AI Training Sessions

# Introduction to Convolutional Neural Networks

by Bartlomiej Borzyszkowski

# AGENDA

➤ Computer Vision - problems

➤ Architecture of CNN

➤ What happens inside?

➤ Popular architectures

➤ Popular datasets

➤ What's next?

_____

➤ Hands-on

# COMPUTER VISION

**Computer vision** is an interdisciplinary field that deals with how computers can be made for gaining high-level understanding from digital images or videos. From the perspective of engineering, it seeks to **automate tasks** that the **human visual system** can do.

*~Wikipedia*

# APPLICATIONS

- agriculture
- **augmented reality**
- **autonomous vehicles**
- **biometrics**
- character recognition
- forensics
- quality inspection
- face recognition
- gesture analysis
- geoscience
- image restoration

- **medical image analysis**
- pollution monitoring
- process control
- remote sensing
- robotics
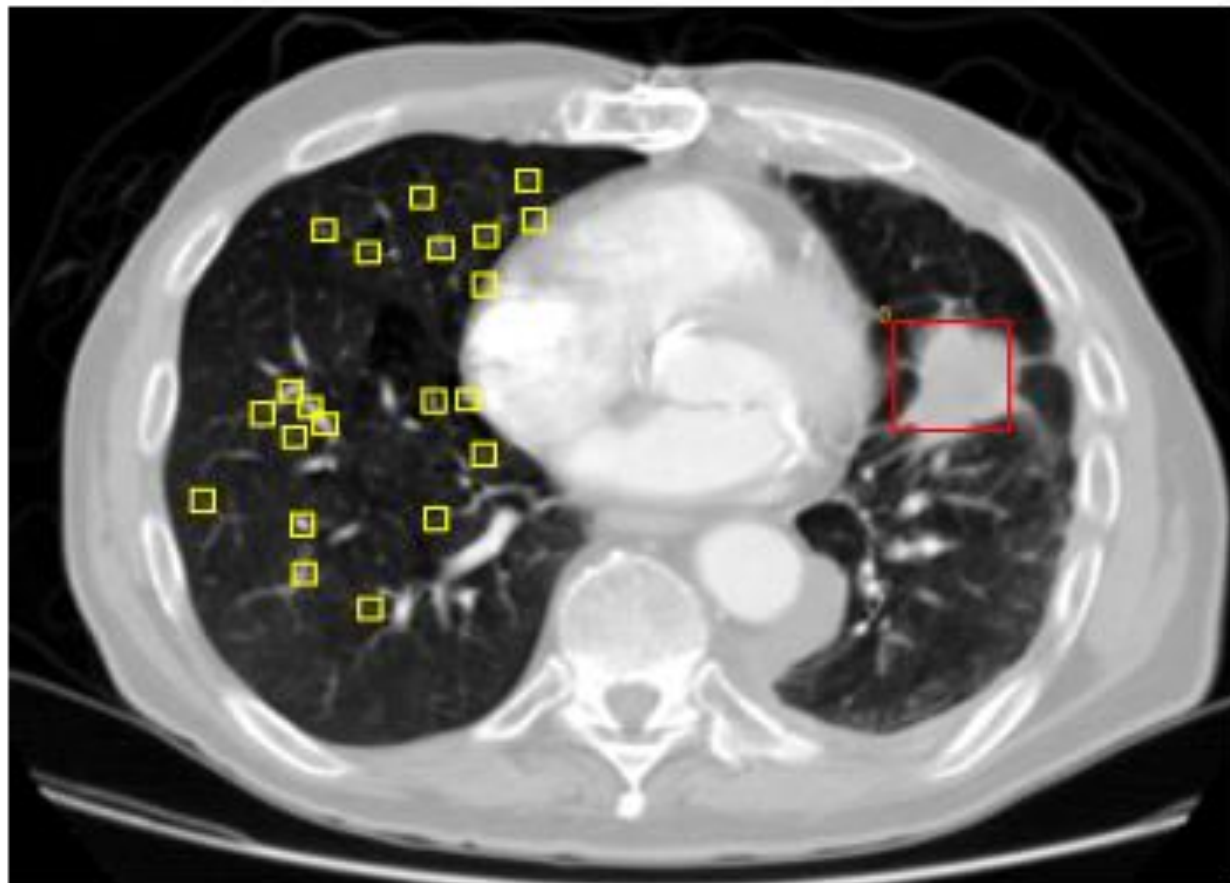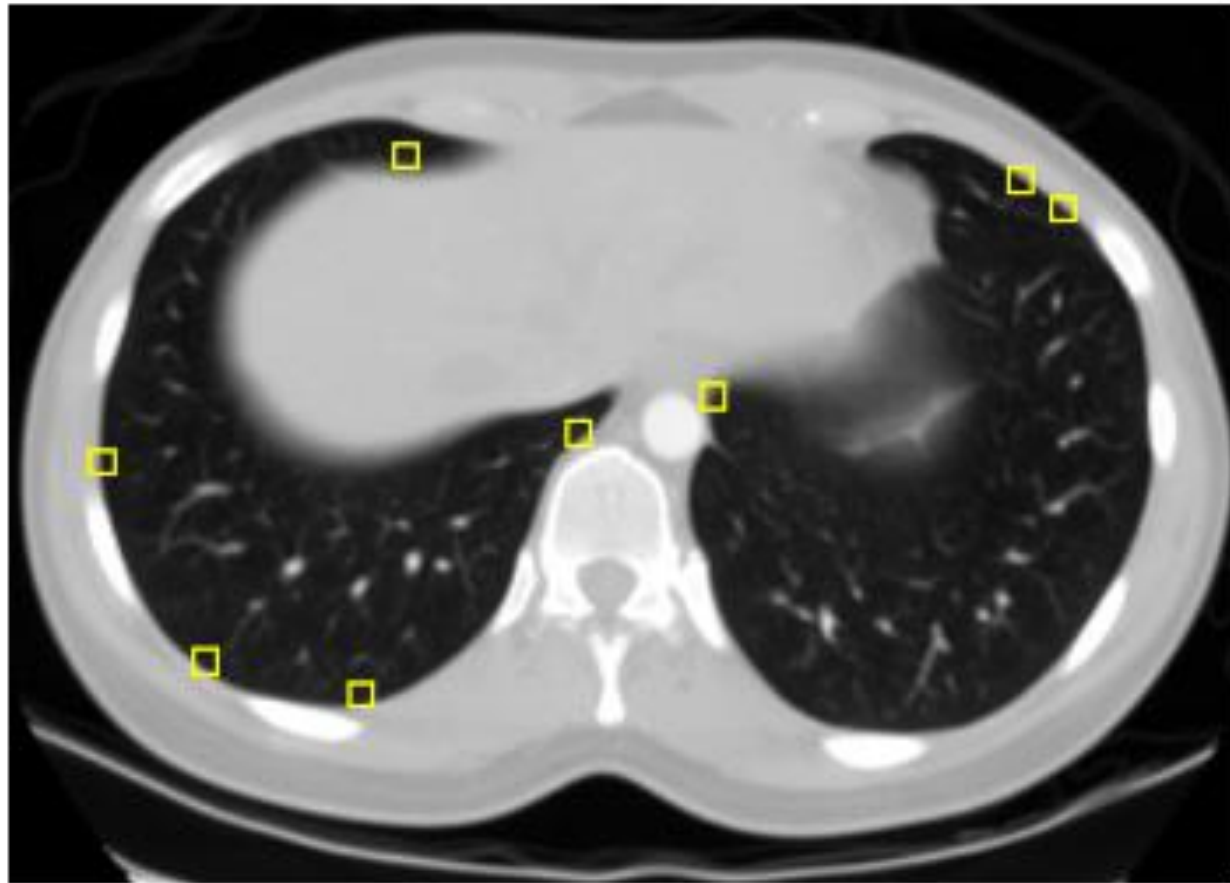- security and surveillance
- transport
- ...

# AUTONOMOUS VEHICLES

➤ Multiple sensors like:

    ➤ Cameras

    ➤ Lidars

    ➤ Ultrasonic sensors

➤ Big players like:

    ➤ Intel & BMW

    ➤ Nvidia & Audi
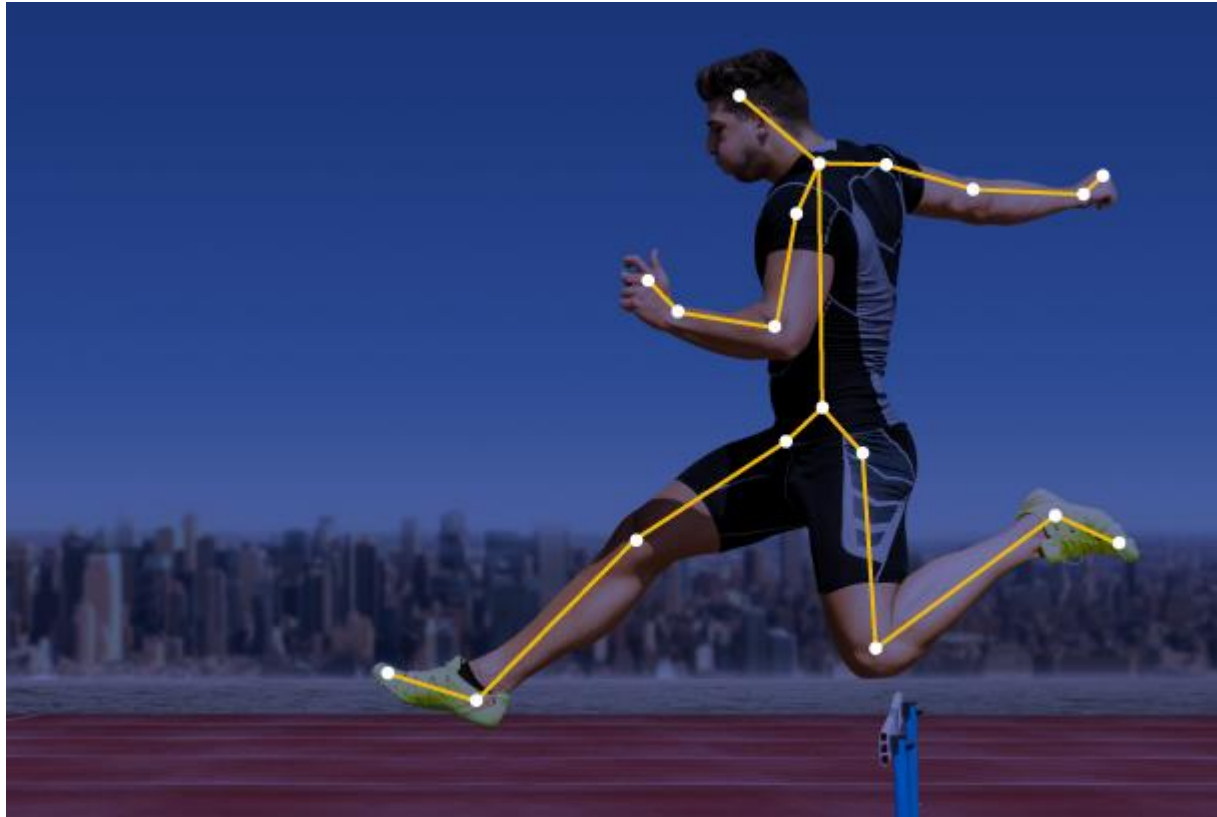
    ➤ Google

    ➤ Tesla

LEFT REARWARD VEHICLE CAMERA

MEDIUM RANGE VEHICLE CAMERA

RIGHT REARWARD VEHICLE CAMERA

MOTION FLOW    LANE LINES    LANE LINES    ROAD FLOW    IN-PATH OBJECTS    ROAD LIGHTS    OBJECTS    ROAD SIGNS

# MEDICINE

- ➤ Assistance for doctors,

- ➤ Cancer detection (Kaggle's Data Science Bowl 2017),

- ➤ Healthcare automation (future?).

# SPORT

➤ Intel Olympic Games!

➤ Statistics for audience

➤ Support for athletes

# ROBOTICS

- ➤ CNNs in manufactures
- ➤ Robots with computer vision
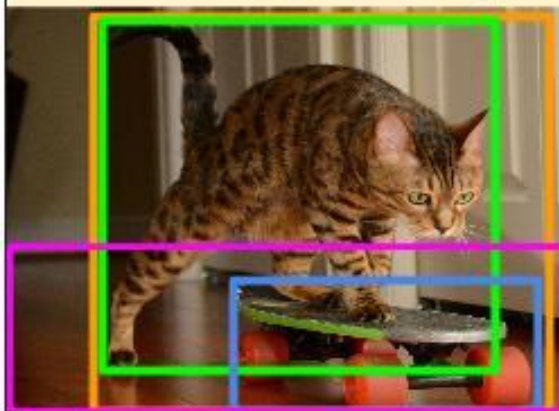- ➤ Reinforcement learning

# WHAT WE CAN DO?

**Classification**

**Classification + Localization**

**Object Detection**

**Instance Segmentation**



CAT

CAT

CAT, DOG, DUCK

CAT, DOG, DUCK
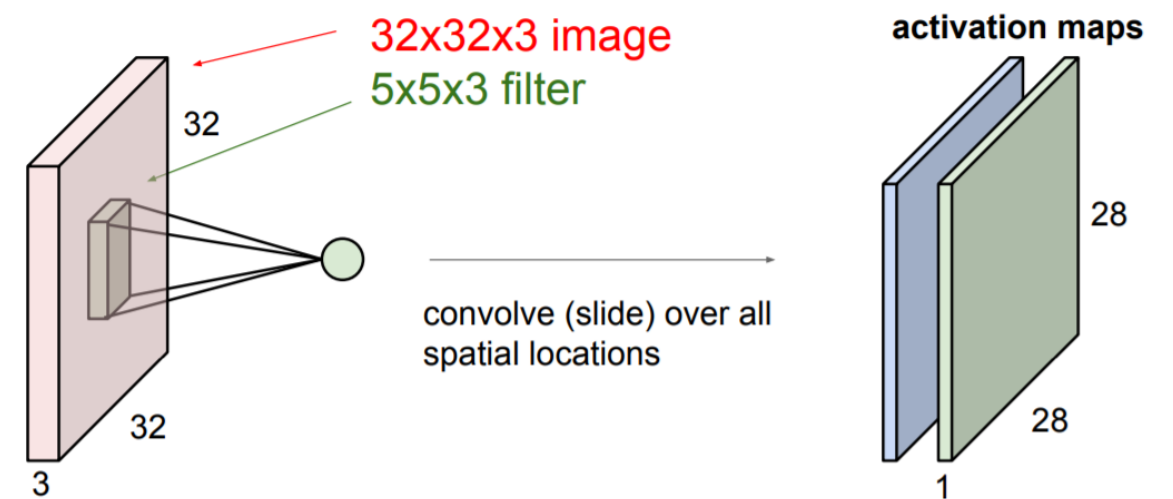
Single object

Multiple objects

# CNN ARCHITECTURE

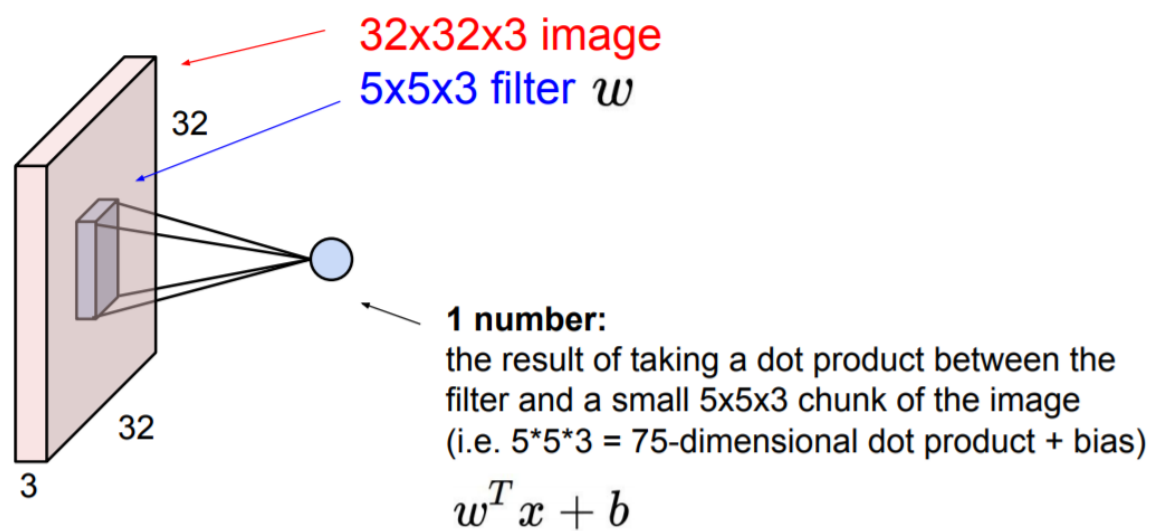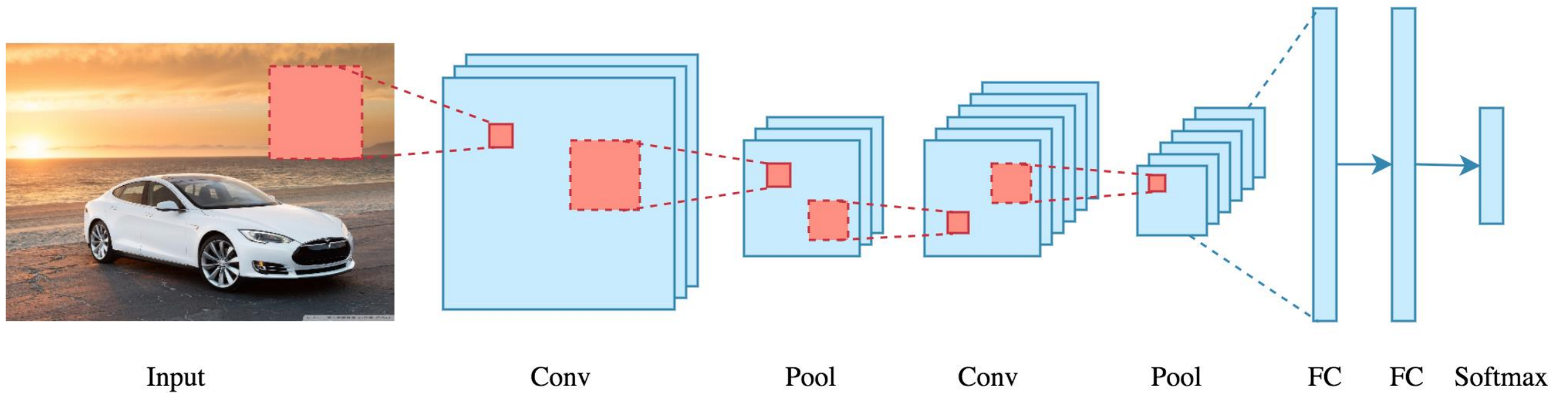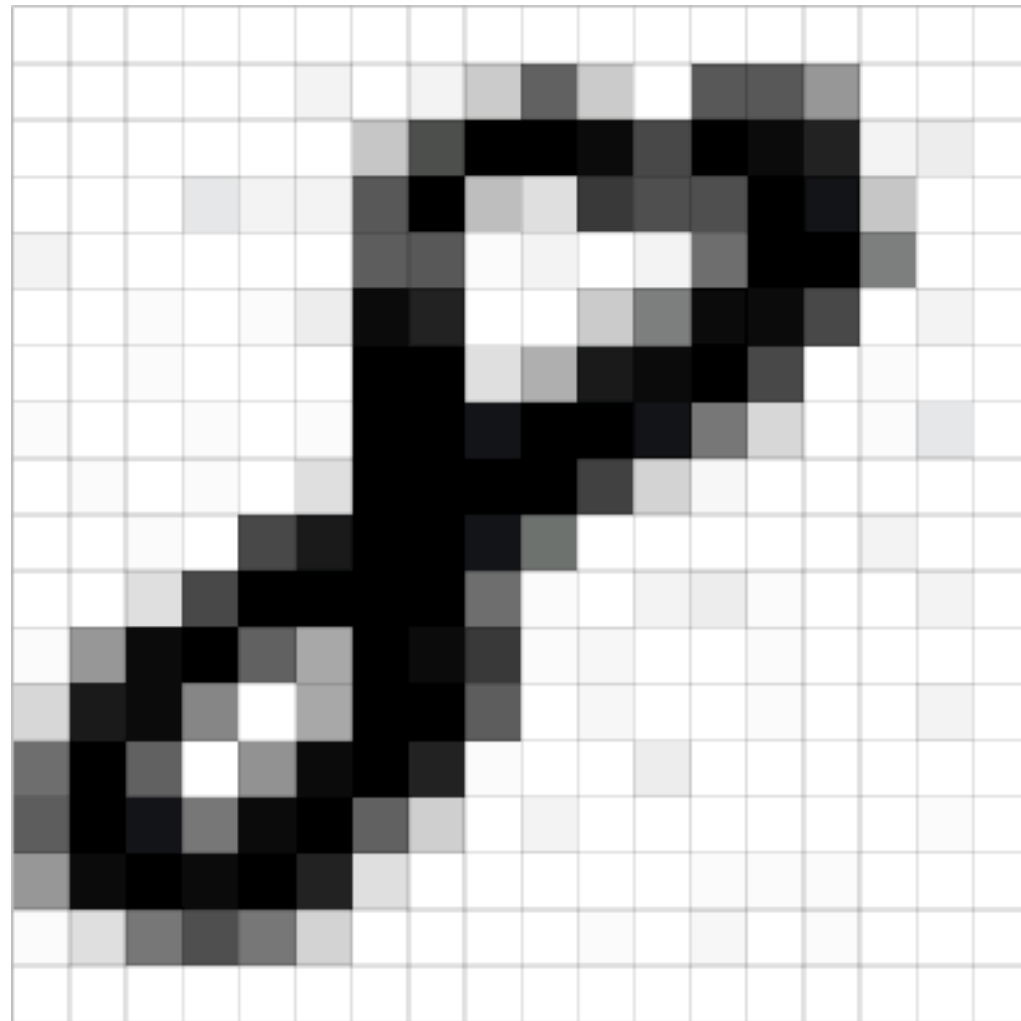What's inside?

# BASIC BLOCKS

➤ Input layer,

➤ Convolution layer,

➤ Activation layer,

➤ Pooling layer,

➤ Fully Connected layer.

# FUNDAMENTAL CONCEPT



Input          Conv          Pool          Conv          Pool          FC    FC    Softmax

32x32x3 image
5x5x3 filter $w$

32

32

3

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

**activation maps**

28

28

1

# WHAT IS OUR INPUT?
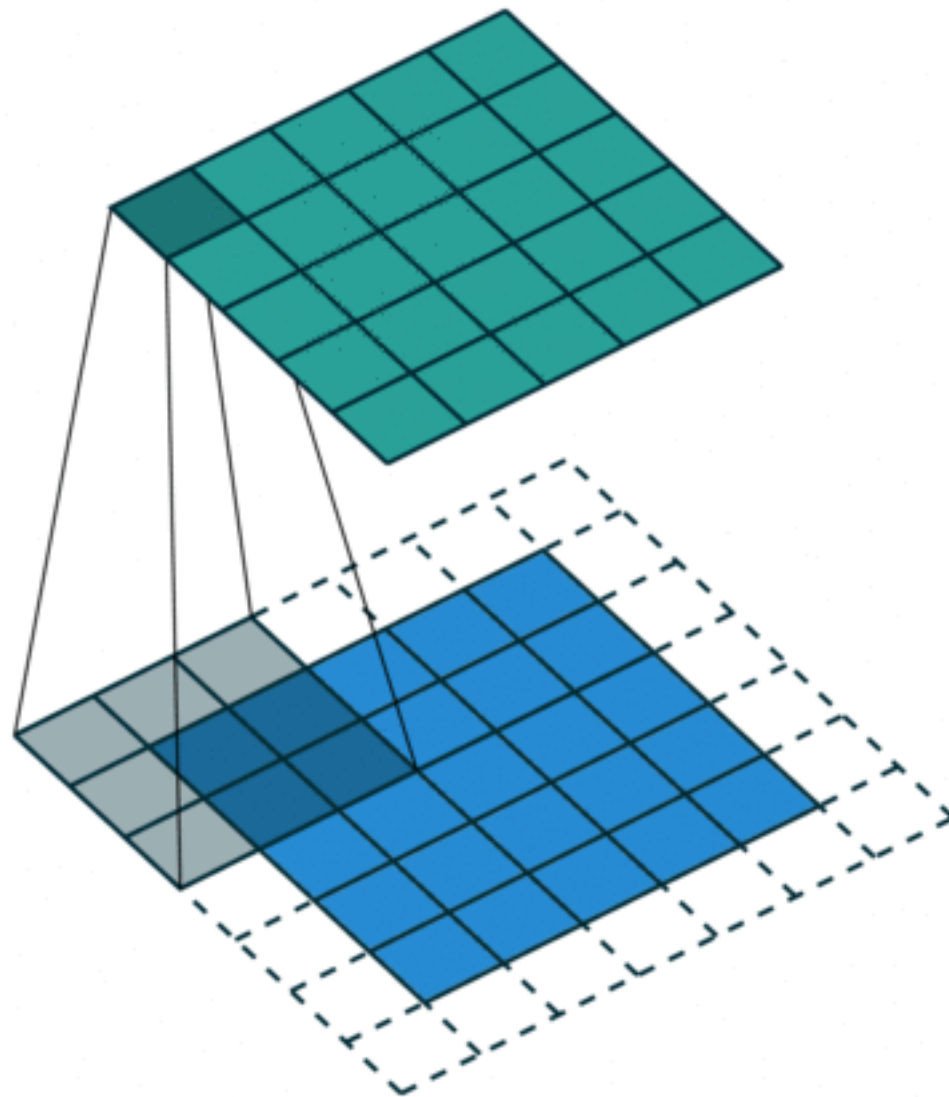
# CONVOLUTION



➤ Creates "**feature maps**",

➤ Apply **filters** on the image,

➤ Move such filter over the image and calculate **feature**,

➤ Follow the **stride** (how many fields it should "jump"),

➤ Is defined by **kernel size** (filter size),
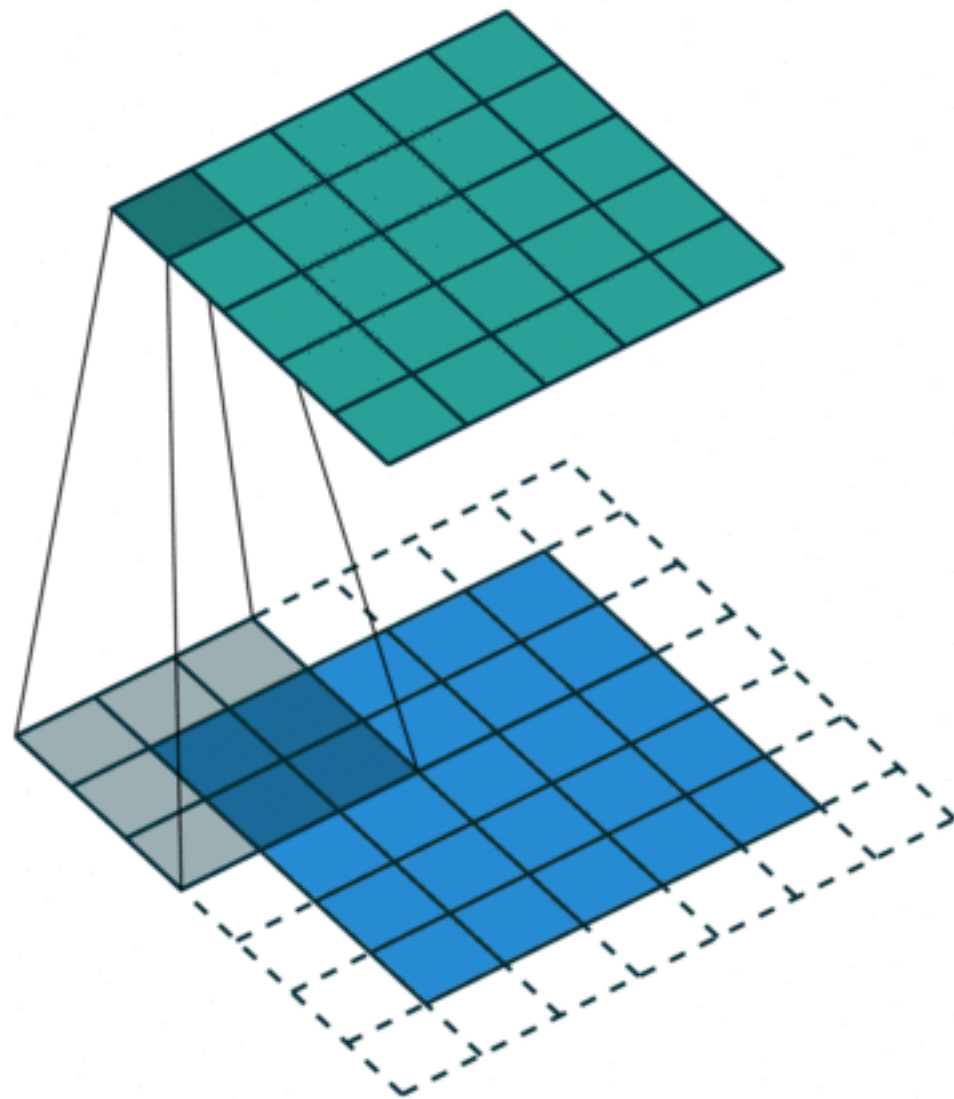
➤ Can use **padding** for bigger receptive field.

Image

Convolved
Feature

Input
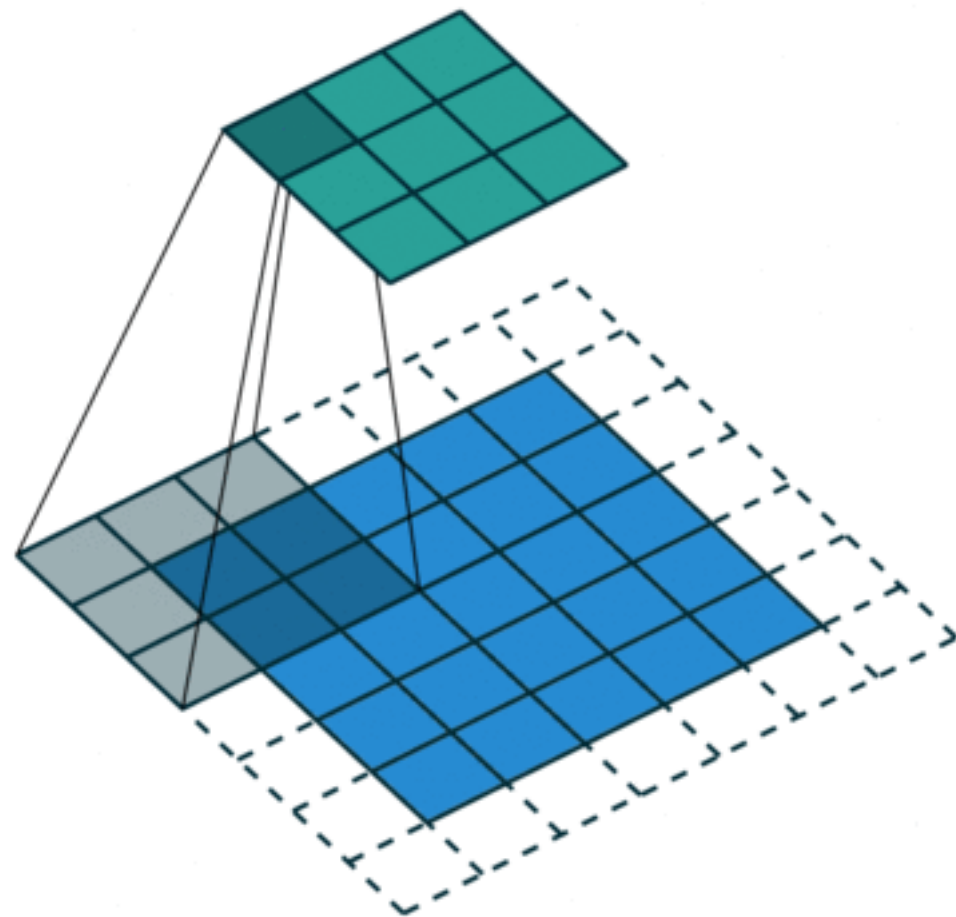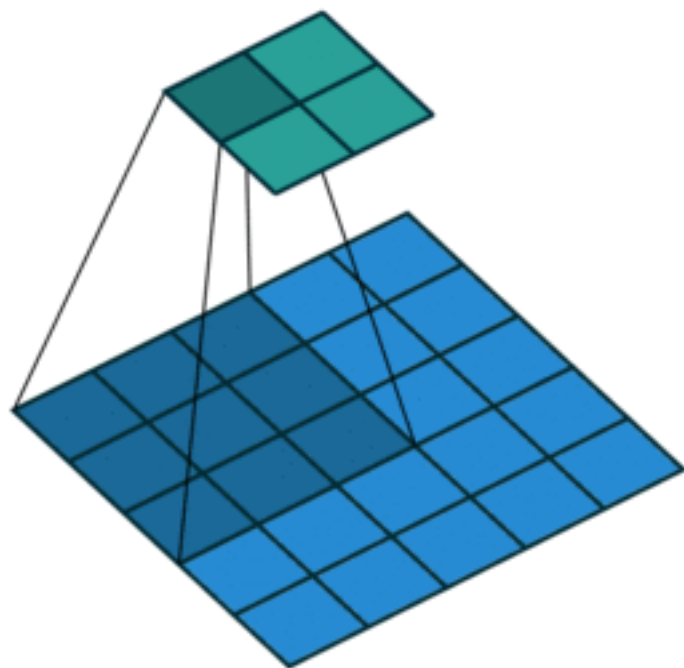
*Stride = 1*

*Stride = 2*
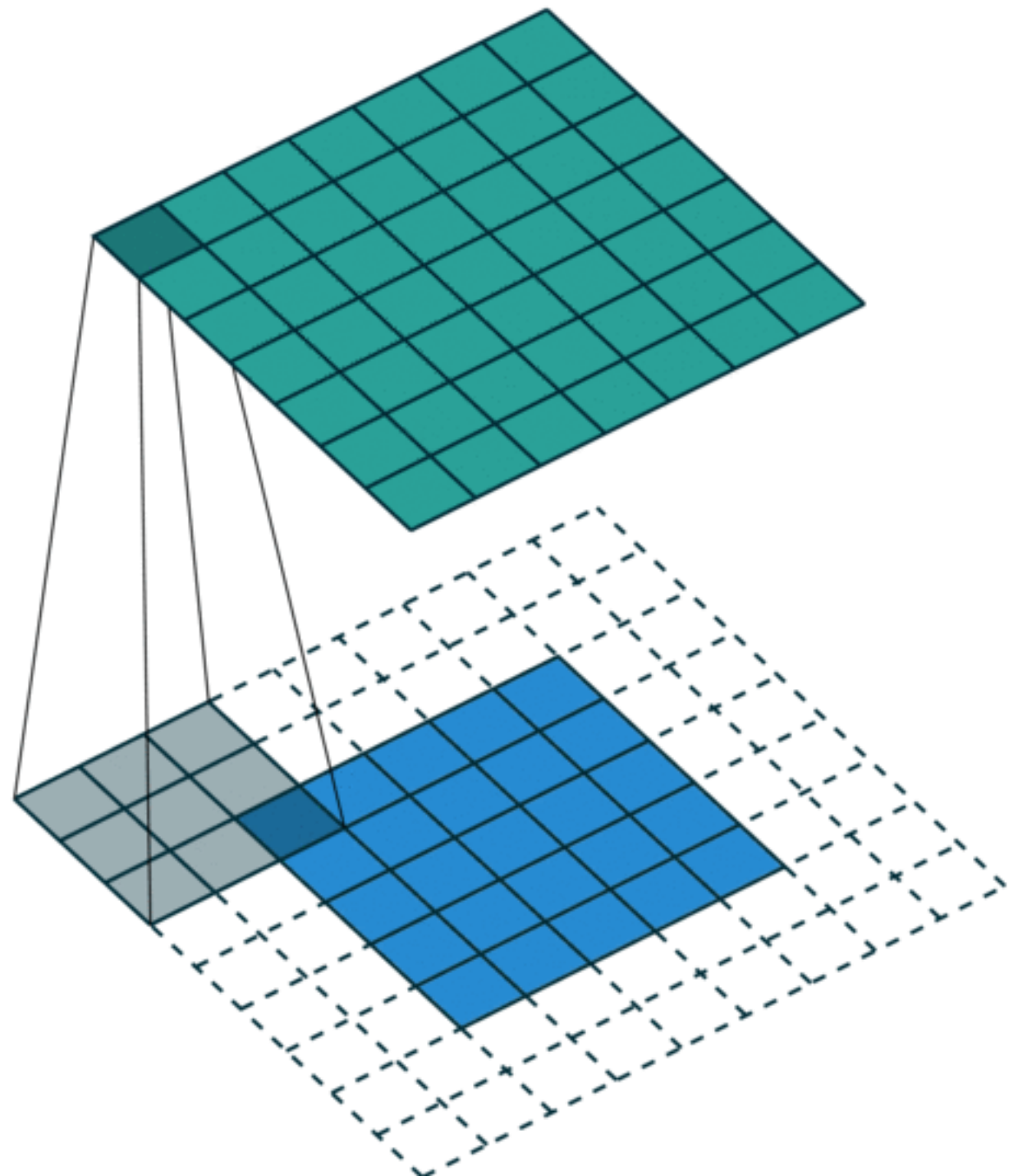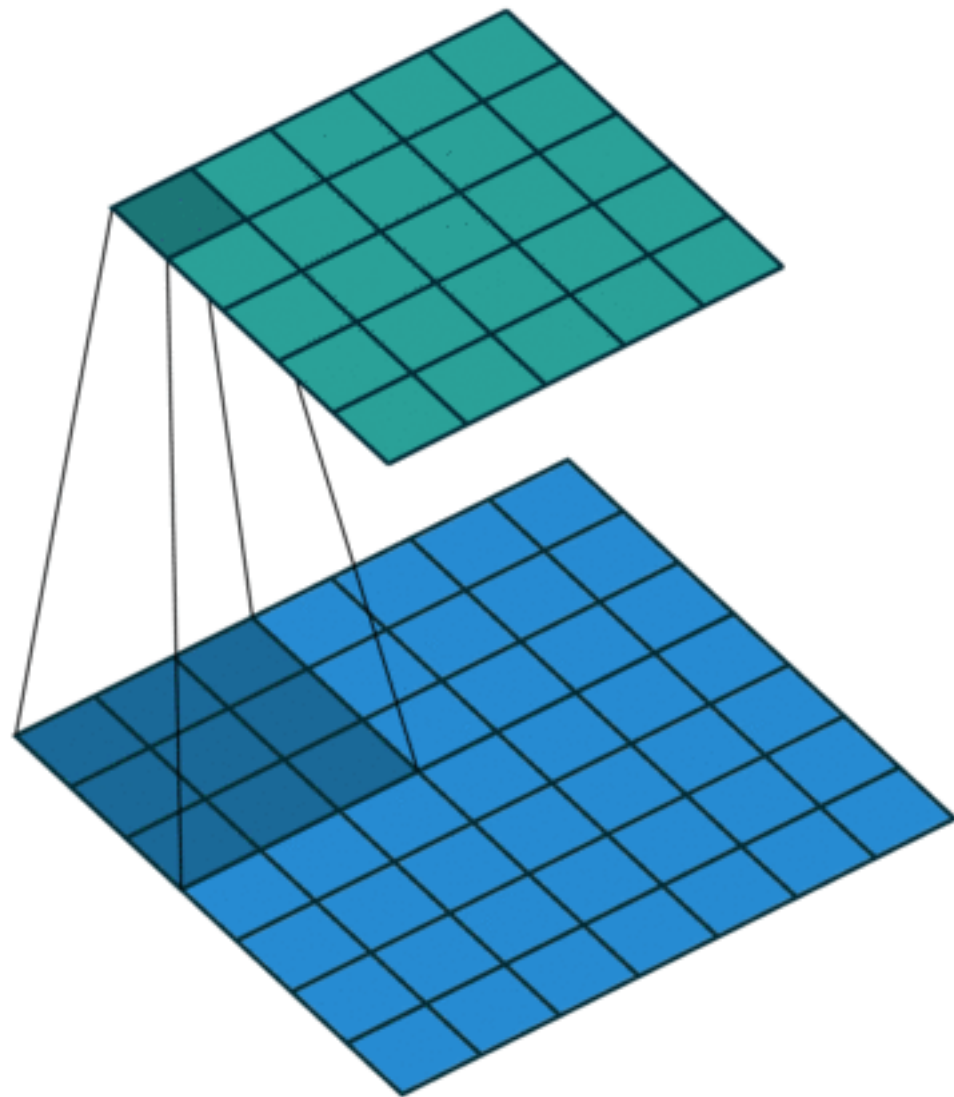
# PADDING

Padding = 0

Padding = 2

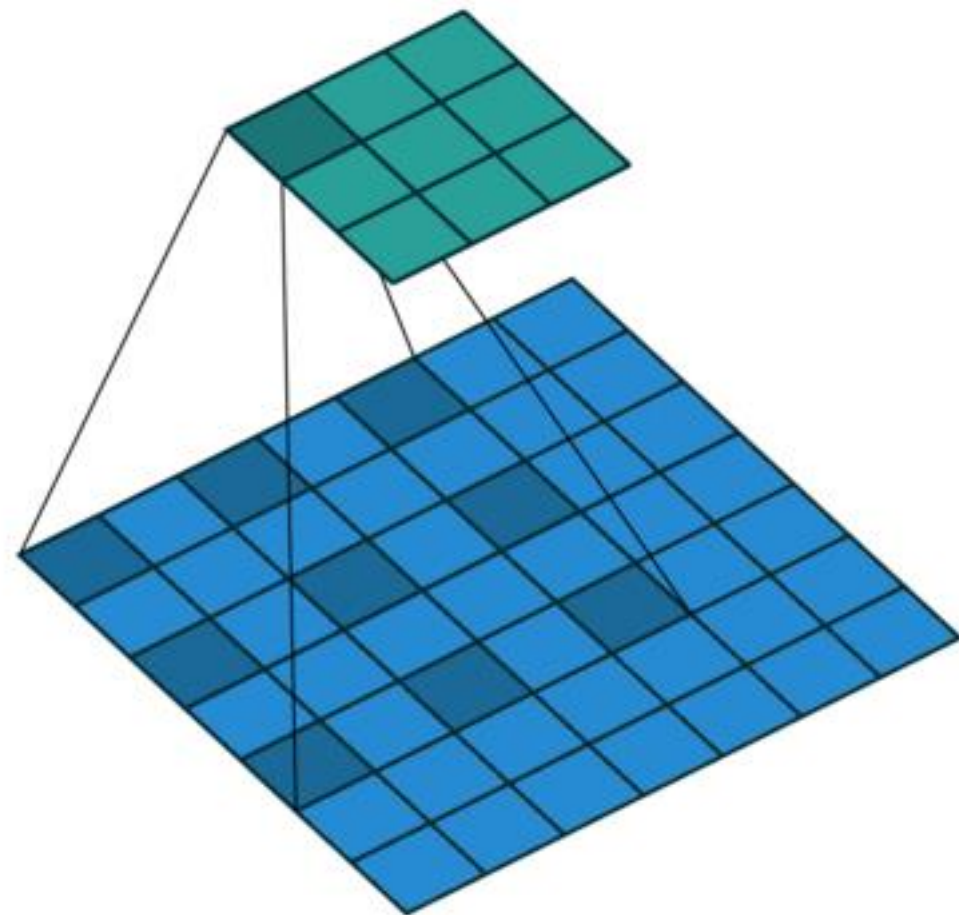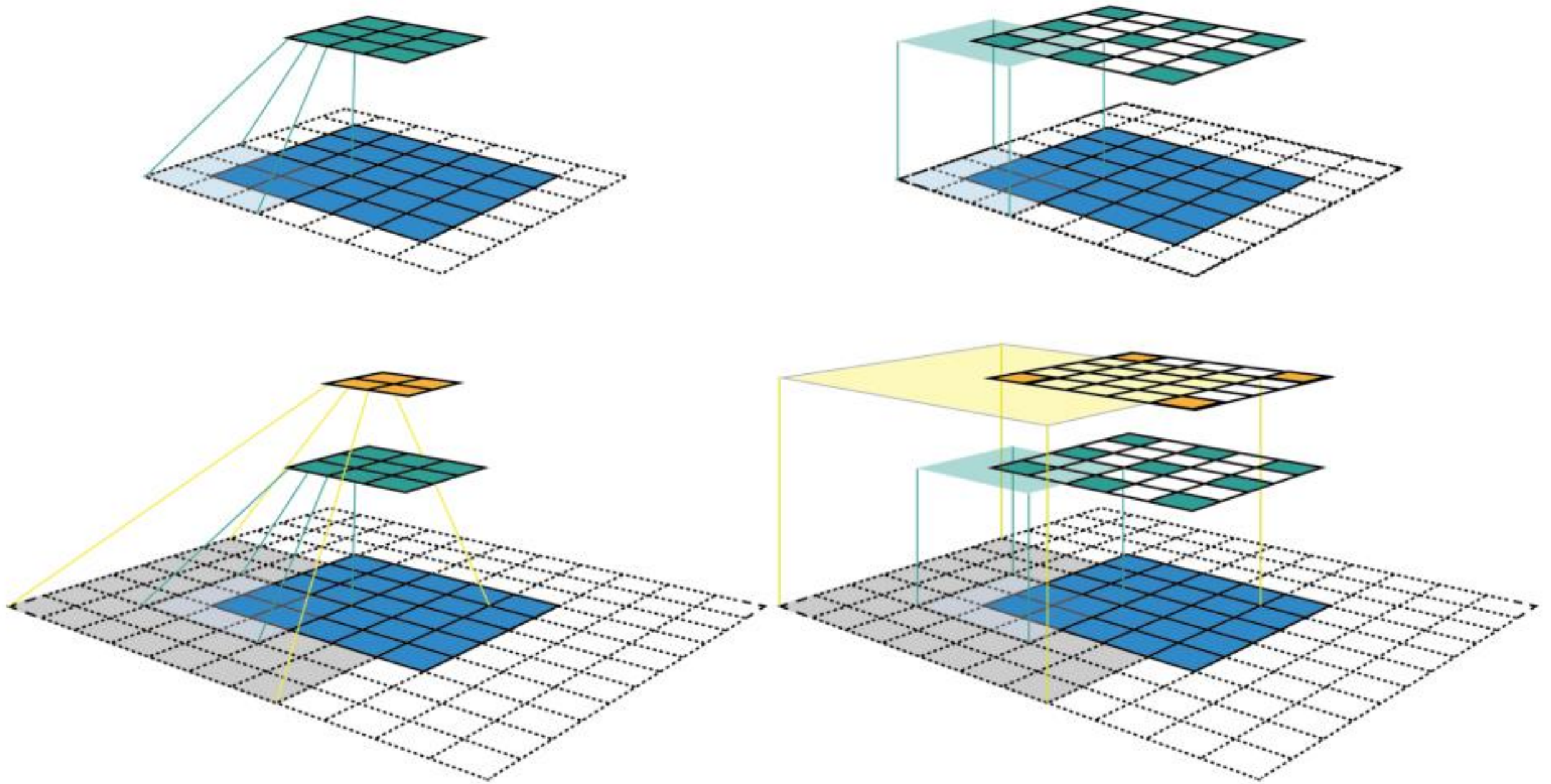# DILATION

Dilation = 1

Dilation = 2

# EQUATION FOR NUMBER OF OUTPUT FEATURES

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

$n_{in}$:　number of input features
$n_{out}$: number of output features
$k$:　　convolution kernel size
$p$:　　convolution padding size
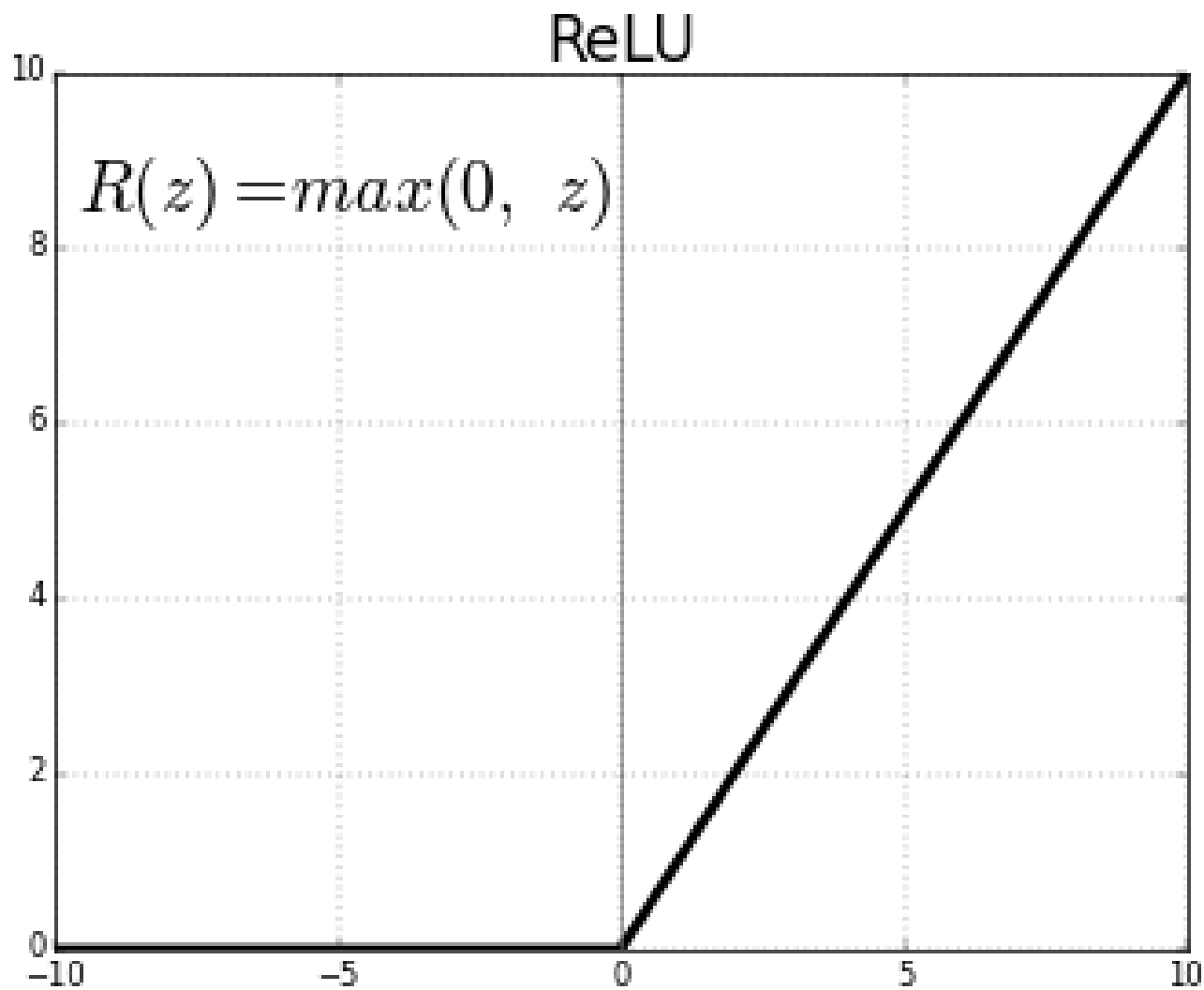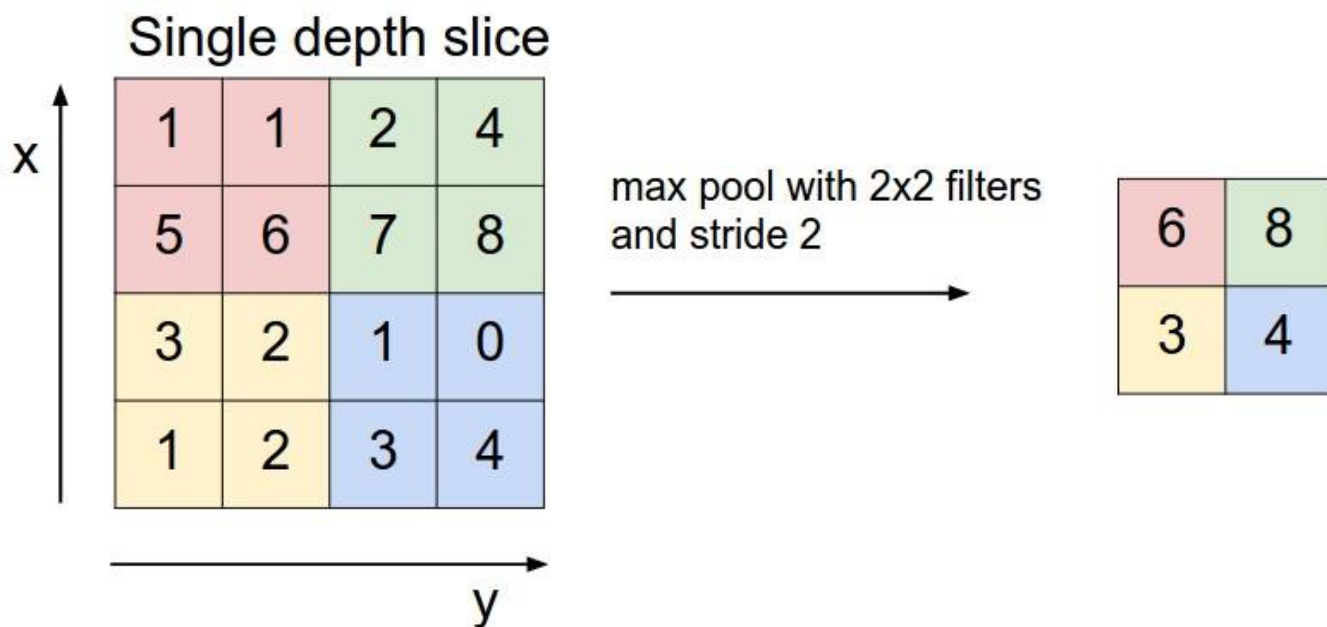$s$:　　 convolution stride size

ReLU

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

➤ **ReLU** is a default way to go,

➤ Some similar layers like:

  ➤ Leaky ReLU,

  ➤ Maxout,

➤ There are many other layers but have many problems (vanishing gradients, etc.).
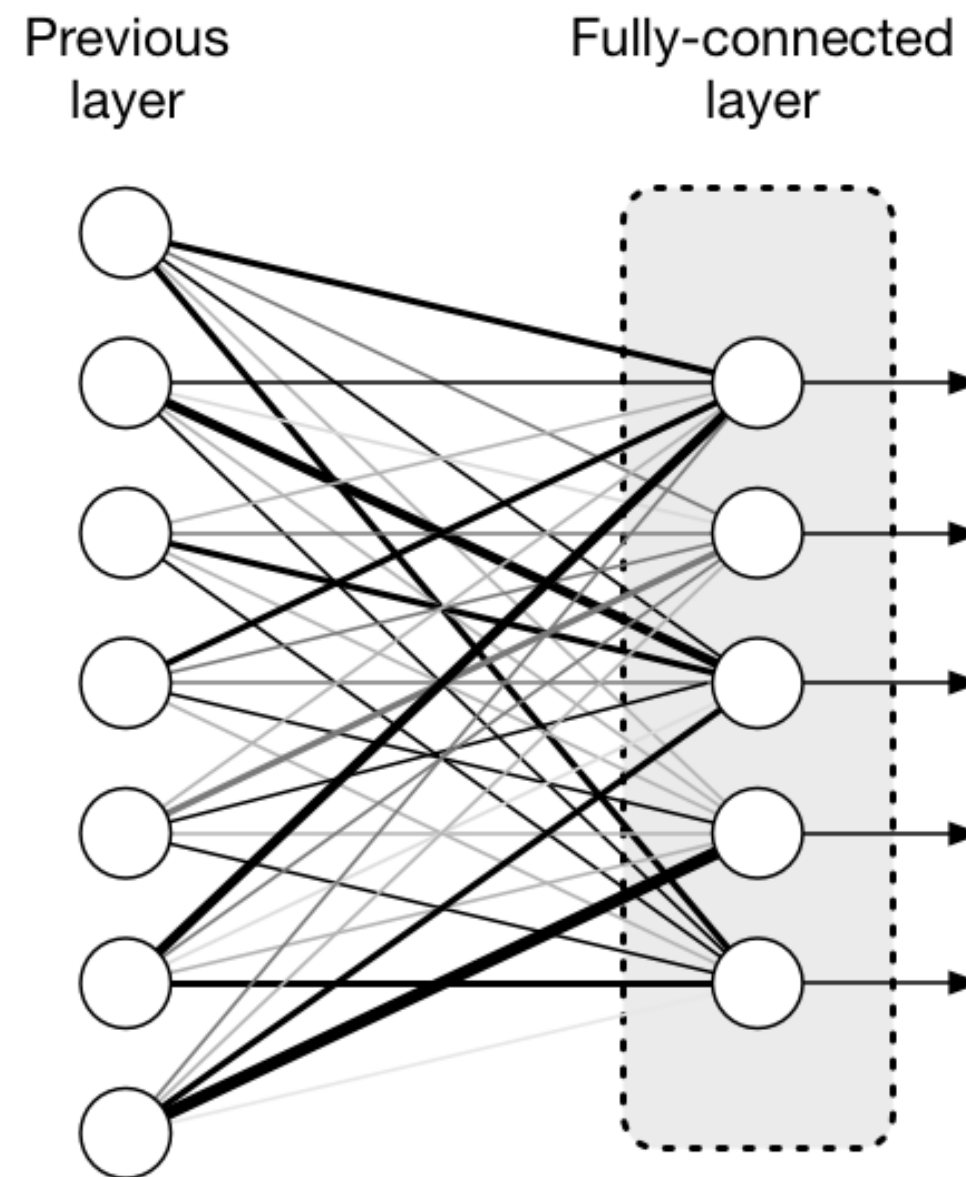
# POOLING



Single depth slice

| x | | | |
|---|---|---|---|
| 1 | 1 | 2 | 4 |
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

y

max pool with 2x2 filters
and stride 2

→

| 6 | 8 |
|---|---|
| 3 | 4 |

➤ Pooling **reduces** spatial space,

➤ **Reduces** amount of parameters,

➤ Reduces overfitting,

➤ A simple **routing** (during back propagation),

➤ Most common: MaxPooling,
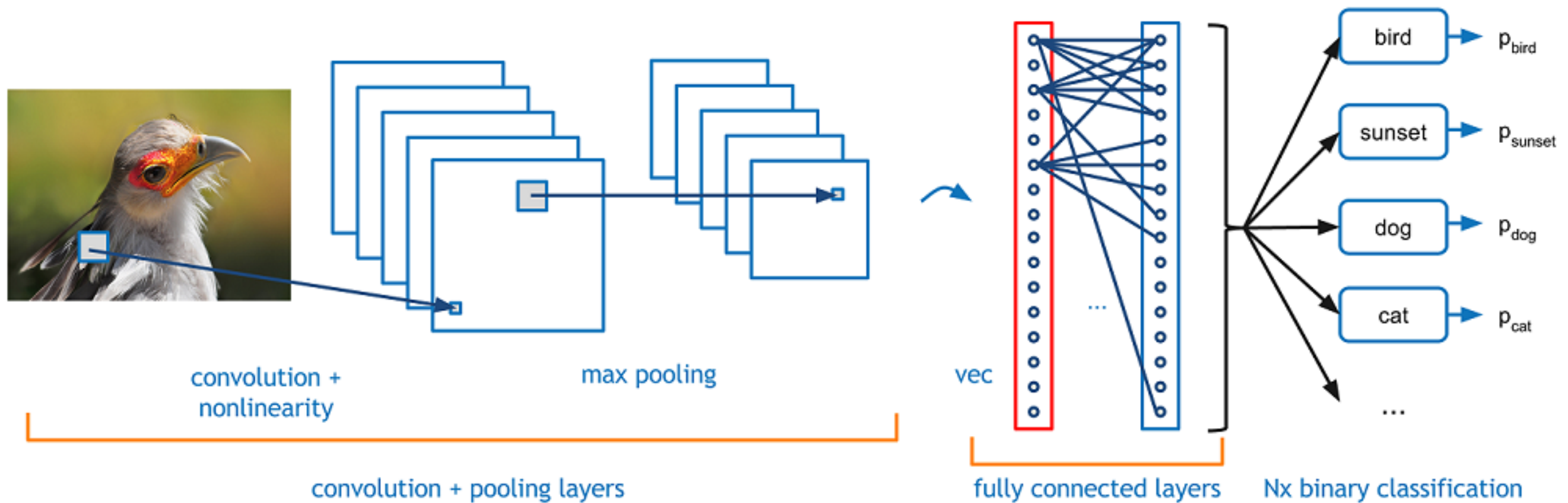
➤ Also: AvgPooling, ...

# FULLY CONNECTED LAYER
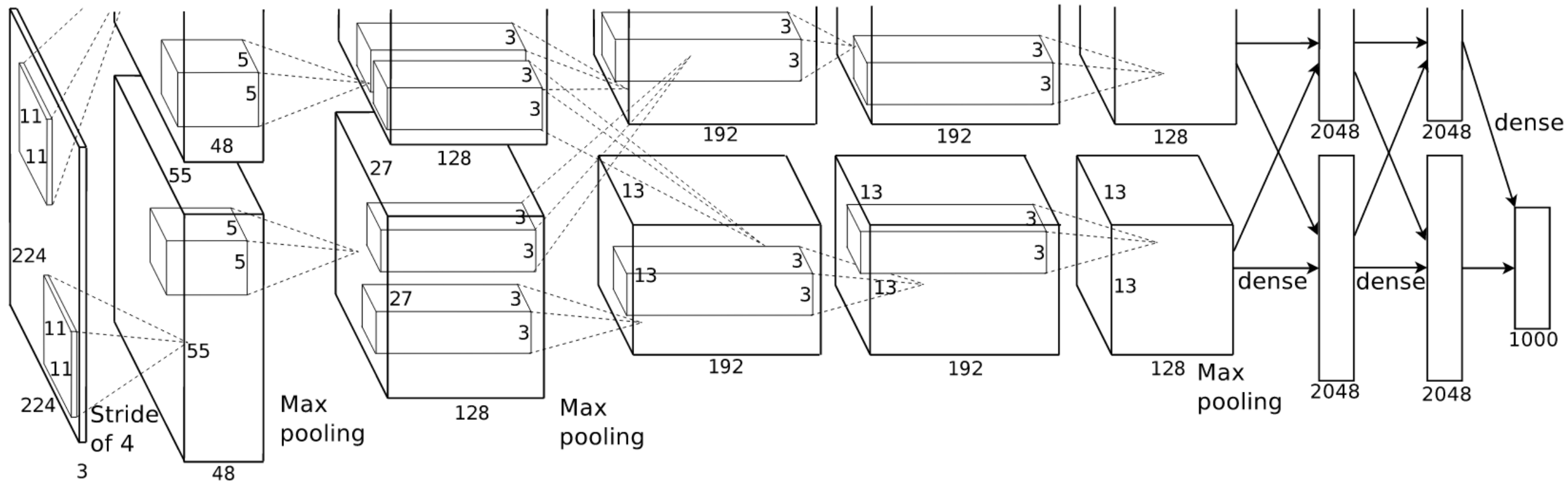


Previous layer

Fully-connected layer

convolution +
nonlinearity

max pooling

vec

convolution + pooling layers

fully connected layers

Nx binary classification

bird → $p_{bird}$

sunset → $p_{sunset}$

dog → $p_{dog}$

cat → $p_{cat}$

...

# POPULAR ARCHITECTURES

...that might be helpful!

$224 \times 224 \times 3$   $224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$   $1 \times 1 \times 1000$

convolution+ReLU
max pooling
fully connected+ReLU
softmax

# GOOGLENET



**Full Inception module**

$\mathcal{F}(\mathbf{x})$

$\mathbf{x}$ identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$

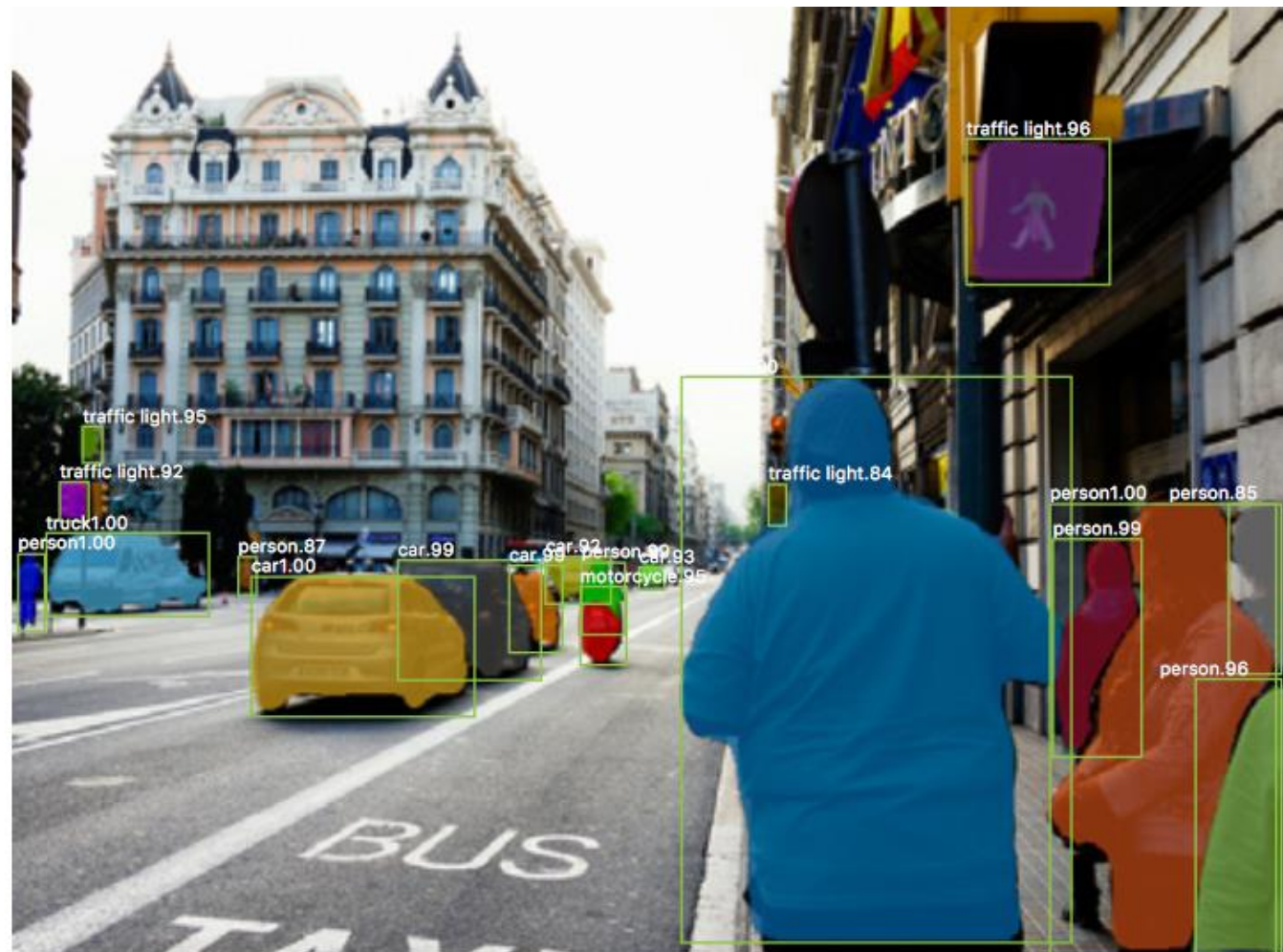# JOURNEY OF CNN IN 3 YEARS

1. R-CNN

2. Fast R-CNN

3. Faster R-CNN

4. Mask R-CNN

Worth reading more!



R-CNN: *Regions with CNN features*

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

R-CNN workflow

# AUTOENCODERS



Original input → Encoder → Compressed representation → Decoder → Reconstructed input

Generative Adversarial Network

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|---|---|---|---|---|---|---|
| conv1 | 112×112 | 7×7, 64, stride 2 | | | | |
| conv2_x | 56×56 | 3×3 max pool, stride 2 | | | | |
| conv2_x | 56×56 | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 64 \\ 3\times3,\ 64 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 64 \\ 3\times3,\ 64 \\ 1\times1,\ 256 \end{bmatrix}\times3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 128 \\ 3\times3,\ 128 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times4$ | $\begin{bmatrix} 1\times1,\ 128 \\ 3\times3,\ 128 \\ 1\times1,\ 512 \end{bmatrix}\times8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 256 \\ 3\times3,\ 256 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times6$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times23$ | $\begin{bmatrix} 1\times1,\ 256 \\ 3\times3,\ 256 \\ 1\times1,\ 1024 \end{bmatrix}\times36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times2$ | $\begin{bmatrix} 3\times3,\ 512 \\ 3\times3,\ 512 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ | $\begin{bmatrix} 1\times1,\ 512 \\ 3\times3,\ 512 \\ 1\times1,\ 2048 \end{bmatrix}\times3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | $1.8\times10^9$ | $3.6\times10^9$ | $3.8\times10^9$ | $7.6\times10^9$ | $11.3\times10^9$ |

# HOW TO DEAL WITH ALL OF THIS?

# AI OPTIMIZED HARDWARE – TRAINING AND INFERENCE



intel Neural Compute Stick 2

## GAUDI™

**Purpose-Built for AI Training**

The only AI processor with Integrated RDMA over Converged Ethernet to provide scalability and lower total cost of ownership.
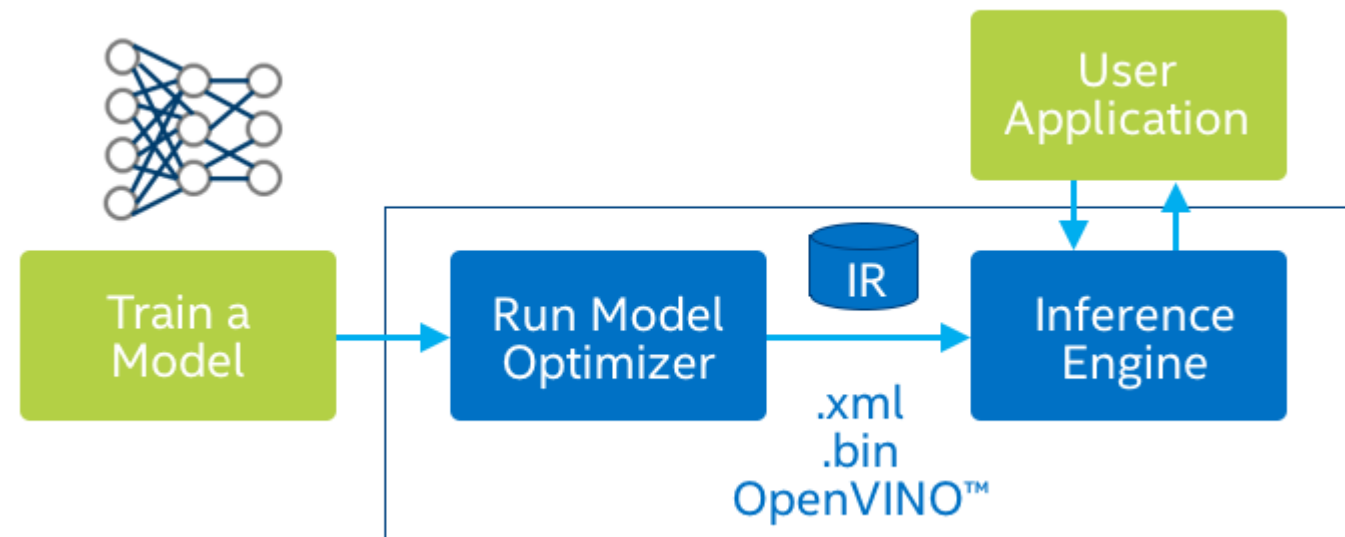
## GOYA™

**Purpose-Built for AI Inference**

15,453 images-per-second throughput on ResNet-50

# OPEN VINO AND MODEL PREPARATION PLATFORM
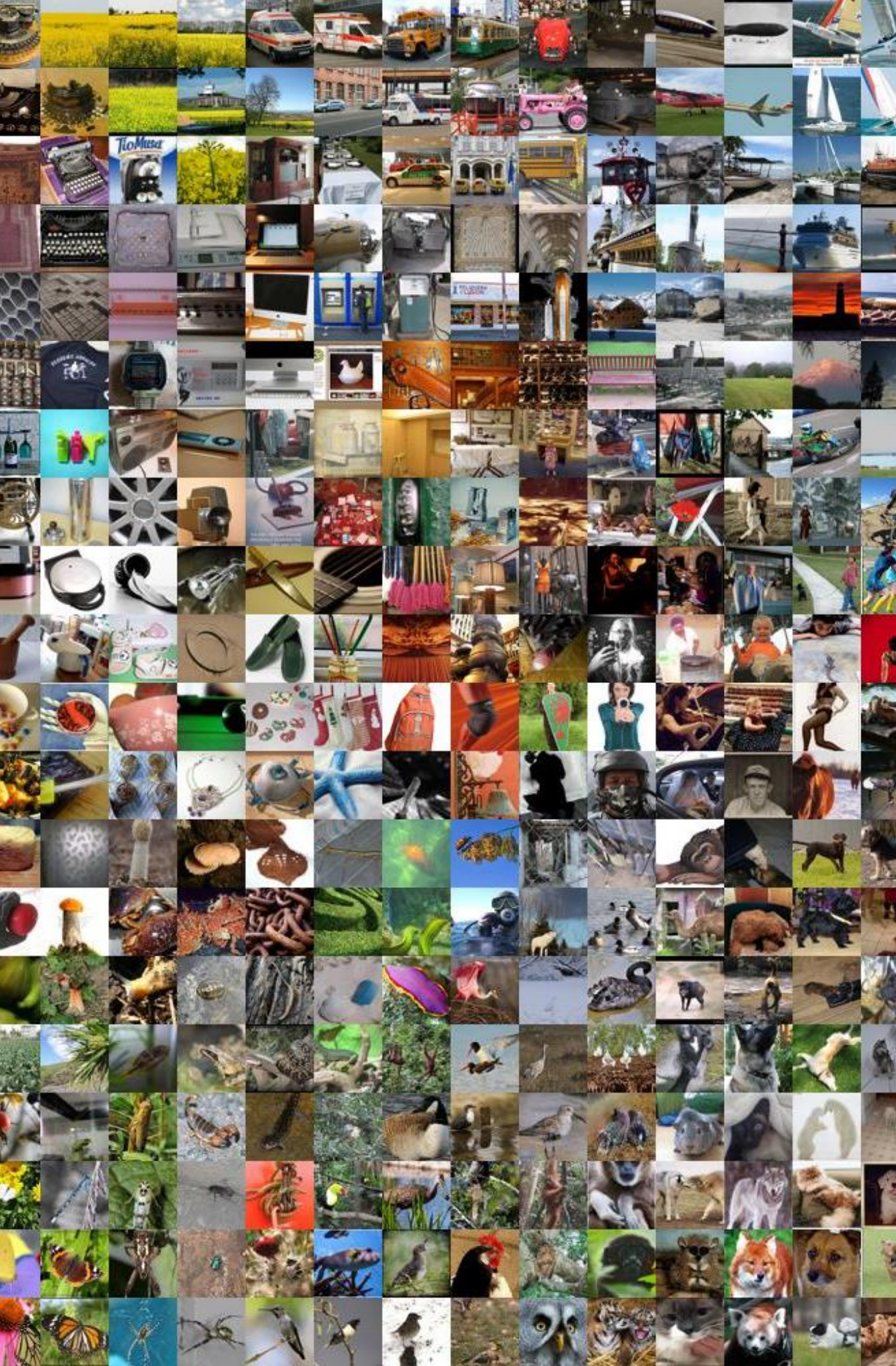
# POPULAR DATASETS

...to train your CNNs!

# IMAGENET

➤ 10 millions hand-labelled images,

➤ 1 million with bounding boxes,

➤ Labels based on WordNet (hierarchical dictionary).

# CIFAR10

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

➤ 60.000 images for training,

➤ 6000 images for testing,

➤ 10 classes,

➤ Also: Version with 100 classes (CIFAR100).

# KAGGLE

➤ Competitions,

➤ 100+ Datasets,

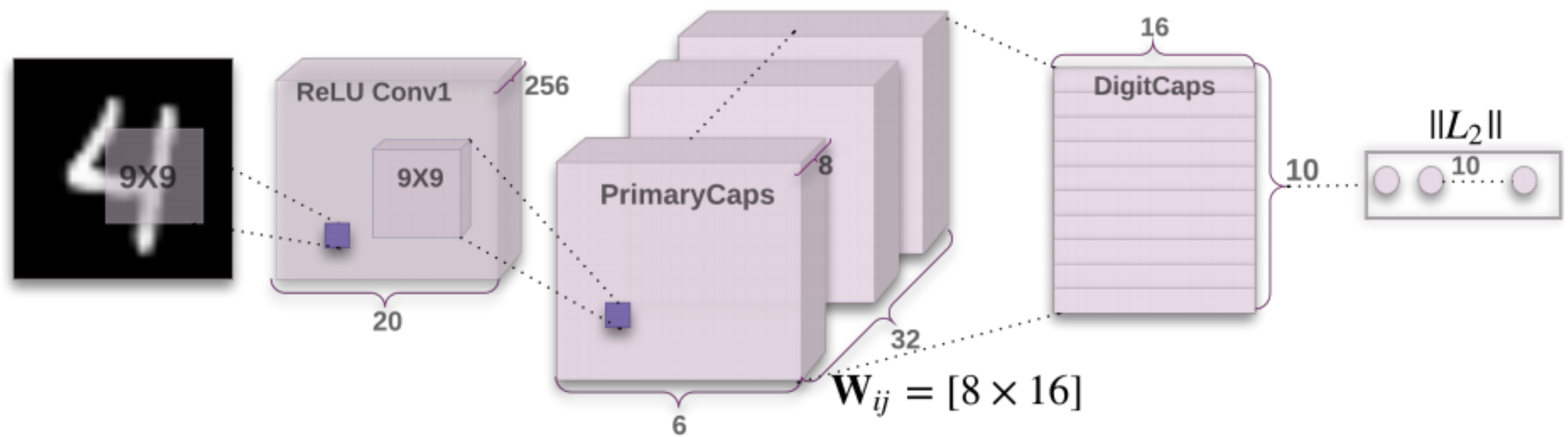➤ Community,

➤ Many code examples,

➤ Many CNNs challenges!

# WHAT'S NEXT?

Future of CNNs...

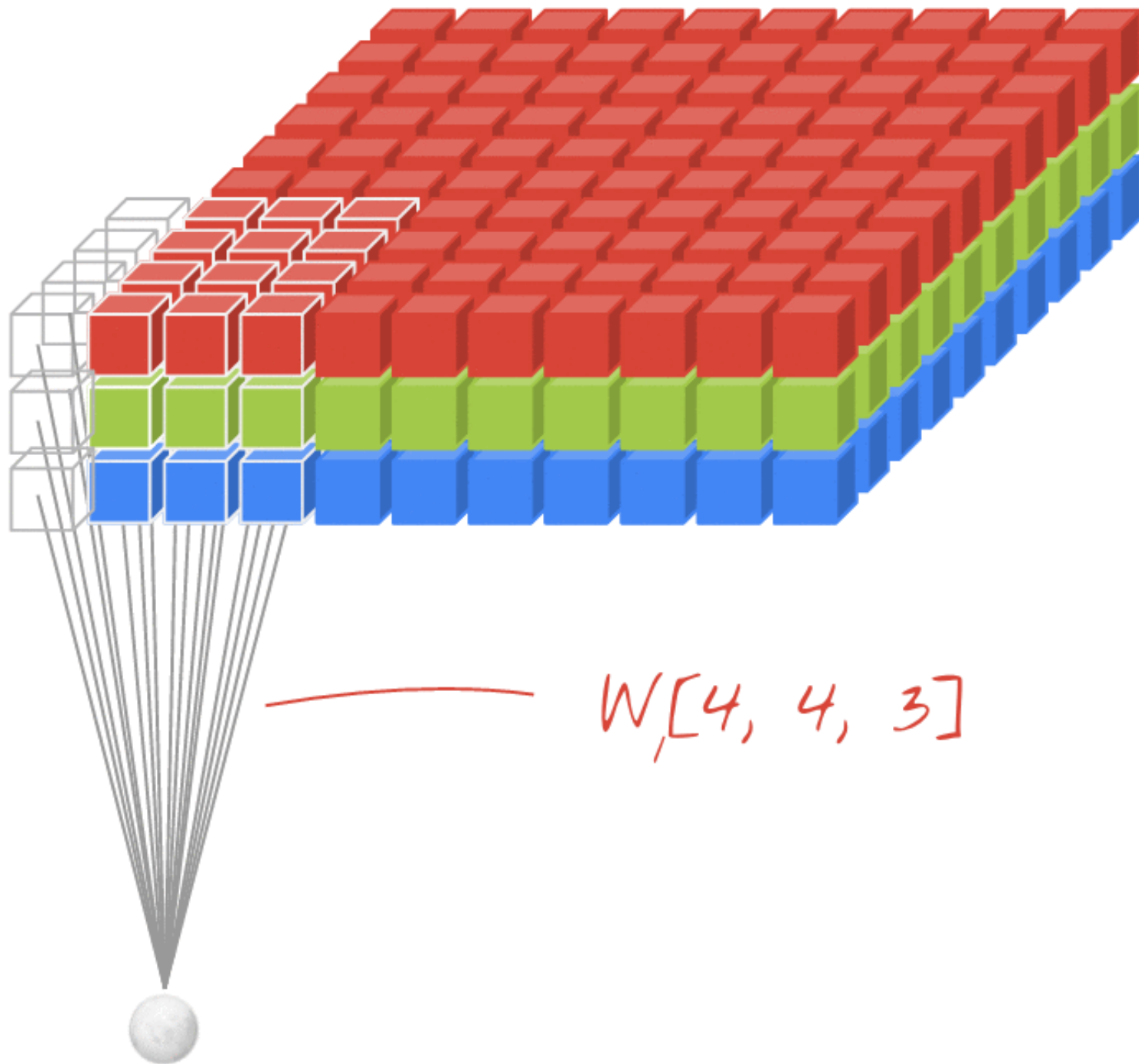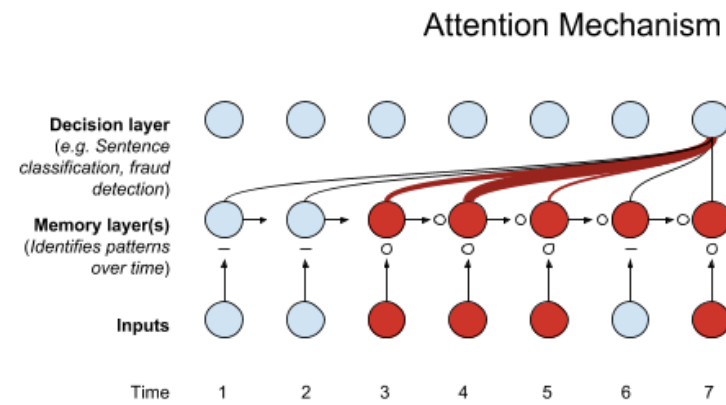$$W[4, 4, 3]$$

**Attention Mechanism**

# THANKS FOR ATTENTION!

Q&A
+
Let's code!