Sensor Fusion for Semantic Segmentation of Urban Scenes

Richard Zhang¹, Stefan A. Candra¹, Kai Vetter¹², Avideh Zakhor¹

¹University of California, Berkeley

²Lawrence Berkeley National Laboratory

Presented at ICRA, May 2015





Goal

- Extract semantic information by fusing multiple modalities
 - Image
 - LiDAR scan
- Challenges
 - Incorporate information from multiple scales
 - Fuse information from multiple modalities with non-overlapping sensor coverage





building , sky , road , vegetation , sidewalk car , pedestrian , cyclist , sign/pole , fence

Top Level Pipeline



Previous Work

	Segmentation	Image Features	Point cloud features	Fusion	
Sengupta et al. ICRA 2013	none	Descriptive	none	majority voting	
He and Upcroft. IROS 2013	Single segmentation image only	Descriptive	none	of image classifications in 3D	

Previous Work

	Segmentation	Image Features	Point cloud features	Fusion	
Sengupta et al. ICRA 2013	none	Descriptive	none	majority voting	
He and Upcroft. IROS 2013	Single segmentation image only	Descriptive	none	classifications in 3D	
Cadena and Košecká. ICRA 2014.	Single segmentation image only	Simple	Simple	Early fusion	

Previous Work

	Segmentation	Image Features	Point cloud features	Fusion
Sengupta et al. ICRA 2013	none	Descriptive	none	majority voting
He and Upcroft. IROS 2013	Single segmentation image only	Descriptive	none	classifications in 3D
Cadena and Košecká. ICRA 2014.	Single segmentation image only	Simple	Simple	Early fusion
Ours	<i>multiple</i> segmentations <i>both</i> domains	Descriptive	Descriptive	Late fusion



Segmentation

- Treating each pixel/point separately is computationally intensive
- Under-segmentation & Oversegmentation errors
- ➔ Multiple segmentations performed to get superpixels/supervoxels







Segmentation

- Treating each pixel/point separately is computationally intensive
- Under-segmentation & Oversegmentation errors
- ➔ Multiple segmentations performed to get superpixels/supervoxels





Feature Extraction

Туре	Name	Dim	Low	High
	Area	1	 ✓ 	✓
	Equivalent Diameter	1	~	✓
Size/Shape	Major/minor axes	2	✓	✓
	Orientation	1	~	~
	Eccentricity	1	✓	✓
Dosition	(x, y) - min, mean, max	6	 ✓ 	√
rostuon	superpixel mask (8x8)	64	✓	✓
Color	rgb+lab (mean, std)	6	✓	 ✓
COIOI	rgb+lab (histogram)	48	✓	√
High-dim	SIFT BoW	400	 ✓ 	
	contextual rgb+lab (mean, std)	6	 ✓ 	
Contextual	contextual rgb+lab (histogram)	48	 ✓ 	
	contextual SIFT BoW	400	 ✓ 	

Туре	Name	Dim	Low	High
	Length proxy - λ_1	1	 ✓ 	✓
Size	Area proxy - $\sqrt{\lambda_1 \lambda_2}$	1	 ✓ 	✓
	Volume proxy - $\sqrt[3]{\lambda_1 \lambda_2 \lambda_3}$	1	✓	✓
	Scatter - λ_3/Λ	1	✓	✓
Shape	Planarity - $(\lambda_2 - \lambda_3)/\Lambda$	1	✓	✓
	Linearity - $(\lambda_1 - \lambda_2)/\Lambda$	1	~	 ✓
Position	$z - z_{gndplane}$ - min, mean, max	3	✓	✓
Orientation	Verticalness - v_{1z}	1	 ✓ 	v
onentation	Horizontalness - $\sqrt{1-v_{1z}^2}$	1	✓	✓
High-dim	Spin image BoW	1000	 ✓ 	

Image superpixel features

Point cloud supervoxel features





Early Fusion



Early Fusion



Early Fusion



Late Fusion



Late Fusion



building , sky , road , vegetation , sidewalk car , pedestrian , cyclist , sign/pole , fence



Point cloud unimodal segmentation



Image unimodal segmentation



Ground Truth

building , sky , road , vegetation , sidewalk car , pedestrian , cyclist , sign/pole , fence



Point cloud unimodal segmentation (projected onto image)



Image unimodal segmentation



Ground Truth

















- Quantitative Results
 - Train and test on KITTI¹ dataset, augmented with additional annotations
 - 252 images across 8 sequences
 - 140 images for training, 112 for testing

	glob	class	bldg	sky	road	veg	side	car	ped	cycl	sgn	fnc
Cadena et al. [3]	84.1%	52.4%	92.5%	95.7%	92.5%	86.3%	51.5%	67.9%	28.6%	4.0%	2.5%	2.3%

- Quantitative Results
 - Train and test on KITTI¹ dataset, augmented with additional annotations
 - 252 images across 8 sequences
 - 140 images for training, 112 for testing

	glob	class	bldg	sky	road	veg	side	car	ped	cycl	sgn	fnc
Cadena et al. [3]	84.1%	52.4%	92.5%	95.7%	92.5%	86.3%	51.5%	67.9%	28.6%	4.0%	2.5%	2.3%
Ours (image only)	83.5%	53.3%	87.5%	92.5%	94.5%	92.5%	34.5%	71.4%	49.0%	3.6%	4.1%	3.3%

- Quantitative Results
 - Train and test on KITTI¹ dataset, augmented with additional annotations
 - 252 images across 8 sequences
 - 140 images for training, 112 for testing

	glob	class	bldg	sky	road	veg	side	car	ped	cycl	sgn	fnc
Cadena et al. [3]	84.1%	52.4%	92.5%	95.7%	92.5%	86.3%	51.5%	67.9%	28.6%	4.0%	2.5%	2.3%
Ours (image only)	83.5%	53.3%	87.5%	92.5%	94.5%	92.5%	34.5%	71.4%	49.0%	3.6%	4.1%	3.3%
Ours (late fused)	88.0%	64.8%	93.5%	92.5%	91.2%	92.0%	69.7%	76.5%	63.7%	10.0%	16.6%	42.2%

- Quantitative Results
 - Train and test on KITTI¹ dataset, augmented with additional annotations
 - 252 images across 8 sequences
 - 140 images for training, 112 for testing

	glob	class	bldg	sky	road	veg	side	car	ped	cycl	sgn	fnc
Cadena et al. [3]	84.1%	52.4%	92.5%	95.7%	92.5%	86.3%	51.5%	67.9%	28.6%	4.0%	2.5%	2.3%
Ours (image only)	83.5%	53.3%	87.5%	92.5%	94.5%	92.5%	34.5%	71.4%	49.0%	3.6%	4.1%	3.3%
Ours (late fused)	88.0%	64.8%	93.5%	92.5%	91.2%	92.0%	69.7%	76.5%	63.7%	10.0%	16.6%	42.2%
Ours (CRF)	89.3%	65.4%	95.0%	92.6%	92.6%	92.8%	73.3%	78.7%	65.1%	7.3%	13.8%	43.2%

¹Geiger, et al. KITTI. CVPR 12

