



Colorization with Convolutional Neural Networks

Division of Computer Science University of Minnesota,
Morris Morris, Minnesota, USA 56267

17 Nov 2018

Yutaro Miyata



Introduction



Pixel: Smallest square that can be displayed on a screen.

RGB: Color used in computer display. Combination of these can generate other colors.

Take Sample Color to find best matching

<https://ieeexplore.ieee.org/document/8014766/references#references>

Introduction

Semantic Segmentation



Person
Bicycle
Background

Classification

Outline

- Convolutional Neural Networks
- U-Net
- Generative Adversarial Network
- Image to Image Translation
- Global and Local Image Priors for Automatic Image Colorization
- Result

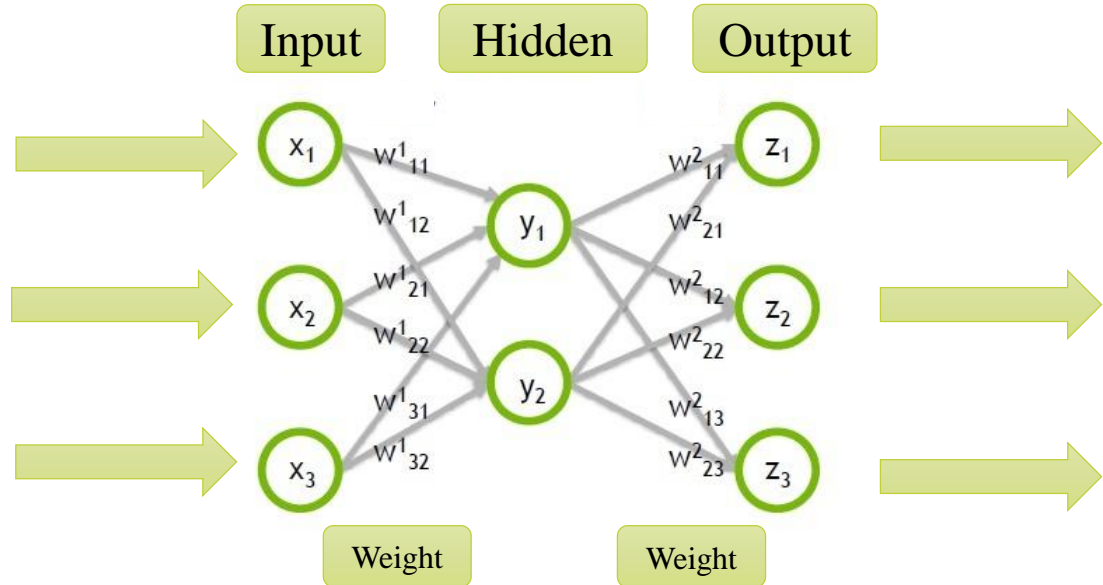
Artificial Neural Network

- What is an Artificial Neural Network?

BW

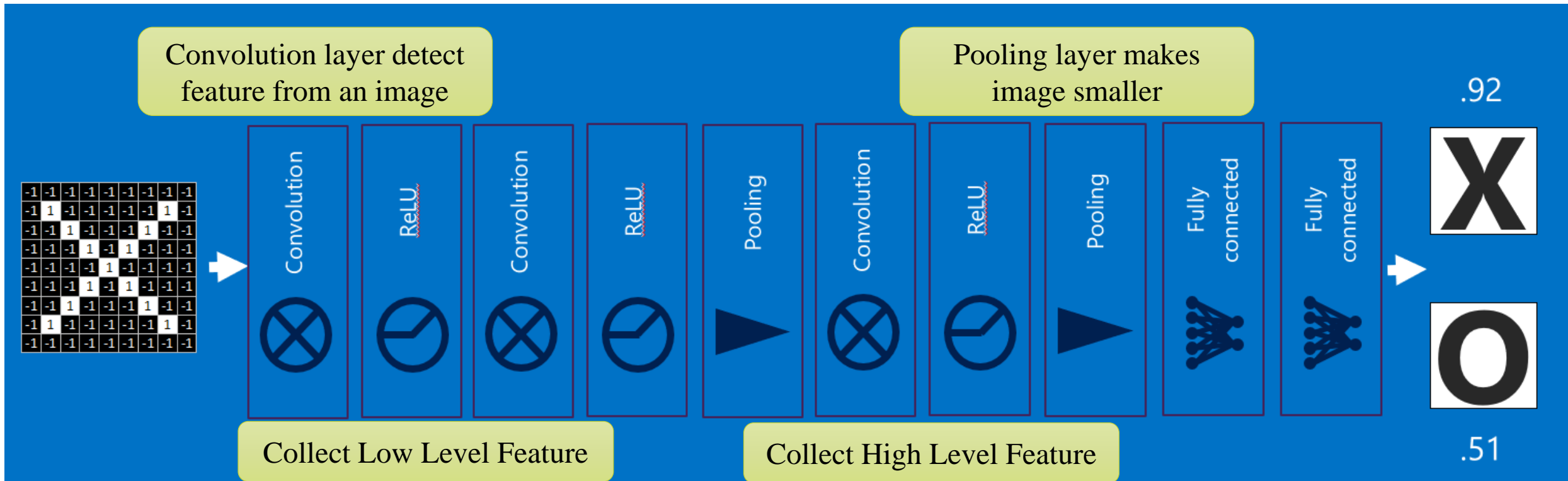
Color

Multilayer Perceptron

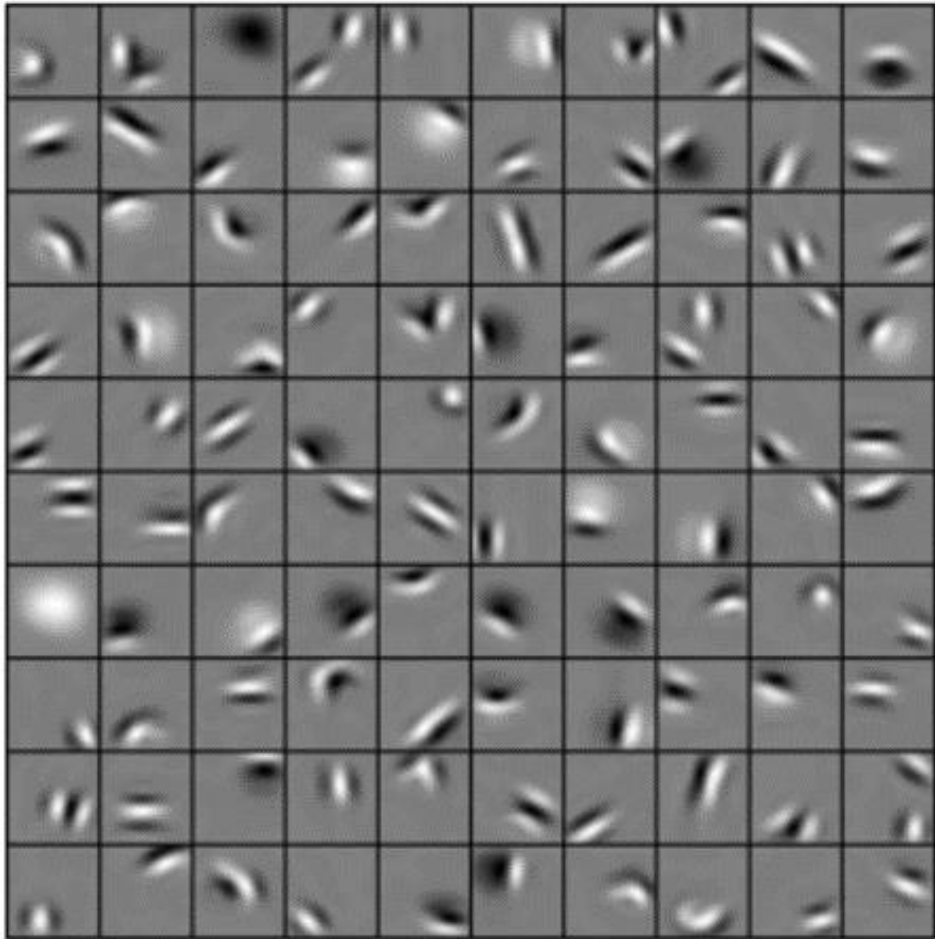


Convolutional Neural Networks

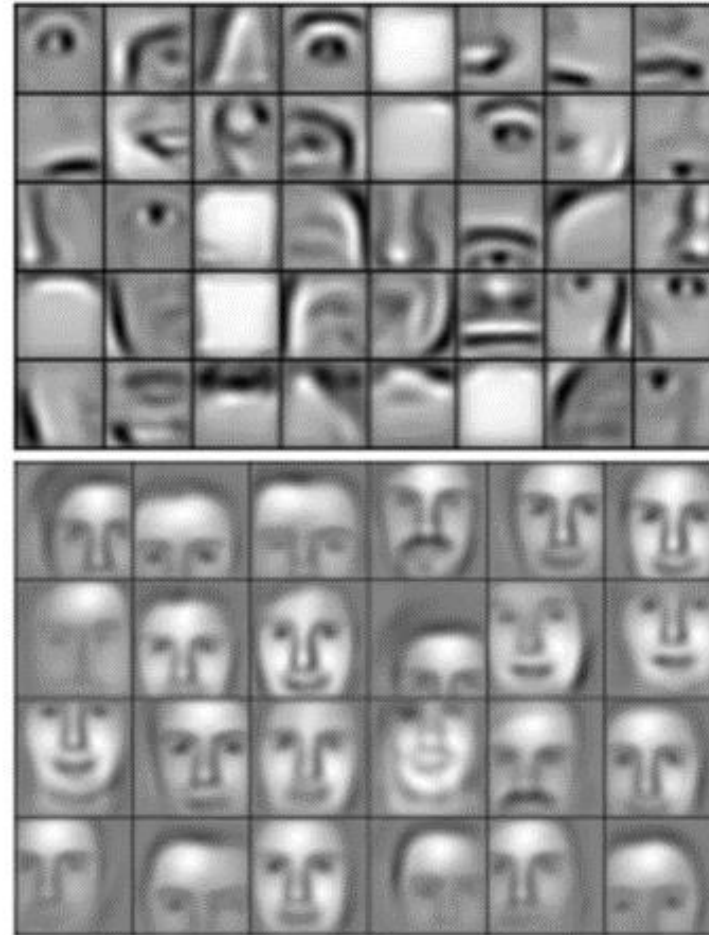
- Structure Of CNNs



Low level feature



High level feature



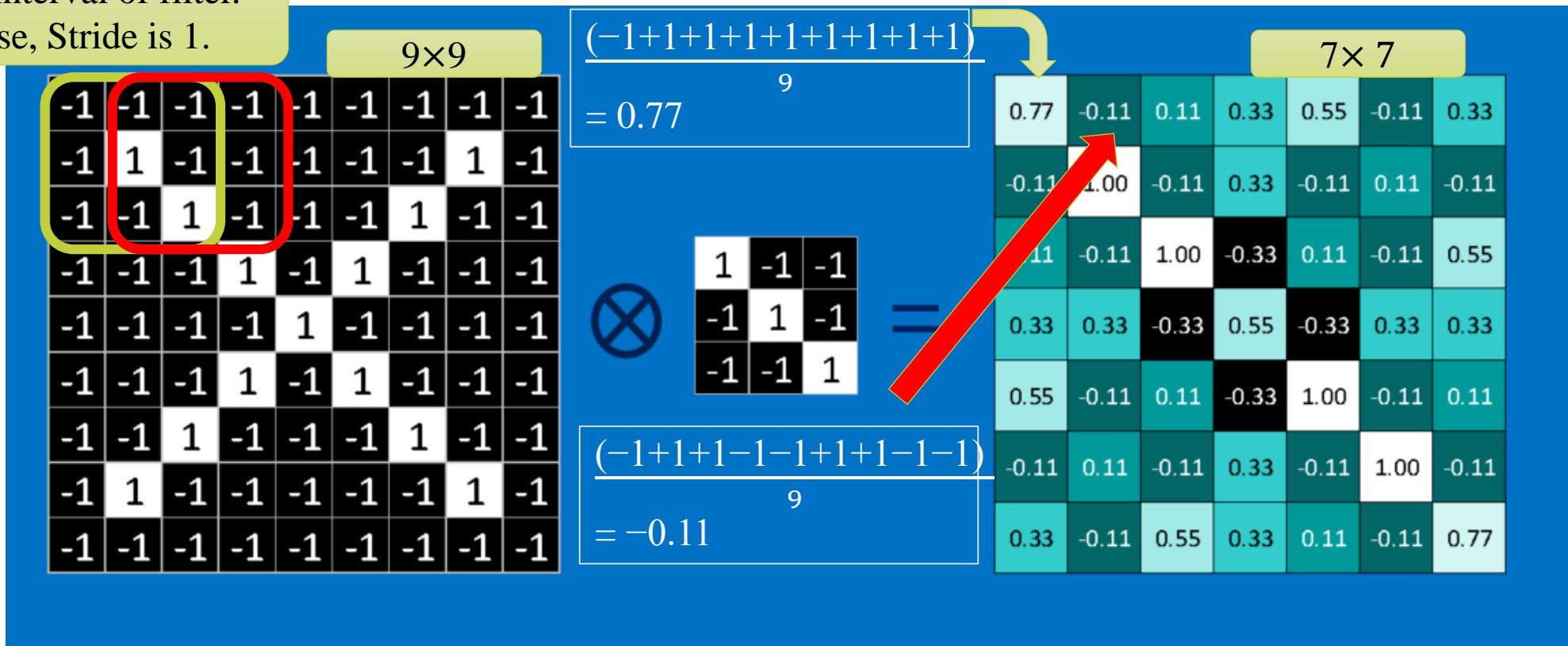
Feature: Pattern which network find from many data set.

Feature map: Certain combination of feature is found in an image.

- Convolutional Layer

- This layer detects many features from image. It applies many filters to image to get feature maps.

Stride is interval of filter.
In this case, Stride is 1.



9×9×1

-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	-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



Filters

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



7×7×3

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-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	-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	-1	1



0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

-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	-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



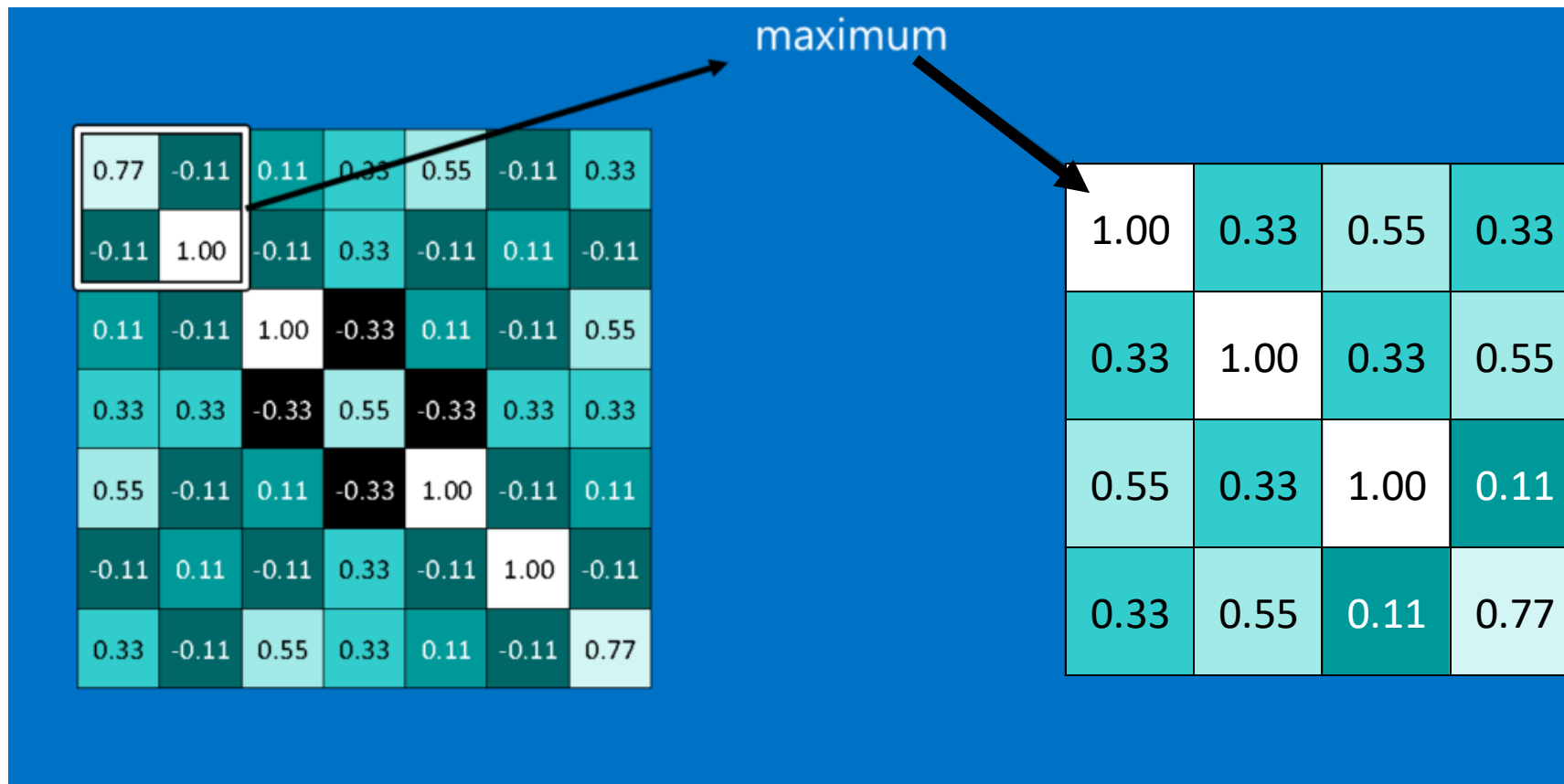
-1	-1	1
-1	1	-1
1	-1	-1



0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33

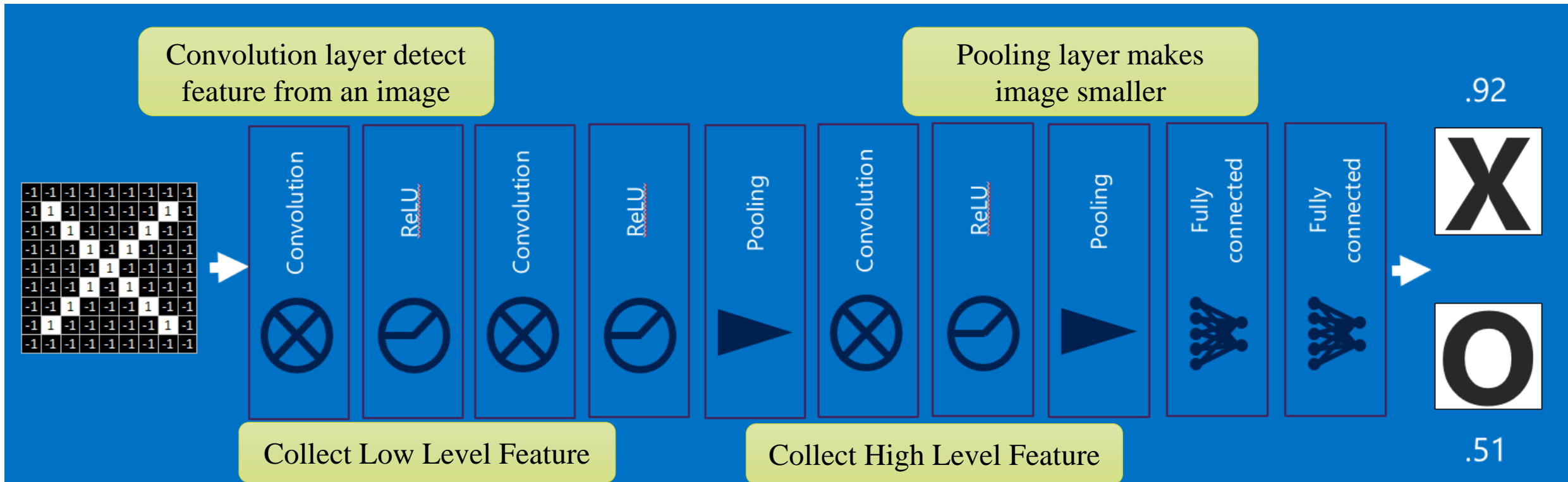
- Pooling Layer

- This layer transform a manageable form to clarify features. After pooling process, an image size will be $\frac{1}{4}$, and a computational cost is reduced.



Convolutional Neural Networks

- Structure Of CNNs
 - CNNs consists of many convolution layer and pooling layer.

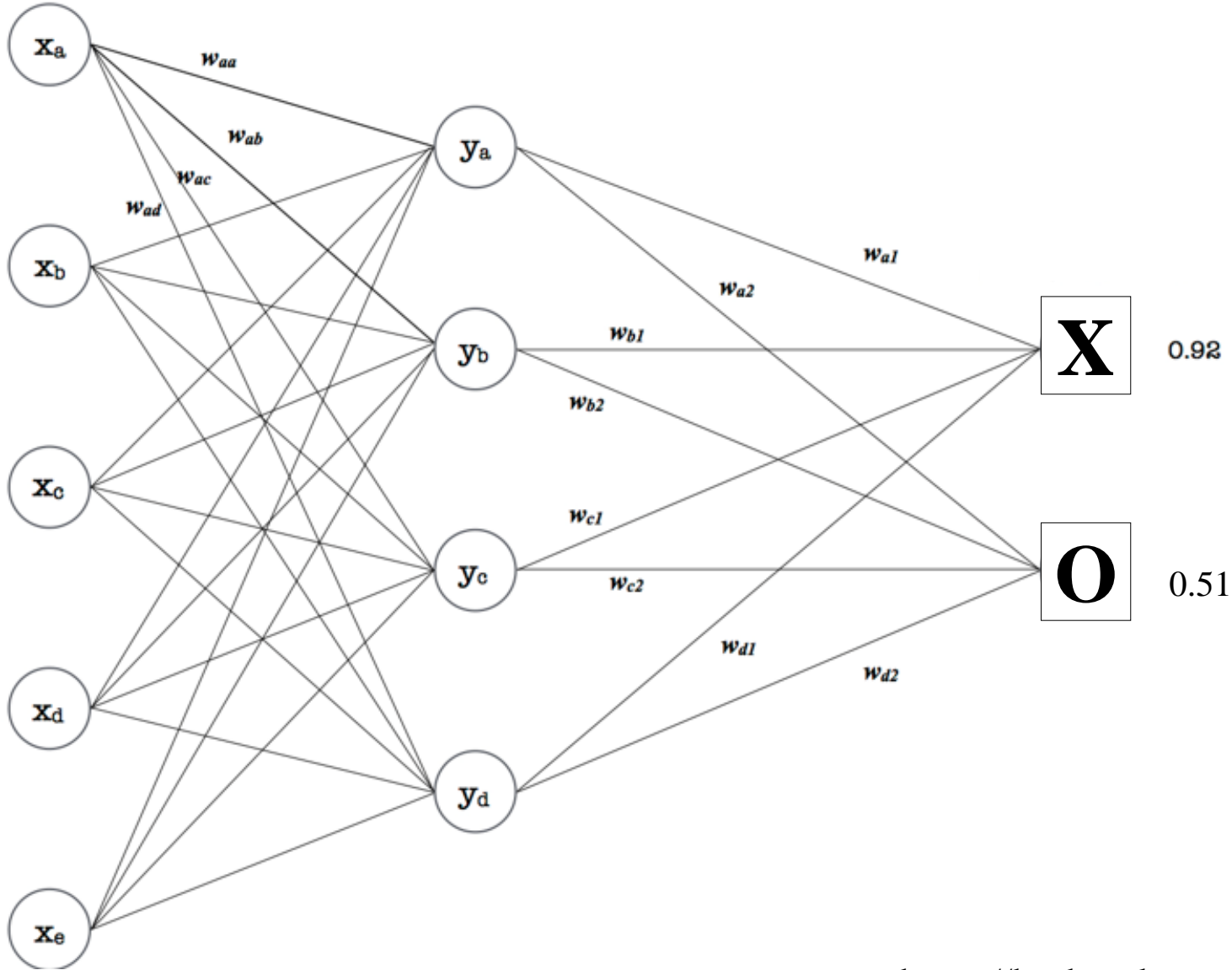


0.9
0.65
0.45
0.87
0.96
0.73
0.23
0.63
0.44
0.89
0.94
0.53

Input Layer

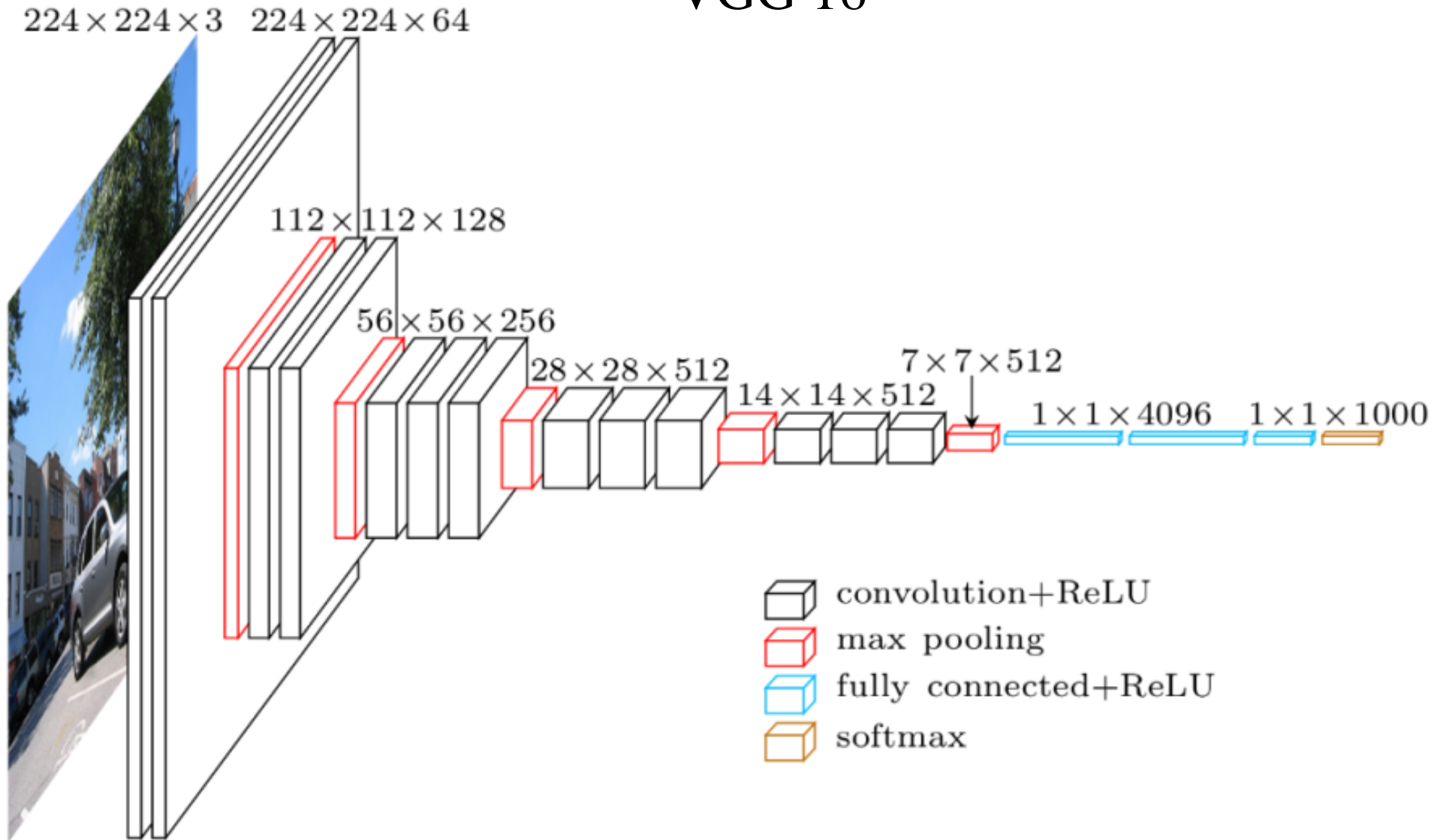
Hidden Layer

Output Layer



<https://hashrocket.com/blog/posts/a-friendly-introduction-to-convolutional-neural-networks>

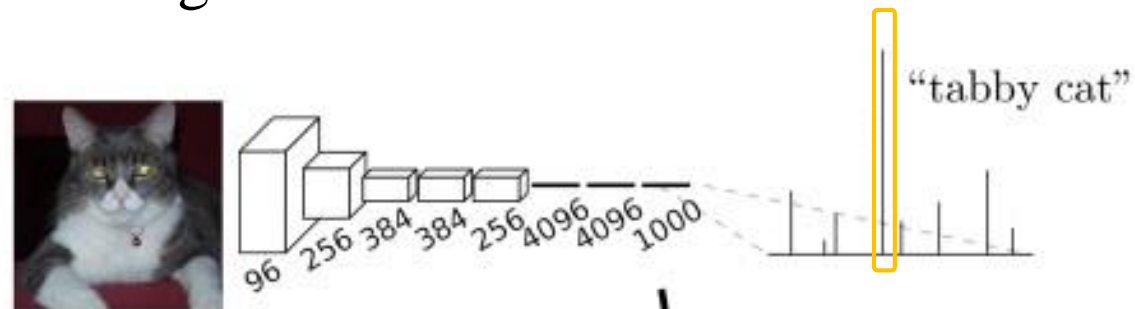
• VGG 16



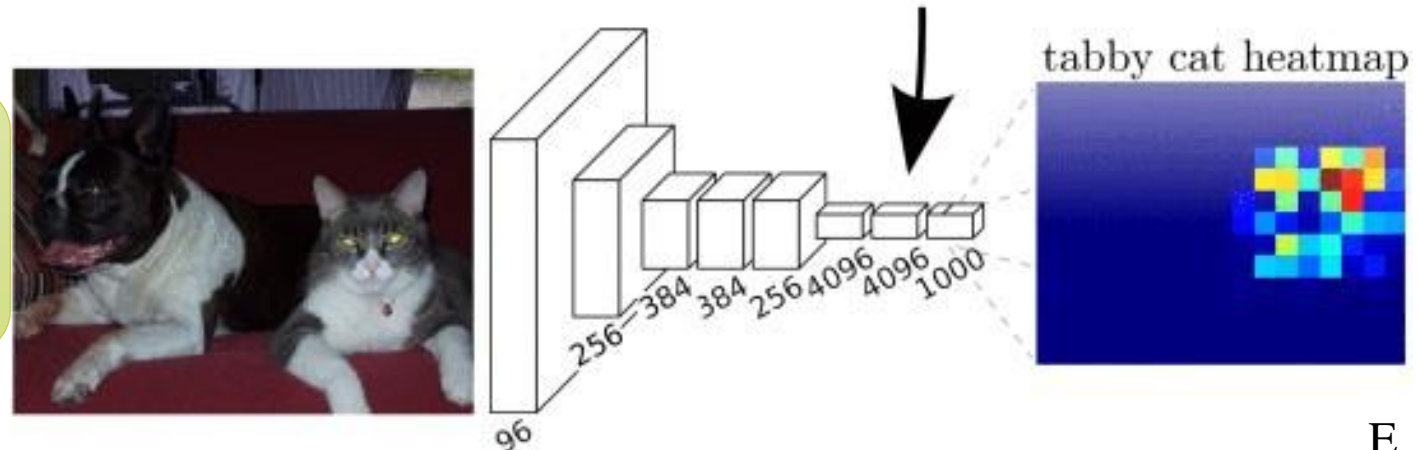
U-Net

- What is U-Net ?
 - U-Net is one method of Semantic Segmentation.

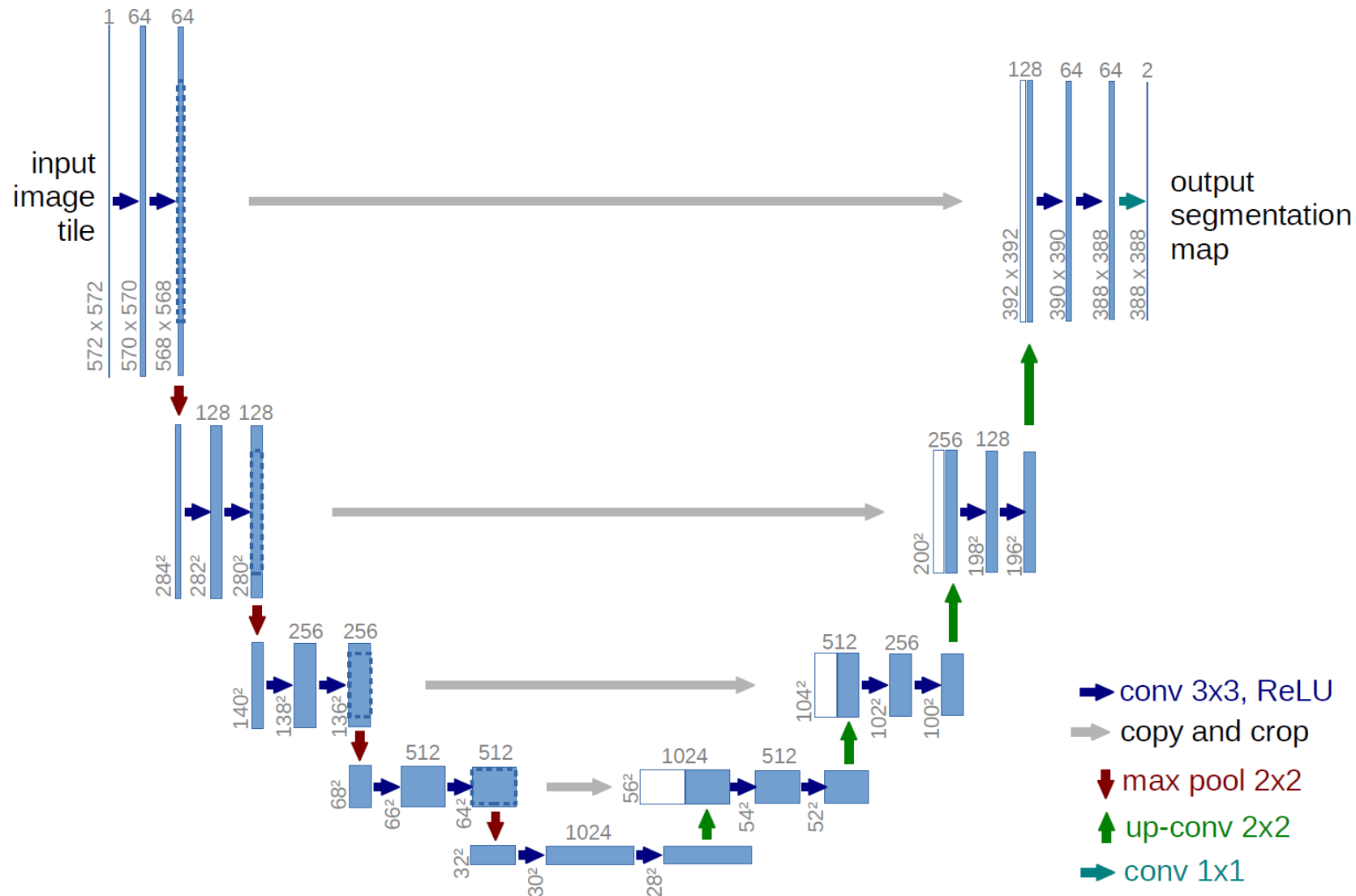
CNNs:
One dimensional vector

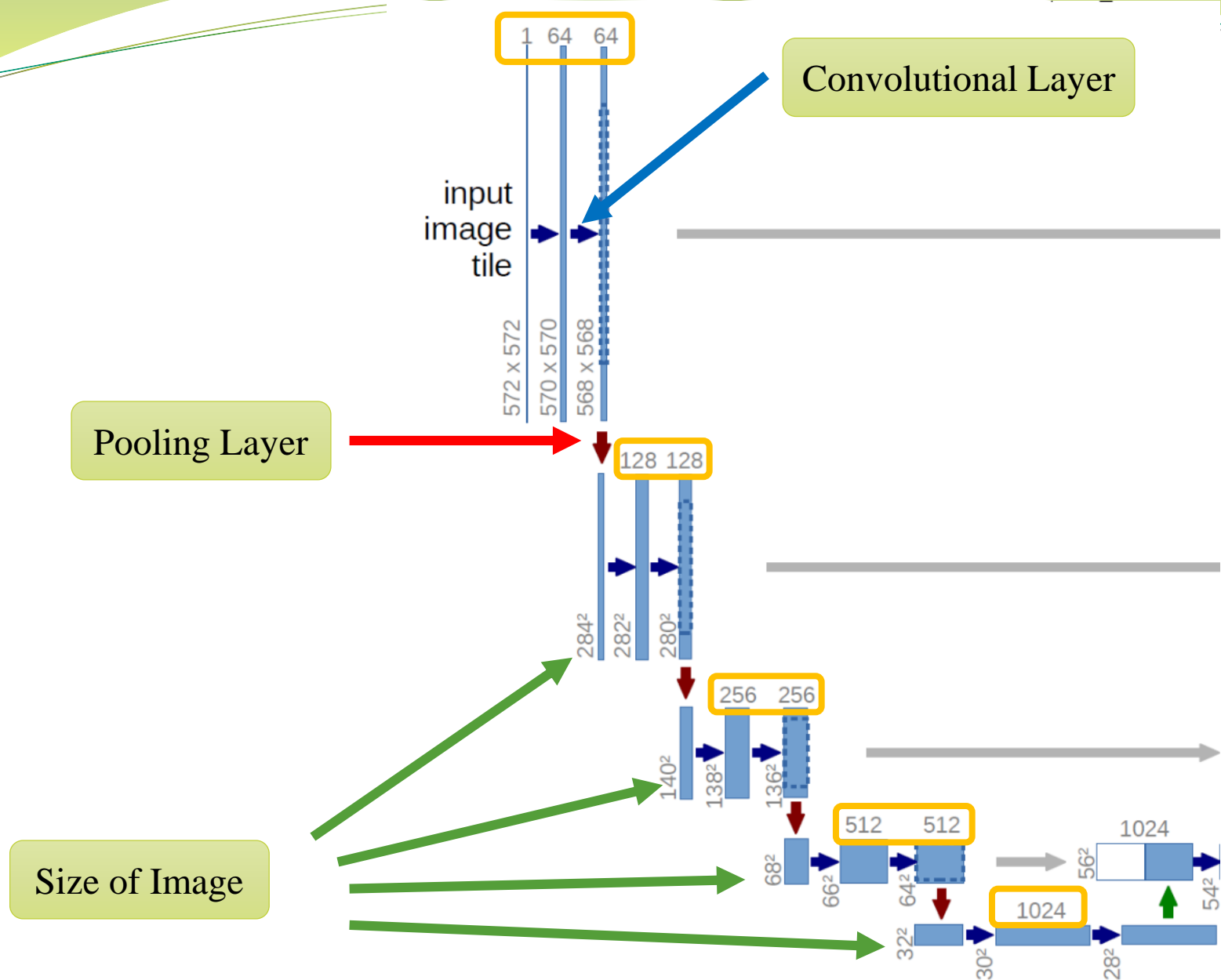


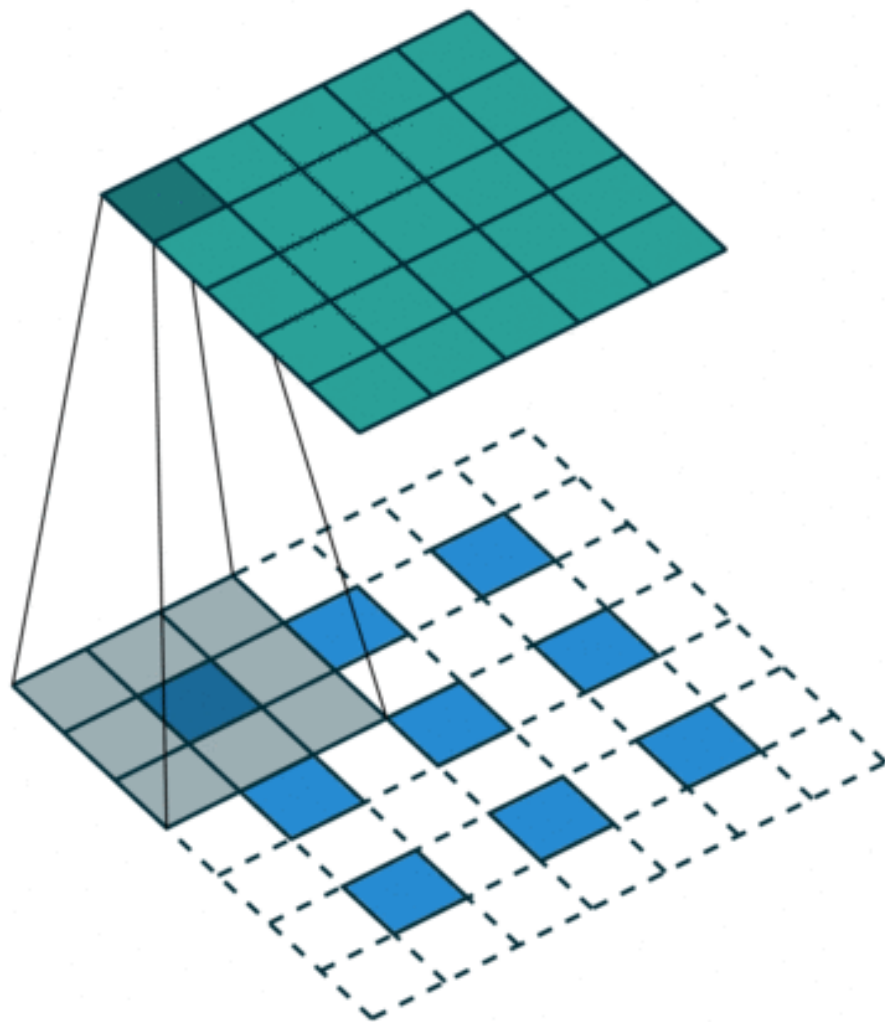
U-Net:
Two dimensional map
($H \times W$)



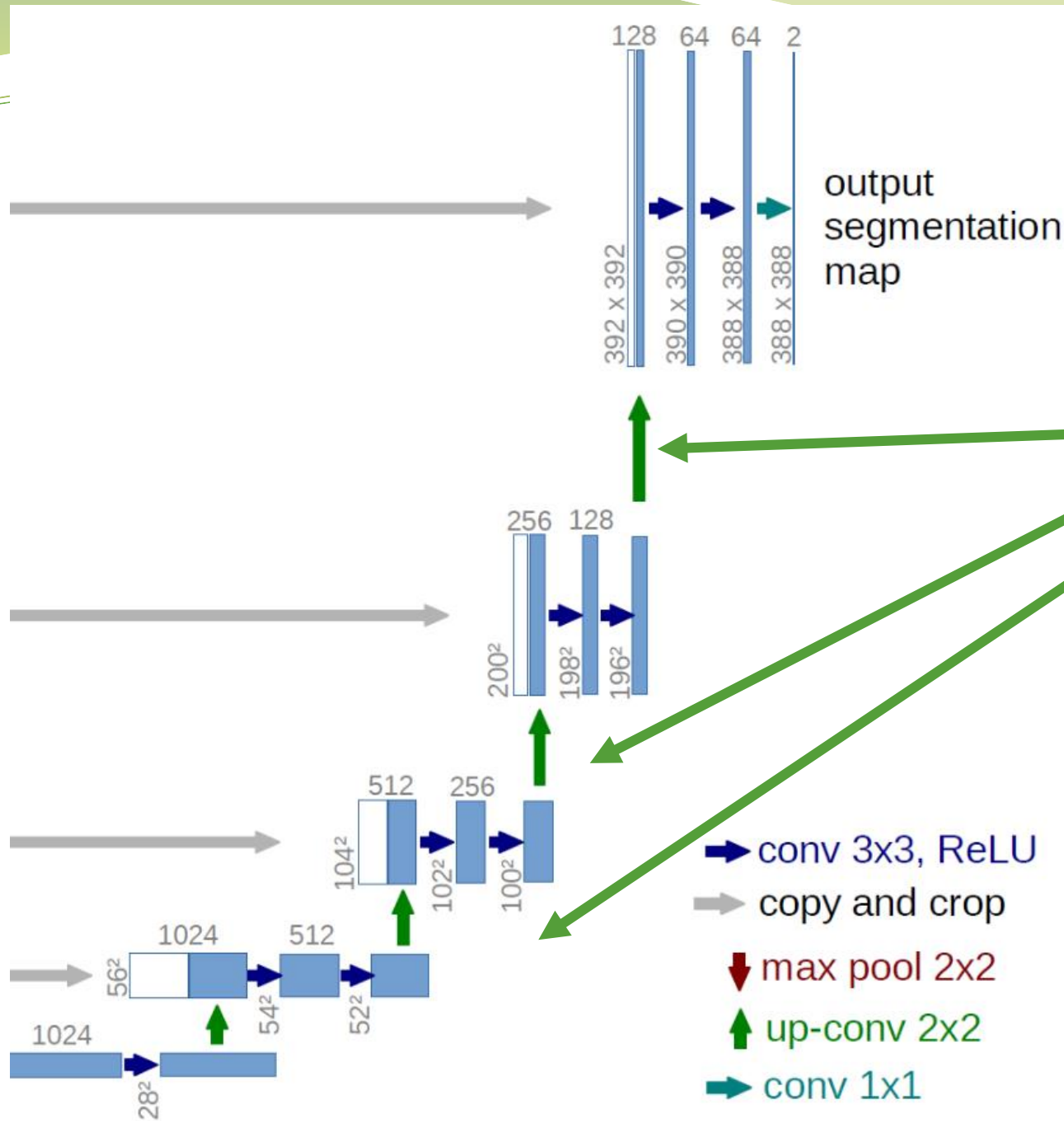
• Structure Of U-Net



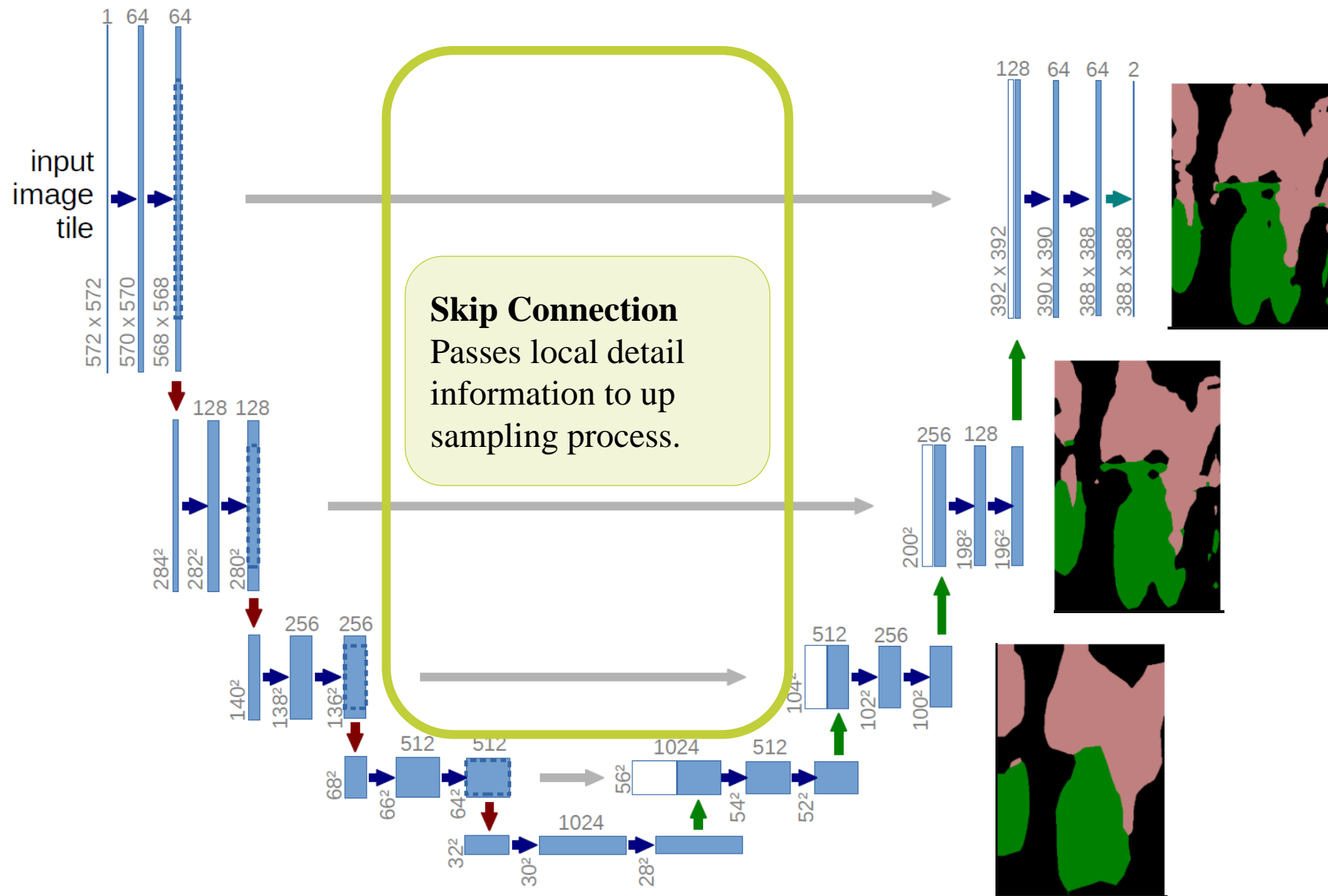




<https://towardsdatascience.com/types-of-convolutions-in-deep-learning-717013397f4d>



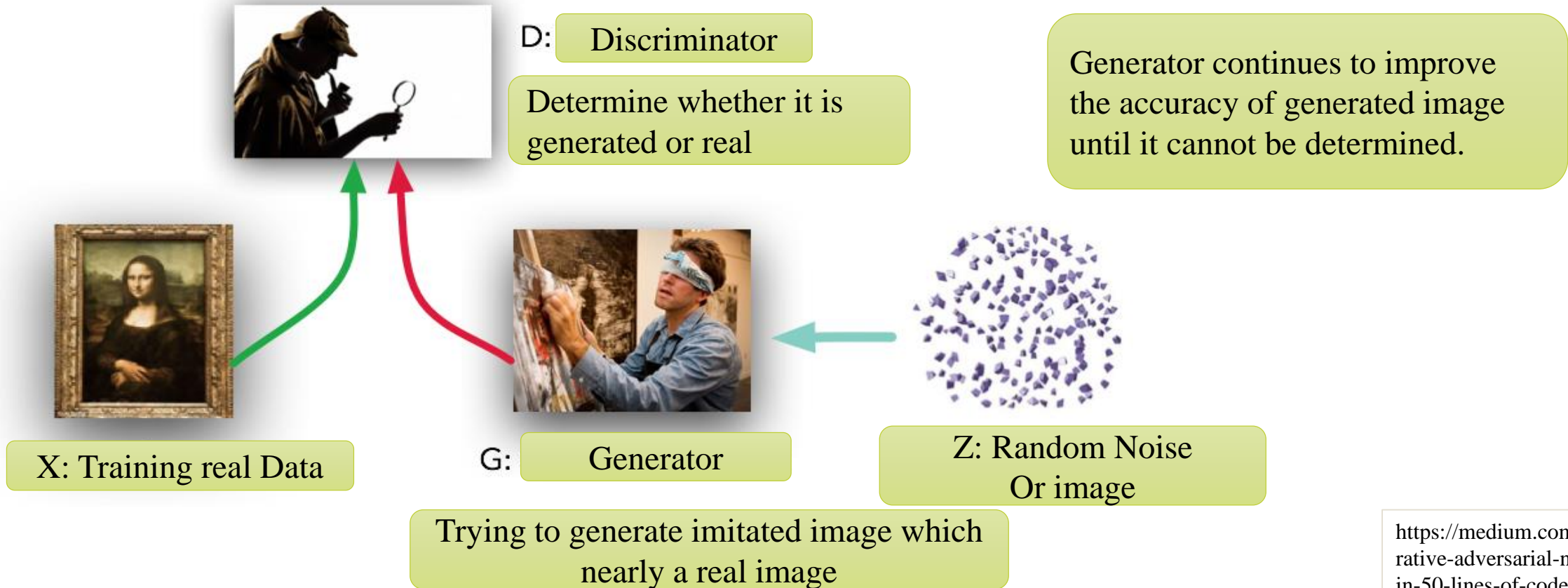
Up Sampling:
Expands feature map
because it became coarse
due to Pooling Layer



Generative Adversarial Network

- Structure of GAN

- The Generative Adversarial Network is generator (G) and discriminator (D) model.



Zebras ↔ Horses



zebra → horse



horse → zebra

Take real image as input

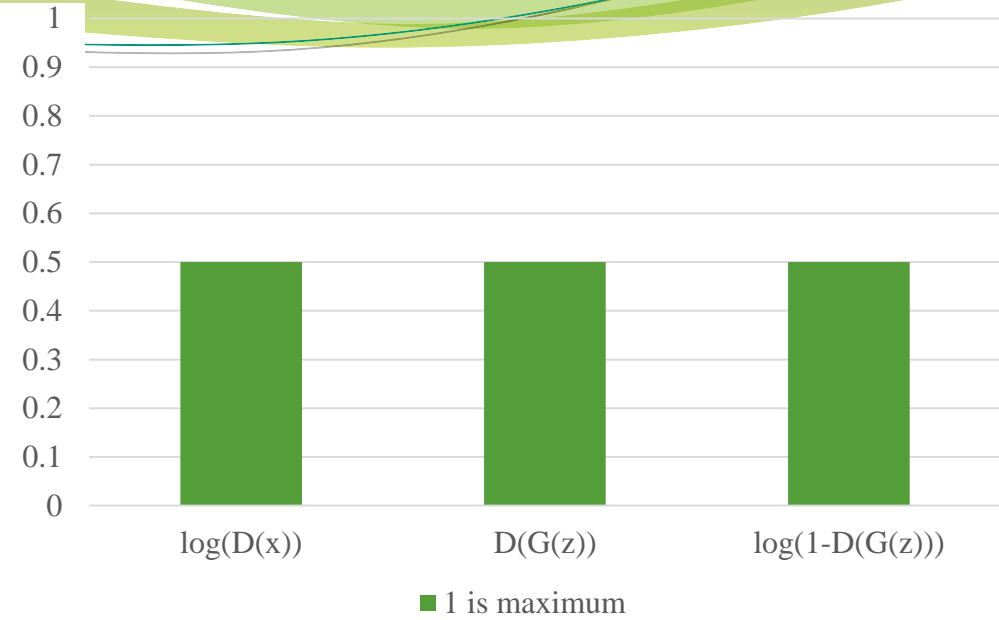
Trying to be bigger(1)

Desired result

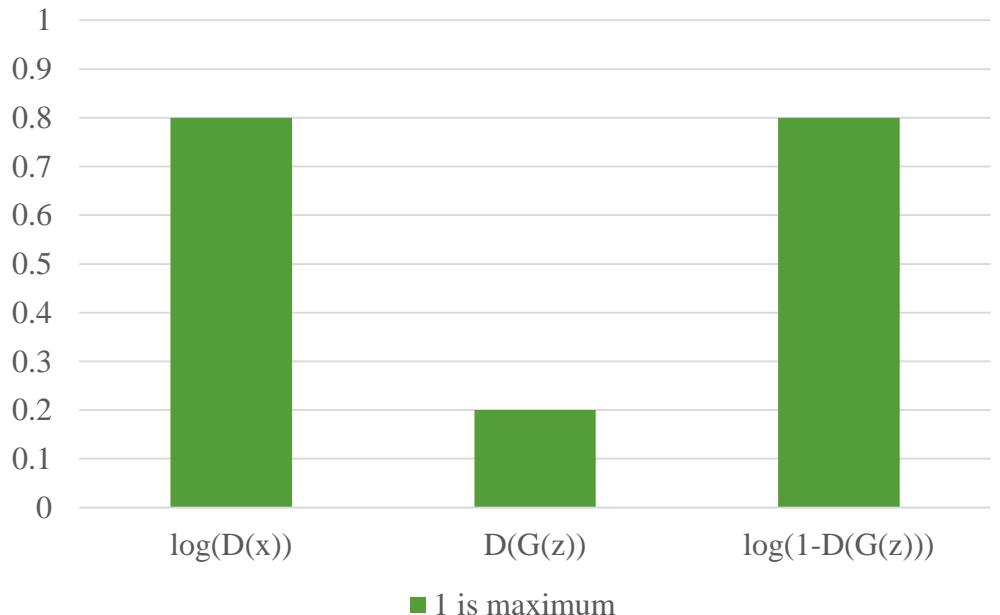
$$\min_G \max_D V(D, G) = E_{x \sim P_{data}(x)} [\log D(x)] + E_{z \sim P_z(z)} [\log(1 - D(G(z)))].$$

Take noise as input

Trying to be smaller (0)



D is well trained



D is not well trained

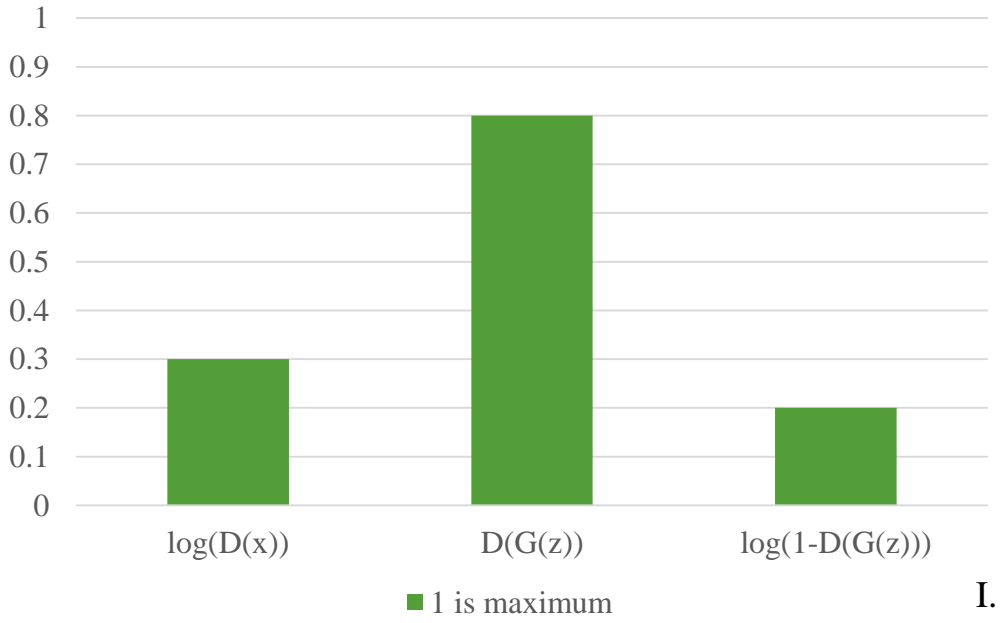


Image to Image Translation

- Prove effectiveness of cGAN and U-Net
- The authors adopt one methods.
 - Amazon Mechanical Turk: Show a generated image and real image to human and evaluate it. Measure the percentage which tester says real.

	Turker	Turker
	Photo to Map	Map to Photo
L1	2.8% ± 1.0%	0.8% ± 0.3%
L1 + cGAN	6.1% ± 1.3%	18.9% ± 2.5%

L1: This is traditional method

Aerial photo to map



Map to aerial photo



Ground truth

L1

cGAN

L1 + cGAN



Effectiveness on Colorization

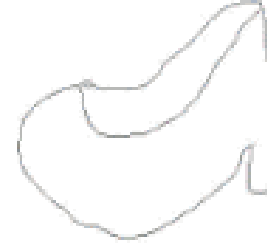
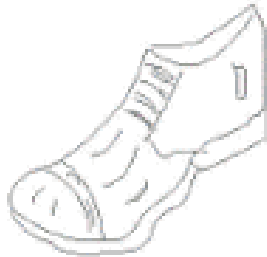
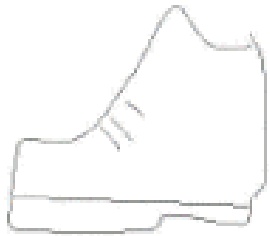
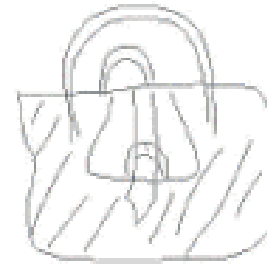
Method	Turkers Labeled real
Zhang et al 2016	27.8% \pm 2.7%
L1 + cGAN	22.5% \pm 1.6%

Another method of colorization which specifically engineered to do well on colorization.



Sketches to Bags Sketches to Shoes

sketches → shoes sketches → bag



Effectiveness of U-Net

L1+cGAN

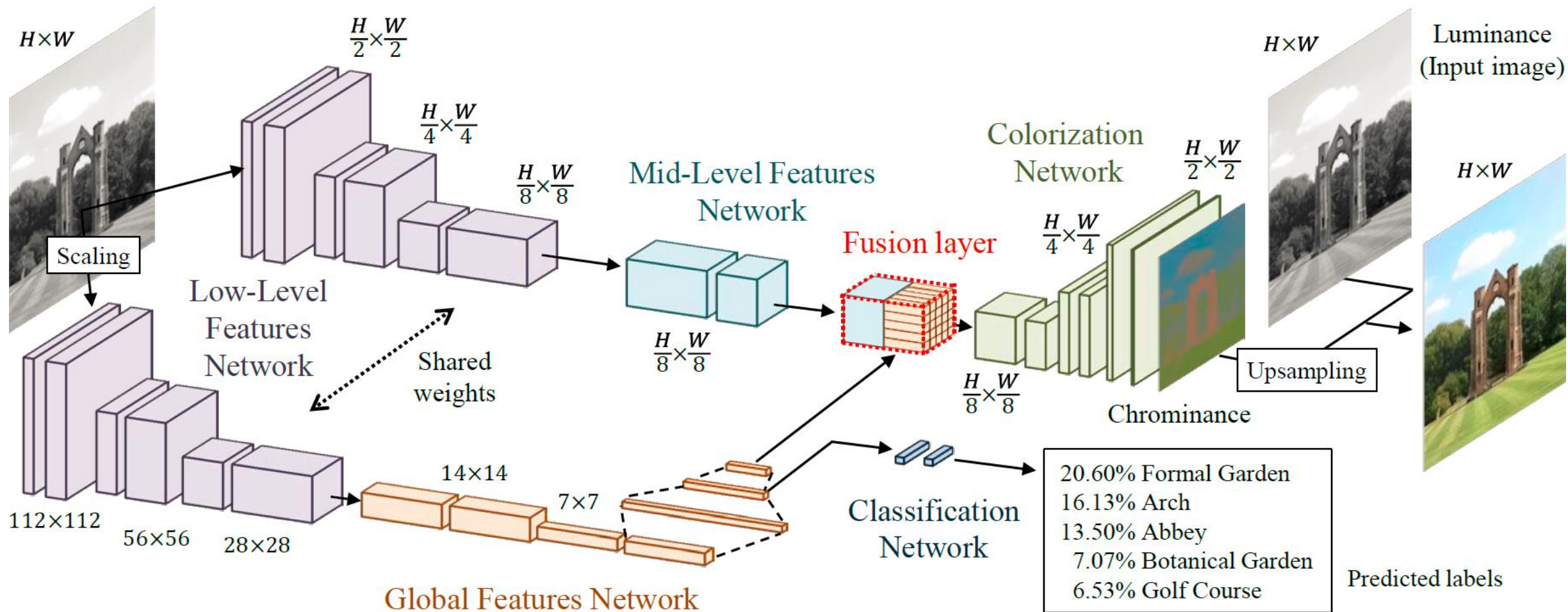
Traditional Encoder
Decoder Network



U-Net



Global and Local Image Priors for Automatic Image Colorization



Middle and Global Level Feature Network

Middle Level Feature Network

It consists of 2 convolution layer to collect middle level feature information

Type	Kernel	Stride	Outputs
conv.	3×3	1×1	512
conv.	3×3	1×1	256

$\frac{H}{4} \times \frac{W}{4}$

$\frac{H}{8} \times \frac{W}{8}$

Mid-Level Features Network

Fusion layer

$\frac{H}{8} \times \frac{W}{8}$

$\frac{H}{8} \times \frac{W}{8}$

Shared weights

Fusion Layer:
Combine global feature(256) and mid level feature(256)

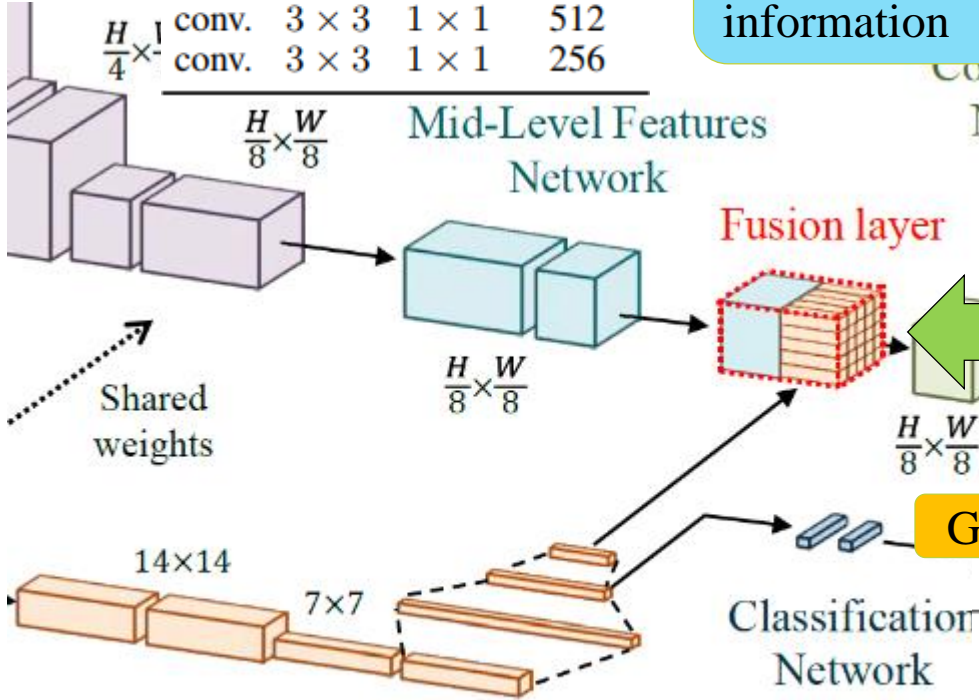
Global Feature Network

Type	Kernel	Stride	Outputs
conv.	3×3	2×2	512
conv.	3×3	1×1	512
conv.	3×3	2×2	512
conv.	3×3	1×1	512
FC	-	-	1024
FC	-	-	512
FC	-	-	256

Classification Network

Global Features Network

It consists of 4 convolution layer and 3 fully connected layer to compute a 256 dimensional vector representation of image

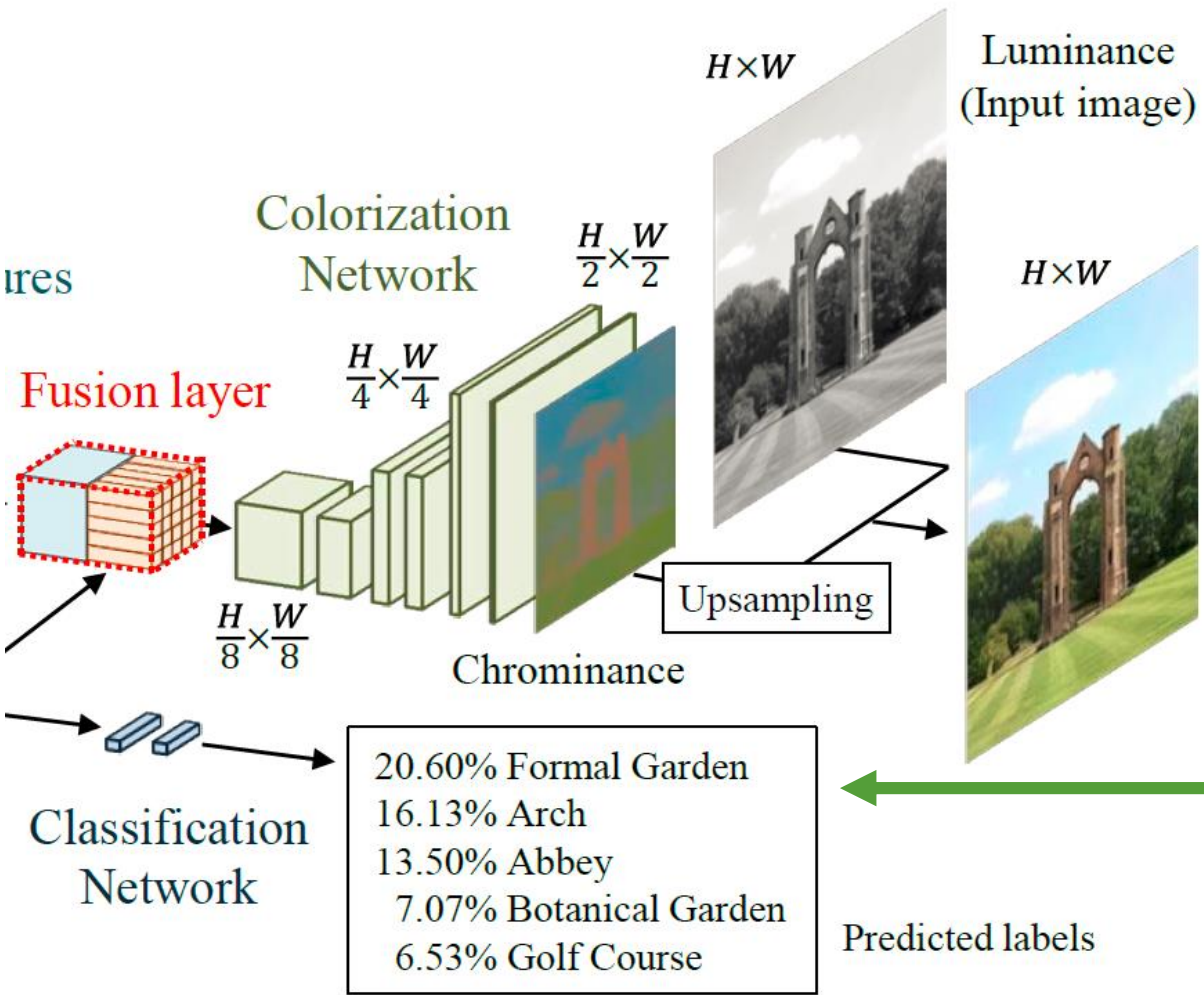


Colorization

Colorization Network

Type	Kernel	Stride	Outputs
fusion	-	-	256
conv.	3×3	1×1	128
upsample	-	-	128
conv.	3×3	1×1	64
conv.	3×3	1×1	64
upsample	-	-	64
conv.	3×3	1×1	32
output	3×3	1×1	2

It consists of convolution layer and up sampling layer to expands feature map until twice as wide and tall.



- 20.60% Formal Garden
- 16.13% Arch
- 13.50% Abbey
- 7.07% Botanical Garden
- 6.53% Golf Course

Predicted labels

Computer learns color and label at same time to understand global context. It guide the training of the global image feature.



Result

- Each method has a field of excellence. For example, global local colorization is specialized for black and white picture colorization. Image to image translation method can be widely used in various fields.
- GAN will be used in various fields because of versatility.



Question?