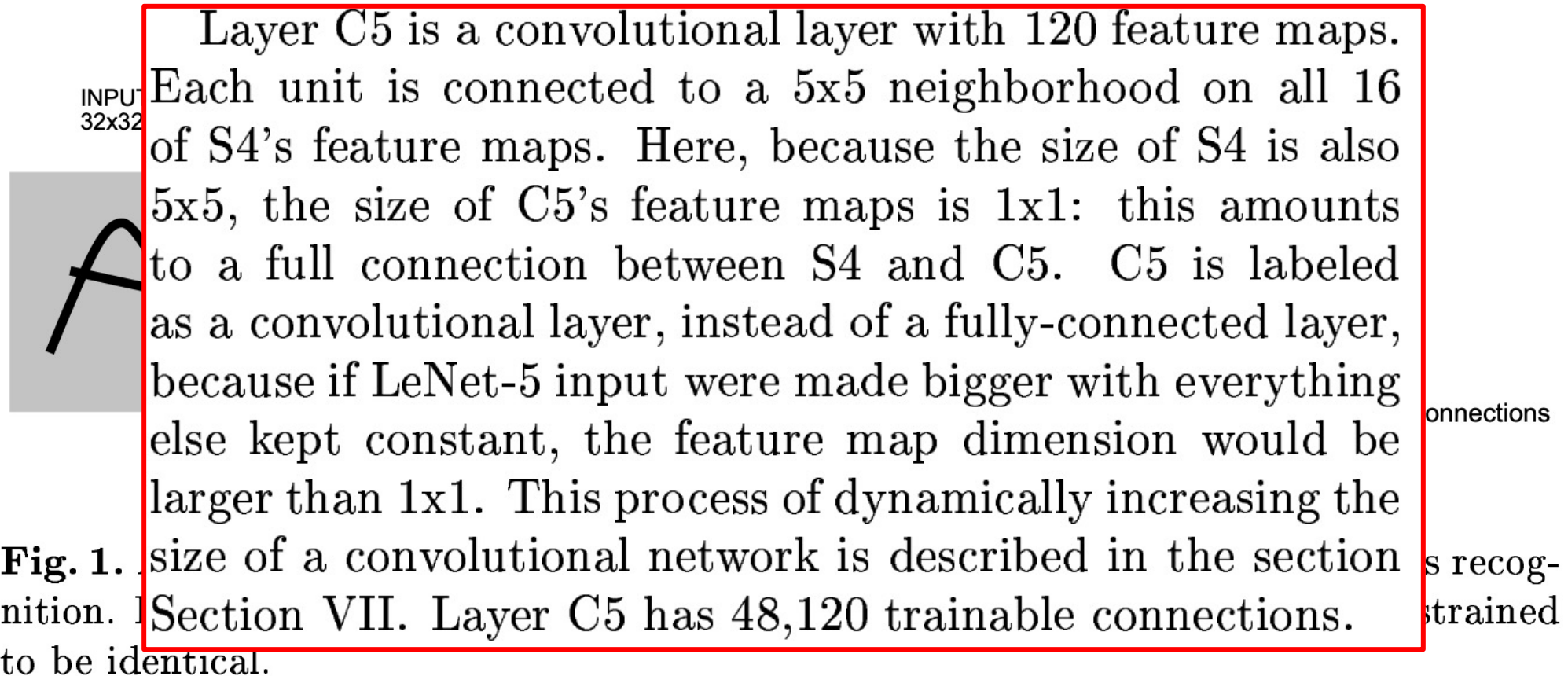


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeNet-5 in 1999

One more **question**:

Where C5 comes from? 16 matrices \rightarrow a 120d vector?

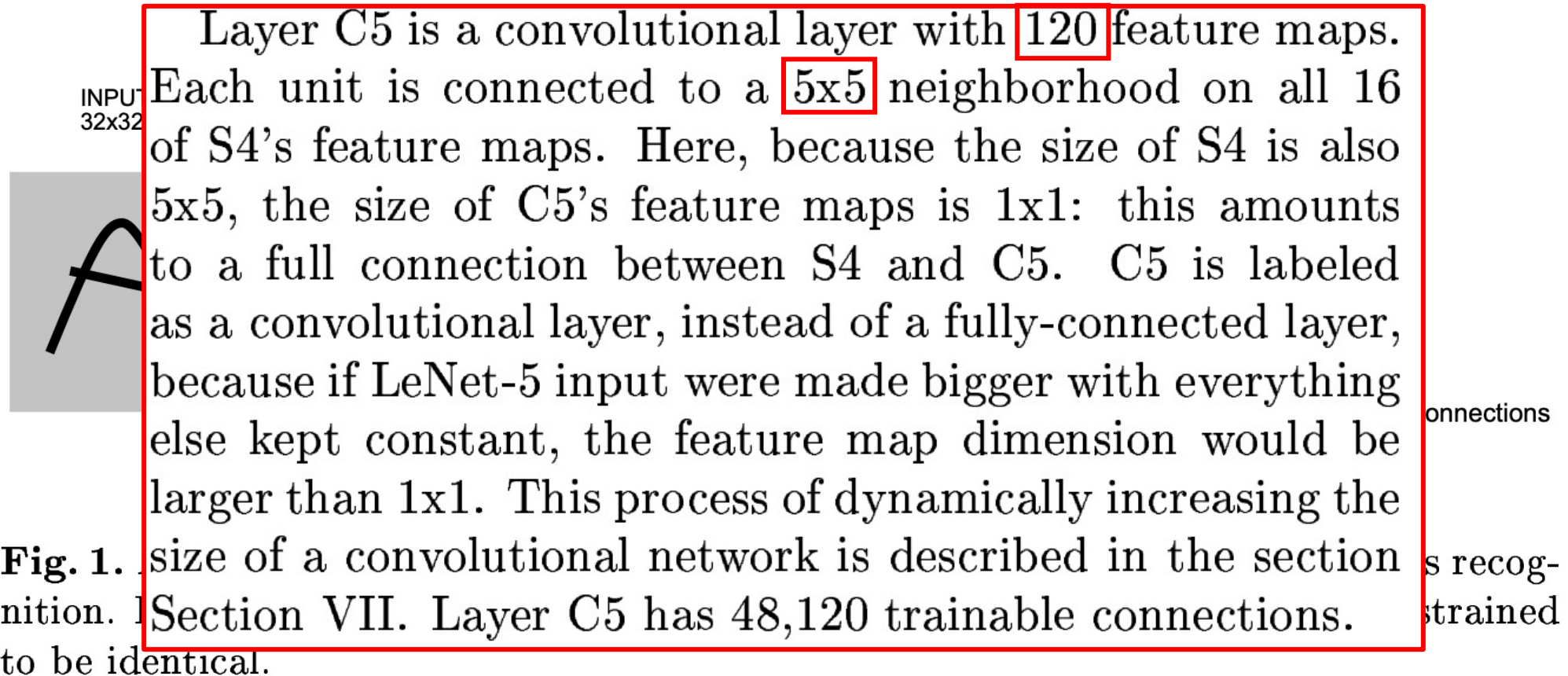


LeCun, Yann, Patrick Haffner, Léon Bottou, and Yoshua Bengio. "Object recognition with gradient-based learning." In *Shape, contour and grouping in computer vision*, pp. 319-345. Springer, Berlin, Heidelberg, 1999.

LeNet-5 in 1999

One more **question**:

Where C5 comes from? 16 matrices \rightarrow a 120d vector?



Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

*

1	0	1
0	1	0
1	0	1

$m=3$

One matrix

→

4	3	4
2	4	3
2	3	4

$n-m+1=3$

One matrix

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=3$

One matrix

→

4	3	4
2	4	3
2	3	4

$n-m+1=3$

One matrix

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

→

$m=5$

One matrix

4	3	4
2	4	3
2	3	4

$n-m+1=3$

One matrix

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

4	3	4
2	4	3
2	3	4

$n-m+1=1$

One matrix

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

Step 1: finding pairs

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

???

$n-m+1=1$

One matrix

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

???

$n-m+1=1$

One matrix

Step 1: finding pairs

1 pair

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

???

$n-m+1=1$

One matrix

Step 1: finding pairs

1 pair

Step 2: elementwise
summation

One input matrix * one filter → one feature matrix

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

???

$n-m+1=1$

One matrix

One input matrix * one filter → one feature matrix

Step 1: finding pairs

1 pair

Step 2: elementwise
summation

Convolution for images (matrices)

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

One matrix

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

*

$m=5$

One matrix

→

13

$n-m+1=1$

One matrix

Step 1: finding pairs

1 pair

Step 2: elementwise summation

One input matrix * one filter → one feature matrix

Operations with convolution layers

- Padding
- Pooling layers for arbitrary input resolution

Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps

Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

$n=5$

*

1	0	1
0	1	0
1	0	1

$m=3$

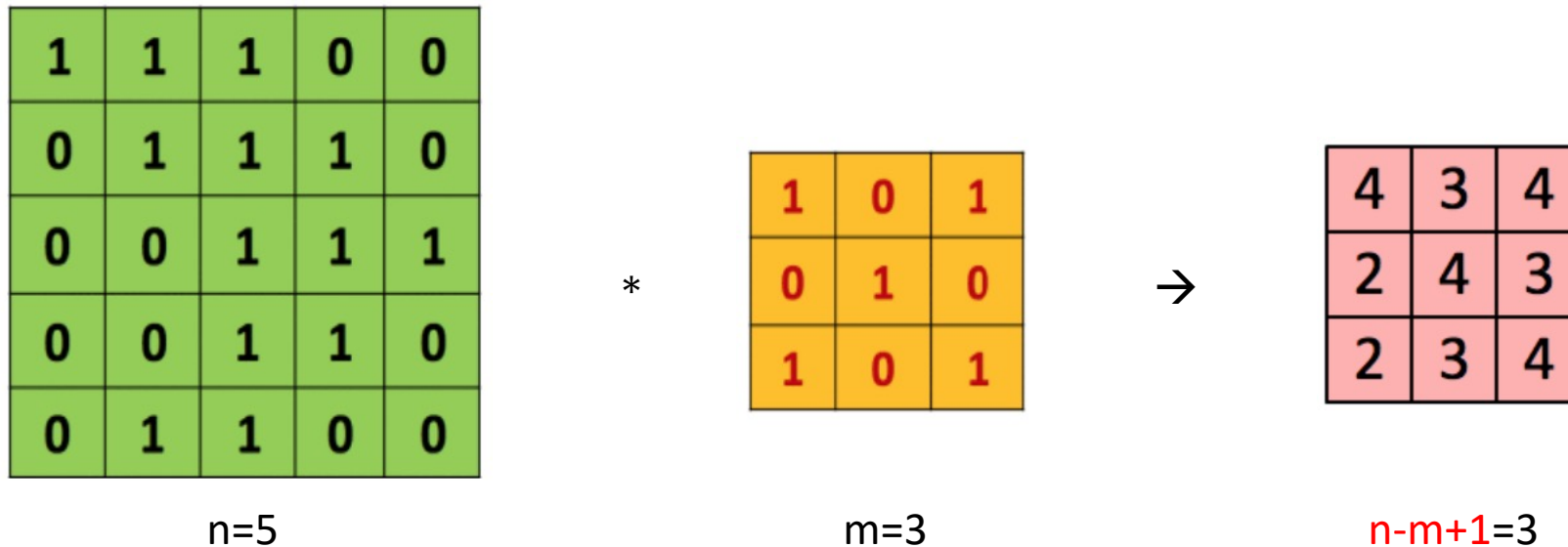
→

4	3	4
2	4	3
2	3	4

$n-m+1=3$

Operations with convolution layers

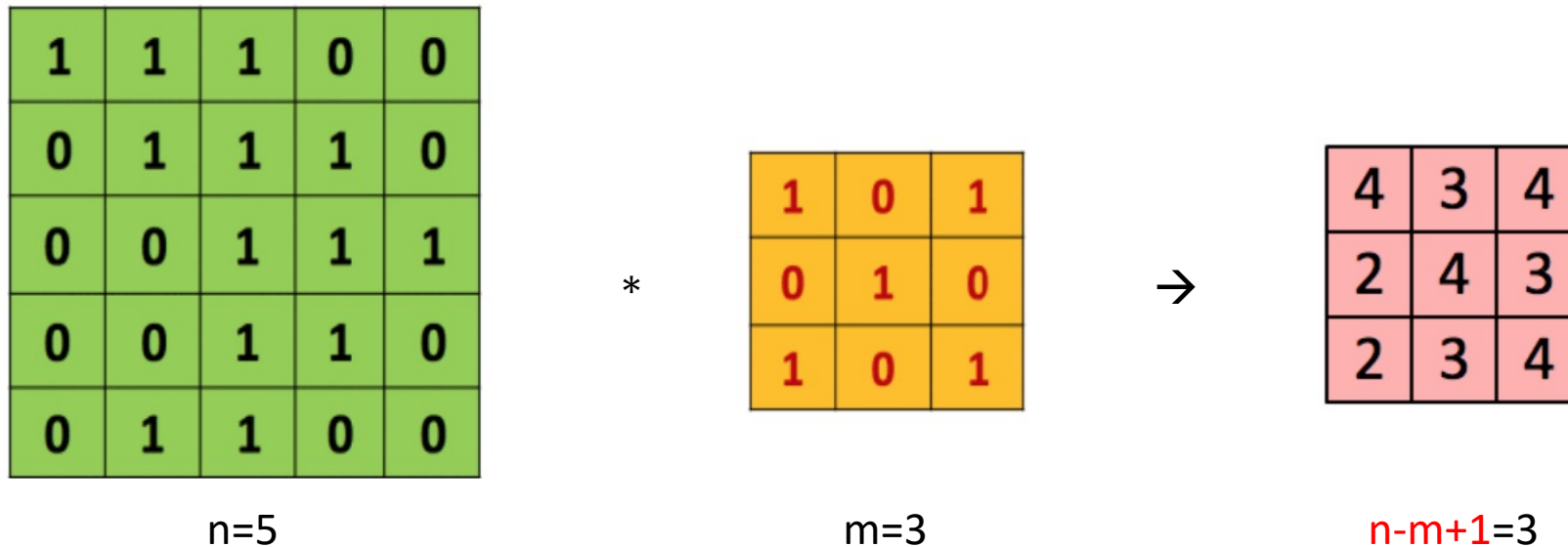
- Padding: convolution operation reduces the size of feature maps



If $m > 1 \rightarrow ??$

Operations with convolution layers

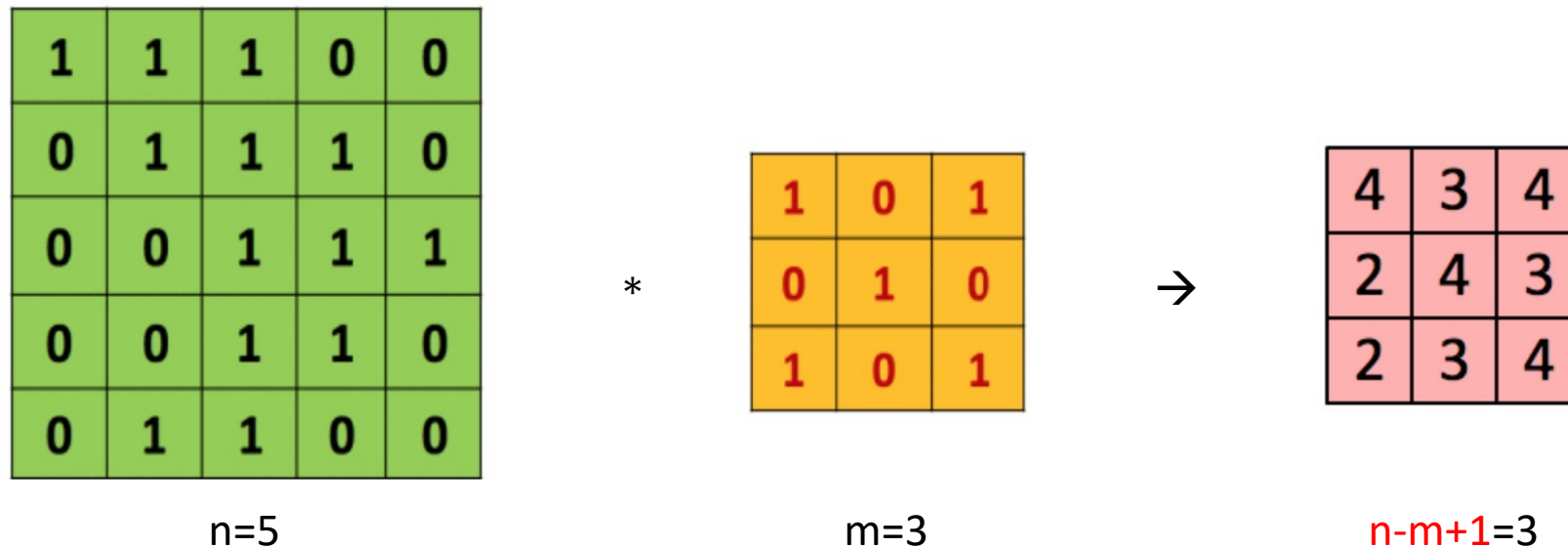
- Padding: convolution operation reduces the size of feature maps



If $m > 1 \rightarrow$ convolution will reduce the dimension

Operations with convolution layers

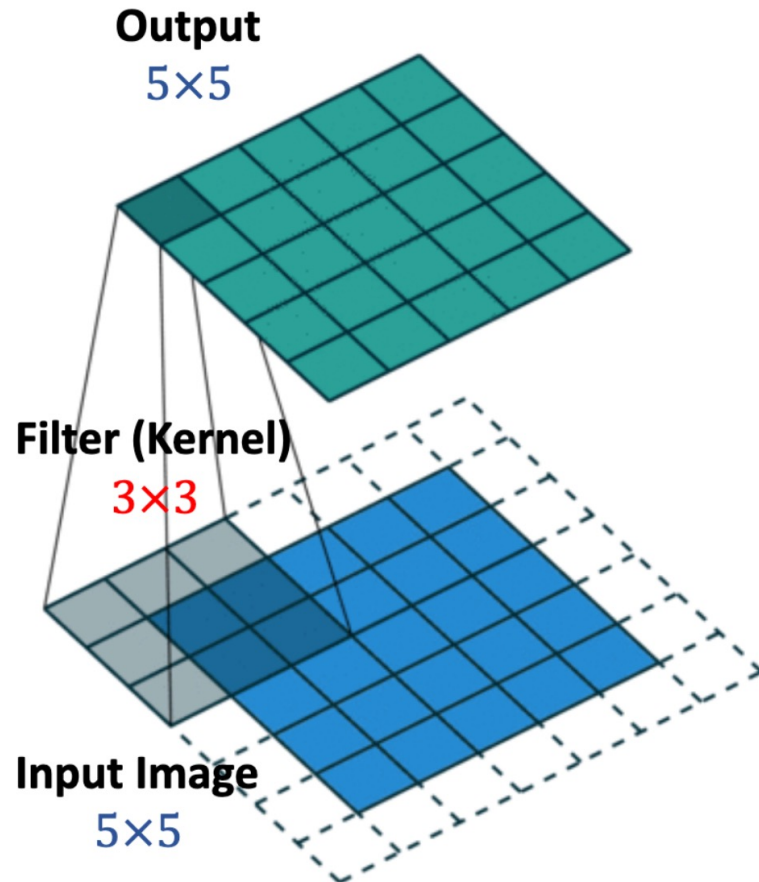
- Padding: convolution operation reduces the size of feature maps



If $m > 1 \rightarrow$ convolution will reduce the dimension
The input resolution introduces a limits of #convolution layers

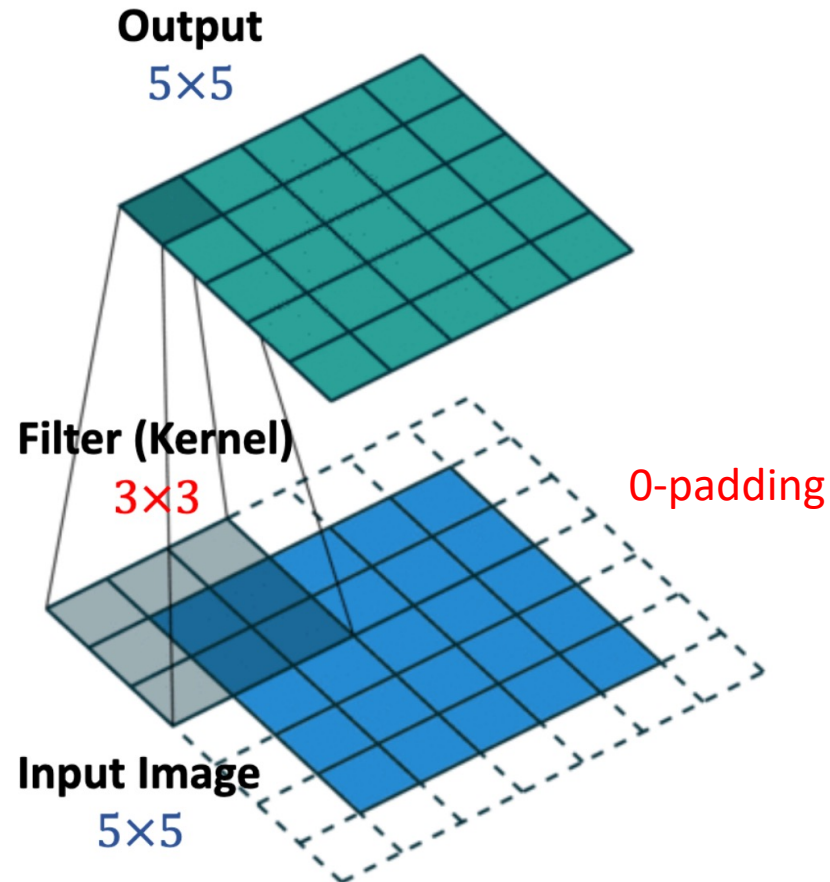
Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps



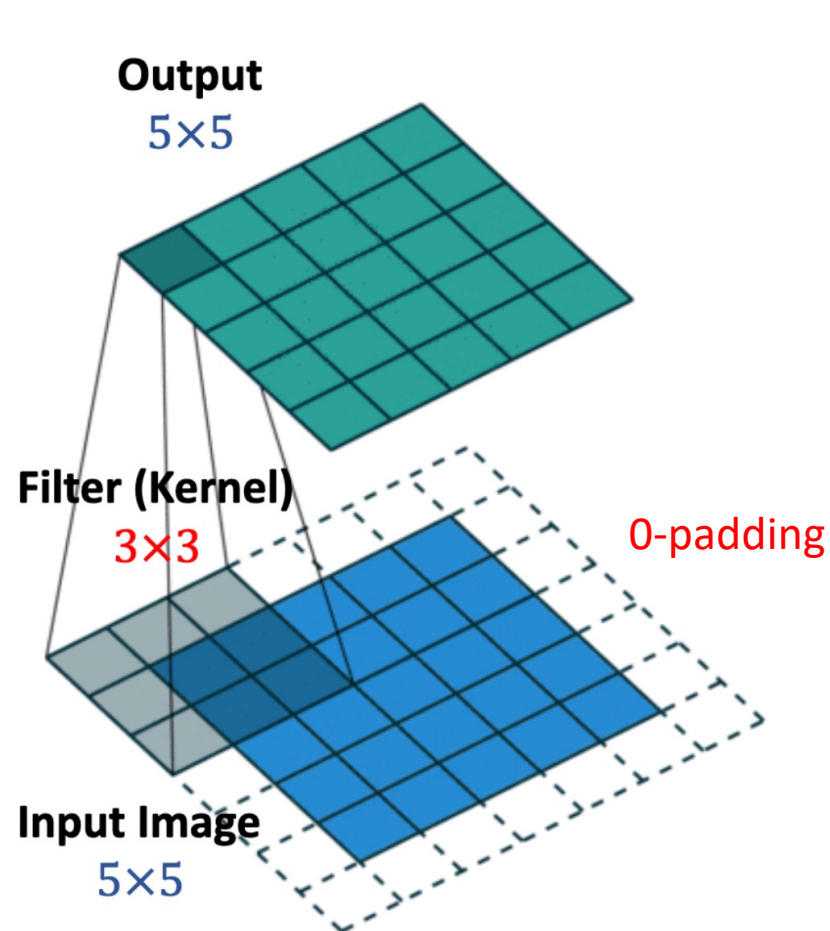
Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps



Operations with convolution layers

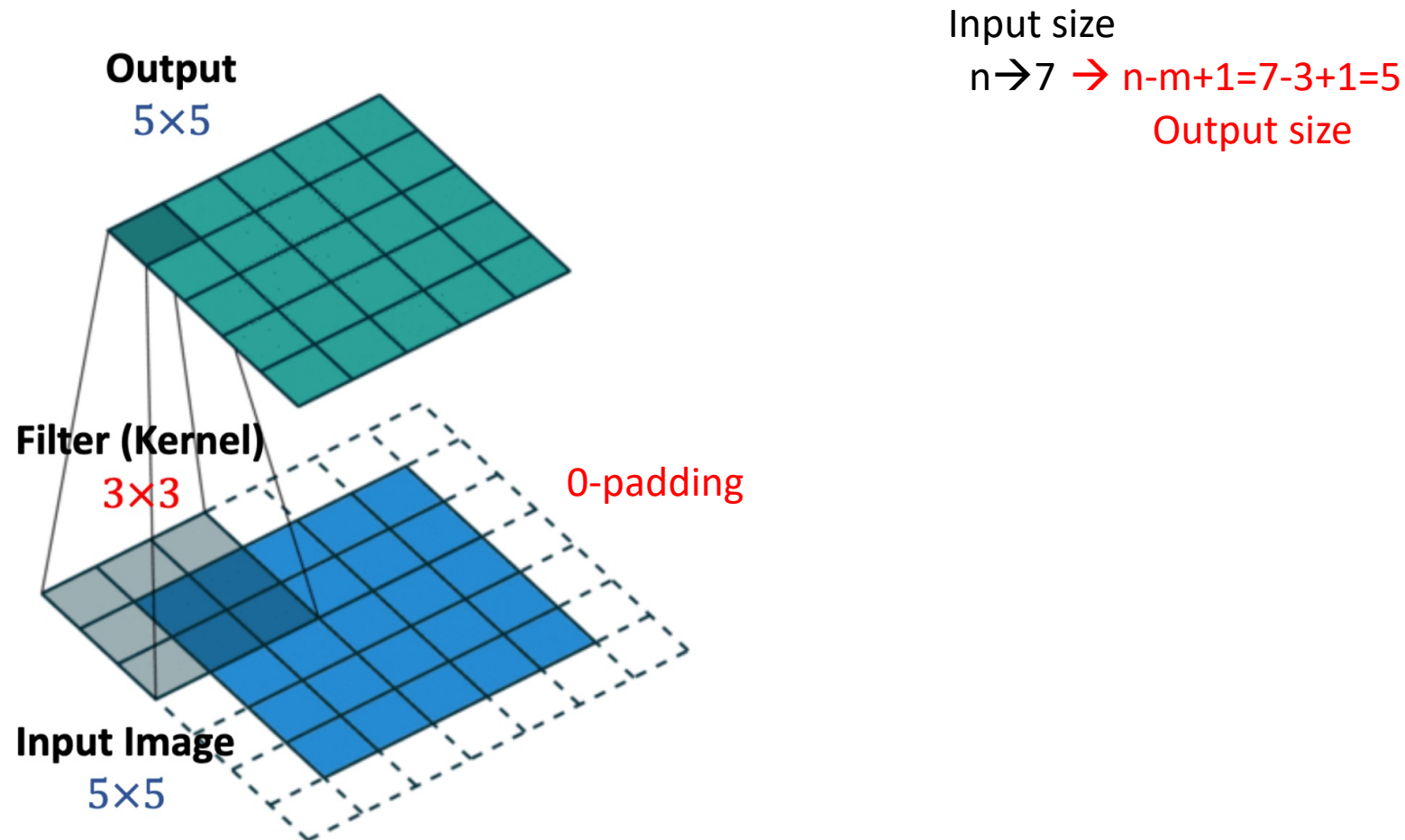
- Padding: convolution operation reduces the size of feature maps



Input size
 $n \rightarrow 7$

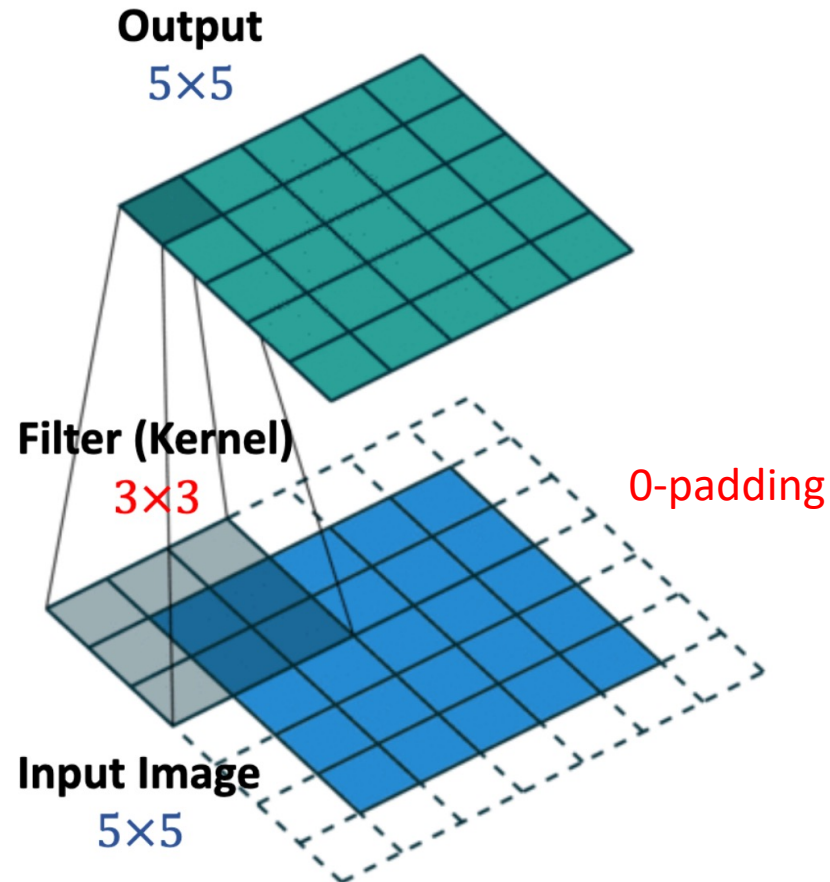
Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps



Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps



$$n \rightarrow 7 \rightarrow n - m + 1 = 7 - 3 + 1 = 5$$

Conclusion:
dimension of feature maps remains the same

Operations with convolution layers

- Padding: convolution operation reduces the size of feature maps
- Pooling layers for an arbitrary input resolution

Input resolution issue

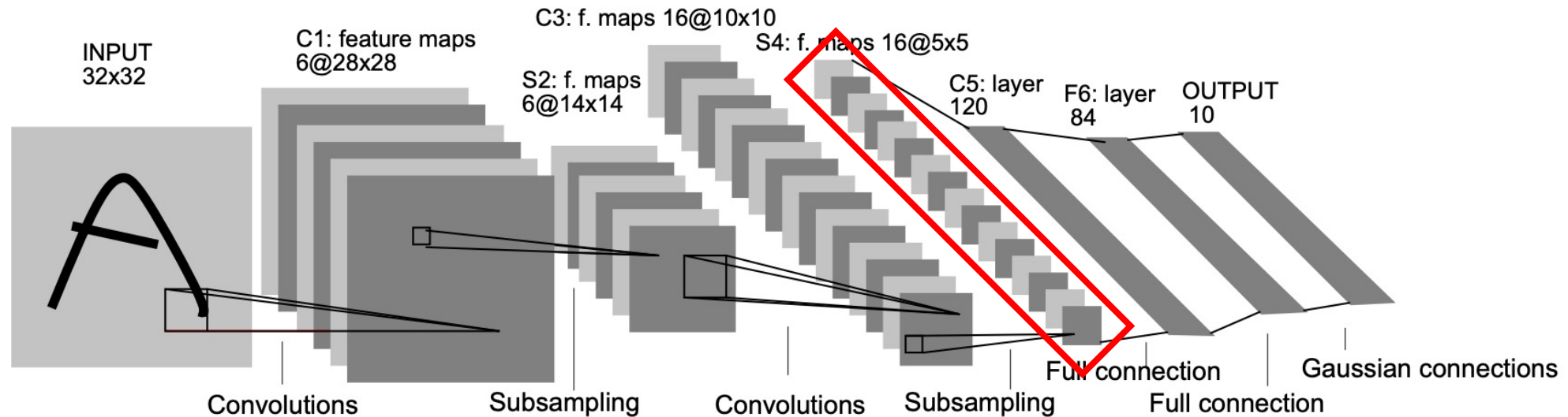


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

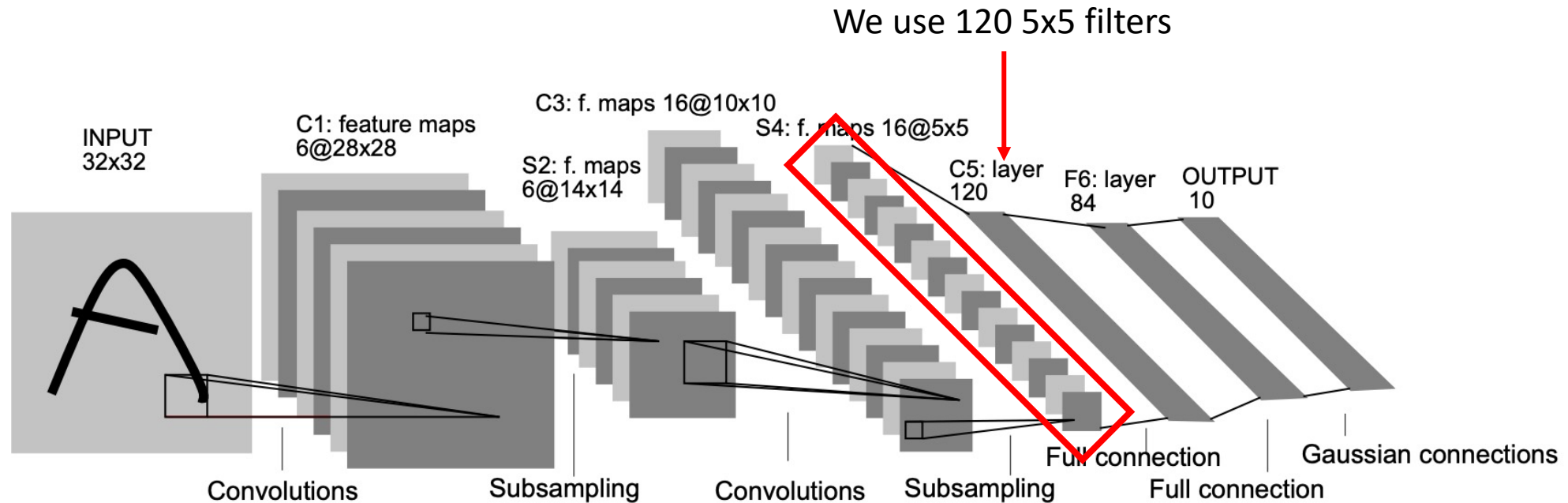


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

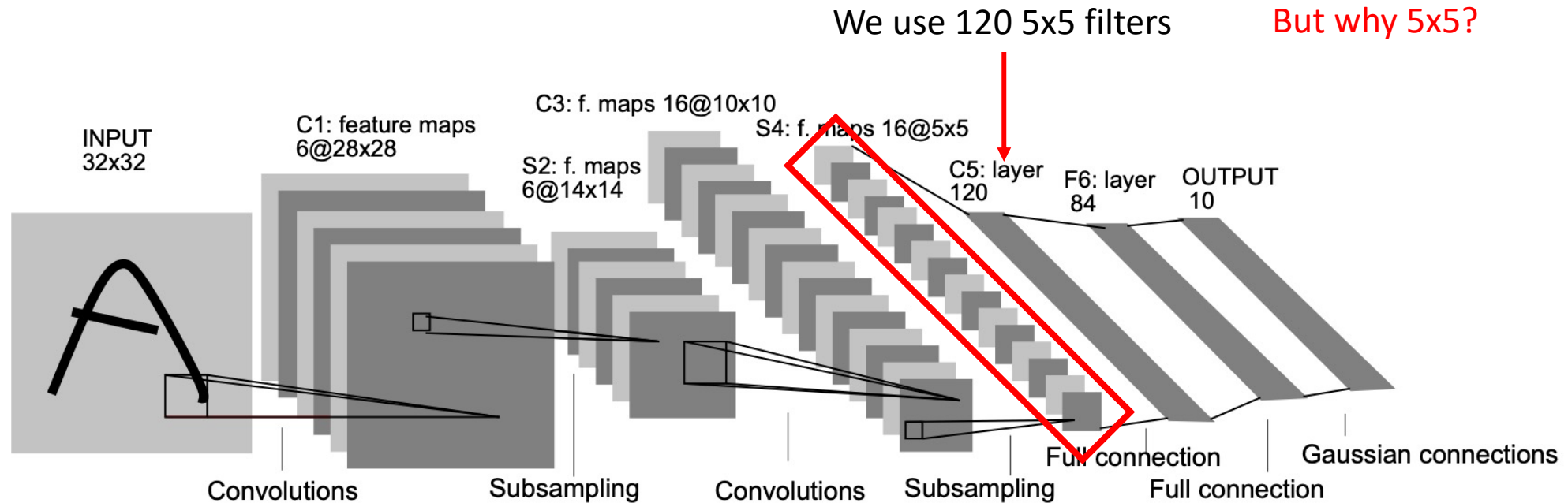


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

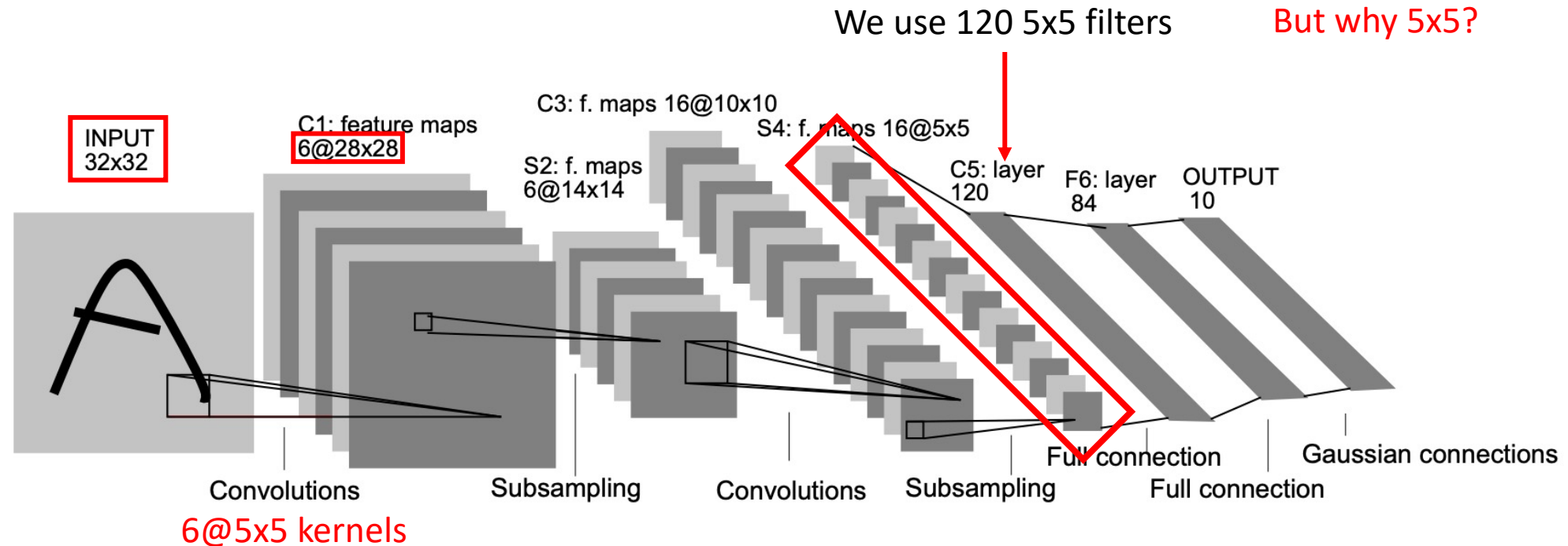


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

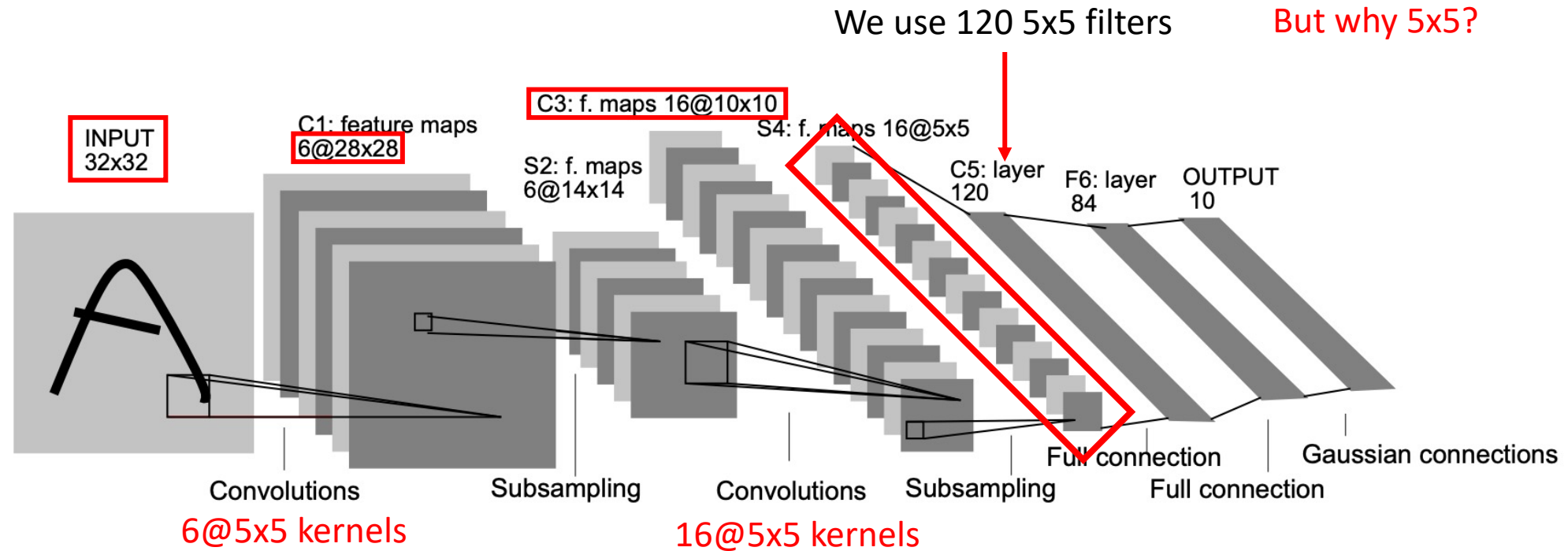


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

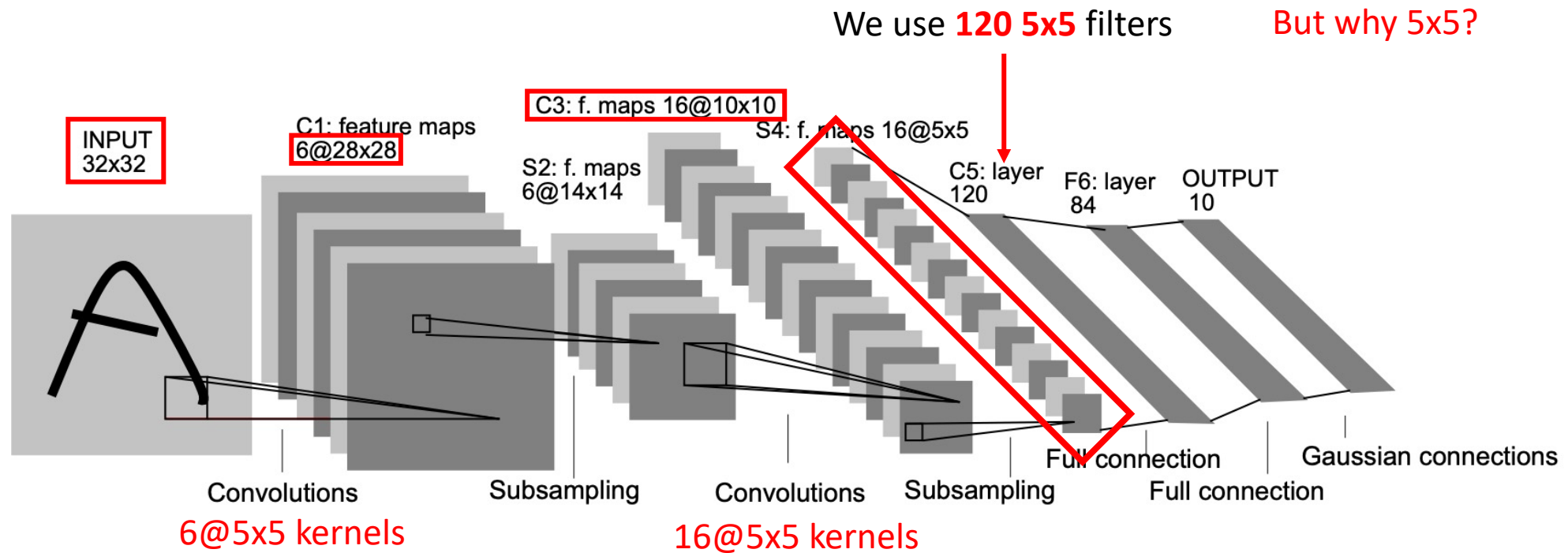


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

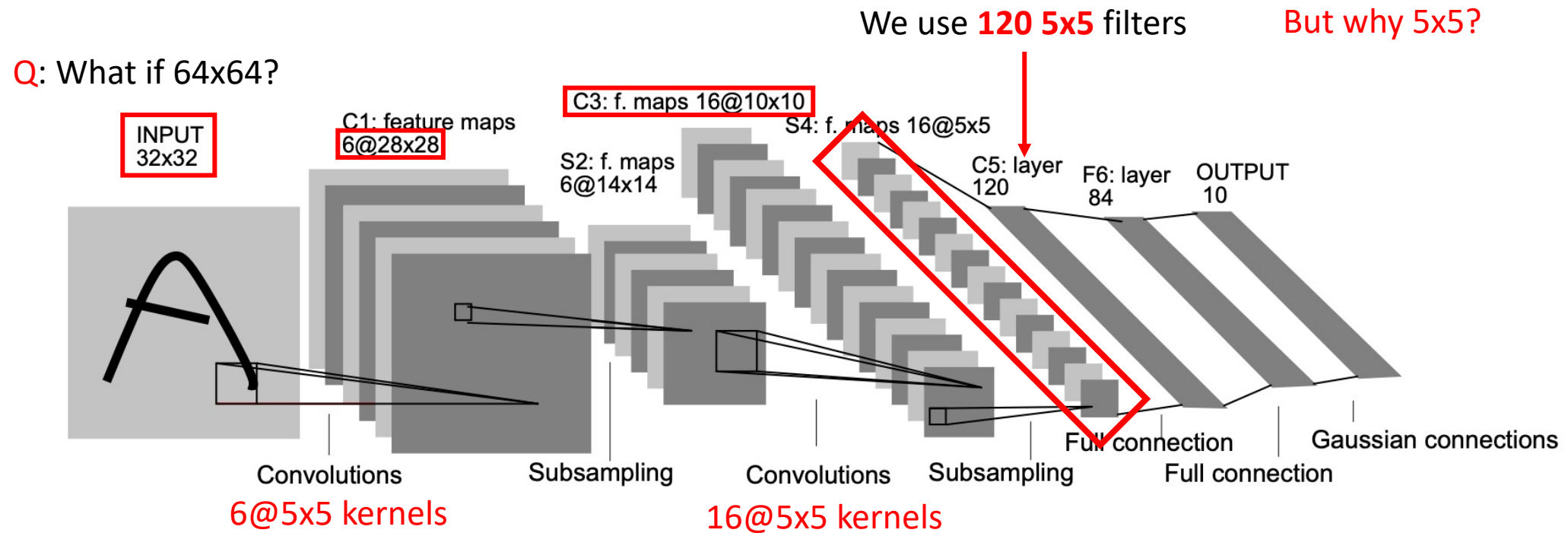
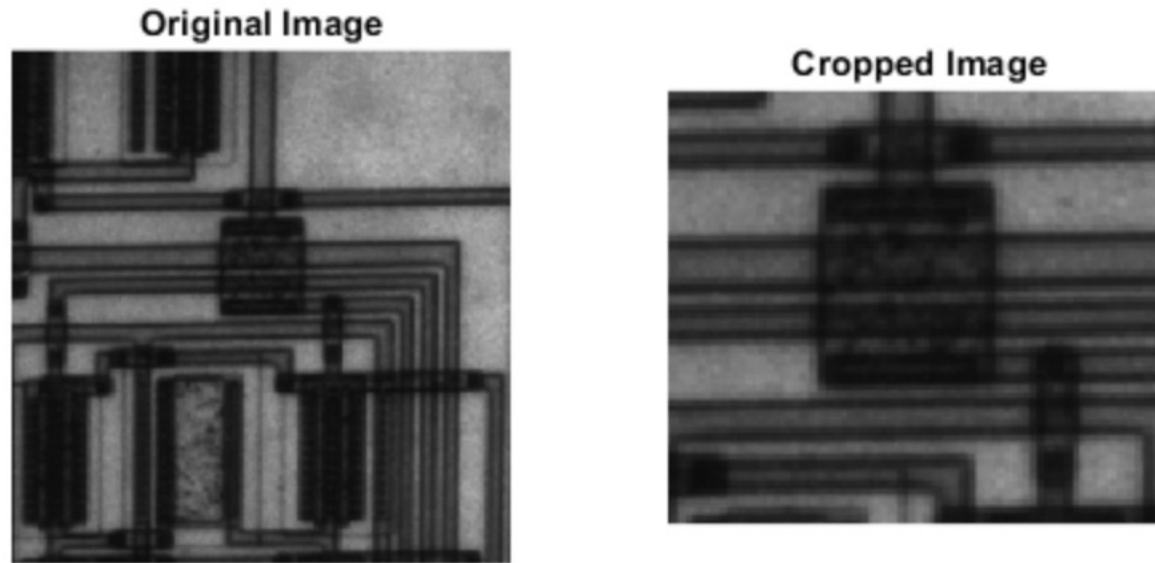


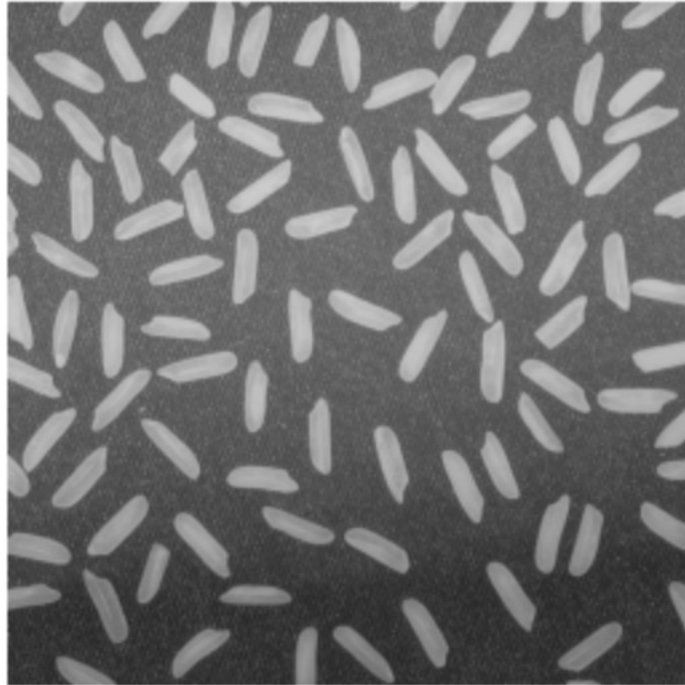
Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

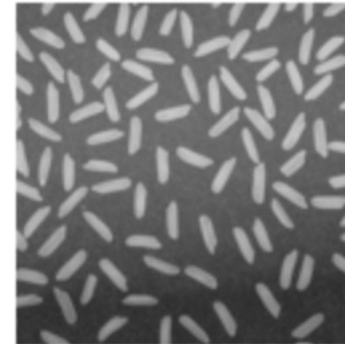


Input resolution issue

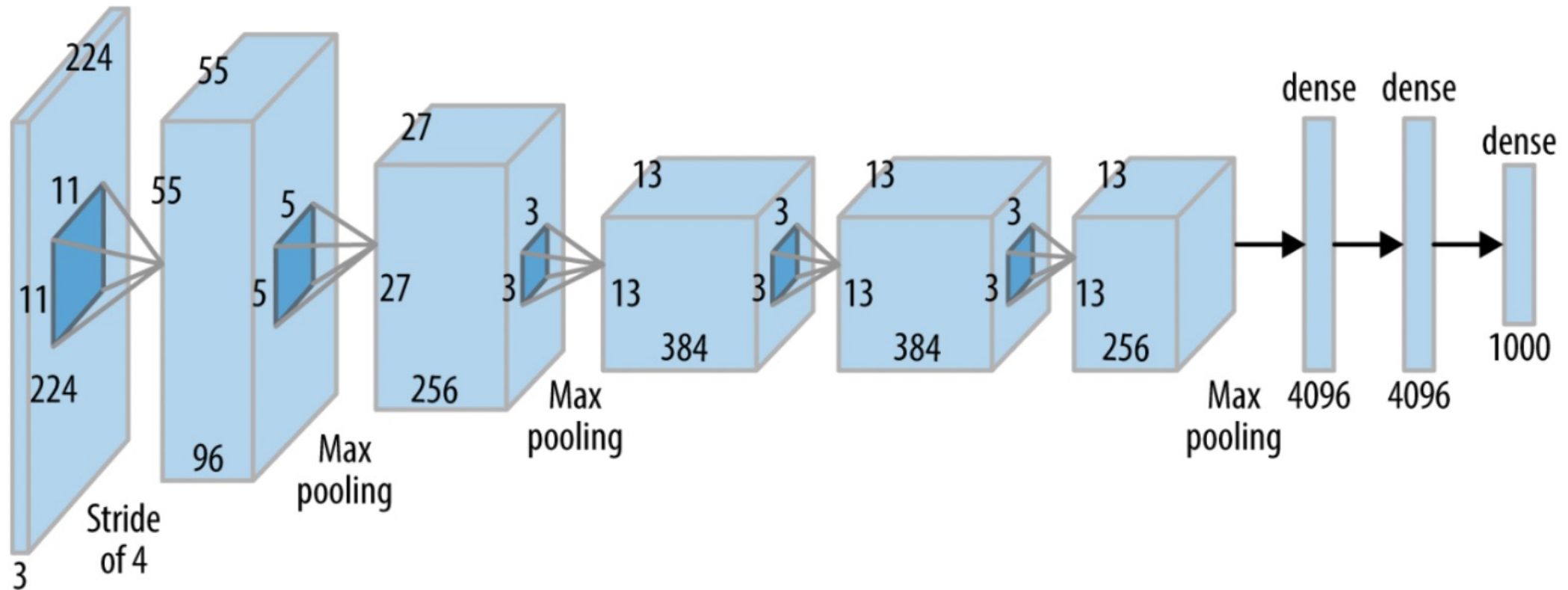
Original Image



Resized Image

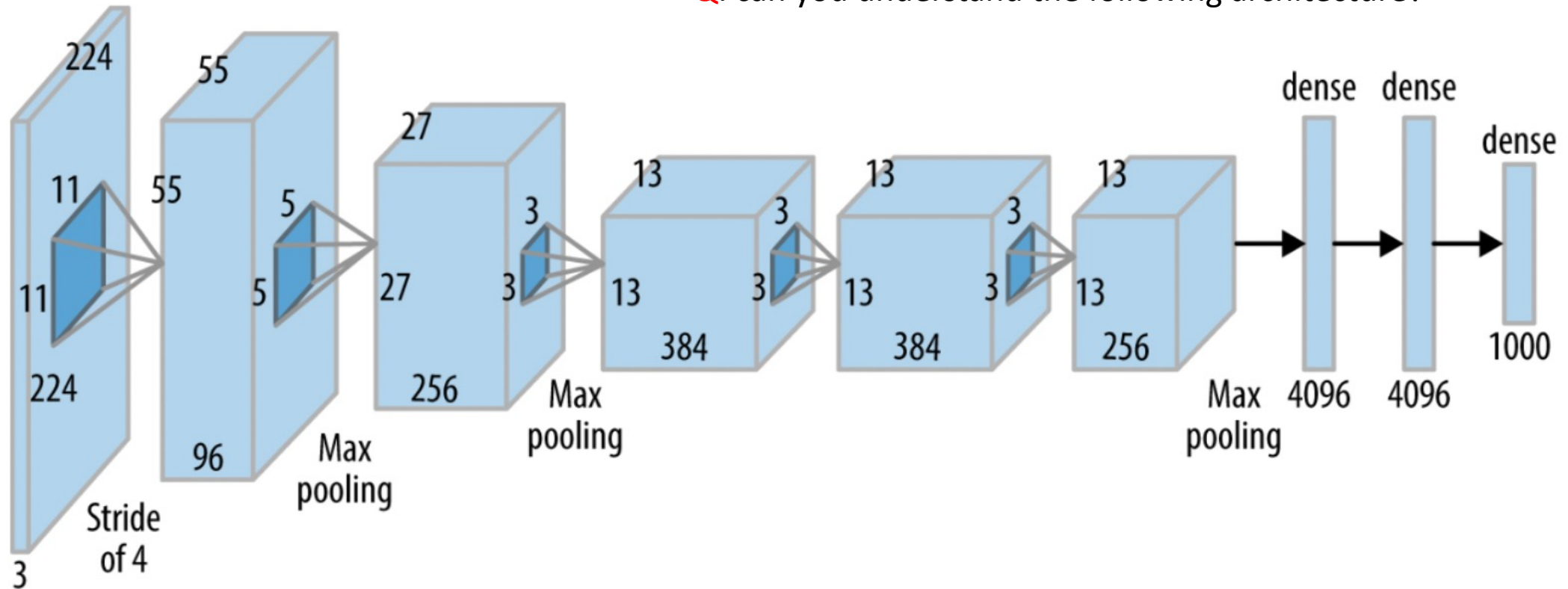


Input resolution issue

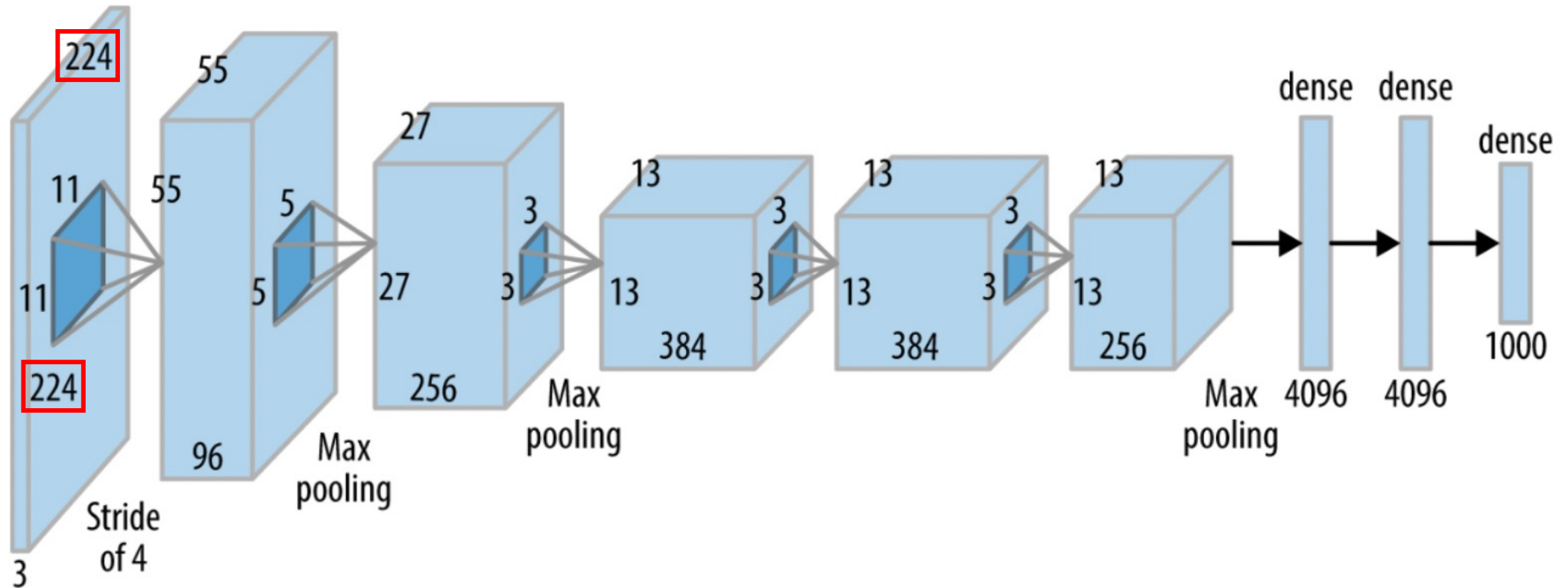


Input resolution issue

Q: can you understand the following architecture?

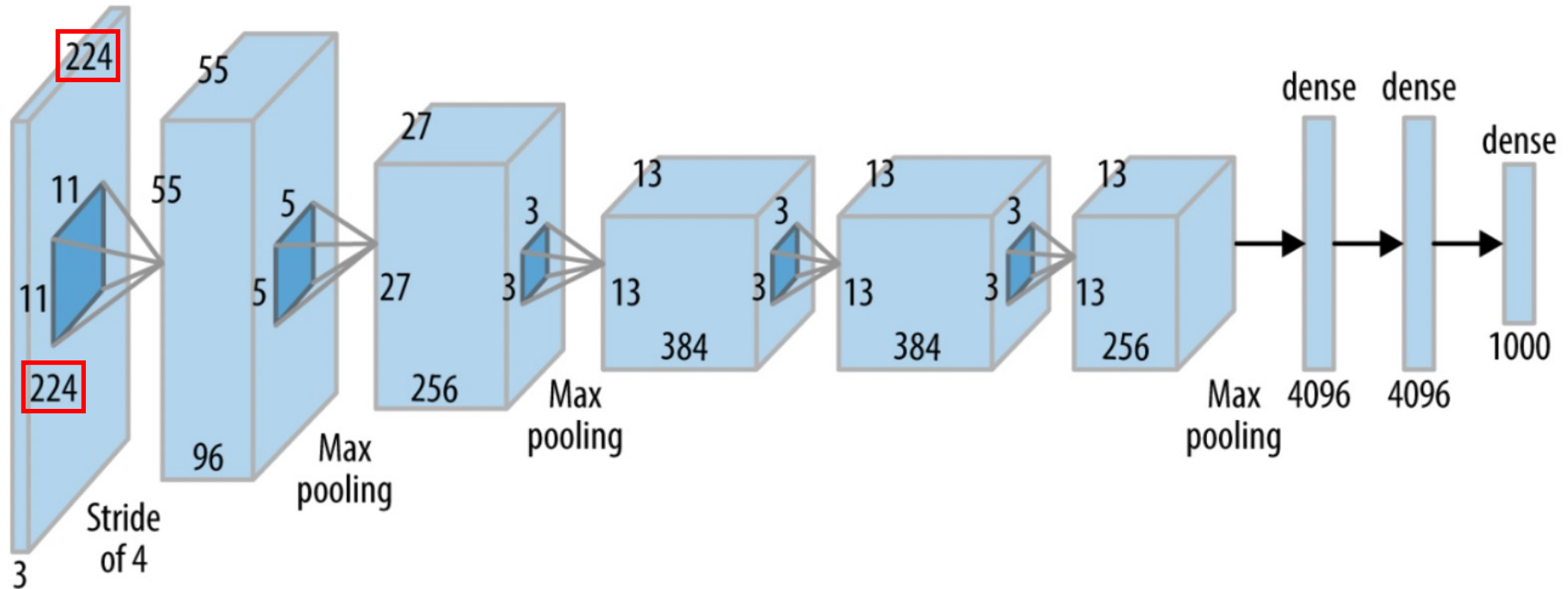


Input resolution issue



Any input image must be 224x224

Input resolution issue



Any input image must be 224x224

Q: how to handle an arbitrary resolution?

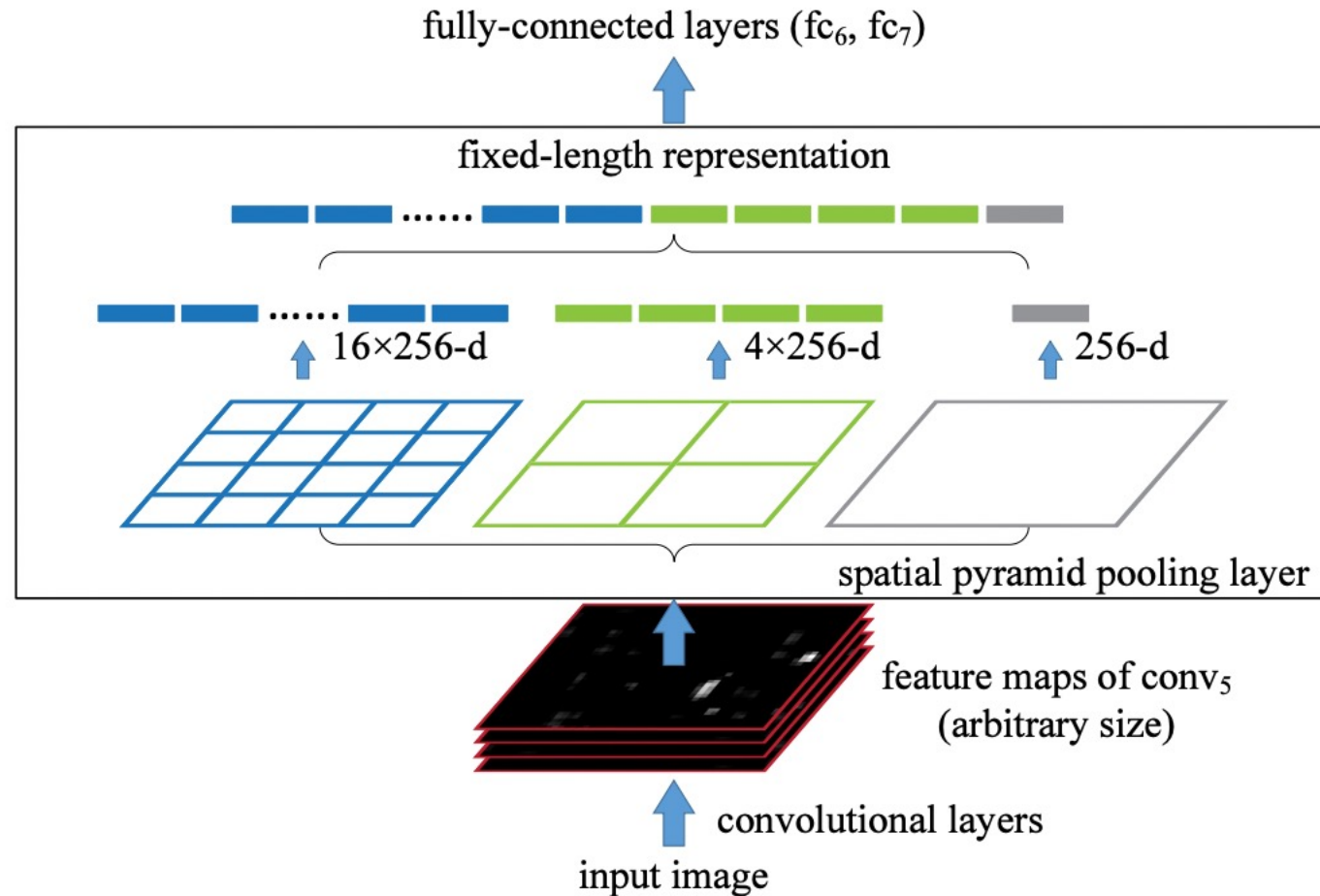
[Alexnet]

Input resolution issue

- Spatial pyramid pooling [pyramid]
- Global average pooling [NIN]
-

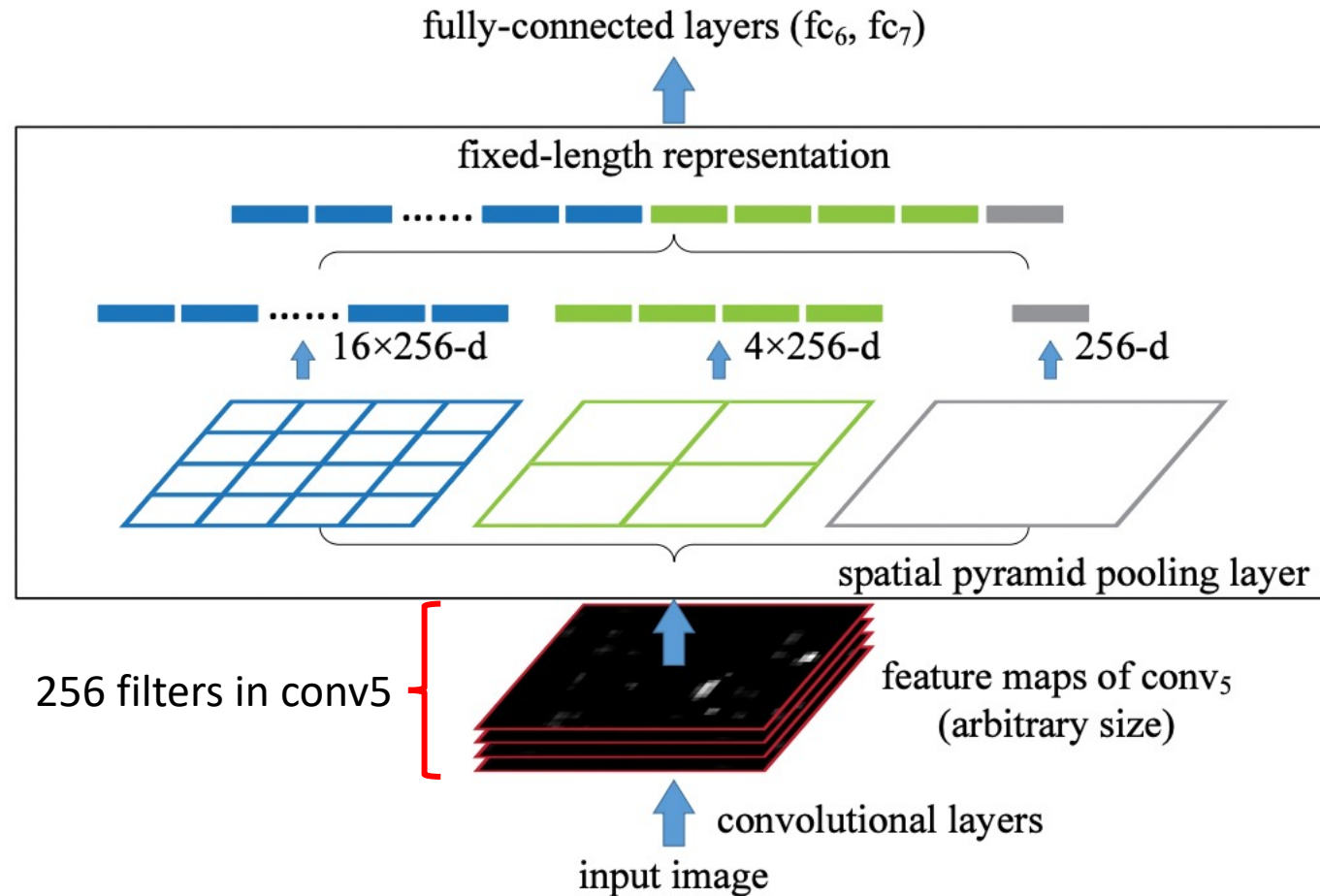
Input resolution issue

- Spatial pyramid pooling



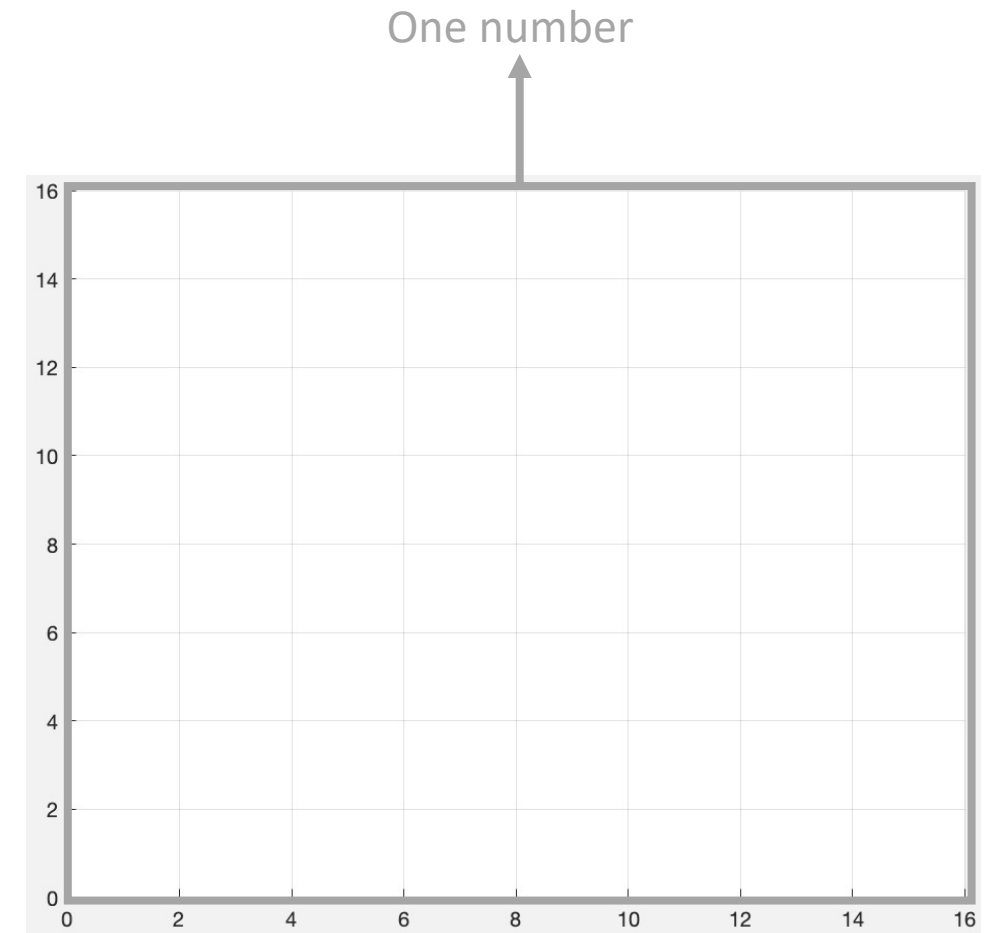
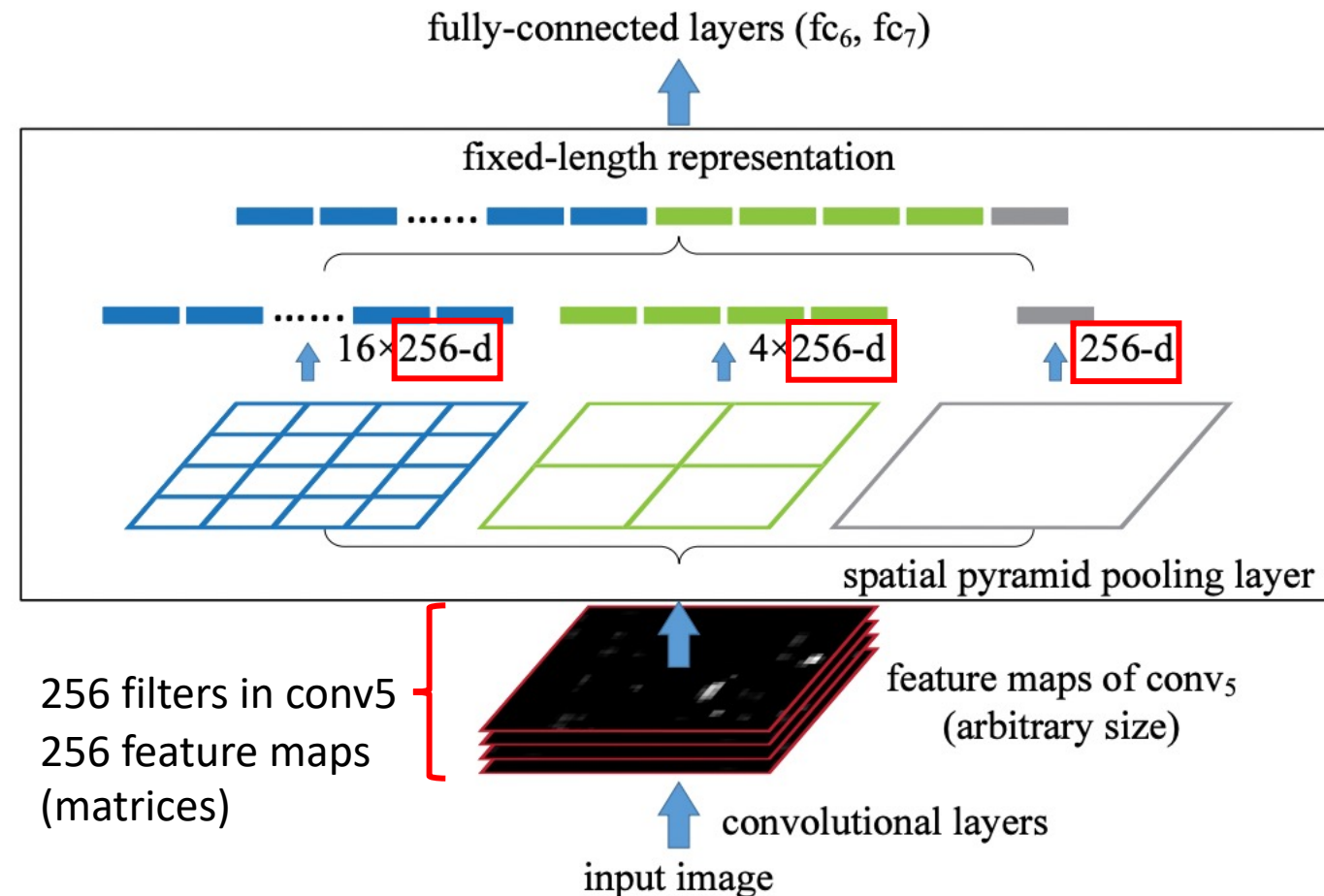
Input resolution issue

- Spatial pyramid pooling



Input resolution issue

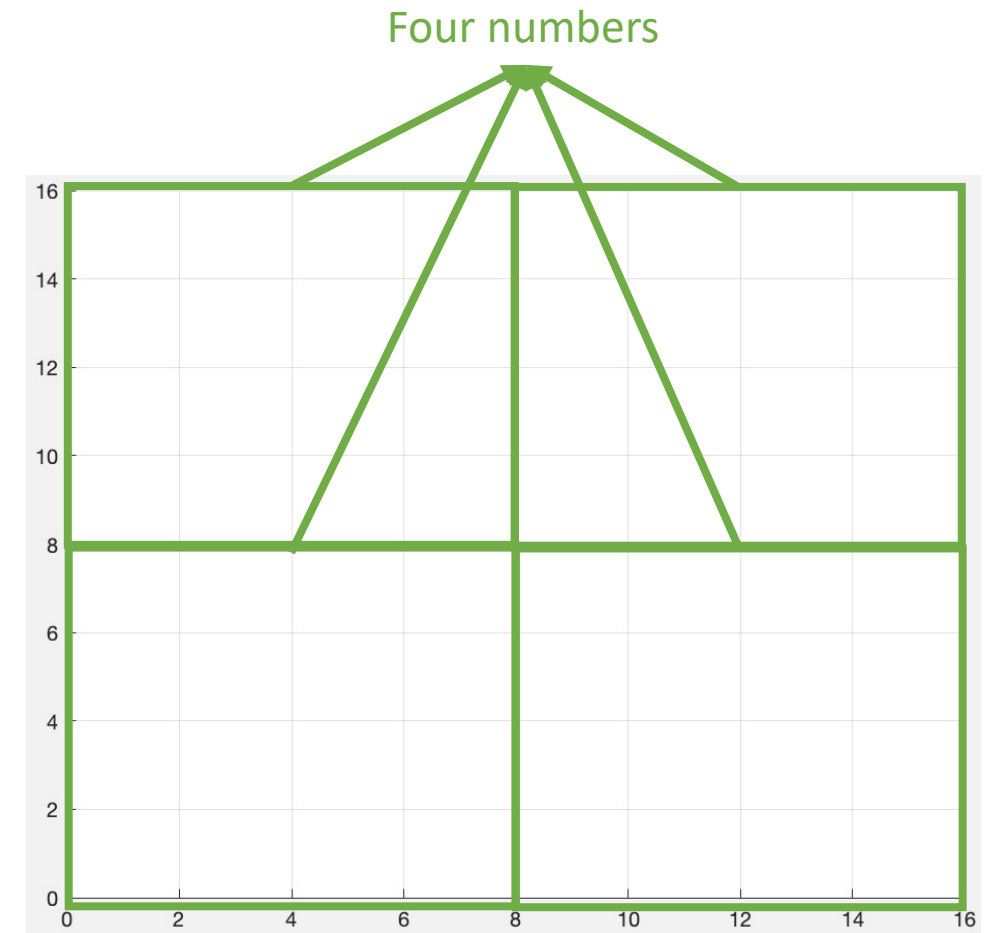
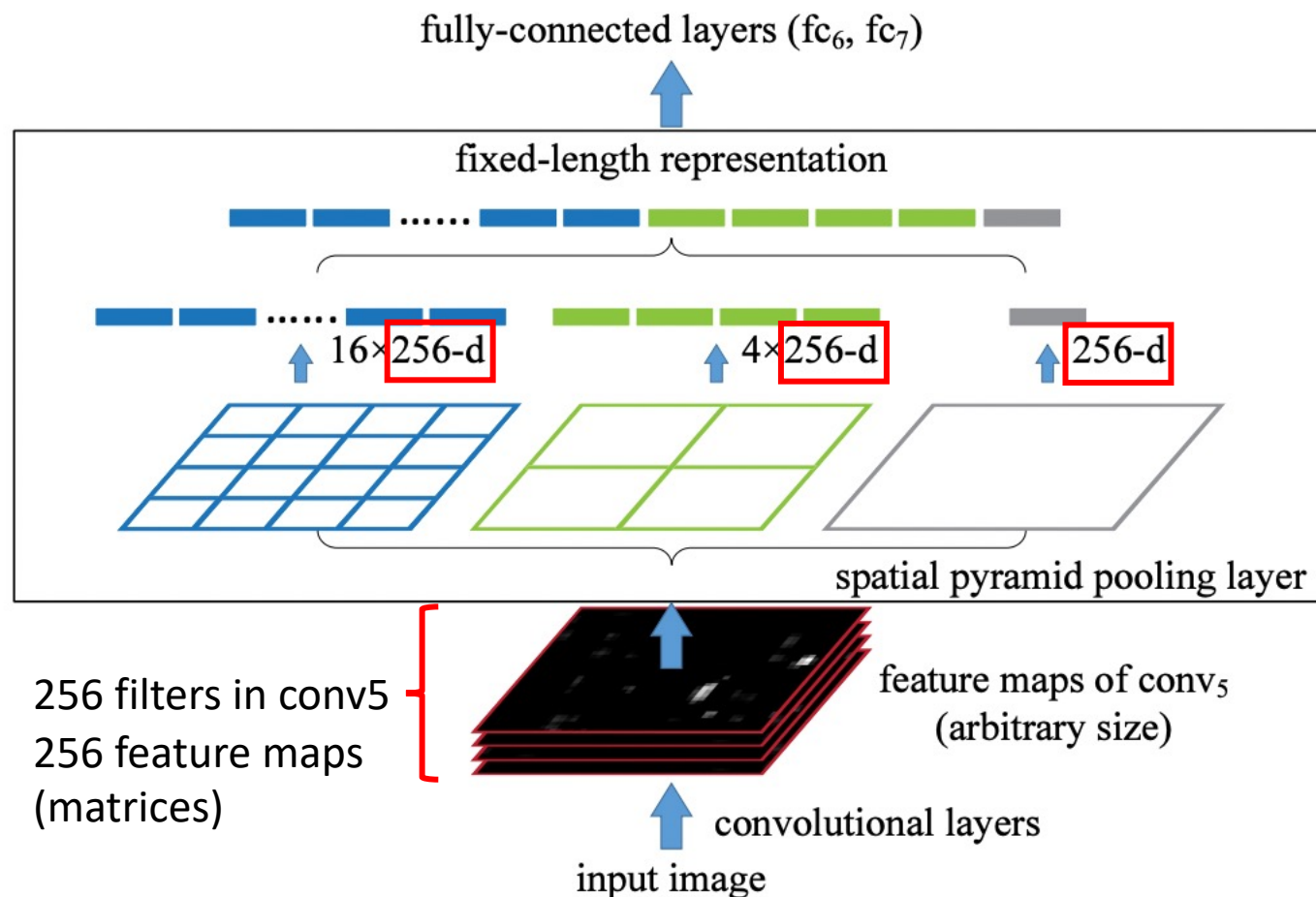
- Spatial pyramid pooling



Some pooling (max/average)

Input resolution issue

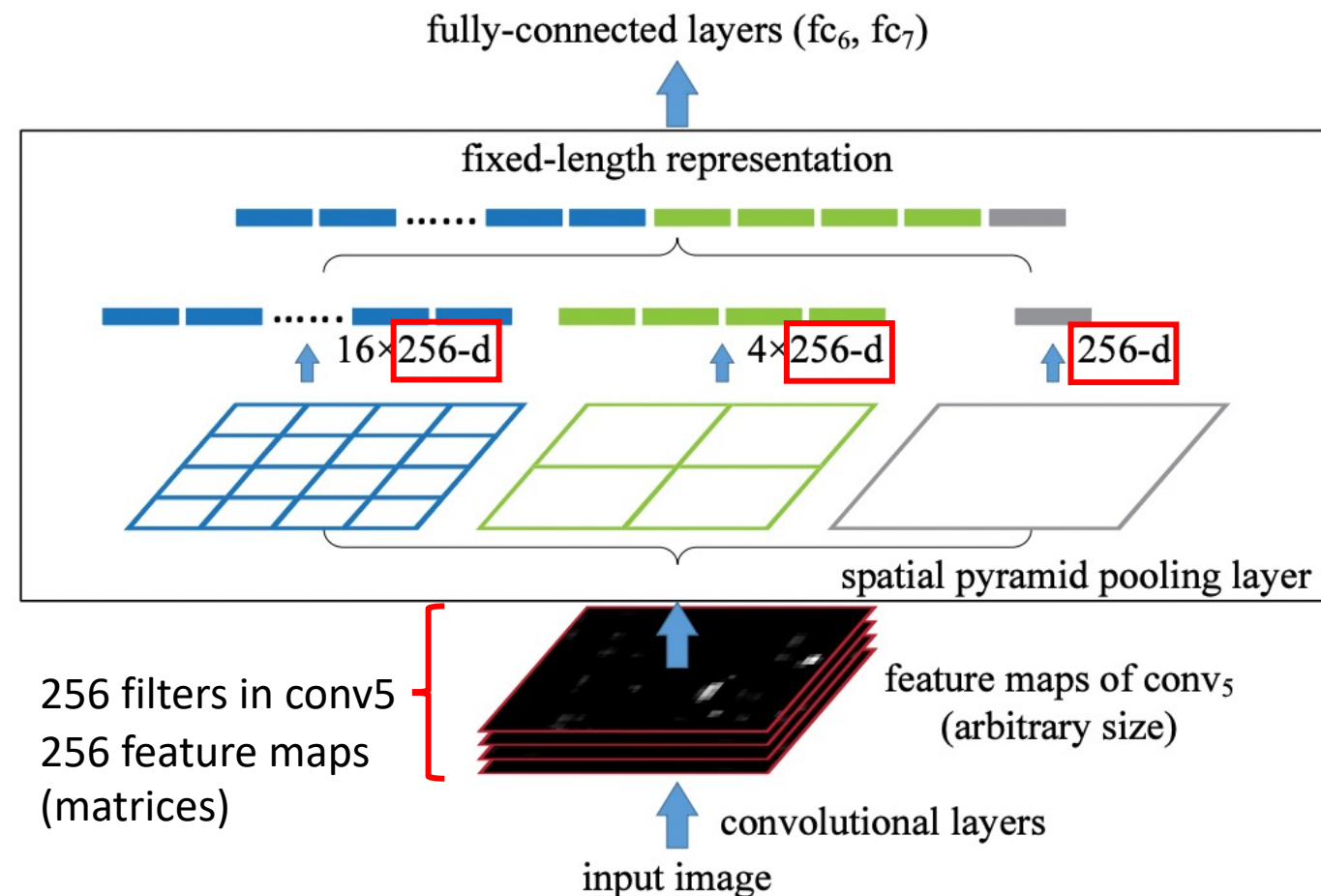
- Spatial pyramid pooling



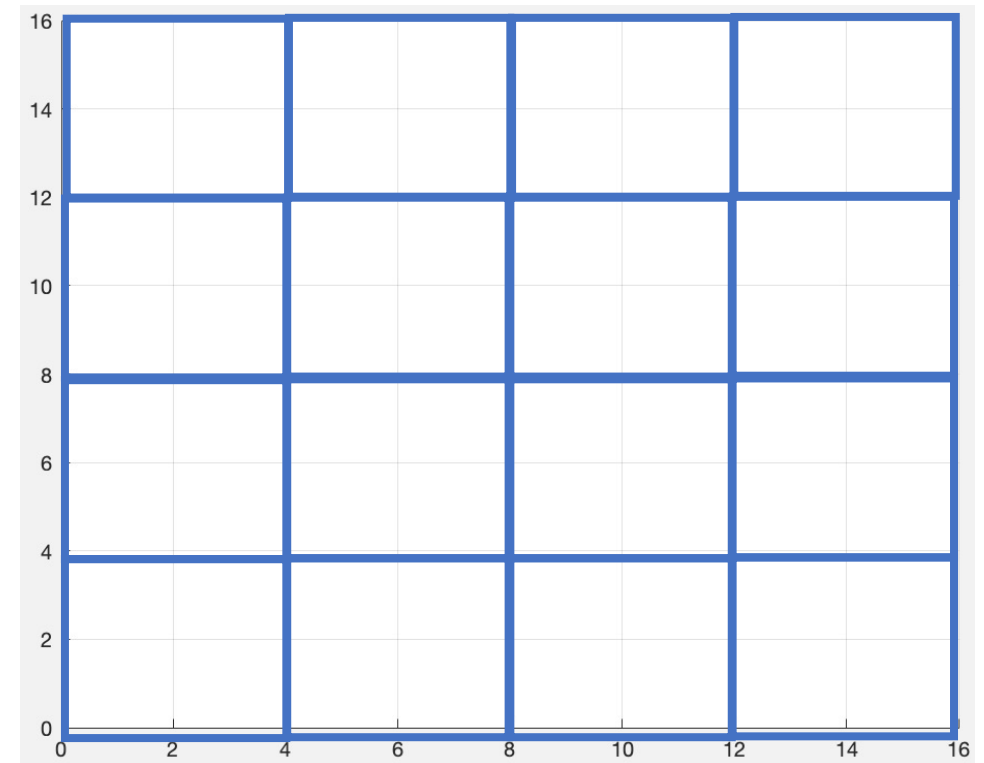
Some pooling (max/average)

Input resolution issue

- Spatial pyramid pooling



16 numbers



Some pooling (max/average)

Input resolution issue

- Spatial pyramid pooling

fully-connected layers (fc_6, fc_7)

fixed-length representation

Concatenation:
(1+4+16) x 256 numbers

16 x 256-d

4 x 256-d

256-d

spatial pyramid pooling layer

256 filters in conv5
256 feature maps
(matrices)

feature maps of conv₅
(arbitrary size)

convolutional layers

input image

Input resolution issue

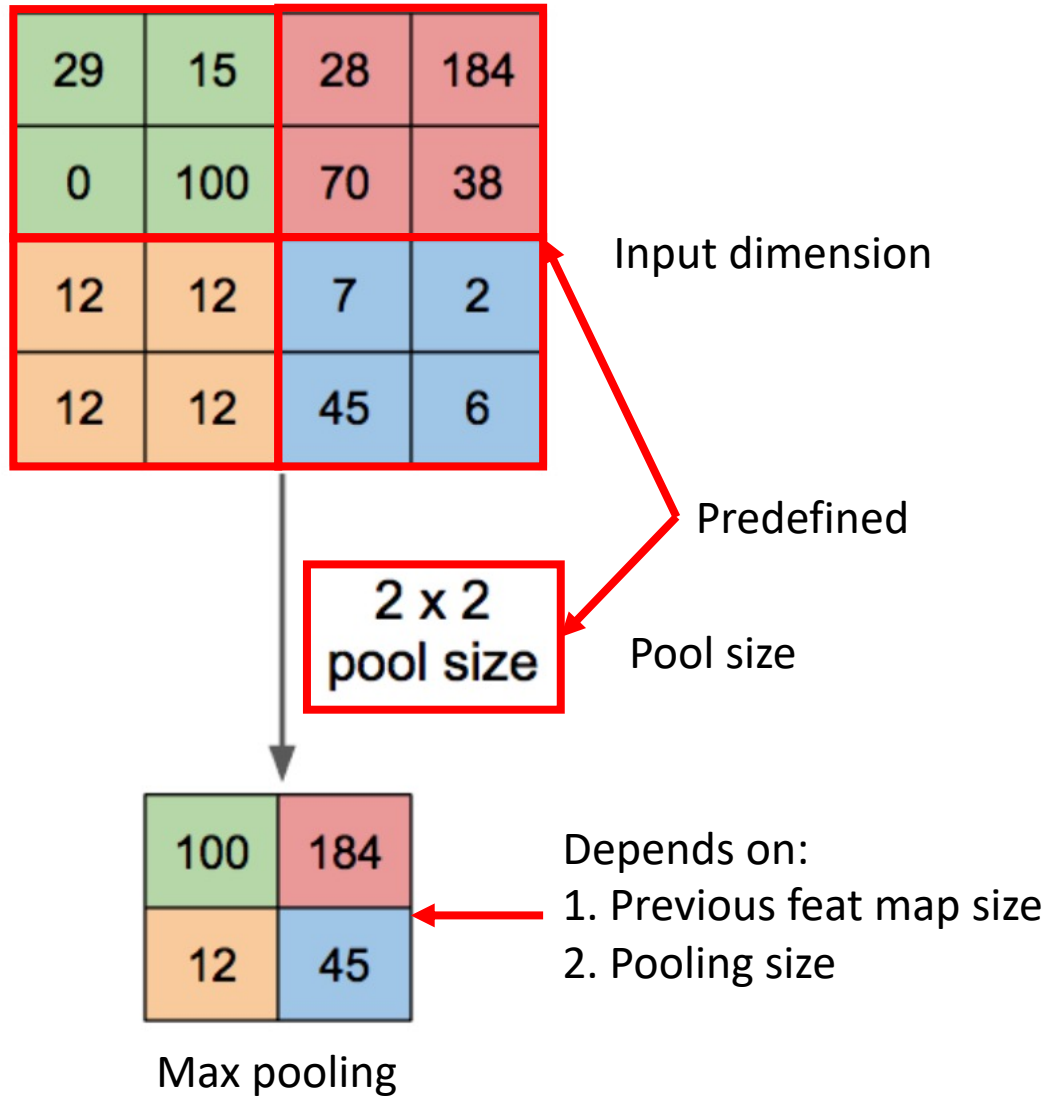
29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

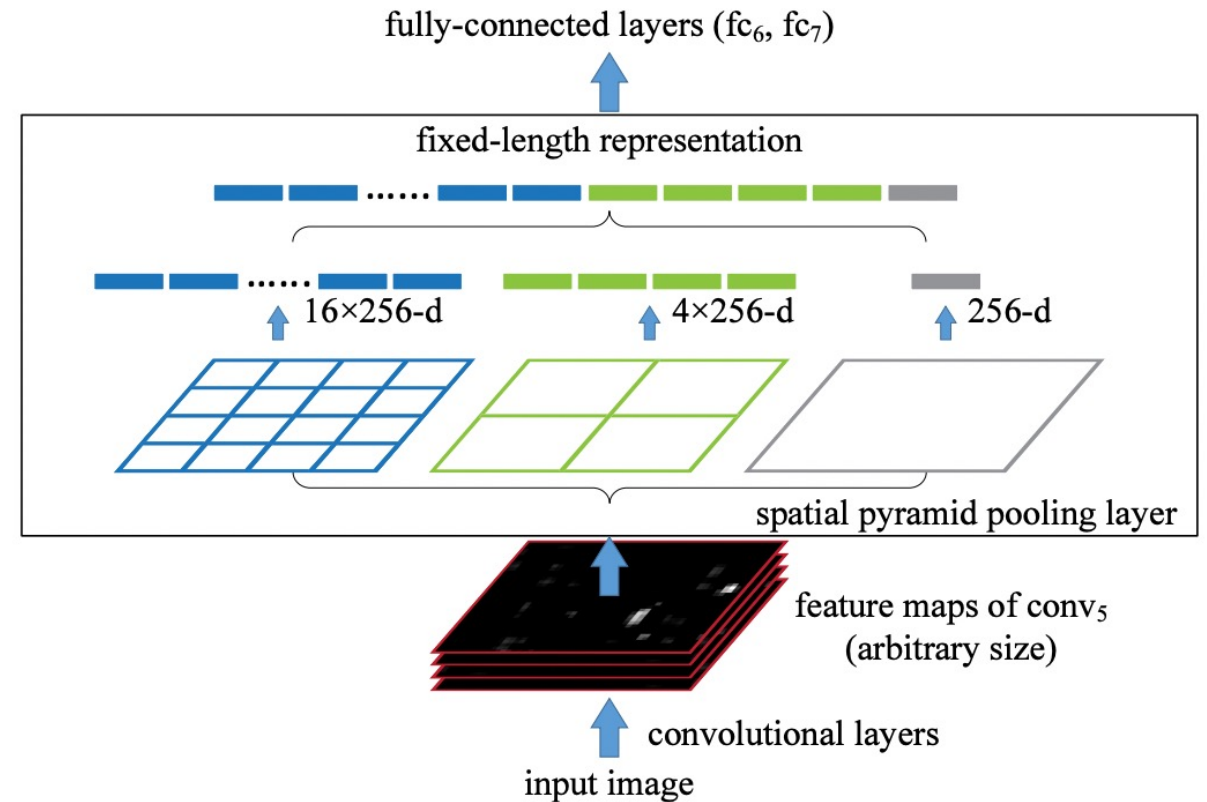
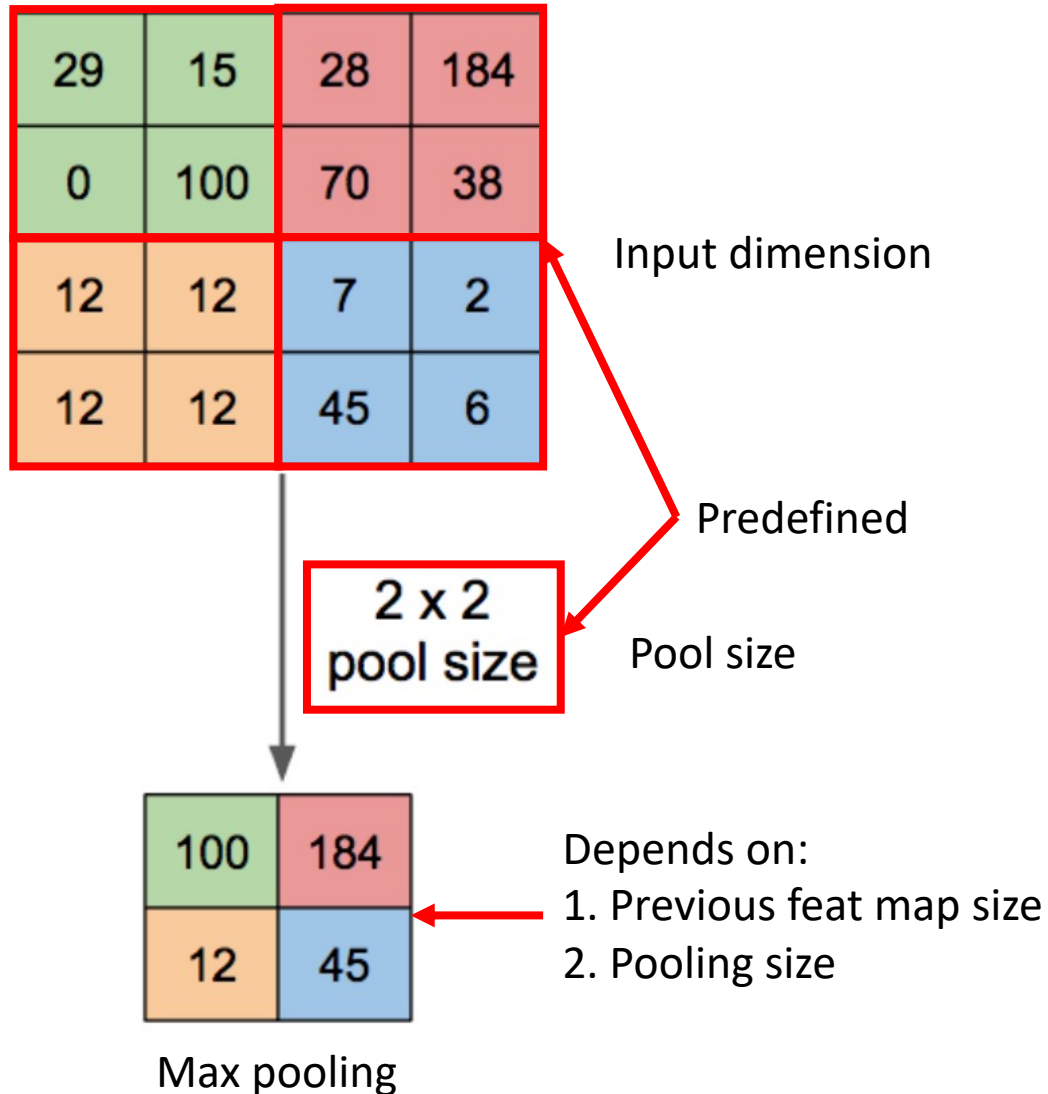
100	184
12	45

Max pooling

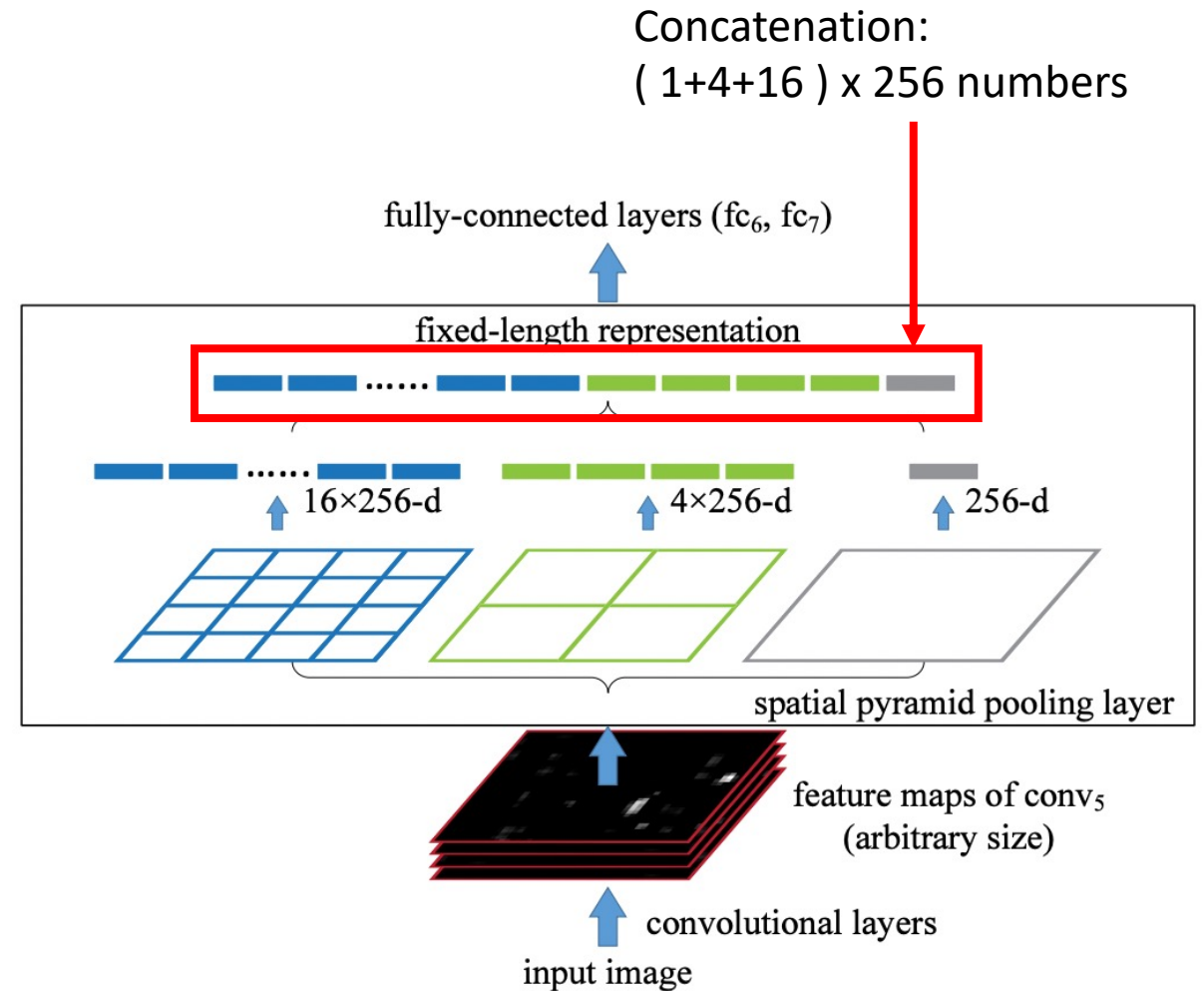
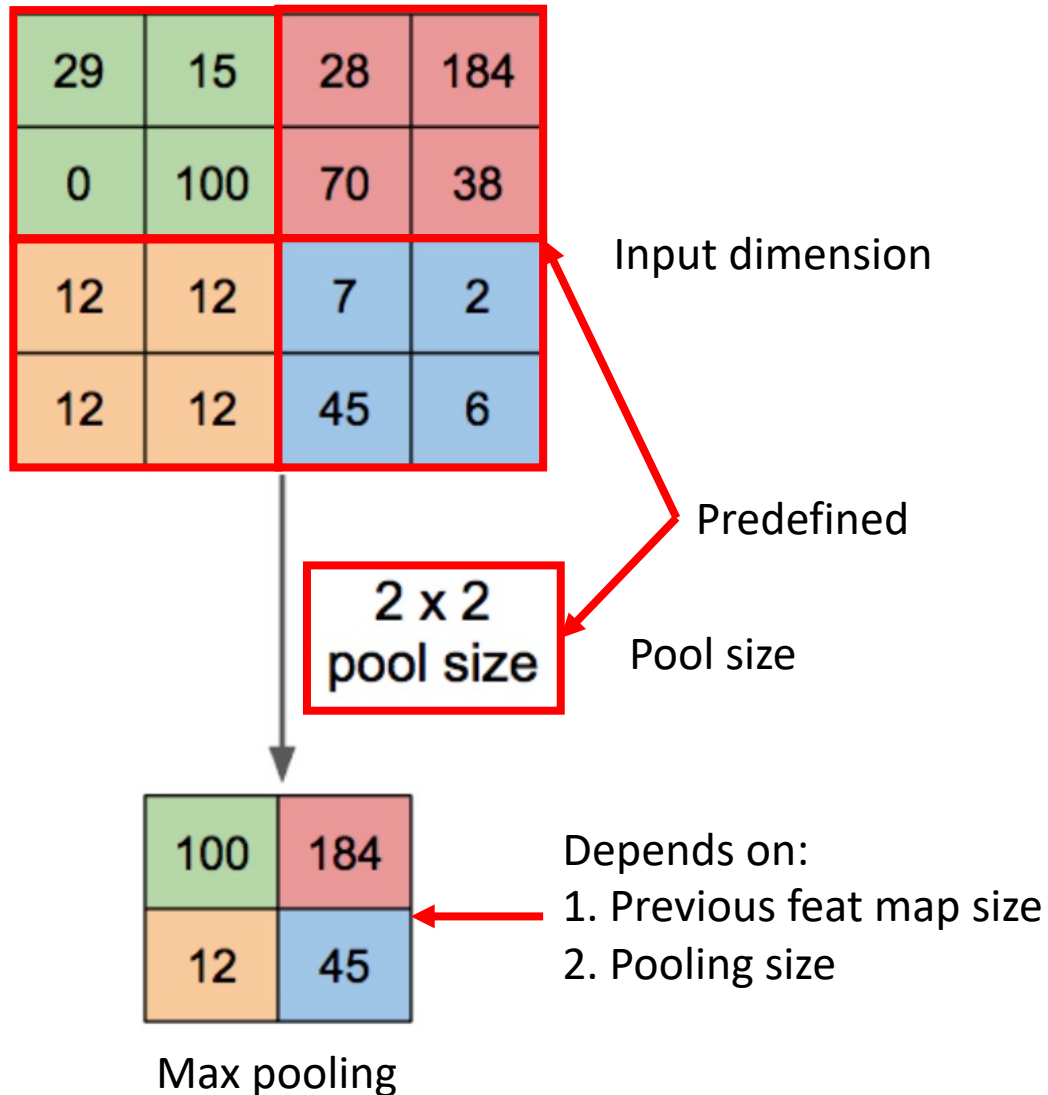
Input resolution issue



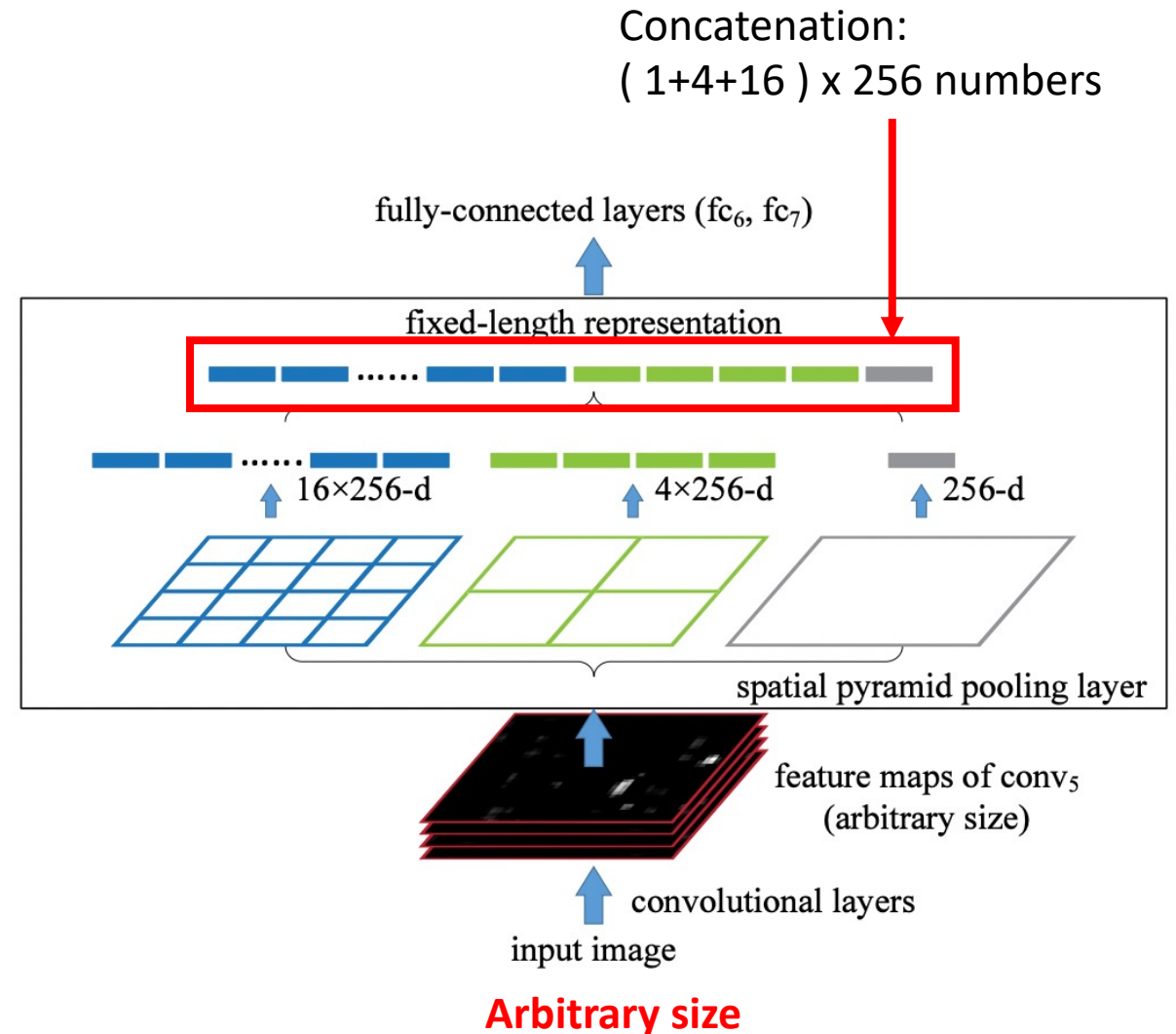
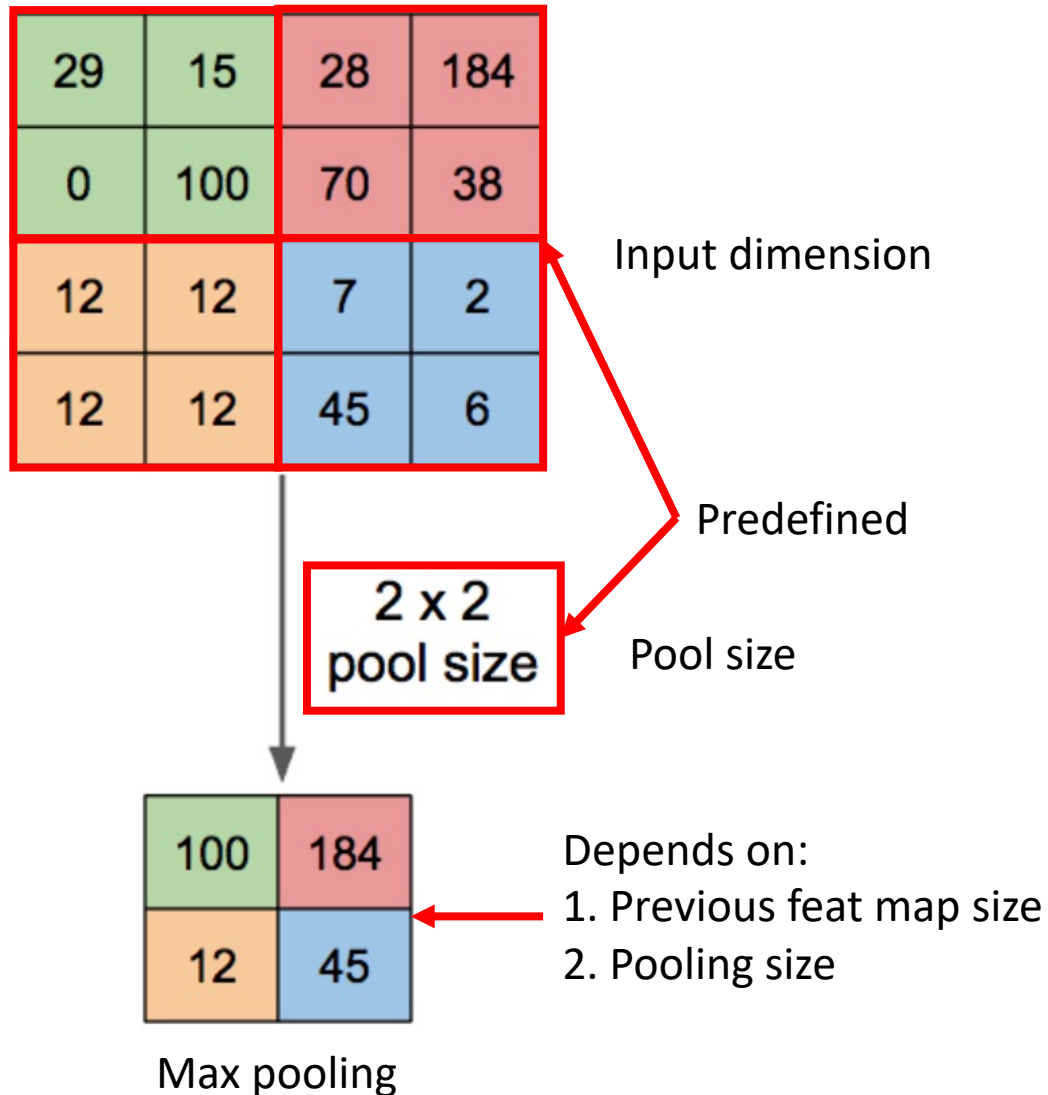
Input resolution issue



Input resolution issue



Input resolution issue



Input resolution issue

- Global average pooling

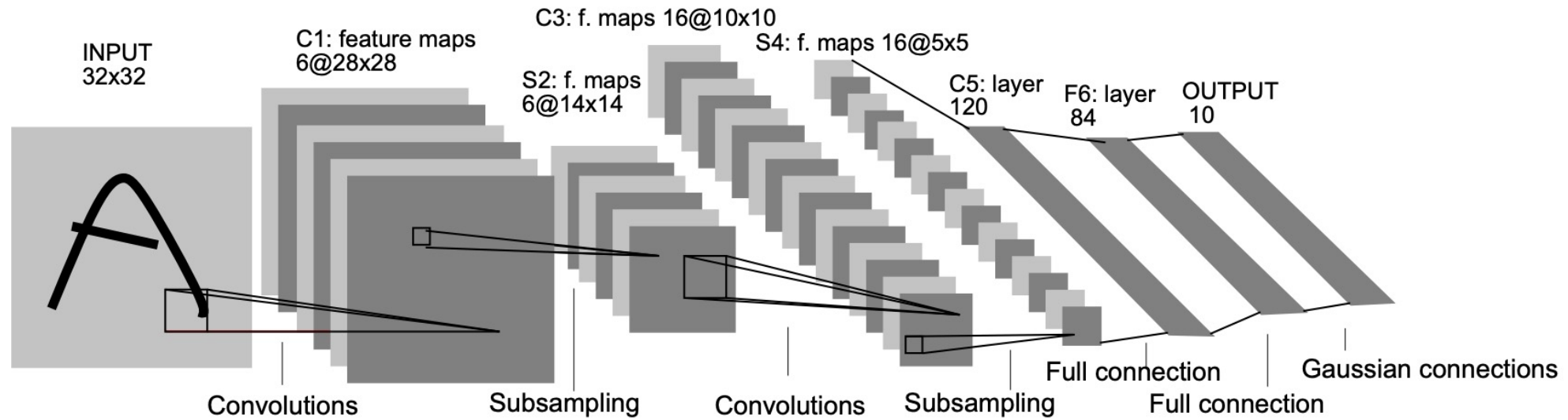


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

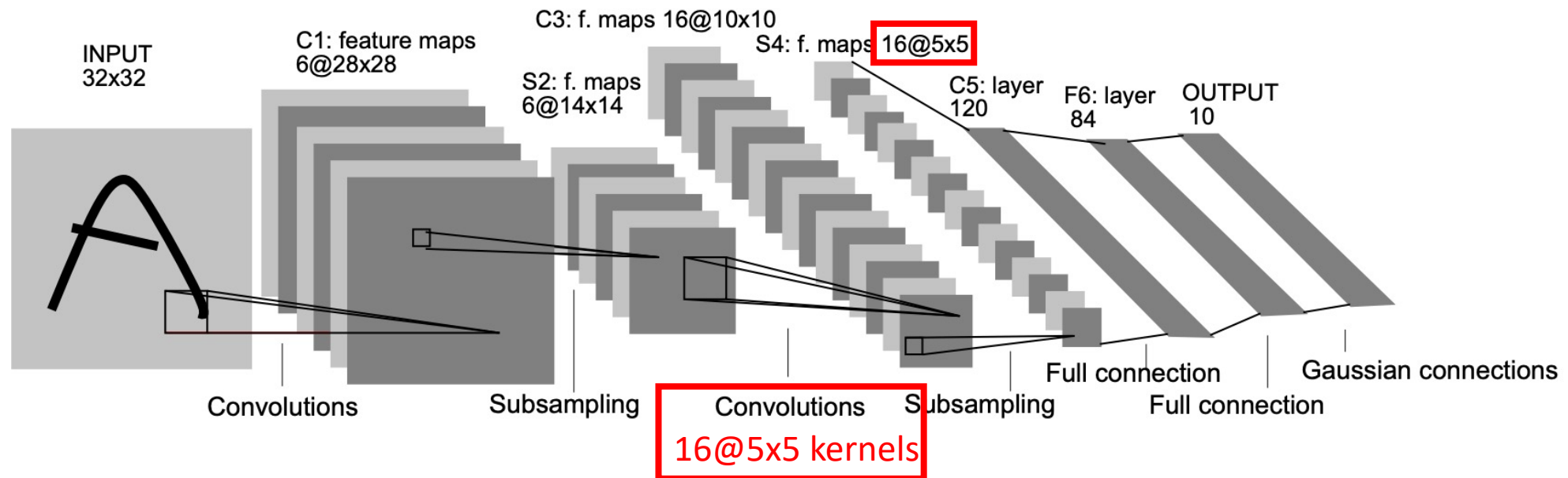


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

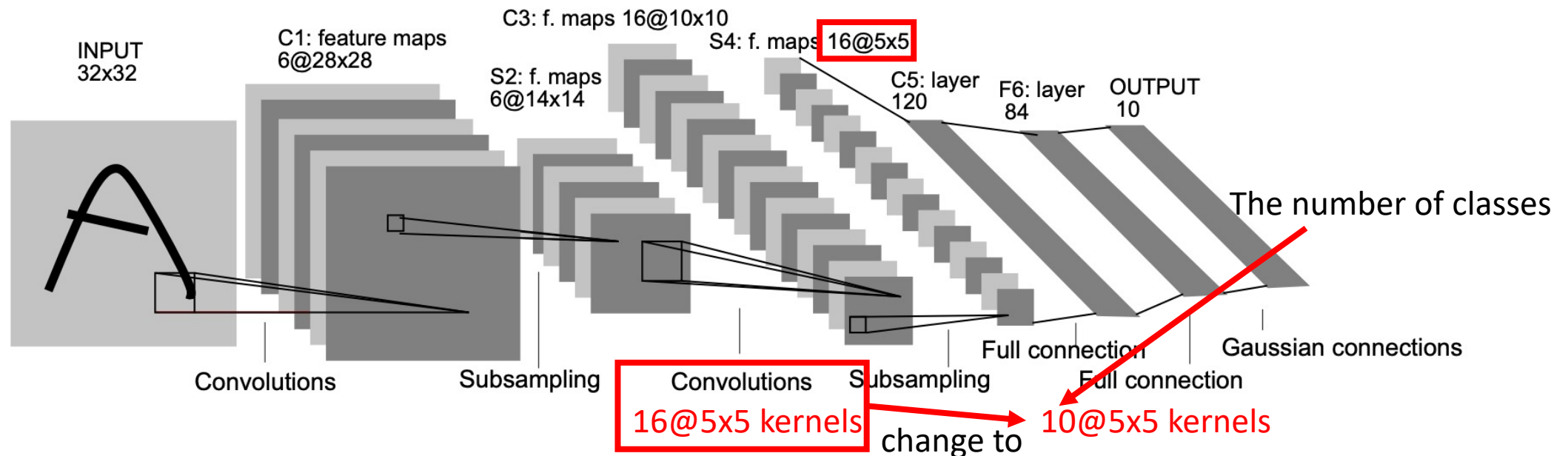


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

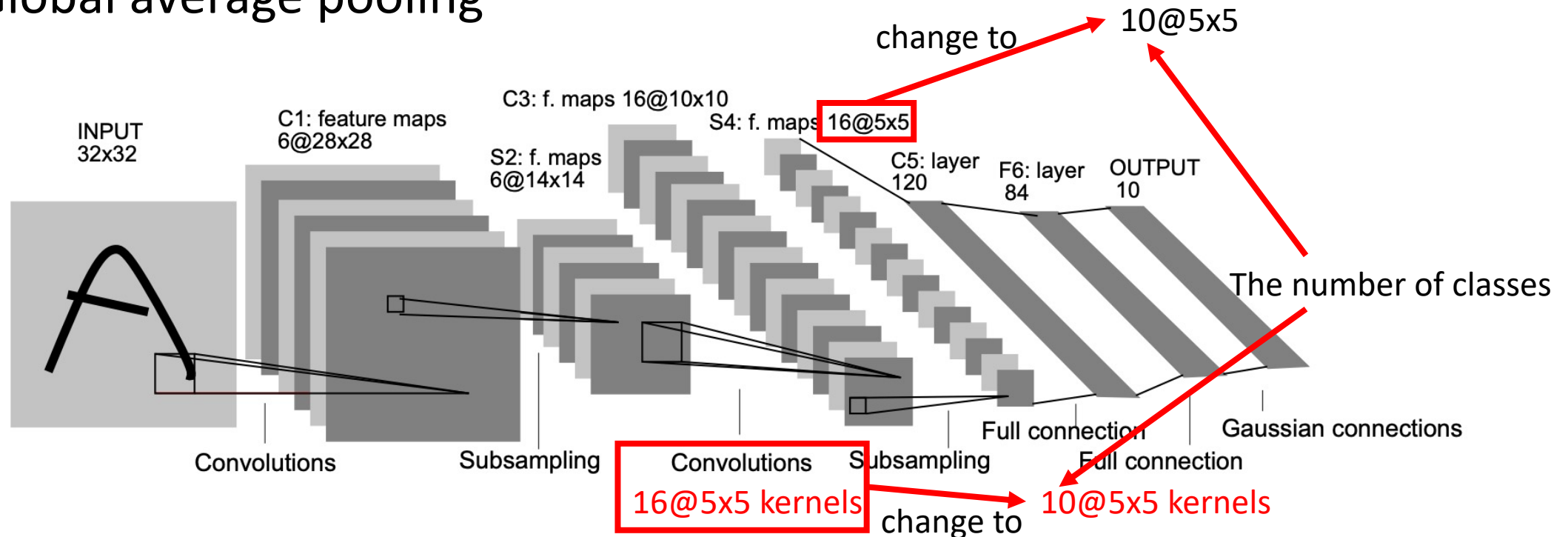


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

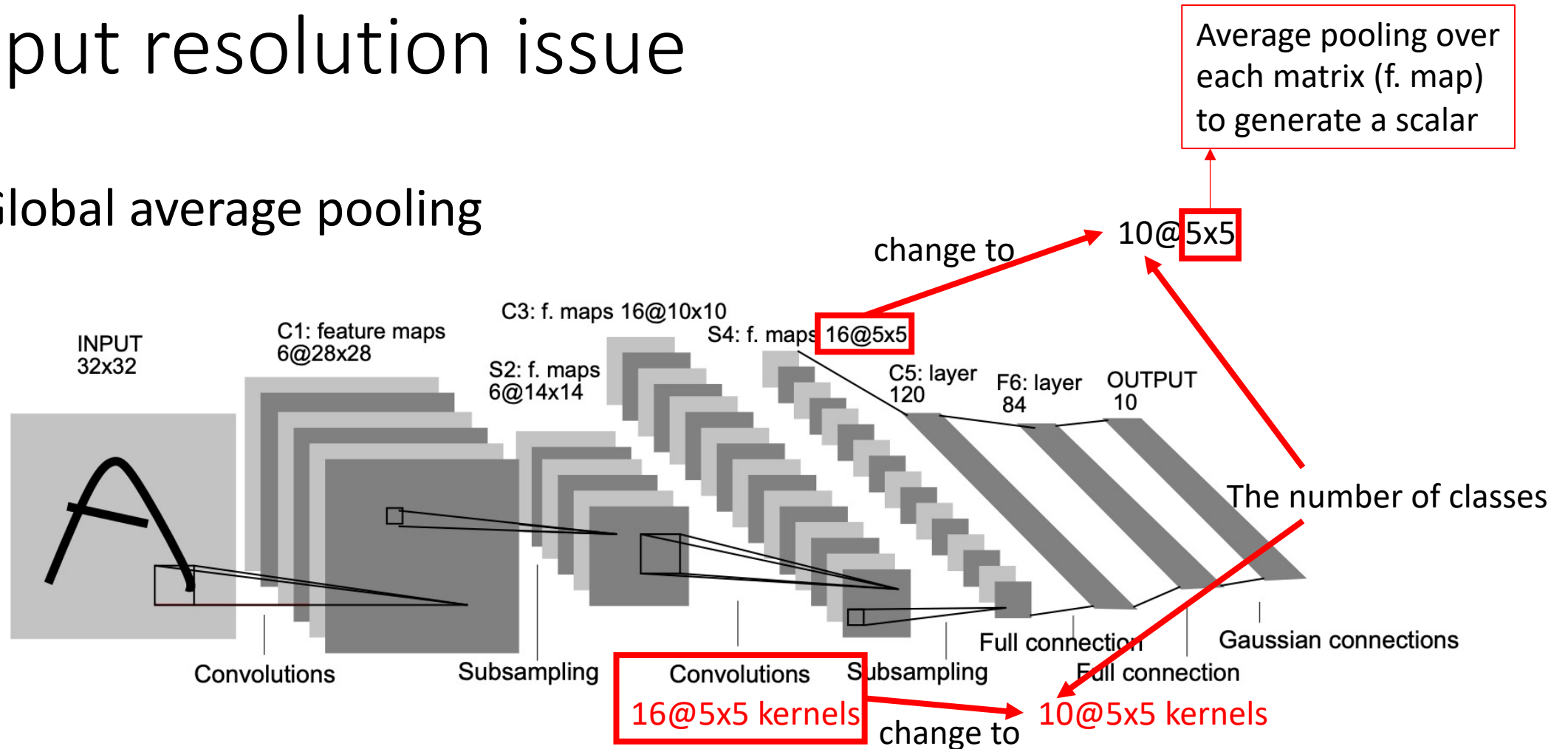


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

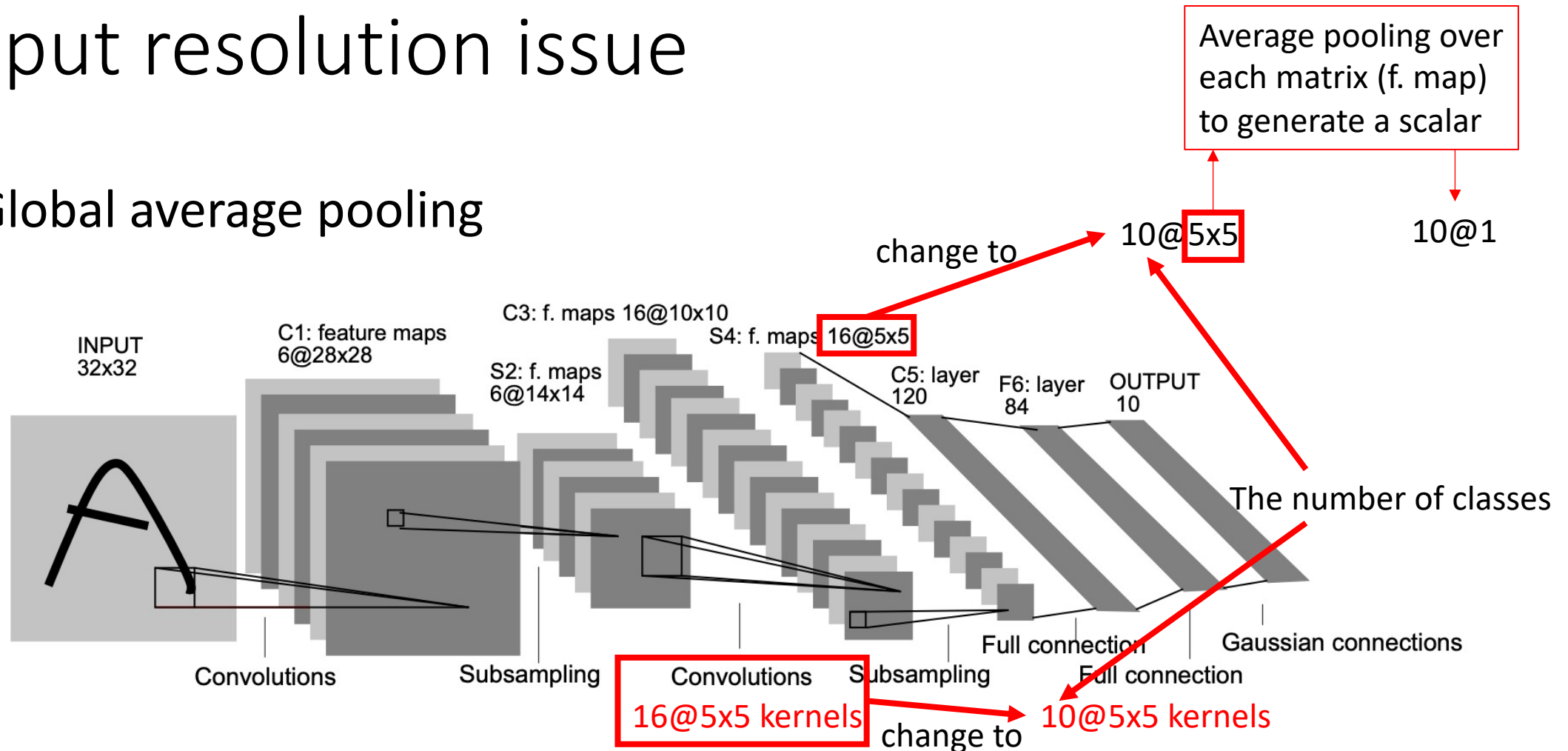


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Input resolution issue

- Global average pooling

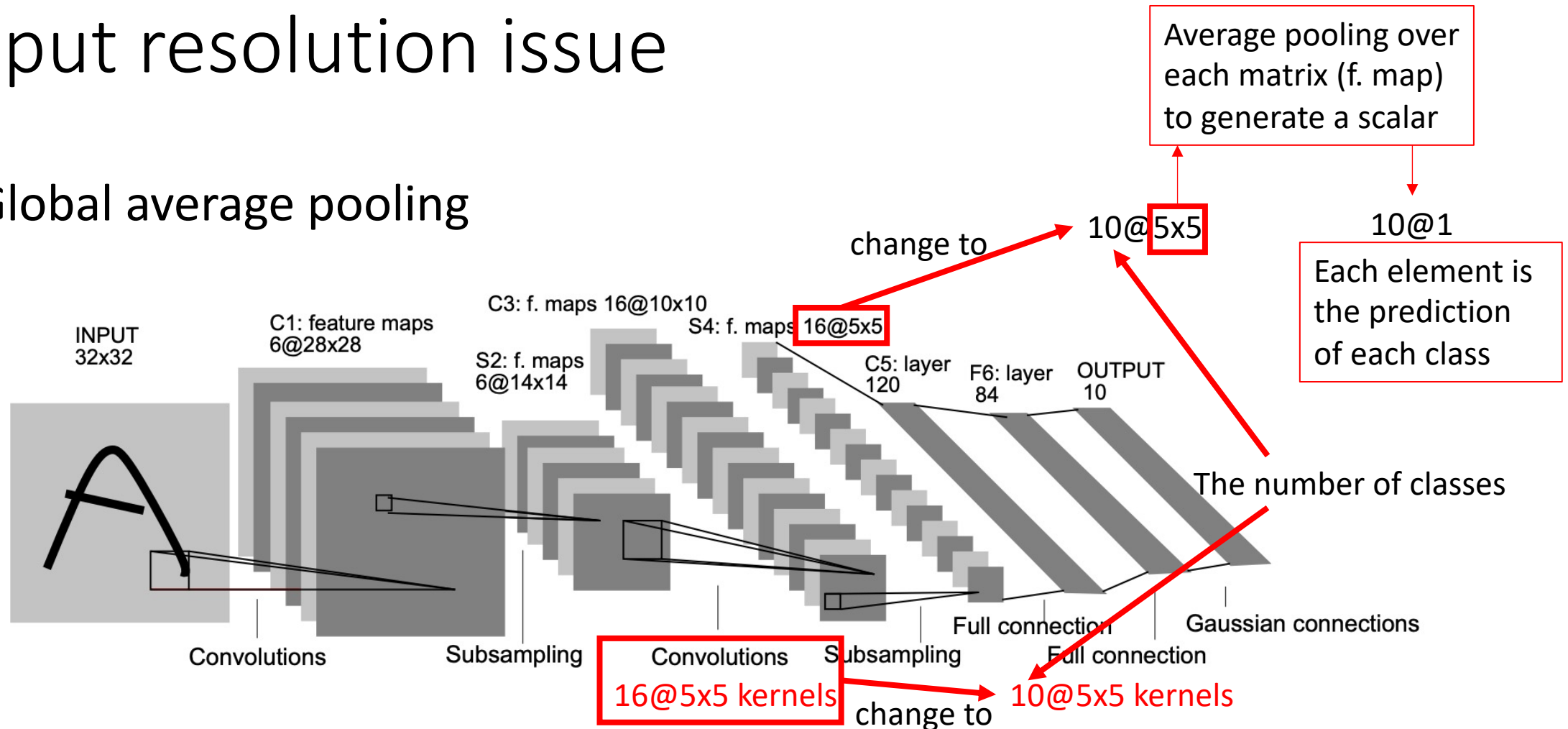


Fig. 1. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

References

- [Alexnet] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* 25 (2012): 1097-1105. Conference proceeding version at <https://papers.nips.cc/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-Abstract.html> or <https://papers.nips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf> (Section 3.5)
- [pyramid] He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Spatial pyramid pooling in deep convolutional networks for visual recognition." *IEEE transactions on pattern analysis and machine intelligence* 37, no. 9 (2015): 1904-1916. ArXiv version at <https://arxiv.org/abs/1406.4729> (Section 2.2)
- [NIN] Lin, Min, Qiang Chen, and Shuicheng Yan. "Network in network." *arXiv preprint arXiv:1312.4400* (2013). ArXiv version at <https://arxiv.org/abs/1312.4400> (Section 3.2)