

# Scene Understanding

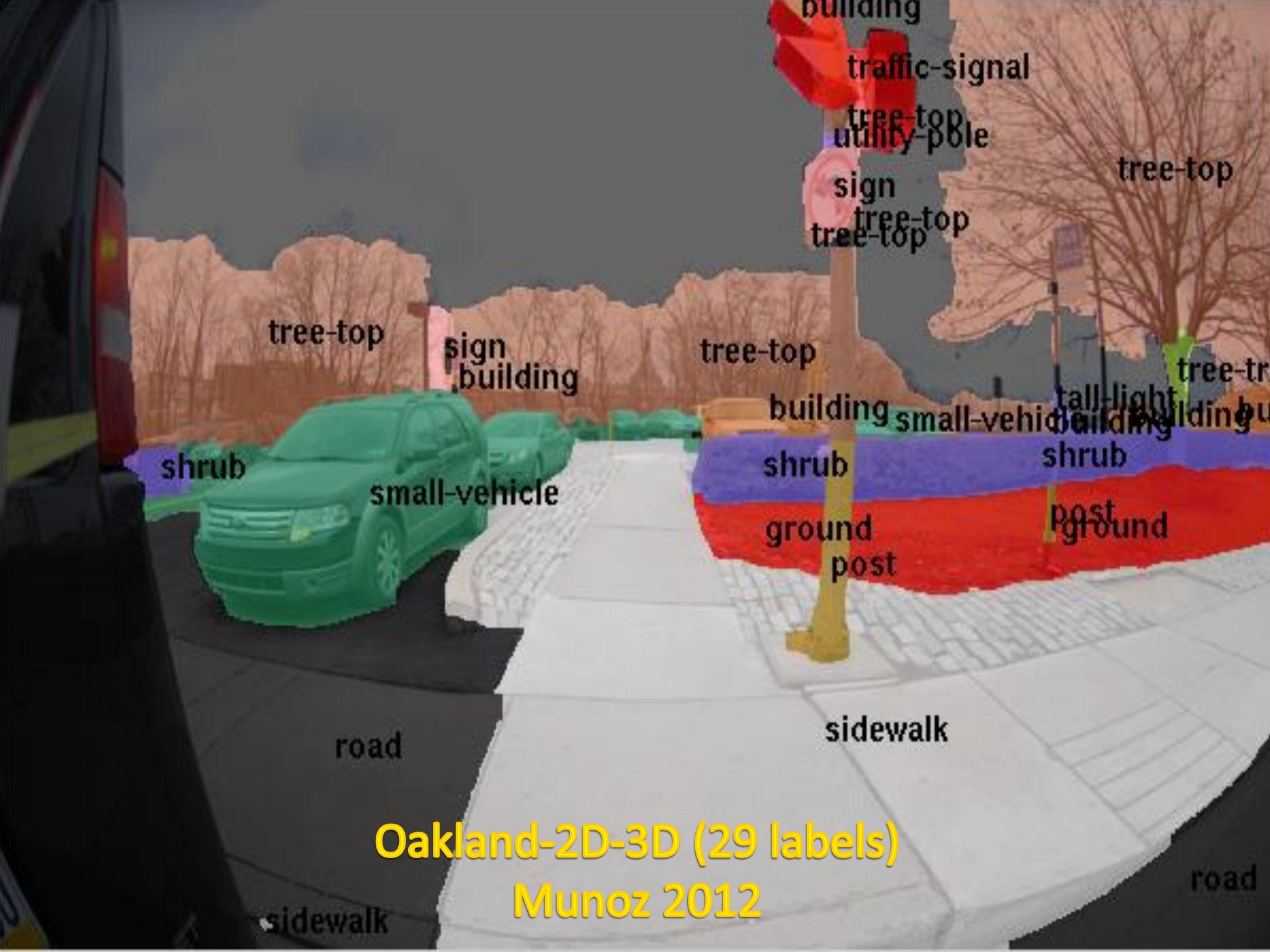


“Now! ... *That* should clear up  
a few things around here!”



Semantic labeling  
Image parsing  
Semantic segmentation





Oakland-2D-3D (29 labels)  
Munoz 2012



CamVid (11 labels)  
Miksik 2013

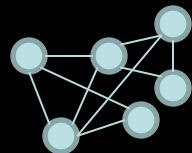
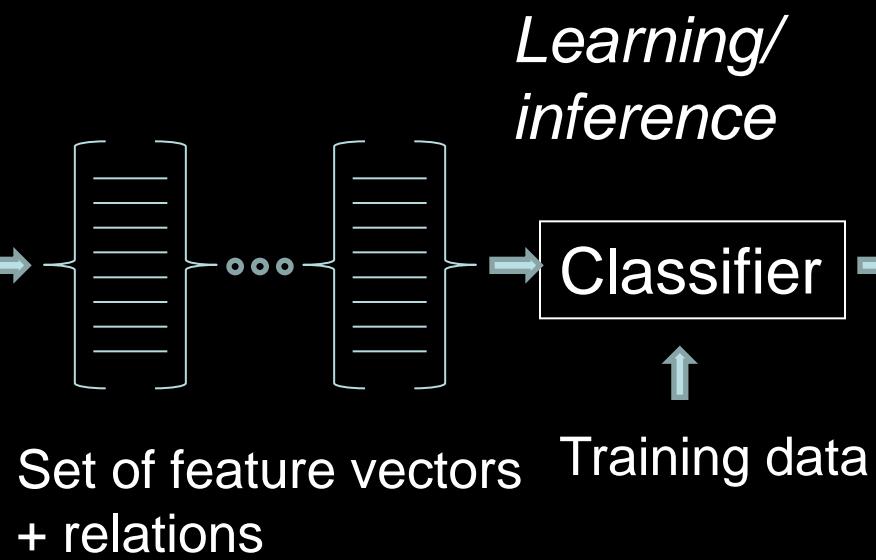
sky	tree	road	sidewalk	building	car
column_pole	pedestrian	bicycle	fence	sign_symbol	

# Outline

- Semantic: Labeling scene regions and objects
- Geometric: Estimating the geometric structure of the scene



Input image



# 1. Independent predictions



$$y^* = \arg \max_y \sum_i \phi(y_i, x)$$

# 1. Independent predictions



$$y^* = \arg \max_y \sum_i \phi(y_i, x)$$

$$P(y|x) = \frac{1}{Z} \prod_i e^{\phi(y_i, x)}$$

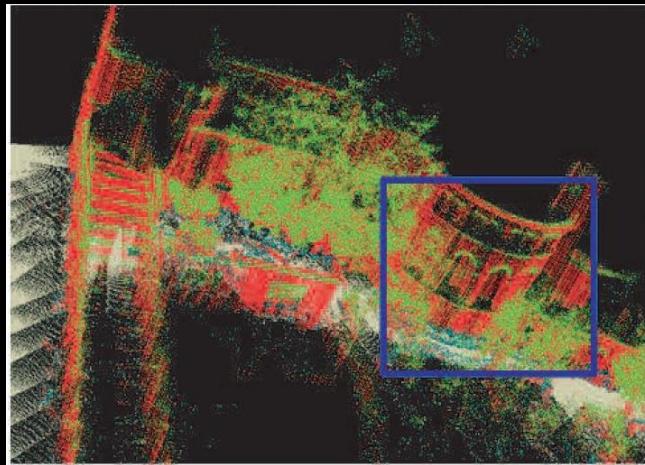
# 1. Independent predictions



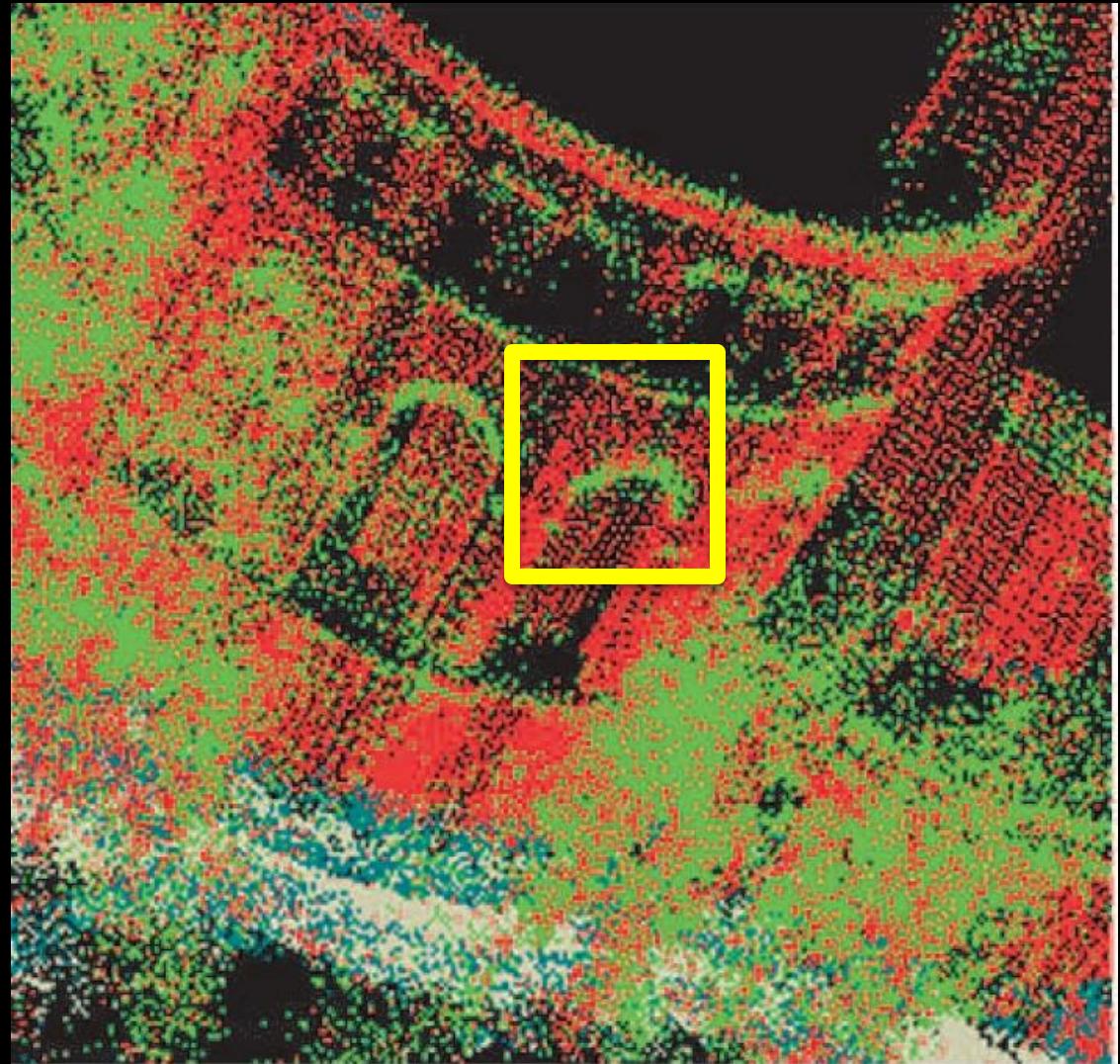
$$y^* = \arg \max_y \sum_i \phi(y_i, x)$$

- Features: local shape, texture, color, etc.
- Predictor: SVM, MaxEnt, etc.

# Needs context



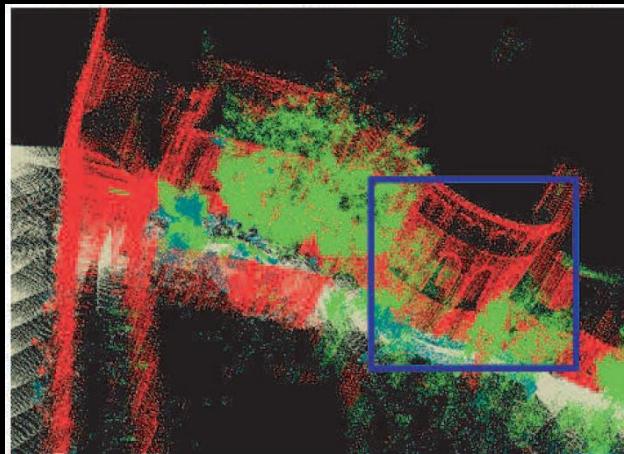
**Buildings**  
**Tree Veg**  
**Shrubs**



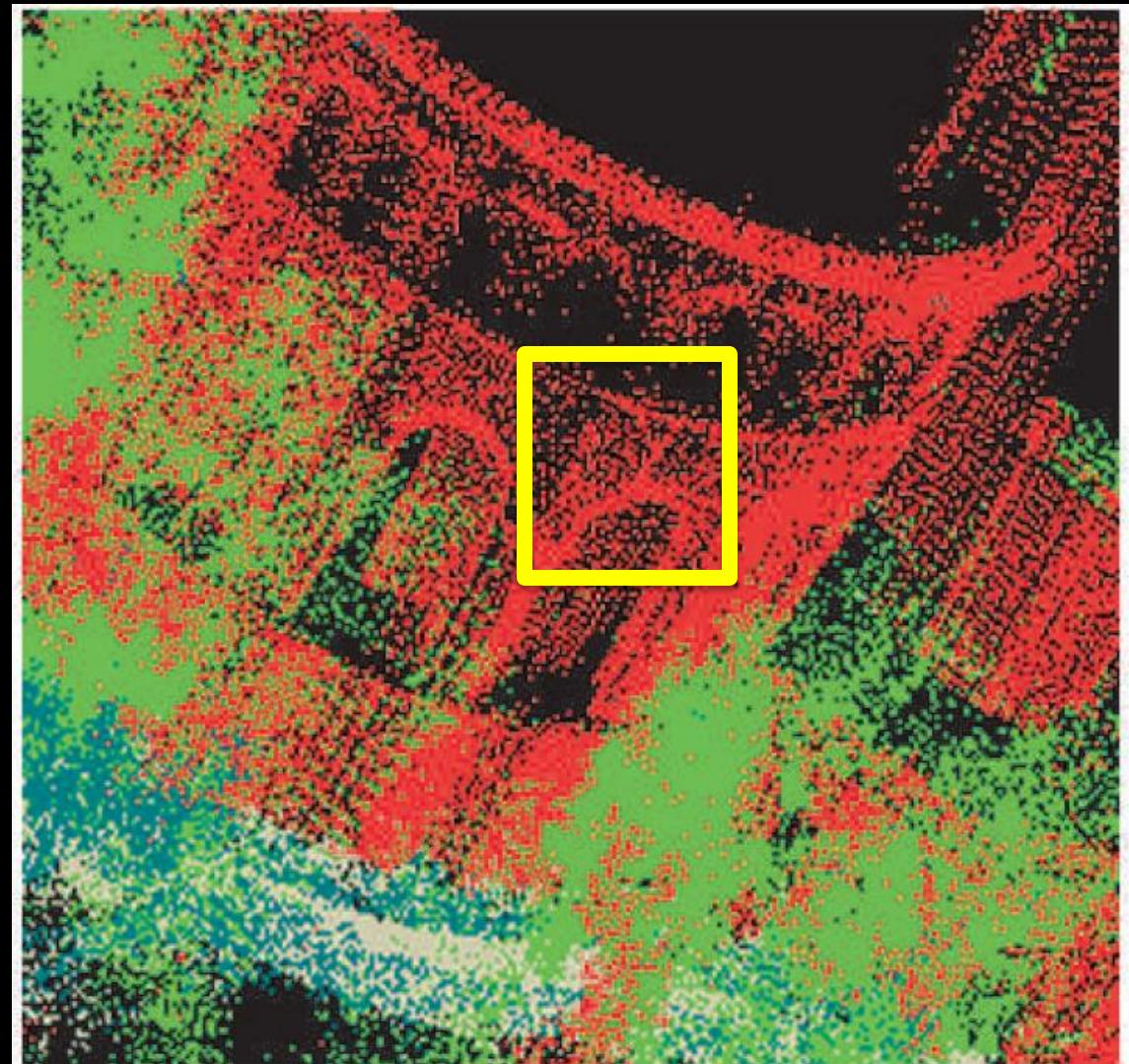
Anguelov 2005

No context

# Needs context



Buildings  
Tree Veg  
Shrubs



With Context

Anguelov 2005

## 2. Pairwise context



$$y^* = \arg \max_y \sum_i \phi(y_i, x) + \sum_{i,j} \phi(y_i, y_j, x)$$

## 2. Pairwise context



$$y^* = \arg \max_y \sum_i \phi(y_i, x) + \sum_{i,j} \phi(y_i, y_j, x)$$
$$P(y|x) = \frac{1}{Z} \prod_i e^{\phi(y_i, x)} \prod_{i,j} e^{\phi(y_i, y_j, x)}$$

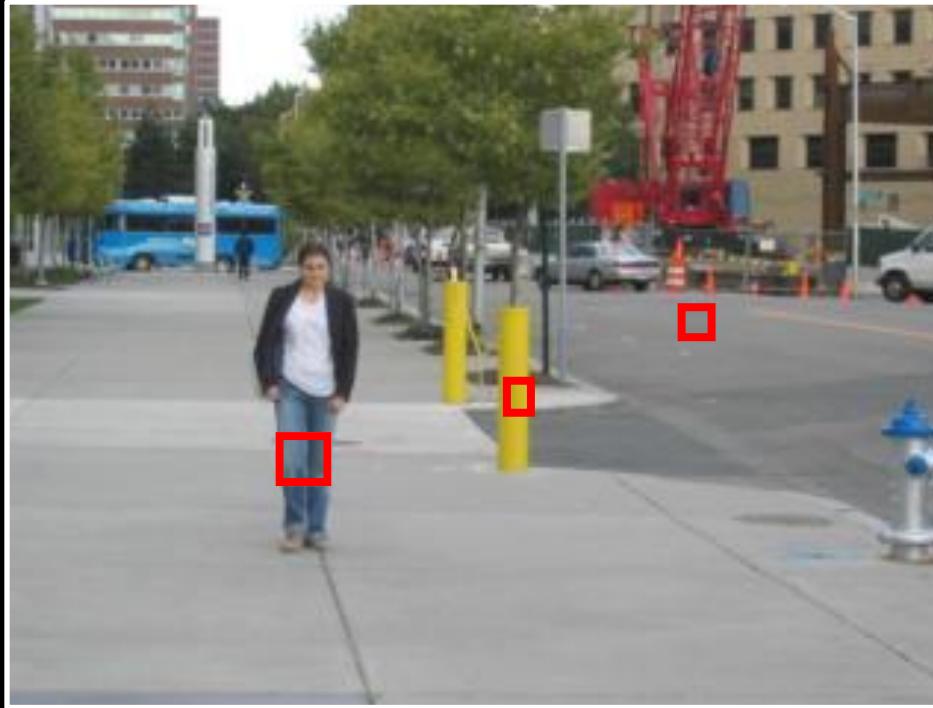
## 2. Pairwise context



$$y^* = \arg \max_y \sum_i \phi(y_i, x) + \sum_{i,j} \phi(y_i, y_j, x)$$

- Inference with BP, MCMC, LP, Graph-cuts
- Approximate inference only

# Difficult from local context



# Easier from broader context



# Ideal regions



Fig. from Tomasz Malisiewicz

# The reality



- How to represent potentials over larger support?
- How to deal with uncertainty on the choice of support regions? Multiple segmentations, hierarchical representations.

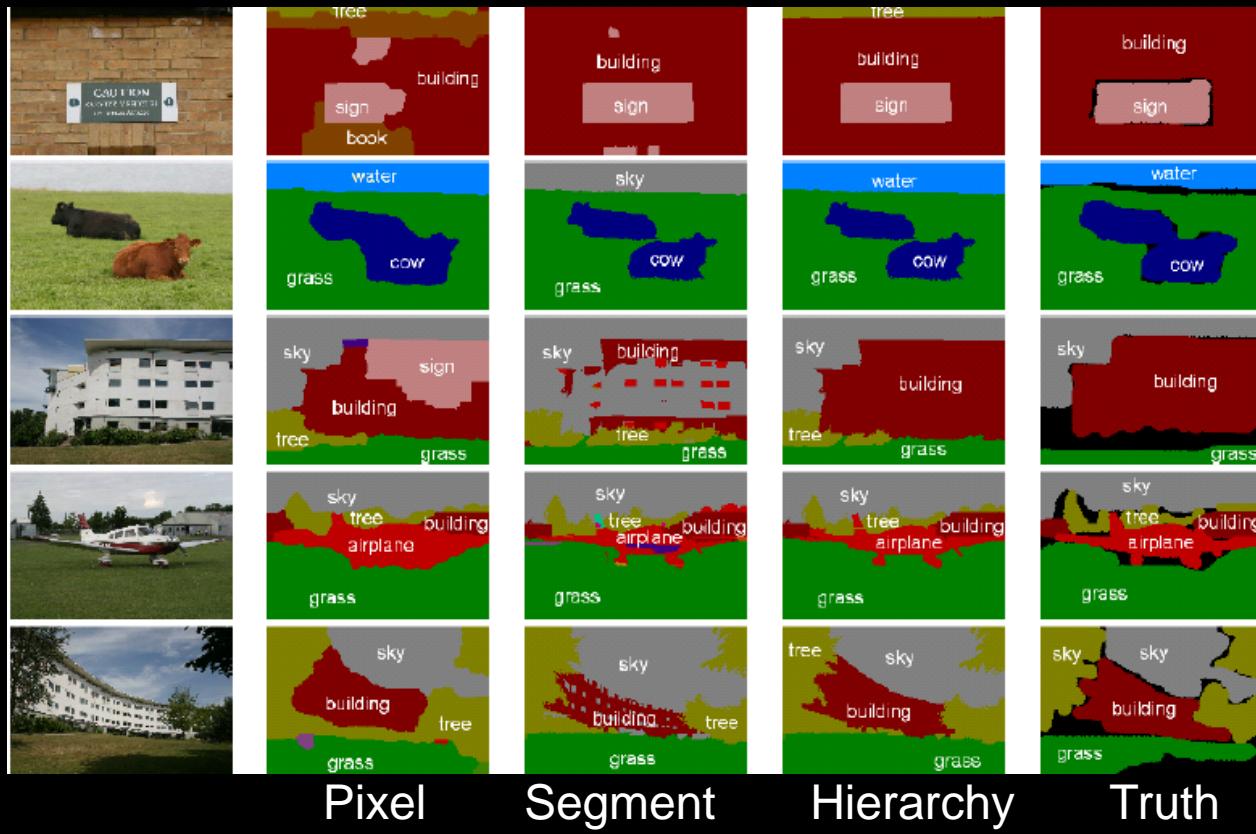
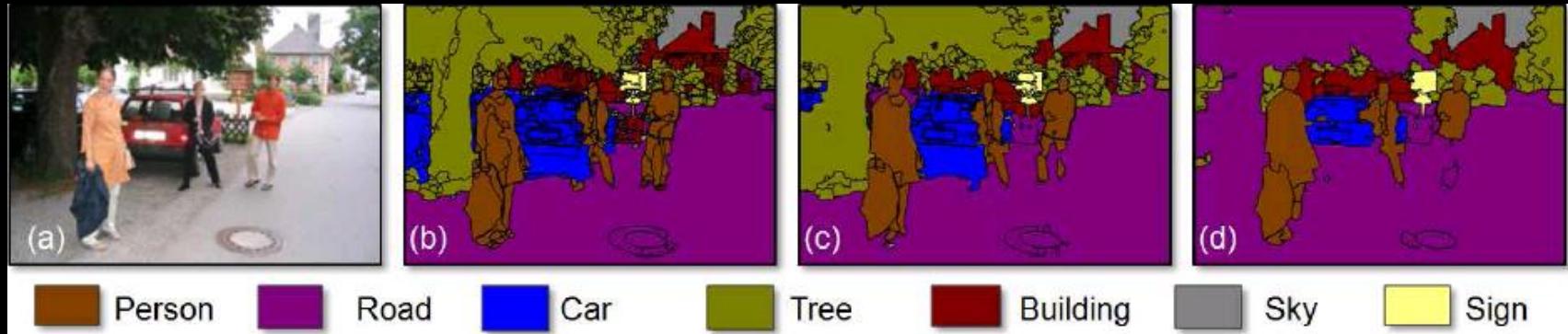
### 3. High-order context



$$y^* = \arg \max_y \sum_i \phi(y_i, x) + \sum_{i,j} \phi(y_i, y_j, x) + \sum_c \phi(y_c, x)$$

Graph cuts [Kohli 2010], MM [Taskar 2004, Munoz 2009], sampling [Gould 2010], ....

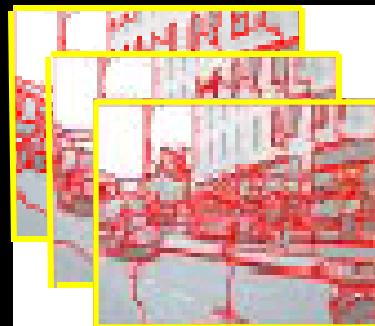
# Example



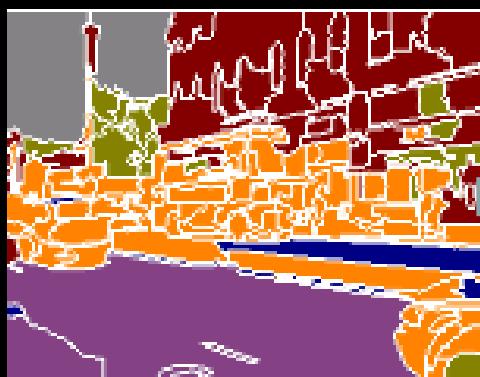
# Approaches

- Set up problem for “exact” solution (e.g., M3N..)
- Approximate solutions
- Sampling
- Decomposition into sequences of simpler problems
- Deep learning

# Example: Sampling



Multiple  
segmentations



Propose move

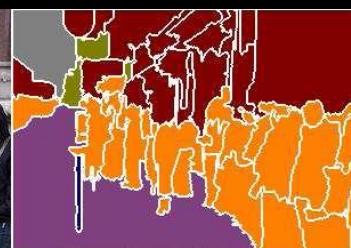
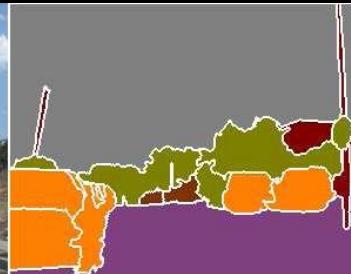
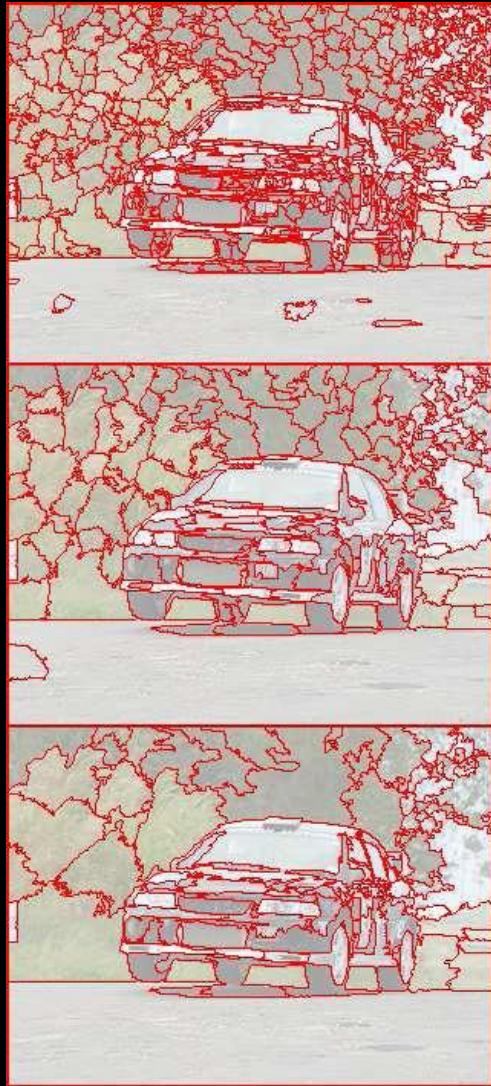
Accept if move  
reduces energy

$$E(y/x)$$

Evaluate energy

[Example from  
Gould et al.]

Gould, S., Fulton, R., Koller, D.: Decomposing a scene into geometric and semantically consistent regions. In: ICCV (2009)  
Stephen Gould, Tianshi Gao and Daphne Koller. Region-based Segmentation and Object Detection. In Advances in Neural  
Information Processing Systems (NIPS),

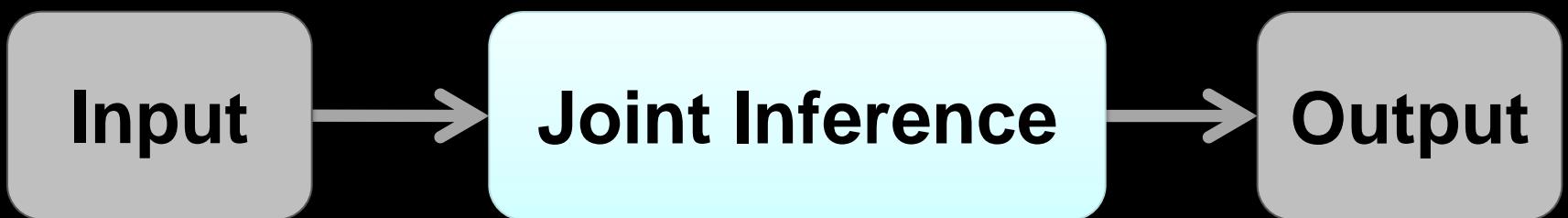
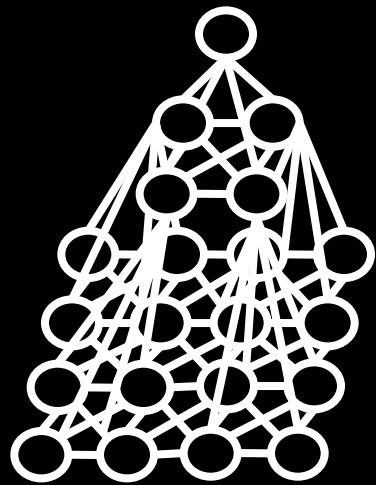


[Example from  
Gould et al.]

sky tree road grass water bldg mntn fg obj.

# Approaches

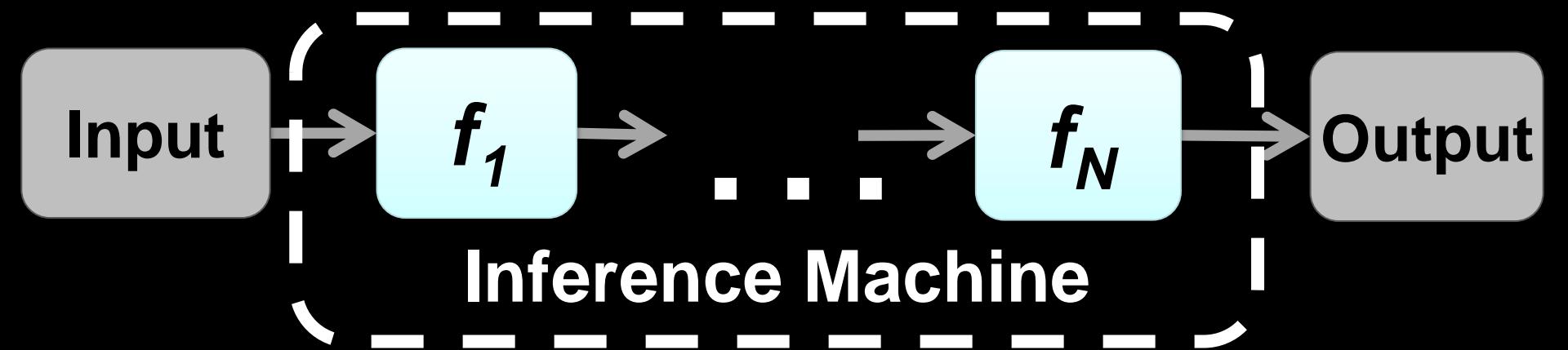
- Set up problem for “exact” solution (e.g., M3N..)
- Approximate solutions
- Sampling
- Decomposition into sequences of simpler problems
- Deep learning



$$y^* = \operatorname{argmax}_y f(x, y)$$

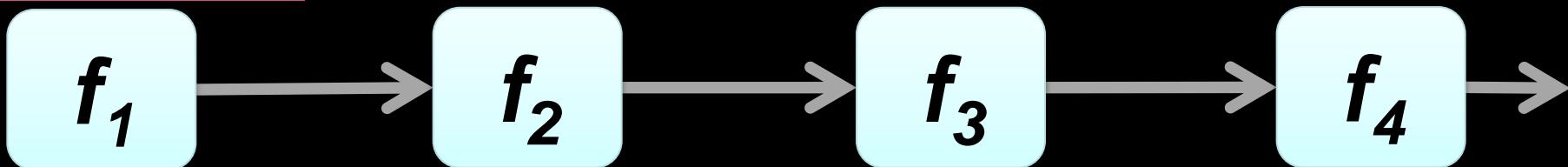
Training: Model the distribution of labels vs. features  
Inference: Find most likely labels given model

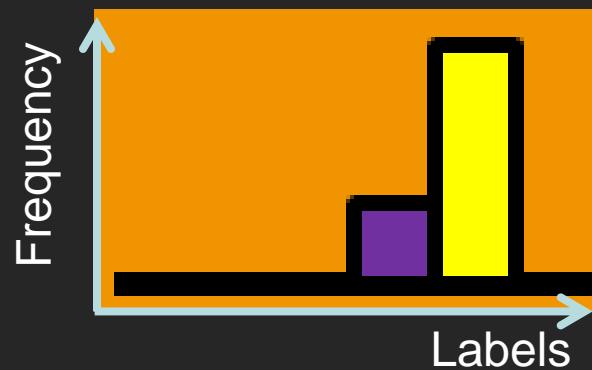
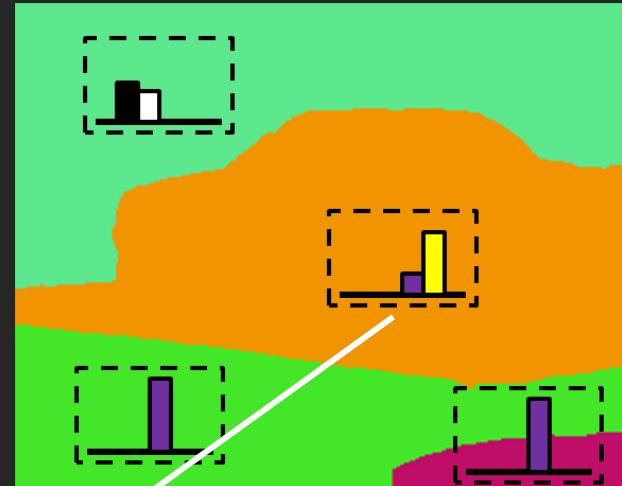
# HIM: Hierarchical Inference Machine

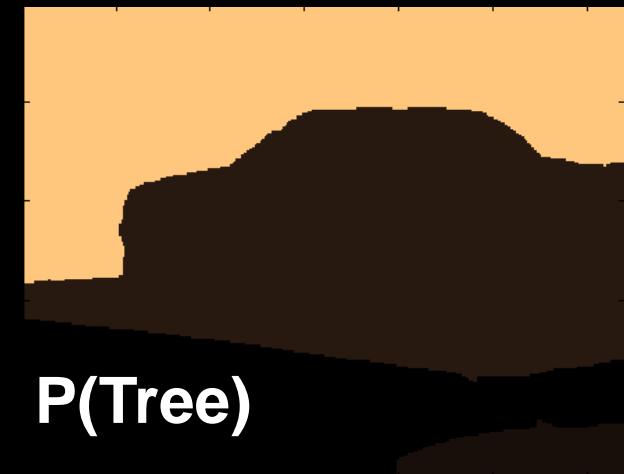
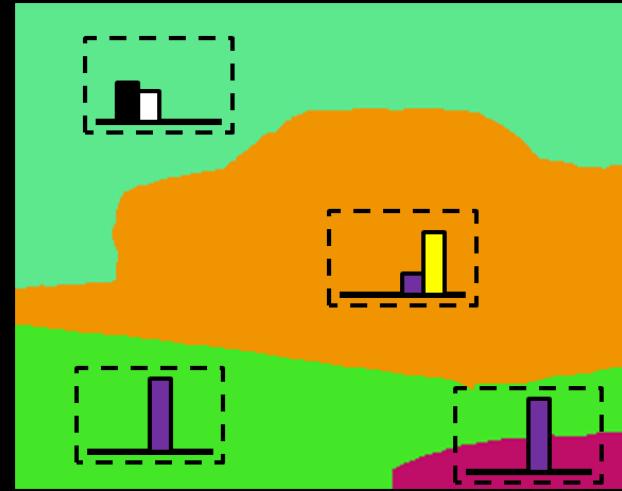


Do not learn a joint distribution/energy

Instead, learn a *procedure* designed to *iteratively* decode the scene







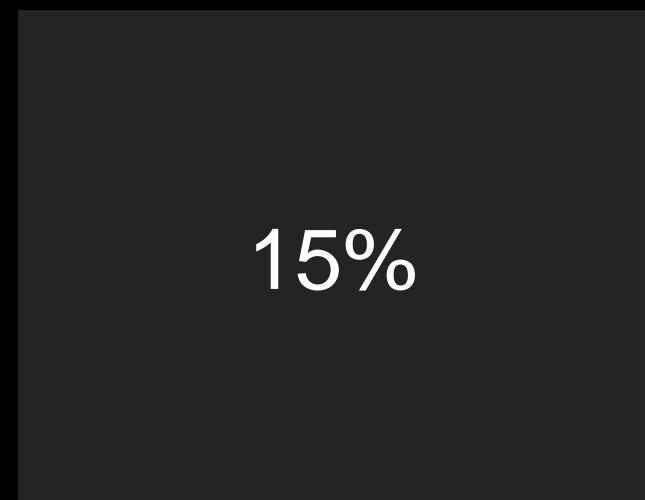
# Level 1/8 Predictions

## Segmentation



# Level 1/8 Predictions

Segmentation P(Foreground)



P(Tree)

P(Building)

P(Road)

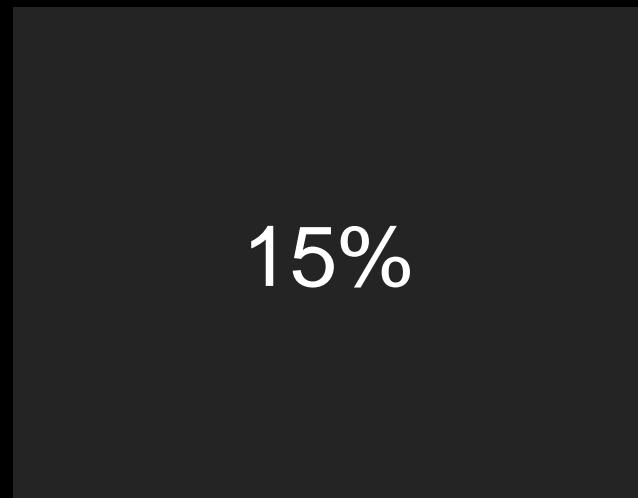
18%

12%

31%

# Level 1/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$

$P(\text{Building})$

$P(\text{Road})$

18%

12%

31%

# Level 2/8 Predictions

Segmentation P(Foreground)



P(Tree)

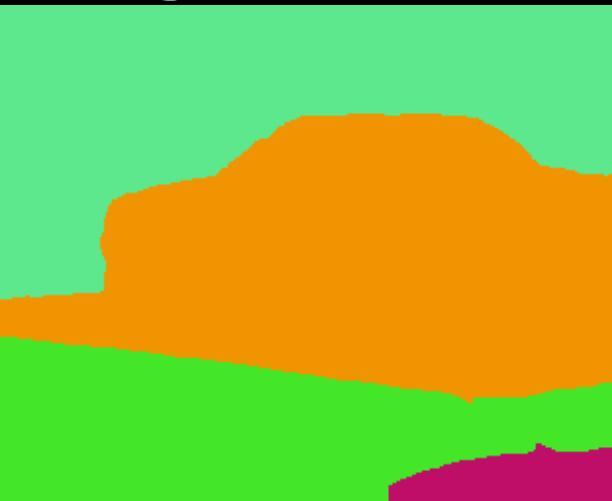
P(Building)

P(Road)



# Level 2/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$



$P(\text{Building})$



$P(\text{Road})$



# Level 3/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$



$P(\text{Building})$



$P(\text{Road})$



# Level 4/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$

$P(\text{Building})$

$P(\text{Road})$



# Level 5/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$



$P(\text{Building})$



$P(\text{Road})$



# Level 6/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$



$P(\text{Building})$



$P(\text{Road})$



# Level 7/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$



$P(\text{Building})$



$P(\text{Road})$



# Level 8/8 Predictions

Current Output   Segmentation    $P(\text{Foreground})$



$P(\text{Tree})$

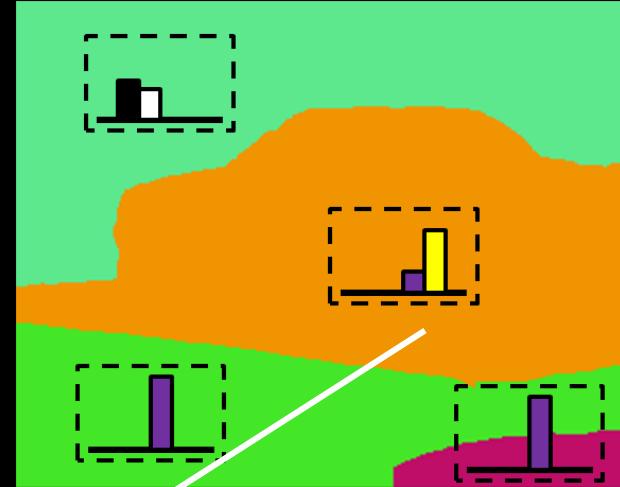


$P(\text{Building})$

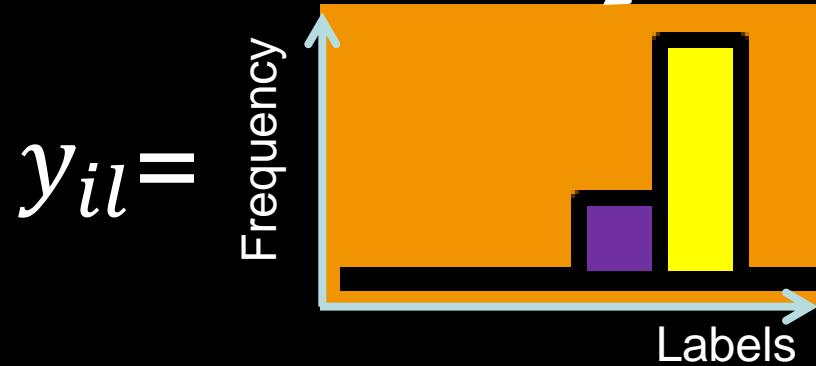


$P(\text{Road})$





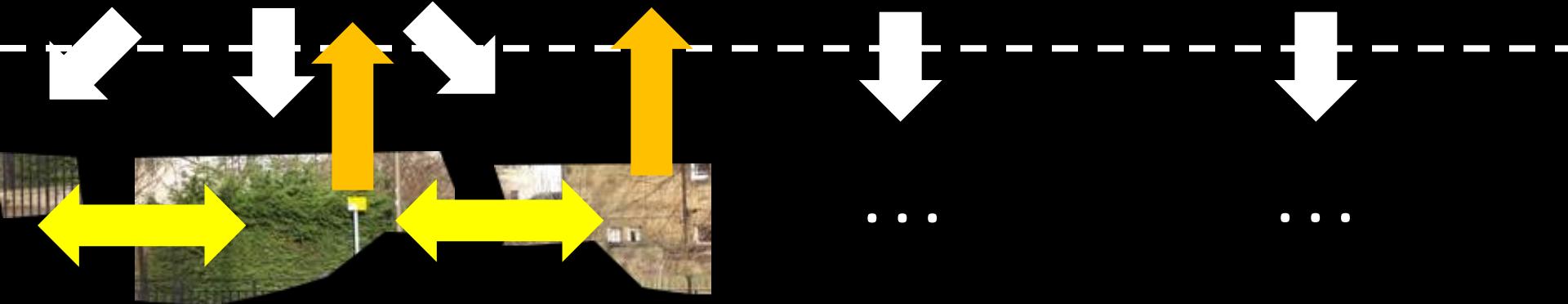
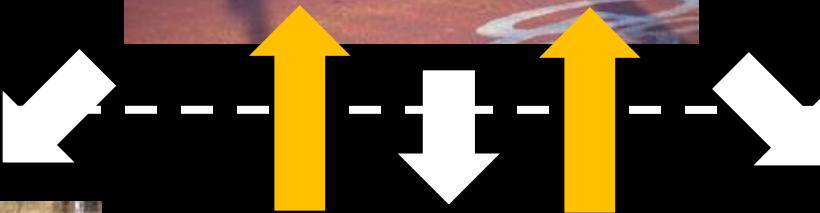
Region  $i$  at level  $l$

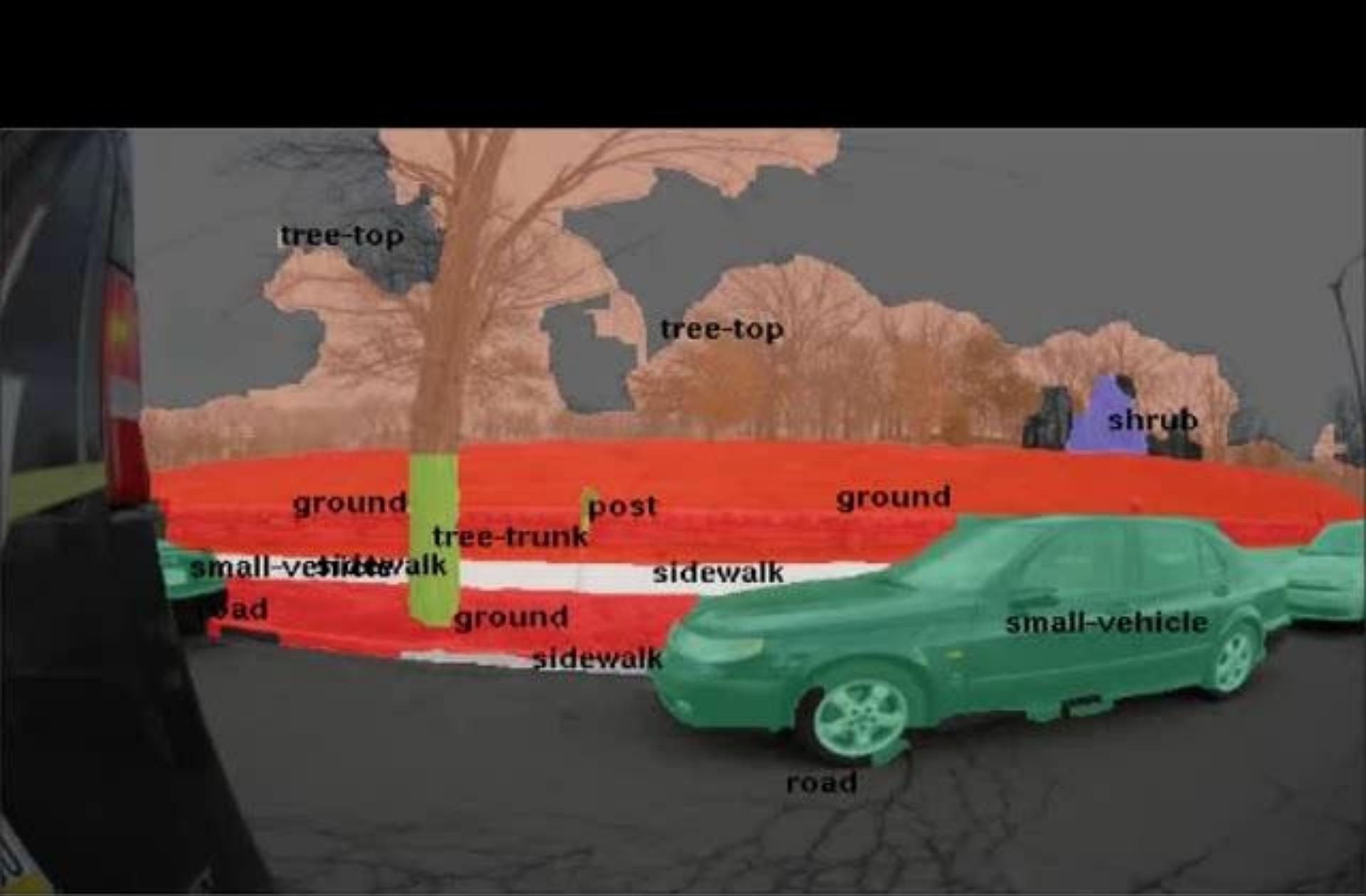


$$y_{il} = q_l(f_{il})$$



$$f_{il} = \begin{bmatrix} \text{Features from region } i \\ \text{Predicted label distribution from parent} \\ \text{Predicted label distributions from neighbors} \end{bmatrix}$$

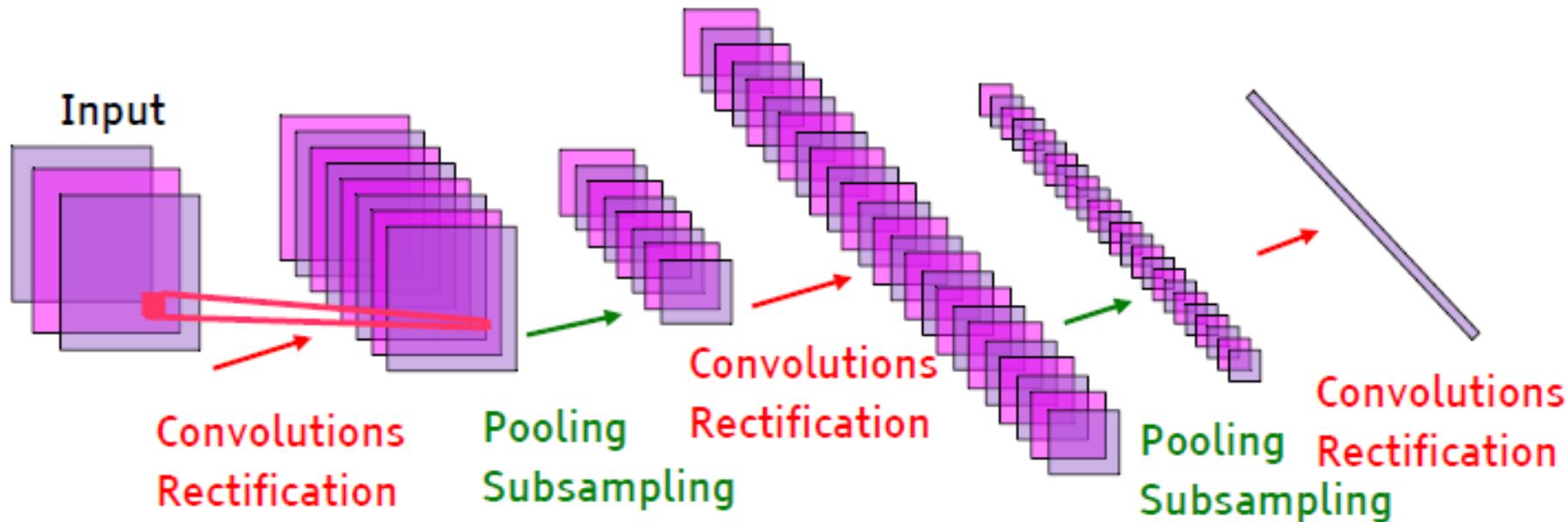




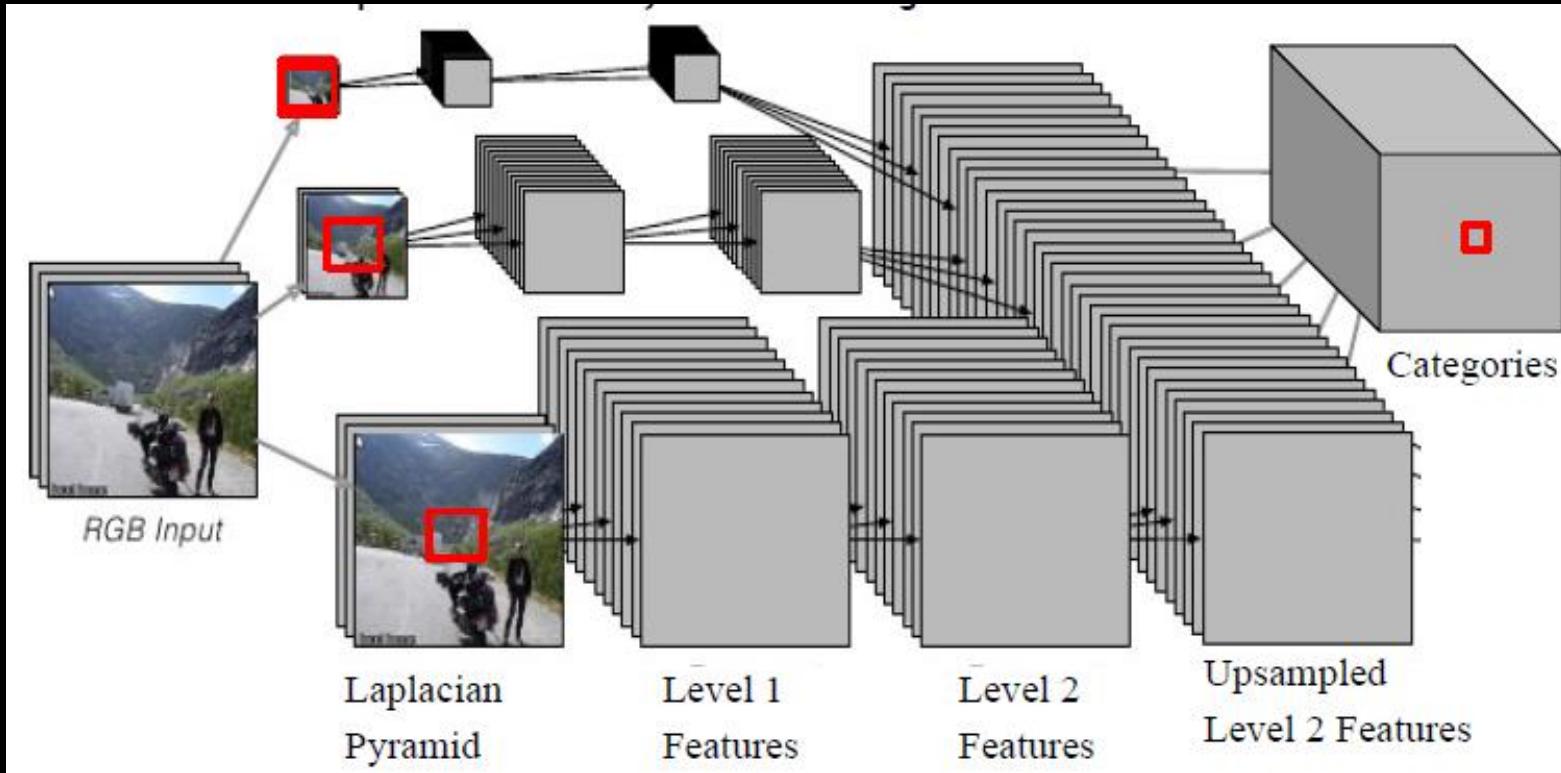
# Approaches

- Set up problem for “exact” solution (e.g., M3N..)
- Approximate solutions
- Sampling
- Decomposition into sequences of simpler problems
- Deep learning

# ConvNets (for detection)



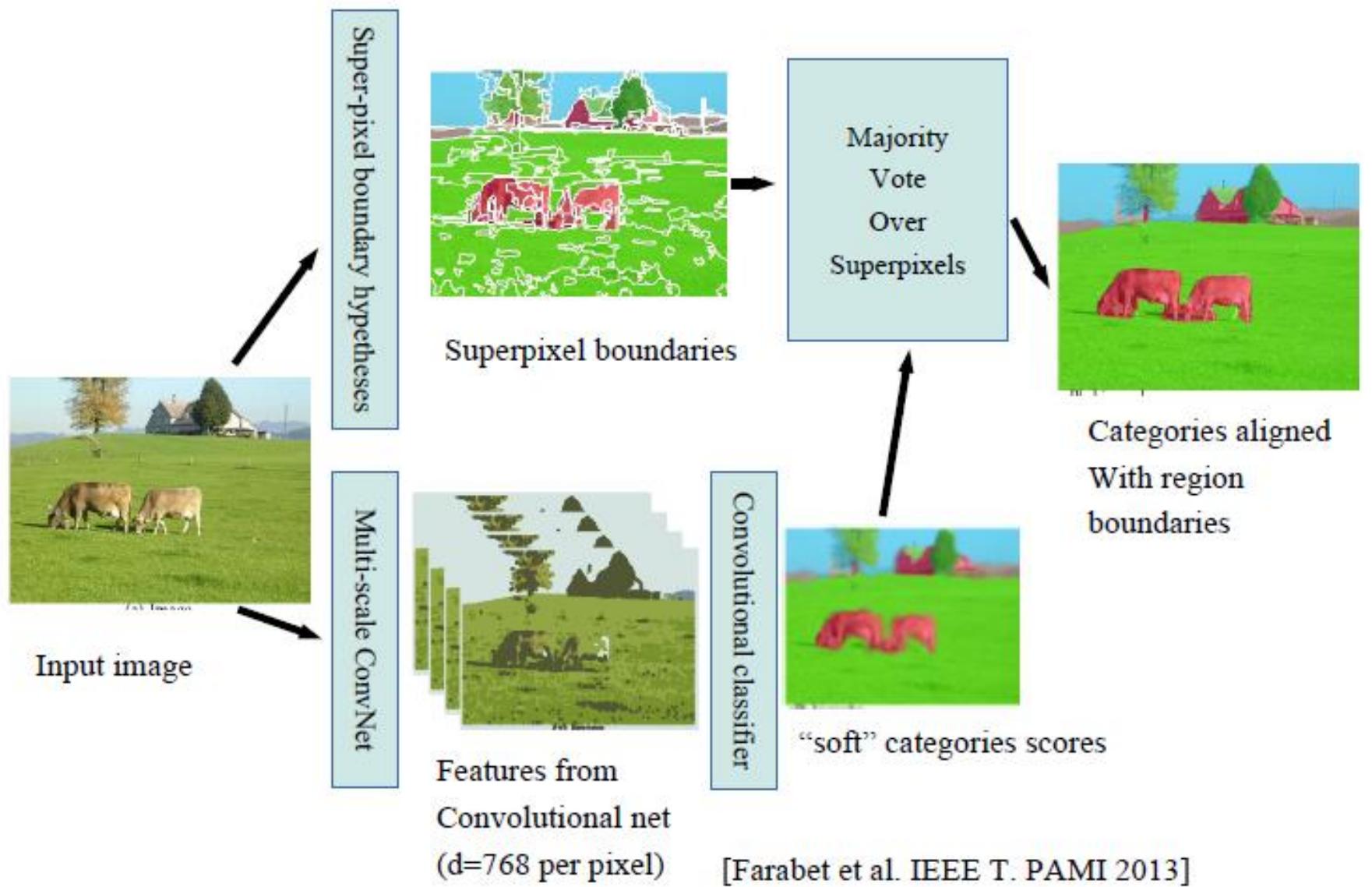
# ConvNets (for semantic labeling)



- Multi-scale features
- How to aggregate over arbitrary support regions?

Slide adapted from Y. Lecun

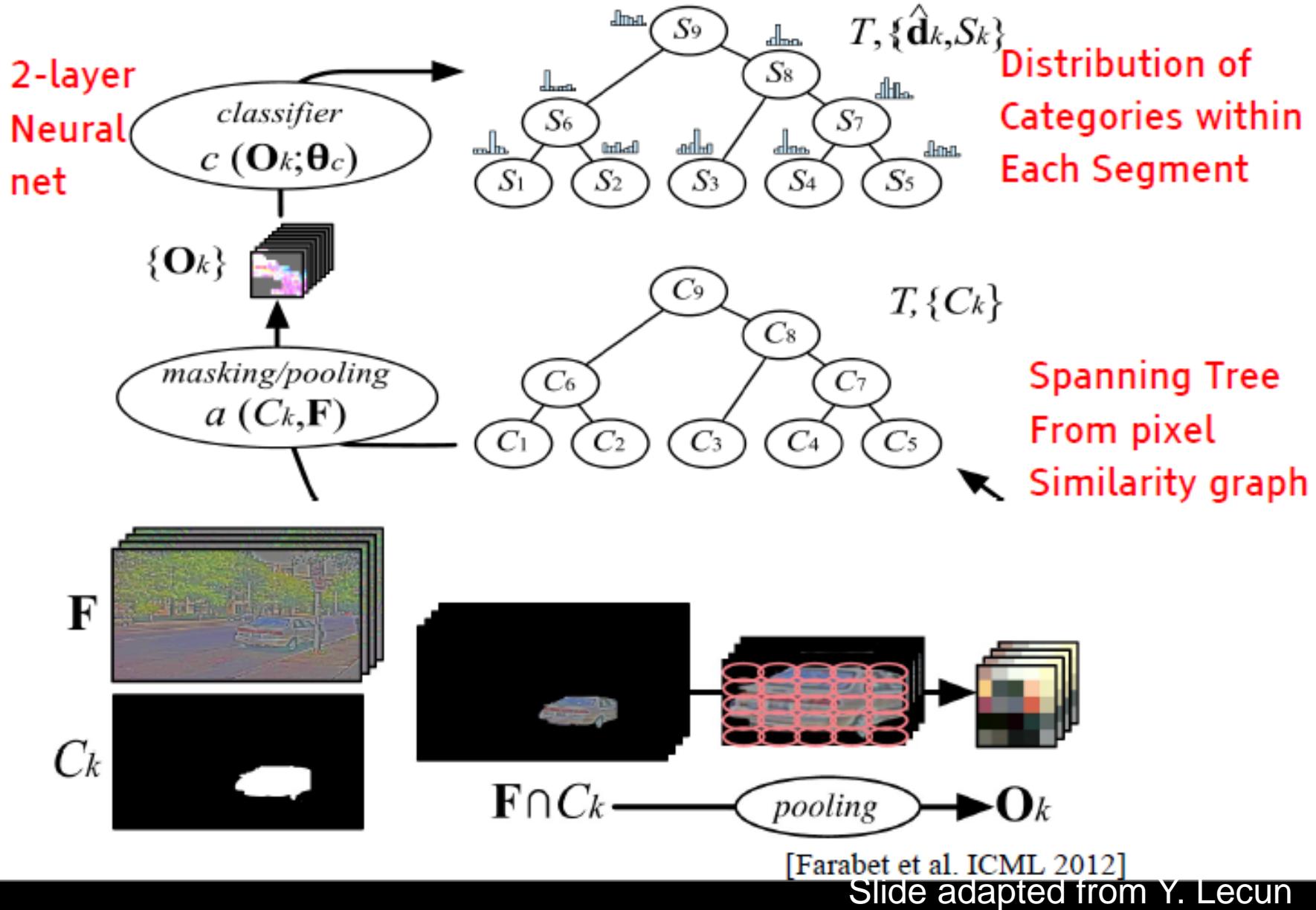
# Method I: Aggregate superpixels



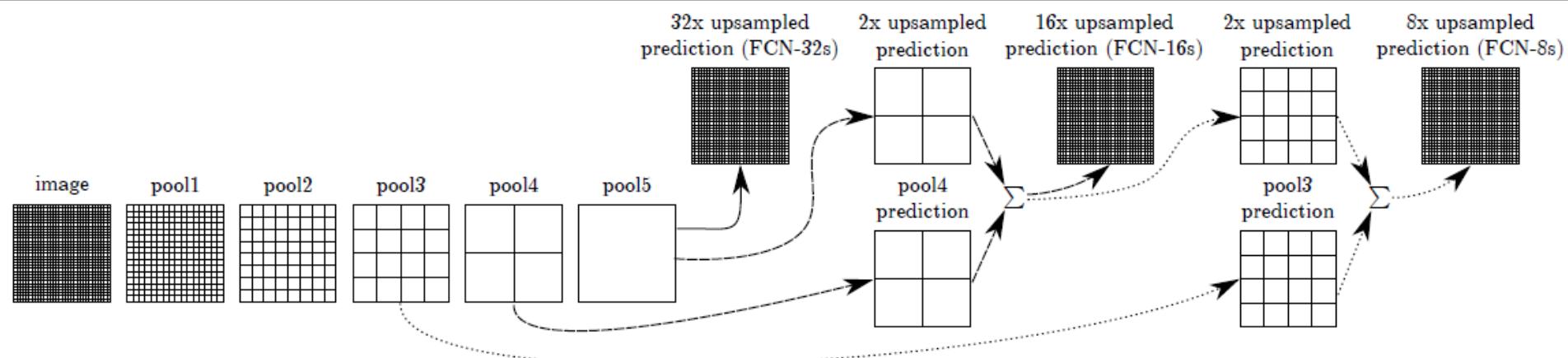
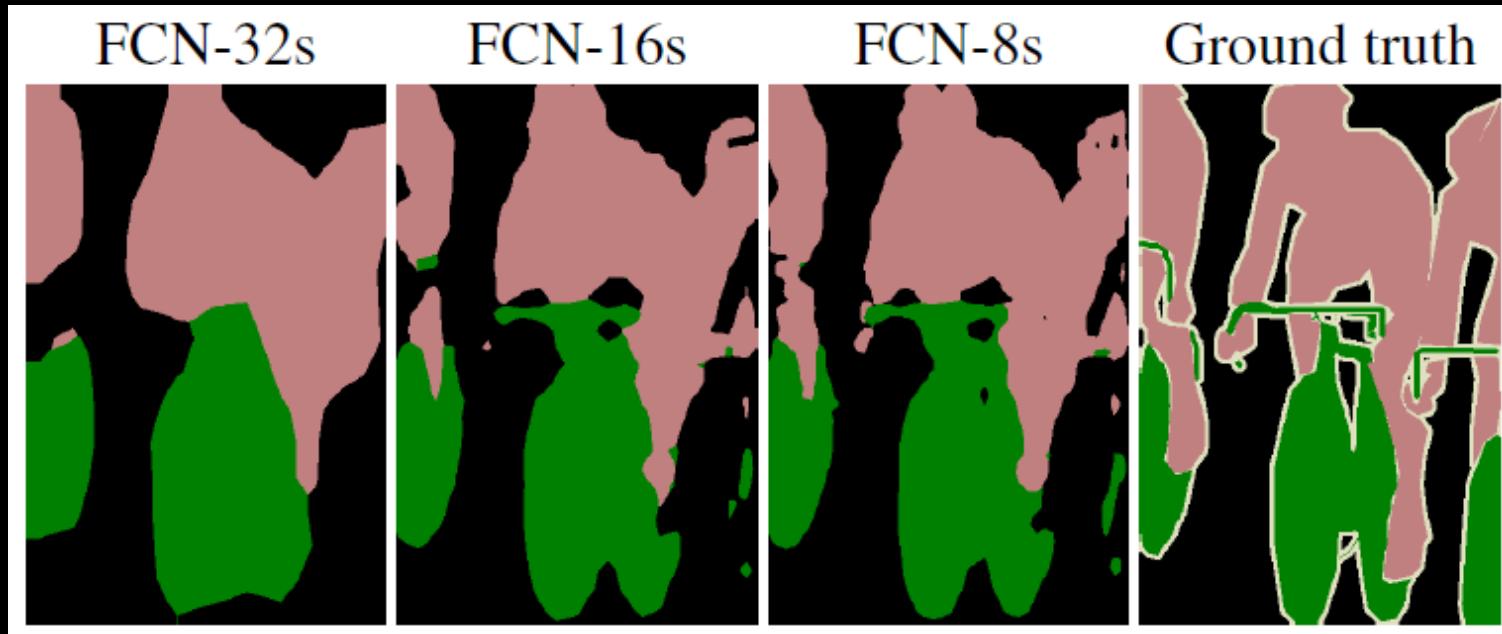
[Farabet et al. IEEE T. PAMI 2013]

Slide adapted from Y. Lecun

# Method II: Tree covering

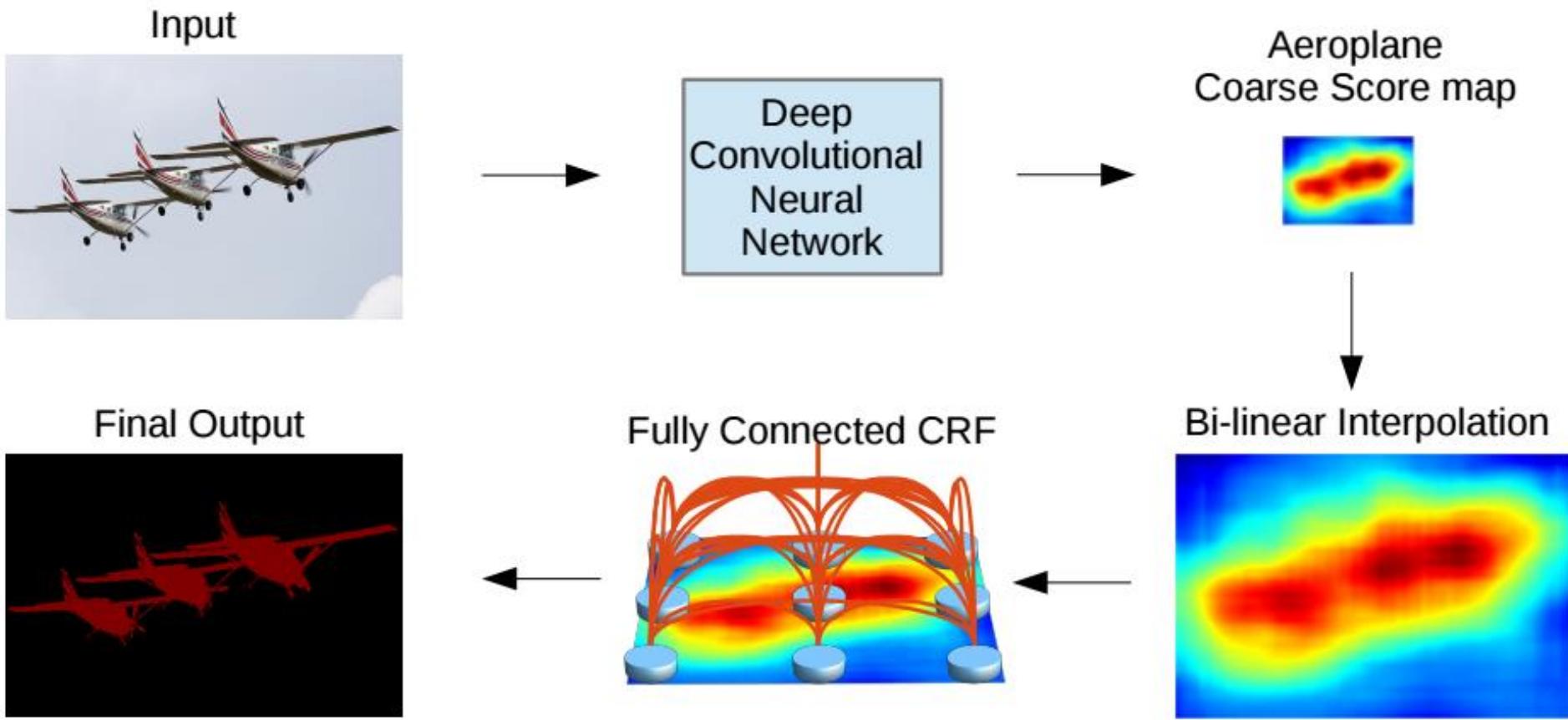


# Method III: (Clever) Sampling



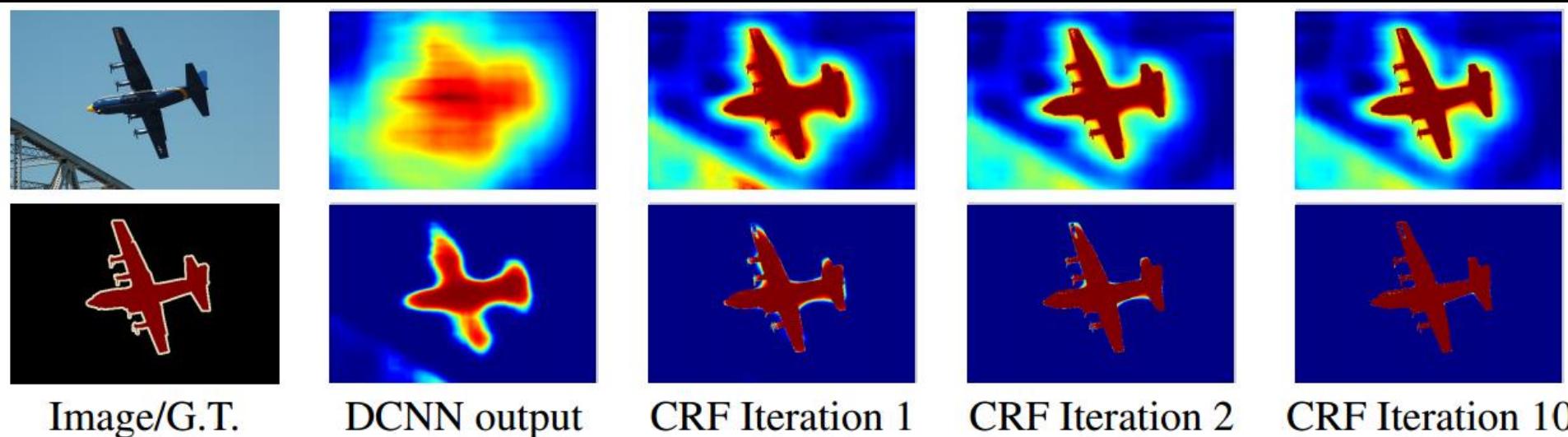
J. Long, E. Shelhamer, T. Darrell. Fully Convolutional Networks for Semantic Segmentation. 2014.

# Method IV: Back to CRFs



L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. Yuille. Semantic Image Segmentation with deep Convolutional Nets and Fully Connected CRFs. ICLR 2015.

# Method IV: Back to CRFs



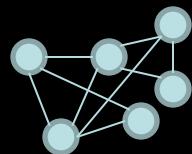
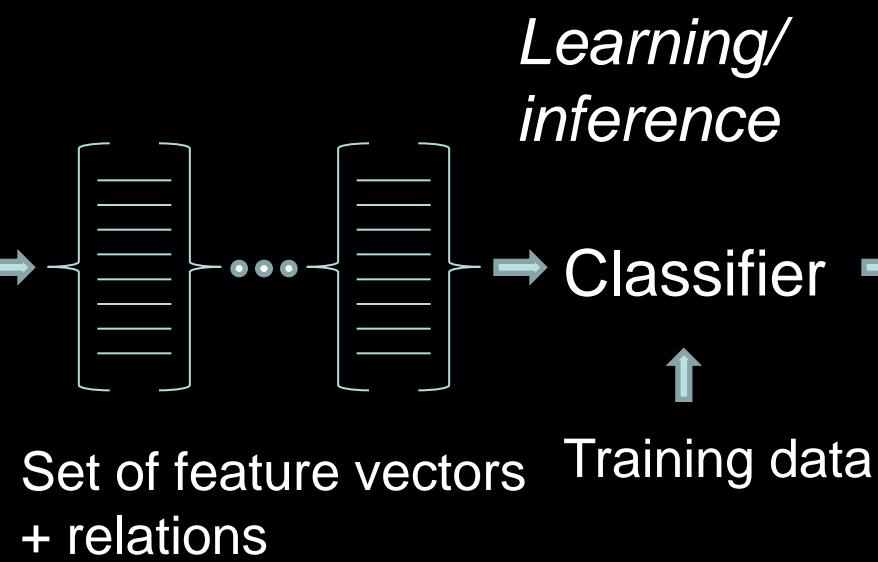
# Method IV: Back to CRFs



L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. Yuille. Semantic Image Segmentation with deep Convolutional Nets and Fully Connected CRFs. ICLR 2015.



Input image

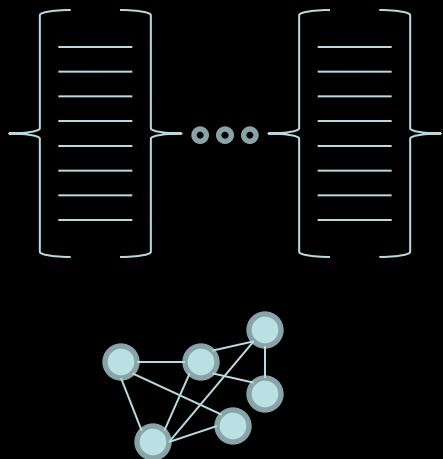


# Outline

- Semantic: Labeling scene regions and objects
- *Geometric: Estimating the geometric structure of the scene*



Input image



Set of feature of  
vectors  
+ additional  
structure (e.g.,  
geometry, relations)

*Reasoning  
Learning/  
inference*

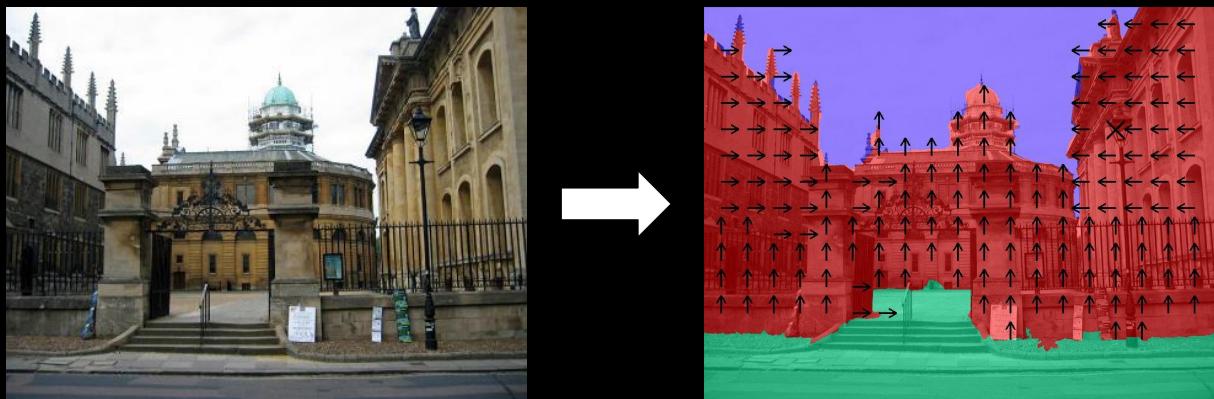


Training data +  
*Geometry, relational  
information, physical  
constraints, domain  
knowledge*

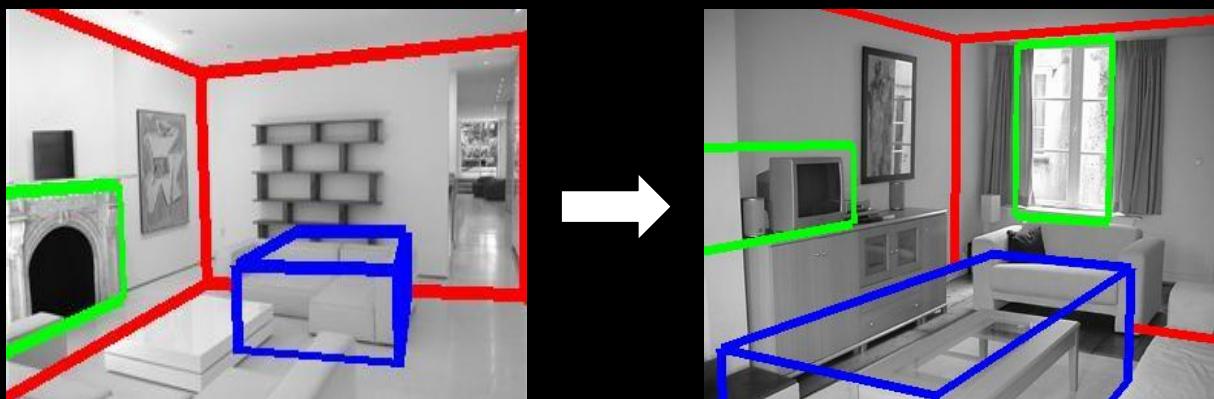


# 3D interpretation from image

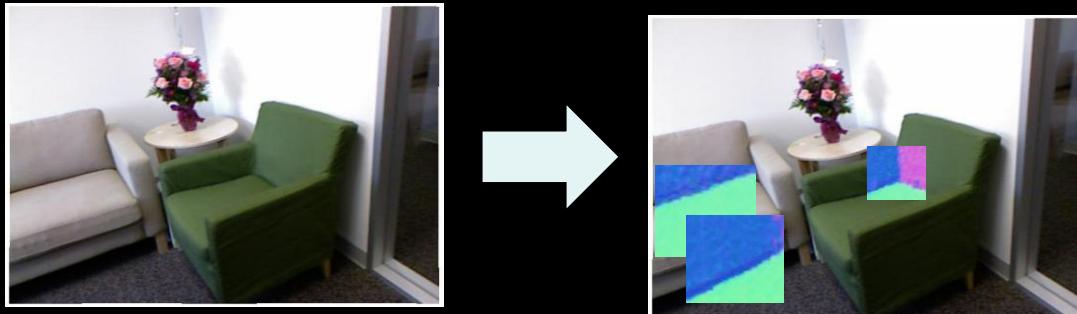
Geometric labels



Volumetric layout



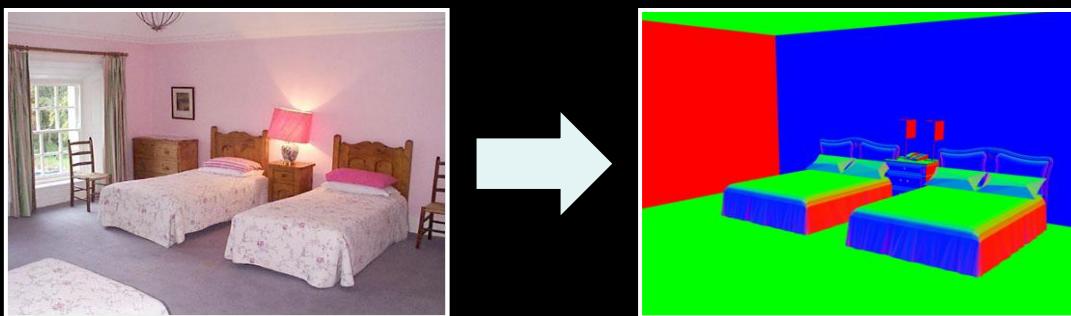
## Sparse primitives



## Dense reconstruction

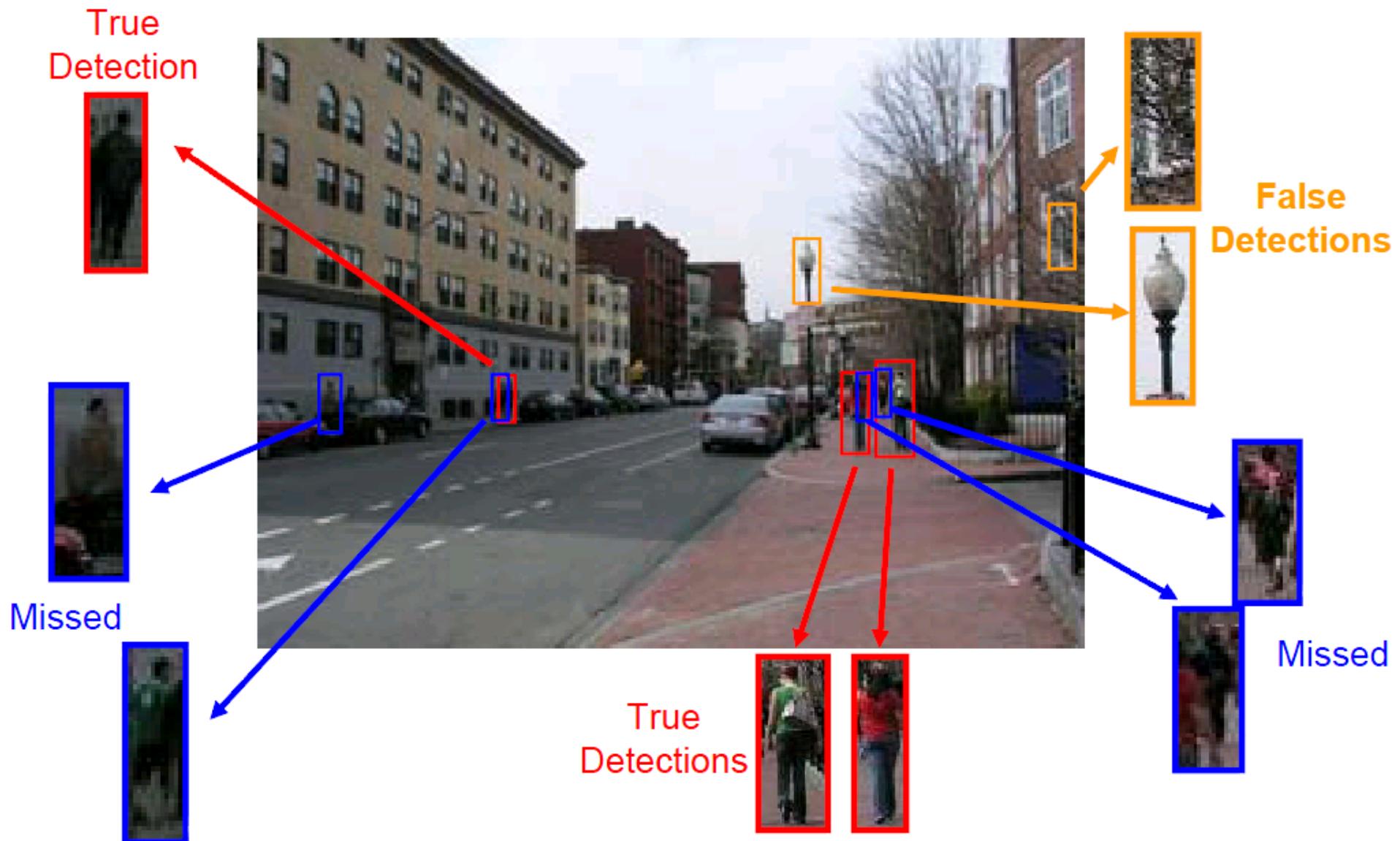


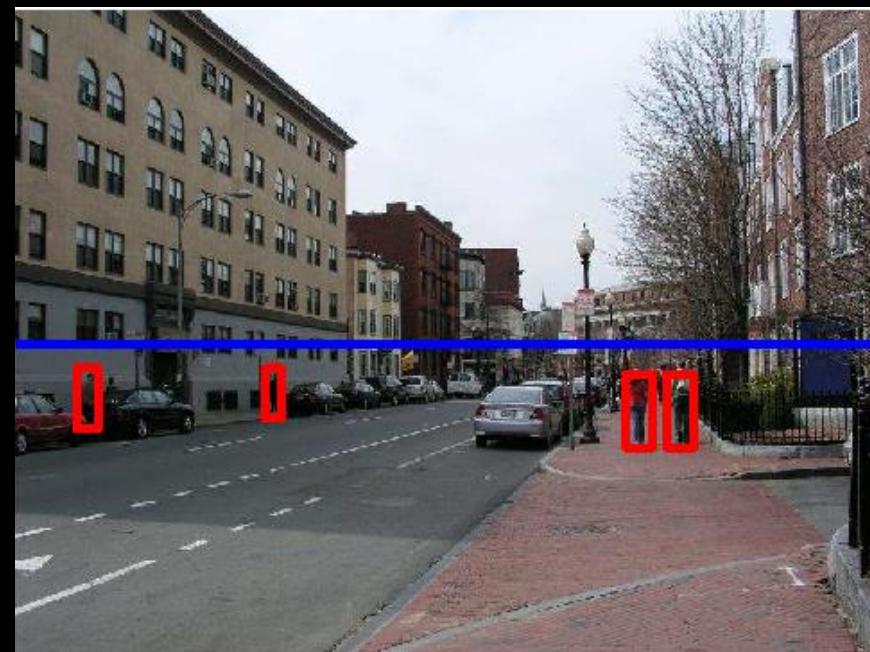
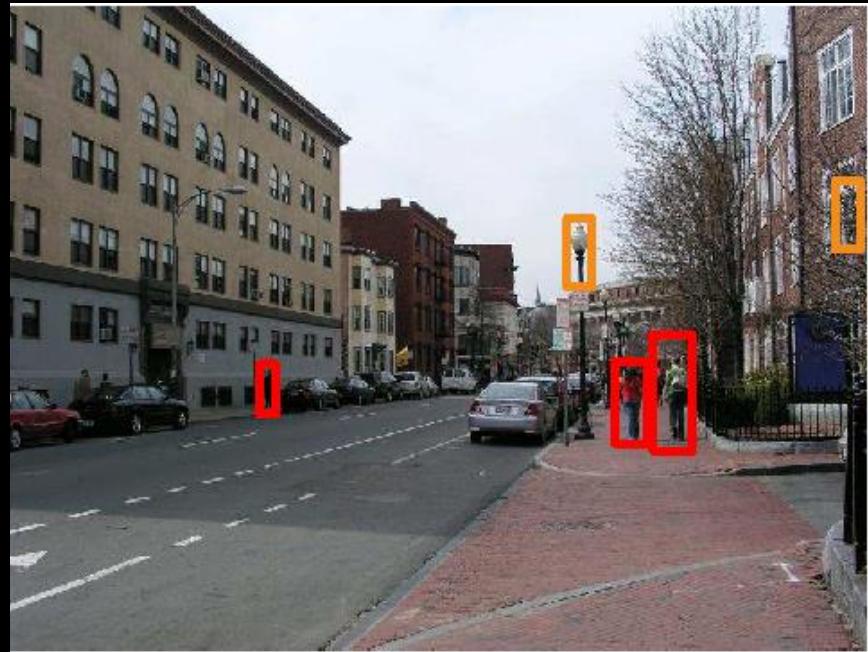
## 3D scene model



# Example applications

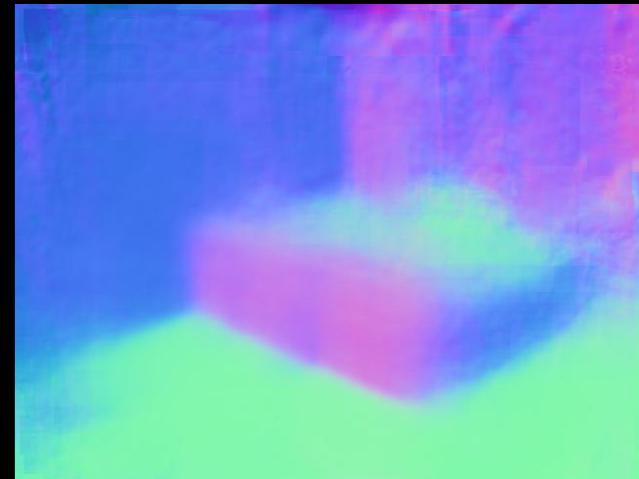
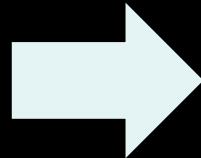
# Informing detectors





D. Hoiem, A. A. Efros, and M. Hebert. *Putting Objects in Perspective*. International Journal of Computer Vision, Vol. 80, No. 1, October, 2008.

# Editing images

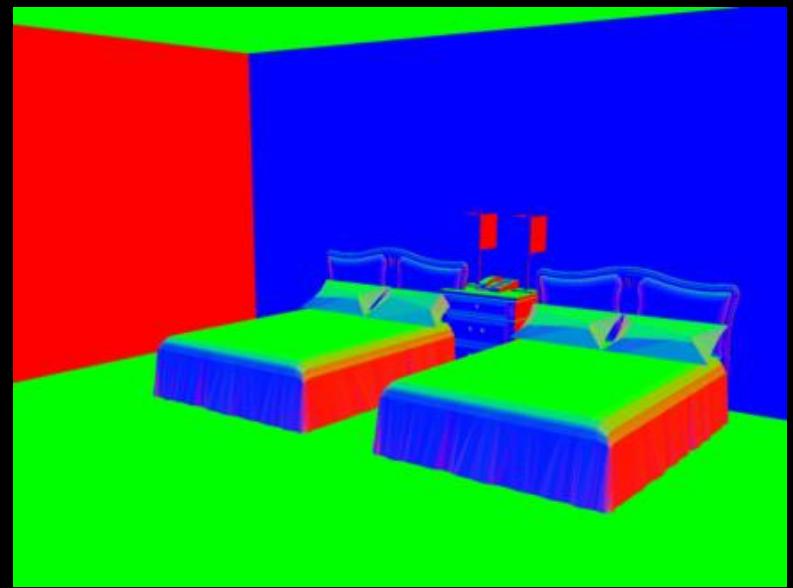
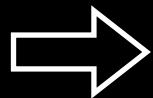




# Predicting actions

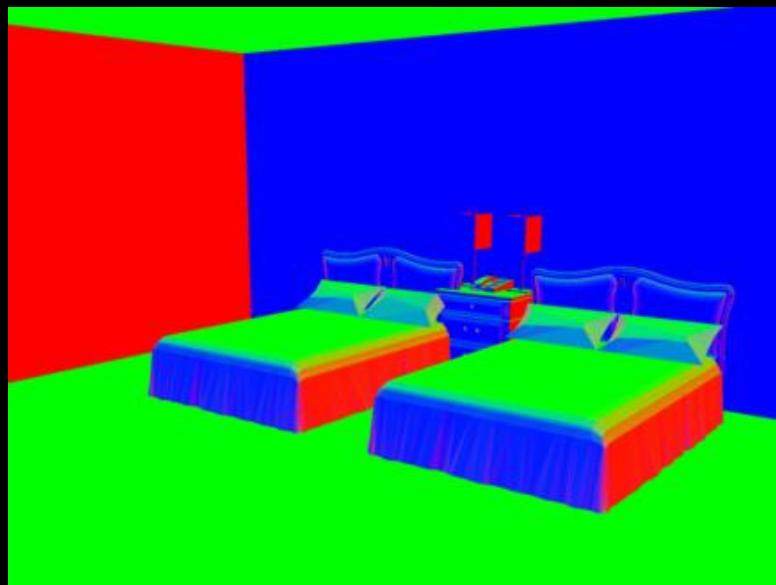


Input Image

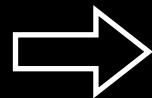


Estimated Geometry

# Application: Affordance Estimation



Estimated Geometry



Predicted Sitting  
Locations

# Application: Affordance Estimation



Sitting Upright



Sitting Reclined

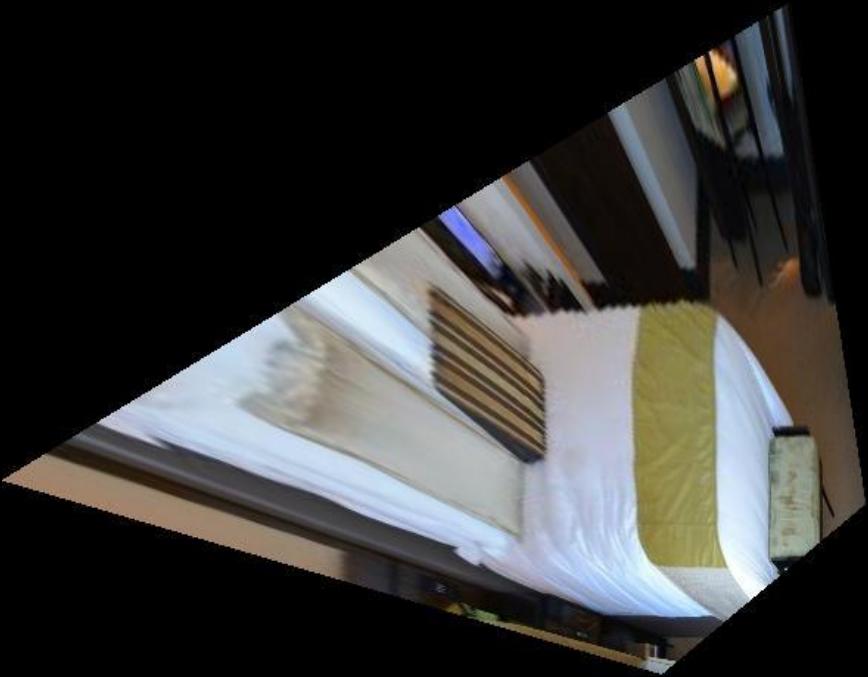


Laying Down



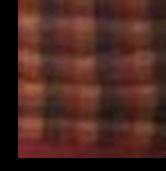
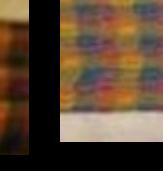
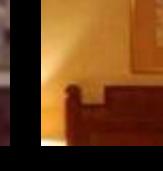
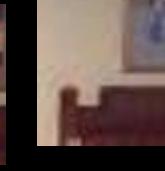
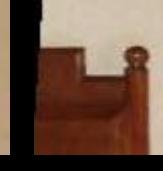
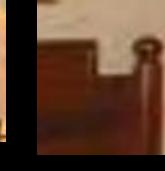
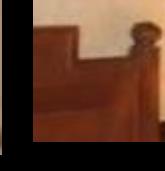
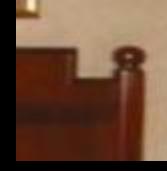
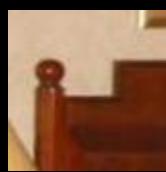
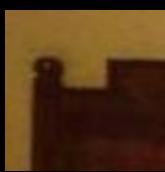
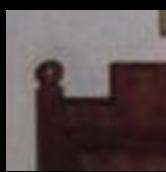
Reaching (4 poses)

# Separating style and structure

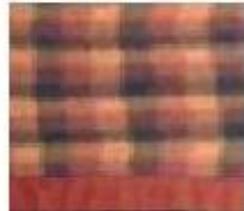


Tenenbaum & Freeman. Separating Style and Content with Bilinear Models. Neural Computation. 2000.

# Casablanca Hotel, New York

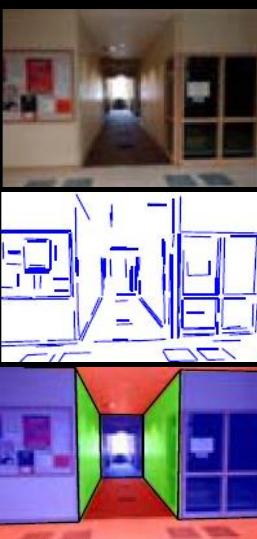
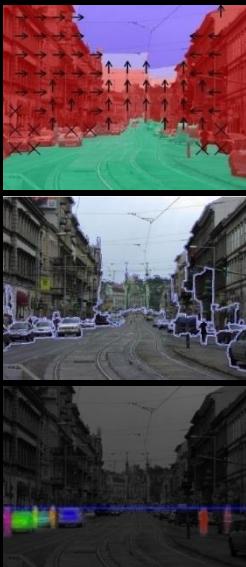






# Outline

## Part 1: Bottom-up Methods for Regions and Boundaries, Global Constraints

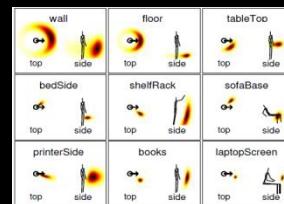
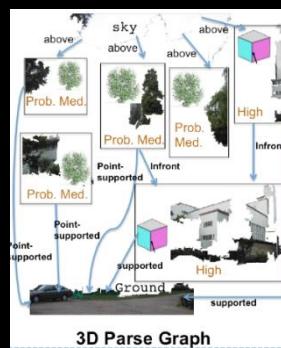


Region labels

+  
Boundaries  
and objects

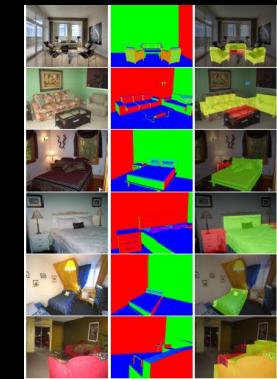
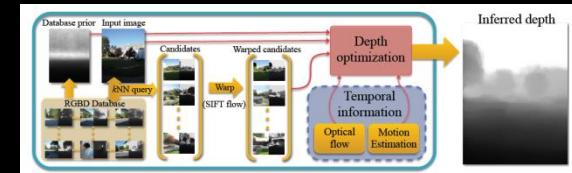
Stronger  
geometric  
constraints  
from domain  
knowledge

## Part 2: Volumetric and Functional Constraints



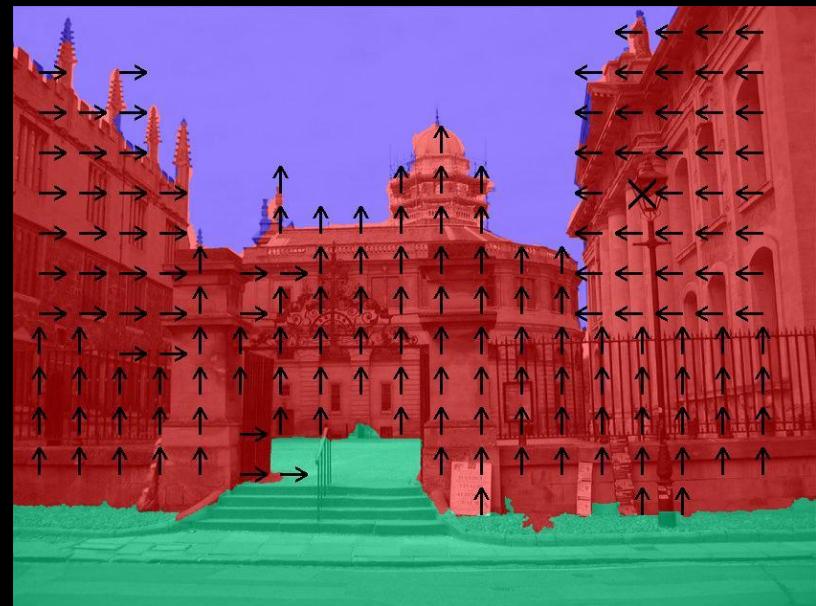
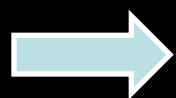
+ physical  
constraints  
+functional  
constraints

## Part 3: Data-driven Models



# Bottom-up estimation

# First attempt: Estimate surface labels



[D. Hoiem, A. A. Efros, and M. Hebert. *Recovering surface layout from an image*. IJCV, 75(1):151–172, 2007]

Multiple  
Classification  
segmentations

Input



...



...



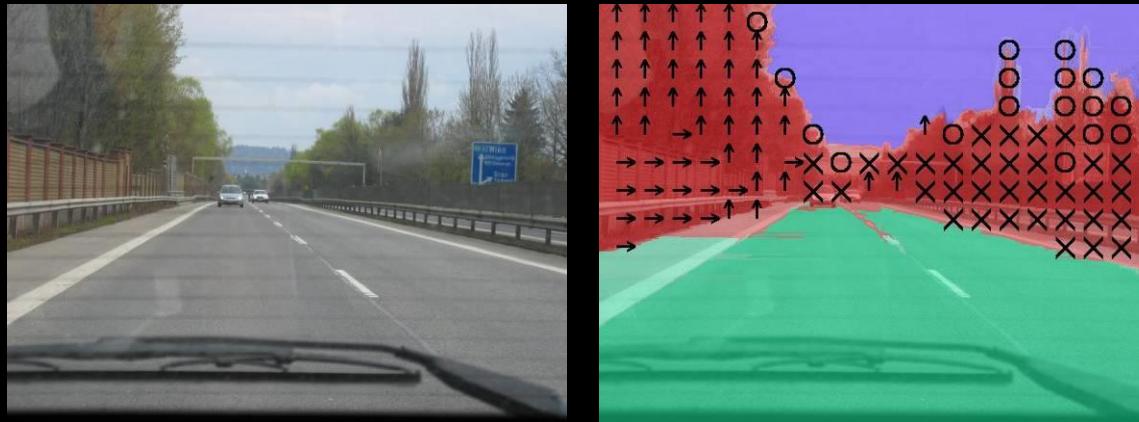
Training data

Classification



[Example from  
Hoiem]

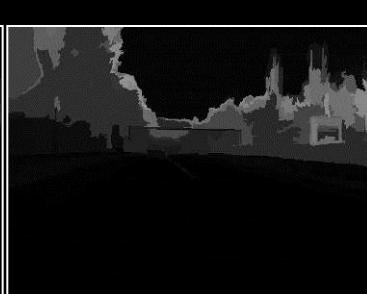
# Example



Support

Vertical

Sky



V-Left

V-Center

V-Right

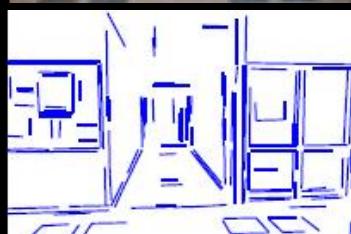
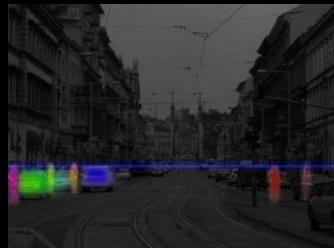
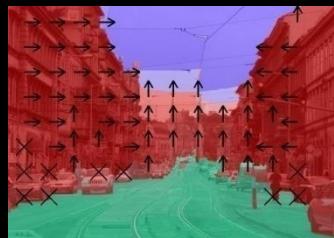
V-Porous

V-Solid

- Is a more precise representation possible?
- For example:
  - We would like to include reasoning about interposition (relations between object relative to a viewpoint induced by occlusion boundaries)
  - We would like to include constraints about object semantics (when known)

# Comments

- Plus:
  - Scene geometry (surface geometry and object relations) estimated from image data
  - Scene geometry used explicitly in scene understanding
- Minus:
  - Still mostly bottom-up classification approach
  - No use of domain constraints or constraints governing the physical world



**Region labels**

**+ Boundaries  
and objects**

**Stronger geometric  
constraints from  
domain knowledge**

**Qualitative**

**More quantitative  
more precise**

# Example

- Using constraints induced by man-made environments in interpreting images
- Examples: Manhattan world, limited vocabulary of object configurations, etc.

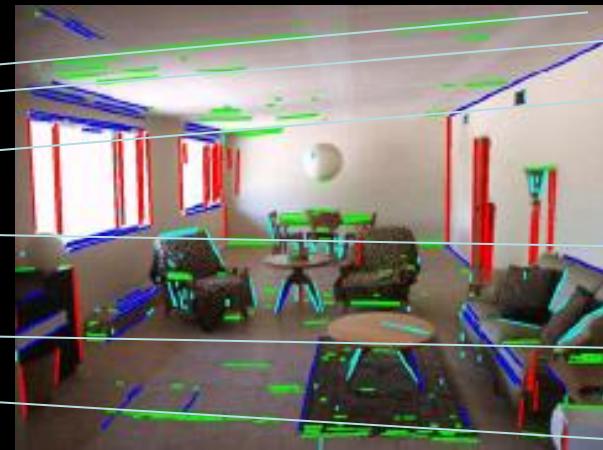


# Constraint: Manhattan world assumption

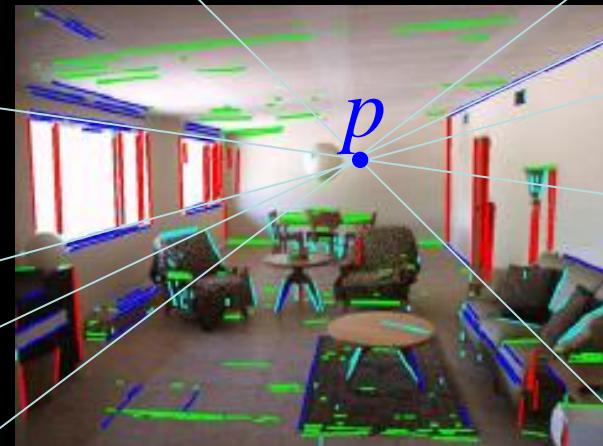
- Three dominant directions corresponding to three “orthogonal” vanishing points



$$n_i = K^{-1}v_i$$



$$n_j \cdot n_i = v_j^T K^{-T} K^{-1} v_i = 0$$



# We need to design 4 things

- Parameterization:  $y$
- Features:  $x$
- Scoring/hypothesis evaluation:  
$$y_o = \operatorname{argmax}_y f(x, y, w)$$
- A way to sample, or generate hypotheses  $y$

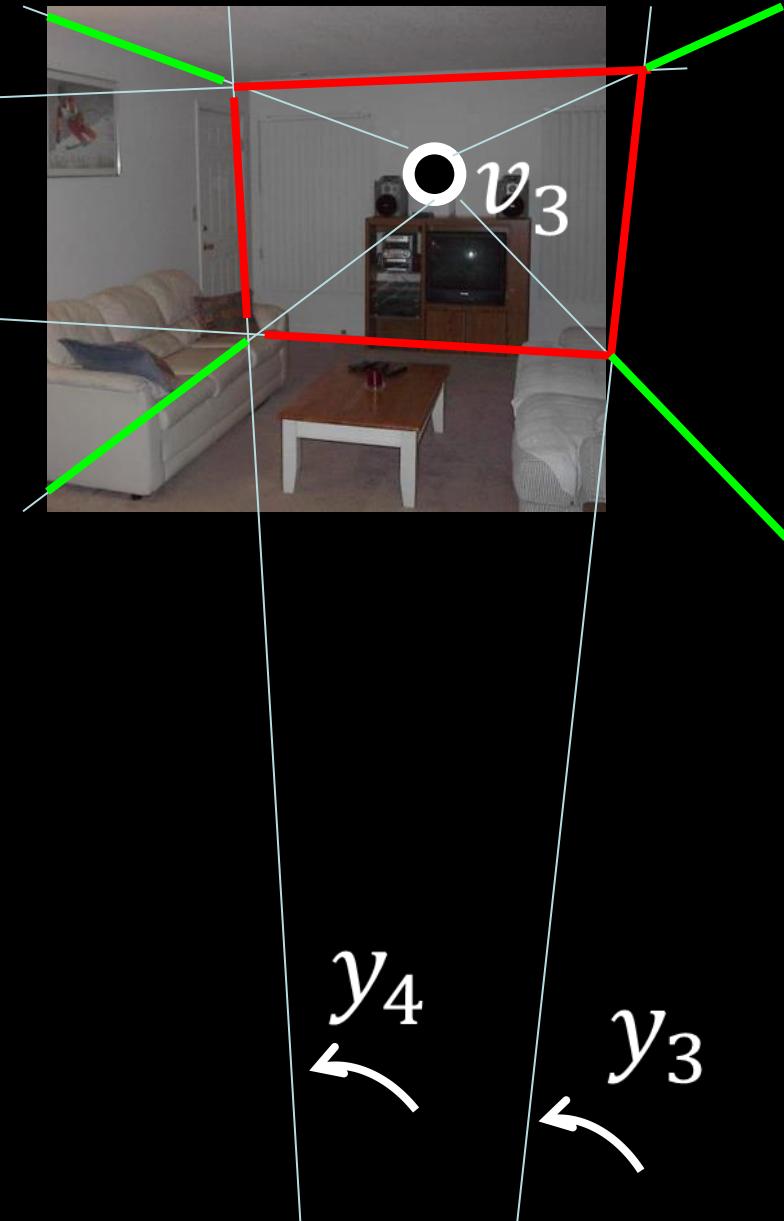
$\circ$

$v_1$

$y_1$

$y_2$

$$y = [y_1 \ y_2 \ y_3 \ y_4]$$



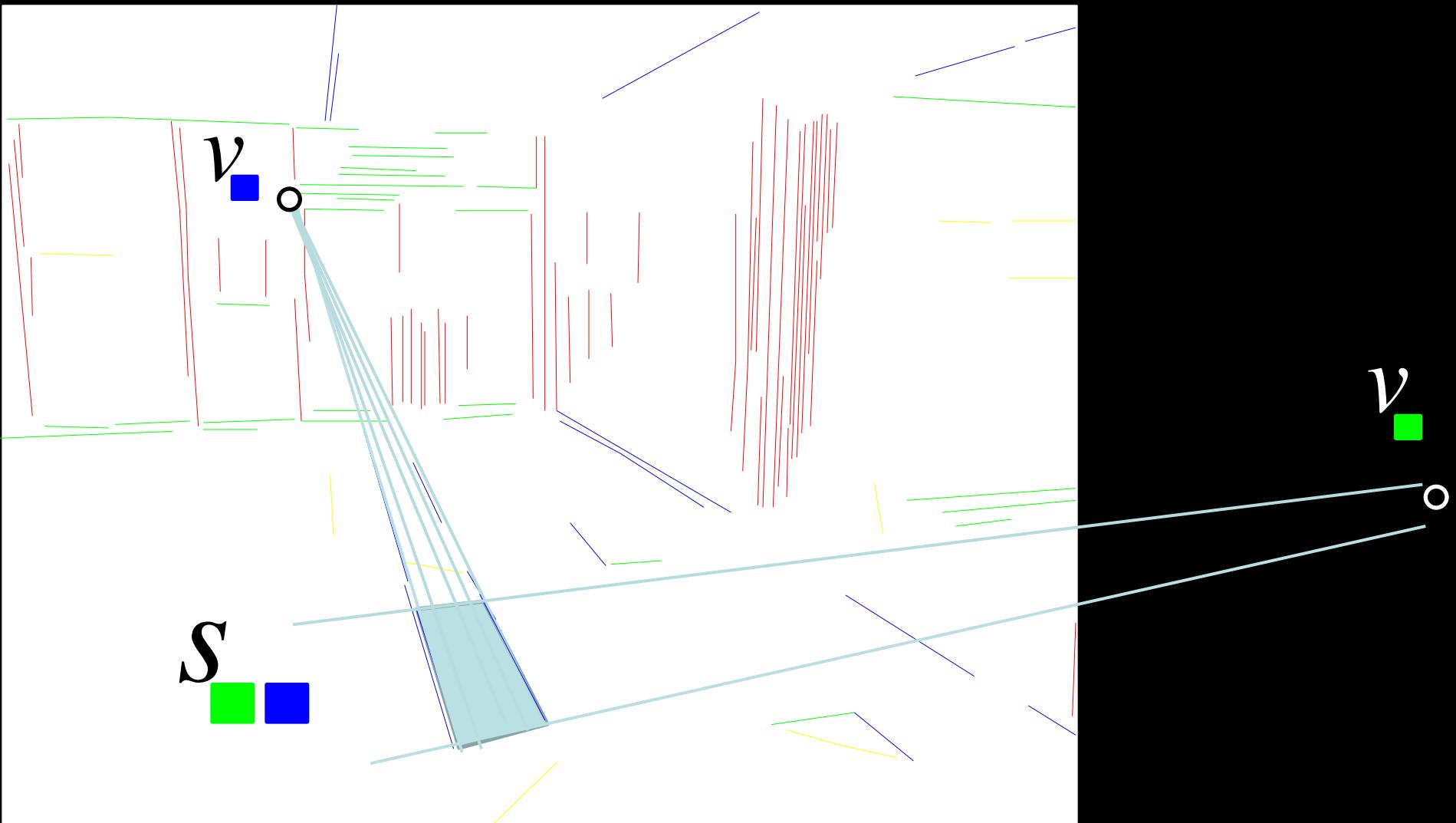
H. Wang, S. Gould, D. Koller. Discriminative learning with latent variables for cluttered indoor scene understanding. ECCV 2010.

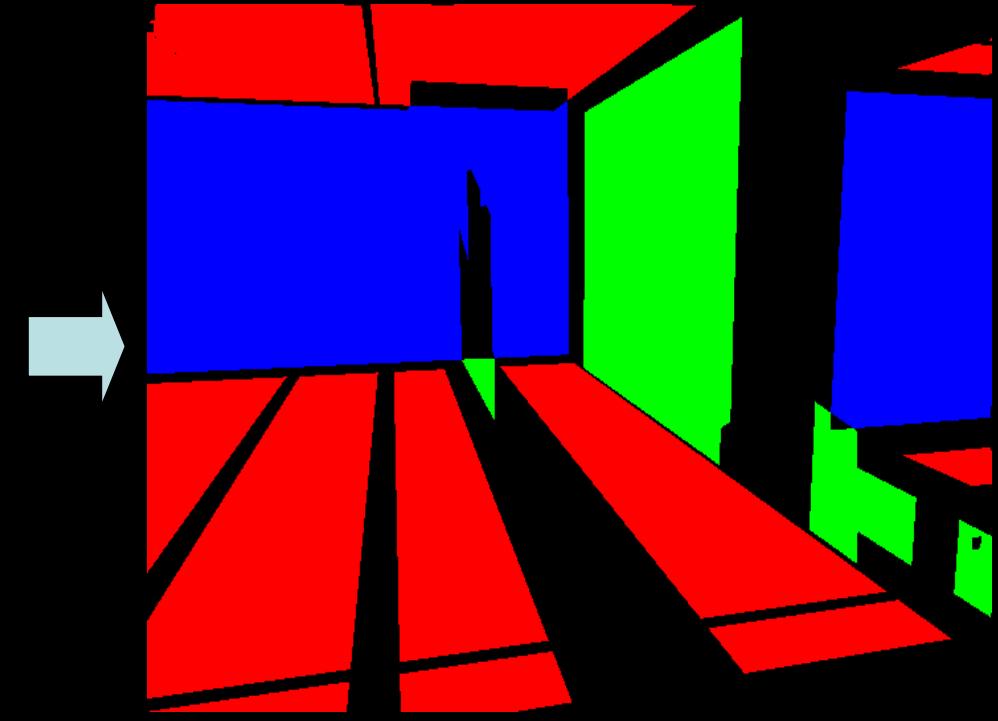
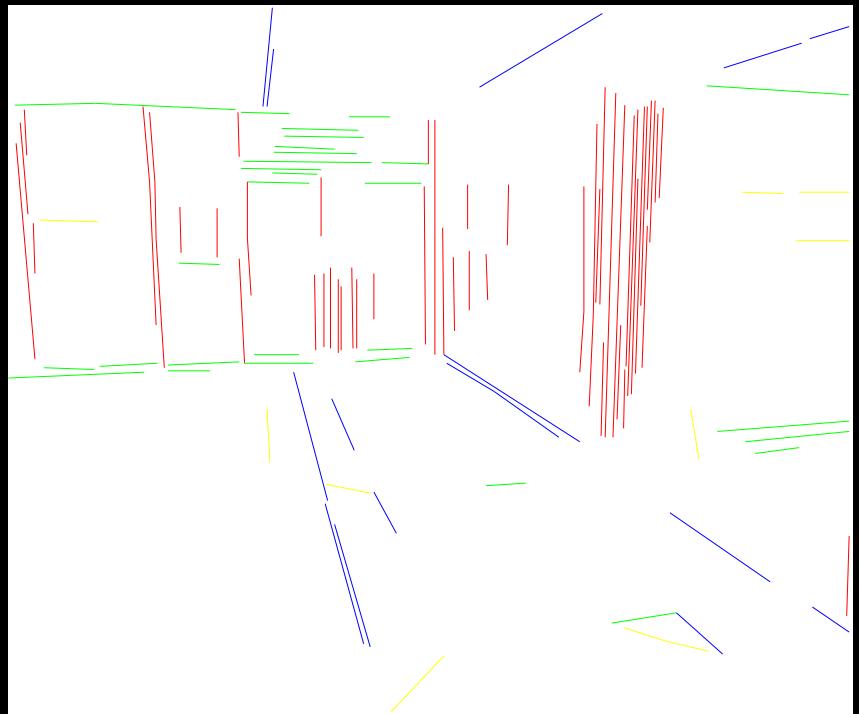
$\circ \ v_2$

# Features

- Features:  $x$ 
  - Surface layout (see earlier)
  - Lines, regions, ...
  - Orientation maps
  - Junctions

# Orientation maps: Sweep algorithm





$p$  is of orientation ■ if  
it is in one  $S_{\blacksquare\blacksquare}$   
It is in one  $S_{\blacksquare\blacksquare}$   
It is in none of  $S_{\blacksquare\blacksquare}$   $S_{\blacksquare\blacksquare}$

# Junctions

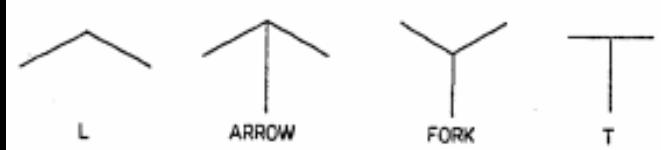


Figure 8 Junction types treated in this paper.

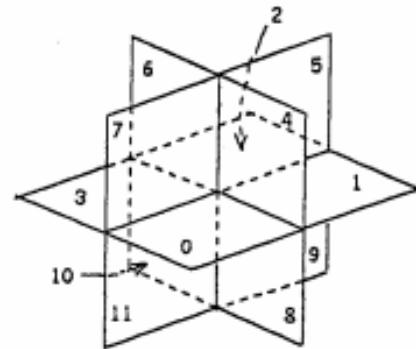
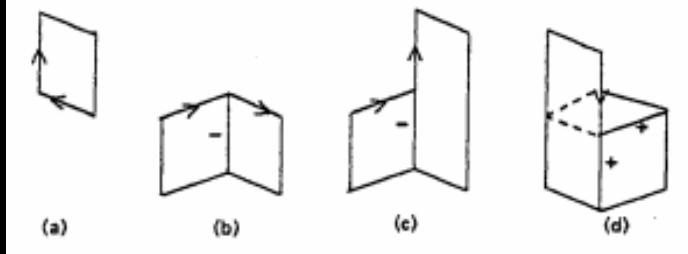
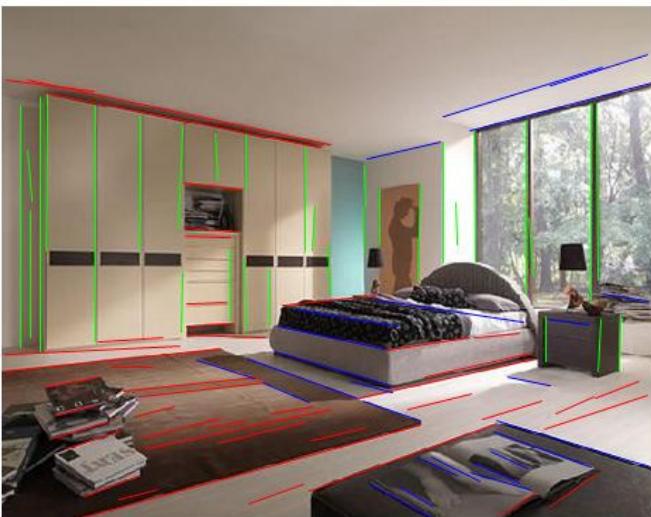


Figure 9 Twelve quadrant planes.

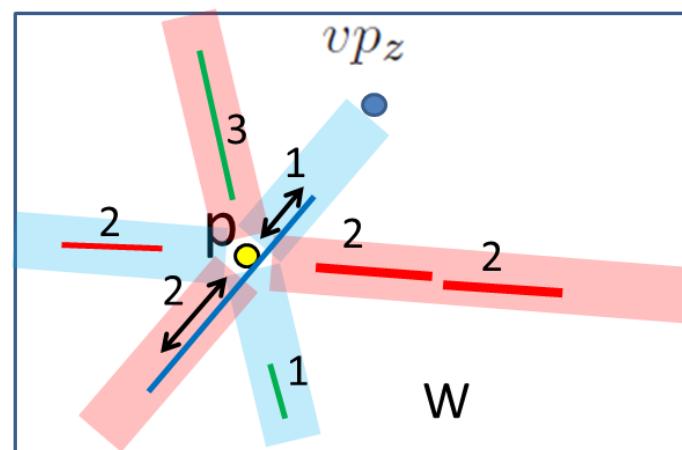


S. Ramalingam, J.K. Pillai, A. Jain, Y. Taguchi.  
Manhattan Junction Catalog for Spatial Reasoning of  
Indoor Scenes. CVPR 2013.

- M. B. Clowes. On seeing things. *AI*, 1971.  
T. Kanade. A theory of origami world. *AI*, 1980.  
D. A. Huffman. Impossible objects as nonsense  
sentences. *Machine Intelligence*, 1971.



vp<sub>x</sub>



vp<sub>y</sub>

L

T

X

Y

W



# Scoring the hypotheses

- Structured prediction

$$y_o = \operatorname{argmax}_y w^T \varphi(x, y)$$

V. Hedau, D. Hoiem, D. Forsyth, “Recovering the Spatial Layout of Cluttered Rooms,” International Conference on Computer Vision (ICCV), 2009.

A.G. Shwing, T. Hazan, M. Pollefeys, R. Urtasun, “Efficient Structured Prediction for 3D Indoor Scene Understanding,” Computer Vision and Pattern Recognition (CVPR), 2012.

Vanishing  
points  
in input  
image

$x$



Score =  
output of  
structured  
predictor

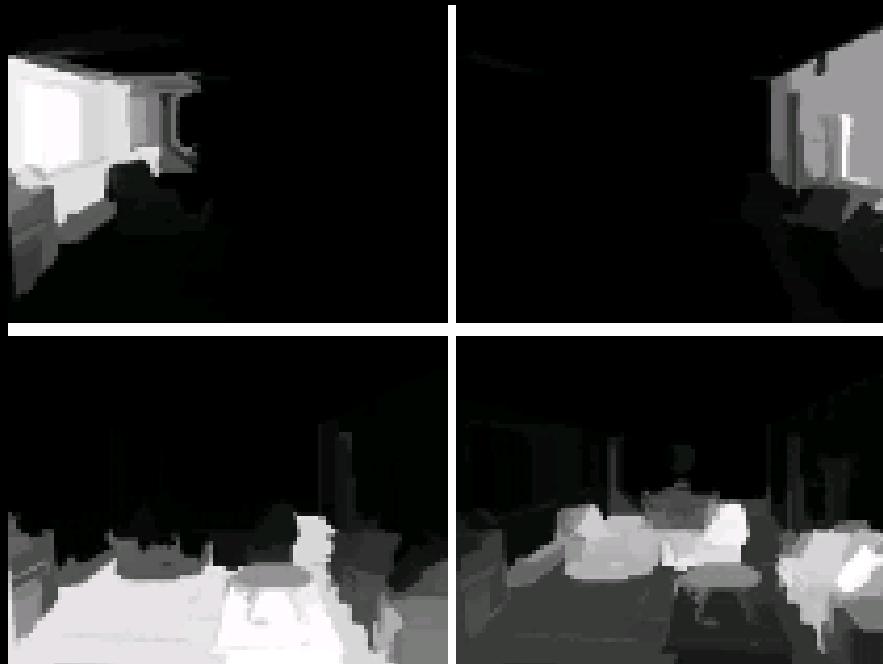
$$\rightarrow f(x, y, w)$$

Many initial  
layout  
hypotheses

$y$



# Data for evaluating $y$



Surface labels	Floor	Left	Middle	Right	Ceiling	Objects
Floor	74/68	0/0	0/1	0/1	0/0	24/30
Left	1/0	75/43	14/44	0/0	1/1	9/12
Middle	1/0	5/2	76/82	4/6	2/1	13/9
Right	1/1	0/0	14/48	73/42	3/2	10/7
Ceiling	0/0	4/3	28/47	2/5	66/45	0/0
Objects	16/12	1/1	5/10	1/2	0/0	76/76

Label confidences  
from classifier

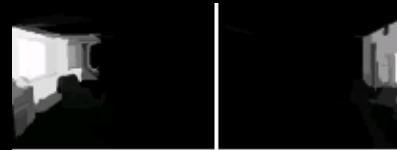
# Definition of mapping function

$$f(x, y, w) = w^T \varphi(x, y)$$

Learned  
weight vector

Feature vector  
measuring  
agreement  
between lines,  
faces and labels

$\varphi(x, y)$ = Relative sum of  
lengths of line segments in  
each face agreeing with  
labels from appearance-  
based classifiers



[Example from Hedau et al.]

# Learning

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

s.t.  $\xi_i \geq 0 \quad \forall i$ , and

$$w^T \psi(x_i, y_i) - w^T \psi(x_i, y) \geq \Delta(y_i, y) - \xi_i,$$

Loss function:

- Distance between centroids of true and estimated faces
- Overlap between true and estimated faces
- Number of missing faces



True:  $y_i$



Estimated:  $y$

$$\Delta(y_i, y)$$

[Example from Hedau et al.]

Lines



Labels



Layout



Lines

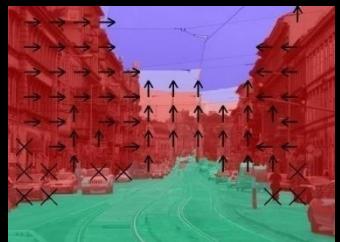


Labels



Layout

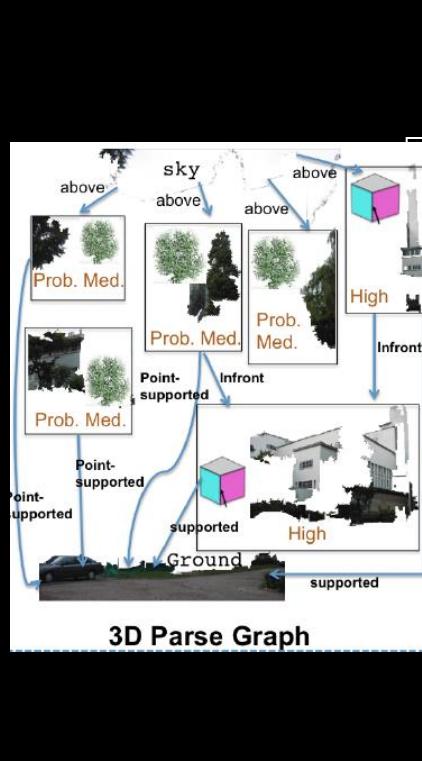
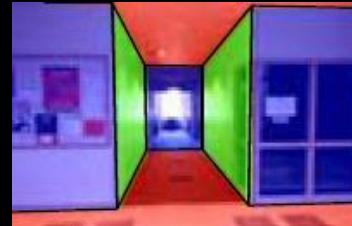
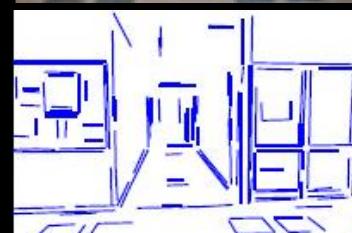




Qualitative

+ Boundaries  
and objects

Stronger geometric  
constraints from  
domain knowledge



+ more constraints

More quantitative  
more precise

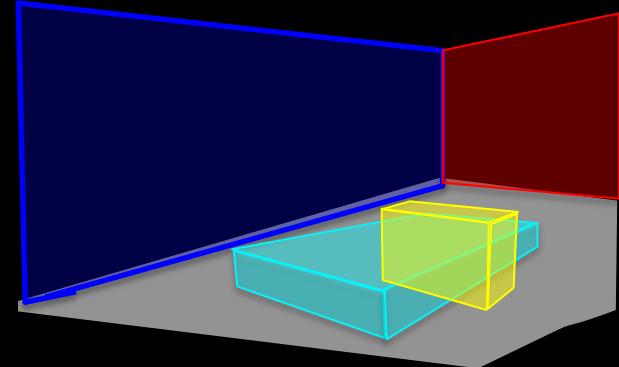
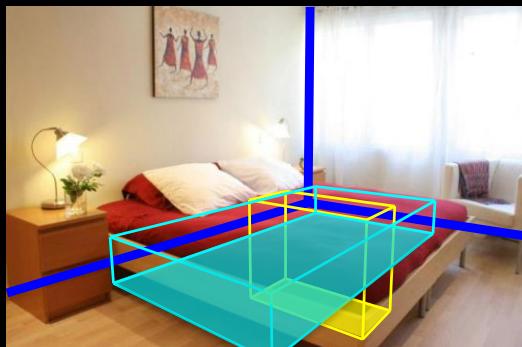
Explicit

# Integrating more constraints

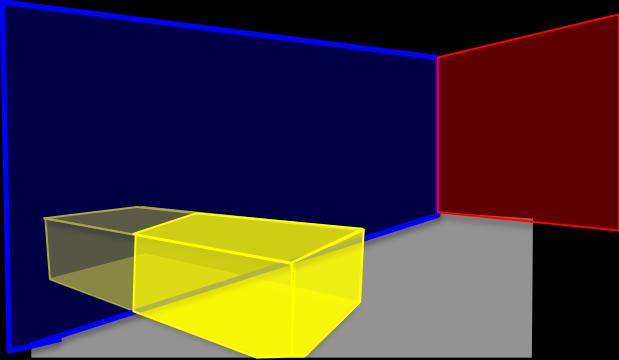
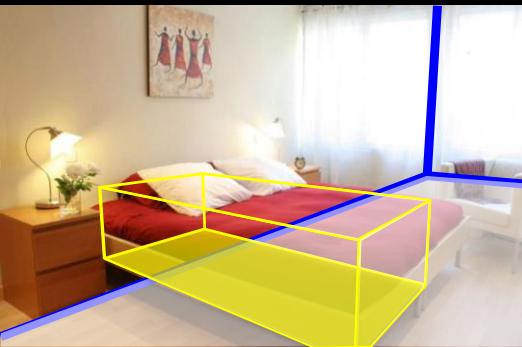
- *Constraints*
  - Volumetric constraints
- *Techniques*
  - Structured prediction

# Constraints: Solid objects must satisfy volumetric/physical constraints

- Finite volume



- Spatial exclusion



- Containment

D. Lee, A. Gupta, M. Hebert, and T. Kanade. Estimating Spatial Layout of Rooms using Volumetric Reasoning about Objects and Surfaces. Advances in Neural Information Processing Systems (NIPS), Vol. 24, 2011.

$$f(x, y)$$

Image

$$y = [r_1 \dots r_n \ o_1 \dots o_m]$$

Labeling: Indicator vector of scene configurations + object hypothesis

Penalty term for incompatible configurations

Compatibility of image data with geometric configuration

$$f(x, y) = \underbrace{w^T \psi(x, y)} + \underbrace{w_\varphi^T \varphi(y)}$$

Features from image (surface labels, vanishing points, etc.)

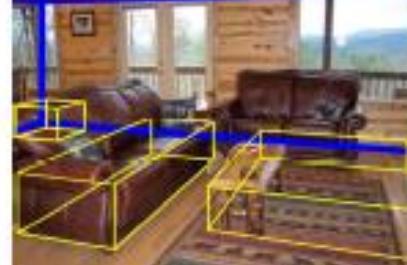
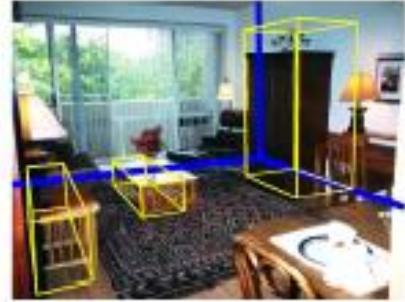
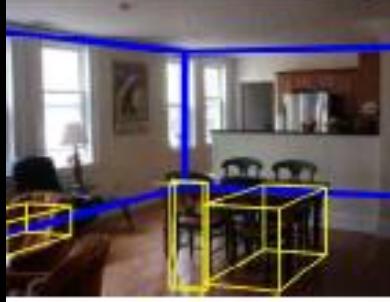
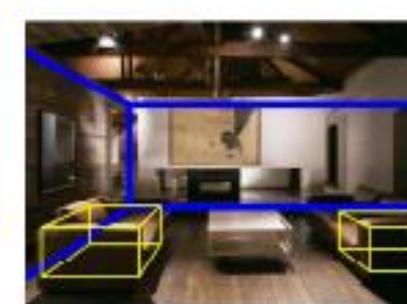
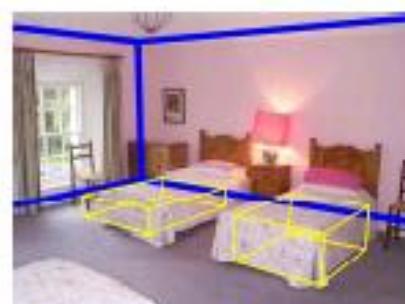
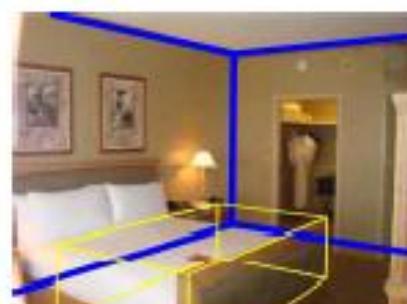
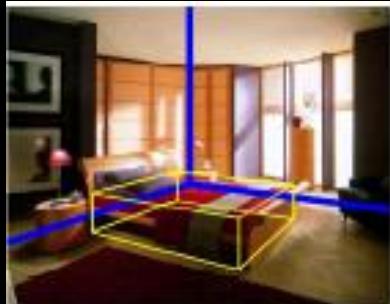
Features of the scene configuration to evaluate constraint violations

$$f(x, y) = w^T \psi(x, y) + w_\varphi^T \varphi(y)$$

- Inference:

$$y^* = \arg \max_y f(x, y)$$

- Training
  - Use structured SVM to estimate  $w$



UIUC 10%

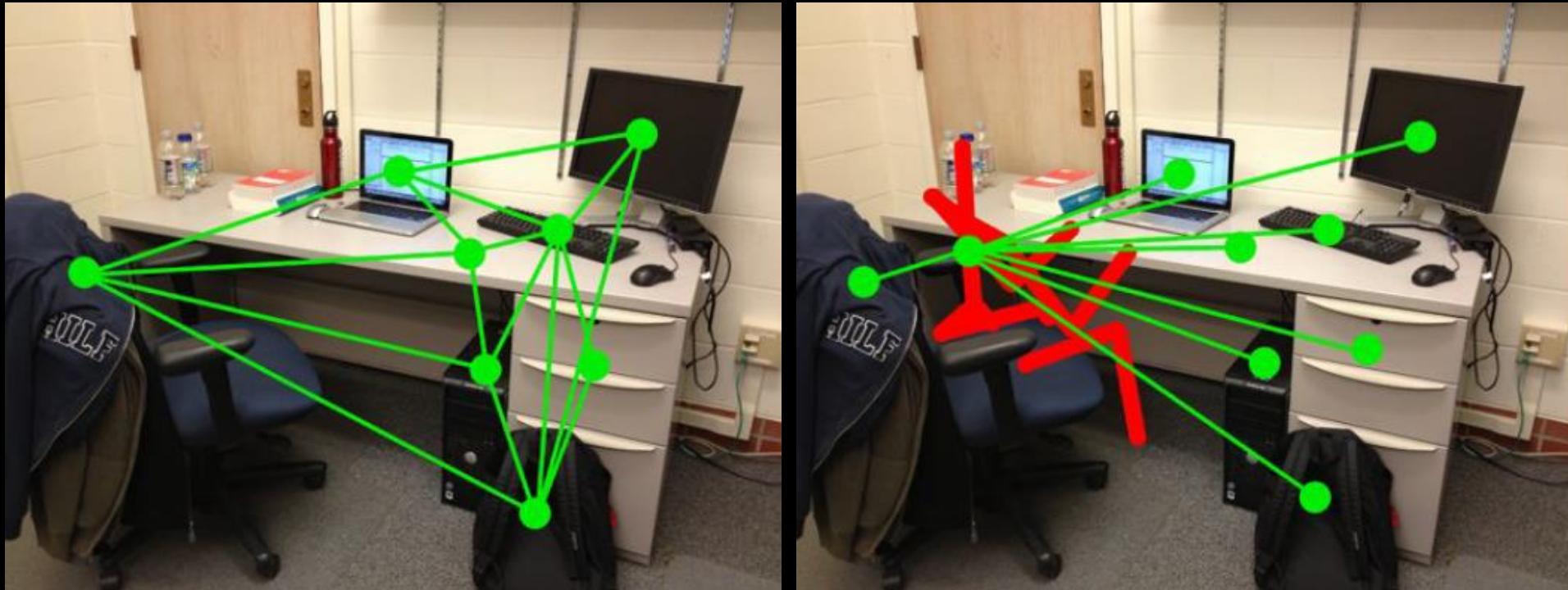
# Integrating more constraints

- *Constraints*
  - Volumetric constraints
- *Techniques*
  - Structured prediction
  - Sampling

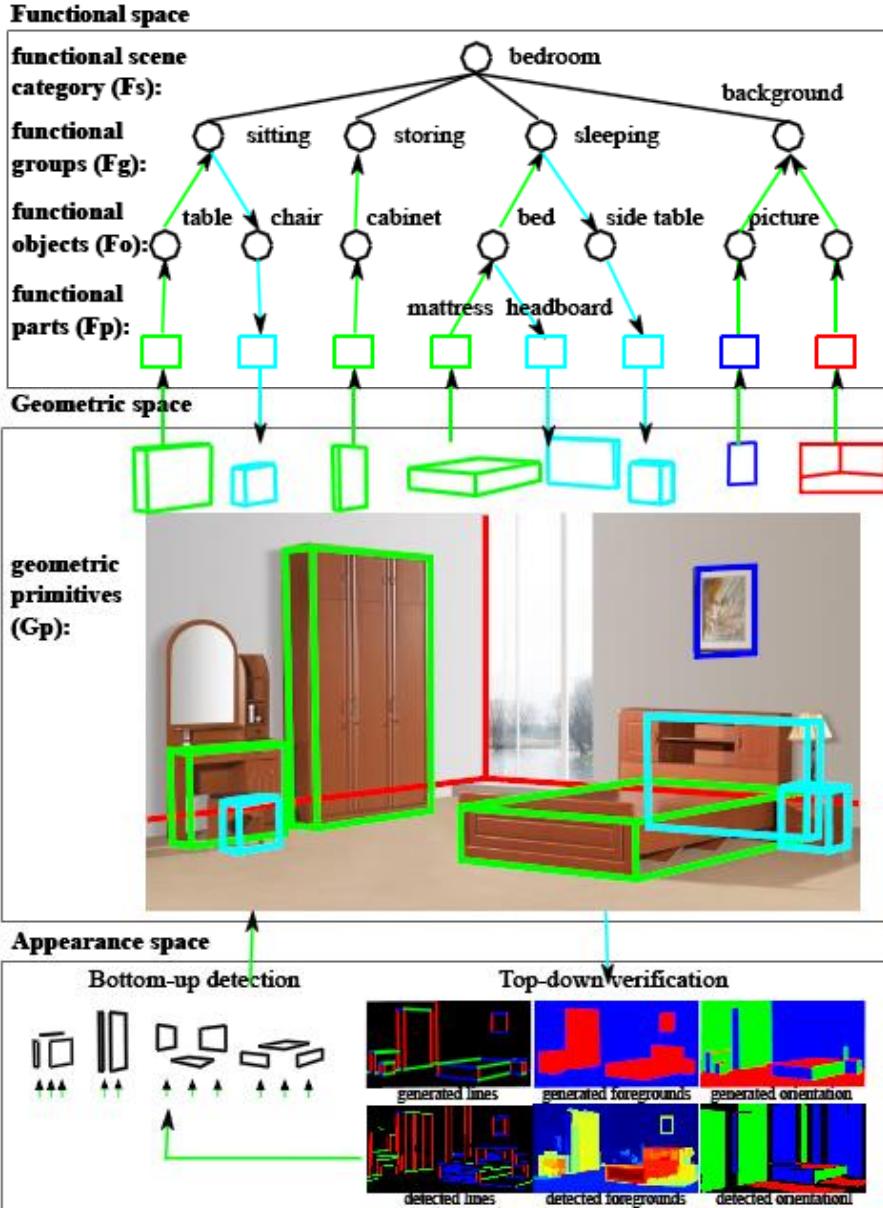
# Integrating more constraints

- *Constraints*
  - Volumetric constraints
  - Physical constraints
  - Relative placement
  - Functional constraints
- *Techniques*
  - Structured prediction
  - Sampling
  - Search through hypothesis space

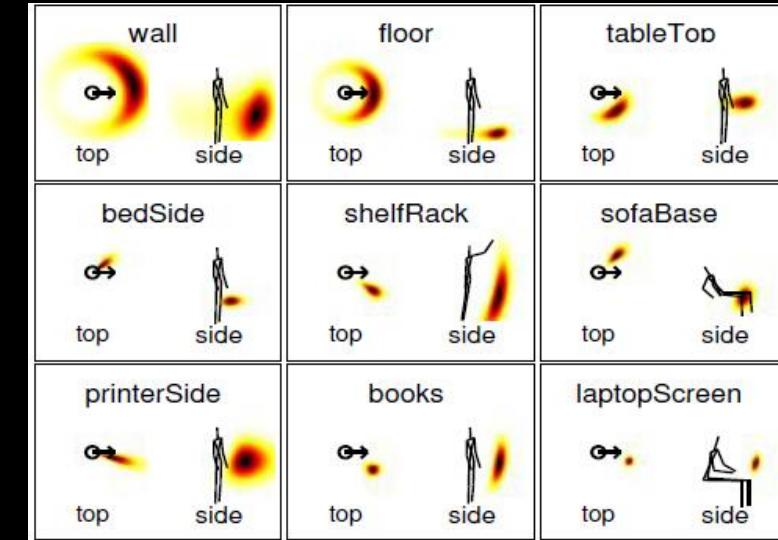
# Even more constraints: Functional



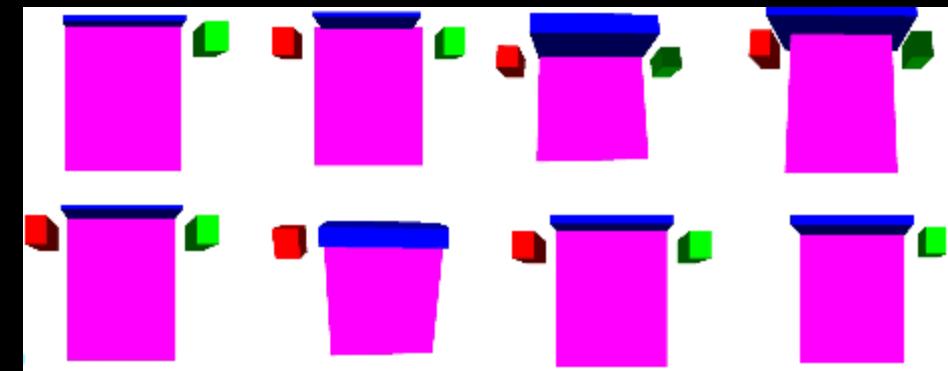
J.J. Gibson. The Theory of Affordances. Lawrence Erlbaum, 1977.



## 1. Structural model function-geometry-appearance



2. Estimate distributions from training data



3. Sample using model

Y.Z. Zhao, S-C. Zhu. Scene Parsing by Integrating Function, Geometry, and Appearance Models. CVPR 2013.  
 Y. Jiang, H. Koppula, A. Saxena. Hallucinated Humans as the Hidden Context for Labeling 3D Scenes. CVPR 2012.

# People as Clutter?



People Occlude the Scene!

# Humans tell a lot about geometry





Timelapse



Timelapse



## Pose Detections

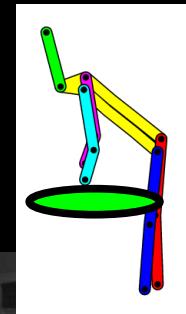
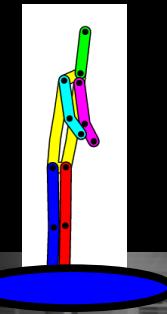
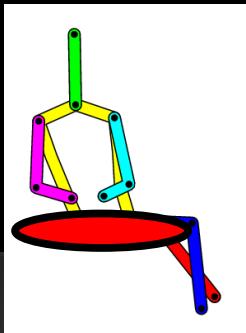
People Watching: Human Actions as a Cue for Single View Geometry. D. Fouhey, V. Delaitre, A. Gupta, A.A. Efros, I. Laptev, J. Sivic. IJCV 2014



Timelapse



Pose Detections



## Estimate Functional Regions from Poses

People Watching: Human Actions as a Cue for Single View Geometry. D. Fouhey, V. Delaitre, A. Gupta, A.A. Efros, I. Laptev, J. Sivic. IJCV 2014



## 3D Room Hypotheses From Appearance

People Watching: Human Actions as a Cue for Single View Geometry. D. Fouhey, V. Delaitre, A. Gupta, A.A. Efros, I. Laptev, J. Sivic. IJCV 2014



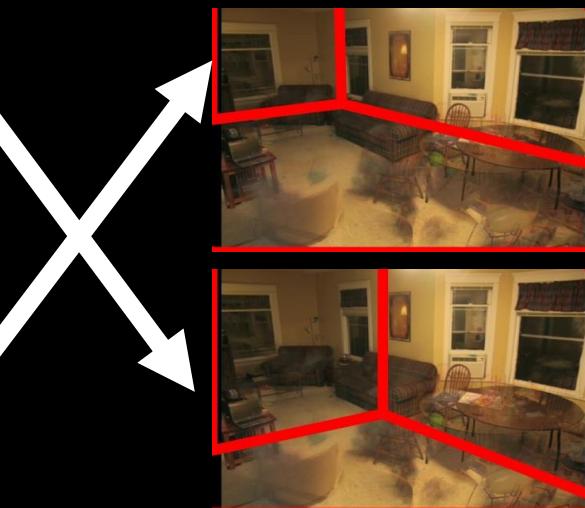
Timelapse



Pose Detections



Functional Regions



#1

#49

Score 3D Room Hypotheses With  
Appearances + Affordances



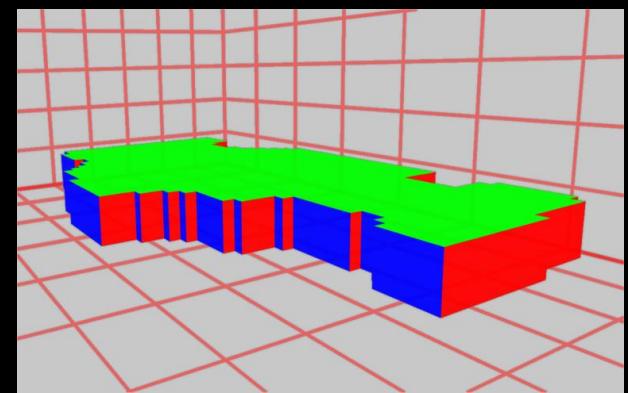
Timelapse



Pose Detections



Functional Regions



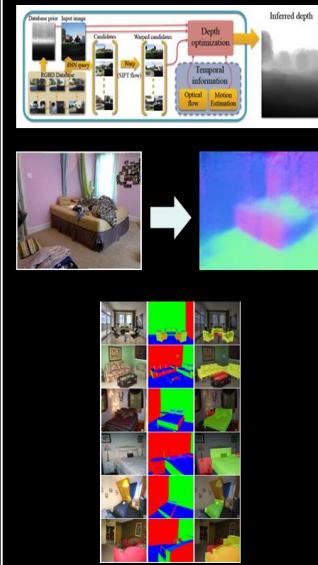
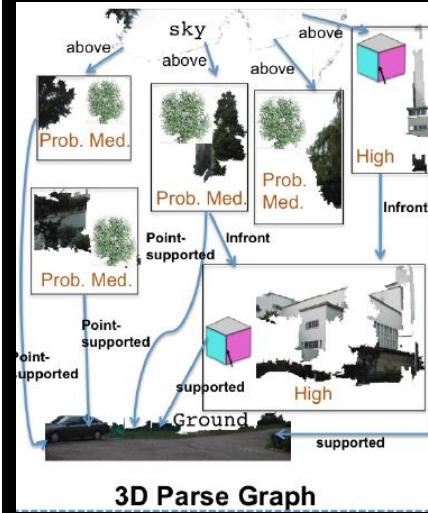
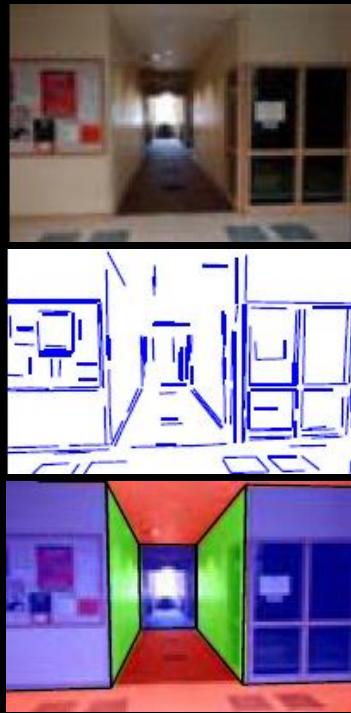
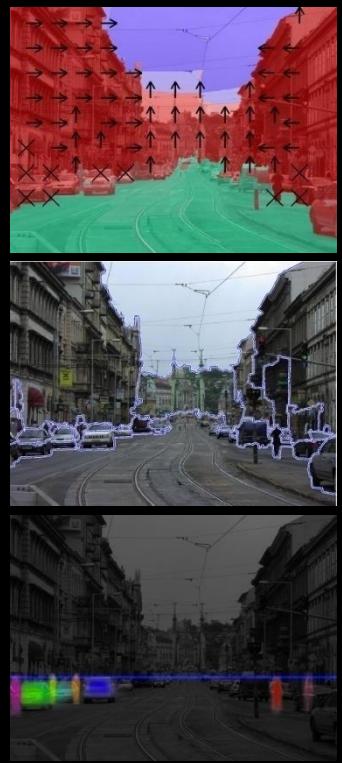
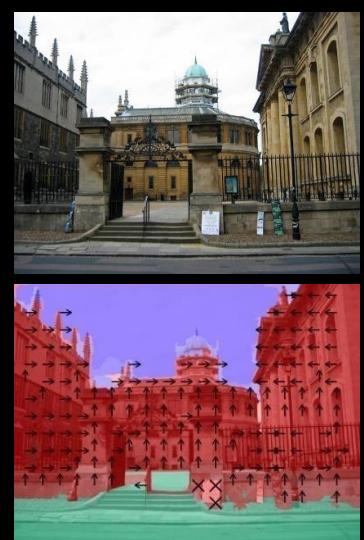
Estimate  
Free-Space



Appearance Alone



Appearance + People



Region labels

+ Boundaries  
and objects

Stronger geometric  
constraints from  
domain knowledge

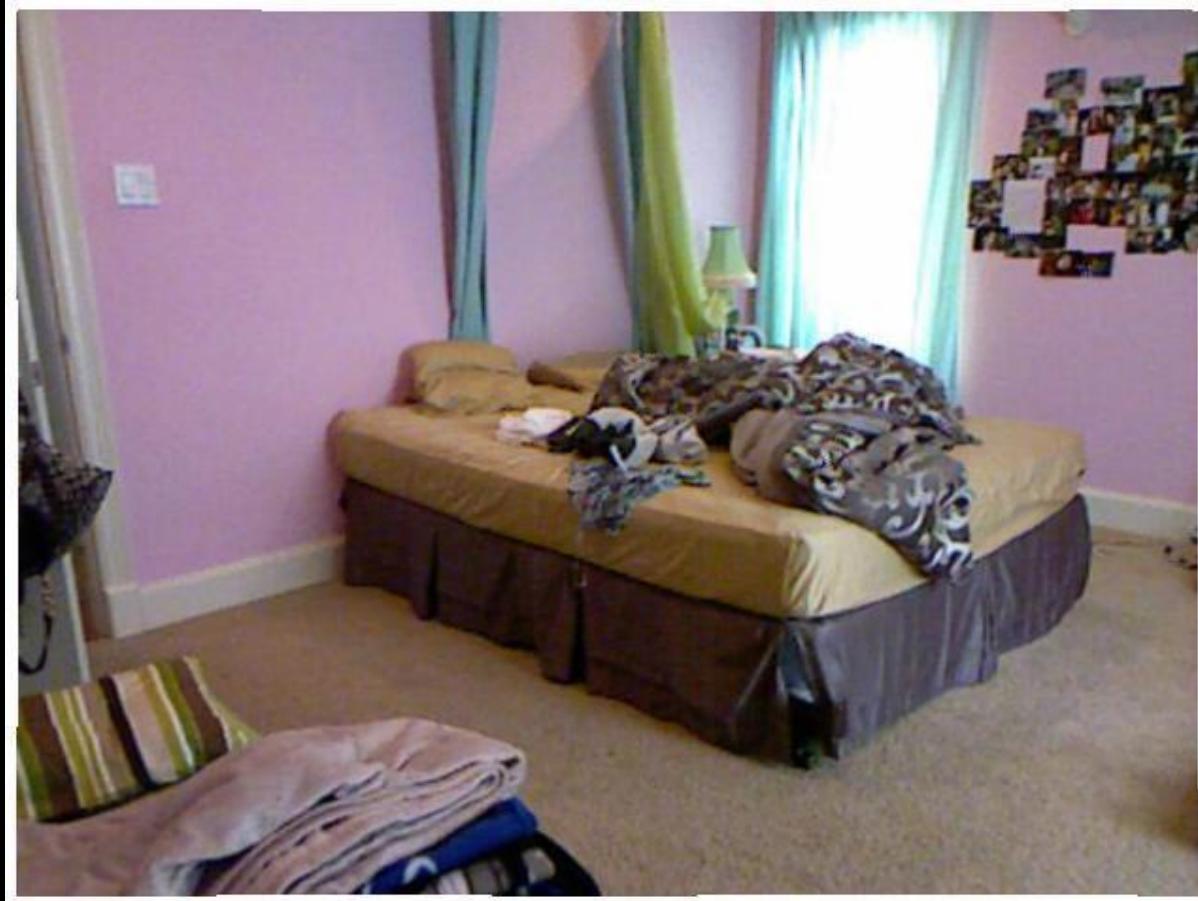
+ more constraints

+ Data-driven  
approaches

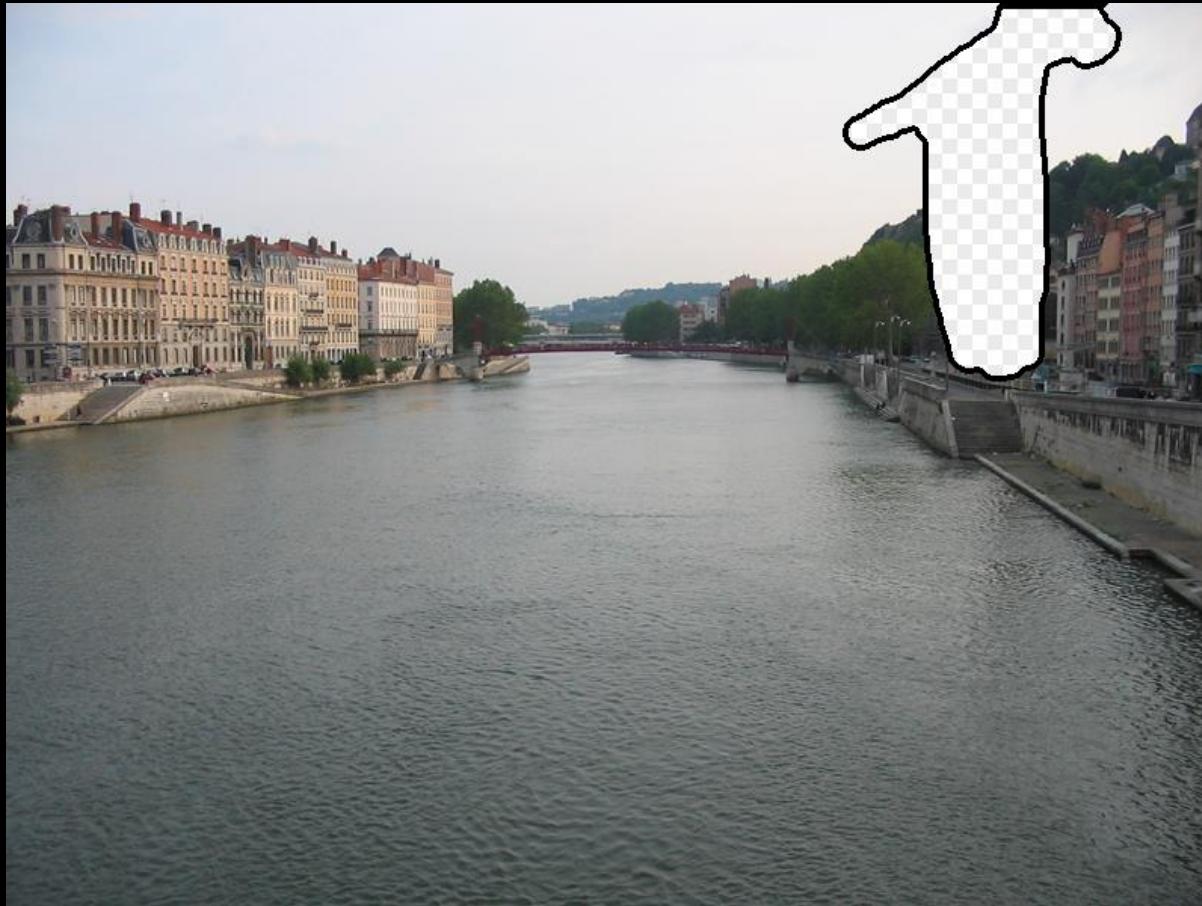
Qualitative

More quantitative  
more precise

# Data-Driven Approaches



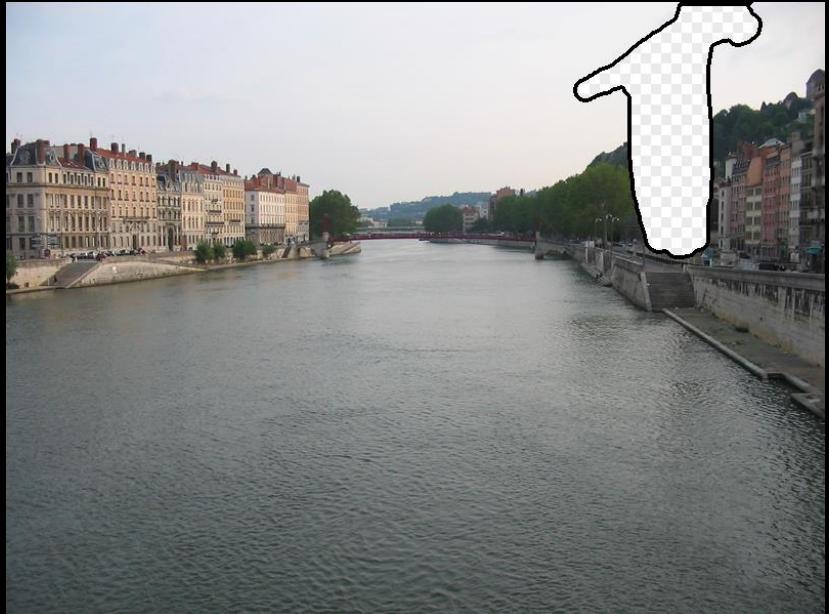
# Data-Driven Interpretation



Every image that can be seen has been seen before (approximately)

Hays and Efros 2007

# Data-Driven Interpretation

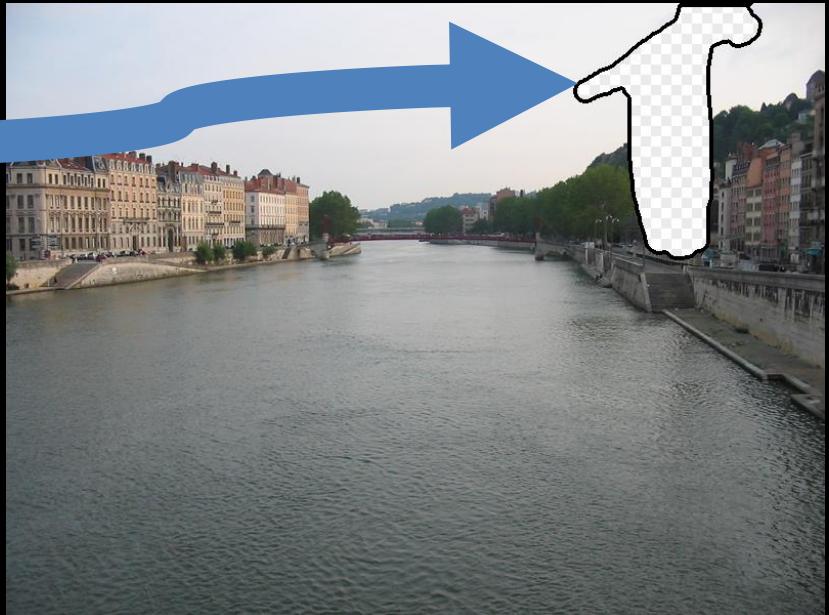
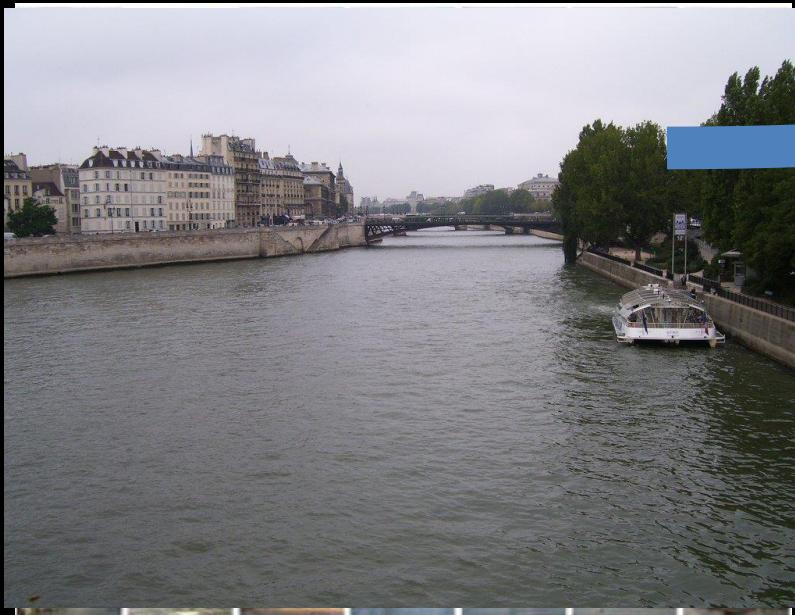


...

Every image that can be seen has been seen before (approximately)

Hays and Efros 2007

# Data-Driven Interpretation



...

Every image that can be seen has been seen before (approximately)

Hays and Efros 2007

# Data-Driven Interpretation



...

Every image that can be seen has been seen before (approximately)

Hays and Efros 2007

# Data-Driven Interpretation

Works well where parametric modeling is hard but where there's data



# Advantages

Volumetric  
Interpretation

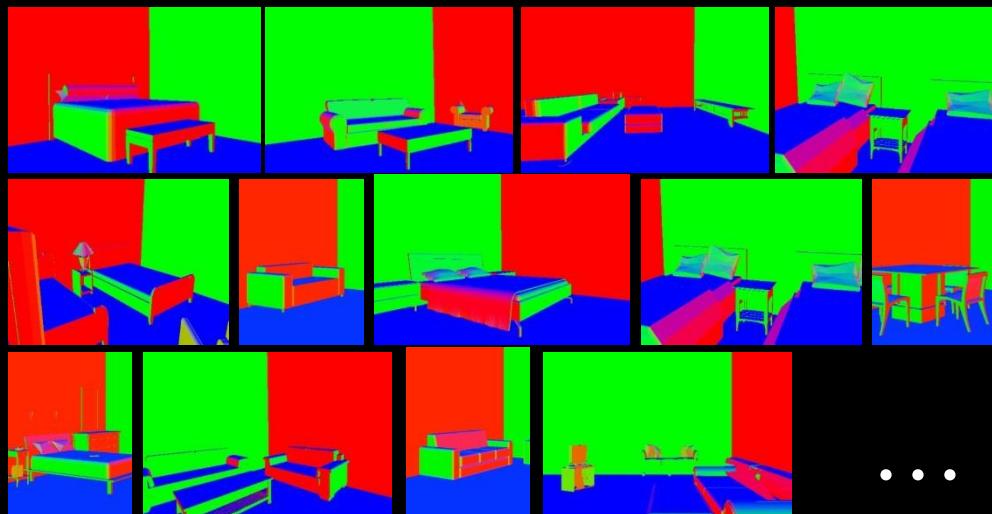


Interpretation by  
3D Models



# Sources

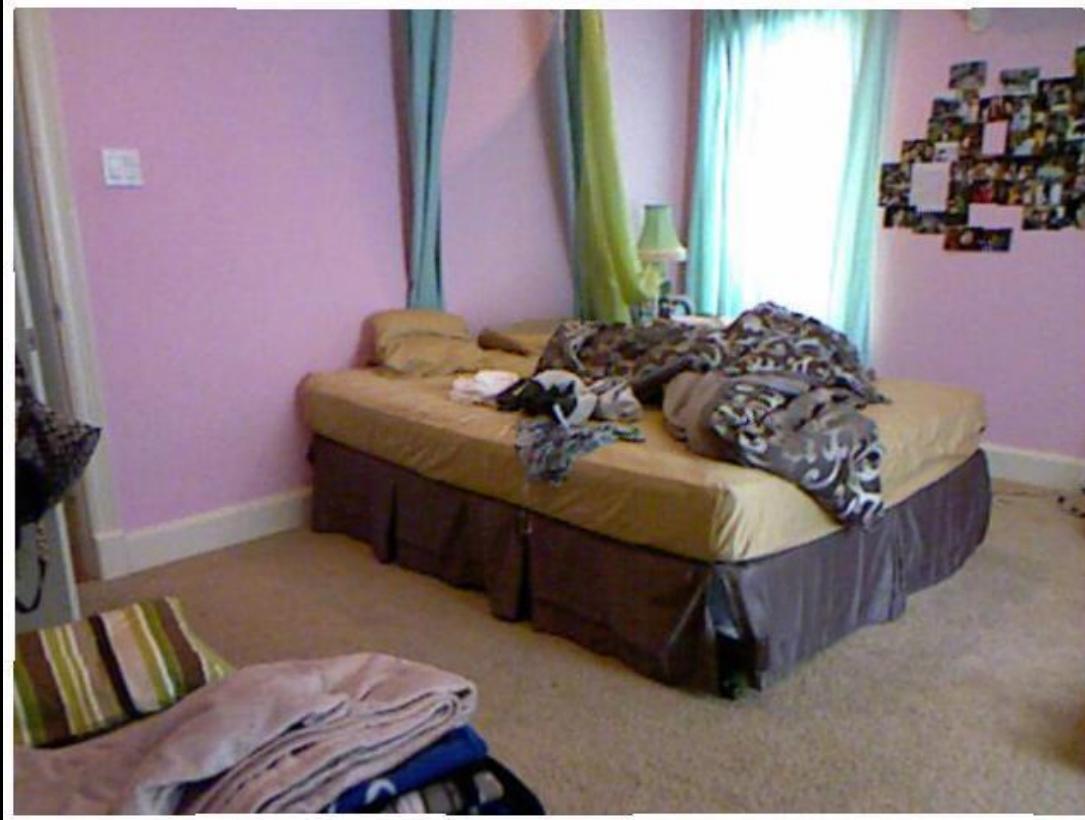
3D Model Databases



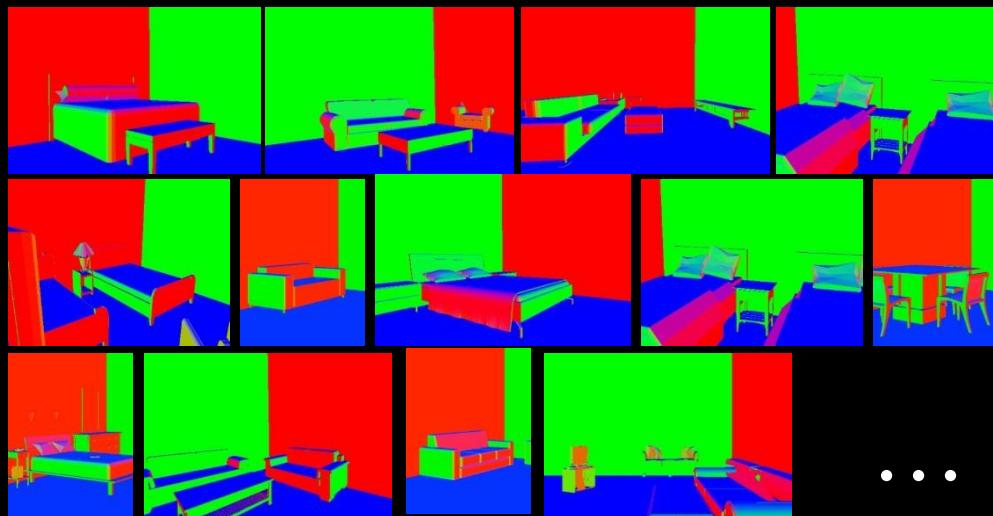
Kinect Databases



# Goal



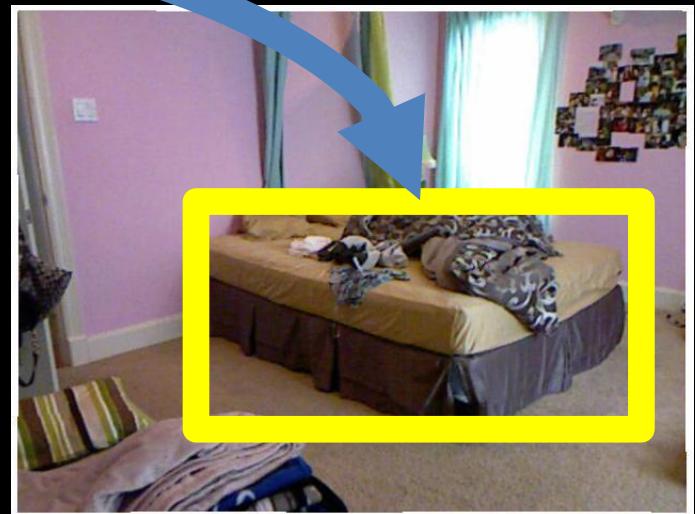
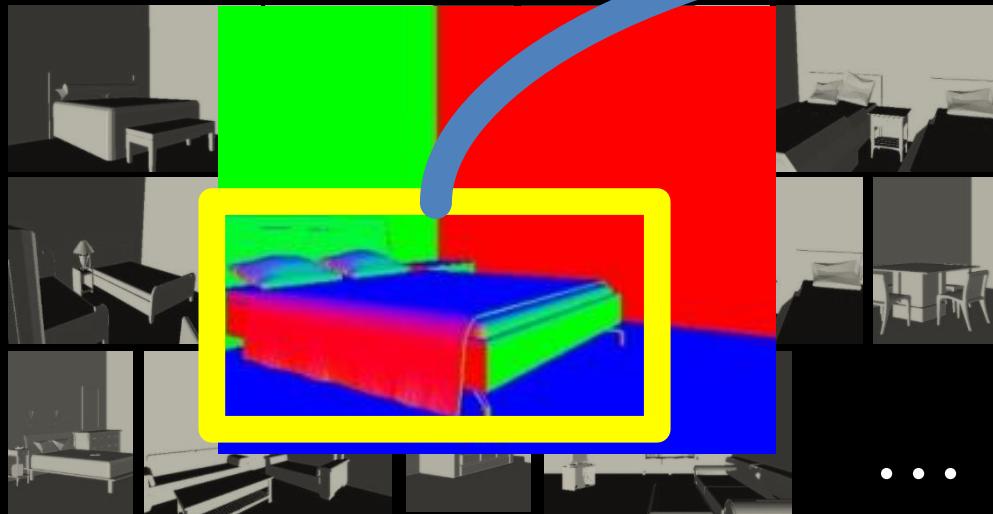
# Goal



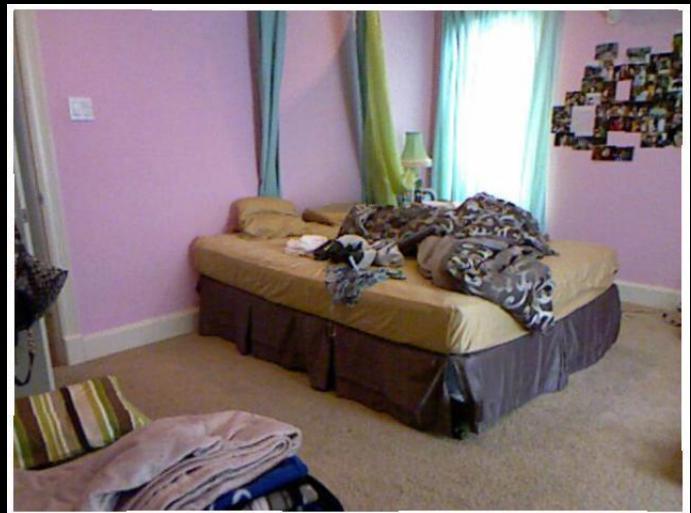
...



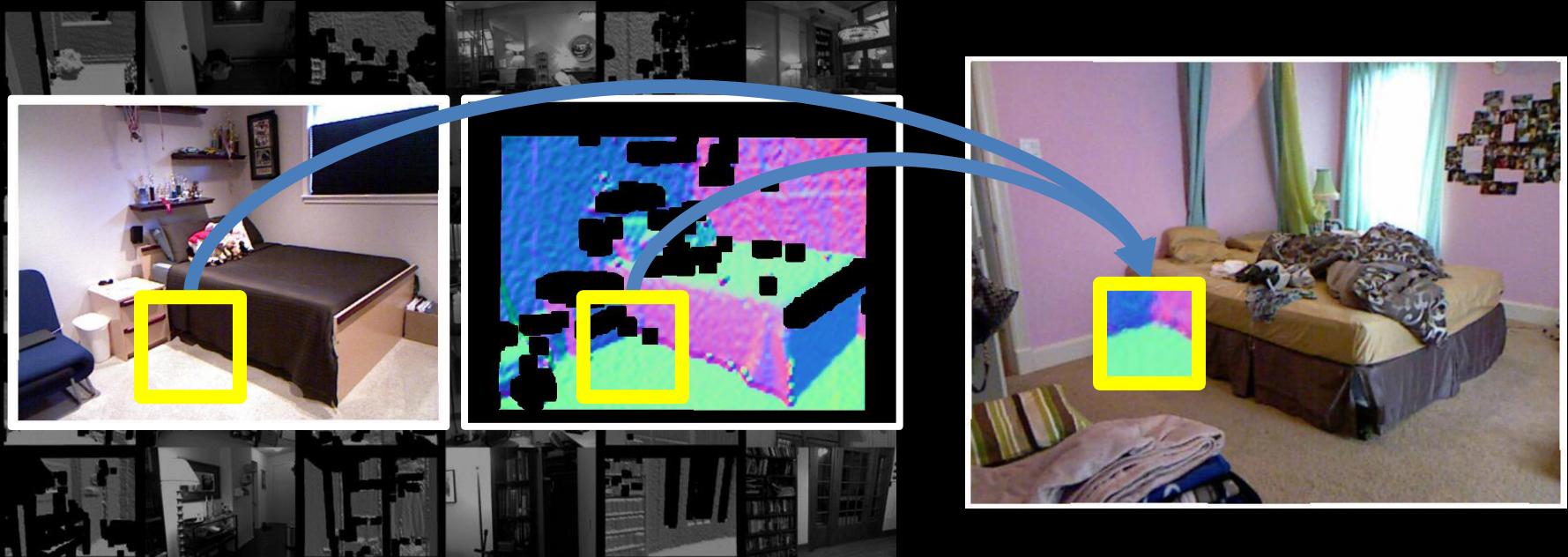
# Goal



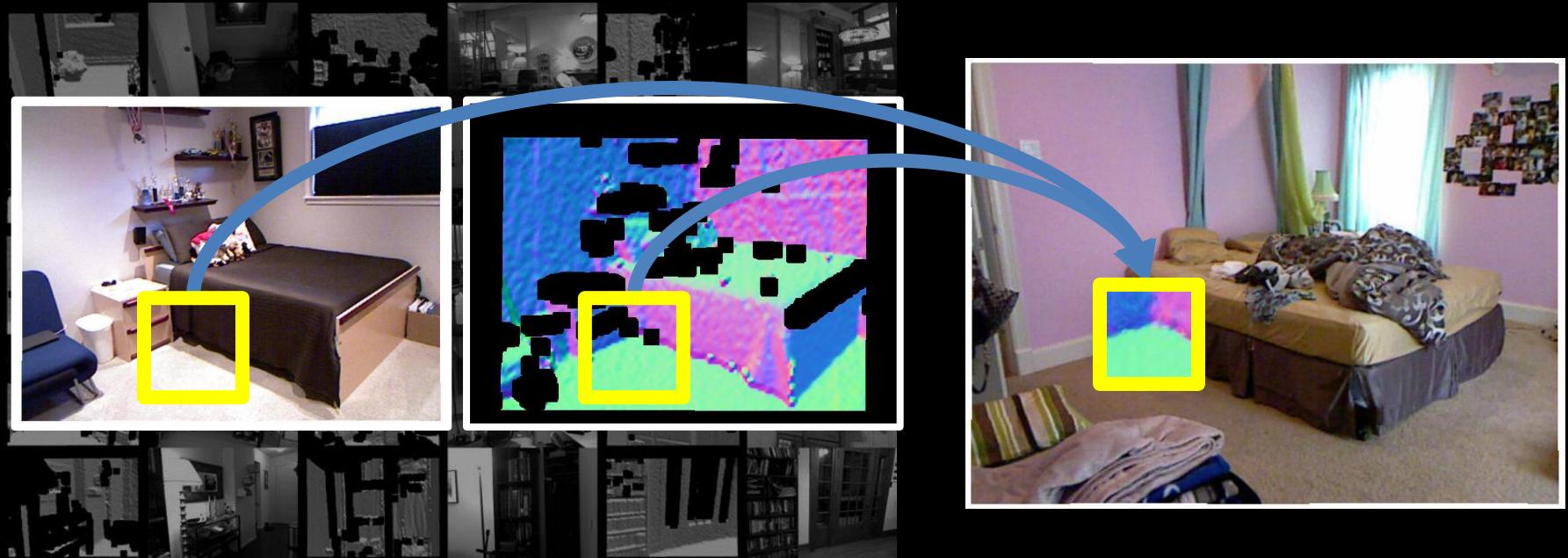
# Goal



# Goal



# Goal



How do you:

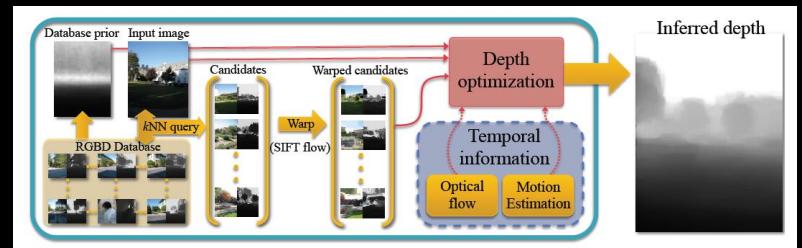
- (a) establish correspondence?
- (b) transfer representations?

# Overview

## 1. How to use 3D models



## 2. How to use depth images



# Why 3D Models

Object Detector



Segmentation

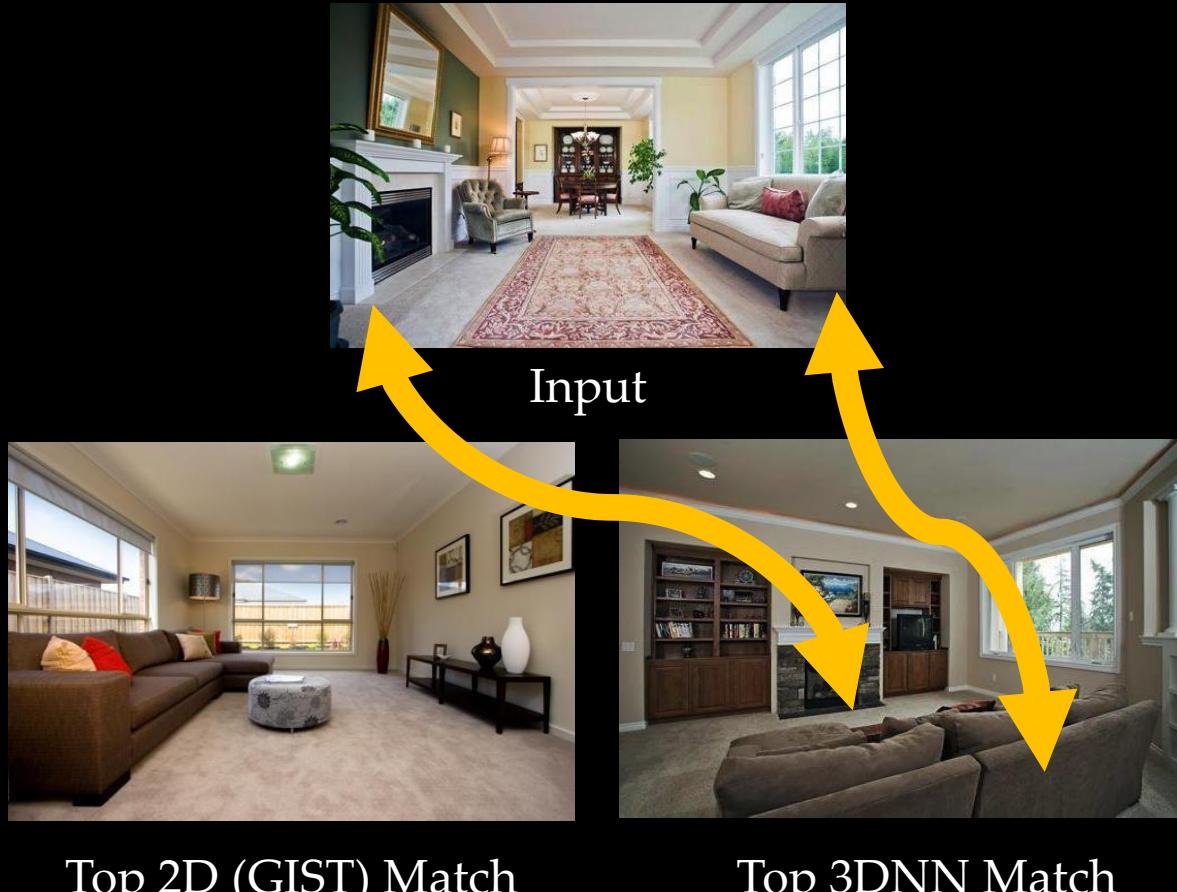


3D Model



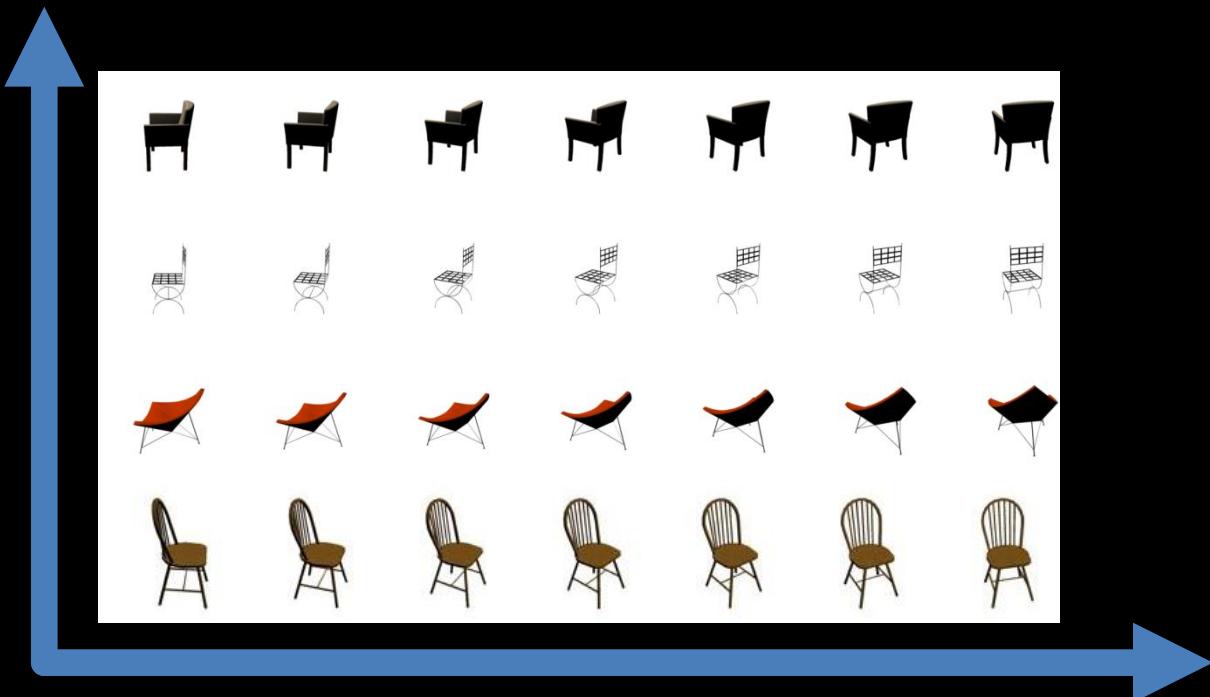
Aubry et al. 2014

# Why 3D Models



# General Approach

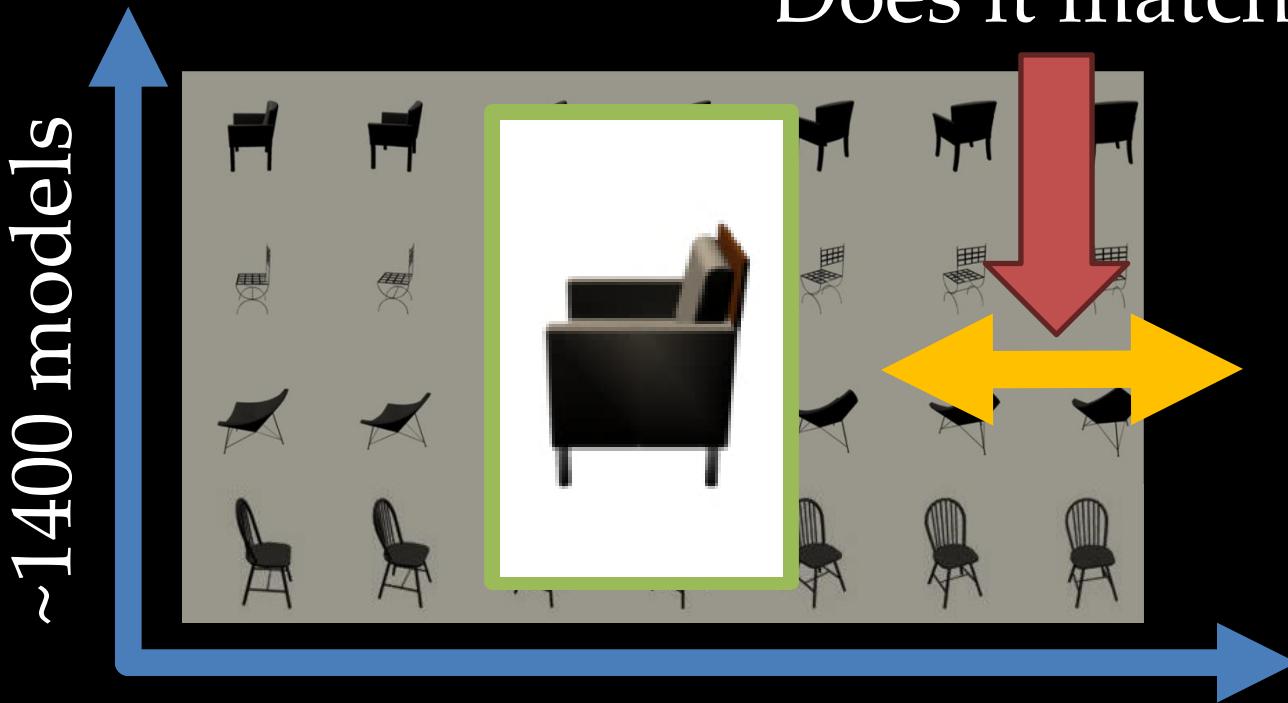
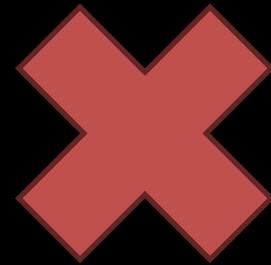
Search over model and viewpoint



Chairs from Aubry et al. 2014

# Primary Question

Does it match?

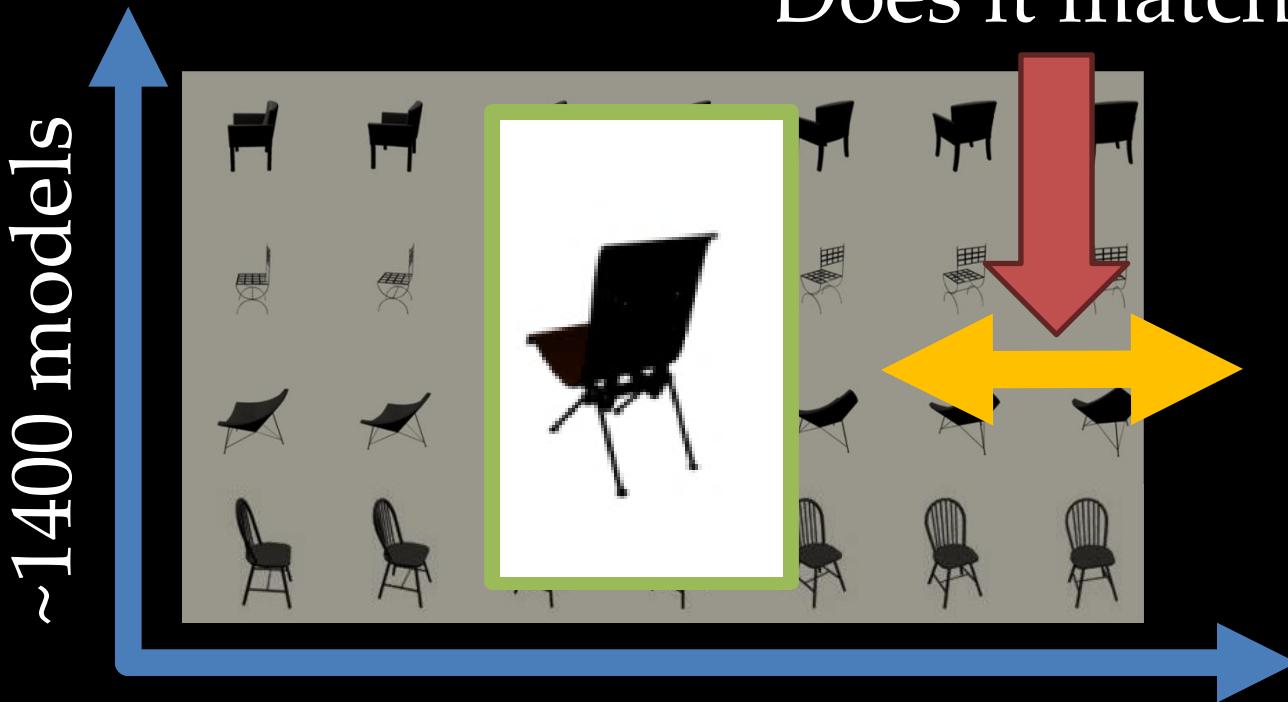
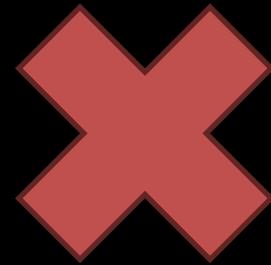


$\sim 60$  viewpoints

Chairs from Aubry et al. 2014

# Primary Question

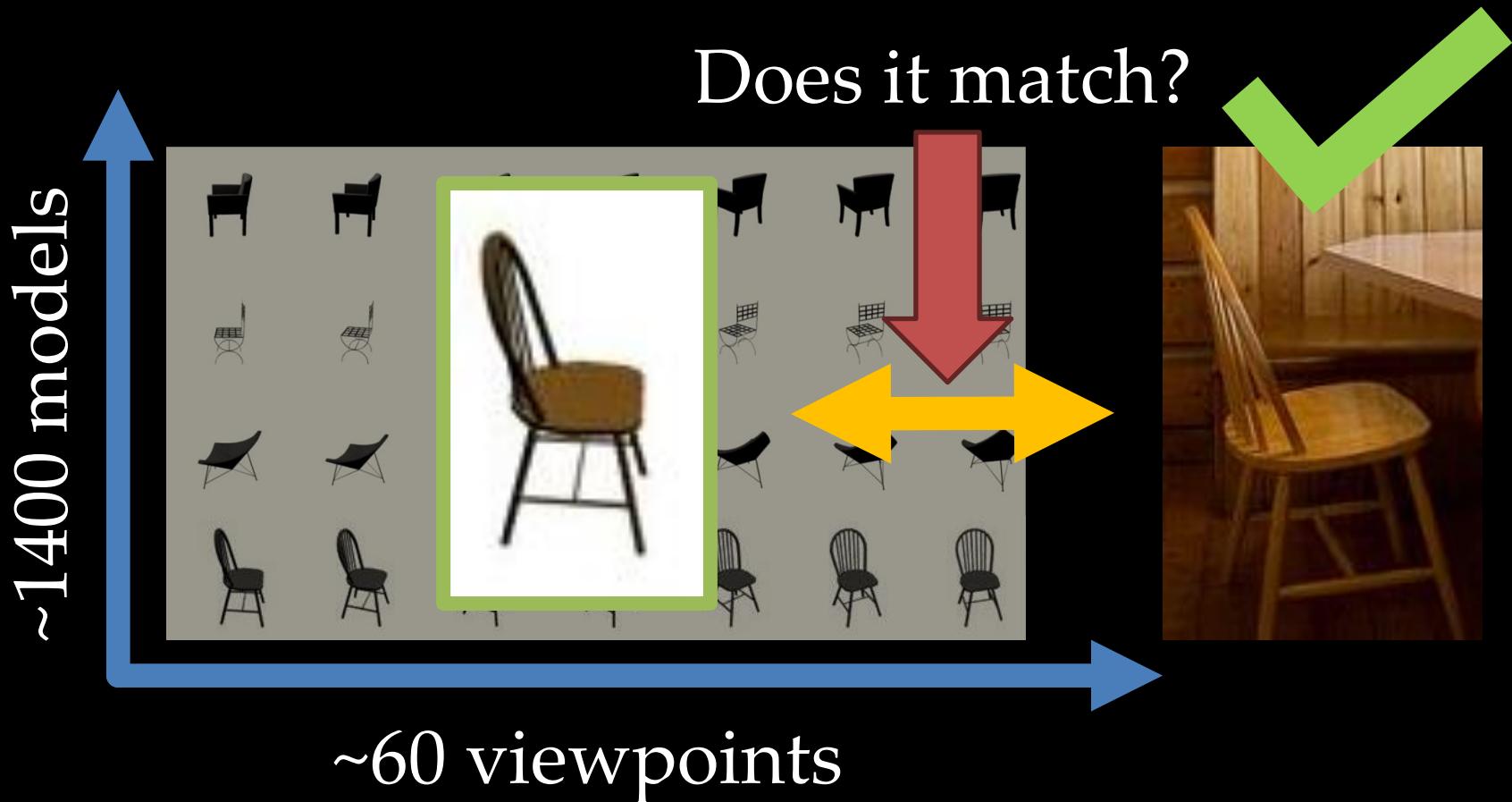
Does it match?



$\sim 60$  viewpoints

Chairs from Aubry et al. 2014

# Primary Question



Chairs from Aubry et al. 2014 141

# Difficulties

Rendered      Natural



Texture

NO

YES

Occlusion

NO

YES

Background

Fake

Natural

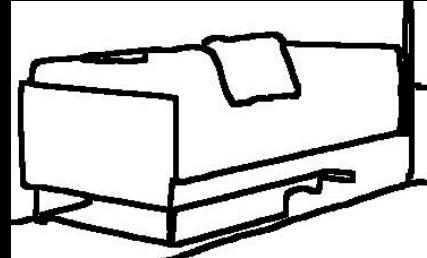
# Cross-Domain Matching

Goal: bring image and model into  
common representation

# Chamfer Matching

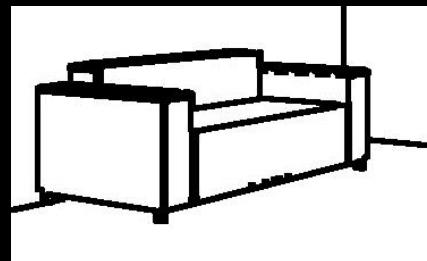
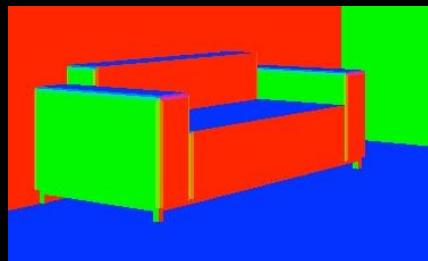
Assumption: edges in 3D are edges in 2D

Image



Match?

3D Model

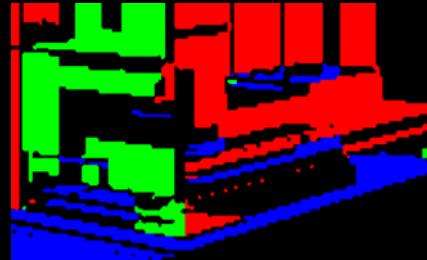


Satkin et al., 2012, 2013, 2014; Lim et al., 2013; Ramnath et al., 2014;

# Domain-Invariant

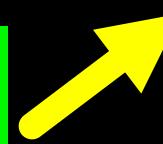
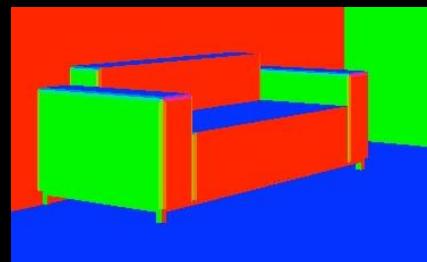
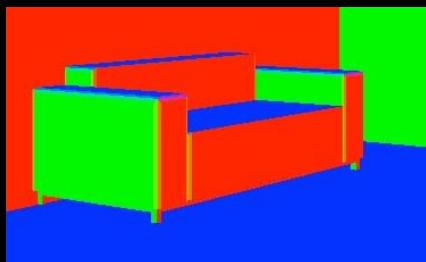
Assumption: can estimate 3D property from 2D

Image



Match?

3D Model

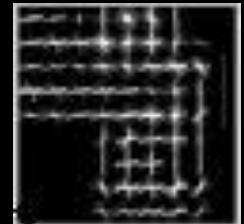
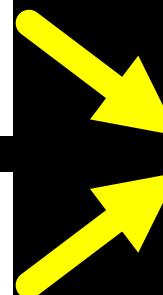


# Domain-invariant “Images”

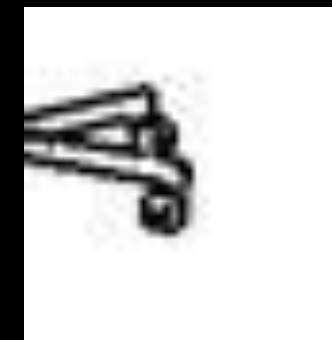
Assumption: edges in 3D are edges in 2D

Apply standard features/techniques

Image



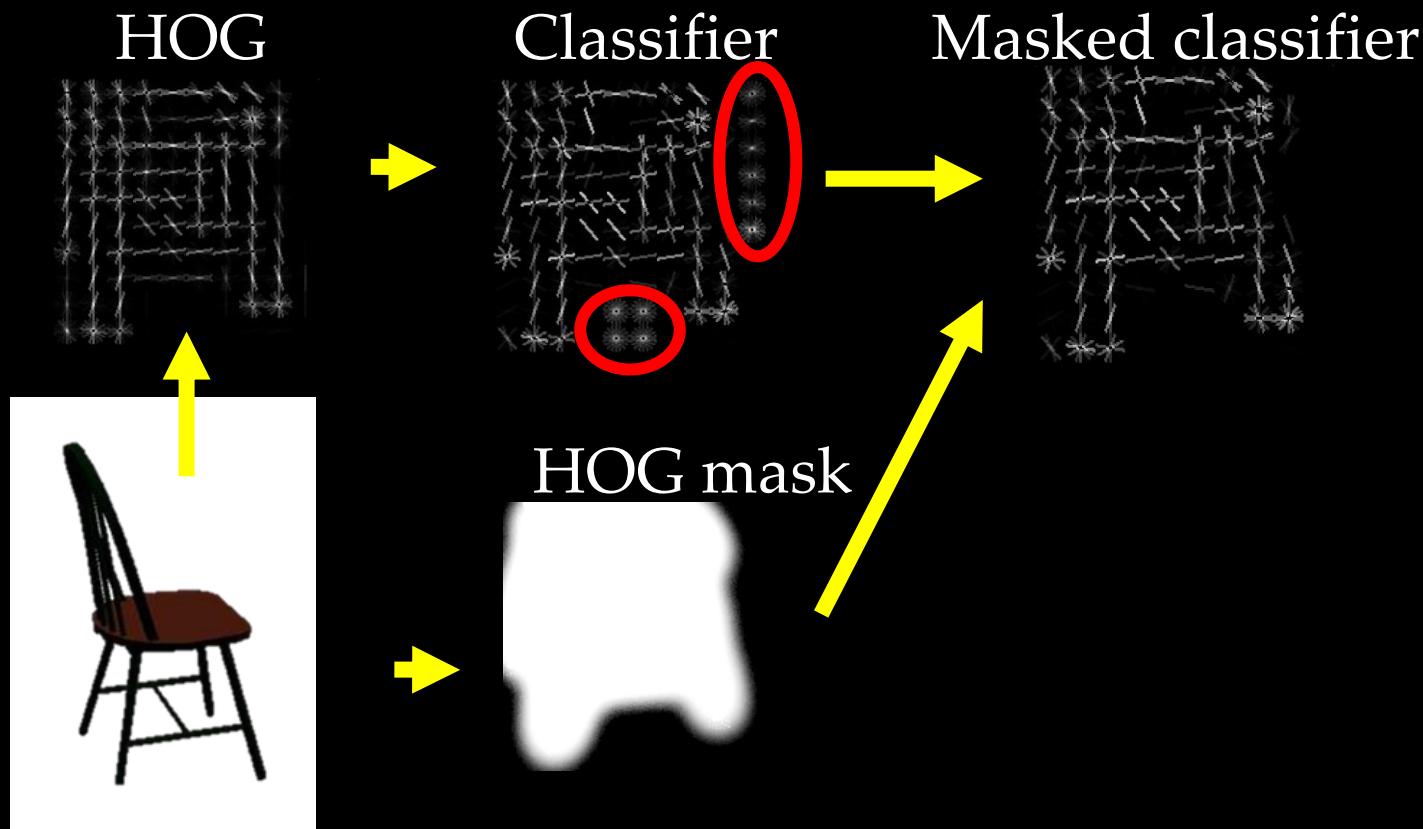
3D Model



Lim et al., 2013

# Masking Features

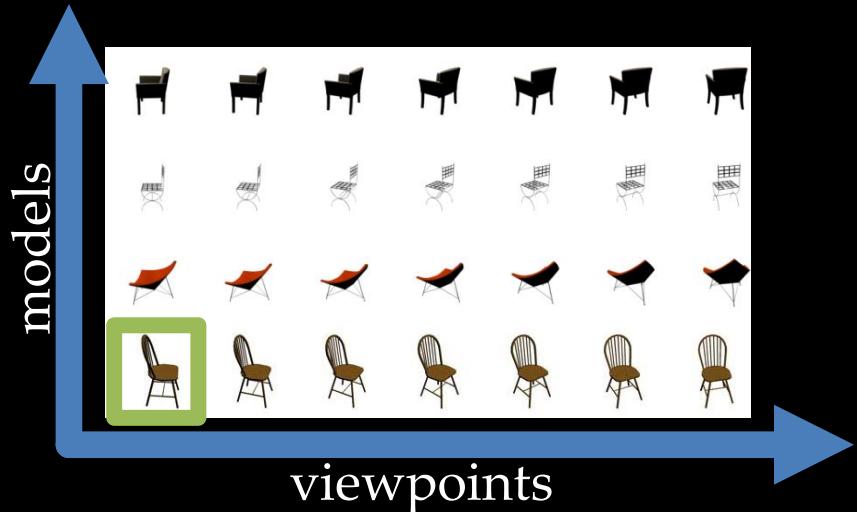
Assumption: only issue is background



Aubry et al., 2014; see also Shrivastava et al., 2011

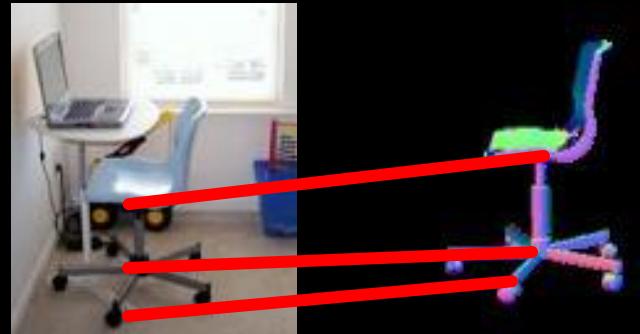
# Searching Hypotheses

Render object parts



Aubry et al., 2014

Matches generate proposals

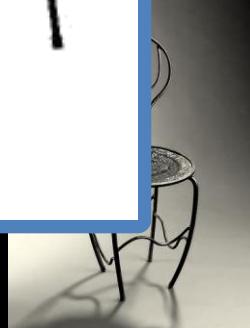


Lim et al., 2013

# Results



# Results



# Results

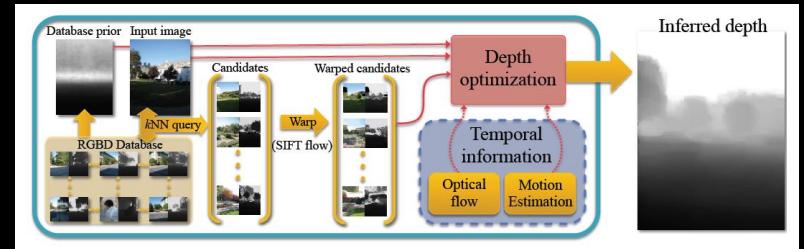
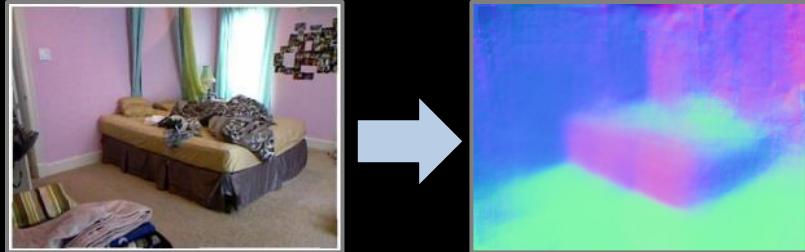


# Overview

## 1. How to use 3D models



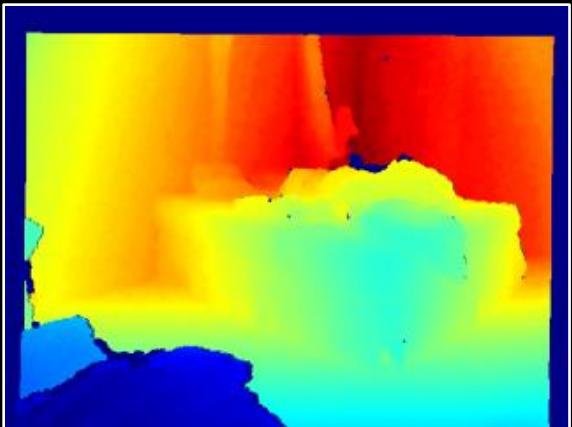
## 2. How to use depth images



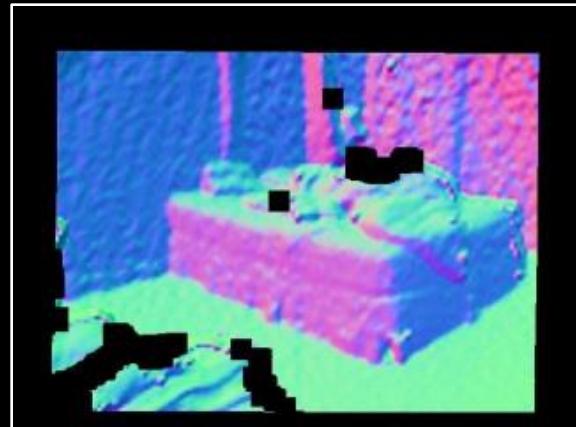
RGB



Depth

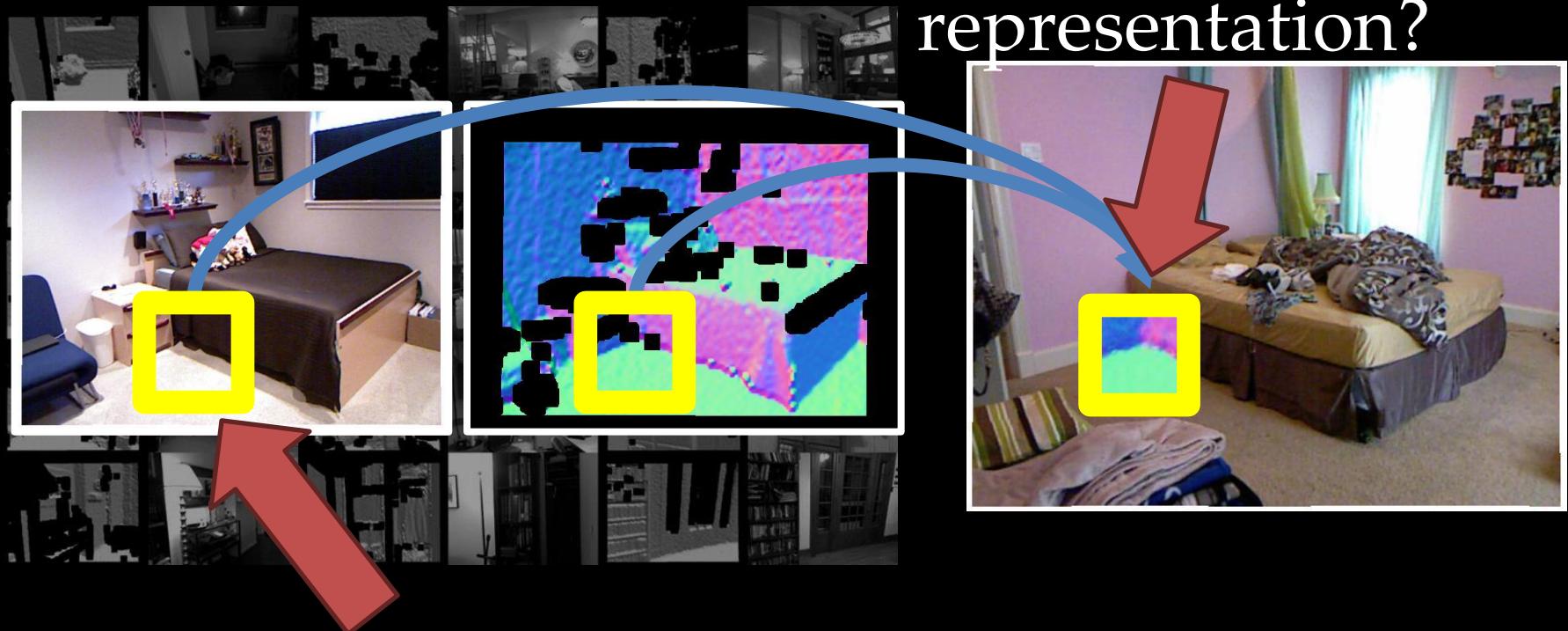


Normals



# General Approach

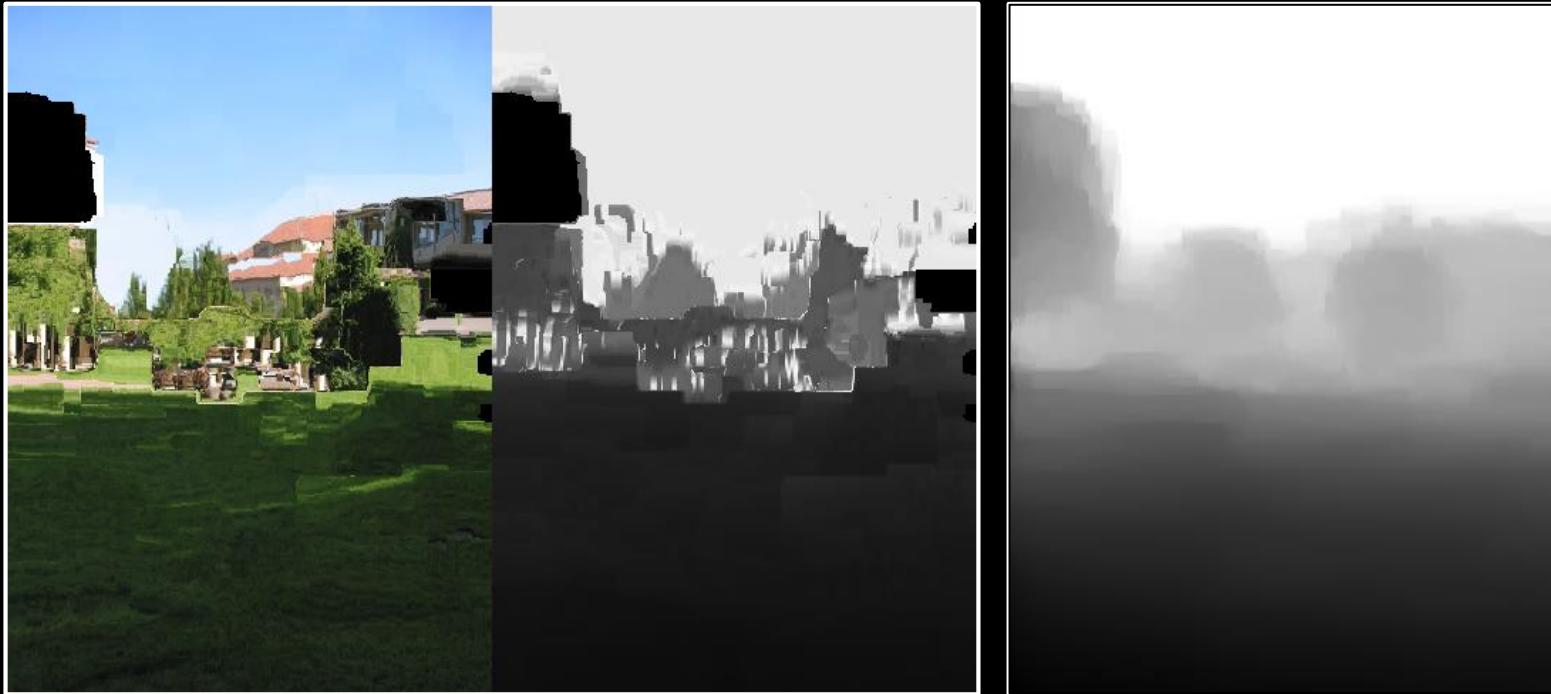
How to transfer  
representation?



How do we get this  
correspondence?

# Two Approaches

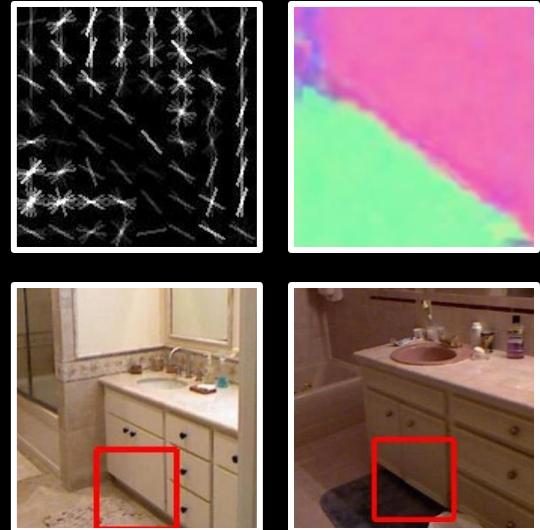
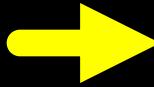
## Data-Driven Alignment



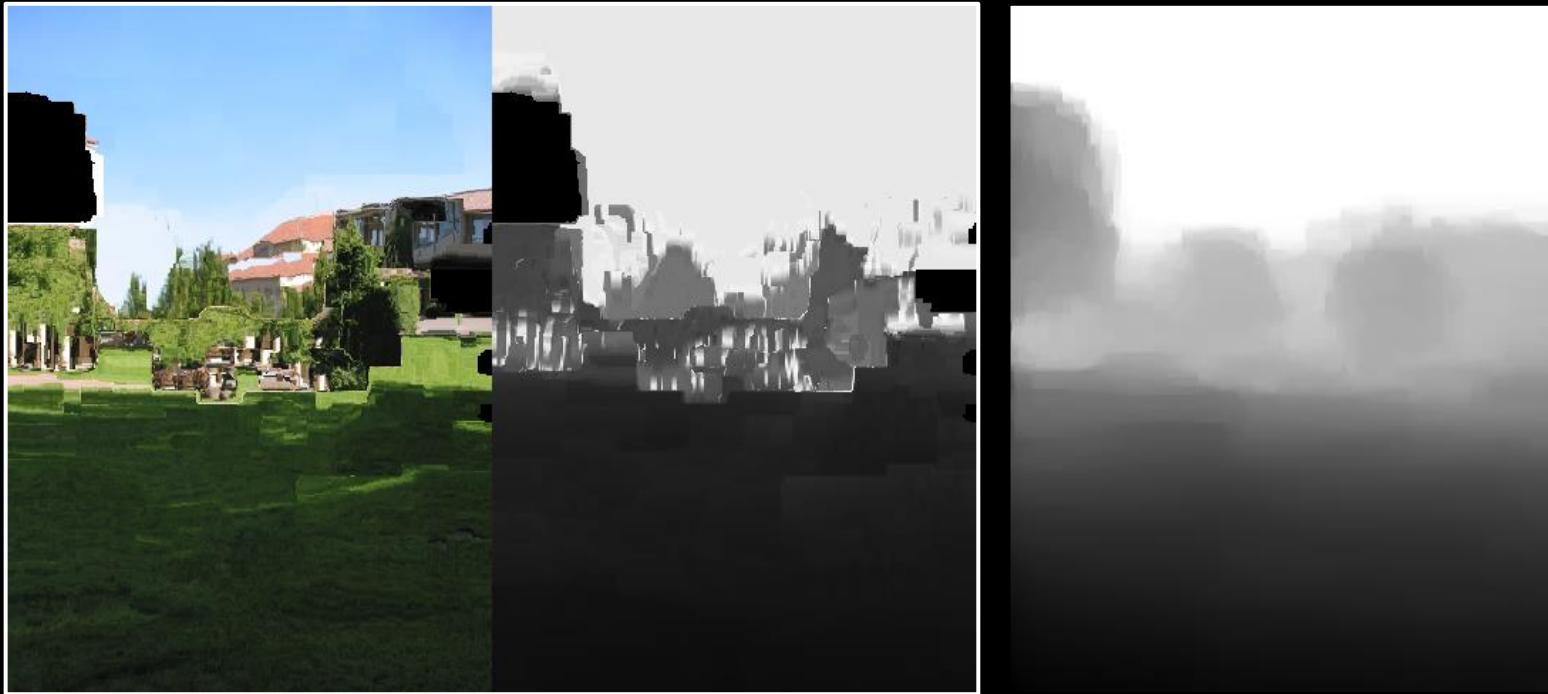
Karsch et al., 2012

# Two Approaches

## Clustering + Detection



# Data-Driven Alignment



Karsch et al., TPAMI 2014

# Finding Correspondences

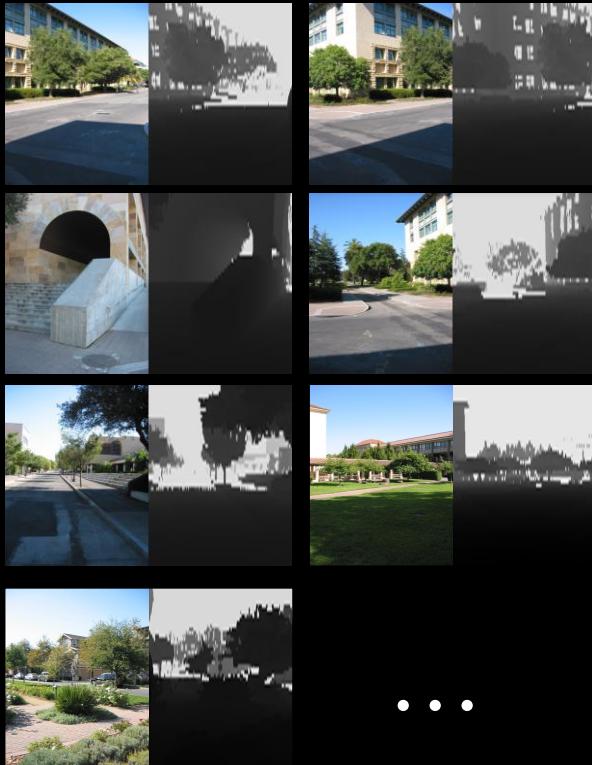
Input



Karsch et al., TPAMI 2014

# Finding Correspondences

Training Set



Input



# Finding Correspondences

Training Set

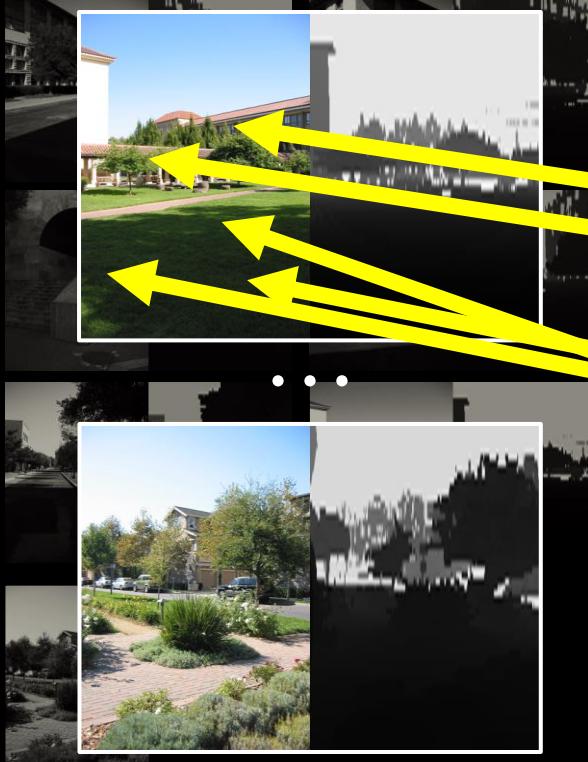


Input



# Finding Correspondences

Training Set

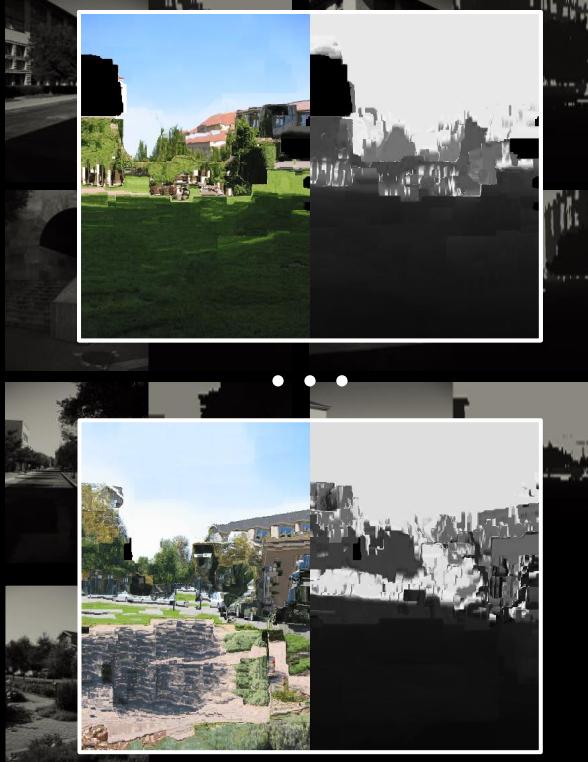


Input



# Finding Correspondences

Training Set



Input



# Finding Correspondences

Candidate 1

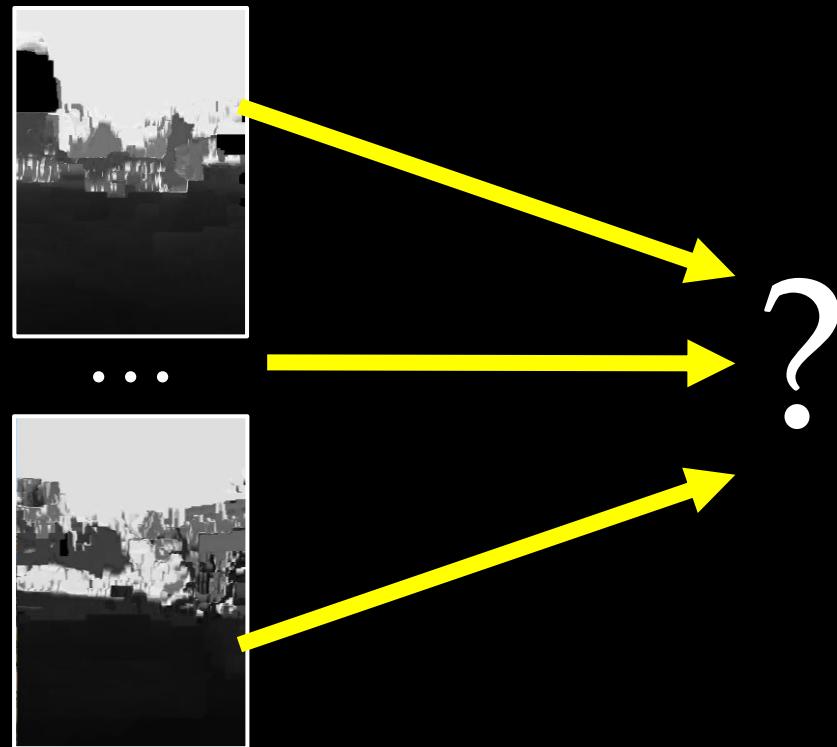


Candidate 2



# Finding Correspondences

Warped Depths



Karsch et al., 2012; see alternate approach from Liu et al., 2014

# Optimizing Depthmaps

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right] \\ + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

$D_i$  - Depth being optimized

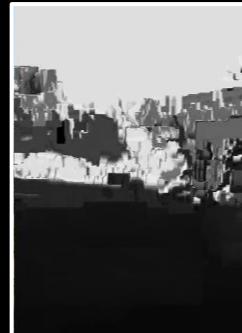
$C_i$  - Warped depth candidate

# Optimizing Depthmaps

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right] \\ + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$



...

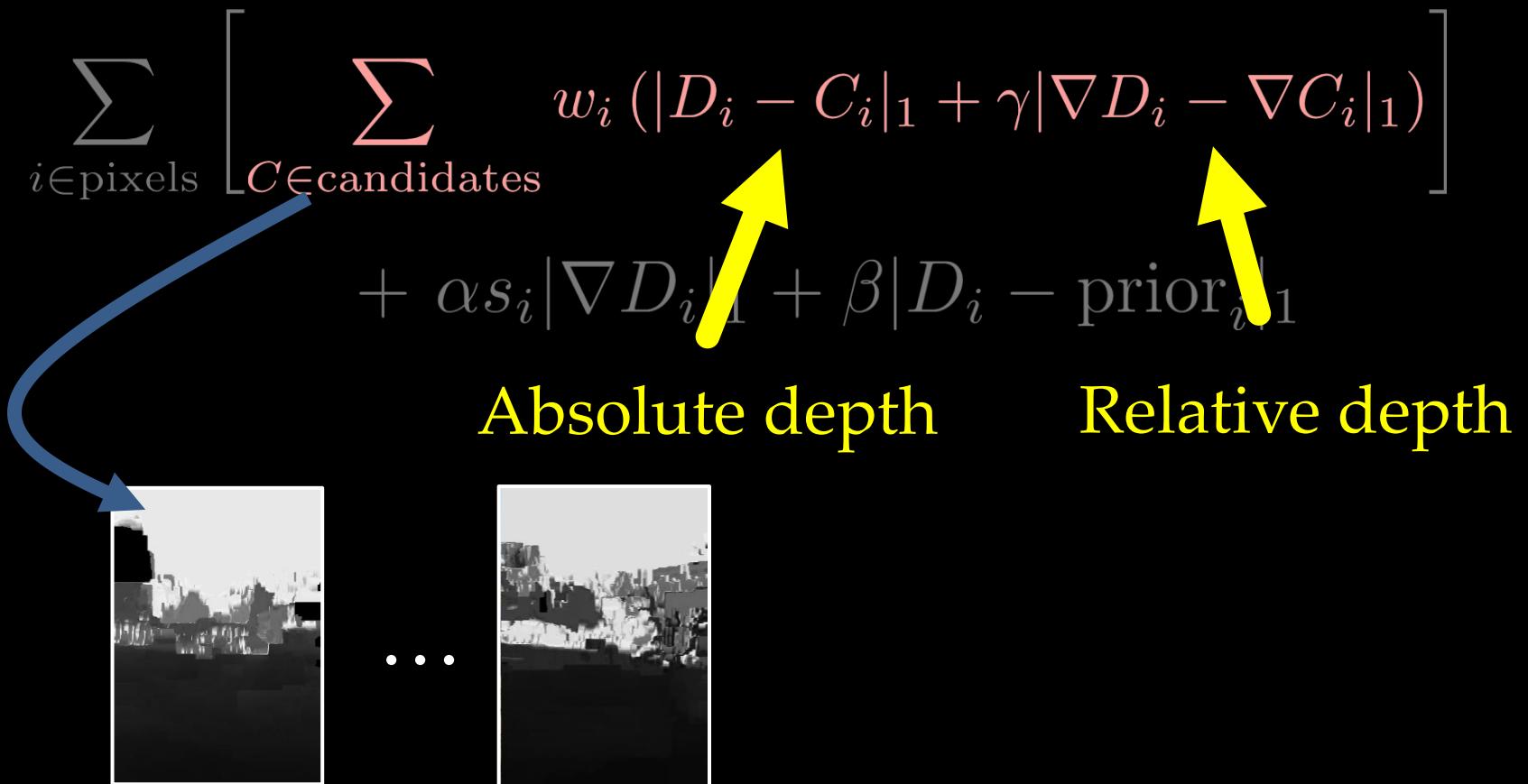


# Optimizing Depthmaps

Enforce depth to match candidates

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right] \\ + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Absolute depth      Relative depth



The diagram illustrates the optimization process. It shows two depthmap candidates for a single pixel. A blue arrow points from the first candidate to the term involving absolute depth in the equation. Two yellow arrows point from the second candidate to the terms involving relative depth and gradient matching in the equation.

# Optimizing Depthmaps

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right] \\ + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Spatial smoothness

# Optimizing Depthmaps

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i (|D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1) \right] \\ + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Match the prior

# Results

Input



True depth



Inferred depth



# Results

Input



True depth



Inferred depth



# Discriminative Clustering + Detection

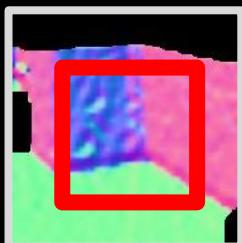
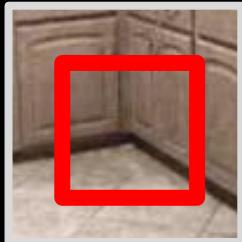


Fouhey et al., 2013, 2104

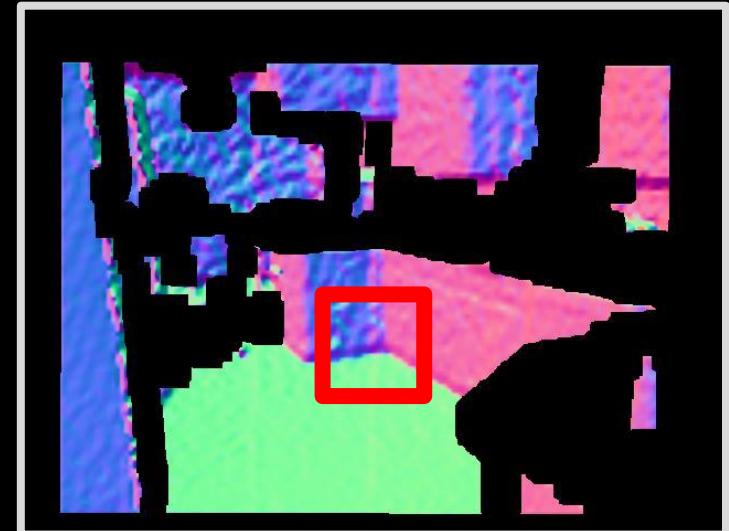
# Goal

Visually  
Discriminative

Geometrically  
Informative



Image



Surface Normals

# Goal

Learn from large-scale RGBD Data

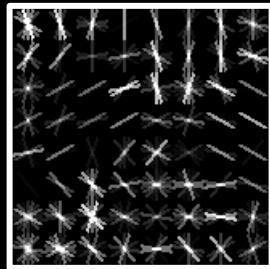


NYU Depth v2, Silberman et al., 2012, 2014

# Approach

Train time: discriminative clustering w/3D

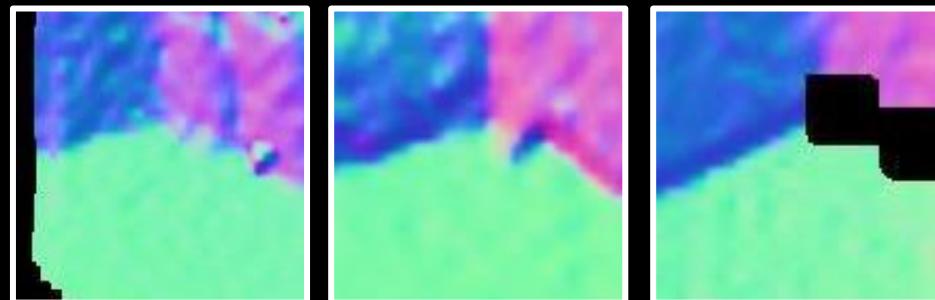
Detector



Instances

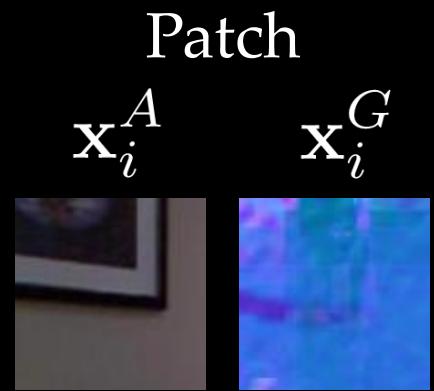
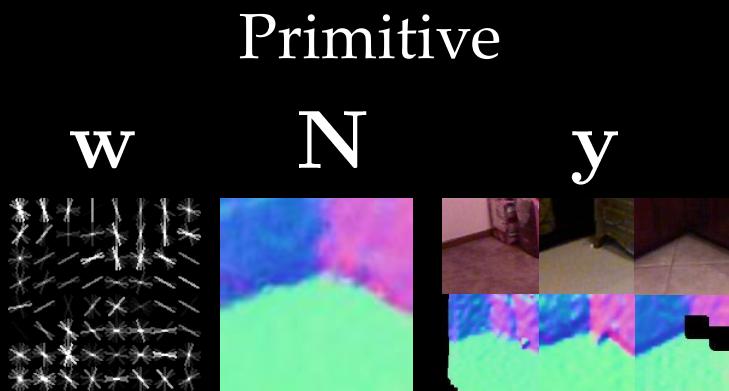


Normals



# Objective

$$\min_{\mathbf{y}, \mathbf{w}, \mathbf{N}} R(\mathbf{w}) + \sum_{i=1}^m \left[ c_2 L(\mathbf{w}, \mathbf{N}, \mathbf{x}_i^A, y_i) + c_1 y_i \Delta(\mathbf{N}, \mathbf{x}_i^G) \right]$$



Fouhey et al., ., ICCV 2013, ECCV 2014

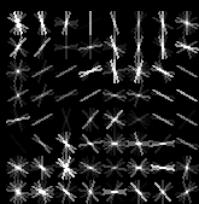
# Objective

Misclassification loss

$$\min_{\mathbf{y}, \mathbf{w}, \mathbf{N}} R(\mathbf{w}) + \sum_{i=1}^m \left[ c_2 L(\mathbf{w}, \mathbf{N}, \mathbf{x}_i^A, y_i) + c_1 y_i \Delta(\mathbf{N}, \mathbf{x}_i^G) \right]$$

Primitive

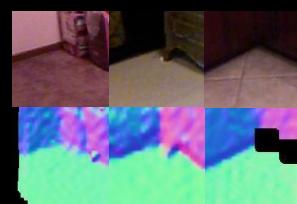
$\mathbf{w}$



$\mathbf{N}$



$\mathbf{y}$



Patch

$\mathbf{x}_i^A$



$\mathbf{x}_i^G$



# Objective

Regularization

$$\min_{\mathbf{y}, \mathbf{w}, \mathbf{N}} R(\mathbf{w}) + \sum_{i=1}^m \left[ c_2 L(\mathbf{w}, \mathbf{N}, \mathbf{x}_i^A, y_i) + c_1 y_i \Delta(\mathbf{N}, \mathbf{x}_i^G) \right]$$

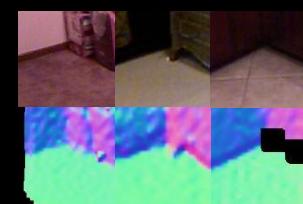
Primitive



$\mathbf{N}$



$\mathbf{y}$



Patch

$\mathbf{x}_i^A$



$\mathbf{x}_i^G$



# Objective

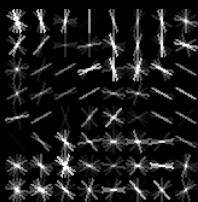
Ensure geometric consistency

$$\min_{\mathbf{y}, \mathbf{w}, \mathbf{N}} R(\mathbf{w}) + \sum_{i=1}^m \left[ c_2 L(\mathbf{w}, \mathbf{N}, \mathbf{x}_i^A, y_i) + c_1 y_i \Delta(\mathbf{N}, \mathbf{x}_i^G) \right]$$



Primitive

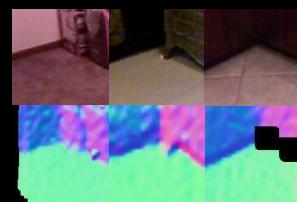
$\mathbf{w}$



$\mathbf{N}$



$\mathbf{y}$



Patch

$\mathbf{x}_i^A$



$\mathbf{x}_i^G$



# Objective

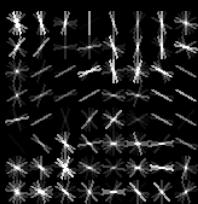
Solved with iterative method similar to  
block-coordinate-descent.



$$\min_{\mathbf{y}, \mathbf{w}, \mathbf{N}} R(\mathbf{w}) + \sum_{i=1}^m \left[ c_2 L(\mathbf{w}, \mathbf{N}, \mathbf{x}_i^A, y_i) + c_1 y_i \Delta(\mathbf{N}, \mathbf{x}_i^G) \right]$$

Primitive

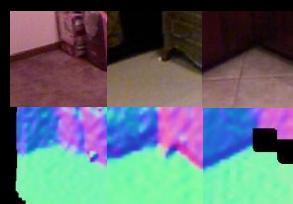
$\mathbf{w}$



$\mathbf{N}$



$\mathbf{y}$



Patch

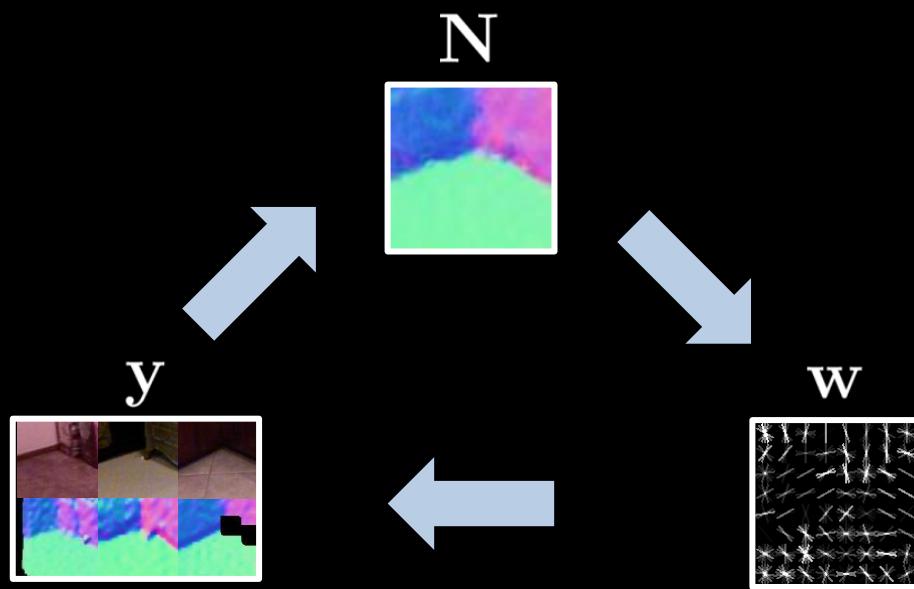
$\mathbf{x}_i^A$



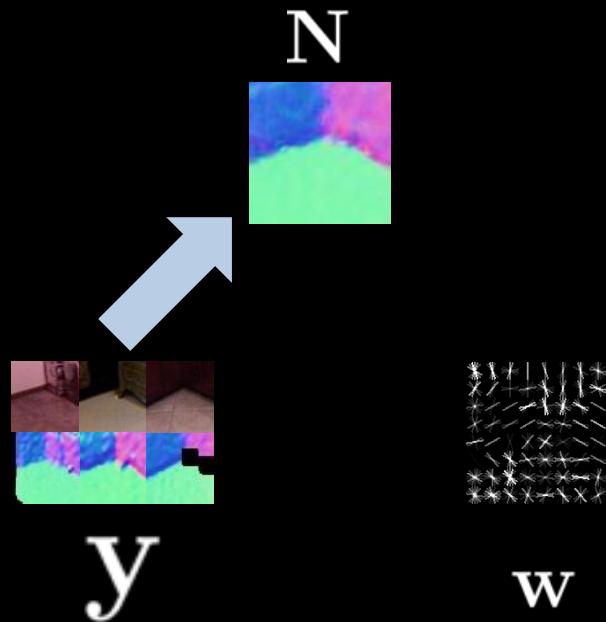
$\mathbf{x}_i^G$



# Iterative Procedure

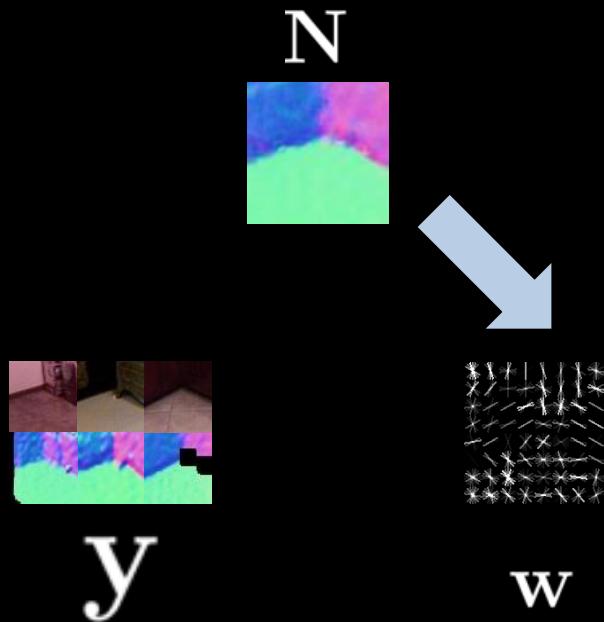


# Iterative Procedure



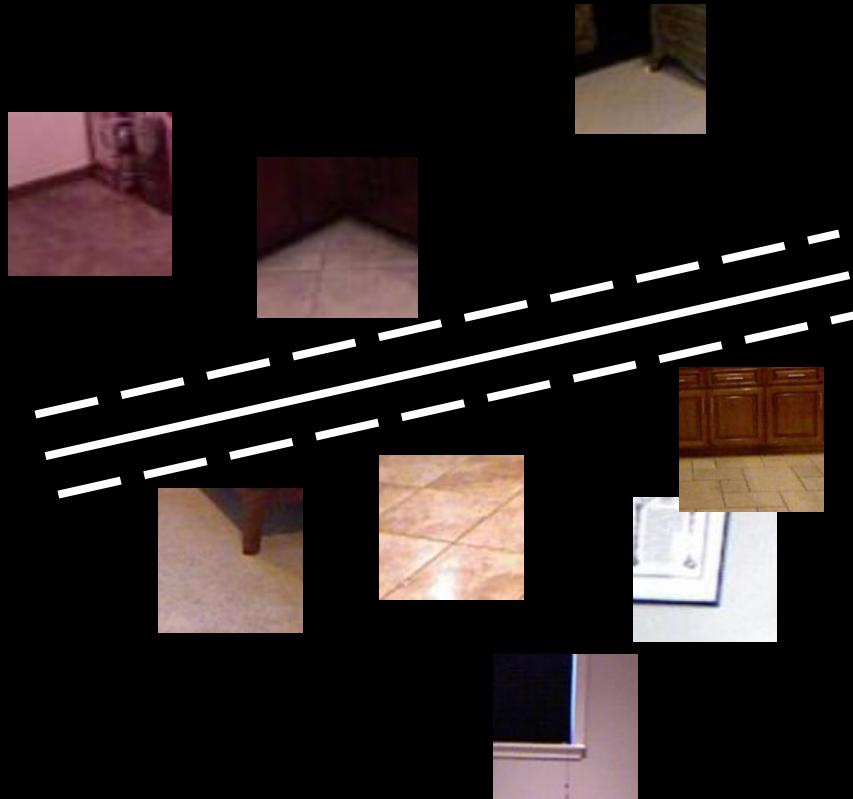
$$\text{Avg} \left( \begin{array}{c} \text{Image} \\ \text{Image} \\ \text{Image} \end{array} \right)$$

# Iterative Procedure

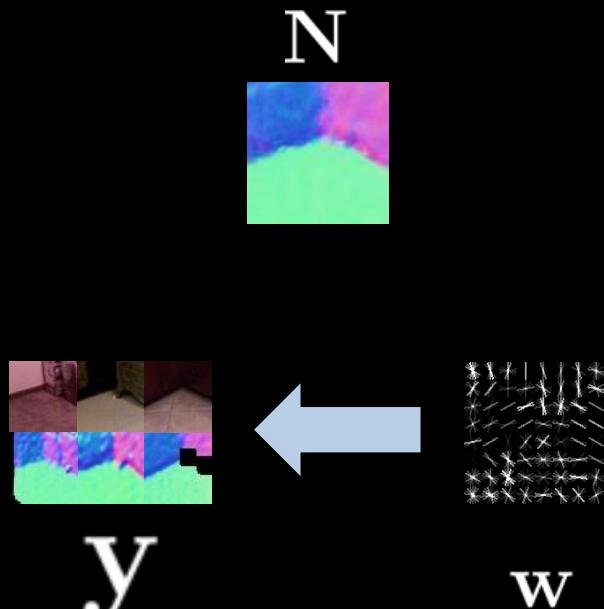


Cluster  
Instances

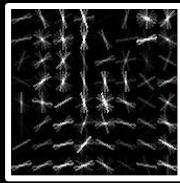
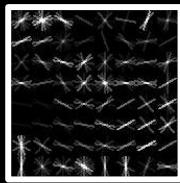
Patches  
Geometrically  
Dissimilar to  $N$



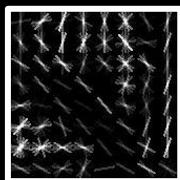
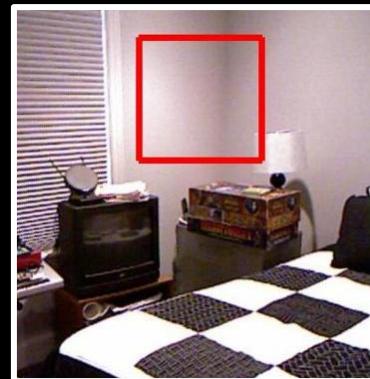
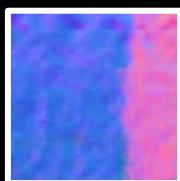
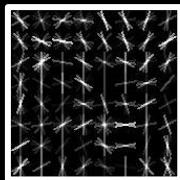
# Iterative Procedure



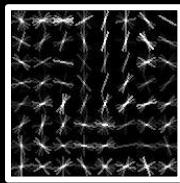
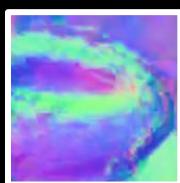
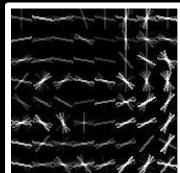
# Primitives



# Primitives

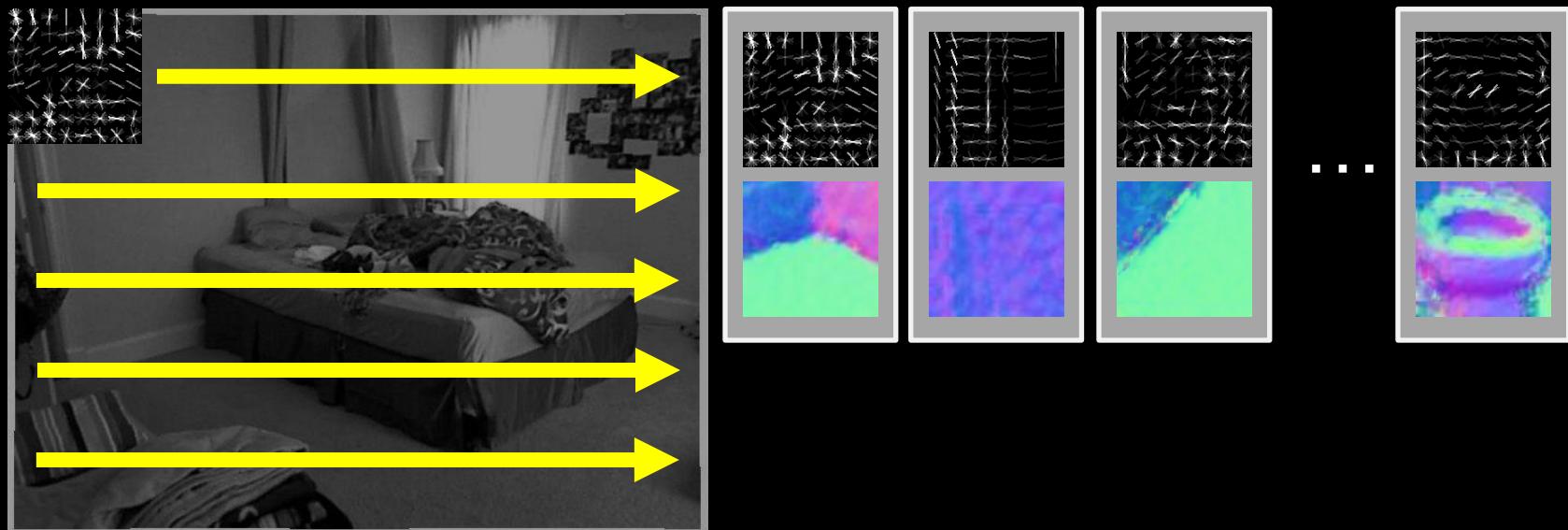


# Primitives

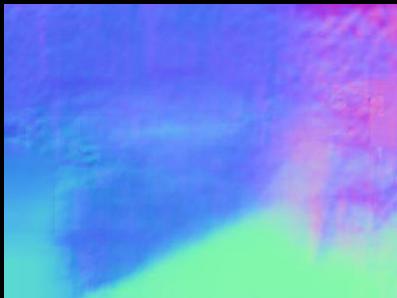
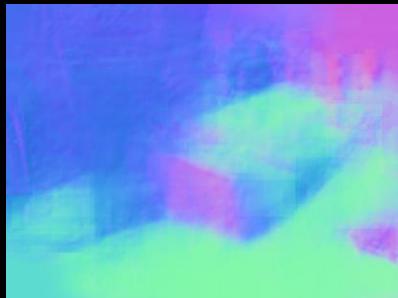
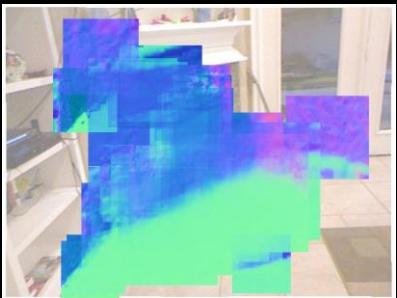
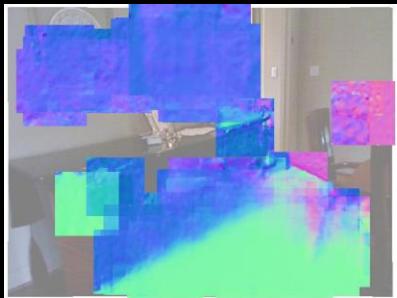
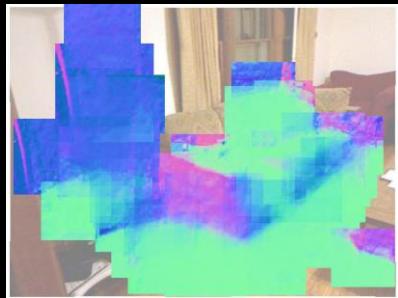


# Test-time Correspondence

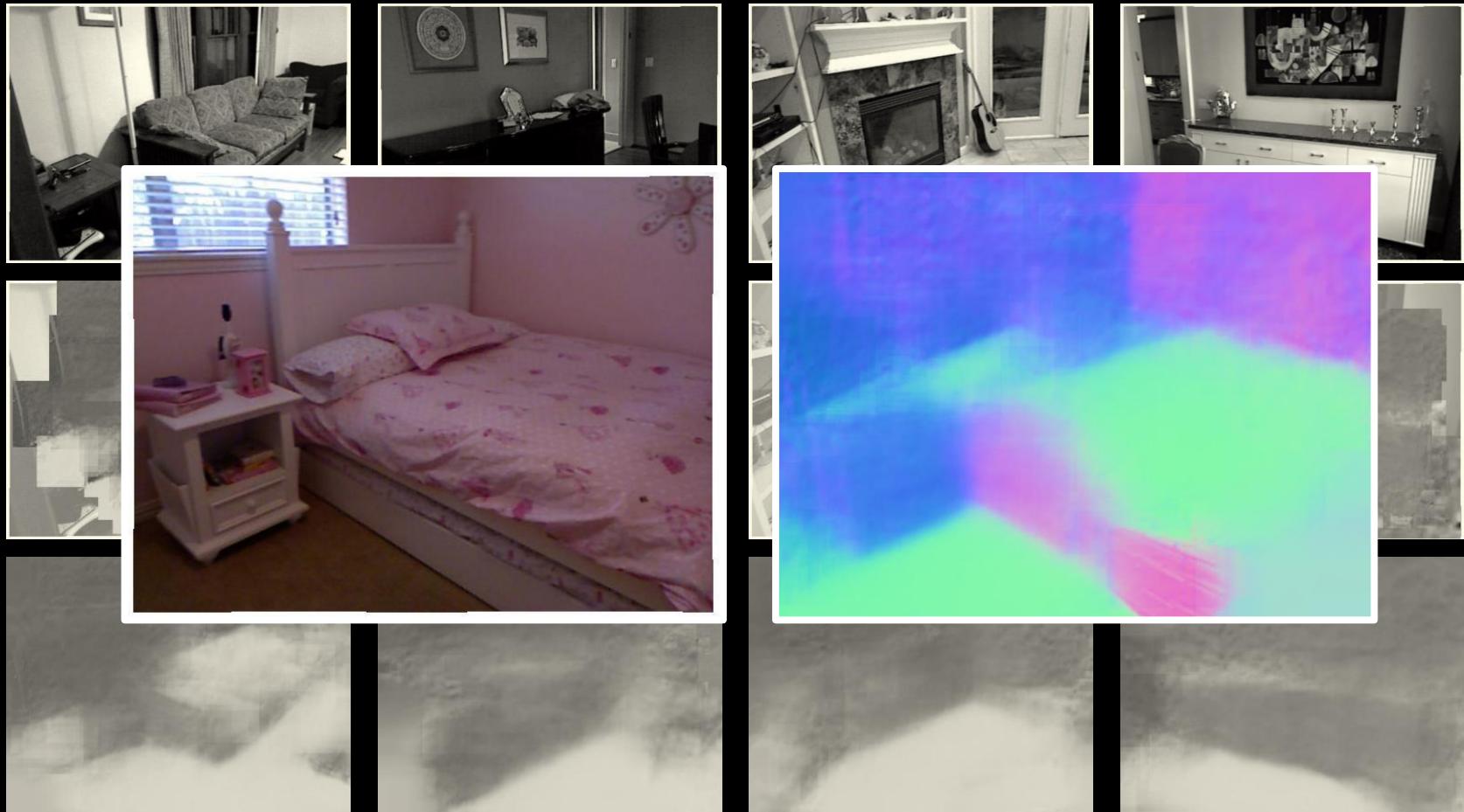
Correspondence via detection



# Results



# Results



Fouhey et al ., ICCV 2013, ECCV 2014

# Conclusions

Introduced Data-Driven 3D Scene Understanding

Full 3D Models



RGBD Data



Two Main Problems:

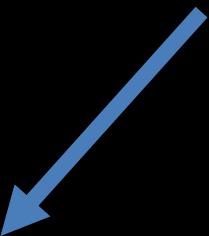
1. Correspondence
2. Representation Transfer

# Available Datasets

- RMRC (NYU + SUN3D)
- NYU v2:
  - 1449 Pairs + semantic labels + raw videos
- SUN3D
  - 415 Sequences in large spaces + raw videos
- Berkeley 3D Object
  - 849 images + bounding boxes
- MSR-V3D
  - 177 sequences

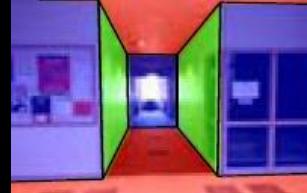
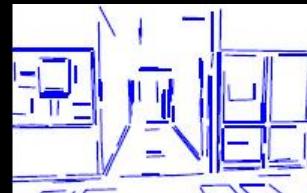
# Available Code

Hoiem et al., Geometric Context,  
Saxena et al., Make 3D

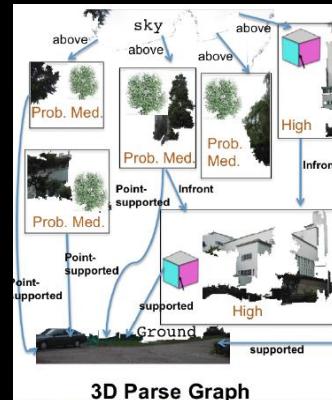


Region labels

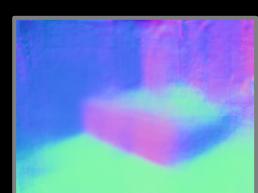
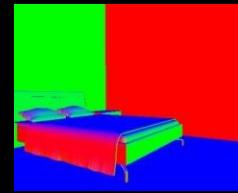
+ Boundaries  
and objects



Stronger geometric  
constraints from  
domain knowledge



Volumetric +  
functional  
constraints



Data-  
driven  
3D

# Available Code

Hoiem et al., Occlusion boundaries

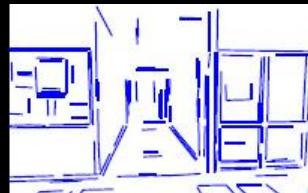
Hoiem et al., Putting objects in perspective



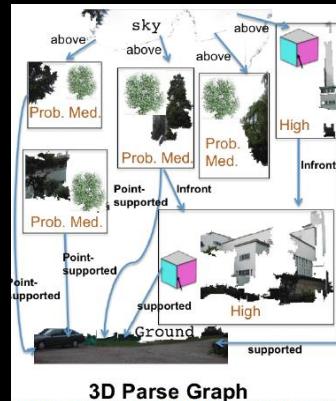
Region labels



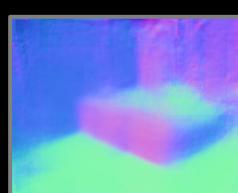
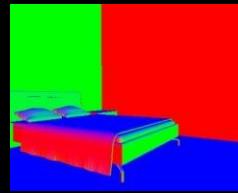
+ Boundaries  
and objects



Stronger geometric  
constraints from  
domain knowledge



Volumetric +  
functional  
constraints



Data-  
driven  
3D

# Available Code

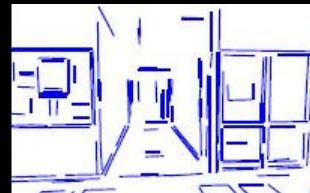
Lee et al., Orientation Maps  
Hedau et al., Room-fitting



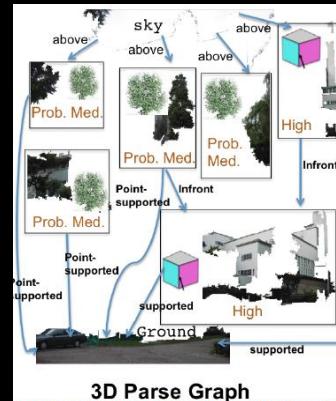
Region labels



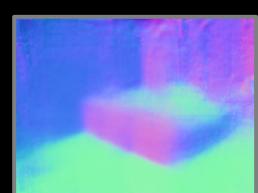
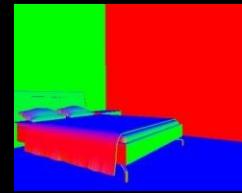
+ Boundaries  
and objects



Stronger geometric  
constraints from  
domain knowledge



Volumetric +  
functional  
constraints



Data-  
driven  
3D

# Available Code

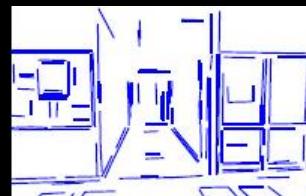
Gupta et al., Blocks World  
Choi et al., Geometric Phrases



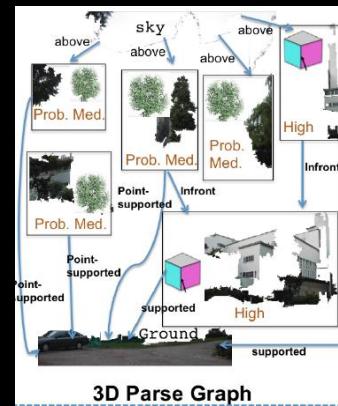
Region labels



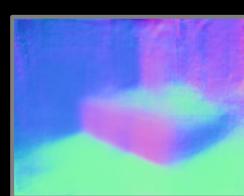
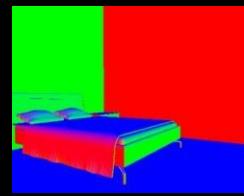
+ Boundaries  
and objects



Stronger geometric  
constraints from  
domain knowledge



Volumetric +  
functional  
constraints



Data-  
driven  
3D

# Available Code

Karsch et al., Depth-Transfer

Fouhey et al., Data-Driven 3D Primitives

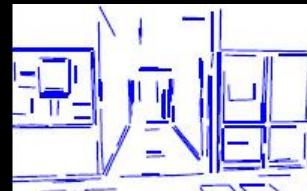
Aubrey et al., Seeing 3D Chairs



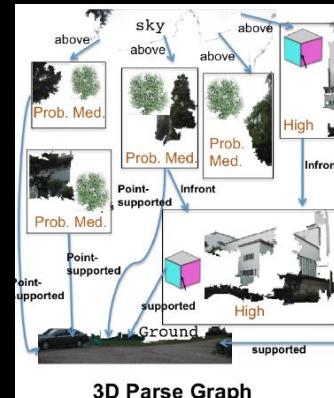
Region labels



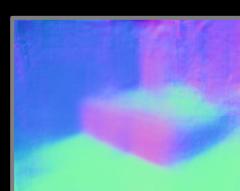
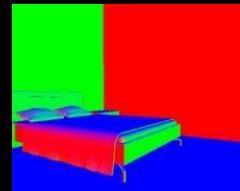
+ Boundaries  
and objects



Stronger geometric  
constraints from  
domain knowledge



Volumetric +  
functional  
constraints



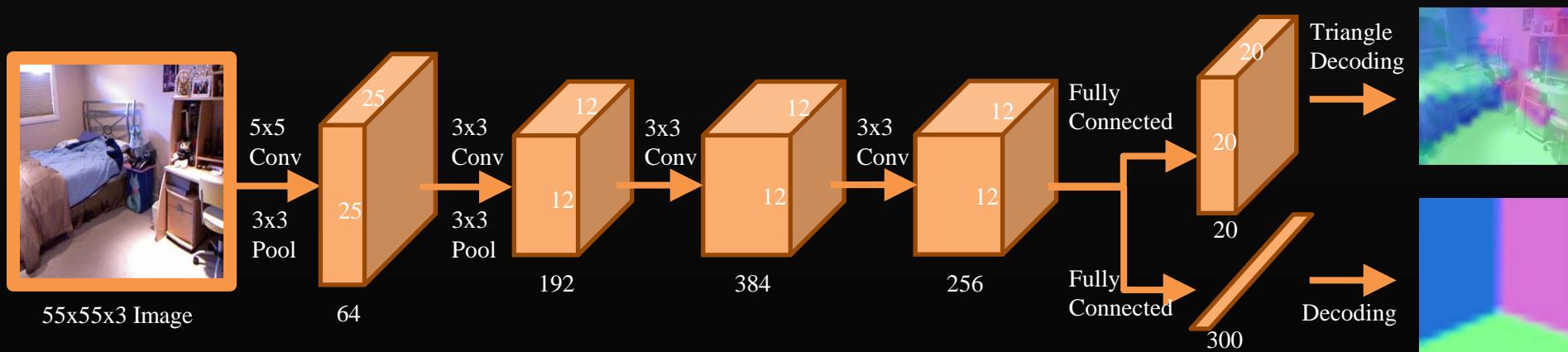
Data-  
driven  
3D

# DEEP LEARNING APPROACH



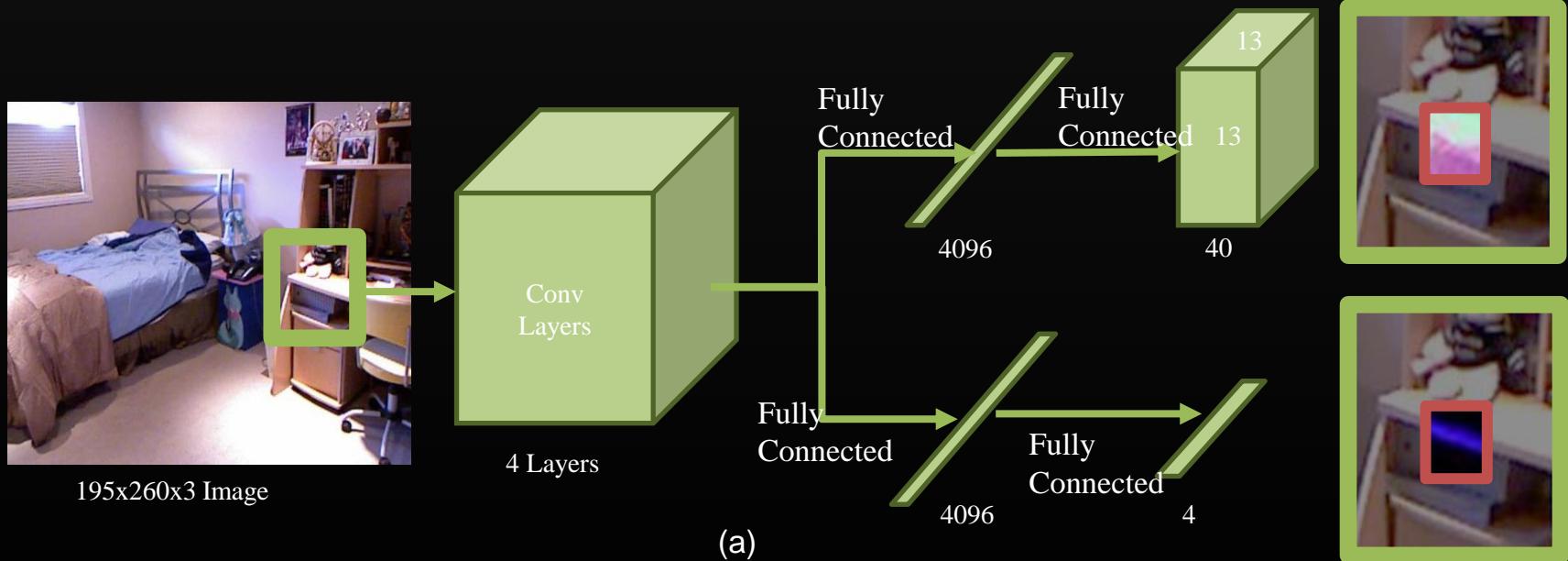
- Possible Approach: Regression to 3 scalars followed by normalization
  - Tends to over-smoothen.
- Formulate regression as classification problem.
  - Cluster output space into code-words: individual class.
- Two networks
  - Capturing global cues and global outputs
  - Capturing local cues and local outputs.

# GLOBAL NETWORK



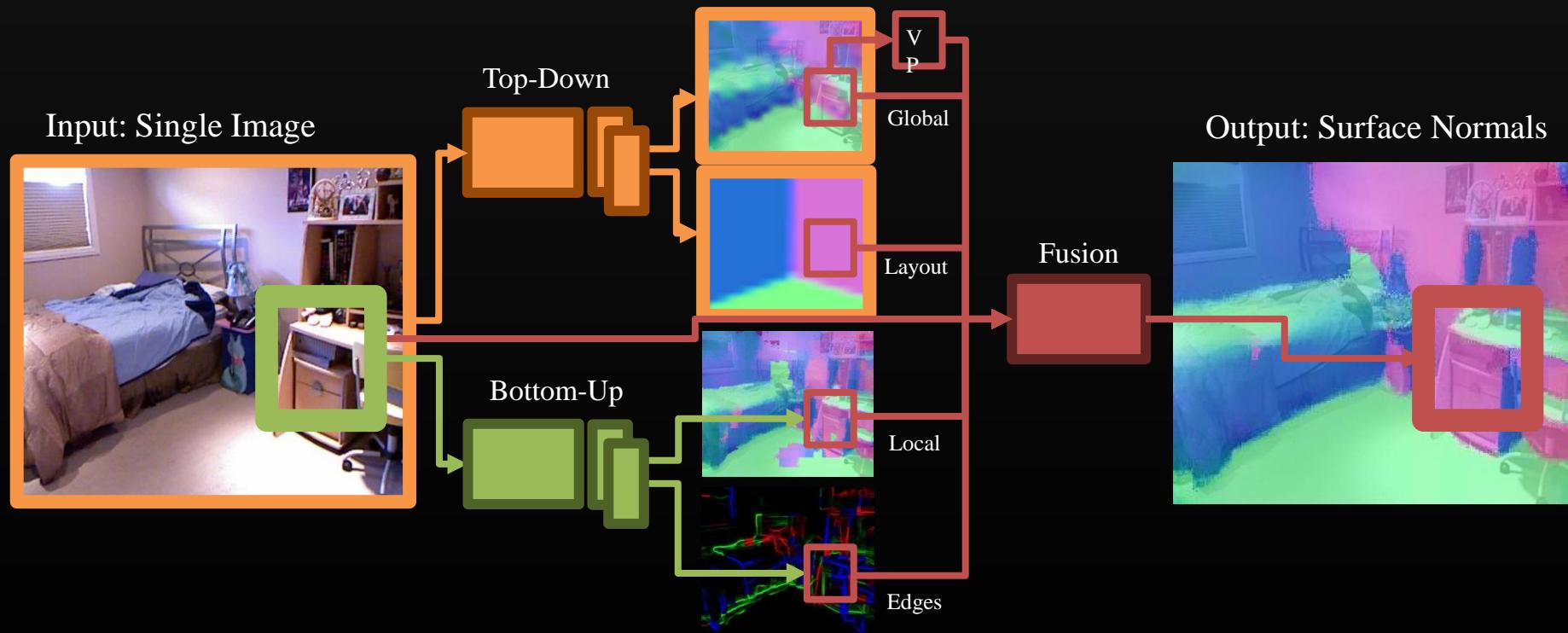
- Predicts  $20 \times 20 \times 20$  Coarse Output
- Room Layout: 1/300 classes
- Dual Prediction leads to better learning

# LOCAL NETWORK

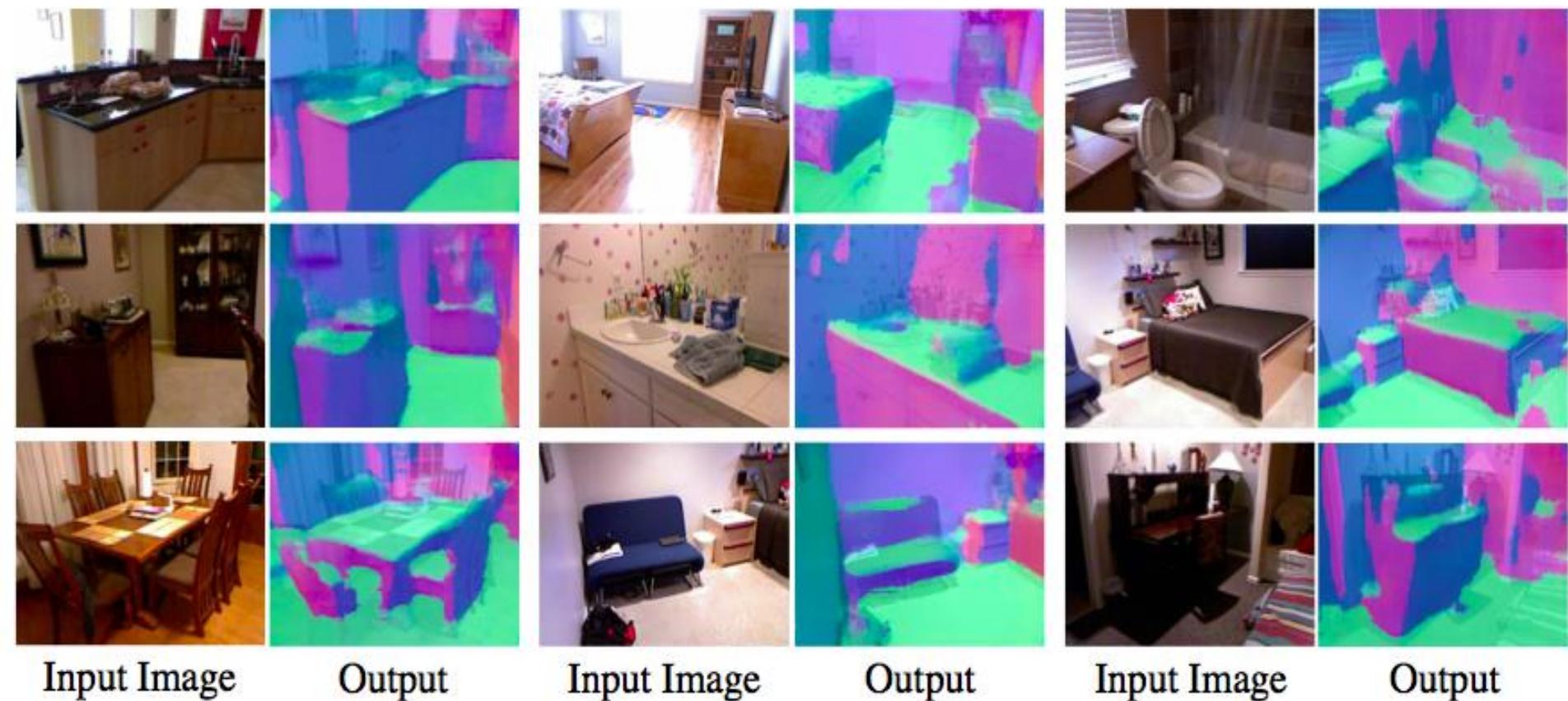


- Input:  $55 \times 55 \times 3$  Pixel Windows
- Predicts  $13 \times 13 \times 40$  Output
- Edge Labels: 4 classes

# PUTTING IT TOGETHER

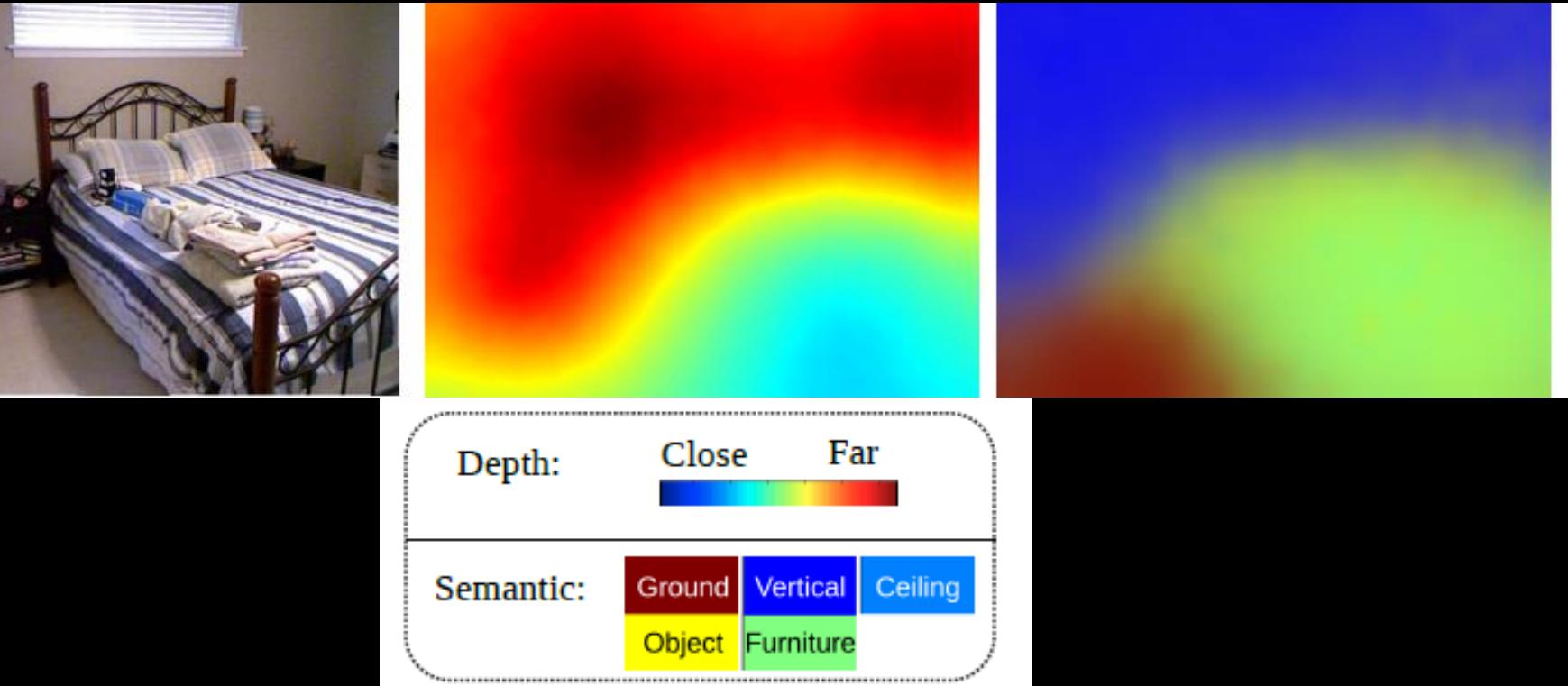


- Train a Separate Fusion Network which takes input all the outputs of Top-Down, Bottom-up Networks



X. Wang, D. Fouhey, A. Gupta. Designing Deep Networks for Surface Normal Estimation. CVPR 2015.

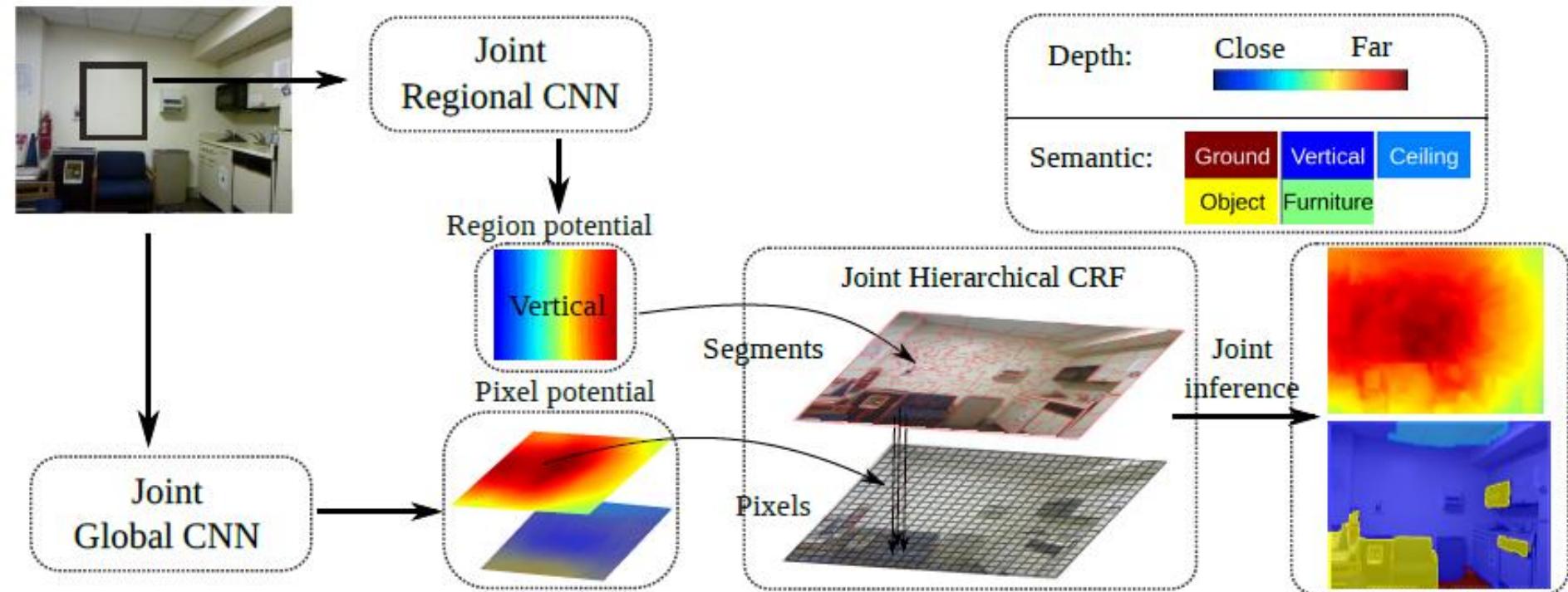
# Closing the loop: Geometric + Semantic (+ CNNs)



$$\text{loss}(\mathcal{X}, \mathcal{X}^*) = \frac{1}{n} \sum_{i=1}^n (\log d_i - \log d_i^*)^2 + \lambda_l \frac{-1}{n} \sum_{i=1}^n \log(P(l_i^*)),$$

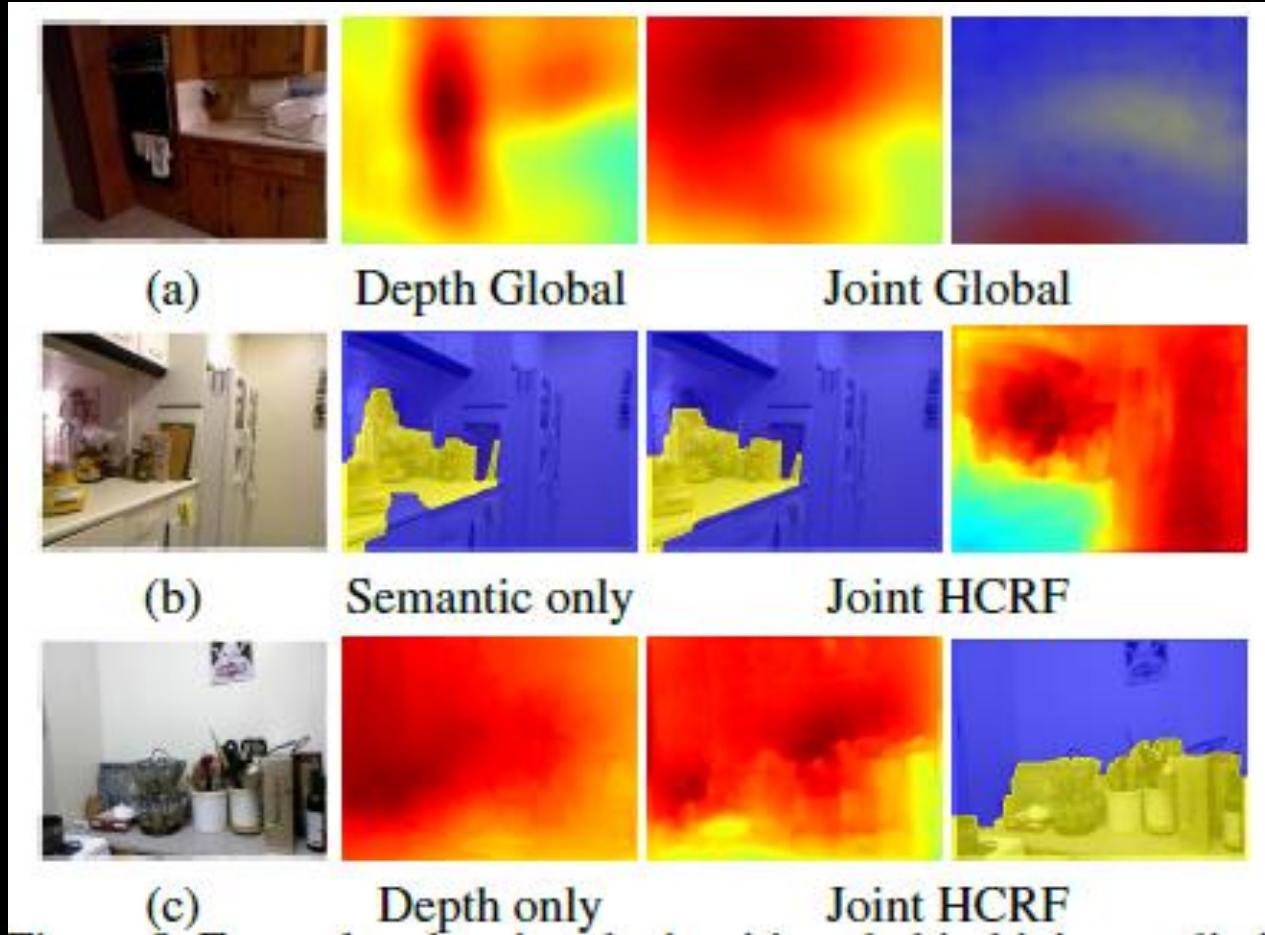
$$\text{and } P(l_i^*) = \exp(z_{i,l_i^*}) / \sum_{l_i} \exp(z_{i,l_i}),$$

# Closing the loop: Geometric + Semantic (+ CNNs)



P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, A. Yuille. Towards Unified Depth and Semantic Prediction from a Single Image. CVPR 2015..

# Closing the loop: Geometric + Semantic (+ CNNs)



P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, A. Yuille. Towards Unified Depth and Semantic Prediction from a Single Image. CVPR 2015..

# Big Questions

- How to estimate geometric properties from an image?
- How to incorporate geometric constraints and which ones?
- How to combine reasoning tools with statistical classification/regression tools?
- How to use large-scale 3D data (3D models, kinect)
- How to combine semantic and geometric representations?
- How to combine with other 3D estimation methods?