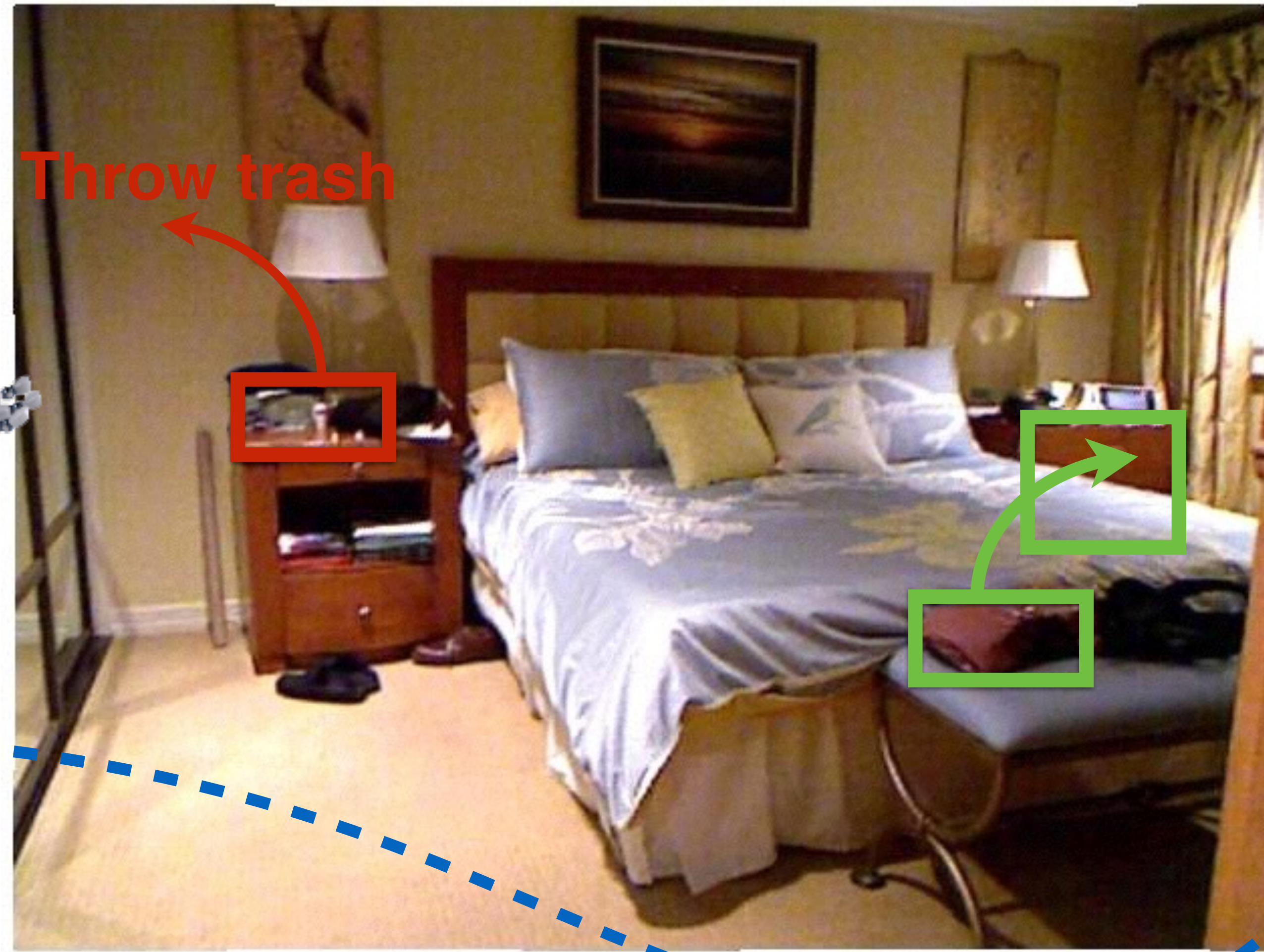
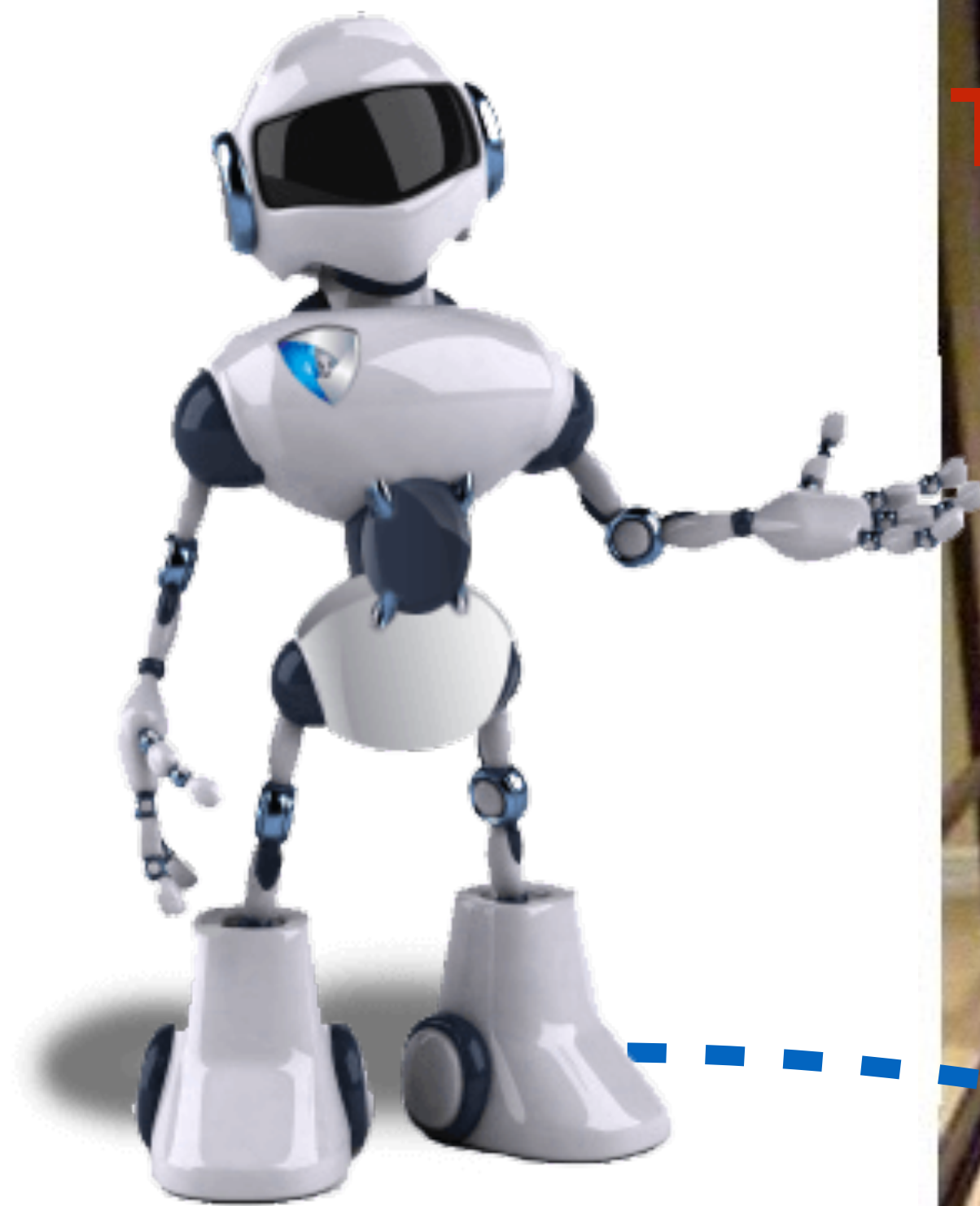


Semantic Scene Completion

from a Single Depth Image

Shuran Song, Fisher Yu, Andy Zeng, Angel X. Chang, Manolis Savva and Thomas Funkhouser
Princeton University

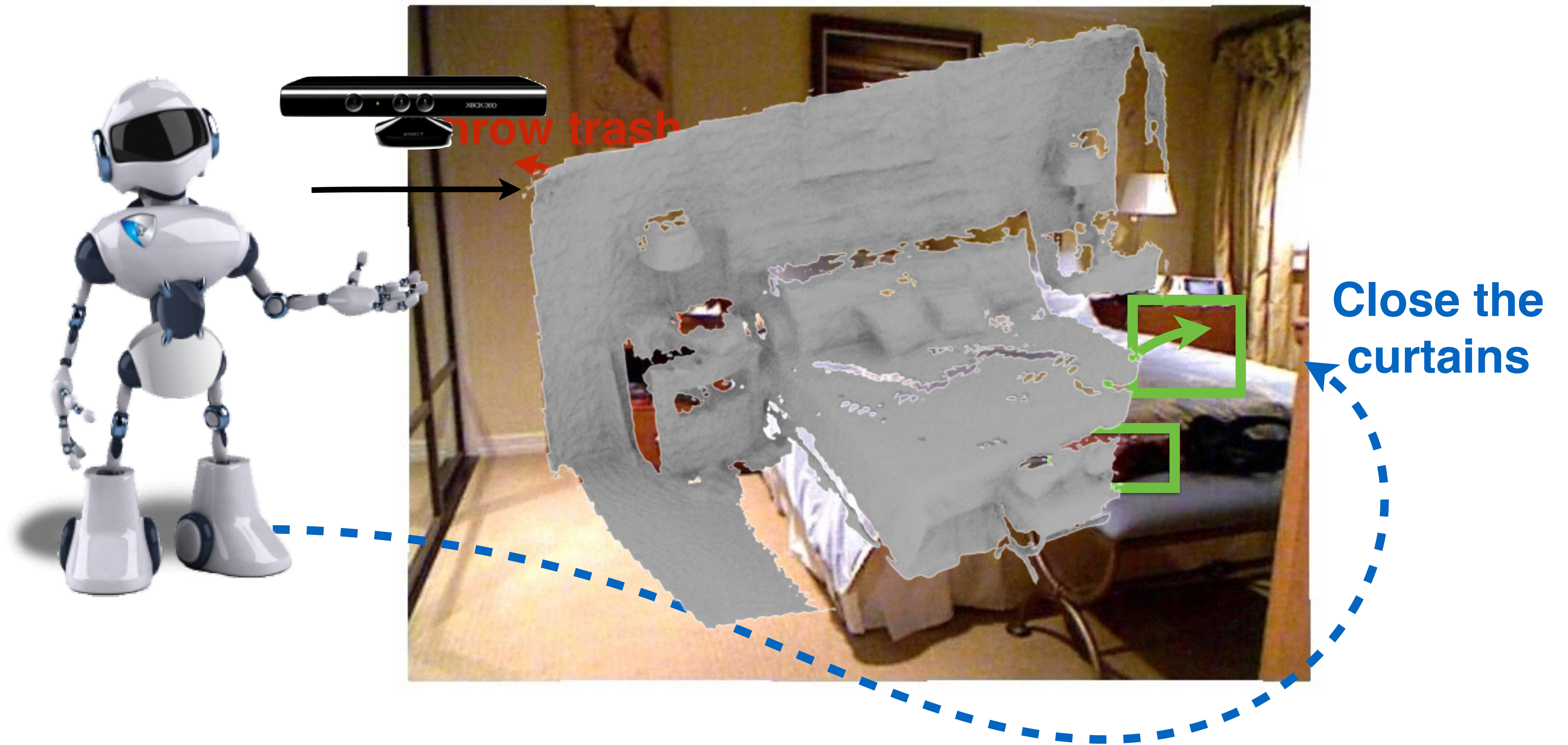
Motivation



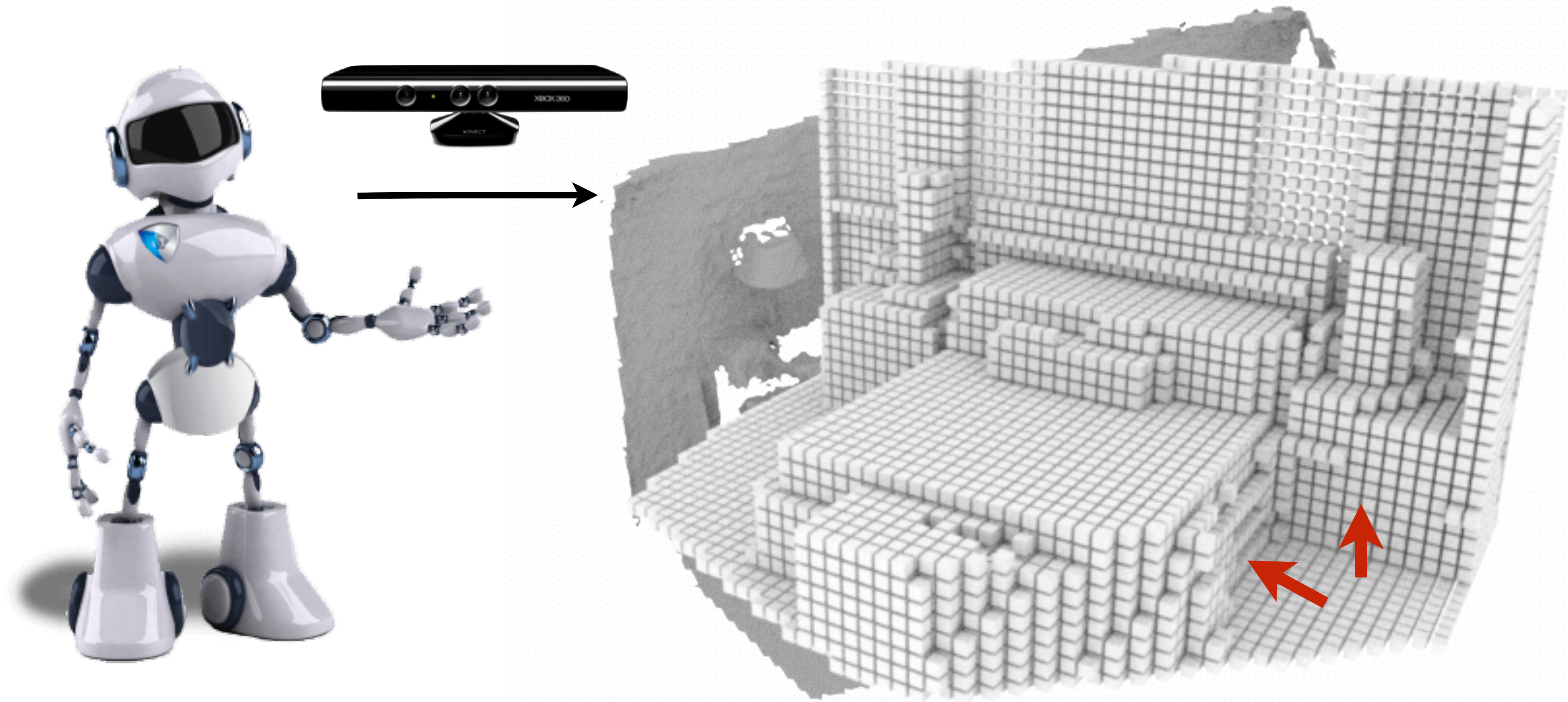
Throw trash

Close the curtains

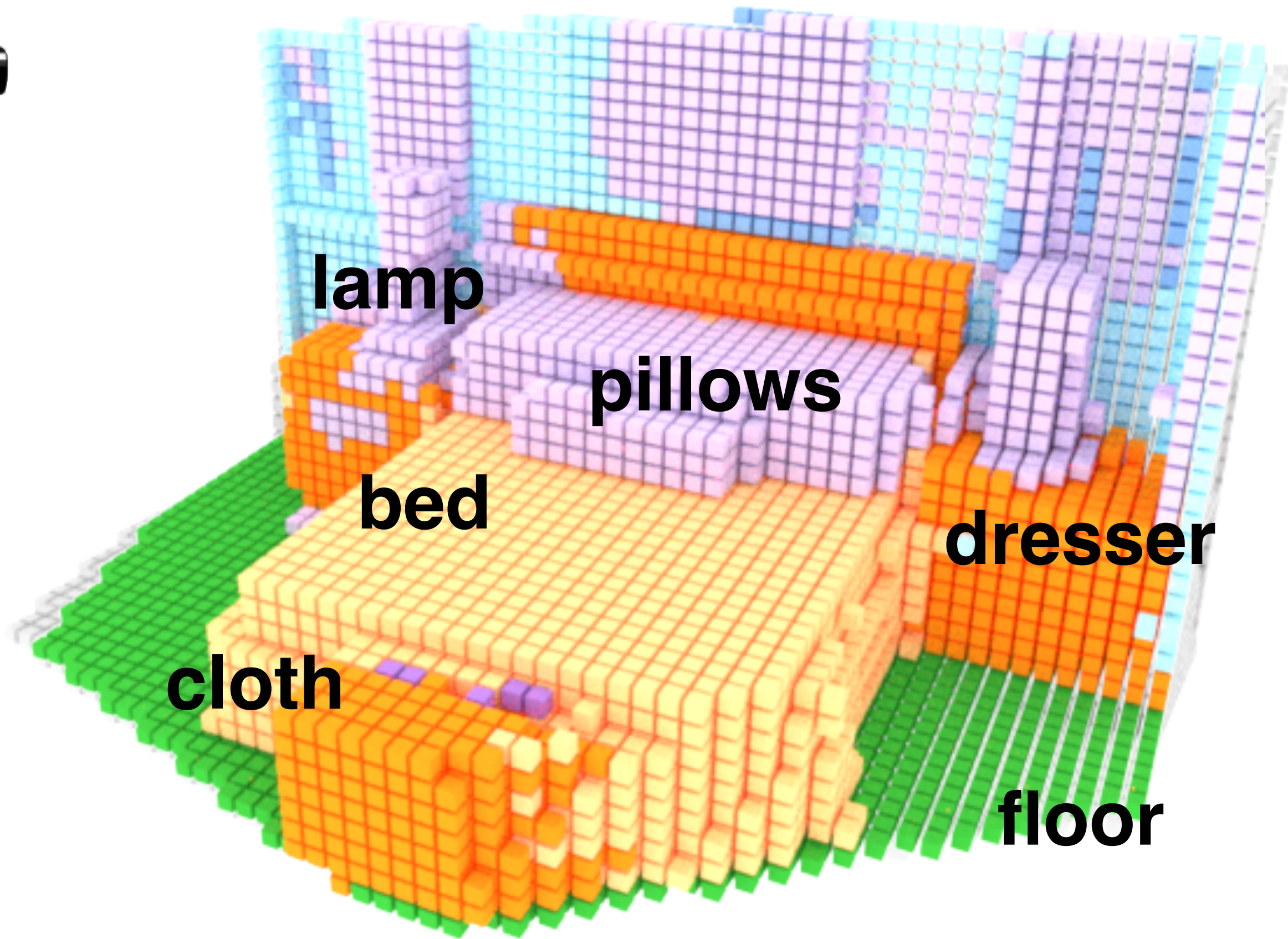
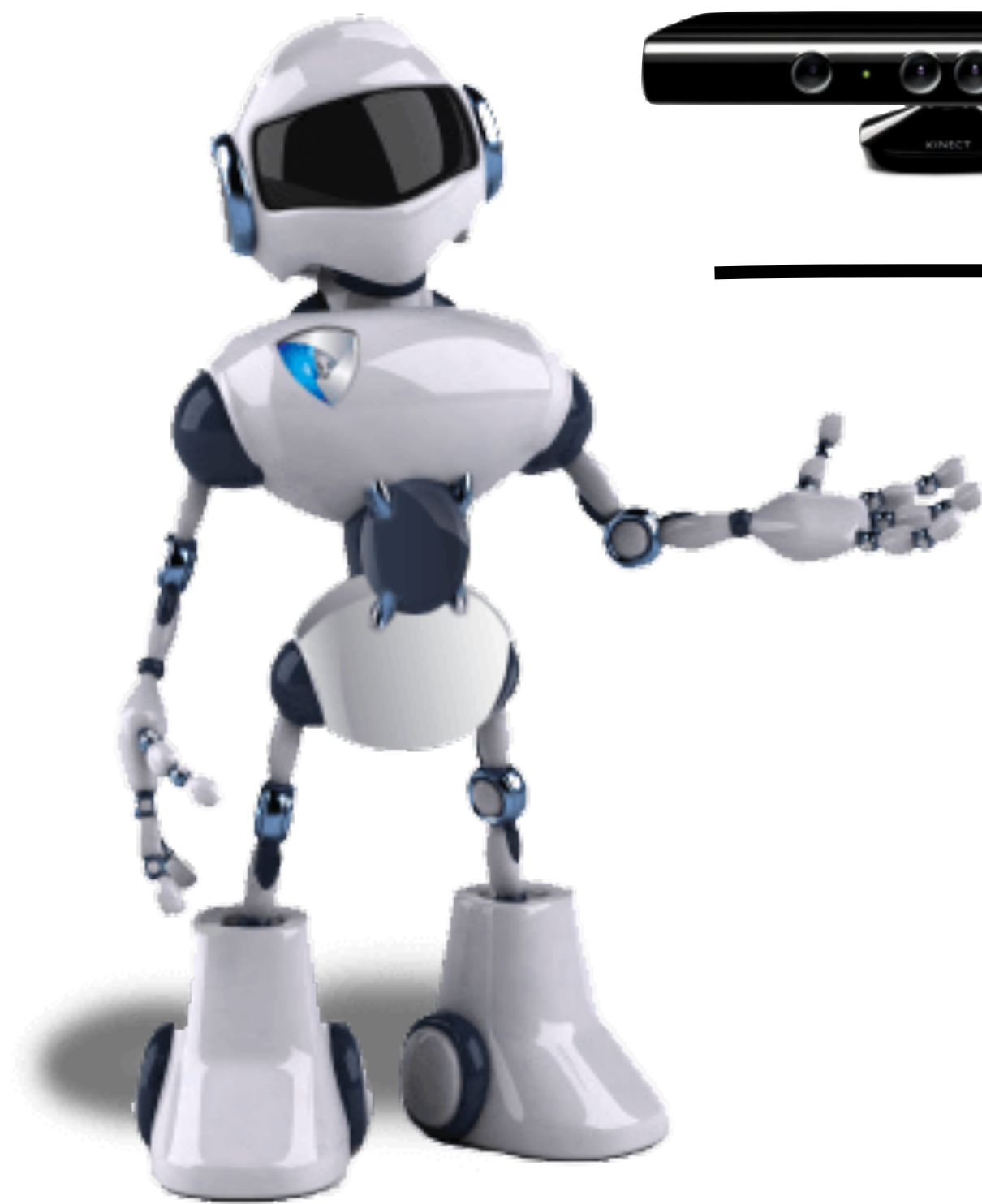
Partial observation



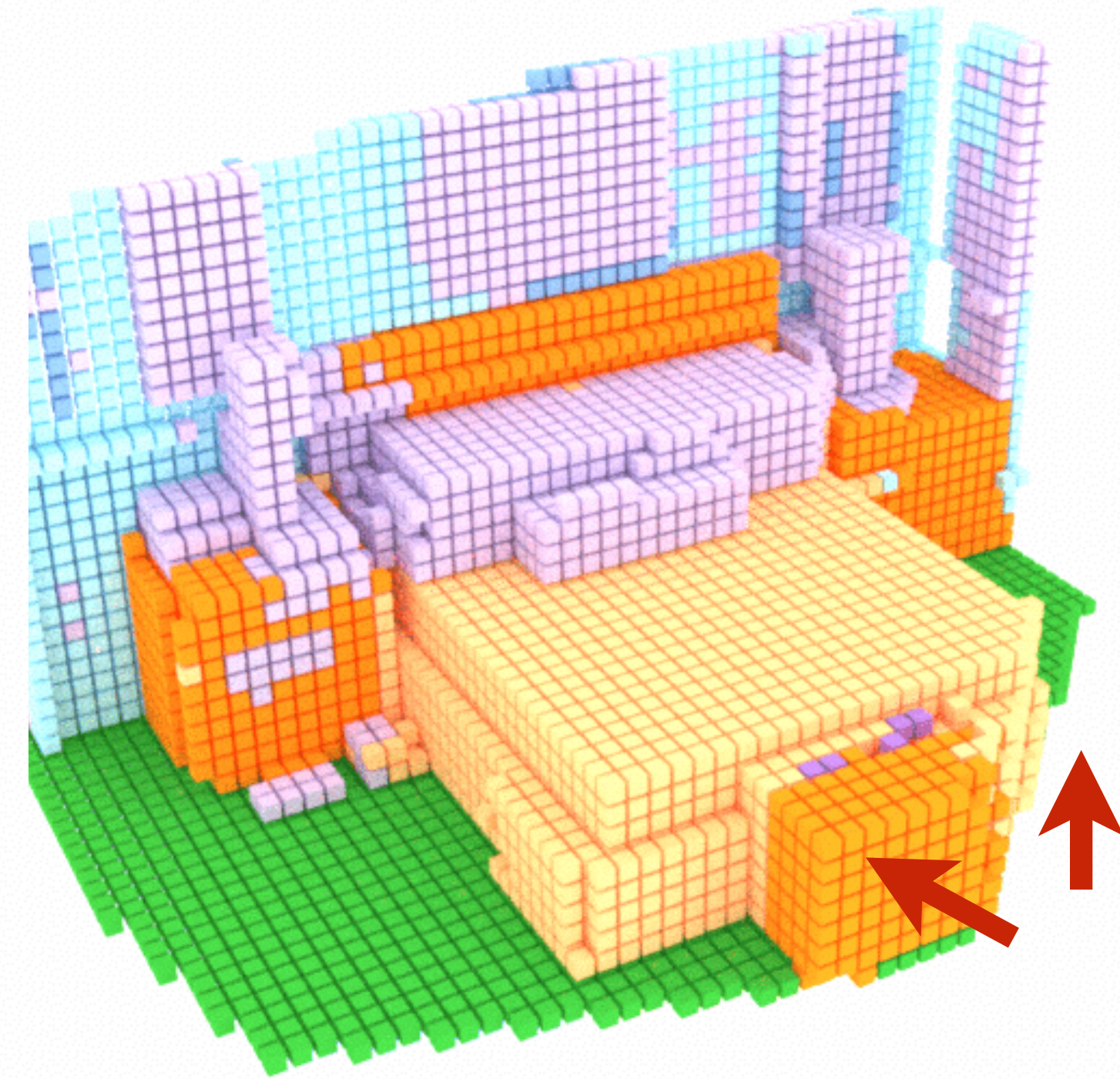
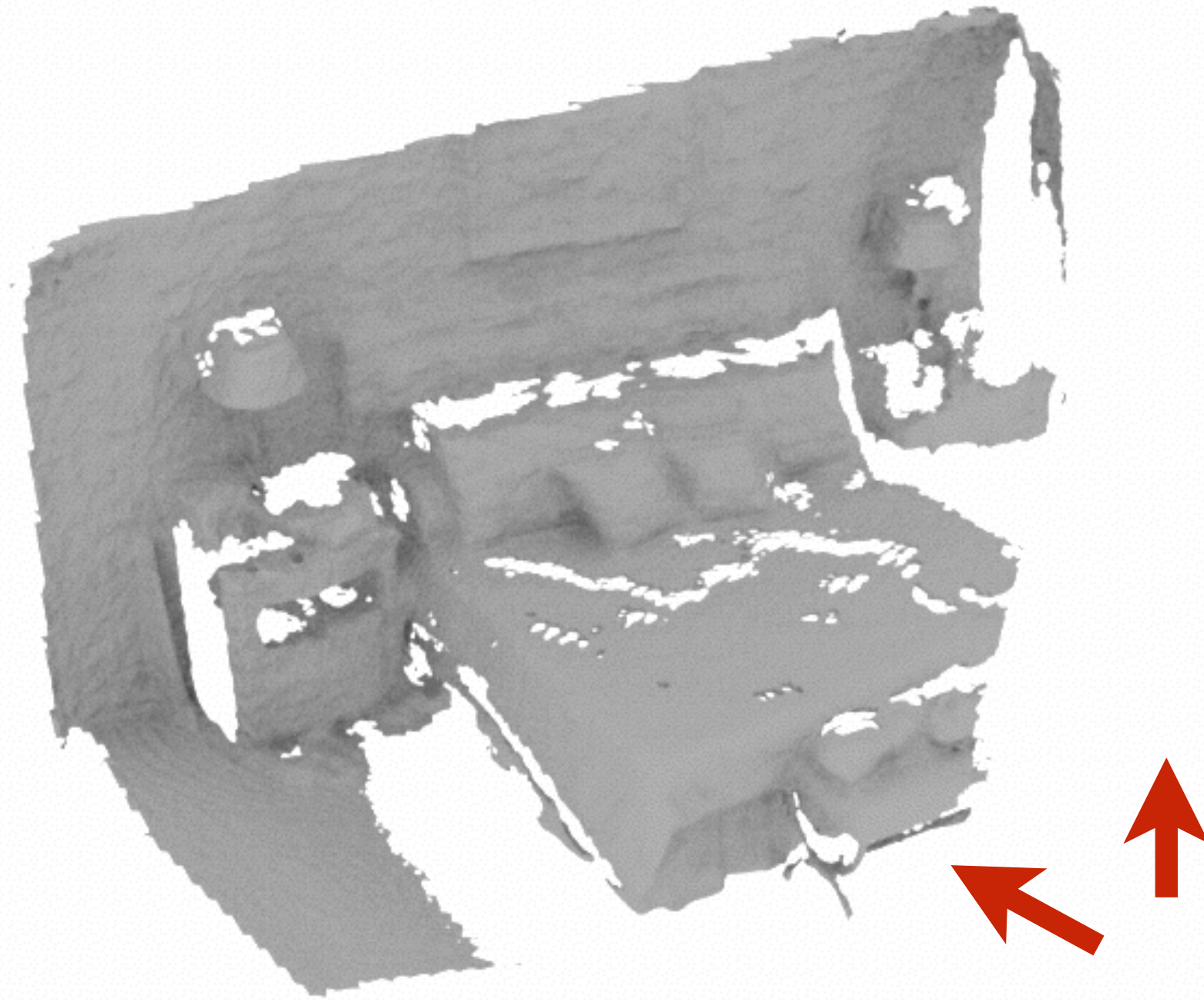
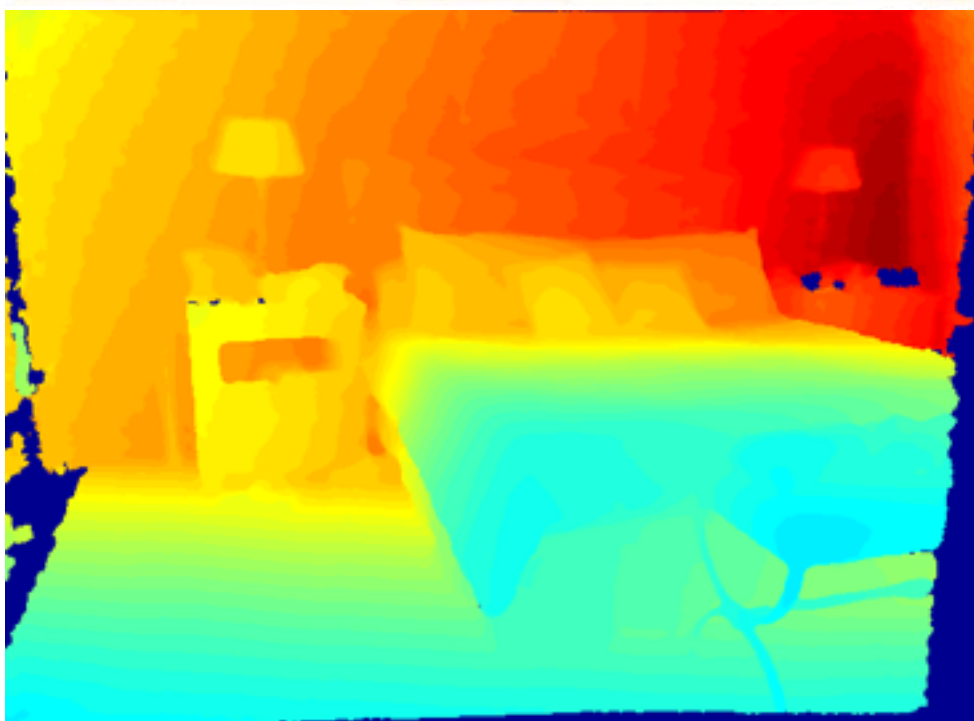
Complete 3D structure



Semantic meaning



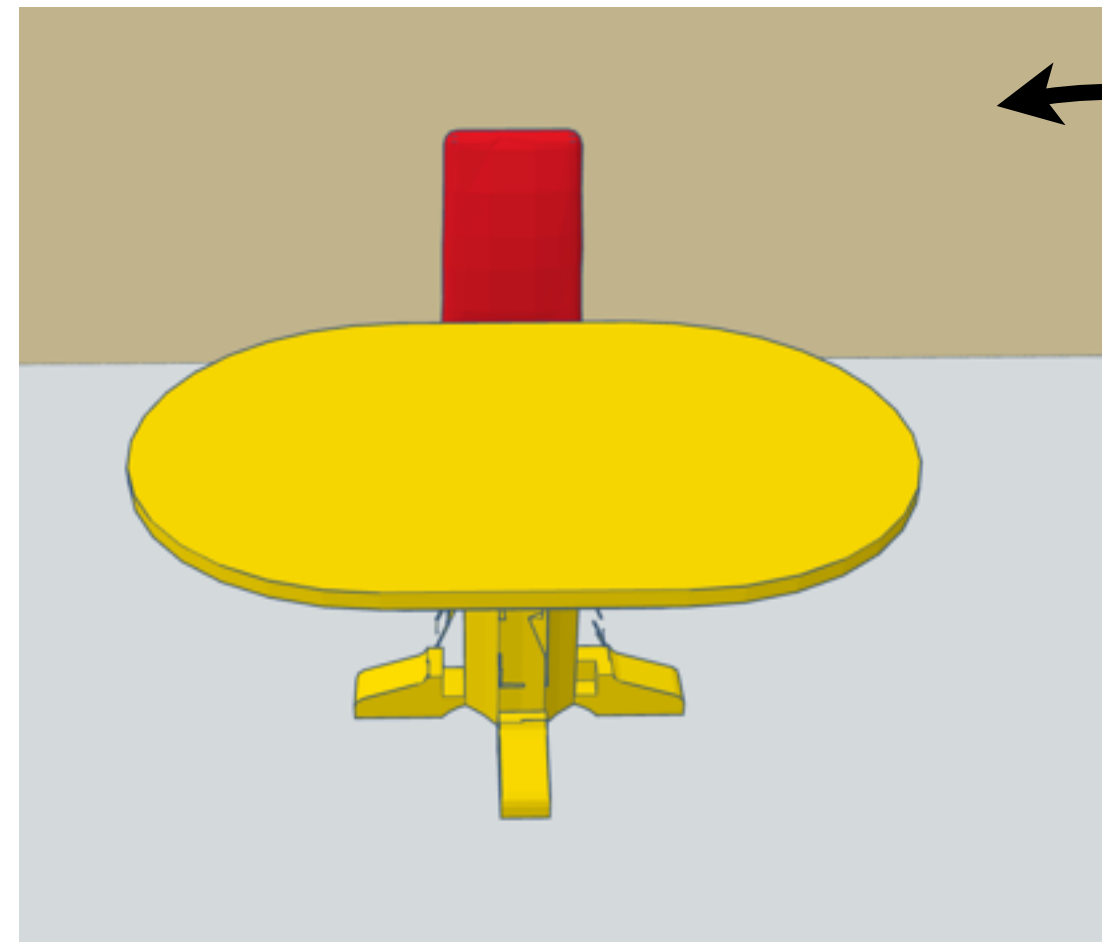
Goal: ~~Semantic Scene Completion~~ Semantic








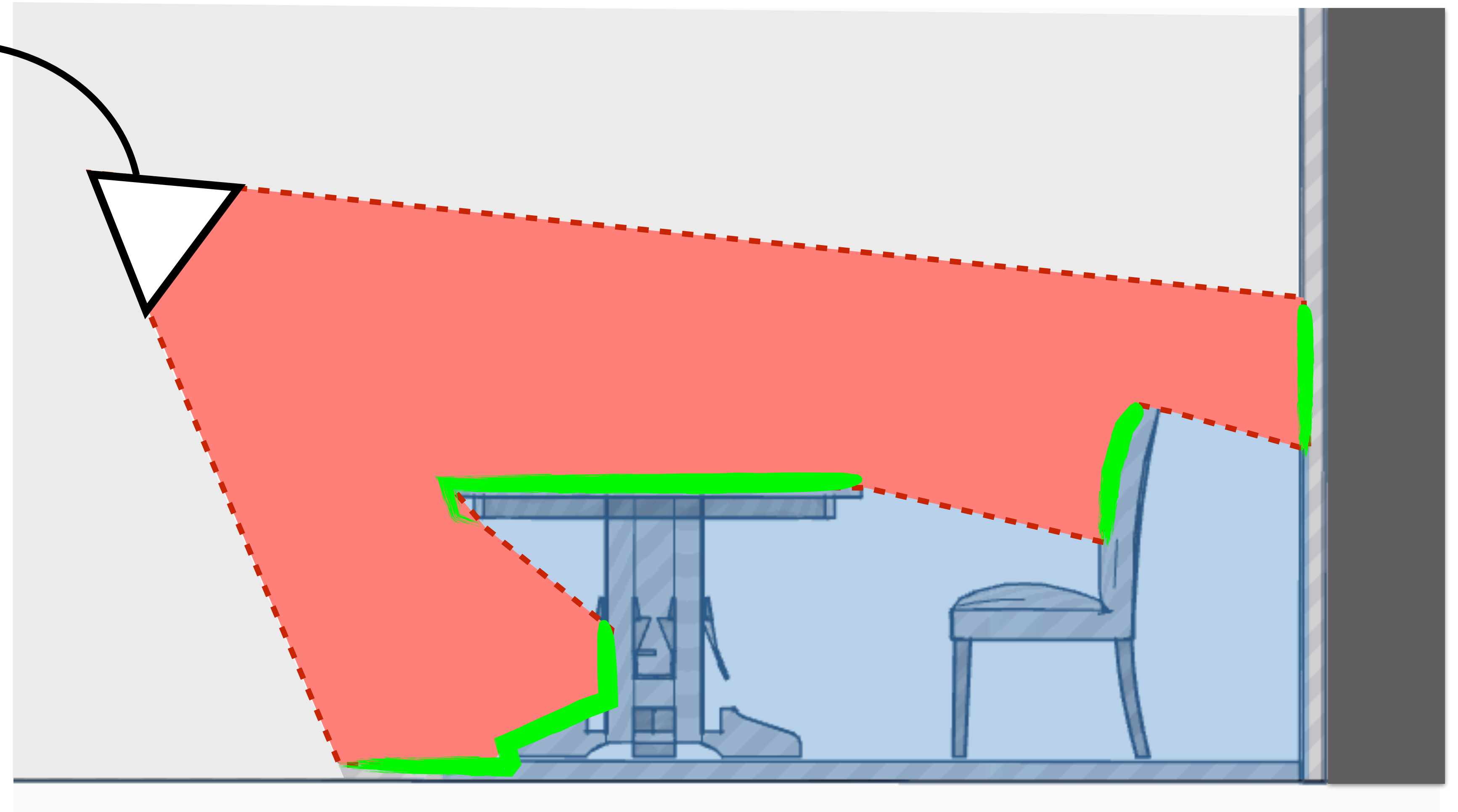
Input:
Single view depth map

Output:
volumetric occupancy + semantic

Problem definition

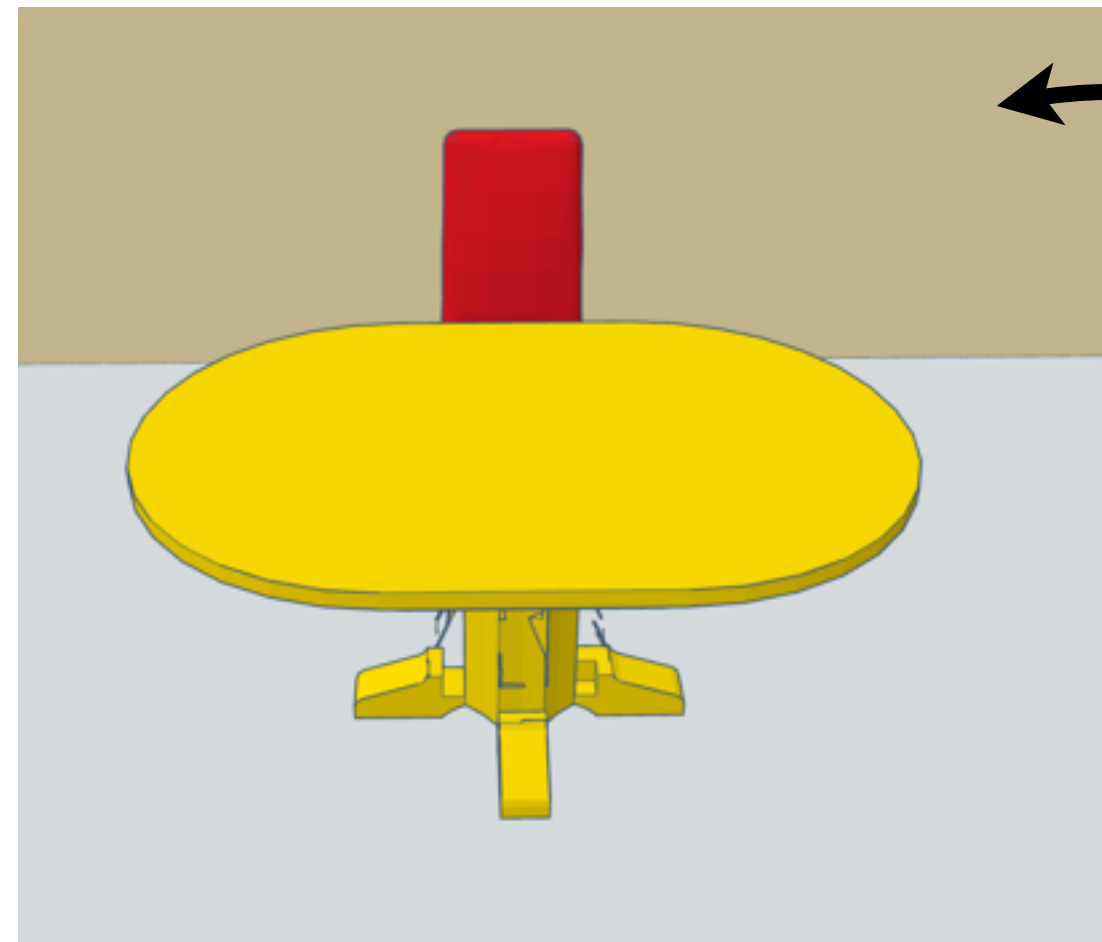


-  visible surface
-  free space
-  occluded space
-  outside view
-  outside room

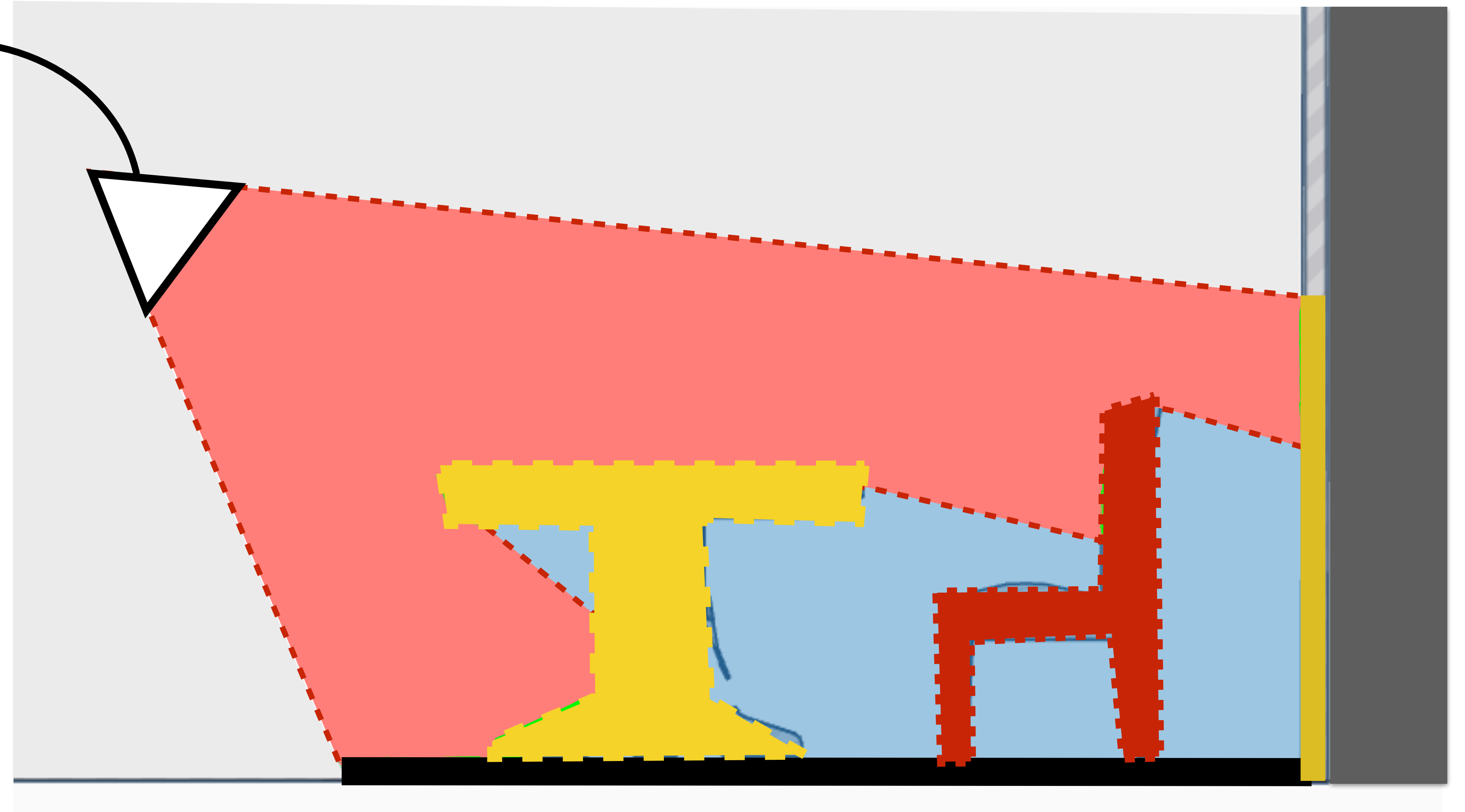


3D Scene

Problem definition

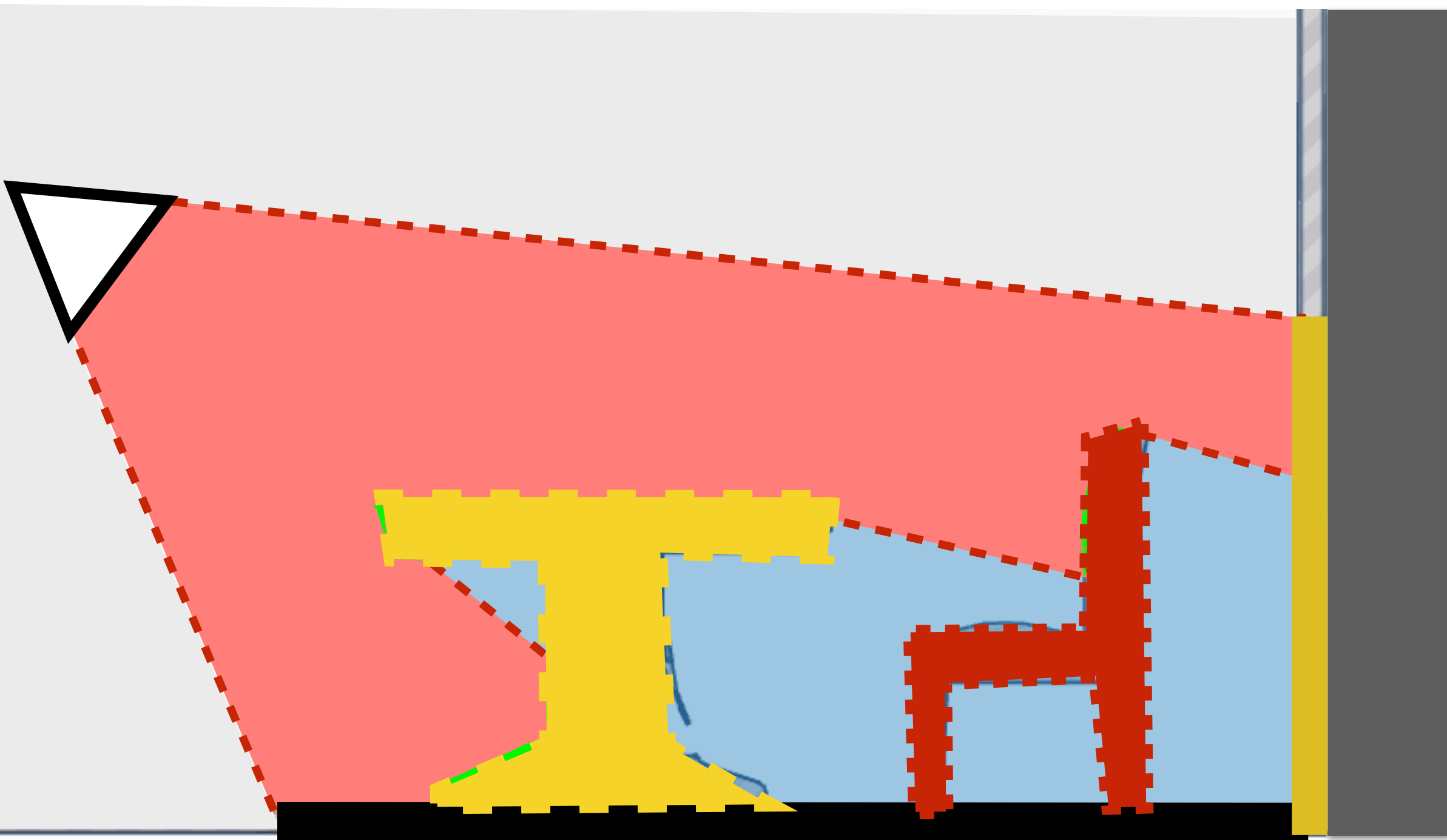


- visible surface
- free space
- occluded space
- outside view
- outside room



3D Scene

Prior work



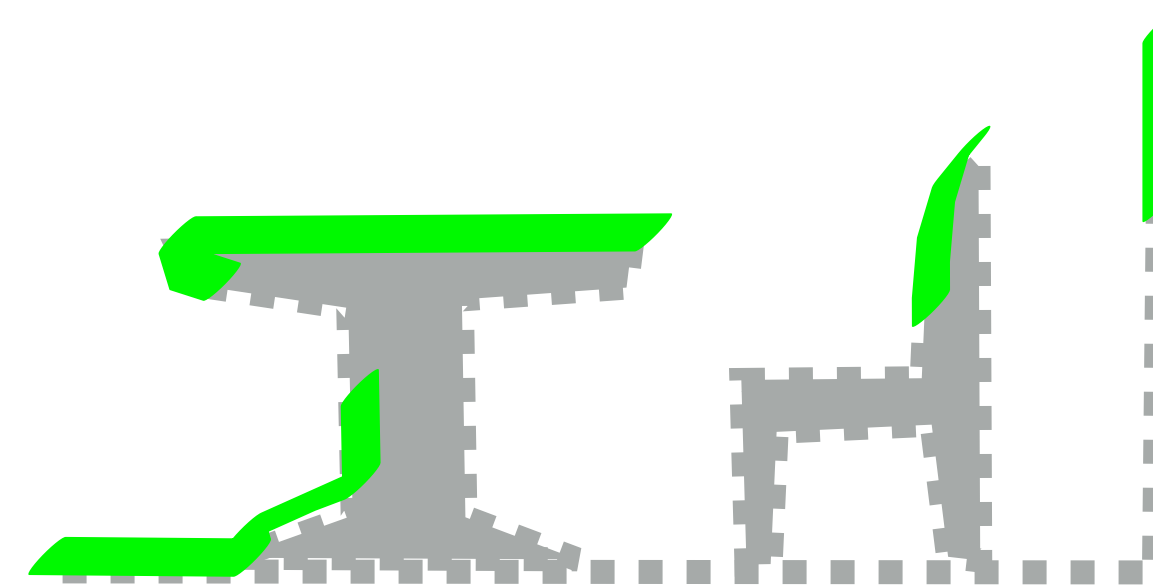
3D Scene



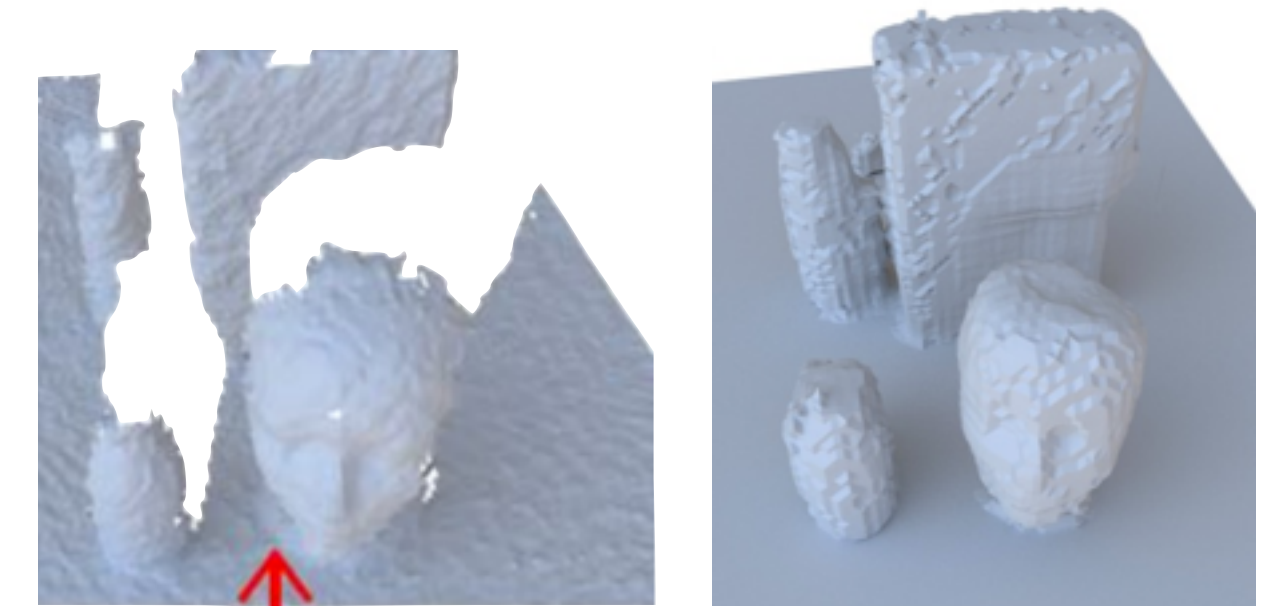
surface segmentation



[Silberman *et al.*]



scene completion



[Firman *et al.*]



semantic scene completion

Object occupancy and semantic

Partial scan of common object



What is this?

What's the complete shape?



a book?

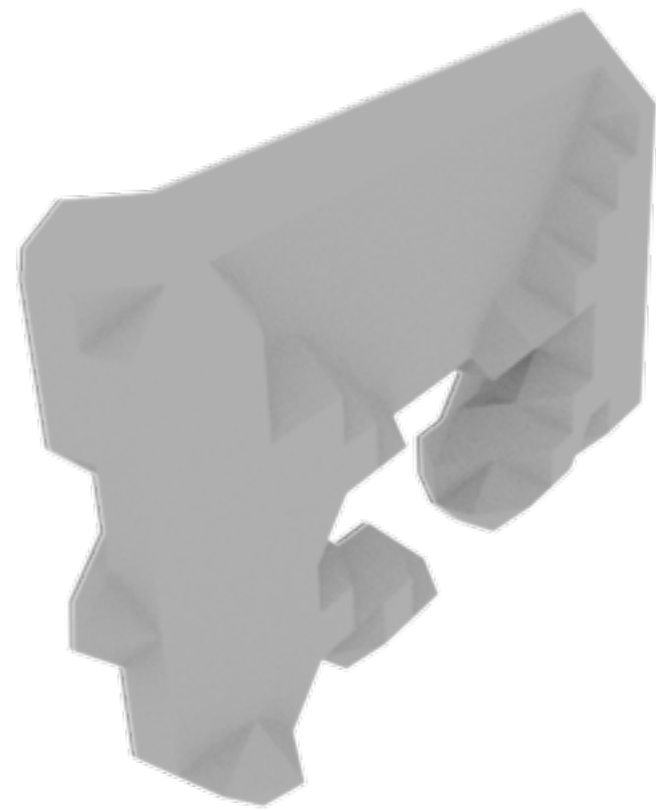


a TV?



a bed?

Object occupancy and semantic



What is this?
What's the complete shape?

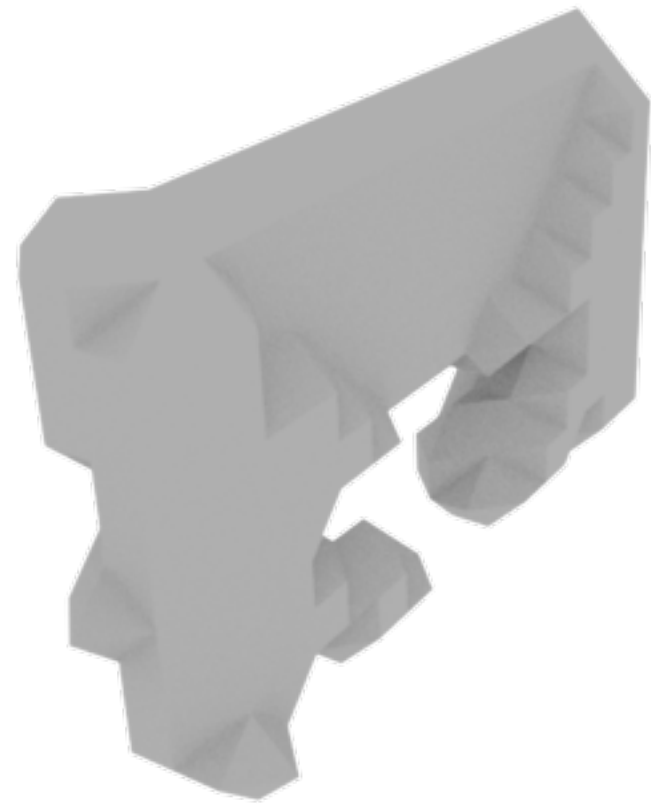
Semantic
meaning

It is part of a chair!

Occupancy
patterns



Object occupancy and semantic



What is this?
What's the complete shape?

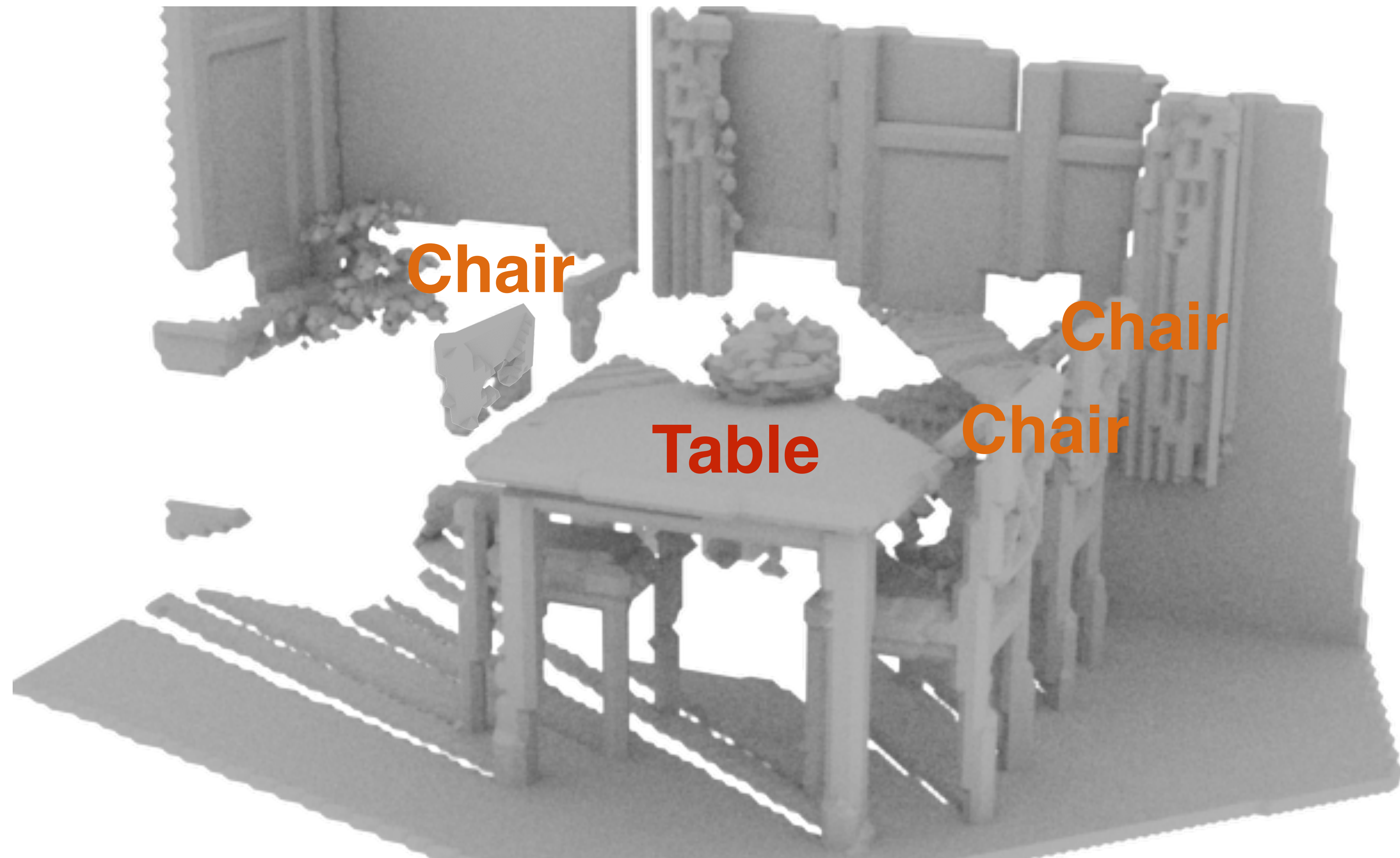
Semantic
meaning

It is part of a chair!

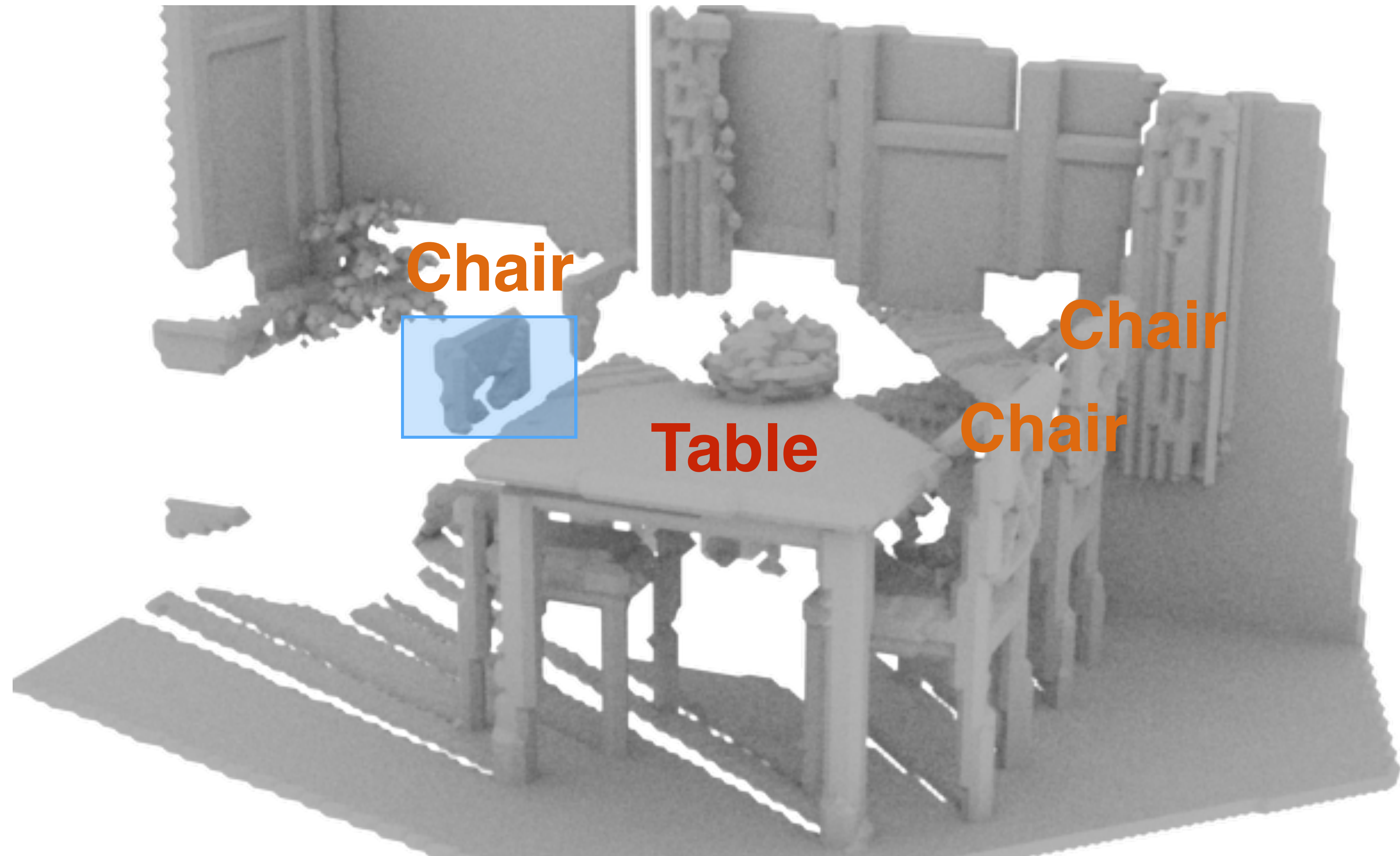
Occupancy
patterns



3D context !



3D context with BIG receptive field!



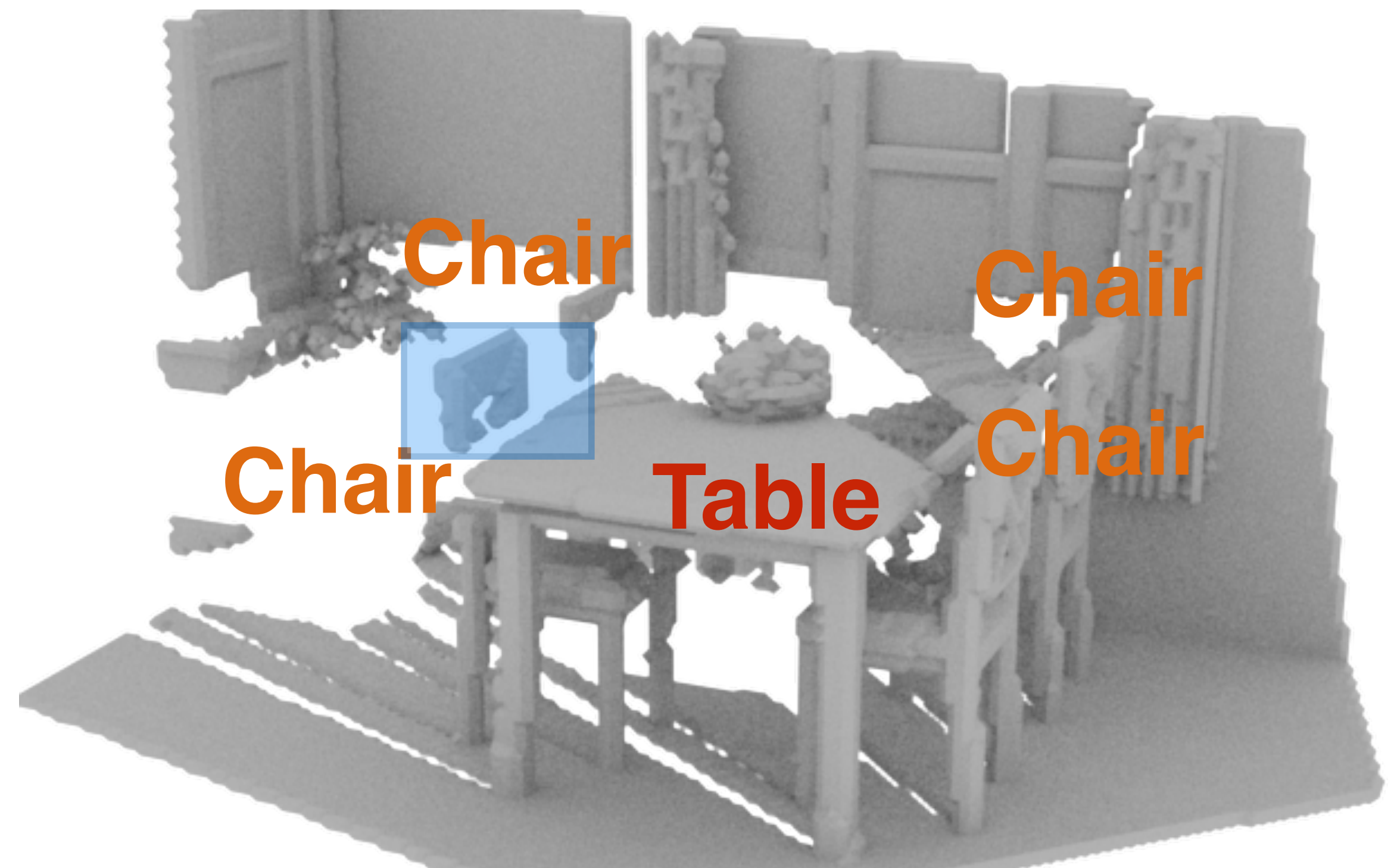
Key ideas:

1. Object occupancy and the identity are tightly intertwined.
2. It is important to capture and understand 3D context with big receptive fields.

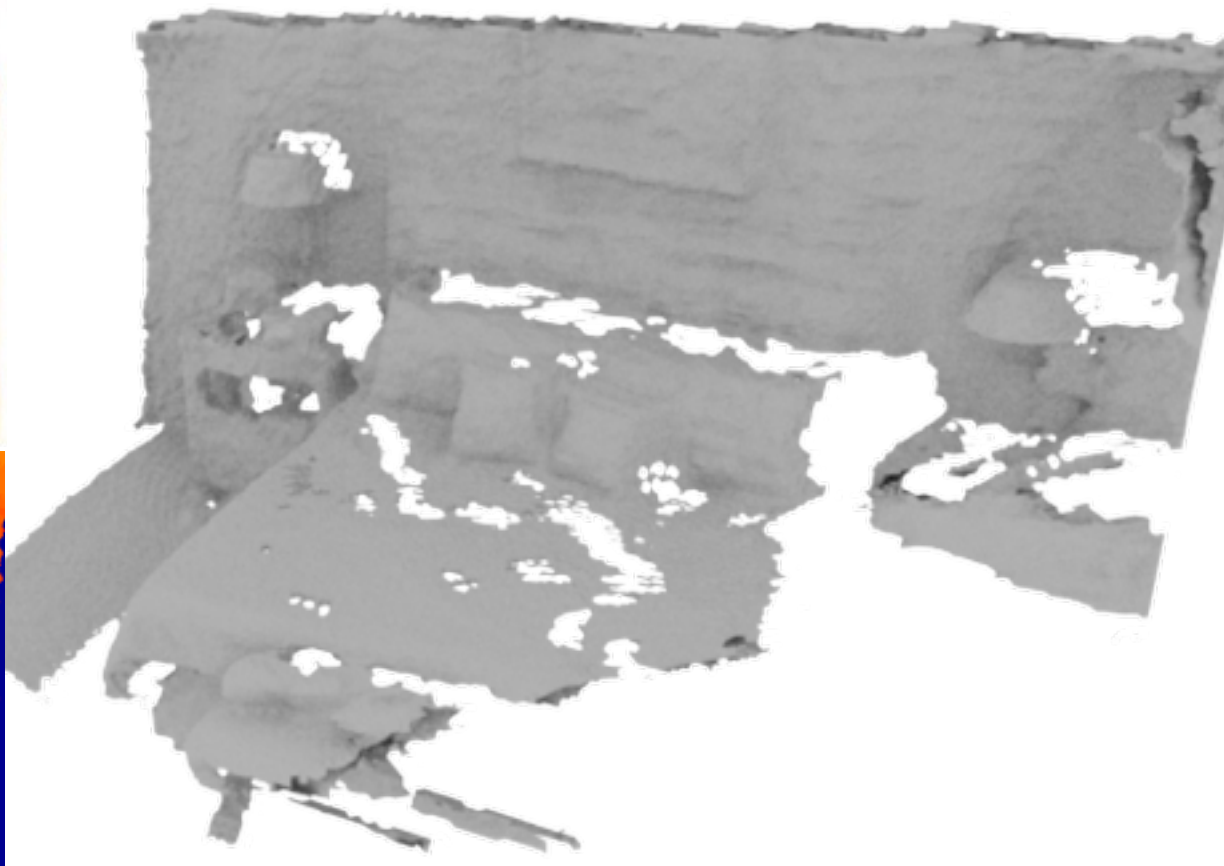
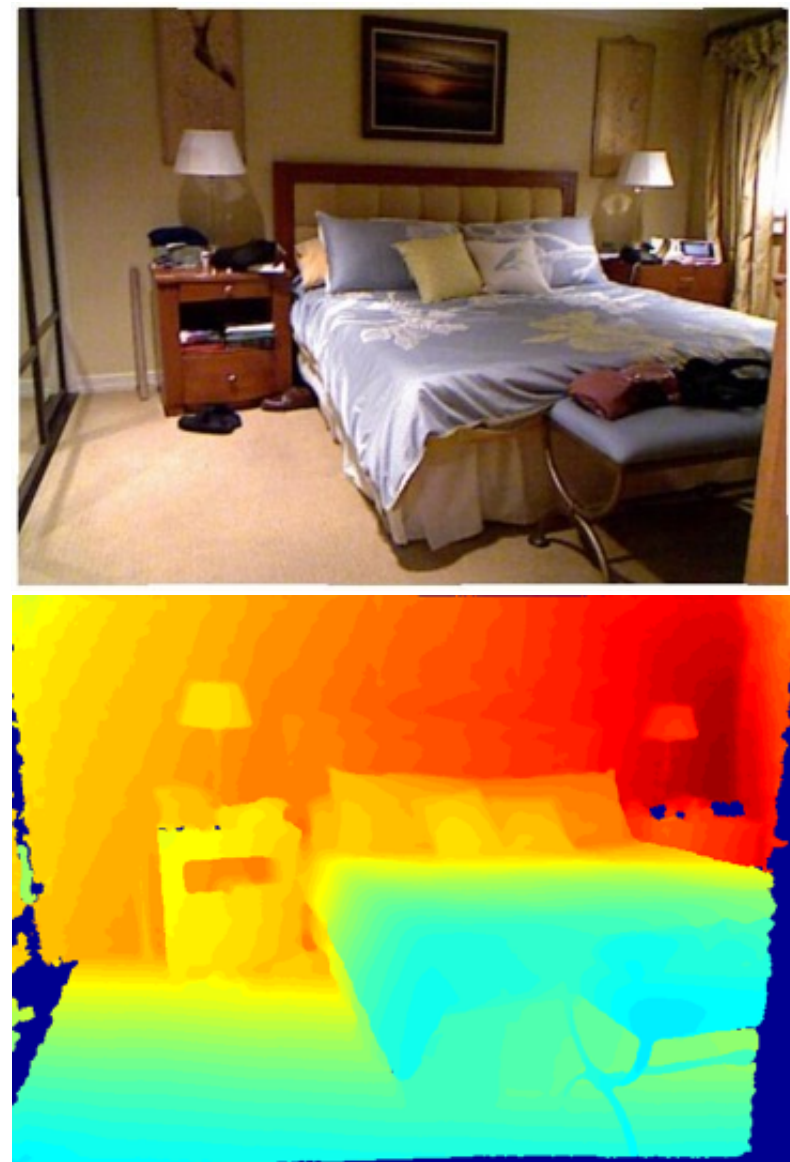
Semantic
meaning

Occupancy
patterns

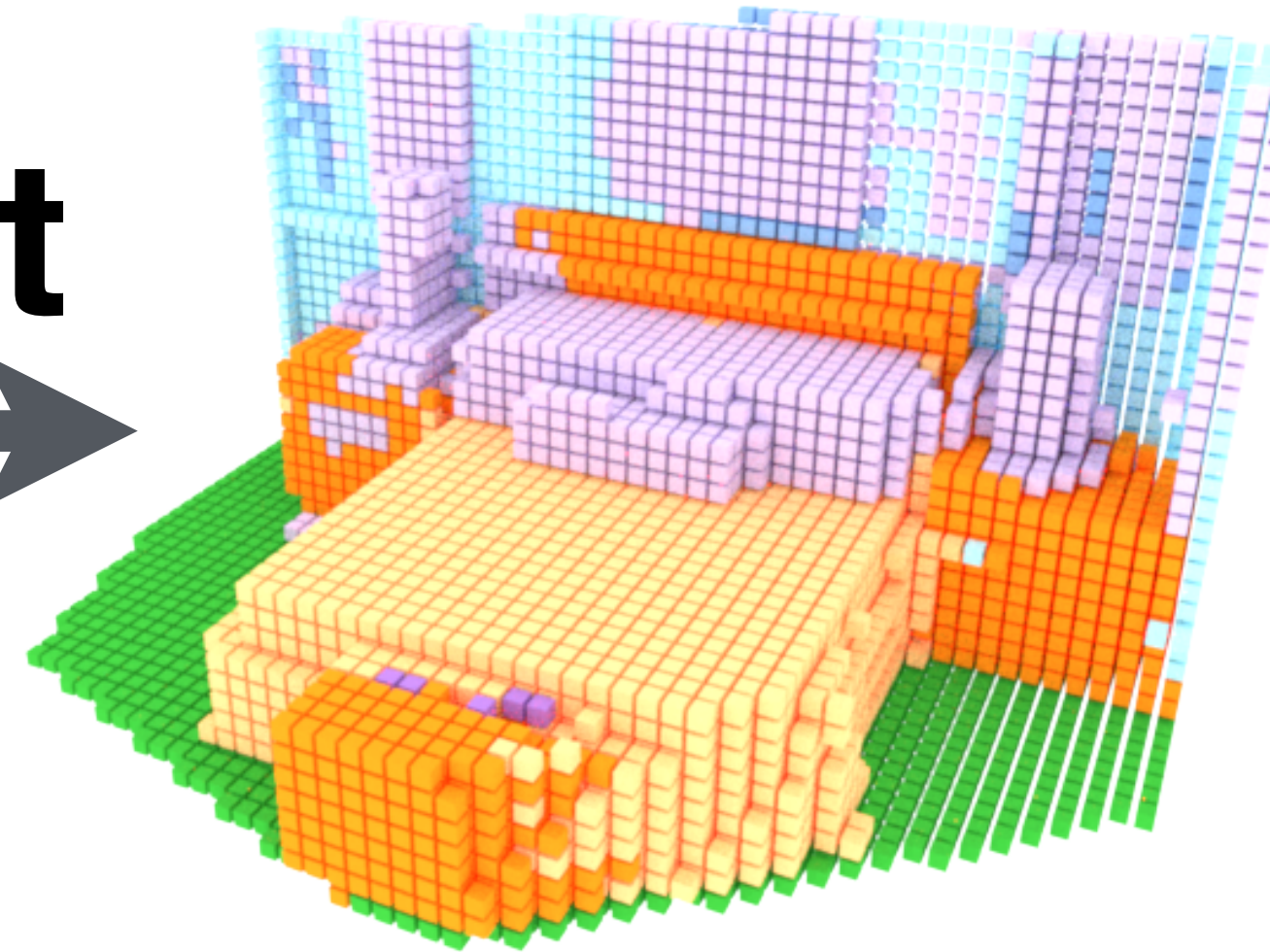
It is part of a chair!



Semantic Scene Completion Network



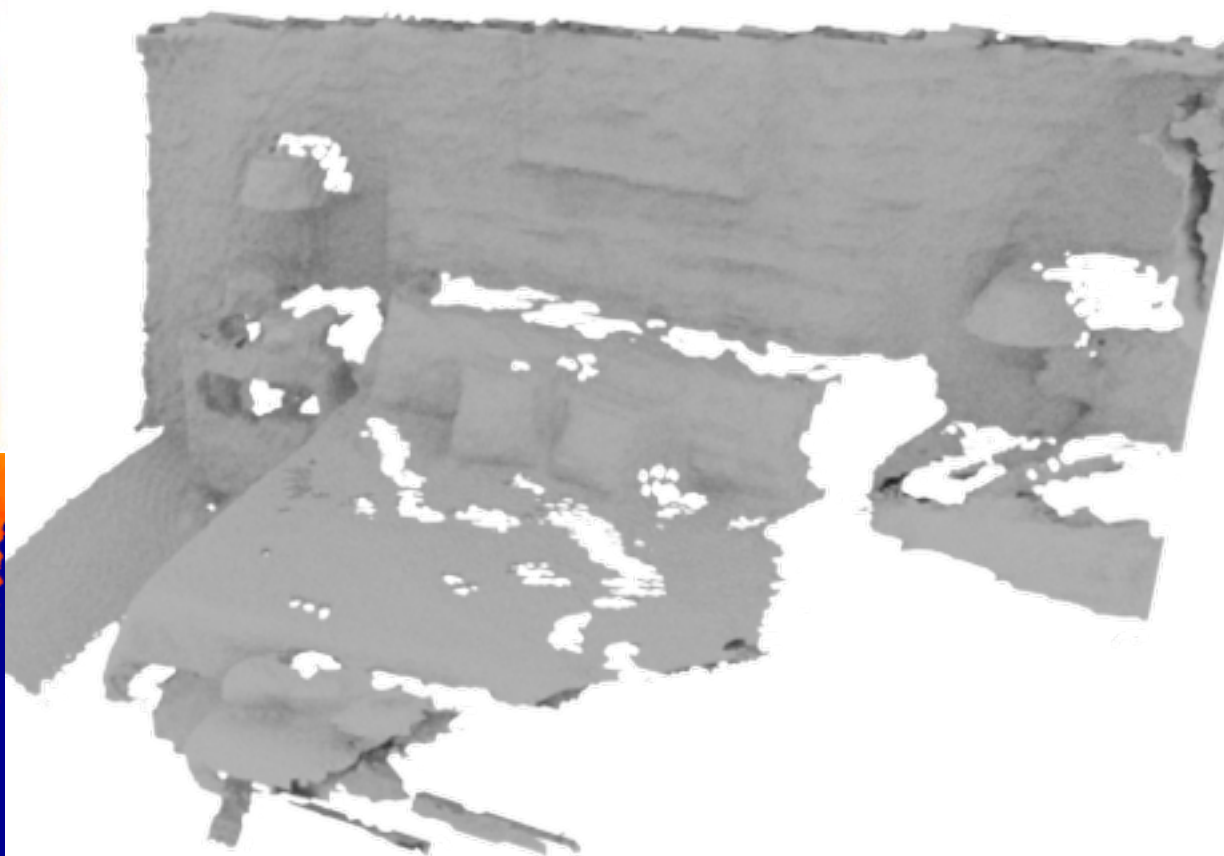
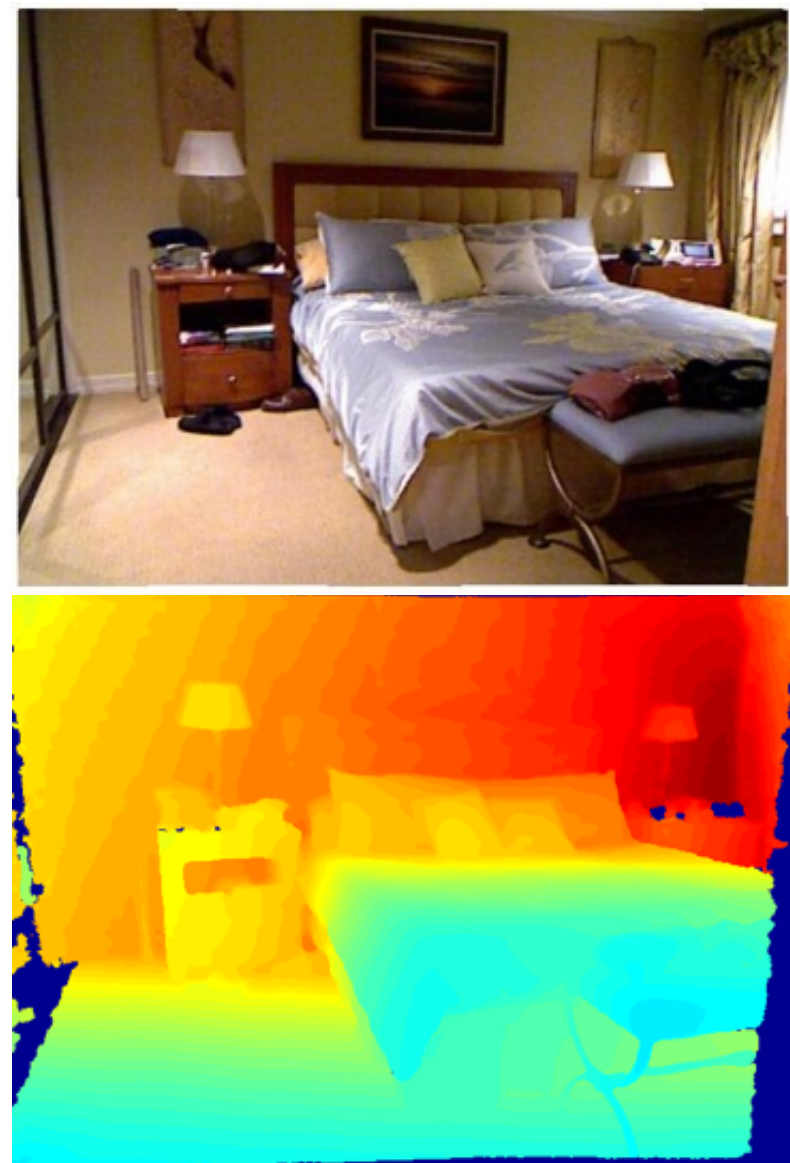
SSCNet



Input: Single view depth map

Output: Semantic scene completion

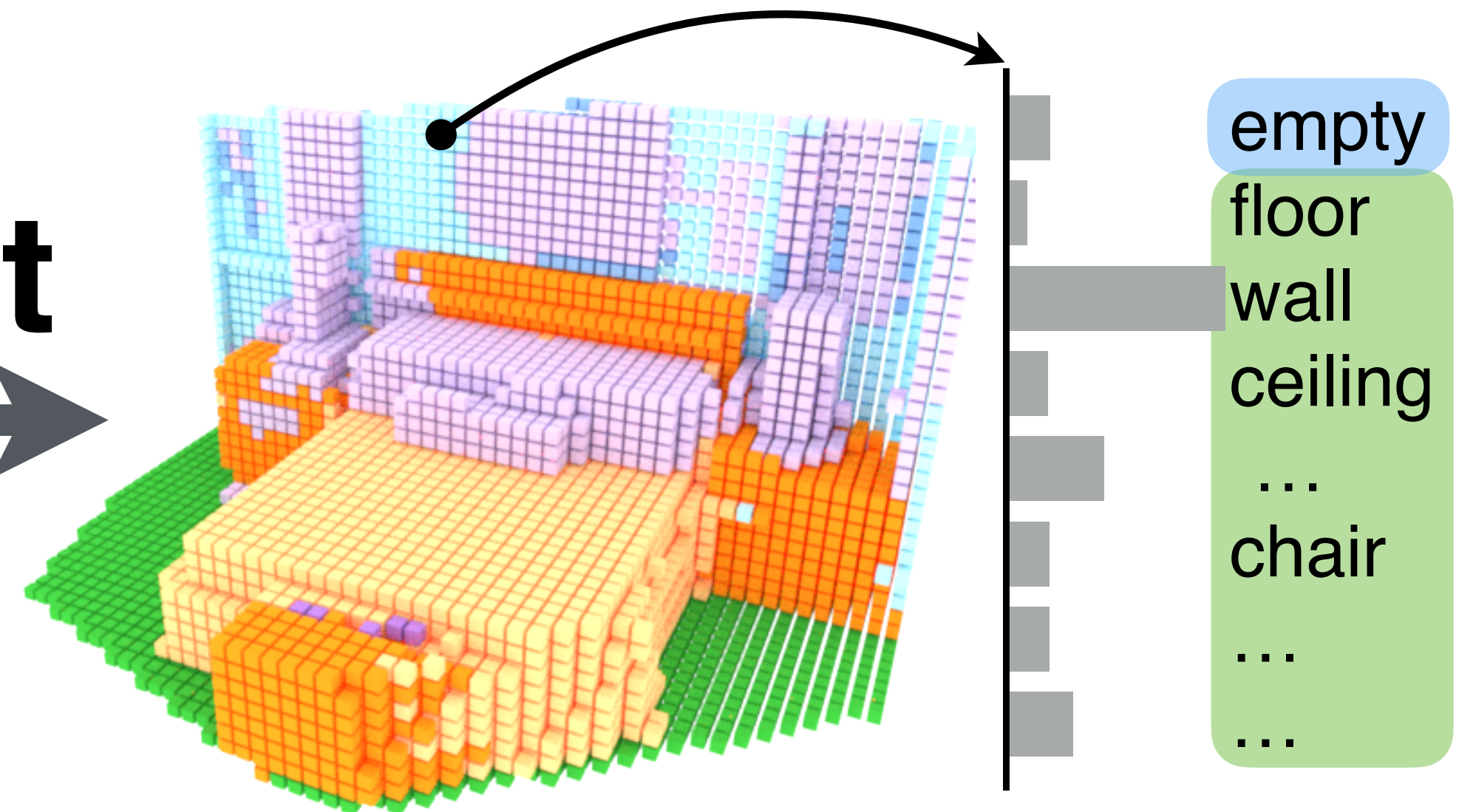
Semantic Scene Completion Network



SSCNet



Prediction: $N+1$ classes

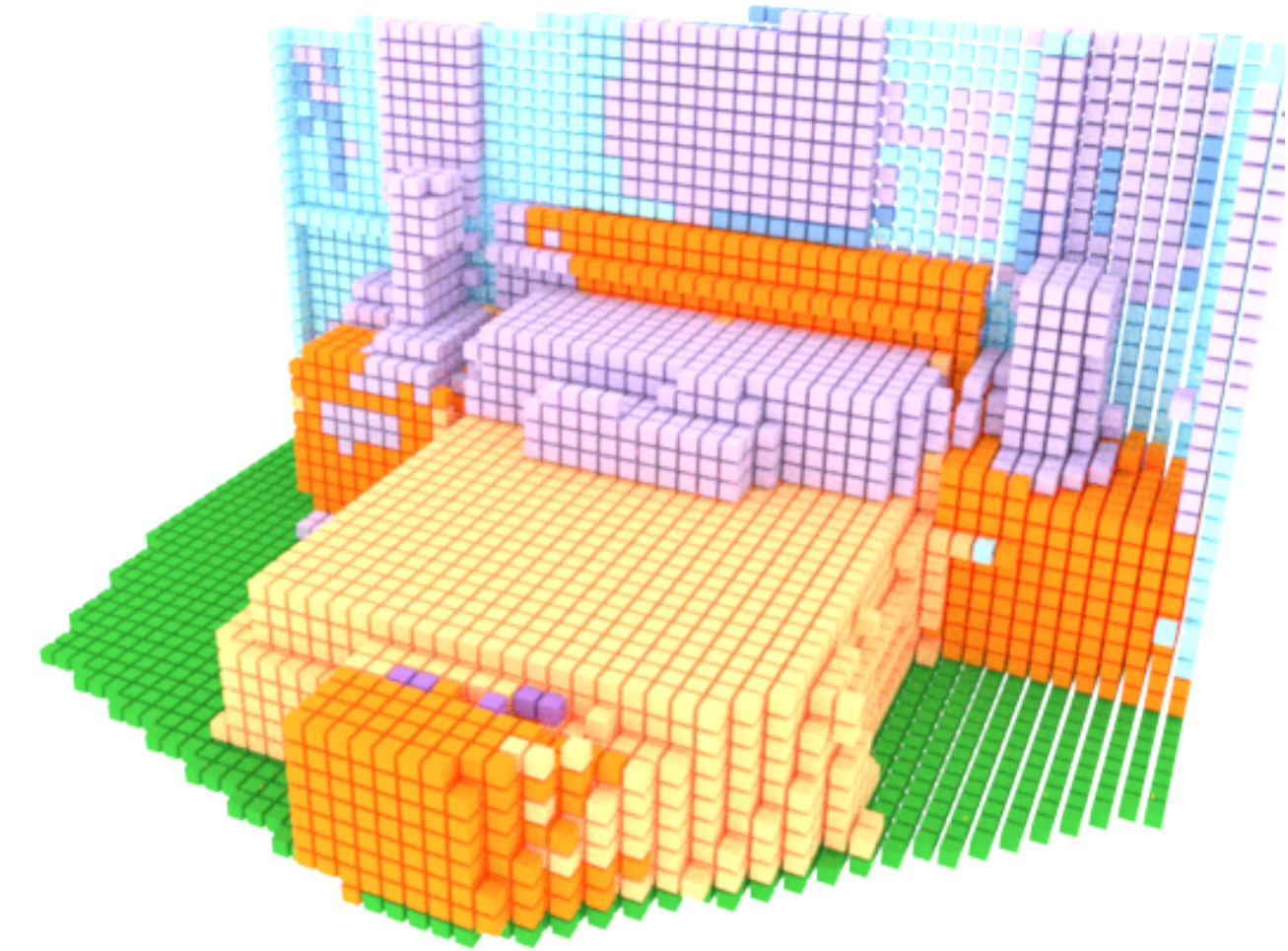
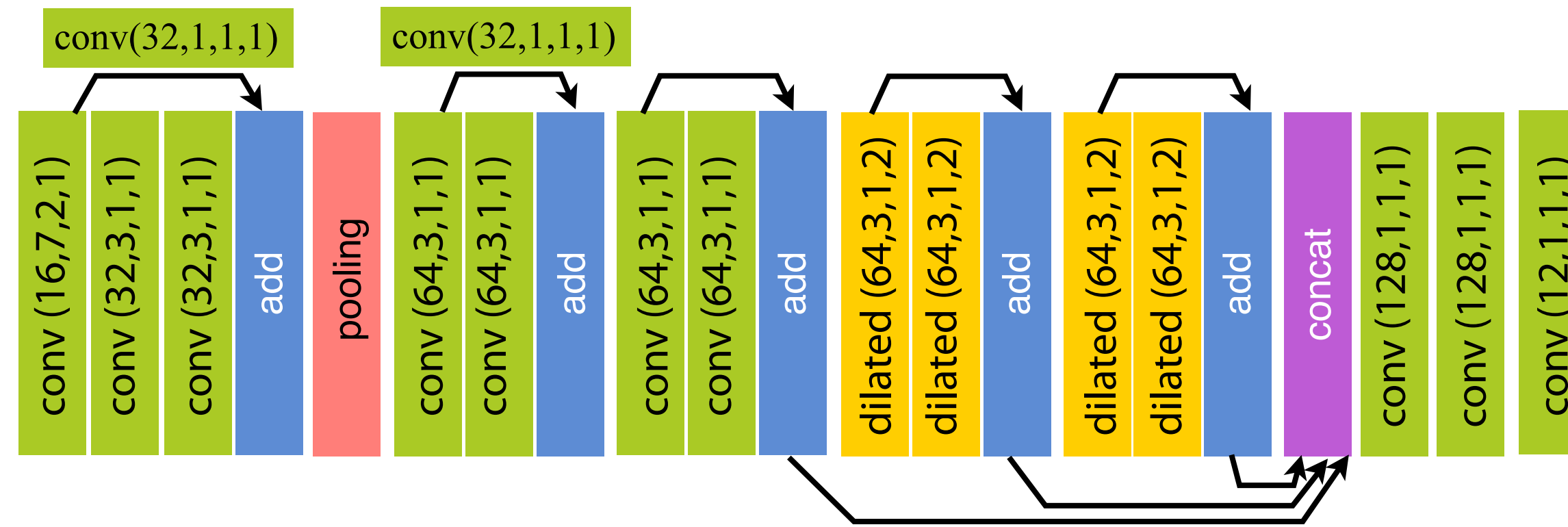
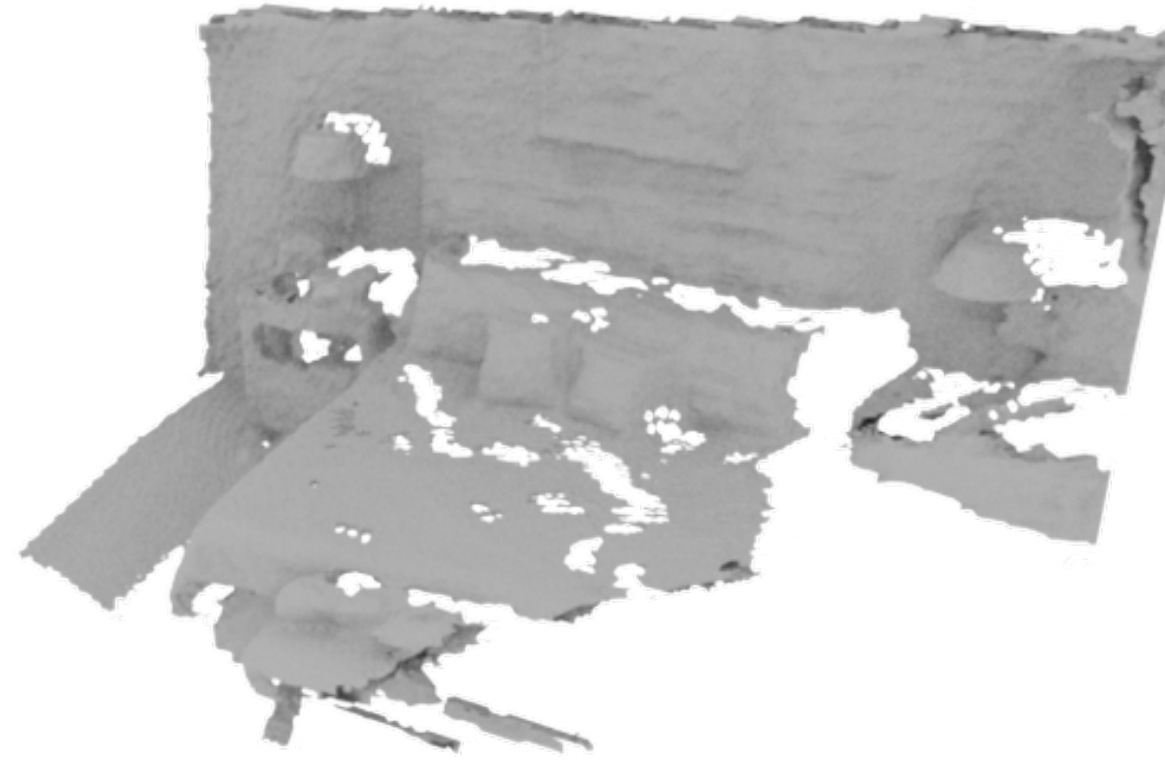


Input: Single view depth map

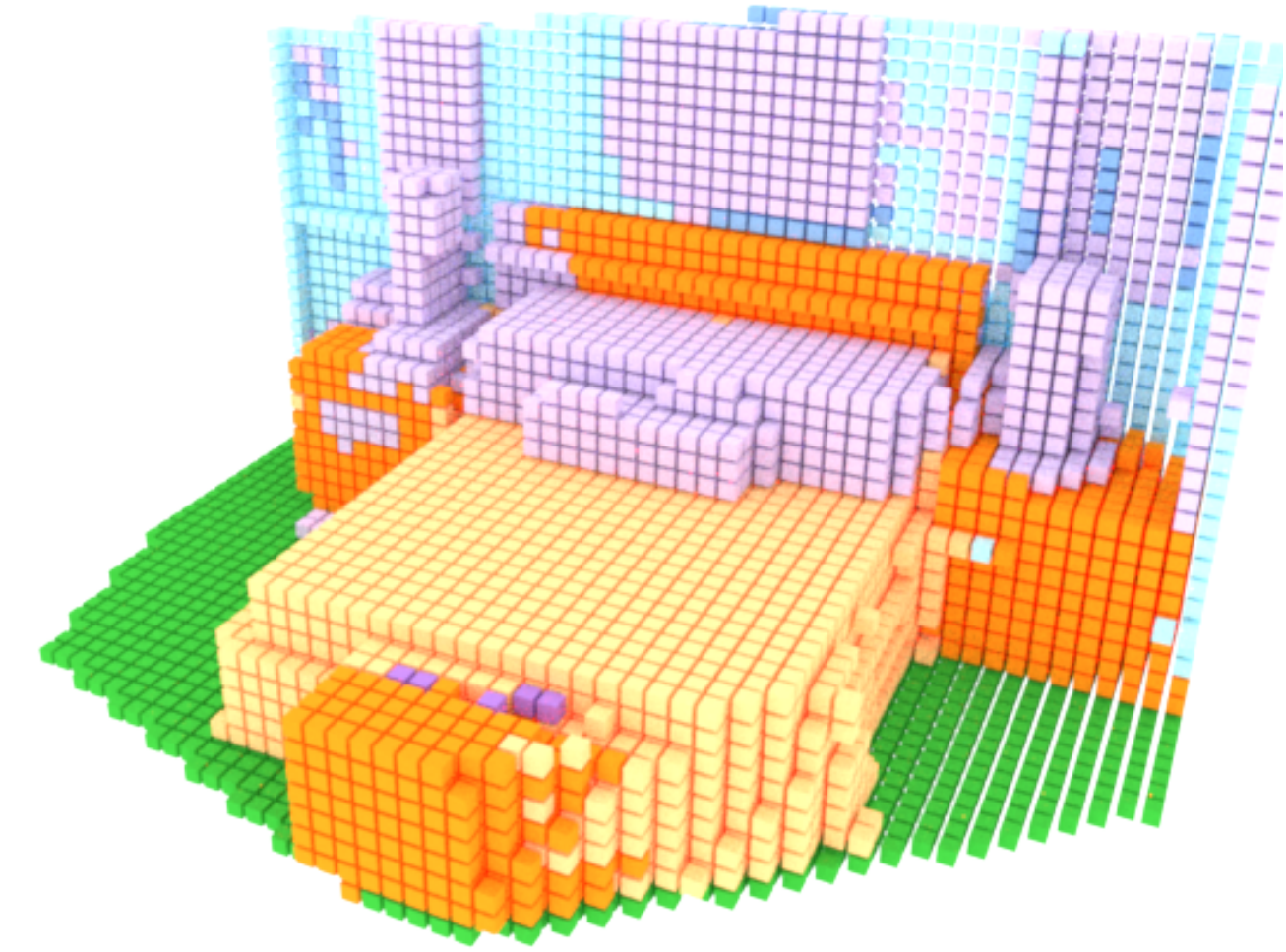
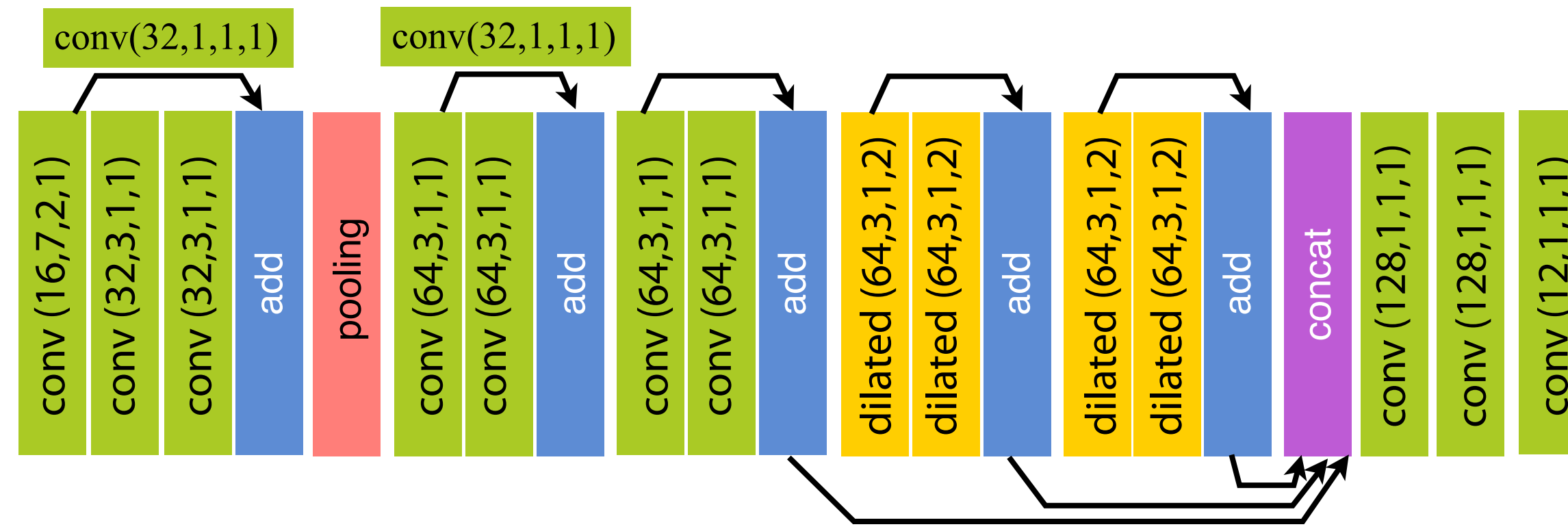
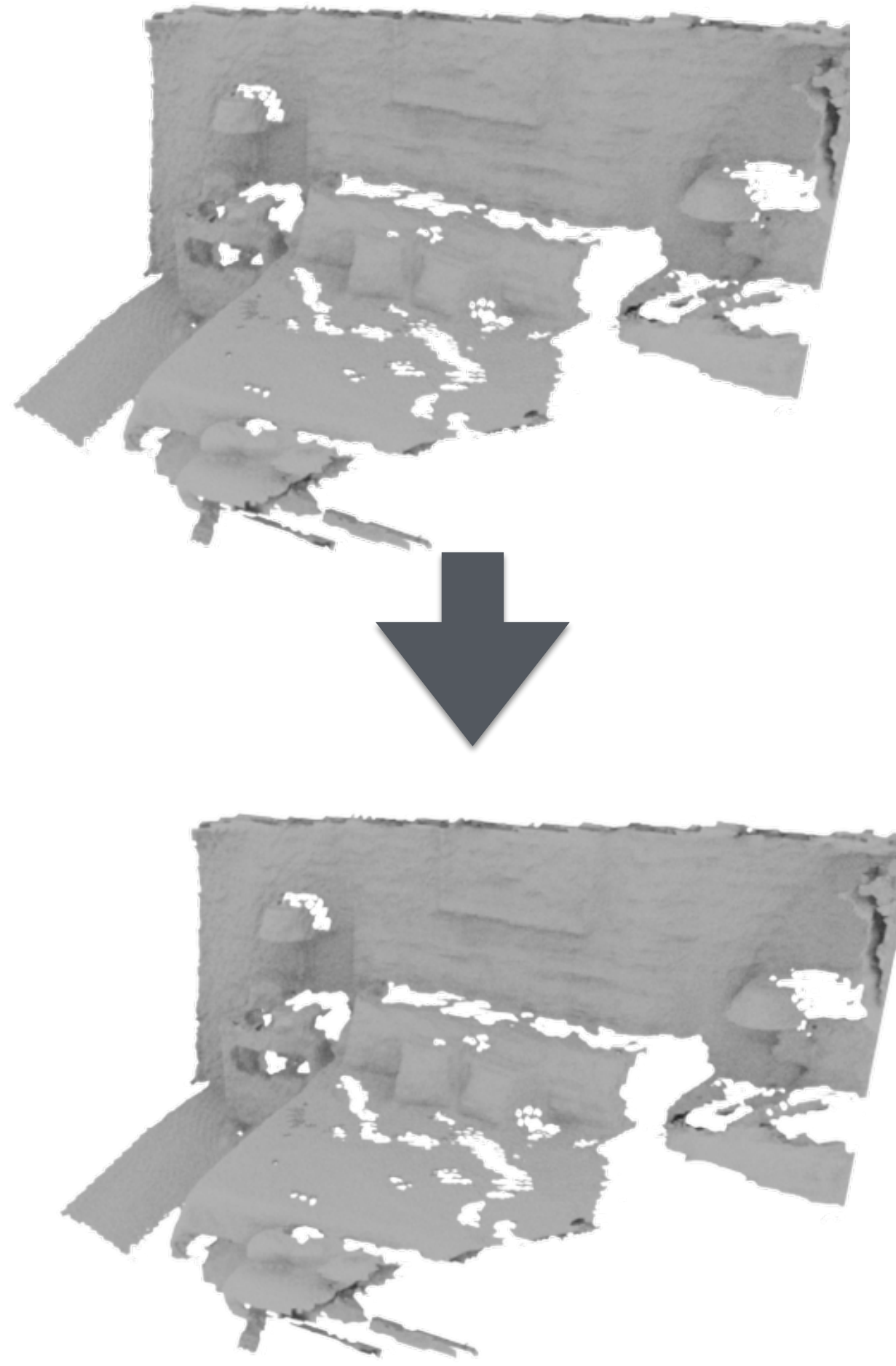
Output: Semantic scene completion

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

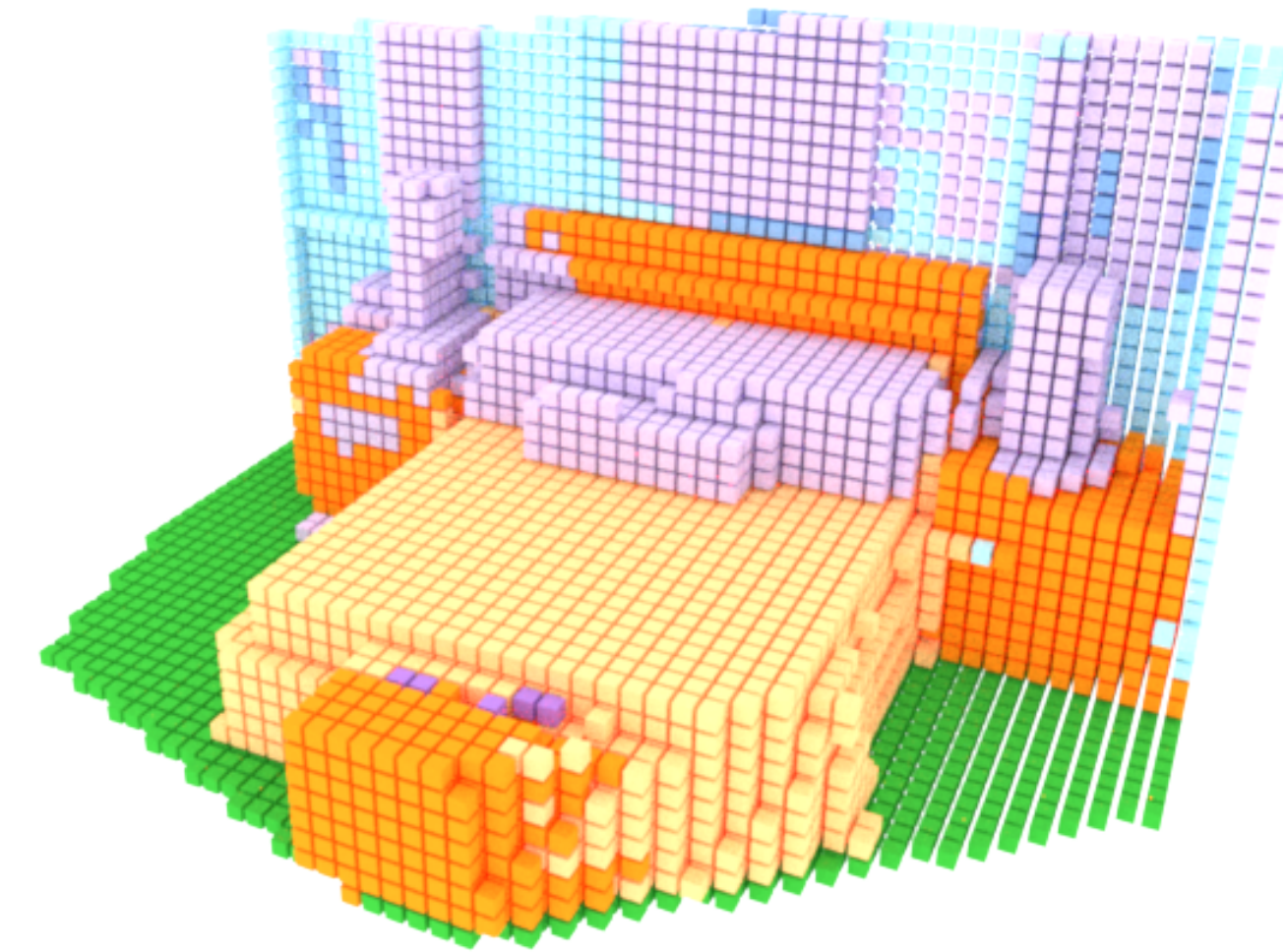
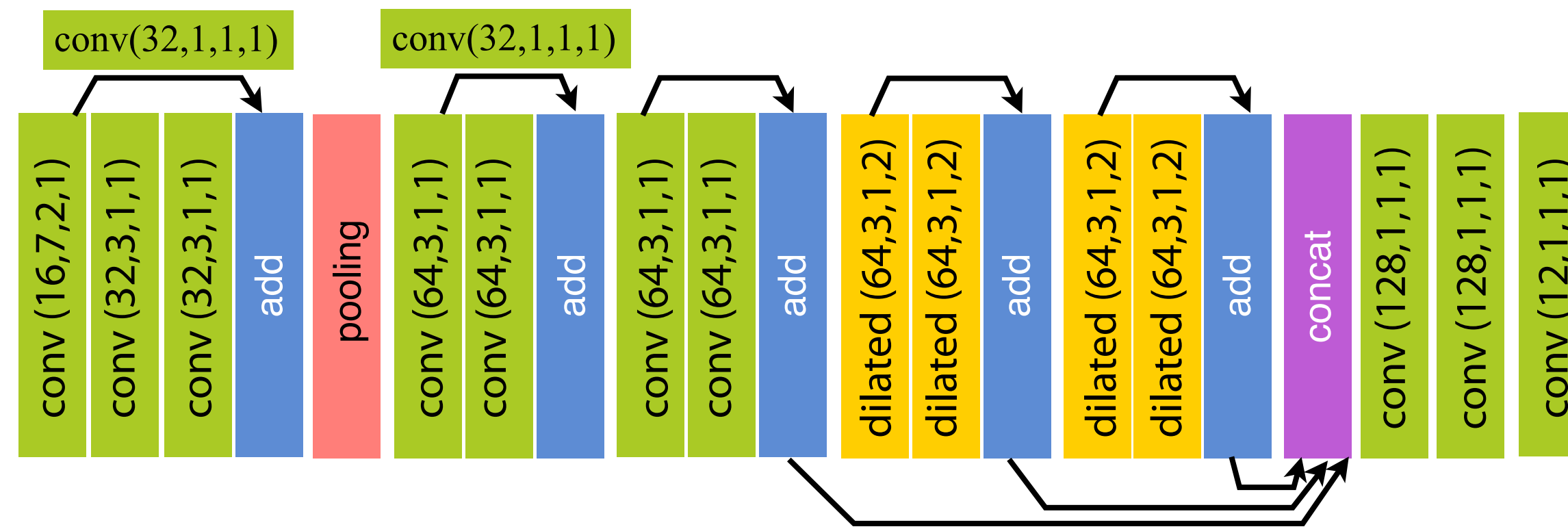
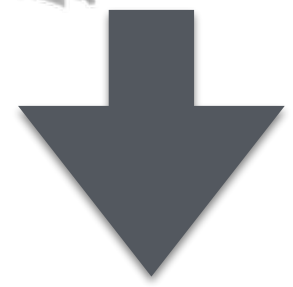
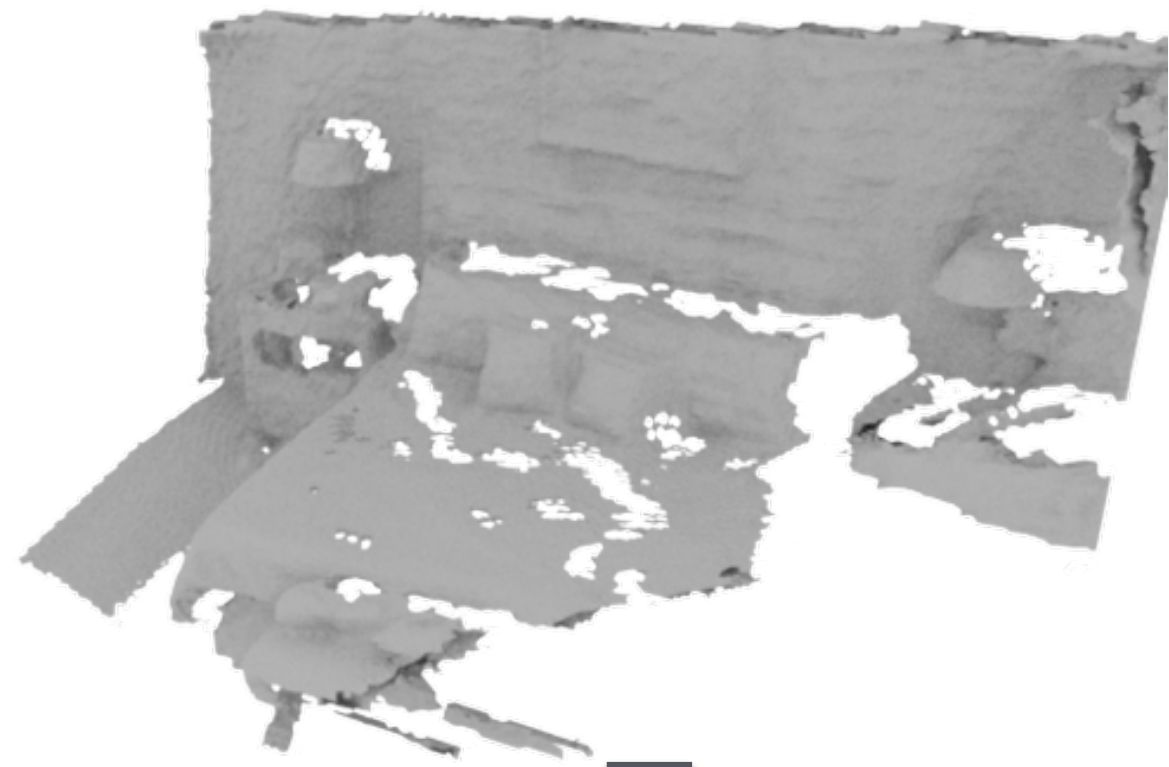
Semantic Scene Completion Network



Semantic Scene Completion Network

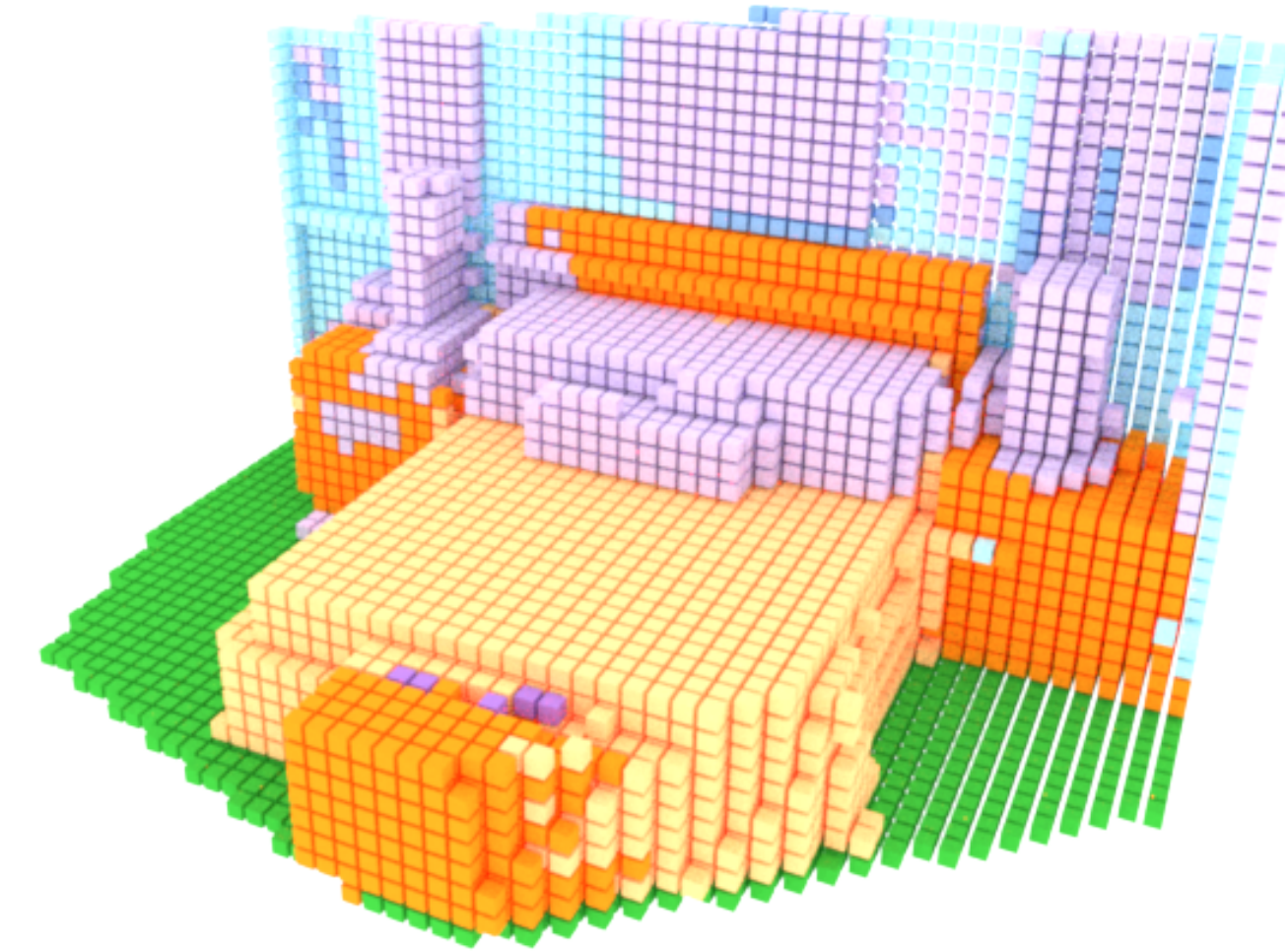
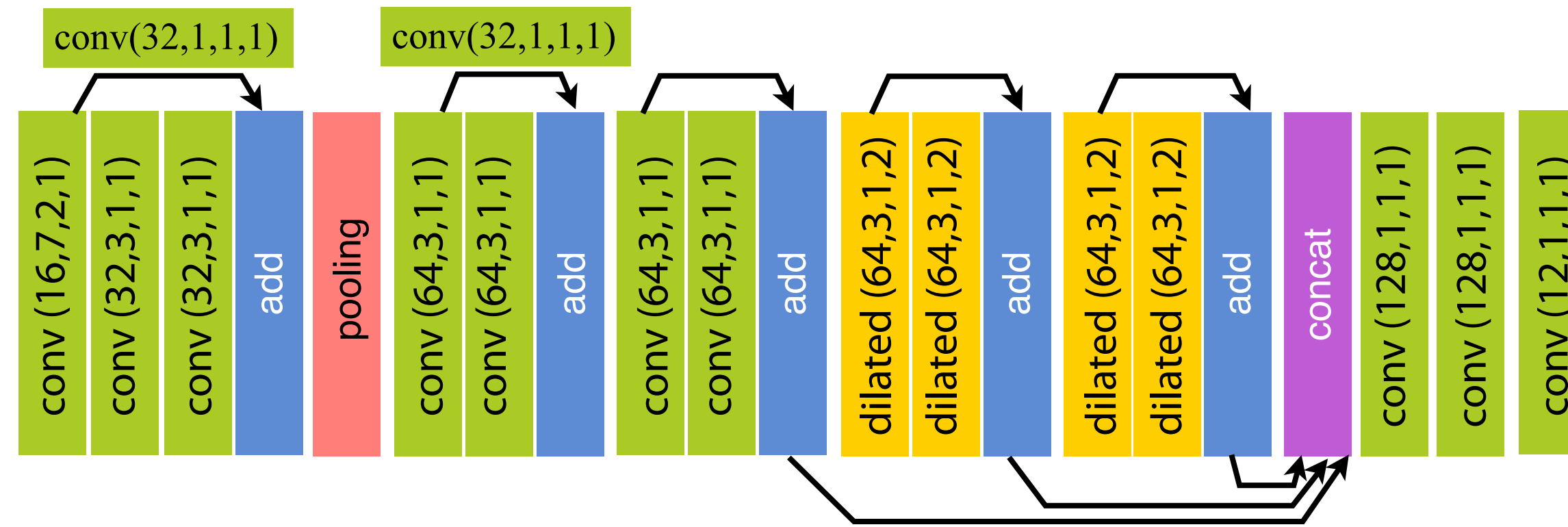
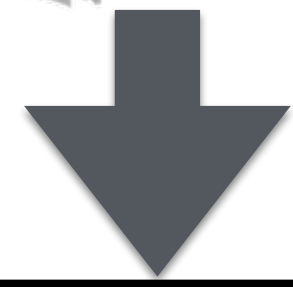
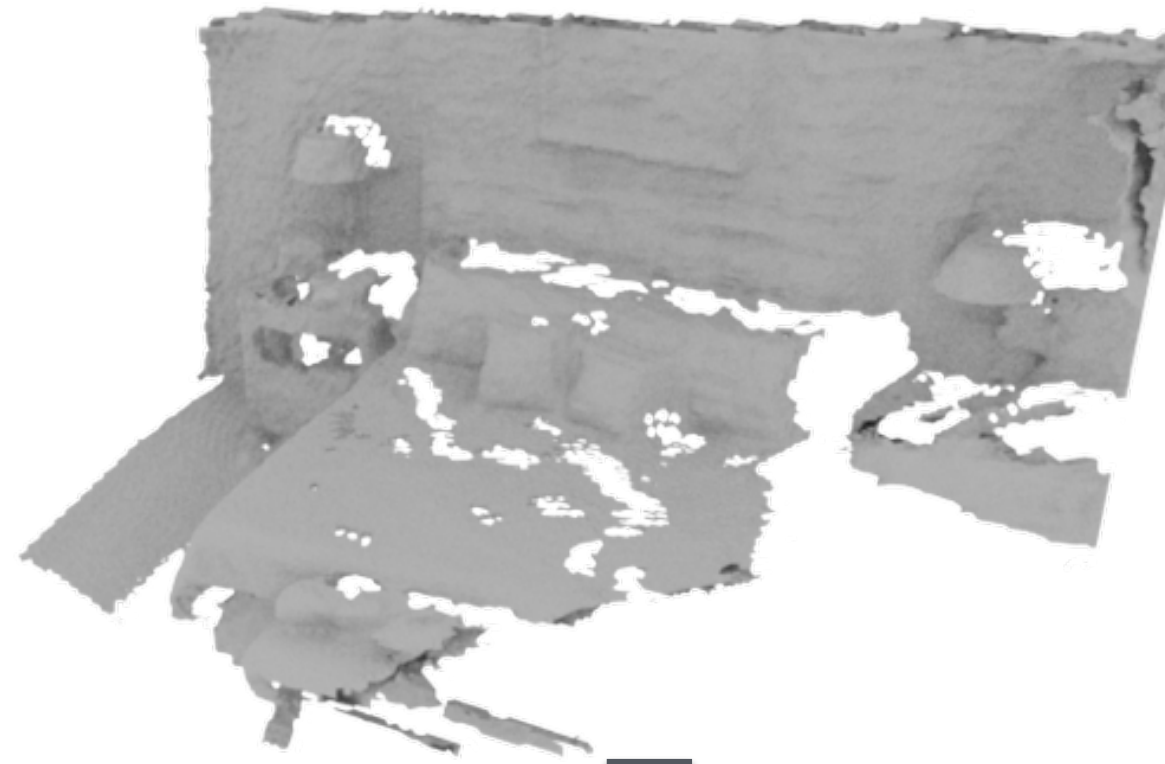


Semantic Scene Completion Network

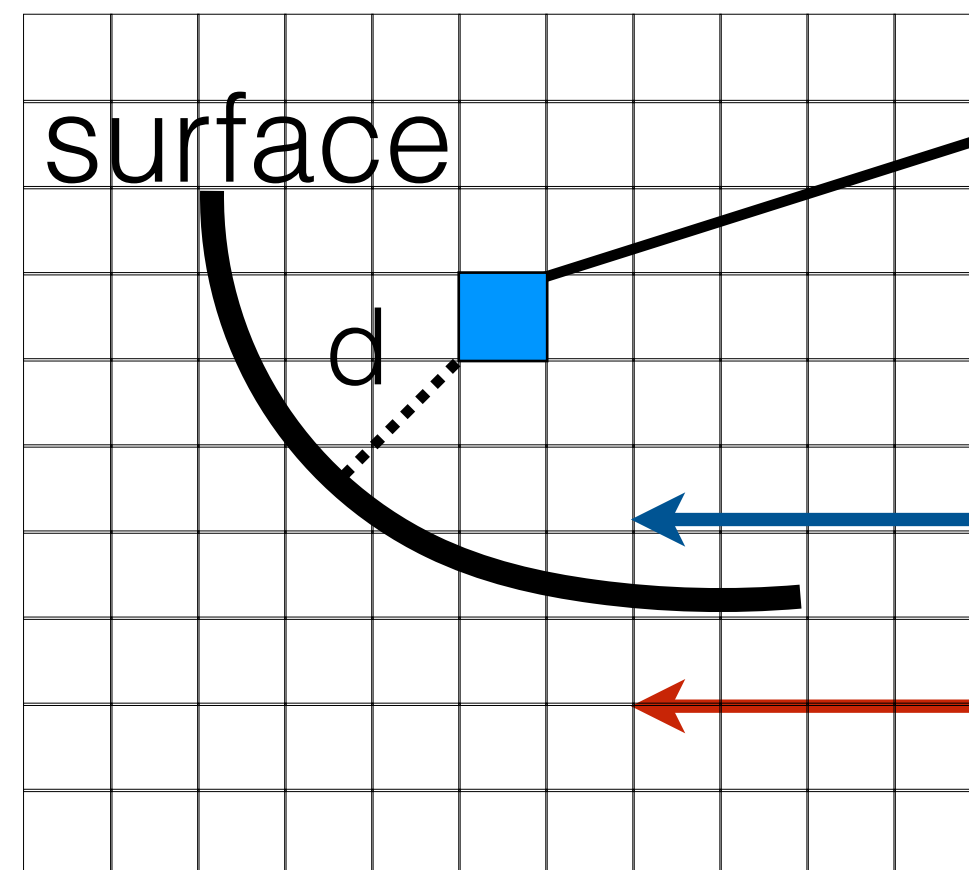


Encode 3D space using flipped TSDF

Semantic Scene Completion Network



Encode 3D space using flipped TSDF



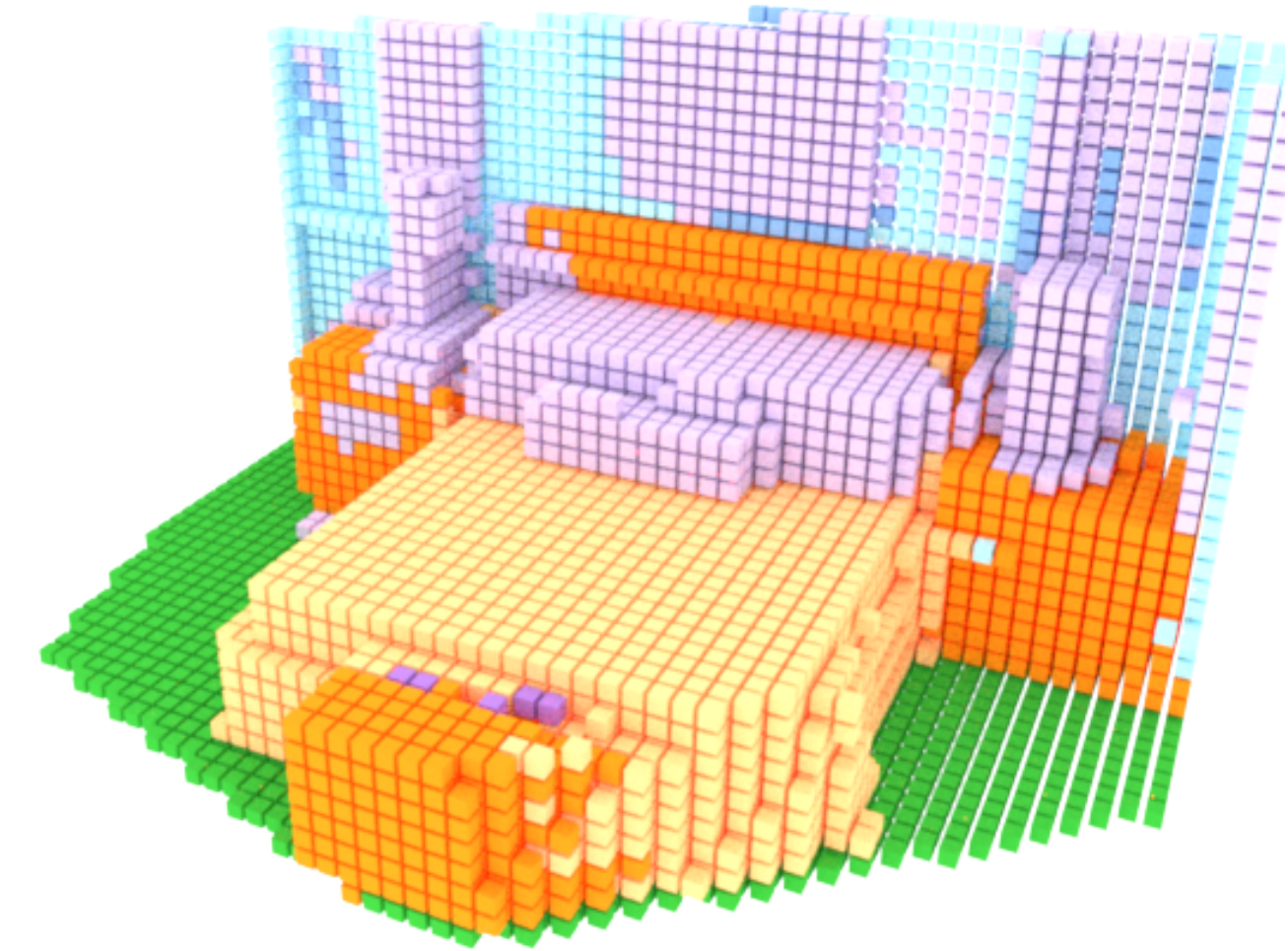
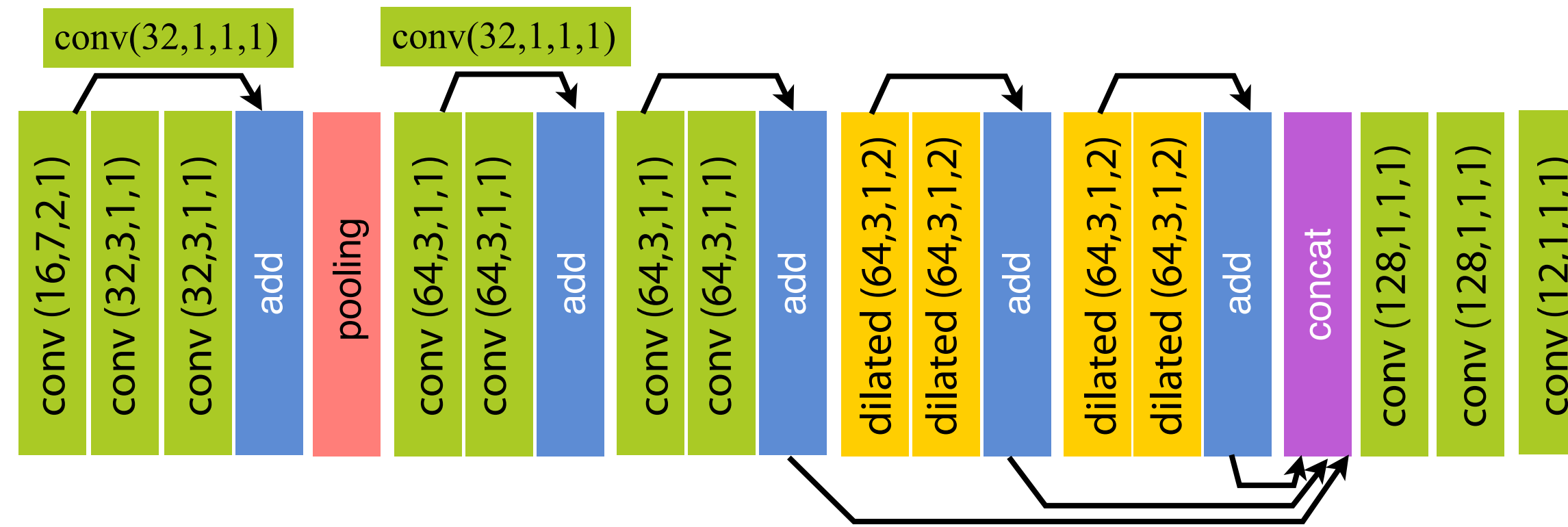
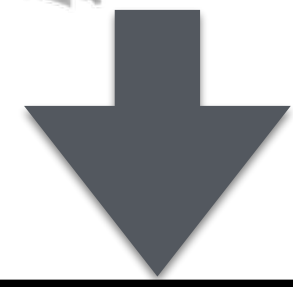
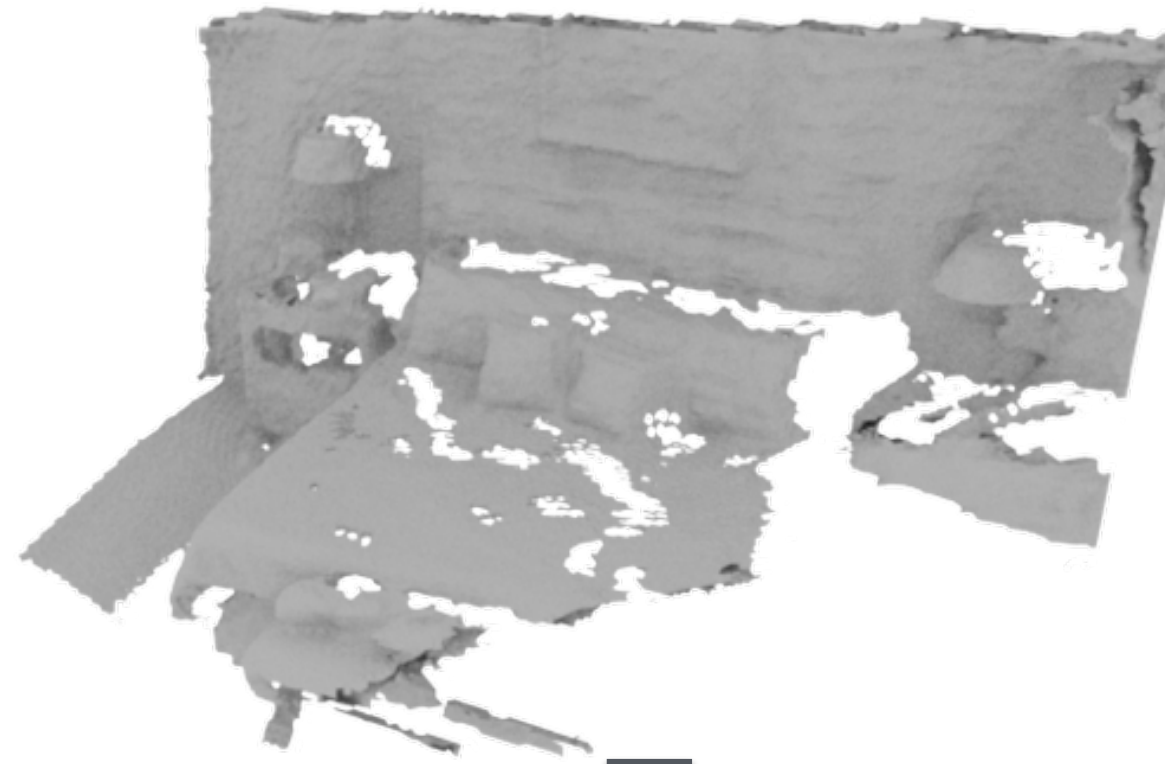
$d = \text{distance to the surface}$

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

Occluded

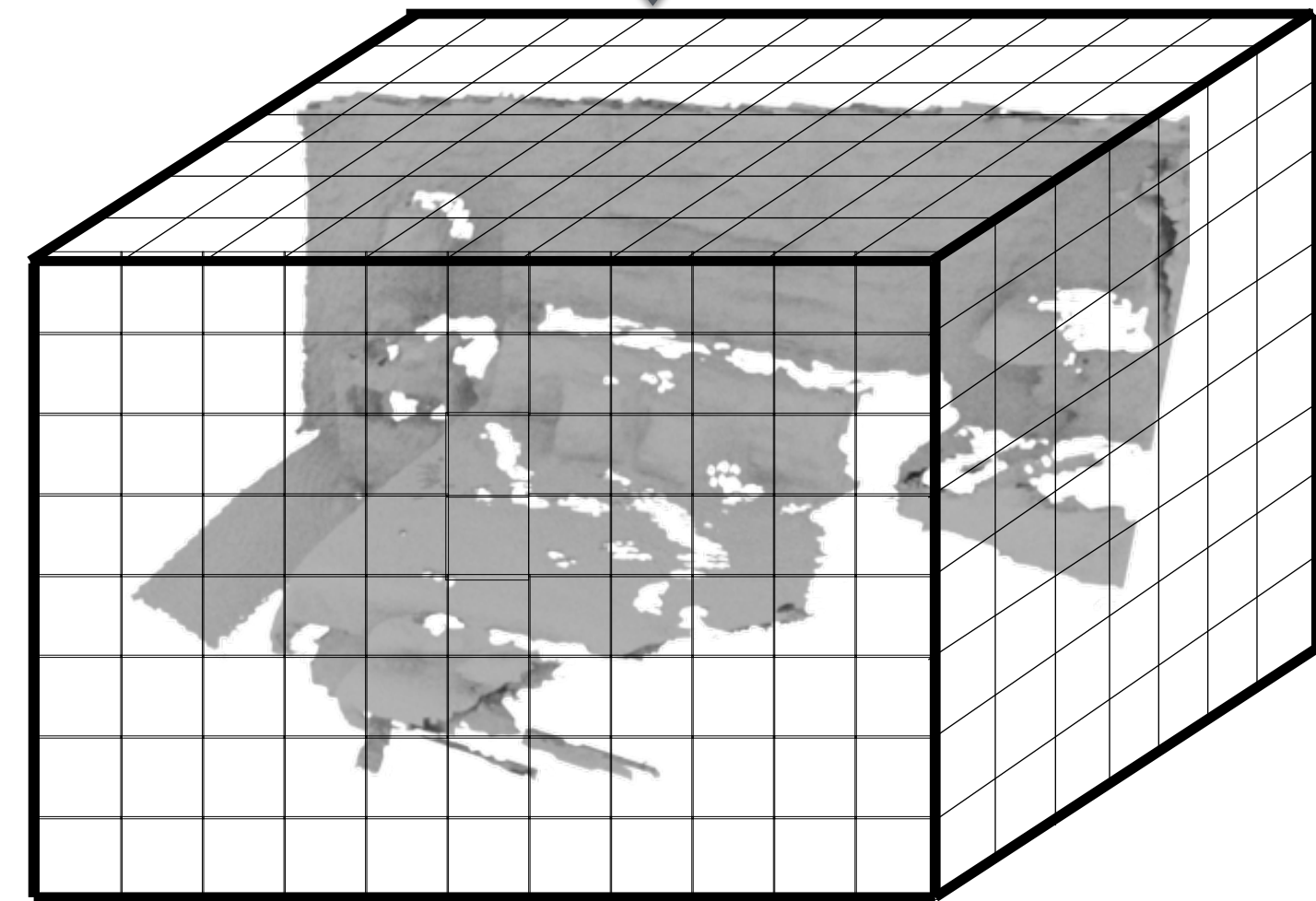
Free space

Semantic Scene Completion Network

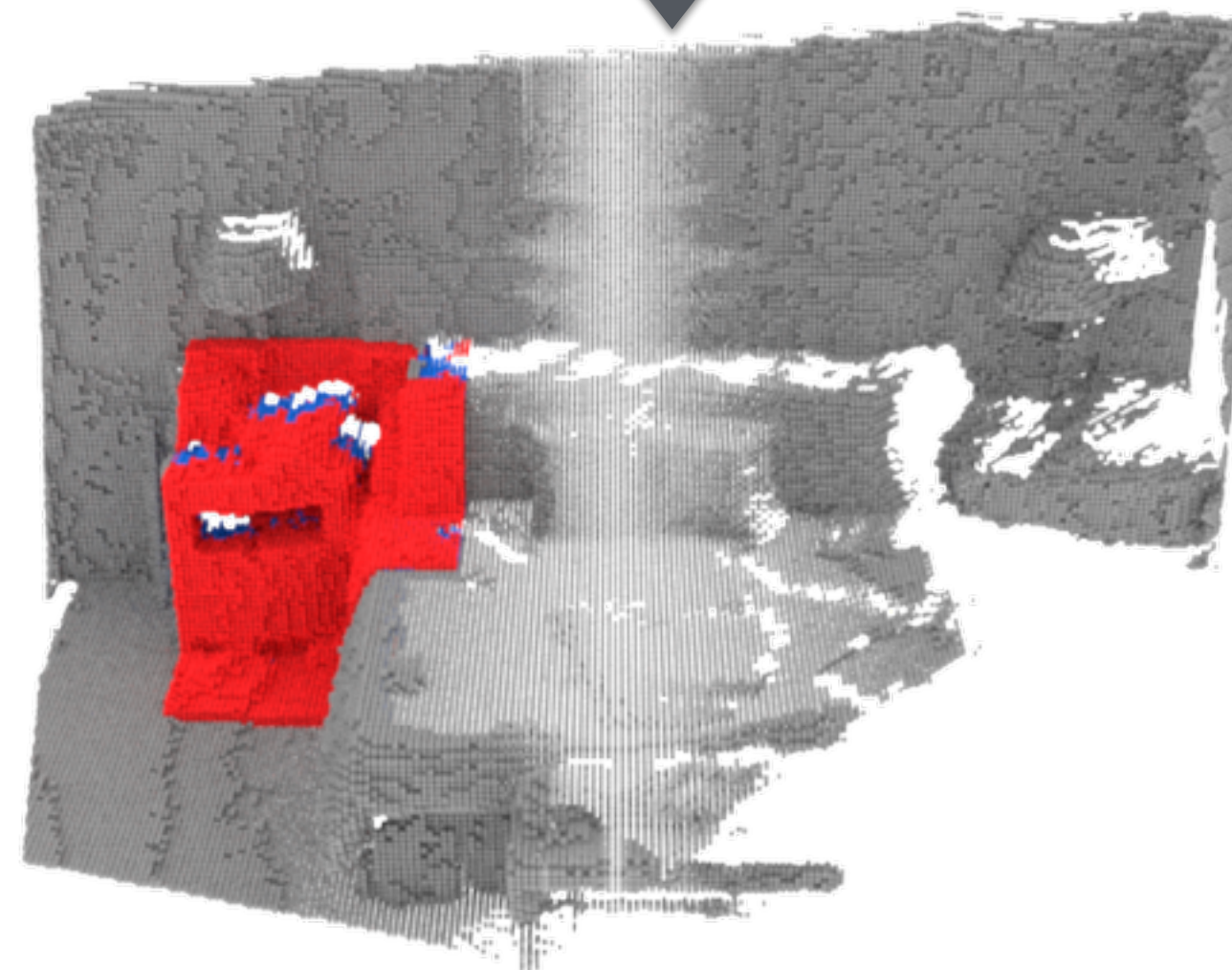
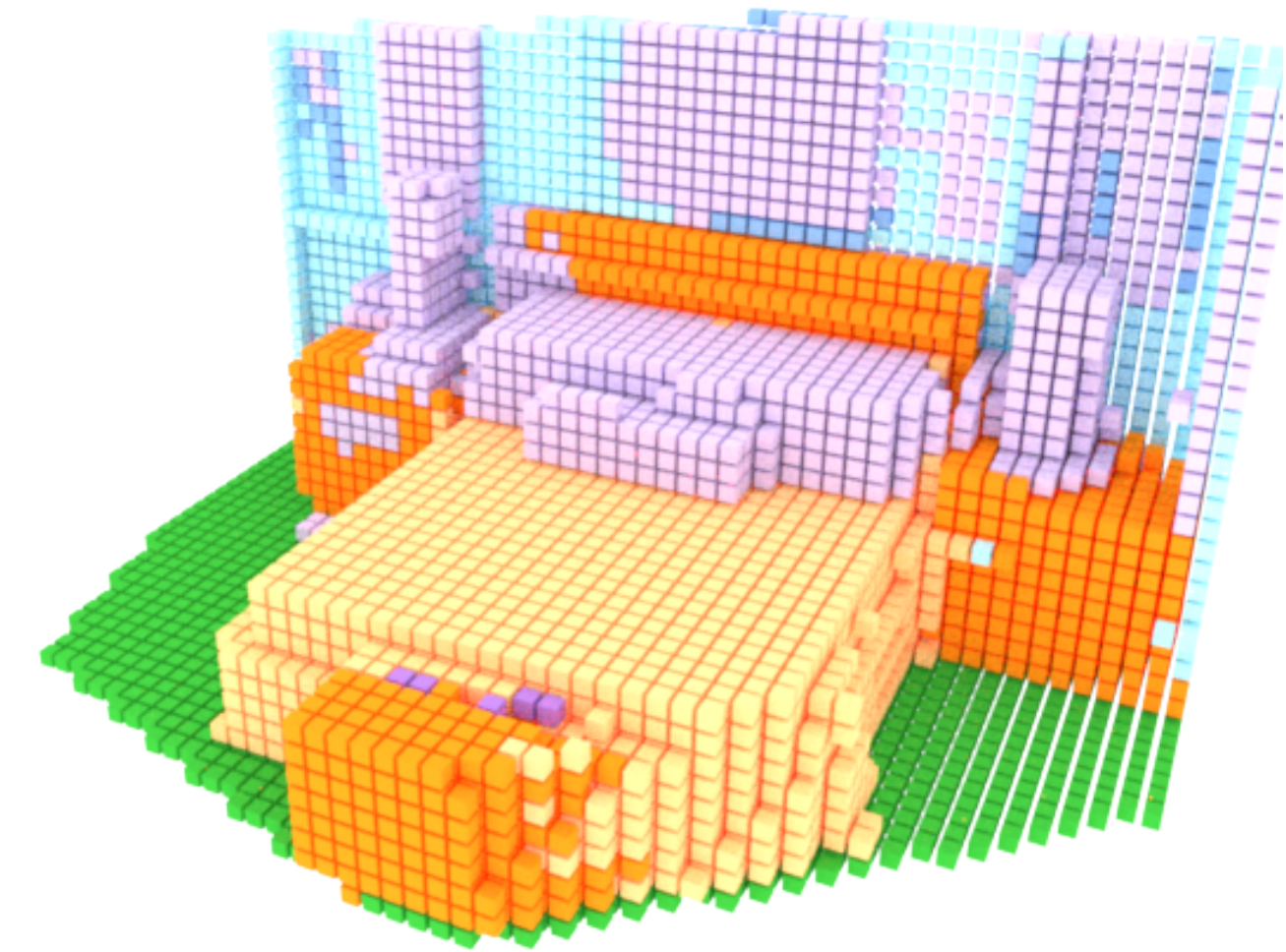
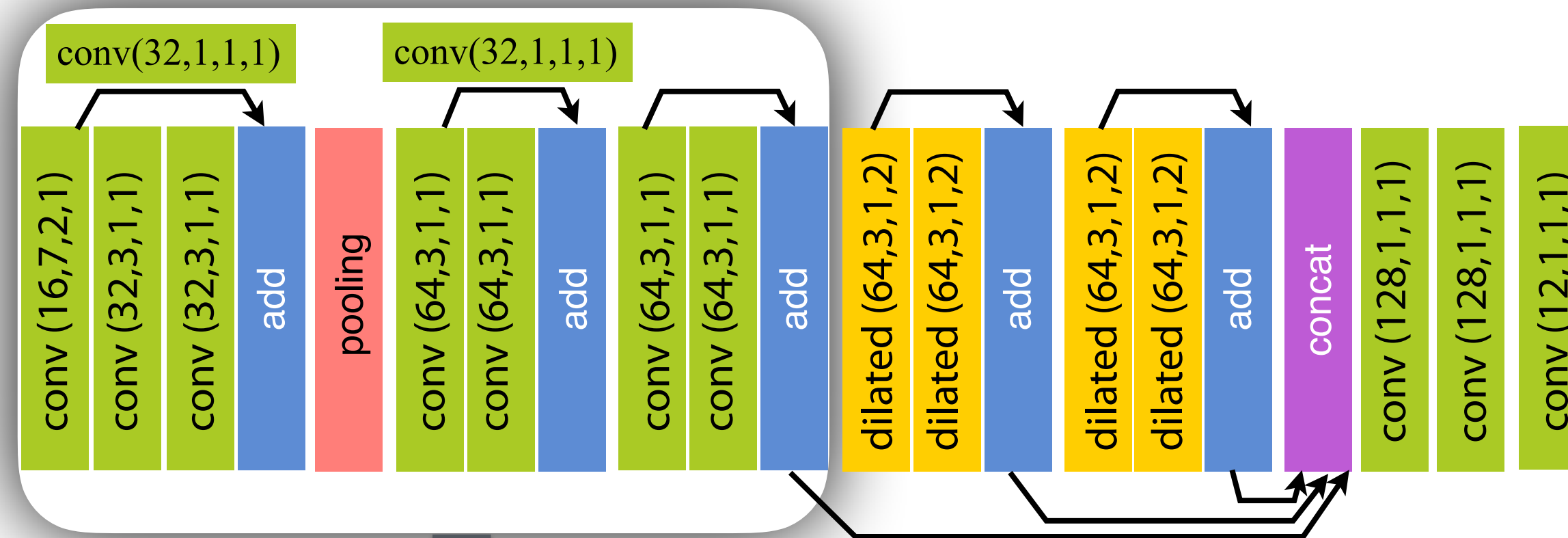
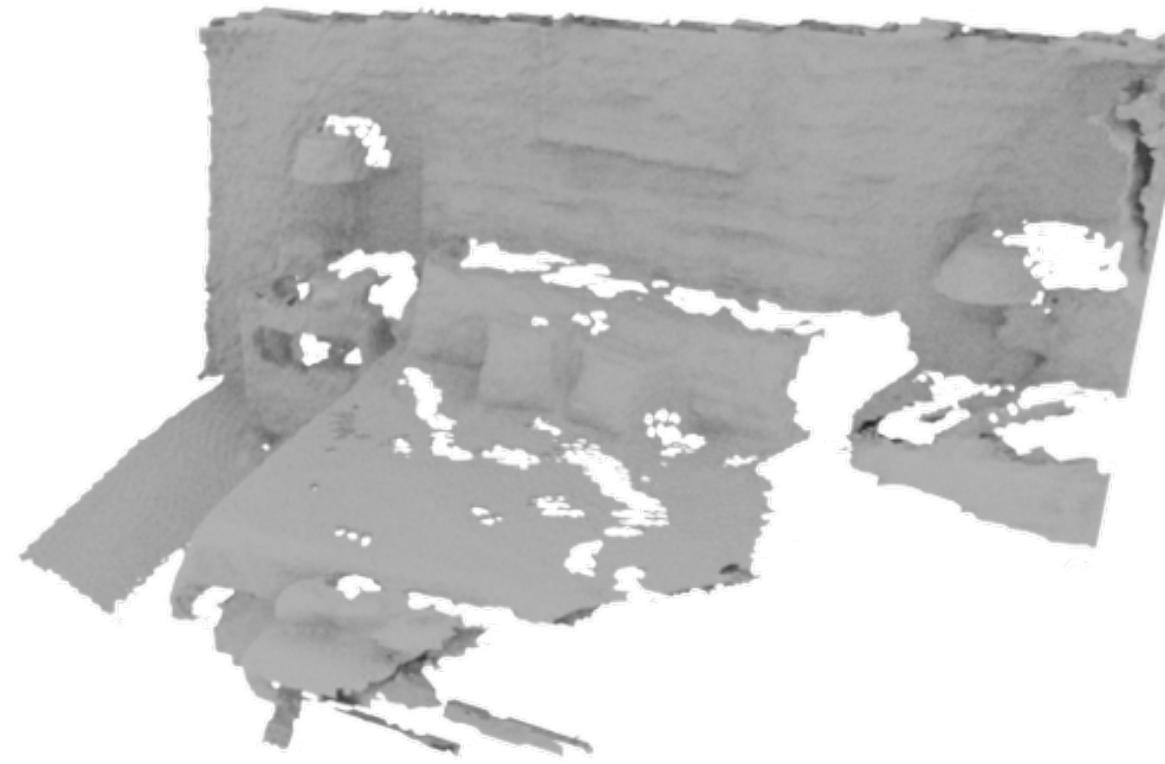


Compare to standard projective TSDF, flipTSDF:

- has less viewpoint dependency
- concentrates the strongest gradient near surface



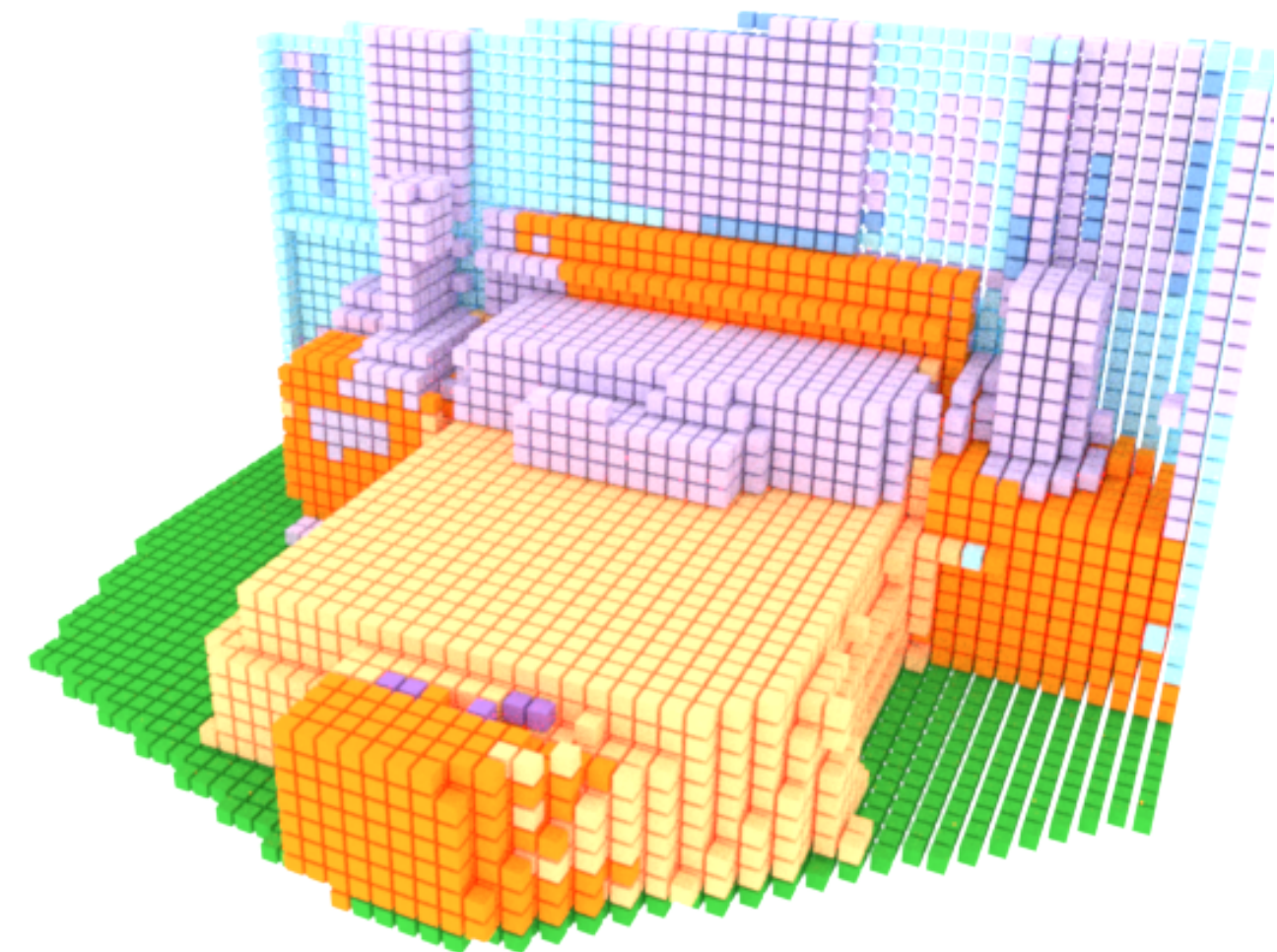
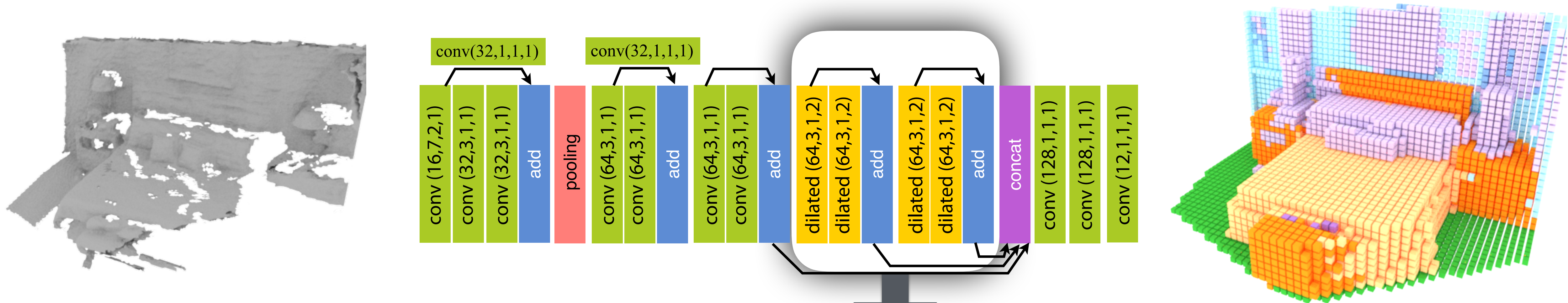
Semantic Scene Completion Network



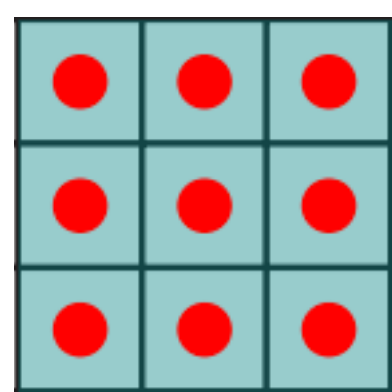
Local geometry

Receptive field: 0.98 m

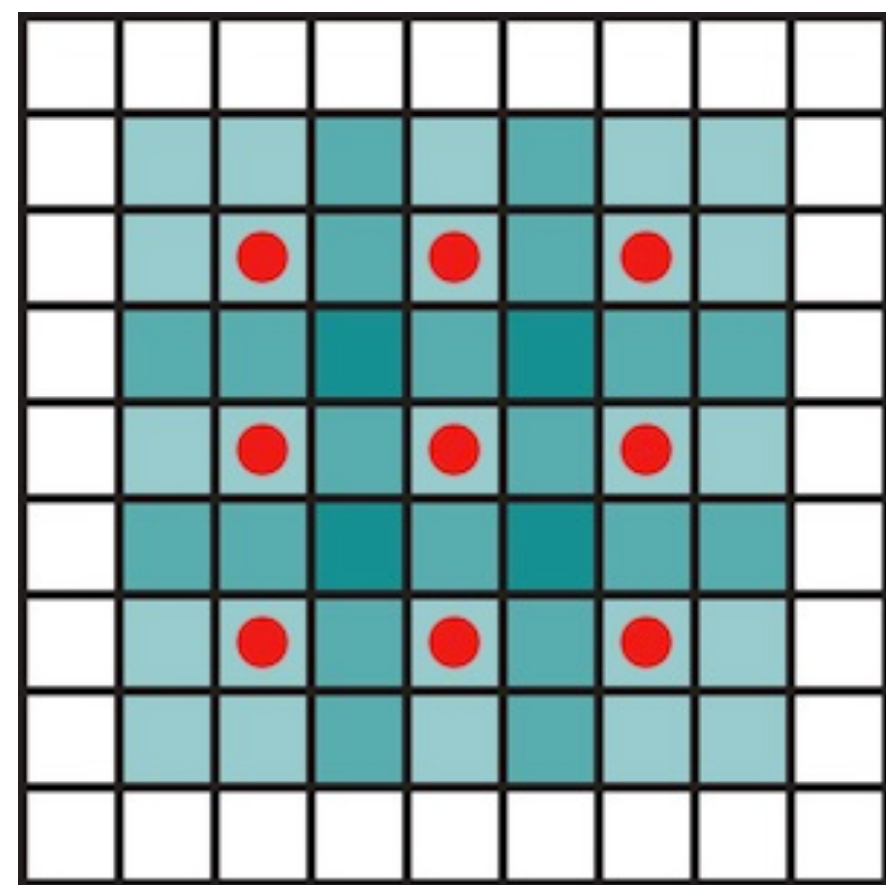
Semantic Scene Completion Network



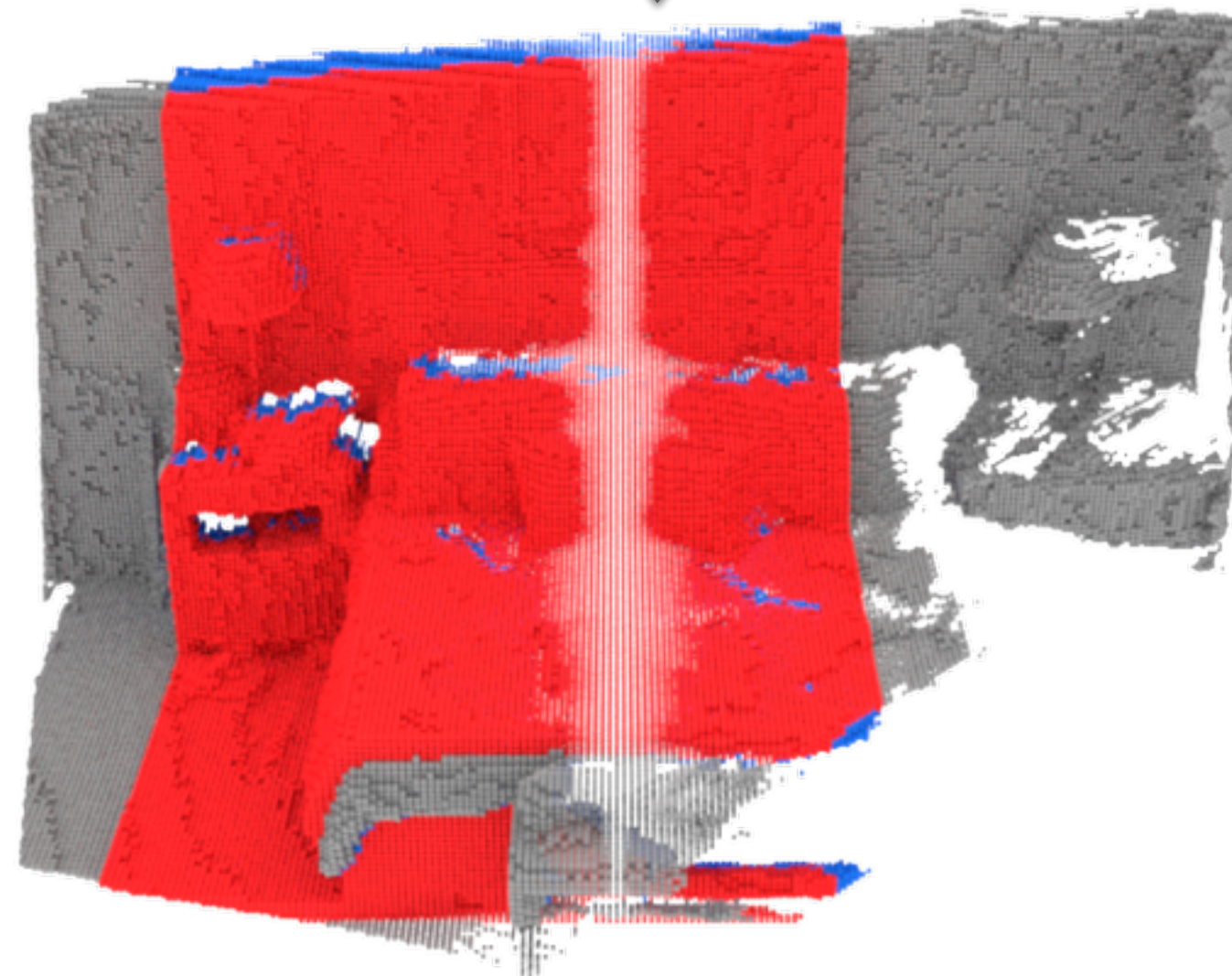
Capturing higher-level
3D context
by big receptive field



Normal kernel

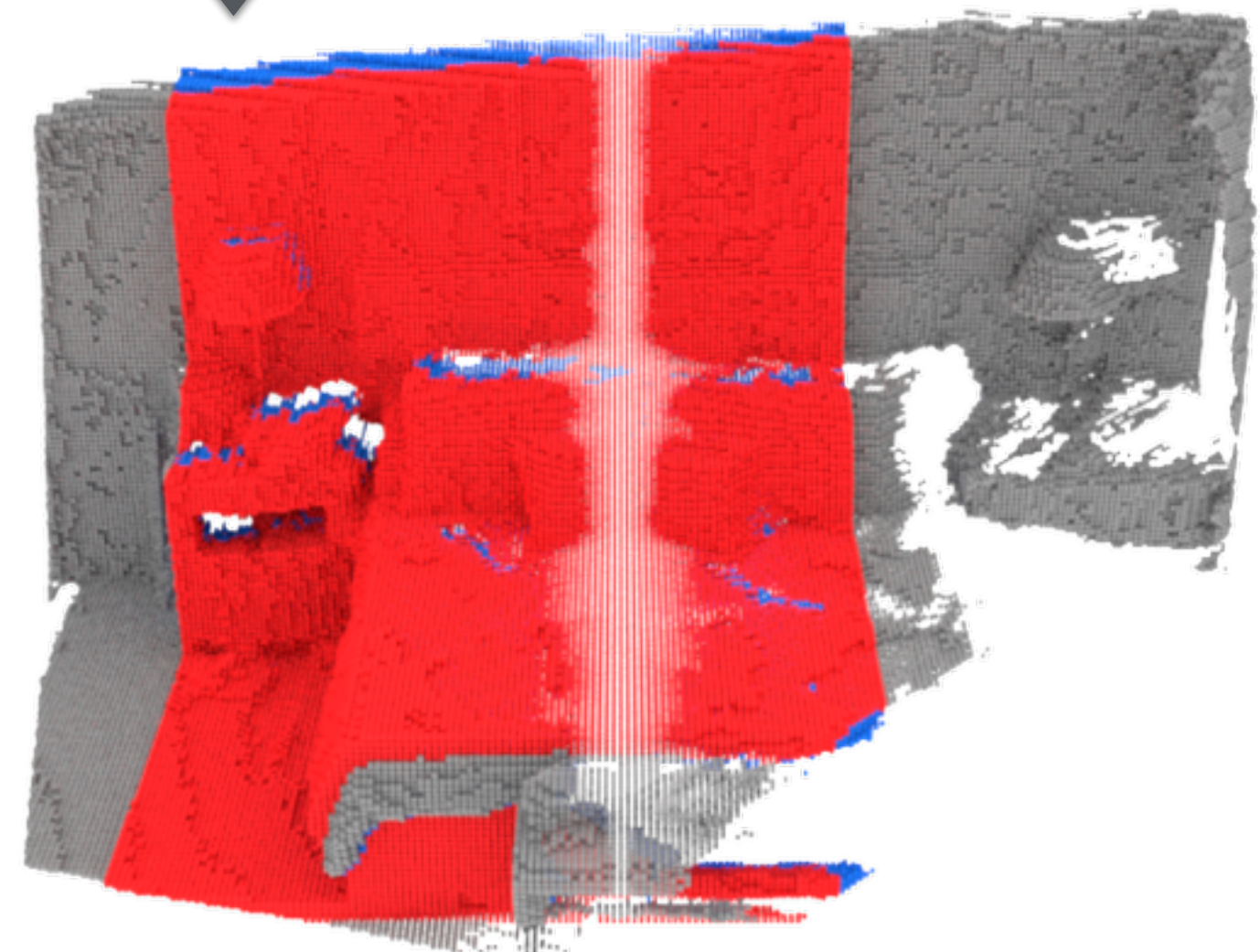
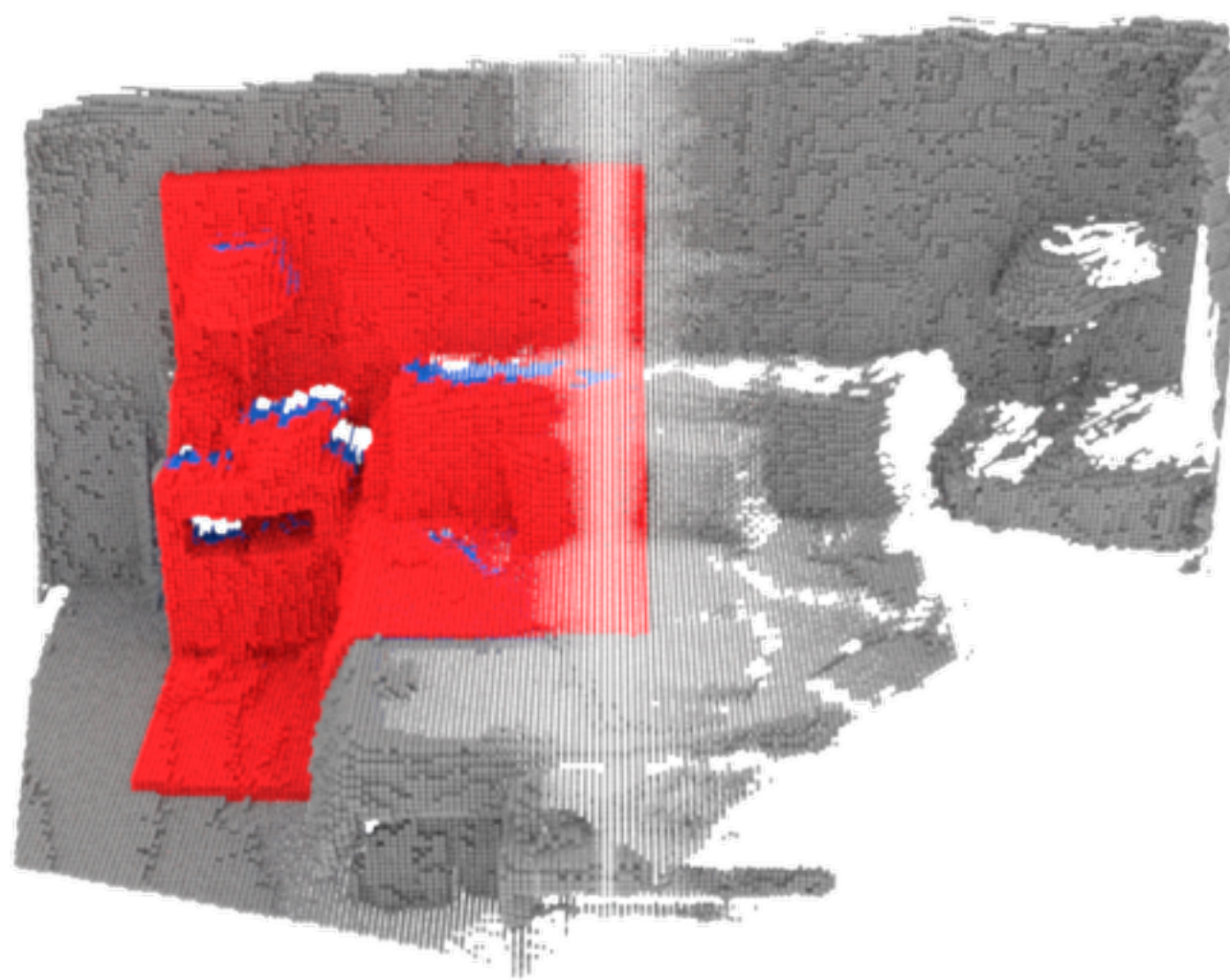
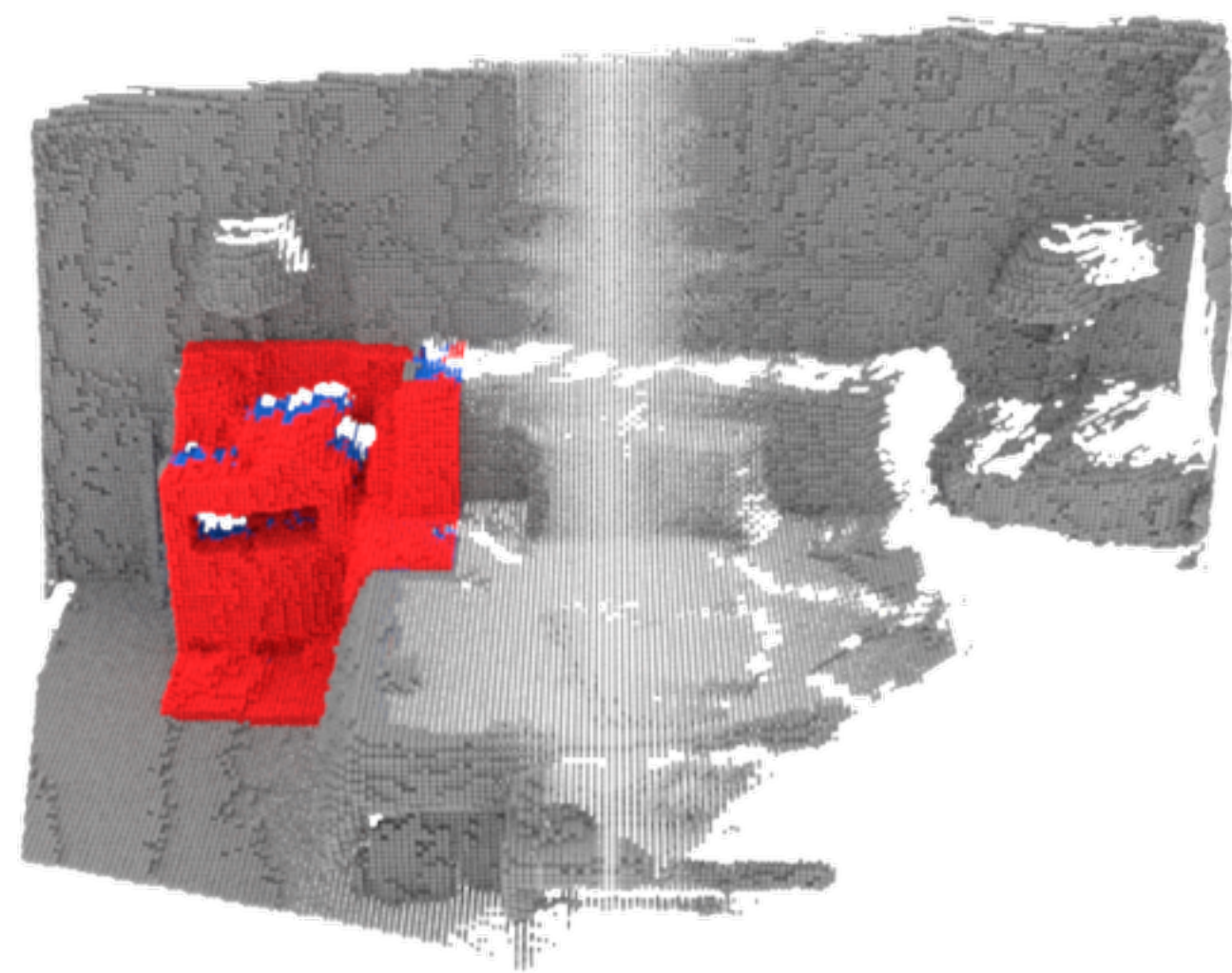
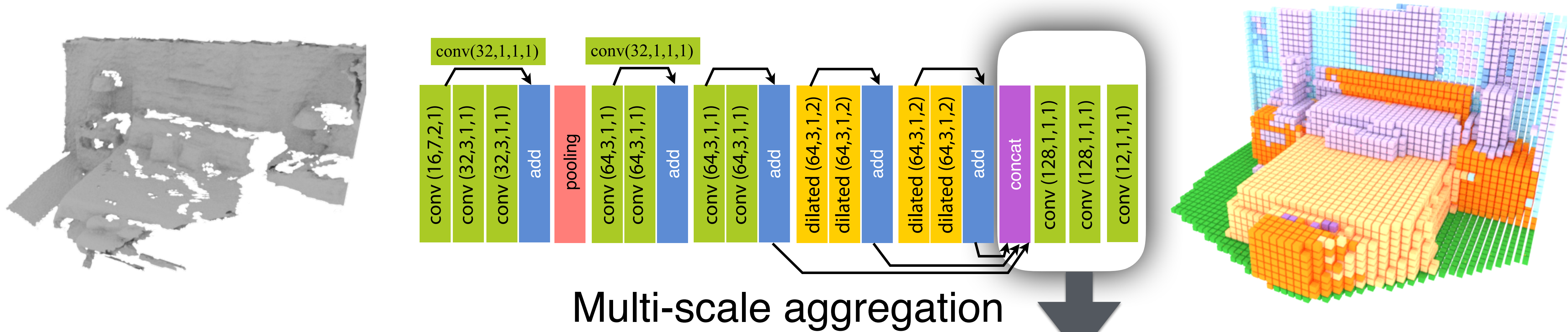


Dilation kernel

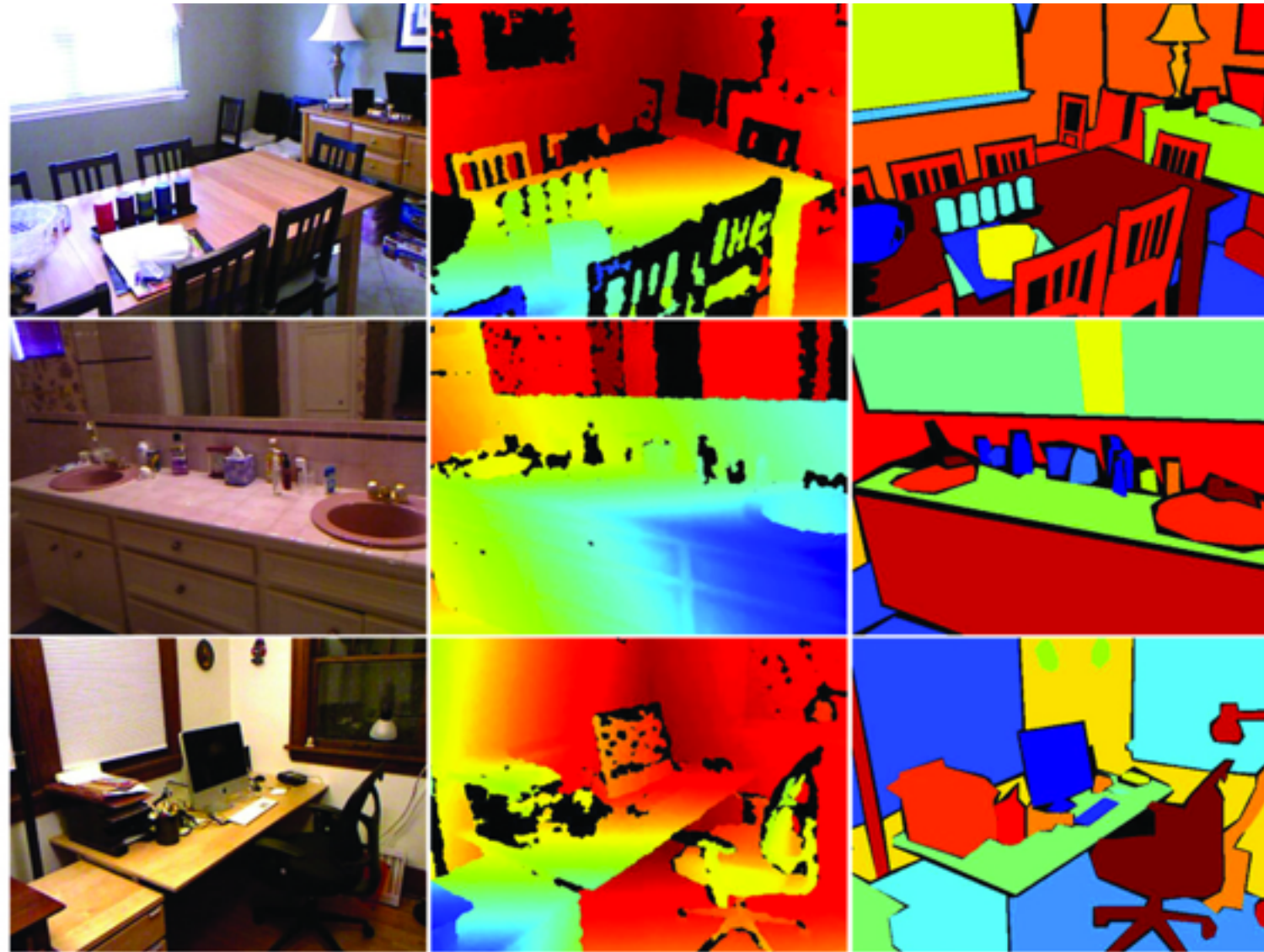


Receptive field: 2.26

Semantic Scene Completion Network

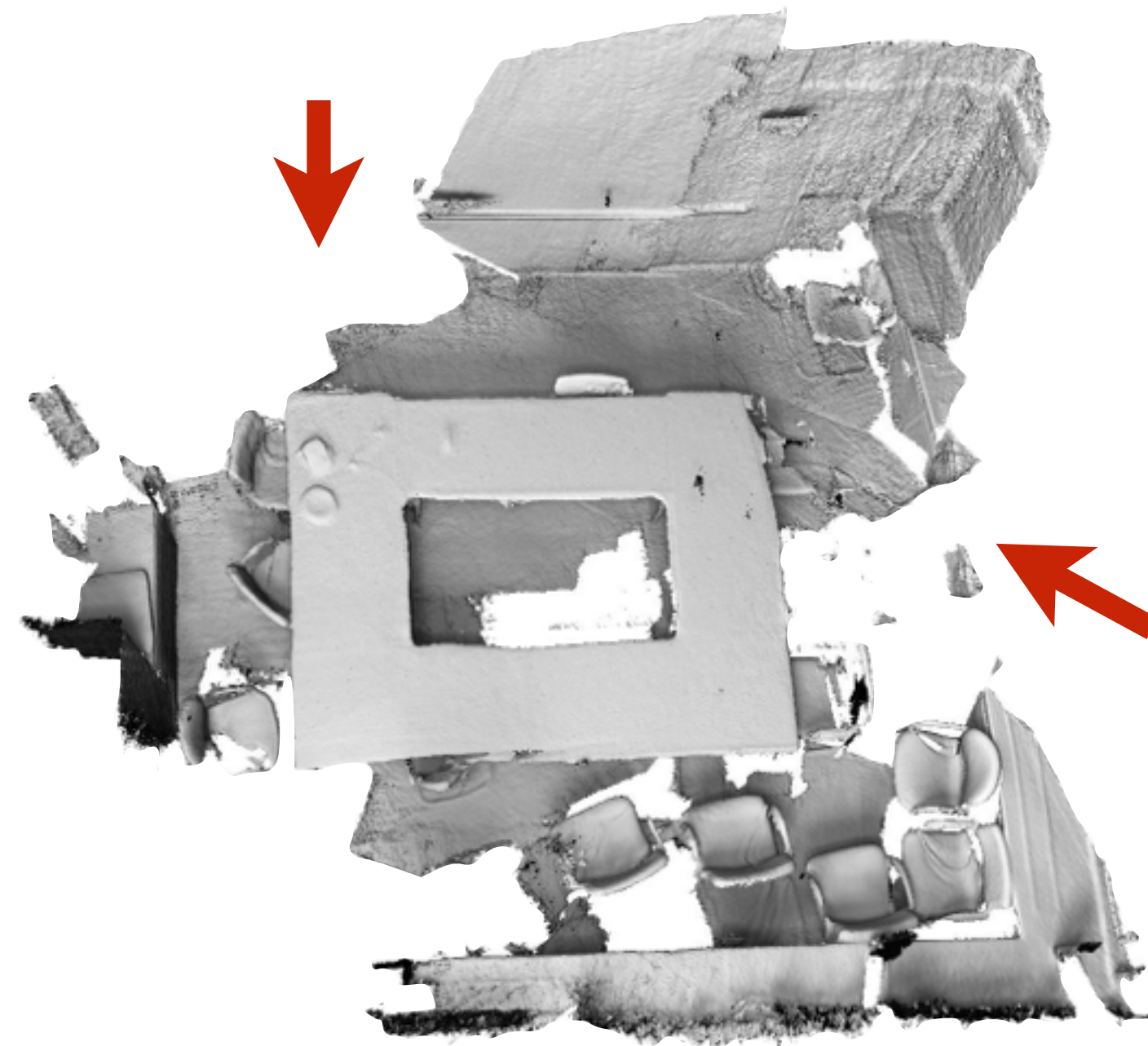


How do we obtain training data ?



Only label visible surface

[Silberman et al.]



Partial observation

[xiao et al.]



Simple scenario

[Firman et al.]

No dense volumetric ground truth with semantic labels for the complete scene.

SUNCG dataset: over 40K houses

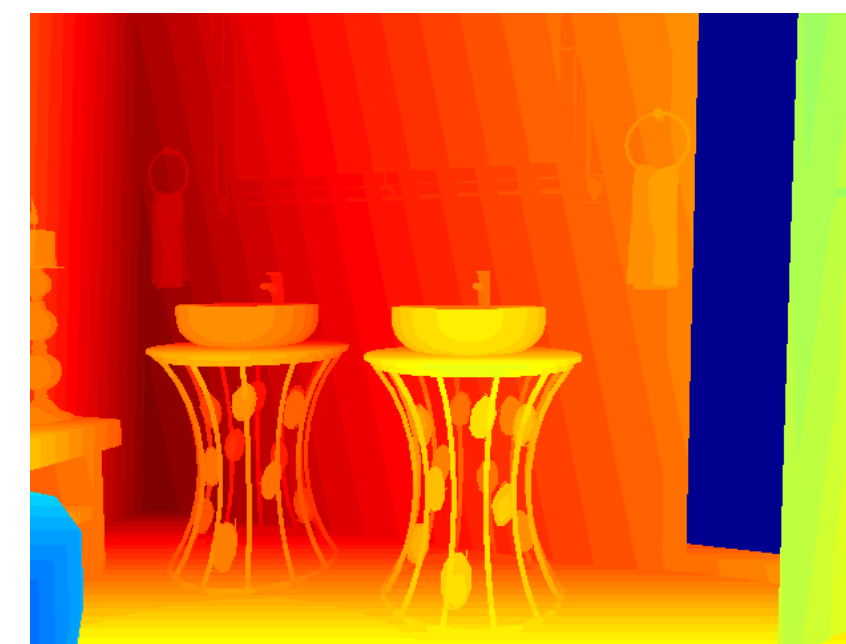
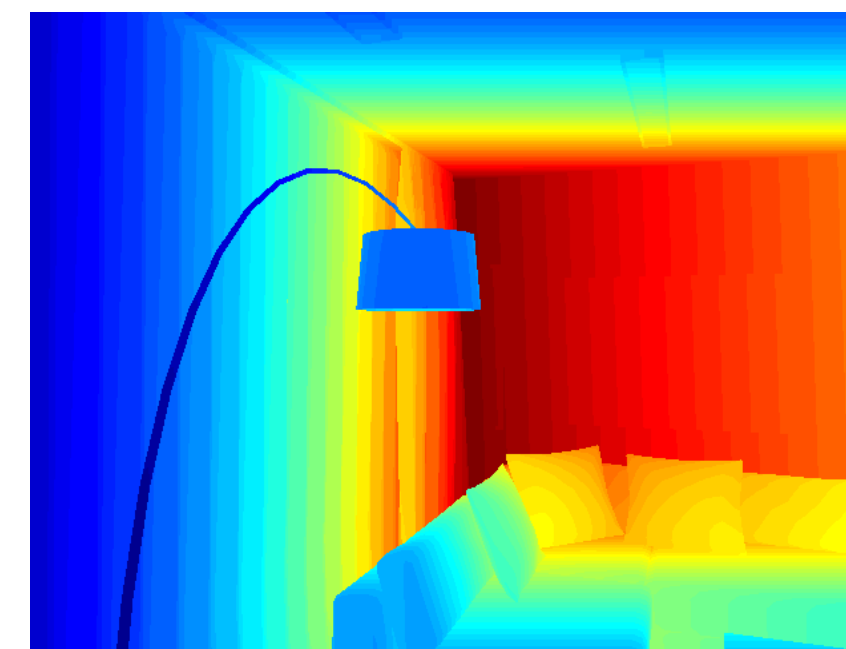
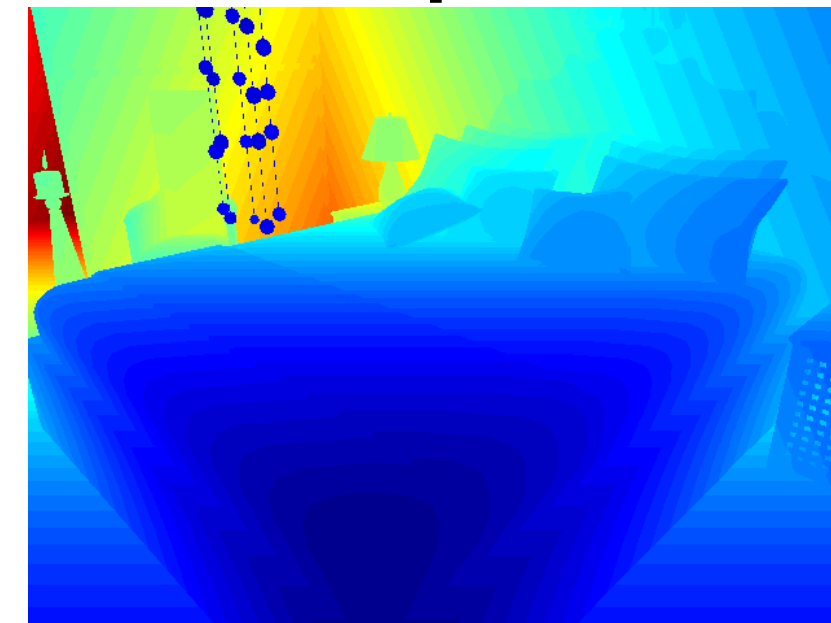


Synthesizing training data

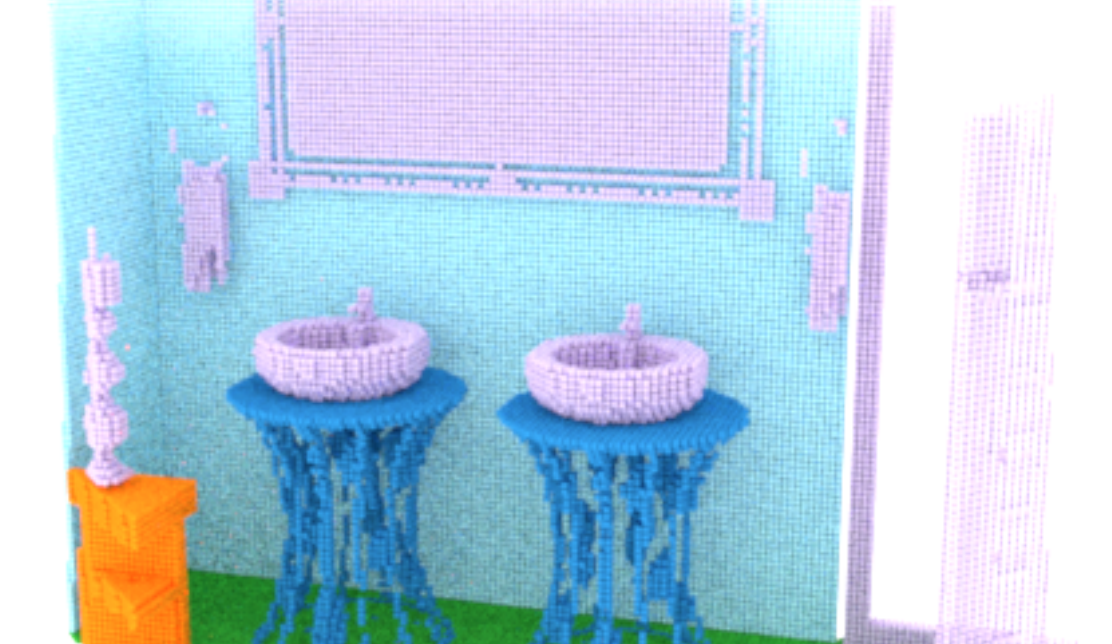
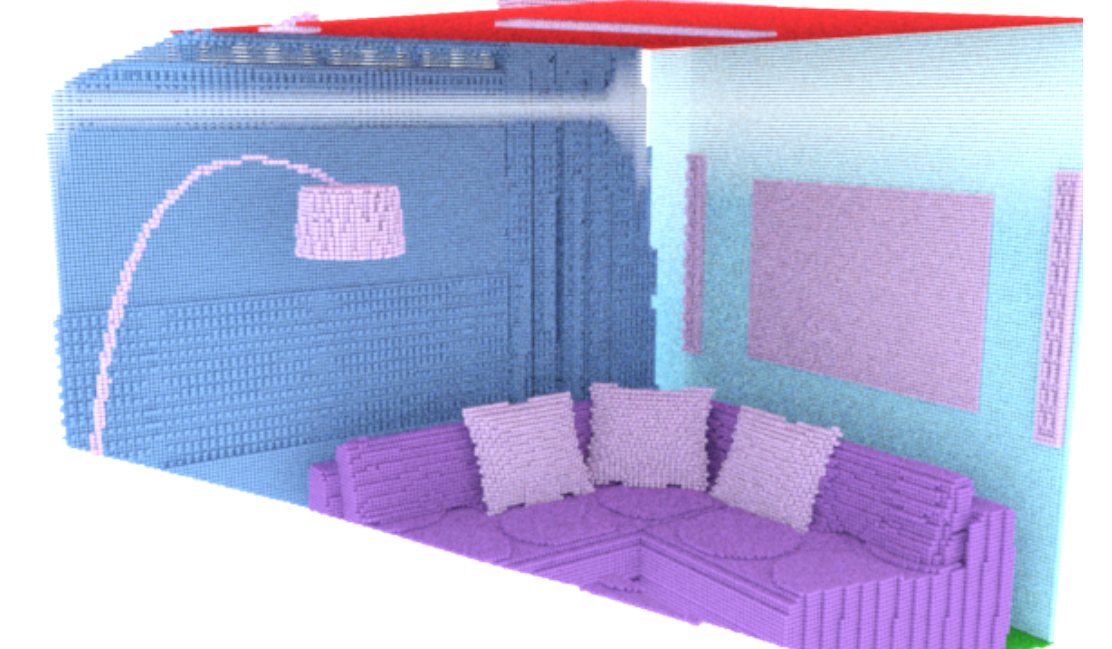
one floor



depth



ground truth



Testing on real-world data

Training on SUNCG



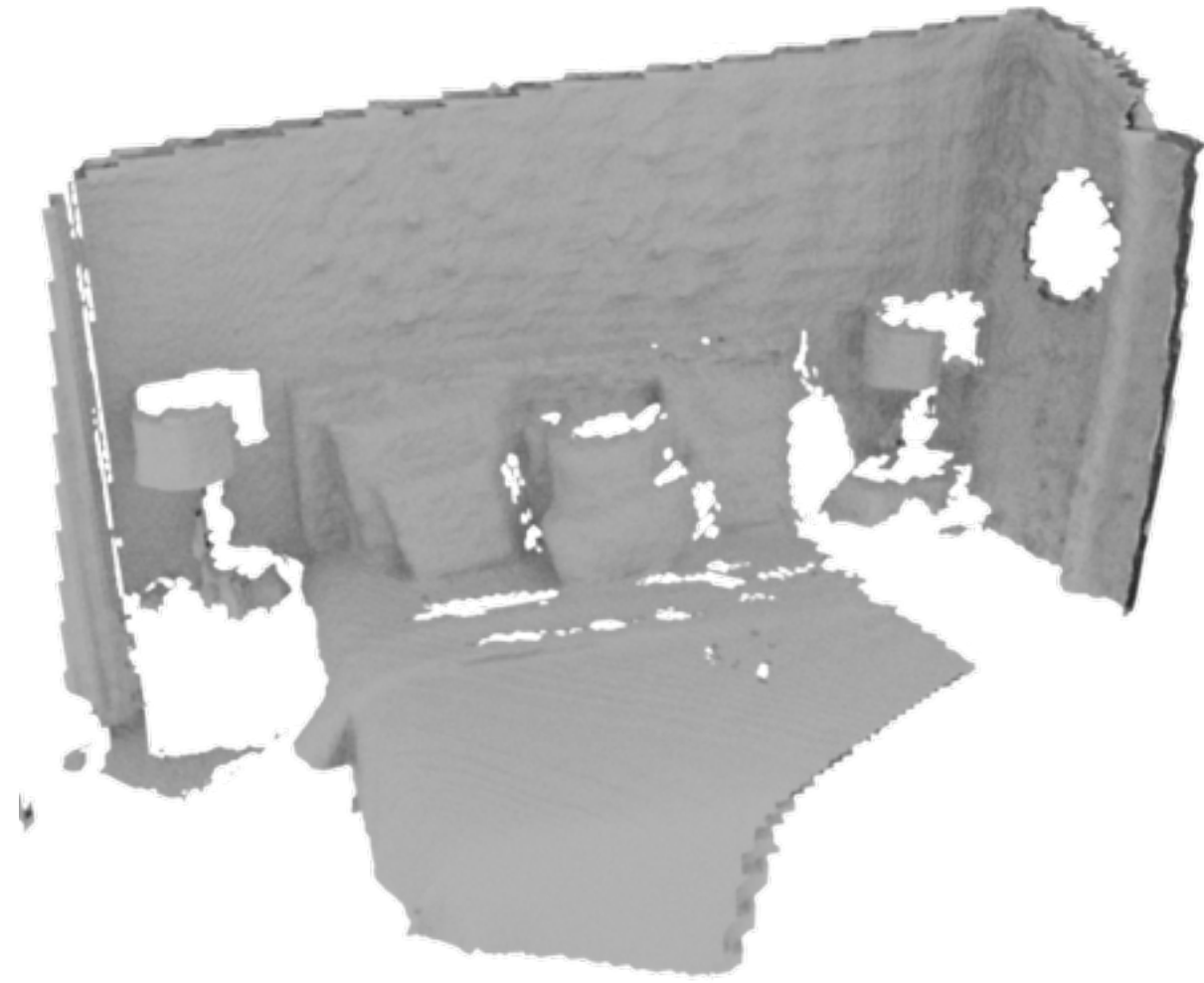
Testing on NYU [1,2]



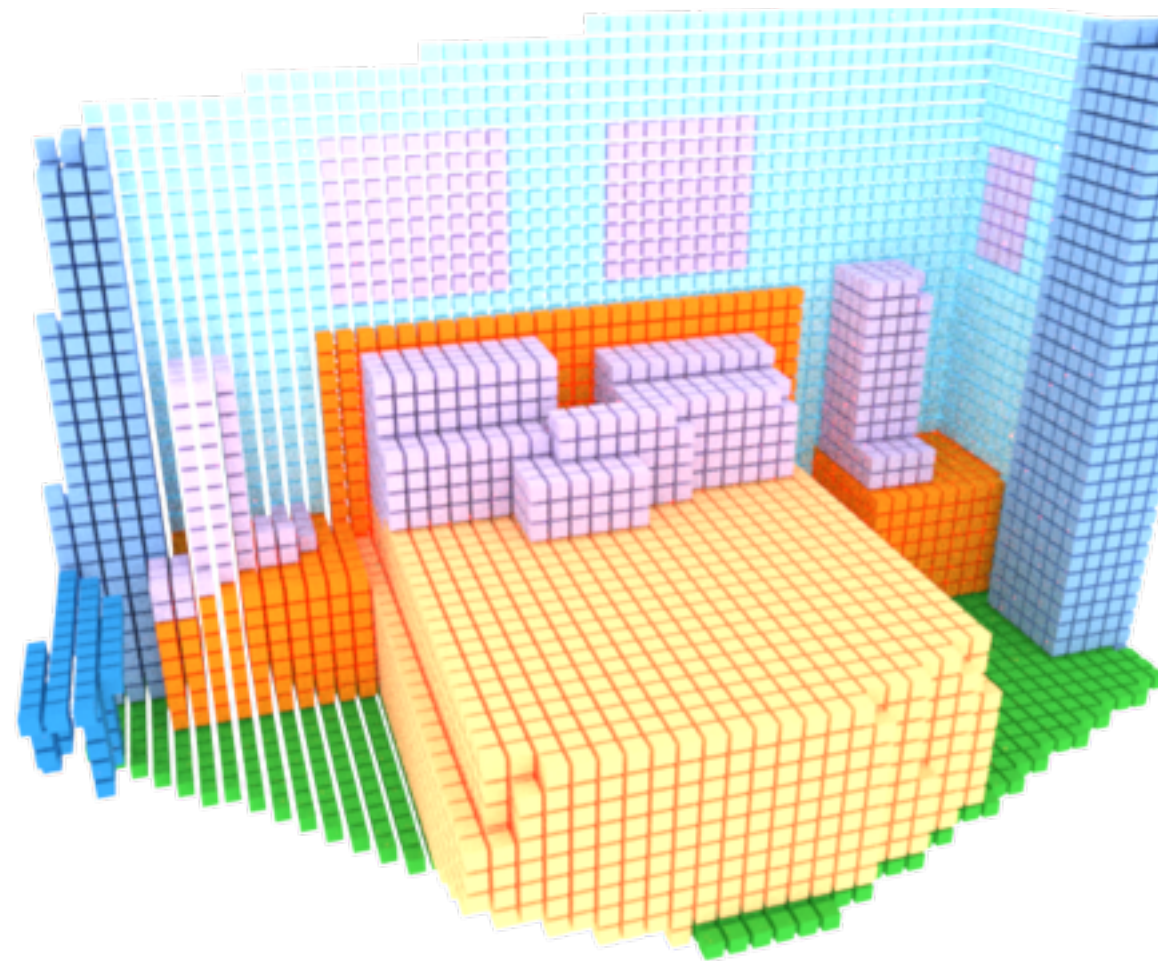
[1] **NYU depth v2**: Nathan Silberman, Pushmeet Kohli, Derek Hoiem, Rob Fergus. Indoor Segmentation and Support Inference from RGBD Images. ECCV 2012

[2] **Ground truth**: Ruiqi Guo, Derek Hoiem. Support surface prediction in indoor scenes. ICCV 2013

Comparison

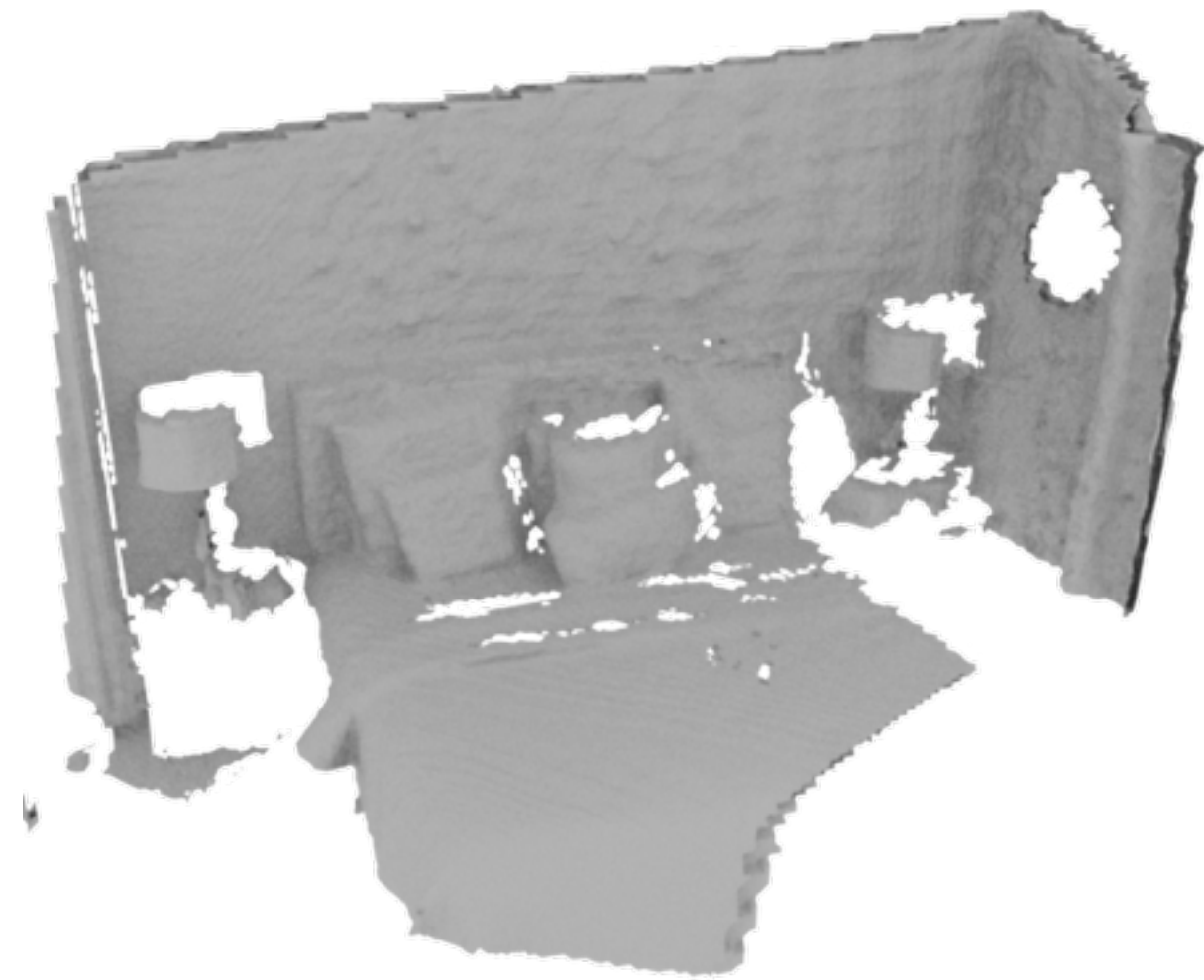


Observed Surface

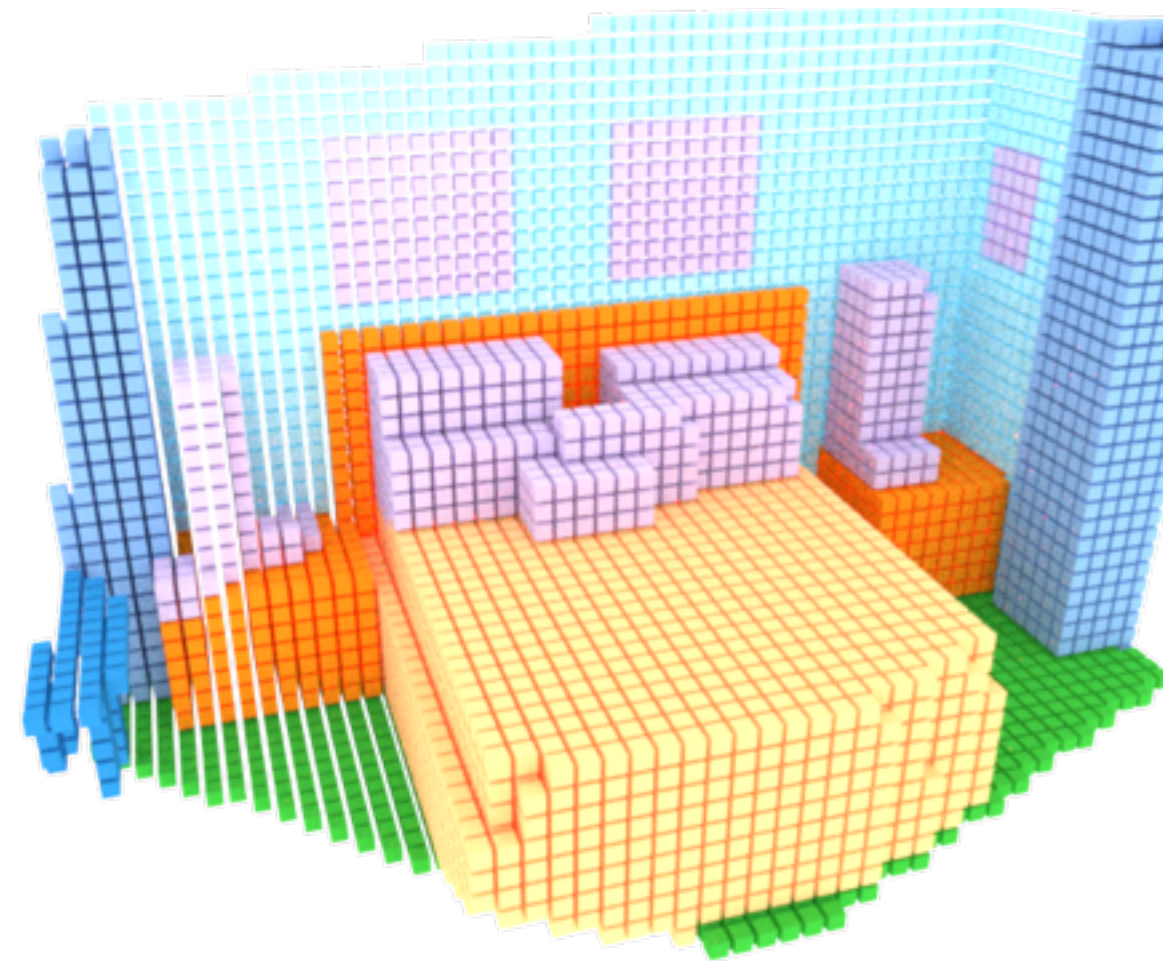


Ground Truth

Comparison

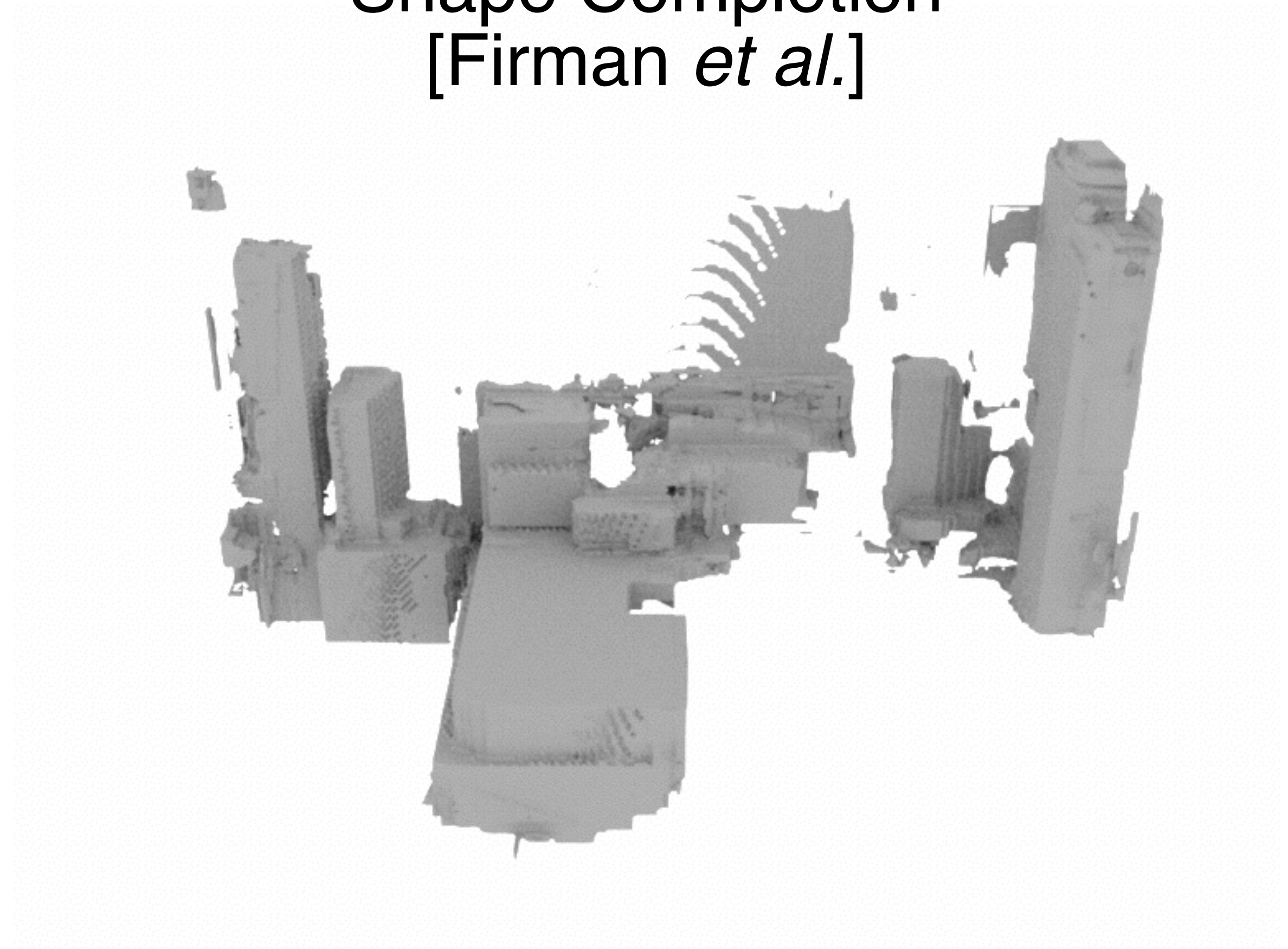


Observed Surface



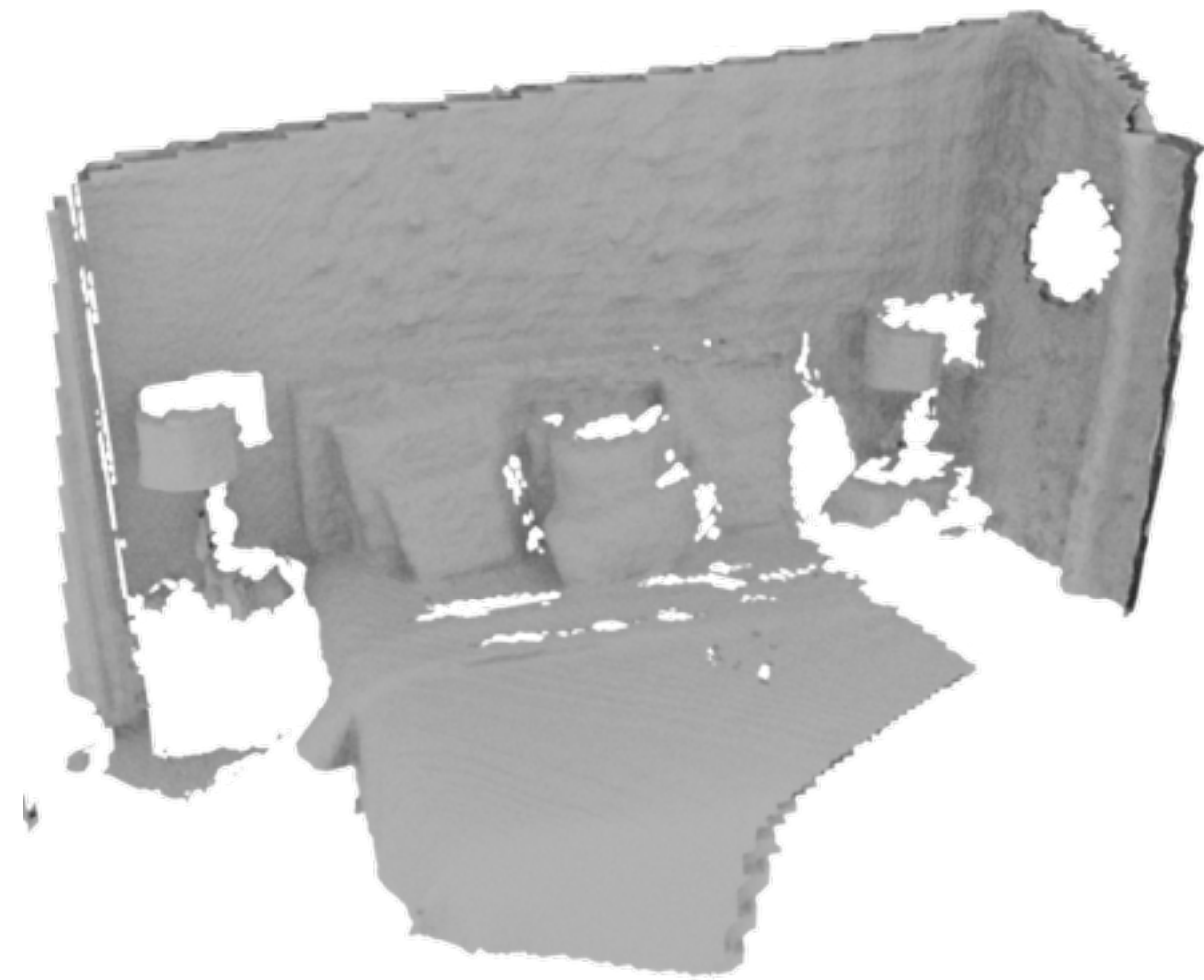
Ground Truth

Shape Completion
[Firman *et al.*]

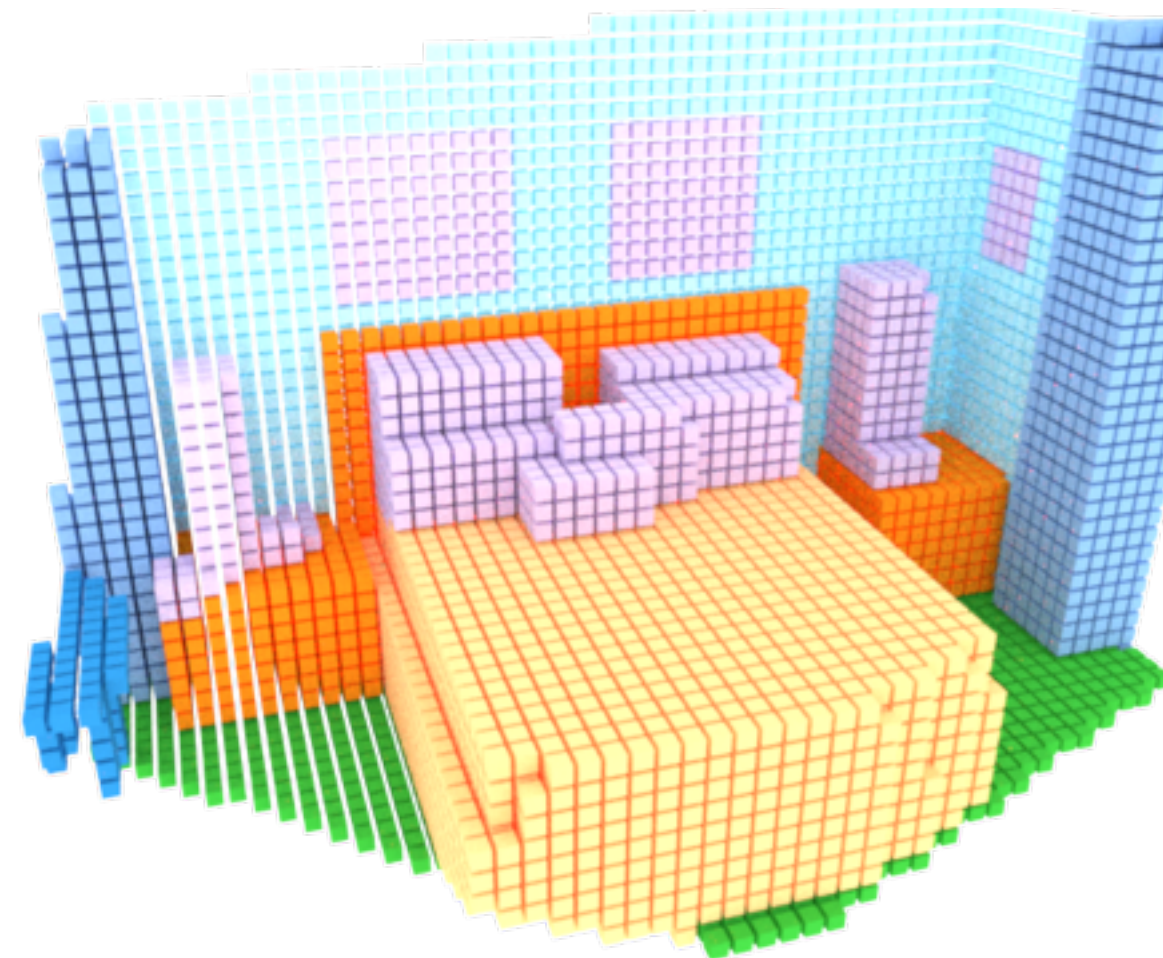


- floor
- wall
- window
- chair
- bed
- sofa
- table
- tvs
- furn.
- objects

Comparison



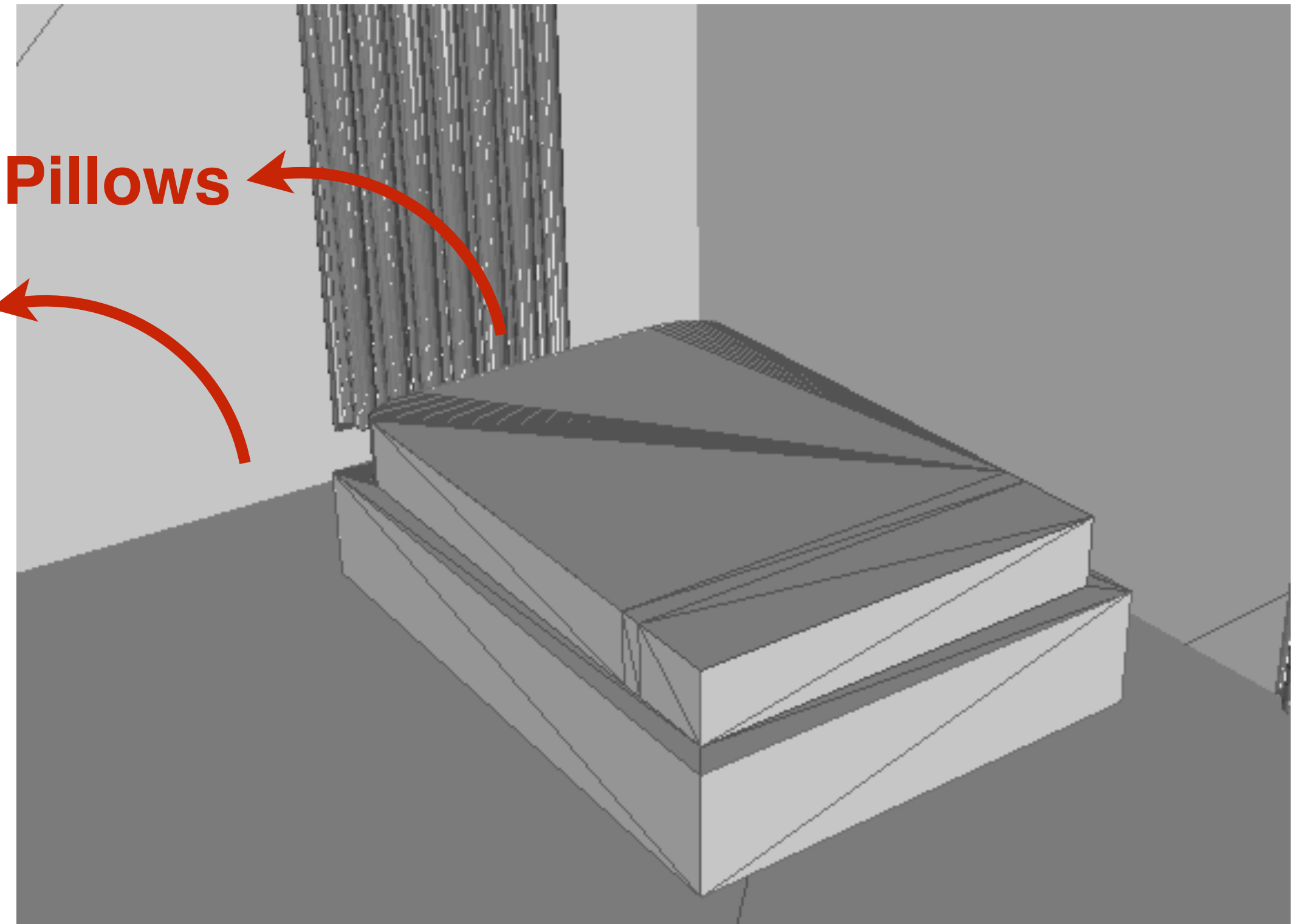
Observed Surface



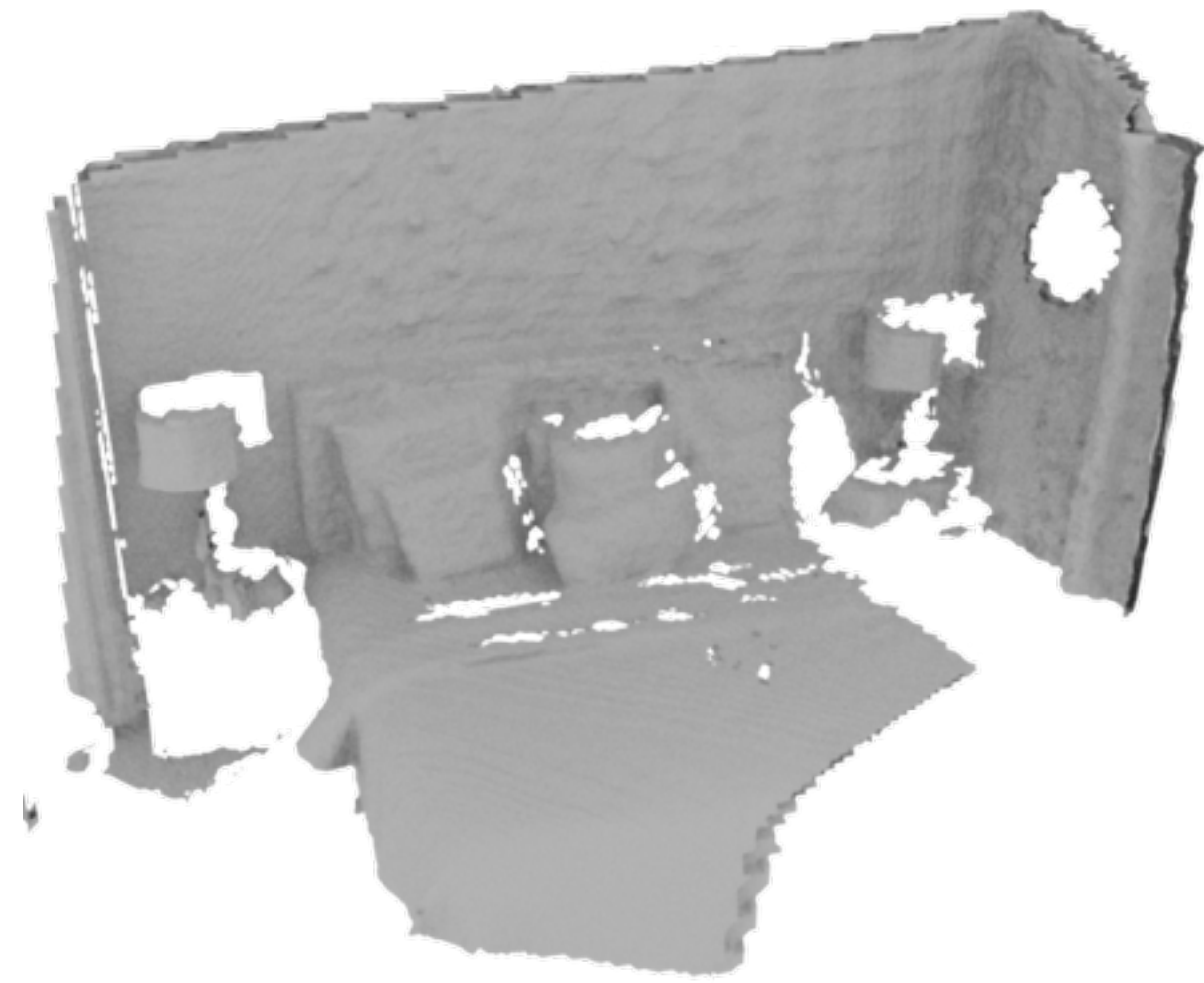
Ground Truth

Model Retrieval+Fitting
[Geiger and Wang]

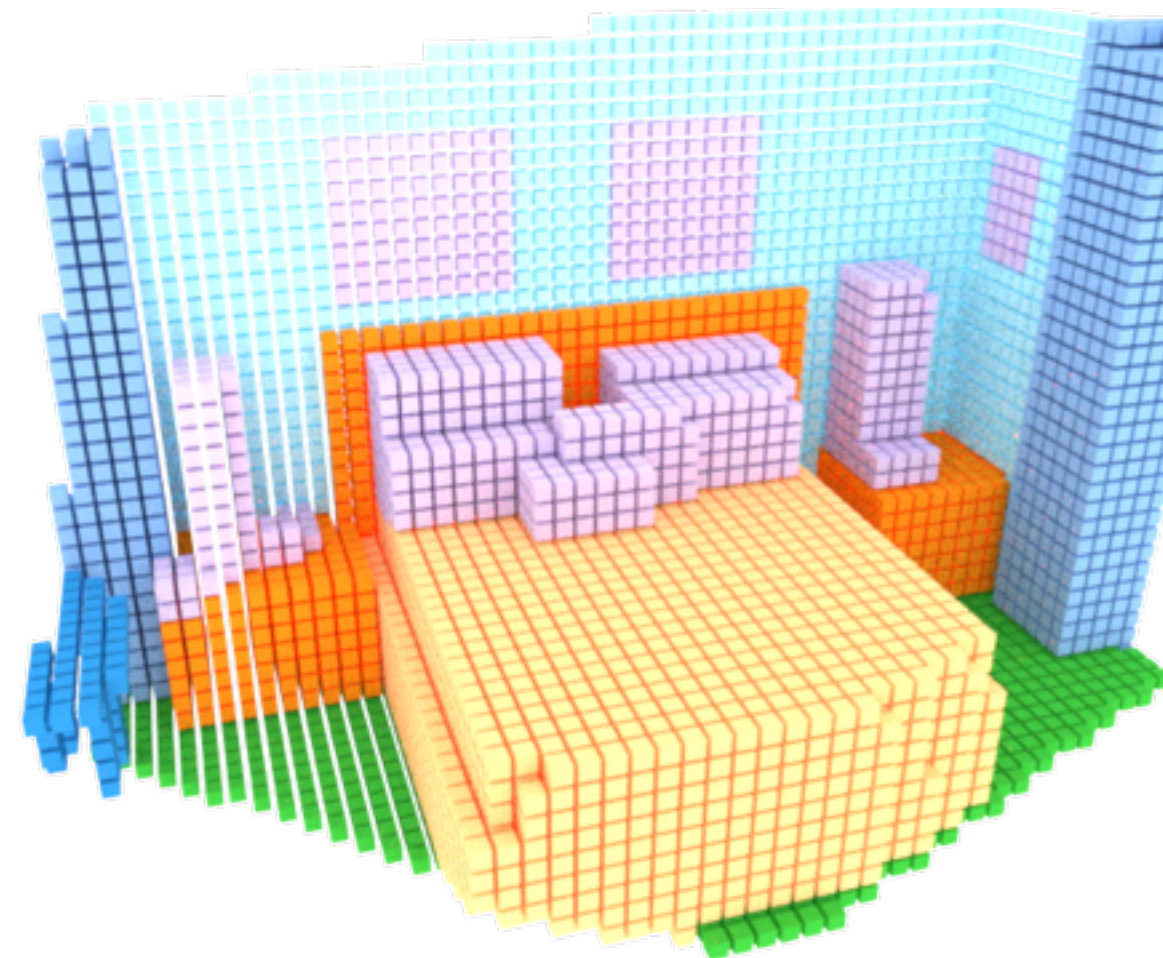
Missing Pillows
Missing Nightstand



Comparison

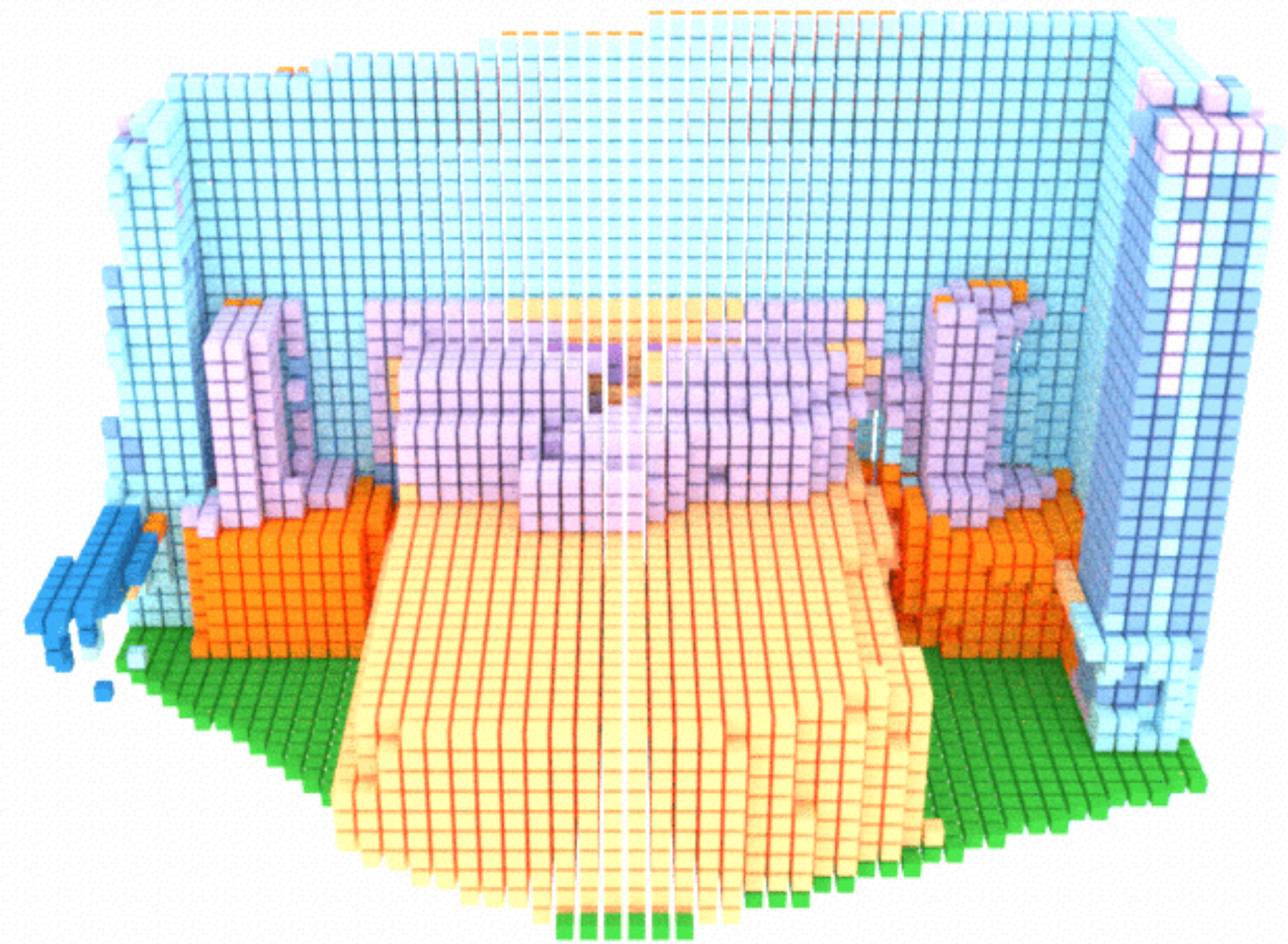


Observed Surface



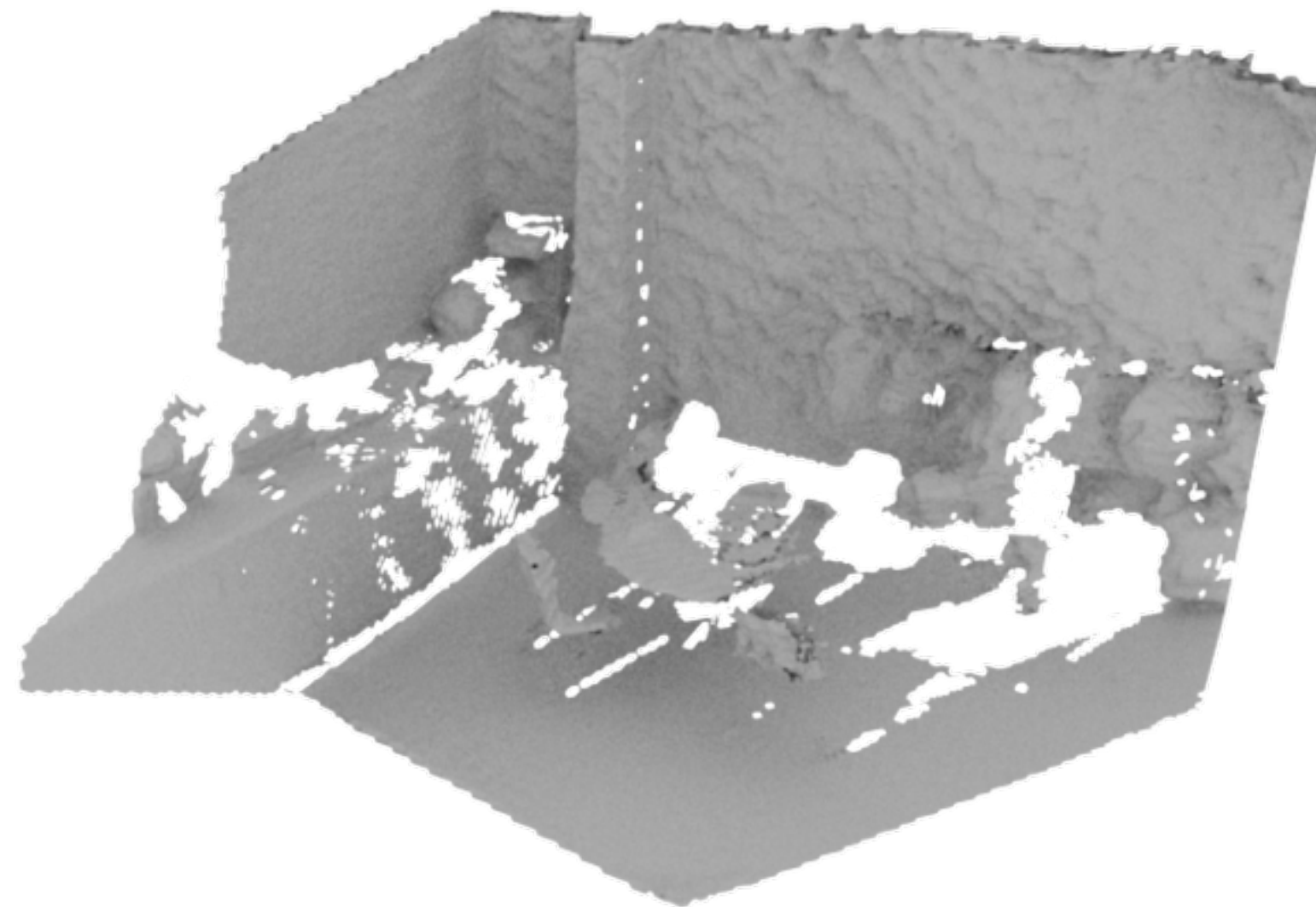
Ground Truth

SSCNet

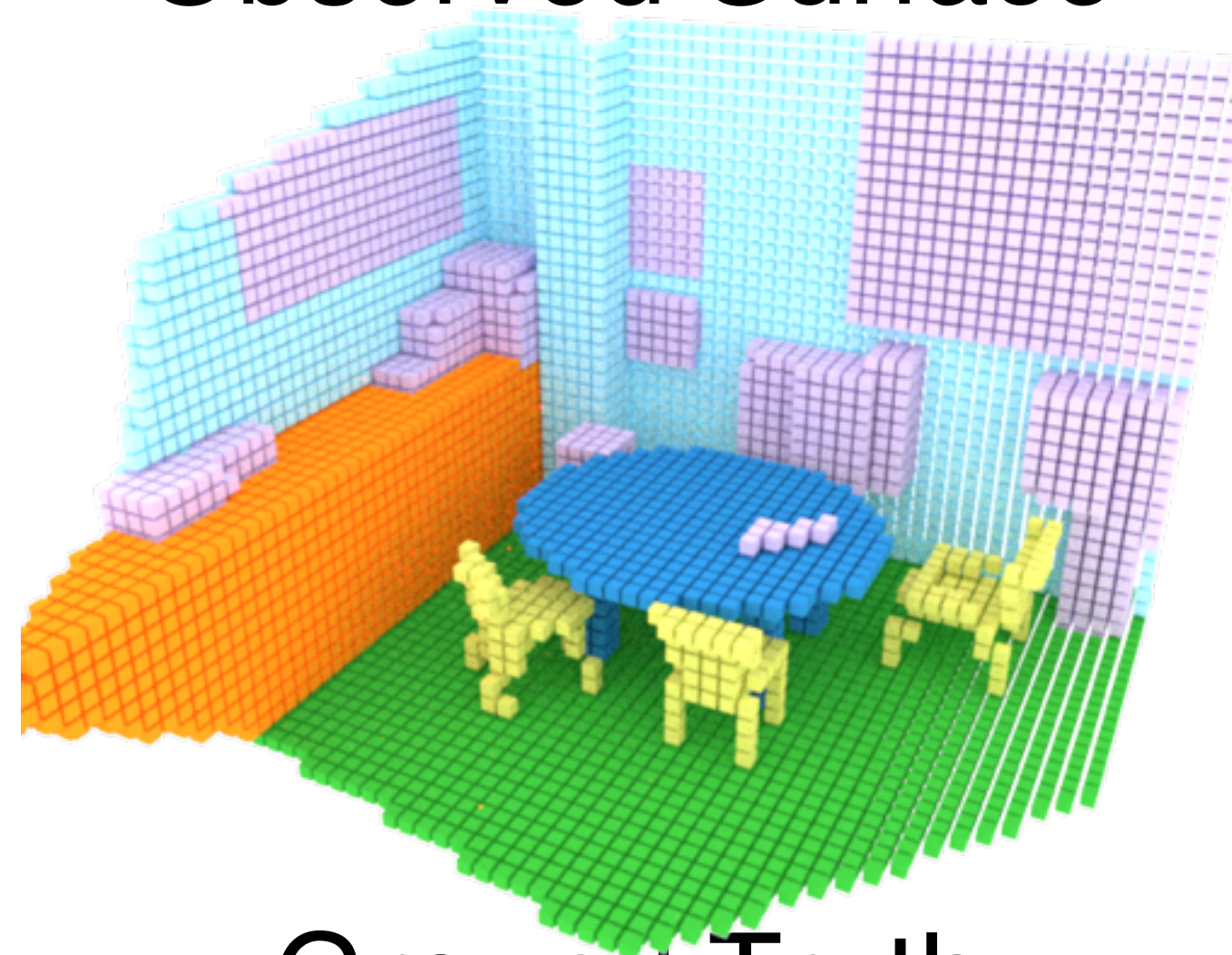


- floor
- wall
- window
- chair
- bed
- sofa
- table
- tvs
- furn.
- objects

Comparison

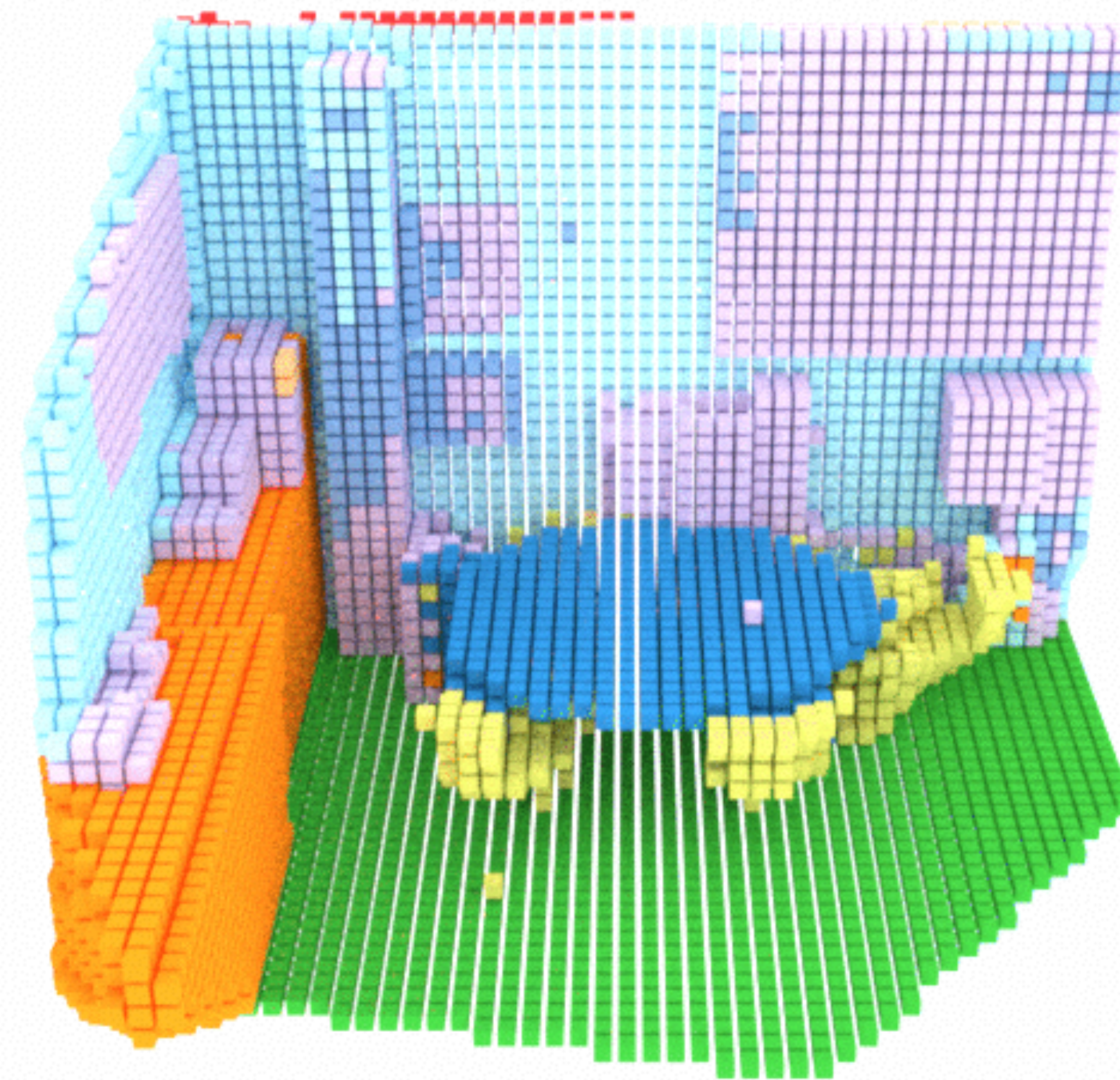


Observed Surface



Ground Truth

SSCNet

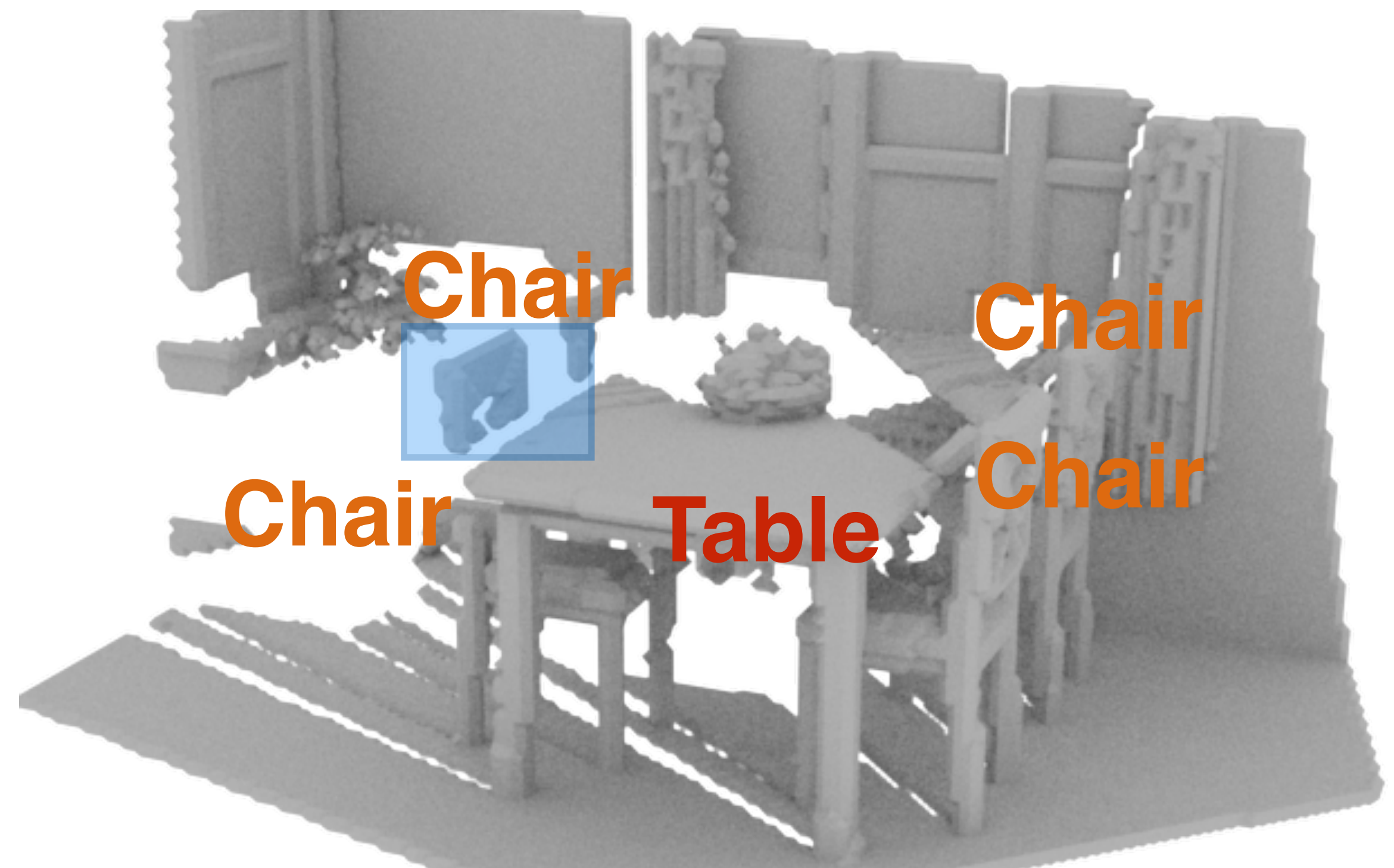
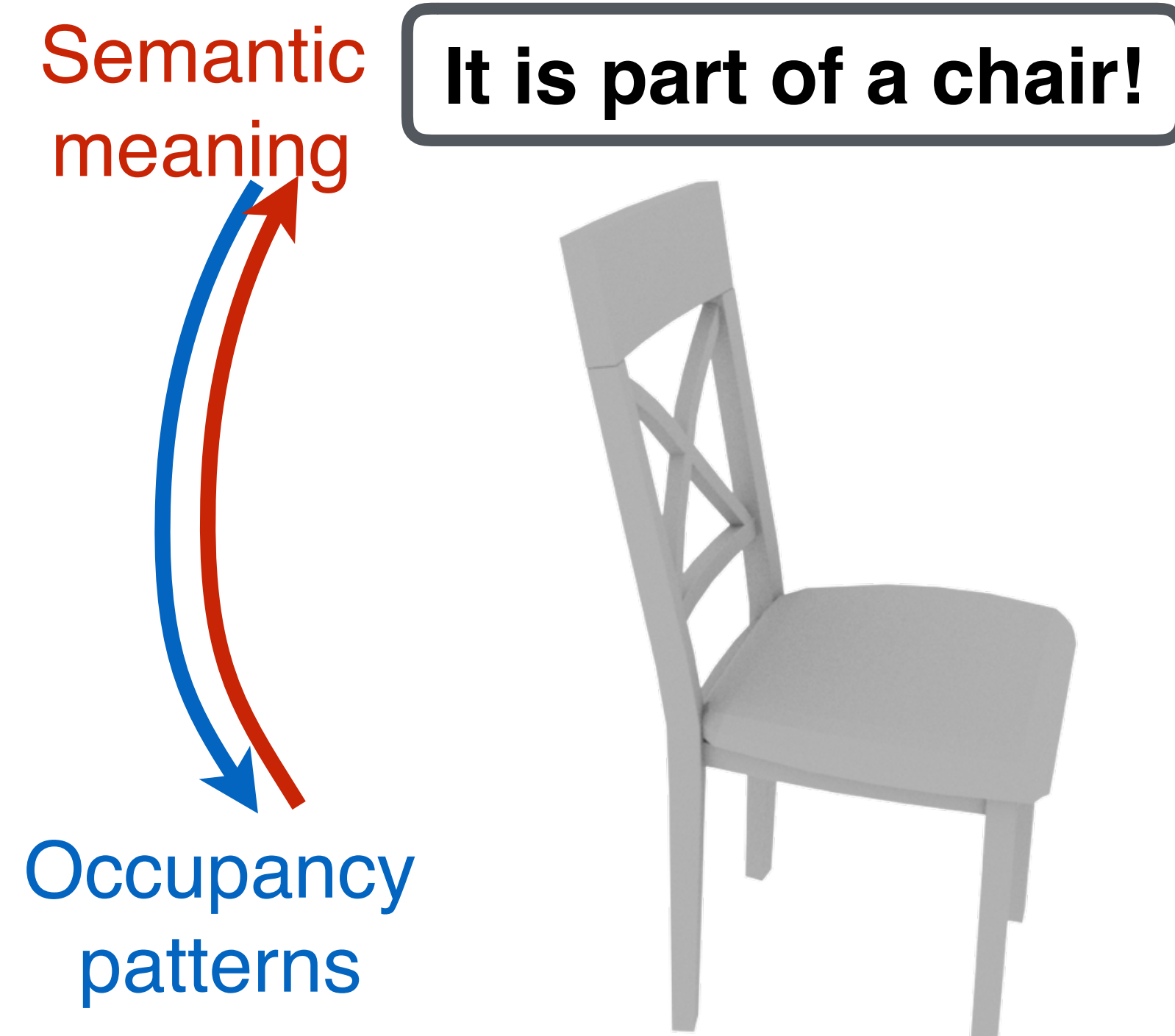


- floor
- wall
- window
- chair
- bed
- sofa
- table
- tv's
- furn.
- objects

Analysis

Key ideas:

1. Object occupancy and the identity are tightly intertwined.
2. It is important to capture and understand 3D context with big receptive fields.



Does joint understanding help?

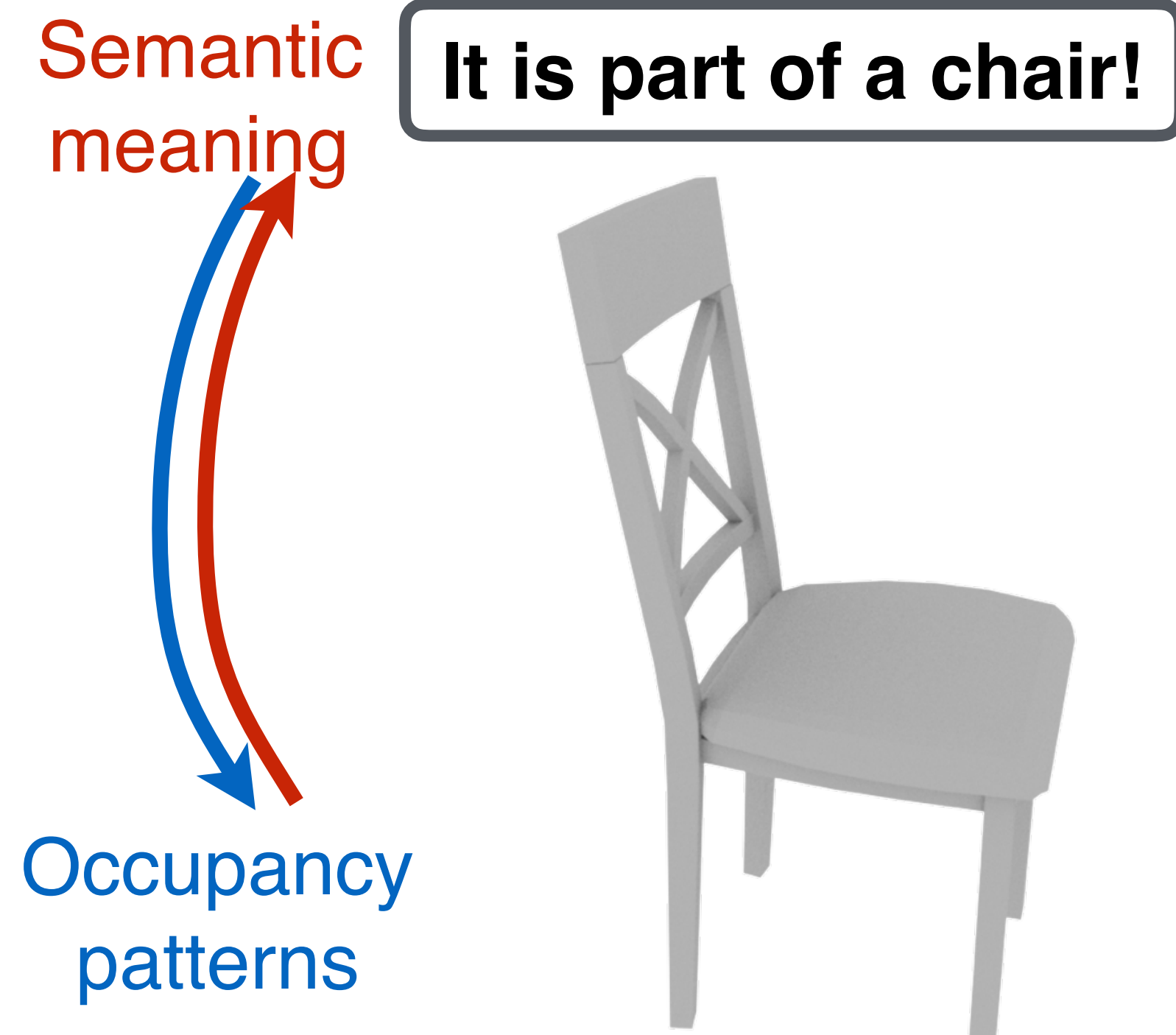
Semantic
meaning

It is part of a chair!



Occupancy
patterns

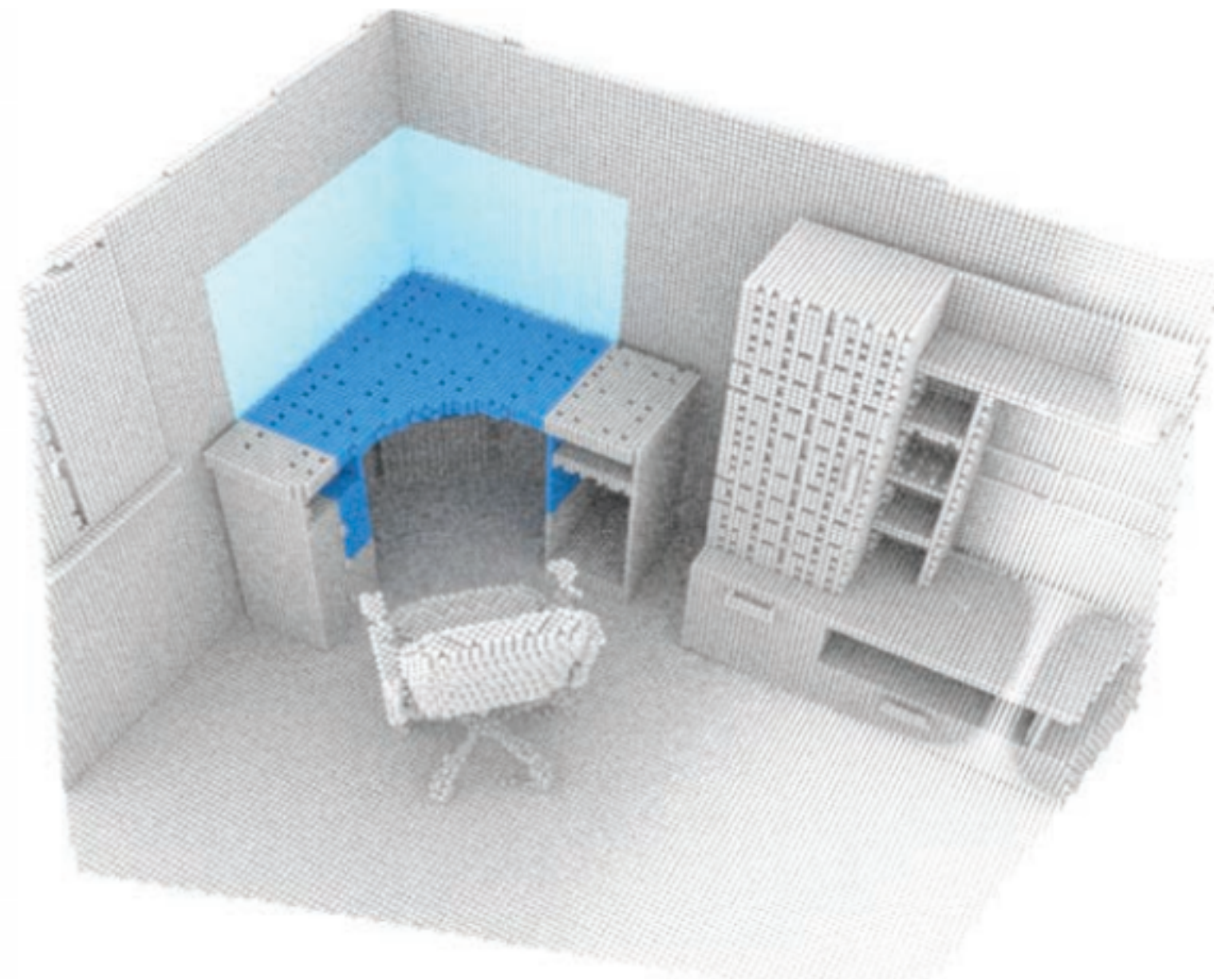
Does joint understanding help?



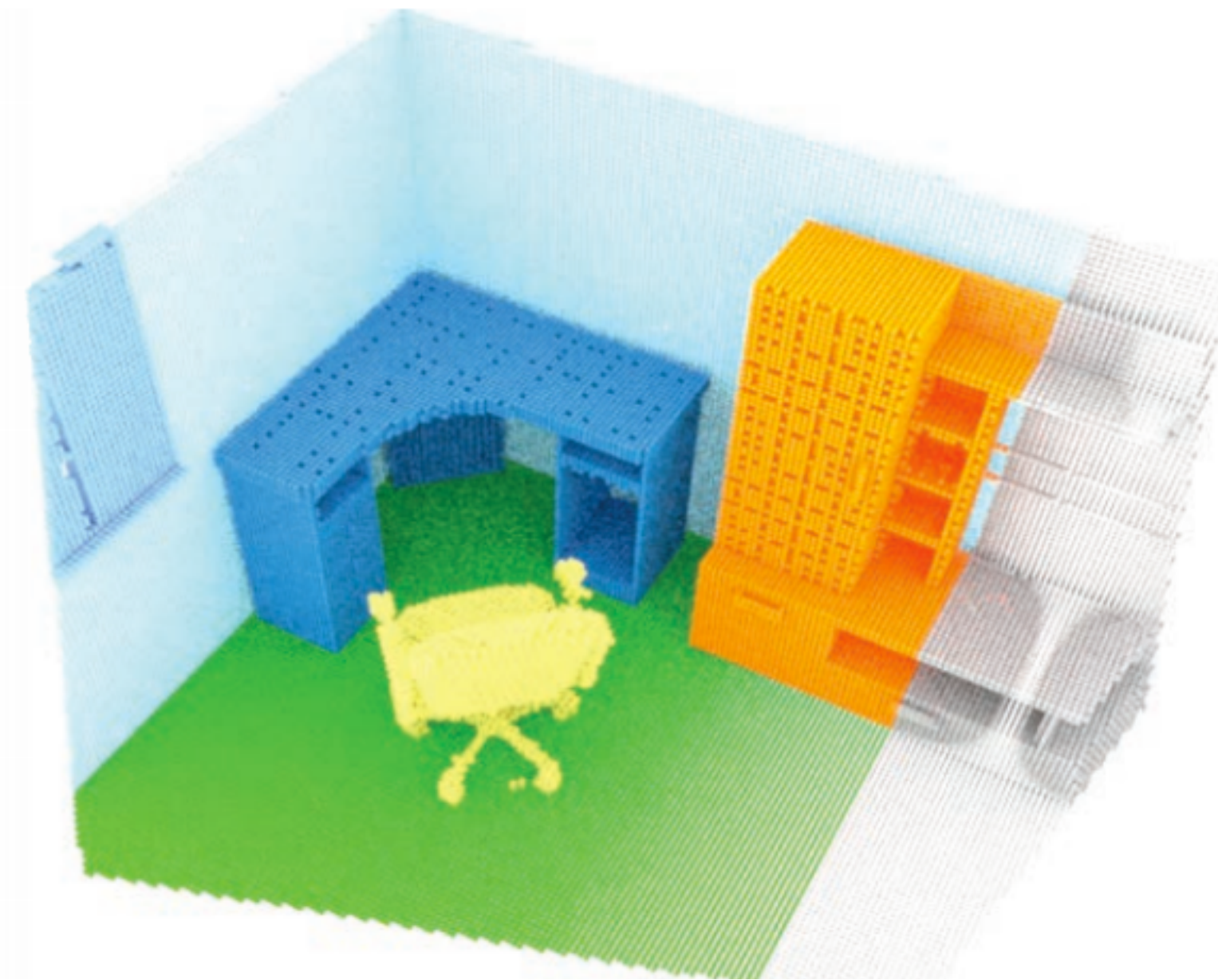
task	scene completion w/o semantics	semantic labeling w/o completion
completion only	64.8	-
semantic only	-	51.2
joint	73.0	54.2

Does a bigger receptive field help?

model/task	Basic	Basic + Dilation
semantic scene completion (IoU %)	38.0	44.3



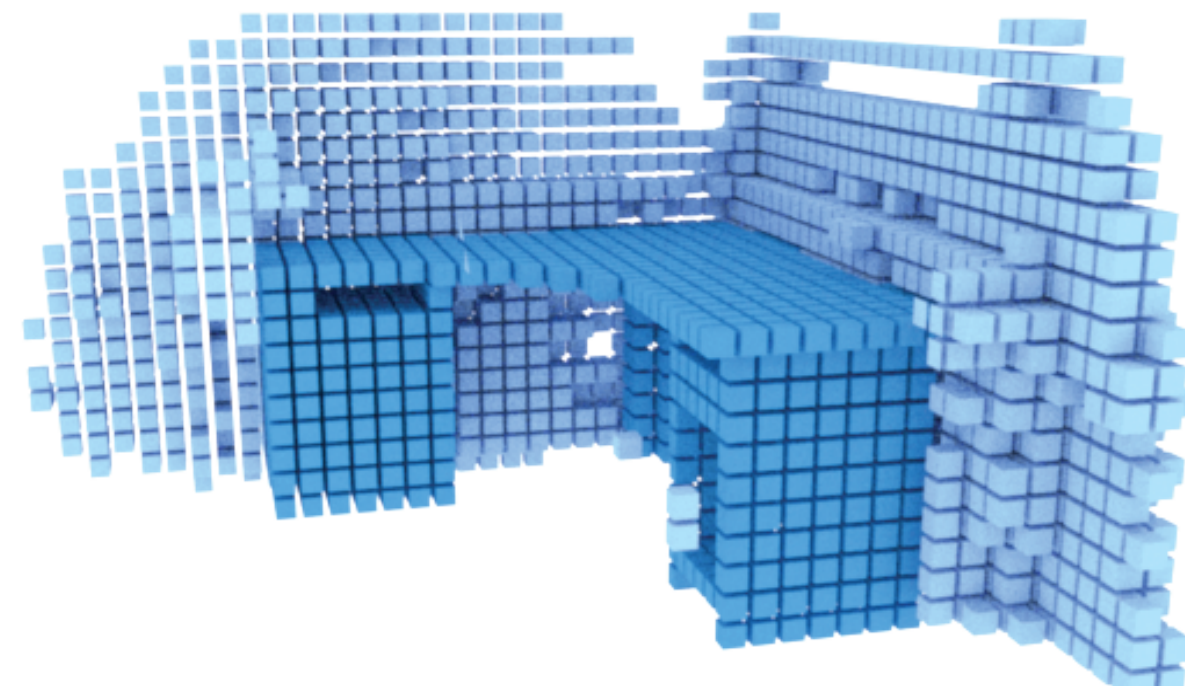
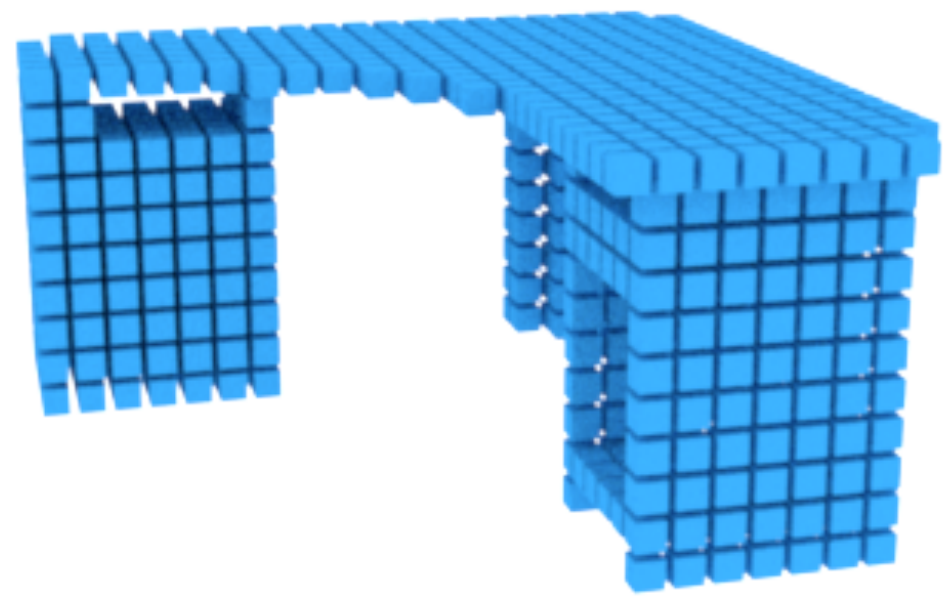
receptive field: 1 meter



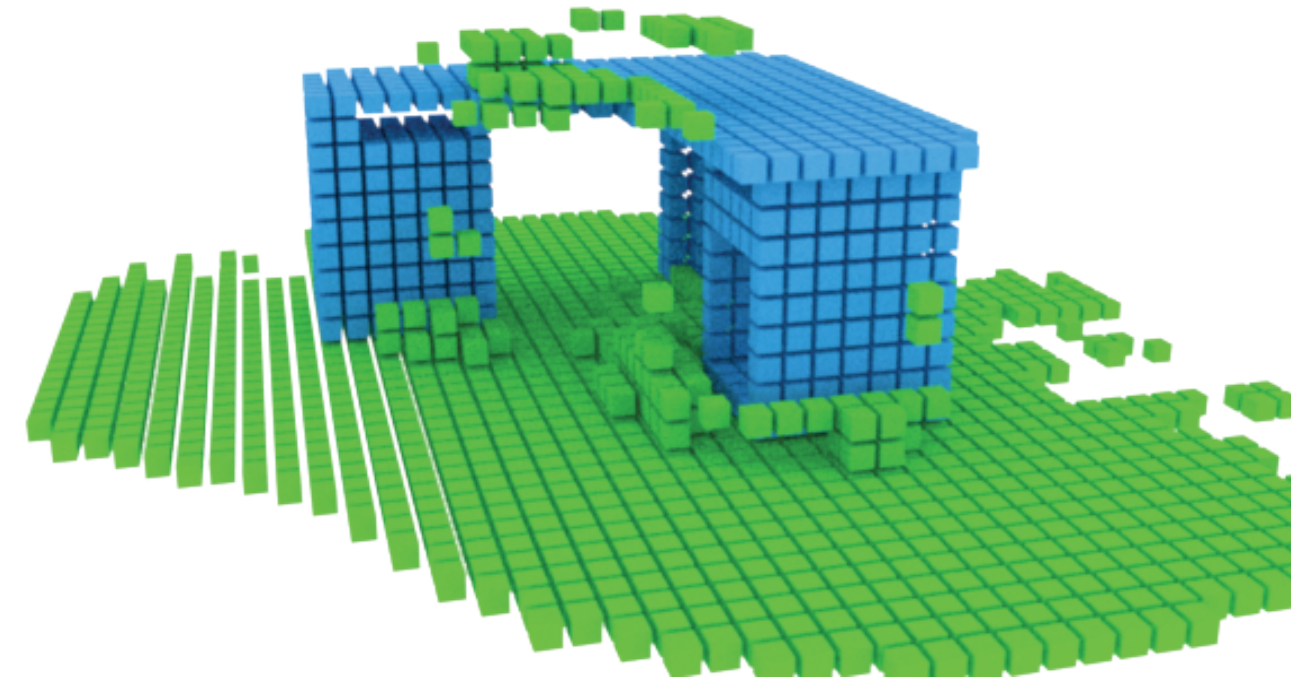
receptive field: 2.26 meter

What 3D context does the network learn?

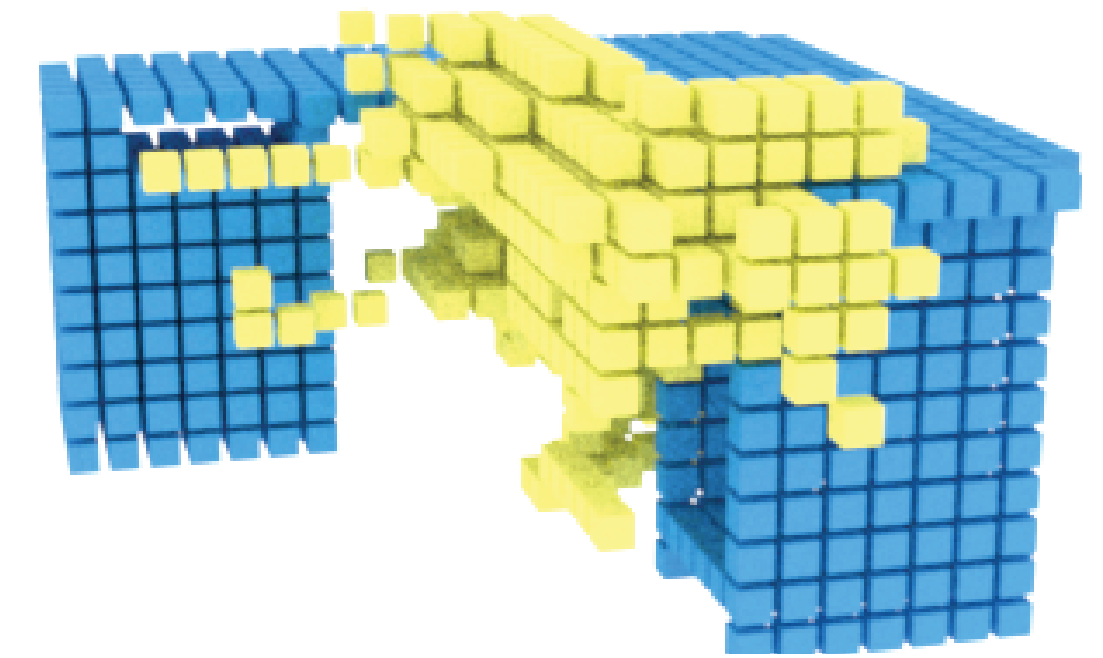
What 3D context does the network learn?



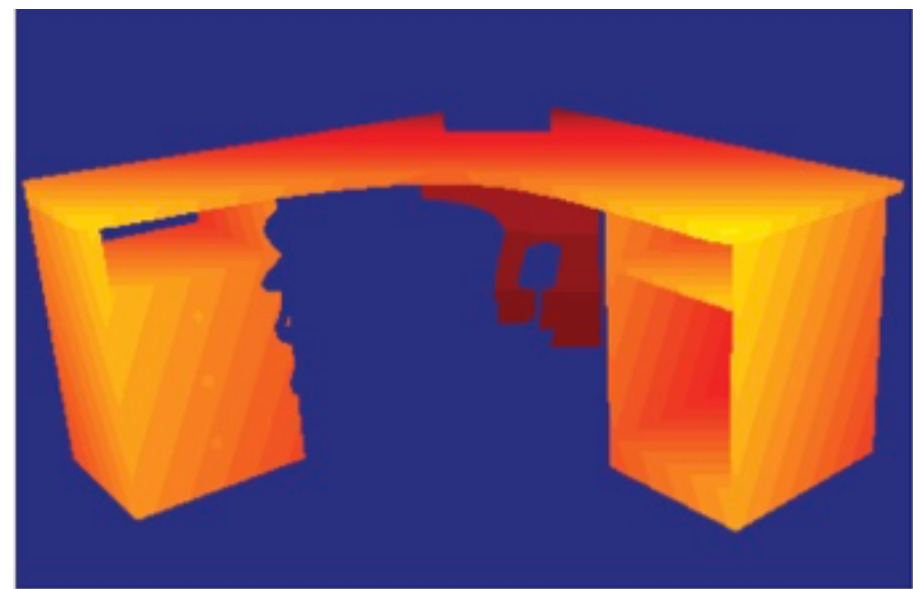
Wall



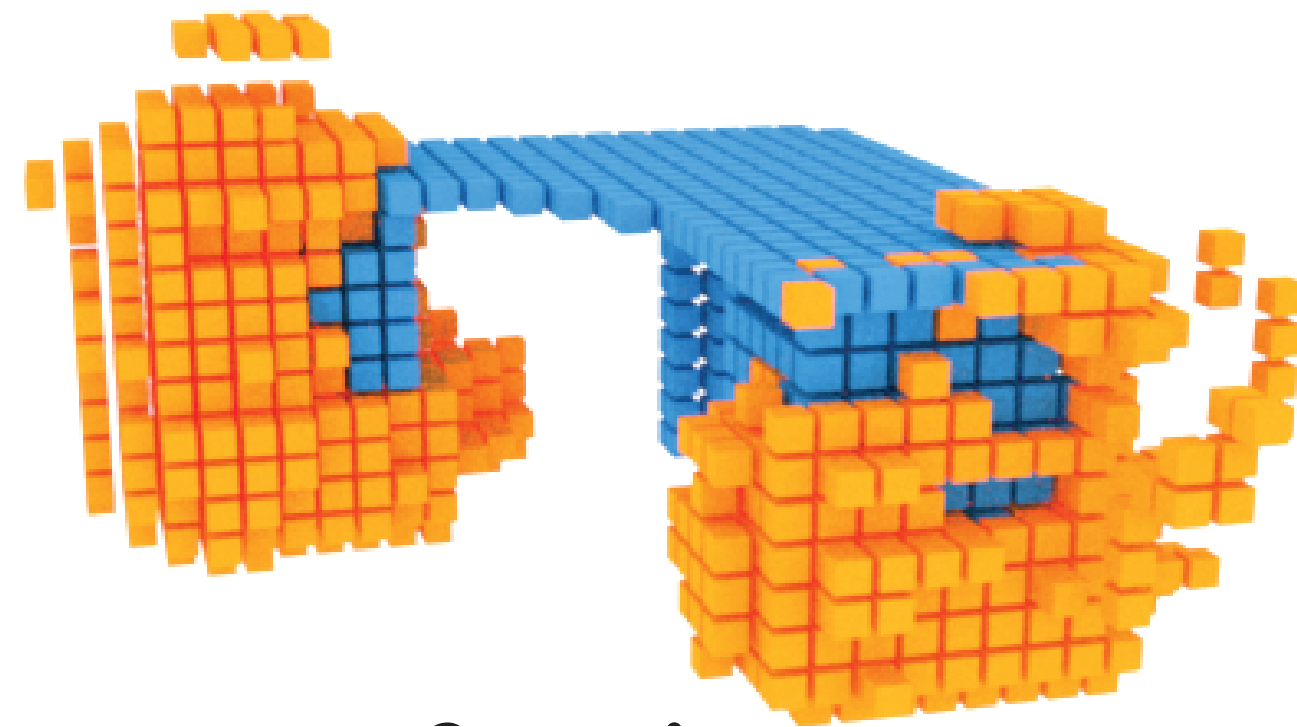
floor



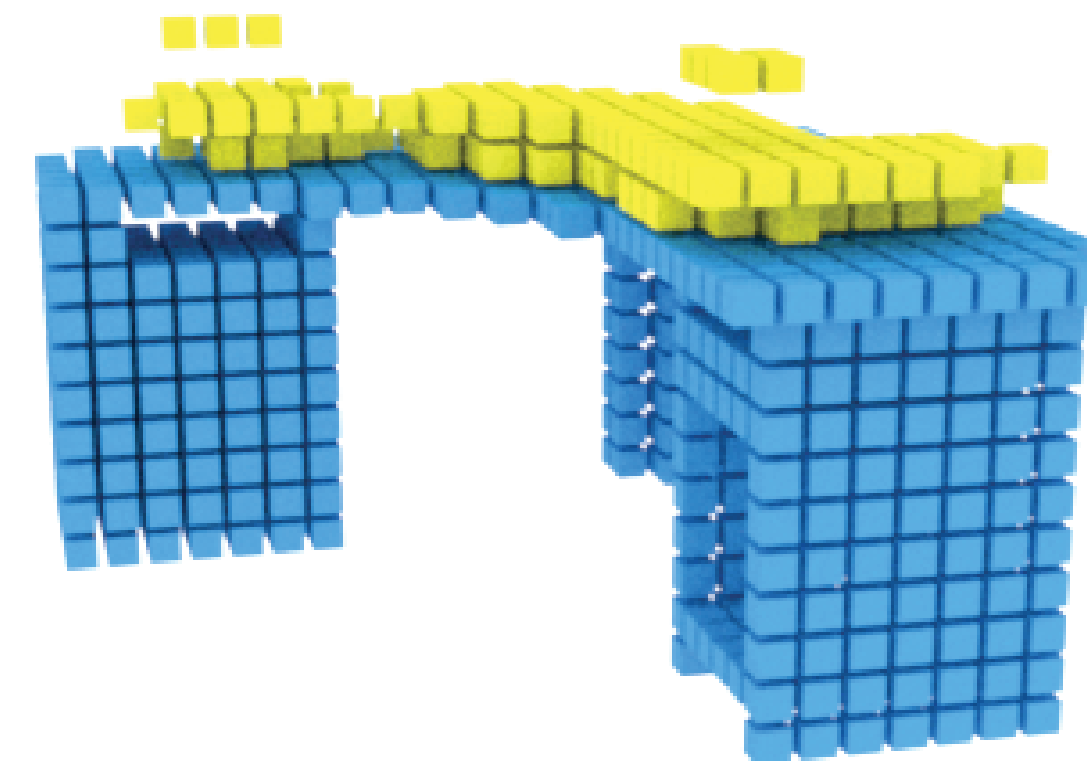
chair



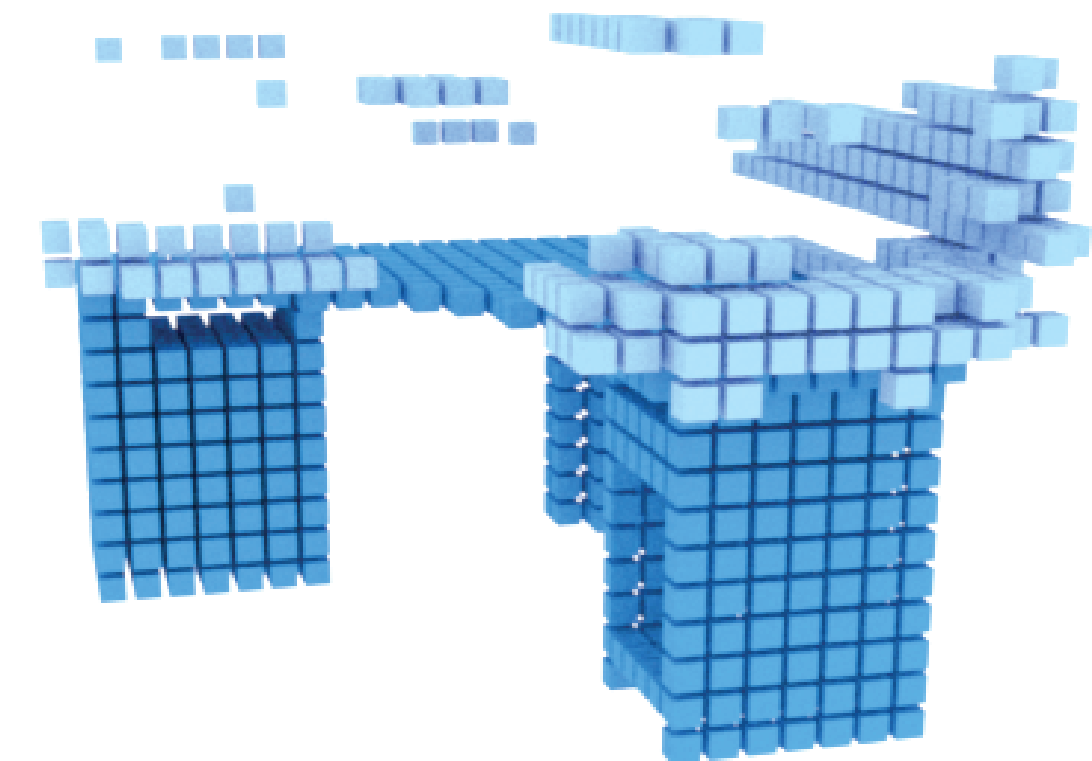
Input



furniture



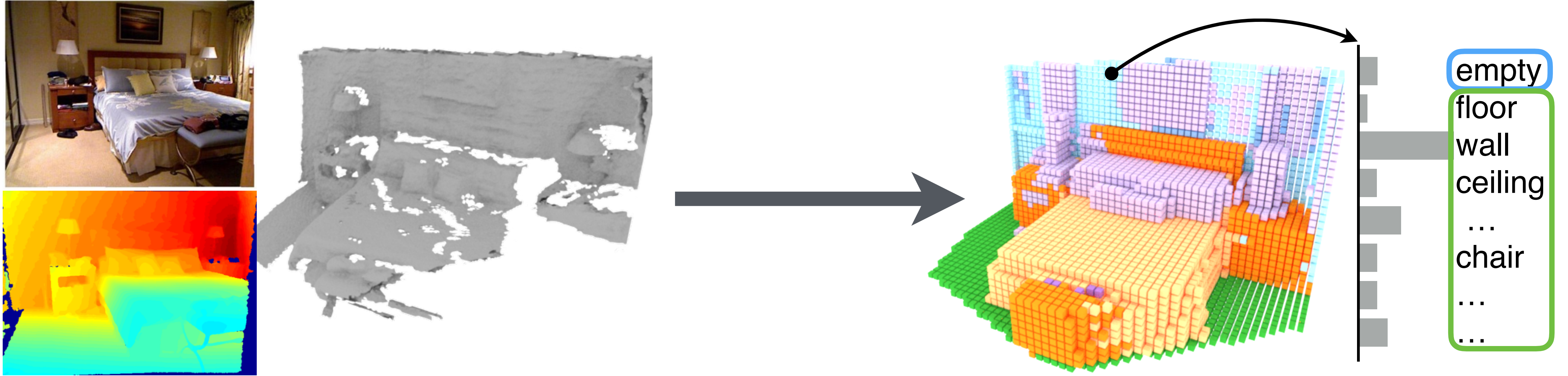
tv/monitor



window

Conclusion

Conclusion

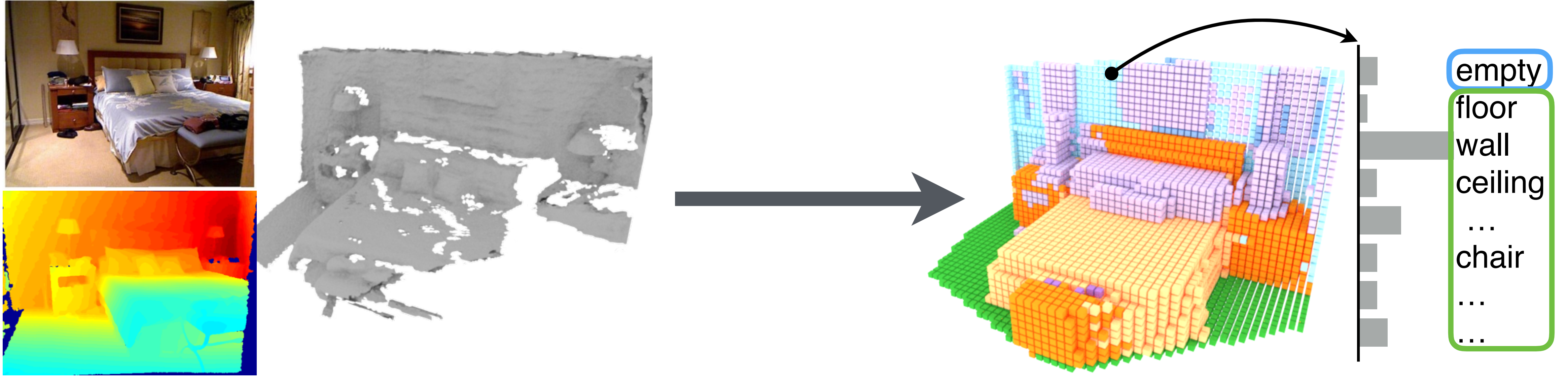


Semantic Scene Completion

- Semantic scene completion network, SSCNet
- A large-scale synthetic scene dataset, SUNCG

Code & Data: sscnet.cs.princeton.edu

Conclusion

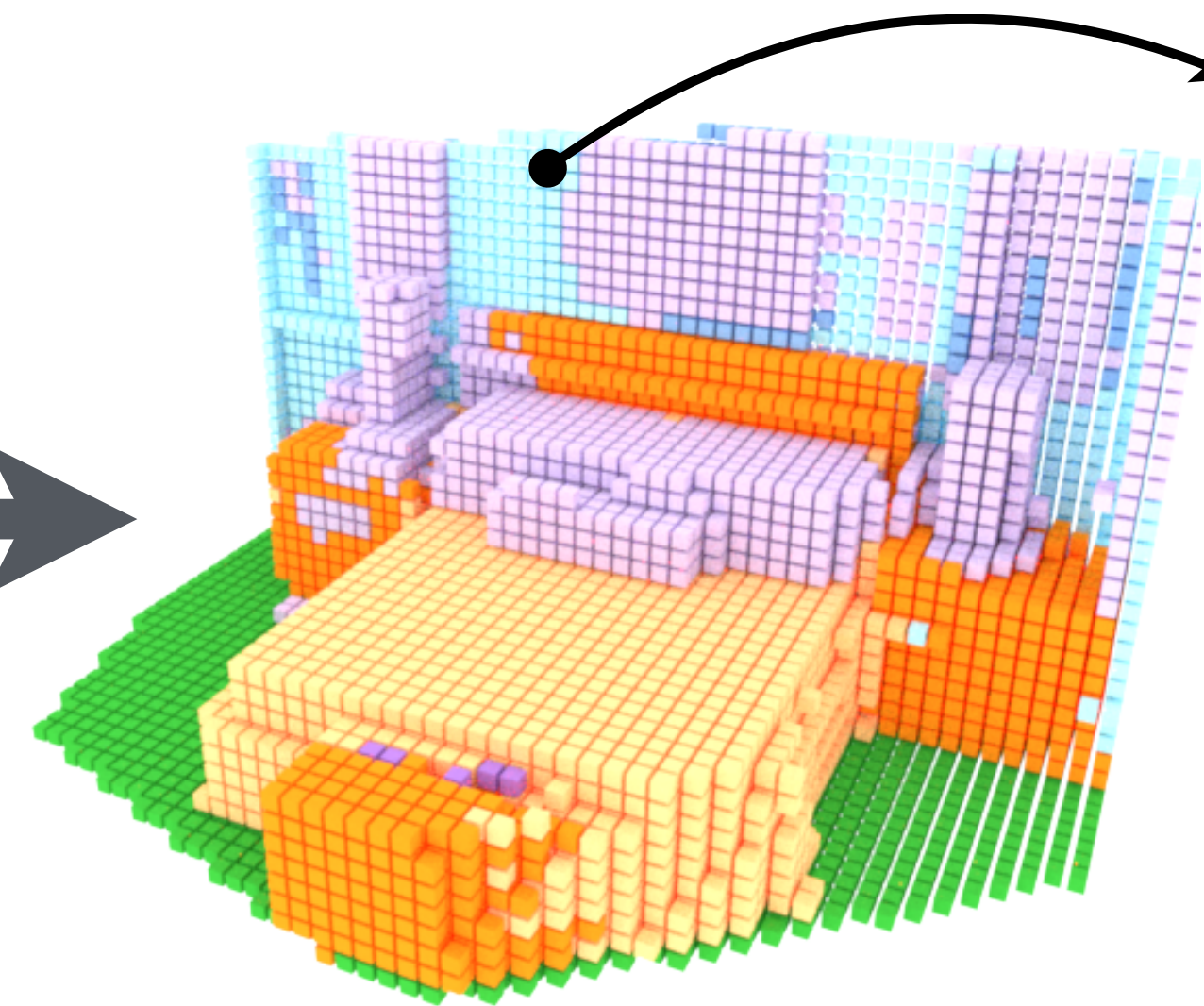
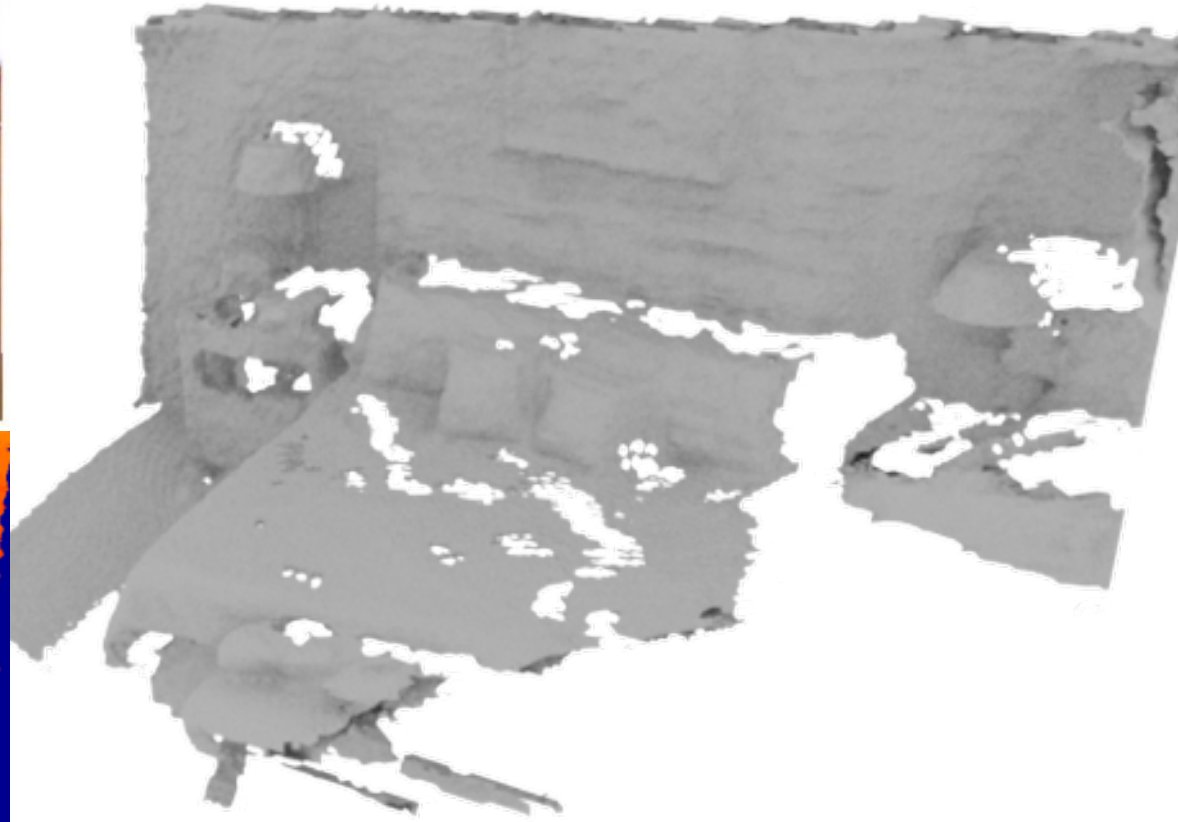
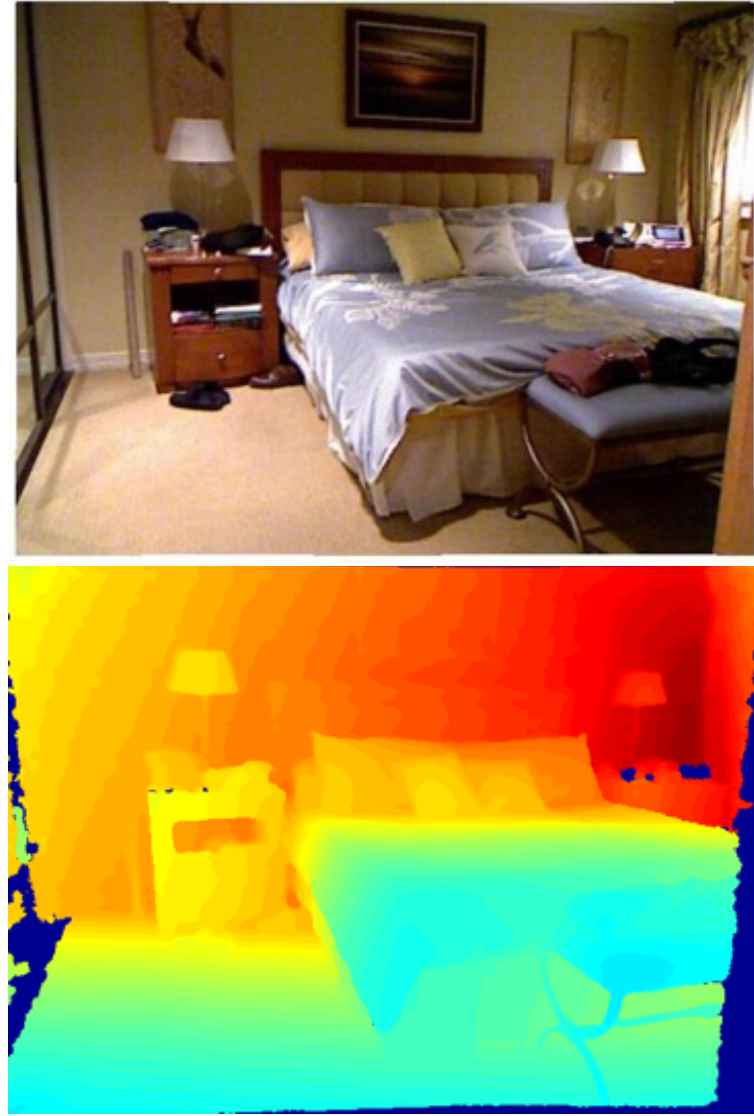


Semantic Scene Completion

- Semantic scene completion network, SSCNet
- A large-scale synthetic scene dataset, SUNCG

Code & Data: sscnet.cs.princeton.edu

Conclusion



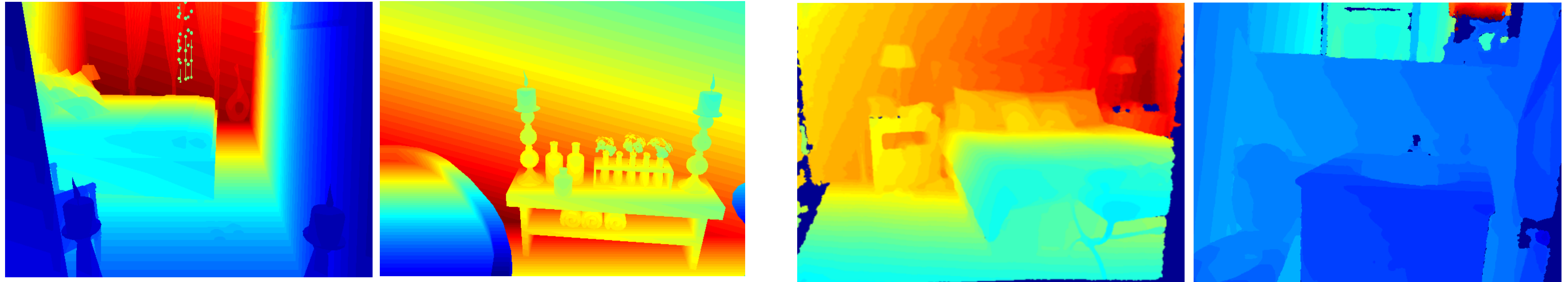
- empty
- floor
- wall
- ceiling
- ...
- chair
- ...
- ...

Semantic Scene Completion

- Semantic scene completion network, SSCNet
- A large-scale synthetic scene dataset, SUNCG

Code & Data: sscnet.cs.princeton.edu

Does synthetic data help?



Rendered depth map

Kinect depth map

Test on NYU	NYU	SUNCG	SUNCG+NYU
semantic scene completion (IoU %)	24.7	20.2	30.5

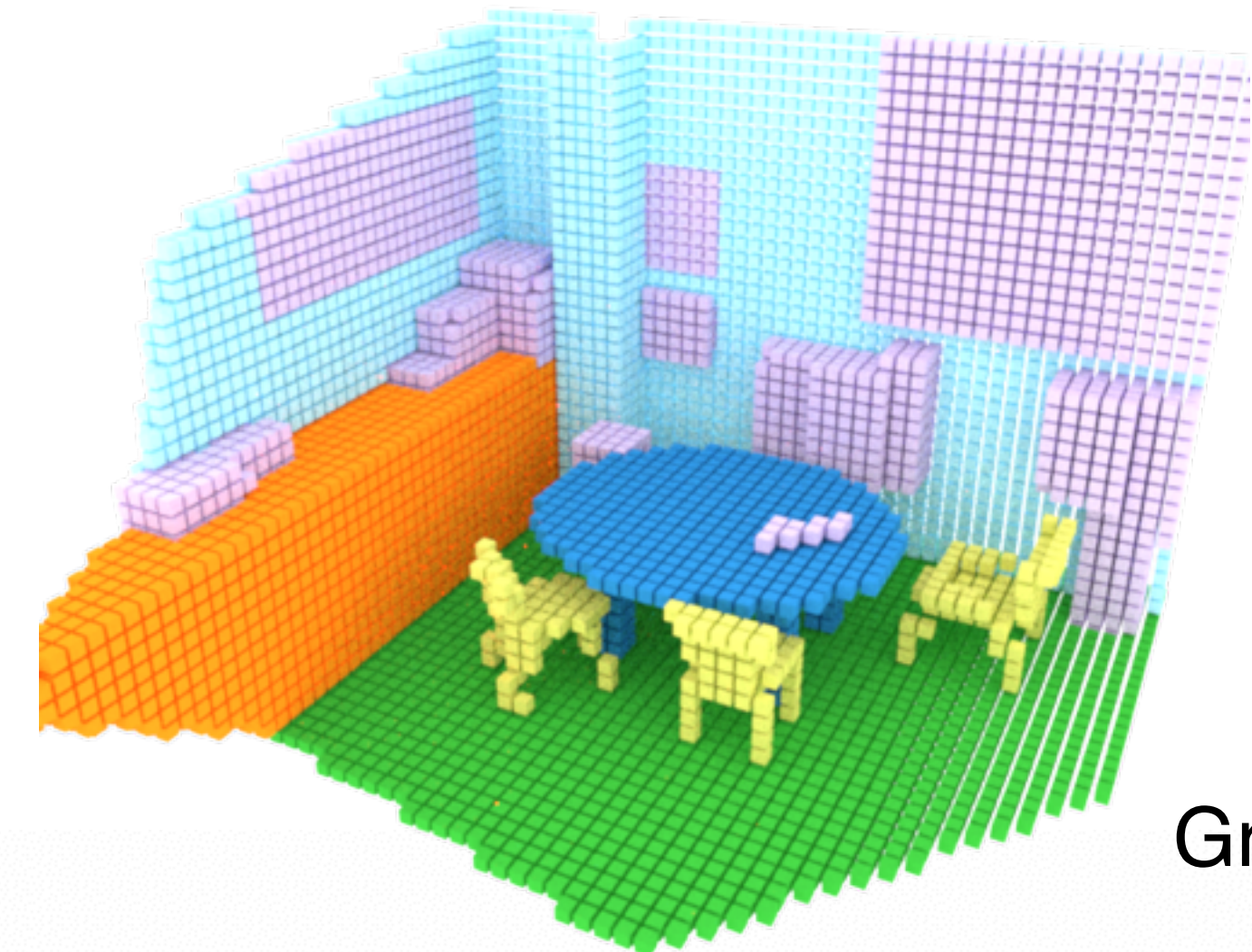
Failure cases

Failure cases: missing fine structures

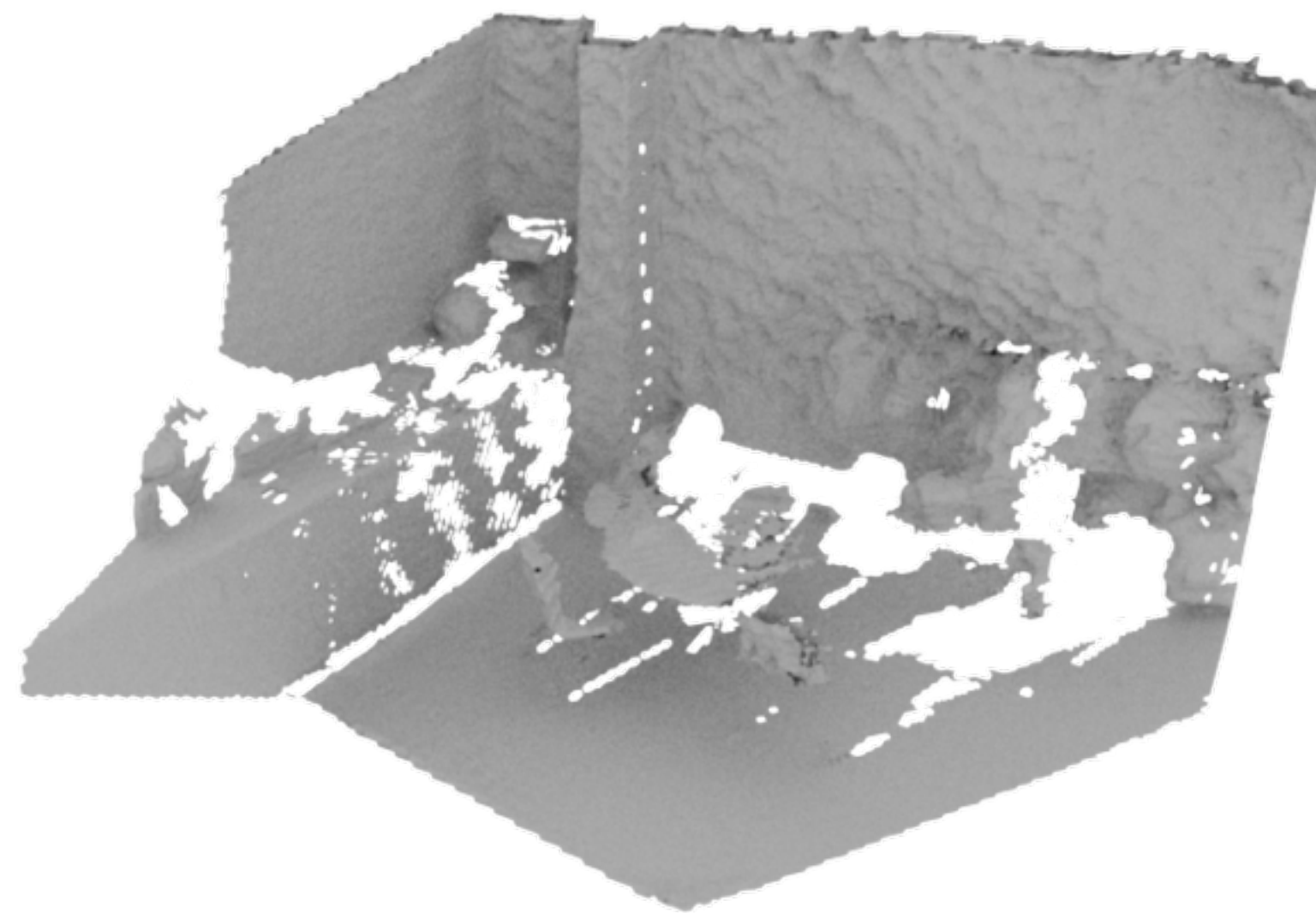
Color Image



Ground Truth

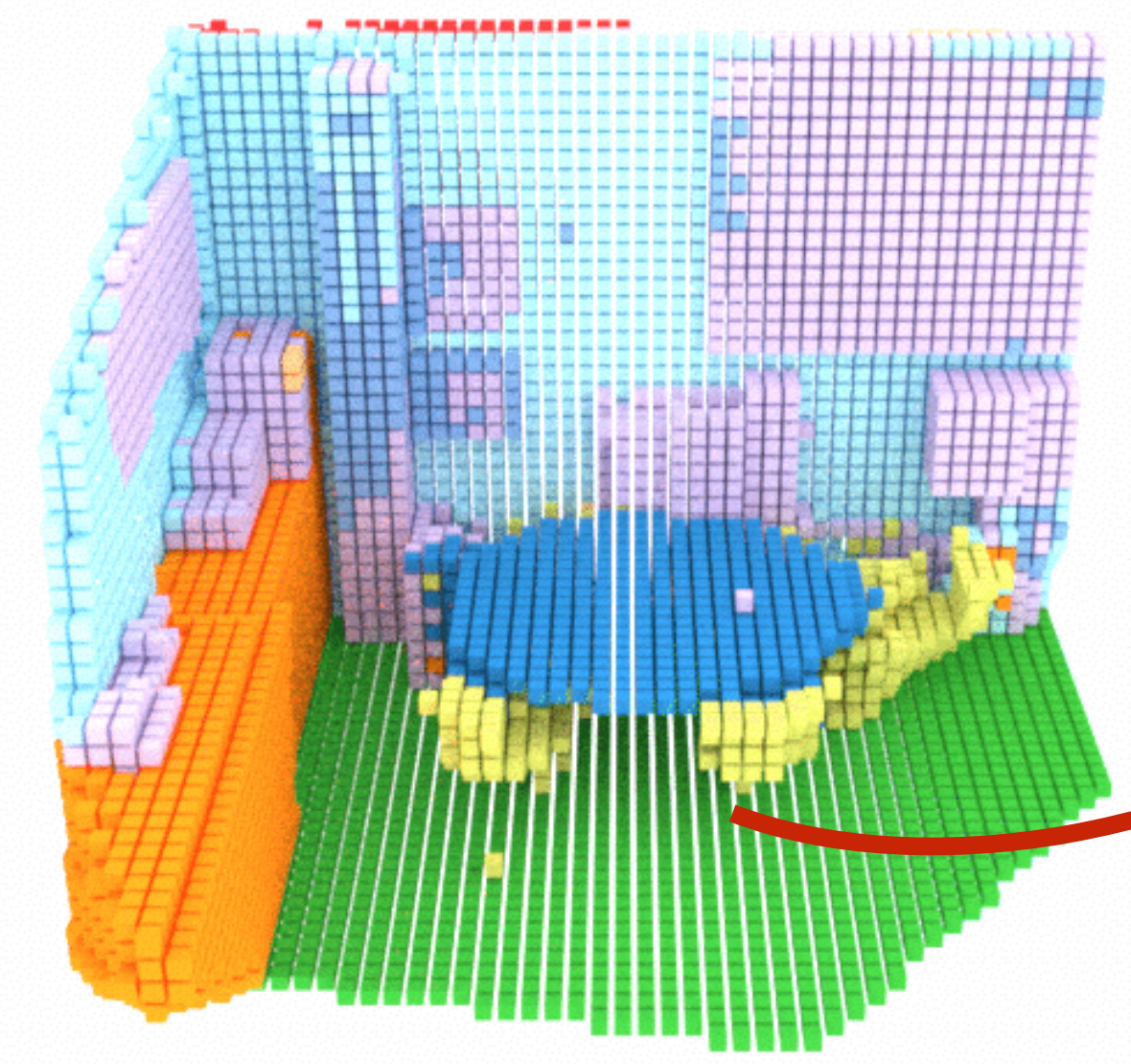


Observed Surface



Missing chair legs

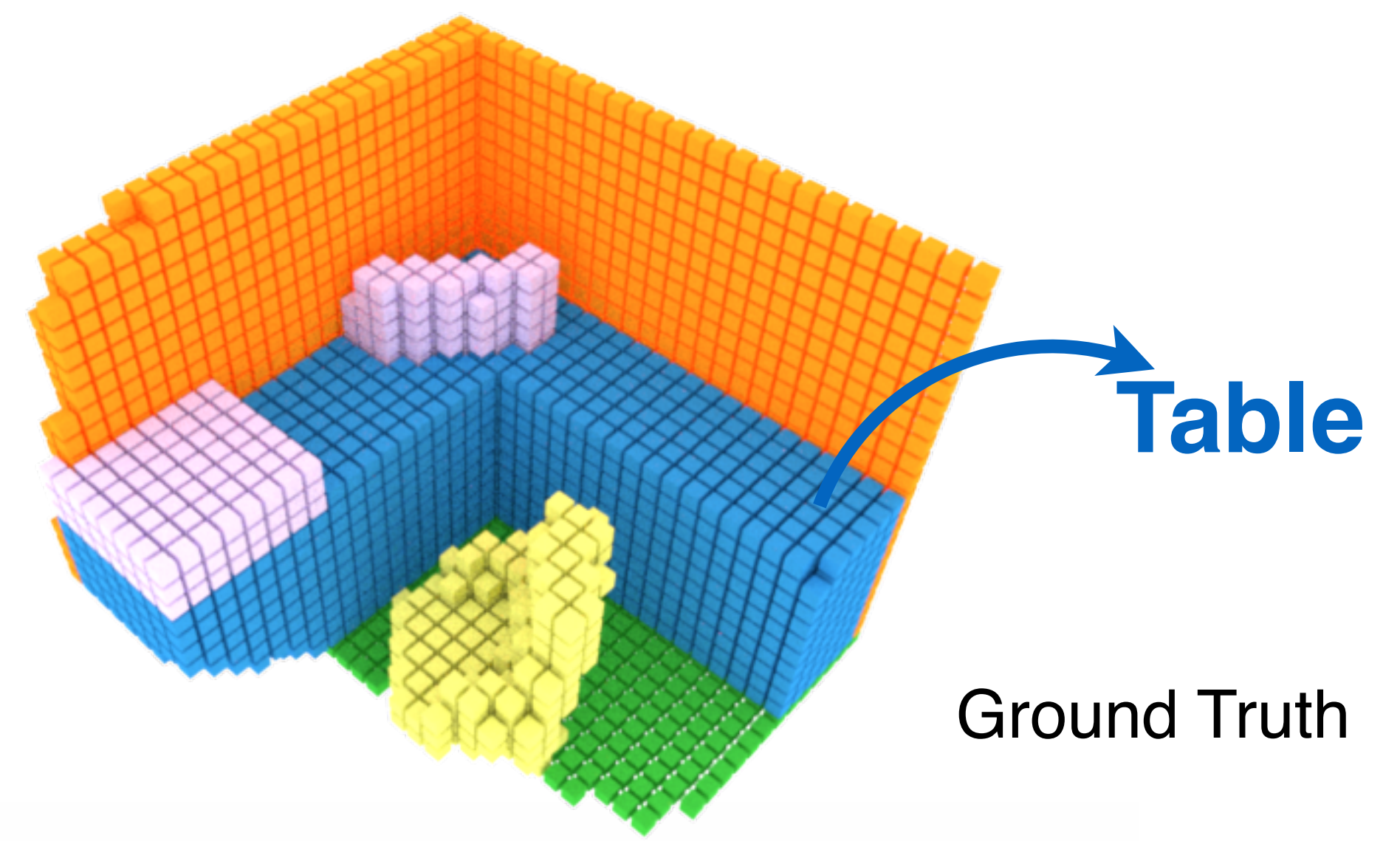
SSCNet prediction



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

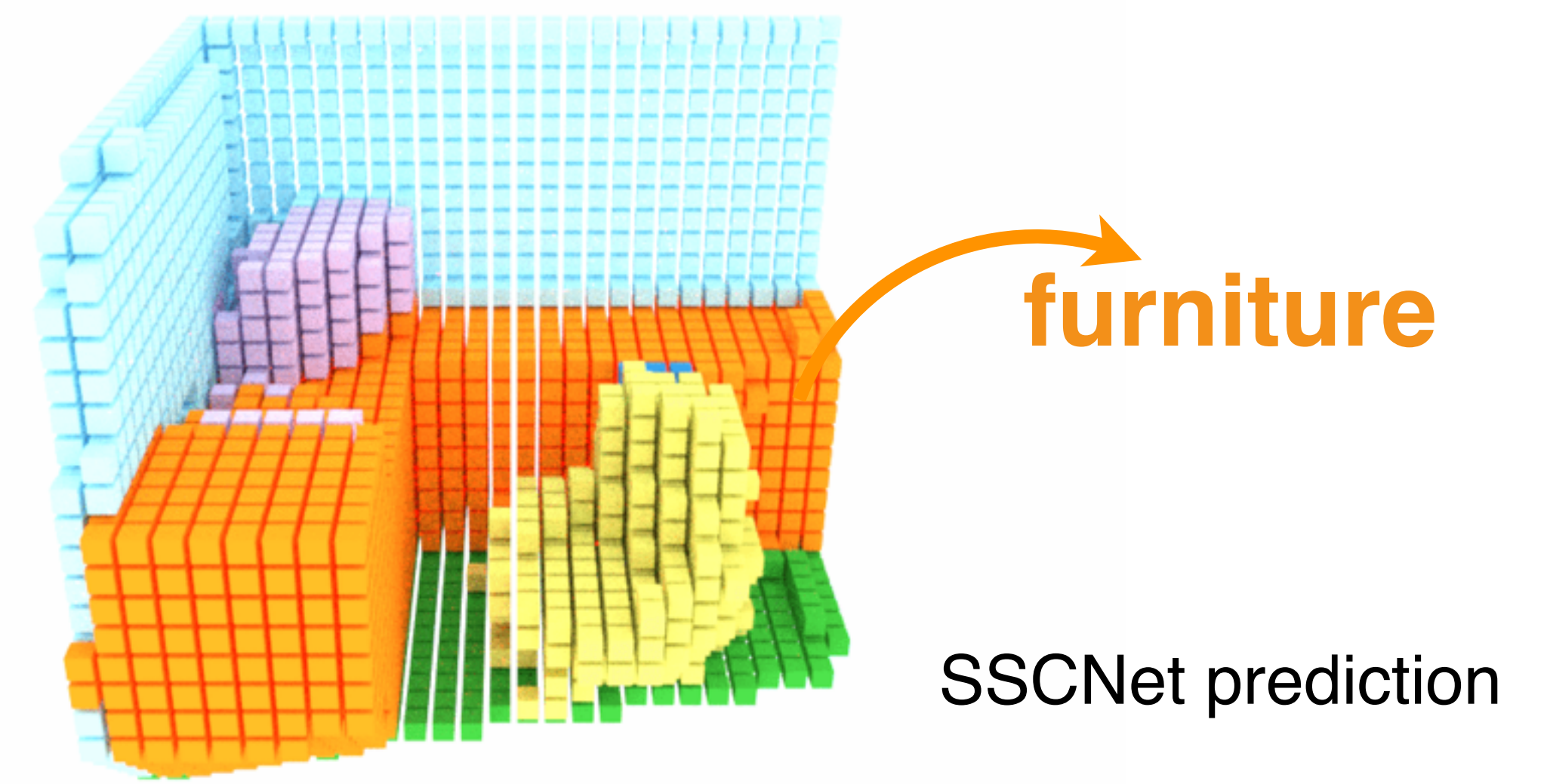
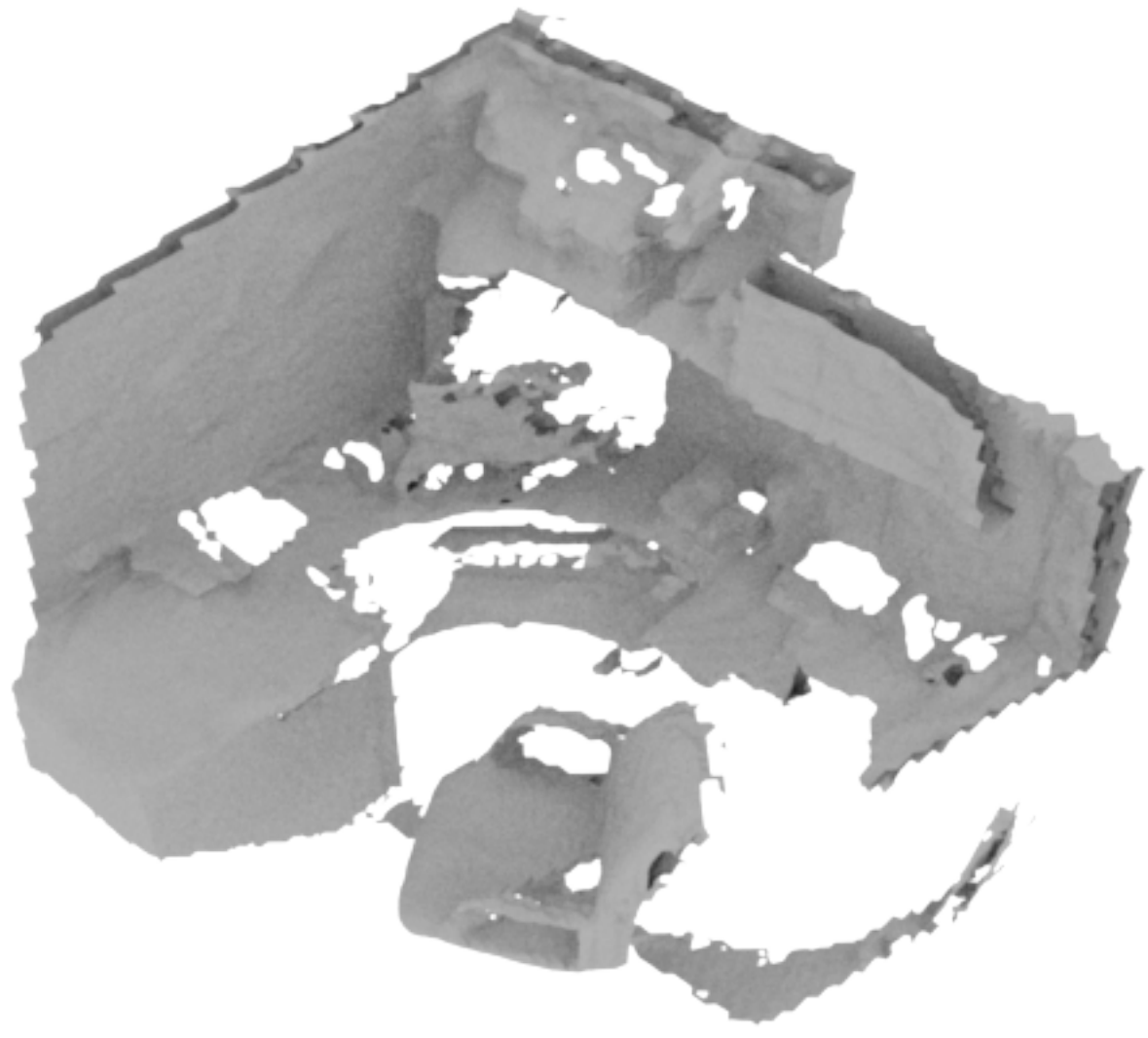
Failure cases: wrong object category

Color Image



Ground Truth

Observed Surface



SSCNet prediction

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

