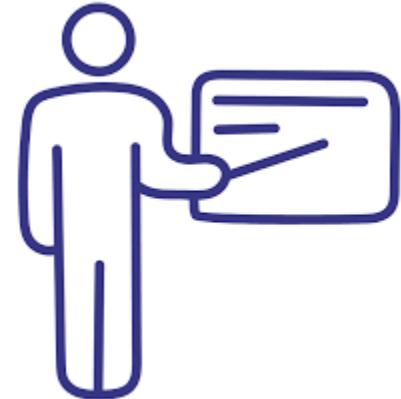


LESSON 23

Image classification
Intelligent vision System,
Object detection,
Segmentation and Recognition,
Machine Learning Hardware



Outline



- Image classification and Intelligent Vision Systems
 - Object detection, Segmentation, Recognition
 - Classical, Deep and Hybrid solutions
 - Use cases
- Hardware for Machine learning



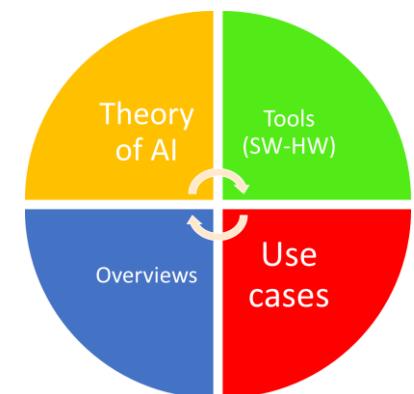
Dog



Cat

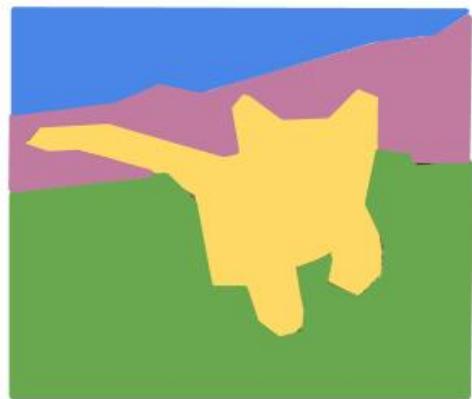
THEORY

Image classification



Intelligent Vision tasks

Semantic Segmentation



Classification + Localization



Object Detection



Instance Segmentation



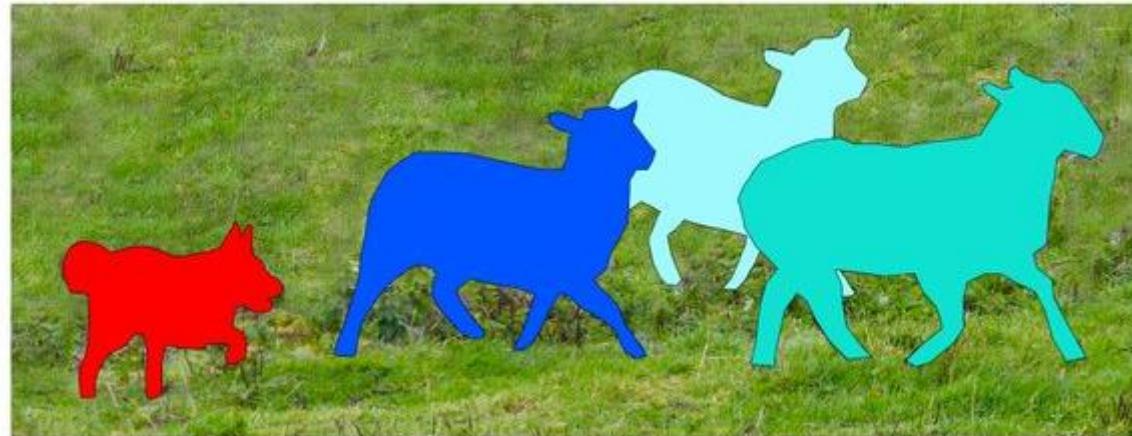
No objects, just pixels

Single Object

Multiple Object

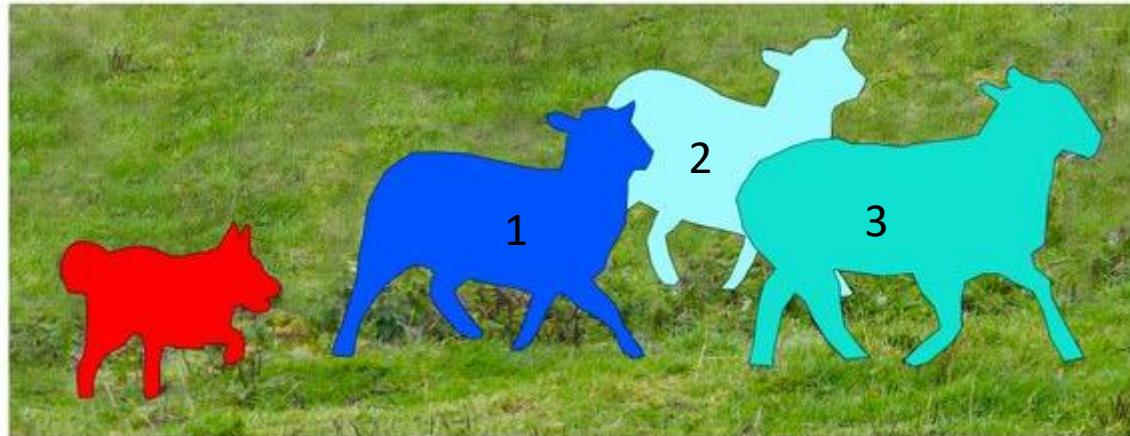
This image is CC0 public domain

A quick test 1 ...



What is this Intelligent vision task?

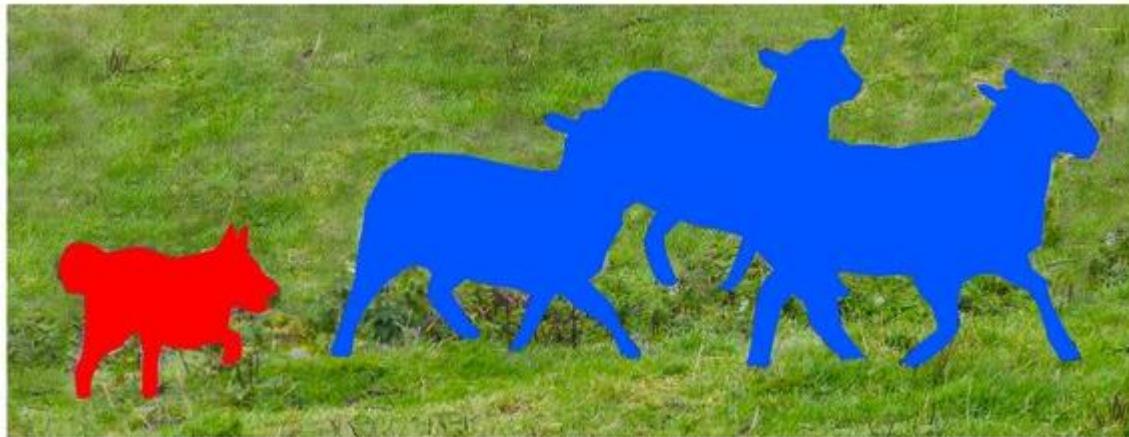
A quick test 1



What is this Intelligent vision task?

INSTANCE SEGMENTATION!

A quick test 2 ...



What is this Intelligent vision task?

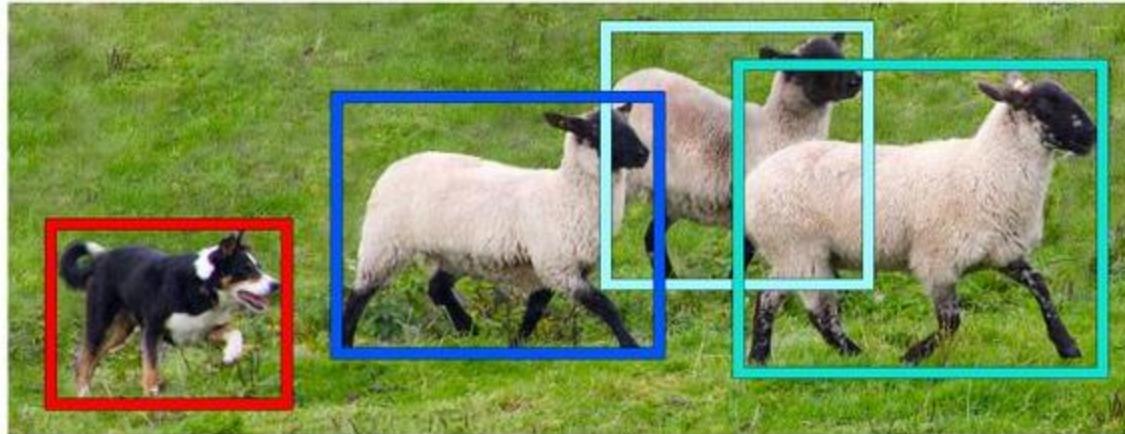
A quick test 2 ...



What is this Intelligent vision task?

SEMANTIC SEGMENTATION!

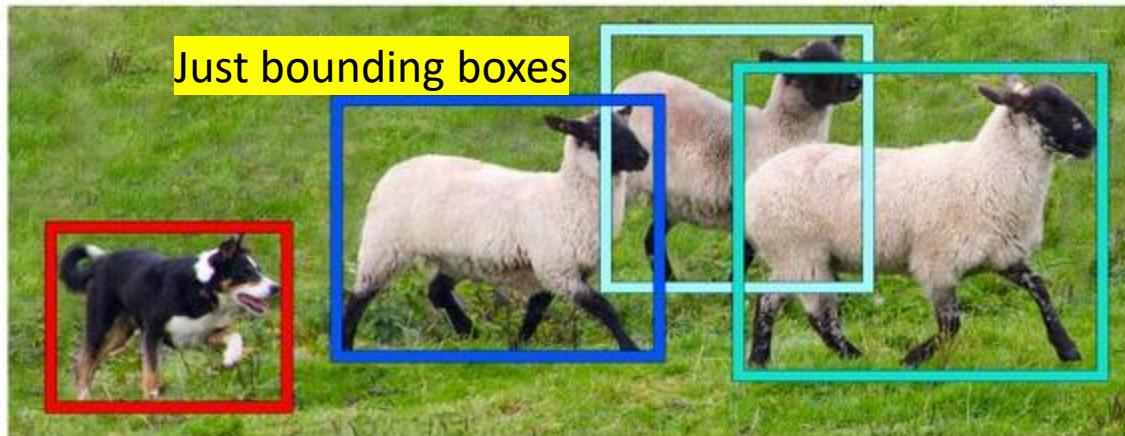
A quick test 3 ...



What is this Intelligent vision task?

A quick test 3 ...

We have just the so-called bounding box



What is this Intelligent vision task?

OBJECT DETECTION

A quick test 4 ...



What Intelligent vision task does it correspond?

A quick test 4 ...

We have just one class in output = Sheep



What is this Intelligent vision task?

IMAGE CLASSIFICATION

Direct image classification



Image classification

«Dog»

No more
information
about the object

Direct image classification: Example with AlexNet

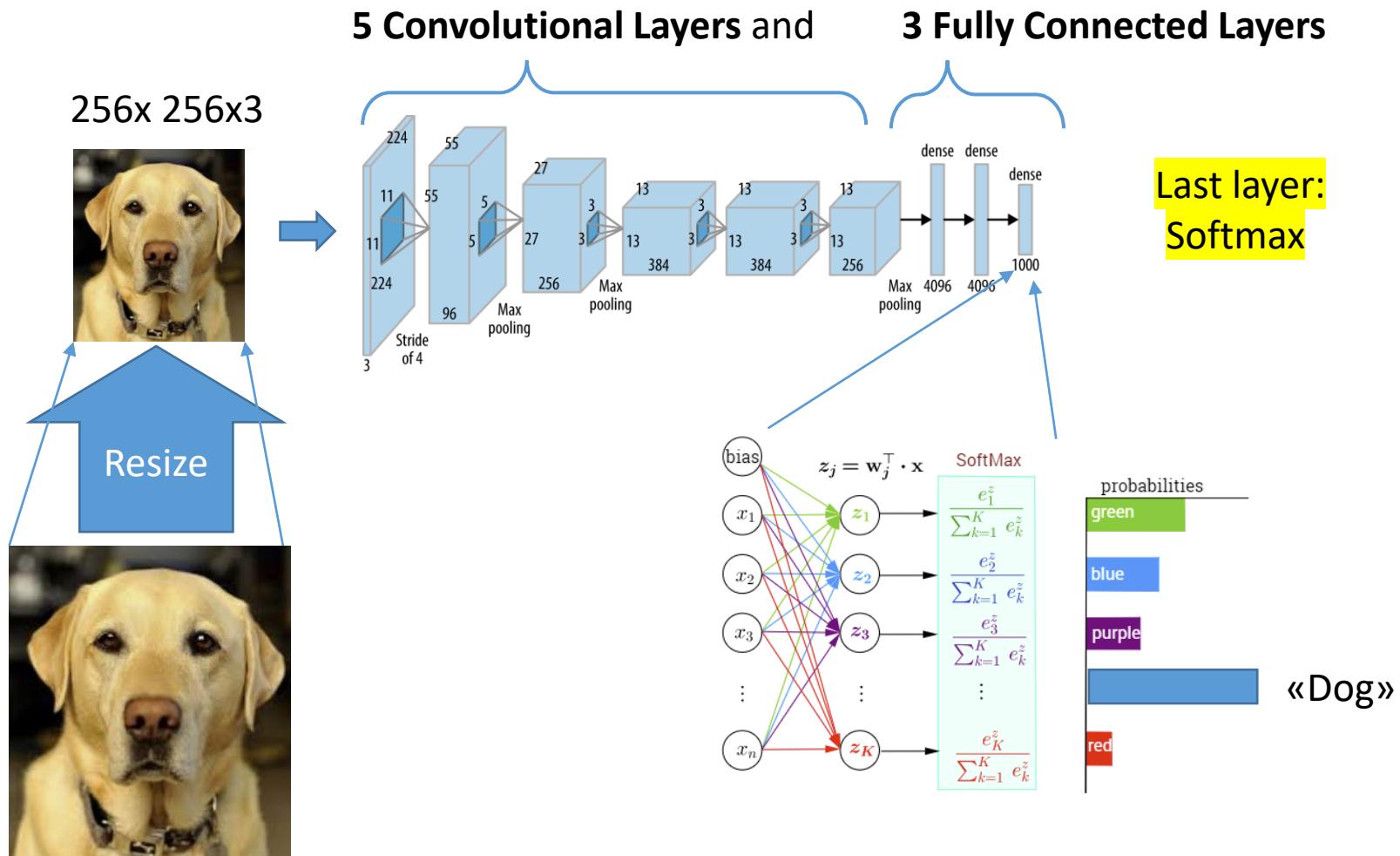


Image classification

Good for training
and fine tuning
(after selection)

EASY SET

red fox (100)



hen-of-the-woods (100)



ibex (100)



goldfinch (100)



flat-coated retriever (100)



tiger (100)



hamster (100)



porcupine (100)



stingray (100)



Blenheim spaniel (100)



HARD SET

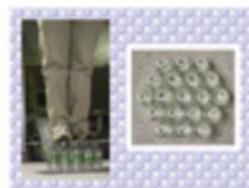
muzzle (71)



hatchet (68)



water bottle (68)



velvet (68)



loupe (66)



hook (66)



spotlight (66)



ladle (65)



restaurant (64)



letter opener (59)



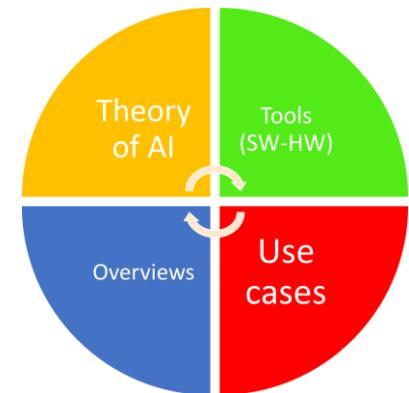
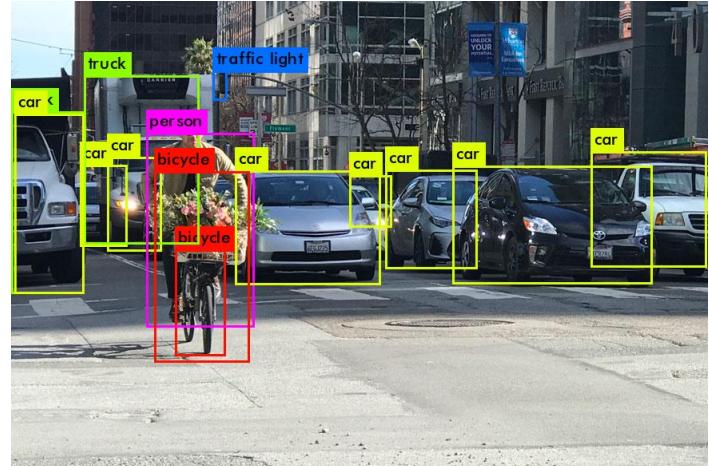
Good for Google...
not for an industrial
application...

But good to train to
extract feature!



THEORY

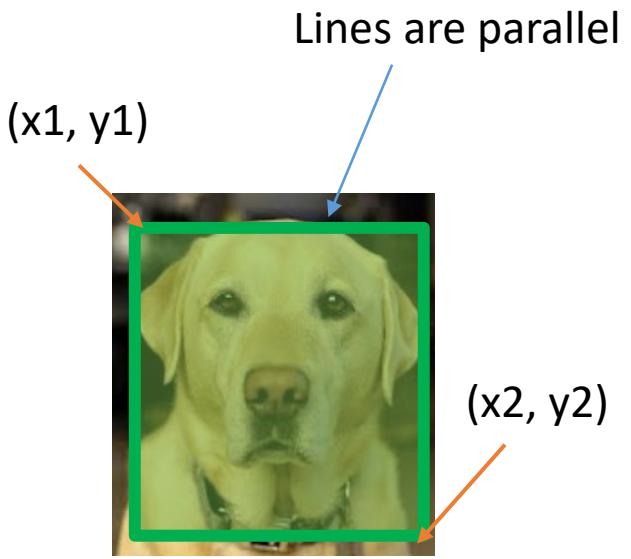
Object detection



Object detection



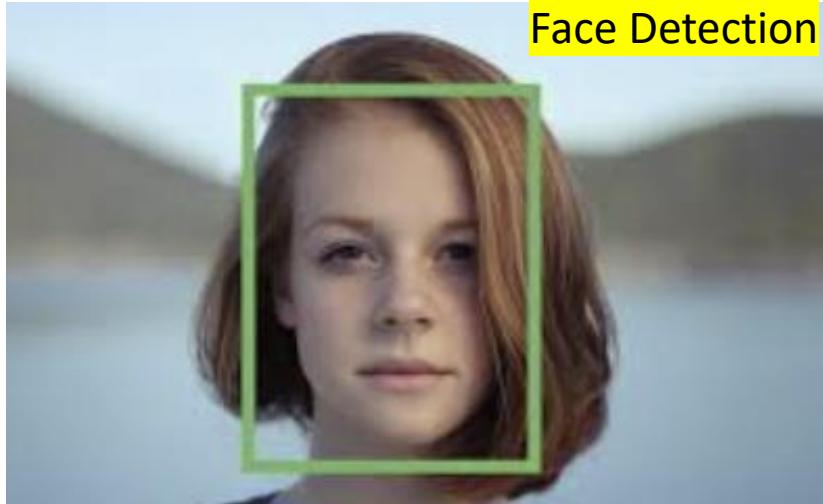
Object detection



Output is the
bounding box
 $\{x_1, y_1, x_2, y_2\}$

It's almost impossible to get measurements

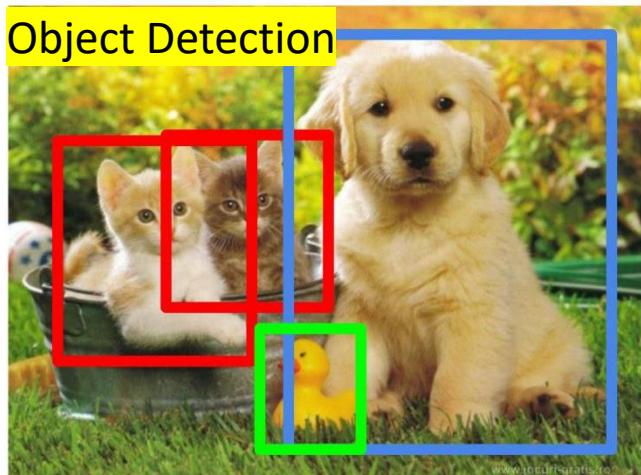
Detection in not segmentation!



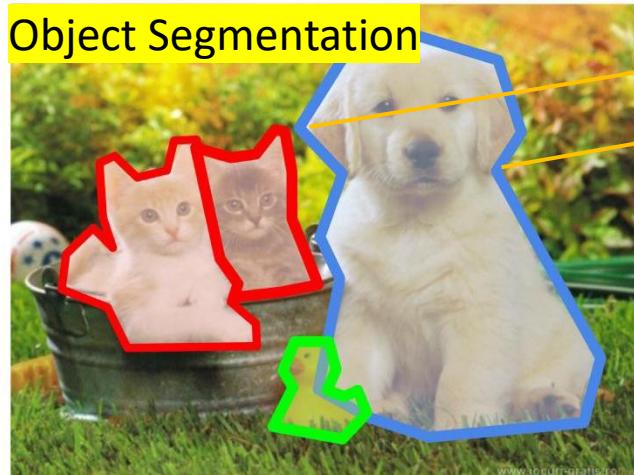
Face Detection



Face Segmentation



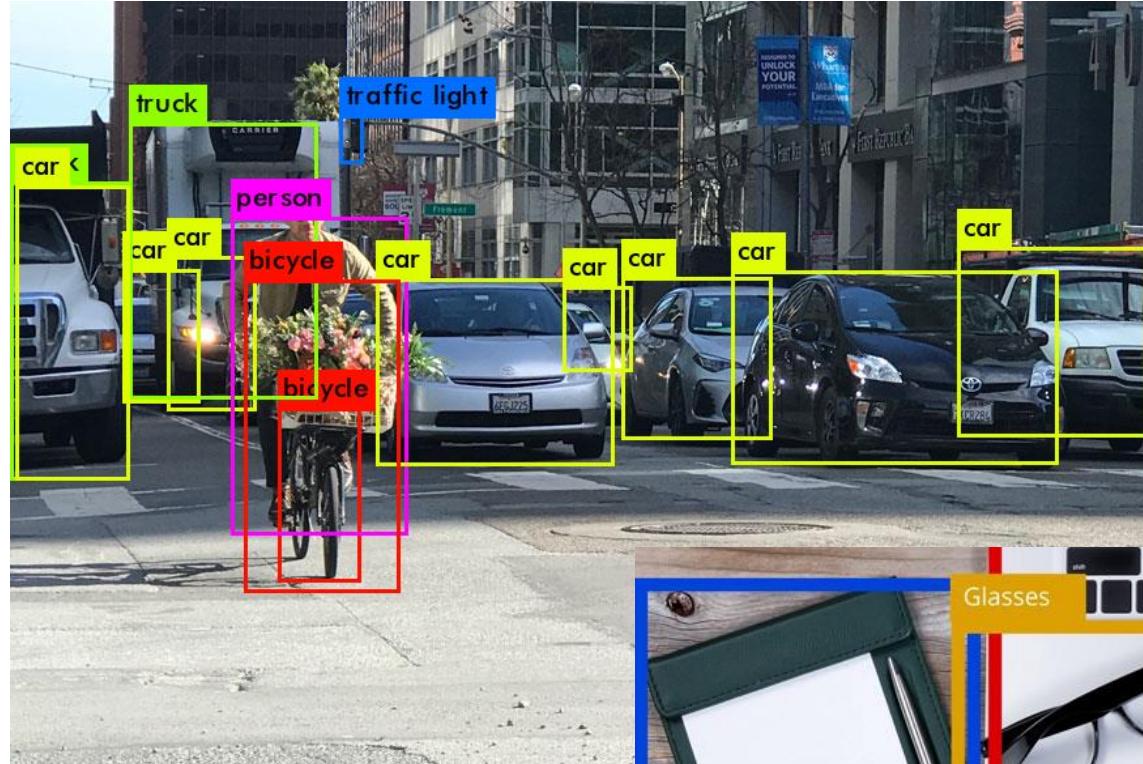
Object Detection



Object Segmentation

Real measures
of the objects!

Multiple object detection



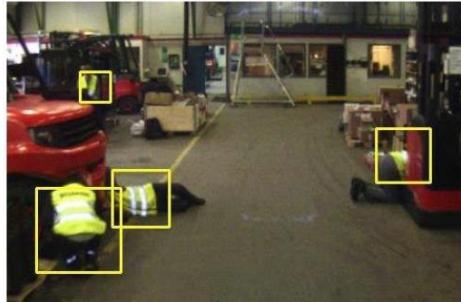
Not an easy task
in complex scenes

Overlapping objects



Detection

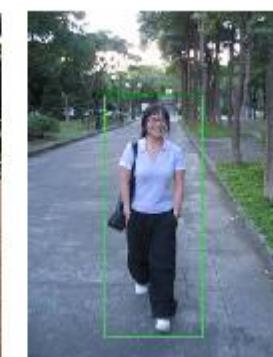
Example: person detection in images



Industrial sites



Surveillance



Articulated body pose estimation: OpenPose

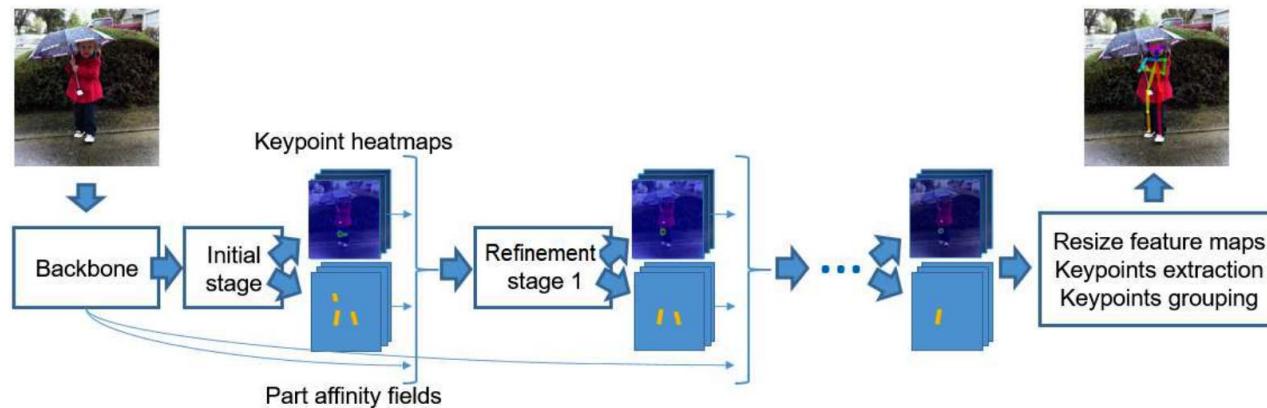
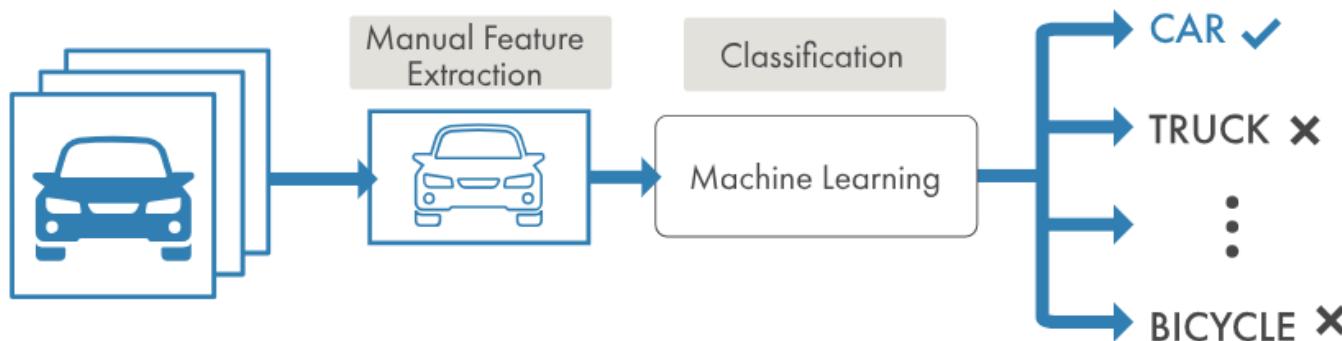


Figure 1: OpenPose pipeline.

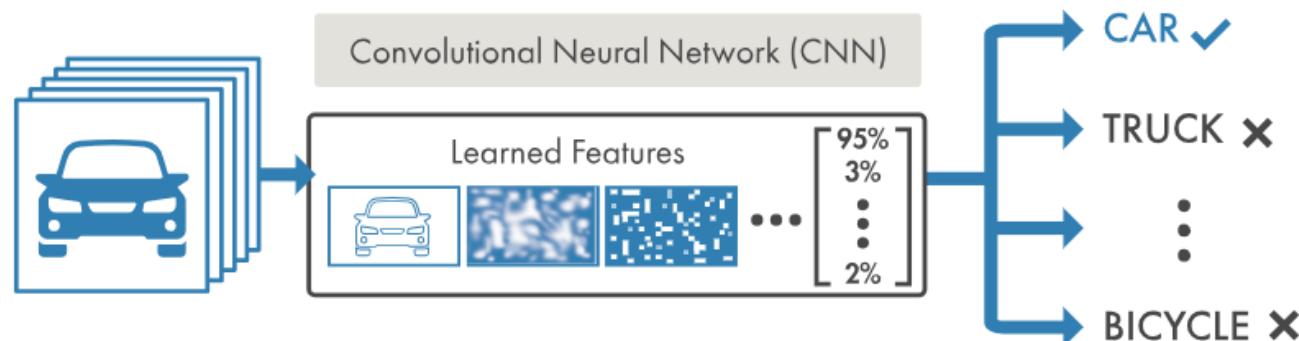
Object Classification/Detection

Classical Machine Learning V.S. Deep Learning

MACHINE LEARNING

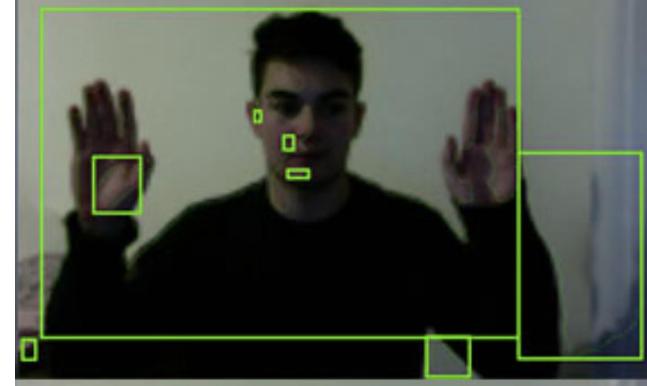


DEEP LEARNING



Hybrid solutions

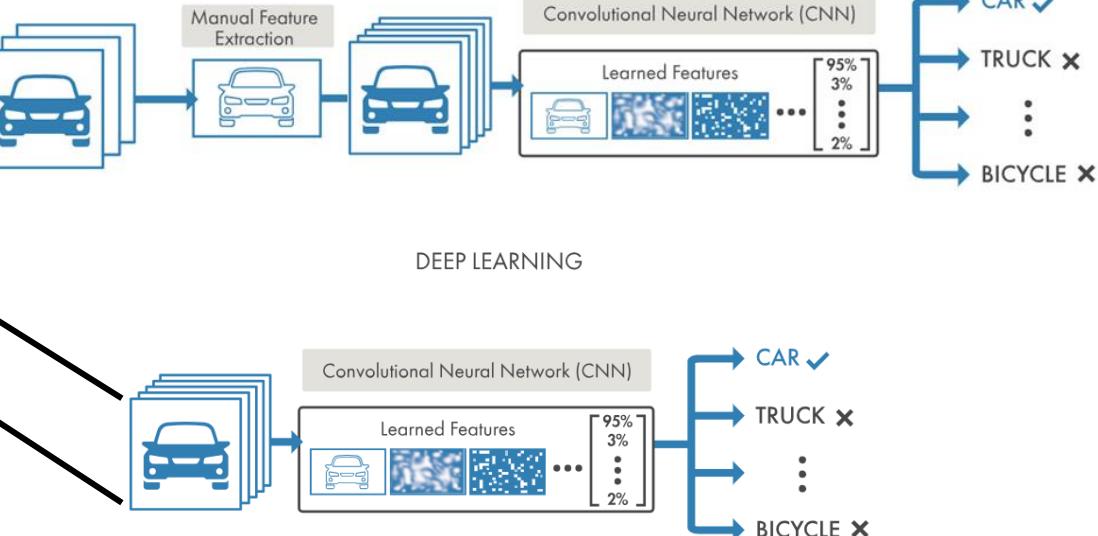
Starting from classical
blob detectors/segmentation
/block processing



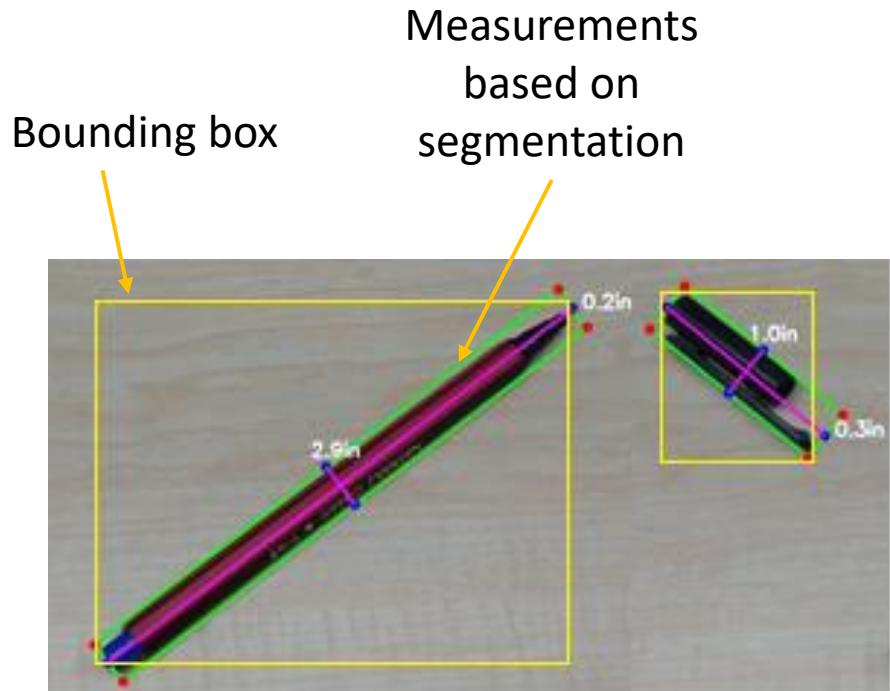
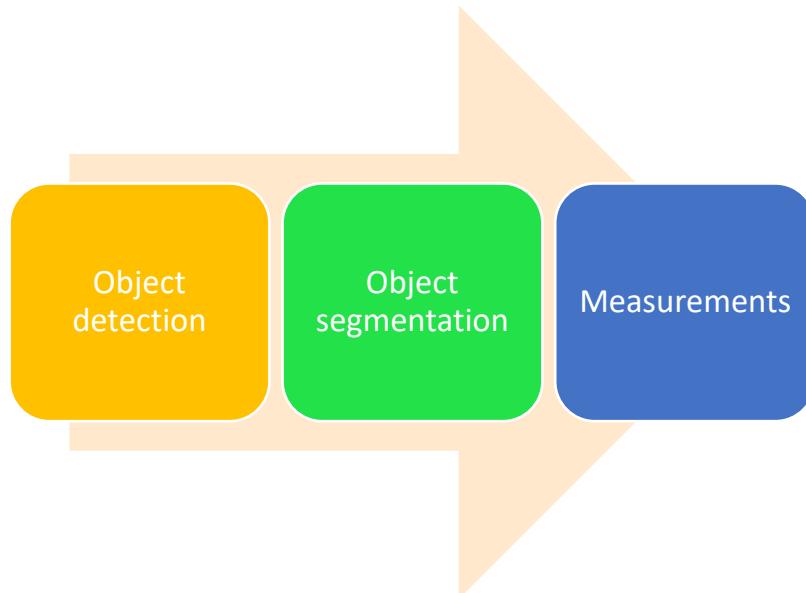
Blob detectors:
where the image responds
to specific kernel convolution
(see “similarity” lessons)



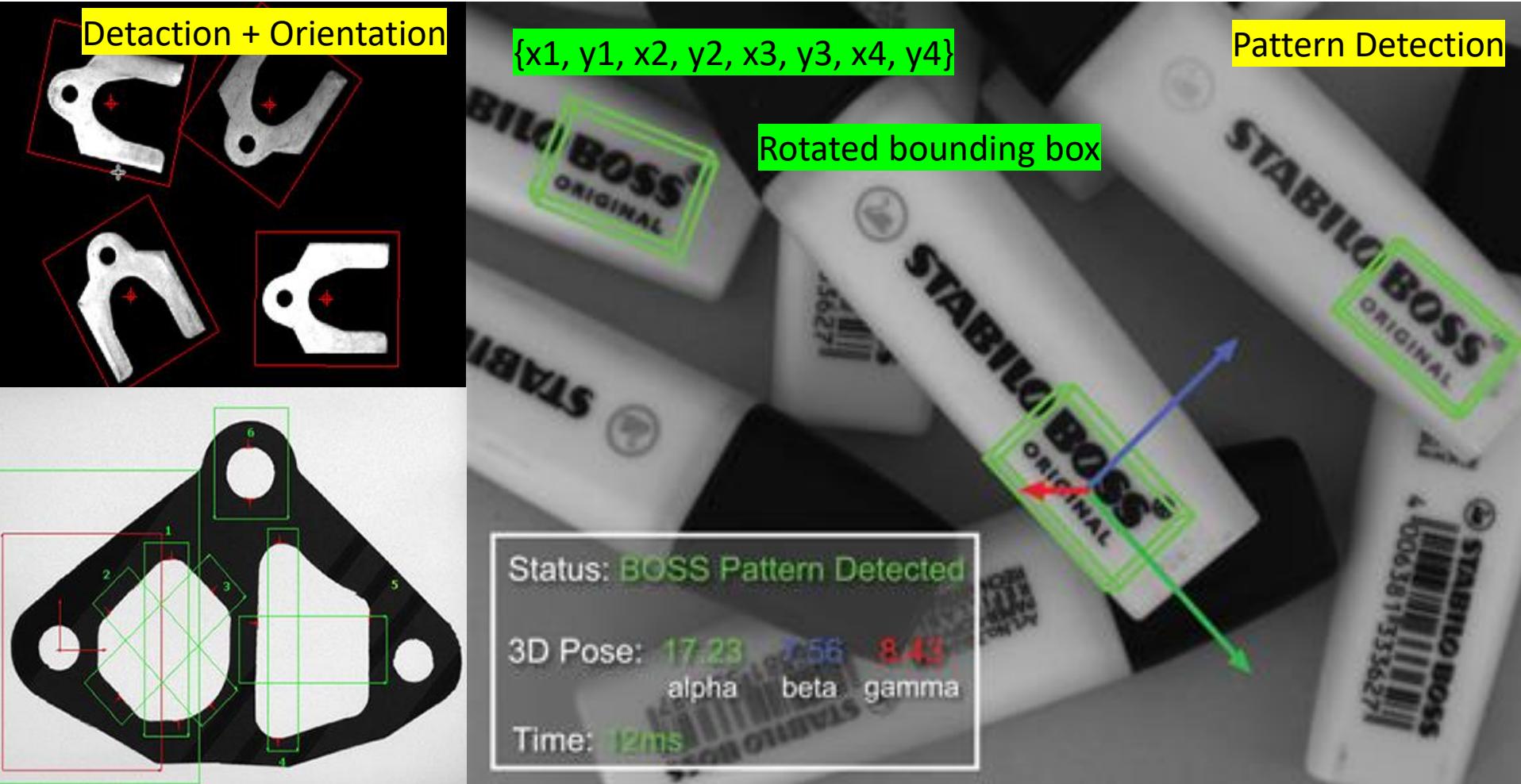
Classical feature extr.



Object Measurement



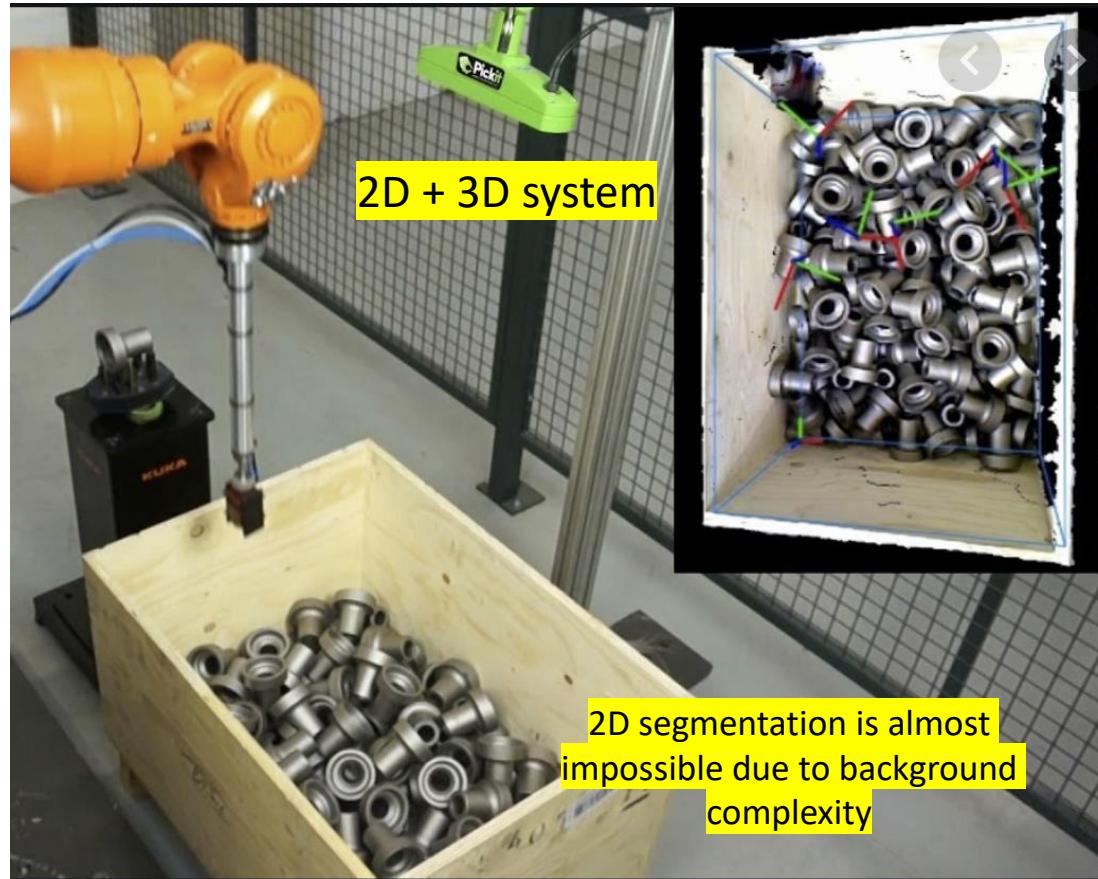
Object measurements in Industrial applications



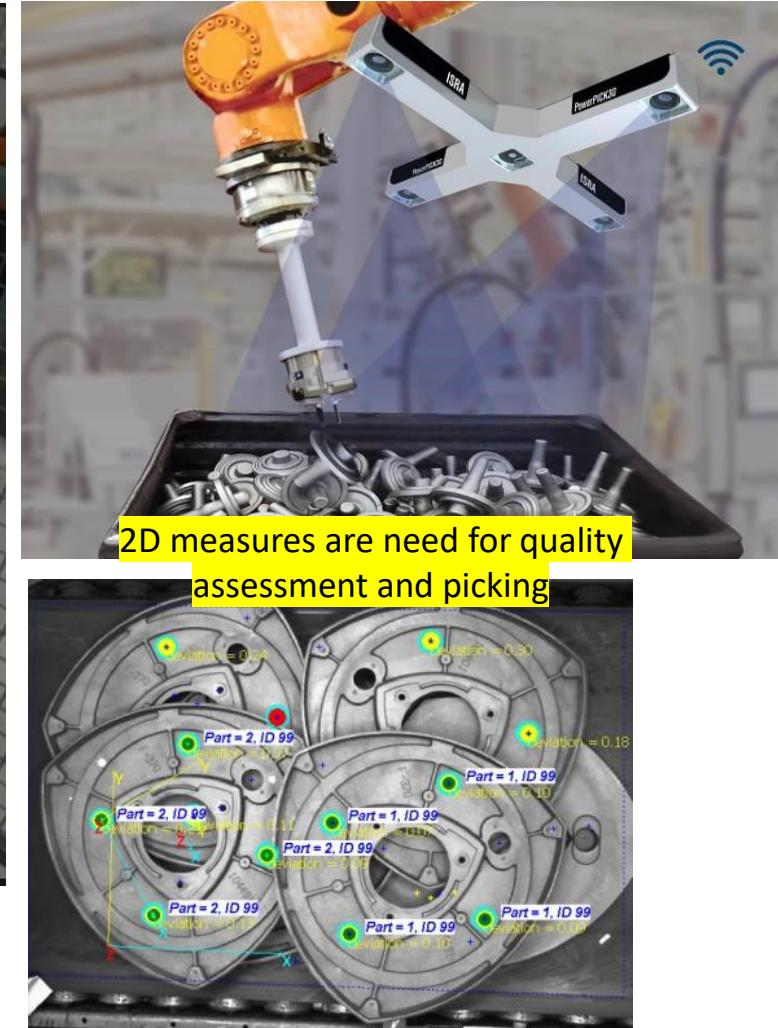
Measurements with traditional pattern matching

Fabio Scotti - Università degli Studi di Milano

Object segmentation and measurement: Bin picking



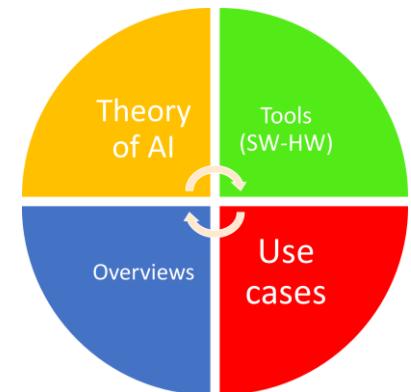
Processing the exact grabbing position with intelligent vision system



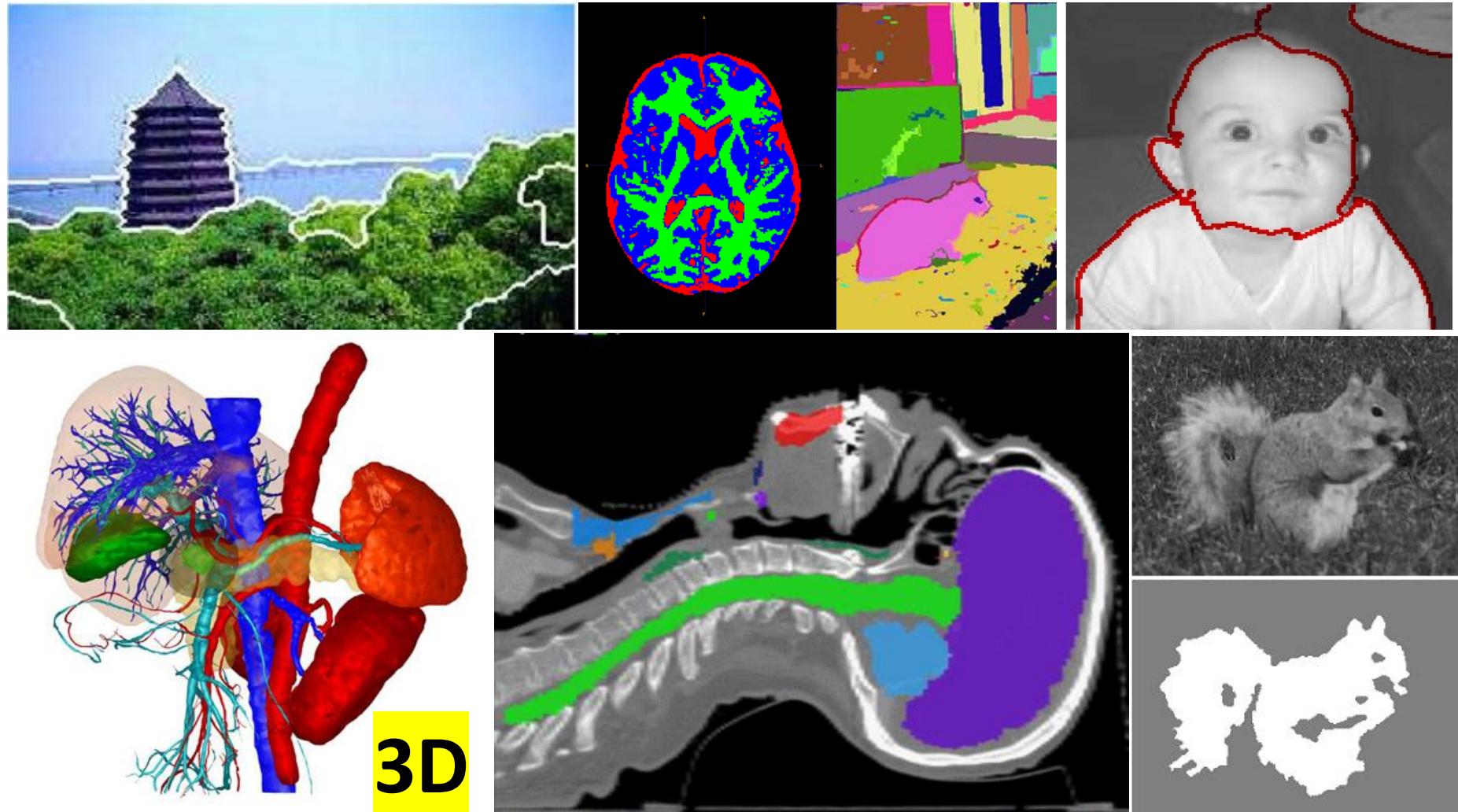


THEORY

Object segmentation



General segmentations



3D

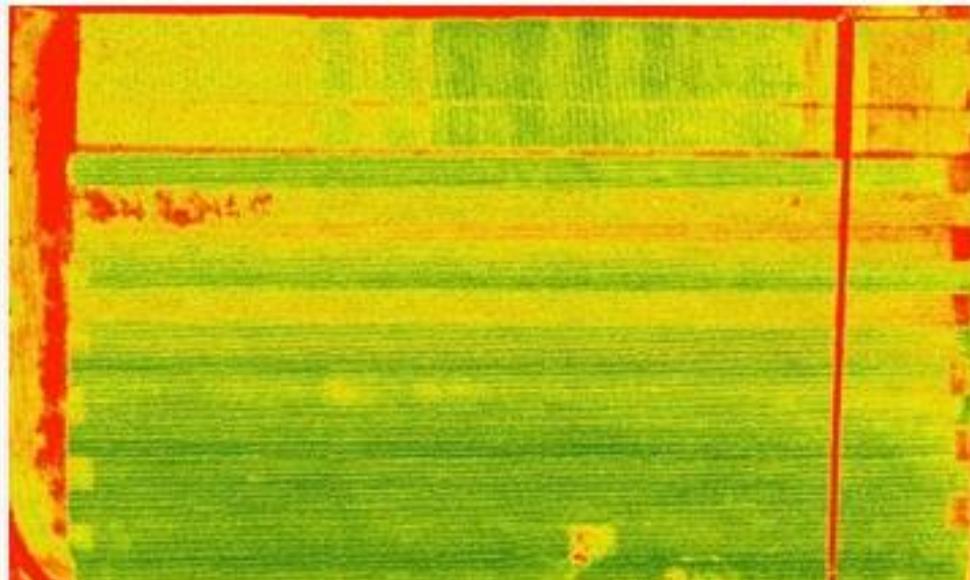
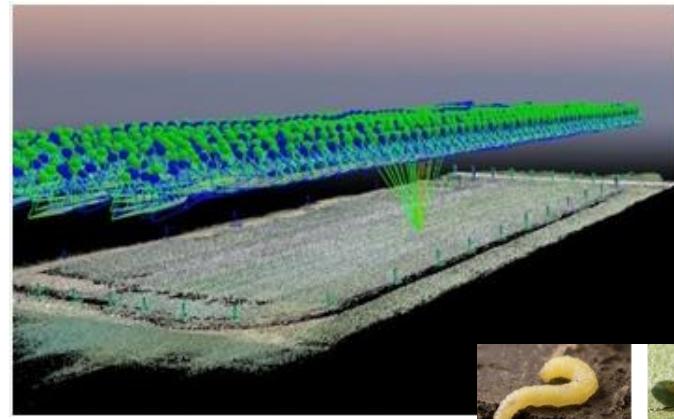
From 3D data it's possible to segment volumes of the objects

Environmental applications: Drones



frames →

Segmentation and measurements

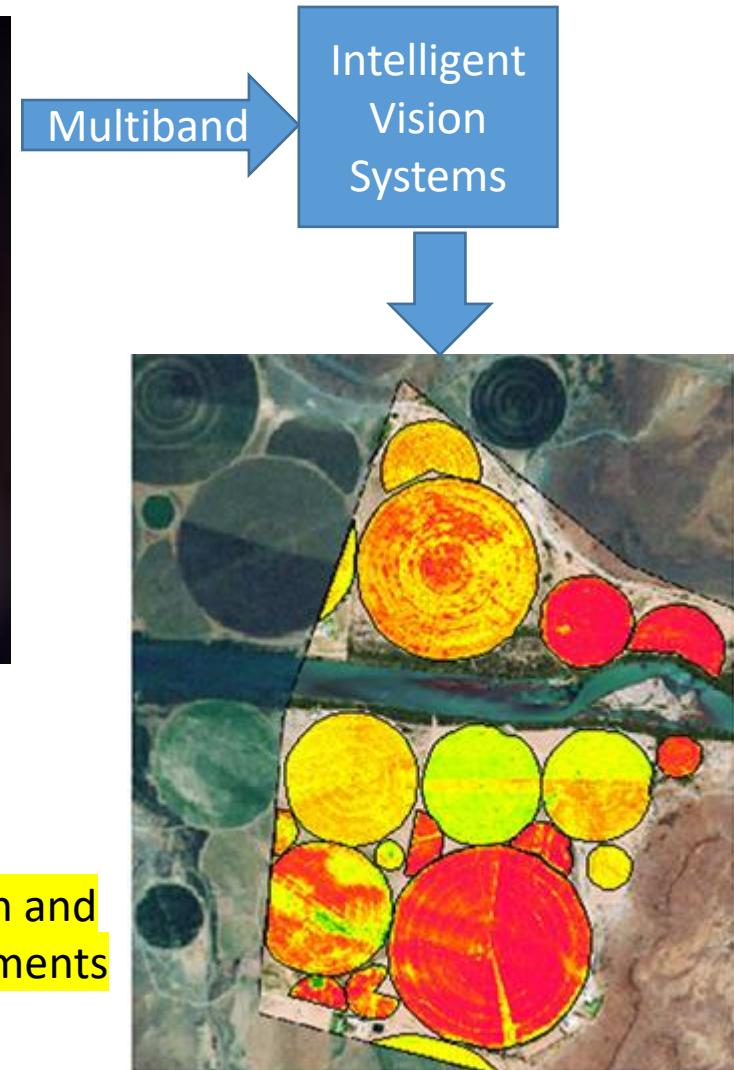


- Example of an outbreak of Corn rootworm.
- The static images and video would help the farmer or crop service vendor identify generally where the damage is occurring.
- Integration with Geographic Information Systems (GIS).

Satellites and airborne Environmental applications



Segmentation and
measurements

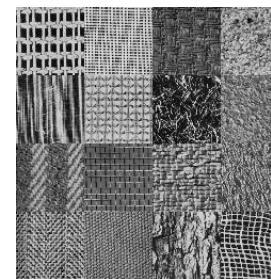
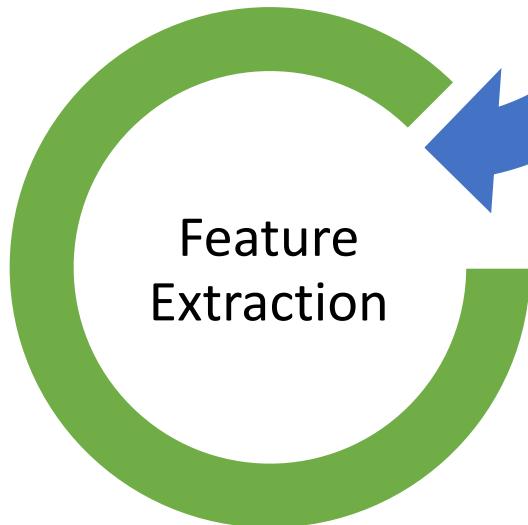
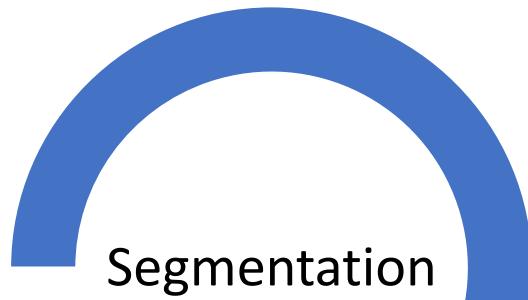
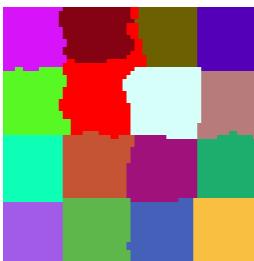


Segmentation “definition”

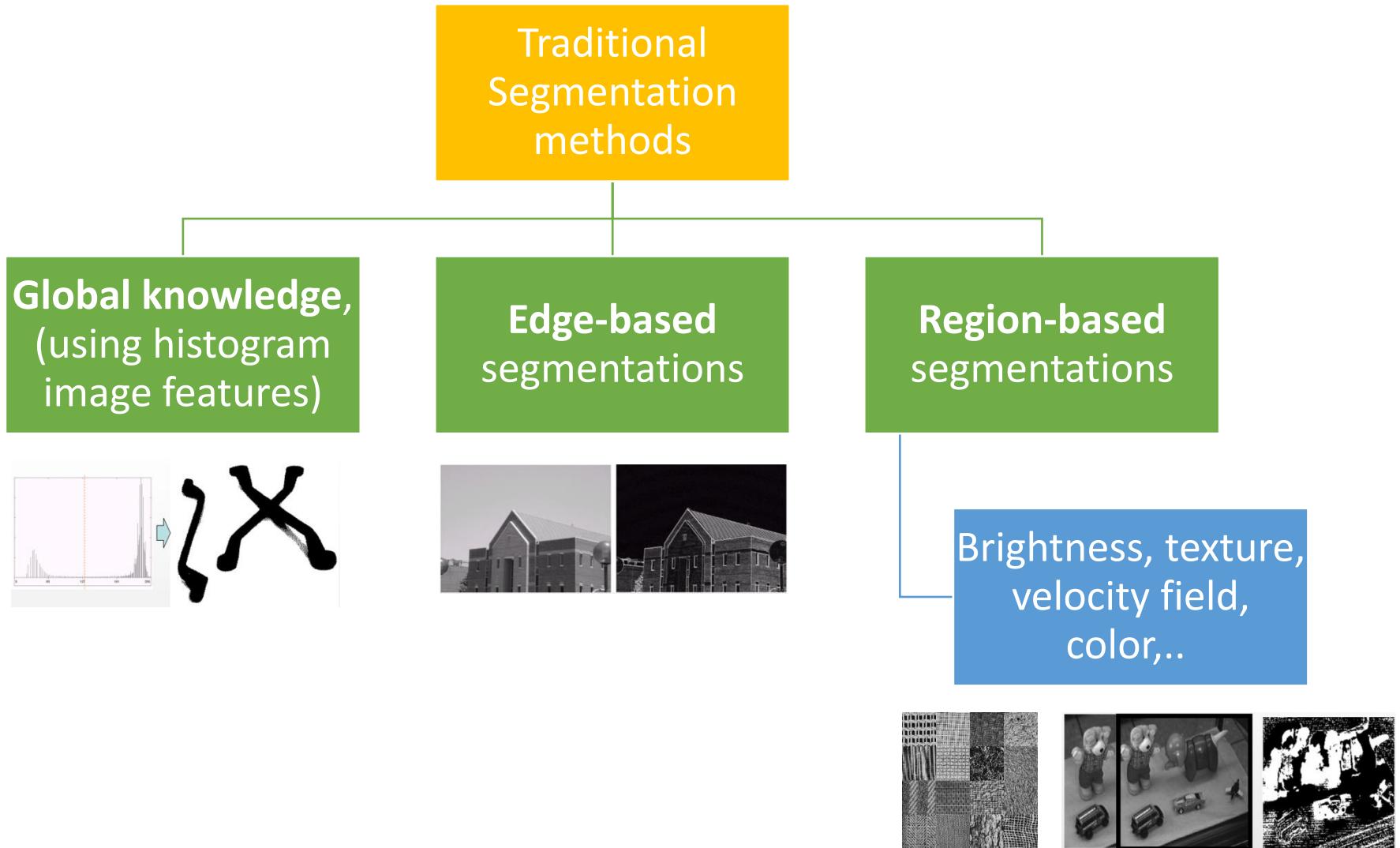
- Segmentation
 - Split or separate an image into regions
 - To facilitate recognition, understanding, and Region Of Interests (ROI) processing
- Ill-defined problem
 - The definition of a region is context-dependent



Two close tasks



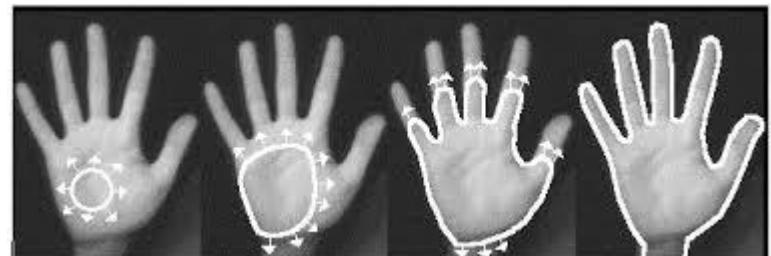
Basic taxonomy



Some classical methods

(for your reference in case of need)

- Thresholding
- Clustering methods
- Compression-based methods
- Histogram-based methods
- Edge detection
- Region-growing methods
- Split-and-merge methods
- Partial differential equation-based methods
 - Parametric
 - Level set
 - Fast Marching
- Graph partitioning methods
- Active Contour Models - “Snakes”
- Watershed transformation
- Model based segmentation
 - Deformable models (Snakes)
- Multi-scale segmentation
 - One-dimensional hierarchical signal segmentation
 - primal sketch
- Semi-automatic segmentation
- Neural networks segmentation

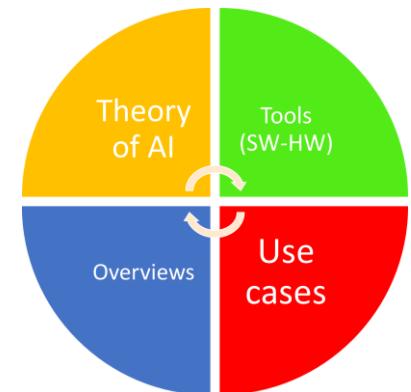


An example of a region-growing method (Snake)



THEORY

Semantic segmentation



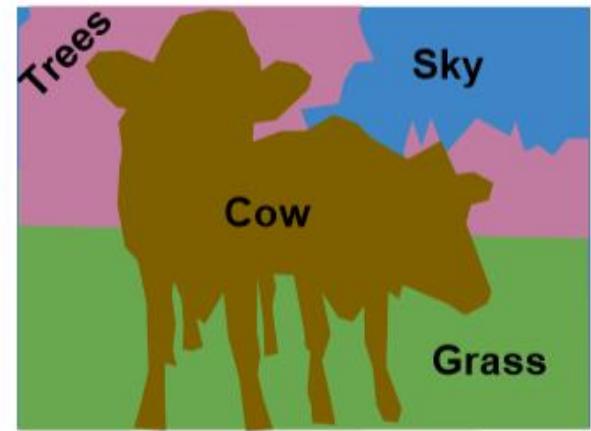
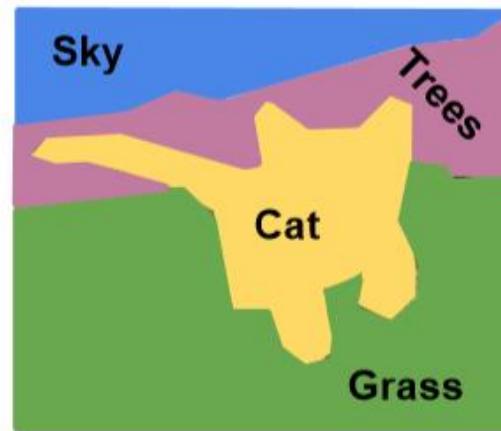
Semantic segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



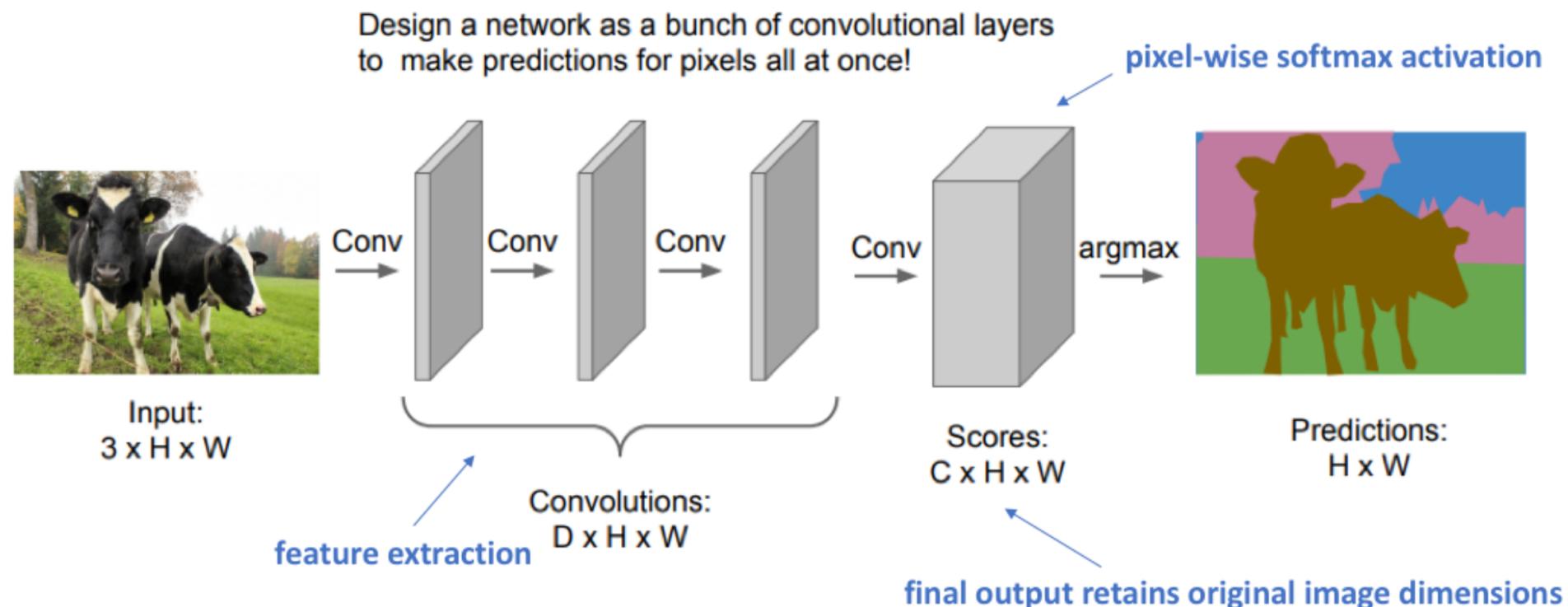
[This image is CC0 public domain](#)



Pixel classification is a good idea but... not feasible

Fully Convolutional Network

Image that every single pixel of the output image is classified with a softmax...



Semantic segmentation

A fully convolutional solution

Fully Convolutional Network with Downsampling and Upsampling

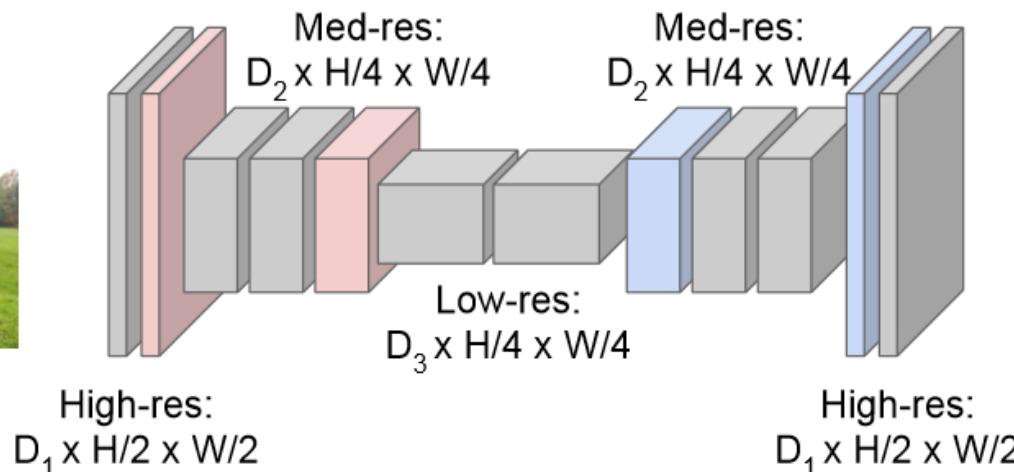
Encoder/decoder structure

Downsampling:
Pooling, strided
convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!

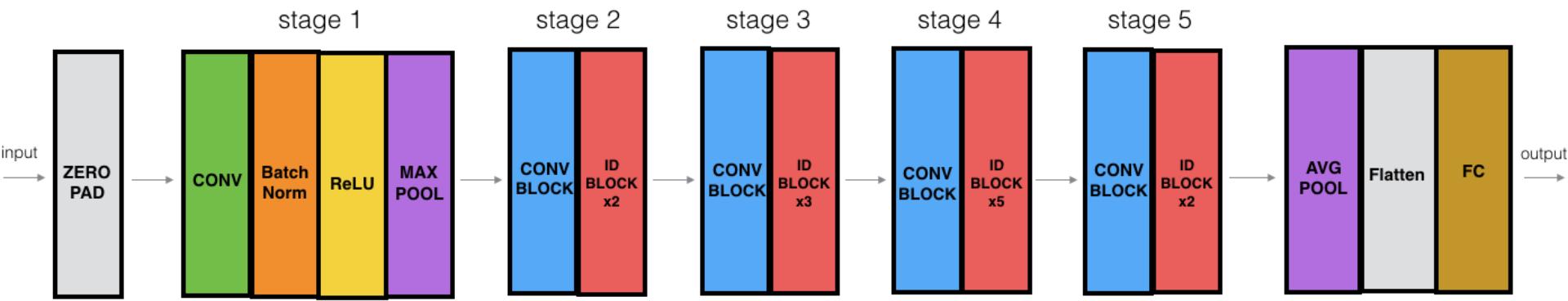


Upsampling:
Unpooling or strided
transpose convolution



Predictions:
 $H \times W$

ResNet-50 for classification

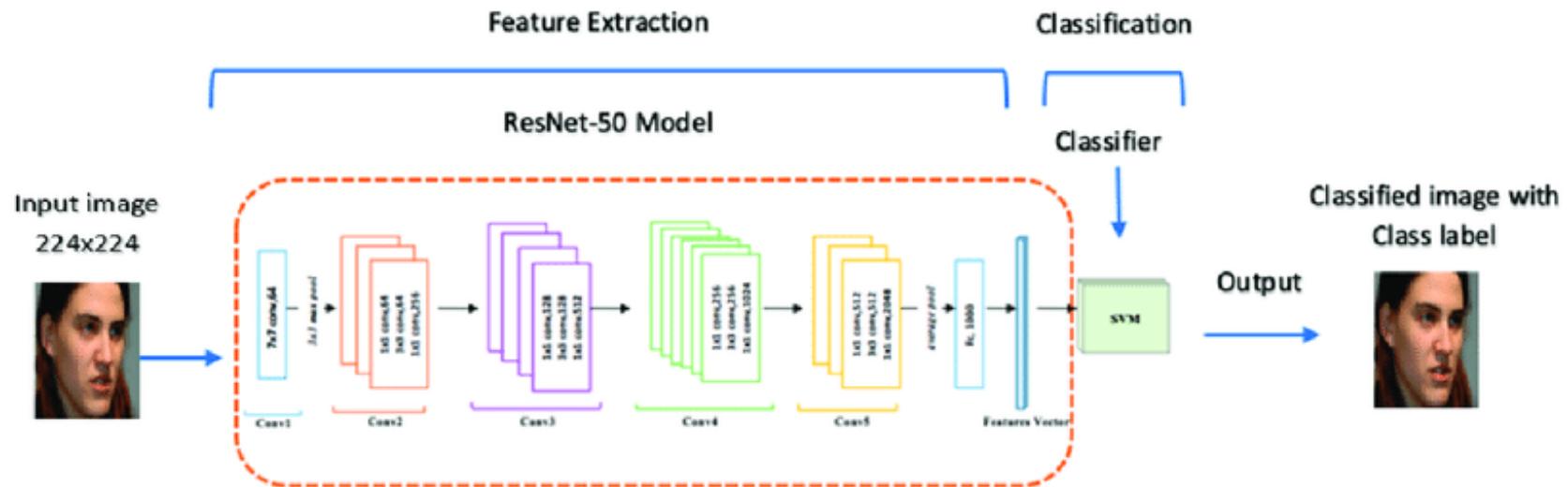


The native network is designed for classification (the output is an integer)!

KERAS: `model.param_count()`
→ 25M parameters

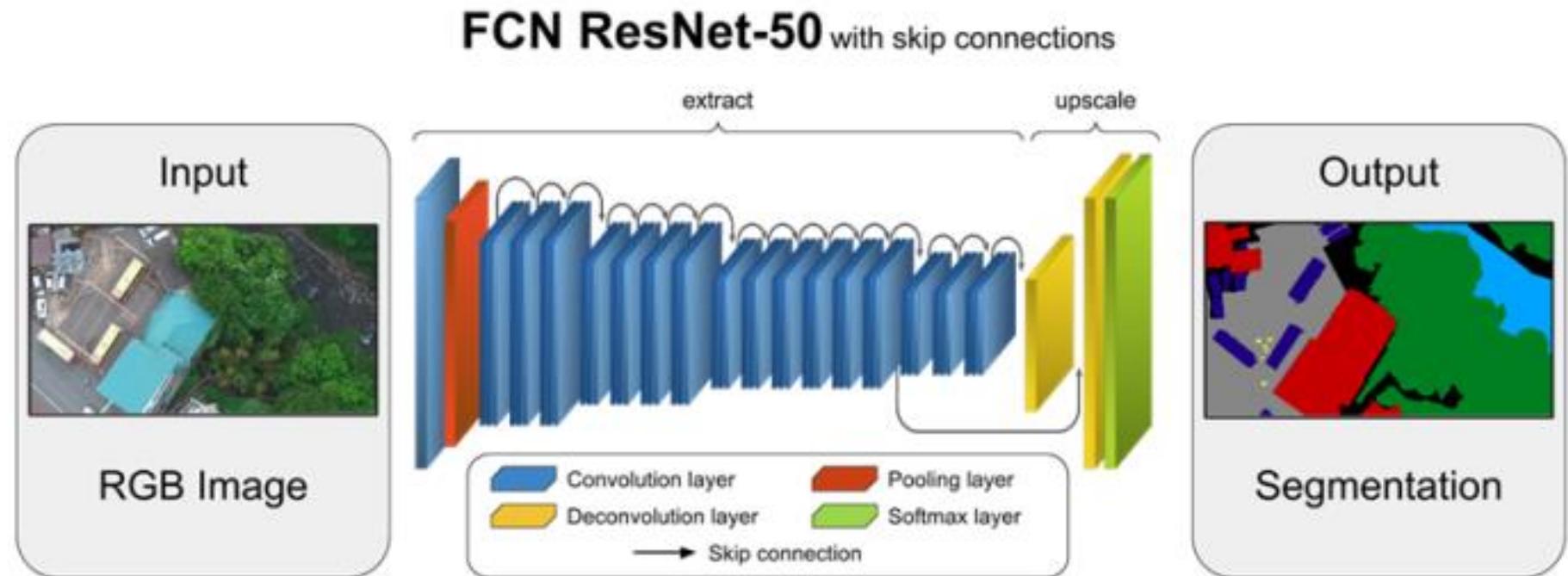
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
				3×3 max pool, stride 2		
conv2_x	56×56	$\left[\begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$	$\left[\begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$
conv3_x	28×28	$\left[\begin{array}{l} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$	$\left[\begin{array}{l} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 4$	$\left[\begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[\begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[\begin{array}{l} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 8$
conv4_x	14×14	$\left[\begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$	$\left[\begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$	$\left[\begin{array}{l} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$	$\left[\begin{array}{l} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 23$	$\left[\begin{array}{l} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 36$
conv5_x	7×7	$\left[\begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[\begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[\begin{array}{l} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

ResNet-50 for feature extraction

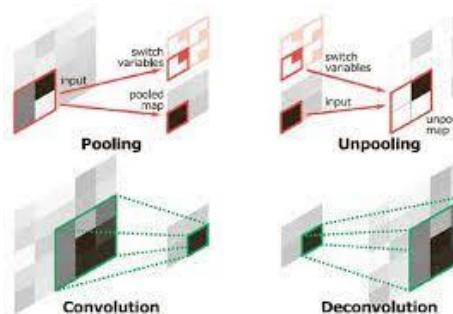


You can use ResNet-50 to process features and
then use your classifier to create a new task

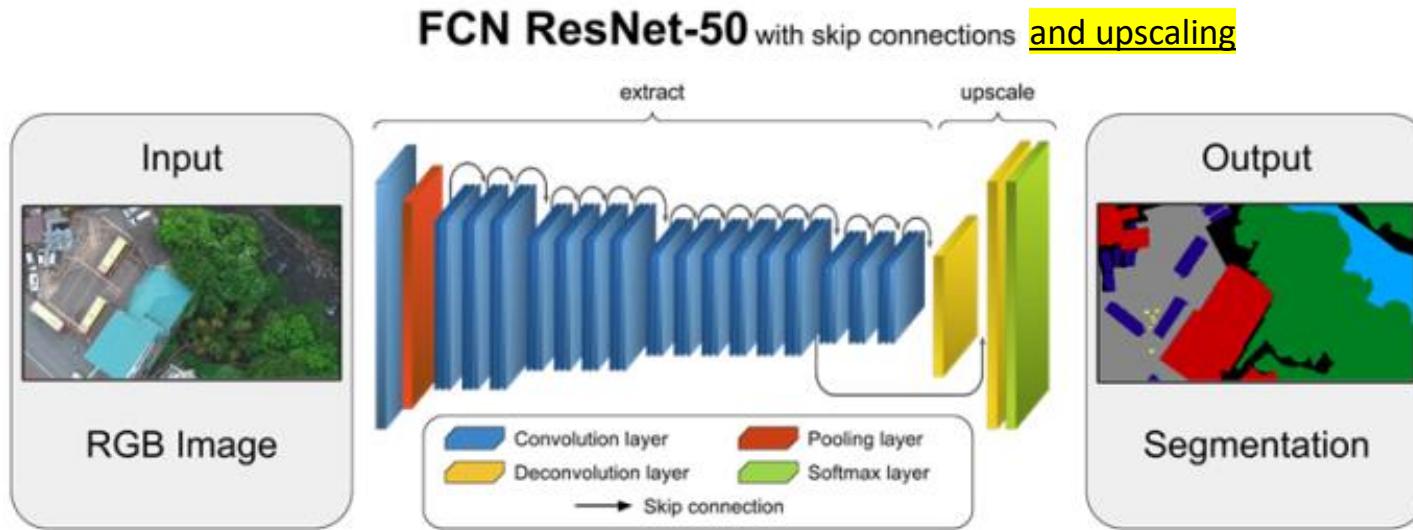
ResNet-50 → ResNet50 + Upscaling trainable semantic segmentation



Deconvolution is the operation inverse to convolution.



ResNet-50 with upscaling for classification? NO!



The native network is designed for **segmentation** (the output is an image)
not for **classification** (the output is an integer)!

Image classification via DL

Classical Pre-Trained Models for Image Classification

- VGG-16
- VGG-19
- ResNet50
- Inceptionv3
- EfficientNet B0-B7

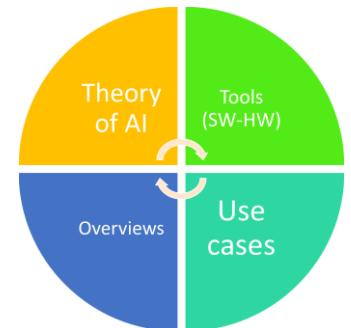
Model	Year	Number of Parameters
VGG-16	2014	138 Million
ResNet-50	2015	25 Million
Inception V3	2015	24 Million
EfficientNetB0	2019	5.3 Million
EfficientNetB7	2019	66 Million



Use cases

Objects recognition from UAVs/drones

Unmanned Aerial Vehicle (UAV)



Drones/UVAs in applications



Objects recognition from UAVs

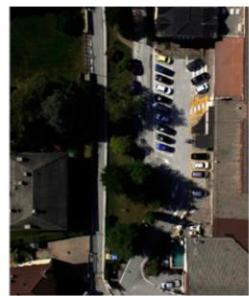


Landing mark

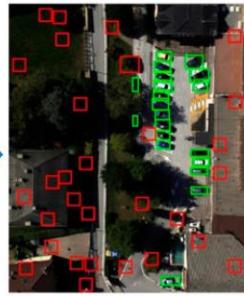
Traffic control

TRAINING

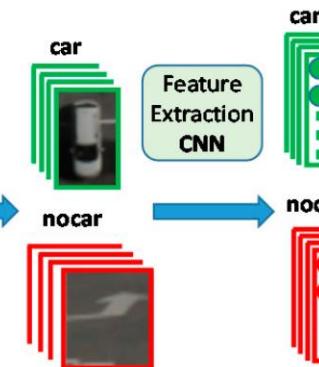
Training Image



Training Images Selection



Training Images Extraction



Classifier Training

SVM

Classifier Testing

DEPLOYMENT



Image Segmentation Mean-Shift



Remove large regions



Feature Extraction CNN



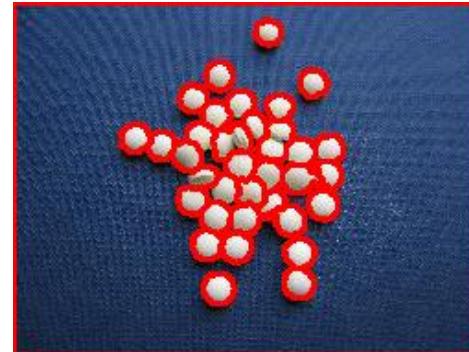
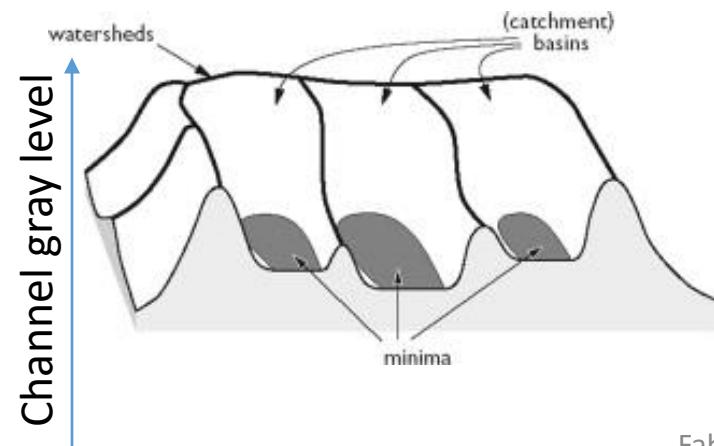
Fine-tuning



TRADITIONAL SEGMENTATION

Obj. Rec. from UAVs: STEP 1 - segmentation

TRADITIONAL SEGMENTATION FOR CANDIDATES

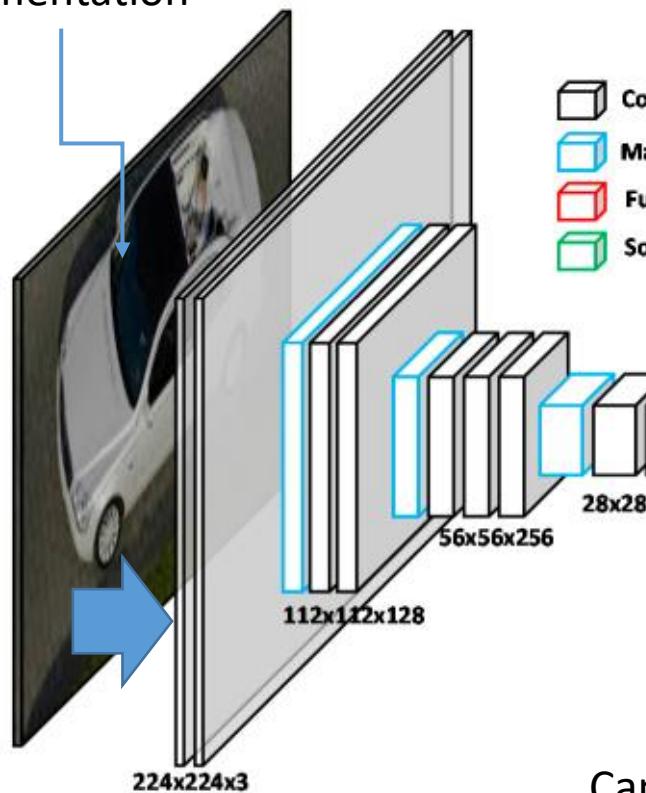


Example:
Watershed Segmentation

Obj. Rec. from UAVs: STEP 2 - CNN design

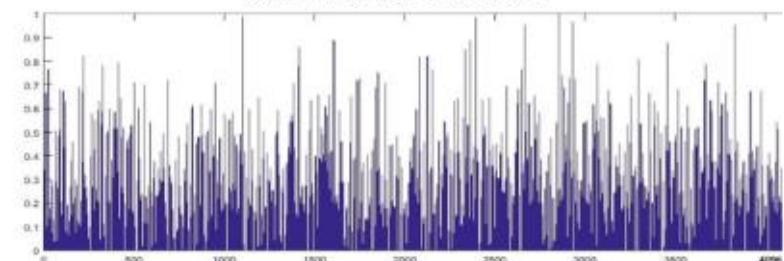
You can use tradition classifier like SVM or deep models

From
segmentation



- Convolutional+ReLU
- Max pooling
- Fully connected + ReLU
- Softmax

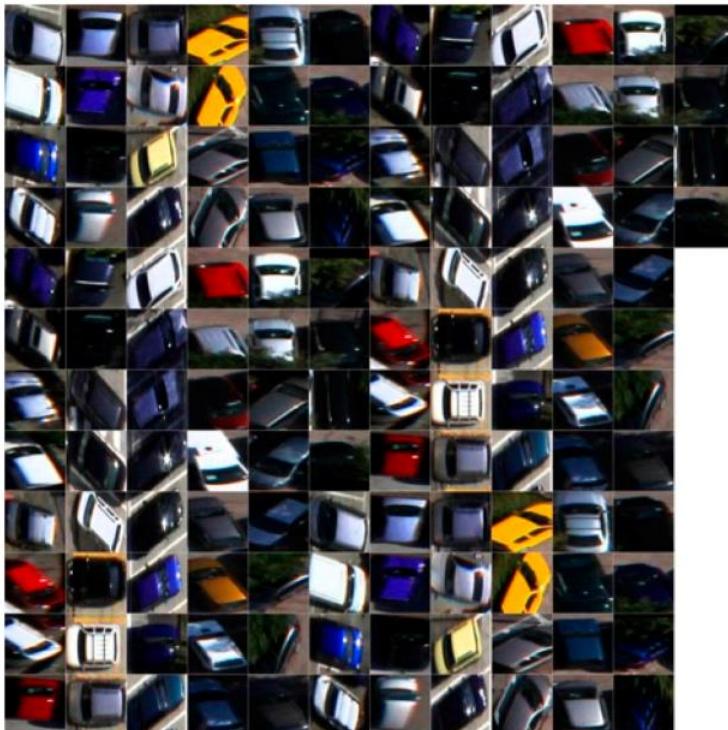
CNN feature vector



Can you reconize the net?

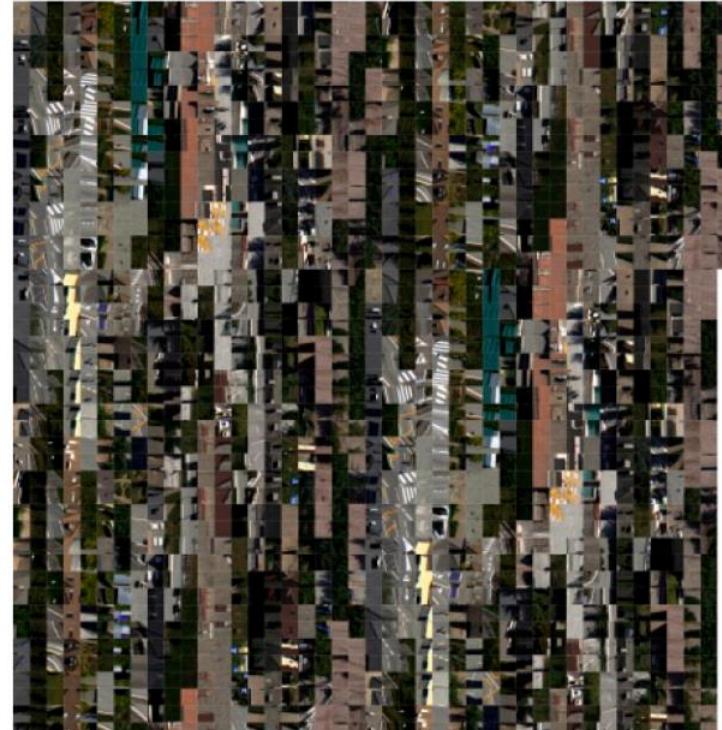
Obj. Rec. from UAVs: CNN learning dataset creation

FINE TUNING

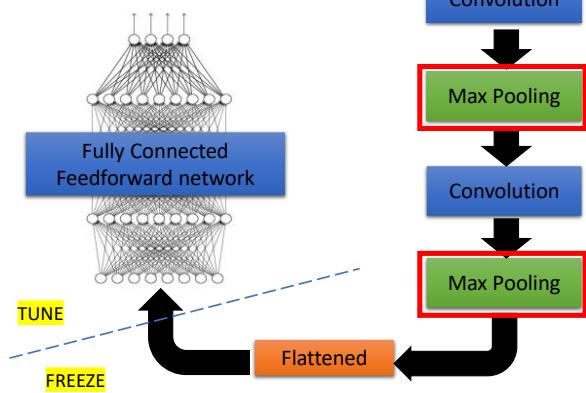


CARS

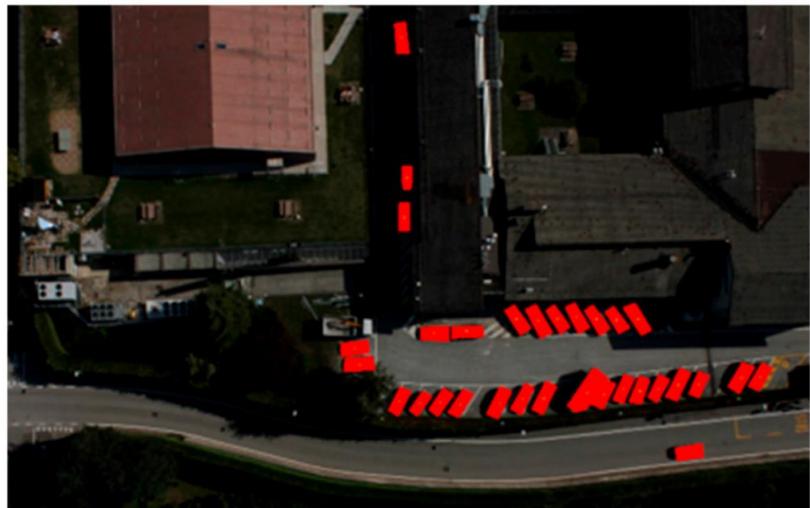
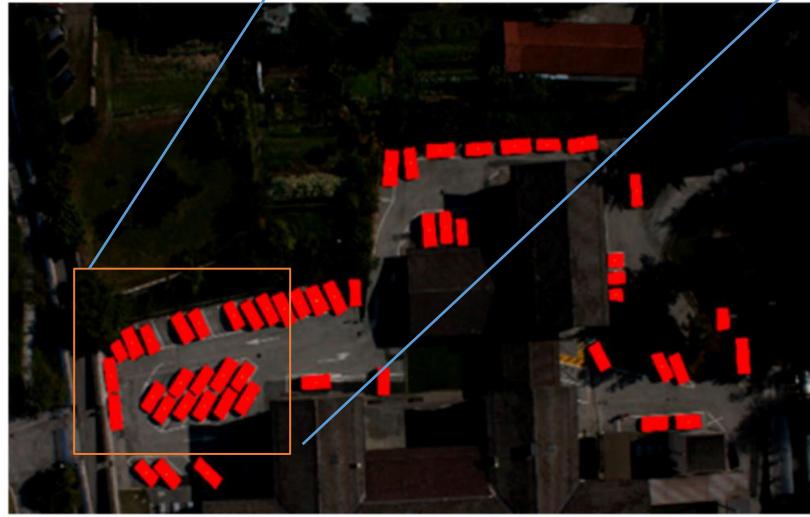
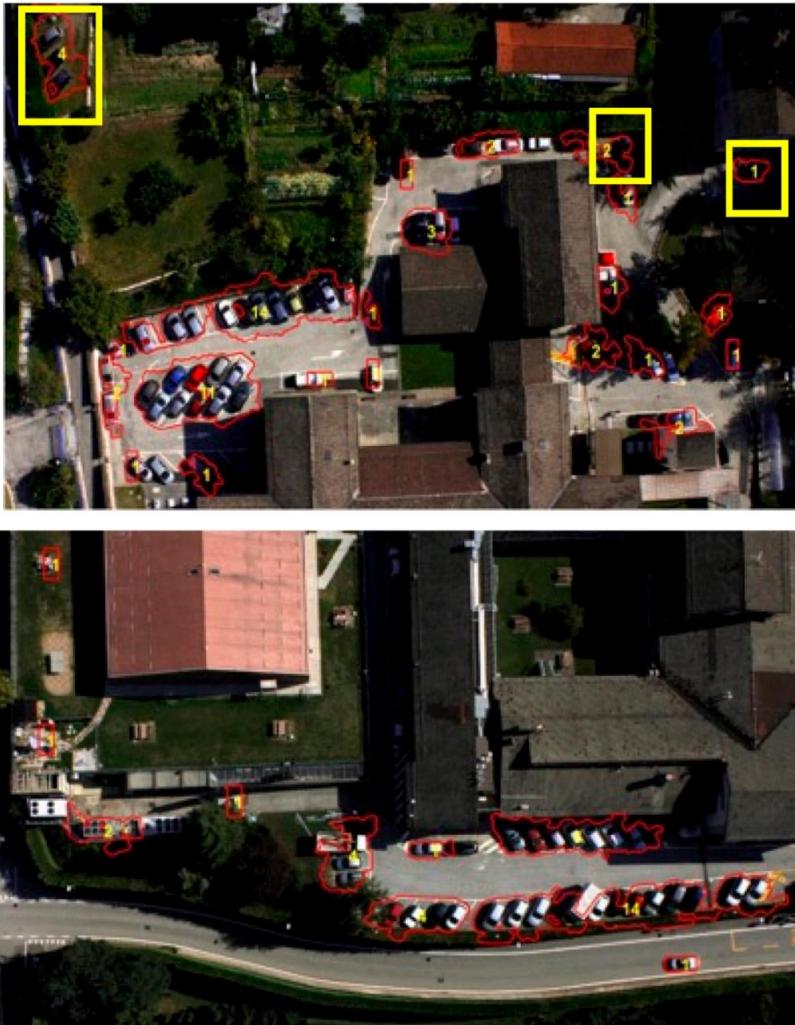
(just a portion)



NO CARS



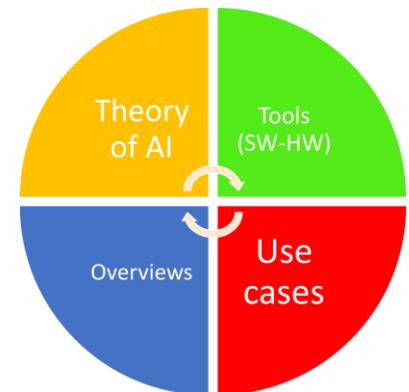
Obj. Rec. from UAVs: results





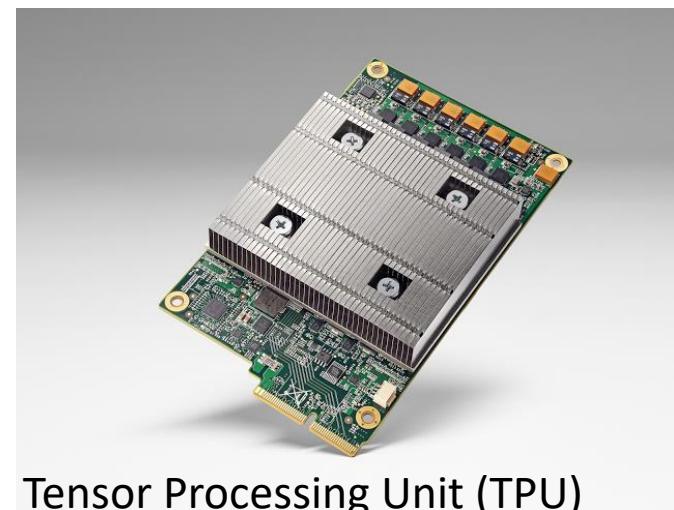
THEORY

Machine Learning HW

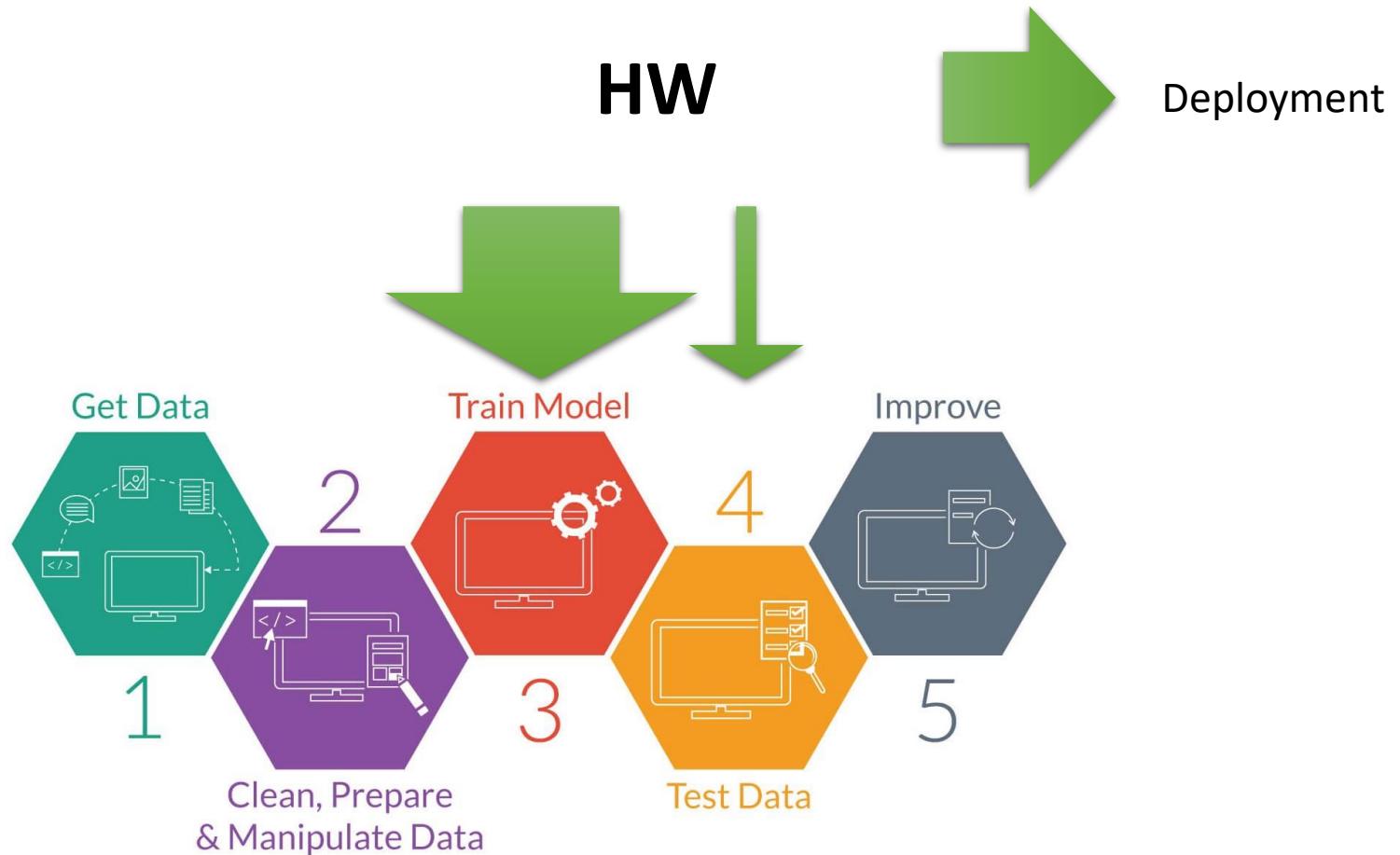


Why Machine Learning Is Possible?

- Mass Storage
 - More data available (TB HD, data centers,...)
- Higher Performance of Computer
 - Larger memory in handling the data
 - Greater computational power for calculating and even online learning

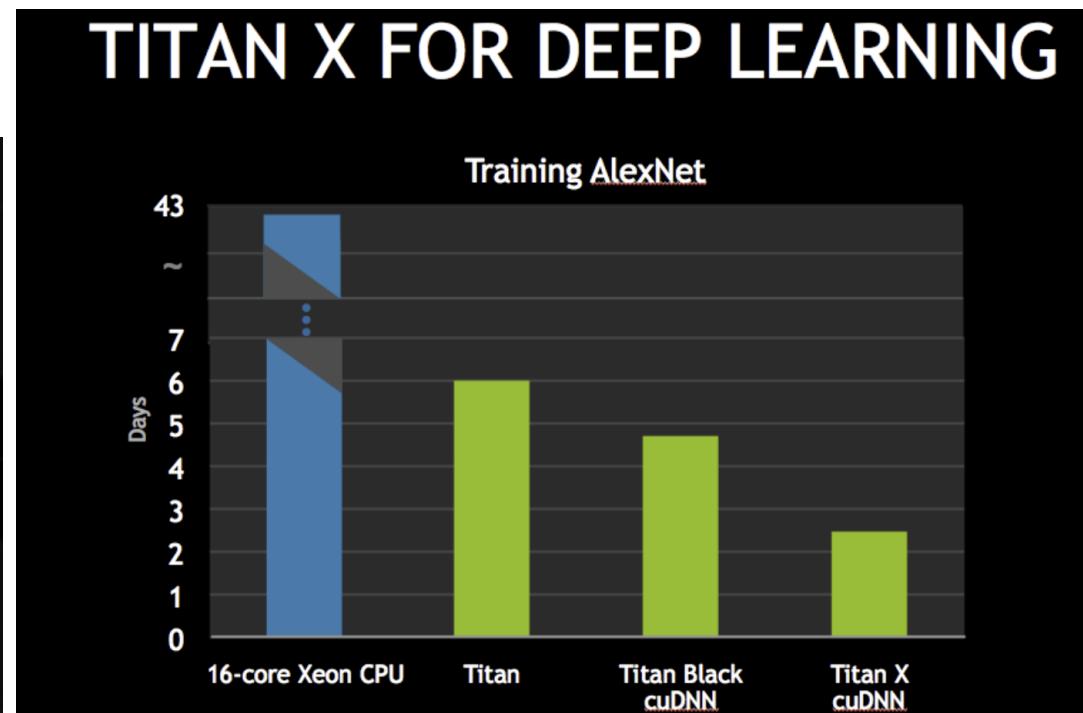
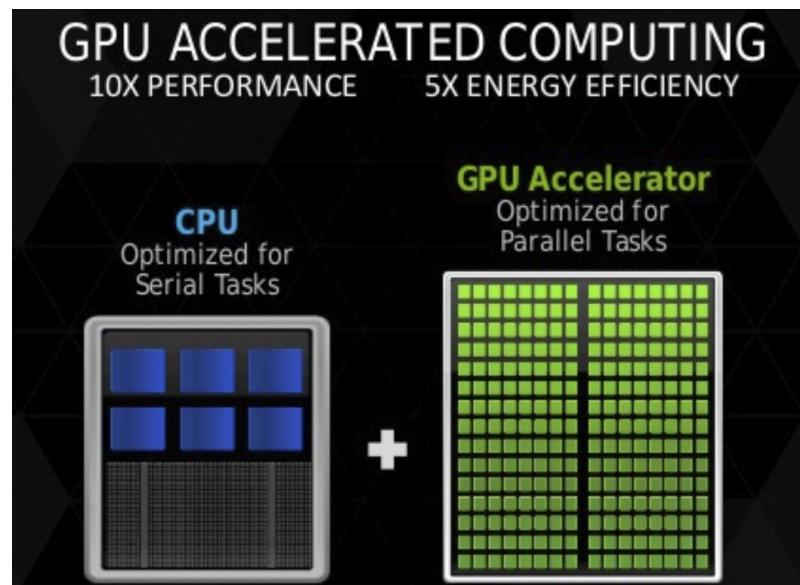


HW: When and where is more relevant?

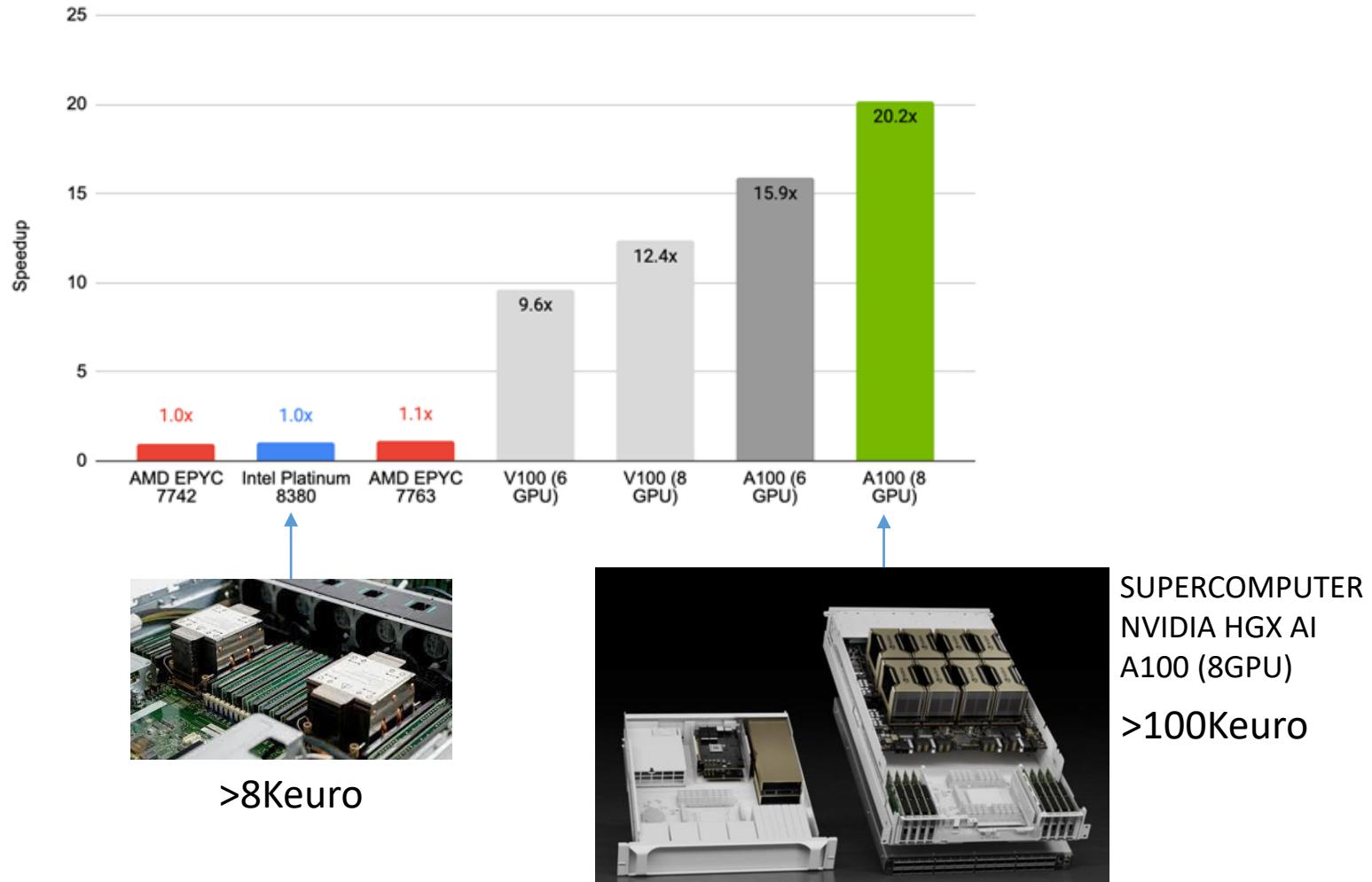


GPU vs CPU

- In brief....you need a GPU for large projects!



Some other data

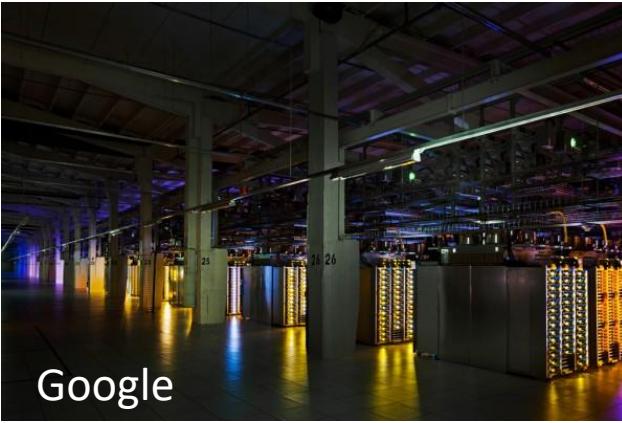


ML Server farm...

- Amazon (AWS)
- Google
- Facebook
- Microsoft (AzureML)
- ...



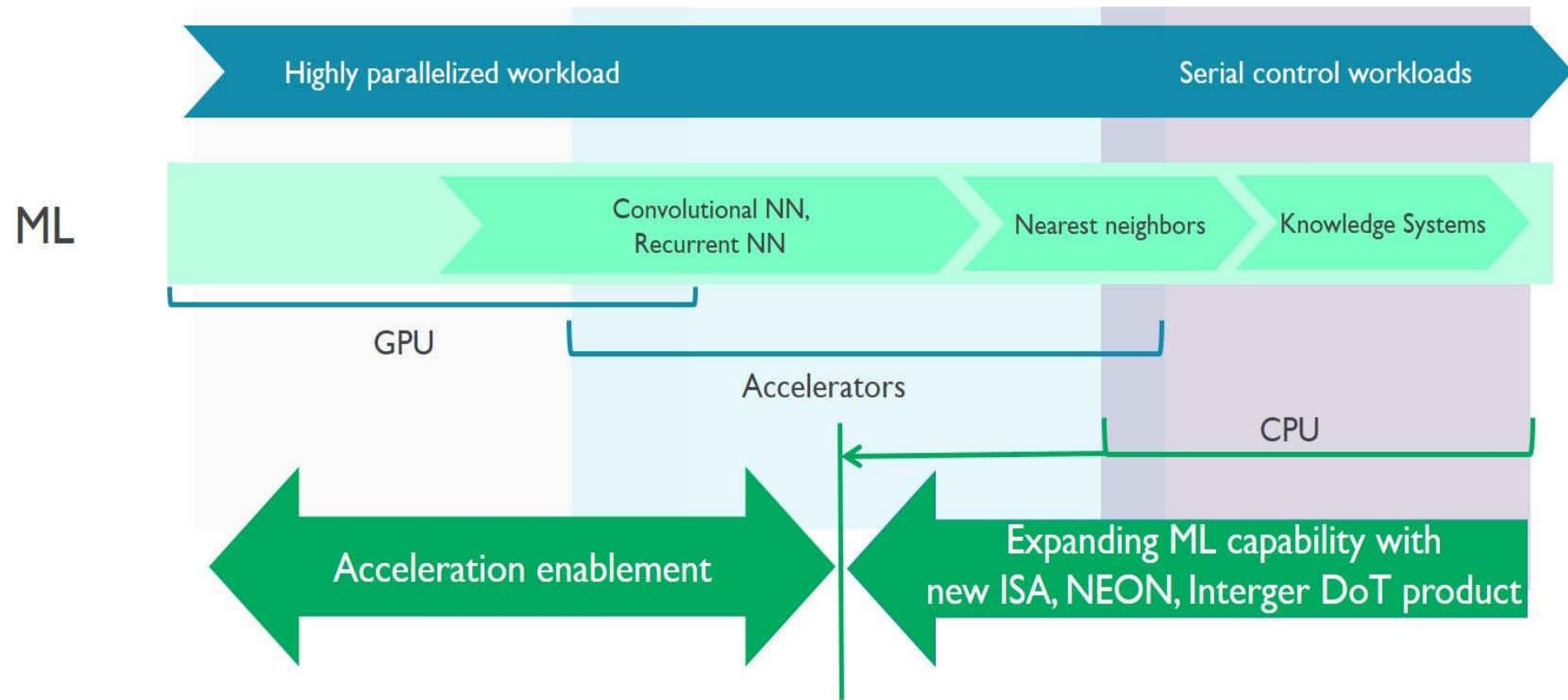
IBM



Google



GPUs: not always needed



ML@CPUs

- Machine Learning on CPUs offer **advantages**
 - Ease of **portability** and use-case flexibility
 - Same code, different applications
 - Market availability at different
 - performances
 - **prices**
 - Deployment across a **wide spectrum of devices** from edge to edge compute and cloud servers

Moving ML to the «edge»

Cloud servers



Training +
inference

Edge Compute
Regional servers



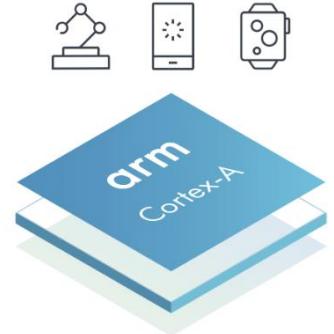
Training +
inference

Edge devices



Sensing, training,
inference & actuation

Example: ML on a ARM proc.

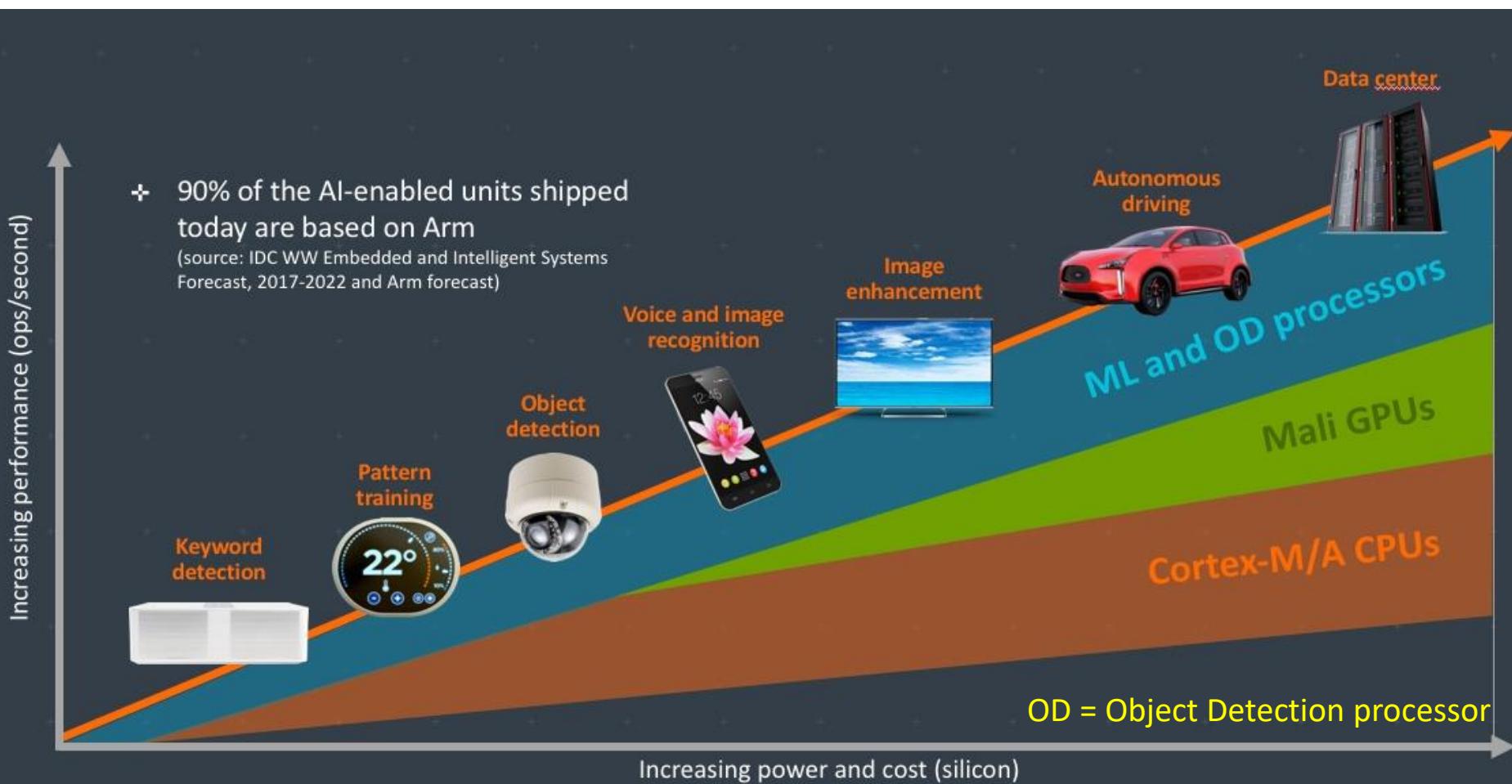


- Arm AI Platform
 - A complete, compute platform including Arm Cortex CPUs and Mali GPUs as well as a new line of highly scalable processors that make the platform versatile enough to scale to any device, from IoT to connected cars and servers.
- Set of functions for
 - ML frameworks like Google's TensorFlow
 - Imaging and vision projects
 - Providing **portable code** that can run across various Arm system configurations.
- Example:
 - the neural networks training considering examples often don't require very high accuracy data, meaning that math calculations can usually be **performed on 16-bit or even 8-bit data**, rather than large 32 or 64-bit entries.
 - The majority of neural network processing uses **8-bit fixed-point matrix multiplication**, the Armv8.2-A architecture introduced support for half-precision (FP16) and integer dot products (INT8) floating point SIMD (single instruction multiple data) NEON instructions to accelerate ML NN processing.

HW for ML

For further reference see

- Proc. Arm Cortex-A77
- GPU Arm Mali-G77
- Chip Arm Machine Learning (ML)

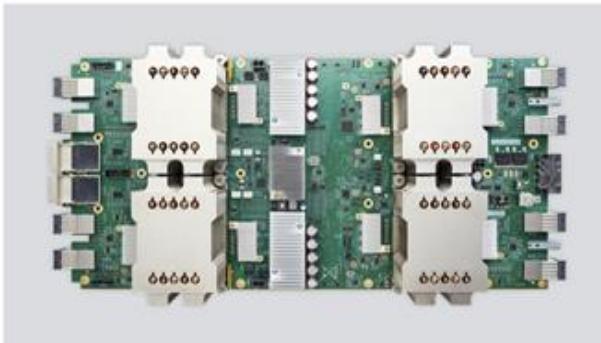


TPU

TPUs in Colab



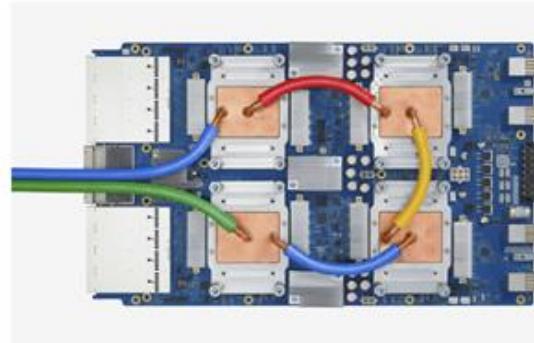
- A tensor processing unit (TPU) is an AI accelerator application-specific integrated circuit (ASIC) developed by Google specifically for neural network machine learning



Cloud TPU v2

180 teraFLOPS

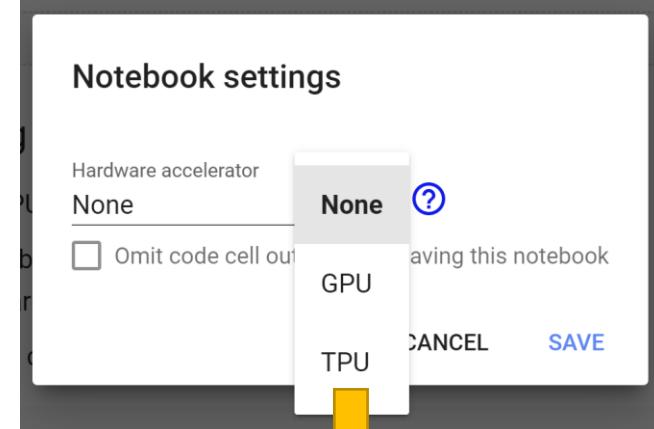
64 GB HBM (High Bandwidth Memory)



Cloud TPU v3

420 teraFLOPS

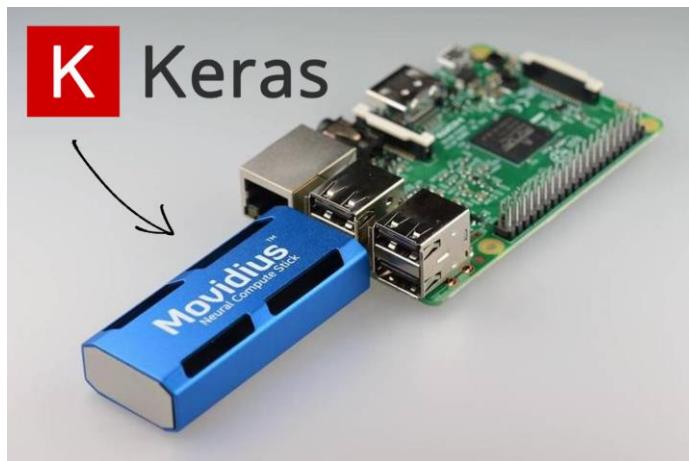
128 GB HBM



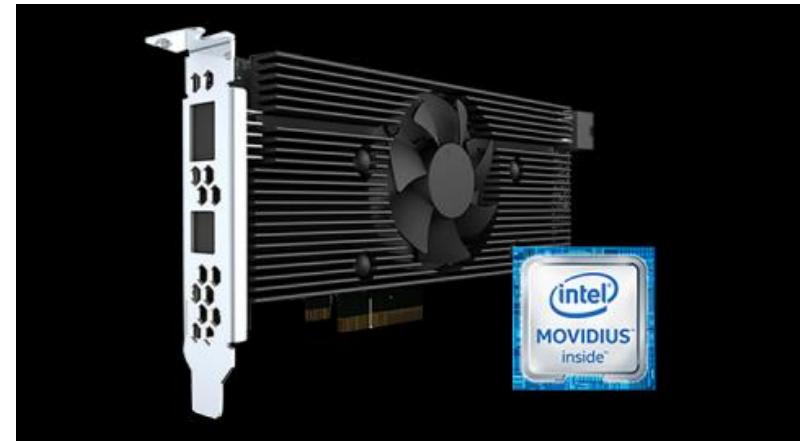
Your CoLab code is running here

Movidius Neural Compute Stick (Intel)

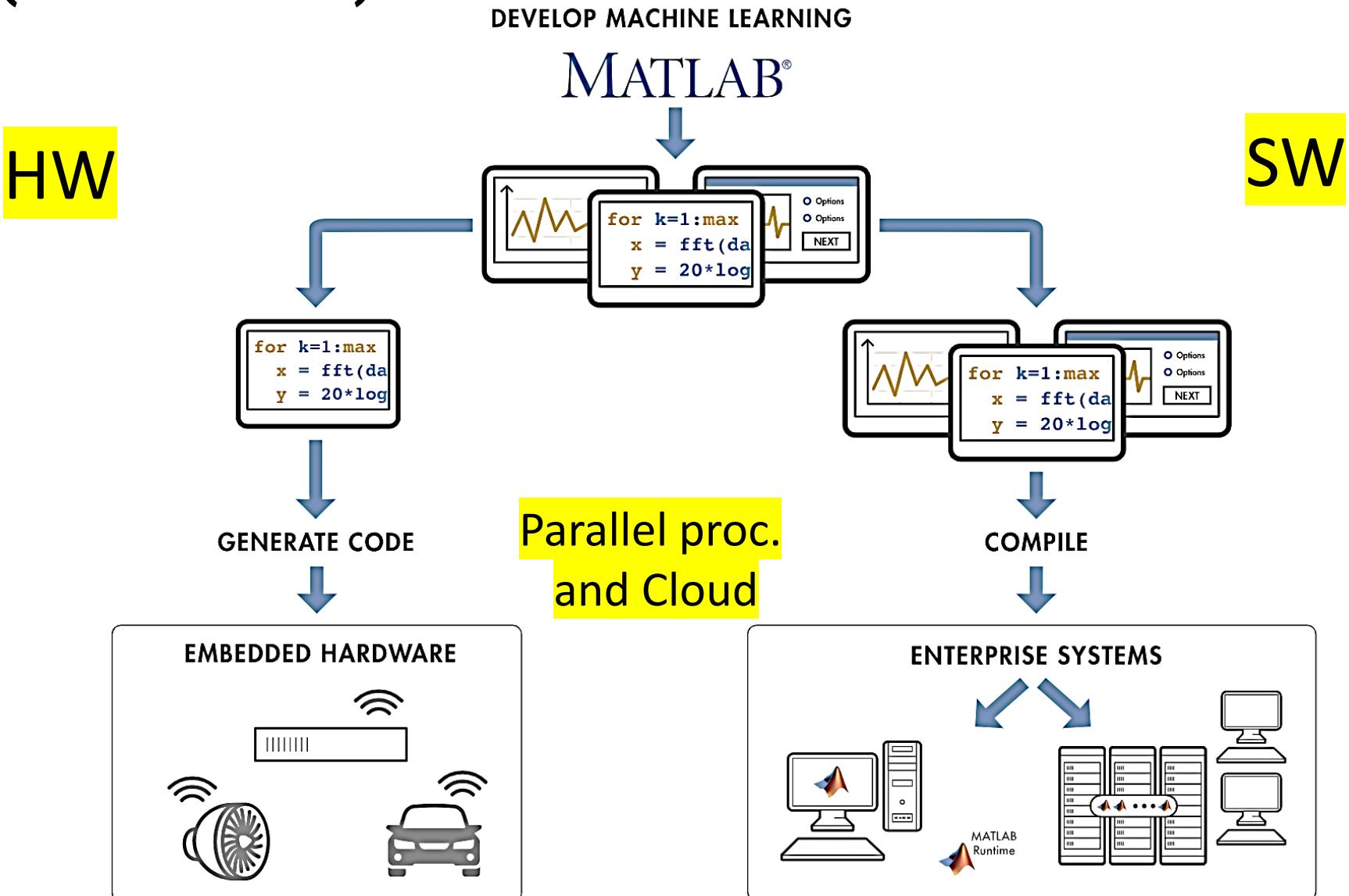
- The Intel Movidius Neural Compute Stick (NCS) is a tiny fanless deep-learning device that can be used to learn AI programming at the edge
- What can you do with a Raspberry PI and a Movidius stick?
<https://developer.movidius.com>



Intel® Movidius™
Vision Processing Unit (VPU)



The Matlab environment (HW+SW)





Main points

- Intelligent Vision Systems
 - Object detection,
Segmentation and Recognition
 - Classical, Deep and Hybrid solutions
 - Use cases
- Machine learning Hardware



Main points

Machine learning: introductions to

- Hardware
- Software and Toolboxes
- Cloud (external resources) – SaaS,
 - Machine Learning as a Service (MLaaS)



Google Cloud



Azure Machine Learning

