# Joint Object Detection and Viewpoint Estimation using CNN features

**Carlos Guindel**, David Martín and José M. Armingol

cguindel@ing.uc3m.es

Intelligent Systems Laboratory · Universidad Carlos III de Madrid

Wien · 28 June 2017

IEEE
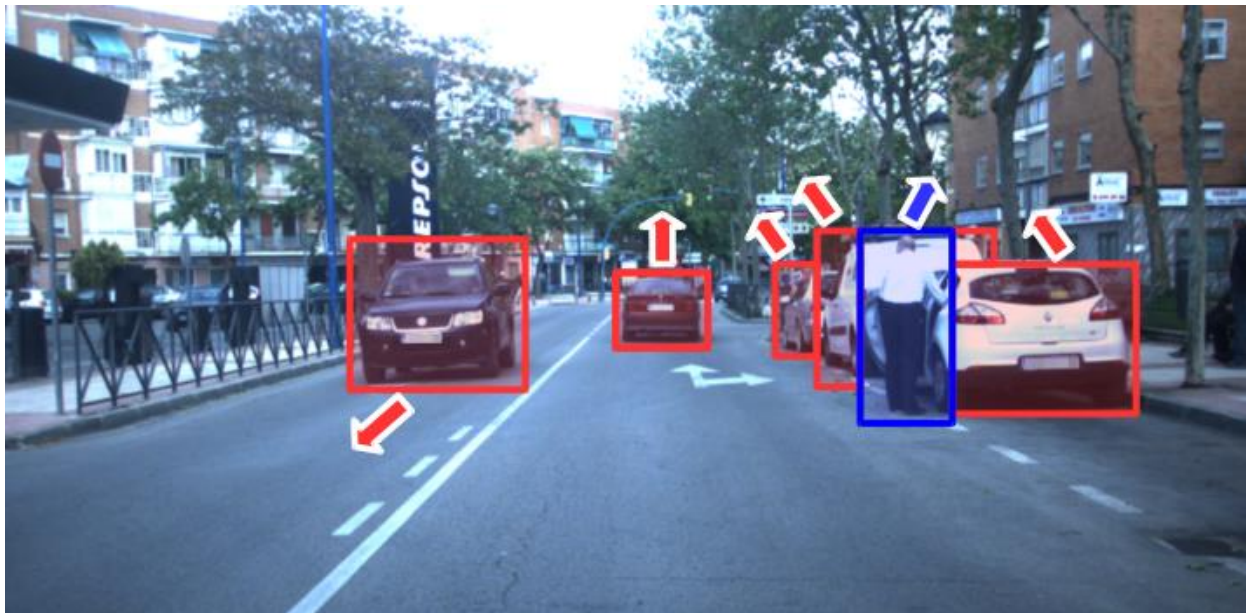ICVES 2017

Intelligent Systems Lab
www.uc3m.es/islab

uc3m

# Outline

- Introduction

- Object detection

- Viewpoint estimation

- Results

- Conclusion

# Outline

- **Introduction**

- Object detection

- Viewpoint estimation

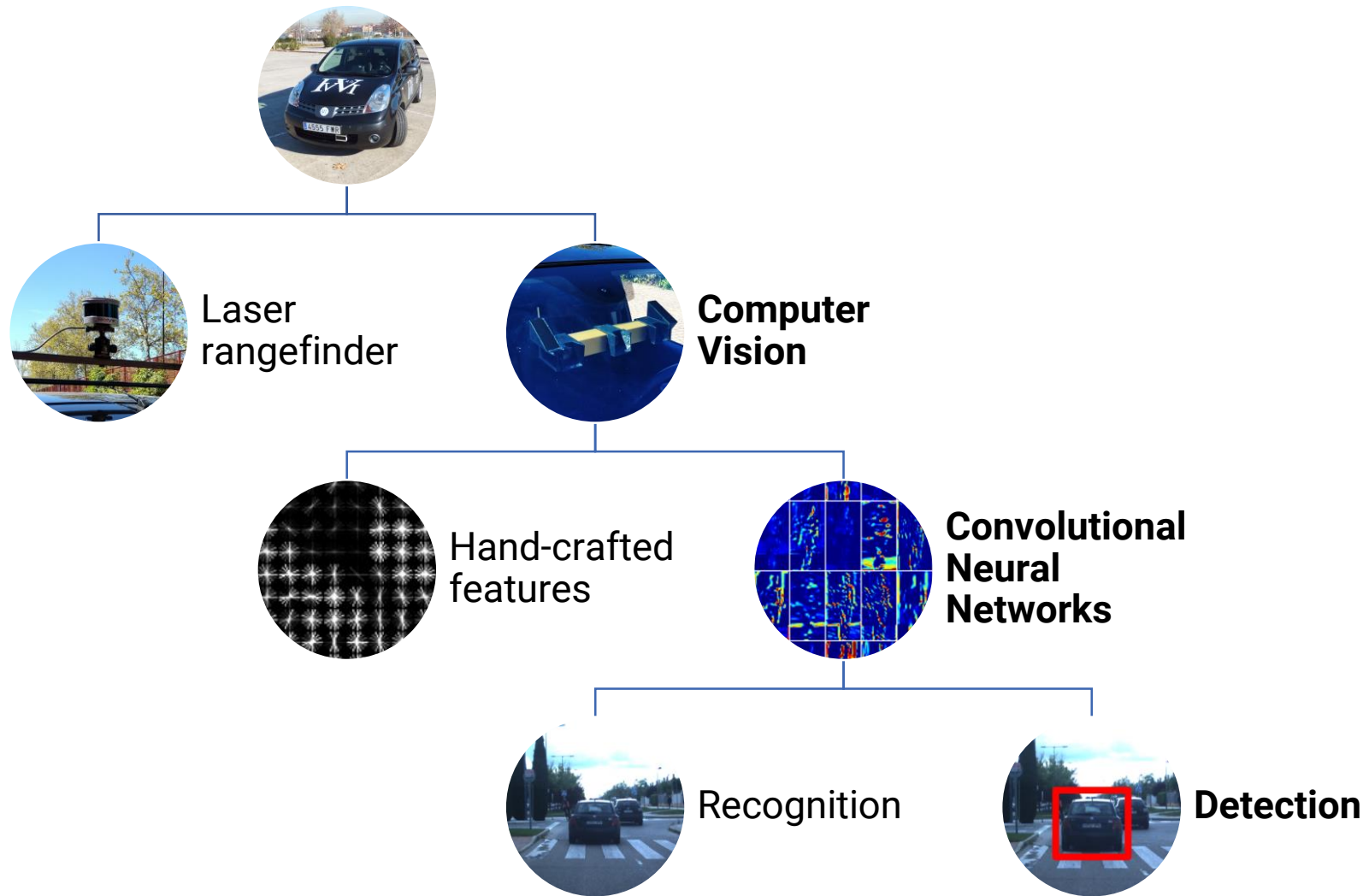- Results

- Conclusion

# Situational awareness for vehicles

- Advanced Driver Assistance Systems (ADAS) and autonomous vehicles rely on a trustable on-board **obstacle detection** module.

- A precise **classification** of the obstacles enables accurate predictions of future traffic situations, including those involving VRU.

- Another significant cue that can be used to anticipate future events is the **orientation** of the objects moving on the ground plane.
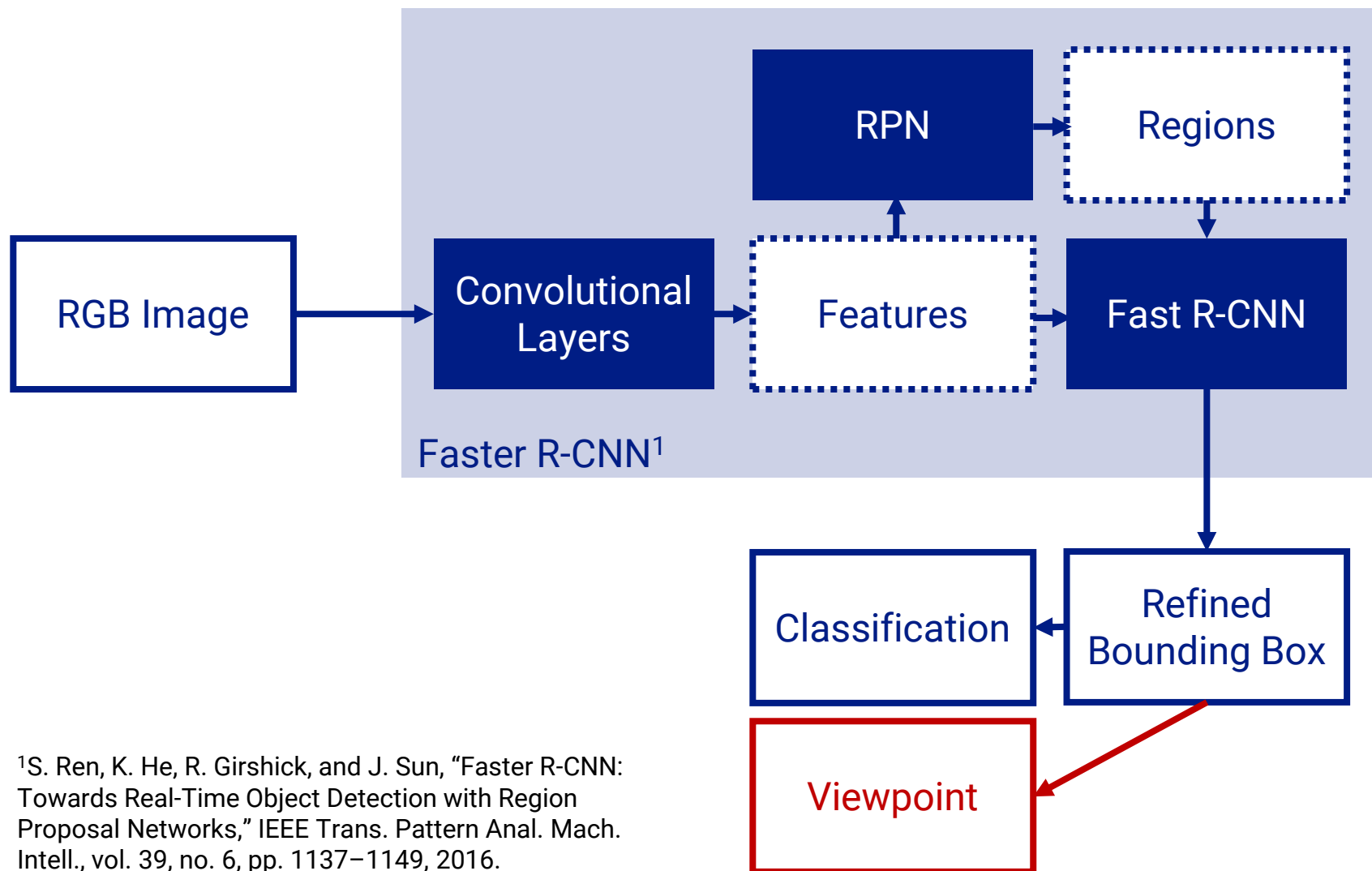
# On-board object detection

Laser rangefinder

**Computer Vision**

Hand-crafted features

**Convolutional Neural Networks**

Recognition

**Detection**

uc3m

# Proposal overview

[1]S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 6, pp. 1137–1149, 2016.

# Outline

- Introduction

- **Object detection**

- Viewpoint estimation

- Results

- Conclusion

# Object detection

- **Traffic environments**

  - Diversity of agents

  - Unstructured environment

- **Faster R-CNN**

  - End-to-end feature learning

  - Highly efficient

  - No prior constraints about the location of objects in the image

  - Meant for more than 21 classes

S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 6, pp. 1137–1149, 2016.
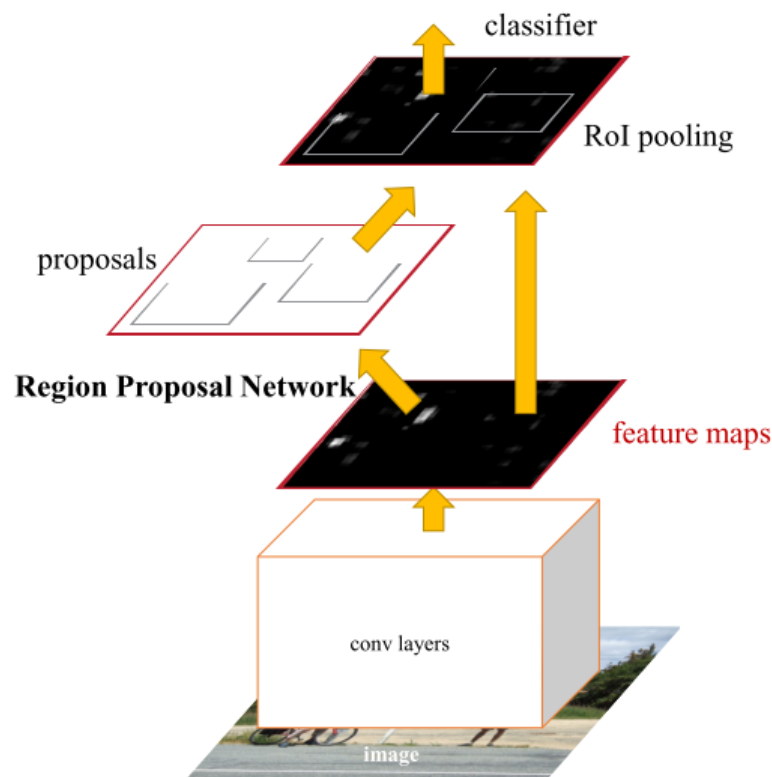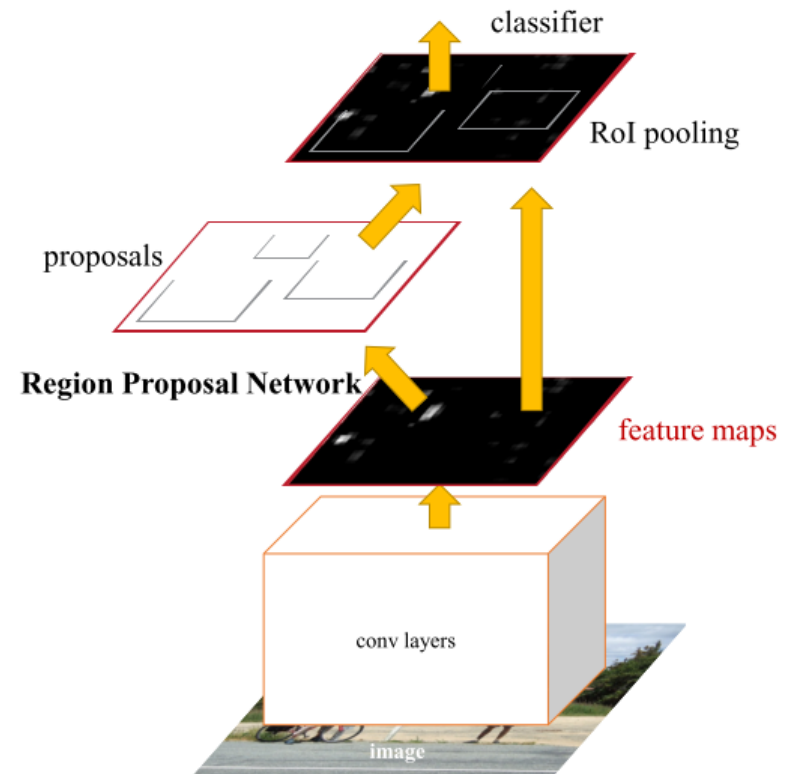
# Faster R-CNN framework

**Parameters are learned through a multi-task loss**

**Conv. features in these regions are pooled for classification**

**A RPN generates proposals wrt. a fixed set of anchors**

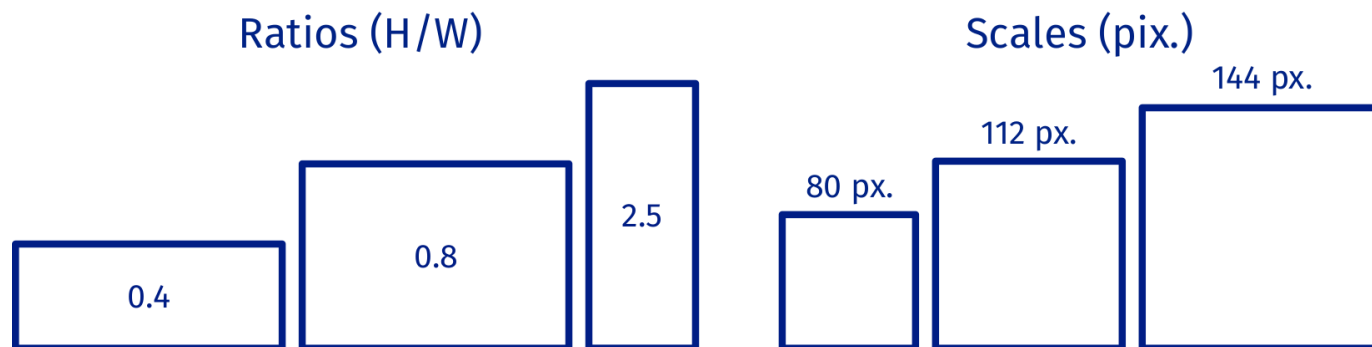**Convolutional features computed only once per image**



S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 6, pp. 1137–1149, 2016.

# Fine-tuning for traffic environments

- Optimized anchors

Ratios (H/W)

0.4

0.8

2.5

Scales (pix.)

80 px.

112 px.

144 px.

- Management of class imbalance
  - Information gain multinomial logistic loss for the **class** inference

$$L = L = \frac{1}{N} \sum_{n=1}^{N} \sum_{k=1}^{K} H_{l_n,k_n} \log(\hat{m}_{n,k_n})$$

Infogain matrix

$$\begin{pmatrix} H_{0,0} & \cdots & 0 \\ 0 & & \vdots \\ \vdots & \ddots & 0 \\ 0 & \cdots & H_{K,K} \end{pmatrix}$$
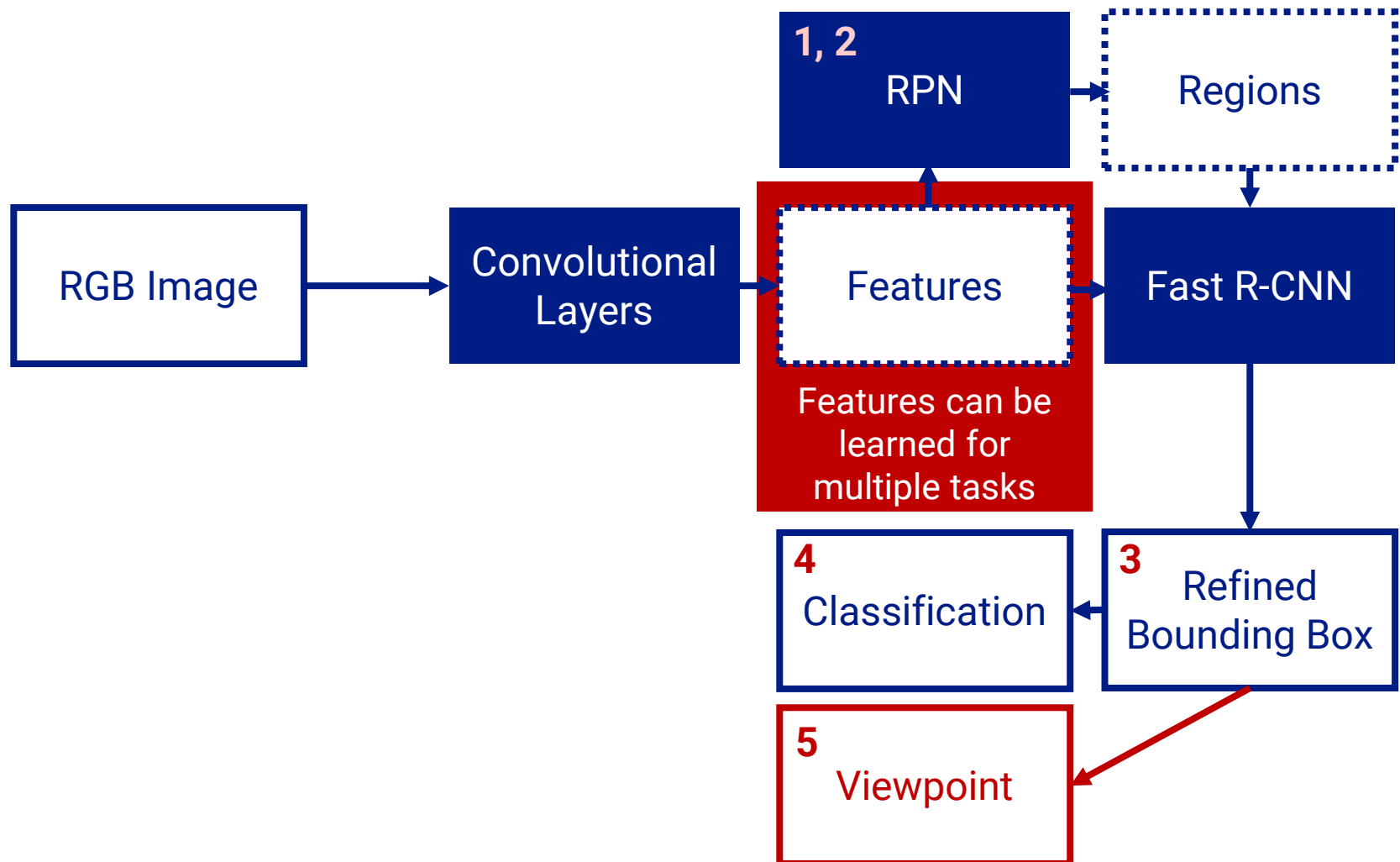
uc3m

# Outline

- Introduction

- Object detection

- **Viewpoint estimation**

- Results

- Conclusion

# Viewpoint estimation

# Discrete viewpoint inference

$N_b$ **angle bins** $\Theta_i \dots \Theta_{N_b}$

$N_b = 8$



$\theta_{i_0}$

**Training:** $\theta_{i_0} \to \Theta_i$

$$\Theta_i = \left\{ \theta \in [0, 2\pi) \;\middle|\; \frac{2\pi}{N_b} \cdot i \leq \theta < \frac{2\pi}{N_b} \cdot (i+1) \right\}$$

**Inference output:** $r \in \Delta^{N_b - 1}$

$$\Delta^N = \left\{ x \in \mathbb{R}^{N+1} \;\middle|\; \sum_{i=1}^{N+1} x_i = 1 \;\wedge\; \forall i : x_i \geq 0 \right\}$$

$r$

$$i^* = \arg\max_i (x_i)$$

Elements of $r$

**Final estimation:** $\Theta_{i^*} \to \hat{\theta}$

$$\hat{\theta} = \frac{\pi(2i^* + 1)}{N_b}$$

# Joint detection and viewpoint estimation

**CNN outputs**

**RPN**

**Fast R-CNN**

For each anchor:

- Objectness

$$a \in \{0, 1\}$$

- Predicted bounding box

$$b = (b_x, b_y, b_w, b_h)$$

For each proposal:

- Class

$$p = (p_0, \ldots, p_K)$$

- Bounding box refinement

$$t^k = (t_x^k, t_y^k, t_w^k, t_h^k) \text{ for } k = 0, \ldots, K$$

Per class

- Viewpoint

$$r^k = (r_0^k, \ldots, r_{N_b}^k) \text{ for } k = 0, \ldots, K$$

Number of angle bins

Number of classes

$N_b \cdot K$ elements

$N_b$

$\cdot K$

uc3m

# Joint detection and viewpoint estimation

**Fast R-CNN**

| B. Box regression | Class | Viewpoint |
| --- | --- | --- |
| Softmax | Softmax | Softmax |
| FC layer | FC layer | FC layer |

Fully connected (FC) layers

Only $N_b \cdot K \cdot 4096$ new weights

Fixed size feat. vector

Conv. design is unchanged

Proposal

Feature map

# Loss function and training

- **Approximate joint** training strategy

- Unweighted muli-task loss with **five** components

**Logistic loss**
for RPN objectness

**Smooth-L1 loss**
for RPN b.box regression

**Smooth-L1 loss**
for b.box regression

$$L = \frac{1}{N_{B_1}} \sum_{j \in B_1} L_{cls}(a_j, u_j) + \frac{1}{N_a} \sum_{j \in B_1} u_j L_{loc}(b_j, b_j^*) +$$

$$\frac{1}{N_{B_2}} \sum_{i \in B_2} L_{inf}(p_i, v_i) + \sum_{i \in B_2} [u \geq 1] L_{loc}(t_i^v, t_i^{v*}) +$$

$$\frac{1}{N_{B_2}} \sum_{i \in B_2} [u \geq 1] L_{cls}(r_i^u, \Theta_i)$$

**Infogain loss**
for class

$$L_{inf}(p_i v_i) = \sum_{n=1}^{N} H_{v_i,k} \log(p_{i,k})$$

$\downarrow$ frequent $\rightarrow \uparrow H_{v_i,k}$

**Logistic loss**
for viewpoint estimation

Only the $N_b$ elements of
the ground-truth class

uc3m

# Outline
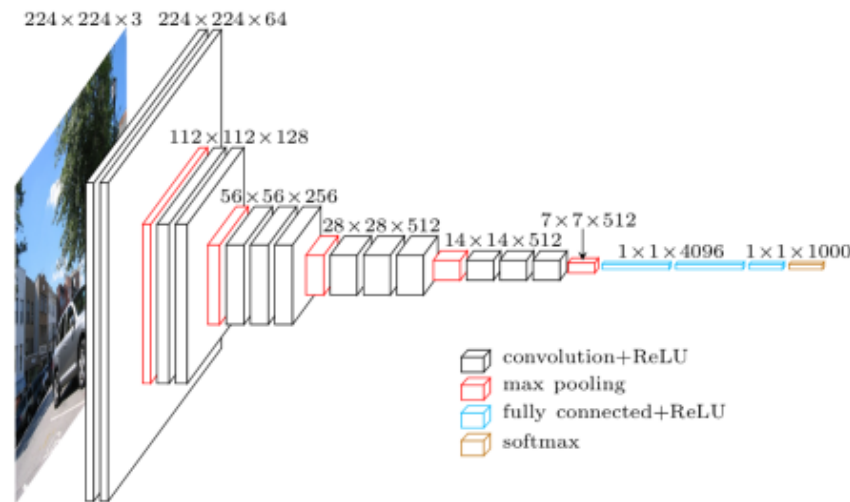
- Introduction

- Object detection

- Viewpoint estimation

- **Results**

- Conclusion

# Experiments

- On the KITTI Vision Benchmark Suite − Object detection dataset

- **Parameters:**

  - **Scale**: 500 px. in height

  - 50k iter. $l_r = 0.001$ **+** 50k iter. $l_r = 0.0001$ **+** 50k iter. $l_r = 0.00001$

  - **VGG16** architecture, initialized with ImageNet weights.



Source: https://blog.heuritech.com/2016/02/29/a-brief-report-of-the-heuritech-deep-learning-meetup-5/

KITTI: A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite," in Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012, pp. 3354–3361.

# Experiments

- $N_b = 8$ (resolution: $\pi/4$ rad)

- Infogain matrix values:

Number of ocurrences of the less frequent class

$$H_{k,k} = 2 \cdot \left(\frac{f_{min}}{f_k}\right)^{\frac{1}{8}}$$

Number of instances of class $k$

**Evaluation criteria**

- Average precision

- Average orientation similarity (performance of **detection + orientation**)

$$AOS = \frac{1}{11} \sum_{r \in \{0, 0.1, .., 1\}} \max_{\tilde{r}: \tilde{r} \geq r} s(\tilde{r}) \qquad s(r) = \frac{1}{|\mathcal{D}(r)|} \sum_{i \in \mathcal{D}(r)} \frac{1 + \cos \Delta_\theta^{(i)}}{2} \delta_i$$

- Minimum overlaps established by the KITTI benchmark

uc3m

# Results (KITTI submission: FRCNN+Or)

| Detection (AP as %) | Easy | Moderate | Hard |
|---|---|---|---|
| **Car** | 89.60 | 78.59 | 68.69 |
| **Pedestrian** | 72.21 | 56.99 | 53.72 |
| **Cyclist** | 68.81 | 55.80 | 50.52 |
| **mAP** | **76.87** | **63.79** | **57,64** |
| *SubCNN* | *84.52* | *77.14* | *64.44* |



The KITTI Vision Benchmark Suite
A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago

home  setup  stereo  flow  scene flow  odometry  object  tracking  road  semantics  raw data  submit results  jobs

Andreas Geiger (MPI Tübingen) | Philip Lenz (KIT) | Christoph Stiller (KIT) | Raquel Urtasun (University of Toronto)

Method

Joint Object Detection and Viewpoint Estimation using CNN features [FRCNN+Or]

Submitted on 7 Jun. 2017 10:58 by
Carlos Guindel (Universidad Carlos III de Madrid)

Running time: 0.1 s
Environment: GPU @ 1.5 Ghz (Python + C/C++)

Method Description:
We enhance the well-known Faster R-CNN method by adding viewpoint inference capabilities. Convolutional features are computed and shared for use in the tasks of region proposal, classification, and orientation estimation, for efficiency reasons.
Parameters:

- Slightly different (generally better) than the ones in the paper:
  - Used the whole KITTI training set
  - Trained only with Car, Pedestrian and Cyclist
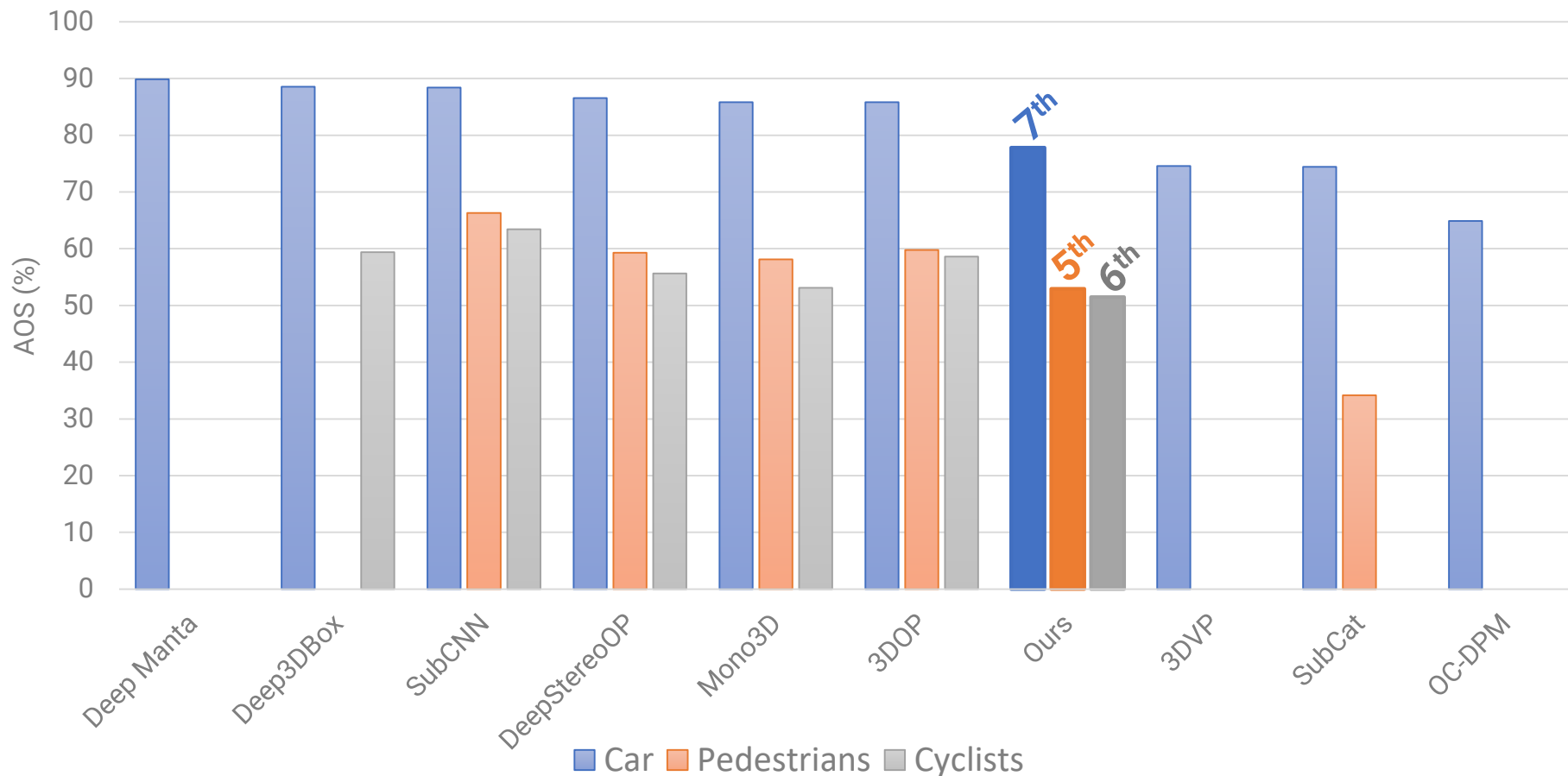  - Non-fixed weights (and bias) at the first convolutional layers

| Det + Or (AOS as %) | Easy | Moderate | Hard |
|---|---|---|---|
| **Car** | 88.93 | 77.8 | 67.87 |
| **Pedestrian** | 67.92 | 52.96 | 49.61 |
| **Cyclist** | 64.90 | 51.47 | 46.48 |
| **mAOS** | **73.92** | **60.74** | **54,65** |
| *SubCNN* | *80.37* | *72.85* | *65.45* |

SubCNN: Y. Xiang, W. Choi, Y. Lin and S. Savarese, "Subcategory-aware Convolutional Neural Networks for Object Proposals and Detection," in IEEE Winter Conference on Applications of Computer Vision (WACV), 2017, pp. 924-933.
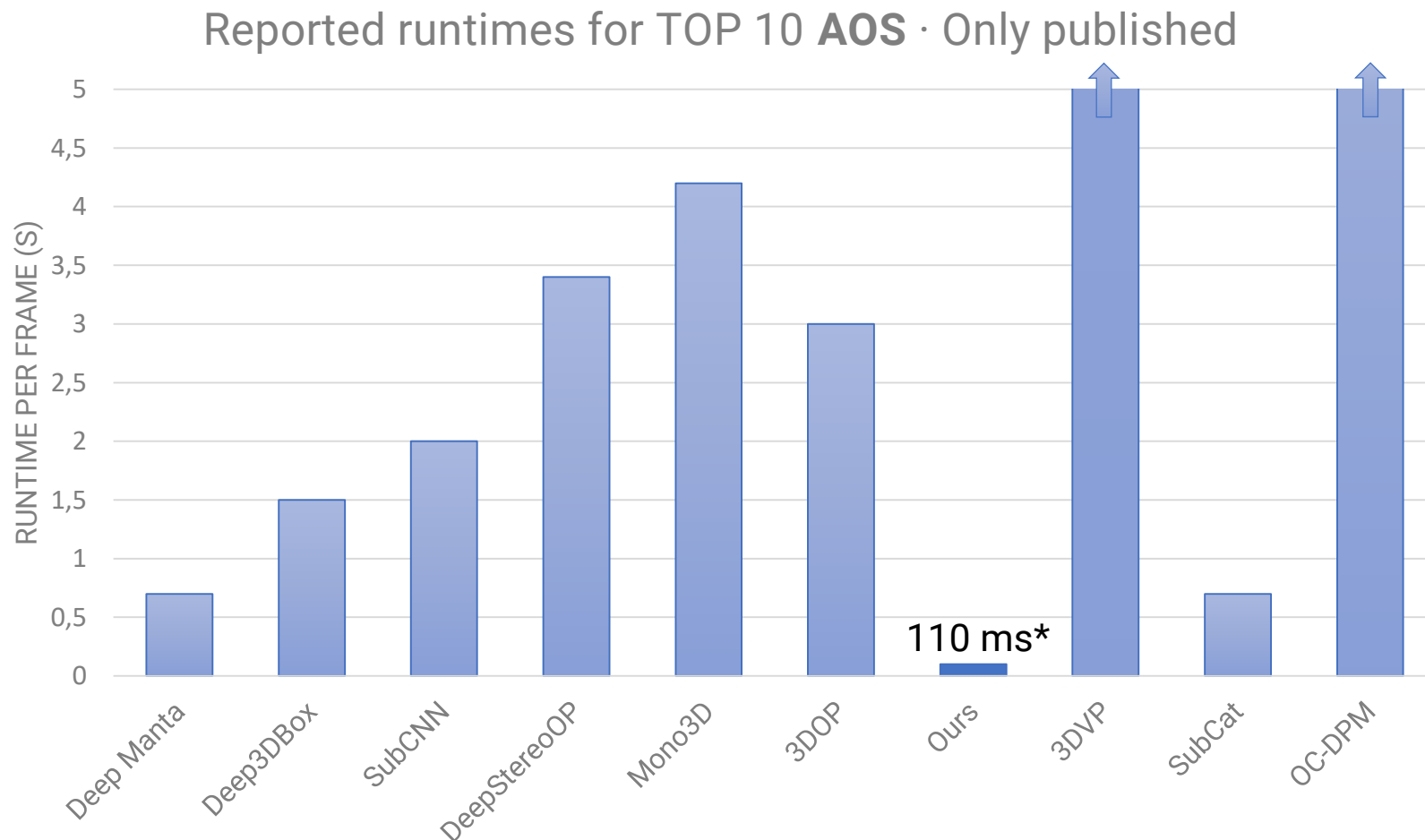
# Comparison with other methods

TOP 10 **AOS** ranking · MODERATE difficulty · Only published methods

Global rankings (including unpublished methods): Car: 11st · Pedestrian: 9th · Cyclist: 10th

# Comparison with other methods
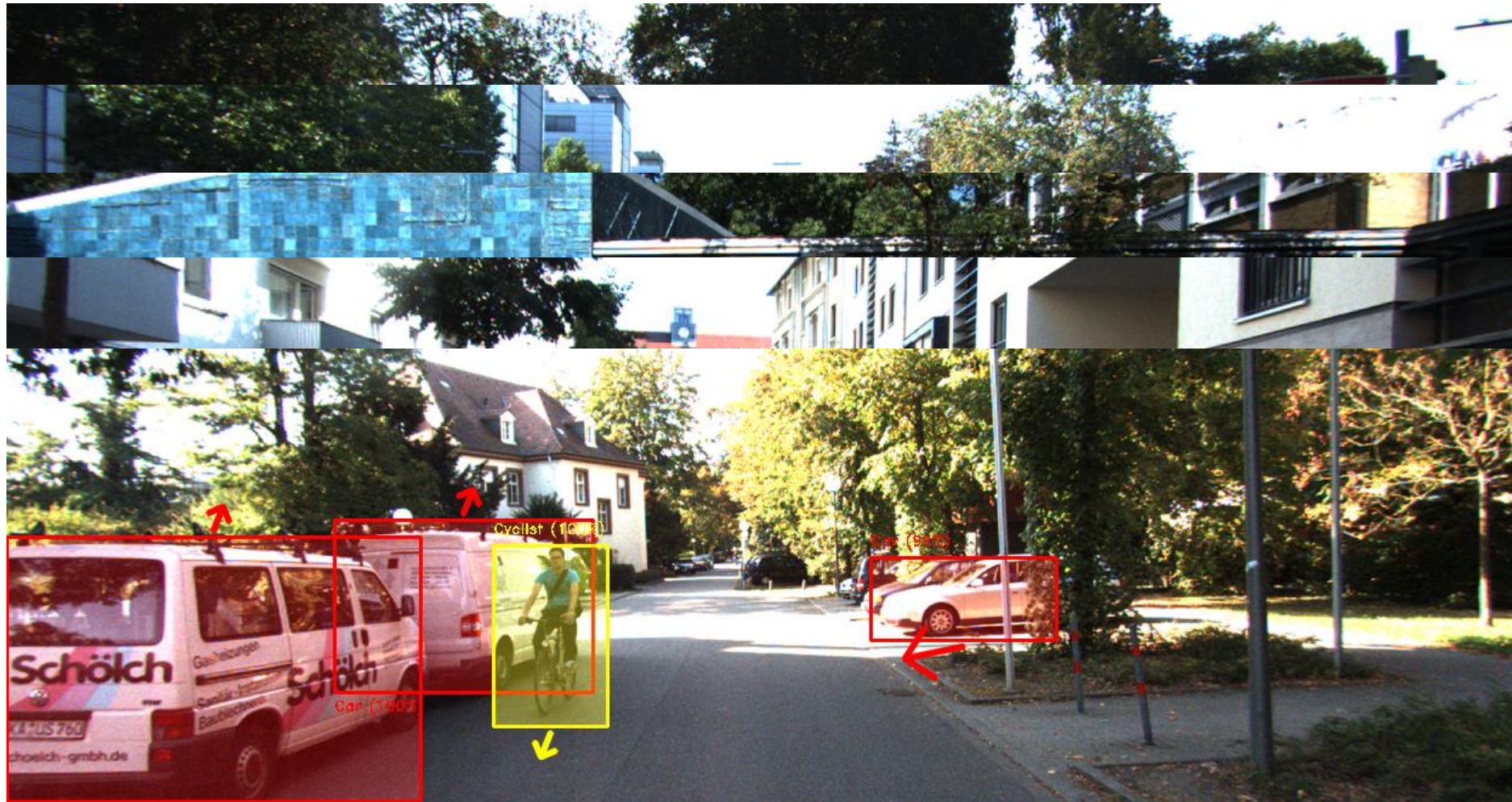
Reported runtimes for TOP 10 **AOS** · Only published

*Average running time using an implementation based on py-faster-rcnn (Python & Caffe) and a NVIDIA **Titan Xp** kindly donated by NVIDIA Corporation

# Sample test images

# Outline

- Introduction

- Object detection

- Viewpoint estimation

- Results

- **Conclusion**

# Conclusion

- Single-image object detection + orientation estimation in traffic scenes.

- The same convolutional features can be successfully used for both tasks.

- Results comparable with non-real-time, sophisticated approaches.

- Orientation is a step towards a real scene understanding.

## Future work

- Fine-grained orientation inference using the **cross-entropy logistic loss**

$$L = -\frac{1}{n} \sum_{n=1}^{n} [p_n \log \hat{p}_n + (1 - p_n)\log(1 - \hat{p}_n)]$$

- Improvements:
  - Network architecture
  - Methods to overcome the fixed-size receptive field

Continuously-evolving code at: https://github.com/cguindel/lsi-faster-rcnn

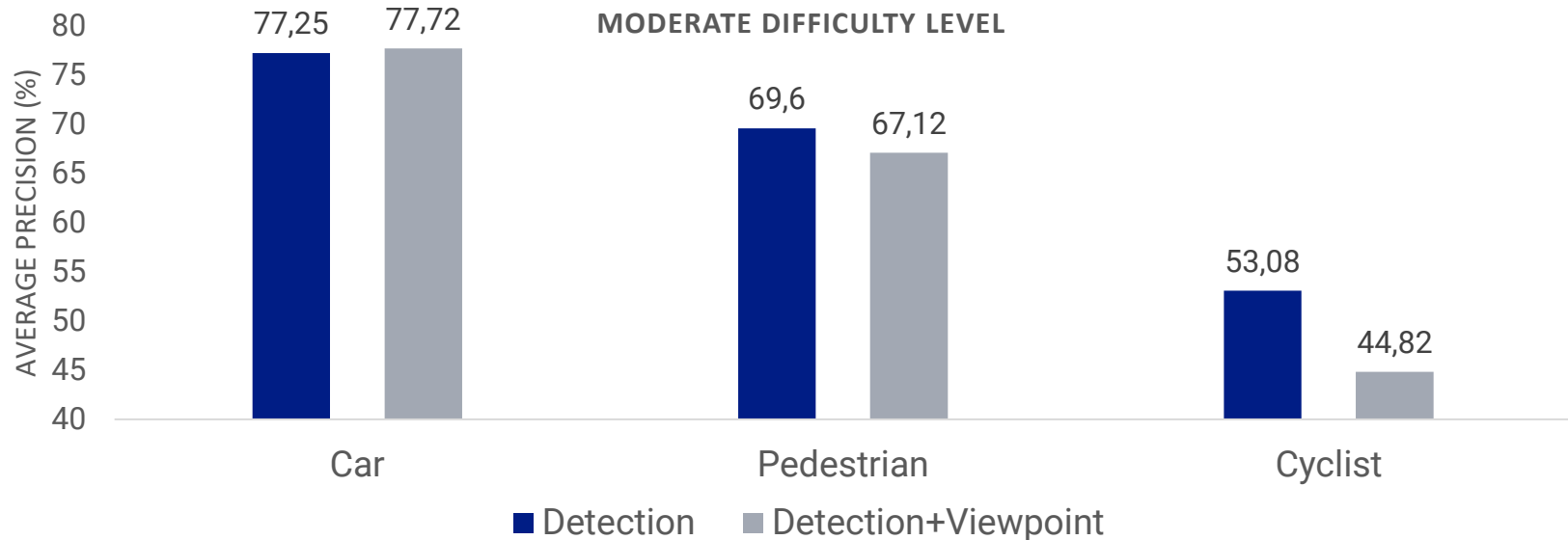# Thank you for your attention!

**Carlos Guindel**

cguindel@ing.uc3m.es

Intelligent Systems Laboratory · Universidad Carlos III de Madrid

Wien · 28 June 2017

# Precision change when introducing viewpoint



**MODERATE DIFFICULTY LEVEL**

Chart showing Average Precision (%):
- Car: Detection 77,25 — Detection+Viewpoint 77,72
- Pedestrian: Detection 69,6 — Detection+Viewpoint 67,12
- Cyclist: Detection 53,08 — Detection+Viewpoint 44,82

Legend: ■ Detection   ■ Detection+Viewpoint

| Detection (ΔAP as %) | Easy | Moderate | Hard |
|---|---|---|---|
| **Car** | +0,72 | +0,47 | +0,26 |
| **Pedestrian** | -1,85 | -2,47 | -3,00 |
| **Cyclist** | -12,16 | -8,25 | -8,11 |