

# **Seeing the Unseen:** Comprehensive 3D Scene Understanding

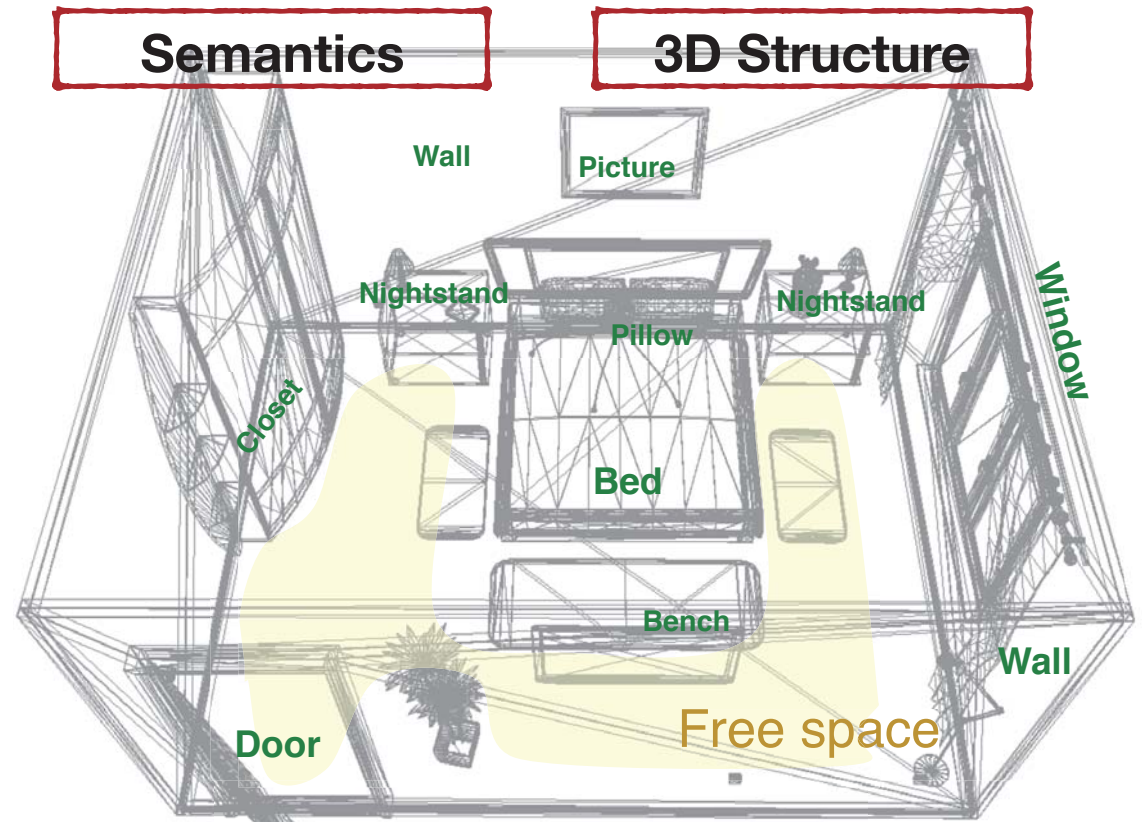
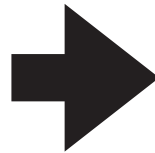
Shuran Song

Princeton —> Google —> Columbia

# Comprehensive 3D Scene Understanding



Partial Observation of the Environment

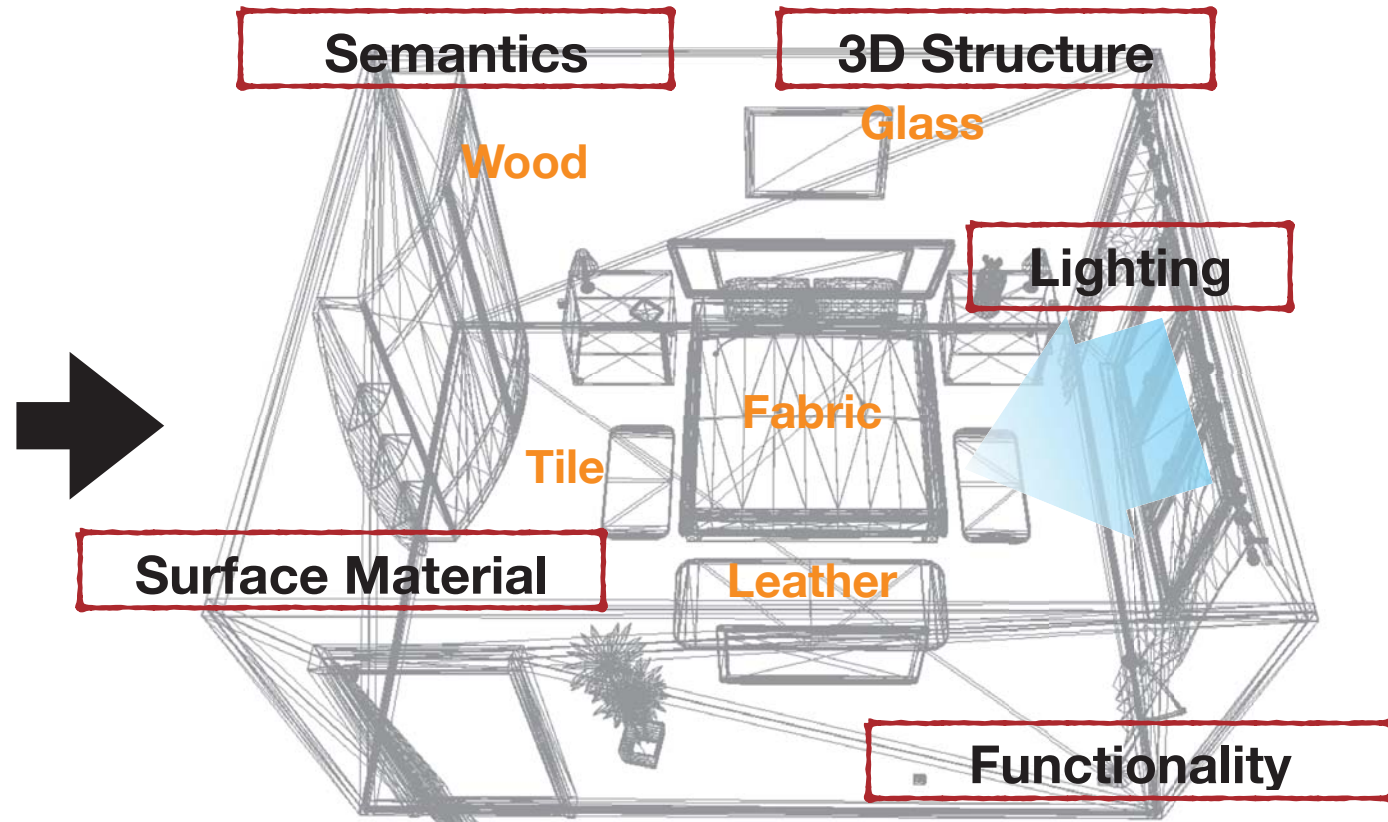


Complete Representation of the 3D Scene

# Comprehensive 3D Scene Understanding



Partial Observation of the Environment



Complete Representation of the 3D Scene

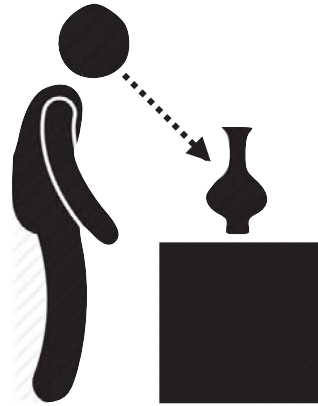
# Challenge: Partial Observation



# Challenge: Partial Observation



**Sensors**



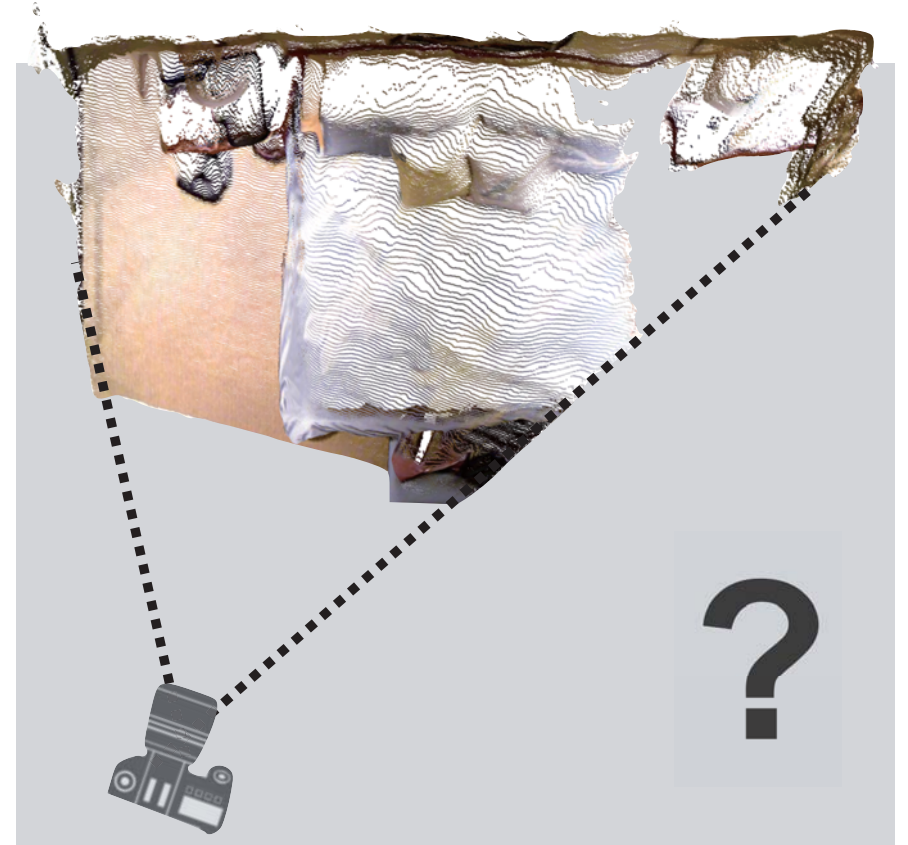
**Partial Observation**



# Challenge: Partial Observation

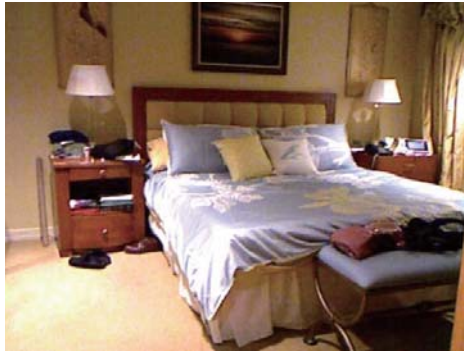


**Occlusion**

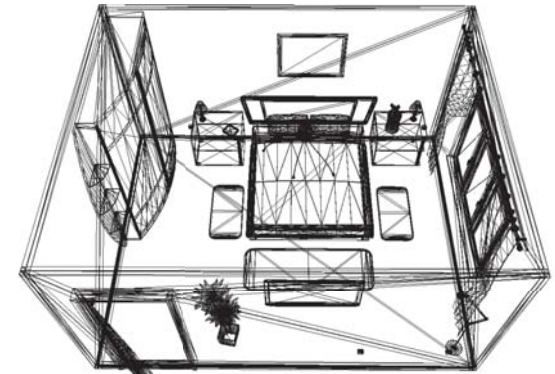


**Limited Camera FOV**

# Comprehensive 3D Scene Understanding



**Partial Observation**

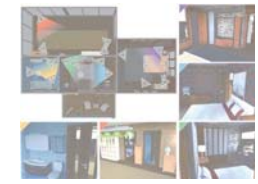
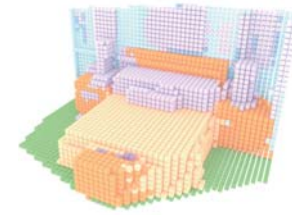
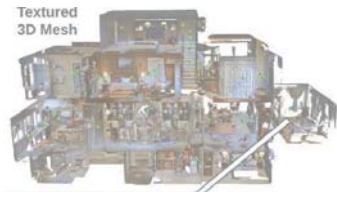


**Complete 3D Scene**



SHAPE  
NET

# Datasets



ShapeNet  
arXiv 2015

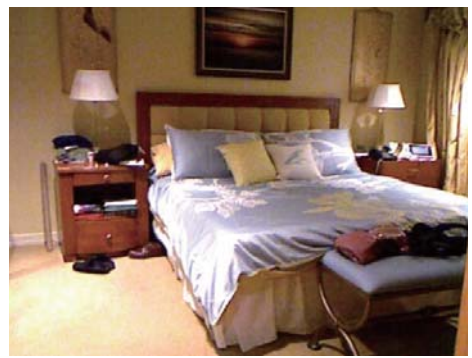
SUN RGB-D  
ECCV 2014

Matterport3D  
3DV 2017

SSCNet  
CVPR 2017

Rendering  
CVPR 2017

3DMatch  
CVPR 2017

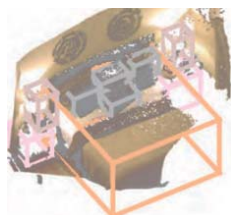
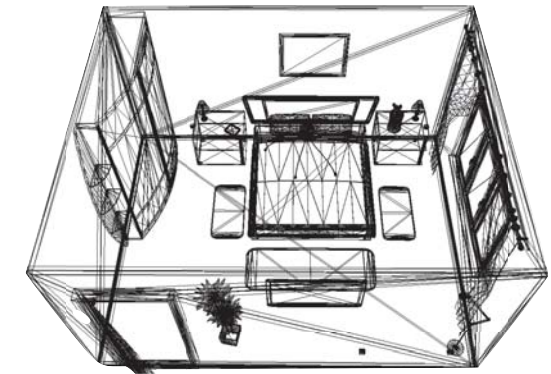


# Algorithms

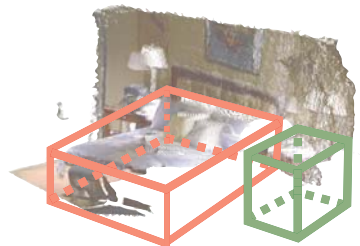
PanoContext  
ECCV 2014

3DShapeNets  
CVPR 2015

RobotInARoom  
arXiv 2015



Deep SlidingShapes  
CVPR 2016



SlidingShapes  
ECCV 2014



LSUN  
arXiv 2016



Im2Pano3D  
CVPR'18



ARC 2017  
ICRA 2017



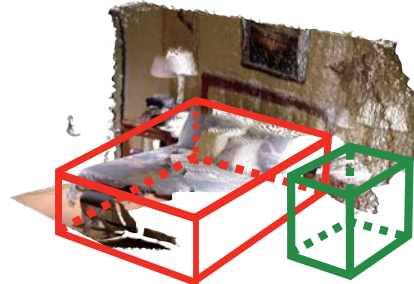
ARC 2018  
ICRA 2018



Tracking  
ICCV 2013

# Applications

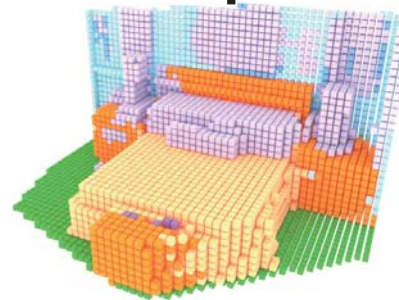
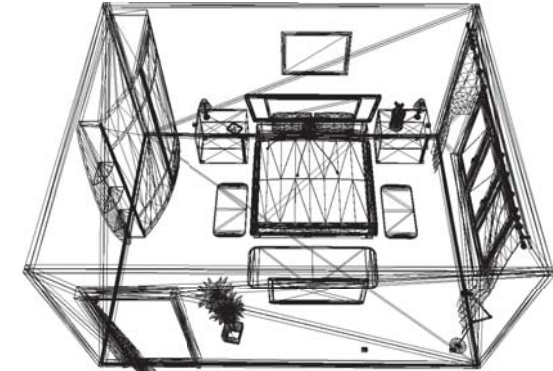
# Comprehensive 3D Scene Understanding



**Amodal 3D**  
[Song and Xiao  
ECCV'14, CVPR'16]



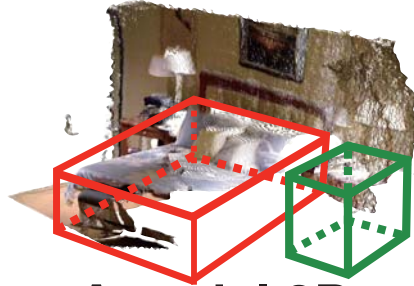
**Beyond FoV**  
[Song et al. CVPR'18]



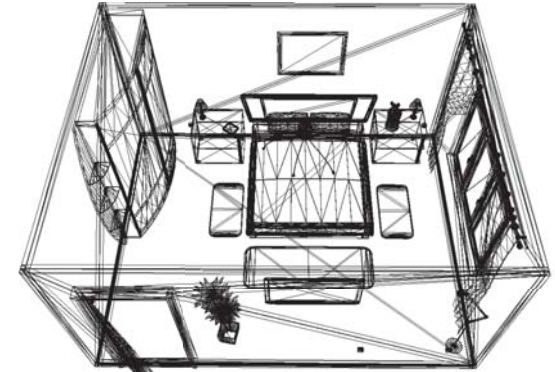
**Higher Fidelity**  
[Song et al. CVPR'17]

- Semantics Category
- 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. Properties ...

# Comprehensive 3D Scene Understanding



**Amodal 3D**  
[Song and Xiao  
ECCV'14, CVPR'16]





# Object Detection

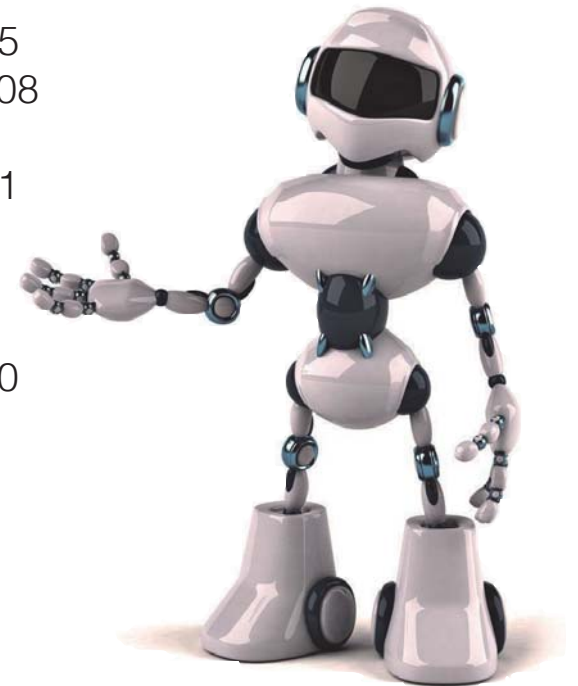
2D Visible (Modal) Surface



Traditional Object Detection Output

Aubry et al. CVPR'14  
Dalal and Triggs CVPR'05  
Felzenszwalb et al. CVPR'08  
Bo et al. CVPR'2011  
Malisiewicz et al. ICCV'11  
Girshick et al. CVPR'14  
Ren et al. NIPS'15  
Girshick, ICCV'15  
Everingham et al. IJCV'10  
He et al. ICCV'17  
Liu et al. ECCV'16  
Erhan et al. CVPR'14  
He et al. ECCV'14  
Szegedy NIPS'13  
...

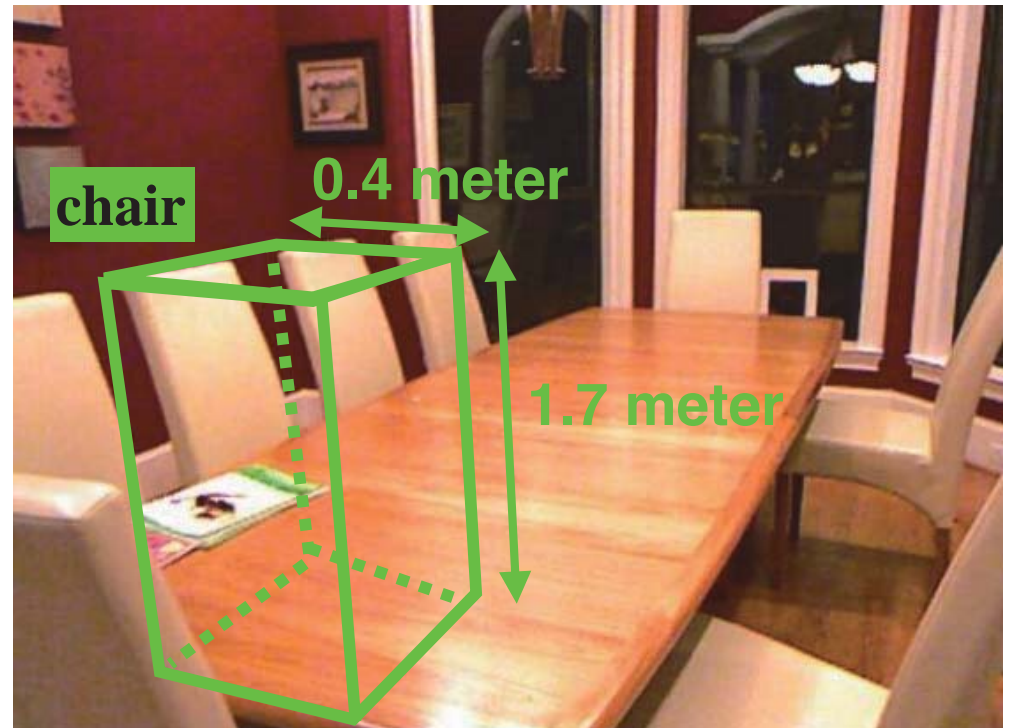
Where to sit?



# Object Detection

2D Visible (Modal) Surface

3D Complete (Amodal) Shape

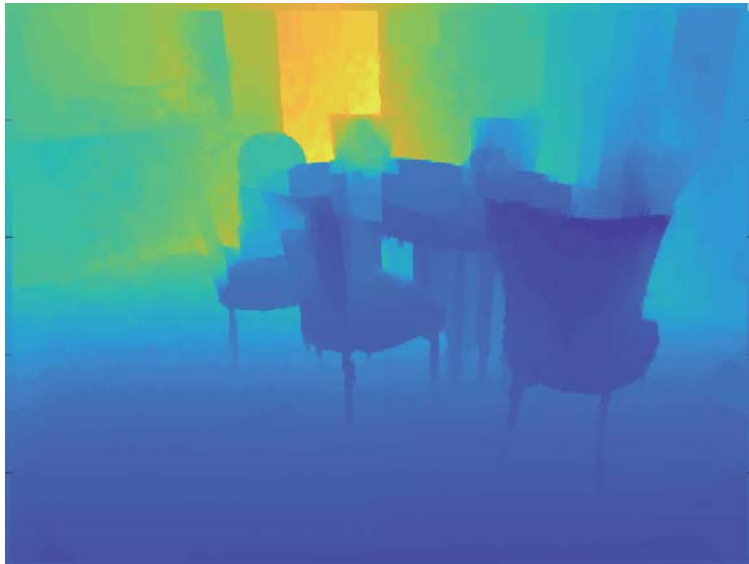


Traditional Object Detection Output

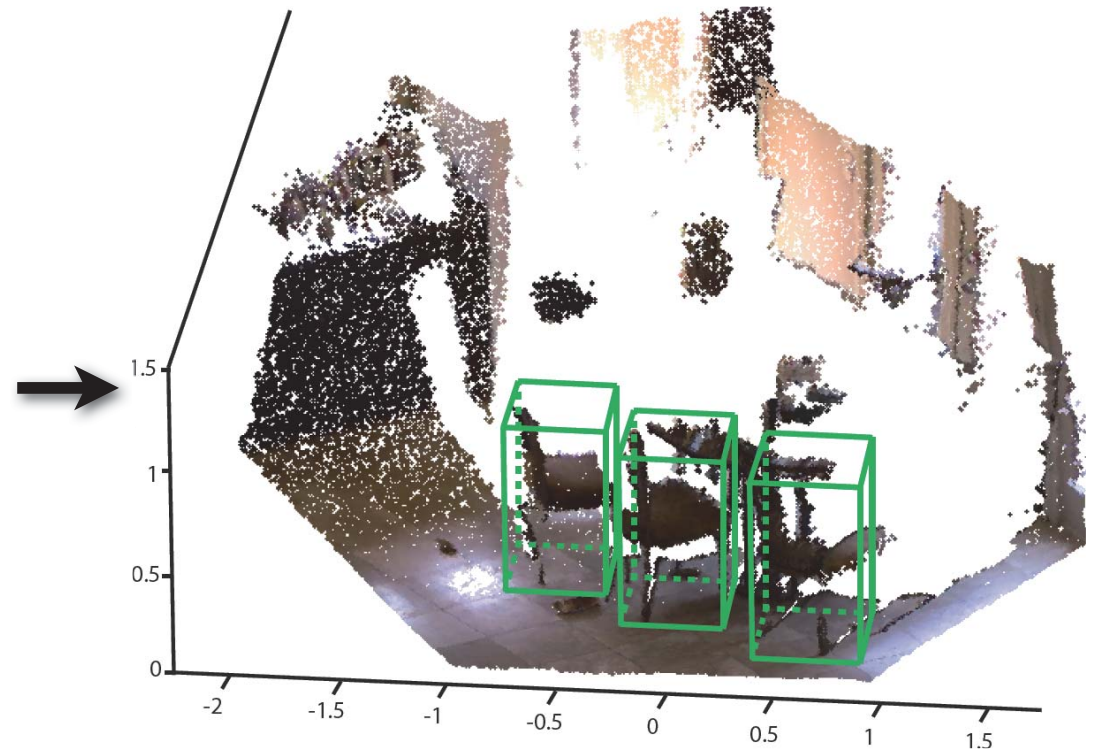
This work

S. Song and J. Xiao, Sliding Shapes for 3D Object Detection in Depth Images, ECCV 2014  
S. Song and J. Xiao, Deep Sliding Shapes for Amodal 3D Object Detection in RGB-D Images, CVPR 2016

# Deep Sliding Shapes

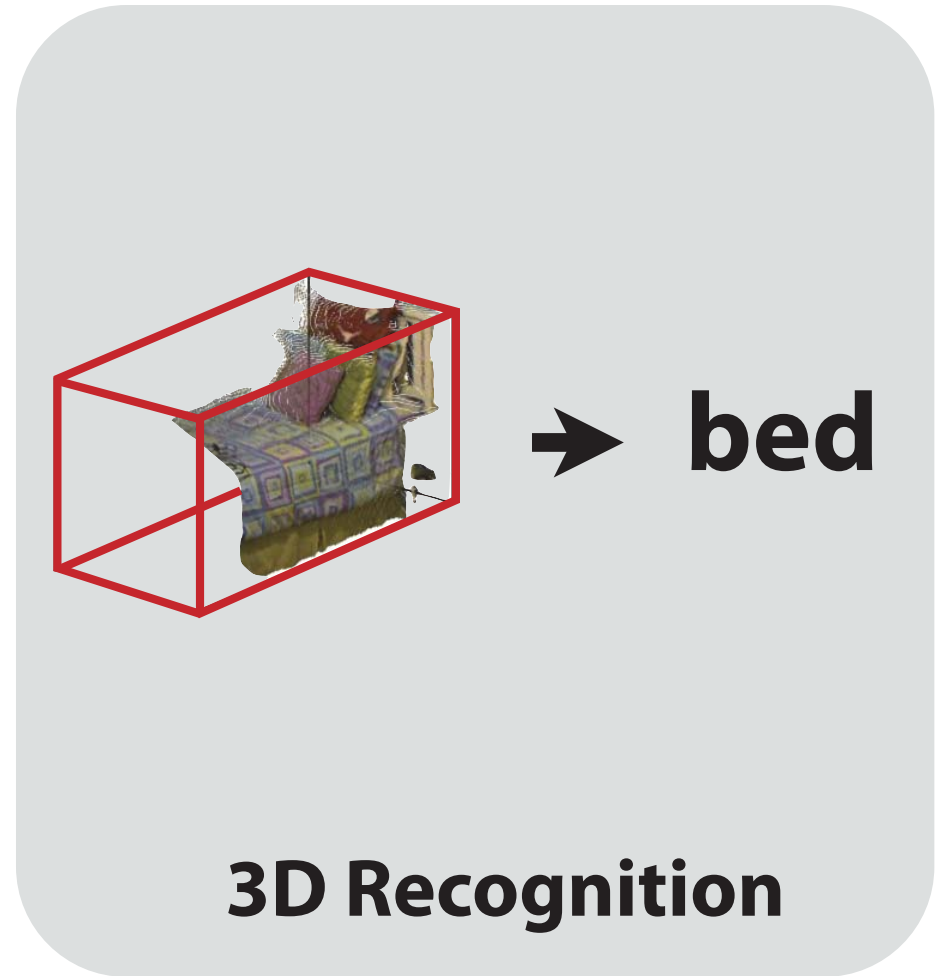
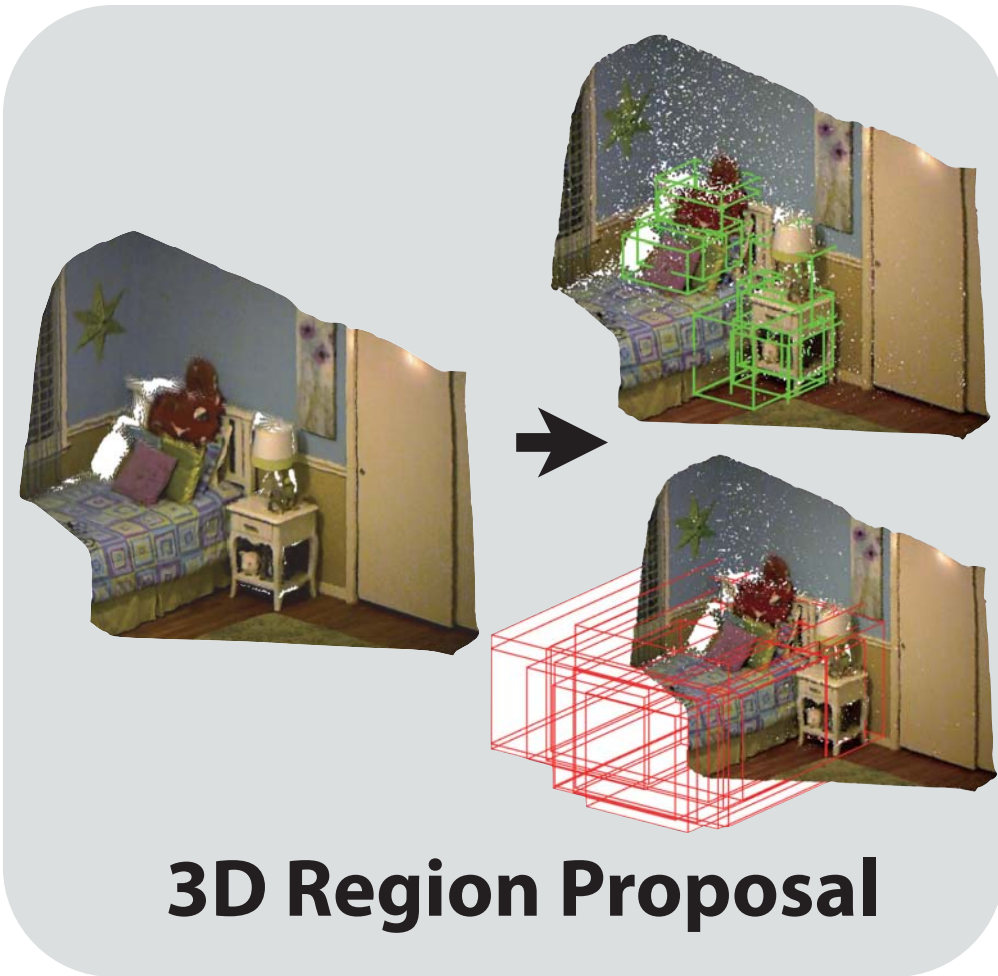


Input: Kinect Depth Image



Output: 3D Bounding Box

# 3D Deep learning

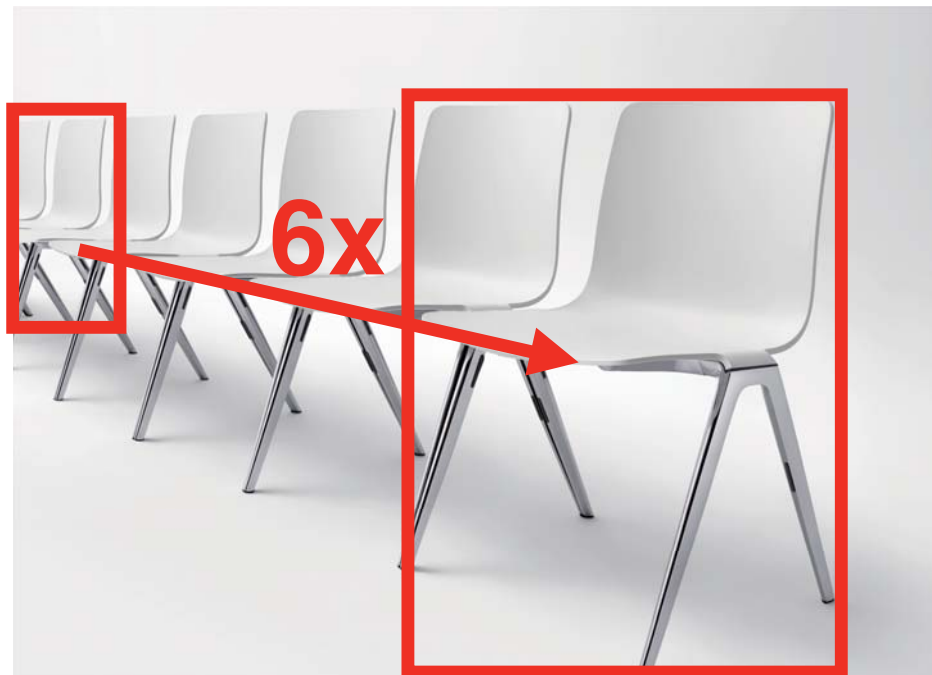




# **Representation: 3D vs. 2D**

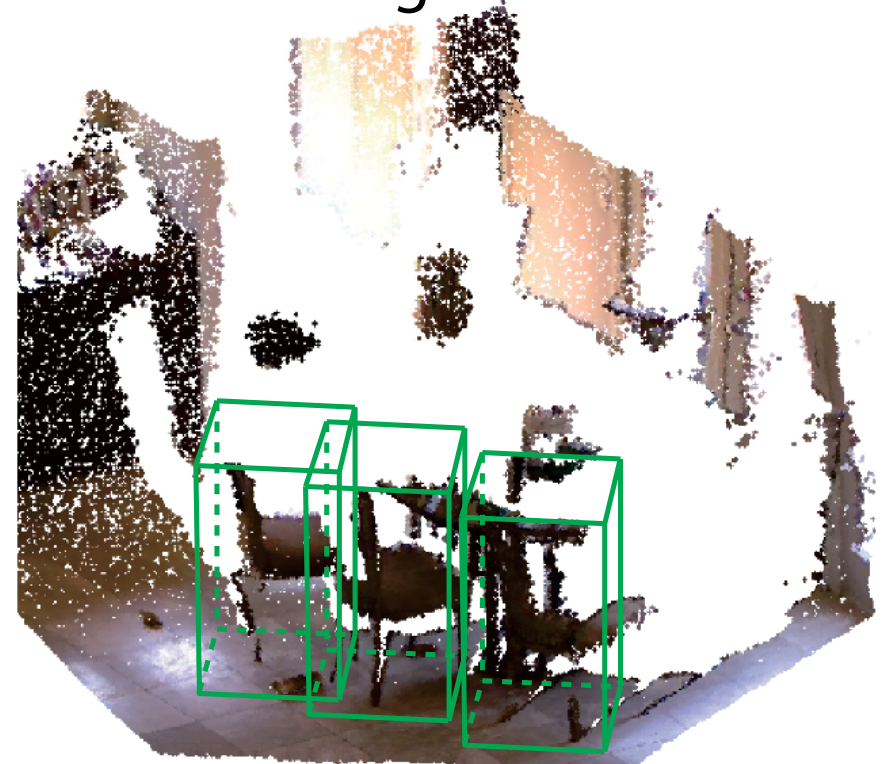
# Advantage: Exploiting Physical Size

2D Sliding Window



Multi-scale searching

3D Sliding Window



Physical size

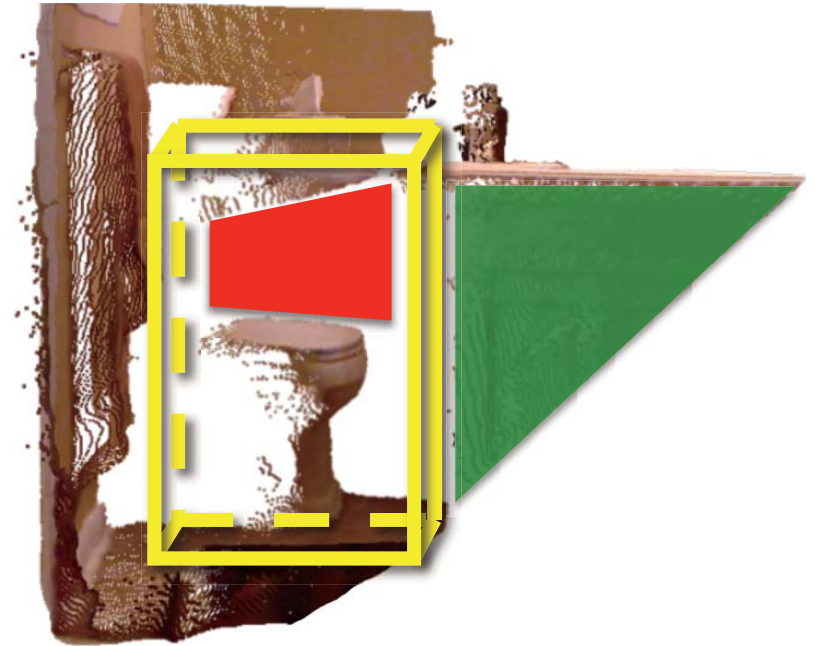


# Advantage: Handling Occlusion

2D Sliding window



3D Sliding window



Using depth, we can know which part is **occluded**.

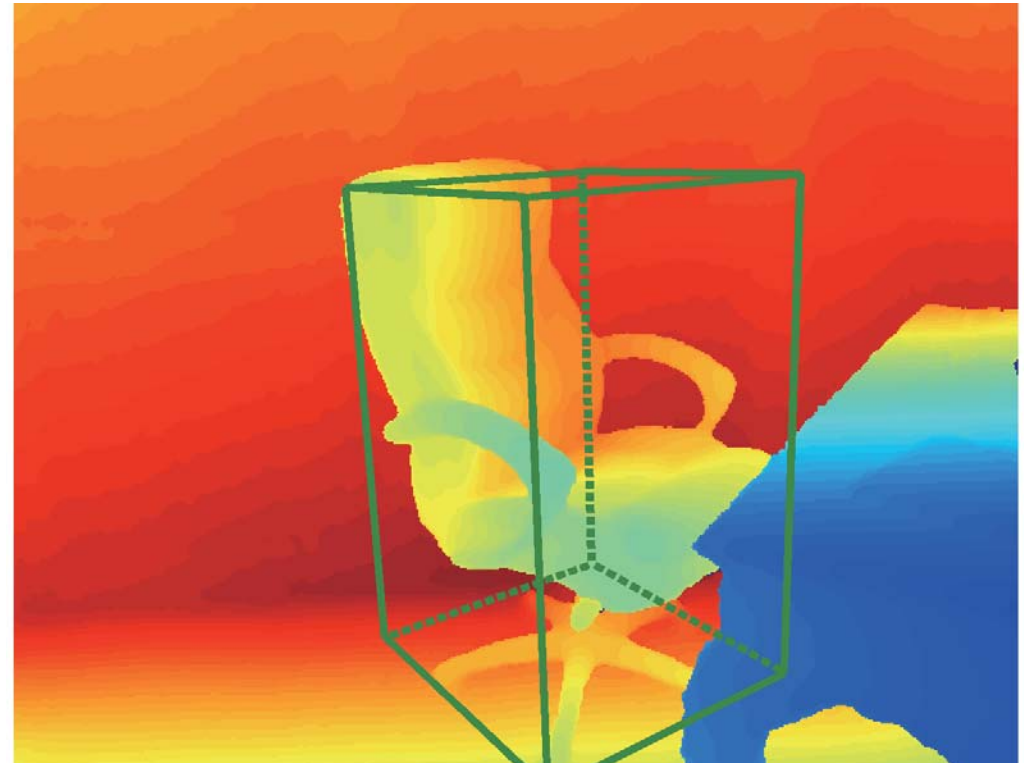
In 3D, we can separate the object from the **occluder**.

# Advantage: Insensitivity to Lighting

Color based detector: miss



Sliding Shapes

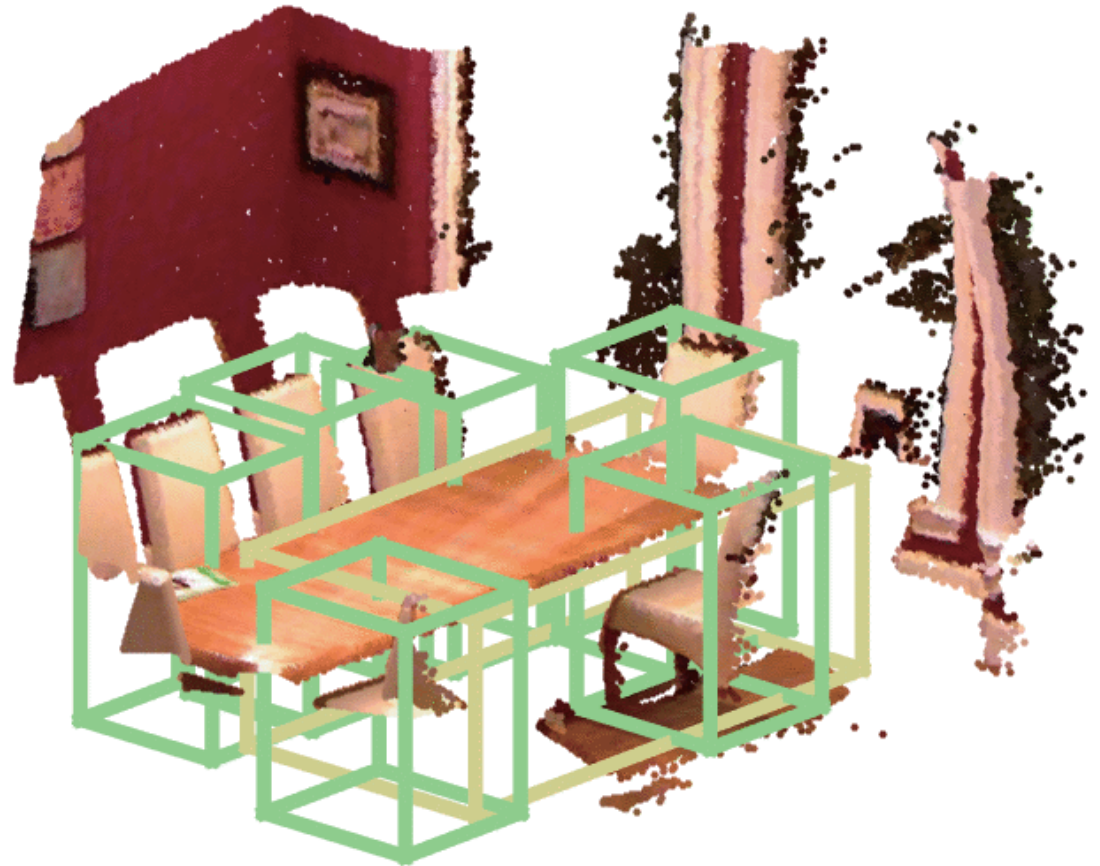
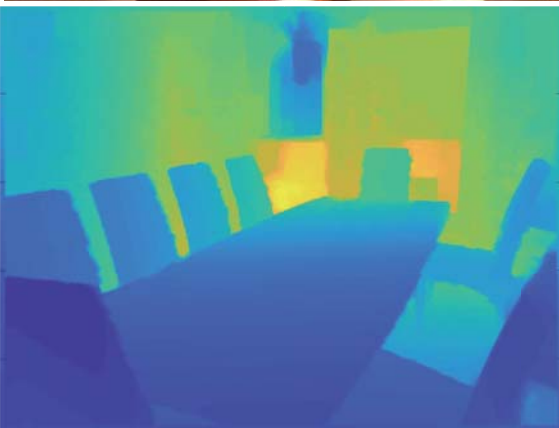


# Results: Deep Sliding Shapes

Color Image



Depth Image



Input: Single RGB-D

Output: 3D Amodal Boxes

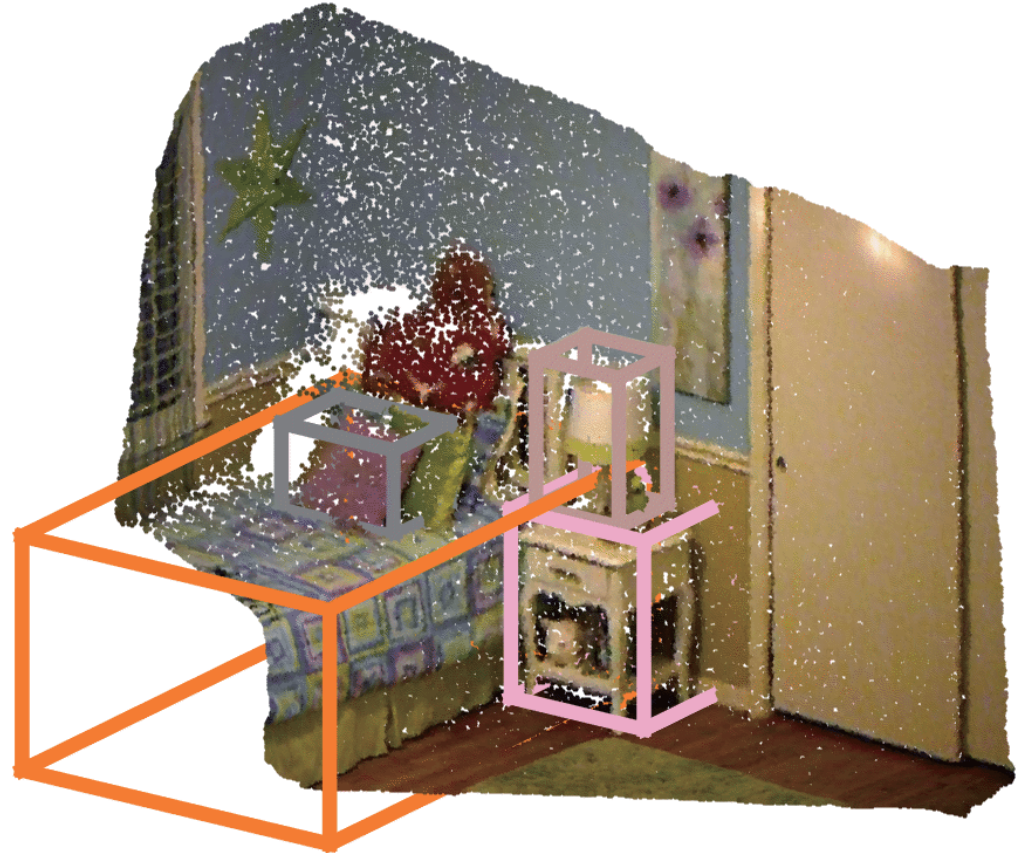
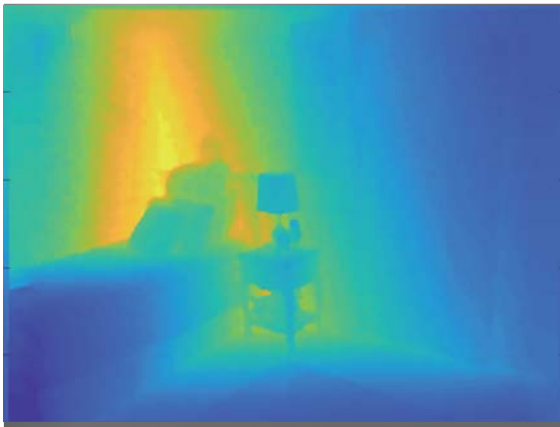


# Results: Deep Sliding Shapes

Color Image



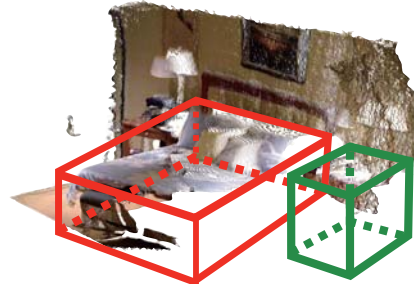
Depth Image



Input: Single RGB-D

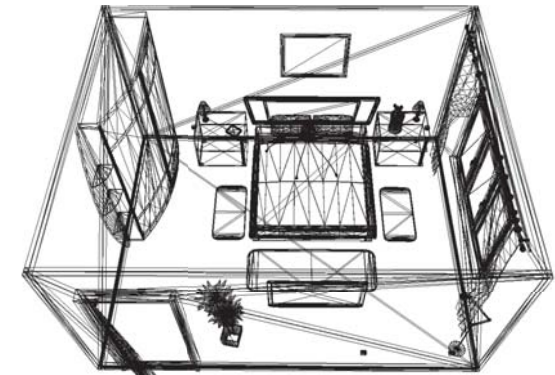
Output: 3D Amodal Boxes

# Data-Driven 3D Scene Understanding



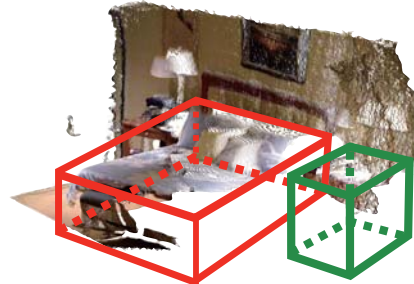
## Amodal 3D

[Song and Xiao  
ECCV'14, CVPR'16]



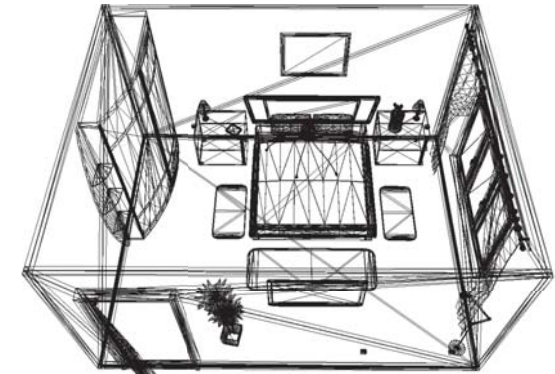
- Semantics Category
- 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. properties ...

# Data-Driven 3D Scene Understanding

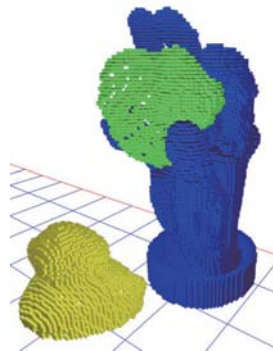


## Amodal 3D

[Song and Xiao  
ECCV'14, CVPR'16]



• Only Boxes, No Detailed Geometry



Not sufficient

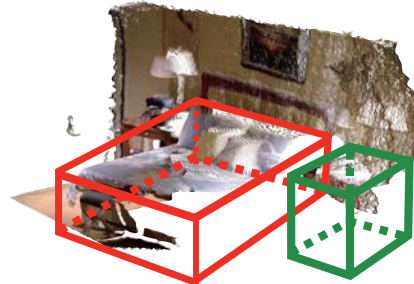
✓ **Semantics Category**

✓ **3D Location, Size**

- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. properties ...



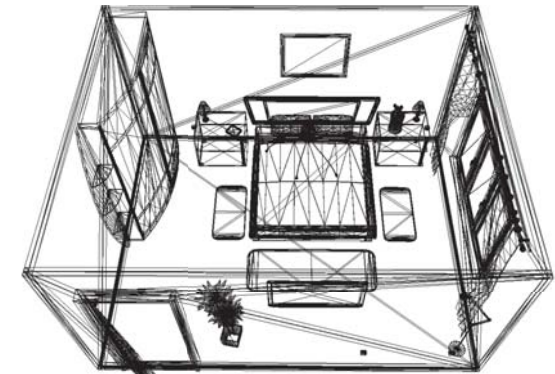
# Data-Driven 3D Scene Understanding



## Amodal 3D

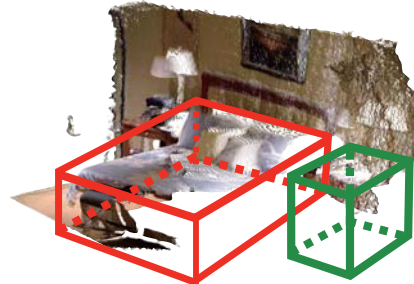
[Song and Xiao  
ECCV'14, CVPR'16]

- Only Boxes, No Detailed Geometry
- Single Object, No Contextual Information



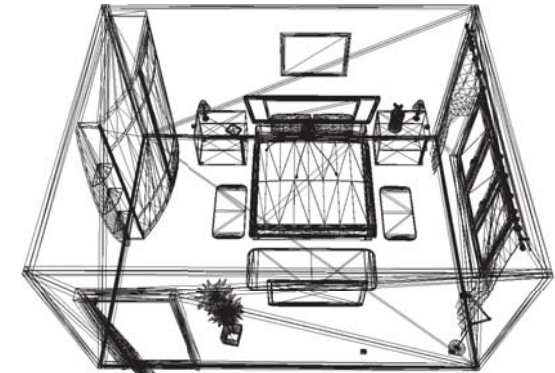
- ✓ **Semantics Category**
- ✓ **3D Location, Size**
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. properties ...

# Data-Driven 3D Scene Understanding



## Amodal 3D

[Song and Xiao  
ECCV'14, CVPR'16]



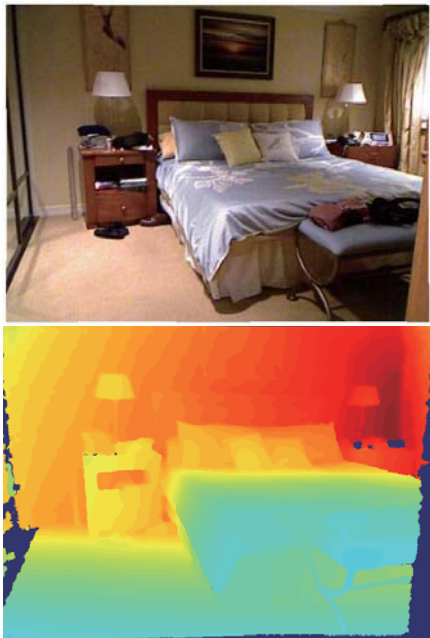
- ✓ Semantics Category
- ✓ 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. properties ...



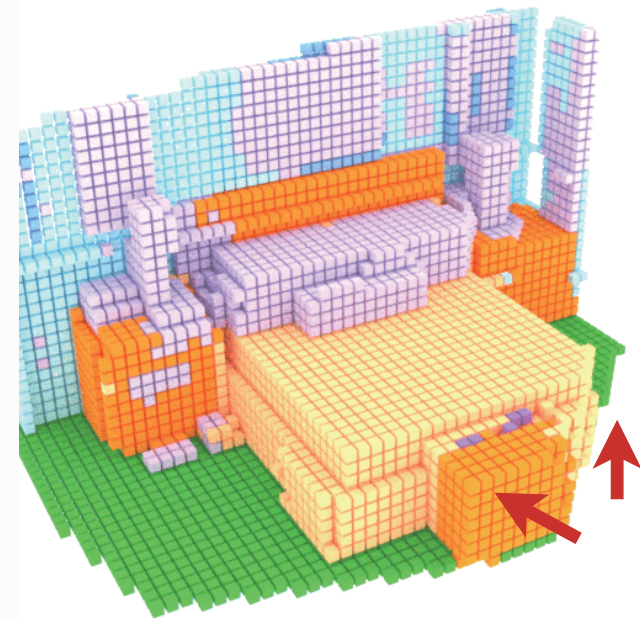
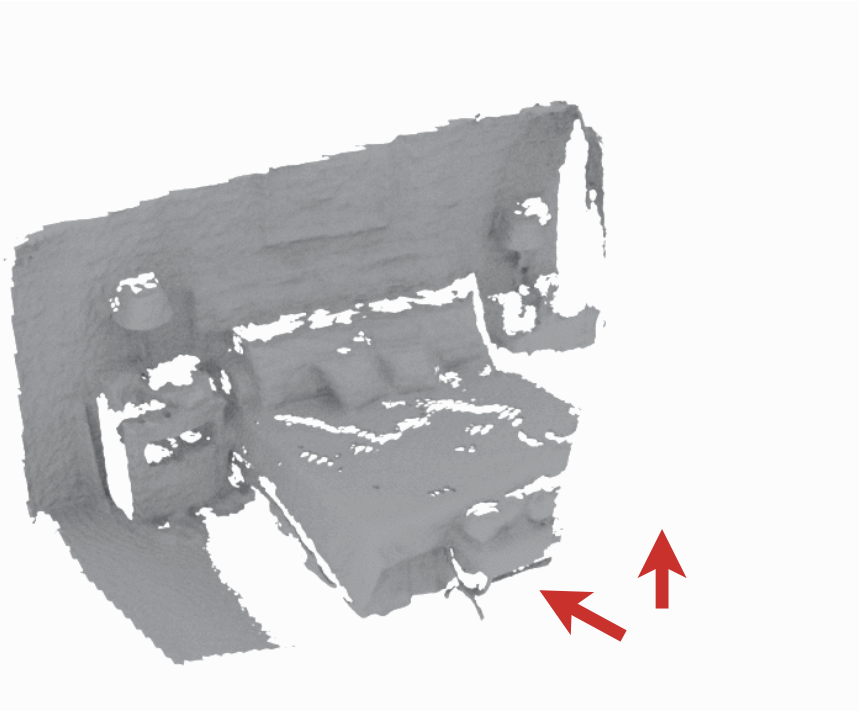
## Higher Fidelity

[Song et al. CVPR'17]

# Completed 3D Scene Completion

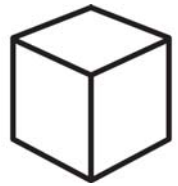


Input:  
Single Depth Map

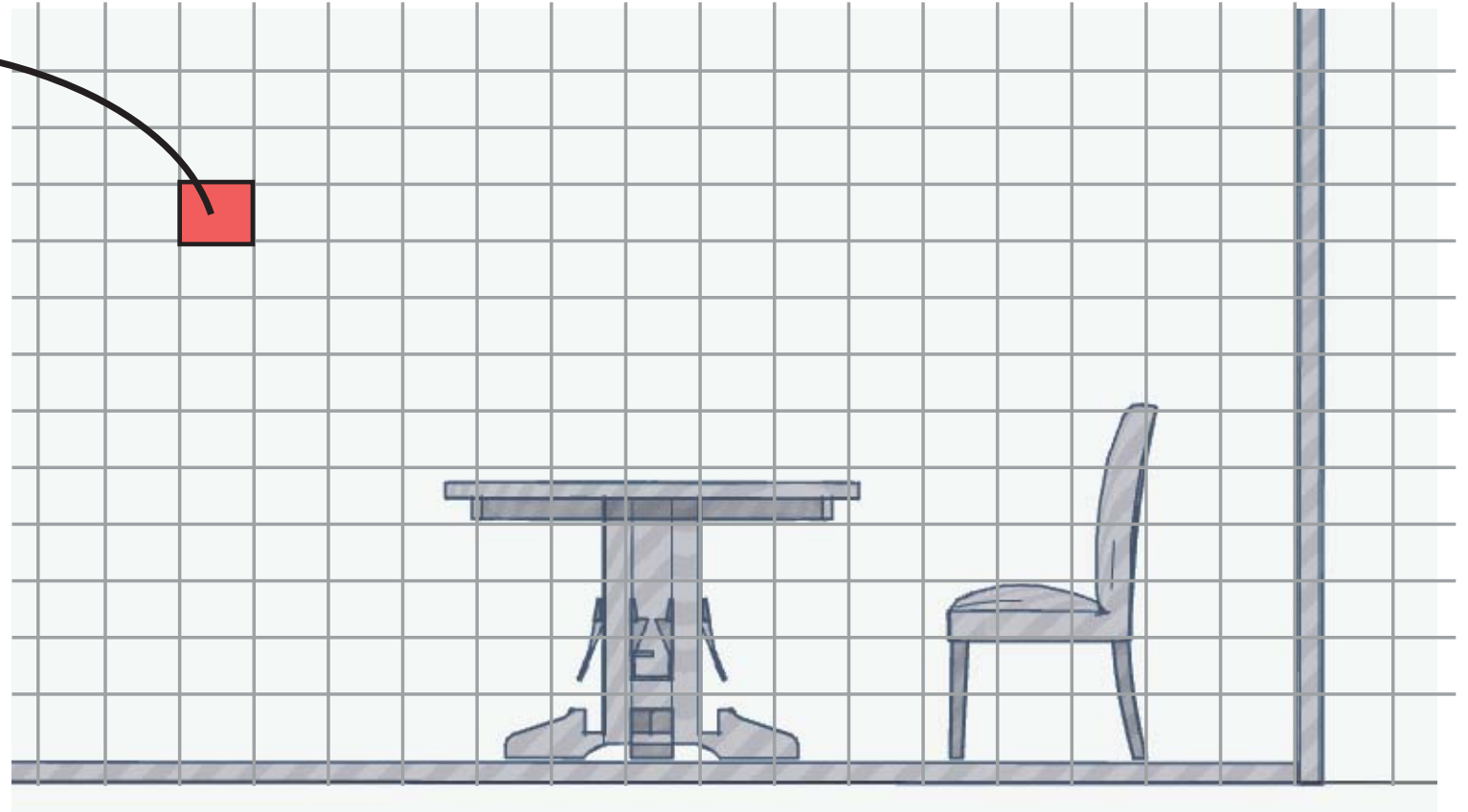


Output:  
Volumetric Occupancy + Semantic

# Problem Definition

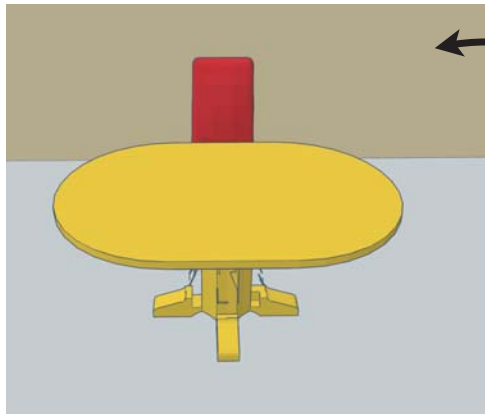






Voxel

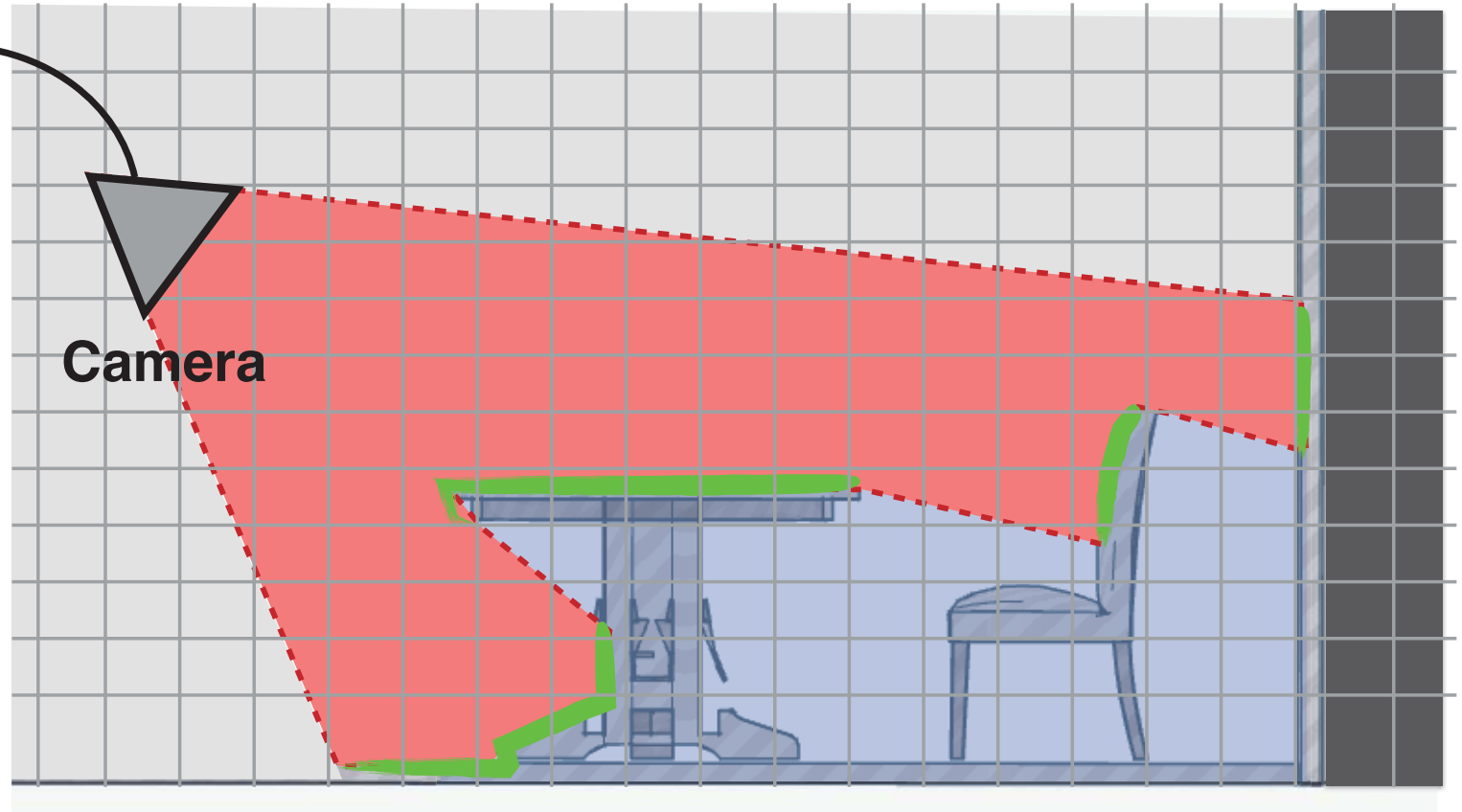


3D Scene

# Problem Definition

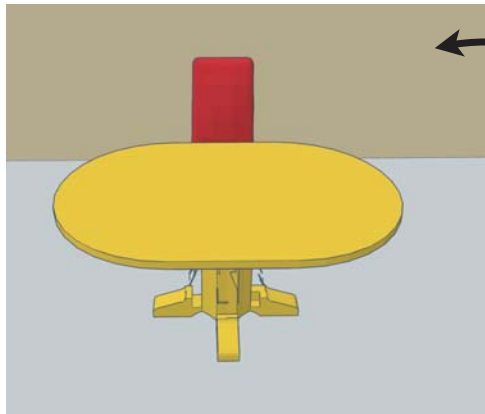


-  visible surface
-  free space
-  occluded space
-  outside view



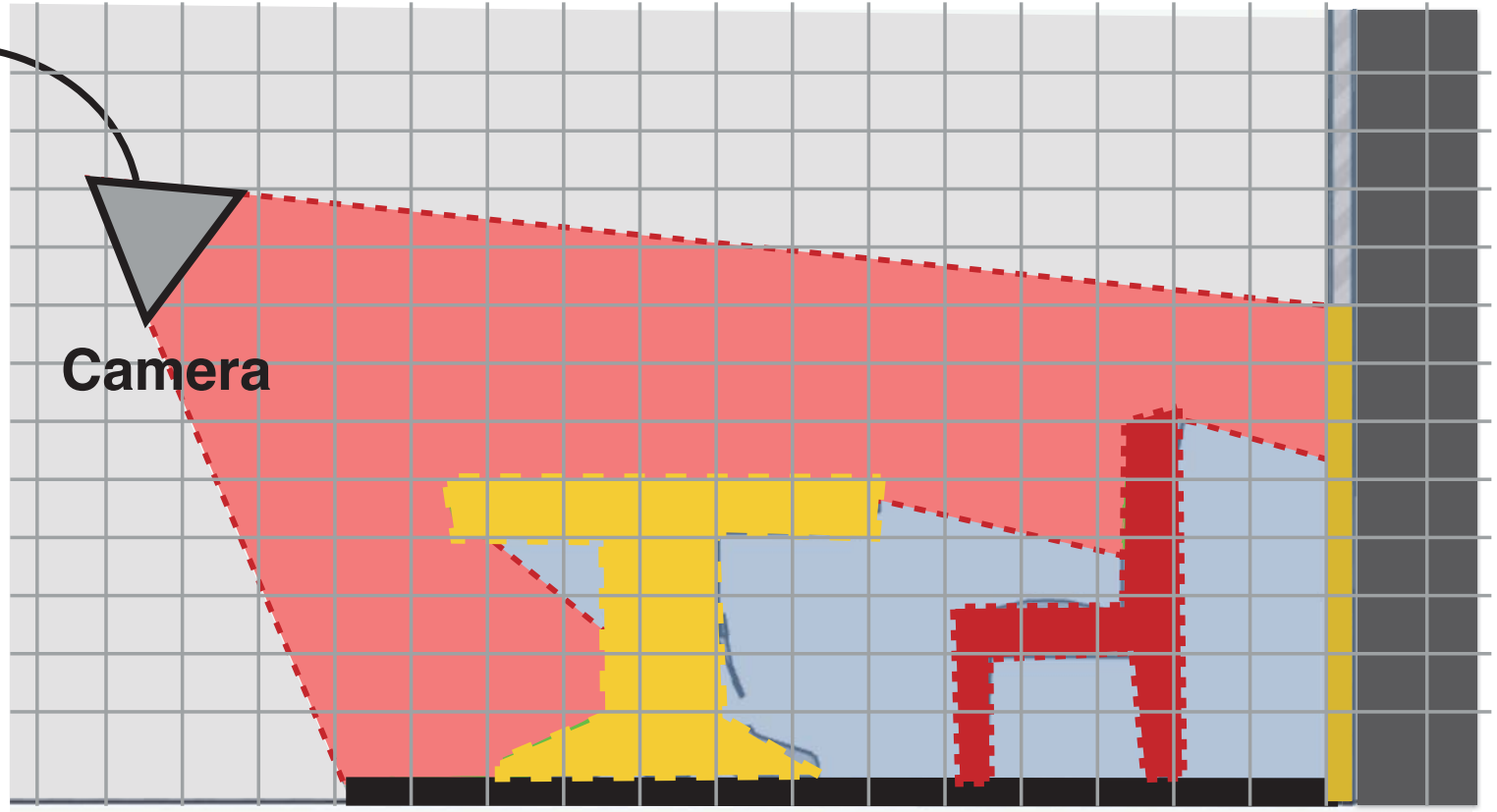
3D Scene

# Problem Definition



Camera

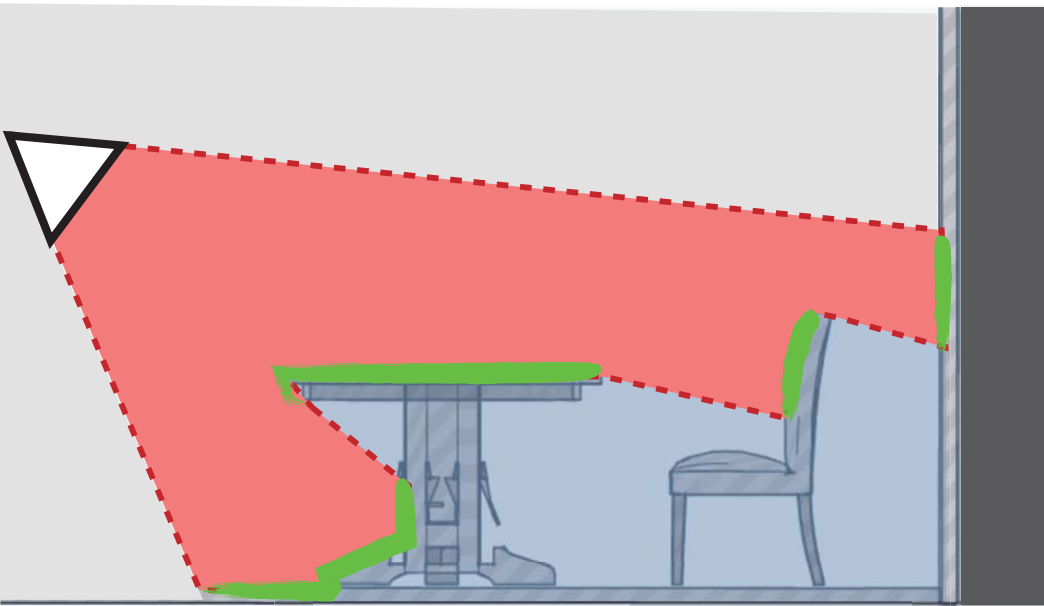
- visible surface
- free space
- occluded space
- outside view



3D Scene



# Problem Definition



3D Scene



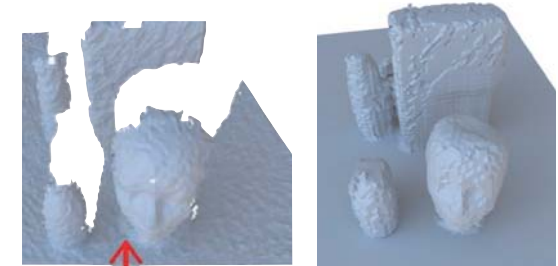
Image Segmentation



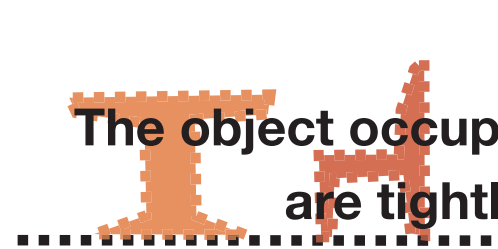
[Long *et al.* CVPR'15]



Shape Completion



[Firman *et al.* CVPR'17]

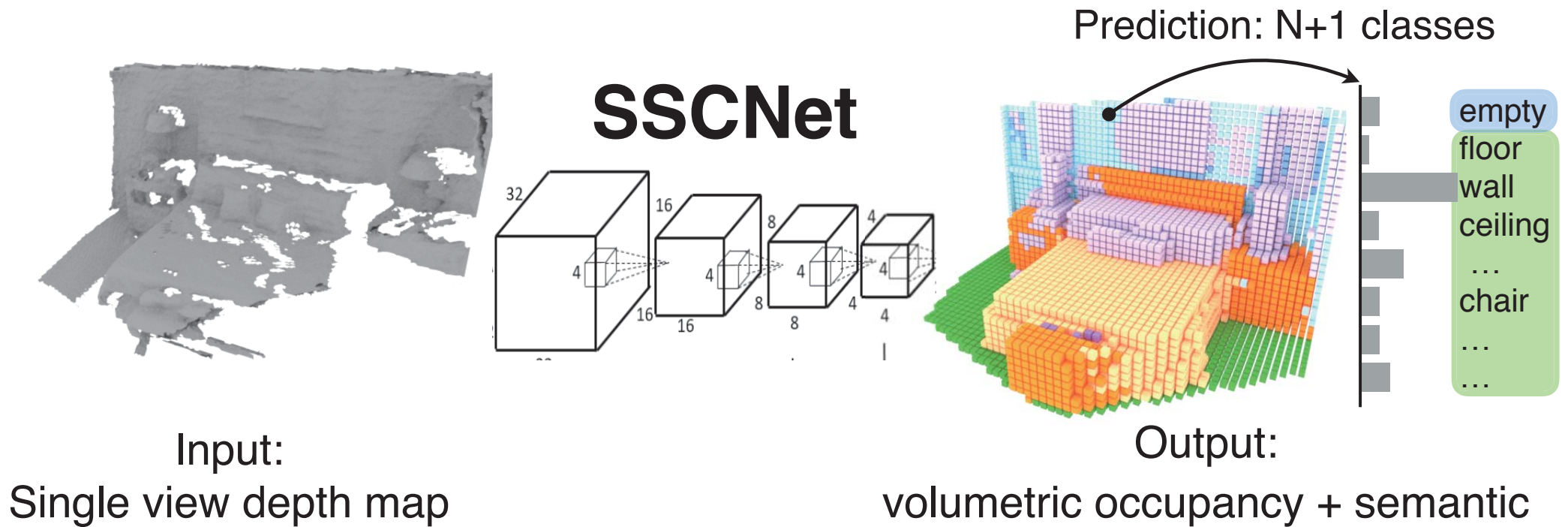


Semantic Scene Completion

**The object occupancy and their identity are tightly intertwined!**

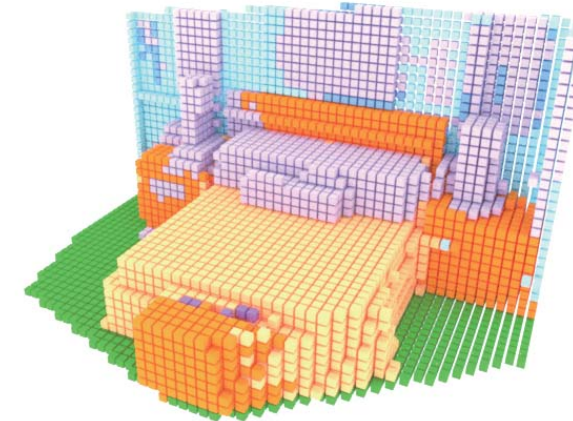
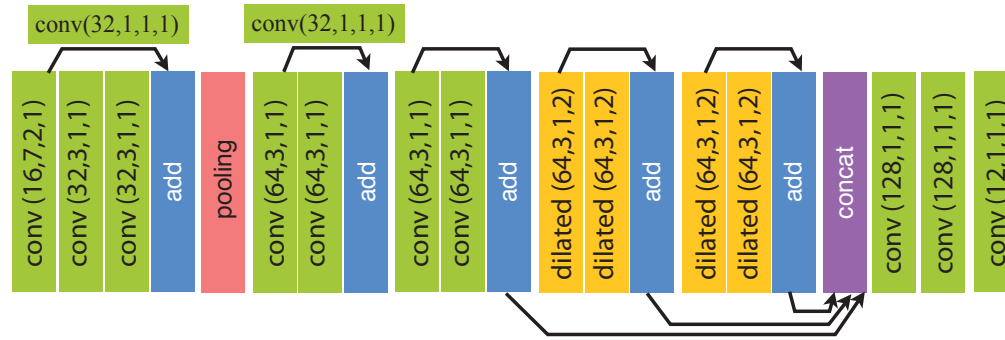
**This paper**

# Semantic Scene Completion Network

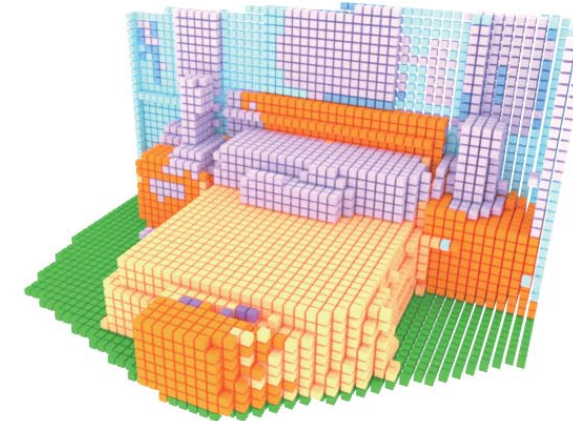
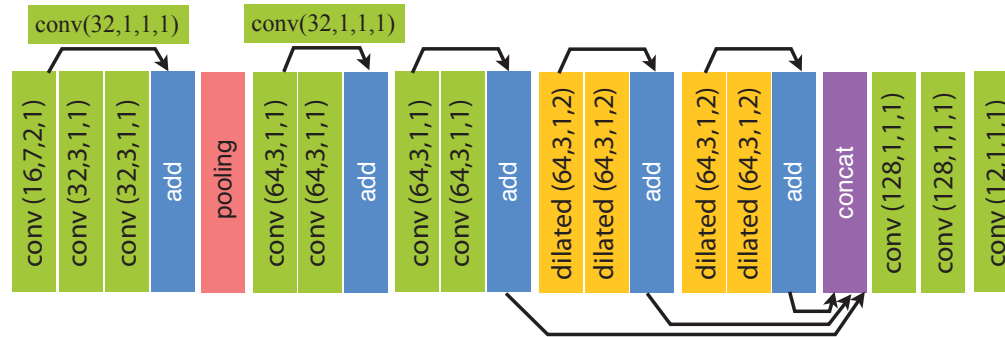
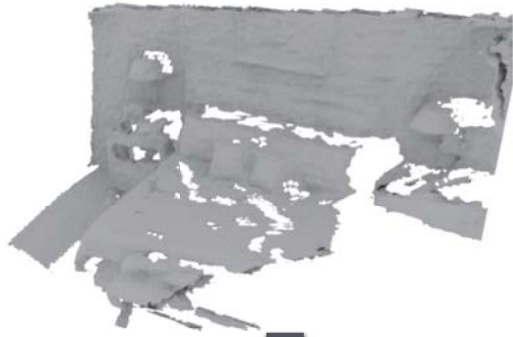


Simultaneously predict voxel occupancy and semantics classes by a single forward pass.

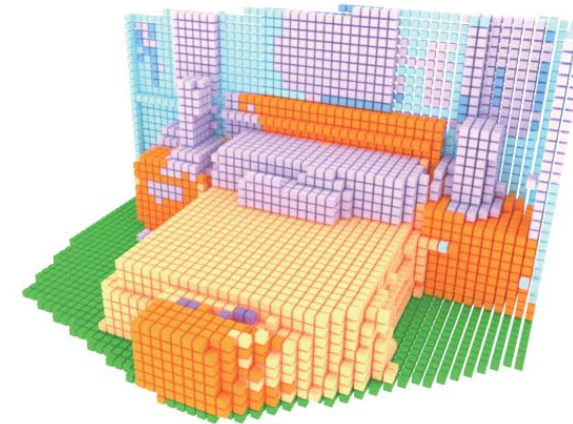
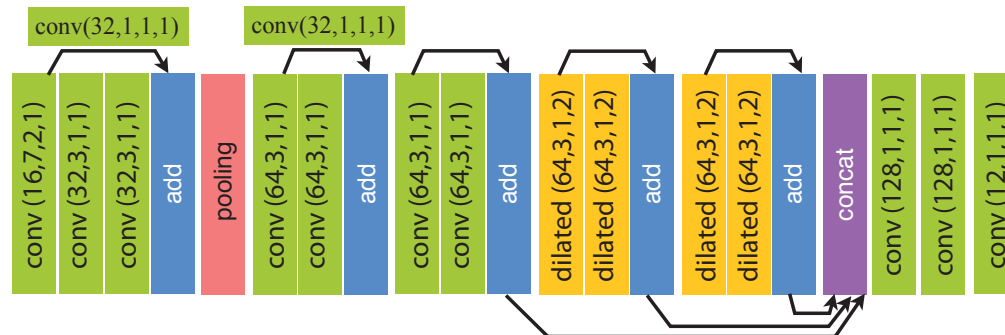
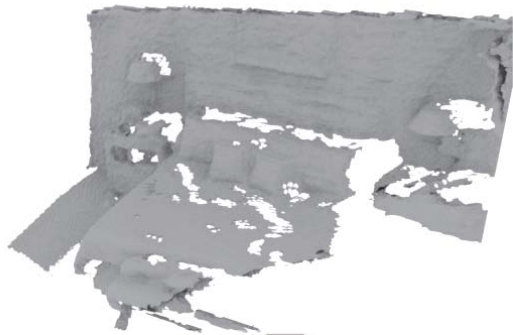
# Semantic Scene Completion Network



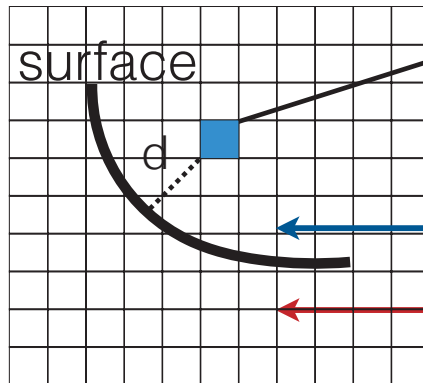
# Semantic Scene Completion Network



# Semantic Scene Completion Network



Encode 3D space using flipped TSDF



$d$  = distance to the surface

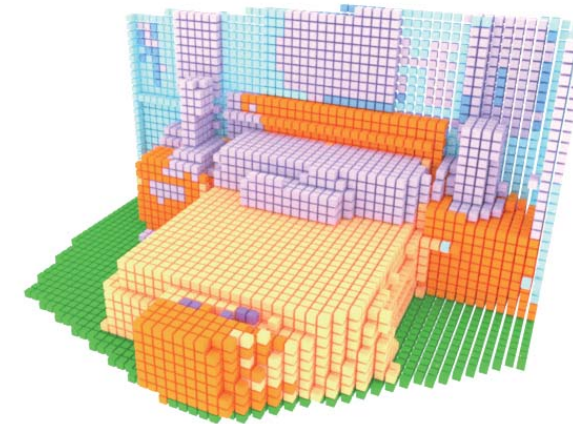
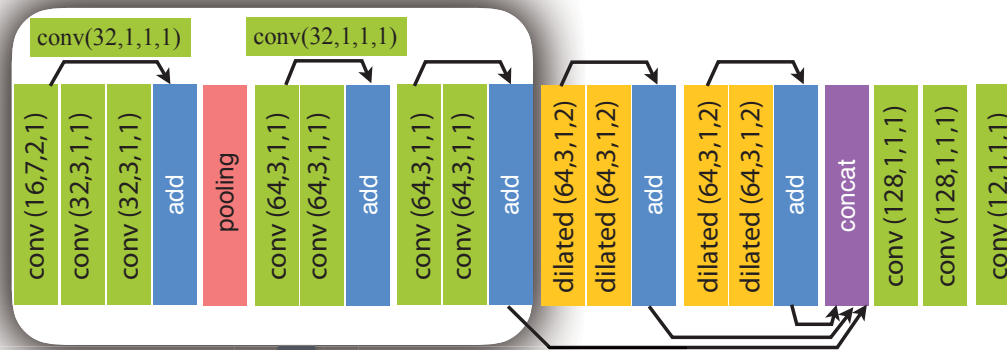
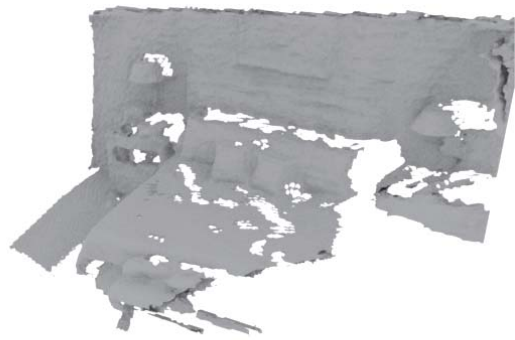
$$\text{flipped TSDF} = \text{sign}(1 - \min(1, d/d_{\text{max}}))$$

Occluded

Free space



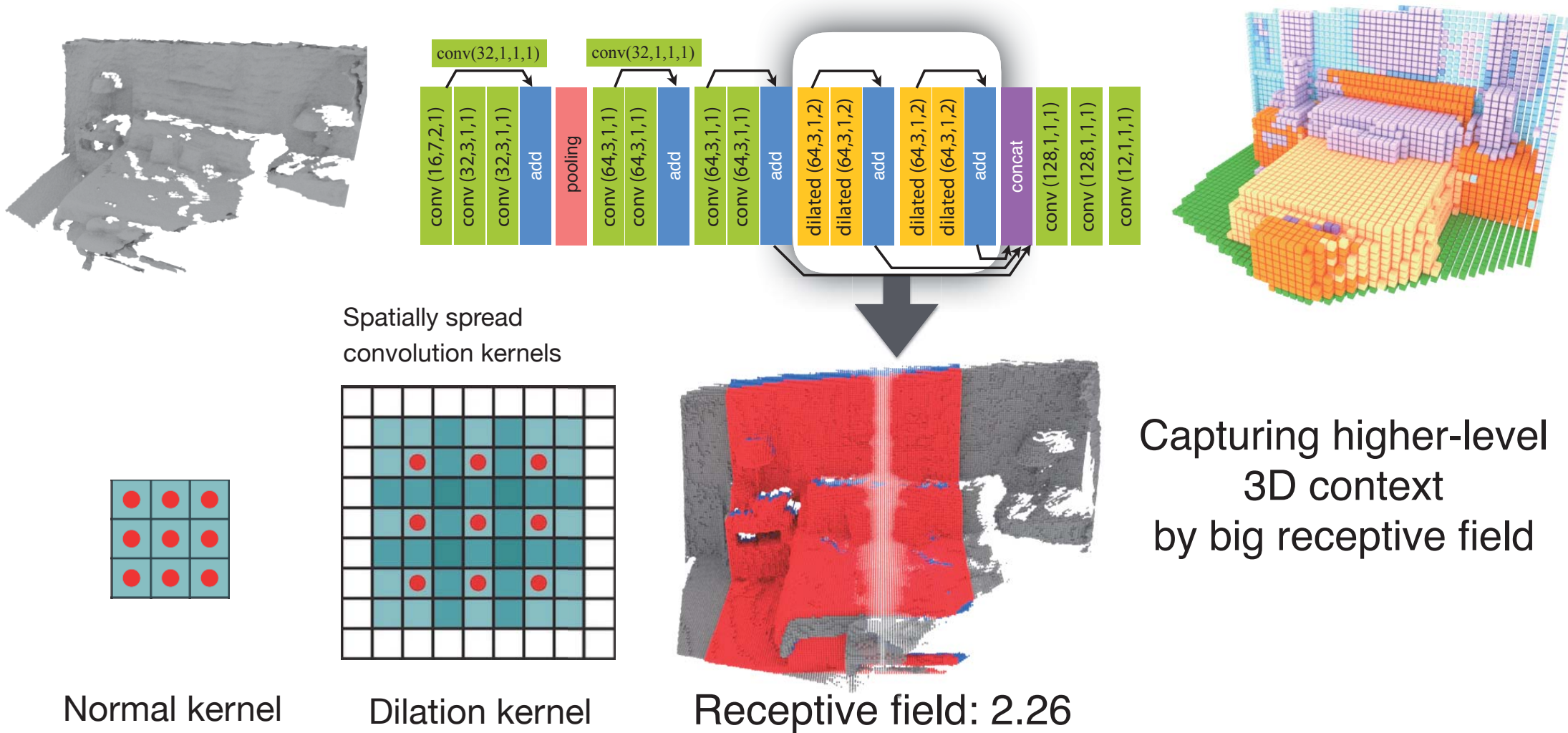
# Semantic Scene Completion Network



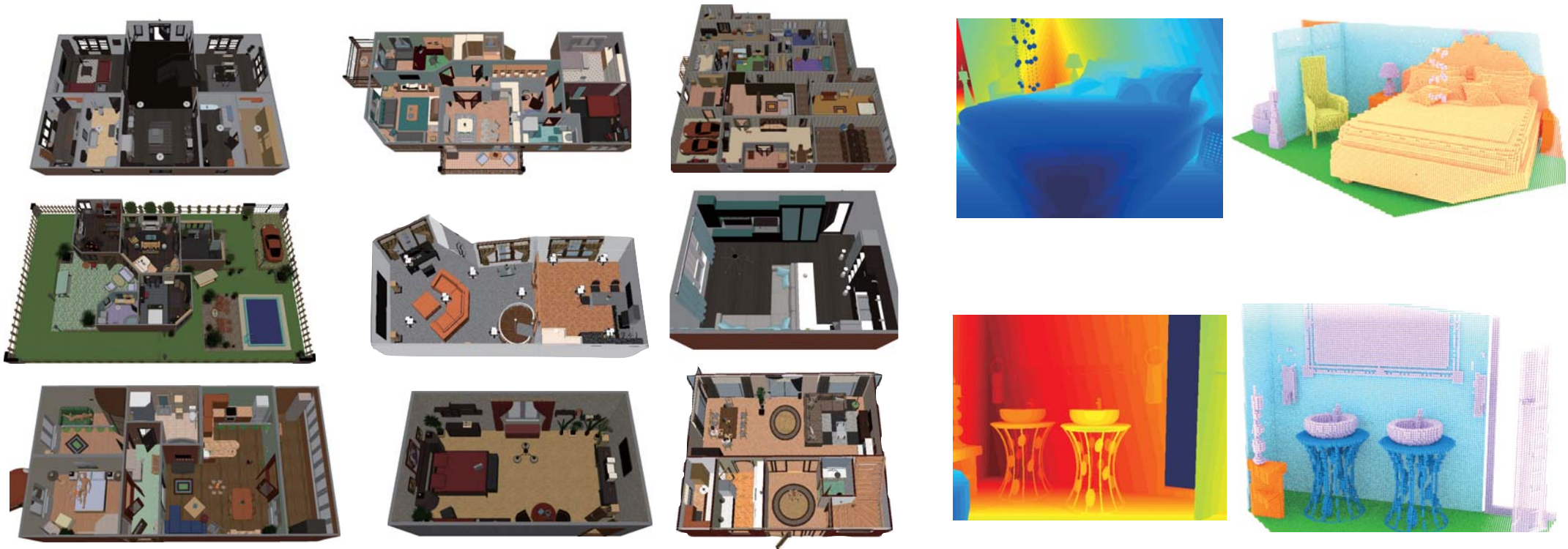
Local geometry

Receptive field: 0.98 m

# Semantic Scene Completion Network



# Train on Synthetic 3D Scenes



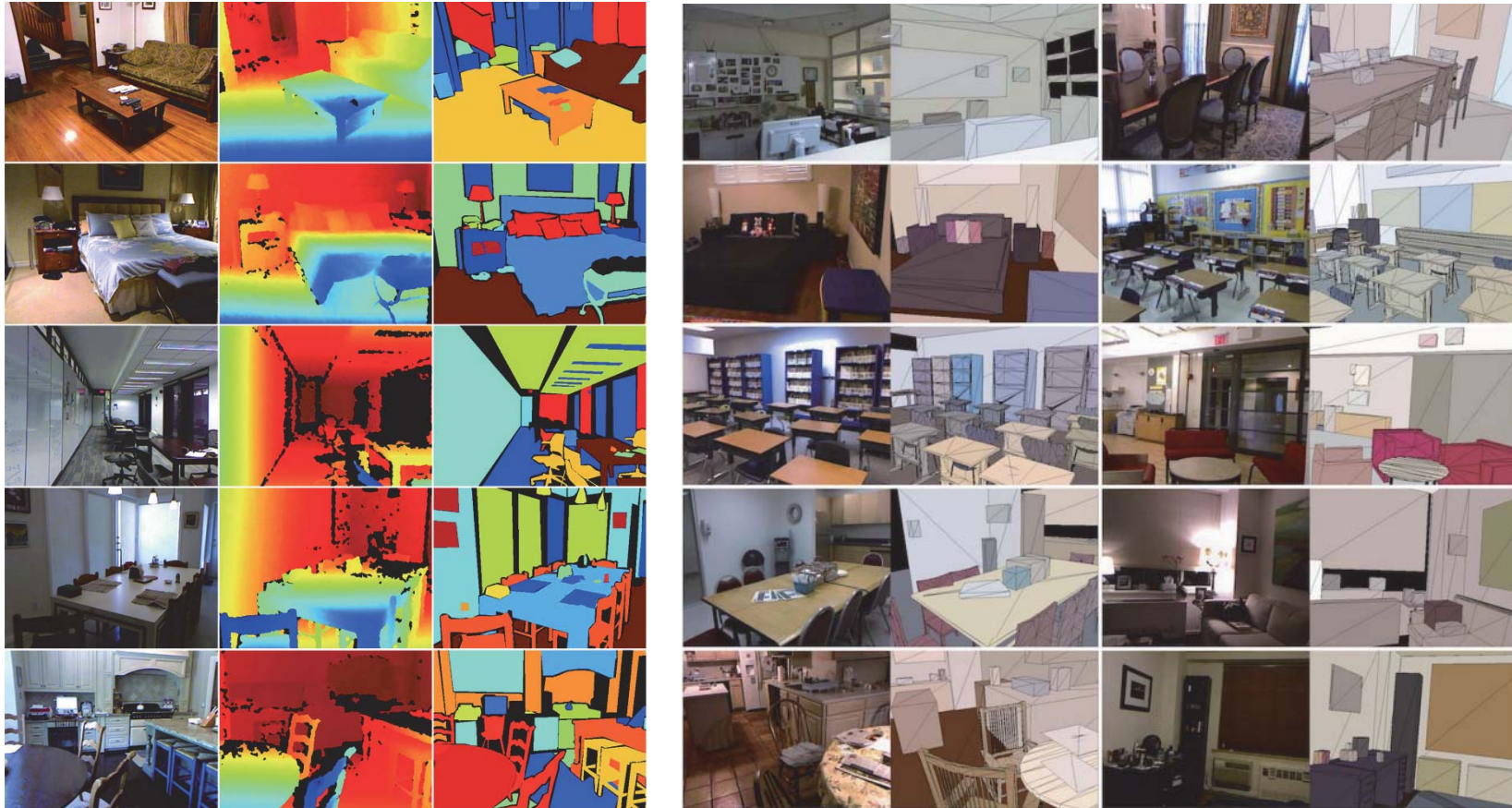
Synthetic Scenes (SUNCG)

Depth

Ground Truth



# Testing on Real-World Data (NYU [1,2])



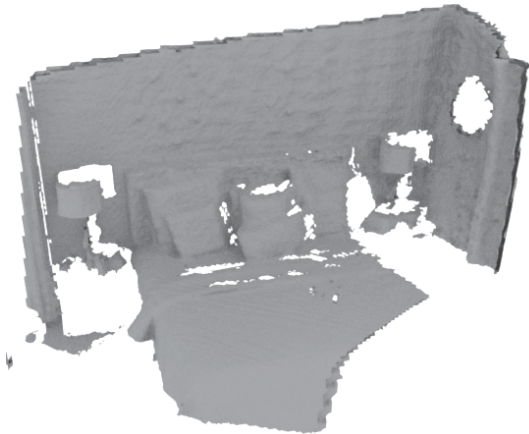
[1] NYU depth v2: Silberman et al. ECCV'12

[2] Ground truth: Guo and Hoiem IJCV'15

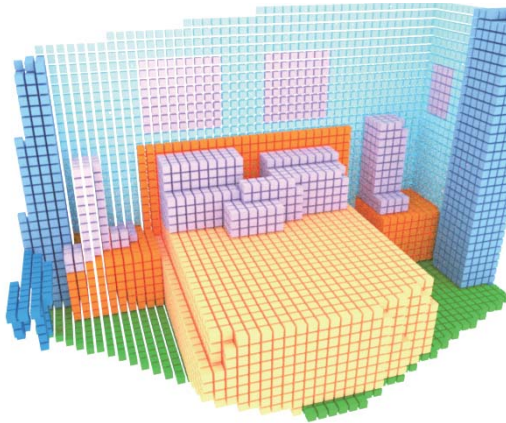
# Comparis on Real-world Dataset



# Comparison on Real-world Dataset

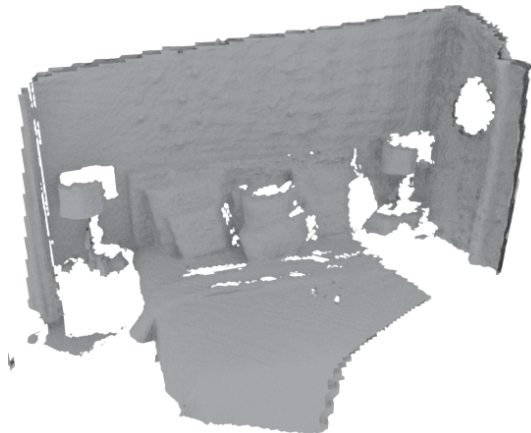


Observed Surface

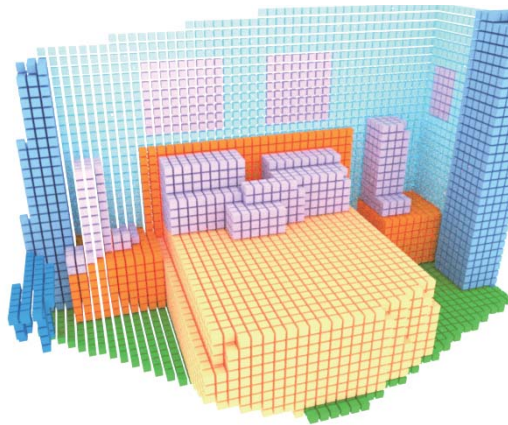


Ground Truth

# Comparison on Real-world Dataset

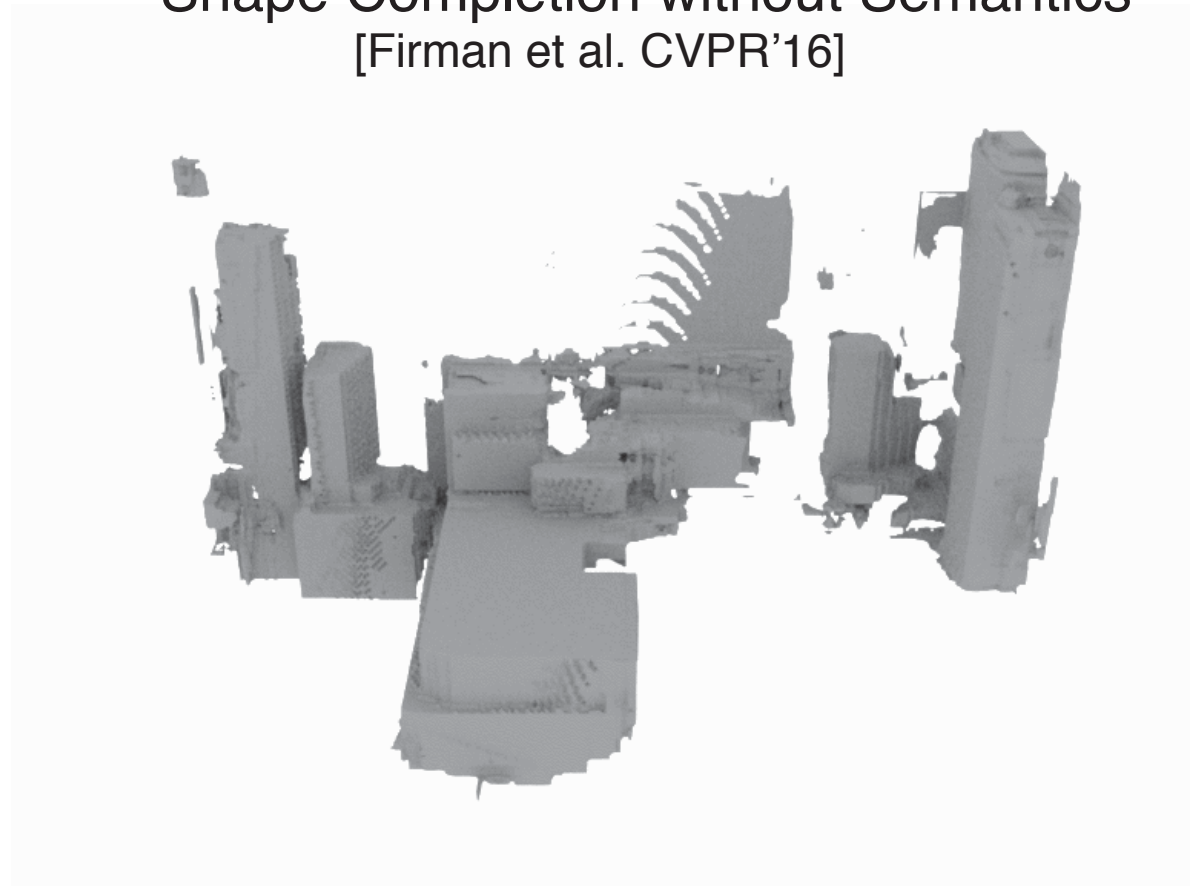


Observed Surface

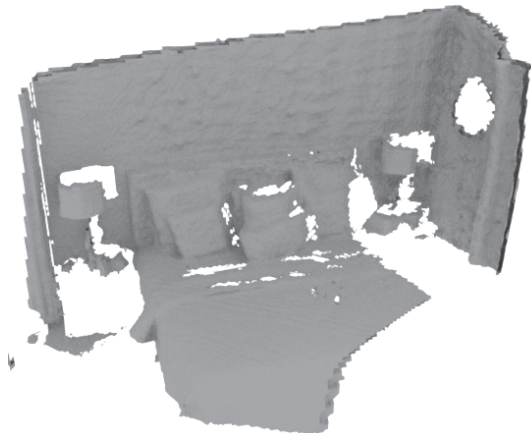


Ground Truth

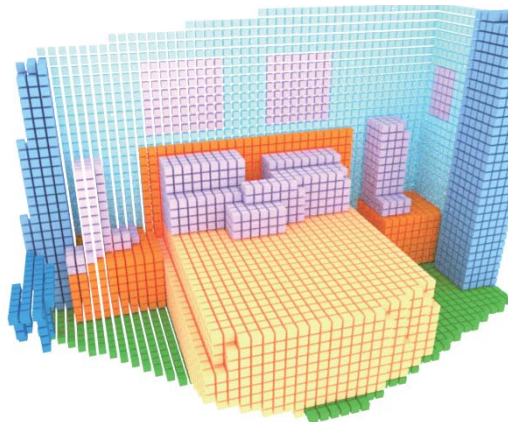
Shape Completion without Semantics  
[Firman et al. CVPR'16]



# Comparison on Real-world Dataset



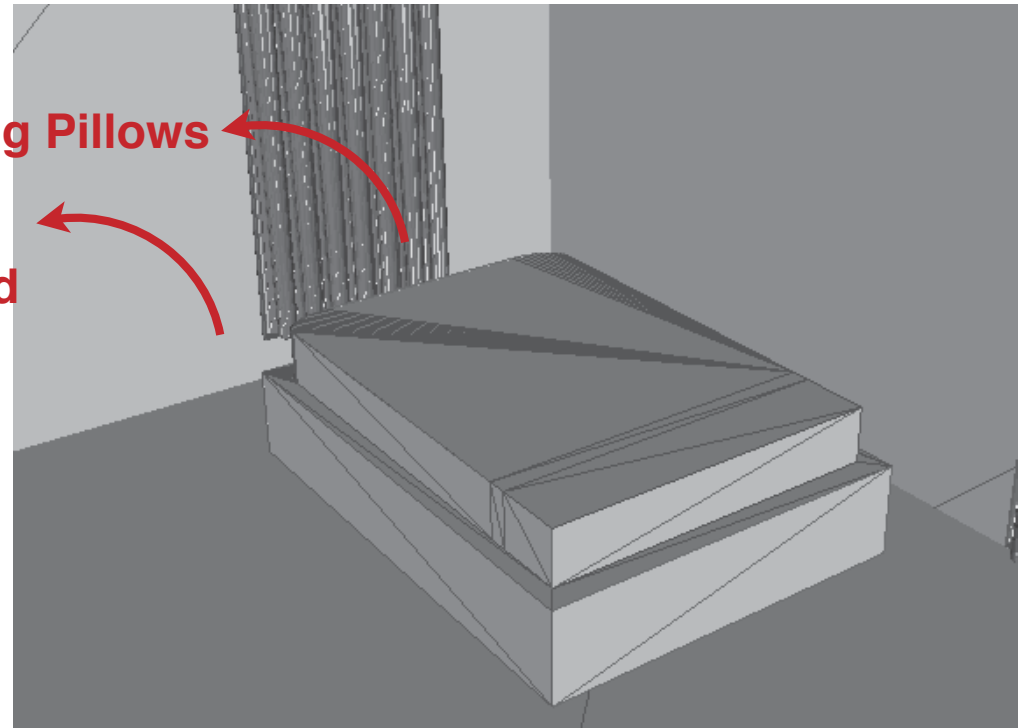
Observed Surface



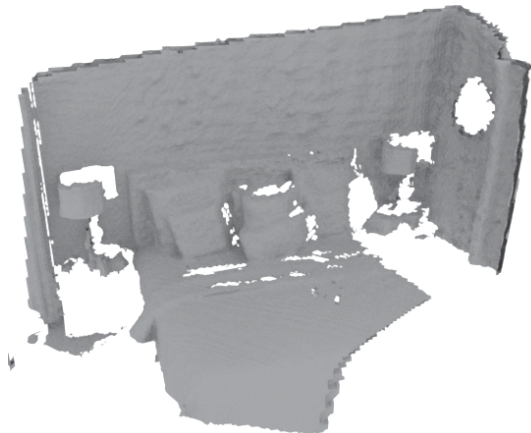
Ground Truth

Model Retrieval+Fitting  
[Geiger and Wang GCPR'15]

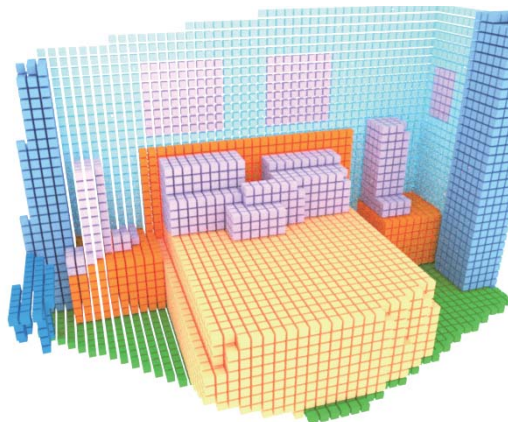
Missing Pillows  
Missing Nightstand



# Comparison on Real-world Dataset

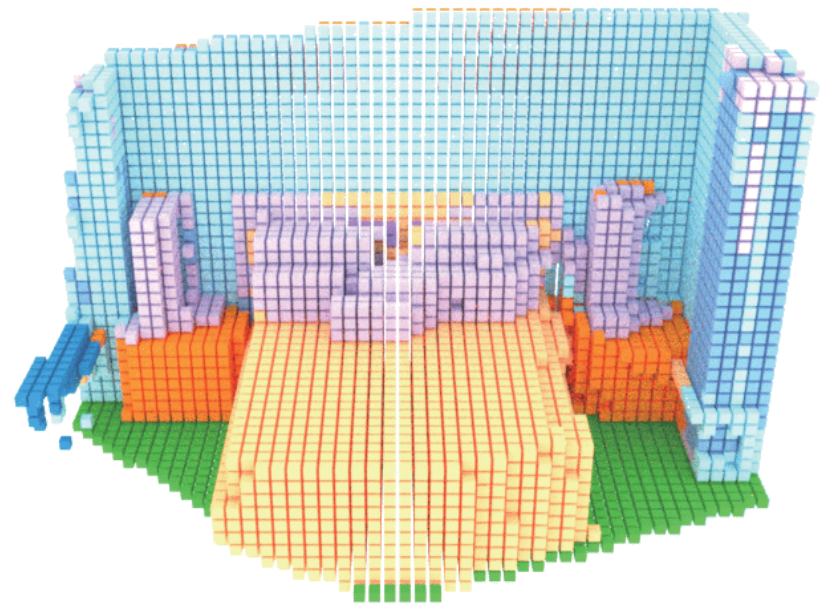


Observed Surface



Ground Truth

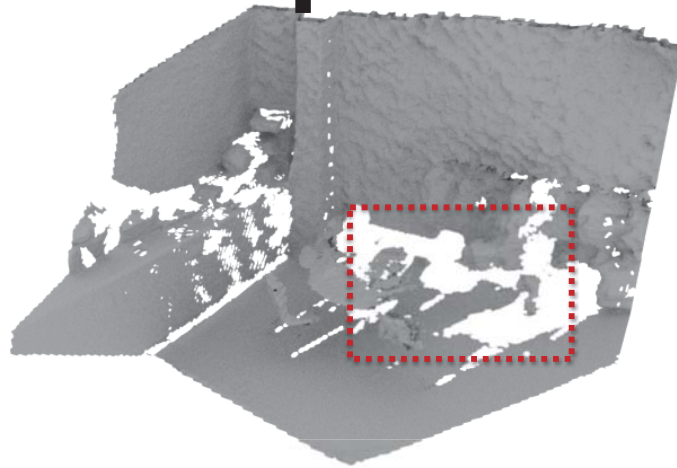
SSCNet  
[Ours]



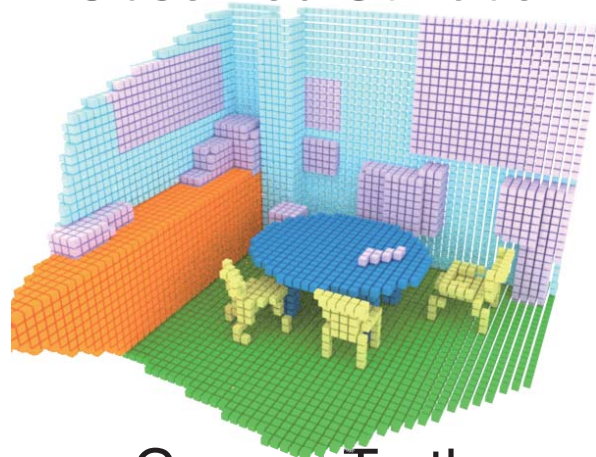
- floor
- wall
- window
- chair
- bed
- sofa
- table
- tv
- furn.
- objects



# Comparison on Real-world Dataset

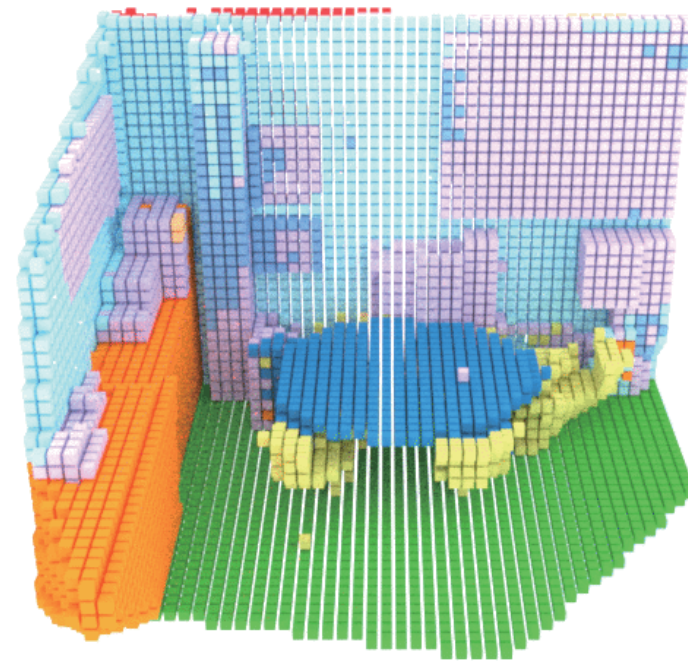


Observed Surface



Ground Truth

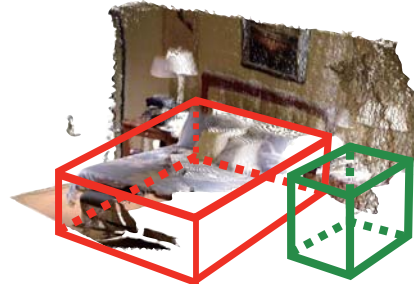
SSCNet  
[Ours]



- floor
- wall
- window
- chair
- bed
- sofa
- table
- tv
- furn.
- objects

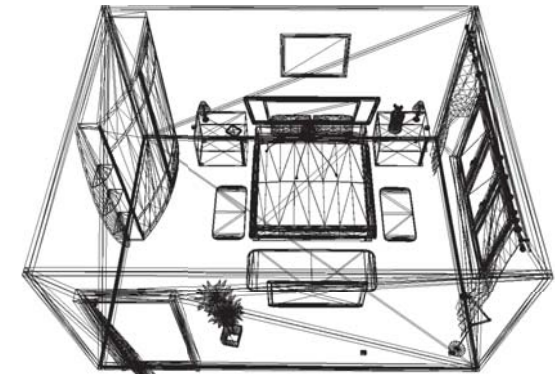


# Data-Driven 3D Scene Understanding



## Amodal 3D

[Song and Xiao  
ECCV'14, CVPR'16]



- ✓ Semantics Category
- ✓ 3D Location, Size
- ✓ Detailed Geometry
- ✓ Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. properties ...

Prediction is limited by  
Camera Field of View



## Higher Fidelity

[Song et al. CVPR'17]

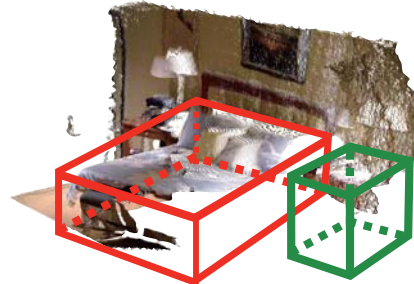
# Limited Camera FoV

Typical camera  
FoV 60 degree

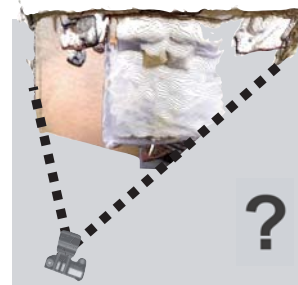
Small Portion of the Scene is  
Observed due to Limited FoV



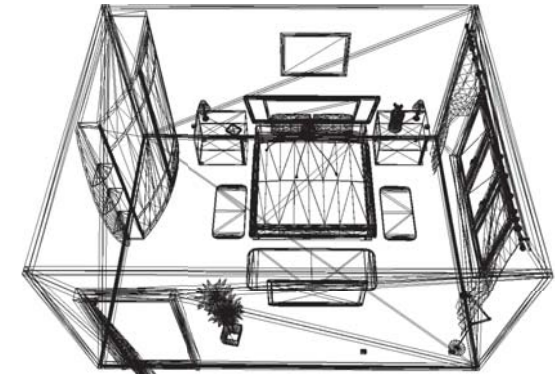
# Data-Driven 3D Scene Understanding



**Amodal 3D**  
[Song and Xiao  
ECCV'14,CVPR'16]



**Beyond FoV**  
[song et al. CVPR'18]



- ✓ Semantics Category
- ✓ 3D Location, Size
- ✓ Detailed Geometry
- ✓ Inter-Object Relationships
- **Not Limited by FoV**
- Action Affordances
- Phys. properties ...



**Higher Fidelity**  
[Song et al. CVPR'17]





# View Extrapolation

Prior work: Predicting Scene Appearance (Only Colored Pixels)

Image Inpainting  
[Pathak et. al]



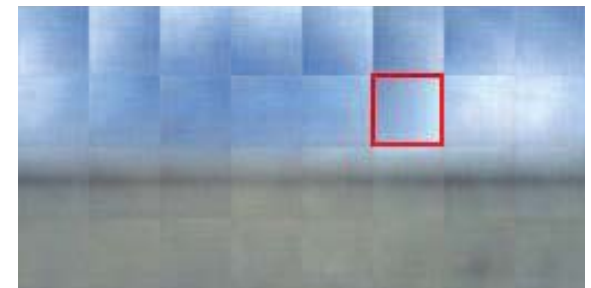
(a) Input context



User-guided view extrapolation [Zhang et al.]

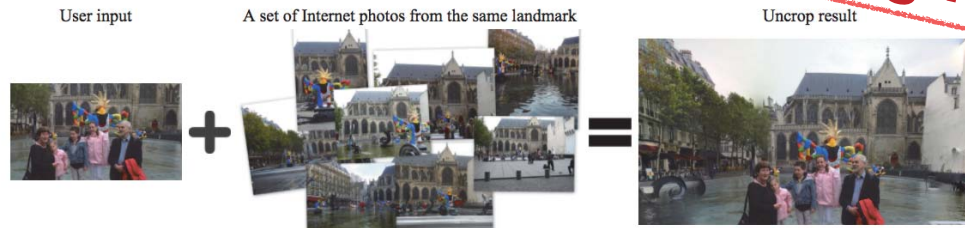


Learning to Look Around  
[Jayaraman and Grauman]



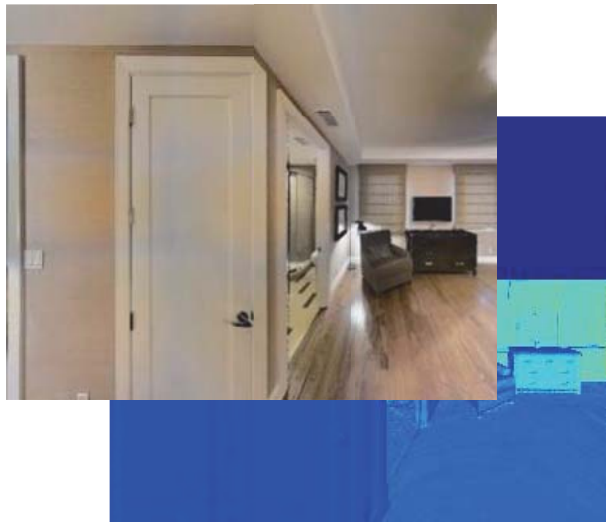
**Hard to be used directly to support high level planning**

Stitching images from the internet

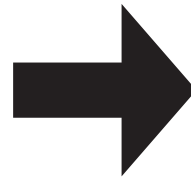




# View Extrapolation



Input: RGB-D images



Output1: 3D Structures



Output2: Semantics

# View Extrapolation

Where can I move?



Output1: 3D Structures

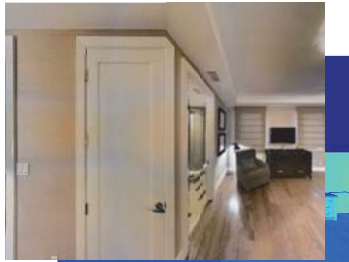
Where should I turn to find a door?



Output2: Semantics

# Semantic-Structure View Extrapolation

Input: RGB-D images



# Semantic-Structure View Extrapolation

Input: RGB-D images



Output: 360° panorama  
with 3D structure & semantics

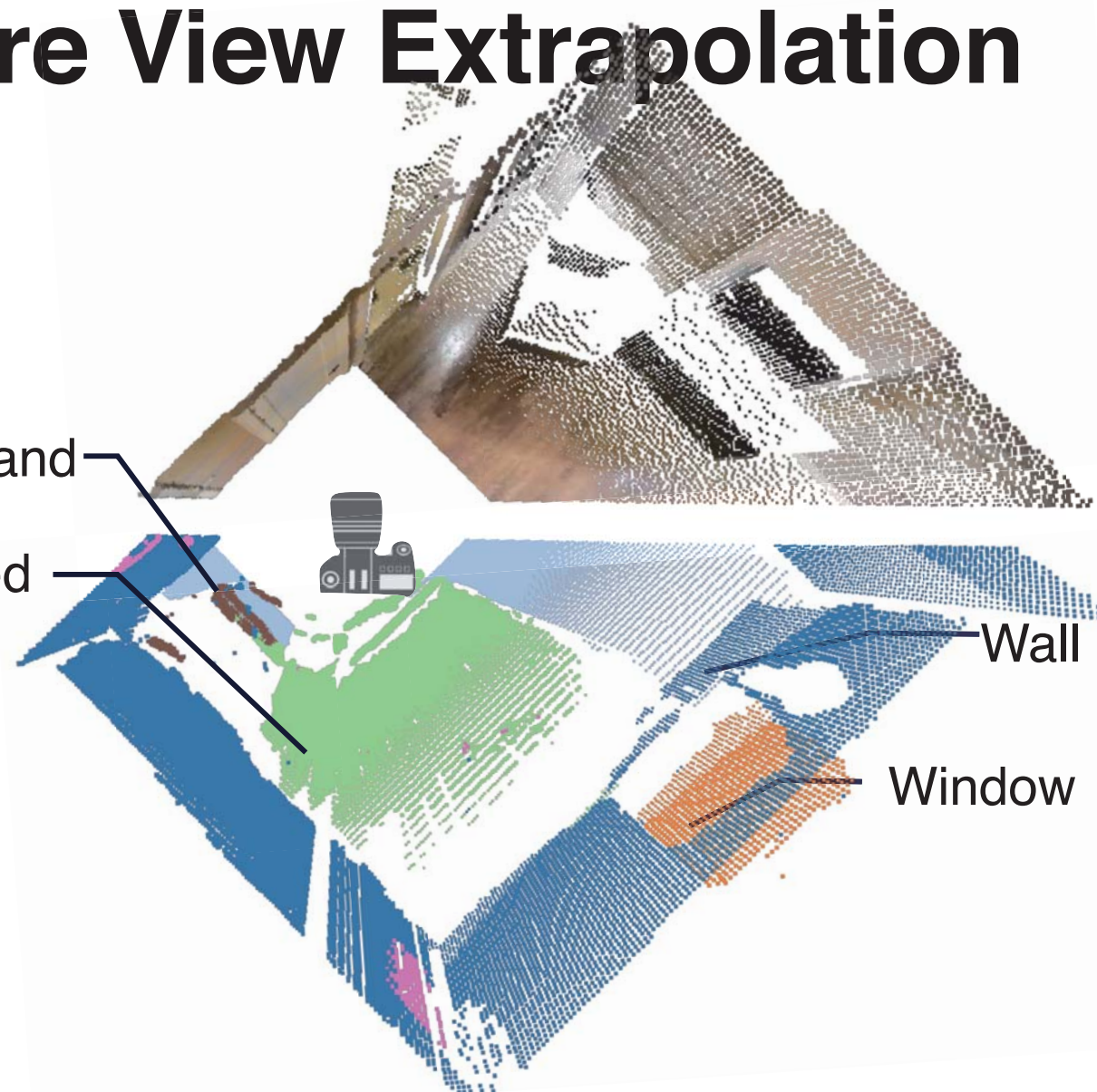


Nightstand

Bed

Wall

Window





# Semantic-Structure View Extrapolation

Input: RGB-D images



Output: 360° panorama with 3D structure & semantics



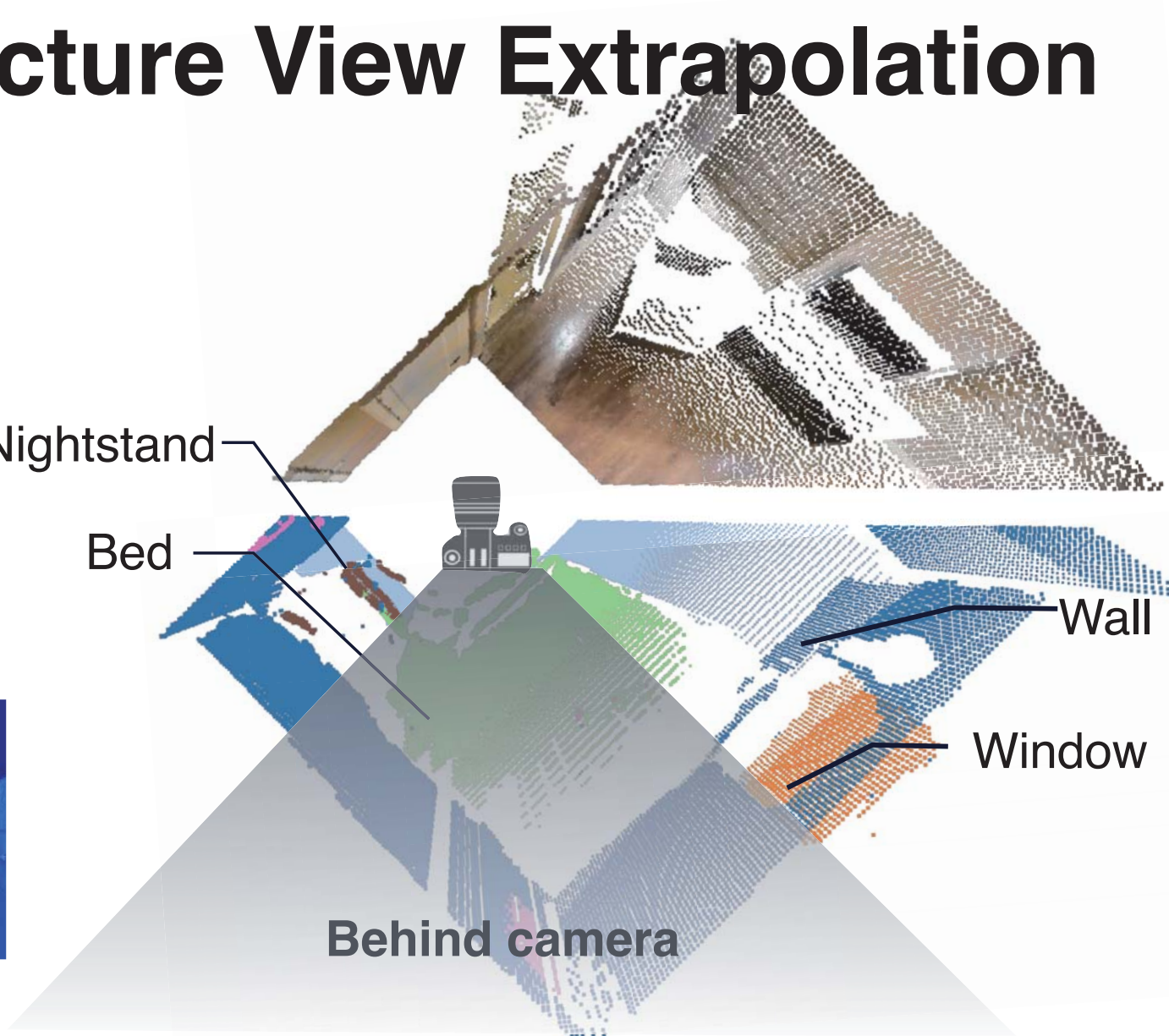
Nightstand

Bed

Wall

Window

Behind camera



# Key idea

**Key idea:** Indoor environments are often **highly structured**.

By learning over the statistics of many typical scenes, the model should be able to leverage **strong contextual cues** inside the image to predict what is beyond the FoV.



# Training data

## 3D House Datasets



### Synthetic Houses (SUNCG):

58,866 RGB-D panoramas

Pre-train



### Real-World Houses (Matterport3D):

5,315 RGB-D panoramas

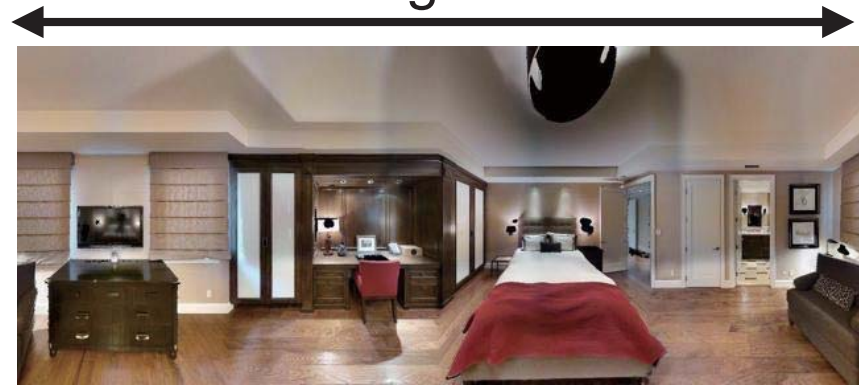
Fine-tune and test

# Data Representation

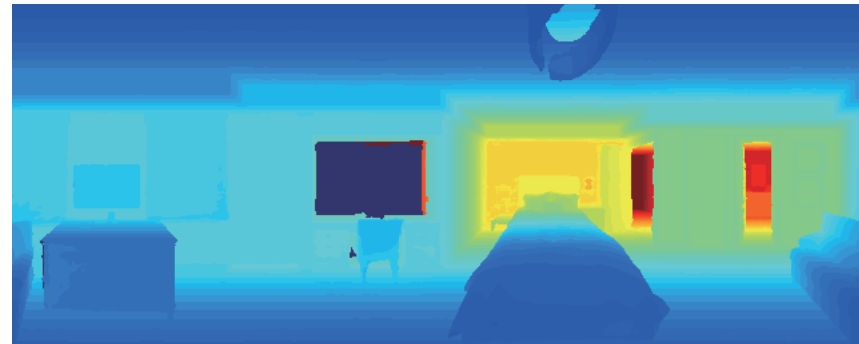


3D Room

360 Degree FoV



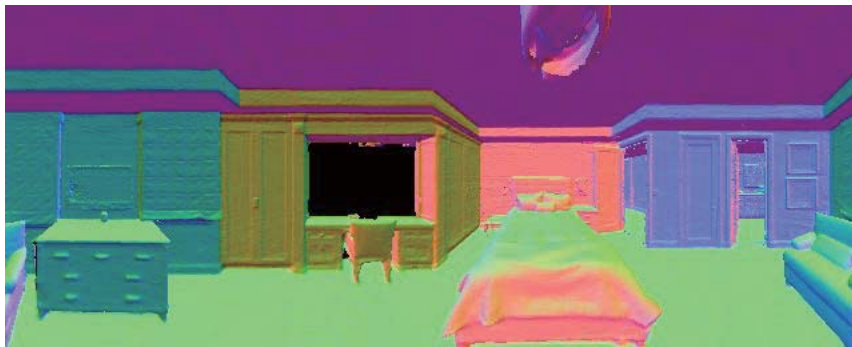
Color Panorama



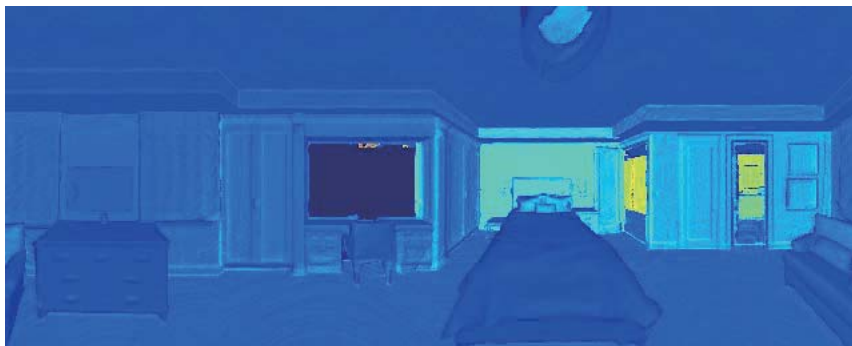
Depth Panorama



# Data Representation

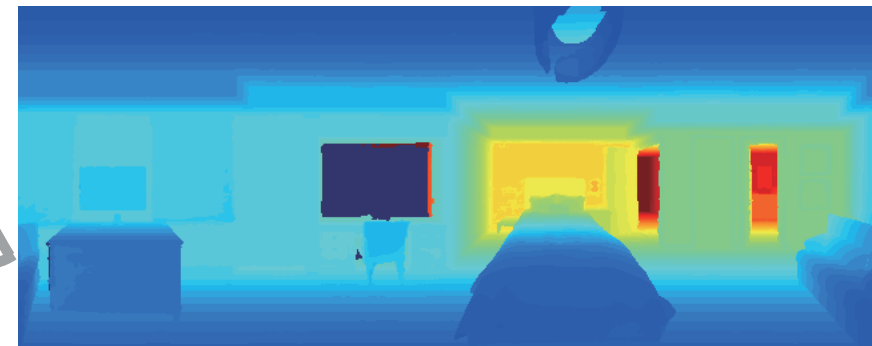
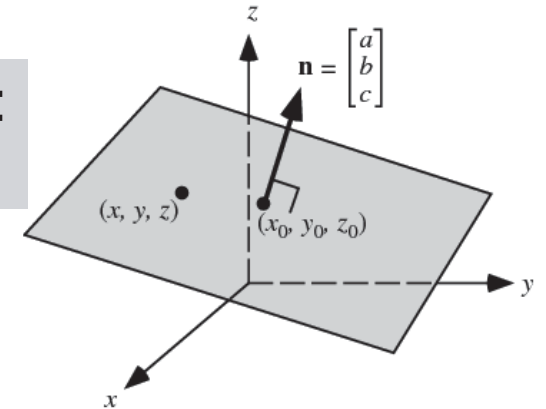


Surface Normal (a,b,c)



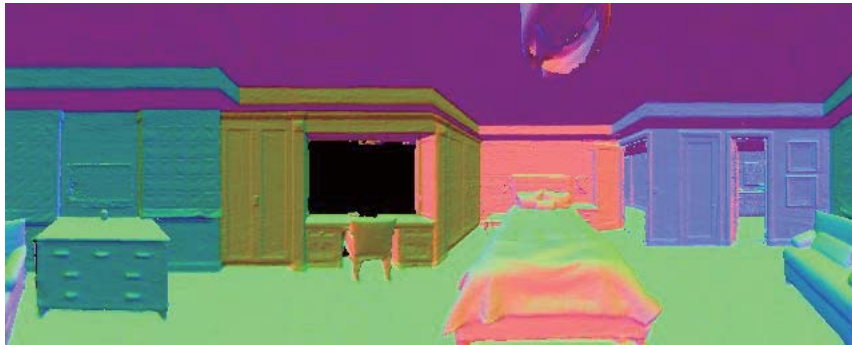
Plane Distance to Origin (p)

Plane Equation:  
 $ax+by+cz-p=0$

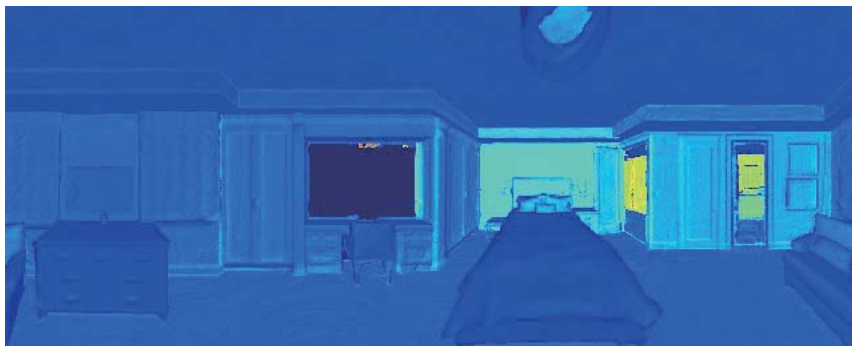


Depth Panorama

# Data Representation

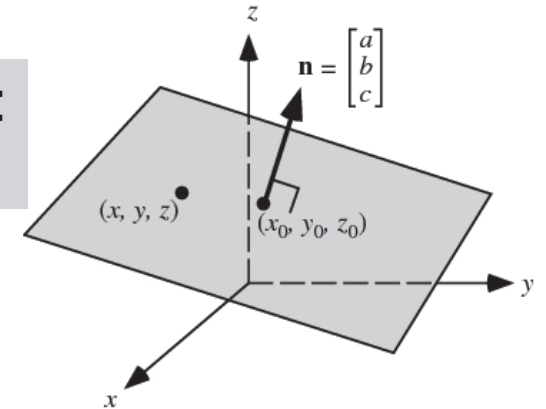


Surface Normal (a,b,c)



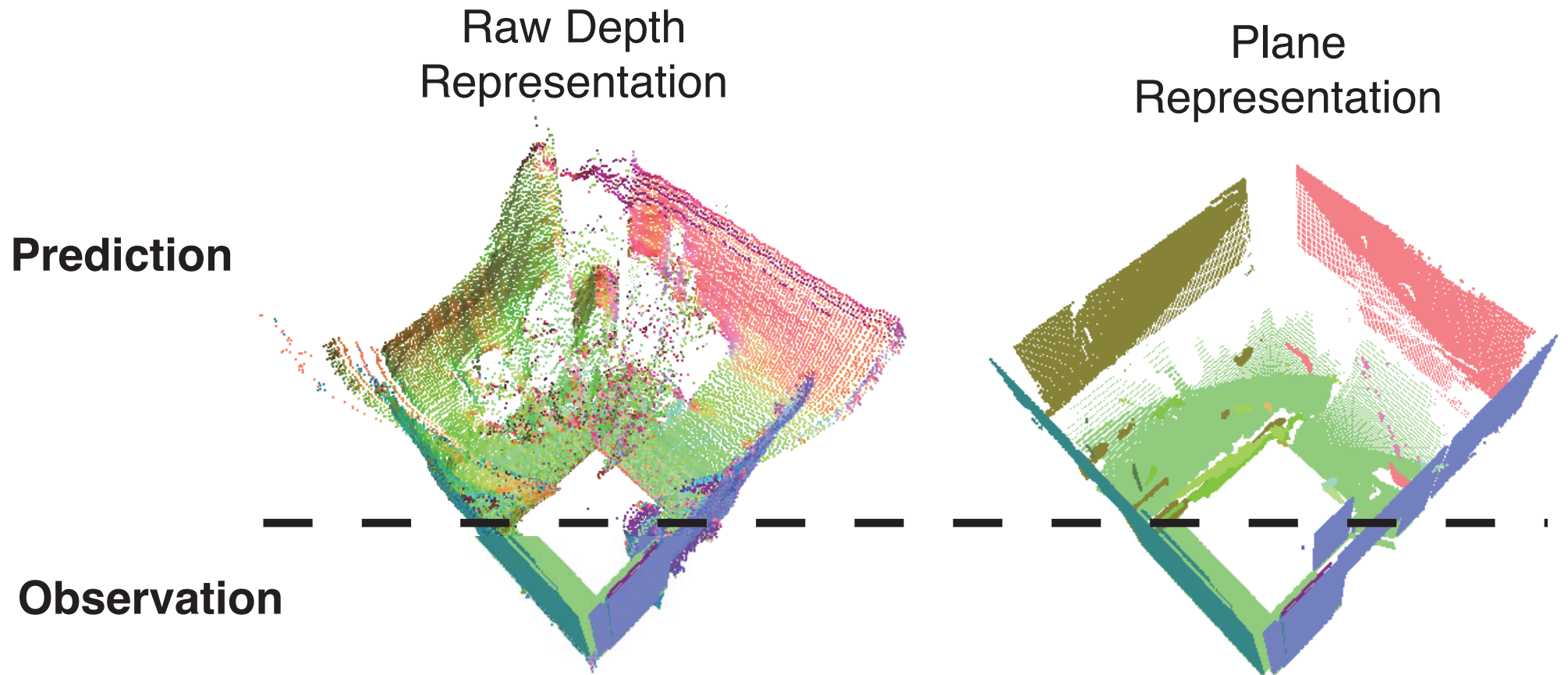
Plane Distance to Origin (p)

Plane Equation:  
 $ax+by+cz-p=0$

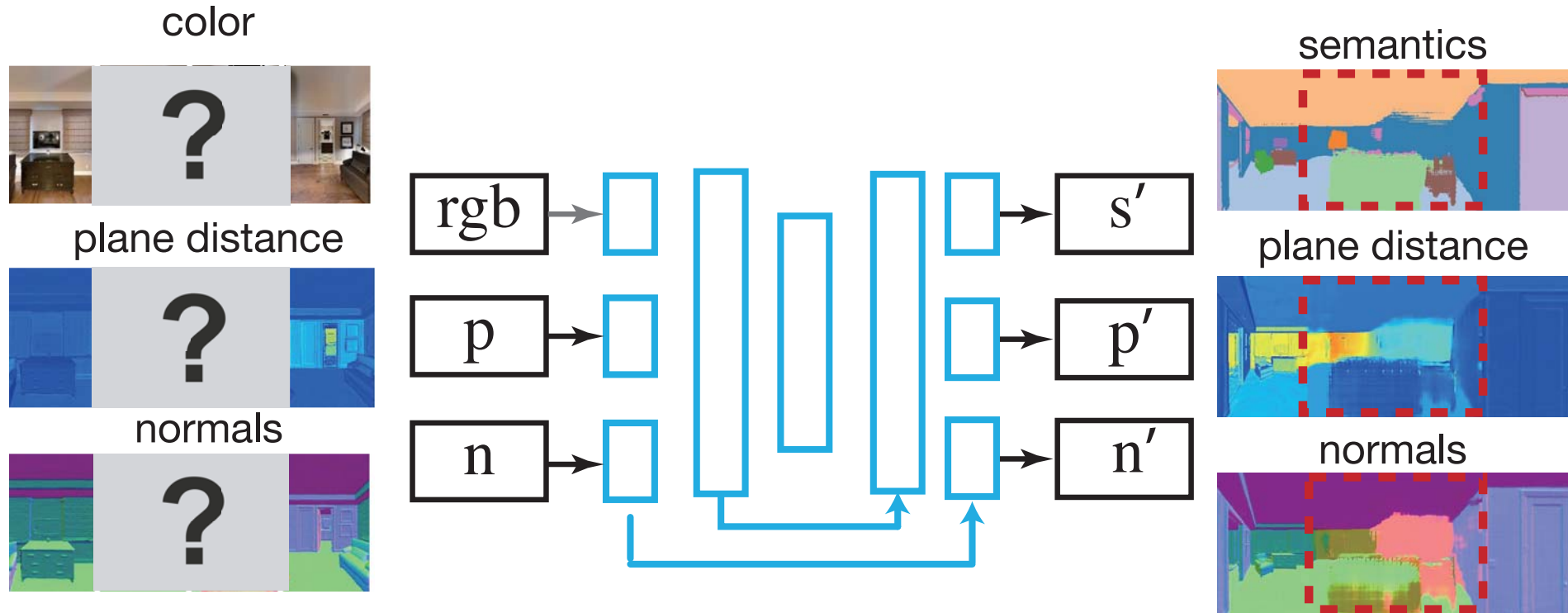


- ✓ Pixels on the same planar surface share the same plane equation.
- ✓ Representation is piecewise constant in a typical indoor environment.

# Data Representation



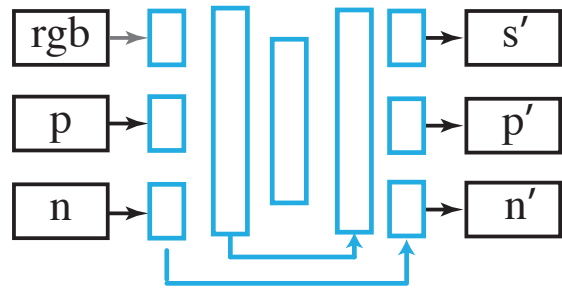
# Im2Pano3D Network



**What training objectives should we use?**



# Training Objectives



# Training Objectives

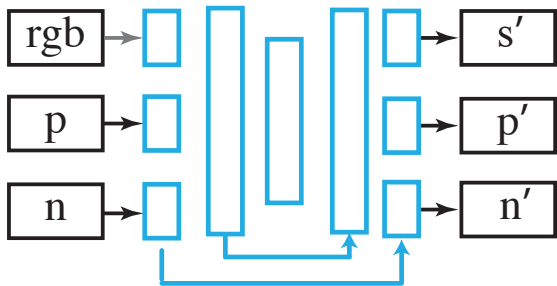
Every Pixel is  
Correct



softmax

L1

cosine



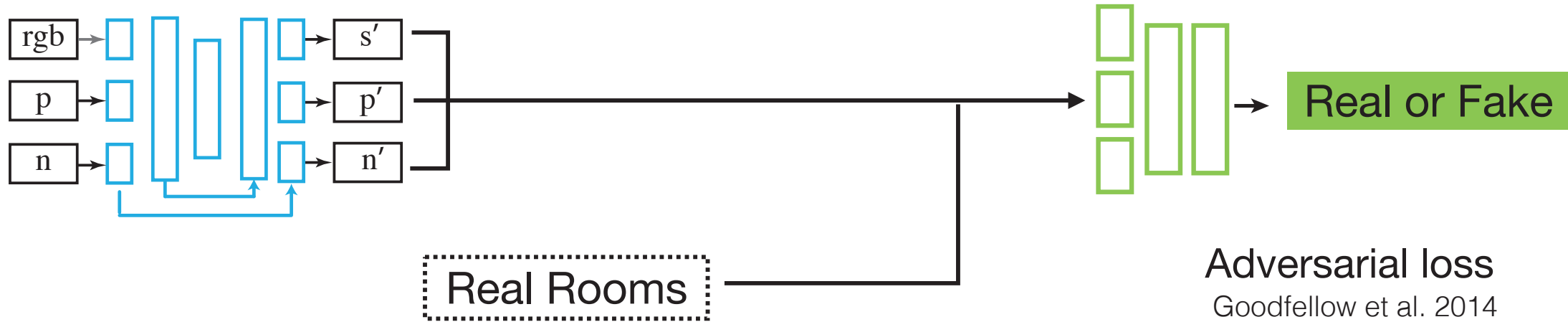
Prediction



Ground truth

# Training Objectives

Prediction is  
Plausible



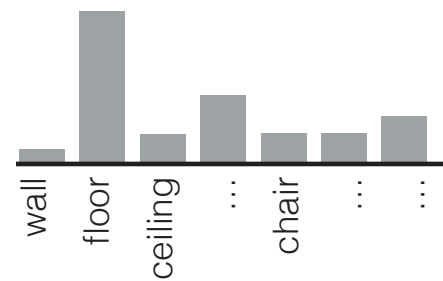
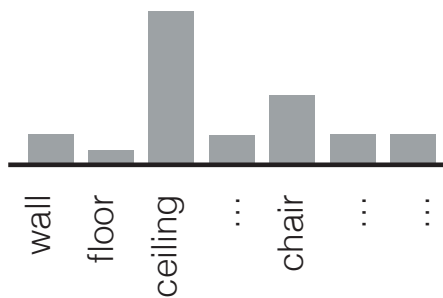
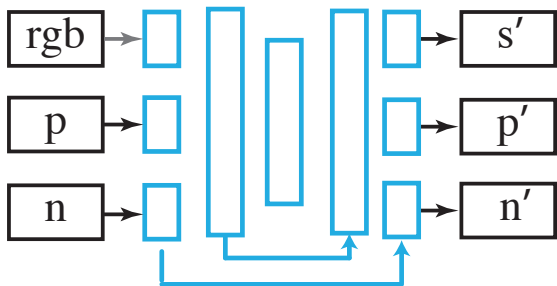
# Training Objectives

Similar Scene  
Attribute



scene category

object distribution



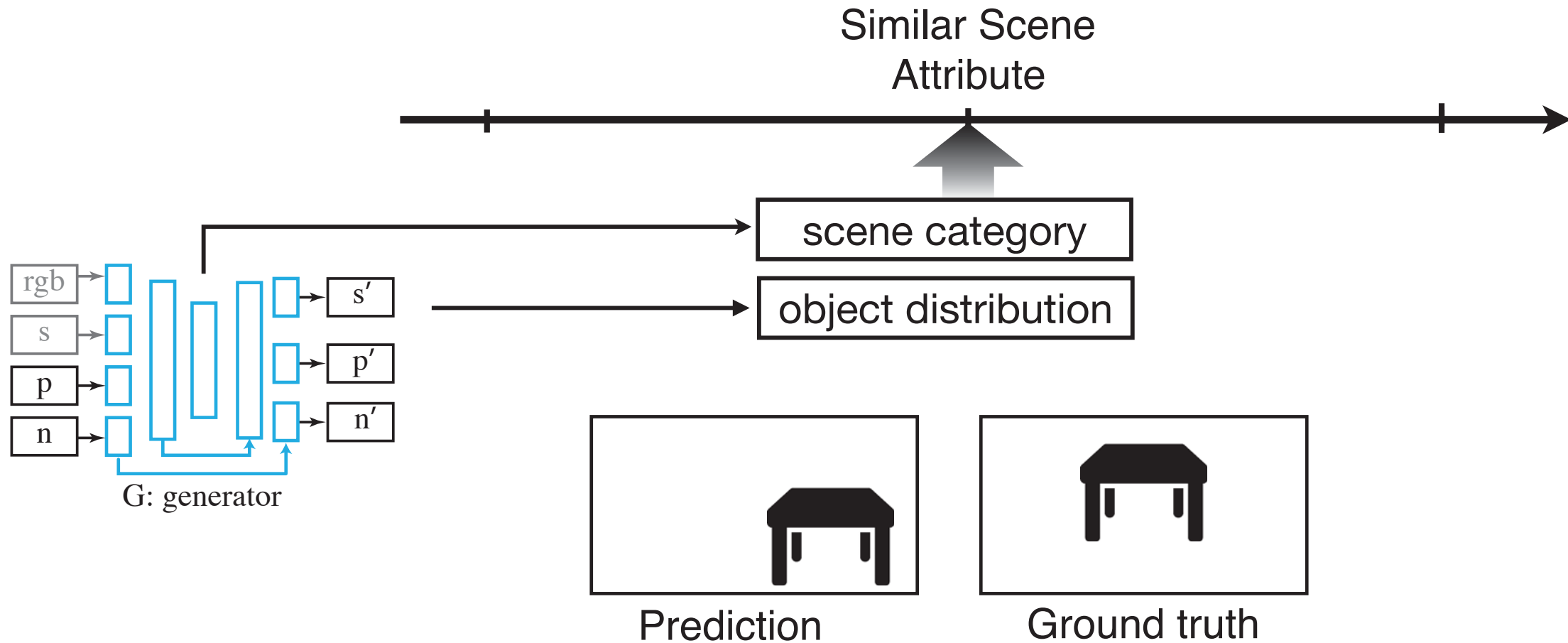
Prediction

$$L_{dis} = \sum_c |y_c - h(x_c)|$$

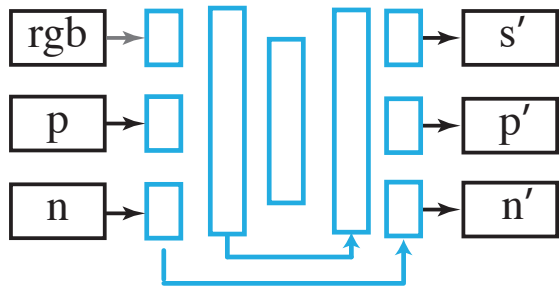
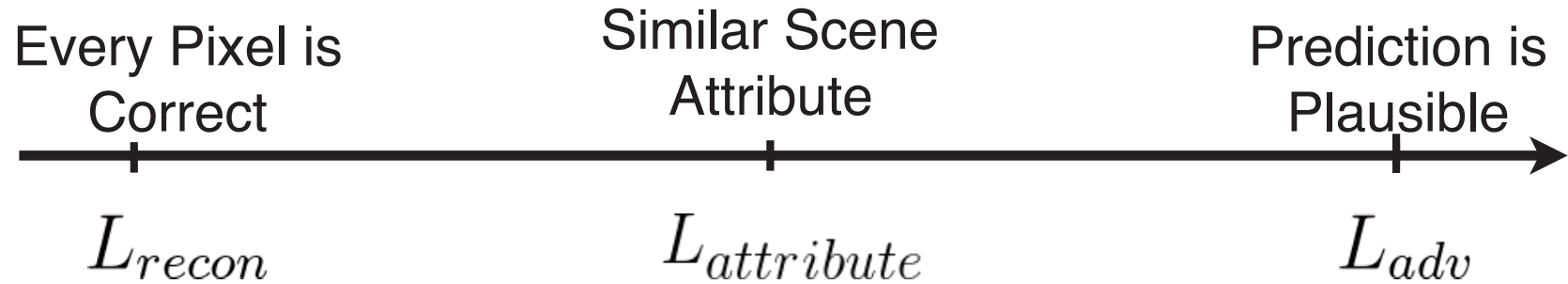
Ground truth



# Training Objectives



# Training Objectives



$$L = \lambda_1 L_{recon} + \lambda_2 L_{attribute} + \lambda_3 L_{adv}$$

# Results

# Results

Input Observation

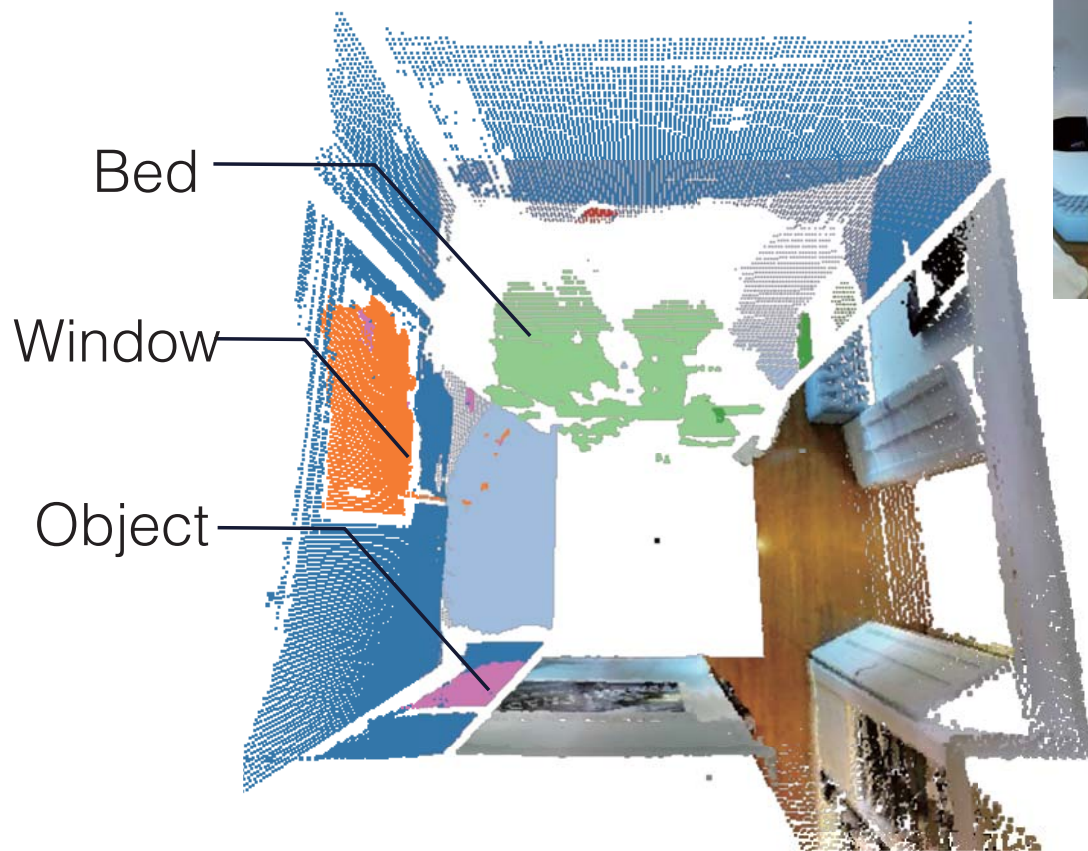


● ceiling ● wall ● floor ● window ● bed ● door ● cabinet ● chair ● sofa ● tv ● table ● object ● furniture



# Results

Prediction



Ground truth

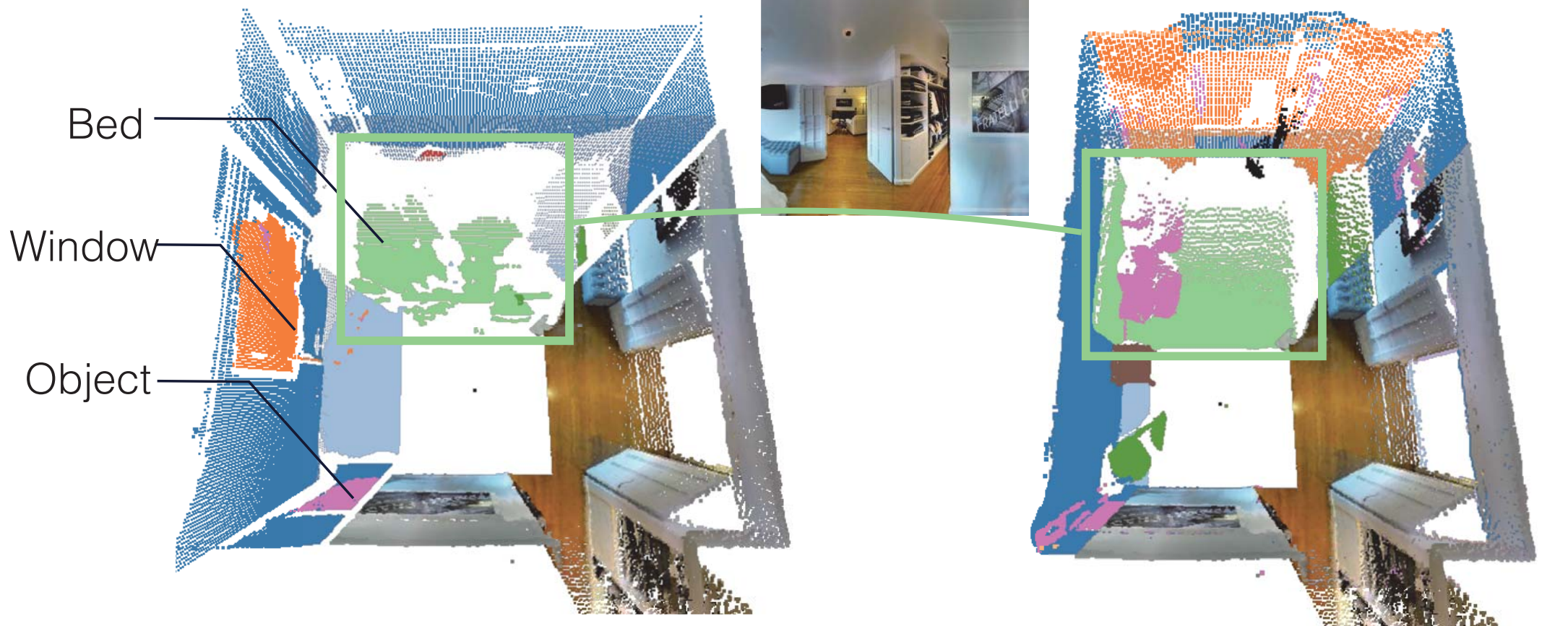


- ceiling ● wall ● floor ● window ● bed ● door ● cabinet ● chair ● sofa ● tv ● table ● object ● furniture

# Results

Prediction

Ground truth

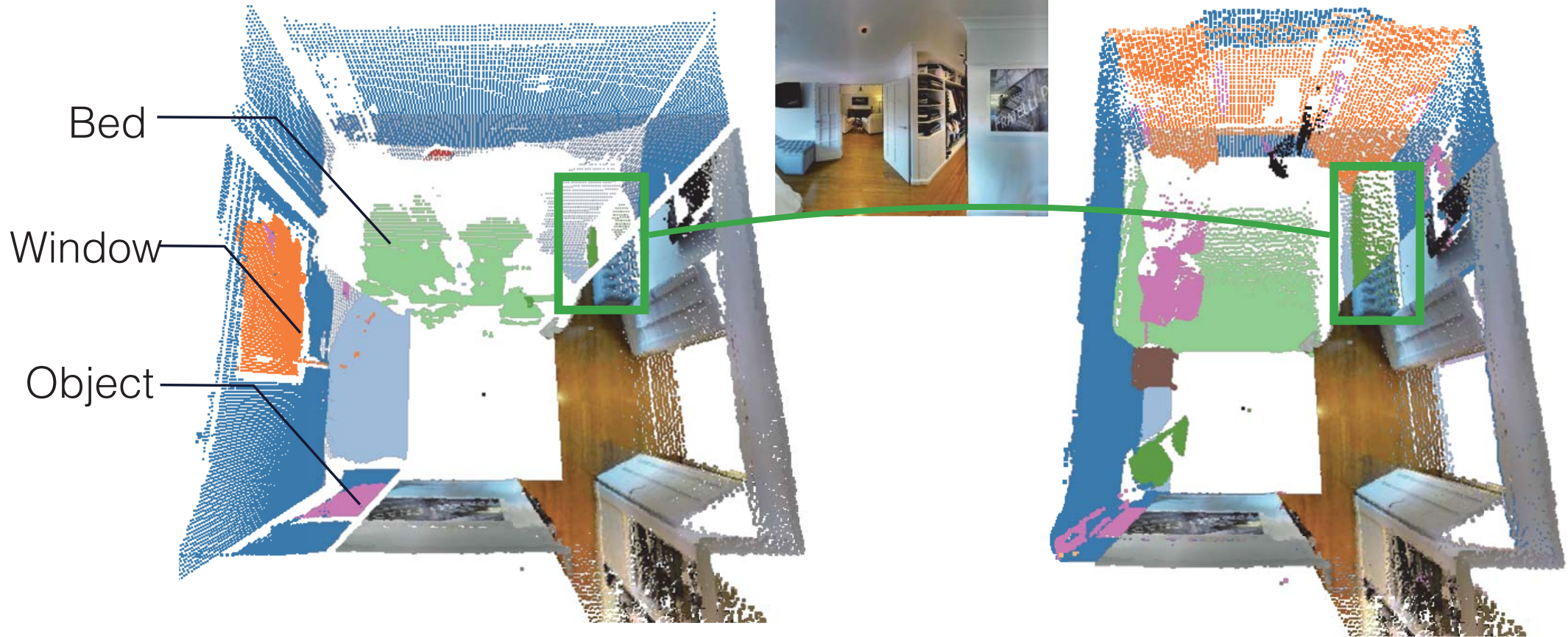




# Results

Prediction

Ground truth



# Results

Prediction

Ground truth

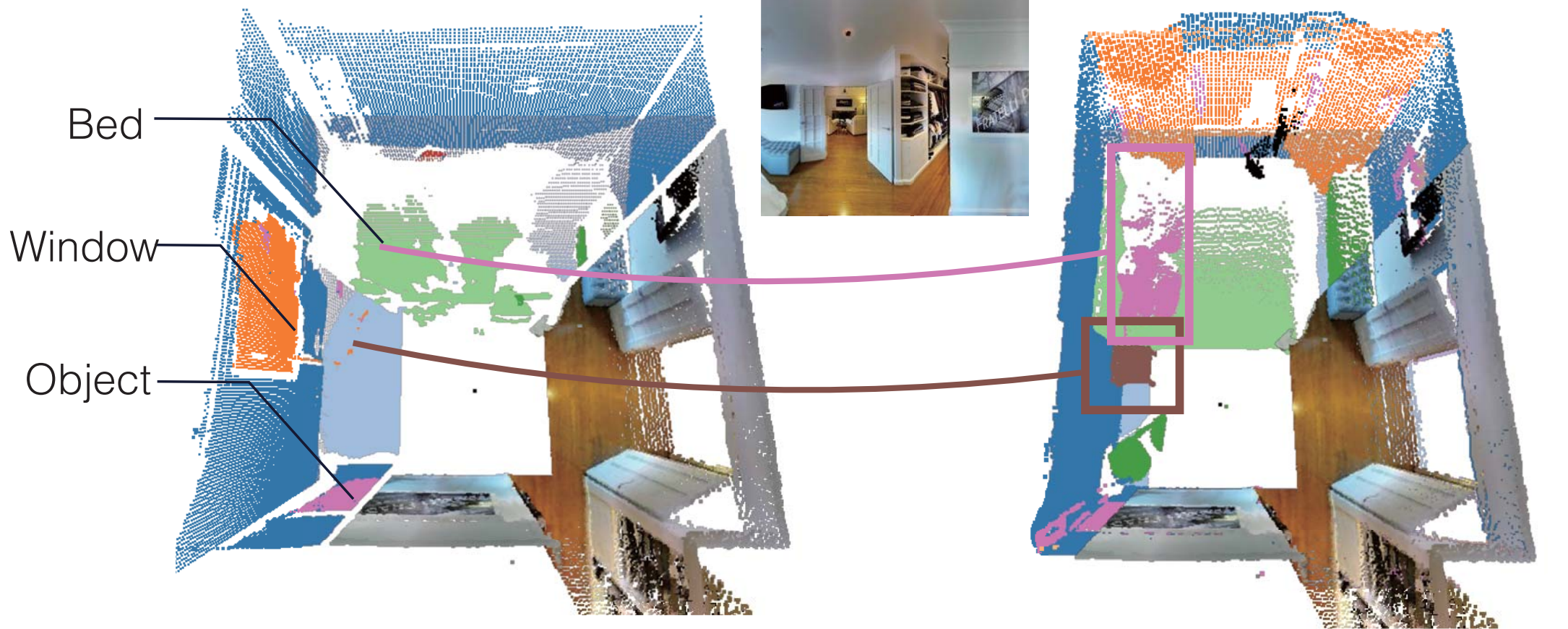




# Results

Prediction

Ground truth

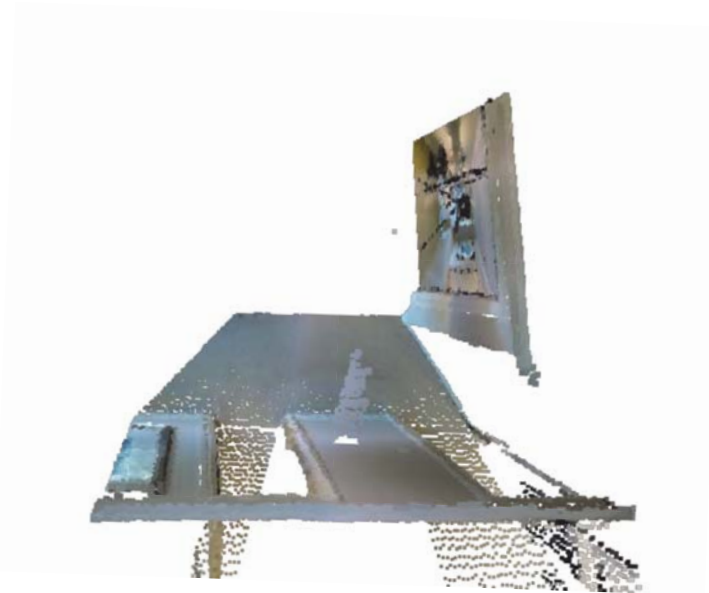


- ceiling
- wall
- floor
- window
- bed
- door
- cabinet
- tv
- table
- object
- furniture



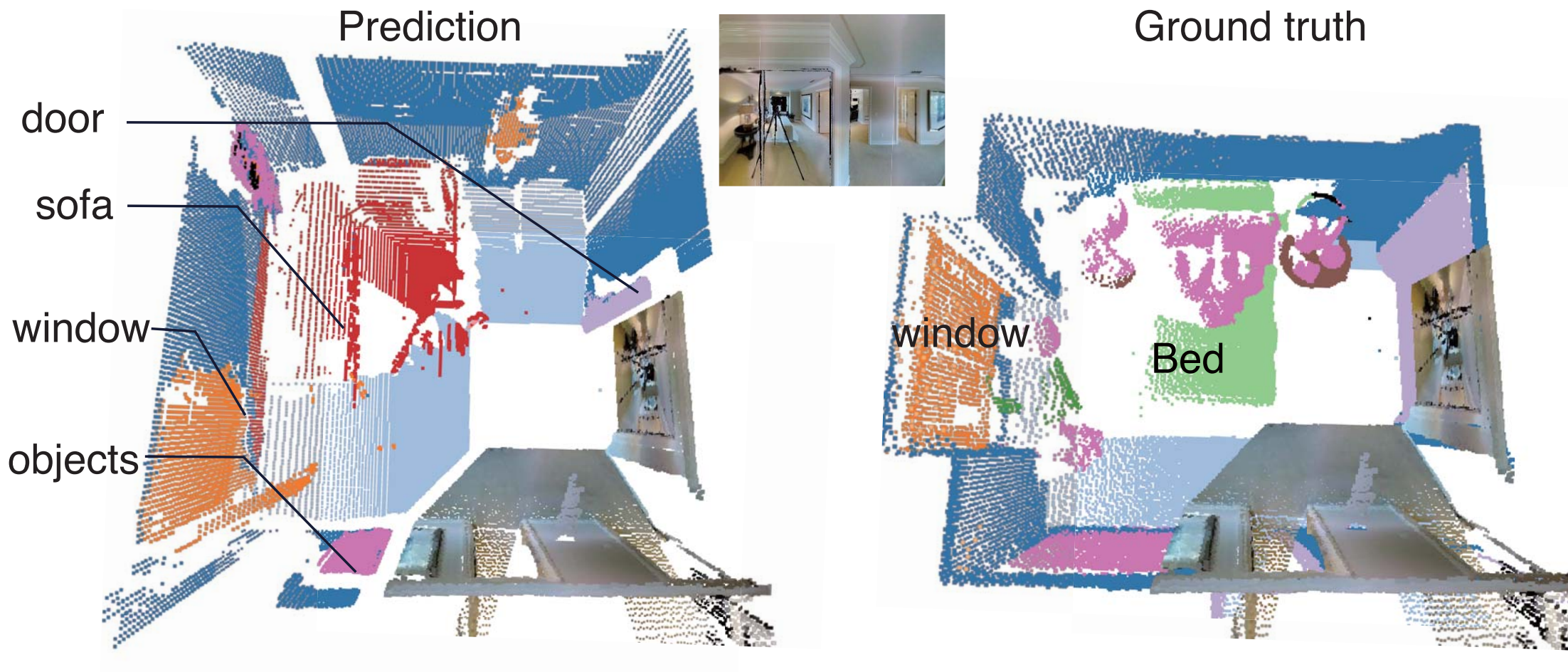
# Results

Input Observation



● ceiling ● wall ● floor ● window ● bed ● door ● cabinet ● chair ● sofa ● tv ● table ● object ● furniture

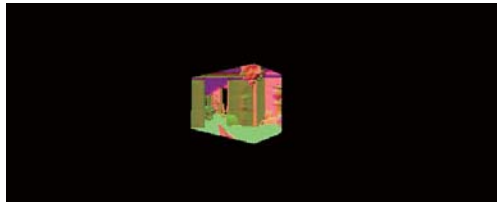
# Results



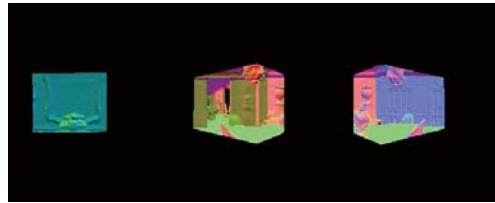
● ceiling ● wall ● floor ● window ● bed ● door ● cabinet ● chair ● sofa ● tv ● table ● object ● furniture

# Camera Configurations in real platforms

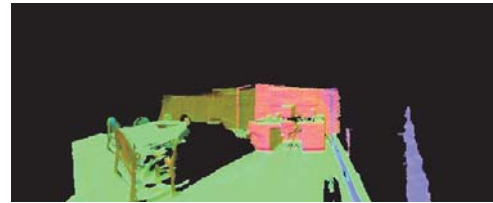
One RGB-D



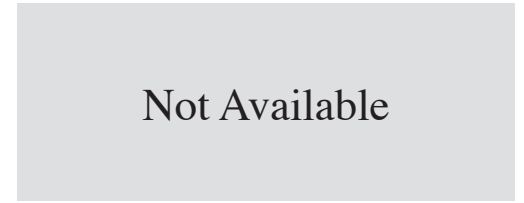
Three RGB-D



One RGB-D+motion



RGB pano



Input

Device



# Camera Configurations

One RGB-D

Three RGB-D

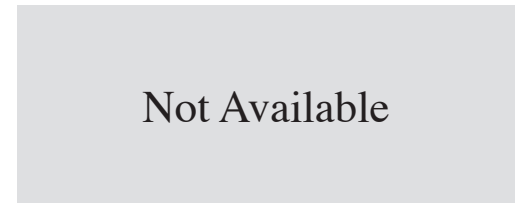
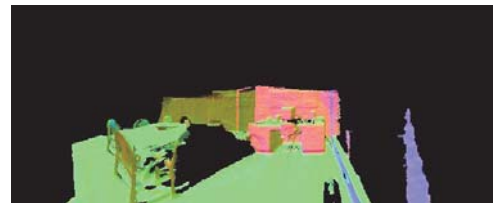
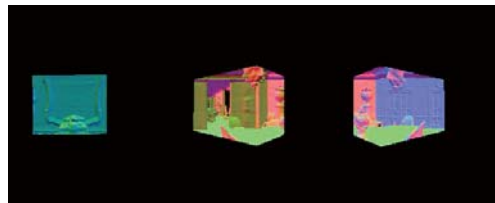
One RGB-D+motion

RGB pano

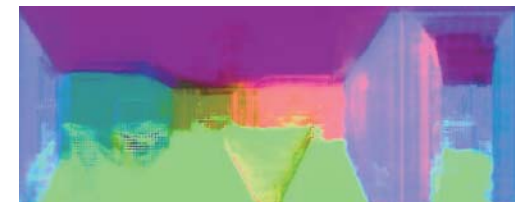
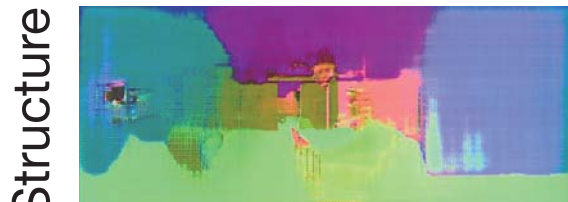
Input



Semantics



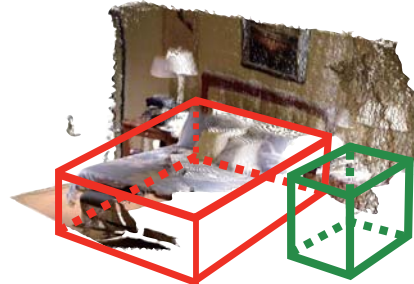
Structure



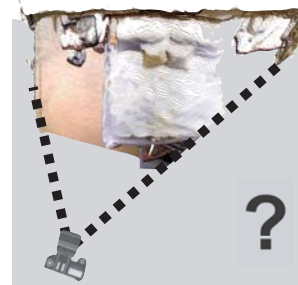
● ceiling ● wall ● floor ● window ● bed ● door ● cabinet ● chair ● sofa ● tv ● table ● object ● furniture



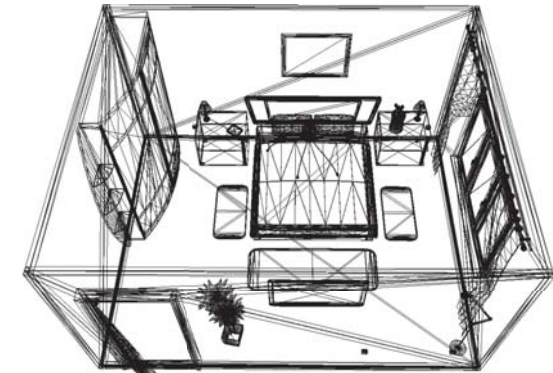
# Advances Towards 3D Scene Understanding



**Amodal 3D**  
[Song and Xiao  
ECCV'14, CVPR'16]



**Beyond FoV**  
[song et al. CVPR'18]



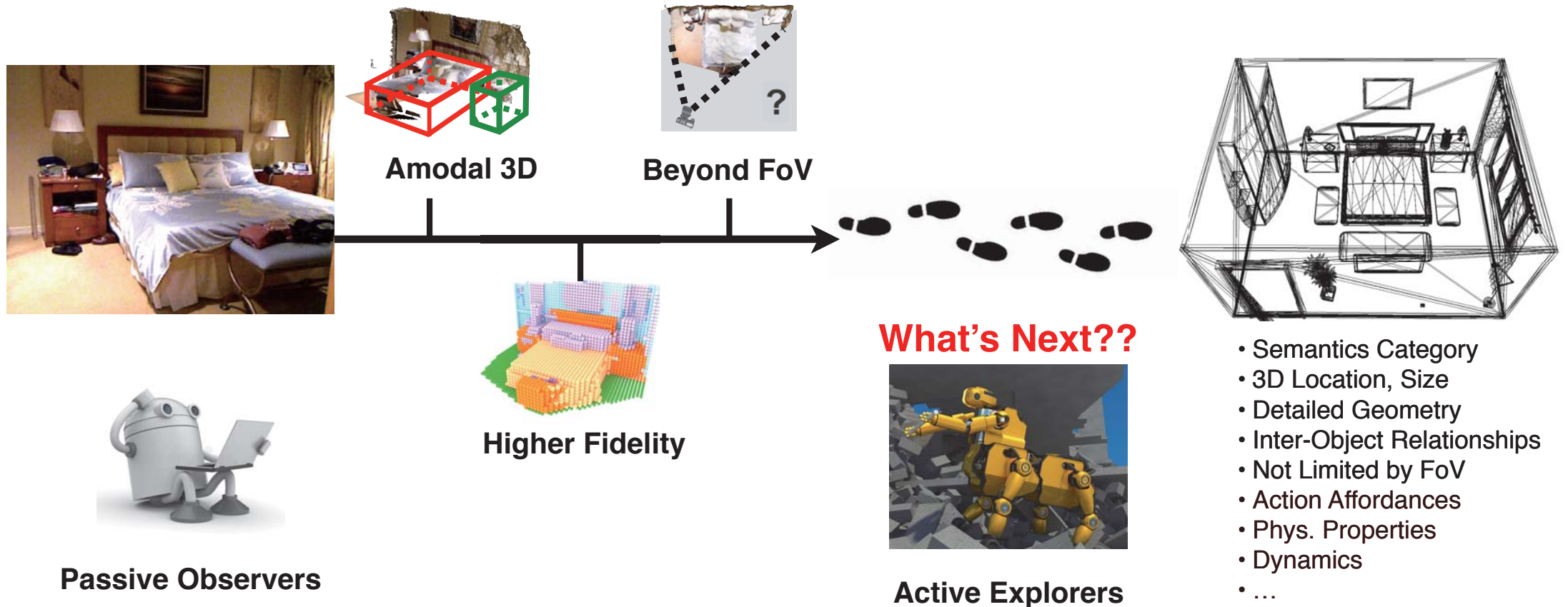
- Semantics Category
- 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. Properties
- Dynamics
- ...



**Higher Fidelity**  
[Song et al. CVPR'17]



# Advances Towards 3D Scene Understanding



# Richer Representation through Interaction

## Active Exploration



Partial Observation

Inference  
(Im2Pano3D)



3D Scene Prior

Guide



Efficient exploration  
+ Most useful observation

Improve



# Richer Representation through Interaction

## Active Exploration



Partial Observation



3D Scene Prior



Efficient exploration  
+ Most useful observation

## Active physical Interaction



**Actions:** Poking, Grasping

**Physical properties:**

Surface material

Friction coefficient

# Richer Representation through Interaction

## Active Exploration



Partial Observation

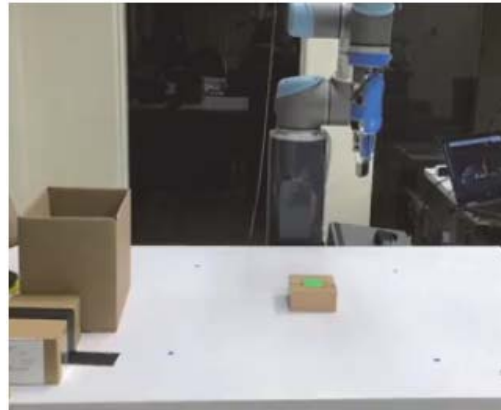


3D Scene Prior



Efficient Exploration  
+ Most useful observation

## Active physical Interaction



**Actions:** Pushing, Grasping

**Physical properties:**  
Surface material  
Friction coefficient

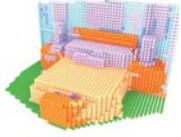
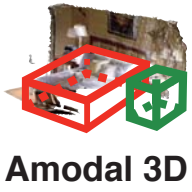


**Actions:** Tossing

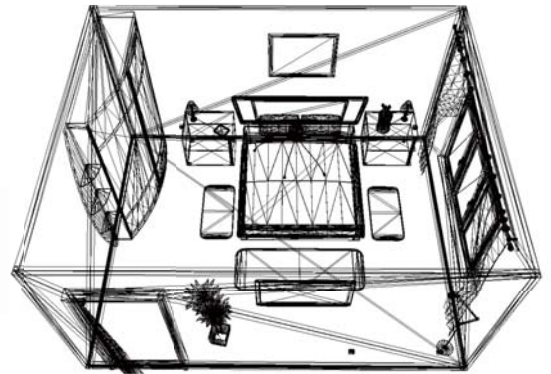
**Physical properties:**  
Mass distribution,  
Aerodynamic



# Comprehensive 3D Scene Understanding



**What's Next??**



- Semantics Category
- 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. Properties
- Dynamics
- ...

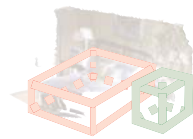


**Passive Observers**



**Active Explorer**

# Comprehensive 3D Scene Understanding



Amodal 3D



Beyond FoV

**We are looking for PhDs and Post-docs!**



**COLUMBIA UNIVERSITY**  
IN THE CITY OF NEW YORK

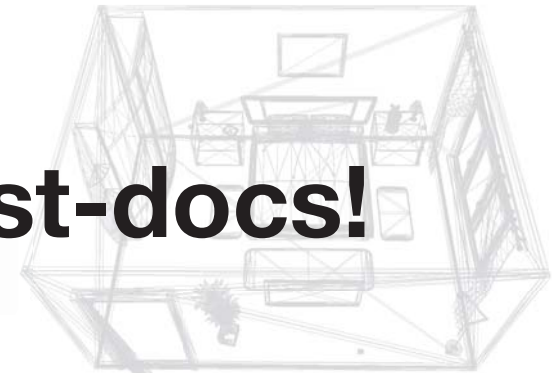
Higher Fidelity

Illumination



Passive Observers

Active Explorer



- Semantics Category
- 3D Location, Size
- Detailed Geometry
- Inter-Object Relationships
- Not Limited by FoV
- Action Affordances
- Phys. Properties
- Dynamics
- ...

# Acknowledgements

## Collaborators

Ferran Alet	Pat Hanrahan	Isabella Morona	Ian Taylor
Maria Bauza	Francois R. Hogan	Prem Qu Nair	Zhirong Wu
Angel Chang	Rachel Holladay	Matthias Nießner	Jianxiong Xiao
Nikhil Chavan Dafle	Qixing Huang	Alberto Rodriguez	Li Yi
Elliott Donlon	Hailin Jin	Eudald Romo	Kuan-Ting Yu
Nima Fazeli	Joon-Young Lee	Silvio Savarese	Fisher Yu
Matthew Fisher	Zimo Li	Manolis Savva	Ersin Yumer
Thomas Funkhouser	Melody Liu	Ari Seff	Andy Zeng
Druck Green	Weber Liu	Hao Su	Linguang Zhang
Leonidas Guibas	Daolin Ma	Orion Taylor	Yinda Zhang

## Funding:

NSF, Google, Intel, Facebook

# Thank You!