

Artificial Intelligence: Deep Learning

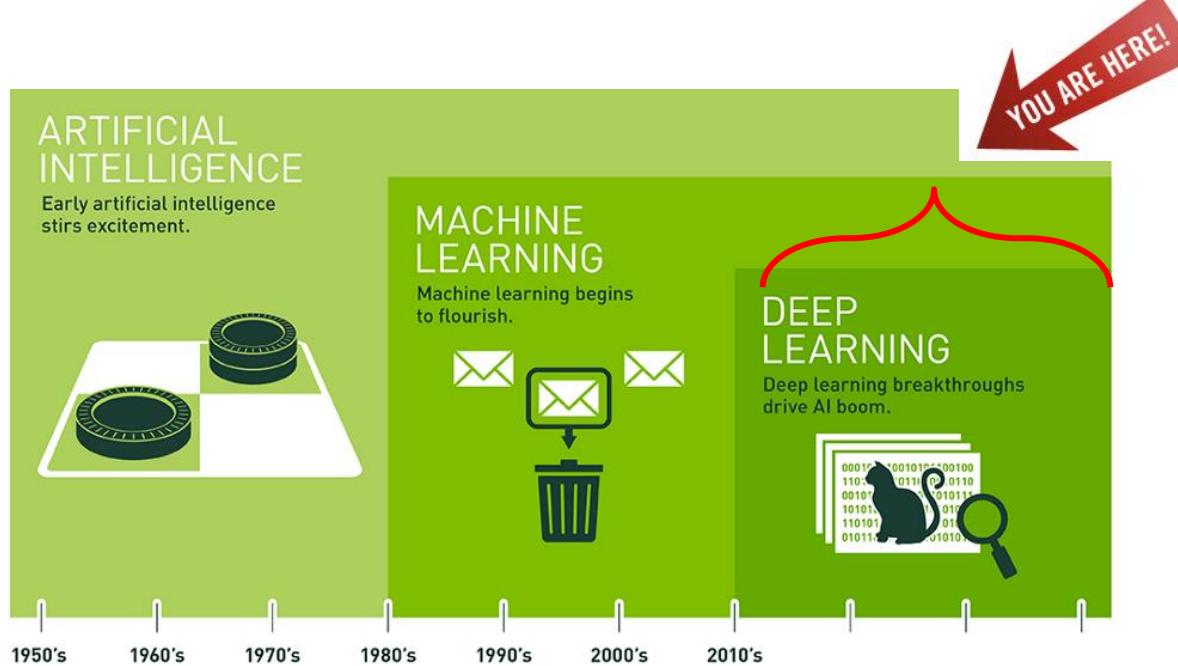
many slides from: Y. Bengio, A. Ng and Y. LeCun

Today

1. Motivation
2. Feature Learning
3. Building Block Design of CNN
 - Convolutional Layer
 - Activation Function
 - Pooling Layer
 - Normalization Layer

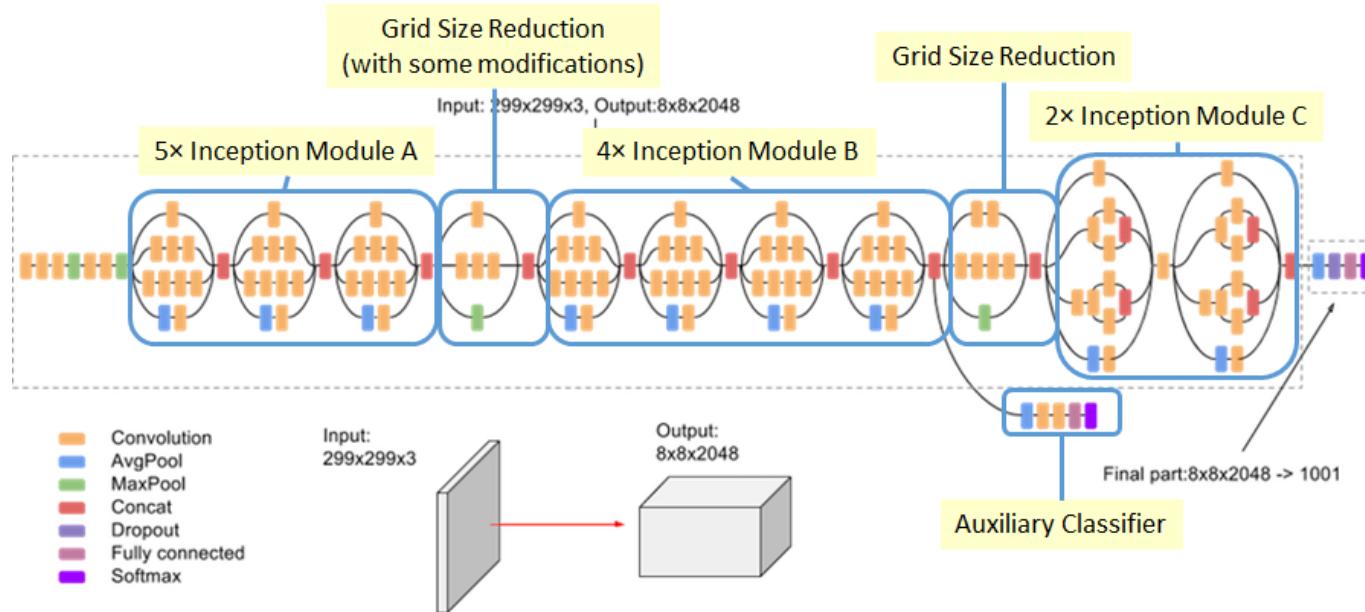


History of AI



What is Really Deep Learning?

1. Good old Neural Networks, with more layers/modules
2. Non-linear, hierarchical, abstract representations of data
3. Flexible models with any input/output type and size
4. Differentiable Functional Programming



The Google "Inception" deep neural network architecture for image recognition (27 layers)
4

Why Deep Learning Now?

1. Better algorithms & understanding
2. Computing power (GPUs, TPUs, ...)
3. Data with labels
4. Open-source tools and models



PYTORCH



Microsoft
CNTK

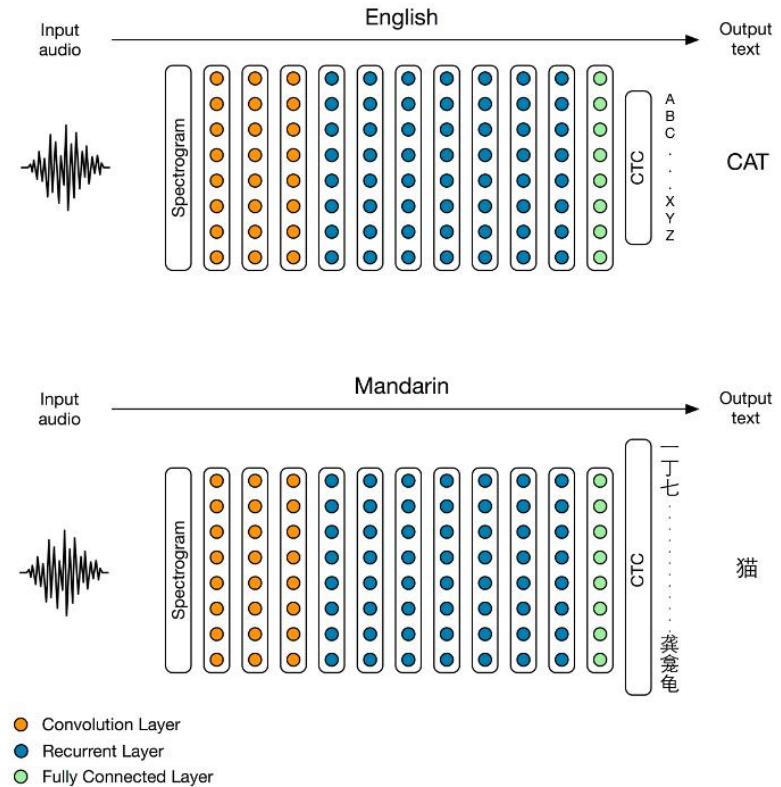
Caffe2

dmlc
mxnet

gensim spaCy

theano

DL Today: Speech-to-Text



[Baidu 2014]

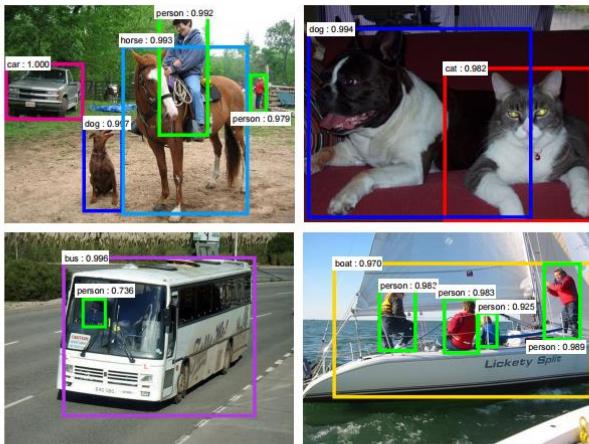
DL Today: Vision



[Krizhevsky 2012]



[Ciresan et al. 2013]



[Faster R-CNN - Ren 2015]

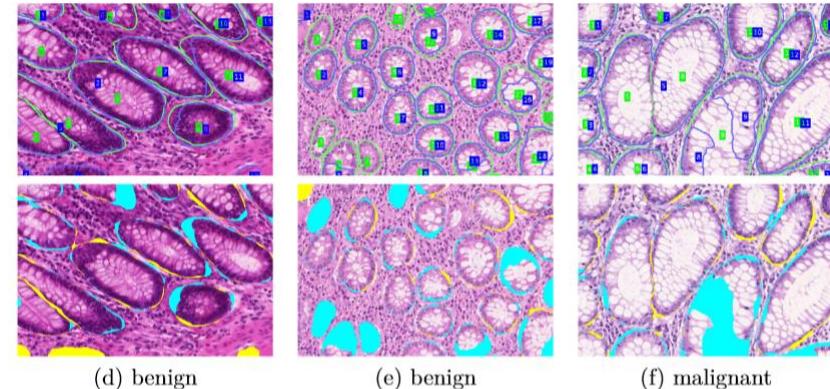


[NVIDIA dev blog]

DL Today: Vision



[Stanford 2017]



[Nvidia Dev Blog 2017]

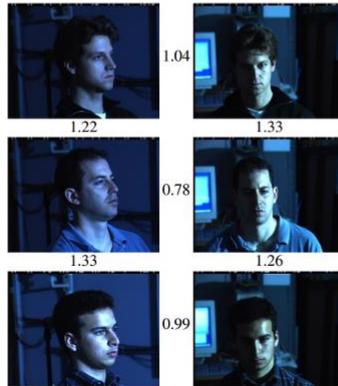
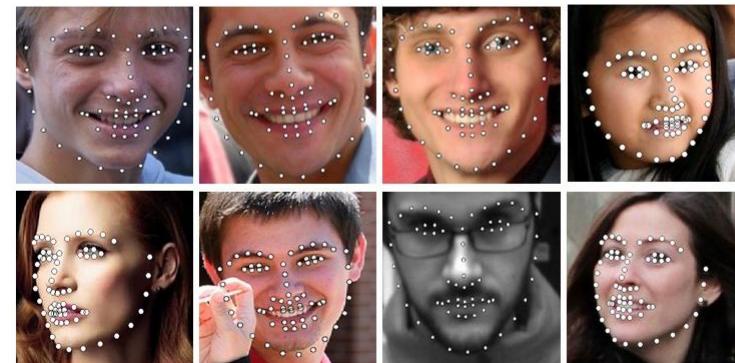


Figure 1. Illumination and Pose invariance.

[FaceNet - Google 2015]



[Facial landmark detection CUHK 2014]

DL Today: NLP

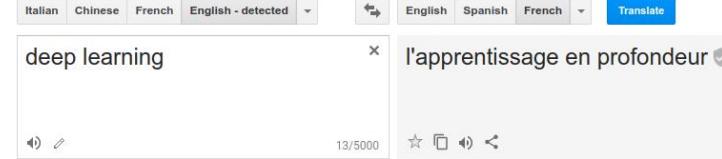
Translate

Italian Chinese French English - detected English Spanish French Translate

deep learning × l'apprentissage en profondeur ✓

13/5000

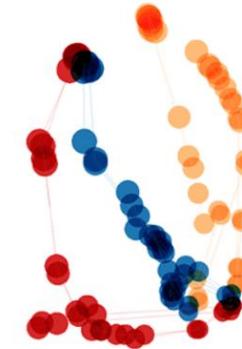
See also
deep, learning



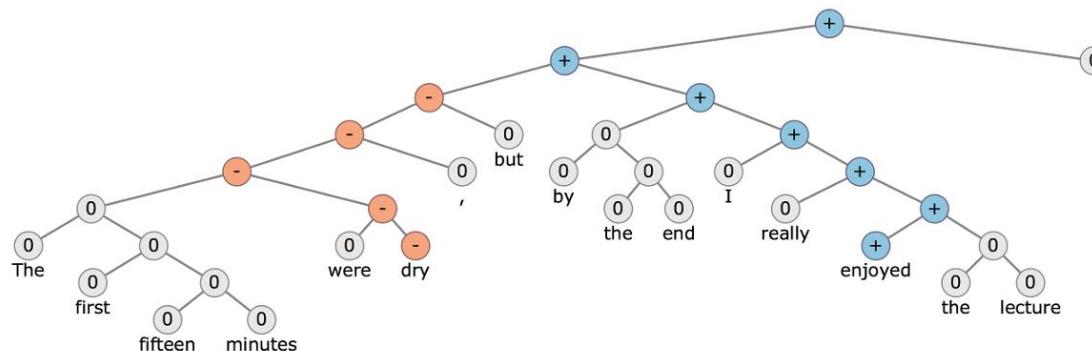
ENGLISH
The stratosphere extends from about 10km to about 50km in altitude.

KOREAN

JAPANESE



[Google Translate System - 2016]



[Socher 2015]

DL Today: NLP



Salit Kulla

11:29 AM •••

to me

Hey, Wynton Marsalis is playing this weekend. Do you have a preference between Saturday and Sunday?

-S

I'm down for either.

Let's do Saturday.

I'm fine with whatever.



Reply



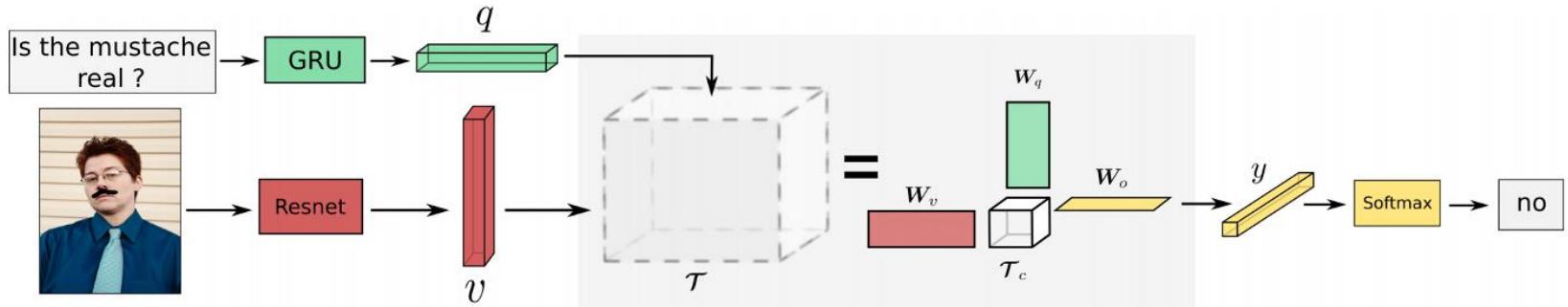
Forward



[Google Inbox Smart Reply]

[Amazon Echo / Alexa]

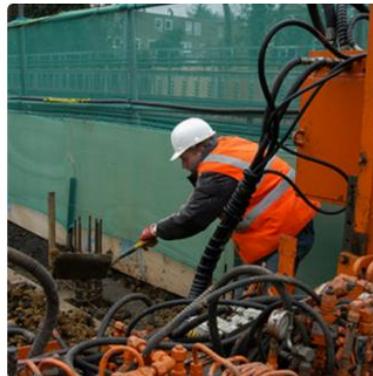
DL Today: Vision+NLP



[VQA - Mutan 2017]



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."

[Karpathy 2015]

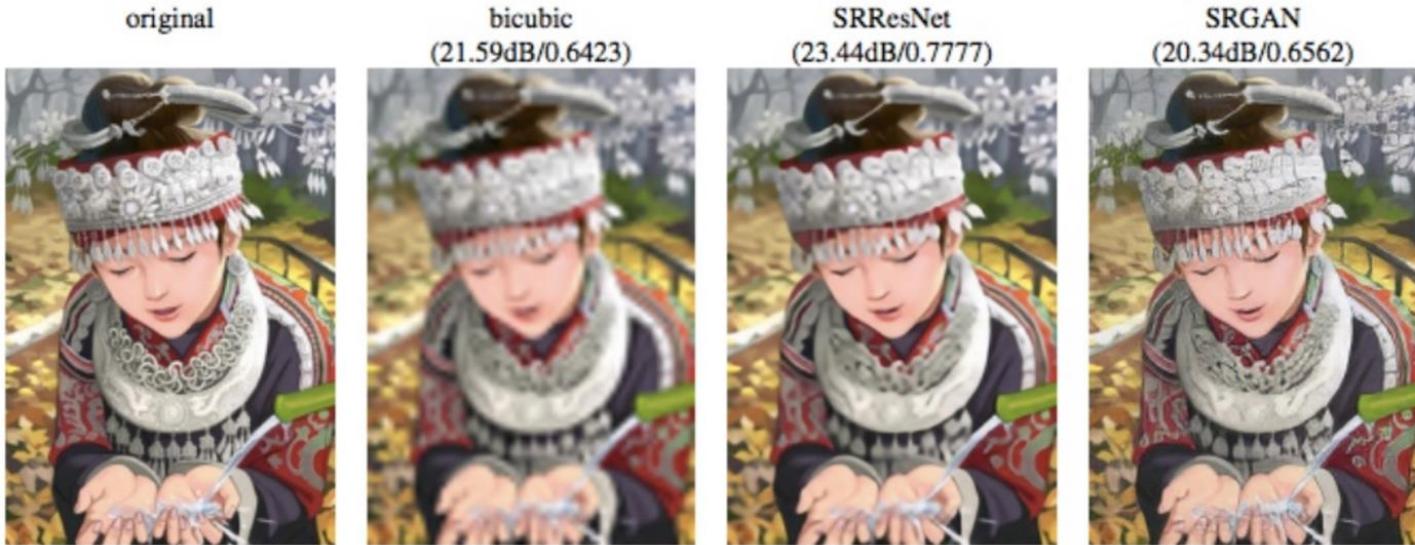
DL Today: Image Translation



[DeepDream 2015]



[Gatys 2015]



[Ledig 2016]

DL Today: Generative Models

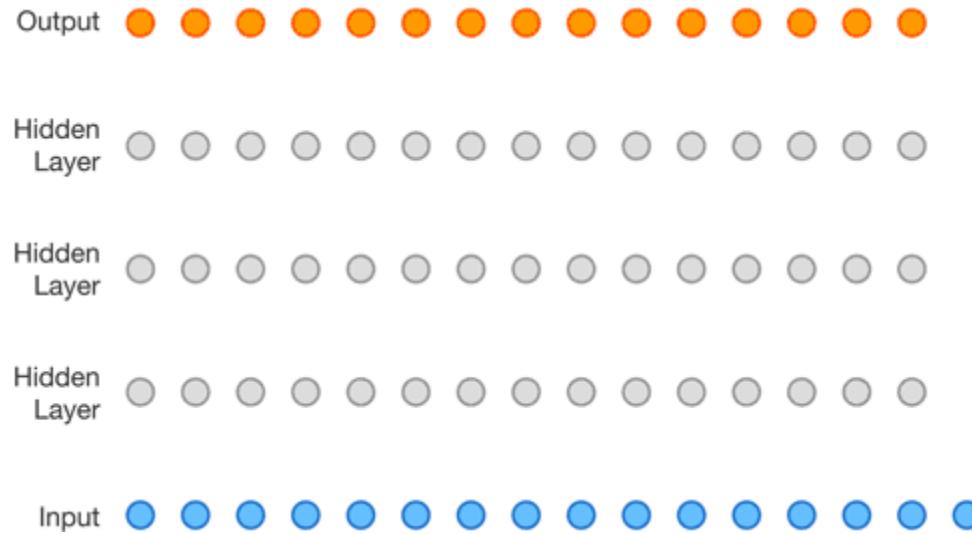


Sampled celebrities [Nvidia 2017]

Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak
Stage-I images							
Stage-II images							

StackGAN v2 [Zhang 2017]

DL Today: Generative Models



Guess which one is generated?



DL Today: Language/Image Models

Open-AI GPT-3, or DALL-E: <https://openai.com/blog/dall-e/>

TEXT PROMPT

an armchair in the shape of an avocado [...]

AI-GENERATED IMAGES



[View more or edit prompt ↓](#)

TEXT PROMPT

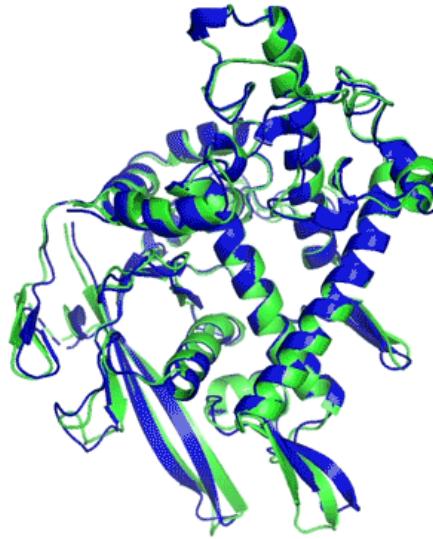
a store front that has the word 'openai' written on it [...]

AI-GENERATED IMAGES

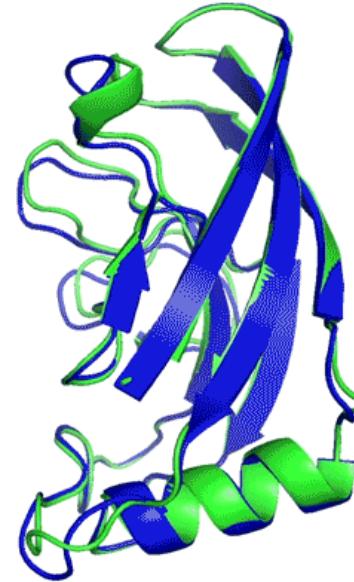


[View more or edit prompt ↓](#)

DL Today: Genomics



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)



T1049 / 6y4f
93.3 GDT
(adhesin tip)

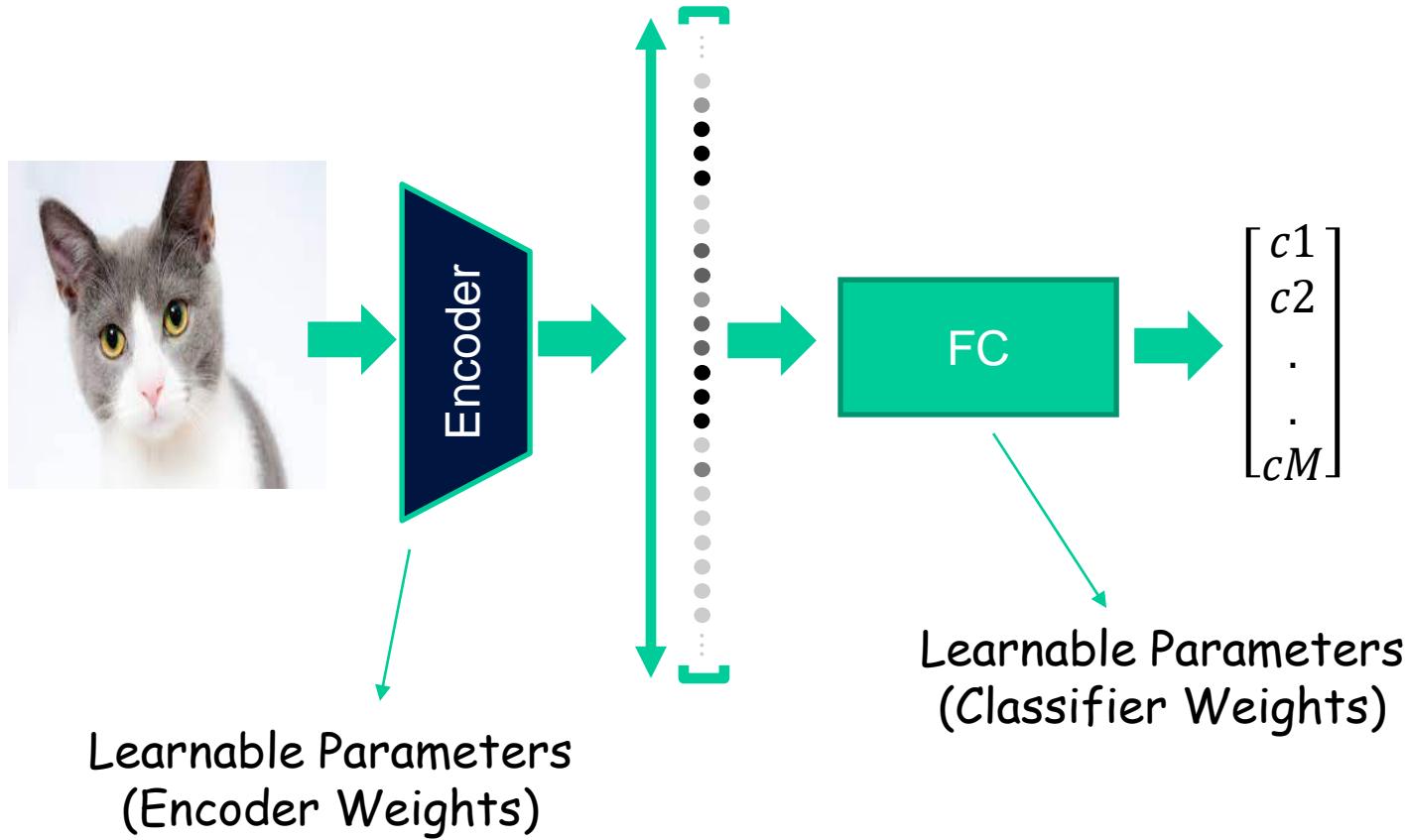
- Experimental result
- Computational prediction

[AlphaFold by DeepMind](#)

Today

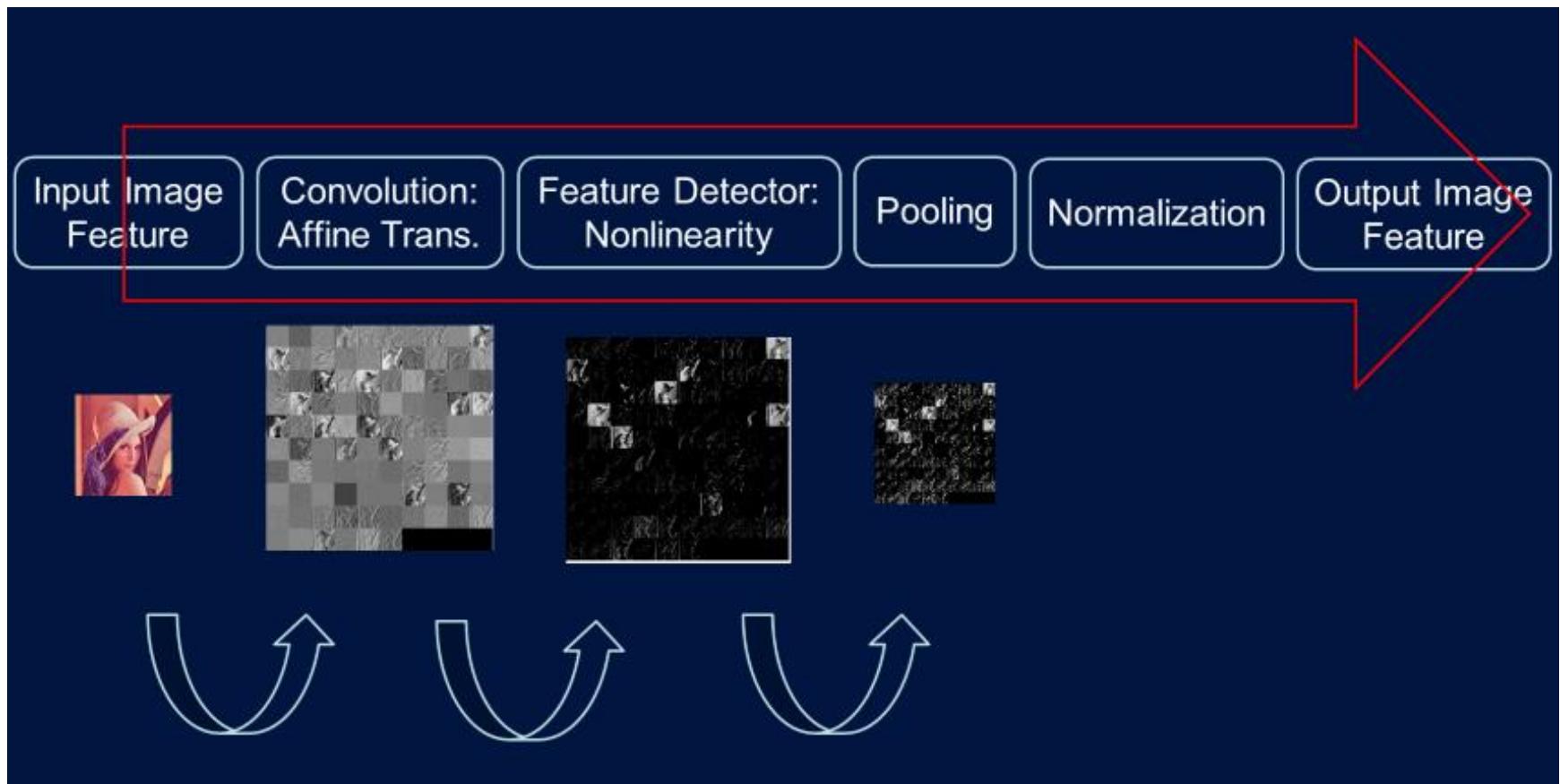
1. Motivation
2. Feature Learning 
3. Building Block Design of CNN
 - Convolutional Layer
 - Activation Function
 - Pooling Layer
 - Normalization Layer

Feature (Encoder) Learning



Encoder: Convolution Layer

Higher level of feature abstraction

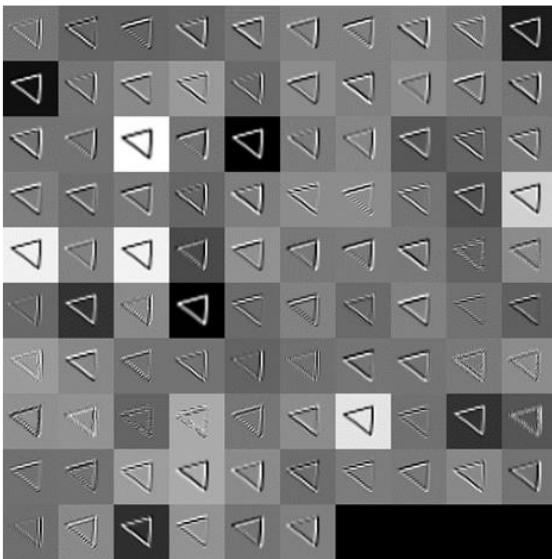
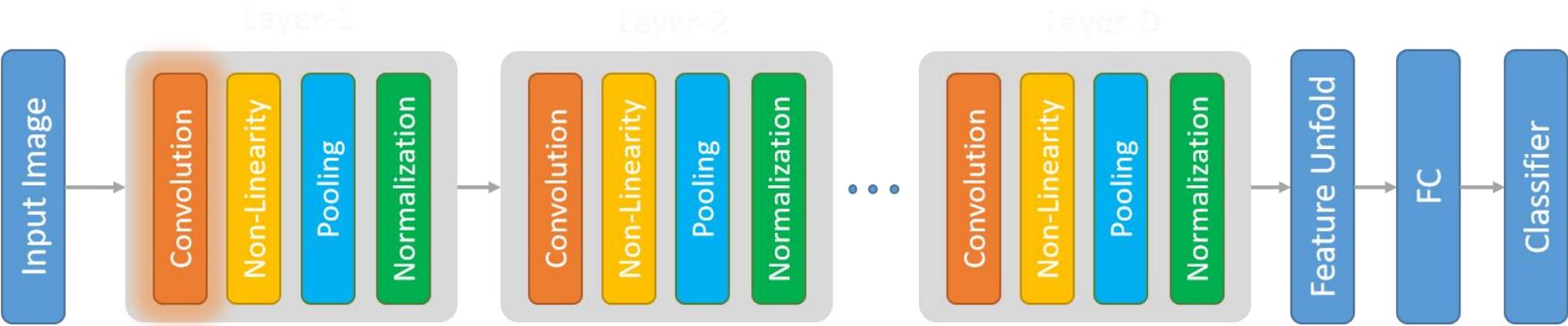


Encoder: Cascade Layer Design

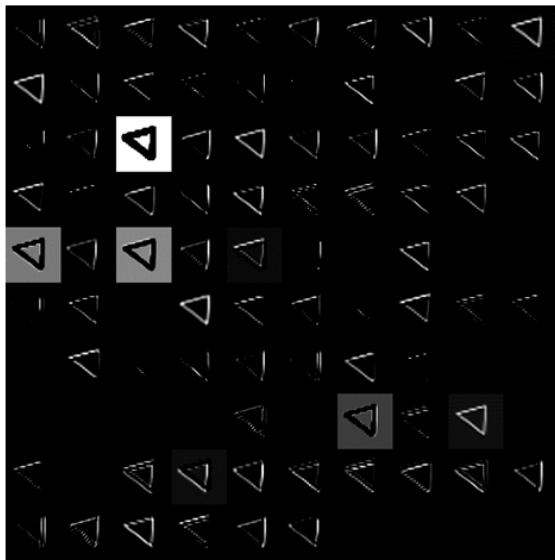


△

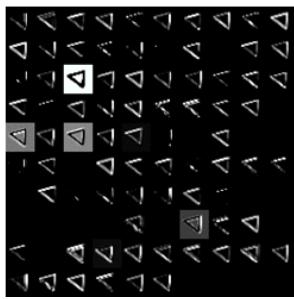
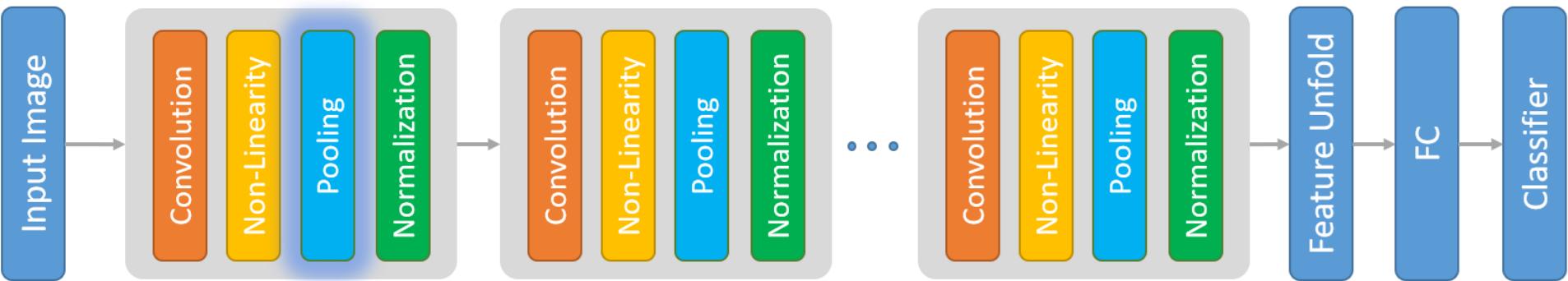
Encoder: Cascade Layer Design



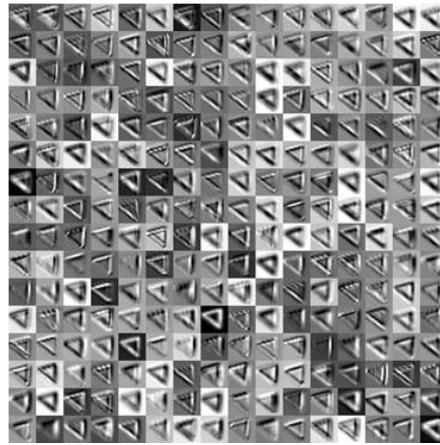
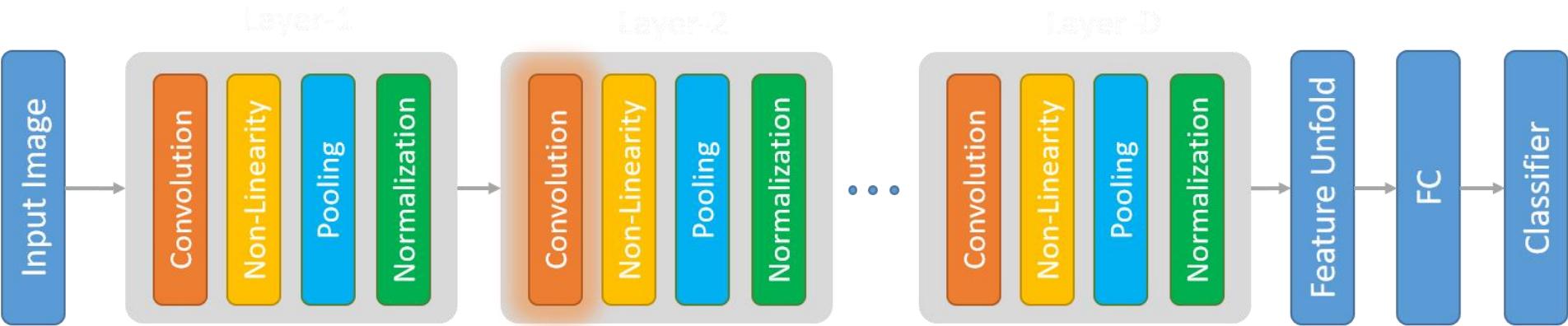
Encoder: Cascade Layer Design



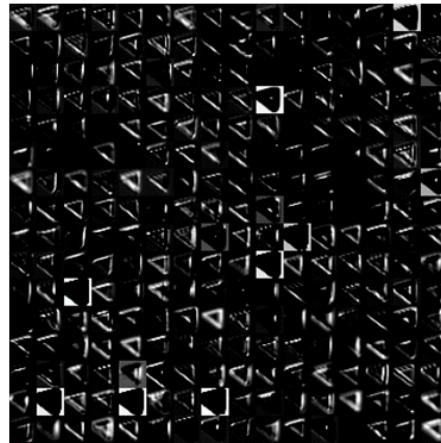
Encoder: Cascade Layer Design



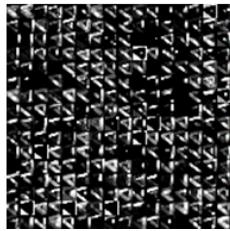
Encoder: Cascade Layer Design



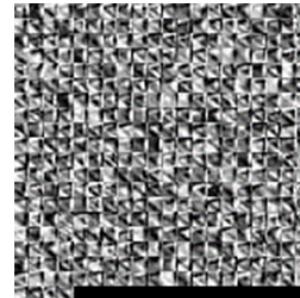
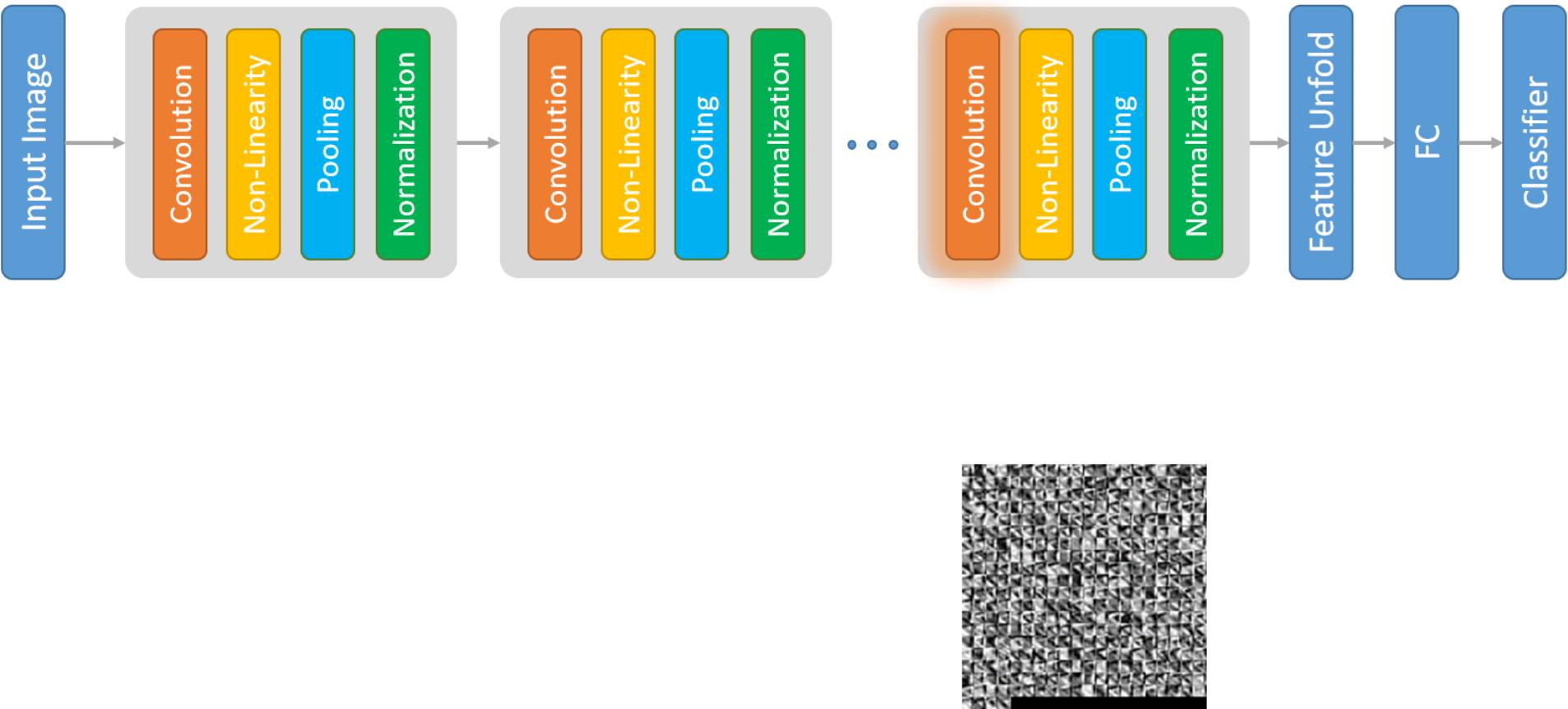
Encoder: Cascade Layer Design



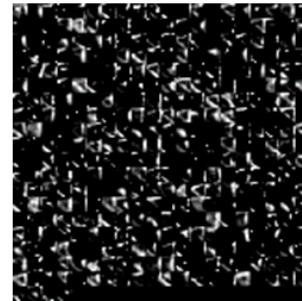
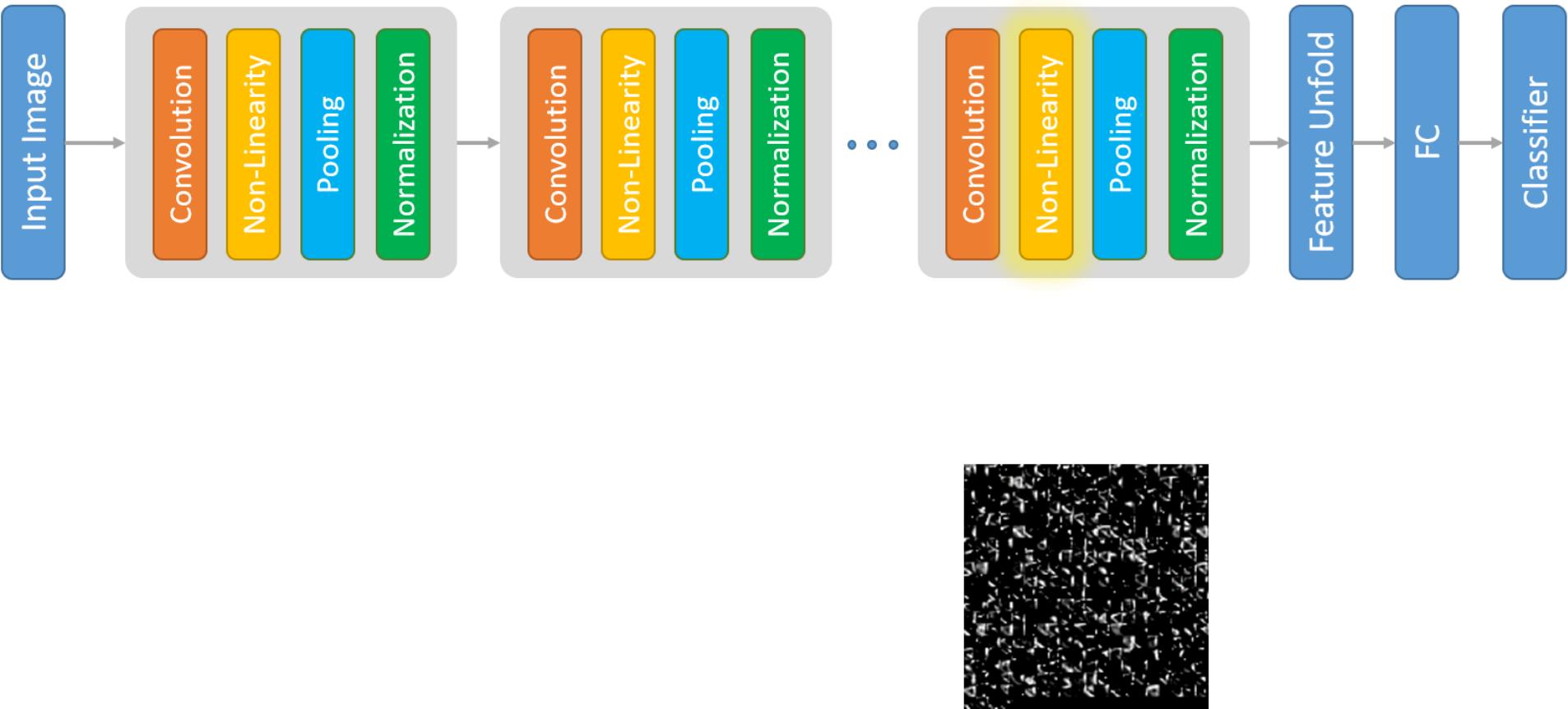
Encoder: Cascade Layer Design



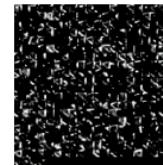
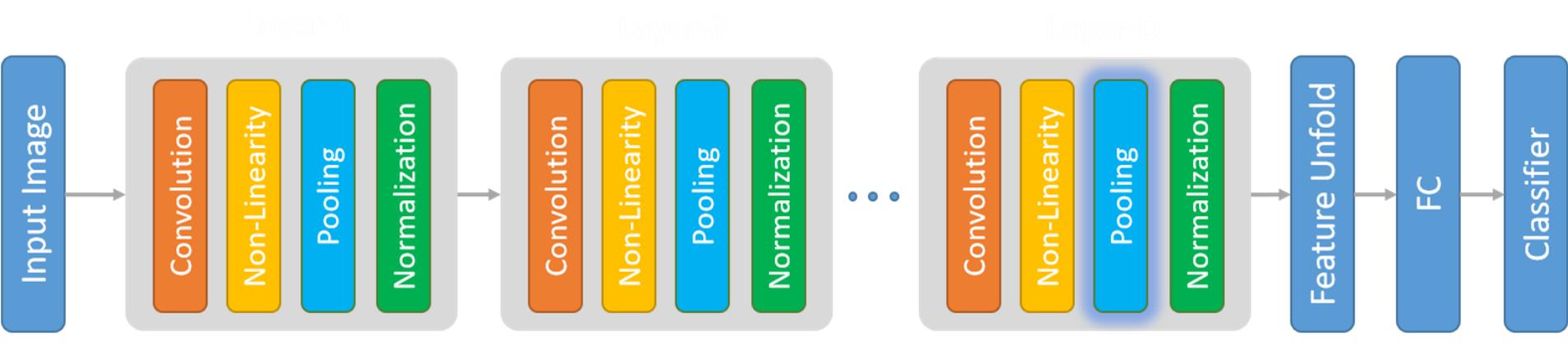
Encoder: Cascade Layer Design



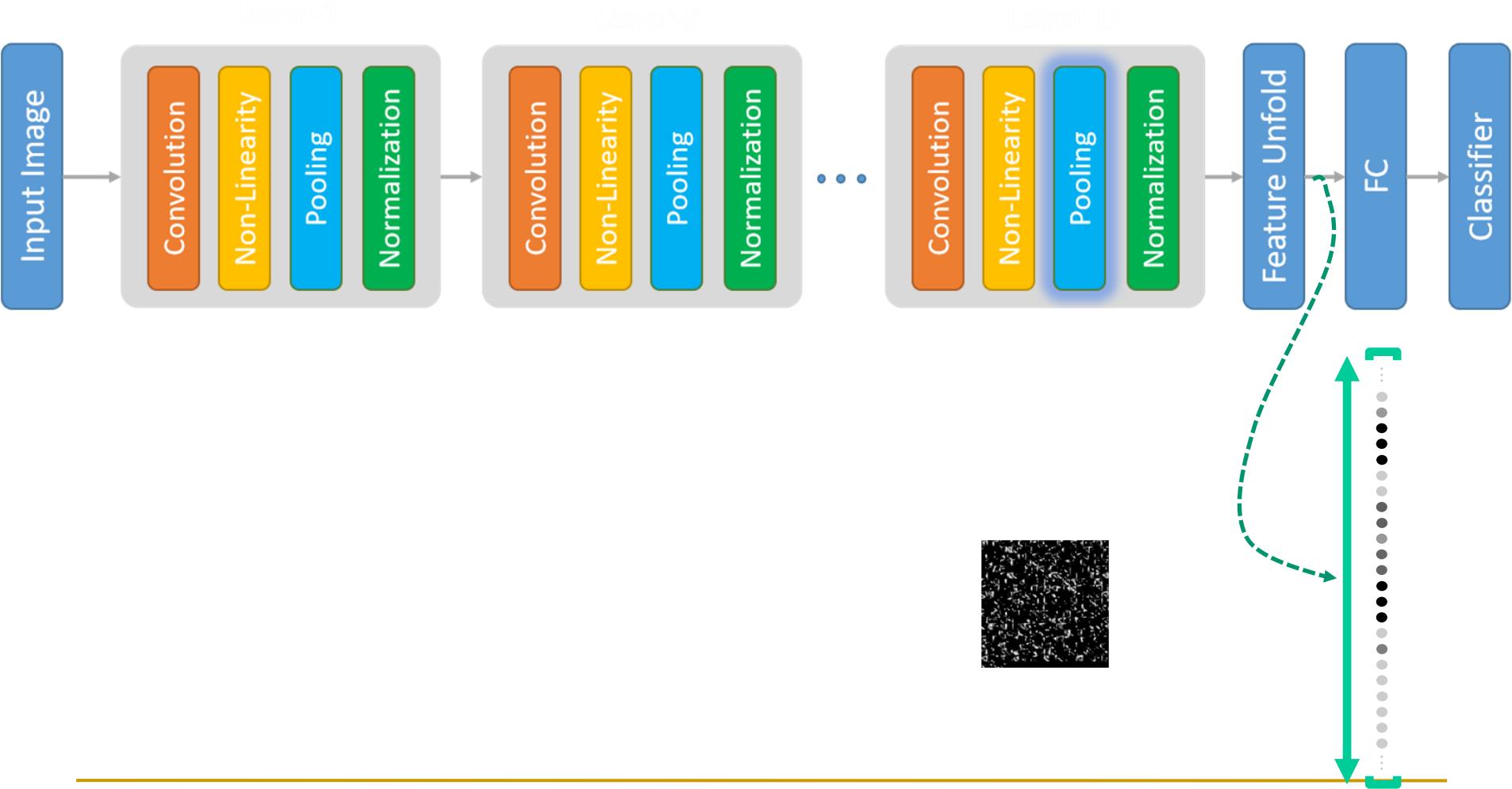
Encoder: Cascade Layer Design



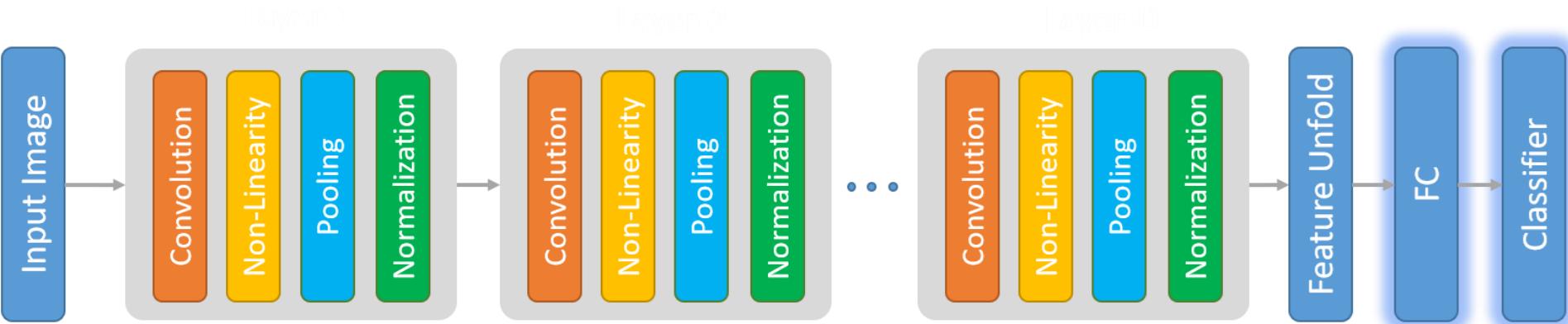
Encoder: Cascade Layer Design



Encoder: Cascade Layer Design



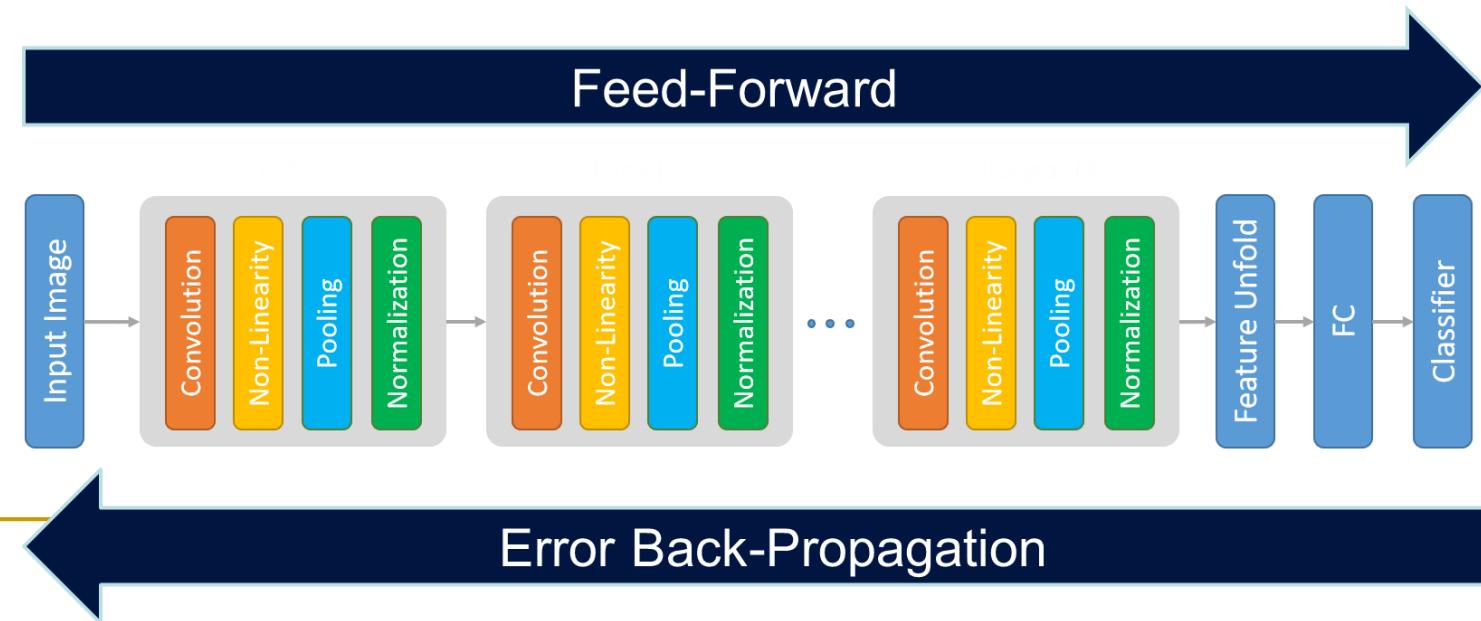
Encoder: Cascade Layer Design



y_1 : Square
 y_2 : Triangle
 y_3 : Circle
⋮
 y_N : Dimond

Encoder: Overall Training

1. Initialize all learnable parameters (e.g. random sampling)
2. Preprocess train image dataset and randomly shuffle them
3. Feed-forward images in mini-batches and compute "Target Predictions"
4. Calculate Error = Abs("Target Labels" - "Target Predictions")
5. Propagate back error gradients and adjust weights using an optimization method (e.g. stochastic gradient descent)
6. Repeat (1)-(4) to converge to a minimal error



Initial Drawbacks

1. Standard backpropagation with sigmoid activation function does not scale well with multiple layers
 - Weight of early layers change too slowly (no learning)
2. Overfitting
 - Large network -> lots of parameters -> increased capacity to "learn by heart"
3. Multilayered ANNs need lots of labeled data
 - Most data is not labeled 

Initial Drawbacks (1)

1. Standard gradient-based backpropagation does not scale well with multiple layers...

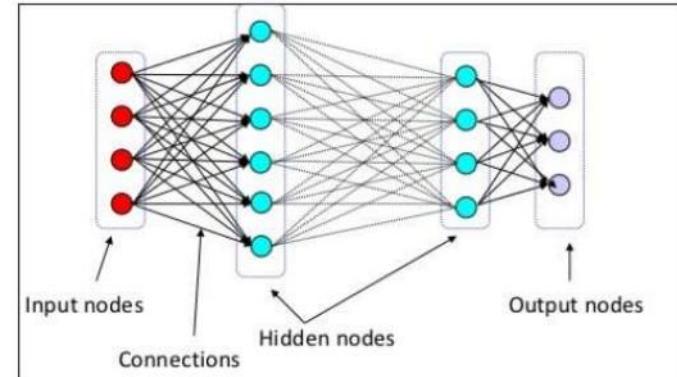
When we multiply the gradients many times (for each layer), it can lead to ...

- a) Vanishing gradient problem:

- ❑ gradients shrink exponentially with the number of layers
- ❑ so weight updates get smaller and smaller
- ❑ and weights of early layers change very slowly and network learns very very slowly

- b) Exploding gradient problem:

- ❑ multiplying gradients could also make them grow exponentially.
- ❑ so weight updates get larger and larger
- ❑ and the weights can become so large as to overflow and result in NaN values



$$\begin{aligned}\delta_h &= g'(x_h) \times Err_h \\ &= O_h (1 - O_h) \times \sum_k (w_{hk} \delta_k)\end{aligned}$$

$$\delta_6 = O_6 (1 - O_6) \times \sum (w_{6,7} \delta_7)$$

$$\delta_5 = O_5 (1 - O_5) \times \sum (w_{5,6} \delta_6)$$

$$\delta_4 = O_4 (1 - O_4) \times \sum (w_{4,5} \delta_5)$$

$$\delta_3 = O_3 (1 - O_3) \times \sum (w_{3,4} \delta_4)$$

...

Initial Drawbacks (1)

To help, we can :

- a) Use other activation functions...
- a) Do "gradient clipping" (i.e. set bounds on the gradients)

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU) [2]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

Initial Drawbacks (2)

2. Overfitting

- Large network → lots of parameters → increased capacity to "learn by heart"

- *Solutions:*

- *Regularization:*

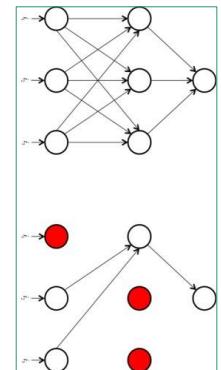
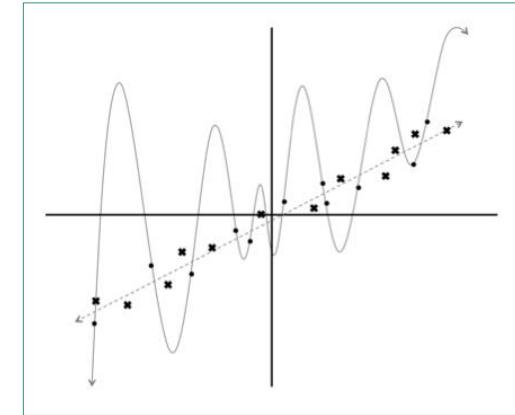
- modify the error function that we minimize to penalize large weights.

$$\frac{1}{2} \sum_{i=0}^n \frac{(T_i - O_i)^2}{n} + \lambda f(w)$$

- where $f(w)$ grows larger as the weights grow larger and λ is the regularization strength

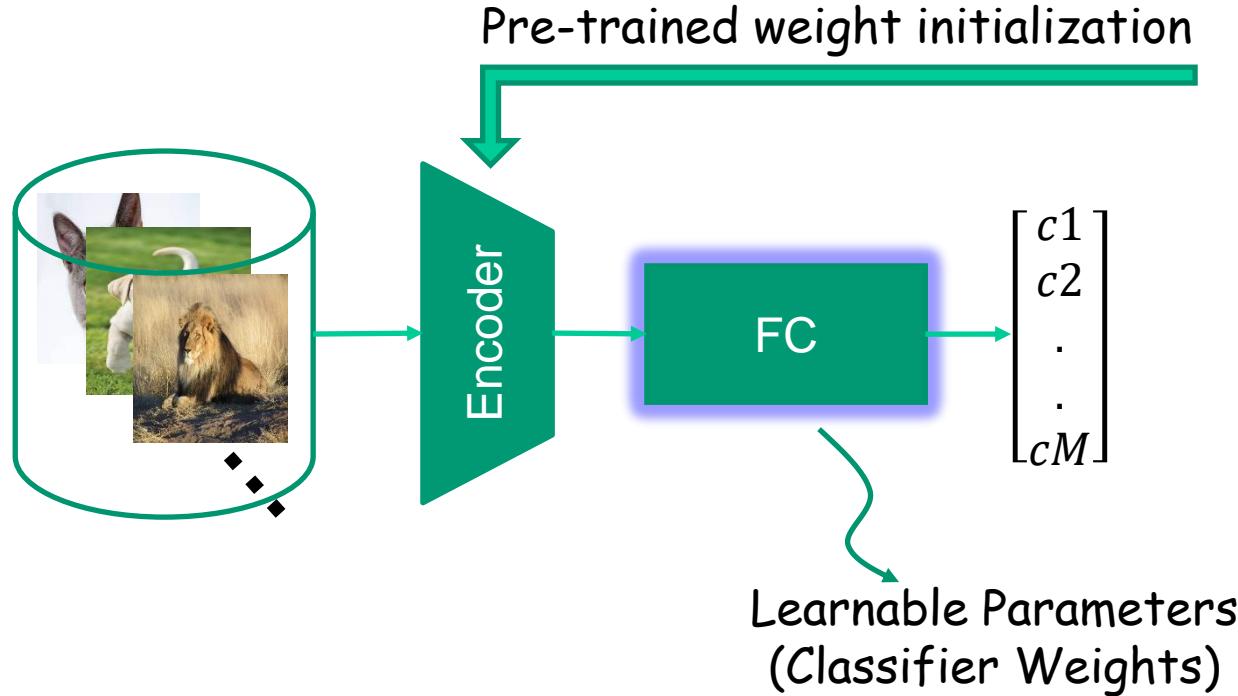
- *Dropout:*

- keep a neuron active with some probability p or setting it to zero otherwise.
 - prevents the network from becoming too dependent on any one neuron.



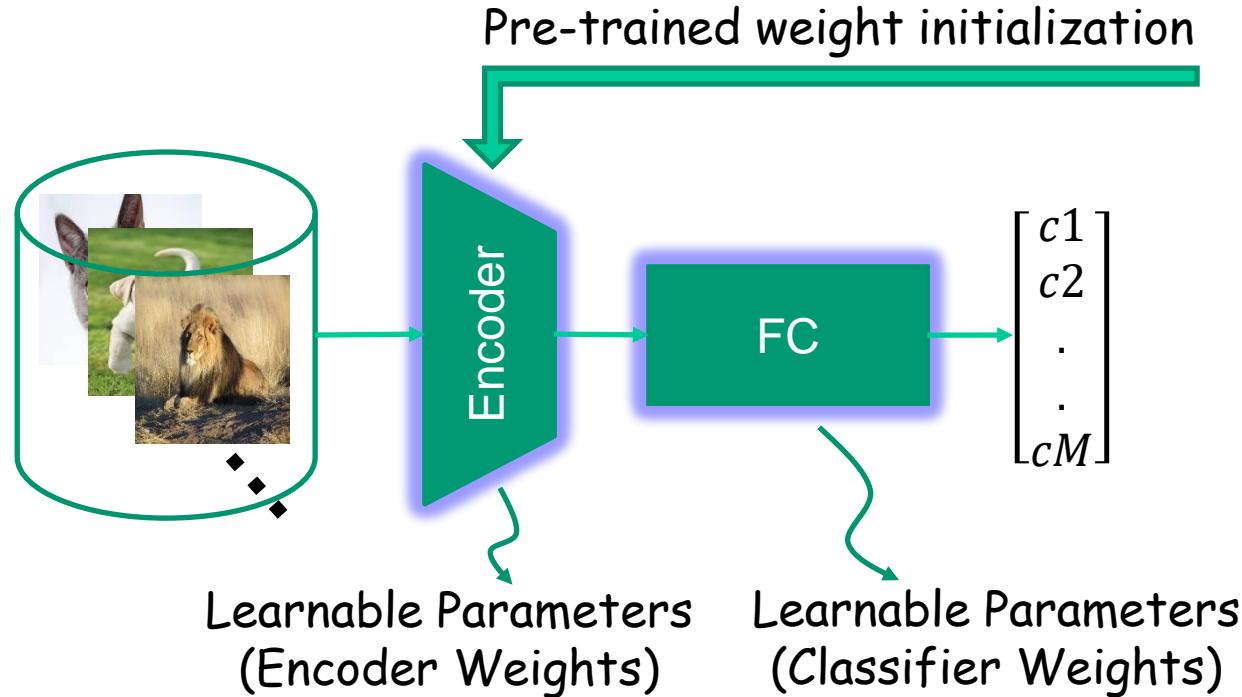
Initial Drawbacks (3)

3. Multilayered ANNs need lots of labeled data for training. To solve the problem:
 - **Transfer Learning:** use pre-trained encoder weights
 - Fine-Tuning



Initial Drawbacks (3)

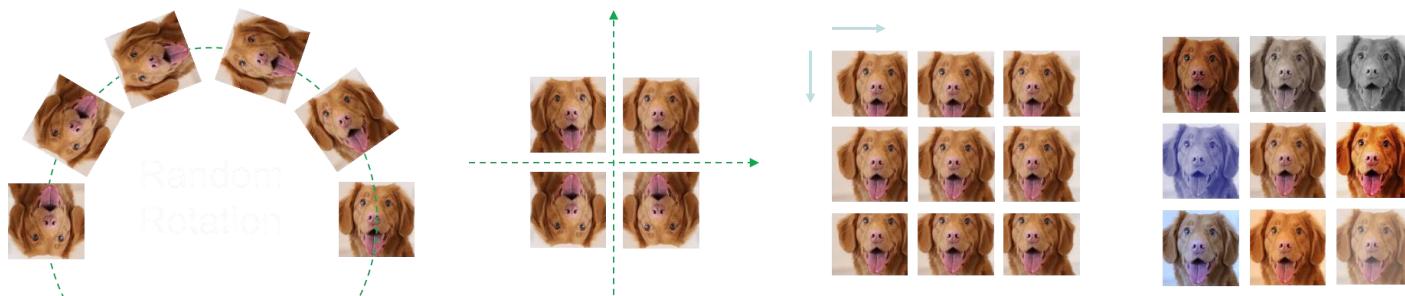
3. Multilayered ANNs need lots of labeled data for training. To solve the problem:
 - **Transfer Learning:** use pre-trained encoder weights
 - Deep-Tuning



Initial Drawbacks (3)

3. Multilayered ANNs need lots of labeled data for training. To solve the problem:

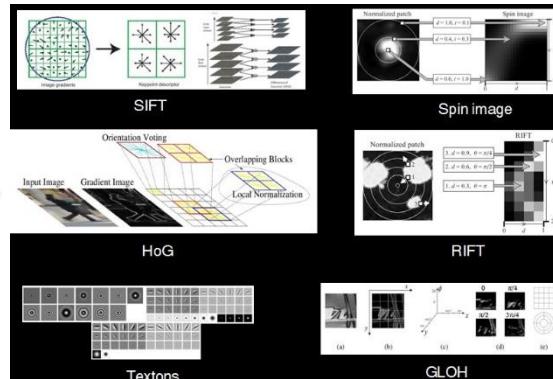
- **Image Augmentation**: randomly perturb representation
 - Shifting
 - Flipping
 - Rotation
 - Skewing
 - Color/illumination perturbation



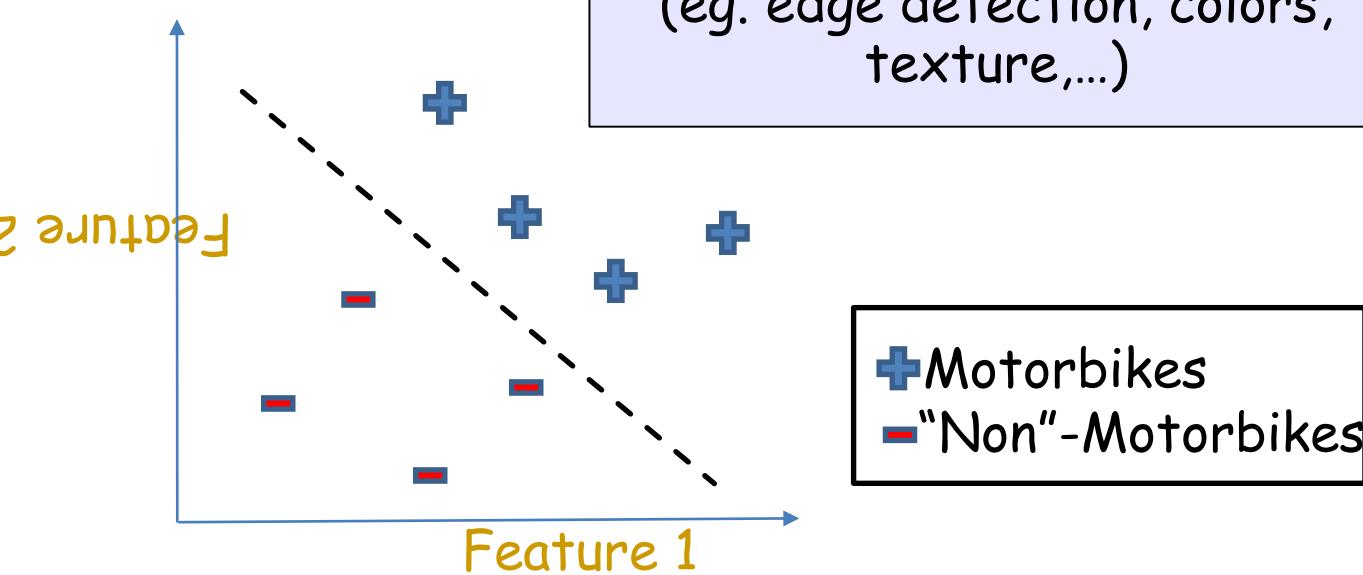
Classic ML



Input



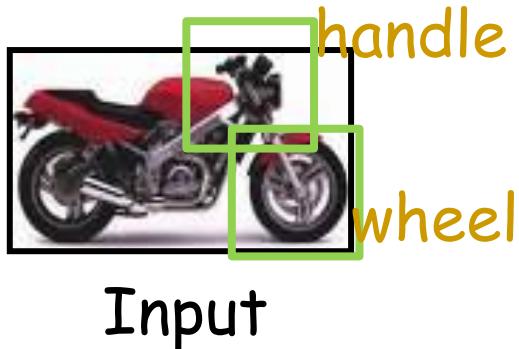
Learning algorithm



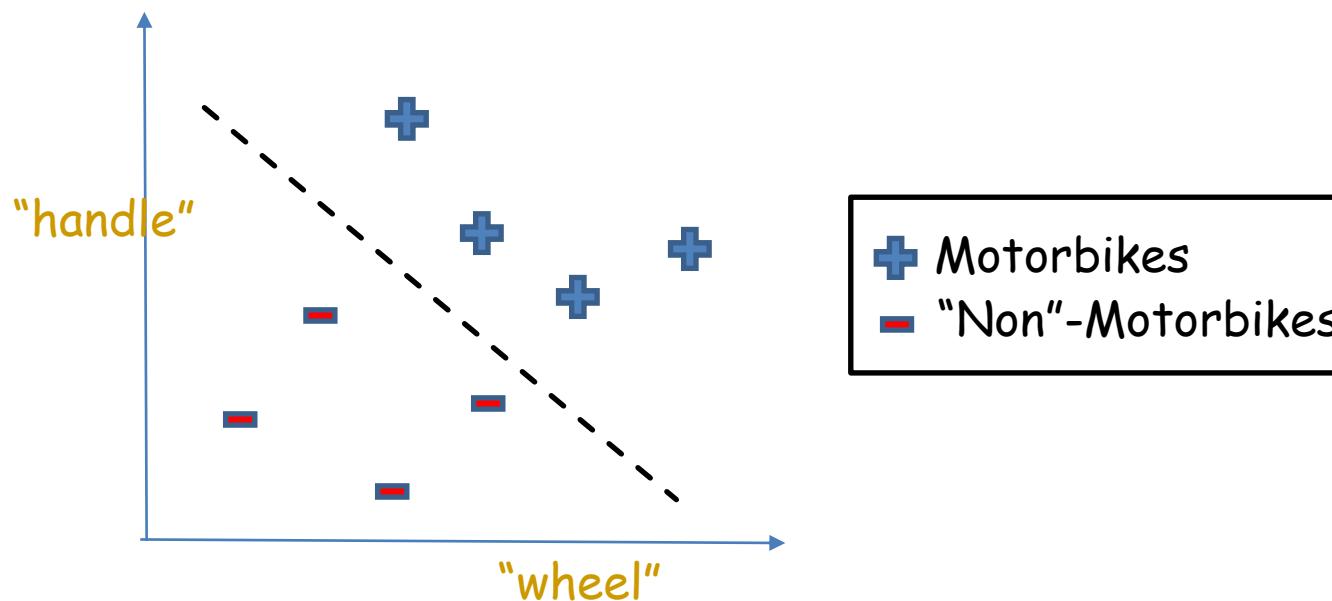
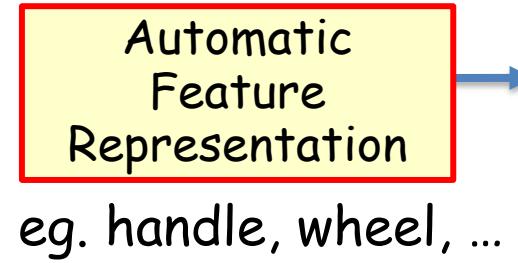
Classic ML,
requires labeled
data and hand-
crafted features

1. Needs expert knowledge
2. Time-consuming and expensive
3. Does not generalize to other domains

Automatic Feature Learning



Input

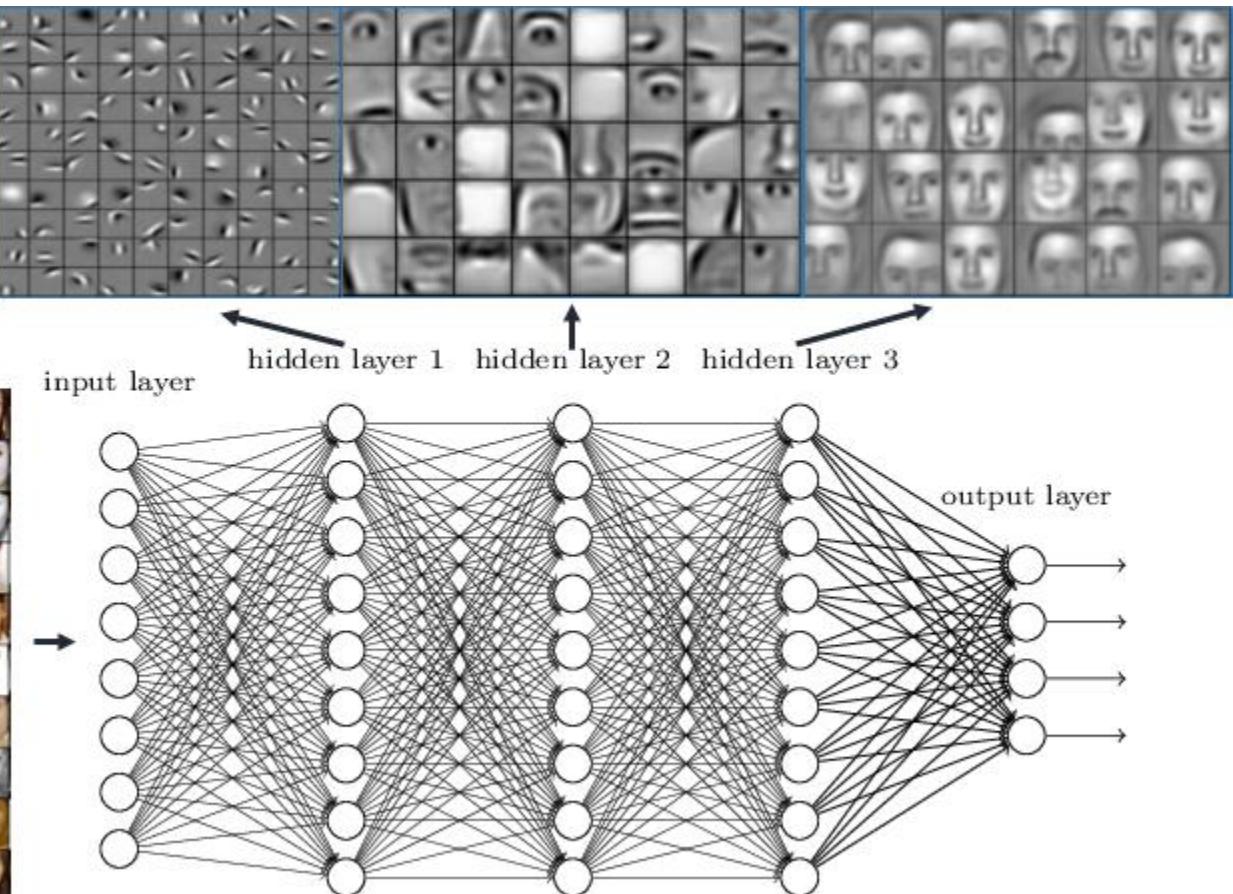


With Automatic Feature Learning:

1. We feed the network the raw data (not feature-curated)
2. The features are learned by the network
3. Features learned can be re-used in similar tasks.

Automatic Feature Learning

Deep neural networks learn hierarchical feature representations



Automatic Feature Learning

Deep Learning = Machine learning algorithms based on learning multiple levels of representation / abstraction.

- Y. Bengio

- Each layer learns more abstract features that are then combined / composed into higher-level features automatically
- Like the human brain ...
 - has many layers of neurons which act as feature detectors
 - detecting more and more abstract features as you go up
- E.g. to classify an image of a cat:
 - Bottom Layers: Edge detectors, curves, corners straight lines
 - Middle Layers: Fur patterns, eyes, ears
 - Higher Layers: Body, head, legs
 - Top Layer: Cat or Dog



Automatic Feature Learning

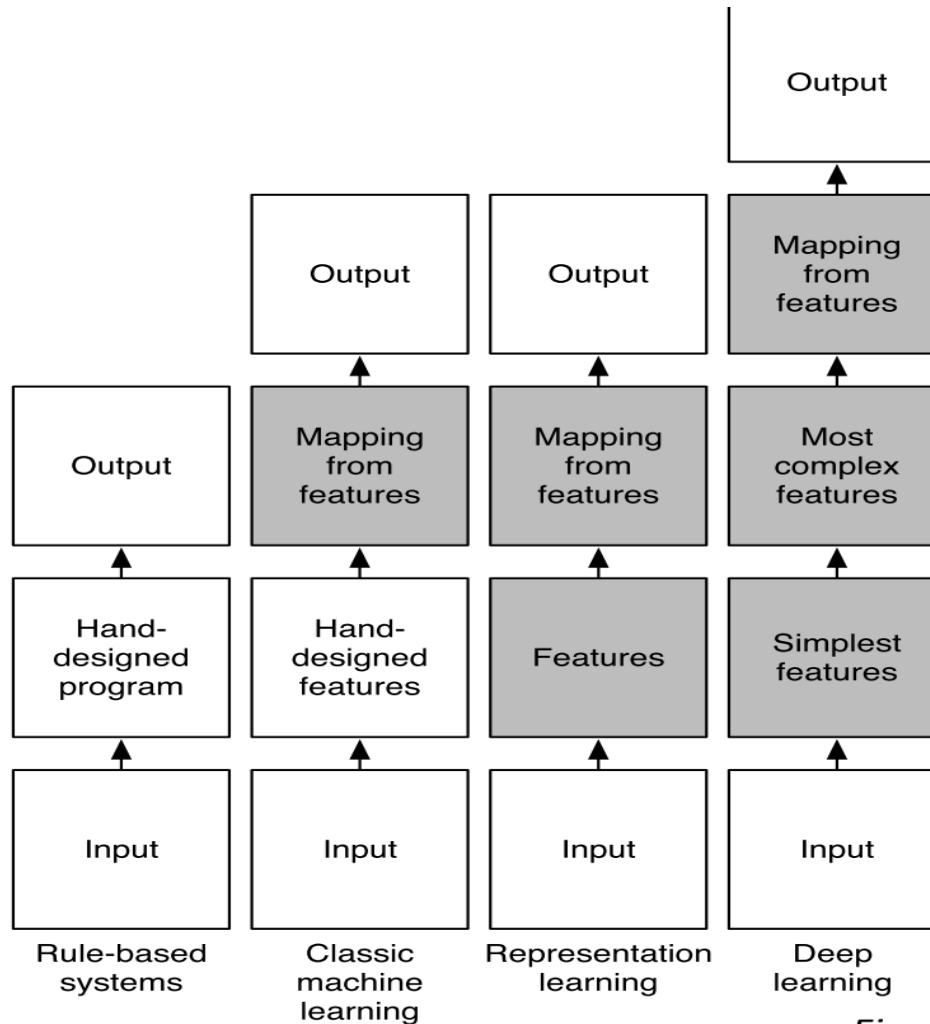
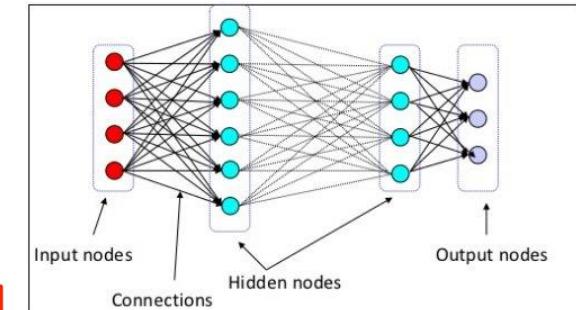
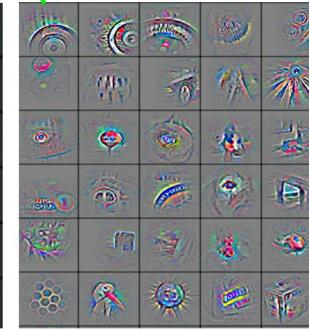
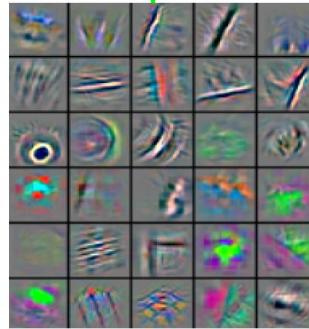
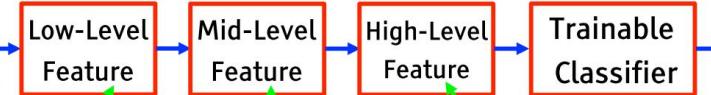


Fig: I. Goodfellow

What Types of Features?

- For image recognition
 - pixel → edge → texton → motif → part → object

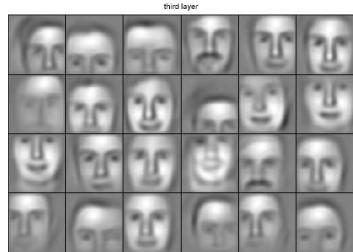


- For NLP
 - character → word → constituents → clause → sentence → discourse
- For speech:
 - sample → spectral band → sound → ... phone → phoneme → word

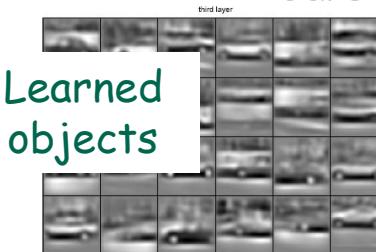
Eg: Learning Image Features

Examples of learned objects parts from object categories

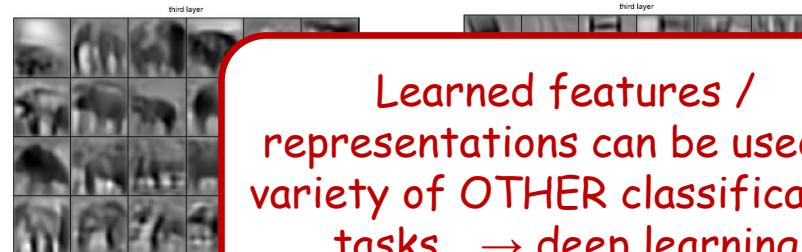
Faces



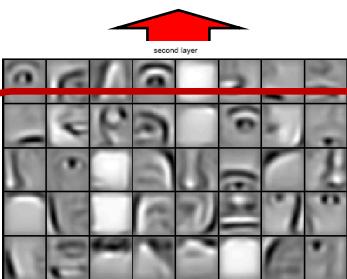
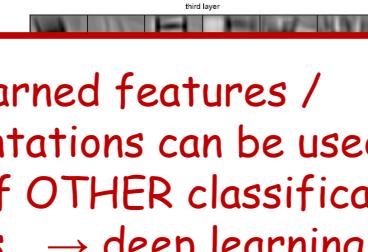
Cars



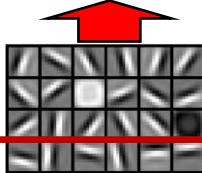
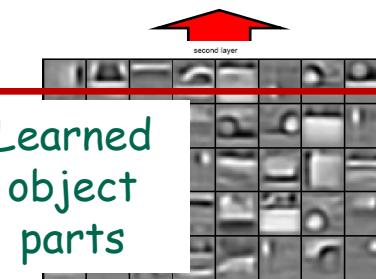
Elephants



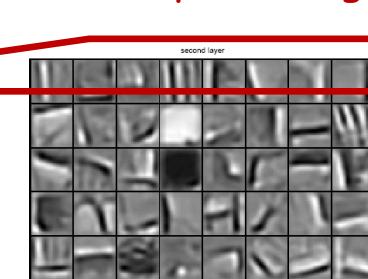
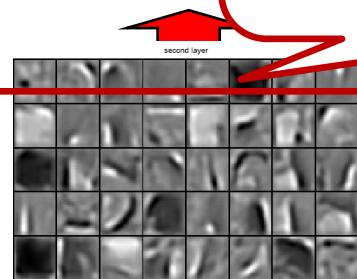
Chairs



Learned object parts



Learned edges



Actual images (pixels)

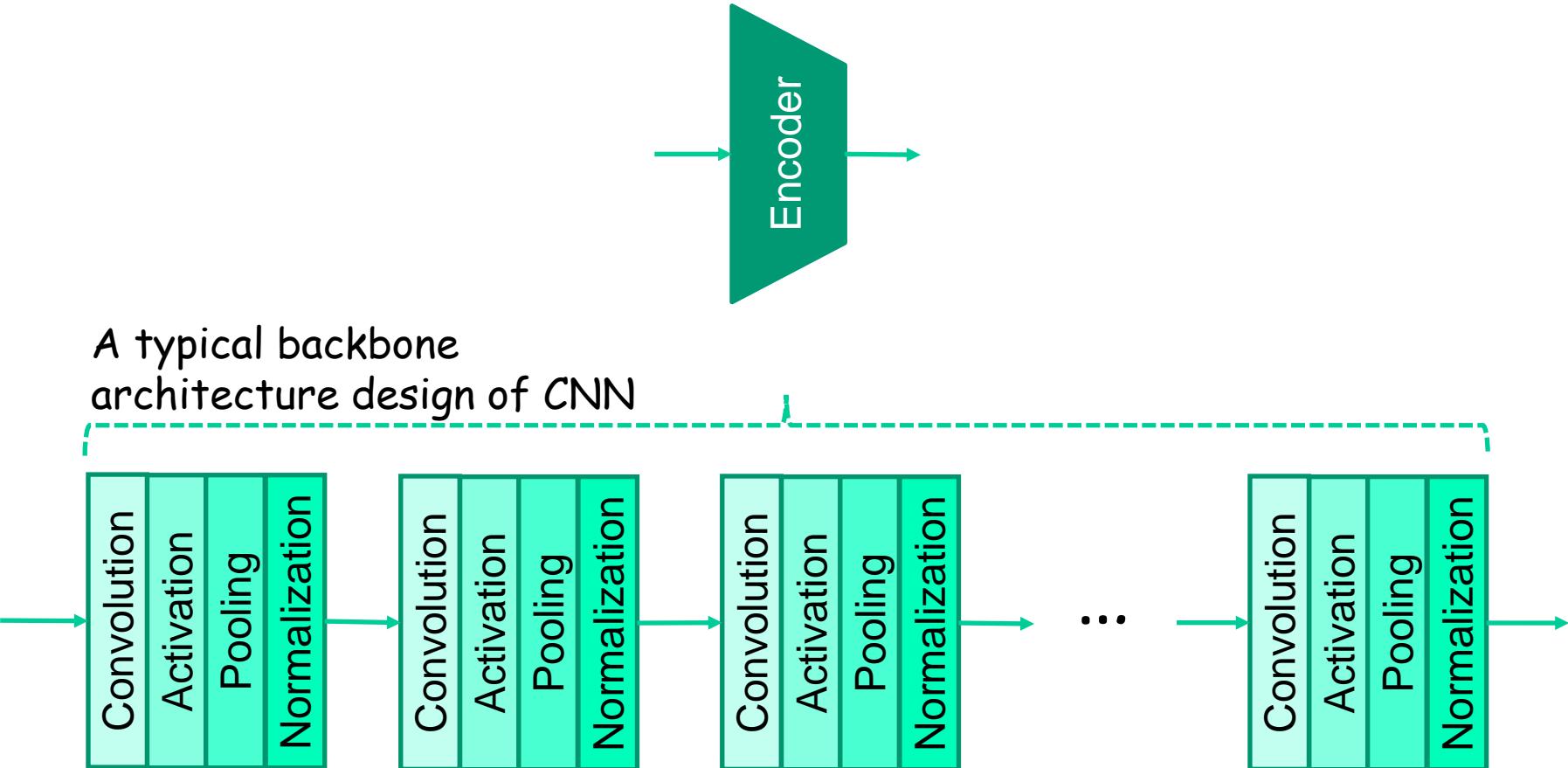


Today

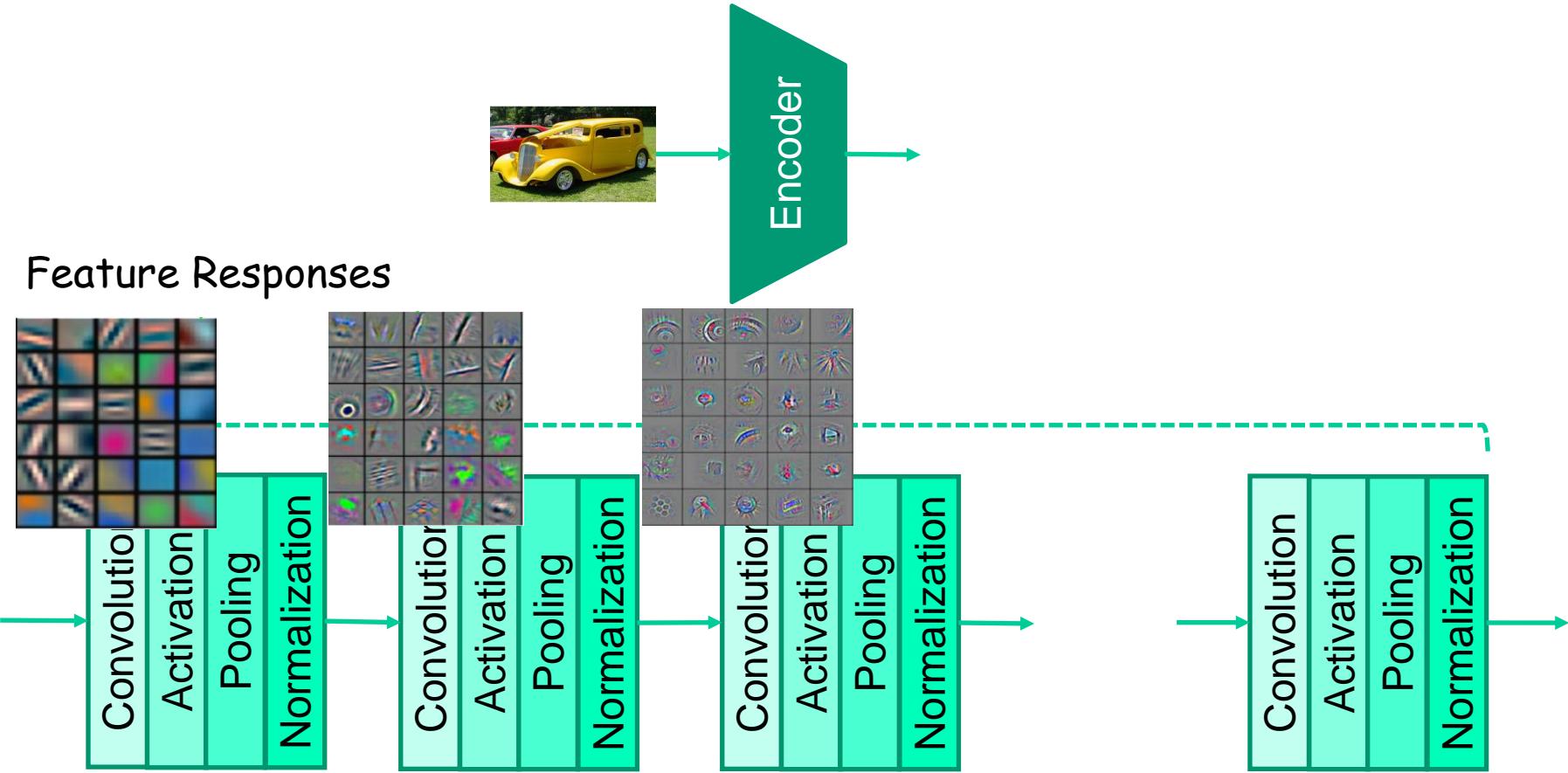
1. Motivation
2. Feature Learning
3. Building Block Design of CNN
 - Convolutional Layer
 - Activation Function
 - Pooling Layer
 - Normalization Layer



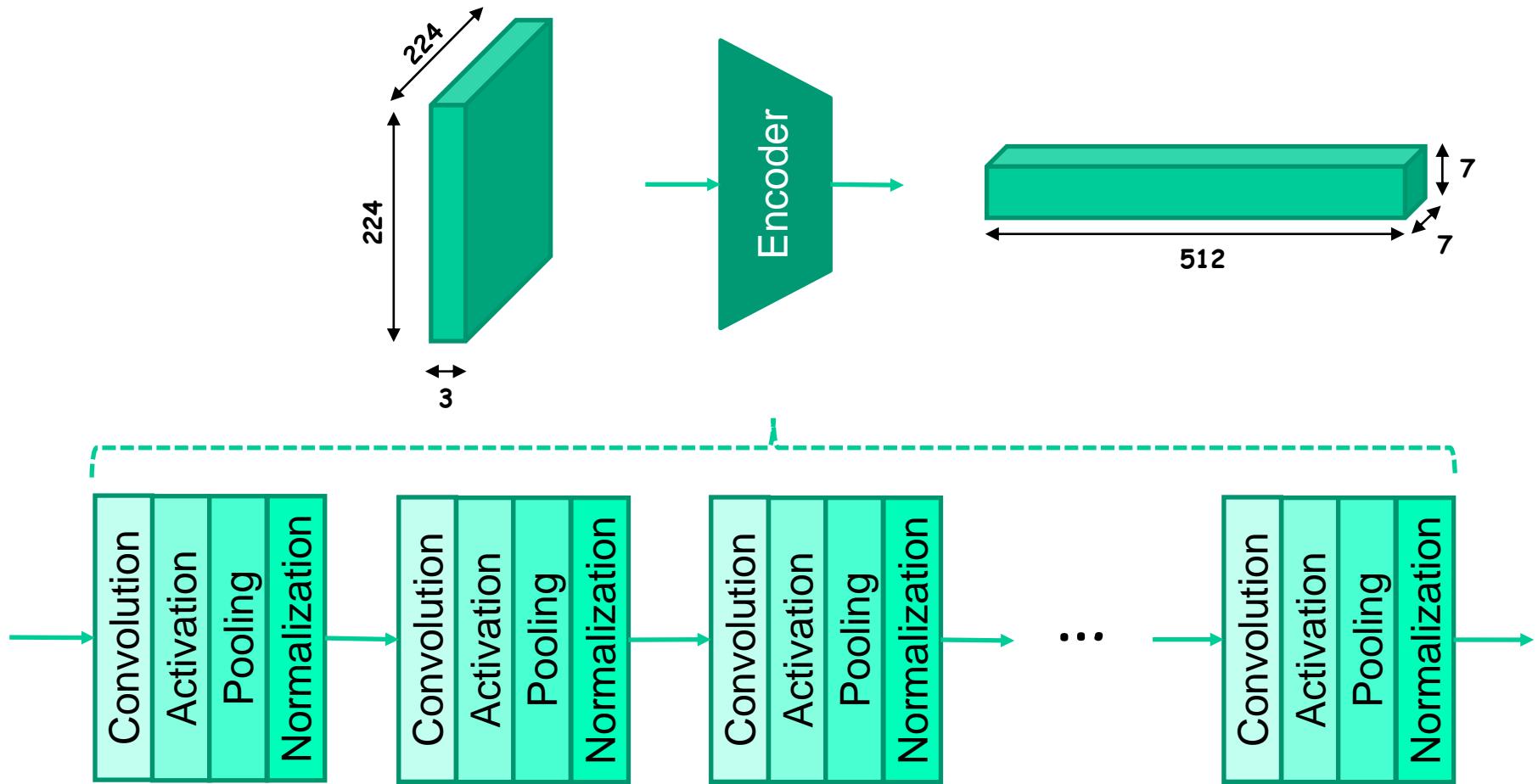
Building Block Design of CNN



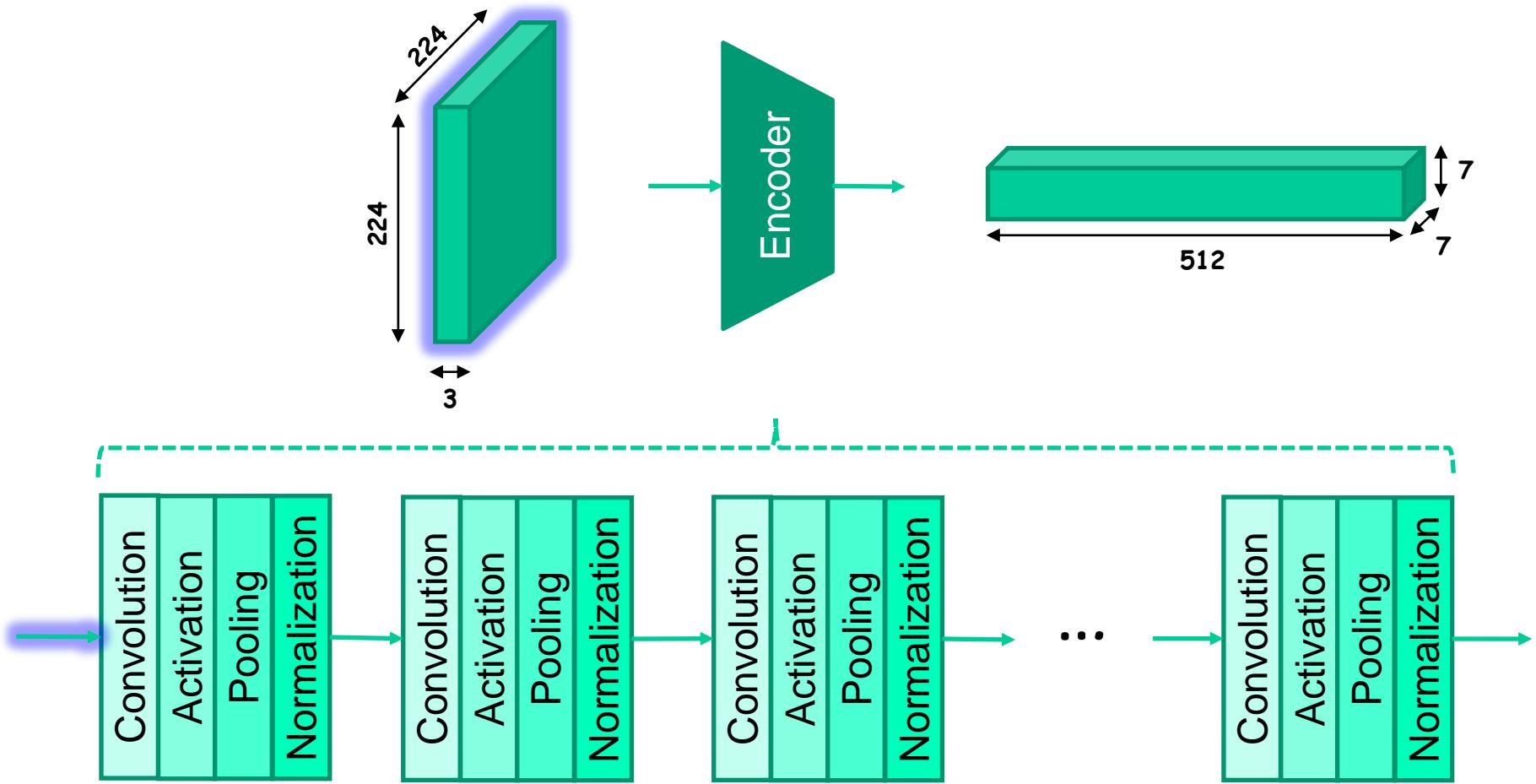
Building Block Design of CNN



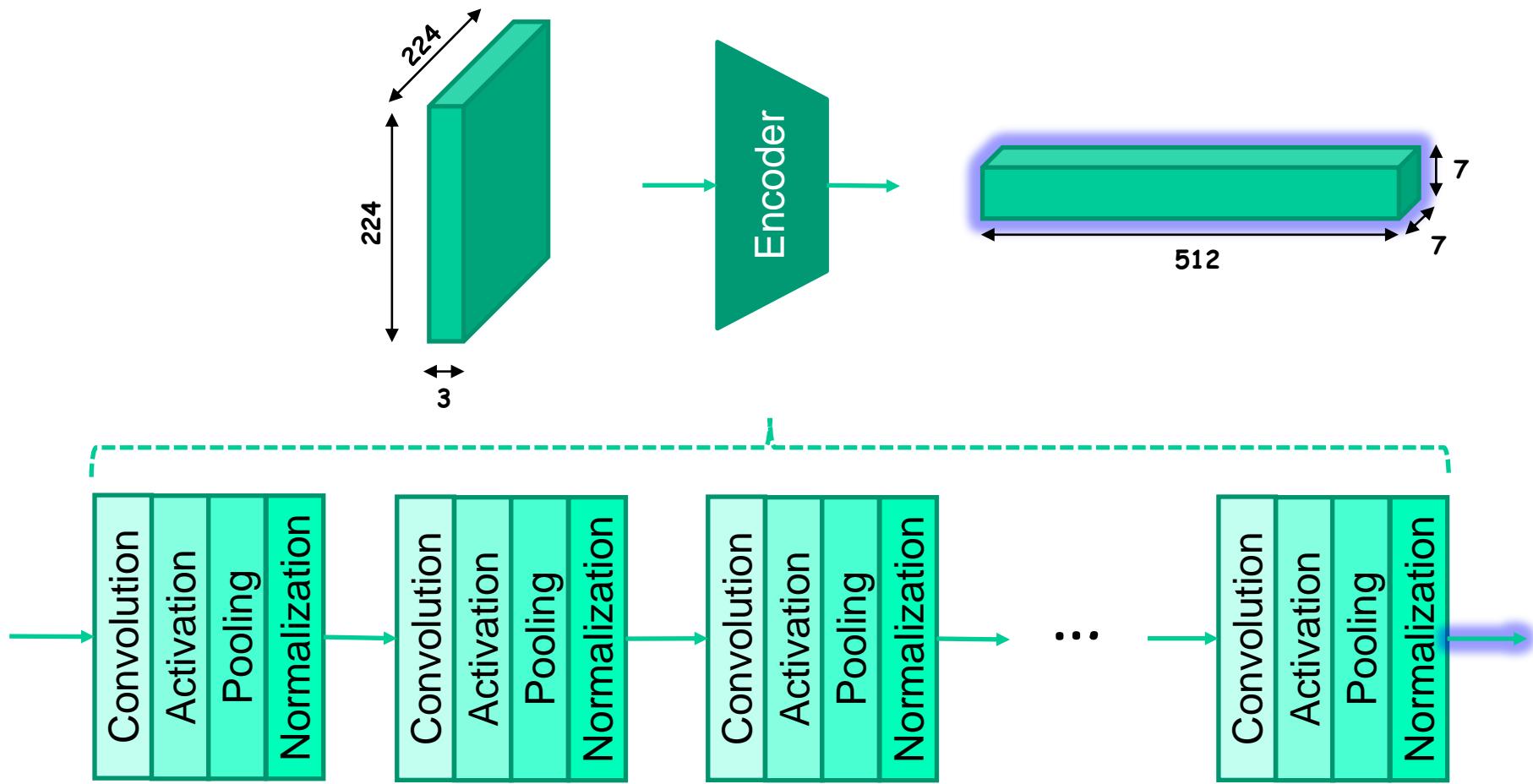
Building Block Design of CNN



Building Block Design of CNN



Building Block Design of CNN



Conv Layer of CNN

What computers 'see': Images as Numbers



Input Image

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	84	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	105	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	105	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	95	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

What you both see

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	106	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	105	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	95	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

Input Image + values

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	106	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	105	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	86	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	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

Pixel intensity values
("pix-el"=picture-element)

Levin Image Processing & Computer Vision

An image is just a matrix of numbers [0,255]. i.e., 1080x1080x3 for an RGB image.

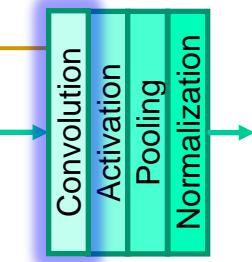
Question: is this Lincoln? Washington? Jefferson? Obama?

How can the computer answer this question?

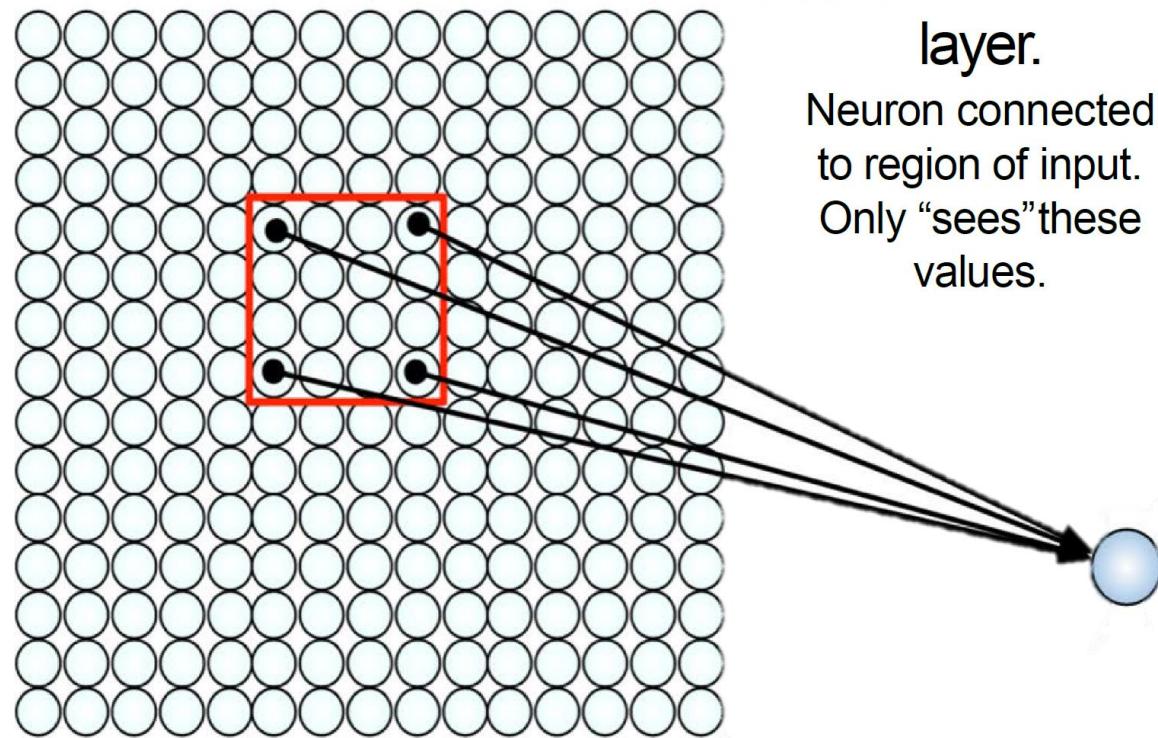
Can I just do classification on the 1,166400-long image vector directly?

No. Instead: exploit image spatial structure. Learn patches. Build them up

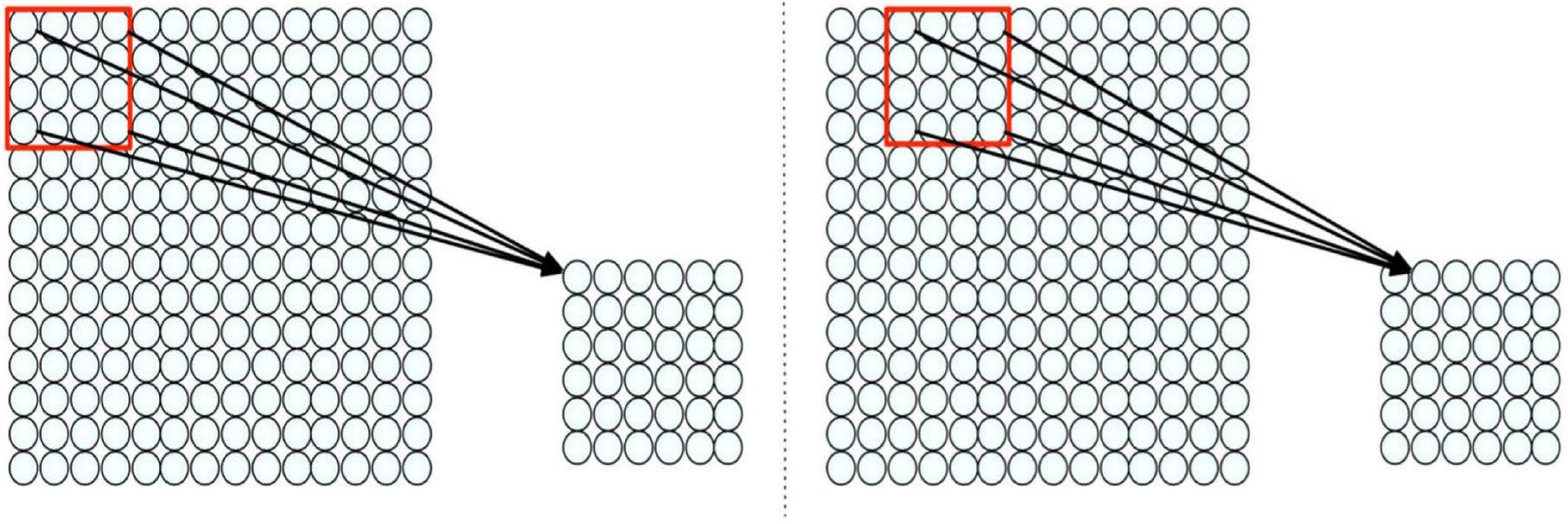
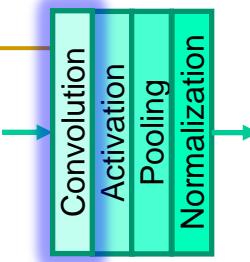
Using Spatial Structure



Input: 2D
image.
Array of pixel
values



Using Spatial Structure

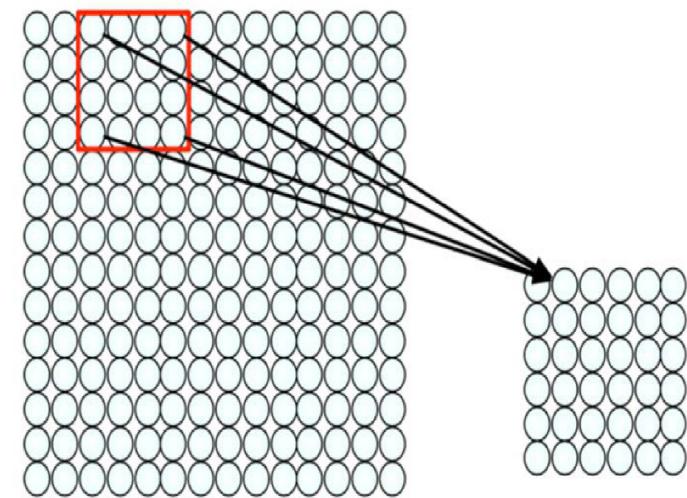
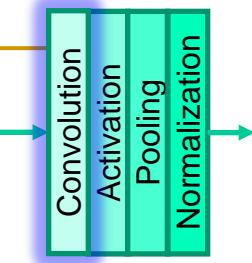


Connect patch in input layer to a single neuron in subsequent layer.

Use a sliding window to define connections.

*How can we **weight** the patch to detect particular features?*

Feature Extraction with Convolution

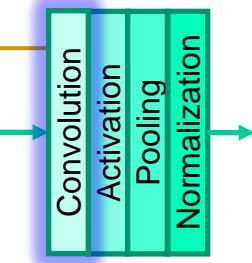


- Filter of size 4×4 : 16 different weights
- Apply this same filter to 4×4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

What is Convolution Operation?



- Elementwise multiplication and addition

Example 5x5 image with binary pixels

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

Example 3x3 filter

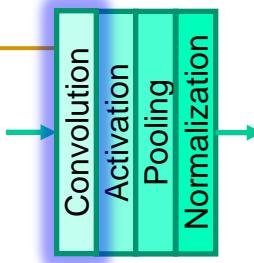
1	0	1
0	1	0
1	0	1

bias

0

- Scanning an image with a “filter”
 - Note: a filter is really just a perceptron, with weights and a bias

What is Convolution Operation?



- Elementwise multiplication and addition

1	0	1
0	1	0
1	0	1

Filter / Kernel

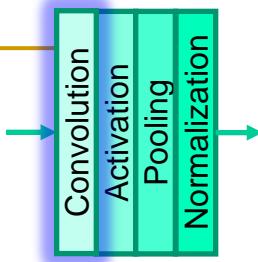
1 x1	1 x0	1 x1	0	0
0 x0	1 x1	1 x0	1	0
0 x1	0 x0	1 x1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

Sliding Window for Convolution



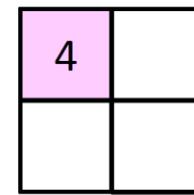
The “Stride” between adjacent scanned locations need not be 1

0
bias

1	0	1
0	1	0
1	0	1

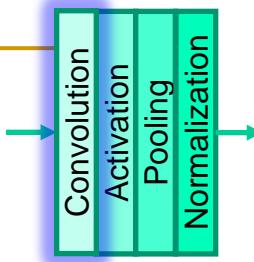
Filter

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



- Scanning an image with a “filter”
 - The filter may proceed by *more* than 1 pixel at a time
 - E.g. with a “stride” of two pixels per shift

Sliding Window for Convolution

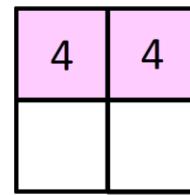


The “Stride” between adjacent scanned locations need not be 1

bias 0
Filter

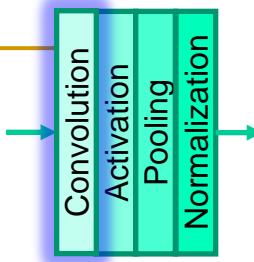
1	0	1
0	1	0
1	0	1

1	1	1 _{x1}	0 _{x0}	0 _{x1}
0	1	1 _{x0}	1 _{x1}	0 _{x0}
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1	1	0
0	1	1	0	0



- Scanning an image with a “filter”
 - The filter may proceed by *more* than 1 pixel at a time
 - E.g. with a “hop” of *two* pixels per shift

Sliding Window for Convolution



The “Stride” between adjacent scanned locations need not be 1

bias
 $\begin{bmatrix} 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$

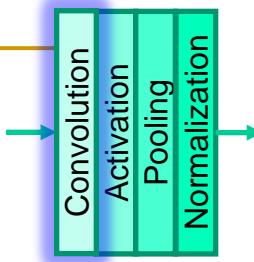
Filter

1	1	1	0	0
0	1	1	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0 _{x0}	0 _{x1}	1 _{x0}	1	0
0 _{x1}	1 _{x0}	1 _{x1}	0	0

4	4
2	

- Scanning an image with a “filter”
 - The filter may proceed by *more* than 1 pixel at a time
 - E.g. with a “hop” of *two* pixels per shift

Sliding Window for Convolution



The “Stride” between adjacent scanned locations need not be 1

bias $\begin{matrix} 0 \\ 1 \ 0 \ 1 \\ 0 \ 1 \ 0 \\ 1 \ 0 \ 1 \end{matrix}$

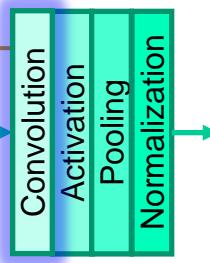
Filter

1	1	1	0	0
0	1	1	1	0
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1 _{x0}	1 _{x1}	0 _{x0}
0	1	1 _{x1}	0 _{x0}	0 _{x1}

4	4
2	4

- Scanning an image with a “filter”
 - The filter may proceed by *more* than 1 pixel at a time
 - E.g. with a “hop” of *two* pixels per shift

Output Feature Map size in Convolution



The size of the convolution

bias 0
 1 0 1
 0 1 0
 1 0 1

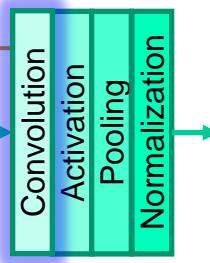
Filter

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



- Image size: 5x5
- Filter: 3x3
- Stride: 2
- Output size = ?

Output Feature Map size in Convolution



The size of the convolution

bias 0
 1 0 1
 0 1 0
 1 0 1

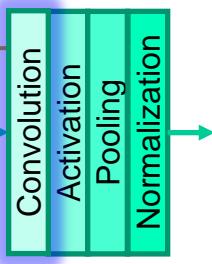
Filter

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

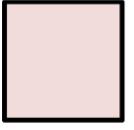
4	4
2	4

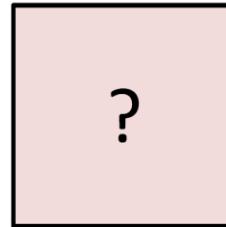
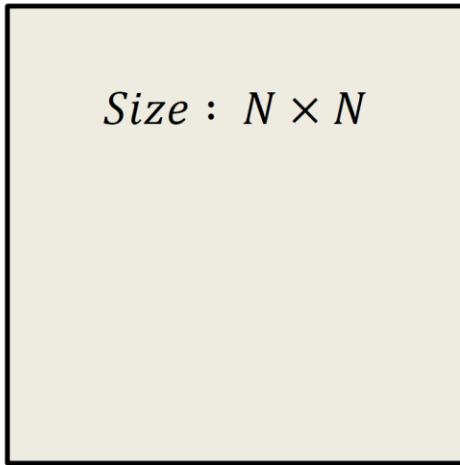
- Image size: 5x5
- Filter: 3x3
- Stride: 2
- Output size = ?

Output Feature Map size in Convolution



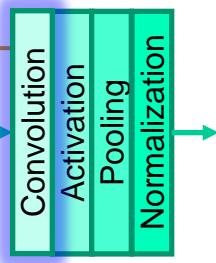
The size of the convolution

bias 0
 $M \times M$
 
Filter



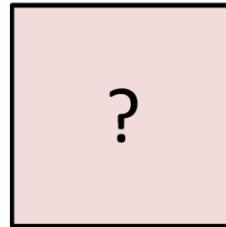
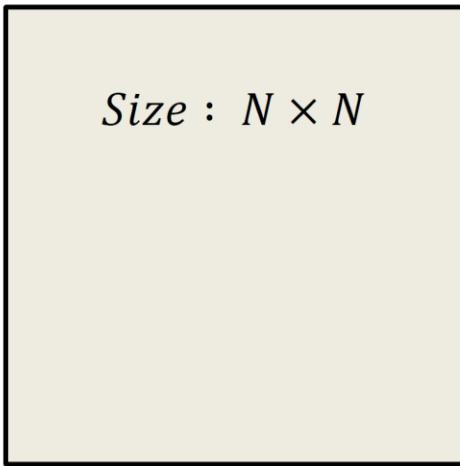
- Image size: $N \times N$
- Filter: $M \times M$
- Stride: 1
- Output size = ?

Output Feature Map size in Convolution



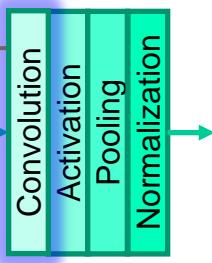
The size of the convolution

bias 0
 $M \times M$
 
Filter

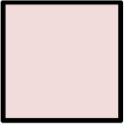


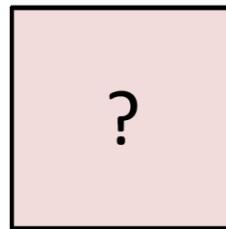
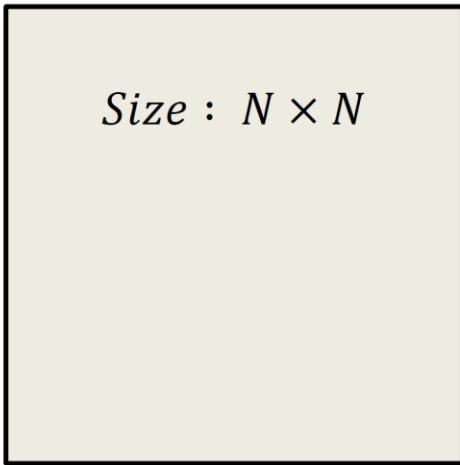
- Image size: $N \times N$
- Filter: $M \times M$
- Stride: S
- Output size = ?

Output Feature Map size in Convolution



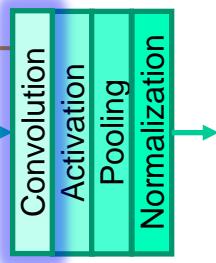
The size of the convolution

bias 0
 $M \times M$
 
Filter



- Image size: $N \times N$
- Filter: $M \times M$
- Stride: S
- Output size (each side) = $\lfloor (N - M)/S \rfloor + 1$
 - Assuming you're not allowed to go beyond the edge of the input

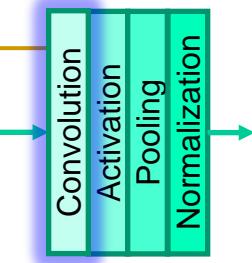
Output Feature Map size in Convolution



Convolution Size

- Simple convolution size pattern:
 - Image size: $N \times N$
 - Filter: $M \times M$
 - Stride: S
 - **Output size (each side)** = $\lfloor (N - M)/S \rfloor + 1$
 - Assuming you're not allowed to go beyond the edge of the input
- Results in a reduction in the output size
 - Even if $S = 1$
 - Sometimes not considered acceptable
 - If there's no active downsampling, through max pooling and/or $S > 1$, then the output map should ideally be the same size as the input

Padding for Convolution



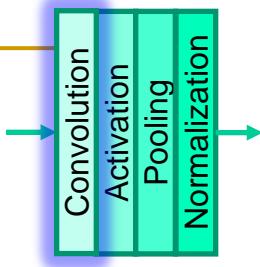
Solution

bias 
Filter

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

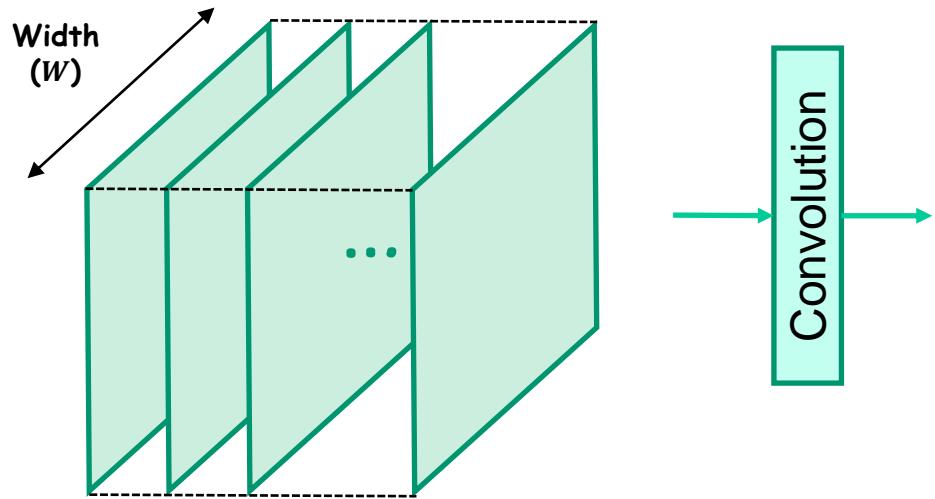
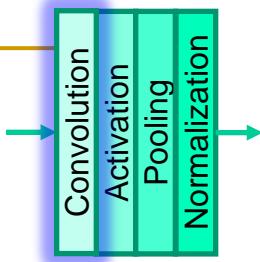
- Zero-pad the input
 - Pad the input image/map all around
 - Pad as symmetrically as possible, such that..
 - **For stride 1, the result of the convolution is the same size as the original image**

Conv Layer of CNN

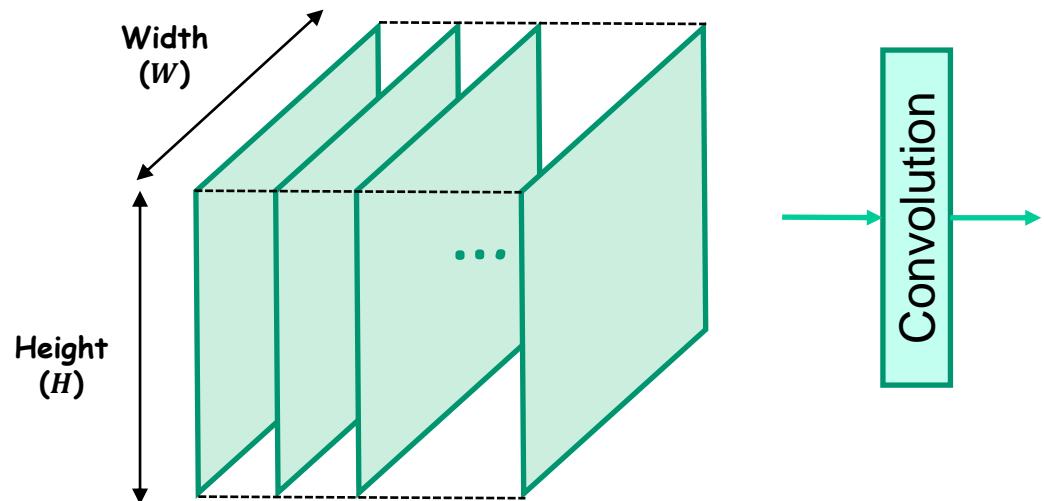
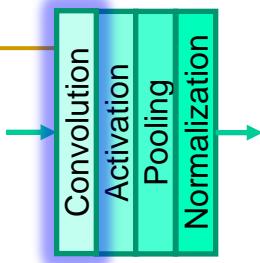


<https://cs231n.github.io/convolutional-networks/>

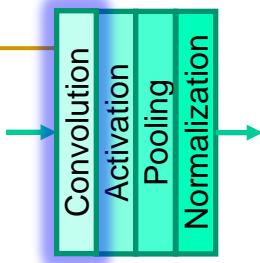
Convolutional Feature Mapping



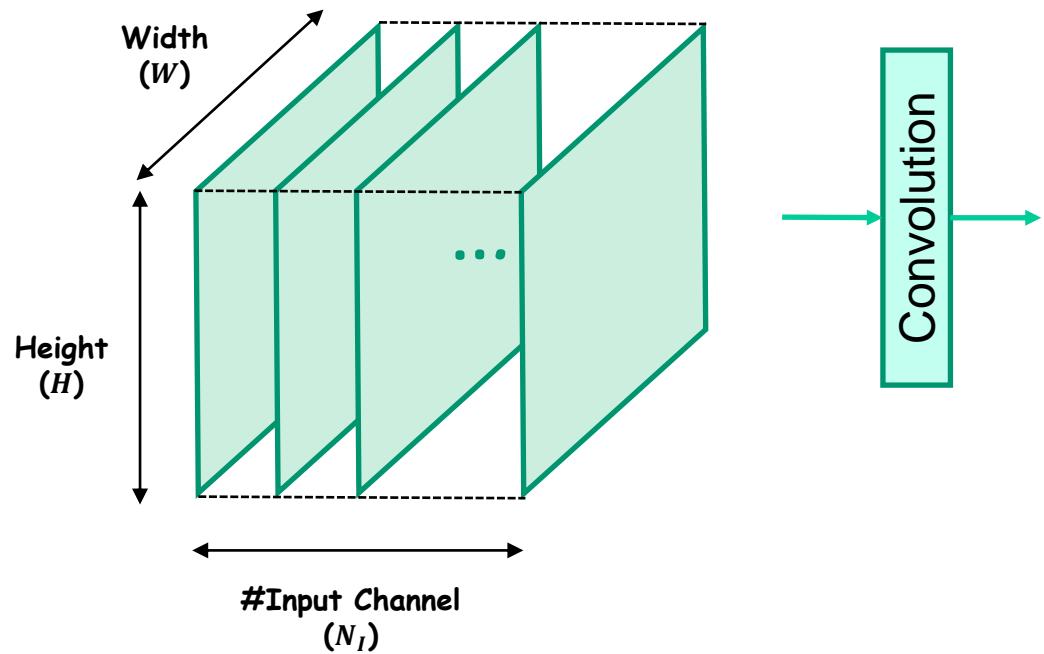
Convolutional Feature Mapping



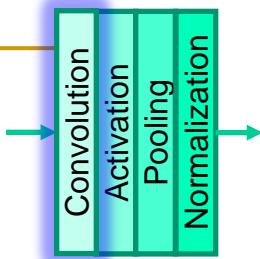
Convolutional Feature Mapping



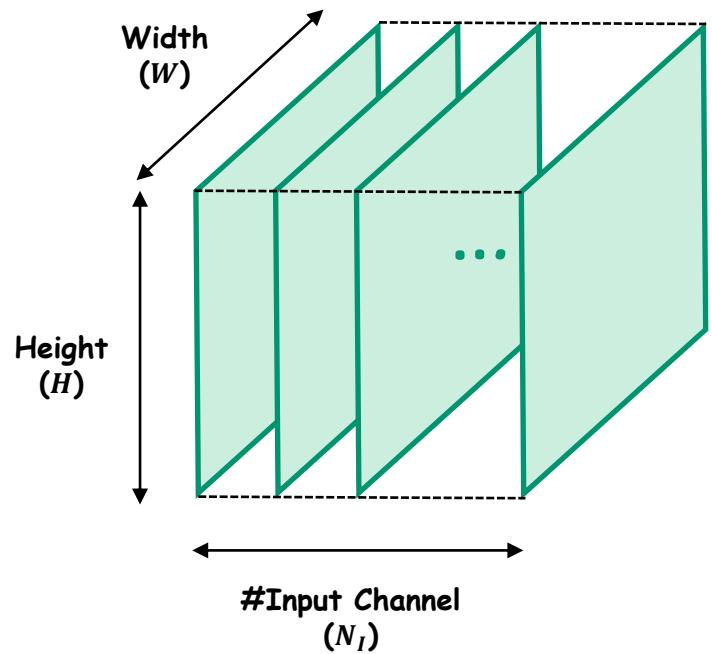
$$F^I \in \mathbb{R}^{H \times W \times N_I}$$



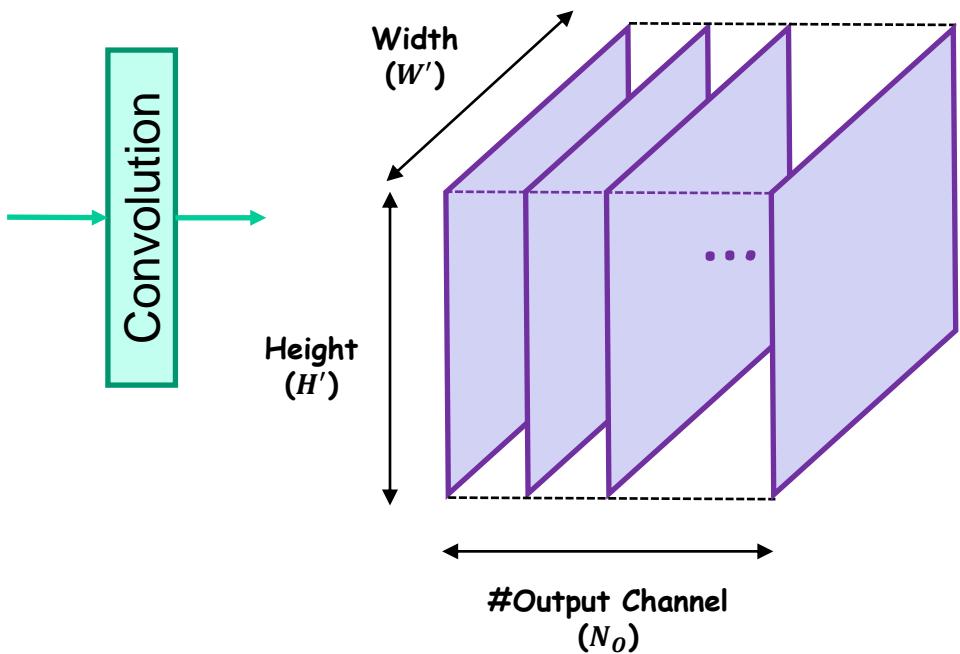
Convolutional Feature Mapping



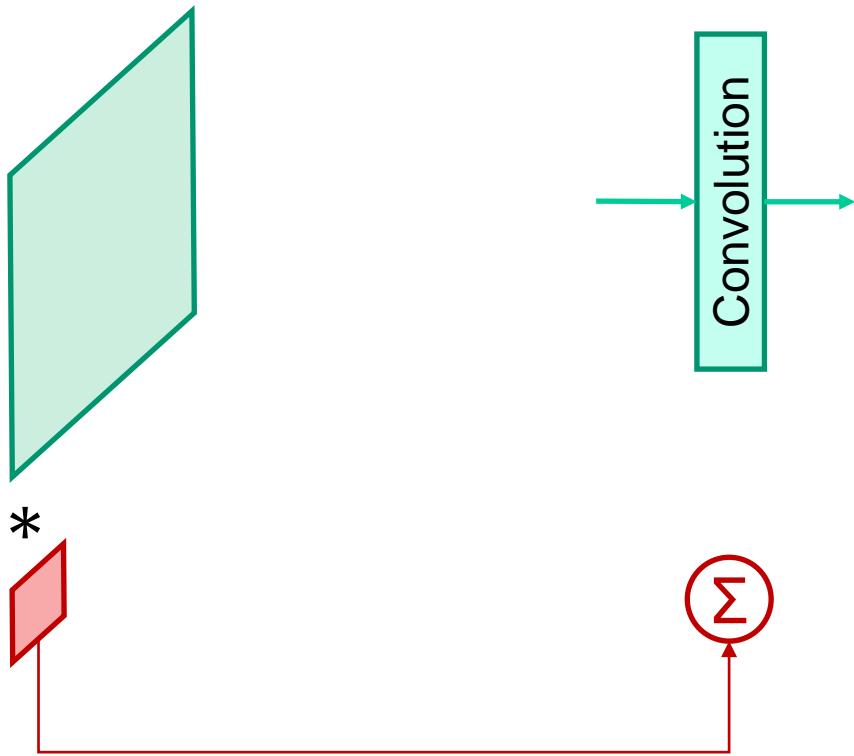
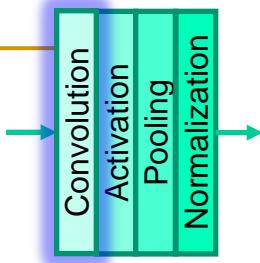
$$F^I \in \mathbb{R}^{H \times W \times N_I}$$



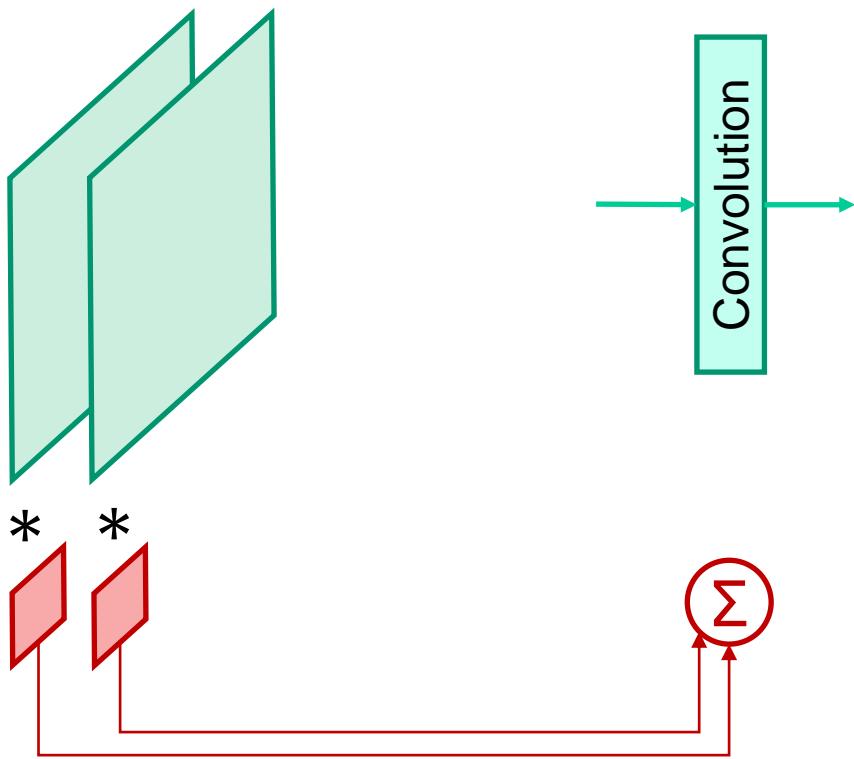
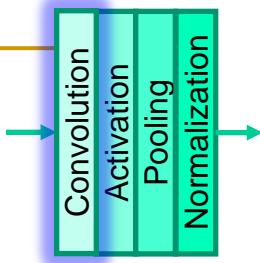
$$F^O \in \mathbb{R}^{H' \times W' \times N_O}$$



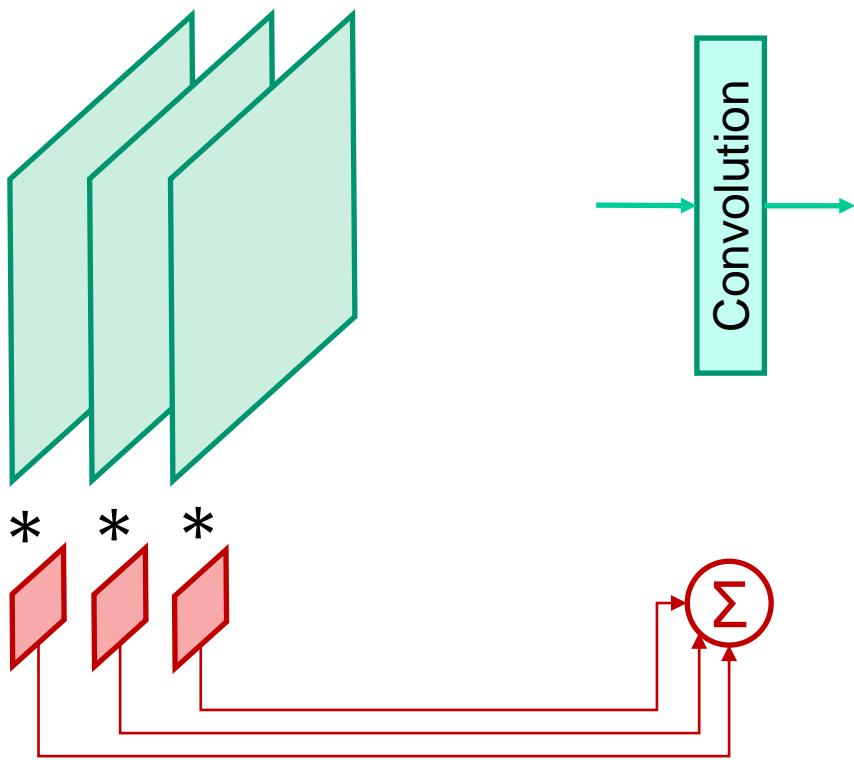
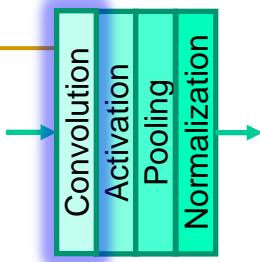
Convolutional Feature Mapping



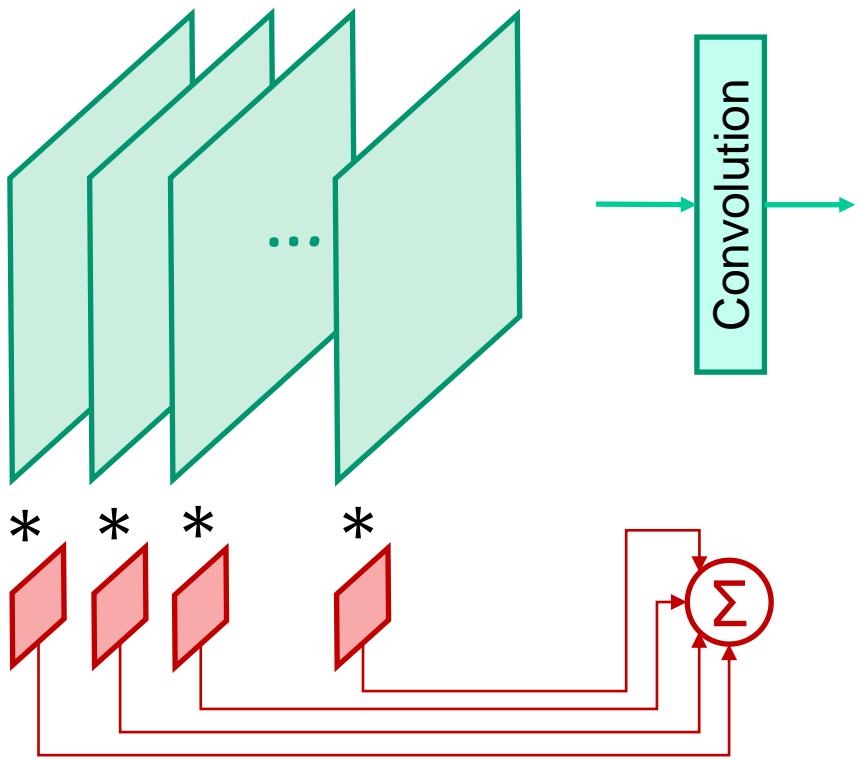
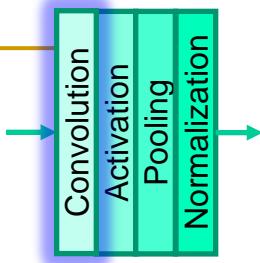
Convolutional Feature Mapping



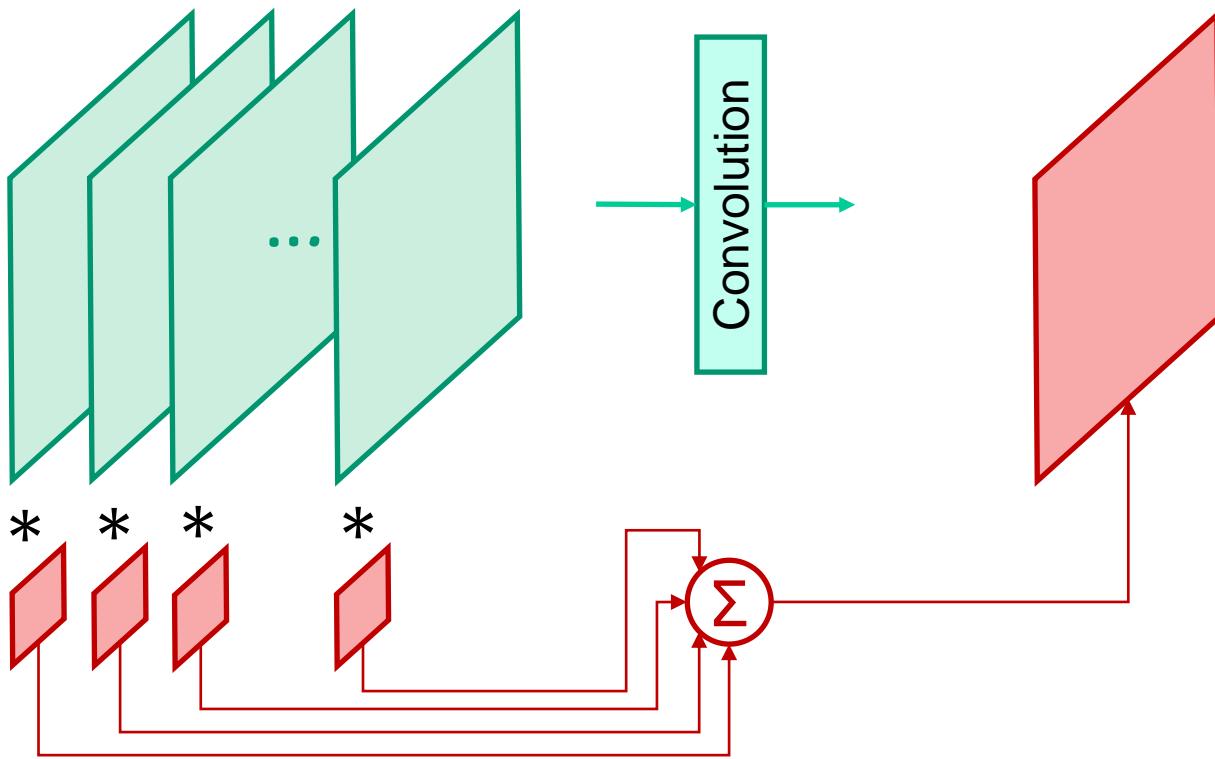
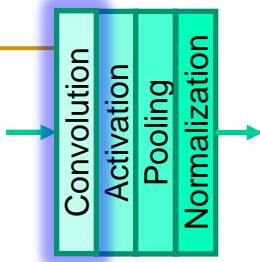
Convolutional Feature Mapping



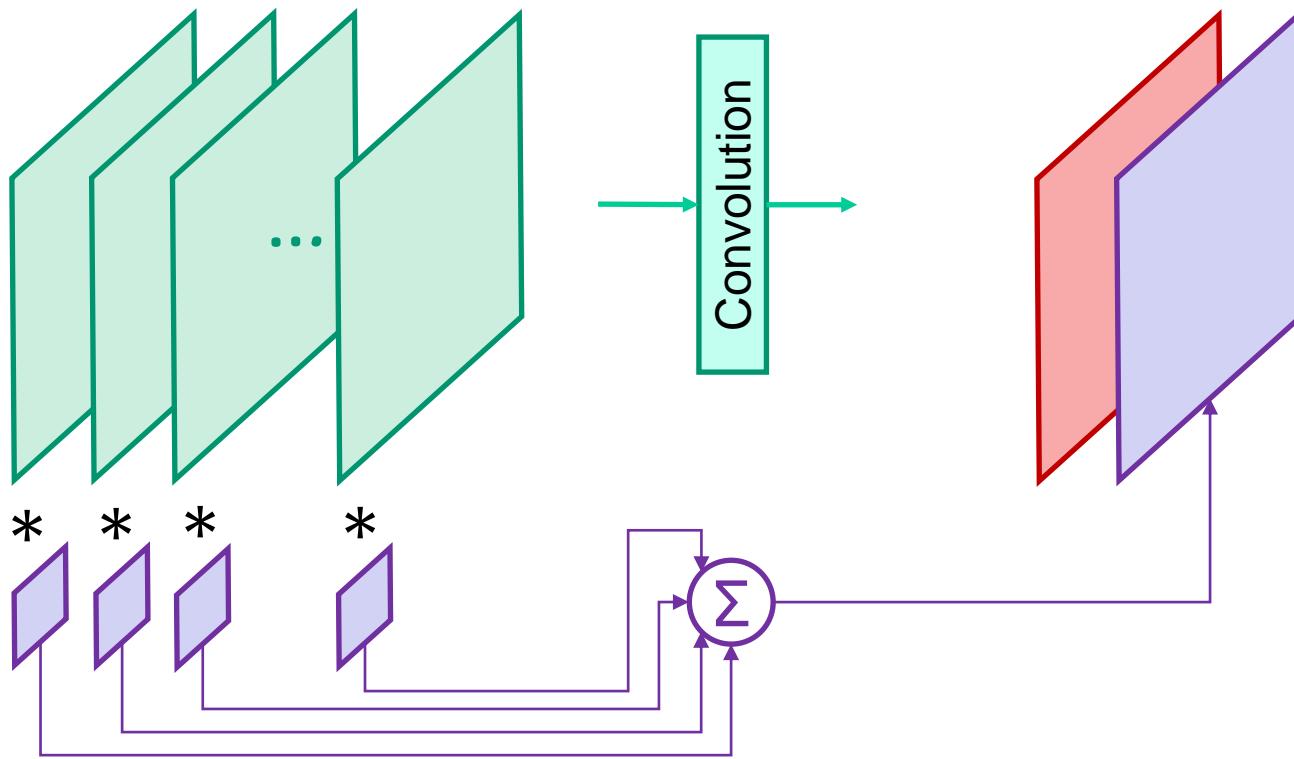
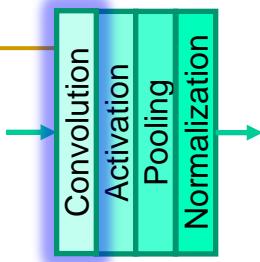
Convolutional Feature Mapping



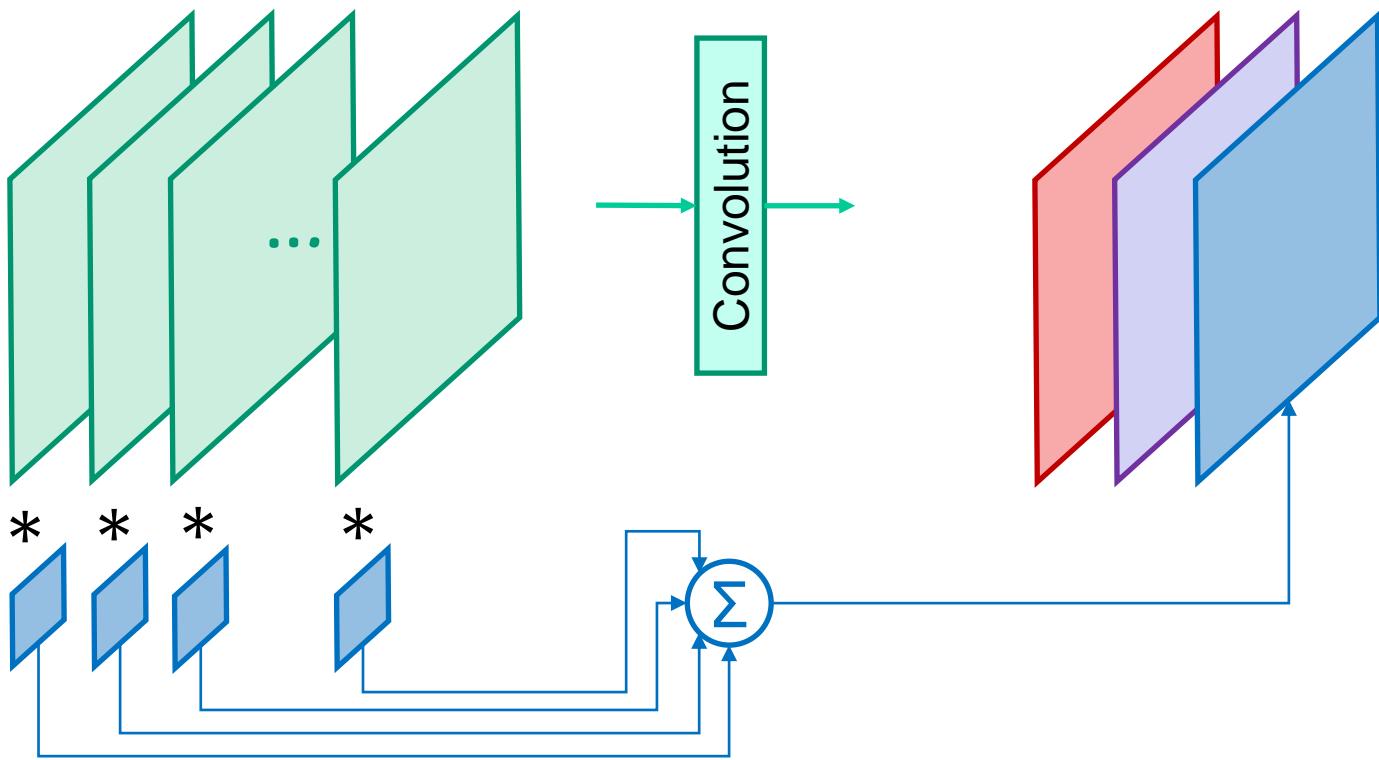
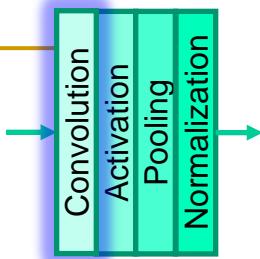
Convolutional Feature Mapping



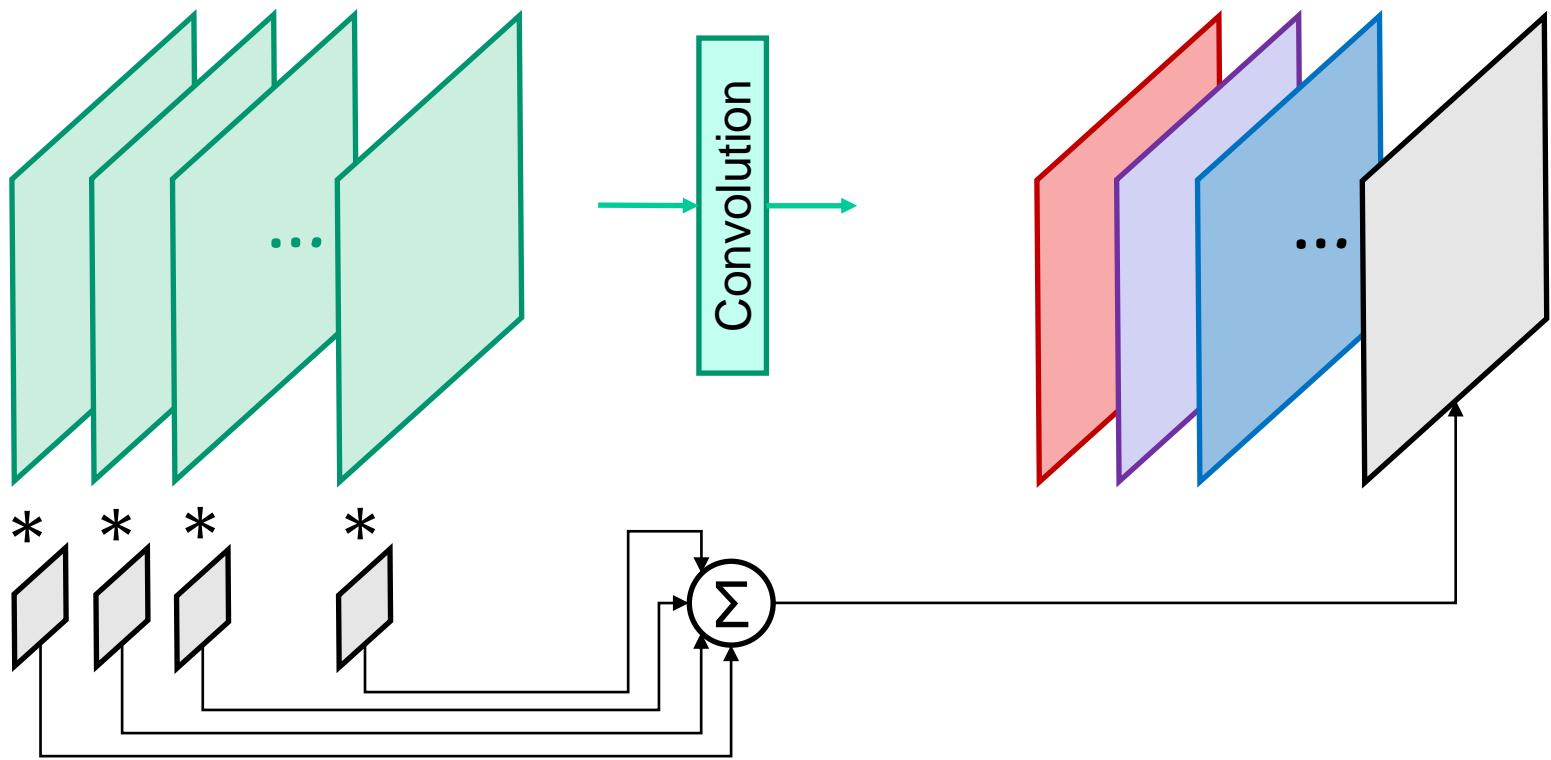
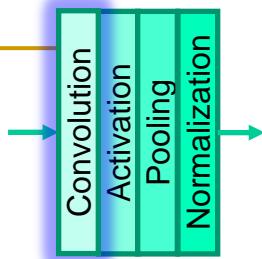
Convolutional Feature Mapping



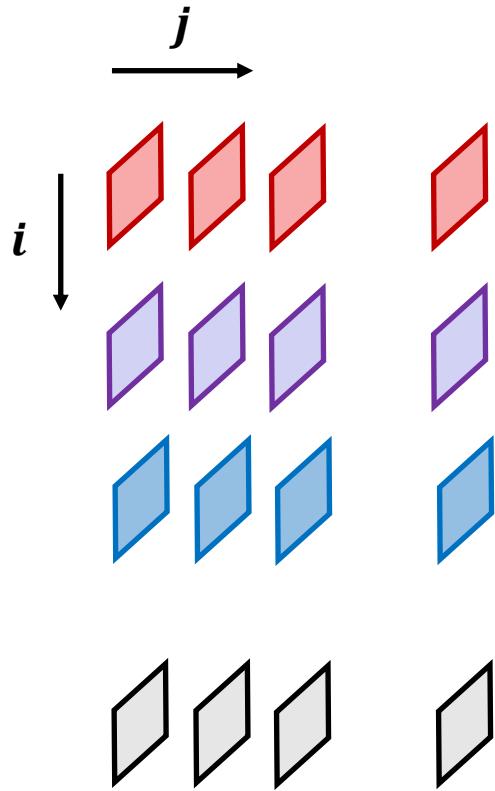
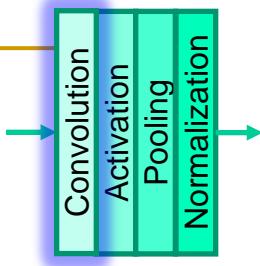
Convolutional Feature Mapping



Convolutional Feature Mapping



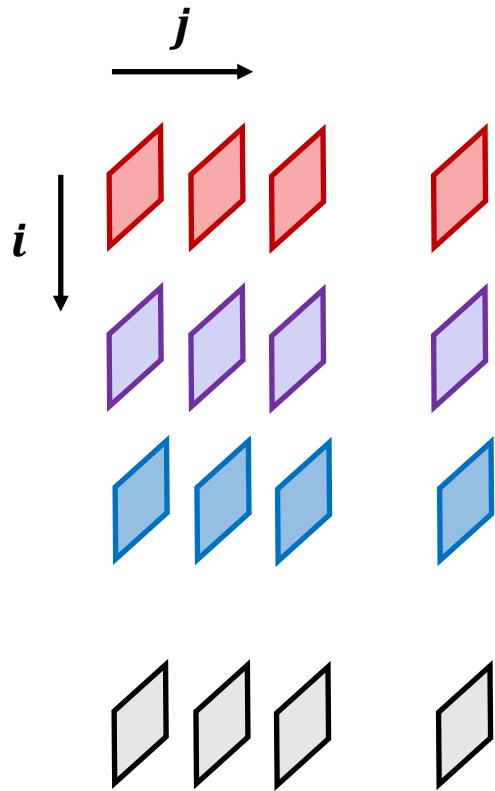
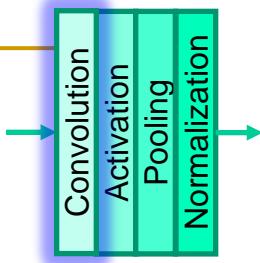
Convolutional Feature Mapping



Convolutional filters (aka kernels) can be represented in Tensor format

$$K \in \mathbb{R}^{h \times w \times N_I \times N_O}$$

Convolutional Feature Mapping

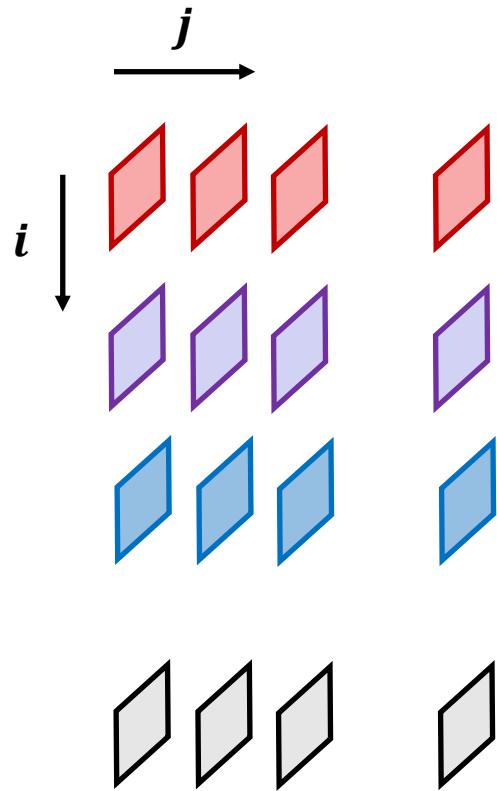
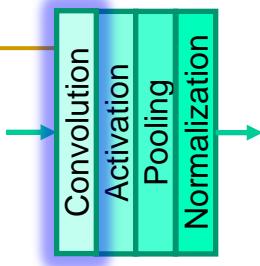


Convolutional filters (aka kernels) can be represented in Tensor format

$$K \in \mathbb{R}^{h \times w \times N_I \times N_O}$$

Ex1. $3 \times 3 \times 64 \times 128$

Convolutional Feature Mapping



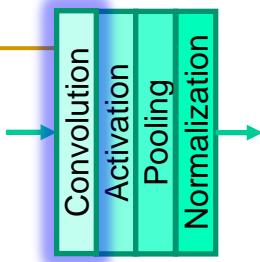
Convolutional filters (aka kernels) can be represented in Tensor format

$$K \in \mathbb{R}^{h \times w \times N_I \times N_O}$$

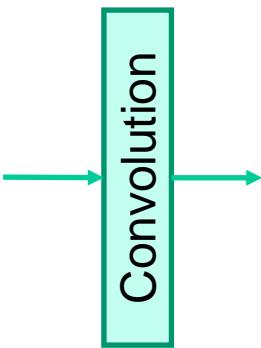
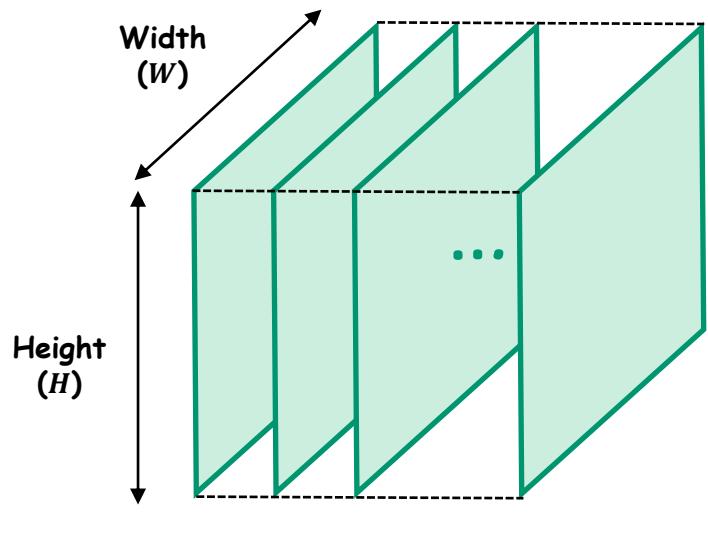
Ex1. $3 \times 3 \times 64 \times 128$

Ex2. $5 \times 5 \times 64 \times 128$

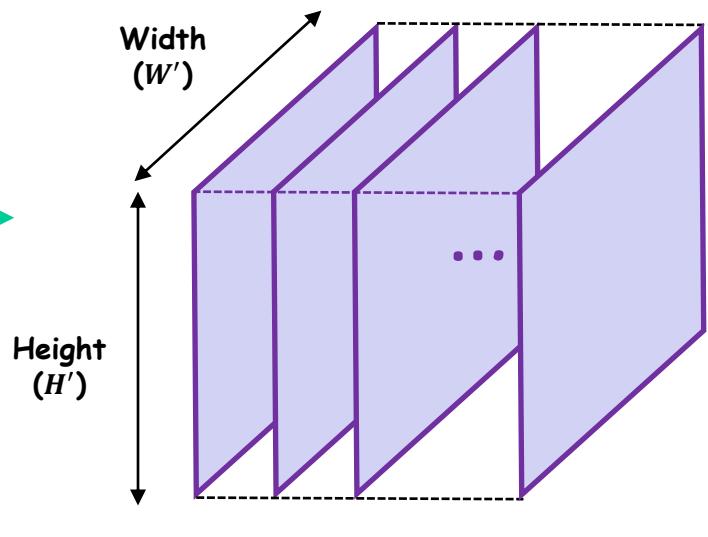
Convolutional Feature Mapping



$$F^I \in \mathbb{R}^{H \times W \times N_I}$$



$$F^O \in \mathbb{R}^{H' \times W' \times N_O}$$

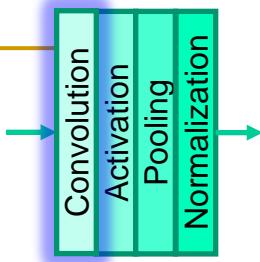


$$\#Input Channel (N_I)$$

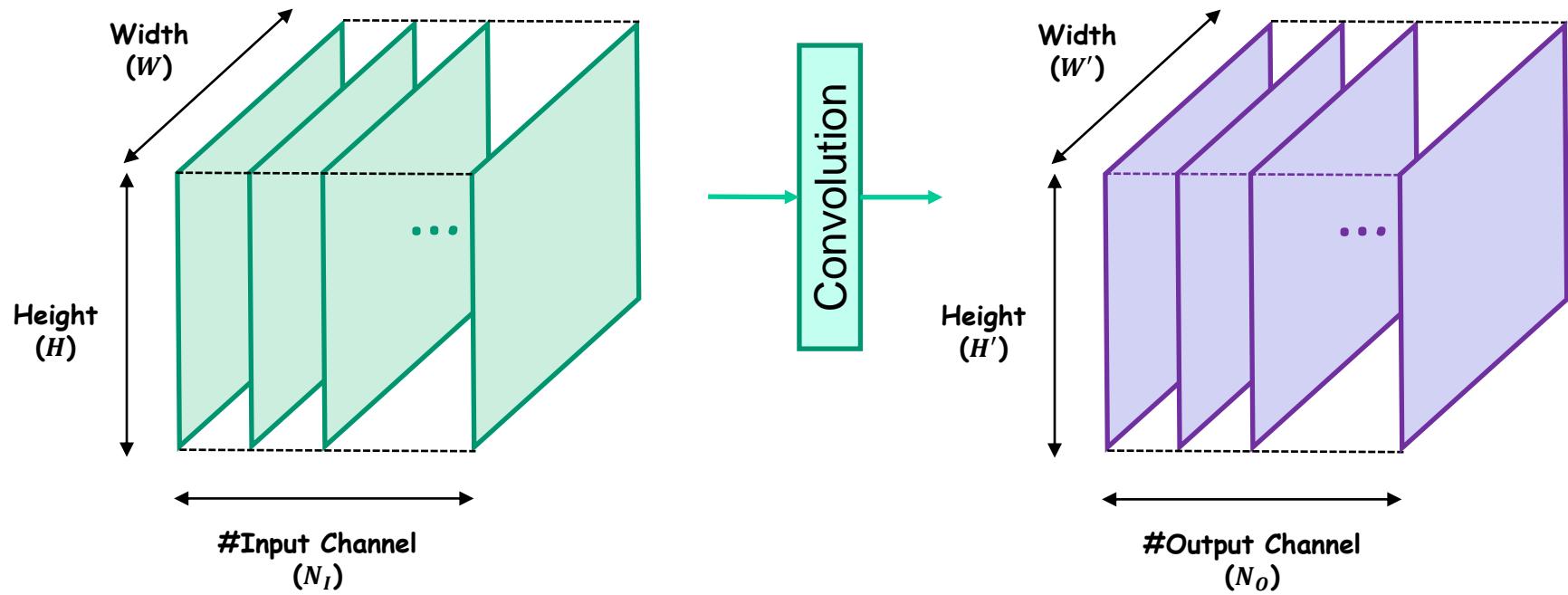
$$\#Output Channel (N_O)$$

$$F^O(:, :, j) = \left[\sum_{i=1}^{N_I} F^I(:, :, i) * K(:, :, i, j) \right] + b_j$$

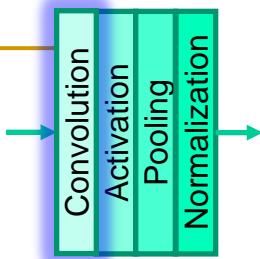
Convolutional Feature Mapping



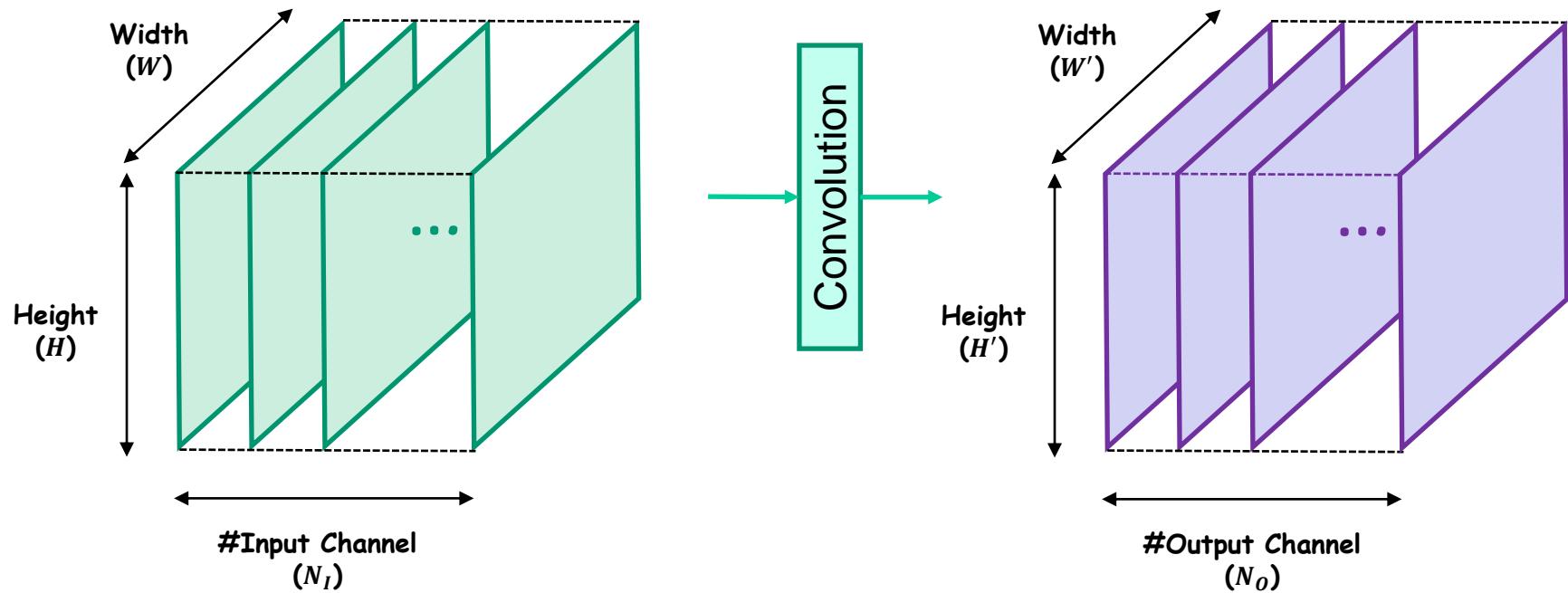
Ex1. Input Feature Map: $56 \times 56 \times 64 \rightarrow$
Output Feature Map: $56 \times 56 \times 128$ (Stride=1)



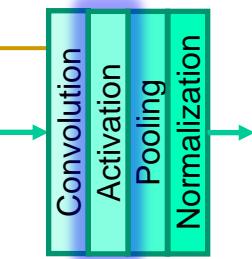
Convolutional Feature Mapping



Ex1. Input Feature Map: $56 \times 56 \times 64 \rightarrow$
Output Feature Map: $28 \times 28 \times 128$ (Stride=2)



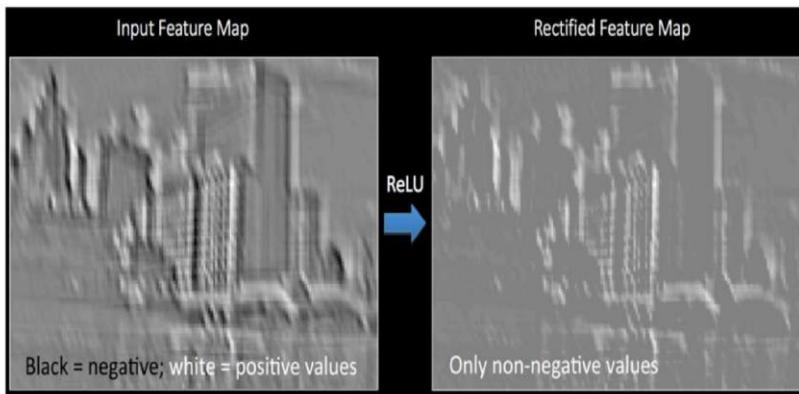
Non-Linear Activation



Introducing Non-Linearity

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero.
- **Non-linear operation**

Rectified Linear Unit
(ReLU)

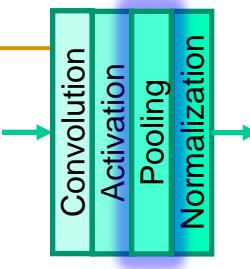


$$g(z) = \max(0, z)$$

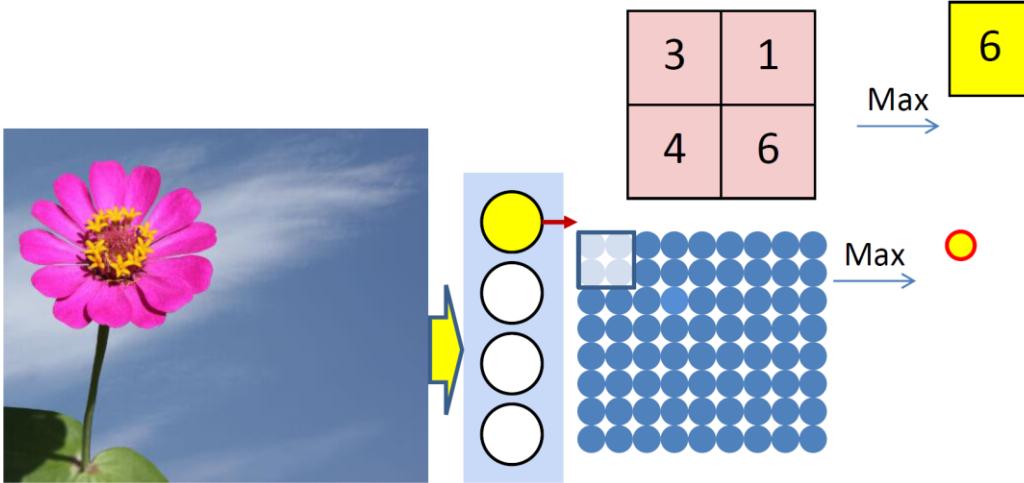
CLASS `torch.nn.ReLU(inplace=False)` [\[SOURCE\]](#)

Karn Intuitive CNNs

Pooling (i.e. down-sampling)

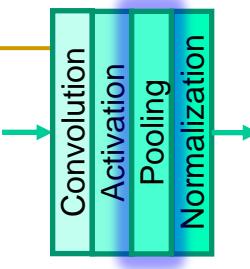


Recall: Max pooling

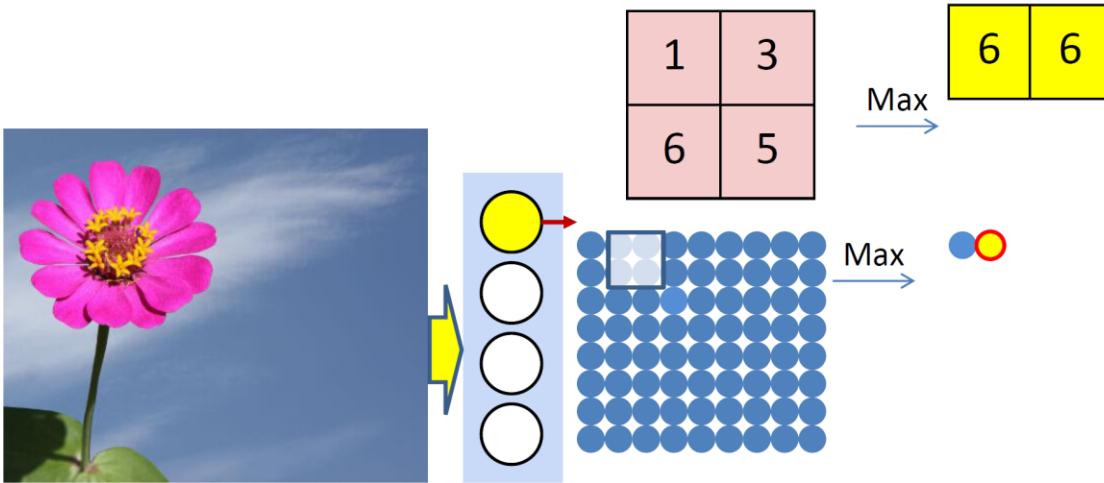


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

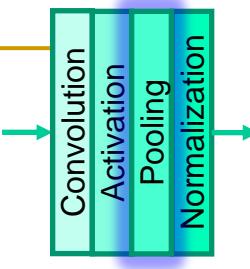


Recall: Max pooling

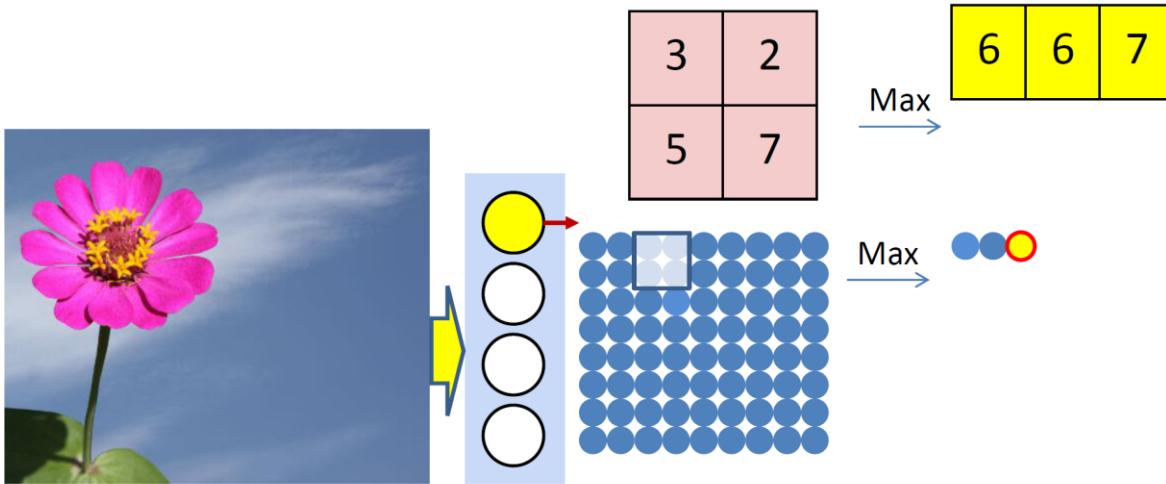


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

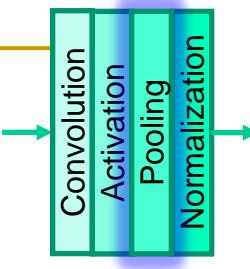


Recall: Max pooling

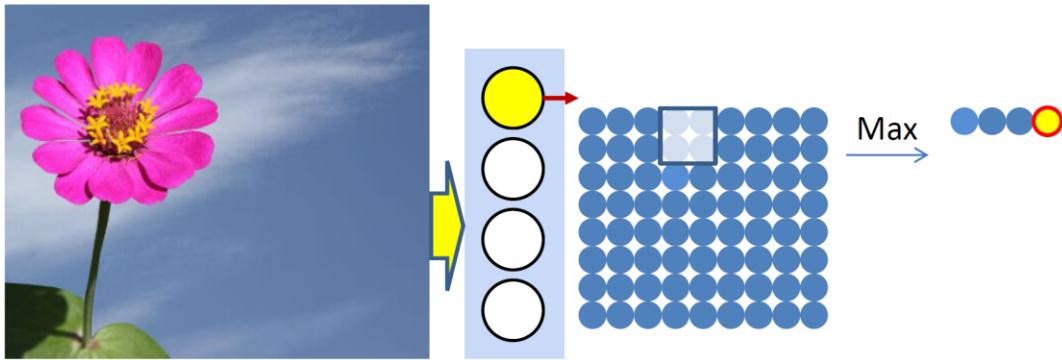


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

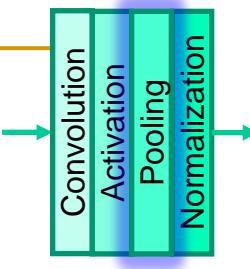


Recall: Max pooling

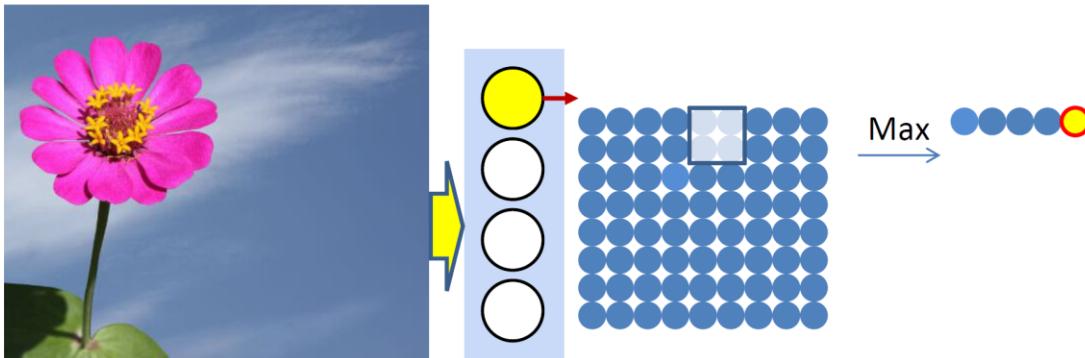


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

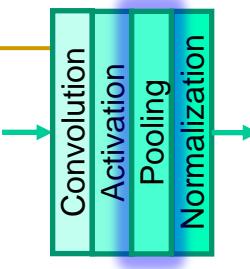


Recall: Max pooling

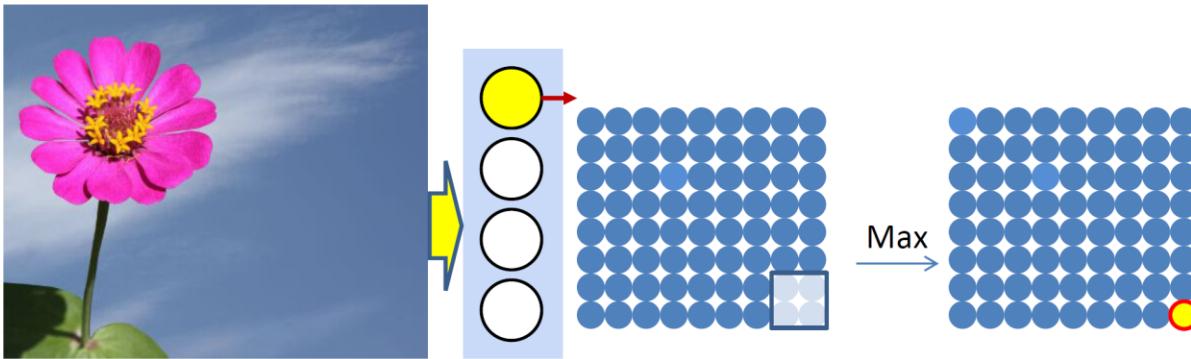


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

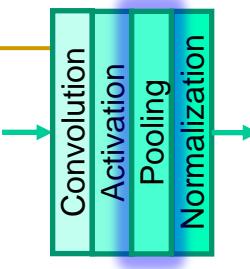


Recall: Max pooling

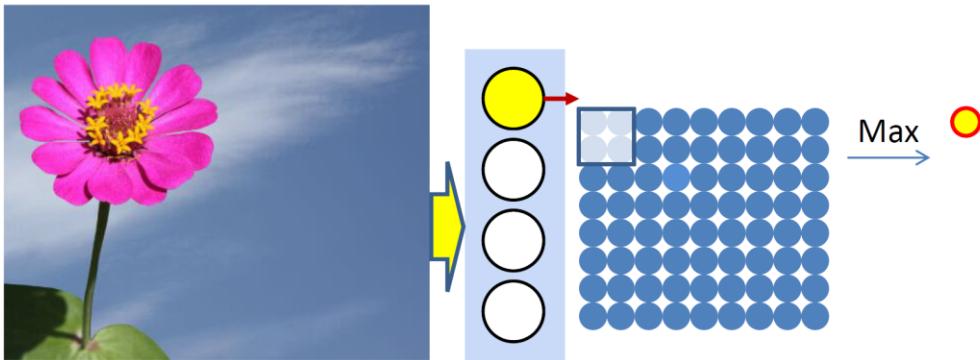


- Max pooling selects the largest from a pool of elements
- Pooling is performed by “scanning” the input

Pooling (i.e. down-sampling)

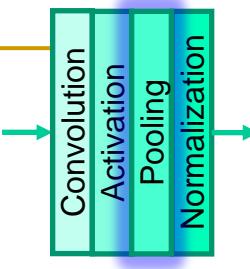


“Strides”

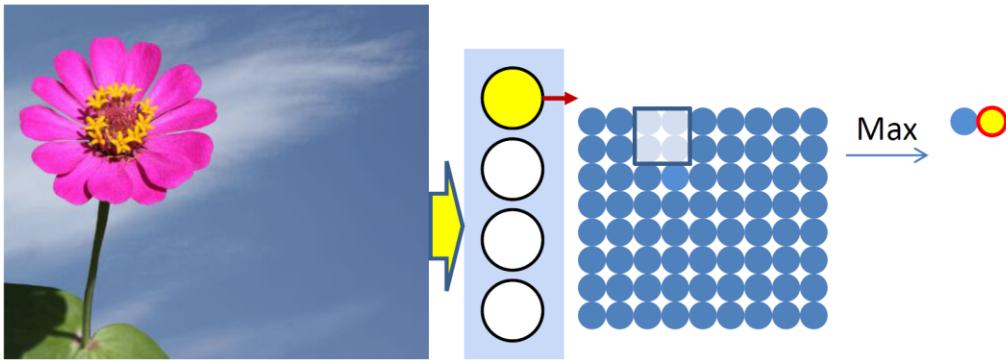


- The “max” operations may “stride” by more than one pixel

Pooling (i.e. down-sampling)

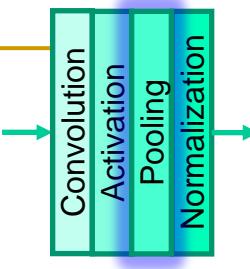


“Strides”

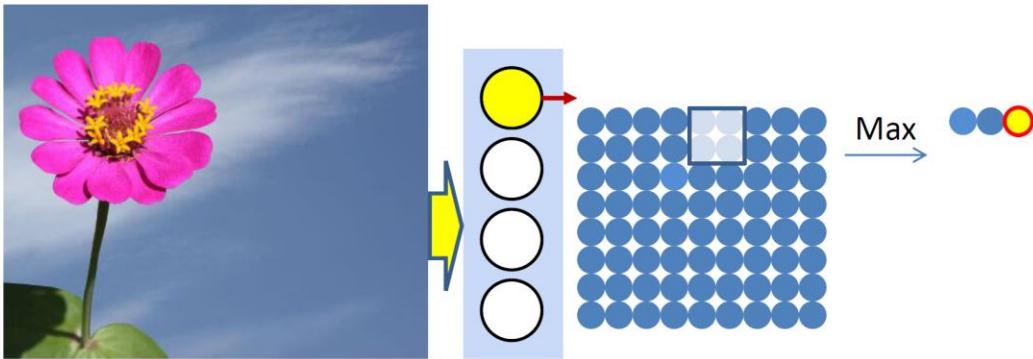


- The “max” operations may “stride” by more than one pixel

Pooling (i.e. down-sampling)

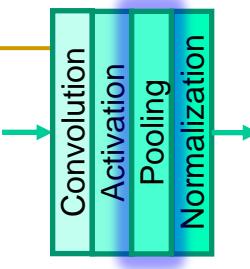


“Strides”

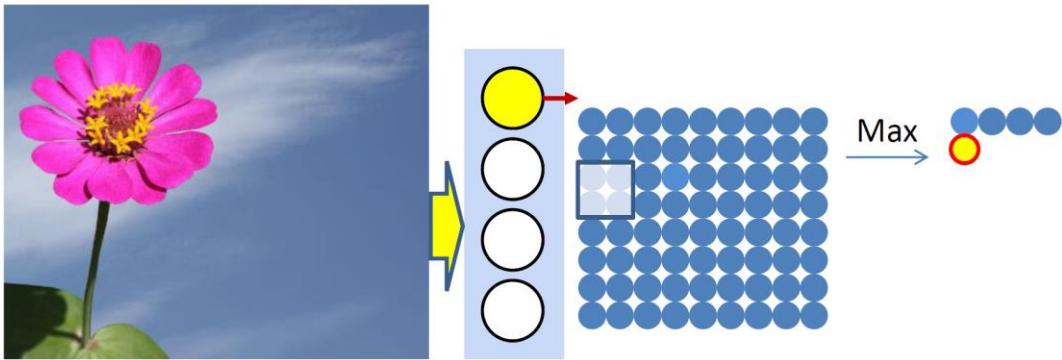


- The “max” operations may “stride” by more than one pixel

Pooling (i.e. down-sampling)

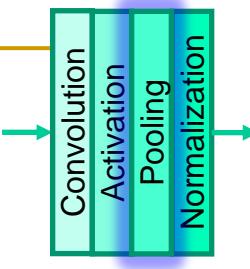


“Strides”

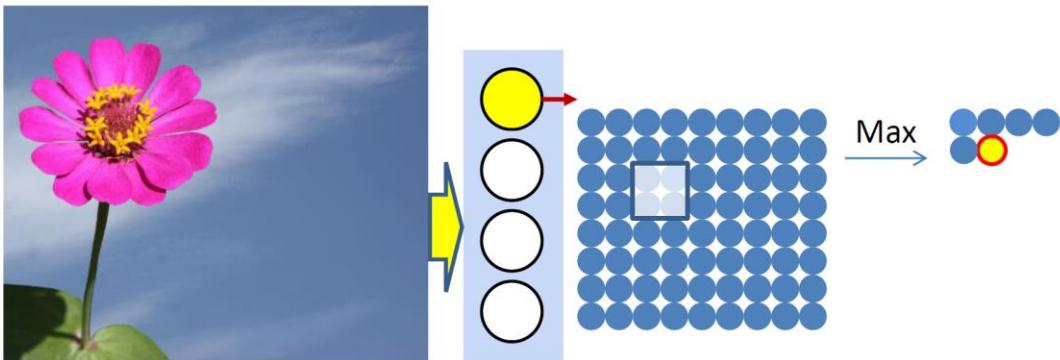


- The “max” operations may “stride” by more than one pixel

Pooling (i.e. down-sampling)

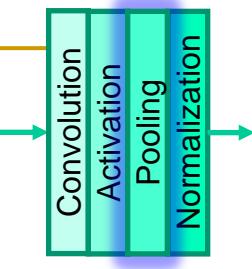


“Strides”

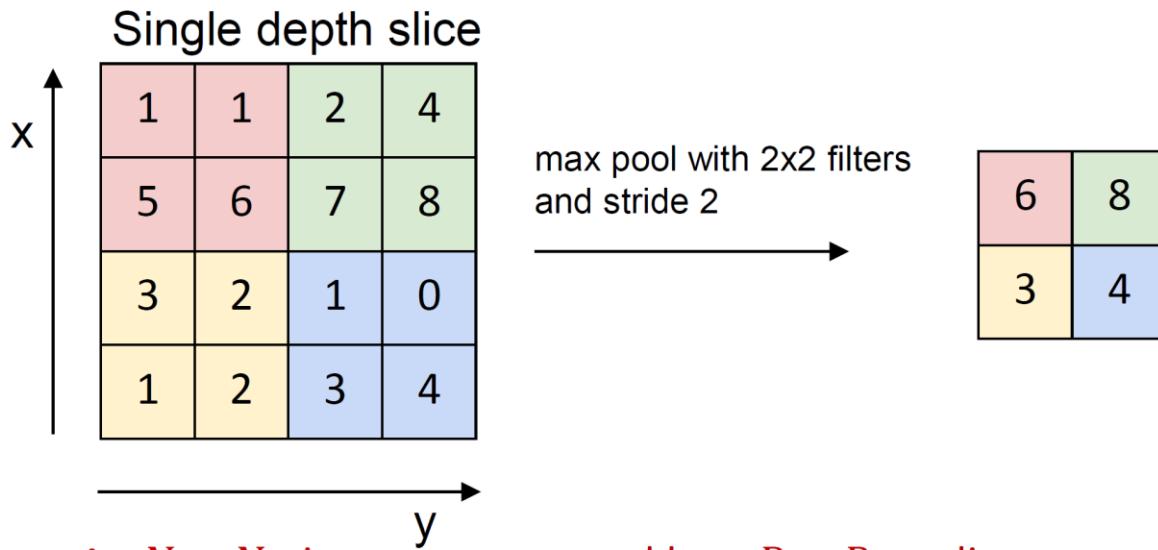


- The “max” operations may “stride” by more than one pixel

Pooling (i.e. down-sampling)

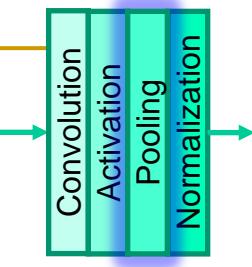


Pooling: Size of output

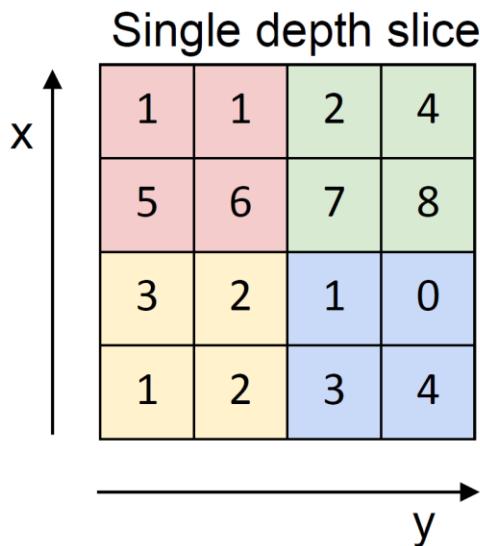


- An $N \times N$ picture compressed by a $P \times P$ pooling filter with stride D results in an output map of side $\lceil (N -$

Pooling (i.e. down-sampling)



Alternative to Max pooling: Mean Pooling

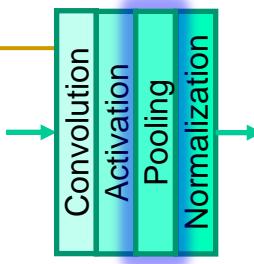


Mean pool with 2x2 filters and stride 2

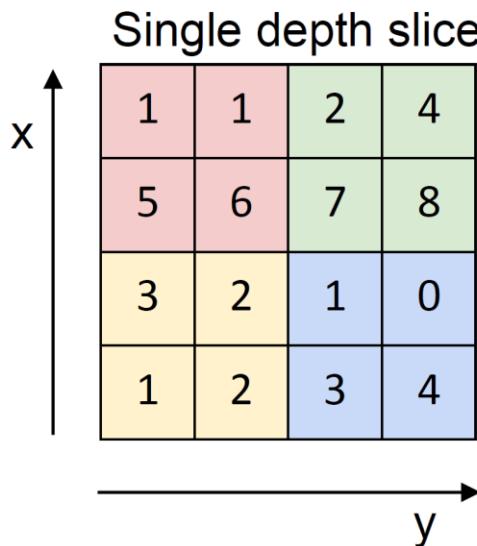
3.25	5.25
2	2

- Compute the mean of the pool, instead of the max

Pooling (i.e. down-sampling)



Alternative to Max pooling: P-norm



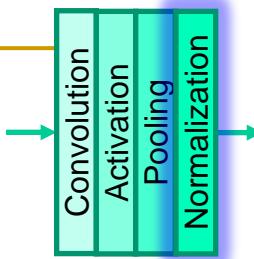
P-norm with 2x2 filters
and stride 2, $p = 5$

$$y = \sqrt[p]{\frac{1}{P^2} \sum_{i,j} x_{ij}^p}$$

4.86	8
2.38	3.16

- Compute a p-norm of the pool

Batch Normalization



- BN is basically Whitening Transformation
- i.e. makes distribution close to standard normal

$$\hat{F} = \frac{F - \mu}{\sigma + \epsilon} \quad \mu = E[F], \quad \sigma = \text{Var}(F)$$

- To bring stochasticity, batch statistics are used

$$\hat{F} = \frac{F - \hat{\mu}}{\hat{\sigma} + \epsilon} \quad \hat{\mu} = \frac{1}{B} \sum_{k=1}^B F_k, \quad \hat{\sigma} = \sqrt{\frac{1}{B} \sum_{k=1}^B (F_k - \hat{\mu})^2}$$

- We make BN layer to learn its scaling/bias adjusting factors

$$\hat{F} = \gamma \frac{F - \hat{\mu}}{\hat{\sigma} + \epsilon} + \beta$$

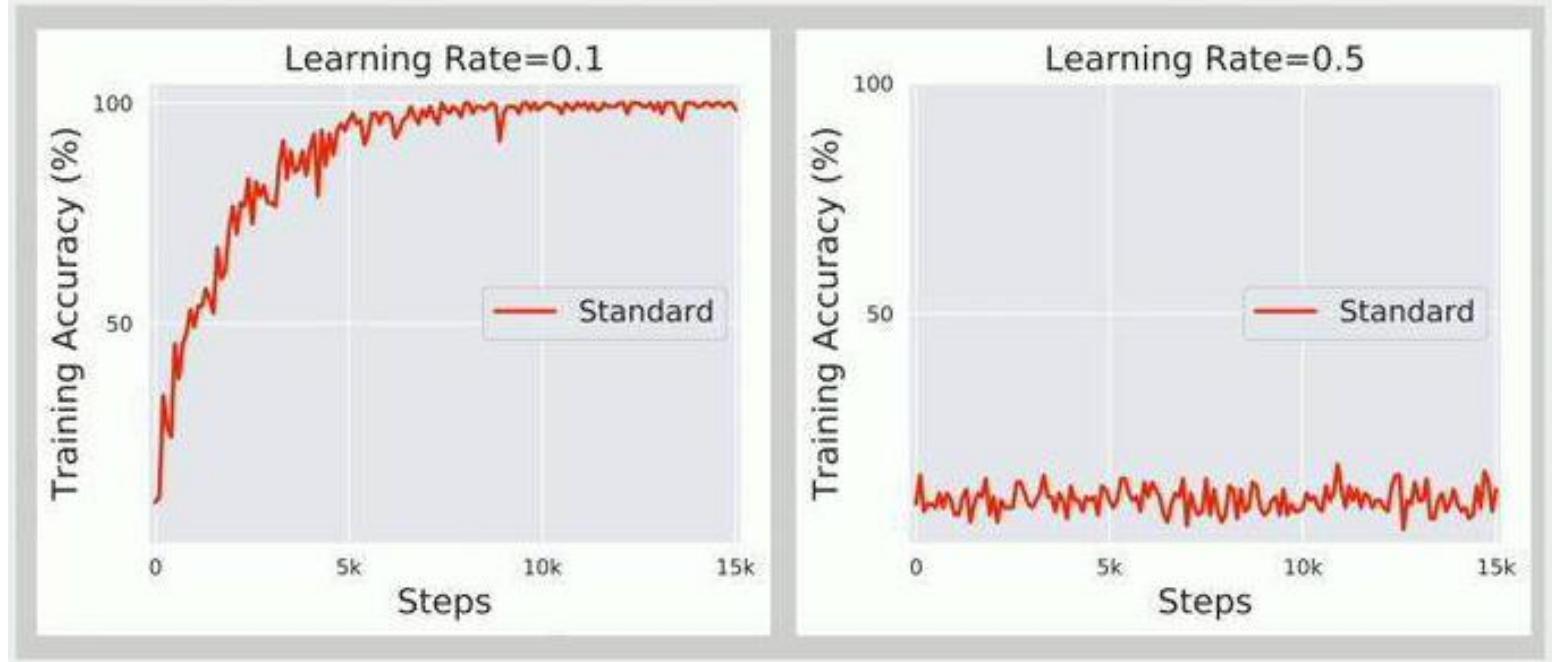
Learnable
parameter: scaling

Learnable
parameter: bias

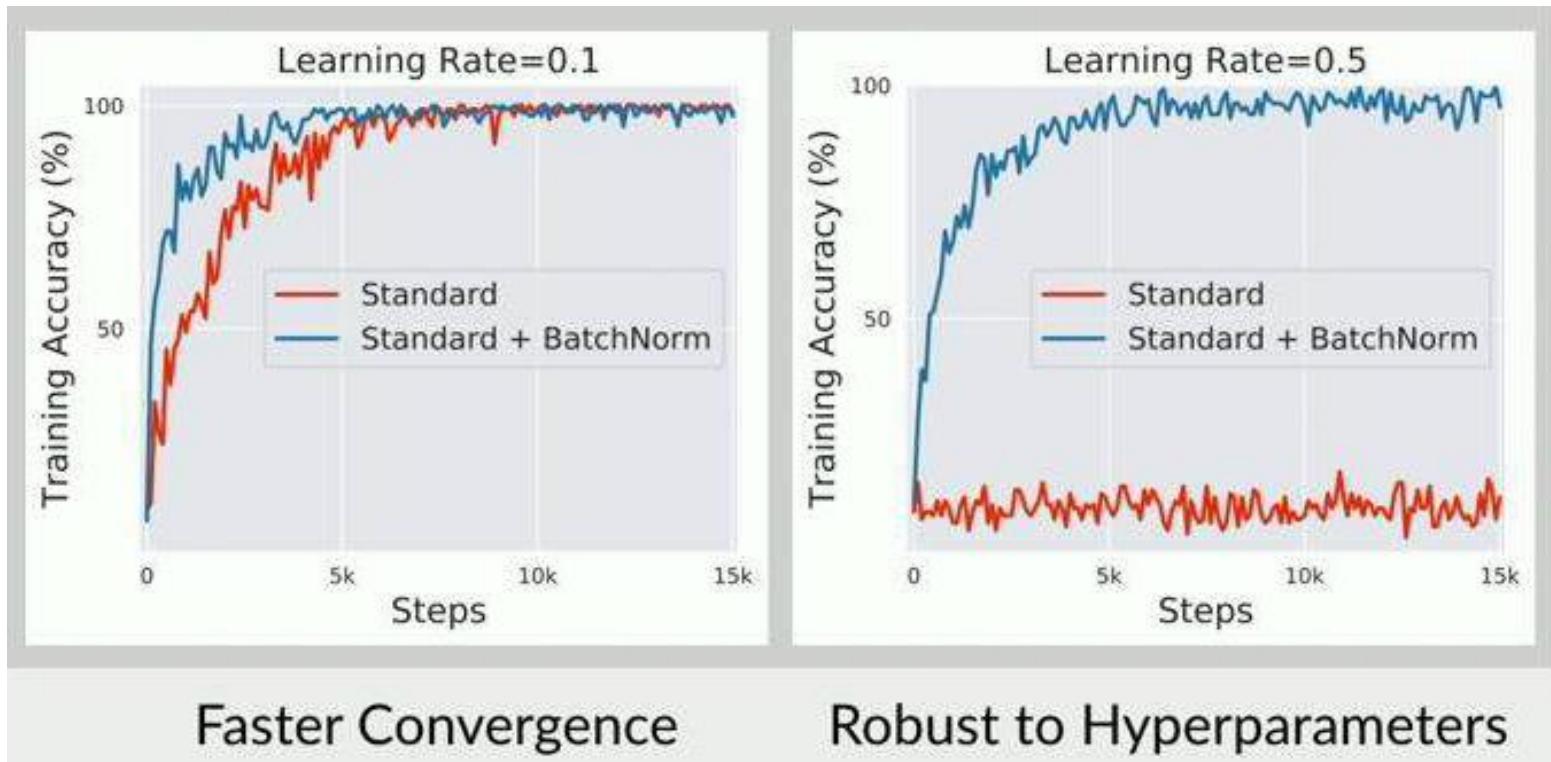
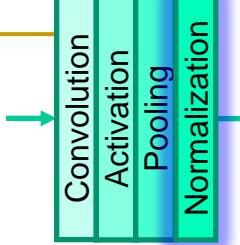
Why Do We Use BN?



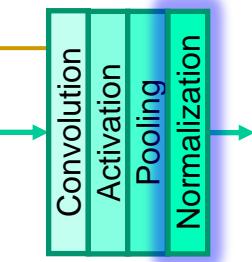
Network without BatchNorm



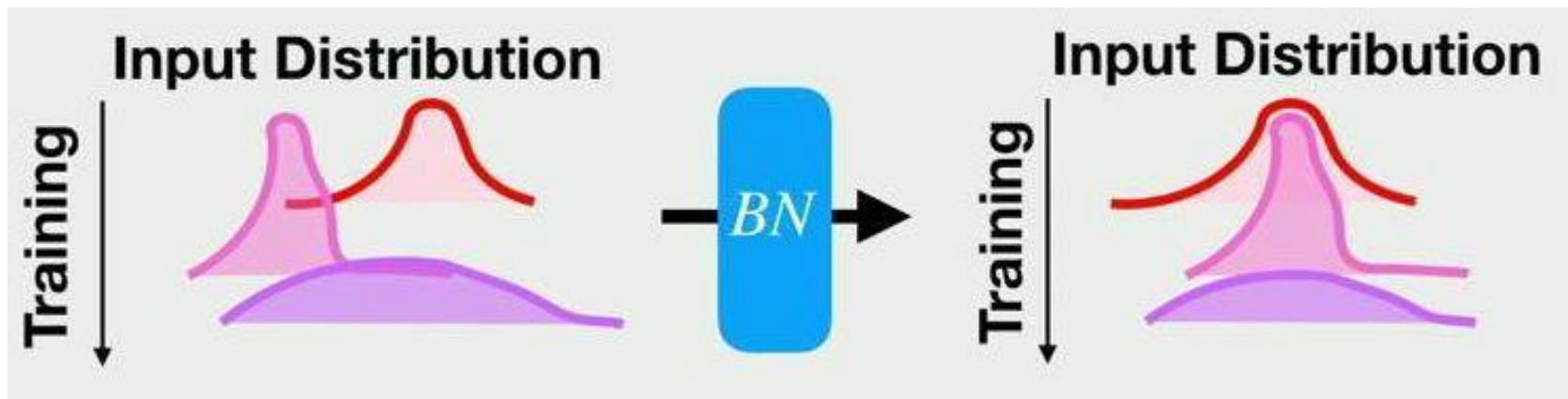
Why Do We Use BN?



How BN Helps?

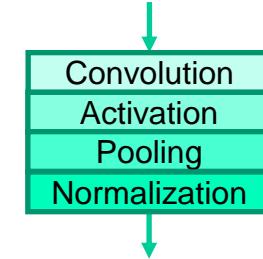


- BN reduces the internal covariate shift



Today

1. Motivation
2. Feature Learning
3. Building Block Design of CNN
 - Convolutional Layer
 - Activation Function
 - Pooling Layer
 - Normalization Layer



YOU ARE HERE!

THE END!