



**ICSCC**  
2021

**1-3 July 2021**  
**Kochi, India**

# **8<sup>th</sup> International Conference on Smart Computing and Communications**



**Paper ID: 314**

**Title: A Comparative Study of Deep Learning-based Depth Estimation Approaches: Application to Smart Mobility**

**Name of Authors: Antoine Mauri - Redouane Khemmar - Benoit Decoux - Tahar Benmoumen - Madjid Haddad - Rémi Boutteau**

**Presented By: Redouane Khemmar**

**Affiliation: Normandie Univ, UNIROUEN, ESIGELEC, IRSEEM**

# Outline

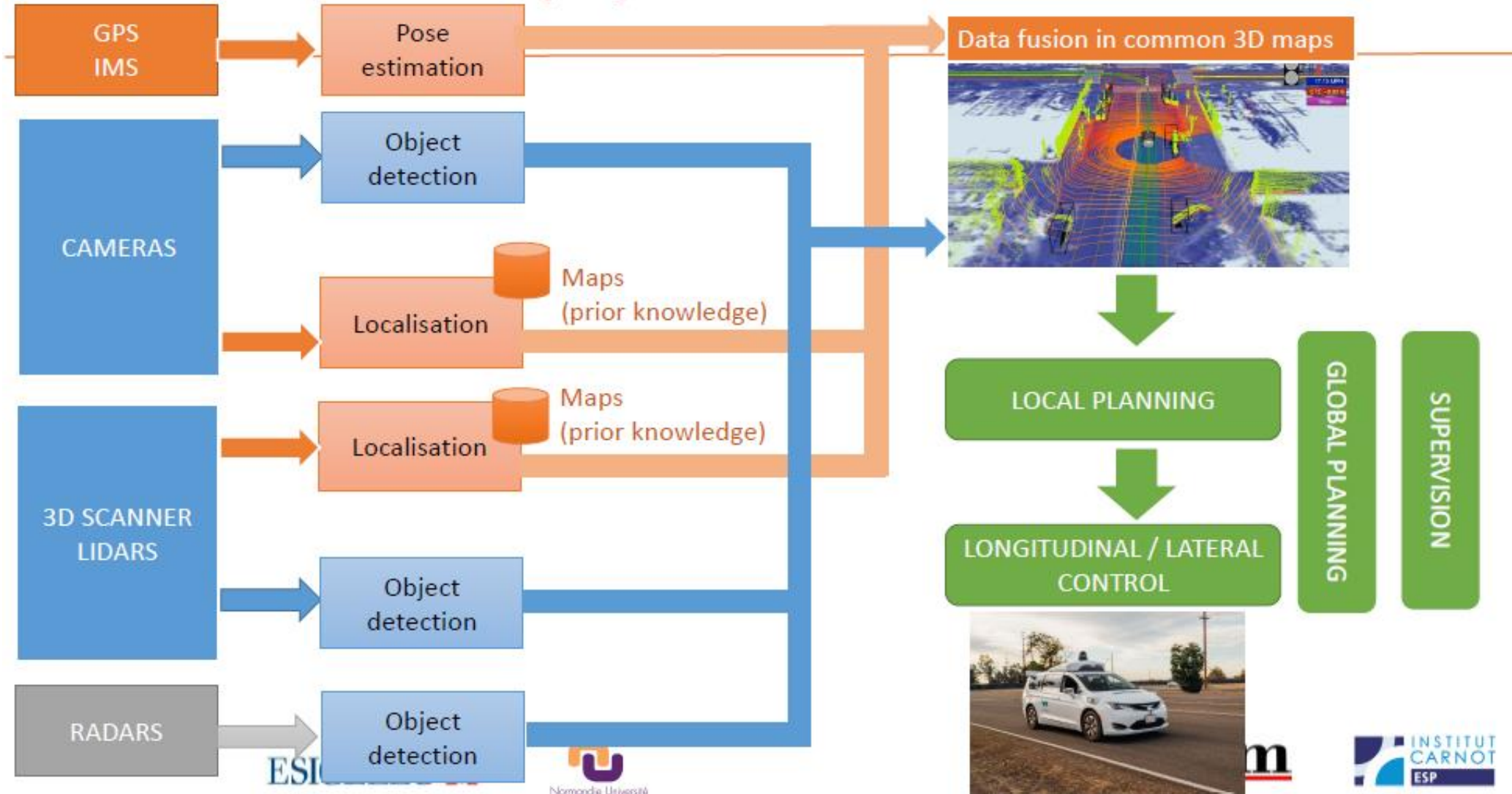
- I. Context & Motivation
- II. Related Work
- III. Evaluated Depth Estimation Methods
- IV. Our New Evaluation Protocols
- V. Experimental Results
- VI. Conclusion and Future Work

# Context & Motivation

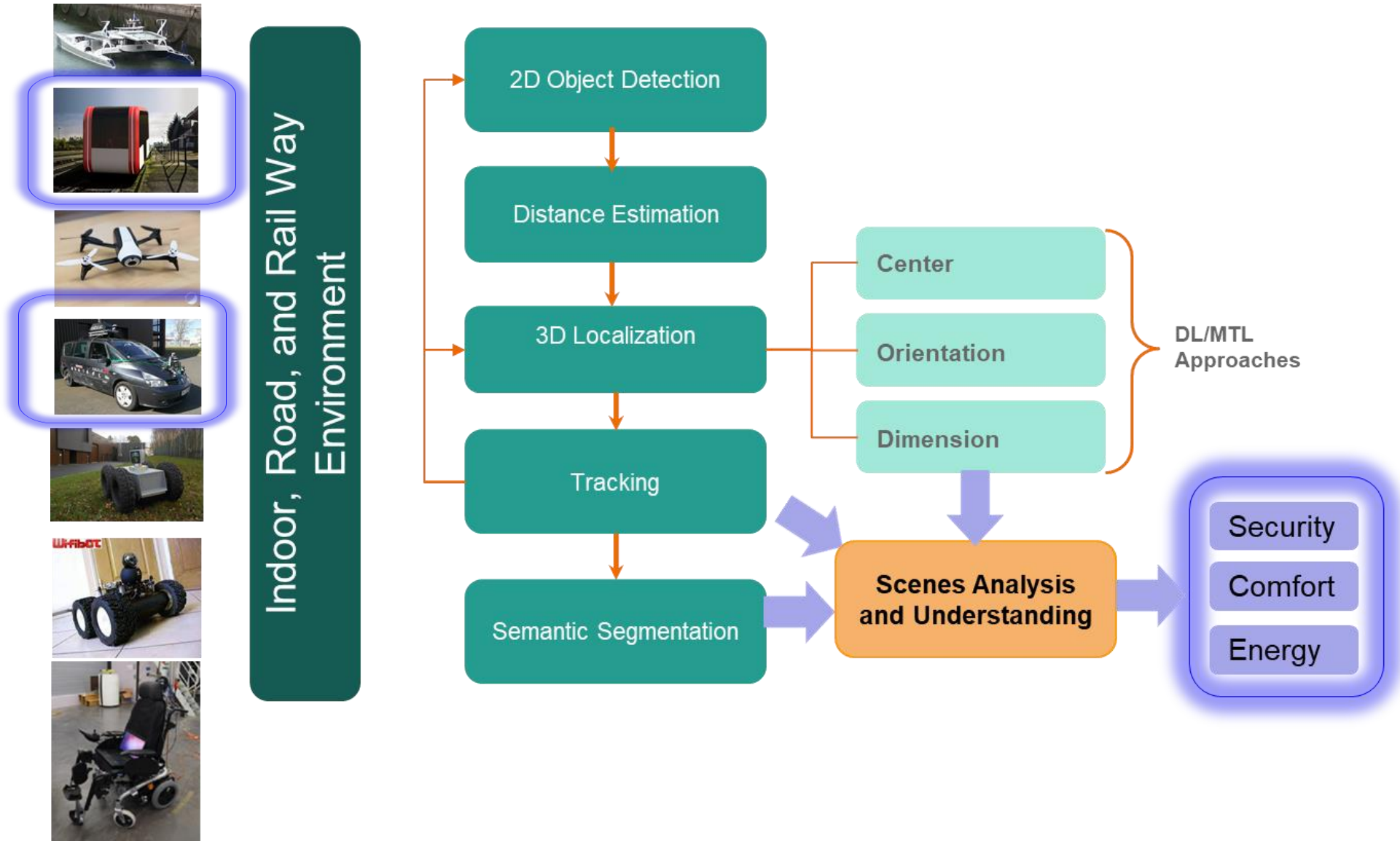
- ▶ Perception of environment is an important aspect for Autonomous Vehicles (AV)
  - Localization of potential obstacles is crucial to prevent collisions for AV
  - Object detection, localization, and tracking are important for scene analysis and understanding
- ▶ CNN based Depth Estimation Methods allow to use **Images from a Monocular or Stereoscopic Camera** to get a Dense Depth Prediction
  - Expensive depth sensors like LiDAR could be replaced by a single camera
- ▶ Few works has been done to compare the performances under realistic environment

# Context & Motivation

## AUTONOMOUS DRIVING (AD) SYSTEM



# Context & Motivation



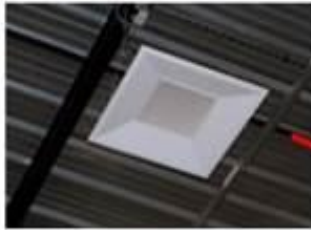
# Context & Motivation

## ▶ Autonomous Vehicle: Acquisition System



# Context & Motivation

## ► Autonomous Vehicle: Autonomous Navigation Laboratory



Eclairage contrôlé  
(tests en luminosité  
dégradée)



Système de capture 3D  
du mouvement VICON



Zone d'expérimentation

Zone « pluie » ;  
test en conditions  
d'environnement  
dégradé



Flotte de robots indoor  
(wifibots)

# Related Work : Depth Evaluation Metrics

## ► Depth error metrics used in our study:

- Relative Error :

$$RE = \frac{1}{N} \sum_i \sum_j \frac{|g_{i,j} - p_{i,j}|}{g_{i,j}}$$

- Squared Relative Error :

$$SRE = \frac{1}{N} \sum_i \sum_j \frac{|g_{i,j} - p_{i,j}|^2}{g_{i,j}}$$

- Root Mean Squared Error :

$$RMSE = \sqrt{\frac{1}{N} \sum_i \sum_j (g_{i,j} - p_{i,j})^2}$$

- Logarithmic Root Mean Squared Error :

$$logRMSE = \sqrt{\frac{1}{N} \sum_i \sum_j (\log(g_{i,j}) - \log(p_{i,j}))^2}$$

p: depth **prediction** of a pixel in the image

g: depth **ground truth** of a pixel in the image

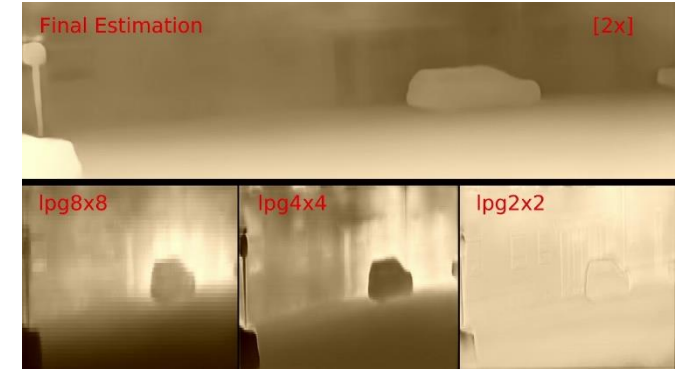
N: total of pixels in the image



# Evaluated Single Image Depth Estimation Methods

## ▶ BTS<sup>1</sup> (Big To Small)

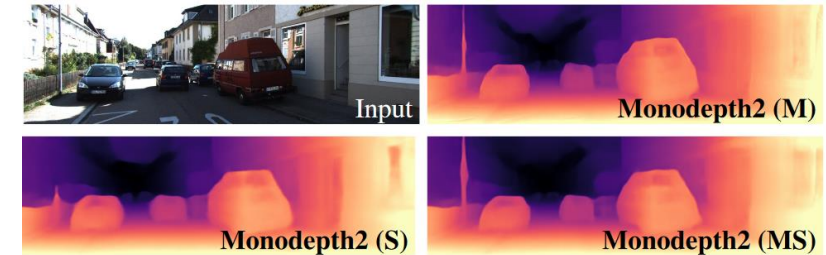
- **Supervised** monocular depth estimation
- Layers located at multiple stages of the decoding phase are used
- Their outputs are combined to **predict depth at full resolution**
- Achieves top precision on [KITTI](#) depth benchmark



BTS Approach

## ▶ Monodepth2<sup>2</sup>

- Can be trained with or without depth supervision
- With self-supervision, the problem of depth estimation is casted into an image reconstruction one
- Self-supervision works with:
  - Monocular image sequences (M)
  - Stereo data (S)
  - Or both of them (MS)



Monodepth2 Approach

1. J. H. Lee, M.-K. Han, D. W. Ko, and I. H. Suh, "From big to small: Multi-scale local planar guidance for monocular depth estimation," arXiv preprint arXiv:1907.10326, 2019.

2. C. Godard, O. Mac Aodha, M. Firman, and G. Brostow, "Digging into self-supervised monocular depth estimation," arXiv preprint arXiv:1806.01260, 2018.

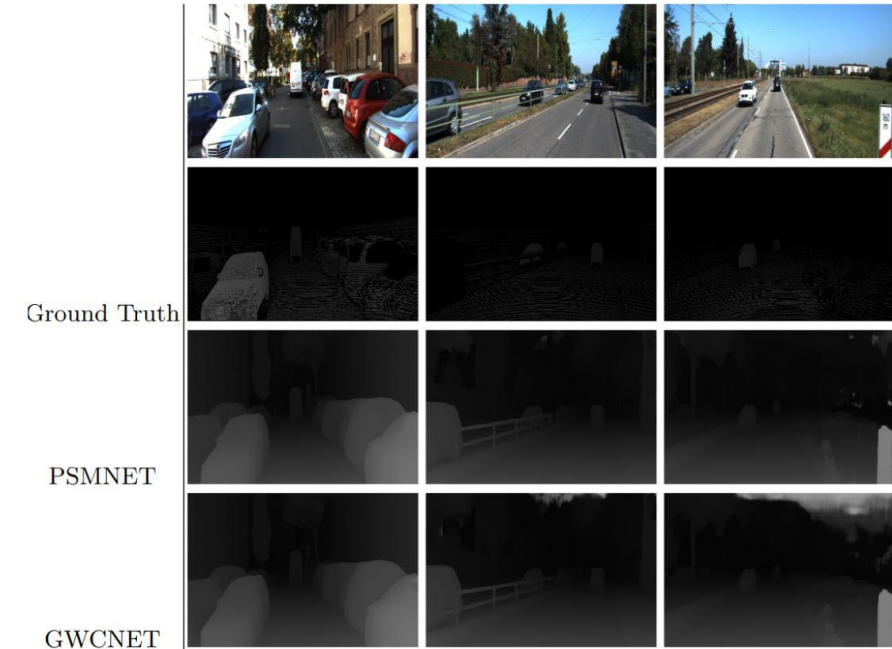
# Evaluated Stereoscopic Image Depth Estimation Methods

## ▶ PSMNet<sup>1</sup>

- **Supervised** stereoscopic depth estimation
- Use **Spatial Pyramid Pooling** to expand the receptive field
- Global contextual informations are extracted by a stacked hourglass 3D **CNN**

## ▶ GwcNet<sup>2</sup>

- **Supervised** stereoscopic depth estimation
- Group wise correlation are used for providing matching features
- Improved stacked hourglass network



PSMNET and GWCNET

1. PSMNET: Pyramid Stereo Matching Network.  
2. GWCNet: Group-wise Correlation Stereo Network

1. J.-R. Chang and Y.-S. Chen, "Pyramid stereo matching network," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5410–5418, 2018.  
2. X. Guo, K. Yang, W. Yang, X. Wang, and H. Li, "Group-wise correlation stereo network," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3273–3282, 2019.

# Experimental Results

## ▶ Dataset : KITTI<sup>1</sup>

- Highly used dataset for road environment with real world traffic
- Data are calibrated and synchronized and come from multiple sensors : Stereo camera, Velodyne Lidar, GPS/IMU navigation system
- Used as benchmark for depth estimation, optical flow, and object detection tasks

## ▶ Results

	RelErr	SqRel	RMSE	logRMSE
Monodepth2	0.115	0.882	4.701	0.190
BTS	<b>0.060</b>	<b>0.249</b>	<b>2.798</b>	<b>0.096</b>

Results of single image-based methods

	RelErr	SqRel	RMSE	logRMSE
GWCNET	<b>0.018</b>	<b>0.048</b>	<b>0.981</b>	<b>0.042</b>
PSMNET	0.032	0.061	1.139	0.056

Results of stereoscopic-based methods

1. A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," The International Journal of Robotics Research, vol. 32, no. 11, pp. 1231–1237, 2013.

# Experimental Results

## ► Input data of our Object Distance

### Evaluation Protocol:

- Input image fed to the depth prediction algorithm
- Disparity map
- Normalized depth map after median scaling
- Object masks from Mask-RCNN



# Experimental Results under KITTI dataset

- ▶ **Monocular Depth Estimation Methods Evaluation: Monotdepth2 vs BTS:**
  - ▶ Depth errors for distance rangers of 10m and up to 80m

Distance ranges	RE		SRE		RMSE		logRMSE		$\alpha_1$		$\alpha_2$		$\alpha_3$	
	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS
0 – 80m	0.115	<b>0.060</b>	0.882	<b>0.249</b>	4.701	<b>2.798</b>	0.190	<b>0.096</b>	0.879	<b>0.955</b>	0.961	<b>0.993</b>	0.982	<b>0.998</b>
0 – 10m	0.102	<b>0.071</b>	0.503	<b>0.188</b>	1.489	<b>0.991</b>	0.141	<b>0.106</b>	0.929	<b>0.959</b>	0.979	<b>0.988</b>	0.99	<b>0.994</b>
10 – 20m	0.116	<b>0.088</b>	0.845	<b>0.395</b>	3.035	<b>2.198</b>	0.18	<b>0.149</b>	0.891	<b>0.924</b>	0.96	<b>0.971</b>	0.979	<b>0.985</b>
20 – 30m	0.168	<b>0.13</b>	1.866	<b>1.055</b>	6.208	<b>4.745</b>	0.261	<b>0.229</b>	0.773	<b>0.836</b>	0.916	<b>0.934</b>	0.957	<b>0.964</b>
30 – 40m	0.196	<b>0.16</b>	2.788	<b>1.945</b>	9.11	<b>7.476</b>	0.307	<b>0.279</b>	0.694	<b>0.764</b>	0.886	<b>0.906</b>	0.942	<b>0.947</b>
40 – 50m	0.209	<b>0.174</b>	3.504	<b>2.64</b>	11.682	<b>10.008</b>	0.318	<b>0.298</b>	0.641	<b>0.725</b>	0.865	<b>0.889</b>	<b>0.943</b>	0.941
50 – 60m	0.221	<b>0.19</b>	4.394	<b>3.739</b>	14.252	<b>12.852</b>	0.332	<b>0.326</b>	0.583	<b>0.675</b>	0.857	<b>0.868</b>	<b>0.927</b>	0.922
60 – 70m	0.212	<b>0.201</b>	4.657	<b>4.584</b>	15.855	<b>15.585</b>	<b>0.325</b>	0.334	0.609	<b>0.619</b>	0.854	<b>0.856</b>	<b>0.93</b>	0.923
70 – 80m	<b>0.181</b>	0.214	<b>4.34</b>	5.454	<b>15.8</b>	18.219	<b>0.284</b>	0.333	<b>0.652</b>	0.548	<b>0.873</b>	0.843	<b>0.945</b>	0.925

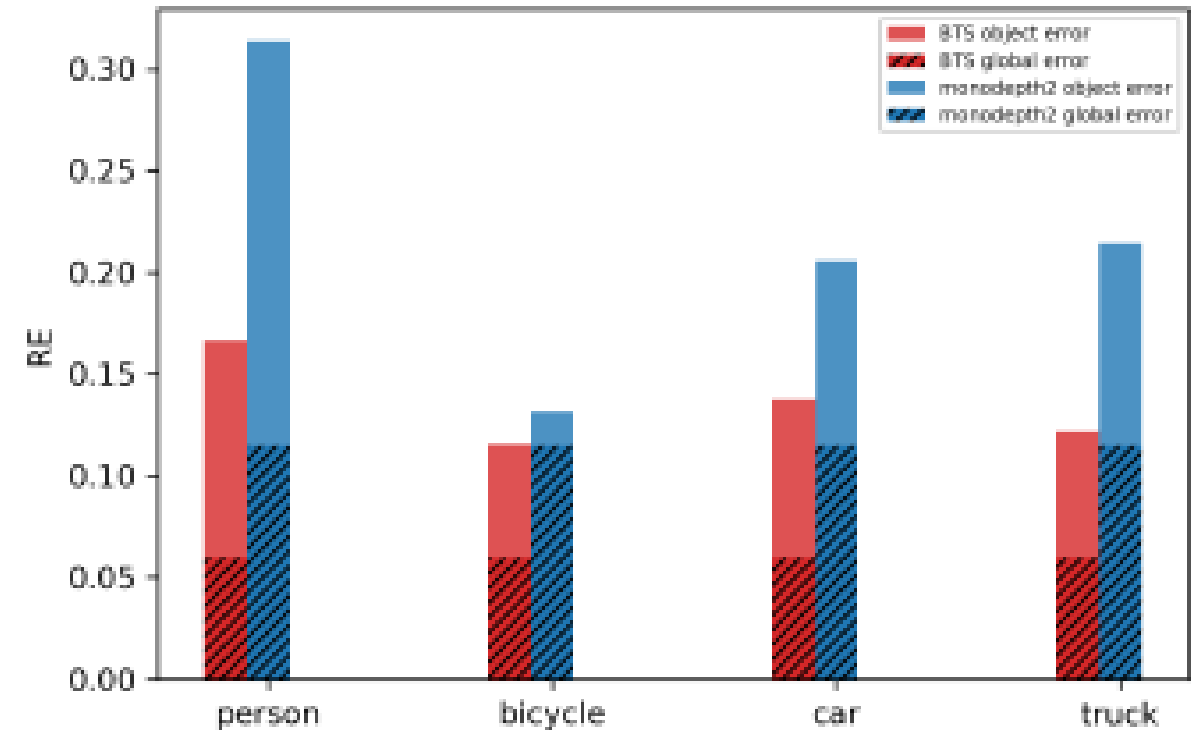
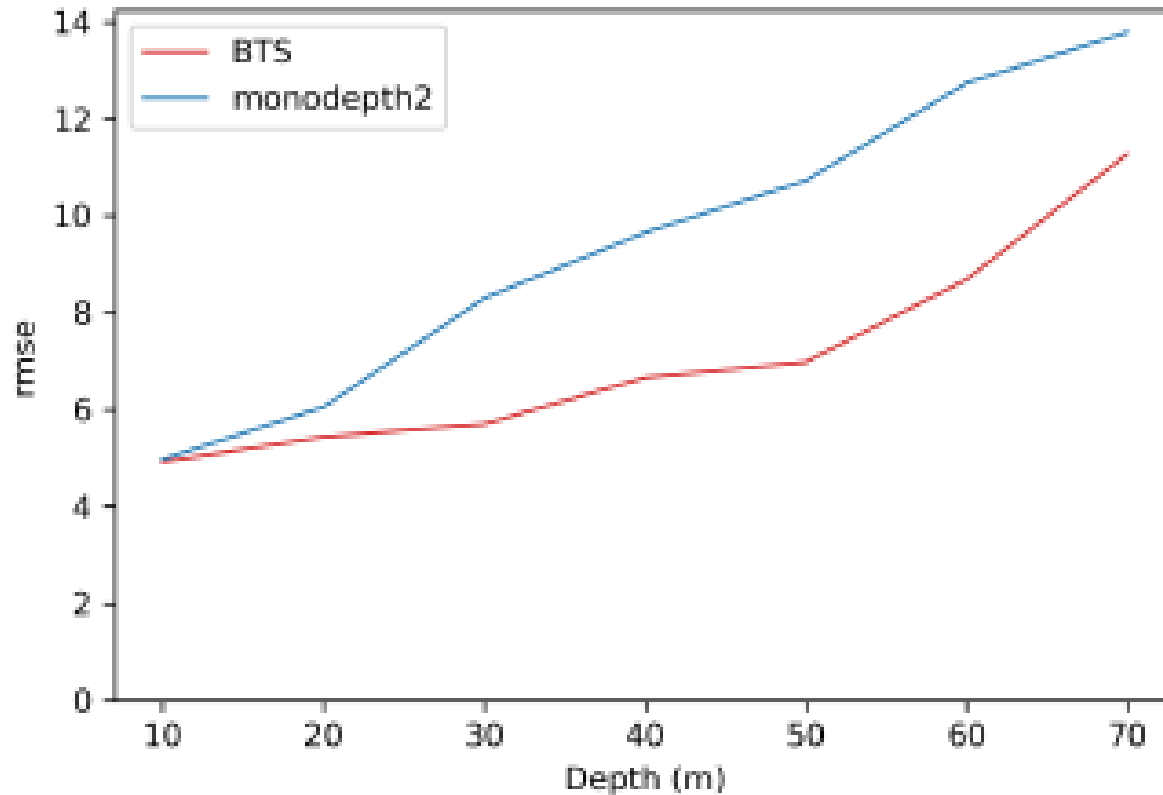
# Experimental Results Under KITTI dataset

- ▶ **Monocular Depth Estimation Methods Evaluation: Monotdepth2 (MD2) vs BTS:**
  - ▶ Depth errors computed for the object classes with enough instance in test split distance rangers of 10m and up to 80m

Object class	RE		SRE		RMSE		logRMSE		$a_1$		$a_2$		$a_3$	
	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS	MD2	BTS
Person	0.314	<b>0.166</b>	5.721	<b>1.786</b>	8.43	<b>5.892</b>	0.326	<b>0.253</b>	0.601	<b>0.772</b>	0.829	<b>0.894</b>	0.92	<b>0.947</b>
Bicycle	0.131	<b>0.116</b>	0.517	<b>0.467</b>	2.81	<b>2.669</b>	0.172	<b>0.163</b>	0.829	<b>0.839</b>	<b>0.964</b>	0.962	0.993	<b>0.994</b>
Car	0.206	<b>0.137</b>	3.132	<b>1.491</b>	7.924	<b>6.052</b>	0.271	<b>0.223</b>	0.773	<b>0.838</b>	0.883	<b>0.922</b>	0.938	<b>0.955</b>
Truck	0.215	<b>0.122</b>	2.769	<b>0.826</b>	6.978	<b>4.523</b>	0.259	<b>0.177</b>	0.694	<b>0.854</b>	0.903	<b>0.969</b>	0.964	<b>0.985</b>

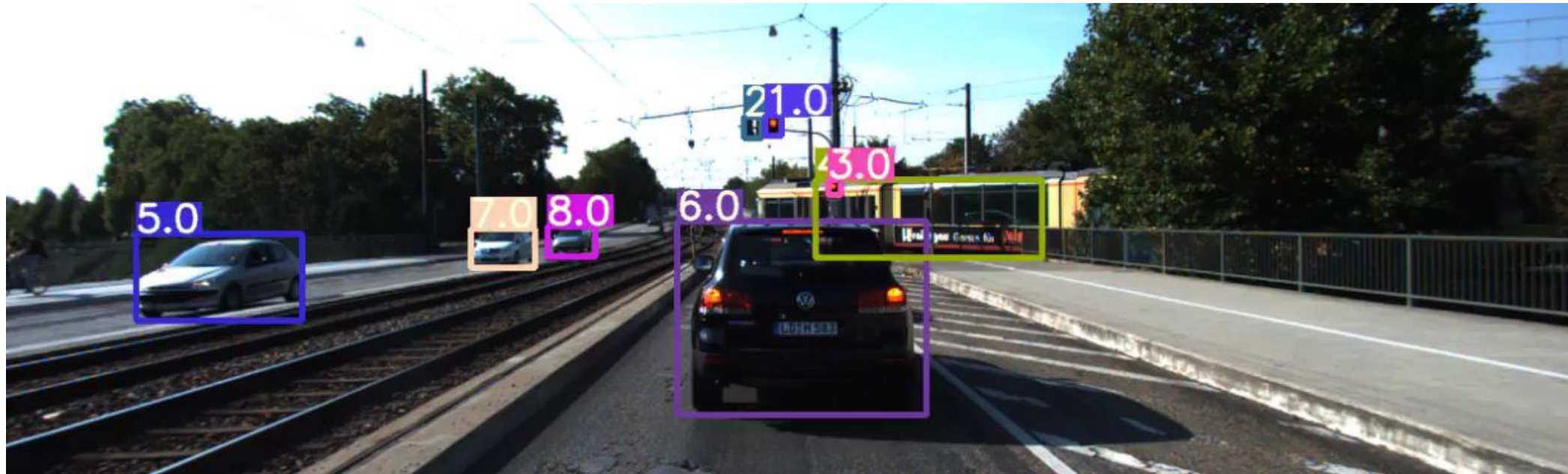
# Experimental Results under KITTI dataset

- ▶ **Quantitative RMSE and RE Results for the car object class over distance ranges on the KITTI dataset.**



# Experimental results

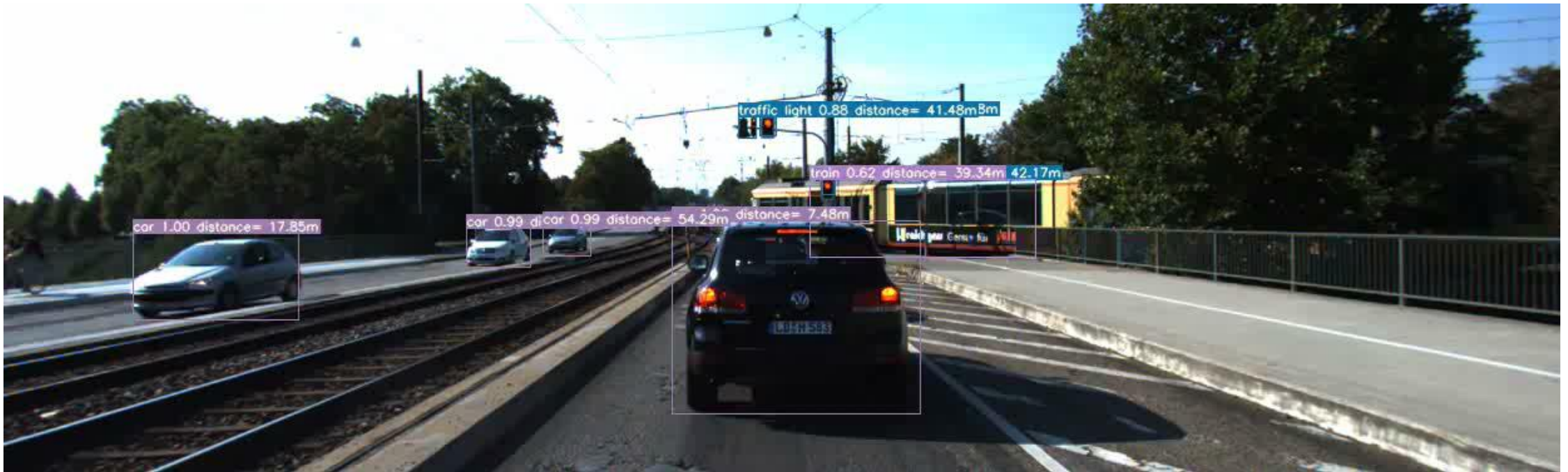
## ► Object Detection





# Experimental results

## ► Distance Estimation of in Traffic Environment under KITTI



# Experimental results

## ► Depth Estimation Evaluation Methods: Evaluation



# Conclusion and Future Work

## ▶ In this work we presented

- An evaluation of State-of-The-Art for both:
  - Stereo and Mono depth estimation methods

## ▶ Our work also showed that

- **BTS** is more accurate than **Monodepth2**
- **GWCnet** outperforms **PSMNet**
- **Stereoscopic** methods have a **greater accuracy** than **monocular-based** methods

## ▶ Future work

- Evaluate depth estimation algorithms on **Rail Environment**
- Acquire our own **Railway Dataset**