

3D ShapeNets: A Deep Representation for Volumetric Shape Modeling

by Wu, Song, Khosla, Yu, Zhang, Tang, Xiao

presented by Abhishek Sinha

3D Shape Prior



3D Shape Prior



3D Shape Prior



Outline

Outline

- Problem

Outline

- Problem
- Motivation

Outline

- Problem
- Motivation
- Desirable Properties for Representation

Outline

- Problem
- Motivation
- Desirable Properties for Representation
- Architecture

Outline

- Problem
- Motivation
- Desirable Properties for Representation
- Architecture
- Dataset

Outline

- Problem
- Motivation
- Desirable Properties for Representation
- Architecture
- Dataset
- Applications

Outline

- Problem
- Motivation
- Desirable Properties for Representation
- Architecture
- Dataset
- Applications
- Extensions

Outline

- Problem
- Motivation
- Desirable Properties for Representation
- Architecture
- Dataset
- Applications
- Extensions
- Discussion Points

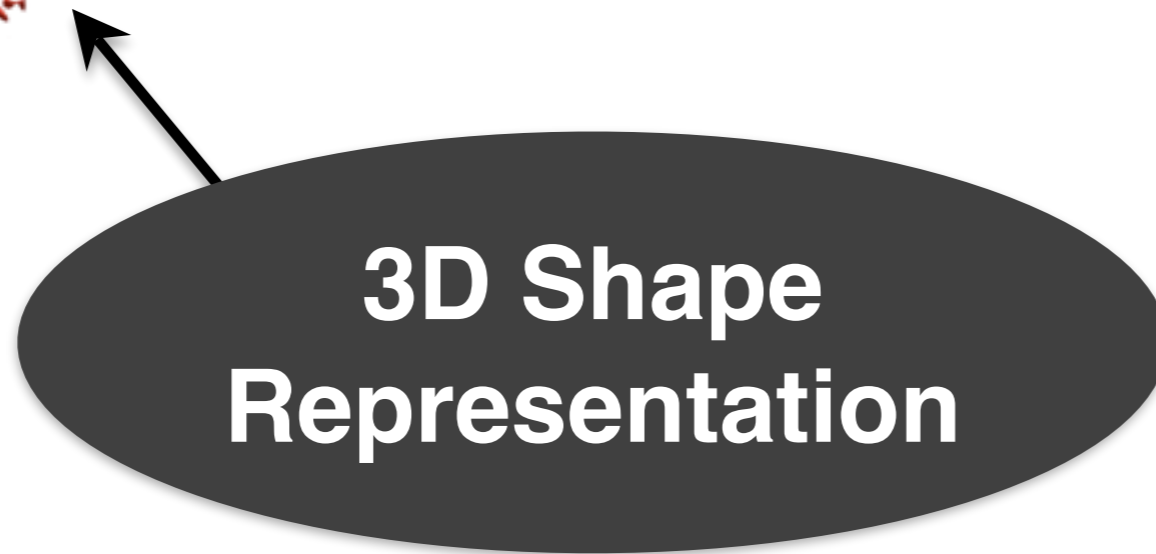
Problem

Learn 'Useful' 3D shape
representations from
images

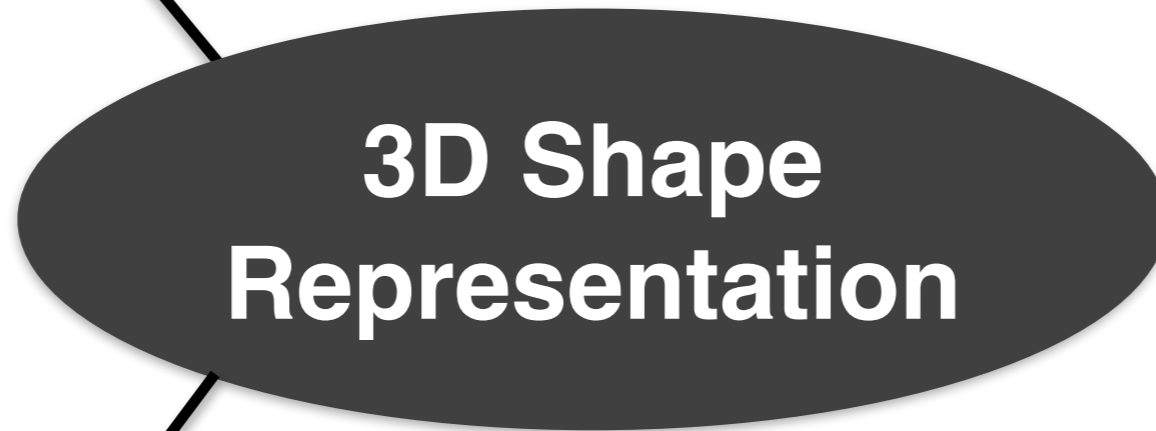
Motivation

3D Shape Representation

Shape Synthesis

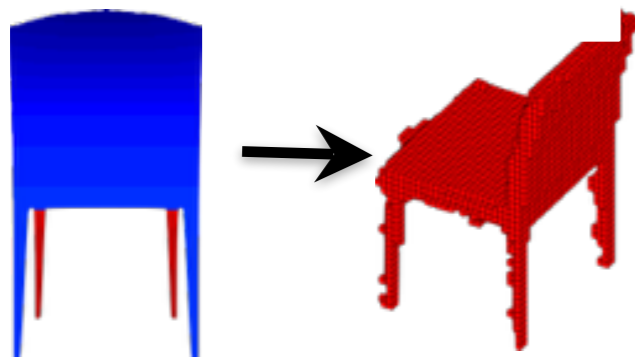


Shape Synthesis



**3D Shape
Representation**

Shape Completion



Shape Synthesis

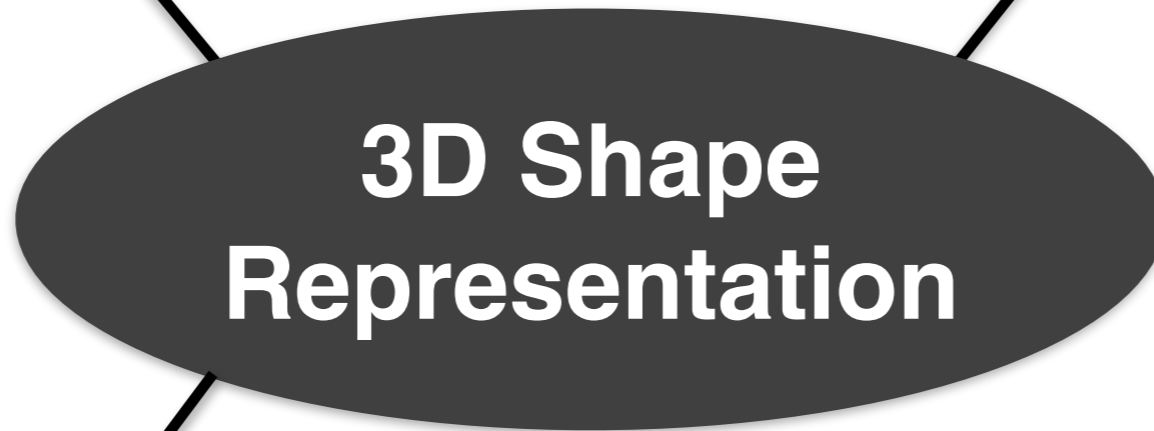


2.5D Object Recognition

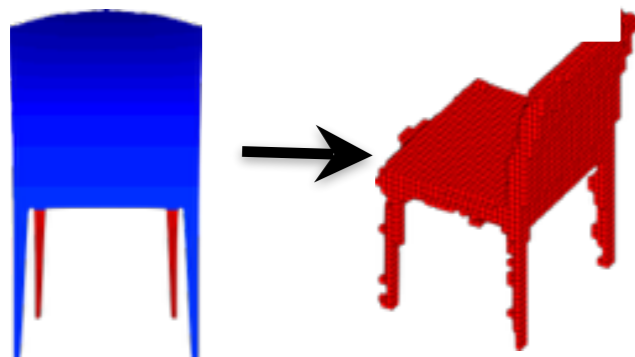


→ person

→ tricycle



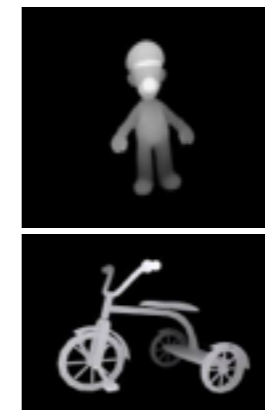
Shape Completion



Shape Synthesis

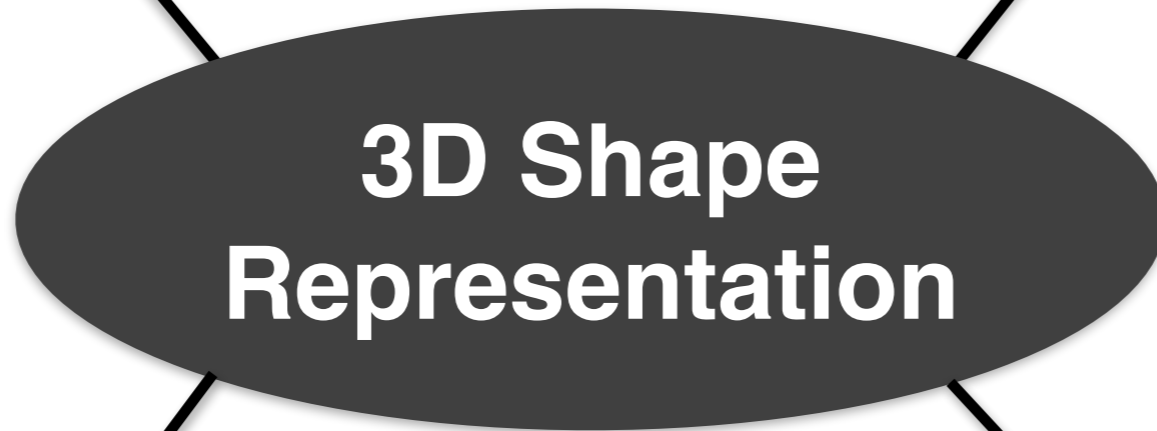


2.5D Object Recognition

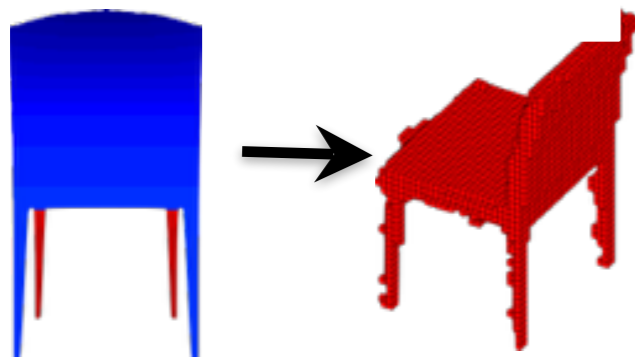


→ person

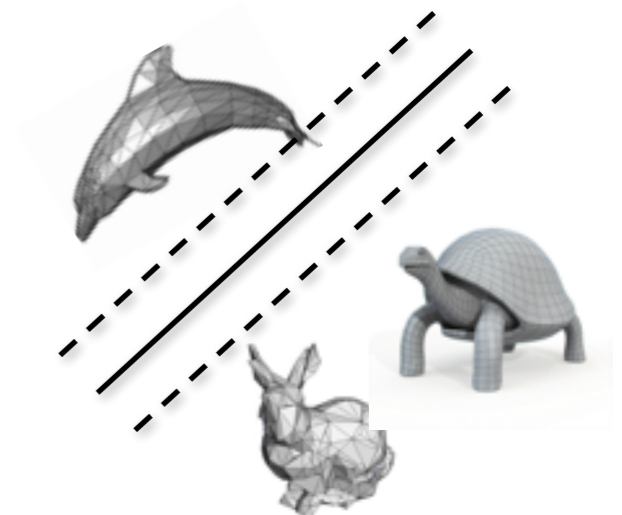
→ tricycle



Shape Completion



Feature Extractor



Desirable Properties

What is a Desirable 3D Shape Representation?

What is a Desirable 3D Shape Representation?

Data-driven

Generic

Compositional

Versatile

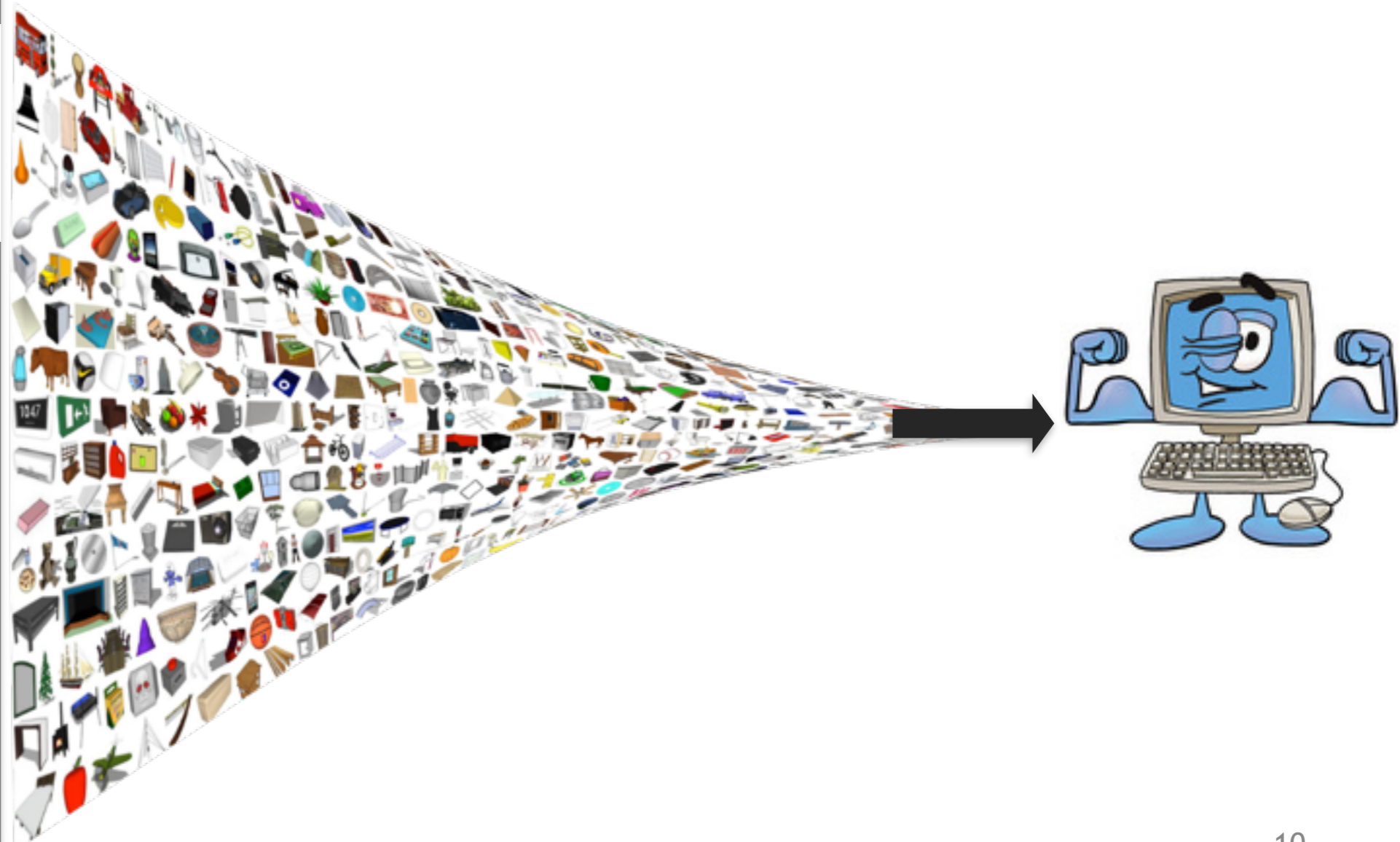
What is a Desirable 3D Shape Representation?

Data-driven

Data-driven

Compositional

Versatile



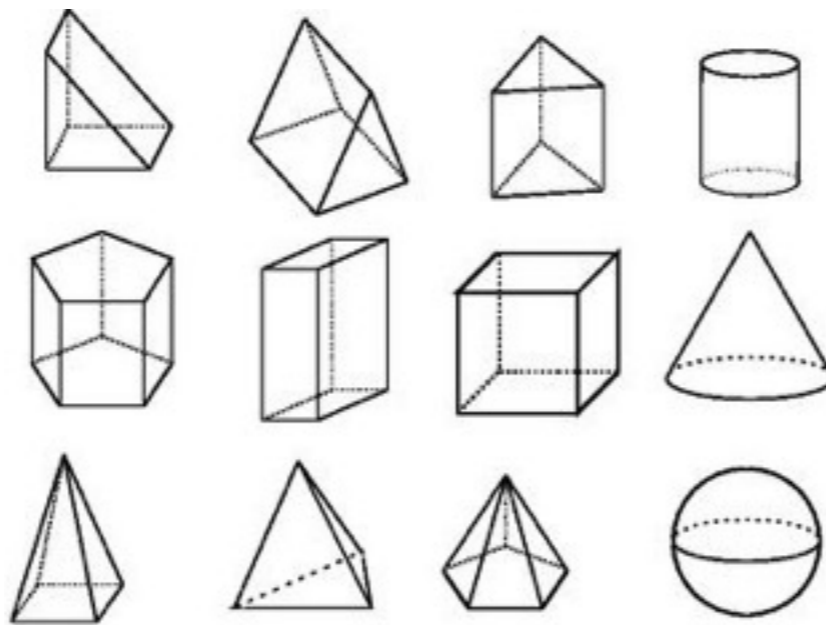
What is a Desirable 3D Shape Representation?

Data-driven

Generic

Generic

Versatile



Simple Shapes



Complex Shapes

What is a Desirable 3D Shape Representation?

Data-driven

Generic

Compositional

Compositiona



building blocks



full object

What is a Desirable 3D Shape Representation?

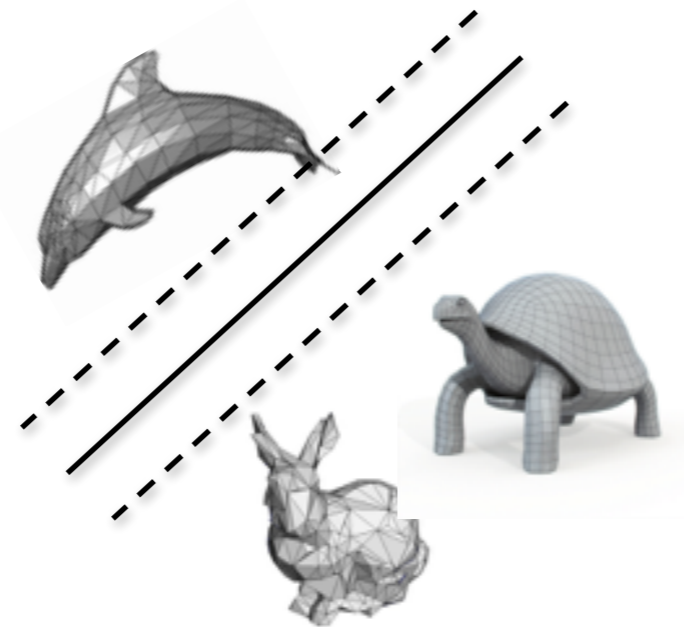
Data-driven

Generic

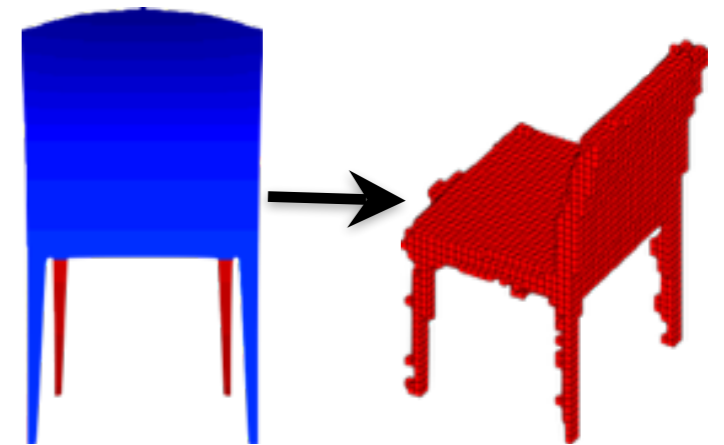
Compositional

Versatile

Versatile



mesh classification



shape completion



→ person



→ tricycle

2.5D object recognition



shape generation

Architecture

3D Deep Learning

3D Deep Learning

3D Shape as Volumetric Representation

3D Deep Learning

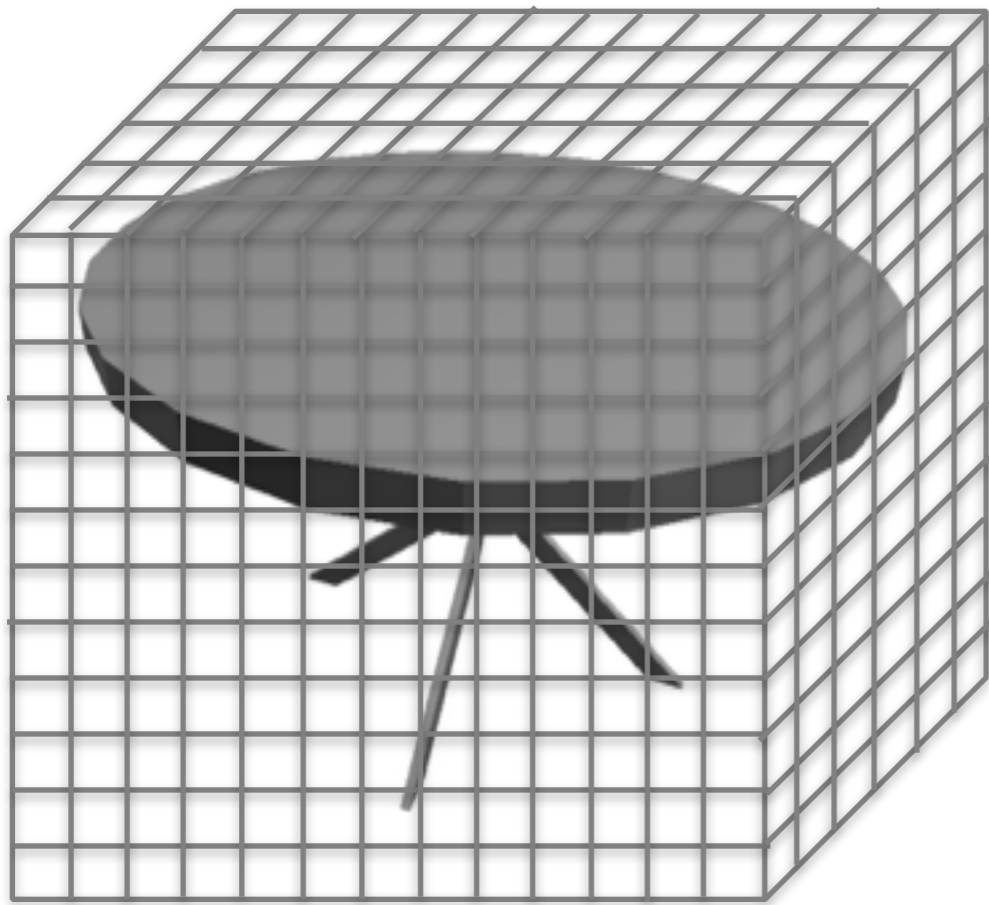
3D Shape as Volumetric Representation



mesh

3D Deep Learning

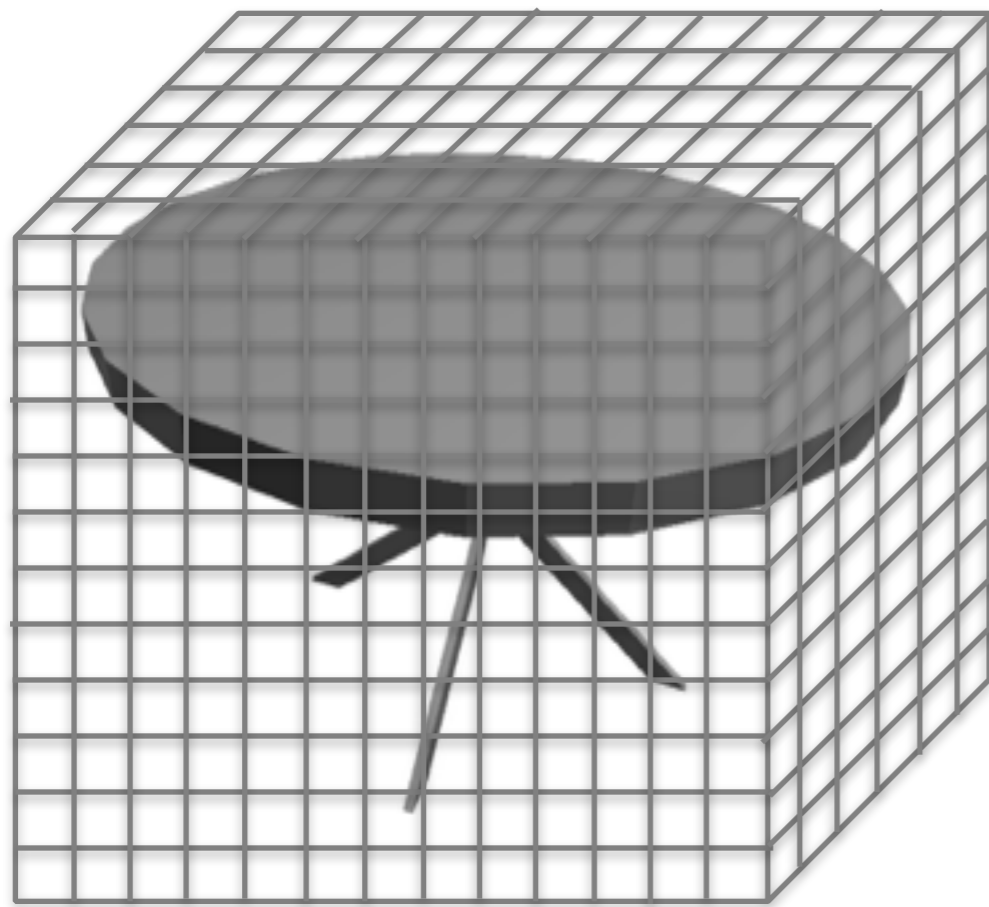
3D Shape as Volumetric Representation



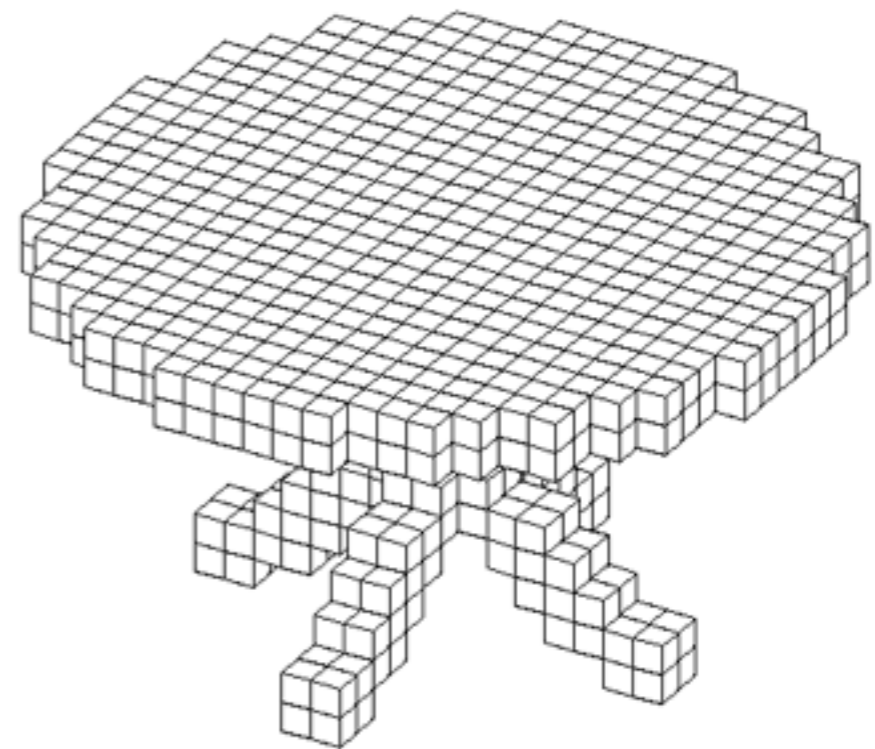
mesh

3D Deep Learning

3D Shape as Volumetric Representation

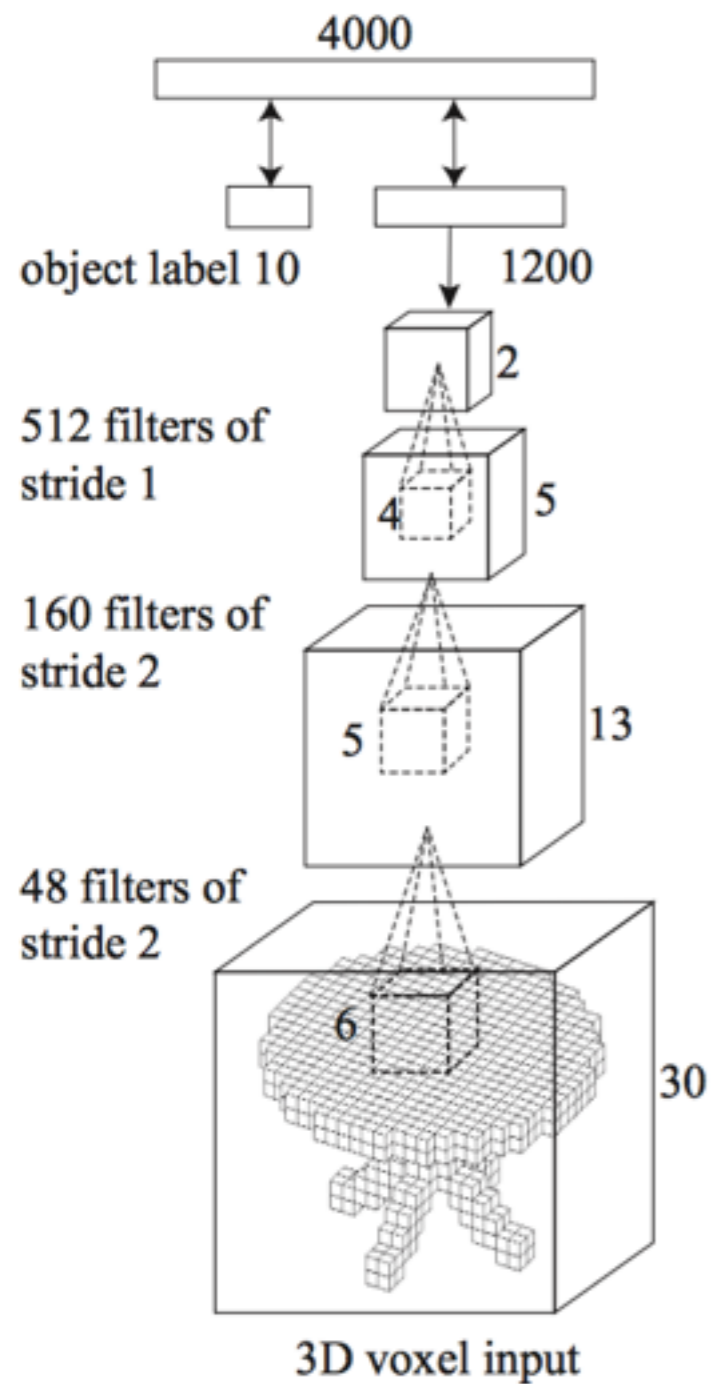


mesh



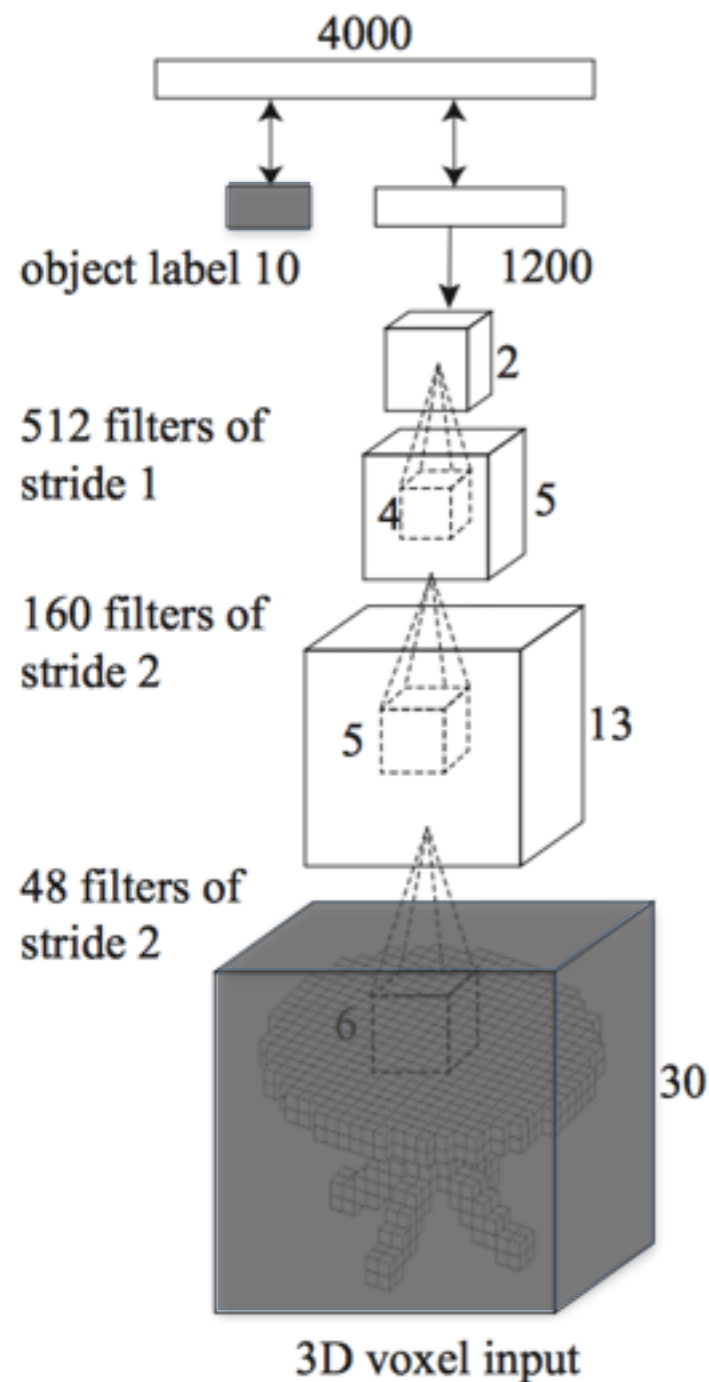
binary voxel

3D ShapeNets



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

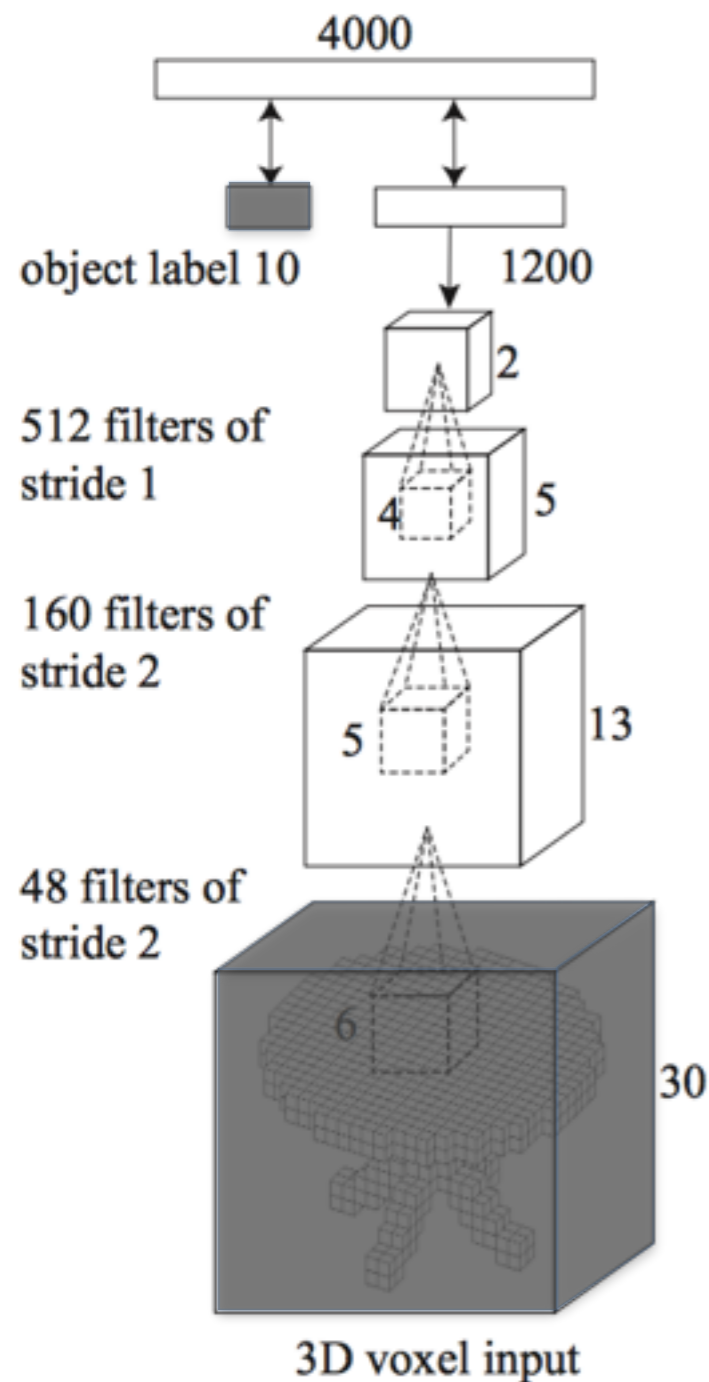
3D ShapeNets



A **Deep Belief Network** is a generative graphical model that describes the distribution of input x over class y .

**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

3D ShapeNets

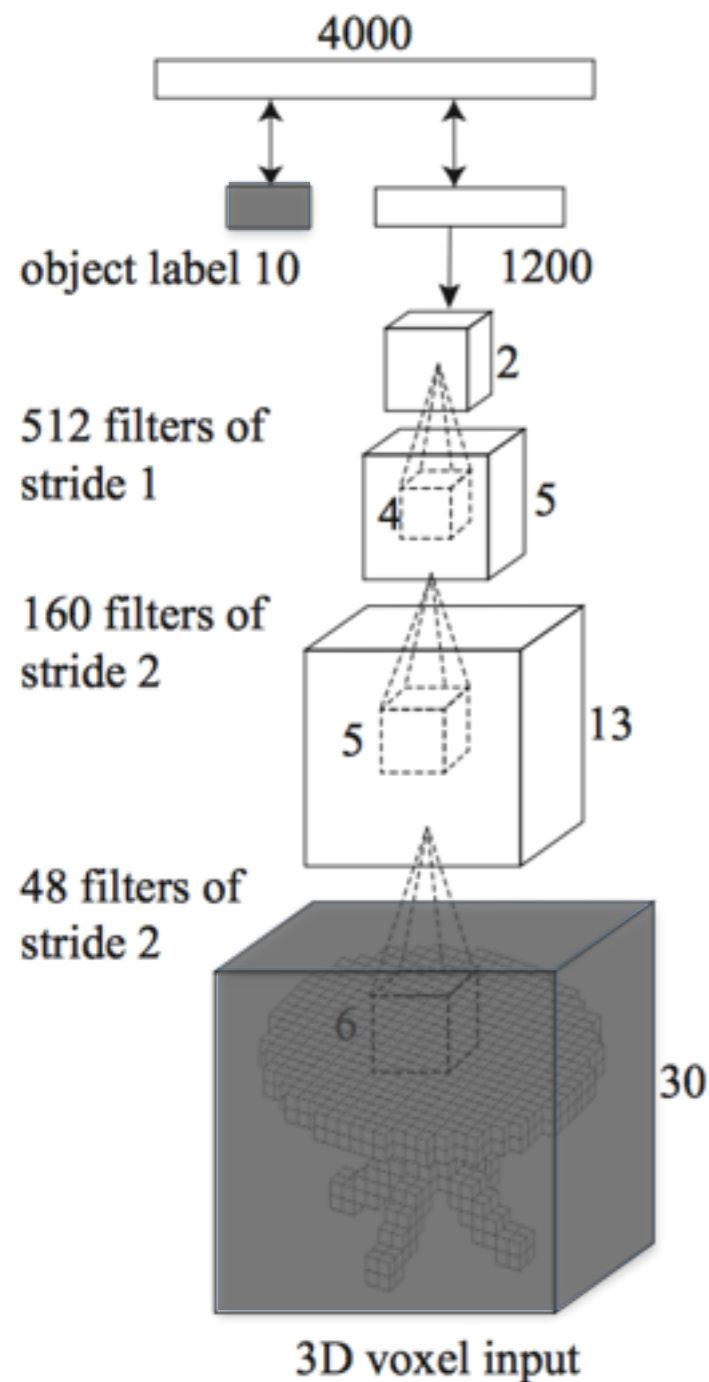


**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

A **Deep Belief Network** is a generative graphical model that describes the distribution of input \mathbf{x} over class y .

- Convolution to enable compositionality
- No pooling to reduce reconstruction error

3D ShapeNets



Convolutional Deep Belief Network $p(\mathbf{x}, y)$

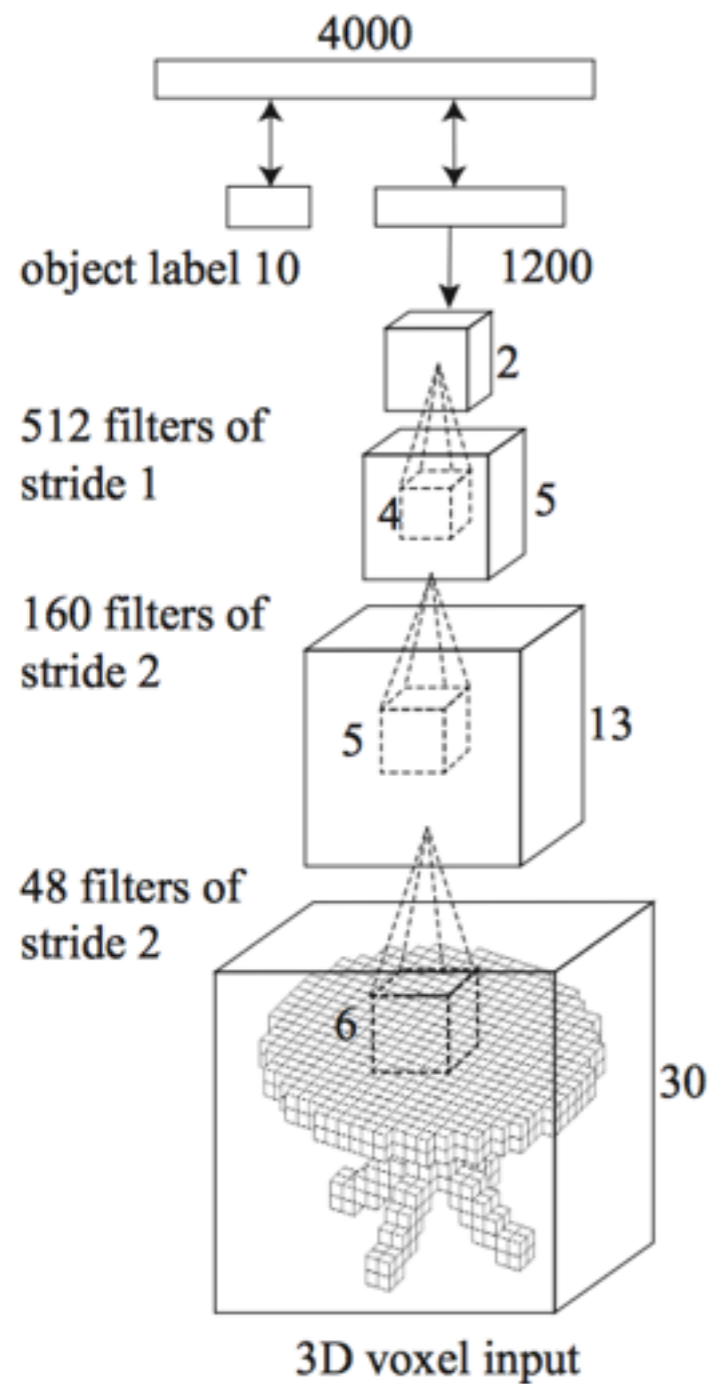
A **Deep Belief Network** is a generative graphical model that describes the distribution of input \mathbf{x} over class y .

- Convolution to enable compositionality
- No pooling to reduce reconstruction error

configurations

Layer 1-3	convolutional RBM
Layer 4	fully connected RBM
Layer 5	multinomial label + Bernoulli feature form an associate memory

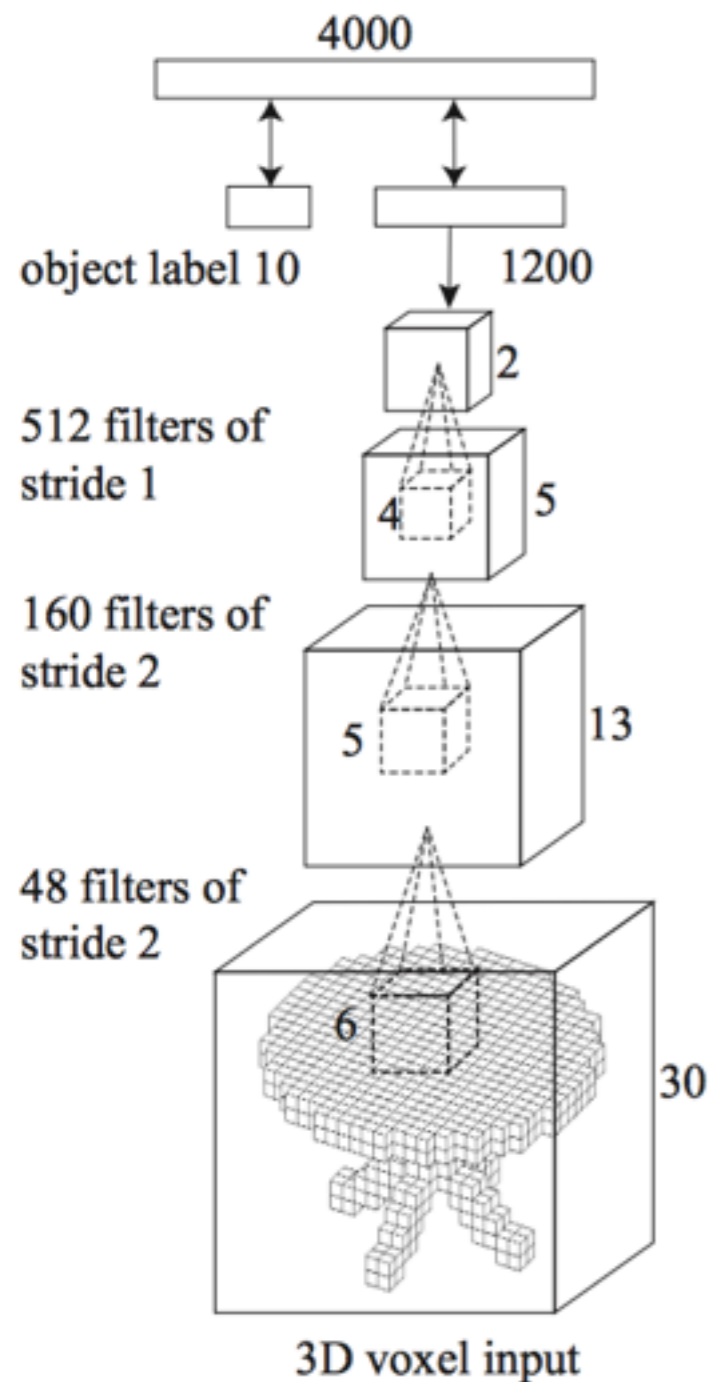
3D ShapeNets



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

3D ShapeNets

3D ShapeNets \neq CNNs

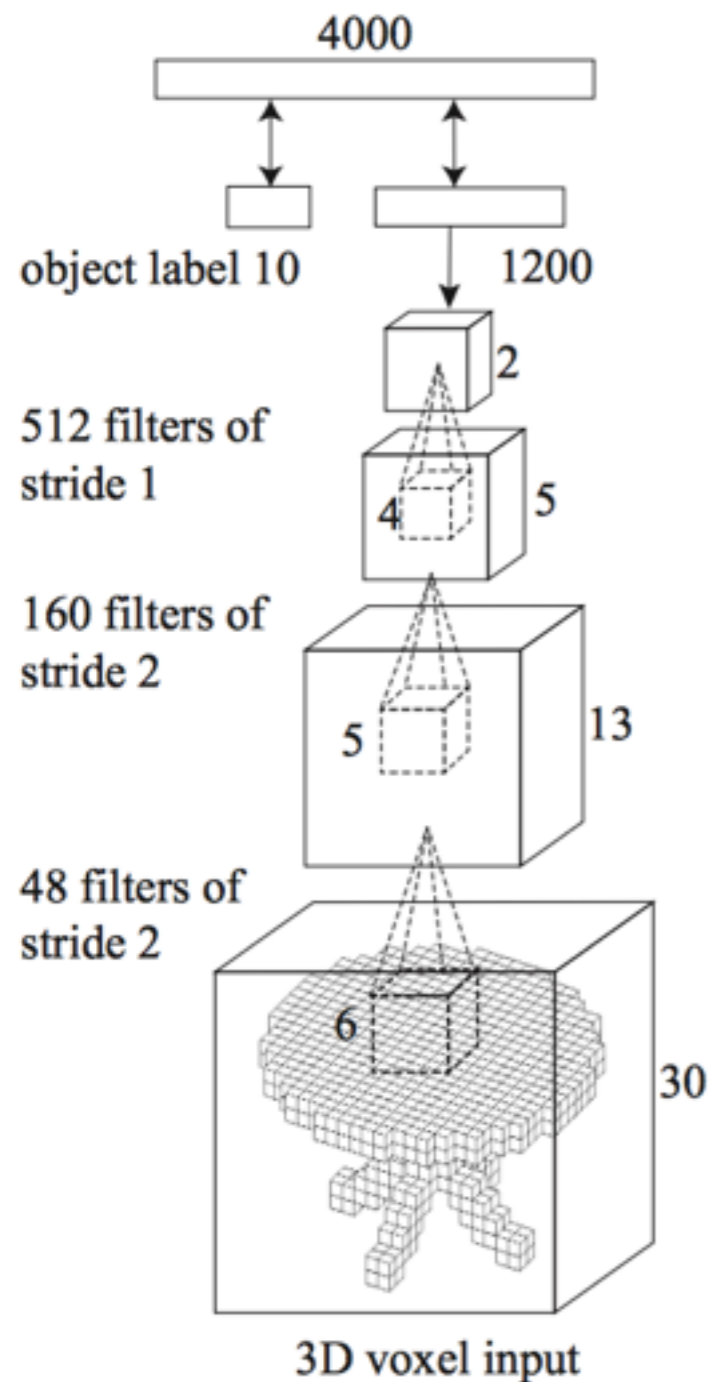


**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

3D ShapeNets

3D ShapeNets \neq CNNs

$$p(y|x)$$



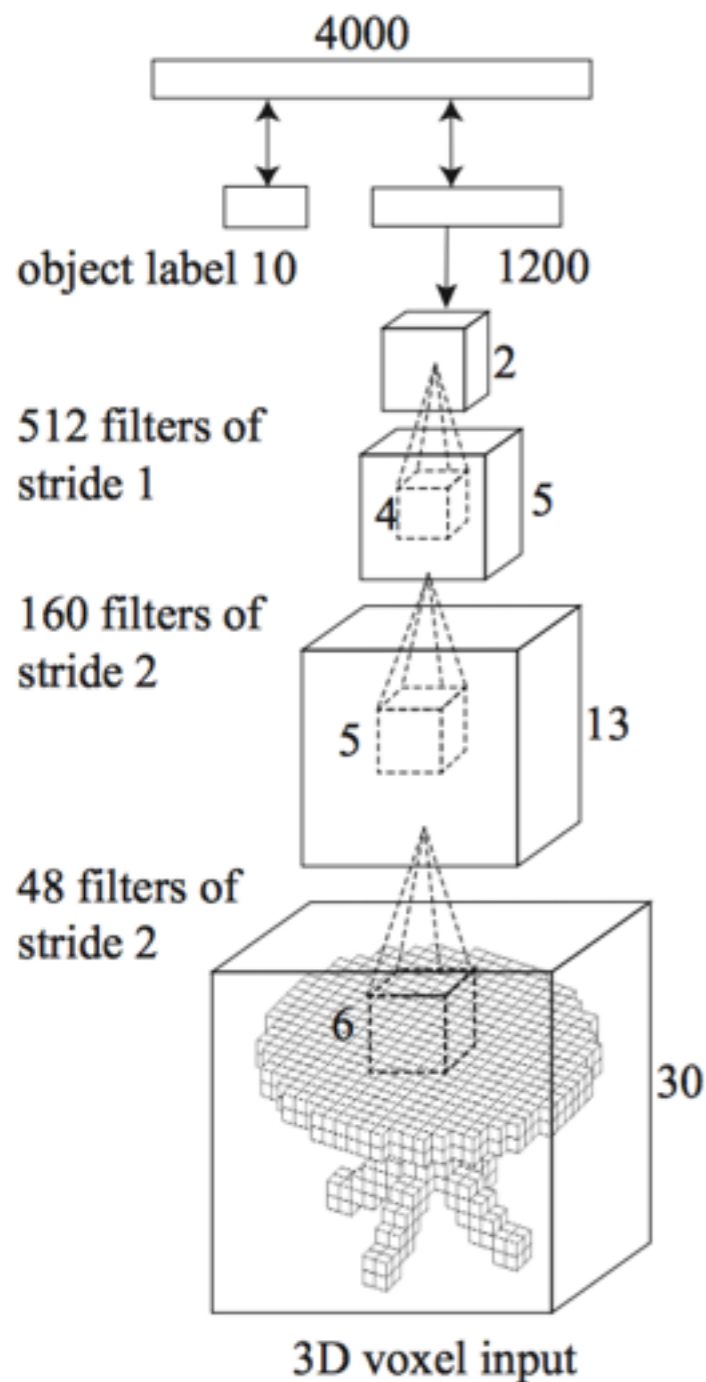
**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

3D ShapeNets

3D ShapeNets \neq CNNs

$$p(x, y)$$

$$p(y|x)$$



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

3D ShapeNets

3D ShapeNets \neq CNNs

$$p(x, y)$$

$$p(y|x)$$

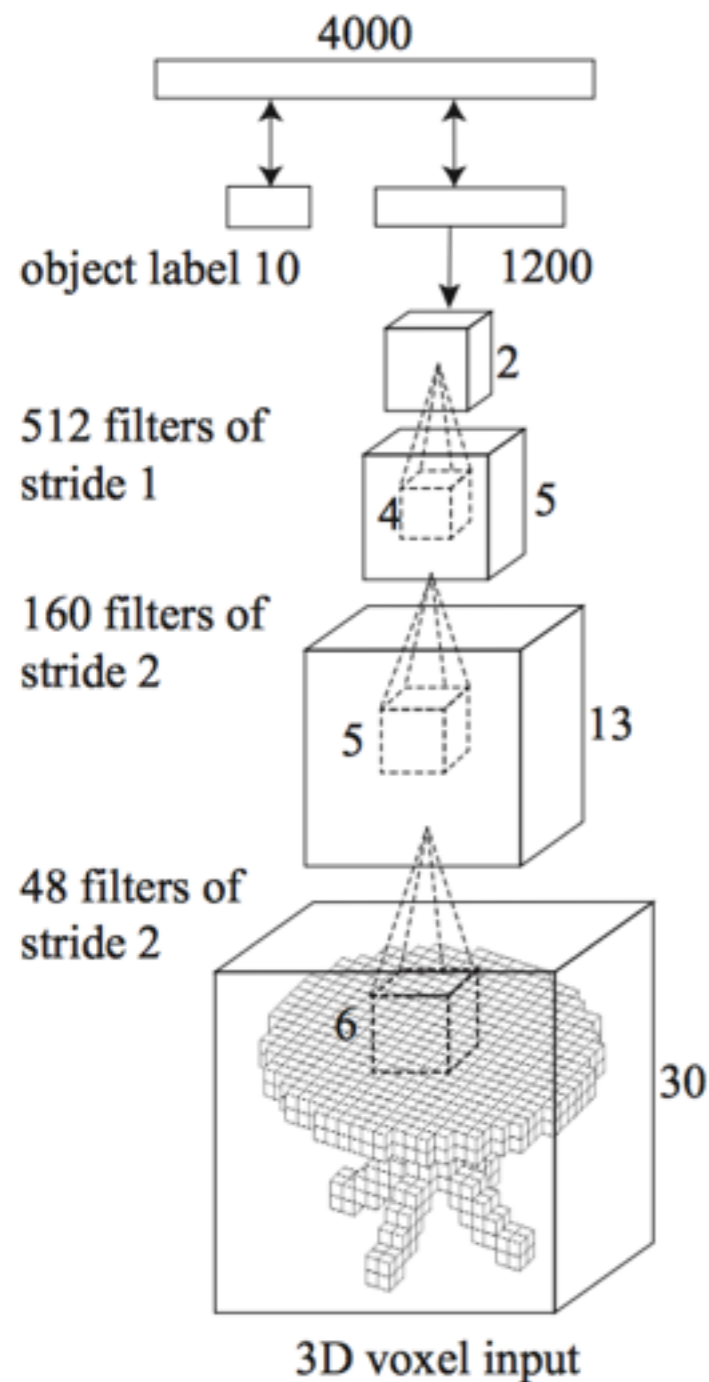


$$p(y|x)$$

discriminative process

$$p(x|y)$$

generative process



Convolutional Deep Belief Network $p(\mathbf{x}, y)$

3D ShapeNets

3D ShapeNets \neq CNNs

$$p(x, y)$$

$$p(y|x)$$

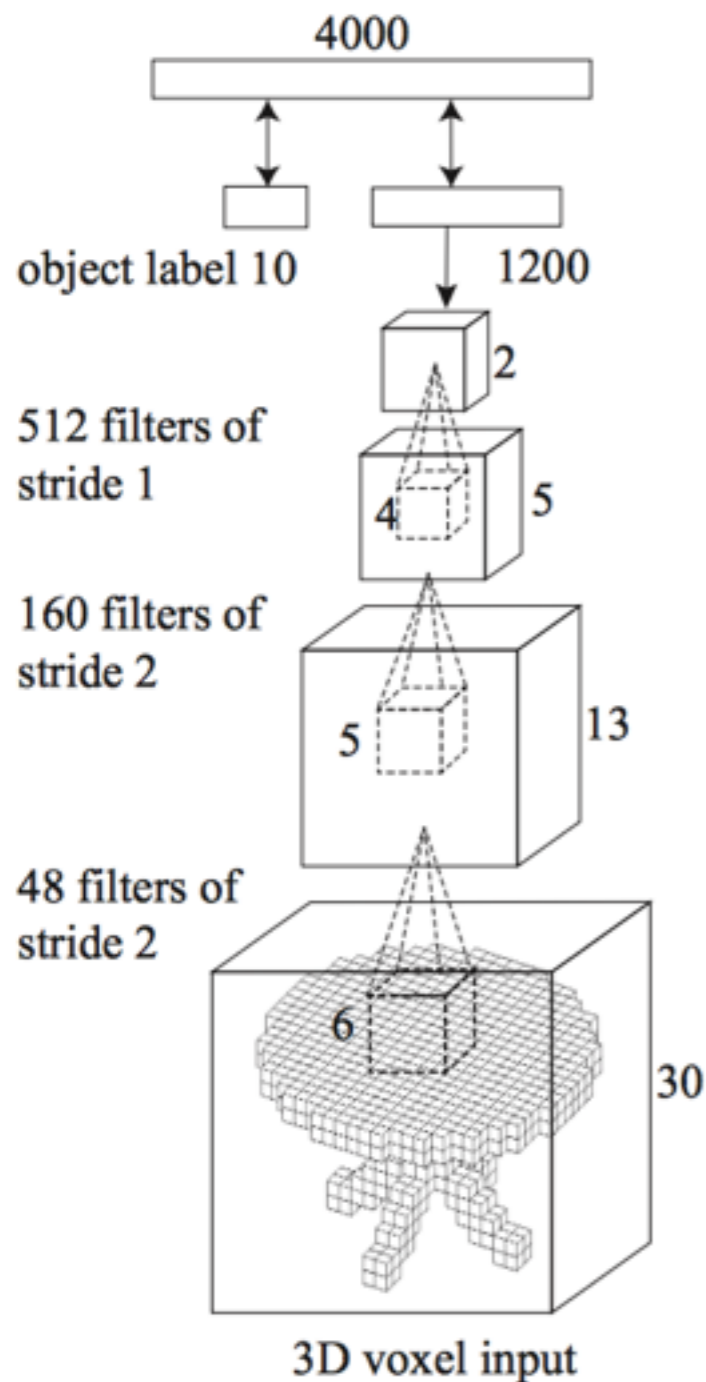


$$p(y|x)$$

discriminative process

$$p(x|y)$$

generative process

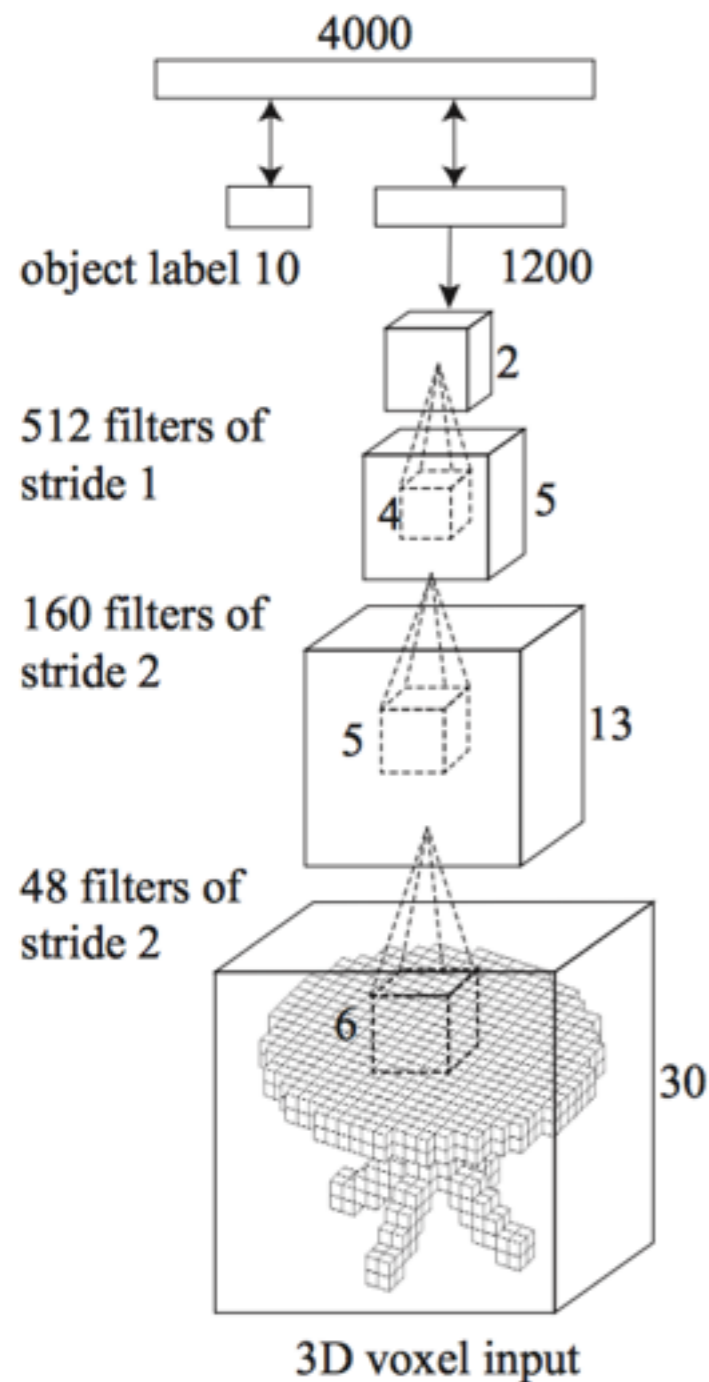


Convolutional Deep Belief Network $p(x, y)$

* 3D ShapeNets can be converted into a CNN, and discriminatively trained with back-propagation.

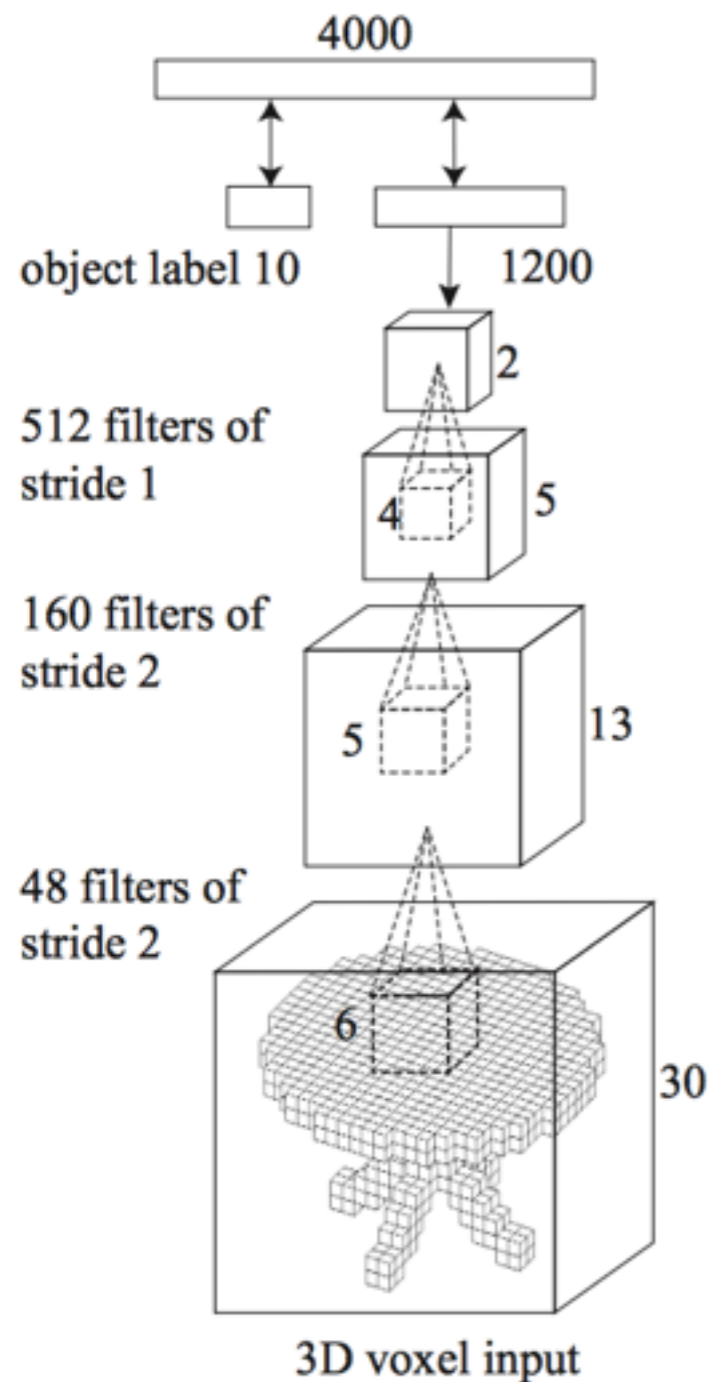
Training

Maximum Likelihood Learning



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

Training



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

Maximum Likelihood Learning

Layer-wise pre-training:

Lower four layers are trained by CD

Last layer is trained by FPCD[1]

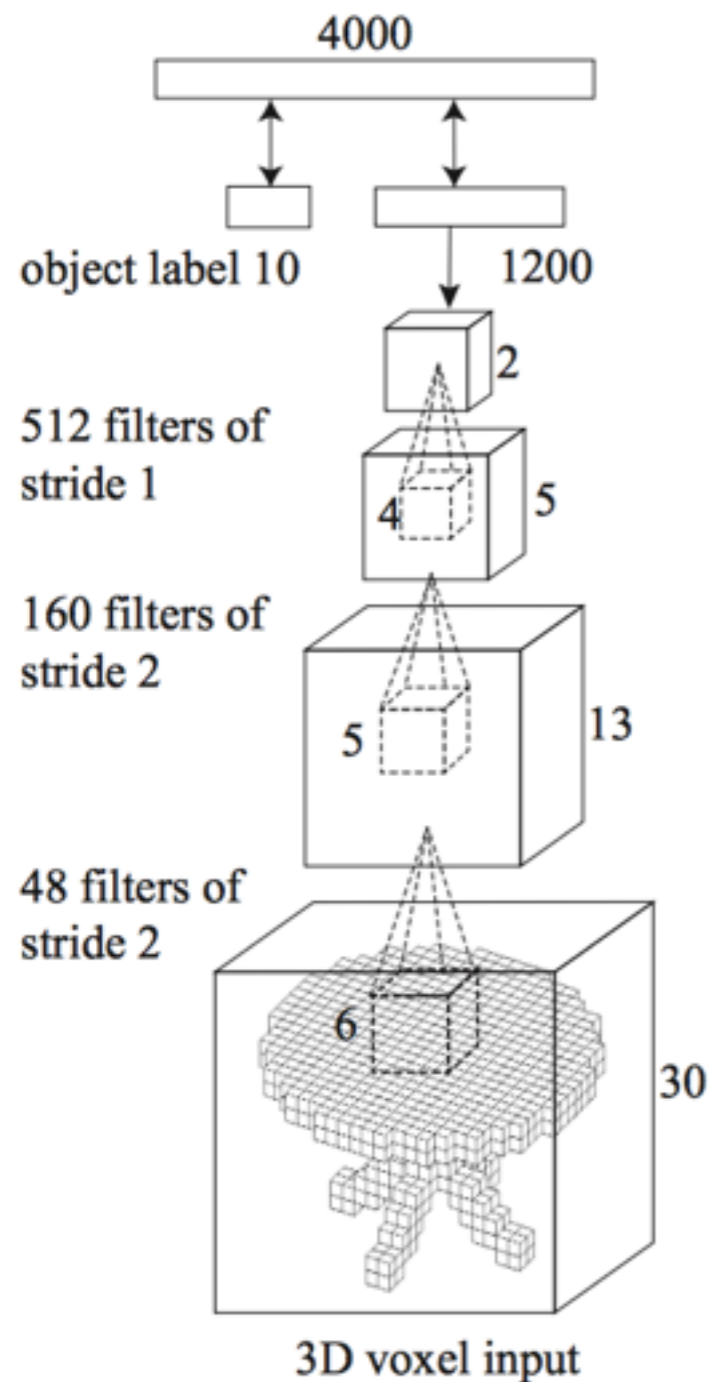
Fine-tuning:

Wake sleep[2] but keep weights tied

[1] Tijmen, et al. "Using fast weights to improve persistent contrastive divergence."

[2] Hinton, et al "A fast learning algorithm for deep belief nets." *Neural computation*

Training



**Convolutional Deep
Belief Network** $p(\mathbf{x}, y)$

Maximum Likelihood Learning

Layer-wise pre-training:

Lower four layers are trained by CD

Last layer is trained by FPCD[1]

Fine-tuning:

Wake sleep[2] but keep weights tied

[1] Tijmen, et al. "Using fast weights to improve persistent contrastive divergence."

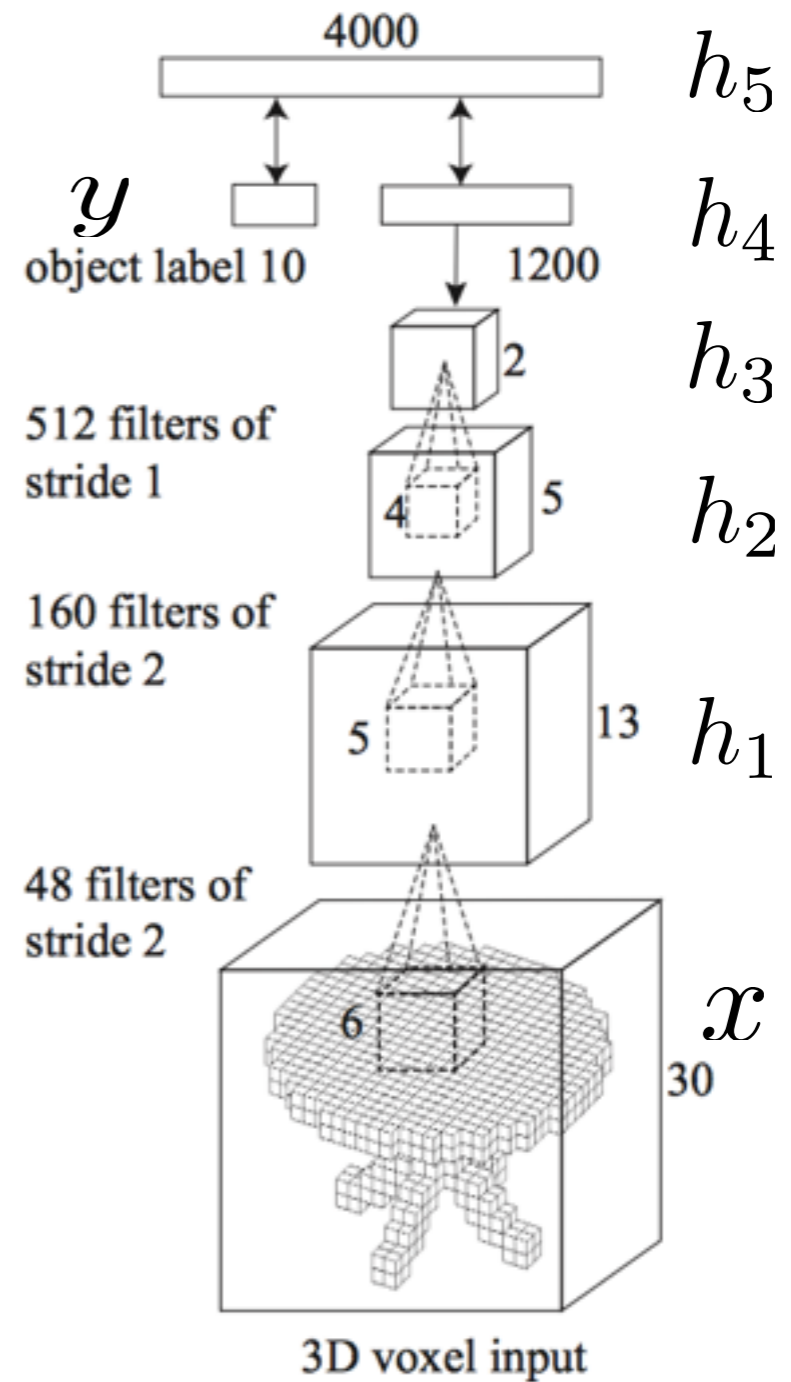
[2] Hinton, et al "A fast learning algorithm for deep belief nets." *Neural computation*

Sampling

generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$



Convolutional Deep Belief Network $p(x, y)$

Sampling

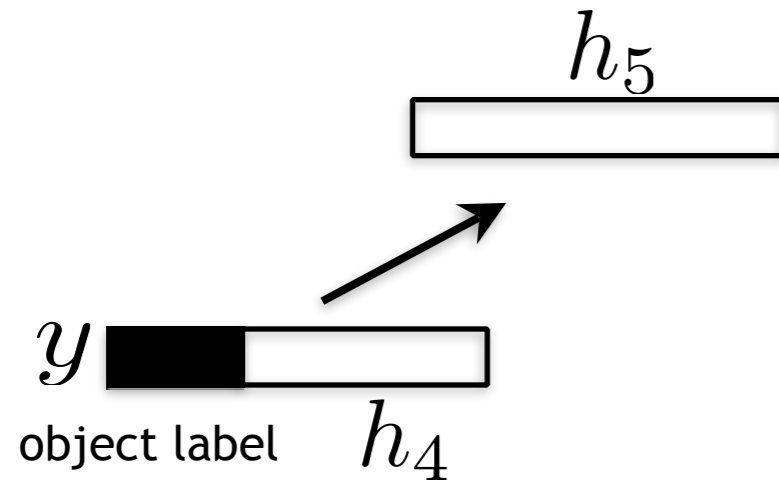
y 
object label h_4

generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$

Sampling

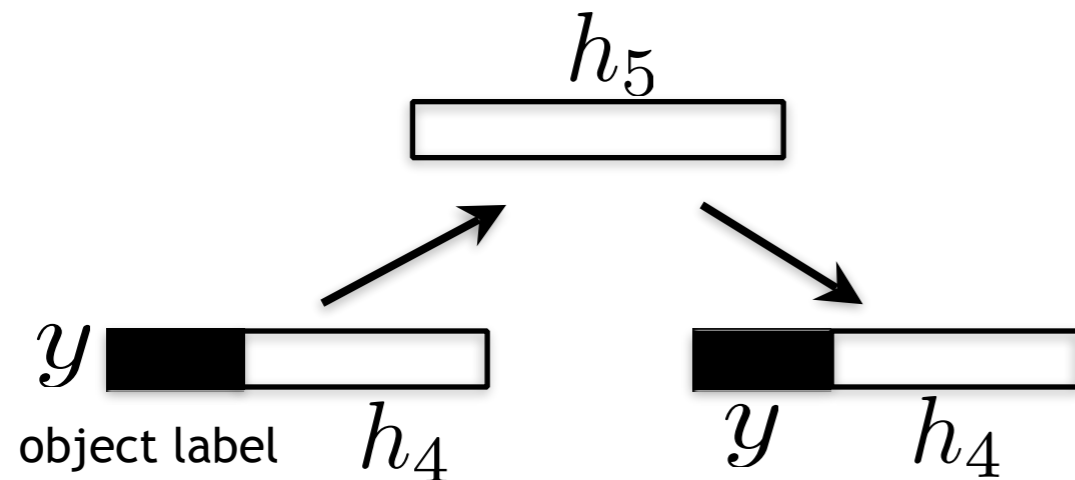


generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$

Sampling

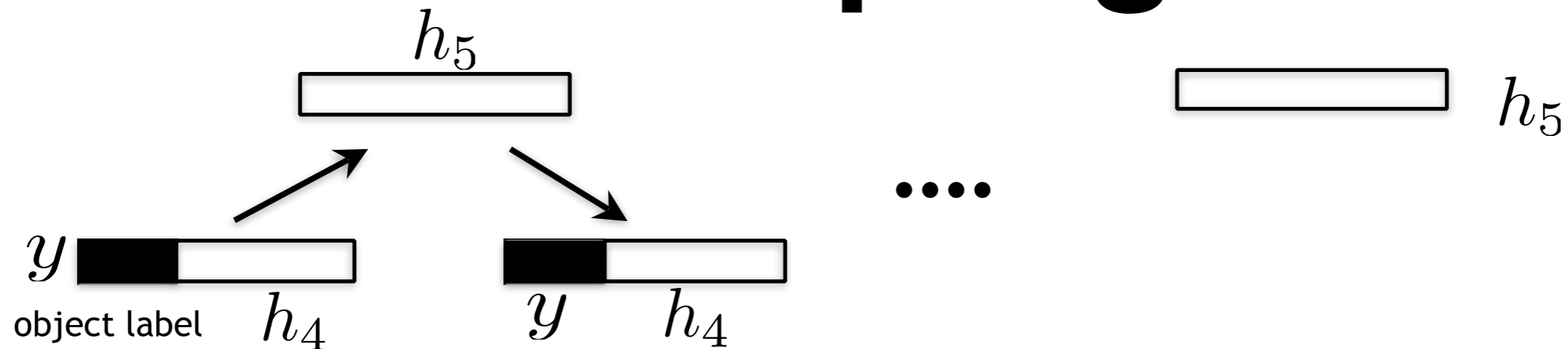


generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$

Sampling

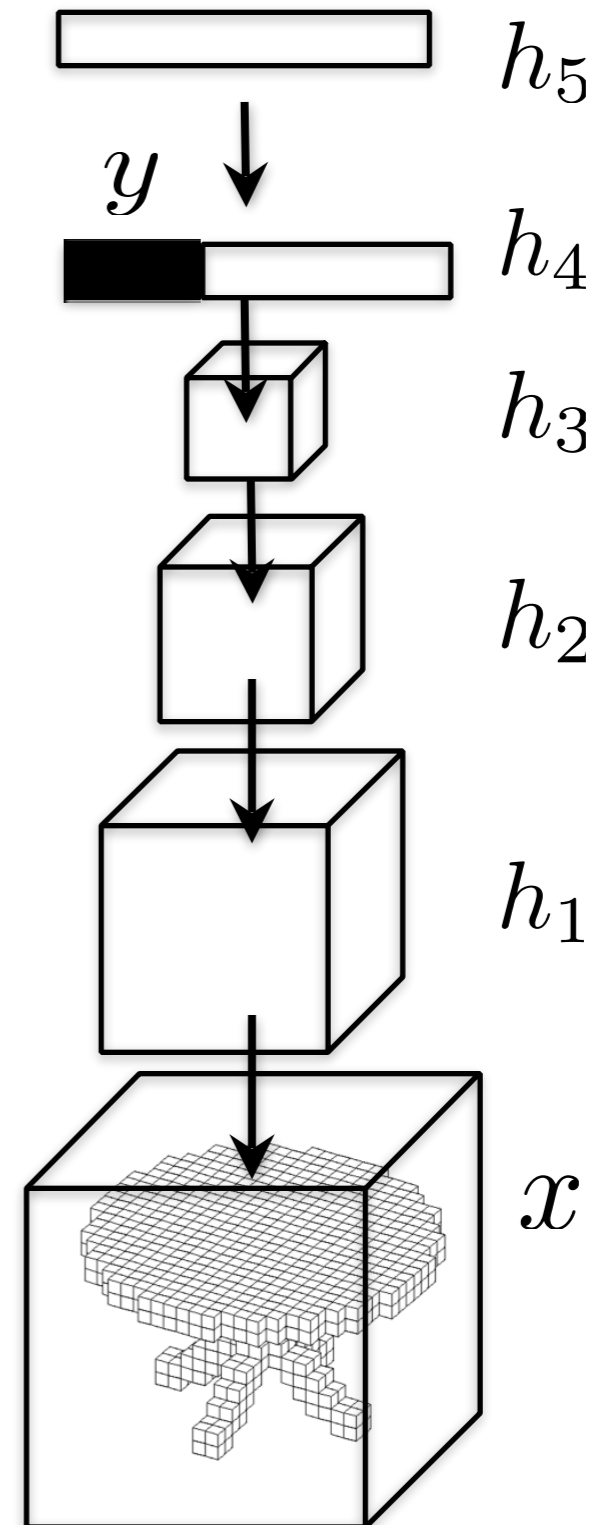
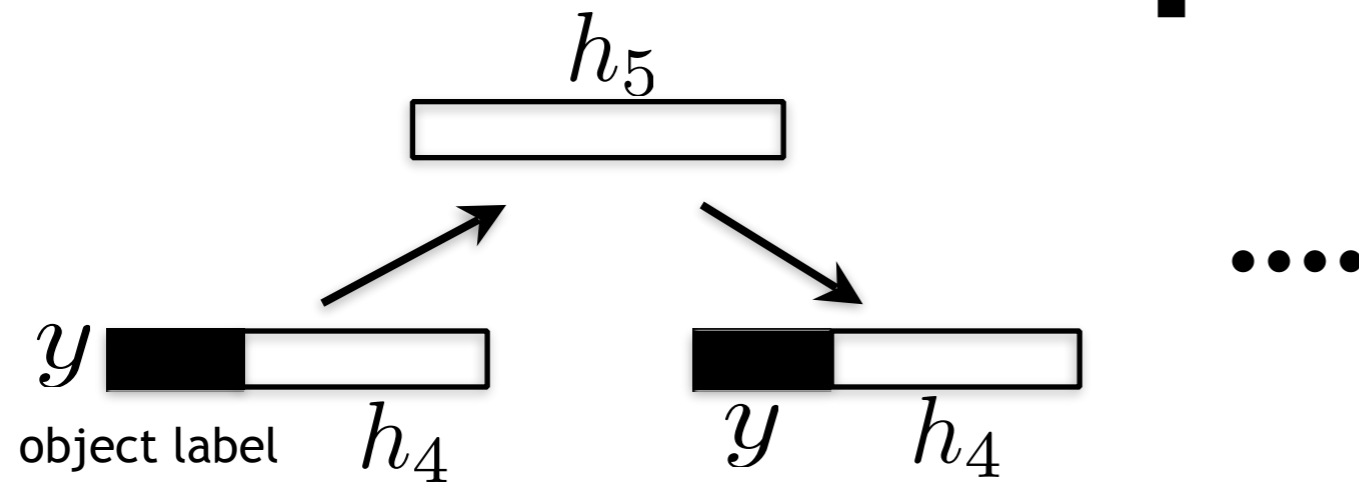


generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$

Sampling



generation process:

Gibbs Sampling

$$p(x, h_1 \dots h_5 | y) = p(h_4, h_5 | y) \cdot \prod_{i=1}^4 p(h_{i-1} | h_i)$$

Dataset

Big 3D Data

Big 3D Data



Query Keyword: common object categories from the SUN database that contain no less than 20 object instances per category



3D Warehouse



yobi 3D

Big 3D Data



Query Keyword: common object categories from the SUN database that contain no less than 20 object instances per category



3D Warehouse



yobi 3D

Instruction **Is this a chair?** Submit (128 images left)

Definition: A separate seat for one person, typically with a back and four legs..



Big 3D Data



151,128 models

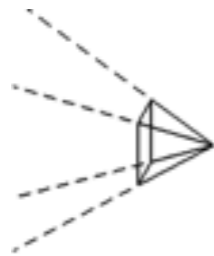
660 categories

Applications

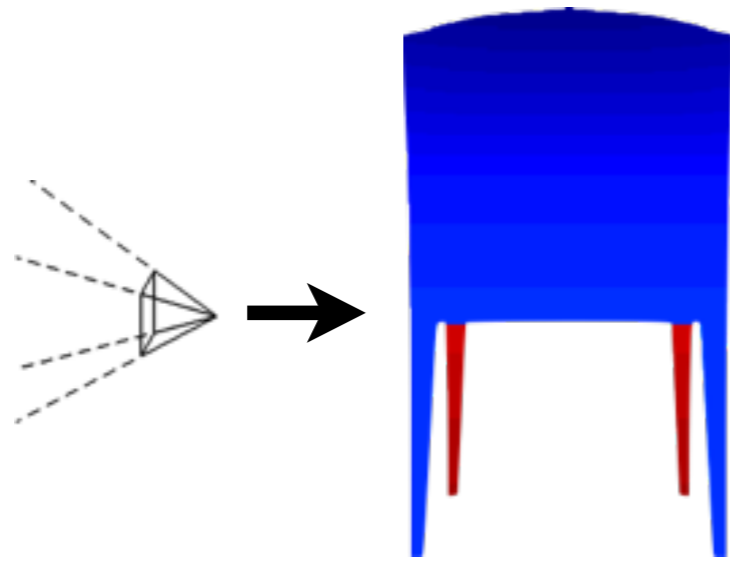
2.5D Completion & Recognition

2.5D Completion & Recognition

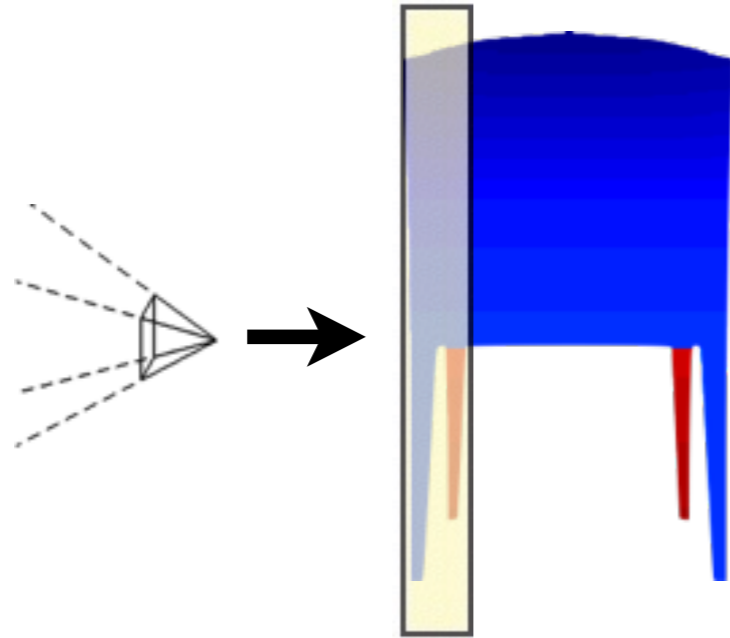
2.5D Completion & Recognition



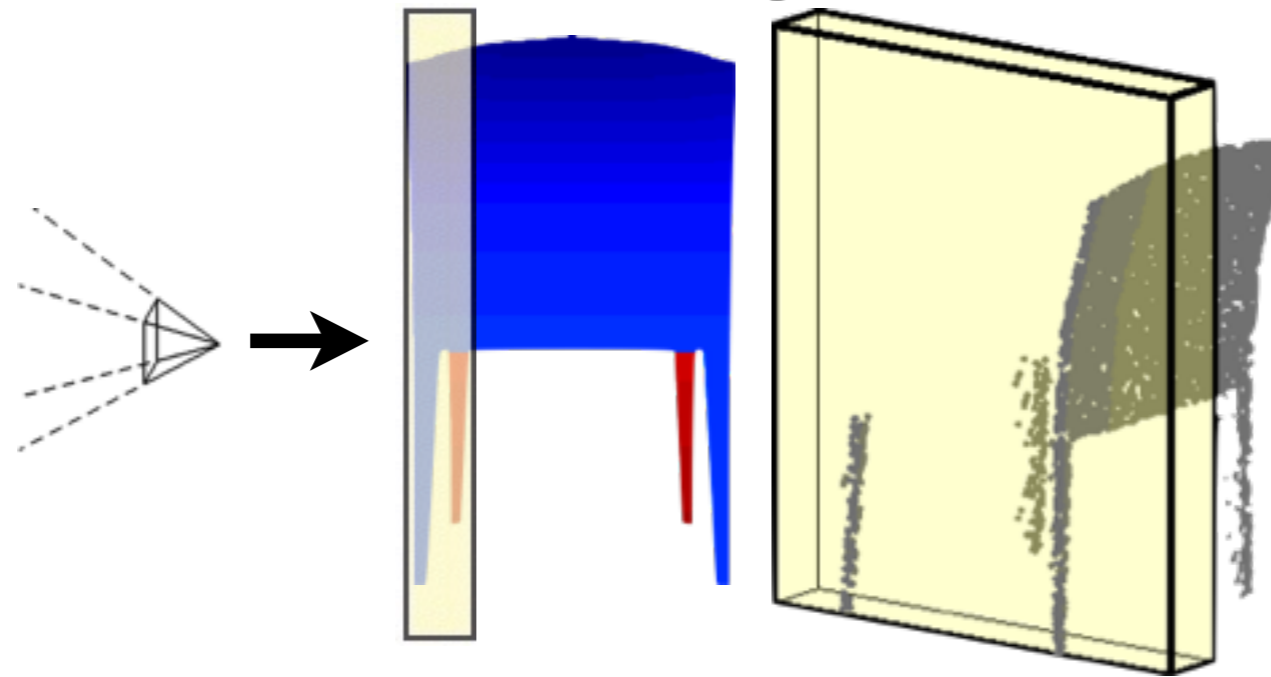
2.5D Completion & Recognition



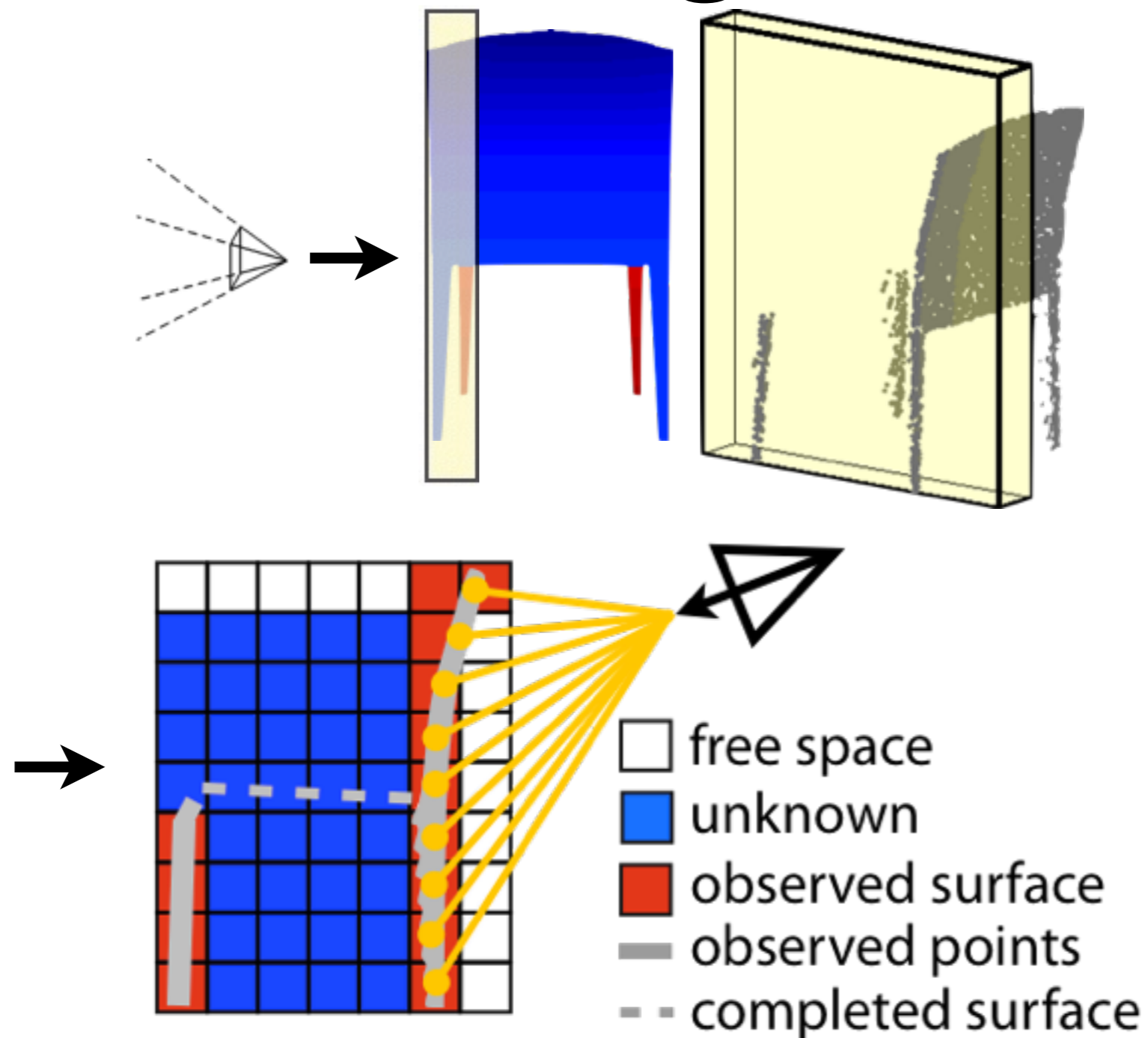
2.5D Completion & Recognition



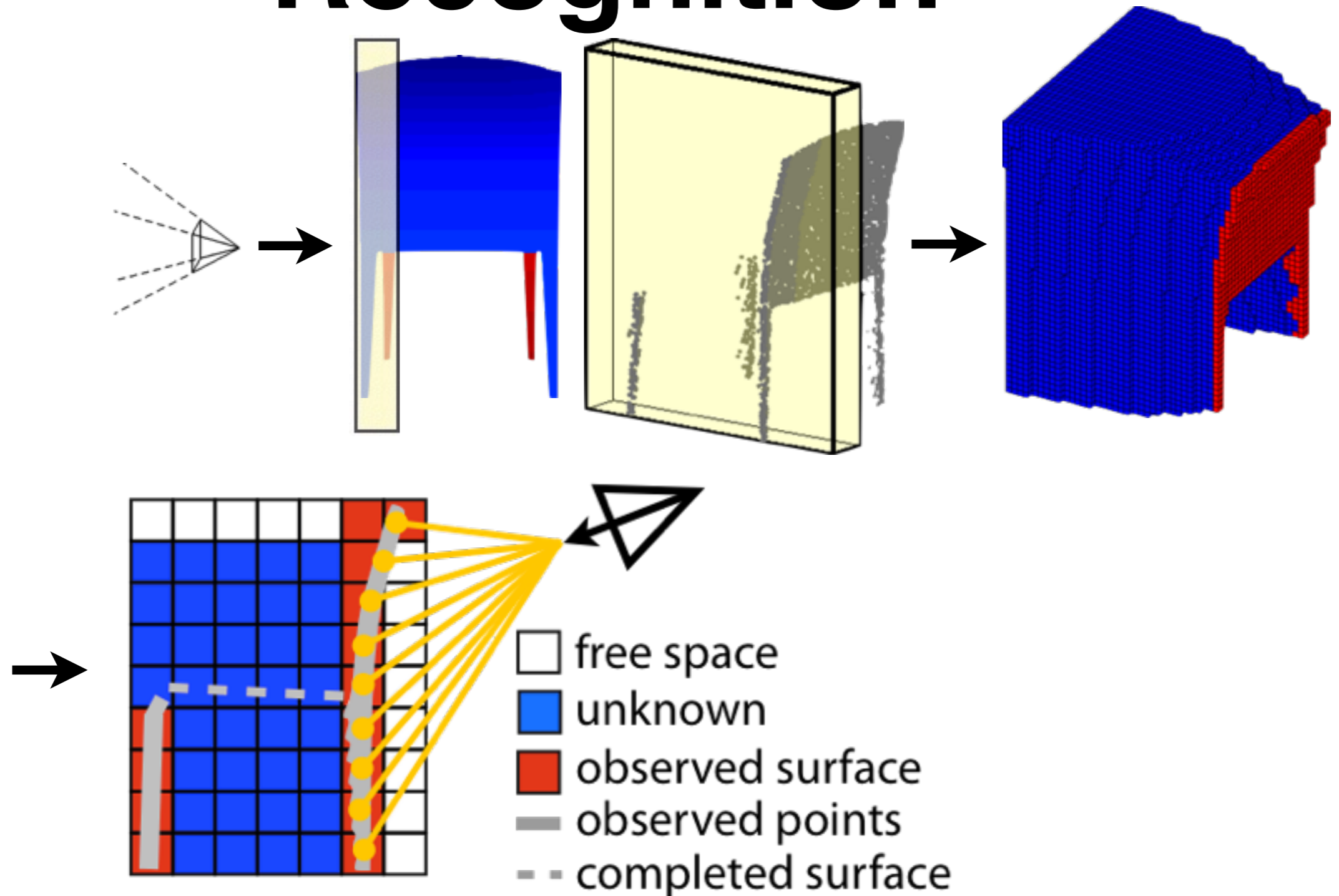
2.5D Completion & Recognition



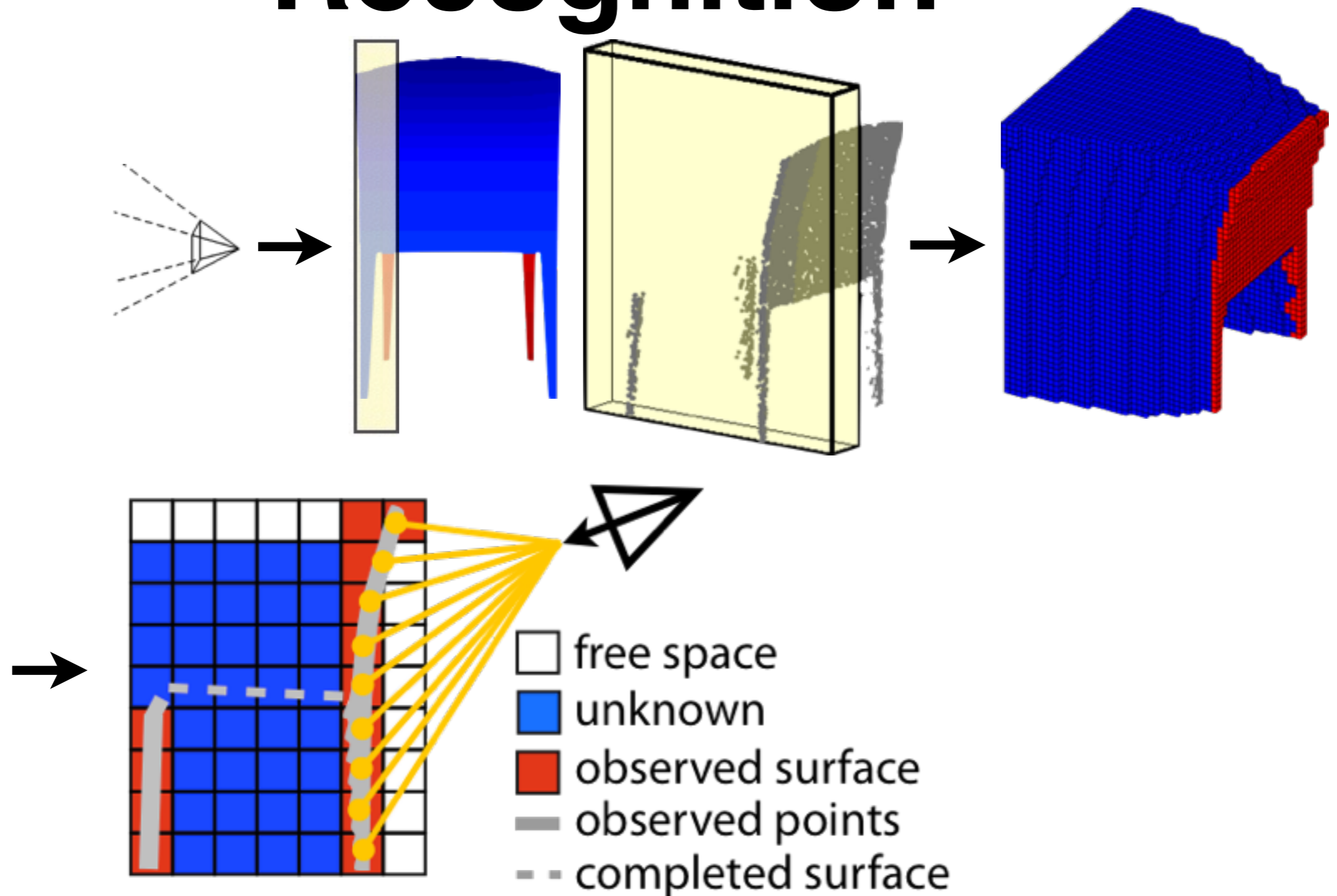
2.5D Completion & Recognition



2.5D Completion & Recognition

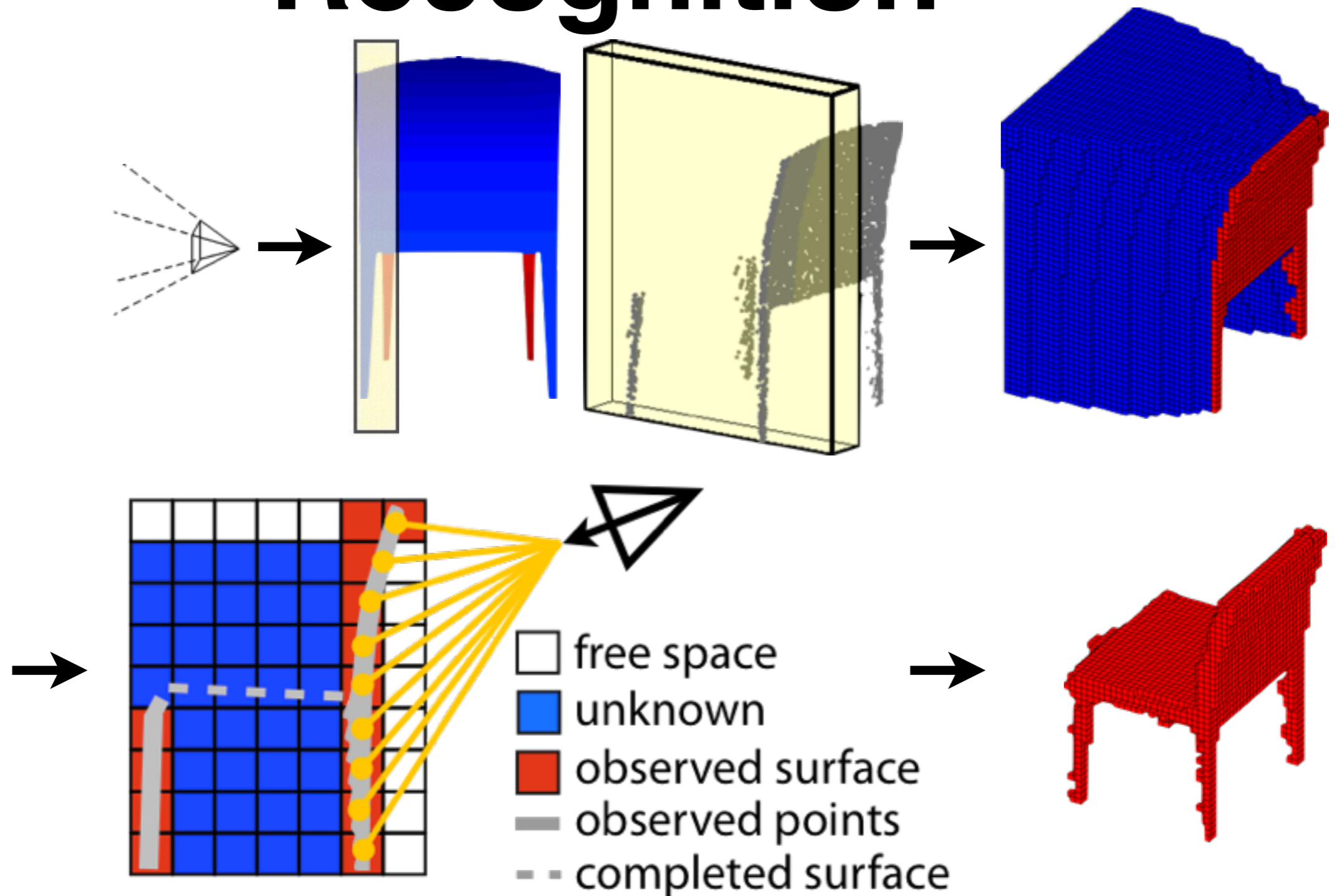


2.5D Completion & Recognition



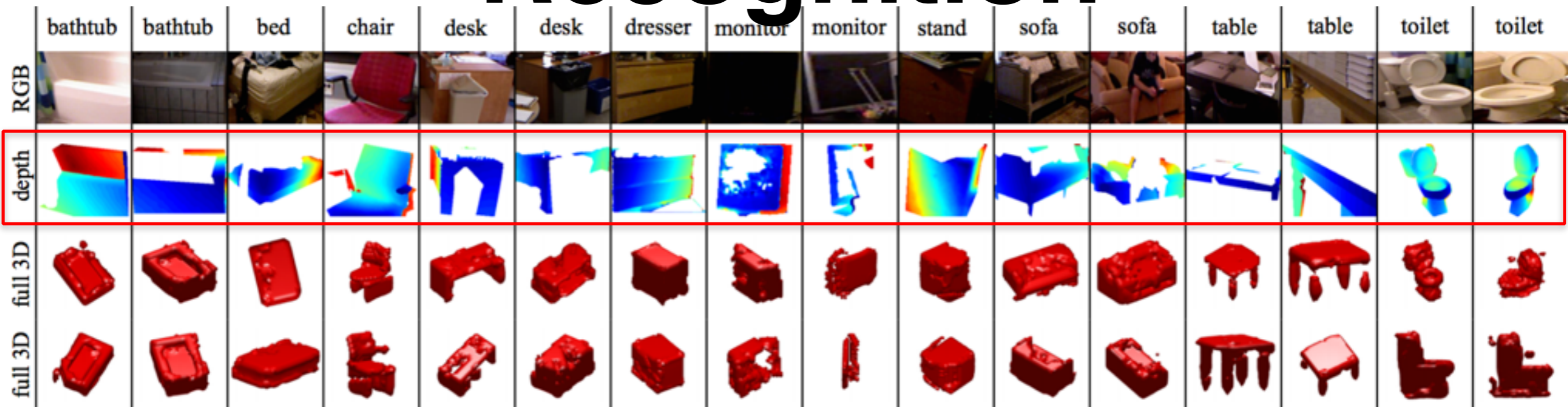
$\mathbf{x} = (\mathbf{x}_u, \mathbf{x}_o)$ $p(y, \mathbf{x}_u | \mathbf{x}_o)$ Gibbs sampling with clamping

2.5D Completion & Recognition



$\mathbf{x} = (\mathbf{x}_u, \mathbf{x}_o)$ $p(y, \mathbf{x}_u | \mathbf{x}_o)$ Gibbs sampling with clamping

2.5D Completion & Recognition



Training on CAD models and no discriminative tuning!

	all
[29] Depth	0.376
NN	0.374
ICP	0.471
3D ShapeNets	0.437
3D ShapeNets fine-tuned	0.579
[29] RGB	0.334
[29] RGBD	0.448

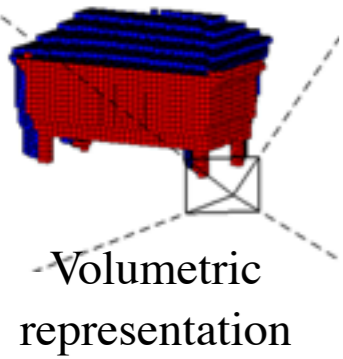
[29] R. Socher, B. Huval, B. Bhat, C. D. Manning, and A. Y. Ng.

Convolutional-recursive deep learning for 3d object classification. In NIPS 2012.

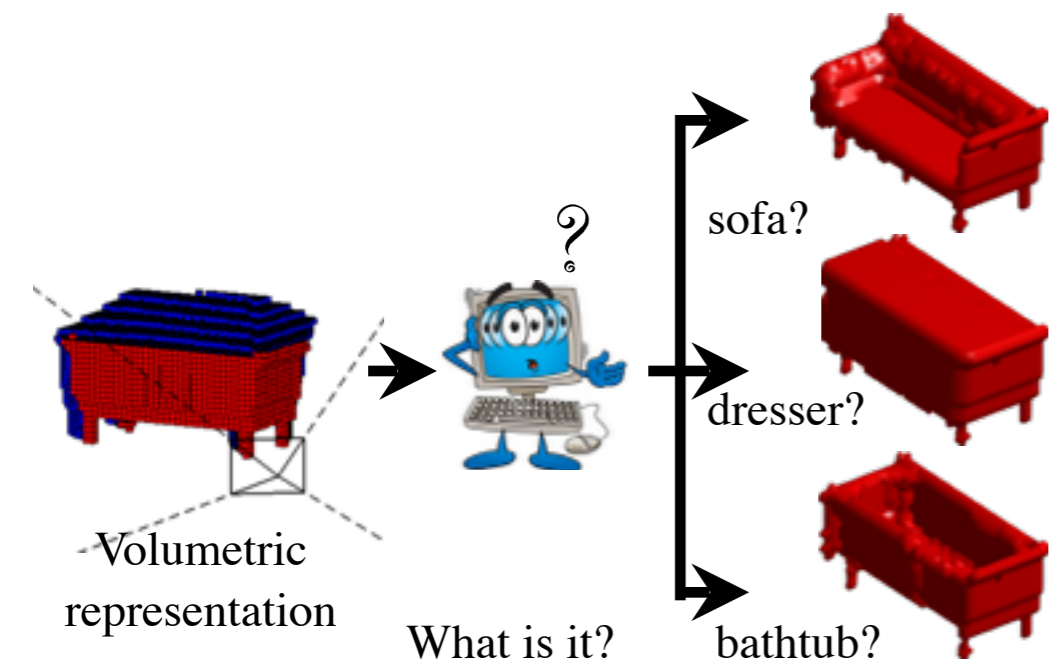
Slide Credit: Wu, Song et al. 3D ShapeNets: A Deep Representation for Volumetric Shape Modeling, CVPR 2015

View Planning for Recognition

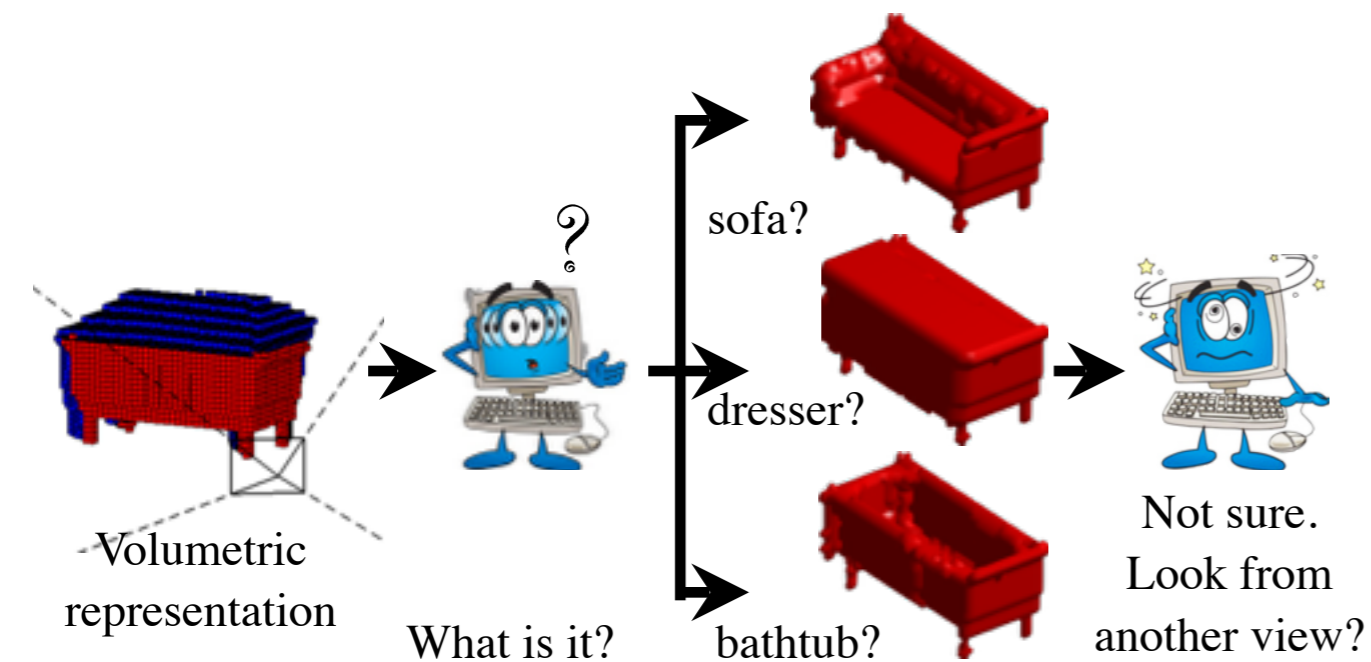
View Planning for Recognition



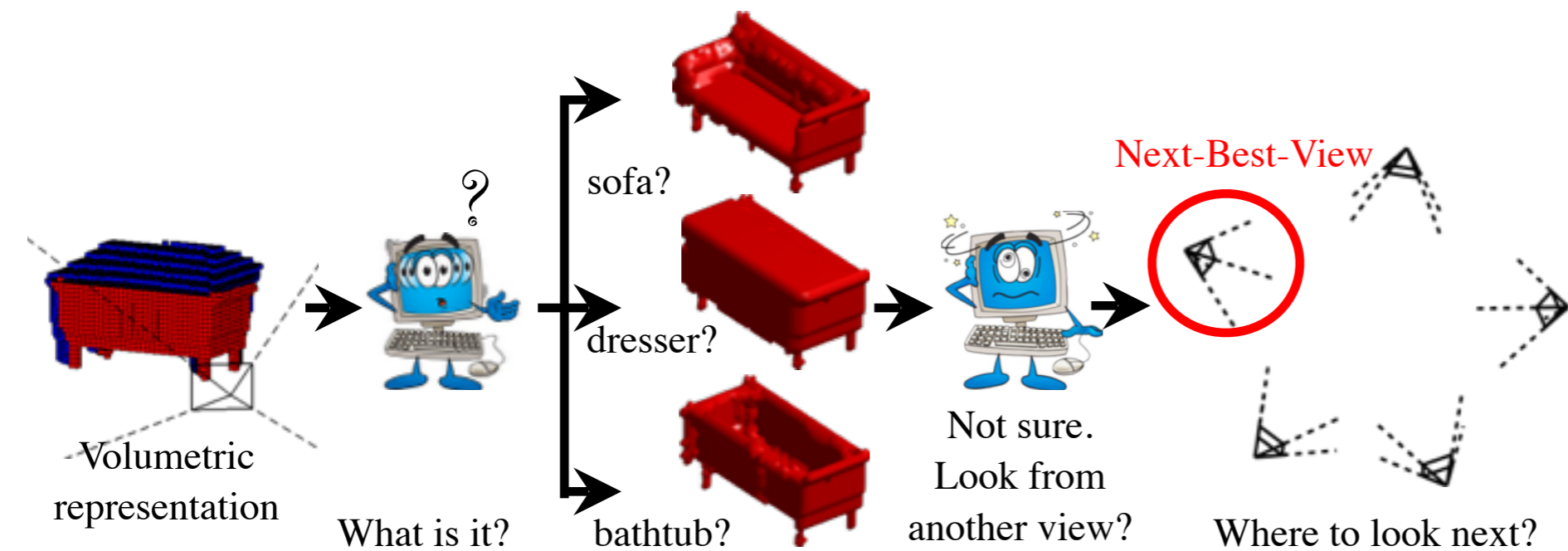
View Planning for Recognition



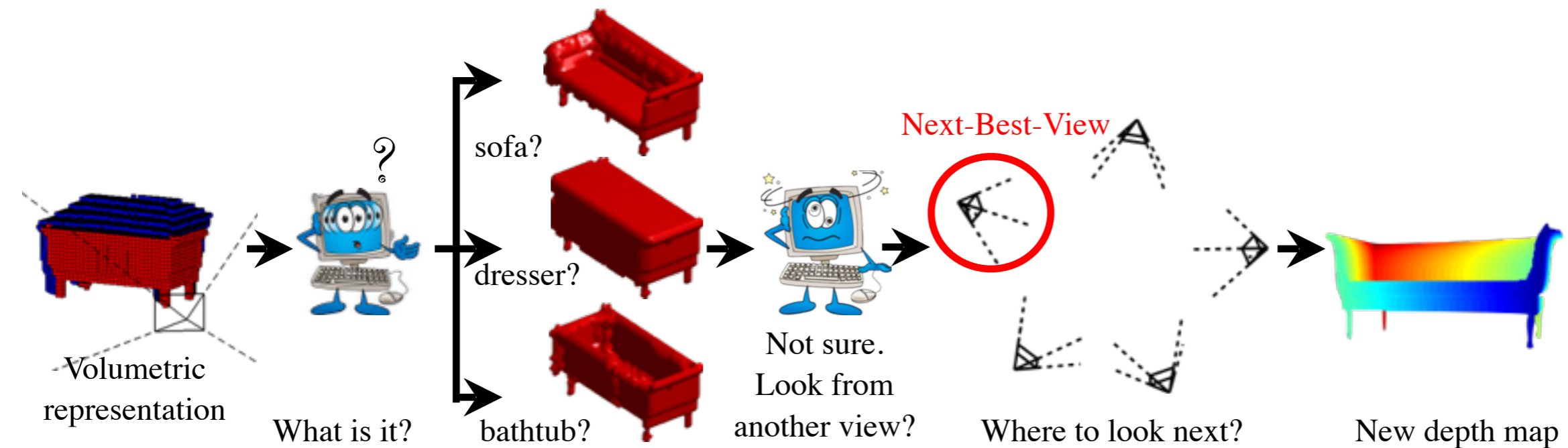
View Planning for Recognition



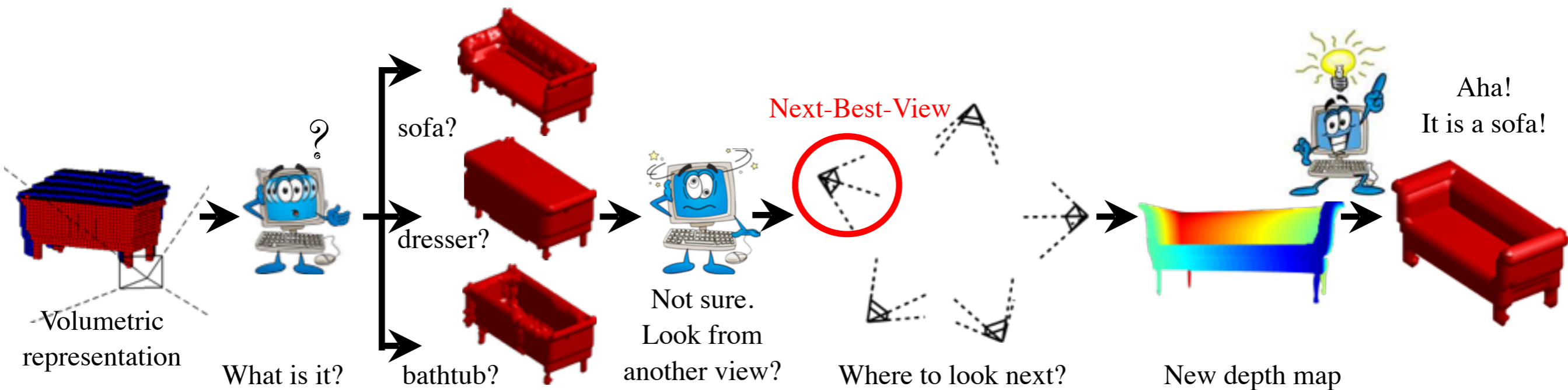
View Planning for Recognition



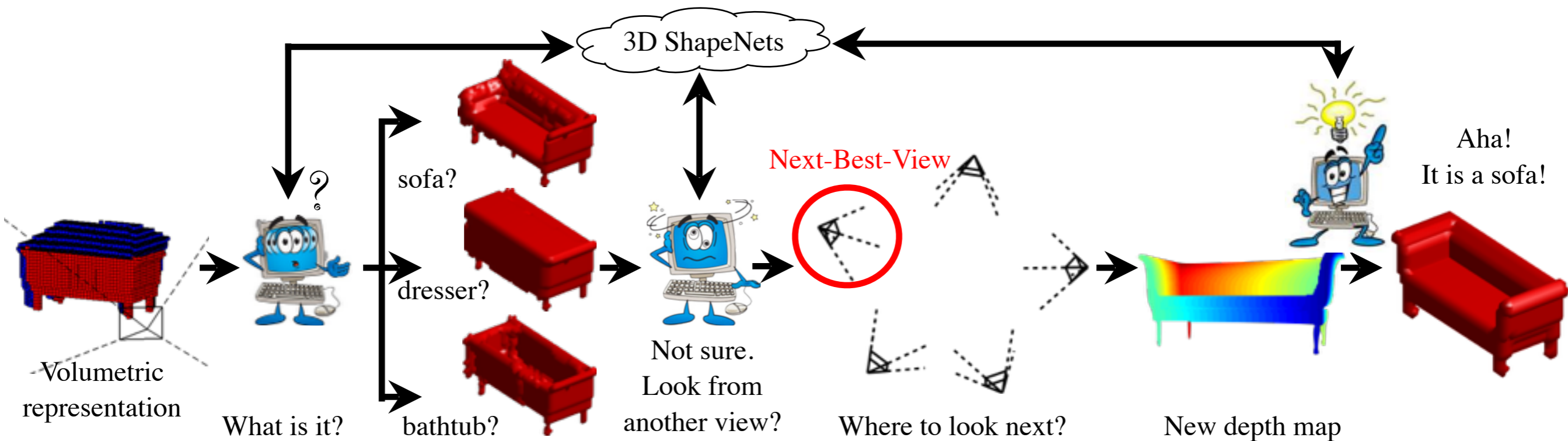
View Planning for Recognition



View Planning for Recognition



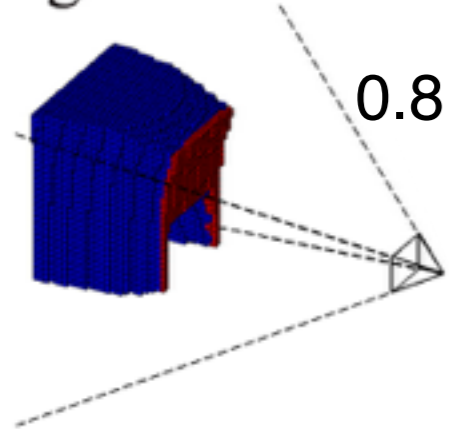
View Planning for Recognition



Deep View Planning

■ observed surface \mathbf{x}_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface \mathbf{x}_n^i □ free space

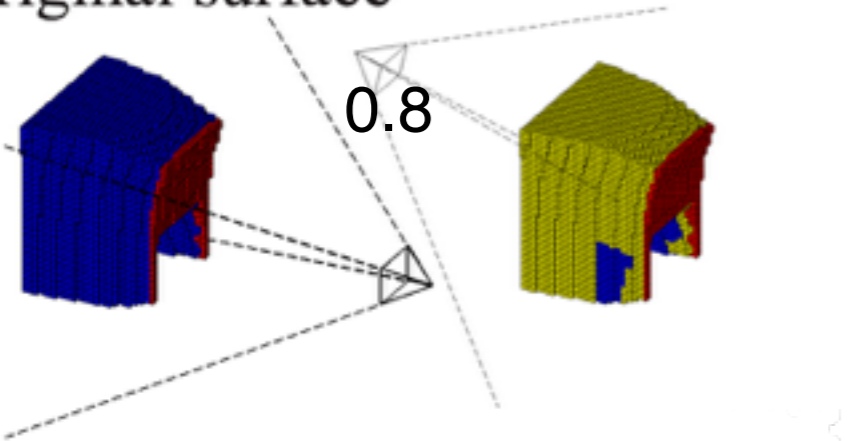
original surface



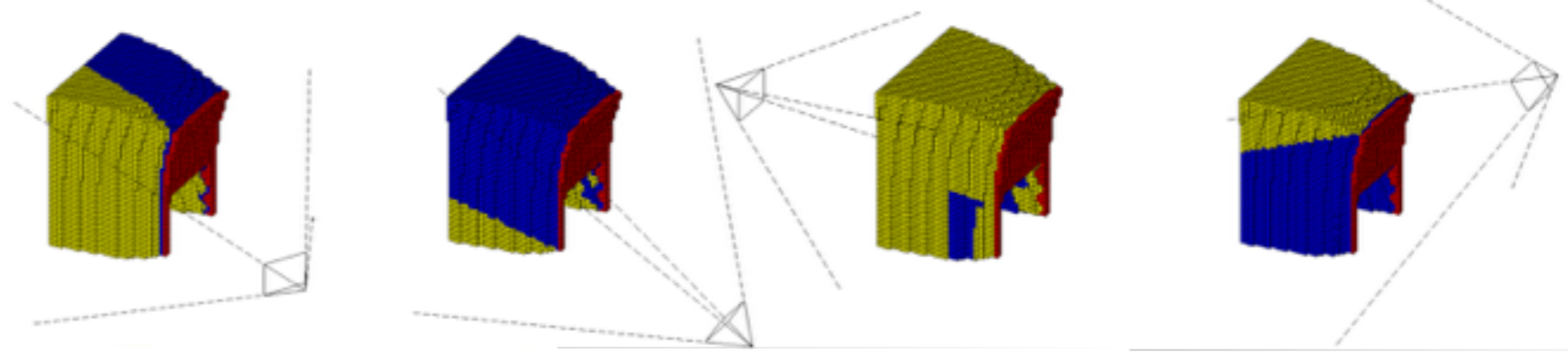
Deep View Planning

■ observed surface x_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface x_n^i □ free space

original surface



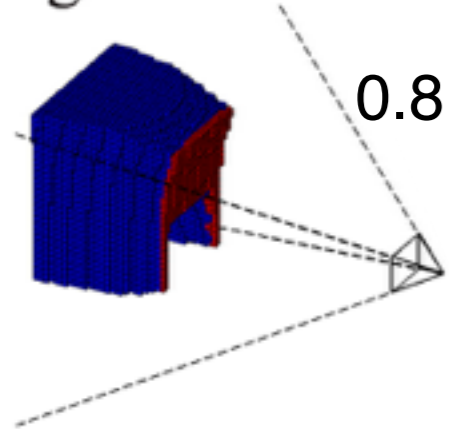
five different next-view candidates



Deep View Planning

■ observed surface \mathbf{x}_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface \mathbf{x}_n^i □ free space

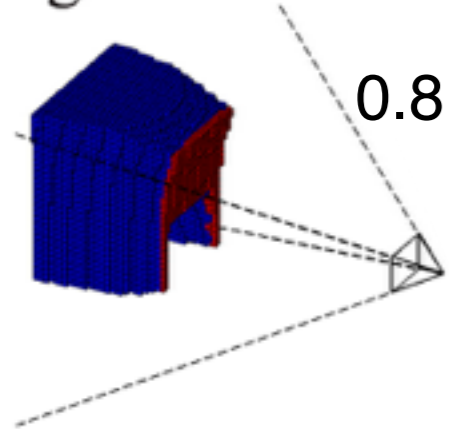
original surface



Deep View Planning

■ observed surface x_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface x_n^i □ free space

original surface



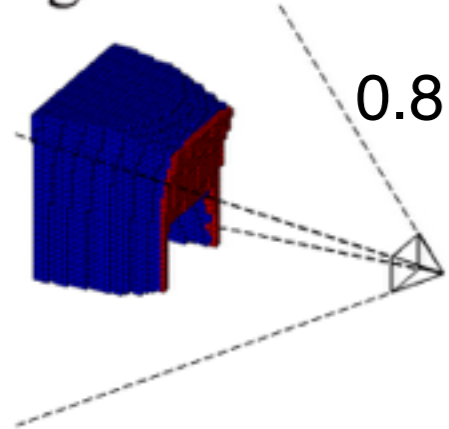
3 possible shapes



Deep View Planning

■ observed surface \mathbf{x}_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface \mathbf{x}_n^i □ free space

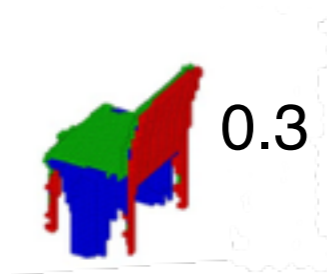
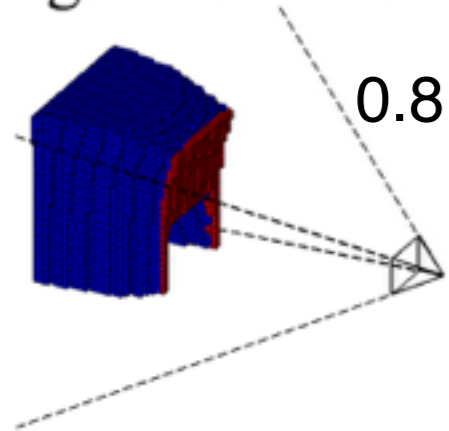
original surface



Deep View Planning

■ observed surface x_o ■ unknown ■ potentially visible voxels in next view ■ newly visible surface x_n^i □ free space

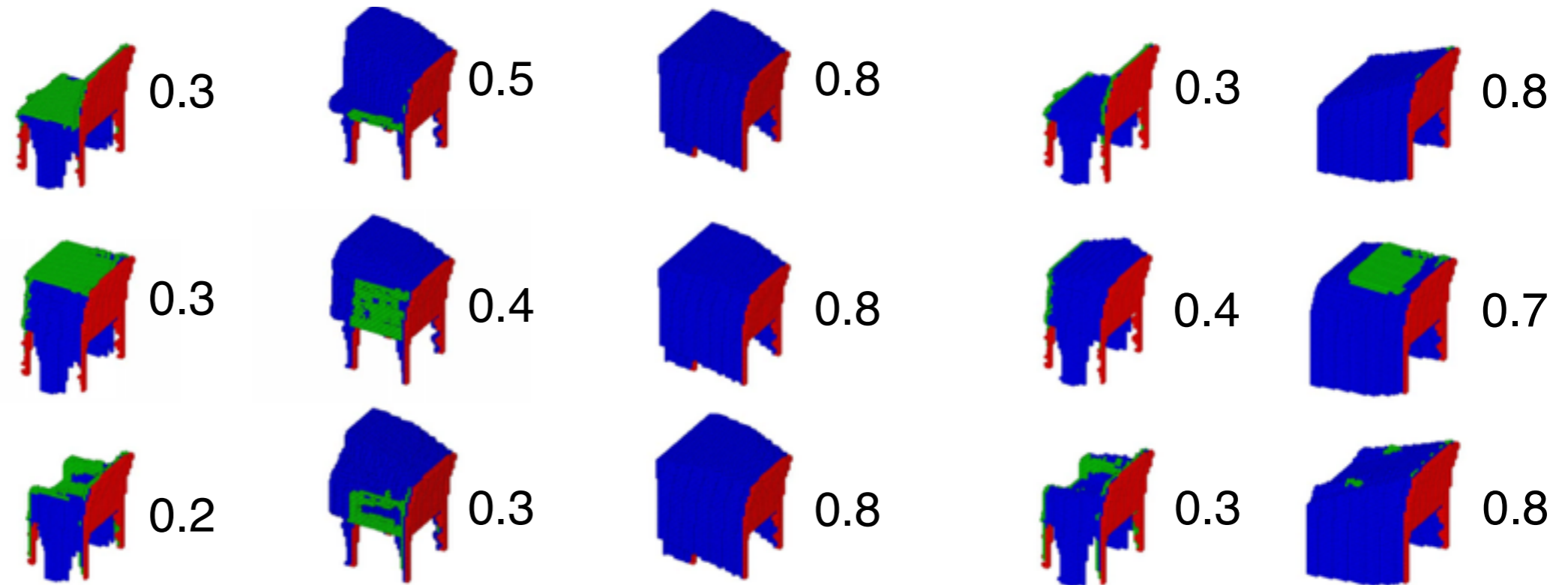
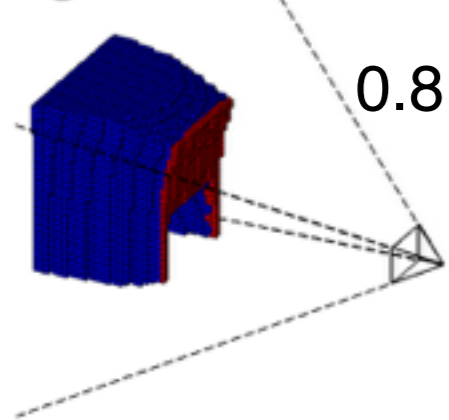
original surface



Deep View Planning

■ observed surface \mathbf{x}_o
■ unknown
 ■ potentially visible voxels in next view
 ■ newly visible surface \mathbf{x}_n^i
 free space

original surface

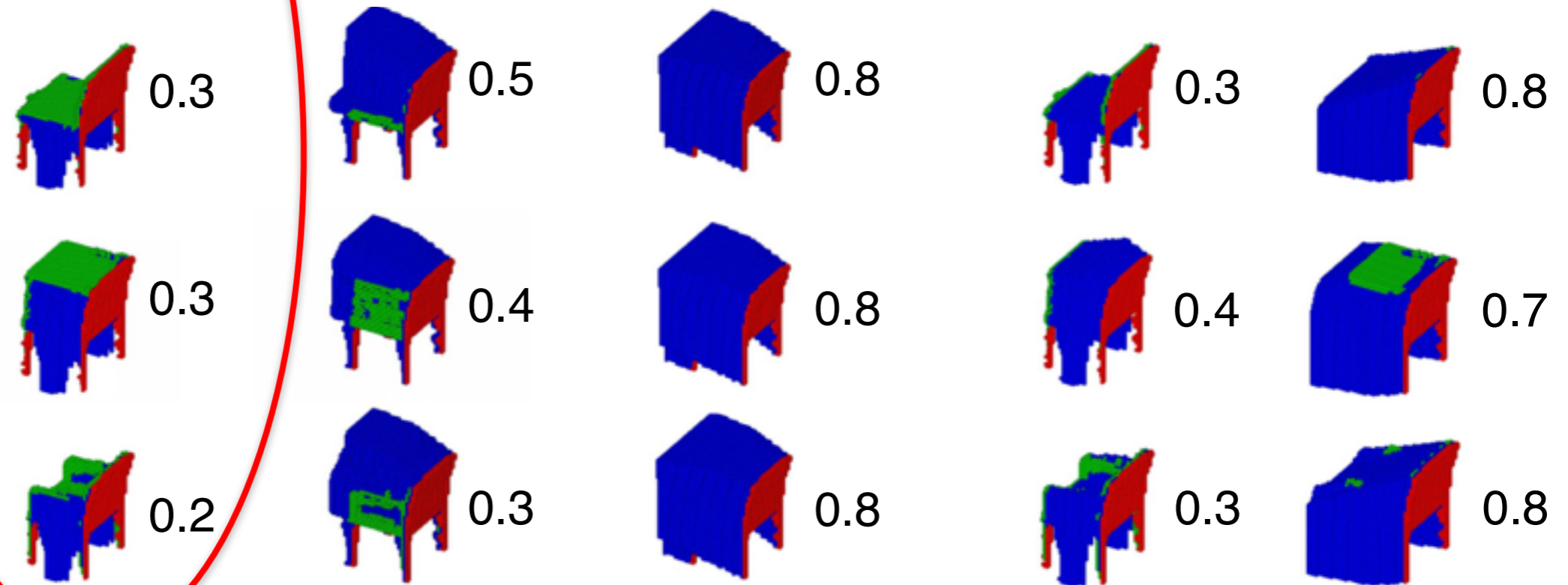
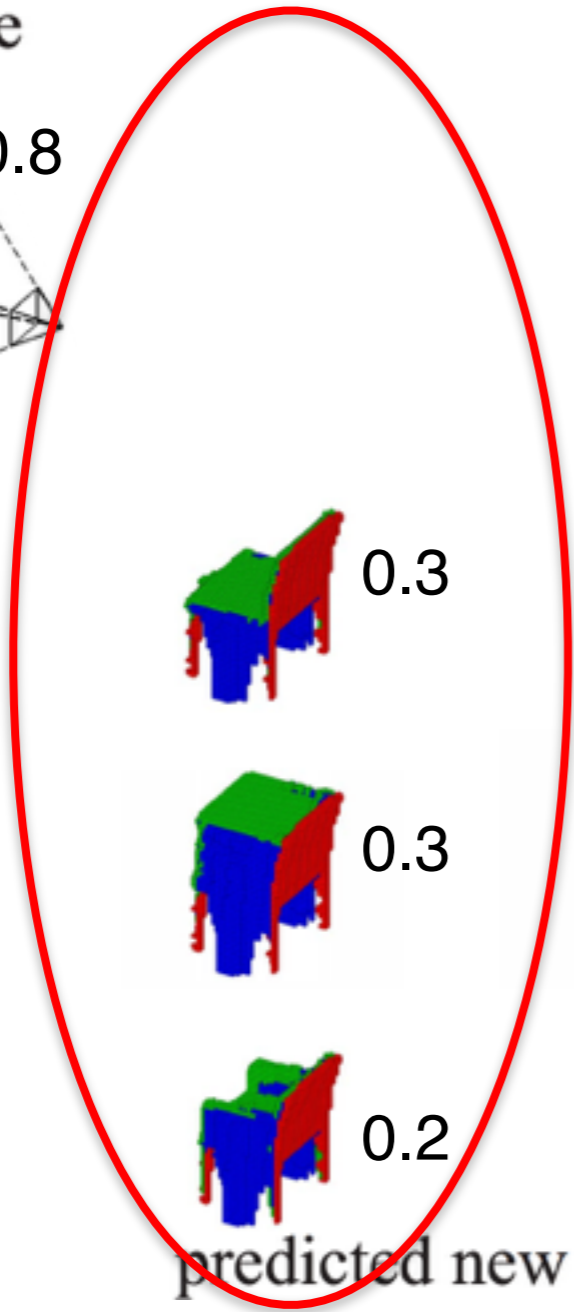
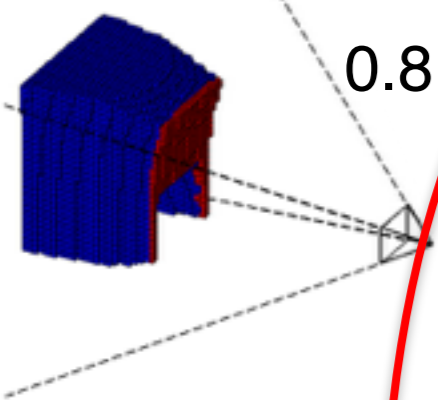


predicted new freespace & visible surface for each shape under each view

Deep View Planning

■ observed surface x_o
■ unknown
 ■ potentially visible voxels in next view
 ■ newly visible surface x_n^i
 free space

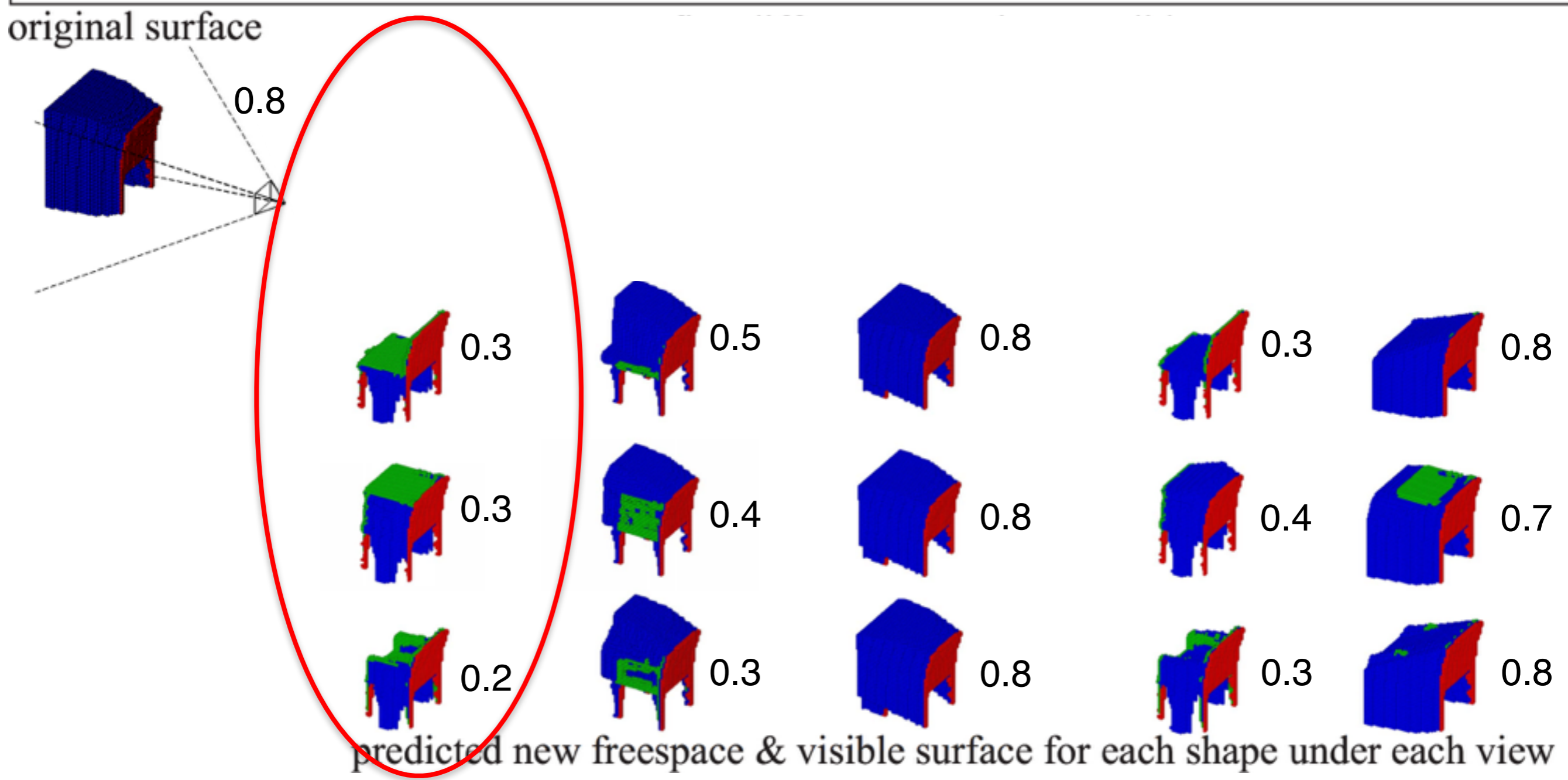
original surface



predicted new freespace & visible surface for each shape under each view

Deep View Planning

■ observed surface \mathbf{x}_o
■ unknown
 ■ potentially visible voxels in next view
 ■ newly visible surface \mathbf{x}_n^i
 free space



Mathematically, this is equivalent to evaluate the conditional mutual information:

$$I(y; \mathbf{x}_n^i | \mathbf{x}_o = x_o)$$

Deep View Planning

	all
Ours	0.80
Max Visibility	0.78
Furthest Away	0.69
Random Selection	0.72

Recognition Accuracy from Two Views.

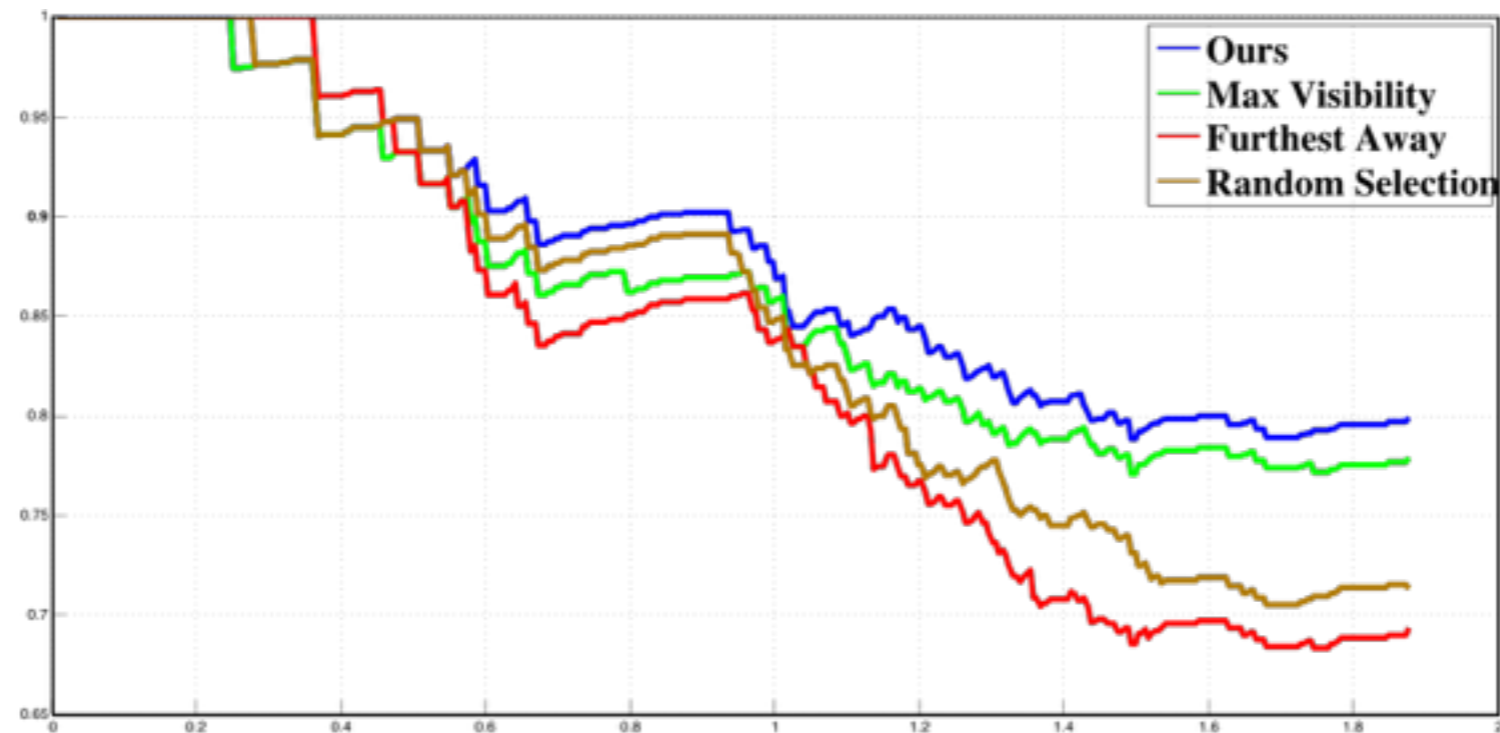
Based on the algorithms' choice, we obtain the actual depth map for the next view and recognize the objects using two views by our 3D ShapeNets.

Deep View Planning

	all
Ours	0.80
Max Visibility	0.78
Furthest Away	0.69
Random Selection	0.72

Recognition Accuracy from Two Views.

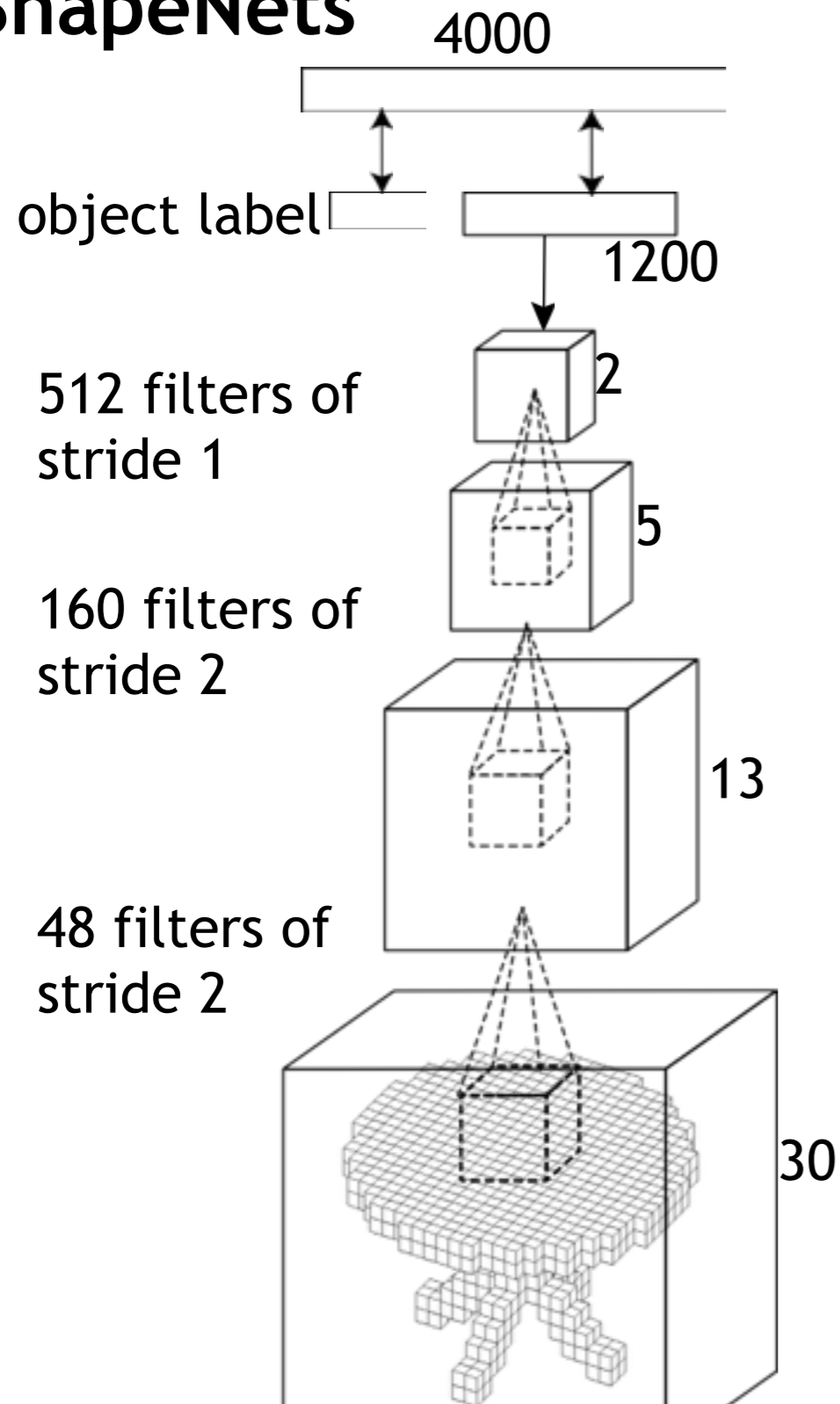
Based on the algorithms' choice, we obtain the actual depth map for the next view and recognize the objects using two views by our 3D ShapeNets.



Our algorithm stands out as the uncertainty of the first view increases 30

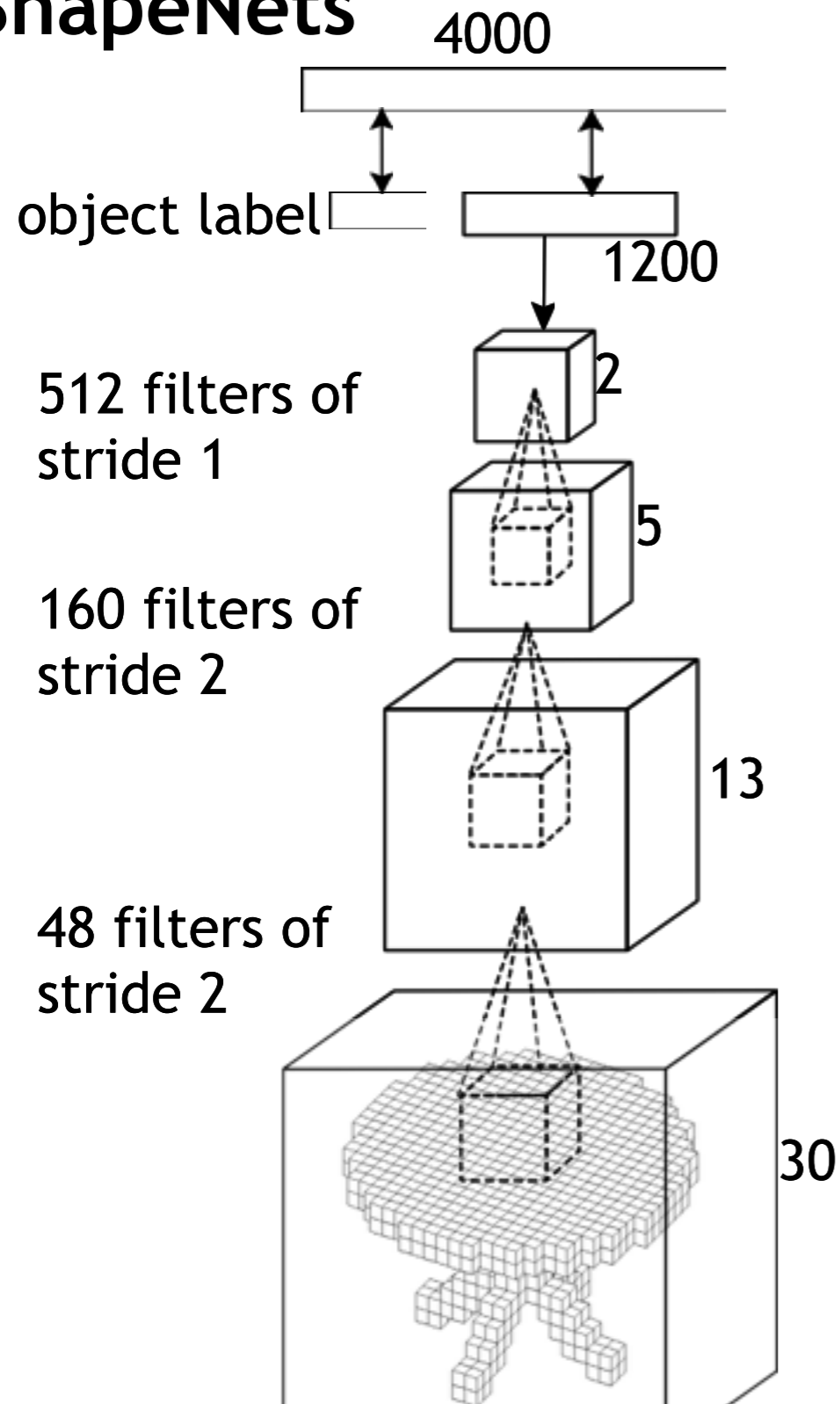
Back Propagation Fine-tuning

3D ShapeNets



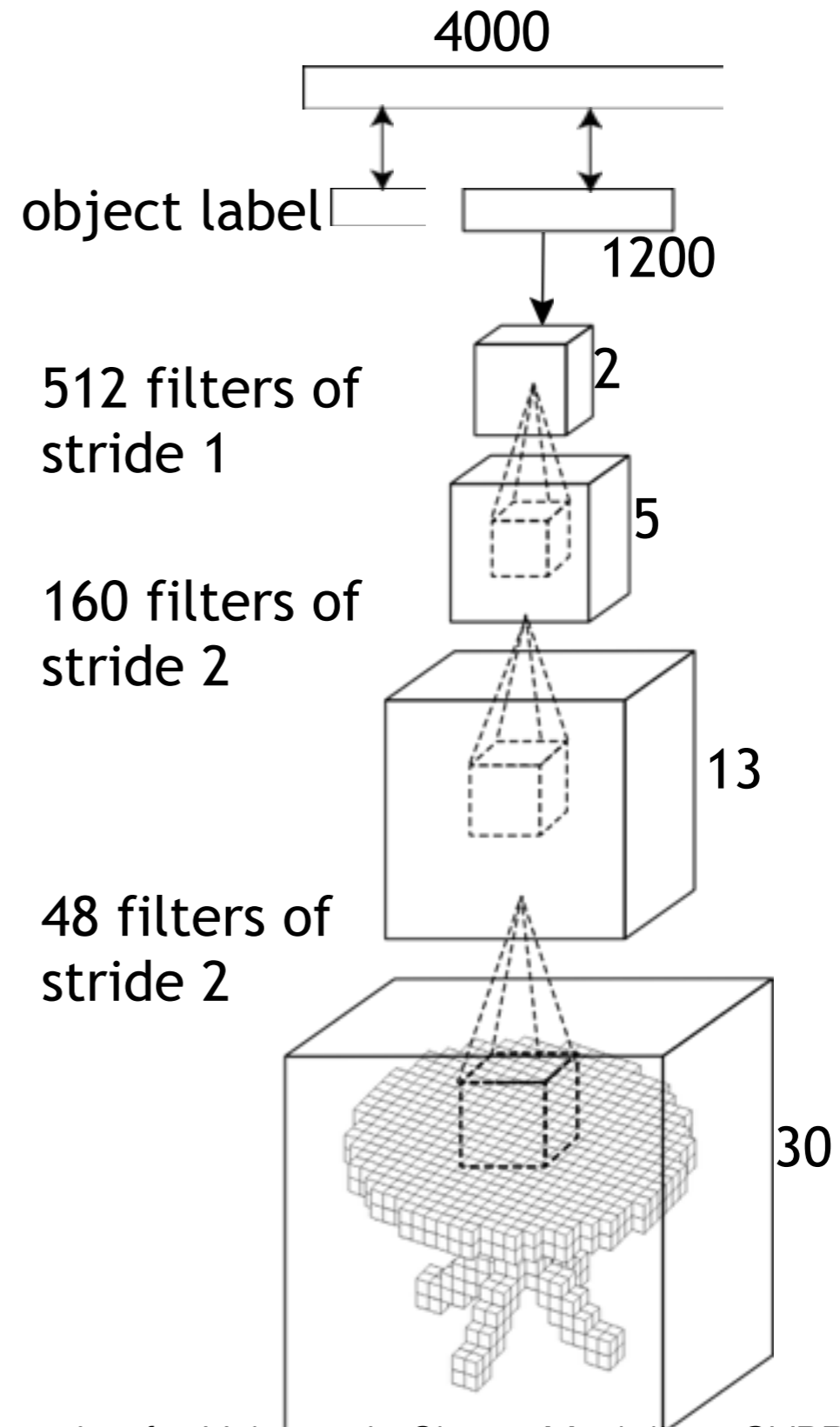
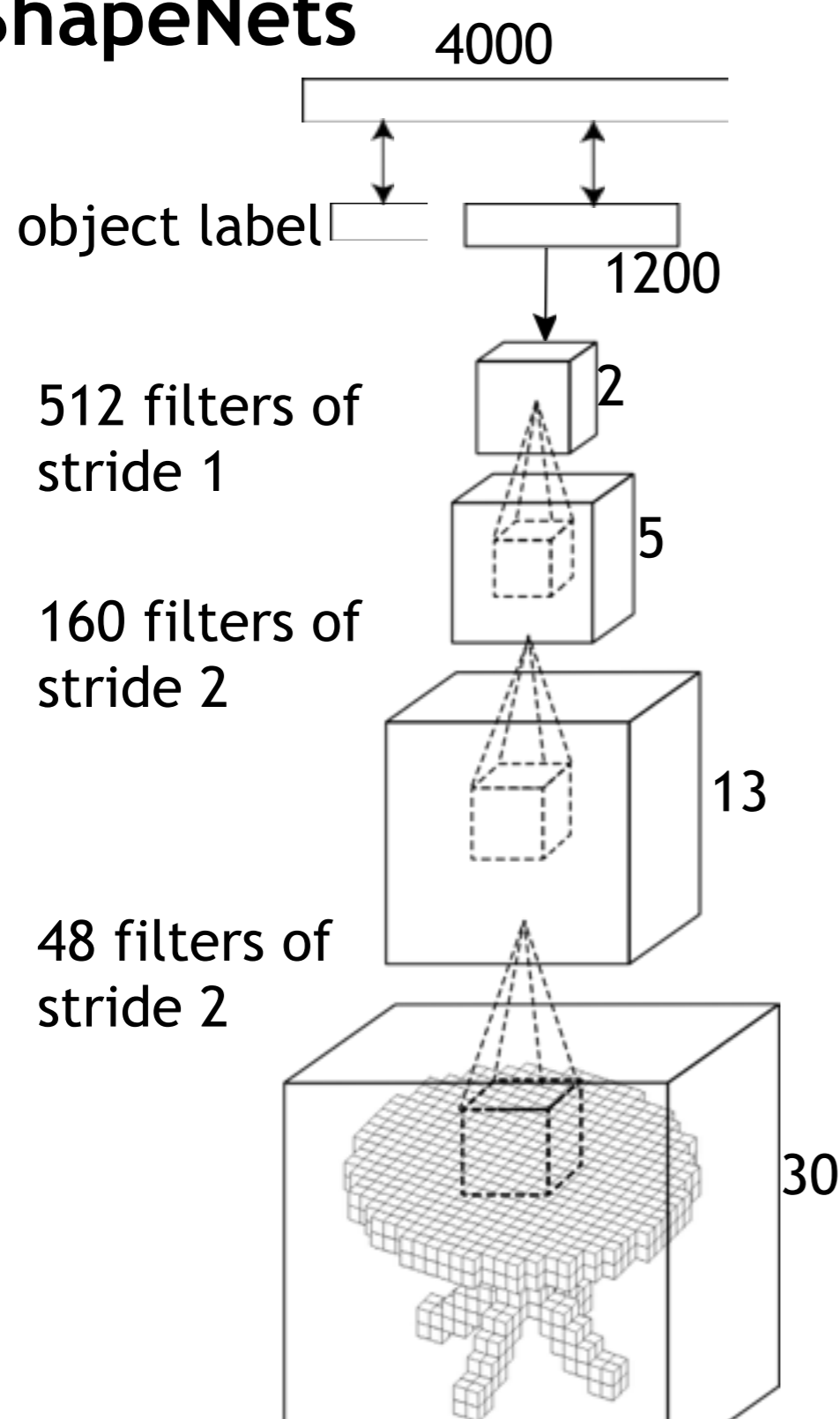
Back Propagation Fine-tuning

3D ShapeNets



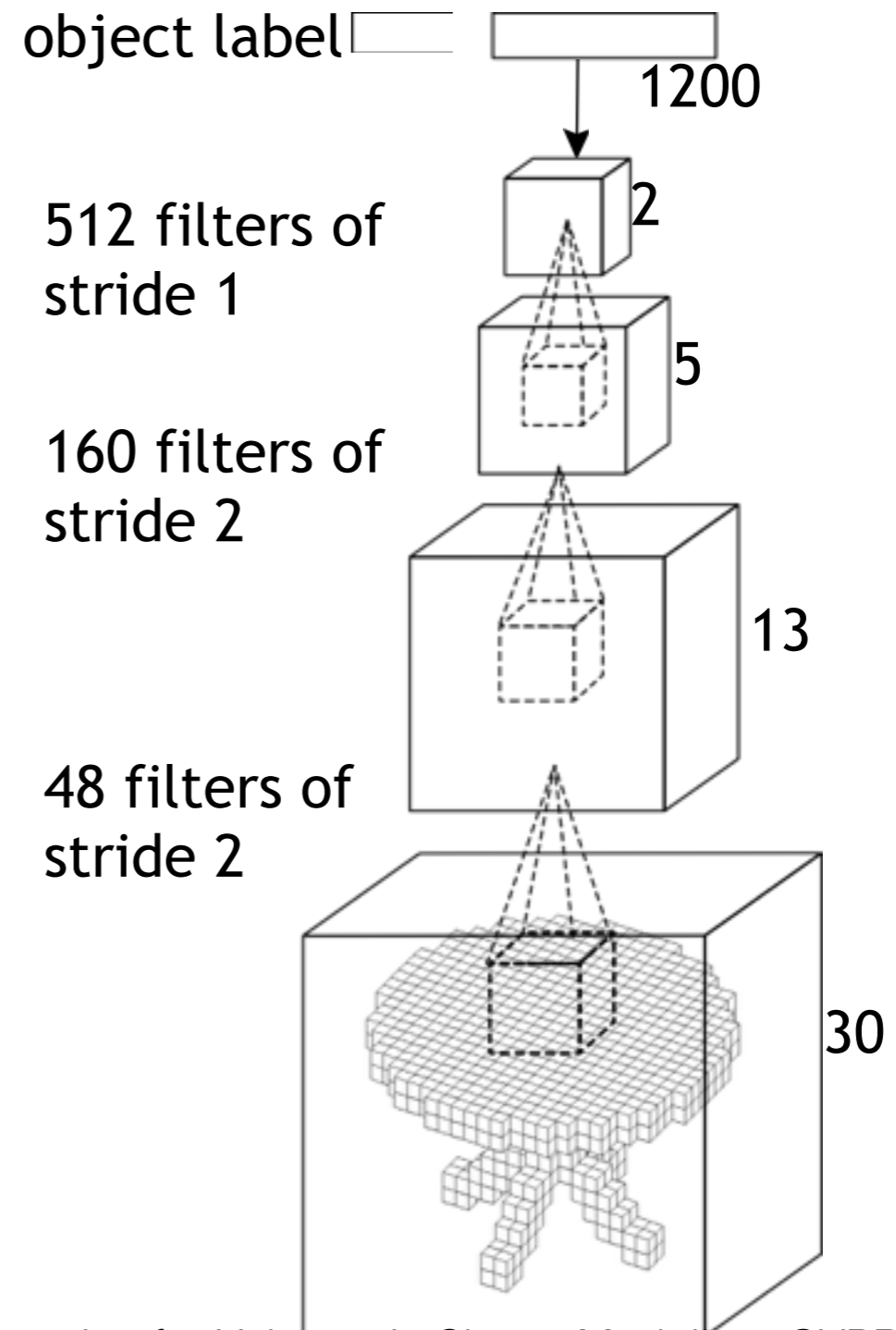
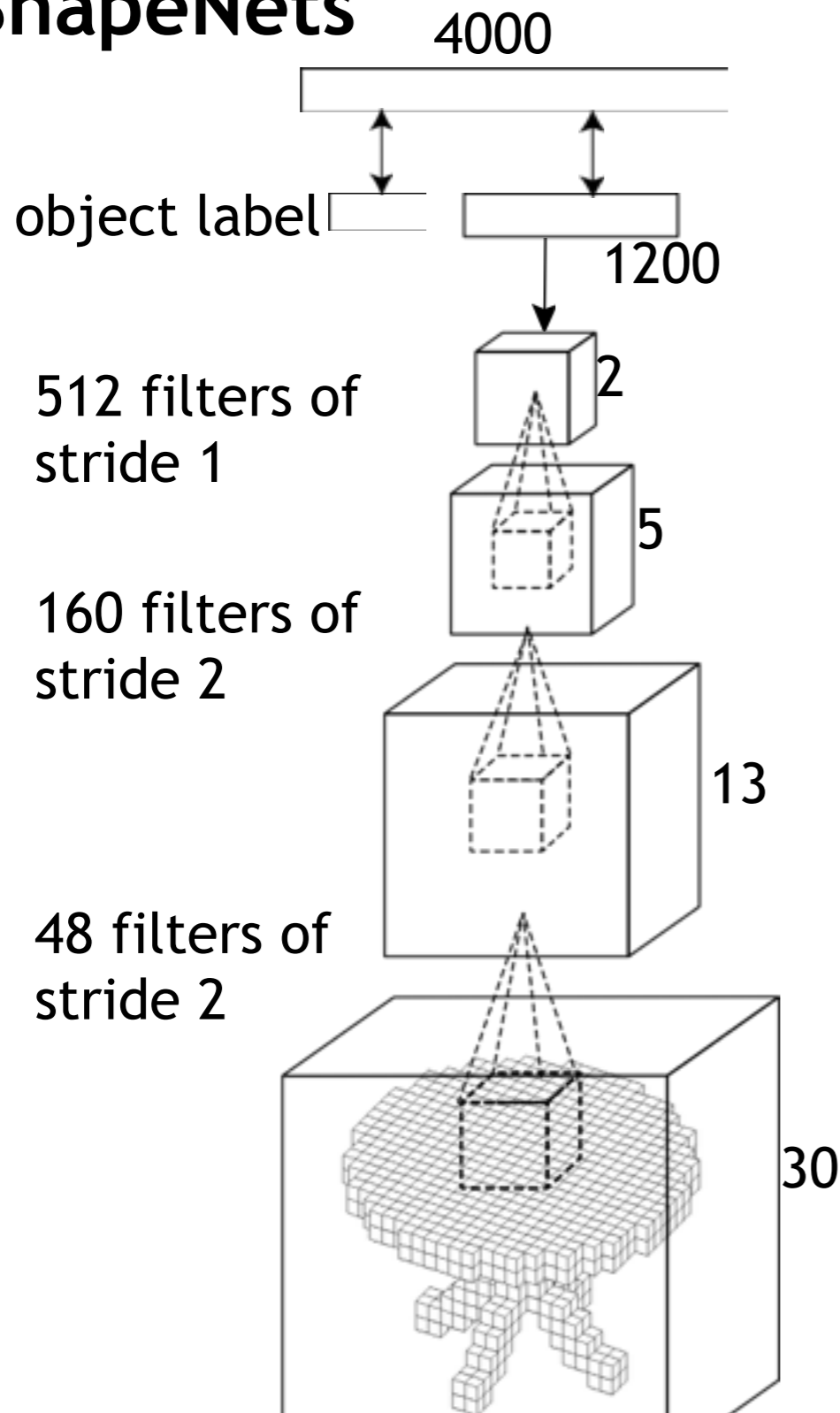
Back Propagation Fine-tuning

3D ShapeNets



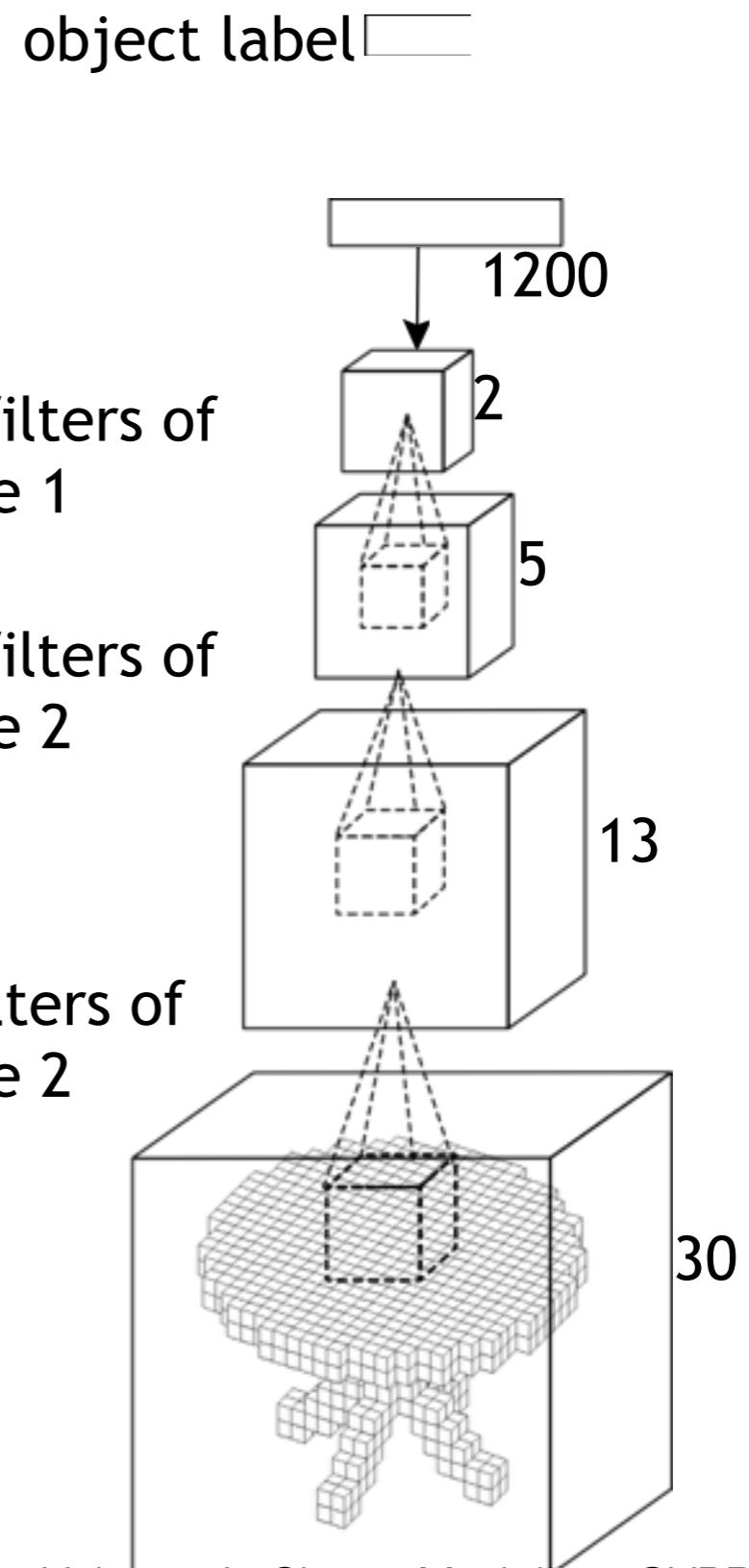
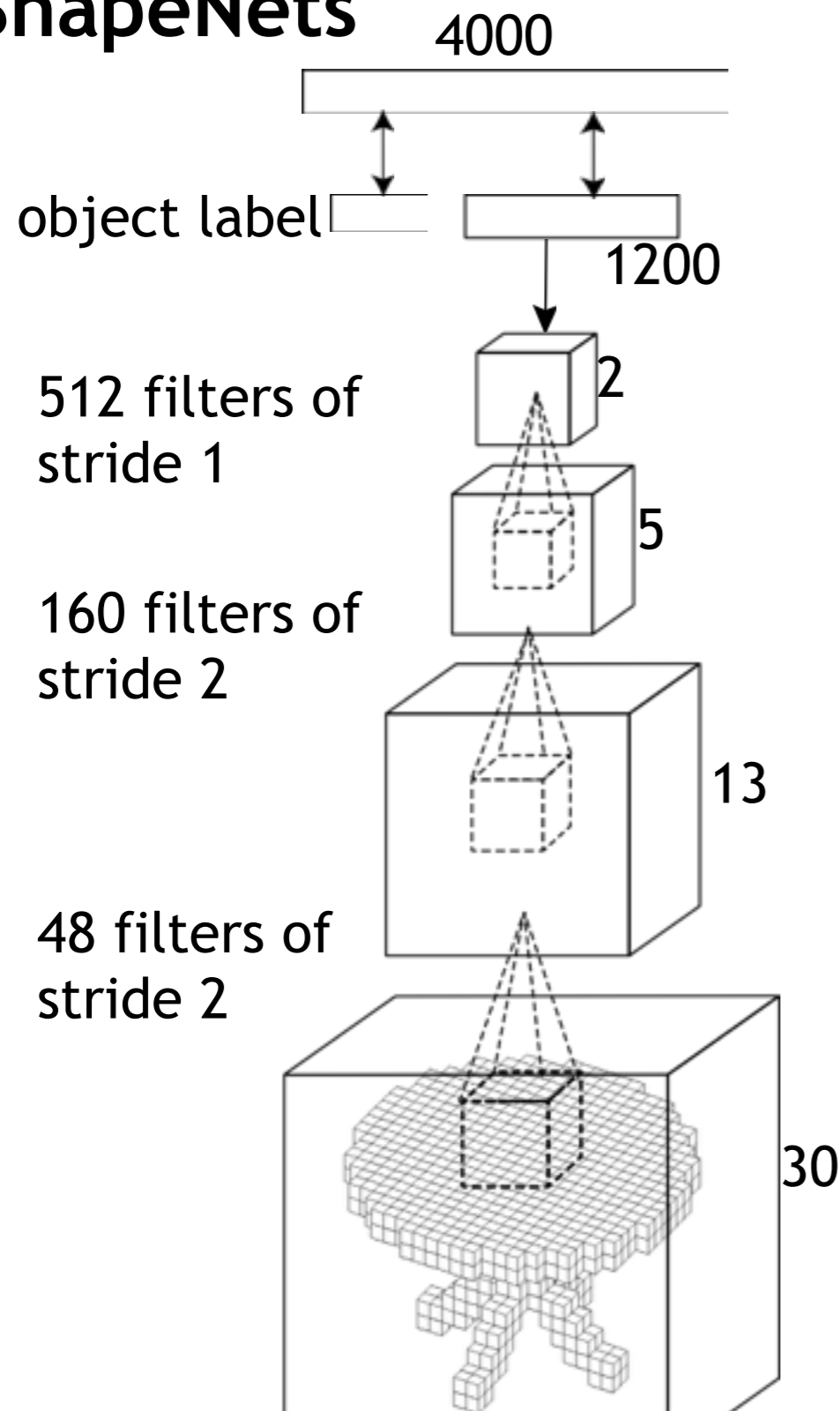
Back Propagation Fine-tuning

3D ShapeNets



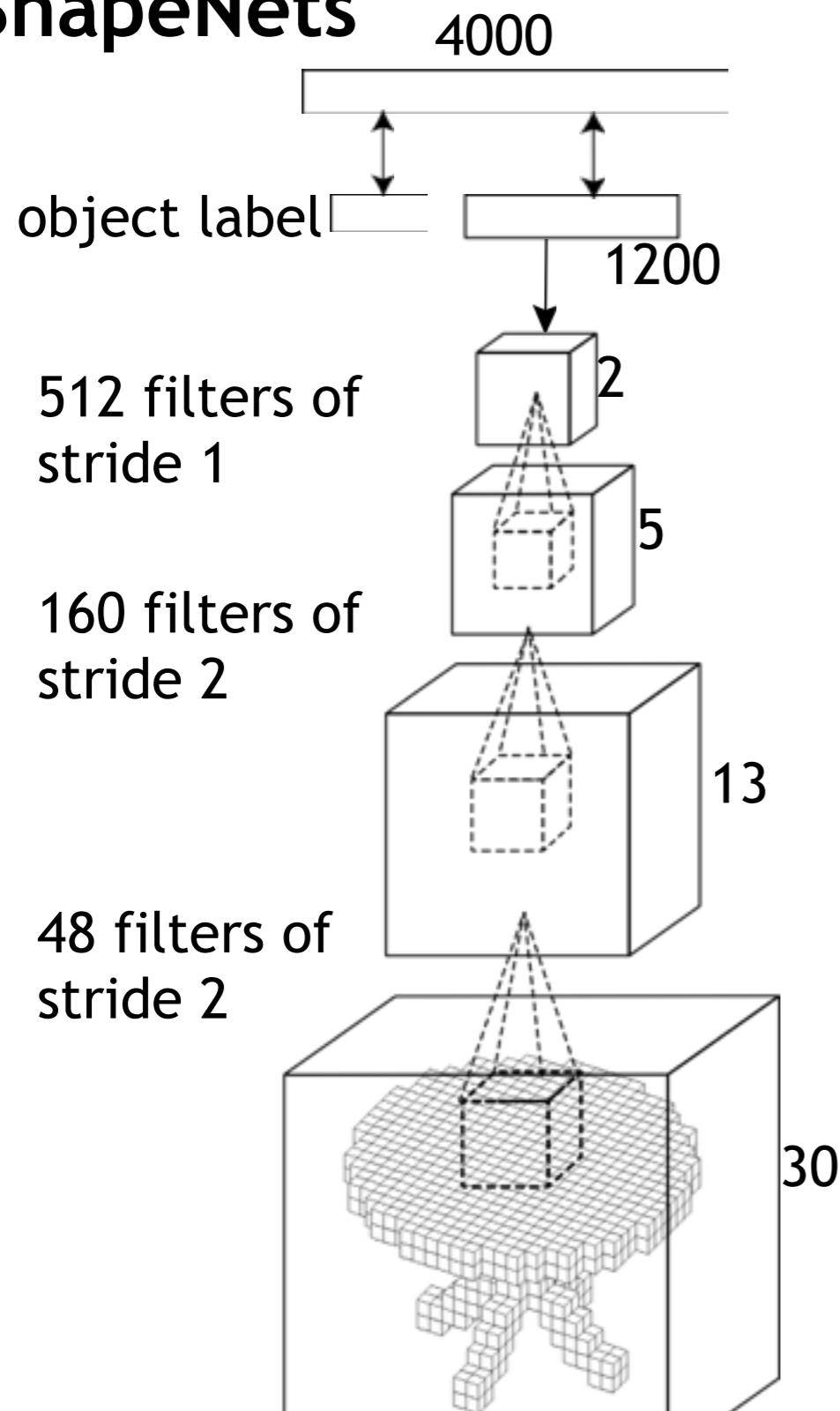
Back Propagation Fine-tuning

3D ShapeNets

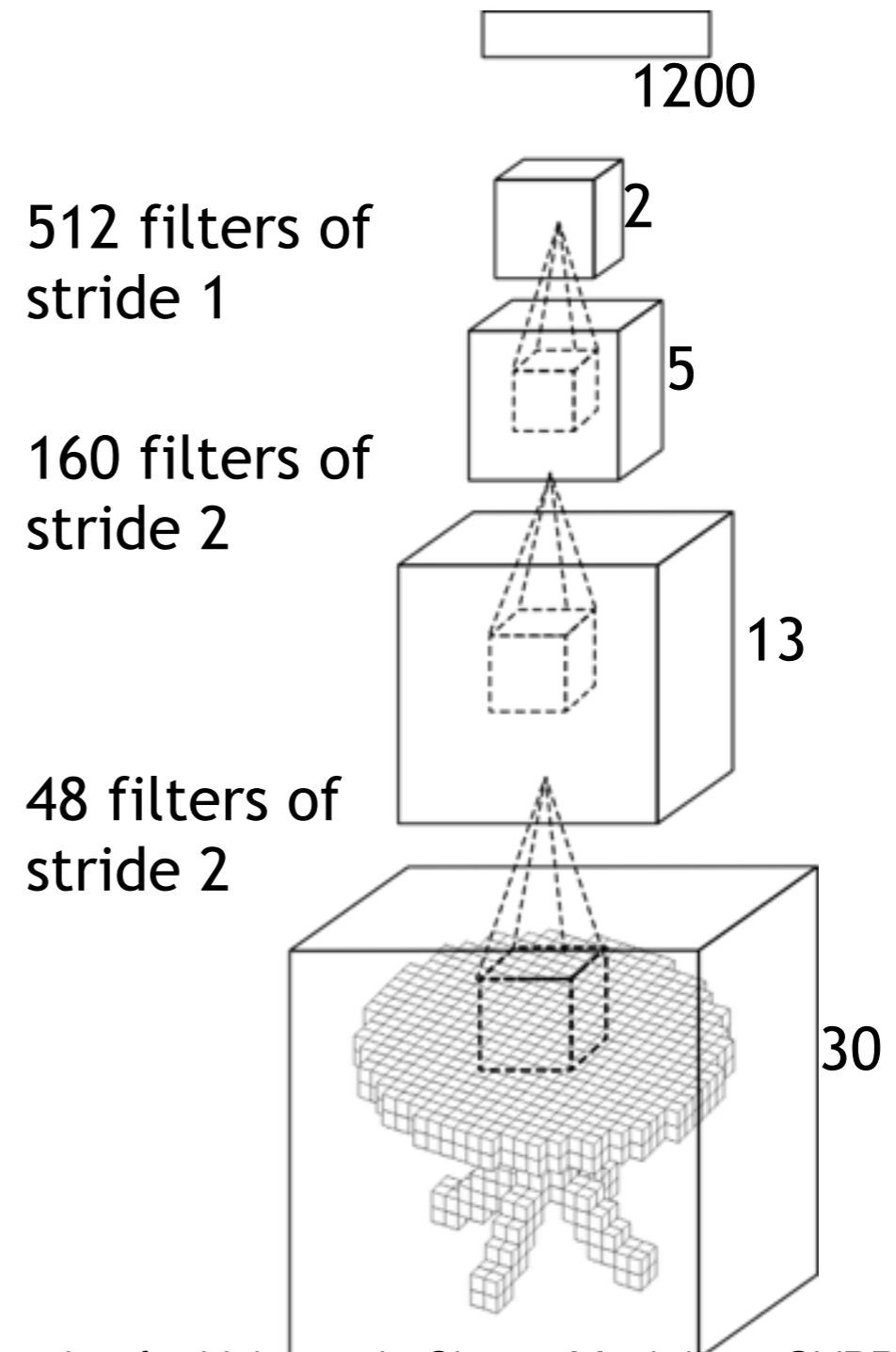


Back Propagation Fine-tuning

3D ShapeNets

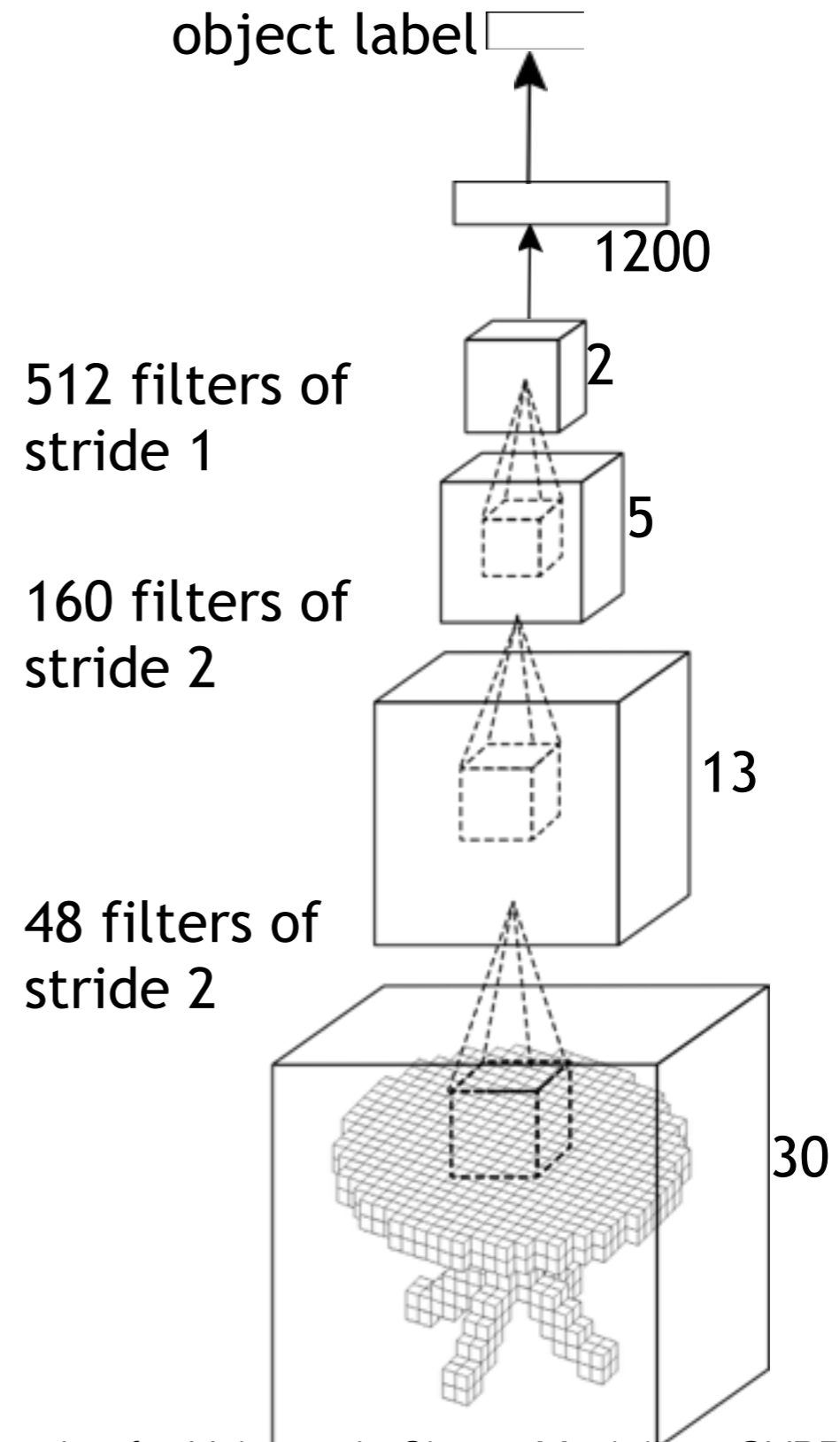
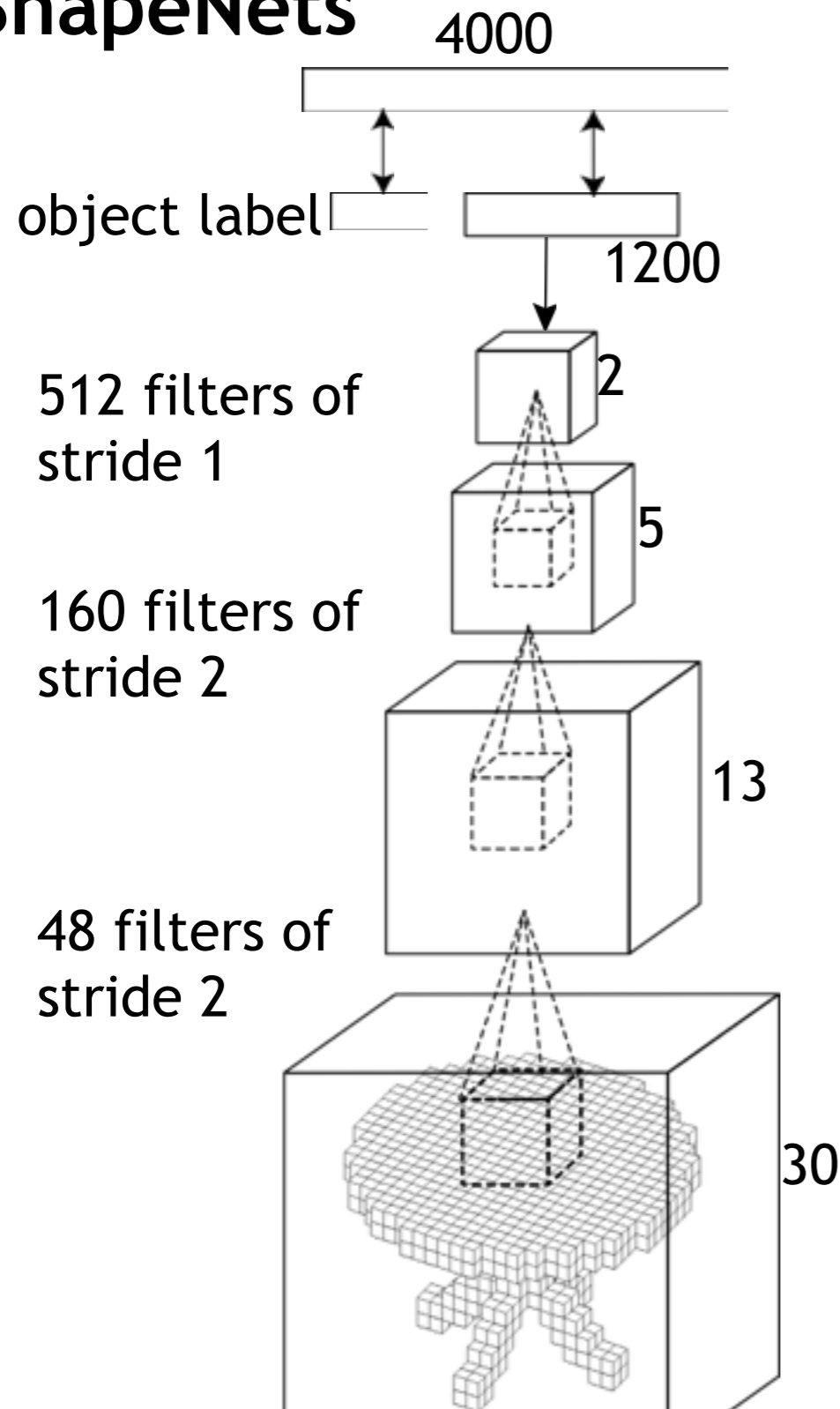


object label



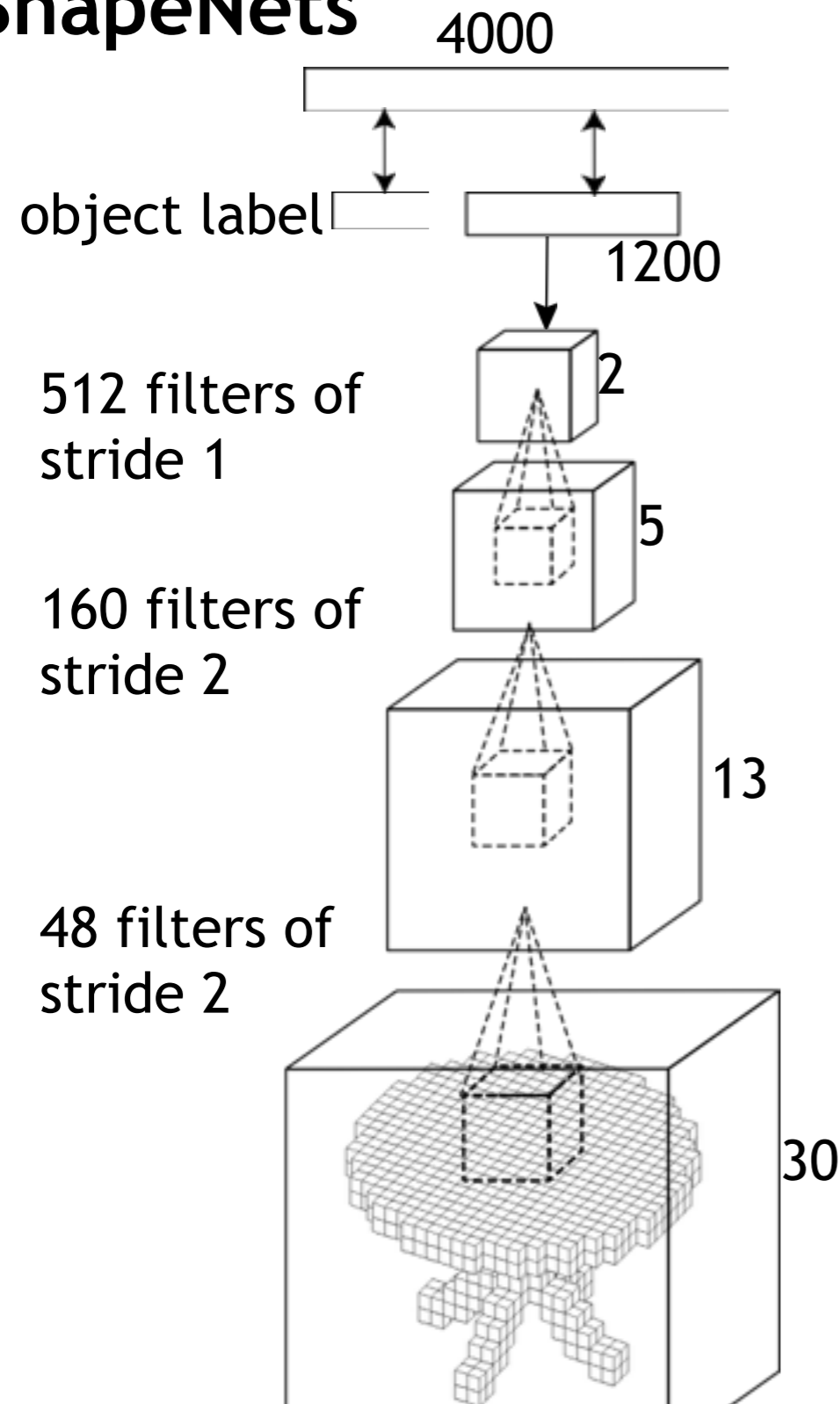
Back Propagation Fine-tuning

3D ShapeNets

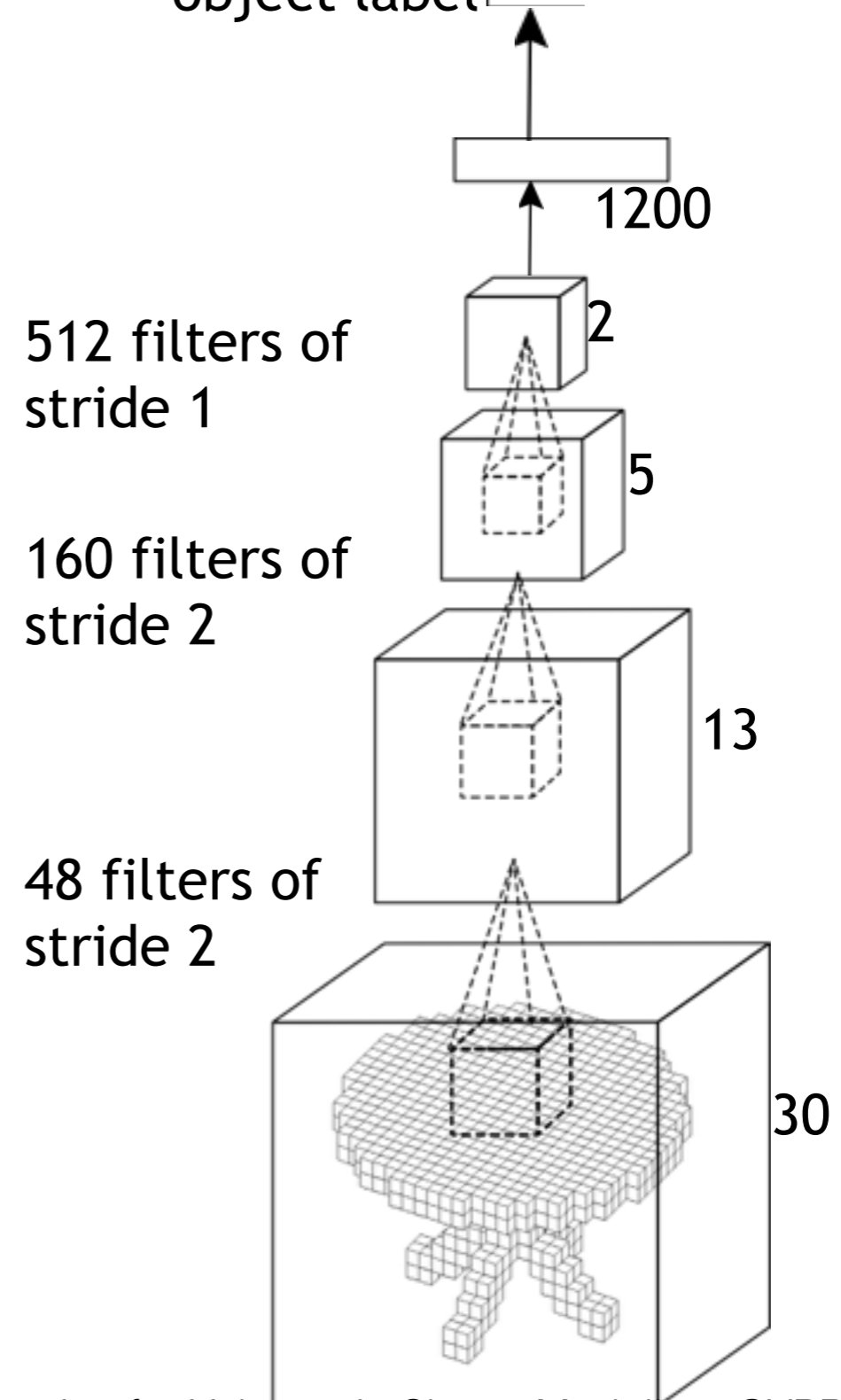


Back Propagation Fine-tuning

3D ShapeNets



3D CNN object label



As a 3D Feature Extractor

As a 3D Feature Extractor

Mesh Classification & Retrieval

10 classes	SPH [18]	LFD [8]	Ours
classification	79.79 %	79.87 %	83.54%
retrieval AUC	45.97%	51.70%	69.28%
retrieval MAP	44.05%	49.82%	68.26%
40 classes	SPH [18]	LFD [8]	Ours
classification	68.23%	75.47%	77.32%
retrieval AUC	34.47%	42.04%	49.94%
retrieval MAP	33.26%	40.91%	49.23%

As a 3D Feature Extractor

Mesh Classification & Retrieval

10 classes	SPH [18]	LFD [8]	Ours
classification	79.79 %	79.87 %	83.54%
retrieval AUC	45.97%	51.70%	69.28%
retrieval MAP	44.05%	49.82%	68.26%
40 classes	SPH [18]	LFD [8]	Ours
classification	68.23%	75.47%	77.32%
retrieval AUC	34.47%	42.04%	49.94%
retrieval MAP	33.26%	40.91%	49.23%

2.5D object recognition

	all
[29] Depth	0.376
NN	0.374
ICP	0.471
3D ShapeNets	0.437
3D ShapeNets fine-tuned	0.579
[29] RGB	0.334
[29] RGBD	0.448

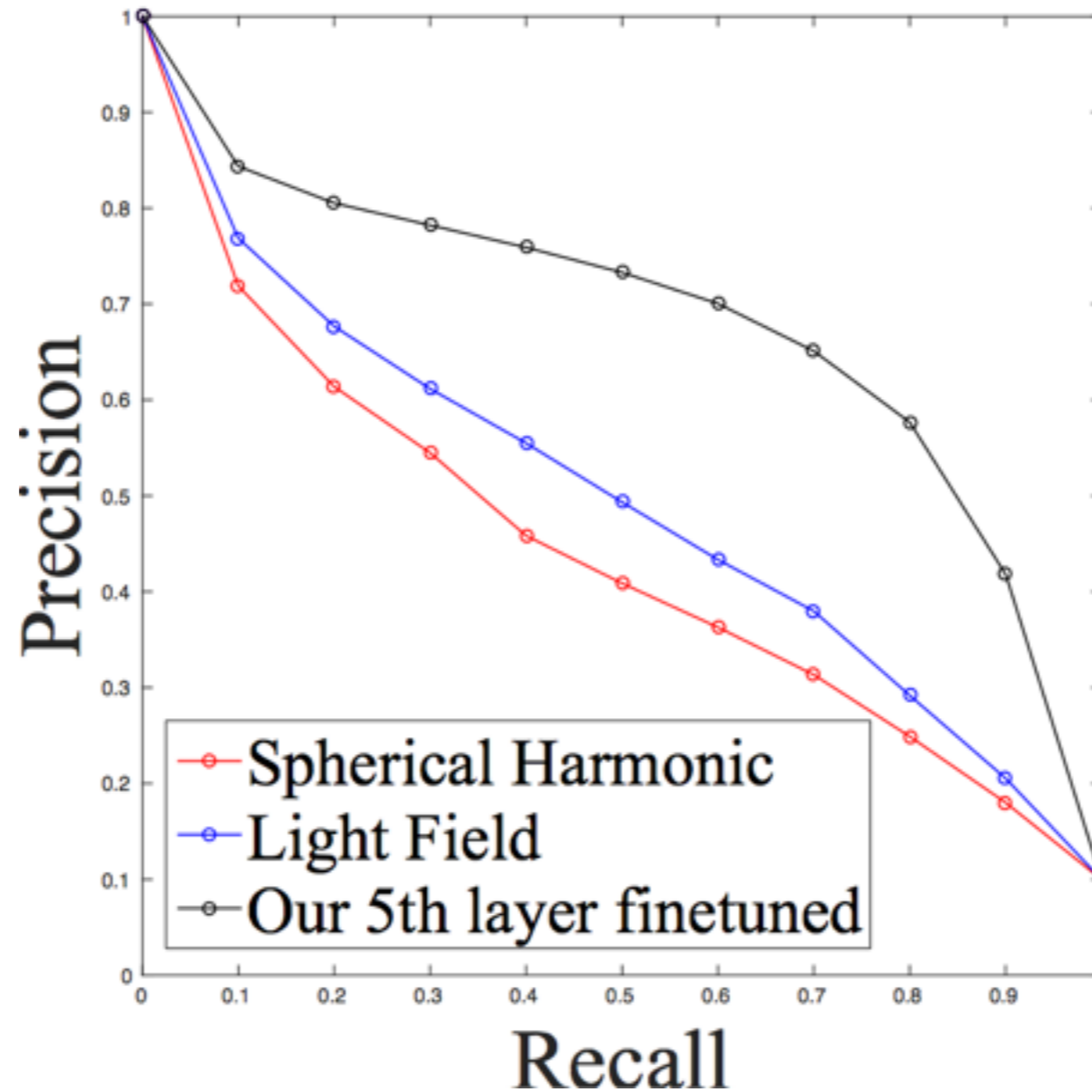
[29] R. Socher, B. Huval, B. Bhat, C. D. Manning, and A. Y. Ng.

Convolutional-recursive deep learning for 3d object classification. In NIPS 2012.

Slide Credit: Wu et al

As a 3D Feature Extractor

10 Classes Results



mesh retrieval

Extensions

- Include RGB information in representation
- 3D Segmentation
- Improve for non-rigid 3D objects

Discussion Points

- Is the network deep enough?
 - $30 \times 30 \times 30 = 27000$ vs $256 \times 256 = 65000$ for Image Net
 - 150K training examples vs millions for Image Net

- Won't removal of max-pooling layers hurt performance on classification tasks?

Algorithm	ModelNet40 Classification (Accuracy)	ModelNet40 Retrieval (mAP)	ModelNet10 Classification (Accuracy)	ModelNet10 Retrieval (mAP)
MVCNN [3]	90.1%	79.5%		
VoxNet [2]	83%		92%	
DeepPano [4]	77.63%	76.81%	85.45%	84.18%
3DShapeNets [1]	77%	49.2%	83.5%	68.3%

[1] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang and J. Xiao. 3D ShapeNets: A Deep Representation for Volumetric Shapes. CVPR2015.

[2] D. Maturana and S. Scherer. VoxNet: A 3D Convolutional Neural Network for Real-Time Object Recognition. IROS2015.

[3] H. Su, S. Maji, E. Kalogerakis, E. Learned-Miller. Multi-view Convolutional Neural Networks for 3D Shape Recognition. ICCV2015.

[4] B Shi, S Bai, Z Zhou, X Bai. DeepPano: Deep Panoramic Representation for 3-D Shape Recognition. Signal Processing Letters 2015.

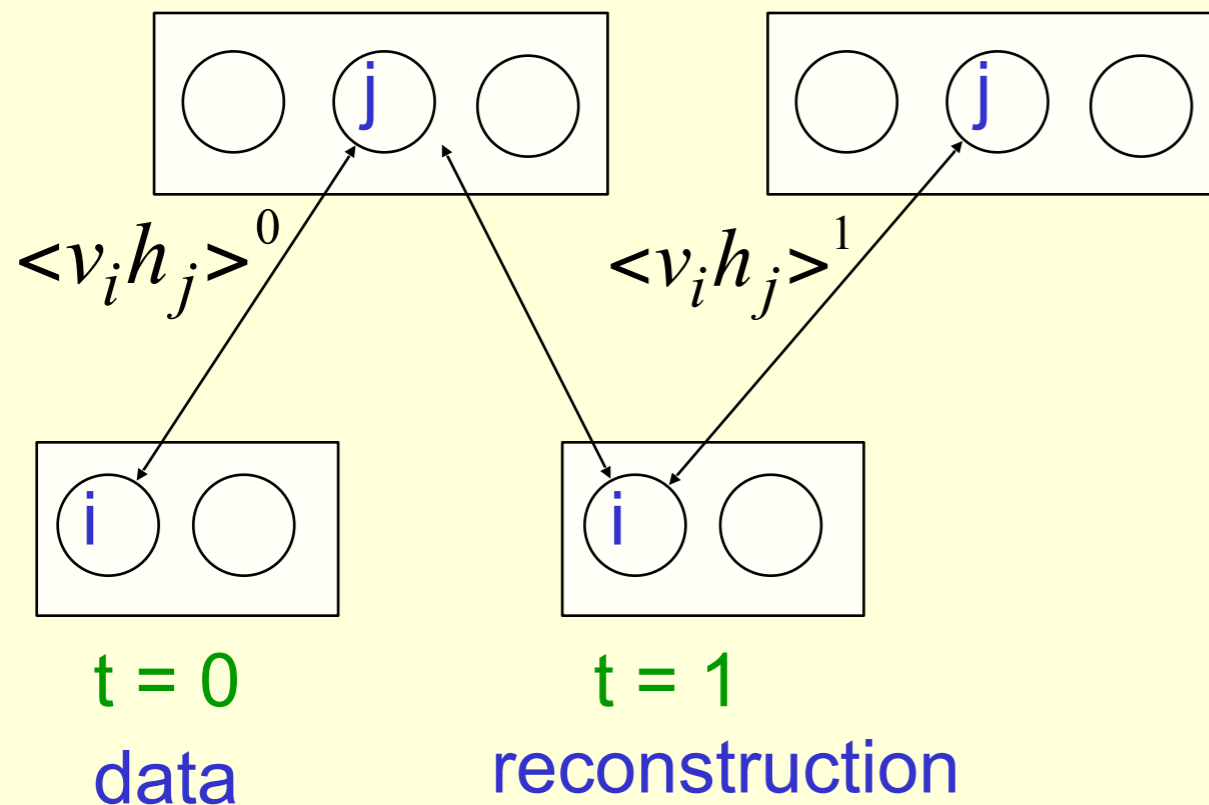
- Any other systems that use binary units with approximate training and inference techniques rather than standard back-prop?
- Hinton, Geoffrey E., Simon Osindero, and Yee-Whye Teh. "A fast learning algorithm for deep belief nets." *Neural computation* 18.7 (2006): 1527-1554
- Salakhutdinov, Ruslan, Andriy Mnih, and Geoffrey Hinton. "Restricted Boltzmann machines for collaborative filtering." *Proceedings of the 24th international conference on Machine learning*. ACM, 2007.

- Are there better ways for representing 3D Shapes. In particular, doesn't the voxel representation have the bottleneck of cubic dependency on grid size?
- Yes. Su, Majhi et al that tries to recognize 3D shapes from multiple 2D views instead of voxel representation and get better results for classification .

- Are there other 3D CAD model datasets
 - 3D Warehouse. <https://3dwarehouse.sketchup.com/>
- Manually removing clutter from 3D CAD models a problem
- Did not address non-rigid objects sufficiently.
 - Even the 40 model classification dataset seemed to contain only 4 non-rigid categories — persons, plant, sofas, curtains.

Appendix

Contrastive divergence learning: A quick way to learn an RBM



Start with a training vector on the visible units.

Update all the hidden units in parallel

Update all the visible units in parallel to get a “reconstruction”.

Update all the hidden units again.

$$\Delta w_{ij} = \varepsilon \left(\langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1 \right)$$

This is not following the gradient of the log likelihood. But it works well.

It is approximately following the gradient of another objective function.

The wake-sleep algorithm for an SBN

- **Wake phase:** Use the recognition weights to perform a bottom-up pass.
 - Train the generative weights to reconstruct activities in each layer from the layer above.
- **Sleep phase:** Use the generative weights to generate samples from the model.
 - Train the recognition weights to reconstruct activities in each layer from the layer below.

