

# Convolutional Neural Networks (Forward)

---

11785 Introduction To Deep Learning (Spring 22) – Recitation 5

Aparajith Srinivasan

# Image Basics

- ❖ Images are nothing but a function of Intensity values  $f(x, y) = \text{Intensity at } (x, y)$  pixel



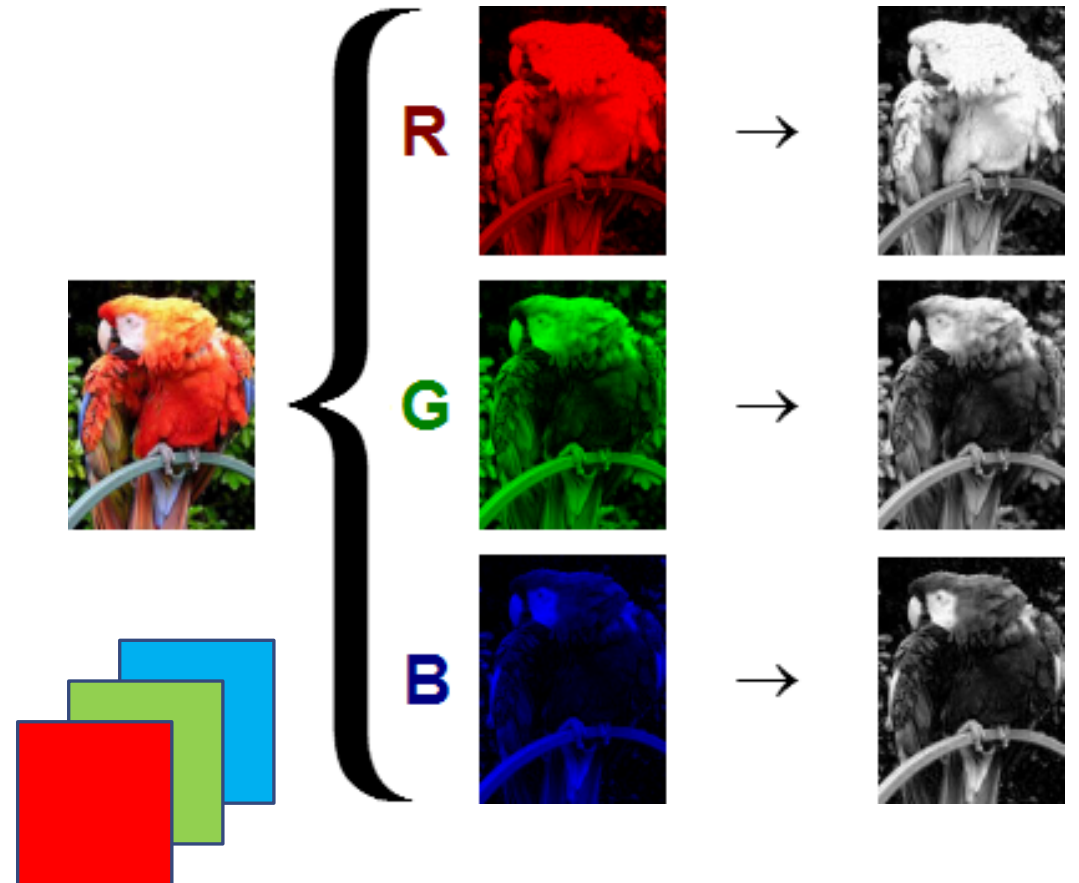
157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	105	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	106	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	96	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

- ❖ Single channel – single 2D Image (M x N) – grayscale

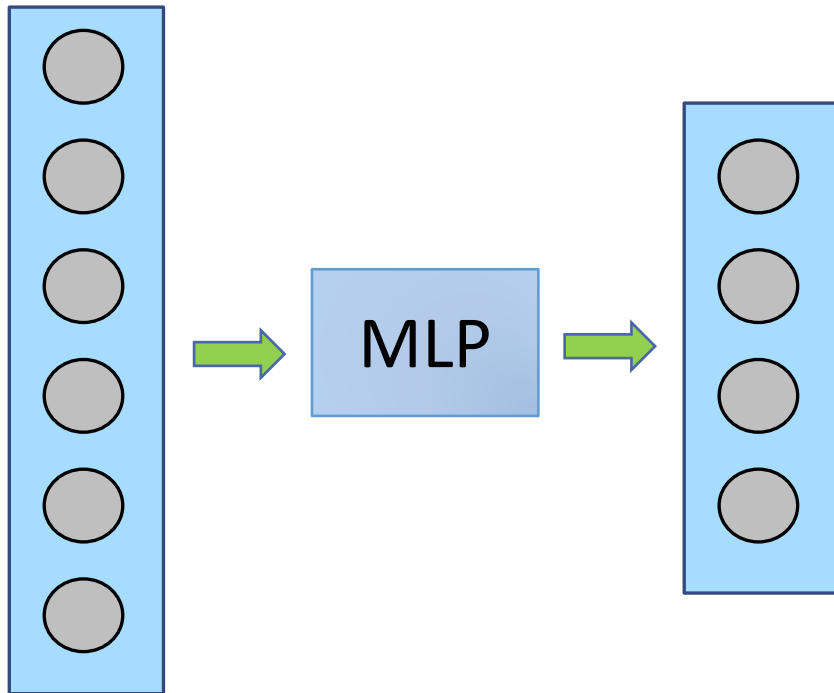
# Image Basics

- ❖ Multichannel images are formed by stacking many grayscale images together
- ❖ Color image:
  - ❑ 3 Channels RGB
  - ❑ Shape:  $(3 \times M \times N)$  or  $(M \times N \times 3)$

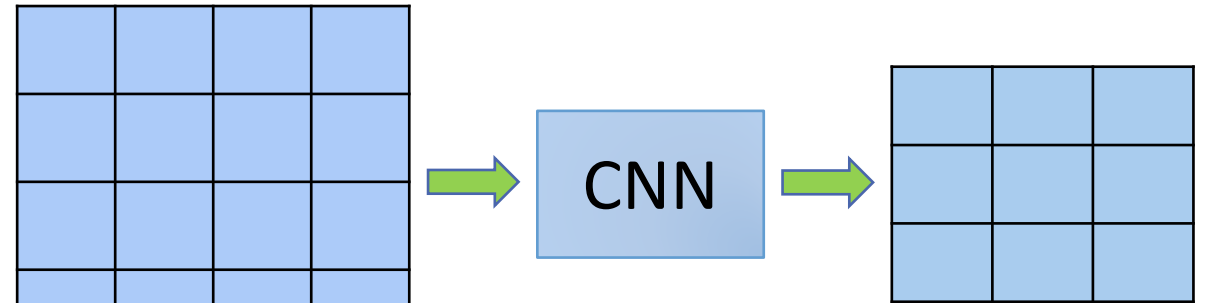


# MLP vs CNN

---



Feature vector to feature vector



Feature map to feature map

# CNN Components

---

- ❖ Any NN will have an input and weights (also bias)
- ❖ Consider the components for a single channel input

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

Input  
A

w11	w12
w21	w22

Kernel  
W

b1
----

Bias  
b

# CNN Components

---

- ❖ Any NN will have an input and weights (also bias)
- ❖ Consider the components for a single channel input

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

Input  
A

w11	w12
w21	w22

Kernel  
W

b1
----

Bias  
b

$$Z = A \otimes W + b$$

# CNN Components

---

- ❖ Any NN will have an input and weights (also bias)
- ❖ Consider the components for a single channel input

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

Input  
A

w11	w12
w21	w22

Kernel  
W

b1
----

Bias  
b

$$Z = A \otimes W + b$$

What is this  $\otimes$  ? 🤔

# CNN Steps

---

- ❖ Nothing but an element wise product and summation



# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

$$z'_{11} = x_{11} * w_{11} + x_{12} * w_{12} + x_{13} * w_{21} + x_{22} * w_{22}$$

z'11		

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

$$z'_{12} = x_{12} * w_{11} + x_{13} * w_{12} + x_{22} * w_{21} + x_{23} * w_{22}$$

z'11	z'12	

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

$$z'_{13} = x_{13} * w_{11} + x_{14} * w_{12} + x_{23} * w_{21} + x_{24} * w_{22}$$

z'11	z'12	z'13

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

z'11	z'12	z'13
z'21		

$$z'_{21} = x_{21} * w_{11} + x_{22} * w_{12} + x_{31} * w_{21} + x_{32} * w_{22}$$

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

$$z'_{22} = x_{22} * w_{11} + x_{23} * w_{12} + x_{32} * w_{21} + x_{33} * w_{22}$$

z'11	z'12	z'13
z'21	z'22	

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

z'11	z'12	z'13
z'21	z'22	z'23
z'31	z'32	z'33

$$z'_{33} = x_{33} * w_{11} + x_{34} * w_{12} + x_{43} * w_{21} + x_{44} * w_{22}$$

# CNN Steps

---

❖ Nothing but an element wise product and summation

x11	x12	x13	x14
x21	x22	x23	x24
x31	x32	x33	x34
x41	x42	x43	x44

w11	w12
w21	w22

z'11	z'12	z'13
z'21	z'22	z'23
z'31	z'32	z'33

$$z'_{33} = x_{33} * w_{11} + x_{34} * w_{12} + x_{43} * w_{21} + x_{44} * w_{22}$$

Kernel crosses 1 pixel at each step, stride = 1

# CNN Steps

---

❖ Nothing but an element wise product and summation

$z'_{11}$	$z'_{12}$	$z'_{13}$
$z'_{21}$	$z'_{22}$	$z'_{23}$
$z'_{31}$	$z'_{32}$	$z'_{33}$

$b_1$

$$z_{ij} = z'_{ij} + b_1$$

$z_{11}$	$z_{12}$	$z_{13}$
$z_{21}$	$z_{22}$	$z_{23}$
$z_{31}$	$z_{32}$	$z_{33}$

Final step is to add bias to all elements



# CNN Steps

---

❖ Note is that the output map size is decreased

# CNN Steps

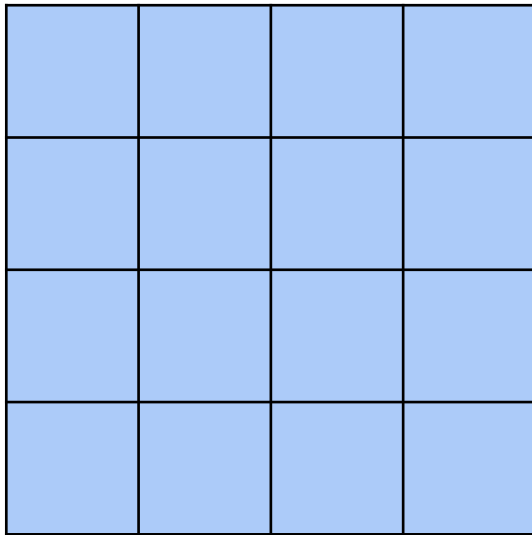
---

- ❖ Note is that the output map size is decreased
- ❖ Output size =  $(\text{Input\_size} - \text{Kernel\_size}) // \text{stride} + 1$

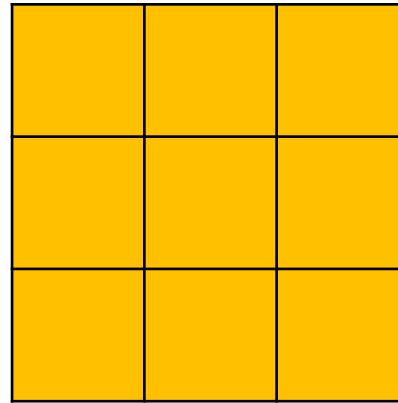
# CNN Steps

---

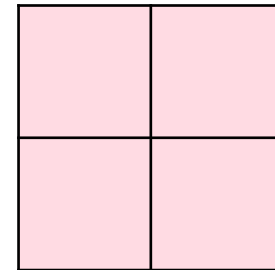
- ❖ Note is that the output map size is decreased
- ❖ Output size =  $(\text{Input\_size} - \text{Kernel\_size}) // \text{stride} + 1$



Input (4x4)



Kernel (3x3)



Output (2x2)  
(stride = 1)

# CNN Steps

---

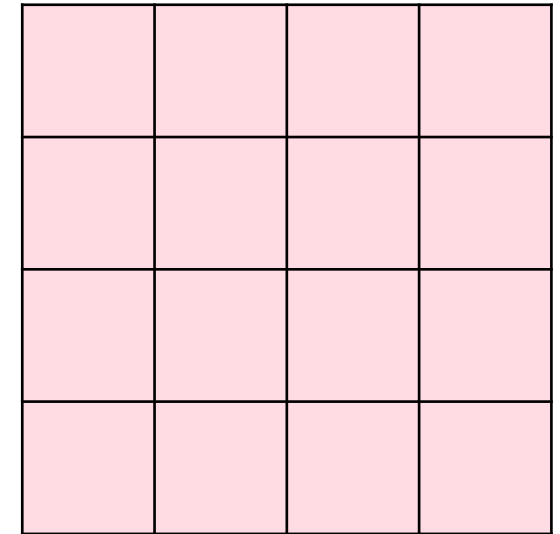
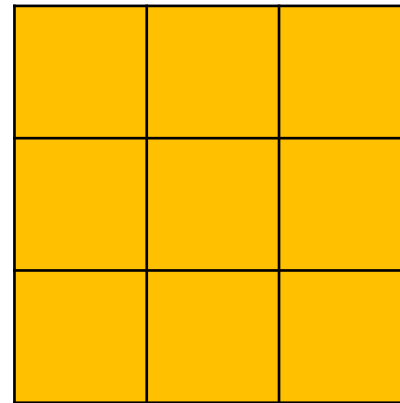
- ❖ Note is that the output map size is decreased
- ❖ Output size =  $(\text{Input\_size} - \text{Kernel\_size}) // \text{stride} + 1$
- ❖ Pad and convolve to conserve size

0	0	0	0	0	0
0					0
0					0
0					0
0					0
0	0	0	0	0	0

# CNN Steps

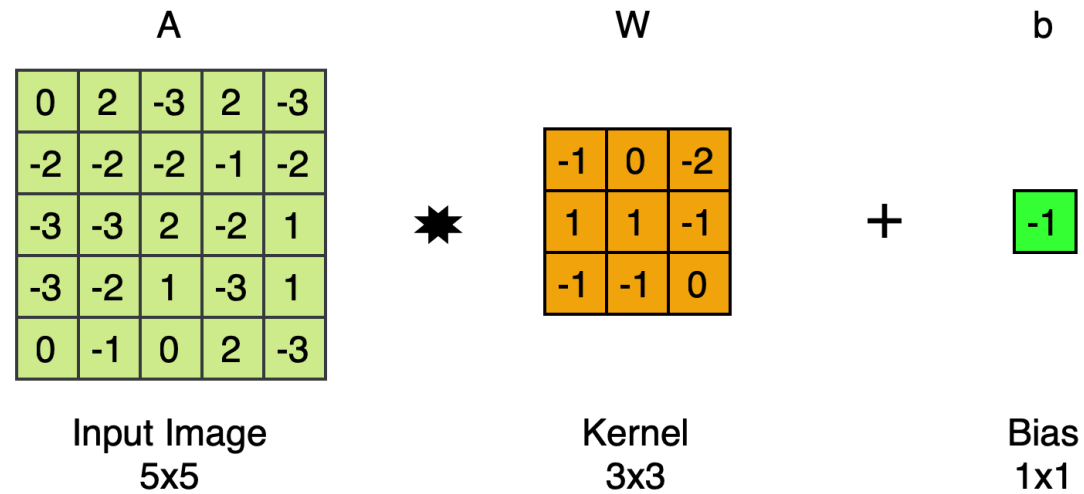
---

0	0	0	0	0	0
0					0
0					0
0					0
0					0
0	0	0	0	0	0



# CNN stride 2 Example

---



# CNN stride 2 Example

0	2	-3	2	-3
-2	-2	-2	-1	-2
-3	-3	2	-2	1
0	-1	0	2	-3
-3	-2	1	-3	1

10	

Step 1

$$\begin{array}{c} \text{A} \\ \begin{array}{ccccc} 0 & 2 & -3 & 2 & -3 \\ -2 & -2 & -2 & -1 & -2 \\ -3 & -3 & 2 & -2 & 1 \\ -3 & -2 & 1 & -3 & 1 \\ 0 & -1 & 0 & 2 & -3 \end{array} \\ \text{Input Image} \\ 5 \times 5 \end{array} \star \begin{array}{c} \text{W} \\ \begin{array}{ccc} -1 & 0 & -2 \\ 1 & 1 & -1 \\ -1 & -1 & 0 \end{array} \\ \text{Kernel} \\ 3 \times 3 \end{array} + \begin{array}{c} \text{b} \\ -1 \\ \text{Bias} \\ 1 \times 1 \end{array}$$

# CNN stride 2 Example

0	2	-3	2	-3
-2	-2	-2	-1	-2
-3	-3	2	-2	1
-3	-2	1	-3	1
0	-1	0	2	-3

10	8

Step 2

$$\begin{array}{c} \text{A} \\ \begin{array}{ccccc} 0 & 2 & -3 & 2 & -3 \\ -2 & -2 & -2 & -1 & -2 \\ -3 & -3 & 2 & -2 & 1 \\ -3 & -2 & 1 & -3 & 1 \\ 0 & -1 & 0 & 2 & -3 \end{array} \\ \text{Input Image} \\ 5 \times 5 \end{array} \star \begin{array}{c} \text{W} \\ \begin{array}{ccc} -1 & 0 & -2 \\ 1 & 1 & -1 \\ -1 & -1 & 0 \end{array} \\ \text{Kernel} \\ 3 \times 3 \end{array} + \begin{array}{c} \text{b} \\ \begin{array}{c} -1 \end{array} \\ \text{Bias} \\ 1 \times 1 \end{array}$$



# CNN stride 2 Example

0	2	-3	2	-3
-2	-2	-2	-1	-2
-3	-3	2	-2	1
-3	-2	1	-3	1
0	-1	0	2	-3

10	8
-6	

Step 3

$$\begin{array}{c} \text{A} \\ \begin{array}{|c|c|c|c|c|} \hline 0 & 2 & -3 & 2 & -3 \\ \hline -2 & -2 & -2 & -1 & -2 \\ \hline -3 & -3 & 2 & -2 & 1 \\ \hline -3 & -2 & 1 & -3 & 1 \\ \hline 0 & -1 & 0 & 2 & -3 \\ \hline \end{array} \\ \text{Input Image} \\ 5 \times 5 \end{array} \star \begin{array}{c} \text{W} \\ \begin{array}{|c|c|c|} \hline -1 & 0 & -2 \\ \hline 1 & 1 & -1 \\ \hline -1 & -1 & 0 \\ \hline \end{array} \\ \text{Kernel} \\ 3 \times 3 \end{array} + \begin{array}{c} \text{b} \\ \begin{array}{|c|} \hline -1 \\ \hline \end{array} \\ \text{Bias} \\ 1 \times 1 \end{array}$$

# CNN stride 2 Example

0	2	-3	2	-3
-2	-2	-2	-1	-2
-3	-3	2	-2	1
-3	-2	1	-3	1
0	-1	0	2	-3

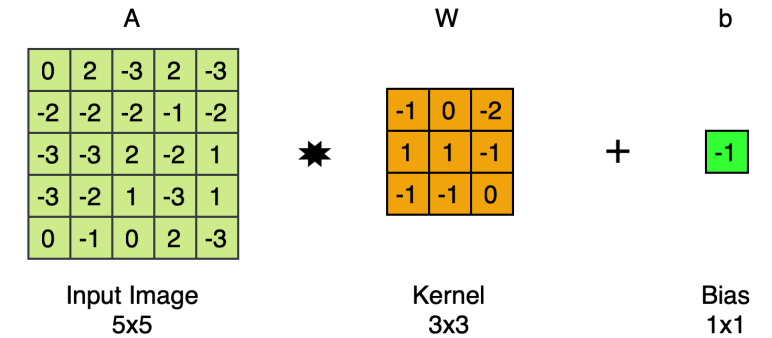
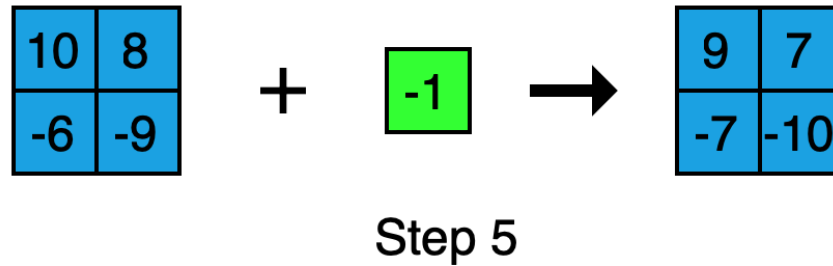
10	8
-6	-9

Step 4

A		W		b																																			
<table><tr><td>0</td><td>2</td><td>-3</td><td>2</td><td>-3</td></tr><tr><td>-2</td><td>-2</td><td>-2</td><td>-1</td><td>-2</td></tr><tr><td>-3</td><td>-3</td><td>2</td><td>-2</td><td>1</td></tr><tr><td>-3</td><td>-2</td><td>1</td><td>-3</td><td>1</td></tr><tr><td>0</td><td>-1</td><td>0</td><td>2</td><td>-3</td></tr></table>	0	2	-3	2	-3	-2	-2	-2	-1	-2	-3	-3	2	-2	1	-3	-2	1	-3	1	0	-1	0	2	-3		<table><tr><td>-1</td><td>0</td><td>-2</td></tr><tr><td>1</td><td>1</td><td>-1</td></tr><tr><td>-1</td><td>-1</td><td>0</td></tr></table>	-1	0	-2	1	1	-1	-1	-1	0		<table><tr><td>-1</td></tr></table>	-1
0	2	-3	2	-3																																			
-2	-2	-2	-1	-2																																			
-3	-3	2	-2	1																																			
-3	-2	1	-3	1																																			
0	-1	0	2	-3																																			
-1	0	-2																																					
1	1	-1																																					
-1	-1	0																																					
-1																																							
Input Image 5x5	$\star$	Kernel 3x3	$+$	Bias 1x1																																			

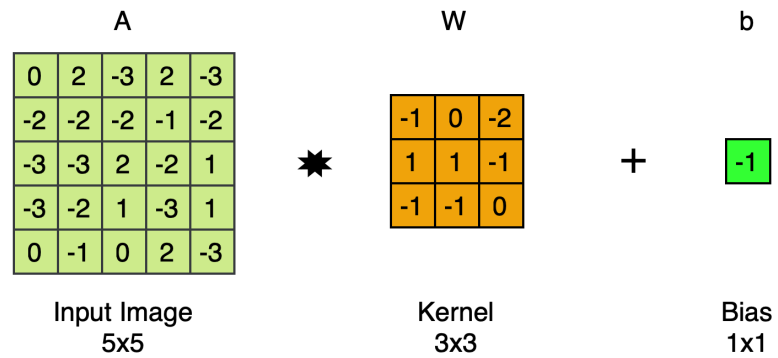
# CNN stride 2 Example

---



# Interpreting stride > 1

---



# Interpreting stride > 1

---

$$\begin{array}{c} \text{A} \\ \begin{array}{|c|c|c|c|c|} \hline 0 & 2 & -3 & 2 & -3 \\ \hline -2 & -2 & -2 & -1 & -2 \\ \hline -3 & -3 & 2 & -2 & 1 \\ \hline -3 & -2 & 1 & -3 & 1 \\ \hline 0 & -1 & 0 & 2 & -3 \\ \hline \end{array} \\ \text{Input Image} \\ 5 \times 5 \end{array} \star \begin{array}{c} \text{W} \\ \begin{array}{|c|c|c|} \hline -1 & 0 & -2 \\ \hline 1 & 1 & -1 \\ \hline -1 & -1 & 0 \\ \hline \end{array} \\ \text{Kernel} \\ 3 \times 3 \end{array} + \begin{array}{c} \text{b} \\ \begin{array}{|c|} \hline -1 \\ \hline \end{array} \\ \text{Bias} \\ 1 \times 1 \end{array}$$

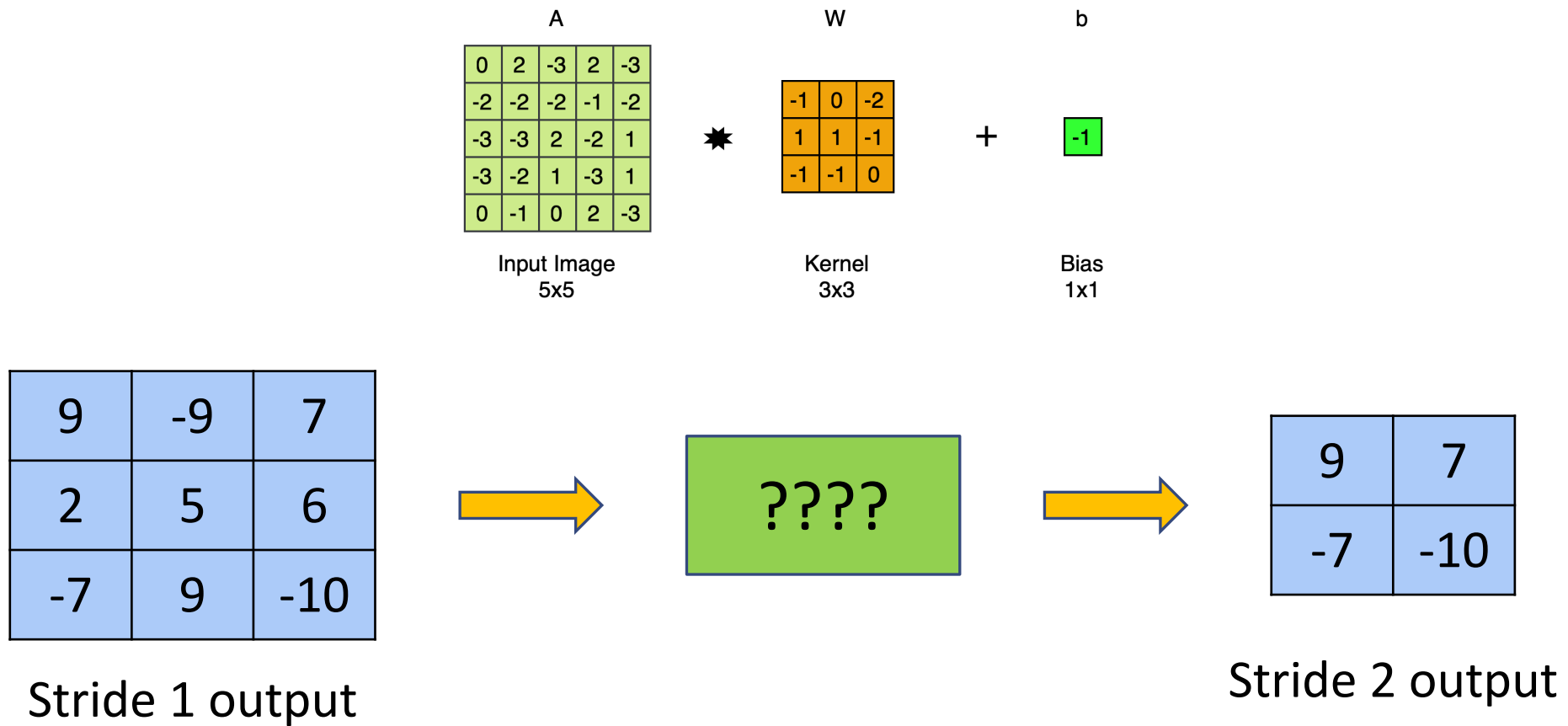
9	-9	7
2	5	6
-7	9	-10

Stride 1 output

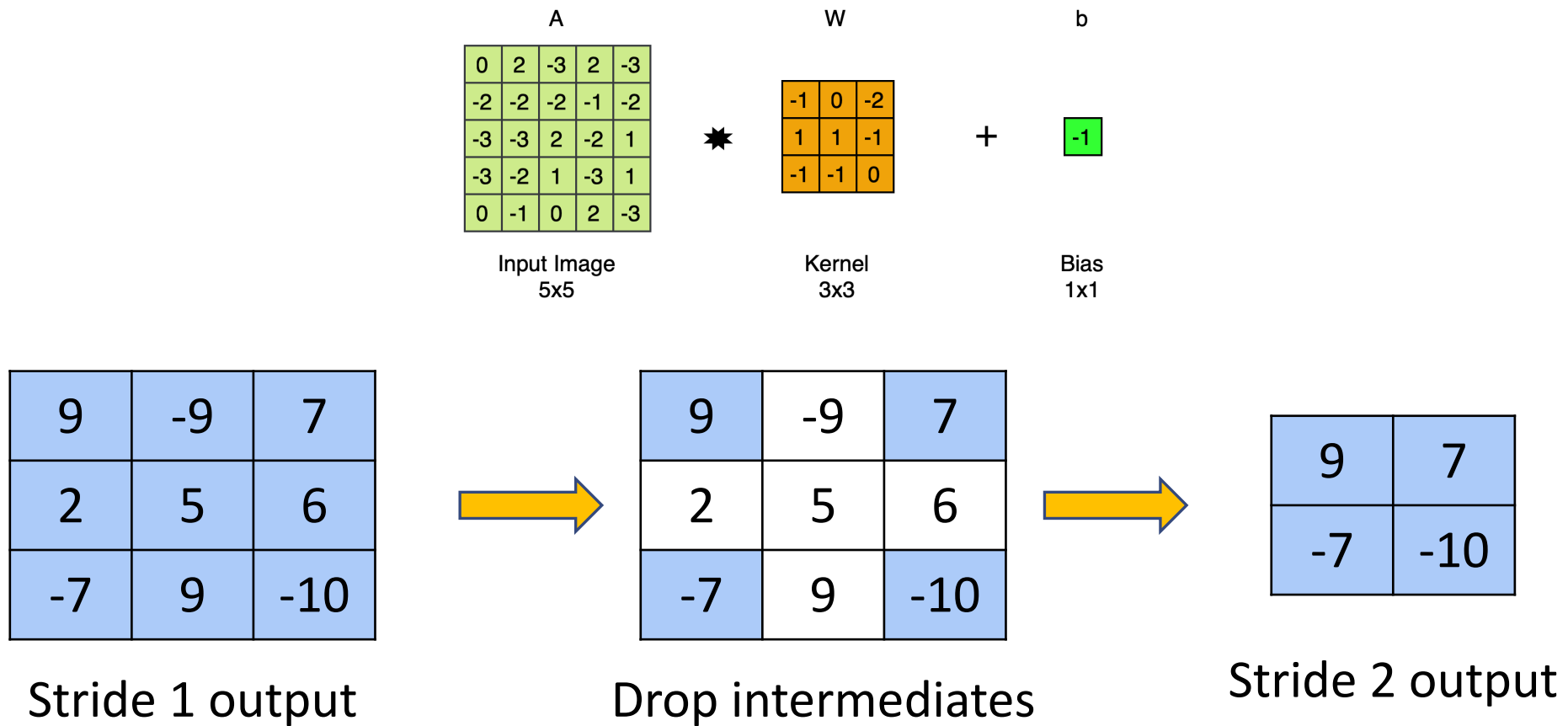
9	7
-7	-10

Stride 2 output

# Interpreting stride > 1



# Interpreting stride > 1



# Multichannel CNN

---

- ❖ When input has multiple channels, kernel has the same number of channels
- ❖ Each channel of the kernel convolves with the corresponding input channel to produce an output channel – all output maps obtained from this convolution are added together
- ❖ 1 filter produces 1 output channel. N filters produce N output channel



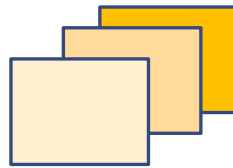
# Multichannel CNN

---

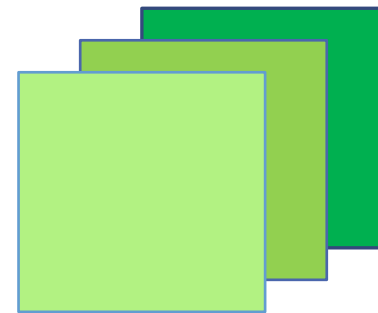
- ❖ When input has multiple channels, kernel has the same number of channels
- ❖ Each channel of the kernel convolves with the corresponding input channel to produce an output channel – all output maps obtained from this convolution are added together
- ❖ 1 filter produces 1 output channel. N filters produce N output channel



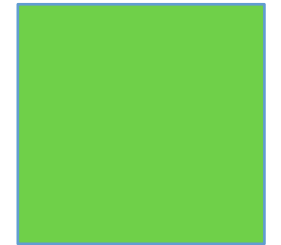
Input



Kernel

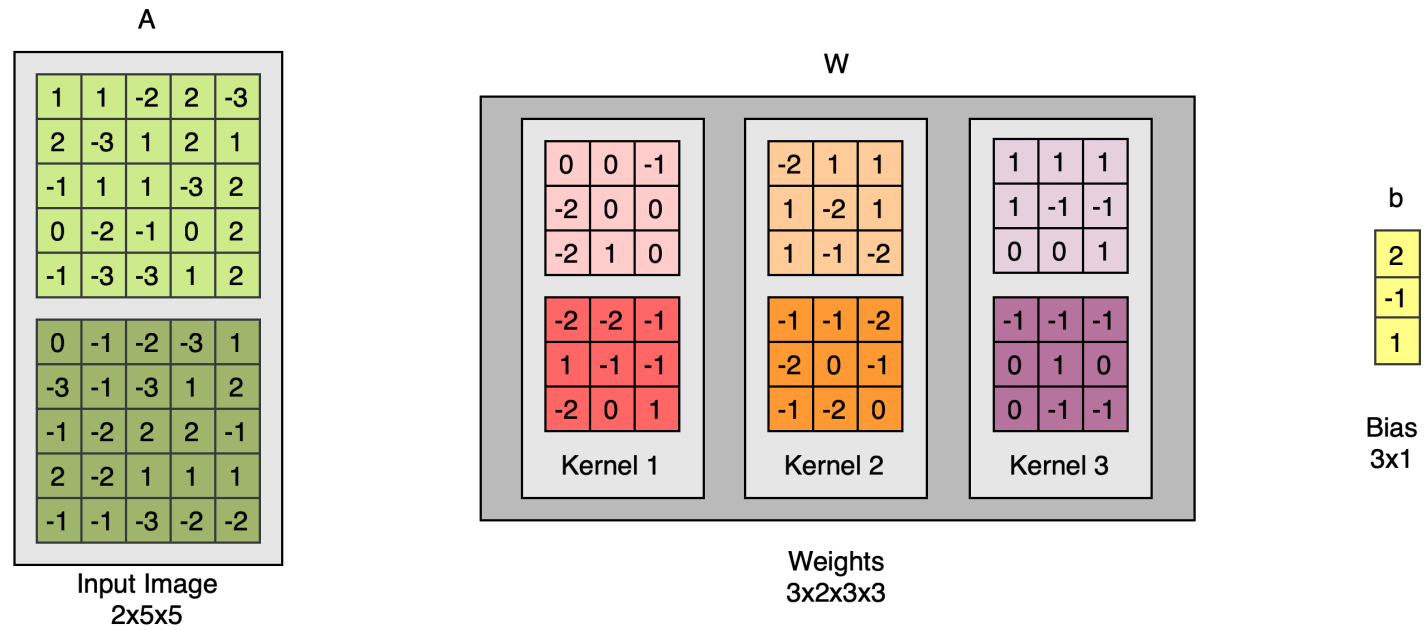


3 maps



Add all maps

# Multichannel CNN



# Multichannel CNN

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

10	

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

10	-6

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

10	-6
4	

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

10	-6
4	3

2
---

12	-4
6	5

A

1	1	-2	2	-3
2	-3	1	2	1
-1	1	1	-3	2
0	-2	-1	0	2
-1	-3	-3	1	2
0	-1	-2	-3	1
-3	-1	-3	1	2
-1	-2	2	2	-1
2	-2	1	1	1
-1	-1	-3	-2	-2

Input Image  
2x5x5

W

0	0	-1
-2	0	0
-2	1	0

Kernel 1

-2	1	1
1	-2	1
1	-1	-2

Kernel 2

1	1	1
1	-1	-1
0	0	1

Kernel 3

Weights  
3x2x3x3

b

2
-1
1

Bias  
3x1

# Multichannel CNN

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

21	

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	2	
-1	-3	-3	1	2		

21	2

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	2	
-1	-3	-3	1	2		

21	2
12	

		0	-1	-2	-3	1
1	1	-2	2	-3	2	
2	-3	1	2	1	-1	
-1	1	1	-3	2	1	
0	-2	-1	0	2	-2	
-1	-3	-3	1	2		

21	2
12	-8

-1

20	1
11	-9

A

1	1	-2	2	-3
2	-3	1	2	1
-1	1	1	-3	2
0	-2	-1	0	2
-1	-3	-3	1	2
0	-1	-2	-3	1
-3	-1	-3	1	2
-1	-2	2	2	-1
2	-2	1	1	1
-1	-1	-3	-2	-2

Input Image  
2x5x5

W

0	0	-1
-2	0	0
-2	1	0
-2	-2	-1
1	-1	-1
-2	0	1

Kernel 1

-2	1	1
1	-2	1
1	-1	-2
-1	-1	-2
-2	0	-1
-1	-2	0

Kernel 2

1	1	1
1	-1	-1
0	0	1
-1	-1	-1
0	1	0
0	-1	-1

Kernel 3

Weights  
3x2x3x3

b  
2  
-1  
1  
Bias  
3x1

# Multichannel CNN

		0	-1	-2	-3	1	
1	1	-2	2	-3	2		
2	-3	1	2	1	-1		
-1	1	1	-3	2	1		
0	-2	-1	0	2	-2		
-1	-3	-3	1	2			

7	

		0	-1	-2	-3	1	
1	1	-2	2	-3	2		
2	-3	1	2	1	-1		
-1	1	1	-3	2	1		
0	-2	-1	0	2	-2		
-1	-3	-3	1	2			

7	1

		0	-1	-2	-3	1	
1	1	-2	2	-3	2		
2	-3	1	2	1	-1		
-1	1	1	-3	2	1		
0	-2	-1	0	2	-2		
-1	-3	-3	1	2			

7	1
4	

		0	-1	-2	-3	1	
1	1	-2	2	-3	2		
2	-3	1	2	1	-1		
-1	1	1	-3	2	1		
0	-2	-1	0	2	-2		
-1	-3	-3	1	2			

7	1
4	1

1
---

8	2
5	2

A

1	1	-2	2	-3
2	-3	1	2	1
-1	1	1	-3	2
0	-2	-1	0	2
-1	-3	-3	1	2
0	-1	-2	-3	1
-3	-1	-3	1	2
-1	-2	2	2	-1
2	-2	1	1	1
-1	-1	-3	-2	-2

Input Image  
2x5x5

W

0	0	-1
-2	0	0
-2	1	0

Kernel 1

-2	1	1
1	-2	1
1	-1	-2

Kernel 2

1	1	1
1	-1	-1
0	0	1

Kernel 3

Weights  
3x2x3x3

b

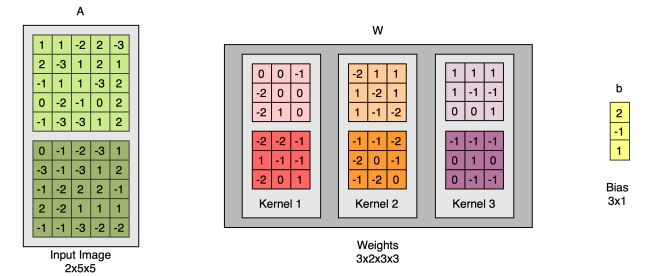
2
-1
1

Bias  
3x1

# Multichannel CNN

			8	2
		20	1	2
12	-4	-9		
6	5			

Output Response  
3x2x2



# Pooling

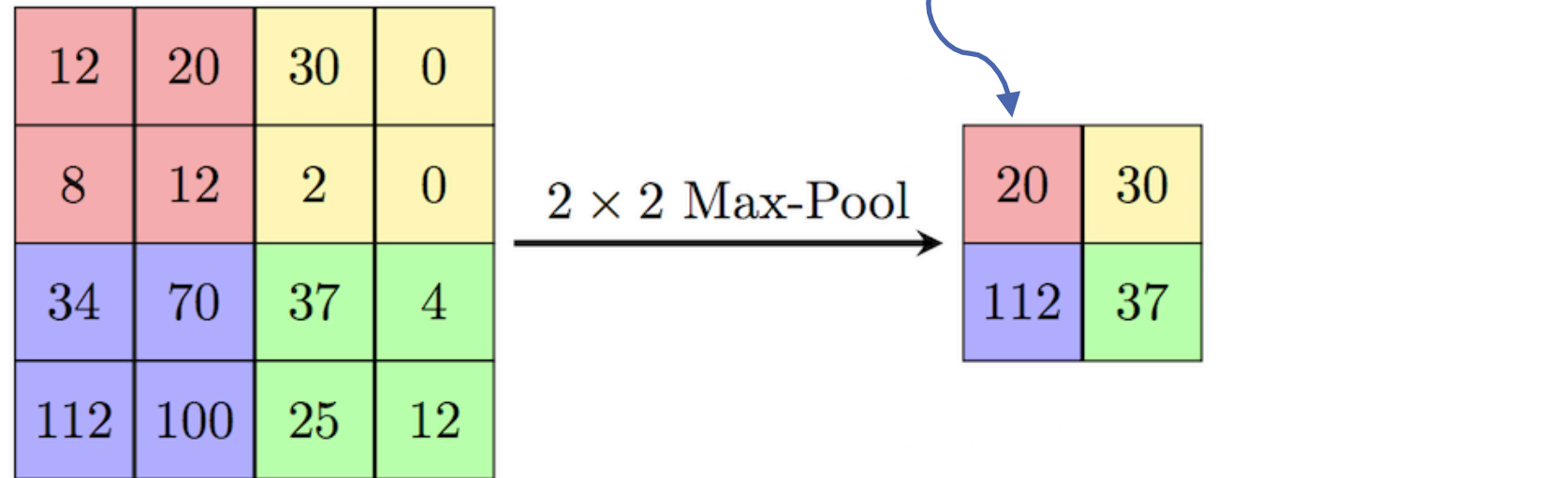
---

- ❖ Step followed by convolution
- ❖ For jitter invariance (also downsamples if stride  $> 1$ )

# Pooling

## ❖ Max Pooling

➤ Kernel\_Size = 2x2; stride = 2





# Pooling

## ❖ Mean Pooling

➤ Kernel\_Size = 2x2; stride = 2

