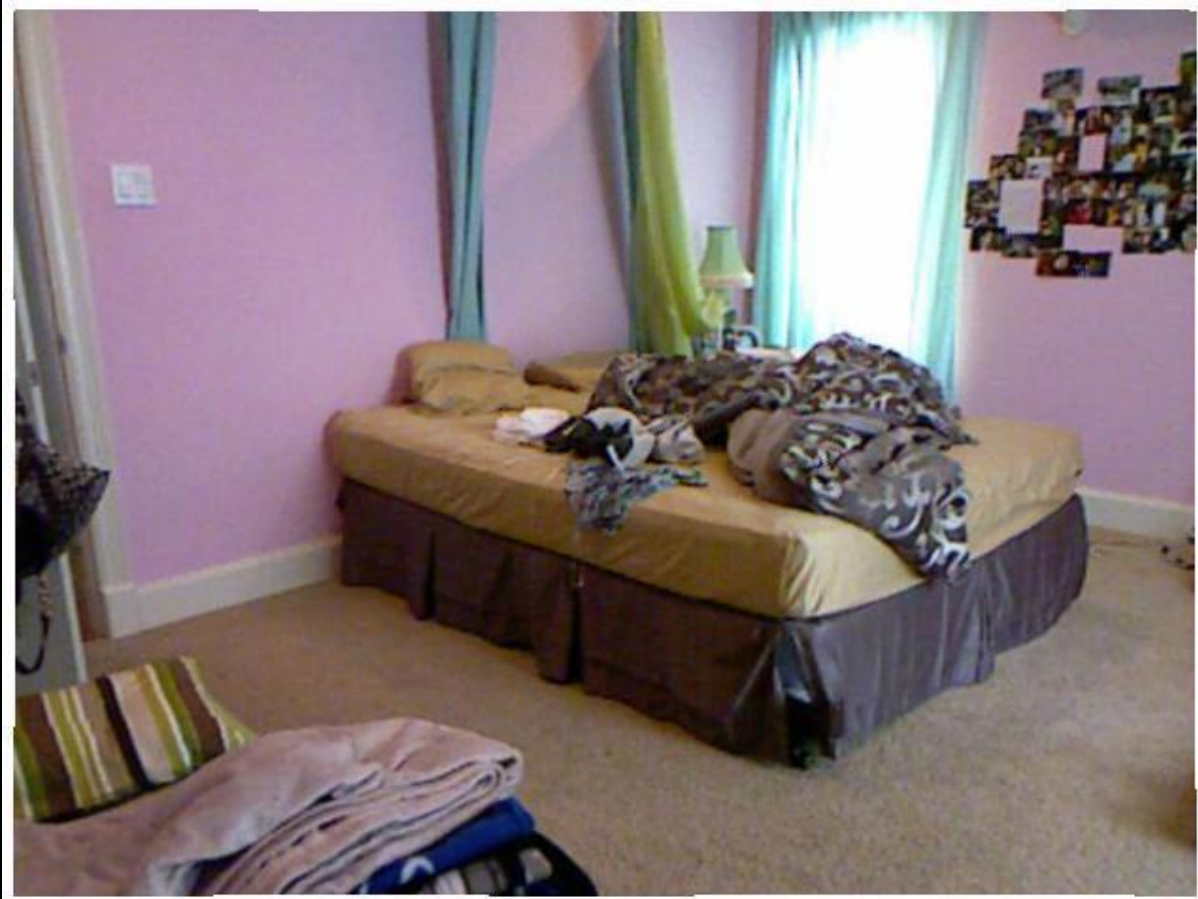# Data-Driven 3D

David Fouhey

# Recap



**Martial**

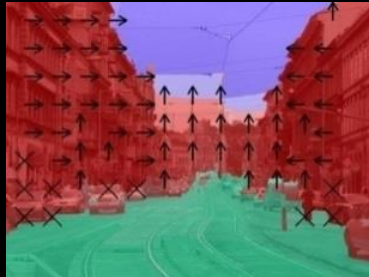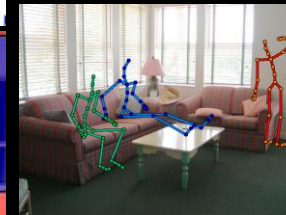**Derek**

**Abhinav**

**David**

3D Parse Graph

**Introduction, Applications, History**

**Region labels +Boundaries +Objects**

**Stronger geometric constraints**

**Volumetric + Functional Constraints**

**Data-Driven 3D**

# Data-Driven Interpretation



Every image that can be seen has been seen before (approximately)
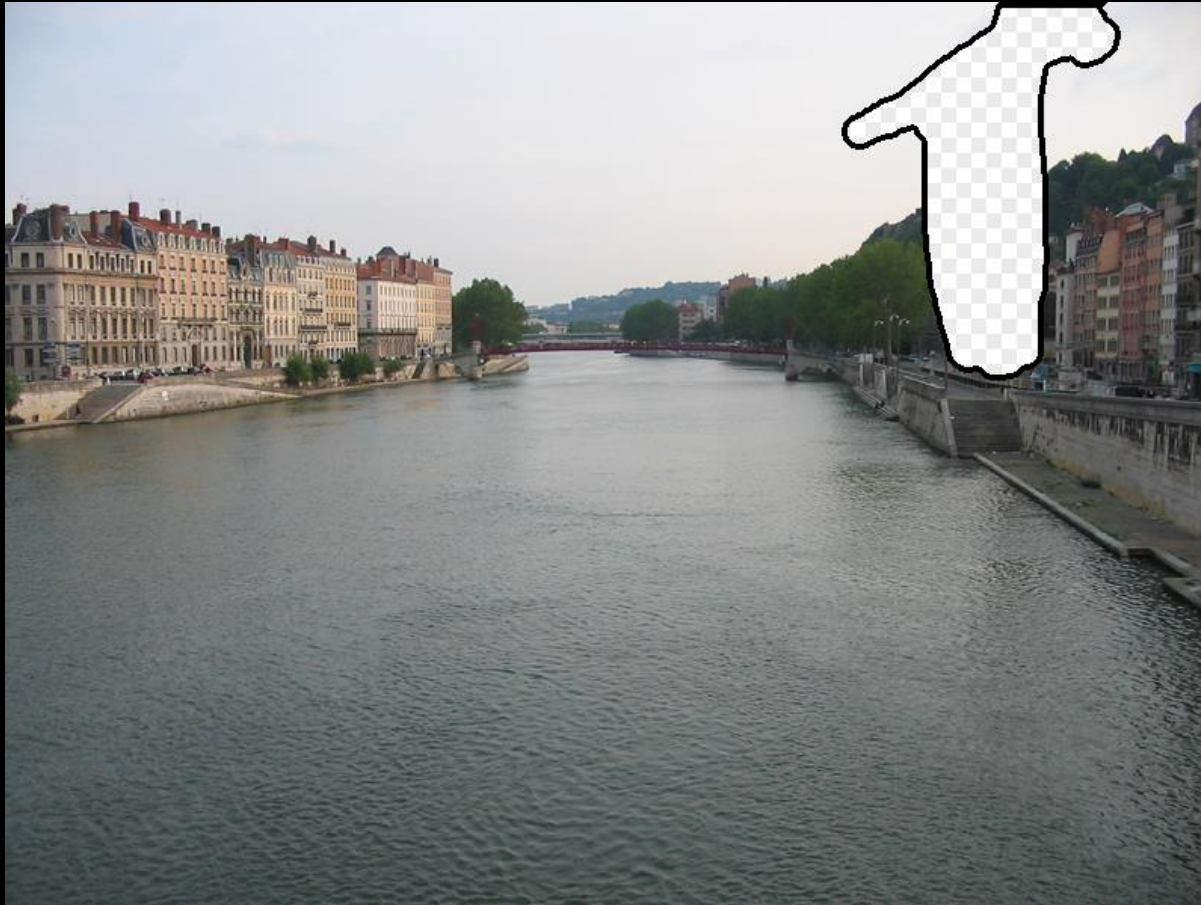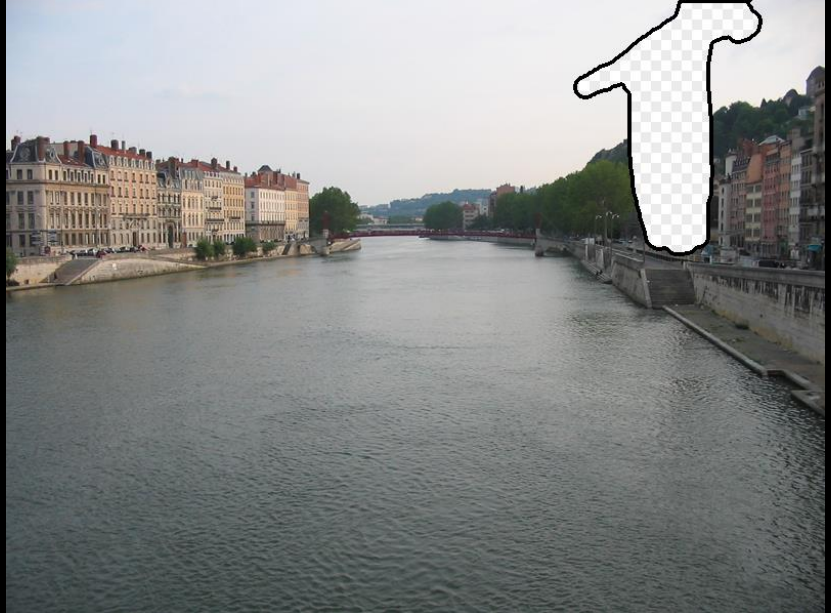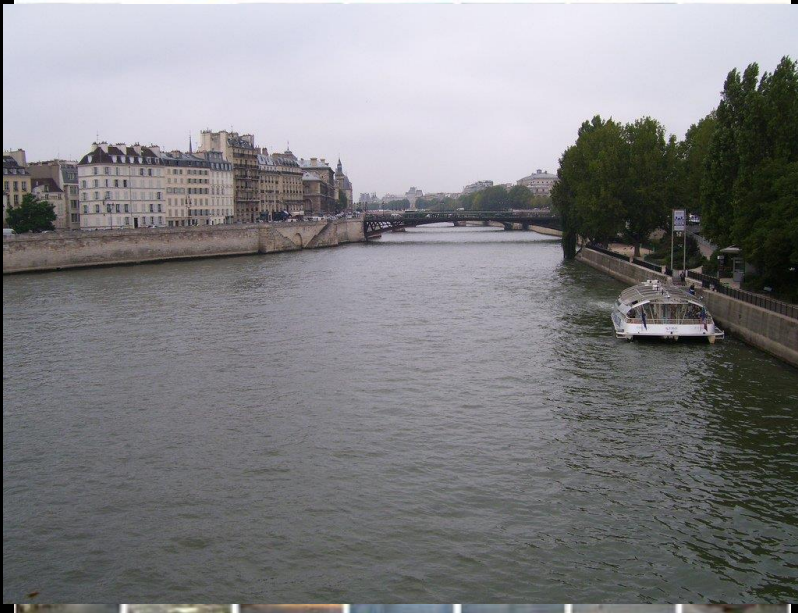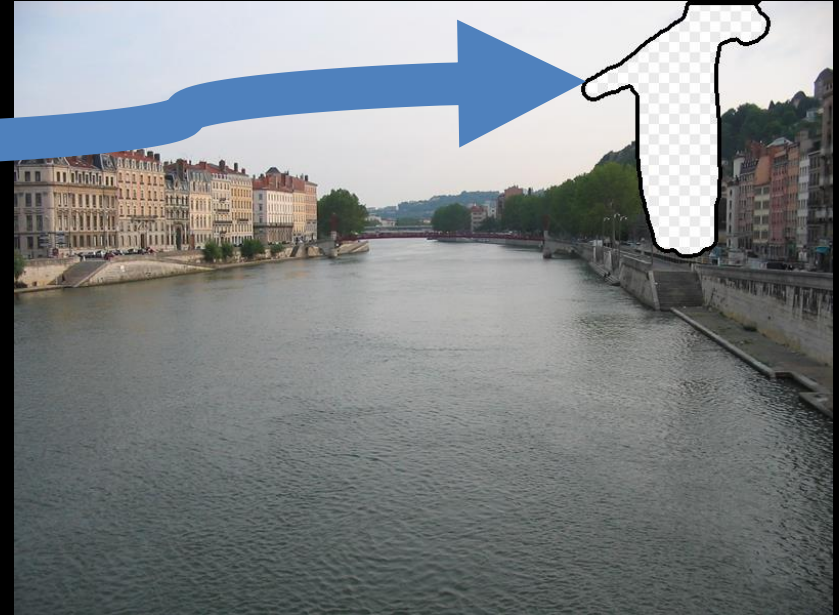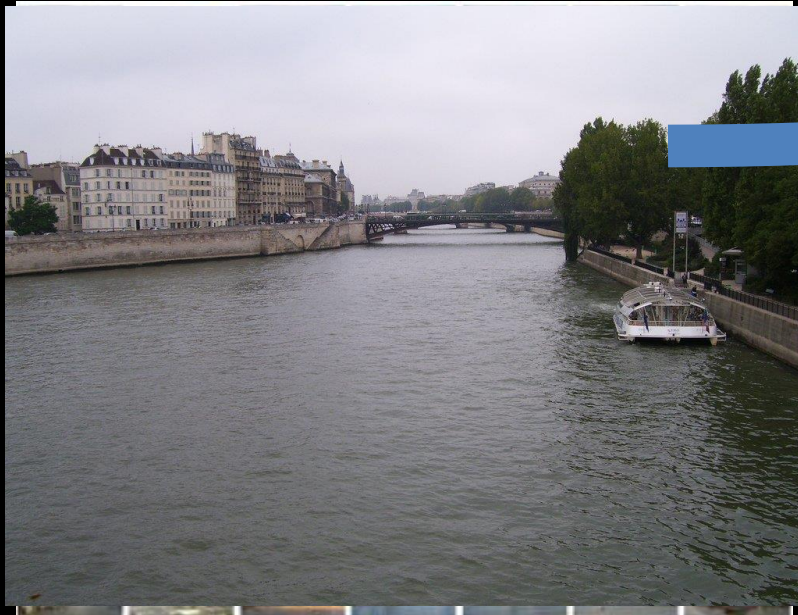
# Data-Driven Interpretation



· · ·

Every image that can be seen has been seen before (approximately)
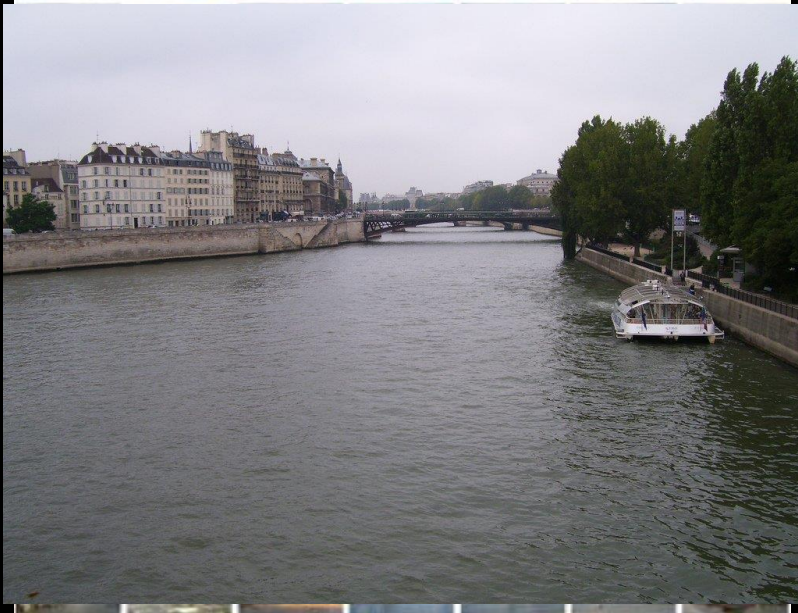
# Data-Driven Interpretation



• • •

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

Hays and Efros 2007    6

# Data-Driven Interpretation



• • •

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

# Data-Driven Interpretation

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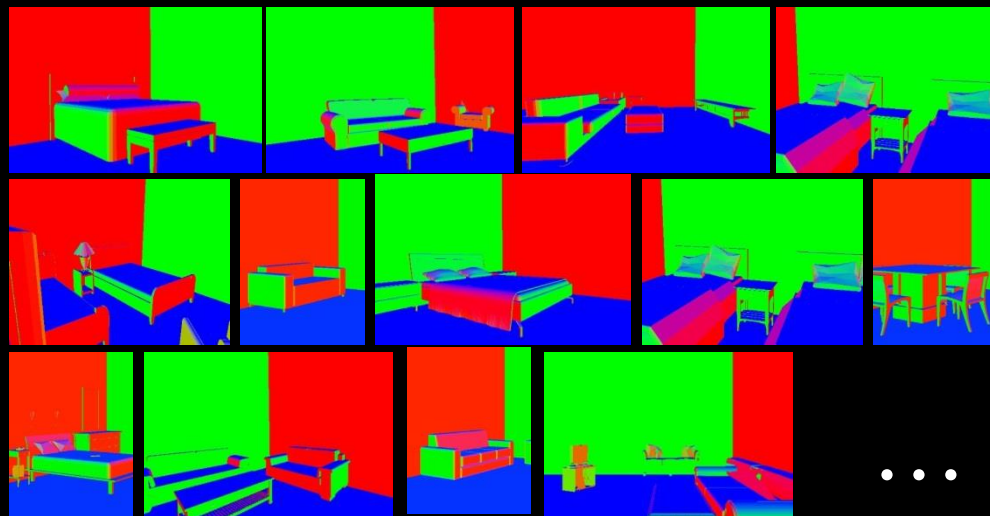
# Advantages

Volumetric
Interpretation
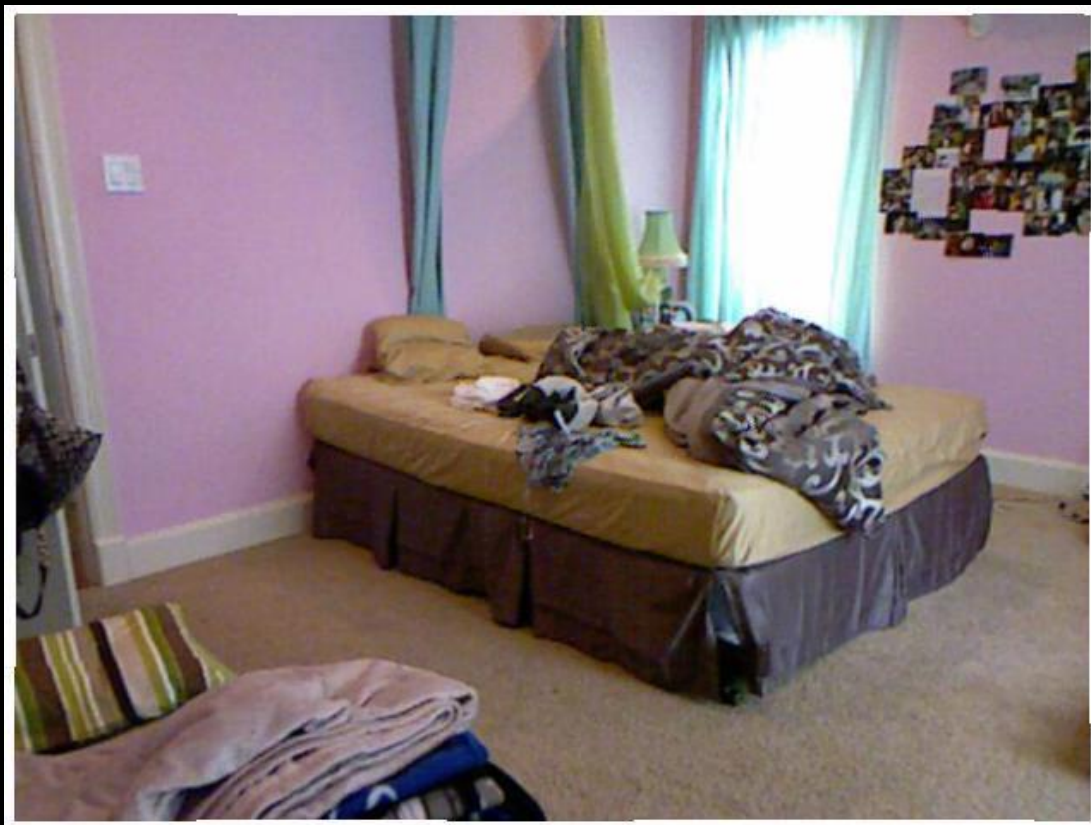
Interpretation by
3D Models
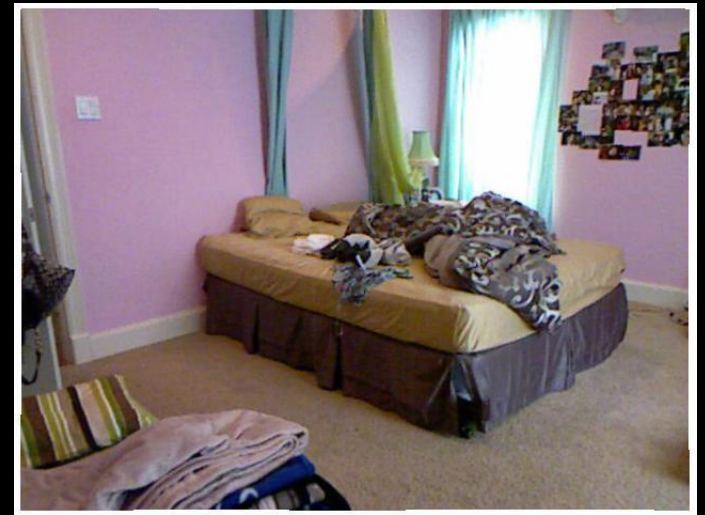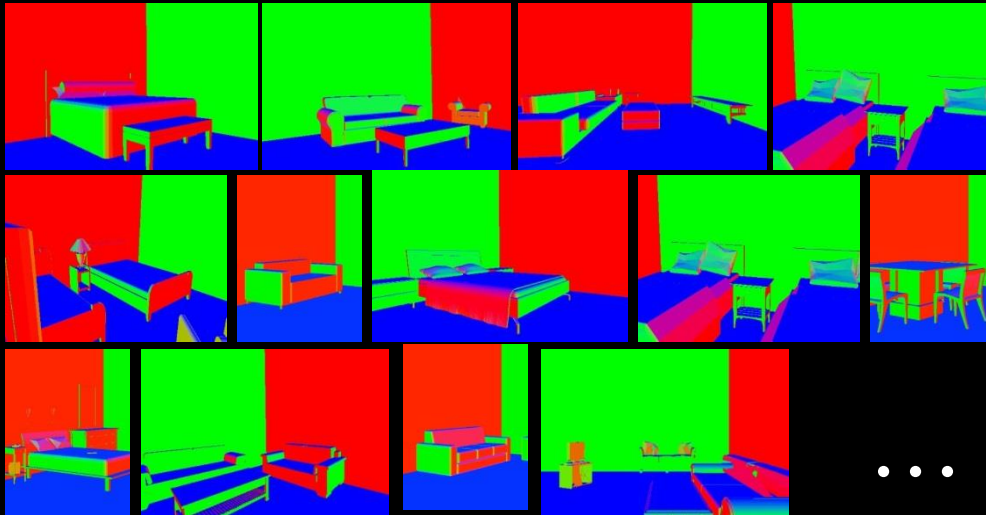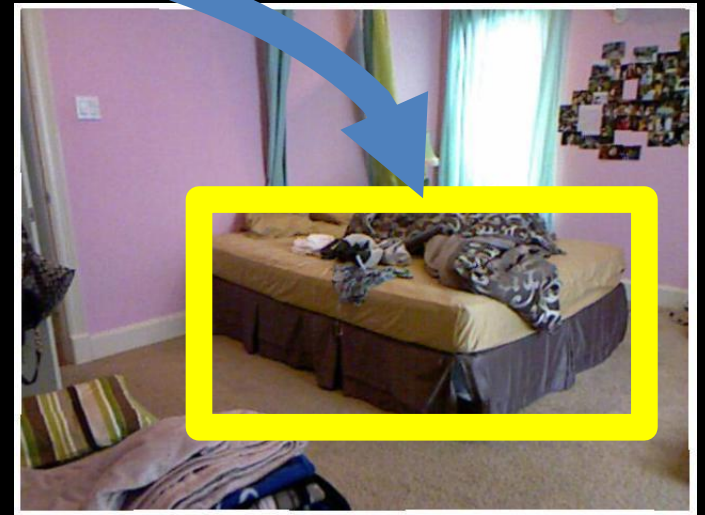
# Sources

## 3D Model Databases



## Kinect Databases
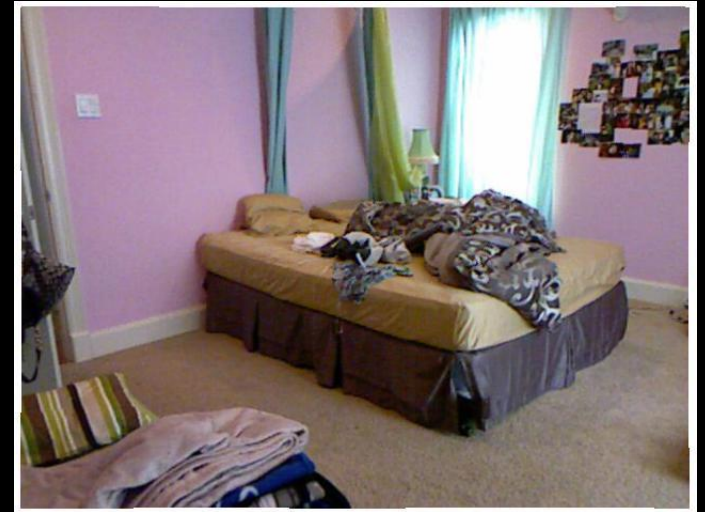


· · ·

· · ·
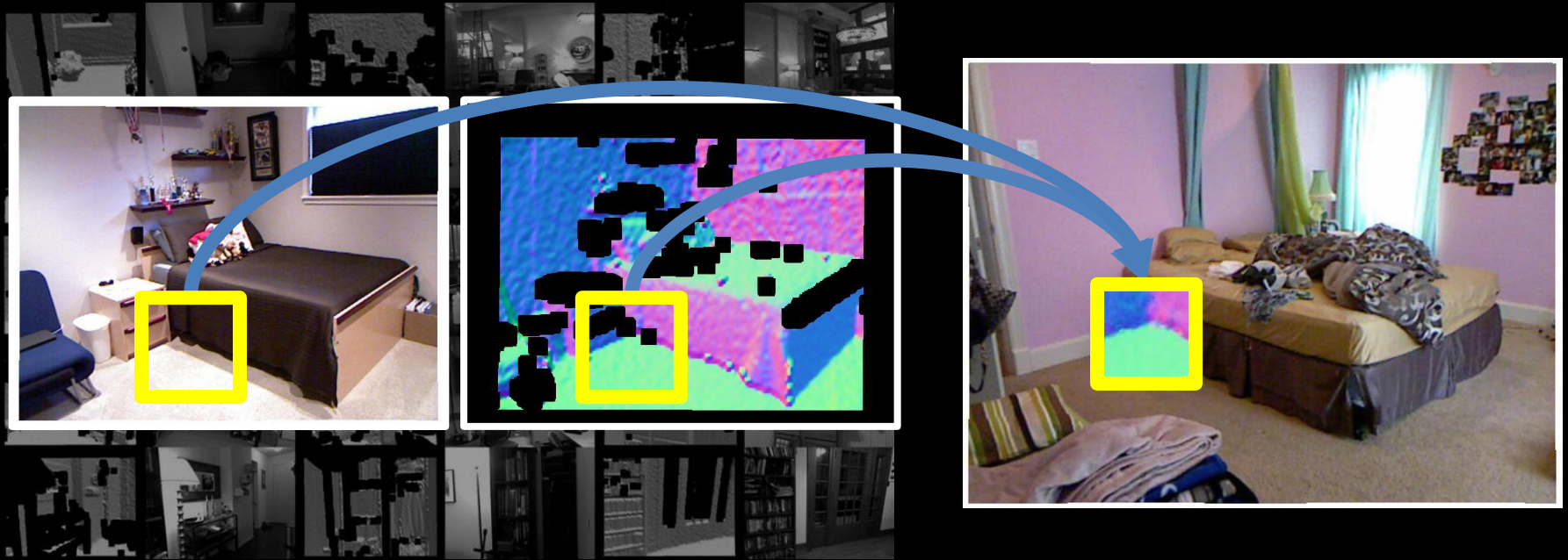
# Goal

# Goal



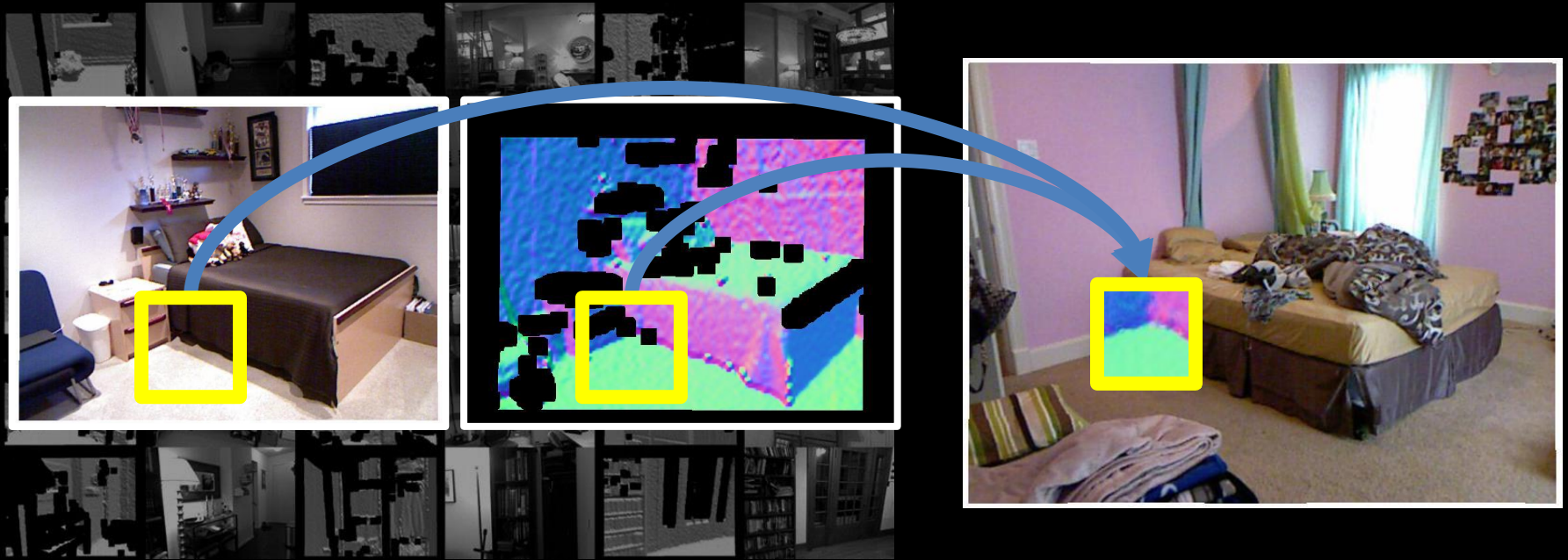. . .

# Goal

# Goal

# Goal

# Goal



How do you:
    (a) establish correspondence?
    (b) transfer representations?

# Overview

## 1. How to use 3D models



## 2. How to use the Kinect

# Why 3D Models

| Object Detector | Segmentation | 3D Model |
|:---:|:---:|:---:|

# Why 3D Models



Input

Top 2D (GIST) Match

Top 3DNN Match

# 3D Models

- Advantages:
  - Full 3D – can be rendered and modified
  - Precise models may exist (e.g., IKEA)

- Disadvantages:
  - No corresponding natural color image (untextured or missing)

# General Approach

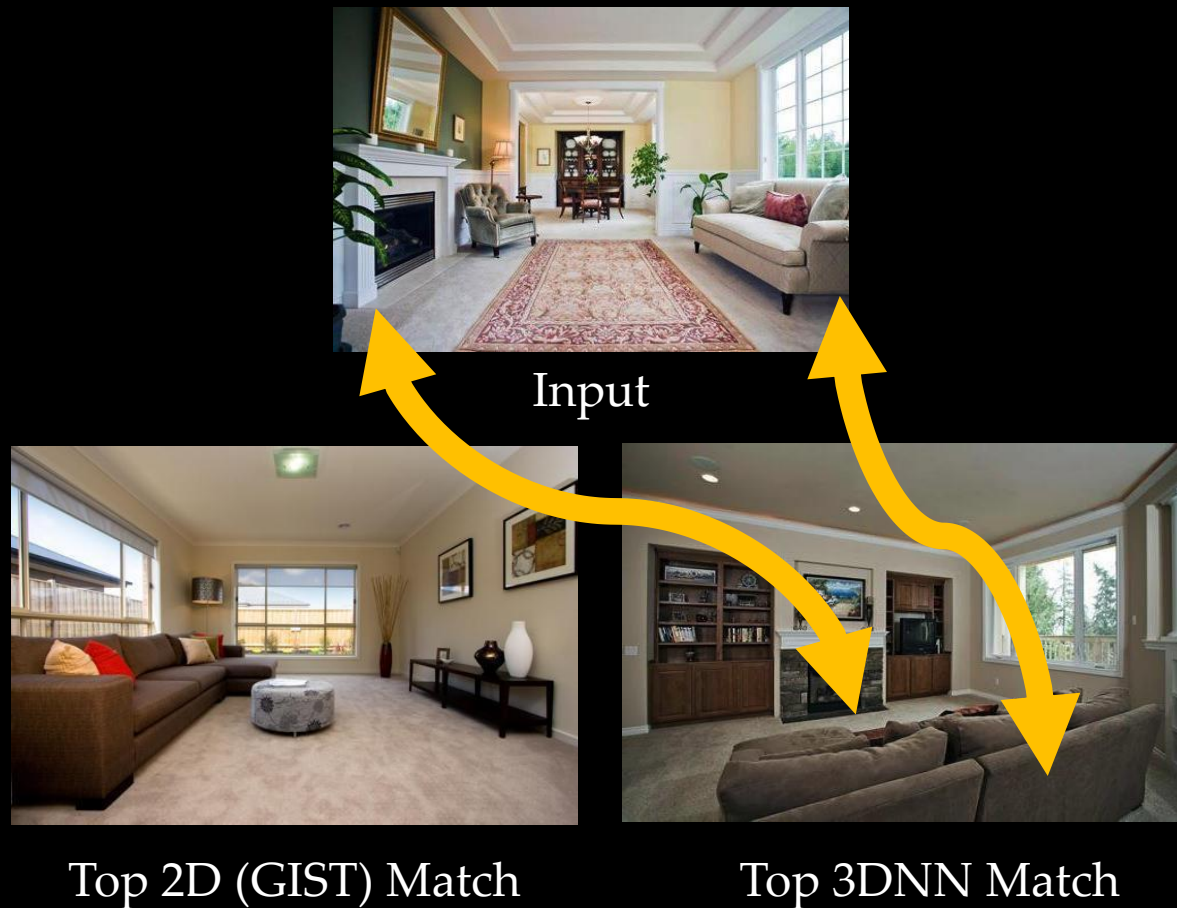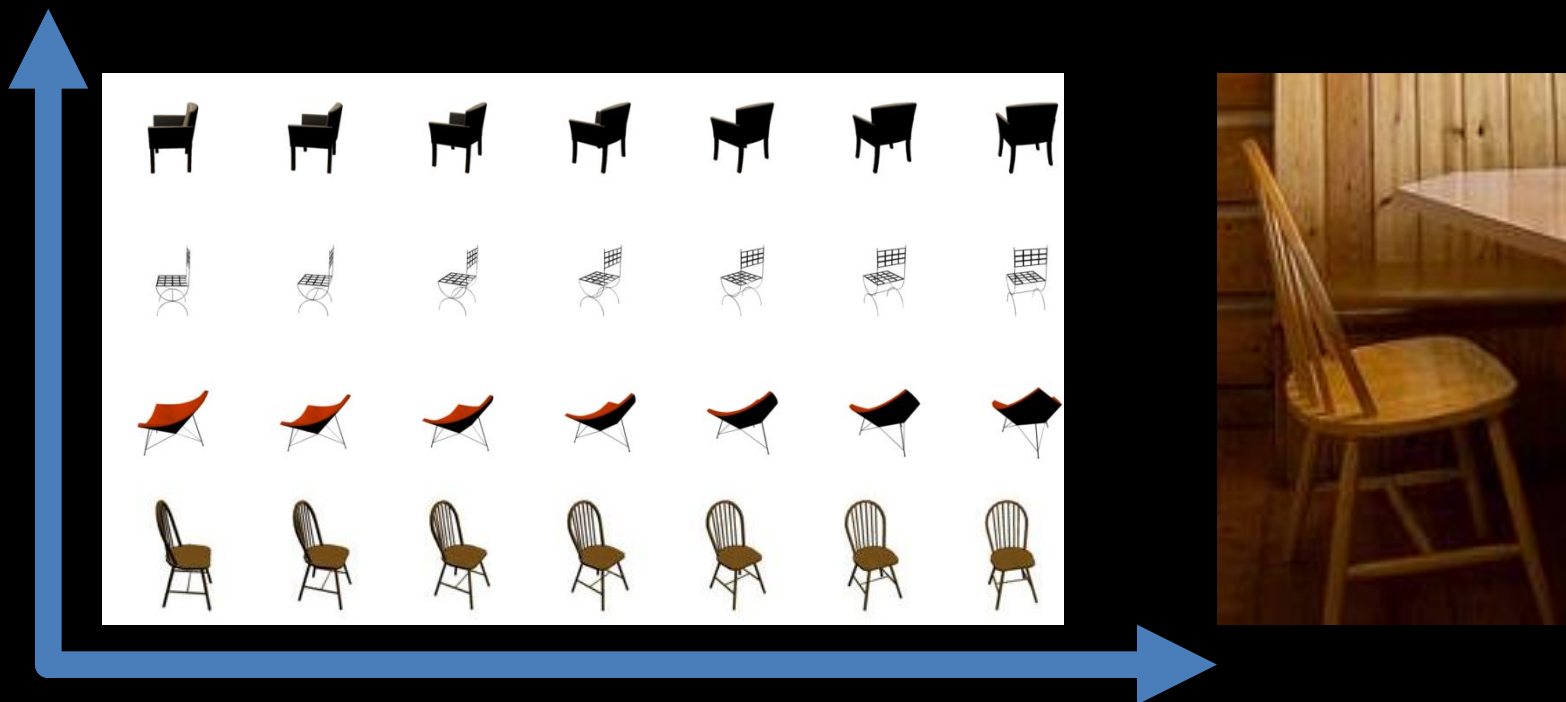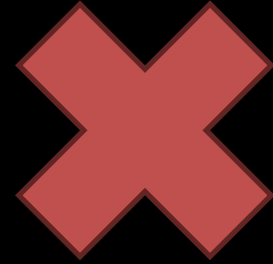Search over model and viewpoint

# Primary Question



Does it match?

~1400 models

~60 viewpoints

# Primary Question



Does it match?
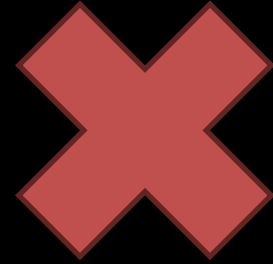
~1400 models

~60 viewpoints

# Primary Question

Does it match?

~1400 models

~60 viewpoints

24

# Difficulties

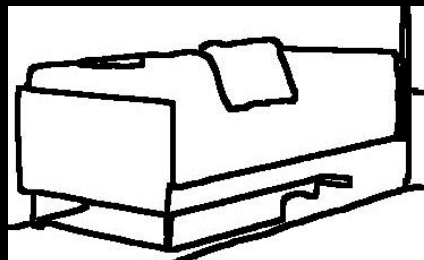|  | Rendered | Natural |
|---|---|---|
| |  |  |
| Texture | NO | YES |
| Occlusion | NO | YES |
| Background | Fake | Natural |

# Cross-Domain Matching

Goal: bring image and model into common representation

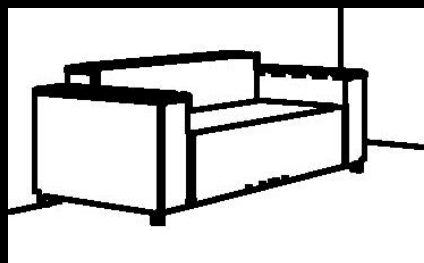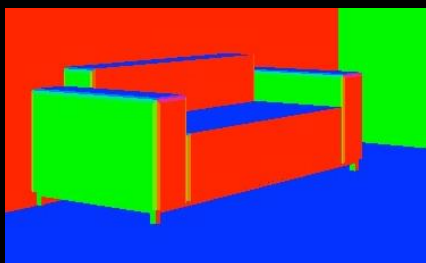# Chamfer Matching

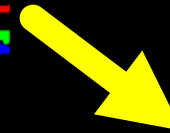Assumption: edges in 3D are edges in 2D
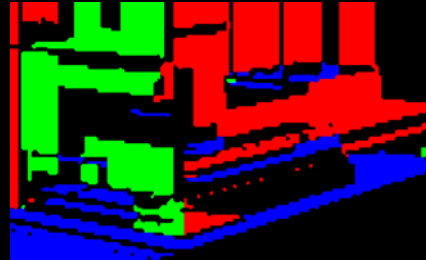


Image

3D Model

Match?

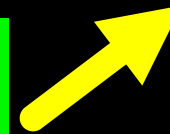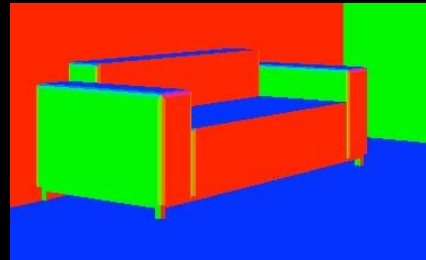Satkin et al., 2012, 2013, 2014; Lim et al., 2013; Ramnath et al., 2014;

# Domain-Invariant

Assumption: can estimate 3D property from 2D

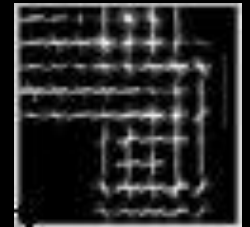Image

3D Model

Match?

Satkin et al., 2012, 2013, 2014;

# Domain-invariant "Images"

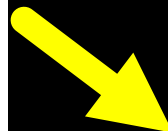Assumption: edges in 3D are edges in 2D
Apply standard features/techniques

Image

3D Model

# Masking Features

Assumption: only issue is background

HOG    Classifier    Masked classifier



HOG mask

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

# Searching Hypotheses

Render object parts



models

viewpoints

Aubry et al., 2014
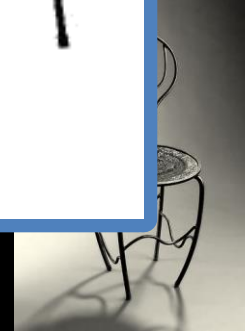
Matches generate proposals



Lim et al., 2013
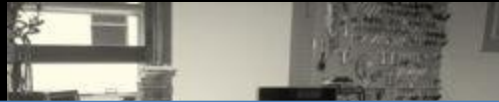
# Results

# Results

# Results

# Issues

What's this?

# Issues

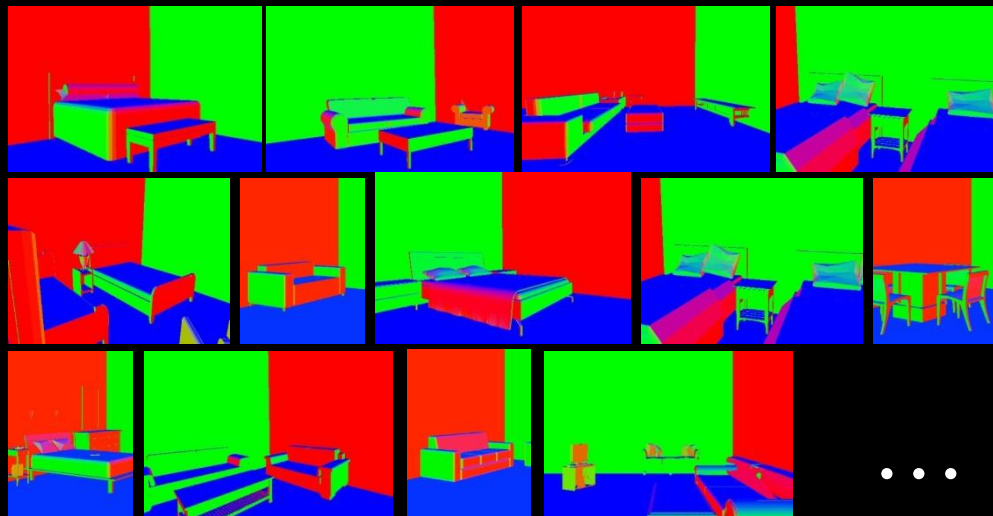Recognition and pose estimation is <u>hard</u>, but made easier by seeing the rest of the room.

# 2D-3D Scene Matching

3D Model Database



Input



· · ·
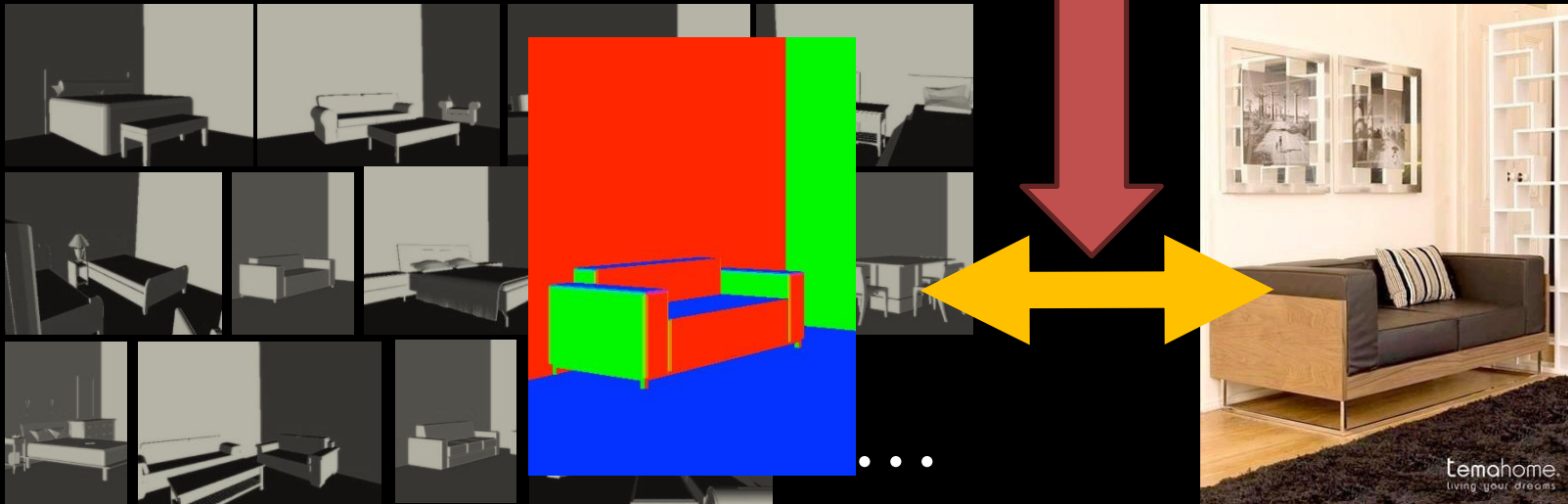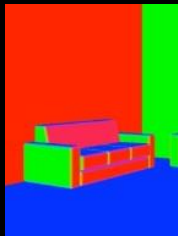
# 2D-3D Scene Matching

3D Model Database   Does it match?   Input



. . .

# Naïve 2D-3D Scene Matching

1K Models

# Naïve 2D-3D Scene Matching

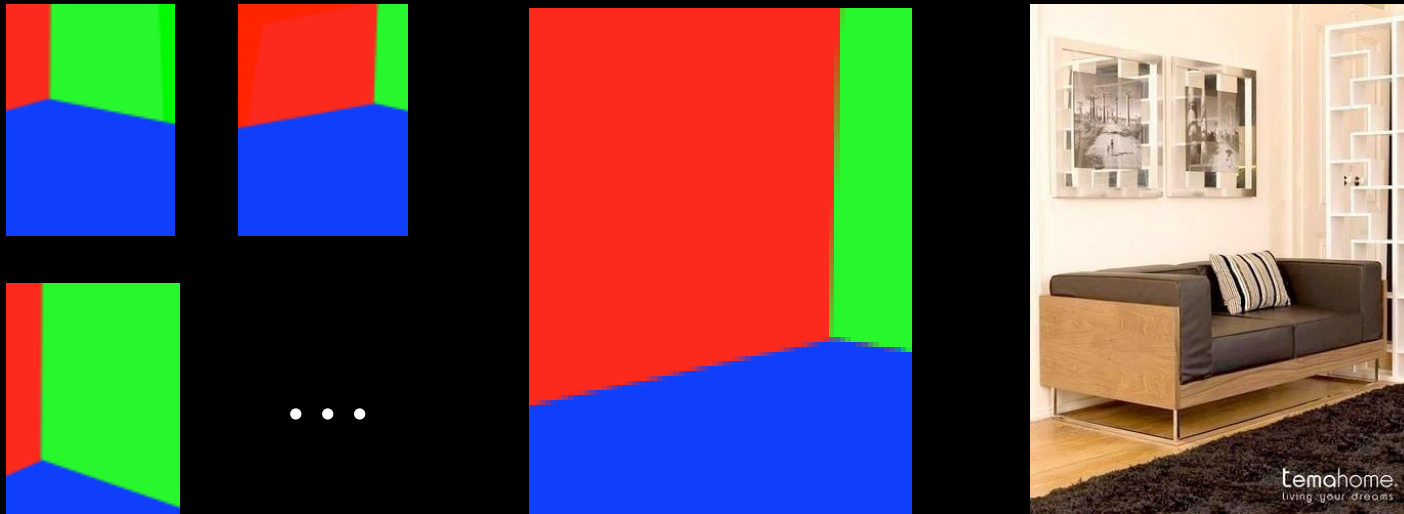1K Models  x 1K Layouts

# Naïve 2D-3D Scene Matching

1K Models x 1K Layouts x 100 rotations

# 2D-3D Scene Matching

Instead: apply what we already know!



Single VP Triplet

Most likely layouts

Model library

Satkin et al., 2012,2013,2014  42

# 2D-3D Scene Matching



$$\mathbf{f}_1 \quad > \quad \mathbf{f}_2$$

Learn w to rank models using ranking svm

Satkin et al., 2012,2013,2014  43

# Pose and Object Sampling

Render+test enables search over hypotheses generated on the fly

# Pose and Object Sampling

On average: 5% gain in accuracy

Initial Estimate

Final Estimate

# Results

Input                         Normals                         Semantics

# Benefits of 3D

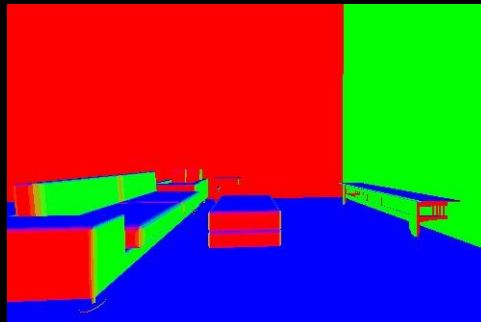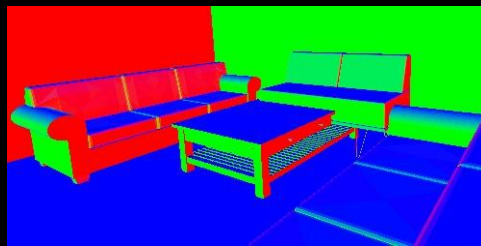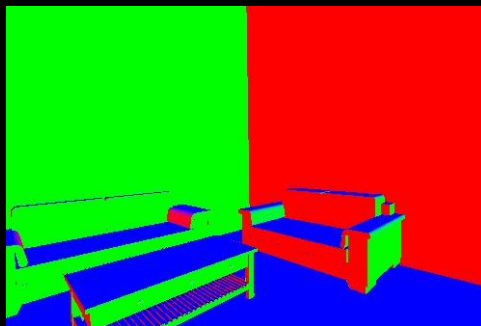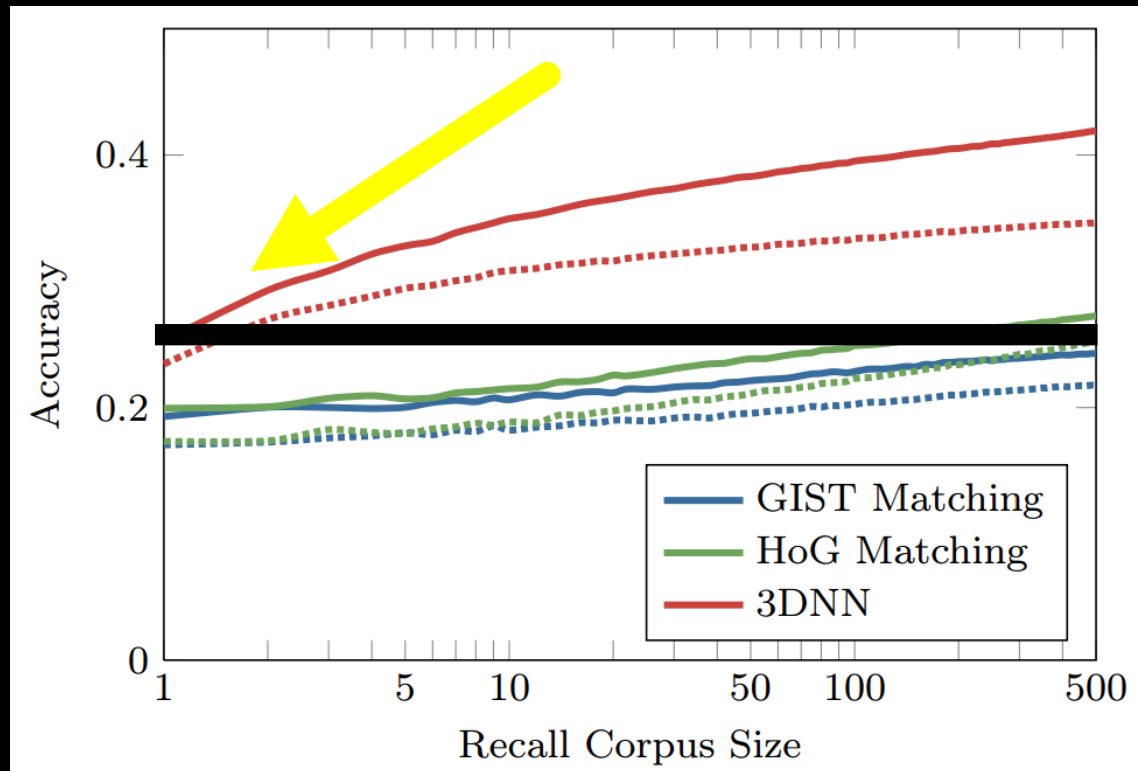## Don't need every viewpoint explicitly!

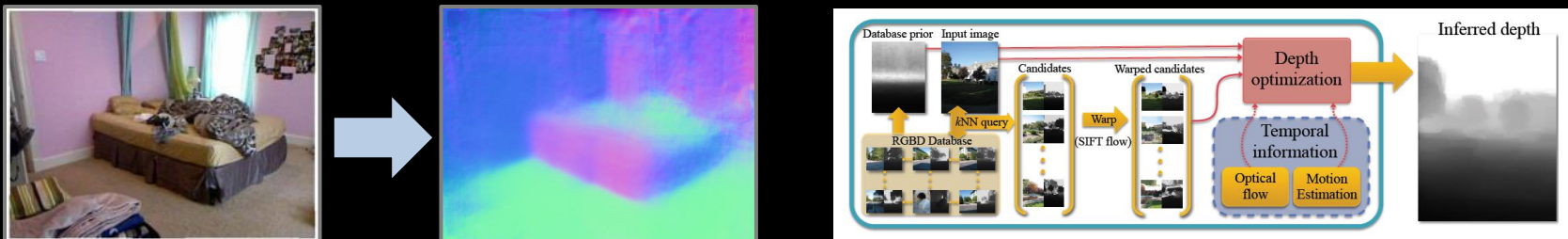# Overview

## 1. How to use 3D models



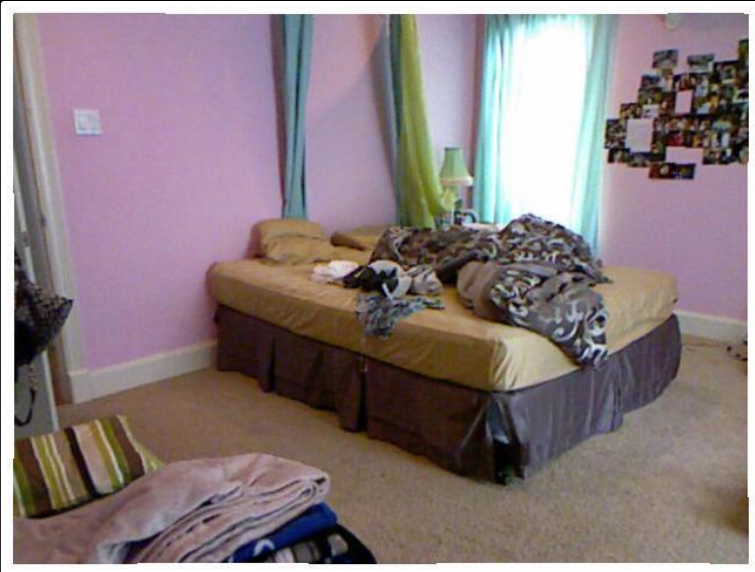## 2. How to use the Kinect

# Kinect Data

RGB

Depth

# Kinect Data

RGB

Depth

Normals

# 2.5D Data

- Advantages:
  - Corresponding natural color image

- Disadvantages:
  - 2.5D (can't render)
  - Missing data, noise
  - Representations can be difficult to transfer

# General Approach



How to transfer representation?

How do we get this correspondence?

# Two Approaches

## Data-Driven Alignment

# Two Approaches

## Clustering + Detection

# Data-Driven Alignment



55

# Finding Correspondences

Input

# 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



61

# Finding Correspondences

Warped Depths



...

?

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

# Optimizing Depthmaps

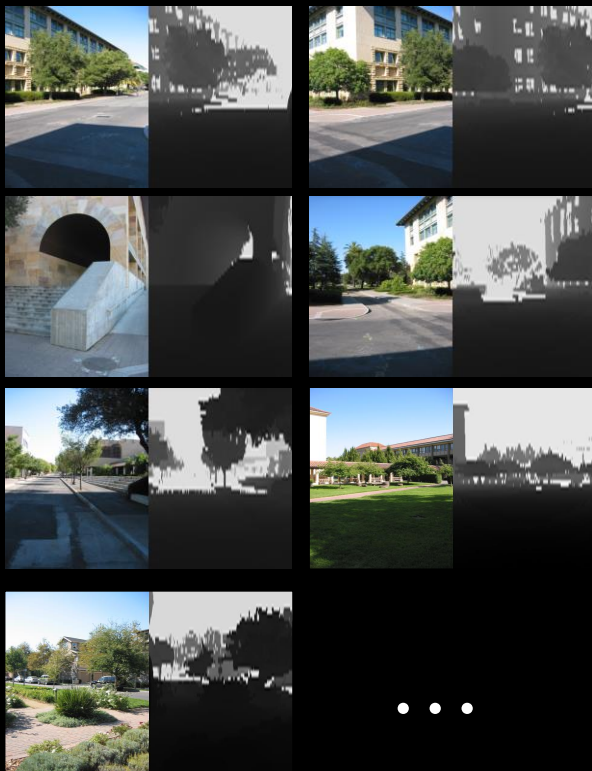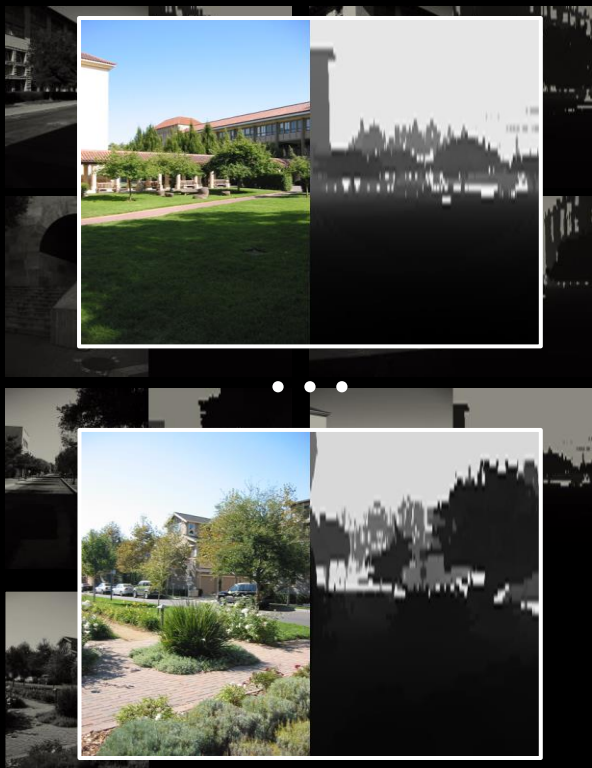$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i \left( |D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1 \right) \right]$$

$$+ \ \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

$D_i$    -Depth being optimized

$C_i$    -Warped depth candidate

# Optimizing Depthmaps

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

$$+ \, \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

...

# Optimizing Depthmaps

Enforce depth to match candidates

$$\sum_{i \in \text{pixels}} \left[ \sum_{C \in \text{candidates}} w_i \left( |D_i - C_i|_1 + \gamma |\nabla D_i - \nabla C_i|_1 \right) \right.$$

$$\left. + \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1 \right]$$

Absolute depth          Relative depth
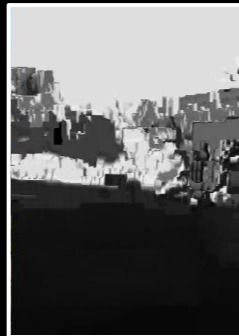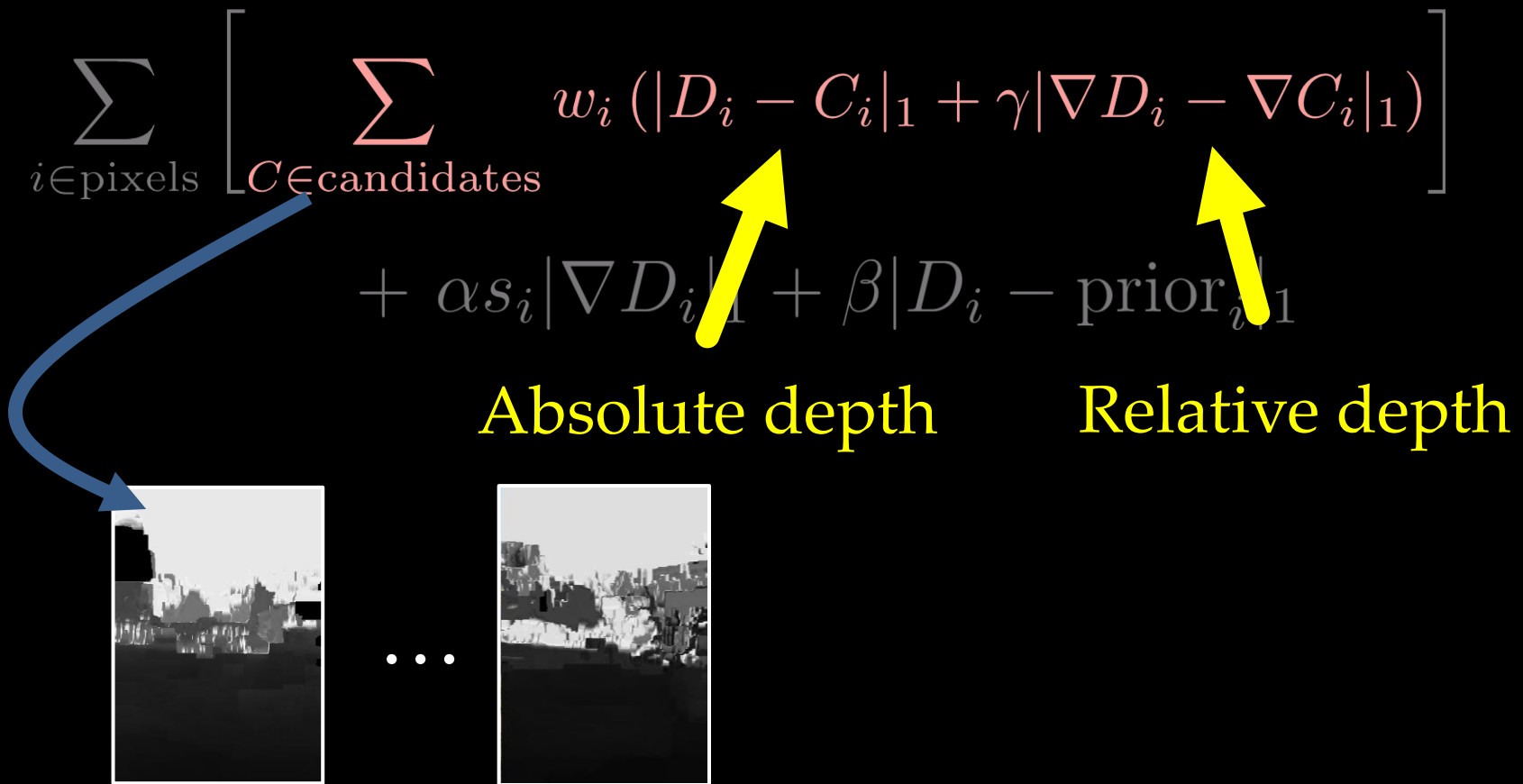
Karsch et al., 2012

# Optimizing Depthmaps

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

$$+ \ \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Spatial smoothness

Karsch et al., 2012

# Optimizing Depthmaps

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

$$+ \; \alpha s_i |\nabla D_i|_1 + \beta |D_i - \text{prior}_i|_1$$

Match the prior

# Results

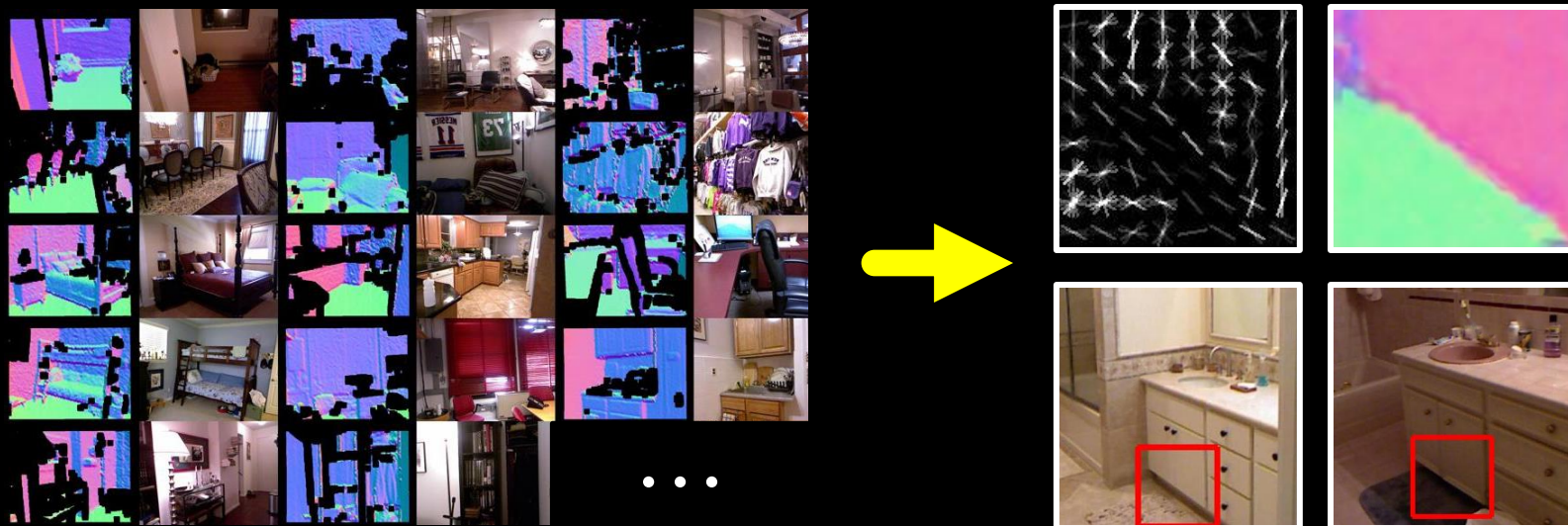Input          True depth          Inferred depth


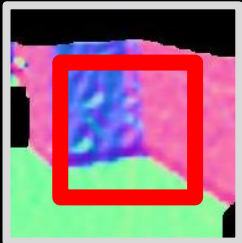
68

# Results

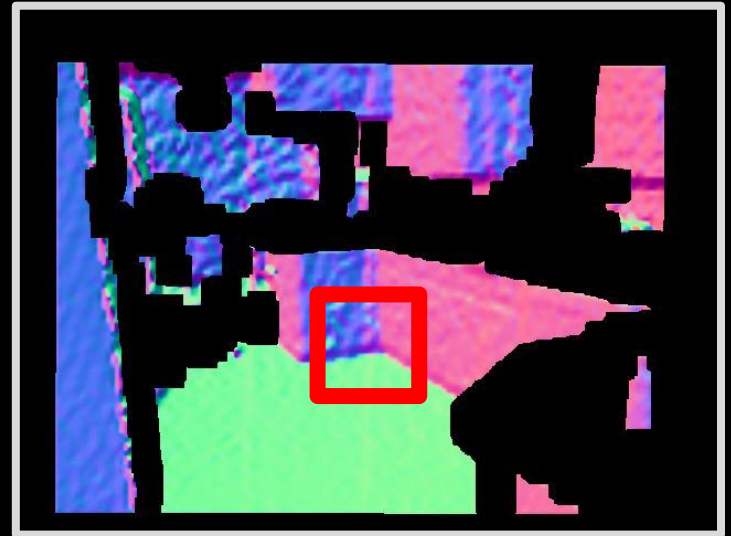Input    True depth    Inferred depth

# Discriminative Clustering + Detection

# Goal

## Visually Discriminative

## Geometrically Informative



Image

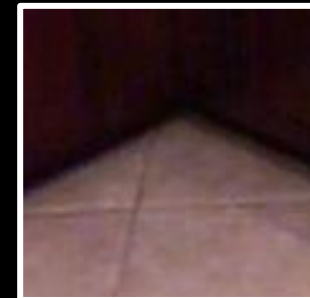Surface Normals

# Goal

## Learn from large-scale RGBD Data
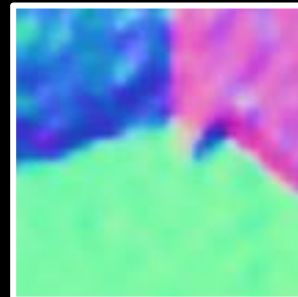
# Approach

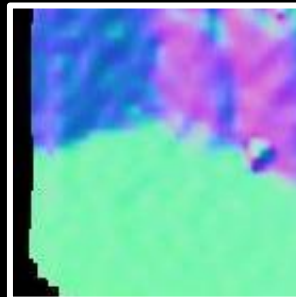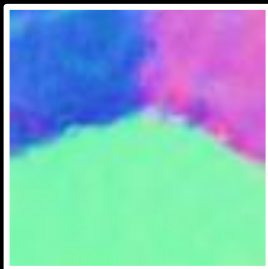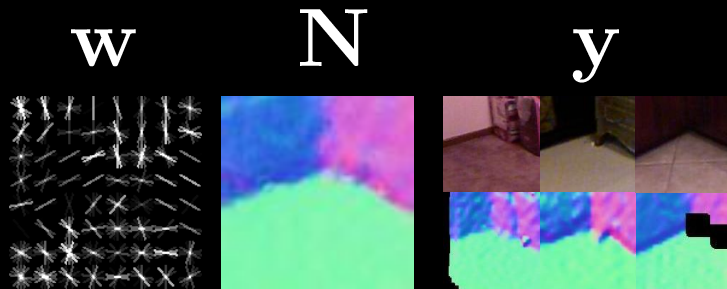## Train time: discriminative clustering w/3D
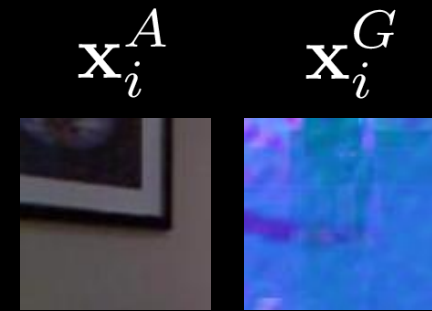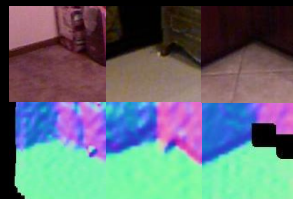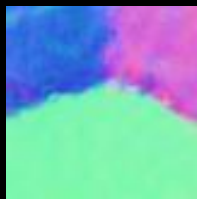
**Detector**

**Instances**

**Normals**

# Objective

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

Primitive

Patch

$\mathbf{w}$ $\quad$ $\mathbf{N}$ $\quad$ $\mathbf{y}$

$\mathbf{x}_i^A$ $\quad$ $\mathbf{x}_i^G$

# 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

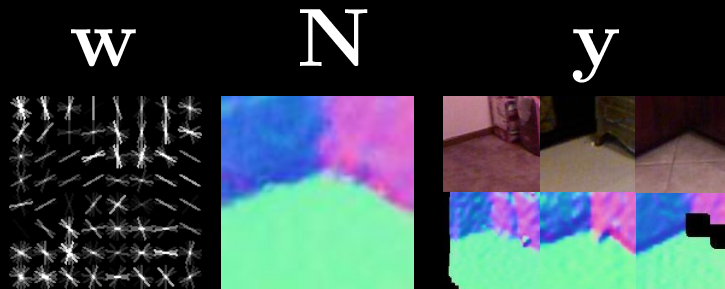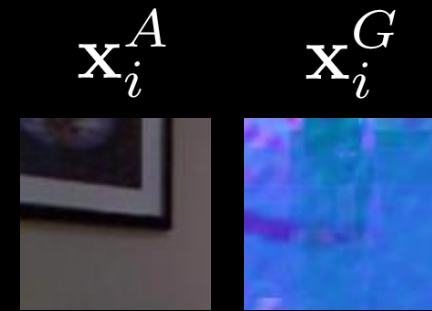$\mathbf{w}$ $\mathbf{N}$ $\mathbf{y}$

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

$\mathbf{w}$    $\mathbf{N}$    $\mathbf{y}$
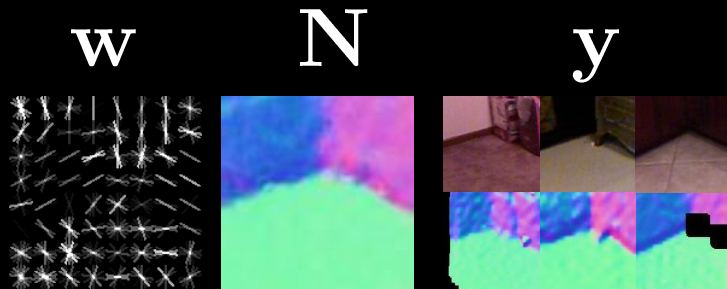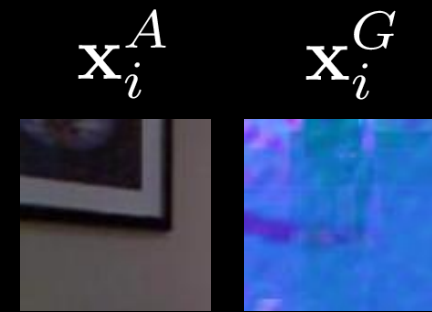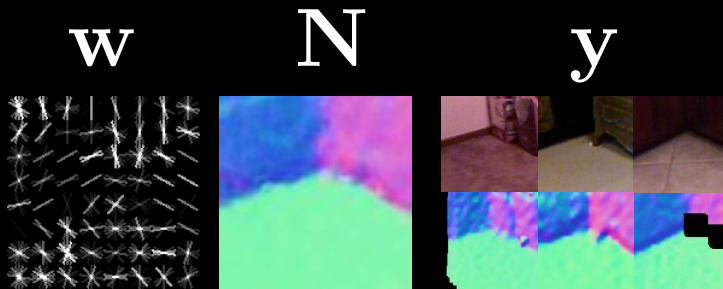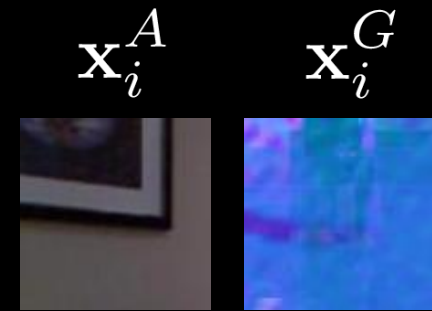
Patch

$\mathbf{x}_i^A$    $\mathbf{x}_i^G$

# Objective

Ensure geometric consistency

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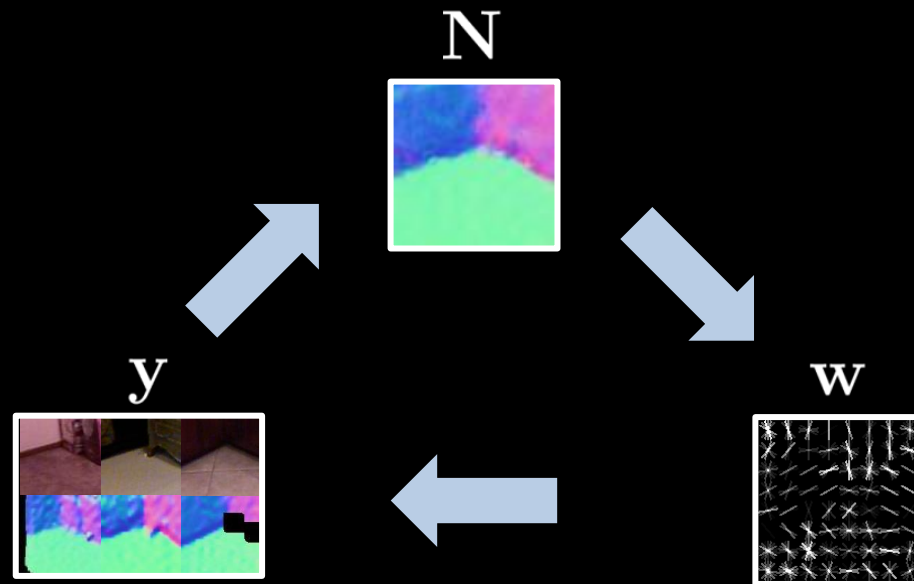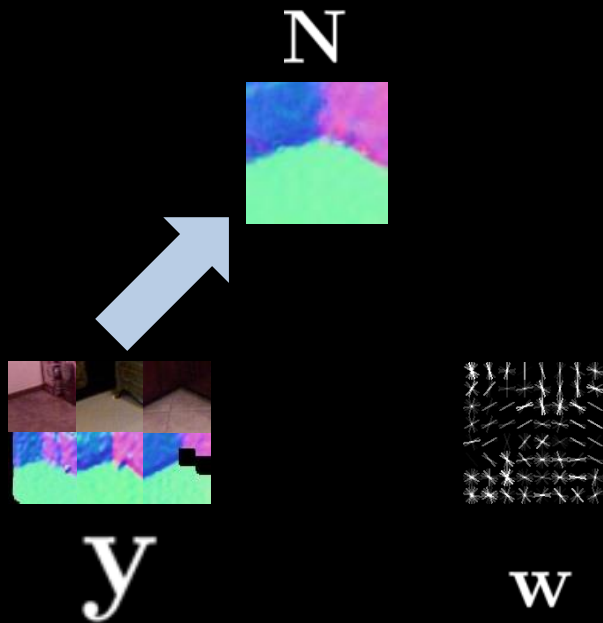
Primitive

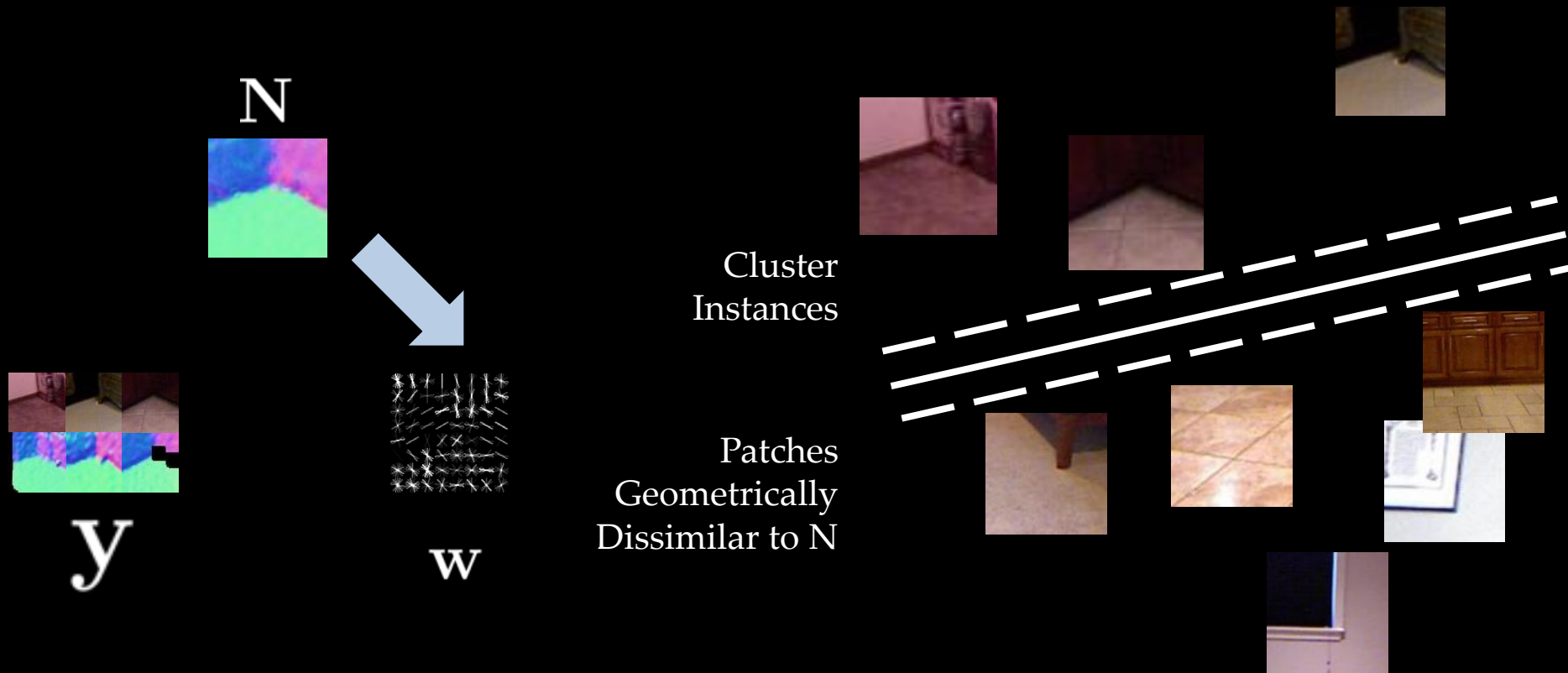$\mathbf{w}$ $\quad$ $\mathbf{N}$ $\quad$ $\mathbf{y}$

Patch

$\mathbf{x}_i^A$ $\quad$ $\mathbf{x}_i^G$

# Objective

Solved with iterative method similar to block-coordinate-descent.
Include min-membership constraint

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

Primitive

$\mathbf{w}$ $\qquad$ $\mathbf{N}$ $\qquad$ $\mathbf{y}$

Patch

$\mathbf{x}_i^A$ $\qquad$ $\mathbf{x}_i^G$

# Iterative Procedure

N

w

y

# Iterative Procedure

N

y          w

$= \mathrm{Avg}\big( \quad \quad \quad \big)$

# Iterative Procedure

N

y → w

Cluster
Instances

Patches
Geometrically
Dissimilar to N

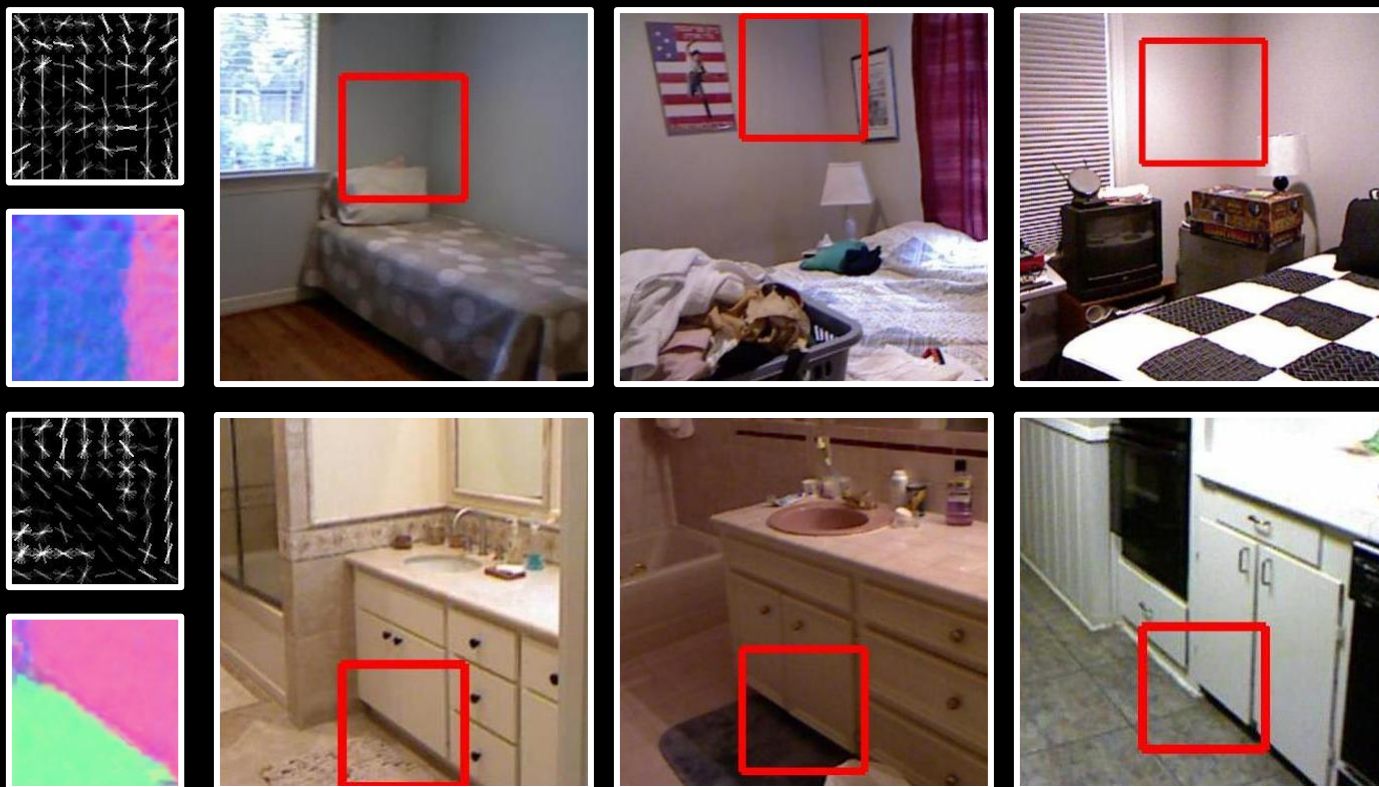# Iterative Procedure



N

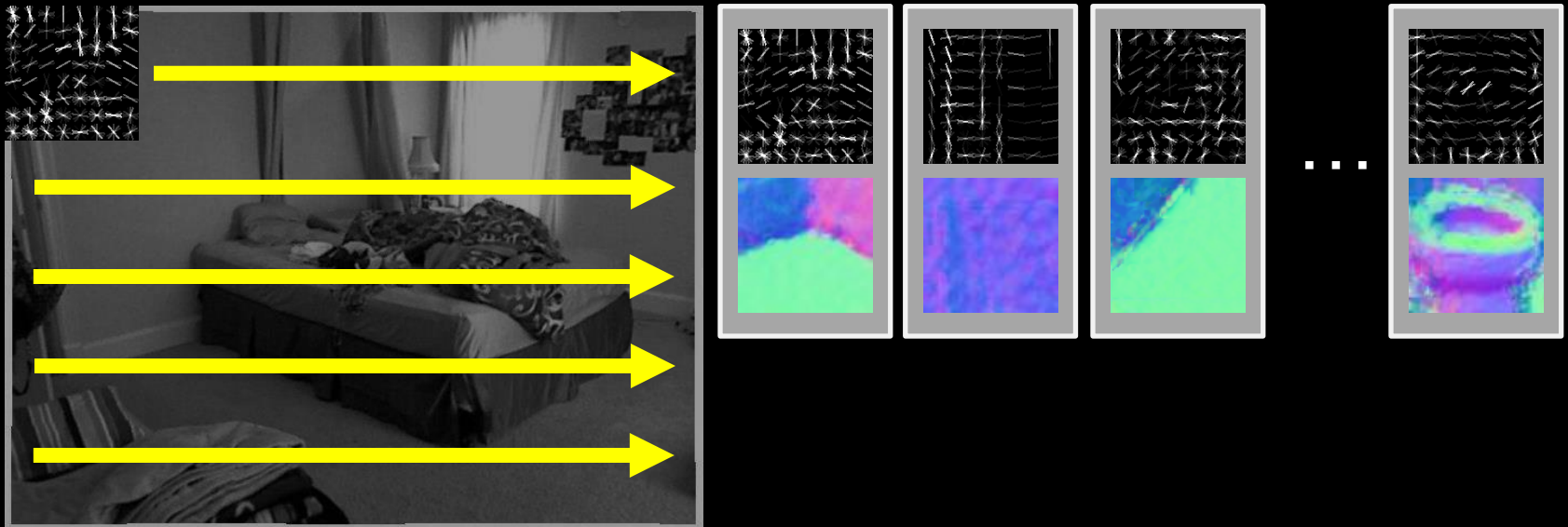y          w

# Primitives
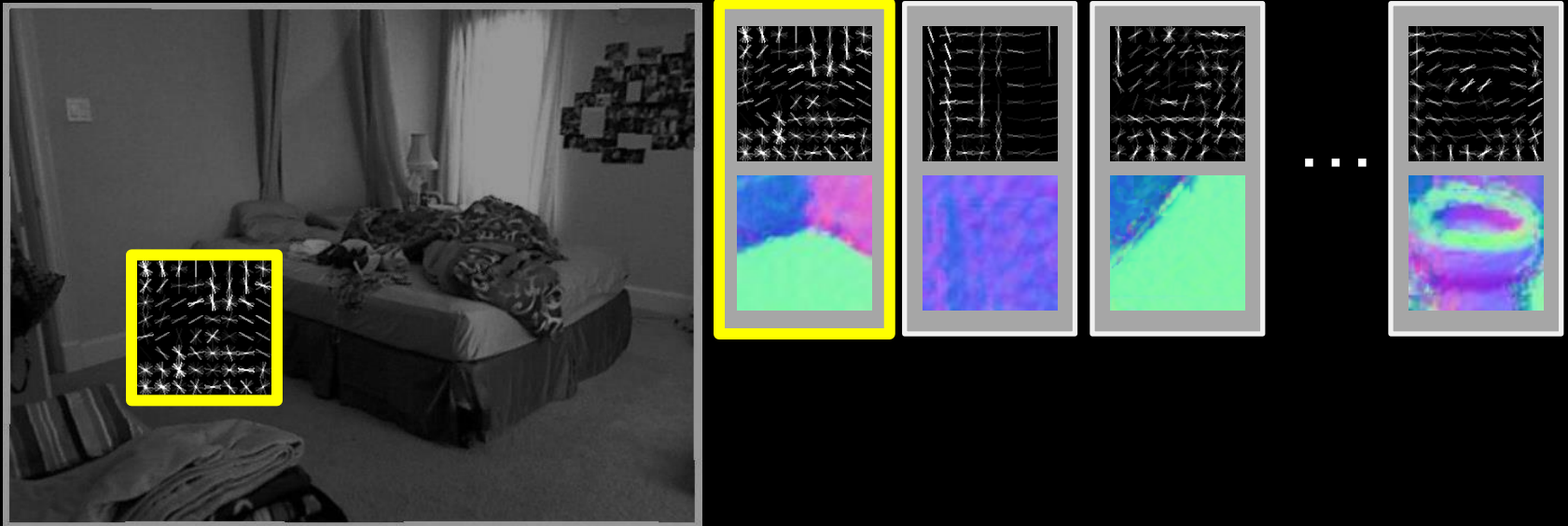
# Primitives

# Primitives
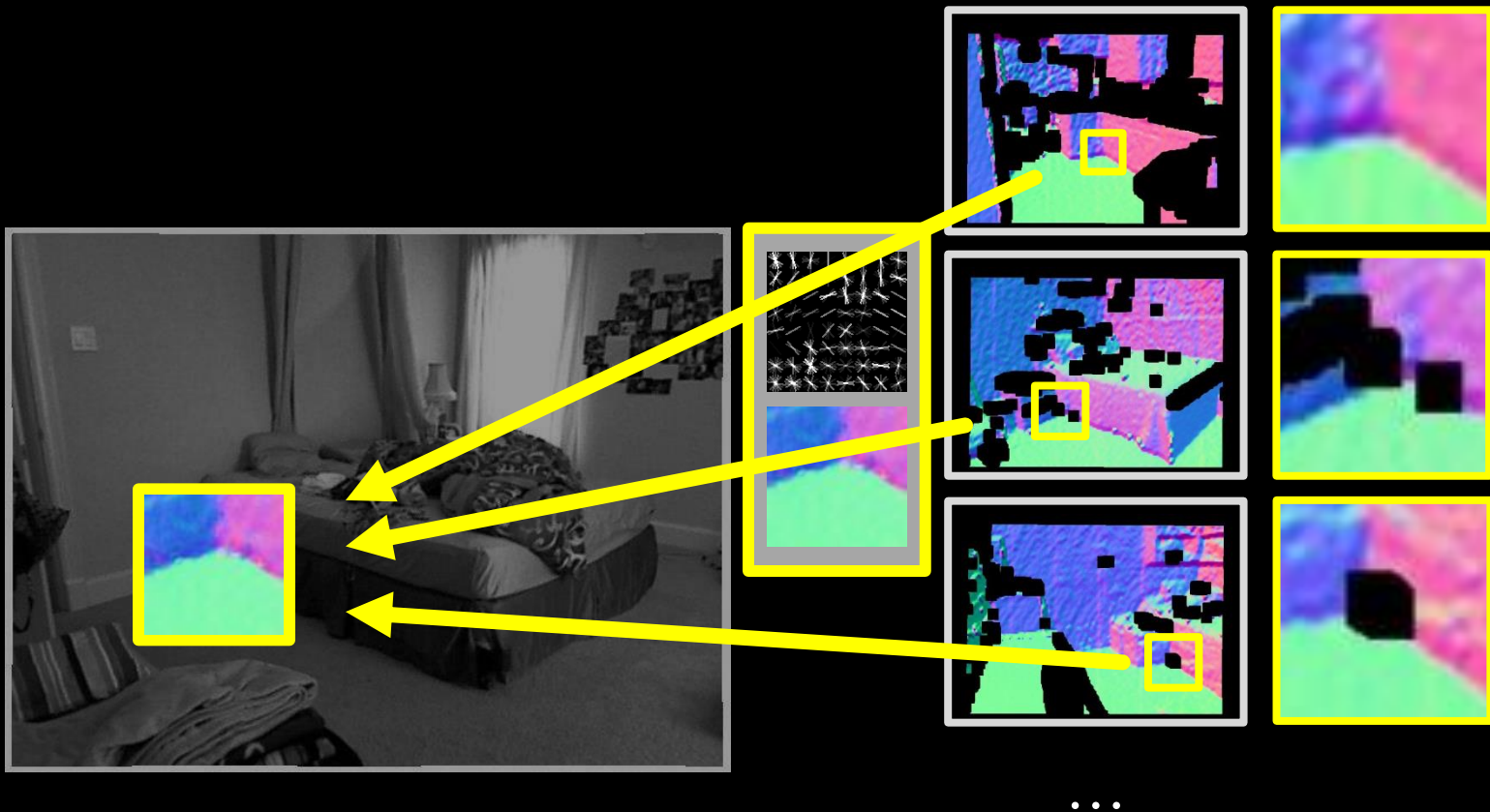


85

# Test-time Correspondence
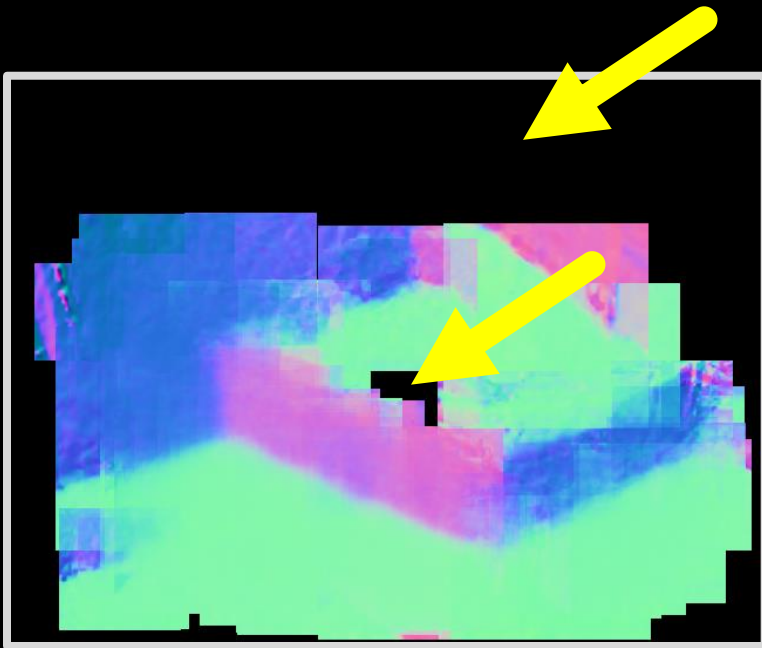
## Correspondence via detection

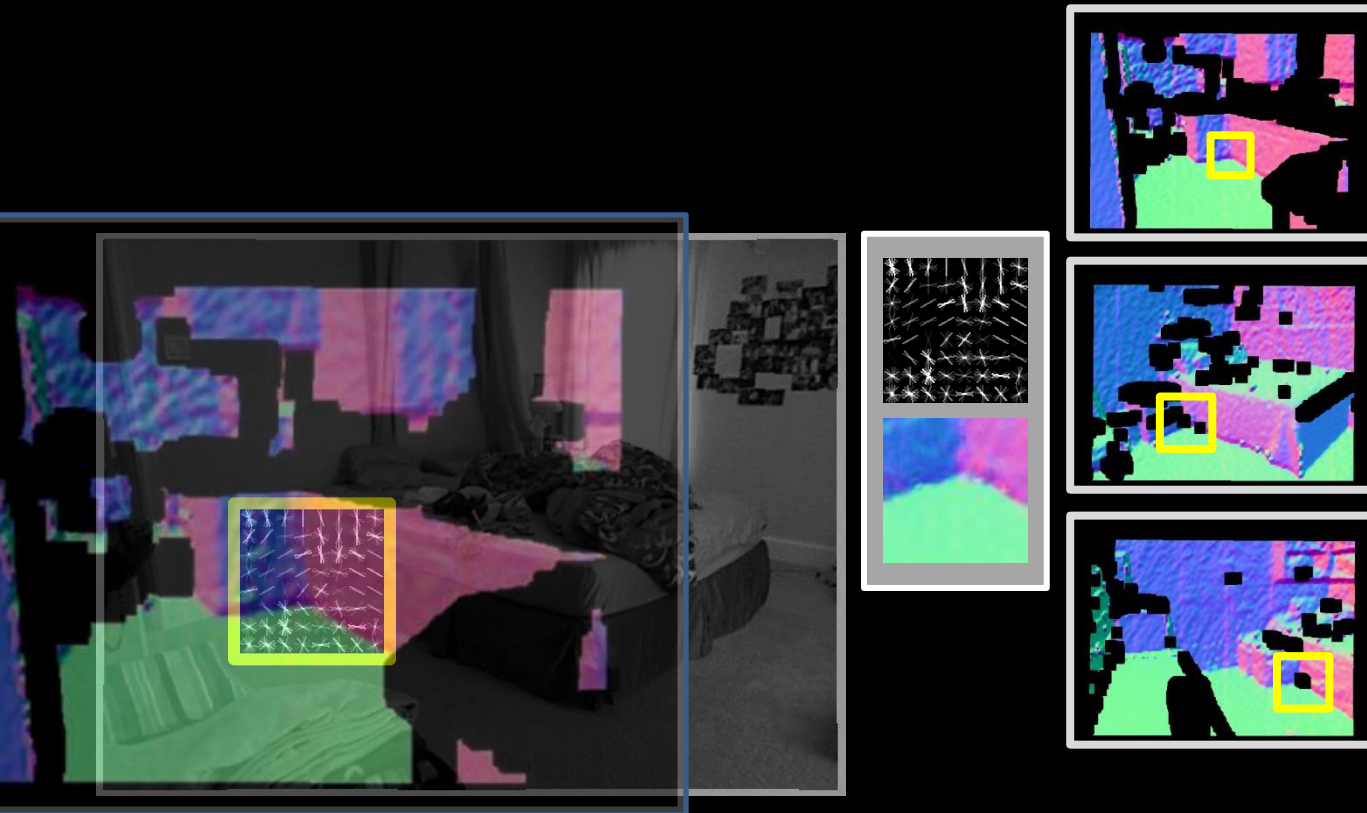# Representation Transfer

# Representation Transfer
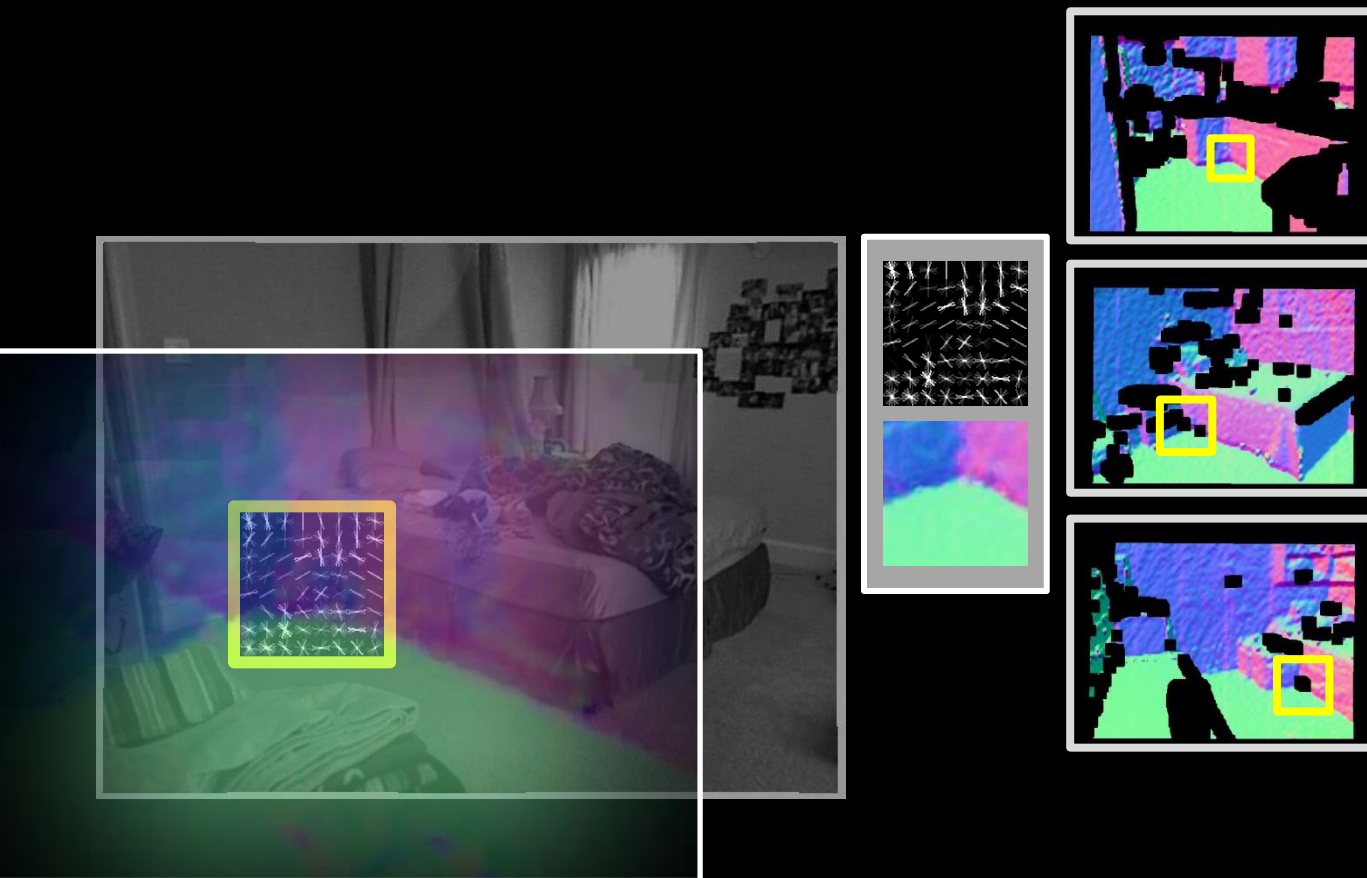


...

88

# Representation Transfer

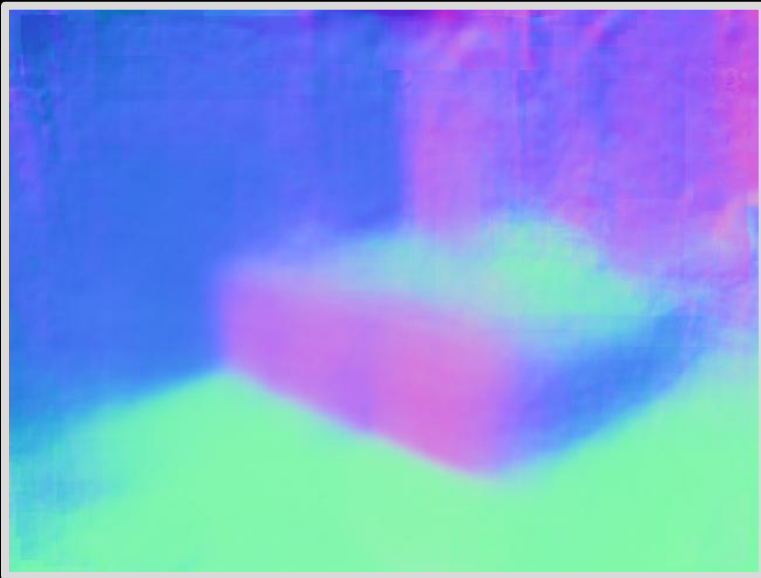Overlaps resolved with averaging

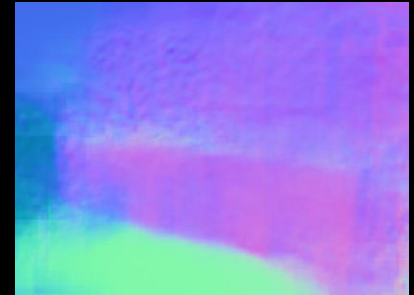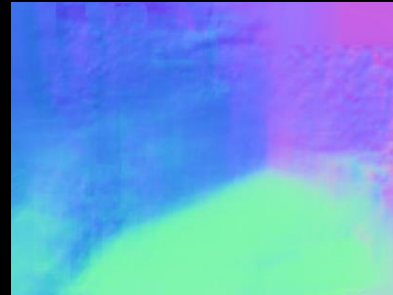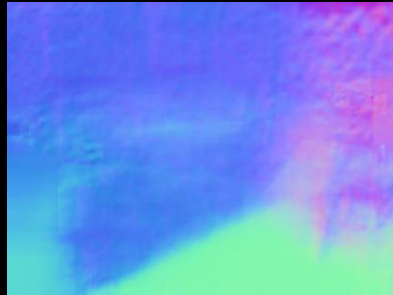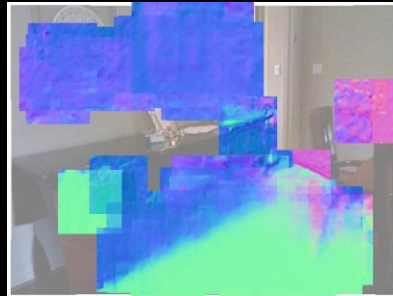# Representation Transfer
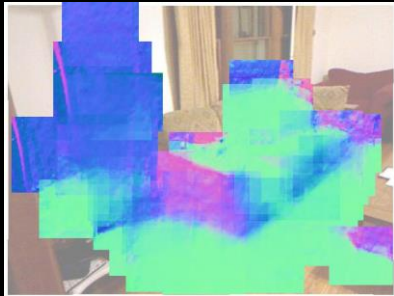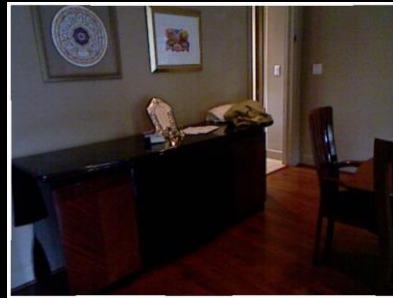
# Representation Transfer
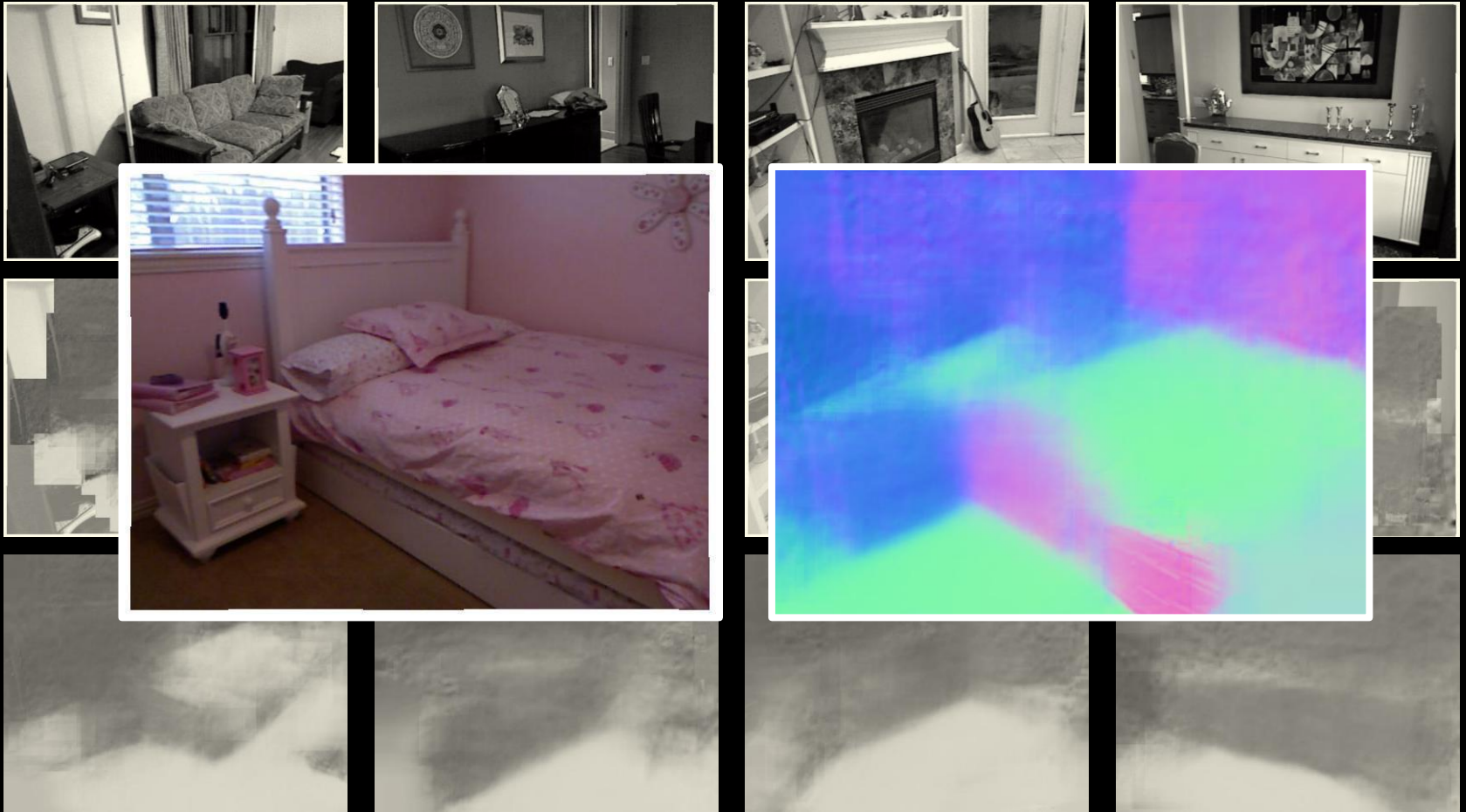
# Representation Transfer

Overlaps resolved with averaging
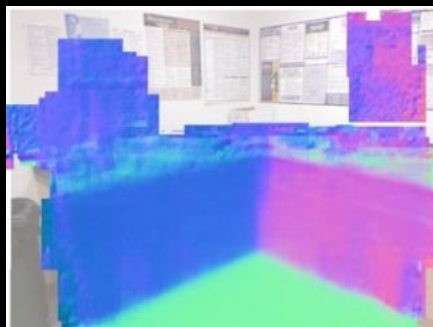
# Results

# Results

# Confidences



Most Confident Result
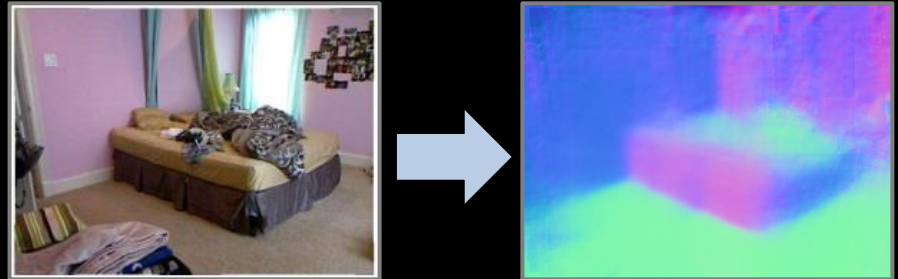
Least Confident Result

# Conclusions

Introduced Data-Driven 3D Scene Understanding

Full 3D Models

RGBD Data



Two Main Problems:

1. Correspondence
2. Representation Transfer

# Future Directions

- How do you get the best of 2.5D and 3D? (see Guo and Hoiem 2013)

- How do you incorporate constraints in data-driven techniques?

# Resources

## (See tutorial website for links + more data/code + slides)

# Survey Books

- D. Hoiem, S. Savarese. *Representations and Techniques for 3D Object Recognition and Scene Interpretation*. Morgan & Claypool, 2011.
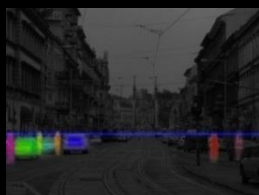
  (link on website)

# Available Kinect 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
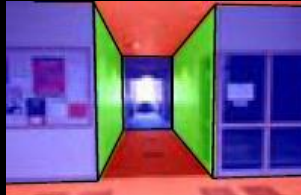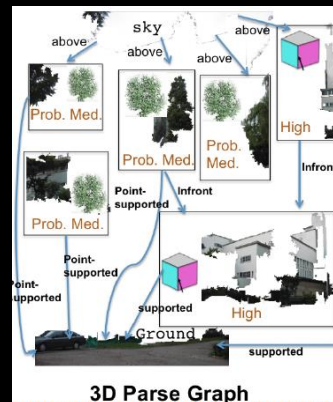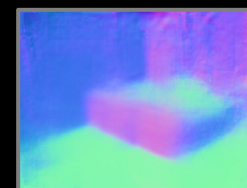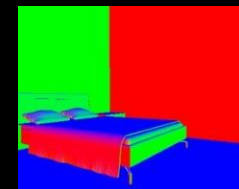
# Available Code



**Region labels**

**+ Boundaries and objects**

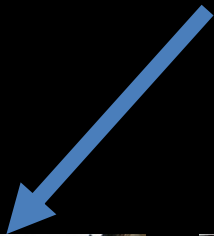**Stronger geometric constraints from domain knowledge**

**Volumetric + functional constraints**

**Data-driven 3D**

# Available Code
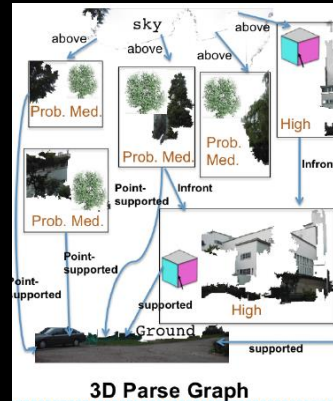
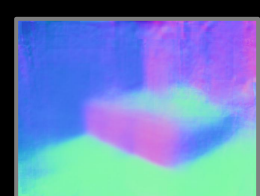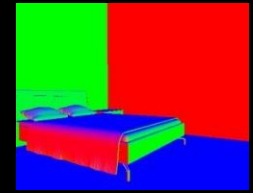Hoiem et al., Geometric Context,
Saxena et al., Make 3D



**Region labels**

**+ Boundaries and objects**

**Stronger geometric constraints from domain knowledge**

**Volumetric + functional constraints**

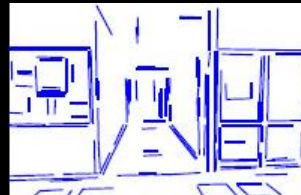**Data-driven 3D**

# Available Code

Hoiem et al., Occlusion boundaries
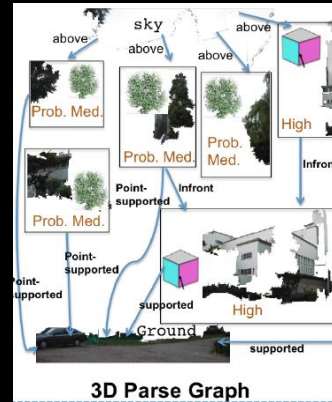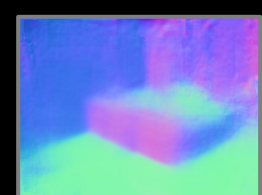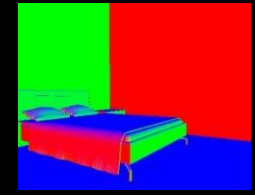Hoiem et al., Putting objects in perspective



**Region labels**

**+ Boundaries and objects**

**Stronger geometric constraints from domain knowledge**

**Volumetric + functional constraints**

**Data-driven 3D**
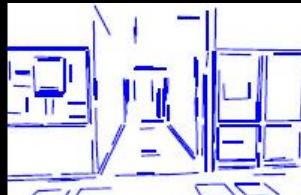
# Available Code

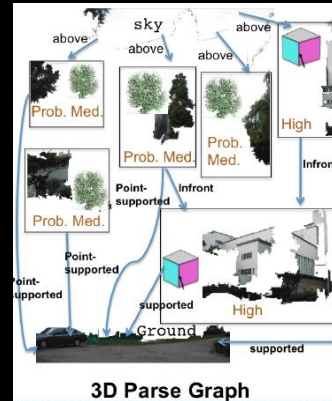Lee et al., Orientation Maps
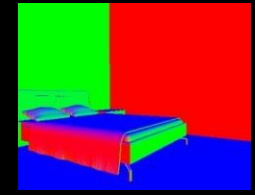Hedau et al., Room-fitting



**Region labels**

**+ Boundaries and objects**

**Stronger geometric constraints from domain knowledge**

**Volumetric + functional constraints**

**Data-driven 3D**

# Available Code

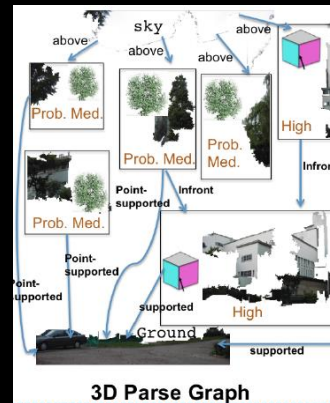Gupta et al., Blocks World
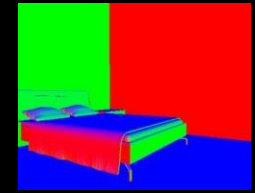Choi et al., Geometric Phrases



**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
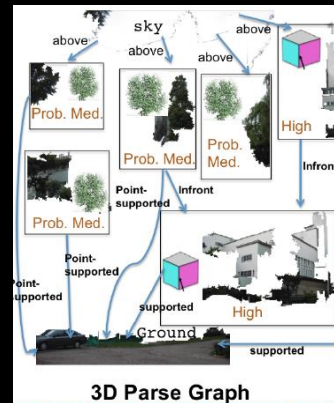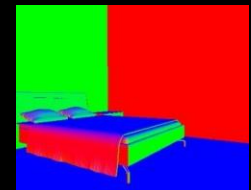Aubrey et al., Seeing 3D Chairs



**Region labels**     **+ Boundaries and objects**     **Stronger geometric constraints from domain knowledge**     **Volumetric + functional constraints**     **Data-driven 3D**
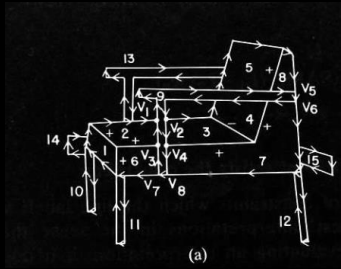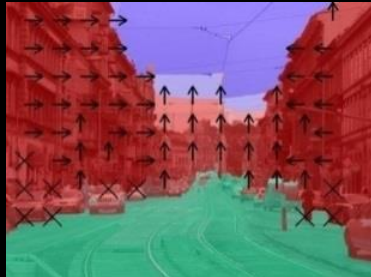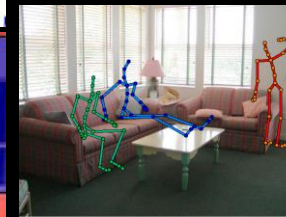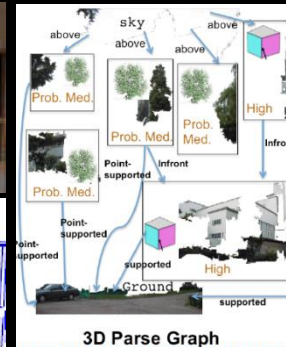
# Thank You



**Martial**

**Derek**

**Abhinav**

**David**

**Introduction,
Applications,
History**

**Region labels
+Boundaries
+Objects**

**Stronger
geometric
constraints**

**Volumetric +
Functional
Constraints**

**Data-Driven 3D**