

# Joint Object Detection and Viewpoint Estimation using CNN features

VII LSI PhD Workshop

---

**Carlos Guindel**

cguindel@ing.uc3m.es

Intelligent Systems Laboratory · Universidad Carlos III de Madrid  
Leganés · 19 June 2017

# Joint Object Detection and Viewpoint Estimation using CNN features

C. Guindel, D. Martín, and J. M. Armingol, “Joint Object Detection and Viewpoint Estimation using CNN features,” in IEEE International Conference on Vehicular Electronics and Safety (ICVES), 2017.

# Outline

---

4

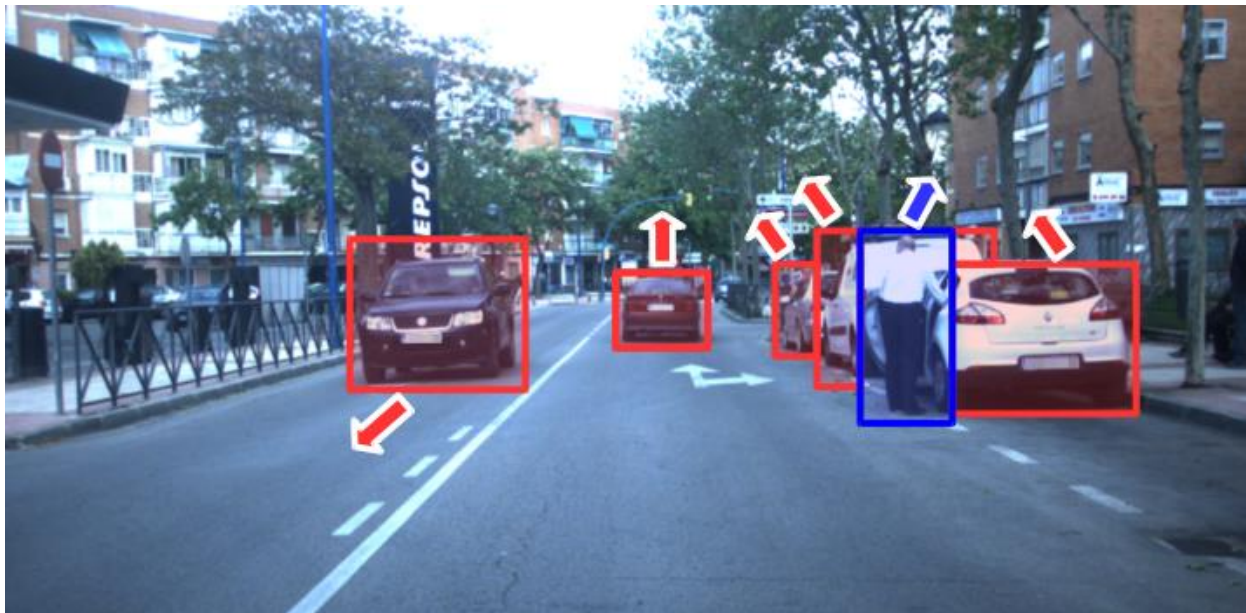
- Introduction
- Object detection
- Viewpoint estimation
- Results
- Conclusion

- **Introduction**
- Object detection
- Viewpoint estimation
- Results
- Conclusion

# Situational awareness for vehicles

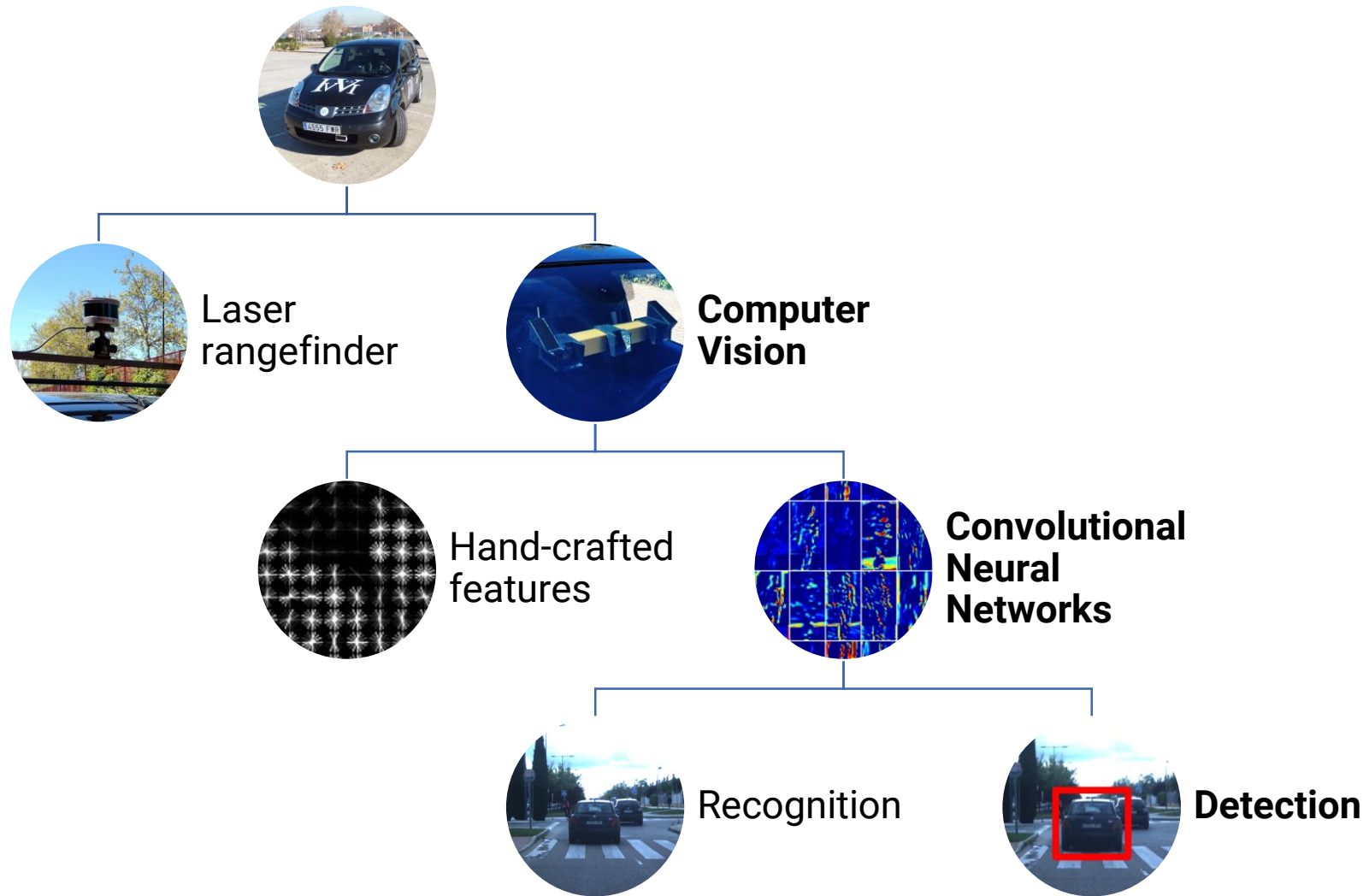
6

- 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.



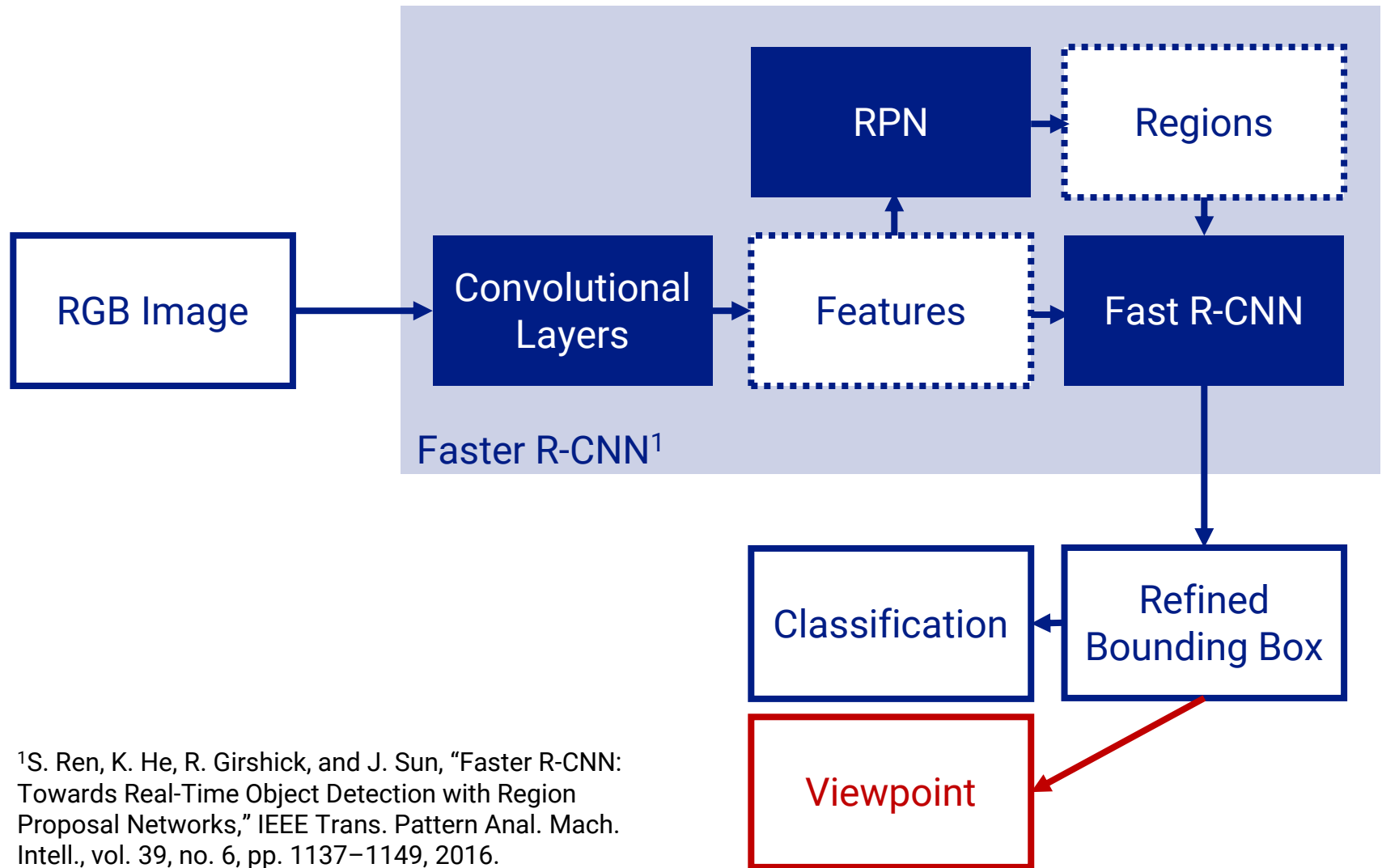
# On-board object detection

7



# Proposal overview

8



<sup>1</sup>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

---

9

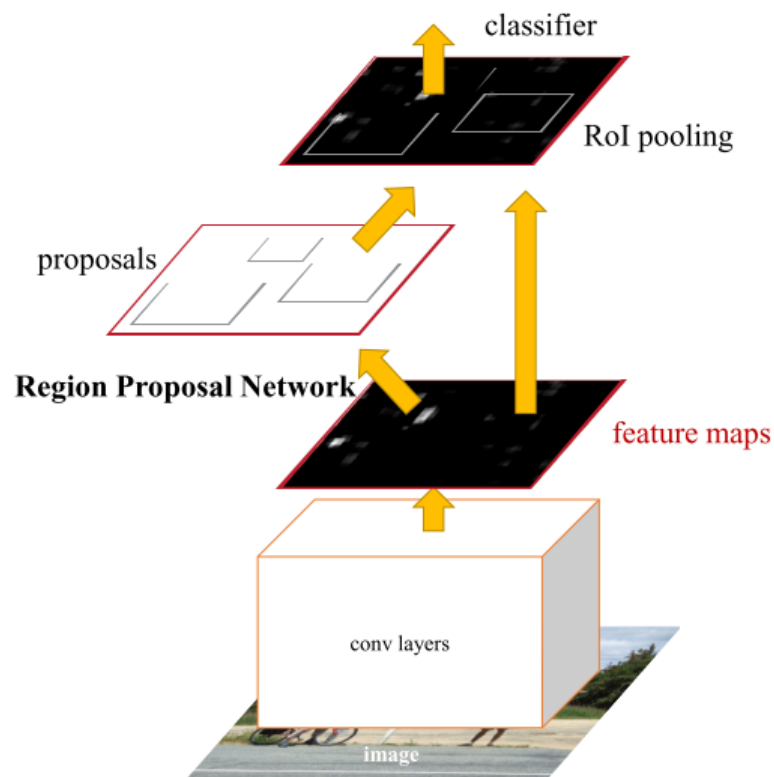
- Introduction
- **Object detection**
- Viewpoint estimation
- Results
- Conclusion



# Object detection

10

- 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

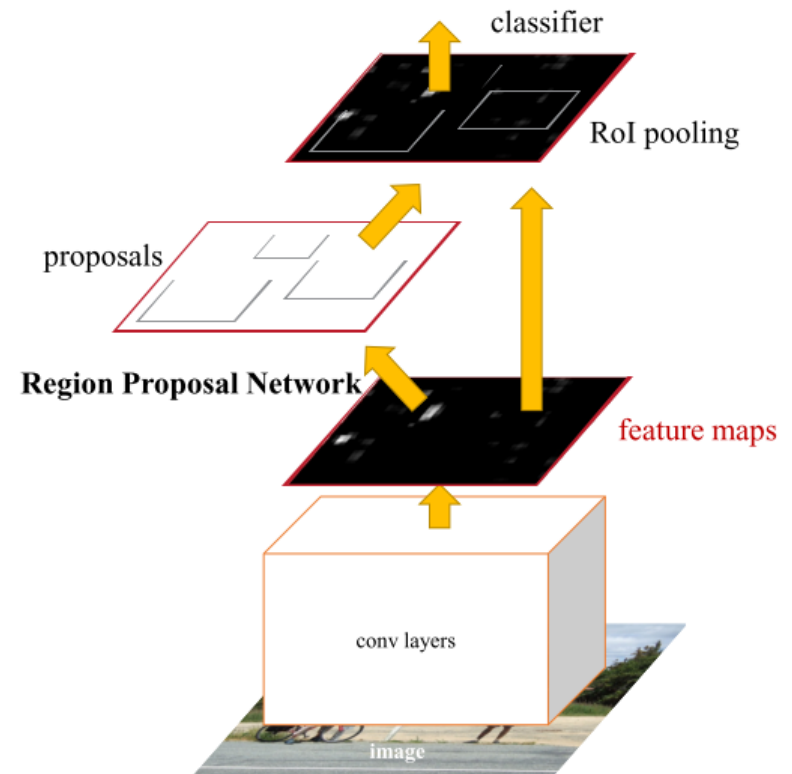
11

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

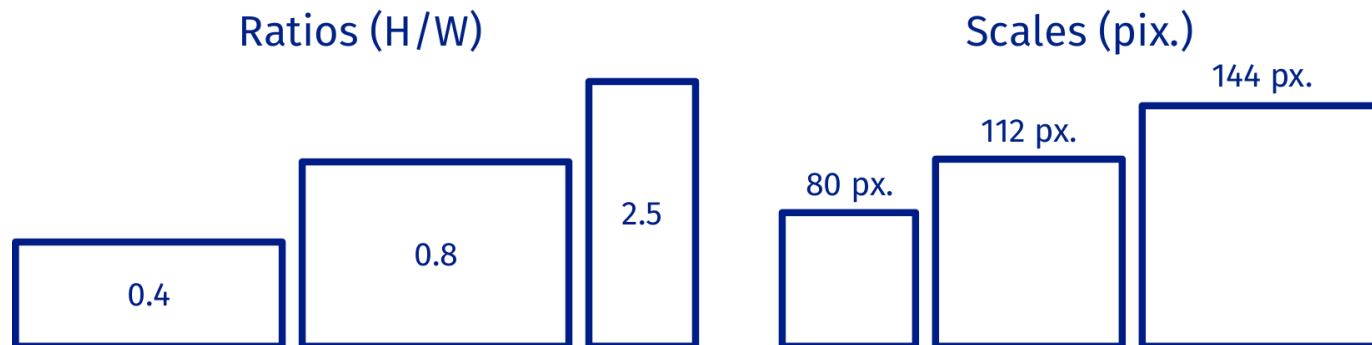


S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 6, pp. 1137–1149, 2016.

# Fine-tuning for traffic environments

12

- Optimized anchors



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

$$L = L_{\text{reg}} - \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K H_{l_n, k_n} \log(\hat{p}_{l_n, k_n})$$

↓

**Infogain matrix**

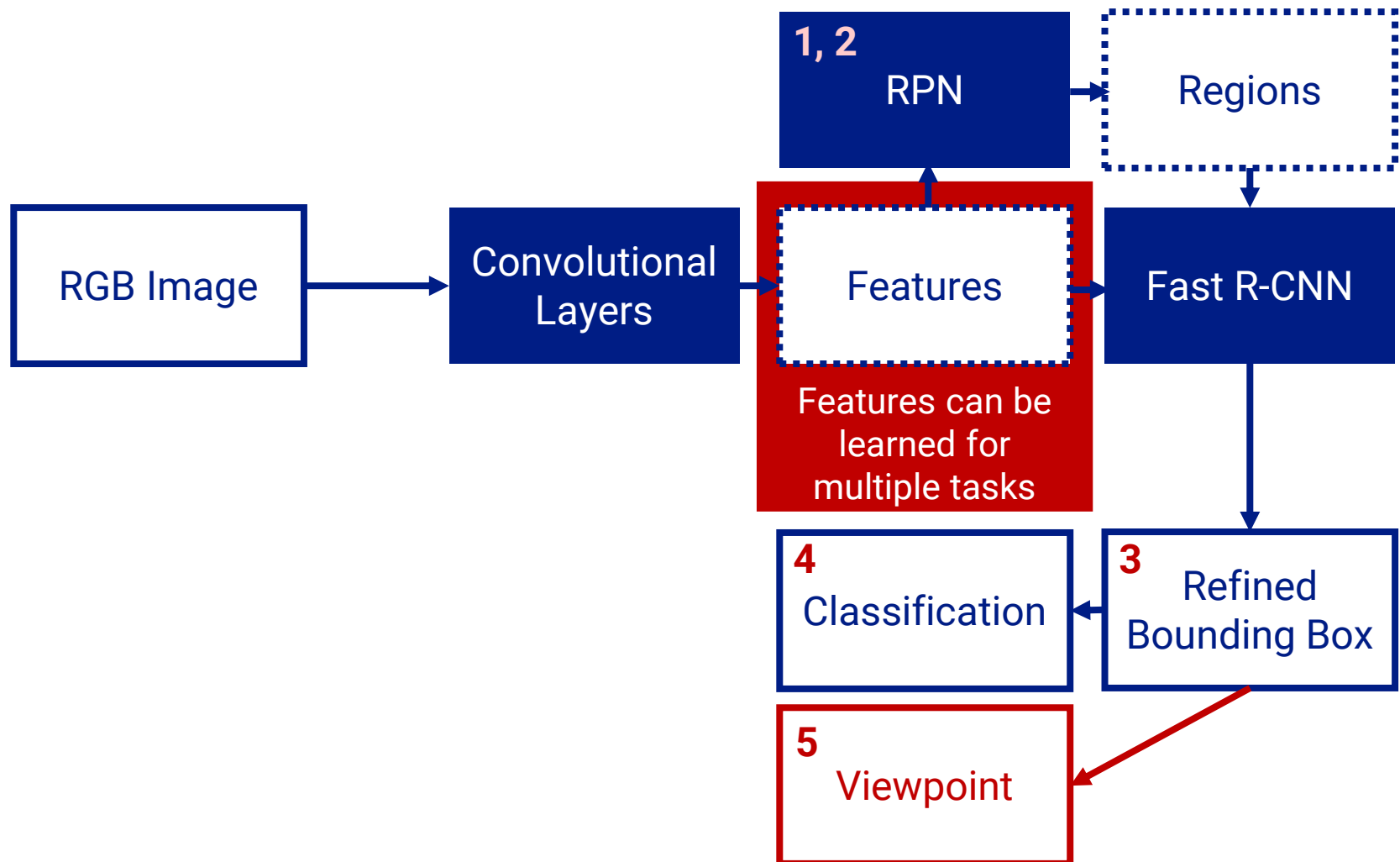
→

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

- Introduction
- Object detection
- **Viewpoint estimation**
- Results
- Conclusion

# Viewpoint estimation

14

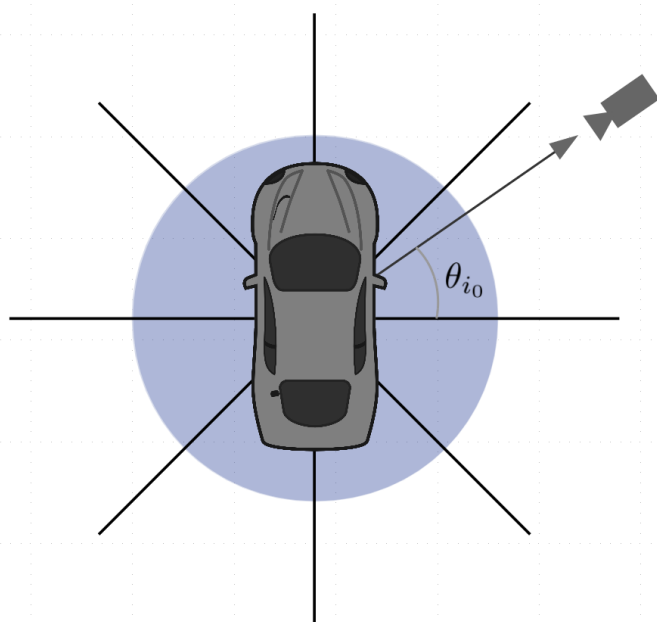


# Discrete viewpoint inference

15

$N_b$  angle bins  $\Theta_1 \dots \Theta_{N_b}$

$N_b = 8$



Training:  $\theta_{i_0} \rightarrow \Theta_i$

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

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

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



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

Elements  
of  $r$

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

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

# Joint detection and viewpoint estimation

16

CNN outputs

RPN

For each anchor:

- Objectness

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

- Predicted bounding box

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

Fast R-CNN

For each proposal:

- Class

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

- Bounding box refinement

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

Per class

- Viewpoint

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



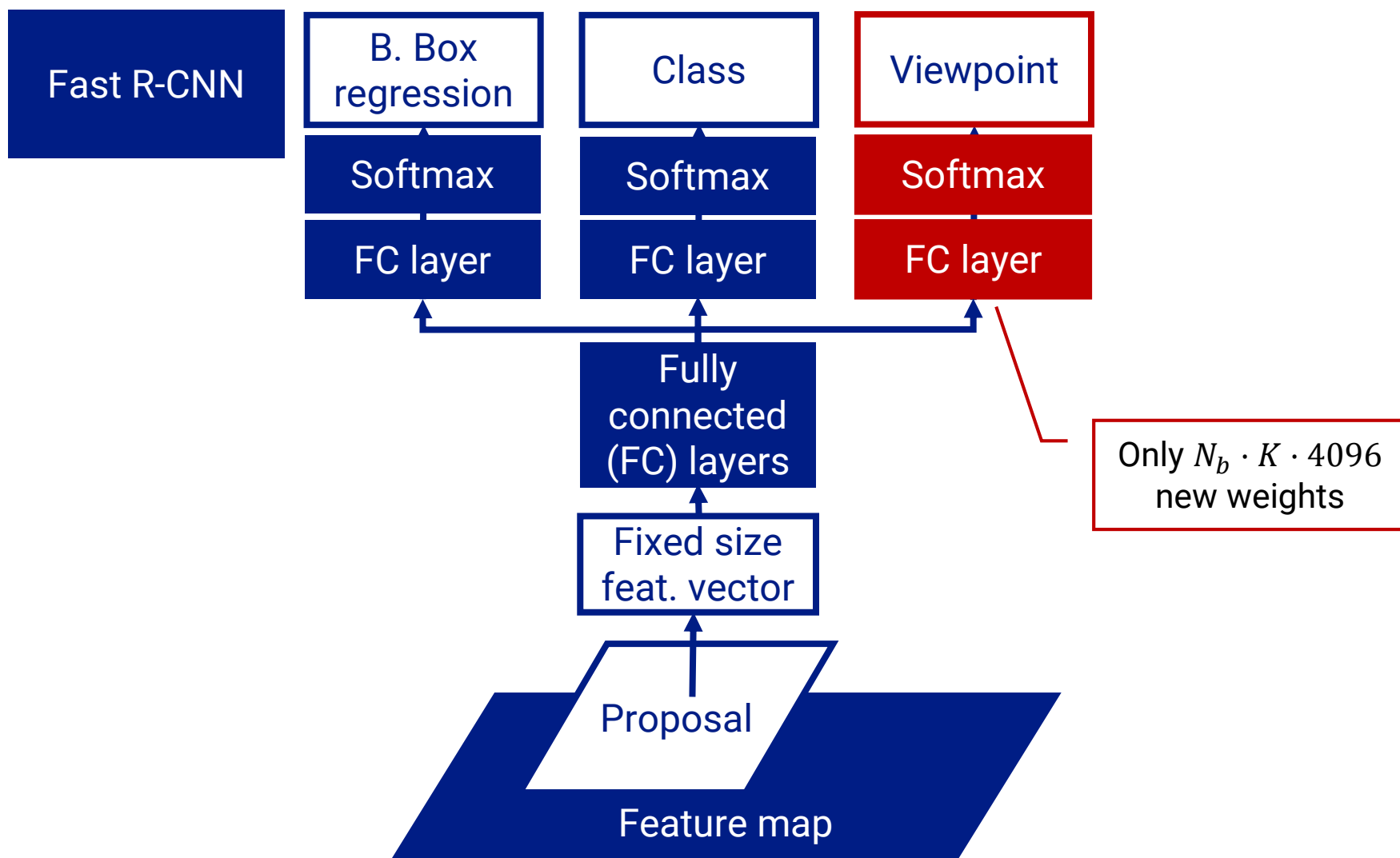
Number of  
angle bins

Number of  
classes

$\cdot K$

# Joint detection and viewpoint estimation

17





# Loss function and training

18

- **Approximate joint** training strategy
- Unweighted multi-task loss with **five** components

**Logistic loss**  
for RPN objectness

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

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

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

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

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

**Infogain loss**  
for class

**Logistic loss**  
for viewpoint estimation

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

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

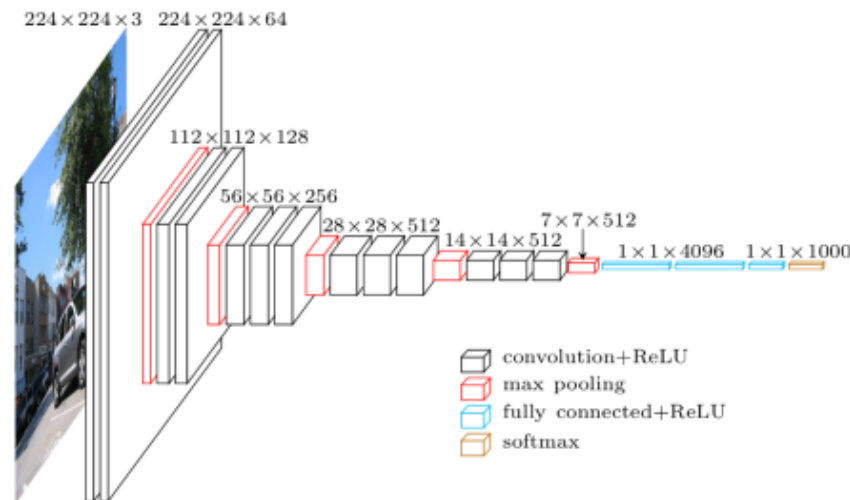
↓ frequent → ↑  $H_{v_i, k}$

- Introduction
- Object detection
- Viewpoint estimation
- **Results**
- Conclusion

# Experiments

20

- On the KITTI Vision Benchmark Suite – Object detection set
- **Training 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.



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.

VGG16 Image: <https://blog.heuritech.com/2016/02/29/a-brief-report-of-the-heuritech-deep-learning-meetup-5/>

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

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

Number of  
occurrences of the  
less frequent class

Number of  
instances of class  $k$




## Evaluation criteria

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

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

- Minimum overlaps established by the KITTI benchmark

**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:**  
Scale=500, 300 RPN proposals, 3 anchor scales and 3 anchor ratios,  $\lambda_N=8$ ,  $\lambda_r=0.001$ , 150k iterations. Running time given for a NVIDIA Titan X Pascal.

**Latex Bibtex:**  
@inproceedings{GuindelICVES,  
author = {Guindel, Carlos and Martin, David and Armingol, Jose M.},  
booktitle = {IEEE International Conference on Vehicular Electronics and Safety (ICVES)},  
title = {{Joint Object Detection and Viewpoint Estimation using CNN features}},  
year = {2017}  
}

# Results (KITTI submission: FRCNN+Or)

23

Detection (AP as %)	Easy	Moderate	Hard
<b>Car</b>	89.60	78.59	68.69
<b>Pedestrian</b>	72.21	56.99	53.72
<b>Cyclist</b>	68.81	55.80	50.52
<b>mAP</b>	<b>76.87</b>	<b>63.79</b>	<b>57.64</b>
<i>SubCNN</i>	84.52	77.14	64.44

- 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

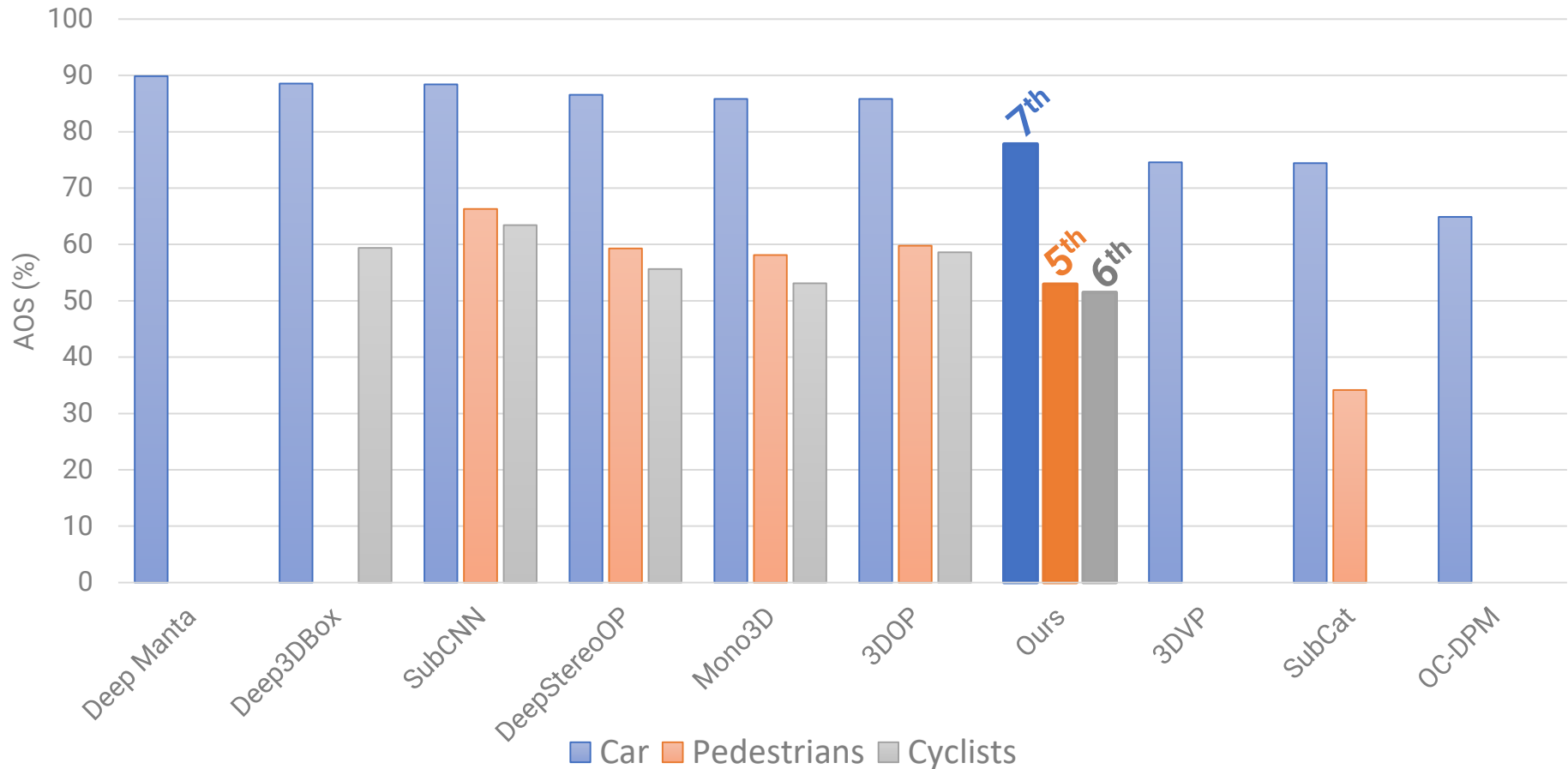
Det + Or (AOS as %)	Easy	Moderate	Hard
<b>Car</b>	88.93	77.8	67.87
<b>Pedestrian</b>	67.92	52.96	49.61
<b>Cyclist</b>	64.90	51.47	46.48
<b>mAOS</b>	<b>73.92</b>	<b>60.74</b>	<b>54.65</b>
<i>SubCNN</i>	80.37	72.85	65.45

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 published KITTI submissions

24

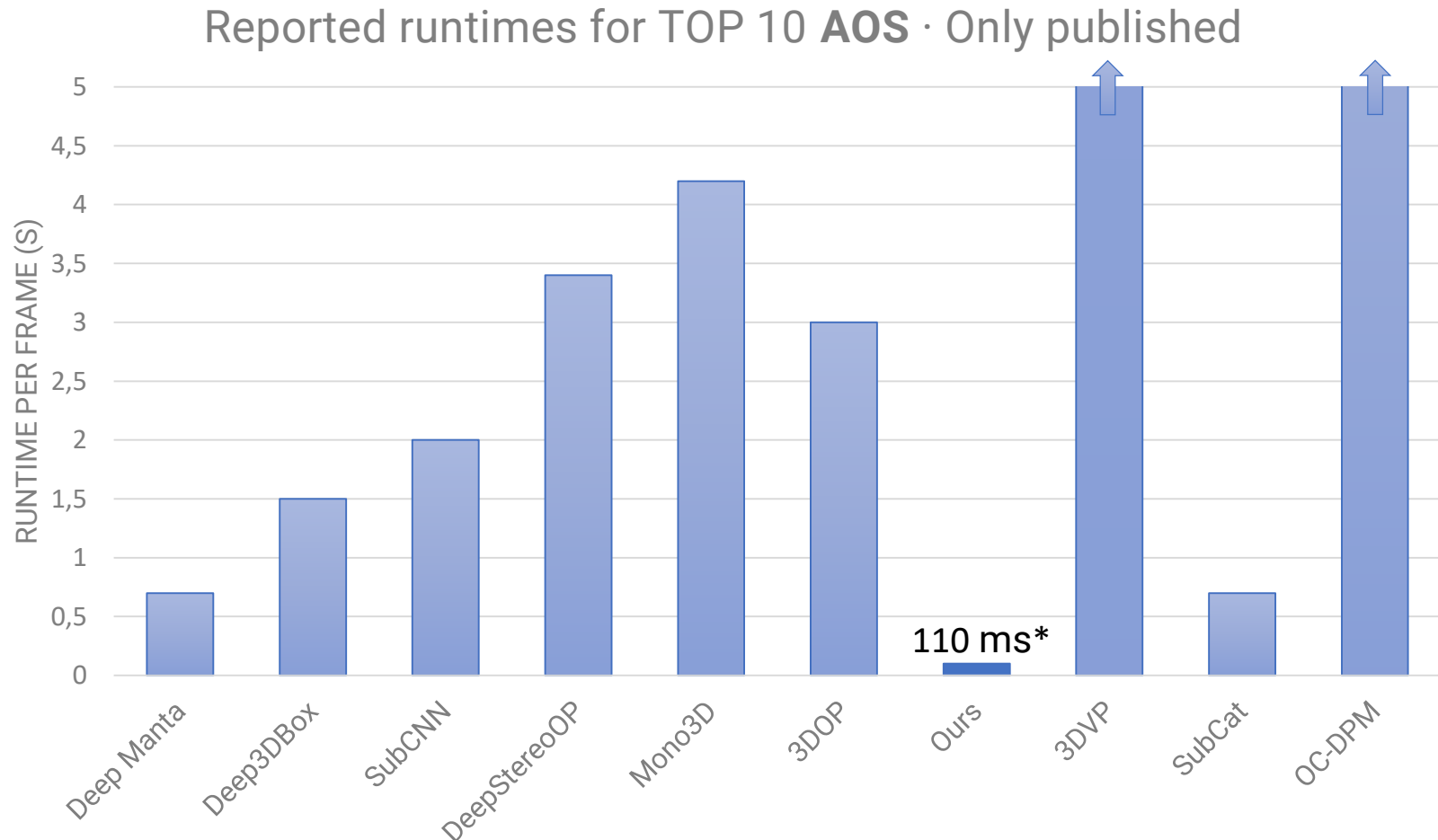
TOP 10 **AOS** ranking · MODERATE difficulty · Only published methods



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

# Comparison with published KITTI submissions

25



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



- Introduction
- Object detection
- Viewpoint estimation
- Results
- **Conclusion**

- Monocular approach for object detection focused on traffic environments.
- Based on a state-of-the-art CNN and adding viewpoint inference.
- Results comparable with non-real-time sophisticated approaches.
- Orientation is a step towards a complete scene understanding.

## Future work

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

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

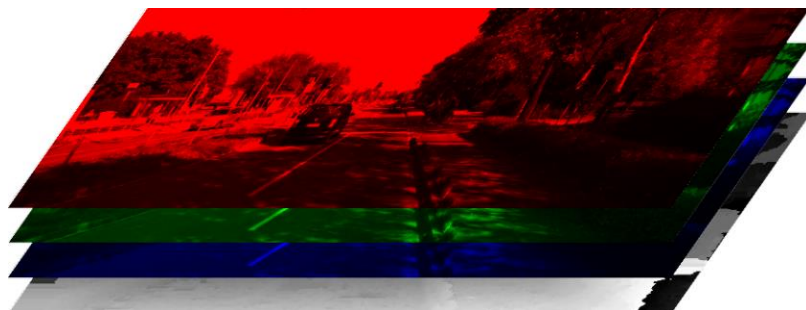
- Improvements:
  - Network architecture
  - Methods to overcome the fixed-size receptive field

# Other Developments

# Additional data for the CNN

30

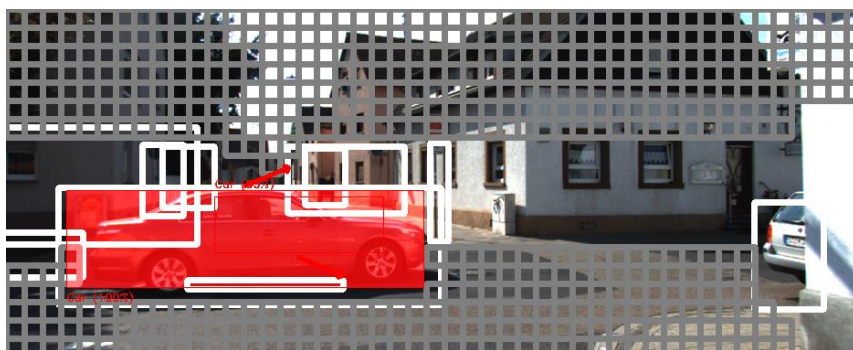
- **Stereo.** Using the disparity map as a fourth channel



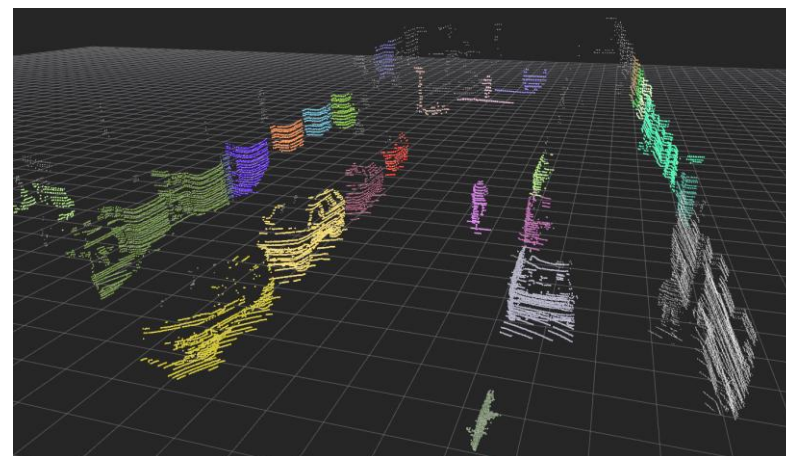
Input	Easy	Moderate	Hard
RGB	79.98	66.46	57.53
RGB+SGM	+0.30	+0.01	+2.52
RGB+DispNet	<b>+1.03</b>	<b>+1.14</b>	<b>+3.57</b>

C. Guindel, D. Martín, and J. M. Armingol, "Stereo Vision-Based Convolutional Networks for Object Detection in Driving Environments," in EUROCAST 2017 - Extended Abstracts, 2017, pp. 288–289.

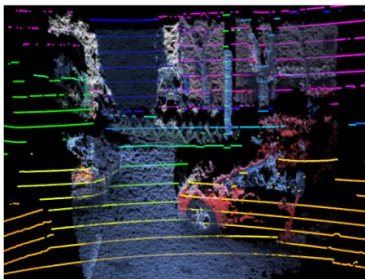
- **Laser.** Classifying object proposals coming from the Velodyne



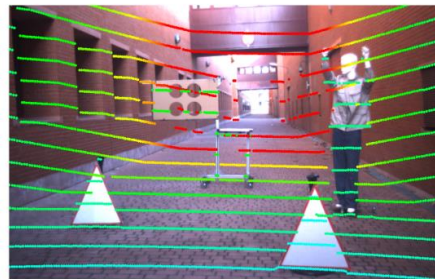
Work in progress!



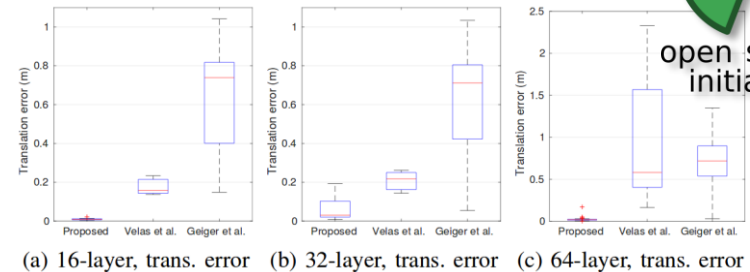
- Lidar-camera calibration. Solving a recurrent problem.



(a)



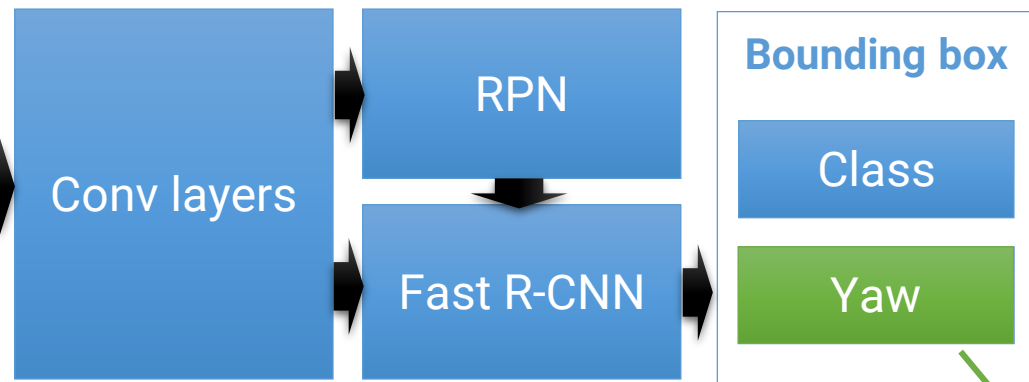
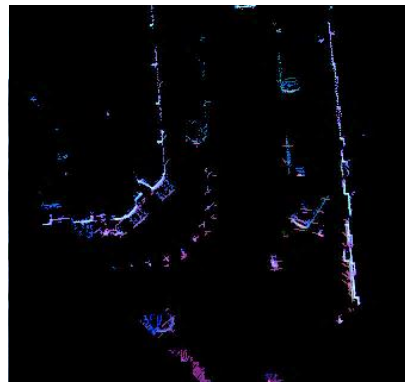
(b)



(a) 16-layer, trans. error (b) 32-layer, trans. error (c) 64-layer, trans. error

**C. Guindel**, J. Beltrán, D. Martín and F. García. "Automatic Extrinsic Calibration for Lidar-Stereo Vehicle Sensor Setups." arXiv:1705.04085 [cs.CV], 2017.

- Didi Challenge. Faster R-CNN applied over image-like inputs (bird-view).



Discrete  
(8 bins)

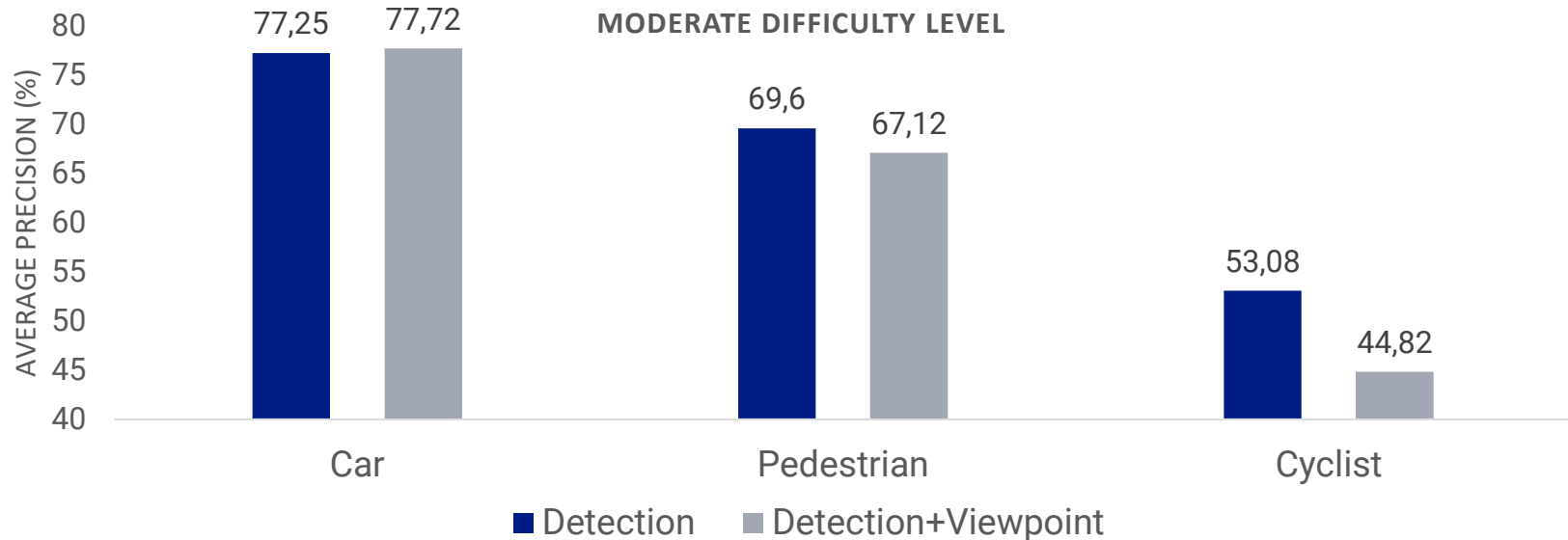
Proudly brought to you by J. Beltrán, D. Cruzado, F. Moreno and me.

# Thank you for your attention!

Questions?

# Precision change when introducing viewpoint

34



Detection ( $\Delta$ AP 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