

# Modeling Traffic Scenes for Intelligent Vehicles using CNN-based Detection and Orientation Estimation

**Carlos Guindel**, David Martín and José María Armingol

Intelligent Systems Laboratory (LSI) · Universidad Carlos III de Madrid

---

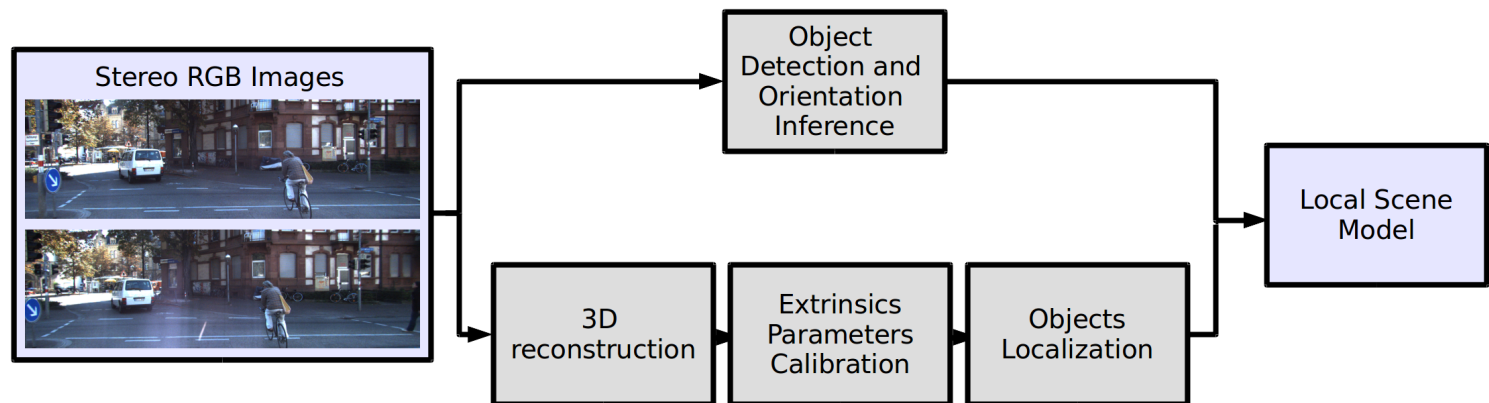


Sevilla · 23 November 2017

# Agenda

2

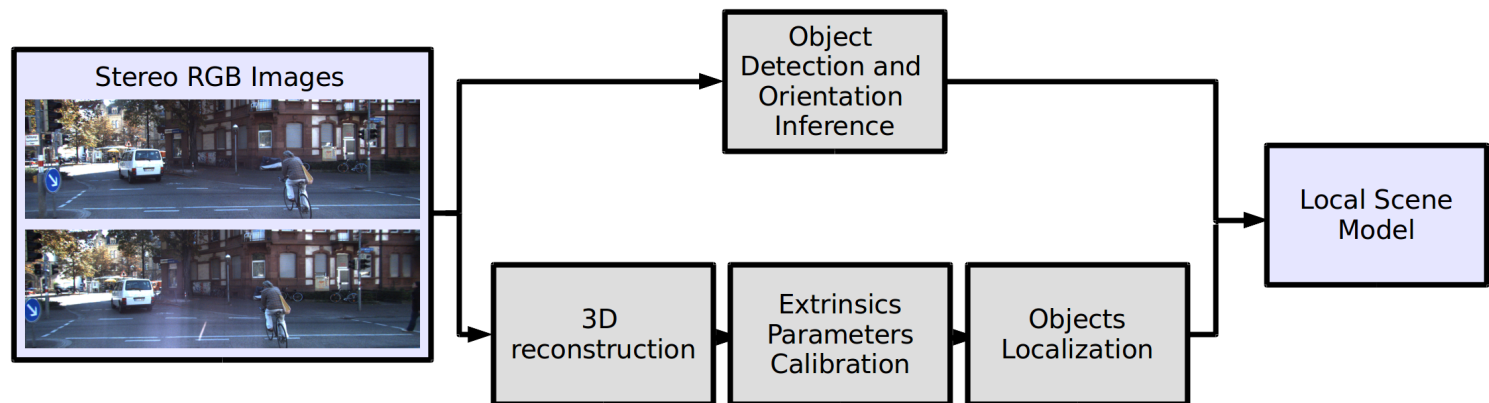
- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion



# Agenda

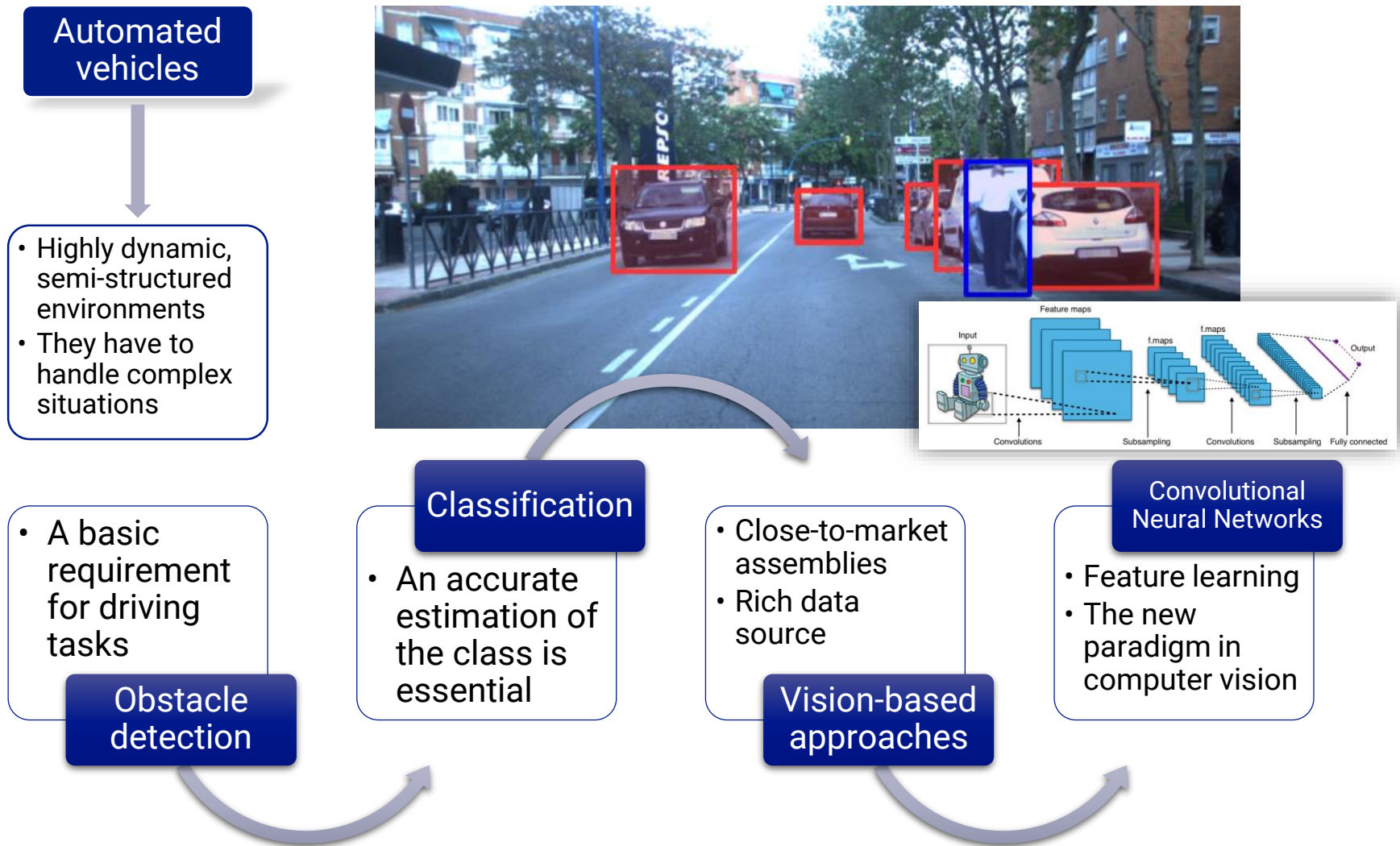
3

- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion



# Introduction

4

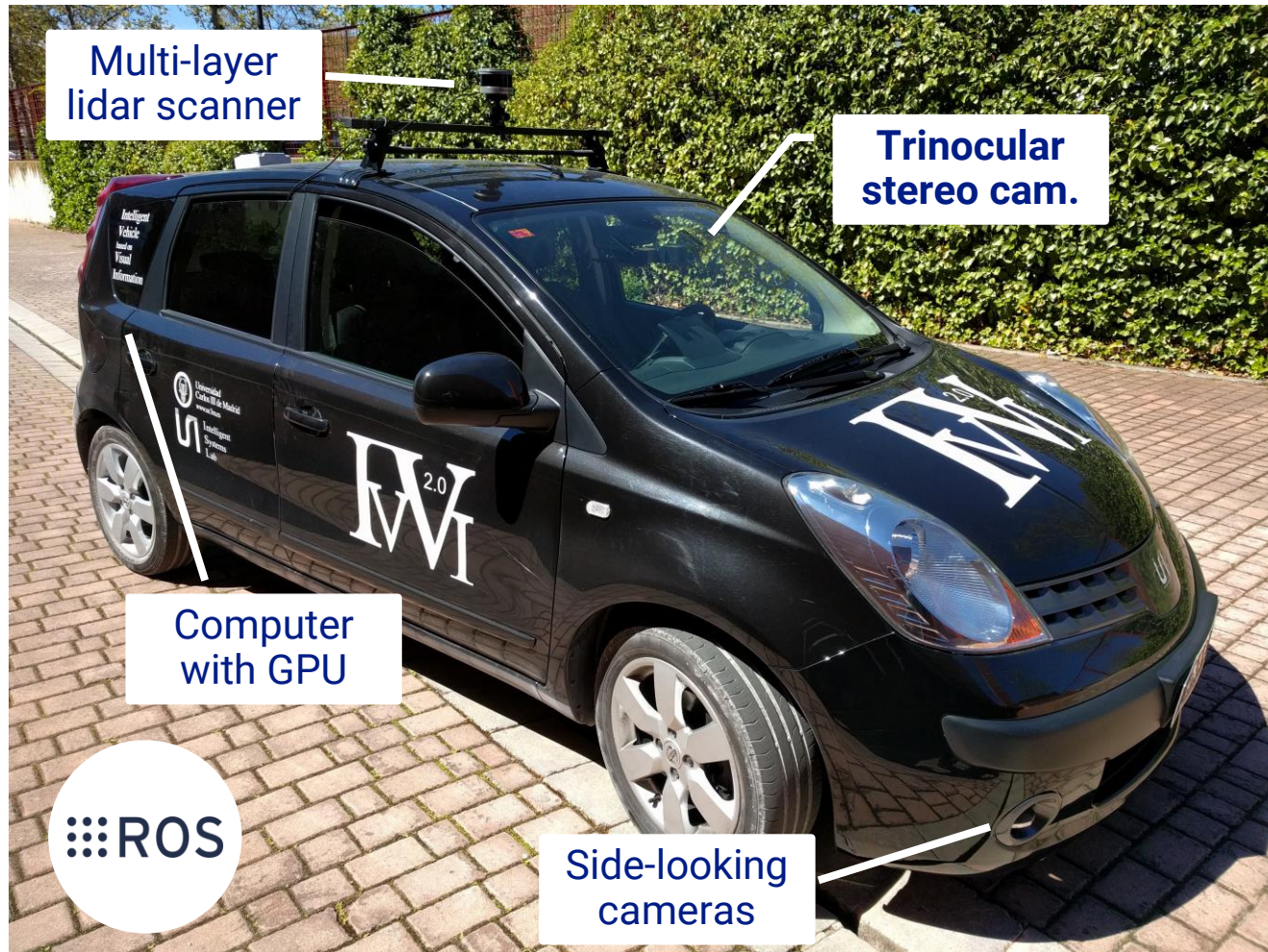




# IVVI 2.0 project

5

## INTELLIGENT VEHICLE BASED ON VISUAL INFORMATION 2.0

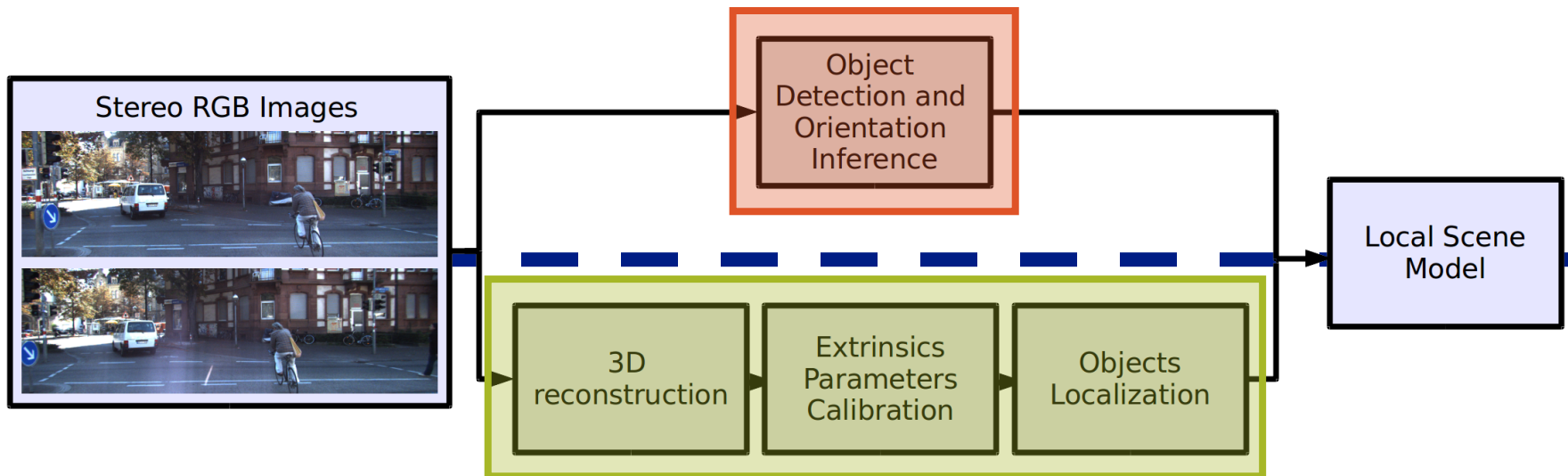


+info:  
[uc3m.es/islab](http://uc3m.es/islab)

# System overview

6

- Two main branches intended to run in parallel
- Obstacle detection
  - Features are extracted exclusively from the left stereo image

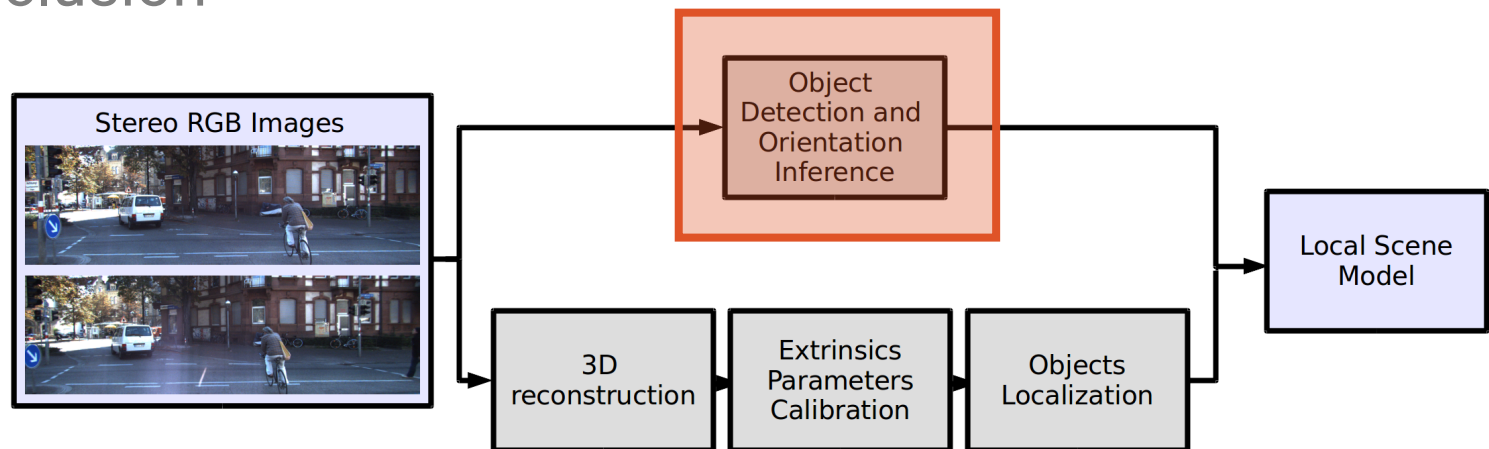


- Scene modeling
  - Stereo-based 3D reconstruction & flat-ground assumption

# Agenda

7

- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion



# Faster R-CNN framework

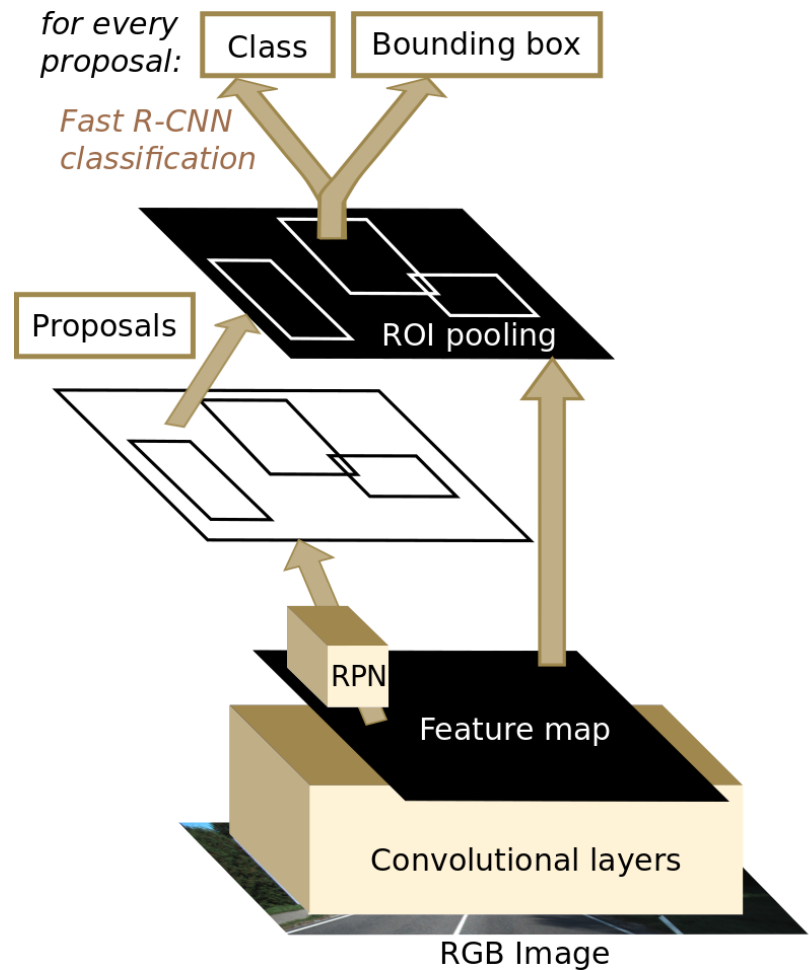
8

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



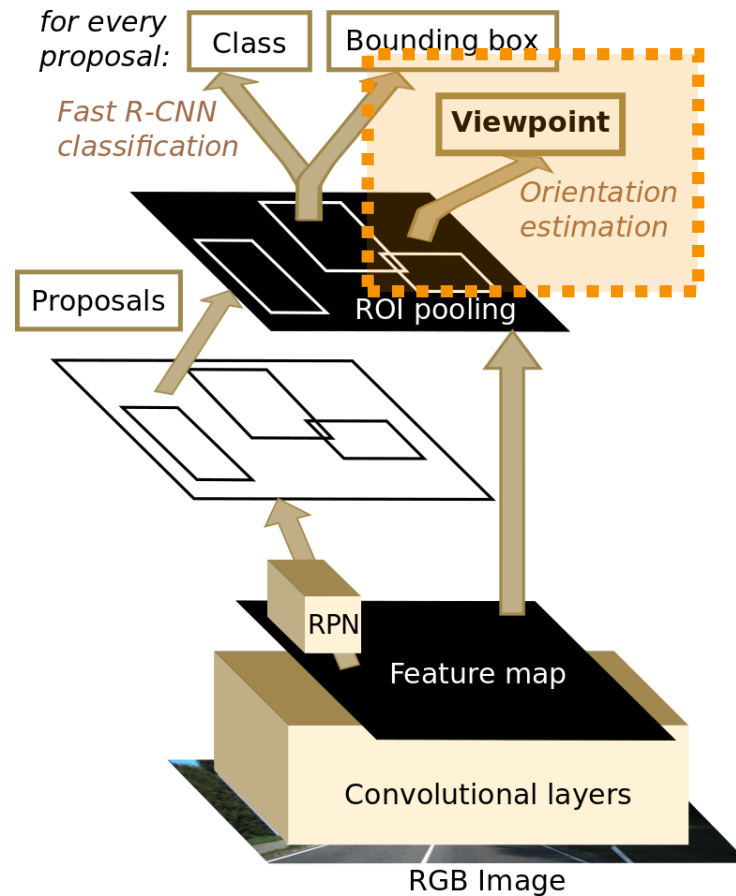
S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 6, pp. 1137–1149, 2016.



# Viewpoint estimation

9

- Faster R-CNN framework was modified to introduce **viewpoint inference**



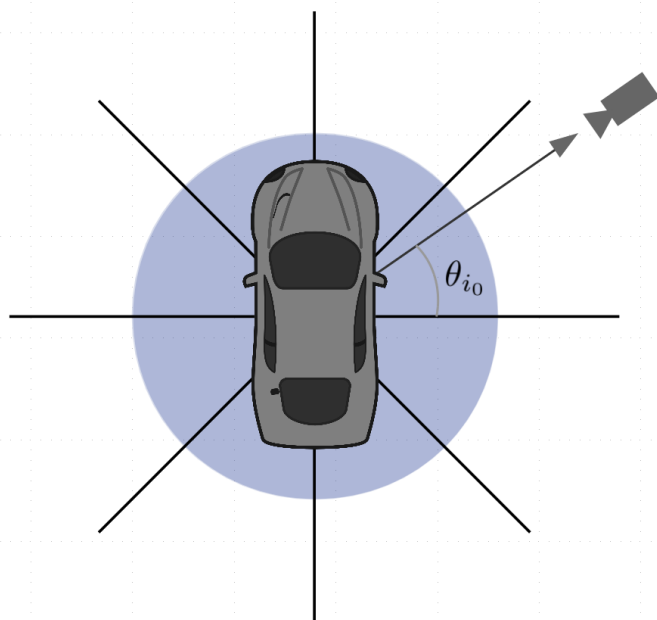
C. Guindel, D. Martin, and J. M. Armingol, "Joint object detection and viewpoint estimation using CNN features," in Proc. of the IEEE International Conference on Vehicular Electronics and Safety (ICVES), 2017, pp. 145–150.

# Discrete viewpoint inference

10

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

$N_b = 8$



- Every object is assigned a bin

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 gives a categorical distribution

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

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



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

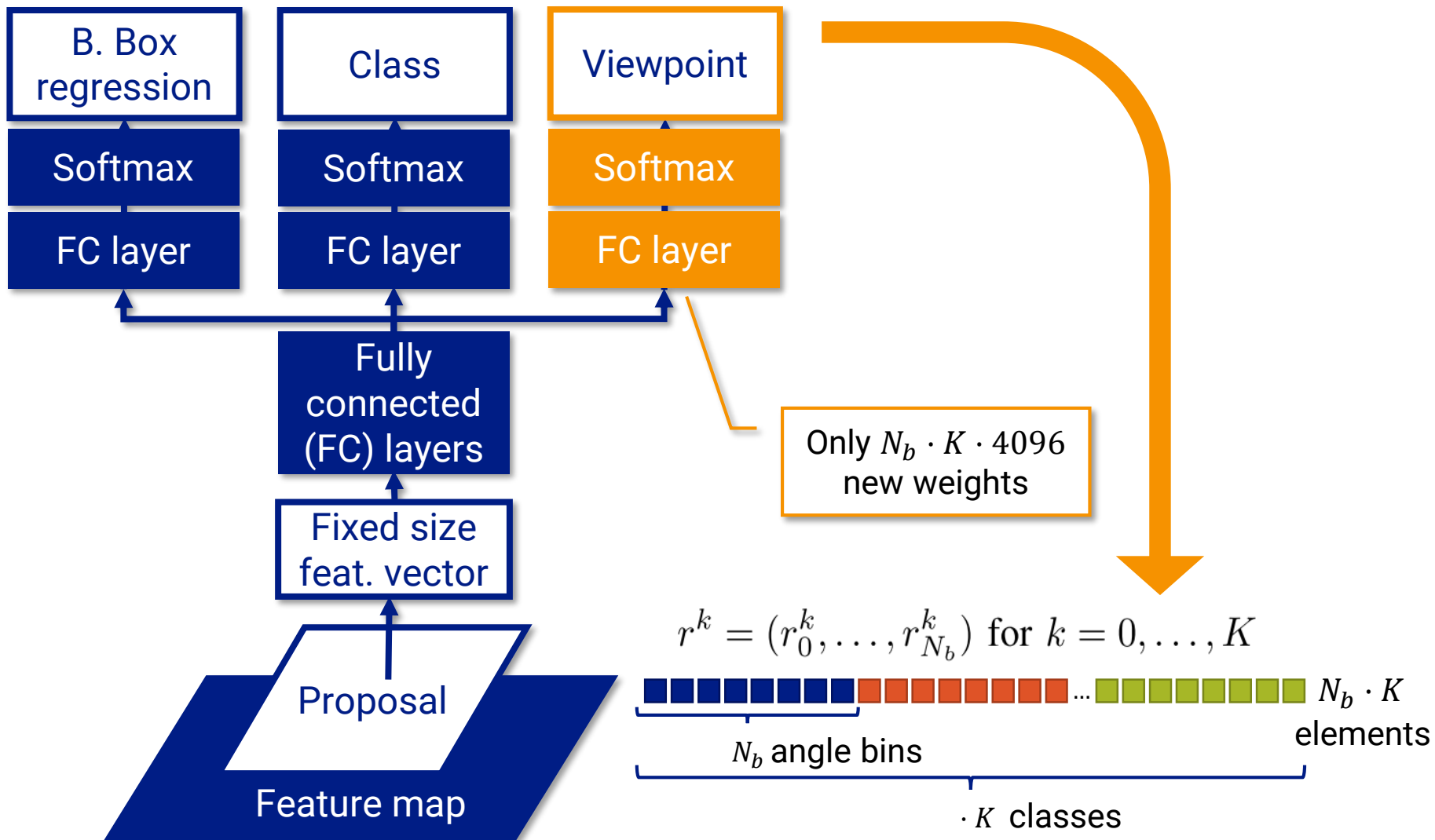
Elements  
of  $r$

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

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

# Joint detection and viewpoint estimation

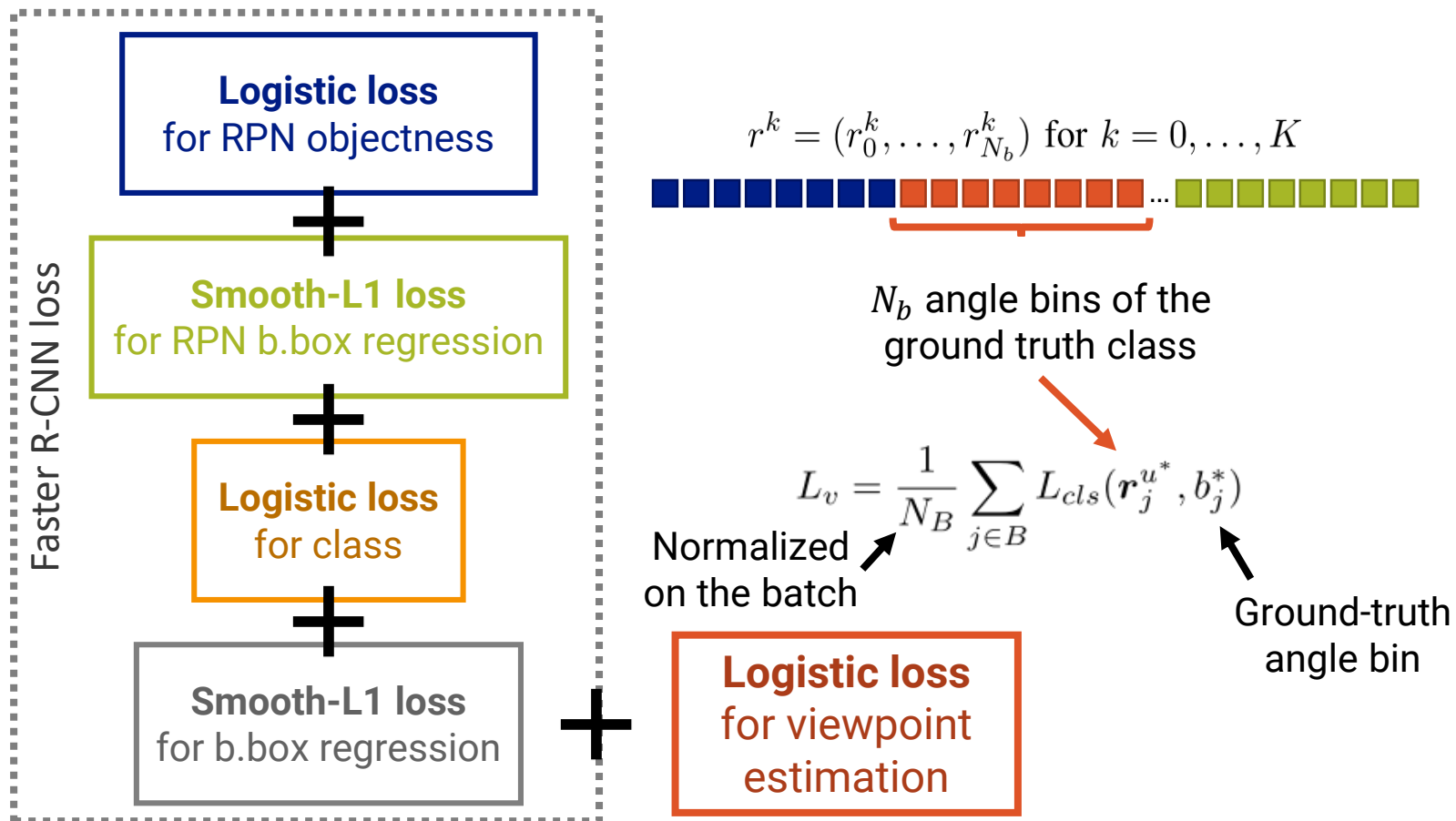
11



# Loss function and training

12

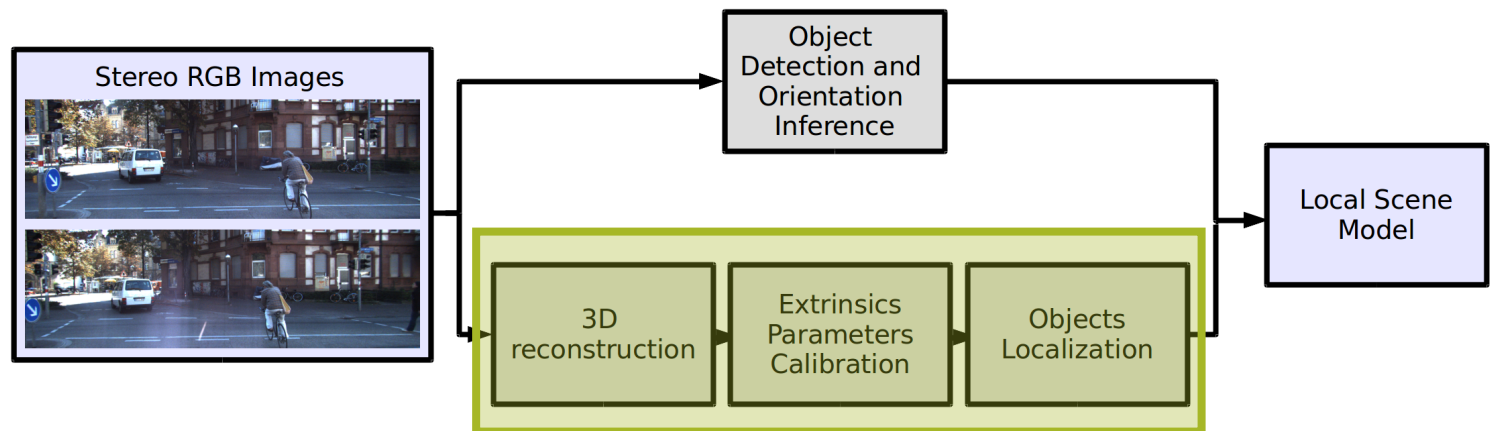
- Unweighted multi-task loss with **five** components



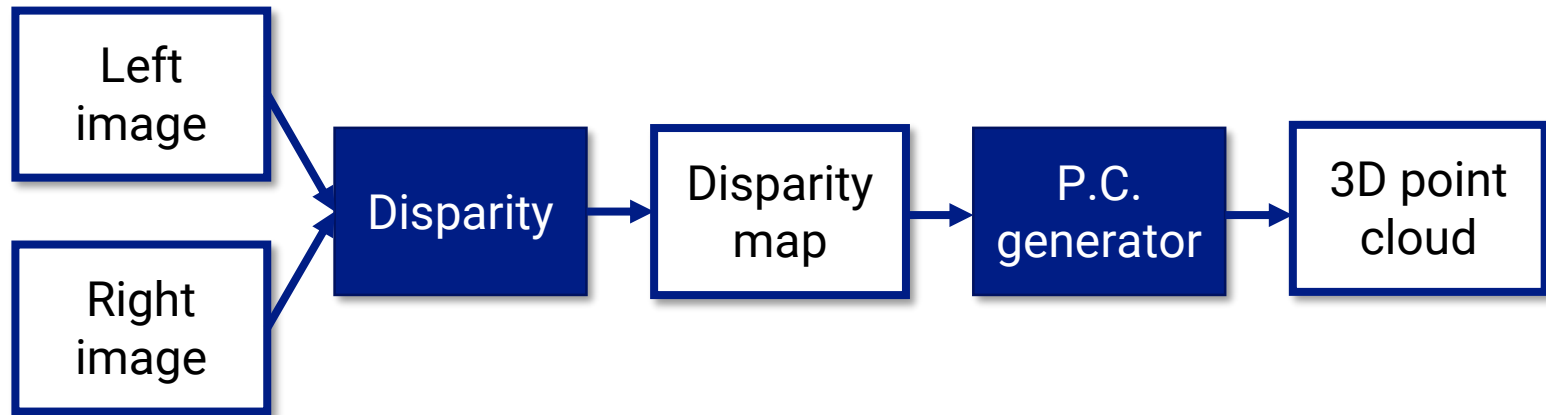
# Agenda

13

- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion

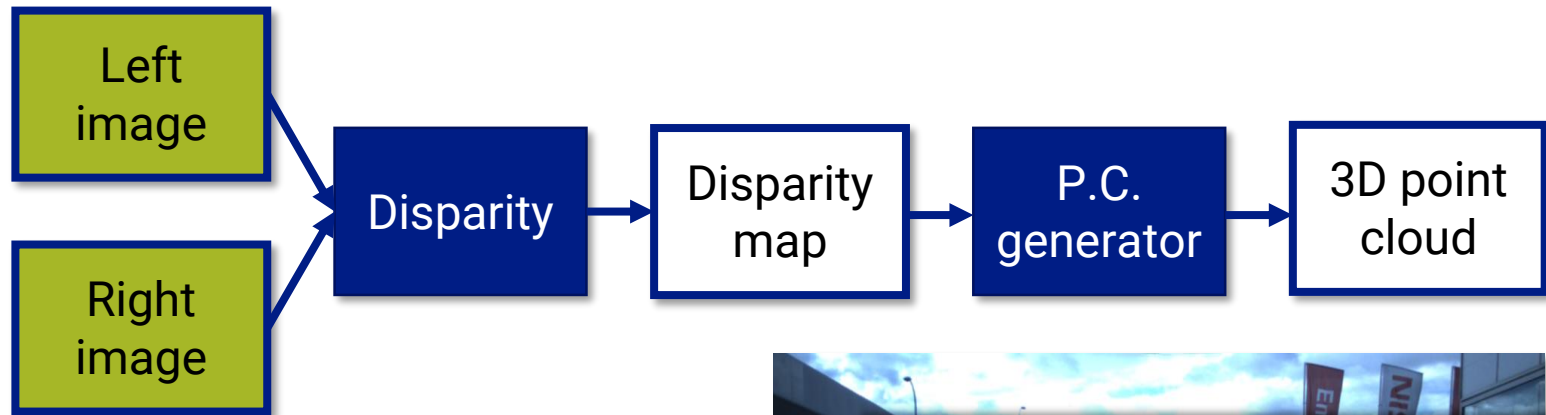


## 3D RECONSTRUCTION

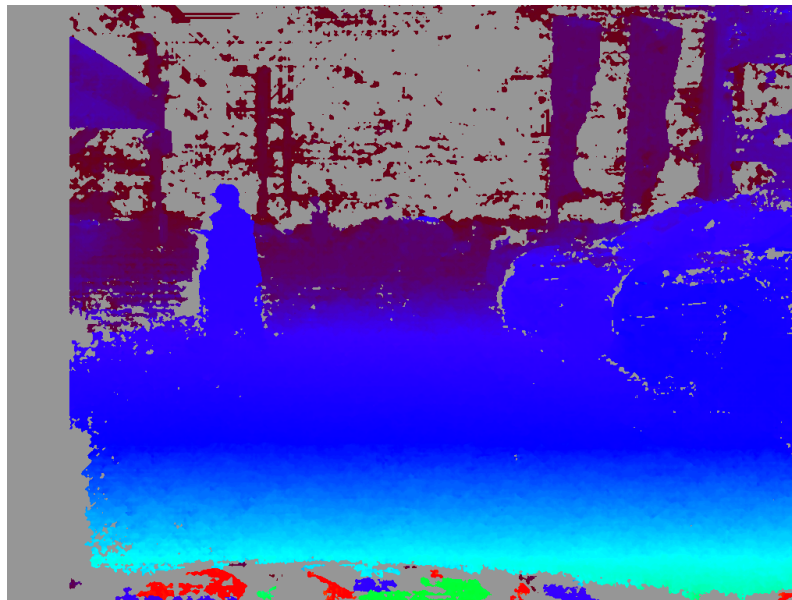
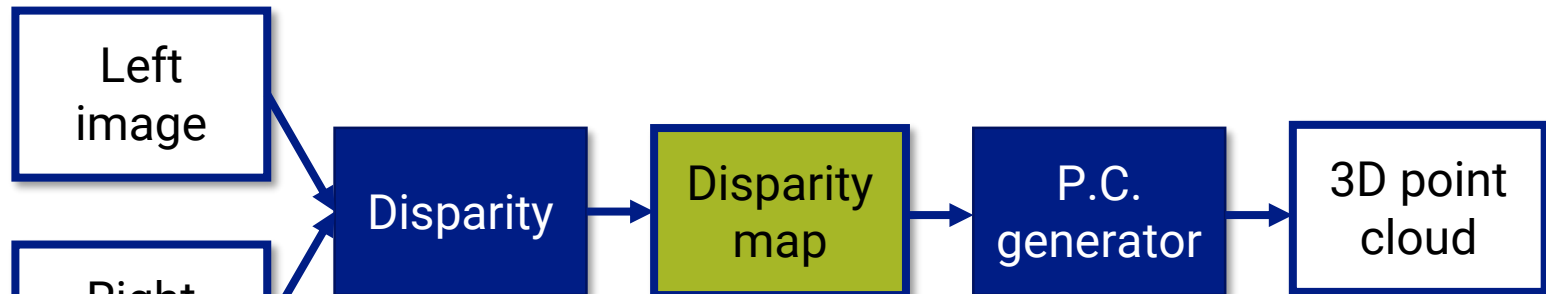




## 3D RECONSTRUCTION

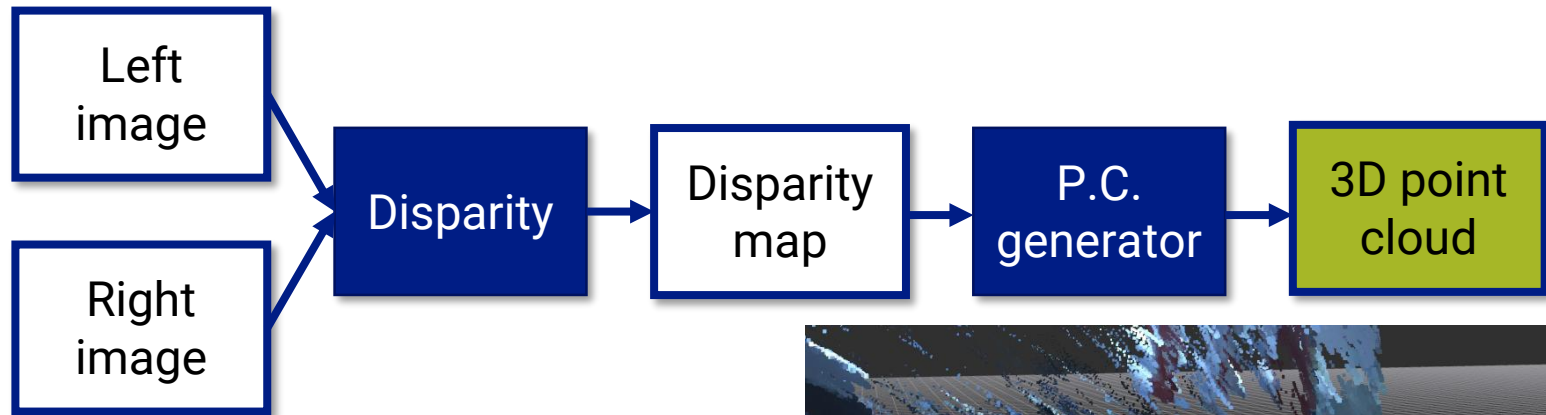


## 3D RECONSTRUCTION



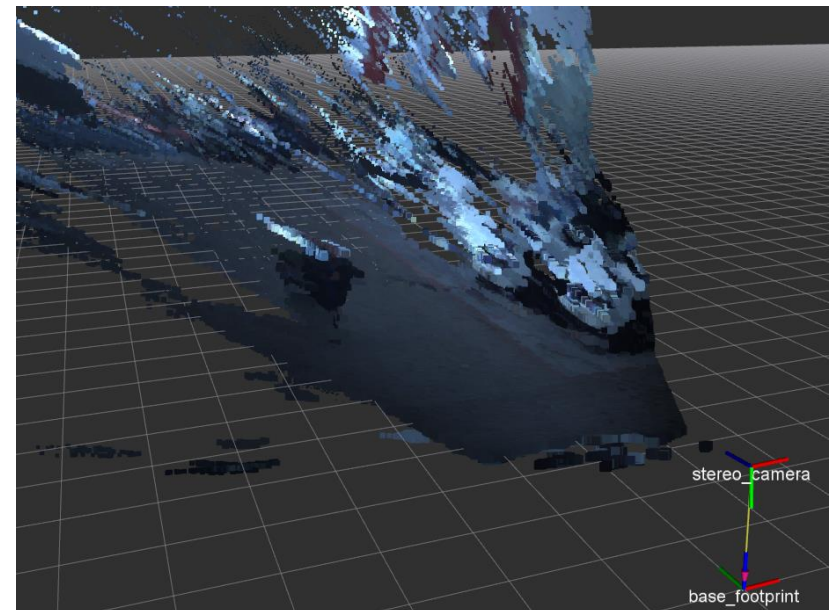
**SGM stereo matching**  
Suitable for environments  
with lack of texture,  
illumination changes, etc.

## 3D RECONSTRUCTION

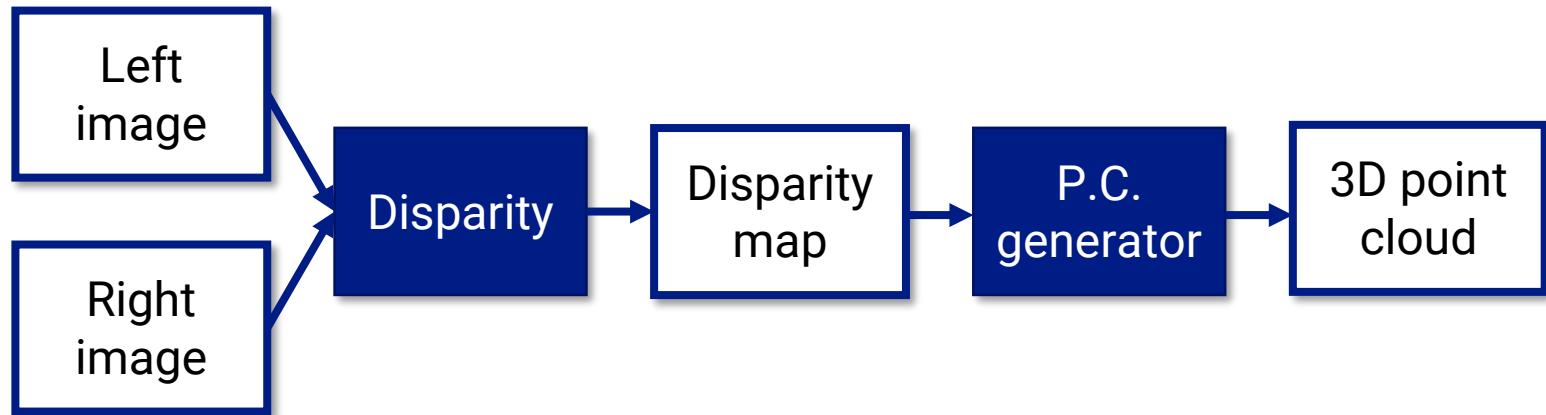


### **Pin-hole + disparity**

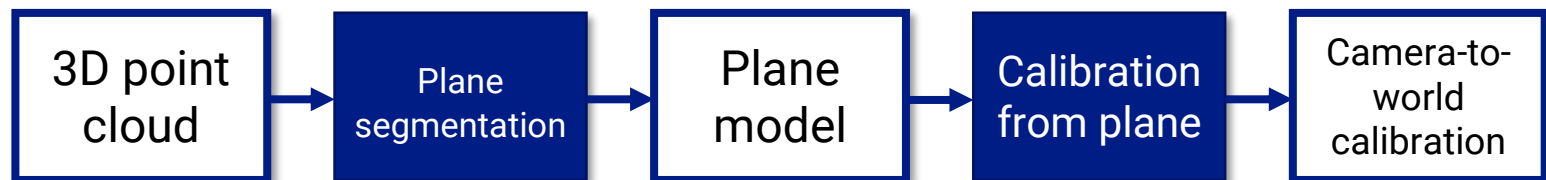
We build a XYZRGB cloud from the left image and the disparity map

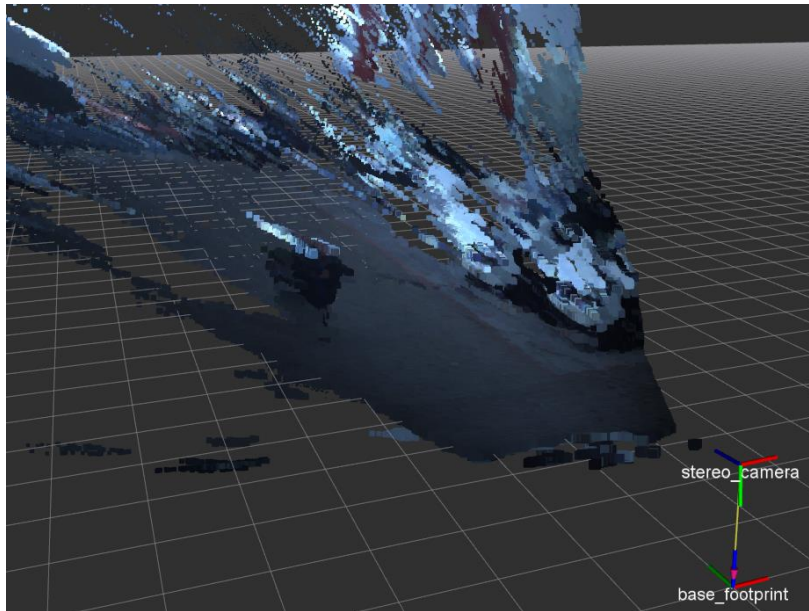


## 3D RECONSTRUCTION



## EXTRINSIC PARAMETERS AUTO-CALIBRATION

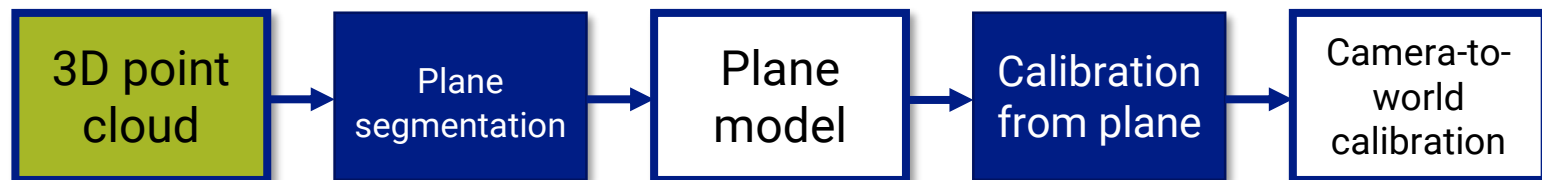




## Voxel grid downsampling

The cloud from the 3D reconstruction pipeline is downsampled (grid size: 20 cm)

## EXTRINSIC PARAMETERS AUTO-CALIBRATION



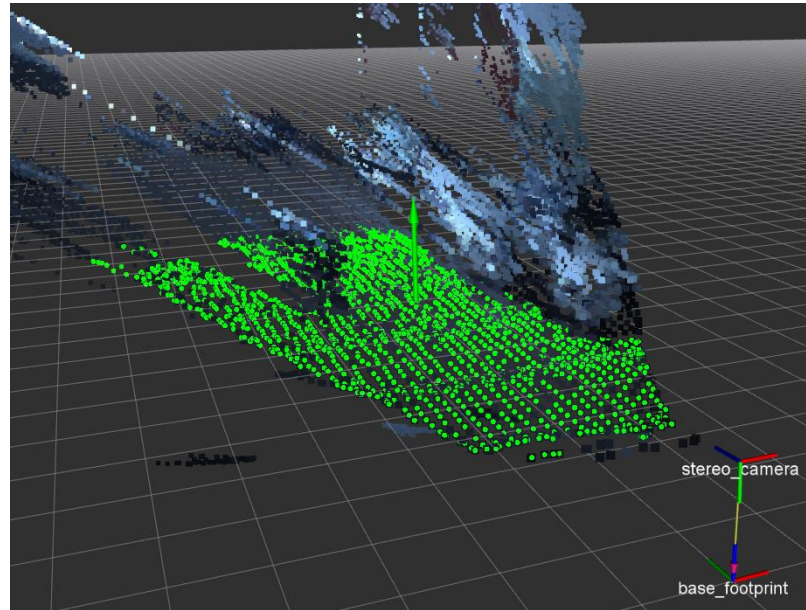
# Scene Modeling

20

## ...Pass through filters

Vertical axis: 0-2 m

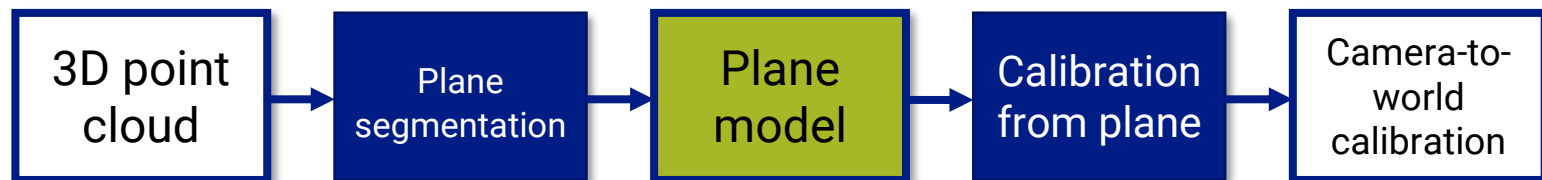
Depth axis: 0-20 m



## Planar segmentation

Using RANSAC with a 10 cm threshold, and a small angular tolerance.

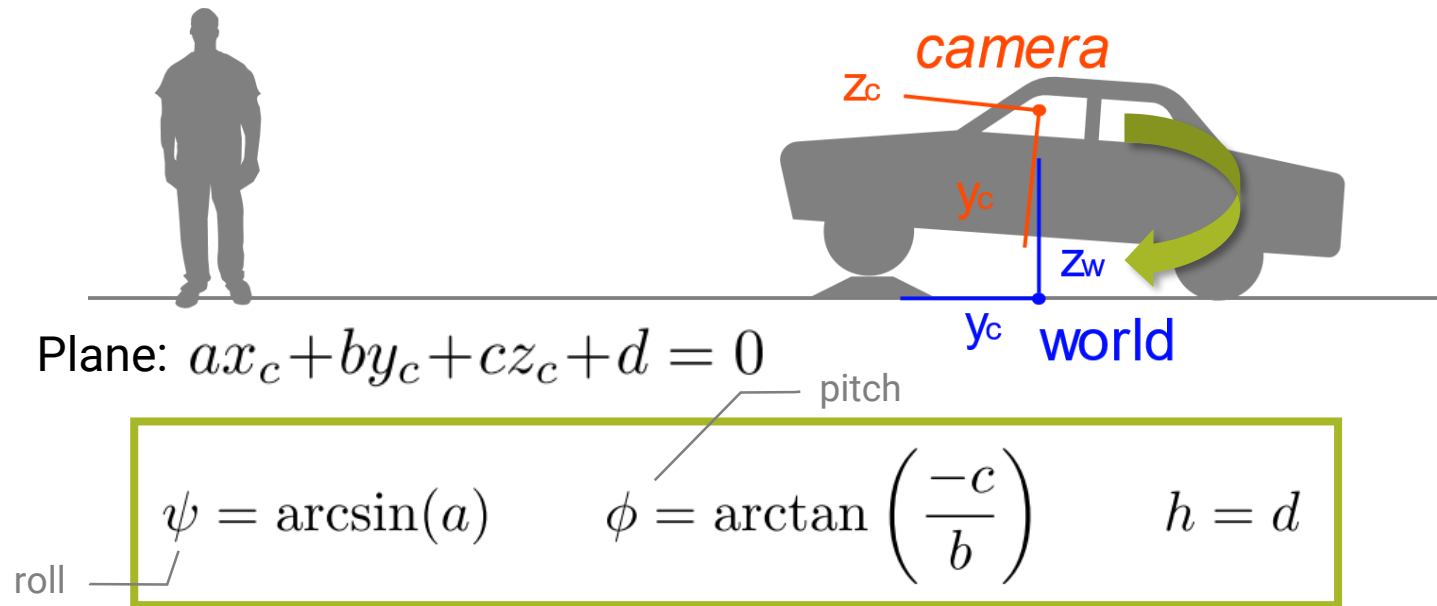
## EXTRINSIC PARAMETERS AUTO-CALIBRATION



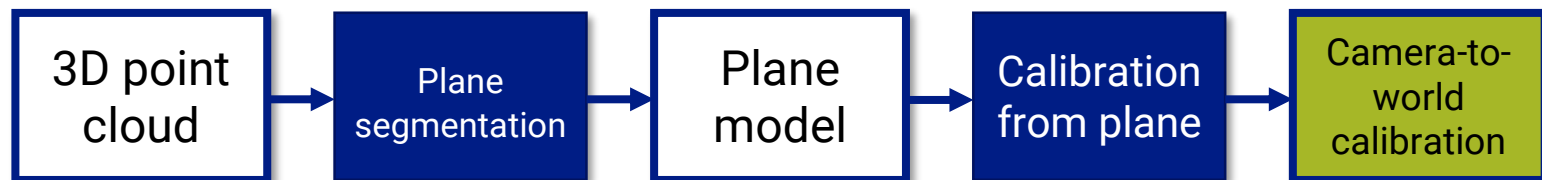


# Scene modeling

21



## EXTRINSIC PARAMETERS AUTO-CALIBRATION



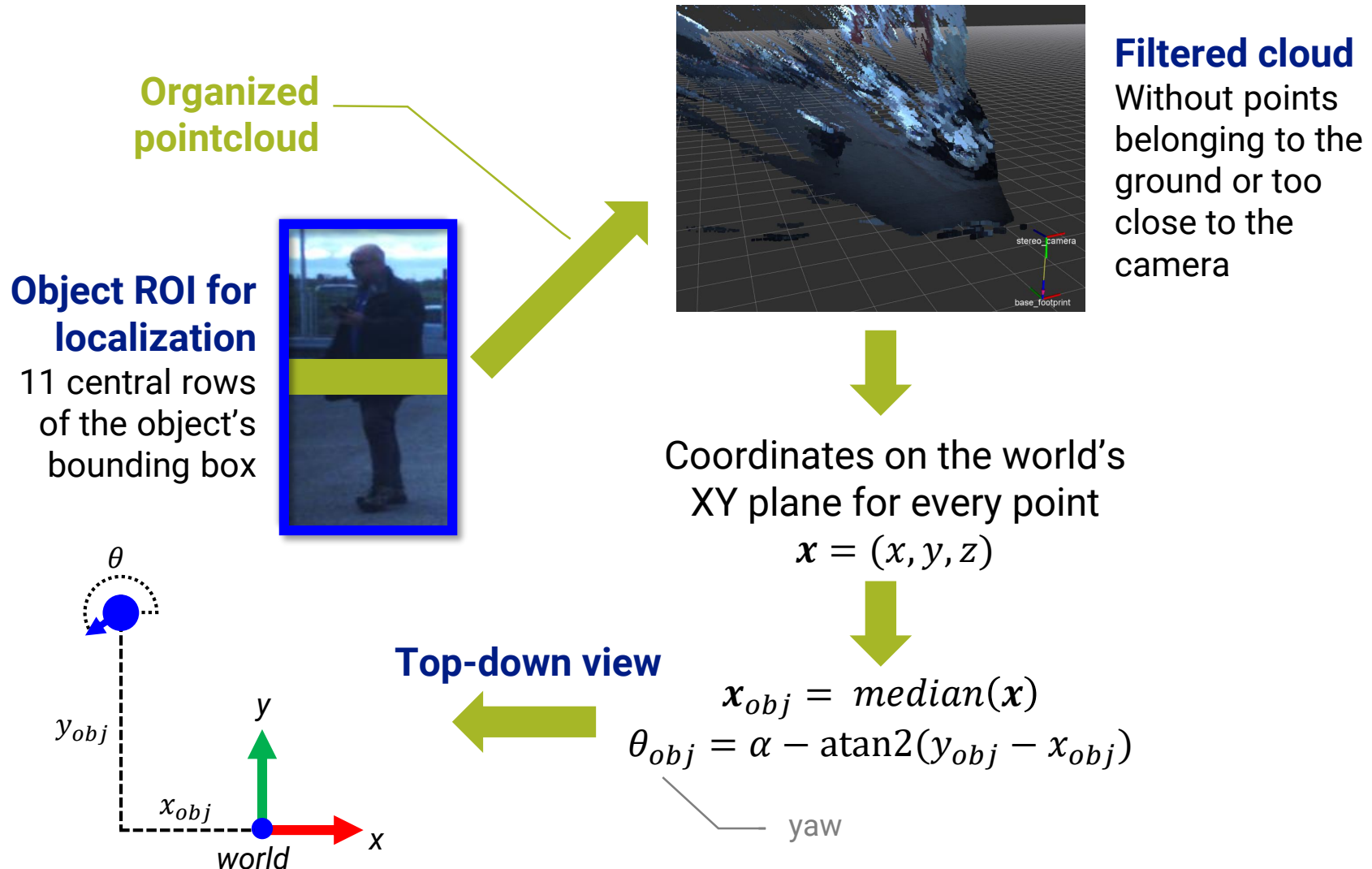
# Object localization

22



# Object localization

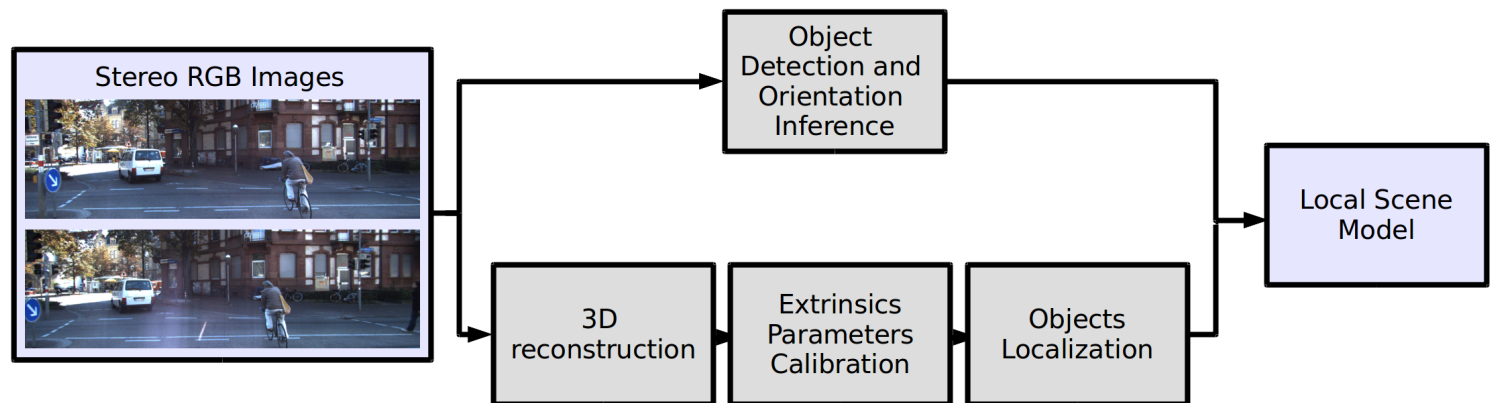
23



# Agenda

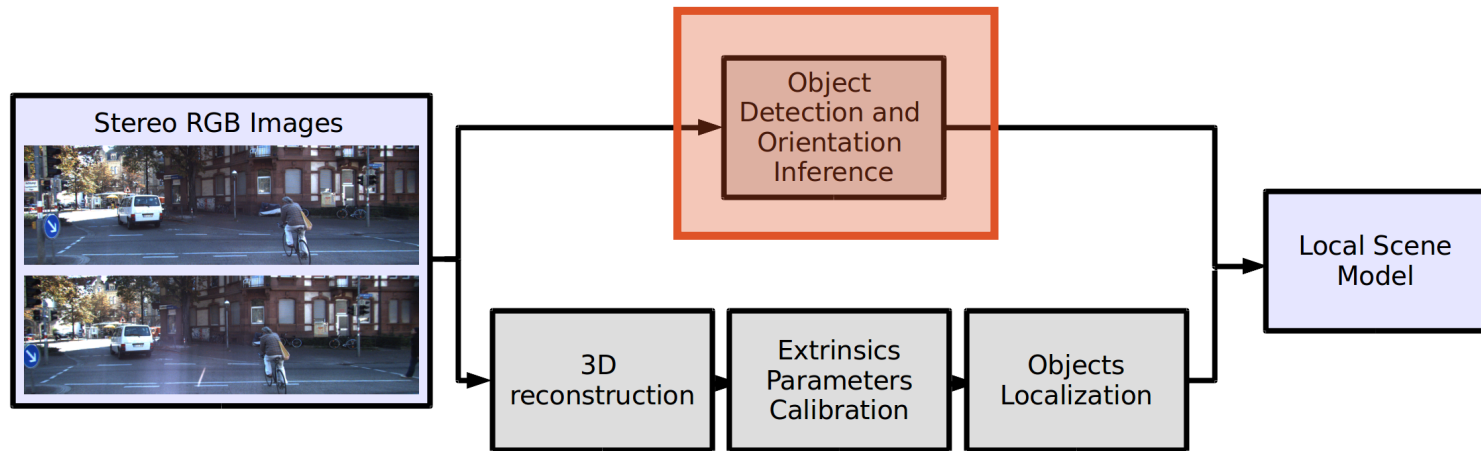
24

- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion



# Results: Detection and viewpoint estimation

25



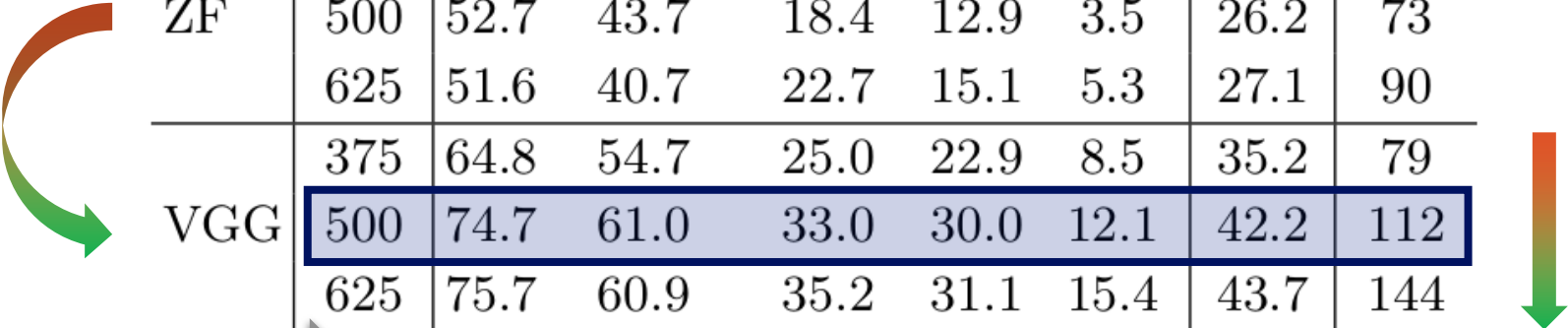
- KITTI Object Detection Benchmark
  - 5,576 images for training and 2,065 for validation
  - Labels for class and orientation available
- Evaluation metric
  - Average Orientation Similarity (AOS)

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

# Results: Detection and viewpoint estimation

26

- Two different architectures:
  - ZF (lightweight) and VGG 16-layer (more complex)
- Three different scales (height in pixels):
  - 375, 500, 625



Net	Scale	Car	Pedest.	Cyclist	Van	Truck	mean	Time (ms)
ZF	375	44.2	35.6	16.1	8.5	3.2	21.5	46
	500	52.7	43.7	18.4	12.9	3.5	26.2	73
	625	51.6	40.7	22.7	15.1	5.3	27.1	90
VGG	375	64.8	54.7	25.0	22.9	8.5	35.2	79
	500	74.7	61.0	33.0	30.0	12.1	42.2	112
	625	75.7	60.9	35.2	31.1	15.4	43.7	144

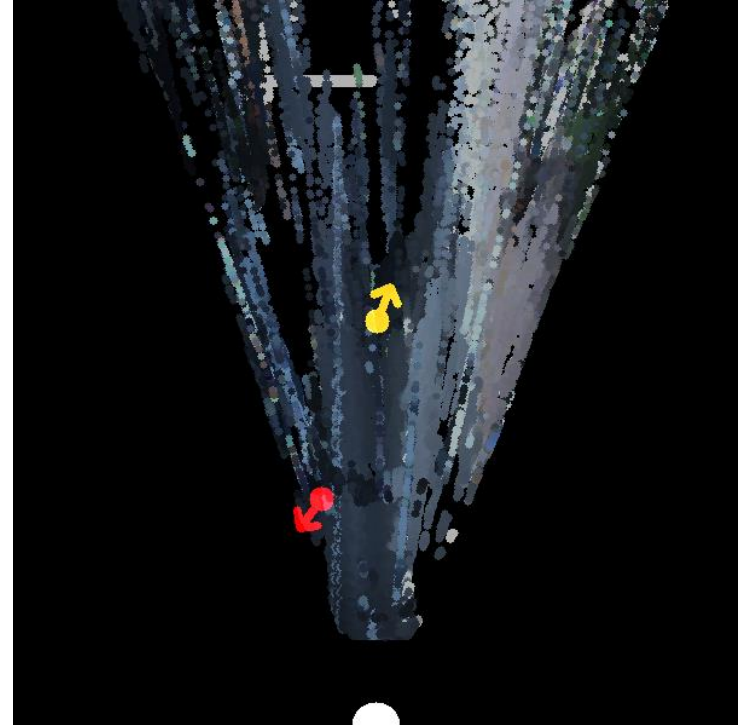
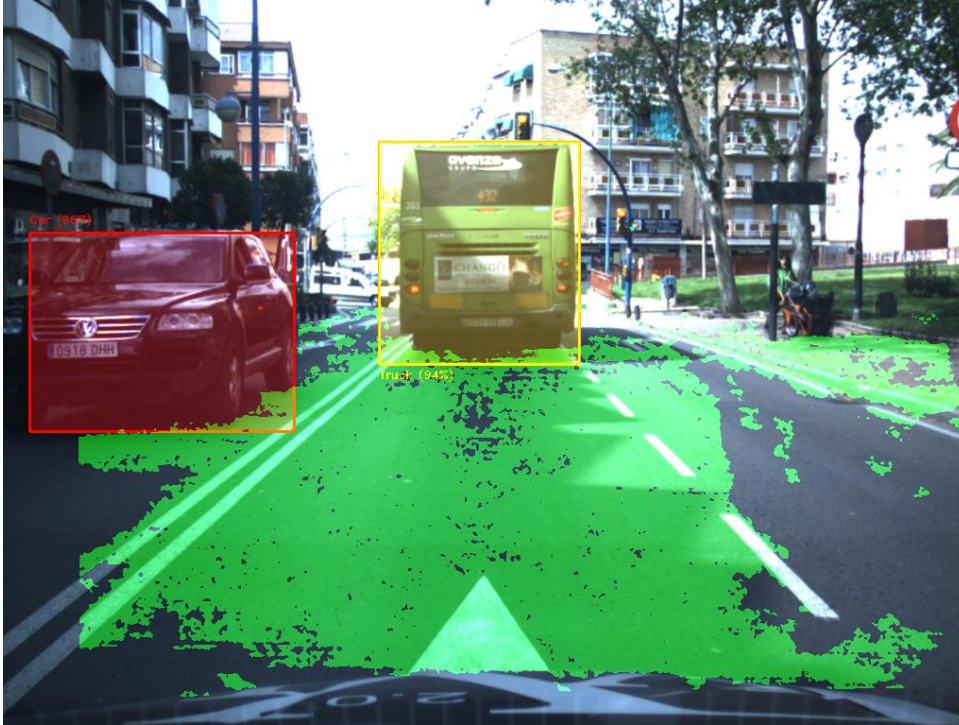
Top-performing  
comparable method  
in the KITTI ranking

88,43	66,28	63,41	N.A.	N.A.	N.A.	2 sec.
-------	-------	-------	------	------	------	--------



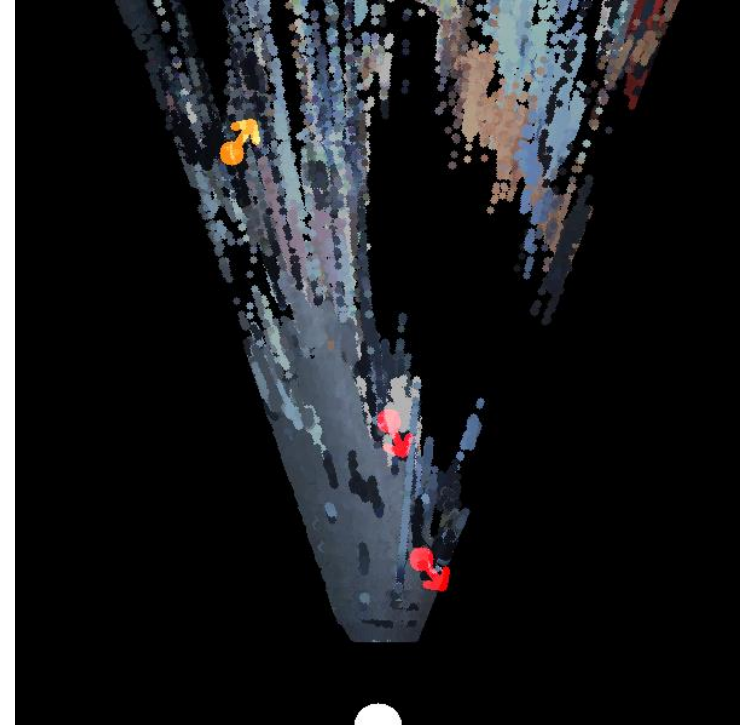
# Results: Scene modeling

27



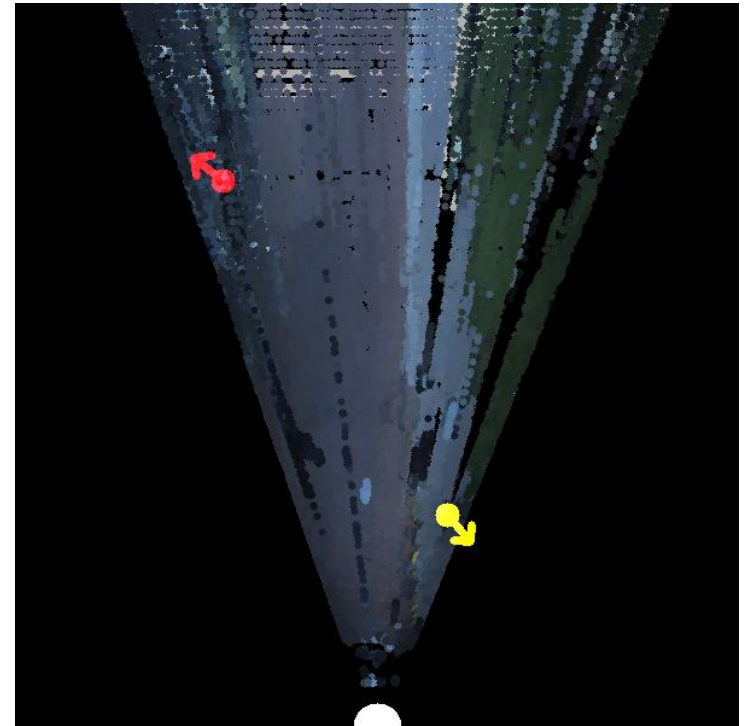
# Results: Scene modeling

28



# Results: Scene modeling

29

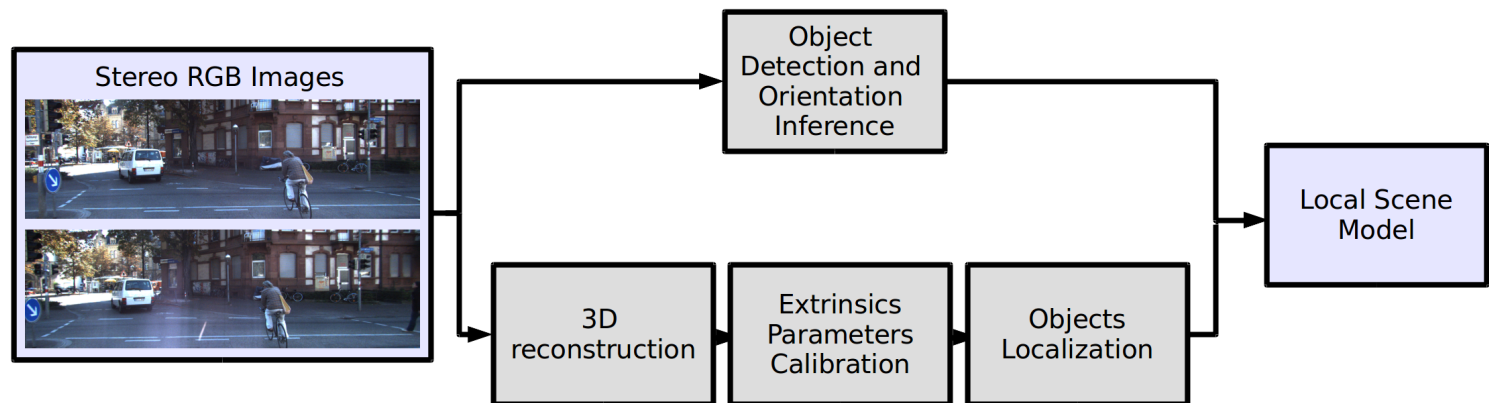




# Agenda

30

- ① Introduction
- ② Obstacle detection
- ③ Scene modeling
- ④ Results
- ⑤ Conclusion



# Conclusion

31

- Towards a full object-based scene understanding
  - CNN-based detection and viewpoint inference
  - Efficient approach: the same set of features is used for all tasks
- Stereo-vision 3D information is included for situation assessment
- Results validate our approach

## Future work

- New categories of traffic elements
- Extension to the time domain
  - Tracking, filtering, etc.
- Including information from other perception modules
  - E.g., semantic segmentation



Code for CNN detection & viewpoints available at  
<https://github.com/cguindel/lsi-faster-rcnn>

# THANKS FOR YOUR ATTENTION

---

23 November 2017

ROBOT'2017 - Third Iberian Robotics Conference



**Intelligent Systems Lab**  
[www.uc3m.es/islab](http://www.uc3m.es/islab)