

Computer Vision

Naeemullah Khan

naeemullah.khan@kaust.edu.sa



جامعة الملك عبد الله
للغعلوم والتكنولوجيا
King Abdullah University of
Science and Technology

KAUST Academy
King Abdullah University of Science and Technology

February 7, 2025

Table of Contents

1. Introduction.
2. Applications of Computer Vision.
3. Image Representation.
4. Convolutional Neural Network (CNN) and its Components.
5. CNN-based Architectures (AlexNet, VGG, InceptionNet, ResNet, EfficientNet, and MobileNet).

Learning Outcomes

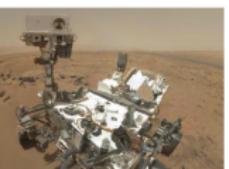
- ▶ Understand the basic concepts of Computer Vision and its real-world applications.
- ▶ Describe how images are represented and processed in a computer.
- ▶ Explain the fundamental building blocks of Convolutional Neural Networks (CNNs).
- ▶ Differentiate between popular CNN architectures (AlexNet, VGG, InceptionNet, ResNet, EfficientNet, and MobileNet) and their key innovations.

Building artificial systems that process, perceive, and reason about visual data

Computer Vision is Everywhere



Left to right:
[Image by Roger H Giesen](#) is licensed under [CC BY 2.0](#)
[Image](#) is CC0 1.0 public domain
[Image](#) is CC0 1.0 public domain
[Image](#) is CC0 1.0 public domain



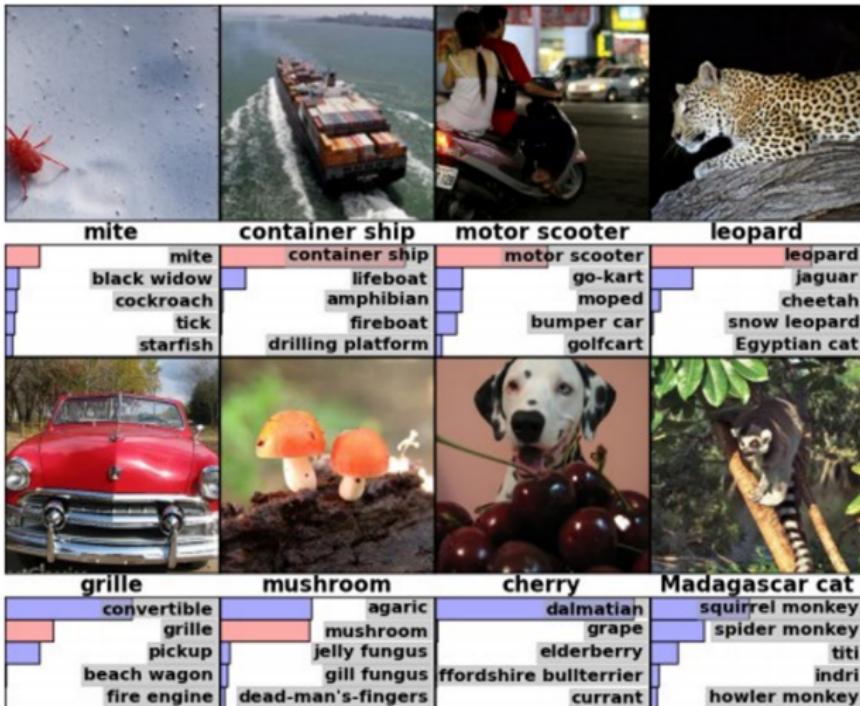
Left to right:
[Image](#) is free to use
[Image](#) is CC0 1.0 public domain
[Image](#) by NASA is licensed under [CC BY 2.0](#)
[Image](#) is CC0 1.0 public domain



Bottom row, left to right:
[Image](#) is CC0 1.0 public domain
[Image](#) by Derek Keats is licensed under [CC BY 2.0](#)
changes made
[Image](#) is public domain
[Image](#) is licensed under [CC BY 2.0](#)
changes made

Some Applications of Computer Vision

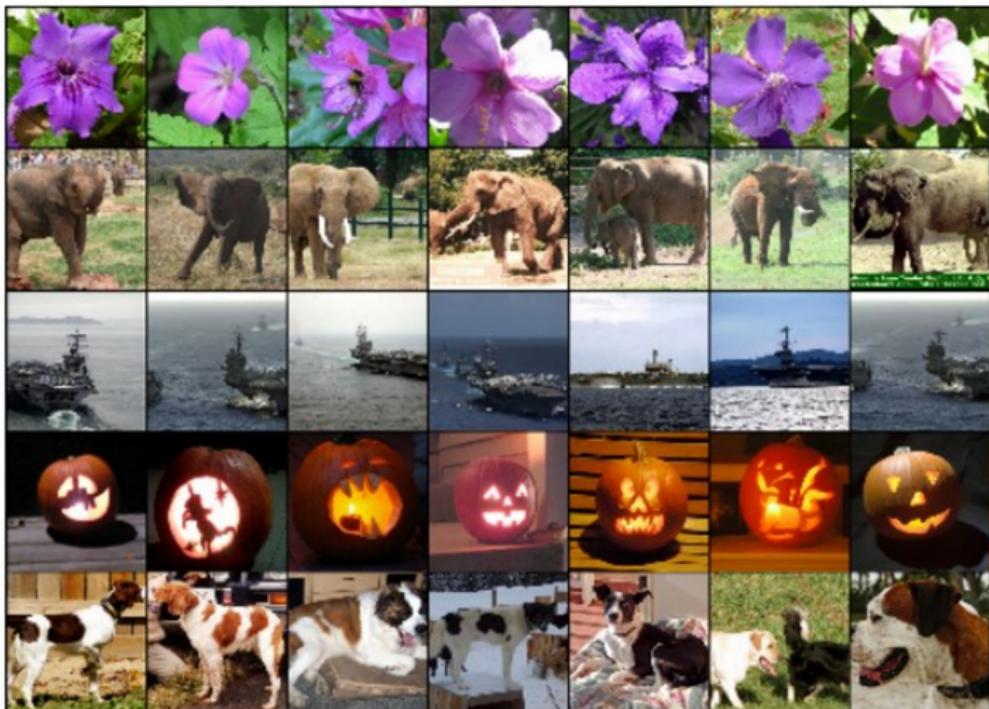
Image Classification



Some Applications (cont.)

Image Retrieval

→
I+ Just
Classification

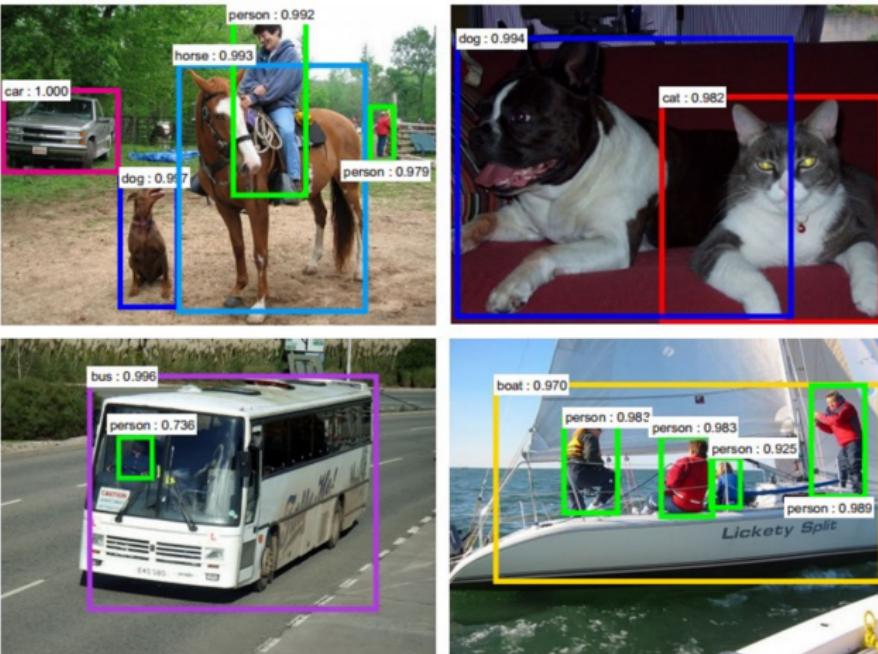


Input

output retrieved from database

Some Applications (cont.)

Object Detection

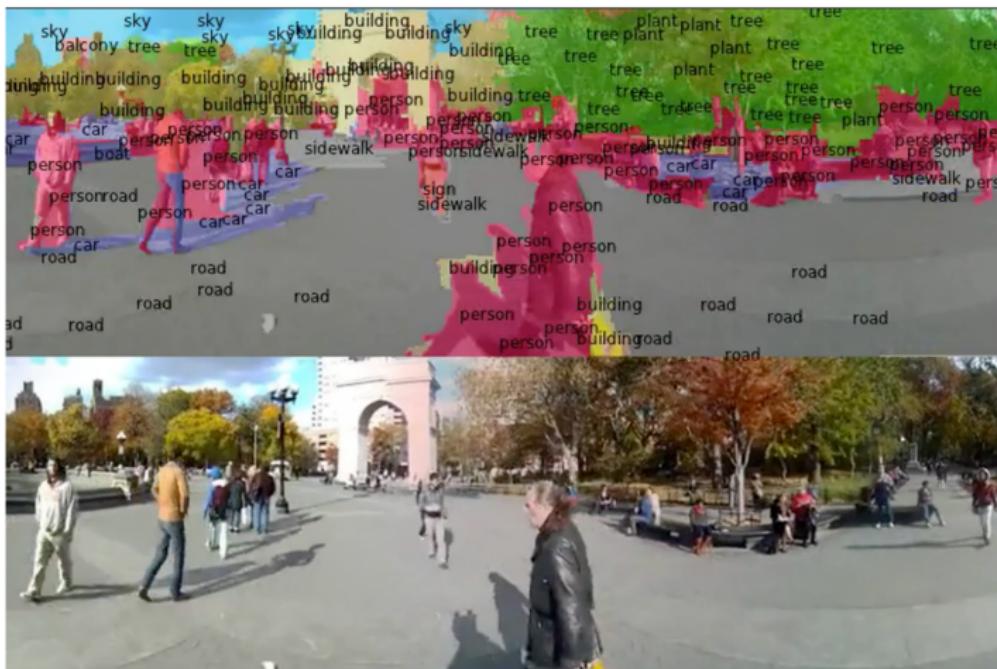


Ren, He, Girshick, and Sun, 2015

Some Applications (cont.)

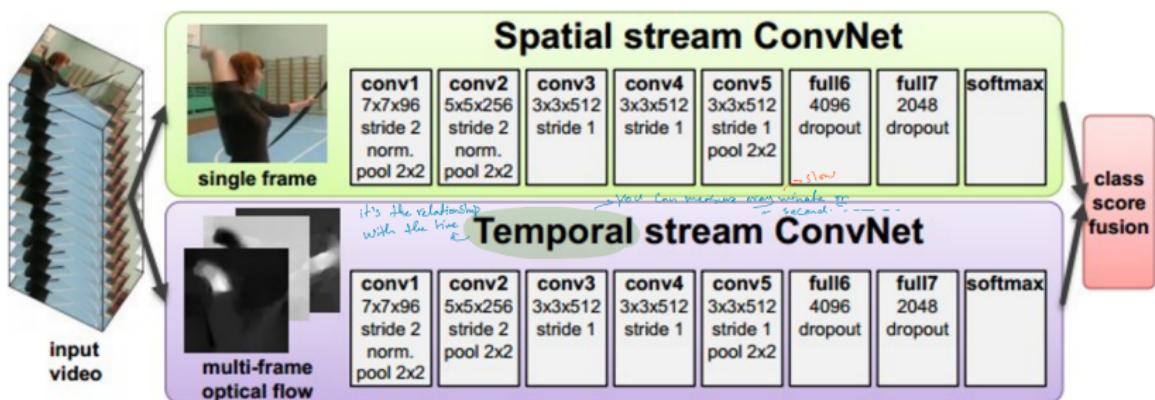
detection II mode

Image Segmentation



Fabaret et al, 2012

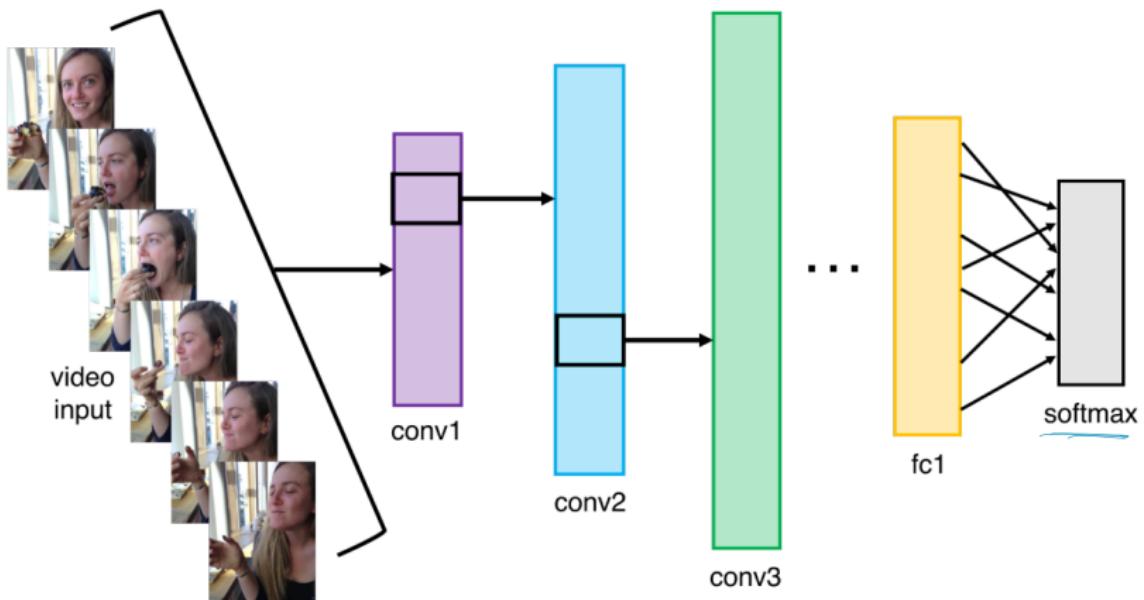
Video Classification



Simonyan et al, 2014

Some Applications (cont.)

Activity Recognition



Some Applications (cont.)



Pose Recognition (Toshev and Szegedy, 2014)



*how do you measure or Count an activity?

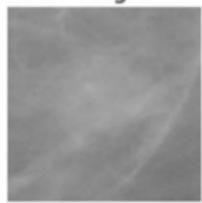
- By change
- By pose estimation → edge detection, find points then the connection

Some Applications (cont.)

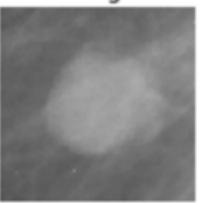
Image classification

Medical Imaging

Benign



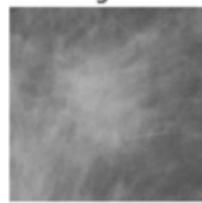
Benign



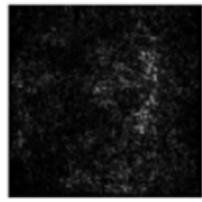
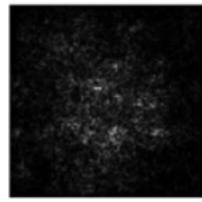
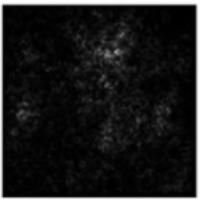
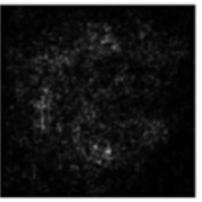
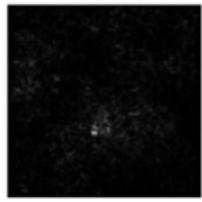
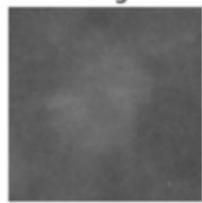
Malignant



Malignant



Benign



Some Applications (cont.)

image \rightarrow Text



A white teddy bear
sitting in the grass

Image Captioning

Vinyals et al, 2015

Karpathy and Fei-Fei, 2015



A man in a baseball
uniform throwing a ball



A woman is holding
a cat in her hand



A man riding a wave
on top of a surfboard



A cat sitting on a
suitcase on the floor



A woman standing on a
beach holding a surfboard

All images are CC0 Public domain:
<http://creativecommons.org/publicdomain/zero/1.0/>
<http://creativecommons.org/licenses/by-sa/1.0/>
<http://creativecommons.org/licenses/by-nd/1.0/>
<http://creativecommons.org/licenses/by-nc/1.0/>
<http://creativecommons.org/licenses/by-nd-nc/1.0/>

Some Applications (cont.)

Text → Image

Image Generation



“Teddy bears working on new AI research underwater with 1990s technology”

DALL-E 2

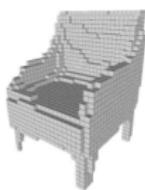
Some Applications (cont.)



Style Transfer

Some Applications (cont.)

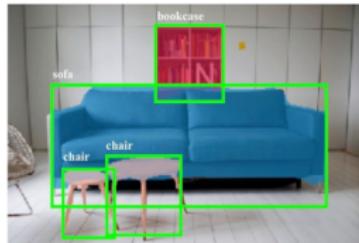
3D Vision



Choy et al., 3D-R2N2: Recurrent Reconstruction Neural Network (2016)



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



Gkioxari et al., "Mesh R-CNN", ICCV 2019

How to represent an image?

- ▶ Images are represented as Matrices with elements in [0, 255]
- ▶ Grayscale images have one channel.



157	153	174	168	150	162	129	161	172	161	155	166					
155	182	163	74	75	62	33	17	110	210	180	164					
180	180	50	14	34	6	10	33	48	106	159	181					
206	106	5	124	131	111	120	204	166	15	56	180					
194	68	137	281	237	239	239	228	227	87	71	201					
172	106	207	238	233	214	220	239	228	98	74	206					
188	88	179	209	186	215	211	158	139	78	20	169					
189	97	166	84	10	168	134	11	31	62	22	148					
199	168	191	193	158	227	178	143	182	106	36	190					
206	174	156	262	236	231	149	178	228	43	95	234					
190	216	116	149	226	187	85	150	79	38	218	241					
190	224	147	108	227	216	127	103	36	101	255	224					
190	214	173	66	103	143	96	80	2	109	249	215					
187	196	236	75	1	81	47	0	6	217	288	211					
183	202	237	146	0	0	12	108	200	128	243	236					
195	206	123	207	177	121	123	200	175	13	95	218					

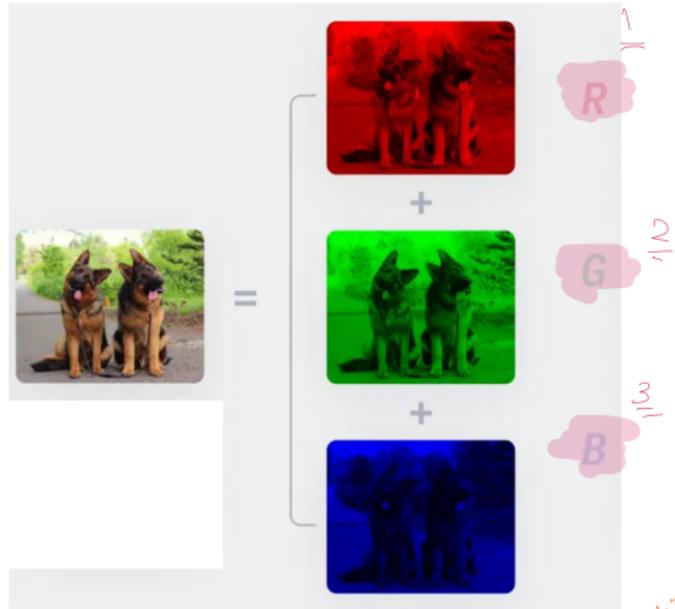
157	153	174	168	150	162	129	161	172	161	155	166					
155	182	163	74	75	62	33	17	110	210	180	164					
180	180	50	14	34	6	10	33	48	106	159	181					
206	109	5	124	131	111	120	204	166	15	56	180					
194	68	137	281	237	239	239	228	227	87	71	201					
172	105	207	233	214	220	239	228	98	74	206						
188	88	179	209	186	215	211	158	139	78	20	169					
189	97	166	84	10	168	134	11	31	62	22	148					
199	168	191	193	158	227	178	143	182	106	36	190					
206	174	156	262	236	231	149	178	228	43	95	234					
190	216	116	149	226	187	85	150	79	38	218	241					
190	224	147	108	227	216	127	103	36	101	255	224					
190	214	173	66	103	143	96	80	2	109	249	215					
187	196	236	75	1	81	47	0	6	217	288	211					
183	202	237	146	0	0	12	108	200	128	243	236					
195	206	123	207	177	121	123	200	175	13	95	218					

⁰<https://www.v7labs.com/blog/image-recognition-guide>

How to represent an image?

- ▶ RGB images have 3 channels.

CMYK
not → HSL or HCB
sure



Normalization:
pick the highest value &
make the rest close to it
between (0,1)

How can we build vision models?

- ▶ So far, we've worked with tabular data, where each sample consists of structured features. We processed this data using Fully Connected Neural Networks.

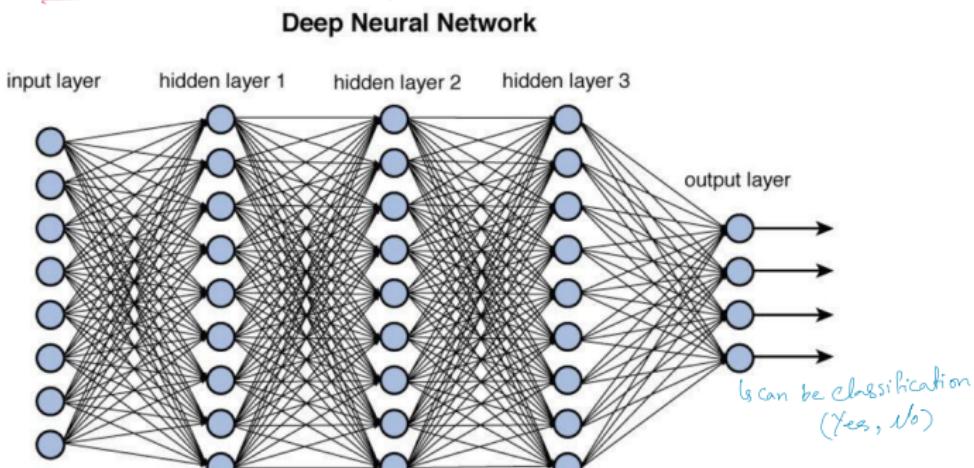


Figure 12.2 Deep network architecture with multiple layers.

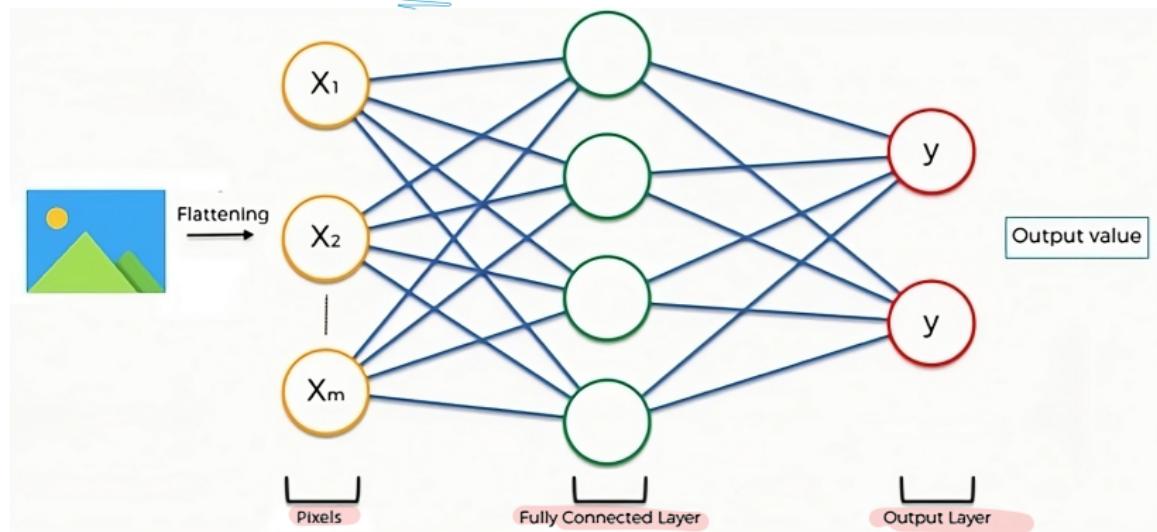
- ▶ But can we apply the same approach to images?

How can we build vision models?

- ▶ Let's try.

How can we build vision models?

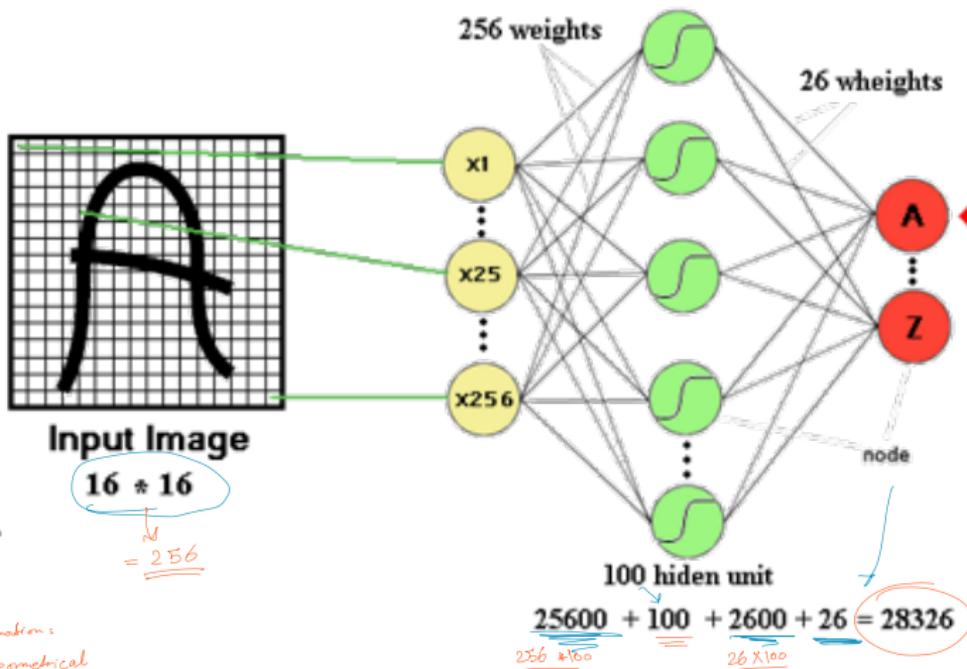
- ▶ Let's consider each pixel as a feature.



- This should work! But...

Drawbacks of Fully-Connected Neural Networks

- ▶ The number of trainable parameters becomes extremely large



Affine transformation:
Scaling any geometrical

Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ Little or no invariance to shifting, scaling, and other forms of distortion
- ▶ In other words, the Fully connected NN cannot recognize that the two images (original and shifted) represent the same object (letter "A").

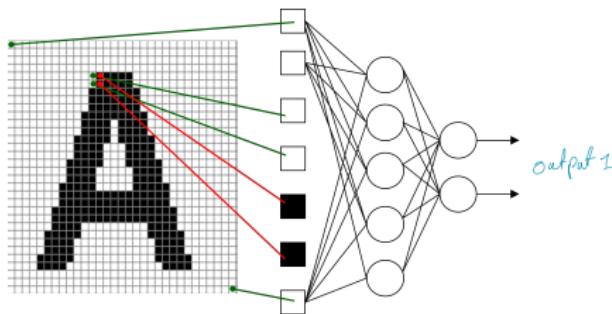


Figure 2: Original "A" character.

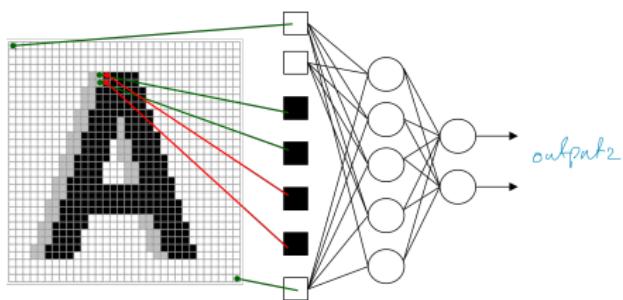


Figure 3: Shifted "A" character.

$$\text{output 1} \neq \text{output 2}$$

Drawbacks of Fully-Connected Neural Networks (cont.)

- ▶ The topology of the input data is completely ignored (it treats pixels as independent features rather than structured patterns.)
- ▶ For a 32×32 image, we have
 - Black and white patterns: $2^{32 \times 32} = 2^{1024}$ pattern.
 - Grayscale patterns: $256^{32 \times 32} = 256^{1024}$ pattern.



What do we need?

► We need a model that can:

- **Find Patterns in Images:** Recognize small features like edges and shapes in different parts of the image.
- **Work Efficiently:** Avoid looking at every pixel separately by focusing on groups of pixels together.
- **Handle Changes:** Still recognize the same object even if it moves, rotates, or looks slightly different.

What do we need?

- ▶ We need a model that can:

- **Find Patterns in Images:** Recognize small features like edges and shapes in different parts of the image.
- **Work Efficiently:** Avoid looking at every pixel separately by focusing on groups of pixels together.
- **Handle Changes:** Still recognize the same object even if it moves, rotates, or looks slightly different.
- We need **Convolutional Neural Networks (CNNs)!**

Convolutional Neural Networks (CNNs)

- ▶ CNNs are neural networks designed for image processing and consist of two main parts:
 - **Feature Extractor:** Automatically learns patterns (features) such as edges, textures, and shapes from the image.
 - **Classifier:** The extracted features are flattened into a vector and passed to a standard neural network (like the ones we used before) for classification.

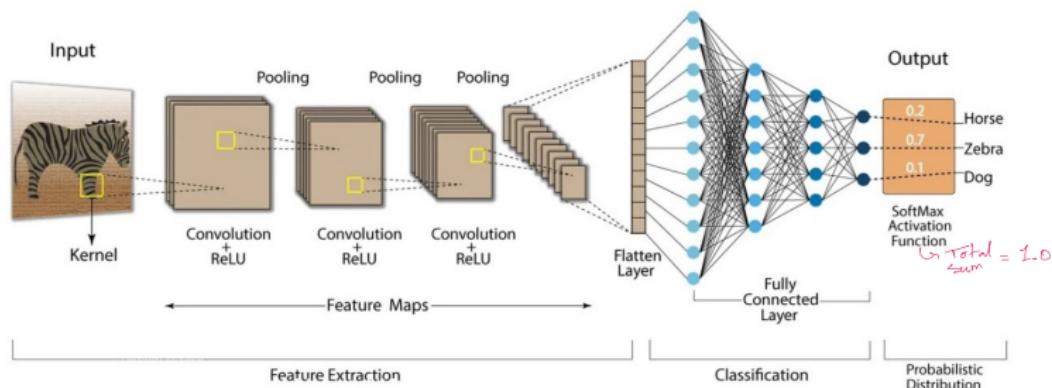


Figure 4: Convolutional Neural Network

Convolutional Neural Networks (CNNs)



- ▶ The feature extractor consists of three essential components:
 - **Convolution Layers:** Detect local patterns such as edges, textures, and shapes by applying filters to the image.
 - **Activation Function (ReLU):** Adds non-linearity.
 - **Pooling Layers:** Reduce the spatial size of extracted features.
- ▶ Let's start with the convolution layer.

How Convolution Works?

- ▶ Let's consider this image.



How Convolution Works? (cont.)

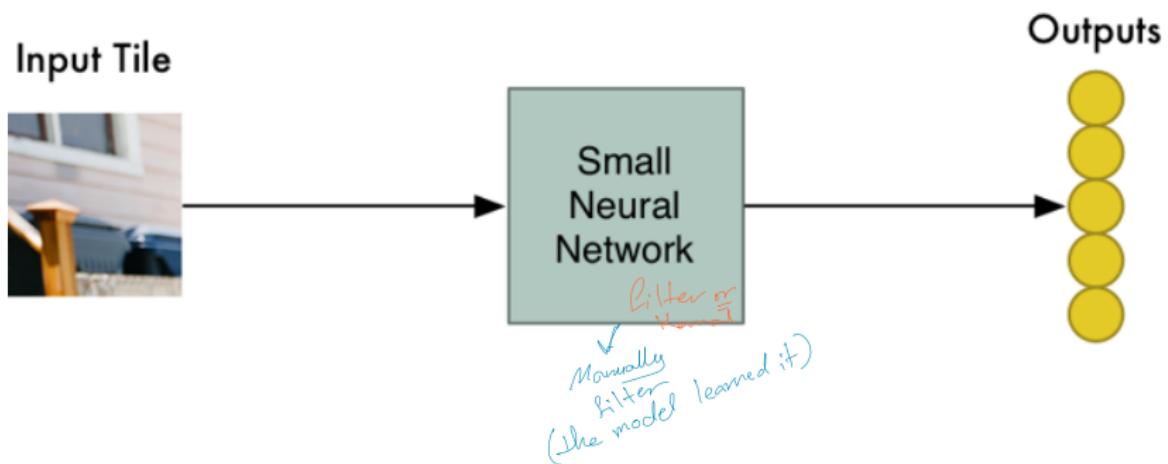
1. The image is divided into small overlapping tiles (regions).



How Convolution Works? (cont.)

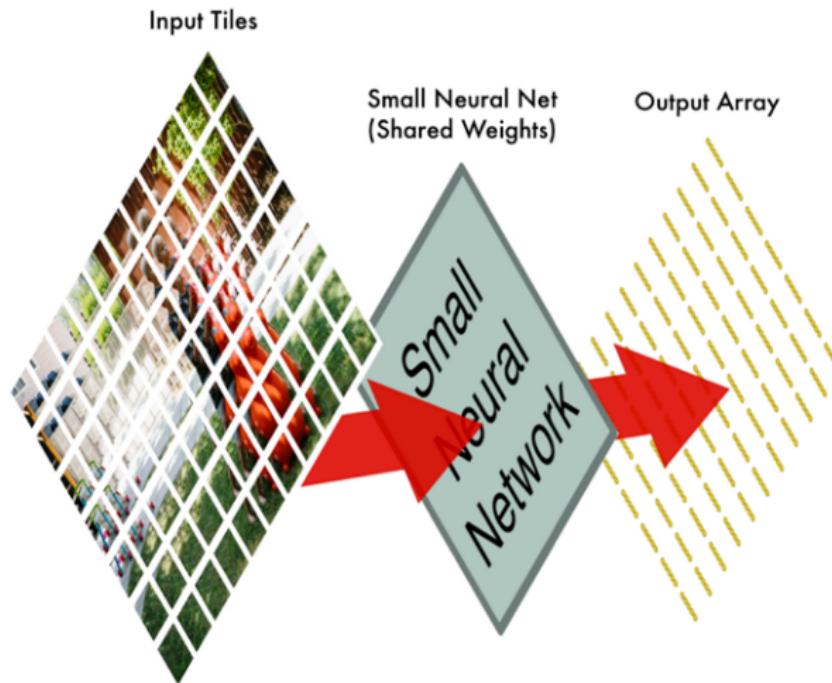
2. We process each tile using the weights matrix of a small neural network. We call this matrix a kernel or filter.

Processing a single tile



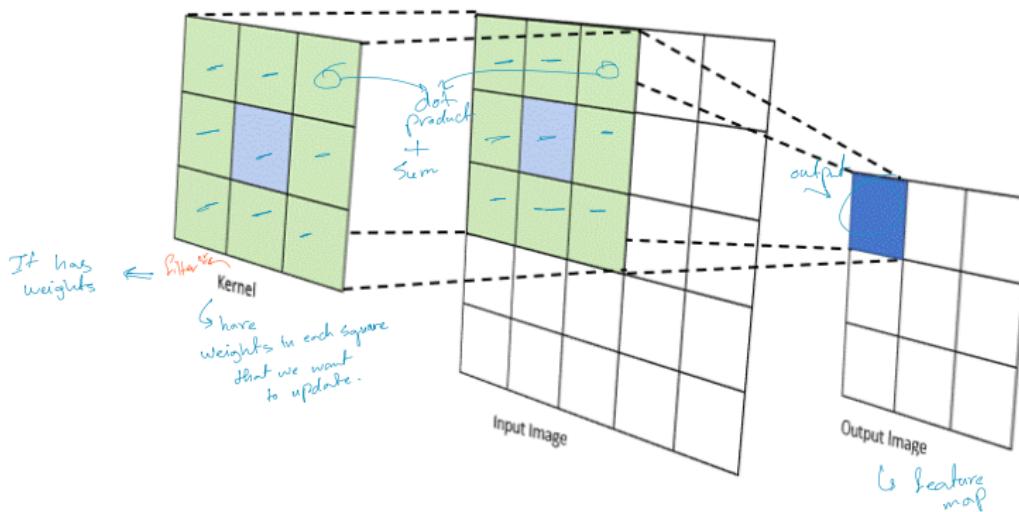
How Convolution Works? (cont.)

3. Finally, the filter slides across the image (with the same weights) and processes all the tiles, creating an output **feature map**.



How Convolution Works? (cont.)

- ▶ What we did in the last step is called the **Convolution Operation**.
 - ▶ We convolved the kernel with the image, which means sliding the kernel over the entire image and computing the **dot product** between the kernel and small regions of the image at each step.



How Convolution Works? (cont.)

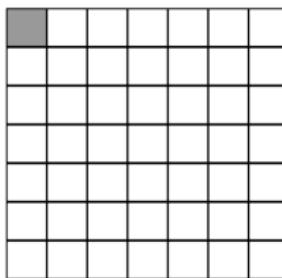
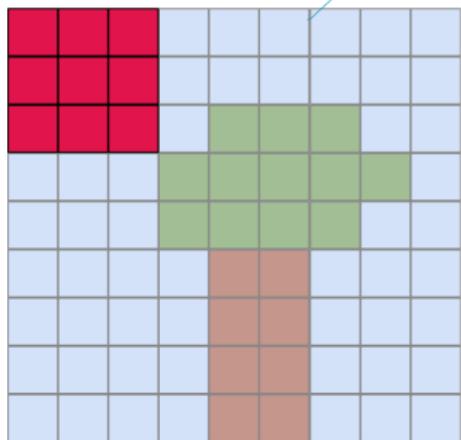
- ▶ This can be represented mathematically by:

$$z = W * x_{i,j} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} W_{ab} X_{(i+a)(j+b)}$$

- z : Output of the convolution at (i, j) .
 - W : Filter (kernel) matrix of size $m \times n$. *Sobel filter : detect edges.*
Blur //
 - $x_{i,j}$: Input value at position (i, j) .
 - W_{ab} : Weight of the filter at (a, b) .
 - $X_{(i+a)(j+b)}$: Input value in the small region at $(i + a, j + b)$.
- You can set it as you like but it must get beyond image size. It can be 1x1*

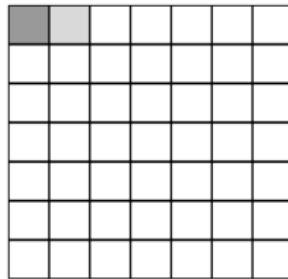
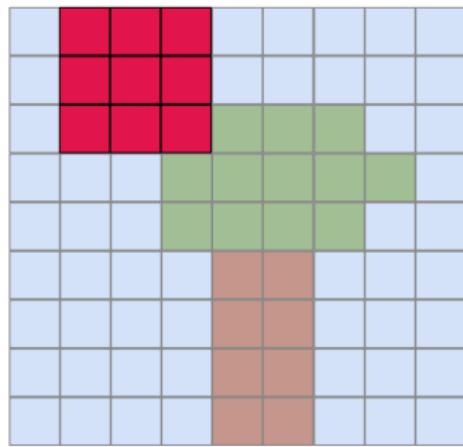
How Convolution Works? (cont.)

You can control
the slides here is
1x1 pixel more
it can be more.



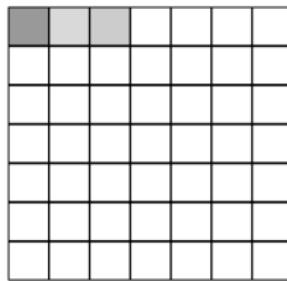
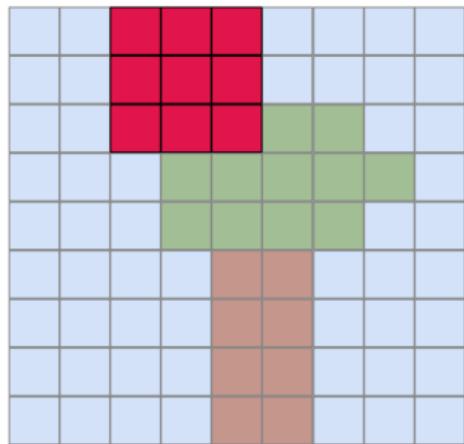
The **kernel** slides across the image and produces an output value at each position

How Convolution Works? (cont.)



The **kernel** slides across the image and produces an output value at each position

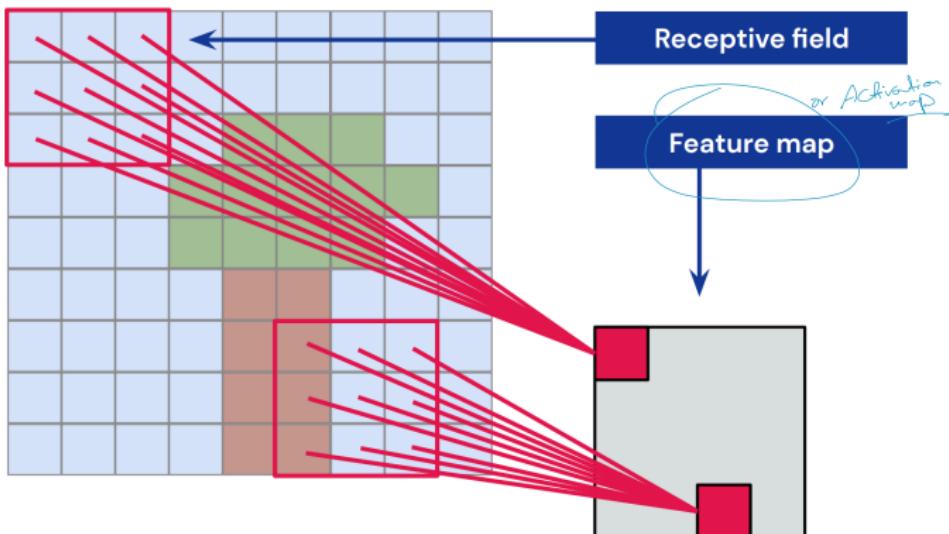
How Convolution Works? (cont.)



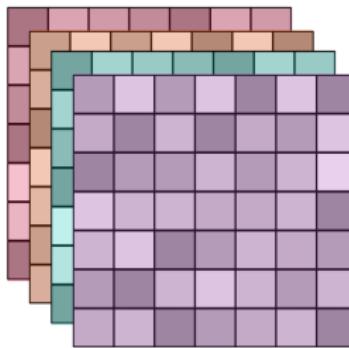
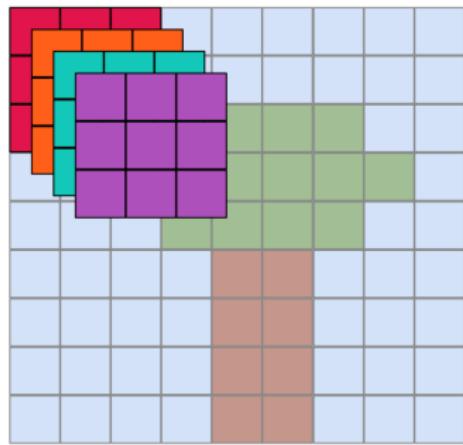
The **kernel** slides across the image and produces an output value at each position

How Convolution Works? (cont.)

- ▶ **Receptive Field:** The region of the input image that the kernel operates on at each step.
 - ▶ **Feature Map:** The features extracted from the input image.



How Convolution Works? (cont.)



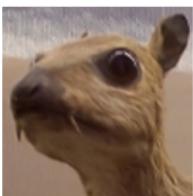
We convolve multiple kernels and obtain multiple feature maps or **channels**

How Convolution Works? (cont.)

- **Filters Example:** Different filters detect patterns like edges, textures, or smoothness, producing corresponding feature maps.

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

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

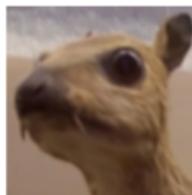


$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

Sharp ↗

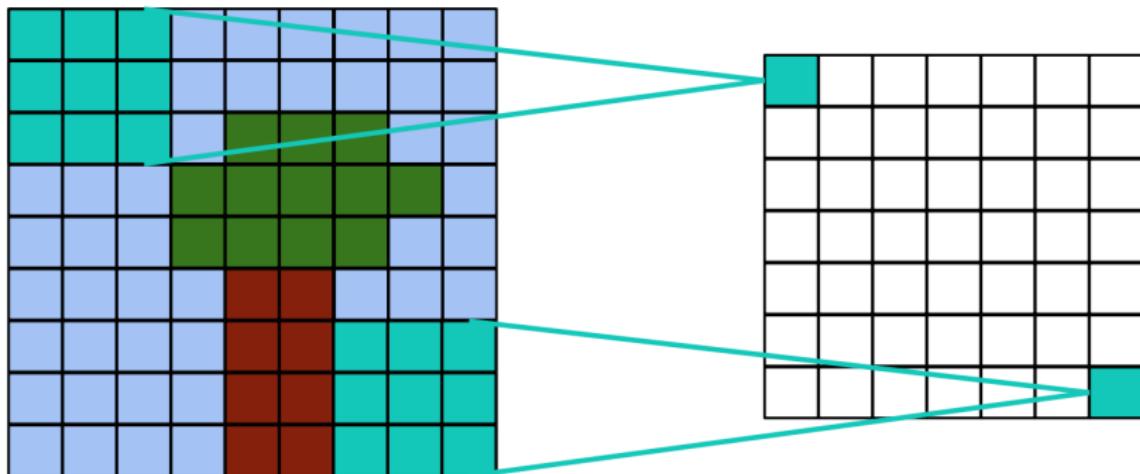


Blur → Smoothness
Filter

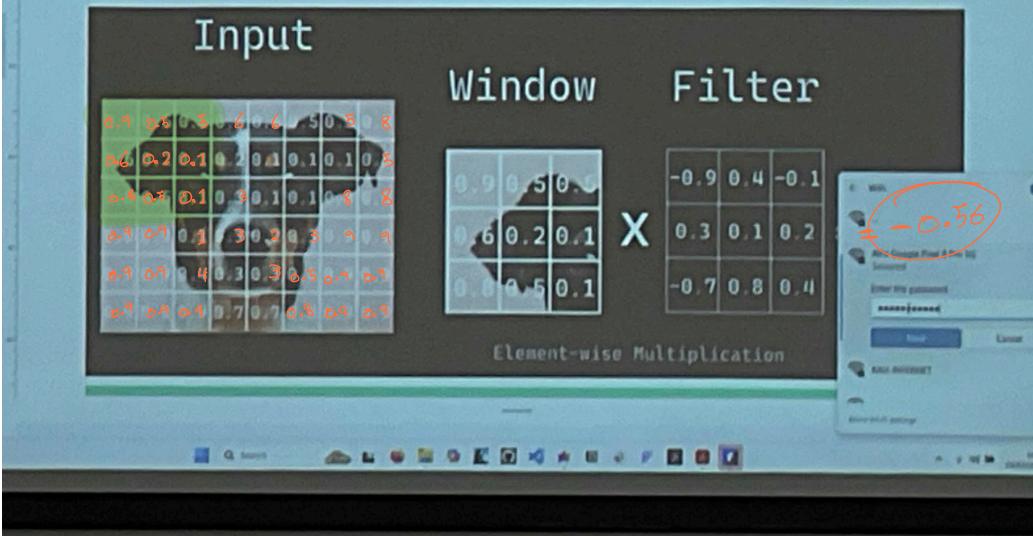
Controlling the Convolution Process

- ▶ Applying Convolution as such reduces the size of the borders.
 - ▶ Sometimes this is not desirable. *but pooling is minimizing more.*

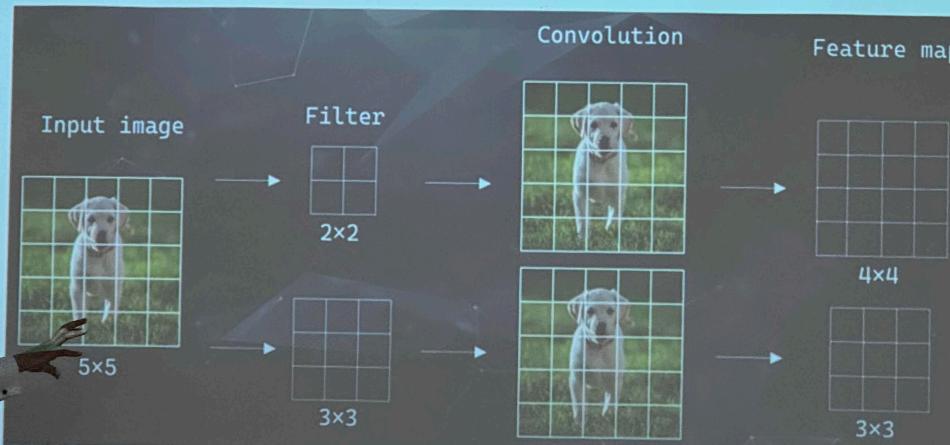
but pooling is minimizing more.



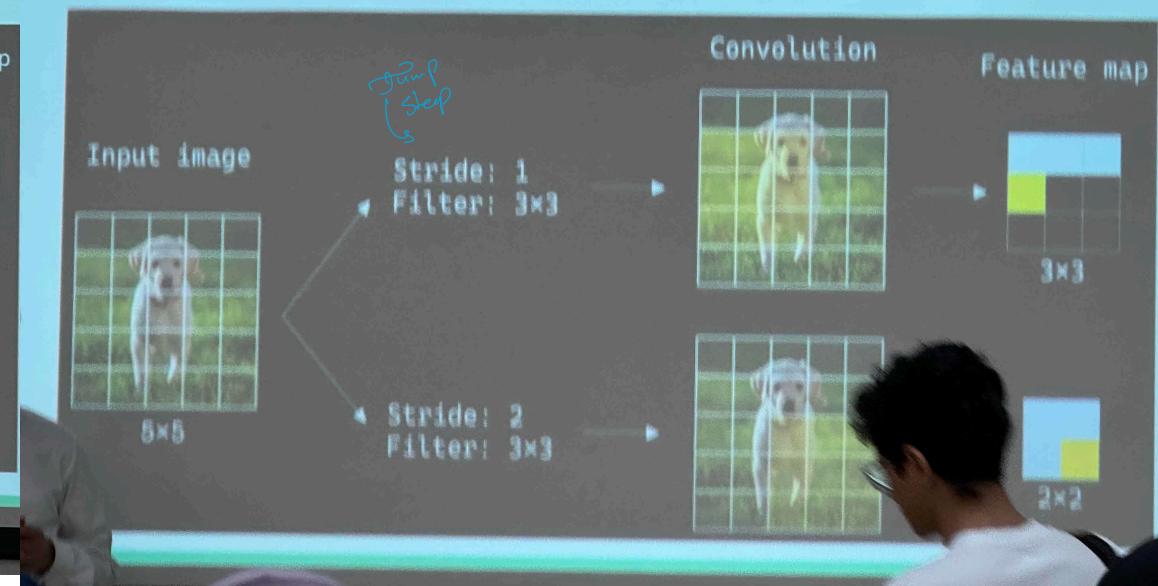
Convolution - Element-wise Multiplication (Scalar Product)



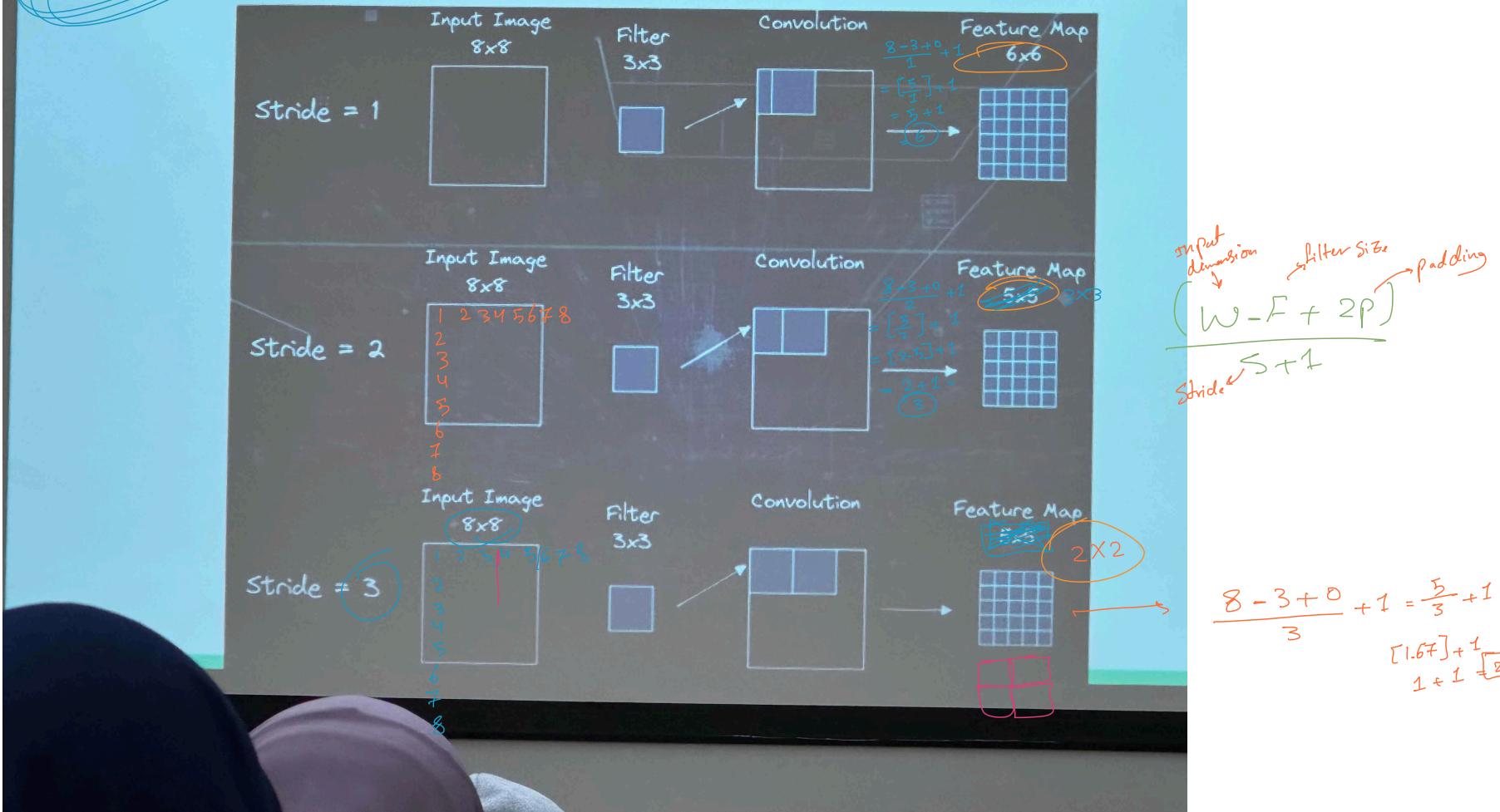
Different Filter Size



Stride - How quickly the filter moves over the image



Stride - Number of pixels the filter moves



~~Important~~
question here

Controlling the Convolution Process (cont.)

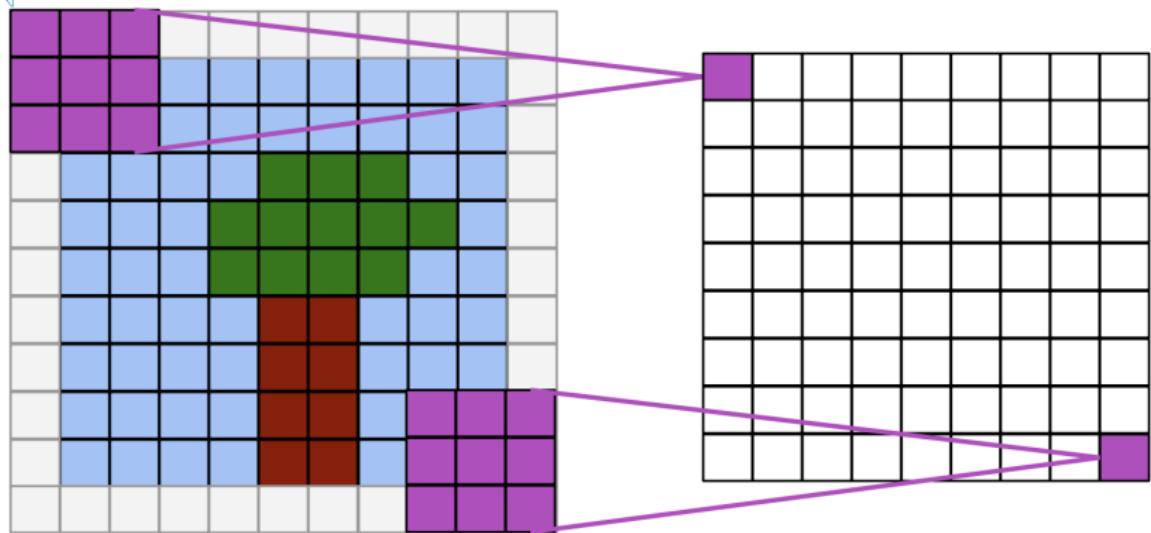
- Solution: We can use **Padding**, which pads the border with zeros.
- it has two types:

1. **Same Convolution:** Padding is added so the output size equals the input size.

get the
same size
in feature
map

مُفْعَلٌ فِي الْجَوَافِرِ
كُلُّهُ وَهُنَّ مُسَعِّدُونَ بِالْمَوْلَى الْمُجَاهِدِ

make output size = input size
useful when images are small dimensions.



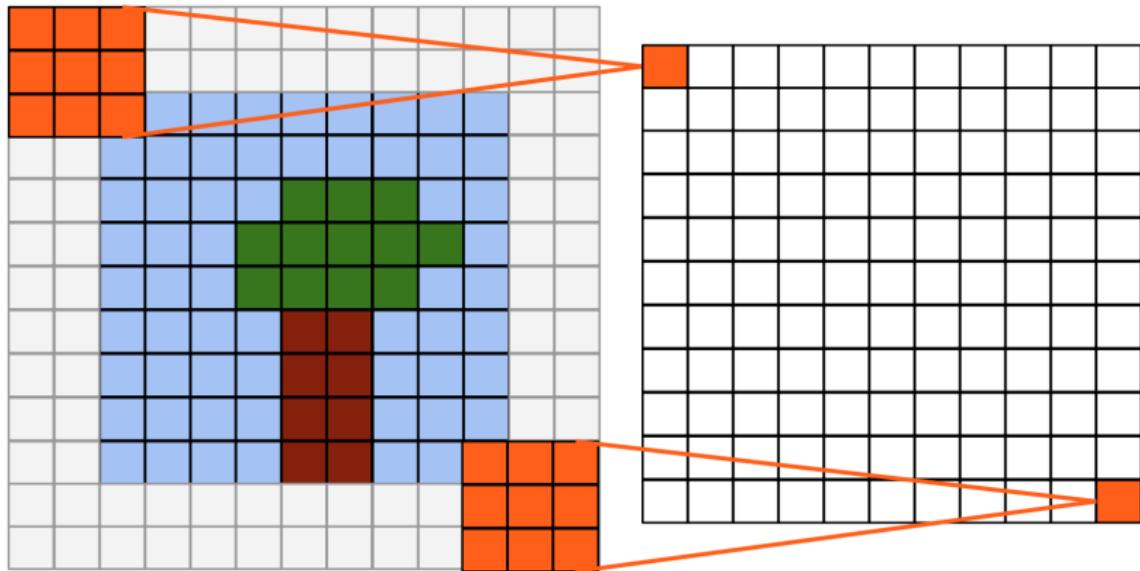
Controlling the Convolution Process (cont.)

Increase the size of feature map

2. **Full Convolution:** Padding input, including edges.
output size = input size

2. **Full Convolution:** Padding is added so the kernel covers the entire input, including edges.

output size = input size + kernel size - 1

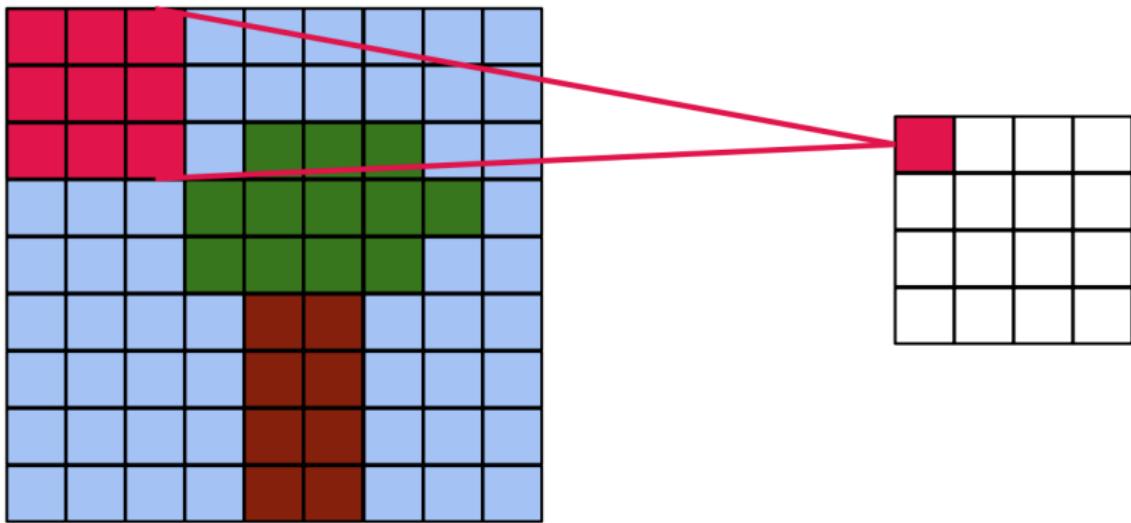


Controlling the Convolution Process (cont.)



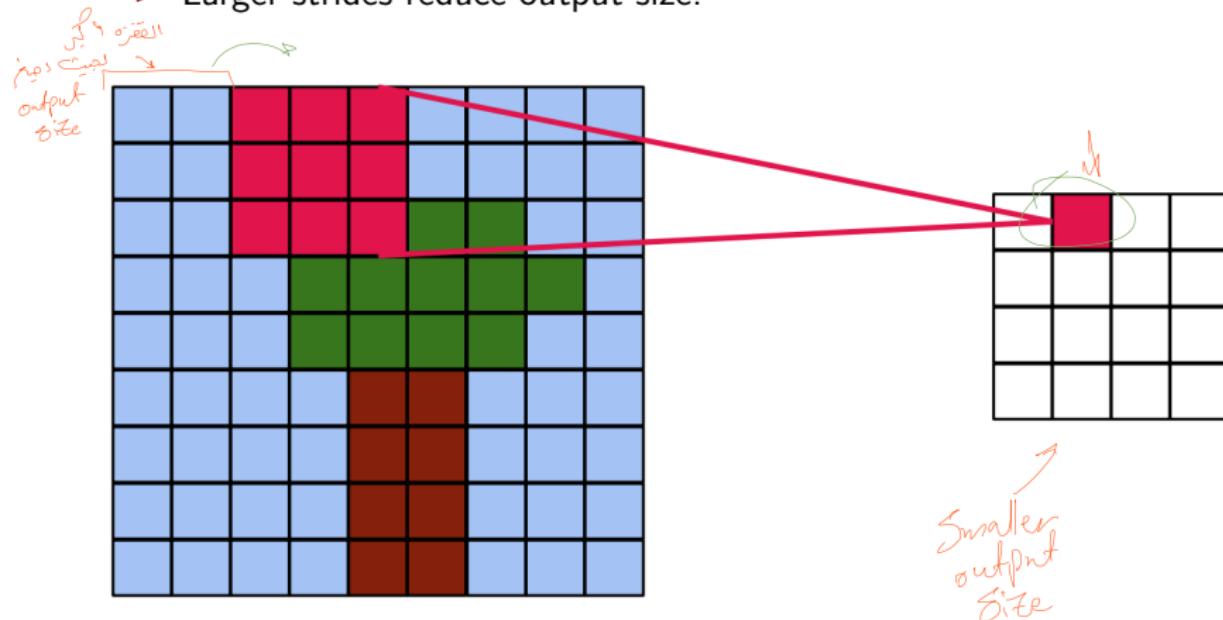
لِعَنْ

- ▶ **Strided Convolution:** Kernel slides along the image with a step > 1
 - ▶ Larger strides reduce output size.



Controlling the Convolution Process (cont.)

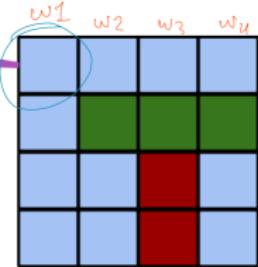
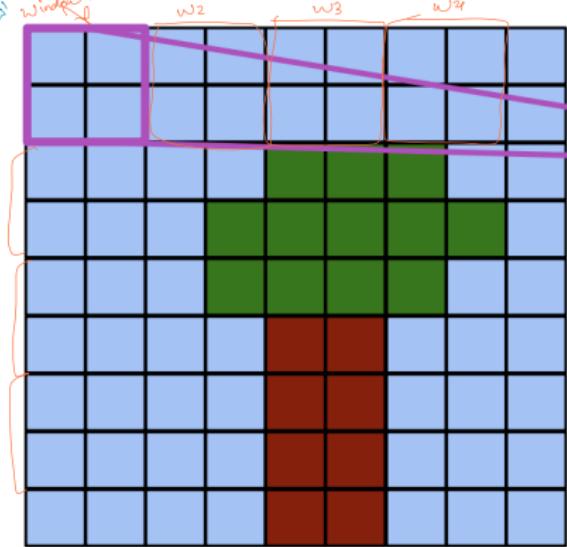
- ▶ **Strided Convolution:** Kernel slides along the image with a step > 1
 - ▶ Larger strides reduce output size.



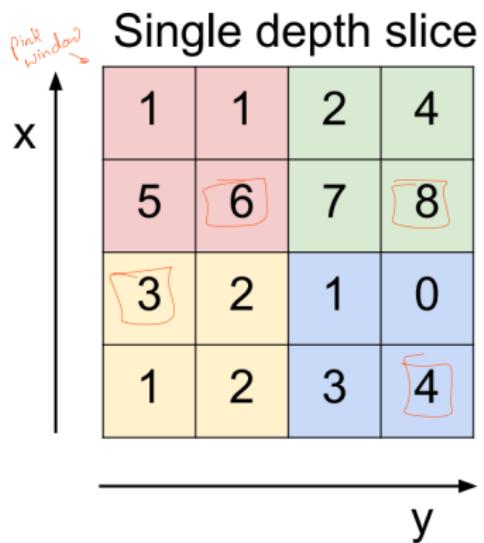
Controlling the Convolution Process (cont.)

Conv layer
ReLU
Batch
Norm
max is it w
mean
for max
Gauss Window
Gauss Window
Gauss Window

▶ **Pooling:** Compute mean or max over small windows to reduce resolution and extract more general features.



Controlling the Convolution Process (cont.)



why max 2,
because it holds the
most of the information
in the single window.

max pool with 2x2 filters
and stride 2

6	8
3	4

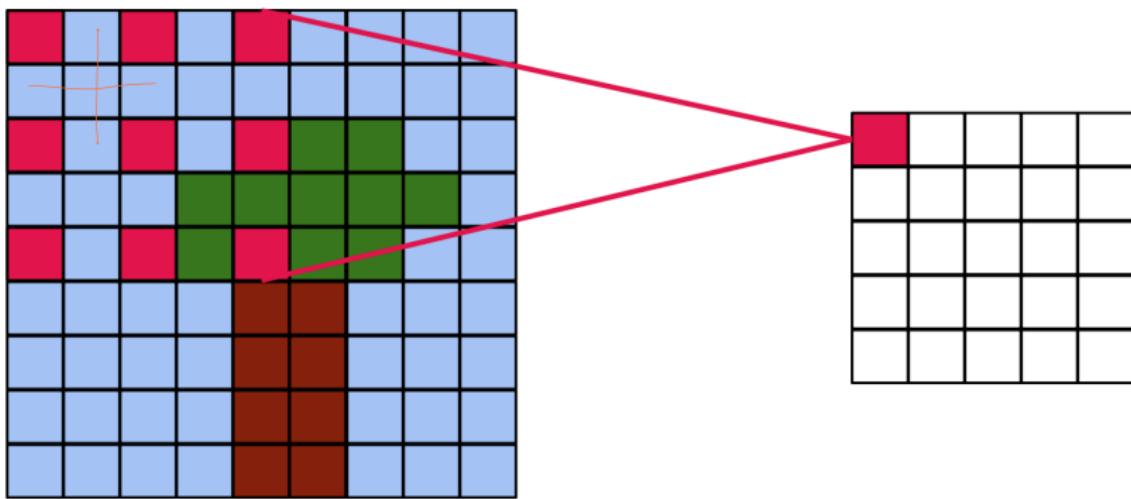
- No learnable parameters
- Introduces spatial invariance

Min Pooling: take the least important features [your network will not learn anything]
network will collapse very fast

Controlling the Convolution Process (cont.)

- ▶ **Dilated Convolution:** Kernel is spread out, step > 1 between kernel elements.
 - ▶ It expands the receptive field without increasing the number of parameters, which makes it efficient.

الفقرة تختلف



CNN Output size

the feature map size

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

n_{in} : number of input features

n_{out} : number of output features

k : convolution kernel size *the window size*

p : convolution padding size

s : convolution stride size

$s=2$ mostly

- ▶ Just like Fully-Connected Neural Networks, we can apply an activation over convolutional layer outputs
- ▶ It helps break linearity
- ▶ For example, Rectified Linear Unit (ReLU): $\sigma(x) = \max(0, x)$

Feature Map

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

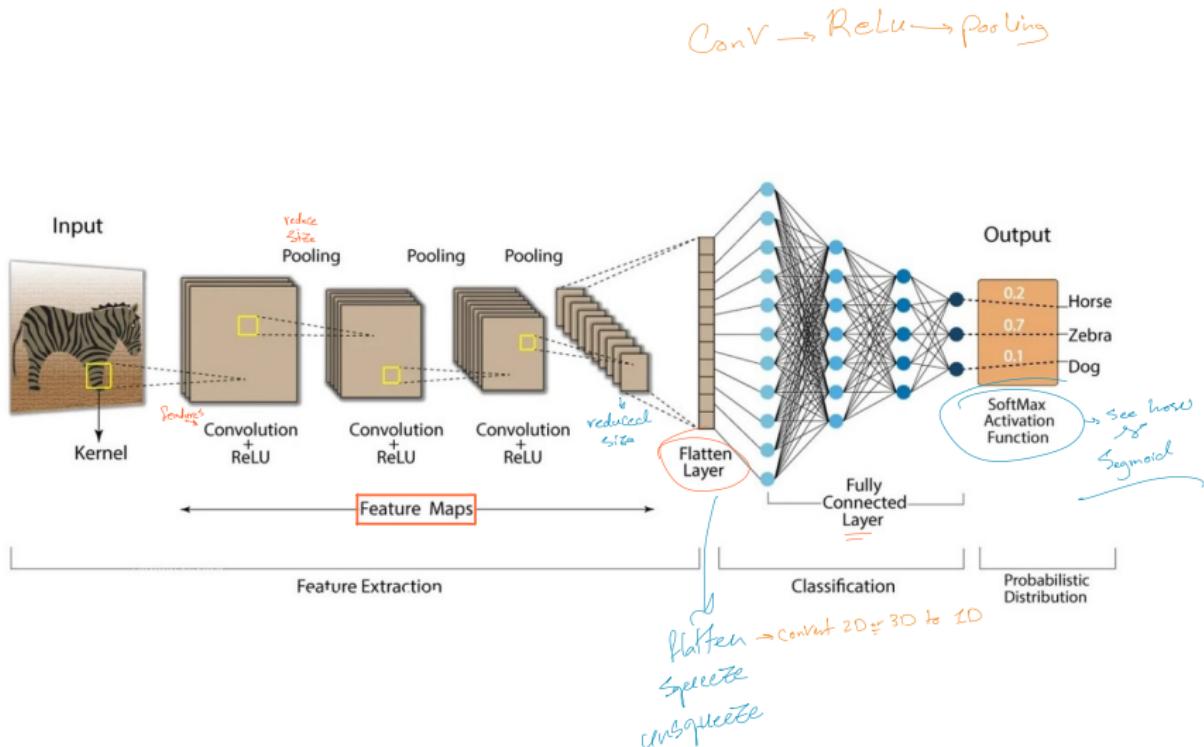
for non-linearity

ReLU Layer

negative numbers \rightarrow because Zeros

9	3	5	0
0	2	0	1
1	3	4	1
3	0	5	1

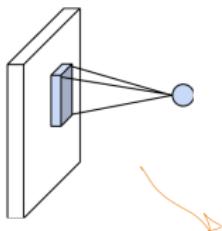
Convolutional Neural Networks



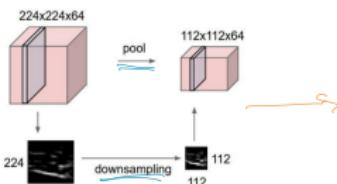
Components of a CNN



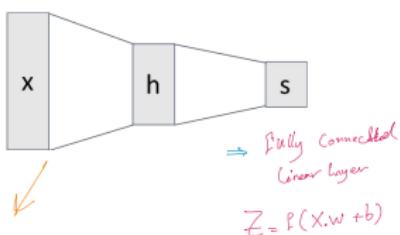
Convolution Layers



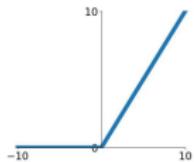
Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

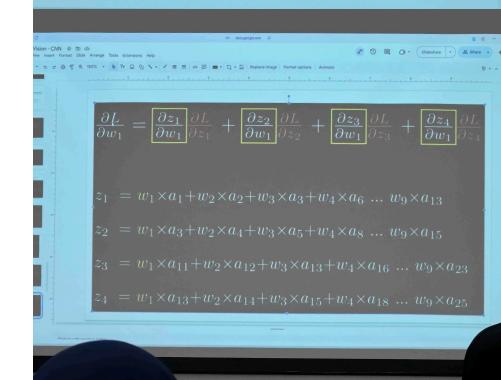
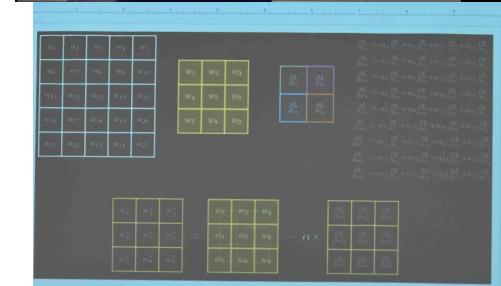
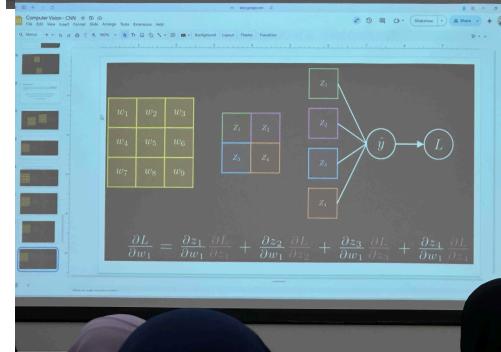
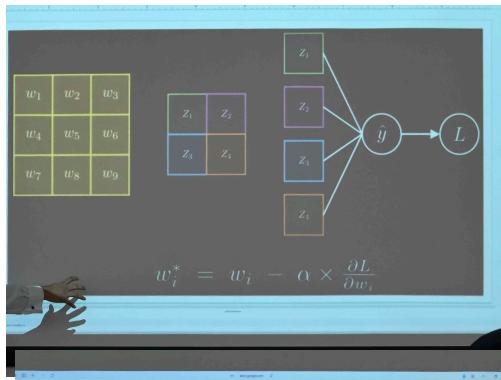
Handwritten notes: alpha, epsilon

See it

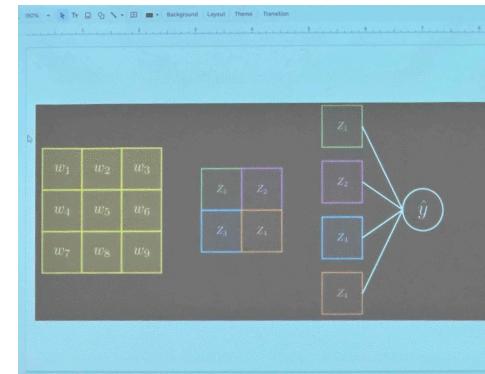
* what we want when we use back propagation:
take the loss and do chain rule to take the derivatives
for each layer

* Back propagation:

- go back with partial derivative
- use chain rule



* Forward propagation



Most Notable CNN-based Architectures

* important
in real-world
no one write CNN
from scratch alone

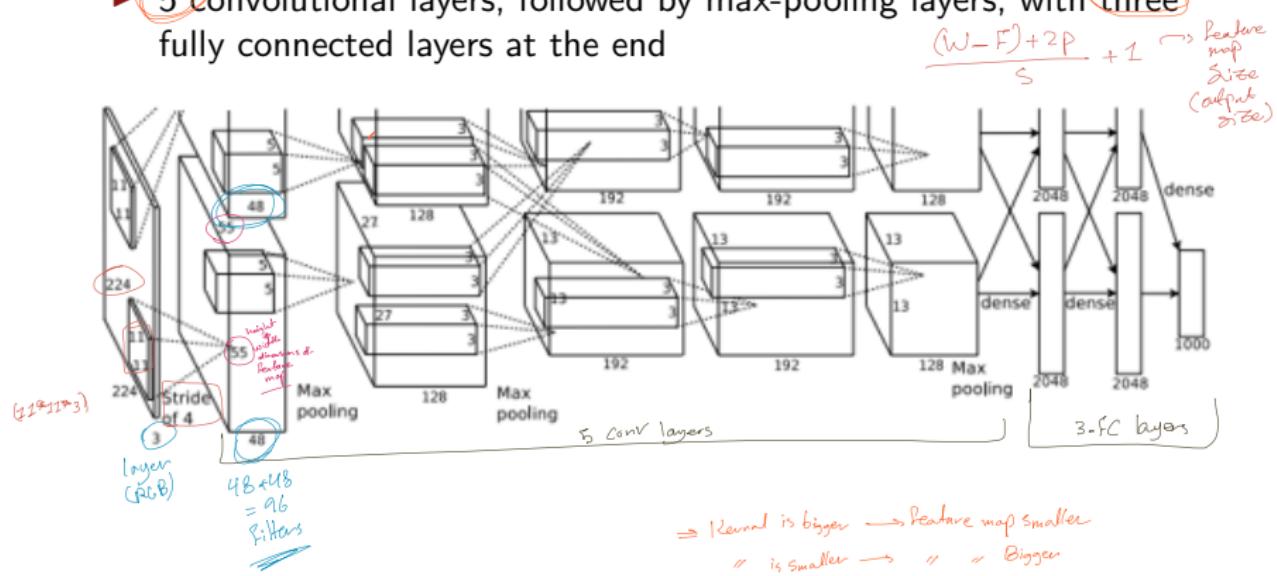
► Over time, researchers built advanced CNN architectures to improve performance and efficiency. These architectures introduced key innovations:

- **AlexNet [Krizhevsky et al. 2012]**: The first CNN to achieve breakthrough performance on image classification.
- **VGGNet [Simonyan and Zisserman, 2014]**: Used very deep networks (up to 19 layers).
- **InceptionNet (GoogLeNet) [Szegedy et al., 2014]**: Used multiple filter sizes per layer (Inception modules).
- **ResNet [He et al., 2015]**: Introduced skip connections for training very deep networks.
- **EfficientNet [Tan and Le, 2019]**: Found a scaling method that simultaneously scales a CNN's depth, width, and resolution optimally using a single scaling coefficient.



Trained on 1000 object categories

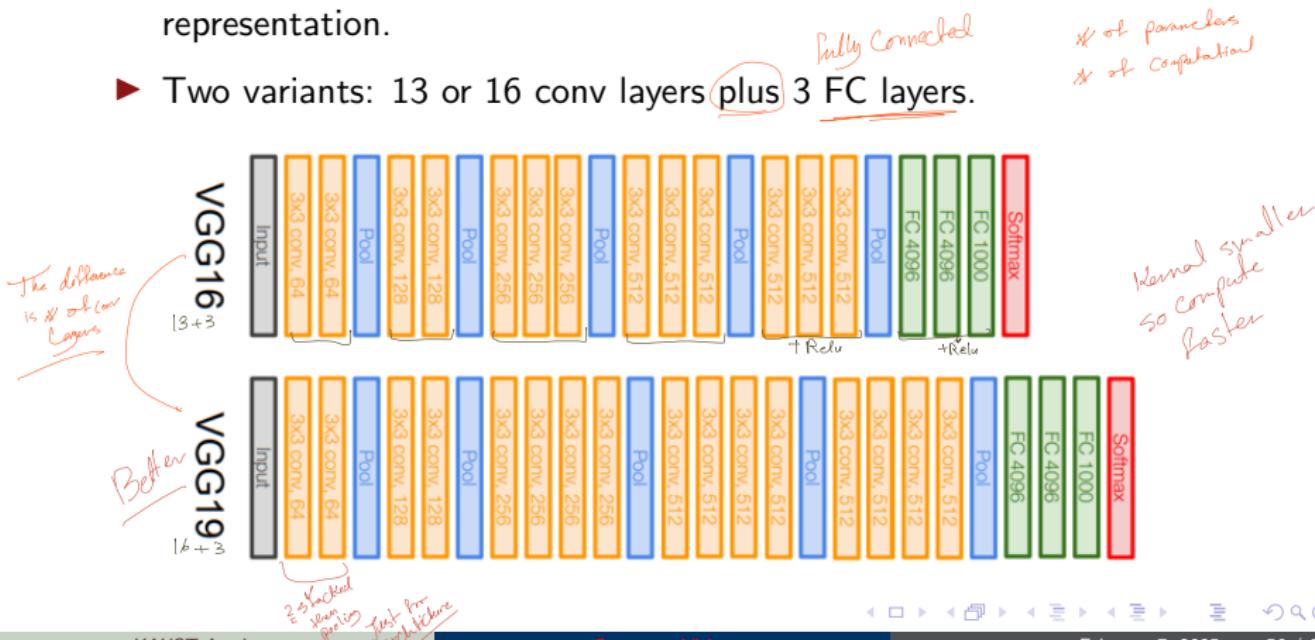
- ▶ First big improvement in image classification.
→ to avoid overfitting
- ▶ Made use of CNN, pooling, dropout, ReLU and training on GPUs.
- ▶ **5 convolutional layers, followed by max-pooling layers; with three fully connected layers at the end**
= 8 layers





► **Improvement over AlexNet:** Uses a deeper network with small filters instead of a shallow network with larger filters.

- A stack of 3x3 conv layers (vs. 7x7 conv) has same receptive field, more non-linearities, fewer parameters and deeper network representation.
- Two variants: 13 or 16 conv layers plus 3 FC layers.



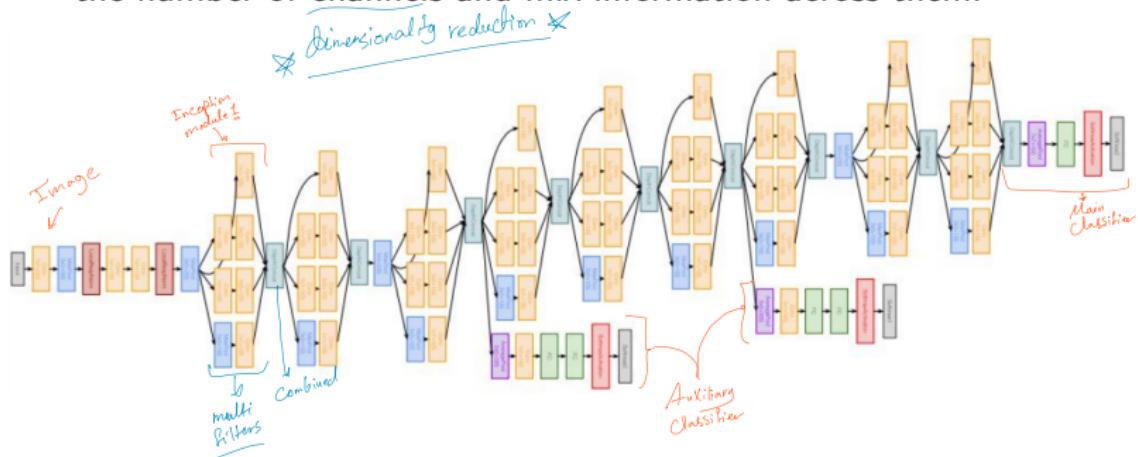
InceptionNet

• Wider rather than deeper.

→ so more parameters

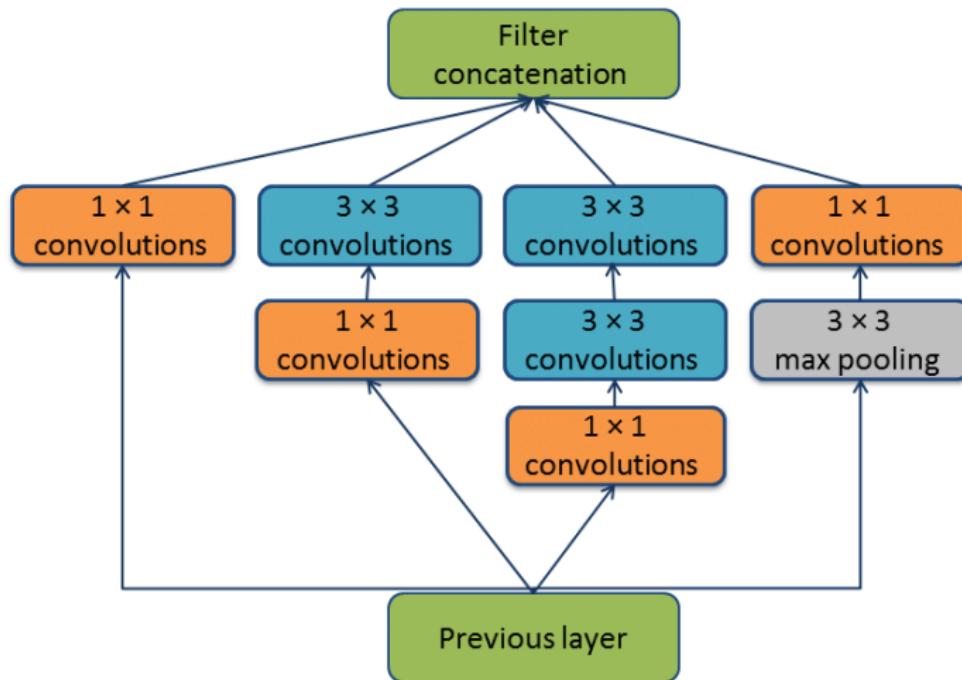
Performance game | Accuracy game
That makes the differences.

- Going Deep: 22 layers
 - Only 5 million parameters! ($12\times$ fewer than AlexNet, $27\times$ fewer than VGGNet).
 - Introduced efficient "Inception module"
 - Introduced "bottleneck" layers that use 1×1 convolutions to reduce the number of channels and mix information across them.



InceptionNet (cont.)

- ▶ **Inception module:** Uses multiple filter sizes (1×1 , 3×3 , 5×5), in parallel, to capture different features, then combines their outputs.

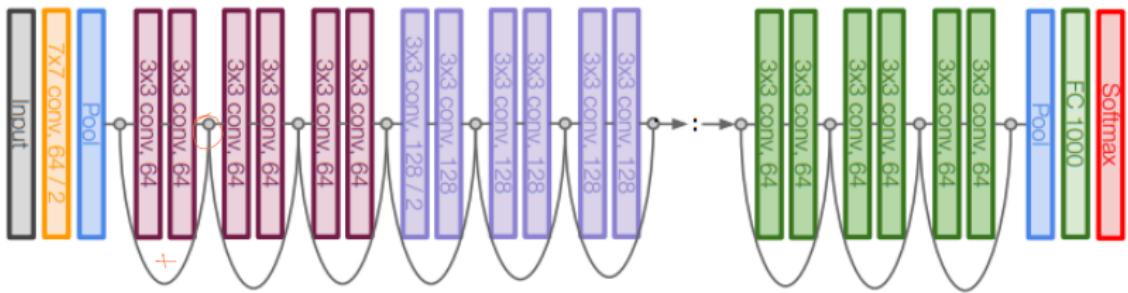


~~↙ better~~

- ▶ **Problem:** Making networks deeper does not always improve accuracy.
- ▶ **Why?** In very deep networks, gradients become extremely small as they move backward through layers, making learning slow or stopping it altogether (**vanishing gradient problem**).
- ▶ **Solution:** Residual Network (ResNet) introduces **skip connections (residuals)**, allowing information to flow more easily.

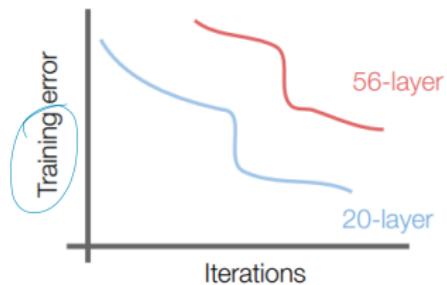
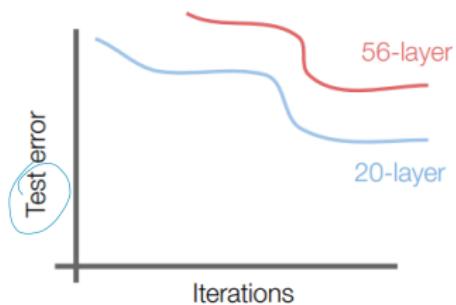


- ▶ Very deep networks using residual connections
 - ▶ 152-layer model for ImageNet
 - ▶ Stacked Residual Blocks
 - ▶ **Residual:** A shortcut connection that helps the network pass information through layers more easily.



٦٠ معلوم قیمه و اینها همچو
نه اسیا در گردنی که \rightarrow So the gradient not vanishing.
این معنی ندارد

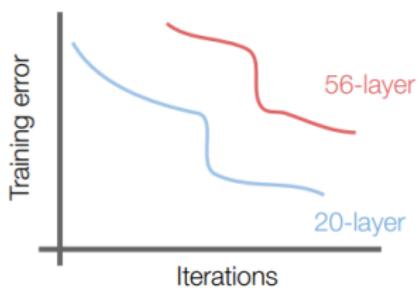
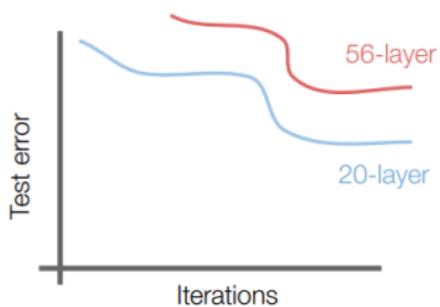
- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



- ▶ 56-layer model performs worse on both test and training error

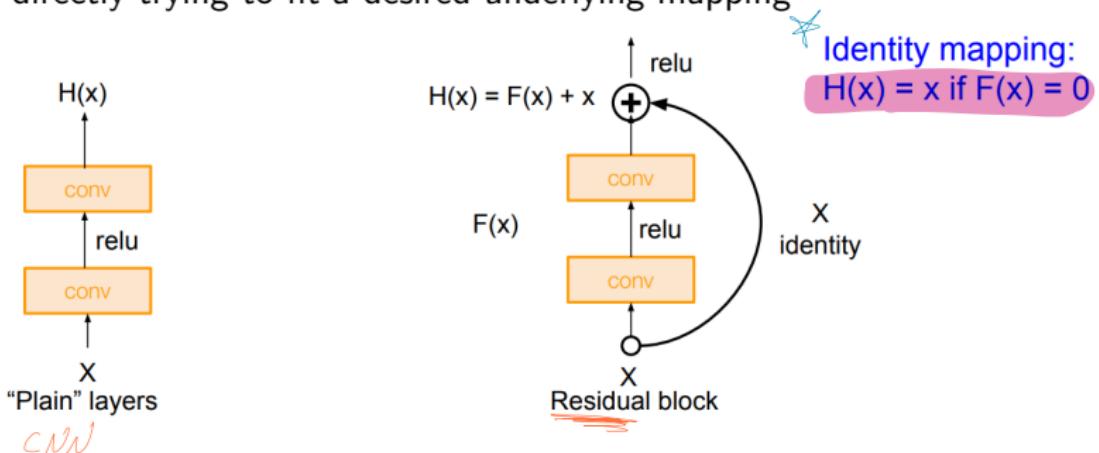
more layers might not be good sometimes because you might not have enough images to train or so.

- ▶ What happens when we continue stacking deeper layers on a "plain" convolutional neural network?



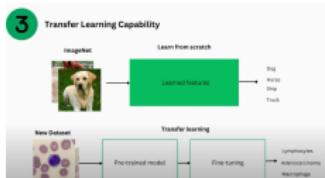
- ▶ 56-layer model performs worse on both test and training error
- ▶ The deeper model performs worse, but it's not caused by overfitting!

- ▶ **Fact:** Deep models have more representation power (more parameters) than shallower models.
- ▶ **Hypothesis:** The problem is an optimization problem, deeper models are harder to optimize
- ▶ **Solution:** Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping

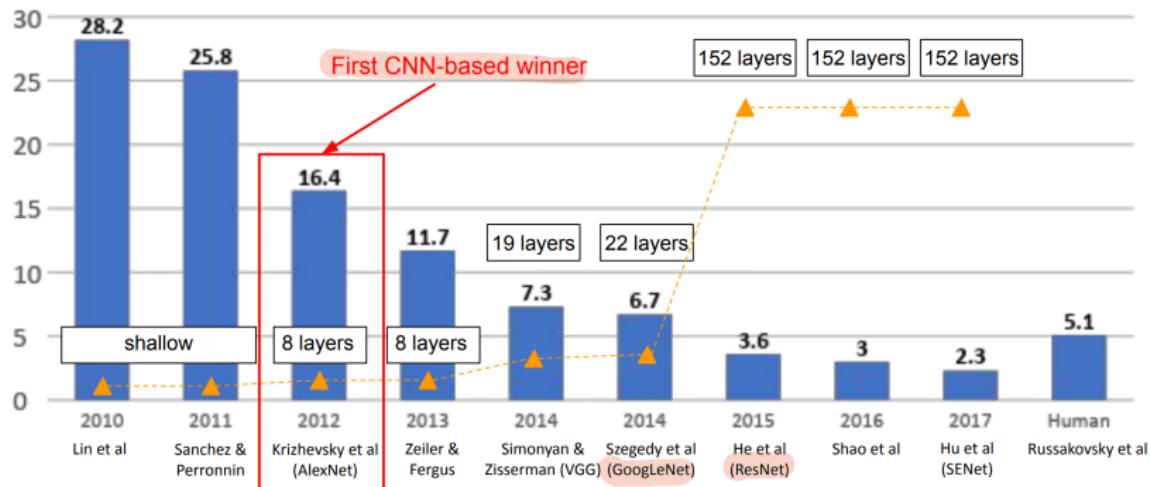




- ▶ The most extensive data for Image Classification , object recognition
- ▶ 3 RGB channels from 0 to 255
 - 15 million labeled in over 22,000 categories
- ▶ 14,197,122 images
- ▶ 1000 classes



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



► **Problem:** To improve accuracy, we can:

- Increase the number of layers (**depth**)
- Increase the number of neurons in each layer (**width**)
- Use higher-resolution images (**resolution**)

But finding the right balance of these three was largely based on trial and error.

► **Solution:** EfficientNet introduced a **compound scaling** method—a mathematical formula to systematically find the optimal balance among **depth**, **width**, and **resolution**.

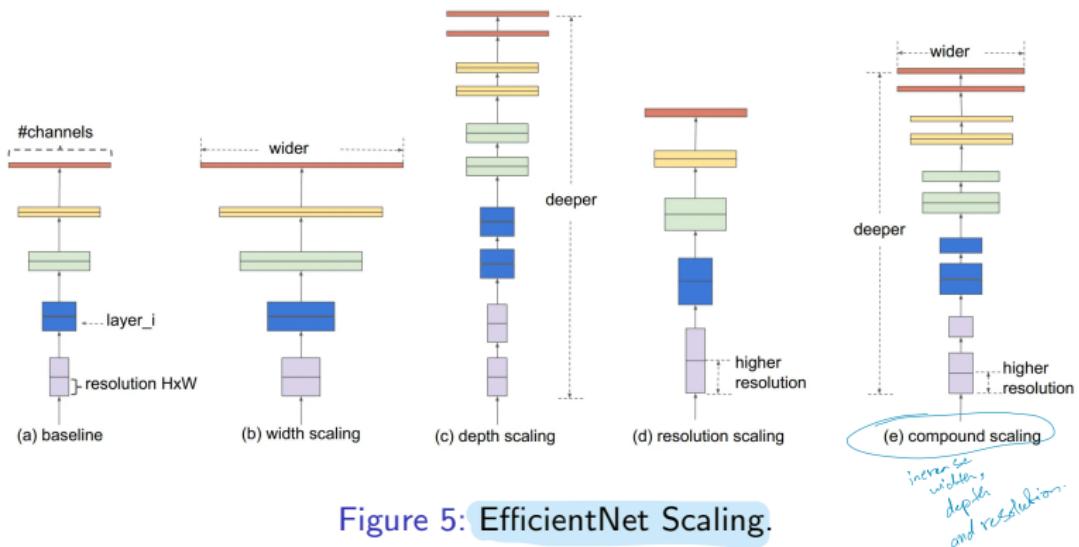


Figure 5: EfficientNet Scaling.

- **Compound Scaling** Uses a single **scaling coefficient** (ϕ) to control:
- **Network Depth** (α^ϕ) → More layers
 - **Network Width** (β^ϕ) → More channels per layer
 - **Input Resolution** (γ^ϕ) → Larger input images
- The goal: find α, β, γ that balance accuracy & efficiency, **then scale up optimally by increasing the global coefficient ϕ .**

- ▶ EfficientNet optimizes depth, width, and resolution using this constraint:

$$\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

- ▶ Why this equation?

- Increasing **depth** (α) increases FLOPs **linearly**.
- Increasing **width** (β) increases FLOPs **quadratically** (β^2).
- Increasing **resolution** (γ) increases FLOPs **quadratically** (γ^2).

- ▶ To double total FLOPs, the three factors must be balanced together.

↳ Floating-point
operator per second
(matrix measures the computational
processing speed of model)
~~with
complexity of
efficiency~~

- ▶ The authors of EfficientNet searched for the best scaling factors on a small baseline model.
- ▶ They found:

$$\underline{\alpha = 1.2},$$

$$\underline{\beta = 1.1},$$

$$\underline{\gamma = 1.15}$$

EfficientNet Scaling: B0 to B7

- The EfficientNet family (B0 to B7) is generated using:

$$\text{Depth} = 1.2^\phi, \quad \text{Width} = 1.1^\phi, \quad \text{Resolution} = 1.15^\phi$$

Model	ϕ	Depth (α^ϕ)	Width (β^ϕ)
B0	0	$1.2^0 = 1.0$	$1.1^0 = 1.0$
B1	1	$1.2^1 = 1.2$	$1.1^1 = 1.1$
B2	2	$1.2^2 = 1.44$	$1.1^2 = 1.21$
B3	3	$1.2^3 = 1.73$	$1.1^3 = 1.33$
B4	4	$1.2^4 = 2.07$	$1.1^4 = 1.46$
B5	5	$1.2^5 = 2.49$	$1.1^5 = 1.61$
B6	6	$1.2^6 = 2.99$	$1.1^6 = 1.77$
B7	7	$1.2^7 = 3.58$	$1.1^7 = 1.94$

$$\begin{aligned} & \alpha \cdot \beta^2 \cdot \text{depth}^2 \approx 2 \\ & \text{depth} \quad \text{width} \quad \text{resolution} \\ & 1.2 \cdot (1.1)^2 \cdot (1.15)^2 \approx 2 \end{aligned}$$

Table 1: Scaling EfficientNet from B0 to B7

- We multiply these scaling factors by the baseline EfficientNet-B0 values and round to the nearest integer to get the new depth, width, and resolution for each model.

- ▶ EfficientNet models achieve state-of-the-art accuracy with significantly fewer parameters and FLOPs.
- ▶ EfficientNet-B7 reaches 84.4% Top-1 and 97.3% Top-5 accuracy on ImageNet.
- ▶ More efficient than previous CNN models—8.4x smaller and 6.1x faster than competitors.

EfficientNet Performance

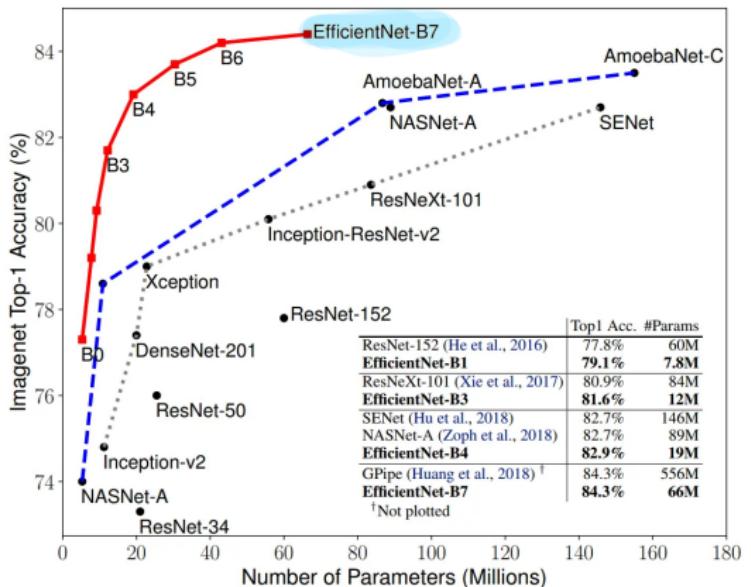


Figure 6: EfficientNet Performance on Imagenet.

▶ Why MobileNets?

- Small-sized models are crucial for mobile and embedded devices.
- MobileNets reduce computational cost and memory usage while maintaining good accuracy.

▶ Key Idea:

- Use **depthwise-separable convolutions** to significantly reduce computation compared to standard convolutions.

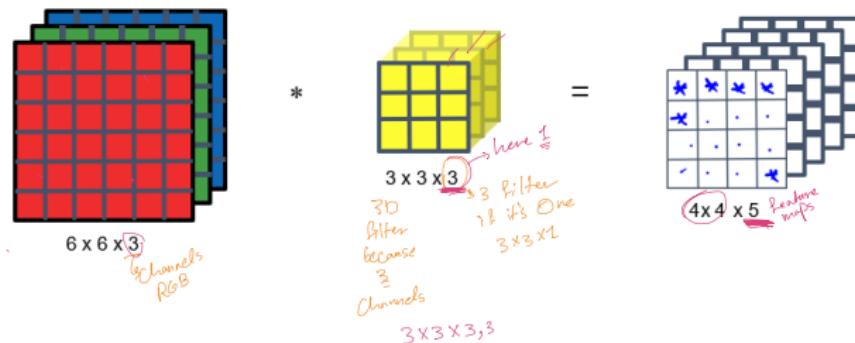
Computational Cost of Convolutions

► Computational cost of standard convolution:

$$\text{Cost} = \# \text{ filter params} \times \# \text{ filter positions} \times \# \text{ filters}$$
$$2160 = 3 \times 3 \times 3 \quad \times \quad \overset{\text{formula}}{4 \times 4} \quad \times \quad 5$$

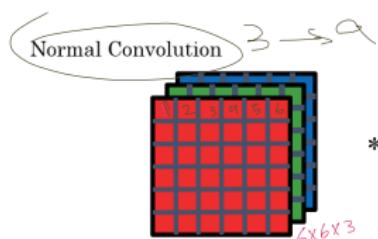
- Filters operate on all input channels, increasing computation significantly.

$$\frac{6 - 3}{1} + 1 = 4$$

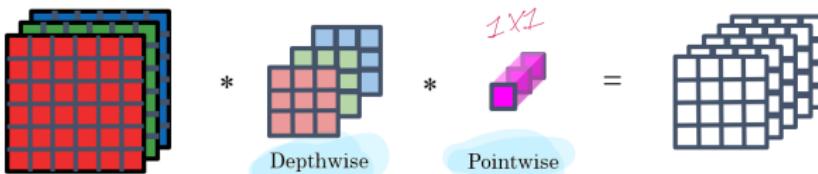


Depthwise-Separable Convolutions

- ▶ Split standard convolution into two steps:
 - **Depthwise Convolution:** Applies a single filter per input channel.
 - **Pointwise Convolution:** Combines outputs from depthwise convolution.
- ▶ **Key Benefit:** Reduces computational cost significantly compared to standard convolution.



Depthwise Separable Convolution



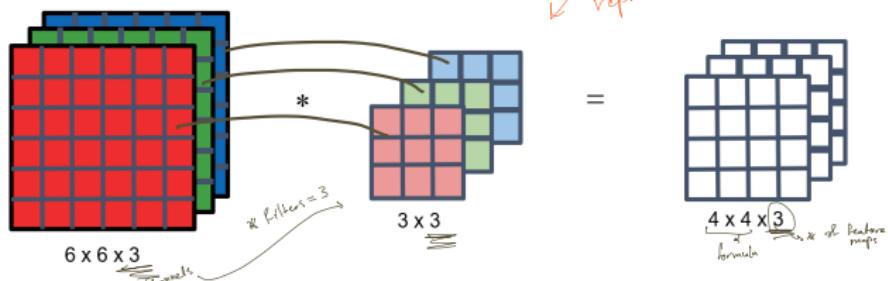
Depthwise Convolution

depthwise

$3 \rightarrow 3$

- Operates on each input channel separately.

depthwise
representation



$$\begin{array}{lcl} \text{Computational cost} & = & \# \text{filter params} \times \# \text{filter positions} \times \# \text{of filters} \\ 482 & = & 3 \times 3 \times 4 \times 4 \times 3 \end{array}$$

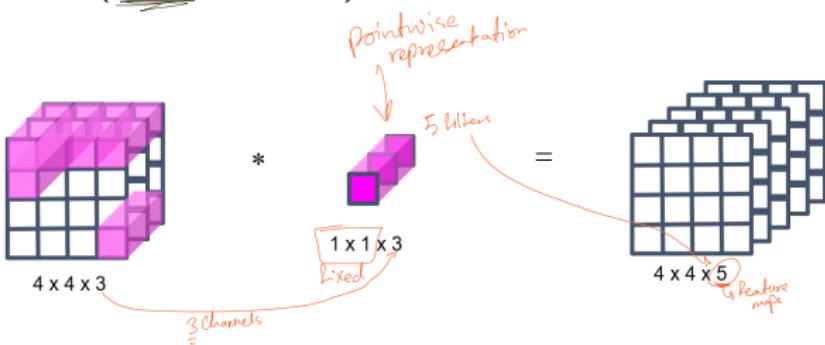
$$\# \text{of channels} = \# \text{of filters} = \# \text{of feature maps}$$

$$\frac{3}{3} = \frac{3}{3}$$

Pointwise Convolution

Pointwise

- Combines outputs from depthwise convolution using 1×1 convolutions (mixes channels).



$$\begin{aligned} \text{Computational cost} &= \# \text{filter params} \times \# \text{filter positions} \times \# \text{of filters} \\ &= 1 \times 1 \times 3 \times 4 \times 4 \times 5 \\ 240 &= 3 \times 16 \times 5 \end{aligned}$$

Cost Summary

Cost of normal convolution 2160

Cost of depthwise separable convolution

$$\begin{array}{l} \text{depthwise} + \text{pointwise} \\ 432 + 240 = 672 \end{array}$$

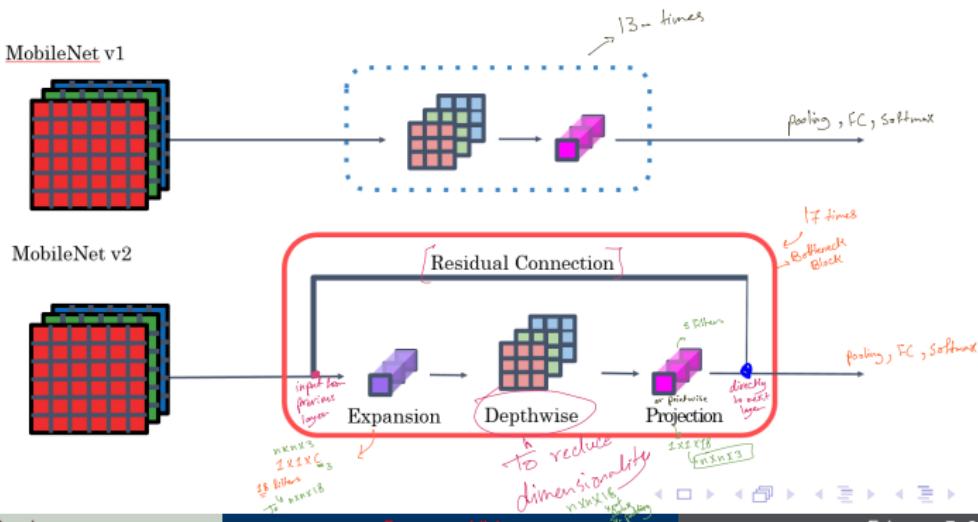
$$\frac{672}{2160} = 0.31 \text{ very efficient cost}$$

Howard et al. 2017 MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications

▶ MobileNet v2:

- Adds **residual connections**.
- Introduces:
 - ▶ **Expansion step**: Expands input dimensions before depthwise convolution.
 - ▶ **Projection step**: Reduces dimensions after processing.

These boxes reduce Computational Cost



These slides have been adapted from

- ▶ Fei-Fei Li, Yunzhu Li & Ruohan Gao, Stanford CS231n: Deep Learning for Computer Vision
- ▶ Assaf Shocher, Shai Bagon, Meirav Galun & Tali Dekel, WAIC DL4CV Deep Learning for Computer Vision: Fundamentals and Applications
- ▶ Justin Johnson, UMich EECS 498.008/598.008: Deep Learning for Computer Vision
- ▶ Sander Dieleman, Deepmind: Deep Learning Lecture Series 2020