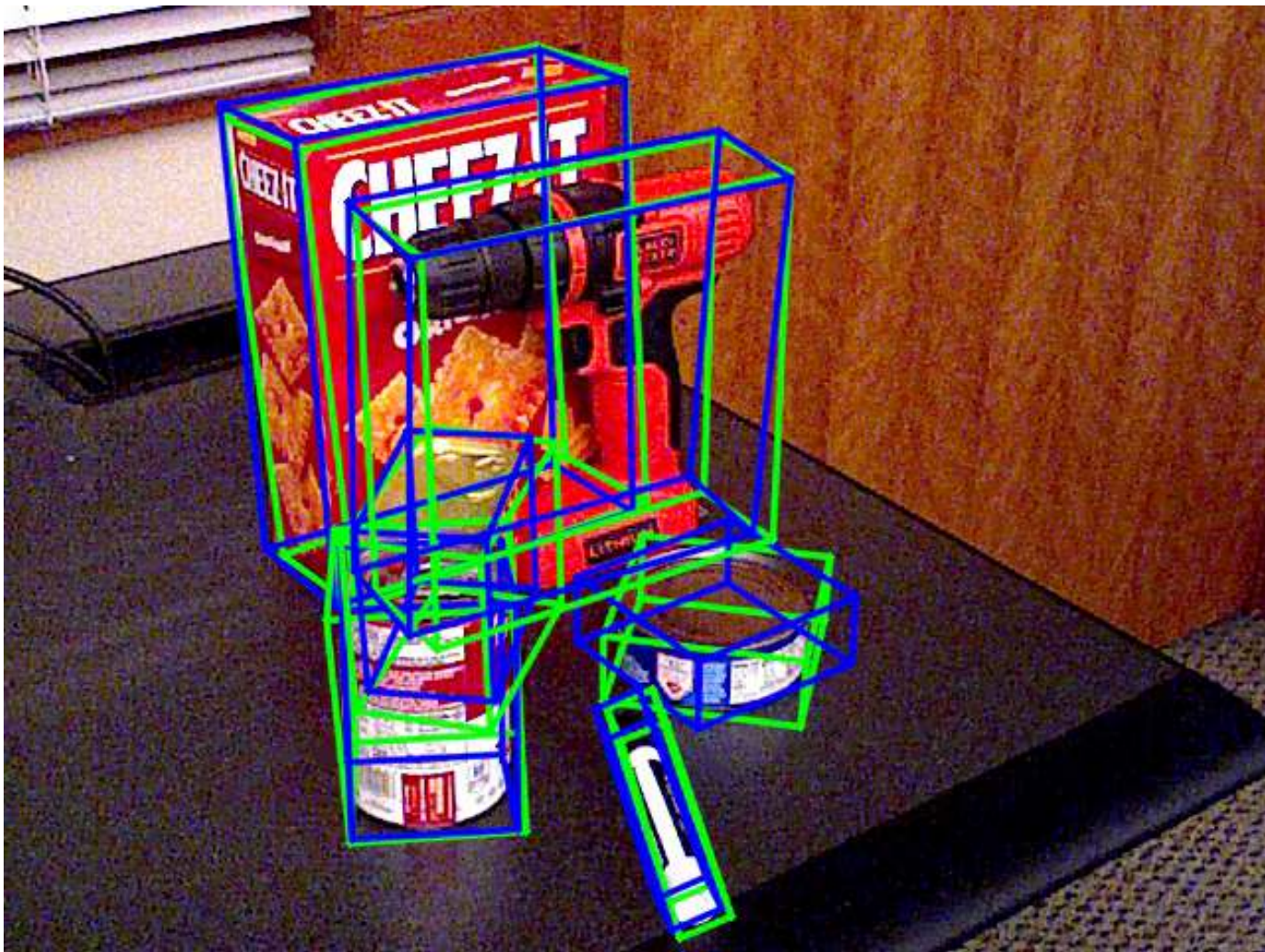
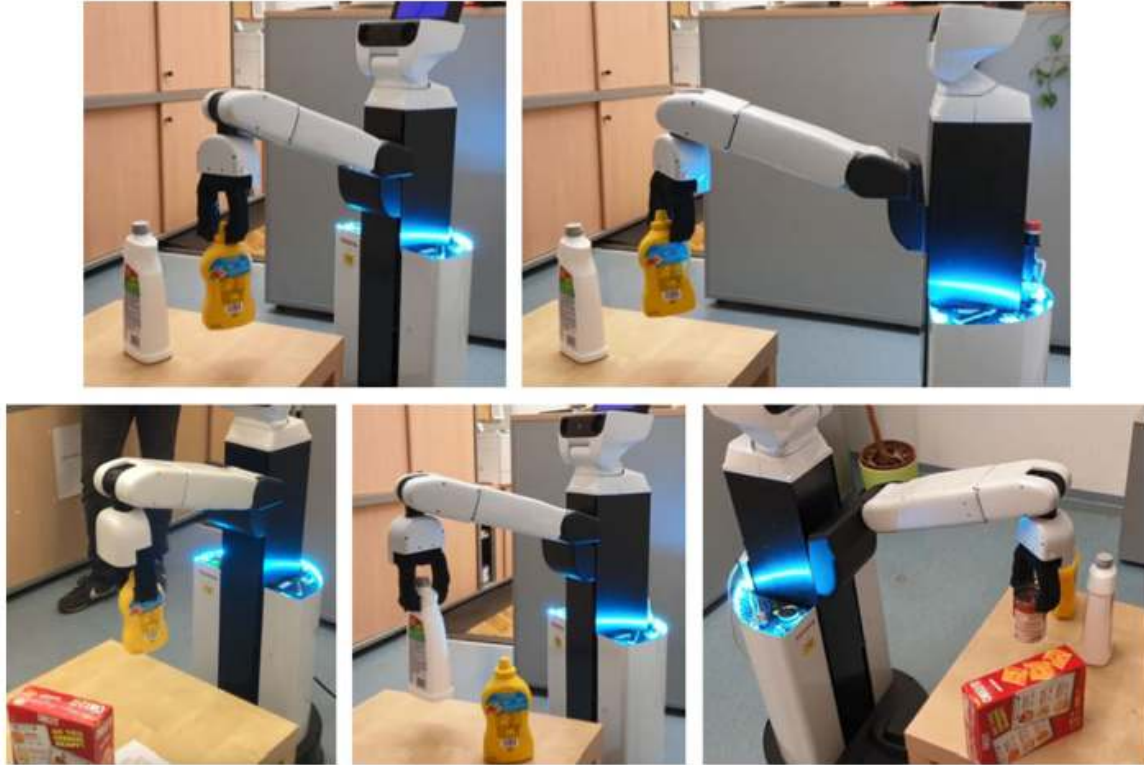


3D Object Detection and Pose Estimation

Vincent Lepetit



possible applications



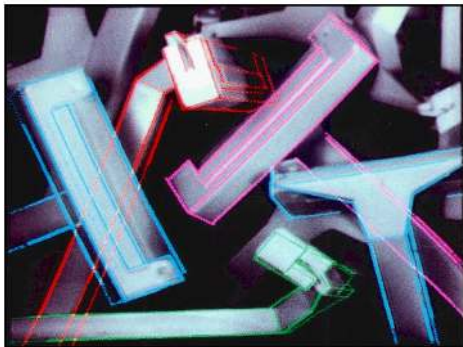
[Vincze et al, 2020]

possible applications

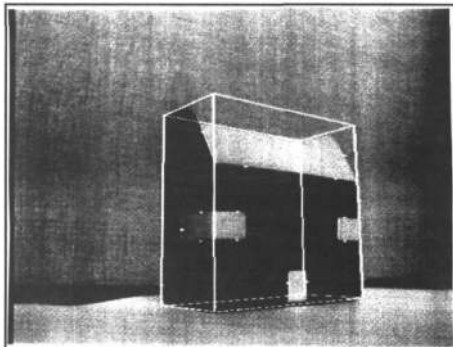


[Petit et al, ISMAR 2013]

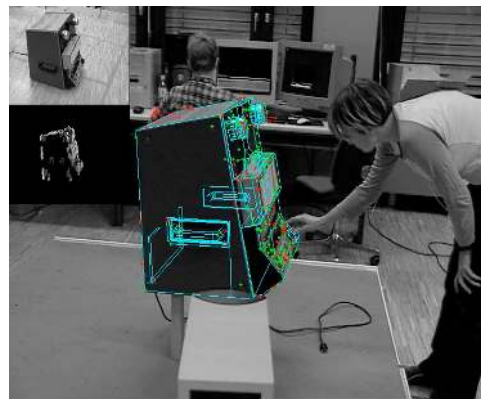
a bit of history



[Lowe, 1987]



[Harris&Stennet, 1990]



[Vacchetti et al, CVPR 2003]



[Lepetit et al, CVPR 2004]

more modern take

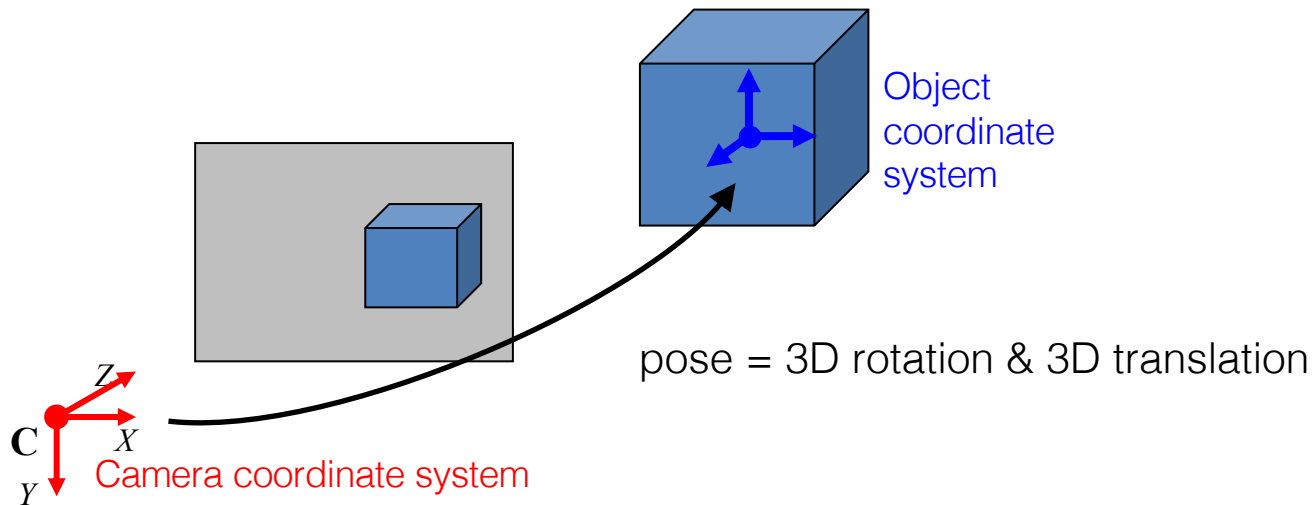
training set: (many) annotated real images



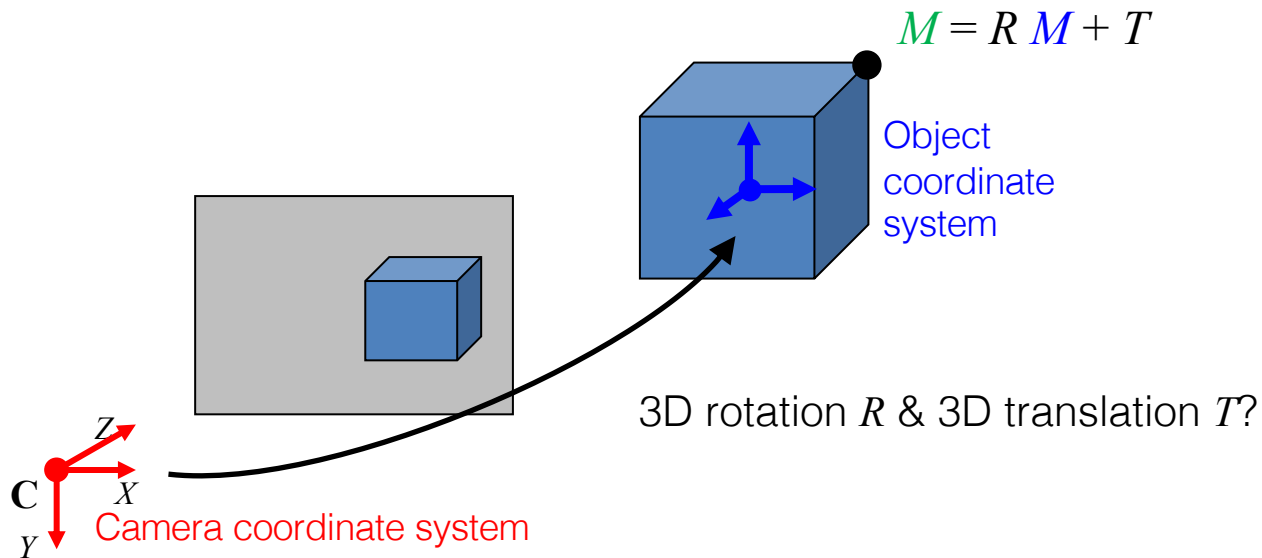
training set: About 200 annotated real images



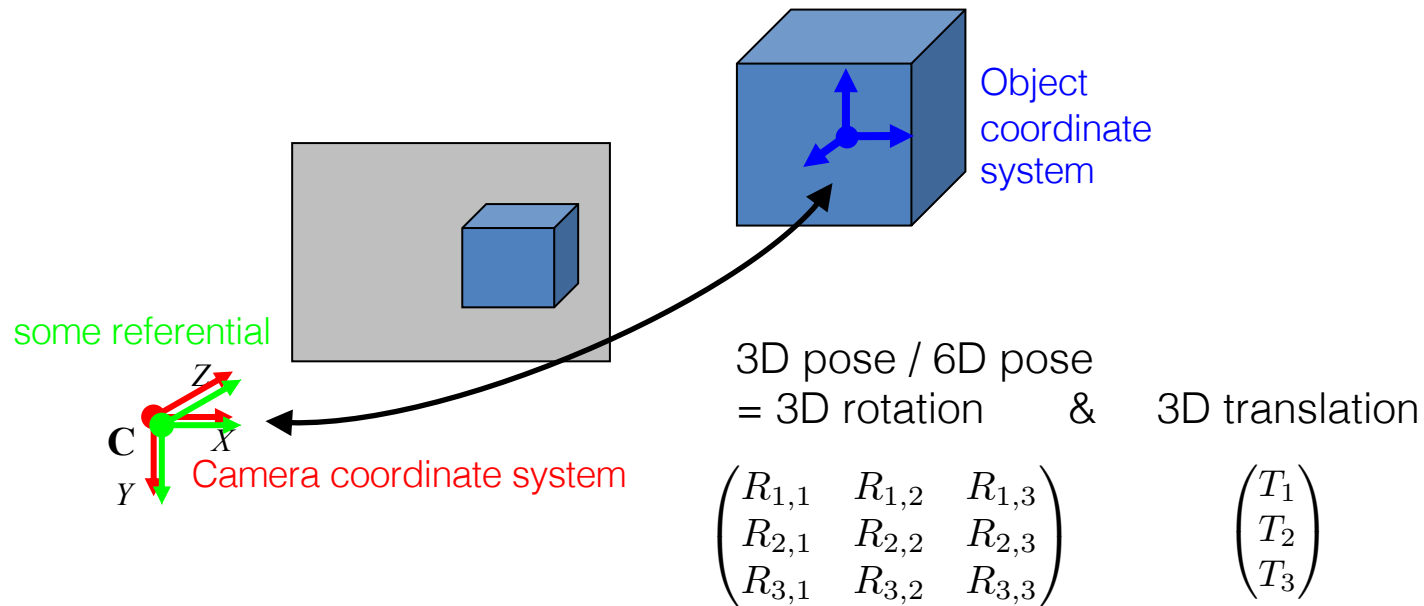
3D Pose / 6D Pose



3D Pose / 6D Pose



3D Pose / 6D Pose



loss for pose prediction

For the 3D translation, simply the Euclidean distance between prediction and ground truth:

$$\mathcal{L}_T = \|T - \hat{T}\|^2$$

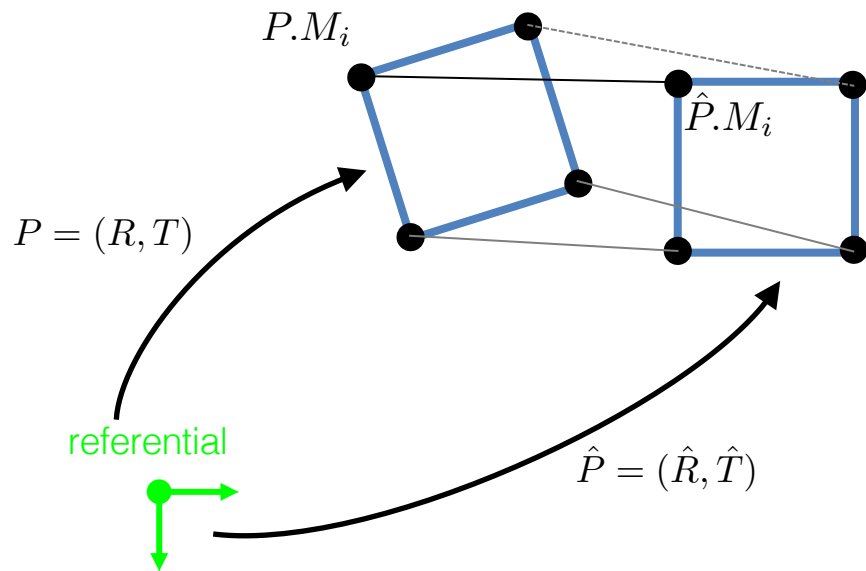
For the 3D rotation, geodesic distance:

$$\begin{aligned}\mathcal{L}_R &= \|\log(R\hat{R}^\top)\|_F \\ &= \cos^{-1}(\text{tr}(R\hat{R}^\top) - 1)/2\end{aligned}$$

For the full 3D pose:

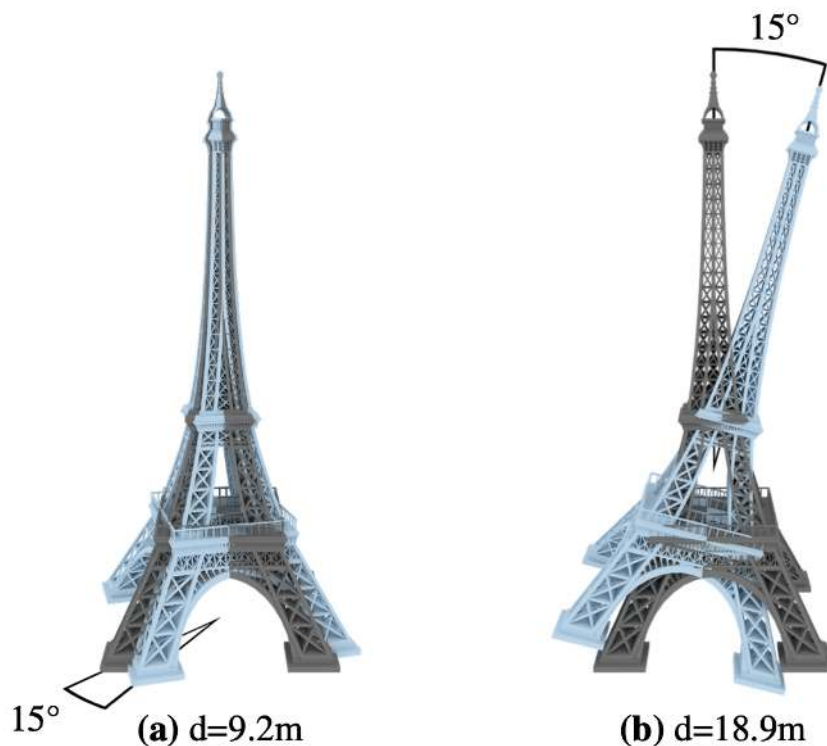
$$\mathcal{L}_{\text{pose}} = \mathcal{L}_T + \gamma \mathcal{L}_R$$

alternative loss for pose prediction



$$\left\{ \begin{array}{l} \mathcal{L} = \sum_i \|P.M_i - \hat{P}.M_i\|^2 \\ P.M_i = RM_i + T \\ \hat{P}.M_i = \hat{R}M_i + \hat{T} \end{array} \right.$$

alternative loss for pose prediction

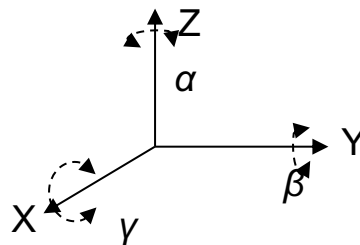


R. Brégier et al. Defining the Pose of Any 3D Rigid Object and an Associated Distance. IJCV, June 2018.

possible parameterizations of the rotation matrix

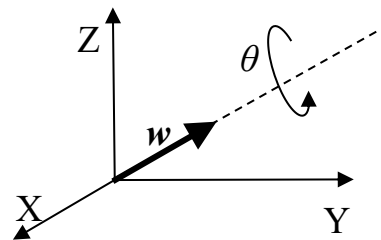
- Directly the rotation matrix (ie 9 values);
- Euler angles (3 values):

$$\mathbf{R} = \begin{bmatrix} \cos \alpha & -\sin \alpha & 0 \\ \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \beta & 0 & \sin \beta \\ 0 & 1 & 0 \\ -\sin \beta & 0 & \cos \beta \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \gamma & -\sin \gamma \\ 0 & \sin \gamma & \cos \gamma \end{bmatrix}$$

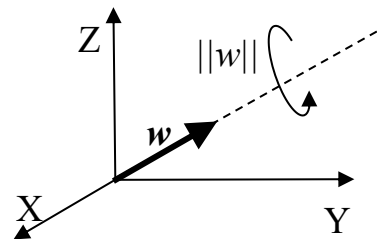


- A unit quaternion (4 values):

$$q = \left(\cos \frac{\theta}{2}, \mathbf{w} \sin \frac{\theta}{2} \right)$$

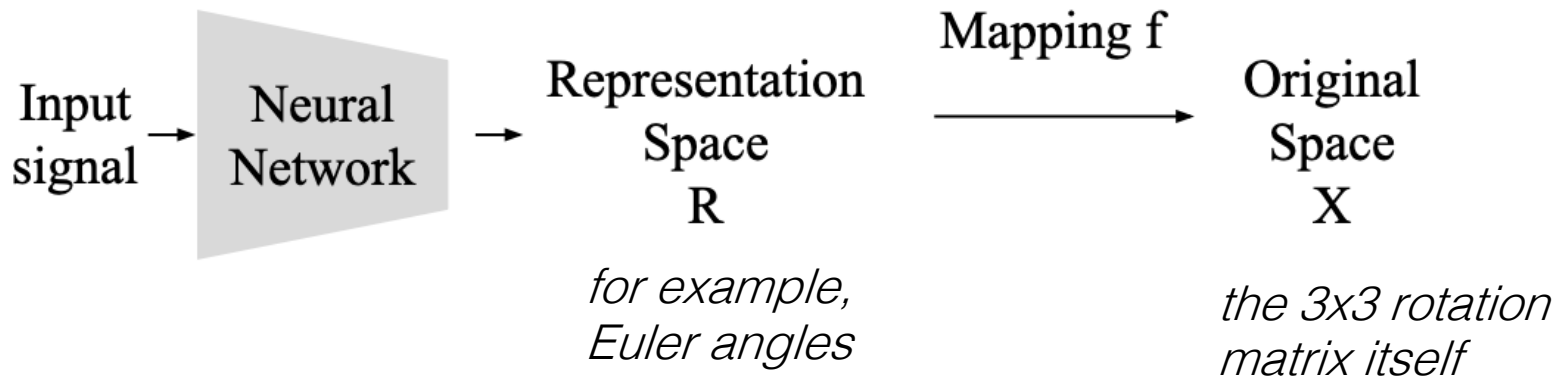


- An Exponential Map (a unit 3-vector, 3 values):

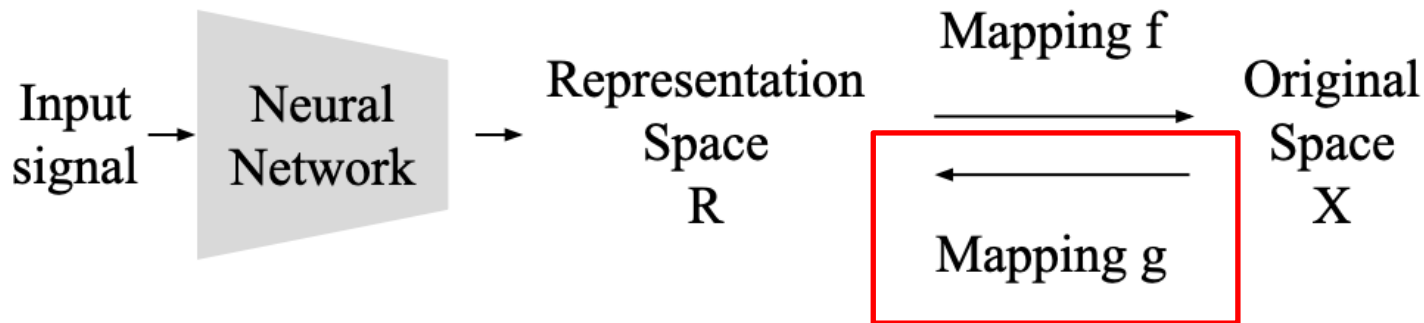


...

the problem with these representations



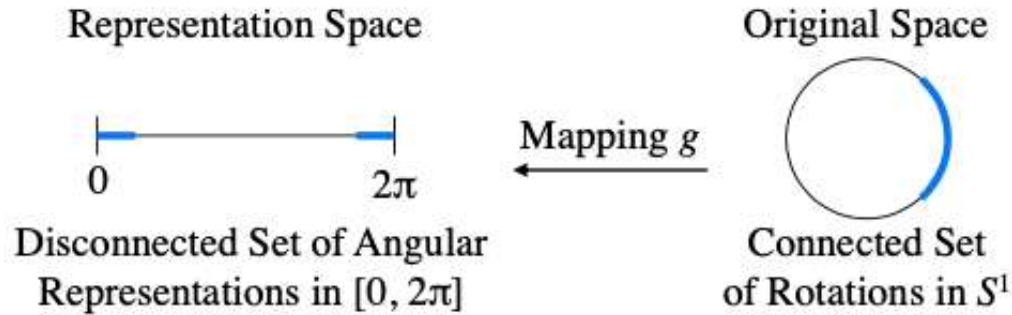
the problem with these representations



Needed for training the network
(in back-propagation).

*Not continuous for these
rotation representations*

discontinuities of g



proposed solution

- 2 3-vectors (6 values): e_1, e_2

$$e'_1 = \frac{e_1}{\|e_1\|_2}$$

$$e'_3 = \frac{e'_1 \wedge e_2}{\|e_2\|_2} \quad R = (e'_1 \quad e'_2 \quad e'_3)$$

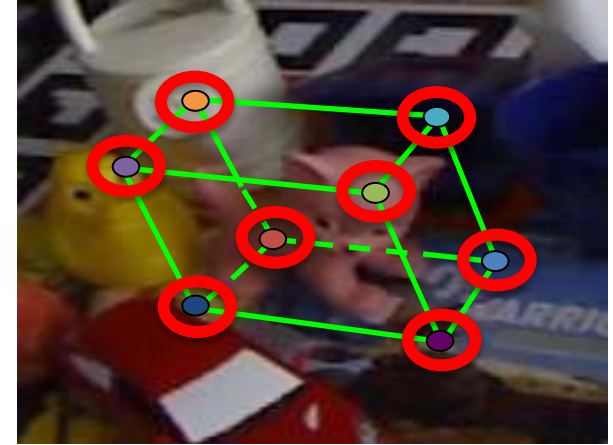
$$e'_2 = e'_3 \wedge e'_1,$$

It is then possible to define a $g(R) = (e_1, e_2)$ function that is continuous.

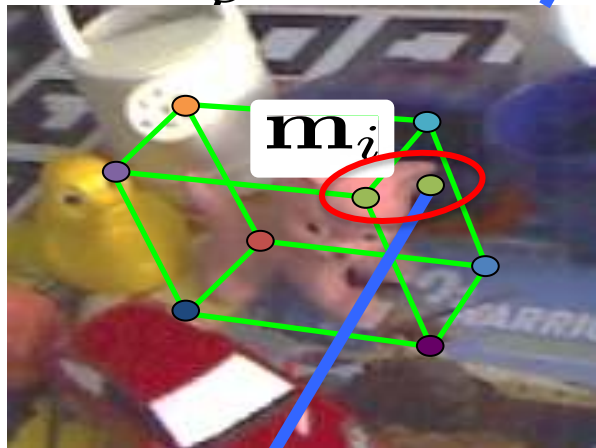
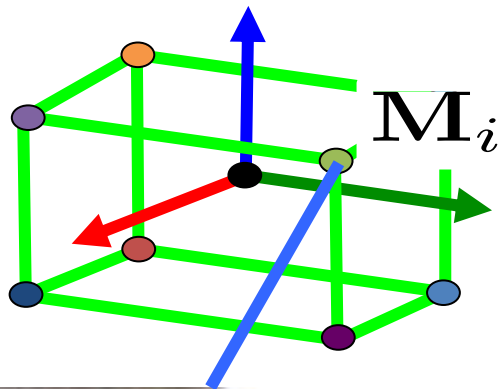
alternative predictions (1)



the 2D projections of the 8 corners of the 3D bounding box

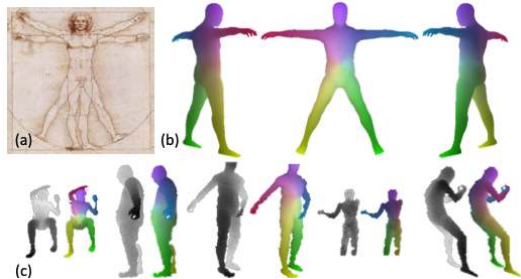


3D pose estimation from correspondences

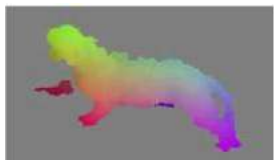


○ Camera center

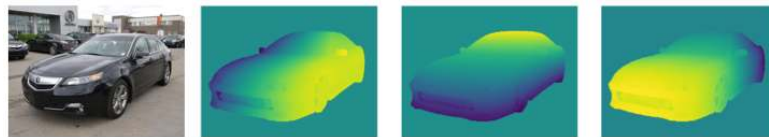
Alternative predictions (2)



Taylor et al. The Vitruvian Manifold: Inferring Dense Correspondences for One-Shot Human Pose Estimation. CVPR 2012.



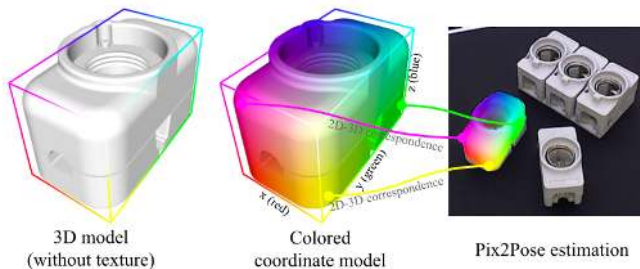
E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother. Learning 6D Object Pose Estimation using 3D Object Coordinates. ECCV 2014.



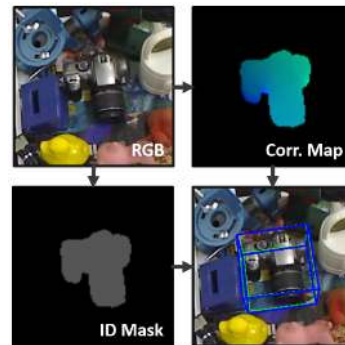
Location Fields. Wang et al., ECCV 2018



Normalized Object Coordinate Space for Category-Level 6D Object Pose and Size Estimation. Wang et al., CVPR 2019.

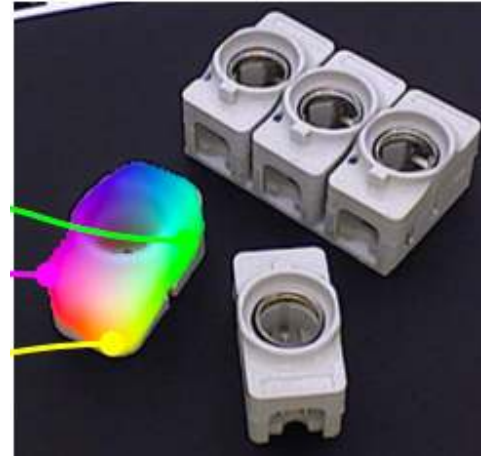
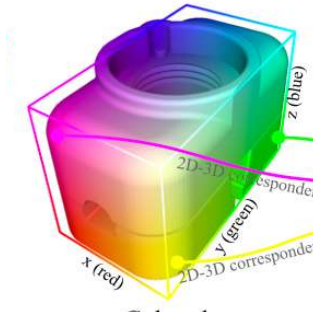


Pix2Pose: Pixel-Wise Coordinate Regression of Objects for 6D Pose Estimation. Park et al., CVPR 2019.



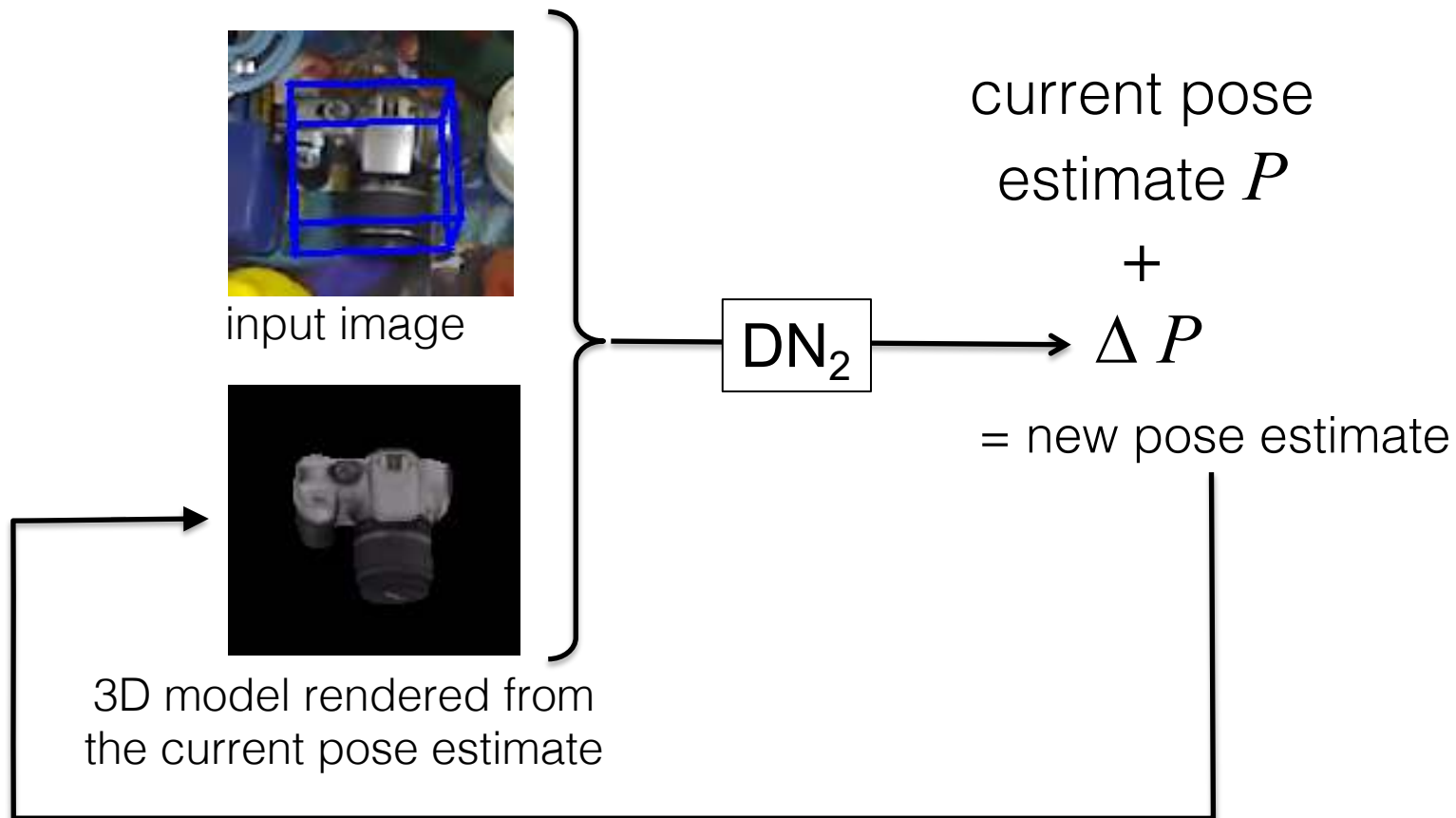
DPOD: 6D Pose Object Detector and Refiner. Zakharov et al. ICCV 2019.

how to use 3D coordinate maps

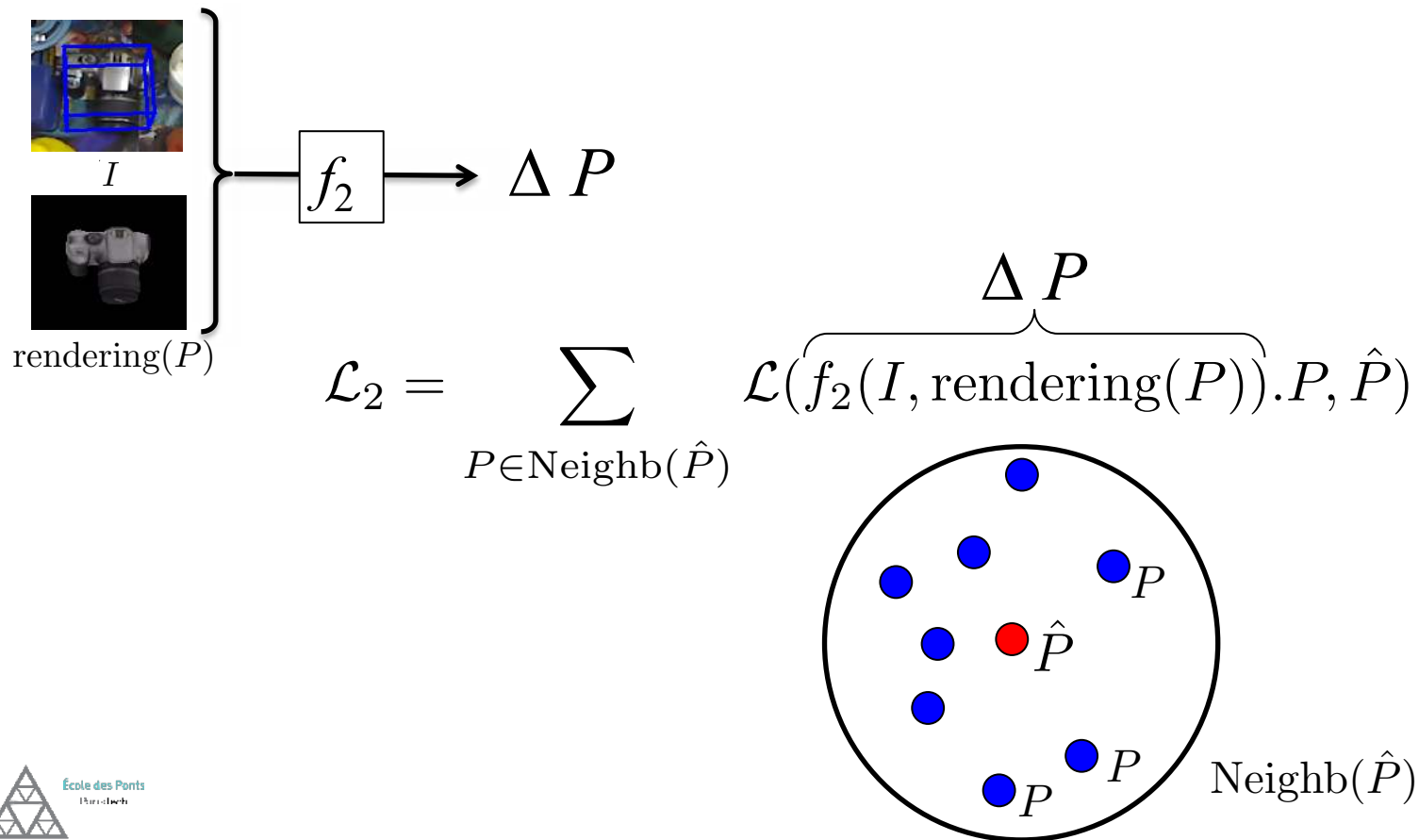


pose refinement

refining the pose



refining the pose, why it works





DeepIM: Decoupled Coordinates



Rot Axis:
[0, 0, 1]
Rot Angle:
0
Trans:
[0]
[0]
[0]



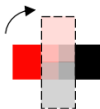
Rot Axis:
[0, 0, 1]
Rot Angle:
90
Trans:
[0]
[0]
[0]



Rot Axis:
[0, 0, 1]
Rot Angle:
90
Trans:
[0]
[0]
[0]



Rot Axis:
[0, 0, 1]
Rot Angle:
90
Trans:
[0]
[0]
[0]



Rot Axis:
[1, 0, 0]
Rot Angle:
90
Trans:
-0.25
0.25
[0]



Rot Axis:
[0, 0, 1]
Rot Angle:
-90
Trans:
[0]
[0]
[0]



Rot Axis:
[0, 0, 1]
Rot Angle:
90
Trans:
[0]
[0]
[0]

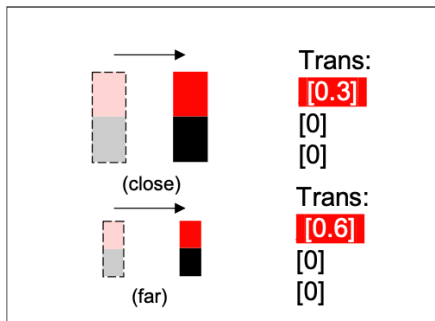
(a) Initial pose

(b) Camera coord.

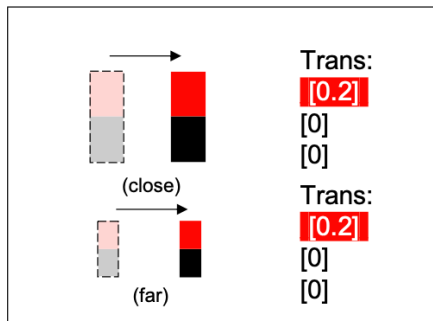
(c) Model coord.

(d) disentangled coord.

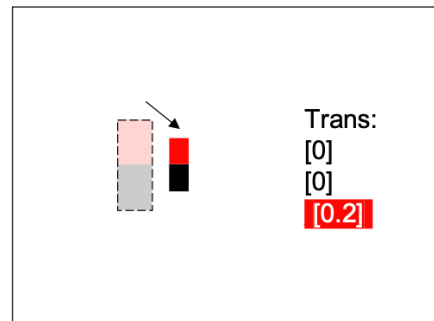
DeepIM: Decoupled Coordinates (T)



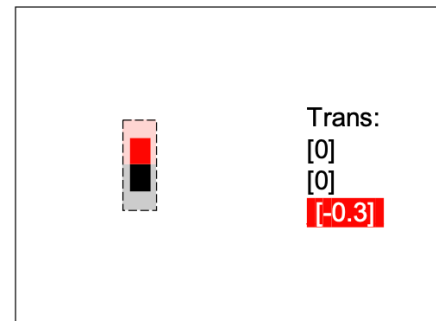
(a) Camera coord. xy-plane translation



(b) Disentangled coord. xy-plane translation

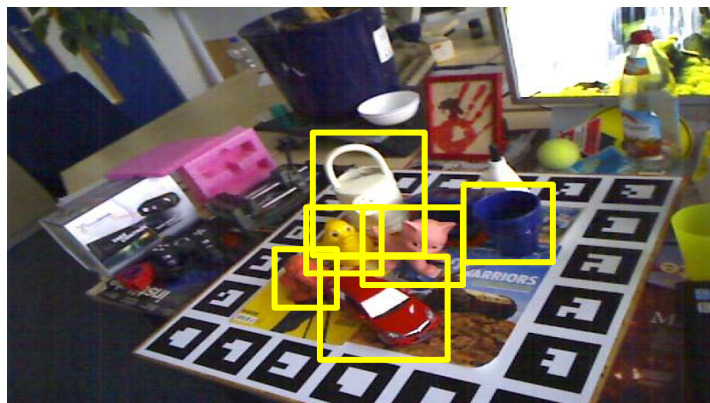
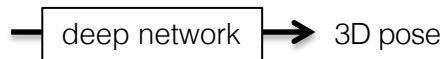


(c) Camera coord. z-axis translation



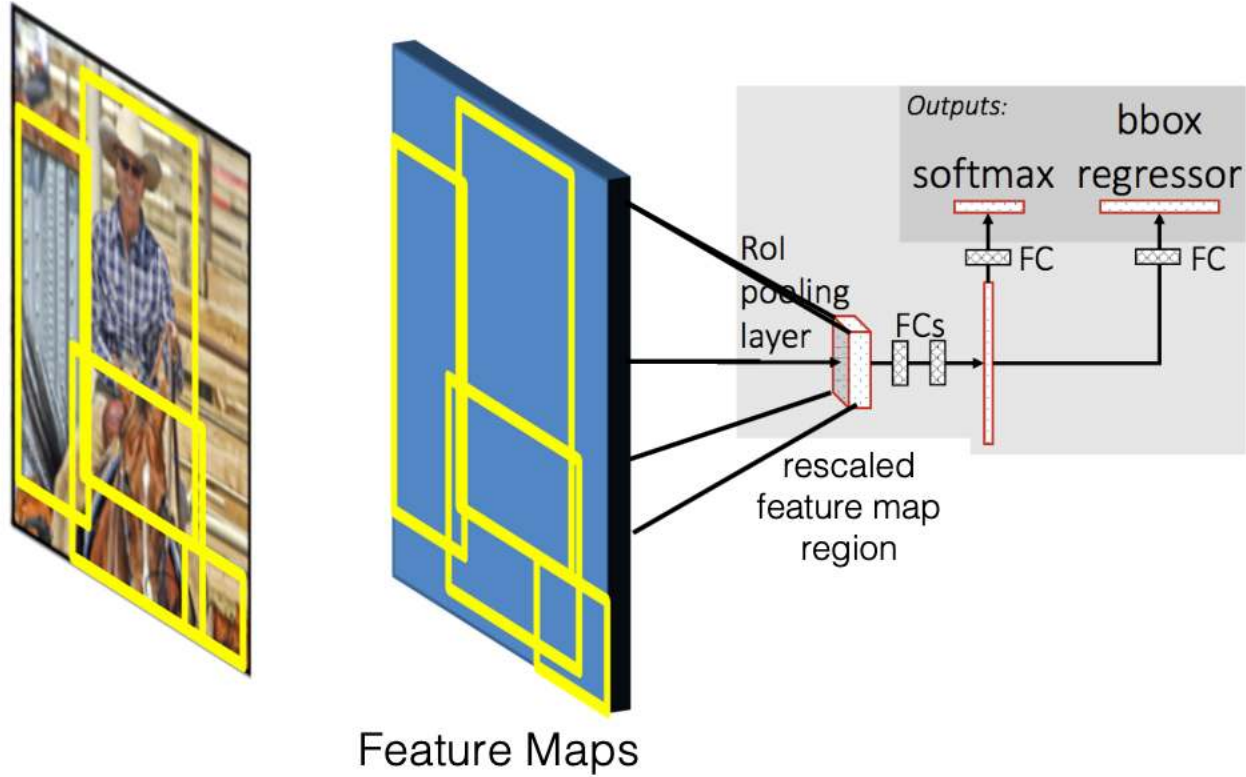
(d) Disentangled coord. z-axis translation

2D detection

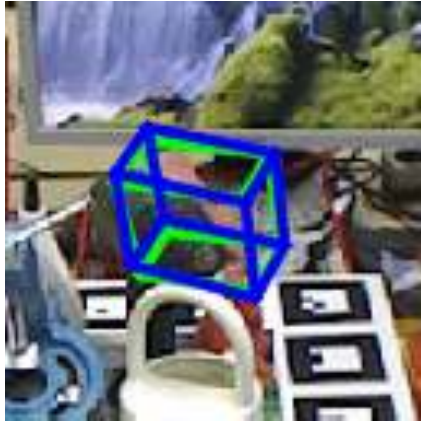


SSD-6D: Making RGB-Based 3D Detection and 6D Pose Estimation Great Again. Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, Nassir Navab. ICCV 2017.

Fast-RCNN / Mask-RCNN / Detectron2



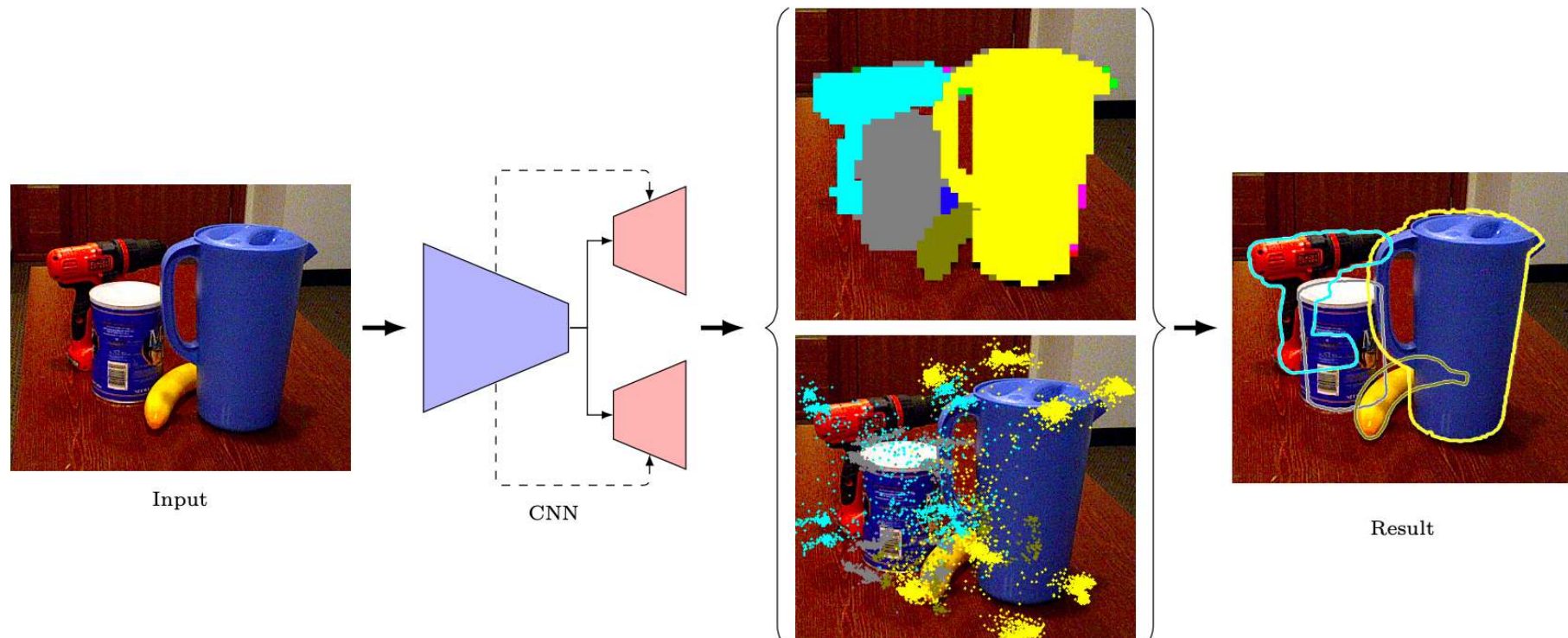
Dealing with Partial Occlusion



Avoid Occlusions in the Input



Voting for the pose



Voting for the pose



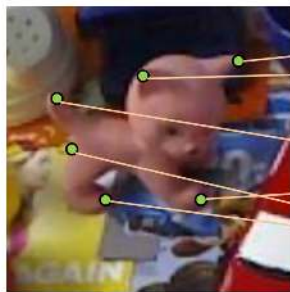
(a) Input image



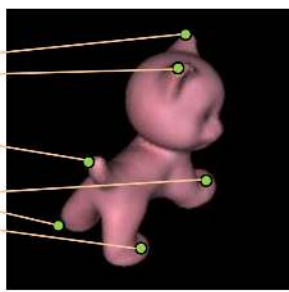
(b) Vectors



(c) Voting



(d) 2D keypoints



(e) 3D keypoints



(f) Aligned model

PVNet: Pixel-wise Voting Network for 6DoF Pose Estimation. Peng et al. CVPR 2019.

Training Set: About 200 Real Images + ...



... Data Augmentation (1)



Data Augmentation and Domain Randomization



Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World. Tobin et al. IROS 2017.

CosyPose: Consistent multi-view multi-object 6D pose estimation. Labbé et al. ECCV 2020.

How Domain Randomization Works

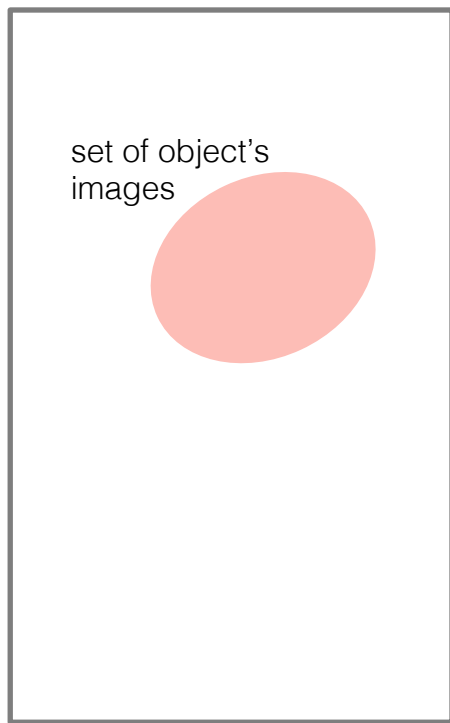


image space

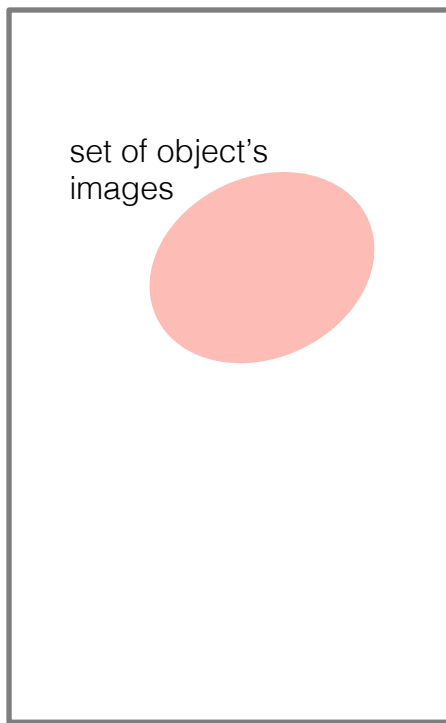


image space

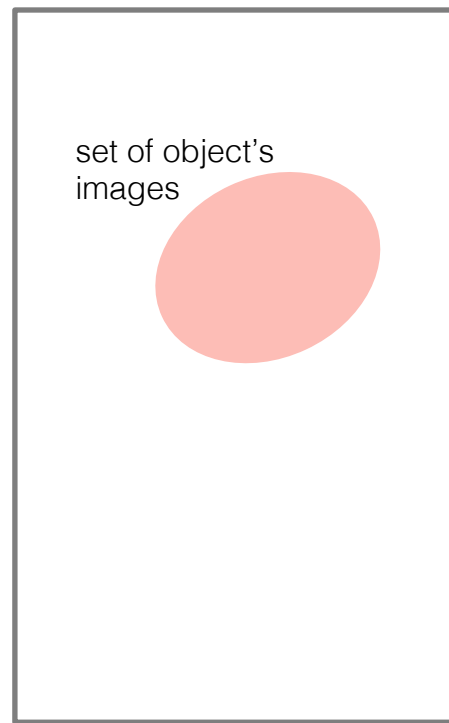


image space

limiting the need for training data and for training time

- Considering object categories;
- Few-shot learning;
- ...

3D Pose Prediction for Object Categories

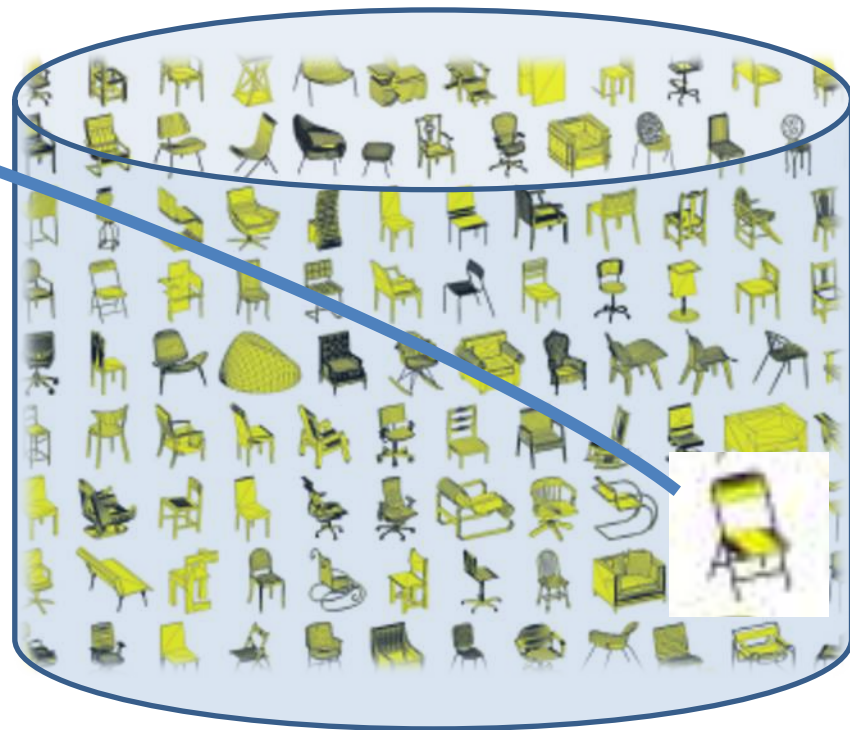


3D Pose Estimation and 3D Model Retrieval for Objects in the Wild. Alexander Grabner, Peter M. Roth, and Vincent Lepetit. CVPR 2018.

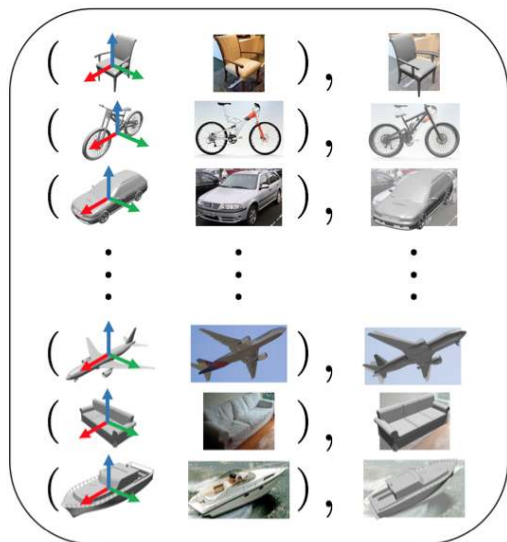


Pix3D: Dataset and Methods for Single-Image 3D Shape Modeling. Sun et al, 2018.

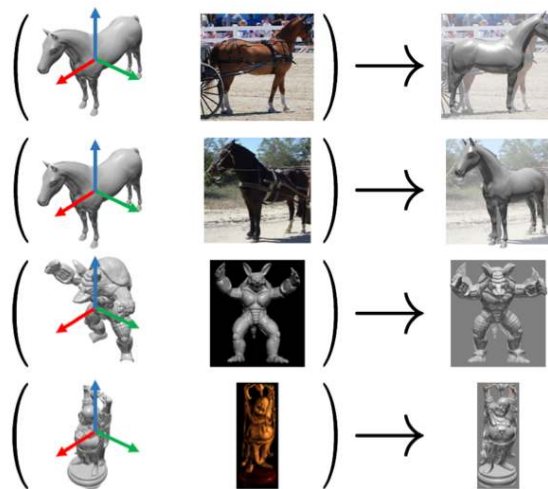
pose-invariant embedding



Pose from Shape



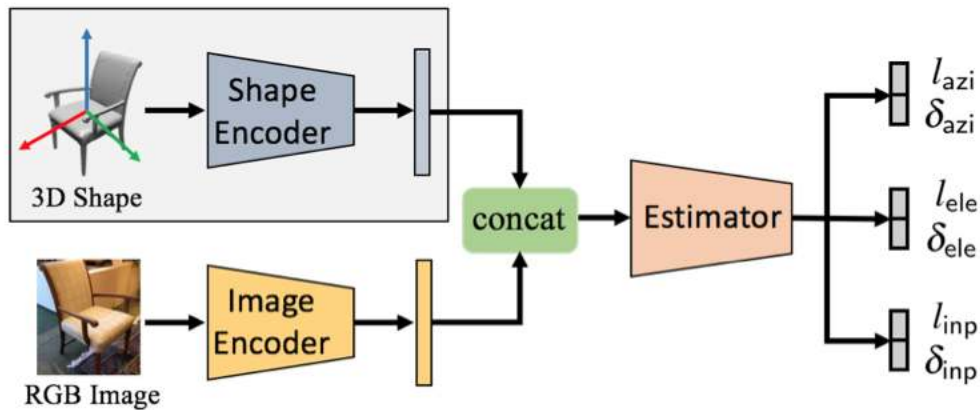
(a) Training with shape and pose



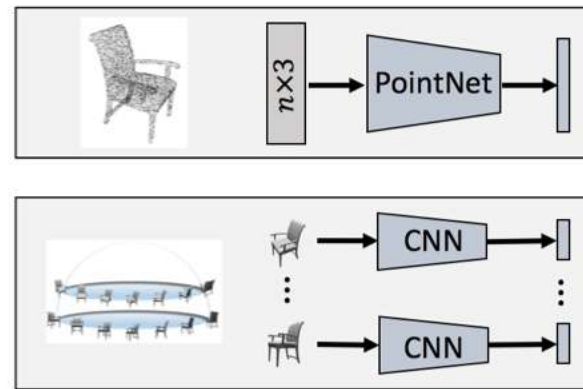
(b) Testing on unseen objects

Pose from Shape: Deep Pose Estimation for Arbitrary 3D Objects. Xiao et al. BMVC 2019.

Pose from Shape



(a) Our pose estimation approach

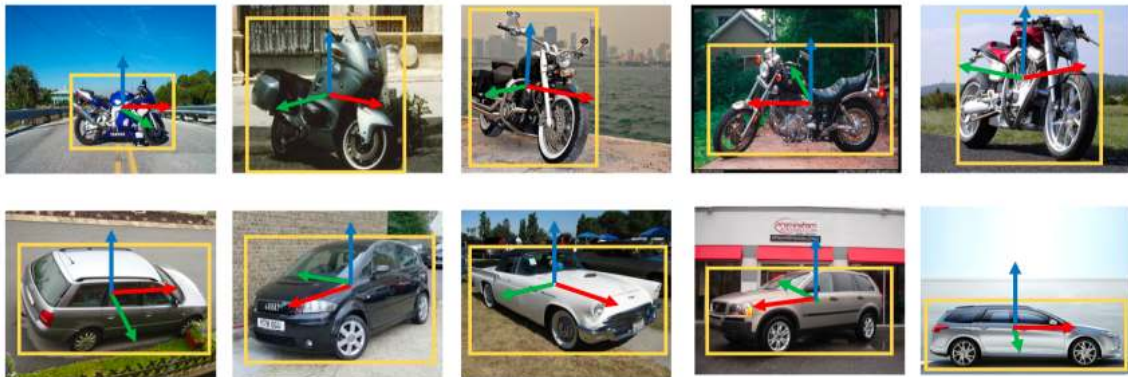


(b) Two possible shape encoders

Pose from Shape: Deep Pose Estimation for Arbitrary 3D Objects. Xiao et al. BMVC 2019.

few-shot learning for 3D scene understanding

Training Examples (Novel Classes)



Optional

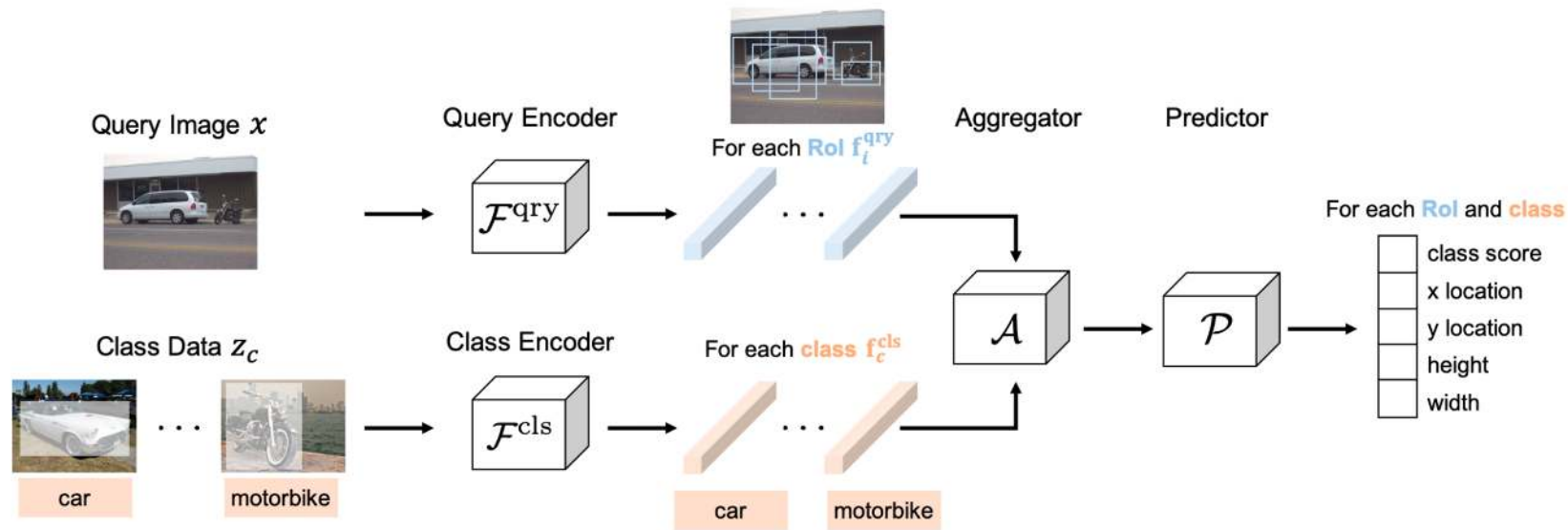


Testing



Few-shot Object Detection and Viewpoint Estimation for Objects in the Wild. Yang Xiao, Vincent Lepetit, Renaud Marlet. arXiv 2020.

few-shot learning for 3D scene understanding

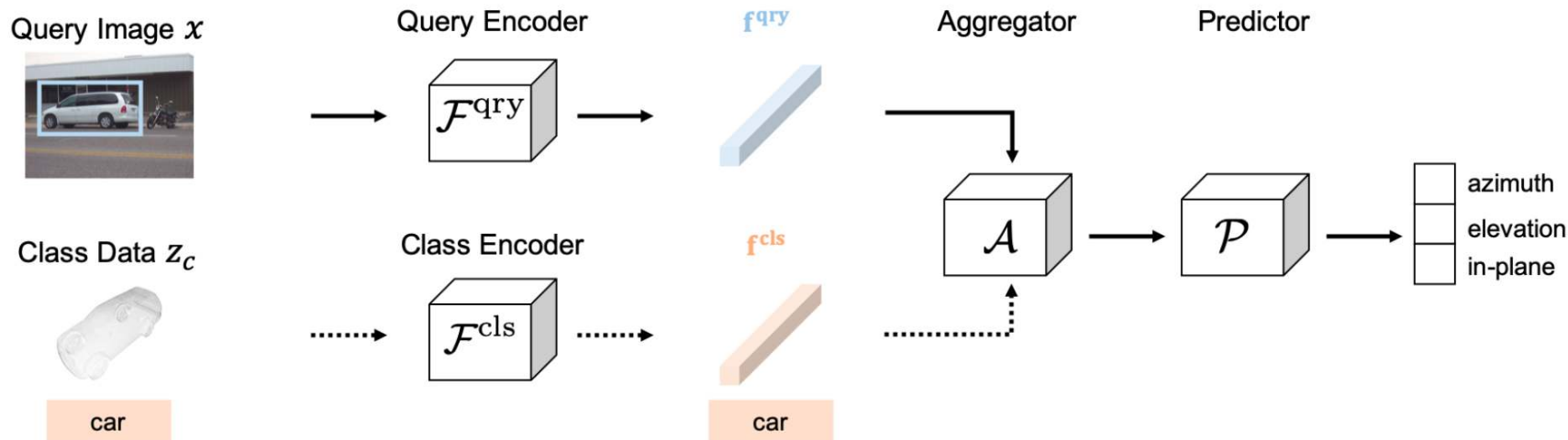


$$\text{cls}_{i,c} = \frac{\alpha \mathcal{A}(f_i^{\text{qry}}, f_c^{\text{cls}})^\top \mathbf{w}_c}{\|\mathcal{A}(f_i^{\text{qry}}, f_c^{\text{cls}})\| \|\mathbf{w}_c\|}$$

$$\mathcal{A}(f^{\text{qry}}, f^{\text{cls}}) = f^{\text{qry}} \otimes f^{\text{cls}}$$

Few-shot Object Detection and Viewpoint Estimation for Objects in the Wild. Yang Xiao, Vincent Lepetit, Renaud Marlet. arXiv 2021.

few-shot learning for 3D scene understanding



$$(\text{azi}, \text{ele}, \text{inp}) = \mathcal{P}(\mathcal{A}(f^{\text{qry}}, f^{\text{cls}}))$$

with $f^{\text{qry}} = \mathcal{F}^{\text{qry}}(\text{crop}(\text{img}(x), \text{box}(x)))$, and
 $f^{\text{cls}} = \mathcal{F}^{\text{cls}}(z_c)$, $c = \text{cls}(x)$