

REPORT DOCUMENTATION PAGE

Form Approved OMB NO. 0704-0188

The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA, 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.  
PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.

1. REPORT DATE (DD-MM-YYYY) 21-12-2014		2. REPORT TYPE Final Report		3. DATES COVERED (From - To) 13-Sep-2011 - 12-Sep-2014	
4. TITLE AND SUBTITLE Final Report: Acquiring Semantically Meaningful Models for Robotic Localization, Mapping and Target Recognition-Research Area 5 Computing Science				5a. CONTRACT NUMBER W911NF-11-1-0476	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER 611102	
6. AUTHORS Jana Kosecka				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAMES AND ADDRESSES George Mason University 4400 University Drive, MSN 4C6 Fairfax, VA 22030 -4422				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS (ES) U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211				10. SPONSOR/MONITOR'S ACRONYM(S) ARO	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S) 60054-CS.13	
12. DISTRIBUTION AVAILABILITY STATEMENT Approved for Public Release; Distribution Unlimited					
13. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.					
14. ABSTRACT The goal of this proposal is to develop novel representations and techniques for localization, mapping and target recognition from videos of indoors and urban outdoors environments. The proposed techniques will facilitate enhanced navigation capabilities by means of visual sensing and enable scalable, long-term navigation and target detection in outdoors and indoors environments. The attained representations will also be applicable towards human-robot interaction, enhancement of human navigational and decision making capabilities and provide compact semantically meaningful summaries of the acquired sensory experience. The					
15. SUBJECT TERMS robot perception, computer vision, semantic labeling, segmentation					
16. SECURITY CLASSIFICATION OF:		17. LIMITATION OF ABSTRACT		15. NUMBER OF PAGES	
a. REPORT UU	b. ABSTRACT UU	c. THIS PAGE UU	UU		
				19a. NAME OF RESPONSIBLE PERSON Jana Kosecka	
				19b. TELEPHONE NUMBER 703-993-1876	



## Report Title

Final Report: Acquiring Semantically Meaningful Models for Robotic Localization, Mapping and Target Recognition-Research Area 5 Computing Science

### ABSTRACT

The goal of this proposal is to develop novel representations and techniques for localization, mapping and target recognition from videos of indoors and urban outdoors environments. The proposed techniques will facilitate enhanced navigation capabilities by means of visual sensing and enable scalable, long-term navigation and target detection in outdoors and indoors environments. The attained representations will also be applicable towards human-robot interaction, enhancement of human navigational and decision making capabilities and provide compact semantically meaningful summaries of the acquired sensory experience. The proposed representations will be governed by principles of compositionality, facilitate bottom-up learning, enable efficient inference and could be adapted to a task at hand.

The main novelty of the approach will be the use both 3D and 2D geometric and photometric cues computed either from video sequence or from novel RGB-D cameras, which provide synchronized video and range data at frame rate. Video poses challenges related to more extreme variations in viewpoint and scale, dramatic changes in lighting and large amount of clutter and occlusions, but also enables computation of 3D structure and motion cues, which can aid segmentation and recognition of object and non-object categories. As a part of this proposal we have developed techniques for semantic labeling of outdoors and indoors environments using photometric and geometric cues from video. The proposed approach is informed by novel features and representations for learning models of objects and non-object categories from video, works effectively with multiple sensing modalities and can be deployed on static frames as well as video in a recursive setting. We have tested the approach extensively on benchmark sequences of indoors and outdoors environments.

---

**Enter List of papers submitted or published that acknowledge ARO support from the start of the project to the date of this printing. List the papers, including journal references, in the following categories:**

**(a) Papers published in peer-reviewed journals (N/A for none)**

<u>Received</u>	<u>Paper</u>
12/18/2014 3.00	A. C. Murillo, G. Singh, J. Kosecka. Localization in urban environments using a panoramic gist descriptor, IEEE Transactions on Robotics, (09 2012): 140. doi:
<b>TOTAL:</b>	<b>1</b>

**Number of Papers published in peer-reviewed journals:**

---

**(b) Papers published in non-peer-reviewed journals (N/A for none)**

<u>Received</u>	<u>Paper</u>
<b>TOTAL:</b>	

**Number of Papers published in non peer-reviewed journals:**

---

**(c) Presentations**

Department of Computer Science, UC Berkeley, Seminar. Semantic Segmentation for Robot Perception, August 8, 2014

RSS Workshop on Multiview Geometry, UC Berkeley, 3D reconstruction and Semantic Parsing for Robotic Systems, July 12, 2014, Invited Speaker

Workshop on Long Term Autonomy, IEEE ICRA 2015, Hong Kong, Semantic Segmentation in the Wild, July 1, 2014, Invited Speaker

Semantic Segmentation for Robot Perception, February 2014, CVPR AC Workshop, University of Maryland

Robot Perception in the Cloud, NSF - ARL Cloud Robotics workshop, Philadelphia, 2013

Semantic Parsing of Street Scenes, Invited talk, May 2012, SUNY Buffalo

**Number of Presentations:** 5.00

---

**Non Peer-Reviewed Conference Proceeding publications (other than abstracts):**

Received

Paper

**TOTAL:**

**Number of Non Peer-Reviewed Conference Proceeding publications (other than abstracts):**

---

**Peer-Reviewed Conference Proceeding publications (other than abstracts):**

<u>Received</u>	<u>Paper</u>
10/01/2012	2.00 Michael Paton, Jana Kosecka. Adaptive RGB-D Localization, 2012 Canadian Conference on Computer and Robot Vision (CRV). 28-MAY-12, Toronto, Ontario, Canada. : ,
10/03/2013	6.00 Cesar Cadena Lerma, Jana Kosecka. Semantic Parsing for Priming Object Detection in RGB-D Scenes , Workshop on Semantic Perception, Mapping and Exploration (in conjunction with IEEE ICRA), 2013. 06-MAY-13, . : ,
10/03/2013	4.00 Jana Kosecka. Detecting Changes in Images of Street Scenes, ACCV, Asian Conference on Computer Vision, 2012. 05-NOV-12, . : ,
10/03/2013	5.00 Gautam Singh, Jana Kosecka. Nonparametric Scene Parsing with Adaptive Feature Relevance and Semantic Context, IEEE Conference on Computer Vision and Pattern Recognition. 21-JUN-13, . : ,
12/11/2014	9.00 Cesar Caden , Jana Kosecka. Recursive Inference for Prediction of Objects in Urban Environments, International Symposium on Robotics Research. 13-DEC-13, . : ,
12/11/2014	10.00 Cesar Cadena, Jana Kosecka. Semantic segmentation with heterogeneous sensor coverages, 2014 IEEE International Conference on Robotics and Automation (ICRA). 31-MAY-14, Hong Kong, China. : ,
12/18/2014	14.00 Md. Alimoor Reza, Jana Kosecka. Object Recognition and Segmentation in Indoor Scenes from RGB-D images, RGB-D Advanced Reasoning with Depth Cameras, Robotics Science and Systems Workshop, RSS . 14-JUL-14, . : ,
<b>TOTAL:</b>	<b>7</b>

**Number of Peer-Reviewed Conference Proceeding publications (other than abstracts):**

---

**(d) Manuscripts**

<u>Received</u>	<u>Paper</u>
12/11/2014	12.00 Cesar Cadena, Jana Kosecka. Semantic Parsing for Priming Object Detection inIndoors RGB-D Scenes <sup>L</sup> , International Journal of Robotics Research (08 2014)
<b>TOTAL:</b>	<b>1</b>

**Number of Manuscripts:**

---

**Books**

Received      Book

**TOTAL:**

Received      Book Chapter

**TOTAL:**

**Patents Submitted**

---

**Patents Awarded**

---

**Awards**

---

**Graduate Students**

<u>NAME</u>	<u>PERCENT SUPPORTED</u>	<u>Discipline</u>
Md. Alimoor Reza	1.00	
George Georgiakis	1.00	
Gautam Singh	0.17	
Xing Zhou	1.00	
Michael Paton	0.08	
<b>FTE Equivalent:</b>	<b>3.25</b>	
<b>Total Number:</b>	<b>5</b>	

### Names of Post Doctorates

<u>NAME</u>	<u>PERCENT SUPPORTED</u>
Cesar Cadena	1.00
<b>FTE Equivalent:</b>	<b>1.00</b>
<b>Total Number:</b>	<b>1</b>

### Names of Faculty Supported

<u>NAME</u>	<u>PERCENT SUPPORTED</u>	National Academy Member
Jana Kosecka	0.11	
<b>FTE Equivalent:</b>	<b>0.11</b>	
<b>Total Number:</b>	<b>1</b>	

### Names of Under Graduate students supported

<u>NAME</u>	<u>PERCENT SUPPORTED</u>	Discipline
Don Ma	0.08	Computer and Computational Sciences
Liam Dan	0.08	Computer and Computational Sciences
<b>FTE Equivalent:</b>	<b>0.16</b>	
<b>Total Number:</b>	<b>2</b>	

### Student Metrics

This section only applies to graduating undergraduates supported by this agreement in this reporting period

The number of undergraduates funded by this agreement who graduated during this period: ..... 2.00

The number of undergraduates funded by this agreement who graduated during this period with a degree in science, mathematics, engineering, or technology fields:..... 2.00

The number of undergraduates funded by your agreement who graduated during this period and will continue to pursue a graduate or Ph.D. degree in science, mathematics, engineering, or technology fields:..... 1.00

Number of graduating undergraduates who achieved a 3.5 GPA to 4.0 (4.0 max scale):..... 2.00

Number of graduating undergraduates funded by a DoD funded Center of Excellence grant for Education, Research and Engineering:..... 0.00

The number of undergraduates funded by your agreement who graduated during this period and intend to work for the Department of Defense ..... 0.00

The number of undergraduates funded by your agreement who graduated during this period and will receive scholarships or fellowships for further studies in science, mathematics, engineering or technology fields:..... 0.00

### Names of Personnel receiving masters degrees

<u>NAME</u>	
Michael Paton	
<b>Total Number:</b>	<b>1</b>

### Names of personnel receiving PHDs

<u>NAME</u>	
<b>Total Number:</b>	

---

**Names of other research staff**

NAME

PERCENT SUPPORTED

**FTE Equivalent:**

**Total Number:**

---

**Sub Contractors (DD882)**

**Inventions (DD882)**

**Scientific Progress**

See the attachement.

**Technology Transfer**

# BAA 1235672148: Final Report

**Project Title: Acquiring semantically meaningful models for robotic localization, mapping and target recognition:**

Jana Košecká, Department of Computer Science George Mason University  
4400 University Drive, Fairfax, VA 22032  
e-mail: kosecka@cs.gmu.edu, phone: 703-993-1876

## 1. Foreword

The main goal of this proposal is to develop novel representations of objects and environments for localization, semantic mapping and target detection. Majority of the research efforts in mapping and visual perception for robotic systems, focused on the problems of localization, map building of doing both jointly know as simultaneous mapping and localization (SLAM) problem. The maps proposed in the past ranged from metric, topological of hybrid representations of the environments. While these models are suitable for navigation tasks, endowing such models with additional semantic information can enable more complex tasks, such as object search or better target/object detection as well as more advanced interactions with humans. Semantic labeling techniques strive to assign different semantic labels to different partitions of the data and use the context of indoors and outdoors environments improve the state of the art of existing visual localization strategies, contextual target detection and recognition and semantic mapping. The focus of our approach is on the development of unified representations, which can be adopted to the task at hand. These representations and framework for learning and inference will be an integral part of perceptual capabilities of a robotic system and will be evaluated using different sensing modalities and different tasks in indoors and outdoors environments.

## 2. Problem Statement

The semantic mapping of the environment requires simultaneous segmentation and categorization of the acquired stream of sensory information. The existing methods typically consider the semantic mapping as the final goal and differ in the number and types of considered semantic categories. We envision semantic understanding of the environment as an on-going process and seek representations which can be refined and adapted depending on the task and robot's interaction with the environment. In this work we propose a novel and efficient method for semantic parsing, which can be adopted to the task at hand and enables localization of objects of interest in indoors environments. For basic mobility tasks we demonstrate how to obtain initial semantic segmentation of the scene into *ground*, *structure*, *furniture* and *props* categories which constitute the first level of hierarchy. Then, we propose a simple and efficient method for predicting locations of objects that based on their size afford a manipulation task. In our experiments we use the publicly available NYU V2 dataset [8] and obtain better or comparable results than the state of the art at the fraction of computational cost. We show the generalization of our approach on two more publicly available datasets.

## 3. Summary of the results

Over the duration of the project we have made several significant contributions in semantic understanding of multimodal sensory data in indoors and outdoors environments. We will briefly summarize them below, while the

additional details can be found in the accompanying publications.

**Priming Object Detection** In [2, 3] we have demonstrated very efficient algorithm which for initial semantic segmentation of the scene into ground, structure, furniture and props categories which constitute the first level of hierarchy. The main technical insights was the use of minimum weight spanning tree approximation of the inference graph, which was computed on 3D depth data and effective and efficient to compute features. These choices enabled us to use well conditioned exact inference techniques for the learning and estimation of the final labeling and yielding improved performance at the fraction of the computational cost on the standard benchmark RGB-D dataset on NYU V2. The initial version of this work was published in a workshop, followed by submission and acceptance of the work to International Journal of Robotics Research [3].

**Recursive Semantic Labeling** In the follow up work we have extended the static semantic parsing to a video setting and proposed a recursive Bayes filter style updating mechanism [1]. In this problem we focused on outdoors environments and exploited widely available exemplars of non-object categories (such as road, buildings, vegetation) and used geometric cues which are indicative of the presence of object boundaries to gather the evidence about objects regardless of their category. We have carried out extensive experiments on videos of urban environments acquired by a moving vehicle and show quantitatively and qualitatively the benefits of our proposal. Another notable feature of the resulting approach was close to real-time performance of the whole system (5 fps), including the feature computation and inference.

**Heterogeneous Coverage** In the previous approaches were were able to compute the semantic labeling for regions of the images and video using only one sensing modality, incorrectly interpolate measurements of other modalities or at best assign semantic labels only to the spatial intersection of coverages of different sensors. In this work we proposed a method for inferring semantic labels Using the previously proposed strategy for inducing the graph structure of Conditional Random Field used for inference, in this work we proposed a novel method for computing the sensor domain dependent potentials. This strategy enabled us to achieve superior semantic segmentation for the regions in the union of spatial coverage of the sensors, while keeping the computational cost of the approach low. The problem is illustrated in Fig. 1. For example with an image sensor note how in column (b) one portion of the car is confused with the ground because their colors are similar. We demonstrated how to combine the visual sensing with the evidence from a 3D laser sensor and mitigate sensor specific perceptual confusers, column (c), but now we are only able to explain a subset of the scene, the spatial intersection coverage, leaving us without output for the car glass and the building in the top portion of the image. With the strategy we introduced, we can take the advantage of both sensor modalities without discarding the non-overlapping zones, column (d) in Fig. 1.

**Finer Grained Semantic Labeling** The previous techniques we have discussed very efficient semantic labeling techniques for small number of semantic categories. This was possible due to efficient features and inference algorithms. In order to obtain better discrimination capabilities for different categories additional features and alternative inference algorithms have to be computed. We proposed to formulate the multi-class object recognition and segmentation in RGB-D data using many binary object-background segmentation, using informative set of features and grouping cues for the small regular superpixels. The main novelty of the proposed approach is the exploitation of the informative depth channel features which indicate presence of depth boundaries, the use of efficient supervised object specific binary segmentation and effective hard negative mining exploiting the object co-occurrence statistics. The binary segmentation is meaningful in the context of robotics applications, where often only the object of interest need to be sought. This yields an efficient and flexible method, which can be easily extended to additional object categories. We report the performance of the approach on NYU-V2 indoors

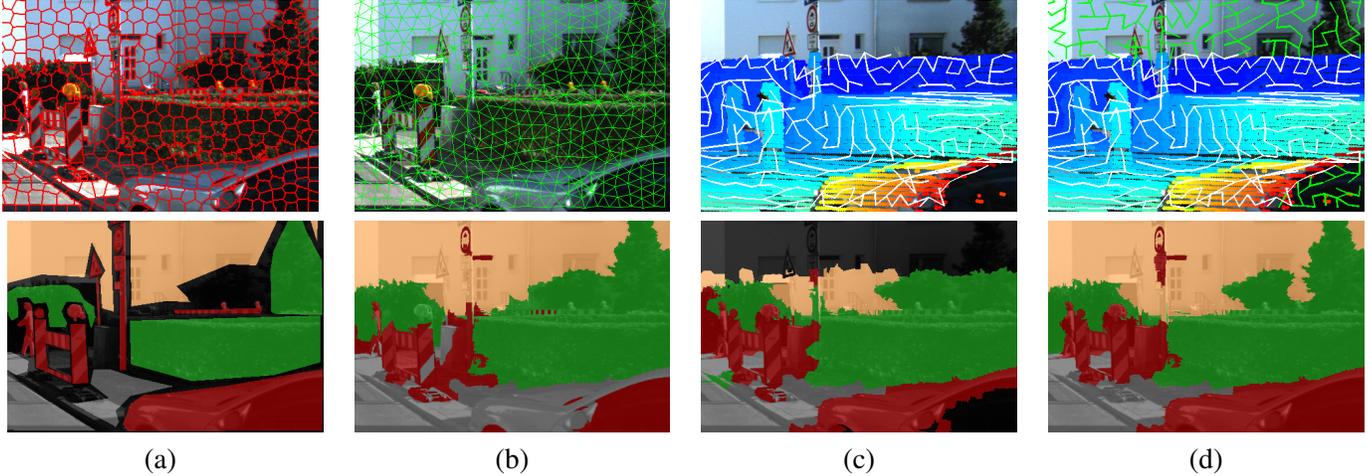


Figure 1. We propose in this work a new approach to semantic parsing, which can seamlessly integrate evidence from multiple sensors with overlapping but possibly different fields of view and account for missing data, while predicting semantic labels over the spatial union of sensors coverages. The semantic segmentation is formulated on a graph, in a manner which depends on sensing modality. First row: (a) over-segmentation over the image; (b) graph induced by the superpixels; (c) the 3D point cloud re-projected on the image with a tree graph structure computed in 3D, and (d) the full graph as proposed here for full scene understanding. In the second row is the semantic segmentation (a) ground truth and results of (b) using the image graph and only visual information; (c) using the 3D graph and visual and 3D information, and finally (d) the result from using a graph for full coverage and all the information. Note the best semantic segmentation achieved over the union of the spatial coverage of both sensors. Color code: ■ ground, ■ objects, ■ building and ■ vegetation.

	Bed	Sofa	Chair	Table	Window	Bookshelf	TV	Bag	Bathtub	Blinds	Books	Box	Cabinet	Clothes	Counter	Curtain	Desks
Silberman[8]	40.00	25.00	32.00	21.00	30.00	23.00	5.70	0.00	0.00	40.00	5.50	0.13	33.00	6.50	33.00	27.00	4.60
Ren[6]	42.00	28.00	33.00	17.00	28.00	17.00	19.00	1.20	7.80	27.00	<b>15.00</b>	3.30	37.00	9.50	39.00	28.00	10.00
Gupta[4]	55.00	<b>44.00</b>	<b>40.00</b>	<b>30.00</b>	<b>33.00</b>	20.00	9.30	0.65	33.00	<b>44.00</b>	4.40	<b>4.80</b>	<b>48.00</b>	6.90	<b>47.00</b>	<b>34.00</b>	10.00
Ours(unary)	50.64	37.44	25.00	19.19	25.93	23.88	26.40	<b>3.28</b>	32.12	29.77	9.17	2.89	27.42	9.79	34.68	25.59	21.04
Ours(CRF)	<b>56.85</b>	42.29	31.44	20.78	30.16	<b>30.29</b>	<b>34.97</b>	3.00	32.95	33.09	10.06	3.99	29.34	10.04	33.82	30.11	<b>23.35</b>

	Door	Dresser	Floor-mat	Lamp	Mirror	Night-stand	Paper	Person	Picture	Pillow	Refrigerator	Shelves	Shower-curtain	Sink	Toilet	Towel	Whiteboard
Silberman[8]	5.90	13.00	7.20	<b>16.00</b>	4.40	6.30	<b>13.00</b>	6.60	36.00	19.00	1.40	3.30	3.60	25.00	27.00	0.11	0.00
Ren[6]	13.00	7.00	20.00	14.00	18.00	9.20	12.00	14.00	32.00	20.00	1.90	6.10	5.40	<b>29.00</b>	35.00	13.00	0.15
Gupta[4]	8.30	22.00	22.00	6.80	19.00	20.00	1.90	16.00	<b>40.00</b>	<b>28.00</b>	15.00	5.10	18.00	26.00	<b>50.00</b>	<b>14.00</b>	37.00
Ours(unary)	14.72	32.35	32.81	6.68	23.09	16.22	7.64	19.54	17.93	16.16	16.86	10.67	25.54	10.98	26.06	7.62	36.25
Ours(CRF)	<b>17.16</b>	<b>35.73</b>	<b>34.19</b>	12.14	<b>27.41</b>	<b>21.54</b>	10.07	<b>30.31</b>	22.21	22.98	<b>20.59</b>	<b>13.46</b>	<b>26.84</b>	11.04	38.65	8.61	<b>37.69</b>

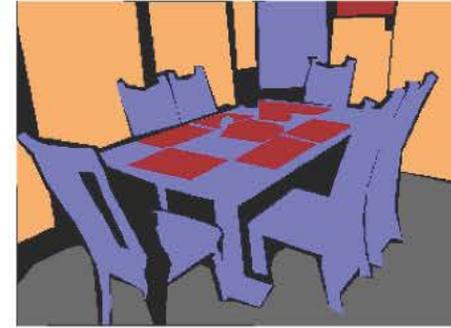
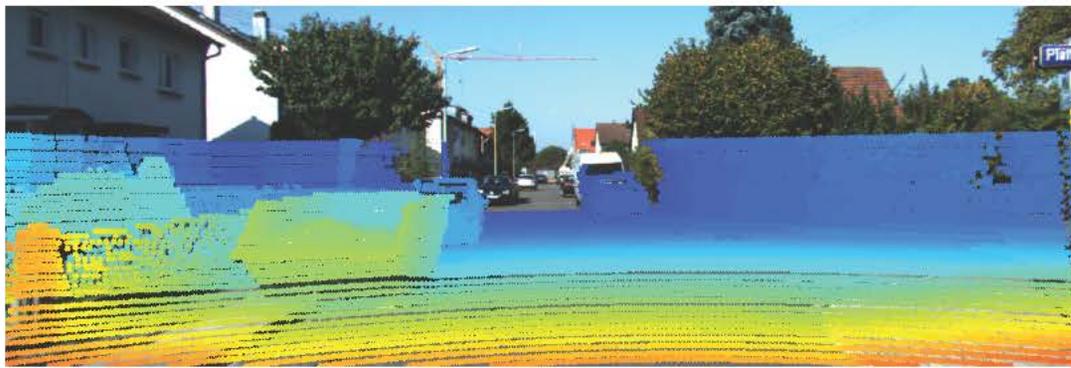
Table 1. Performance on the NYUD-V2 dataset in Jaccard index.

dataset and demonstrate improvement in the global and average accuracy compared to the state of the art methods. The brief summary of the results, highlighting two different performance measures in different categories can be seen in Tables 1. More details of the proposed methods can be found in [7] and the follow up submission, which is currently under review.

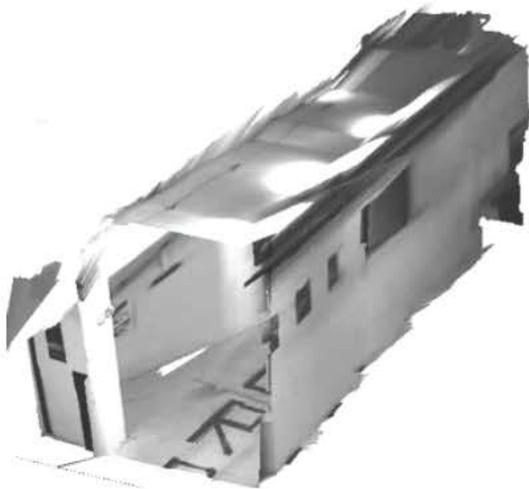
## References

- [1] C. Cadena and J. Košecka. Recursive inference for prediction of objects in urban environments. *International Symposium on Robotics Research, Singapore*, 2013. 2
- [2] C. Cadena and J. Košecka. Semantic parsing for priming object detection in RGB-D scenes. *International Conference on Robotics Automation (ICRA) - 3rd Workshop on Semantic Perception, Mapping and Exploration (SPME)*, 2013. 2

- [3] C. Cadena and J. Košecka. Semantic parsing for priming object detection in RGB-D scenes. *International Journal of Robotics Research*, 2014. 2
- [4] S. Gupta, P. Arbelaez, and J. Malik. Perceptual organization and recognition of indoor scenes from RGB-D images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 3
- [5] M. Paton and J. Košecka. Dynamic RGB-D mapping. In *George Mason University Technical Report TR-GMU-2012-001*, January 2012.
- [6] X. Ren, L. Bo, and D. Fox. RGB-(D) scene labeling: Features and algorithms. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 3
- [7] M. A. Reza and J. Košecka. Object Recognition and Segmentation in Indoor Scenes from RGB-D Images. *RGB-D Advanced Reasoning with Depth Cameras Workshop, RSS*, 2014. 3
- [8] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision (ECCV)*, 2012. 1, 3



# 3D Reconstruction and Semantic Parsing Robotic Systems



Jana Kosecka  
MVGRO, RSS July 2014

# Robotic Perception

- Methods for modeling 3D geometry of the environment and semantic concepts
- Seamless processing of video streams
- Recursive and multi-view setting
- Efficient Inference

# Robotic Tasks

Navigation

Mapping  
Localization

Manipulation

Human  
Interaction

# Representations

Road, floor, free space  
landmarks, locations

Structural Obj.  
(doors, ...)

Specific Obj.  
(traffic signs, cups,  
bins, ...)

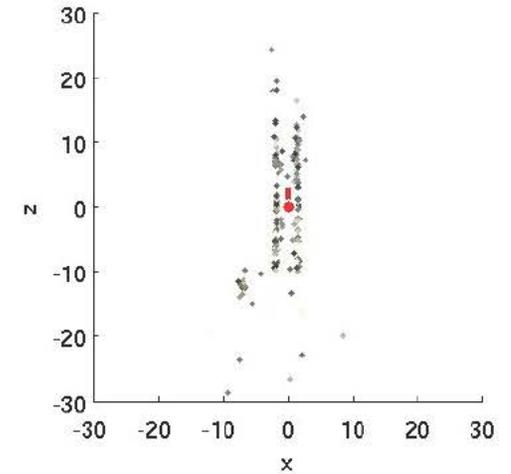
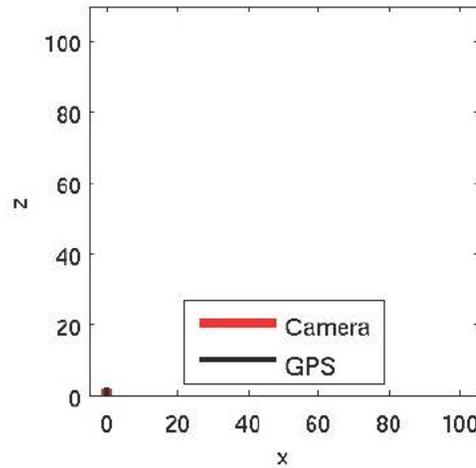
Dynamic Obj, People  
(cars, bikes, ...)

- Representation which can be computed efficiently and are reusable for multiple tasks, extended efficiently to new semantic concepts

# Tasks

# Representations

Mapping  
Localization



- Visual odometry – linear algorithm adopted to 360 FOV
- No need for bundle adjustment
- Guided sampling from entire FOV for RANSAC

# Tasks

# Representations

Mapping  
Localization

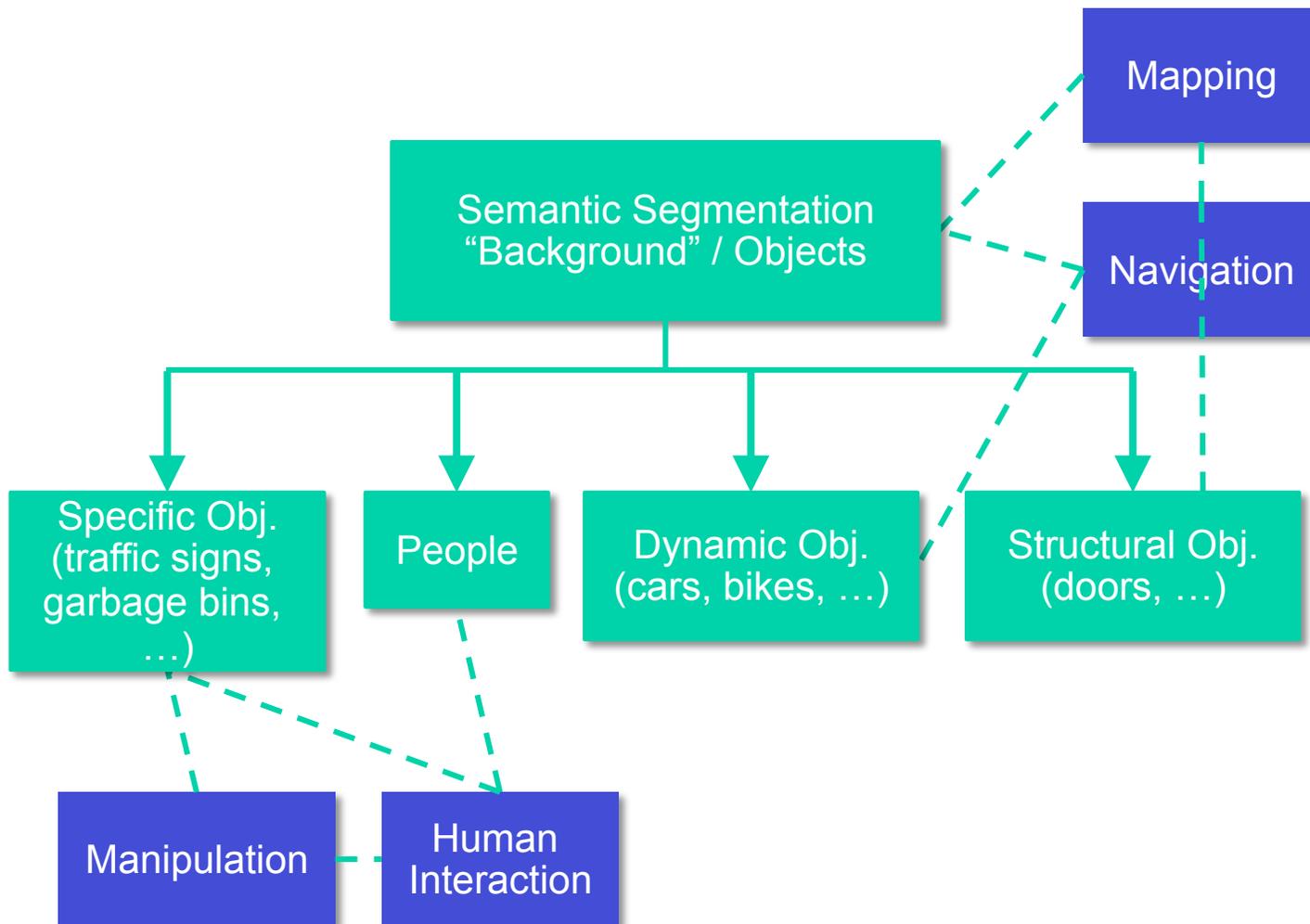
- Point features tracking
- Recovery of relative motion, visual odometry
- Loop closure
- Environment models, sparse clouds of points
- Often sufficient for navigation
- Not sufficient for navigation free space/obstacles
- Location recognition
- Data assoc. with large variations of appearance
- Semantic information for object recognition, human robot interaction
- Obtain high quality denser env. models
- Associate useful semantic inform. with regions/volumes in 3D space

# 3D geometry, Features, Semantics

- 3D geometry: move from sparse feature points to dense models, reason about surfaces etc.
- Overcome challenging, matching and correspondence problems
- Semantic Categorization: learning and inference in a multi-view/recursive setting **Learning and inference**
- Multi-view reconstruction using super-pixels
- Task Dependent Sematic Parsing

# Semantic Hierarchy

Strongly depends on the current task



# images → 3D model



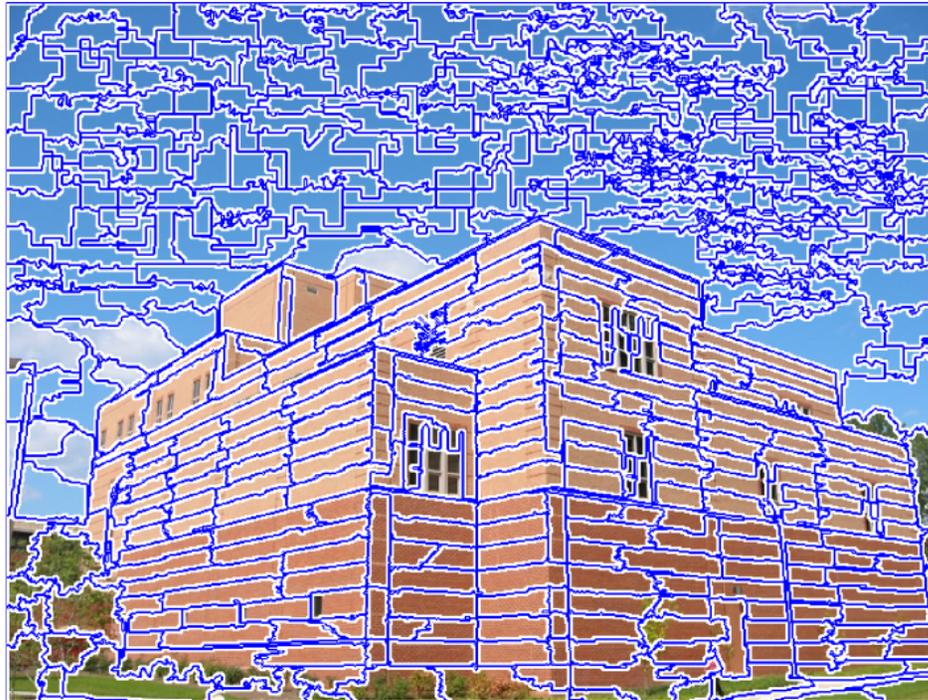
## Challenges of indoors and outdoors environment

- ✗ no or repetitive texture
- ✗ illumination, scale, viewpoint changes, occlusions ...



# Multi-view Superpixel Reconstruction

- pre-segment the reference image into superpixels

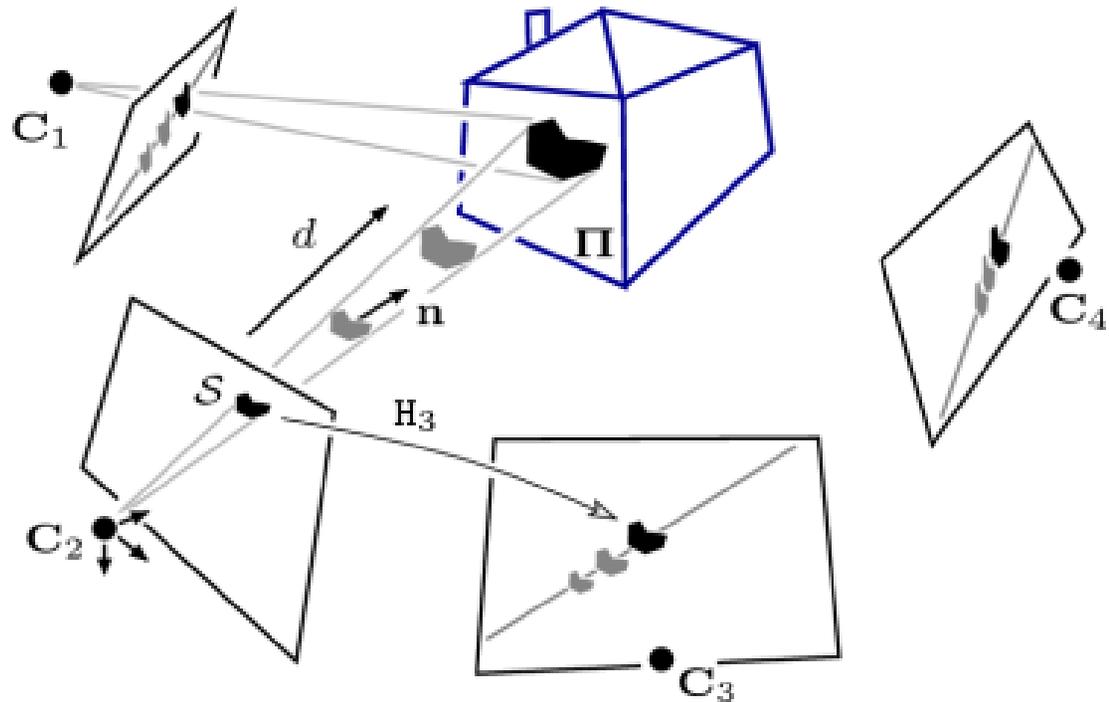


(Felzenszwalb & Huttenlocher, IJCV'04)

- large support area -> more robust measures

**Goal:** to find depth and normal for each superpixel,  $\mathbf{\Pi} = [\mathbf{n}^\top \ d]^\top$

**Assumption:** each super-pixel corresponds to a planar patch in 3D



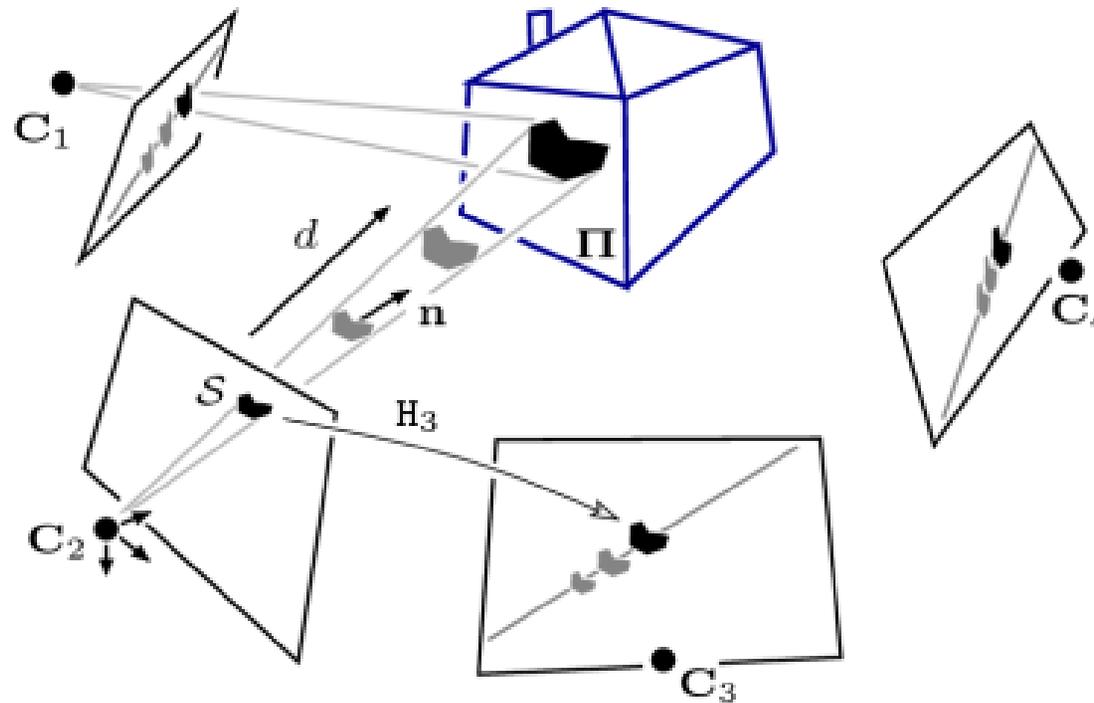
✓ *Known* camera projective matrices  $P_k = K_k [R_k \ t_k]$

- searching for MAP of MRF as a labelling problem

$$\operatorname{argmin}_{\{\mathbf{n}, d\}^{|S|}} \left[ \sum_{\{s\}} E_{photo} + \lambda \sum_{\{s, s'\}} E_{smooth} \right]$$

× **Intractable** over all depths and normals !

# Plane sweeping



✓ Plane induced homography

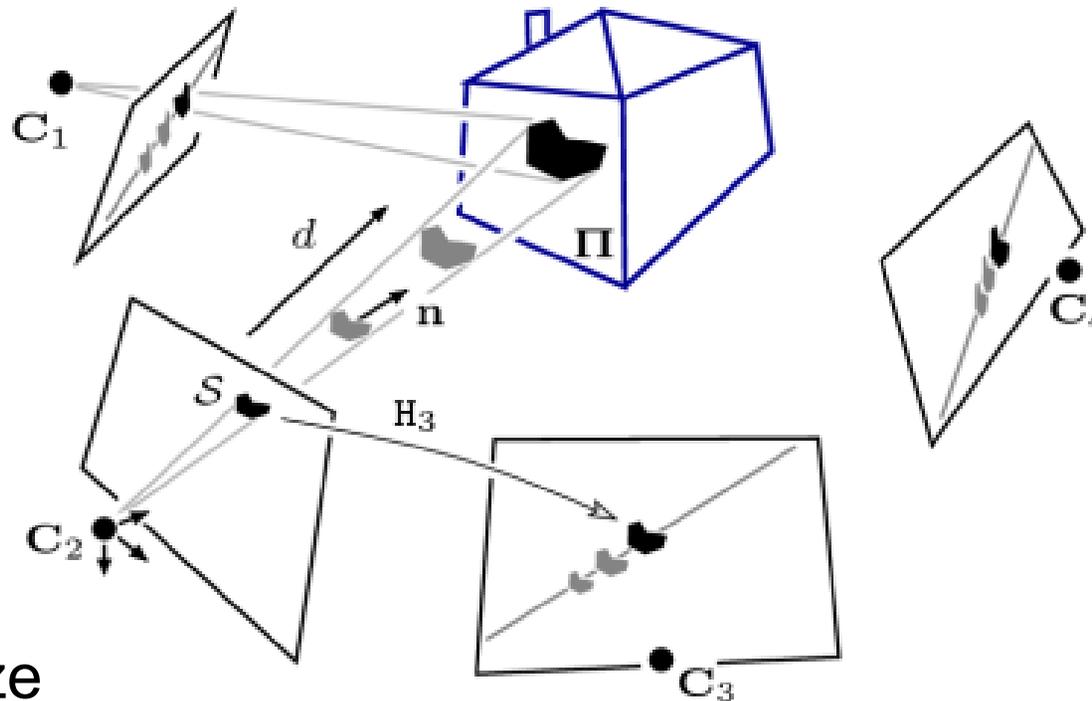
$$H_k(\mathbf{\Pi}, P_k, K_{ref}) = K_k \left( R_k - \mathbf{t}_k \mathbf{n}^\top / d \right) K_{ref}^{-1}$$

**Known** cameras ...  $P_k = K_k [R_k \ \mathbf{t}_k]$

**Unknown** plane ....  $\mathbf{\Pi} = [\mathbf{n}^\top \ d]^\top$

## Restriction of depths

- For each plane normal sweep along depth range and remember only **depth candidates** !

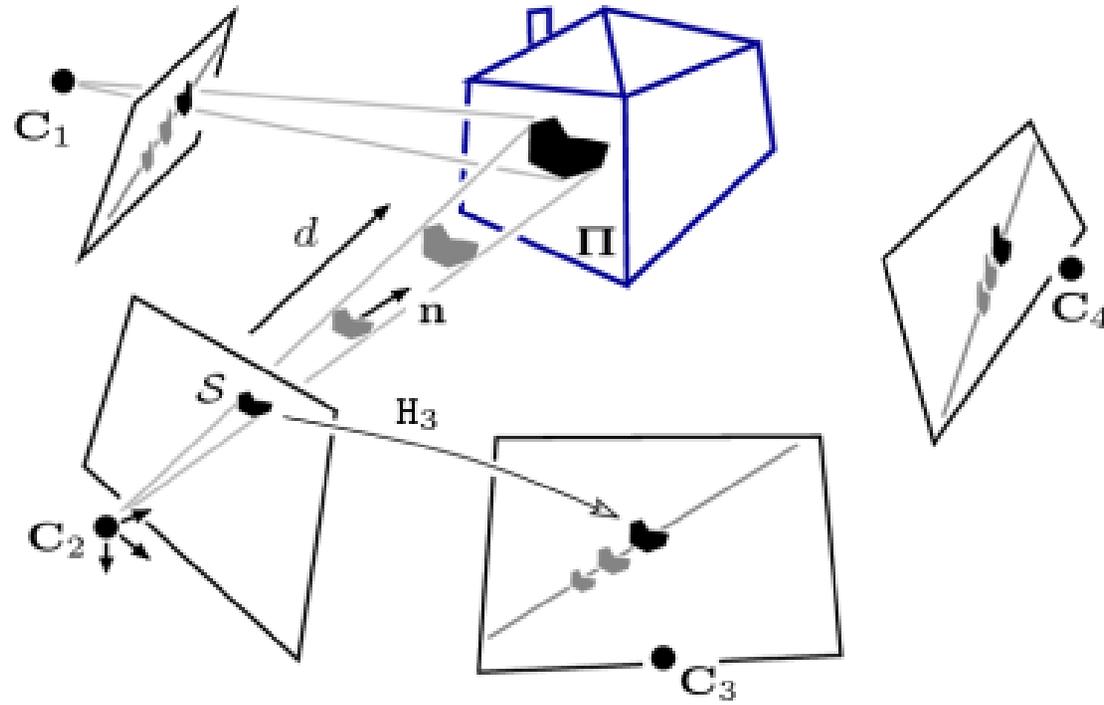


- and optimize

$$\operatorname{argmin}_{\mathcal{P}} \left[ \sum_s E_{photo} + \lambda_1 \sum_s E_{geom} + \lambda_2 \sum_{\{s,s'\}} E_{norm} + \lambda_3 \sum_{\{s,s'\}} E_{depth} \right]$$

over the depth candidates and dominant normals !

# Plane-induced homography



one depth, one normal

# Sweeping



- photometrically normalize all splx projections

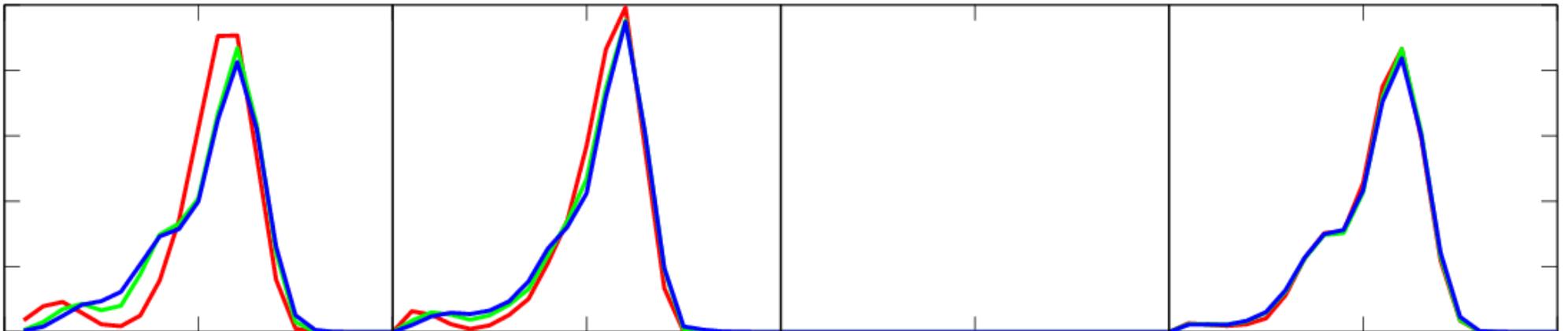


$$\mathbf{c}_k = \frac{[\mu^R \ \mu^G]^\top}{\mu^R + \mu^G + \mu^B}$$

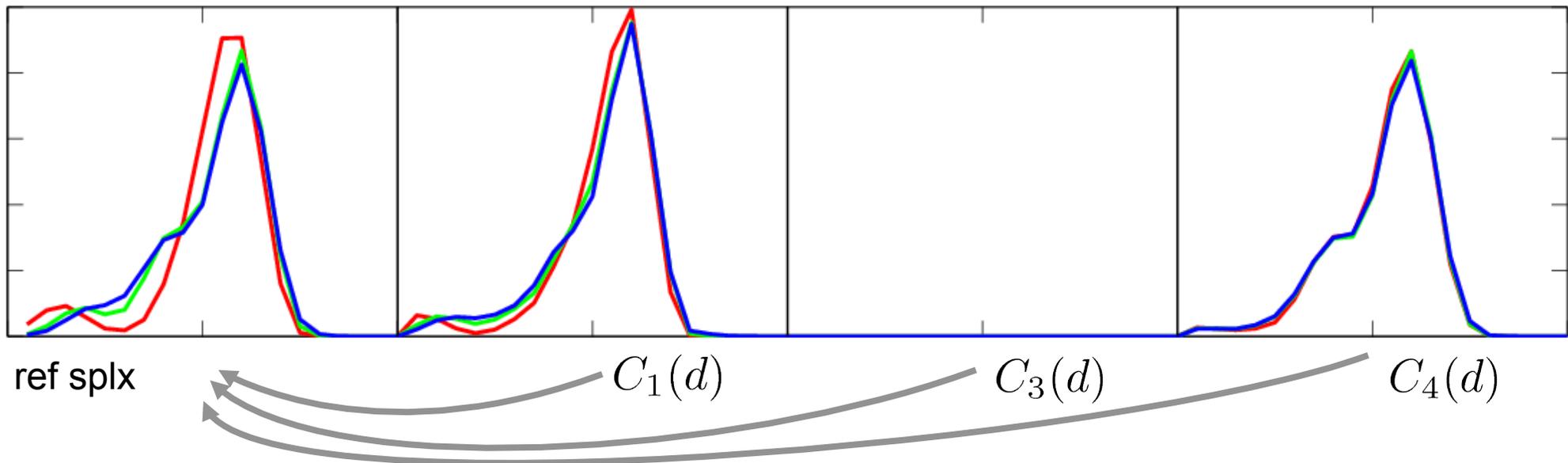
chromacity vector

$$\begin{aligned} \mu^{R,G,B} &= 0 \\ \sigma^{R,G,B} &= 1 \end{aligned}$$

- compute normalized histograms



# Photometric measure

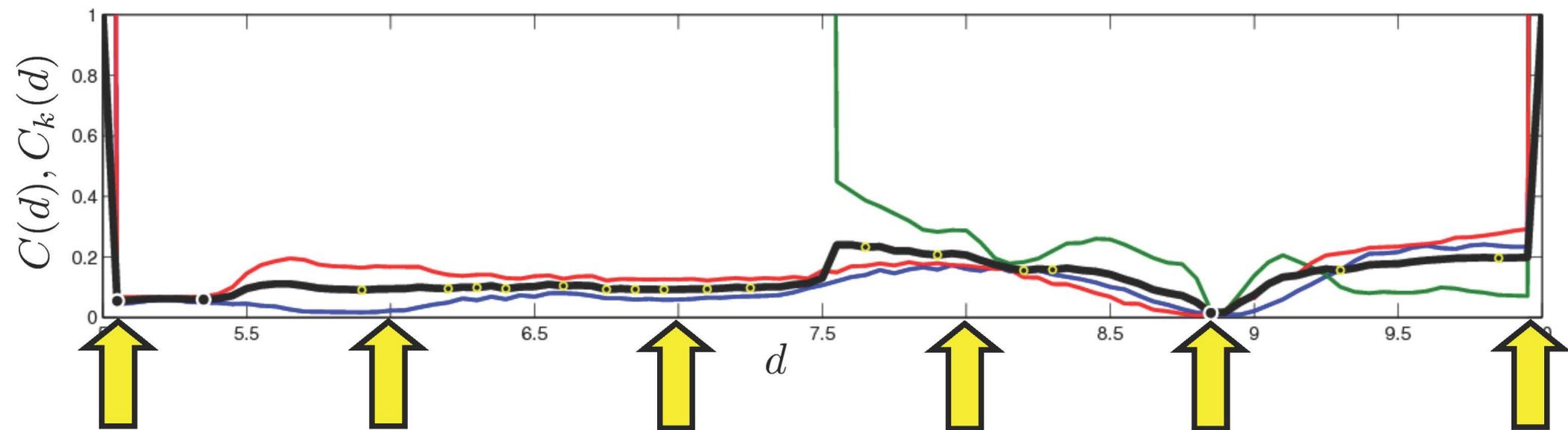
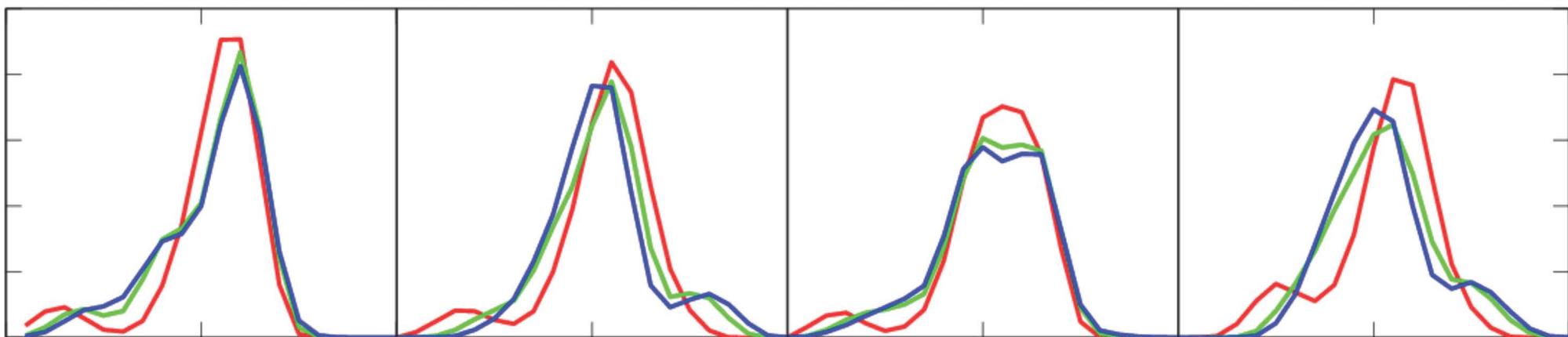


- photometric measure for each splx projection – histogram difference and chromaticity for each reference view/view pair

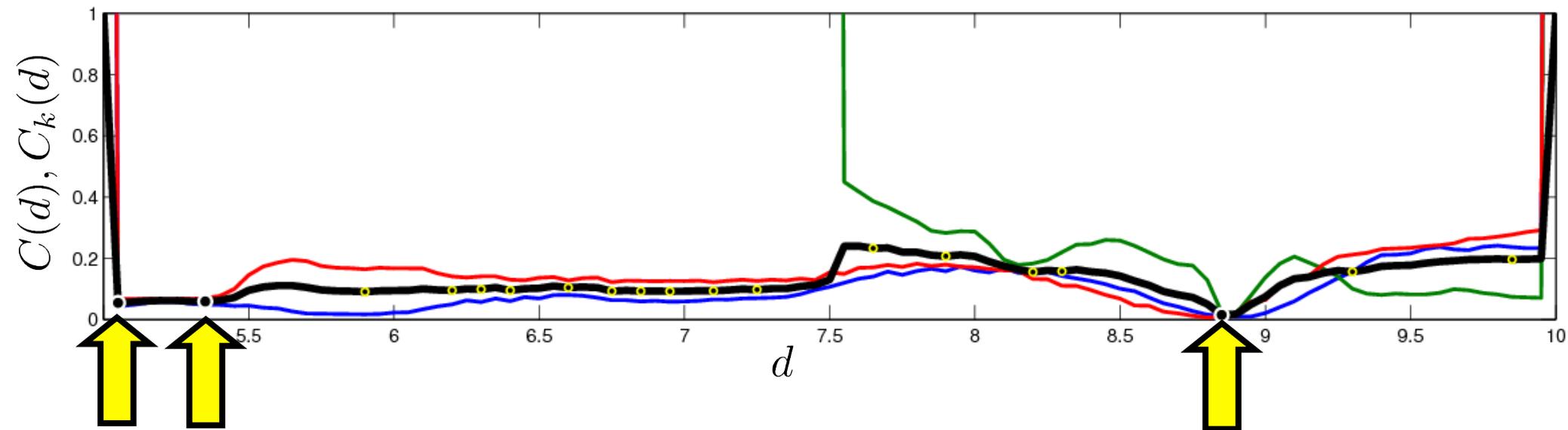
$$C_k(d) = \chi_k^2 + \alpha \|\mathbf{c}_k - \mathbf{c}_{ref}\|^2$$

- composite photometric measure for all views

$$C(d) = \frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} C_k(d)$$



# Depth candidates



- For each normal store  $N$ -best minima as depth candidates in matrix

$$T_{(:,j)}^s = \begin{bmatrix} d \\ i \\ C(d) \end{bmatrix}$$



- Photo-consistency term

$$\operatorname{argmin}_{\mathcal{P}} \left[ \sum_s E_{photo} + \lambda_1 \sum_s E_{geom} + \lambda_2 \sum_{\{s,s'\}} E_{norm} + \lambda_3 \sum_{\{s,s'\}} E_{depth} \right]$$

# Labelling Problem - Energy terms

- Nodes of the graph – superpixels
- Edges of the graph – induced by neighboring superpixels
- Typical pair wise MRF

$$\arg \min_{\mathcal{P}} \sum_s E(s) + \sum_{(s,s')} E(s, s')$$

- In our case:

$$\operatorname{argmin}_{\mathcal{P}} \left[ \sum_s E_{photo} + \lambda \sum_s E_{geom} + \lambda_2 \sum_{\{s,s'\}} E_{norm} + \lambda_3 \sum_{\{s,s'\}} E_{depth} \right]$$

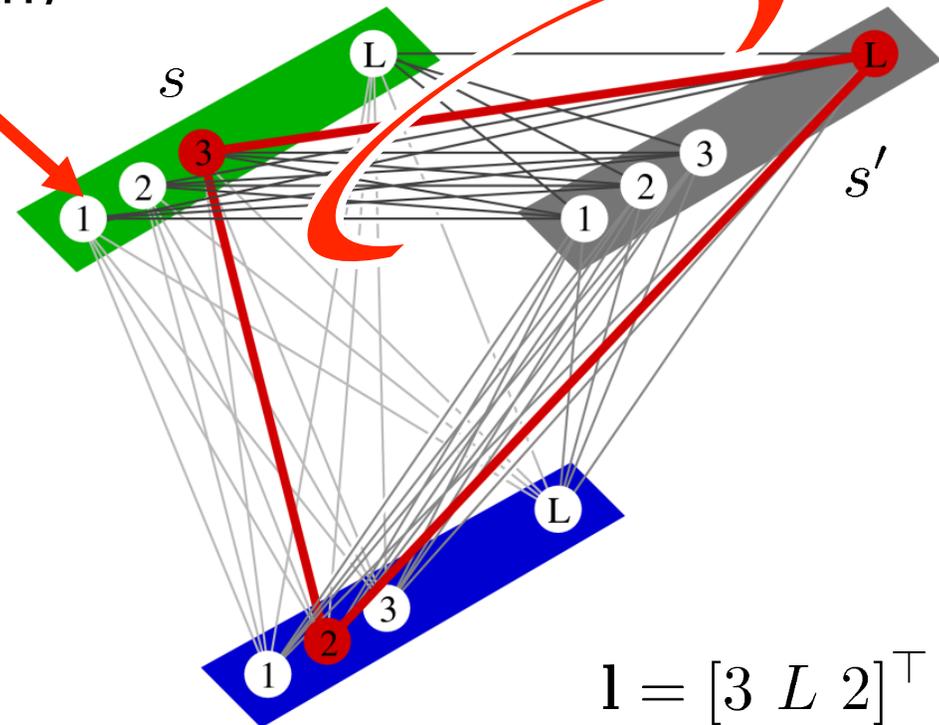
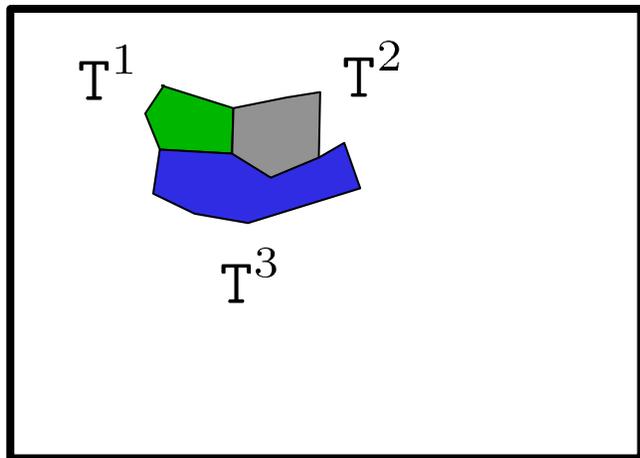
# Labeling

$$\operatorname{argmin}_{\mathcal{P}} \left[ \sum_s E_{photo} + \lambda_1 \sum_s E_{geom} + \lambda_2 \sum_{\{s,s'\}} E_{norm} + \lambda_3 \sum_{\{s,s'\}} E_{depth} \right]$$

$$\mathcal{P} = \text{vector } \mathbf{1}_{S \times 1} \quad E(s, l_s)$$

$$E(s, s', l_s, l_{s'})$$

- **Labeling:** assign a label to each splx.
- plane hypothesis (normal + depth)

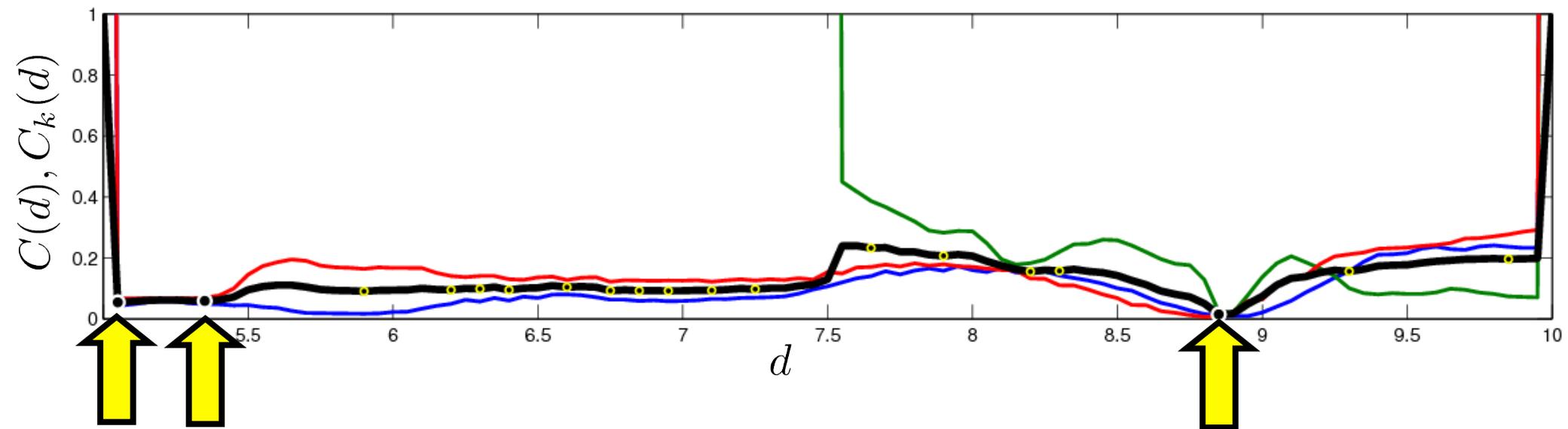


- ✓ available solvers: Kolmogorov PAMI'06, Werner PAMI'07

# Photoconsistency term

$$E_{photo}(s, l_s)$$

Color histogram and chromacity differences

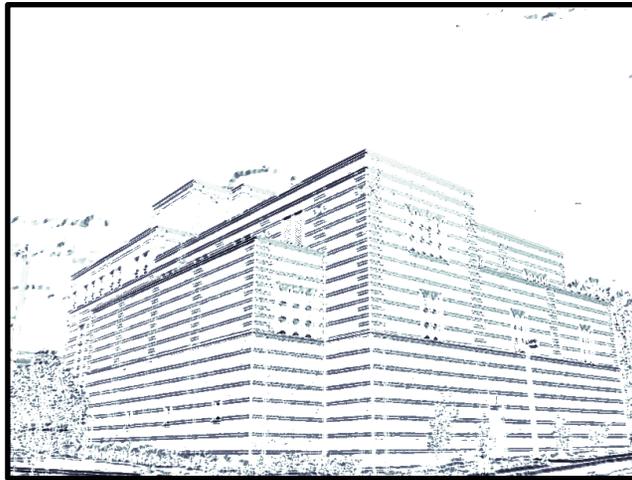


$$\mathbb{T}_{(:,j)}^s = \begin{bmatrix} d \\ i \\ C(d) \end{bmatrix}$$

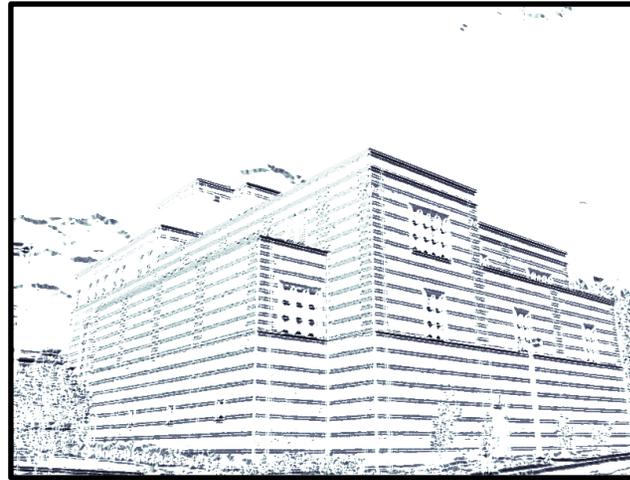
# Geometric term

- ✓ Splx boundaries are usually consistent with dominant directions

gradient mixture model (Coughlan & Yuille NC'03)



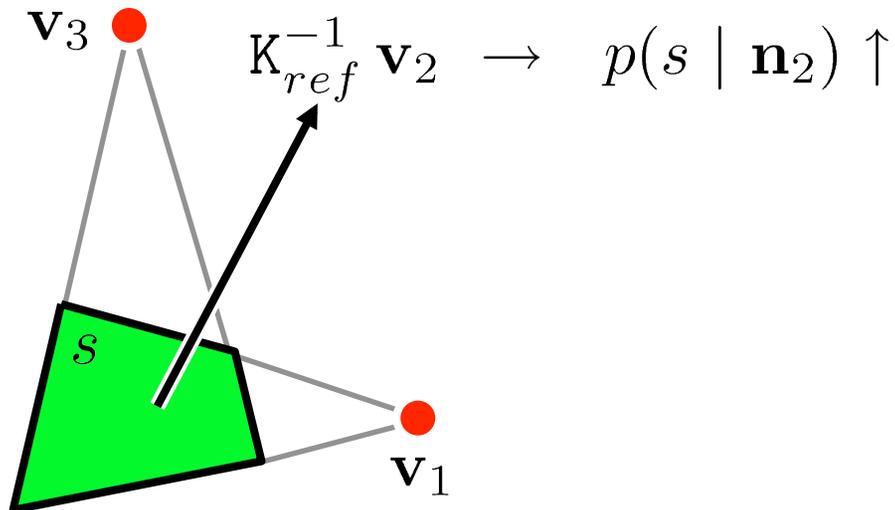
$$p(\mathbf{u} \mid \mathbf{v}_1)$$



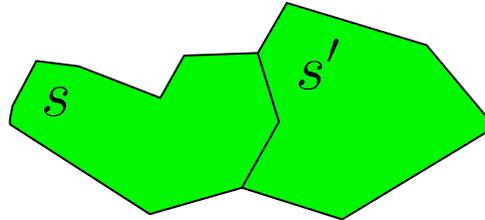
$$p(\mathbf{u} \mid \mathbf{v}_2)$$



$$p(\mathbf{u} \mid \mathbf{v}_3)$$



# Pair wise normal term



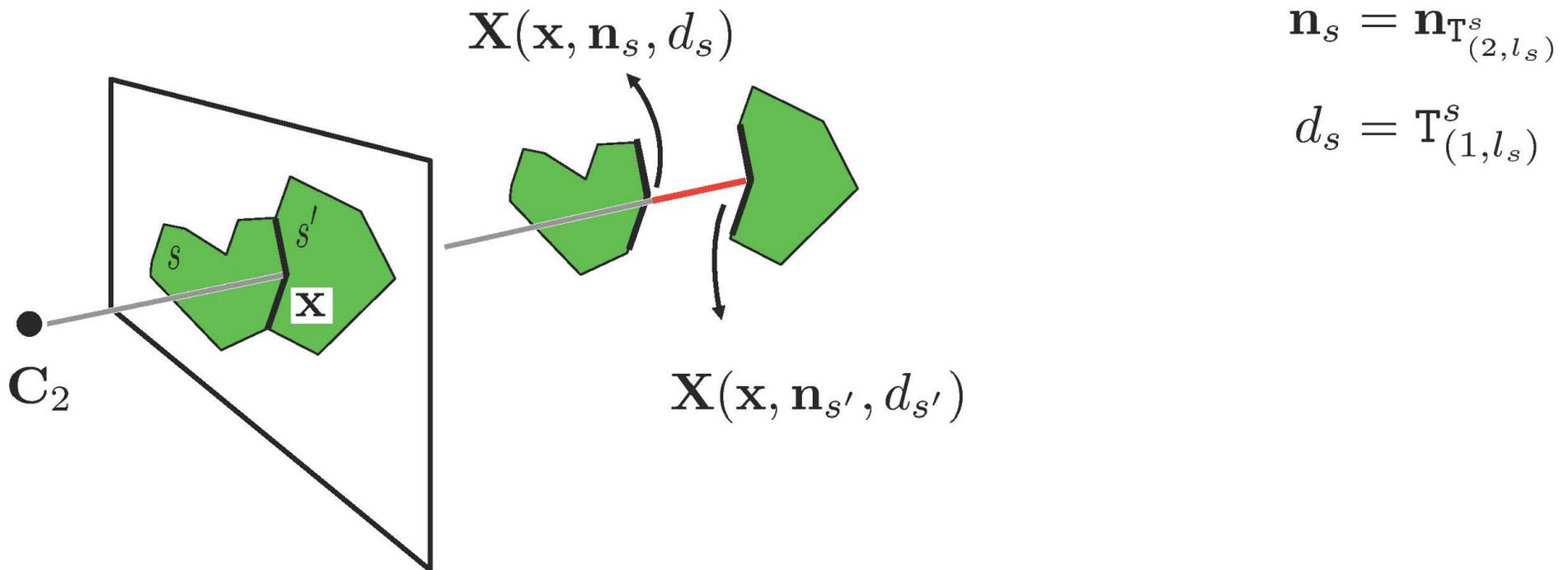
- Force neighboring splxs to have same normal

$$E_{norm}(s, s', l_s, l_{s'}) = \delta(\mathbb{T}_{(2, l_s)}^s \neq \mathbb{T}_{(2, l_{s'})}^{s'})$$

Ising prior

# Pair wise depth term

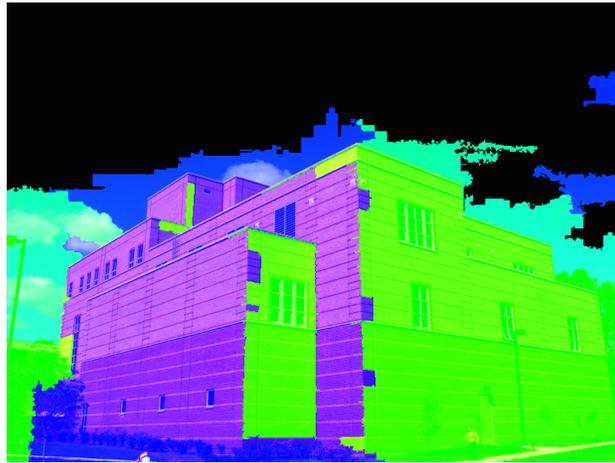
- Force neighboring splxs to touch in 3D



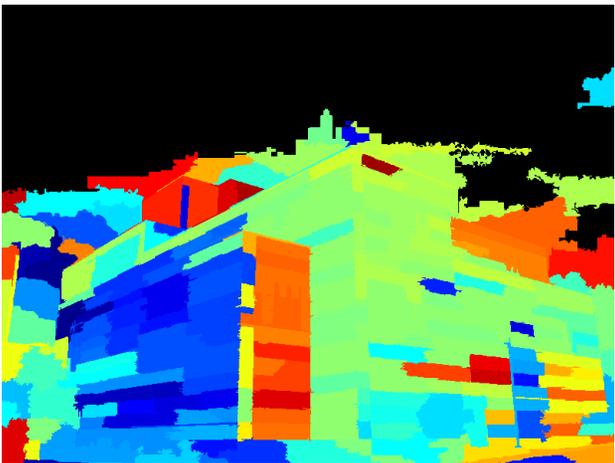
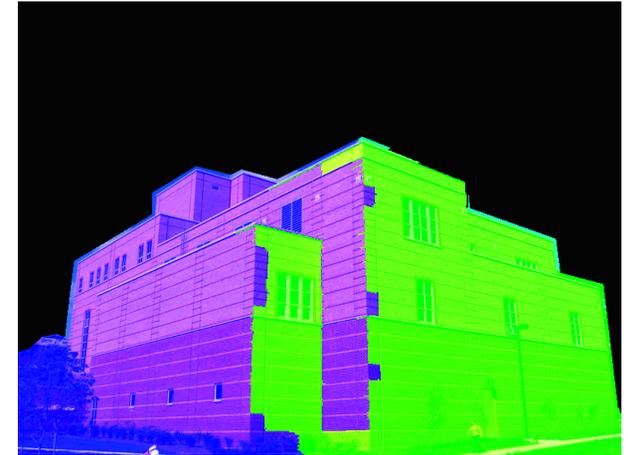
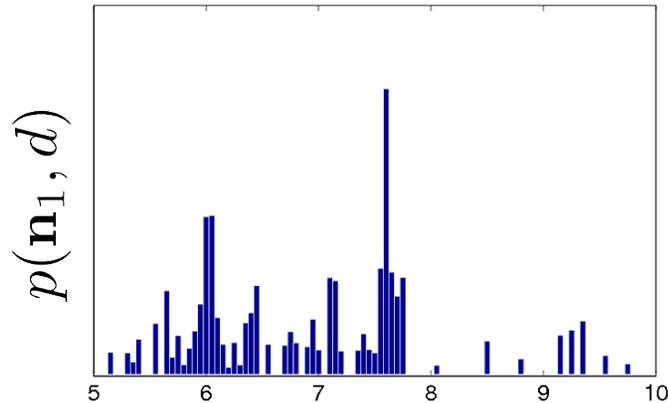
$$E_{depth}(s, s', l_s, l_{s'}) = \min \left( \underset{\mathbf{x} \in \mathcal{S}_s \cap \mathcal{S}_{s'}}{\text{med}} \frac{\|\mathbf{X}(\mathbf{x}, \mathbf{n}_s, d_s) - \mathbf{X}(\mathbf{x}, \mathbf{n}_{s'}, d_{s'})\|}{\|\mathbf{X}(\mathbf{x}, \mathbf{n}_s, d_s)\|}, \Theta_2 \right)$$

# Utilizing plane priors

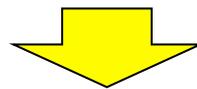
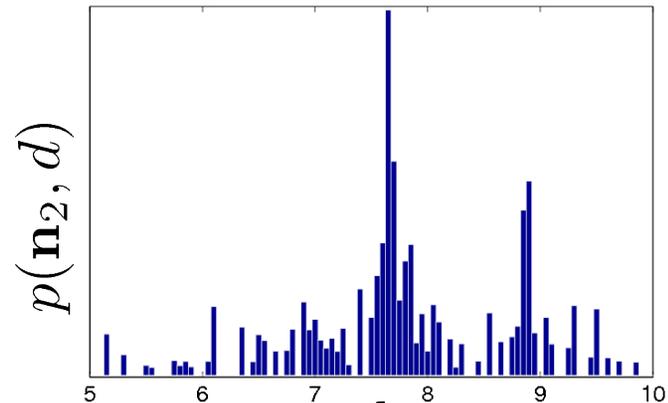
- Compute MAP estimate (run the MRF graph solver once)



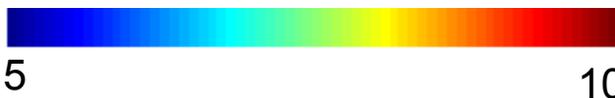
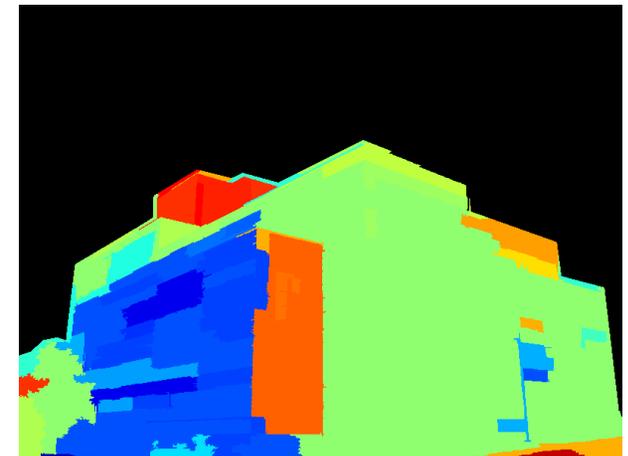
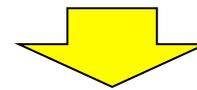
normals



depth map



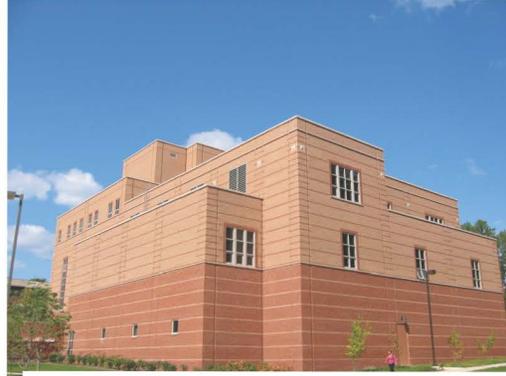
Bayesian formulation



5 10

$$C_2(d) = C(d) - 2\sigma^2 \log p(\mathbf{\Pi}_s)$$

# GMU building



3D model

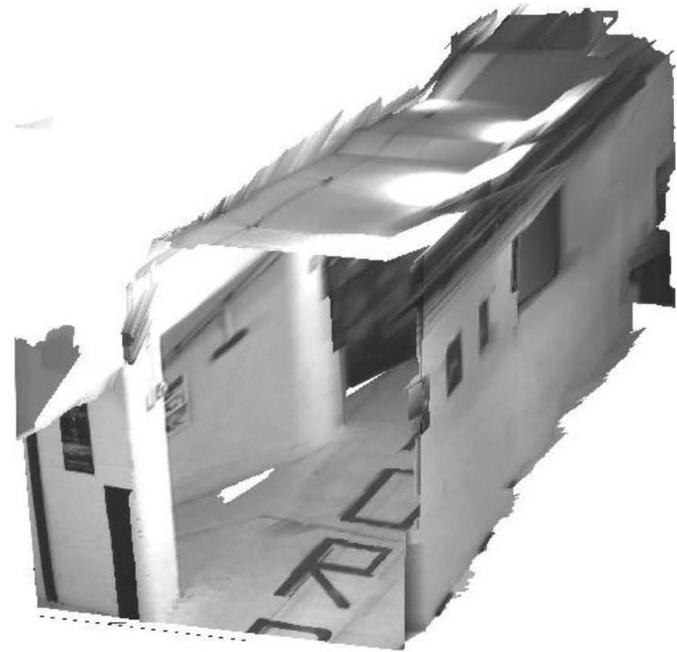
# 3D Reconstruction of street scenes



# Oxford corridor



using 6 images



3D model

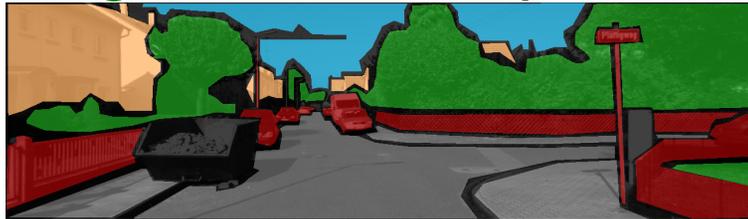
# Multi-view Superpixel Reconstruction

- Superpixel 3D reconstruction of one indoors - outdoors scenes
  - ✓ less computational complexity
  - ✓ alternative photometric/similarity measures
  - ✓ piecewise planar surfaces can be favorably handled
- Extensions to general mostly planar scenes, real-time settings (Gallup'09)
- Integration of 3D reconstruction with recognition

# Associating Semantic Information

Segmentation and categorization of different partitions of sensory data using geometric and appearance cues

- **Navigation** - free space, occupied space



- **Localization/Place recognition** - static portions of environment



- **Object search**, human detection - generate hypotheses about presence of objects

- Object Categorization  
Object Instance recognition

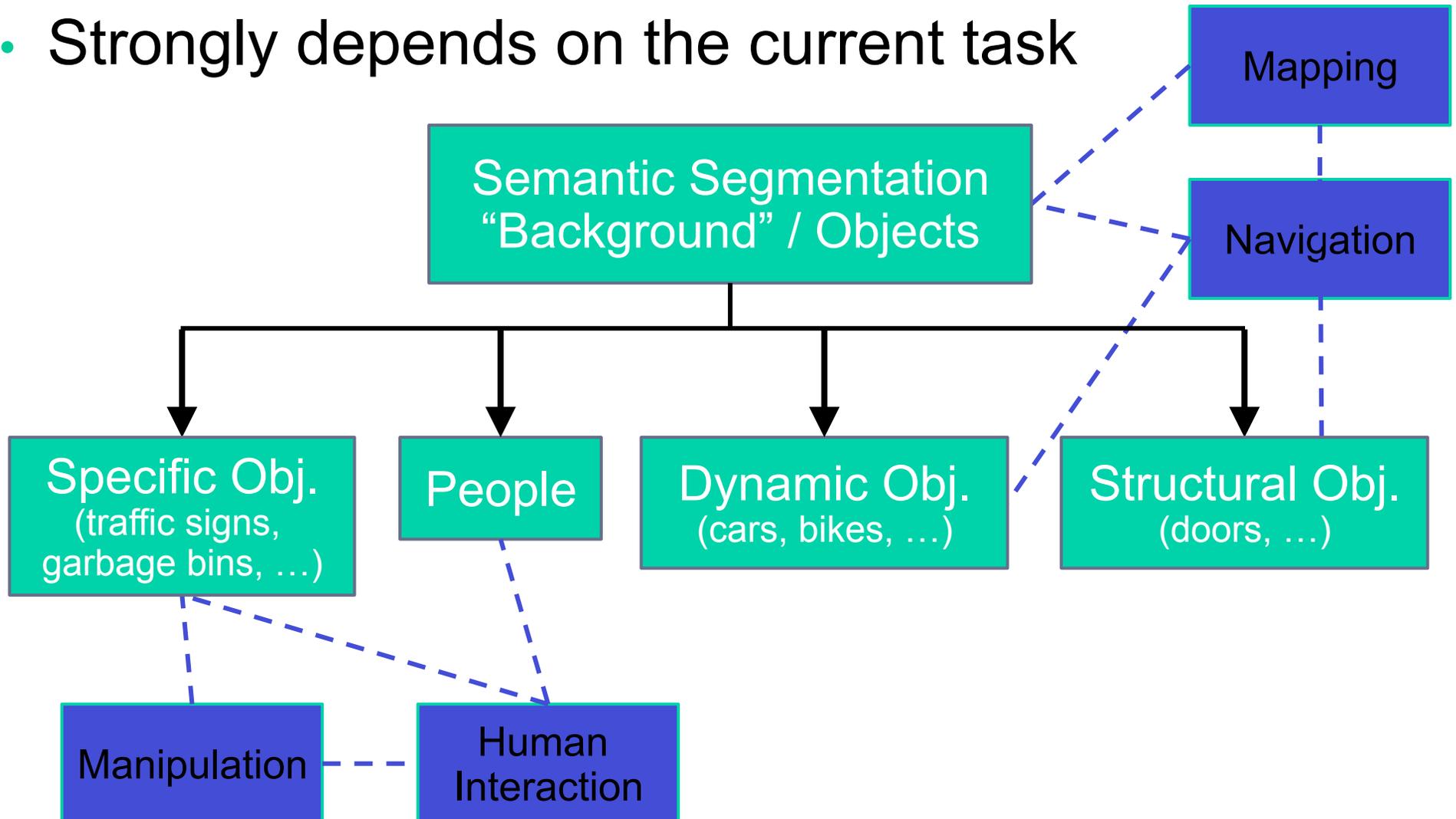


# Semantic Segmentation

1. Task Dependent Semantic Hierarchy
2. Multi-view and Recursive formulation
3. Capability of handling missing 3D data (full sensor's FOVs coverage)

# Proposed Semantic Hierarchy

- Coarse to fine manner
- Strongly depends on the current task



# Proposed Semantic Hierarchy

## “Background” and object categories

- Non-object categories specific to types of scenes
- we can assume to be present (almost) always
- mostly static or slow changing
- Props/Objects

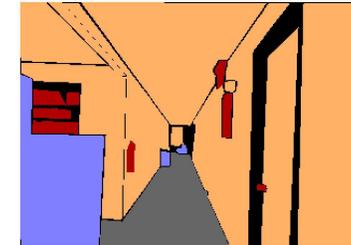
KITTI dataset, Geiger et al. IJRR  
2013



Ground  
Buildings  
Vegetation  
Sky  
Objects

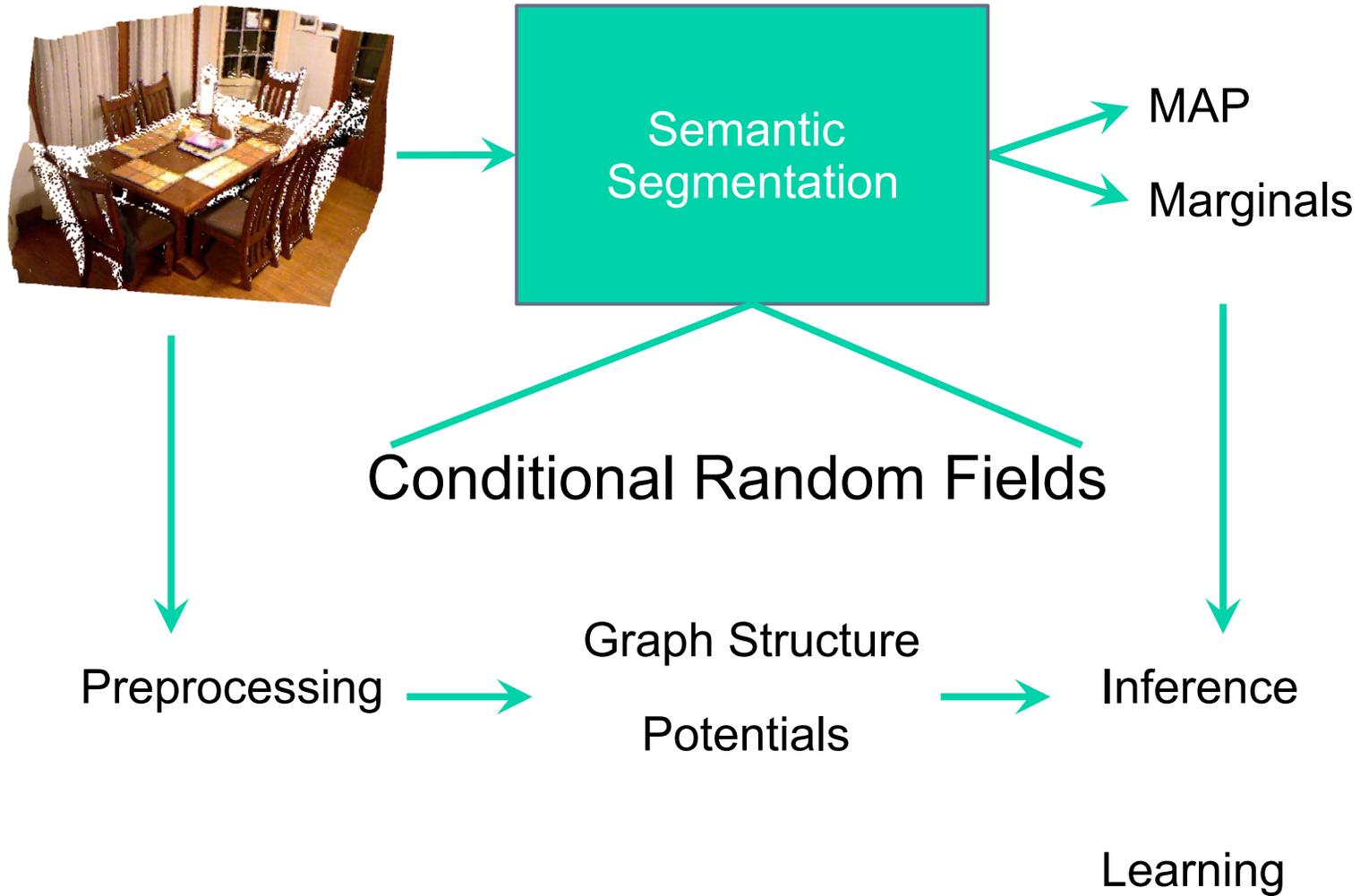
Mapped from Sengupta et al. ICRA  
2013

NYU V2 dataset, Silberman et al. ECCV  
2012



Ground  
Structure  
Furniture  
Props

# Our formulation



# Preprocessing: Over-segmentation

## SLIC superpixels

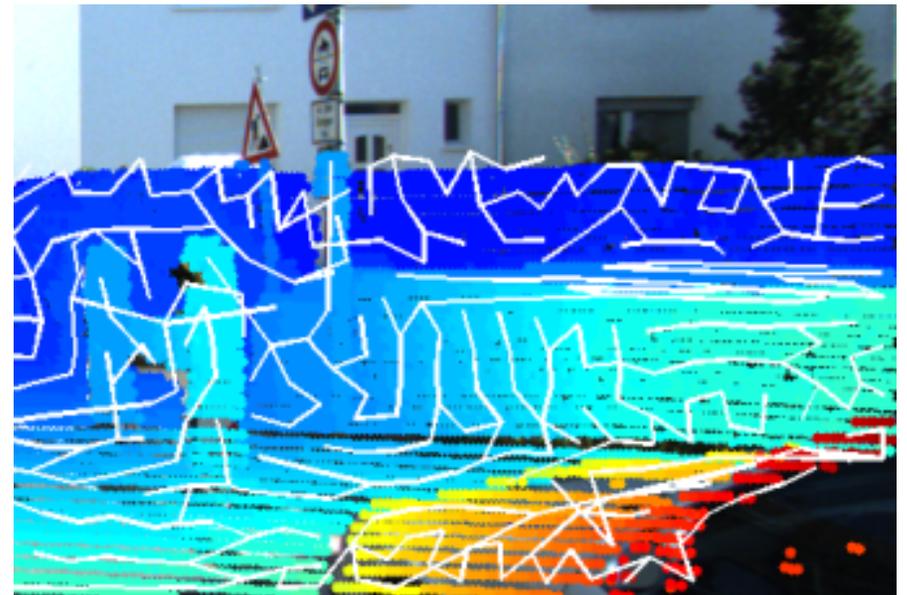
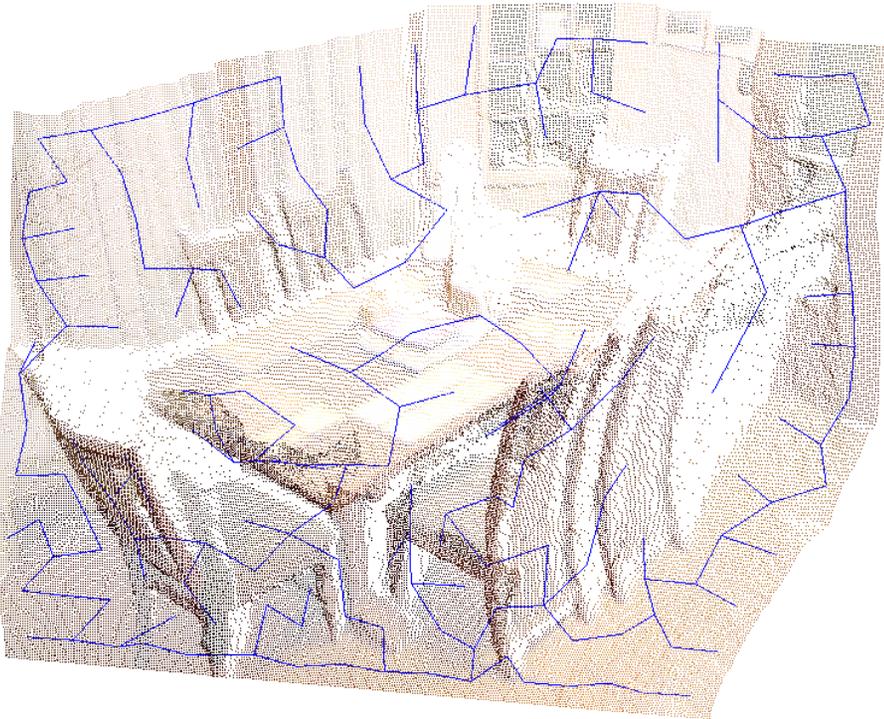


Preserve contours - Regularity - Efficiency

R. Achanta et al. *SLIC superpixels compared to state-of-the-art superpixel methods*, PAMI, 2012.

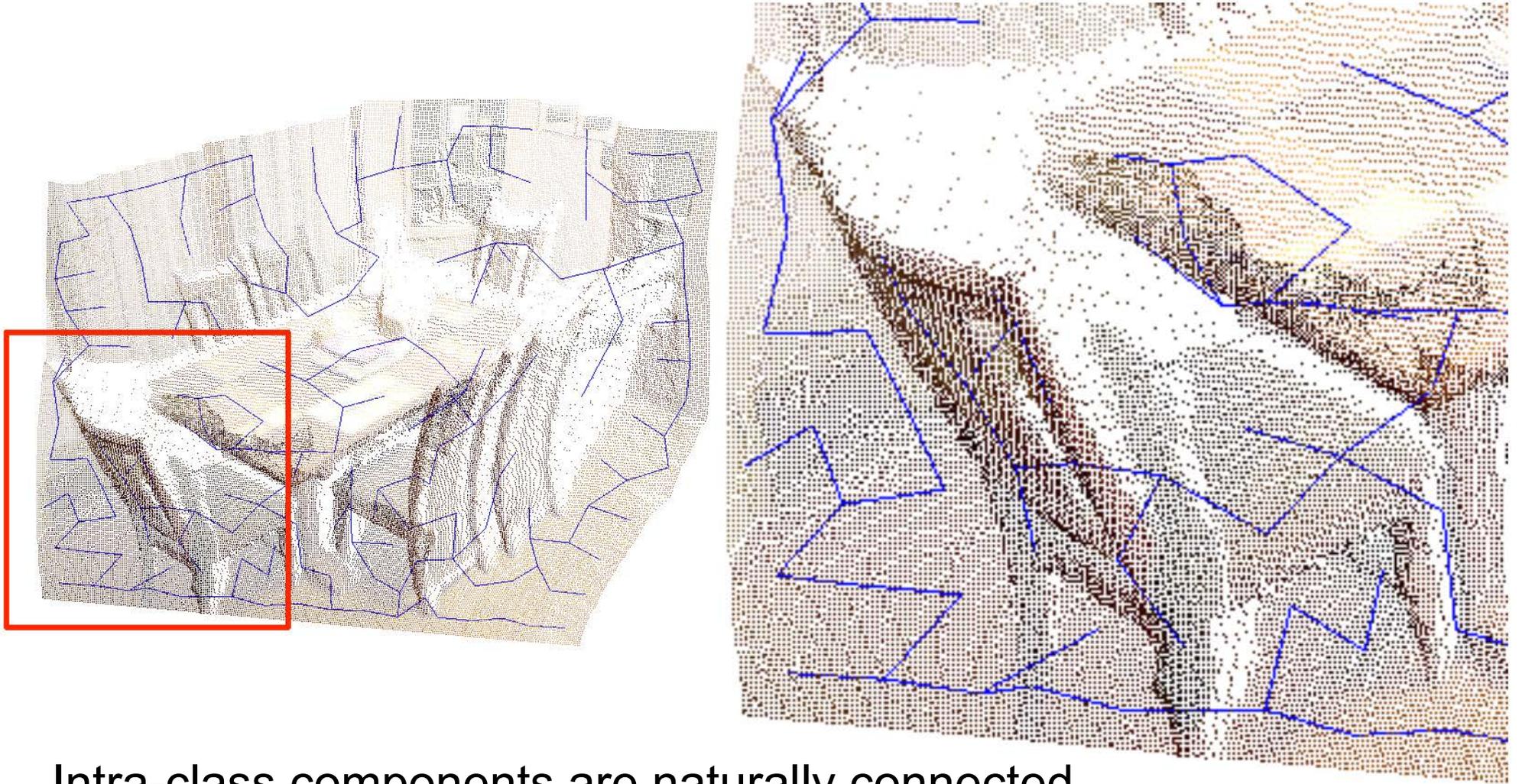
# Graph Structure: Our choice

Minimum Spanning Tree  
Over 3D



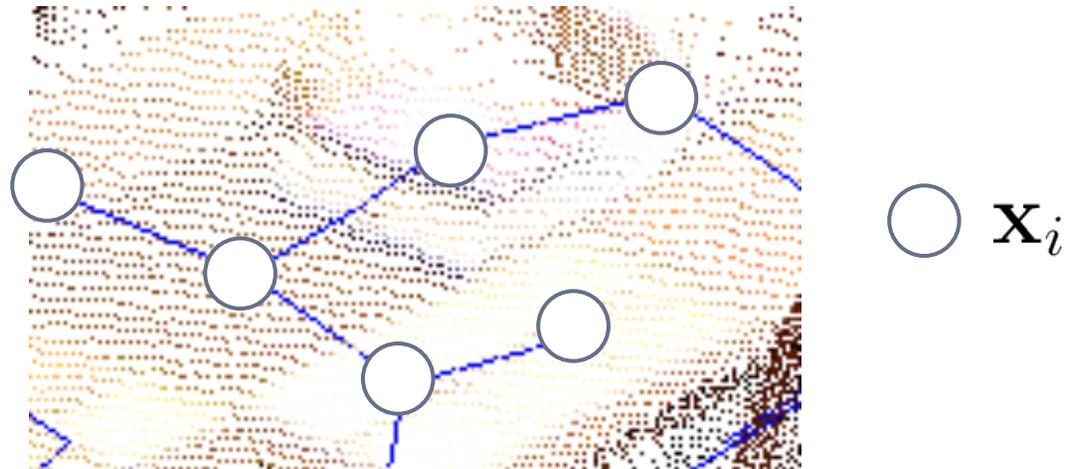
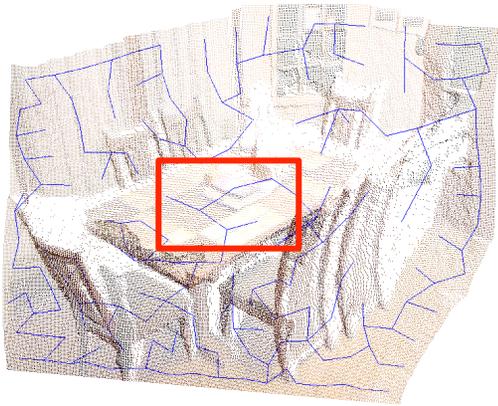
Edges are determined by the MST over 3D distances between superpixels' centroids

# Graph Structure: Our choice



Intra-class components are naturally connected

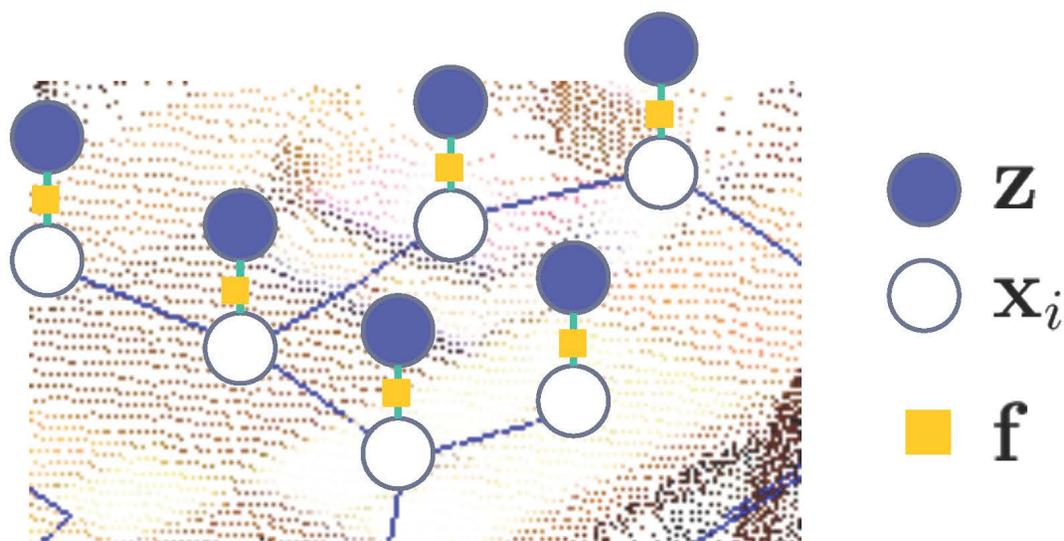
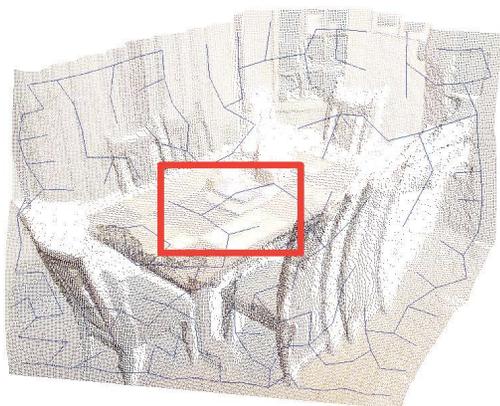
# Formulation as Pairwise CRFs



Directly models  $p(\mathbf{x}|\mathbf{z})$

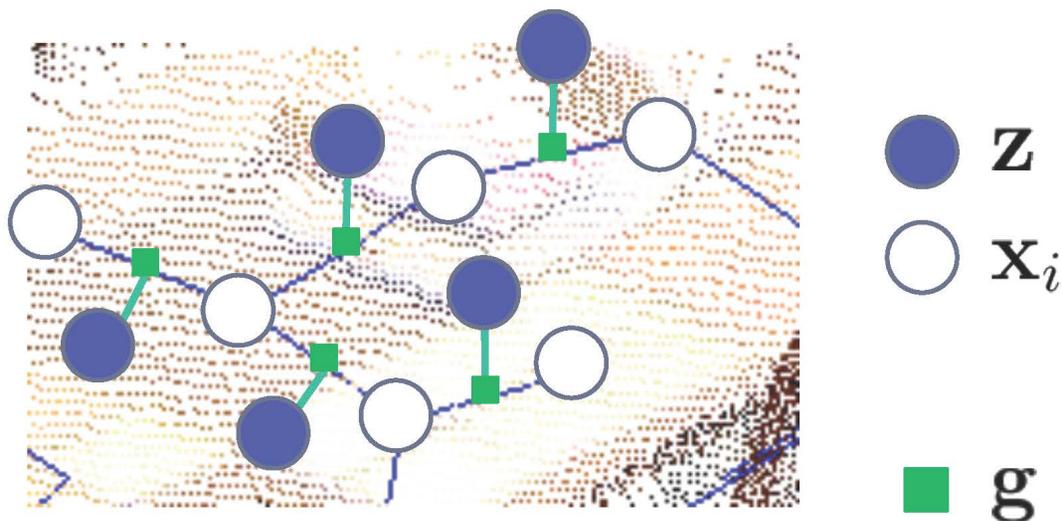
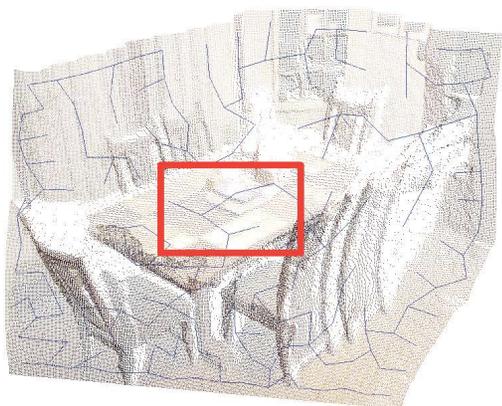
$$p(\mathbf{x}|\mathbf{z}) = \frac{1}{Z(\mathbf{z})} \exp \left( \sum_{i \in \mathcal{N}} \mathbf{w}_u^T \mathbf{f}(\mathbf{x}_i, \mathbf{z}) + \sum_{i, j \in \mathcal{E}} \mathbf{w}_p^T \mathbf{g}(\mathbf{x}_{i, j}, \mathbf{z}) \right)$$

# Potentials: Pairwise CRFs



$$p(\mathbf{x}|\mathbf{z}) = \frac{1}{Z(\mathbf{z})} \exp \left( \sum_{i \in \mathcal{N}} \mathbf{w}_u^T \mathbf{f}(\mathbf{x}_i, \mathbf{z}) + \sum_{i, j \in \mathcal{E}} \mathbf{w}_p^T \mathbf{g}(\mathbf{x}_{i, j}, \mathbf{z}) \right)$$

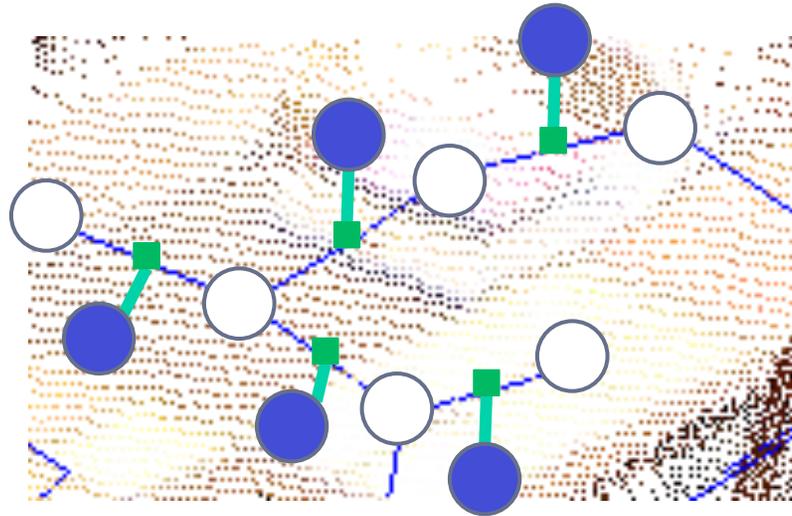
# Potentials: Pairwise CRFs



$$p(\mathbf{x}|\mathbf{z}) = \frac{1}{Z(\mathbf{z})} \exp \left( \sum_{i \in \mathcal{N}} \mathbf{w}_u^T \mathbf{f}(\mathbf{x}_i, \mathbf{z}) + \sum_{i, j \in \mathcal{E}} \mathbf{w}_p^T \mathbf{g}(\mathbf{x}_{i, j}, \mathbf{z}) \right)$$

# Potentials: pairwise

- Favor (penalize) same class for nodes close (far) in Lab color
- Favor (penalize) different classes for nodes far (close) in Lab color



$$\blacksquare \quad g(\mathbf{x}_{i,j}, \mathbf{z}) = \begin{cases} 1 - \exp(-\|c_i - c_j\|_2) & \rightarrow l_i = l_j \\ \exp(-\|c_i - c_j\|_2) & \rightarrow l_i \neq l_j \end{cases}$$

$c$  : Lab color

Same form with 3D positions

# Potentials: unary

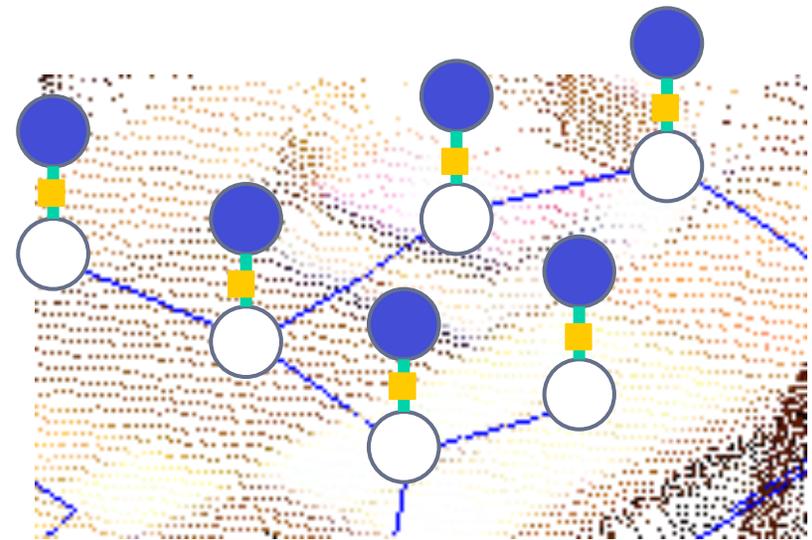
unary (local) potential  
using a k-NN classifier

$$\blacksquare f(\mathbf{x}_i, \mathbf{z}) = -\log P_i(\mathbf{x}_i | \mathbf{z})$$

$$P_i(\mathbf{x}_i = l_j | \mathbf{z}) = \eta \frac{f(l_j) \bar{F}(l_j)}{\bar{f}(l_j) F(l_j)}$$

$f(l_j)$  frequency of label  $j$  in a k-NN query

$F(l_j)$  frequency of label  $j$  the database



The database is a kd-tree of features from training data

# Features

Indoors (15D)

Outdoors(21D)

## From Image

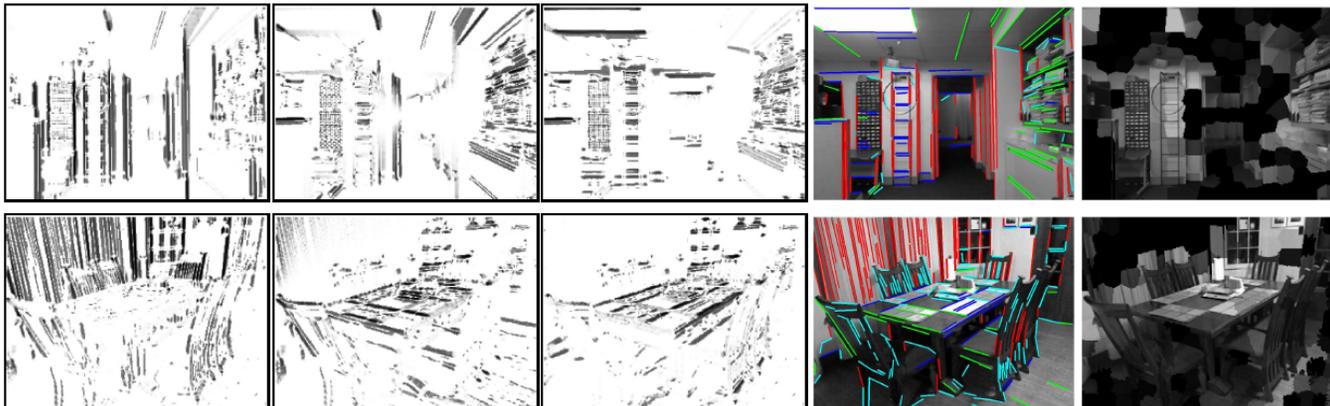
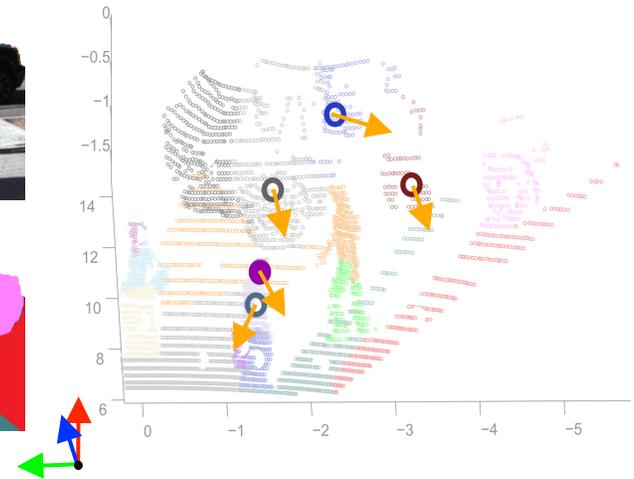
Lab-color: mean and std	6D	6D
RGB-color: mean and std		6D
vertical centroid location	1D	1D
entropy from vanishing points	1D	

## From 3D

3-D centroid position	2D	3D
differences on depth: mean and std	2D	2D
local planarity	1D	1D
neighboring planarity	1D	1D
vertical orientation	1D	1D

# Features

- From 3D
  - mean and std of differences on depth
  - local planarity
  - neighboring planarity
  - vertical orientation
- From Image:
  - entropy from vanishing points



$$H_s = - \sum_{j=1}^4 h g_s(y = j) \log (h g_s(y = j))$$

# Inference

- We use belief propagation:
  - Exact results in MAP/marginals
  - Efficient computation, in  $\mathcal{O}(nm^2)$

# Learning

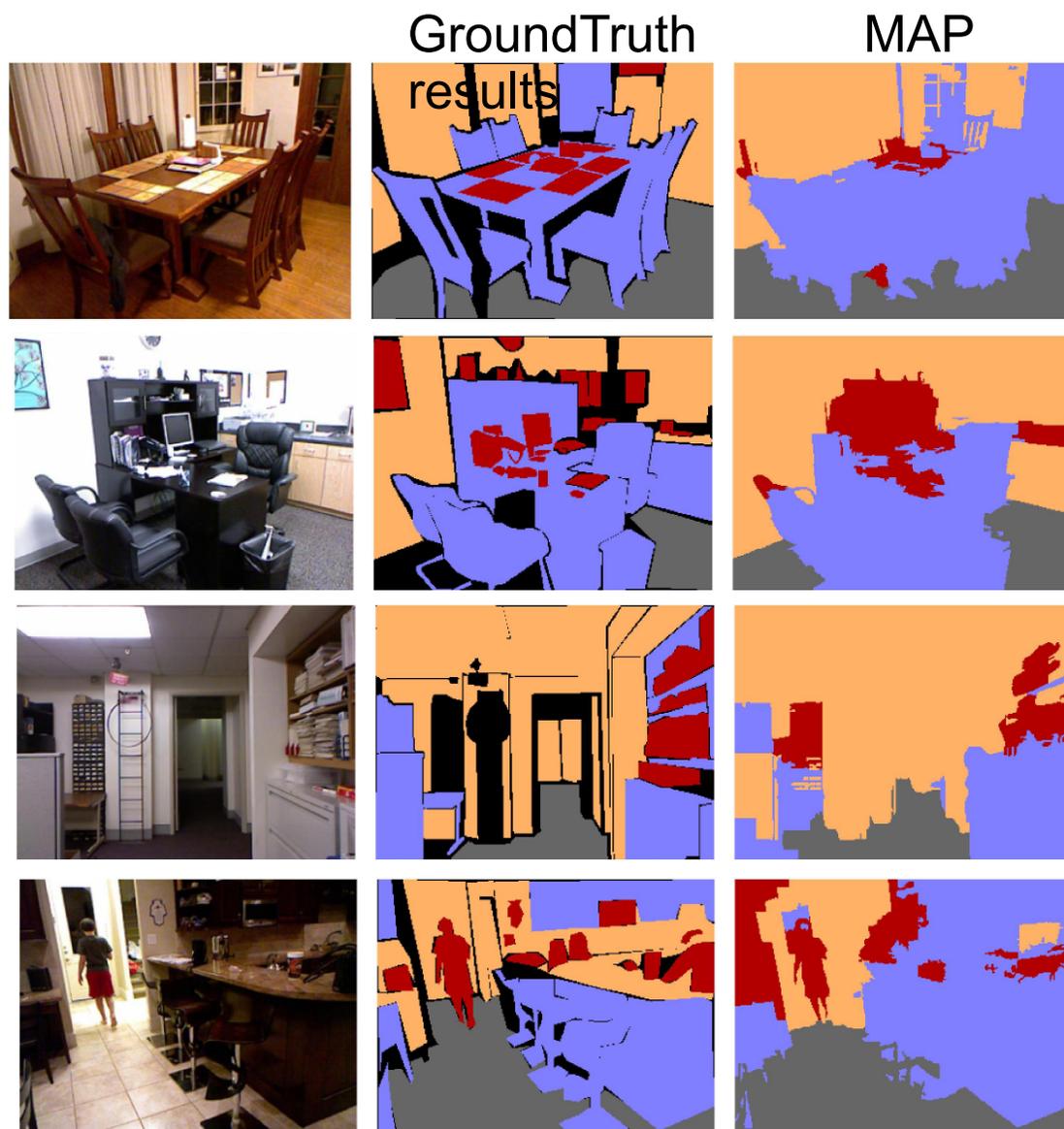
- Maximum Likelihood Estimation

- To learn  $[\mathbf{w}_u, \mathbf{w}_p]$

Tree graph structure:  
good convergence

# Results: NYU-Depth v2 Dataset

Qualitative comparison:



Ground  
Structure

Furniture

Props

# Results: NYU-Depth v2 Dataset

Quantitative comparison:

- Recall accuracy in pixel-wise percentage:

	Ground	Furniture	Props	Structure	Average	Global
CRF-MST-kNN	88.4	64.1	30.5	78.6	<b>65.4</b>	<b>67.2</b>
only Image Feat.	63.2	47.5	24.5	73.6	52.2	56.1
only 3D Feat.	<b>89.5</b>	<b>70.0</b>	16.9	79.4	62.7	65.8
data-term (kNN)	87.3	60.6	33.7	74.8	64.1	64.9
Silberman et al. 2012	68	<b>70</b>	<b>42</b>	59	59.6	58.6
Couprie et al. 2013	87.3	45.3	35.5	<b>86.1</b>	63.5	64.5

# Results: KITTI Dataset

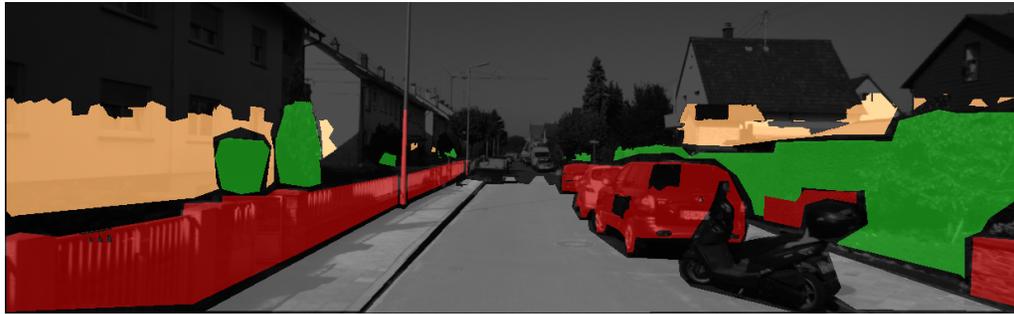
- Recall accuracy in pixel-wise percentage:

	Ground	Objects	Building	Vegetation	Average	Global
CRF-MST-kNN	<b>97.3</b>	82.9	<b>82.8</b>	86.9	<b>87.5</b>	<b>88.4</b>
only Image Feat.	96.8	49.2	64.6	<b>95.5</b>	76.5	76.8
only 3D Feat.	95.9	<b>84.2</b>	80.5	46.7	76.8	78.8
data-term (kNN)	96.8	75.9	80.7	77.6	82.8	83.5

# Results: KITTI Dataset

Ground Truth

MAP results



Ground  
Objects

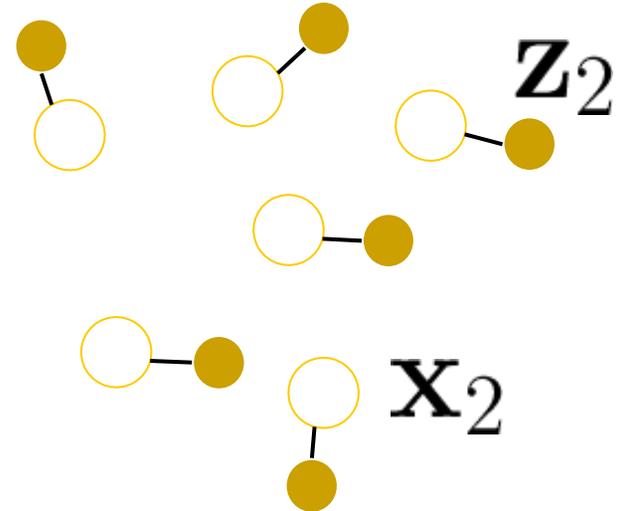
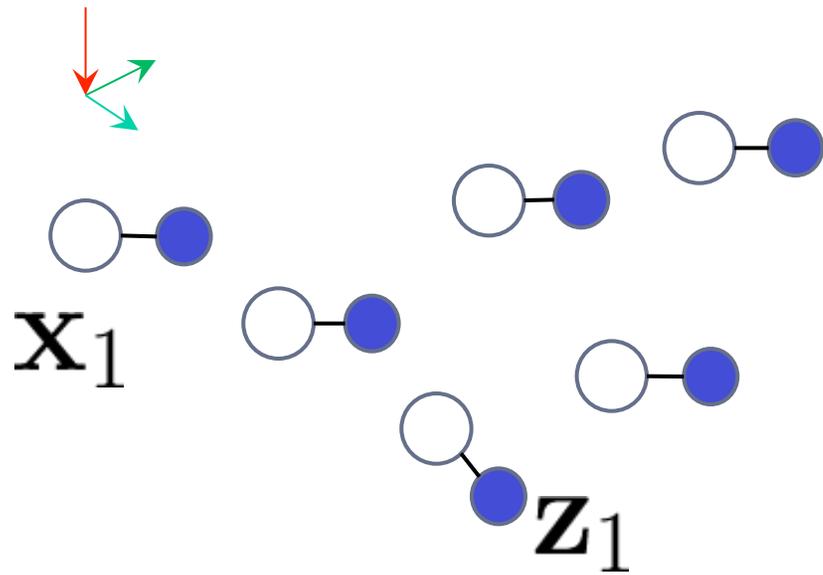
Buildings

Vegetation

# Overview

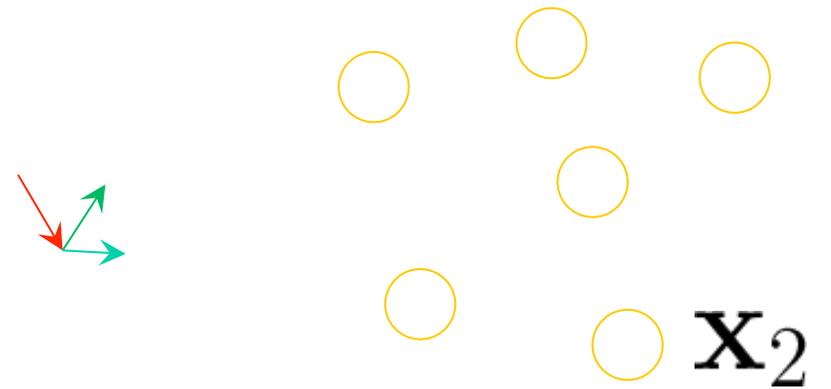
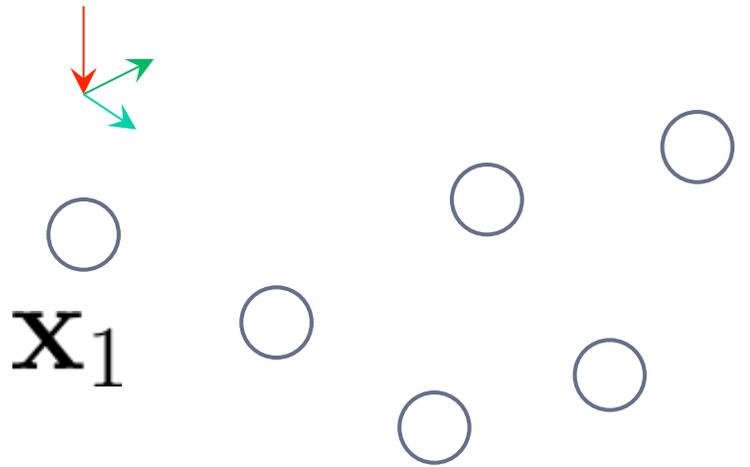
1. Task Dependent Semantic Hierarchy
- 2. Multi-view and Recursive formulation**
3. Capability of handling missing 3D data (full sensor's FOVs coverage)

# Multi-view:

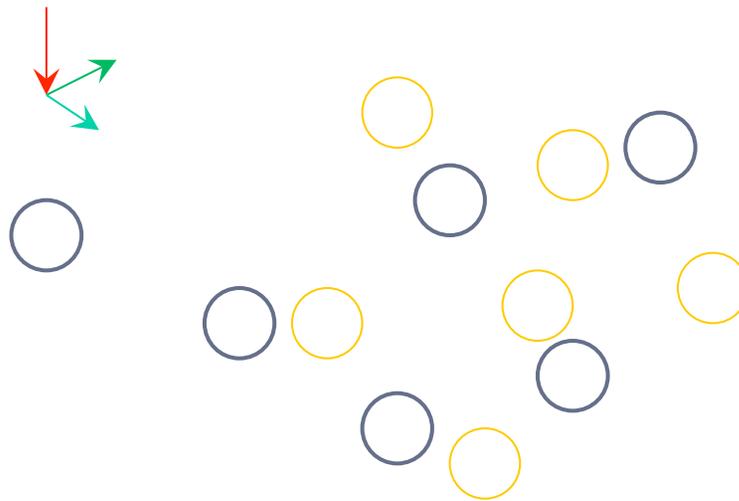


Different views/nodes in their local reference frame – mean 3D position

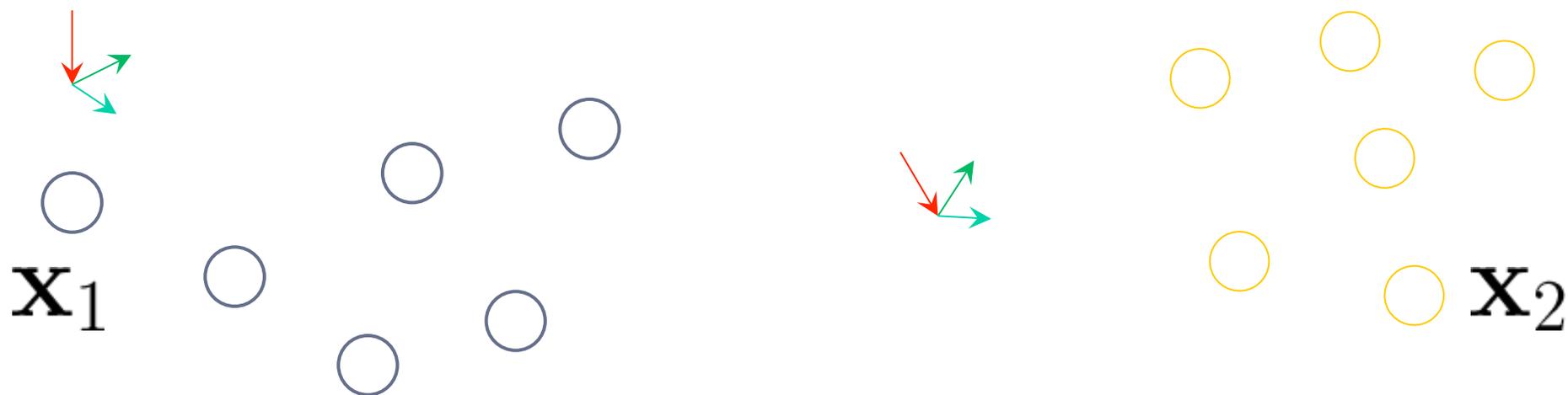
# Multi-view:



Use relative pose to align the nodes to the same reference frame

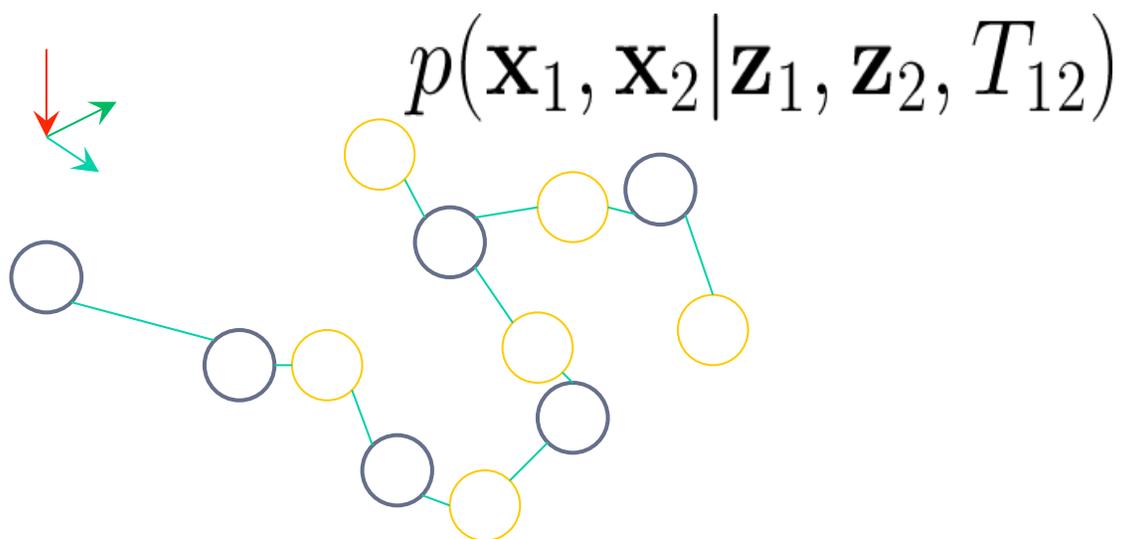


# Multi-view:

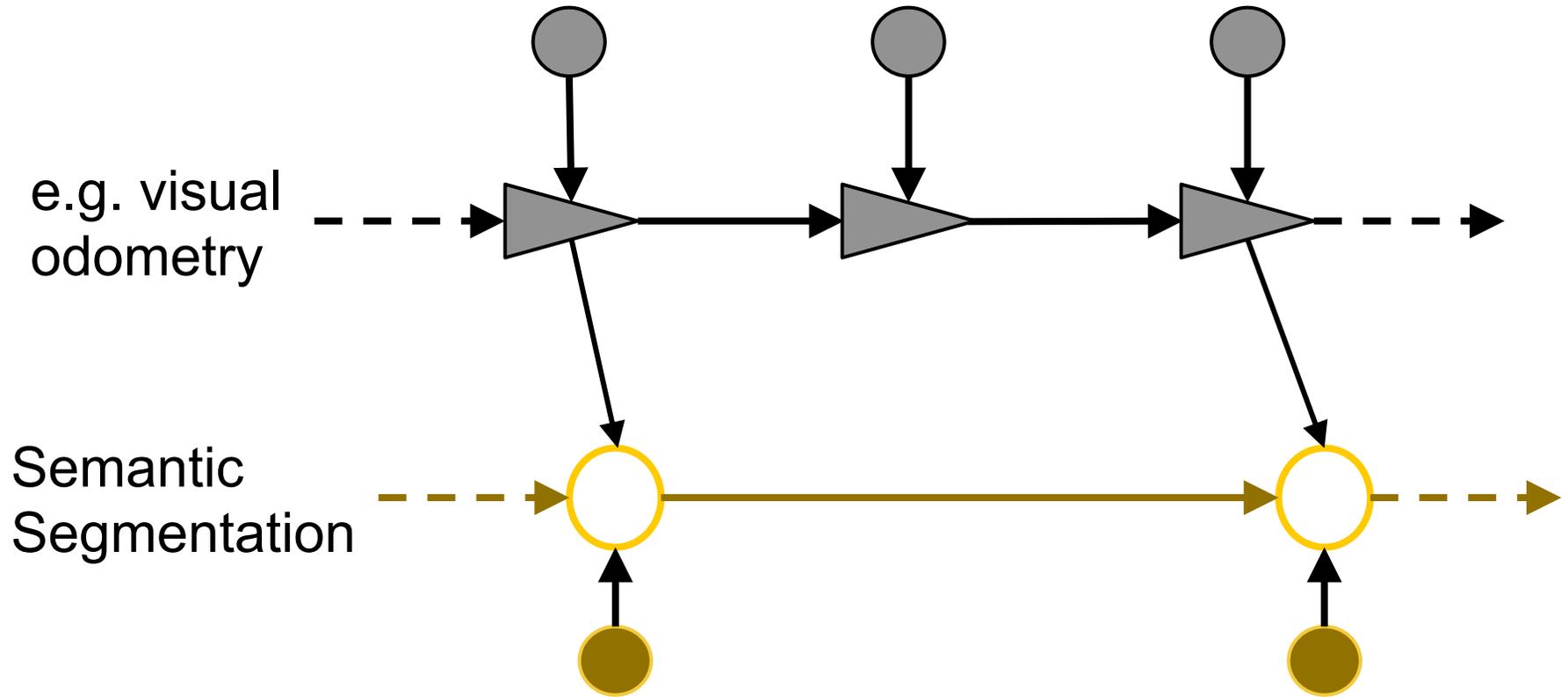


MST in the common  
reference frame

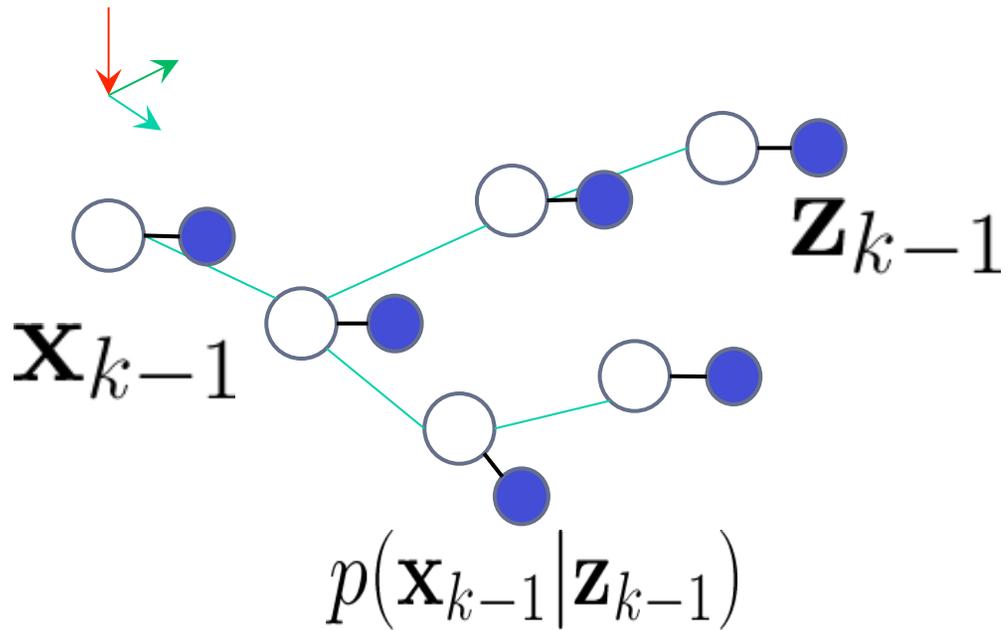
No need to find hard  
correspondences



# Recursive Inference



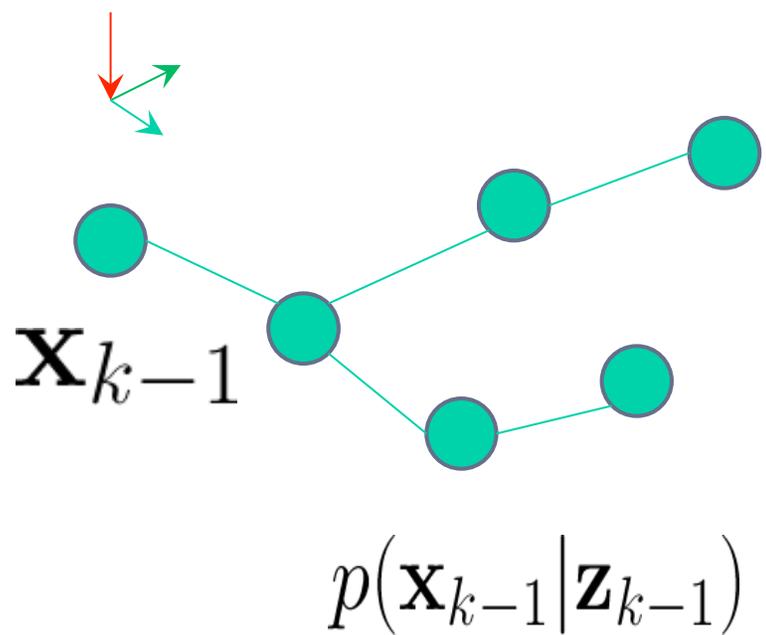
# Video sequences:



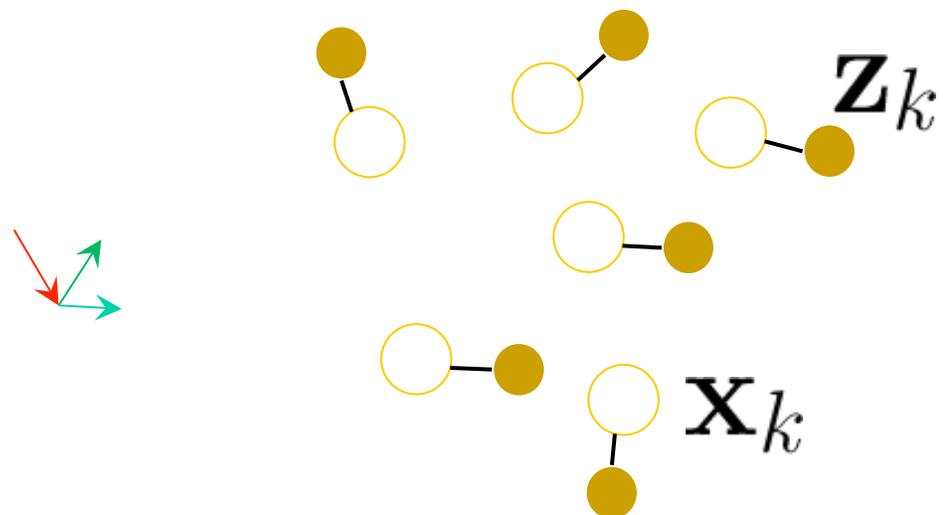
On-line operation

Infer the marginals at  
 $k-1$

# Recursive Inference

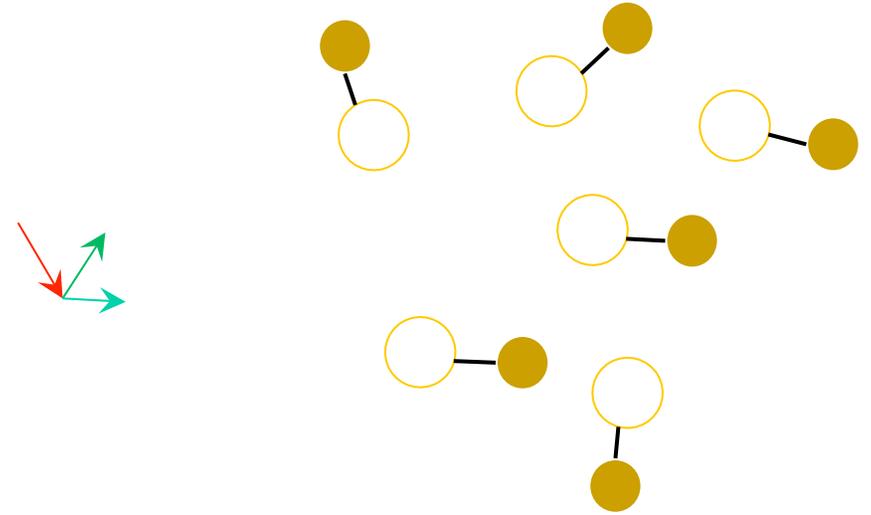
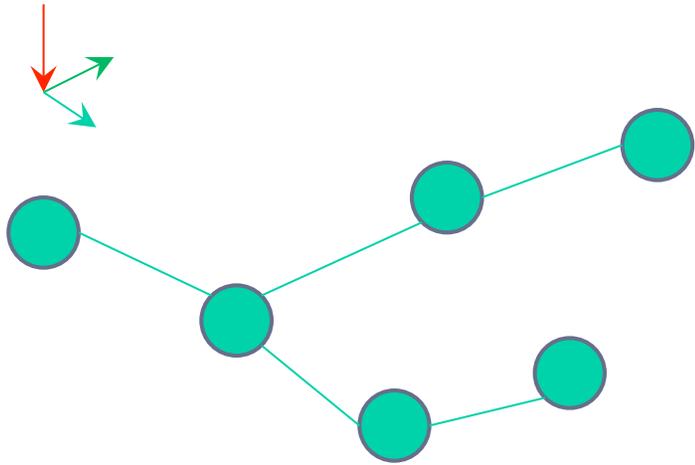


Marginals at  $k-1$



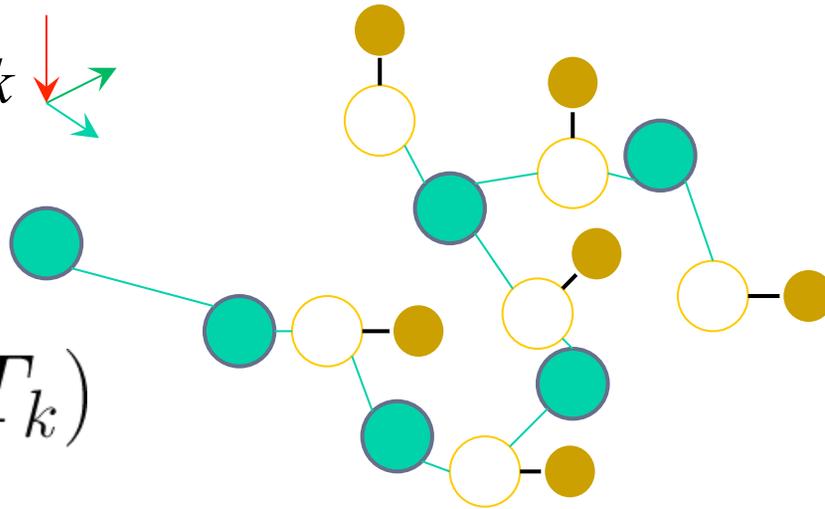
Sensing at time  $k$

# Recursive Inference



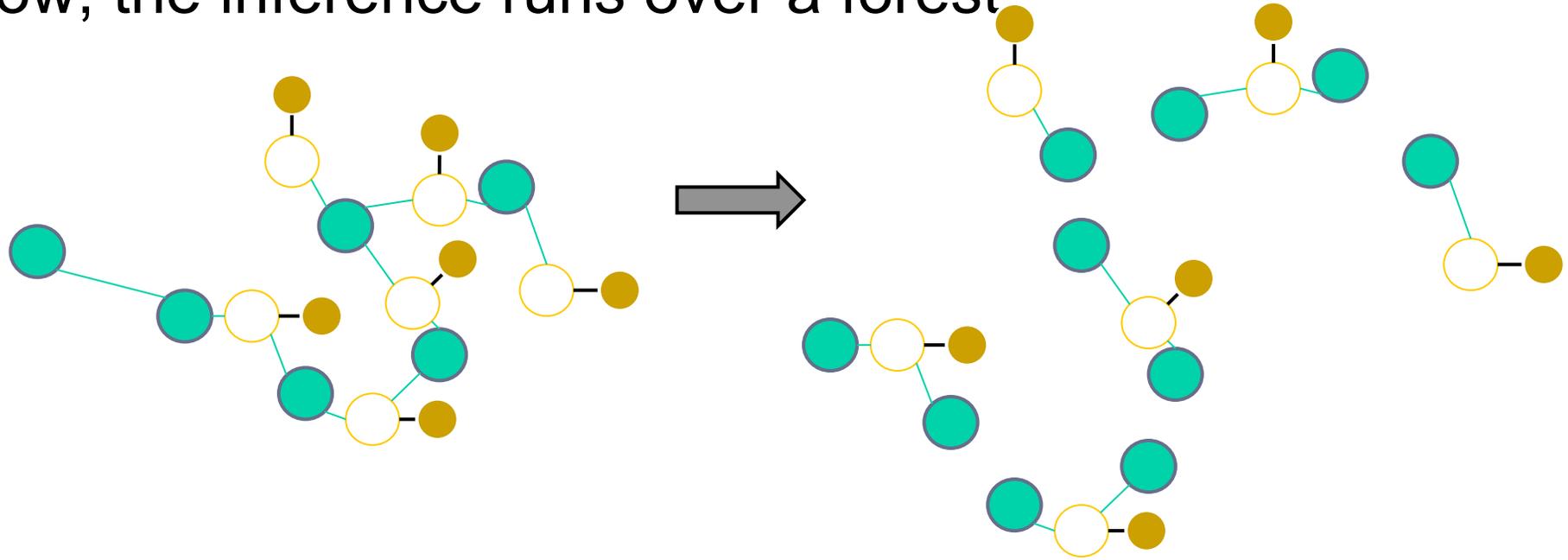
Infer marginals at time  $k$

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{z}_k, T_k)$$



# Recursive Inference

- Now, the inference runs over a forest



F1 KITTI Dataset	Ground	Objects	Building	Vegetation	Time MST	Time BP
Single View	0.977	0.854	0.870	0.811	21 ms	164 ms
Recursive Inference	0.977	0.853	0.879	0.809	57 ms	69 ms

# Single view vs. Recursive mode

Original  
video



Single view  
mode



$$I_m \times p(\mathbf{x} = \text{objects})$$

Recursive  
mode



$$I_m \times p(\mathbf{x} = \text{objects})$$

Sequence:  
2011\_09\_29/2011\_09\_29\_drive\_0071

# Test in Dynamic Street - KITTI



Removing the object class  
from the reprojected  
pointcloud:

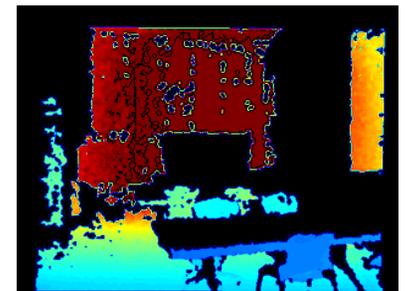
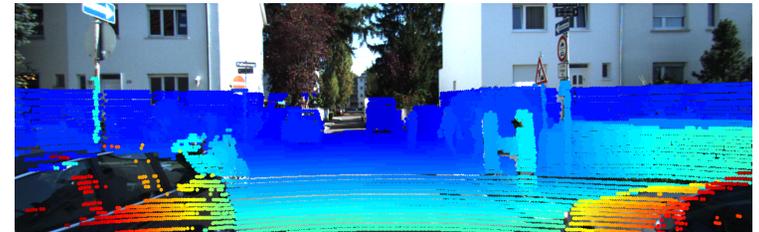


# Overview

1. Task Dependent Semantic Hierarchy
2. Multi-view and Recursive formulation
3. Capability of handling missing 3D data (full sensor's FOVs coverage)

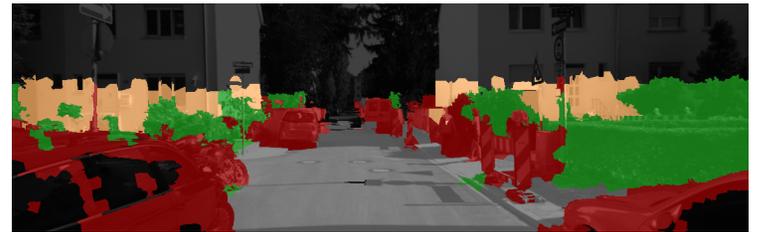
# Multiple Sensor Modalities

- Different fields of view, missing 3D data
- Every sensor suffers from specific blind spots, e.g.
  - Laser: limited range, specular surfaces
  - Vision: low light conditions
  - Depth (Kinect): natural light, specularities
- Every modality suffers from different sources of ambiguities

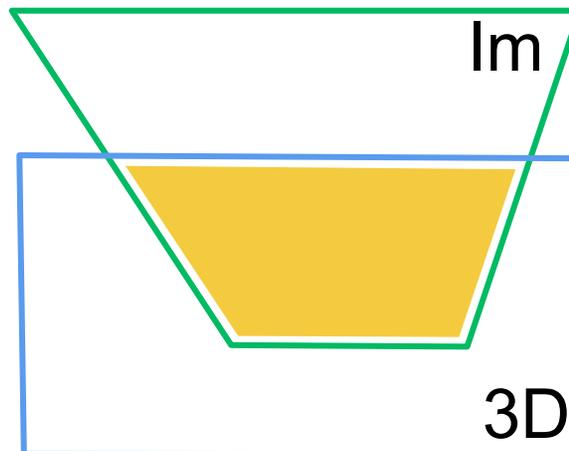


# Multiple Sensor Modalities

- So far, we have dealt with the ‘semantic’ ambiguities fusing image and 3D sensors

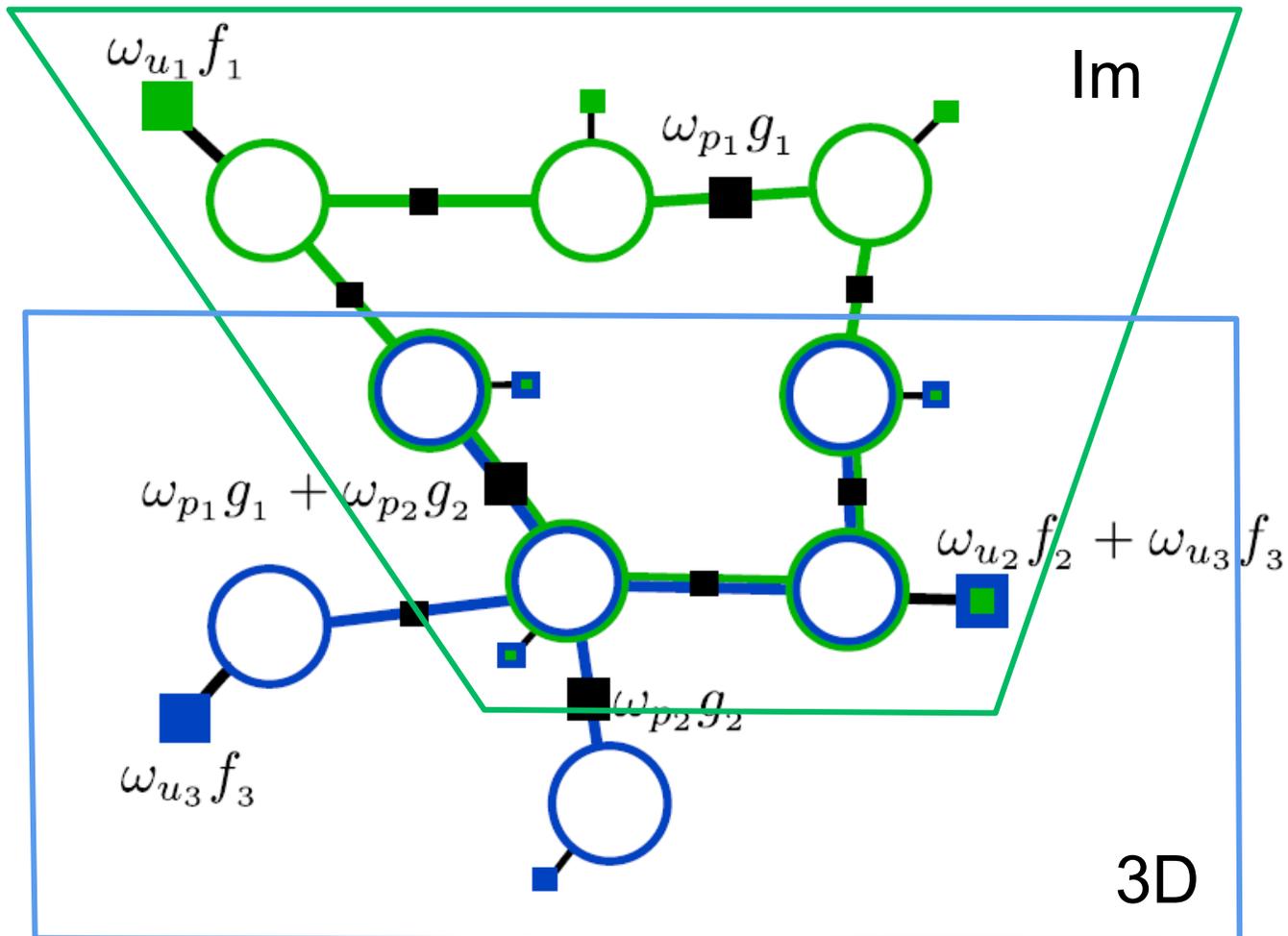


- But still, only over the common spatial coverage, and without handling missing data



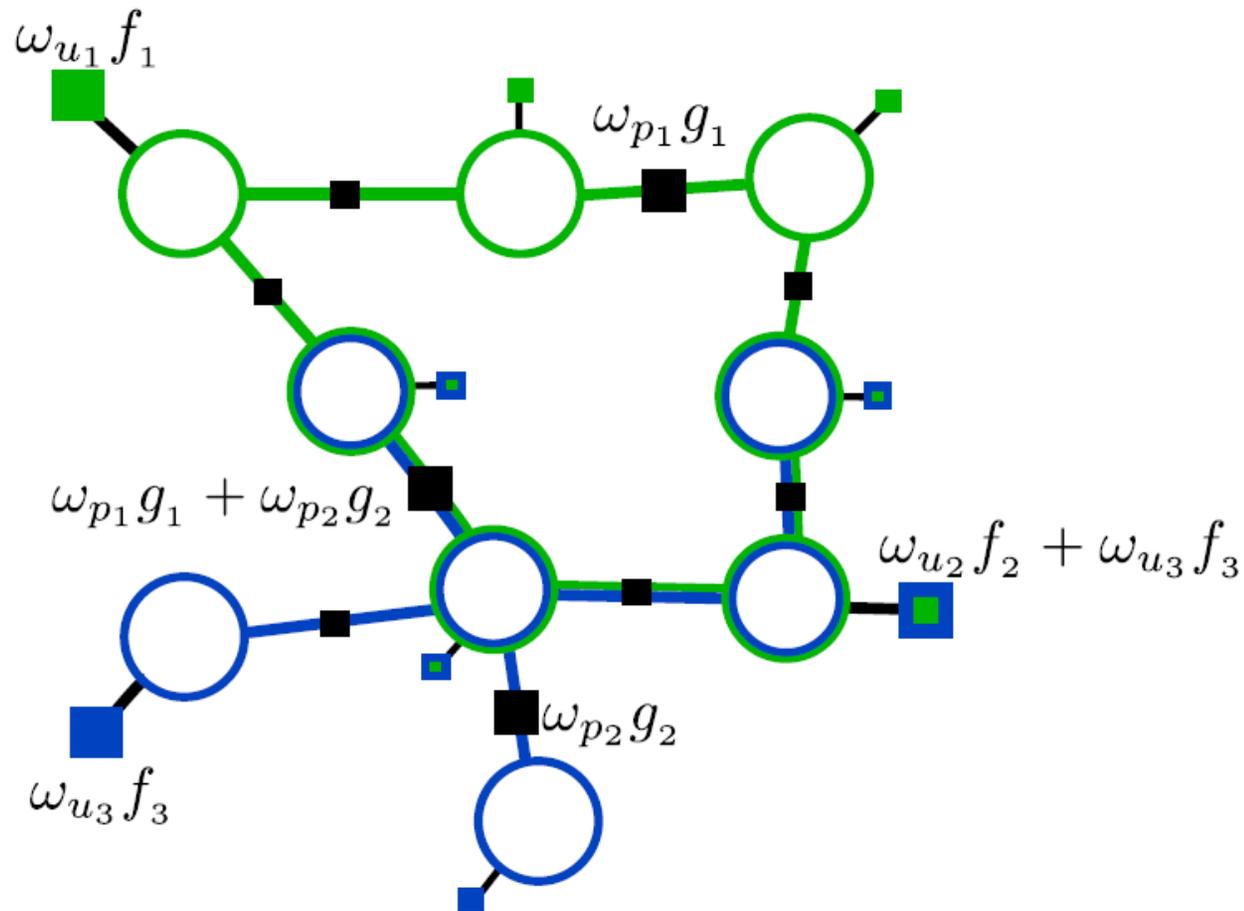
# Full Coverage - Formulation

- Define a graph structure for full coverage



# Full Coverage - Formulation

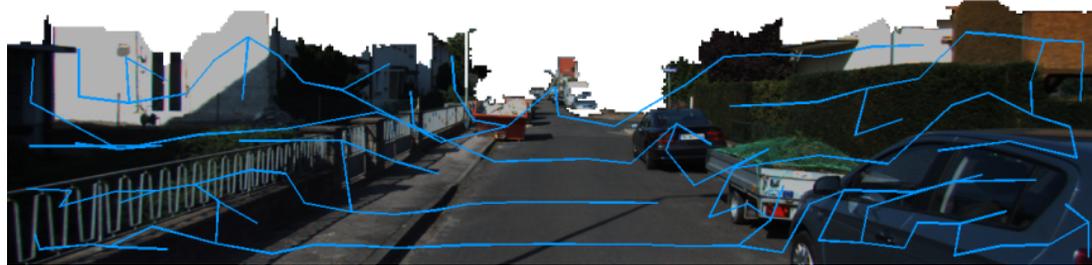
- Define a graph structure for full coverage
- Domain based potentials



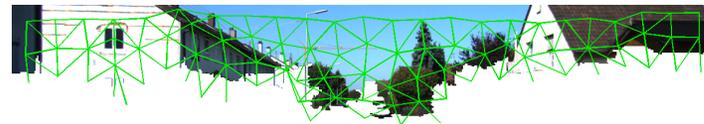
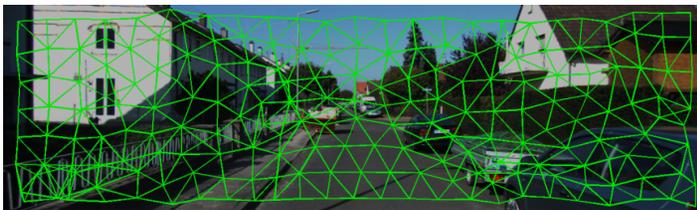
# Graph Structure

- Start with the MST over 3D distances as before

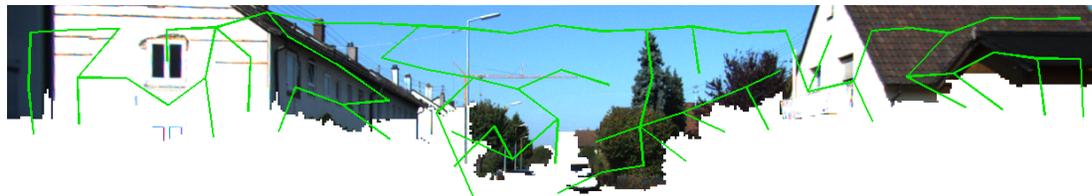
3D



- Look for the image graph on superpixels without 3D

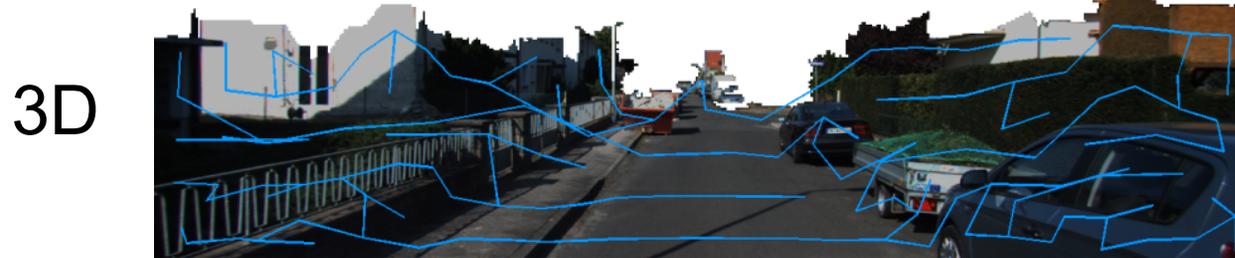


- And compute the MST over Lab-color distances

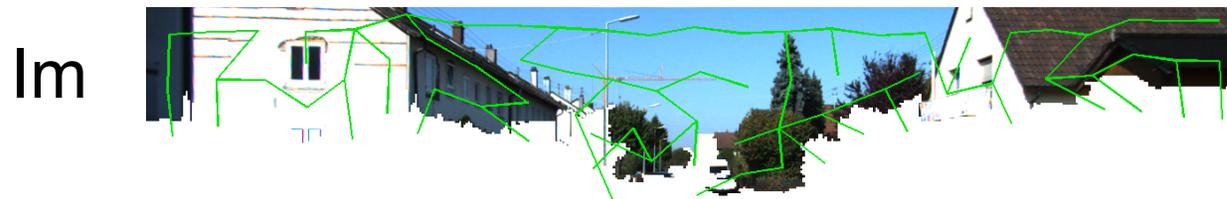


# Graph Structure

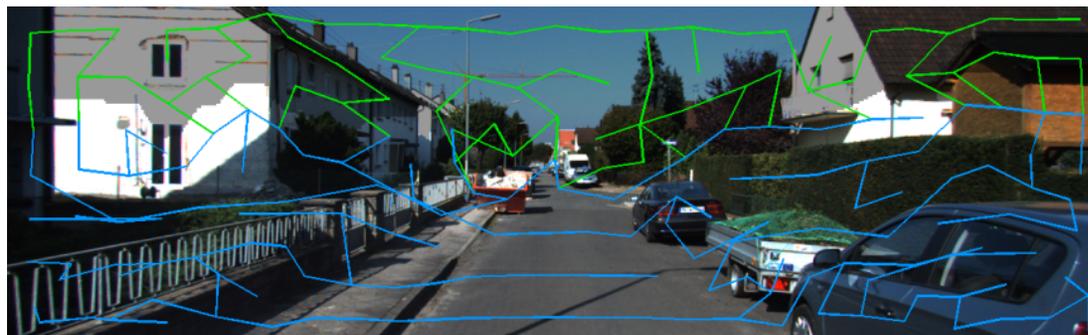
- Sub-graph over 3D



- Sub-graph over Image

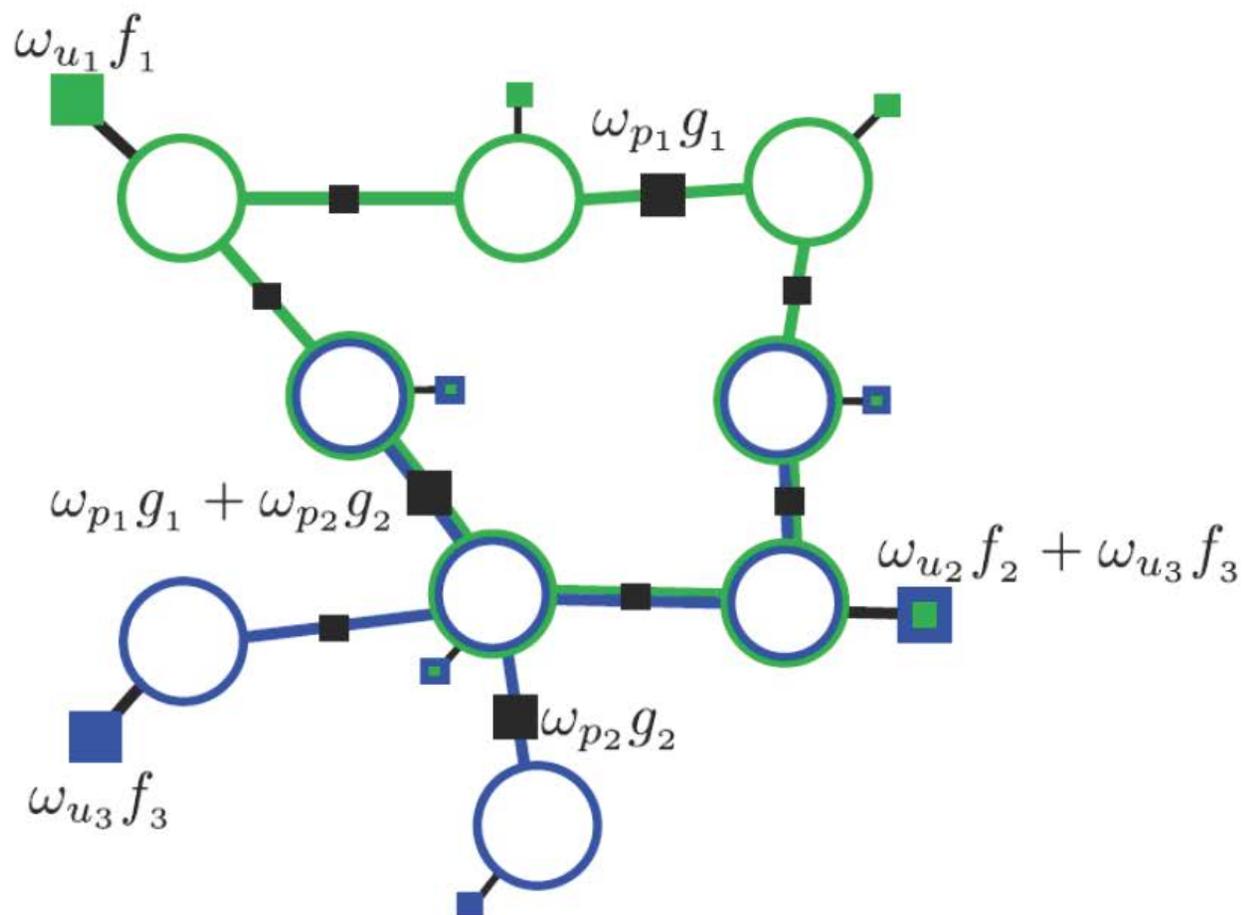


- Results in a graph for full coverage



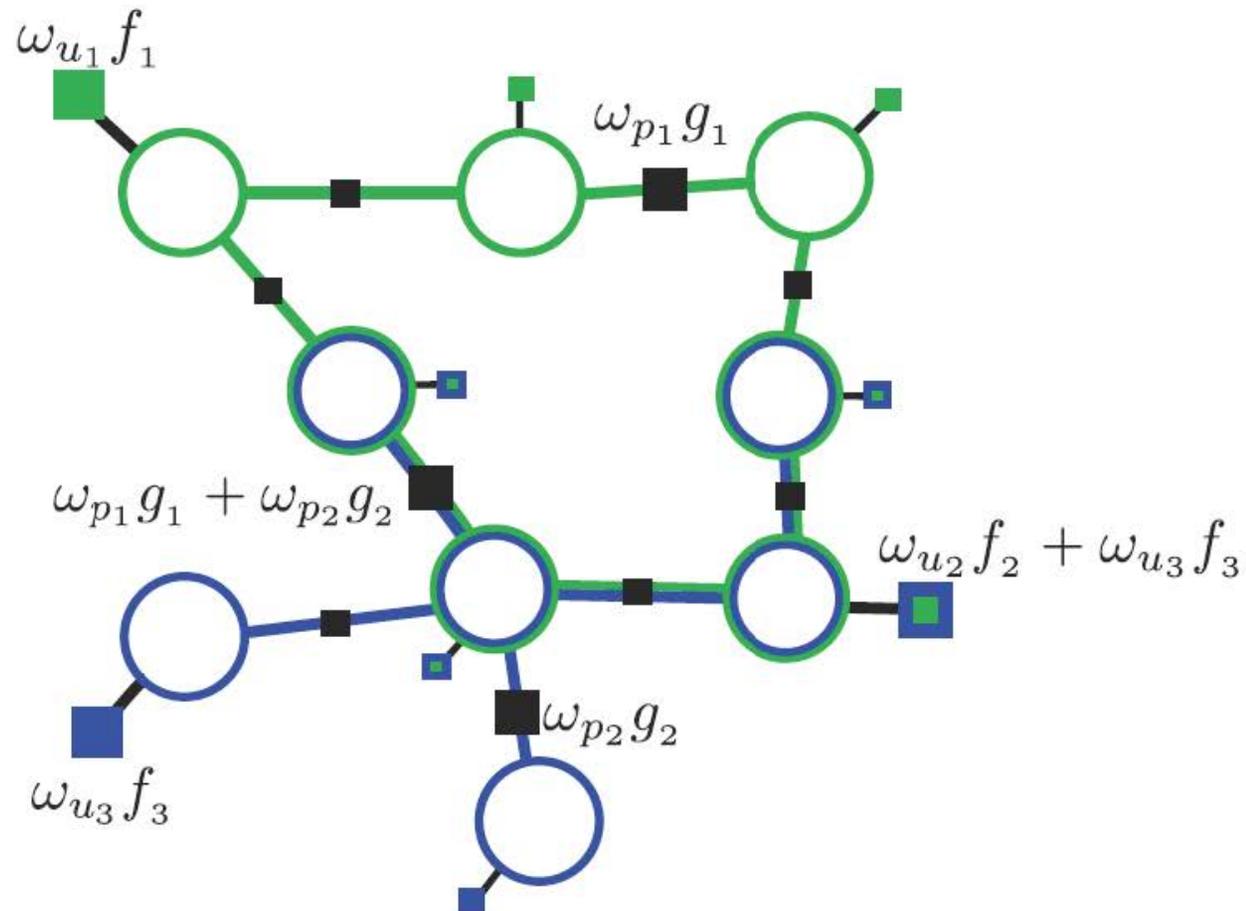
# CRFs Formulation

$$p(\mathbf{x}|\mathbf{z}) = \frac{1}{Z(\mathbf{z})} \exp \left( \sum_{i \in \mathcal{N}} \mathbf{w}_u^T \mathbf{f}(\mathbf{x}_i, \mathbf{z}) + \sum_{i,j \in \mathcal{E}} \mathbf{w}_p^T \mathbf{g}(\mathbf{x}_{i,j}, \mathbf{z}) \right)$$



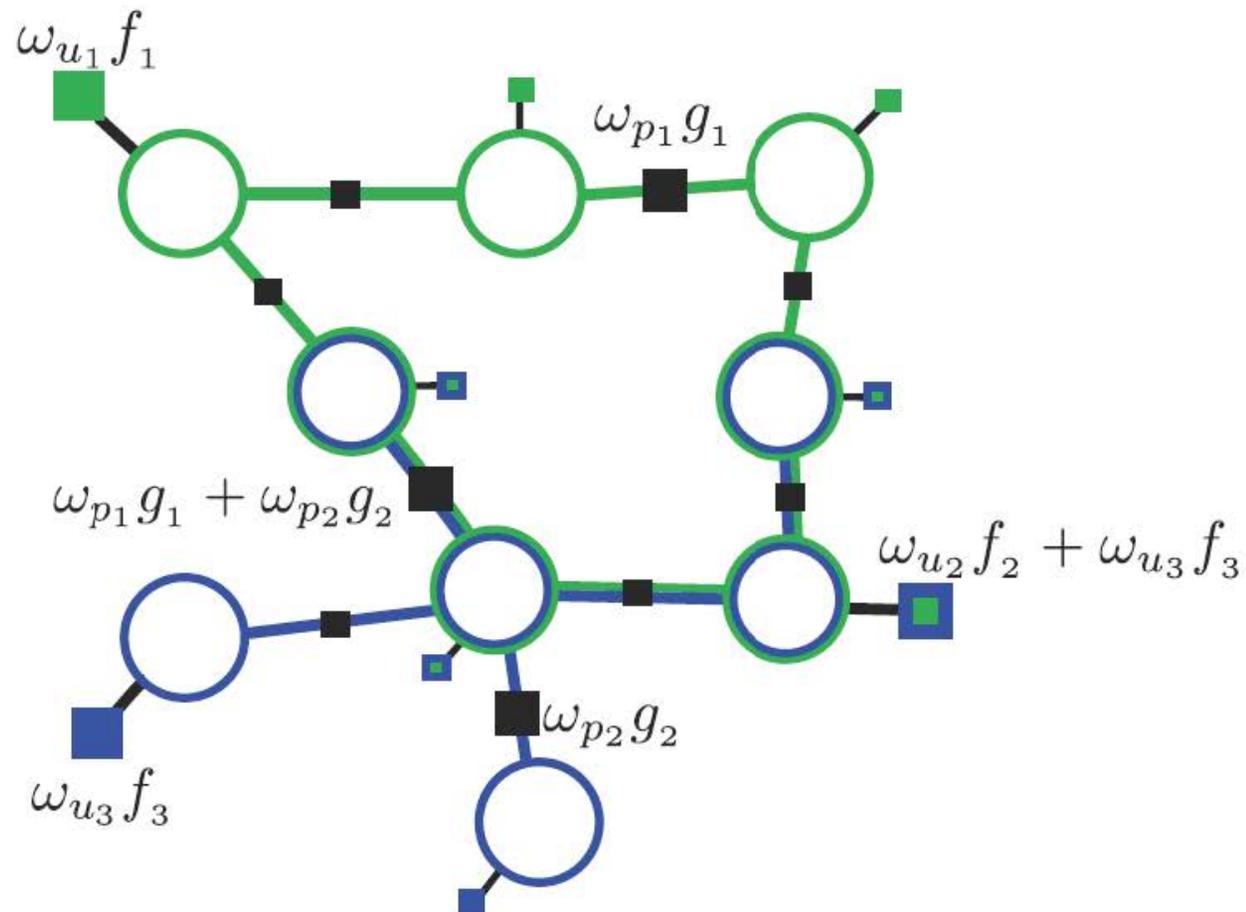
# Unary Potentials

$$\sum_{i \in \mathcal{N}_{Im \setminus 3D}} \omega_{u_1} f_1(\mathbf{x}_i, \mathbf{z}_{Im}) + \sum_{i \in \mathcal{N}_{Im \cap 3D}} \omega_{u_2} f_2(\mathbf{x}_i, \mathbf{z}_{Im, 3D}) + \sum_{i \in \mathcal{N}_{3D}} \omega_{u_3} f_3(\mathbf{x}_i, \mathbf{z}_{3D})$$

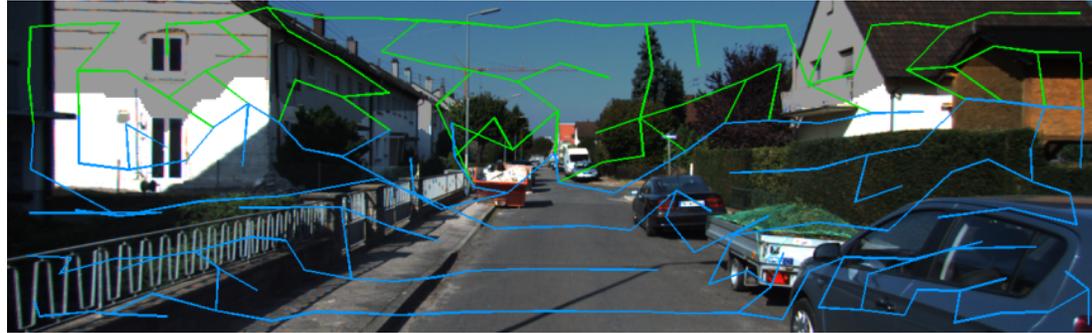


# Pairwise Potentials

$$\sum_{i,j \in \mathcal{E}_{Im}} \omega_{p_1} g_1(\mathbf{x}_{i,j}, \mathbf{z}_{Im}) + \sum_{i,j \in \mathcal{E}_{3D}} \omega_{p_2} g_2(\mathbf{x}_{i,j}, \mathbf{z}_{3D})$$



# Inference and Learning



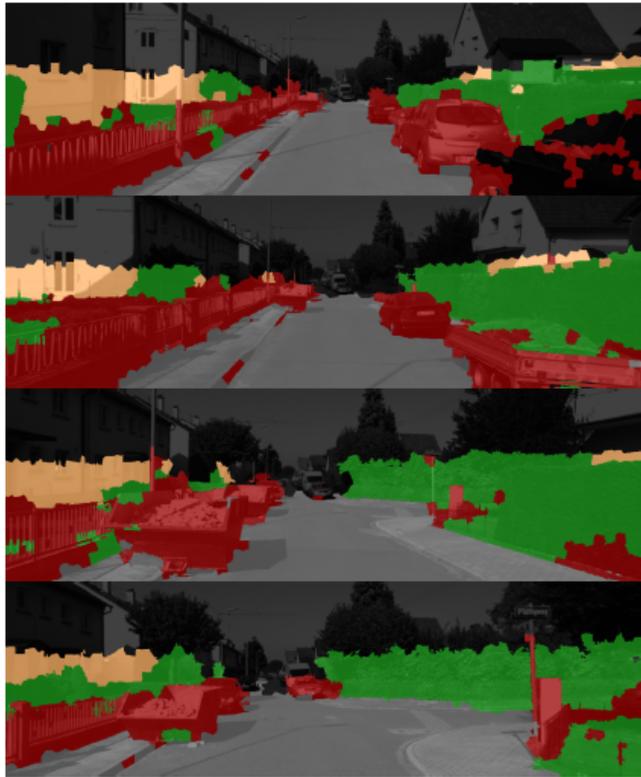
- This graph is not a tree
- We use Loopy Belief Propagation for Inference
- Our experiments always have converged in less than 5 iterations
- Learning with MPL

# Results: KITTI Dataset

Im



Im and 3D



Im or 3D



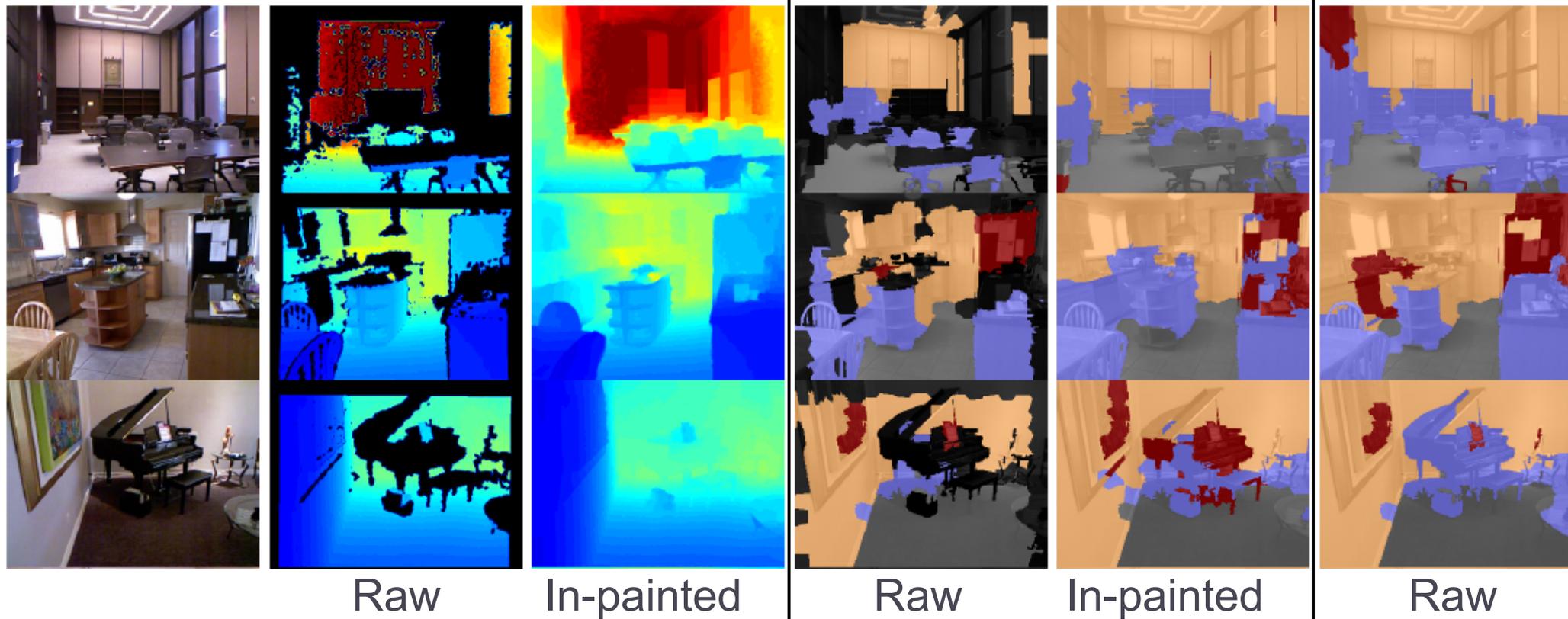
# Results: KITTI Dataset

- Recall accuracy in pixel-wise percentage:

	ground		objects					building	vegetation	sky	Average	Global	Coverage
	road	pavem.	car	fence	post	people	sign.						
Image only													
Sengupta et al. [21]	98.3	91.3	93.9	48.5	49.3	—	—	97.0	93.4	—	81.7	88.4	—
CRF-Im	97.8		61.1					87.4	94.6	97.6	87.7	85.5	100
CRF-Im $\cap$ 3D <sup>a</sup> [3]	97.3		82.9					82.8	86.9	—	87.5	88.4	60.1
CRF-Im $\cup$ 3D	96.6		83.6					86.1	94.3	97.2	91.6	90.1	100

Best performance in average and global accuracies

# Results: NYU Depth V2 Dataset



# Results: NYU Depth V2 Dataset

- Recall accuracy in pixel-wise percentage:

	ground	furniture	props	structure	Average	Global	In-painting	Coverage
Silberman et al. [22]	68	70	42	59	59.6	58.6	Required	100
Couprie et al. [5]	87.3	45.3	35.5	86.1	63.5	64.5	Required	100
CRF-Im $\cap$ 3D [4]	88.4	64.1	30.5	78.6	65.4	67.2	Required	100
CRF-Im $\cap$ 3D raw-depth <sup>b</sup>	88.5	69.0	23.1	78.6	64.8	67.4	No	74.6
CRF-Im $\cup$ 3D	87.9	63.8	27.1	79.7	64.3	67.0	No	100

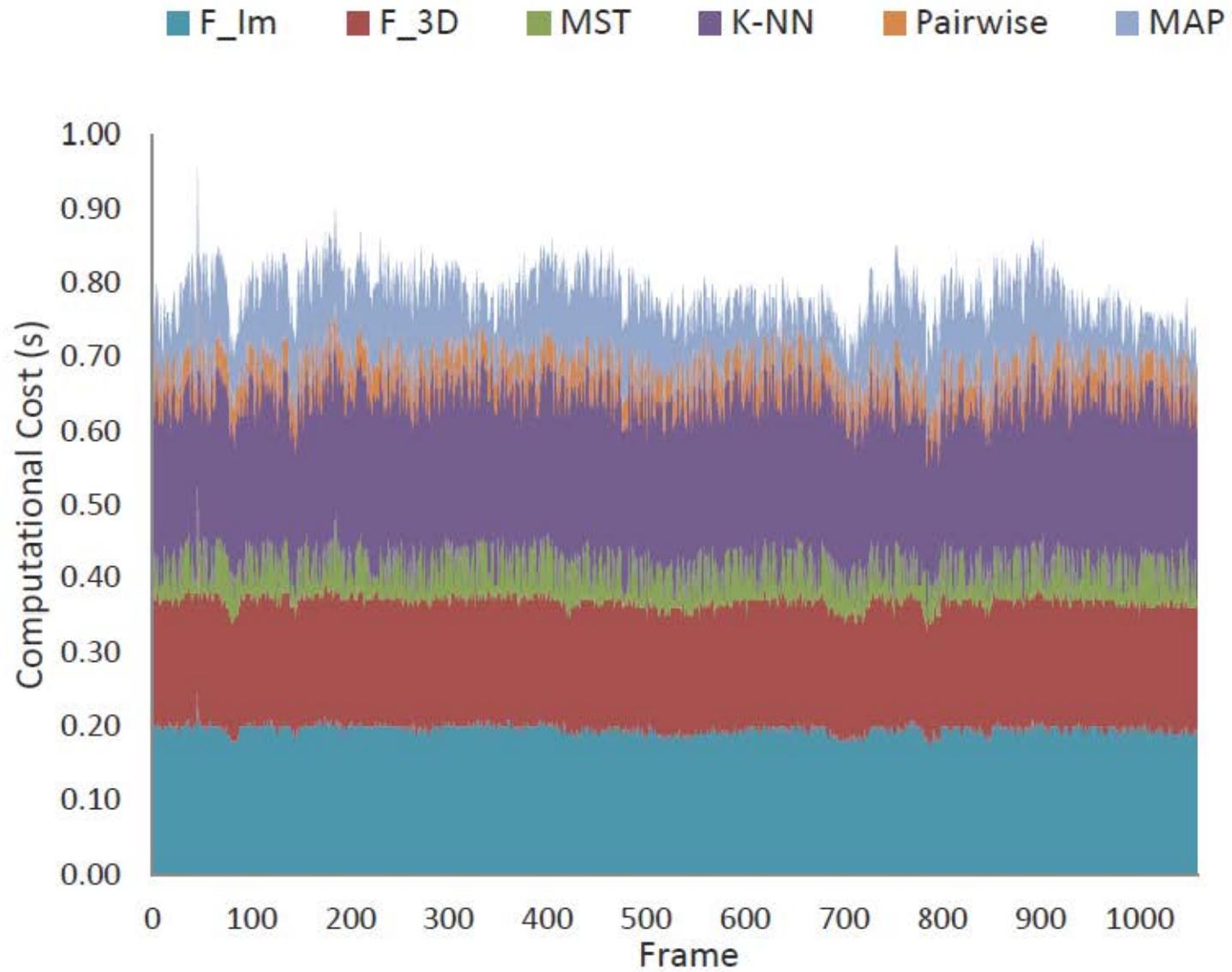
State of the art performance in full coverage

without In-painting (15s)

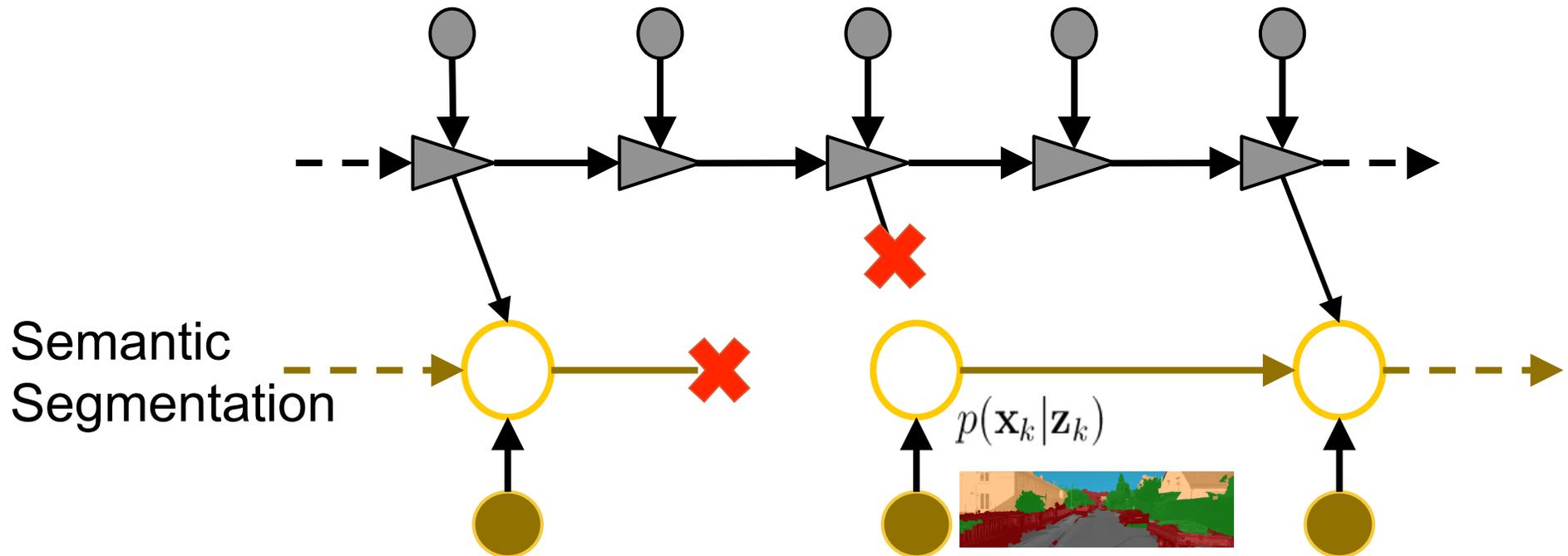
# Overview

1. Proposed Semantic Hierarchy
2. Basic Formulation: Graph, Appearance and 3D data
3. Multiview and Recursive formulation
4. Full sensor's FOVs coverage
5. Conclusions

# Timing

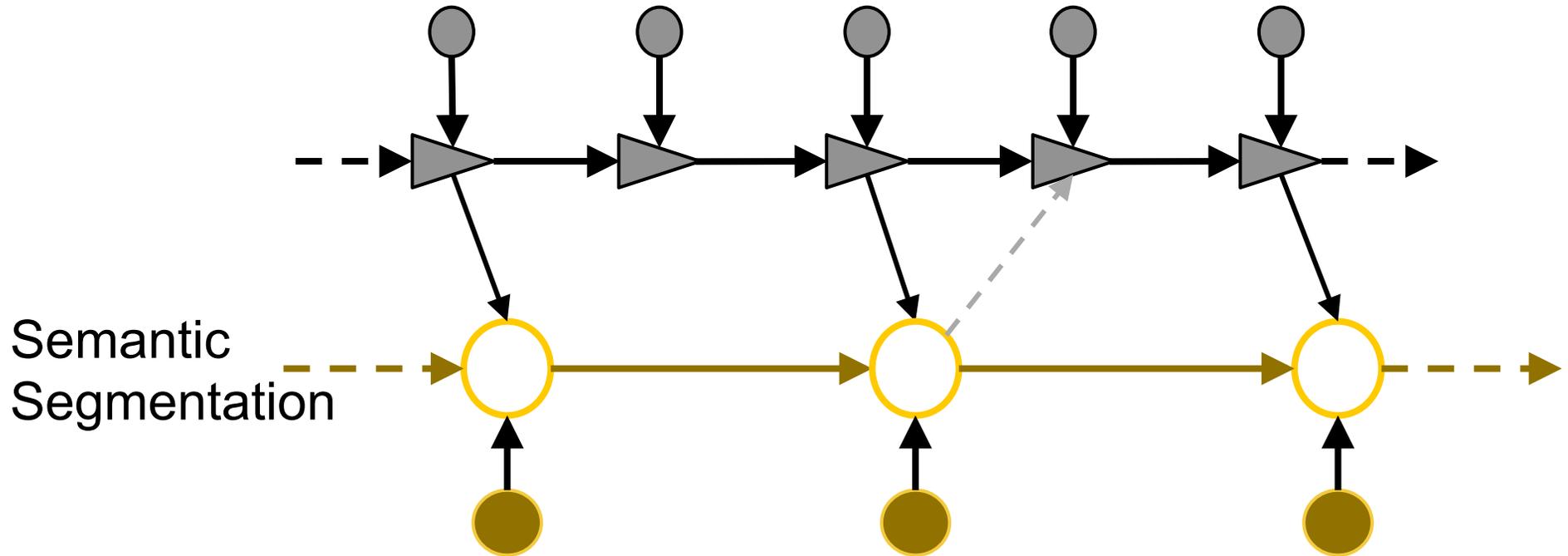


# Semantic Segmentation for Robotic Systems



Single view and Video

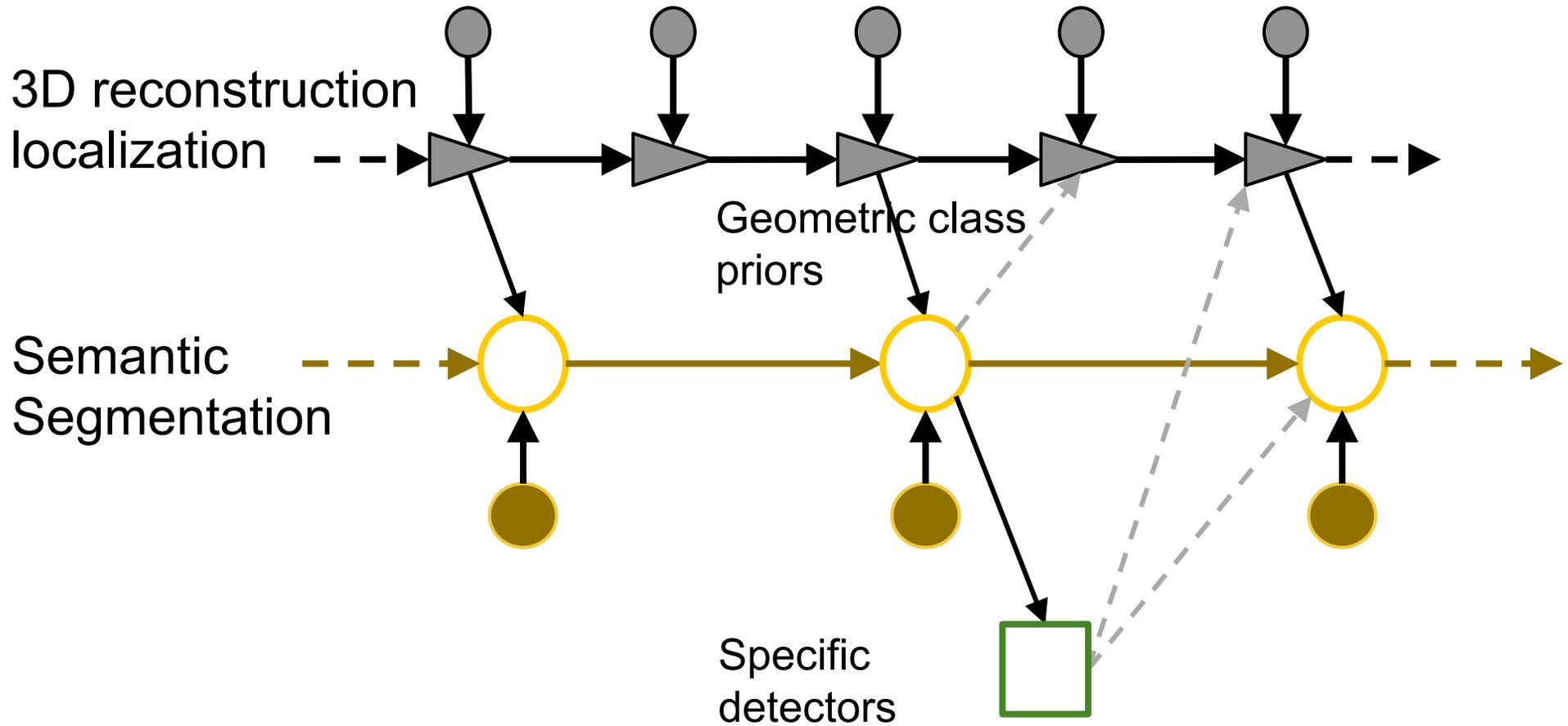
# Semantic Segmentation for Robotic Systems



Semantic  
Segmentation

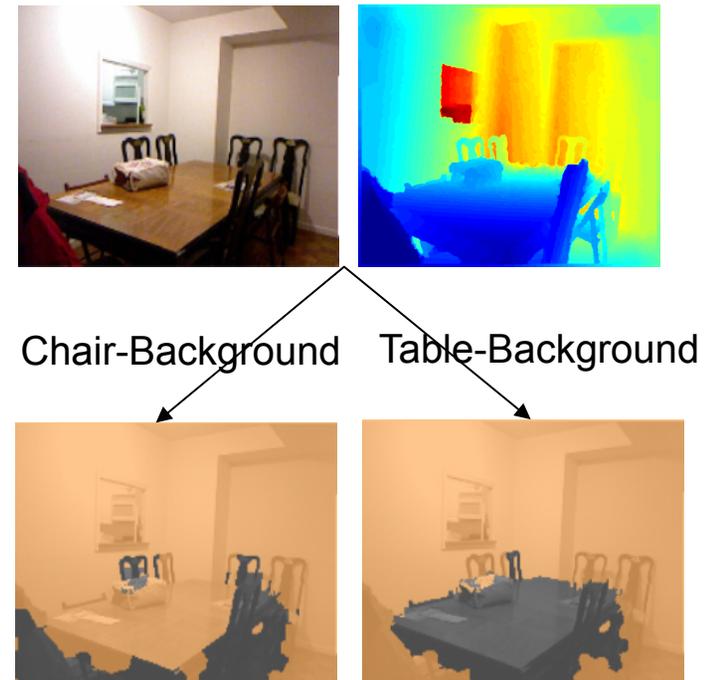
Feedback of priming  
classes

# Semantic Segmentation for Robotic Systems



# Semantic Refinement

- Refinement of the object category
- Binary **object-of-interest vs background** segmentation task
- Endow superpixels with richer features



- Learned one-vs-all AdaBoost per object.
- Equal proportion of positive-negatives ex.
- Negative mining to select samples from other objects that co-occur with the object of interest

# Object Level Segmentation Jaccard Index

	Bed	Sofa	Chair	Table	Window	Bookshelf	TV	Bag	Bathub	Blinds	Books	Box	Cabinet	Clothes	Counter	Curtain	Desks
Silberman <sup>[5]</sup>	40.00	25.00	32.00	21.00	30.00	23.00	5.70	0.00	0.00	40.00	5.50	0.13	33.00	6.50	33.00	27.00	4.60
Ren <sup>[4]</sup>	42.00	28.00	33.00	17.00	28.00	17.00	19.00	1.20	7.80	27.00	<b>15.00</b>	3.30	37.00	9.50	39.00	28.00	10.00
Gupta <sup>[2]</sup>	55.00	<b>44.00</b>	<b>40.00</b>	<b>30.00</b>	<b>33.00</b>	20.00	9.30	0.65	33.00	<b>44.00</b>	4.40	<b>4.80</b>	<b>48.00</b>	6.90	<b>47.00</b>	<b>34.00</b>	10.00
Ours(unary)	50.64	37.44	25.00	19.19	25.93	23.88	26.40	<b>3.28</b>	32.12	29.77	9.17	2.89	27.42	9.79	34.68	25.59	21.04
Ours(CRF)	<b>56.85</b>	42.29	31.44	20.78	30.16	<b>30.29</b>	<b>34.97</b>	3.00	32.95	33.09	10.06	3.99	29.34	10.04	33.82	30.11	<b>23.35</b>
	Door	Dresser	Floor-mat	Lamp	Mirror	Night-stand	Paper	Person	Picture	Pillow	Refridgerator	Shelves	Shower-curtain	Sink	Toilet	Towel	Whiteboard
Silberman <sup>[5]</sup>	5.90	13.00	7.20	<b>16.00</b>	4.40	6.30	<b>13.00</b>	6.60	36.00	19.00	1.40	3.30	3.60	25.00	27.00	0.11	0.00
Ren <sup>[4]</sup>	13.00	7.00	20.00	14.00	18.00	9.20	12.00	14.00	32.00	20.00	1.90	6.10	5.40	<b>29.00</b>	35.00	13.00	0.15
Gupta <sup>[2]</sup>	8.30	22.00	22.00	6.80	19.00	20.00	1.90	16.00	<b>40.00</b>	<b>28.00</b>	15.00	5.10	18.00	26.00	<b>50.00</b>	<b>14.00</b>	37.00
Ours(unary)	14.72	32.35	32.81	6.68	23.09	16.22	7.64	19.54	17.93	16.16	16.86	10.67	25.54	10.98	26.06	7.62	36.25
Ours(CRF)	<b>17.16</b>	<b>35.73</b>	<b>34.19</b>	12.14	<b>27.41</b>	<b>21.54</b>	10.07	<b>30.31</b>	22.21	22.98	<b>20.59</b>	<b>13.46</b>	<b>26.84</b>	11.04	38.65	8.61	<b>37.69</b>

## References

- S. Gupta, P. Arbelaez, and J. Malik. Perceptual organization and recognition of indoor scenes from RGB-D images. (CVPR) 2013.
- X. Ren, L. Bo, and D. Fox. RGB-(D) scene labeling: Features and algorithms. (CVPR), 2012.
- N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgb-d images. (ECCV), 2012.

Silber.[5]	15.12
RenFox[4]	17.99
Gupta [2]	23.92
Ours	<b>24.92</b>

Table: Summary of results in mean Jaccard Index metric.

# Conclusions

Computational efficient approach for semantic segmentation

We see it as the first stage of a scalable semantic understanding for mobile robots

Our approach effectively uses 3D and Images cues

Both 3D reconstruction and semantic segmentation formulated on the same graph induced by superpixels

We exploited the versatility and flexibility of CRFs to connect and use different sensory modalities for full coverage



Cesar Cadena, University of Adelaide

Branislav Micusik, currently at ARC Austria

Md. Alimoor Reza, GMU

Thank you !

Multi-view Superpixel Stereo in Man-Made Environments

B. Micusik and J. Kosecka

International Journal of Computer Vision, 2010, Vol 89, Number 1, 106-119.

Recursive Inference for Prediction of Objects in Urban Environments

Cadena C. and Kosecka J.

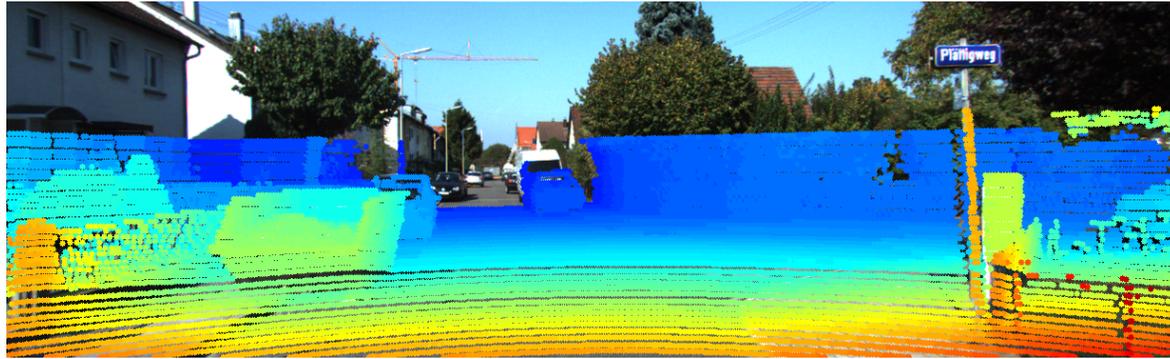
International Symposium on Robotics Research. December 2013, Singapore

Object Recognition and Segmentation in RGB-D scenes.

Md. Alimoor Reza and Kosecka J.

RGB-D workshop, RSS 2014, Berkeley

Supported by the US Army Research Office Grant W911NF-1110476



# Semantic Parsing in Indoors and Outdoors Environments

Jana Kosecka

George Mason University



Joint work with C. Cadena Lerma, G. Singh, Md. Alimoor Reza  
Supported by the Army Research Office Grant W911NF-1110476

# Semantic Parsing of Open Environments

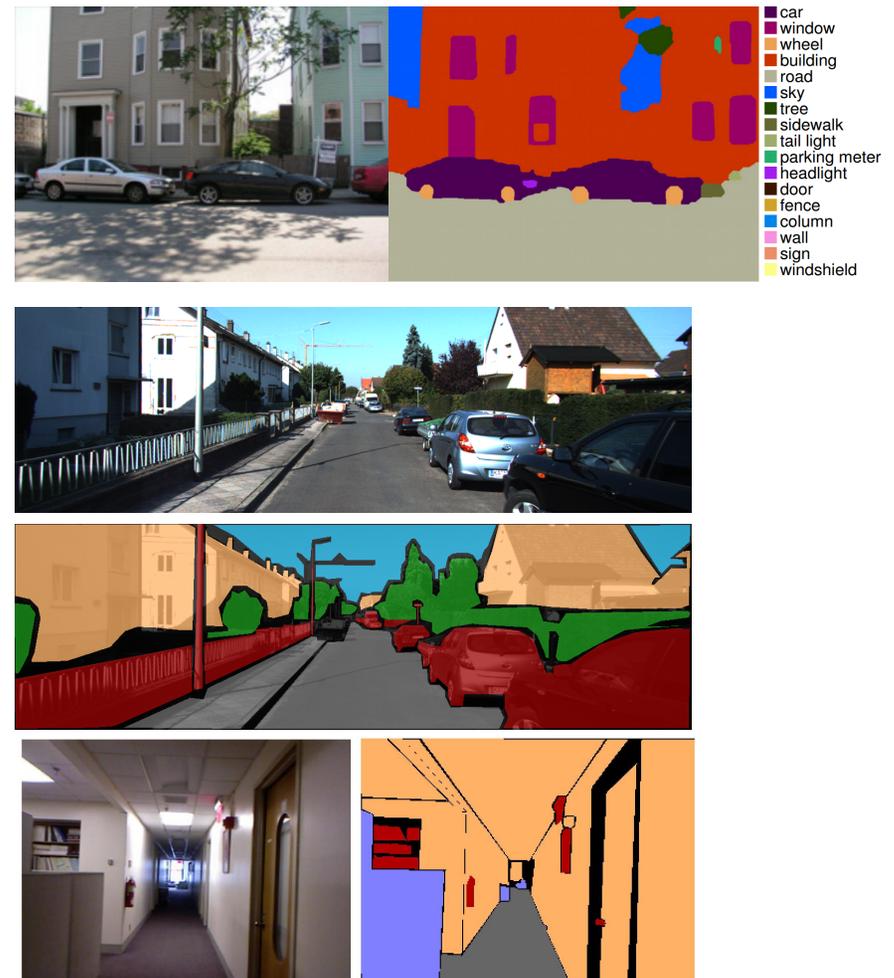
Simultaneous segmentation and categorization of partitions of sensory data into background and object categories



# Semantic Parsing of Open Environments

Endowing environment models with semantic information can enhance robustness and sophistication of robotic tasks

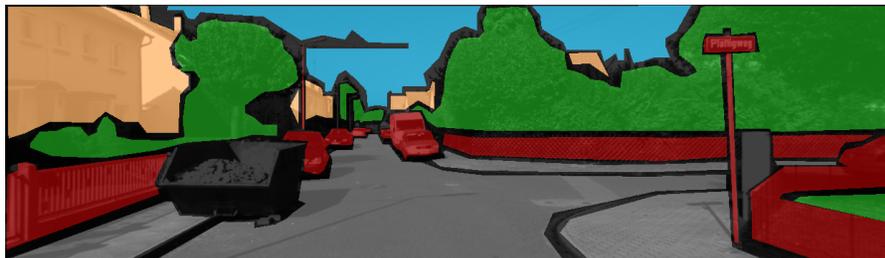
- In **Computer Vision**, often in single view setting, expensive preprocessing, learning and inference
- **Robotics** requires processing of video and multiple sensor modalities
- Semantic categories can be constrained by the tasks
- Incremental re-usable representations



# Semantic Segmentation

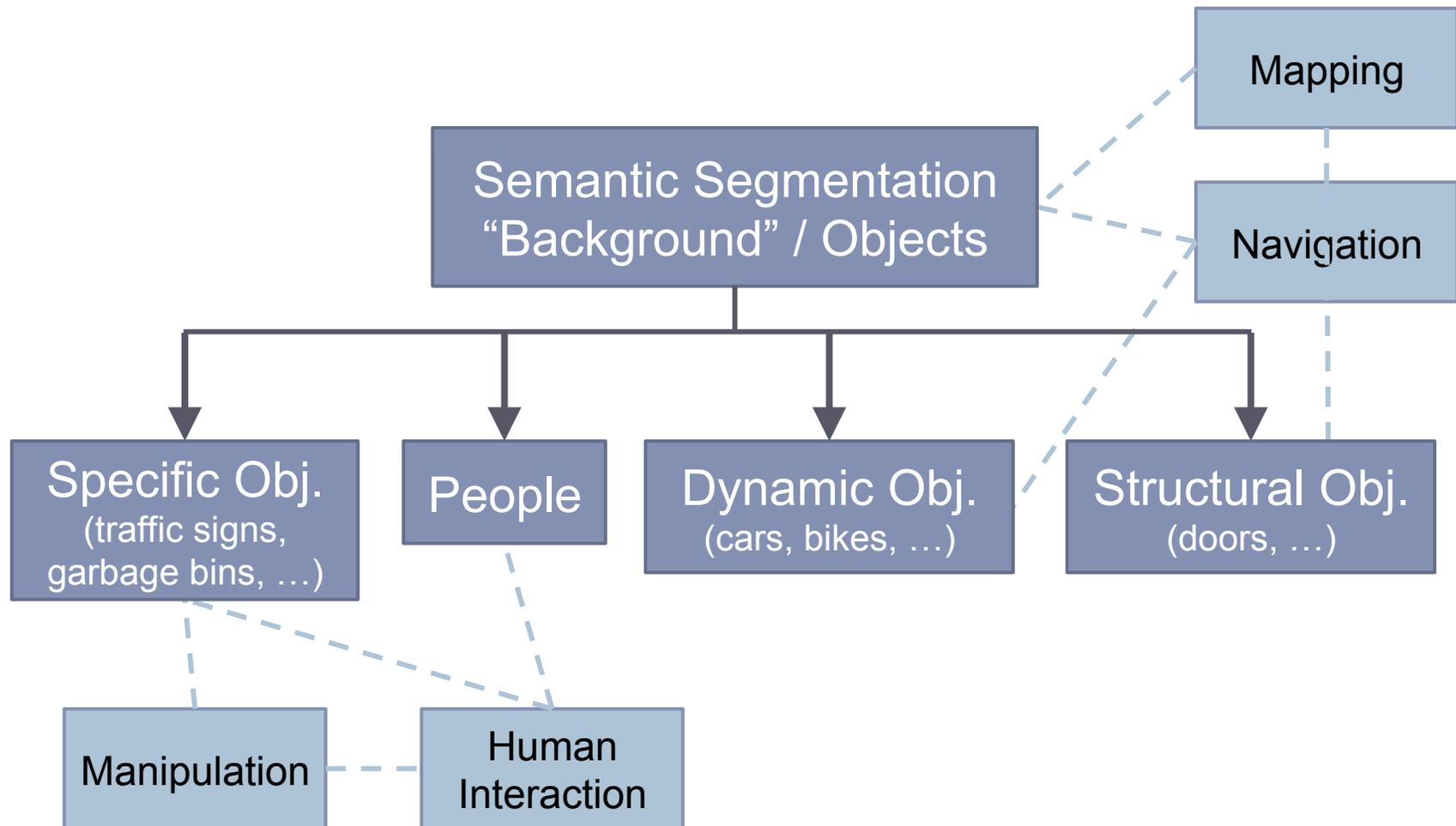
## Task constraints

- Navigation (road, free space)
- Localization (landmarks)
- Manipulation (object categorization, object search)
- Reasoning about static and dynamic object instances



Ground  
Truth

# Semantic Concepts & Tasks



# Semantic Hierarchy

## “Background” and object categories

- Non-object categories specific to types of scenes
- we can assume to be present (almost) always
- mostly static or slow changing
- Generic objects share some characteristics

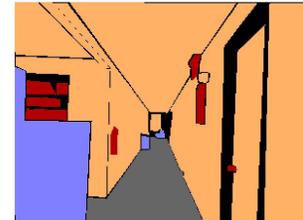
KITTI dataset, Geiger et al. IJRR 2013



Ground  
Buildings  
Vegetation  
Sky  
Objects

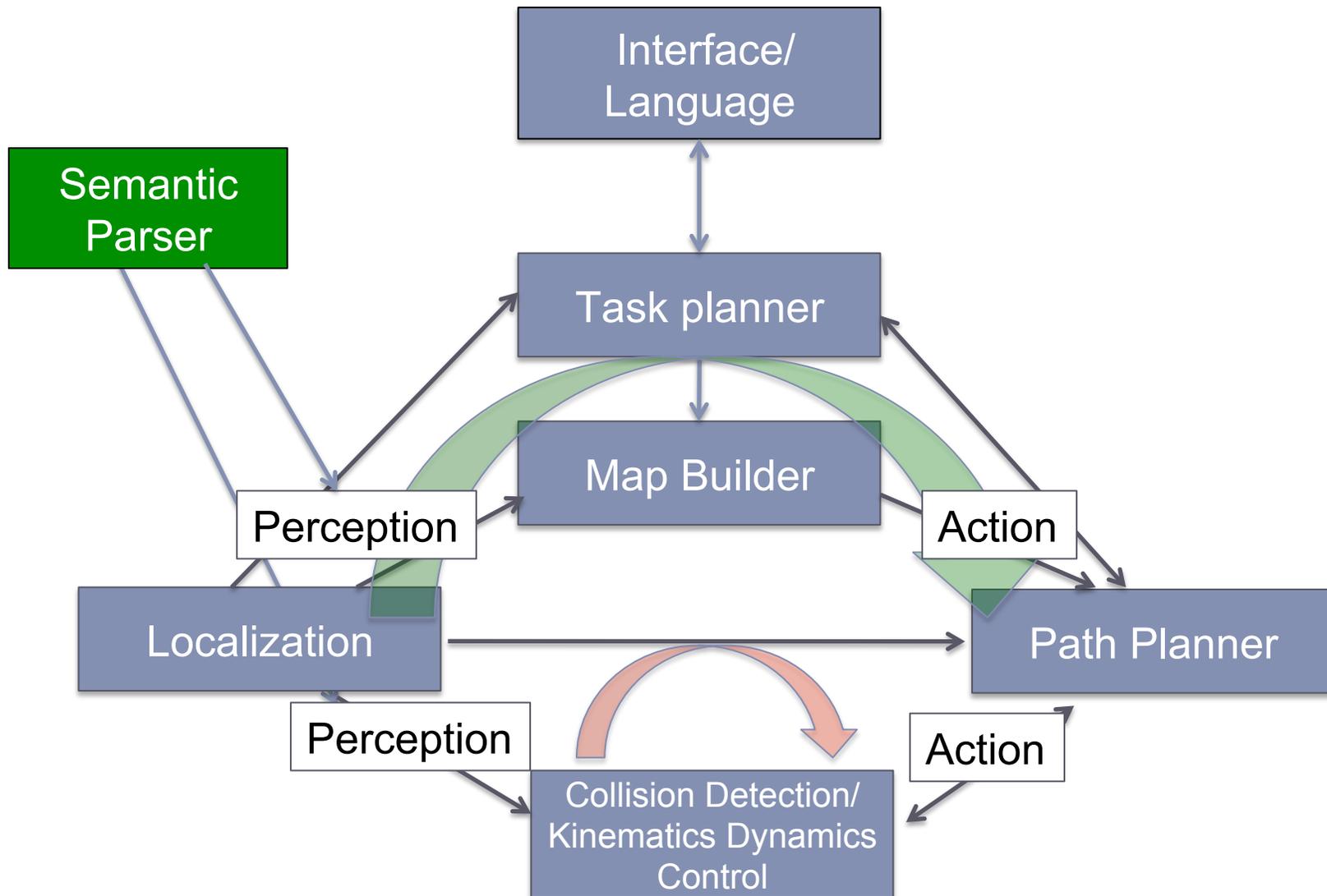
Mapped from Sengupta et al. ICRA 2013

NYU V2 dataset, Silberman et al. ECCV 2012

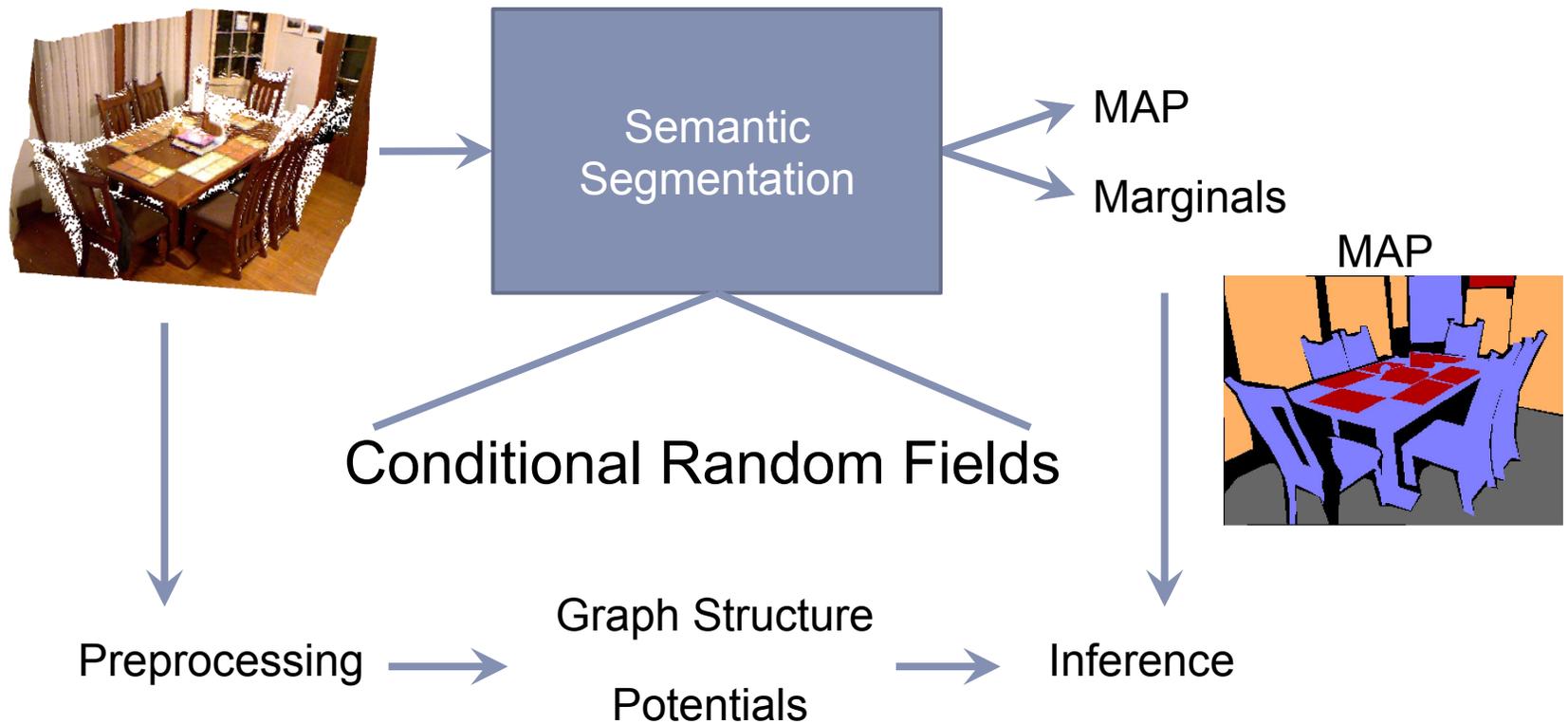


Ground  
Structure  
Furniture  
Props

# Different time scales



# Approach



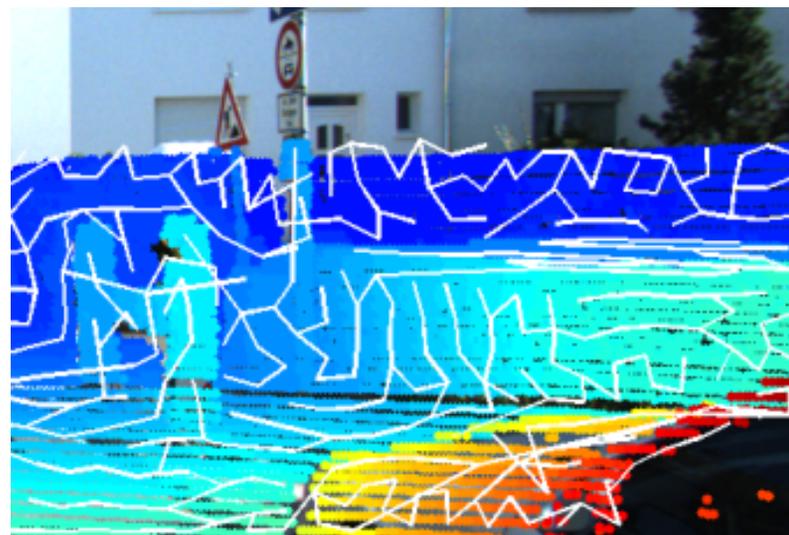
Previous methods suffer from :

→ Expensive over-segmentation, Expensive features Expensive Inference

Learning

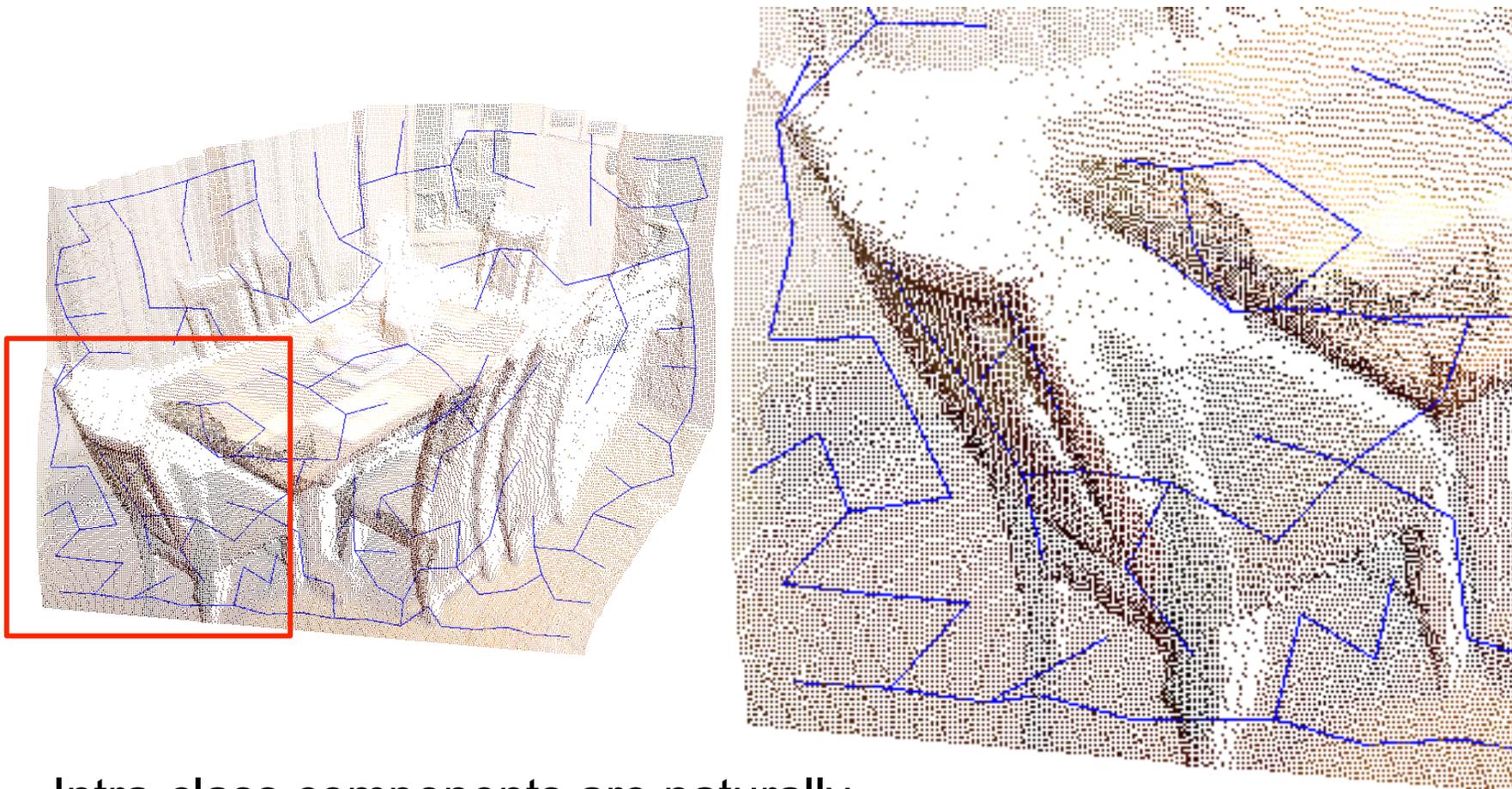
# Graph Structure:

## Minimum Spanning Tree Over 3D



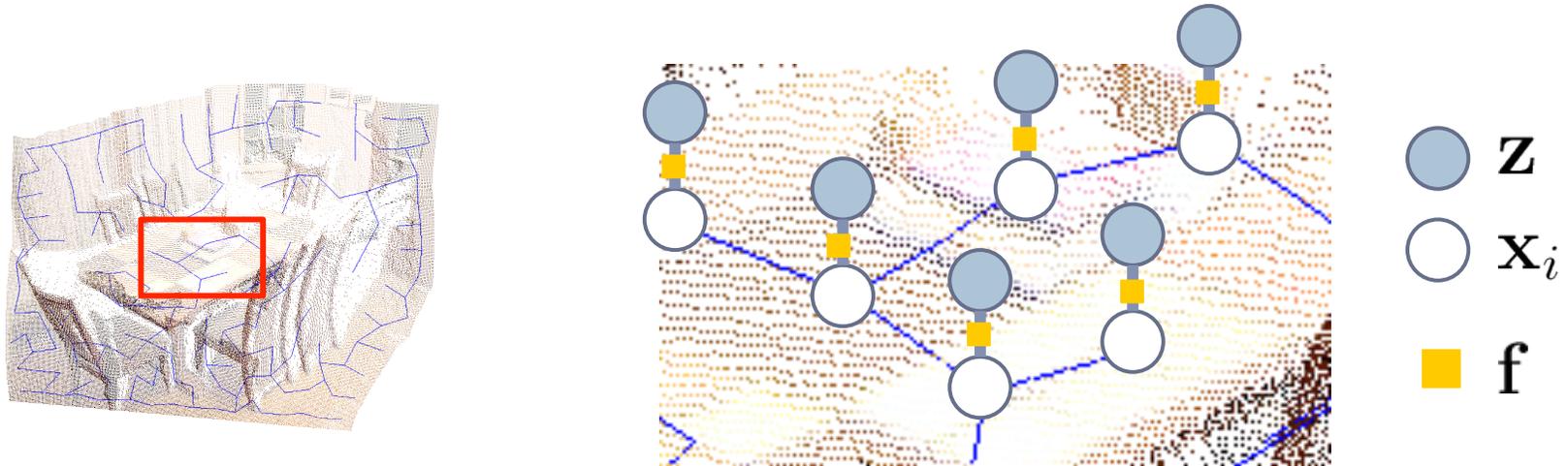
- SLIC superpixels, preserve contours, regularity, efficiency
- Edges are determined by the MST over 3D distances between superpixels' centroids

# Graph Structure: Our choice MST over 3D



Intra-class components are naturally connected

# Potentials: Pairwise CRFs



$$p(\mathbf{x}|\mathbf{z}) = \frac{1}{Z(\mathbf{z})} \exp \left( \sum_{i \in \mathcal{N}} \mathbf{w}_u^T \mathbf{f}(\mathbf{x}_i, \mathbf{z}) + \sum_{i, j \in \mathcal{E}} \mathbf{w}_p^T \mathbf{g}(\mathbf{x}_i, \mathbf{x}_j, \mathbf{z}) \right)$$

MAP inference in CRF:

compute most likely labels  $\mathbf{x}$  given observations  $\mathbf{z}$

# Potentials: unary & pairwise

unary (local) potential  
using a k-NN classifier

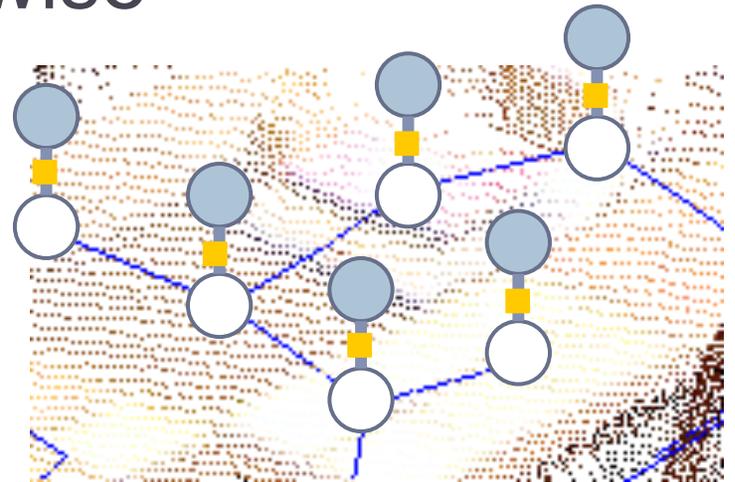
$$\blacksquare f(\mathbf{x}_i, \mathbf{z}) = -\log P_i(\mathbf{x}_i | \mathbf{z})$$

$$P_i(\mathbf{x}_i = l_j | \mathbf{z}) = \eta \frac{f(l_j) \bar{F}(l_j)}{\bar{f}(l_j) F(l_j)}$$

$f(l_j)$  frequency of label  $j$  in a k-NN query

$F(l_j)$  frequency of label  $j$  the database

$$\blacksquare g(\mathbf{x}_{i,j}, \mathbf{z}) = \begin{cases} 1 - \exp(-\|c_i - c_j\|_2) & \rightarrow l_i = l_j \\ \exp(-\|c_i - c_j\|_2) & \rightarrow l_i \neq l_j \end{cases}$$

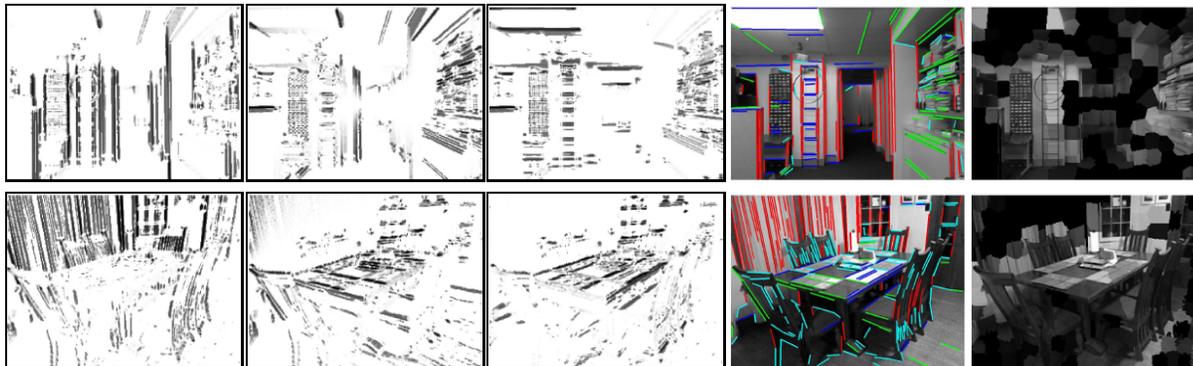
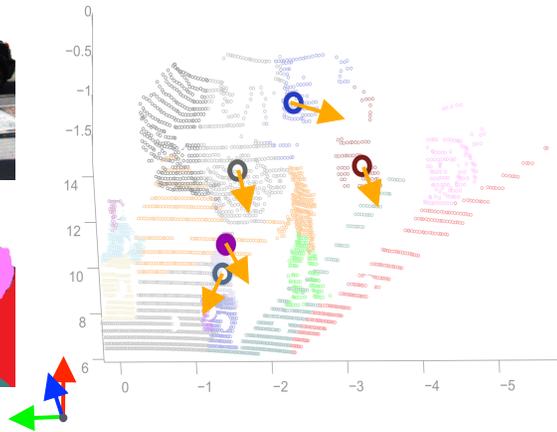


# Features

indoors (15D)

outdoors(6D)

- From 3D
  - mean and std of differences on depth
  - local planarity
  - neighboring planarity
  - vertical orientation
- From Image:
  - entropy from vanishing points



$$H_s = - \sum_{j=1}^4 h_{g_s}(y = j) \log(h_{g_s}(y = j))$$

# Inference

- We use belief propagation:
  - Exact results in MAP/marginals
  - Efficient computation, in  $\mathcal{O}(nm^2)$

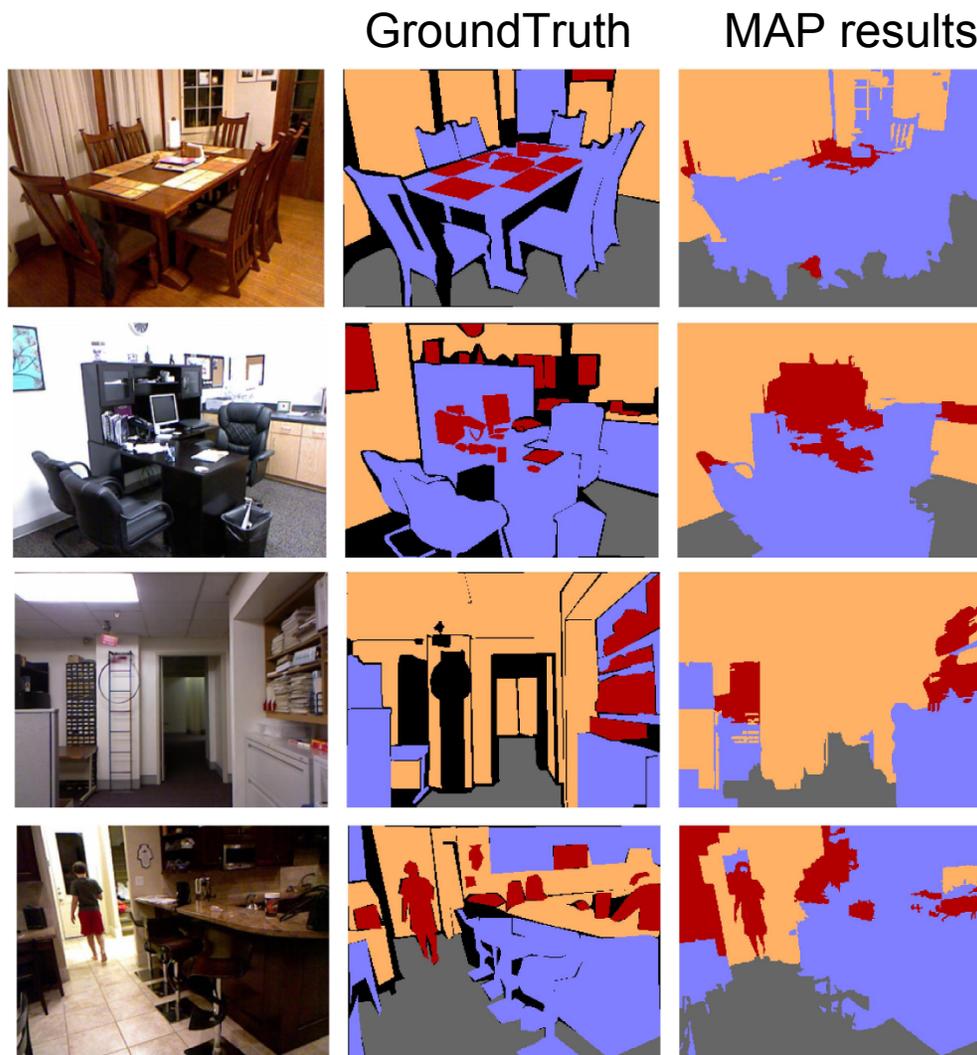
# Learning

- Maximum Likelihood Estimation
  - To learn  $[\mathbf{w}_u, \mathbf{w}_p]$

Tree graph structure:  
good convergence

# Results: NYU-Depth v2 Dataset

Qualitative comparison:



Ground  
Structure  
Furniture  
Props

# Results: NYU-Depth v2 Dataset

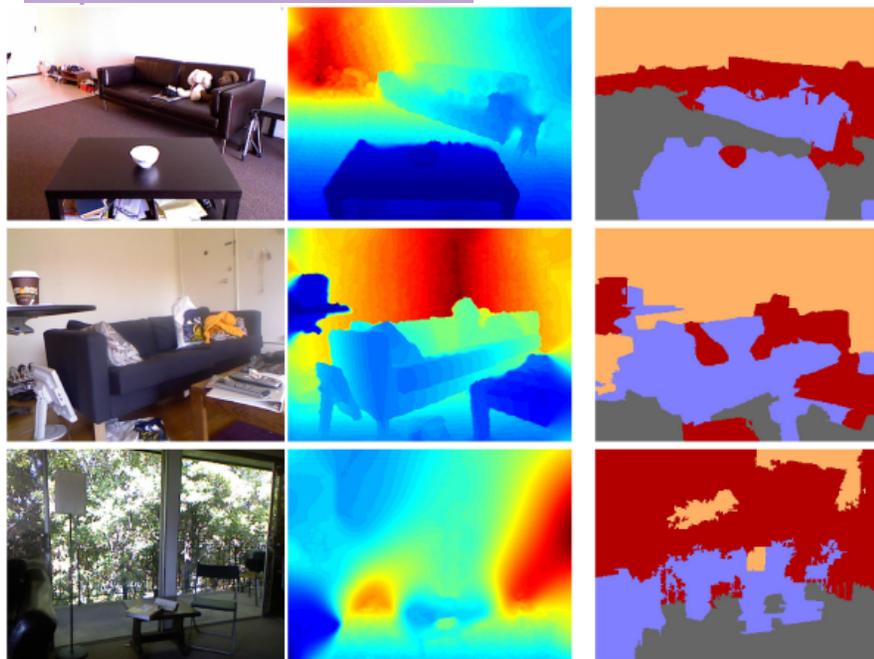
Quantitative comparison:

- Recall accuracy in pixel-wise percentage:

	Ground	Furniture	Props	Structure	Average	Global
CRF-MST-kNN	88.4	64.1	30.5	78.6	<b>65.4</b>	<b>67.2</b>
only Image Feat.	63.2	47.5	24.5	73.6	52.2	56.1
only 3D Feat.	<b>89.5</b>	<b>70.0</b>	16.9	79.4	62.7	65.8
data-term (kNN)	87.3	60.6	33.7	74.8	64.1	64.9
Silberman et al. 2012	68	<b>70</b>	<b>42</b>	59	59.6	58.6
Couprie et al. 2013	87.3	45.3	35.5	<b>86.1</b>	63.5	64.5

# Results: Independent datasets

B3DO: Berkeley 3-D Object Dataset  
Janoch and Karayev, 2011  
<http://kinectdata.com/>

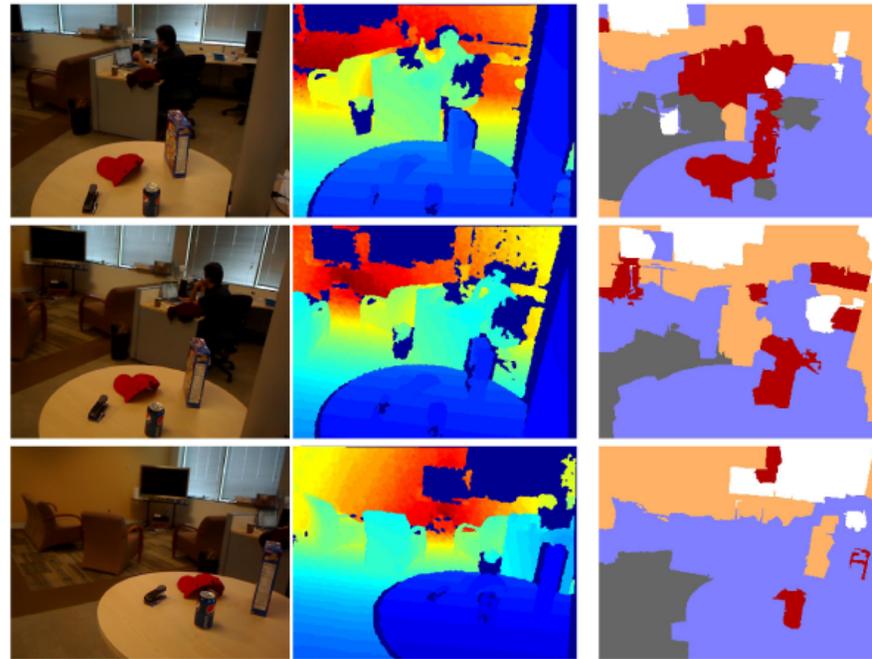


RGB  
results

Depth

MAP

RGB-D Object Dataset  
K. Lai, L. Bo, X. Ren and D. Fox, 2011  
<http://www.cs.washington.edu/rgbd-dataset/>



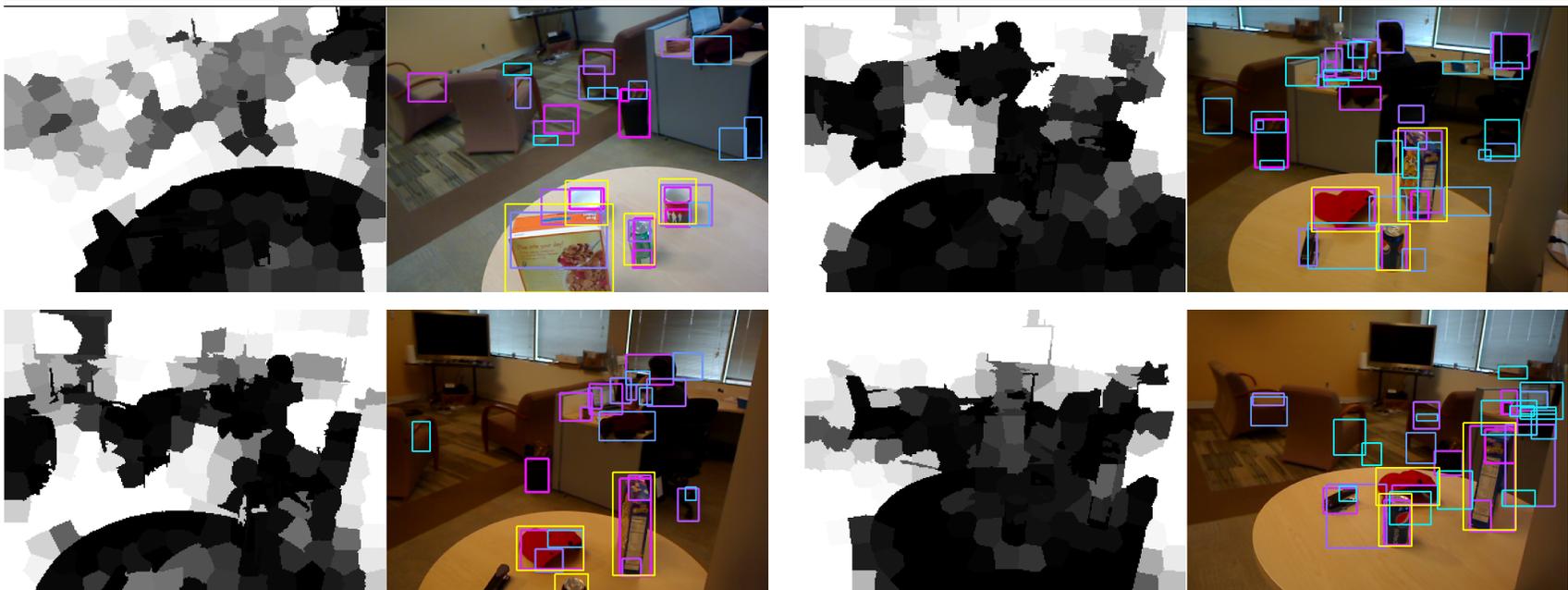
RGB  
results

Depth

MAP

# Generating Object Hypotheses

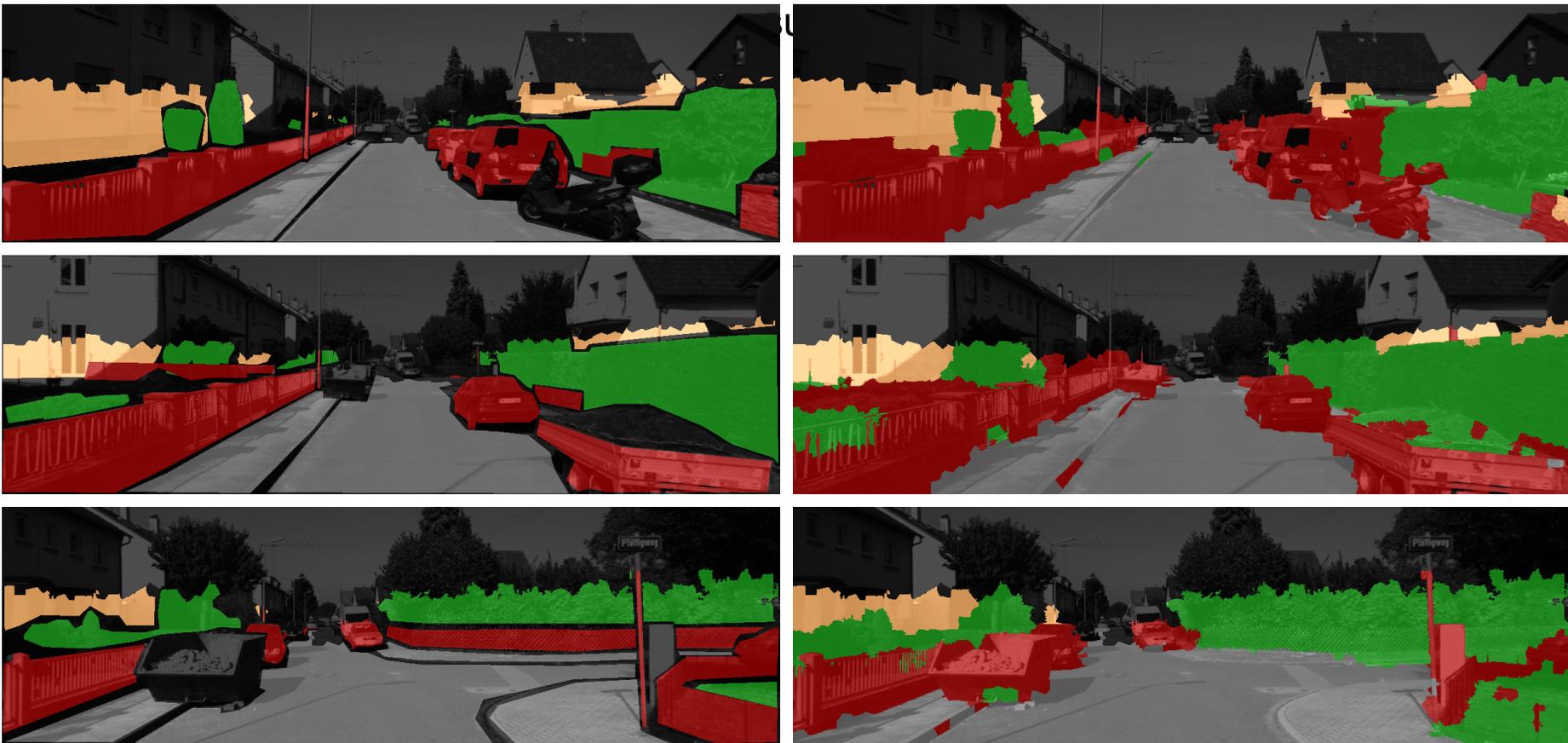
- Given the knowledge of the size of object to be manipulated, priming object detection
- Generate object hypotheses using the prob. map



# Results: KITTI Dataset

Ground Truth

MAP



Ground

Buildings

Vegetation

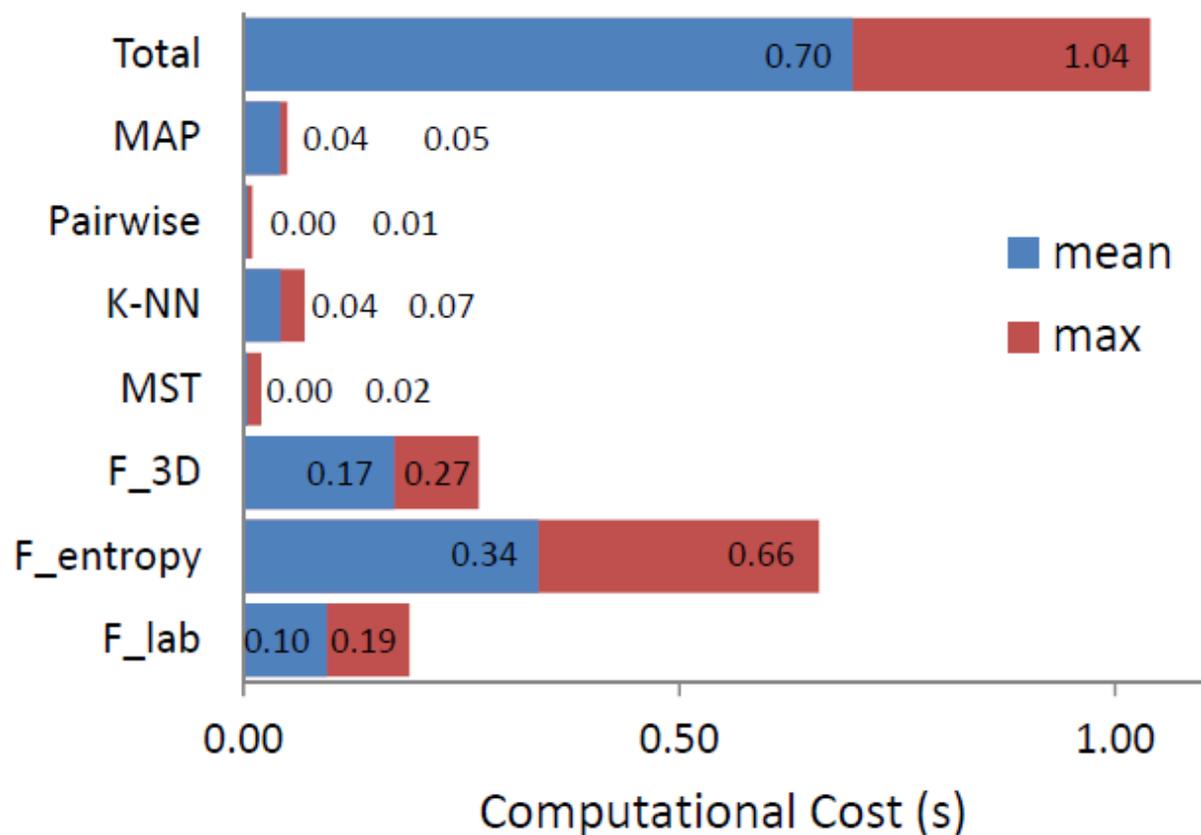
Objects

# Results: KITTI Dataset

- Recall accuracy in pixel-wise percentage:

	Ground	Objects	Building	Vegetation	Average	Global
CRF-MST-kNN	<b>97.3</b>	82.9	<b>82.8</b>	86.9	<b>87.5</b>	<b>88.4</b>
only Image Feat.	96.8	49.2	64.6	<b>95.5</b>	76.5	76.8
only 3D Feat.	95.9	<b>84.2</b>	80.5	46.7	76.8	78.8
data-term (kNN)	96.8	75.9	80.7	77.6	82.8	83.5

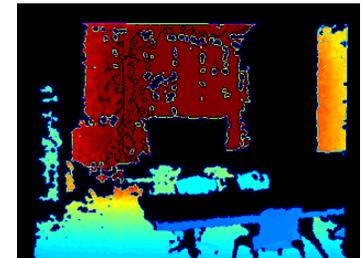
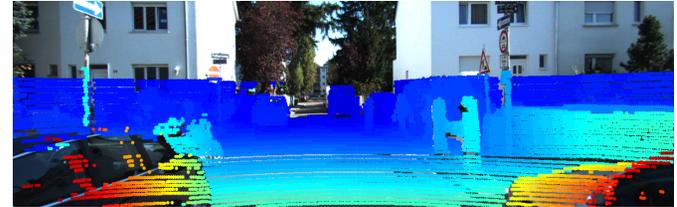
# Results: NYU-Depth v2 Dataset



Initial implementation in C++ with SLIC in GPU 5fps

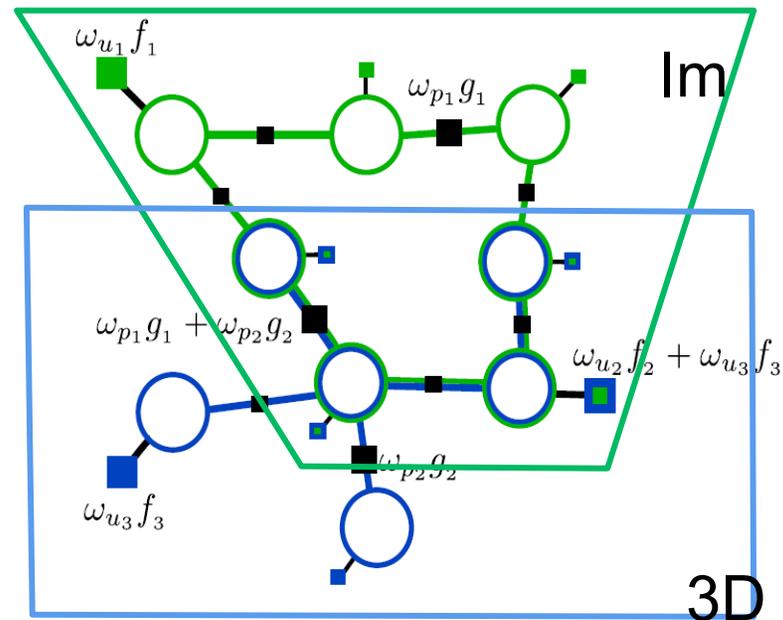
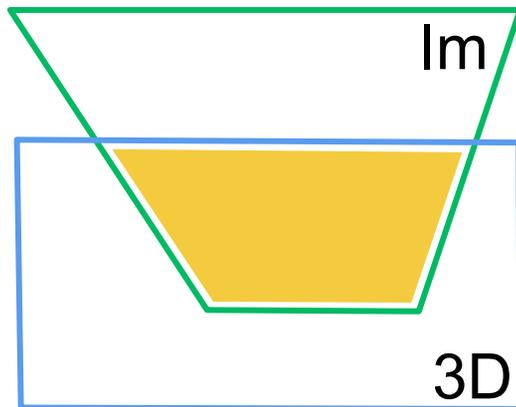
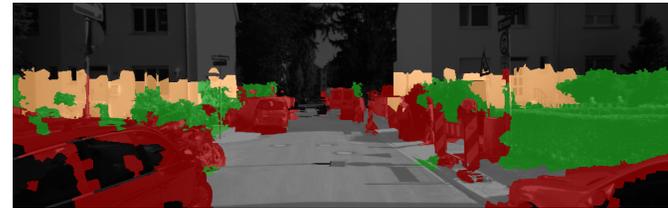
# Multiple Sensor Modalities

- Different fields of view
- Every sensor suffers from specific blind spots, e.g.
  - Laser: limited range, specular surfaces
  - Vision: low light conditions
  - Depth (Kinect): natural light, specularities
- Every modality suffers from different sources of ambiguities



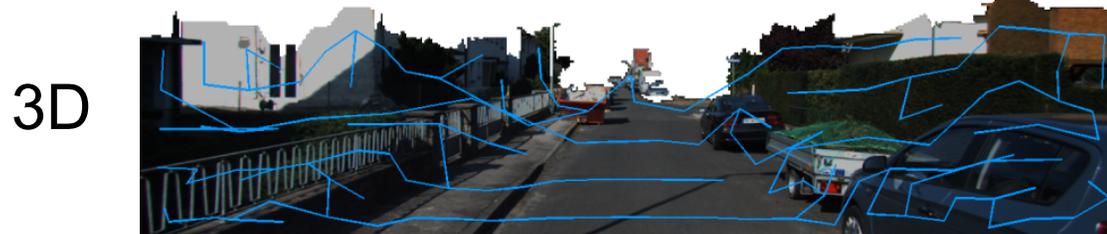
# Multiple Sensor Modalities

- Previously ‘semantic’ ambiguities fusing image and 3D sensor
- Common spatial coverage
- Without handling missing data
- Extension to union of FOV’s

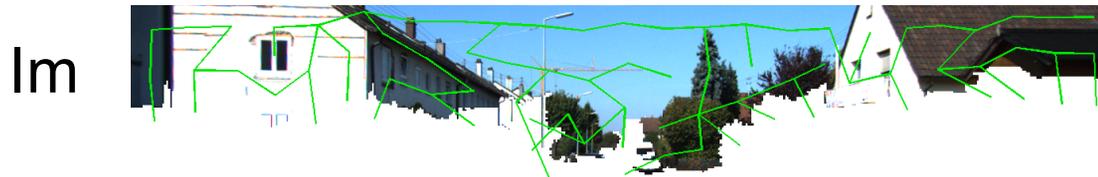


# Graph Structure

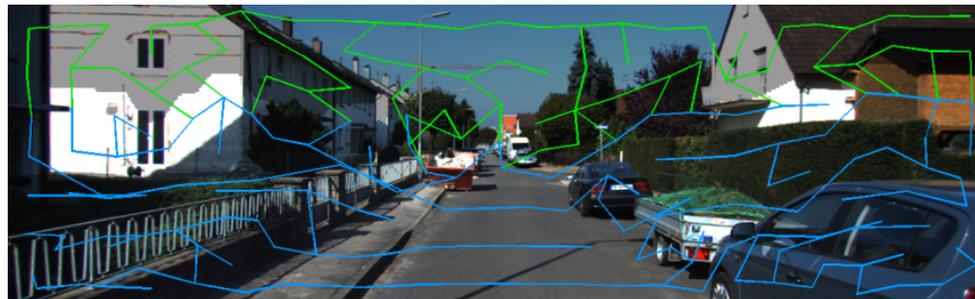
- Sub-graph over 3D



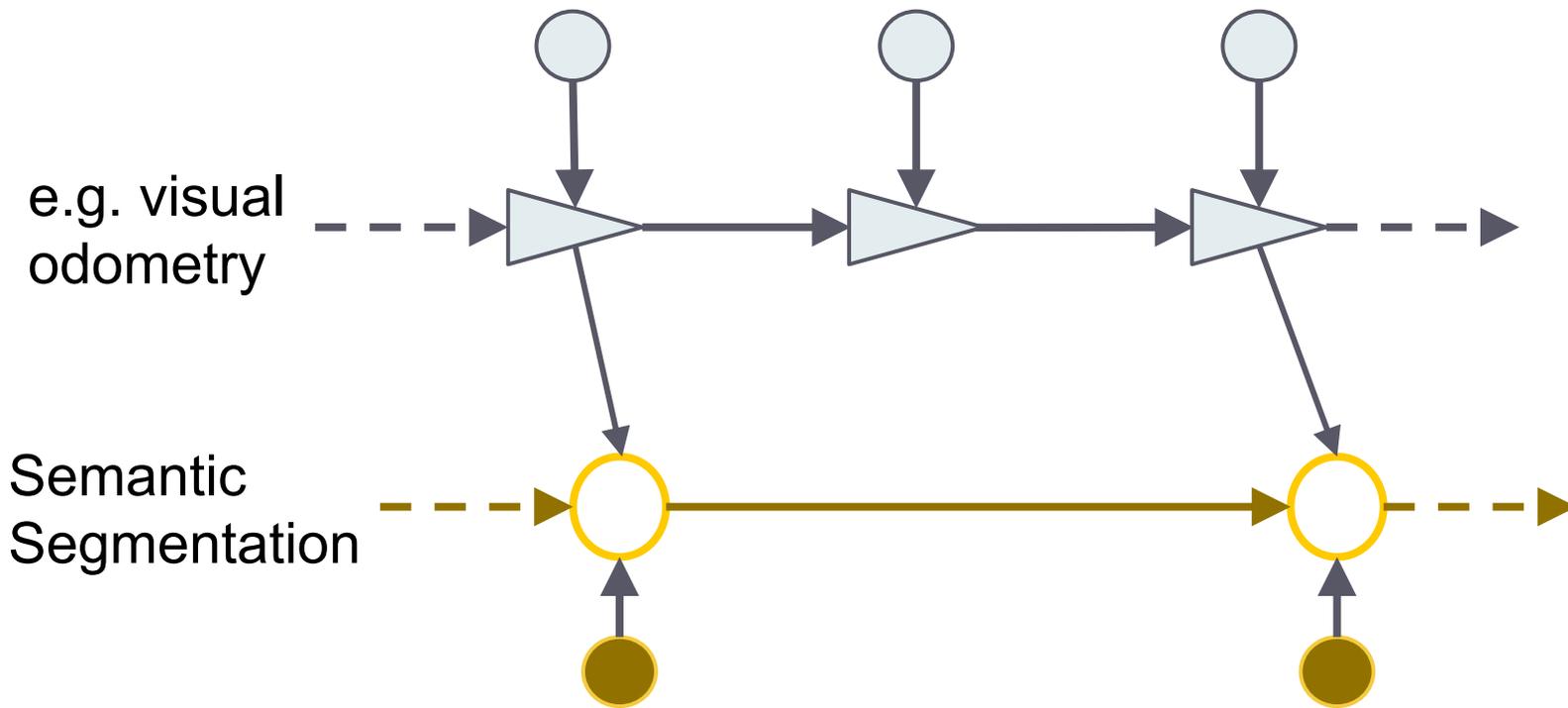
- Sub-graph over Image



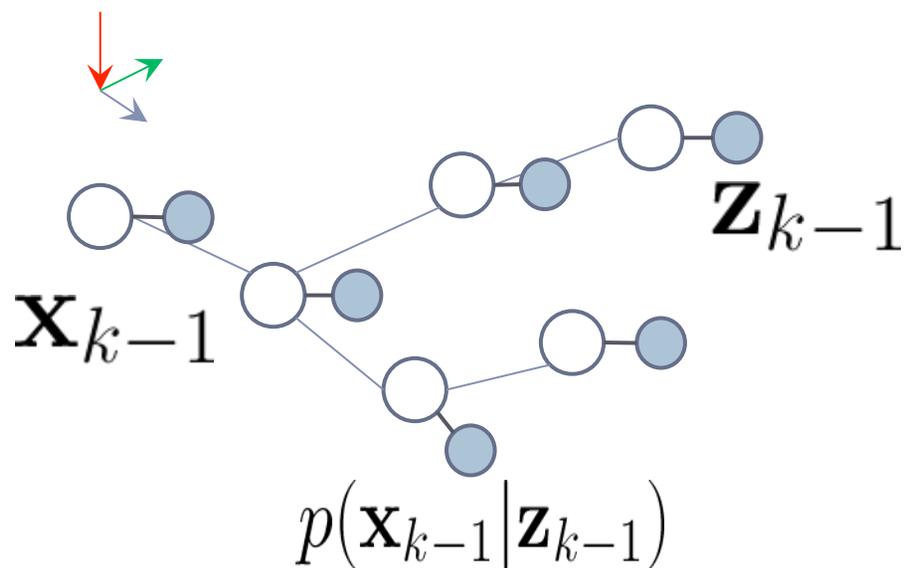
- Results in a graph for full coverage



# Recursive Inference



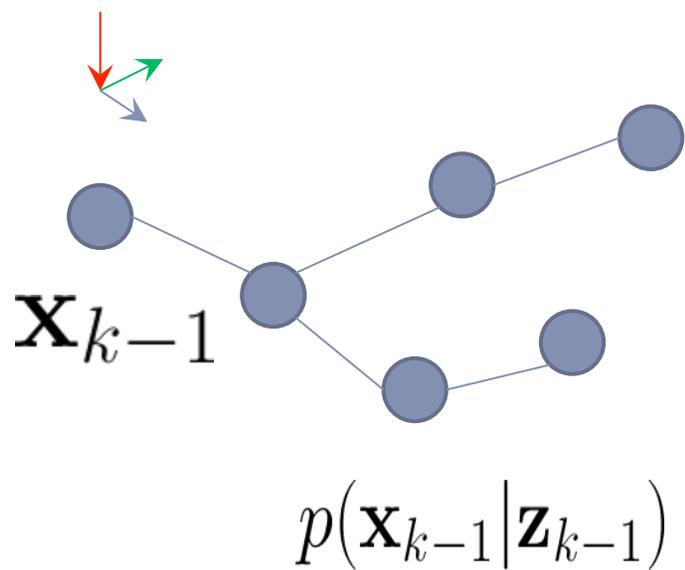
# Video sequences:



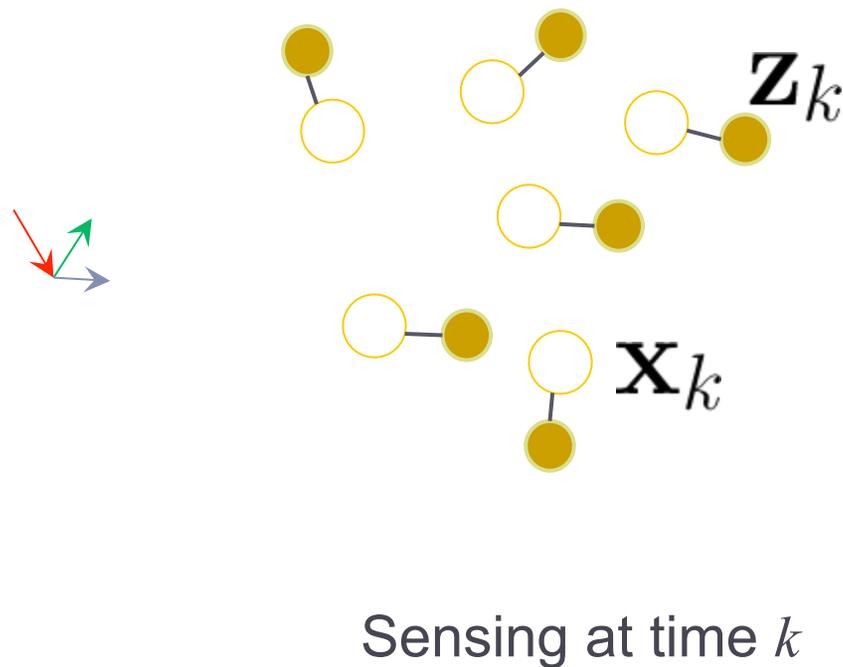
On-line operation

Infer the marginals at  
 $k-1$

# Recursive Inference

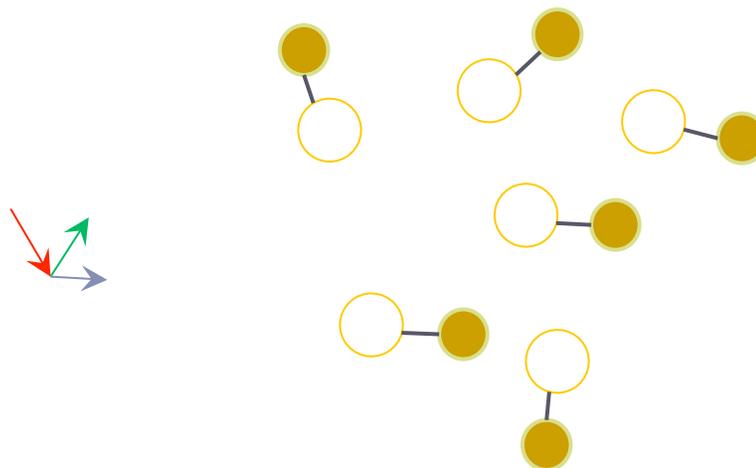
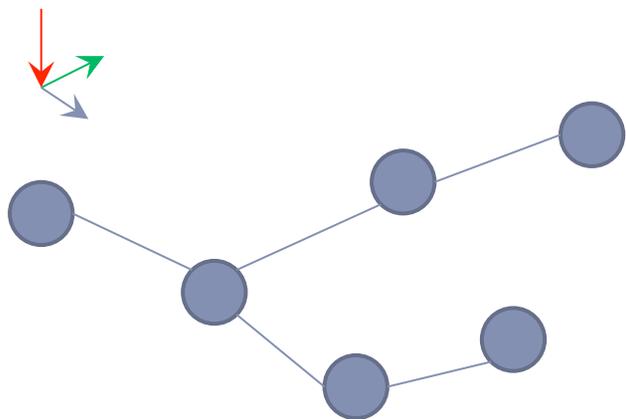


Marginals at  $k-1$



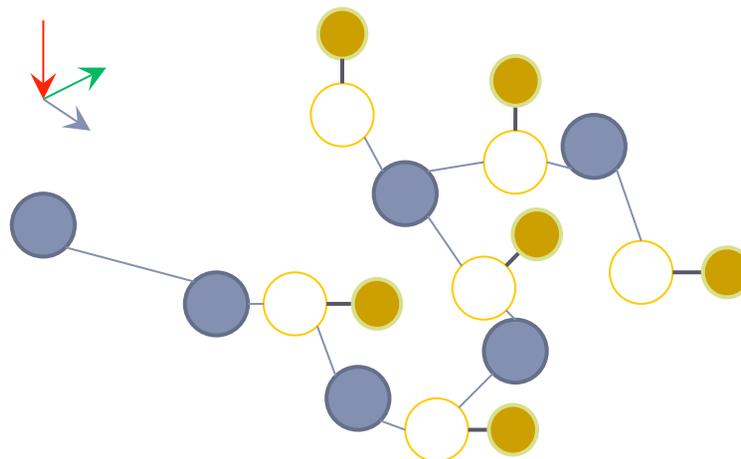
Sensing at time  $k$

# Recursive Inference



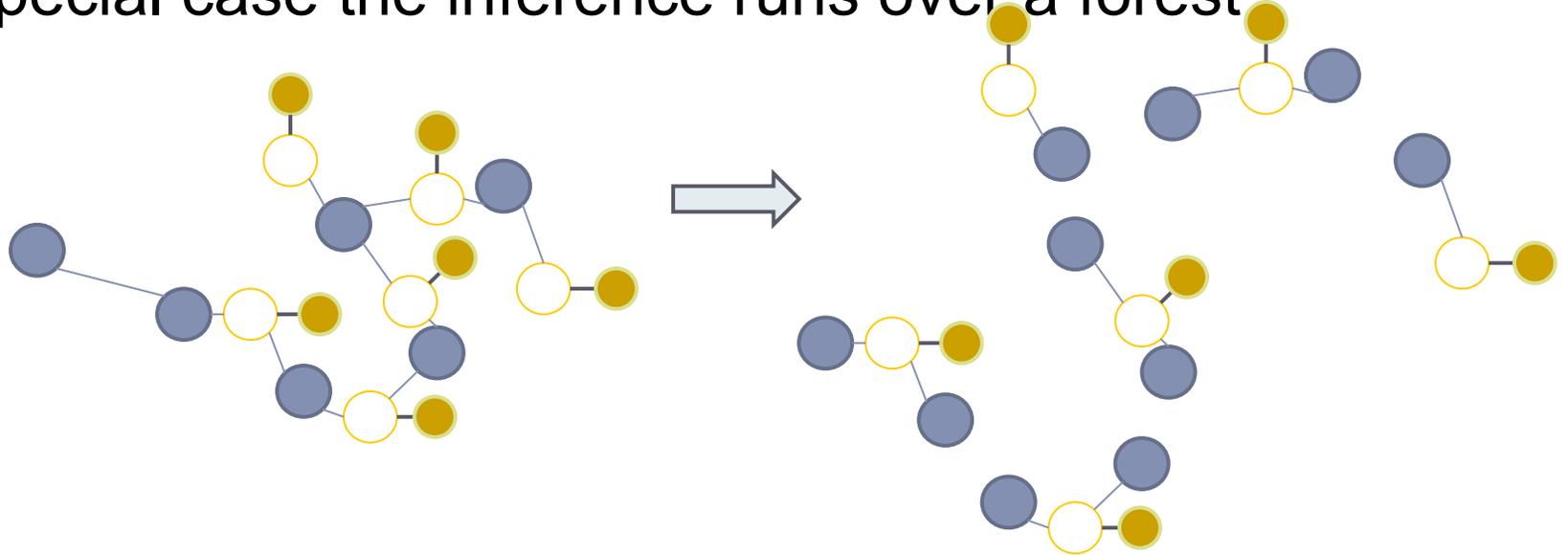
Infer marginals at time  $k$

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{z}_k, T_k)$$



# Recursive Inference

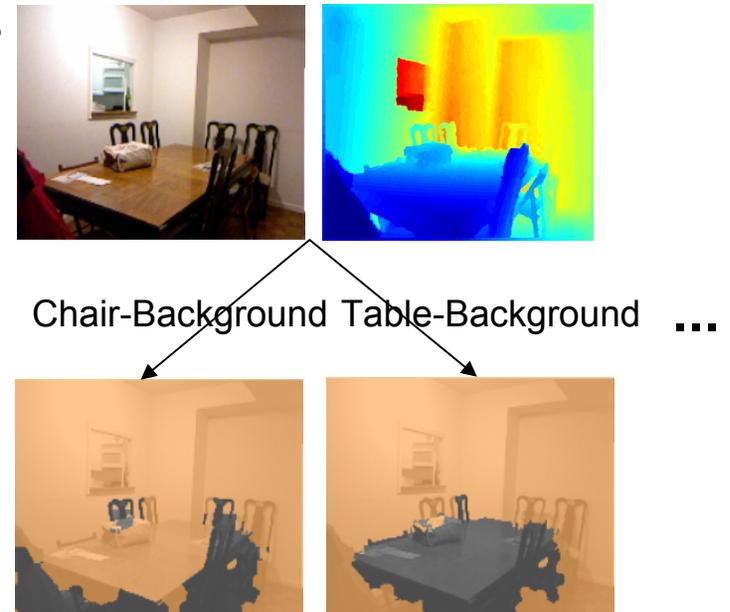
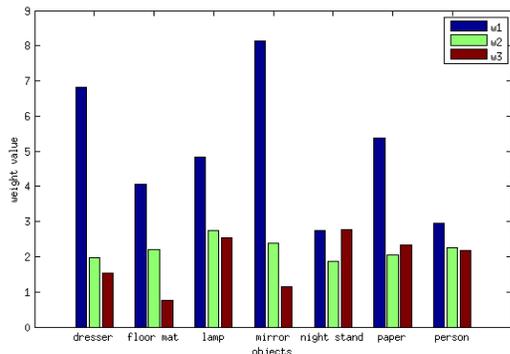
- Special case the inference runs over a forest



F1 KITTI Dataset	Ground	Objects	Building	Vegetation	Time MST	Time BP
Single View	0.977	0.854	0.870	0.811	21 ms	164 ms
Recursive Inference	0.977	0.853	0.879	0.809	57 ms	69 ms

# Finer Grained Categorization

- We formulate the problem of recognition and segmentation of objects in indoor scenes as a binary **object-of-interest vs background** segmentation task. Learn per category binary segmentations
- **Our choices:**
  - Enrich set of features
  - Low level per category grouping rules are learned in CRF setting

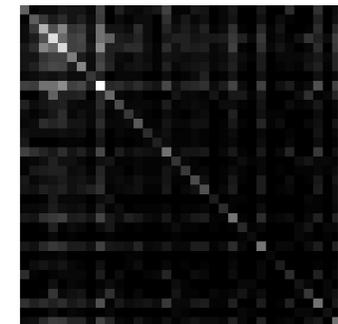


# Fine grained categorization

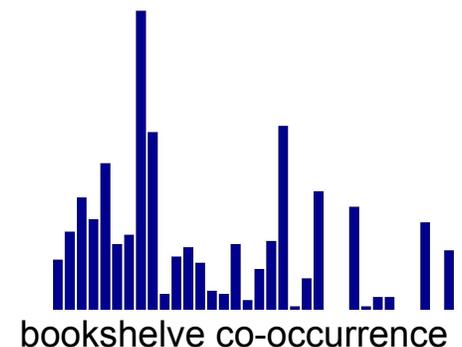
- Extend the set of features:
  - Color, Texture histograms Histograms C1, T1
  - Geometric Features (previous) G1
  - Generic Features G2 (planarity features, alignment with respect to gravity, orientation context)

Descriptor type	description
C1	75 dim. histogram HSV values
T1	240 dim. histogram
G1	11 dim. descriptor of geometric features
G2	60 dim. descriptor of generic features

- Adaboost classifier with Decision trees
- Exploit hard negative mining using context
- Sampling negative example proportional to
- Co-occurrence



object co-occurrence



bookshelve co-occurrence

# Finer grained categorization

- Object recognition and segmentation
- Train per class object-background models CRF's
- Evaluation in terms of per class segmentation accuracy, using Jaccard Index

$$JI = \frac{P \cap G}{P \cup G}$$

	Mean JI		Bed	Sofa	Chair	Table	Window	Books	TV
Silberman[15]	15.12								
Ren[13]	17.99								
Gupta[5]	23.92	Coupric[4]	38.4	24.6	34.1	10.2	15.9	13.7	6.0
		Hermans[6]	68.4	28.5	41.9	27.1	46.1	45.4	38.4
Ours	<b>24.99</b>	Ours	<b>87.8</b>	<b>86.5</b>	<b>83.1</b>	<b>78.3</b>	<b>78.5</b>	<b>73.8</b>	<b>82.4</b>

- Improves state of the art performance, very efficient (computational bottleneck feature computation)

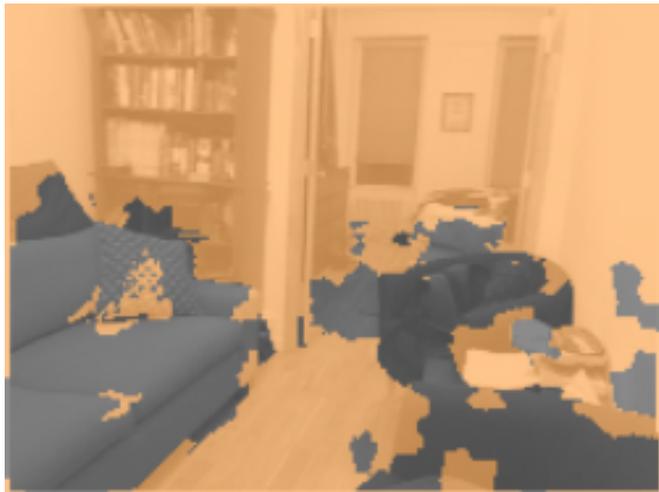
# More detailed per category results

	Bed	Sofa	Chair	Table	Window	Bookshelf	TV	Bag	Bathub	Blinds	Books	Box	Cabinet	Clothes	Counter	Curtain	Desks
Silberman[15]	40.00	25.00	32.00	21.00	30.00	23.00	5.70	0.00	0.00	40.00	5.50	0.13	33.00	6.50	33.00	27.00	4.60
Ren[13]	42.00	28.00	33.00	17.00	28.00	17.00	19.00	1.20	7.80	27.00	<b>15.00</b>	3.30	37.00	9.50	39.00	28.00	10.00
Gupta[5]	55.00	<b>44.00</b>	<b>40.00</b>	<b>30.00</b>	<b>33.00</b>	20.00	9.30	0.65	33.00	<b>44.00</b>	4.40	<b>4.80</b>	<b>48.00</b>	6.90	<b>47.00</b>	<b>34.00</b>	10.00
Ours(unary)	50.64	37.44	25.00	19.19	25.93	23.88	26.40	<b>3.28</b>	32.12	29.77	9.17	2.89	27.42	9.79	34.68	25.59	21.04
Ours(CRF)	<b>56.03</b>	42.51	30.59	20.89	30.21	<b>29.98</b>	<b>34.94</b>	3.13	<b>33.94</b>	34.77	11.01	3.70	29.56	<b>10.68</b>	34.02	30.28	<b>23.02</b>
	Door	Dresser	Floor-mat	Lamp	Mirror	Night-stand	Paper	Person	Picture	Pillow	Refridgerator	Shelves	Shower-curtain	Sink	Toilet	Towel	Whiteboard
Silberman[15]	5.90	13.00	7.20	<b>16.00</b>	4.40	6.30	<b>13.00</b>	6.60	36.00	19.00	1.40	3.30	3.60	25.00	27.00	0.11	0.00
Ren[13]	13.00	7.00	20.00	14.00	18.00	9.20	12.00	14.00	32.00	20.00	1.90	6.10	5.40	<b>29.00</b>	35.00	13.00	0.15
Gupta[5]	8.30	22.00	22.00	6.80	19.00	20.00	1.90	16.00	<b>40.00</b>	<b>28.00</b>	15.00	5.10	18.00	26.00	<b>50.00</b>	<b>14.00</b>	37.00
Ours(unary)	14.72	32.35	32.81	6.68	23.09	16.22	7.64	19.54	17.93	16.16	16.86	10.67	25.54	10.98	26.06	7.62	36.25
Ours(CRF)	<b>17.18</b>	<b>35.80</b>	<b>34.02</b>	11.17	<b>26.66</b>	<b>20.65</b>	10.29	<b>29.60</b>	21.91	22.00	<b>21.84</b>	<b>13.26</b>	<b>27.49</b>	12.11	39.37	9.71	<b>37.29</b>

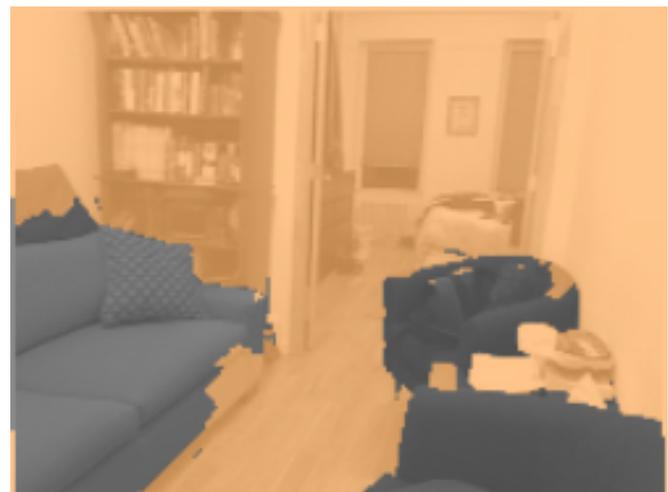
# Results



Ground Truth



Data Term



CRF

# Results



Ground Truth

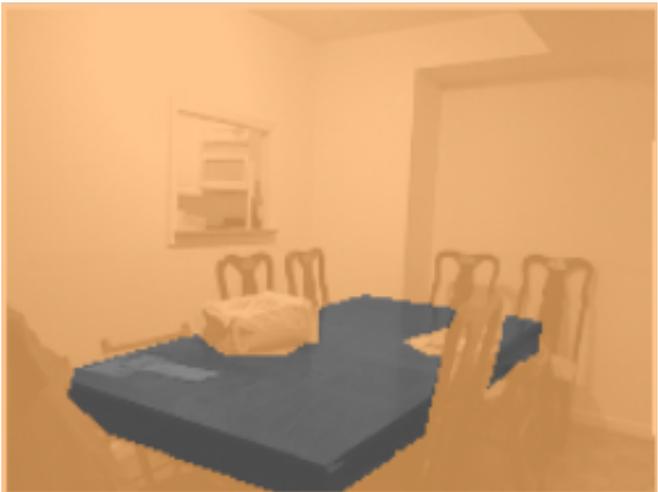


Data Term

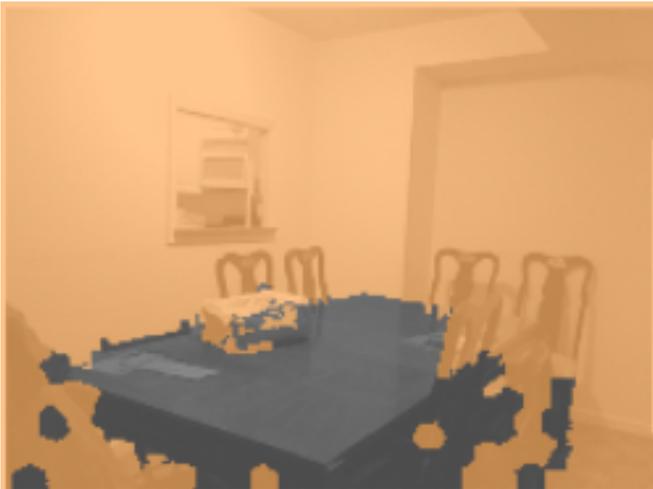


CRF

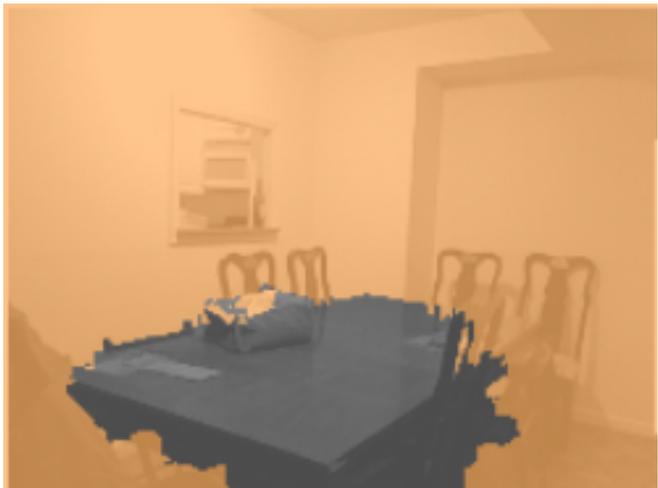
# Results



Ground Truth



Data Term

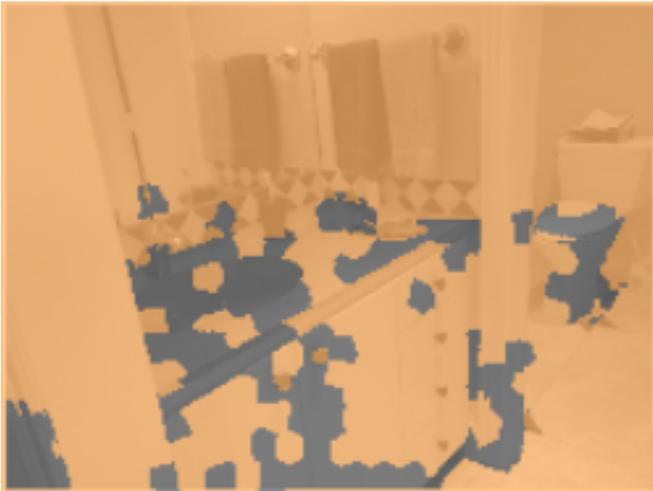


CRF

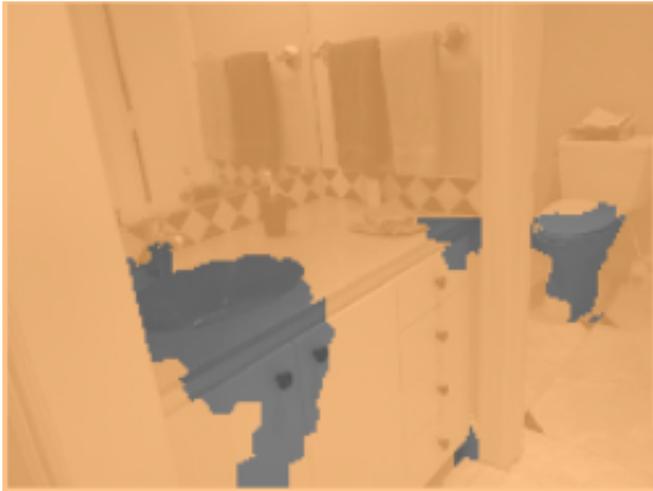
# Results



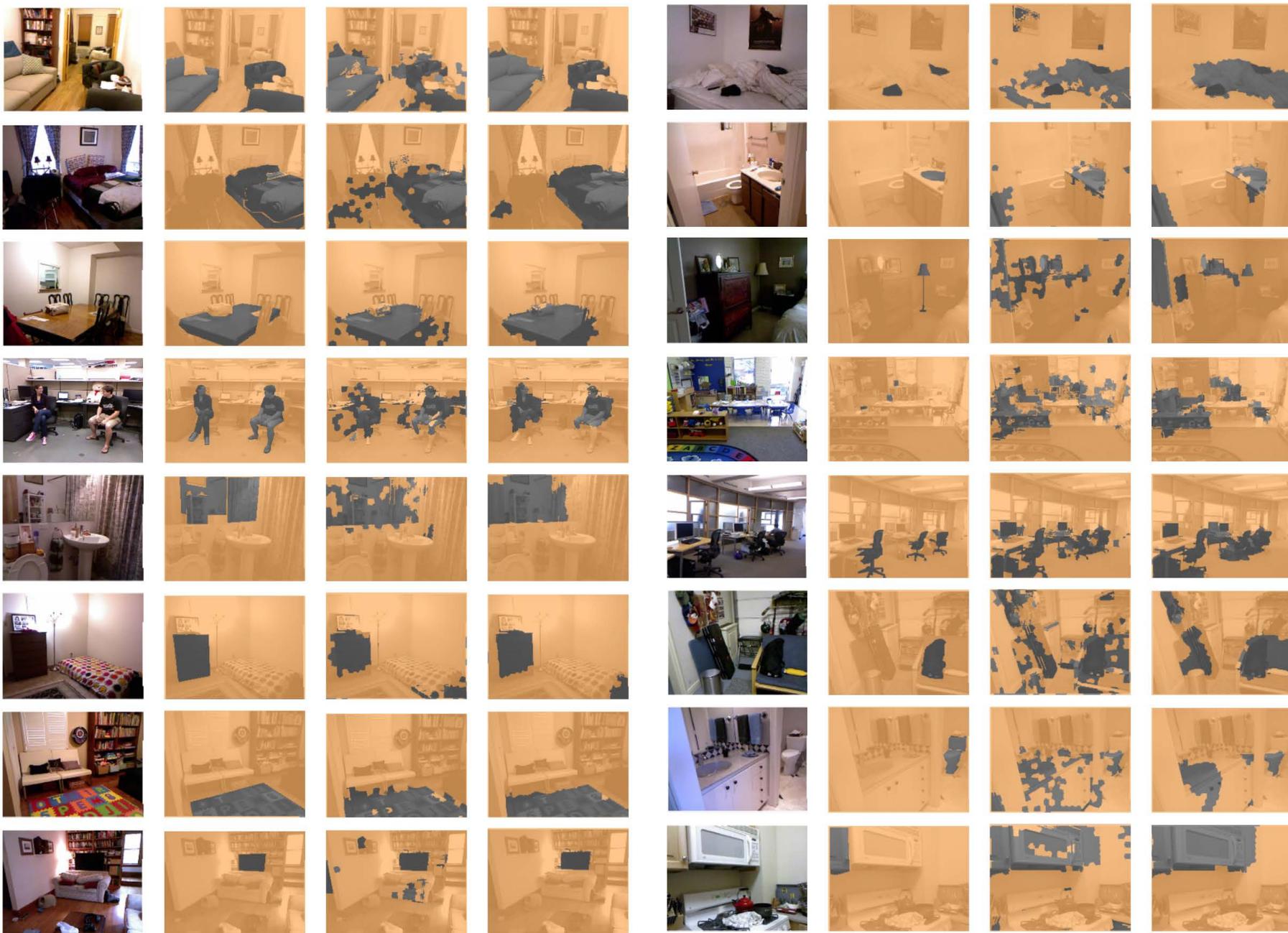
Ground Truth



Data Term



CRF



# Conclusions

- Computationally efficient approach for semantic segmentation, effective use of image and 3D cues
- Proposed Semantic Hierarchy: Background/Objects
- CRF Framework, Efficient exact inference on trees in 3D
- Recursive setting and multiple sensing modalities
- Refining Semantic Hierarchy for Objects

## Life-long Semantic Mapping

- **Reusable Representations** of sensory streams, which will generalize across different environments
- New semantic concepts can be learned incrementally, **fine grained semantic categories**
- Tightly couple localization, reconstruction, mapping