

Section2

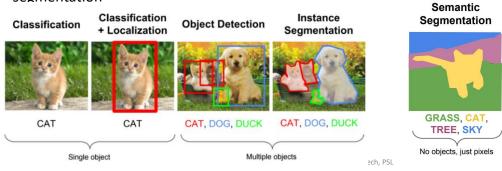
Semantic Segmentation, object detection, instance segmentation

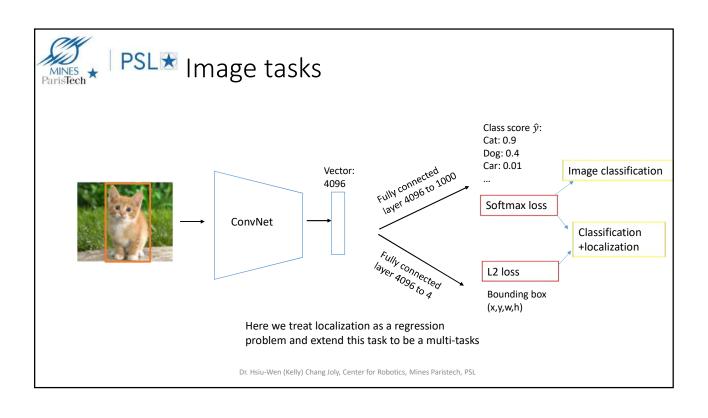
Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL

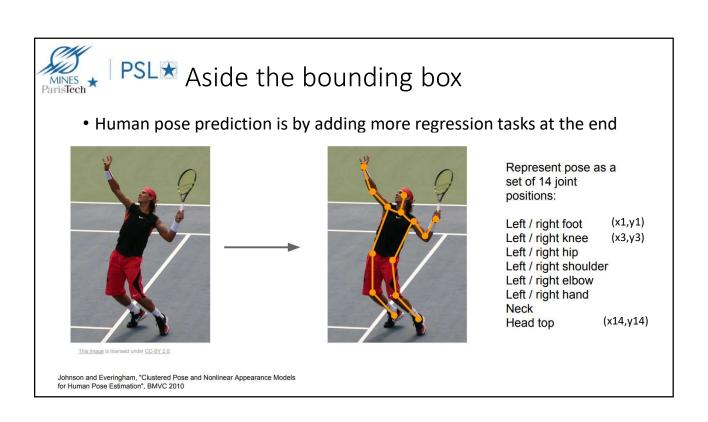


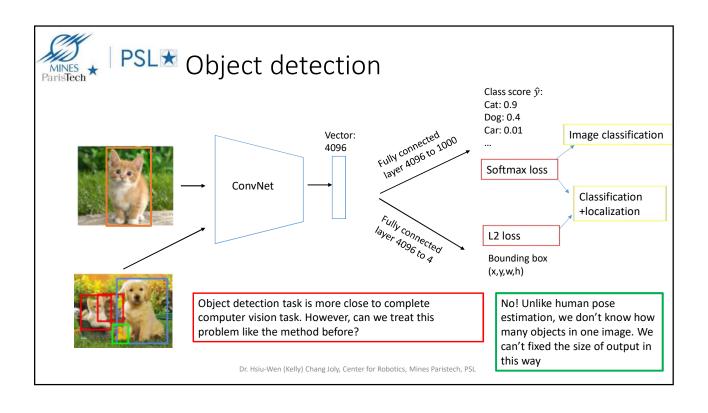
PSL Computer vision tasks

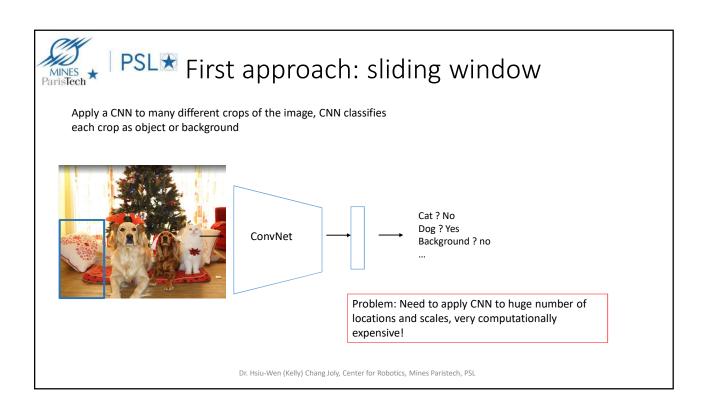
- Computer vision tasks are the main engines to drive the development of deep learning. Start with the image classification to more complete task like semantic segmentation and tracking task. We will see what we have solved in the past and what are the challenges that need to be solved.
- Image classification → Object detection → Instance segmentation → Semantic segmentation











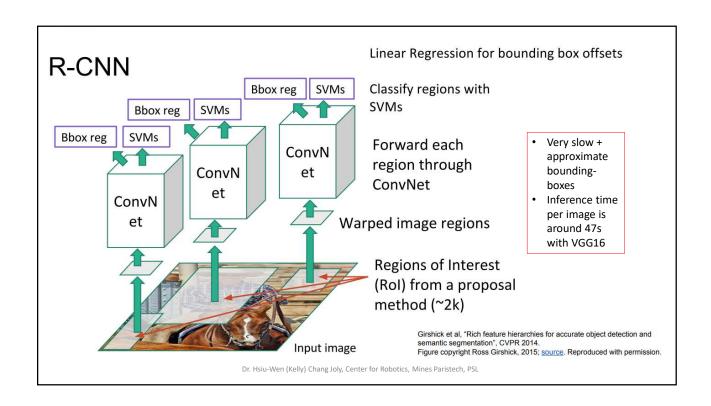


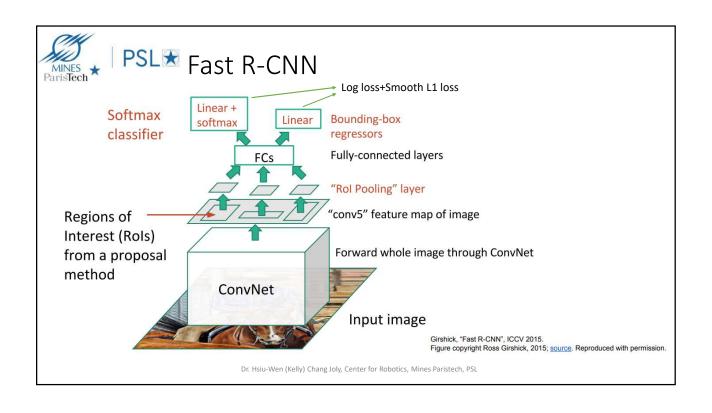
PSL First approach: region proposals

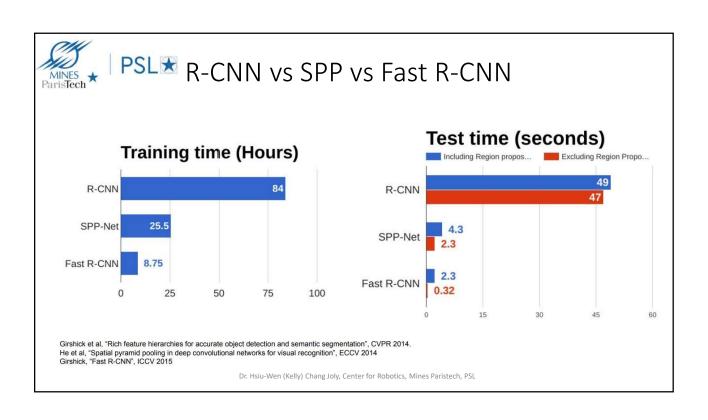
- Find "blobby" image regions that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 1000 region proposals in a few seconds on CPU

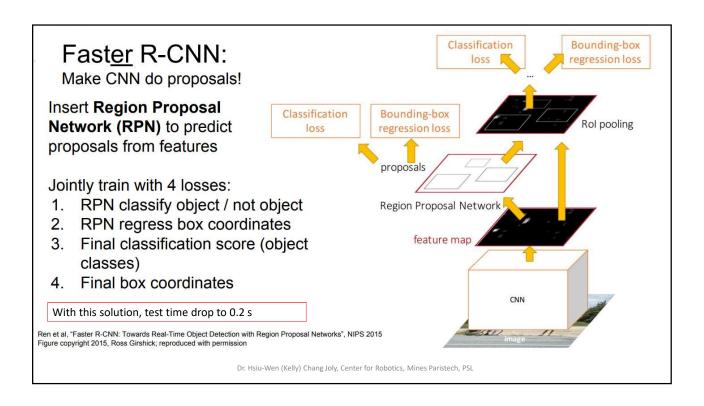








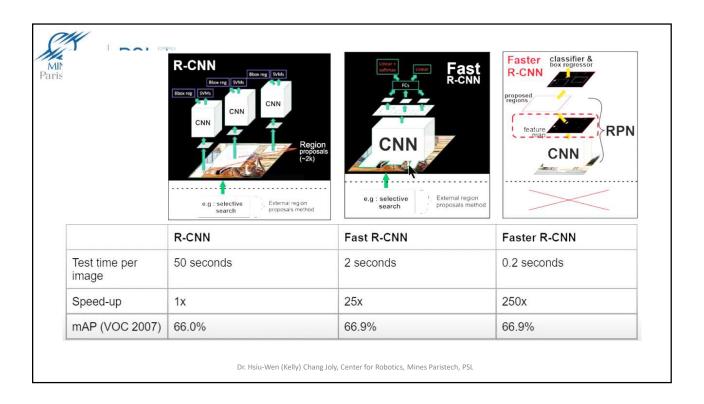






- History
 - R-CNN: Selective search → Cropped Image → CNN
 - Fast R-CNN: Selective search → Crop feature map of CNN
 - Faster R-CNN: CNN → Region-Proposal Network → Crop feature map of CNN
- Best performances, but longest run-time
- End-to-end, multi-task loss

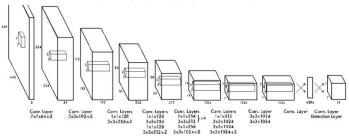
[https://github.com/endernewton/tf-faster-rcnn]





PSLM Detection without proposals: YOLO

You Only Look Once



- Modified GoogleNet/Inception
- Super fast (21~155 fps)
- Finds objects in image grids in parallel
- Only slightly worse performance than Faster R-CNN

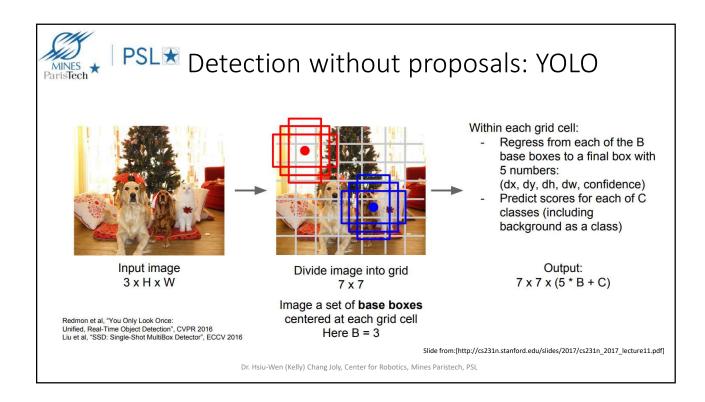


PSLM YOLO limitations and v2 improvements

- Non-maximal suppression (NMS):
 - Since YOLO uses 7x7 grid then if an object occupies more than one grid this object may be detected in more than one grid
 - For each class (cars, pedestrians, cats,....) do:
 - 1-Discard all boxes with confidence C<C -threshold (for example C<0.5)

• 5-Start again from step (3) until all remaining predictions are checked.

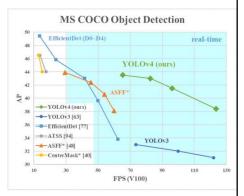
- 2- Sort the predictions starting form the highest confidence C.
- 3-Choose the box with the highest ${\bf C}$ and output it as a prediction .
- 4-Discard any box with IOU>IOU-threshold with the box in the previous step .
- Limitations
 - It can only detect maximum 49 objects
 - · Groups of small objects
 - · Unusual aspect ratios
 - · High error of localization
- YOLOv2 (2016):
 - add Batch Normalization
 - Increase input image from 224x224 to 448x448
 - Add anchor boxes:multi-object prediction per grid cell
 - · Backbone: Darnet19





PSL™ YOLO families

- YOLO9000 (2017):
 - It is a real-time framework for detecting more than 9000 object categories by jointly optimizing detection and classification. During training, they mix images from both detection and classification datasets.
- YOLOv3 (2018):
 - Backbone: Darknet53
 - better performance for small objects
 - worse performance on medium and larger size objects.
 - YOLOv3 predicts boxes at 3 different scales
- YOLOv4 (Alexey Bochkovskiy et al.,2020):
 - bag of freebies: CutMix and Mosaic data augmentation, DropBlock regularization, Class label smoothing
 - bag of specials: Mish activation, Cross-stage partial connections (CSP), Multi- input weighted residual connections (MiWRC)
- YOLOv5 (no published paper)



The speed and accuracy of YOLO v4 (source: YOLO v4 paper)

Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL★ YOLO versions

YOLO (darknet) - https://pjreddie.com/darknet/yolov1/ (C++)

YOLO v2 (darknet) - https://pjreddie.com/darknet/yolov2/ (C++)

- Better and faster - 91 fps for 288 x 288

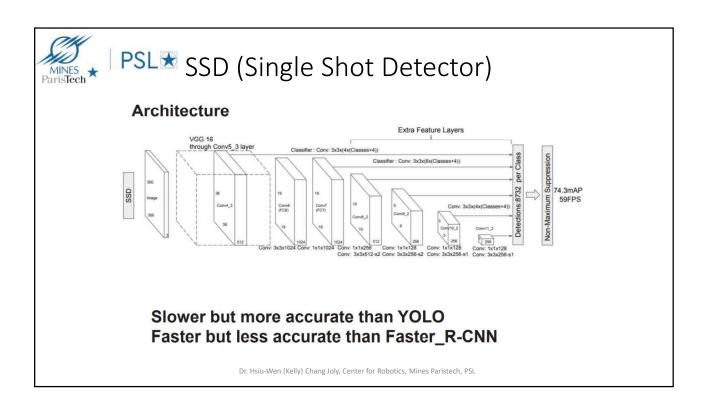
YOLO v3 (darknet) - https://pireddie.com/darknet/yolo/ (C++)

https://github.com/ultralytics/yolov3 (Pytorch)

YOLO (caffe) - https://github.com/xingwangsfu/caffe-yolo

YOLO (tensorflow) - https://github.com/thtrieu/darkflow

YOLO v5 (pytorch) - https://github.com/ultralytics/yolov5





PSLM SSD available source code versions

SSD (caffe) - https://github.com/weiliu89/caffe/tree/ssd

SSD (tensorflow) - https://github.com/balancap/SSD-Tensorflow

SSD (pytorch) - https://github.com/amdegroot/ssd.pytorch

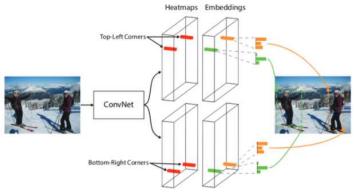


Hei Law, Jia Deng "CornerNet: Detecting Objects as Paired Keypoints", European Conference on Computer Vision (ECCV), 2018

- Detecting Objects as Paired Key points: the top-left corner and the bottom-right corner
- No anchor boxes
- Proposed corner pooling

 that helps the network
 better localize corners
- Backbone: Hourglass-104

$$L = L_{det} + \alpha L_{pull} + \beta L_{push} + \gamma L_{off} \quad --$$



 L_{det} is the detection loss, which is responsible for proper corner detection and is a variant of focal loss;

 L_{pull} is the grouping loss, used to pull corners of one object together; $L_{push},$ is, on the opposite, used to separate corners of different objects; L_{off} is the smooth L1 loss used for offset correction; and $\alpha,\,\beta$ and γ are parameters, which are set to 0.1, 0.1 and 1 respectively.

Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



- Base Network
 - VGG16
 - ResNet-101
 - Inception V2
 - Inception V3
 - Inception
 - ResNet
 - MobileNet

- Object Detection architecture
 - Faster R-CNN
 - R-FCN
 - SSD
 - YOLO

Recommended reading:

Huang et al, "Speed/accuracy trade-offs for modern convolutional object detectors", CVPR 2017

R-FCN: Dai et al, "R-FCN: Object Detection via Region-based Fully Convolutional Networks", NIPS 2016
Inception-V2: Ioffe and Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal
Covariate Shift", ICML 2015

Inception V3: Szegedy et al, "Rethinking the Inception Architecture for Computer Vision", arXiv 2016 Inception ResNet: Szegedy et al, "Inception-V4, Inception-ResNet and the Impact of Residual Connections on Learning", arXiv 2016

MobileNet: Howard et al, "Efficient Convolutional Neural Networks for Mobile Vision Applications", arXiv 2017

istech, PSL

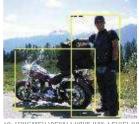


PSLM Training sets for Visual objects detection

- Training a visual objects detector requires a training set contain images WITH BOUNDINGBOXES (or even mask) ANNOTATION
- Two main « reference » training sets of this type:
 - Pascal VOC (Visual Object Class) http://host.robots.ox.ac.uk/pascal/VOC/

• Coco (Common Objects in Context) [more classes + MASK annotations] http://cocodataset.org/

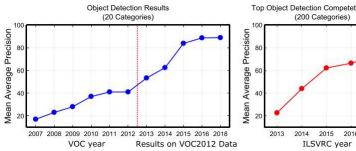


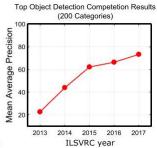






PSL An overview of recent object detection performance





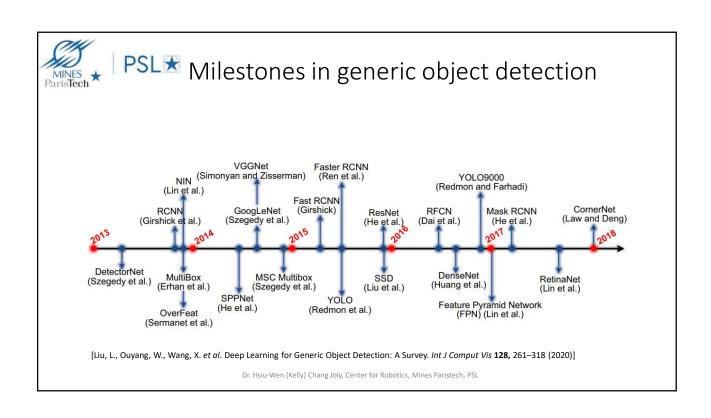
Turning point at 2012: Deep learning achieved record breaking Image Classification Result

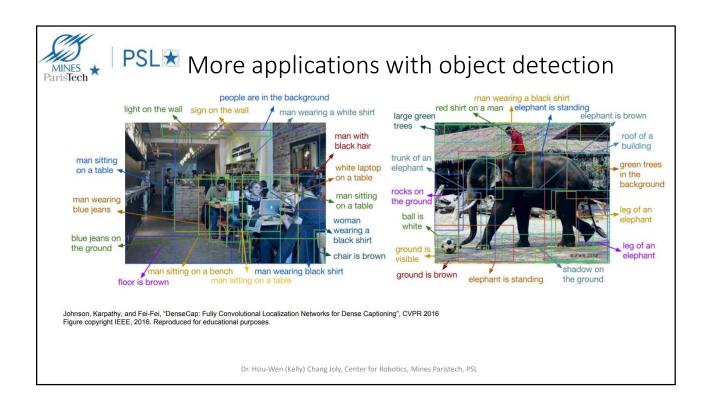


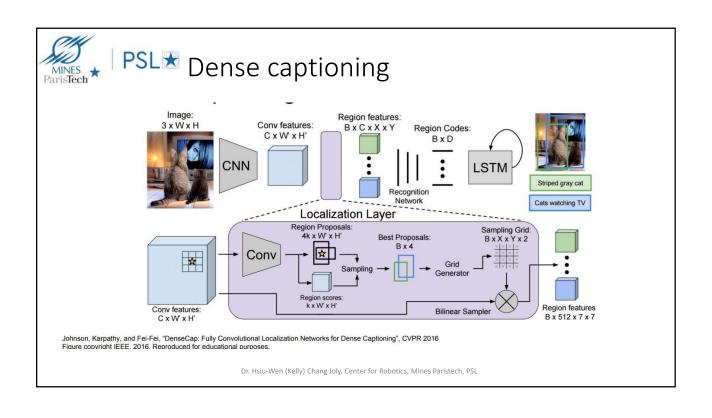
PSL The difficulties for image classification

- (a-h) Changes in appearance of the same class with variations in imaging conditions
- (i) There is an astonishing variation in what is meant to be a single object class
- (j) In contrast, the four images in j appear very similar, but in fact are from four different object classes.
- Most images are from ImageNet (Russakovsky et al. 2015) and MS COCO (Lin et al. 2014)



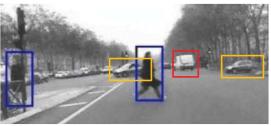








PSLM Drawbacks of object detections





- Problem for objects without sharp boundaries (trees, ...) or very dense group of objects (crowd of pedestrians, ...)
- Only « compact » objects are categorized (what about « road », « sidewalk », « building », ...?)

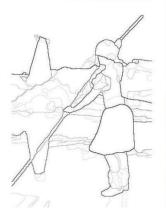
Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL What is image segmentation

• Identify groups of continuous pixels that go together





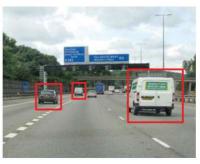


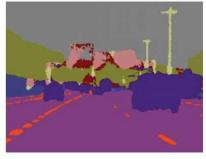


| PSL★ Advantage of Semantic (full) segmentation

- One single semantic segmentation →all interesting object categories (cars, pedestrians, signs, etc...) and categorization of whole image
- Can also categorize non-compact areas (road, sky, buildings, trees, traffic lanes...)







Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL Many approaches for image segmentation

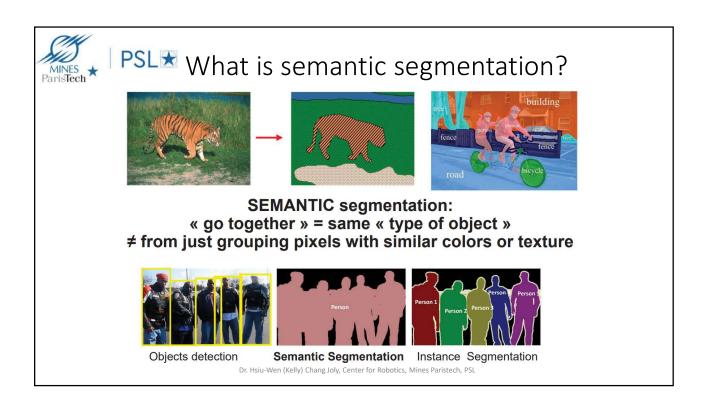
- Clustering (K-means, GMM, MeanShift, ...)
- Graph-based (graph-cuts)

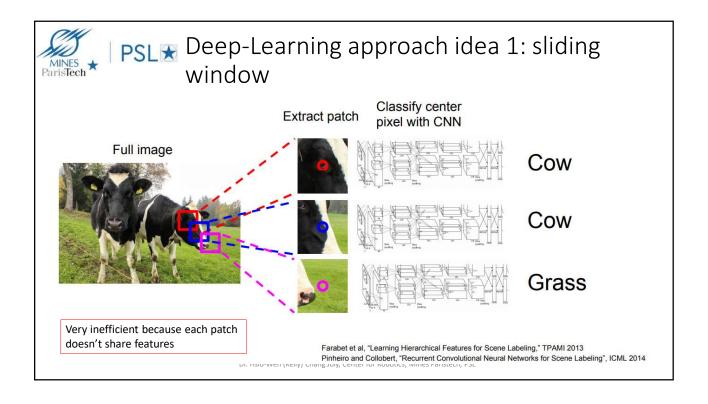


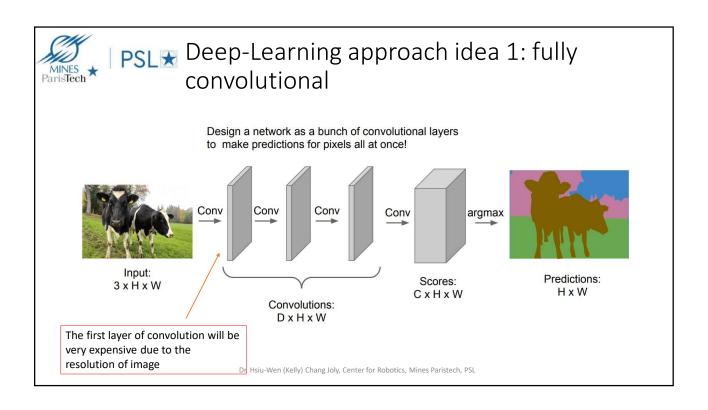


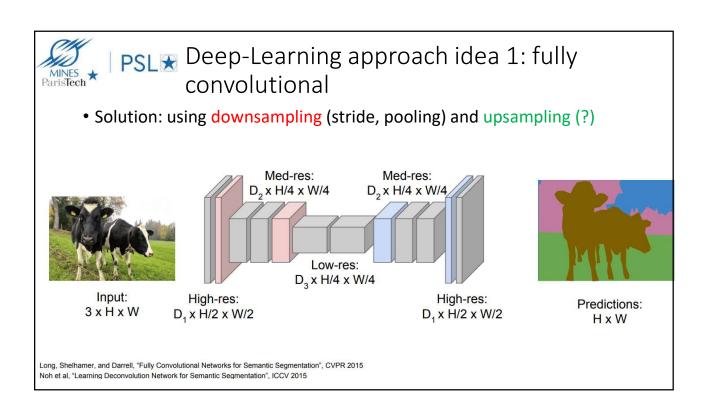
- Node (vertex) for every pixelEdge between pairs of pixels, (p,q)
- Affinity weight w_{pq} for each edge

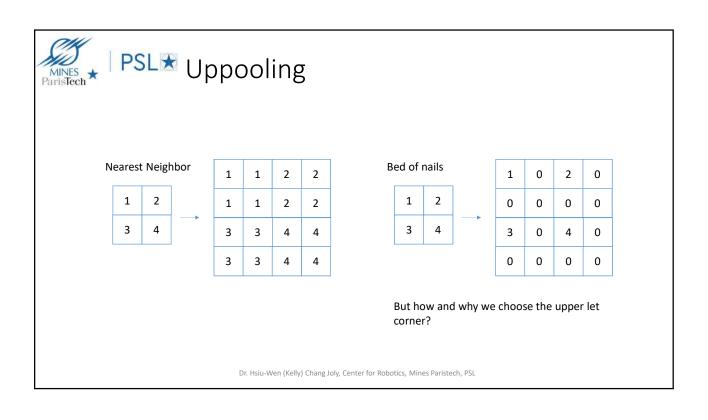
 w_{pq} measures similarity
 Similarity is inversely proportional to difference (in color and position...)
- Mathematical Morphology (watershed, etc...)
- Energy minimization (Conditional Random Fields)
- Deep-Learning

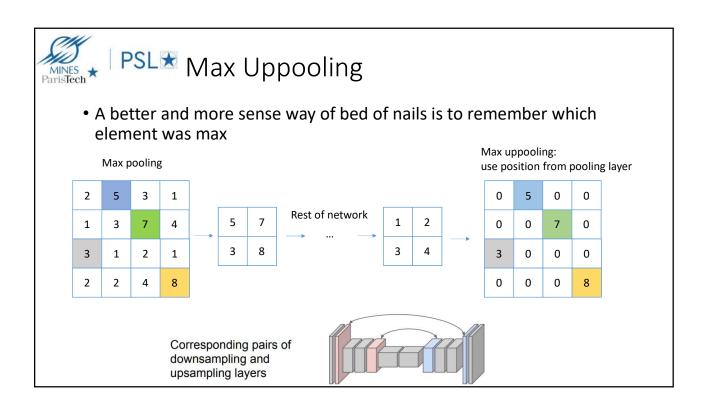


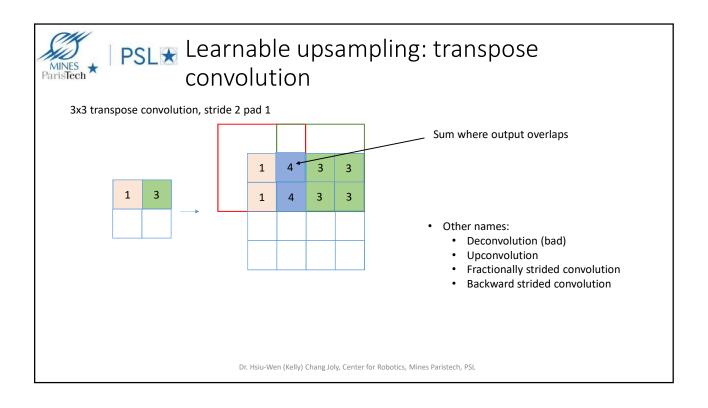


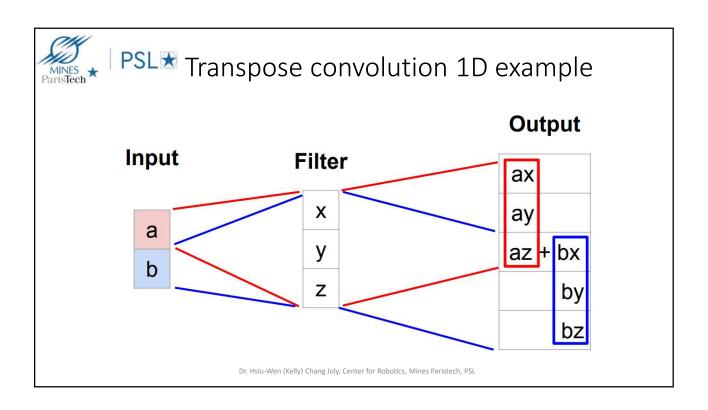














PSL Convolution as matrix operation

We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & z & 0 & 0 & 0 \\ 0 & x & y & z & 0 & 0 \\ 0 & 0 & x & y & z & 0 \\ 0 & 0 & 0 & x & y & z \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

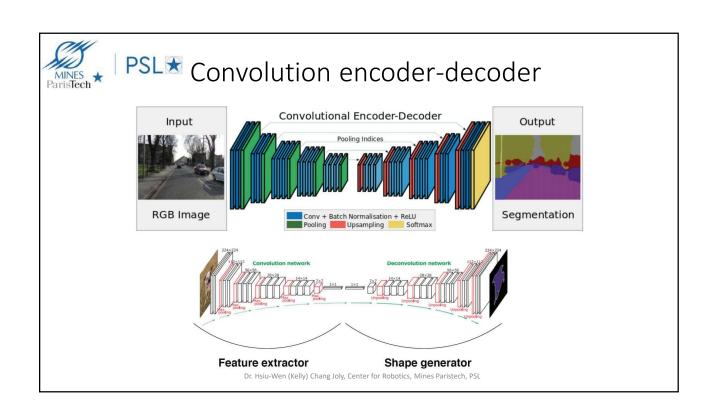
Example: 1D conv, kernel size=3, stride=1, padding=1

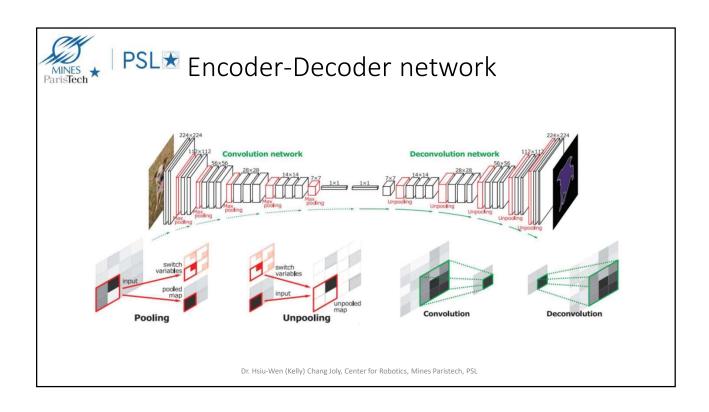
Convolution transpose multiplies by the transpose of the same matrix:

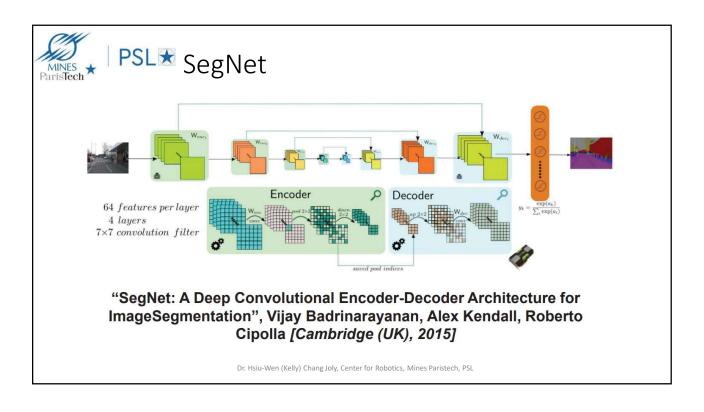
$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix} \qquad \begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

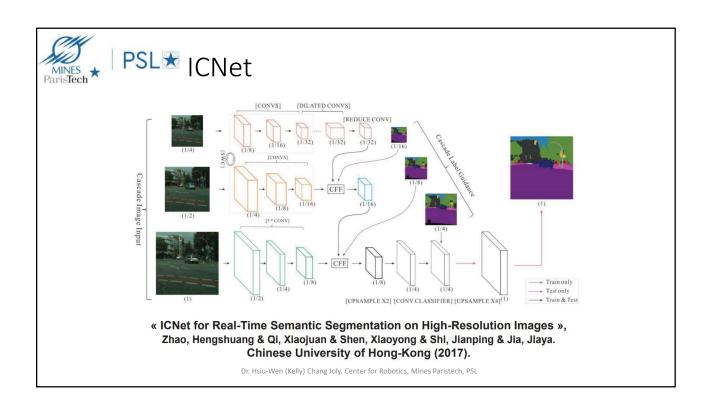
When stride=1, convolution transpose is just a regular convolution (with different padding rules)





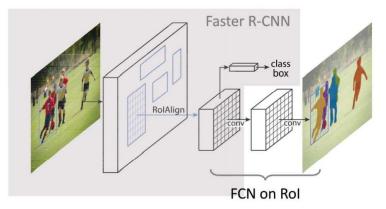








Mask R-CNN = Faster R-CNN with FCN on Rols



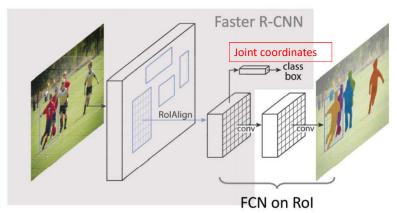
Mask R-CNN architecture extract detailed contours and shape of objects instead of just bounding-boxes

Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL Addition function of Mask R-CNN

Mask R-CNN = Faster R-CNN with FCN on Rols



Mask R-CNN architecture extract detailed contours and shape of objects instead of just bounding-boxes



PSL Results from Mask R-CNN





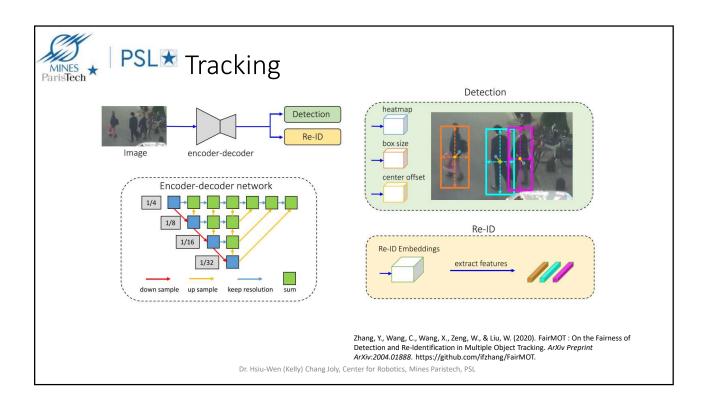


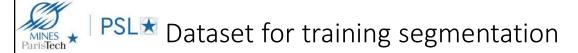
Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSLM And many other competitors

- 2015: U-Net (Keras)
 - https://github.com/zhixuhao/unet
- RefineNet (2016)
- DeepLab (Caffe)
 - https://github.com/Robotertechnik/DeepLab
- DeepLabv3 (Tensorflow)
 - https://github.com/NanqingD/DeepLabV3-Tensorflow





- It is very expensive to create dataset for training segmentation
- Synthetic images are popular in this case. They are very real.



Example from SYNTHIA

http://synthia-dataset.net

oi. Insid-vveir (keily) Chang Joly, Center for Robotics, Ivilles Falistech, FSL



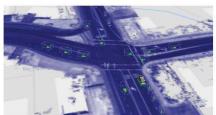
Interest of synthetic images for Machine-Learning in IV applications

- Possible to generate as many as needed at nearly no cost (in particular compared to recording while driving)
- Easy to generate controlled variability in environment, luminosity conditions, scenarii, etc + also images « dangerous situations »
- NO NEED FOR MANUAL LABELLING: ground truth (ie target value) for classifiers, localizers, and semantic segmentation provided automatically

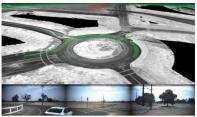
Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL Simulators dedicated to Autonomous Vehicles



Scenario-builing with CarCraft by Google/Waymo



Simulation of a virtual scenario in XView by Google/Waymo



, PSL



PSL CARLA open-source urban driving simulator

 Still few driving simulators adapted for DL and RL, and best ones not totally mature

Simulateur	GTA	DeepDrive.io	AirSim	CARLA[1]
Flexibilité		++	++	++
Variété	++		_	+
Complexité/Réalisme	++		_	_
Objets mobiles	++			+
Vitesse éxecution		+	+	+
Multi-agent		-	-	++

→ Choice of CARLA

[1] A. Dosovitskiy: CARLA: An Open Urban Driving Simulator (2017)

Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



PSL® Synthetic images use in ML/DL for IV

- Initial training of a classifier / segmented / controller only on simulated images / videos / scenarios
- Possible to then adaptation to real-world by fine-tuning on REAL images/video datasets
- Cheaper / more extensive testing than on real-world videos
- REINFORCEMENT LEARNING in simulation!



PSL ■ We ask computer to do it?!

- Is it possible we ask the computer to create the dataset for us?
- The answer becomes more and more clear to us in the past few year
- GameGAN:
 - generative model that learns to visually imitate a desired game by ingesting screenplay and keyboard actions during training
 - CameGAN

 Memory

 Action a
 Random z
 Noise
 Noise
 Memory m_{t-1}
 Image
 Engine

 Rendering
 Engine

 Image
 X_{t+1}

- Dall-e (OpenAI, 2021):
 - creates images from text captions for a wide range of concepts expressible in natural language.
 - https://openai.com/blog/dall-e/

Dr. Hsiu-Wen (Kelly) Chang Joly, Center for Robotics, Mines Paristech, PSL



Question?

hsiu-wen.chang_joly@mines-paristech.fr