

ANALYSIS OF THE INFLUENCE OF TRAINING DATA ON ROAD USER DETECTION

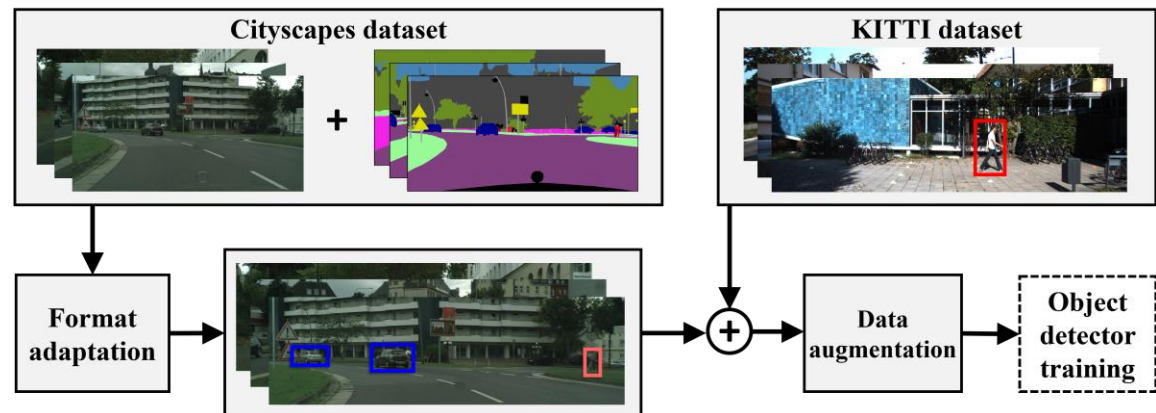
Carlos Guindel, David Martín, José María Armingol, and Christoph Stiller



20th IEEE International Conference on Vehicular Electronics and Safety
Madrid · 12 September 2018

Agenda

- Motivation and goals
- Experimental setup
- Analysis
- Conclusion



Motivation

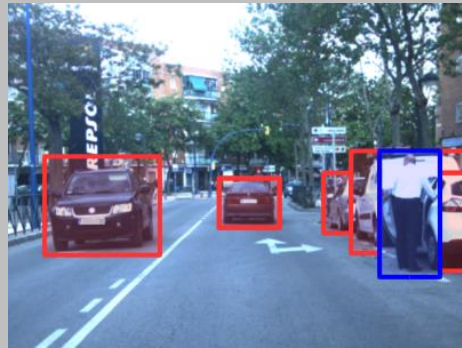
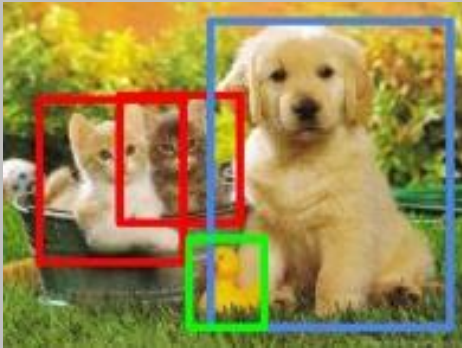
Object detection



Computer vision



Autonomous driving



**Instance
segmentation
(e.g., Mask R-CNN)**

**Deep
Learning**



Data



Motivation

Object detection



Computer vision



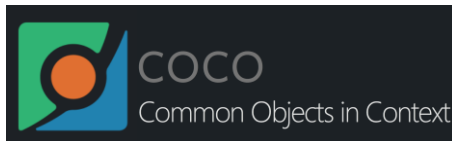
Autonomous driving

Deep
Learning



IMAGENET

450 000+ images
200 categories



200 000+ images
80 categories

...

**The KITTI Vision
Benchmark Suite**

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago

7 481 images
9 categories



CITYSCAPES
DATASET

2 975 images
10 categories

Data



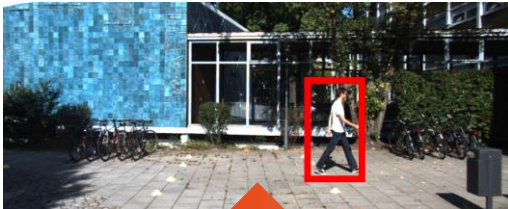
Motivation

Object detection

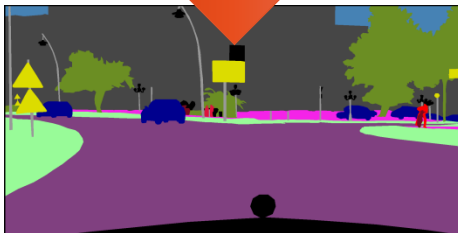


Autonomous driving

Deep
Learning



Different
labels



The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago

7 481 images
9 categories



CITYSCAPES
DATASET

2 975 images
10 categories

Data



Goals



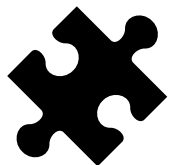
Research is often narrow-focused on the development of new architectures and models.

R-FCN, SSD; ResNet, Inception, MobileNet,...

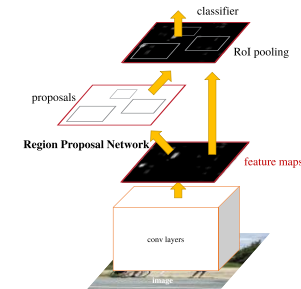
Instead, we investigate:



The improvement provided by introducing additional samples into the training process.



The possibility of using heterogeneous labels in a multi-task learning method.



Faster R-CNN
State-of-the-art object
detection meta-architecture

Datasets



The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago

Training

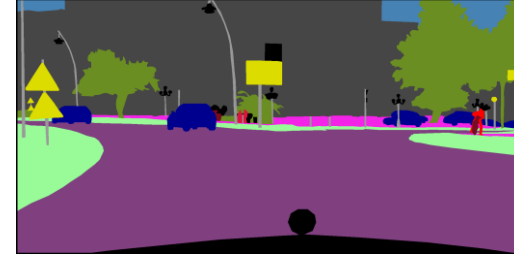


3,712 images

Validation



3,769 images



CITYSCAPES
DATASET

Training



2,975 images

X. Chen, K. Kundu, Y. Zhu, A. Berneshawi, H. Ma, S. Fidler, and R. Urtasun, "3D Object Proposals for Accurate Object Class Detection," in Advances in Neural Information Processing Systems (NIPS), 2015, pp. 424–432.

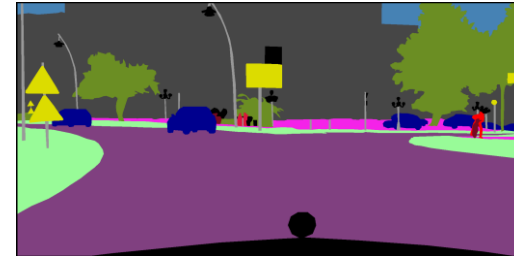
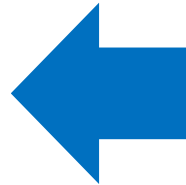
category	KITTI			Cityscapes
	train	val	total	total
Car	10 753	10 963	21 716	21 637
Pedestrian	2 104	2 172	4 276	15 788
Cyclist	594	600	1 194	1 481

Adaptation



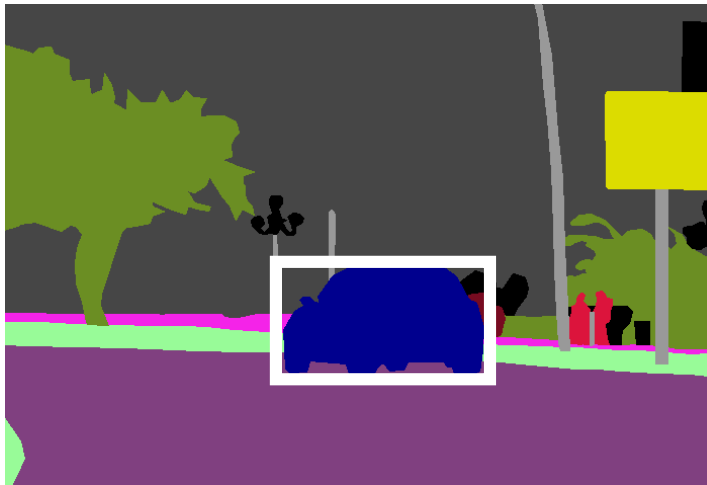
The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago



 **CITYSCAPES**
DATASET

① Semantic labeling to bounding boxes



Instantiable objects

Minimum enclosing box

Categories:

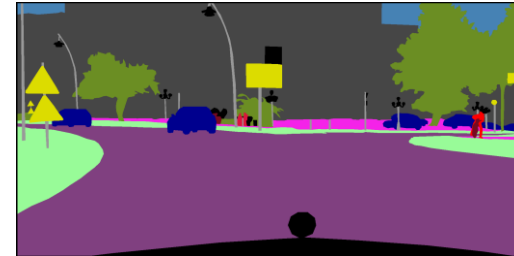
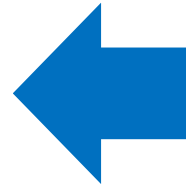
- Person → Pedestrian
- Rider + Bicycle → Cyclist

Adaptation



The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago



 **CITYSCAPES**
DATASET

② Occlusion & truncation



Occlusion

Cityscapes labels contain foreground-background ordering

$$\text{Degree of occlusion} = \frac{\text{intersection}}{\text{background}}$$

Truncation

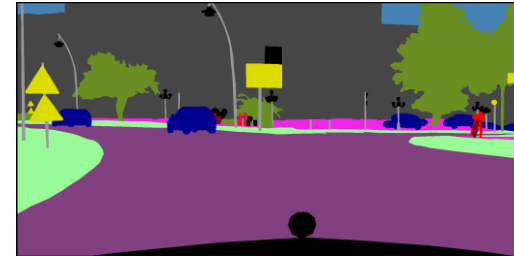
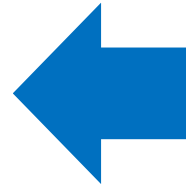
Whenever any of the sides of the b. box coincides with the image boundaries

Adaptation



The KITTI Vision
Benchmark Suite

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago



 CITYSCAPES
DATASET

③ Resolution/FOV



2048 × 1024



2048 × 620

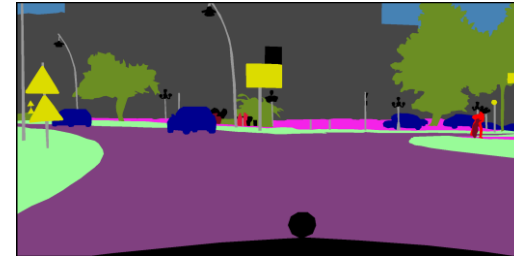
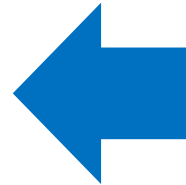
Removes the hood and the
Mercedes-Benz emblem

Adaptation



The KITTI Vision
Benchmark Suite

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago



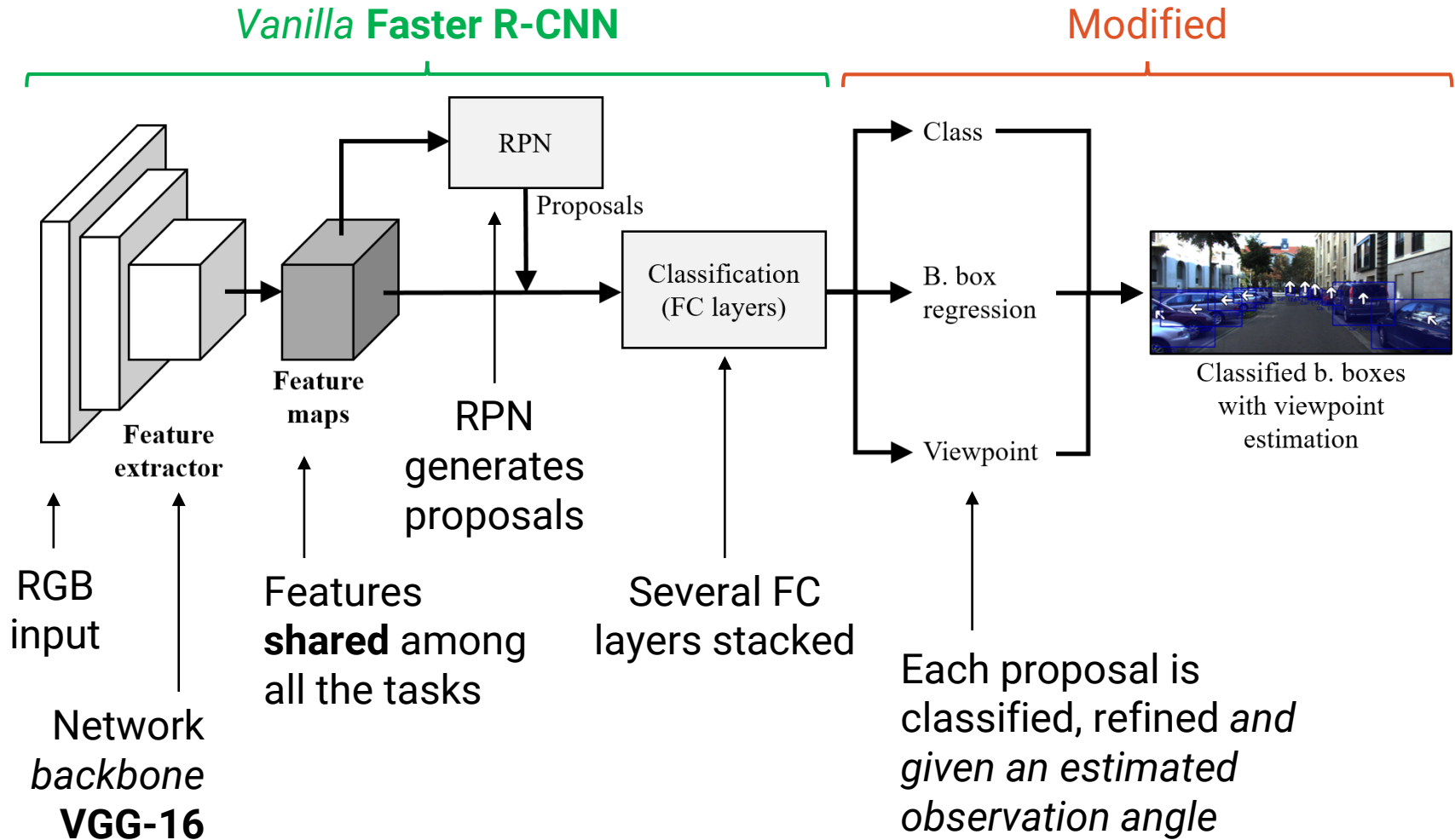
CITYSCAPES
DATASET

④ Difficulty levels

We ignore samples not meeting the KITTI's **Hard** level requirements:

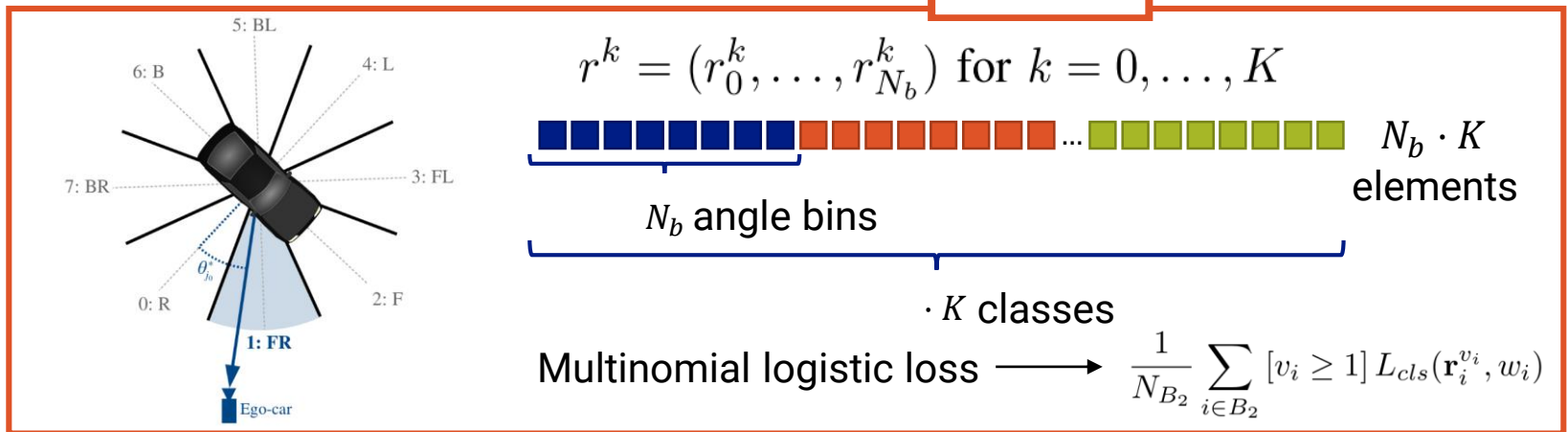
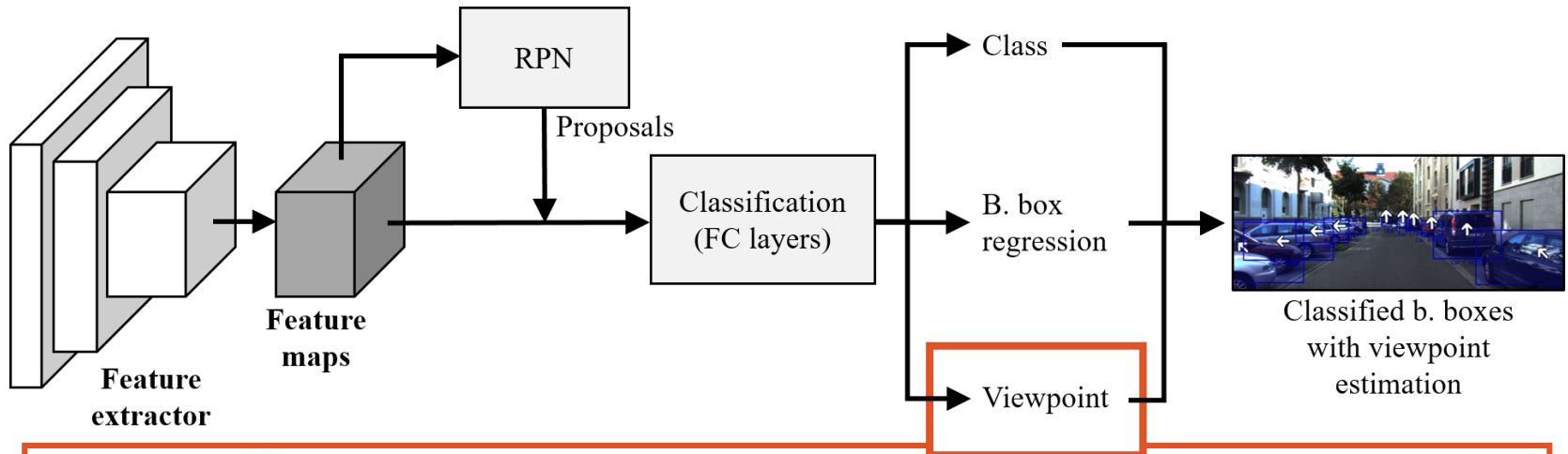
- Larger than **25 pixels**
- Max **occlusion**: “Difficult to see”: level 2 (KITTI) or 75% (Cityscapes)
- Maximum **truncation**: 50% (KITTI) or no truncation (Cityscapes)

Object detection method



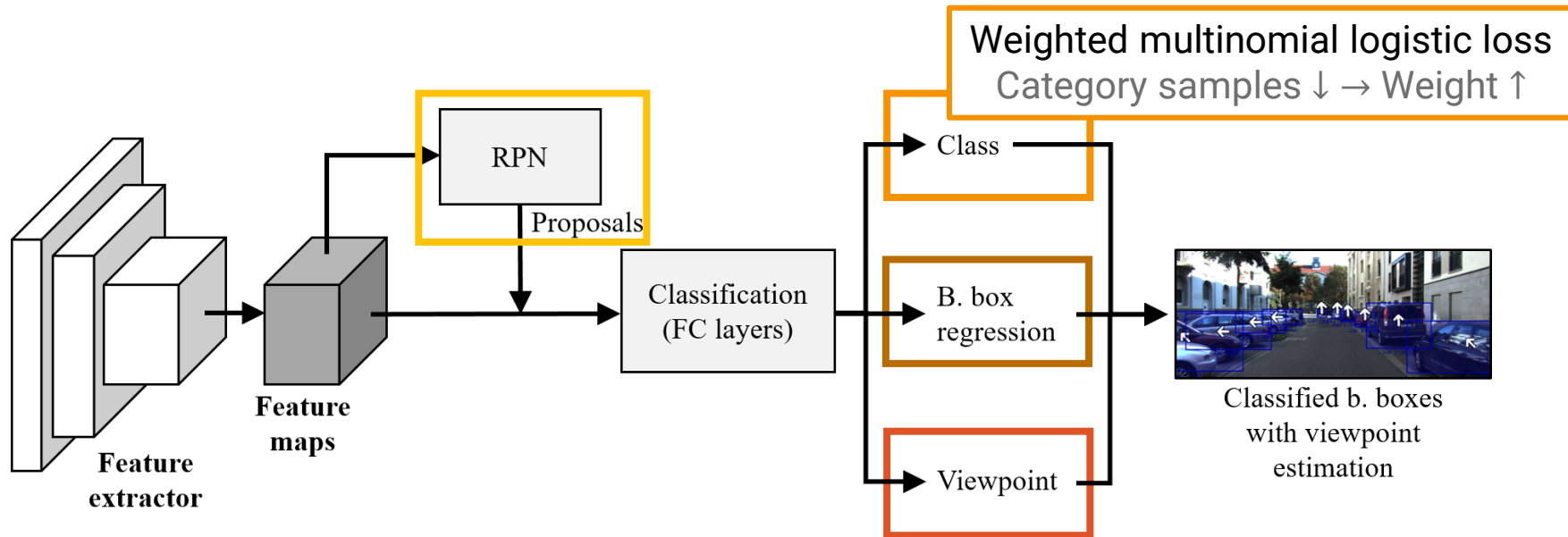
C. Guindel, D. Martin, and J. M. Armingol, "Fast Joint Object Detection and Viewpoint Estimation for Traffic Scene Understanding," accepted for publication in: Intelligent Transportation Systems Magazine.

Object detection method



C. Guindel, D. Martin, and J. M. Armingol, "Fast Joint Object Detection and Viewpoint Estimation for Traffic Scene Understanding," accepted for publication in: Intelligent Transportation Systems Magazine.

Multi-task loss

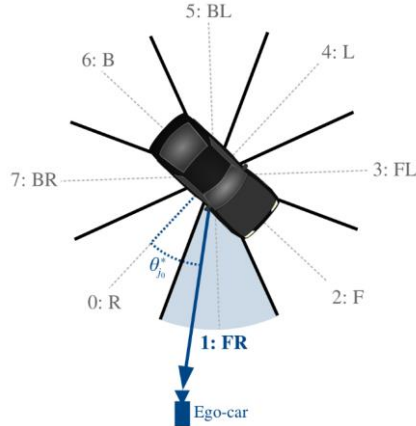
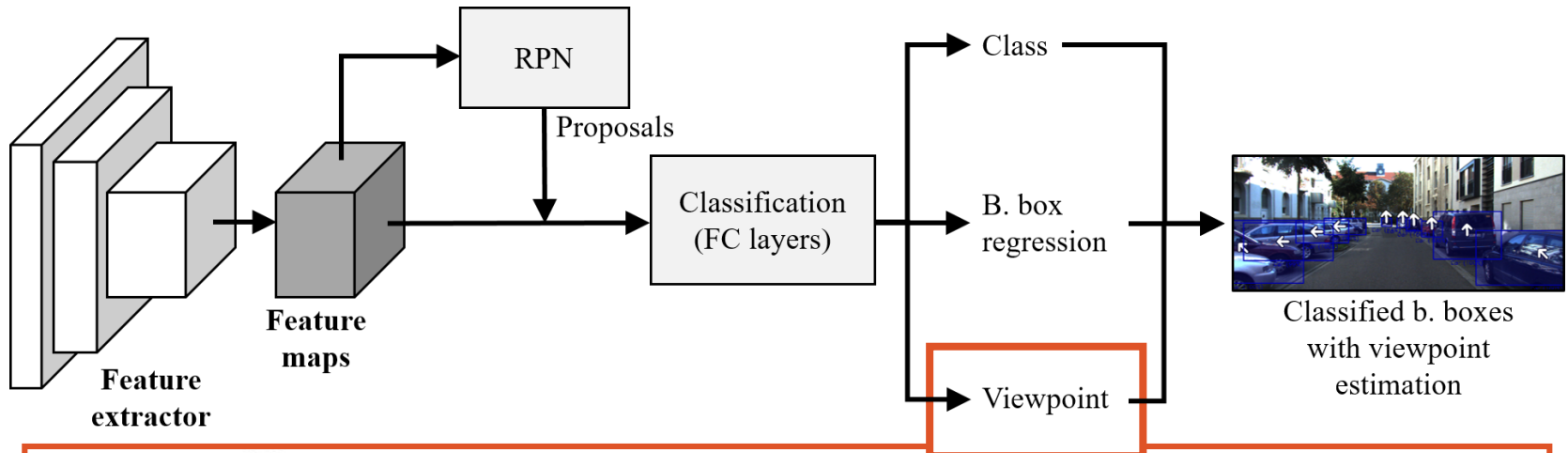


$$\text{Loss} = \text{Objectness} + \text{Proposal regr.} + \text{Class} + \text{B.box regr.} + \text{Viewpoint}$$

$$\text{Multinomial logistic loss} \longrightarrow \frac{1}{N_{B_2}} \sum_{i \in B_2} [v_i \geq 1] L_{cls}(\mathbf{r}_i^{v_i}, w_i)$$

C. Guindel, D. Martin, and J. M. Armingol, "Fast Joint Object Detection and Viewpoint Estimation for Traffic Scene Understanding," accepted for publication in: Intelligent Transportation Systems Magazine.

Multi-task loss



The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago

Viewpoint loss $\neq 0$

Multinomial logistic loss $\longrightarrow \frac{1}{N_{B_2}} \sum_{i \in B_2} [v_i \geq 1] L_{cls}(\mathbf{r}_i^{v_i}, w_i)$



CITYSCAPES
DATASET

No viewpoint/orientation annotations

Viewpoint loss = 0

Assessment method

Evaluation metrics

Average precision (AP)



Assess object detection

$$AP = \frac{1}{N} \sum_r \tilde{p}(r) \quad \tilde{p}(r) = \max_{\tilde{r}: \tilde{r} > r} p(\tilde{r})$$

Precision Recall

Average orientation similarity (AOS)



Assess object detection
AND orientation

$$AOS = \frac{1}{N} \sum_r \tilde{s}(r) \quad \tilde{s}(r) = \max_{\tilde{r}: \tilde{r} > r} s(\tilde{r})$$

Cosine similarity Recall

Common Testbed



The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago

Validation set



3,769 images

Training Parameters

- 1- image batch
- SGD, initial lr = 0.001
- Step decay schedule
0.1 × every 50k iterations
- 80k iterations


Experiment 1: Combined datasets

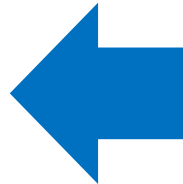


**The KITTI Vision
Benchmark Suite**

A project of Karlsruhe Institute of Technology
and Toyota Technological Institute at Chicago

Training

 3,712 images



**Add
training
samples**




Adapted



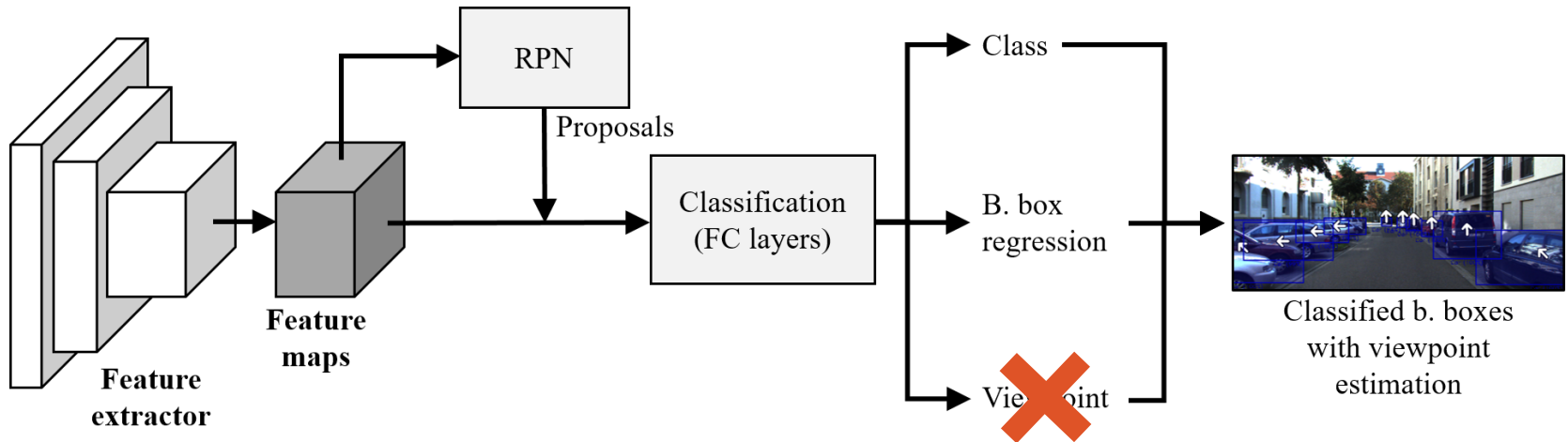
**CITYSCAPES
DATASET**

Training

 2,975 images

- In each training iteration, an image (single batch) is randomly chosen from a mix of both datasets.
- Tests without (a) and with (b) viewpoint estimation branch

Experiment 1: Combined datasets (a)



Detection

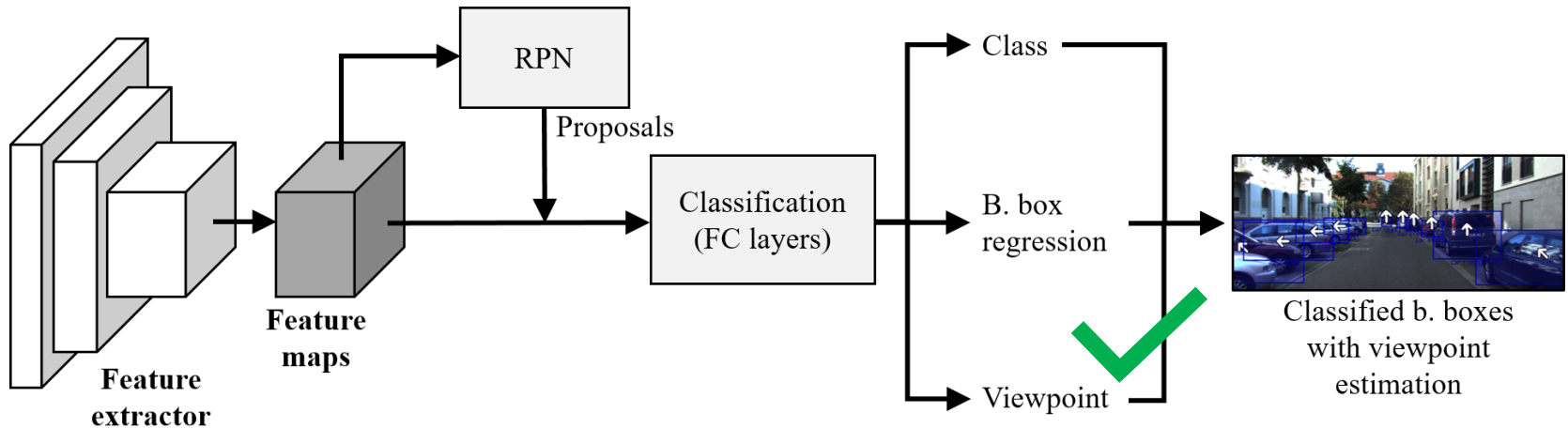
category	tr. data C	Easy	Mod.	Hard
Car	KITTI	90.05	79.32	70.04
	Cityscapes	81.37	63.66	53.47
	KITTI + CS	90.31	84.94	70.33
Pedestrian	KITTI	75.80	67.17	58.58
	Cityscapes	72.00	63.92	55.33
	KITTI + CS	77.77	68.72	60.05
Cyclist	KITTI	77.47	56.96	54.64
	Cityscapes	63.09	50.14	46.85
	KITTI + CS	82.90	62.50	58.05

+5.62 AP

+1.55 AP

+5.54 AP

Experiment 1: Combined datasets (b)



Detection +
Orientation

Two alternative strategies:

- ① Pick images from a **KITTI+Cityscapes** mix
Viewpoint = 0 when a Cityscapes sample is chosen
- ② Pre-train with **Cityscapes**, fine-tune with **KITTI**

Train without viewpoint branch and transfer weights to the complete model

Experiment 1: Combined datasets (b)



Detection +
Orientation

① Pick images from a **KITTI+Cityscapes** mix

② Pre-train with **Cityscapes**, fine-tune with **KITTI**

category	tr. data C	Detection (AP)			Orientation (AOS)			
		Easy	Mod.	Hard	Easy	Mod.	Hard	
Car	KITTI	90.01	79.03	69.67	88.26	77.35	67.97	+5.6 AP +1.6 AOS
	① KITTI + CS	90.39	84.59	70.21	88.68	82.79	68.57	
	② KITTI (w. CS pret.)	90.33	86.16	70.58	88.63	84.43	69.01	
Pedestrian	KITTI	71.19	64.05	55.75	65.31	57.62	50.01	+3.9 AP +2 AOS
	① KITTI + CS	76.32	67.98	59.11	67.83	59.65	51.69	
	② KITTI (w. CS pret.)	74.54	66.01	57.68	67.33	59.01	51.52	
Cyclist	KITTI	77.33	54.87	52.89	69.73	48.79	47.06	+13.6 AP +12.4 AOS
	① KITTI + CS	86.11	68.49	63.46	77.66	61.23	56.83	
	② KITTI (w. CS pret.)	83.18	60.37	57.35	75.55	54.36	51.74	

Experiment 2: Can we get rid of ImageNet?



Pre-training with ImageNet (generalist dataset) generates good initial weights



init.	tr. data	Detection (mAP)			Orientation (mAOS)		
		Easy	Mod.	Hard	Easy	Mod.	Hard
Yes	KITTI	79.51	65.98	59.44	74.43	61.25	55.02
No	K. + CS	53.80	42.99	37.25	47.93	39.13	33.00

-22.99 mAP

-22.12 mAOS

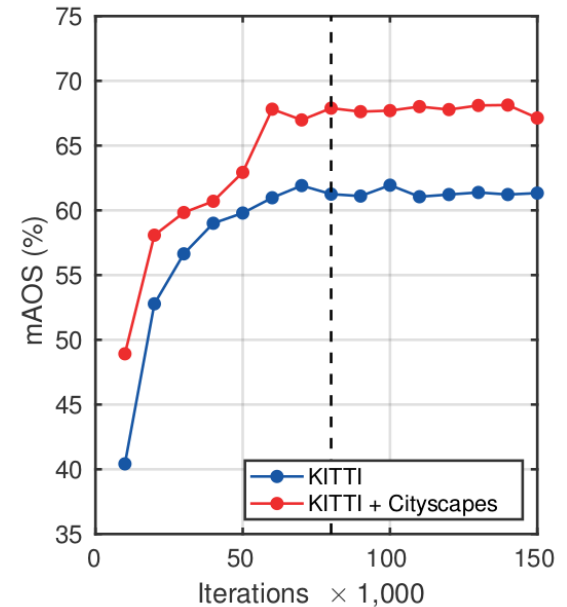
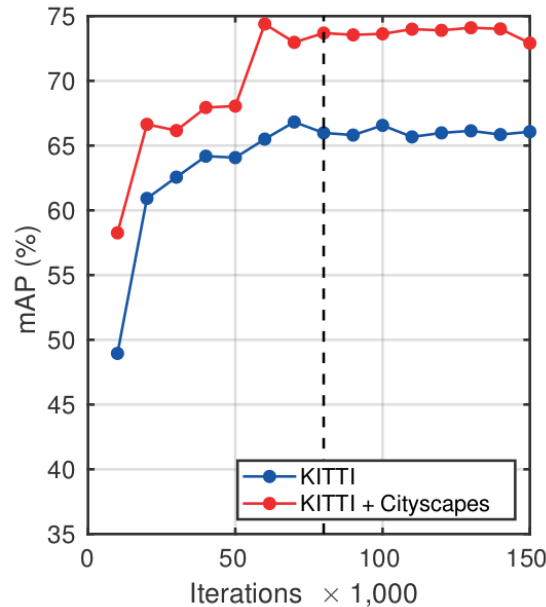


Initialization with a large dataset is still an essential requirement to achieve a proper generalization ability

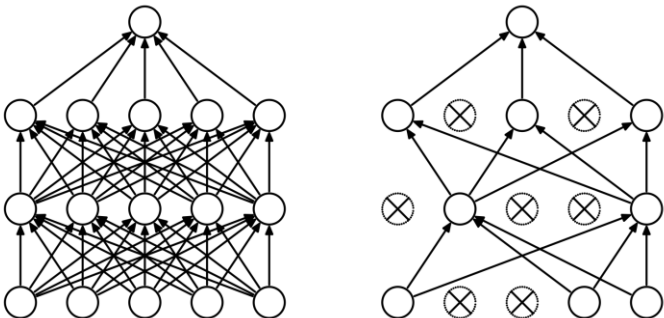
Experiment 3: Overfitting

Performance on the validation set vs # of iterations

✓ No symptoms of overfitting



Dropout? $p = 0.5$



dropout	Detection (mAP)			Orientation (mAOS)		
	Easy	Mod.	Hard	Easy	Mod.	Hard
No	79.51	65.98	59.44	74.43	61.25	55.02
Yes	79.20	65.34	58.43	73.77	60.73	54.16

-0.64 mAP

-0.52 mAOS

✗ No apparent benefit

Experiment 4: Missing labels



CITYSCAPES
DATASET



Detection +
Orientation
+6.64 mAOS



Orientation
?

**Observation angle
annotations**



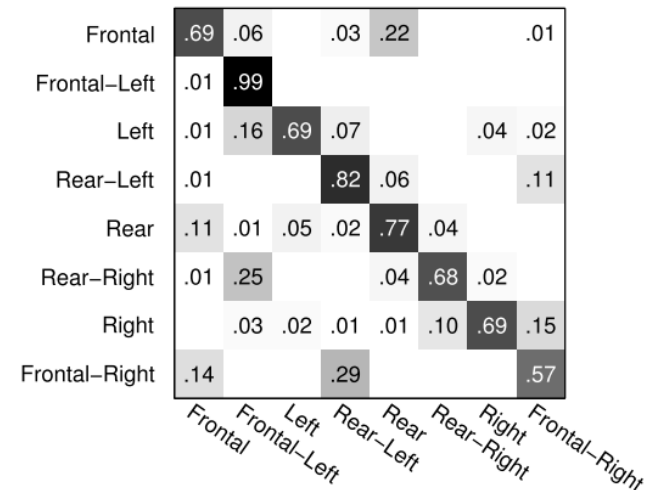
**No observation
angle annotations**



Additional measure

Mean precision in pose estimation
(MPPE)

- Discrete viewpoint estimation is a classification problem
- MPPE is the mean of the elements on the main diagonal of the confusion matrix (correctly predicted viewpoint bins)



R. J. López-Sastre, T. Tuytelaars, and S. Savarese, "Deformable part models revisited: A performance evaluation for object category pose estimation," in Proc. IEEE International Conference on Computer Vision (ICCV), 2011, pp. 1052–1059.

Experiment 4: Missing labels



CITYSCAPES
DATASET



Detection +
Orientation
+6.64 mAOS



Orientation
+1,52 MPPE

**Observation angle
annotations** ✓

**No observation
angle annotations** ✗

category	tr. data	Easy	Mod.	Hard
Car	KITTI	92.24	80.93	69.29
	KITTI + CS	92.13	83.40	71.72
Pedestrian	KITTI	59.03	51.02	43.71
	KITTI + CS	57.74	51.35	44.41
Cyclist	KITTI	70.71	49.84	49.00
	KITTI + CS	64.95	51.60	49.11

+2.47 MPPE

+0.33 MPPE

+1.76 MPPE

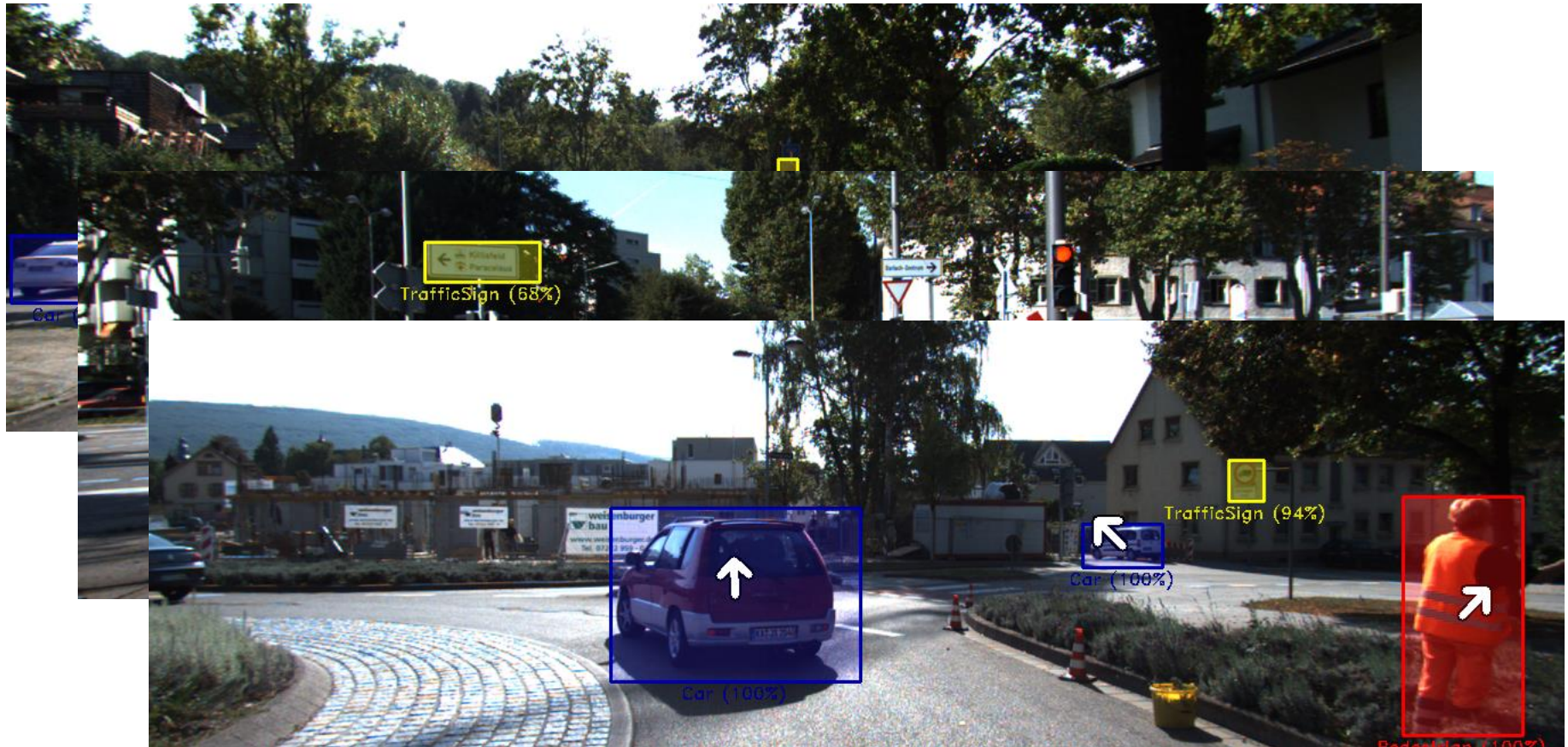
**Missing labels does
not hurt orientation
estimation
performance** ✓

Experiment 5: Mixed labels

Including all categories

Car, Truck, Pedestrian, Cyclist, Train, Traffic Sign

Cityscapes-only



Experiment 6: Data augmentation

Horizontal flip + texture augmentations

<https://github.com/aleju/imgaug>

Choose a random subset of them, between 0 and 4

Add

$[-40, 40]$



Multiplication

$[0.5, 1.5]$



Gaussian noise

$\mathcal{N}(0, 5.1^2)$



Saturation

$[-20, 20]$ (H, S)



Experiment 6: Data augmentation

Two separated experiments

① Only **KITTI**

To assess the overall effect of the augmentation techniques

② **Cityscapes + KITTI**

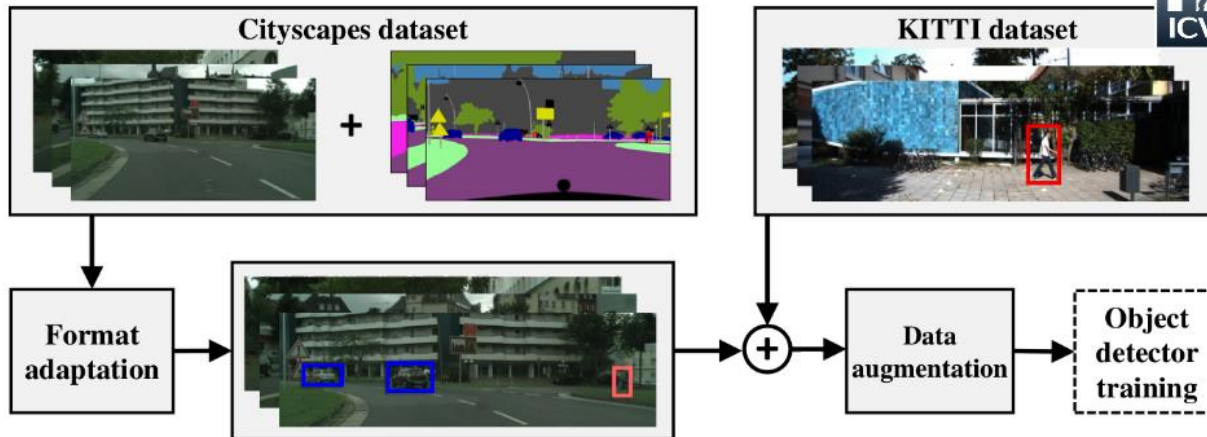
Augmentation could help mitigate the difference between both sets of images

	tr. data	aug.	Detection (mAP)			Orientation (mAOS)			
			Easy	Mod.	Hard	Easy	Mod.	Hard	
①	K.	No	79.51	65.98	59.44	74.43	61.25	55.02	✗ No apparent benefit
		Yes	80.39	65.87	58.96	74.56	61.00	54.38	
②	K. + CS	No	84.27	73.69	64.26	78.06	67.89	59.03	~ Limited benefit
		Yes	83.96	74.14	65.16	77.95	68.09	59.59	

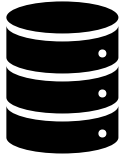
Comparison

Analysis of the Influence of Training Data on Road User Detection

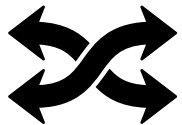
C. Guindel, D. Martín, J. M. Armingol and C. Stiller



Conclusion



Modestly enhancing the training data can lead to notable improvements on the results obtained by a CNN object detector



The variability introduced by Cityscapes samples can achieve a non-negligible improvement, even when evaluated on the KITTI dataset



Results pave the way for future works taking advantage of multiple data sources

THANK YOU



20th IEEE International Conference on Vehicular Electronics and Safety
Madrid · 12 September 2018