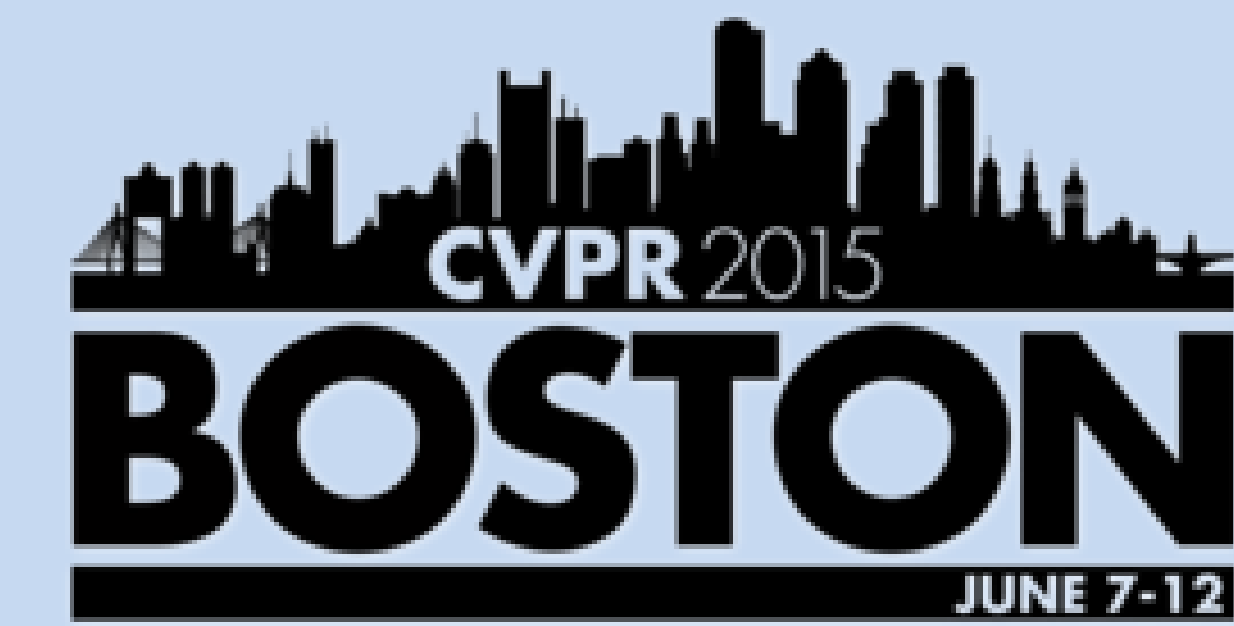


Separating Objects and Clutter in Indoor Scenes

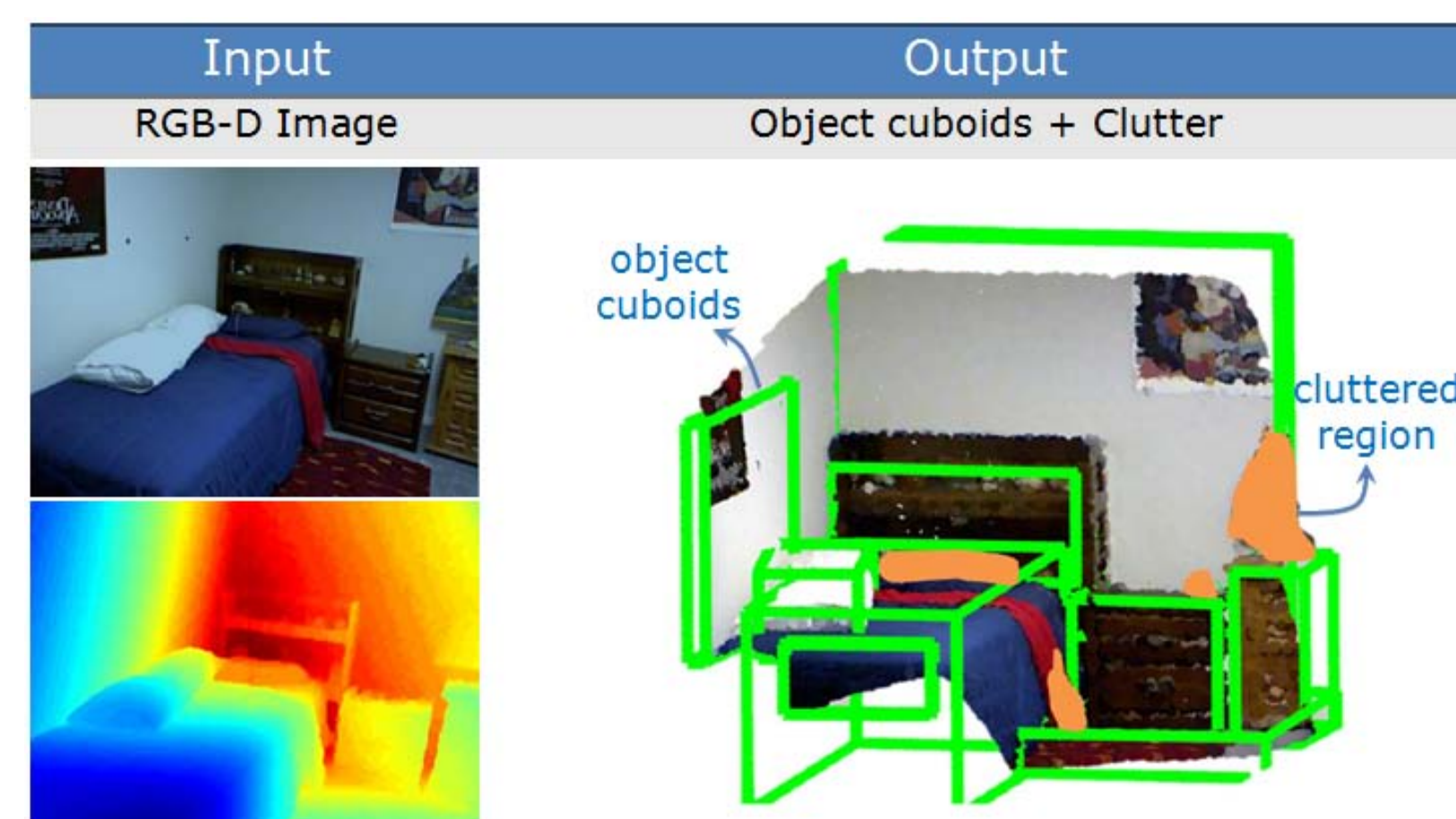
Salman H. Khan¹, Xuming He², Mohammed Bennamoun¹, Ferdous Sohel¹, Roberto Togneri³

¹ School of CSSE, ³ School of EECE, University of Western Australia, ² NICTA, Australian National University



Introduction

Problem Definition:



Highlights:

- We perform 'fine-grained structure categorization' by predicting all the major objects and structures and simultaneously labeling the cluttered regions.
- The proposed CRF model incorporates a rich set of local appearance, geometric features and interactions between the scene elements.
- We take a structural learning approach with a loss of 3D localisation to estimate the model parameters from a large annotated RGBD dataset, and a mixed integer linear programming formulation for inference.

Our Approach

- The goal is to describe an RGBD image with an optimal set of cuboids and pixel-level labeling of cluttered regions.
- We represent an indoor scene as an overlay of the cluttered regions (modeled as local surfaces) and the non-cluttered regions (modeled using 3D cuboids).
- Our method for initial cuboid hypothesis generation is based on a bottom-up clustering and fitting procedure, which generates both 'object cuboids' and 'scene bounding cuboids'.

- We build a CRF model on the superpixel clutter variables (\mathbf{m}) and the object variables (\mathbf{c}) to describe the properties of clutter, objects and their relationship in the scene.

$$E(\mathbf{m}, \mathbf{c}|\mathcal{I}) = E_{obj}(\mathbf{c}) + E_{sp}(\mathbf{m}) + E_{com}(\mathbf{m}, \mathbf{c})$$

- The first term is defined as a combination of unary and pair-wise terms on cuboids:

$$E_{obj}(\mathbf{c}) = \sum_{k=1}^K [\psi_{obj}^u(c_k) + \psi_{obj}^h(c_k)] + \sum_{i < j} \psi_{obj}^p(c_i, c_j)$$

Cuboid feature set include:

- Volumetric occupancy feature, color consistency feature, normal consistency feature, tightness feature, support feature, geometric plausibility, cuboid size feature, cuboid MDL potential.

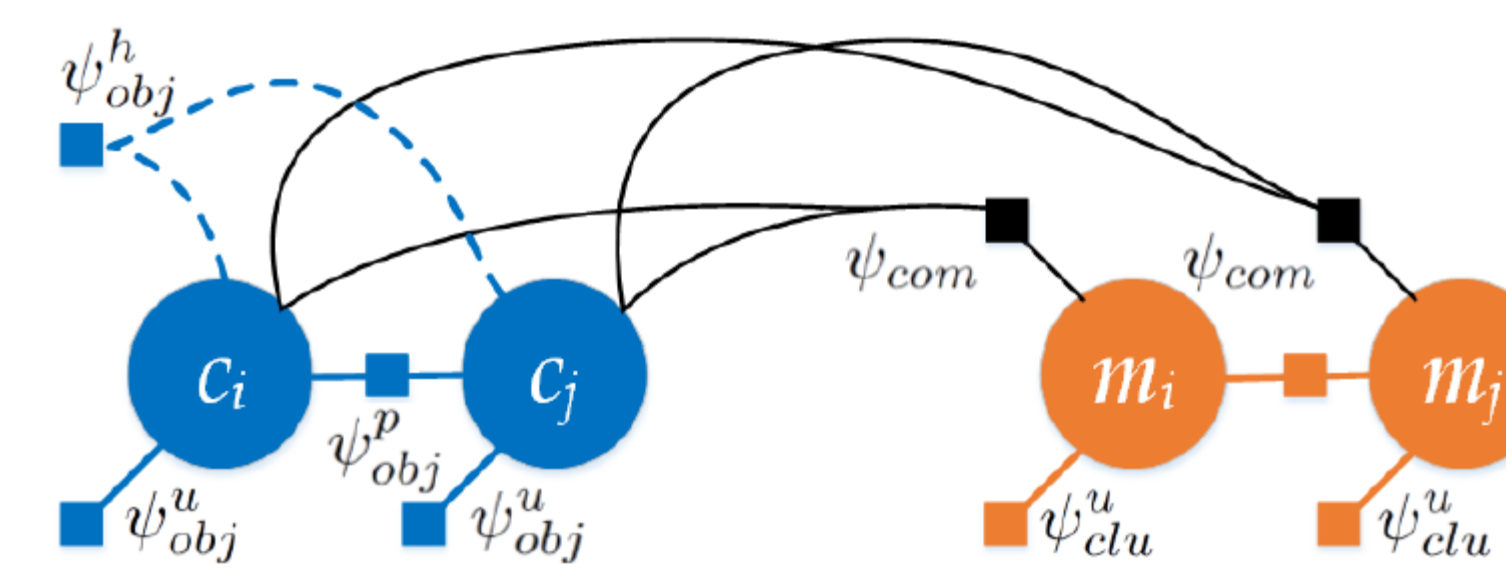


Figure: Graph structure representation for the potentials defined on the object cuboids and the cluttered/non-cluttered regions.

- View obstruction potential, cuboid intersection potential.

- Similarly at the superpixel level, the second term consists of two potential functions:

$$E_{sp}(\mathbf{m}) = \sum_{j=1}^J \psi_{sp}^u(m_j) + \sum_{(i,j) \in N_s} \psi_{sp}^p(m_i, m_j)$$

- The unary potential is based on the rich kernel descriptors (Bo et al., 2011).
- The pairwise potential is a contrast based Potts model.

- The third term is the compatibility constraint which enforces the consistency of the cuboid activations and the superpixel labeling:

$$E_{com}(\mathbf{m}, \mathbf{c}) = \sum_{j=1}^J \psi_{com}(m_j, \mathbf{c})$$

It consists of two terms:

- superpixel membership potential
- superpixel cuboid occlusion potential

Inference

- We minimize the CRF energy to get optimal labeling:

$$\{\mathbf{m}^*, \mathbf{c}^*\} = \underset{\mathbf{m}, \mathbf{c}}{\operatorname{argmin}} E(\mathbf{m}, \mathbf{c}|\mathcal{I}).$$

- We adopt the relaxation method in [Geiger'11, Jiang'13] and transform the minimization into a Mixed Integer Linear Program (MILP) with linear constraints.

$$\begin{aligned} & \min_{\mathbf{m}, \mathbf{c}, \mathbf{x}, \mathbf{y}, \mathbf{w}, \mathbf{z}} E(\mathbf{m}, \mathbf{c}, \mathbf{x}, \mathbf{y}, \mathbf{w}, \mathbf{z}|\mathcal{I}) \\ & \text{s.t. linear inequality constraints in } \mathcal{LC}, \\ & m_j, c_k \in \{0, 1\}, \quad \forall j, k \\ & w_{i,j}, x_{i,j}, y_{i,j}, z_{j,k} \geq 0, \quad \forall i, j, k \end{aligned}$$

- The MILP formulation can be solved much faster compared to the original ILP, using the branch and bound method.

- On average, 81948% variables are involved in each inference and the final MILP gap is zero for 98.5% of the cases on the whole dataset.

	Small gap	Large gap	Cuts	LP relax.
Time (sec)	1.84 ± 31%	1.31 ± 24%	0.45 ± 13%	0.001 ± 0.4%
Det. Rate	26.8%	26.1%	24.4%	19.9%

Table: Inference running time comparisons for variants of MILP formulation.

Cuboid Detection Results



Figure: Comparison of our results (3rd row) with the state of the art technique [Jiang et al., 2013] (2nd row) and Ground Truth (1st row). The detected cuboids are shown in color on top of the original RGB image.

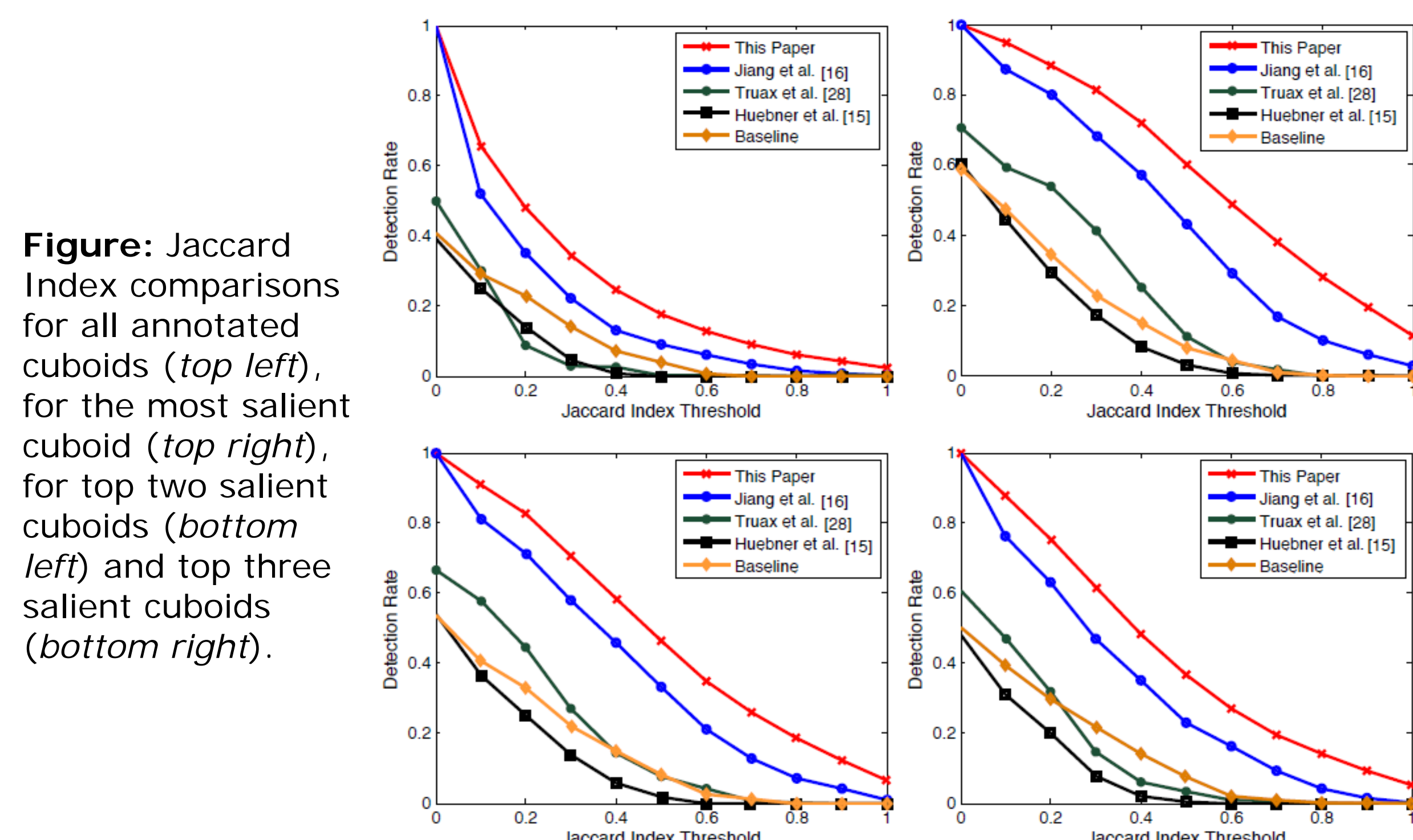


Figure: An ablation study on the model potentials/features for the cuboid detection task at the 40% JI threshold.

Method	Accuracy
Unary cuboid cost of Jiang [16]	6.5%
Our unary cuboid cost only	8.8%
Our unary + pairwise cuboid cost only	19.4%
Our full model	26.1%



Figure: Examples of detection errors.

Parameter Learning

- For parameter learning, we apply the structural SVM framework with margin re-scaling and 3D localization loss given by:

$$\Delta(\mathbf{t}^{(n)}, \mathbf{t}) = \sum_i \left(1 - \frac{|o_i^{(n)} \cap o_i|}{|o_i^{(n)} \cup o_i|} \right)$$

- The formulation uses the cutting plane algorithm [Joachims'07] to search the optimal parameter setting.

Segmentation Results

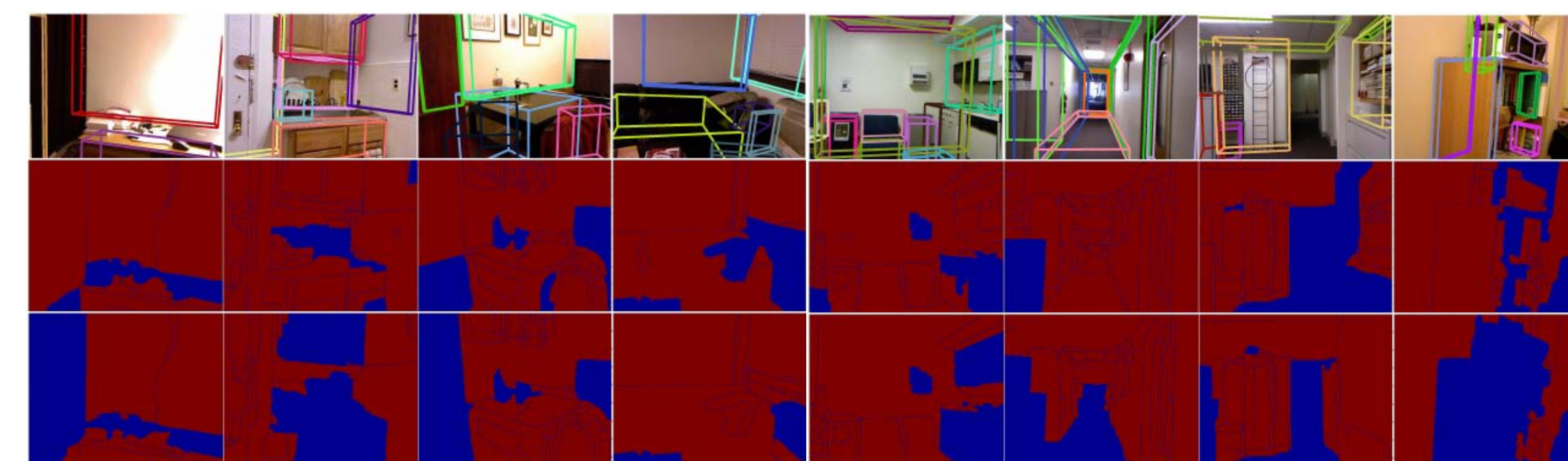


Figure: Our method is able to accurately detect cuboids in the case of cluttered indoor scenes (1st row). The 2nd and 3rd rows show our clutter labelling and the ground-truth labelling on superpixels, respectively. In the bottom two rows, red color represents non-clutter while blue color represents clutter.

Method	Precision	Recall	F-Score
Superpixel unary only	0.43 ± 13%	0.45 ± 11%	0.44 ± 16%
Unary + pairwise	0.46 ± 12%	0.48 ± 10%	0.47 ± 16%
Full model (all classes)	0.65 ± 9%	0.68 ± 8%	0.66 ± 12%
Full model (only object classes)	0.75 ± 6%	0.71 ± 8%	0.73 ± 10%

Table: Evaluation on Clutter/Non-Clutter Segmentation Task. Precision signifies the accuracy of clutter classification.

Eval. Criterion	CPMC [6]		This Paper	
	Prc.	Rec.	Prc.	Rec.
Most salient obj.	0.83 ± 11%	0.79 ± 12%	0.85 ± 15%	0.82 ± 15%
Top 2 salient obj.	0.77 ± 12%	0.73 ± 14%	0.81 ± 16%	0.79 ± 16%
Top 3 salient obj.	0.69 ± 15%	0.66 ± 17%	0.79 ± 21%	0.76 ± 19%
All objects	0.54 ± 17%	0.51 ± 20%	0.73 ± 23%	0.69 ± 21%

Table: Evaluation on Foreground/Background Segmentation Task. Precision signifies the accuracy of clutter classification.