

Basics on DNNs and Kernel Computation

Keypoints in Lecture 1

- Technology related issues
 - End of Moore's law, Dennard Scaling, ...
- Turing Tariff



- Domain Specific Architectures
 - ASIP, GPU, FPGA, ASIC

Workshop Agenda

- Lecture 1: Domain Specific Architectures
- **Lecture 2: Kernel computation**
- Lecture 3: Data-flow techniques
- Lecture 4: DNN accelerators architectures

Agenda – Lecture 2

- Deep Neural Networks
- Types of layers
- Memory and communication traffic
- Kernel computation
- Quantitative analysis

Artificial Intelligence

"The
John

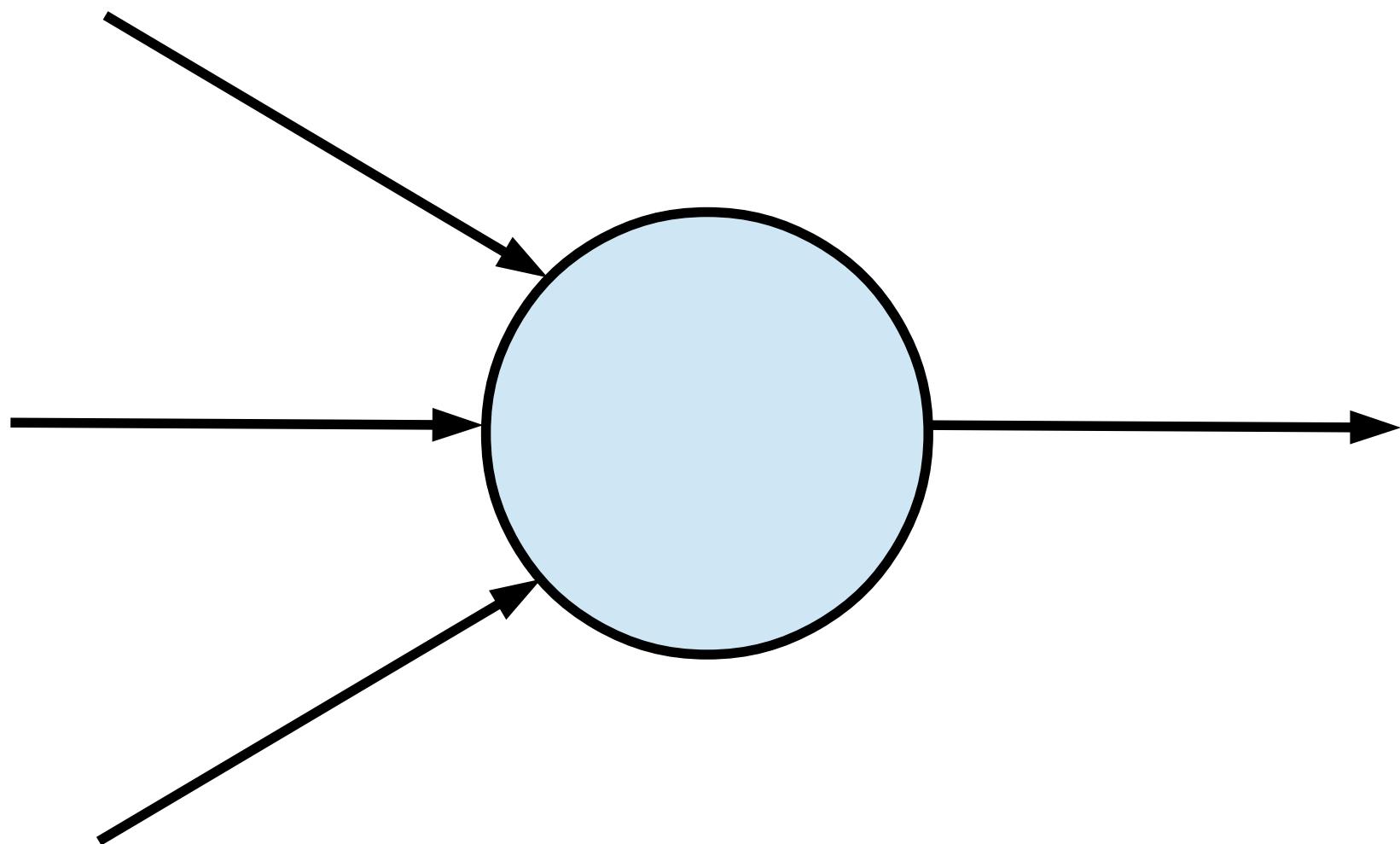
Machine Learning

"Fie
bei
Art

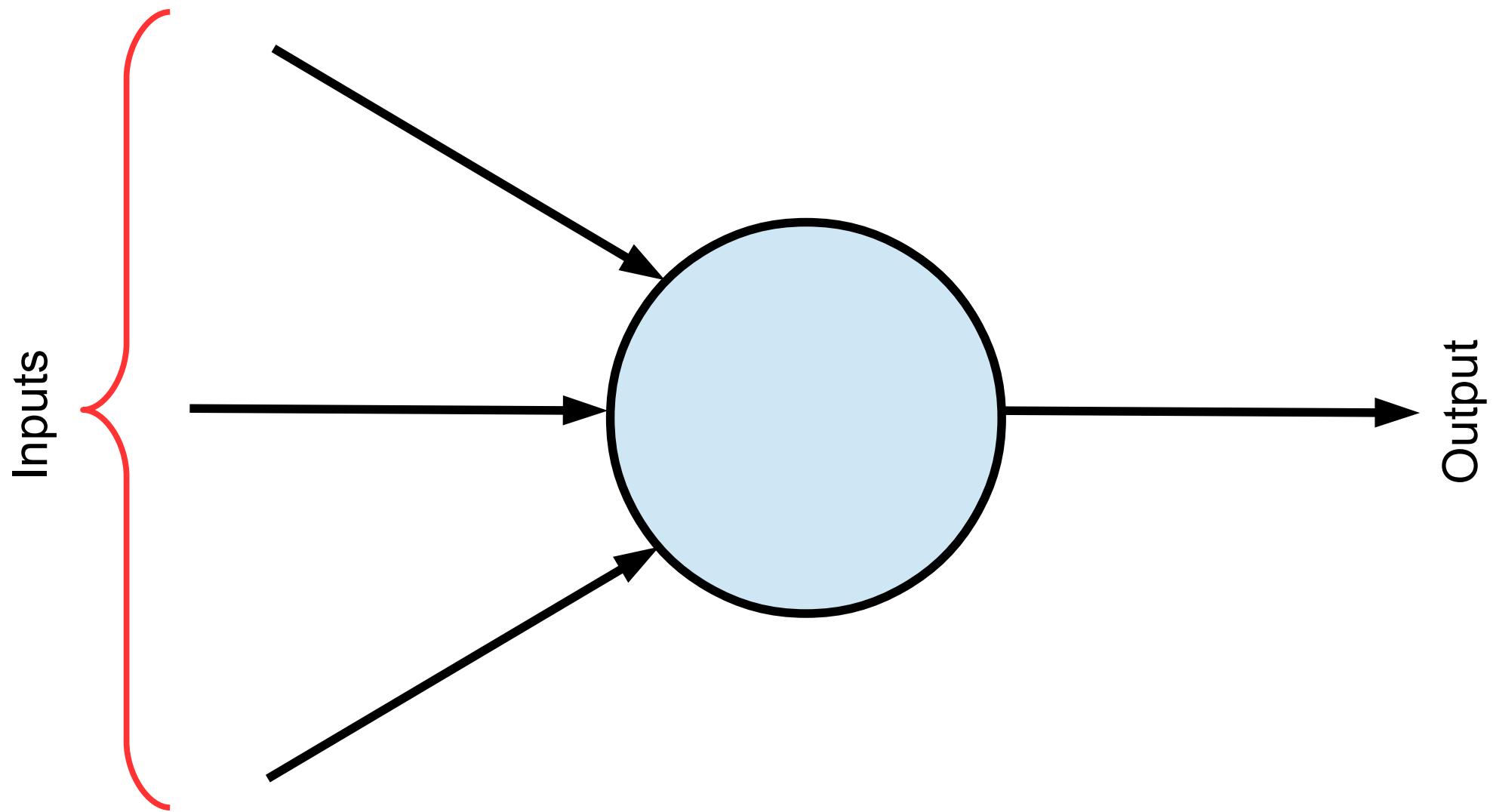
Brain-Inspired
*An algorithmic
our underst*

**Neural
Networks**

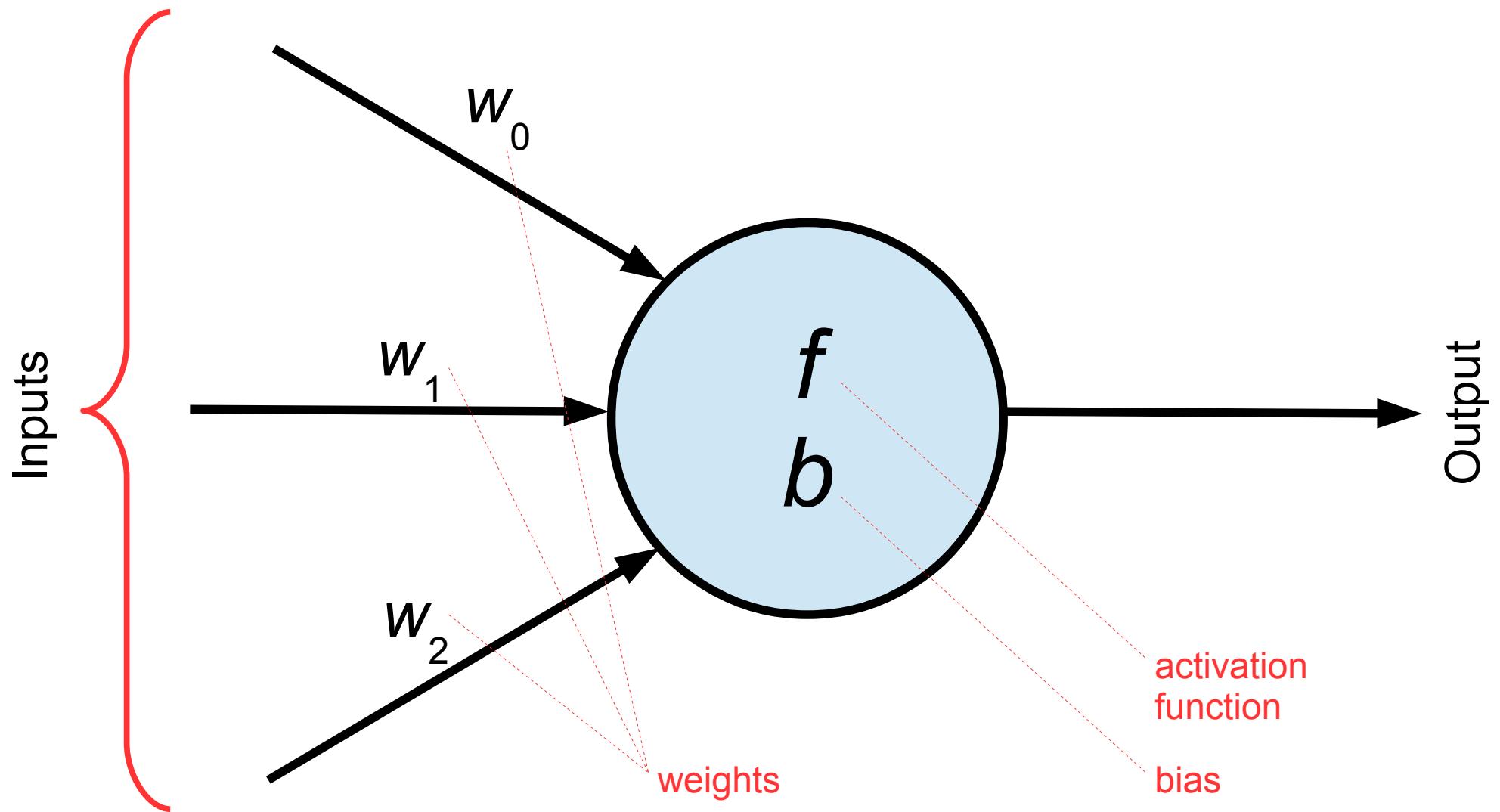
Neuron



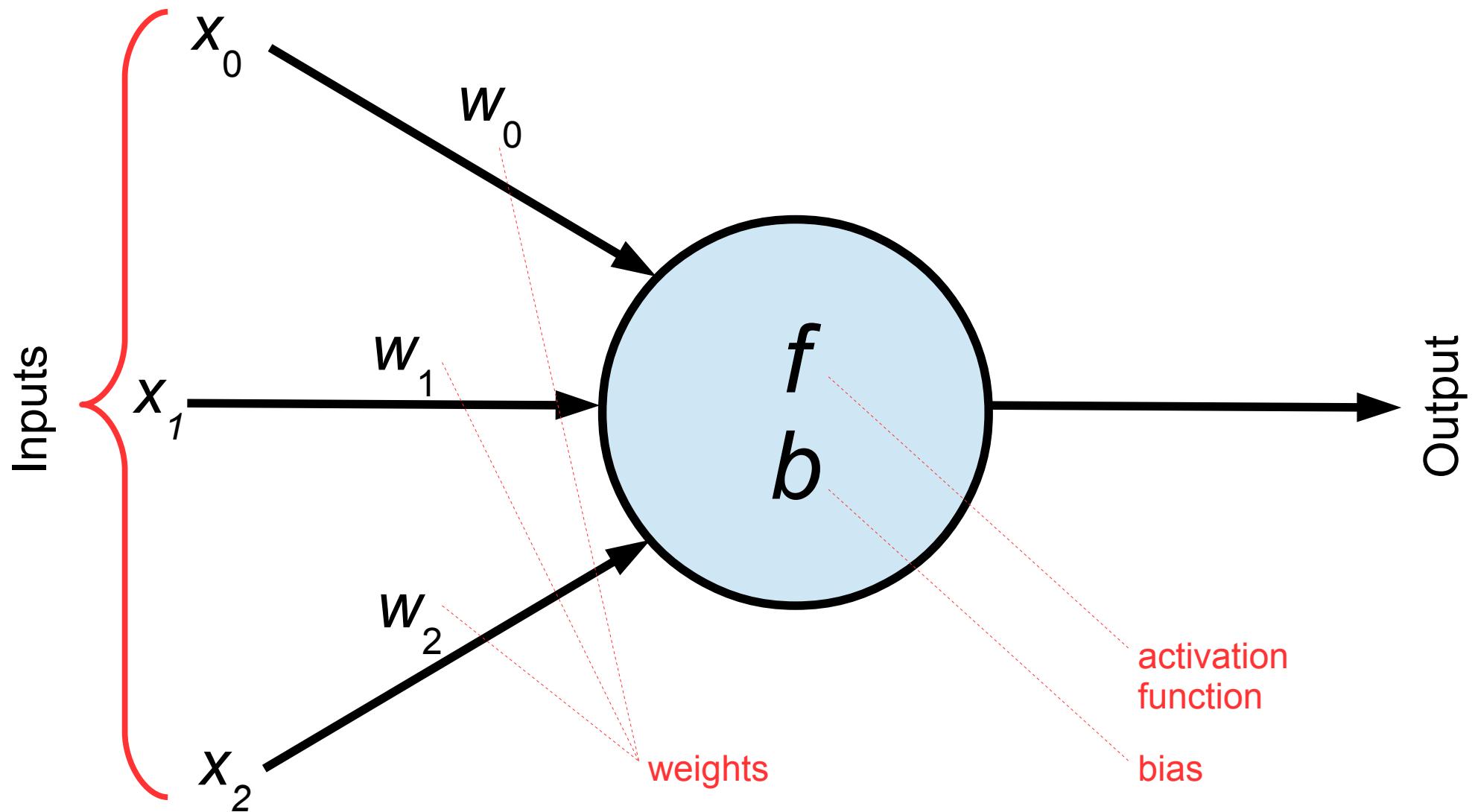
Neuron



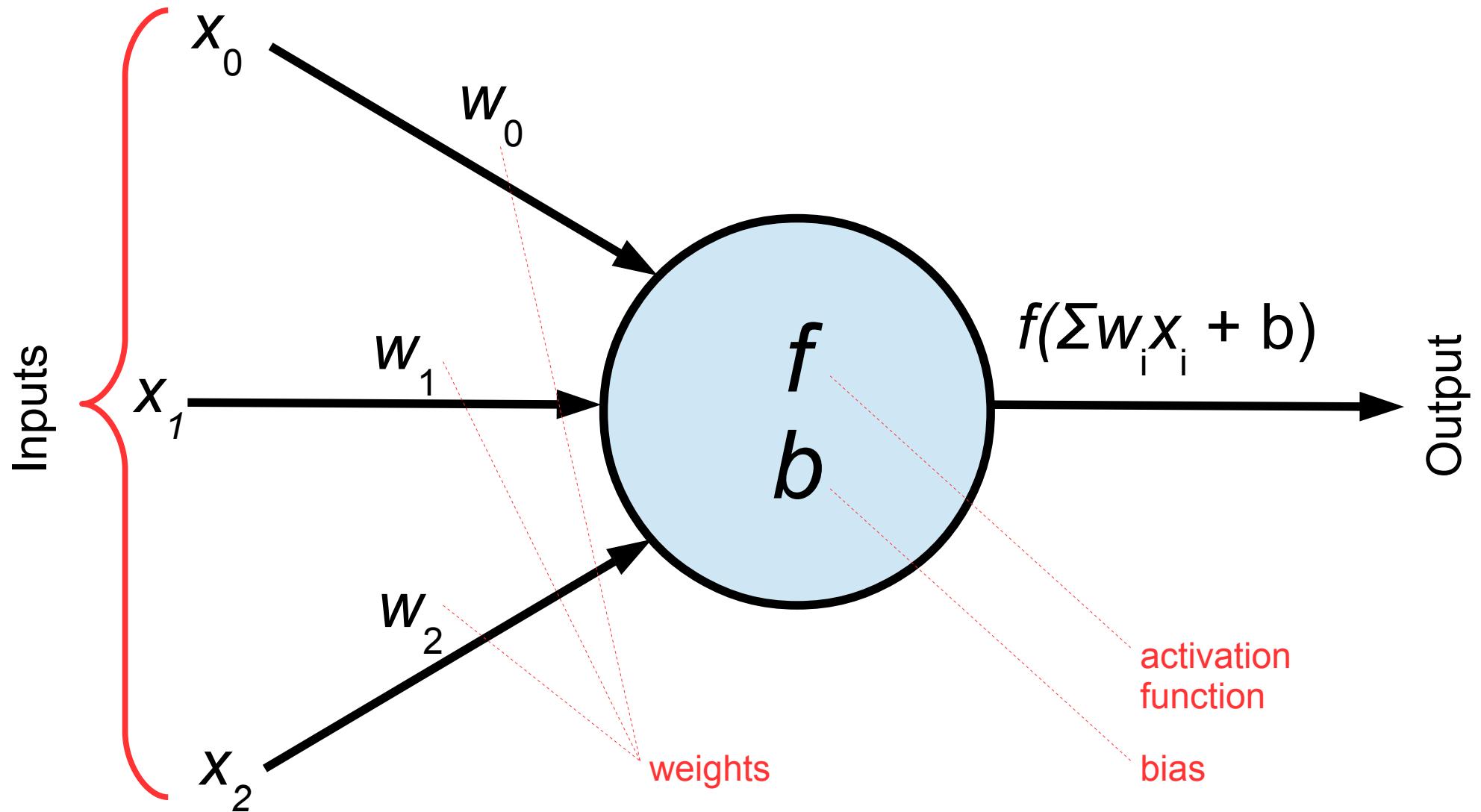
Neuron



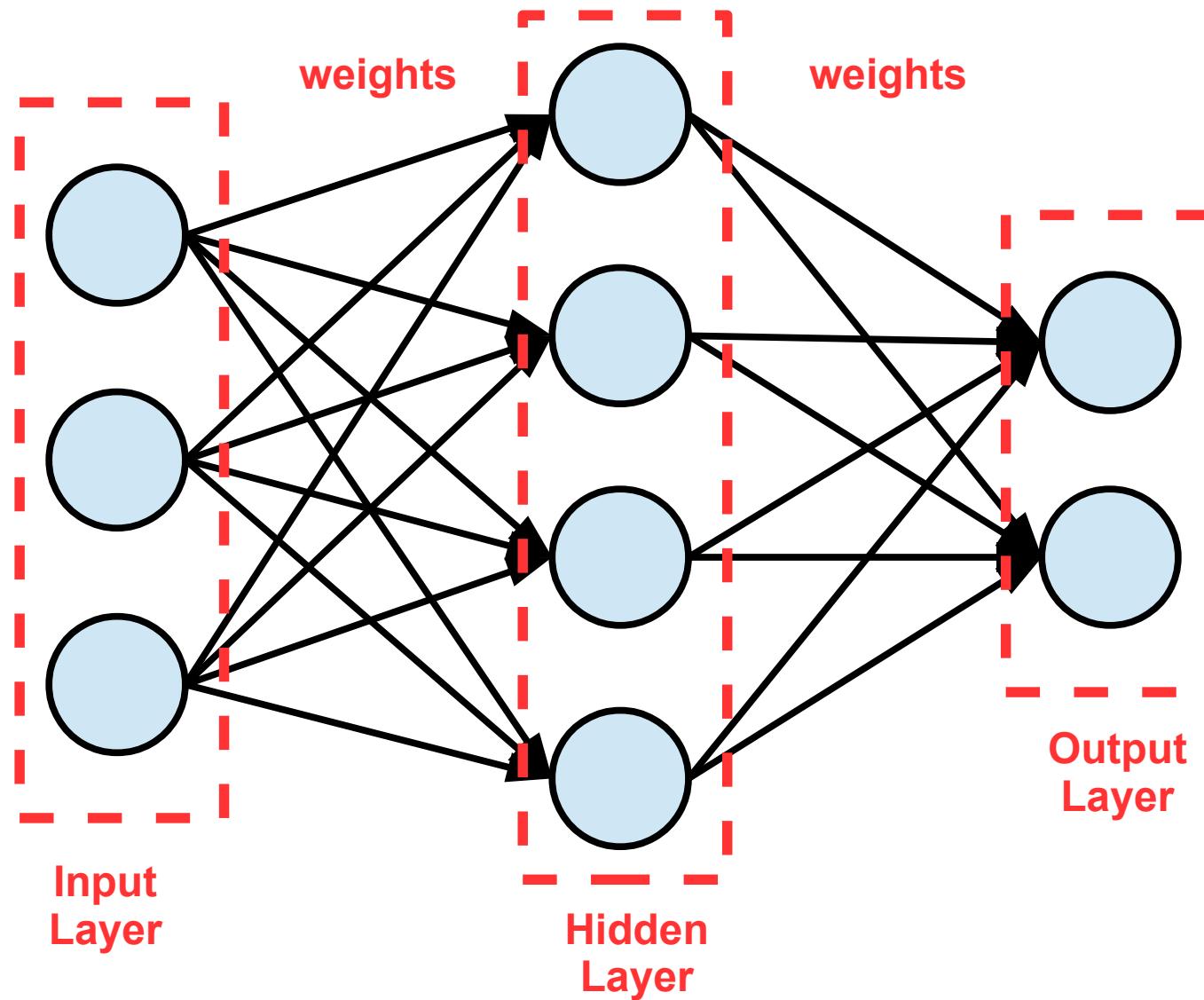
Neuron



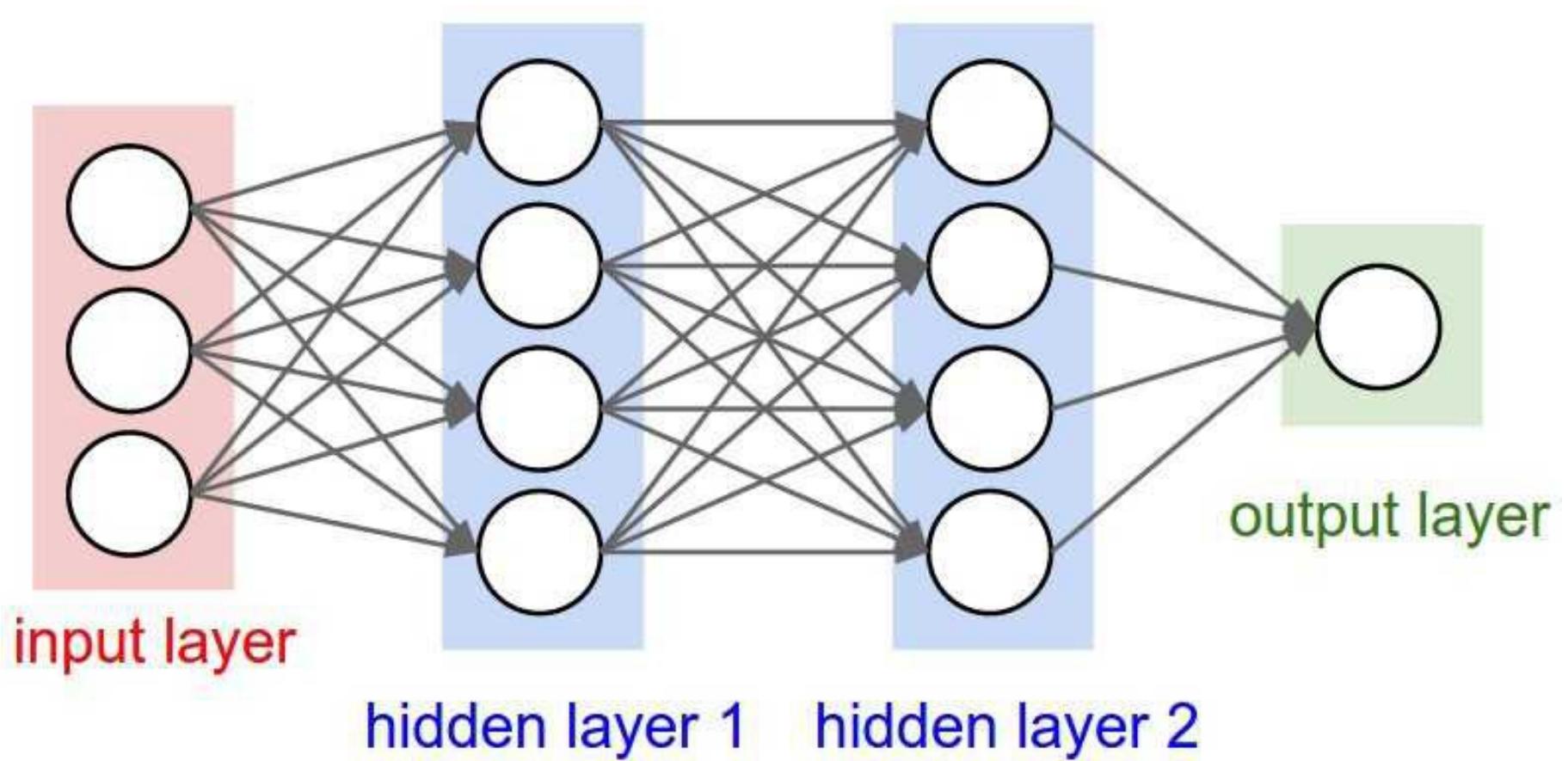
Neuron



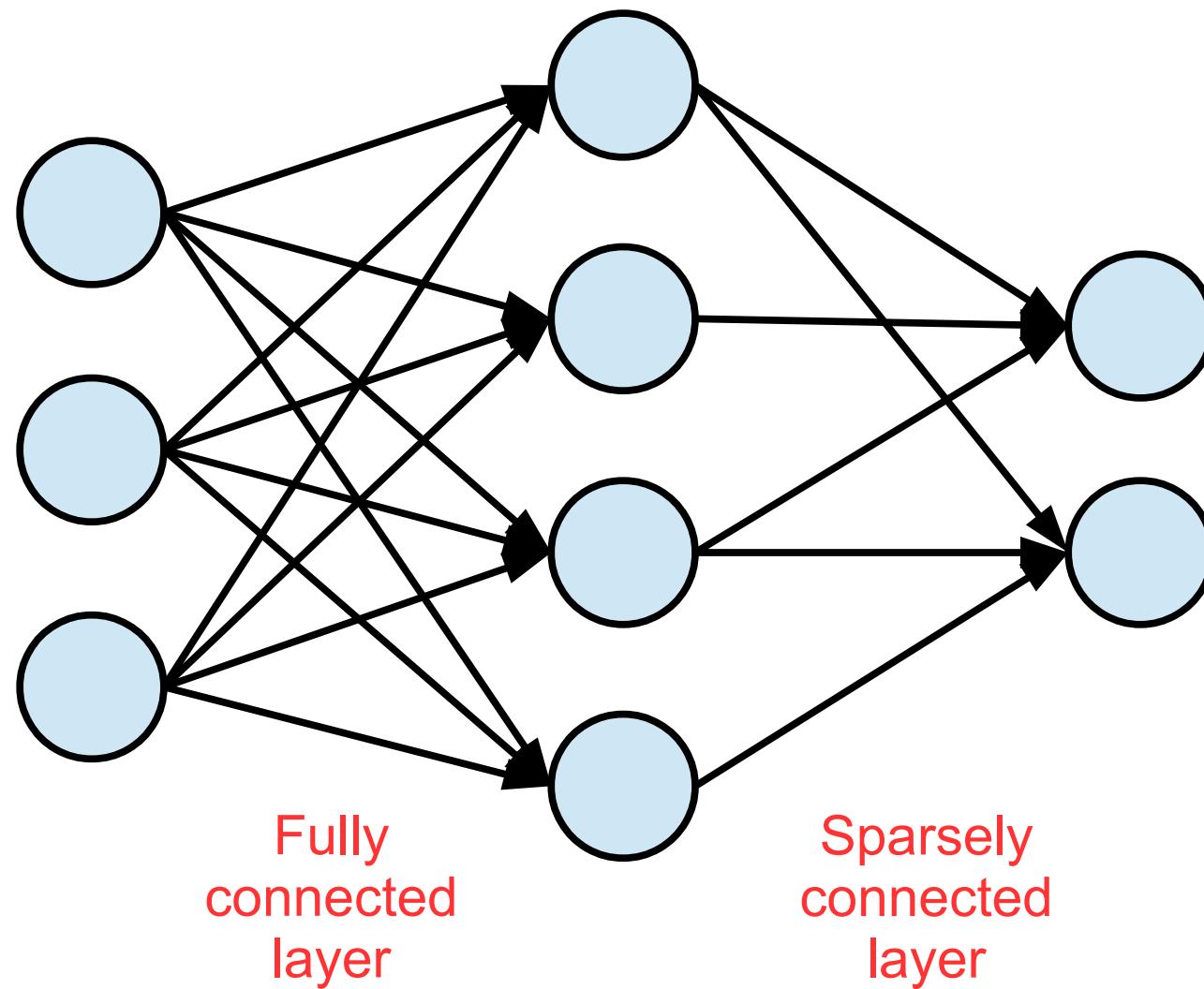
Terminology



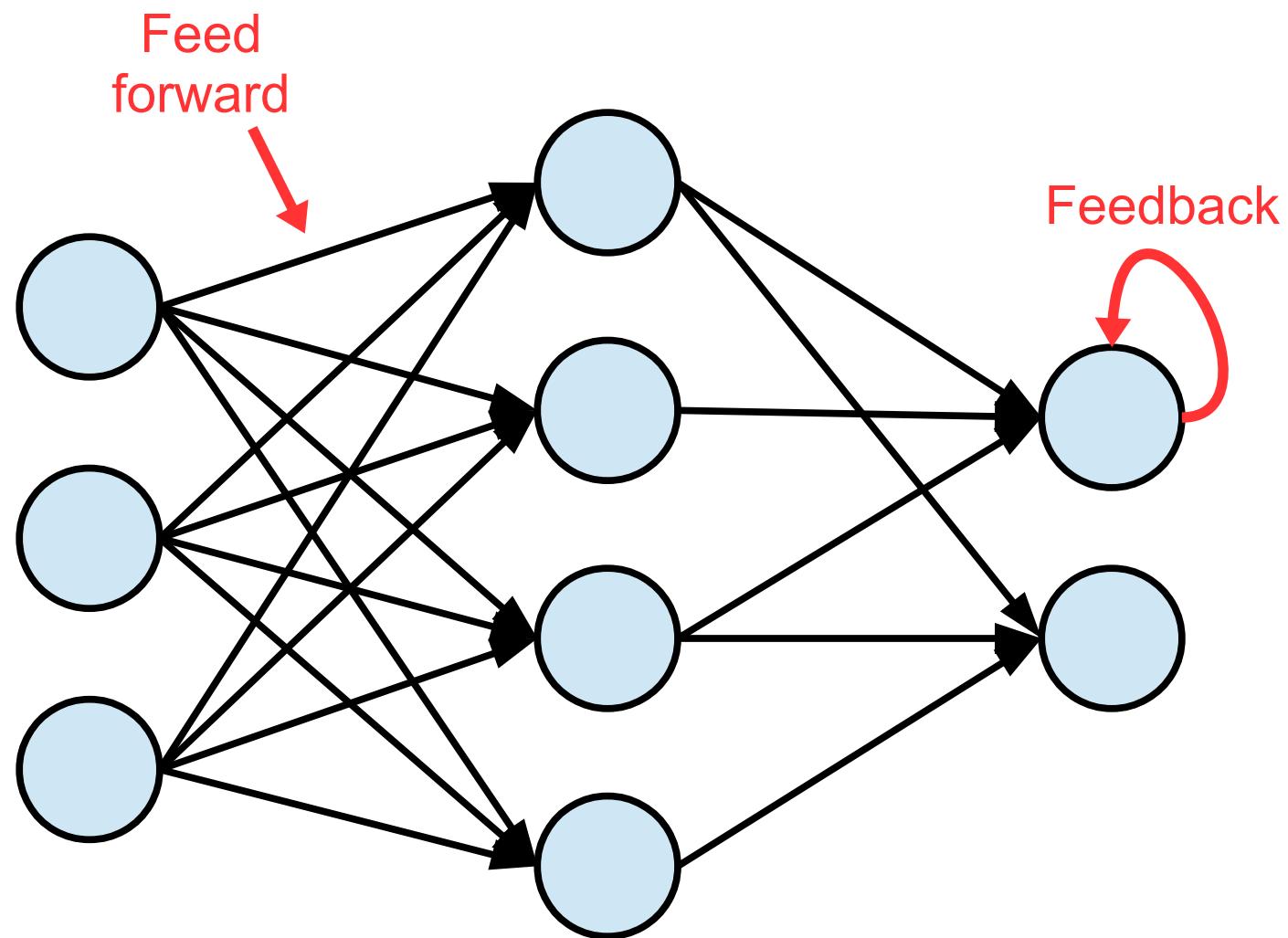
Terminology



Terminology



Terminology



Artificial Intelligence

"The
John

Machine Learning

"Fie
bei
Art

Brain-Inspired

An algorithm
our underst

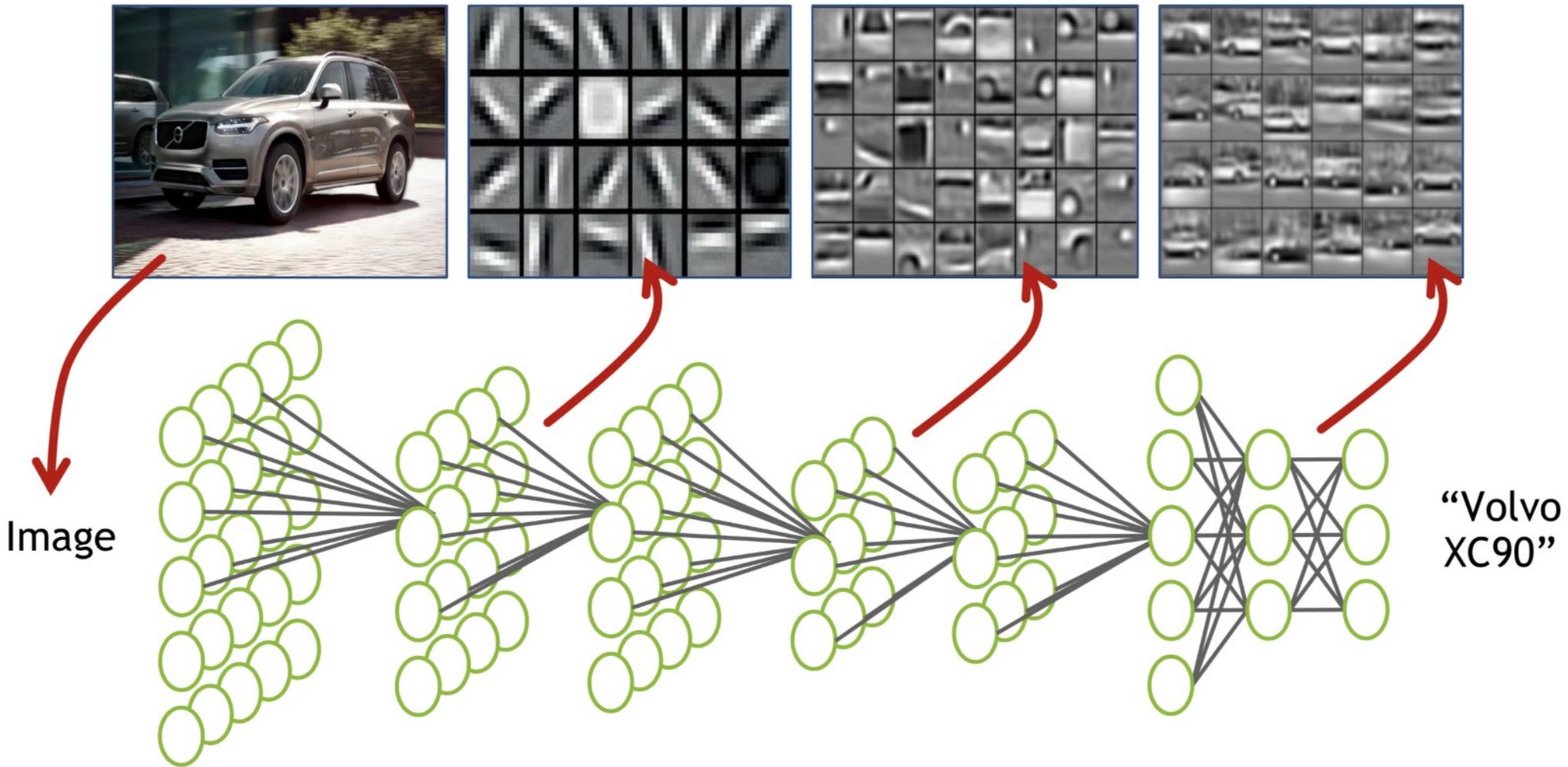
Neural Networks

Deep Learning

hout

rom

Deep Neural Network



[H. Lee *et al.*, Unsupervised learning of hierarchical representations with convolutional deep belief networks. Commun. ACM 2011]

Types of DNNs

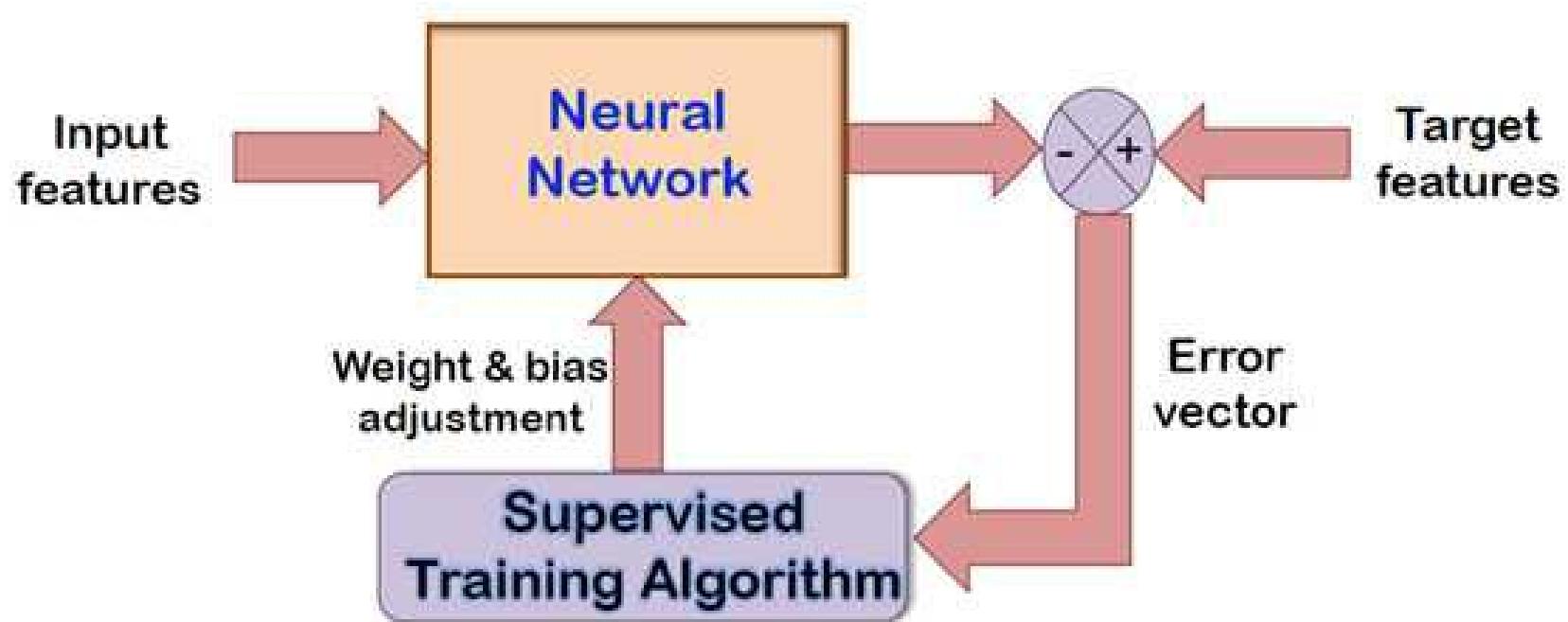
- Fully-Connected NN
 - feed forward, a.k.a. multilayer perceptron (MLP)
- Convolutional NN (CNN)
 - feed forward, sparsely-connected w/ weight sharing
- Recurrent NN (RNN)
 - feedback
- Long Short-Term Memory (LSTM)
 - feedback + storage

The Two Phases

- Training
- Inference

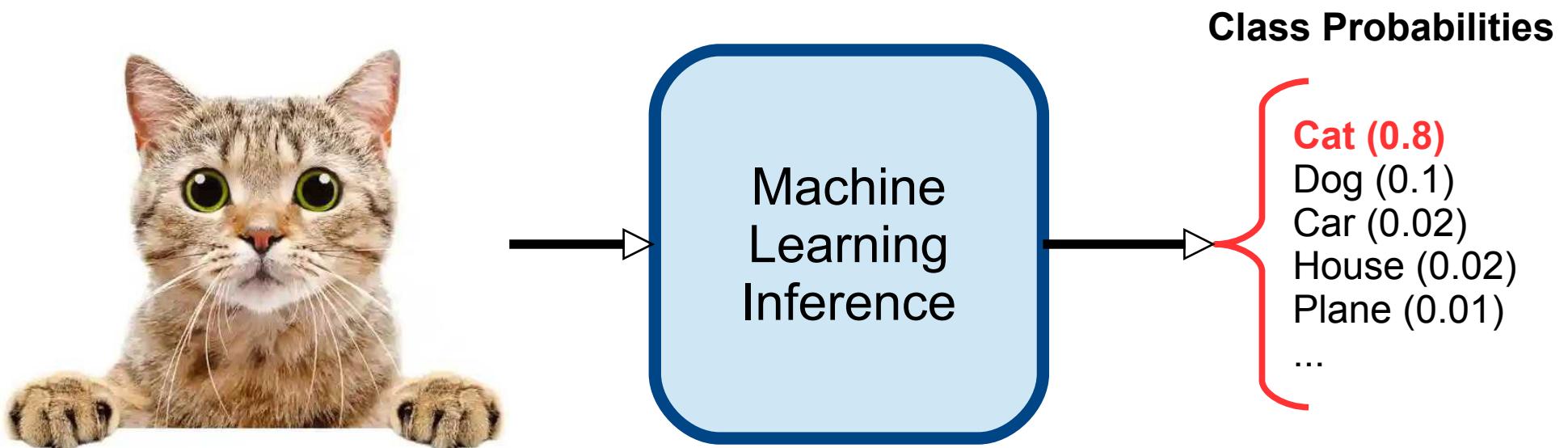
Training

- Determine weights (and biases)
 - Supervised
 - Training set has inputs and outputs, *i.e.*, labeled



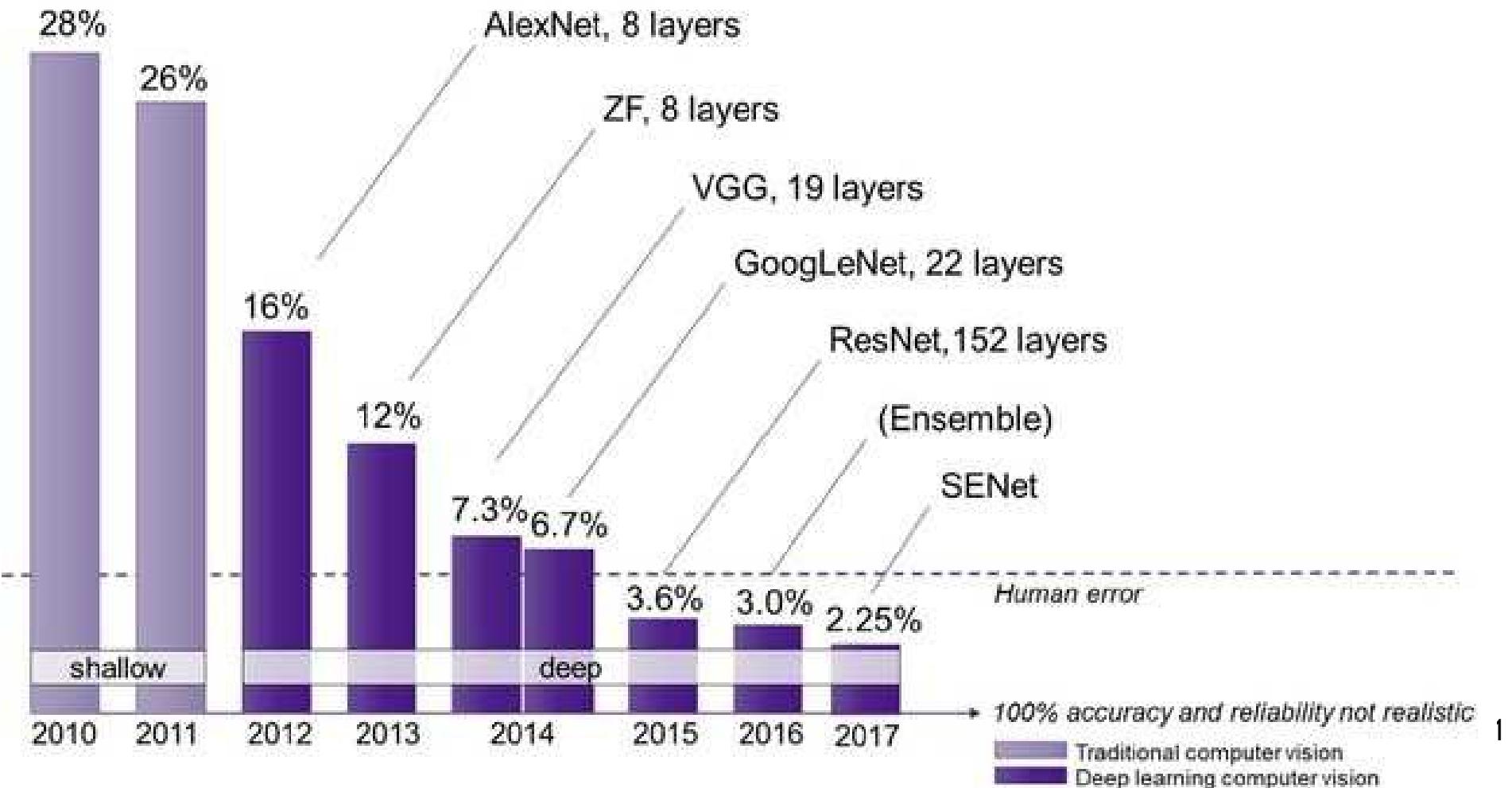
Inference

- Apply weights to determine output



ImageNet Challenge

- Training 1.2M images, 1K object categories



Applications of DNNs

- Image and Video
 - 70% to 80% of internet traffic is video (Cisco Visual Networking Index 2021)
 - 800 Mhours/day video surveillance in 2016
 - 7x increase in 2021 (VNI Complete Forecast Highlights)
 - DNN based CV to extract meaningful information
 - Image classification
 - Object localization and detection
 - Image segmentation
 - Action recognition

Applications of DNNs

- Speech and Language
 - Speech recognition
 - Machine translation
 - Natural language processing
 - Audio generation

Applications of DNNs

- Medicine and Health Care
- Game Play
- Robotics
- ...

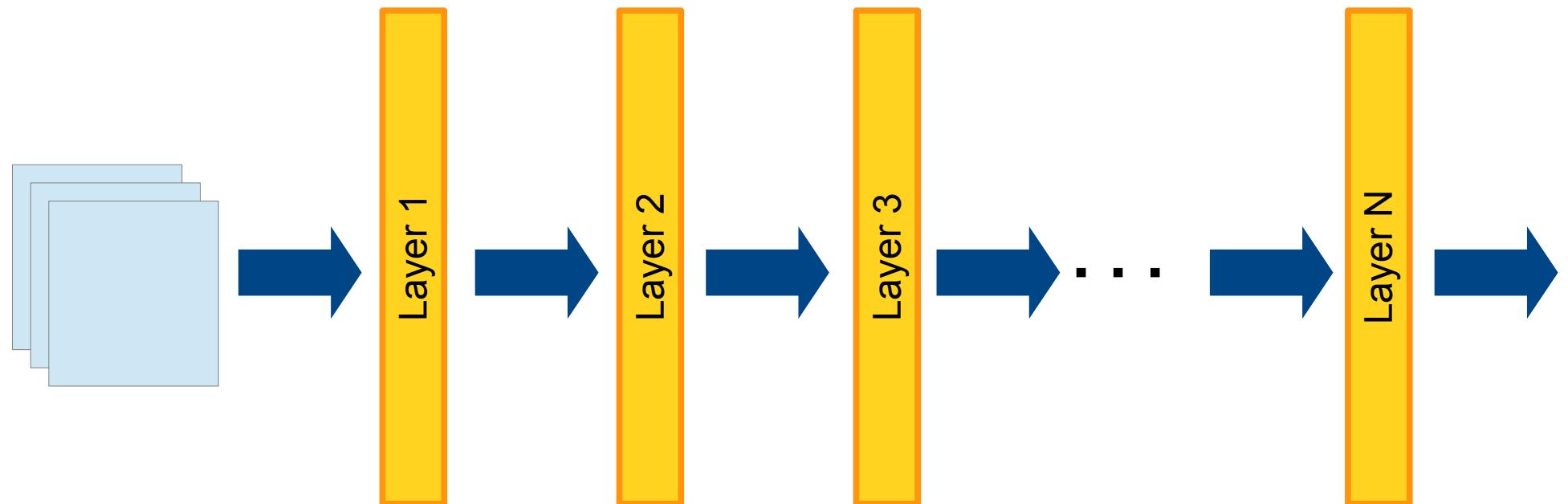
Embedded vs. Cloud

- Training
 - Typical performed in the cloud
- Inference
 - Trend is moving to the edge
 - Privacy/security
 - Latency
 - Communication energy

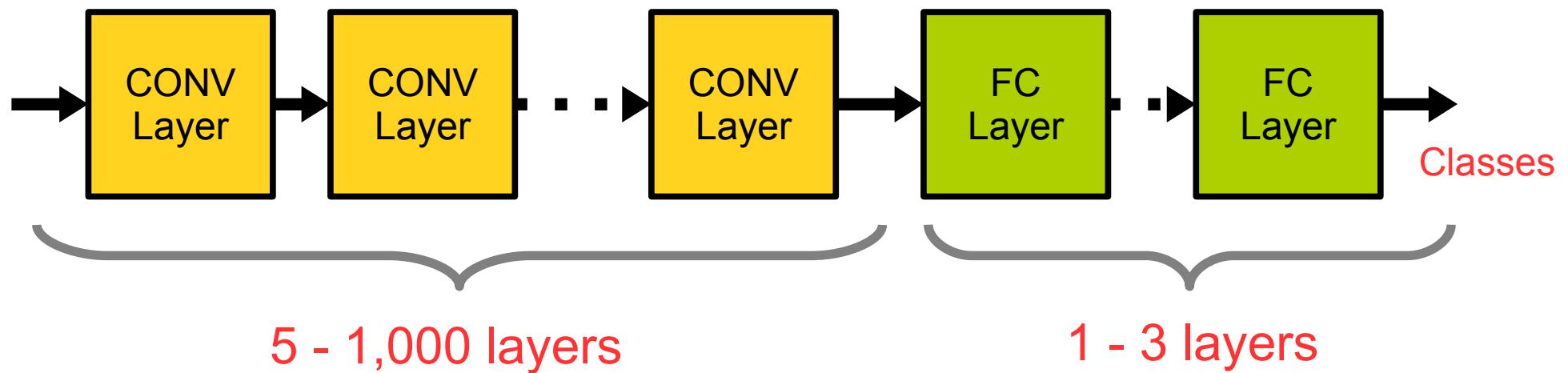
Agenda

- Deep Neural Networks
- Types of layers
- Memory and communication traffic
- Kernel computation
- Quantitative analysis

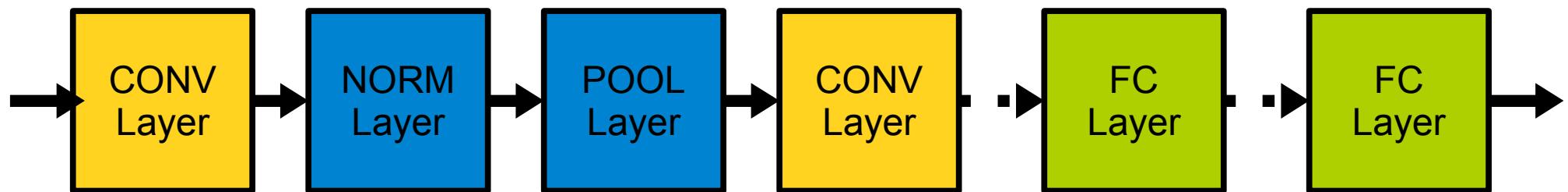
Deep Neural Network



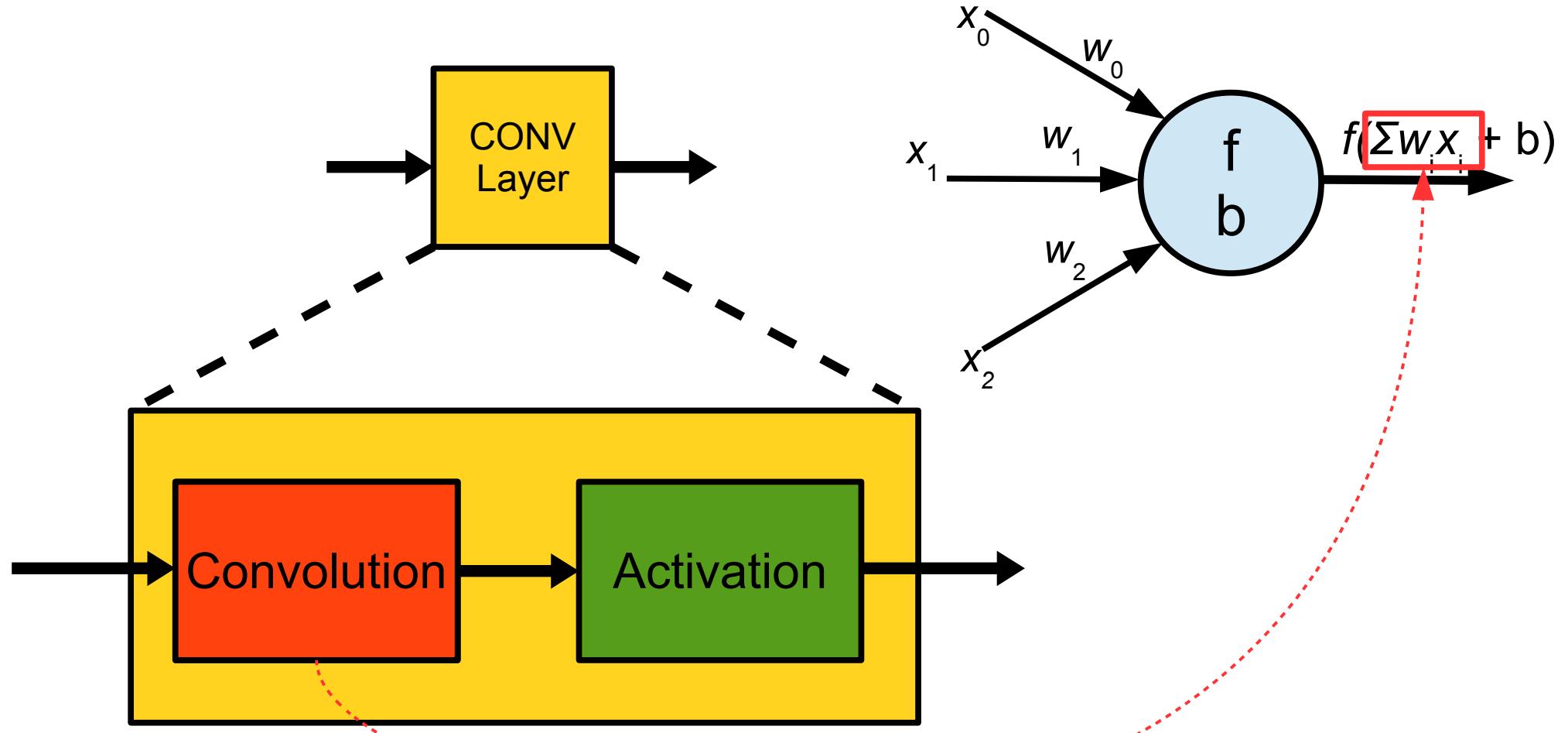
Deep Convolutional Neural Networks



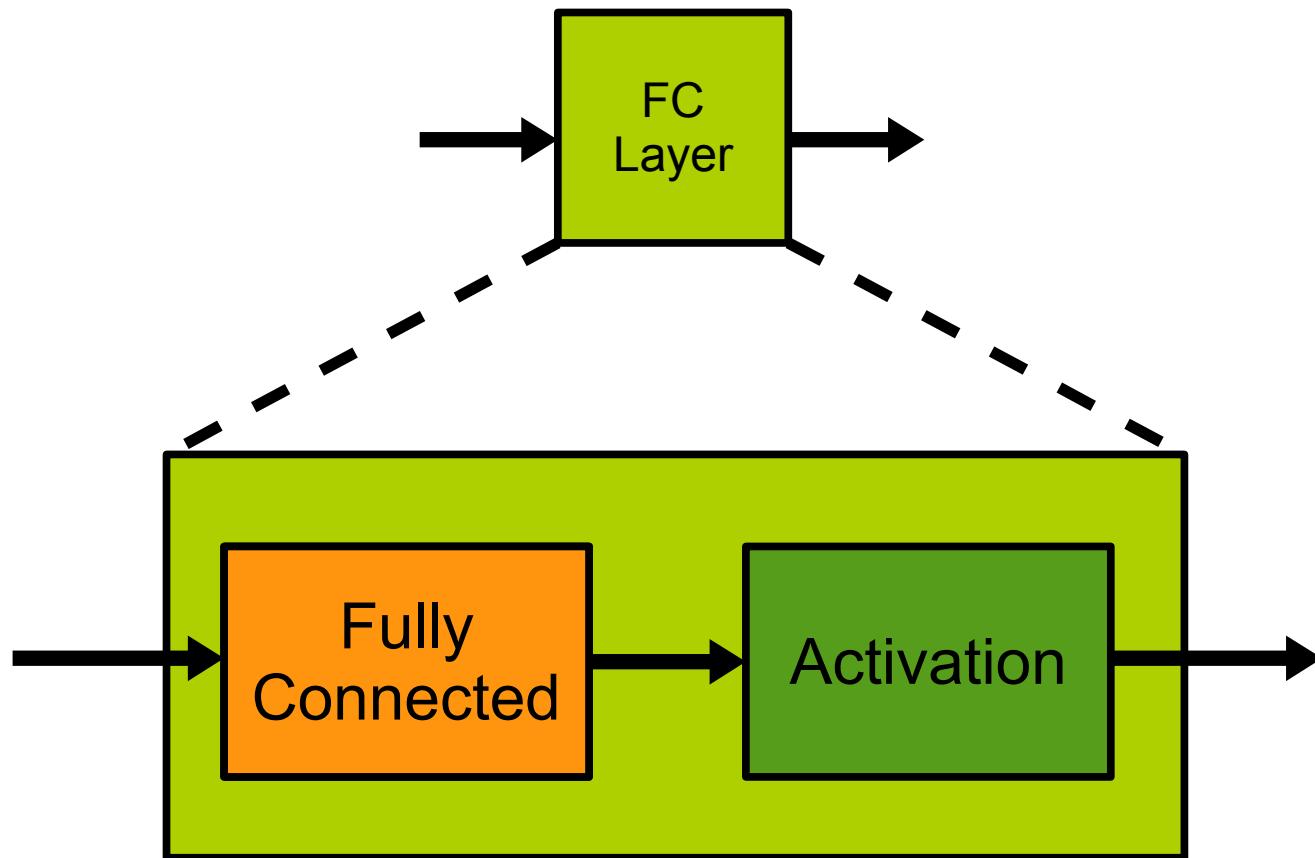
Deep Convolutional Neural Networks – Optional Layers



Deep Convolutional Neural Networks

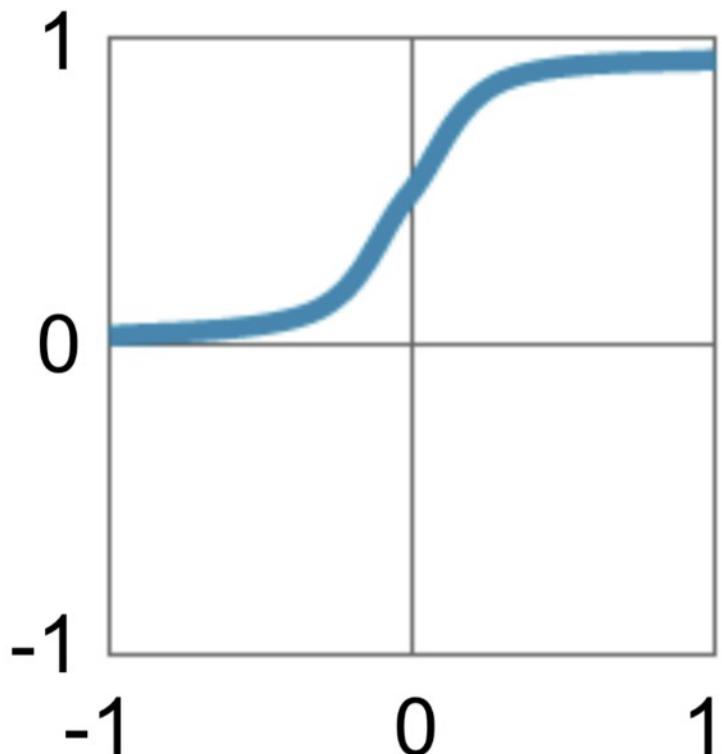


Deep Convolutional Neural Networks



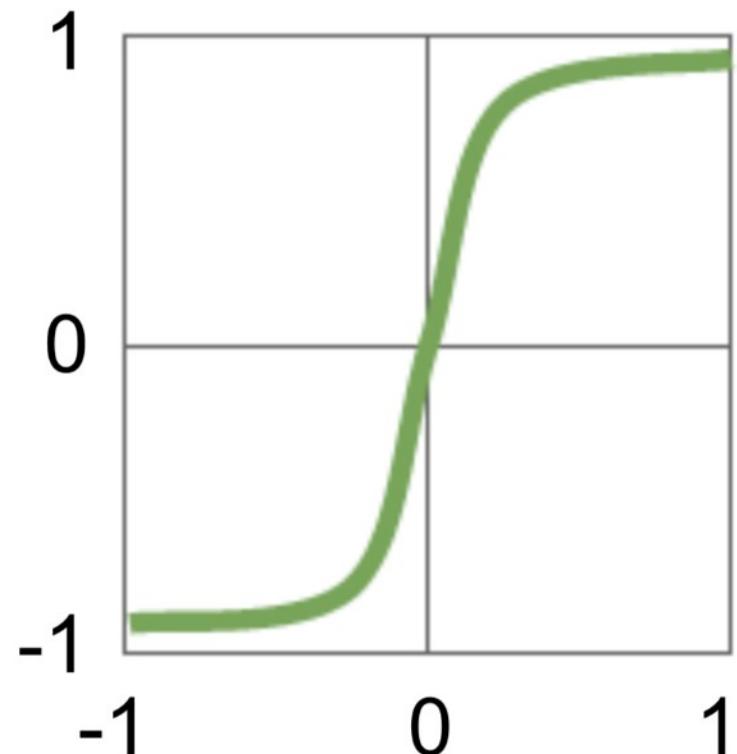
Activation -- Traditional

Sigmoid



$$y = 1/(1+e^{-x})$$

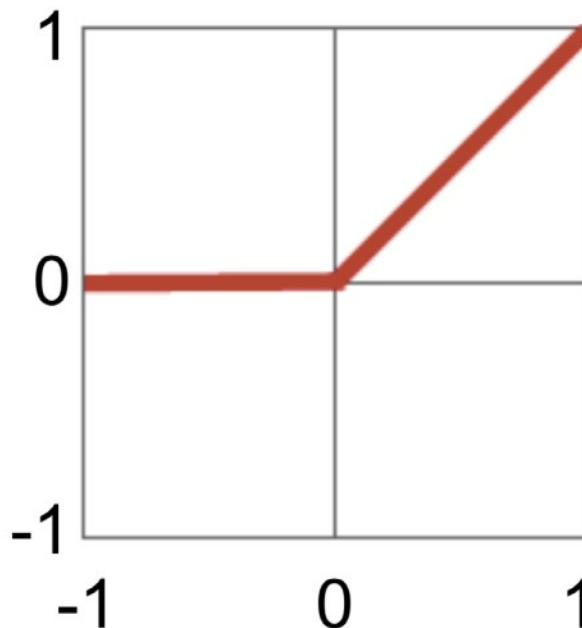
Hyperbolic Tangent



$$y = (e^x - e^{-x}) / (e^x + e^{-x})$$

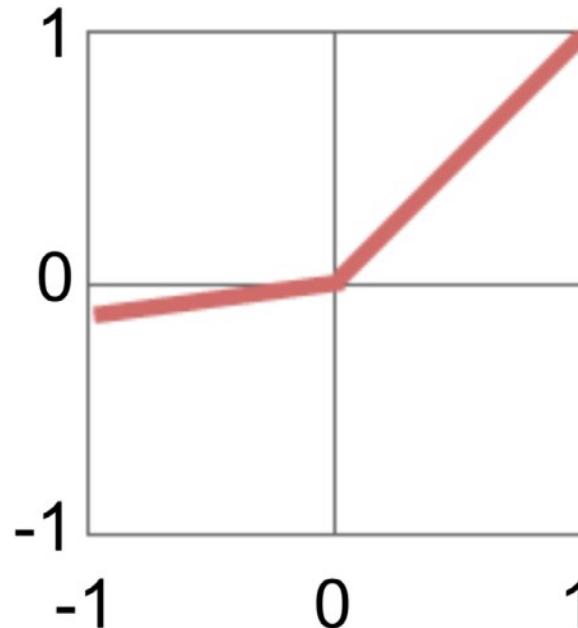
Activation Functions -- Modern

Rectified Linear Unit
(ReLU)



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

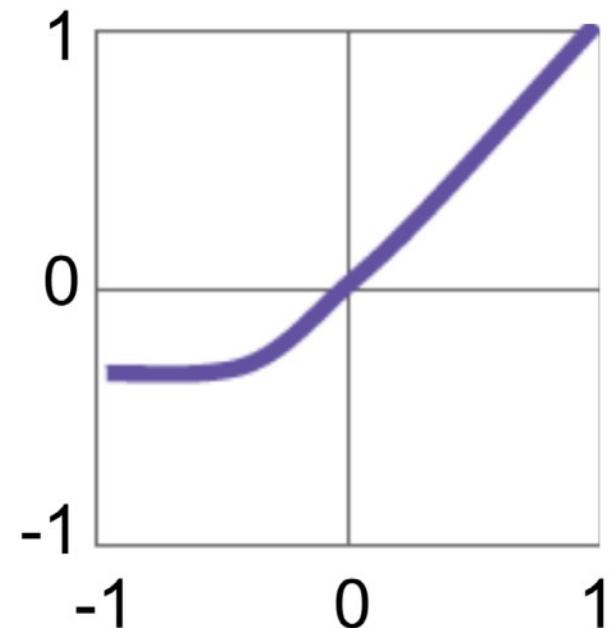
Leaky ReLU



$$y = \max(\alpha x, x)$$

$\alpha = \text{small const. (e.g. 0.1)}$

Exponential LU

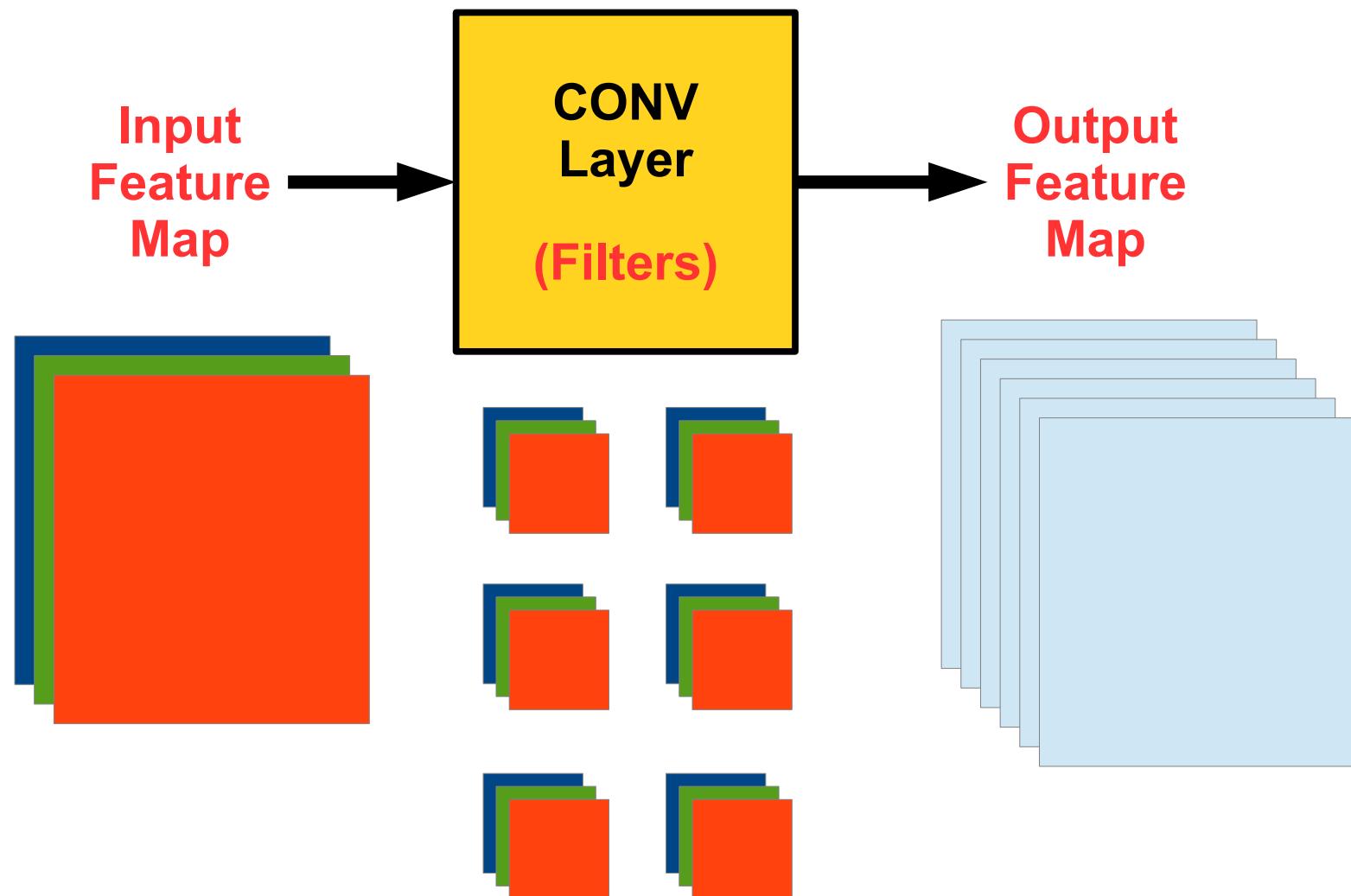


$$y = \begin{cases} x, & x \geq 0 \\ \alpha(e^x - 1), & x < 0 \end{cases}$$

Convolution

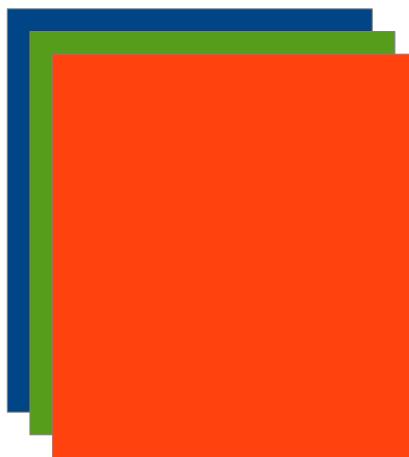
- Convolutions account for more than 90% of overall computation
- Dominate runtime and energy consumption

Convolution

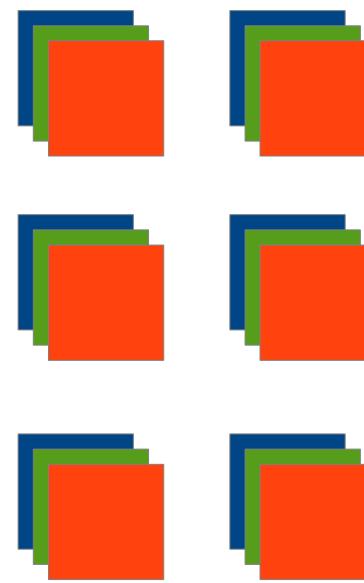


Convolution

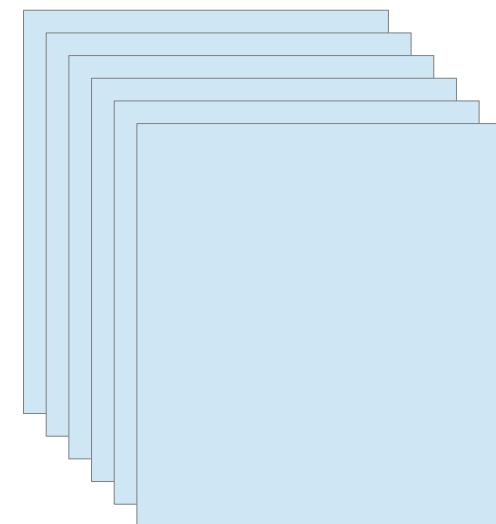
Input Feature Map



Filters



Output Feature Map

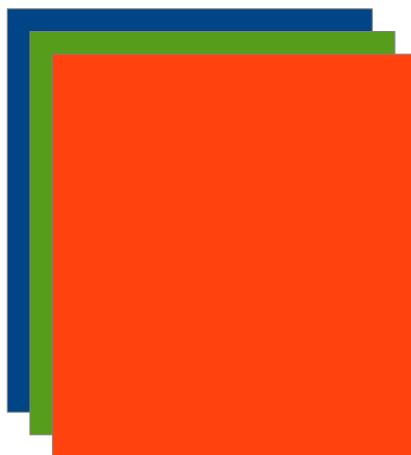


$W \times H \times 6$

$W_f \times H_f \times 3$

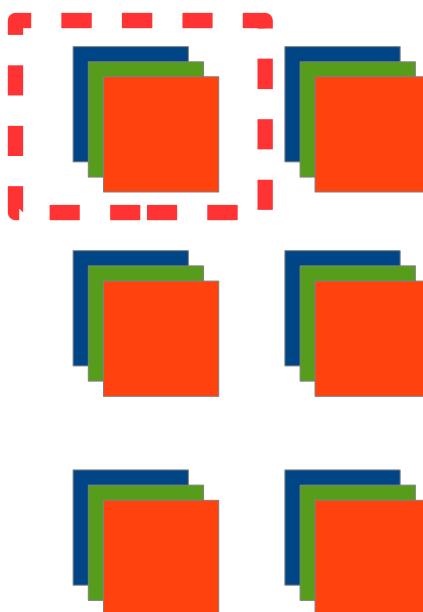
Convolution

Input Feature Map



$W \times H \times 3$

Filters

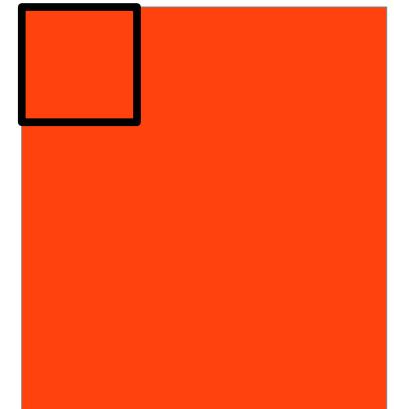
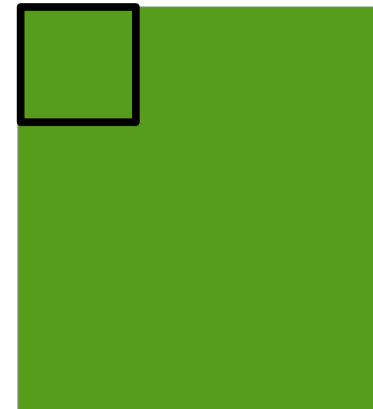
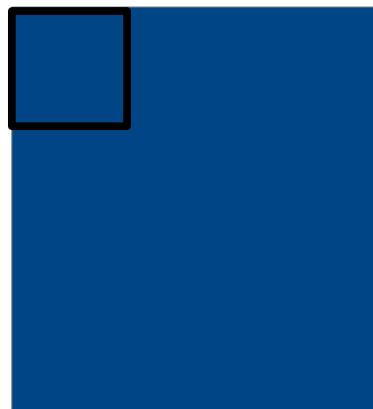


$W_f \times H_f \times 3$

Convolution

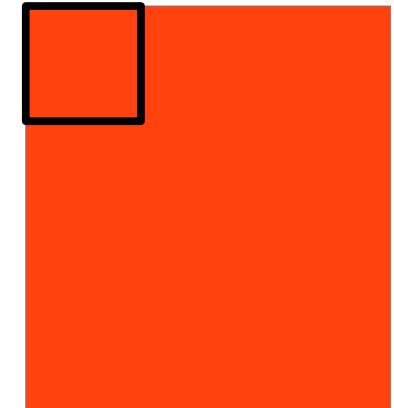
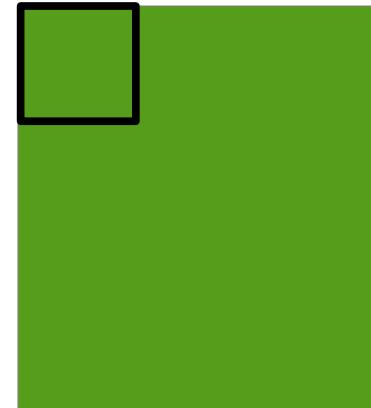
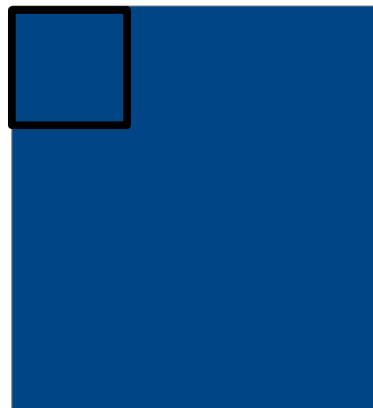


Convolution

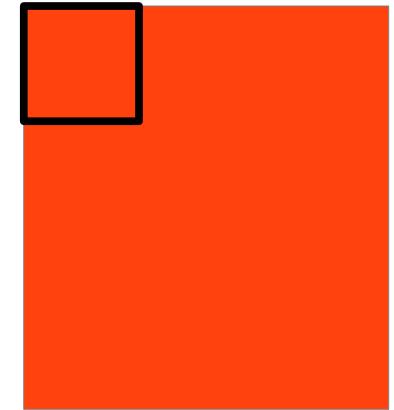
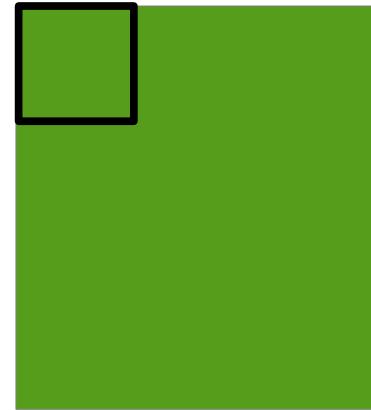
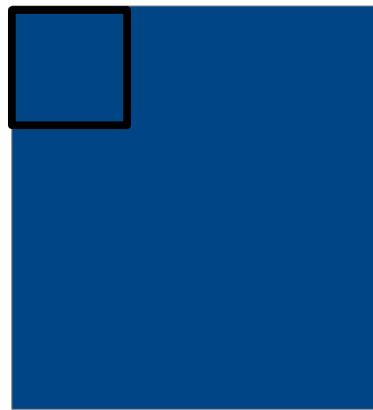


$$\begin{aligned} & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \end{aligned}$$

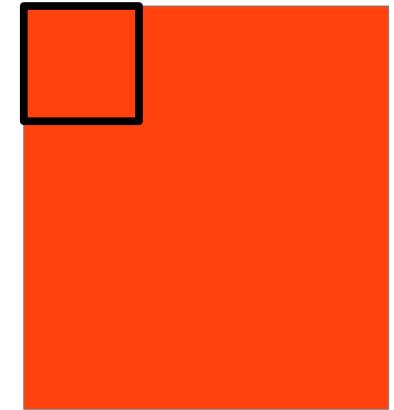
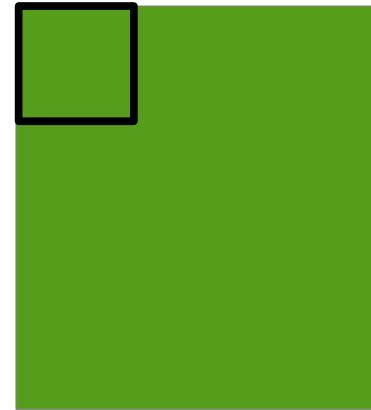
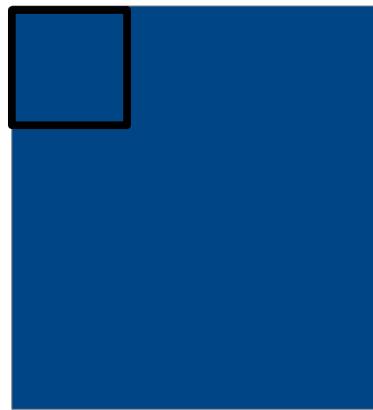
Convolution


$$\begin{aligned} & f1[0,0,\textcolor{blue}{\blacksquare}] * I[0,0,\textcolor{blue}{\blacksquare}] + f1[1,0,\textcolor{blue}{\blacksquare}] * I[1,0,\textcolor{blue}{\blacksquare}] + \dots + f1[Wf-1,0,\textcolor{blue}{\blacksquare}] * I[Wf-1,0,\textcolor{blue}{\blacksquare}] + \\ & f1[0,1,\textcolor{blue}{\blacksquare}] * I[0,1,\textcolor{blue}{\blacksquare}] + f1[1,1,\textcolor{blue}{\blacksquare}] * I[1,1,\textcolor{blue}{\blacksquare}] + \dots + f1[Wf-1,1,\textcolor{blue}{\blacksquare}] * I[Wf-1,1,\textcolor{blue}{\blacksquare}] + \\ & + \dots + \\ & f1[0,Hf-1,\textcolor{blue}{\blacksquare}] * I[0,Hf-1,\textcolor{blue}{\blacksquare}] + f1[1,Hf-1,\textcolor{blue}{\blacksquare}] * I[1,Hf-1,\textcolor{blue}{\blacksquare}] + \dots + f1[Wf-1,Hf-1,\textcolor{blue}{\blacksquare}] * I[Wf-1,Hf-1,\textcolor{blue}{\blacksquare}] + \\ & f1[0,0,\textcolor{green}{\blacksquare}] * I[0,0,\textcolor{green}{\blacksquare}] + f1[1,0,\textcolor{green}{\blacksquare}] * I[1,0,\textcolor{green}{\blacksquare}] + \dots + f1[Wf-1,0,\textcolor{green}{\blacksquare}] * I[Wf-1,0,\textcolor{green}{\blacksquare}] + \\ & f1[0,1,\textcolor{green}{\blacksquare}] * I[0,1,\textcolor{green}{\blacksquare}] + f1[1,1,\textcolor{green}{\blacksquare}] * I[1,1,\textcolor{green}{\blacksquare}] + \dots + f1[Wf-1,1,\textcolor{green}{\blacksquare}] * I[Wf-1,1,\textcolor{green}{\blacksquare}] + \\ & + \dots + \\ & f1[0,Hf-1,\textcolor{green}{\blacksquare}] * I[0,Hf-1,\textcolor{green}{\blacksquare}] + f1[1,Hf-1,\textcolor{green}{\blacksquare}] * I[1,Hf-1,\textcolor{green}{\blacksquare}] + \dots + f1[Wf-1,Hf-1,\textcolor{green}{\blacksquare}] * I[Wf-1,Hf-1,\textcolor{green}{\blacksquare}] + \end{aligned}$$

Convolution


$$\begin{aligned} & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \\ & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \\ & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \end{aligned}$$

Convolution


$$\begin{aligned} & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \\ & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \\ & f1[0,0,\blacksquare] * I[0,0,\blacksquare] + f1[1,0,\blacksquare] * I[1,0,\blacksquare] + \dots + f1[Wf-1,0,\blacksquare] * I[Wf-1,0,\blacksquare] + \\ & f1[0,1,\blacksquare] * I[0,1,\blacksquare] + f1[1,1,\blacksquare] * I[1,1,\blacksquare] + \dots + f1[Wf-1,1,\blacksquare] * I[Wf-1,1,\blacksquare] + \\ & + \dots + \\ & f1[0,Hf-1,\blacksquare] * I[0,Hf-1,\blacksquare] + f1[1,Hf-1,\blacksquare] * I[1,Hf-1,\blacksquare] + \dots + f1[Wf-1,Hf-1,\blacksquare] * I[Wf-1,Hf-1,\blacksquare] + \text{Bias}^{4^2} \end{aligned}$$

Convolution

$f1[0,0,\text{blue}]*I[0,0,\text{blue}] + f1[1,0,\text{blue}]*I[1,0,\text{blue}] + \dots + f1[Wf-1,0,\text{blue}]*I[Wf-1,0,\text{blue}] +$
 $f1[0,1,\text{blue}]*I[0,1,\text{blue}] + f1[1,1,\text{blue}]*I[1,1,\text{blue}] + \dots + f1[Wf-1,1,\text{blue}]*I[Wf-1,1,\text{blue}] +$
+ ... +
 $f1[0,Hf-1,\text{blue}]*I[0,Hf-1,\text{blue}] + f1[1,Hf-1,\text{blue}]*I[1,Hf-1,\text{blue}] + \dots + f1[Wf-1,Hf-1,\text{blue}]*I[Wf-1,Hf-1,\text{blue}] +$
 $f1[0,0,\text{green}]*I[0,0,\text{green}] + f1[1,0,\text{green}]*I[1,0,\text{green}] + \dots + f1[Wf-1,0,\text{green}]*I[Wf-1,0,\text{green}] +$
 $f1[0,1,\text{green}]*I[0,1,\text{green}] + f1[1,1,\text{green}]*I[1,1,\text{green}] + \dots + f1[Wf-1,1,\text{green}]*I[Wf-1,1,\text{green}] +$
+ ... +
 $f1[0,Hf-1,\text{green}]*I[0,Hf-1,\text{green}] + f1[1,Hf-1,\text{green}]*I[1,Hf-1,\text{green}] + \dots + f1[Wf-1,Hf-1,\text{green}]*I[Wf-1,Hf-1,\text{green}] +$
 $f1[0,0,\text{red}]*I[0,0,\text{red}] + f1[1,0,\text{red}]*I[1,0,\text{red}] + \dots + f1[Wf-1,0,\text{red}]*I[Wf-1,0,\text{red}] +$
 $f1[0,1,\text{red}]*I[0,1,\text{red}] + f1[1,1,\text{red}]*I[1,1,\text{red}] + \dots + f1[Wf-1,1,\text{red}]*I[Wf-1,1,\text{red}] +$
+ ... +
 $f1[0,Hf-1,\text{red}]*I[0,Hf-1,\text{red}] + f1[1,Hf-1,\text{red}]*I[1,Hf-1,\text{red}] + \dots + f1[Wf-1,Hf-1,\text{red}]*I[Wf-1,Hf-1,\text{red}] + \text{Bias1}$

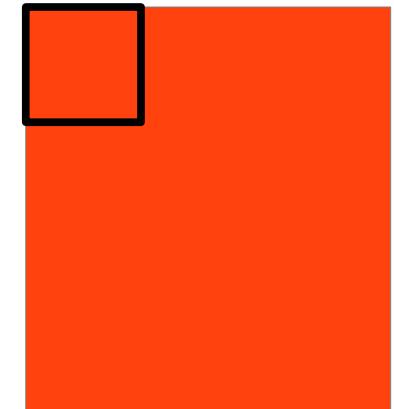
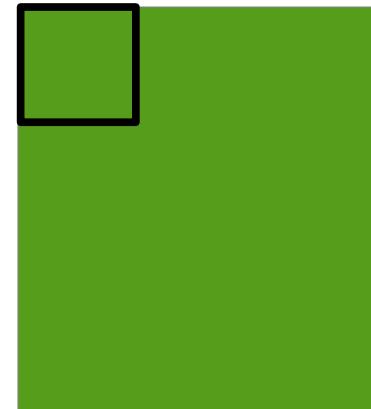
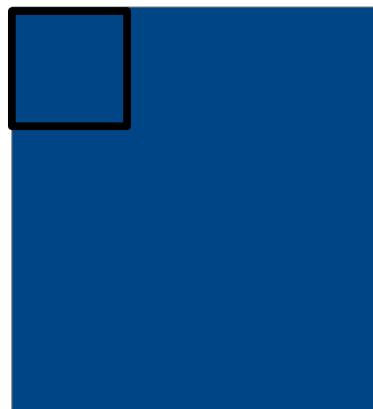


Activation Function

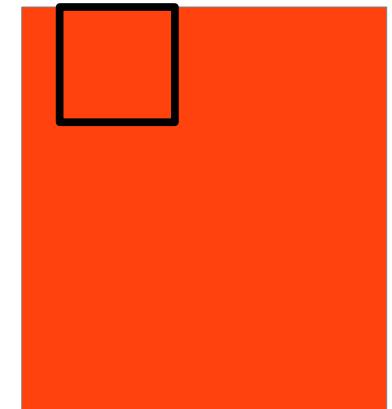
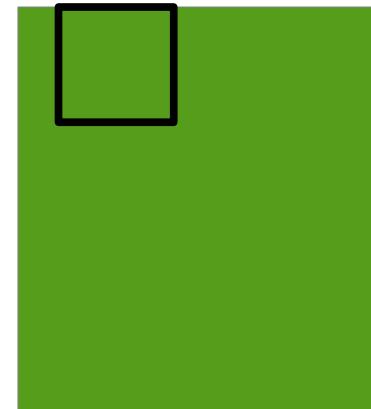
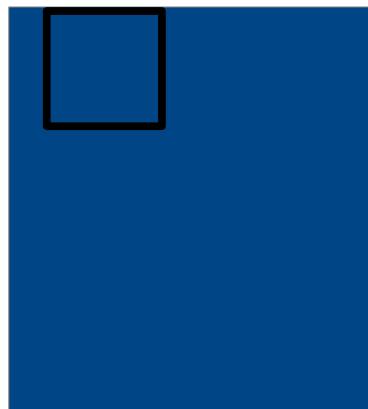


1x1

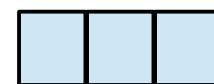
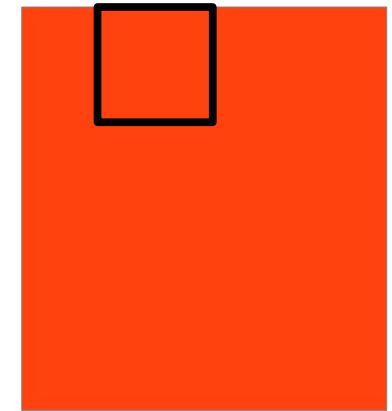
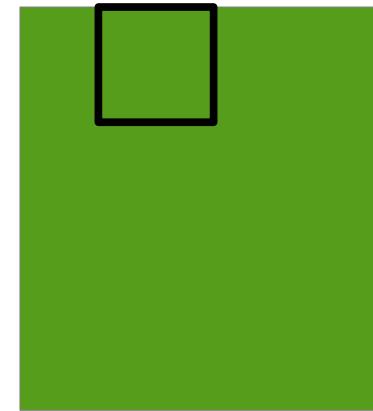
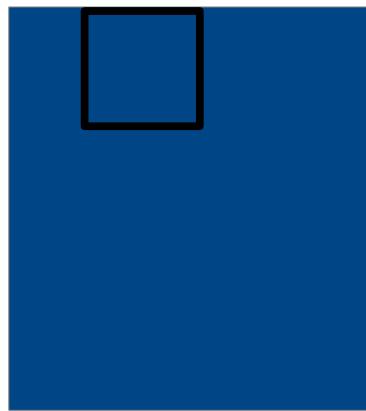
Convolution



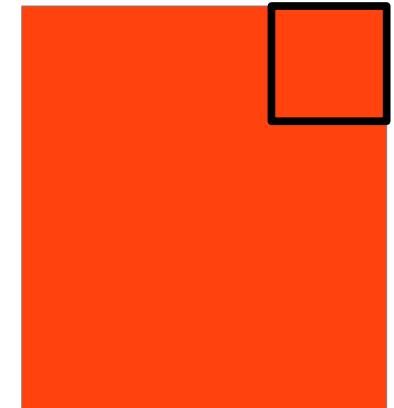
Convolution



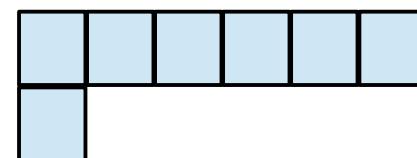
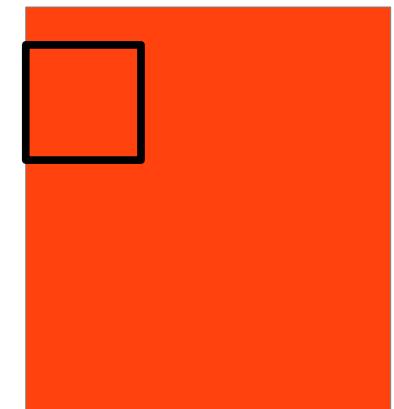
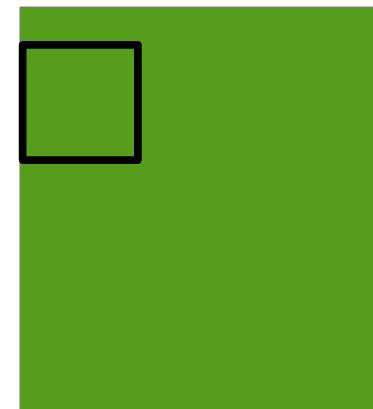
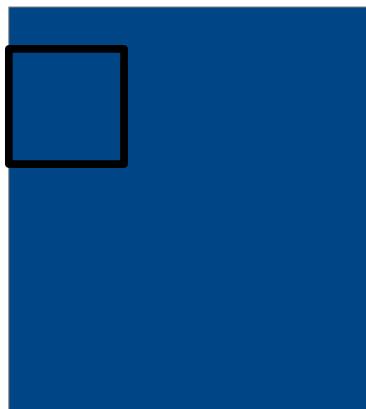
Convolution



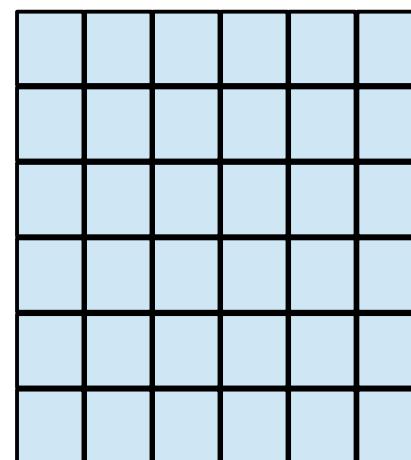
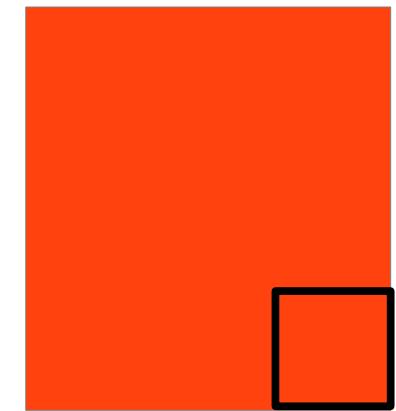
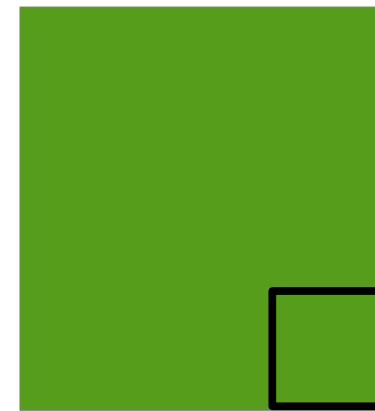
Convolution



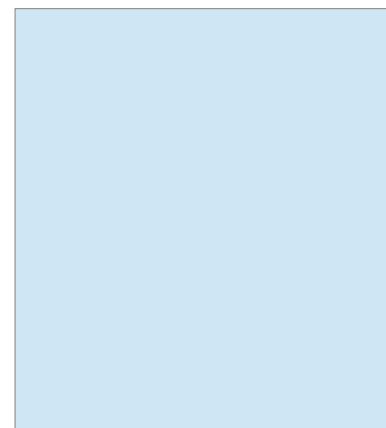
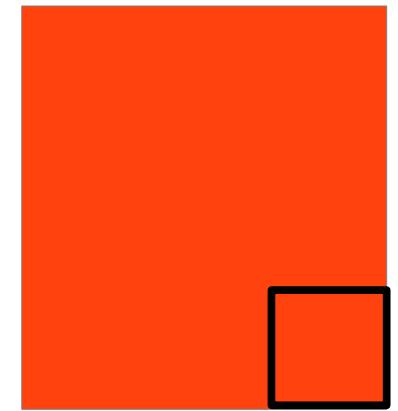
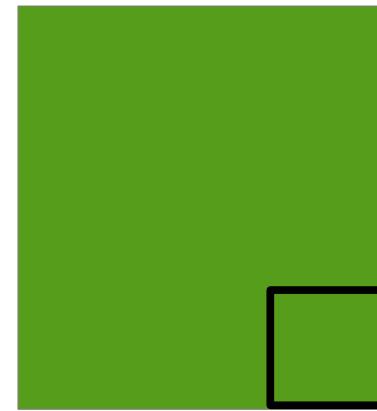
Convolution



Convolution

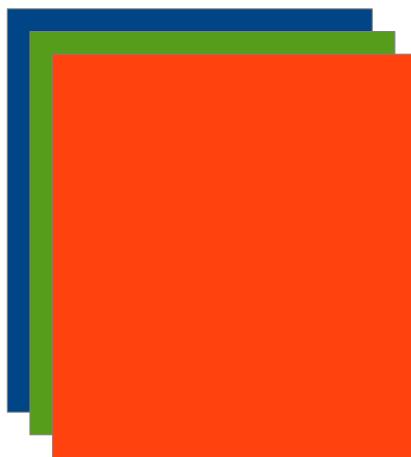


Convolution



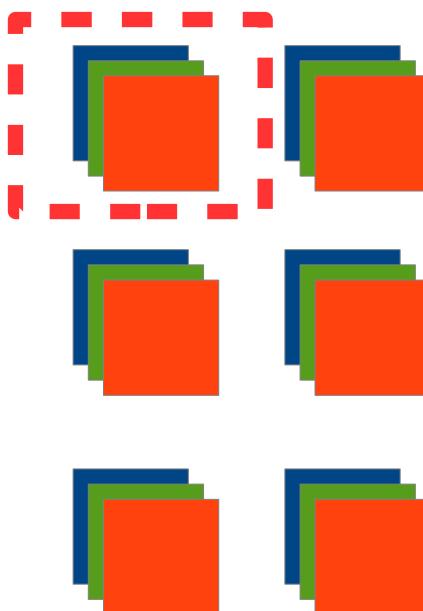
Convolution

Input Feature Map



$W \times H \times 3$

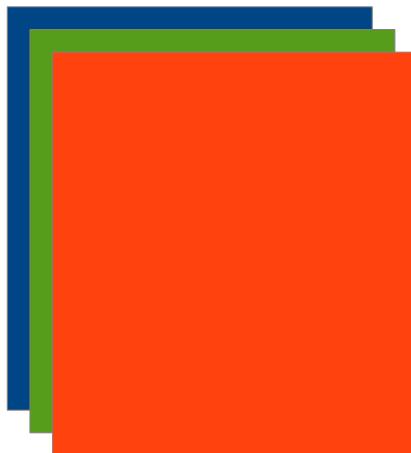
Filters



$W_f \times H_f \times 3$

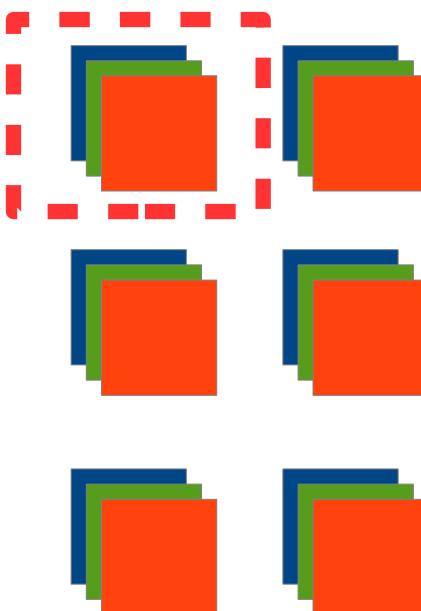
Convolution

Input Feature Map

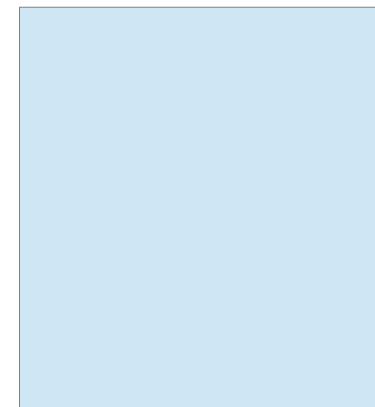


$W \times H \times 3$

Filters



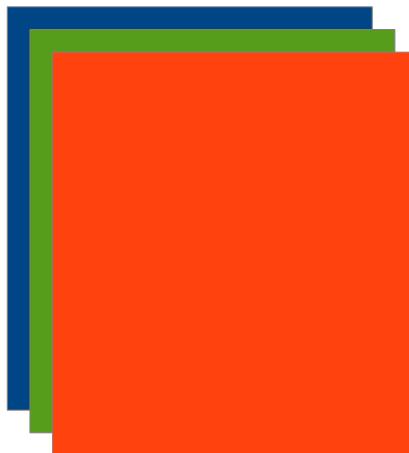
Output Feature Map



$W_f \times H_f \times 3$

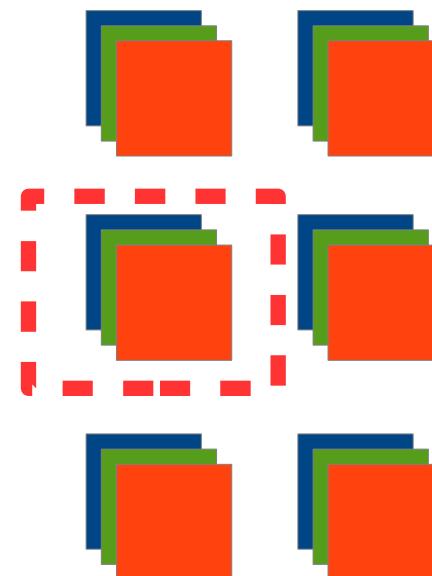
Convolution

Input Feature Map

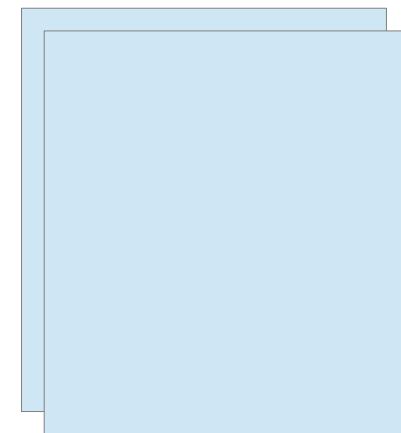


$W \times H \times 3$

Filters



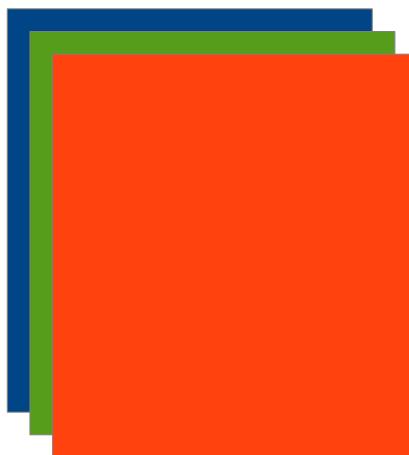
Output Feature Map



$W_f \times H_f \times 3$

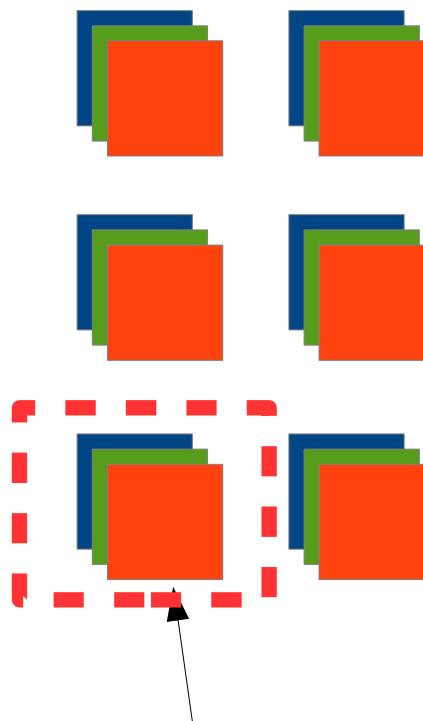
Convolution

Input Feature Map



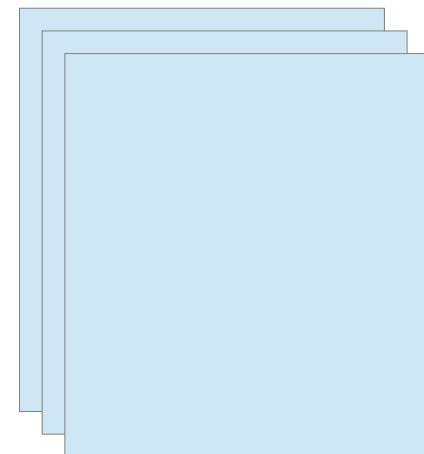
$W \times H \times 3$

Filters



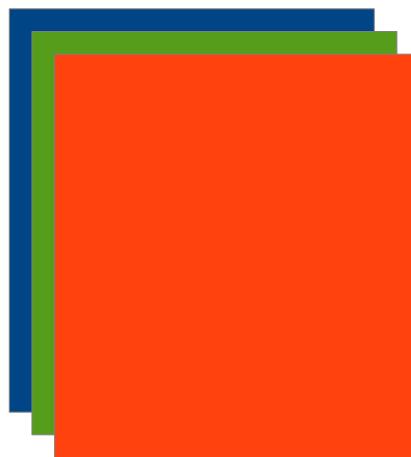
$W_f \times H_f \times 3$

Output Feature Map

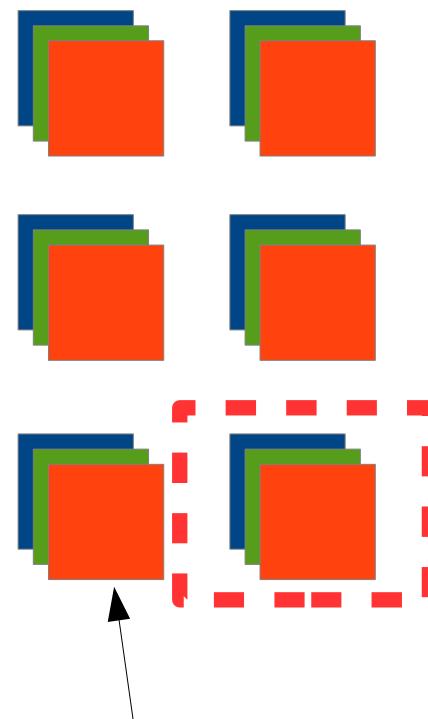


Convolution

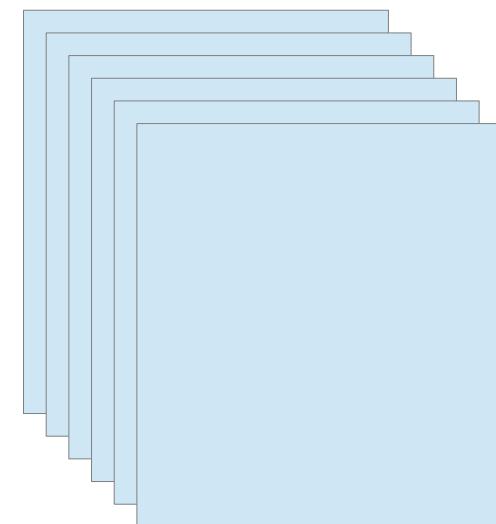
Input Feature Map



Filters

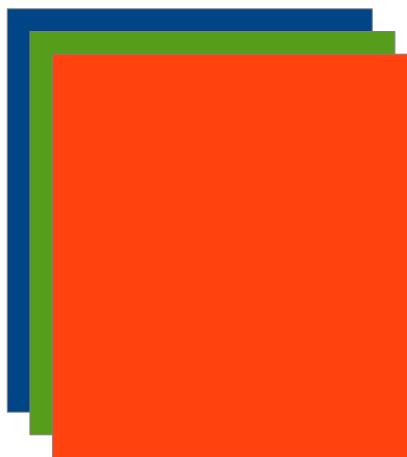


Output Feature Map

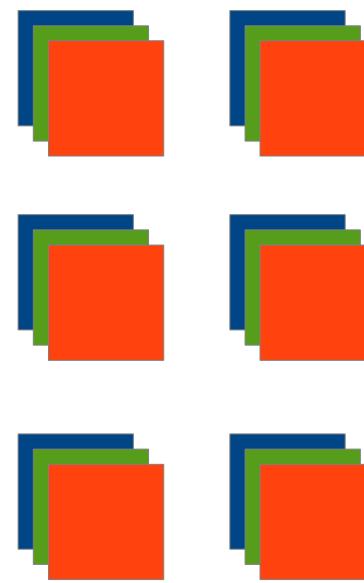


Convolution

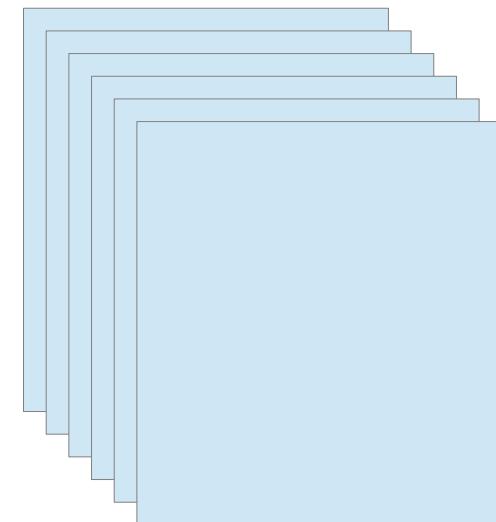
Input Feature Map



Filters



Output Feature Map



$W \times H \times 3$

$W_f \times H_f \times 3$

$W \times H \times 6$

CONV Layer Computation

Output fmaps (O) **Input fmaps (I)** **Filter weights (W)**

Biases (B)

$$\underline{\mathbf{O}[n][m][x][y]} = \text{Activation}(\underline{\mathbf{B}[m]} + \sum_{i=0}^{R-1} \sum_{j=0}^{S-1} \sum_{k=0}^{C-1} \underline{\mathbf{I}[n][k][Ux+i][Uy+j]} \times \underline{\mathbf{W}[m][k][i][j]}),$$

$$0 \leq n < N, 0 \leq m < M, 0 \leq y < E, 0 \leq x < F,$$

$$E = (H - R + U)/U, F = (W - S + U)/U.$$

Shape Parameter	Description
N	fmap batch size
M	# of filters / # of output fmap channels
C	# of input fmap/filter channels
H/W	input fmap height/width
R/S	filter height/width
E/F	output fmap height/width
U	convolution stride

CONV Layer Implementation

- Naïve 7-layer for-loop implementation

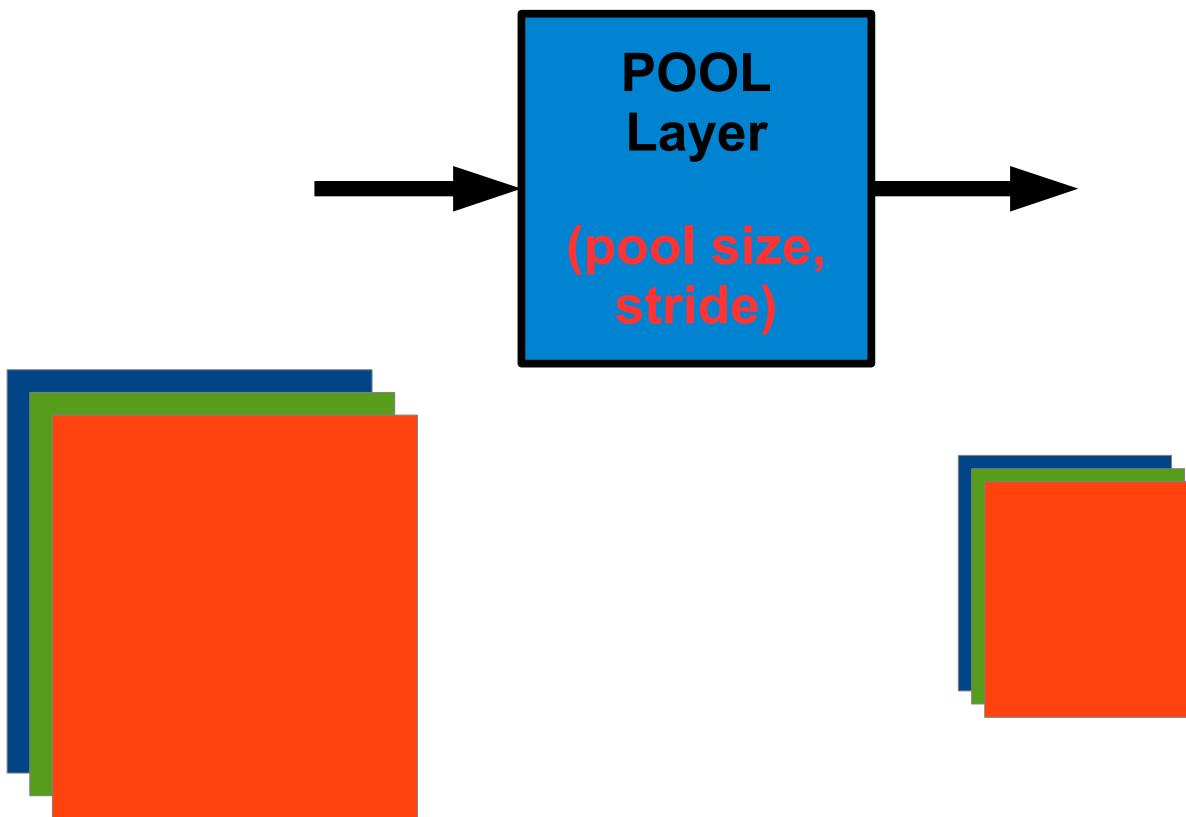
```

for (n=0; n<N; n++) {
    for (m=0; m<M; m++) {
        for (x=0; x<F; x++) {
            for (y=0; y<E; y++) {
                o[n][m][x][y] = B[m];
                for (i=0; i<R; i++) {
                    for (j=0; j<S; j++) {
                        for (k=0; k<C; k++) {
                            o[n][m][x][y] += I[n][k][Ux+i][Uy+j] * W[m][k][i][j];
                        }
                    }
                }
                o[n][m][x][y] = Activation(o[n][m][x][y]);
            }
        }
    }
}

```

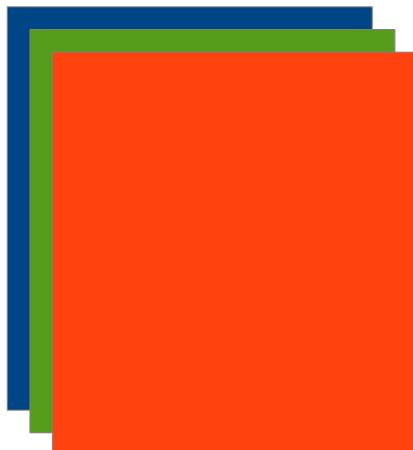
} for each output fmap value

Pooling



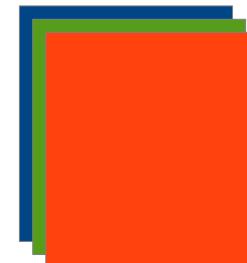
Pooling (Max/Avg)

Input Feature Map



$W \times H \times 3$

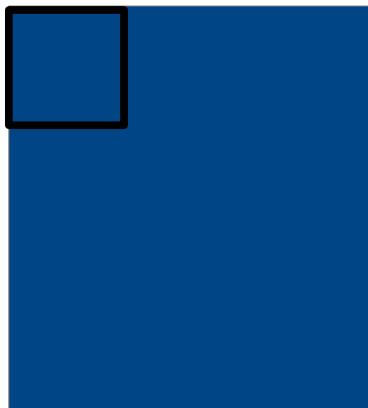
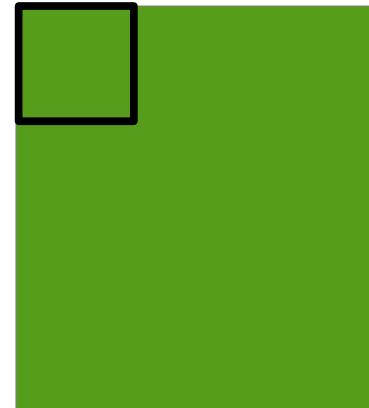
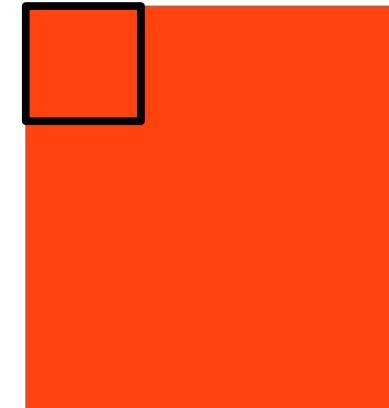
Output Feature Map



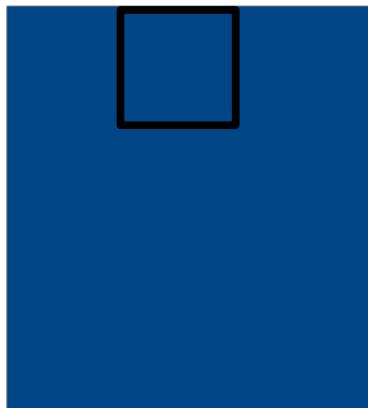
Pool size 2

$W/2 \times H/2 \times 3$

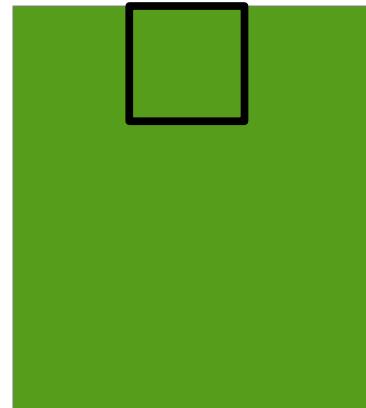
Max Pool


$$\text{Max}(I[0,0,\blacksquare], I[1,0,\blacksquare], \\ I[0,1,\blacksquare], I[1,1,\blacksquare])$$

$$\text{Max}(I[0,0,\blacksquare], I[1,0,\blacksquare], \\ I[0,1,\blacksquare], I[1,1,\blacksquare])$$

$$\text{Max}(I[0,0,\blacksquare], I[1,0,\blacksquare], \\ I[0,1,\blacksquare], I[1,1,\blacksquare])$$

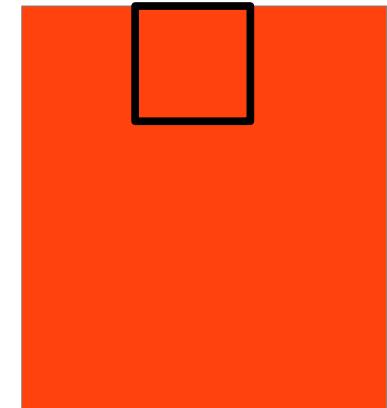

Max Pool



$\text{Max}(I[2,0, \blacksquare], I[3,0, \blacksquare],
I[2,1, \blacksquare], I[3,1, \blacksquare])$



$\text{Max}(I[2,0, \blacksquare], I[3,0, \blacksquare],
I[2,1, \blacksquare], I[3,1, \blacksquare])$

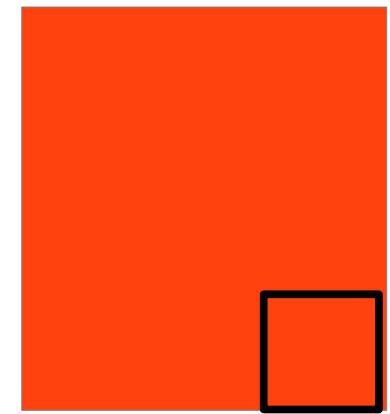
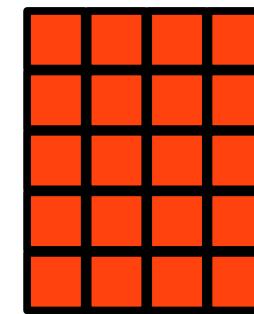
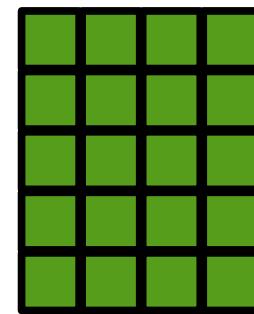
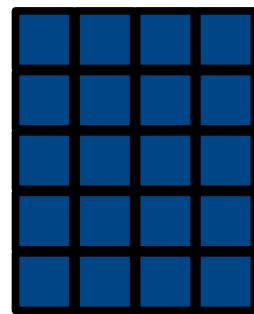


$\text{Max}(I[2,0, \blacksquare], I[3,0, \blacksquare],
I[2,1, \blacksquare], I[3,1, \blacksquare])$

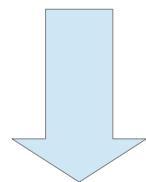
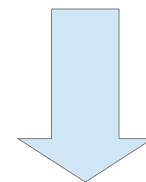
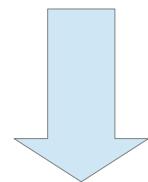
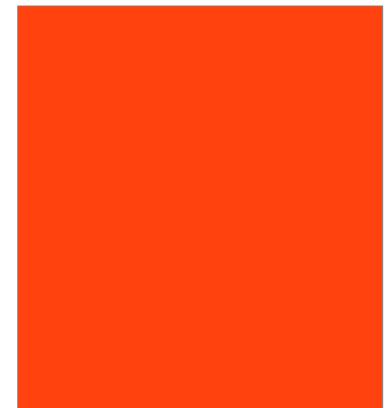


Max Pool


$$\text{Max}(I[2,0,\blacksquare], I[3,0,\blacksquare], \\ I[2,1,\blacksquare], I[3,1,\blacksquare])$$

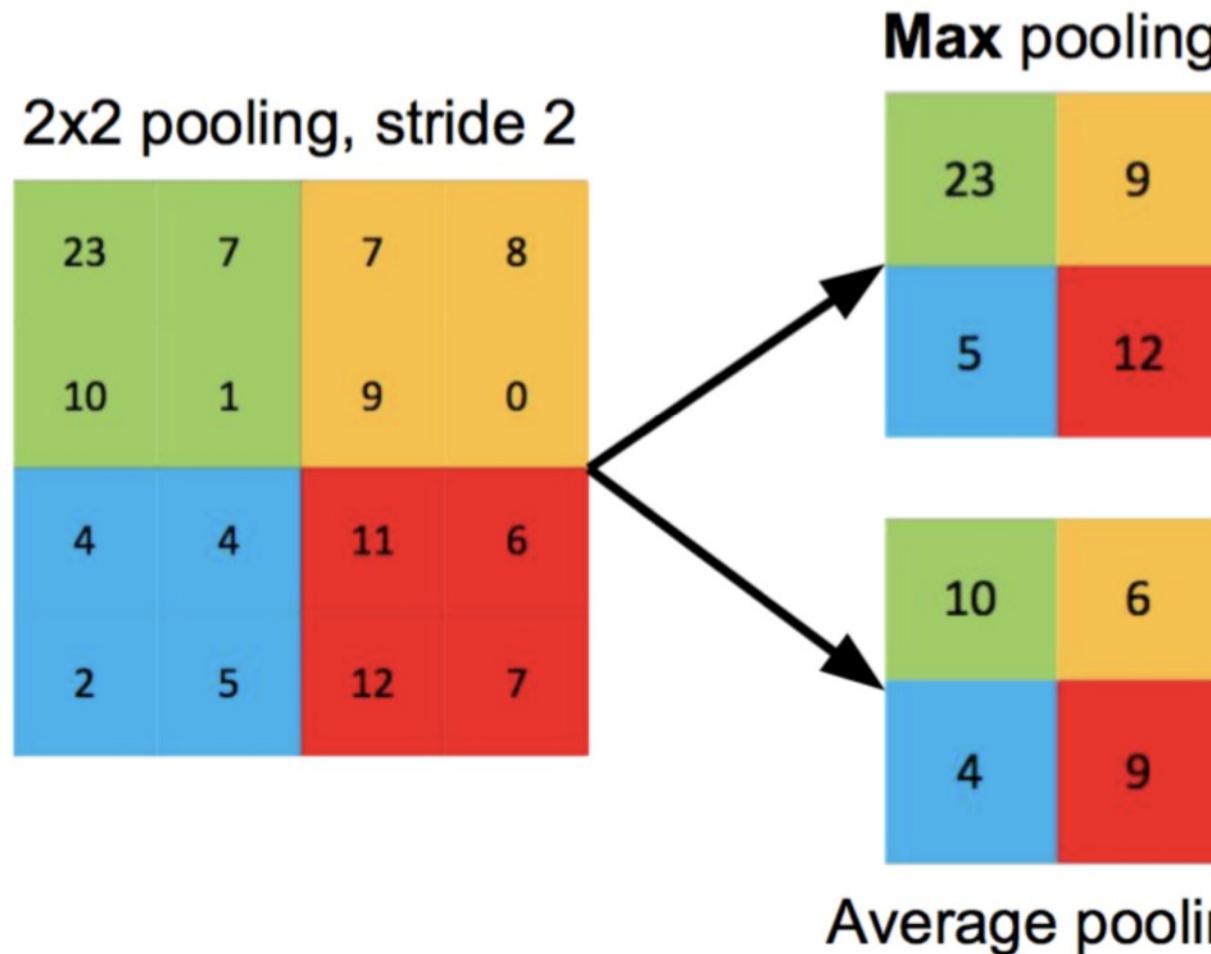
$$\text{Max}(I[2,0,\blacksquare], I[3,0,\blacksquare], \\ I[2,1,\blacksquare], I[3,1,\blacksquare])$$

$$\text{Max}(I[2,0,\blacksquare], I[3,0,\blacksquare], \\ I[2,1,\blacksquare], I[3,1,\blacksquare])$$


Max Pool

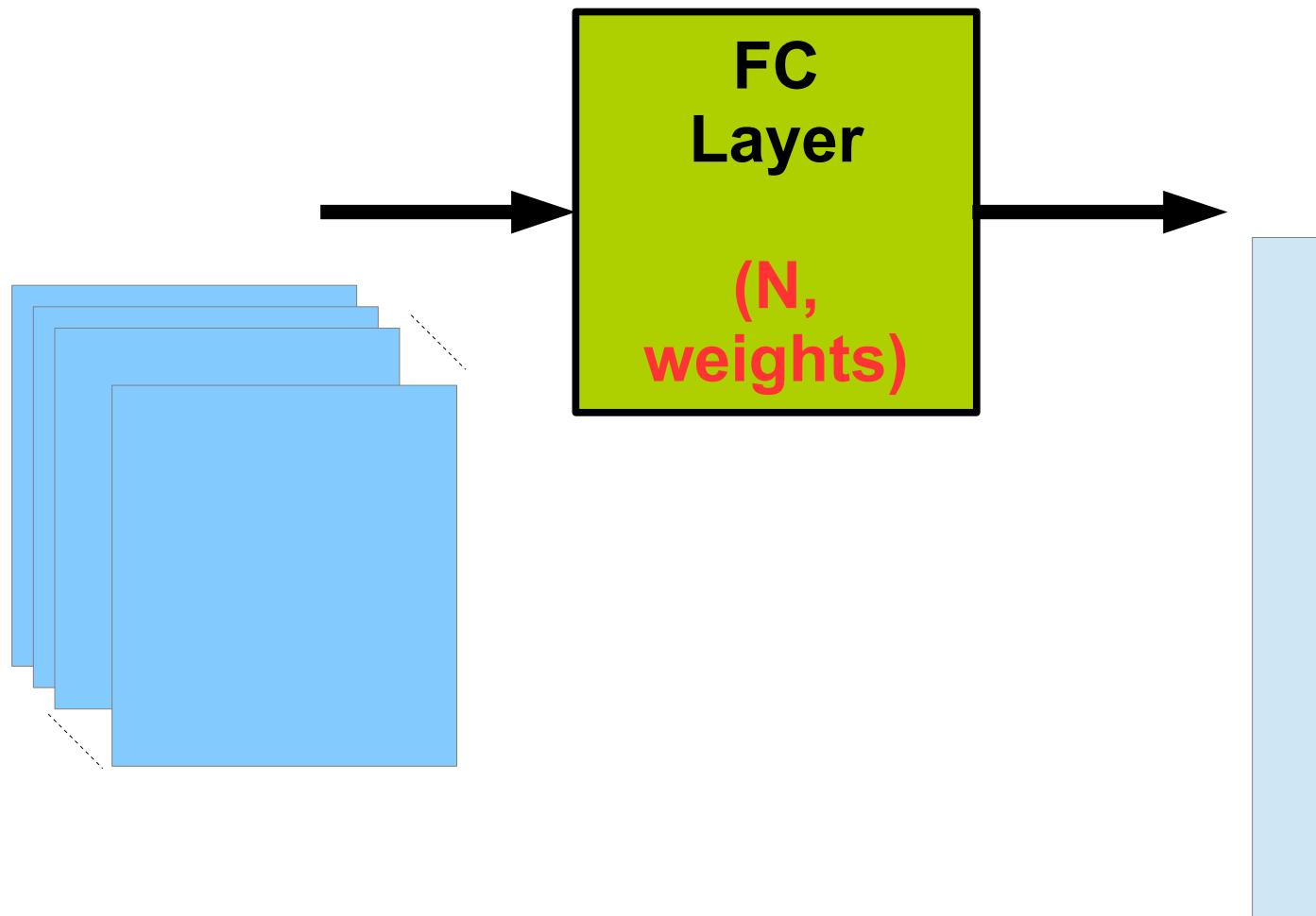


Pooling Layer

- Reduce resolution of each channel independently
- Overlapping or non-overlapping depending on stride

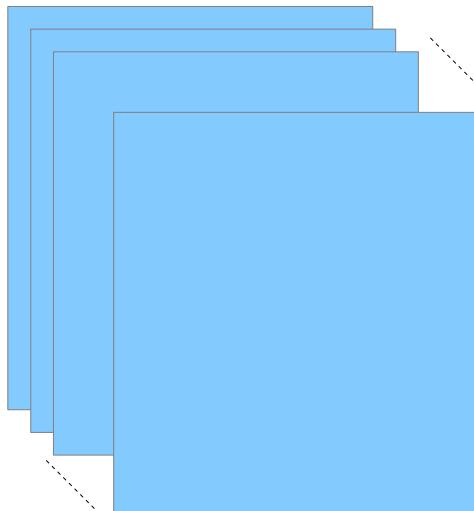


Fully Connected Layer



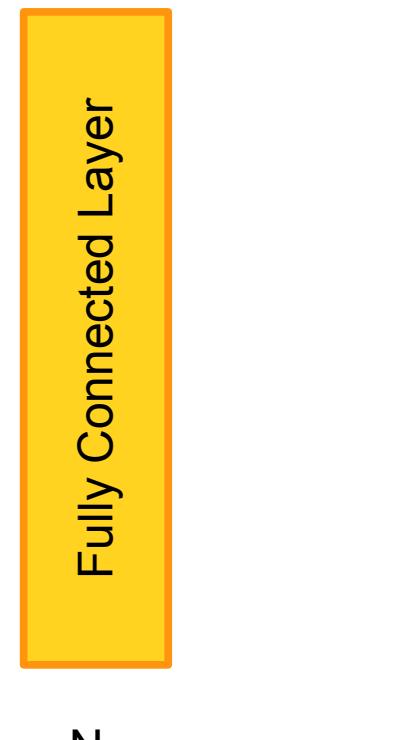
Fully Connected Layer

Input Feature Map



$W \times H \times C$

Output Feature Map



Z

$N \times 1$

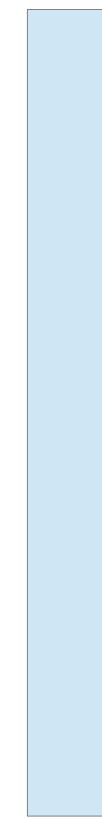
Fully Connected Layer

Input Feature Map



$(W \times H \times C) \times 1$

Output Feature Map



$N \times 1$

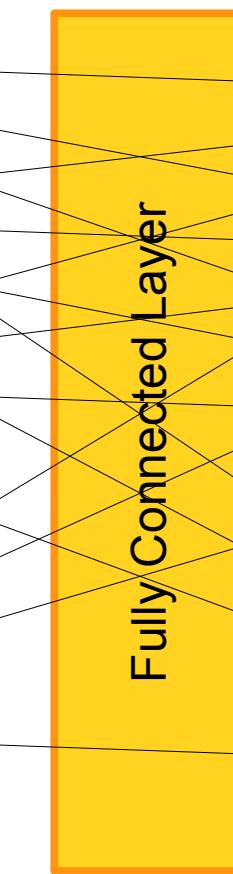


N

Fully Connected Layer

Input Feature Map

$(W \times H \times C) \times 1$



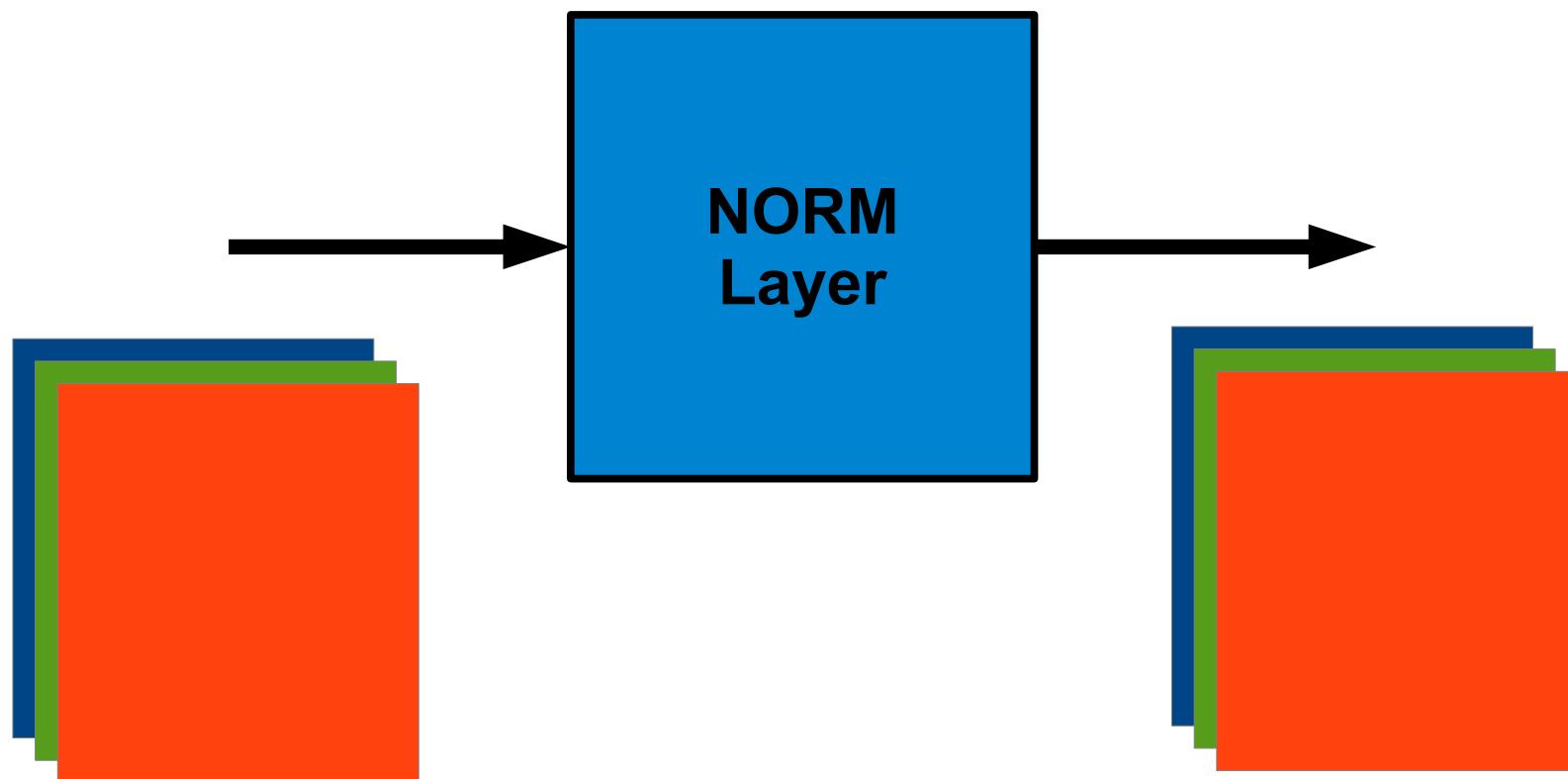
N

Output Feature Map

$N \times 1$

$(W \times H \times C \times N)$ weights

Normalization Layer



Normalization Layer

- Batch Normalization (BN)
 - Normalize activations towards mean 0 and stdev 1
 - Believed to be key to getting high accuracy and faster training on very deep neural networks

[S. Loffe & C. Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. ICML 2015]

BN Layer Implementation

- The normalized value is further scaled and shifted, the parameters of which are learned from training

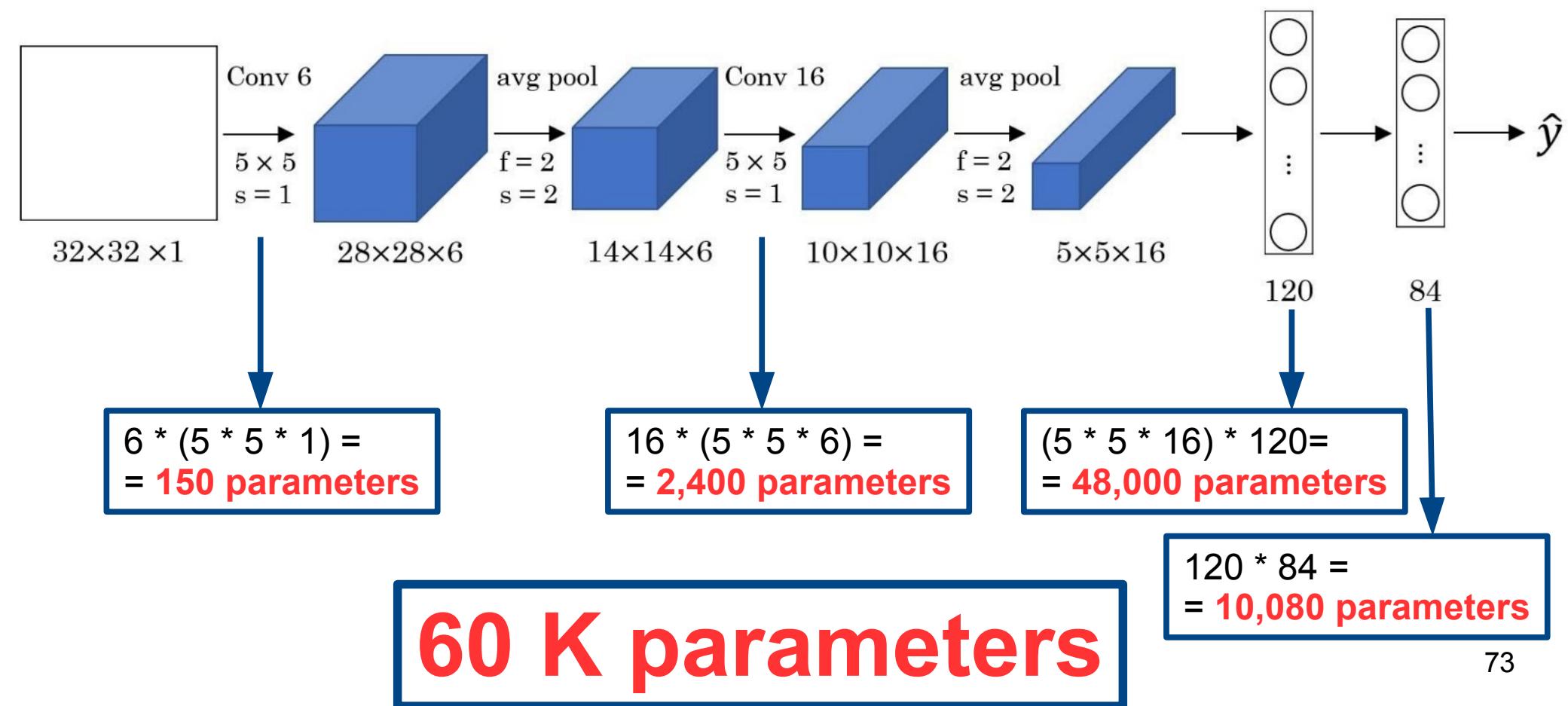
$$y = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta$$

Annotations for the BN formula:

- data mean**: μ (red arrow pointing down)
- data std. dev.**: σ (red arrow pointing up)
- learned scale factor**: γ (blue arrow pointing down)
- learned shift factor**: β (blue arrow pointing up)
- small const. to avoid numerical problems: ϵ

LeNet-5

LeCun et al., 1998. Gradient-based learning applied to document recognition



LeNet-5

LeCun et al., 1998. Gradient-based learning applied to document recognition

- MNIST

- 28x28 pixels (B&W)
- 10 Classes
- 60,000 Training
- 10,000 Testing



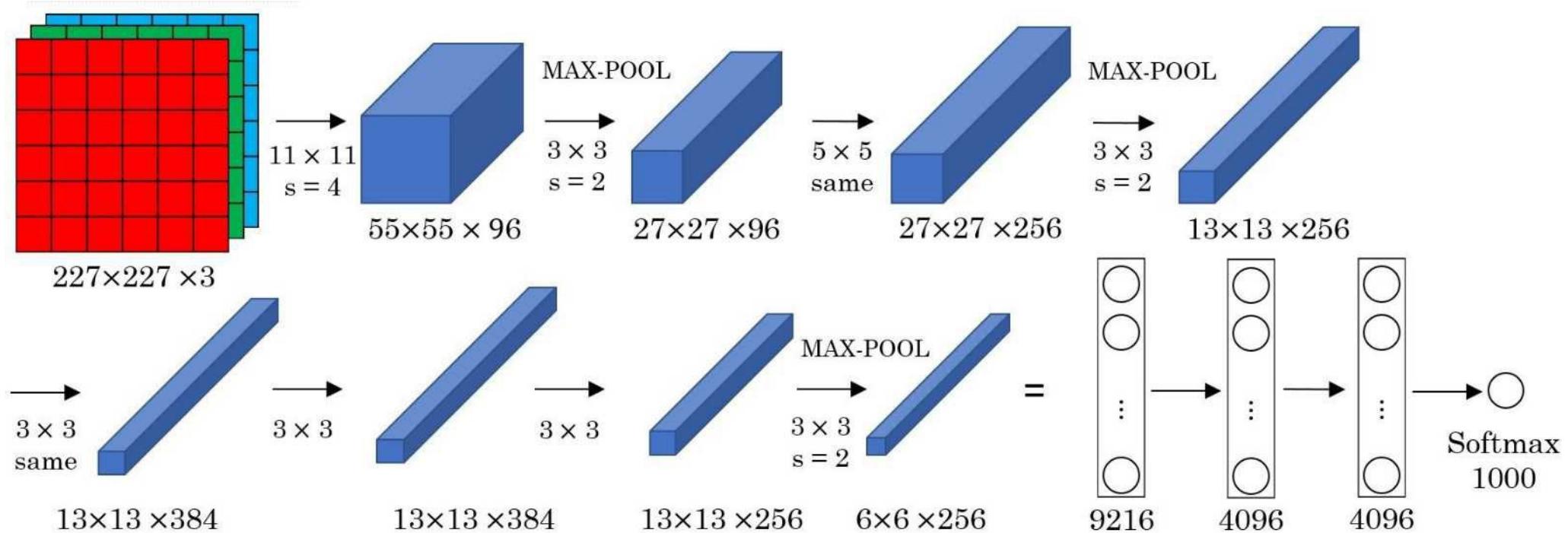
LeNet-5

LeCun et al., 1998. Gradient-based learning applied to document recognition

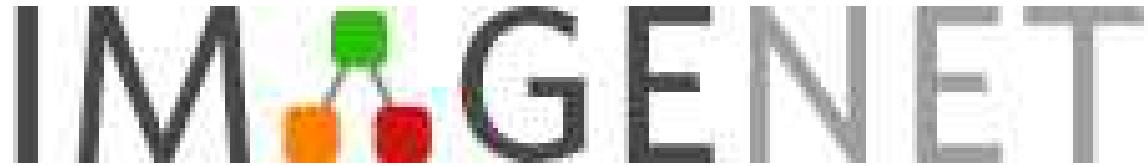
<http://yann.lecun.com/exdb/lenet/>

AlexNet

Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks



~60 M parameters



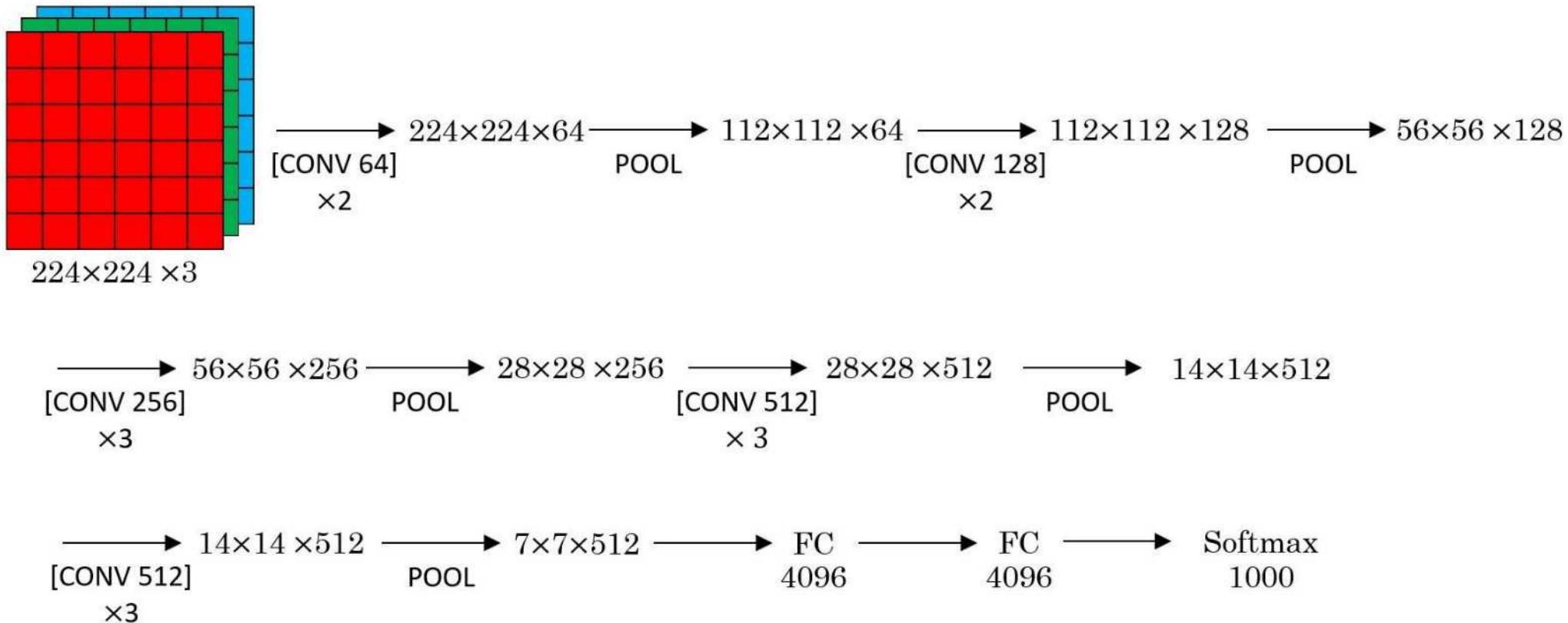
- <http://www.image-net.org/challenges/LSVRC/>
- ~256x256 pixels (color)
- 1000 Classes
- 1.3M Training
- 100,000 Testing

VGG-16

Simonyan & Zisserman 2015. Very deep convolutional networks for large-scale image recognition

CONV = 3×3 filter, $s = 1$, same

MAX-POOL = 2×2 , $s = 2$



~138 M parameters

Evolution of DNN Architectures

Year	Model name	Accuracy (%)	Size (MB)	MACs (G)	#Params (M)
2012	AlexNet	56.48	237.89	0.72	61.10
2014	VGG19	72.38	461.10	0.40	20.08
2014	Inception_v3	79.35	5212.60	1.53	6.25
2016	ResNet18	76.82	721.60	0.56	11.22
2016	ResNet50	78.07	4233.60	1.30	23.70
2016	ResNet101	77.35	6409.60	2.52	42.70
2016	ResNet152	78.28	9097.60	3.74	58.34
2018	DenseNet121	78.46	5025.60	0.90	7.05
2018	DenseNet169	75.56	6111.60	1.07	12.64
2018	DenseNet201	76.87	7947.60	1.38	18.28
2020	COVID-Net-Large*	-na-	270.88	0.13	33.85

Variants and Expedients

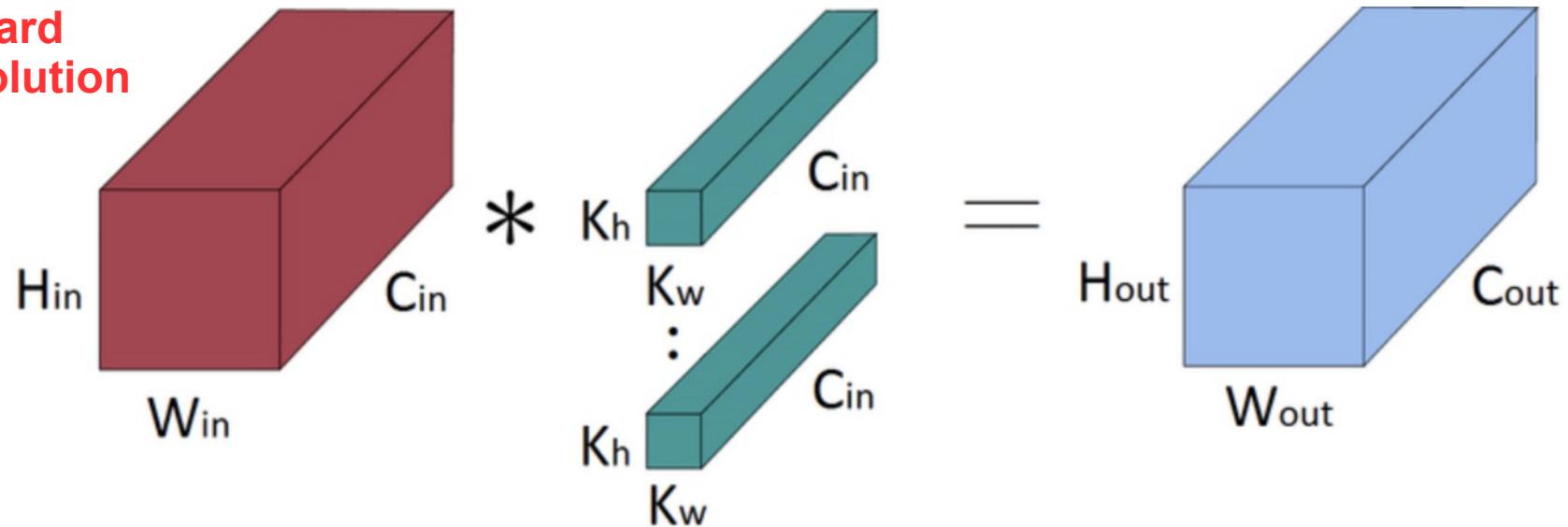
- Grouped Convolution
- Stacked Filter
- Inception Module
- 1x1 Filter
- Residual Connection

Grouped Convolution

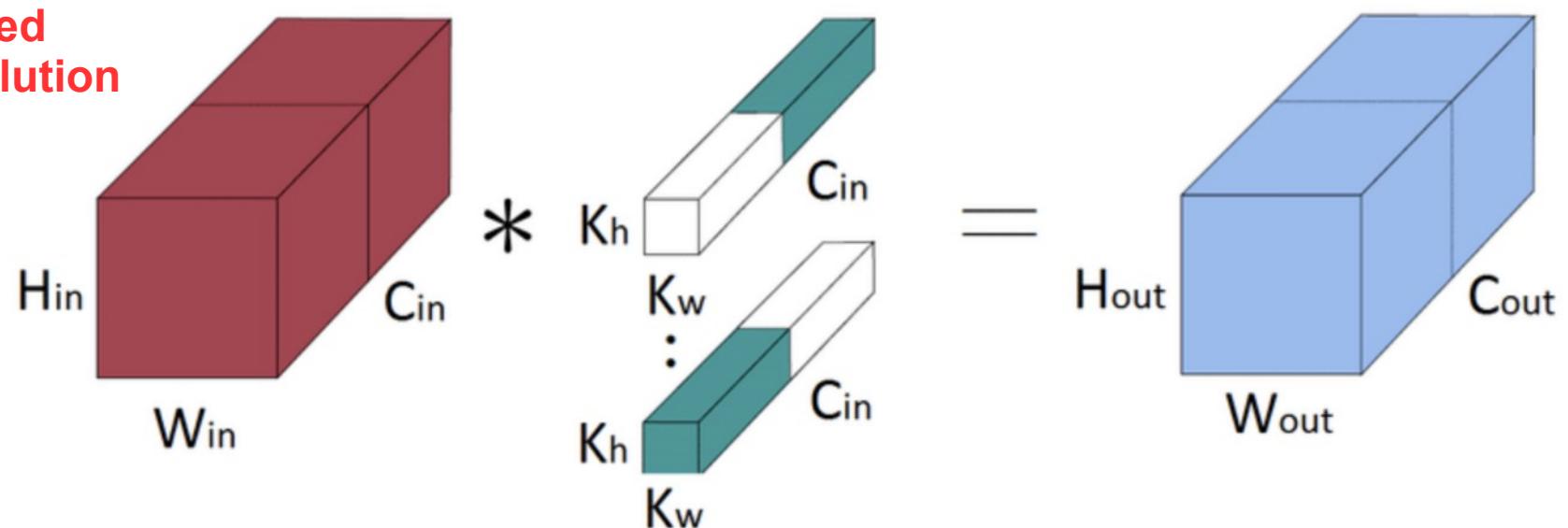
- To reduce the amount of weights
 - e.g., used in AlexNet
- Channels of the Ifmap are split such that filters can have less channels

Grouped Convolution

Standard
Convolution

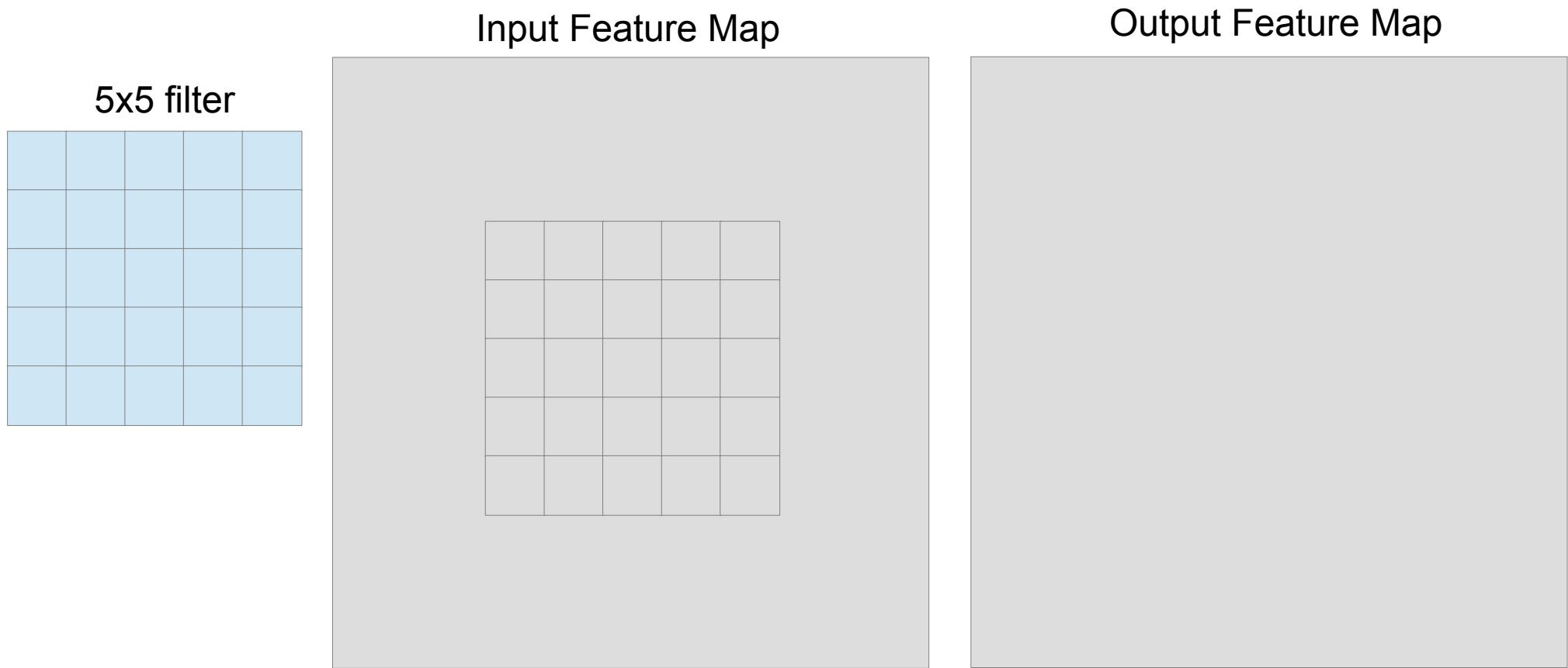


Grouped
Convolution



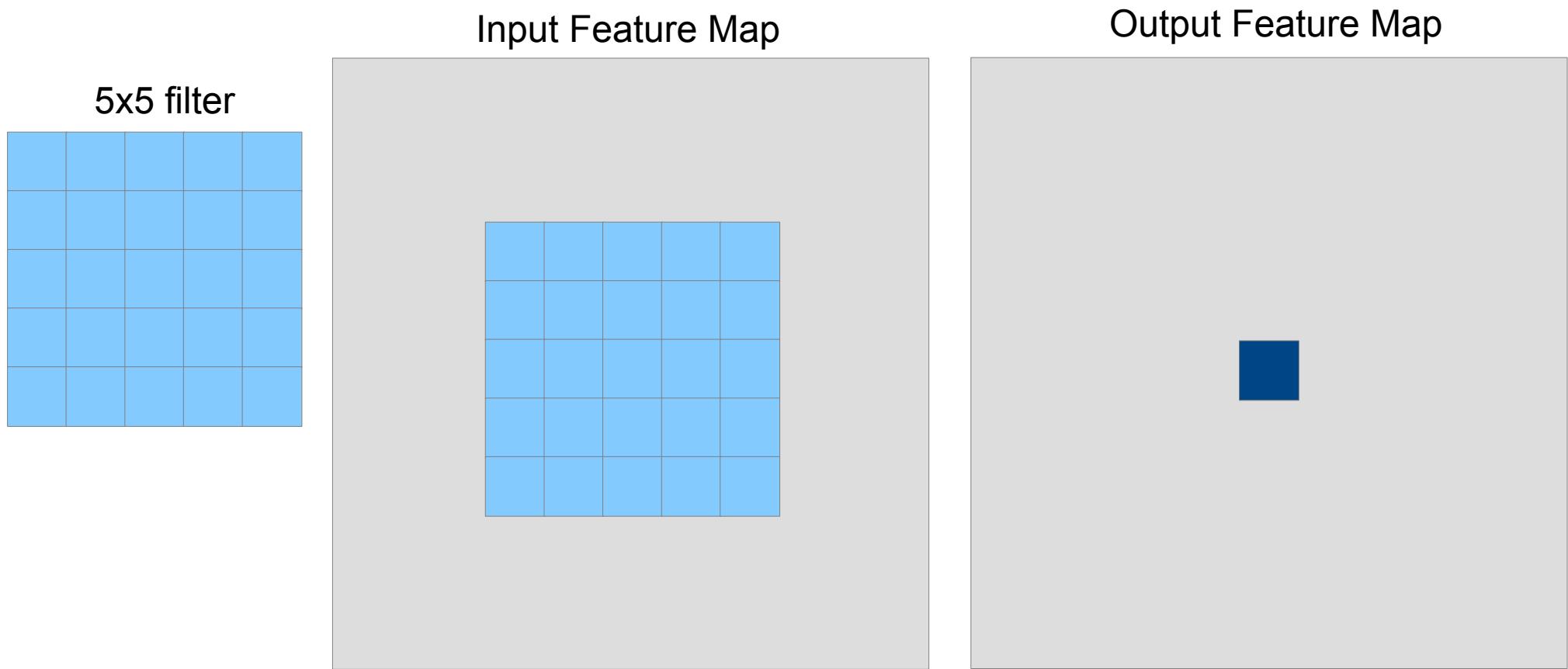
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16



Stacked Filter

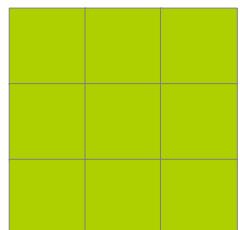
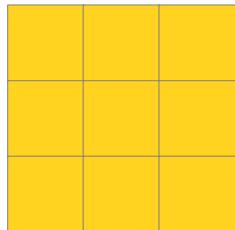
- Large filters built from multiple small filters
 - e.g., used in VGG-16



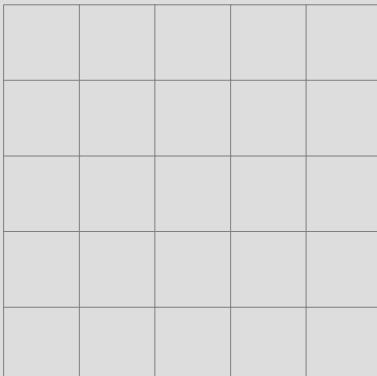
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16

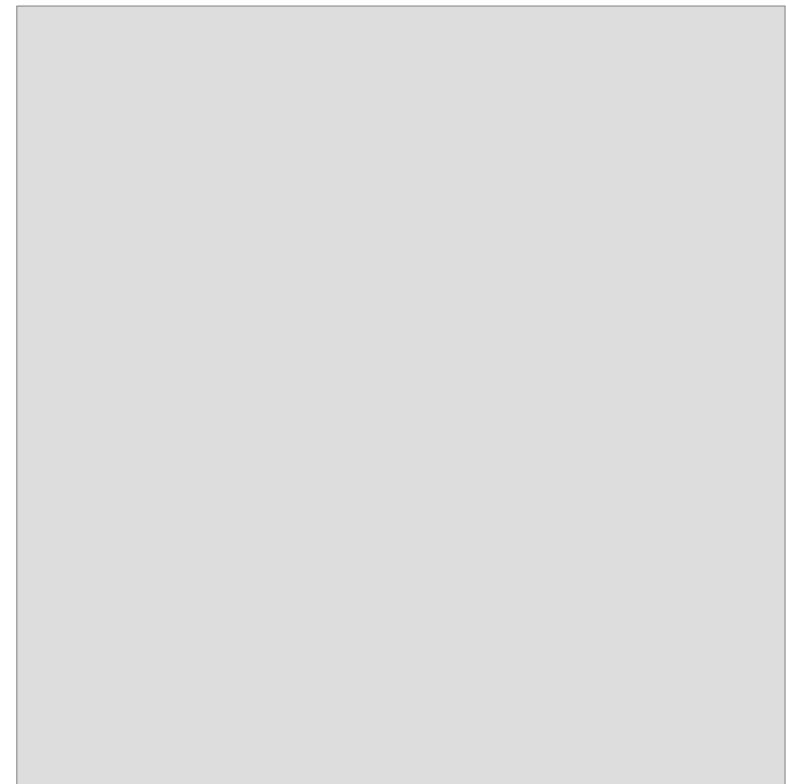
3x3 filters



Input Feature Map

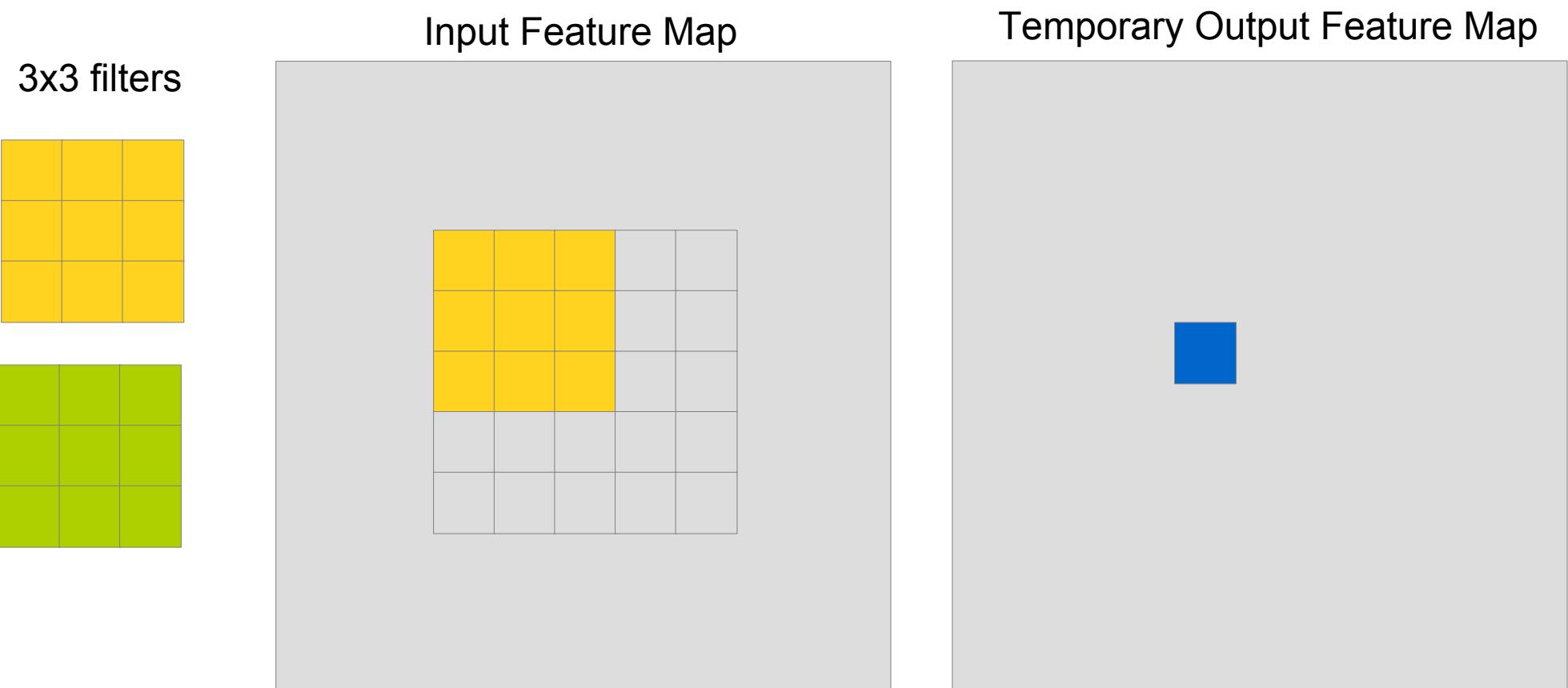


Temporary Output Feature Map



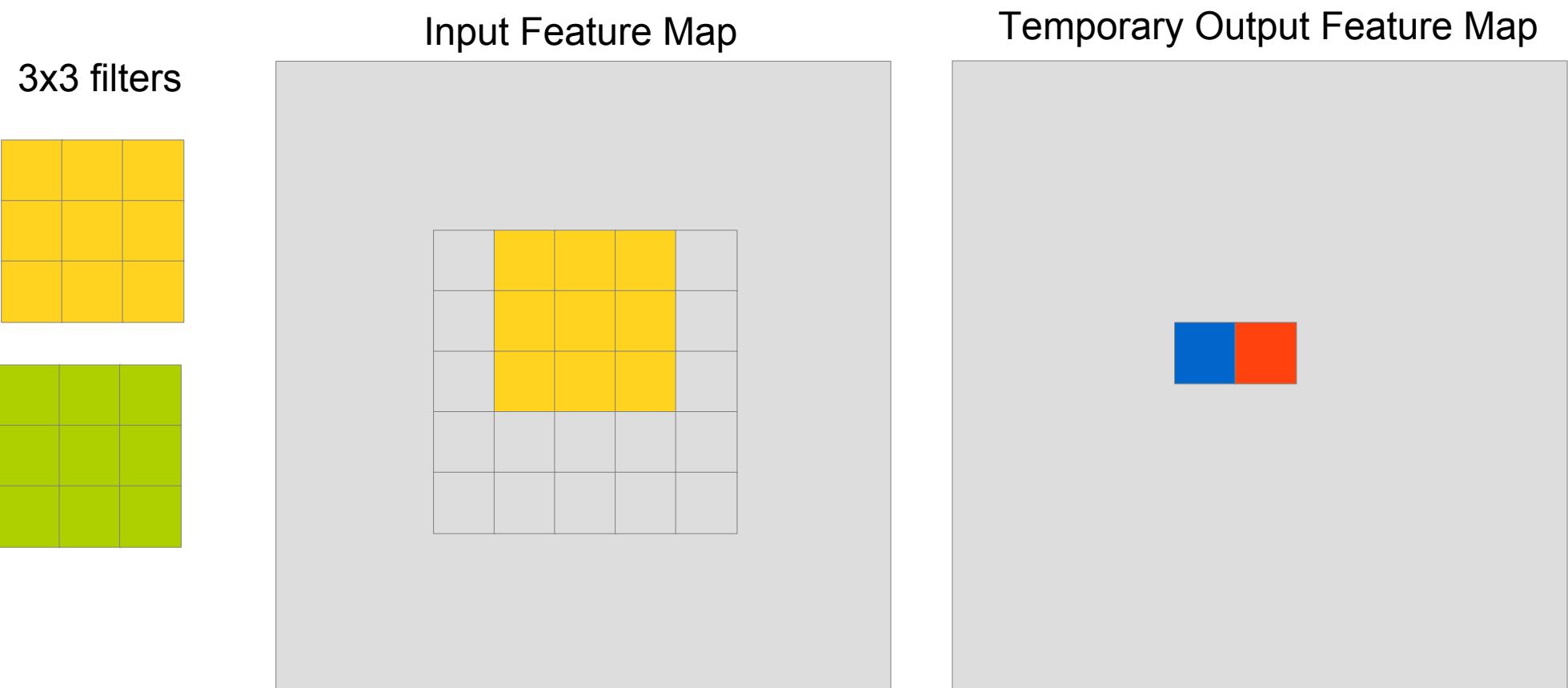
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16



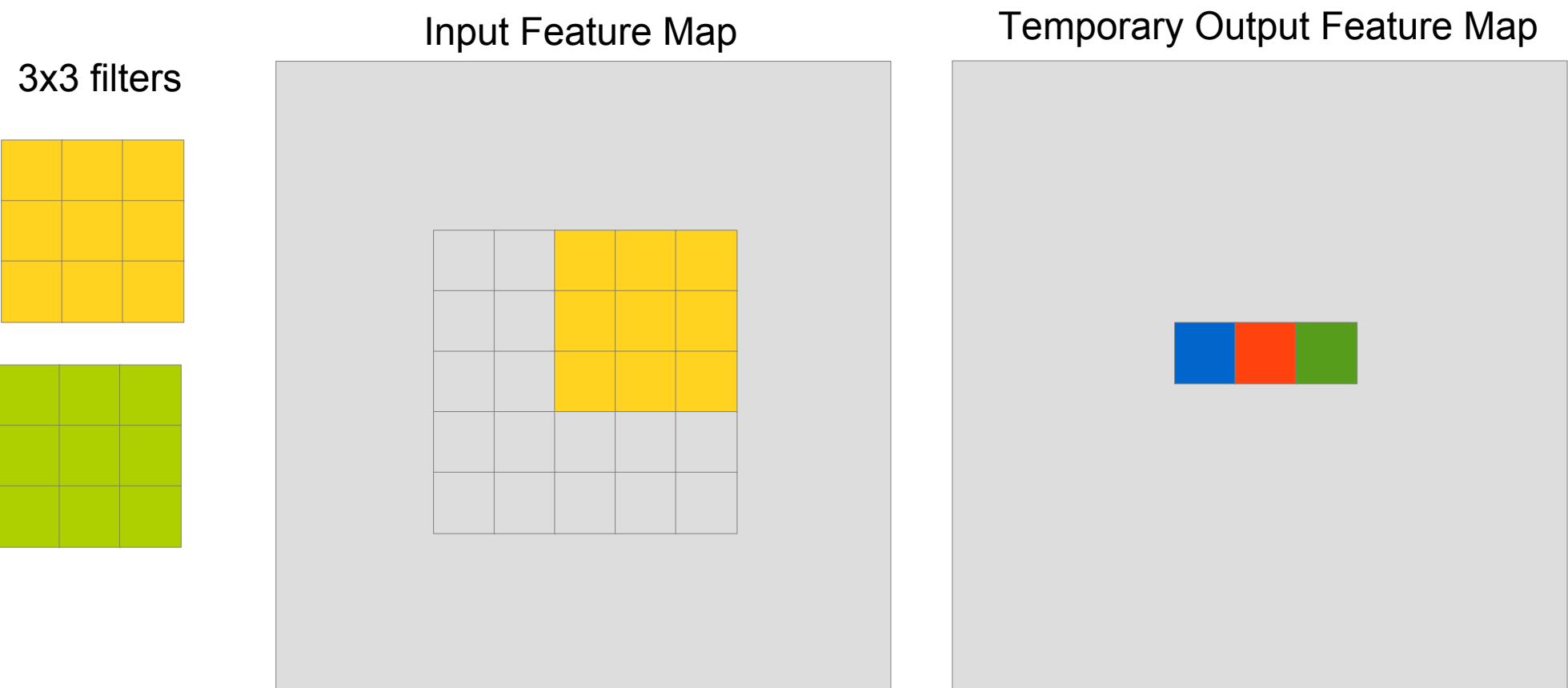
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16



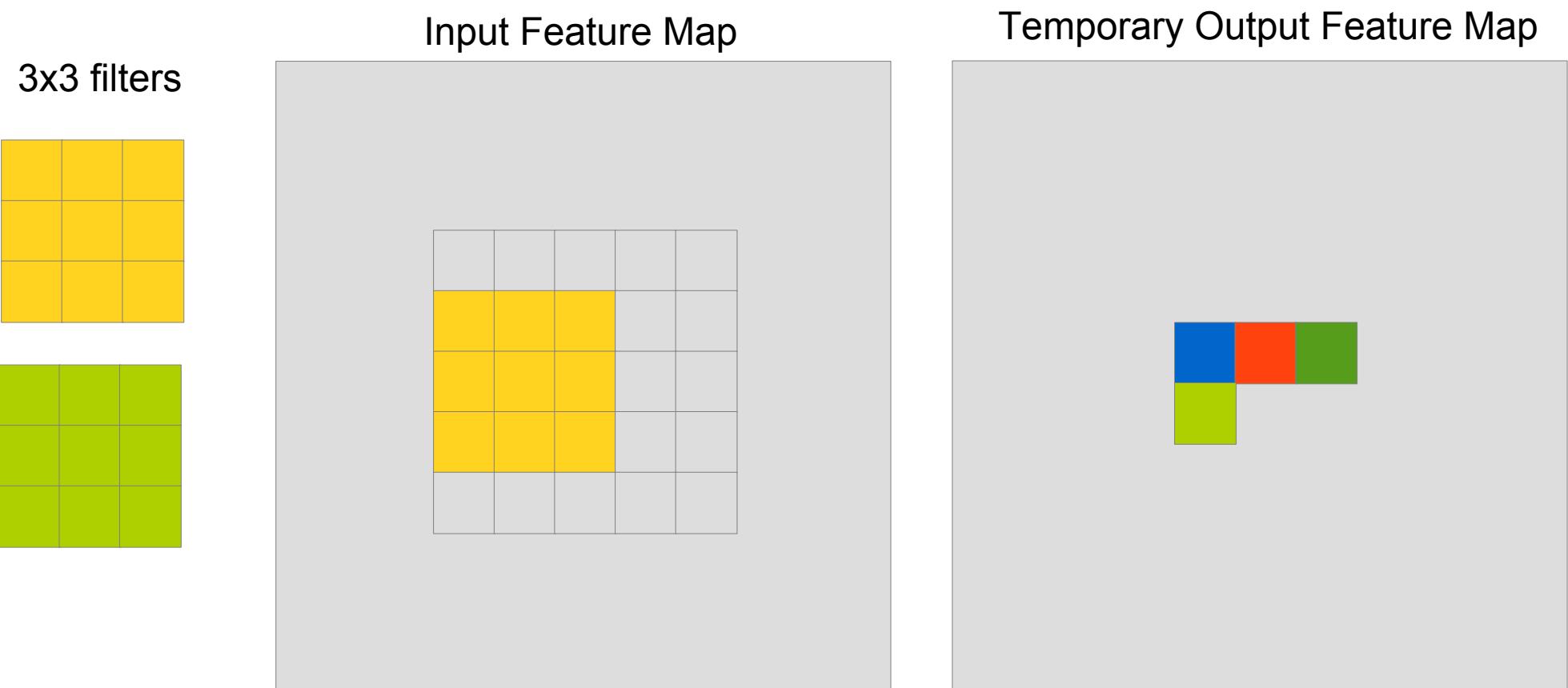
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16



Stacked Filter

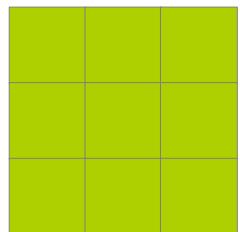
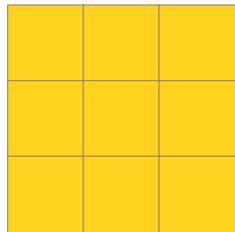
- Large filters built from multiple small filters
 - e.g., used in VGG-16



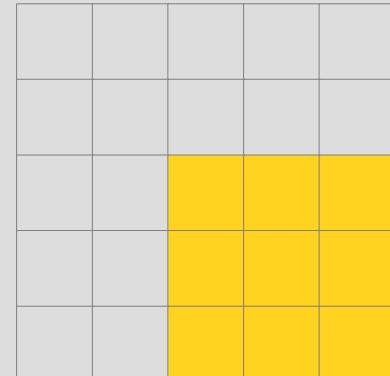
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16

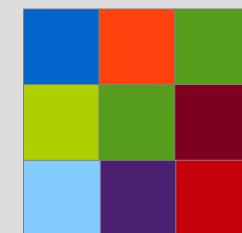
3x3 filters



Input Feature Map



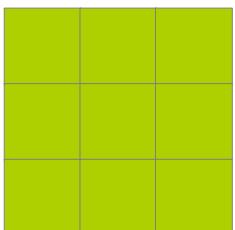
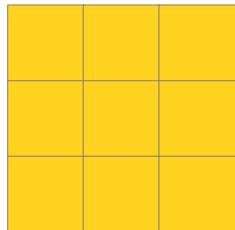
Temporary Output Feature Map



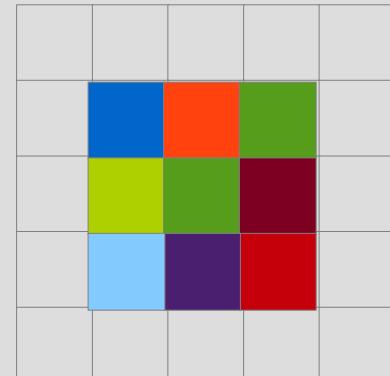
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16

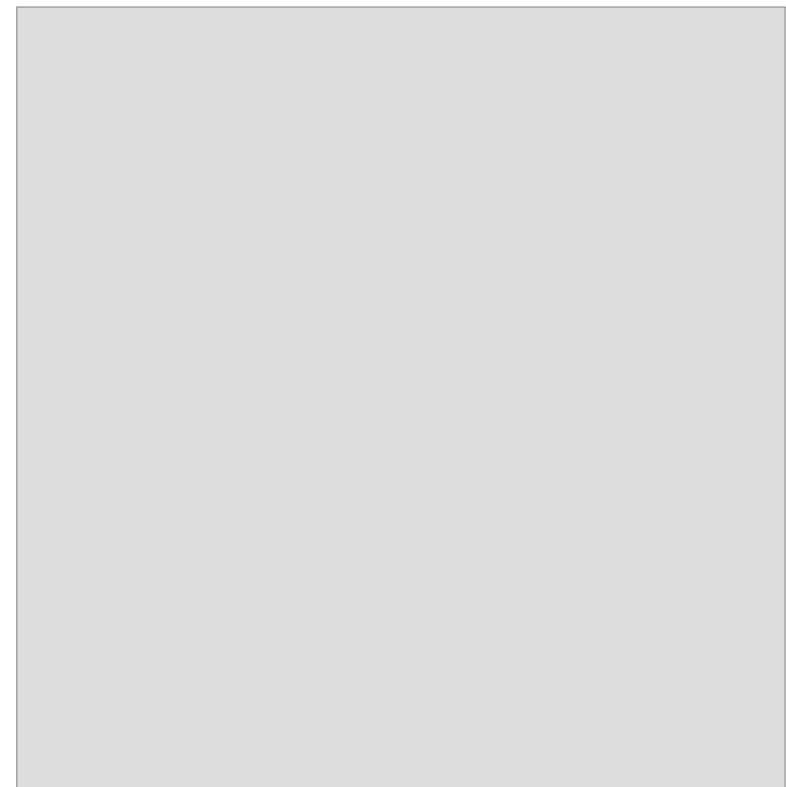
3x3 filters



Temporary Output Feature Map

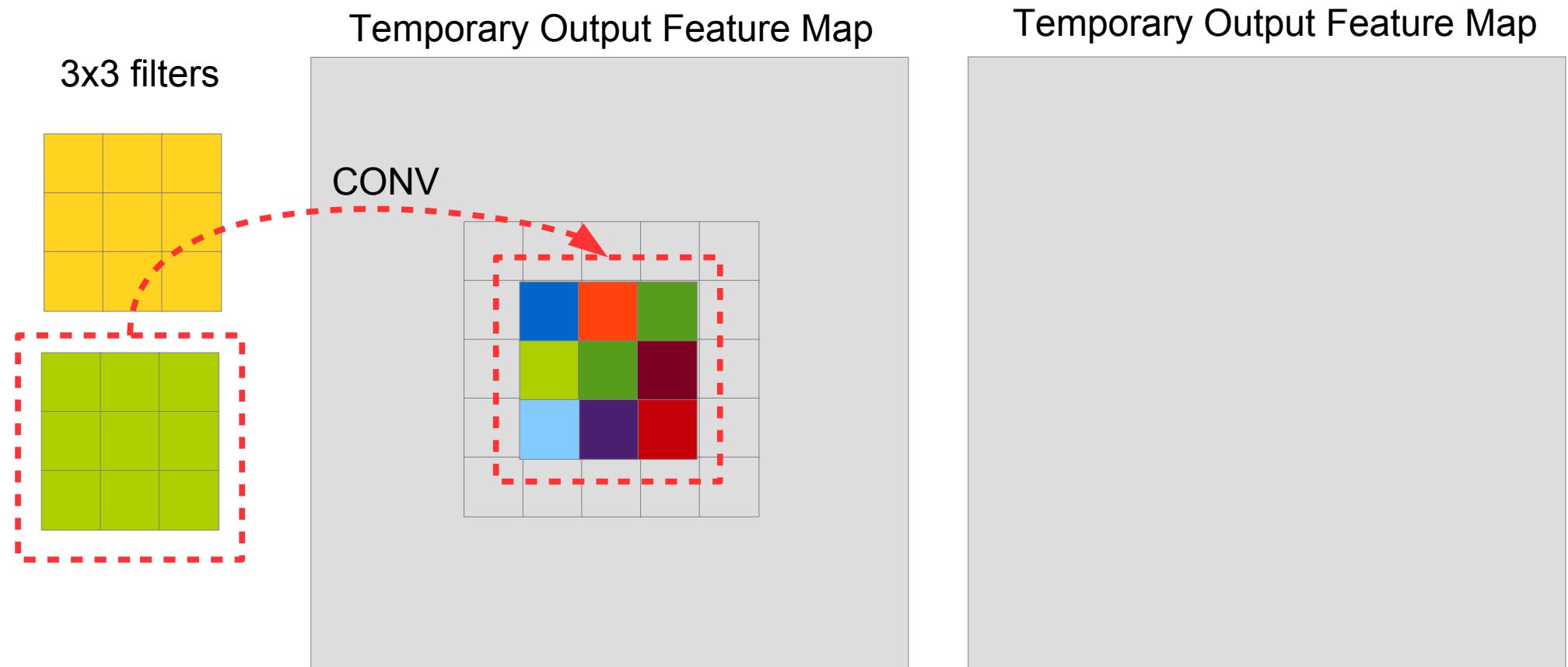


Temporary Output Feature Map



Stacked Filter

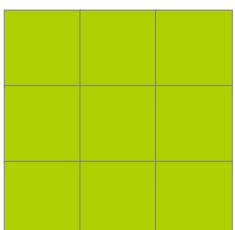
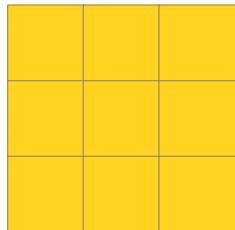
- Large filters built from multiple small filters
 - e.g., used in VGG-16



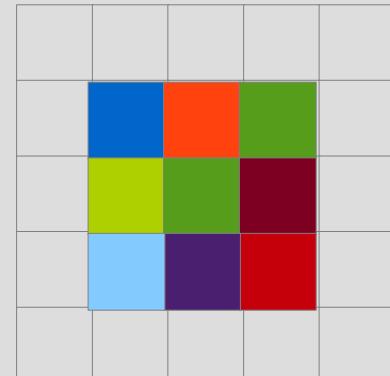
Stacked Filter

- Large filters built from multiple small filters
 - e.g., used in VGG-16

3x3 filters



Temporary Output Feature Map

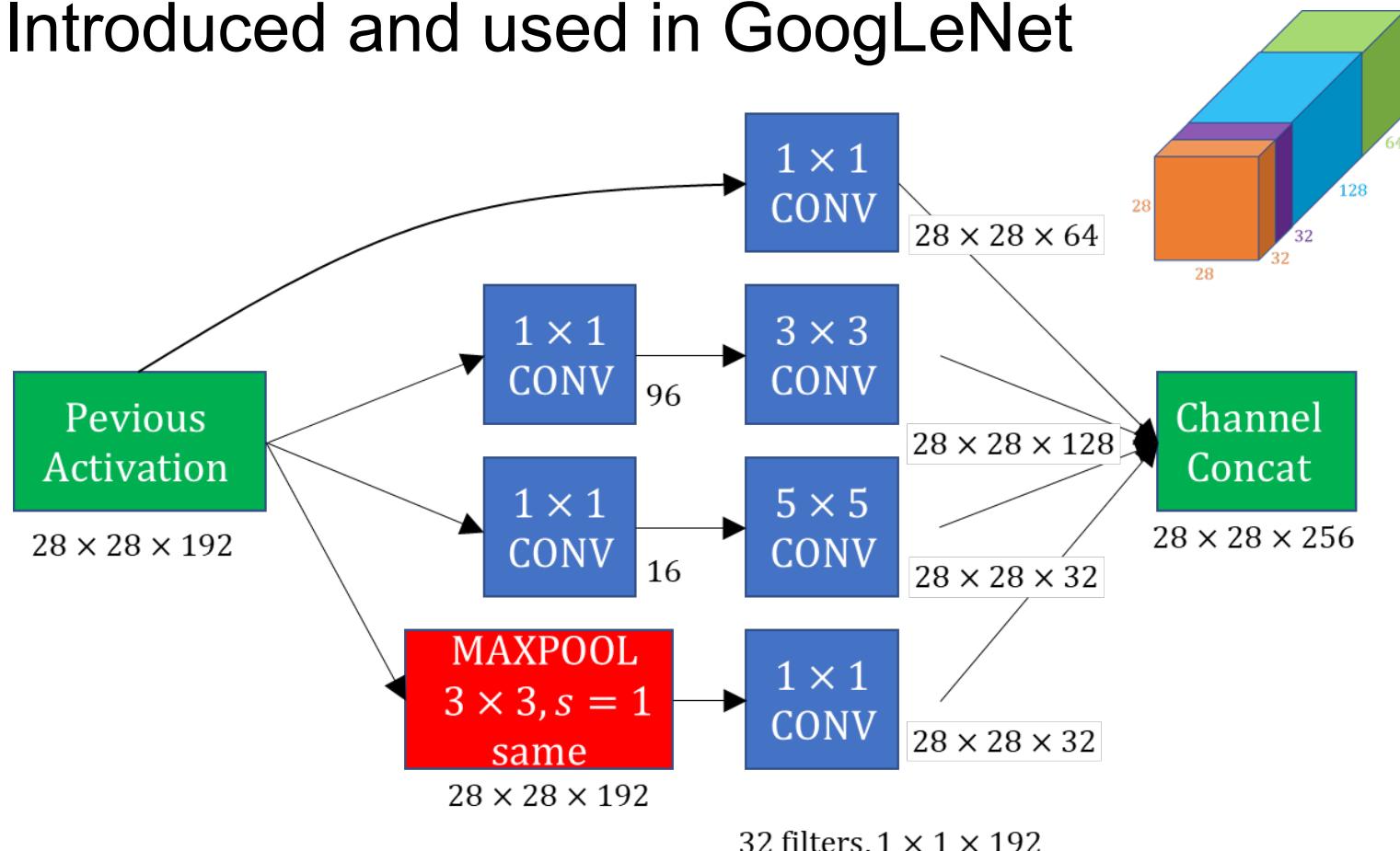


Temporary Output Feature Map



Inception Module

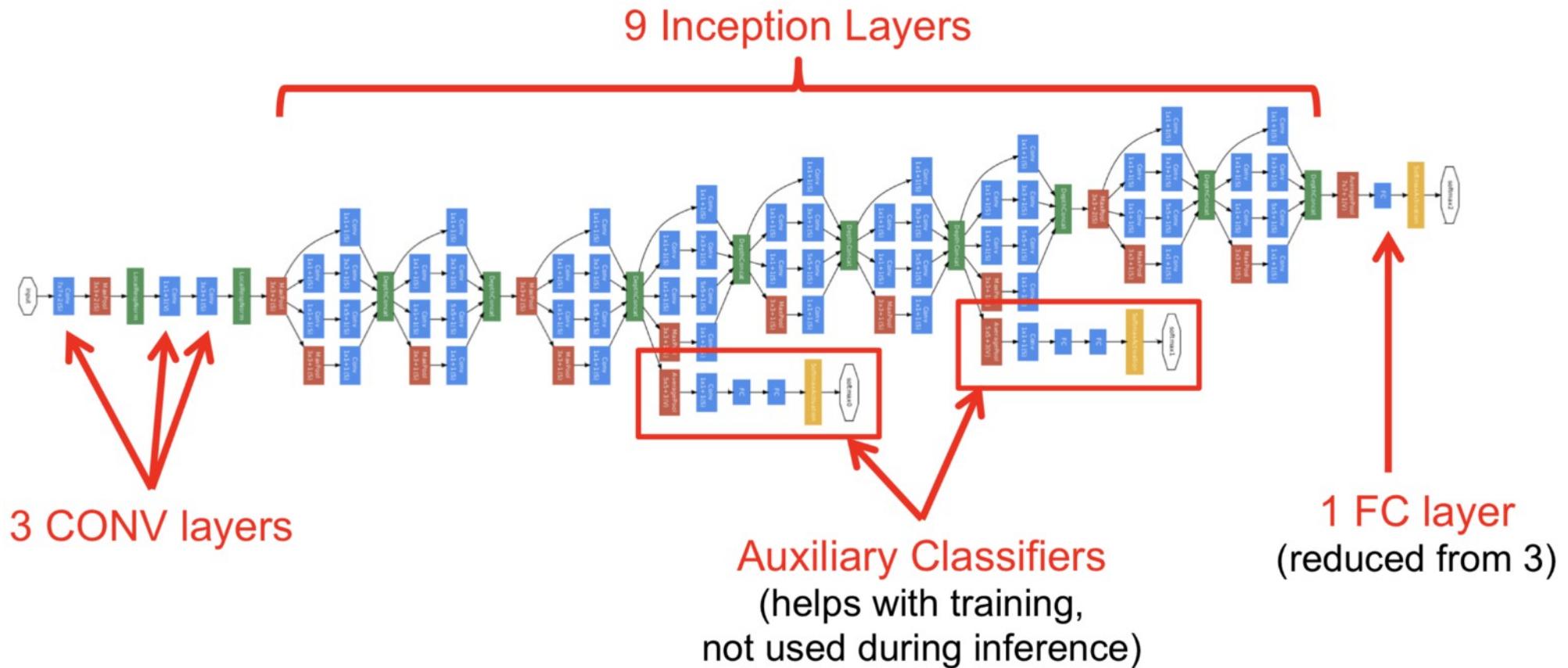
- Use multiple filter sizes for processing the input at multiple scales
 - Introduced and used in GoogLeNet



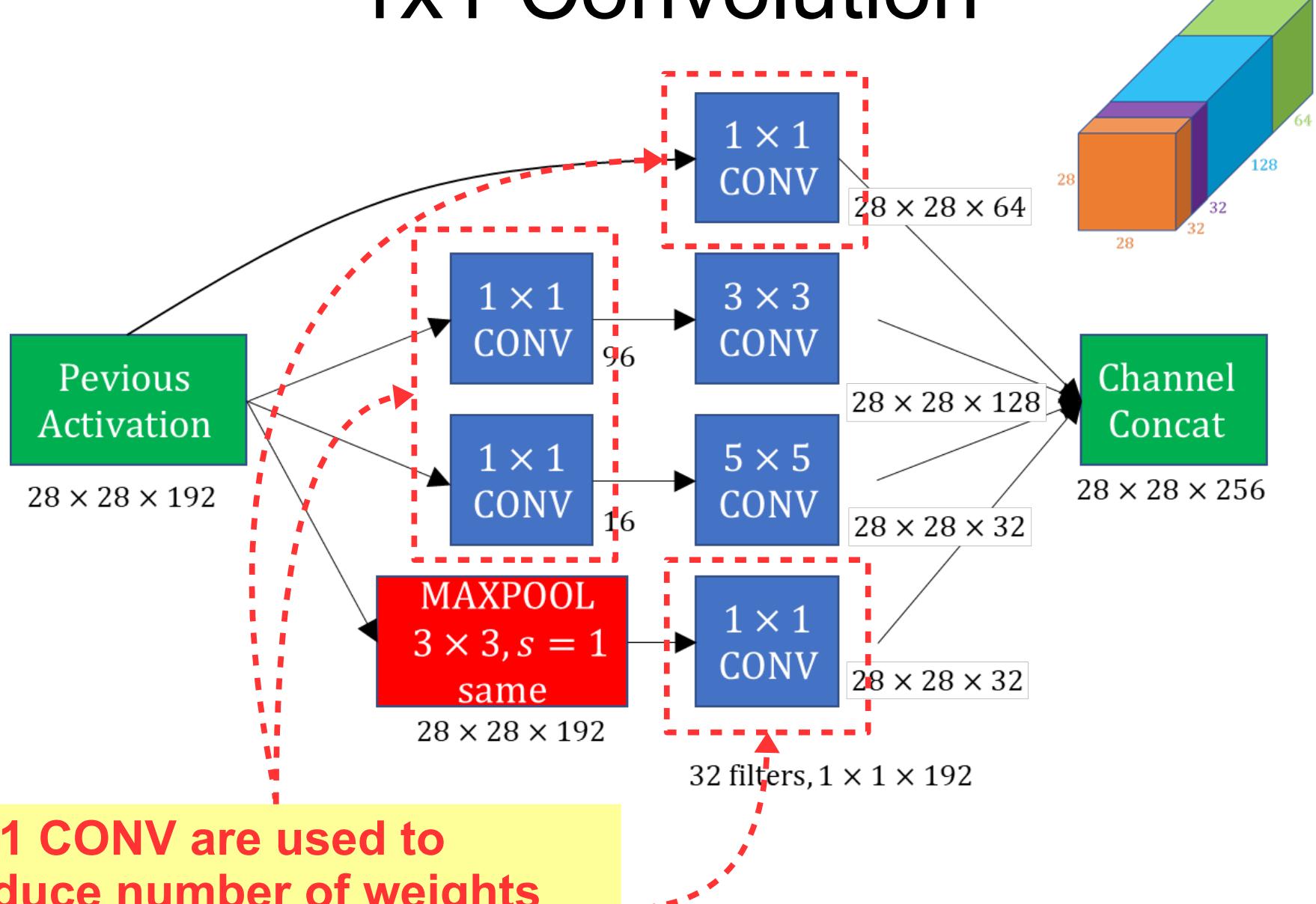
GoogLeNet/Inception (v1)

(Also, v2, v3 and v4)
ILSVRC14 Winner

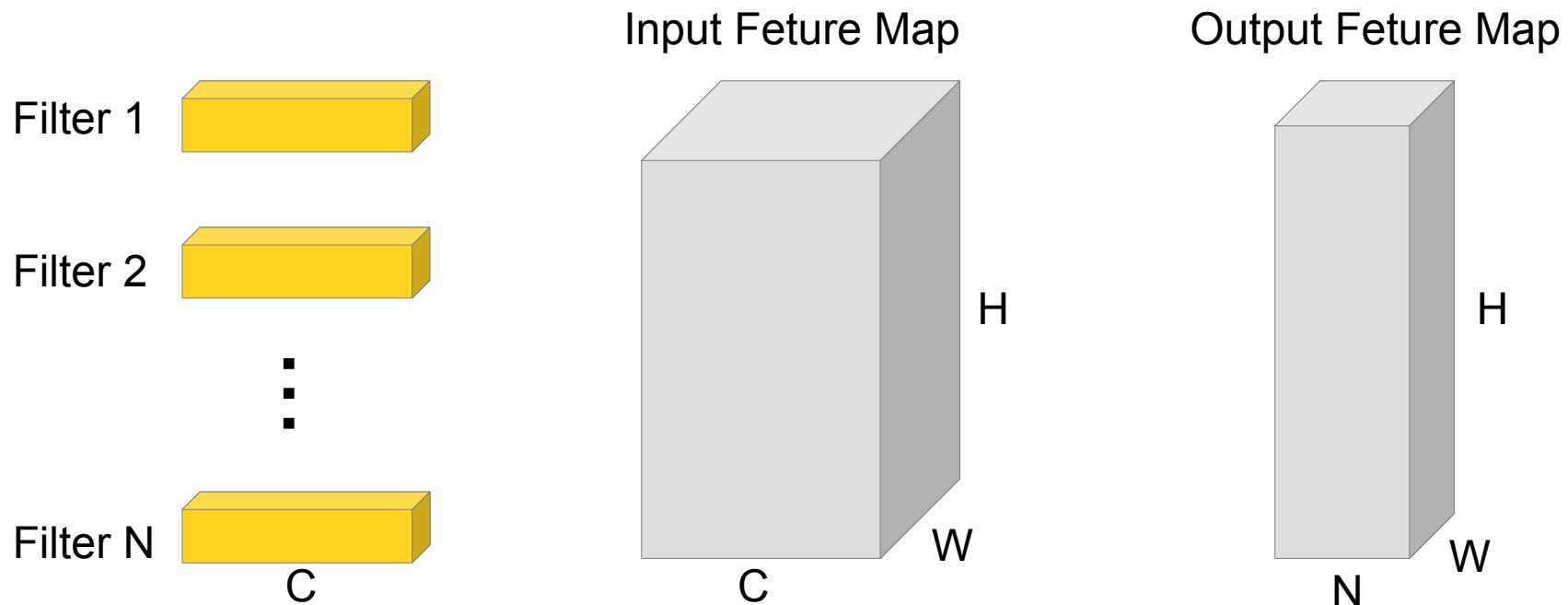
- 57 layers (21 CONV, 1 FC)
- 7M parameters, 1.43G MACs



1x1 Convolution



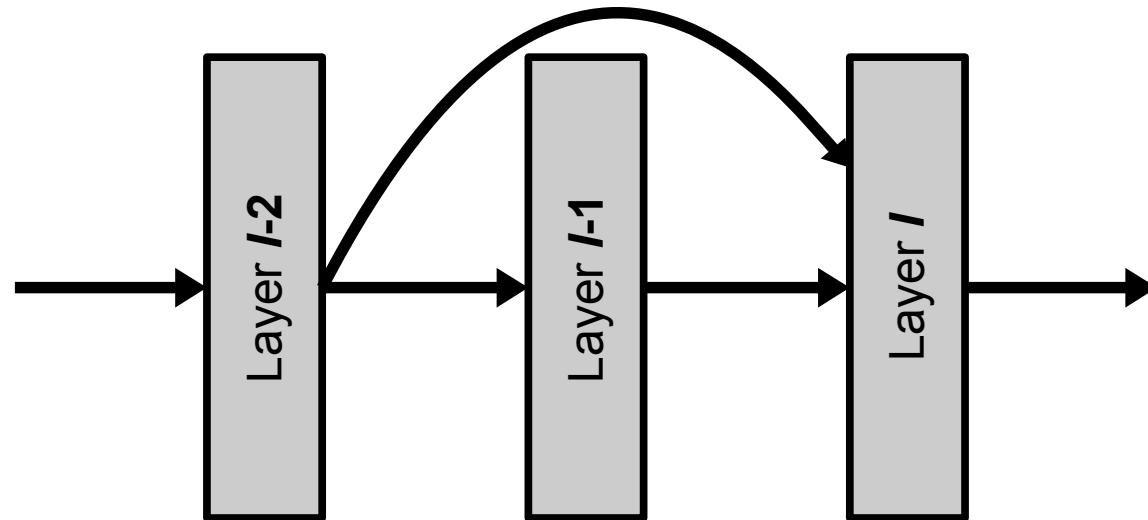
1X1 Convolution



- Capture cross-channel correlation
- Does not capture spatial correlation
- Reduces the number of channels in the next layer
 - If $N < C$

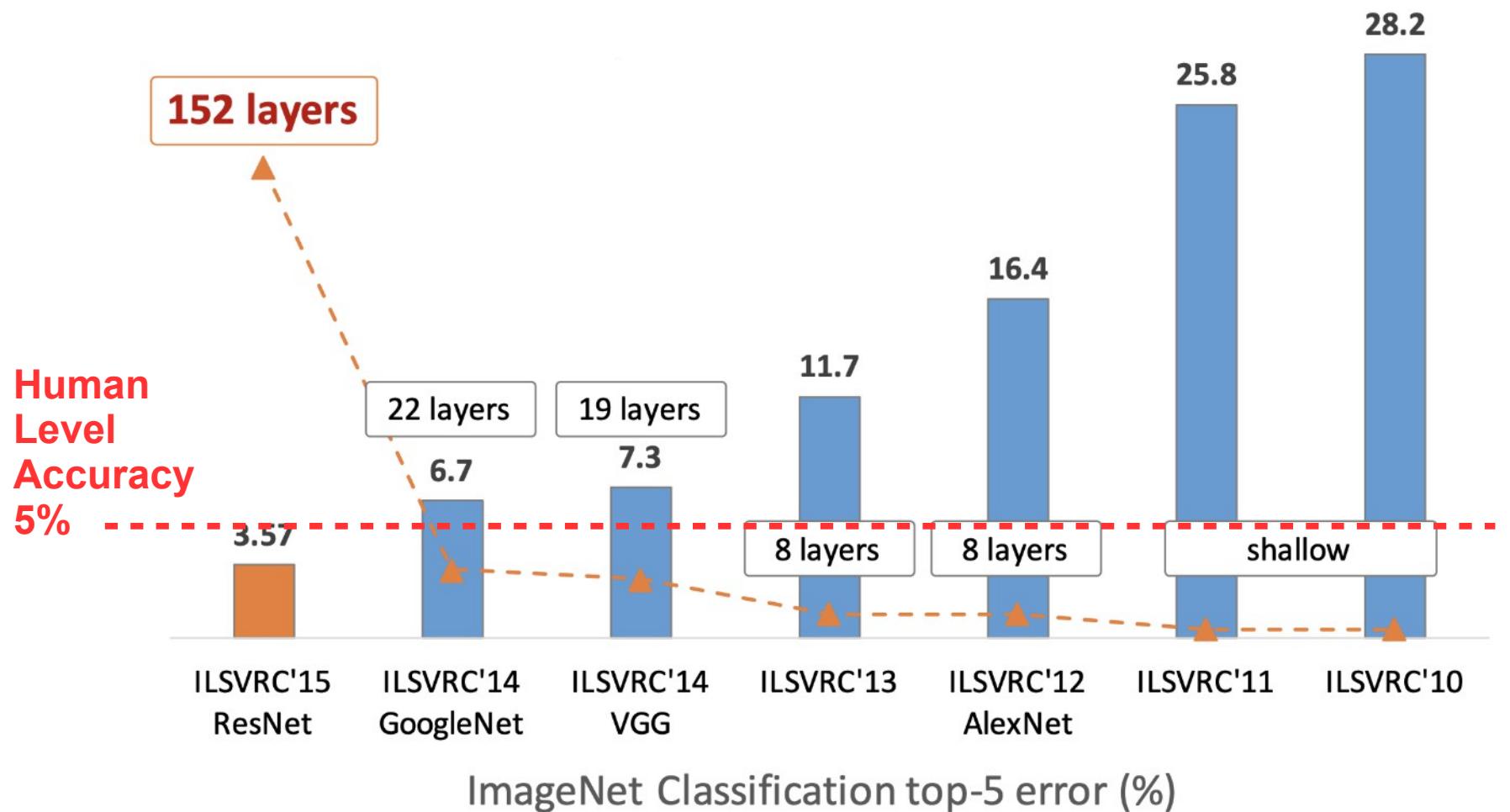
Residual Connection

- Feed-forward connection that connects to layers beyond the immediate next layer
- Improve the training of DNN with many layers
 - Address the *vanishing gradient* problem
- Introduced and used in ResNet



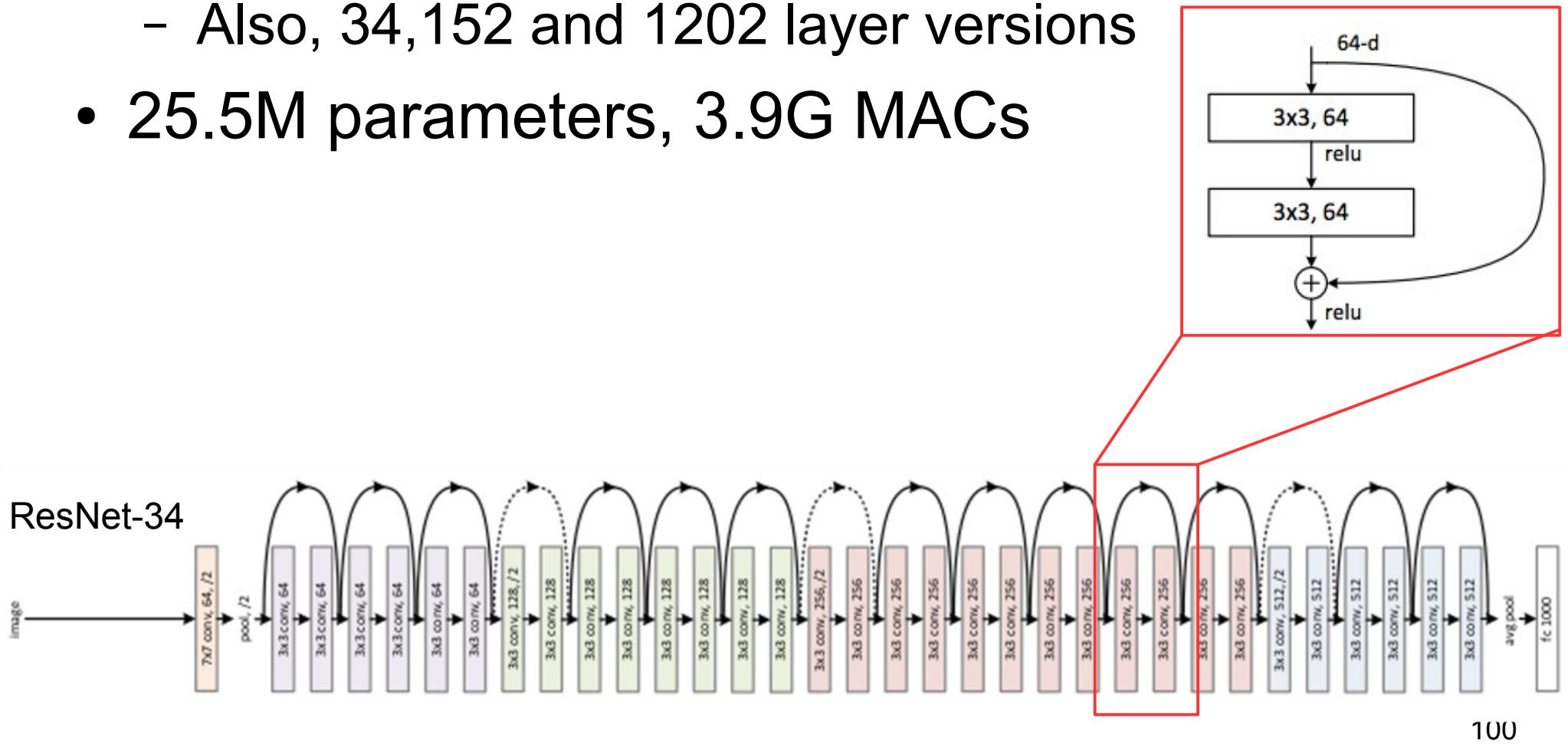
ResNet

- ILSVRC15 Winner (ResNet-152)
 - Better than human level accuracy!



ResNet-50

- 50 layers (49 CONV, 1 FC)
 - Also, 34, 152 and 1202 layer versions
 - 25.5M parameters, 3.9G MACs



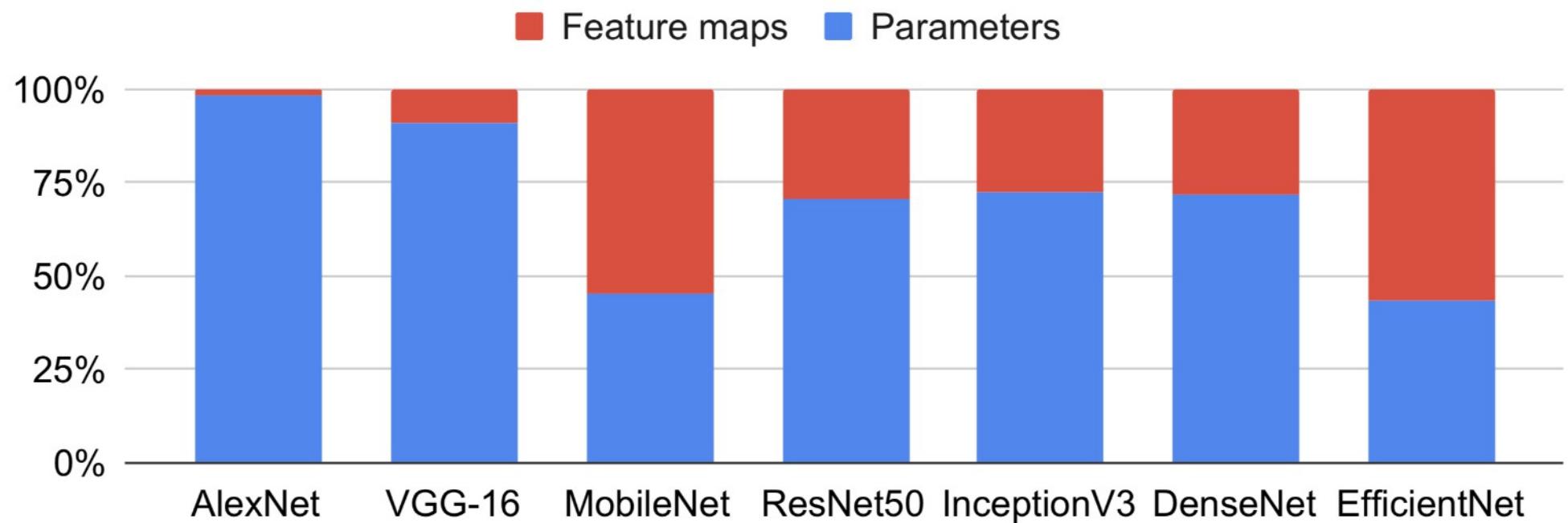
Trends

- Increase of the depth
- Increase in the number of filter shape
 - Flexibility is important
- Most of the computation on CONV layers
- Most of the parameters in FC layers
 - Even if FC layers have been replacing by CONV

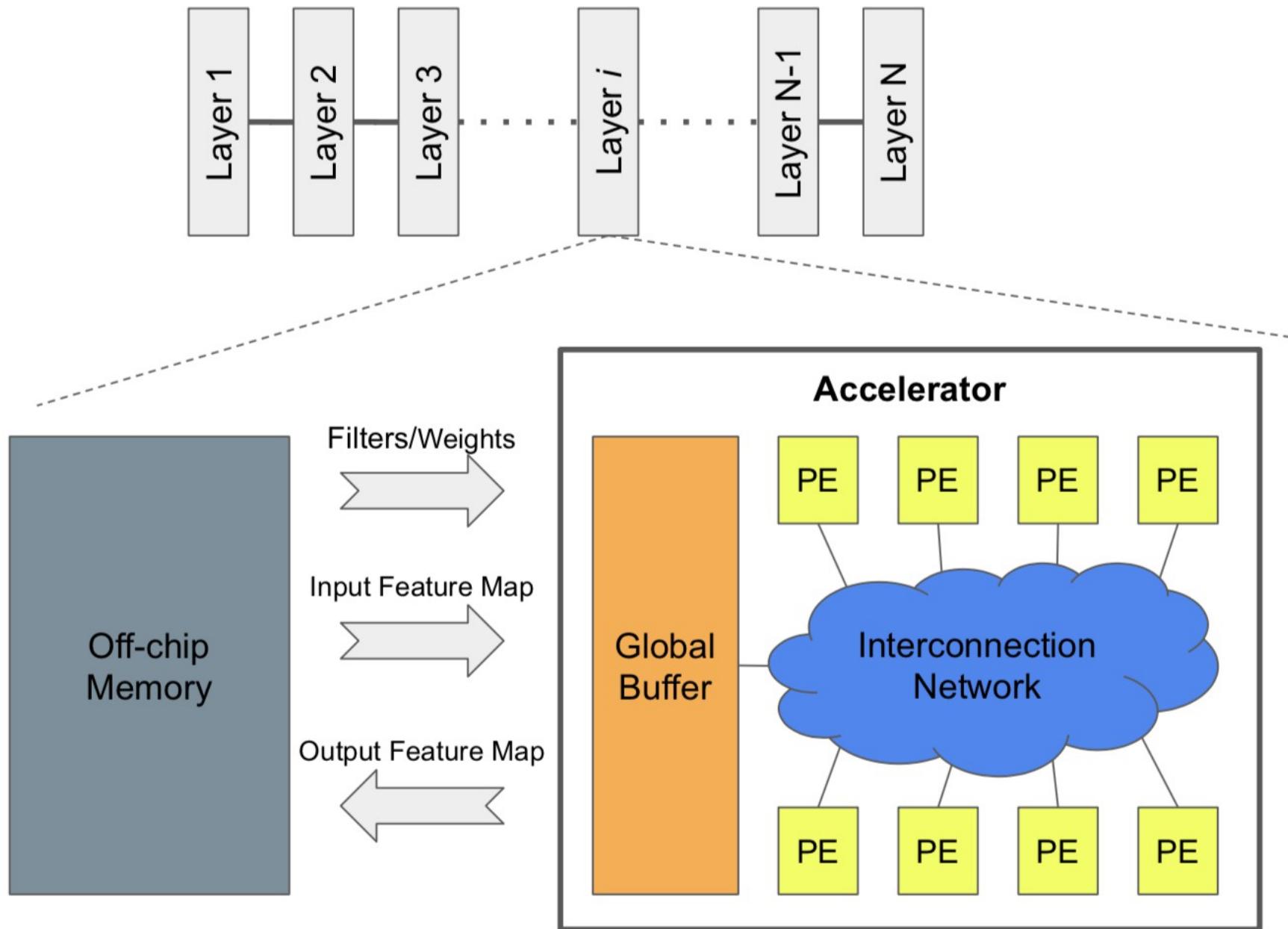
Agenda

- Deep Neural Networks
- Types of layers
- Memory and communication traffic
- Kernel computation
- Quantitative analysis

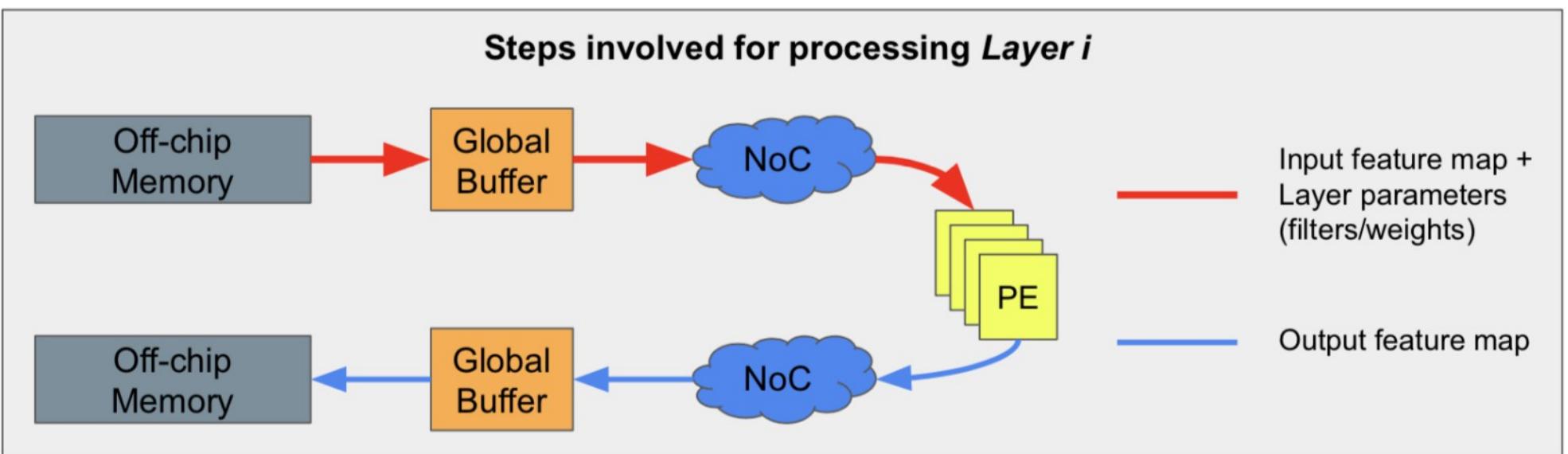
Feature Maps vs. Parameters



DNN Layer Computation



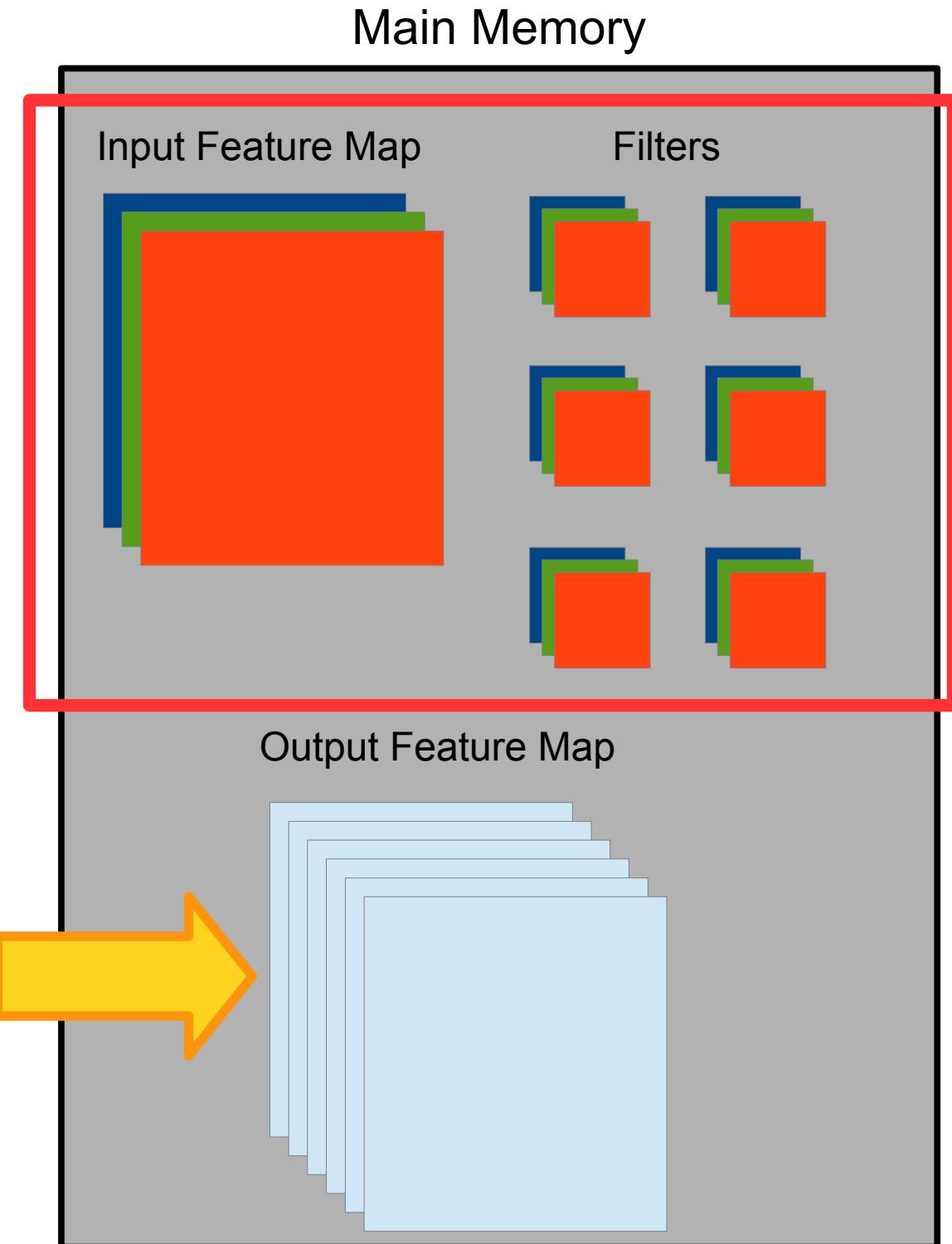
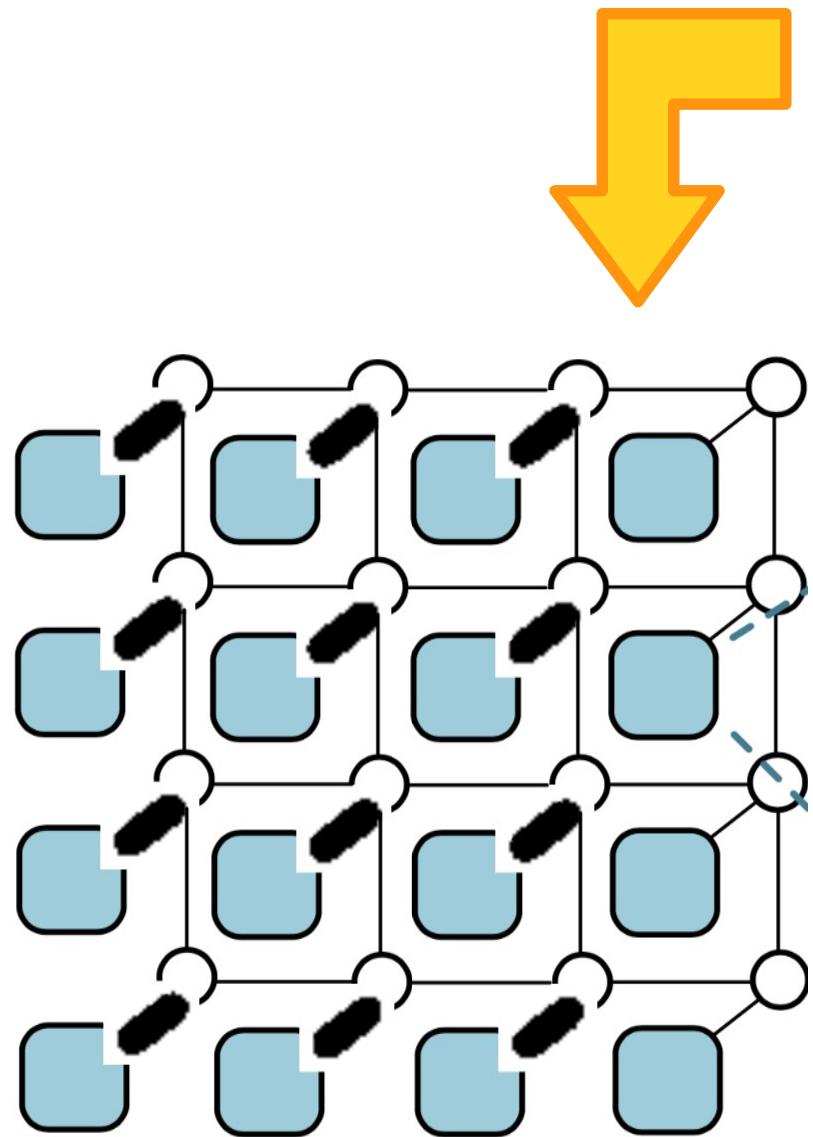
DNN Layer Computation



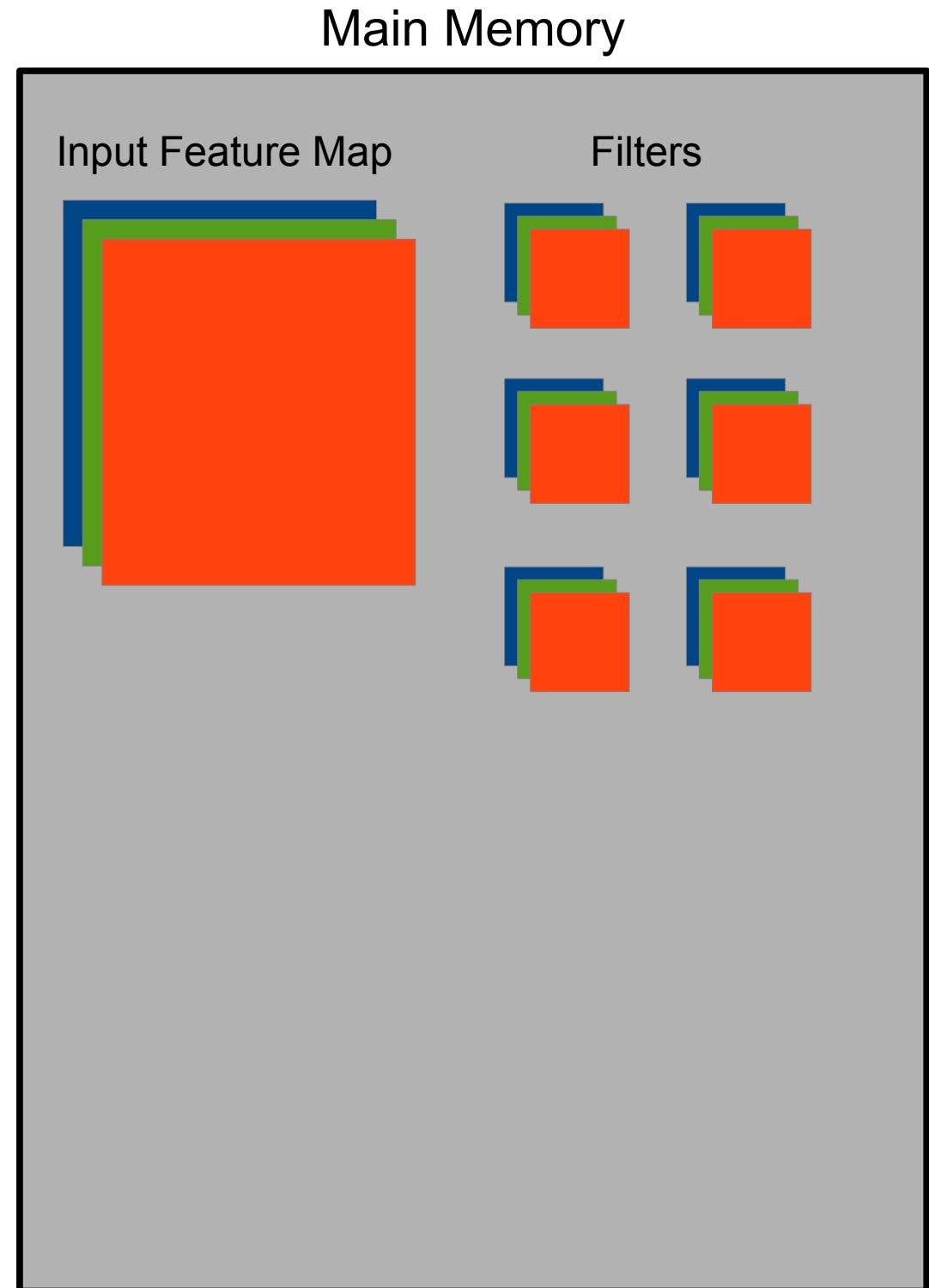
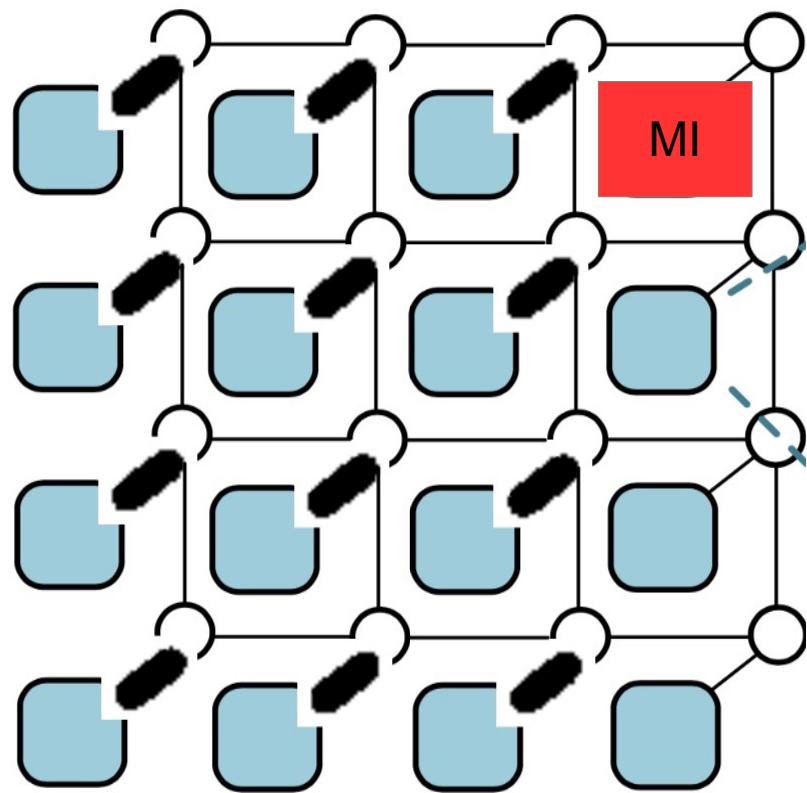
Induced Traffic

- From/to the memory
- Into the accelerator

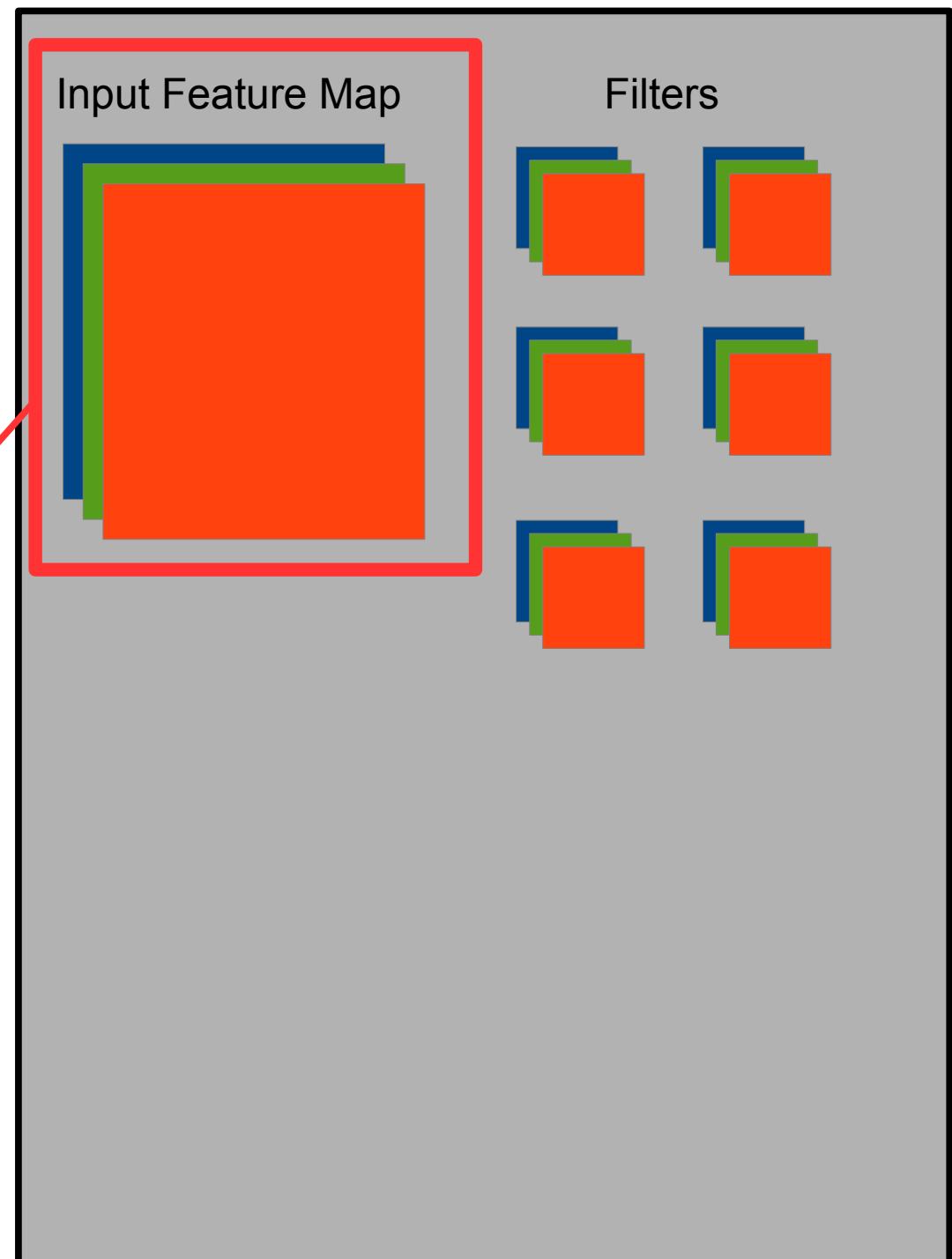
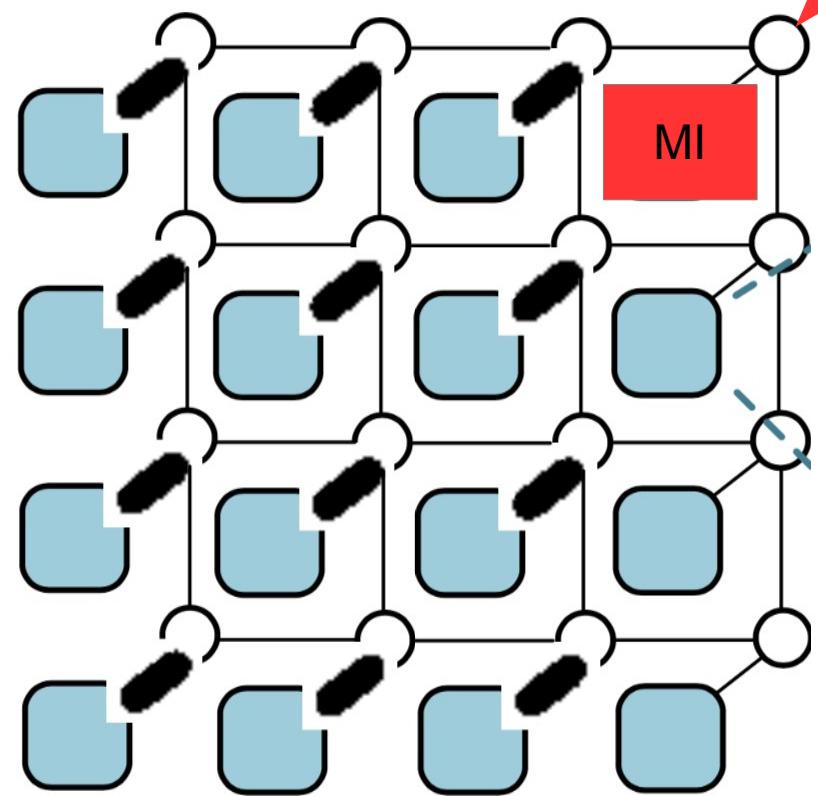
Convolution



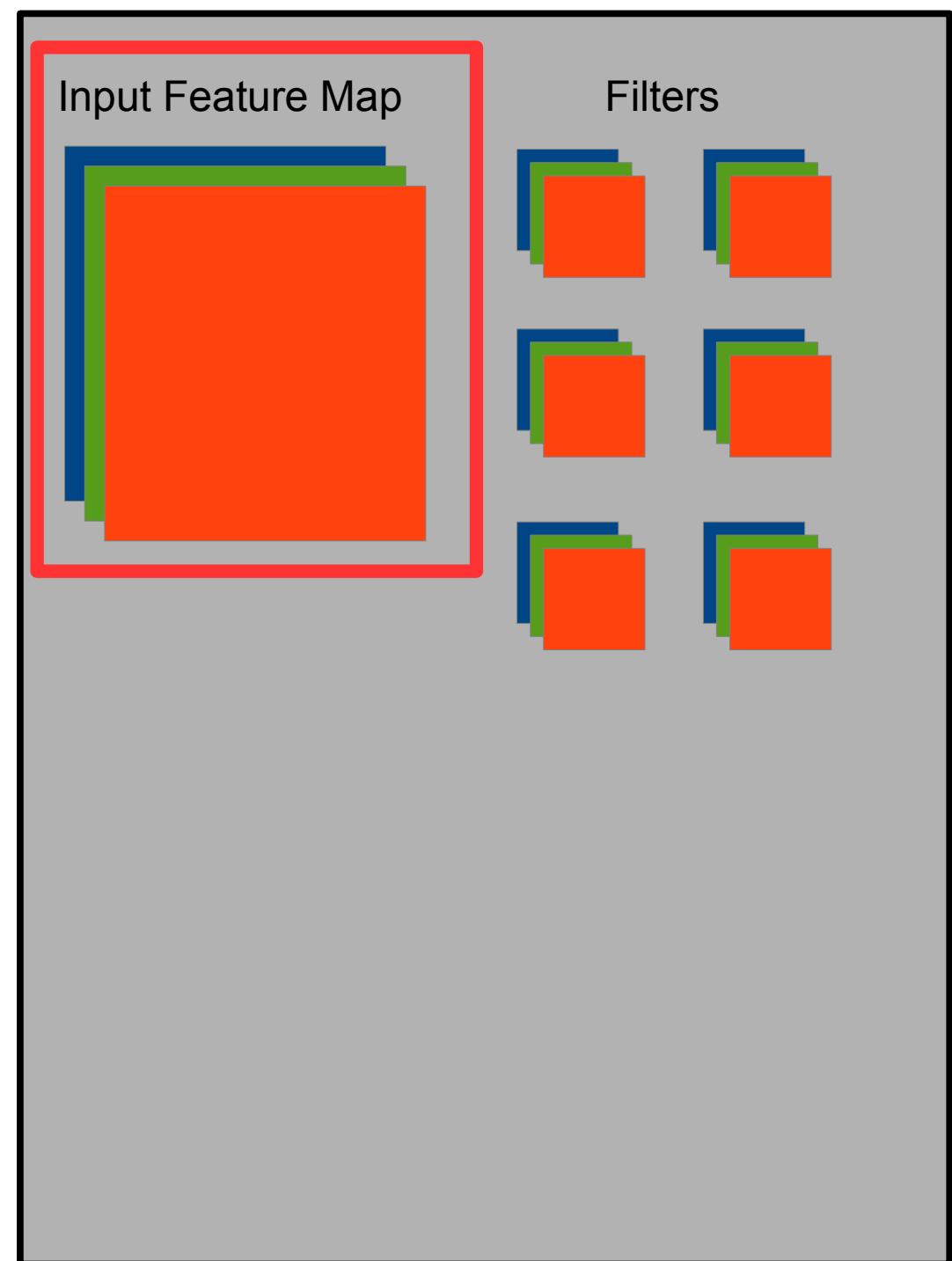
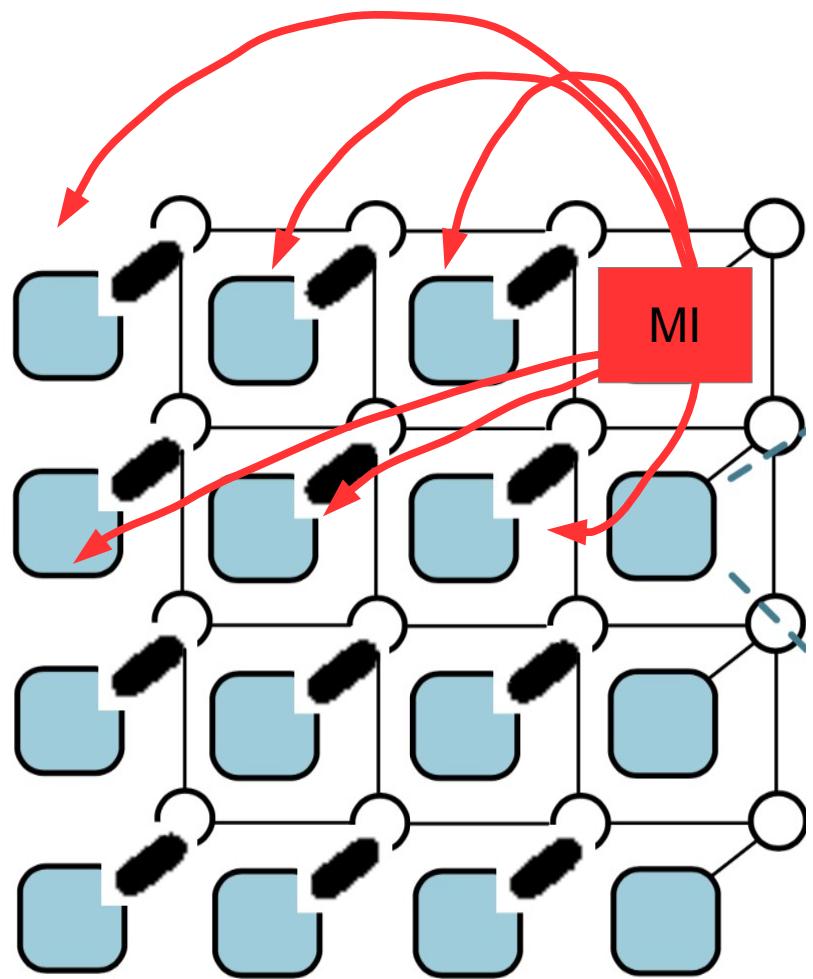
Convolution



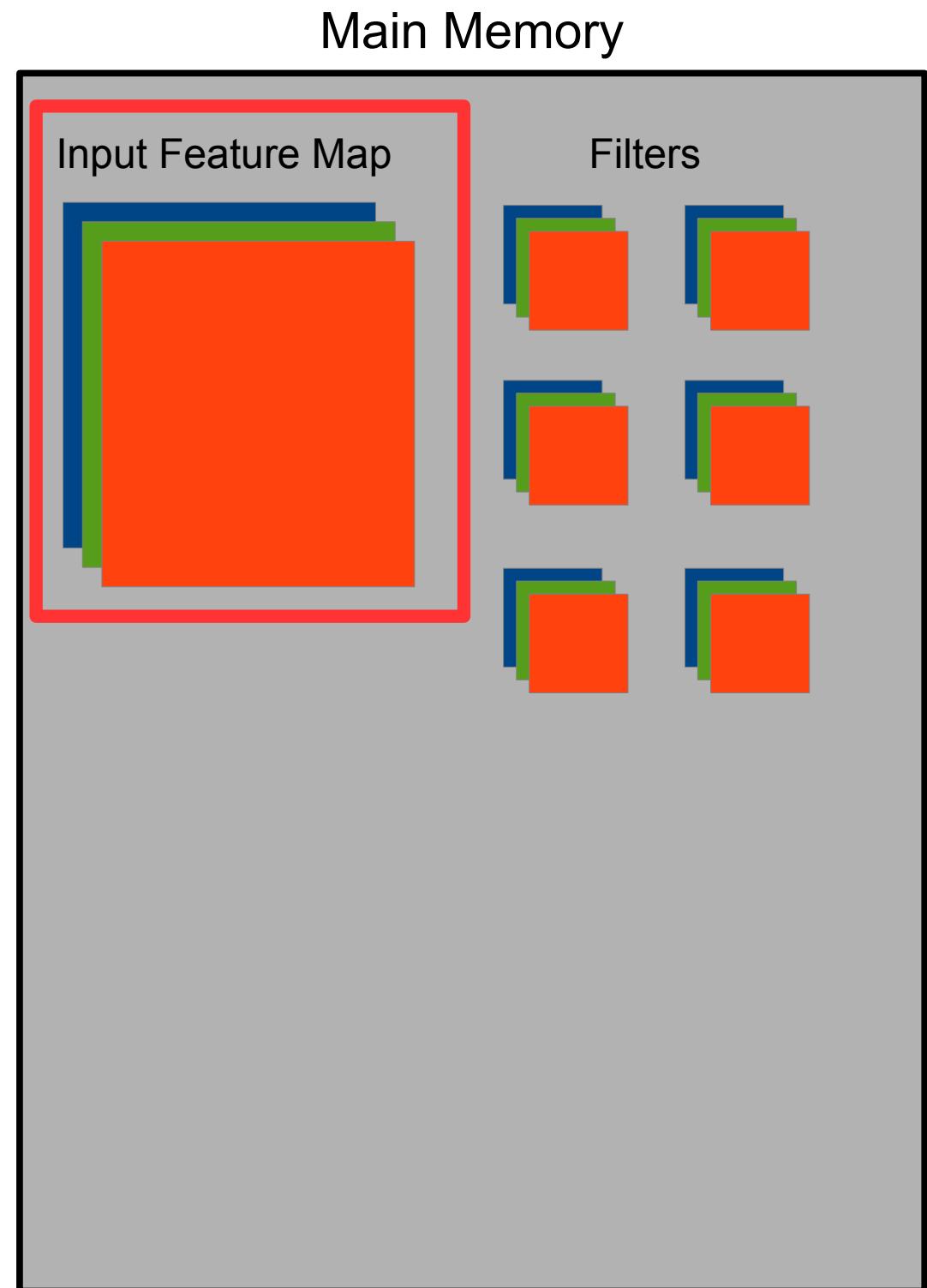
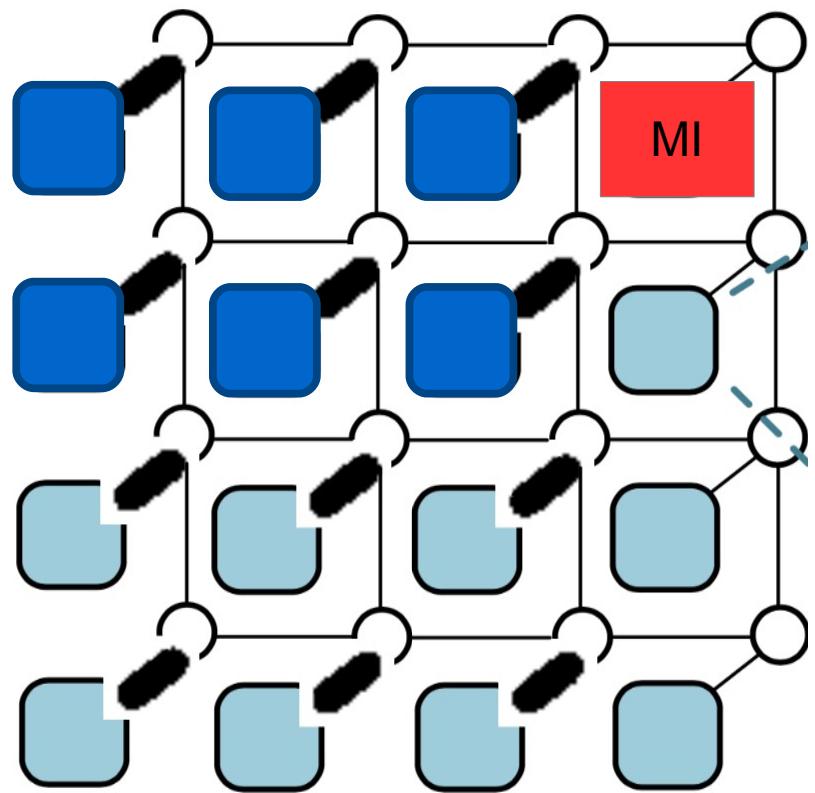
Convolution



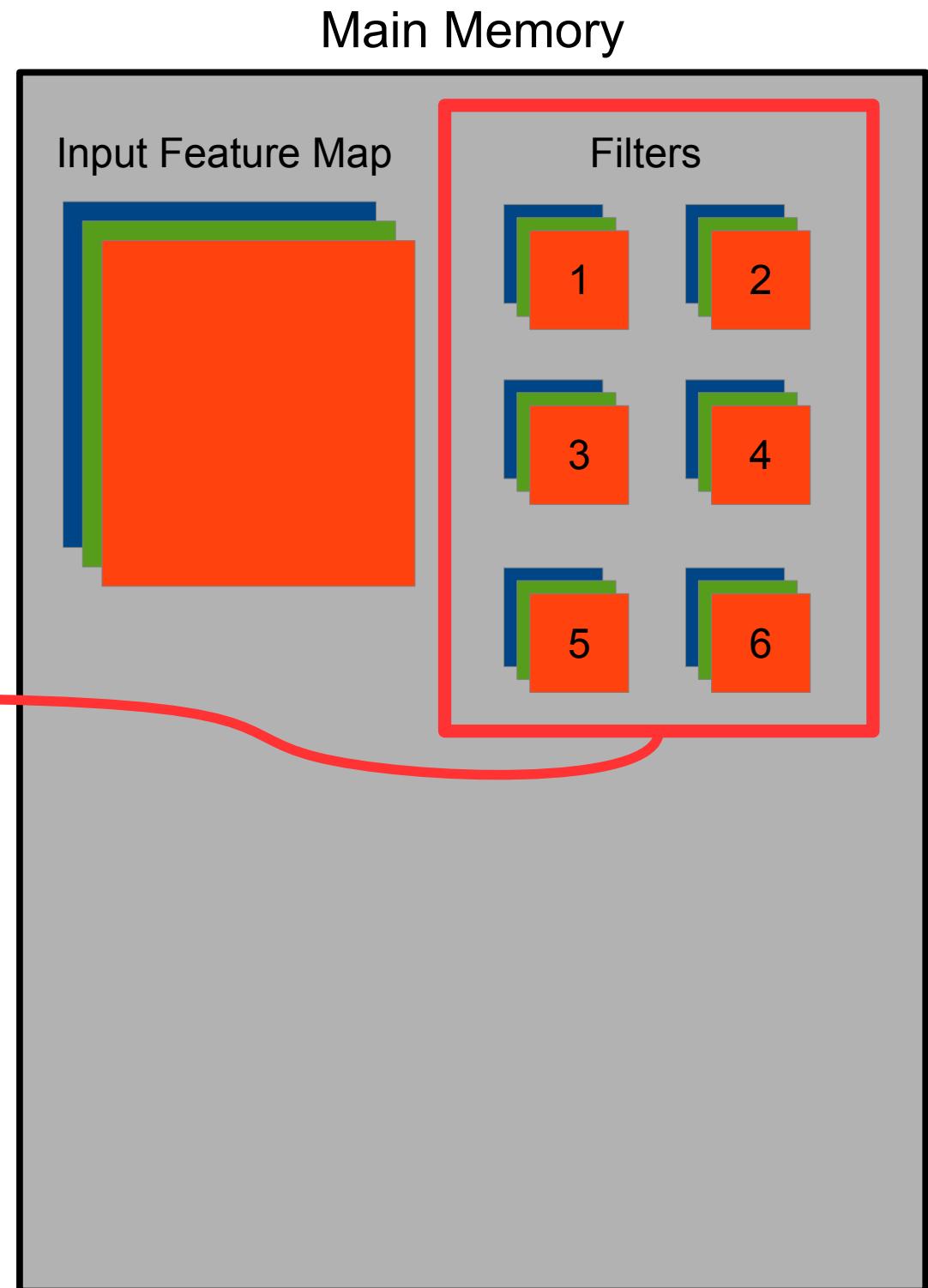
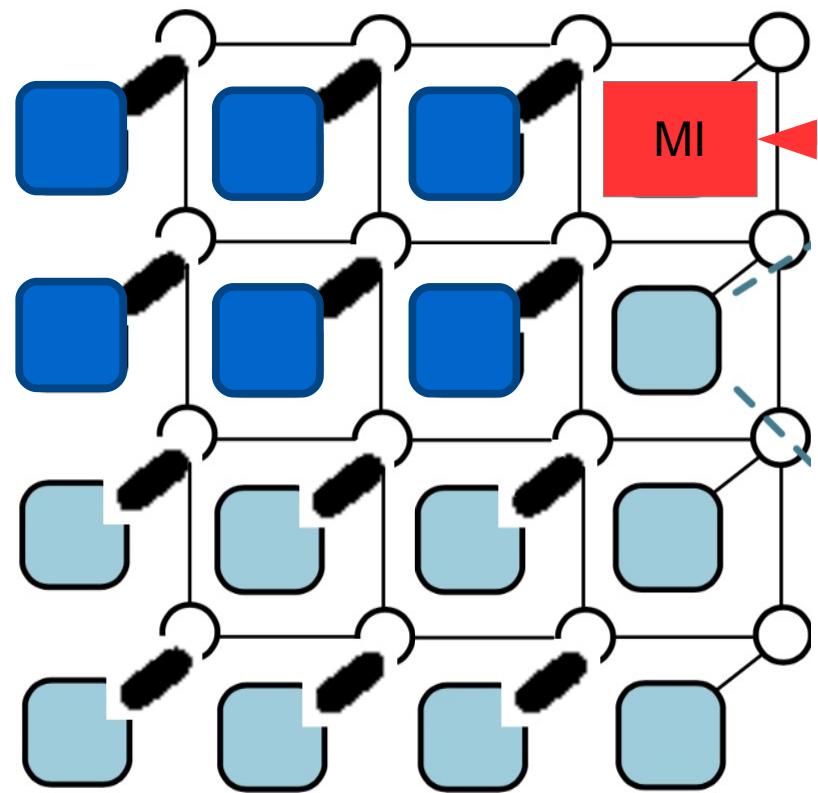
Convolution



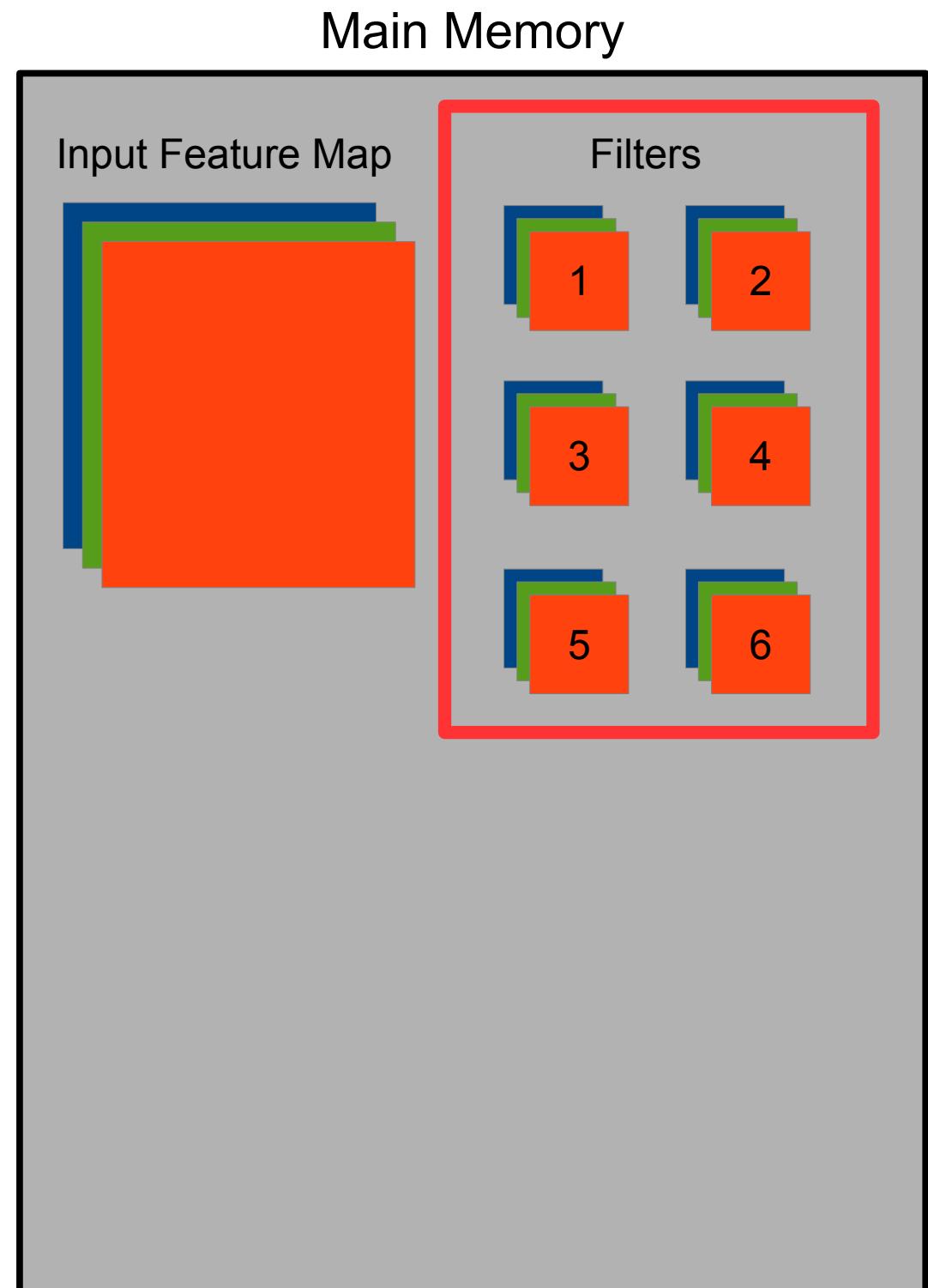
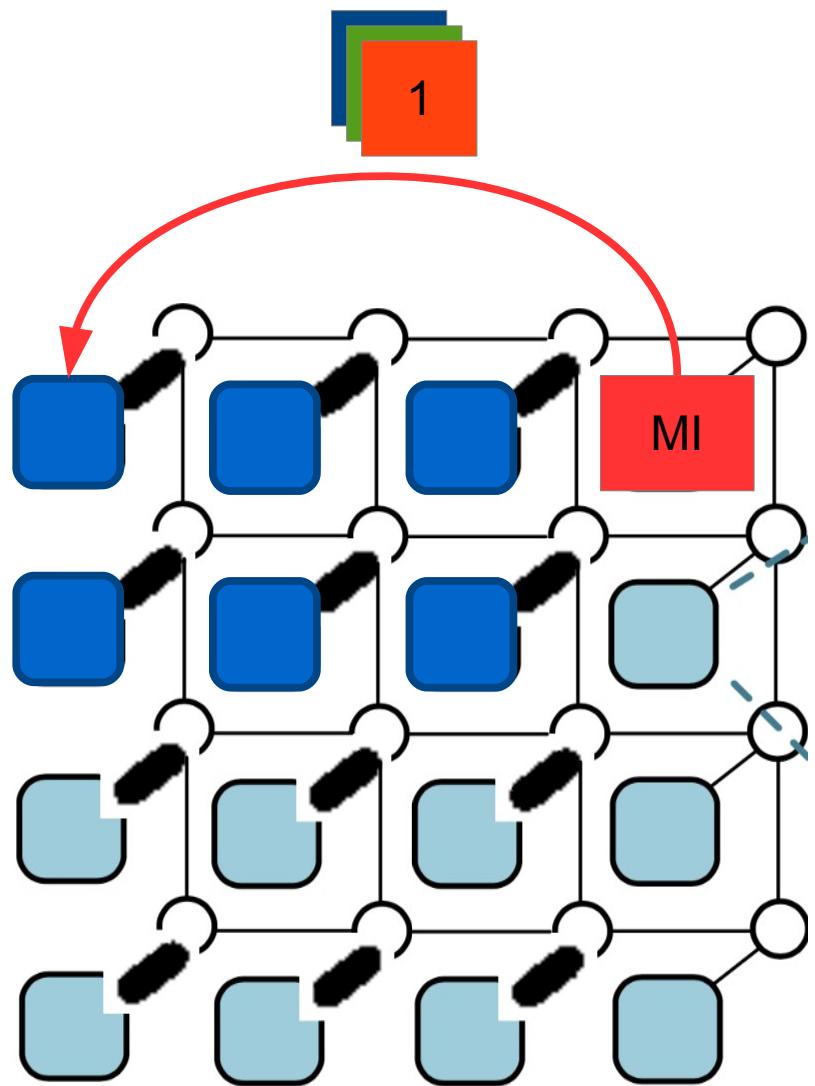
Convolution



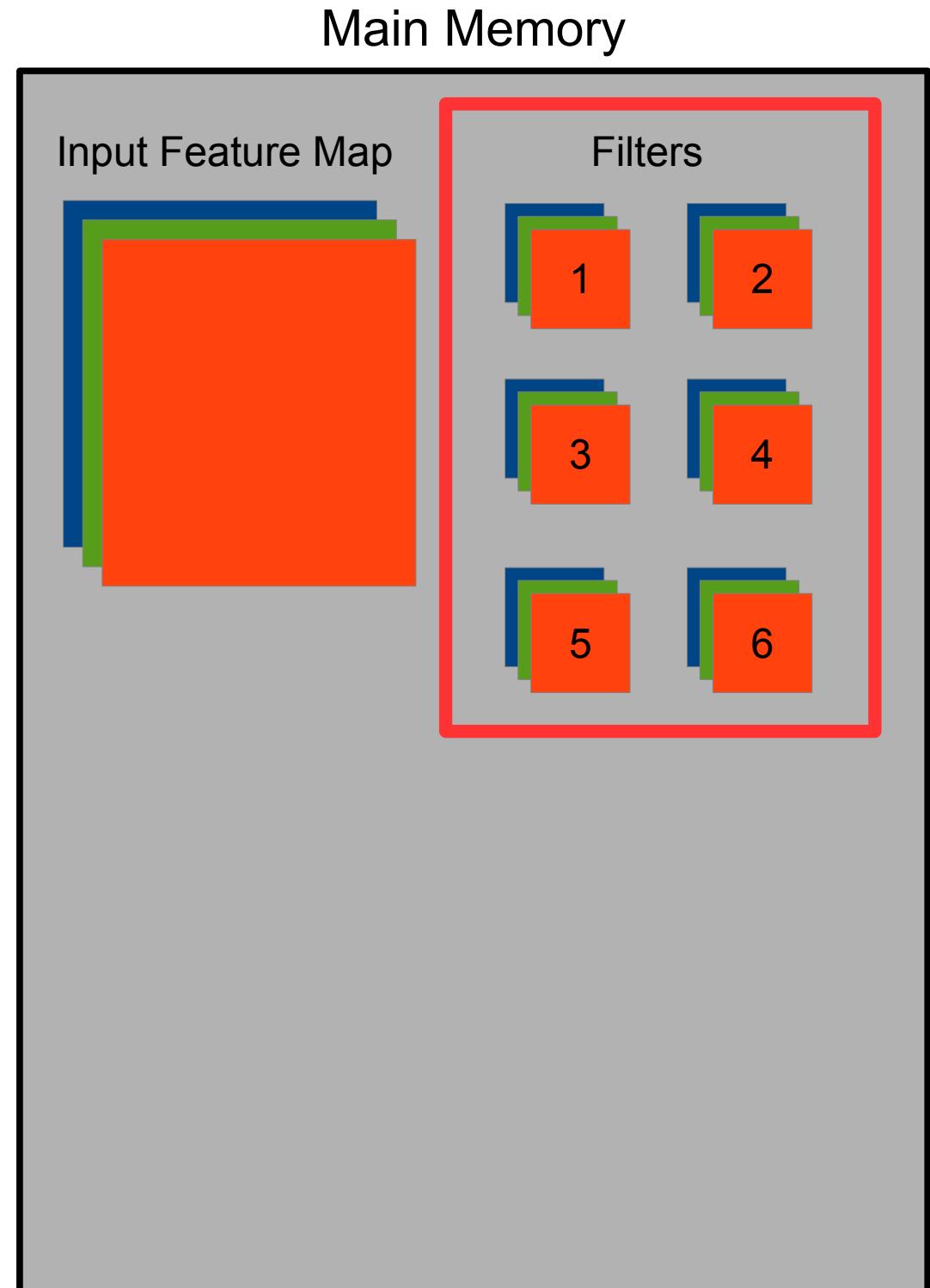
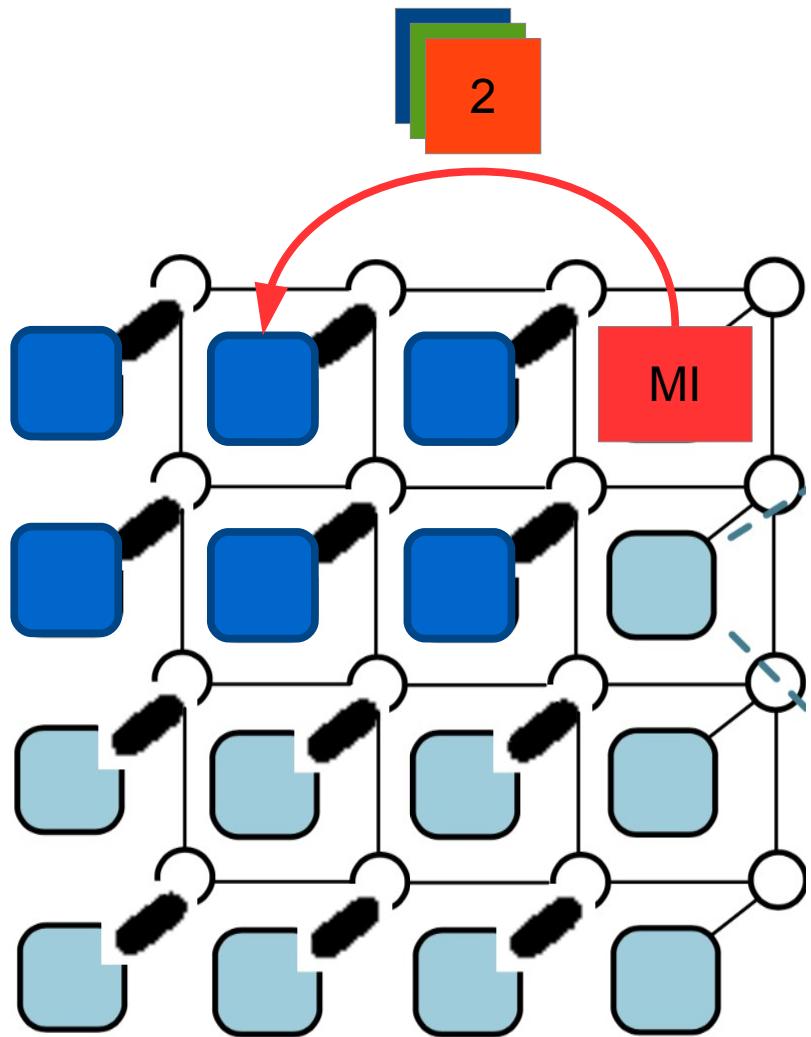
Convolution



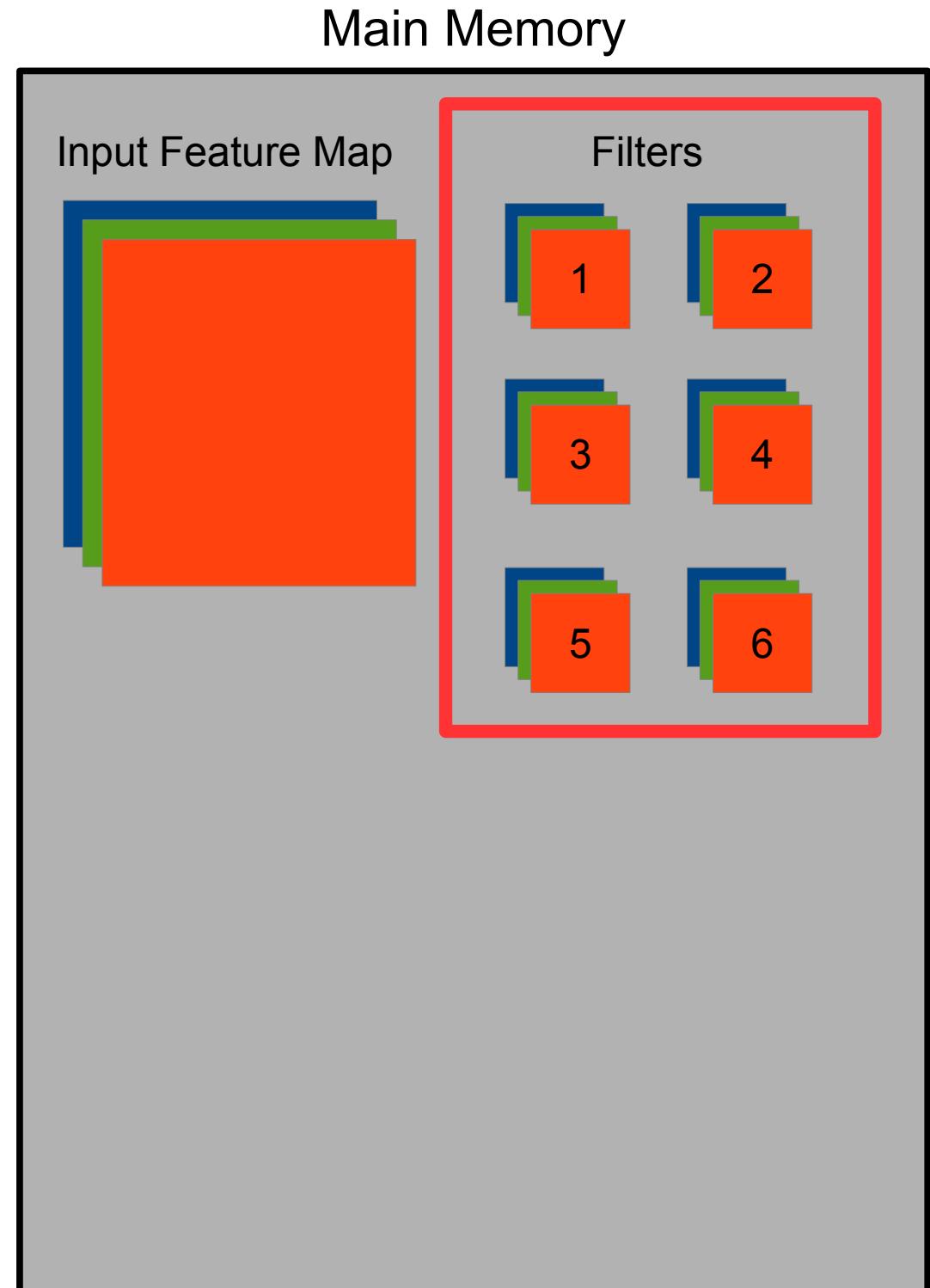
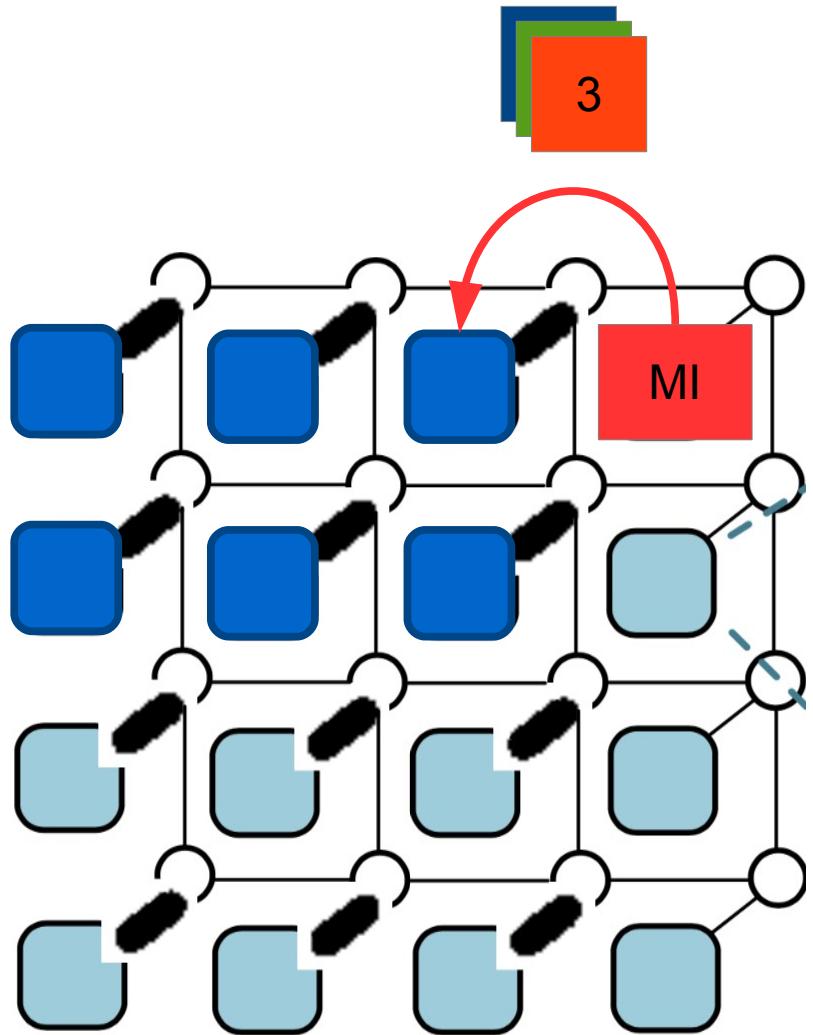
Convolution



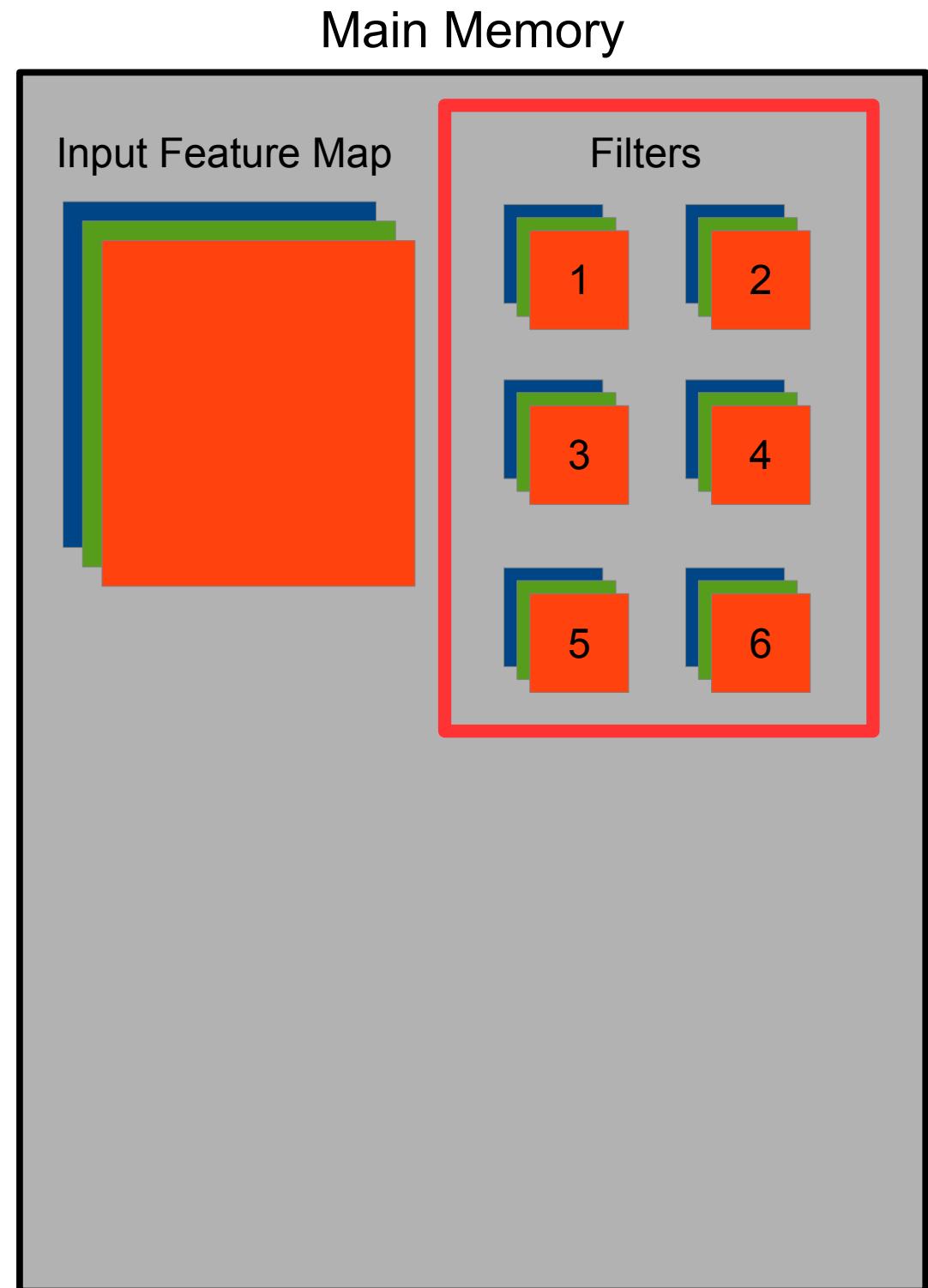
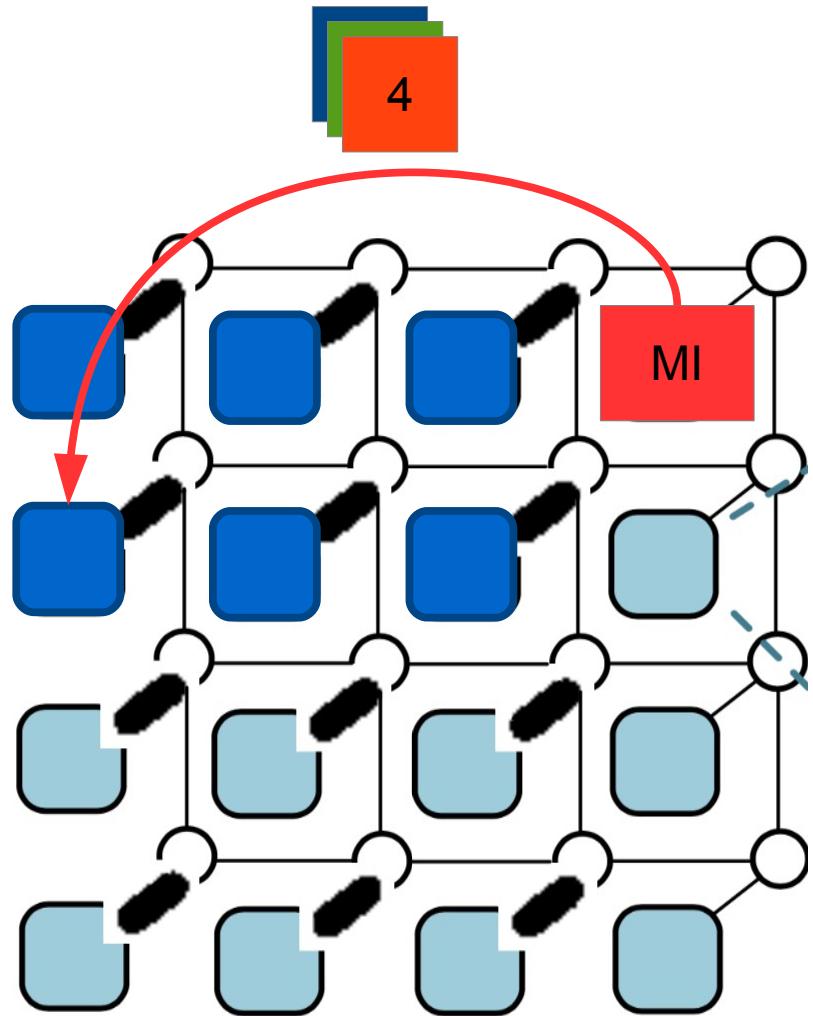
Convolution



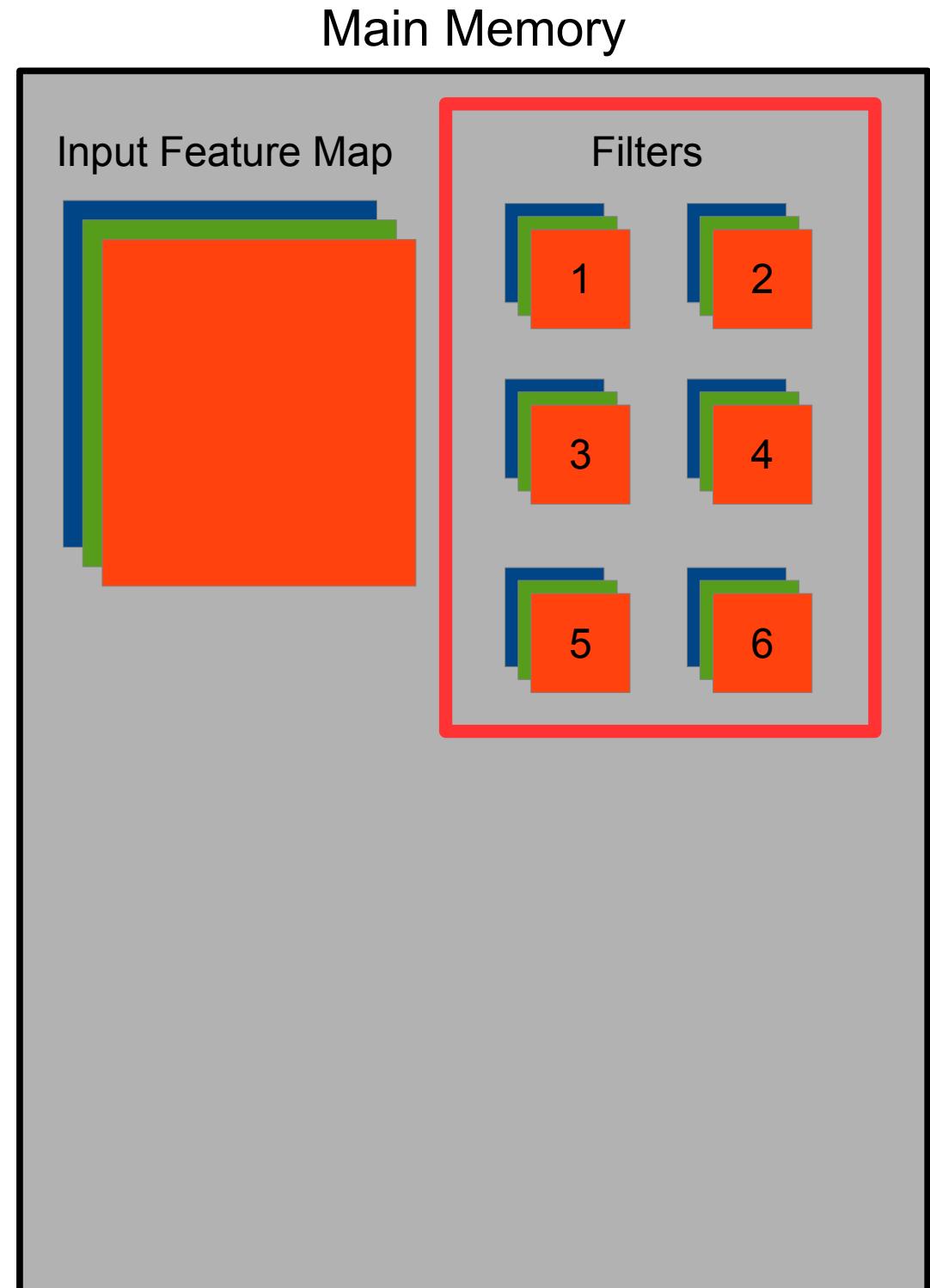
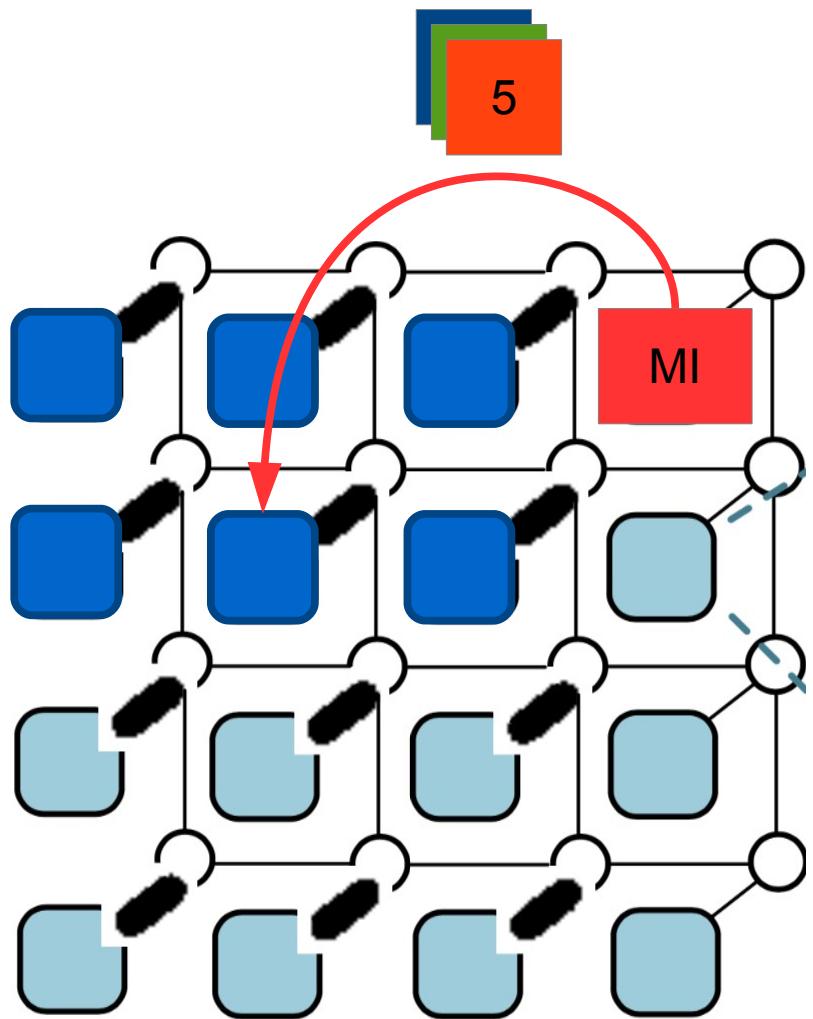
Convolution



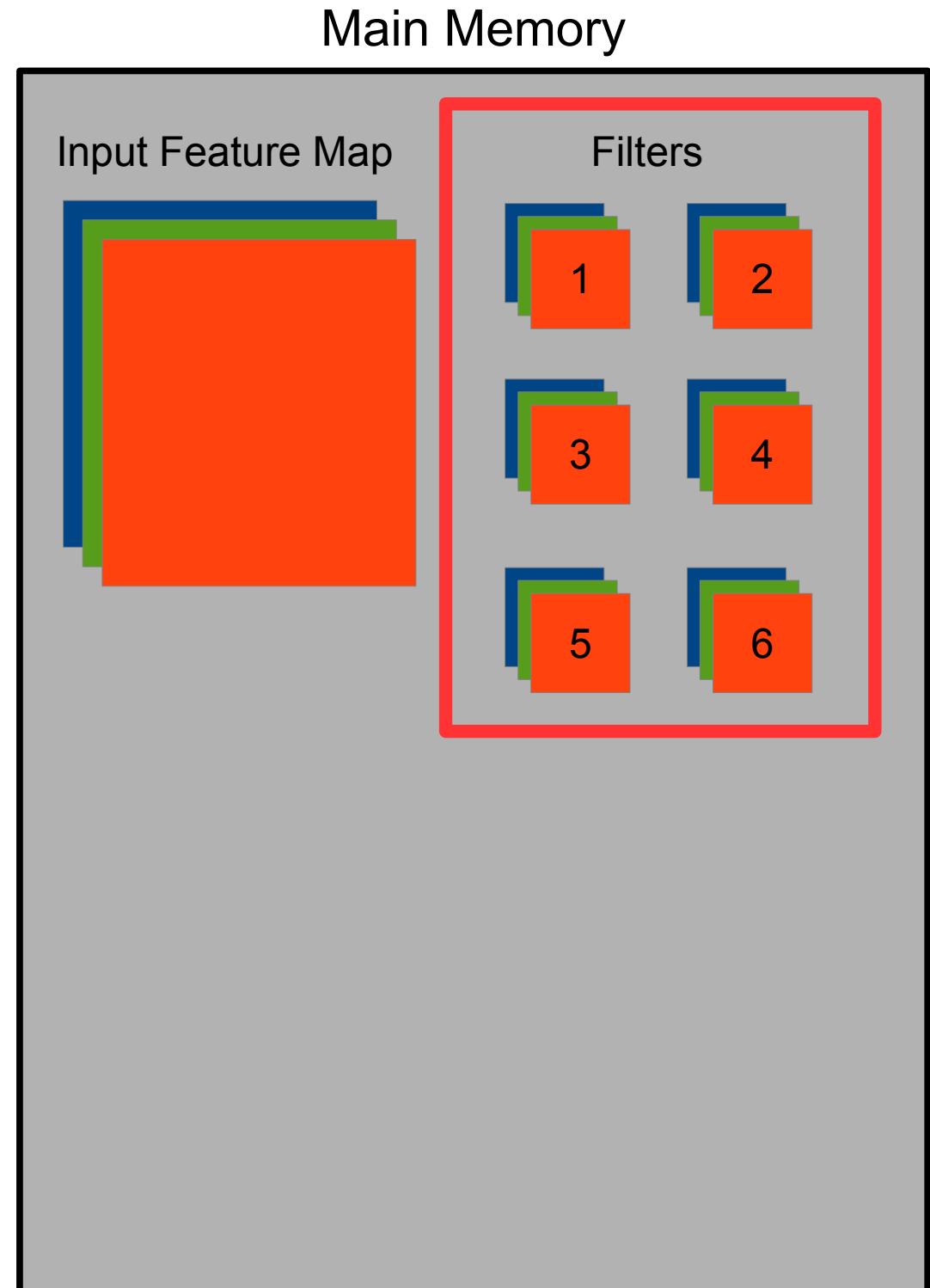
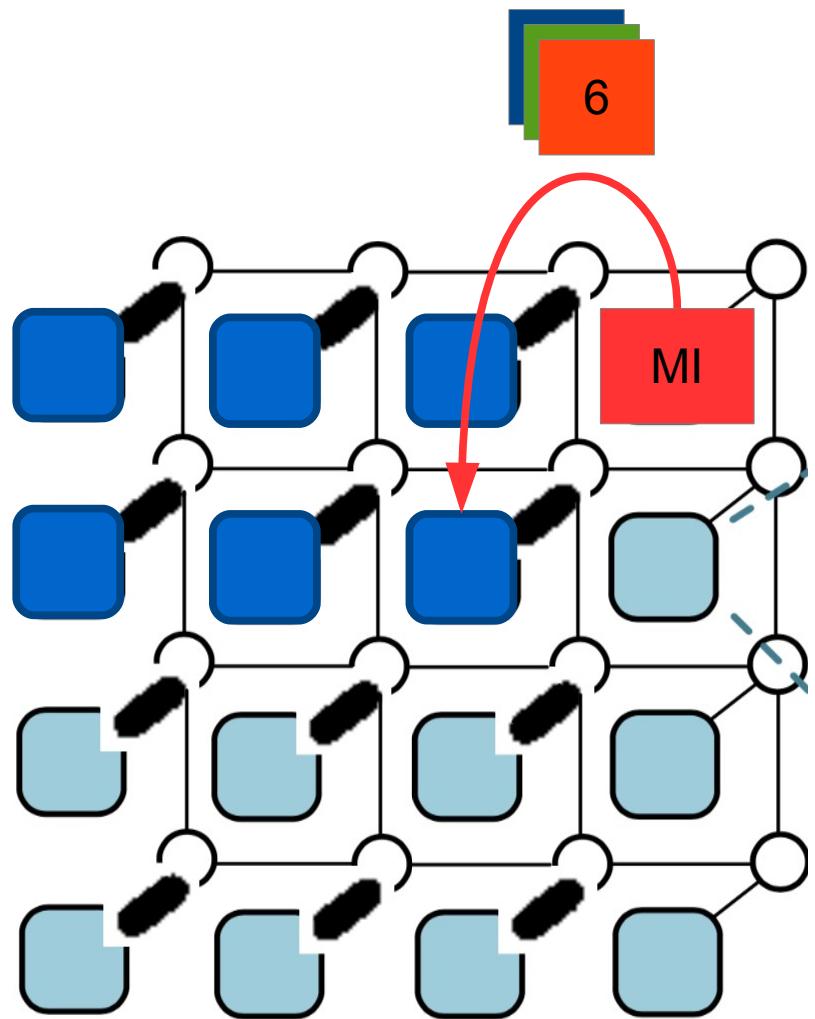
Convolution



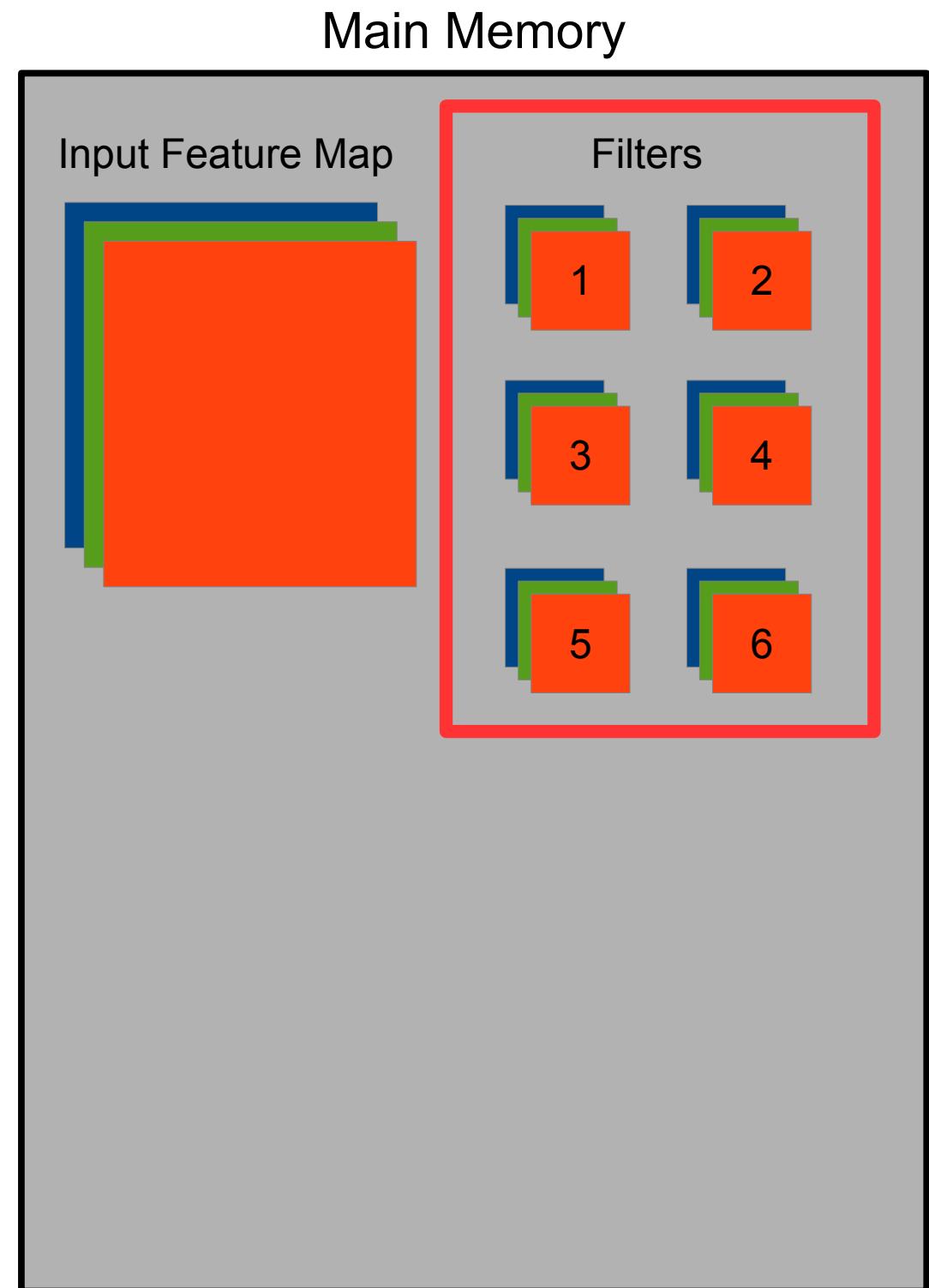
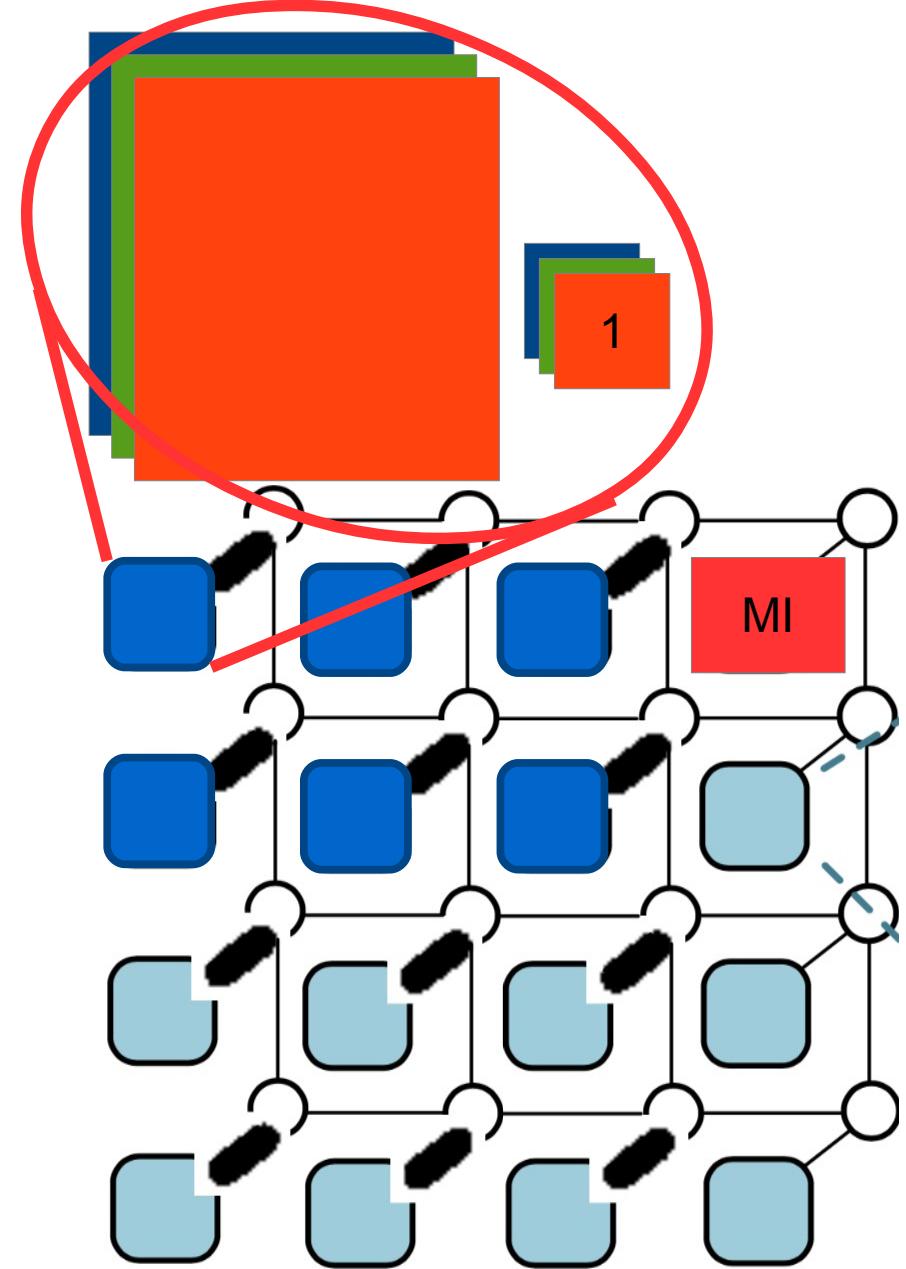
Convolution



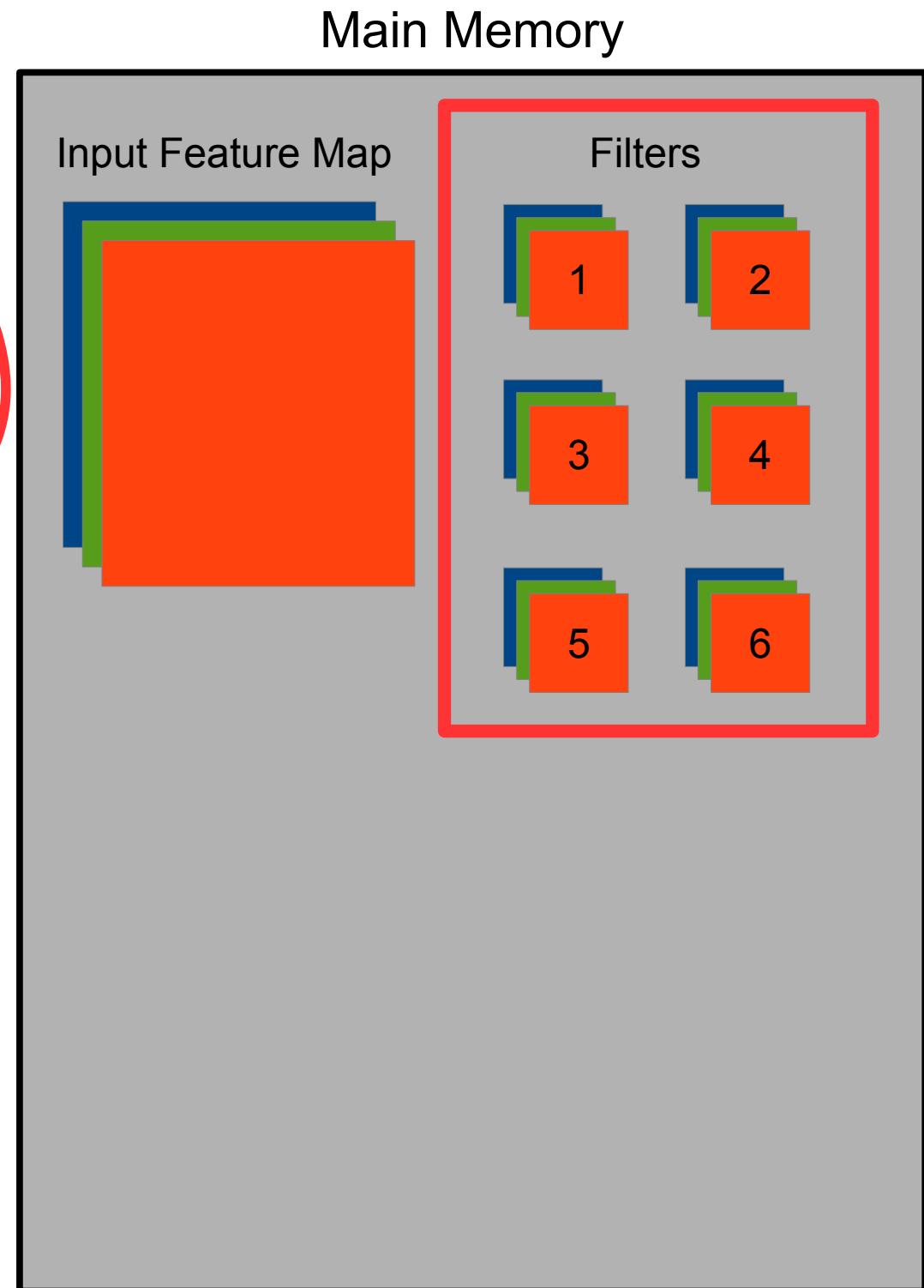
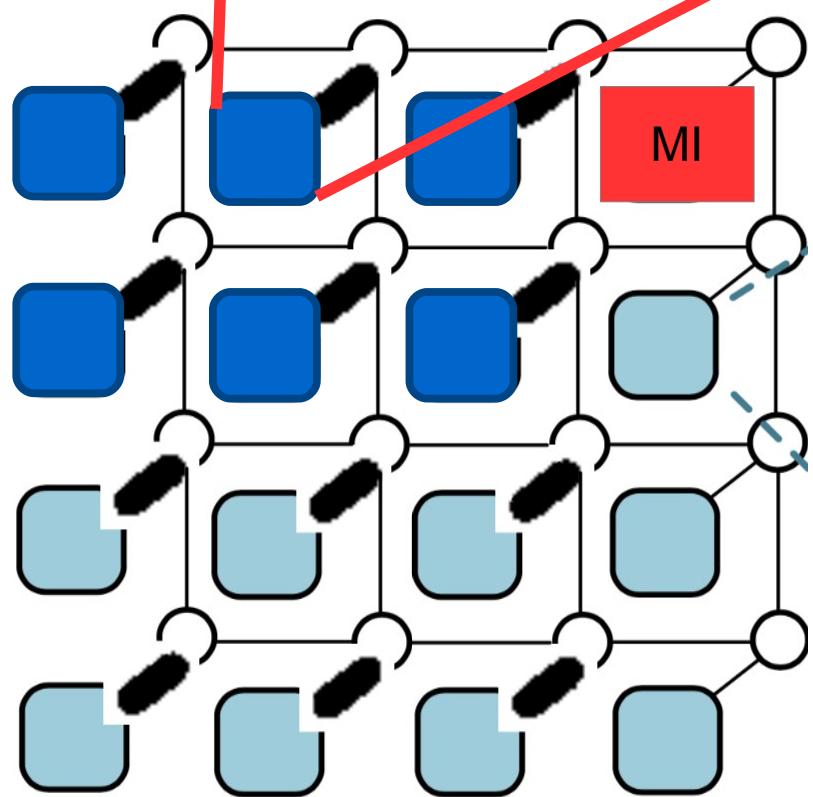
Convolution



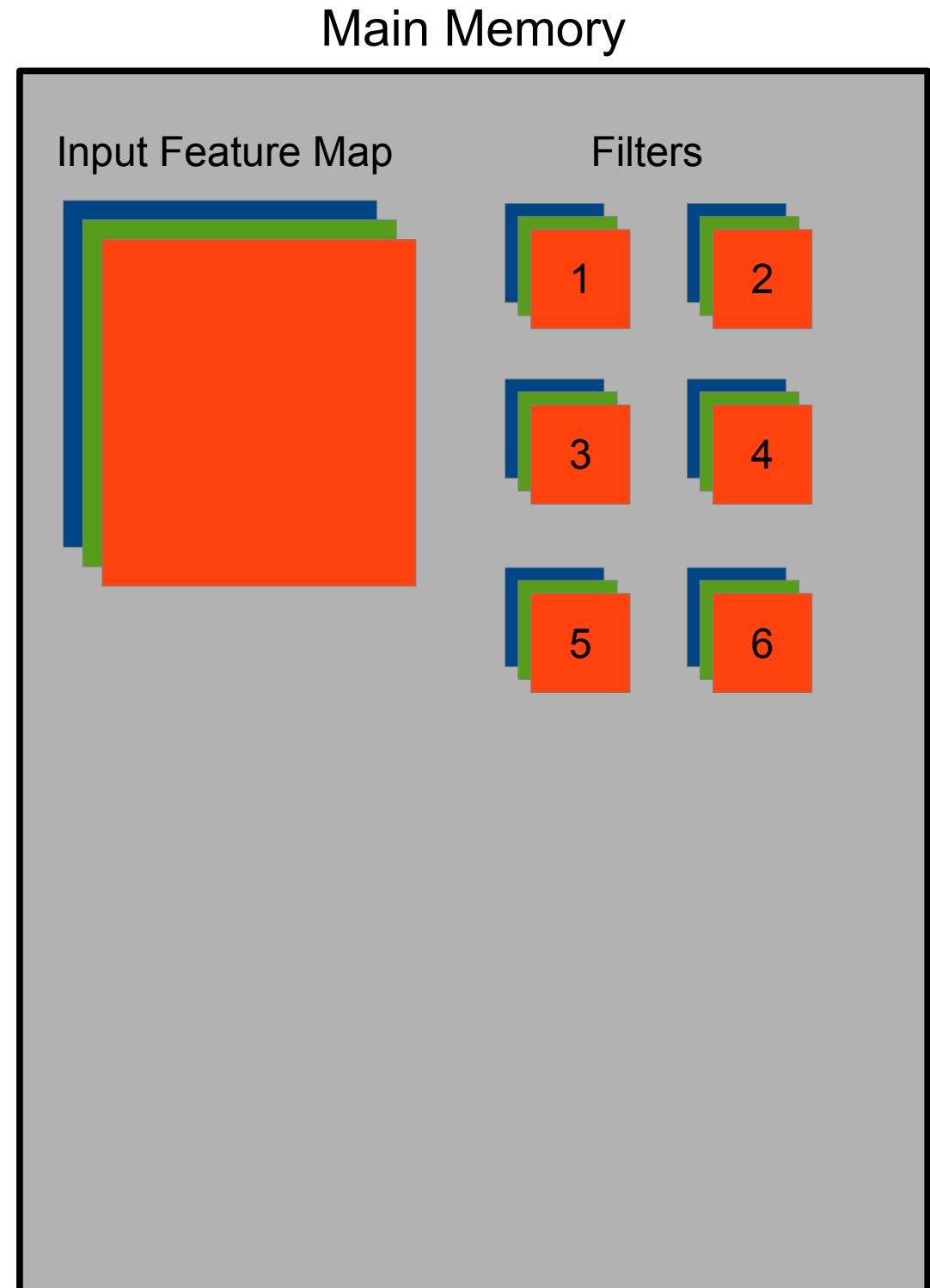
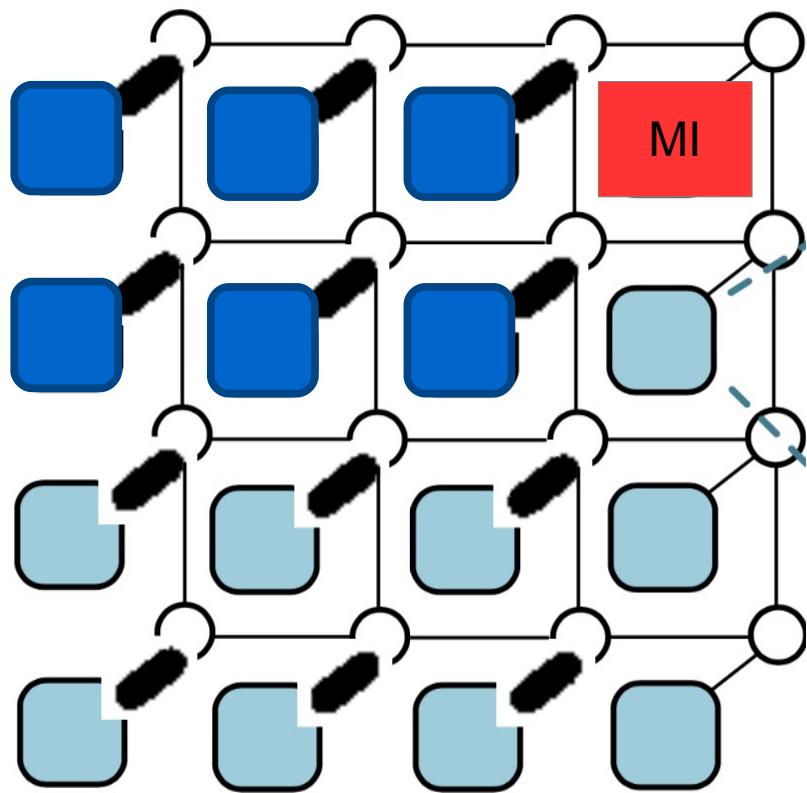
Convolution



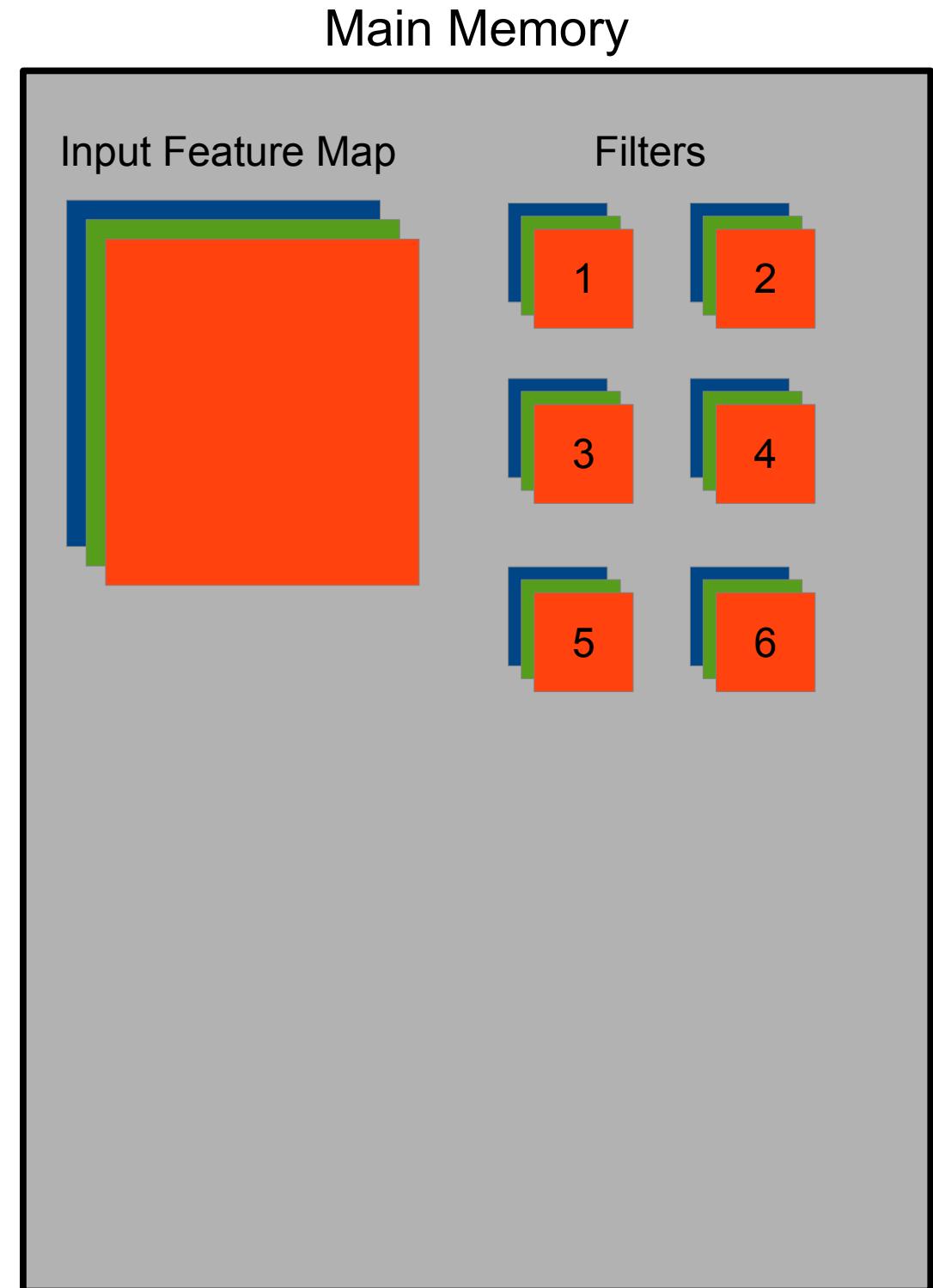
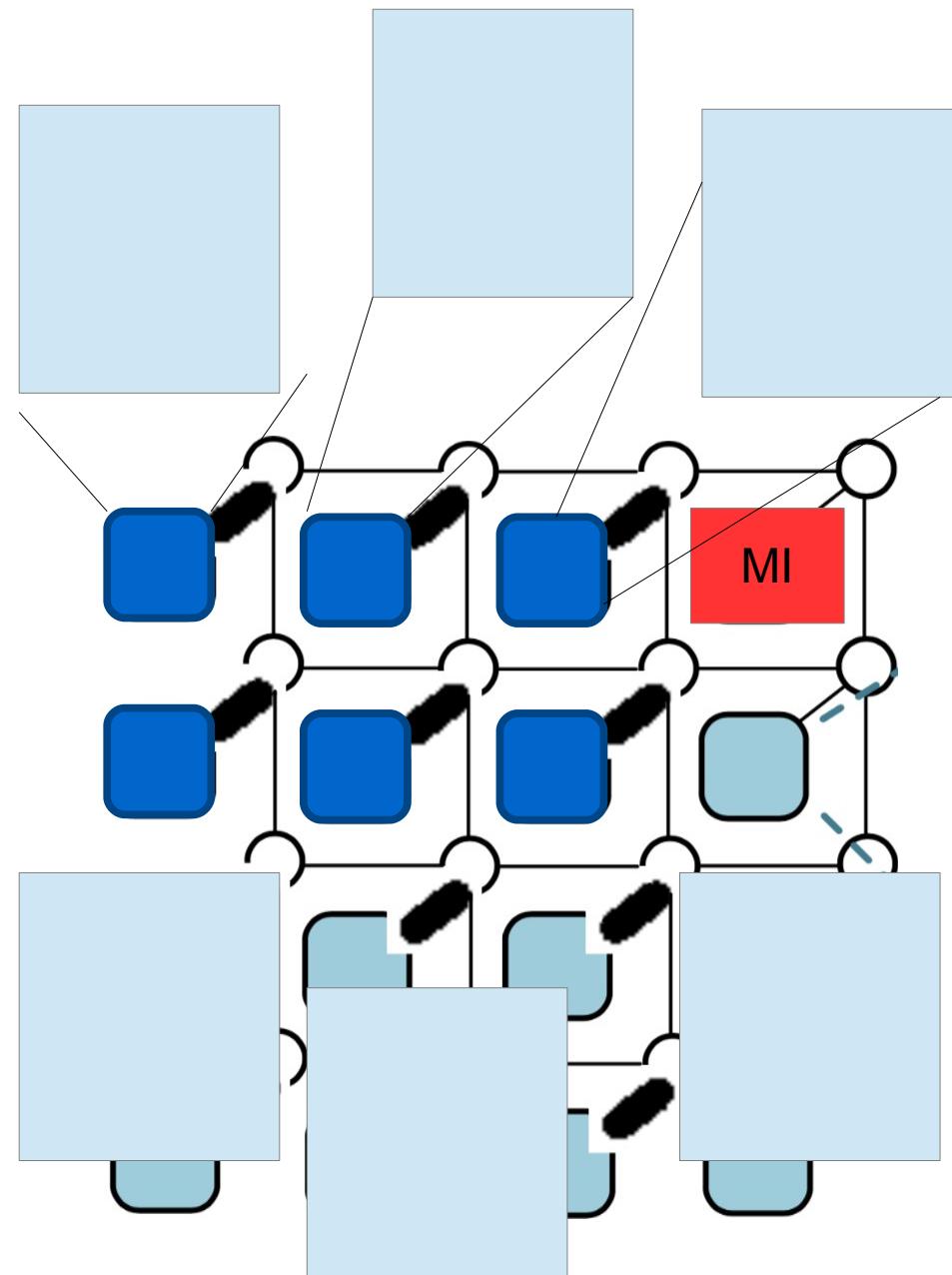
Convolution



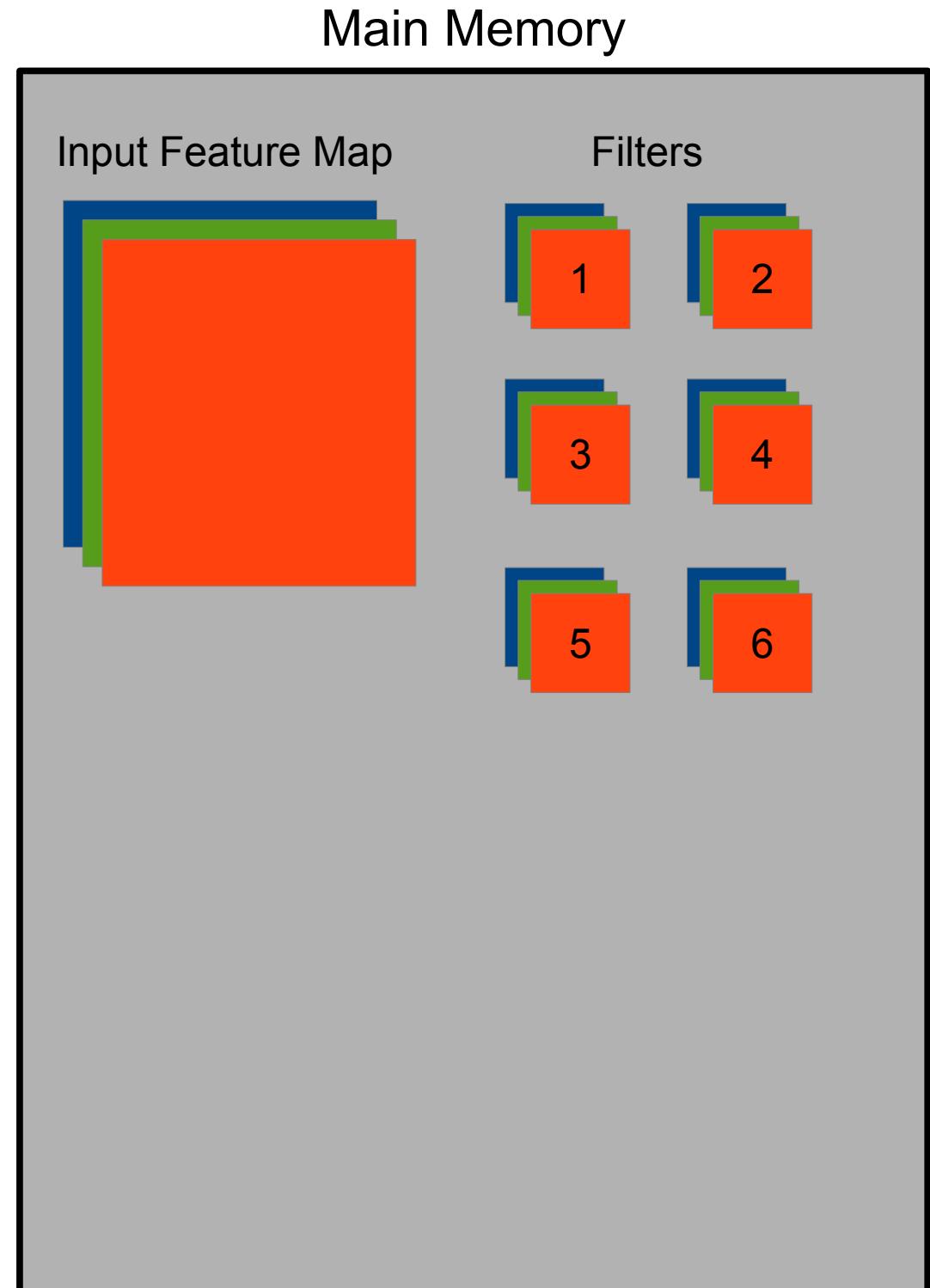
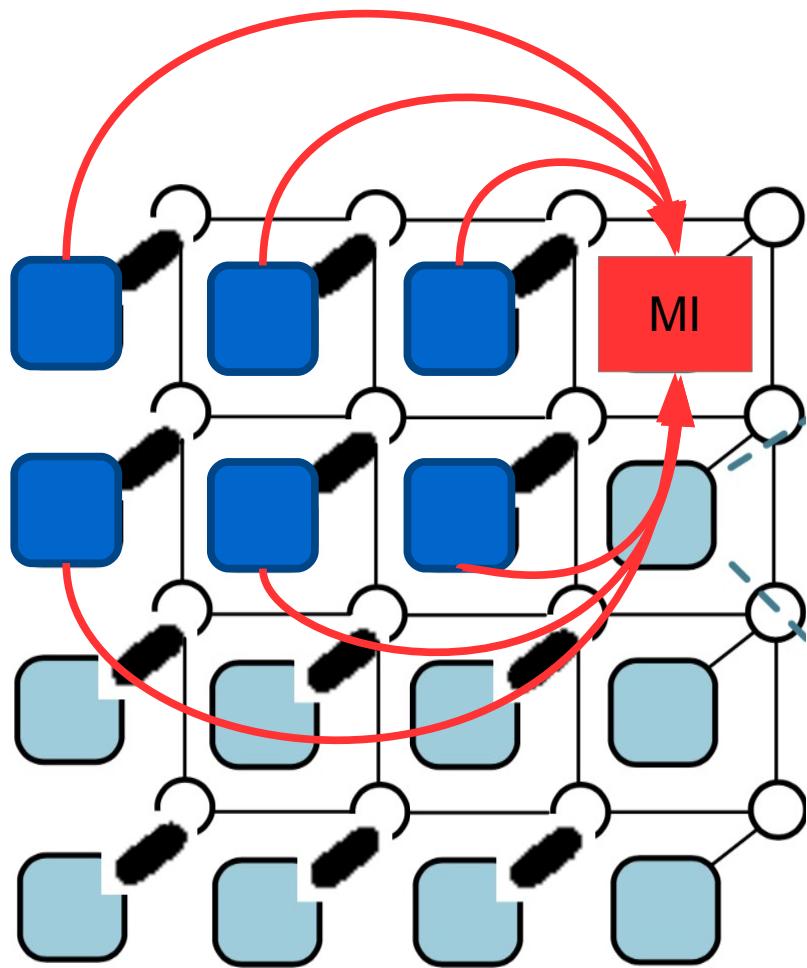
Convolution



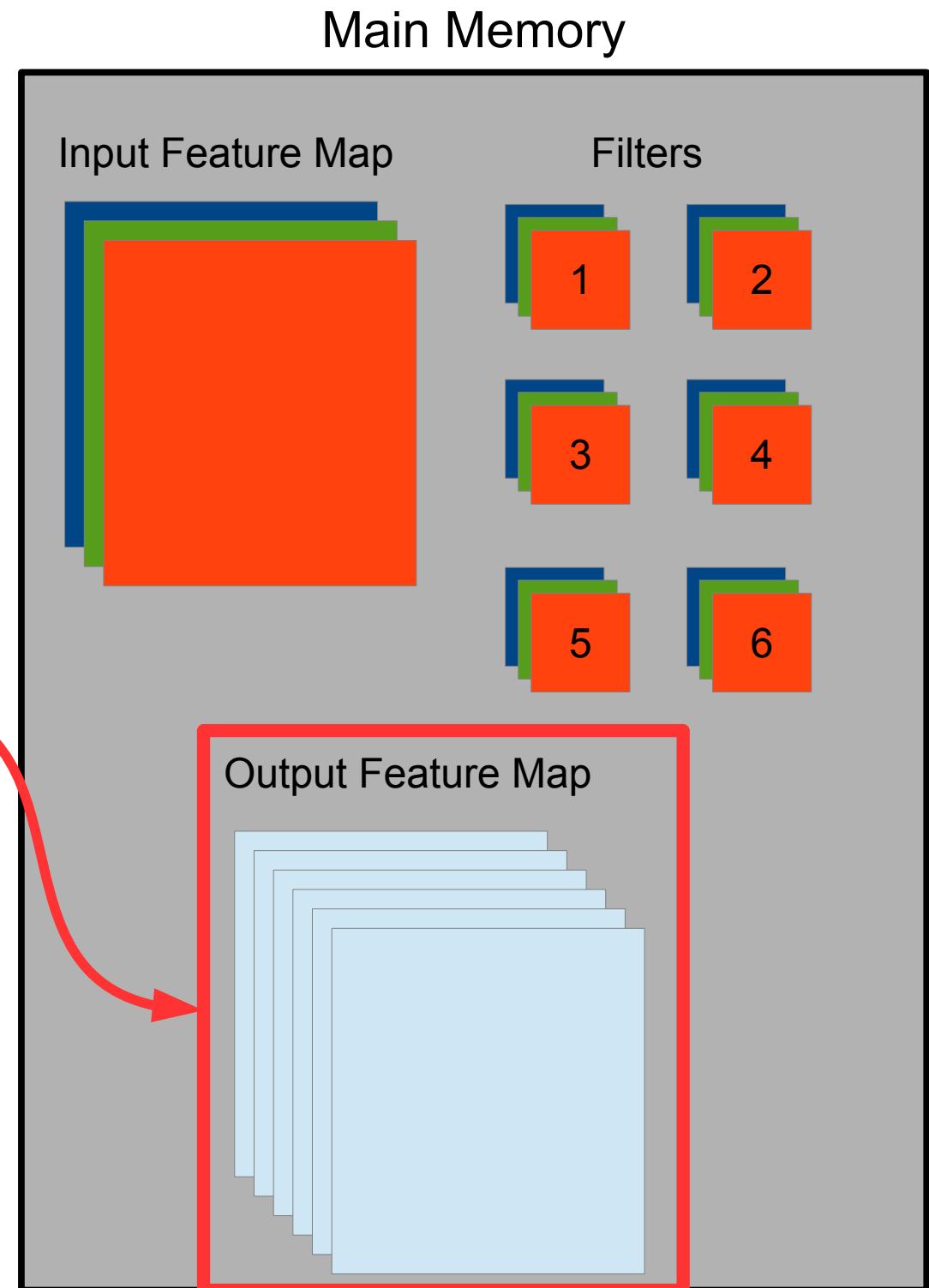
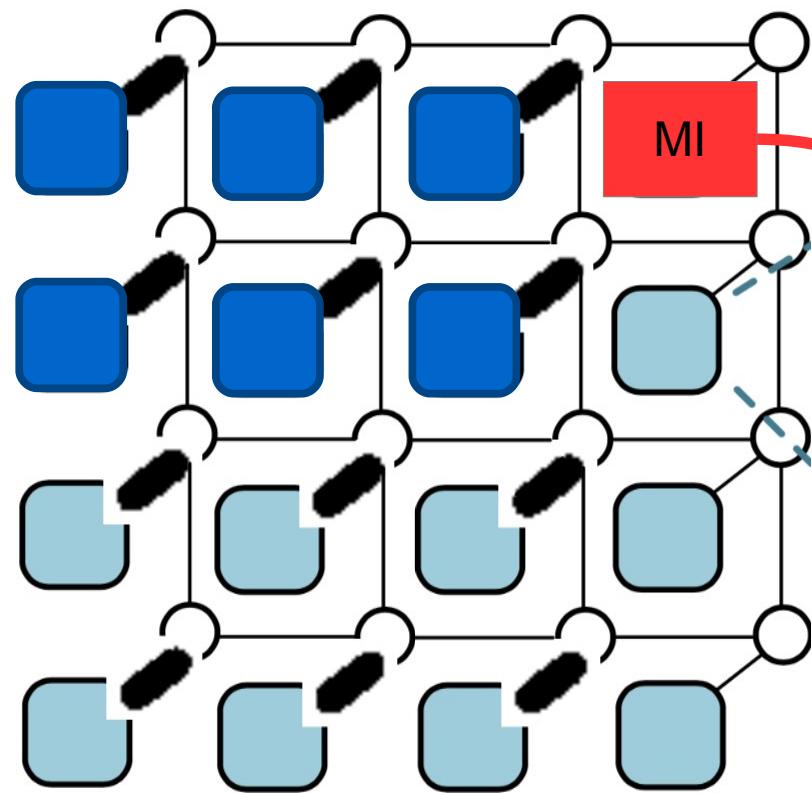
Convolution



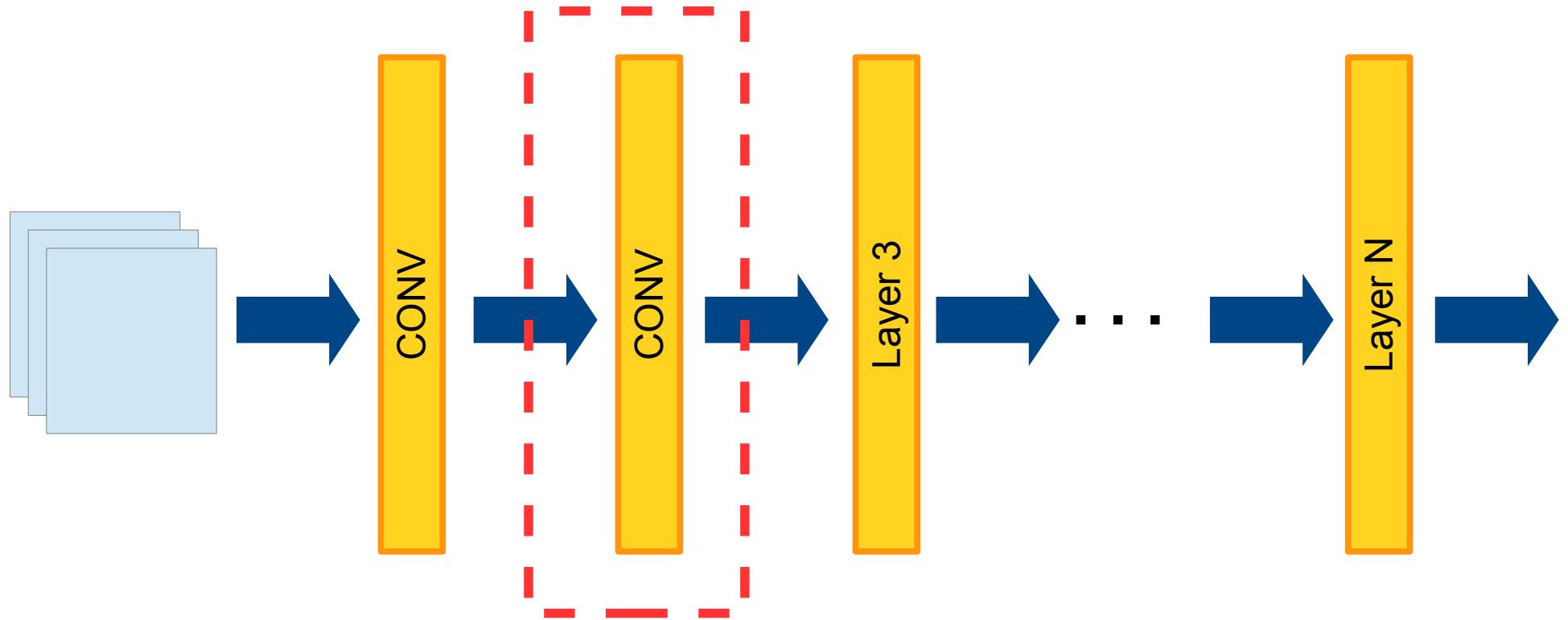
Convolution



Convolution

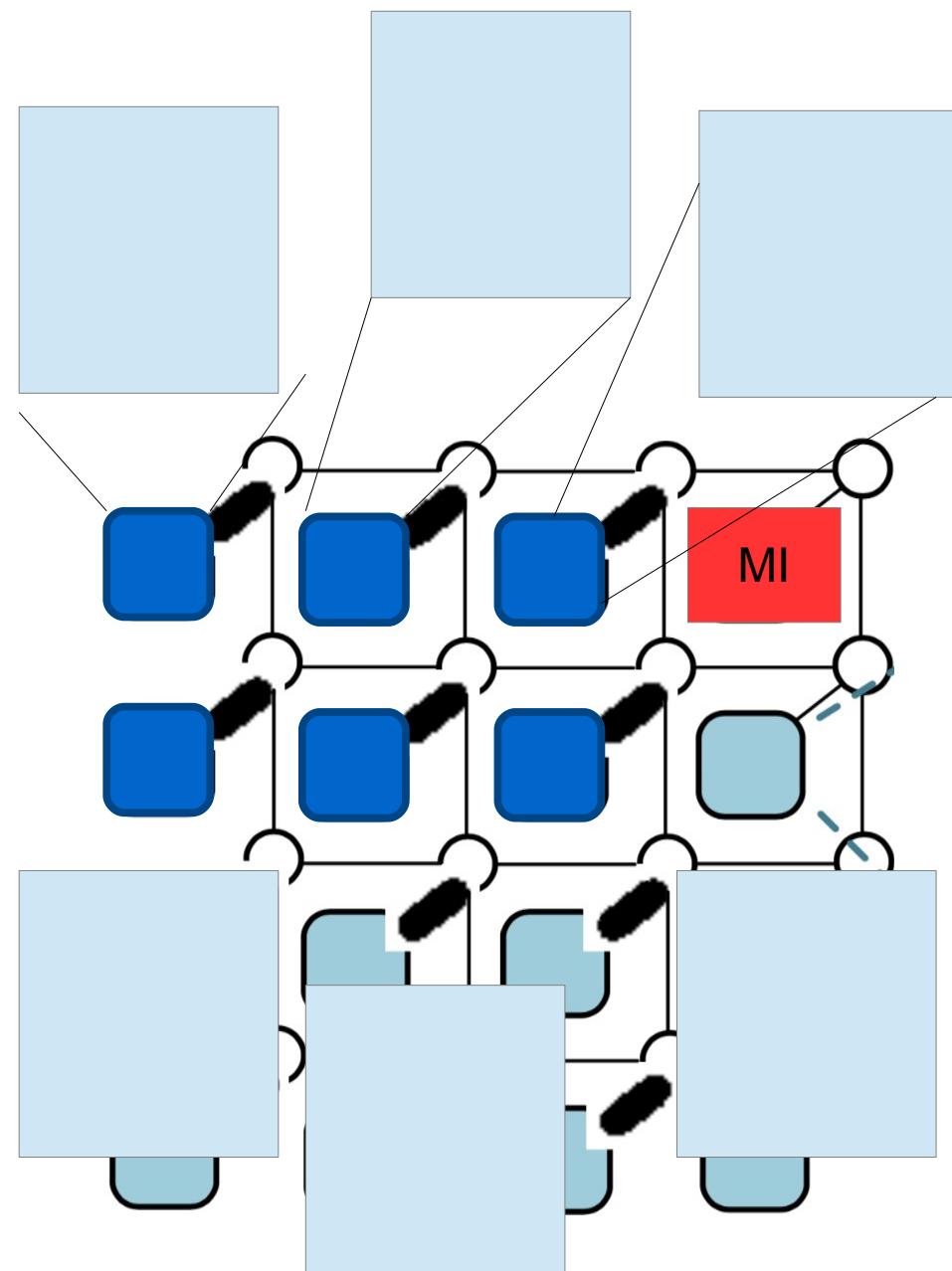


Example



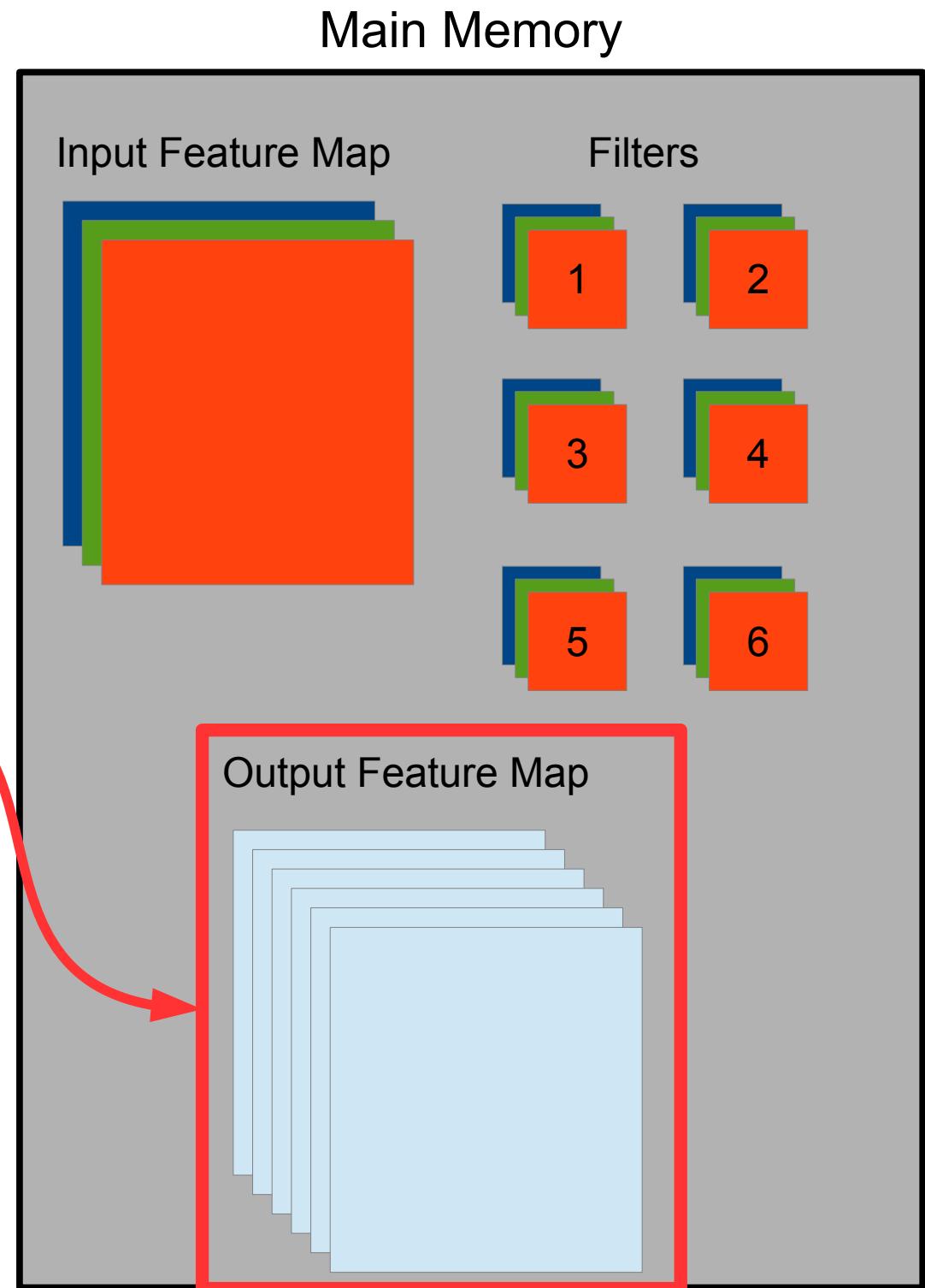
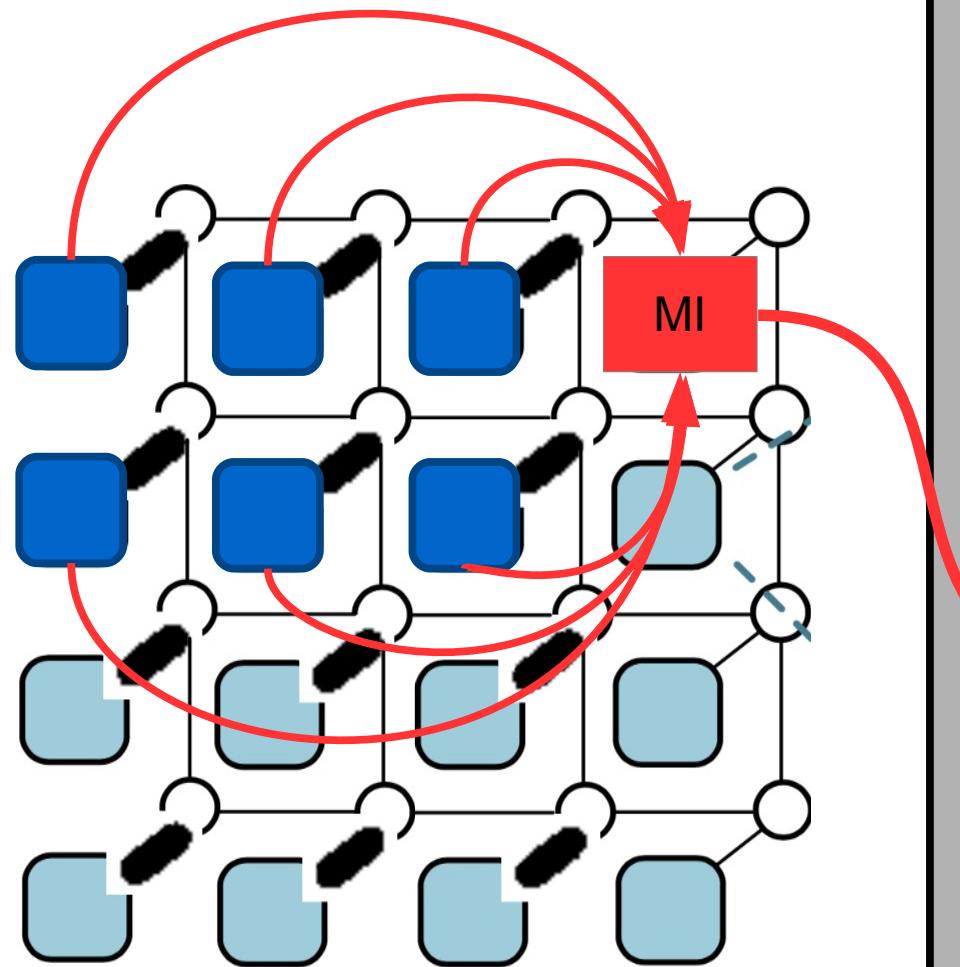
- The output feature map of layer 1 corresponds to the input feature map for layer2

Convolution

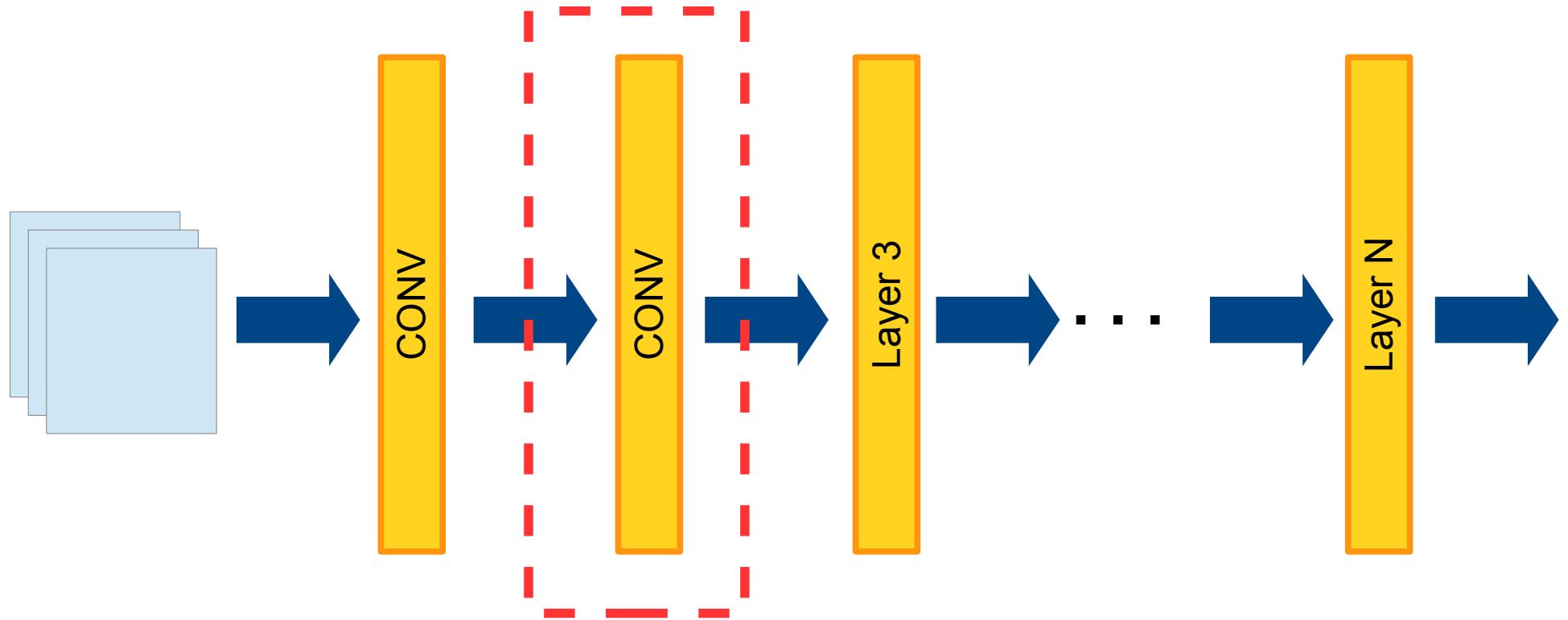


- The input feature map for the next layer is spread over the PE used in the previous layer

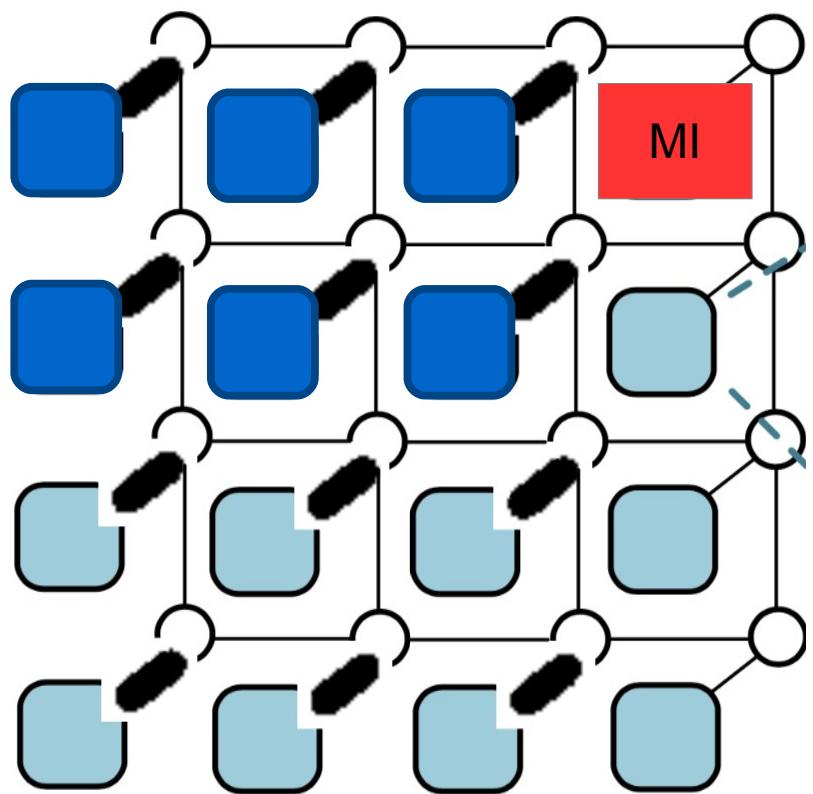
- No need to store back the output feature map into main memory



Convolution

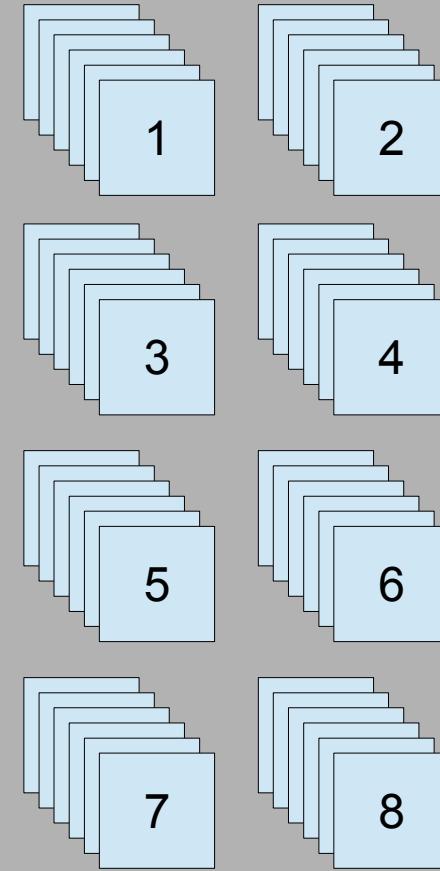


Convolution



Main Memory

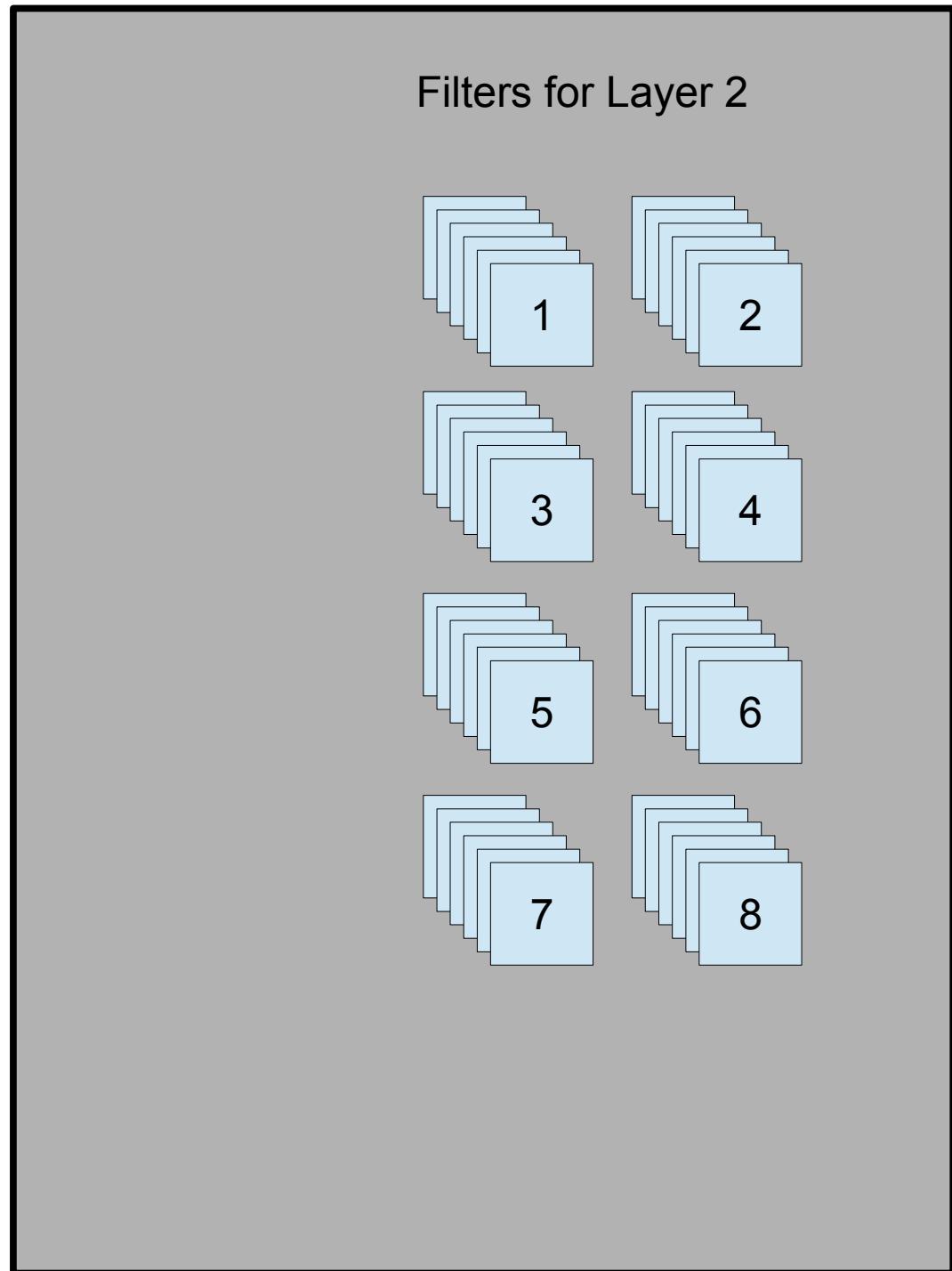
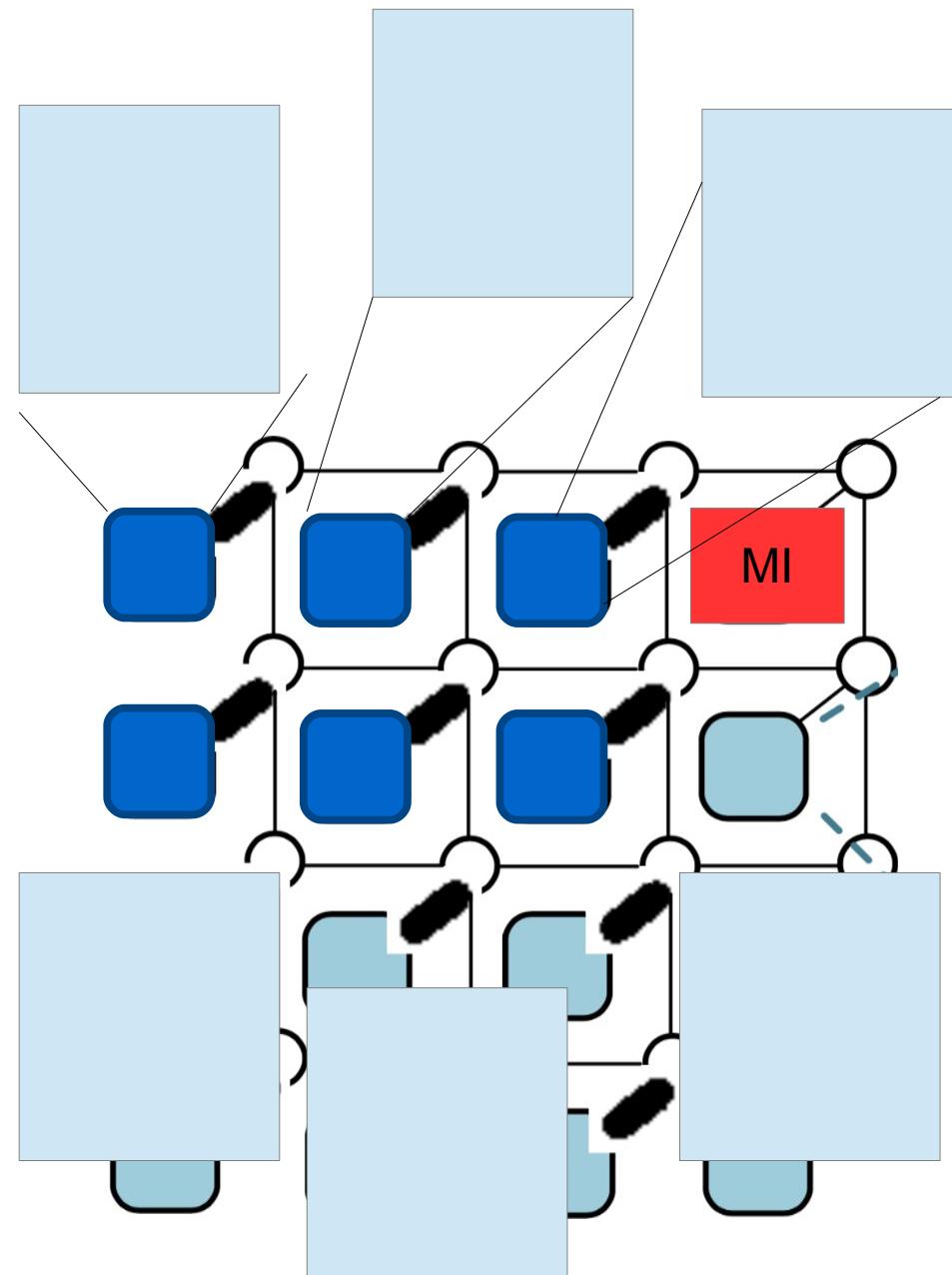
Filters for Layer 2



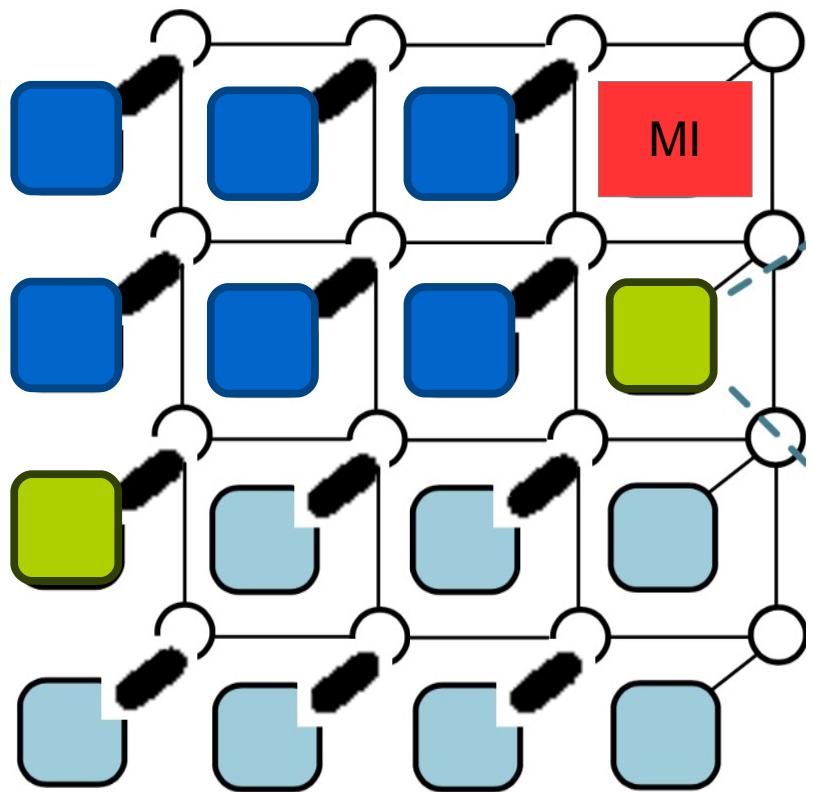
Convolution

Main Memory

Filters for Layer 2

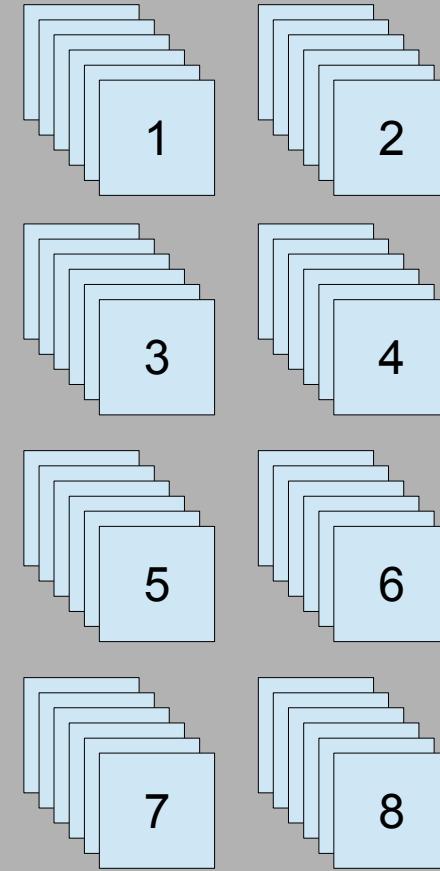


Convolution

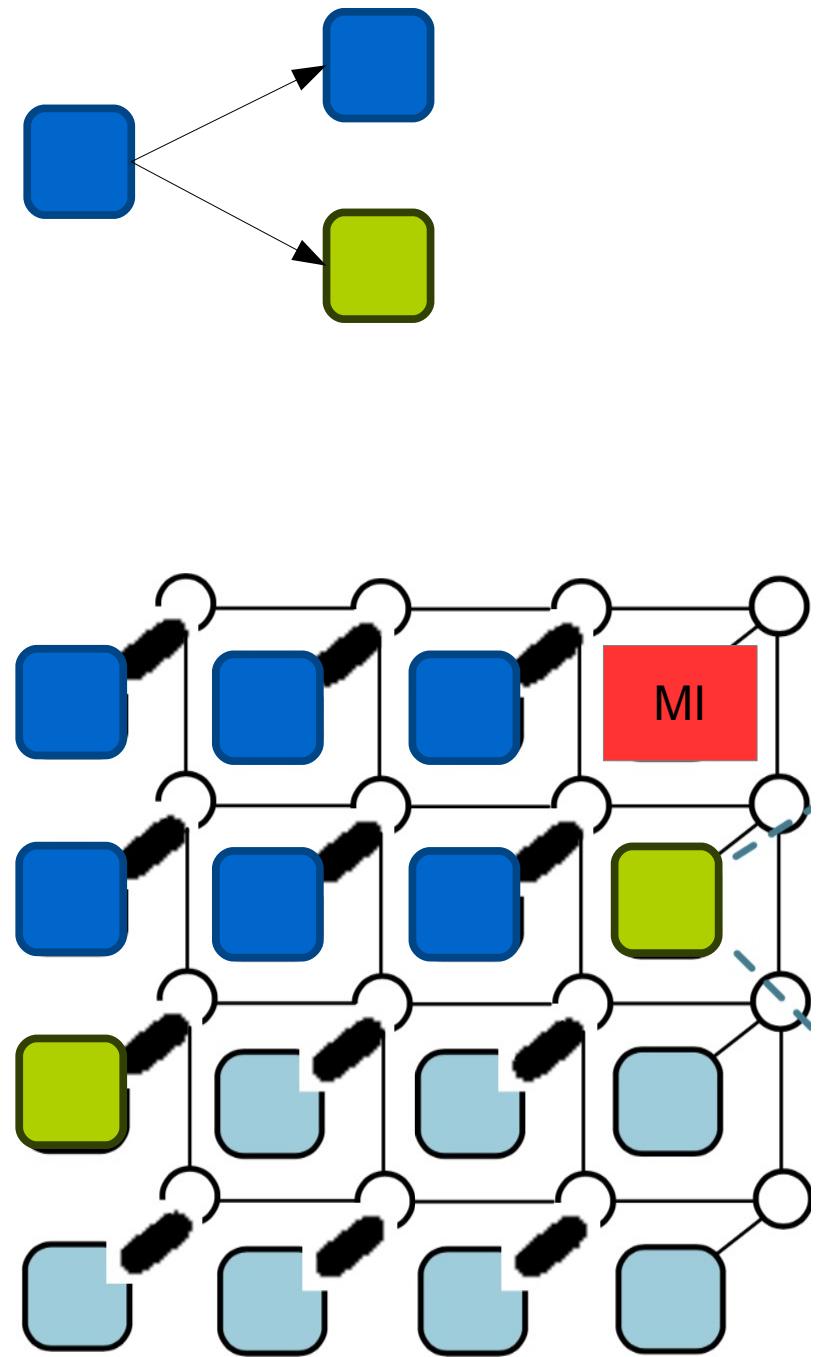


Main Memory

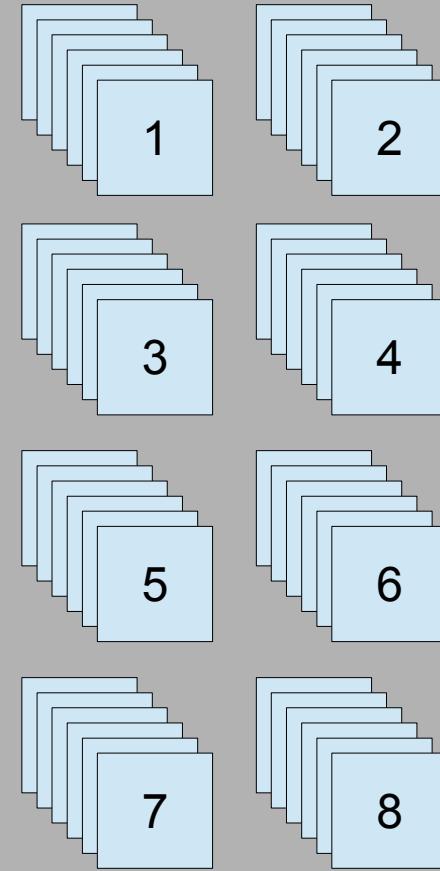
Filters for Layer 2



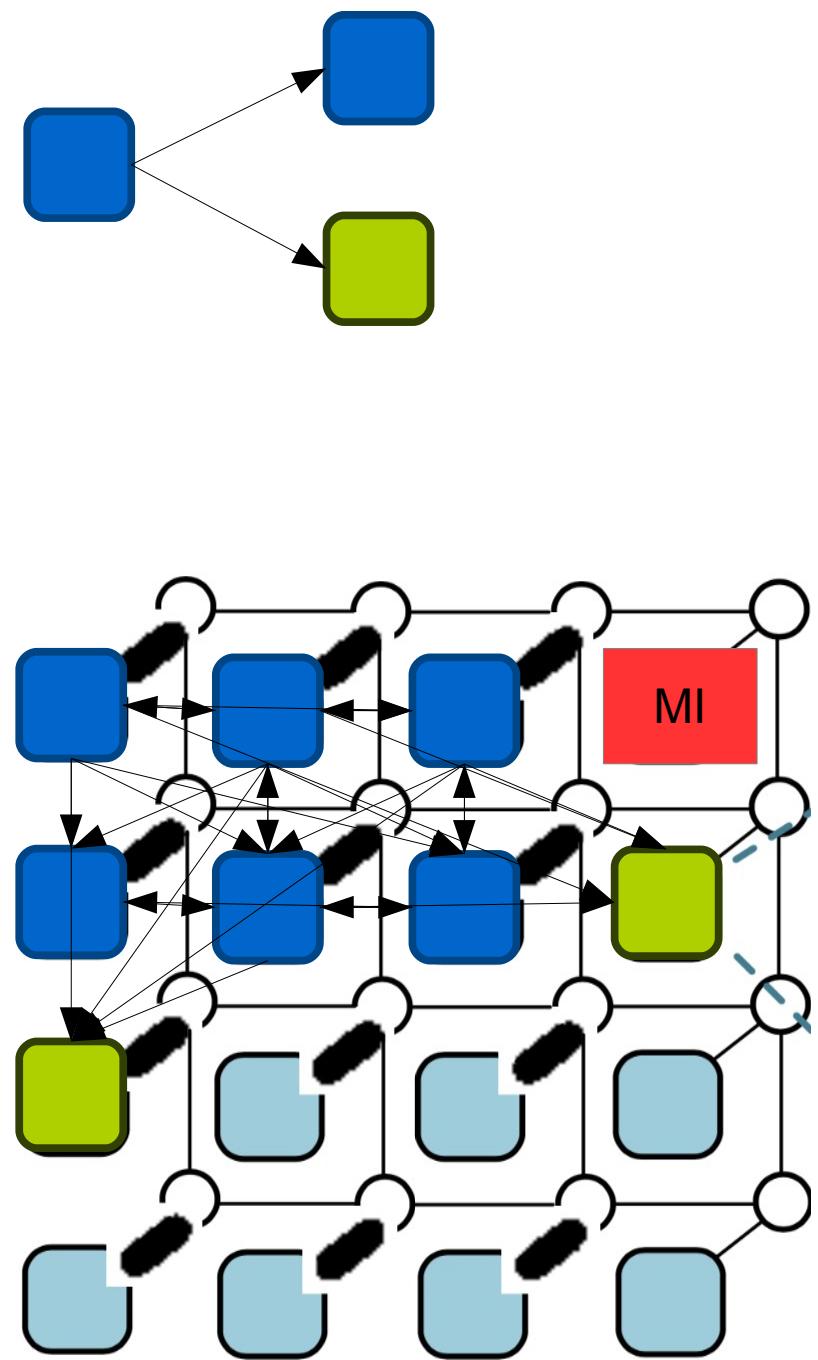
Main Memory



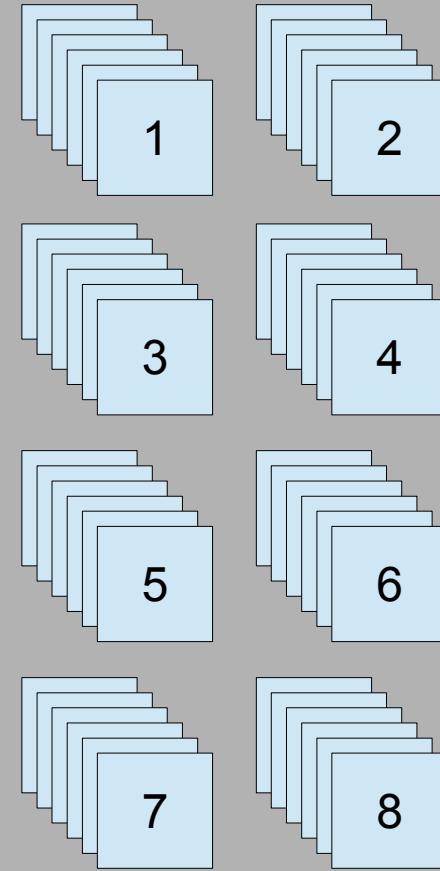
Filters for Layer 2



Main Memory

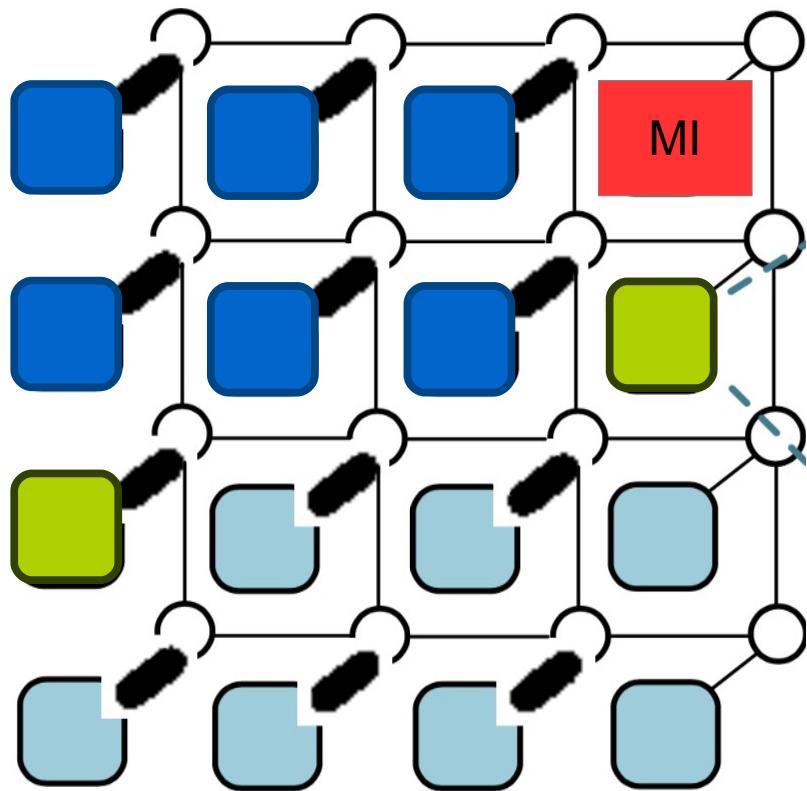


Filters for Layer 2



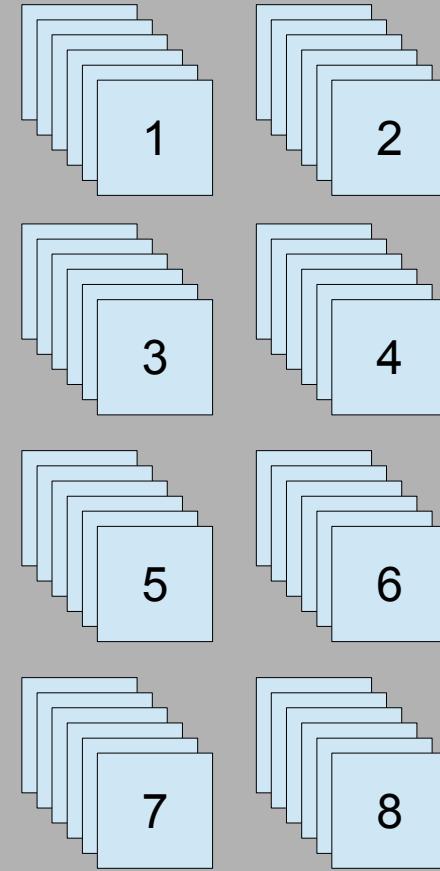
Convolution

- All PE have the input feature map

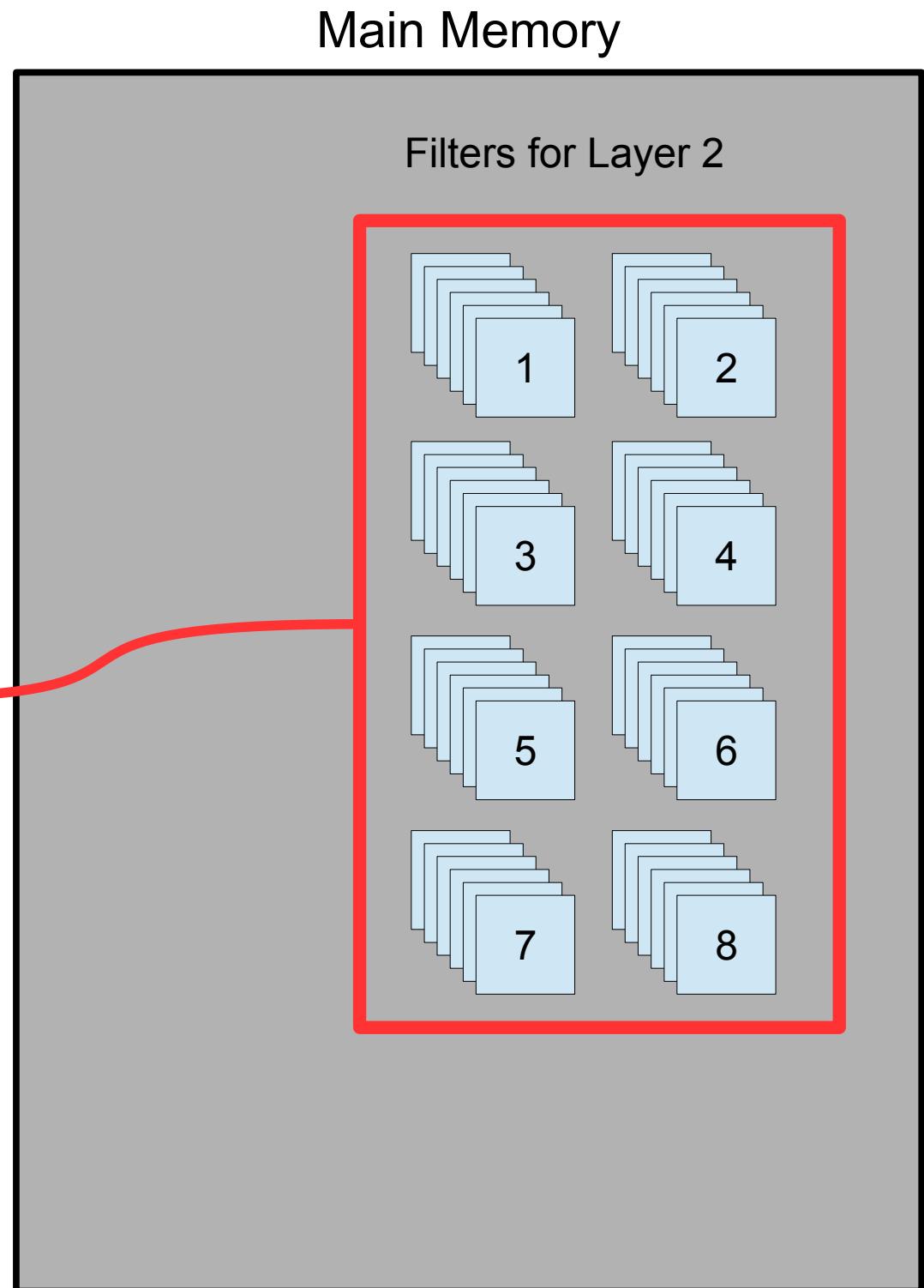
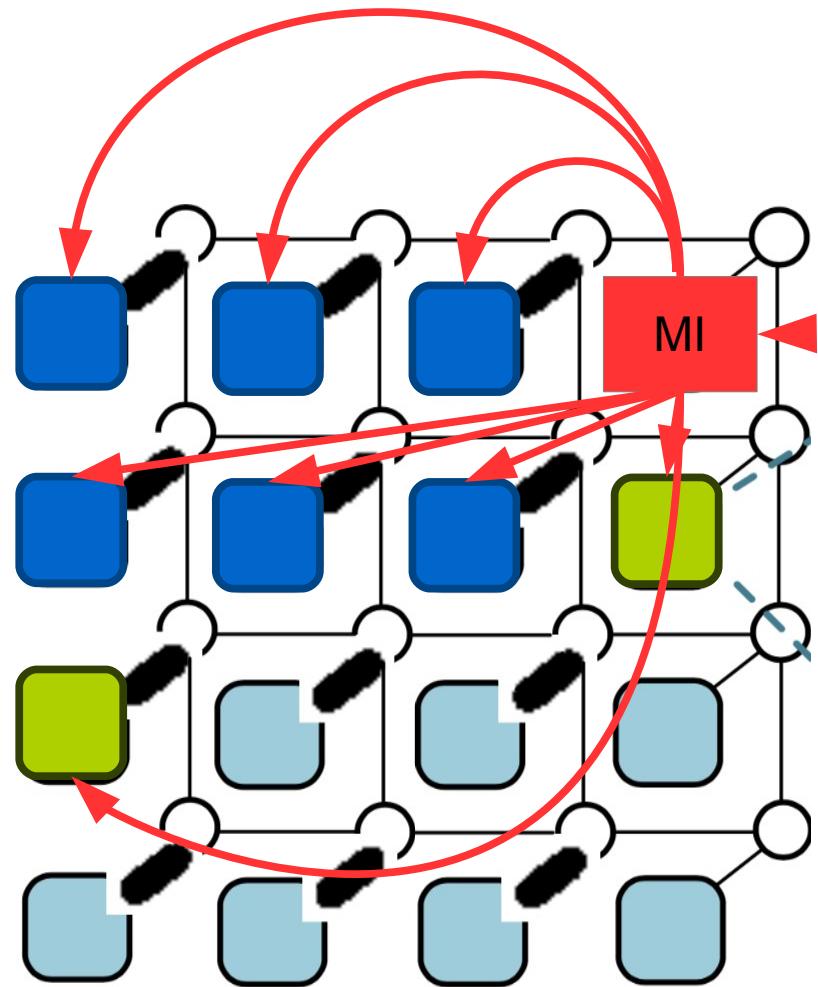


Main Memory

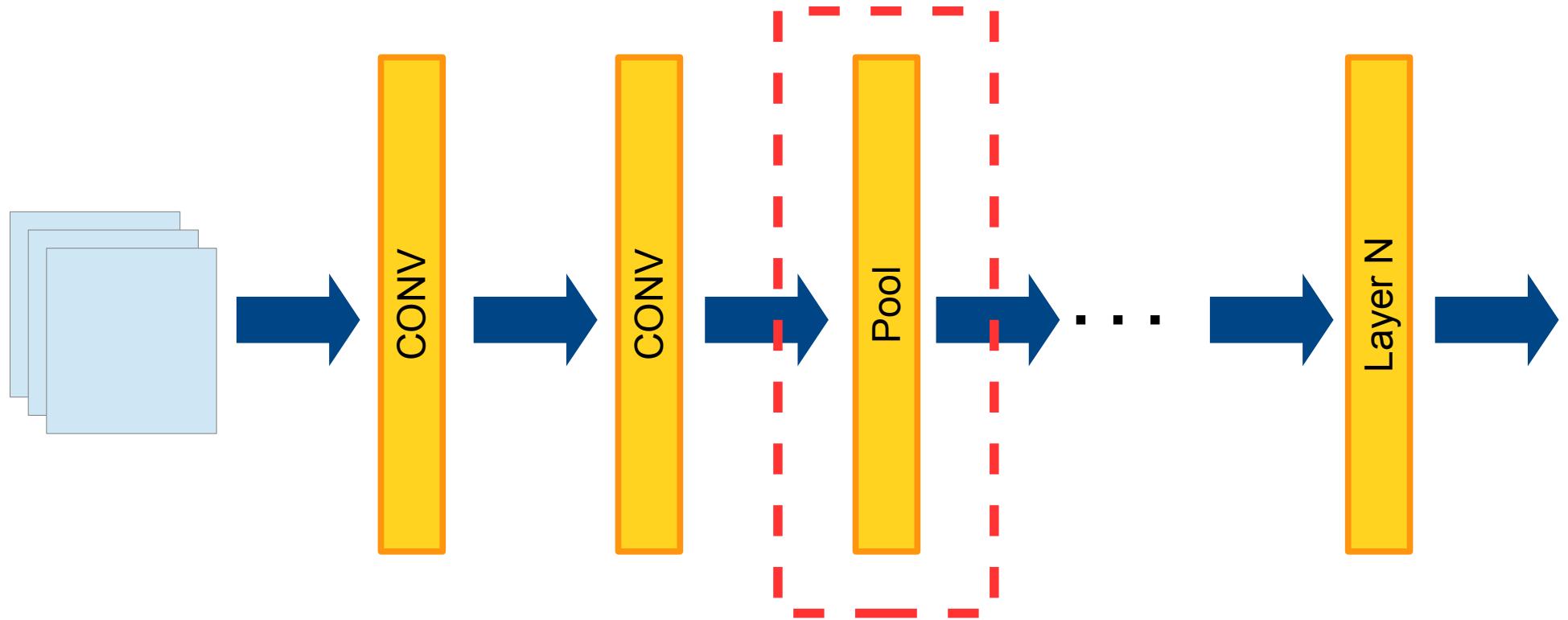
Filters for Layer 2



Convolution

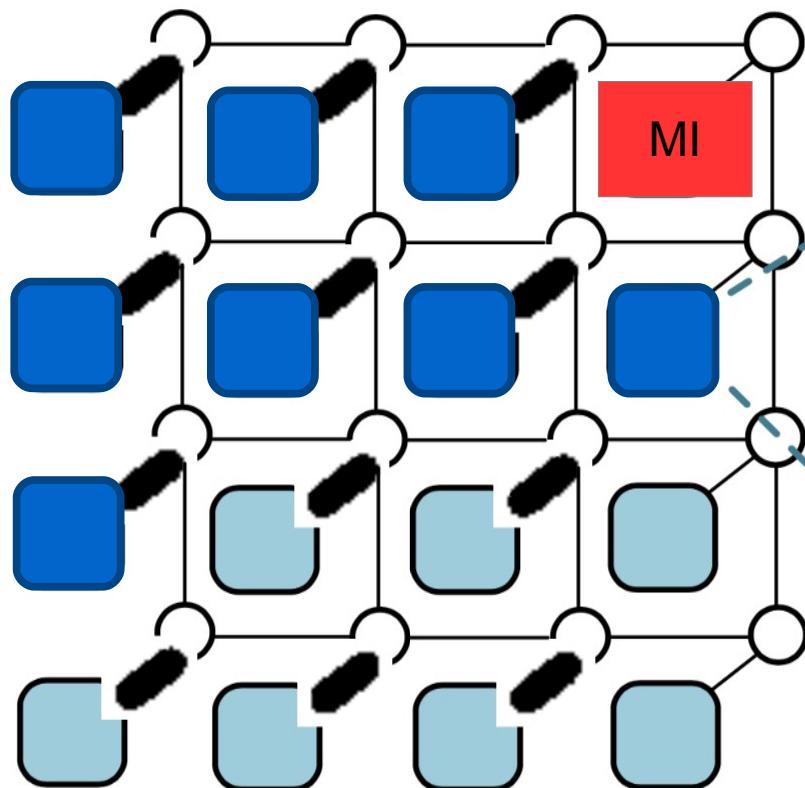


Max/Avg Pool



Avg/Max Pool

- Input feature map spread over the PE used in the previous layer
- Each PE computes its feature map channel



No memory traffic

Fully Connection

Input Feature Map



$(W \times H \times C) \times 1$

Output Feature Map

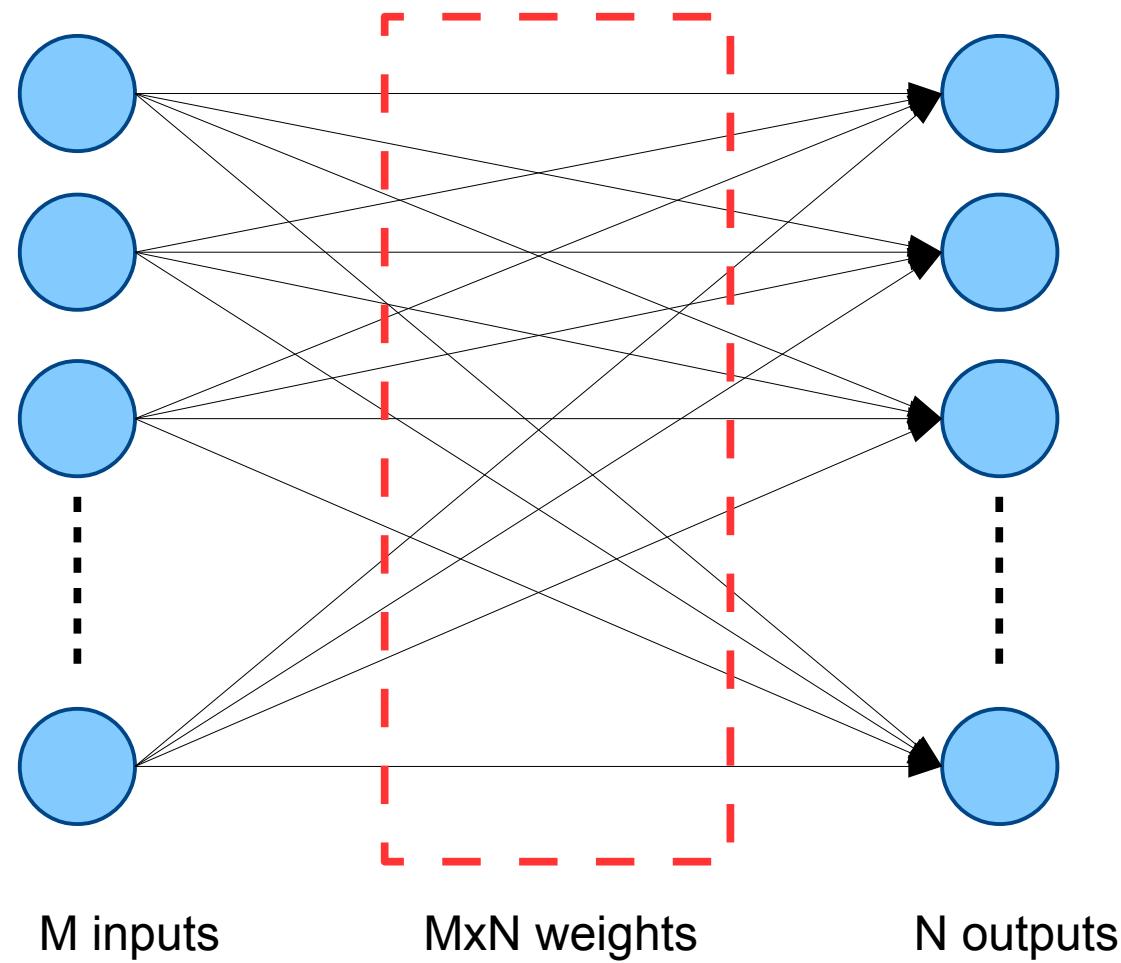


$N \times 1$

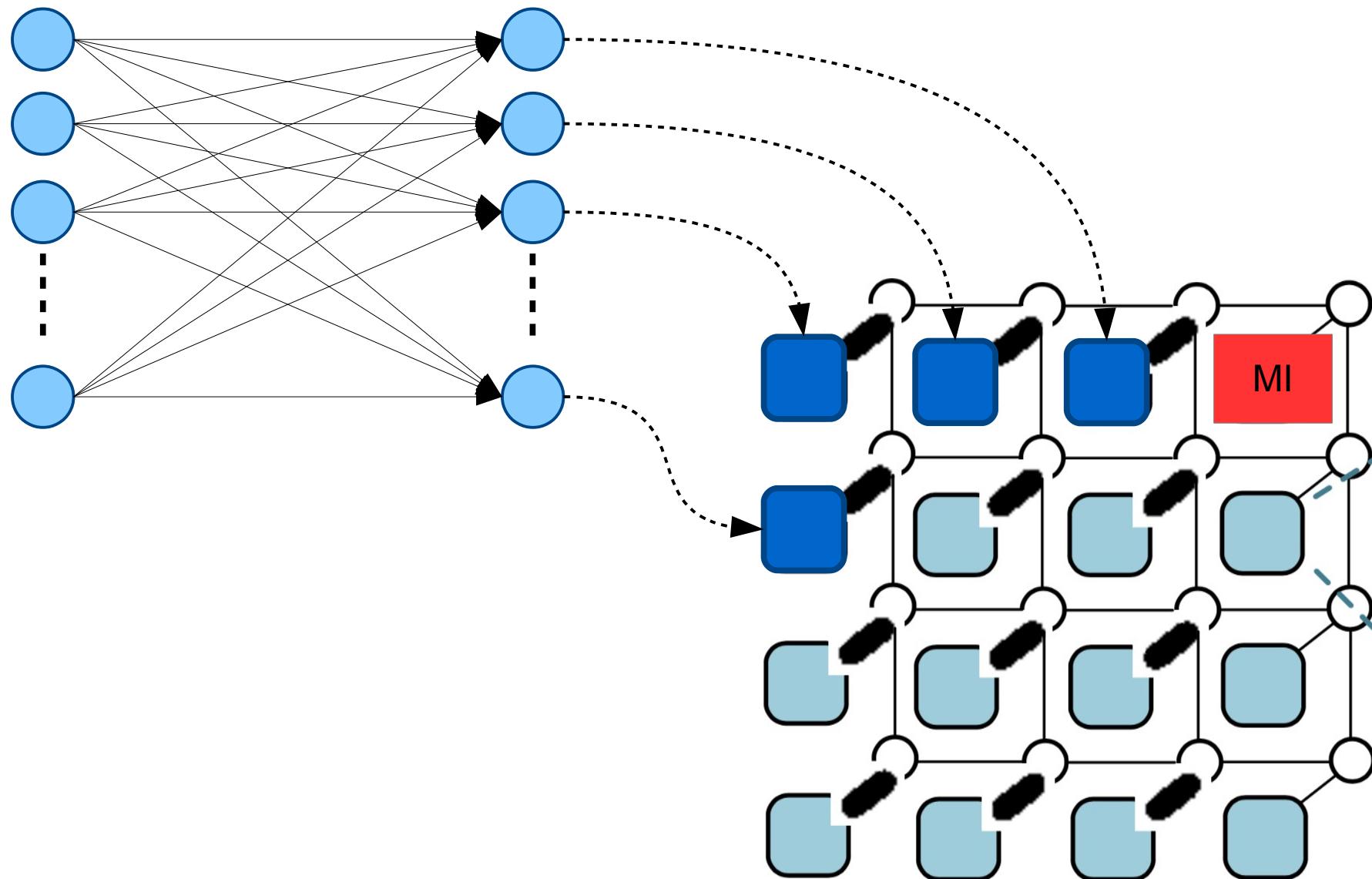
Fully Connected Layer

N

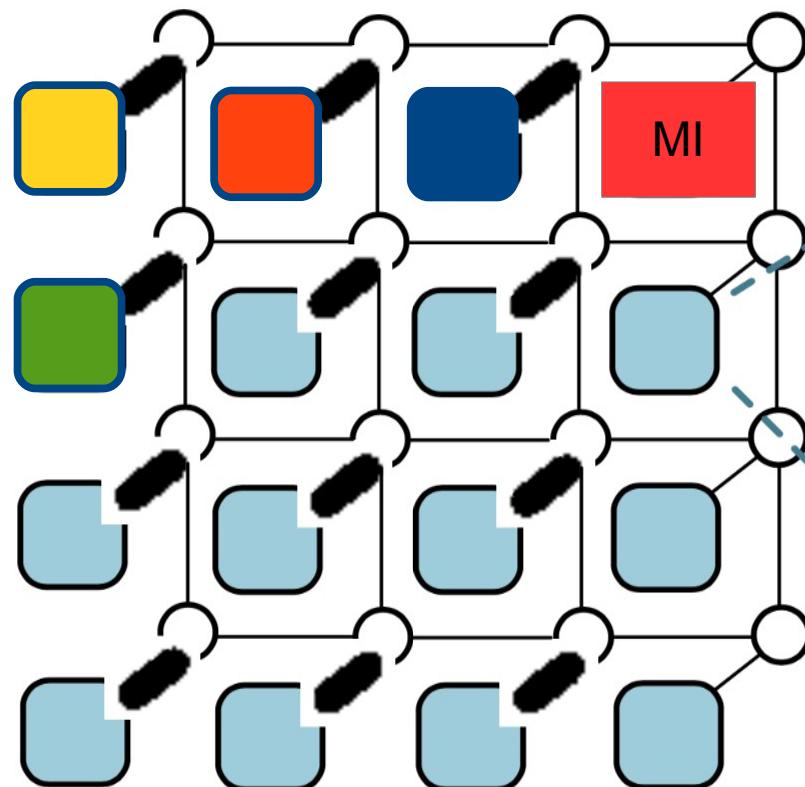
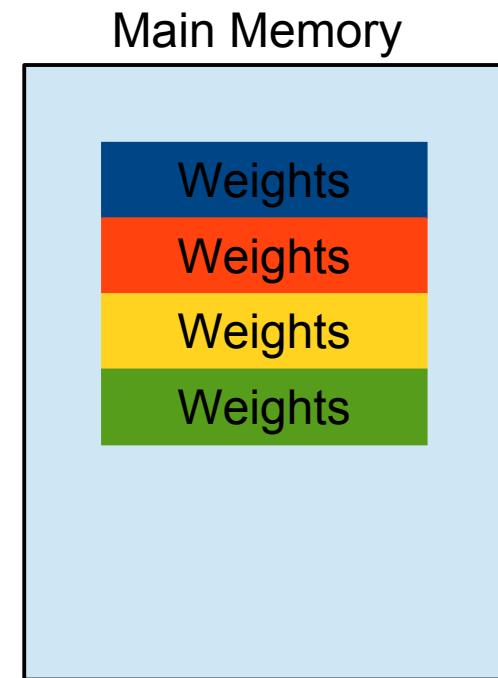
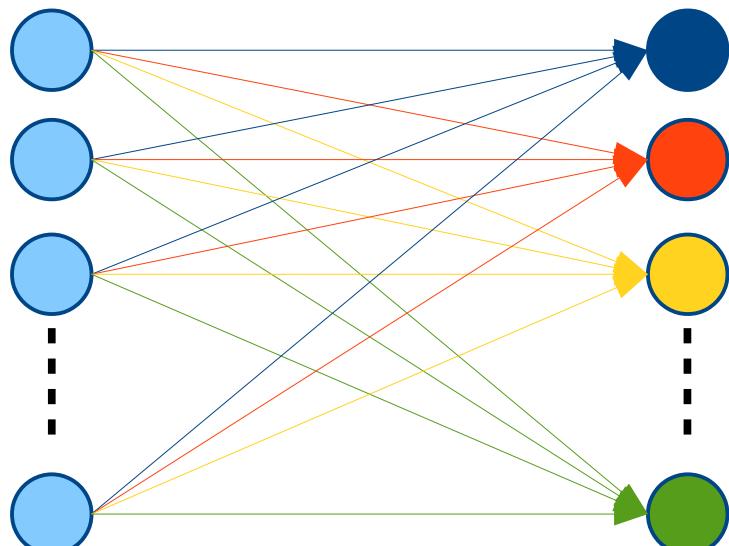
Fully Connection



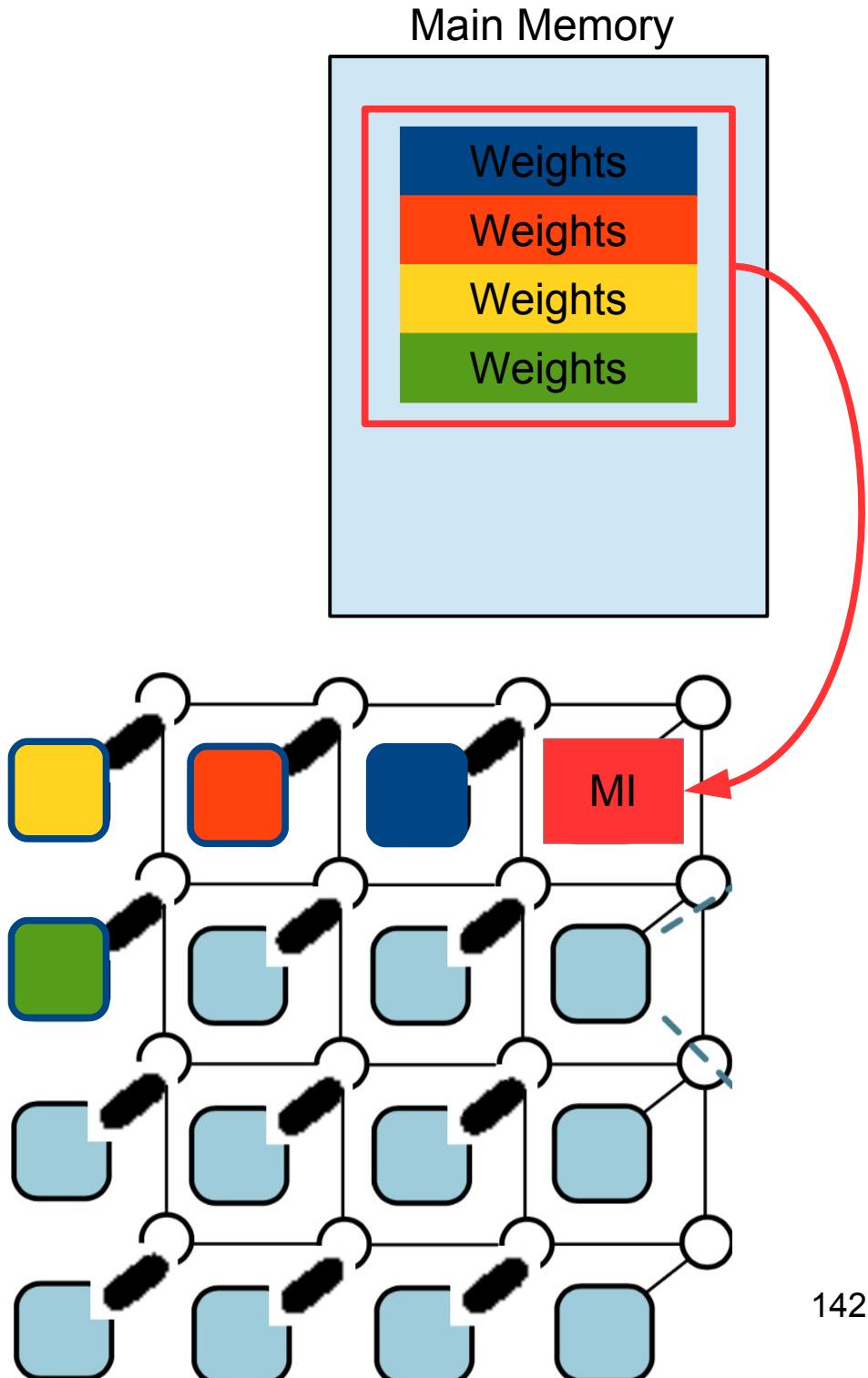
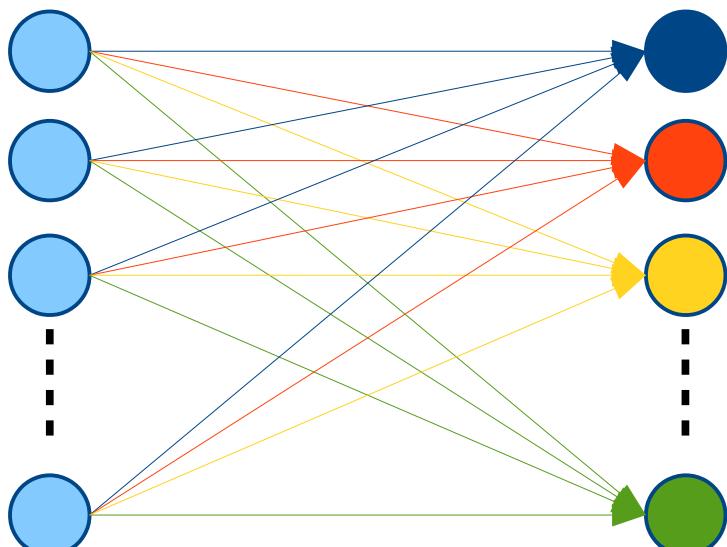
Fully Connection



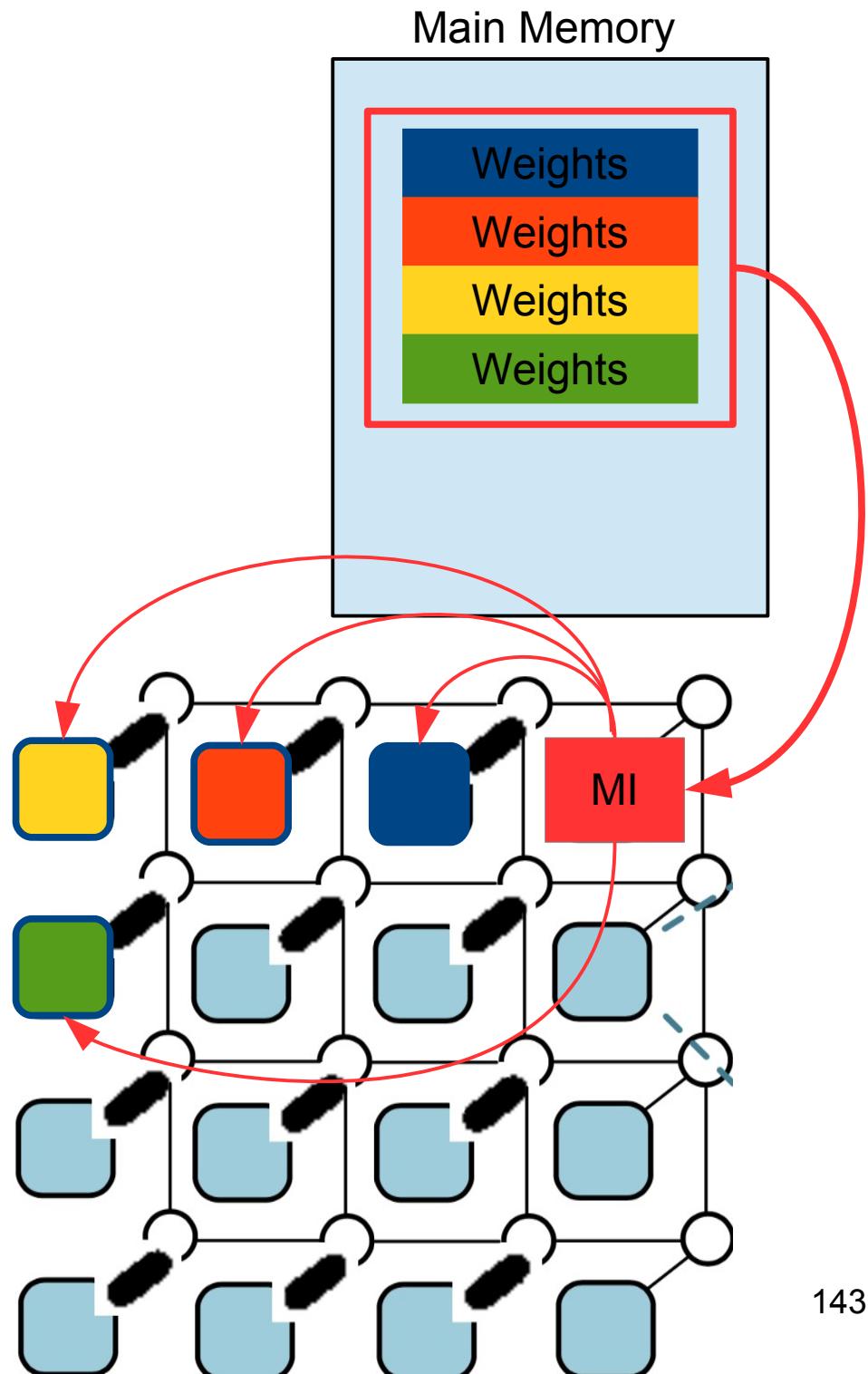
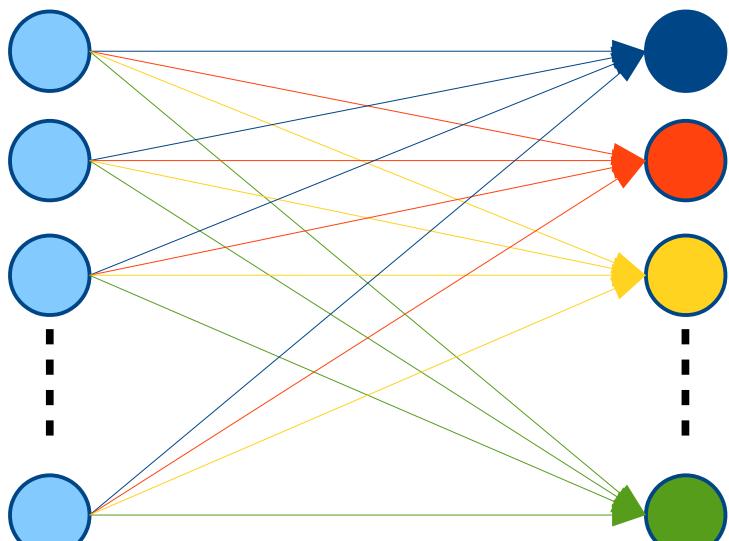
Fully Connection



Fully Connection



Fully Connection



Memory/Comm Traffic Summary

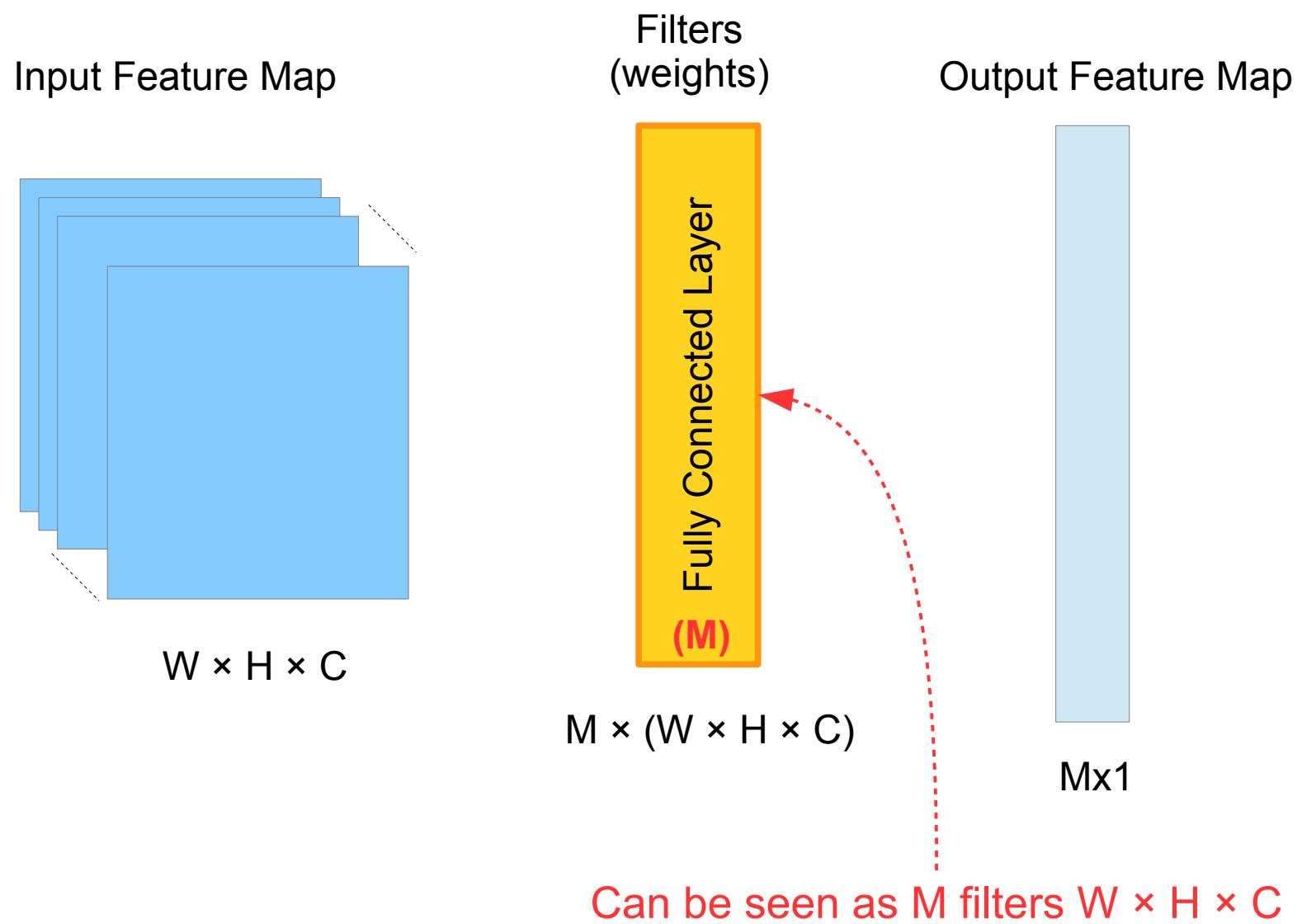
- High traffic volume to/from memory
- MI is an hot-spot point
- Different traffic types
 - Unicast
 - Multicast
 - Broadcast
- Communication network plays an important role

[S. M. Nabavinejad, *et al.* “An Overview of Efficient Interconnection Networks for Deep Neural Network Accelerators”. IEEE J. Emerg. Sel. Topics Circuits Syst. 2020]

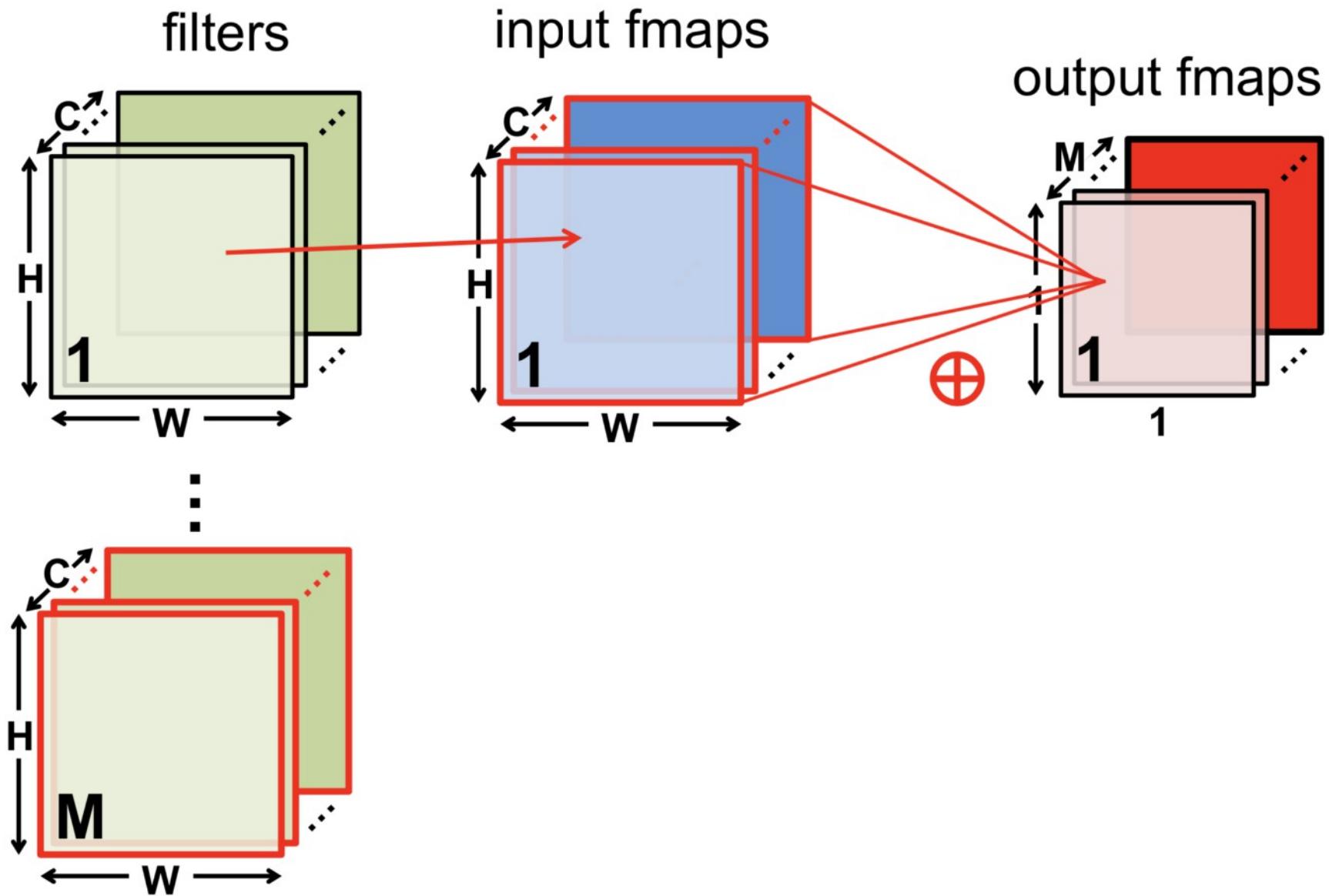
Agenda

- Deep Neural Networks
- Types of layers
- Memory and communication traffic
- Kernel computation
- Quantitative analysis

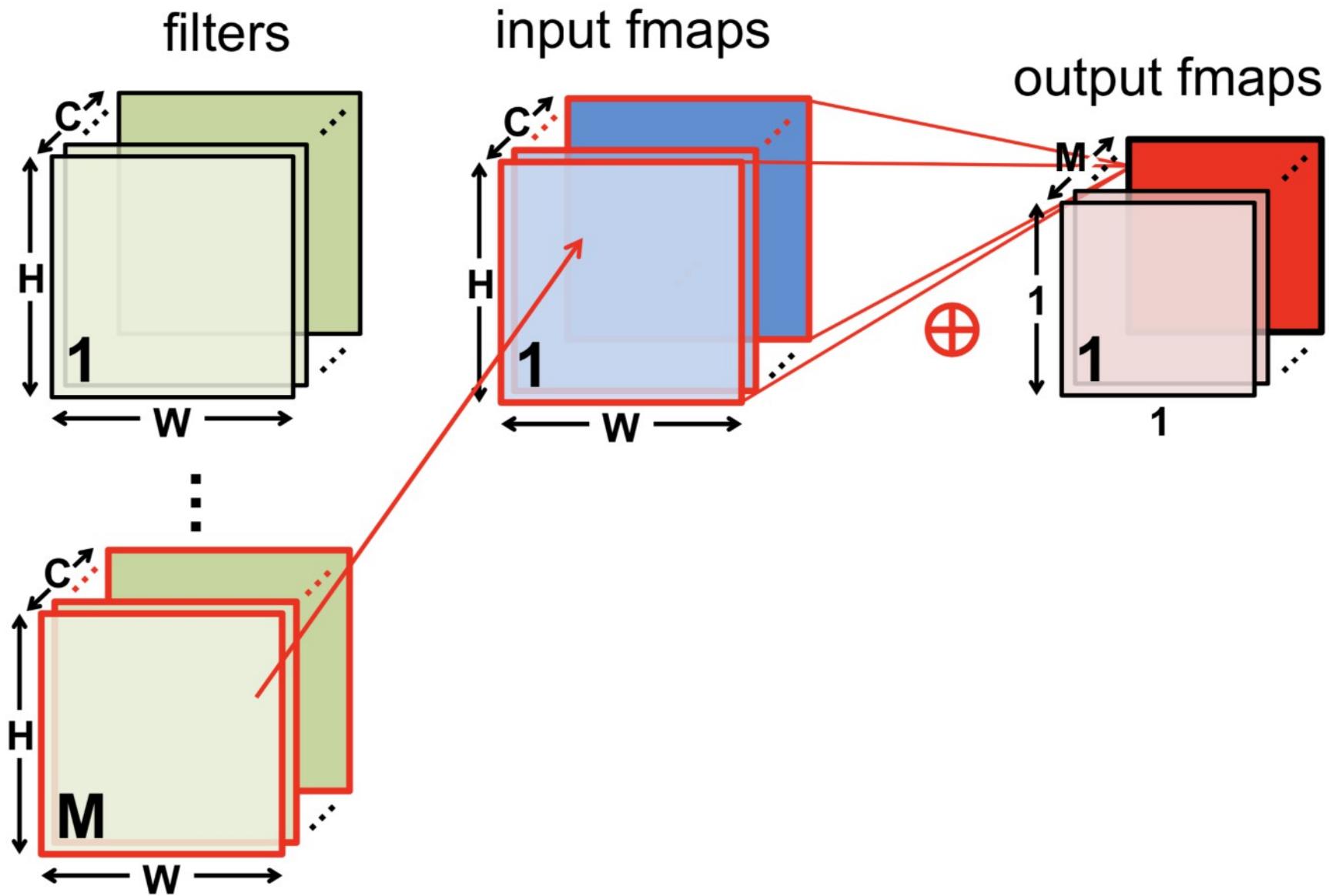
Fully Connected Layer



Fully Connected Layer

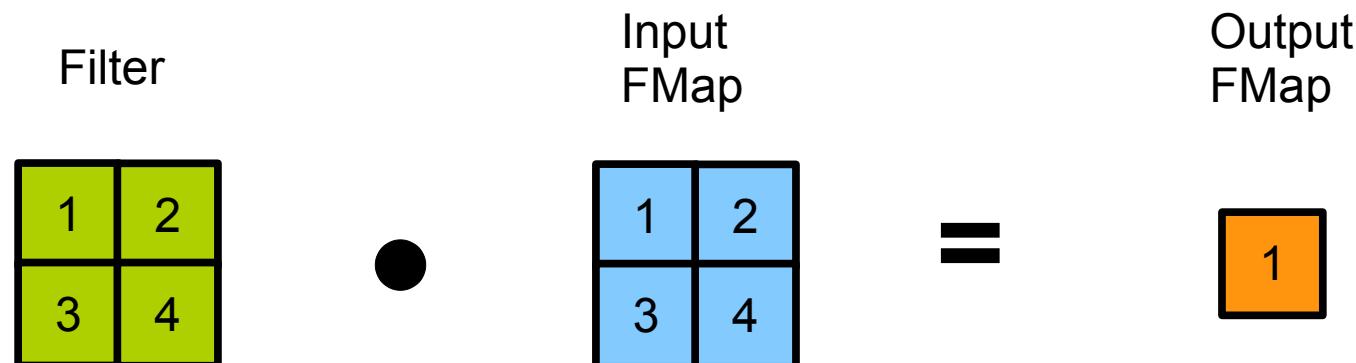


Fully Connected Layer



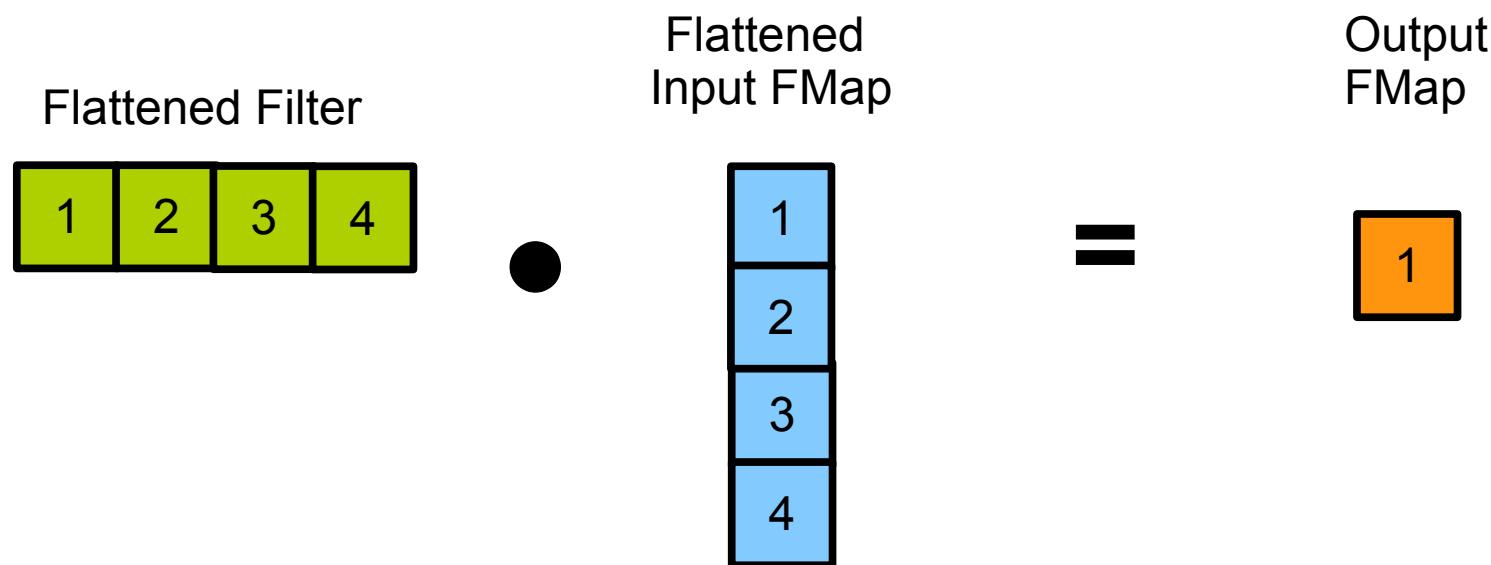
Kernel Operation

- Computing the Output Fmap **point-by-point**
 - 2D Dot Product



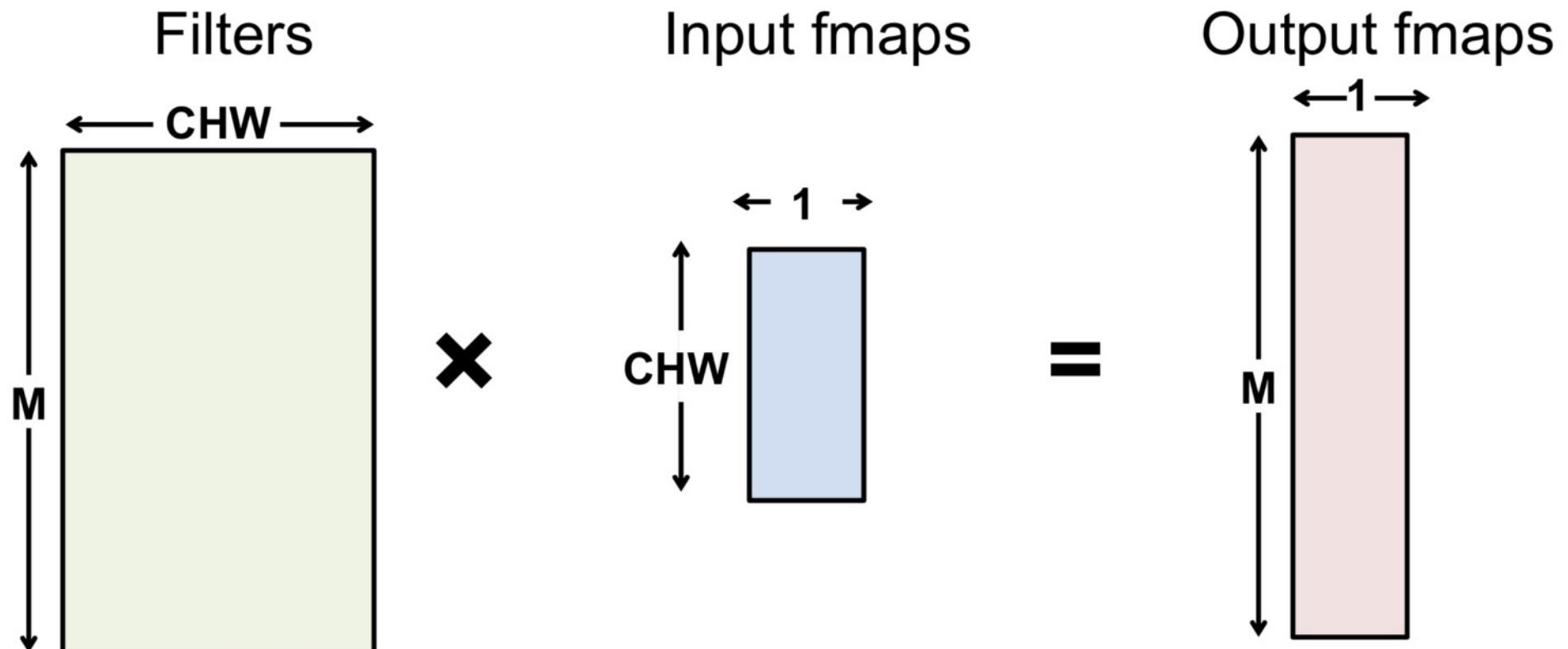
Kernel Operation

- 2D Dot Product can be converted to 1D Dot Product



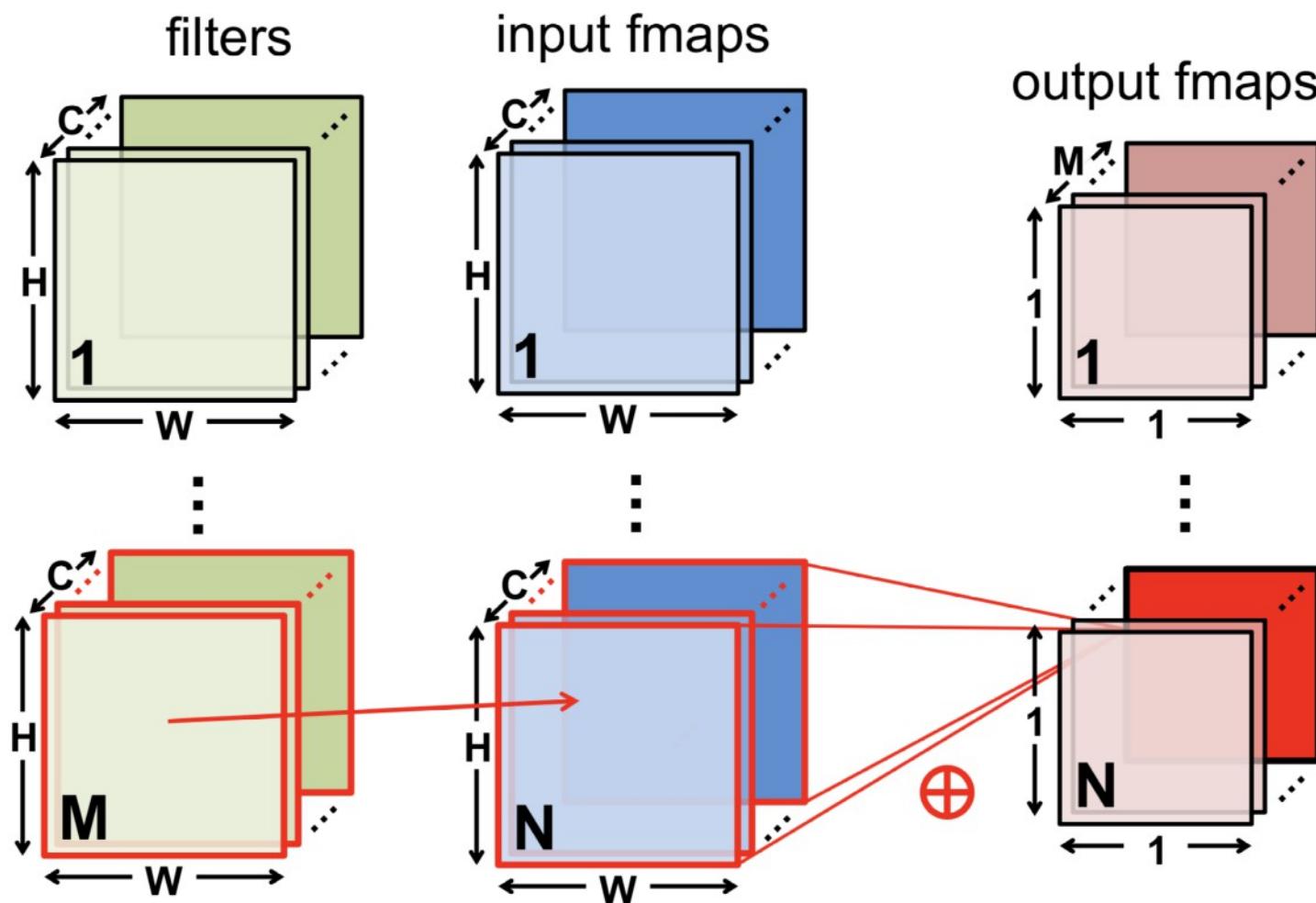
Fully Connected Layer

- Compute the Output fmap **all at once**
 - Matrix-Vector Multiply



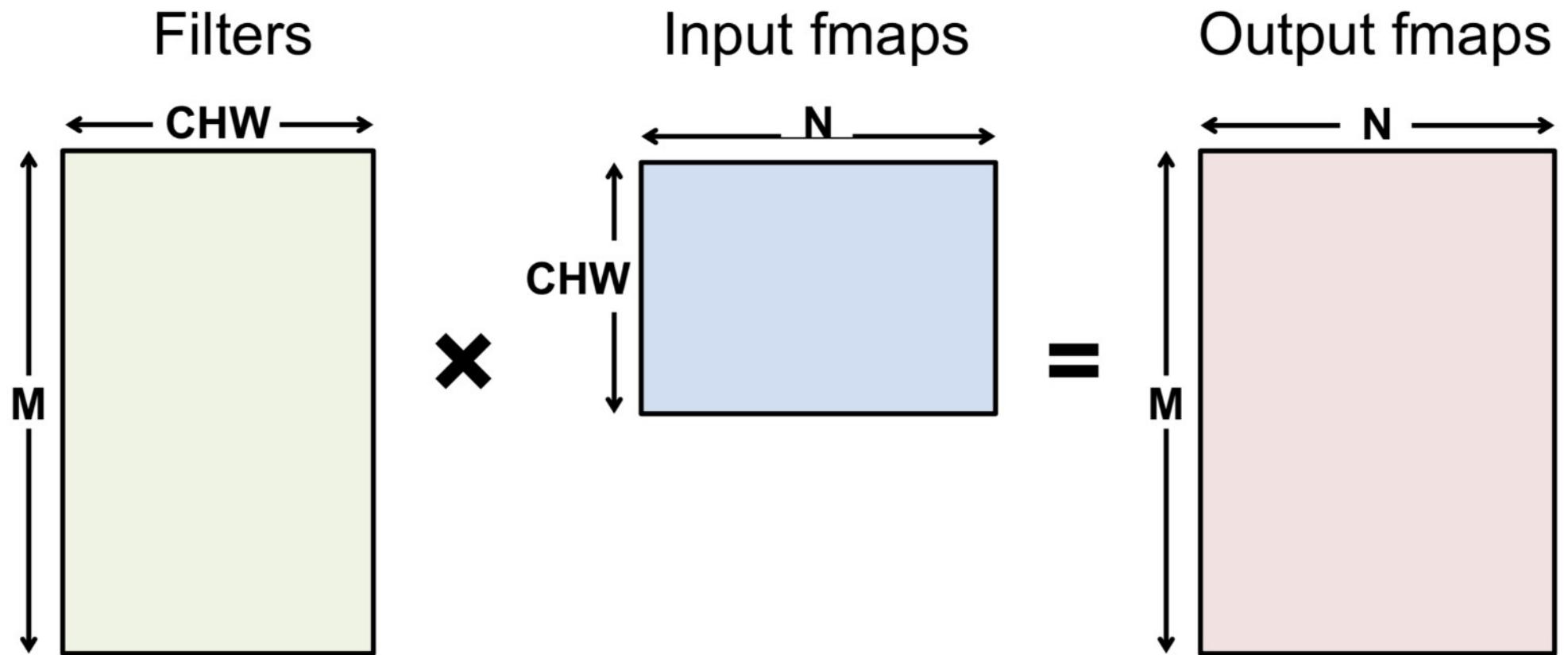
Fully Connected Layer

- Batch size > 1



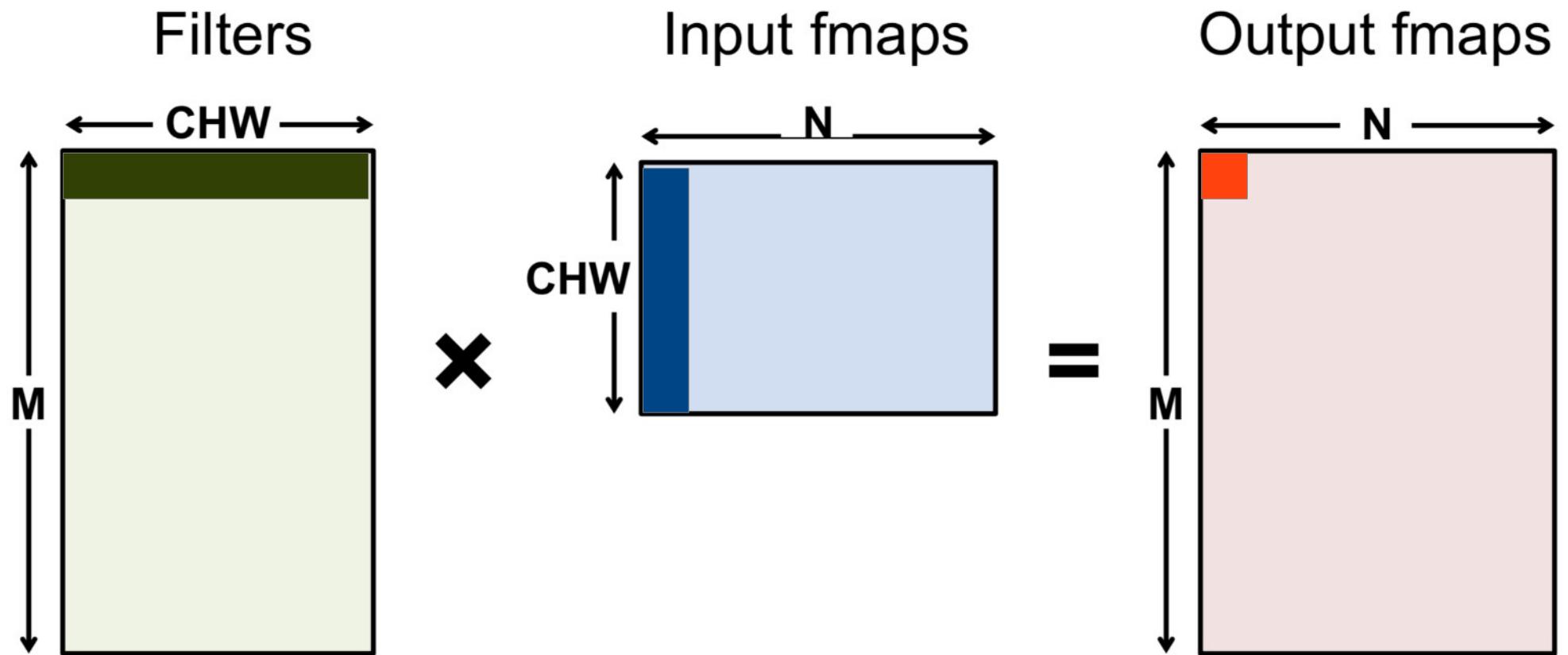
Fully Connected Layer

- Batch size > 1



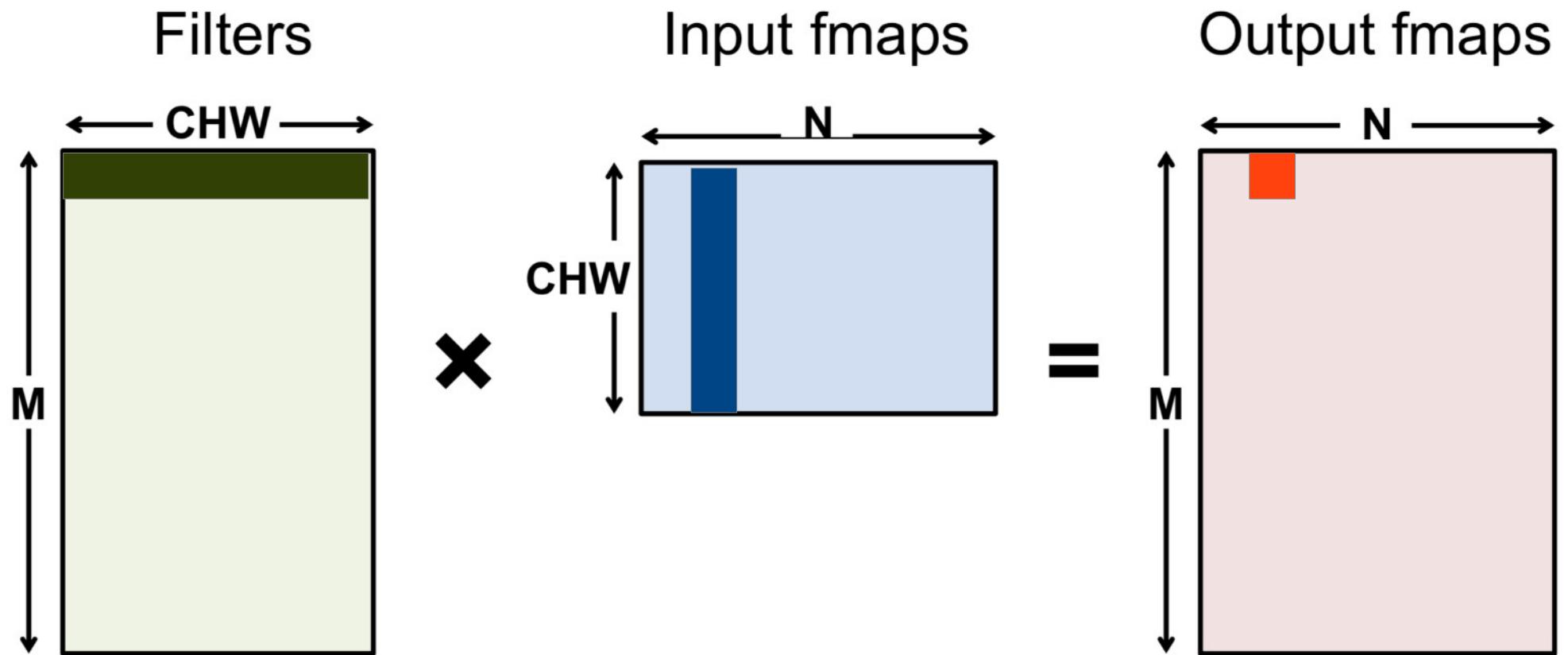
Fully Connected Layer

- Batch size > 1



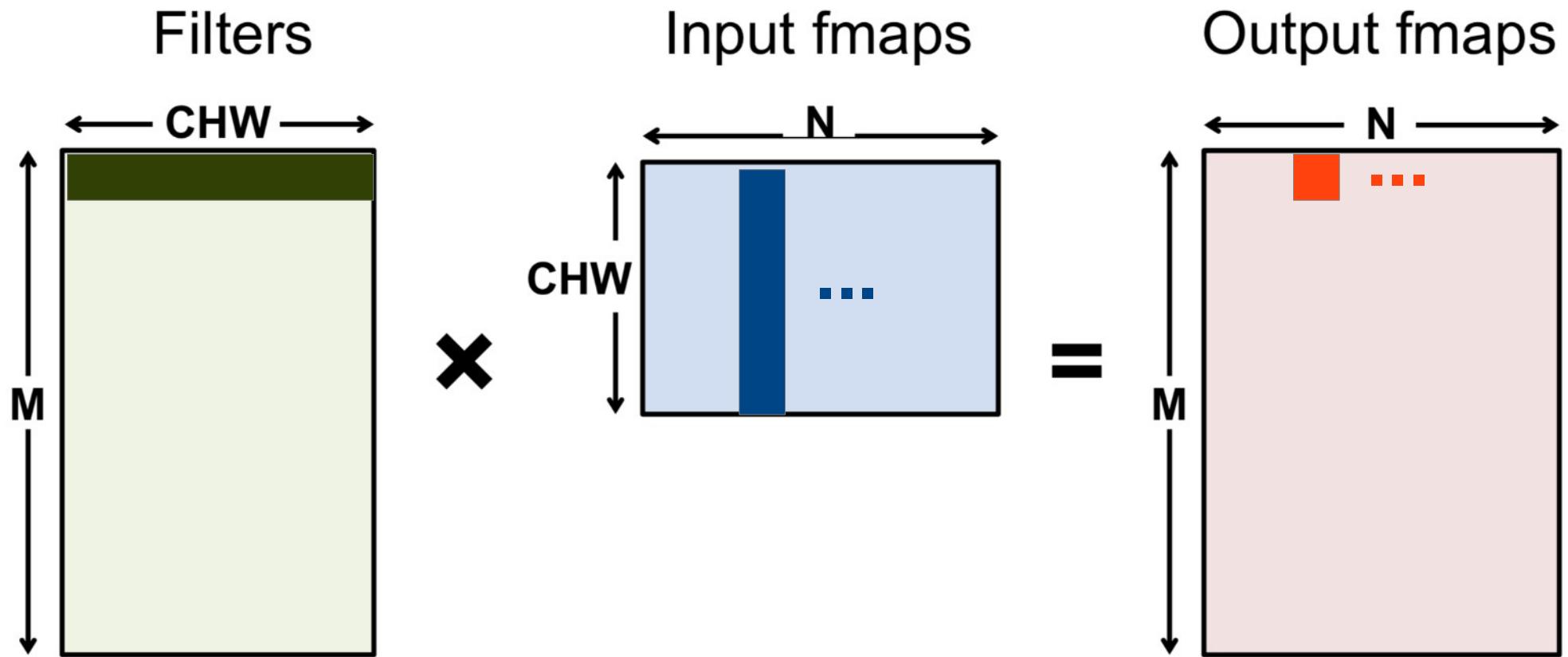
Fully Connected Layer

- Batch size > 1



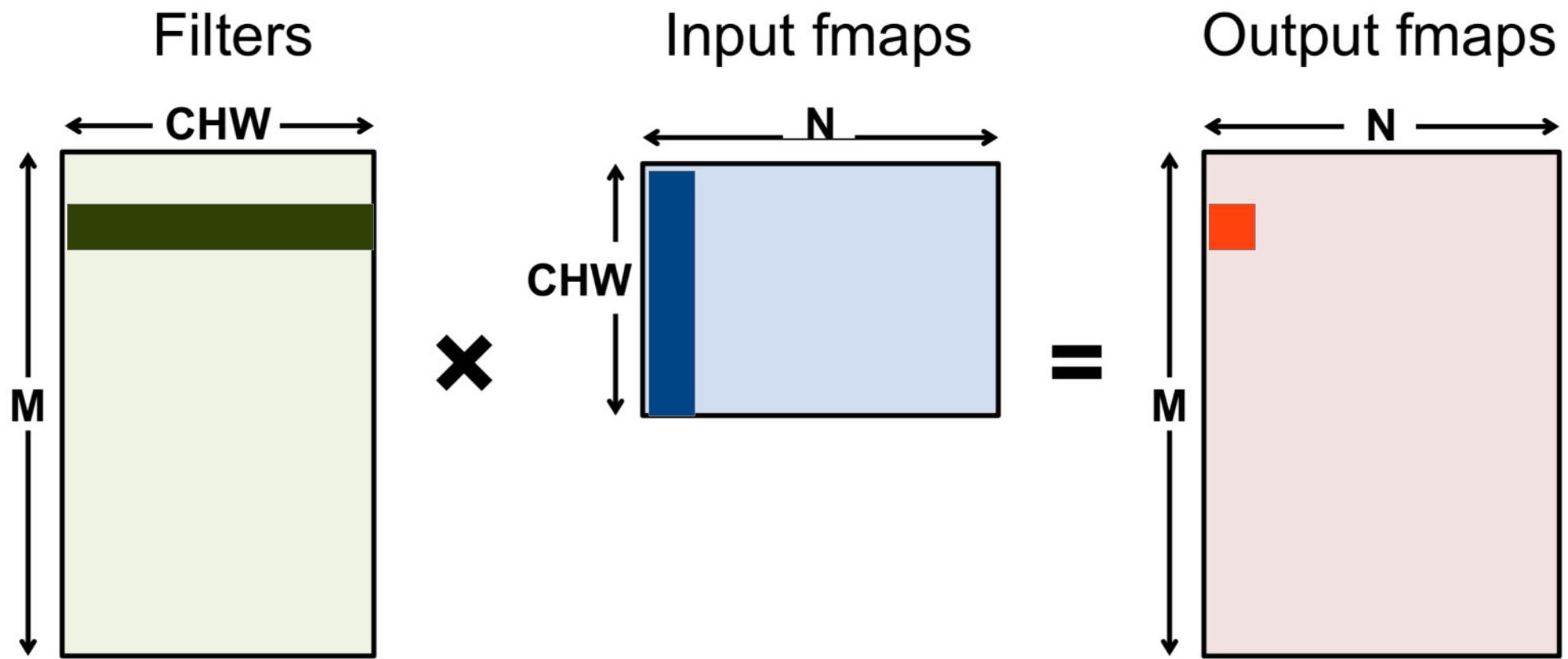
Fully Connected Layer

- Batch size > 1



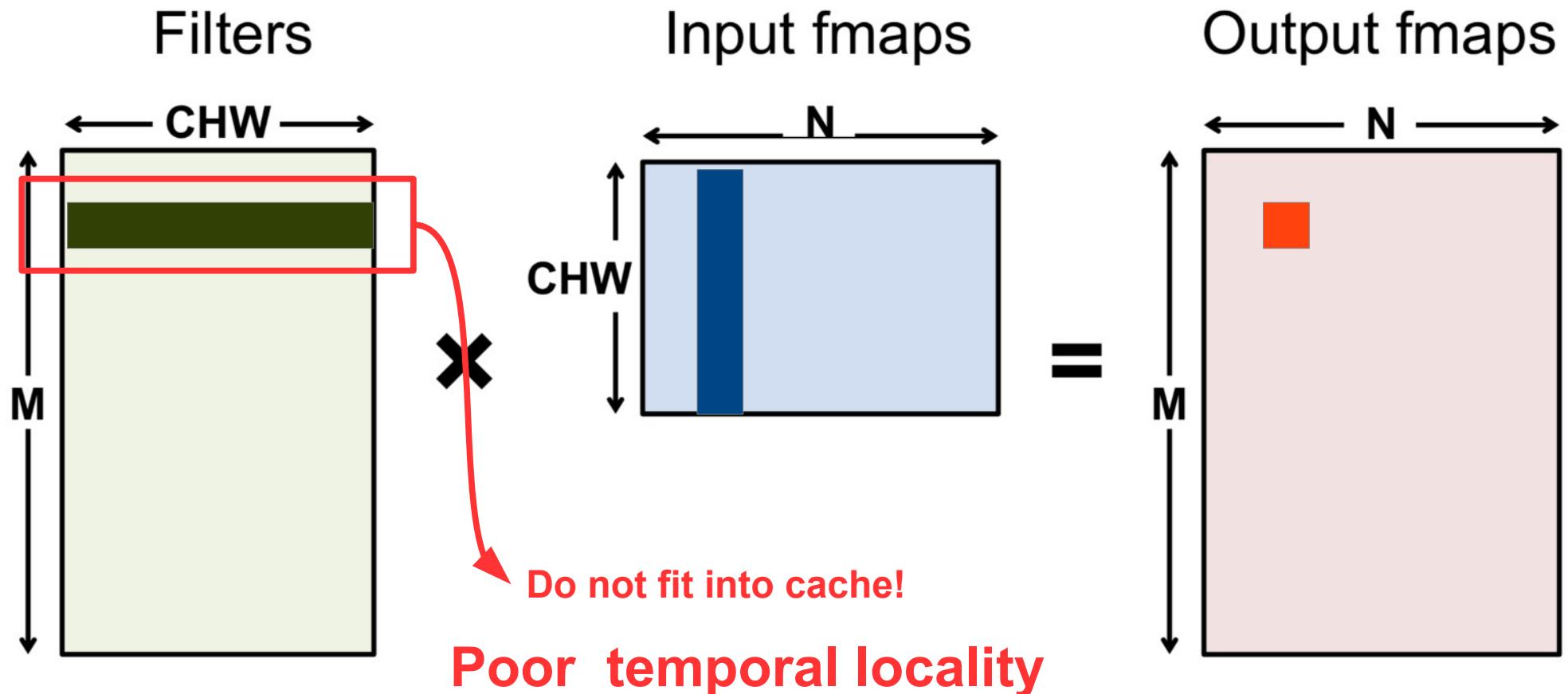
Fully Connected Layer

- Batch size > 1



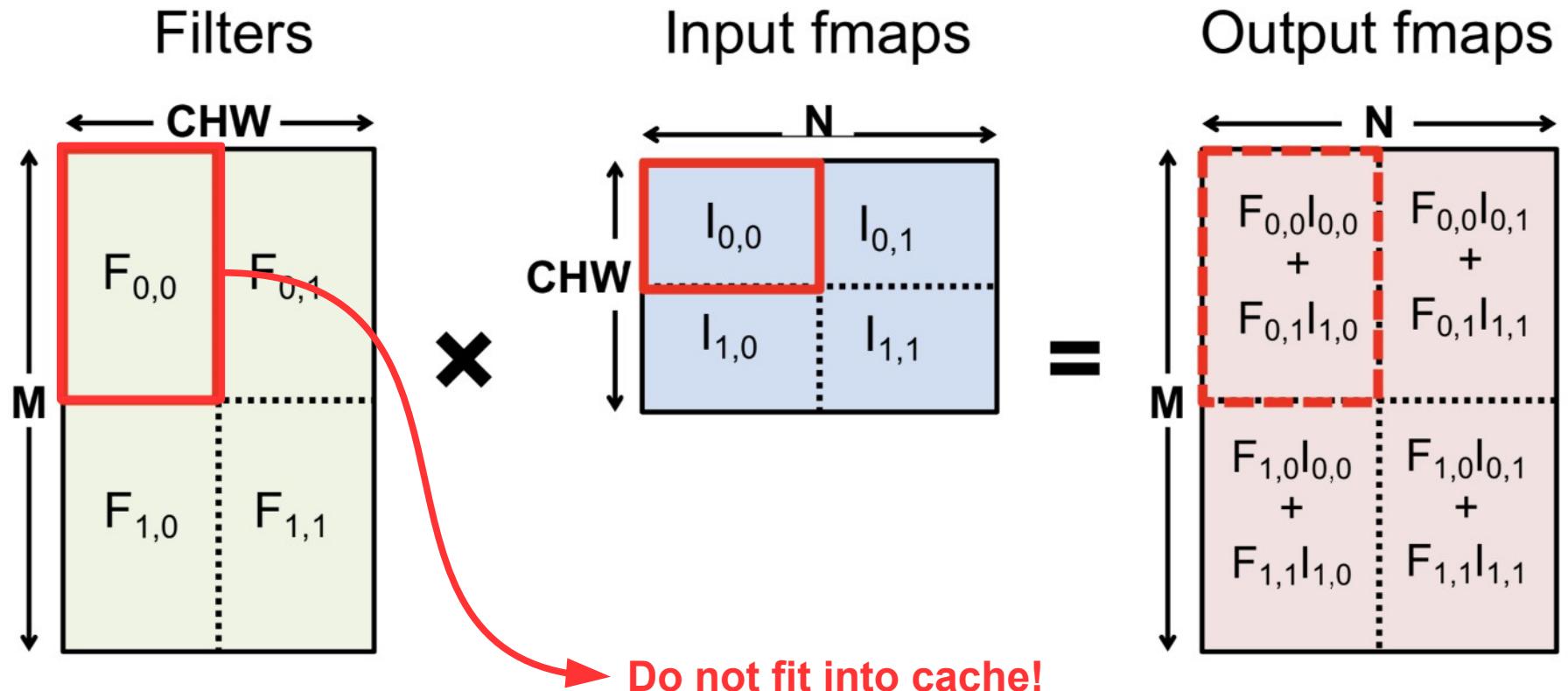
Fully Connected Layer

- Batch size > 1



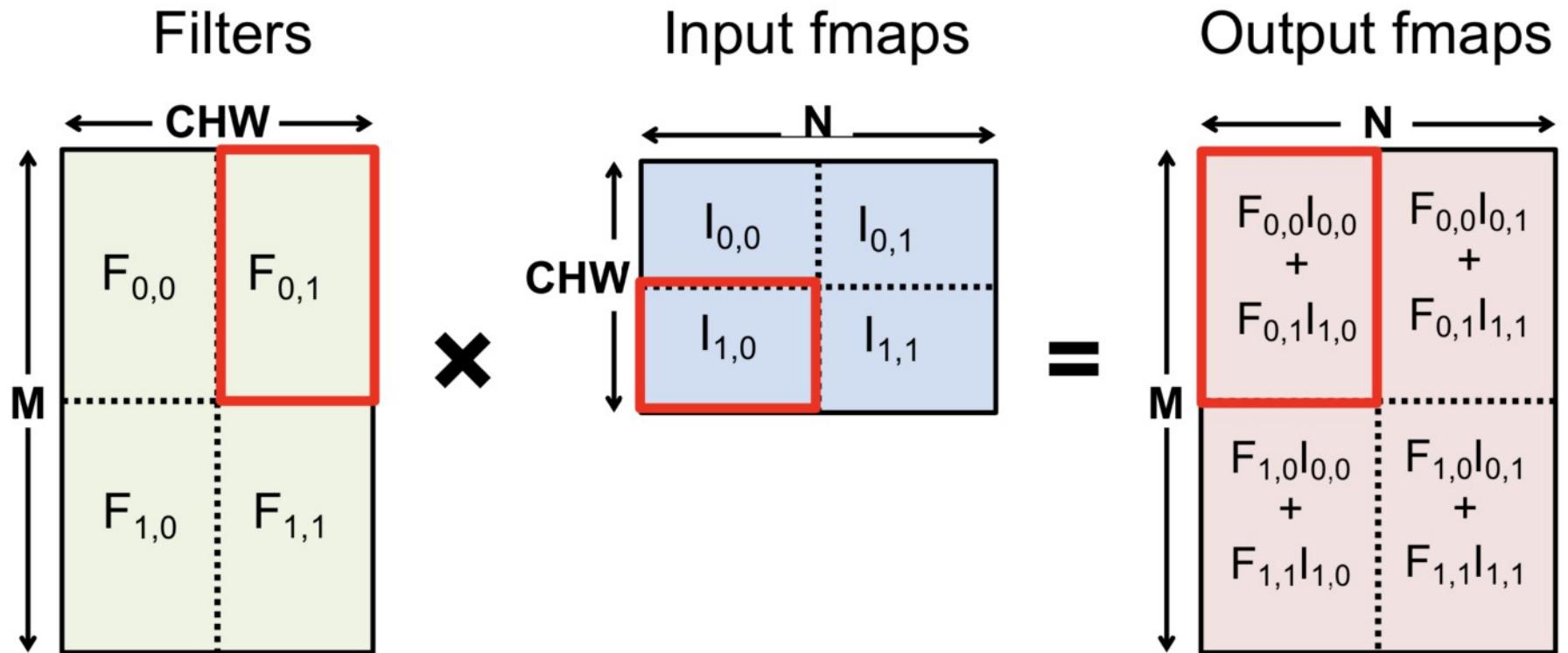
Tiling

- Matrix multiply tiled to fit in cache
- Computation ordered to maximize data reuse



Tiling

- Matrix multiply tiled to fit in cache
- Computation ordered to maximize data reuse



Implementation

- Matrix Multiplication
 - CPU: OpenBLAS, Intel MKL, ...
 - GPU: cuBLAS, cuDNN, ...
- Library will note shape of the matrix multiply and select implementation optimized for that shape
- Optimization involves proper tiling to storage hierarchy

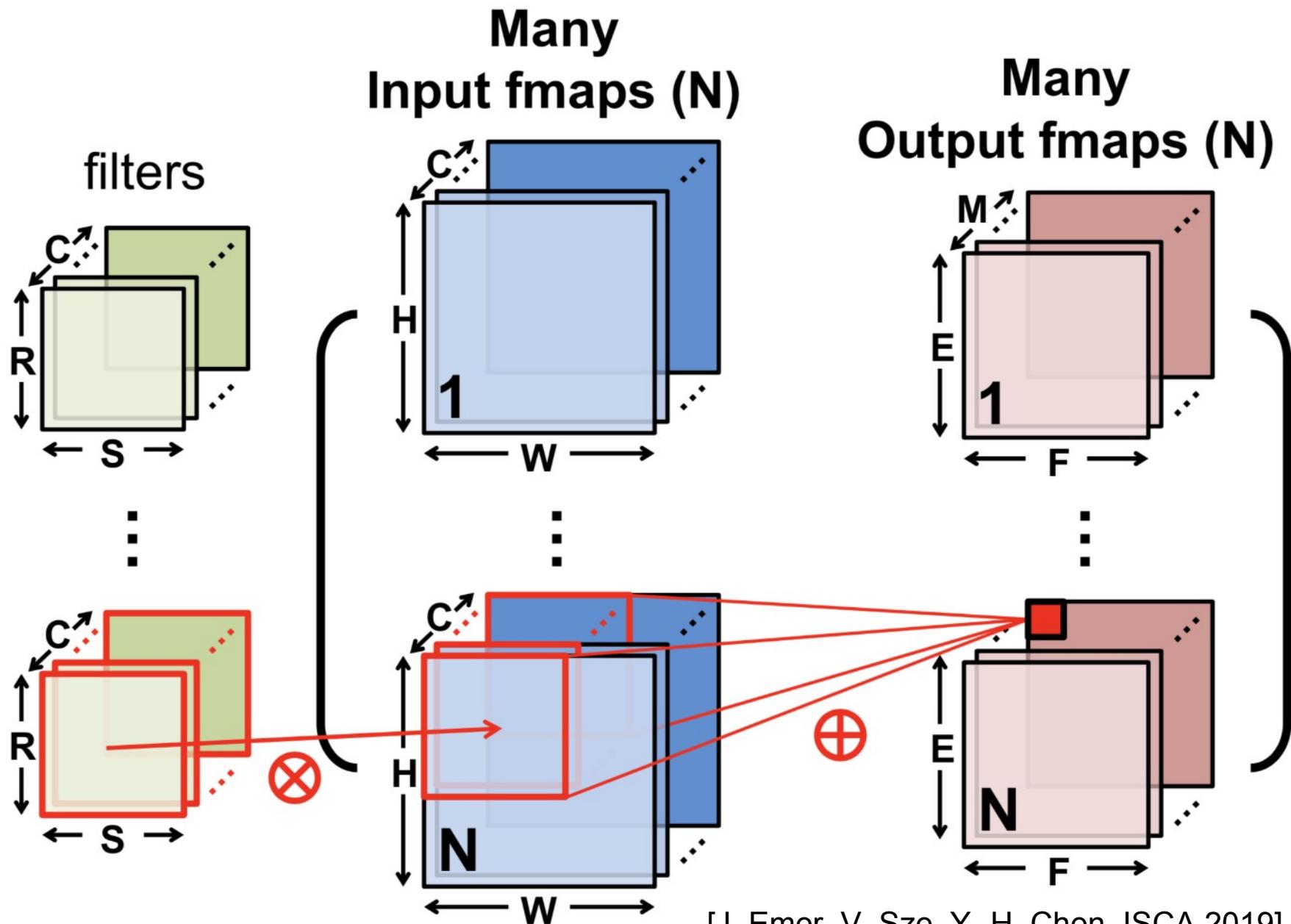
GV100 – Tensor Core

$$\mathbf{D} = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix}_{\text{FP16 or FP32}} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix}_{\text{FP16}} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}_{\text{FP16 or FP32}}$$

- Matrix Multiply Accumulate (HMMA)

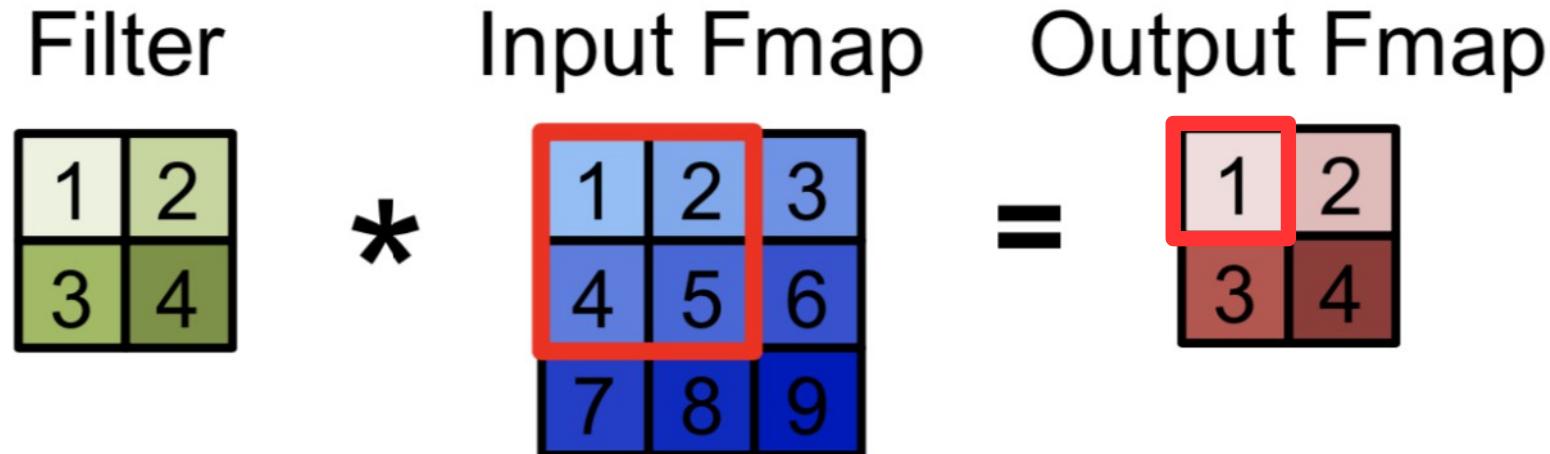
- FP16 operands → Inputs 48 / Outputs 16
 - Multiplies → 64
 - Adds → 64

Convolution Layer

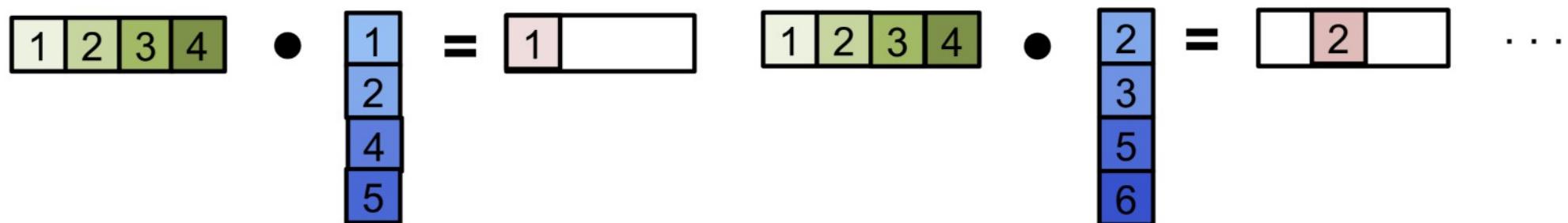


Convolution by 1D Products

Convolution



1D Product (flattened)



Convolution by Matrix Multiplication

Filter

1	2
3	4

Input Fmap

1	2	3
4	5	6
7	8	9

Output Fmap

1	2
3	4

*

=

Convolution:



Matrix Multiply (by Toeplitz Matrix)

1	2	3	4
---	---	---	---

x

1	2	4	5
2	3	5	6
4	5	7	8
5	6	8	9

1	2	3	4
---	---	---	---

Convolution by Matrix Multiplication

Filter

1	2
3	4

Input Fmap

1	2	3
4	5	6
7	8	9

Output Fmap

1	2
3	4

*

=

Convolution:



Matrix Multiply (by Toeplitz Matrix)

1	2	3	4
---	---	---	---

x

1	2	4	5
2	3	5	6
4	5	7	8
5	6	8	9

1	2	3	4
---	---	---	---

Convolution by Matrix Multiplication

Filter

1	2
3	4

*

Input Fmap

1	2	3
4	5	6
7	8	9

Output Fmap

1	2
3	4

Convolution:



Matrix Multiply (by Toeplitz Matrix)

1	2	3	4
---	---	---	---

×

1	2	4	5
2	3	5	6
4	5	7	8
5	6	8	9

1	2	3	4
---	---	---	---

Convolution by Matrix Multiplication

Filter

1	2
3	4

Input Fmap

1	2	3
4	5	6
7	8	9

Output Fmap

1	2
3	4

Convolution:



Matrix Multiply (by Toeplitz Matrix)

1	2	3	4
---	---	---	---

×

1	2	4	5
2	3	5	6
4	5	7	8
5	6	8	9

1	2	3	4
---	---	---	---

Convolution by Matrix Multiplication

Filter

1	2
3	4

Input Fmap

1	2	3
4	5	6
7	8	9

Output Fmap

1	2
3	4

*

=

Convolution:



Matrix Multiply (by Toeplitz Matrix)

1	2	3	4
---	---	---	---

×

1	2	4	5
2	3	5	6
4	5	7	8
5	6	8	9

1	2	3	4
---	---	---	---

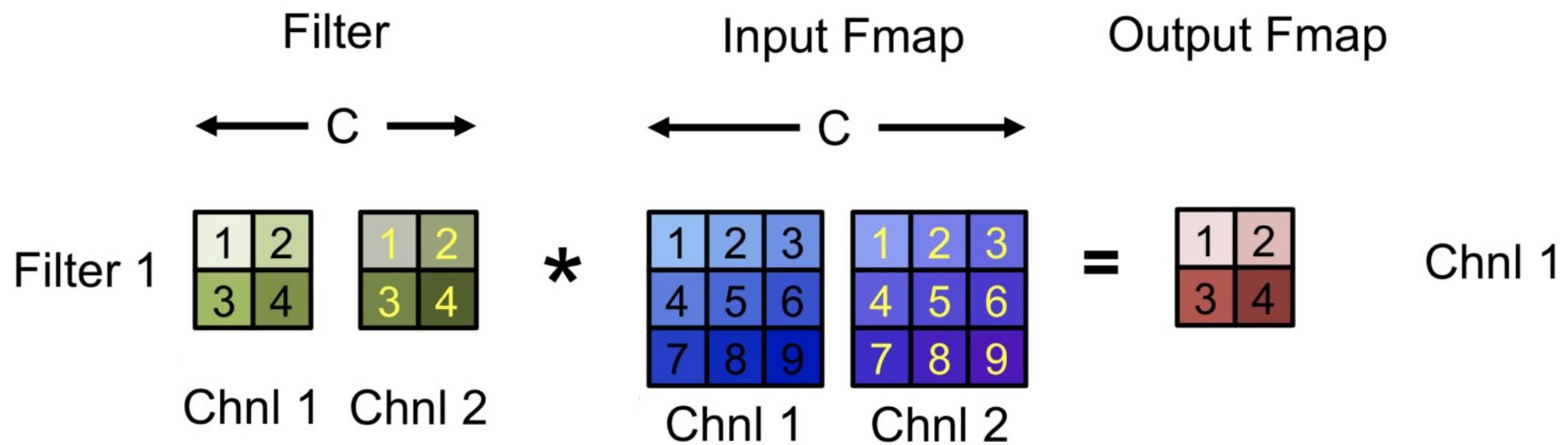
Convolution by Matrix Multiplication

$$\begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 4 & 5 \\ 2 & 3 & 5 & 6 \\ 4 & 5 & 7 & 8 \\ 5 & 6 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$

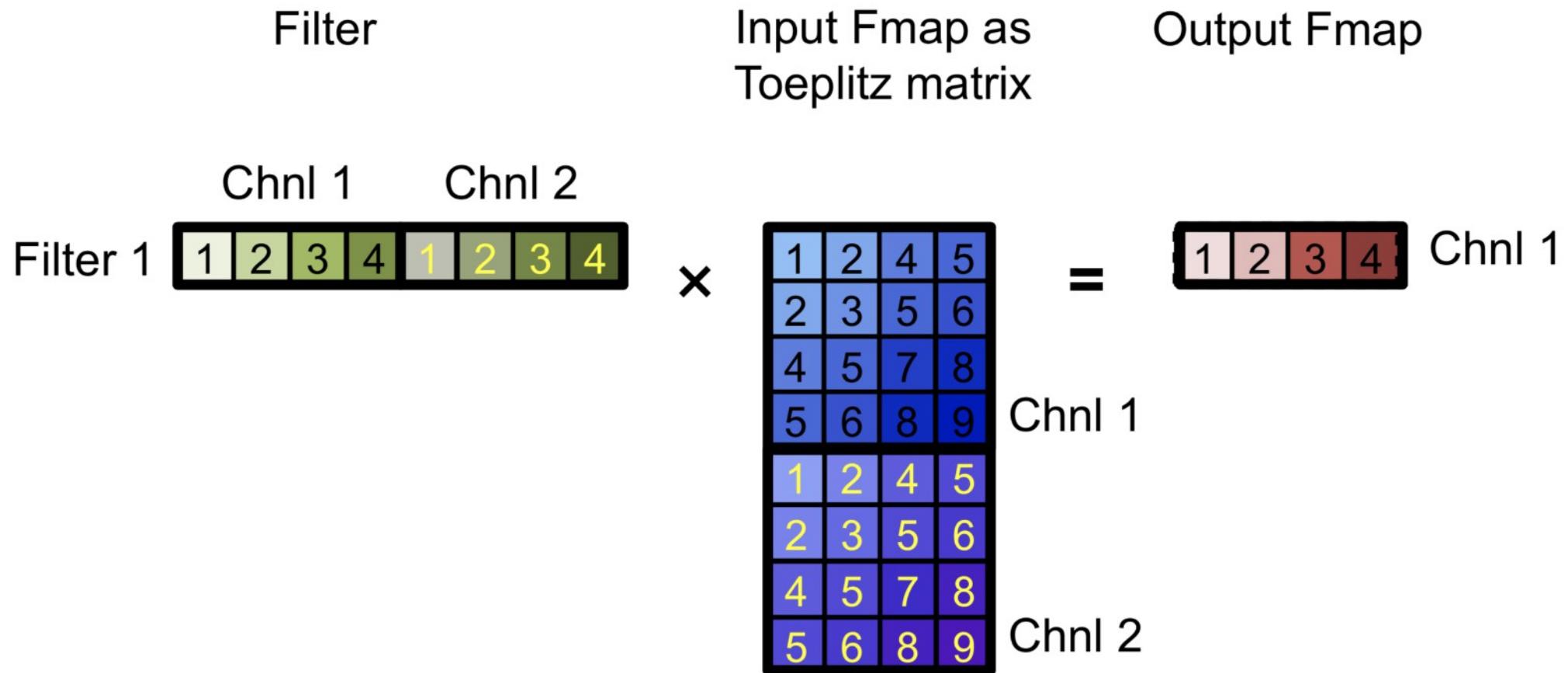
The diagram illustrates the convolution operation using matrix multiplication. On the left, a 1x4 input vector [1 2 3 4] is multiplied by a 4x4 kernel matrix. The kernel matrix has its elements colored in a repeating pattern: row 1: green, yellow, red, blue; row 2: blue, green, red, orange; row 3: orange, blue, green, yellow; row 4: red, orange, yellow, blue. The result of the multiplication is shown on the right as a 1x4 output vector [1 2 3 4], where each element is the sum of the products of the corresponding input element and the kernel elements. For example, the first element of the output is $1 \cdot 1 + 2 \cdot 2 + 3 \cdot 4 + 4 \cdot 5 = 30$.

- Data is repeated
 - Memory vs. Computation trade-off

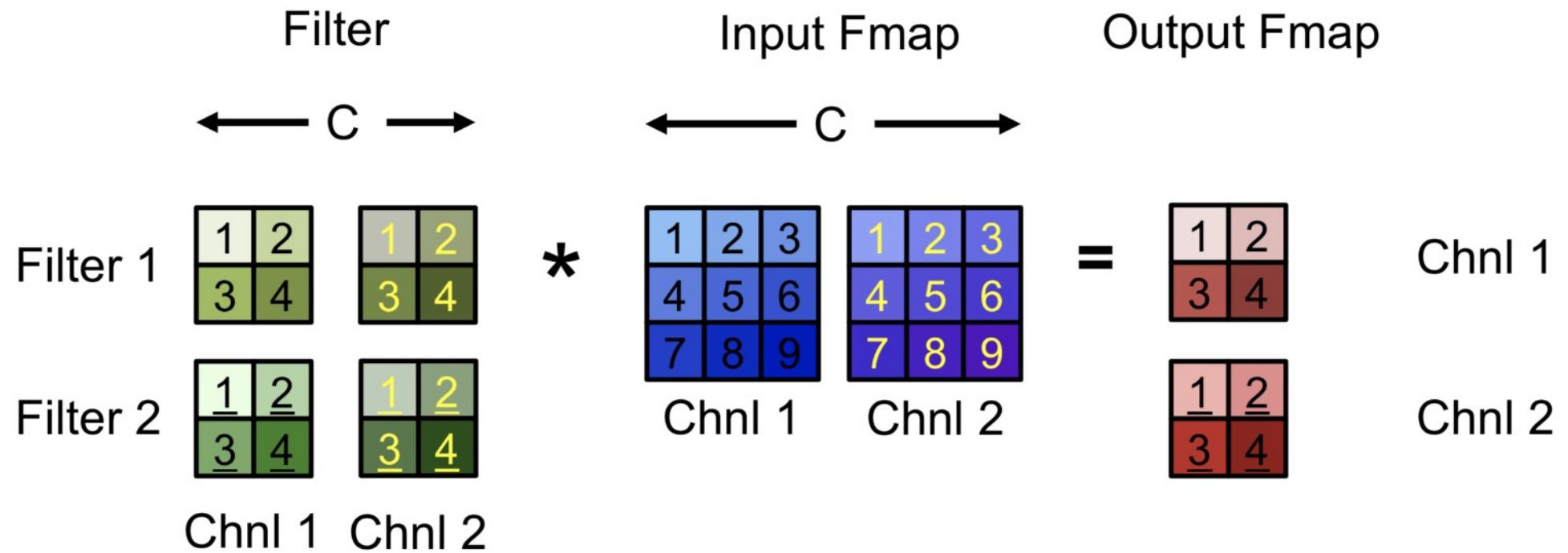
Convolution with Multiple Channels



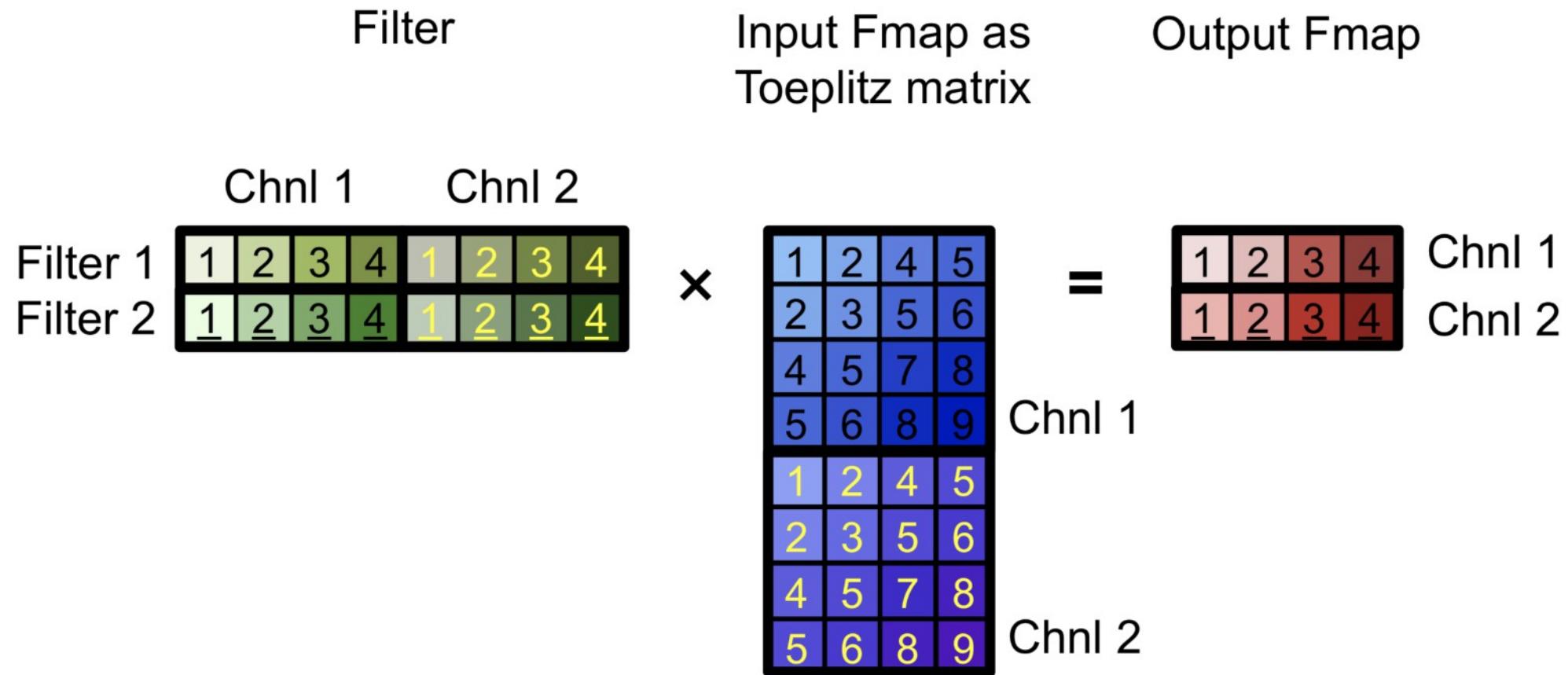
Convolution by Matrix Multiplication



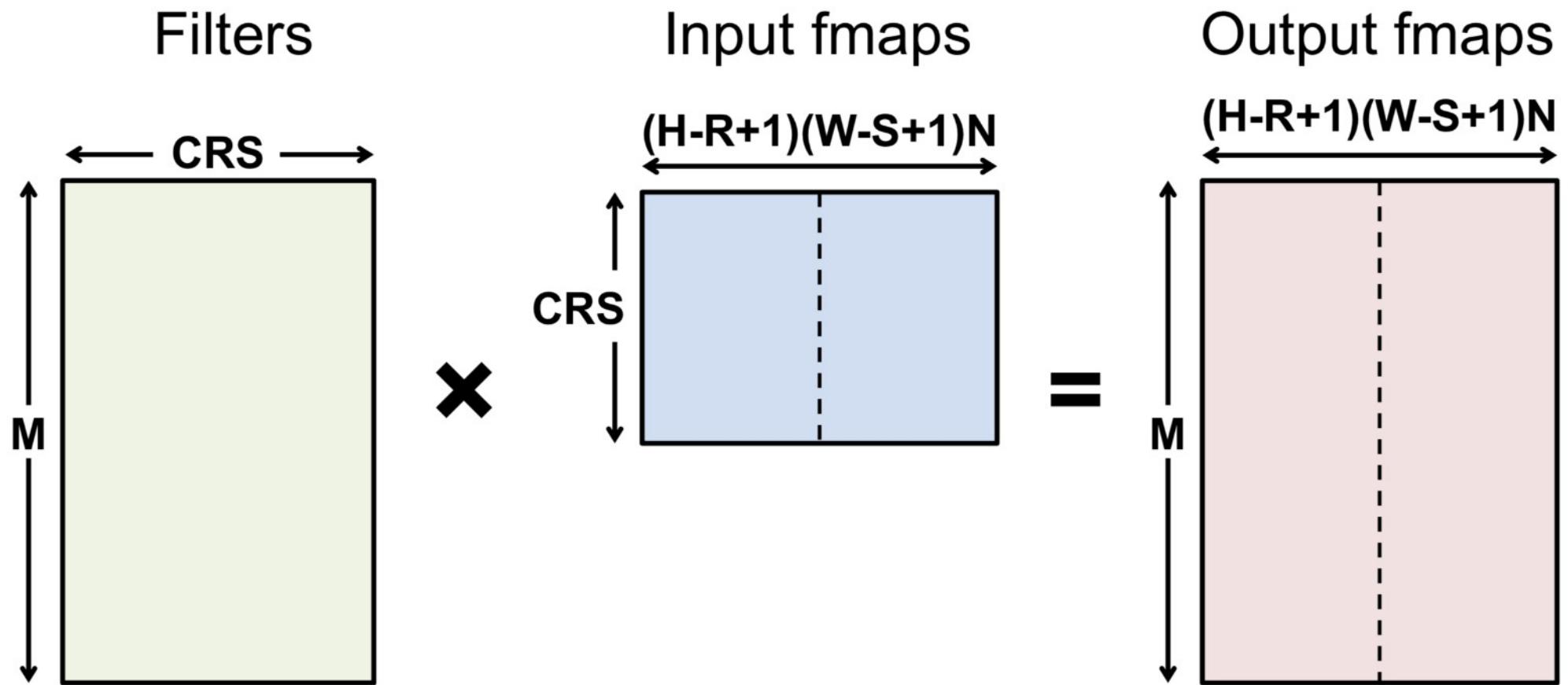
Convolution with Multiple Filters



Convolution by Matrix Multiplication



Convolution with Batch size > 1



Computation Transformations

- Obtaining the same results with less computation
 - Gauss's multiplication
 - Strassen
 - Winograd
 - FFT

Gauss's Multiplication

$$(a + bi)(c + di) = (ac - bd) + (bc + ad)i.$$

4 multiplications + 3 additions



$$k_1 = c \cdot (a + b)$$

$$k_2 = a \cdot (d - c)$$

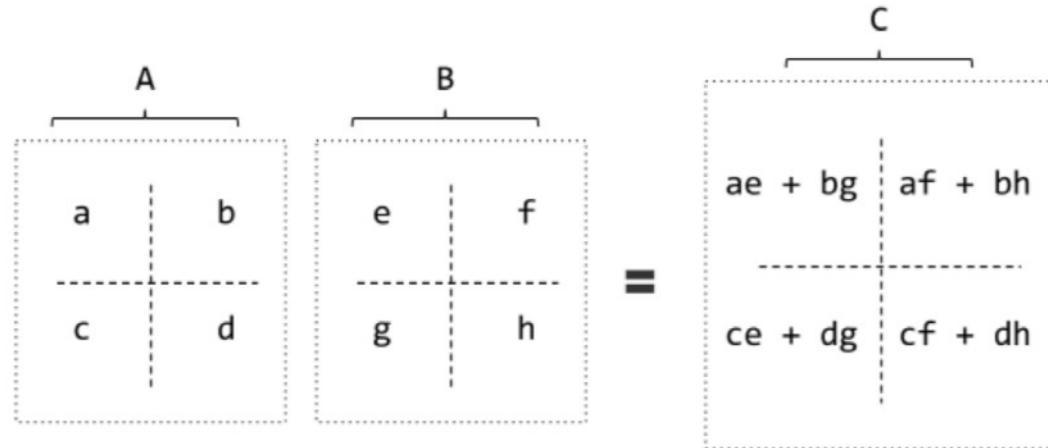
3 multiplications + 5 additions

$$k_3 = b \cdot (c + d)$$

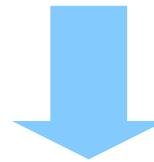
$$\text{Real part} = k_1 - k_3$$

$$\text{Imaginary part} = k_1 + k_2.$$

Strassen



8 multiplications + 4 additions



$$\begin{aligned} P1 &= a(f - h) \\ P2 &= (a + b)h \\ P3 &= (c + d)e \\ P4 &= d(g - e) \end{aligned}$$

$$\begin{aligned} P5 &= (a + d)(e + h) \\ P6 &= (b - d)(g + h) \\ P7 &= (a - c)(e + f) \end{aligned}$$

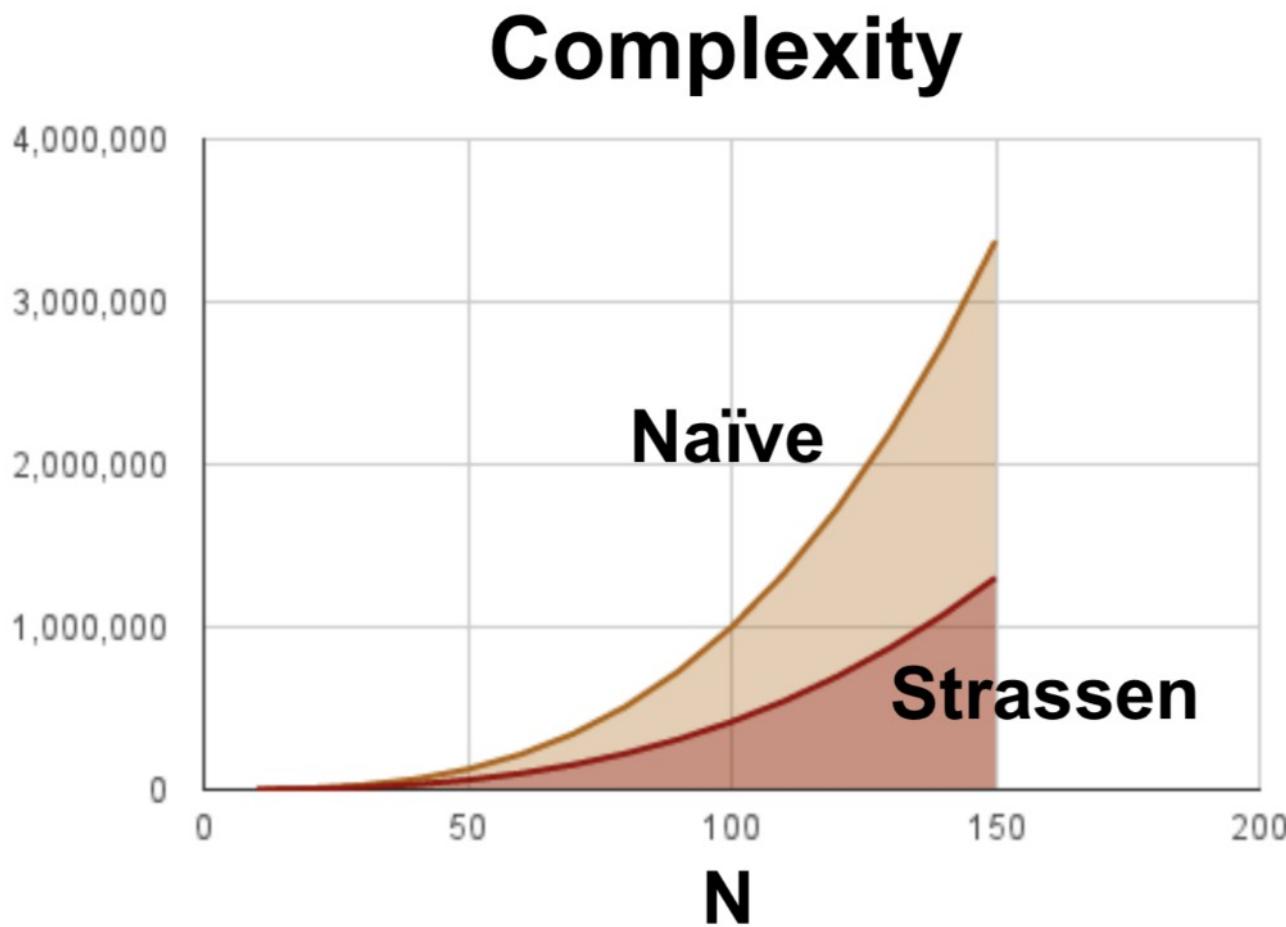
$$AB =$$

$$\begin{bmatrix} P5 + P4 - P2 + P6 & P1 + P2 \\ P3 + P4 & P1 + P5 - P3 - P7 \end{bmatrix}$$

7 multiplications + 18 additions

Strassen

- Reduce the complexity of matrix multiplication from $O(N^3)$ to $O(N^{2.807})$



- ...but
- Reduced numerical stability
 - Requires significantly more memory

Winograd 1D

- Used for convolutions

$$input = [d_0 \ d_1 \ d_2 \ d_3] \quad filter = [g_0 \ g_1 \ g_2]$$

$$output = \begin{bmatrix} d_0 \times g_0 + d_1 \times g_1 + d_2 \times g_2 \\ d_1 \times g_0 + d_2 \times g_1 + d_3 \times g_2 \end{bmatrix}$$

Original:

6 multiplications, 4 additions

Winograd 1D

- Used for convolutions

$$\begin{array}{c} \text{input} \\ \left[\begin{matrix} d_0 & d_1 & d_2 \\ d_1 & d_2 & d_3 \end{matrix} \right] \end{array} \begin{array}{c} \text{filter} \\ \left[\begin{matrix} g_0 \\ g_1 \\ g_2 \end{matrix} \right] \end{array} = \left[\begin{matrix} m_1 + m_2 + m_3 \\ m_2 - m_3 - m_4 \end{matrix} \right]$$

$$m_1 = (d_0 - d_2)g_0$$

$$m_4 = (d_1 - d_3)g_2$$

$$m_2 = (d_1 + d_2) \frac{g_0 + g_1 + g_2}{2}$$

$$m_3 = (d_2 - d_1) \frac{g_0 - g_1 + g_2}{2}$$

Winograd:

4 multiplications, 12 additions, 2 shifts

Filter weights are constant. Only need to be performed once!

Winograd 1D – Observation

$$\begin{array}{c} \text{input} \\ \left[\begin{matrix} d_0 & d_1 & d_2 \\ d_1 & d_2 & d_3 \end{matrix} \right] \end{array} \quad \begin{array}{c} \text{filter} \\ \left[\begin{matrix} g_0 \\ g_1 \\ g_2 \end{matrix} \right] \end{array} = \left[\begin{matrix} m_1 + m_2 + m_3 \\ m_2 - m_3 - m_4 \end{matrix} \right]$$

- Works on a small part of input
 - 4 elements in this example
- An integer convolution requires the application of Winograd on a tile-by-tile basis
- **Reduction of multiplications depends on filter/tile size**

Winograd 2D

Filter			Input Fmap				Output Fmap		
g_{00}	g_{01}	g_{02}	\ast				y_{00}	y_{01}	
g_{10}	g_{11}	g_{12}	d_{00}	d_{01}	d_{02}	d_{03}	$=$	y_{10}	y_{11}
g_{20}	g_{21}	g_{22}	d_{10}	d_{11}	d_{12}	d_{13}			
			d_{20}	d_{21}	d_{22}	d_{23}			
			d_{30}	d_{31}	d_{32}	d_{33}			

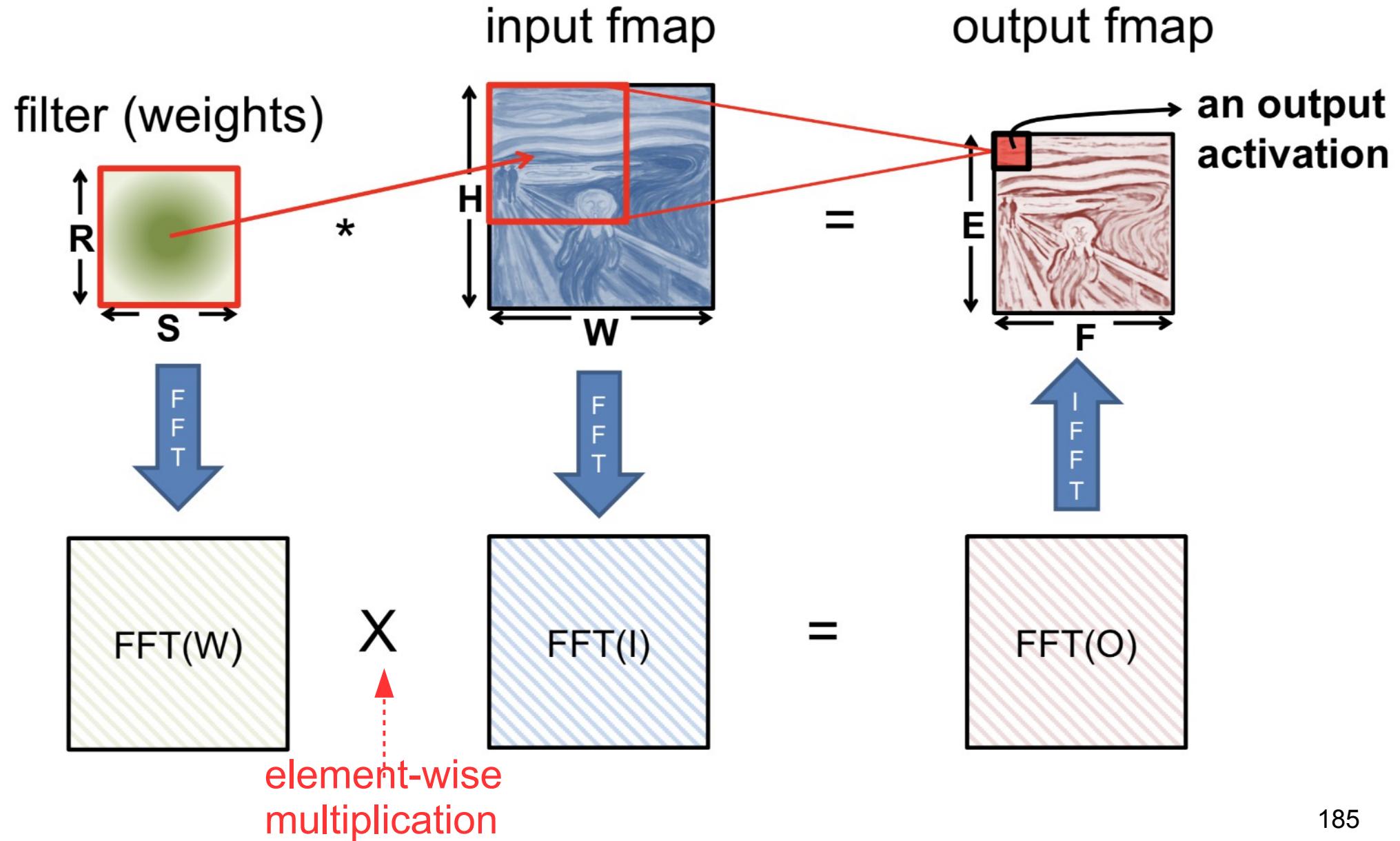
Original:
Winograd:

36 multiplications
16 multiplications → Speedup 2.25

Winograd Summary

- Pros
 - Optimized computation for convolutions
 - Significantly reduce multiplies
- Cons
 - Each filter size (and output size) is a different computation

Fast Fourier Transform



FFT Cost

- Complexity
 - Convolution $O(RSEF)$
 - Multiplication $O(EF \log_2 EF)$
- Computational benefit of FFT decreases with decreasing size of filter
- Large storage capacity and bandwidth
- Coefficients in the frequency domain are complex

FFT Optimization and Trade-offs

- FFT of real matrix is symmetric
 - Save 1/2 of the computations
- Filters can be pre-computed and stored
 - In frequency domain they are much larger than in space domain

Agenda

- Deep Neural Networks
- Types of layers
- Memory and communication traffic
- Kernel computation
- Quantitative analysis

Key Metrics and Design Objectives

- Accuracy
- Throughput and Latency
- Energy
- ...

Throughput

- Number of execution of a task that can be completed in a given time period
 - e.g, infereces per second

Latency

- Time between the beginning of a task and its completion
 - e.g., inference latency

Fallacies

- If throughput increases latency decreases and viceversa
- Example
 - Want throughput of 30 fps
 - With a batch size of 100 frames → Amortize overhead due to load weights → 100 frames in 1s (more than 30 fps required)
 - What about latency for a single frame?
 - Results after filling up the batch
 - Batch filled up in 3.3 sec
 - Latency > 3.3 sec! Not acceptable for real-time applications (e.g., high-speed navigation)

Throughput – Quantitative Analysis

$$\frac{\text{inferences}}{\text{second}} = \frac{\text{operations}}{\text{second}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}}$$

Depends on

- DNN hardware
- DNN model

Depends on

- DNN model

Throughput – Quantitative Analysis

$$\frac{\text{operations}}{\text{second}} = \frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}}$$

Peak throughput
(aggregate for **all** PEs)

Throughput – Quantitative Analysis

$$\frac{\text{operations}}{\text{second}} = \left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right) \times \text{number of PEs}$$

Peak throughput
(aggregate for **single** PEs)

Throughput – Quantitative Analysis

$$\frac{\text{operations}}{\text{second}} = \left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right) \times \text{number of PEs} \times \text{utilization of PEs}$$

A red dotted box highlights the term $\frac{1}{\frac{\text{cycles}}{\text{operation}}}$. A red arrow points from the text "Peak throughput (aggregate for single PEs)" below to this highlighted term.

Peak throughput
(aggregate for **single** PEs)

Throughput – Quantitative Analysis

$$\frac{\text{operations}}{\text{second}} = \left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right) \times \text{number of PEs} \times \text{utilization of PEs}$$

Peak throughput of a PE Parallelism Degradation due to the inability to utilize the PE

Throughput – Quantitative Analysis

$$\frac{\text{operations}}{\text{second}} = \left(\frac{1}{\frac{\text{cycles}}{\text{operation}}} \times \frac{\text{cycles}}{\text{second}} \right) \times \text{number of PEs} \times \text{utilization of PEs}$$

Decrease the cycles for performing an operation (e.g., non-pipelined multi-cycle MAC vs. pipelined MAC)

Increase the clock frequency (e.g., micro-architectural improvement)

- If area is fixed
- Reducing PE area
 - Trade-off with storage area → impact on PEs utilization

Mapping and data-flow technique

Throughput – Quantitative Analysis

$$\text{utilization of PEs} = \frac{\text{number of active PEs}}{\text{number of PEs}} \times \text{utilization of active PEs}$$

Ability to distribute the workload among PEs

How efficiently active PEs process the workload

The equation illustrates the components of throughput analysis. The first part, the fraction, represents the ability to distribute work among available processing elements. The second part, the multiplication by utilization, represents how effectively those active elements process the assigned work.

Throughput – Quantitative Analysis

$$utilization\ of\ PEs = \frac{number\ of\ active\ PEs}{number\ of\ PEs} \times utilization\ of\ active\ PEs$$

- Number of PEs that receive work
- **Flexibility** of the architecture plays a key role
 - On-chip network able to support different layer shapes
- **Mapping**
 - Placement and scheduling in place and time of operations
 - Delivery of the data to PEs

Throughput – Quantitative Analysis

$$\text{utilization of PEs} = \frac{\text{number of active PEs}}{\text{number of PEs}} \times \text{utilization of active PEs}$$

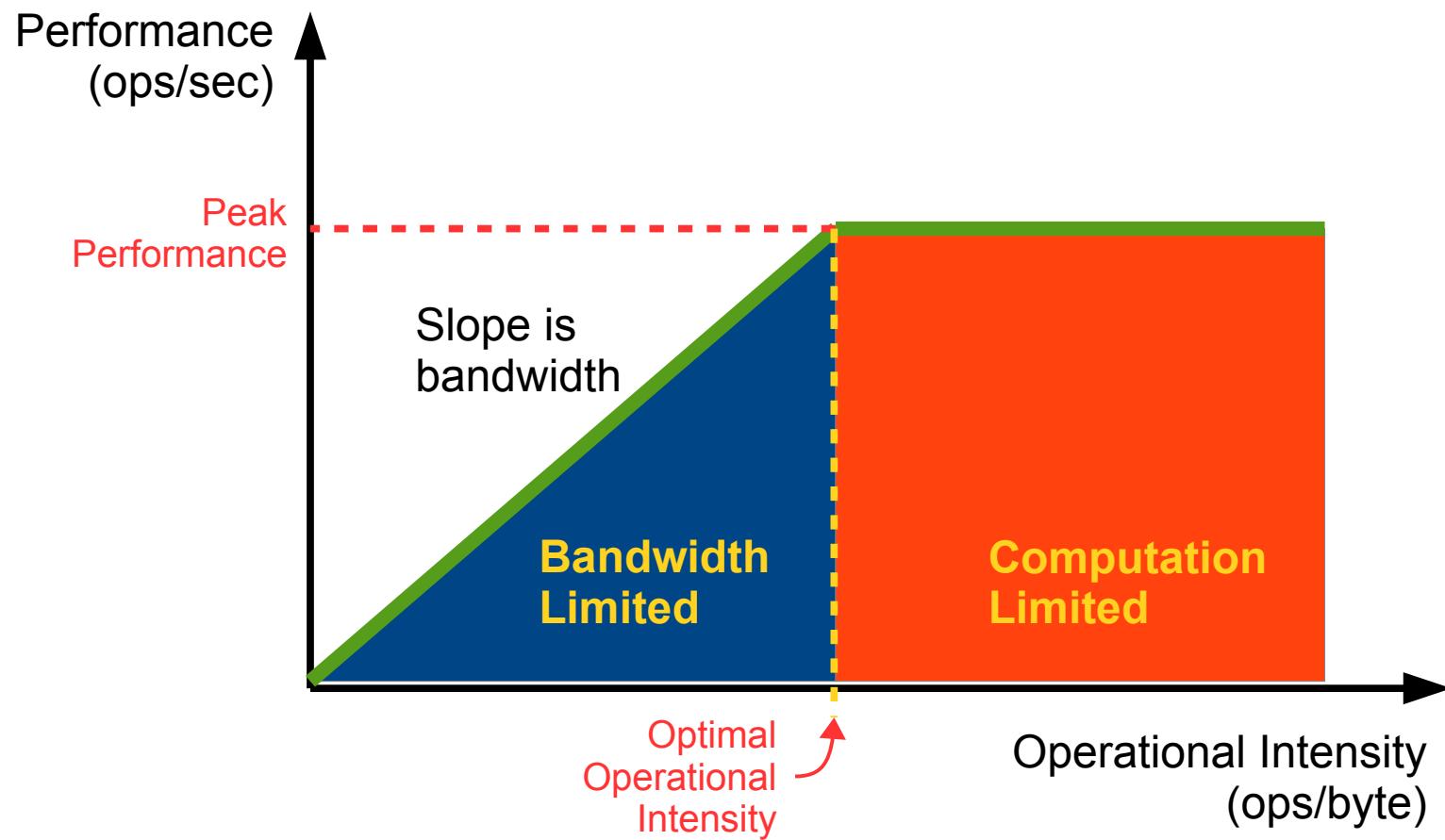
- Timely arrival of operands to PEs
 - PEs don't become idle while waiting for data to arrive
- Causes that affect utilization
 - Bandwidth and latency of memory and communication
 - Affected by the amount of data reuse
 - Data-flow techniques
 - Use of large batch size
 - Imbalance of work allocated across PEs
 - *E.g.*, due to sparsity

Throughput – Quantitative Analysis

$$\text{utilization of PEs} = \frac{\text{number of active PEs}}{\text{number of PEs}} \times \text{utilization of active PEs}$$

- Interplay PEs utilization vs. number of PEs
 - Use local memory as buffer for avoid PEs waiting for data → increase PE size → reduce the number of PEs

Roofline Model



Throughput – Other Observations

$$\frac{\text{inferences}}{\text{second}} = \frac{\text{operations}}{\text{second}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}}$$

Depends on

- DNN hardware
- DNN model

Depends on

- DNN model

- Efficient Network Architectures
 - DNN models with efficient layer shapes
 - Reduce MAC operations → improve ops/inference
 - **...but** wide range of layer shapes → might result in poor utilization of PEs → reduce ops/second

Throughput – Other Observations

$$\frac{\text{inferences}}{\text{second}} = \boxed{\frac{\text{operations}}{\text{second}}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}}$$

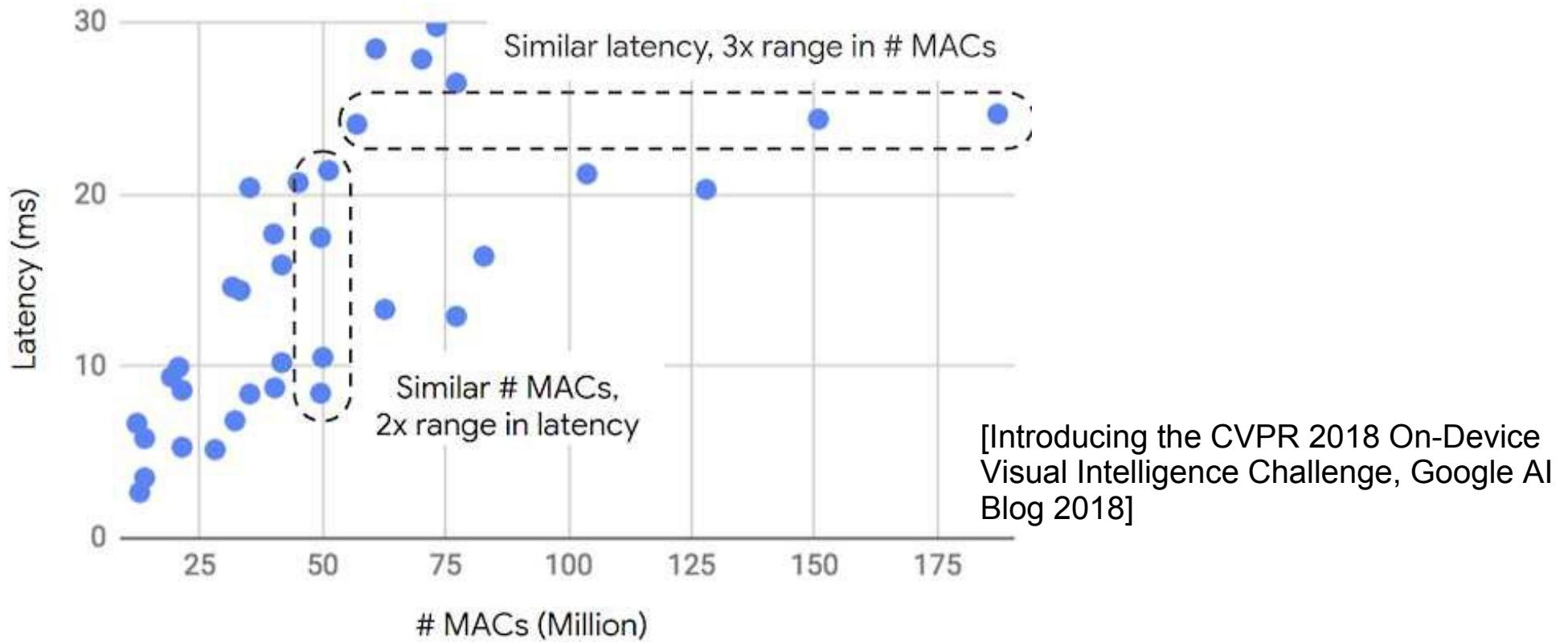
- Not all the operations are the same
 - e.g., multiply by 0 → *ineffectual operations*
 - HW exploiting ineffectual operations → increase ops/sec
 - **...but** overhead in HW → increase of the critical path and/or decrease the number of PEs if area fixed

Throughput – Other Observations

- Designing HW and DNN models that support reduced precision
 - Increase the number of operations per second
 - Required memory bandwidth decreases → Increase PE utilization
 - PE area decreases → Increase the number of PEs
 - **...but**, if multiple level of precision need to be supported → HW overhead → Increase of the critical path, etc.

Throughput – Summary

- Number of MAC operations alone **is not sufficient** for evaluating throughput (and latency)



- Need for designing DNN models with HW in the loop

Energy

- Importance of energy in DNN processing
 - Edge vs. Cloud processing
 - Privacy, latency, comm. BW limitation...
- Energy efficiency
 - Amount of work that can be completed with for unit of energy

Power vs. Energy

- Power consumption
 - Energy consumed per unit time
 - Related to heat dissipation
 - Important when DNN processed in the cloud
 - Data centers have stringent cooling costs
 - Important in wearable devices
 - Form factor limited by cooling mechanisms

Power vs. Energy

- Energy consumption
 - Integral of the power
 - Important in mobile terminals
 - Battery has finite energy
 - Determines life-time of a device before recharging
 - Affects form factor/weights of the device

Energy – Quantitative Analysis

$$\frac{\text{inferences}}{\text{second}} \leq \max\left(\frac{\text{joules}}{\text{second}}\right) \times \frac{\text{inferences}}{\text{joule}}$$

Thermal Design Power (TPD)
Maximum power consumption

If inference per joule increases, inference per second can be increased

Energy – Quantitative Analysis

$$\frac{\text{inferences}}{\text{joule}} = \frac{\text{operations}}{\text{joule}} \times \frac{1}{\frac{\text{operations}}{\text{inference}}}$$

Depends on

- DNN hardware
- DNN model

Energy – Quantitative Analysis

- Energy per operation (in general)

$$\frac{\text{joules}}{\text{operation}} = \alpha \times C \times V_{dd}^2$$

- Two components

$$Energy_{total} = Energy_{data} + Energy_{MAC}$$

Dominant contribution

Energy – Quantitative Analysis

Operation	Energy (pJ)	Operation	Energy (pJ)
32 bit int ADD	0.1	32 bit float MULT	3.7
32 bit float ADD	0.9	32 bit SRAM memory	5.0
32 bit register file	1.0	32 bit NoC hop ^a	13.4
32 bit int MULT	3.1	32 bit DRAM memory	640.0

45 nm CMOS process

[M. Horowitz, “Energy table for 45nm process,” Stanford VLSI Wiki.]

Energy – Quantitative Analysis

$$Energy_{total} = Energy_{data} + Energy_{MAC}$$

- Reducing $Energy_{data}$
 - Improving memory system
 - Improving communication system
 - Improving dataflow technique

Energy – Quantitative Analysis

$$Energy_{total} = Energy_{data} + Energy_{MAC}$$

- Reducing $Energy_{MAC}$
 - Reducing precision
 - e.g., 8-bit add requires 10× less energy than 32-bit add
 - Energy vs. Accuracy trade-off
 - For instruction based systems (GP, ASIP, GPU)
 - Use large aggregate instructions (SIMD/Vector)
 - e.g., HMMA, IMMA Nvidia Volta V100 GPU

Summary

- Common layers in DNN
 - CONV, FC, Pooling, expedients
- Kernel computation
 - Layer computation led to matrix multiplication
- Induced traffic
 - Communication/Memory play a key role
- Quantitative analysis
 - MACs/Parameters not a good metric for performance/energy
 - Need of designing with HW in the loop