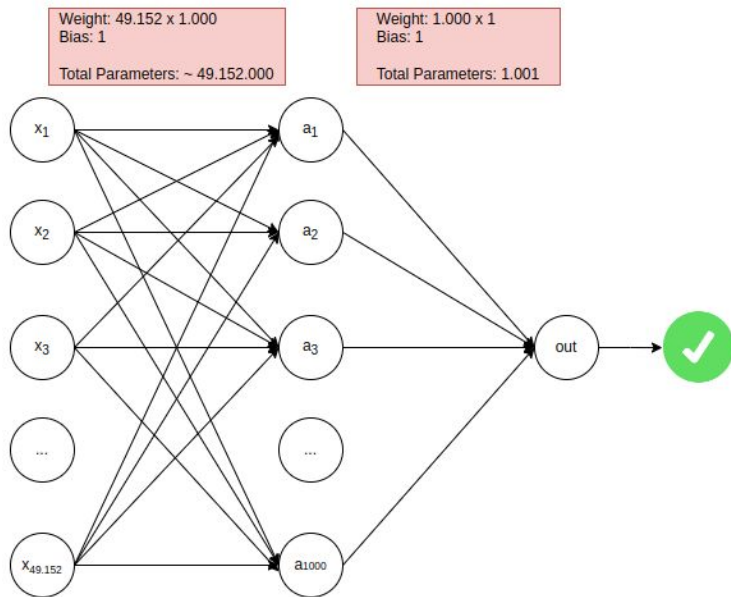

Day 2: Introduction to Convolutional Neural Networks (CNNs)

Lesson of Content

1. What is the problem with Neural Network ?
2. Image processing
3. Convolutional Operation
4. Some blocks
5. ImageNet challenge
6. Some Architecture
7. CNN Applications

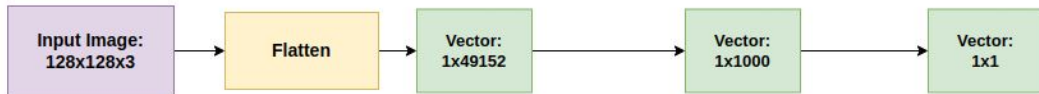
1. What problem with Neural Network ?



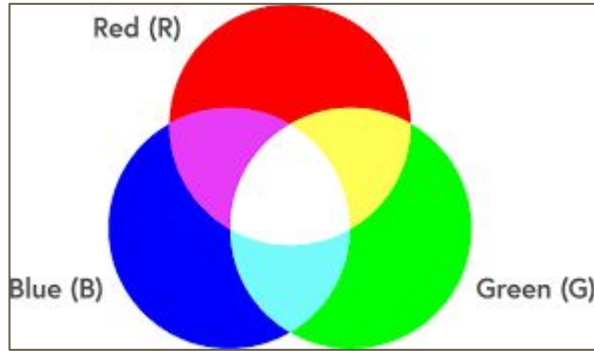
Với ảnh kích thước nhỏ:
 $128 \times 128 \times 3$

Số lớp Hidden: 1

=> Mạng tương đối đơn giản
mà số parameters đã đến gần
50 M. (Too big)



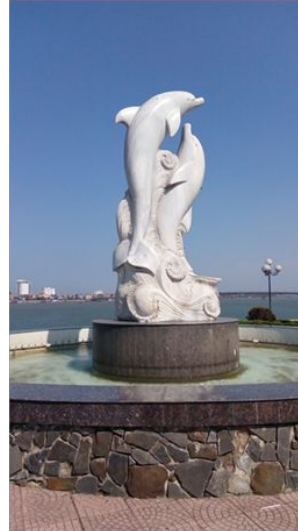
2. Image Processing



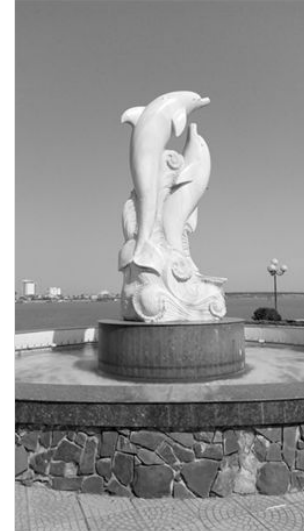
- **Color:** Được tạo từ ba kênh màu.
Shape: **32 x 80 x 3**
- **GrayScale:** Được tạo từ 1 kênh màu range từ 0 -> 255.
Shape: **32 x 80 x 1**
- **Black-and-white:** Được bởi hai giá trị 1 và 0.
Shape : **32 x 80 x 1**

80

32



Color

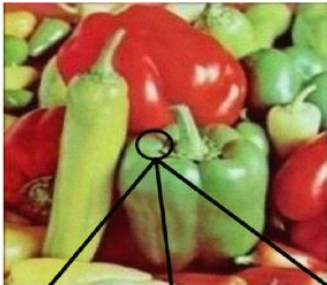


Grayscale



Black-and-white

2. Image Processing



240	241	241
240	237	238
239	240	240
238	237	240
240	240	239
239	240	240

Red

207	199	196
183	163	195
183	166	184
176	172	181
184	167	176
182	180	170

Green

234	231	225
223	213	225
219	211	195
176	205	189
168	141	117
160	142	117

Blue



0.1216	0.1255	0.1059
0.1176	0.1176	0.1137
0.1020	0.1020	0.1059
0.1490	0.0980	0.0902
0.5020	0.4196	0.2941
0.6392	0.6431	0.6510
0.7255	0.6667	0.6353
0.6824	0.7137	0.6863
0.6784	0.7373	0.7373
0.6980	0.7176	0.7176
0.7255	0.7216	0.7098

Element-wise multiplication matrix

Đây là một phép toán được sử dụng nhiều nhất trong CNNs.

<table border="1"><tr><td>a_1</td><td>a_2</td><td>a_3</td></tr><tr><td>a_4</td><td>a_5</td><td>a_6</td></tr><tr><td>a_7</td><td>a_8</td><td>a_9</td></tr></table>	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9	$*$	<table border="1"><tr><td>b_1</td><td>b_2</td><td>b_3</td></tr><tr><td>b_4</td><td>b_5</td><td>b_6</td></tr><tr><td>b_7</td><td>b_8</td><td>b_9</td></tr></table>	b_1	b_2	b_3	b_4	b_5	b_6	b_7	b_8	b_9	$=$	<table border="1"><tr><td>a_1b_1</td><td>a_2b_2</td><td>a_3b_3</td></tr><tr><td>a_4b_4</td><td>a_5b_5</td><td>a_6b_6</td></tr><tr><td>a_7b_7</td><td>a_8b_8</td><td>a_9b_9</td></tr></table>	a_1b_1	a_2b_2	a_3b_3	a_4b_4	a_5b_5	a_6b_6	a_7b_7	a_8b_8	a_9b_9
a_1	a_2	a_3																													
a_4	a_5	a_6																													
a_7	a_8	a_9																													
b_1	b_2	b_3																													
b_4	b_5	b_6																													
b_7	b_8	b_9																													
a_1b_1	a_2b_2	a_3b_3																													
a_4b_4	a_5b_5	a_6b_6																													
a_7b_7	a_8b_8	a_9b_9																													

3. Convolutional Operation

Input

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Kernel

0	-1	0
-1	5	-1
0	-1	0

★

Bias

+ 5

Output

11	12
15	16

=

★ Cross-Correlation

3. Convolutional Operation

Input: $256 \times 256 \times 3$

Kernel: $4 \times 4 \times 3$, có 5 cái kernel,
Padding = 0, Stride = 2

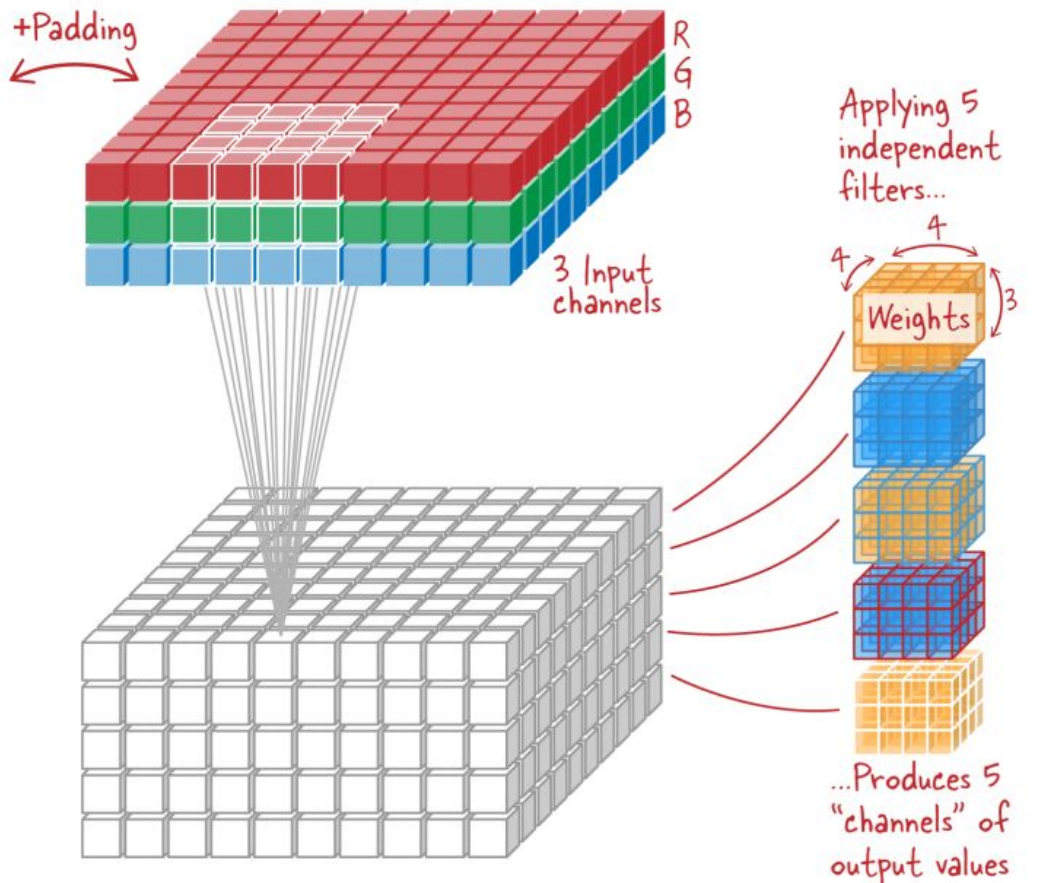
Output: $127 \times 127 \times 5$

- **Trích xuất đặc trưng:** Giúp học và phát hiện được các đặc trưng của dữ liệu như (Màu, góc, cạnh, ...).

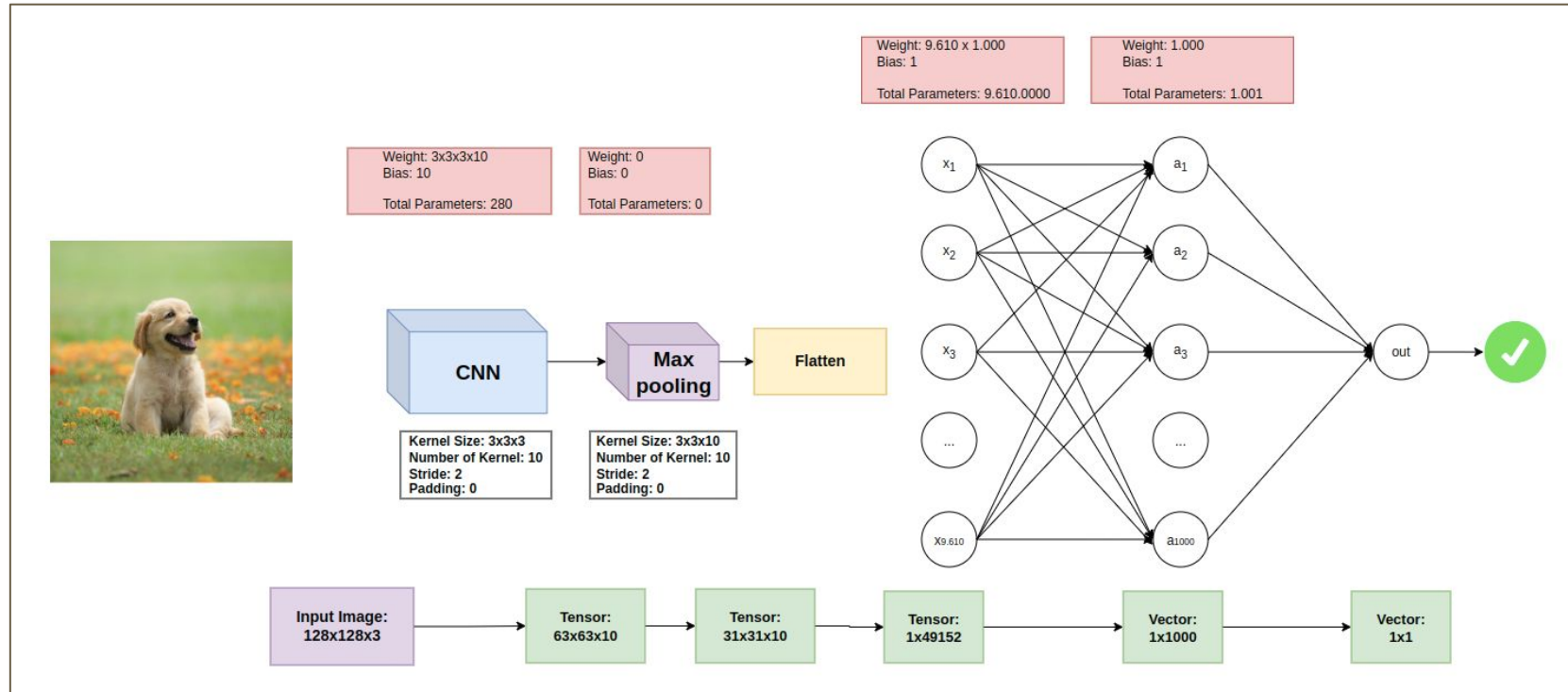
- **Hiểu được cấu trúc không gian:** của dữ liệu.

- **Tăng tính phi tuyến:** CNN có khả năng phát hiện các đặc trưng phức tạp -> nó có thể tạo thành các hàm phức tạp.

- **Giảm số lượng tham số:** Giúp inference nhanh và training nhanh (Do CNN không phụ thuộc vào kích thước ảnh đầu vào mà phụ thuộc vào kernel size)



3. Convolutional Operation



Total Parameter:

- **Neural Network: 50M**
- **CNNs: 10M**

3. Convolutional Operation

Input

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Kernel

0	-1	0
-1	5	-1
0	-1	0

Bias

+ 5

Output

11	12
15	16

★ Cross-Correlation

$$1 \times 0 + 5 \times -1 + 9 \times 0 + 2 \times -1 + 6 \times 5 + 10 \times -1 + 3 \times 0 + 7 \times -1 + 11 \times 0 + 5 = 11$$

3. Convolutional Operation

Input

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Kernel

0	-1	0
-1	5	-1
0	-1	0

Bias

+ 5

Output

11	12
15	16

★ Cross-Correlation

$$2 \times 0 + 6 \times -1 + 10 \times 0 + 3 \times -1 + 7 \times 5 + 11 \times -1 + 4 \times 0 + 8 \times -1 + 12 \times 0 + 5 = 12$$

3. Convolutional Operation

Input

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Kernel

0	-1	0
-1	5	-1
0	-1	0

Bias

+ 5

Output

11	12
15	16

★ Cross-Correlation

$$5 \times 0 + 9 \times -1 + 13 \times 0 + 6 \times -1 + 10 \times 5 + 14 \times -1 + 7 \times 0 + 11 \times -1 + 15 \times 0 + 5 = 15$$

3. Convolutional Operation

Input

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Kernel

0	-1	0
-1	5	-1
0	-1	0

★

Bias

+ 5

Output

11	12
15	16

★ Cross-Correlation

$$6 \times 0 + 10 \times -1 + 14 \times 0 + 7 \times -1 + 11 \times 5 + 15 \times -1 + 8 \times 0 + 12 \times -1 + 16 \times 0 + 5 = 16$$

3. Convolutional Operation

- Input: $(N, C_{in}, H_{in}, W_{in})$ or (C_{in}, H_{in}, W_{in})
- Output: $(N, C_{out}, H_{out}, W_{out})$ or $(C_{out}, H_{out}, W_{out})$, where

$$H_{out} = \left\lfloor \frac{H_{in} + 2 \times \text{padding}[0] - \text{dilation}[0] \times (\text{kernel_size}[0] - 1) - 1}{\text{stride}[0]} + 1 \right\rfloor$$

$$W_{out} = \left\lfloor \frac{W_{in} + 2 \times \text{padding}[1] - \text{dilation}[1] \times (\text{kernel_size}[1] - 1) - 1}{\text{stride}[1]} + 1 \right\rfloor$$

source: <https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html>

N: Số lượng sample

C_{in/out}: Số lượng Channel ngõ vào/ra

H, W: Lần lượt là Width và Height

Padding: Kích thước phần được thêm vào (Theo chiều Width và chiều Height)

Stride: Kích thước trượt của Kernel (Theo chiều Width và chiều Height)

Kernel_size: Kích thước của Kernel

Pooling

8	7	5	3
12	9	5	7
13	2	10	3
9	4	5	14

2x2 pooling,
stride 2

Max pooling

12	7
13	14



Average pooling

9	5
7	8

4. Some blocks

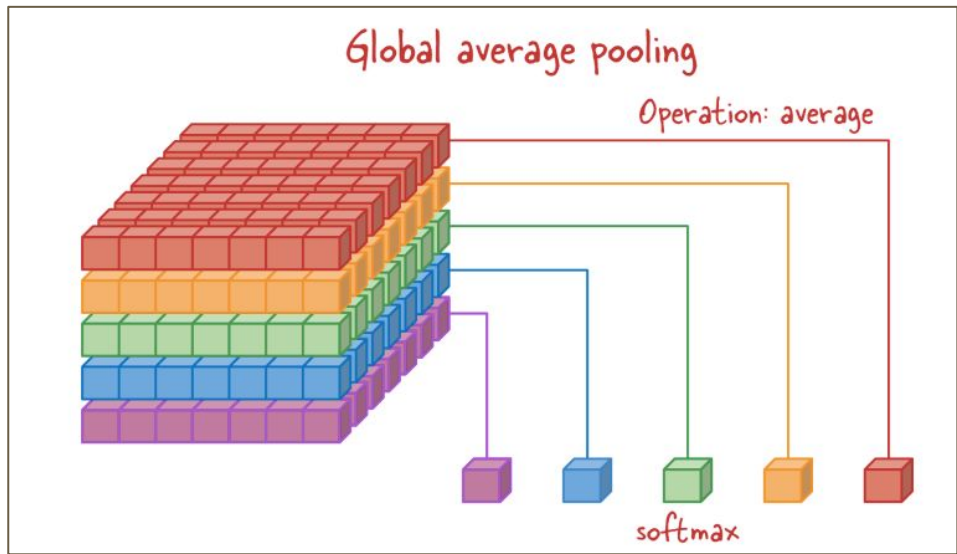
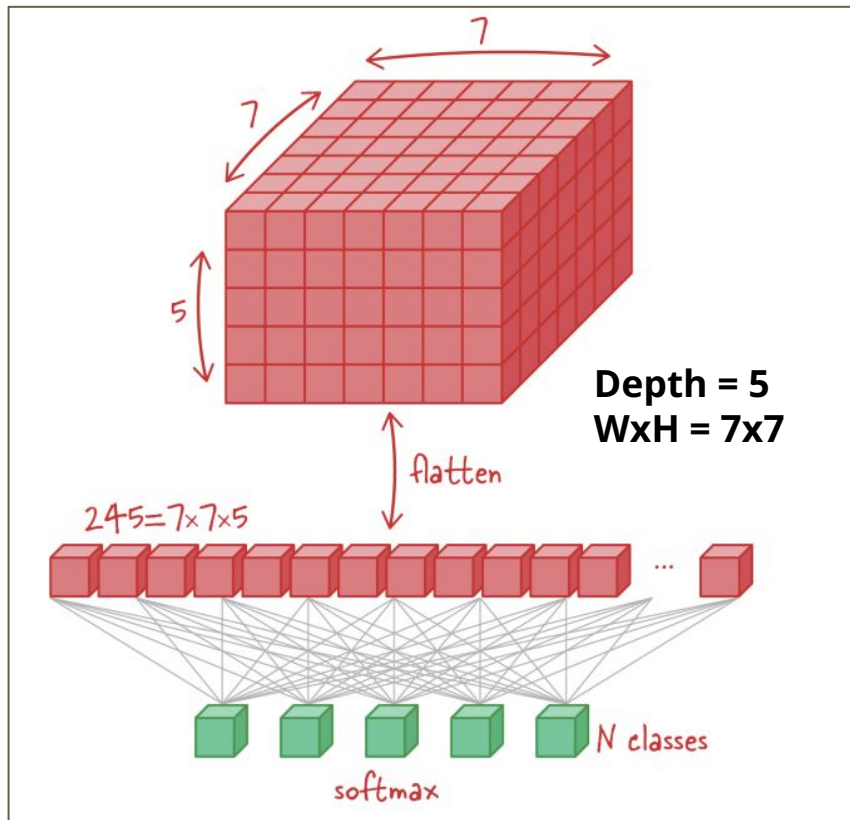
Giảm kích thước không giảm : Giảm số lượng tham số và phép tính, giúp mô hình nhanh và ít bộ nhớ hơn.

Trích xuất đặc trưng cục bộ

Tăng tính phi tuyến: Giúp thêm một chút tính phi tuyến vào mô hình.

Mất thông tin: Mỗi lần pooling mất một số thông tin. Nếu dùng MaxPooling thì chỉ giữ những pixel có giá trị max, còn AveragePooling thì giữ lại giá trị trung bình.

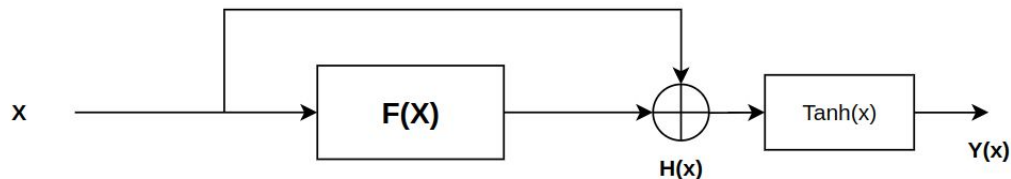
Flatten and Global Average Pooling (GAP)



Flatten thường dùng cho mạng CNN ít lớp (LeNet, AlexNet)

GAP thường dùng cho mạng CNN hiện đại như ResNet

Skip Connection (Residual Connection)



Giảm Vanishing gradient

Cho phép mạng sâu hơn: Minh chứng cho thấy mạng ResNet thường có số lớp sâu nhưng vẫn hiệu quả.

Khả năng hội tụ nhanh hơn: Trên thực nghiệm cho thấy mô hình có skip connection sẽ hội tụ nhanh hơn.

$$H(x) = F(x) + x$$

$$Y(x) = \tanh(H(x))$$

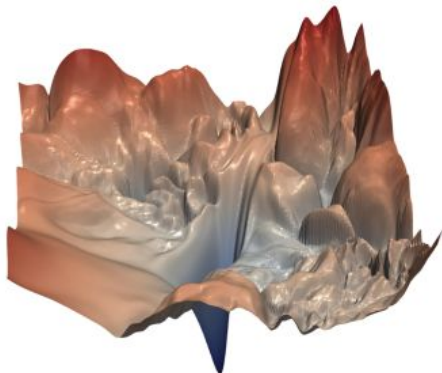
$$\frac{dY}{dx} = \frac{dY}{dH} \cdot \frac{dH}{dx}$$

$$\frac{d \tanh(z)}{dz} = 1 - \tanh^2(z)$$

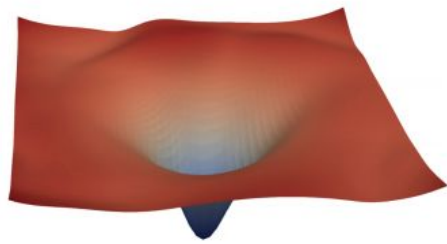
$$\frac{dY}{dH} = 1 - \tanh^2(H(x))$$

$$\frac{dH}{dx} = \frac{dF}{dx} + 1$$

$$\frac{dY}{dx} = (1 - \tanh^2(H(x))) \cdot (\frac{dF}{dx} + 1)$$



(a) without skip connections

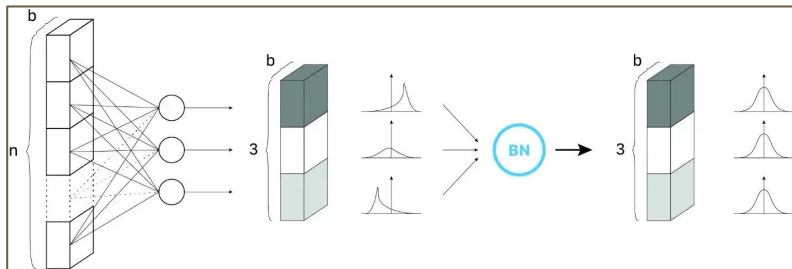


(b) with skip connections

Batch Normalization

- Thường đặt ở sau hàm kích hoạt để chuẩn hóa đưa phổ dữ liệu về vùng bão hòa. Giúp **giảm Vanishing/ Exploding Gradient**.

- **Đồng nhất phân phối**: Giúp việc huấn luyện nhanh và ổn định hơn.



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

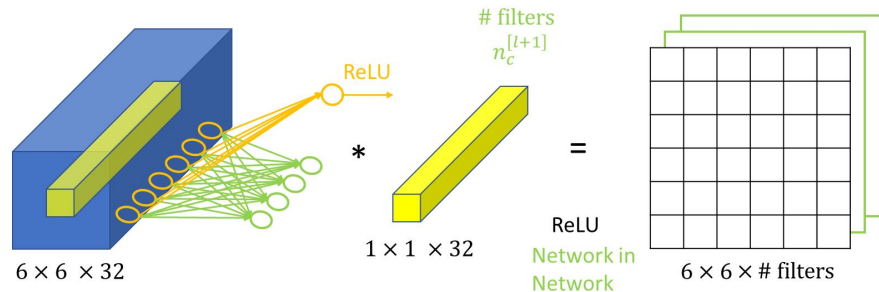
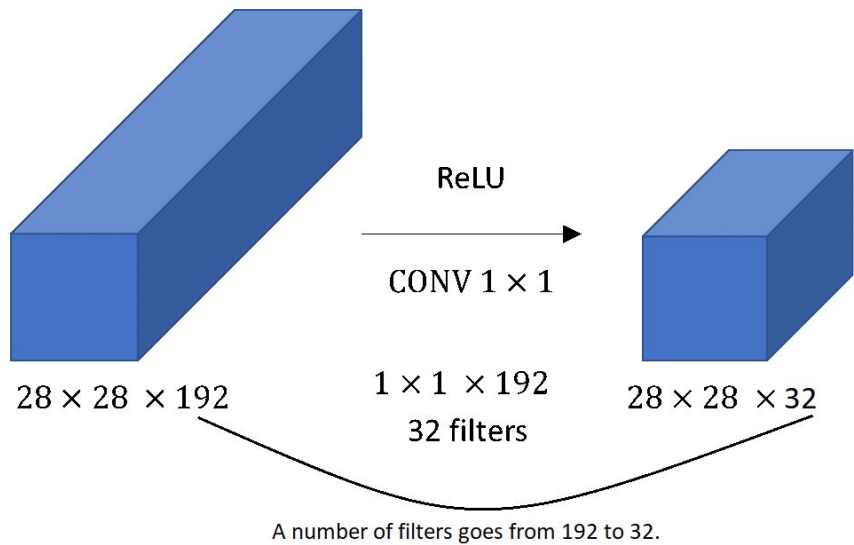
$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

1x1 Convolutions (Pointwise Convolutions)



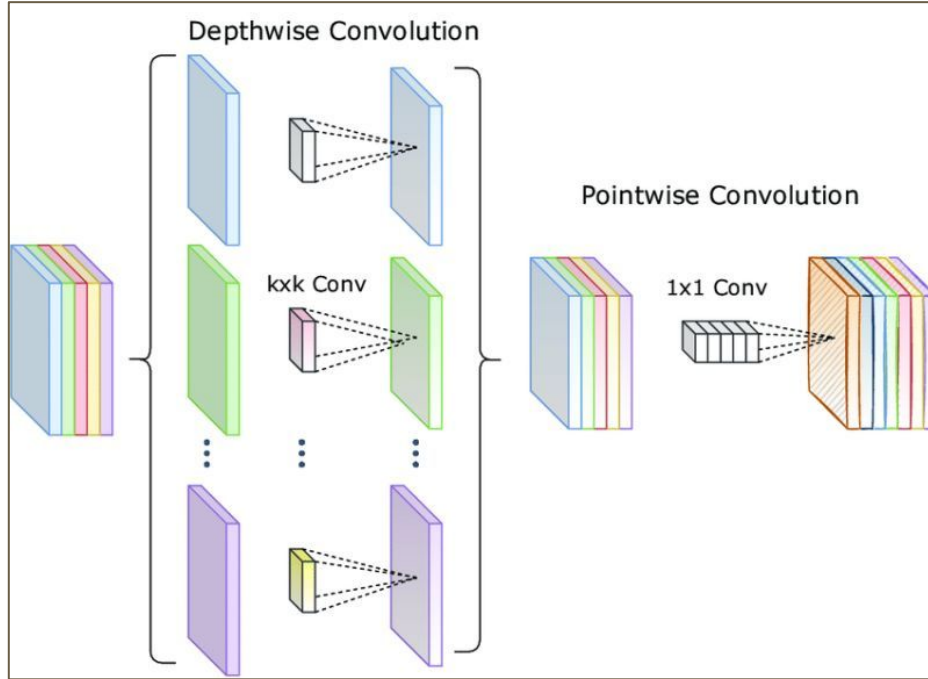
Thường dùng để **Giảm/Tăng chiều** của một feature map mà không thay đổi chiều cao, chiều rộng.

- **Tính toán phi không gian:** Giúp mô hình có sự học hiểu giữa các kênh của tích chập.

- **Giảm chi phí tính toán:** Thường đặt 1x1 Conv trước khi vào 3x3 hoặc 5x5 Conv , giúp giảm số lượng tham số tính toán.

- **Tăng cường tính phi tuyến:** Giúp mô hình có những biểu diễn phức tạp hơn.

Depthwise Convolutions



Depth-Separable Convolutions

Giảm số lượng tham số: Mỗi kênh đầu vào có một kernel riêng, nên số lượng tham số giảm đáng kể.

Giảm tính toán: Số lượng tính toán chỉ phụ thuộc vào số lượng channel.

Tích hợp thông tin đặc trưng và không gian (
Nếu dùng theo sau là Pointwise Convolutions)

Không thể học mối quan hệ phức tạp như CNN truyền thống.

Các architecture chứa loại này: **Xception, MobileNet, EfficientNet, SqueezeNet**

4. Some blocks

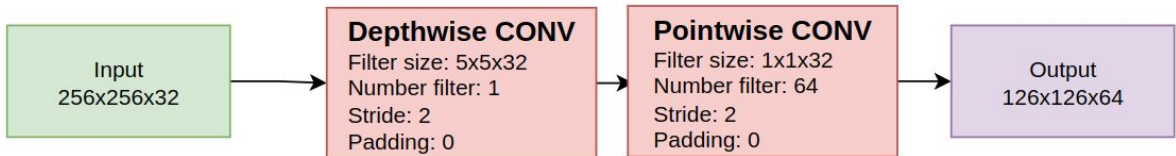
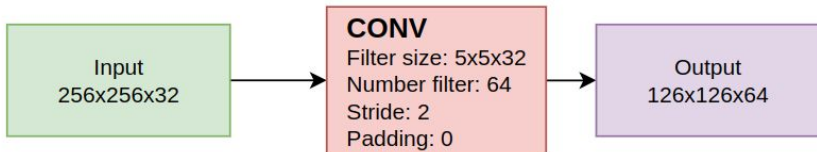
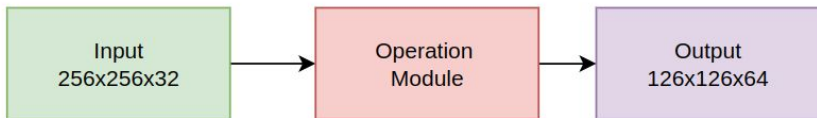
CONV:

Parameters:

$$(5 \times 5 \times 32 + 1) \times 64 = 51.264$$

Computational Complexity:

$$(126 \times 126) \times (5 \times 5 \times 32) \times 64 = 812.851.200 \text{ (Phép tính)}$$



Depth-Separable Convolutions:

Parameters:

$$w1 = 5 \times 5 \times 32 = 800$$

$$w2 = (1 \times 1 \times 32 + 1) \times 64 = 2.112$$

$$w1 + w2 = 2.912$$

Computational Complexity:

$$t1 = (126 \times 126) \times (5 \times 5 \times 32) \times 1 = 12.700.800$$

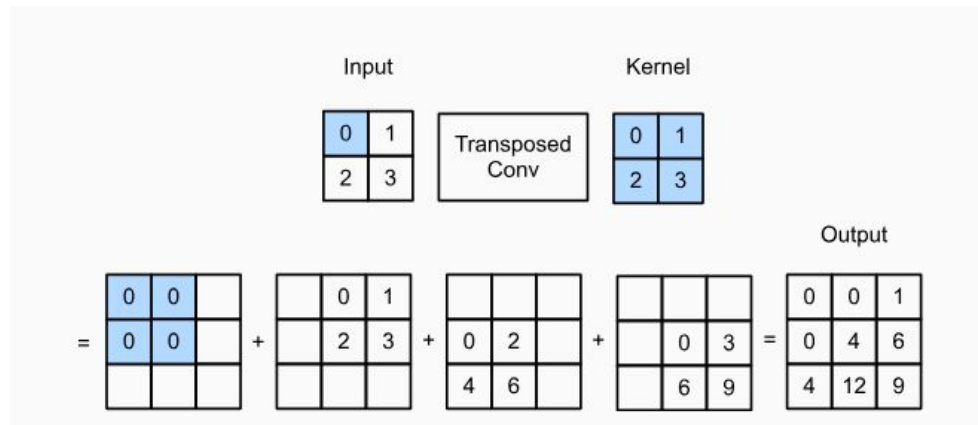
$$t2 = (126 \times 126) \times (1 \times 1 \times 32) \times 64 = 32.514.048$$

$$t1 + t2 = 45.214.848 \text{ (Phép tính)}$$

Giảm số lượng tính toán đi: 17.9 lần

Giảm số lượng tham số đi: 17.6 lần

Transposed Convolution



Thường dùng trong Image Segmentation , Generative Adversarial Networks (GANs) và Super Resolution.

- **Khôi phục thông tin:** Do các ConV thường làm mất đi thông tin, thì TConv sẽ là khôi phục thông tin của feature đã mất.

- **Phân đoạn hình ảnh:** Nó là bước trung gian từ 1 pixel tăng lên rồi giảm về lại kích thước ban đầu, từ đó đưa ra dự đoán mức độ pixel

- **Tăng kích thước feature map đầu ra.**

Vấn đề **checkerboard pattern** cần lưu ý.

```

1 import torch
2 from torch import nn
3 X = torch.tensor([[[0.0, 1.0], [2.0, 3.0]]])
4 K = torch.tensor([[[0.0, 1.0], [2.0, 3.0]]])
5
6 X, K = X.reshape(1, 1, 2, 2), K.reshape(1, 1, 2, 2)
7 tconv = nn.ConvTranspose2d(1, 1, kernel_size=2, bias=False)
8 tconv.weight.data = K
9 tconv(X)

```

[8] ✓ 0.0s

... tensor([[[[0., 0., 1.],
[0., 4., 6.],
[4., 12., 9.]]]], grad_fn=<ConvolutionBackward0>)

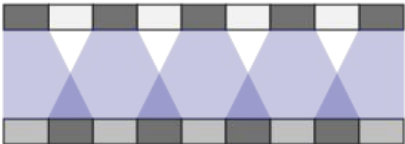
Checkerboard pattern



Xảy ra do việc sử dụng stride lớn trong TConv. Khi tính toán output nhiều vùng được cộng chồng lên nhau.

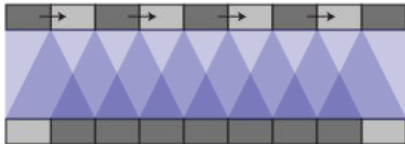
Solution:

- Thường dùng Kernel là số chẵn.
- Kết hợp thêm một lớp ConV ngay sau khi TConv giúp làm "mịn" hình ảnh và giảm thiểu sự không đồng đều.



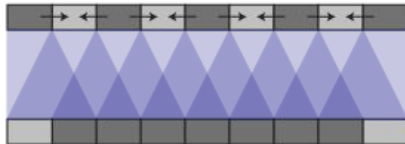
$$\begin{bmatrix} a & c \\ & b \\ & a & c \\ & & b \end{bmatrix}$$

Deconvolution



$$\begin{bmatrix} a+b & c \\ & a & b+c \\ & & a+b & c \\ & & & a & b+c \end{bmatrix}$$

NN-Resize Convolution



$$\begin{bmatrix} a+\frac{1}{2}b & \frac{1}{2}b+c \\ \frac{1}{2}a & \frac{1}{2}a+b+\frac{1}{2}c & \frac{1}{2}c \\ & a+\frac{1}{2}b & \frac{1}{2}b+c \\ & \frac{1}{2}a & \frac{1}{2}a+b+\frac{1}{2}c & \frac{1}{2}c \end{bmatrix}$$

Bilinear-Resize Convolution

5. ImageNet Challenge

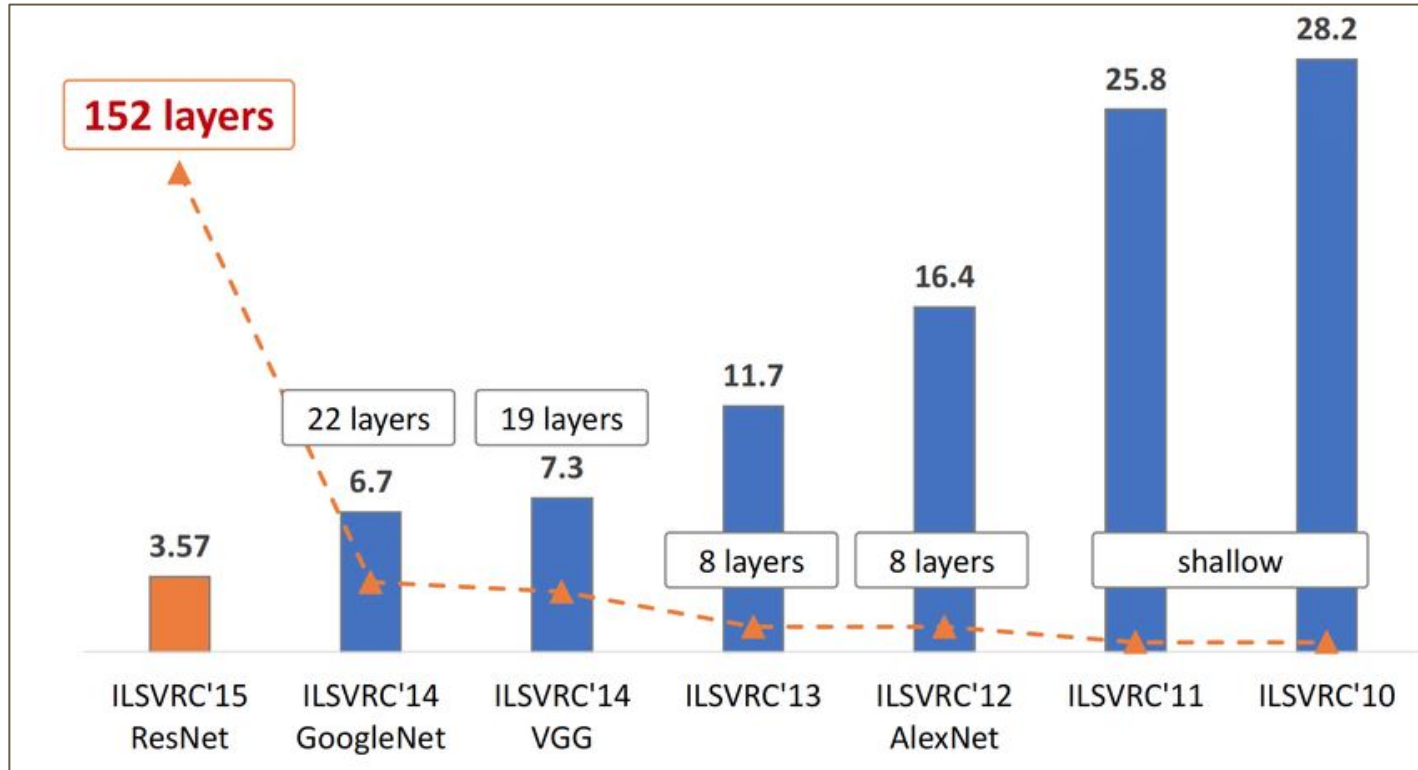


~14 million labeled images, 20k classes
Images gathered from Internet
Human labels via Amazon MTurk

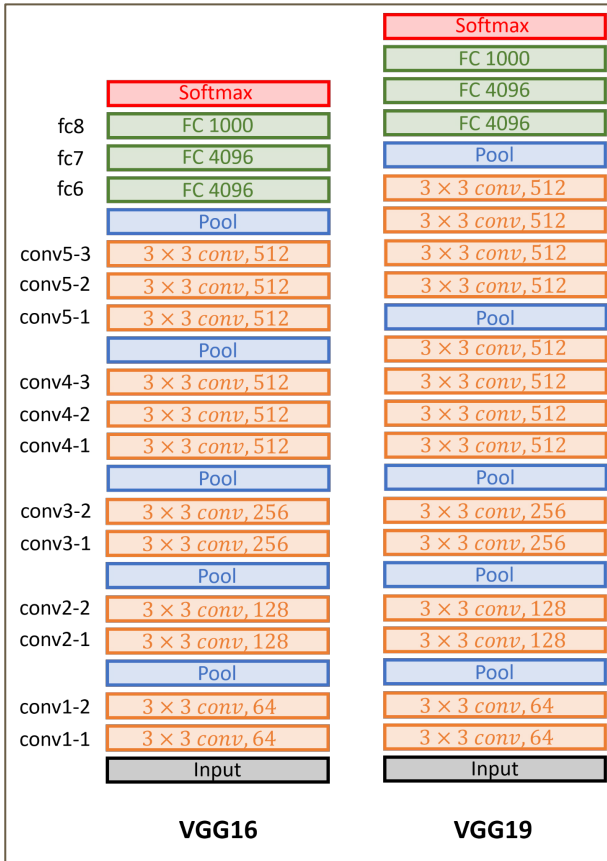
**ImageNet Large-Scale Visual Recognition Challenge
(ILSVRC): 1.2 million training images, 1000 classes**

Link: www.image-net.org/challenges/LSVRC/

Winner in year 2015



VGG-19



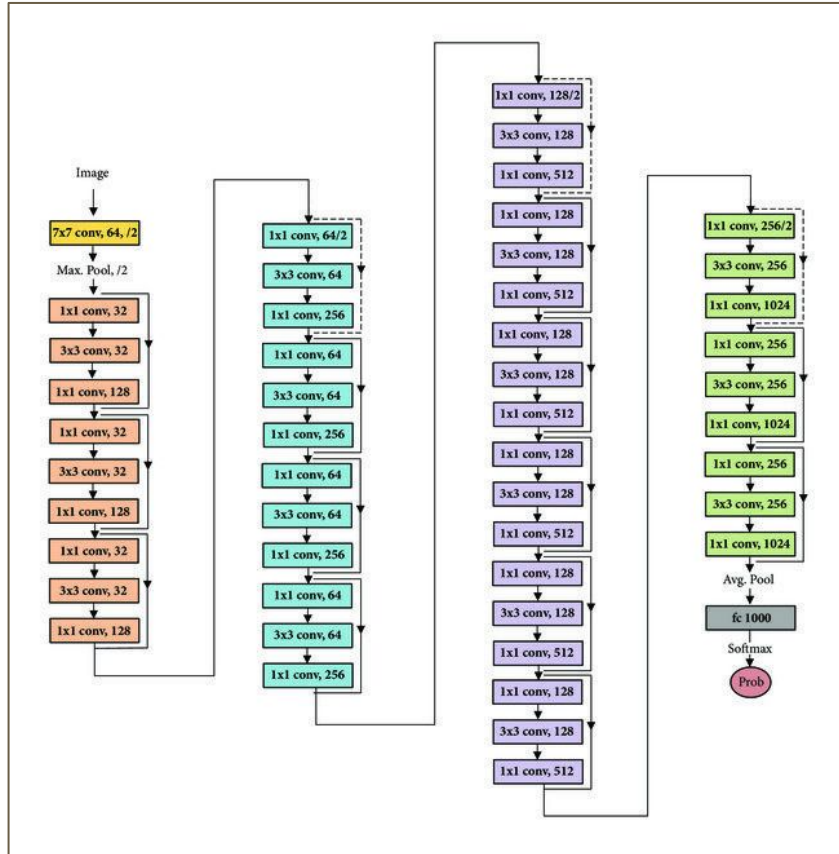
Model	Parameters (excl. classif. head ^a)	ImageNet accuracy	104 flowers F1 score ^b (fine-tuning)	104 flowers F1 score (trained from scratch)
VGG19	20M	71%	88% precision: 89%, recall: 88%	N/A ^c
Previous best for comparison:				
AlexNet	3.7M	60%		39% precision: 44%, recall: 38%

- Dễ hiểu và triển khai

- **Số lượng tham số lớn** (Do lạm dụng nhiều Convolutional và Pooling) dẫn đến train và dự đoán lâu
- **Vẫn có thể bị vanishing gradient**
- **Dễ bị overfit với tập huấn luyện nhỏ** do số lượng tham số lớn.

6. Some Architecture

ResNet-50



Model	Parameters (excl. classif. head ^a)	ImageNet accuracy	104 flowers F1 score ^b (fine-tuning)	104 flowers F1 score (trained from scratch)
ResNet50	23M	75%	94% prec.: 95%, recall: 94%	73% prec.: 76%, recall: 72%
Previous best for comparison:				
InceptionV3	22M	78%	95% prec.: 95%, recall: 94%	
SqueezeNet, 24 layers	2.7M			76% prec.: 77%, recall: 75%

- **Mạng sâu và hiệu quả:** Nhờ sử dụng **Residual blocks** giúp không bị vanishing gradient.
- **Mạng sâu hơn ít tham số** (so với VGG): Nhờ việc sử dụng nhiều Convolutional 1x1.
- **Số lượng tham số lớn** (Do lạm dụng nhiều CNN và Pooling) dẫn đến train và dự đoán lâu
- **Đễ bị overfit với tập huấn luyện nhỏ** do số lượng tham số lớn.

Applications of CNN



Computer Vision and related application



Natural Language Processing



Object Detection and Segmentation



Image Classification



Speech Recognition



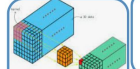
Video Processing



Low-resolution images



Limited Resource system



CNN for various dimensional data



Object Counting

7. CNN Applications

