

CS395T paper review

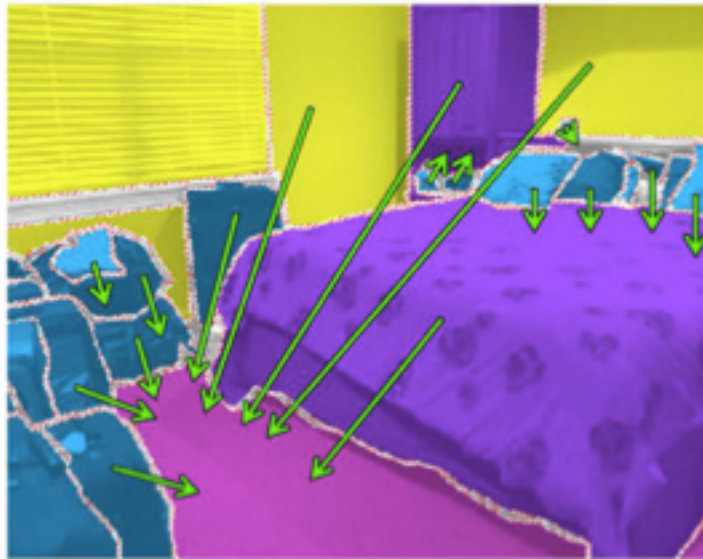
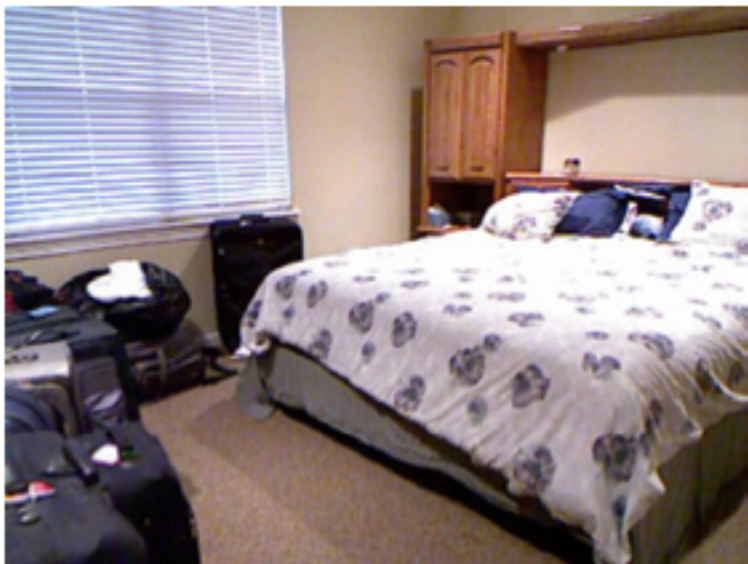
# Indoor Segmentation and Support Inference from RGBD Images

Chao Jia

Sep 28 2012

# Introduction

- What do we want -- Indoor scene parsing
  - Segmentation and labeling
  - **Support relationships**



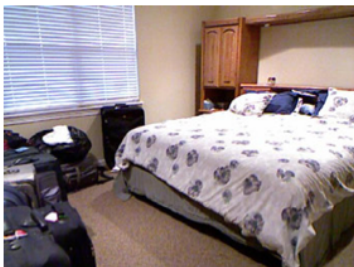
Different colors show different kinds of objects;

Support relationships help understand the scene and interact with scene elements.

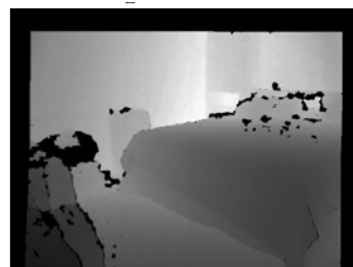
# Introduction

- What do we have

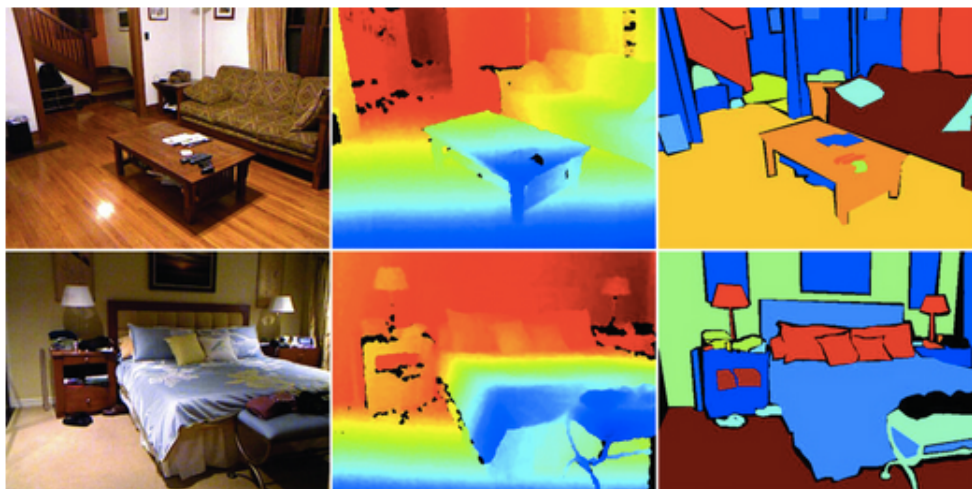
- Color image



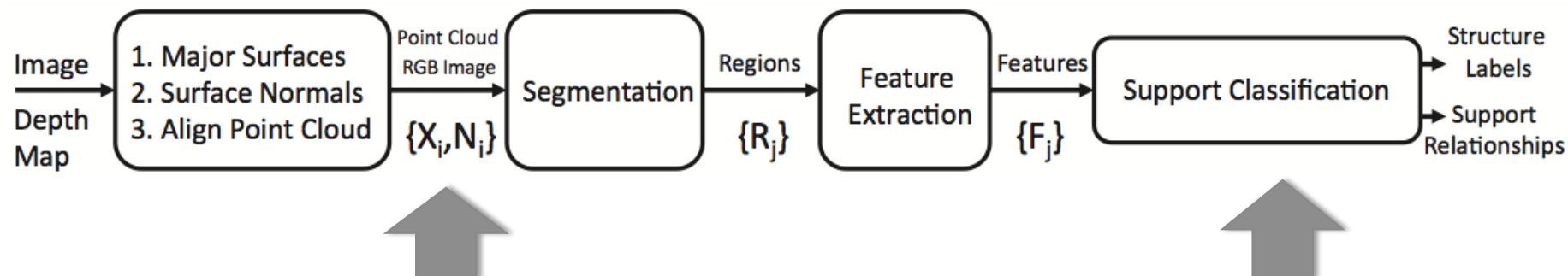
- Depth image (3D coordinates)



- How 3D cues can best inform a structured 3D interpretation
- Dataset with 1449 densely labeled images

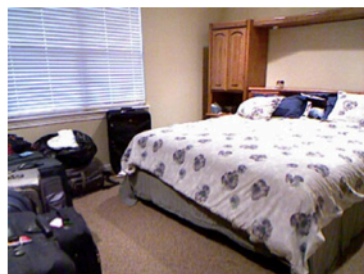


# General Steps



How 3D cues help  
scene interpretation

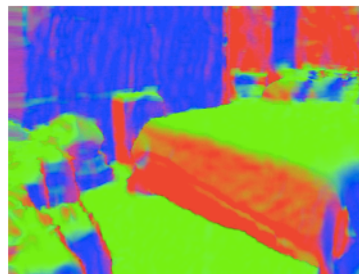
Integer programming  
formulation



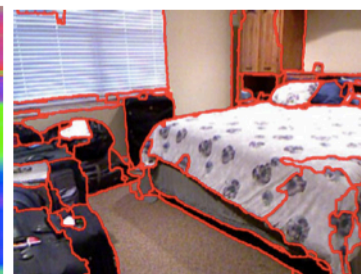
Input RGB



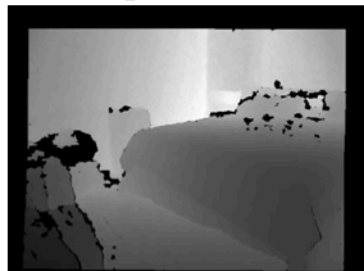
Surface Normals



Aligned Normals



Segmentation



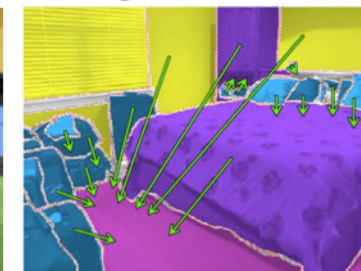
Input Depth



Inpainted Depth



3D Planes

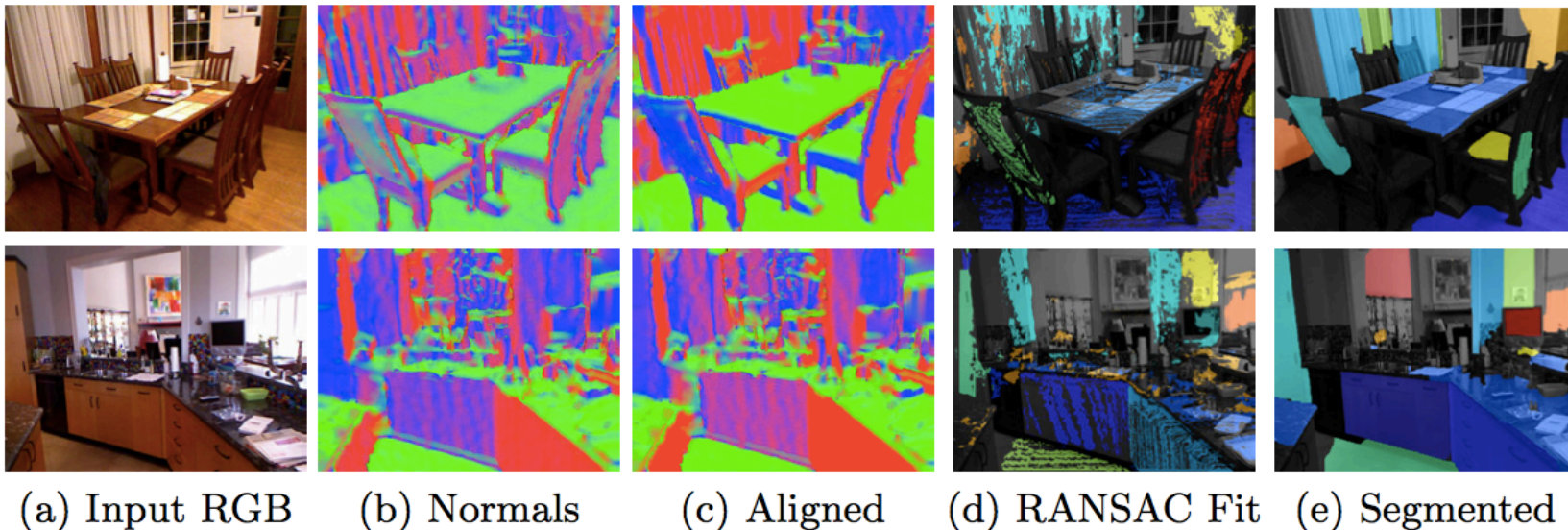


Support Relations



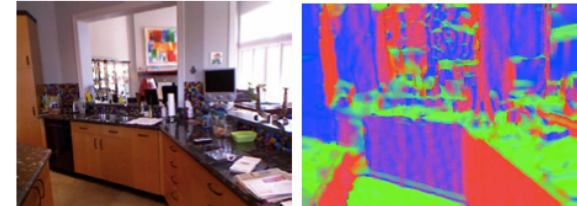
# Scene Structure Modeling

- Align the room with the 3 principle directions
  - Compute 3D lines and surface normals
  - Find the most probable X-Y-Z axis
- Segment the visible regions into 3D planes
  - Propose 3D planes using RANSAC
  - Segment the image into the proposed planes



# Aligning to Room Coordinates

- Preparation using 3D coordinates
  - Straight line segments
  - 3D surface normals at each pixel
- Propose candidates (100-200)
  - All the straight 3D lines
  - Mean-shift modes of surface normals
- Search for the most probable X-Y-Z triple
  - Random sample a triple, compute the score
- Choose the triple with highest score
- Warp the image to align with principle directions

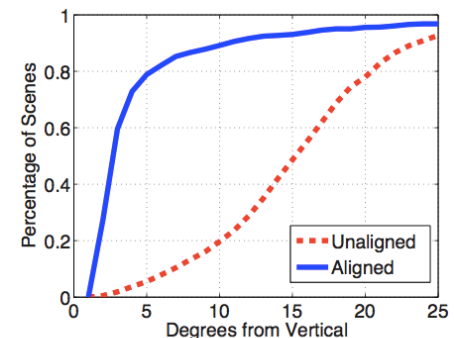


(a) Input RGB

(b) Normals

Manhattan  
world  
assumption

$$S(v_1, v_2, v_3) = \sum_{j=1}^3 \left[ \frac{w_N}{N_N} \sum_i^{N_N} \exp\left(-\frac{(\mathbf{N}_i \cdot \mathbf{v}_j)^2}{\sigma^2}\right) + \frac{w_L}{N_L} \sum_i^{N_L} \exp\left(-\frac{(\mathbf{L}_i \cdot \mathbf{v}_j)^2}{\sigma^2}\right) \right]$$



**Fig. 3.** Alignment of Floors

# Proposing and Segmenting Planes

- Generating potential planes
  - Sample the grid of pixel and propose planes (>2500 inliers)
- Assign each pixel a label to a certain plane
  - Latent variables to infer: plane label
  - Observable variables: 3D coordinates, RGB intensities, **surface normals**
  - Conditional random field modeling solved by graph cuts

$$E(\mathbf{data}, \mathbf{y}) = \alpha_i \left[ \sum_i f_{3d}(\mathbf{X}_i, y_i) + f_{norm}(\mathbf{N}_i, y_i) \right] + \sum_{i,j \in \mathcal{N}_8} f_{pair}(y_i, y_j, \mathbf{I})$$

3D coordinates

unary term

surface normals

pairwise term

RGB intensities

# Proposing and Segmenting Planes

- Unary term  $f_{3d}(\mathbf{X}_i, y_i) + f_{norm}(\mathbf{N}_i, y_i)$ 
  - Geometrically validate the labels

$$-\log \frac{Pr(dist|inlier)}{Pr(dist|outlier)}$$

from RANSAC  
plane proposing

- Pairwise term  $f_{pair}(y_i, y_j, \mathbf{I})$

$$\mathbf{1}(y_i \neq y_j) \exp \left( -(\beta_1 + \beta_2 \|\mathbf{I}_i - \mathbf{I}_j\|^2) \right)$$

smoothness weighed by  
RGB intensity difference

# Segmentation

- Oversegmentation into superpixels
  - Boundaries detection from **RGB intensities**
  - Force consistency with **3D planes regions**
- Iterative merging of regions
  - Regions with minimum boundary strength are merged
  - Boundary strength:  $P(y_i \neq y_j | \mathbf{x}_{ij}^s)$ 
    - Trained boosted decision tree classifier
    - $y$ : labels of regions
    - $x$ : paired regions features



# Segmentation

- Paired region features
  - **RGB features:** crucial for nearby or touching objects
  - **3D features (plane labels, surface normals, depth):** help differentiate between texture and object edges



# Modeling Support Relationships

- Variables to infer for each region (  $R$  regions in total)
  - $S_i \in \{1 \dots R, h, g\}$  the support region
    - supported by other regions
    - supported by an invisible region
    - not supported (ground)
  - $T_i \in \{0, 1\}$  supported from below/behind
  - $M_i \in \{1, 2, 3, 4\}$  structure class
    - 1: Ground
    - 2: Furniture (large objects that cannot be carried)
    - 3: Prop (small objects that can be easily carried)
    - 4: Structure (walls, ceiling, columns)

# Modeling Support Relationships

- Energy minimization

$$\{\mathbf{S}^*, \mathbf{T}^*, \mathbf{M}^*\} = \arg \max_{\mathbf{S}, \mathbf{T}, \mathbf{M}} P(\mathbf{S}, \mathbf{T}, \mathbf{M} | I) = \arg \min_{\mathbf{S}, \mathbf{T}, \mathbf{M}} E(\mathbf{S}, \mathbf{T}, \mathbf{M} | I)$$

- Factorize posterior distribution

$$P(\mathbf{S}, \mathbf{T}, \mathbf{M} | I) \propto \underbrace{\prod_{i=1}^R P(I | S_i, T_i) P(I | M_i)}_{\text{likelihood + factorization}} \underbrace{P(\mathbf{S}, \mathbf{T}, \mathbf{M})}_{\text{Prior}}$$

- Final problem

$$E(\mathbf{S}, \mathbf{T}, \mathbf{M}) = - \sum_{i=1}^R \underbrace{\log(D_s(F_{i,S_i}^s | S_i, T_i) + \log(D_m(F_i^m | M_i)))}_{\text{likelihood + factorization}} + \underbrace{E_P(\mathbf{S}, \mathbf{T}, \mathbf{M})}_{\text{Prior}}$$

# Modeling Support Relationships

- Prior term

$$E_P(\mathbf{S}, \mathbf{T}, \mathbf{M}) = \sum_{i=1}^R \psi_{TC}(M_i, M_{S_i}, T_i) + \psi_{SC}(S_i, T_i) + \psi_{GC}(S_i, M_i) + \psi_{GGC}(\mathbf{M})$$

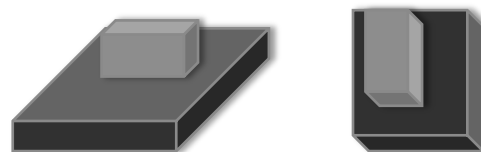
- Transition prior (supporting relationship between two structure classes)

$$\psi_{TC}(M_i, M_{S_i}, T_i) \propto -\log \frac{\sum_{z \in \text{supportLabels}} \mathbb{1}[z = [M_i, M_{S_i}, T_i]]}{\sum_{z \in \text{supportLabels}} \mathbb{1}[z = [M_i, *, T_i]]}$$

which combination is more likely

- Support consistency (between 3D structure and support relationship)

$$\psi_{SC}(S_i, T_i) = \begin{cases} (H_i^b - H_{S_i}^t)^2 & \text{if } T_i = 0 \\ V(i, S_i)^2 & \text{if } T_i = 1 \end{cases}$$



- Global ground consistency

$$\psi_{GGC}(\mathbf{M}) = \sum_{i=1}^R \sum_{j=1}^R \begin{cases} \kappa & \text{if } M_i = 1 \wedge H_i^b > H_j^b \\ 0 & \text{otherwise,} \end{cases}$$

Everything is above floor


- Ground consistency

$$\psi_{GC}(S_i, M_i) = \begin{cases} \infty & \text{if } S_i = g \text{ and } M_i \neq 1 \\ 0 & \text{else} \end{cases}$$

No support for floor

# Modeling Support Relationships

- Likelihood term

$$-\sum_{i=1}^R \log(D_s(F_{i,S_i}^s | S_i, T_i) + \log(D_m(F_i^m | M_i))$$


support relation classifier

structure classifier

- $F_{i,S_i}^S$  support features

*proximity, containment, characteristics of supporting objects, absolute 3D locations of candidate objects*

- $F_i^M$  structure features

*SIFT features, color histogram, ... (object classification)*

- Classifiers trained by logistic regression



# Modeling Support Relationships

- Introduce Boolean indicator variables:

$$\begin{aligned}
 & \arg \min_{\mathbf{s}, \mathbf{m}, \mathbf{w}} \sum_{i,j} \theta_{i,j}^s s_{i,j} + \sum_{i,u} \theta_{i,u}^m m_{i,u} + \sum_{i,j,u,v} \theta_{i,j,u,v}^w w_{i,j}^{u,v} \\
 & \text{s.t. } \sum_j s_{i,j} = 1, \quad \sum_u m_{i,u} = 1 \quad \forall i \\
 & \quad \sum_{j,u,v} w_{i,j}^{u,v} = 1, \quad \forall i \\
 & \quad s_{i,2R'+1} = m_{i,1}, \quad \forall i \\
 & \quad \sum_{u,v} w_{i,j}^{u,v} = s_{i,j}, \quad \forall u, v \\
 & \quad \sum_{j,v} w_{i,j}^{u,v} \leq m_{i,u}, \quad \forall i, u \\
 & \quad s_{i,j}, m_{i,u}, w_{i,j}^{u,v} \in \{0, 1\}, \quad \forall i, j, u, v
 \end{aligned}$$

- Problem is linearized !
- Integer programming  $\rightarrow$  relax the integrality constraints

# Experiments

- Segmentation evaluation
  - measured as average overlap over ground truth regions for best-matching segmented region

Features	Weighted Score	Unweighted Score
RGB Only	52.5	48.7
Depth Only	55.9	47.3
RGBD	62.7	52.7
RGBD + Support	63.4	53.7
RGBD + Support + Structure classes	63.9	54.1

# Support Relationships Evaluation

- Evaluate proposed inference model against
  - Image plane rules  
(no structure class assignment)
  - Structure class rules  
(class assignment by trained classifier)
  - Support classifier  
(no structure class assignment; infer the support relationship between every pair of regions)
- Metric
  - Percentage of correct supports

# Support Relationships Evaluation

Predicting Support Relationships				
Region Source	Ground Truth		Segmentation	
Algorithm	Type Agnostic	Type Aware	Type Agnostic	Type Aware
Image Plane Rules	63.9	50.7	22.1	19.4
Structure Class Rules	72.0	57.7	45.8	41.4
Support Classifier	70.1	63.4	45.8	37.1
Energy Min (LP)	75.9	72.6	55.1	54.5

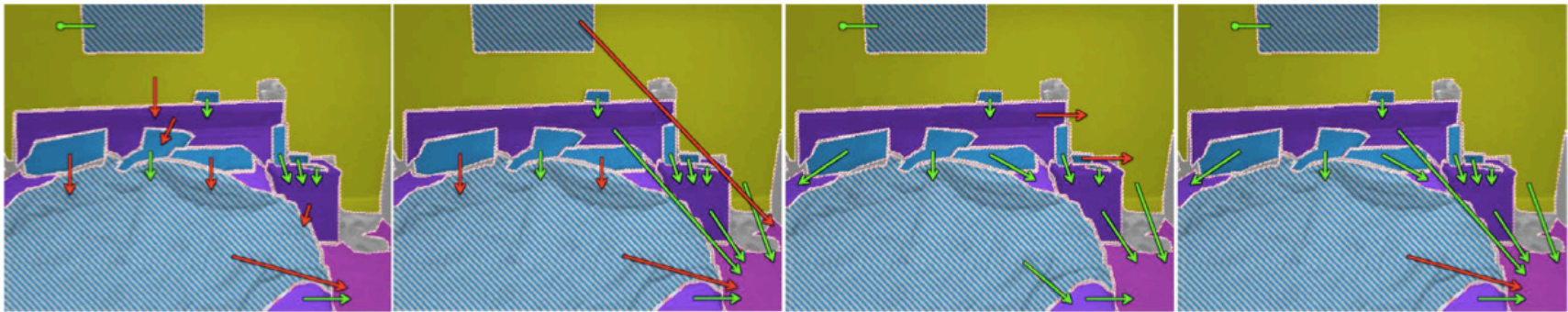


Image Plane Rules

Structure Class Rules

Support Classifier

Energy Minimization

# Experiments

- Structure class prediction evaluation
  - only slightly better than local classification

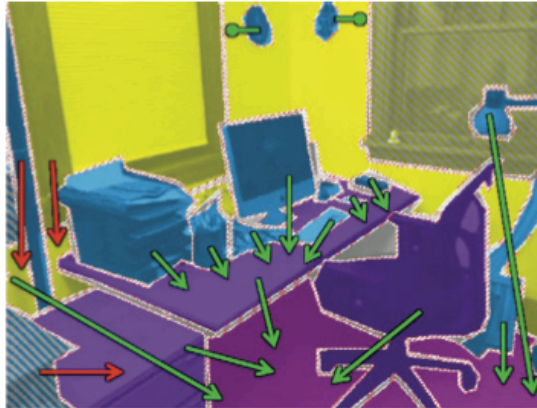
Predicting Structure Classes				
	Overall		Mean Class	
Algorithm	G. T.	Seg.	G. T.	Seg.
Classifier	79.9	58.7	79.2	59.0
Energy Min (LP)	80.3	58.6	80.3	59.6

Labels	Ground	Furniture	Prop	Structure
	.68	.28	.02	.02
	.04	.70	.14	.12
	.03	.43	.42	.12
	.01	.24	.14	.59
Predictions				



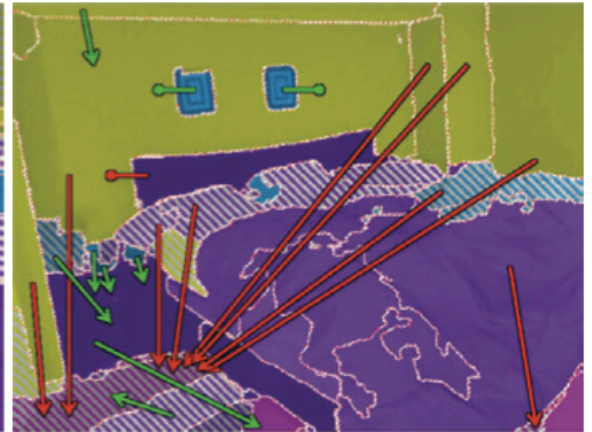
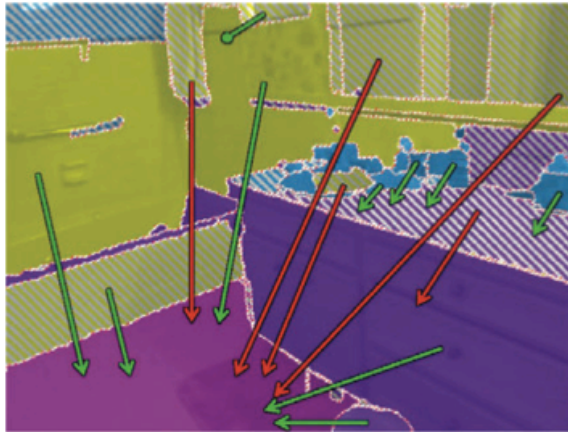
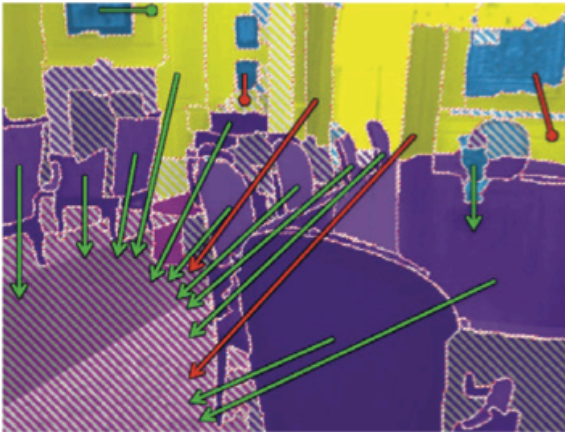
# More results

- Using ground-truth segmentation



# More results

- Using proposed segmentation



# Summary

- Pros
  - 3D features (planes, surface normals, 3D coordinates) help segmentation and support relationship inference
  - Globally infer the support relationships with high accuracy (50% - 70%)
- Cons
  - Too many functions based on training ---- training time and training data size
  - What is a good factorization of the posterior distribution in inference of support relationships ---- Are structure class features and support features really separable ?
  - Should we consider more kinds of objects instead of just props (to make features more distinguishable) ?