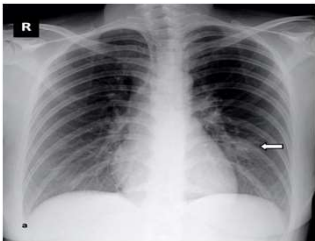
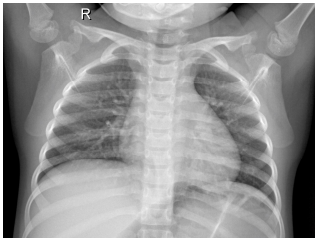


Convolutional Neural Networks

Traditional approach

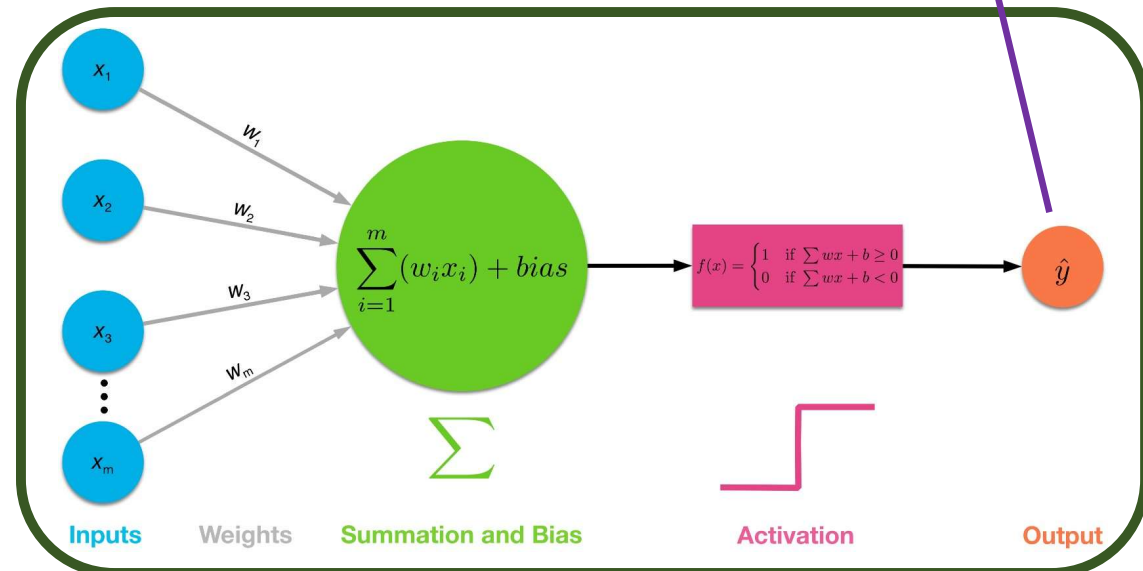


- Extract features, feed to the neural network

- Opacity Patterns
- Lesion Localization
-
- Pleural Effusion

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} = \begin{bmatrix} 2 \\ 1.5 \\ \dots \\ 8 \end{bmatrix}$$

0 – Healthy
1 – Corona



Wouldn't it be nice if we could feed the image directly ?

... and let the “learning process” figure out which features to extract ?

What is an Image?

- Data in the form of matrix (Rows and Columns) consisting of Pixels



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	156	24	0	0	6	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5
0	0	23	113	215	255	250	248	255	255	248	118	14	12	0	0
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	14	1	0	6	6	0	0	0

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	4	97	255	255	255	248	252	255	244	255	182	10	0	4
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	156	24	0	0	6	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5
0	0	23	113	215	255	250	248	255	255	248	118	14	12	0	0
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	14	1	0	6	6	0	0	0

<https://mozanunal.com/images/pixel.png>



(a)

Color Image



(b)

Grayscale Image

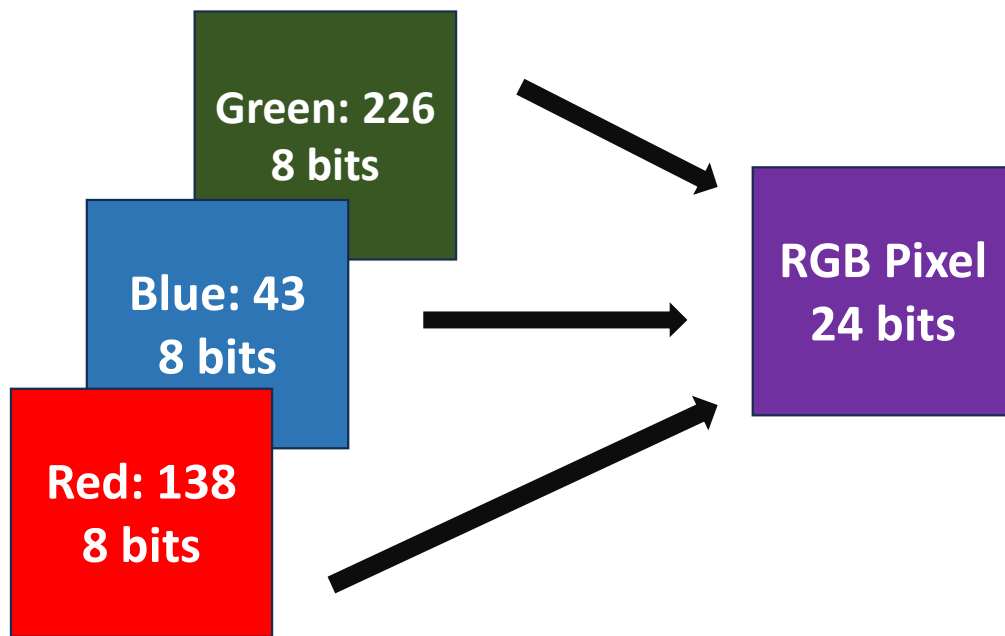


(c)

Binary Image

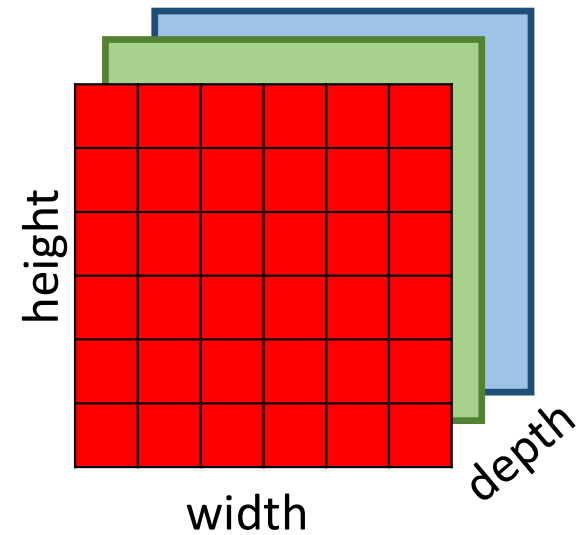
<https://www.researchgate.net/profile/Sanskriti-Patel-2>

Color Image



Width \times Height \times Depth

Depth: [Red, Green, Blue]





28 × 28 pixels

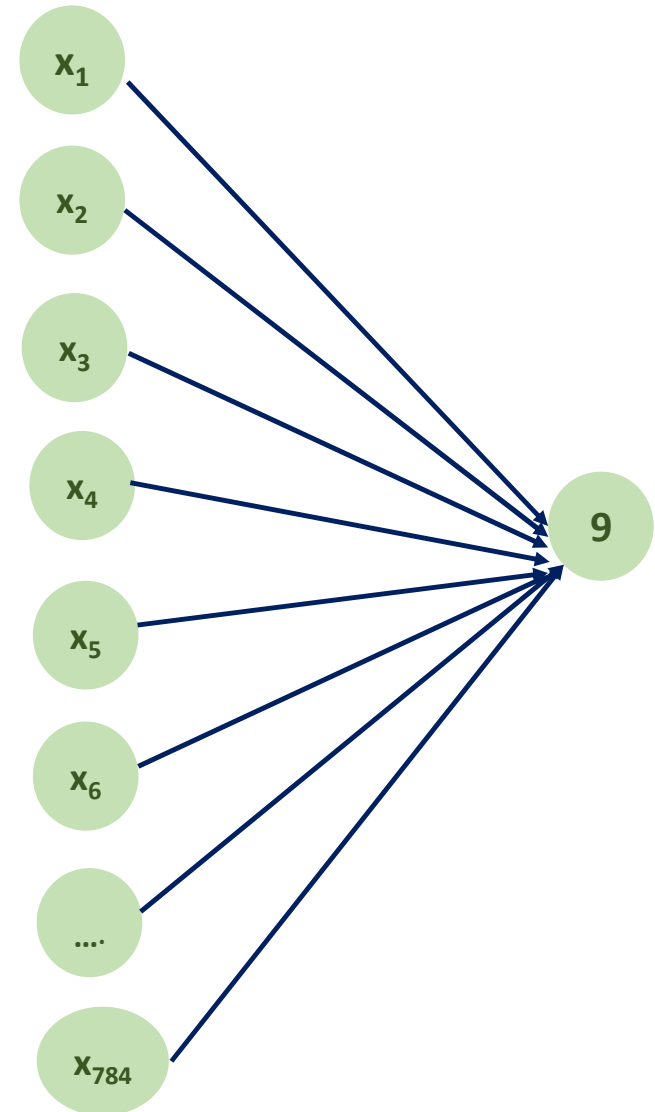


28 × 28 matrix

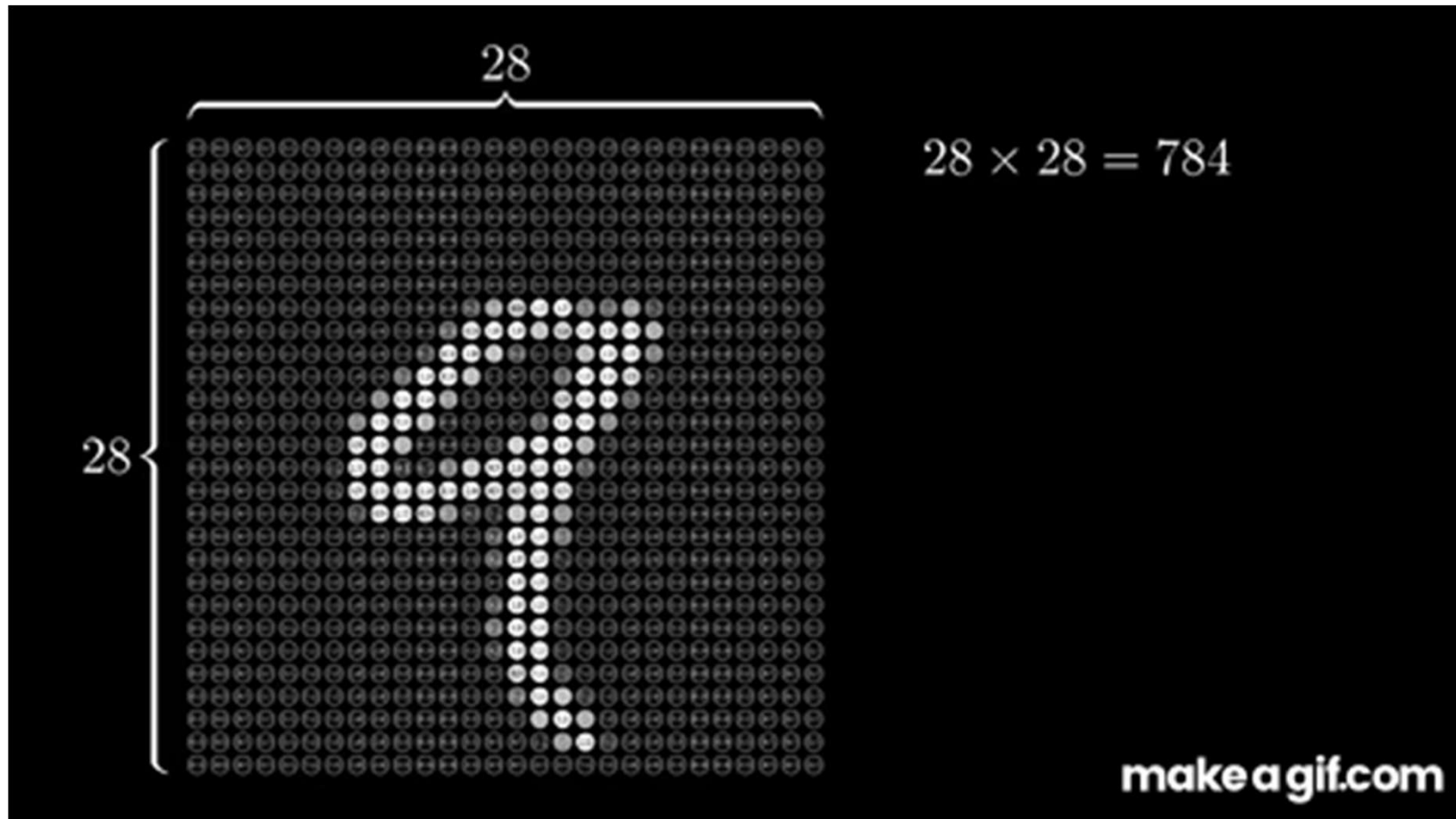
0	0	0	0	0	0	-	0
0	87	240	210	24	45	-	0
0	47	0	28	183	195	-	0
0	35	0	99	167	210	-	0
0	12	32	76	101	201	-	0
0	109	145	230	254	123	-	0
-	-	-	-	-	-	-	-
0	0	0	0	0	0	-	0

784 × 1 matrix

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 0 \\ 0 \\ 87 \\ 240 \\ 210 \\ 24 \\ 45 \\ \dots \\ 0 \\ 0 \\ 0 \\ 47 \\ 0 \\ 28 \\ \dots \\ \dots \\ 0 \end{pmatrix}$$

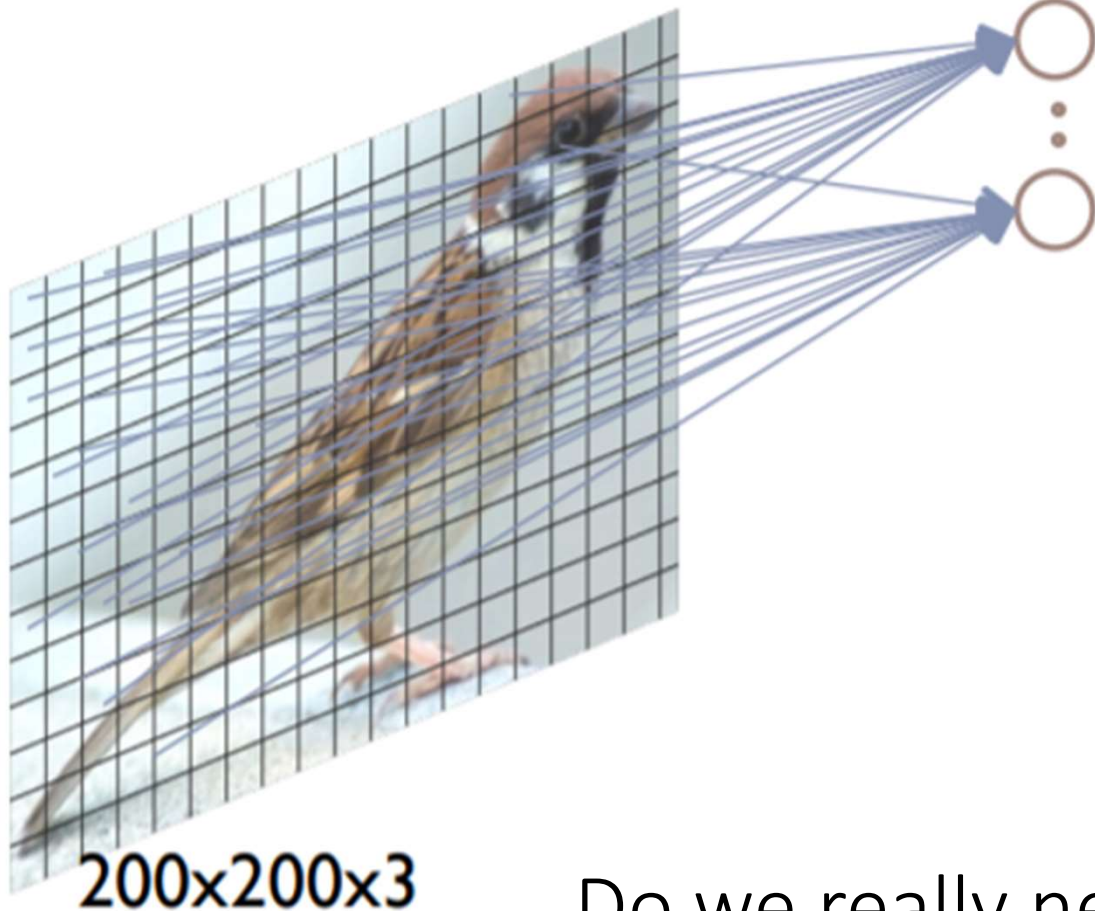


Flatten image, Feed to NN



<https://www.youtube.com/watch?v=aircAruvnKk>

For colour images?



- Input Layer: 120000
- # Hidden Units: 100000
- # Params to train: 12 billion
- Need huge training data to prevent overfitting
- Will not perform well, even if the image is shifted by one pixel
- No correlation exists among pixels

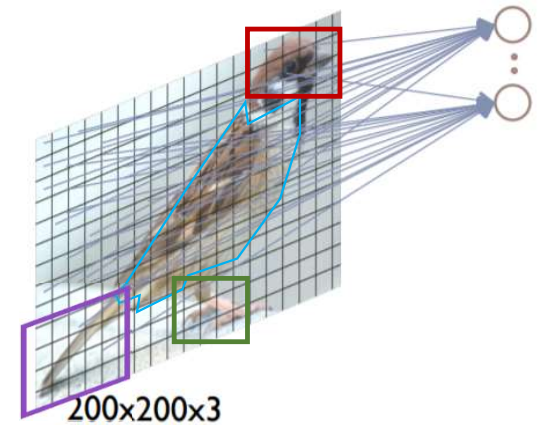
Do we really need full connectivity ?

Solutions?

- Huge number of parameters to train
 - Reduce the number of parameters, # of input nodes
- Need huge training data to prevent overfitting
- Will not perform well, even if the image is shifted by one pixel
 - Tolerate small shifts in where the pixels are in the image
- No correlation exists among pixels
 - Take advantage of the correlations that we observe in complex images

Hierarchical Combination

- Images are 2D.
- Assumption: Object image = combination of **2D image patterns**

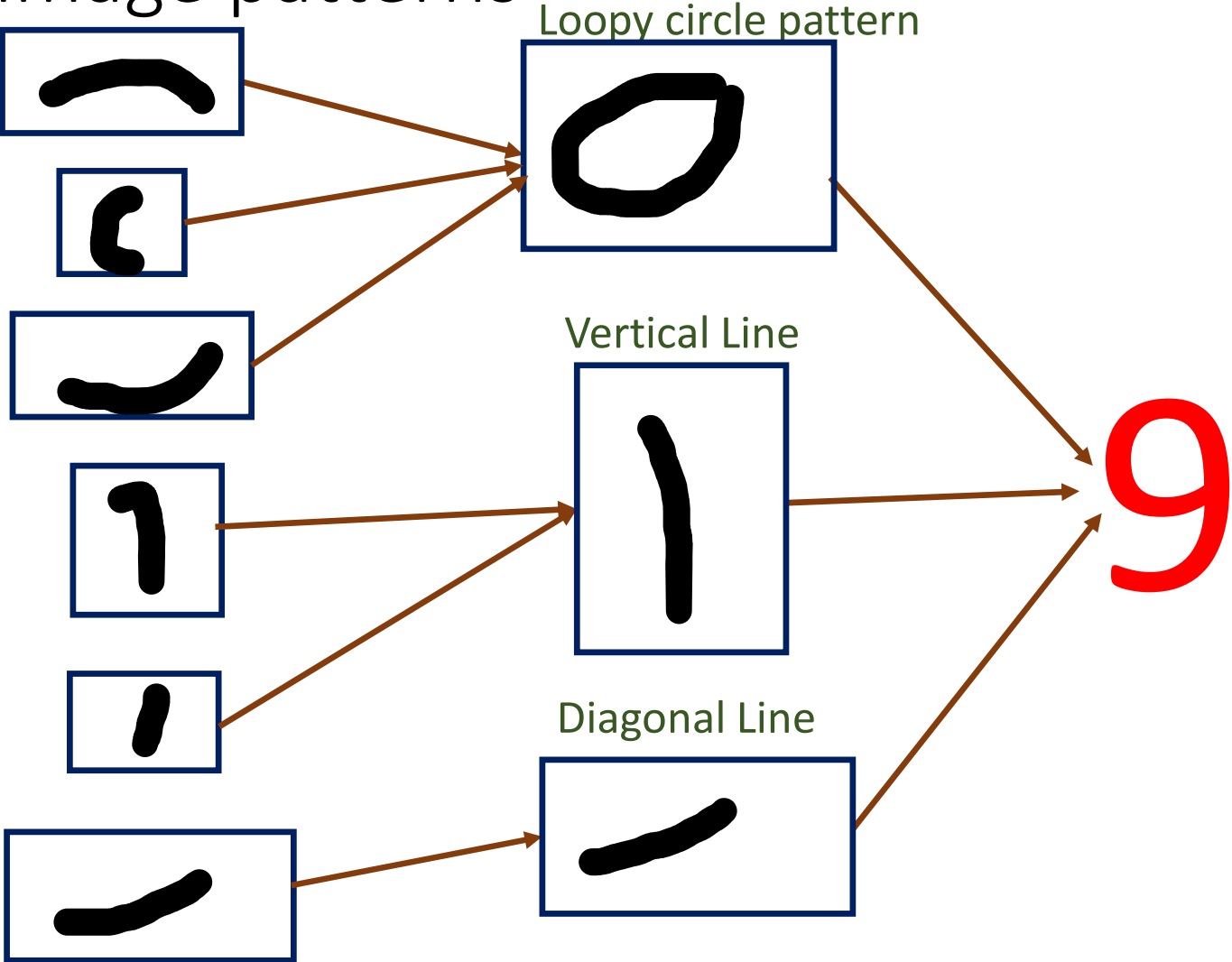
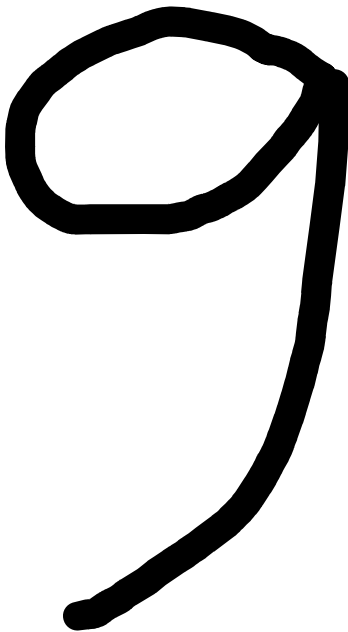


hierarchical combination of **2D image patterns**

Goal:

1. Determine 2D image patterns [2D image patterns are smaller than image]
2. Map small image patterns → Large image patterns
3. Map larger image patterns → Target

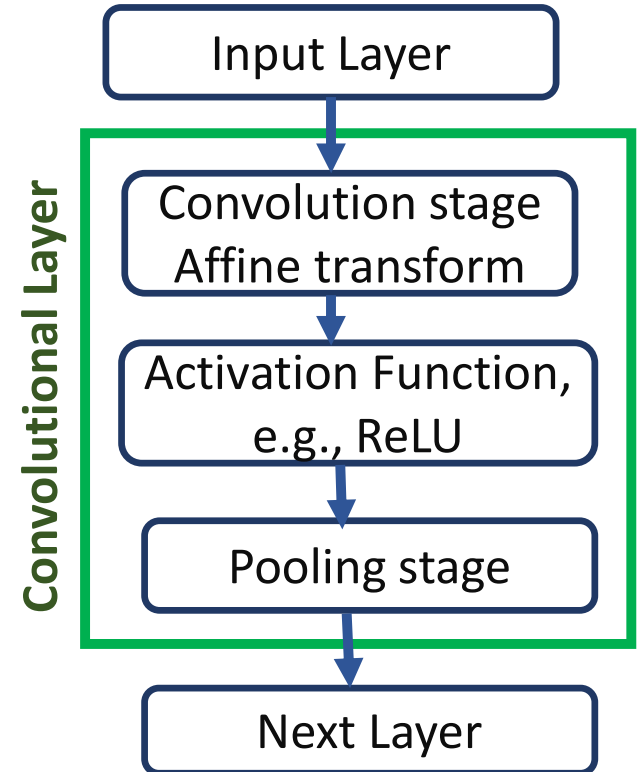
Determining 2D image patterns



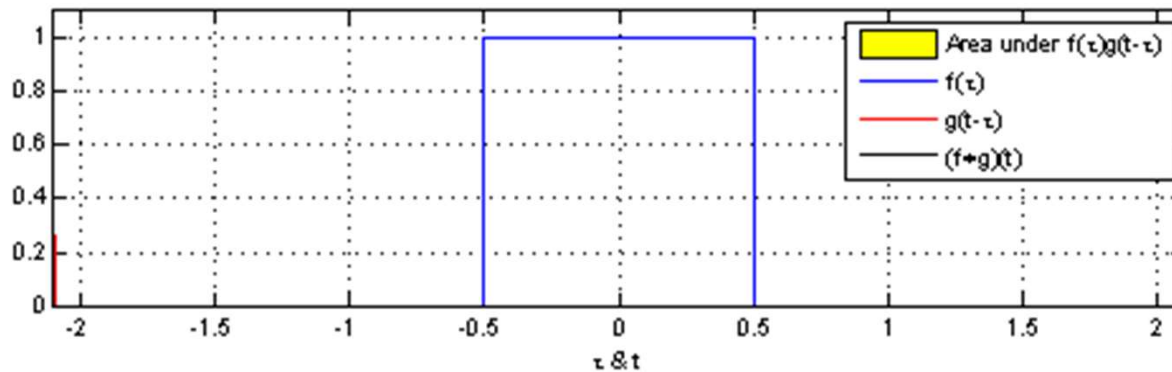
Convolutional Neural Networks (CNN)

- ✓ CNN are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers
- Convolution is a mathematical operation on two functions (f and g), that produces a third function ($f * g$) that expresses how the shape of one is modified by the other

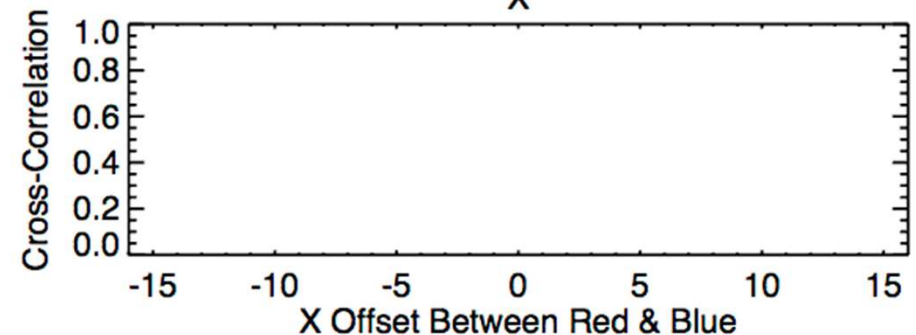
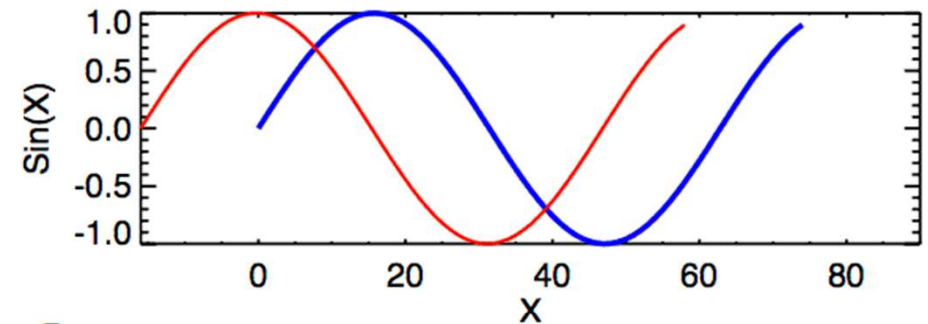
$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau$$



Convolution



- Cross Correlation: It is a measure of similarity of two series as a function of the displacement of one relative to another.



Convolutional Stage

Step 1: Apply a **Filter** to the input Image

- A filter is just a smaller square that is commonly N_{odd} pixels by N_{odd} pixels

-1	1	1	1	-1
-1	1	-1	1	-1
-1	1	1	1	-1
-1	-1	-1	1	-1
-1	-1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	-1	-1

3 pixels by 3 pixels

1/9	1/9	1/9
1/9	-1/9	1/9
1/9	1/9	1/9

- ✓ Start with random values.
- ✓ The intensity of each pixel in the filter/kernel determined(learned) by Backpropagation.

To apply the Filter to the input image, we overlay the Filter onto the image

Convolutions

-1	1	1	1	1	-1
-1	1	1	-1	1	-1
-1	1	1	1	1	-1
-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1

$\ast \frac{1}{9}$

1	1	1
1	-1	1
1	1	1

=

-0.11		

$$(-1 \ast 1 + 1 \ast 1 + 1 \ast 1 + (-1) \ast 1 + 1 \ast (-1) + (-1) \ast 1 + (-1) \ast 1 + 1 \ast 1 + 1 \ast 1) / 9 = -1/9 = -0.11$$

Convolutions

-1	1 1	1 1	1 1	-1
-1	1 1	-1 -1	1 1	-1
-1	1 1	1 1	1 1	-1
-1	-1	-1	1	-1
-1	-1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	-1	-1

$\ast \frac{1}{9}$

1	1	1
1	-1	1
1	1	1

=

-0.11	1	

-1	1	1	1	-1
-1	1	-1	1	-1
-1	1	1	1	-1
-1	-1	-1	1	-1
-1	-1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	-1	-1



1/9	1/9	1/9
1/9	-1/9	1/9
1/9	1/9	1/9

=

-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33

Feature Map

- ✓ Filters are nothing but the feature detectors. A 2-D map is created where certain features appear in the input.
- ✓ Filters are location invariant. They can detect the loopy patterns in any location of the image.

Filters

- Different values of the filter matrix will produce different Feature Maps for the same input image — **Multiple filters, multiple 'feature maps'**



- The Convolution operation captures the local dependencies in the original image

Padding: Valid and Same

- The size of the output is smaller than the input. To maintain the dimension of output as in input, we use padding. Padding is a process of adding zeros to the input matrix symmetrically

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

*

$1/9$	$1/9$	$1/9$
$1/9$	$-1/9$	$1/9$
$1/9$	$1/9$	$1/9$

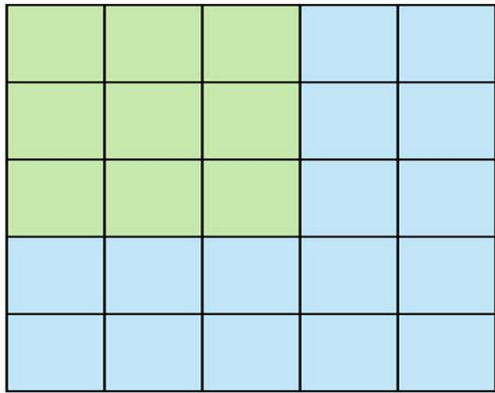
=

0.22	-0.22	0.22	-0.22	0.22
0.22	-0.11	1	-0.11	0.22
0	-0.55	0.11	-0.33	0.22
-0.22	-0.33	0.33	-0.33	0
-0.44	-0.22	-0.11	-0.22	-0.22
-0.22	-0.33	-0.33	-0.33	-0.22
0	0	-0.22	0	-0.22

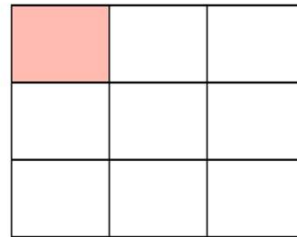
- Valid: no padding
- Same: Pad so that the output size is same as the input size

Stride

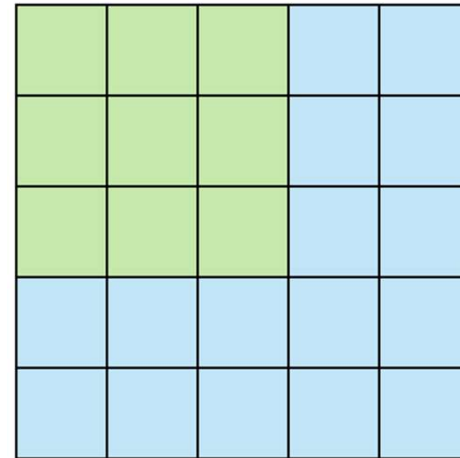
- Stride denotes how many steps we are moving in each step during convolution.



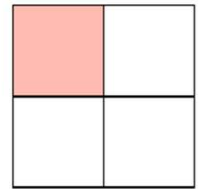
Stride 1



Feature Map



Stride 2

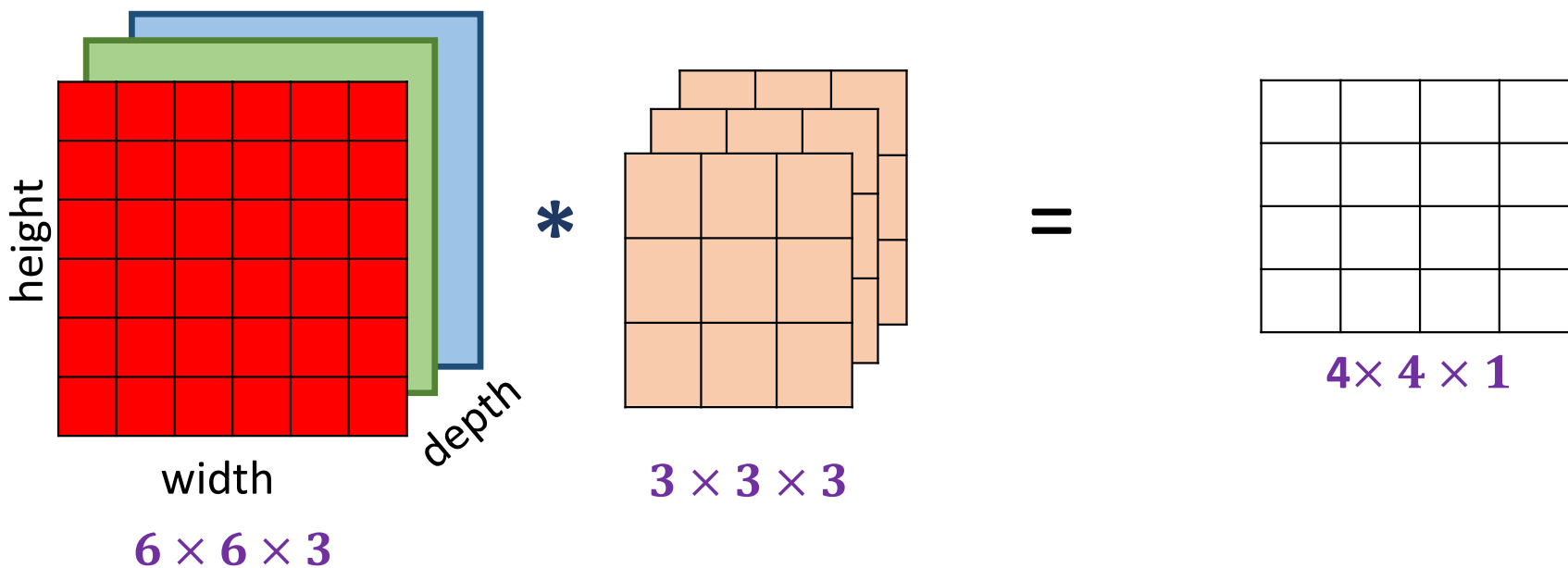


Feature Map

Solutions Obtained

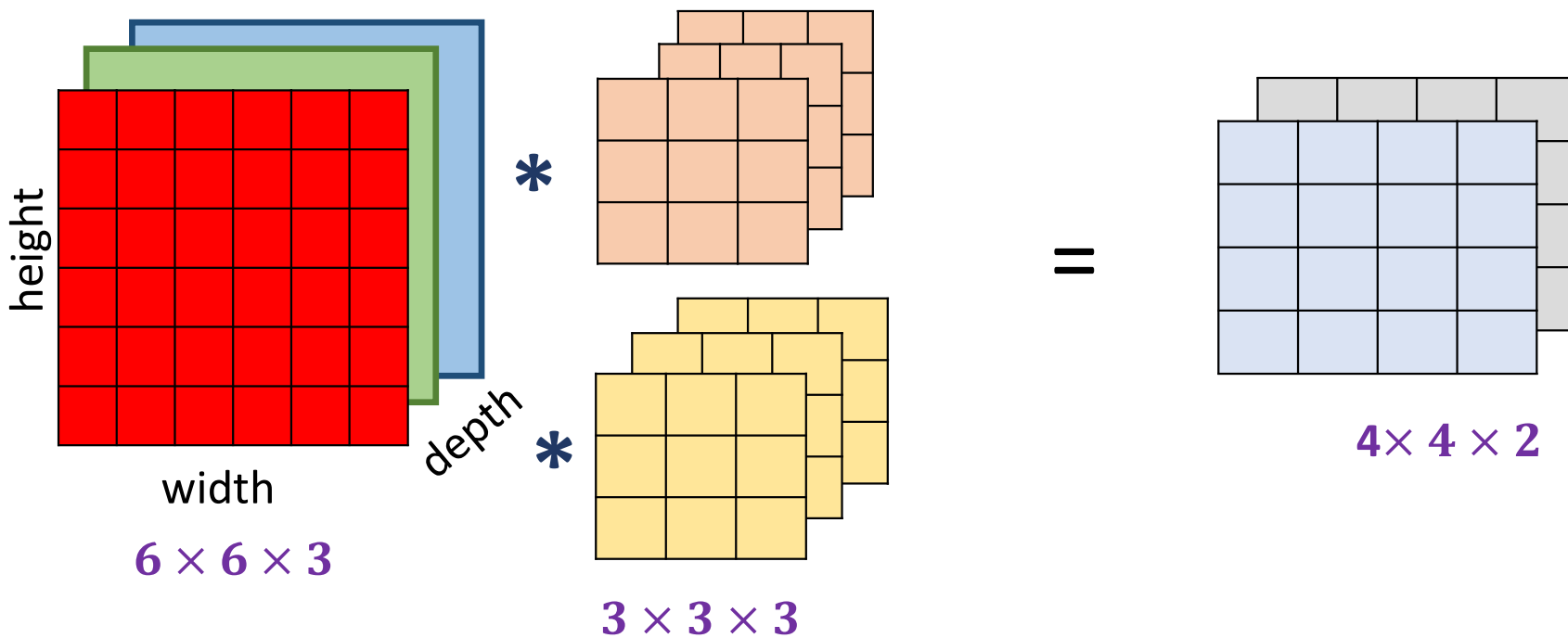
- In feed forward neural network, every output unit interacts with every input unit.
- CNN, typically have sparse connectivity (sparse weights). Each member of the kernel is used at every position of the input.
- This is accomplished by making the kernel smaller than the input.
- Reduce the number of parameters, where instead of learning a separate set of parameters for every location, we learn only one set – Parameter Sharing.
- Hyperparameters: Filter size, and number of filters.

Convolutions on RGB images



If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

Convolutions on RGB images



Summary of convolutions

f^l = filter size

s^l = stride

p^l = padding

$n_c^{[l]}$ = number of filters

Layer l is a convolution layer

- Input: $n_H^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]} : 28 \times 28 \times 3$
- Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} : 5 \times 5 \times 3$
- Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]} : 5 \times 5 \times 3 \times 10$
- Output Size: $n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$
 - : $24 \times 24 \times 10$ [with valid padding]
 - : $28 \times 28 \times 10$ [with same padding]

$$n_{H/w}^{[l]} = \left\lceil \frac{n_{H/w}^{[l-1]} + 2p^l - f^l}{s} + 1 \right\rceil$$
$$= \left\lceil \frac{28 + 2 \times 0 - 5}{1} + 1 \right\rceil = 24$$

Convolutional Neural Networks (CNN)

❑ First Stage: Convolution

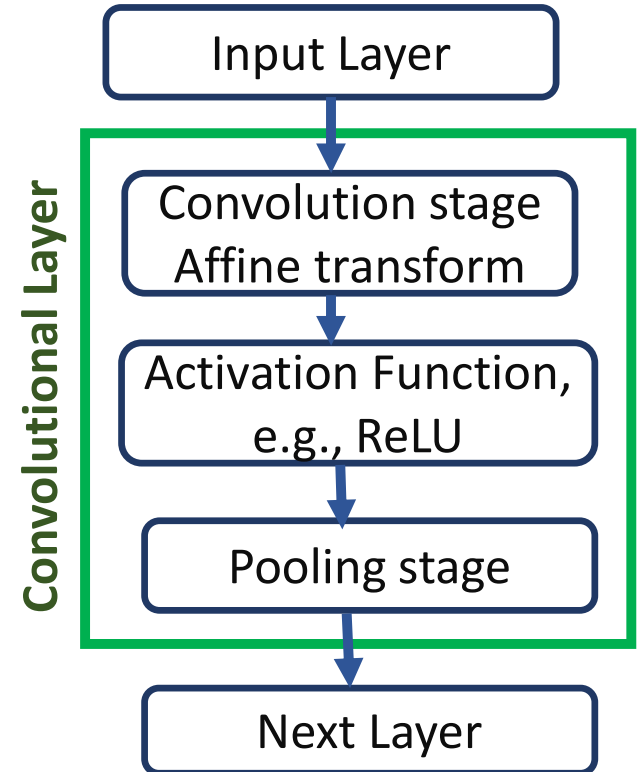
- The layer performs several convolutions in parallel to produce a set of pre-activations

❑ Second Stage: Non-linearity

- Each pre-activation is run through a nonlinear activation function

❑ Third Stage: Pooling

- Retaining the statistics of the data



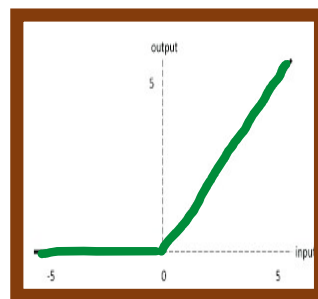
Activation Function

Feature Map

-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33



ReLU Activation



Feature Map, Post ReLU

0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0

Pooling

- Pooling is used to reduce the spatial volume of input image after convolution
- **Motivation:** We care about presence of features, not their exact location !
 - Dimensionality Reduction
 - Prevents overfitting

0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0

Max Pooling

1	1
0.33	0.33
0.33	0.33
0	0

Average Pooling

0.28	0.28
0.22	0.22
0.16	0.16
0	0

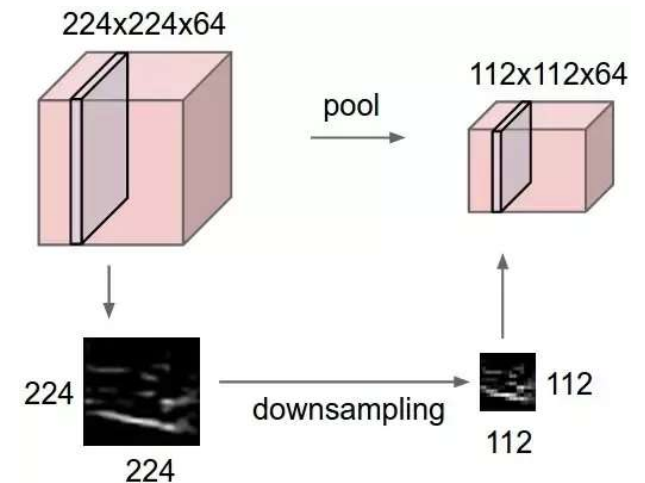
Summary of Pooling

Hyperparameters:

- ✓ f : filter size
- ✓ s : stride
- ✓ Max or Average pooling

Note: Pooling does not learn any parameters

$$n_H \times n_W \times n_c \xrightarrow{\text{Pooling}} \left[\frac{n_H - f}{s} + 1 \right] \times \left[\frac{n_W - f}{s} + 1 \right] \times n_c$$



Summary

-1	1	1	1	-1
-1	1	-1	1	-1
-1	1	1	1	-1
-1	-1	-1	1	-1
-1	-1	-1	1	-1
-1	-1	1	-1	-1
-1	1	-1	-1	-1

*

Loopy pattern filter

1/9	1/9	1/9
1/9	-1/9	1/9
1/9	1/9	1/9



-0.11	1	-0.11
-0.55	0.11	-0.33
-0.33	0.33	-0.33
-0.22	-0.11	-0.22
-0.33	-0.33	-0.33



ReLU
Activation

0	1	0
0	0.11	0
0	0.33	0
0	0	0
0	0	0



Max
Pooling

1	1
0.33	0.33
0.33	0.33
0	0

Summary

1	1	1	-1	-1
1	-1	1	-1	-1
1	1	1	-1	-1
-1	-1	1	-1	-1
-1	-1	1	-1	-1
-1	1	-1	-1	-1
1	-1	-1	-1	-1

*

Loopy pattern filter

1/9	1/9	1/9
1/9	-1/9	1/9
1/9	1/9	1/9



1	-0.11	-0.11
0.11	-0.33	0.33
0.33	-0.33	-0.33
-0.11	-0.55	-0.33
-0.55	-0.33	-0.55



ReLU
Activation

1	0	0
0.11	0	0.33
0.33	0	0
0	0	0
0	0	0

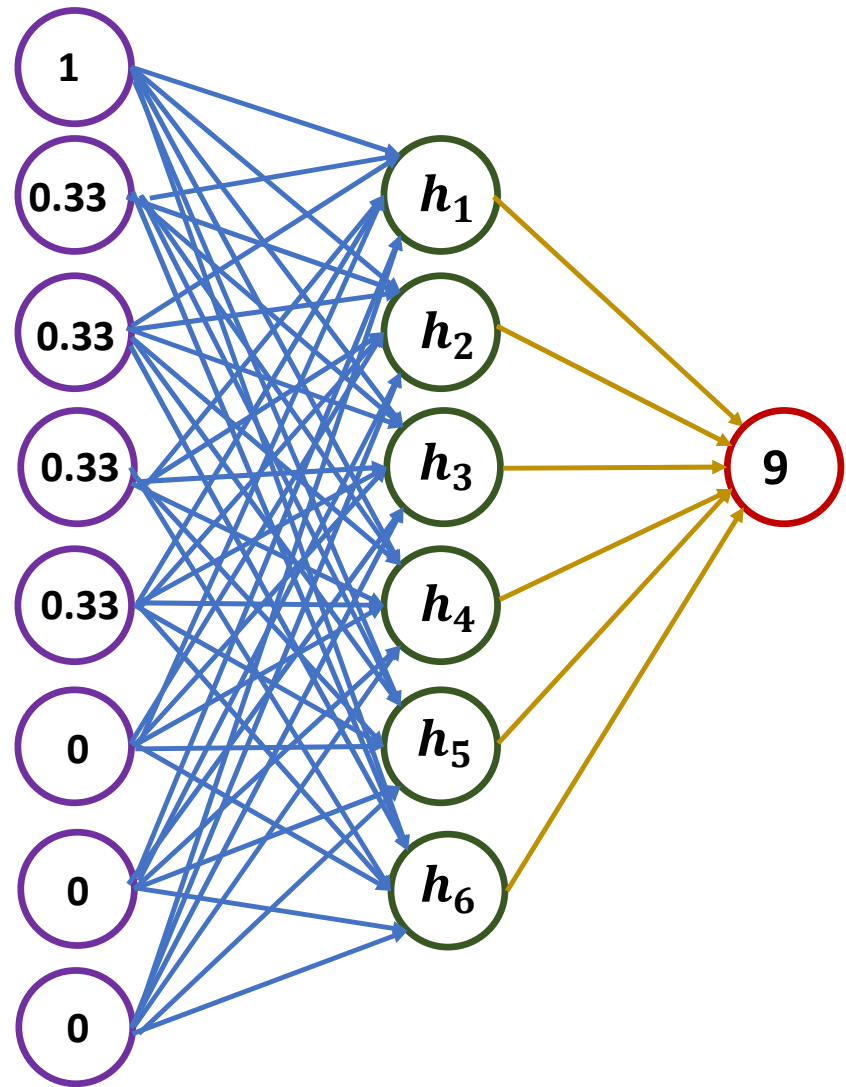


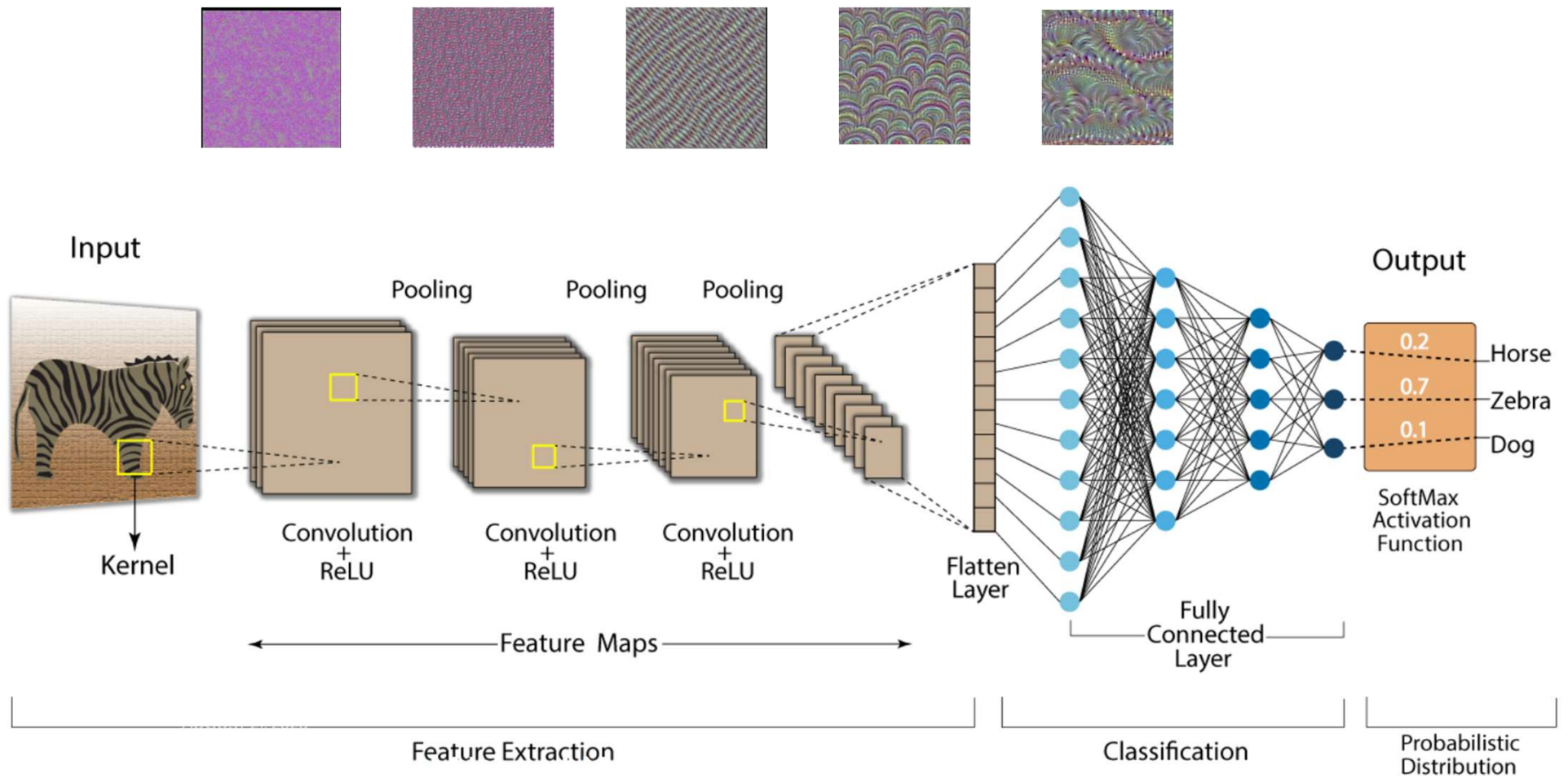
Max
Pooling

1	0.33
0.33	0.33
0.33	0
0	0

Classification

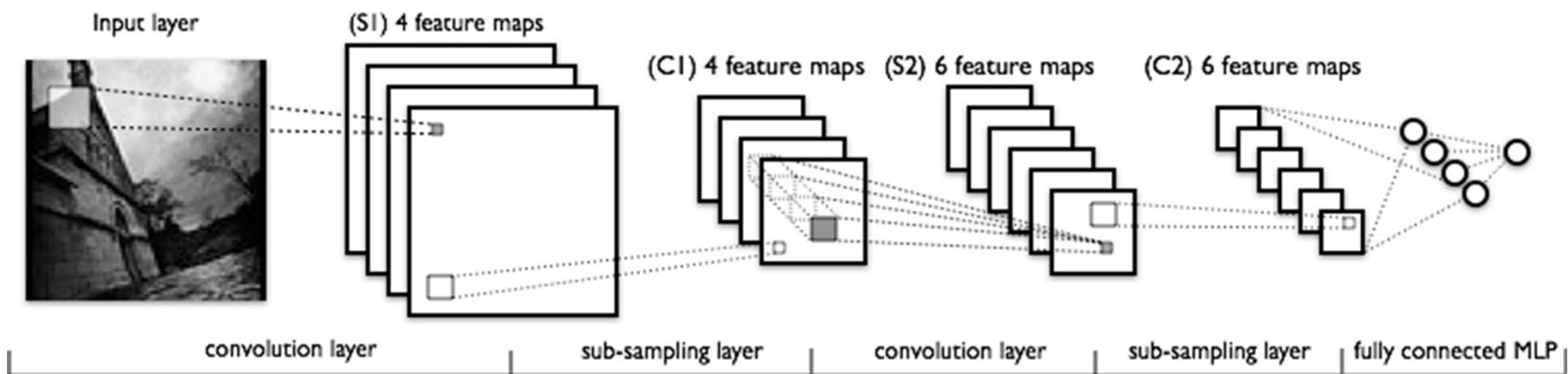
1	0.33
0.33	0.33
0.33	0
0	0





Convolution Neural Network Architecture

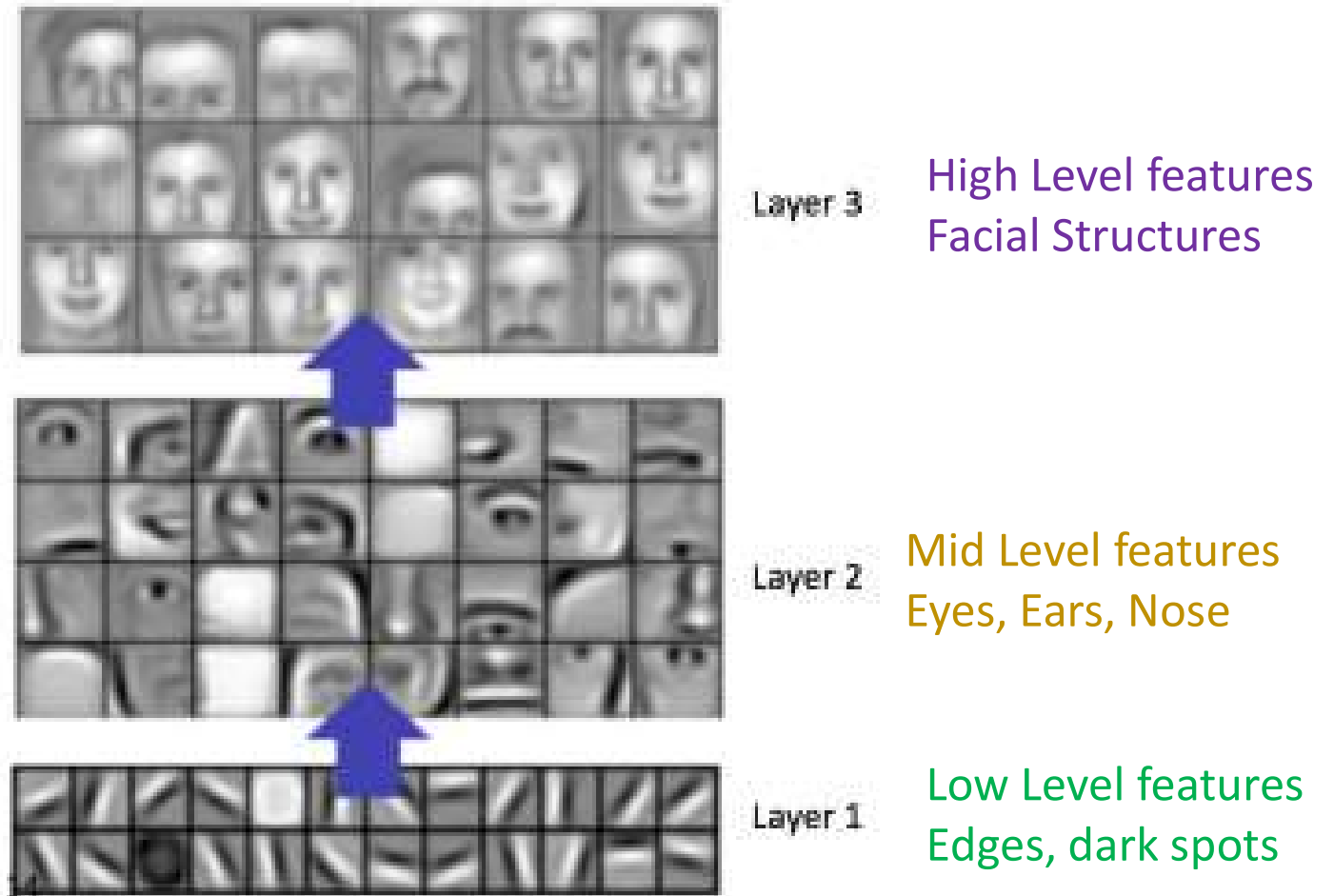
Convolutional Neural Networks (CNN)



- Learn features in input image through convolution
- Introduce non-linearity through activation function (real-world data is non-linear!)
- Reduce dimensionality and preserve spatial invariance with pooling
- Fully connected layer uses these features for classifying input image

Learned Features

Learn **hierarchy of features** directly from data (rather than hand-engineering them)



<http://web.eecs.umich.edu/~honglak/icml09-ConvolutionalDeepBeliefNetworks.pdf>

Summary: Layer Patterns

INPUT → **[CONV → RELU]*N** → **POOL?***M → **[FC → RELU]*K** → **CLASSIFIER**

* -> Repetition
? -> Optional

Common ConvNet Architectures:

- INPUT → FC → CLASSIFIER
- INPUT → CONV → RELU → FC → CLASSIFIER
- INPUT → [CONV → RELU → POOL]*2 → FC → RELU → CLASSIFIER
- INPUT → [CONV → RELU → CONV → RELU → POOL]*3 → [FC → RELU]*2 → CLASSIFIER