

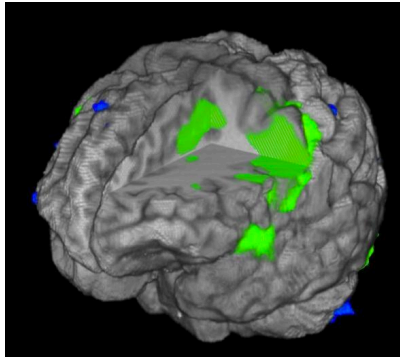
UC San Diego

**Lecture 12:**  
**Deep Learning on  
Volumetric Representation**

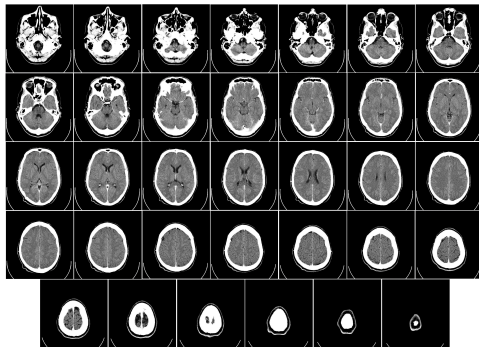
Instructor: Hao Su

Feb 19, 2018

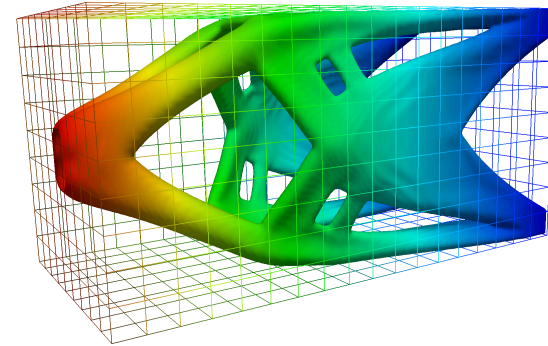
# Popular 3D volumetric data



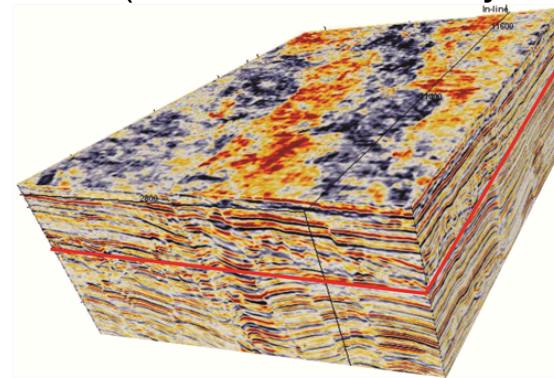
fMRI



CT



Manufacturing  
(finite-element analysis)



Geology

# 3D volumetric representations

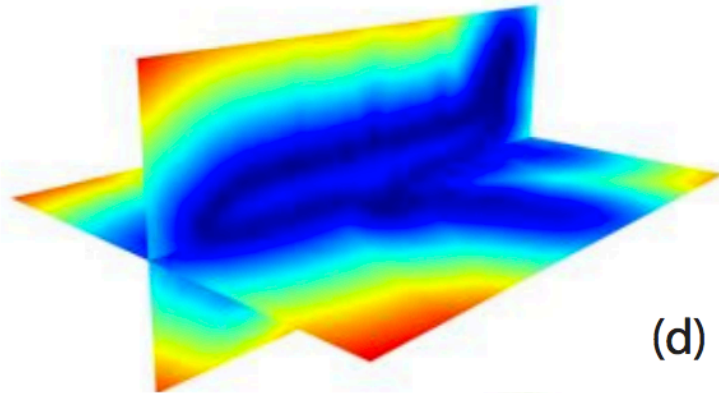
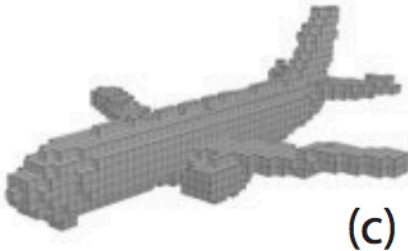
Mesh



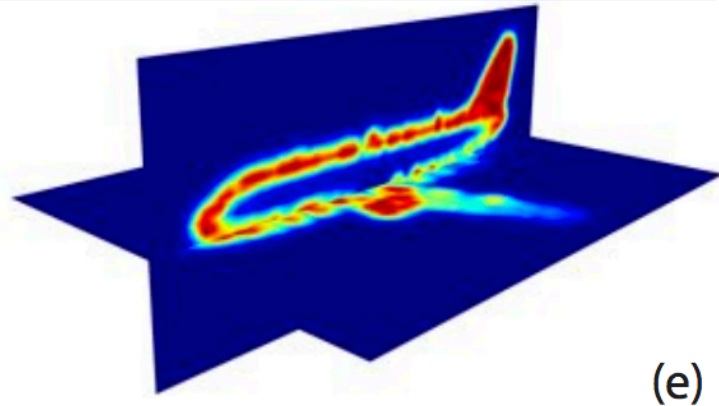
Point cloud



Volumetric occupancy



Volumetric distance field



Volumetric Gaussian distance field

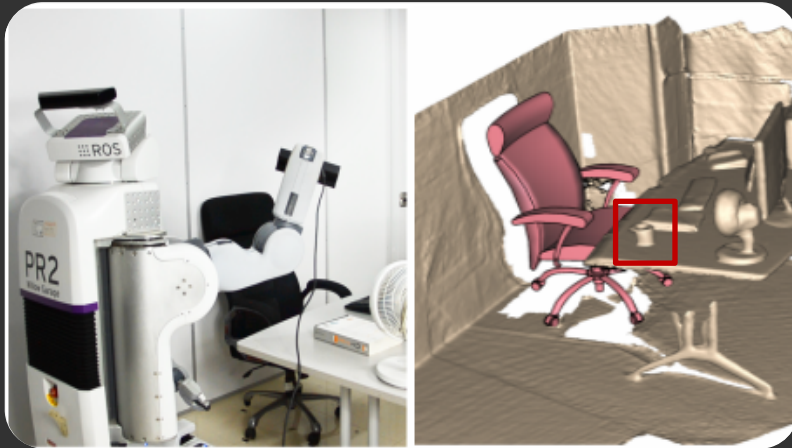
UC San Diego

# Shape Analysis

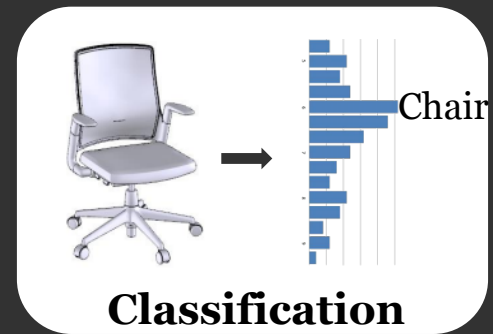


# 3D Shape Analysis

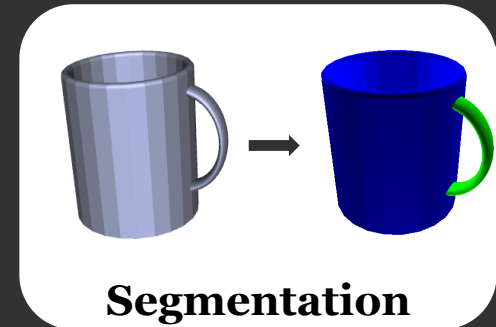
---



Robotics



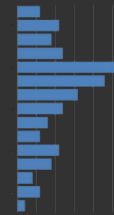
**Classification**



**Segmentation**

# CNN for 3D Shape Analysis

---



Chair

# Goal

---

- General
- Efficient
- Effective



3D data

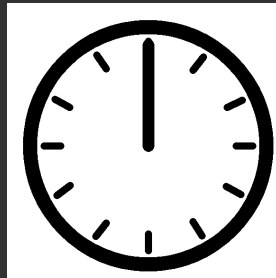


CNN  
structure

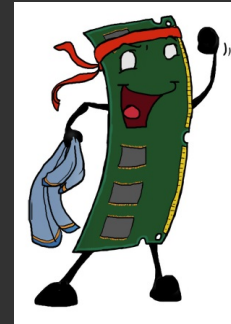
# Goal

---

- General
- Efficient
- Effective



Time cost



Memory cost

# Goal

---

- General
- Efficient
- Effective

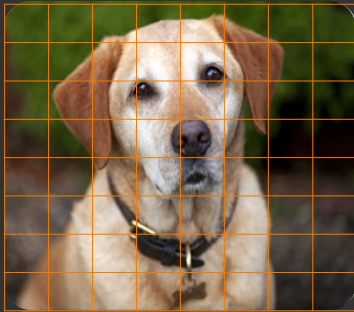


Performanc  
e

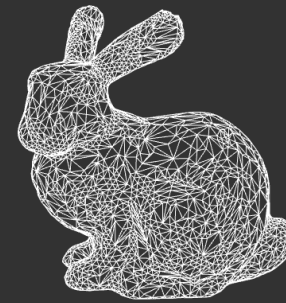
# Key Challenge

---

- A 3D shape representation for efficient CNN on GPU
  - 2D Regular grid
  - Irregular 3D shape

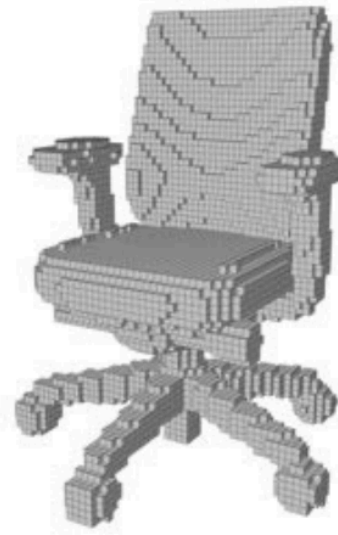
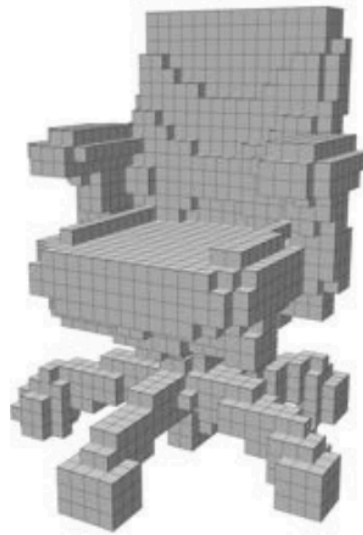


Image



Mesh

# The sparsity characteristic of 3D data



$$\frac{\#occupied\ grid}{\#total\ grid}$$

Occupancy:

10.41%

5.09%

2.41%

Resolution:

32

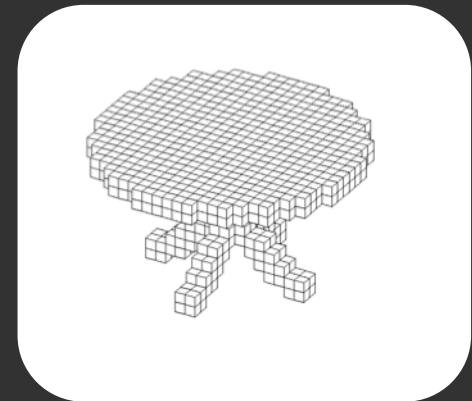
64

128

# Full Voxel based Solutions

---

- Related work: [Wu et al. 2015], [Maturana and Scherer 2015], ...
- **General**: intuitive extension of images ✓
- **Efficient**:  $O(N^3)$  ✗

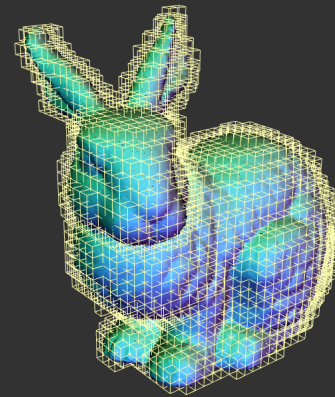
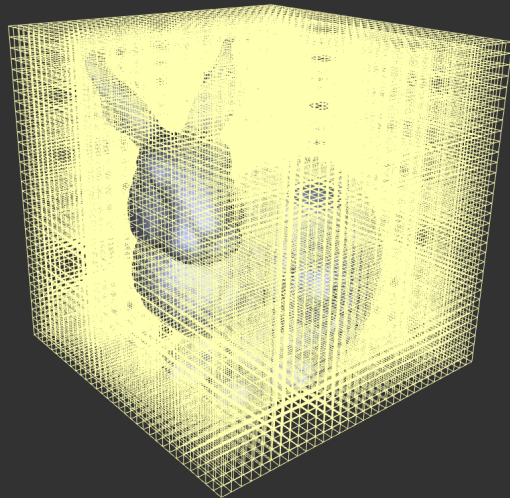




# Key Idea

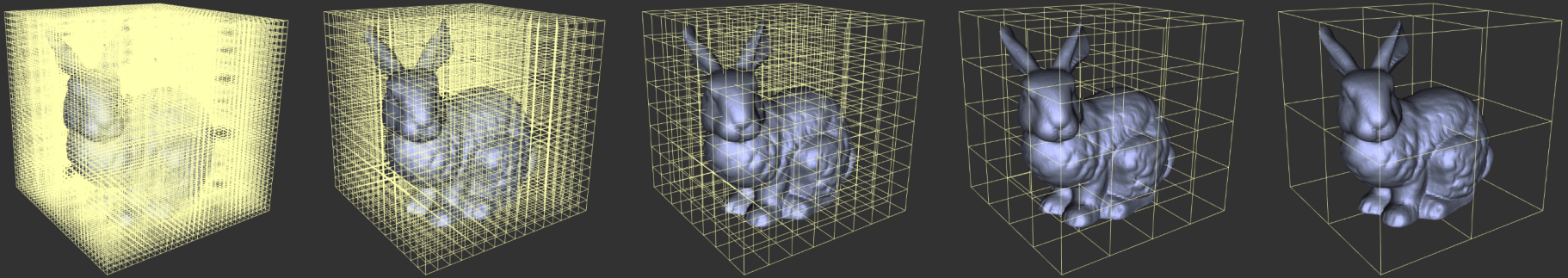
---

- Store the sparse surface signals
- Constrain the computation near the surface

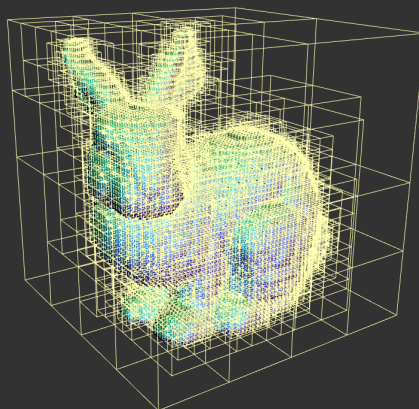
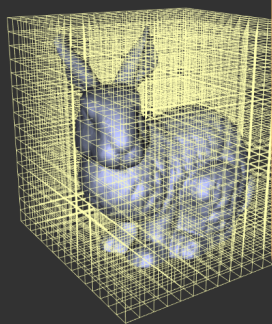
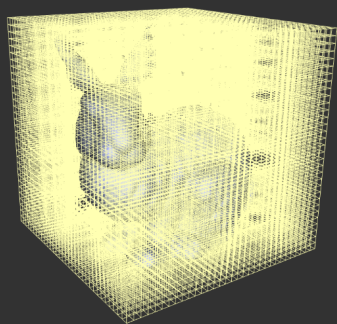


# Solution: Octree based CNN (O-CNN)

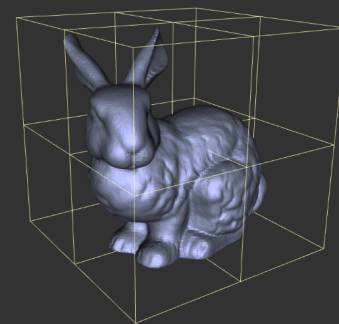
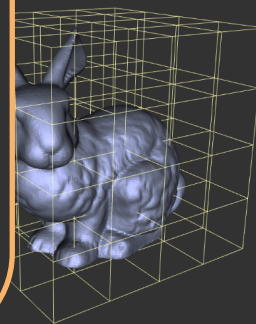
---



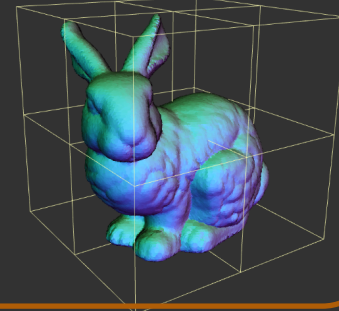
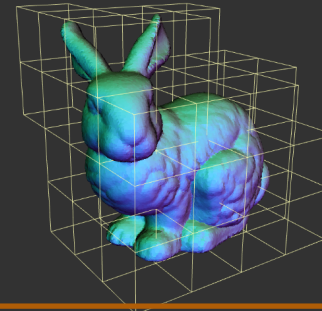
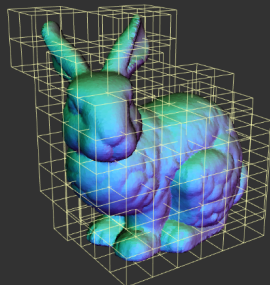
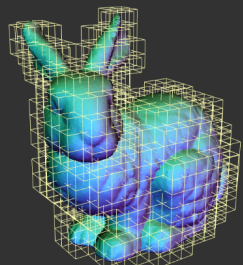
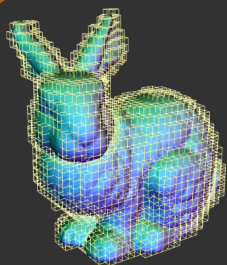
# Solution: Octree CNN (O-CNN)



**Octree**

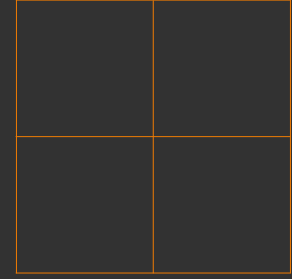
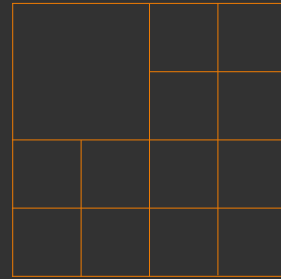
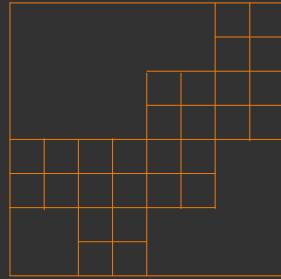
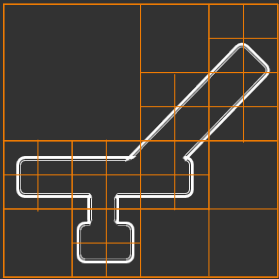


*CONV → POOL → CONV → POOL → CONV → POOL → CONV → POOL → CONV*



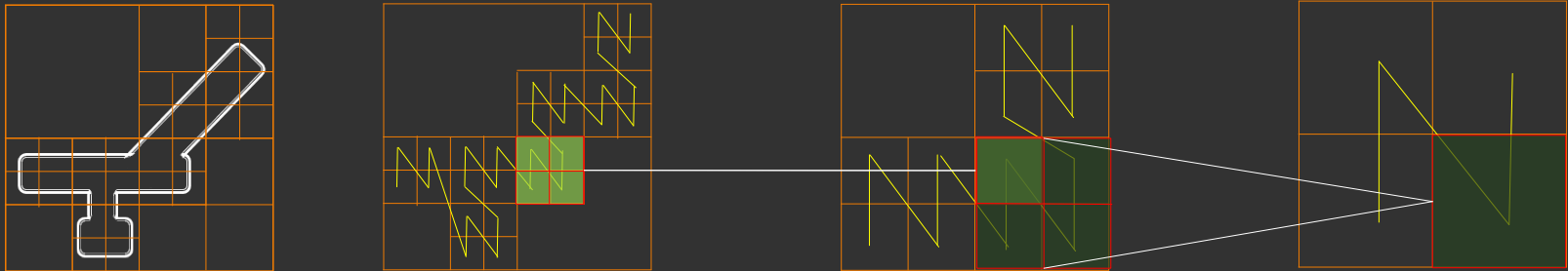
# Octree Data Structure

---



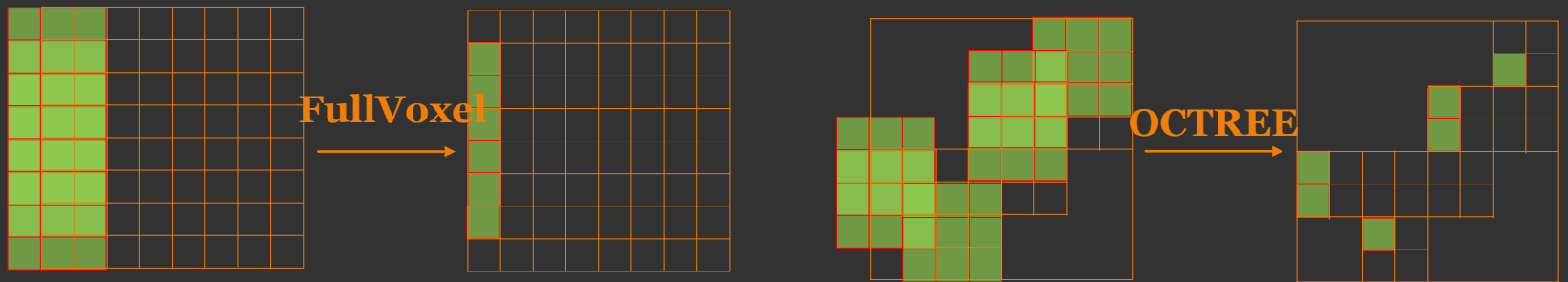
# Octree Data Structure

---



# Convolution on Octree

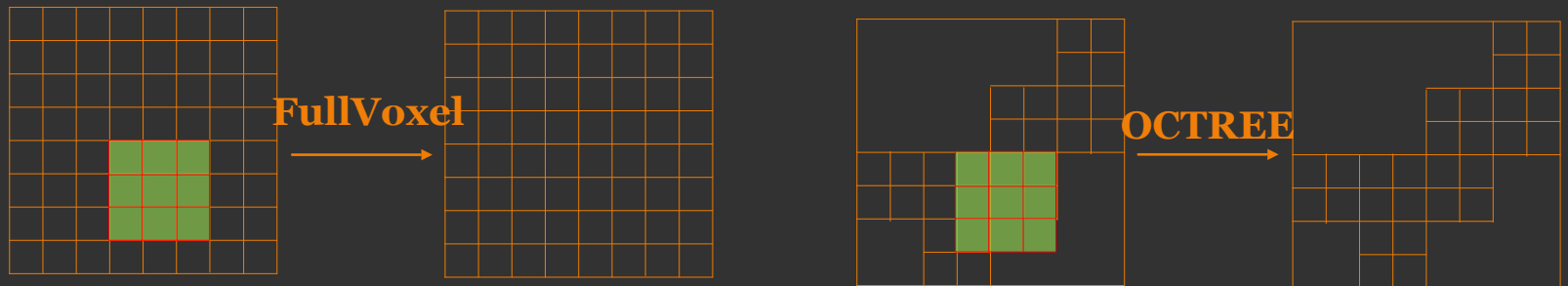
---



# Convolution on Octree

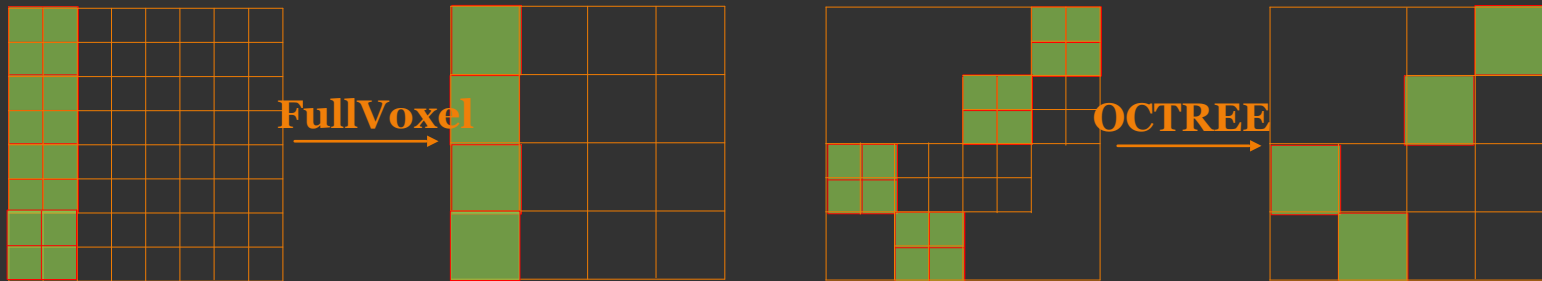
---

- Neighborhood searching: Hash table



# Pooling on Octree

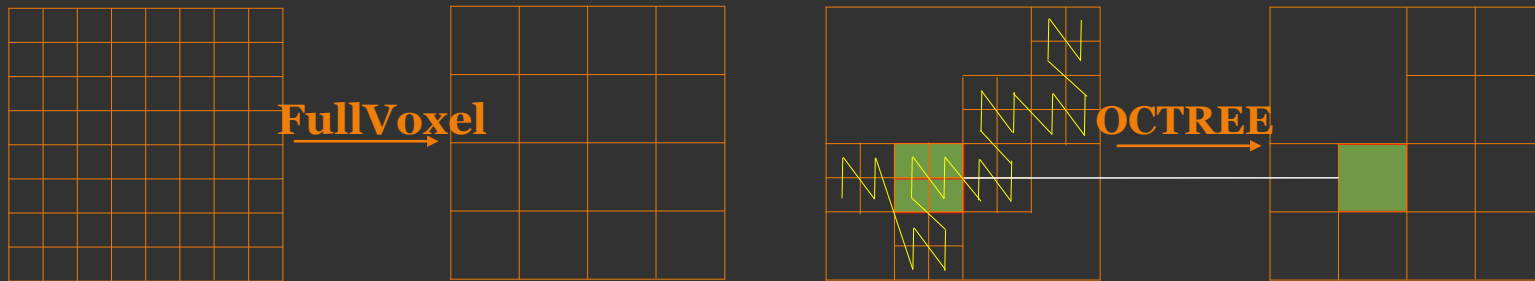
---





# Pooling on Octree

---



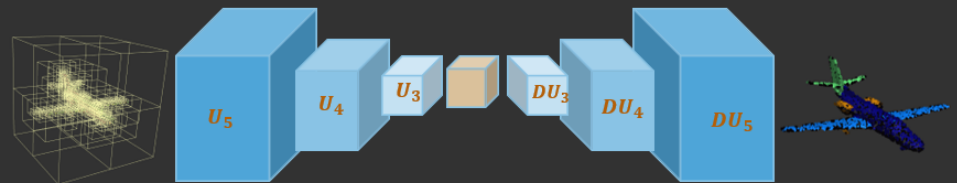
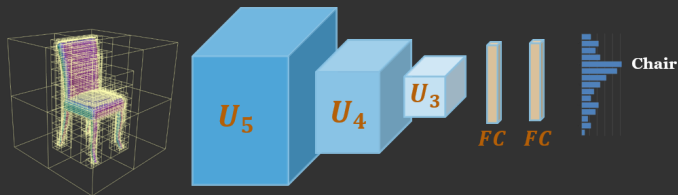
# Other CNN Operations on Octree

---

- Convolution with stride  $> 1$
- Deconvolution and un-pooling
  - Inverse operations of convolution and pooling
- Support most CNN architectures for images
  - LeNet [Lecun et al. 1998], GoogLeNet [Szegedy et al. 2015], ResNet [He et al. 2016], DeconvNet [Noh et al. 2015], FCN [Long et al. 2015] ...

# O-CNN for Shape Analysis

- Shape classification and retrieval
  - LeNet [Lecun et al. 1998]
- Shape segmentation
  - DeconvNet [Noh et al. 2015] + DenseCRF [Krähenbühl and Pfister 2011]



# Efficiency of O-CNN

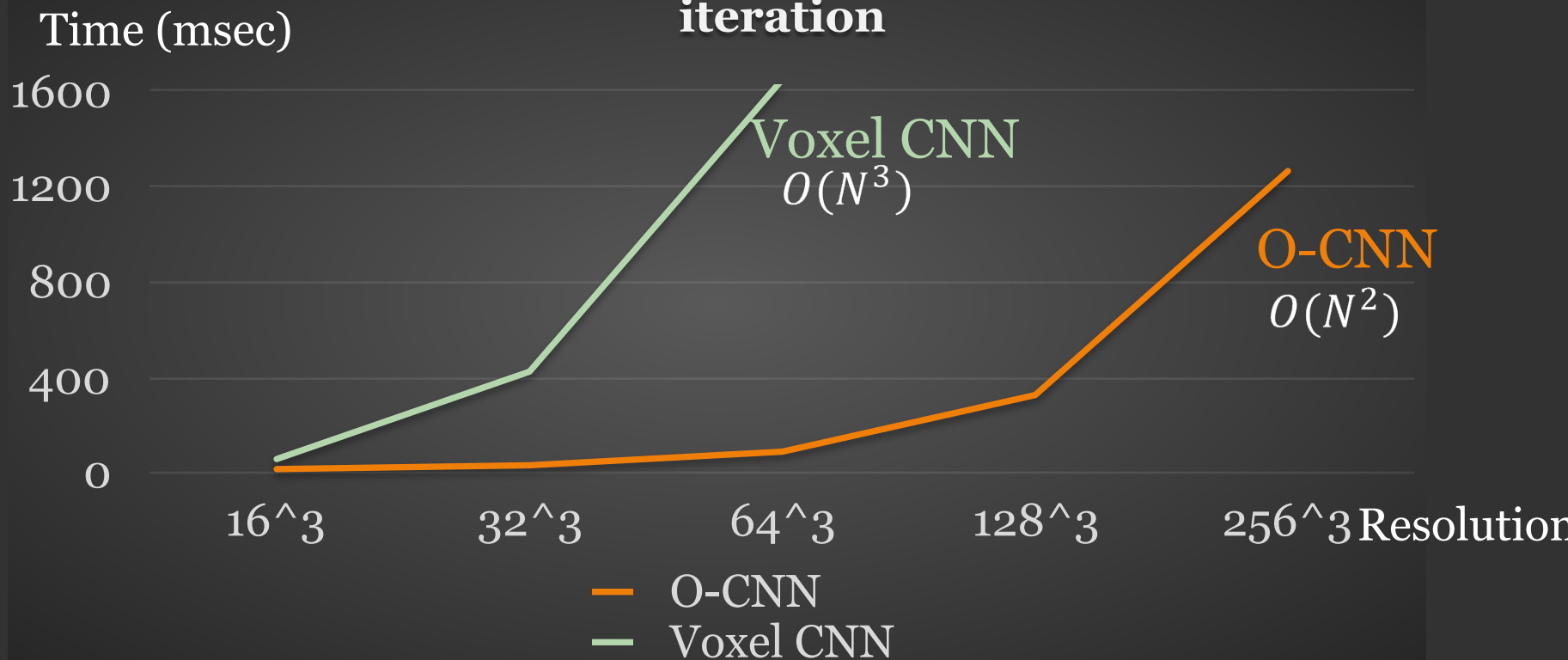
---

- **O-CNN** vs. full voxel CNN
  - Geforce 1080 GPU (8GB); Batch size 32

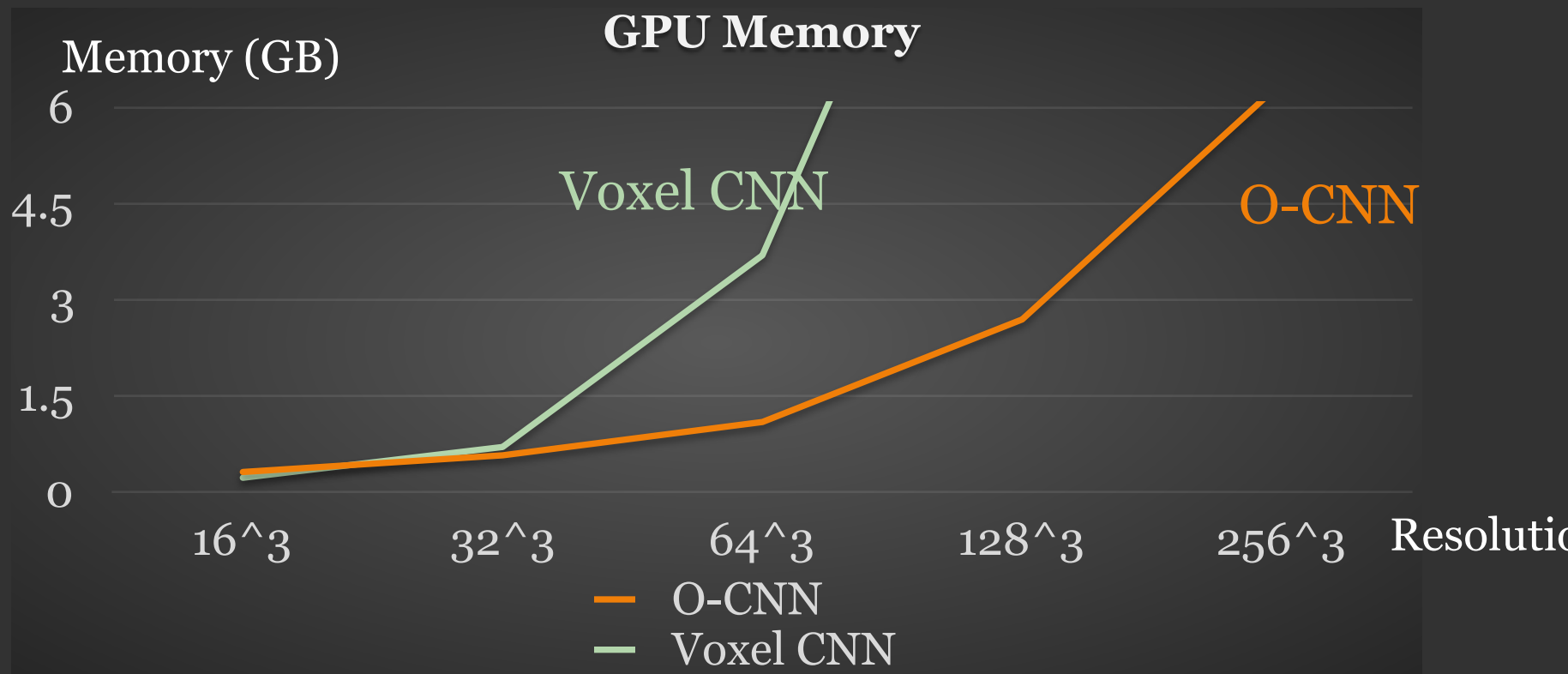


# Computational Efficiency

Average time for each forward and backward iteration



# Memory Efficiency





# Results – Classification

- **Task:** recognize the shape category
- **Dataset:** Princeton ModelNet40, 12311 3D models, 40 categories
- **Evaluation metric:** classification accuracy



Network	without voting
VoxNet ( $32^3$ )	82.0%
Geometry image	83.9%
SubVolSup ( $32^3$ )	87.2%
FPNN ( $64^3$ )	87.5%
FPNN+normal( $64^3$ )	88.4%
PointNet	89.2%
VRN ( $32^3$ )	89.0%
O-CNN(3)	85.5%
O-CNN(4)	88.3%
O-CNN(5)	89.6%
O-CNN(6)	<b>89.9%</b>
O-CNN(7)	89.5%
O-CNN(8)	89.6%

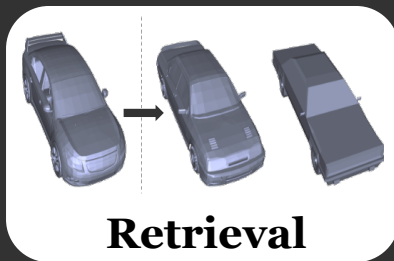
$8^3$   
 ↓  
 $256^3$

O-CNN(8)	88.9%
O-CNN(Δ)	88.2%
O-CNN(ϕ)	88.8%
O-CNN(2)	88.9%
O-CNN(1)	88.3%

# Results – Shape Retrieval



- **Task:** Given a query shape, retrieve similar shapes from the database
- **Dataset:** ShapeNet55 Core, 51190 3D models, 55 categories
- **Evaluation metric:** precision, recall, mAP, F-score, and NDCG

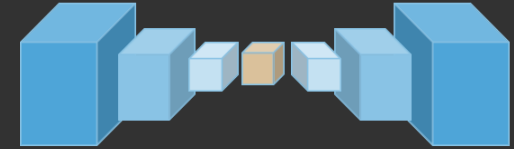


Method	P@N	R@N	F1@N	mAP	NDCG@N
Tatsuma_DB	0.427	0.689	0.472	0.728	0.875
Wang_CCMLT	0.718	0.350	0.391	0.823	0.886
Li_ViewAggr	0.508	<b>0.868</b>	0.582	0.829	0.904
Bai_GIFT	0.706	0.695	0.689	0.825	0.896
Su MVCNN	0.770	0.770	0.764	0.873	0.899
O-CNN(5)	0.768	0.769	0.763	0.871	0.904
O-CNN(6)	<b>0.778</b>	0.782	<b>0.776</b>	<b>0.875</b>	<b>0.905</b>

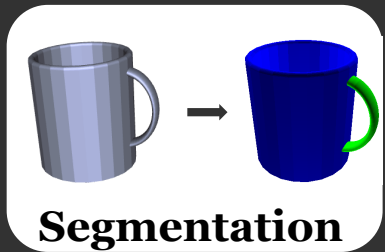
O-CNN(6)	0.778	0.782	0.776	0.875	0.905
O-CNN(6)	0.778	0.782	0.776	0.875	0.905



# Results – Segmentation



- **Task:** Segment a 3D shape into semantic parts
- **Dataset:** dataset from [Yi et al. 2016], 16881 models, 2~6 parts
- **Evaluation metric:** Intersection over Union



	mean	plane	bag	cap	car	chair	e.ph.	guitar	knife
# shapes		2690	76	55	898	3758	69	787	392
[Yi et al. 2016]	81.4	81.0	78.4	77.7	75.7	87.6	61.9	92.0	85.4
PointNet [Qi et al. 2017]	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9
SpecCNN [Yi et al. 2017]	84.7	81.6	81.7	81.9	75.2	90.2	74.9	93.0	86.1
O-CNN(5)	85.2	84.2	86.9	84.6	74.1	90.8	81.4	91.3	87.0
O-CNN(6)	<b>85.9</b>	<b>85.5</b>	<b>87.1</b>	<b>84.7</b>	<b>77.0</b>	<b>91.1</b>	<b>85.1</b>	<b>91.9</b>	<b>87.4</b>
O-CNN(6)	82.8	82.2	84.1	84.4	74.0	91.1	82.1	91.8	84.4
O-CNN(6)	82.5	84.5	80.8	84.0	74.1	90.8	81.4	91.3	84.0

# Conclusion

---

- Key idea
  - Store sparse surface signal
  - Constrain the computation near surface
- Octree based 3D CNNs
  - General, efficient, and effective



Code and data online

<http://wang-ps.github.io/O-CNN>

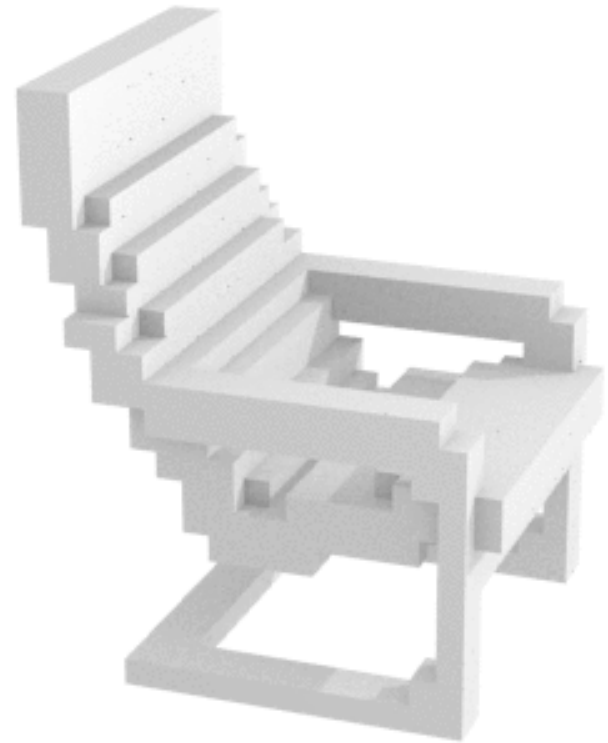
UC San Diego

# Shape Reconstruction

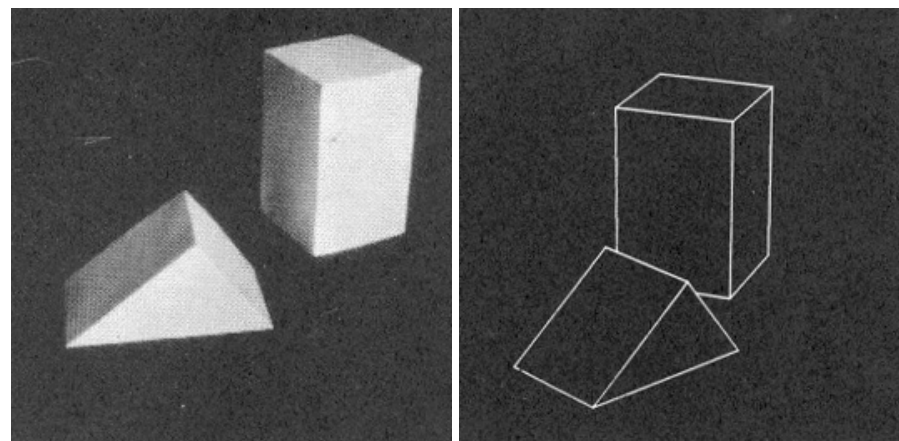
How do we learn to perceive 3D ?



How do we learn to perceive 3D ?



# Single-view Reconstruction

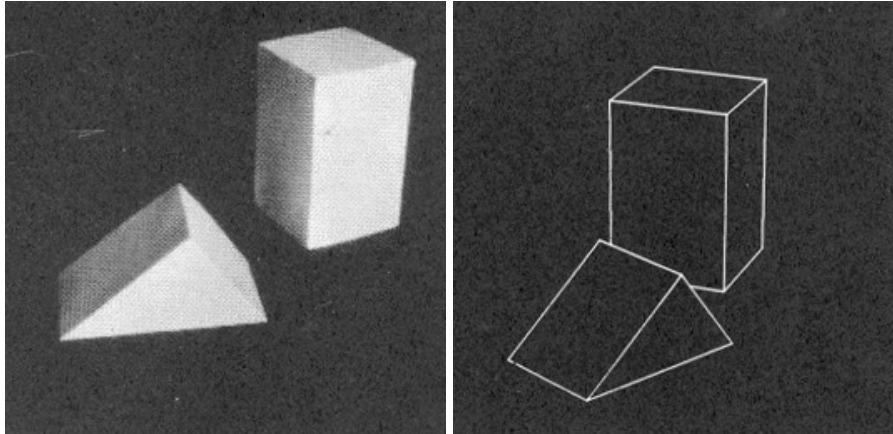


Roberts. PhD Thesis, MIT. 1963

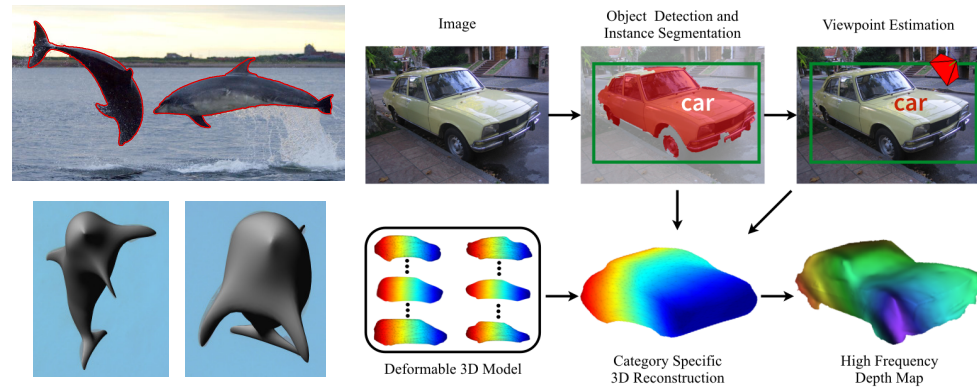
Unsupervised



# Single-view Reconstruction

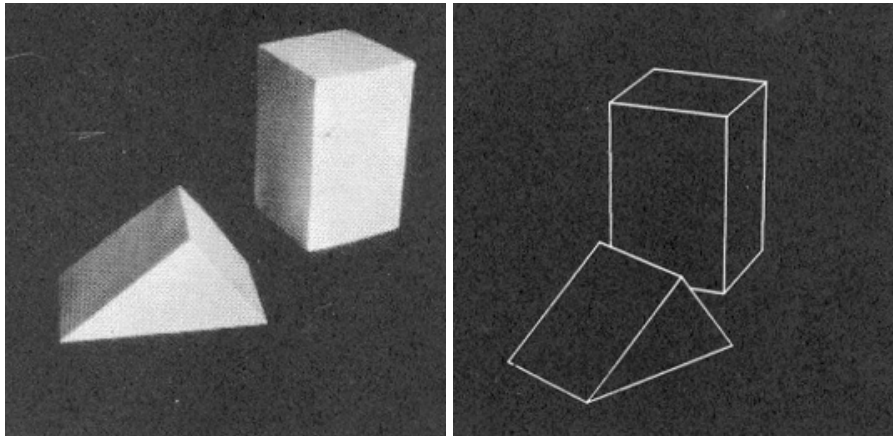


Roberts. PhD Thesis, MIT. 1963  
Unsupervised



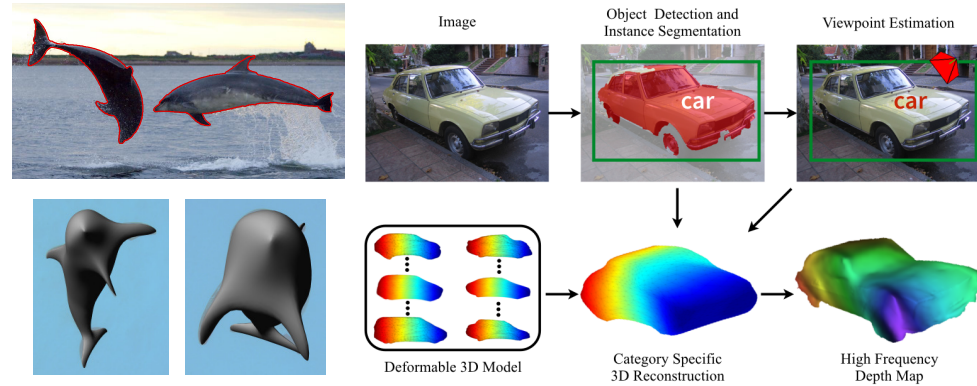
Cashman & Fitzgibbon, PAMI 2013  
Kar et al., CVPR 2015  
Supervision : Masks + Pose

# Single-view Reconstruction



Roberts. PhD Thesis, MIT. 1963

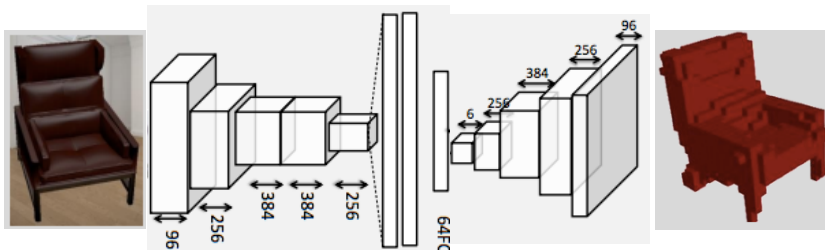
Unsupervised



Cashman & Fitzgibbon, PAMI 2013

Kar et al., CVPR 2015

Supervision : Masks + Pose



Loss



Choy et al., Girdhar et al.

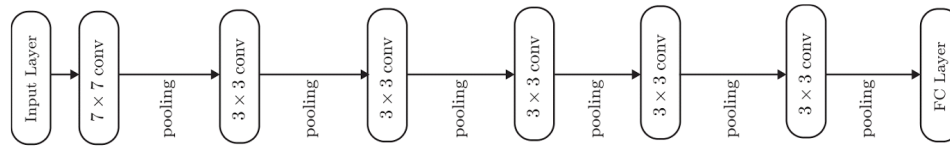
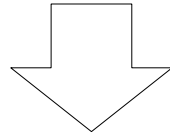
ECCV 2016

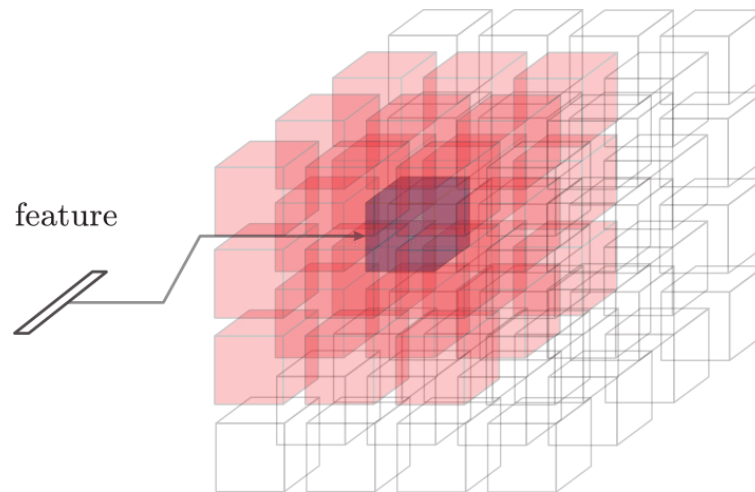
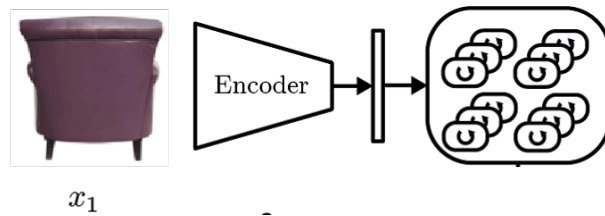
Supervision : Ground-truth 3D



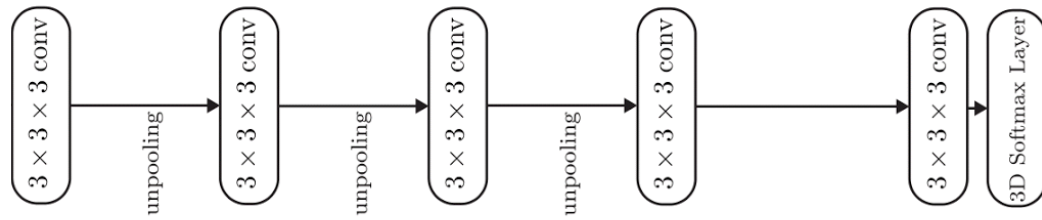
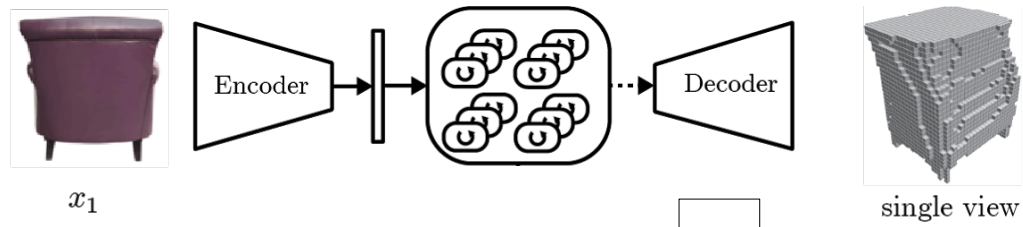


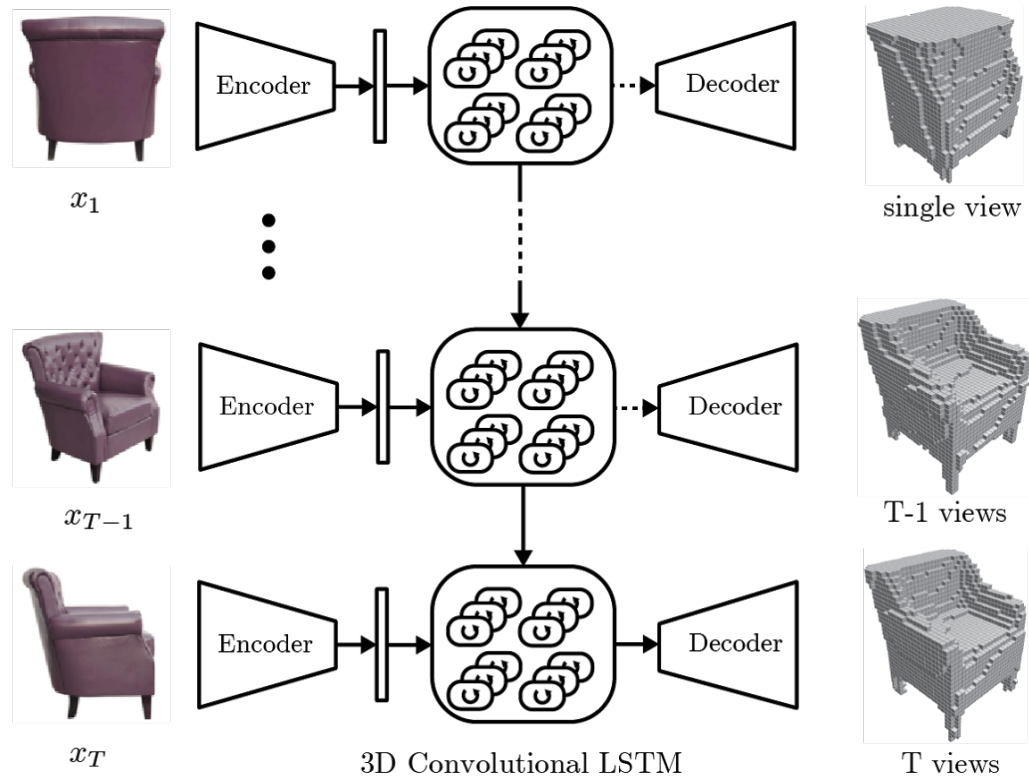
$x_1$



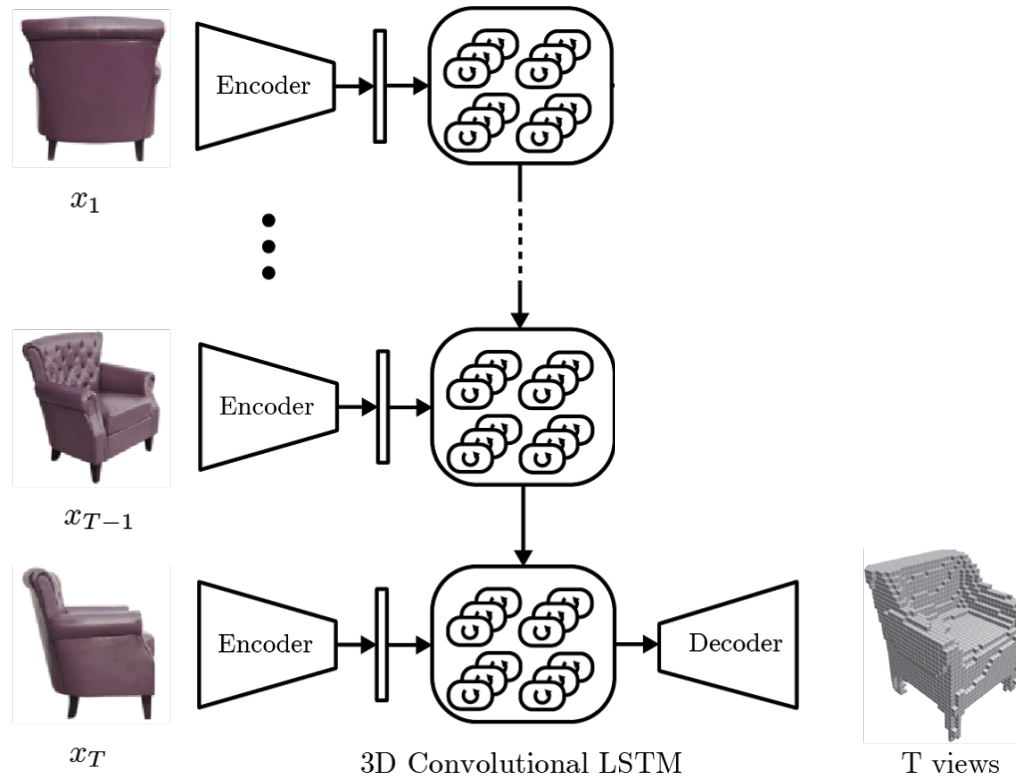


## 3D Convolutional LSTM

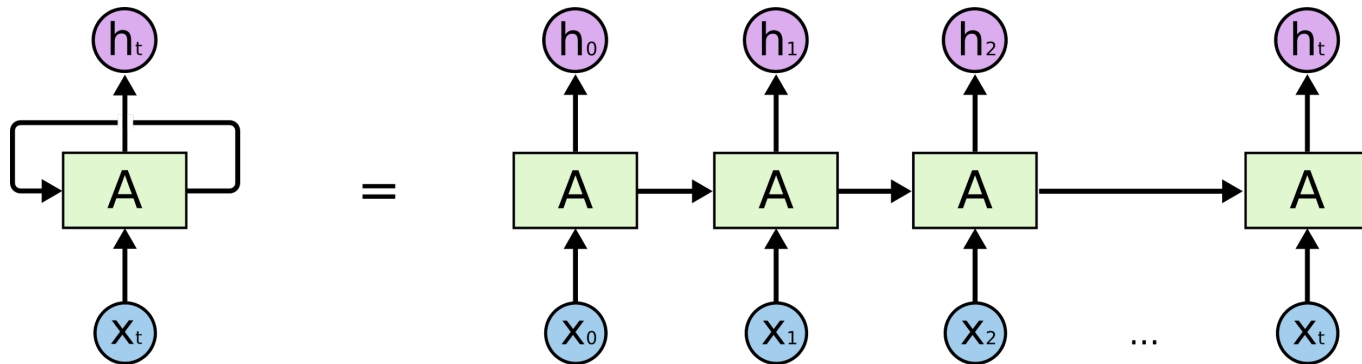




It is possible to aggregate information from multiple views

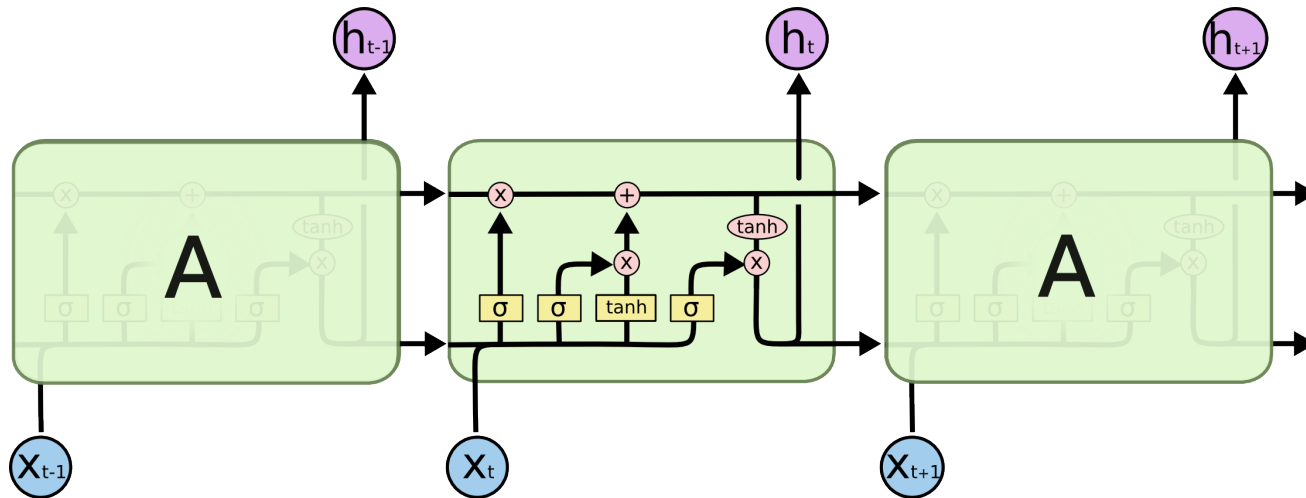


# Recurrent Neural Network



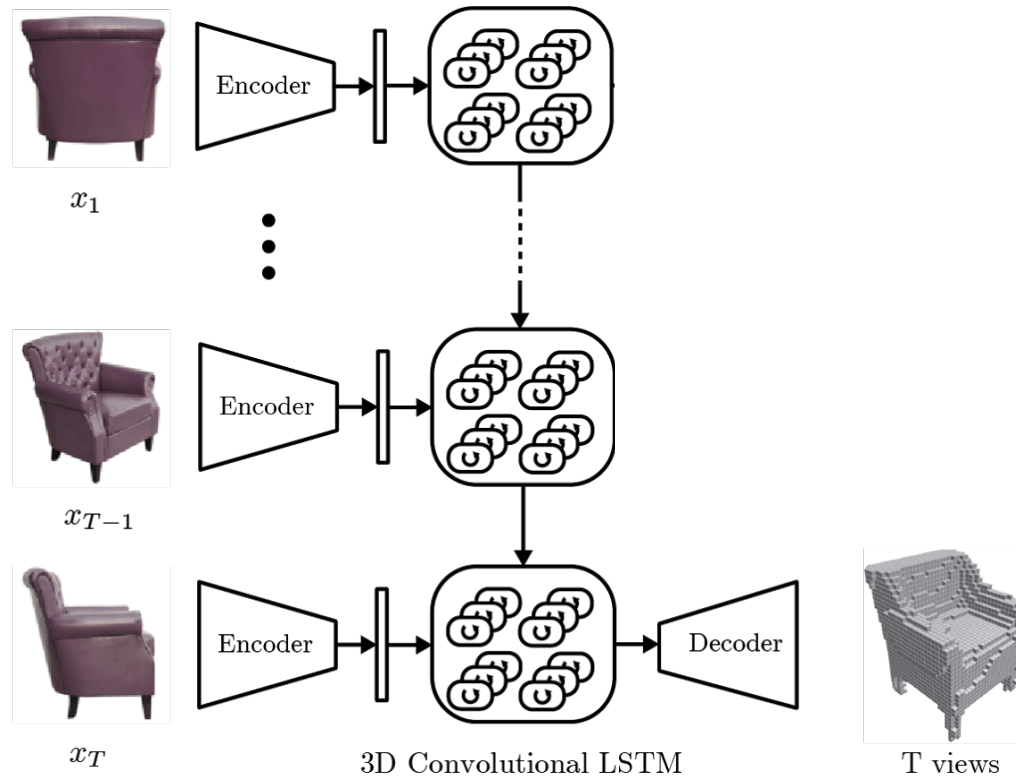
[Christopher Olah] Understanding LSTM Networks, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

# Long Short Term Memory



[Christopher Olah] Understanding LSTM Networks, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

It is possible to aggregate information from multiple views





# Training

- ShapeNet
  - 50k CAD models
  - Render from arbitrary views
  - Random number of images w/ random order
  - Random background, translation

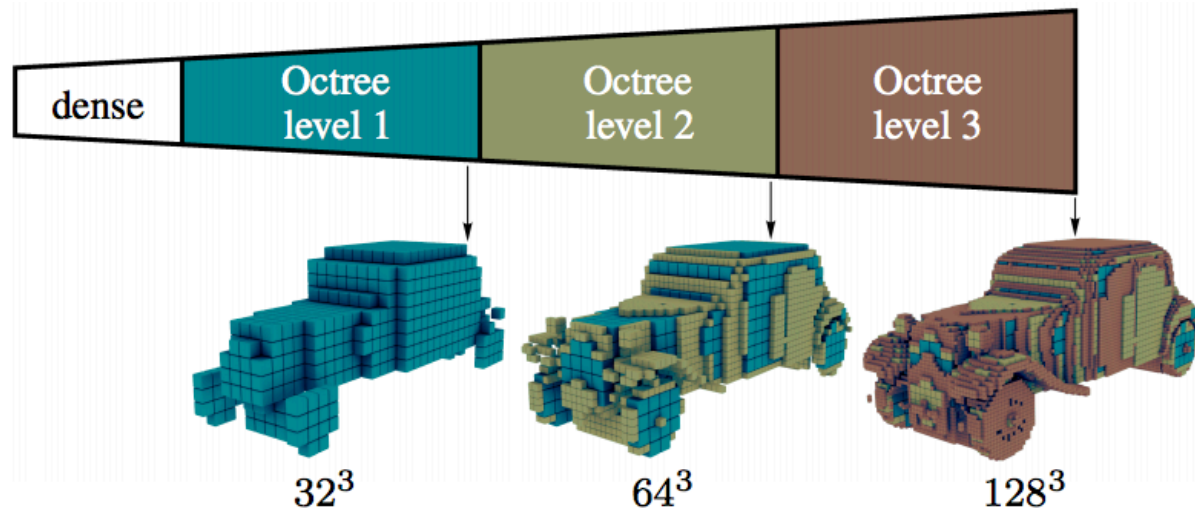
$$L(\mathcal{X}, y) = \sum_{i,j,k} y_{(i,j,k)} \log(p_{(i,j,k)}) + (1 - y_{(i,j,k)}) \log(1 - p_{(i,j,k)})$$

- Voxel-wise cross entropy loss





# Towards higher spatial resolution

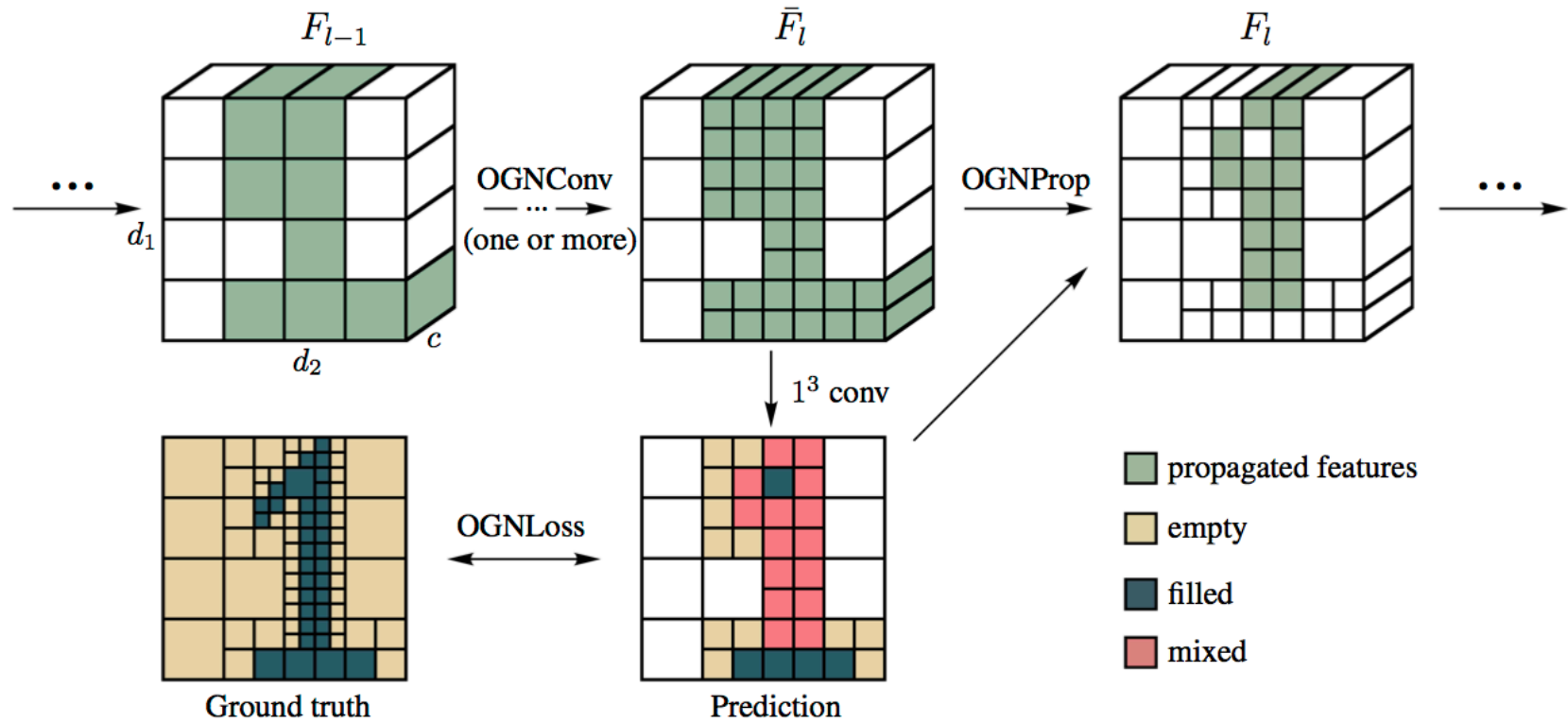


Maxim Tatarchenko, Alexey Dosovitskiy, Thomas Brox

**“Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs”**

*arxiv (March, 2017)*

# Progressive voxel refinement



# Results

Input

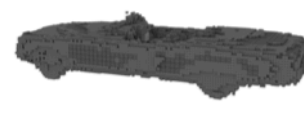
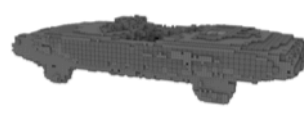
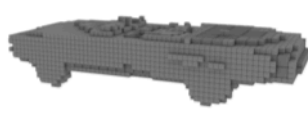
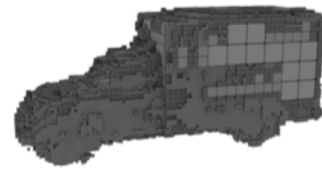
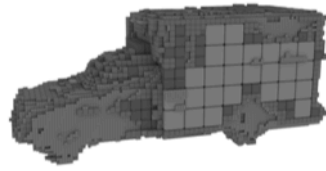
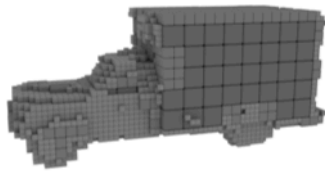
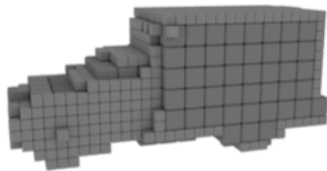
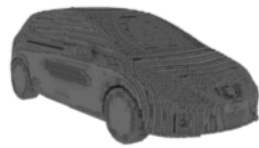
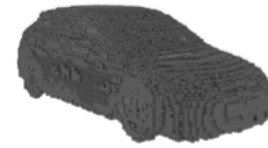
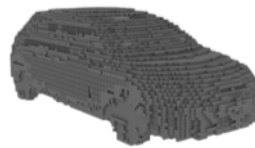
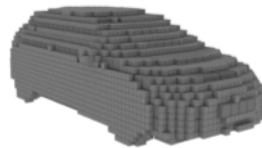
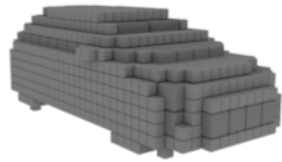
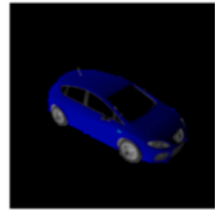
$32^3$

$64^3$

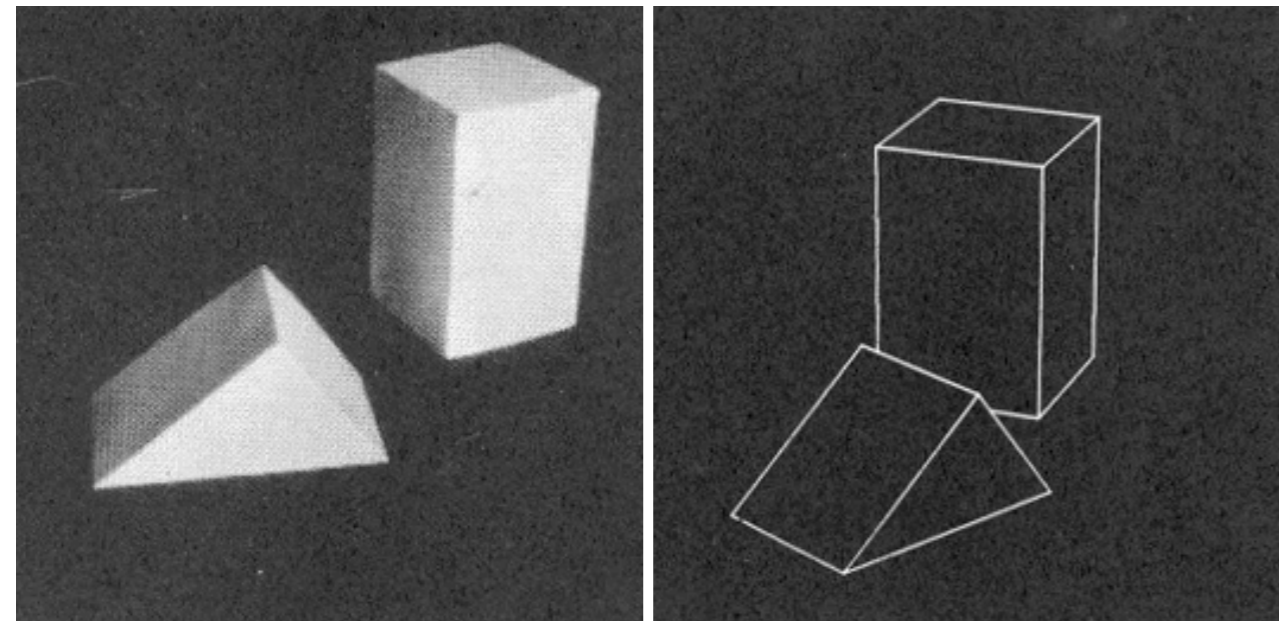
$128^3$

$256^3$

GT  $256^3$

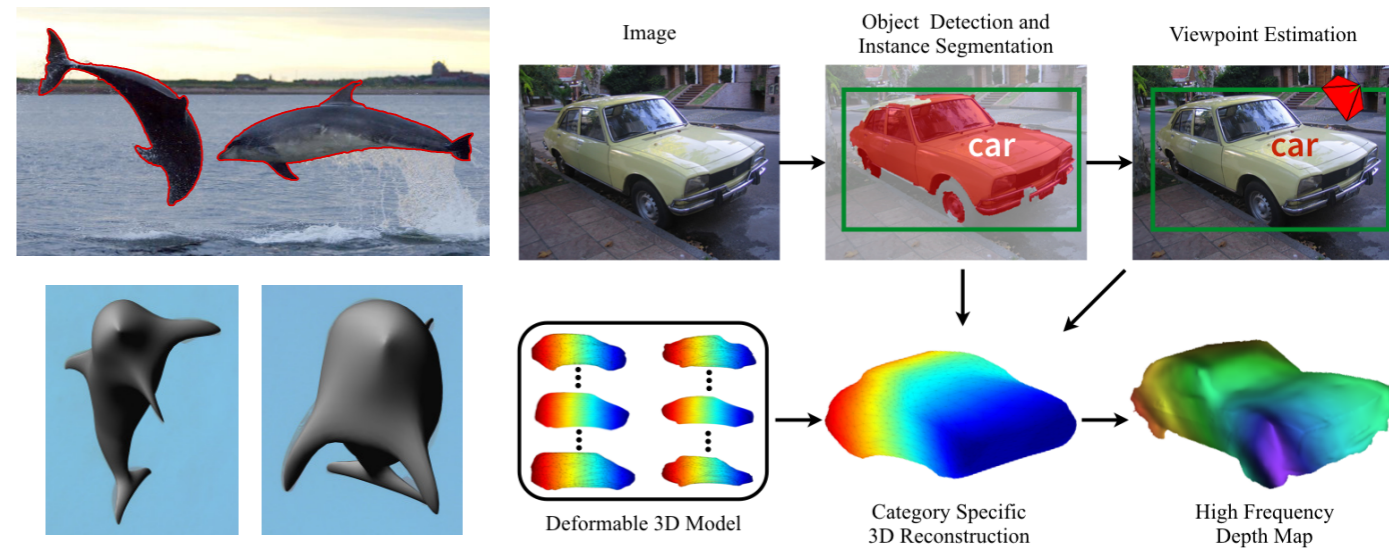


# Single-view Reconstruction



Roberts. PhD Thesis, MIT. 1963

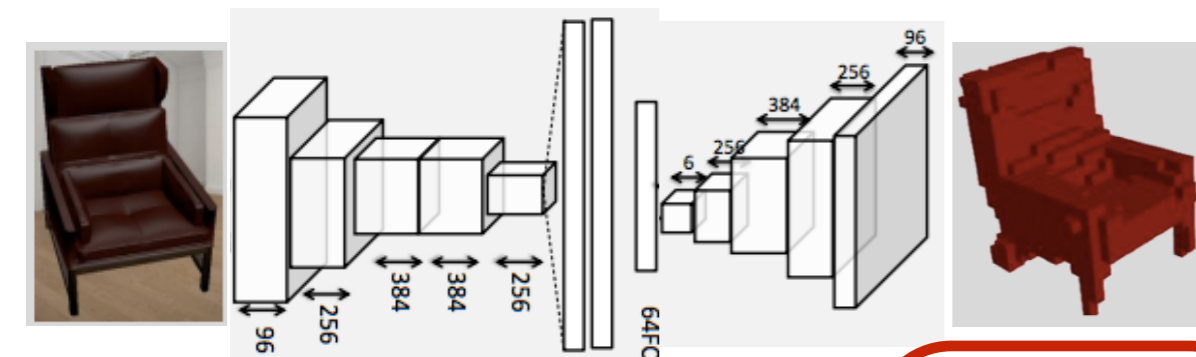
Unsupervised



Cashman & Fitzgibbon, PAMI 2013

Kar et al., CVPR 2015

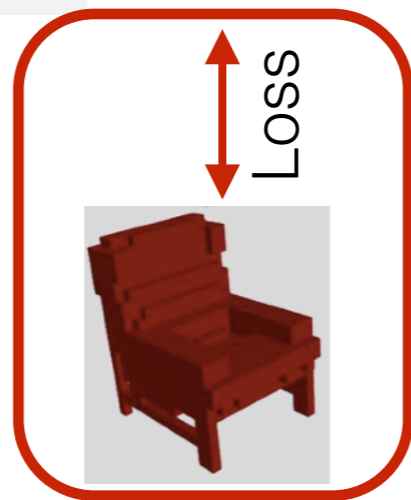
Supervision : Masks + Pose



Choy et al., Girdhar et al.

ECCV 2016

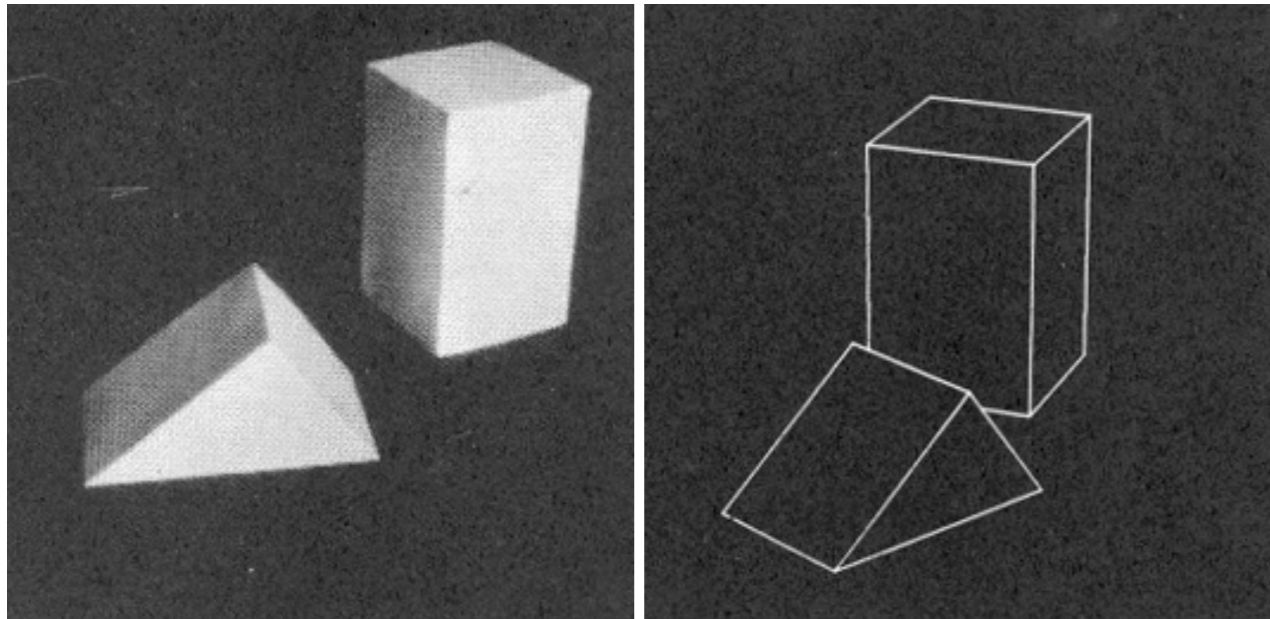
Supervision : Ground-truth 3D



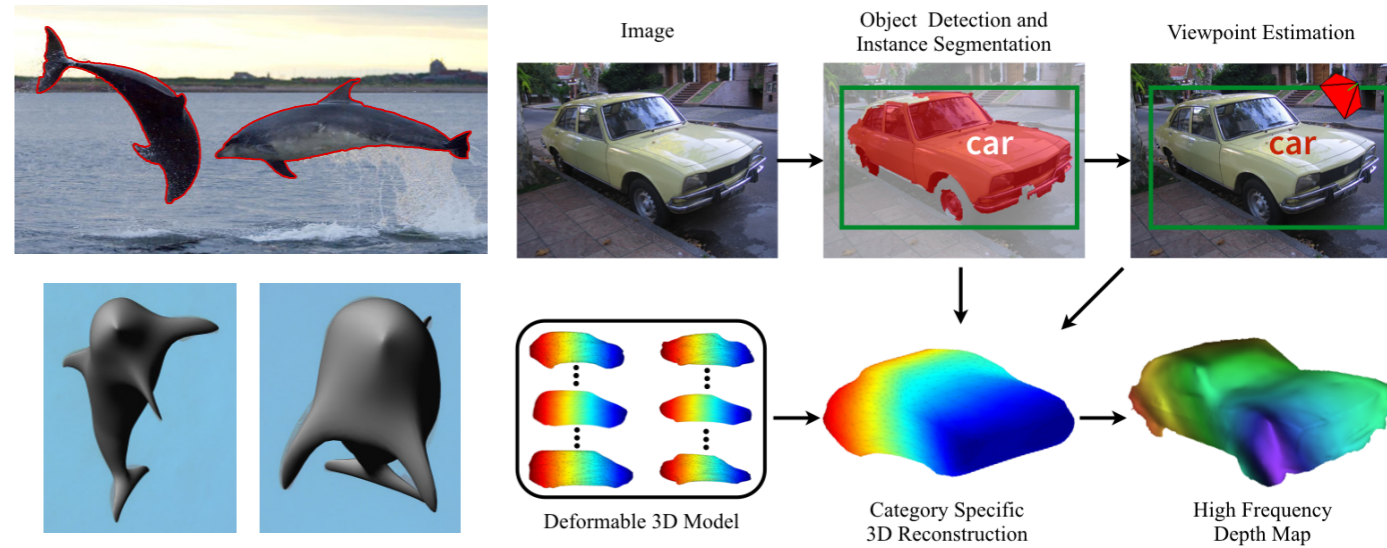
But we don't have ground-truth 3D !



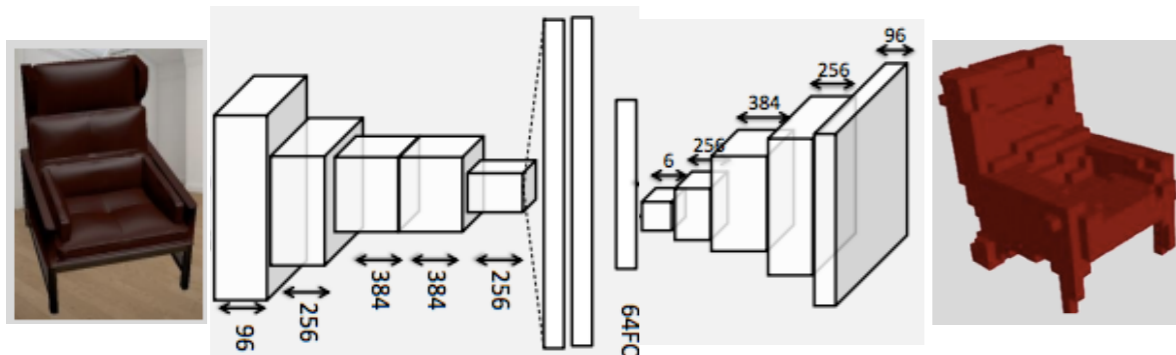
# Single-view Reconstruction



Roberts. PhD Thesis, MIT. 1963  
Unsupervised



Cashman & Fitzgibbon, PAMI 2013  
Kar et al., CVPR 2015  
Supervision : Masks + Pose



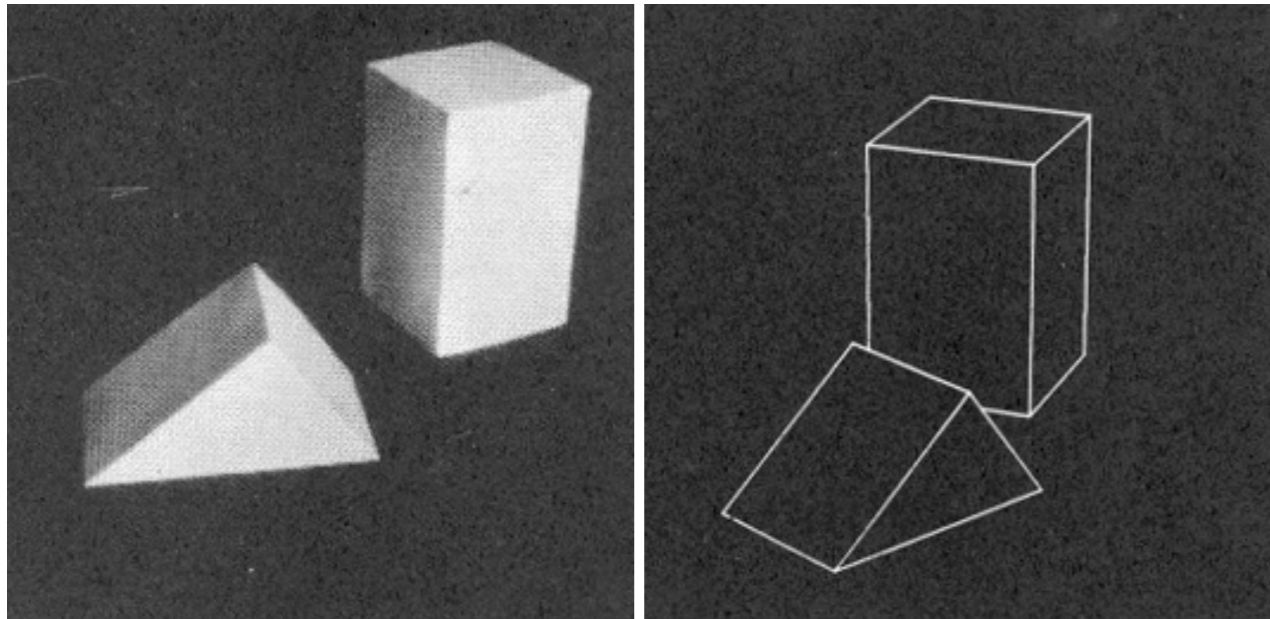
SSOT

Choy et al., Girdhar et al.  
ECCV 2016  
Supervision : Ground-truth 3D



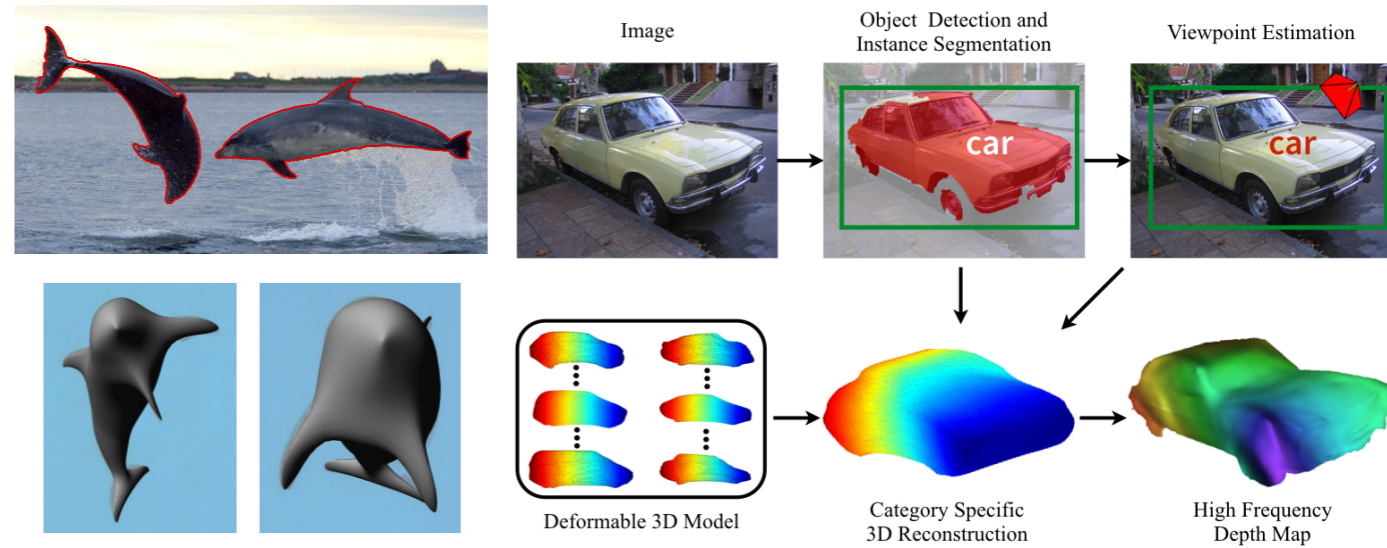


# Single-view Reconstruction



Roberts. PhD Thesis, MIT. 1963

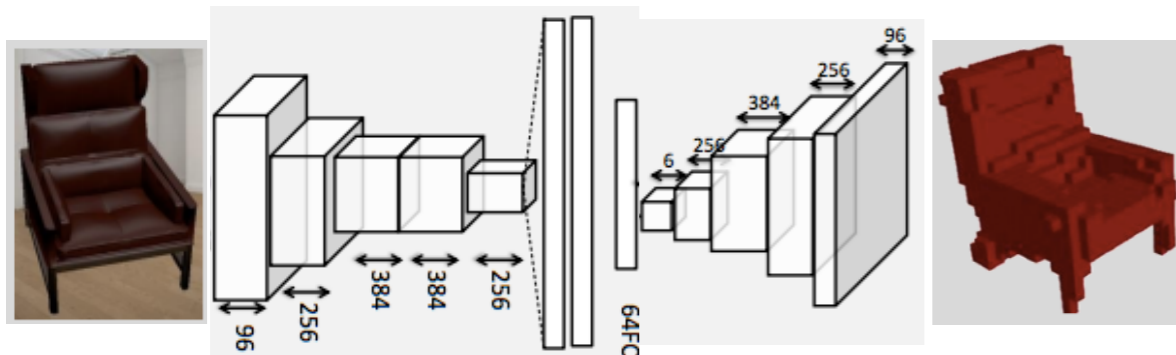
Unsupervised



Cashman & Fitzgibbon, PAMI 2013

Kar et al., CVPR 2015

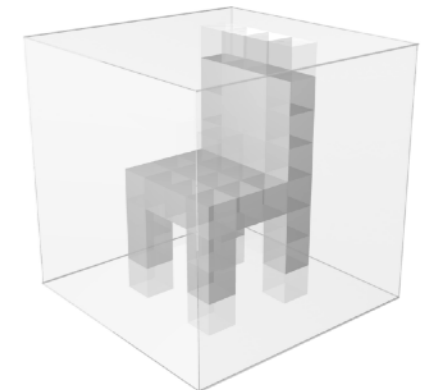
Supervision : Masks + Pose



Choy et al., Girdhar et al.

ECCV 2016

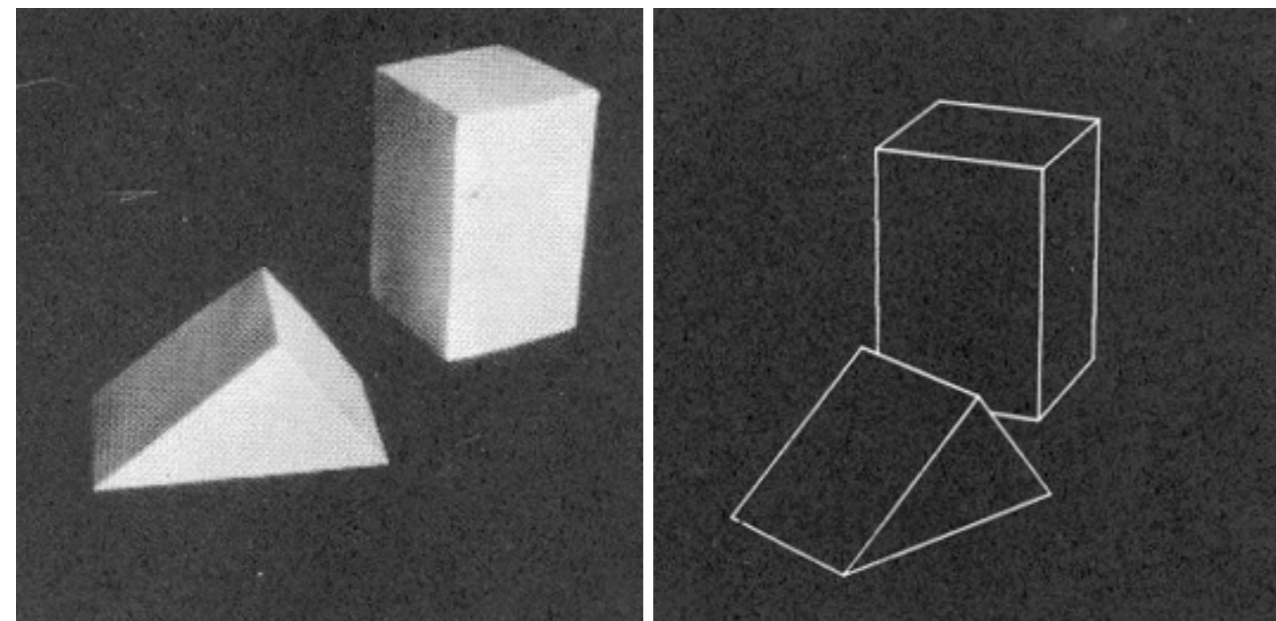
Supervision : Ground-truth 3D



Multi-view Supervision for Single-view Reconstruction via Differentiable Ray Consistency  
Tulsiani et al.

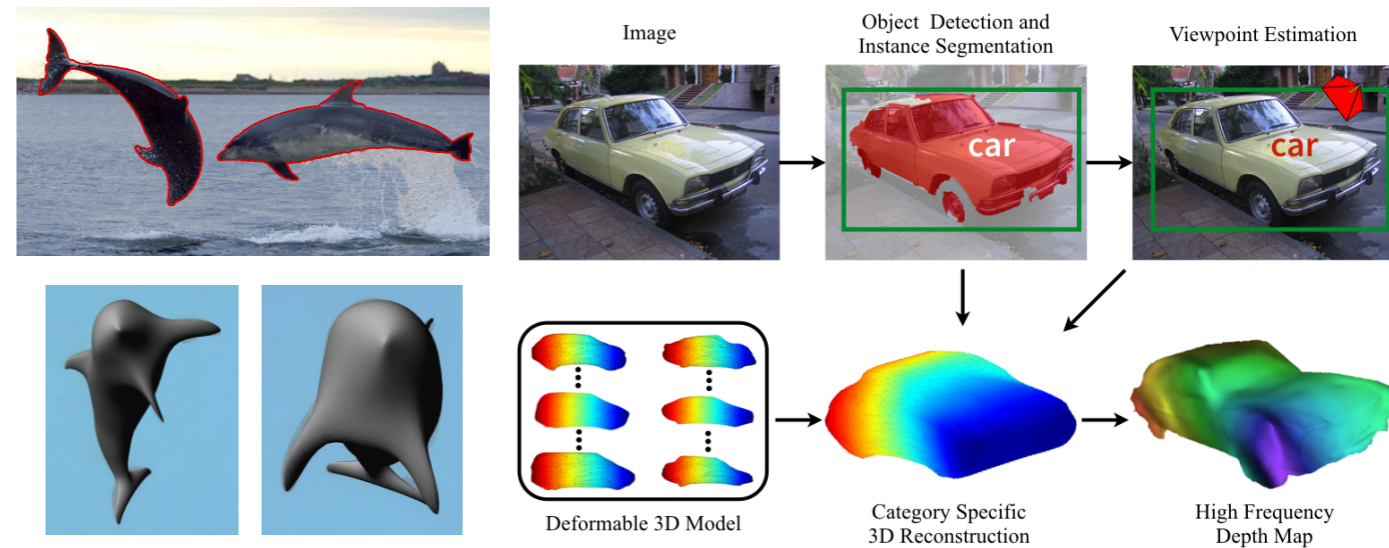
Supervision : Multi-view

# Single-view Reconstruction



Roberts. PhD Thesis, MIT. 1963

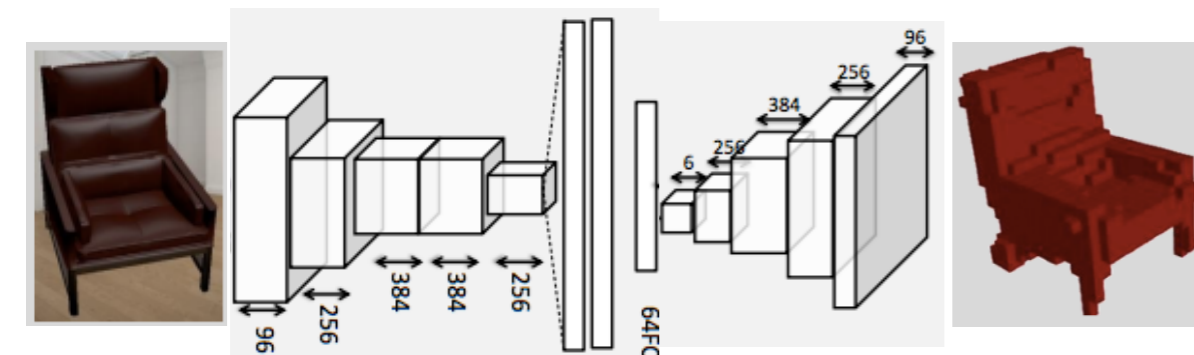
Unsupervised



Cashman & Fitzgibbon, PAMI 2013

Kar et al., CVPR 2015

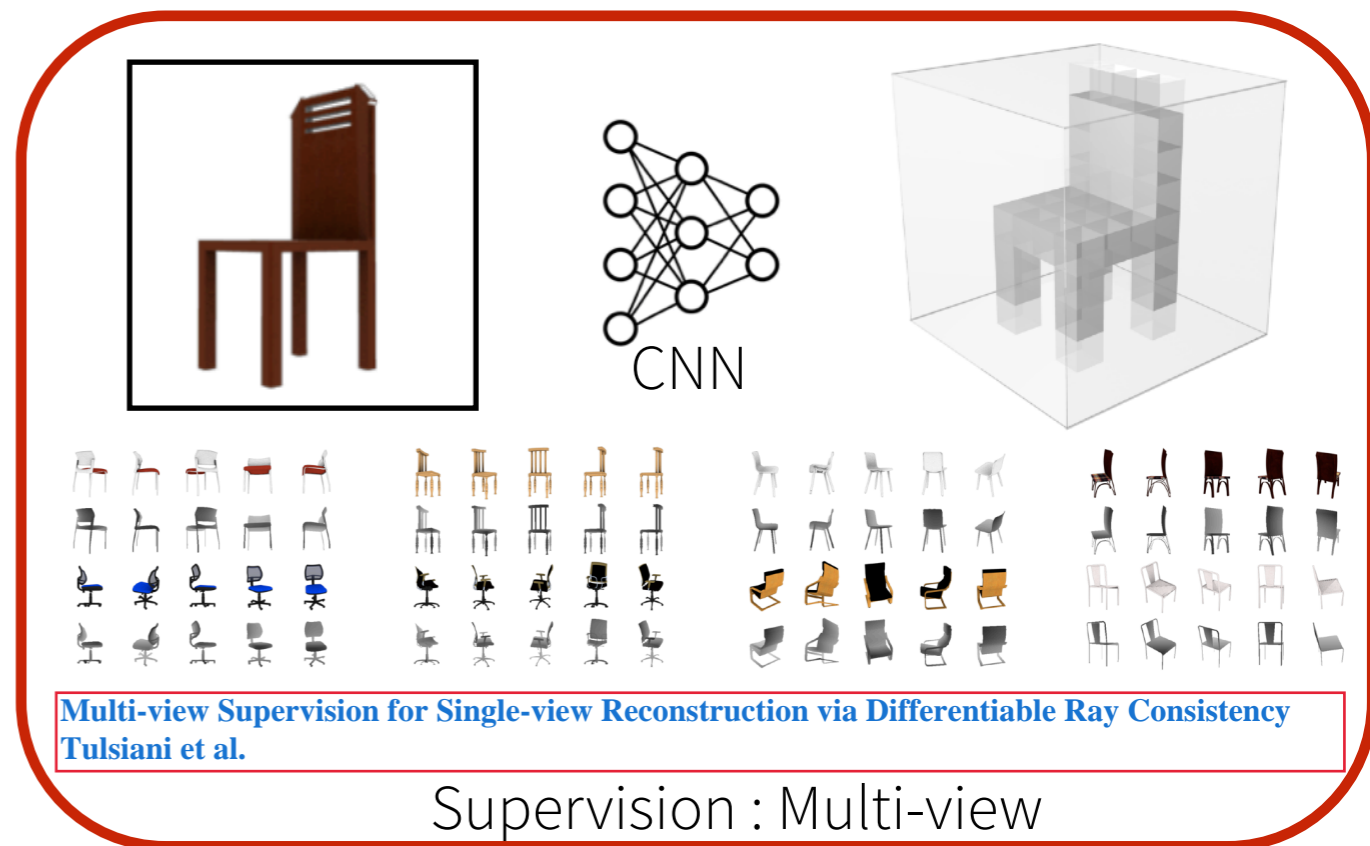
Supervision : Masks + Pose



Choy et al., Girdhar et al.

ECCV 2016

Supervision : Ground-truth 3D

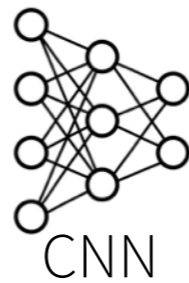


Supervision : Multi-view

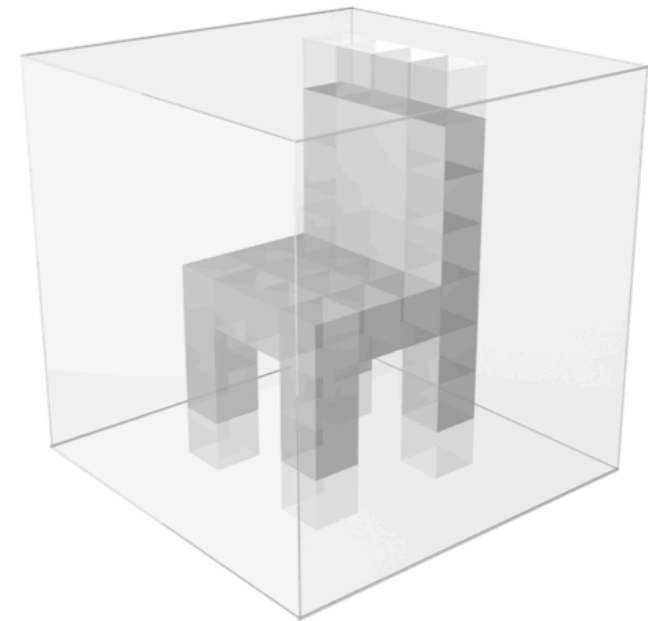
# How to use Multi-view Supervision ?



Input Image



CNN

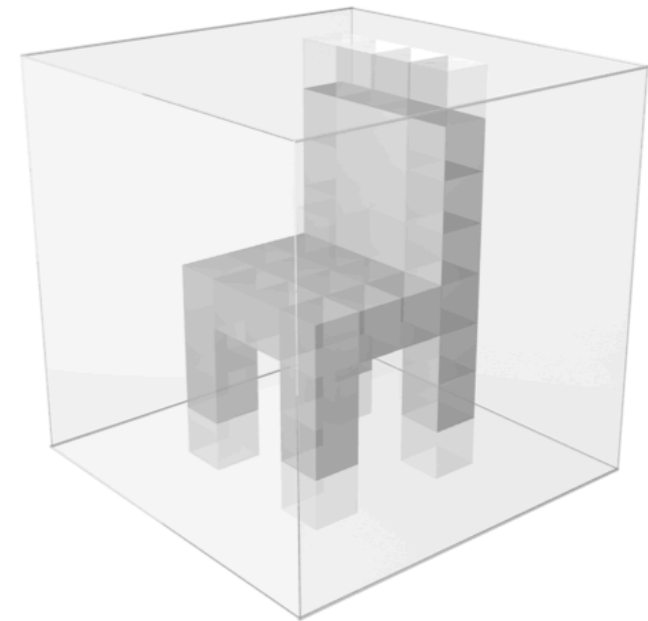




# How to use Multi-view Supervision ?

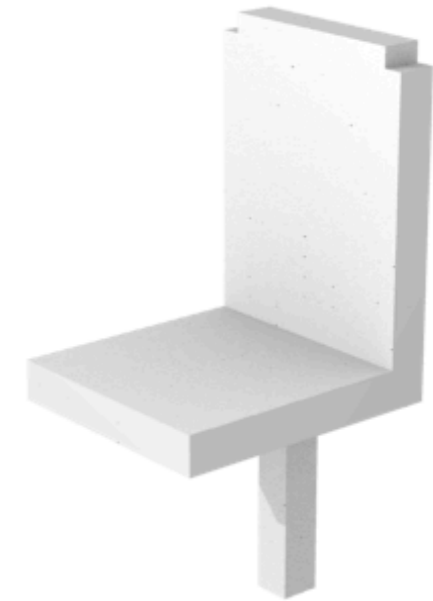
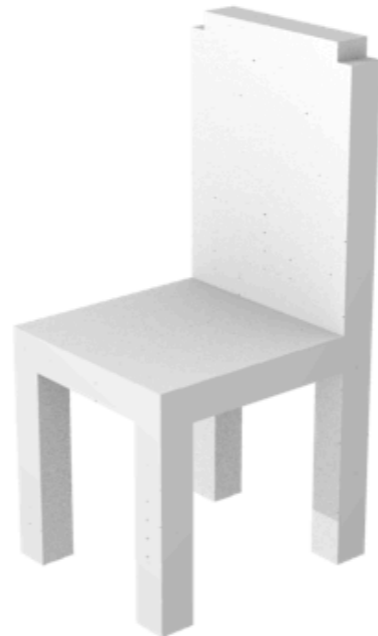
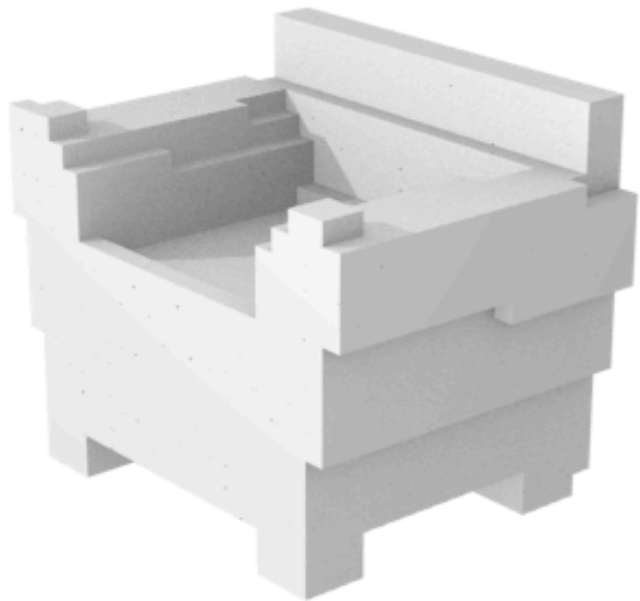


Input Image



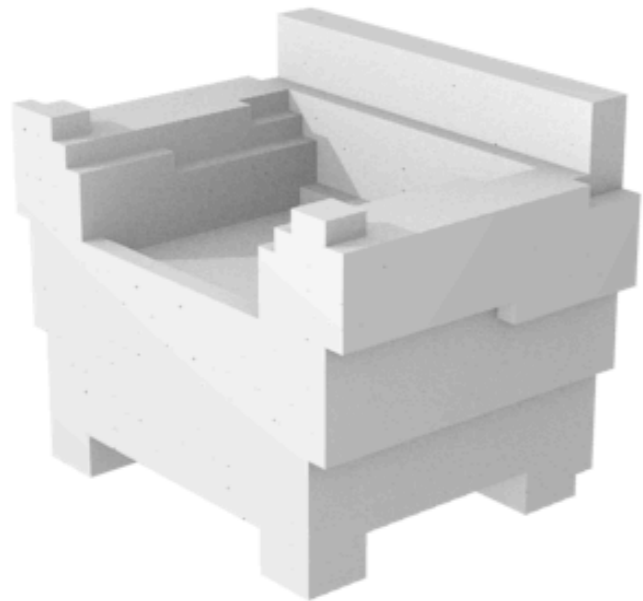
Observation **O**  
from camera **C**

# How to use Multi-view Supervision ?

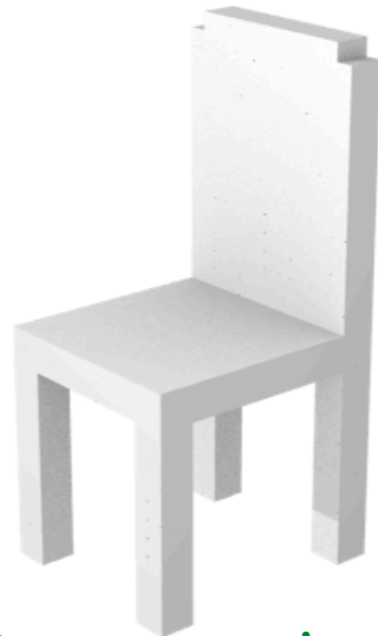


Observation **O**  
from camera **C**

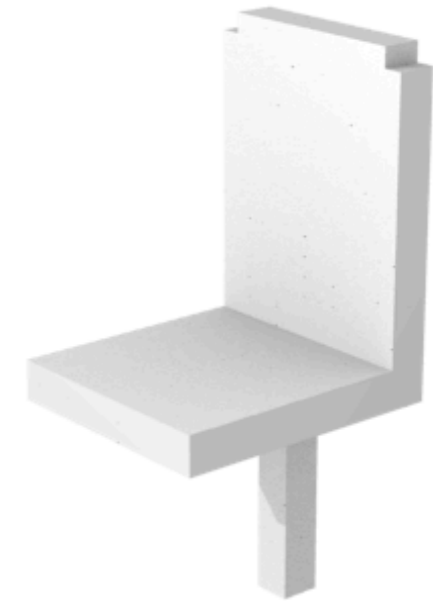
# How to use Multi-view Supervision ?



Geometrically  
Inconsistent



Geometrically  
Consistent



Geometrically  
Inconsistent

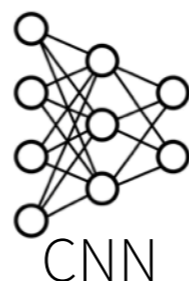


Observation **O**  
from camera **C**

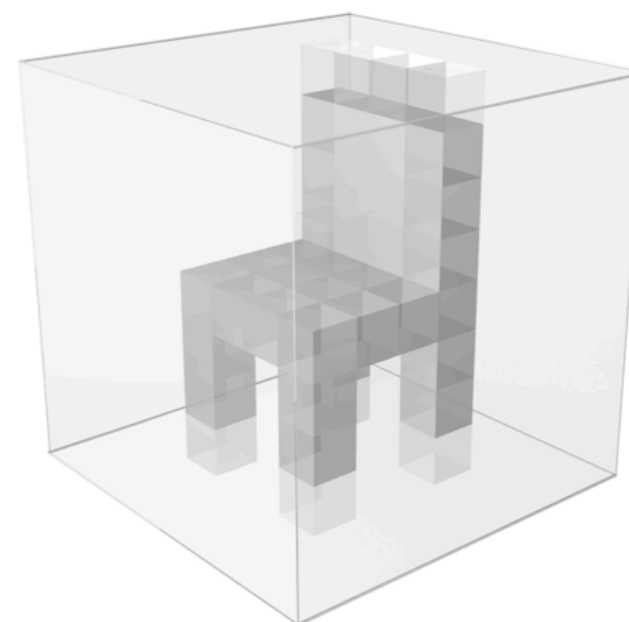
# Learning via Geometric Consistency



Input Image



CNN

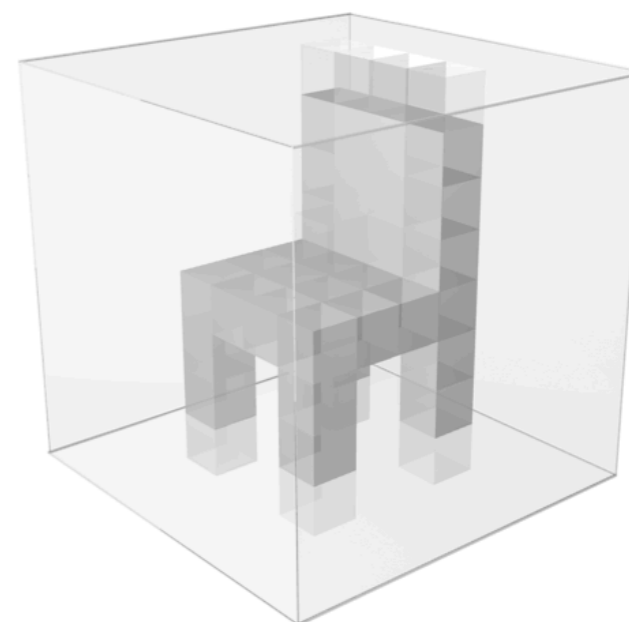


Observation **O**  
from camera **C**

# Learning via Geometric Consistency



Input Image



Observation  $\mathbf{O}$   
from camera  $\mathbf{C}$



$$L(\text{Image of Chair}, \text{3D Voxel Chair})$$

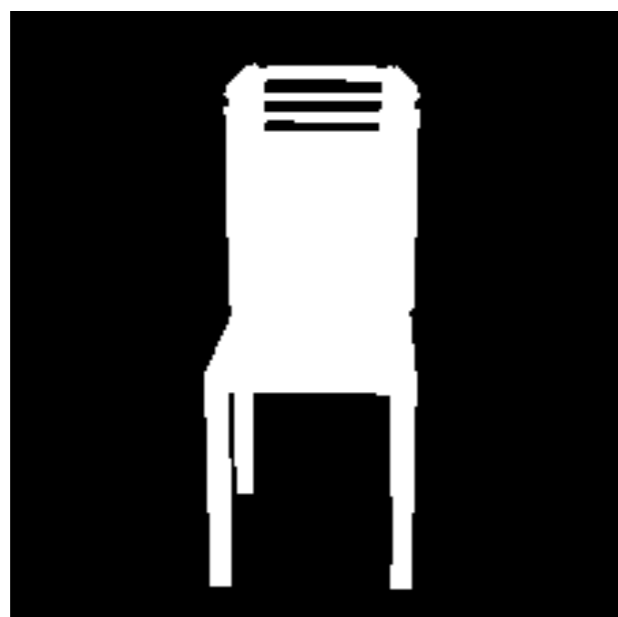
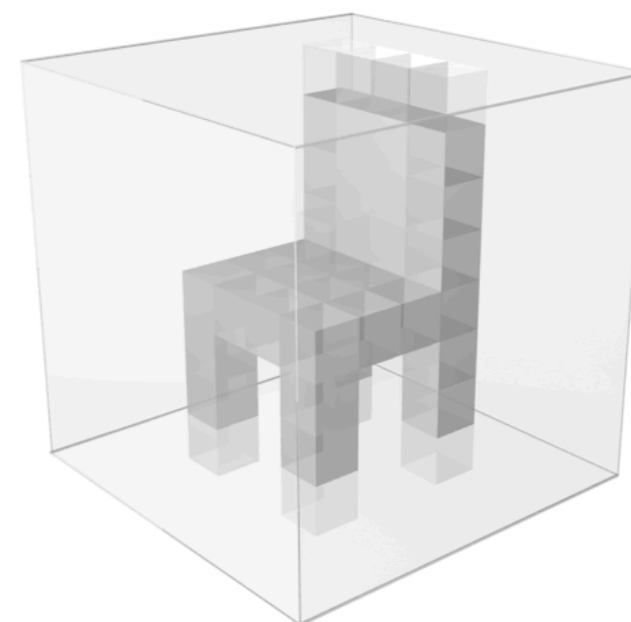
Geometric  
Consistency Loss



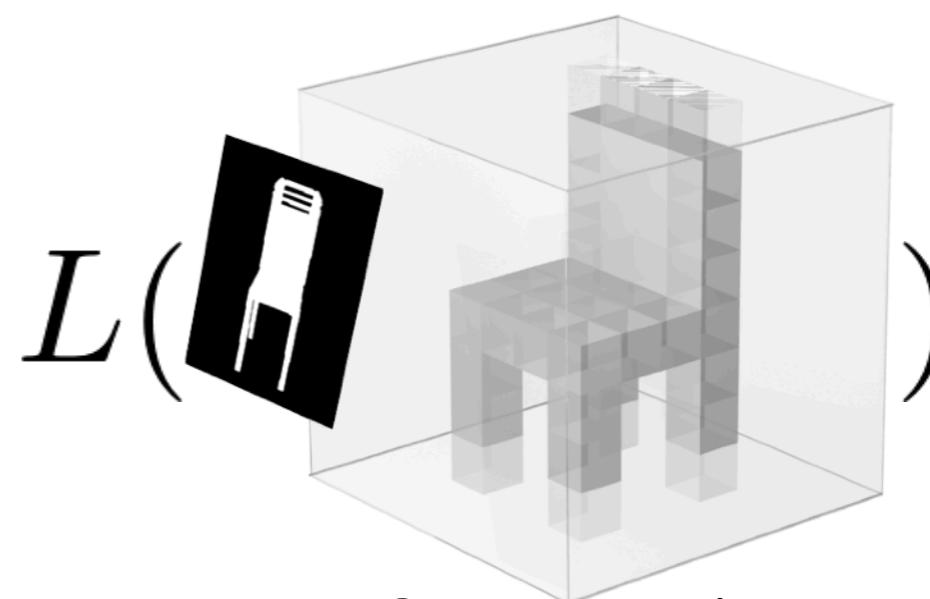
# Learning via Geometric Consistency



Input Image



Observation **O**  
from camera **C**

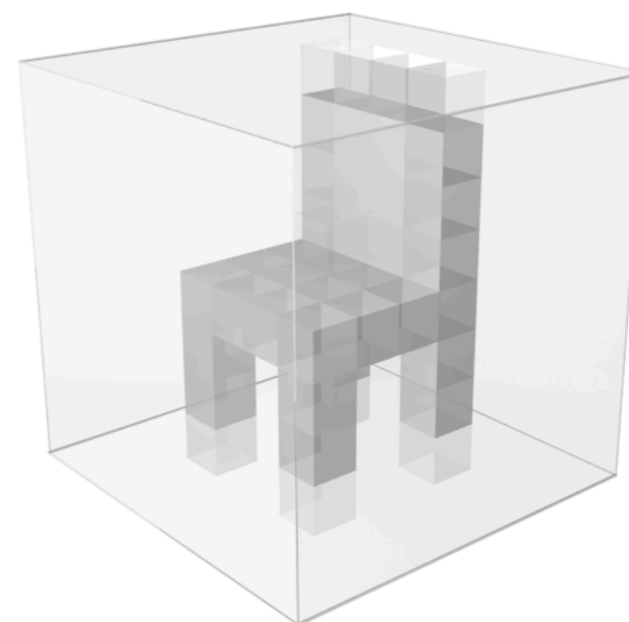
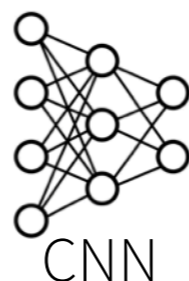


Geometric  
Consistency Loss

# Learning via Geometric Consistency



Input Image



$$L(\text{Input Image}, \text{3D Voxel Chair})$$

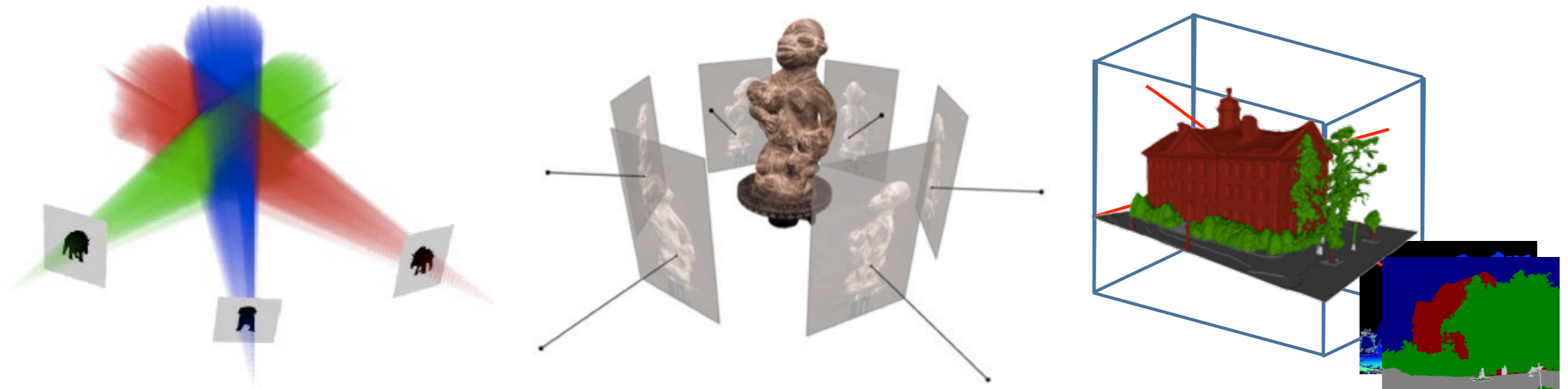
Geometric

Consistency Loss

Observation **O**  
from camera **C**

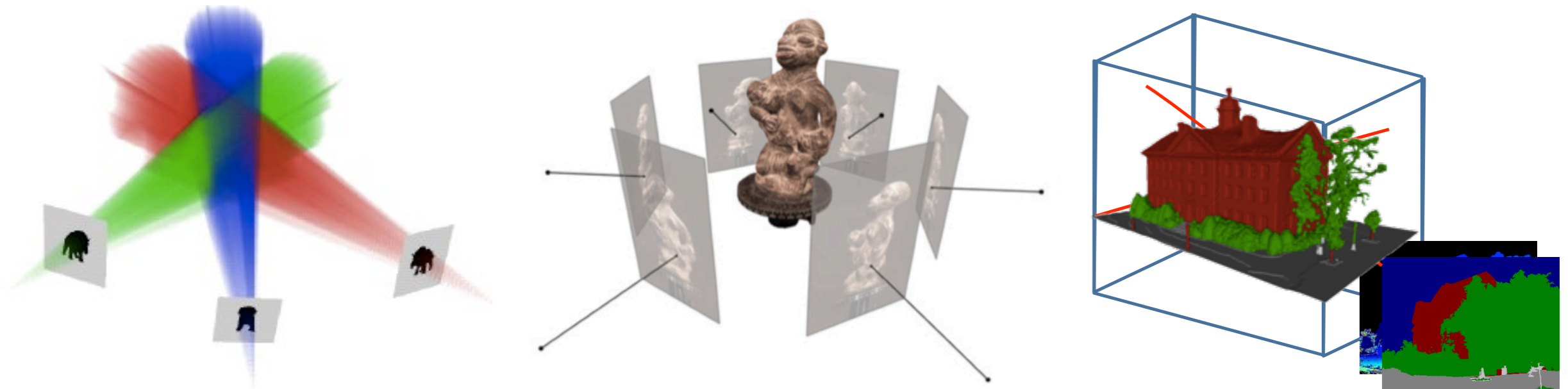
# 3D from Geometric Consistency

# 3D from Geometric Consistency

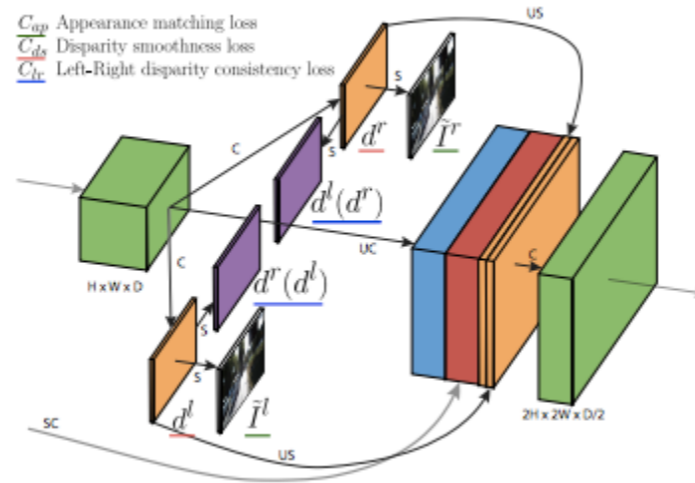
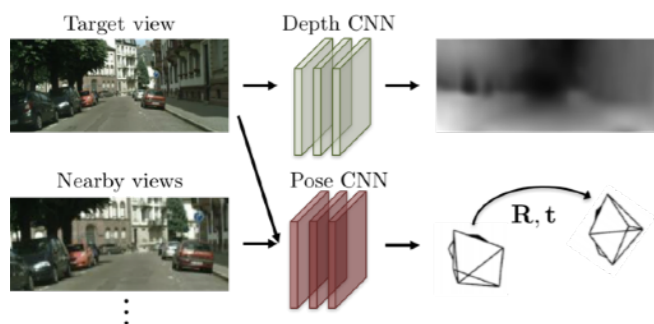
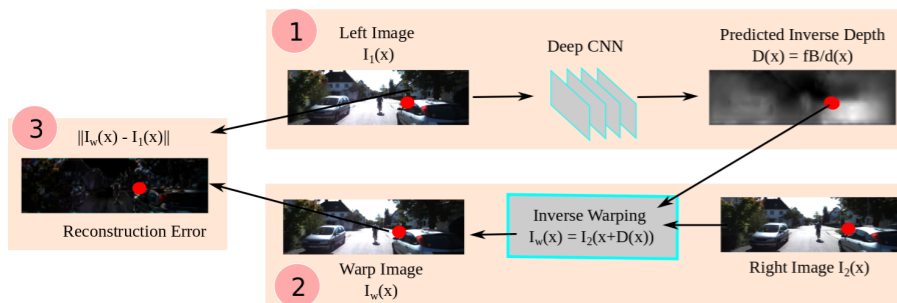


Space Carving, Multi-view Stereo, Multi-view Reconstruction

# 3D from Geometric Consistency



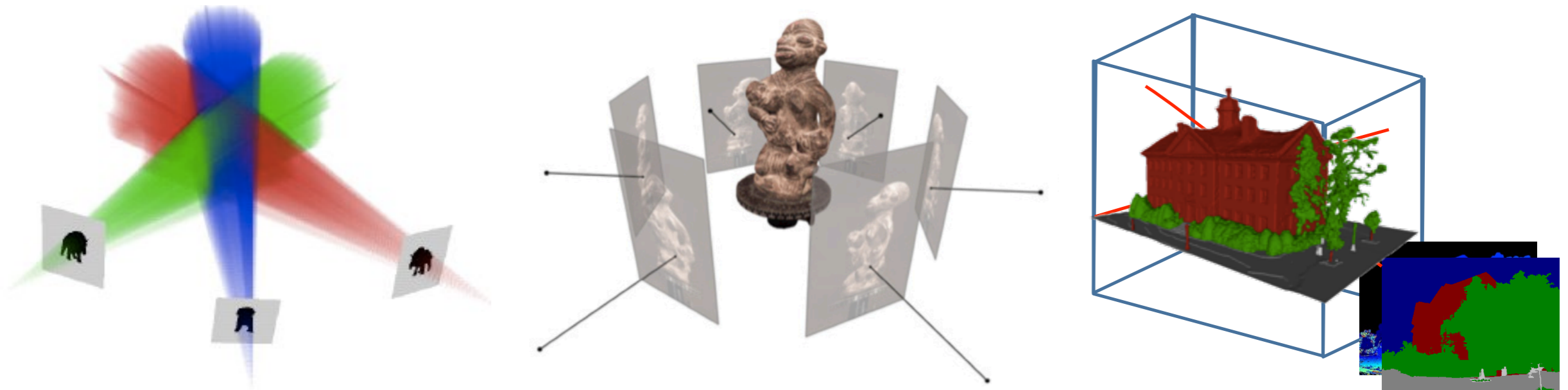
Space Carving, Multi-view Stereo, Multi-view Reconstruction



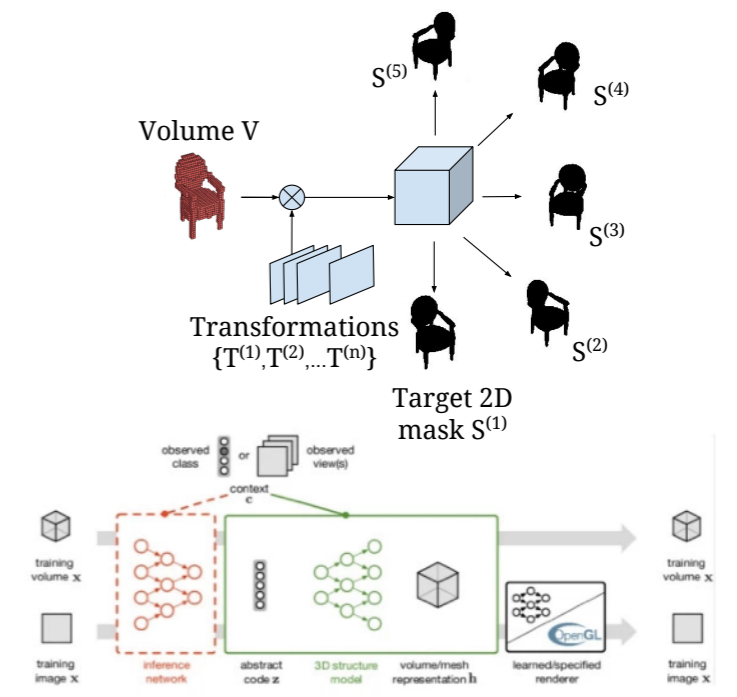
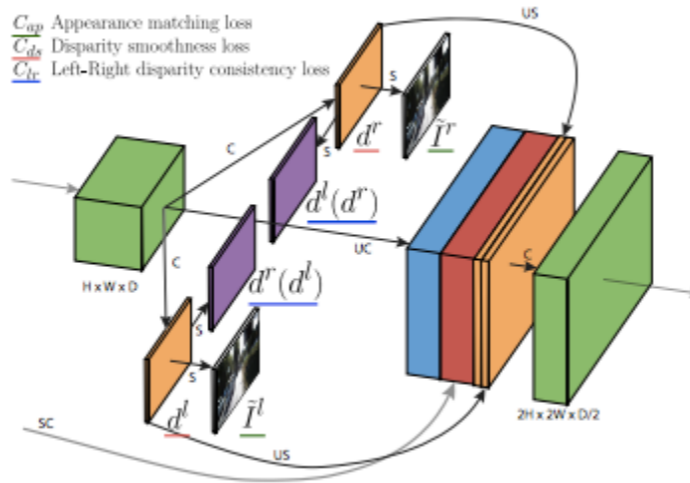
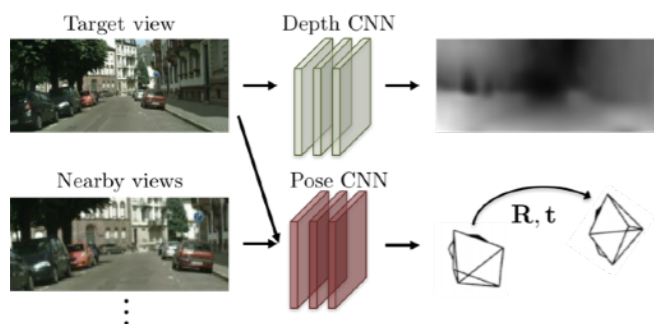
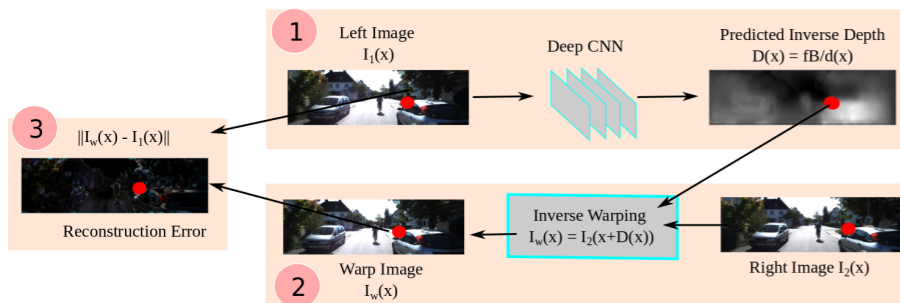
Garg et. al. ECCV 16

Godard et. al., Zhou et. al., CVPR 17

# 3D from Geometric Consistency



## Space Carving, Multi-view Stereo, Multi-view Reconstruction



Garg et. al. ECCV 16  
 Godard et. al., Zhou et. al., CVPR 17

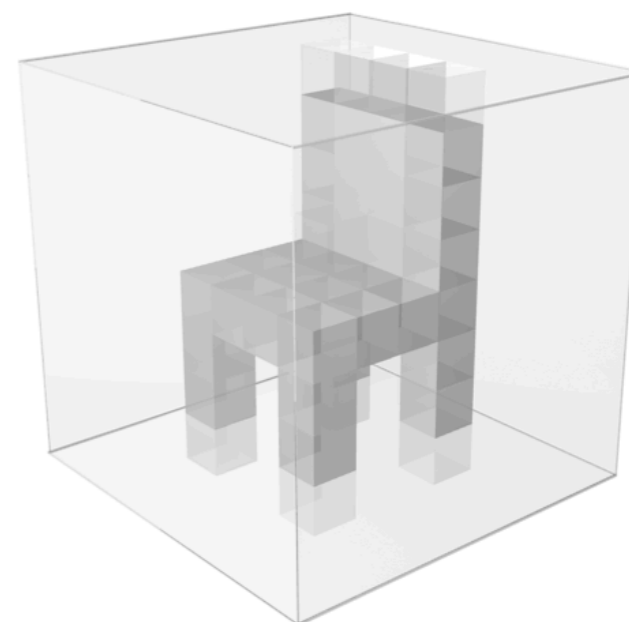
Yan et. al., Rezende et. al.  
 NIPS 16



# Learning via Geometric Consistency



Input Image



$$L(\text{Input Image}, \text{3D Voxel Chair})$$

Geometric

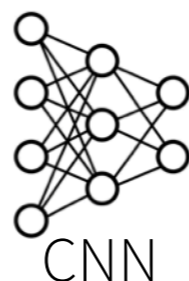
Consistency Loss

Observation **O**  
from camera **C**

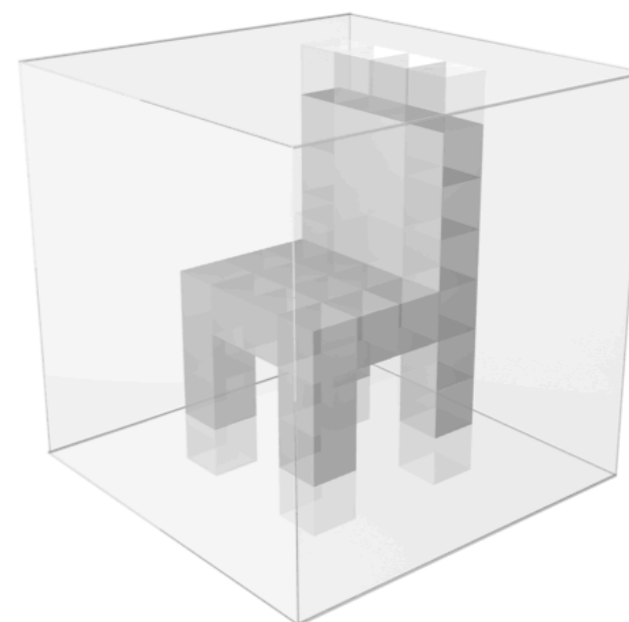
# Learning via Geometric Consistency



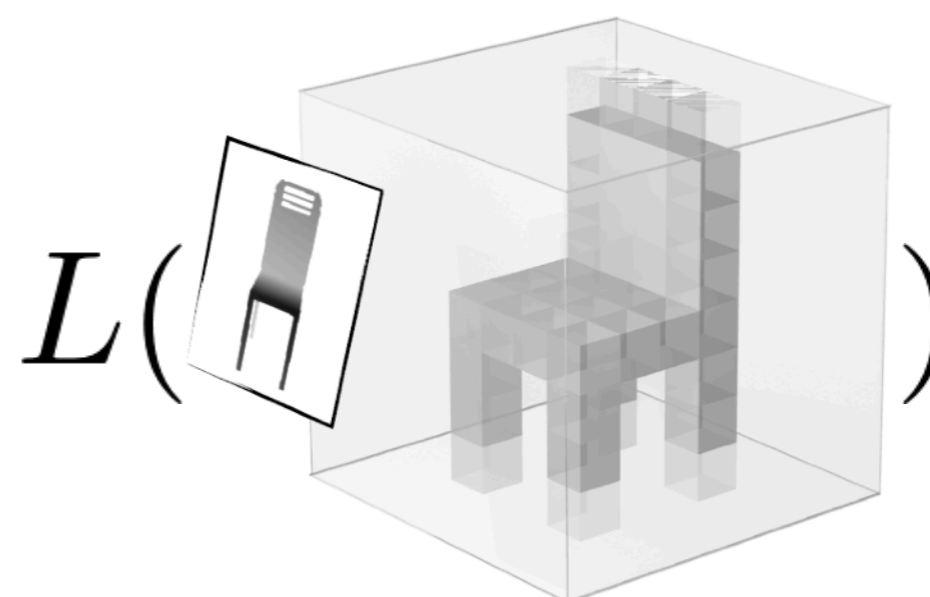
Input Image



CNN



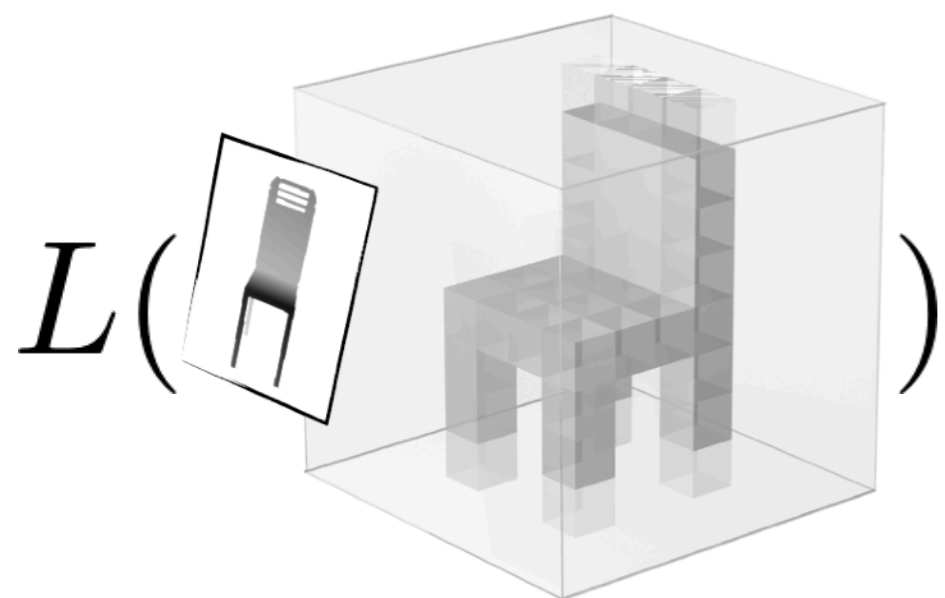
Observation **O**  
from camera **C**



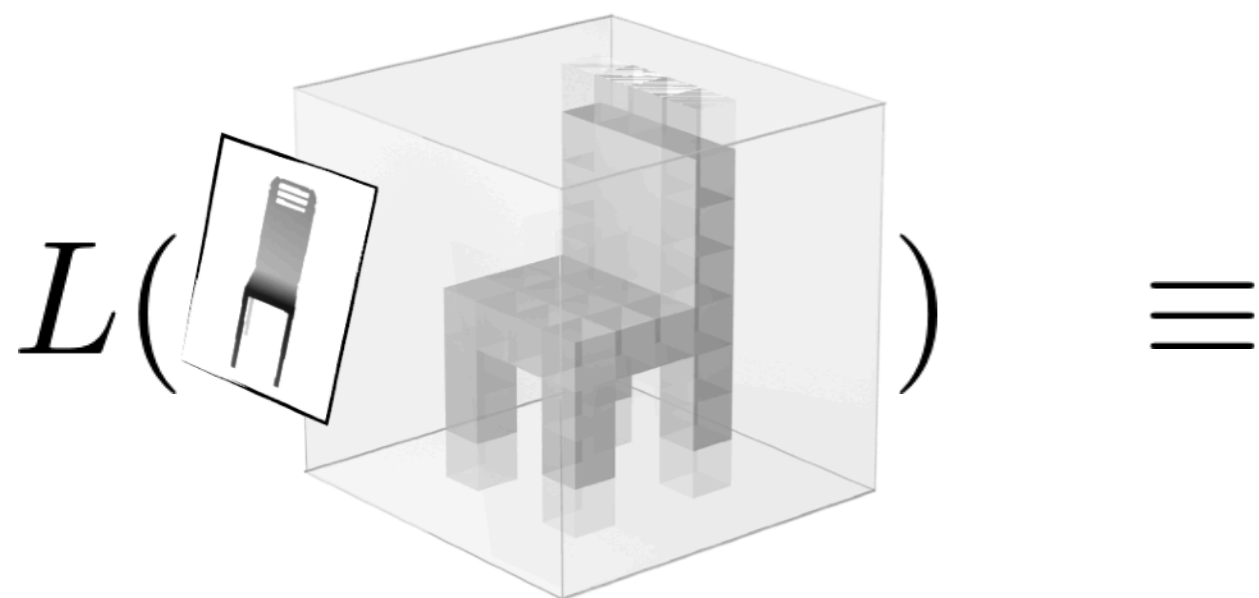
Geometric  
Consistency Loss



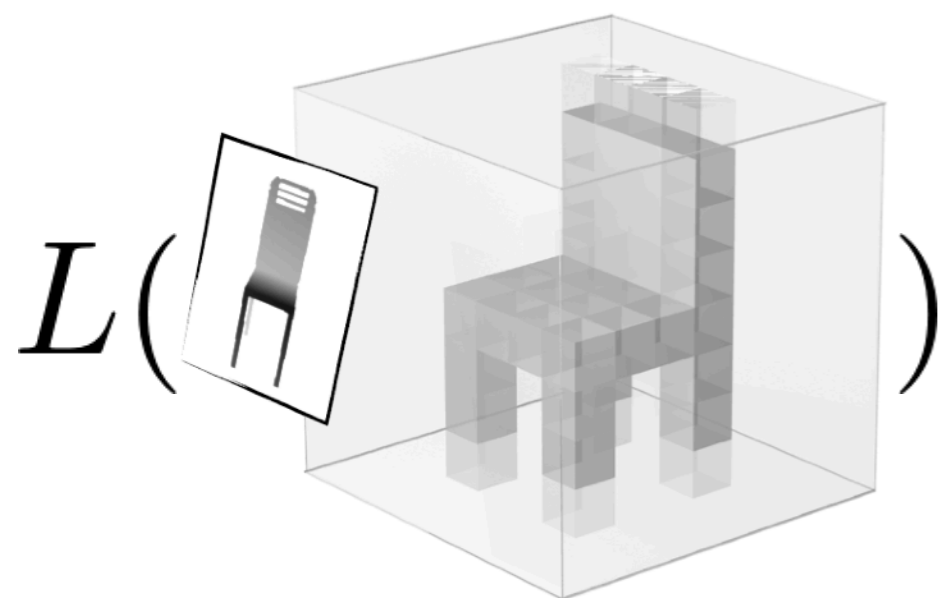
# View Consistency as Ray Consistency



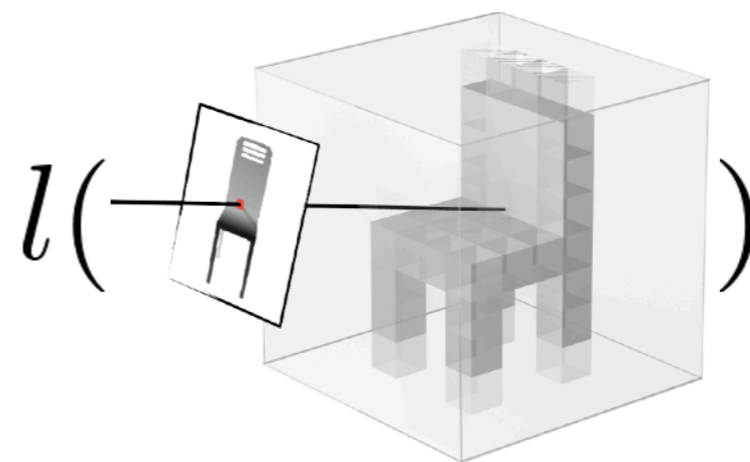
# View Consistency as Ray Consistency



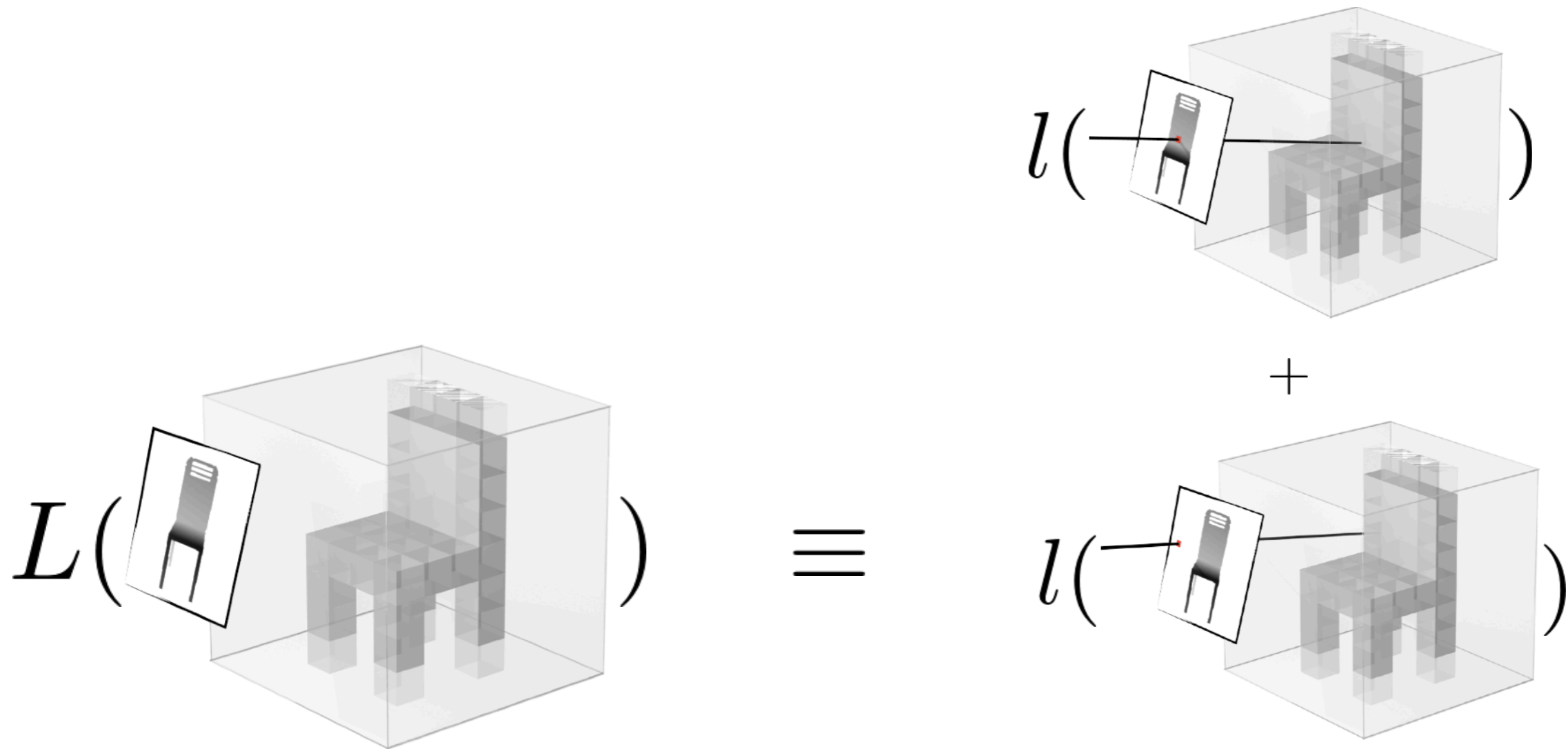
# View Consistency as Ray Consistency



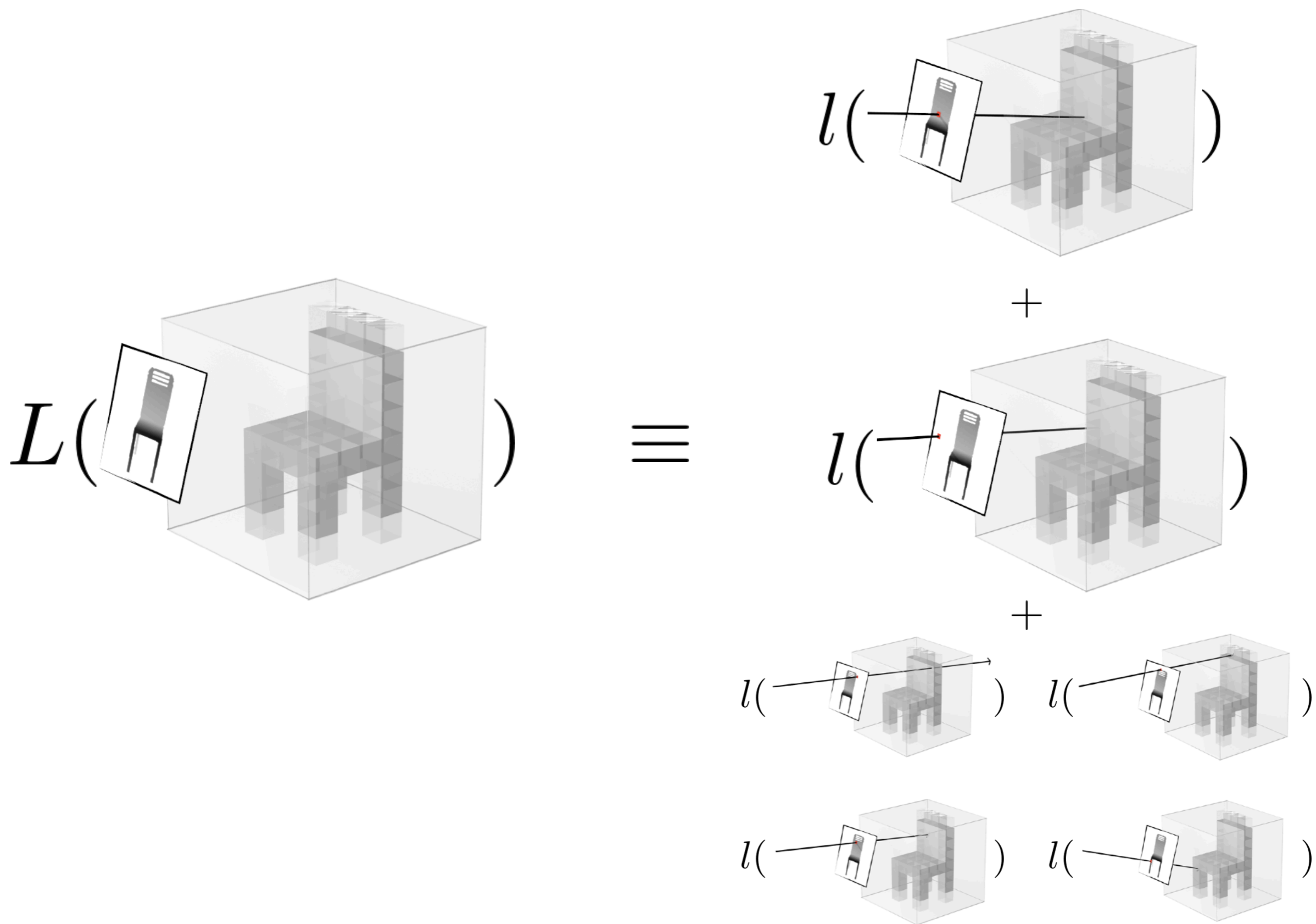
$\equiv$



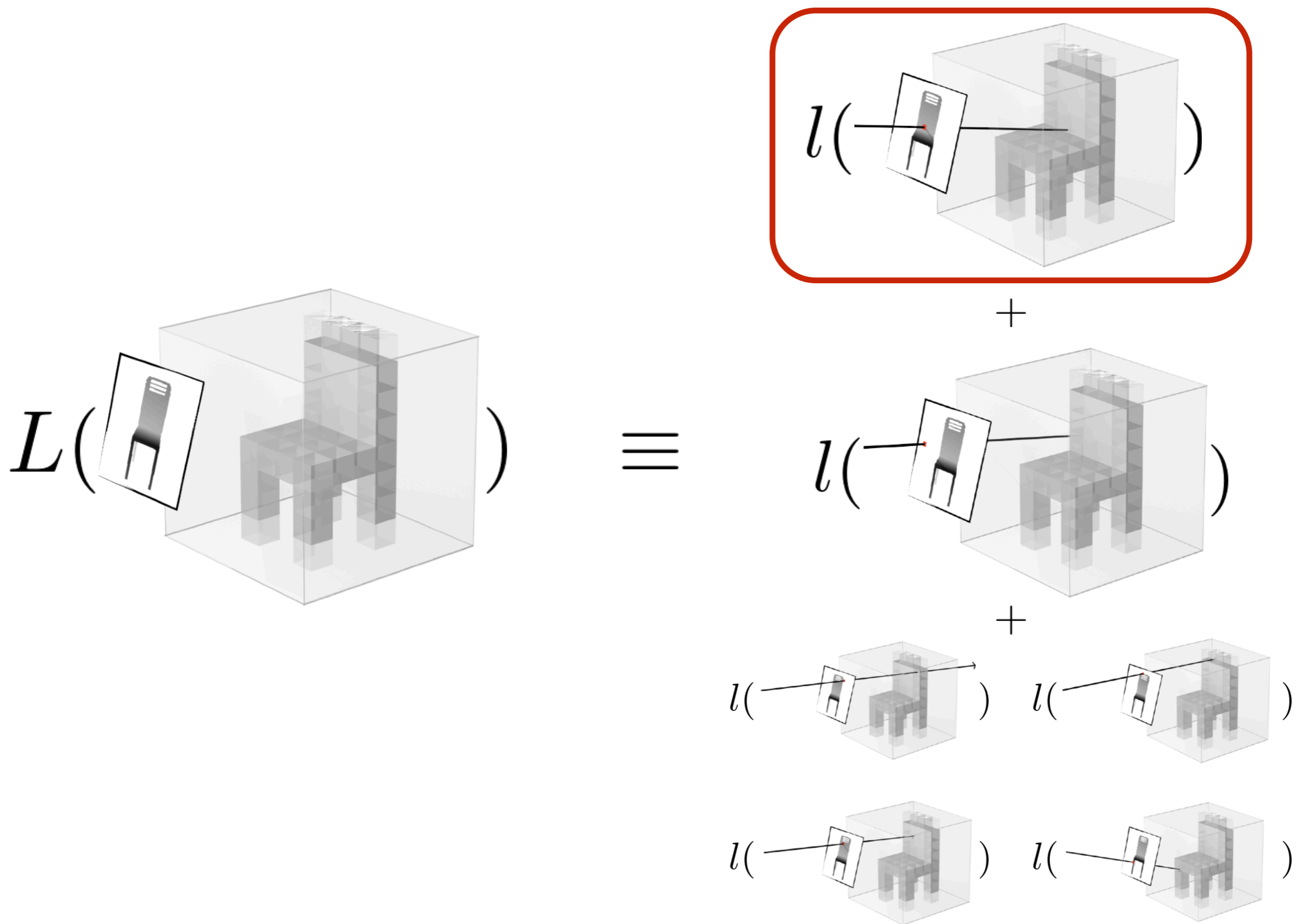
# View Consistency as Ray Consistency



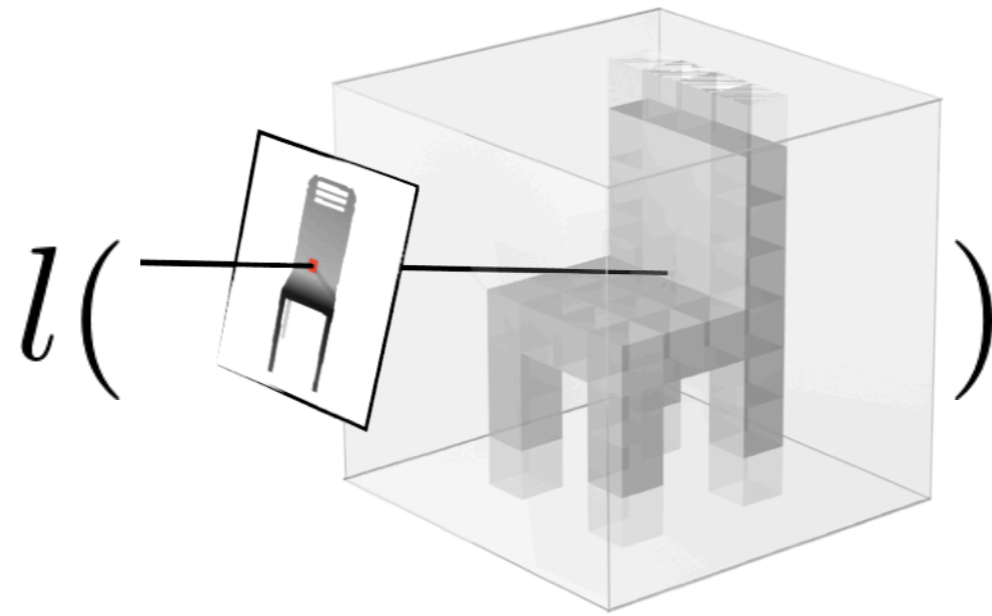
# View Consistency as Ray Consistency



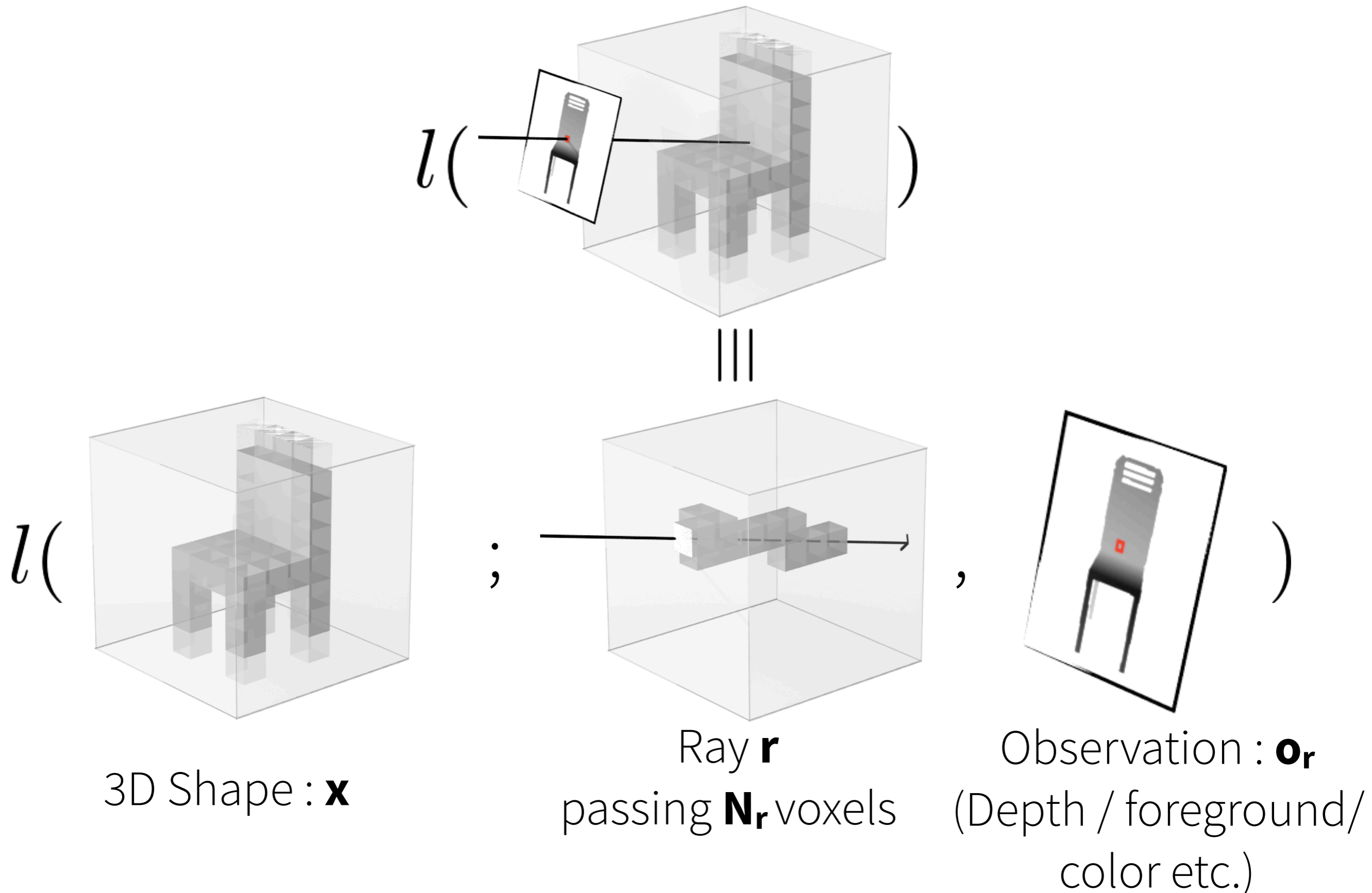
# View Consistency as Ray Consistency



# Differentiable Ray Consistency

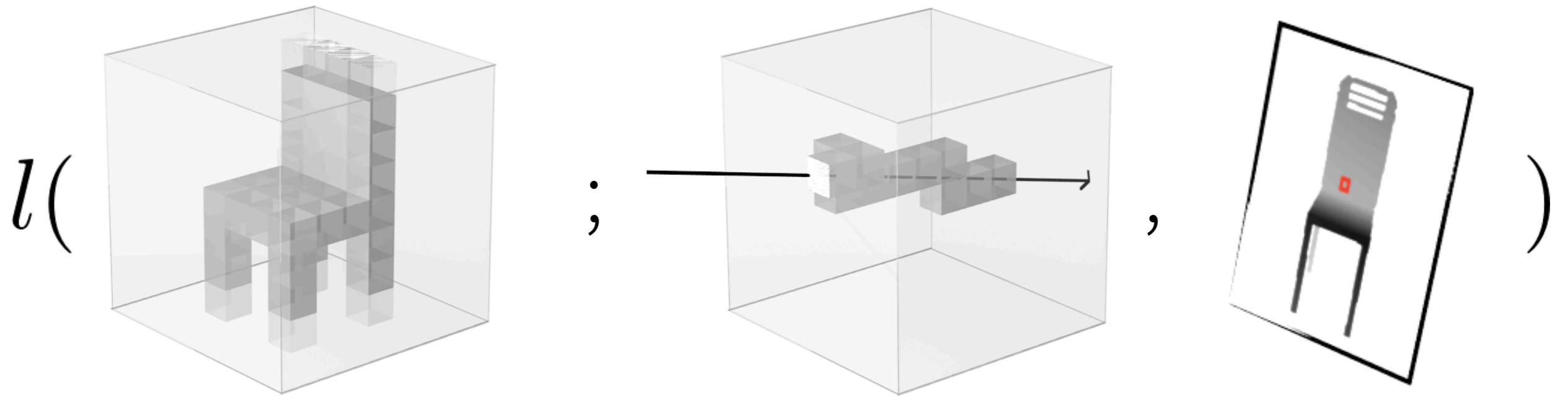


# Differentiable Ray Consistency

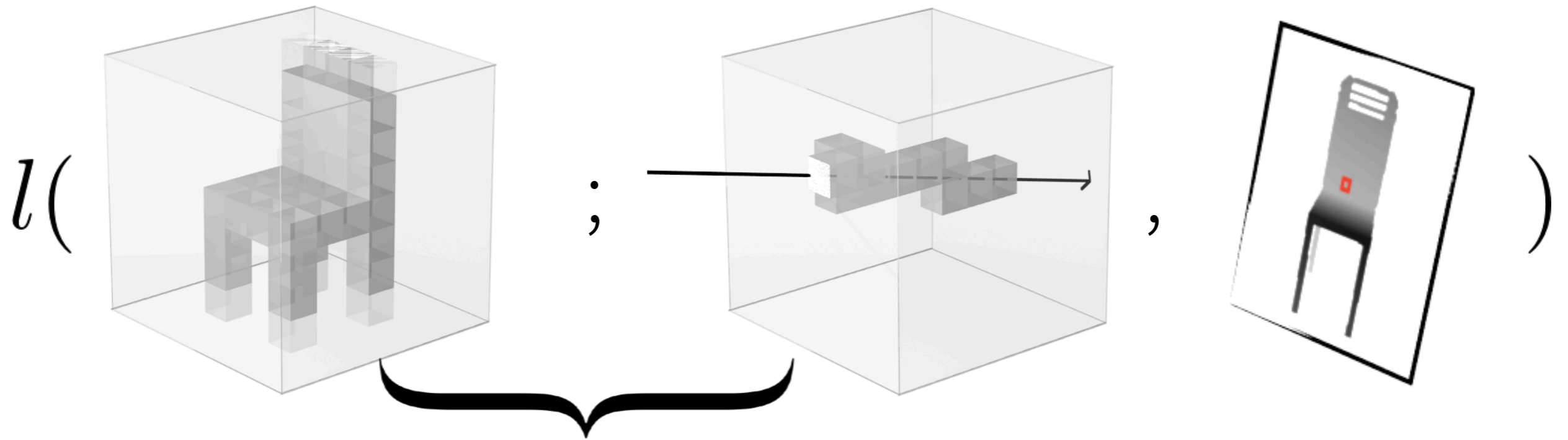




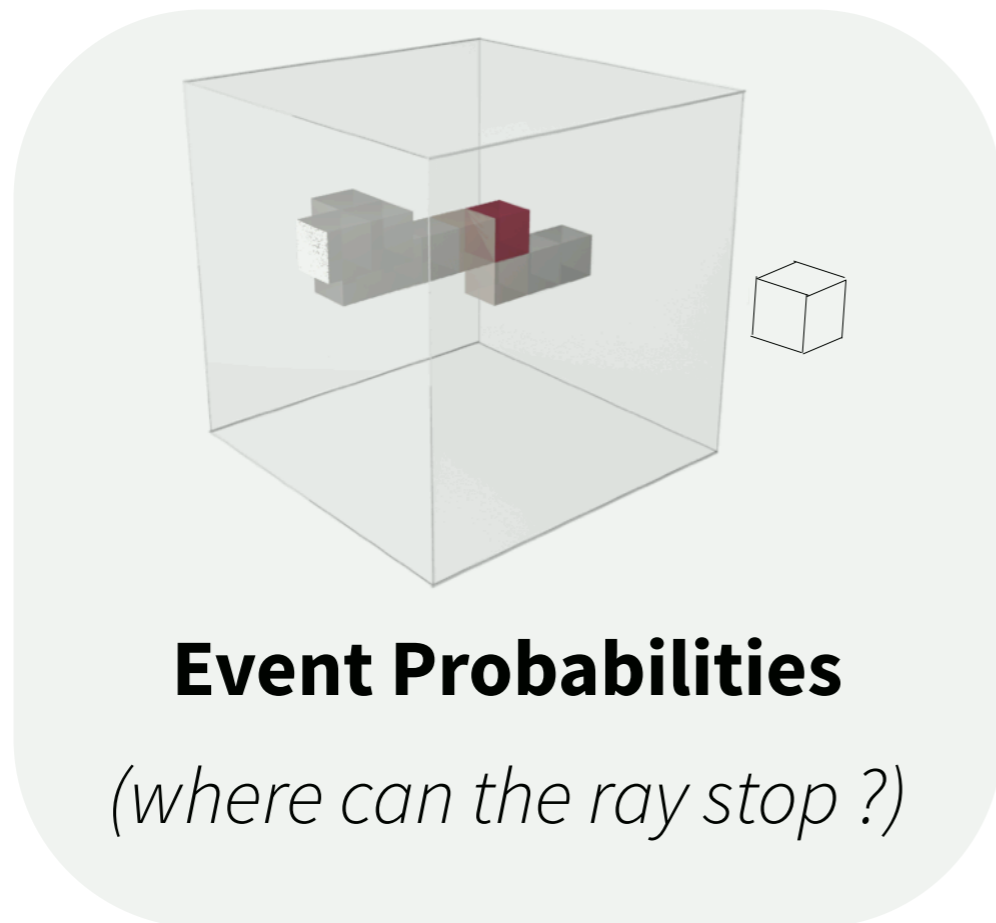
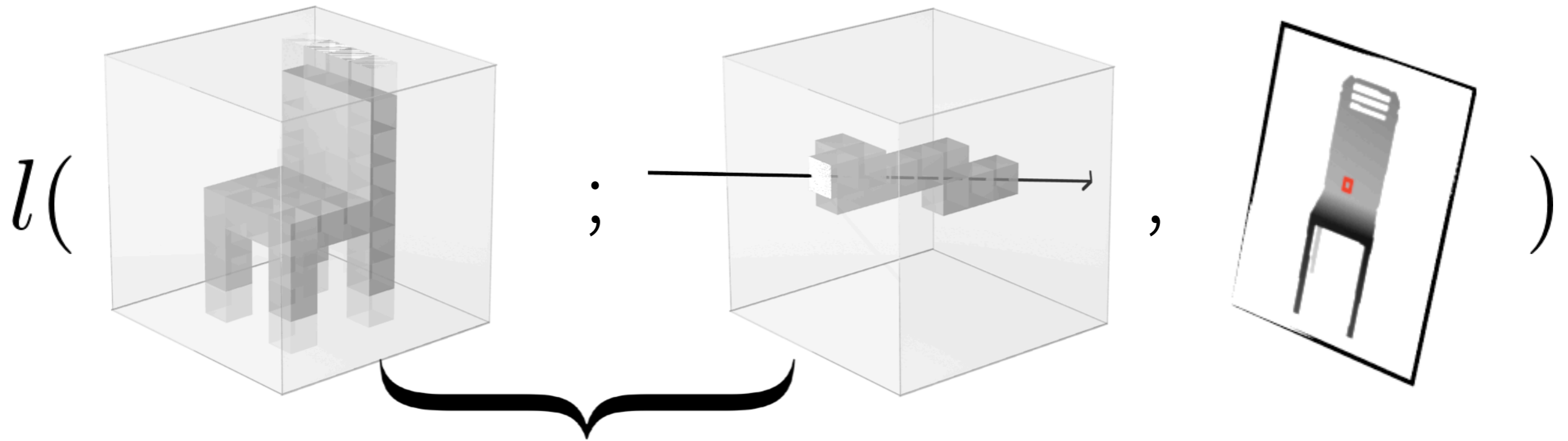
# Differentiable Ray Consistency



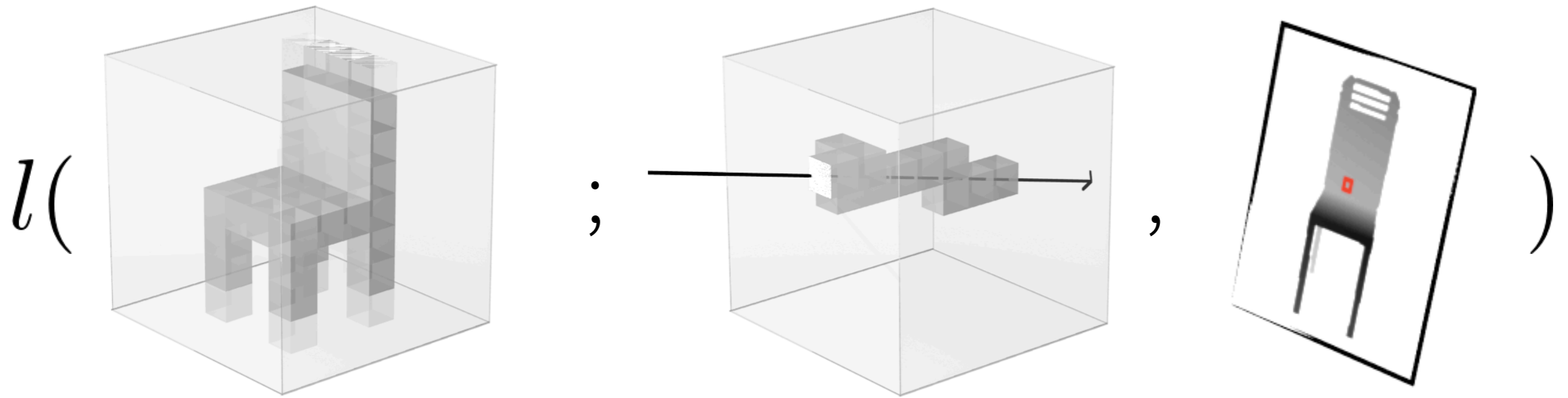
# Differentiable Ray Consistency



# Differentiable Ray Consistency



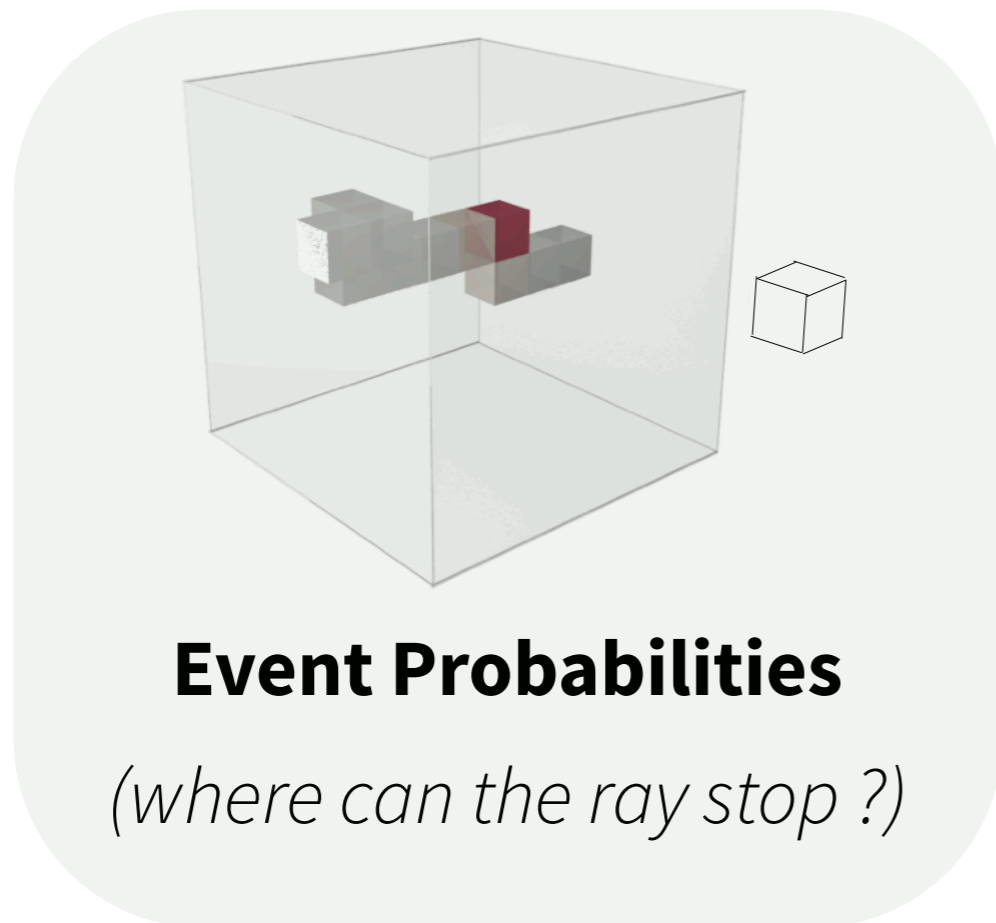
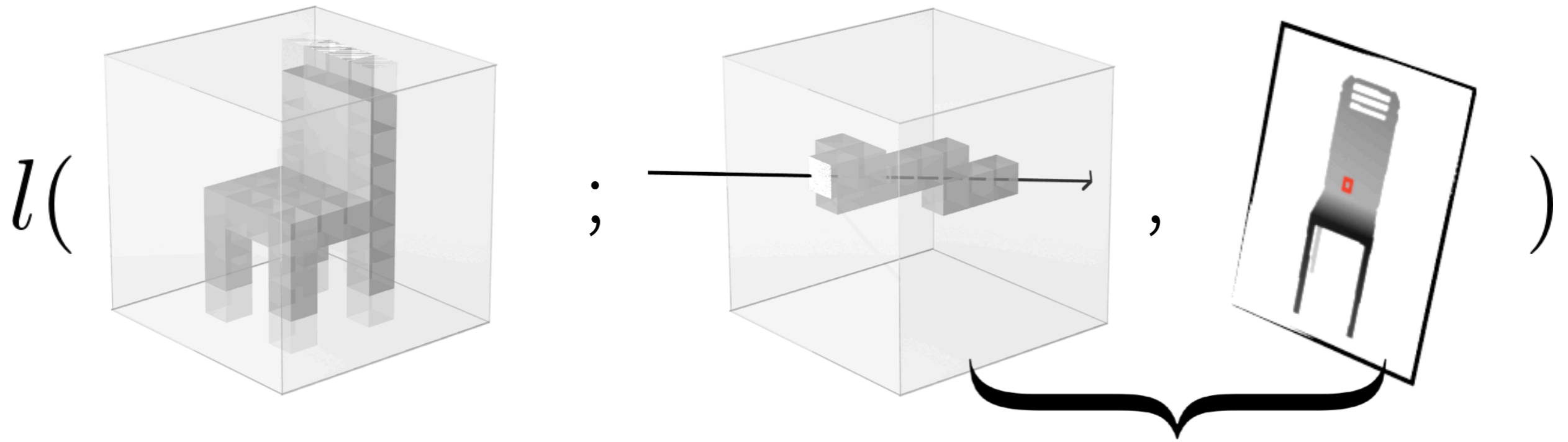
# Differentiable Ray Consistency



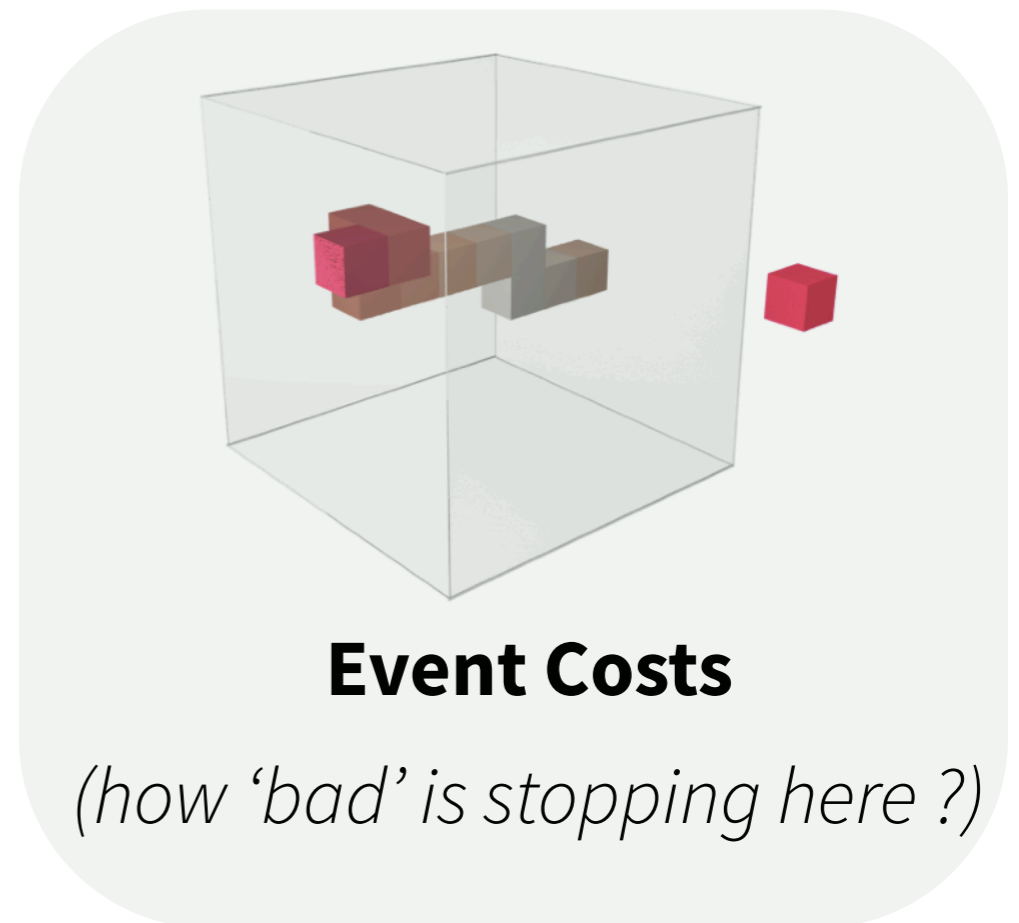
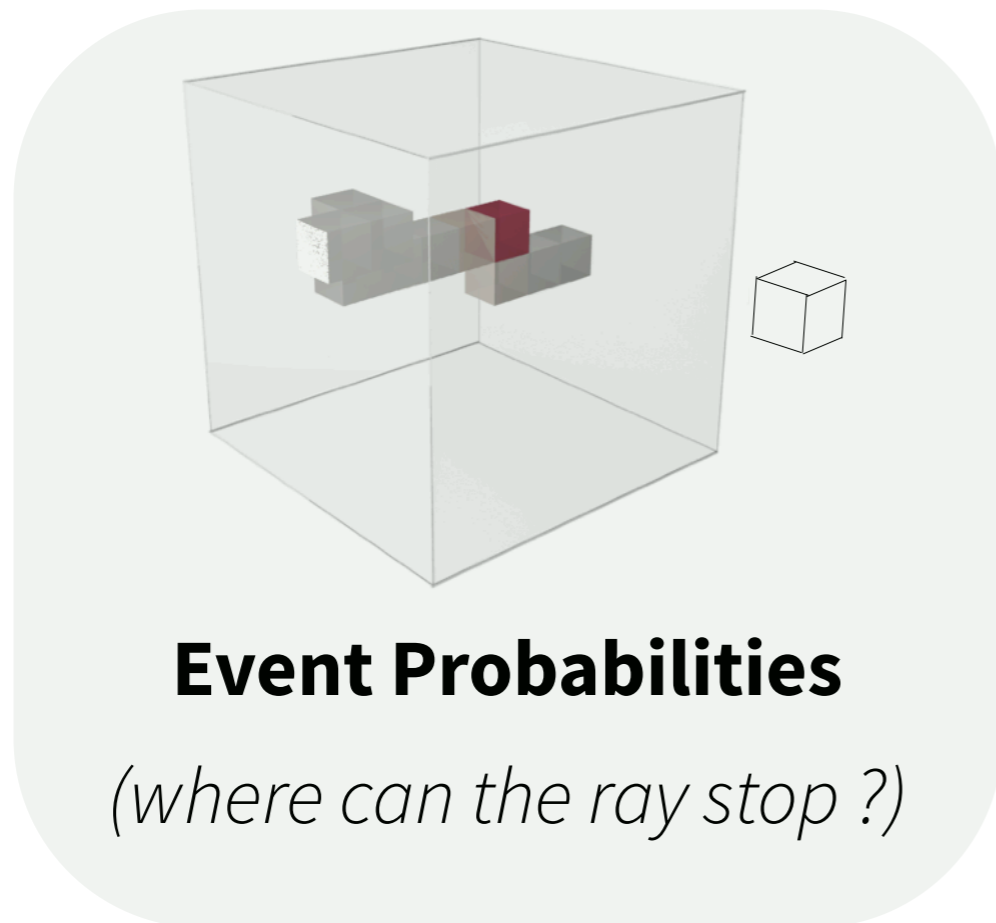
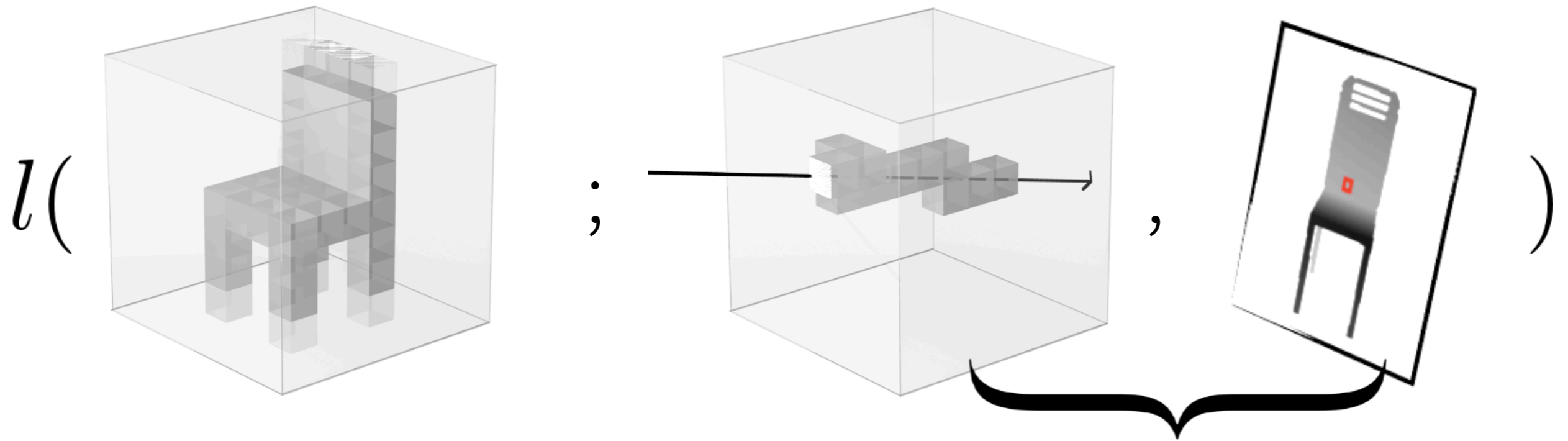
The diagram shows a 3D scene with a horizontal bar of blocks. One block is highlighted in red. A ray is shown passing through the blocks. To the right of the scene is a small white cube. Below the scene, the text reads:

**Event Probabilities**  
*(where can the ray stop ?)*

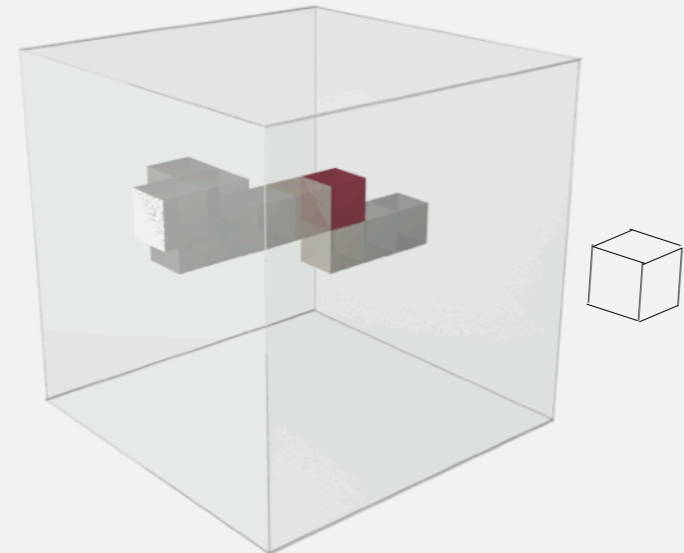
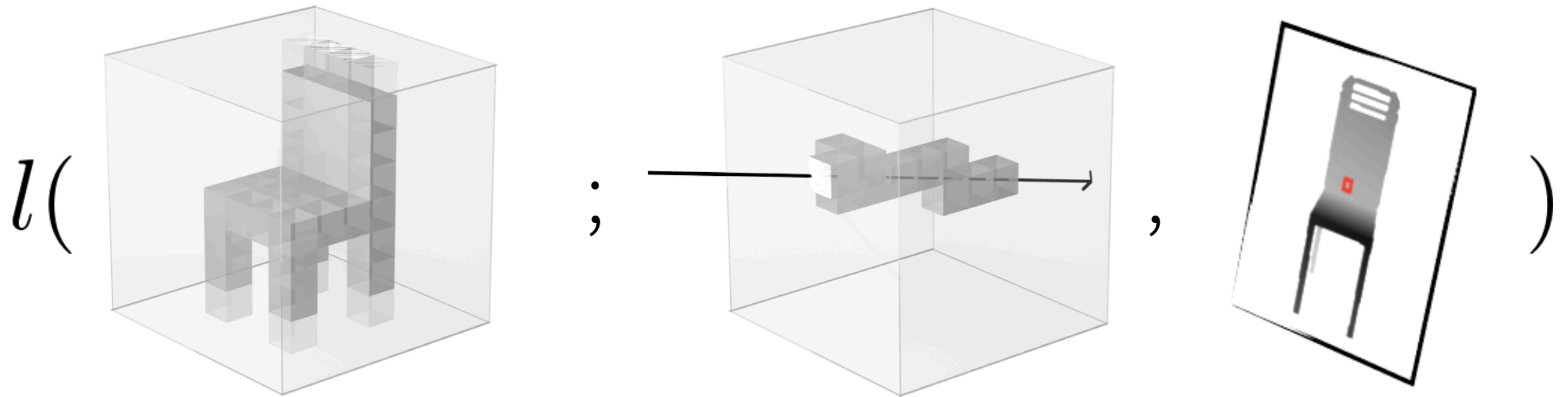
# Differentiable Ray Consistency



# Differentiable Ray Consistency



# Differentiable Ray Consistency

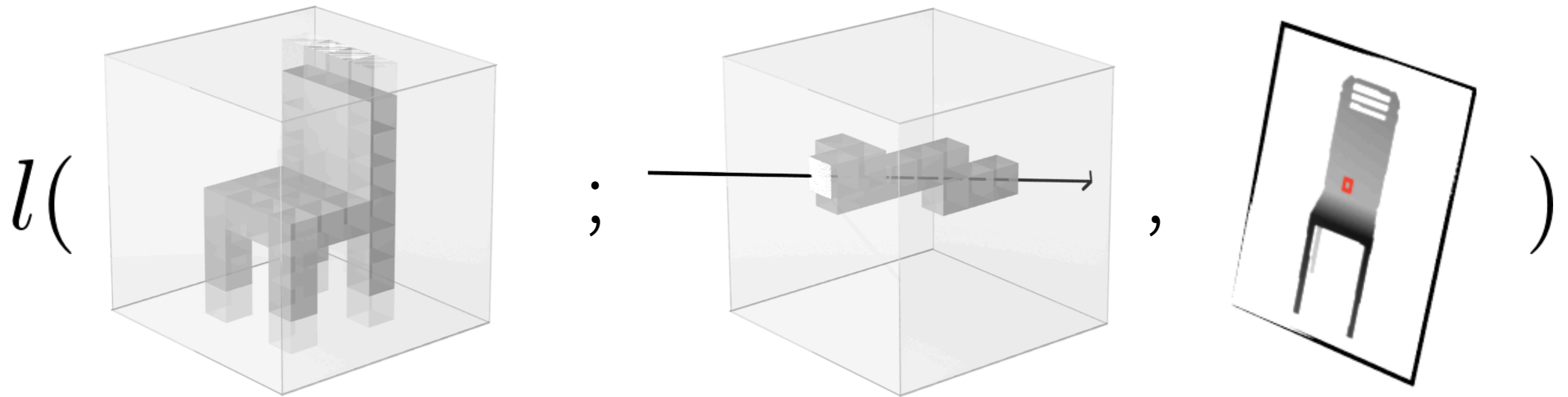


**Event Probabilities**  
*(where can the ray stop ?)*

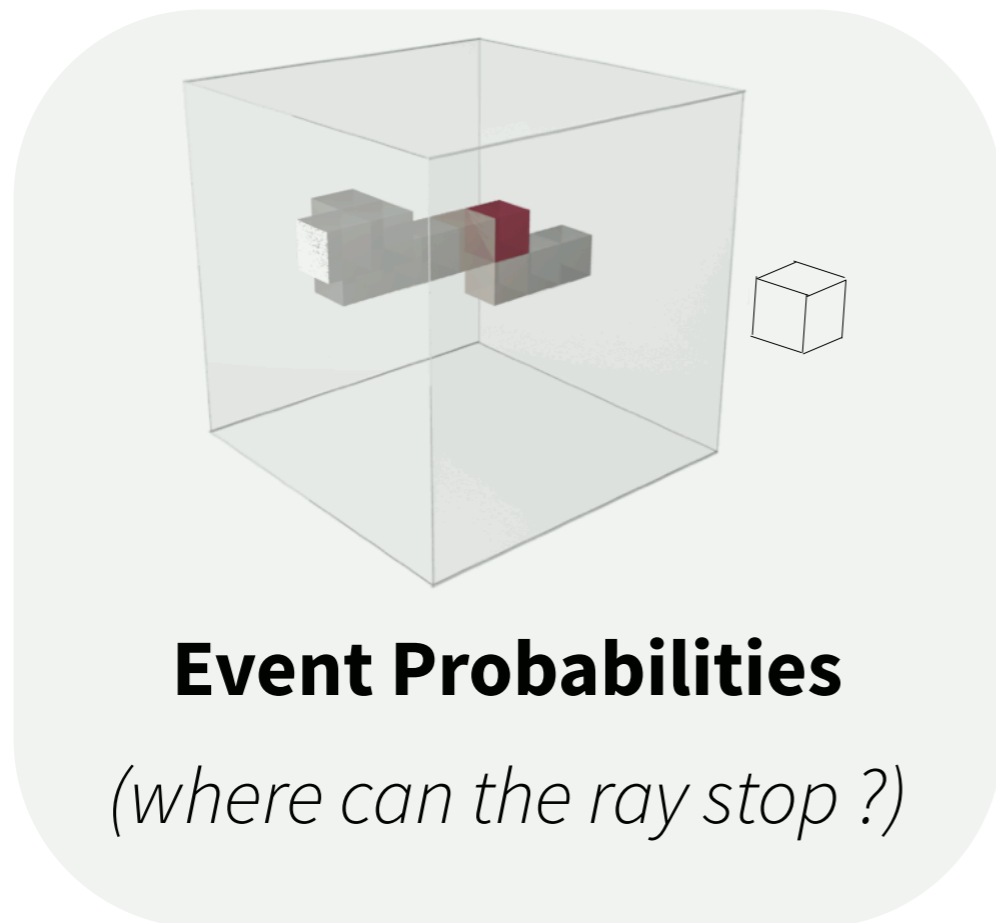


**Event Costs**  
*(how 'bad' is stopping here ?)*

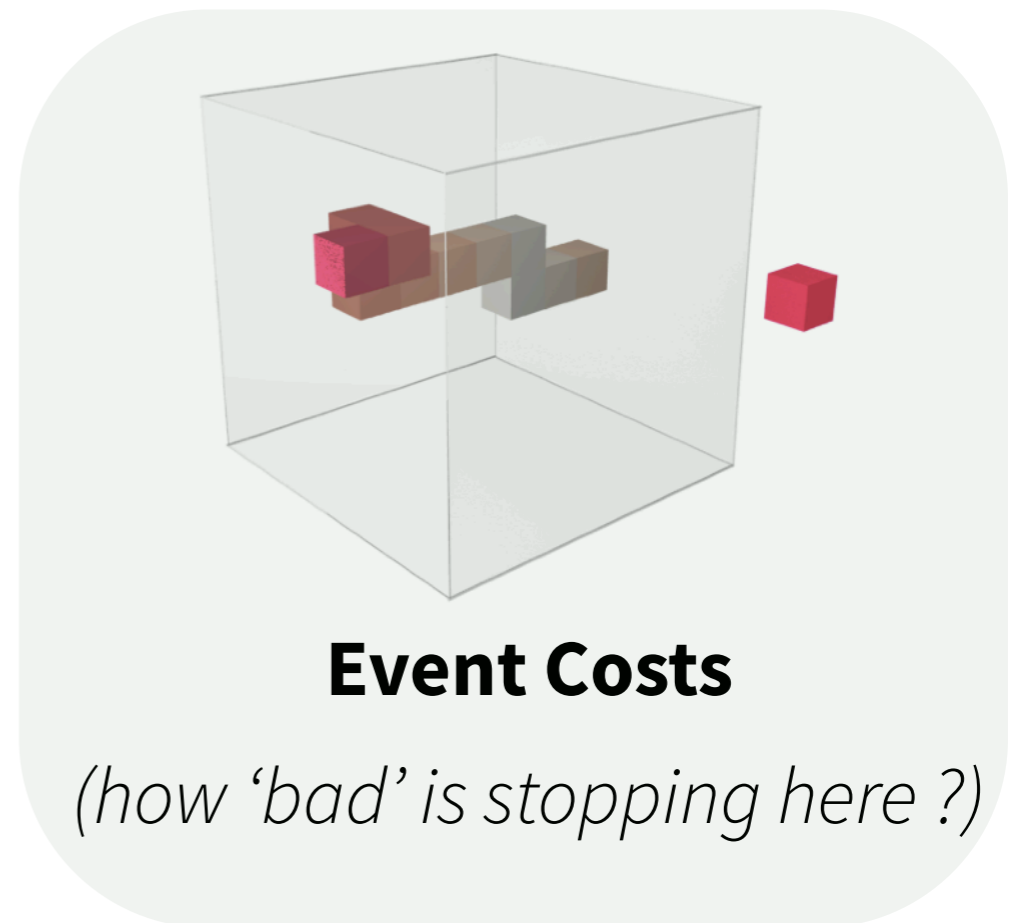
# Differentiable Ray Consistency



$\Sigma($

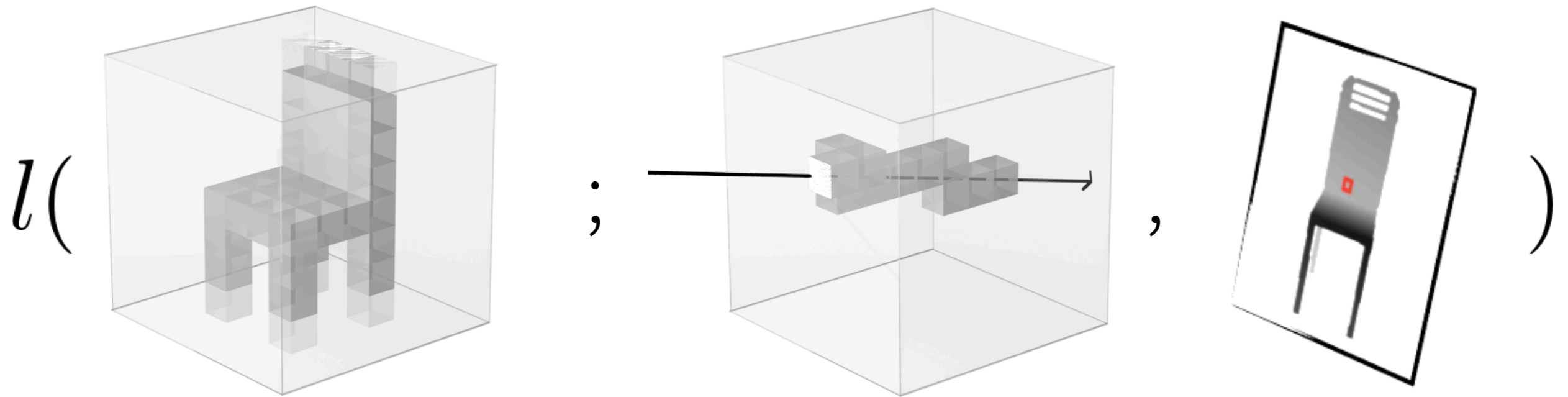


$\odot$

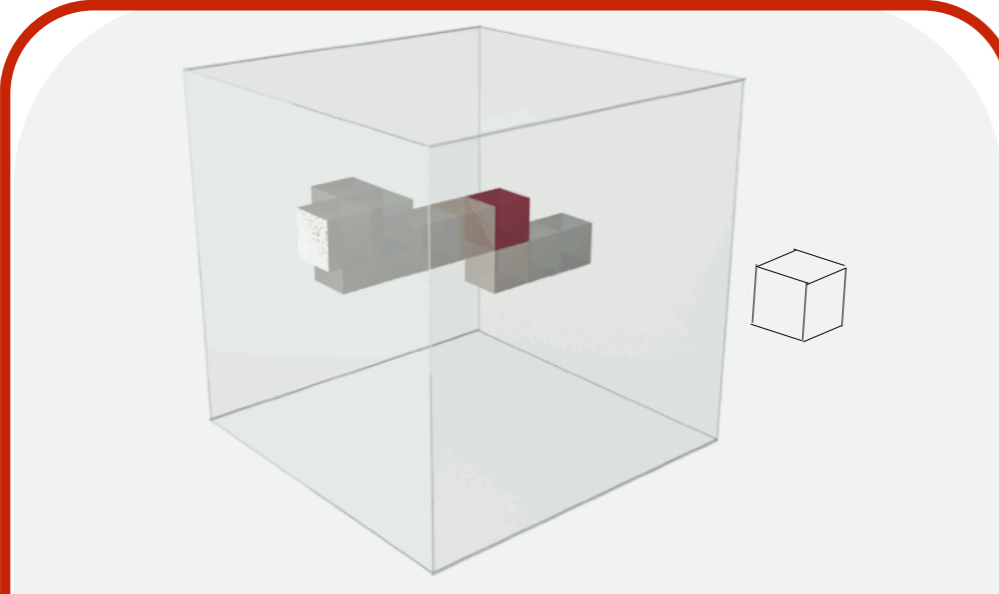




# Differentiable Ray Consistency



$\Sigma($



A 3D scene with a ray and a semi-transparent cube at the intersection point. A small white cube is shown to the right of the scene. The entire scene is enclosed in a rounded rectangle with a red border.

**Event Probabilities**  
*(where can the ray stop ?)*

$\odot$

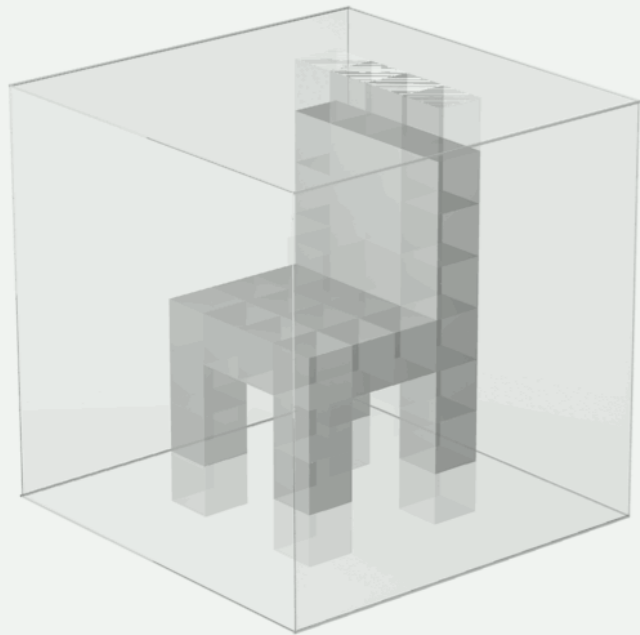


A 3D scene with a ray and a semi-transparent cube at the intersection point. A small red cube is shown to the right of the scene.

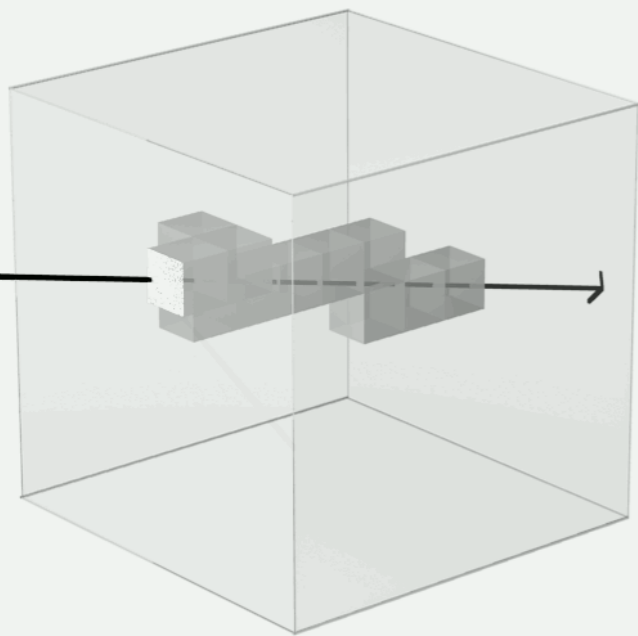
**Event Costs**  
*(how 'bad' is stopping here ?)*

$)$

# Probabilistic Ray Tracing



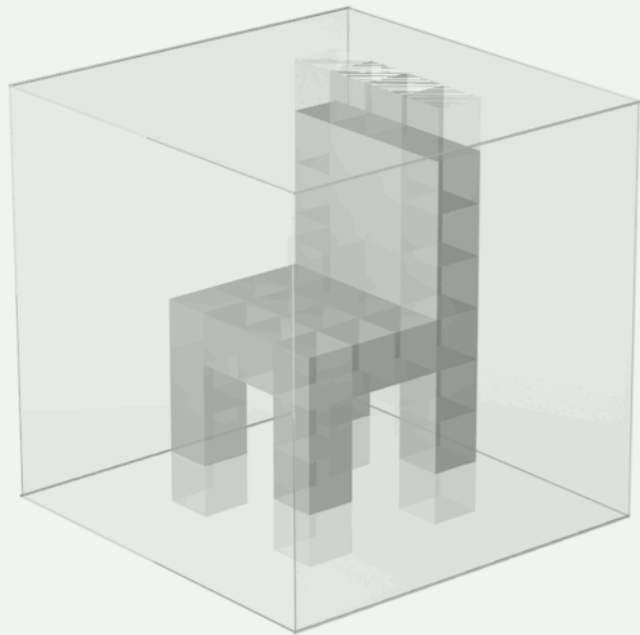
**x**



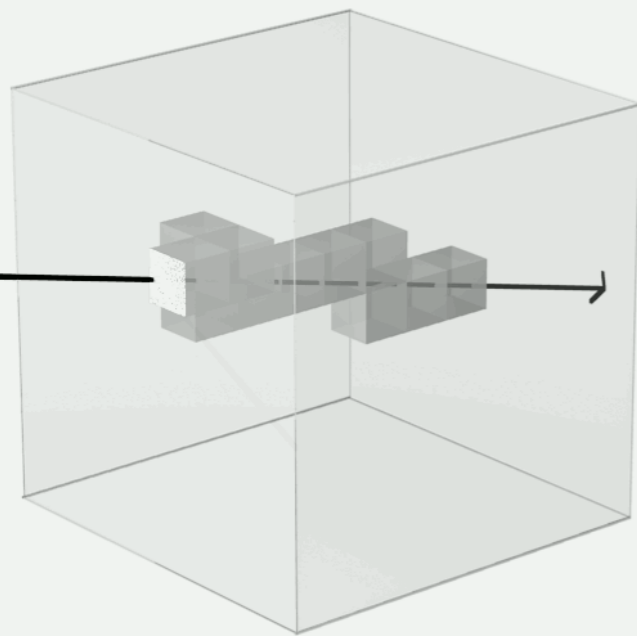
**r**

Possible Events

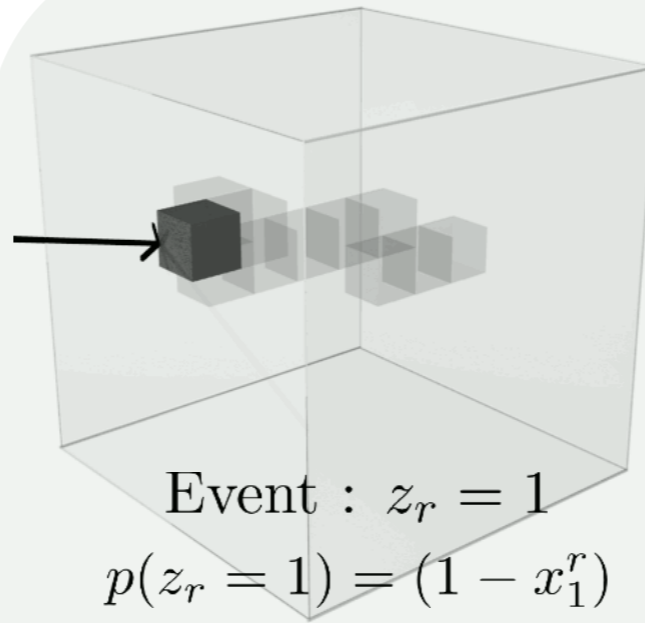
# Probabilistic Ray Tracing



**x**



**r**

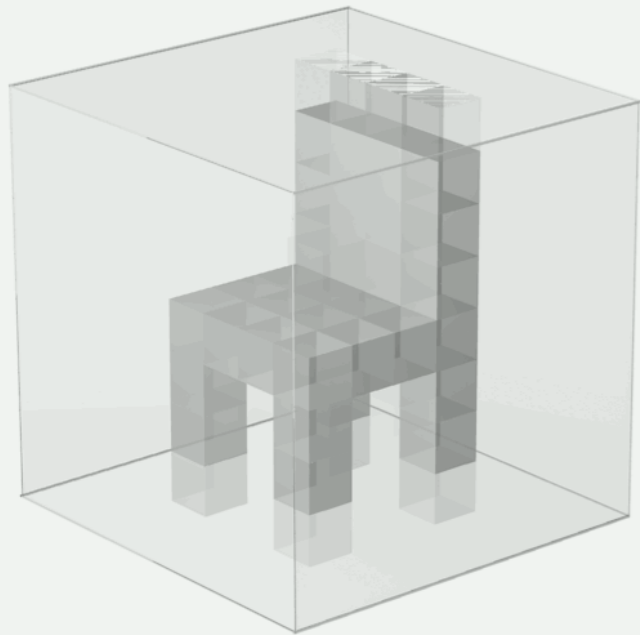


Event :  $z_r = 1$

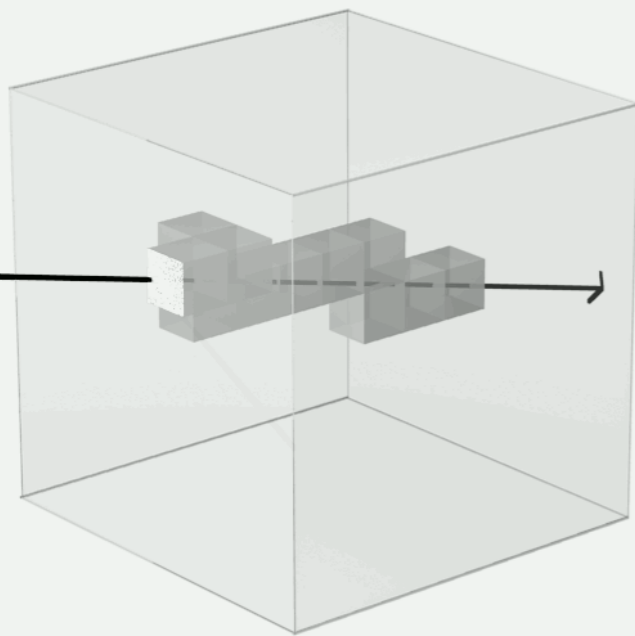
$$p(z_r = 1) = (1 - x_1^r)$$

Possible Events

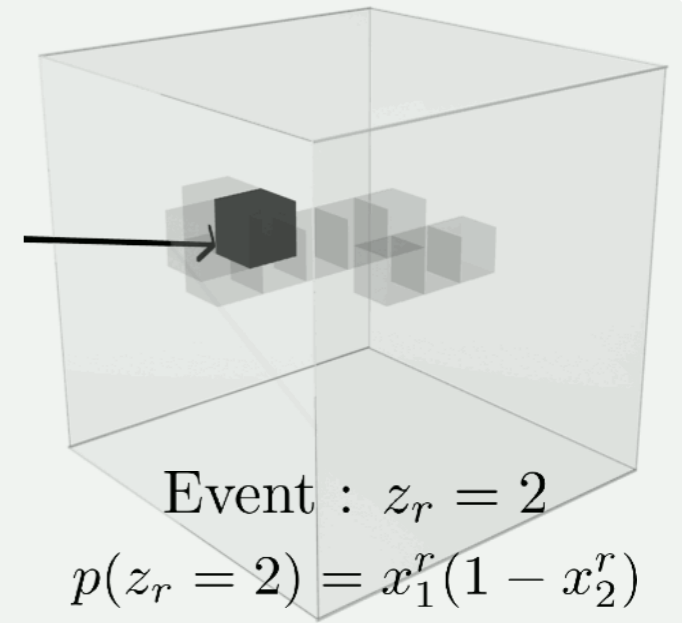
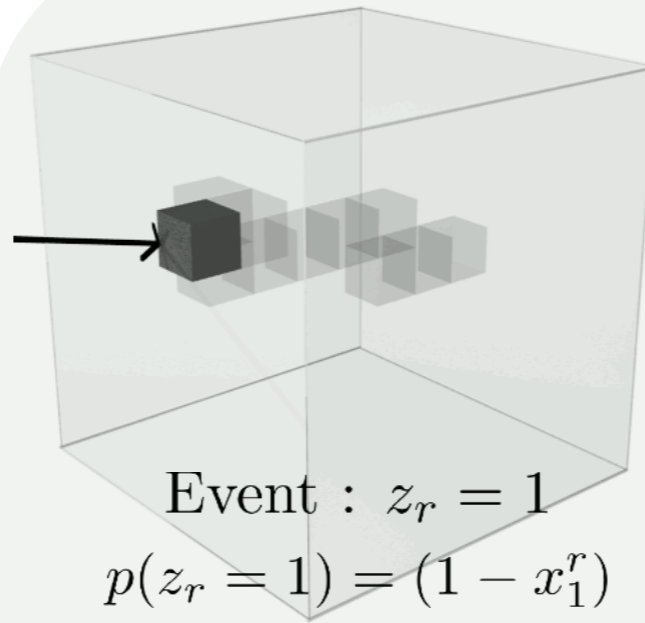
# Probabilistic Ray Tracing



**x**

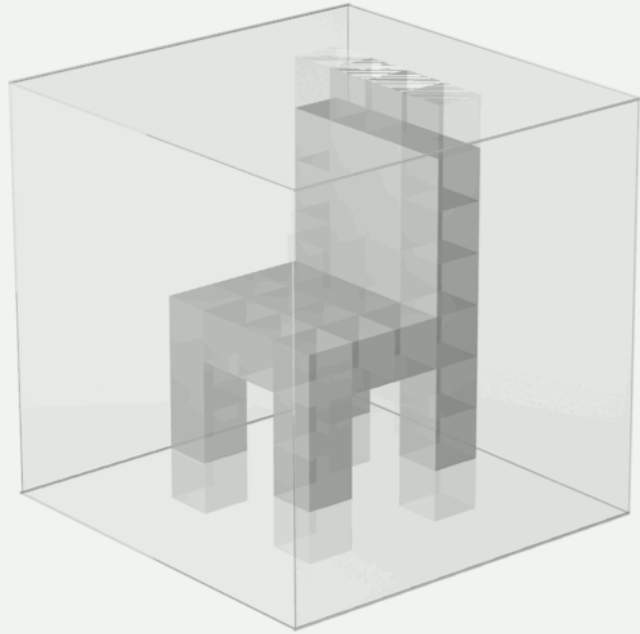


**r**

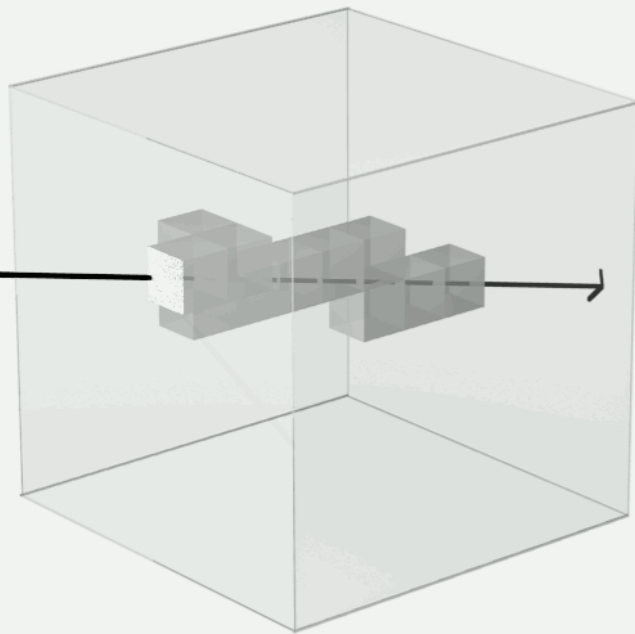


Possible Events

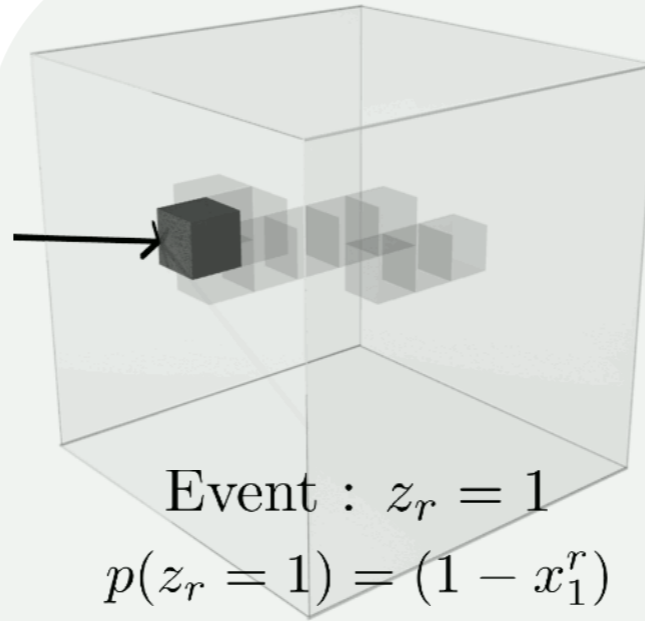
# Probabilistic Ray Tracing



**x**

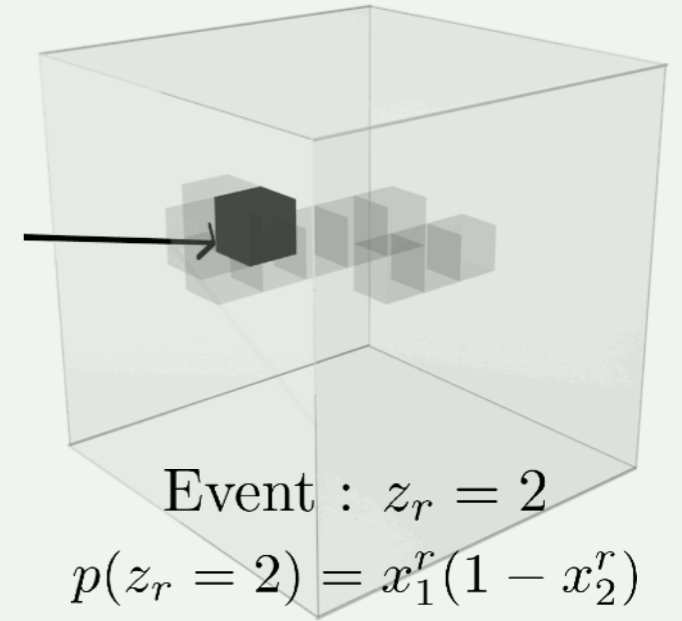


**r**



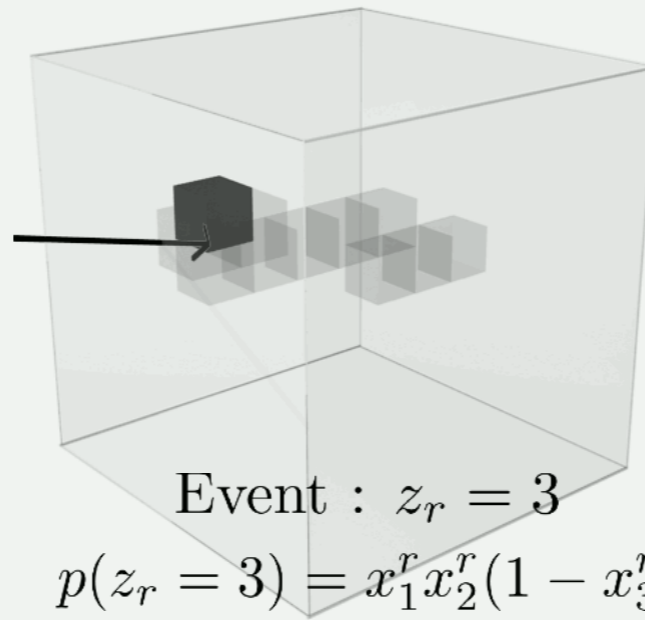
Event :  $z_r = 1$

$$p(z_r = 1) = (1 - x_1^r)$$



Event :  $z_r = 2$

$$p(z_r = 2) = x_1^r(1 - x_2^r)$$

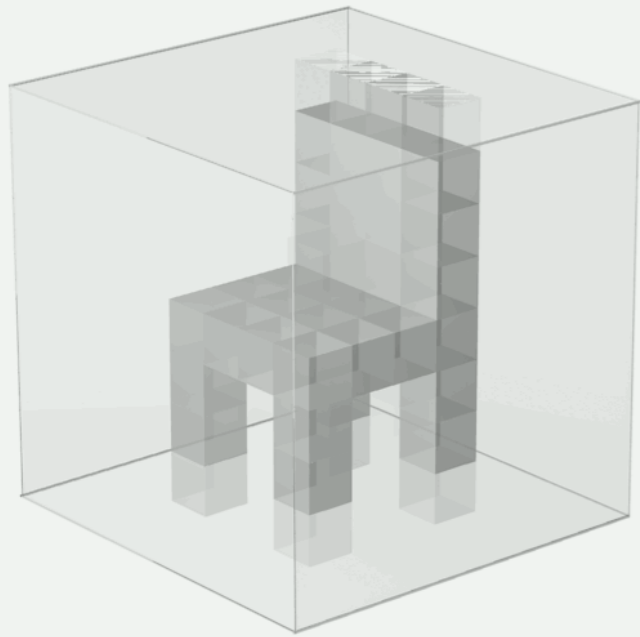


Event :  $z_r = 3$

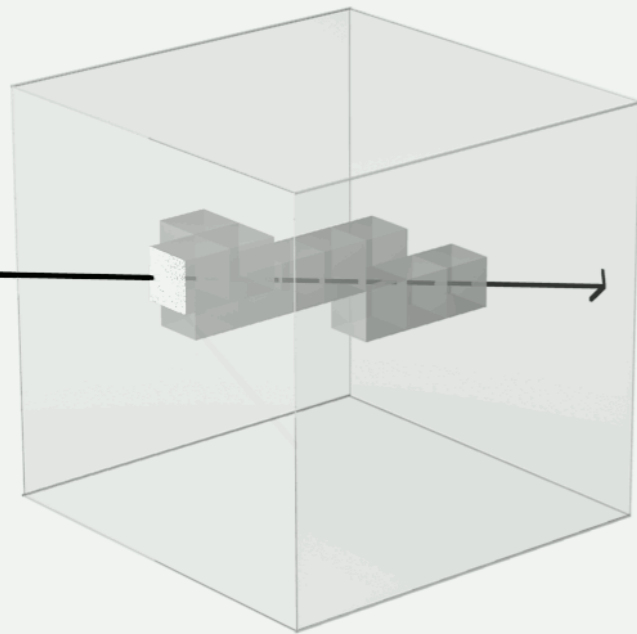
$$p(z_r = 3) = x_1^r x_2^r (1 - x_3^r)$$

Possible Events

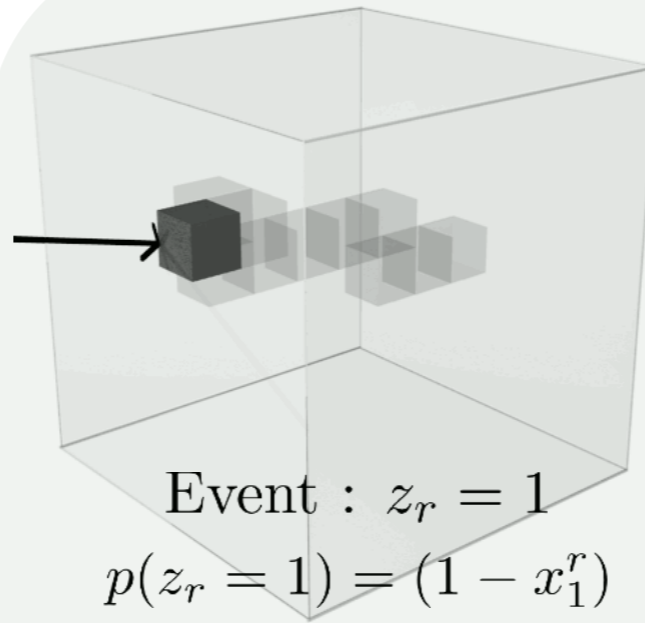
# Probabilistic Ray Tracing



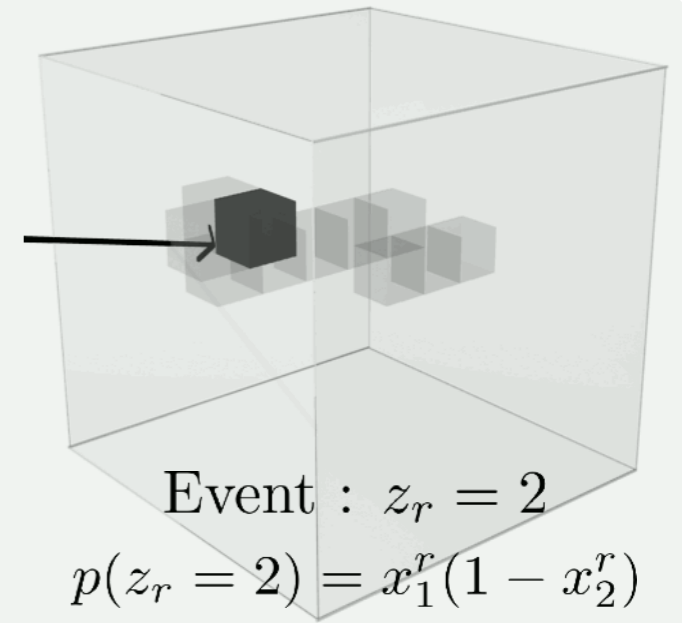
**x**



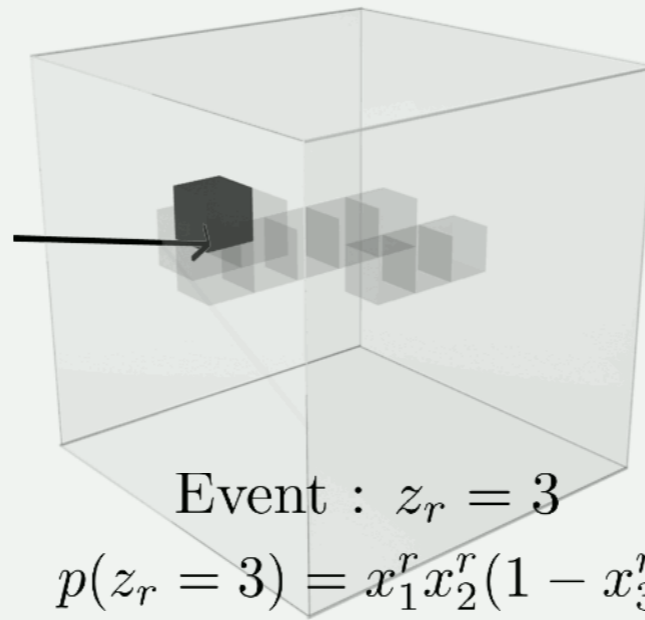
**r**



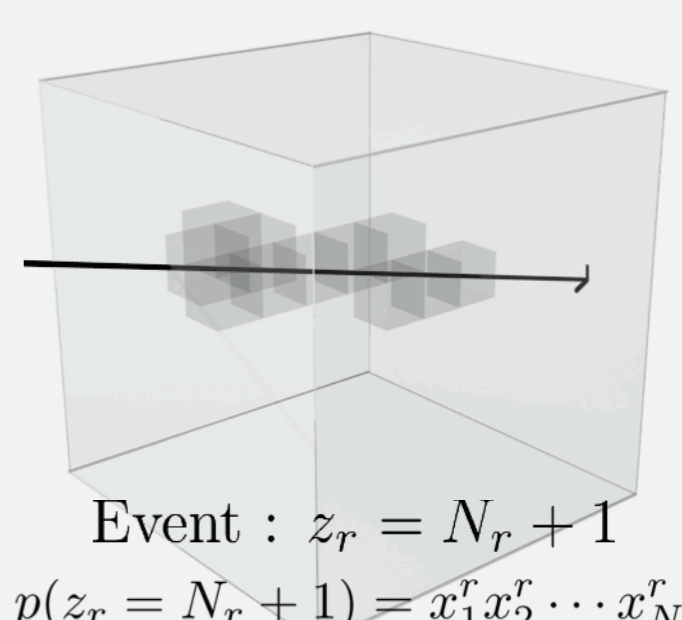
Event :  $z_r = 1$   
 $p(z_r = 1) = (1 - x_1^r)$



Event :  $z_r = 2$   
 $p(z_r = 2) = x_1^r(1 - x_2^r)$



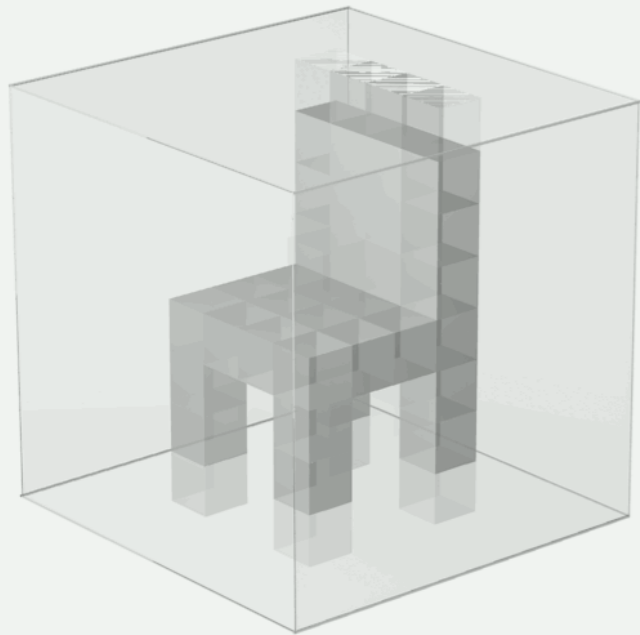
Event :  $z_r = 3$   
 $p(z_r = 3) = x_1^r x_2^r (1 - x_3^r)$



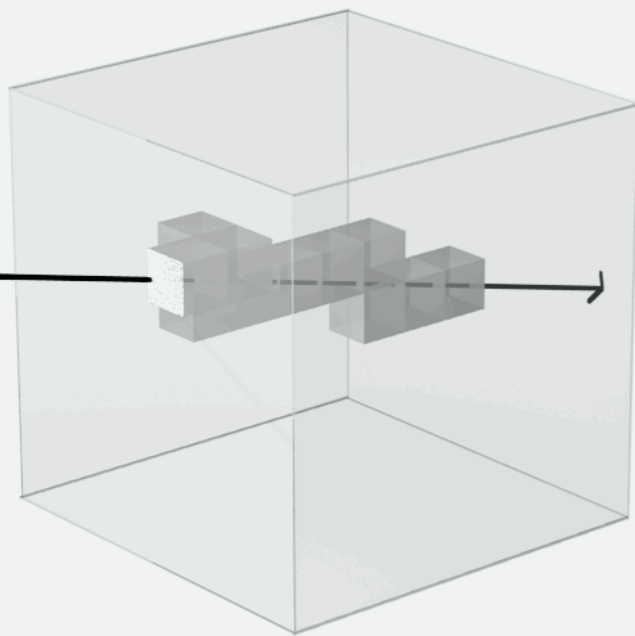
Event :  $z_r = N_r + 1$   
 $p(z_r = N_r + 1) = x_1^r x_2^r \cdots x_{N_r}^r$

Possible Events

# Probabilistic Ray Tracing



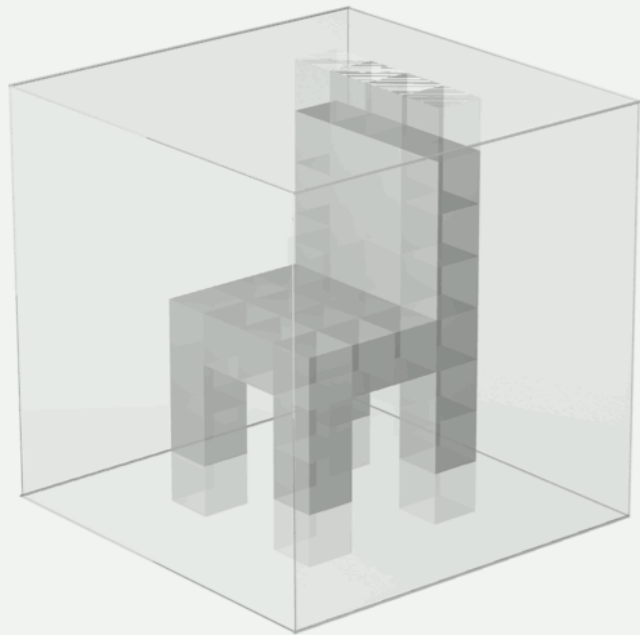
**x**



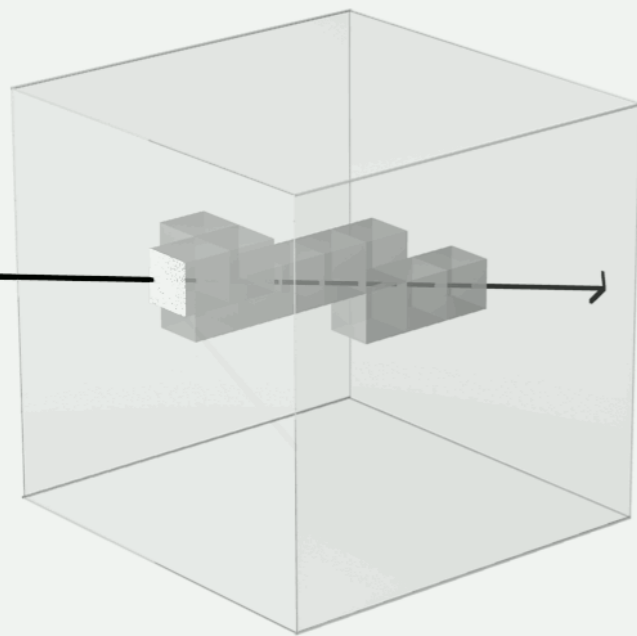
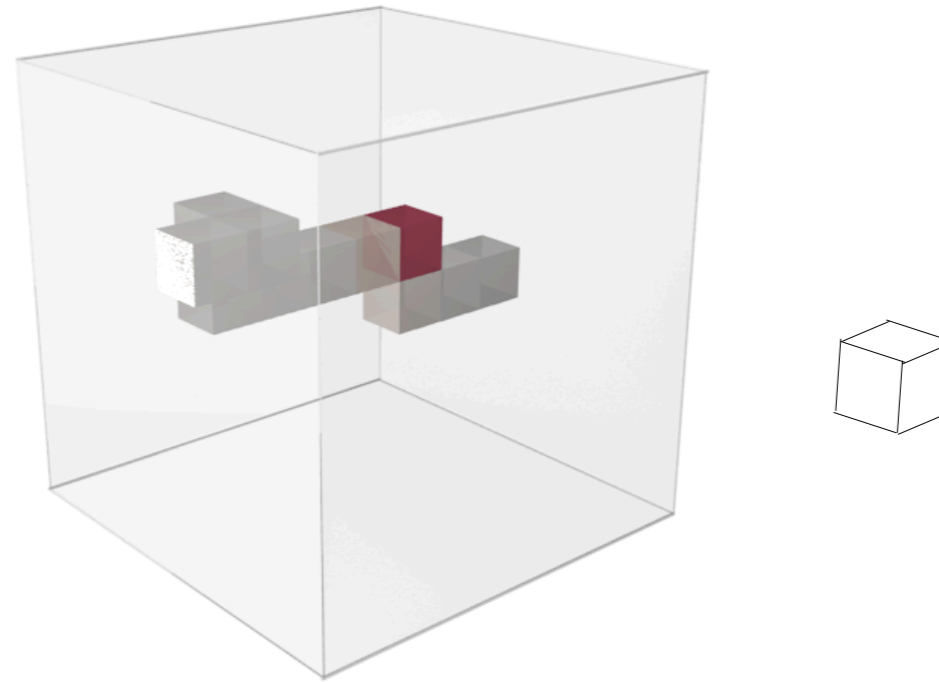
**r**

$$p(z_r = i) = \begin{cases} (1 - x_i^r) \prod_{j=1}^{i-1} x_j^r, & \text{if } i \leq N_r \\ \prod_{j=1}^{N_r} x_j^r, & \text{if } i = N_r + 1 \end{cases}$$

# Probabilistic Ray Tracing



**x**

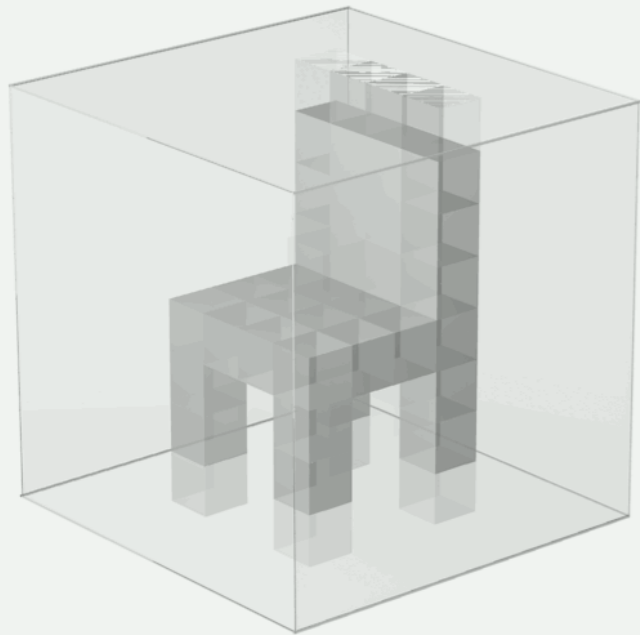


**r**

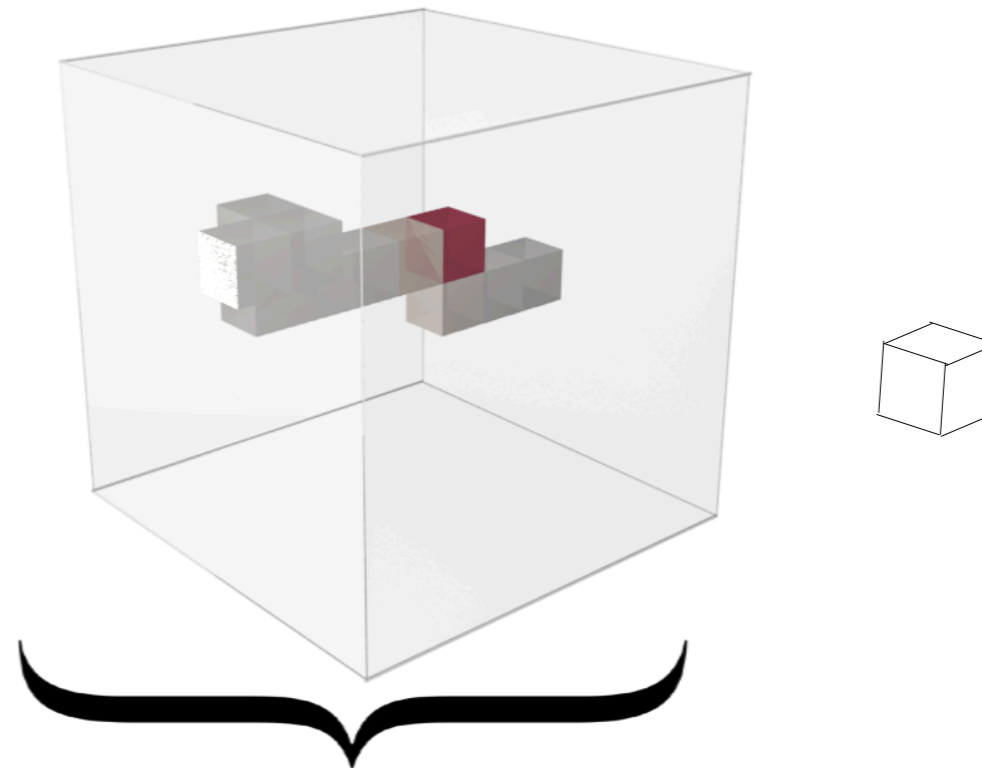
$$p(z_r = i) = \begin{cases} (1 - x_i^r) \prod_{j=1}^{i-1} x_j^r, & \text{if } i \leq N_r \\ \prod_{j=1}^{N_r} x_j^r, & \text{if } i = N_r + 1 \end{cases}$$



# Probabilistic Ray Tracing

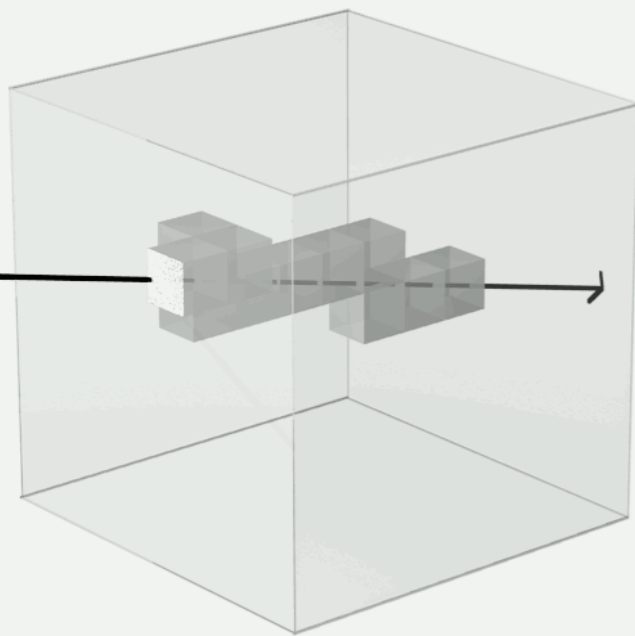


**x**



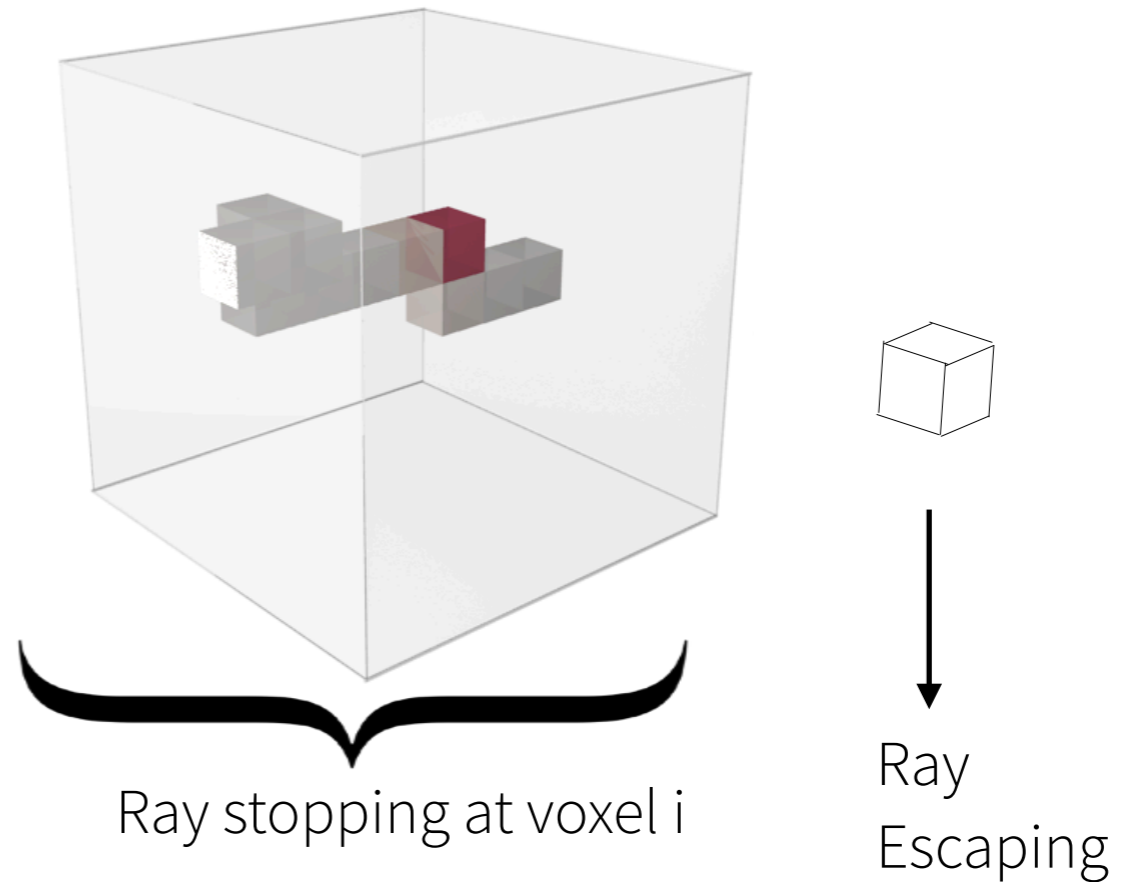
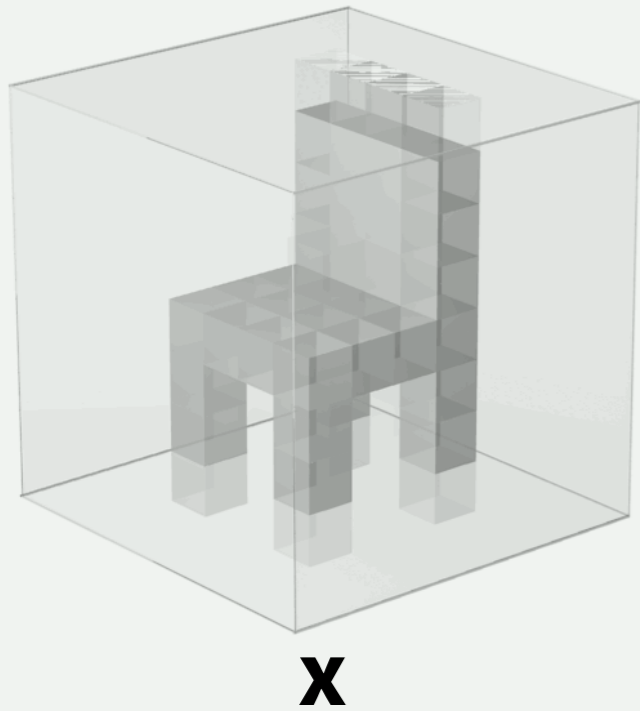
Ray stopping at voxel  $i$

$$p(z_r = i) = \begin{cases} (1 - x_i^r) \prod_{j=1}^{i-1} x_j^r, & \text{if } i \leq N_r \\ \prod_{j=1}^{N_r} x_j^r, & \text{if } i = N_r + 1 \end{cases}$$

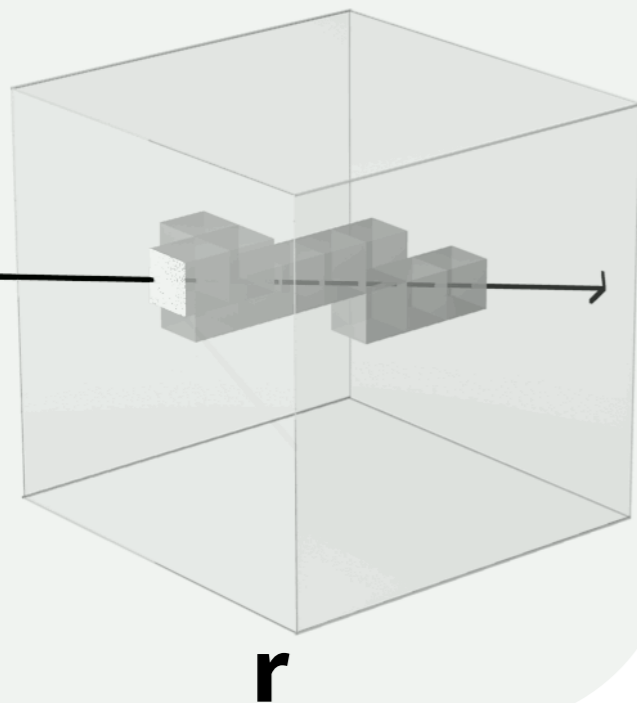


**r**

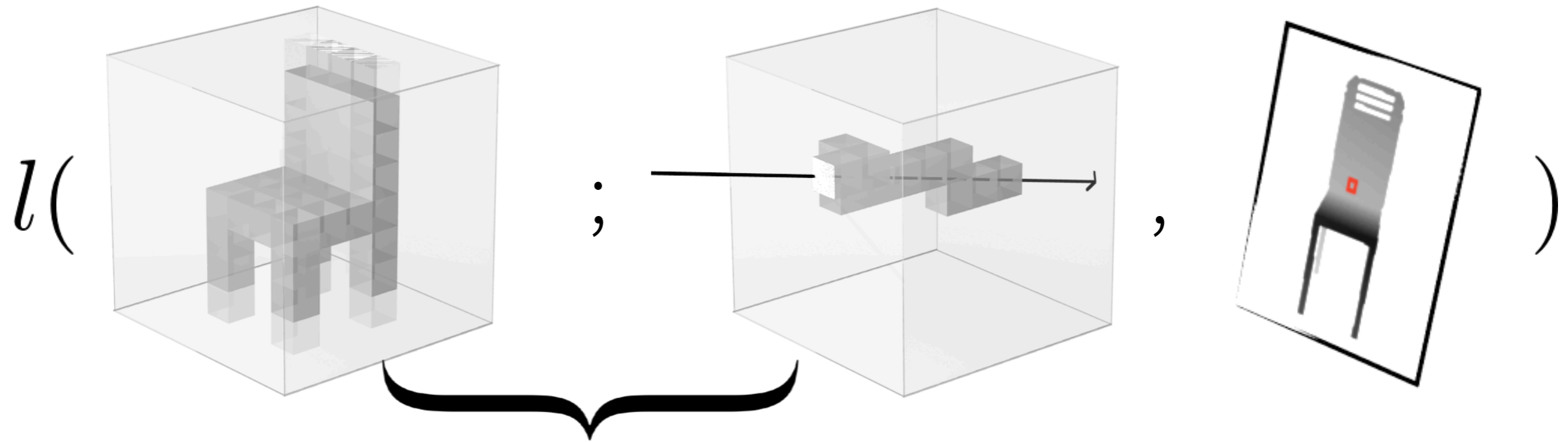
# Probabilistic Ray Tracing



$$p(z_r = i) = \begin{cases} (1 - x_i^r) \prod_{j=1}^{i-1} x_j^r, & \text{if } i \leq N_r \\ \prod_{j=1}^{N_r} x_j^r, & \text{if } i = N_r + 1 \end{cases}$$



# Differentiable Ray Consistency



$\Sigma$  (

A 3D scene with a ray path. A red cube is highlighted on the ray path. A small white cube is shown to the right of the scene.

**Event Probabilities**  
*(where can the ray stop ?)*

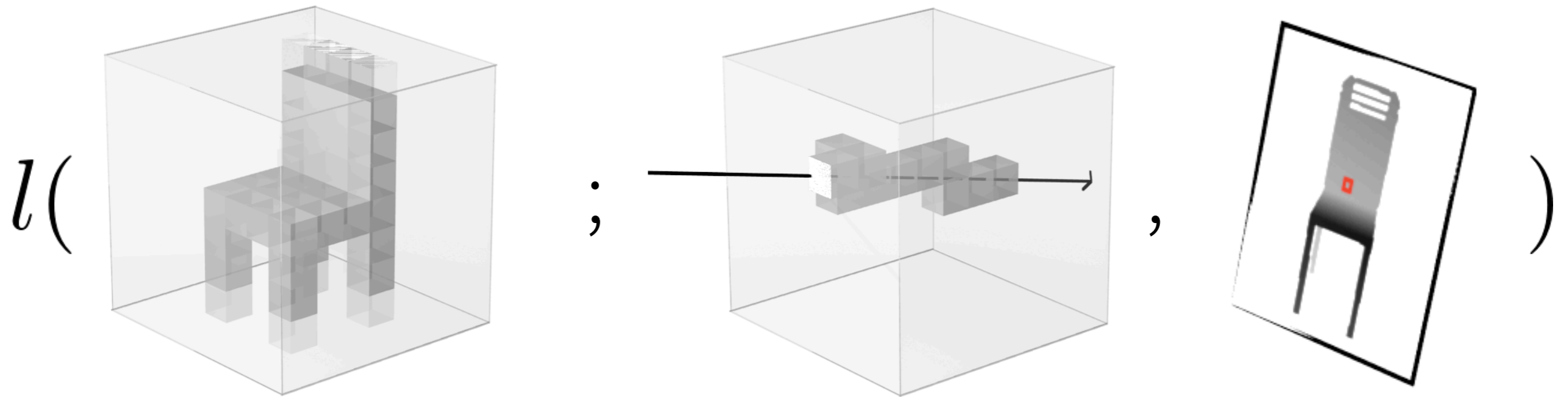
$\odot$

A 3D scene with a ray path. A red cube is highlighted on the ray path. A small red cube is shown to the right of the scene.

**Event Costs**  
*(how 'bad' is stopping here ?)*

)

# Differentiable Ray Consistency



$\Sigma$  (

A 3D scene with a ray passing through a chair. A small white cube is shown next to the ray's path, representing a potential stopping point. The entire diagram is enclosed in a rounded rectangle with a red border.

**Event Probabilities**  
*(where can the ray stop ?)*

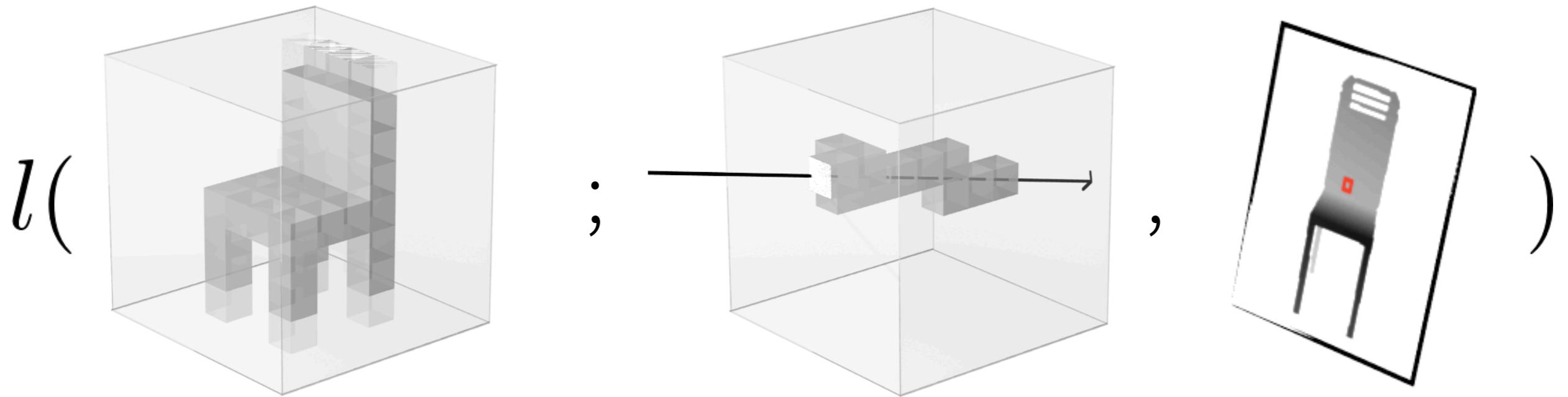
$\odot$

A 3D scene with a ray passing through a chair. A small red cube is shown next to the ray's path, representing a cost associated with stopping at that location. The entire diagram is enclosed in a rounded rectangle with a light gray border.

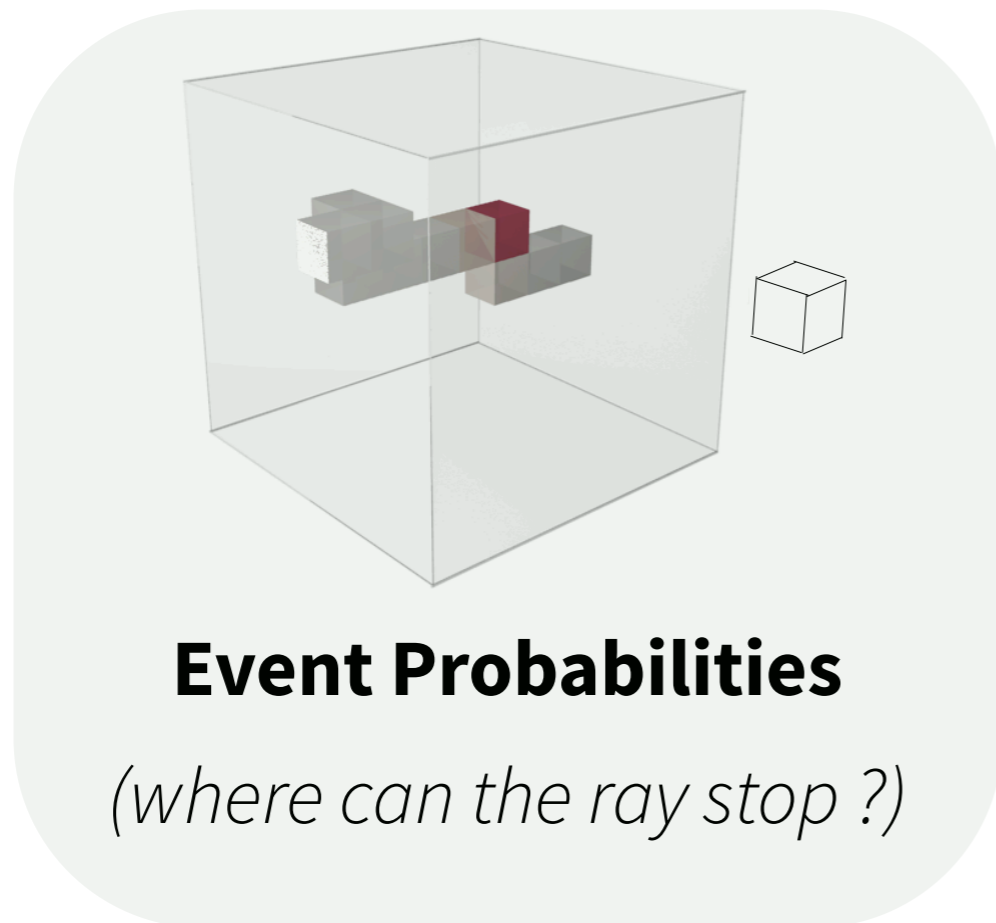
**Event Costs**  
*(how 'bad' is stopping here ?)*

)

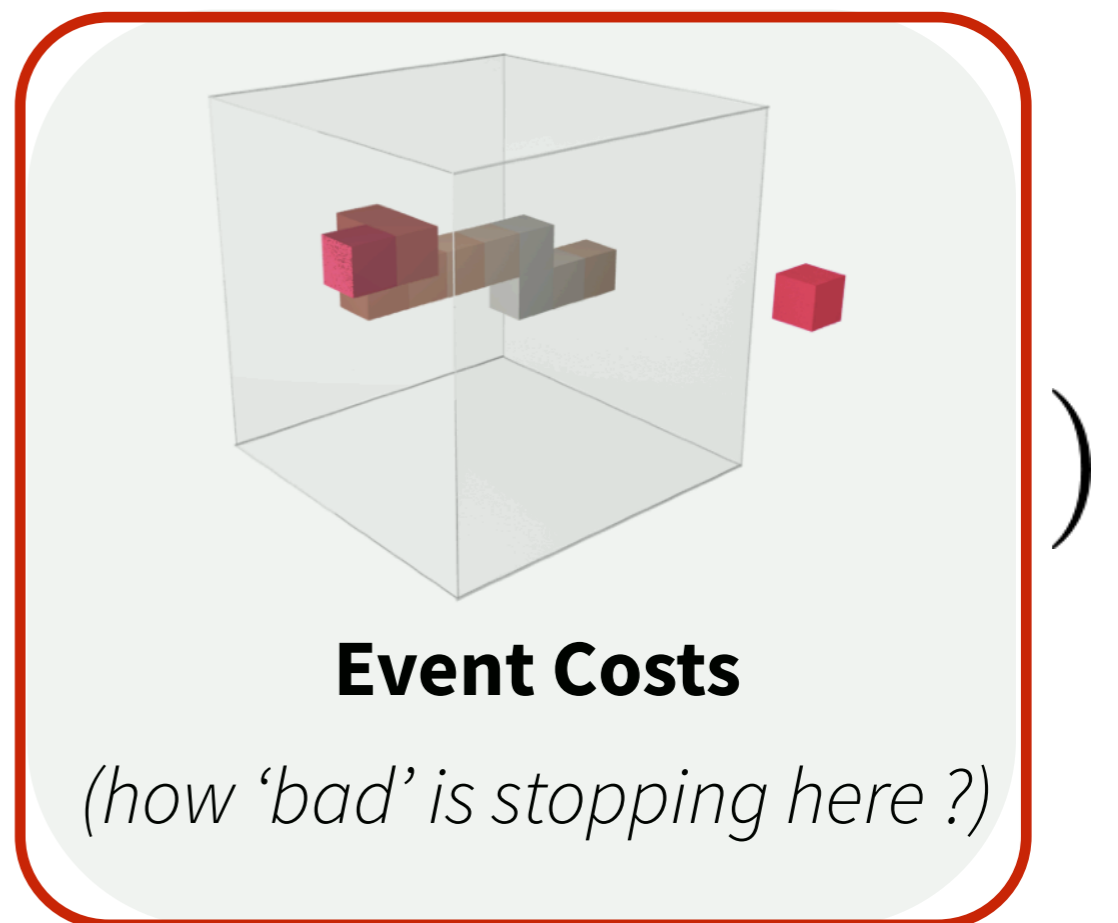
# Differentiable Ray Consistency



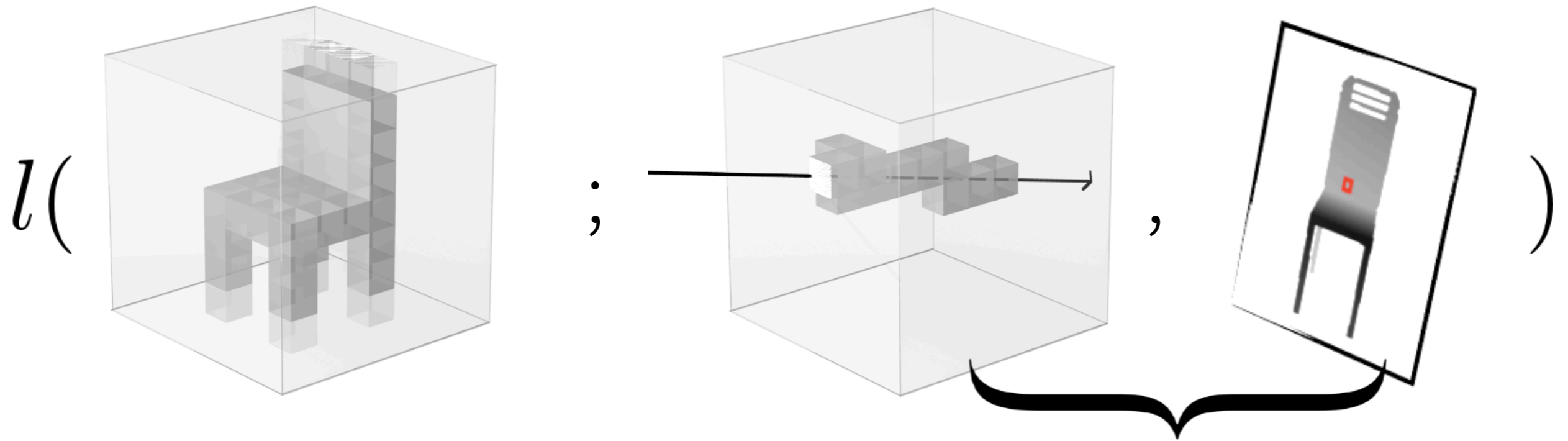
$\Sigma($



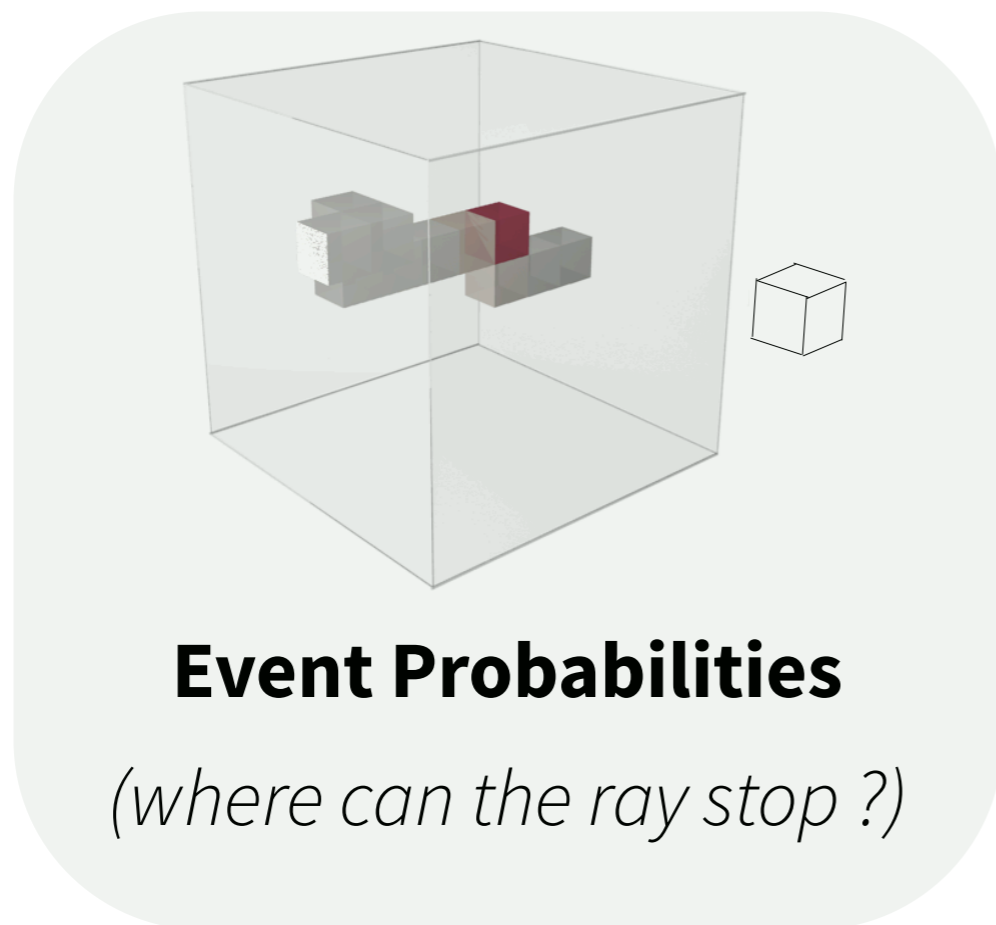
$\odot$



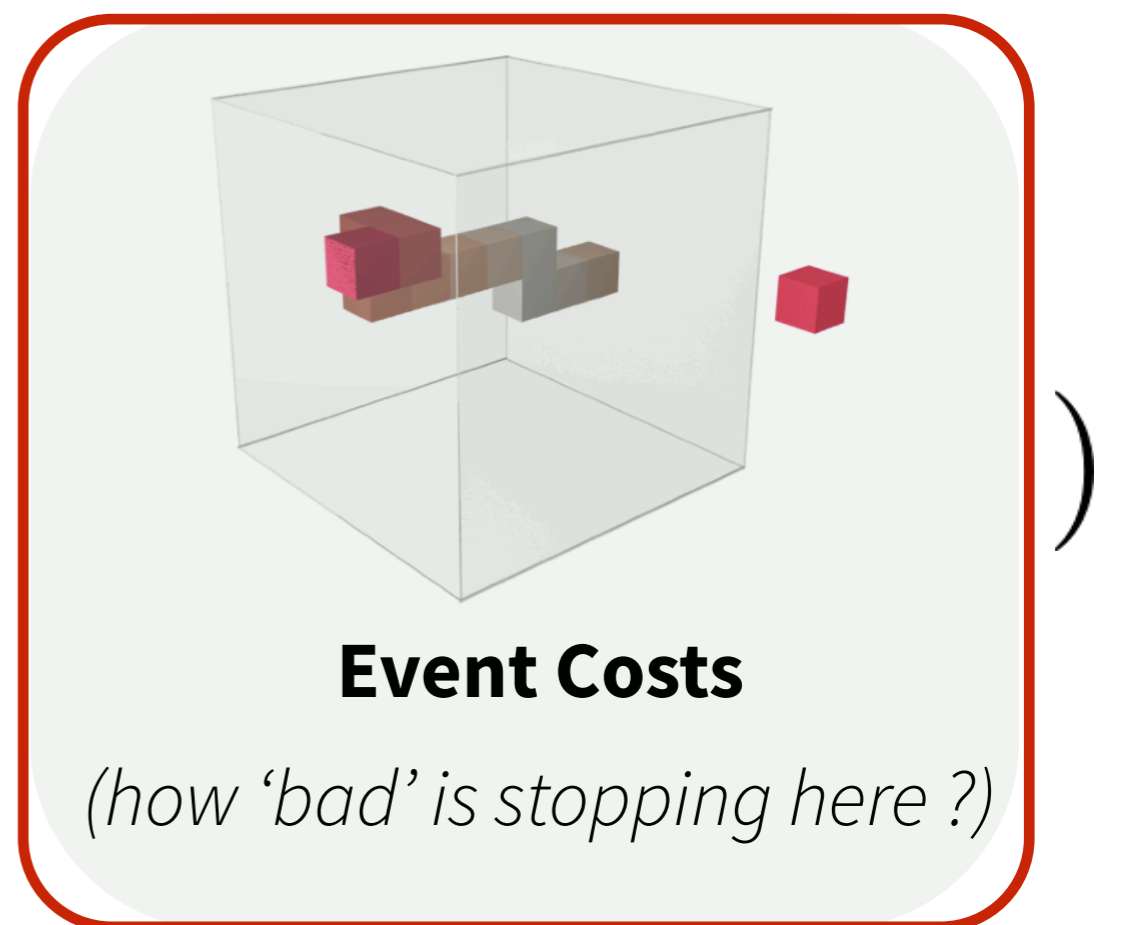
# Differentiable Ray Consistency



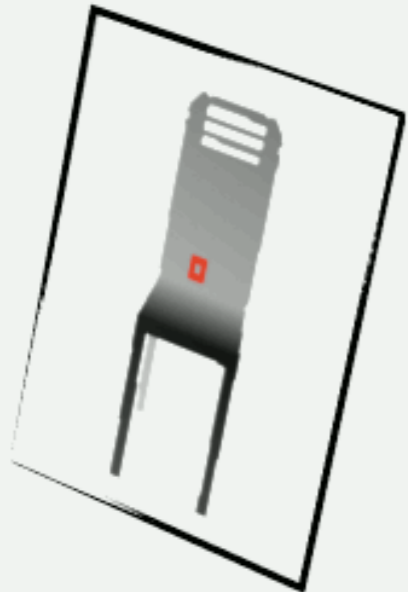
$\Sigma$  (



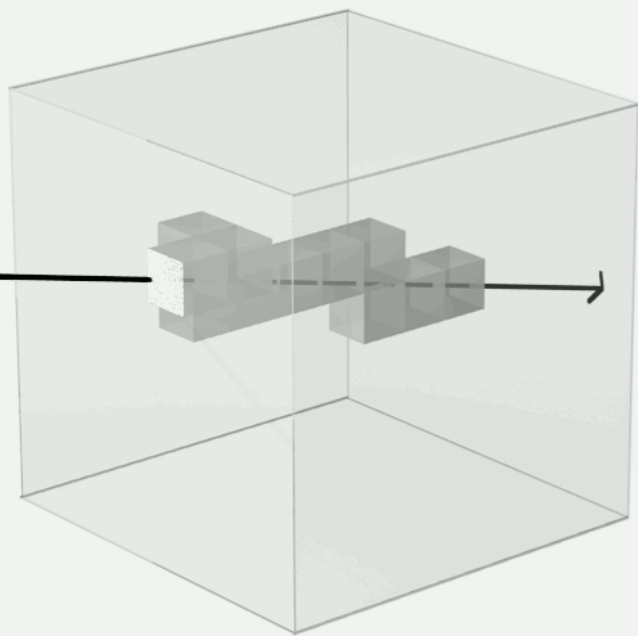
$\odot$



# Event Costs



$\mathbf{o}_r$  (*Depth*)



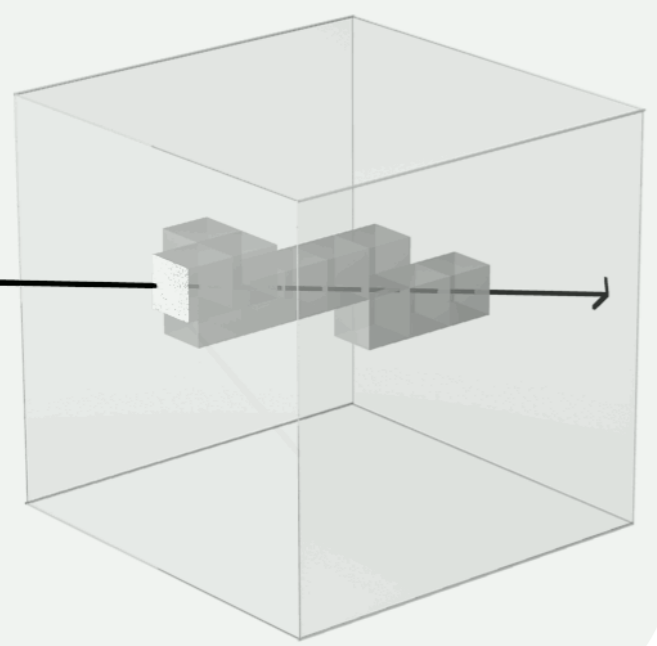
$\mathbf{r}$

How inconsistent is each event w.r.t  $\mathbf{o}_r$  ?

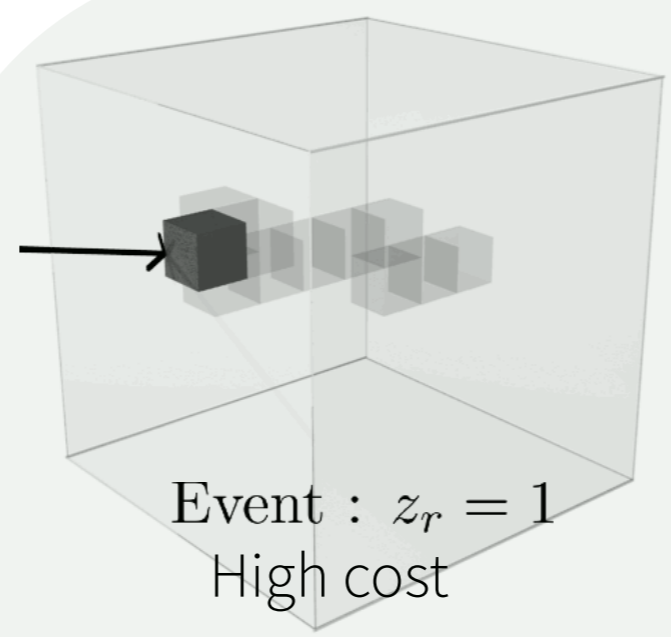
# Event Costs



$\mathbf{o}_r$  (Depth)



$\mathbf{r}$



Event :  $z_r = 1$   
High cost

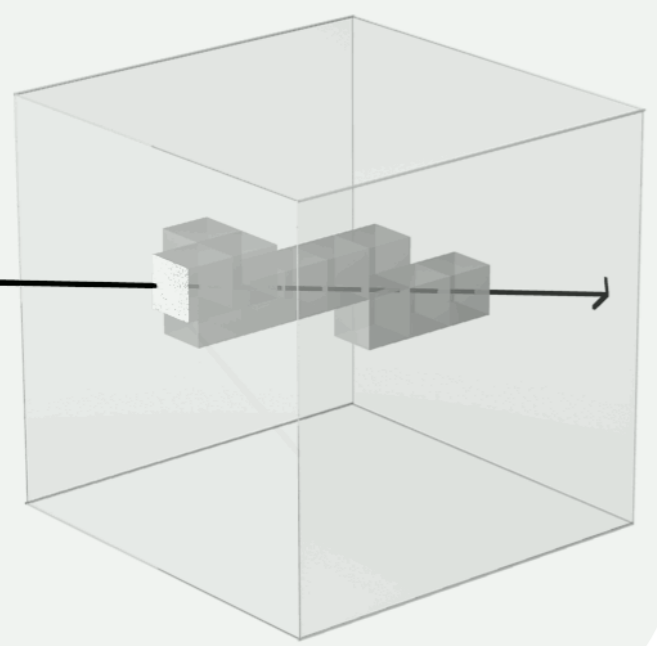
How inconsistent is each event w.r.t  $\mathbf{o}_r$  ?



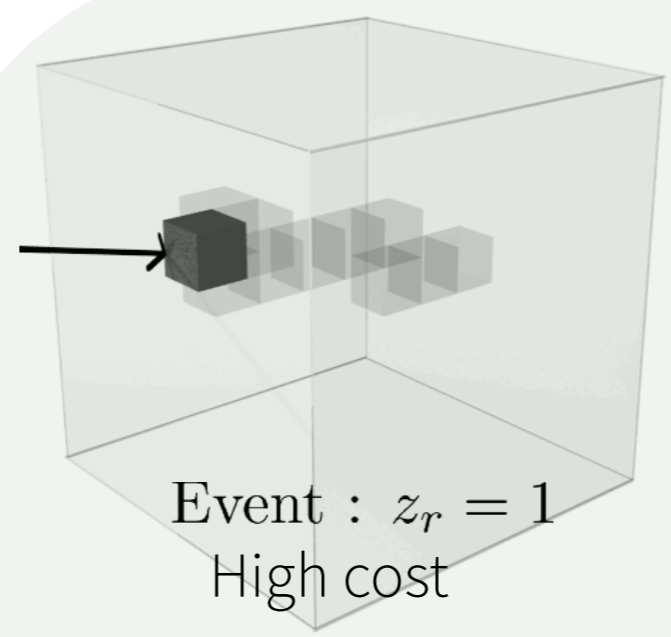
# Event Costs



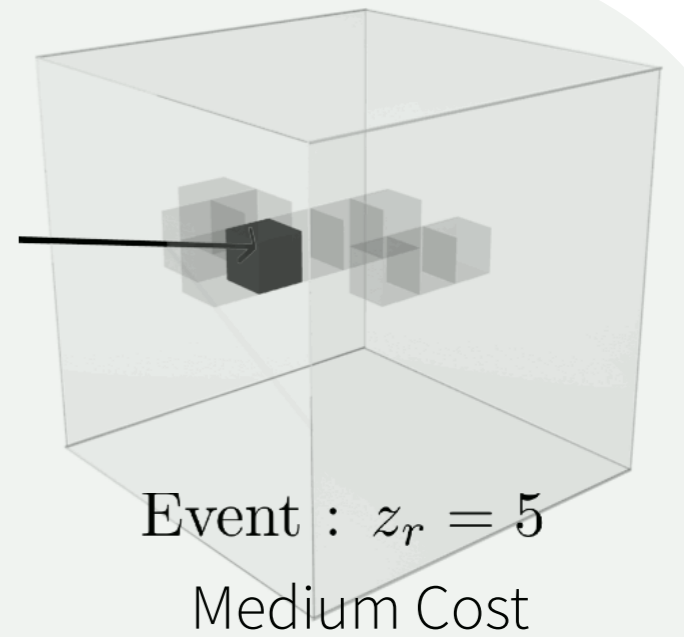
$\mathbf{O}_r$  (*Depth*)



$\mathbf{r}$



Event :  $z_r = 1$   
High cost



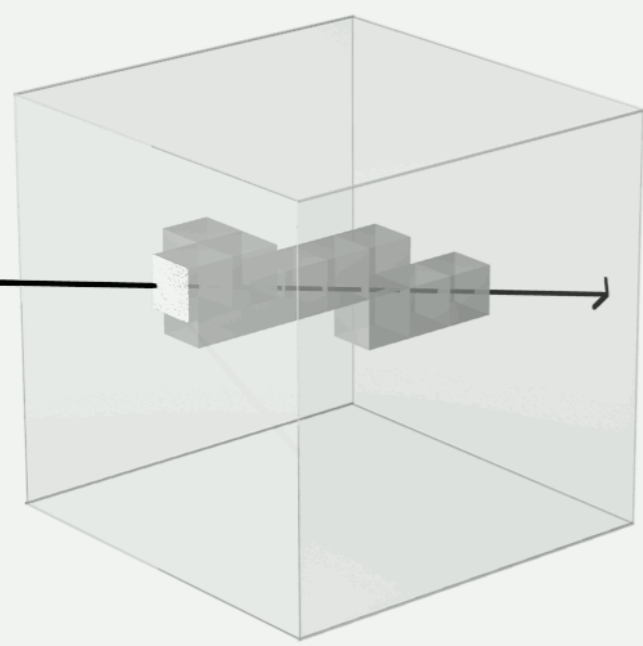
Event :  $z_r = 5$   
Medium Cost

How inconsistent is each event w.r.t  $\mathbf{O}_r$  ?

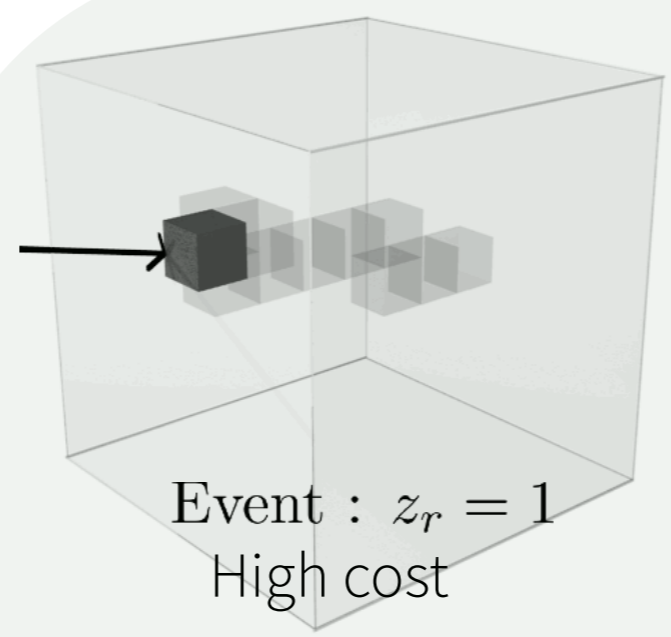
# Event Costs



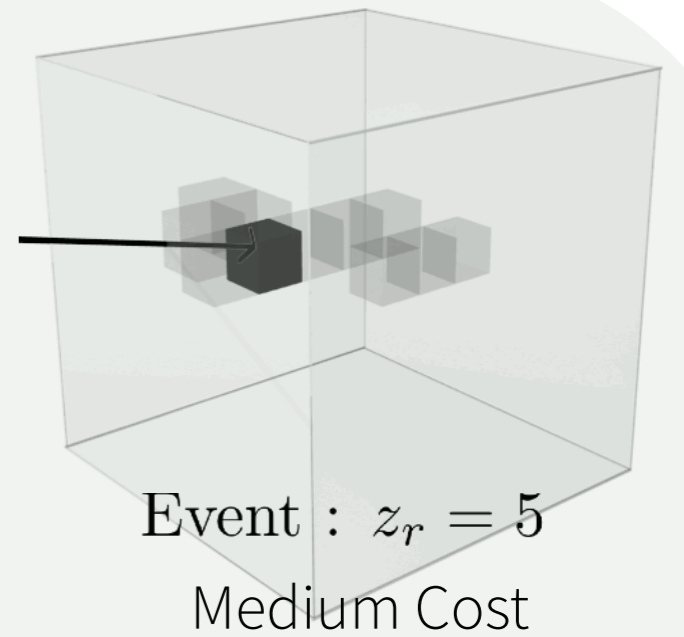
$\mathbf{O}_r$  (Depth)



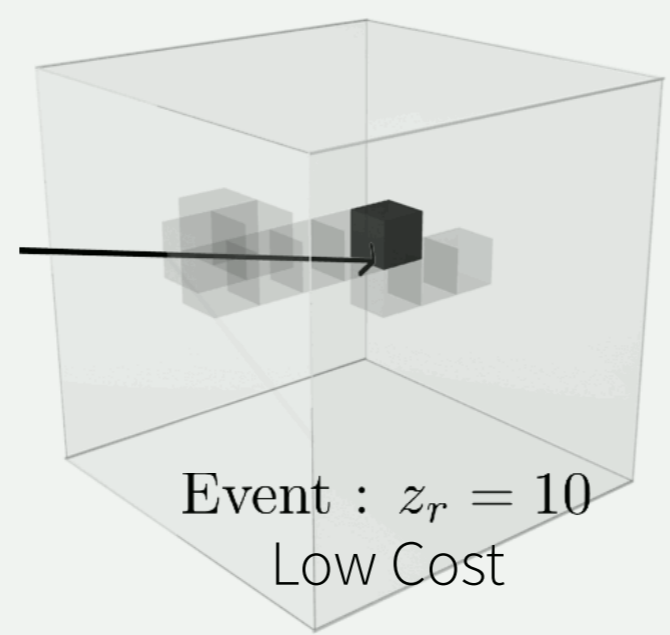
$\mathbf{r}$



Event :  $z_r = 1$   
High cost



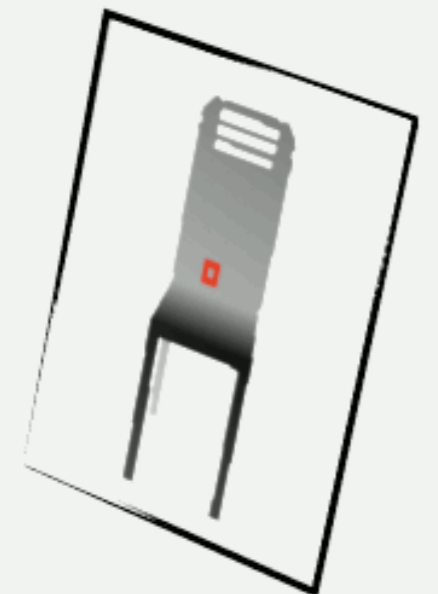
Event :  $z_r = 5$   
Medium Cost



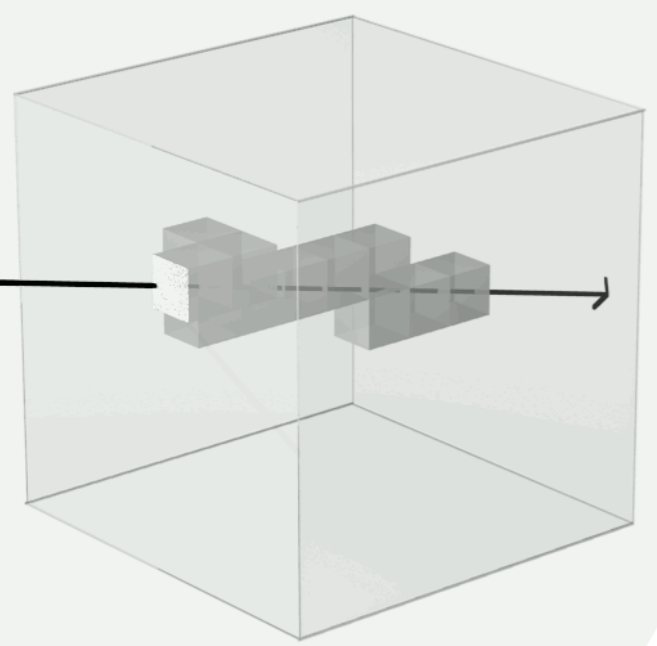
Event :  $z_r = 10$   
Low Cost

How inconsistent is each event w.r.t  $\mathbf{O}_r$  ?

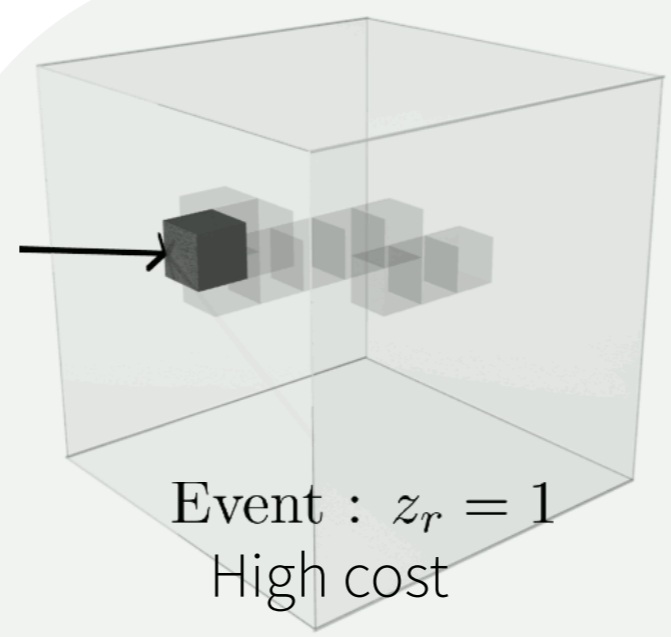
# Event Costs



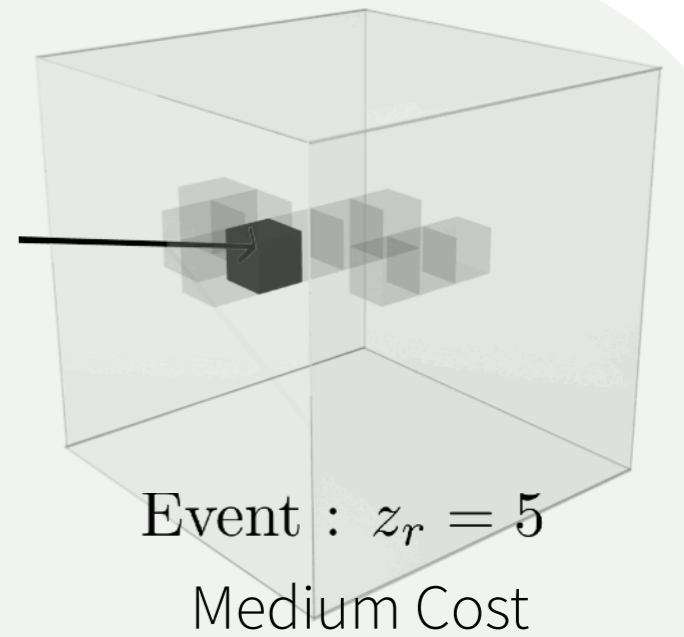
$\mathbf{O}_r$  (Depth)



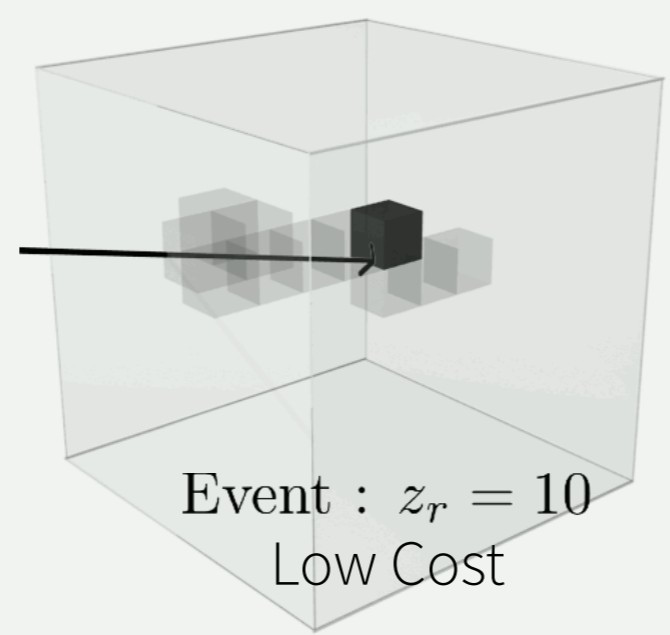
$\mathbf{r}$



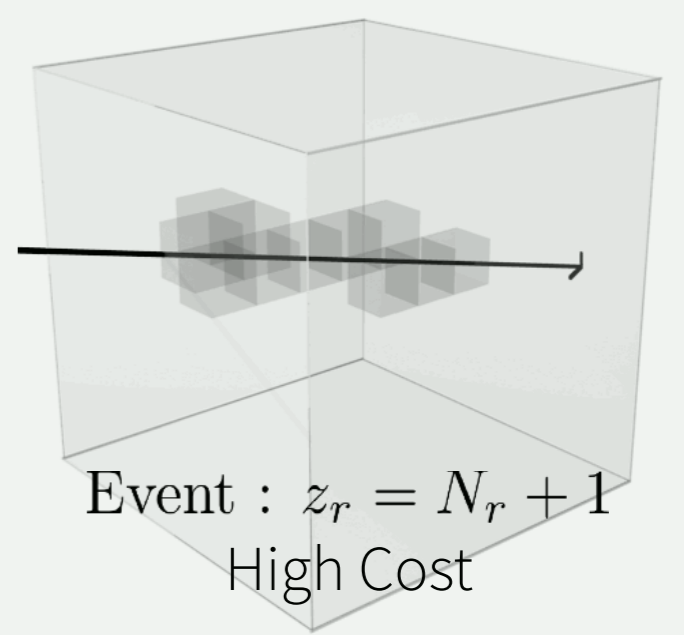
Event :  $z_r = 1$   
High cost



Event :  $z_r = 5$   
Medium Cost



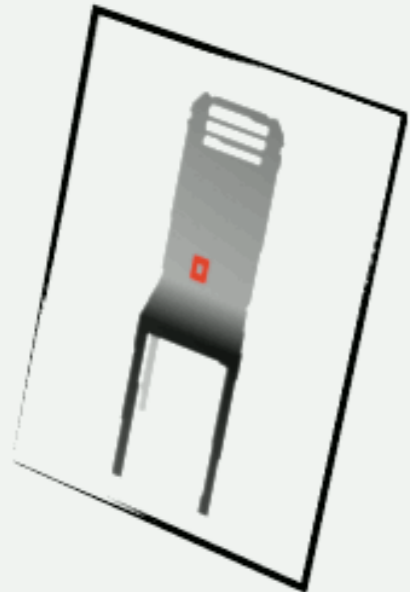
Event :  $z_r = 10$   
Low Cost



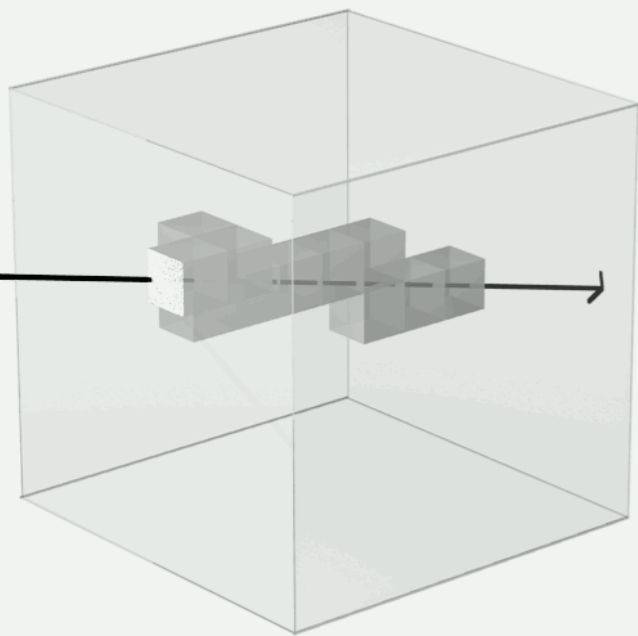
Event :  $z_r = N_r + 1$   
High Cost

How inconsistent is each event w.r.t  $\mathbf{O}_r$  ?

# Event Costs



$\mathbf{O}_r$  (*Depth*)



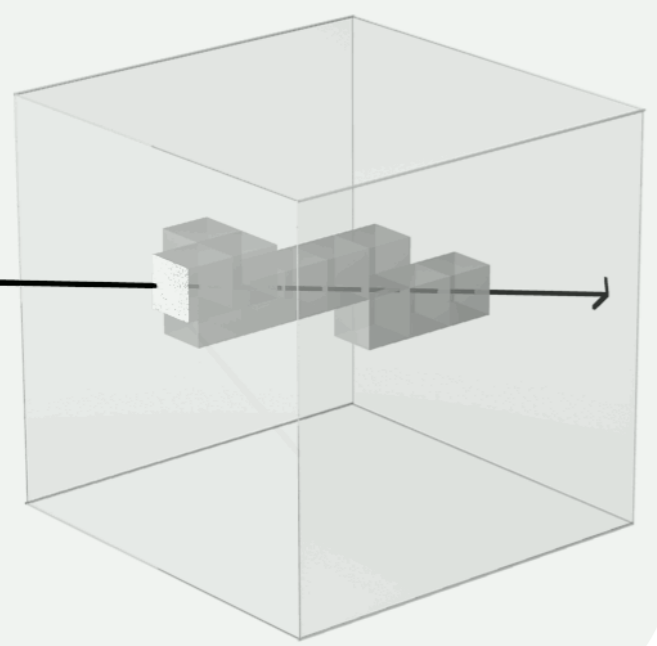
$\mathbf{r}$

$$\psi_r^{\text{depth}}(i) = |d_i^r - d_r|$$

# Event Costs



$\mathbf{O}_r$  (*Depth*)



$\mathbf{r}$

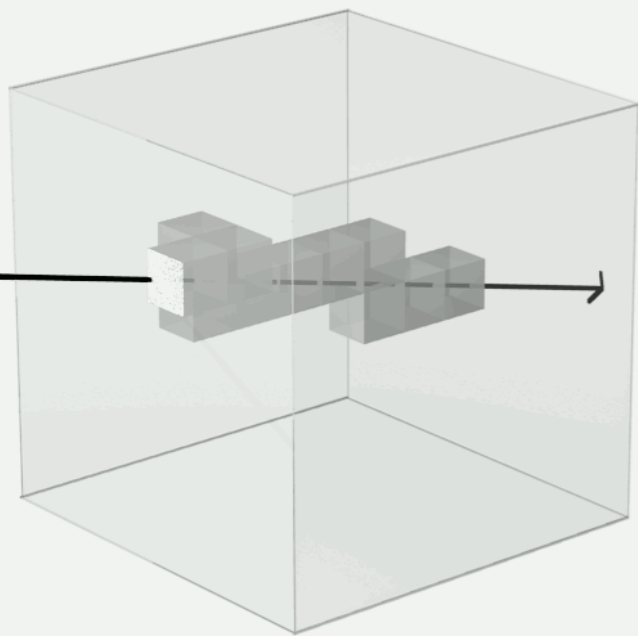
$$\psi_r^{depth}(i) = |d_i^r - \boxed{d_r}|$$

Observed  
Depth

# Event Costs



$\mathbf{O}_r$  (*Depth*)



$\mathbf{r}$

$$\psi_r^{depth}(i) = |d_i^r - d_r|$$

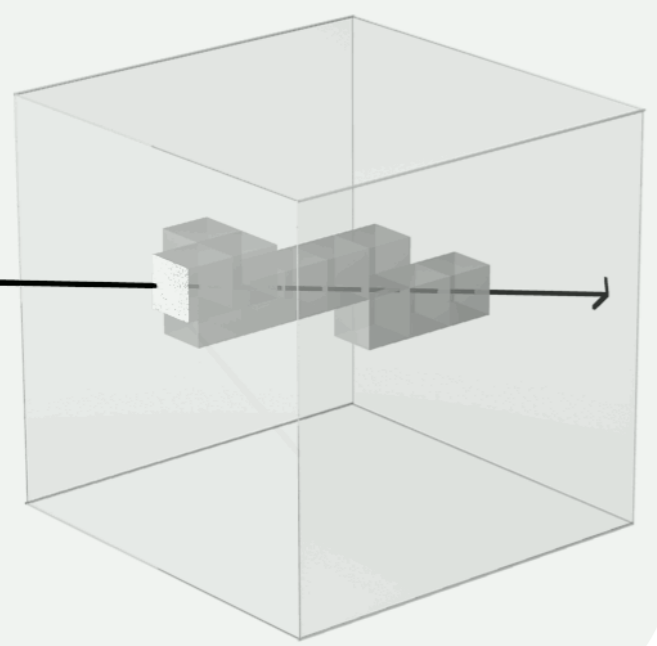
Depth under  
event  $i$

Observed  
Depth

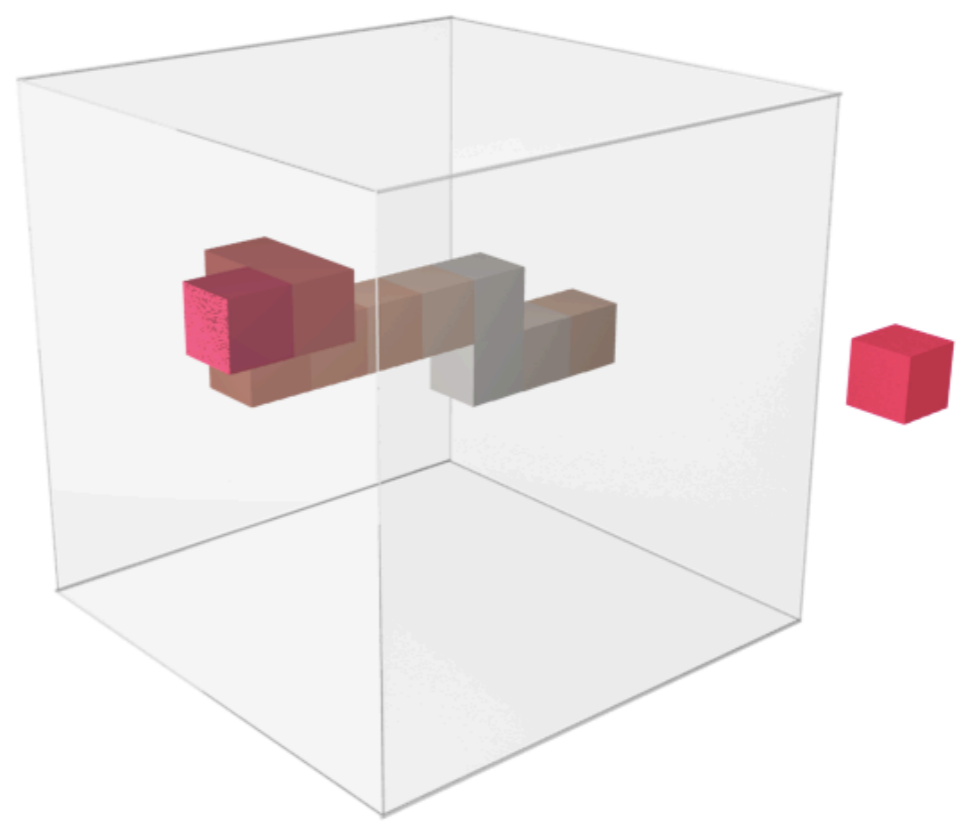
# Event Costs



$\mathbf{O}_r$  (*Depth*)



$\mathbf{r}$



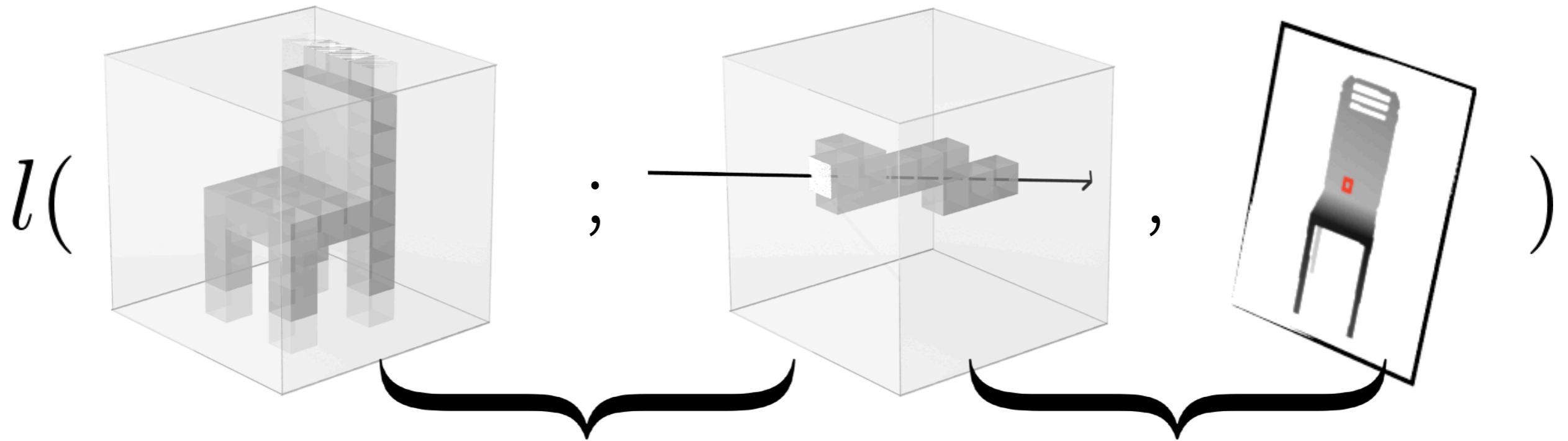
Event Costs

$$\psi_r^{depth}(i) = |d_i^r - d_r|$$

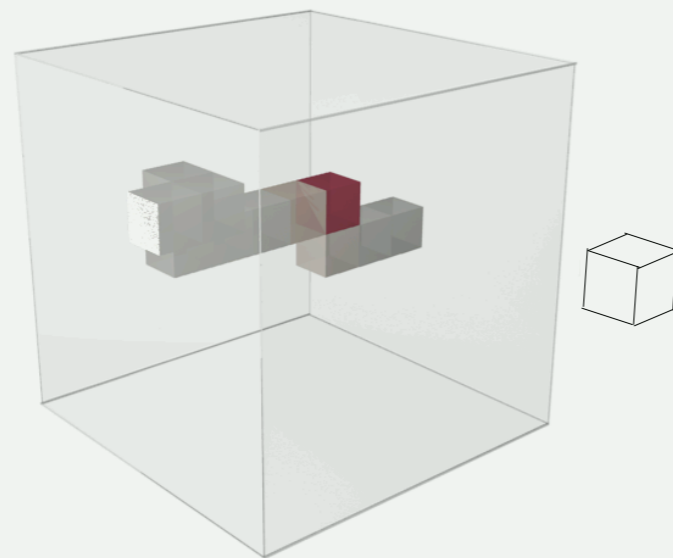
Depth under event  $i$

Observed Depth

# Differentiable Ray Consistency



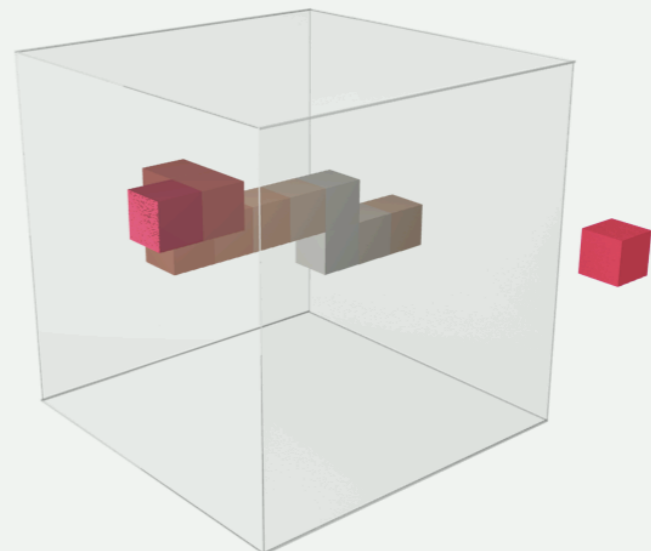
$\Sigma($



**Event Probabilities**

*(where can the ray stop ?)*

$\odot$

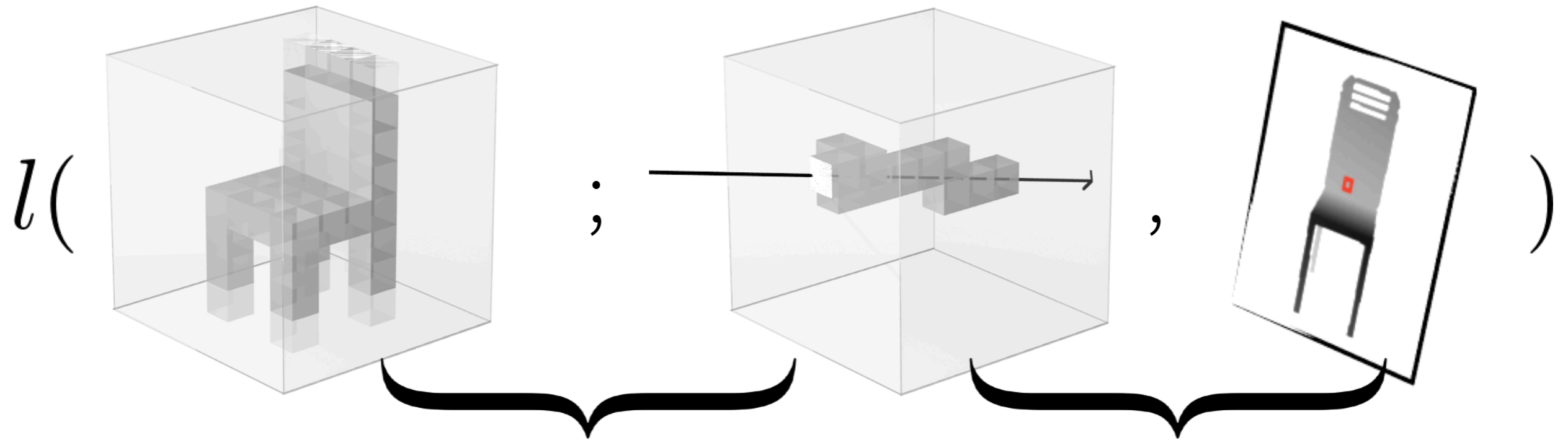


**Event Costs**

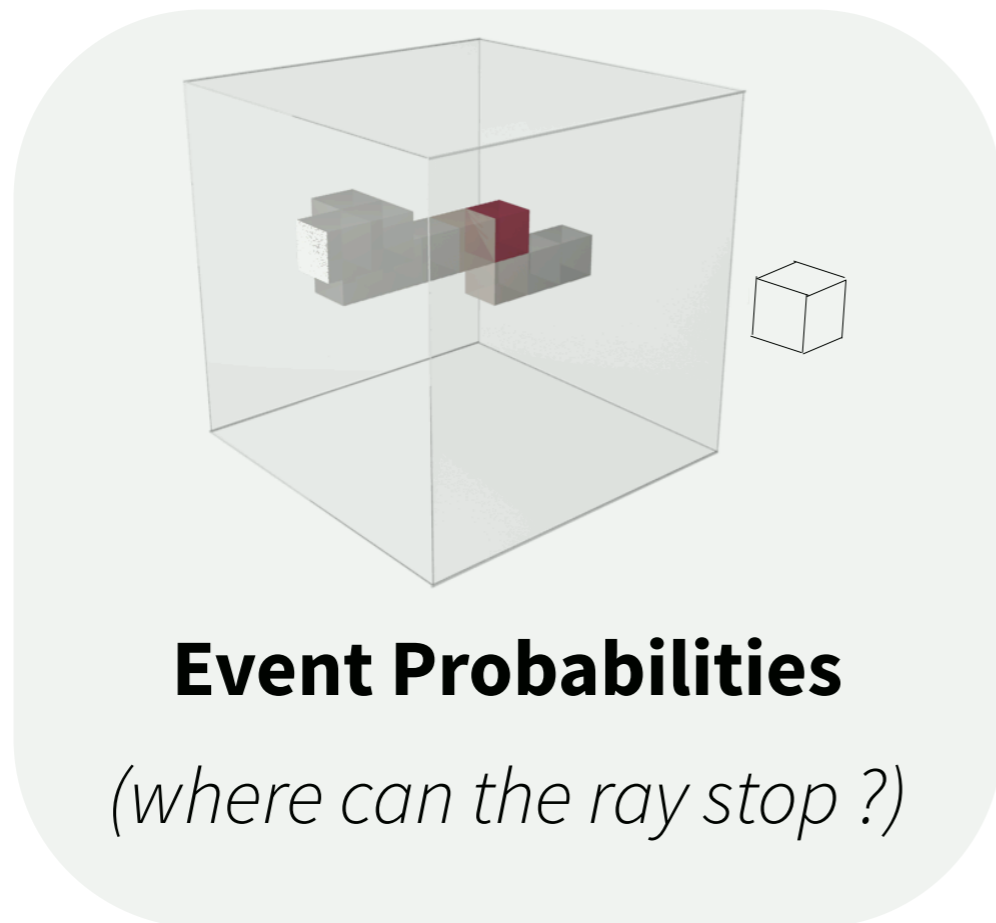
*(how 'bad' is stopping here ?)*



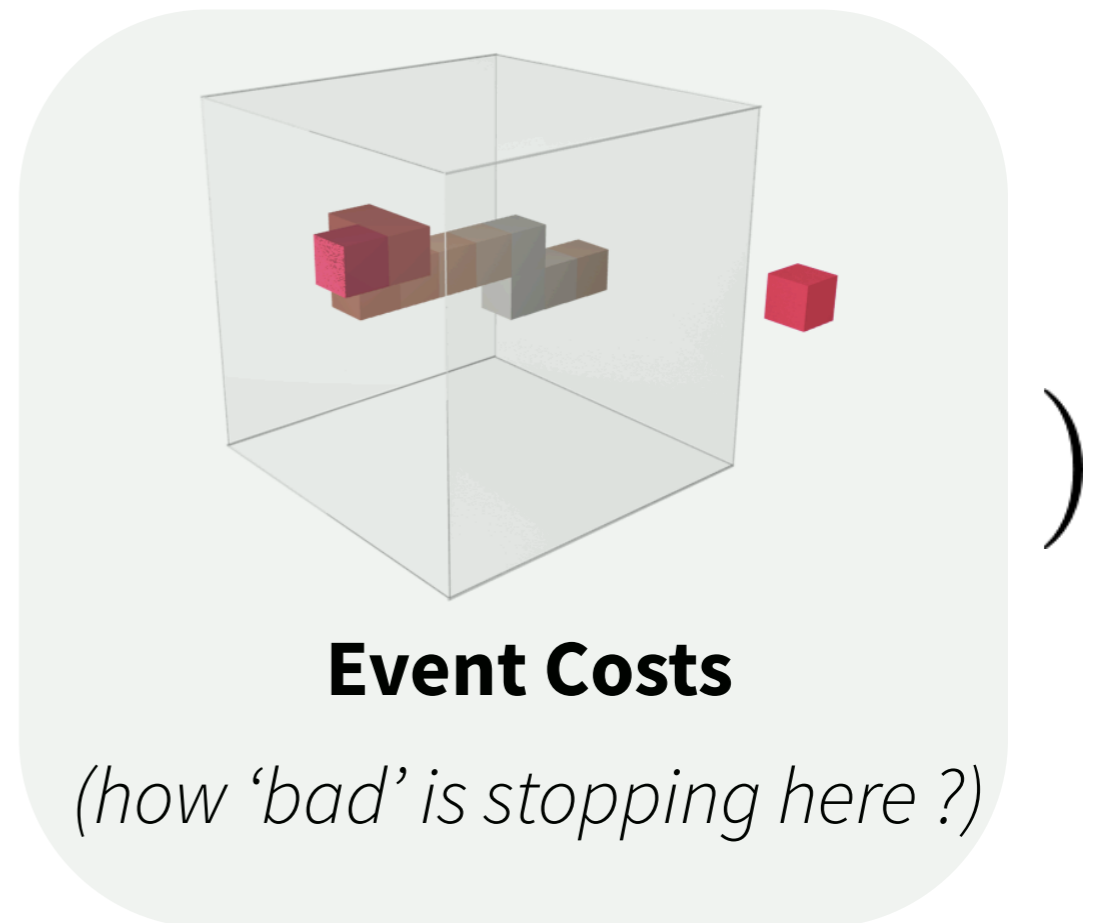
# Differentiable Ray Consistency



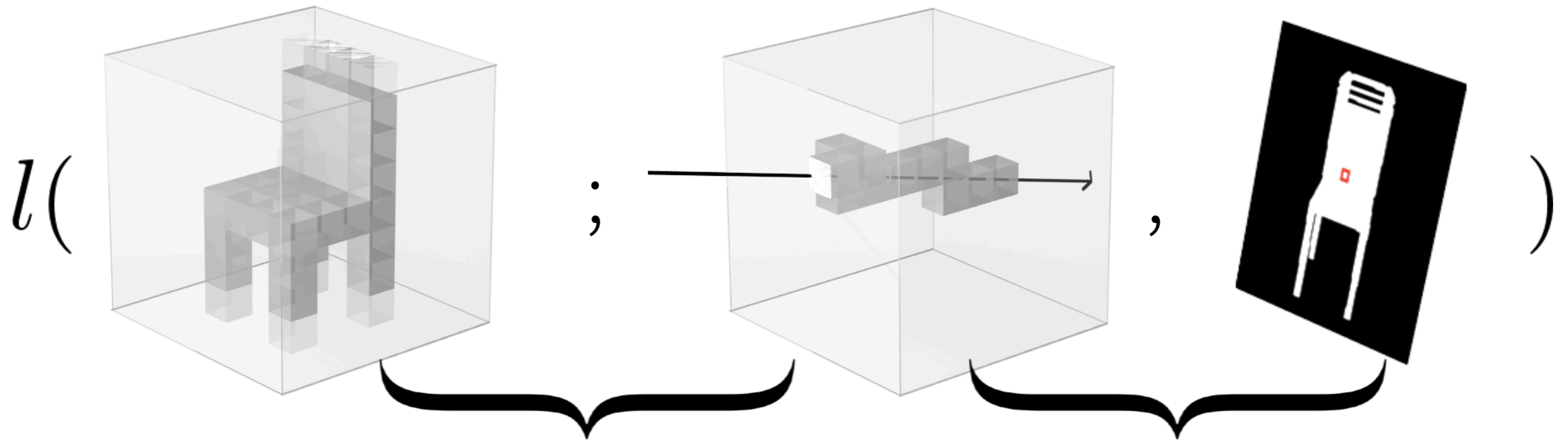
$\Sigma$  (



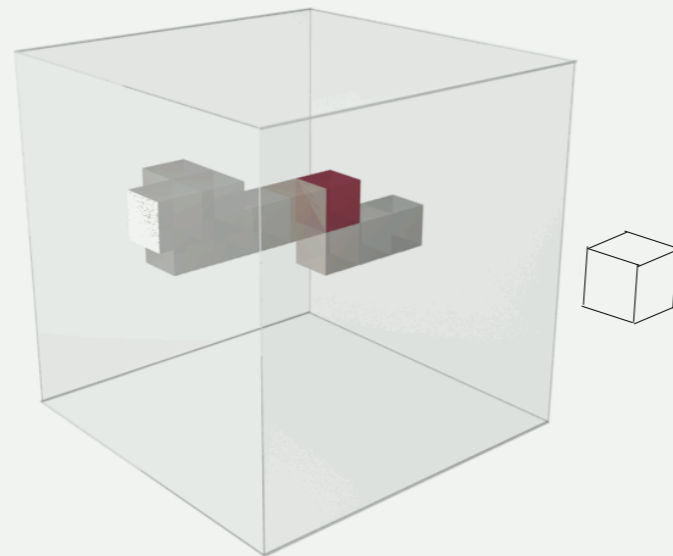
$\odot$



# Mask Observation

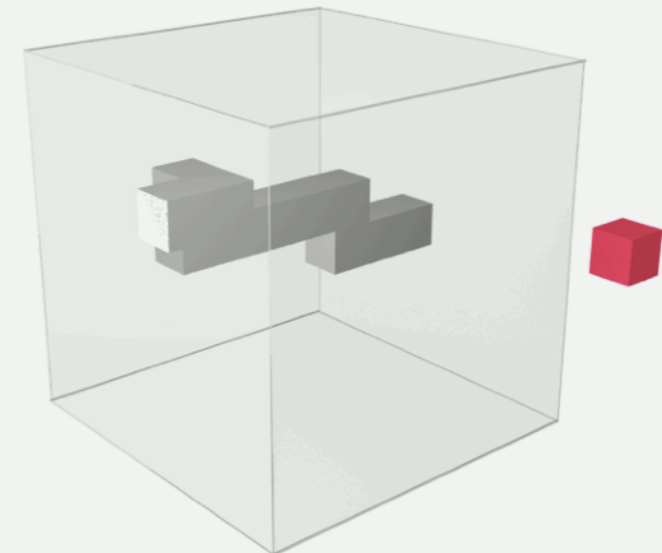


$\Sigma($



**Event Probabilities**

*(where can the ray stop ?)*

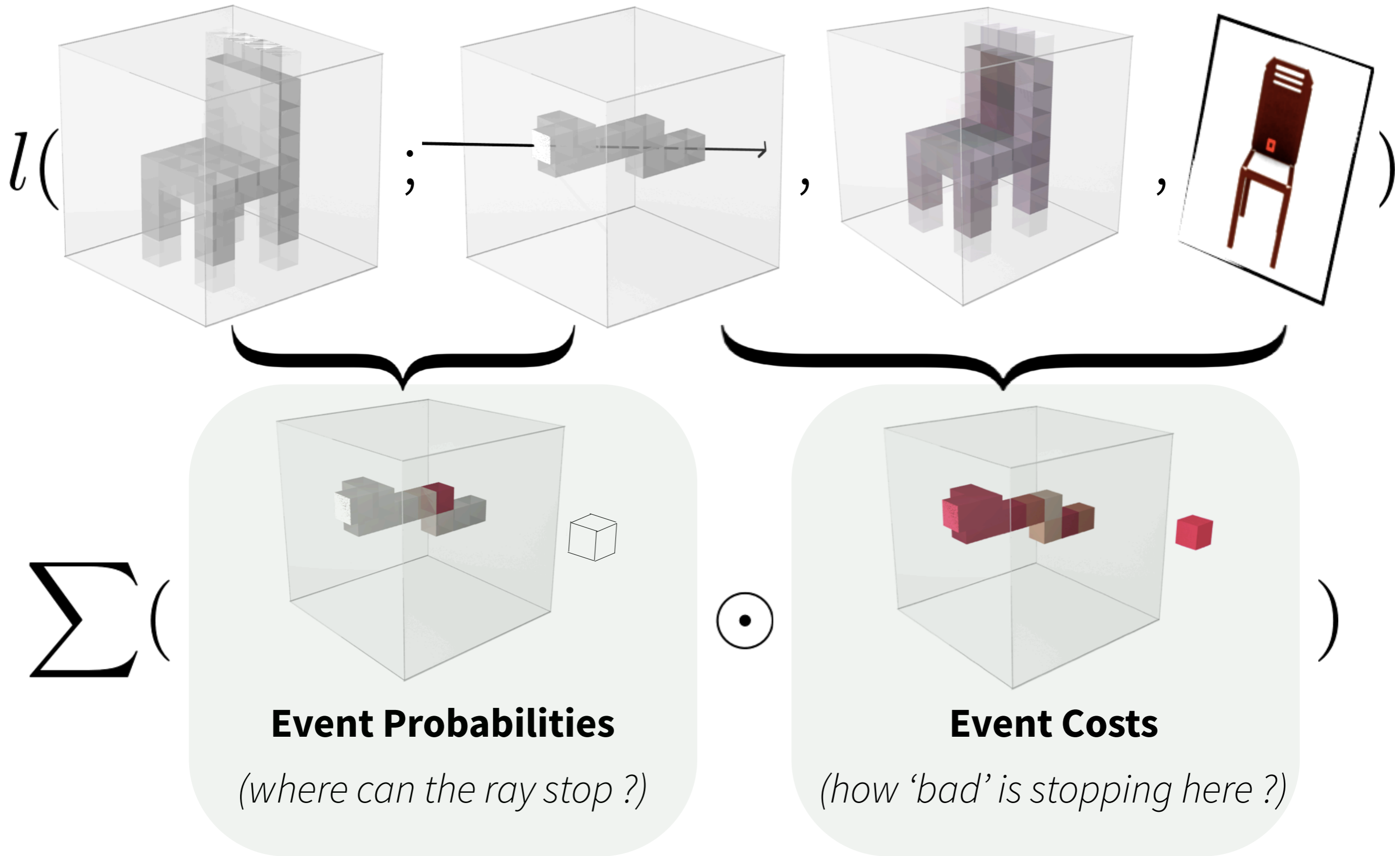


**Event Costs**

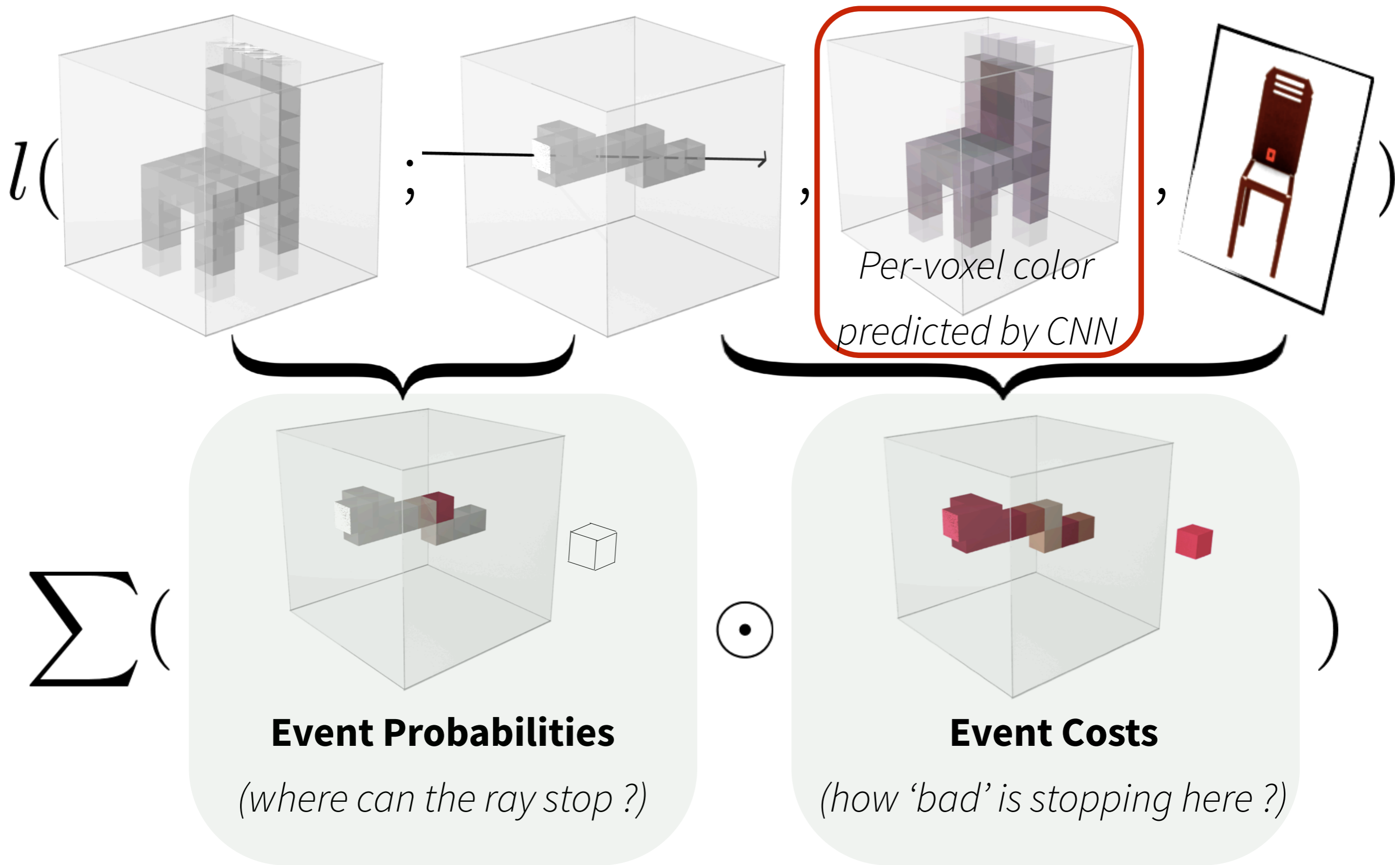
*(how 'bad' is stopping here ?)*

)

# Color Observation



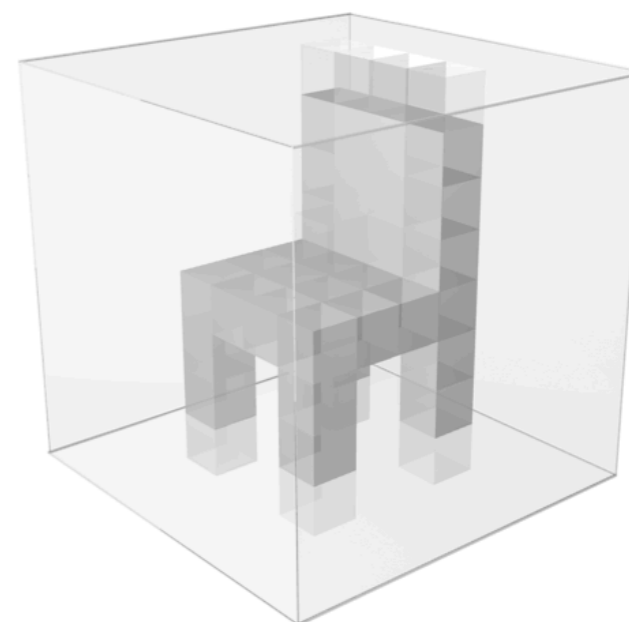
# Color Observation



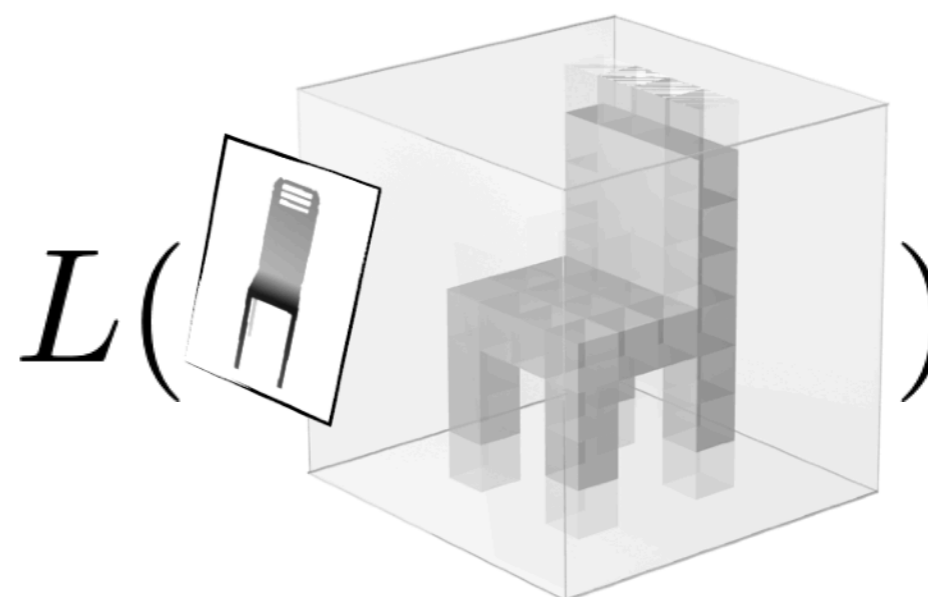
# Learning via Geometric Consistency



Input Image



Observation  $\mathbf{O}$   
from camera  $\mathbf{C}$

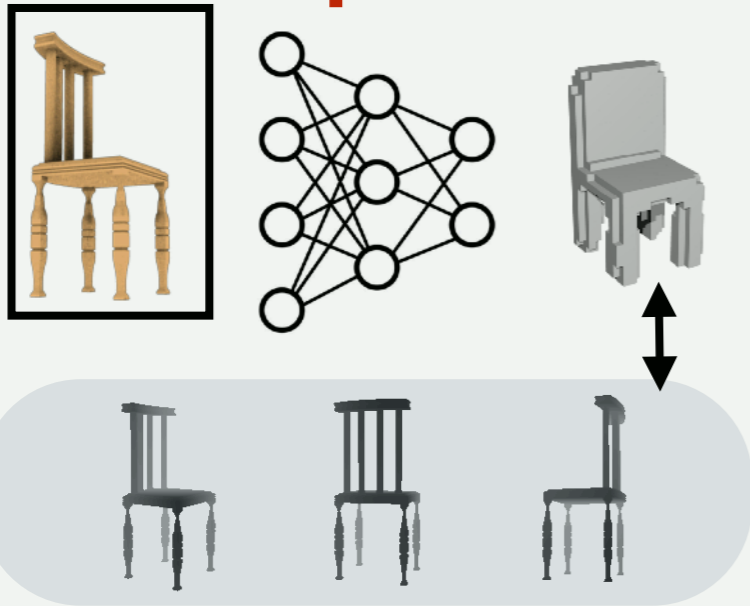


Geometric  
Consistency Loss

# Learning Single-view Reconstruction

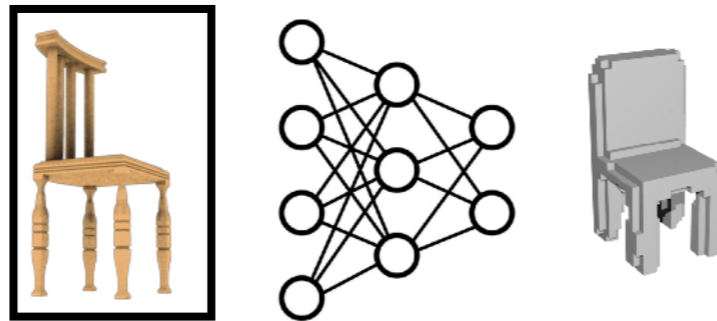
# Learning Single-view Reconstruction

## ShapeNet



Supervision : Pose + Depth/Mask

# Experiments - ShapeNet



Input

GT

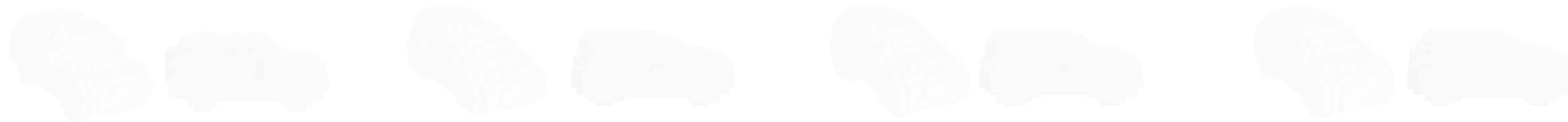
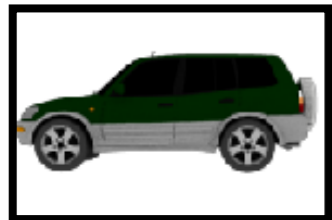
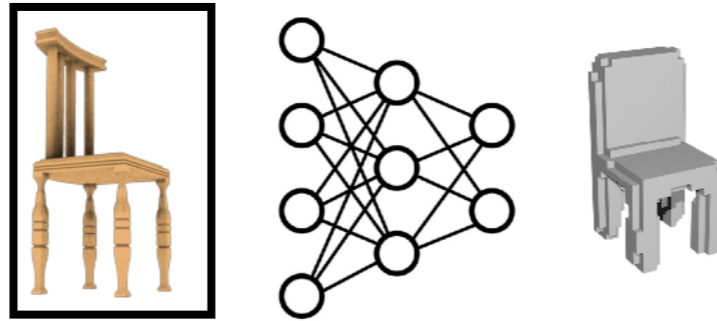
DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision



# Experiments - ShapeNet



Input

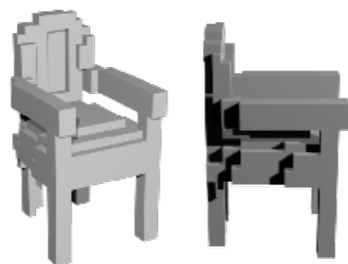
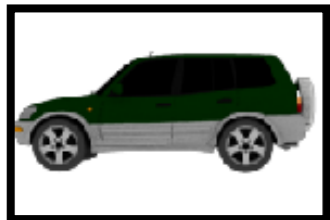
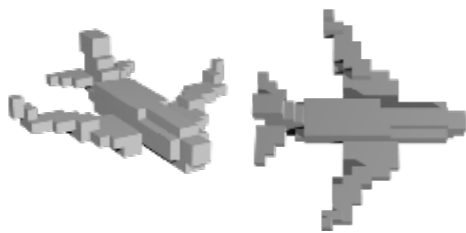
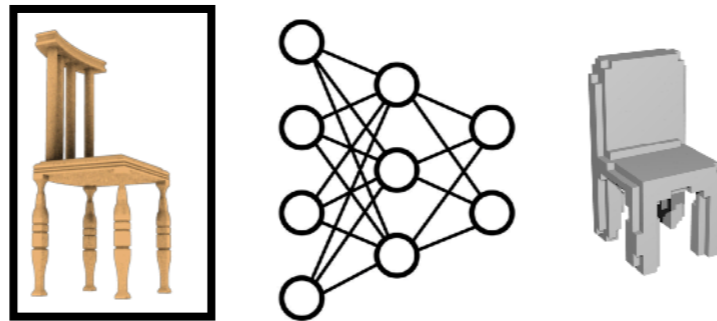
GT

DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision

# Experiments - ShapeNet



Input

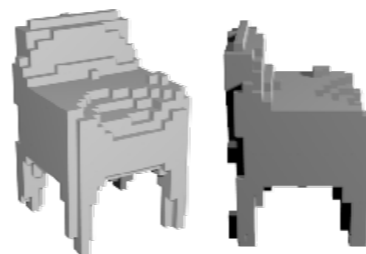
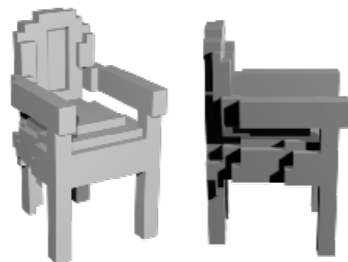
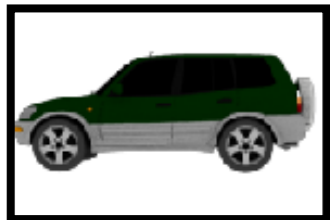
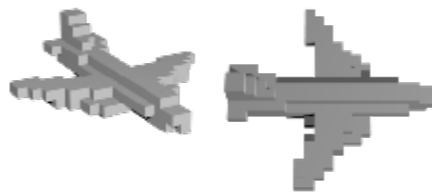
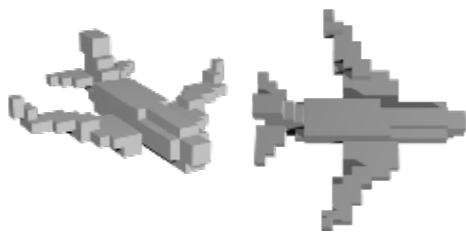
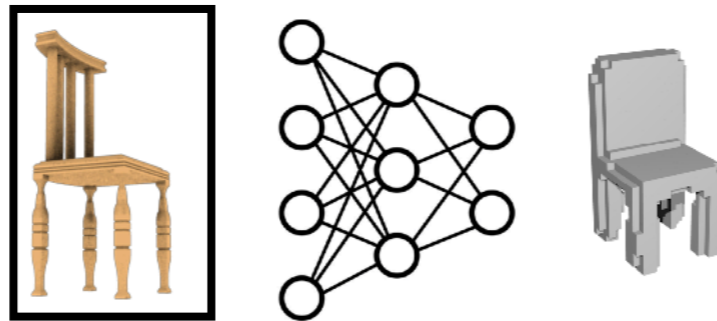
GT

DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision

# Experiments - ShapeNet



Input

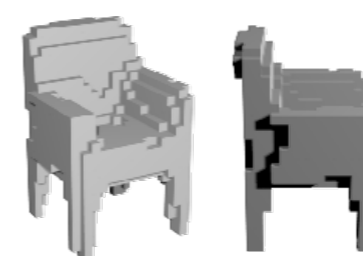
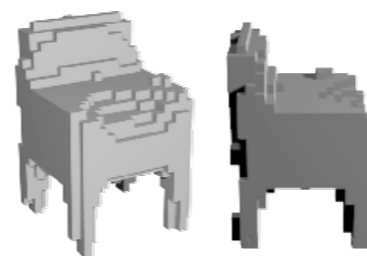
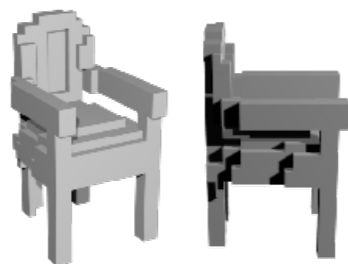
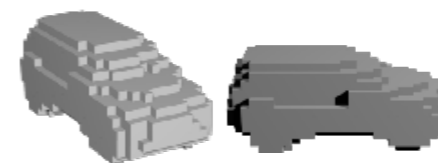
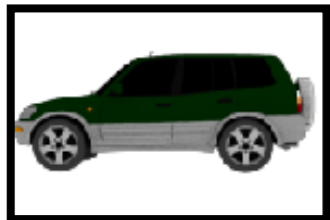
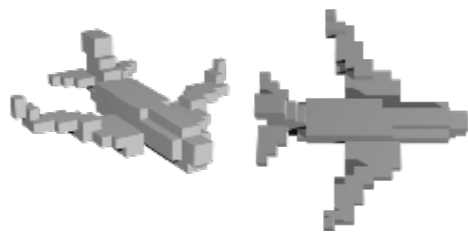
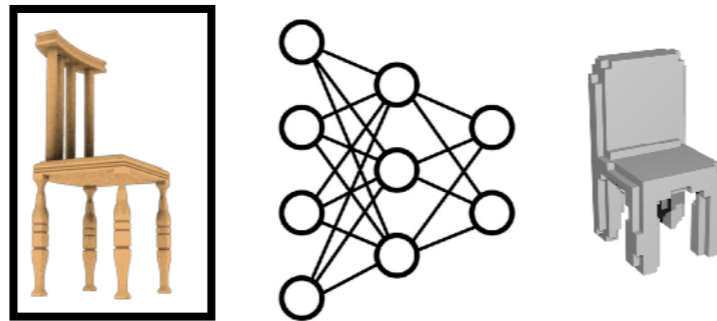
GT

DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision

# Experiments - ShapeNet



Input

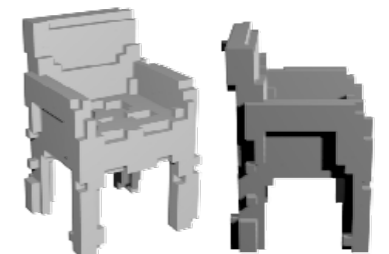
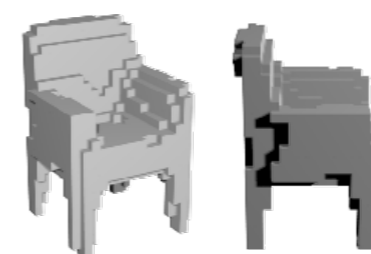
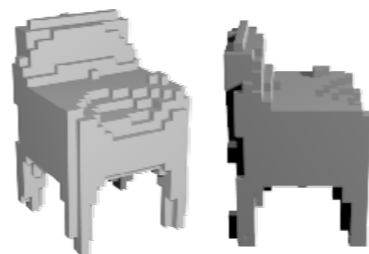
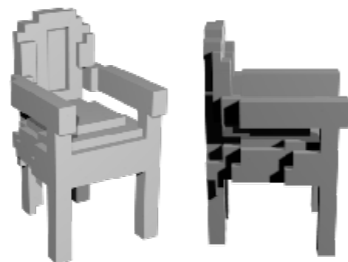
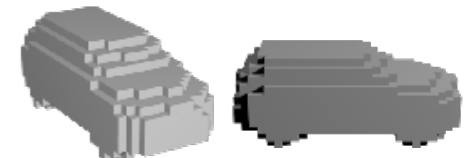
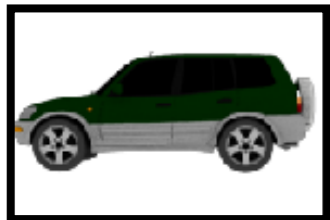
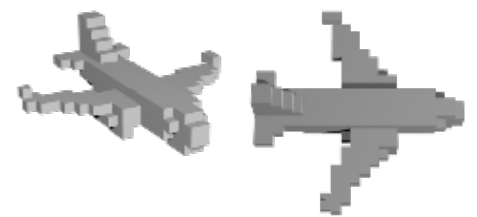
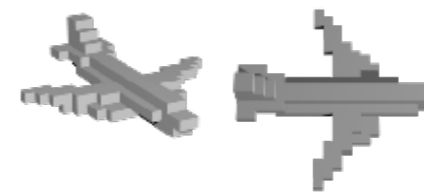
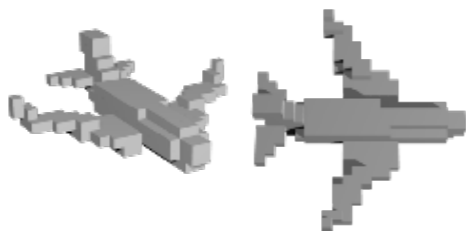
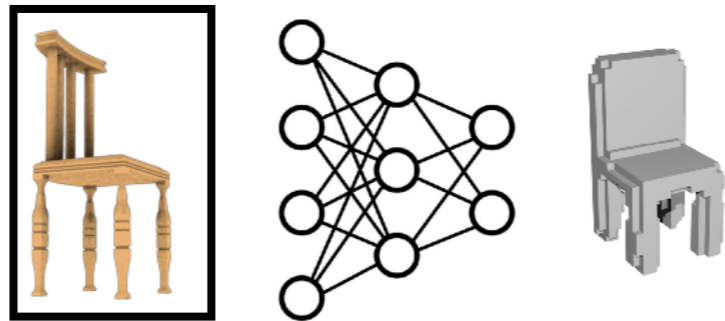
GT

DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision

# Experiments - ShapeNet



Input

GT

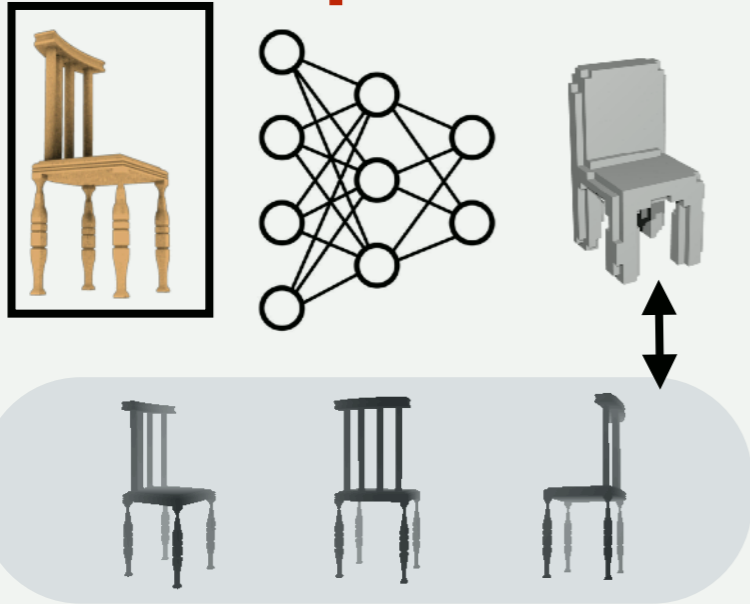
DRC  
(Mask)

DRC  
(Depth)

3D  
Supervision

# Learning Single-view Reconstruction

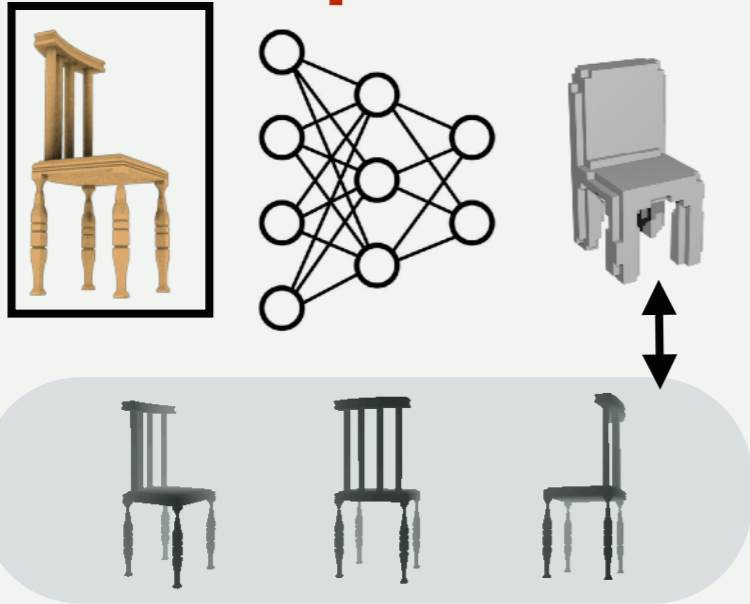
## ShapeNet



Supervision : Pose + Depth/Mask

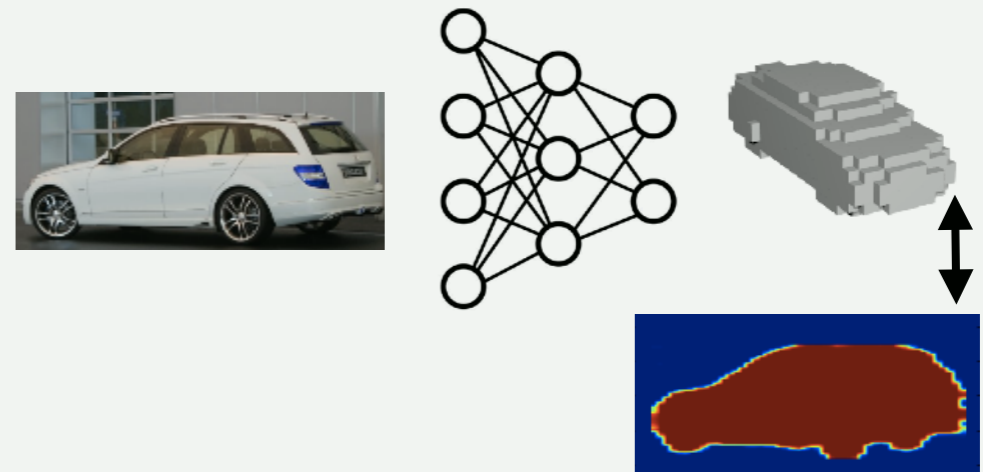
# Learning Single-view Reconstruction

## ShapeNet



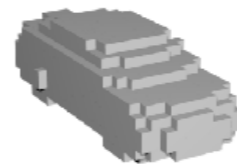
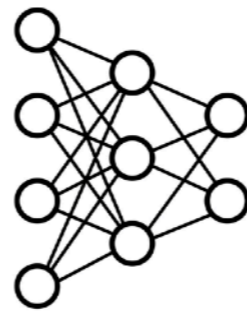
Supervision : Pose + Depth/Mask

## PASCAL VOC



Supervision : Pose + Mask

# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

DRC  
(Pascal)

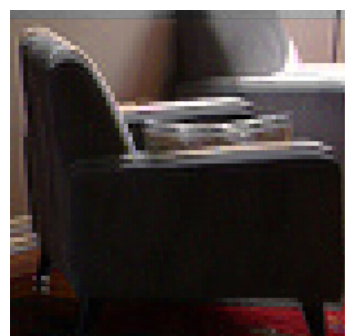
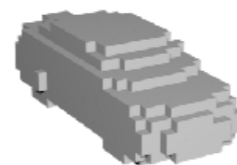
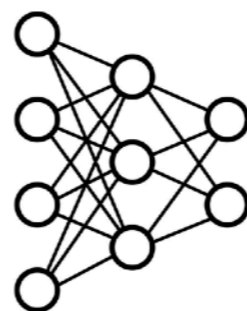
SNet 3D

DRC  
(Joint)

'Ground-  
Truth'



# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

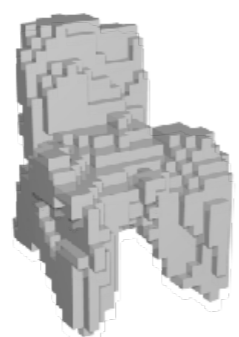
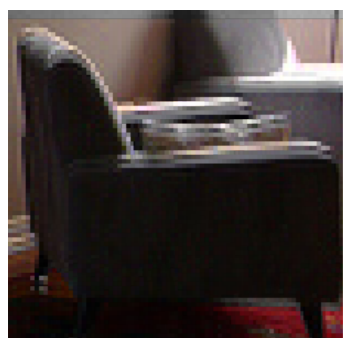
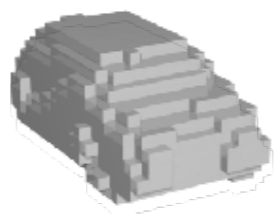
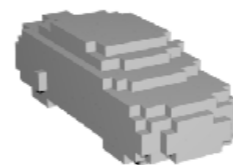
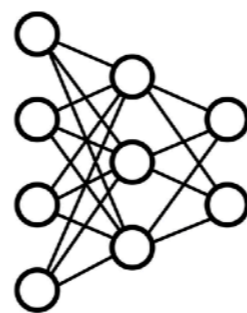
DRC  
(Pascal)

SNet 3D

DRC  
(Joint)

'Ground-  
Truth'

# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

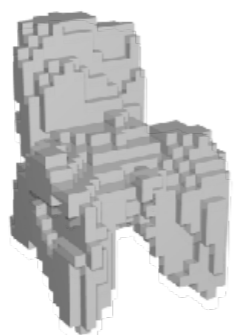
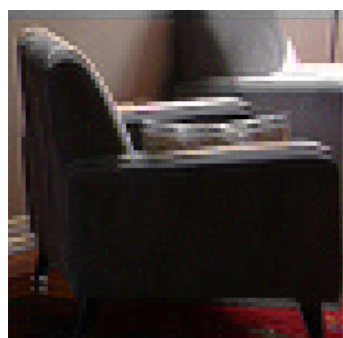
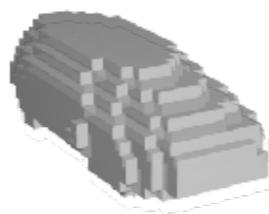
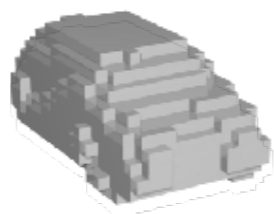
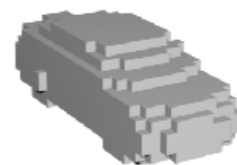
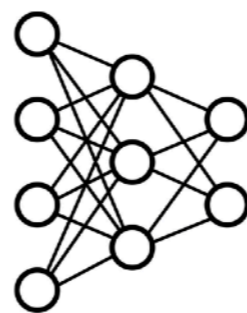
DRC  
(Pascal)

SNet 3D

DRC  
(Joint)

'Ground-  
Truth'

# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

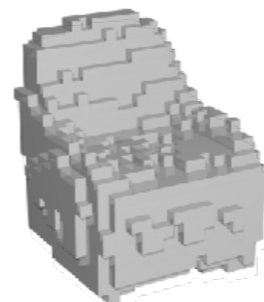
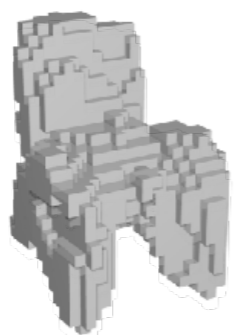
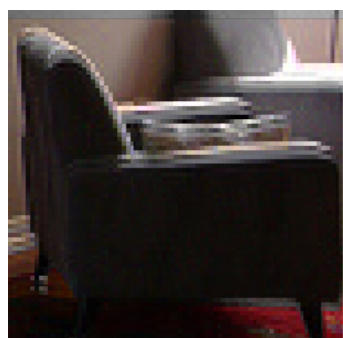
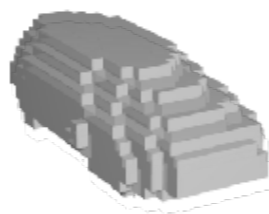
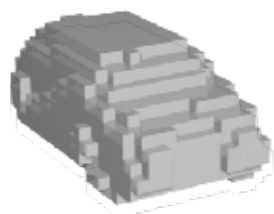
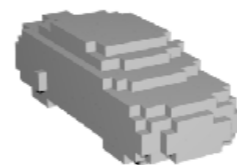
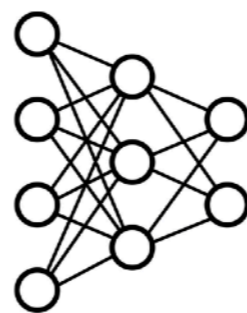
DRC  
(Pascal)

SNet 3D

DRC  
(Joint)

'Ground-  
Truth'

# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

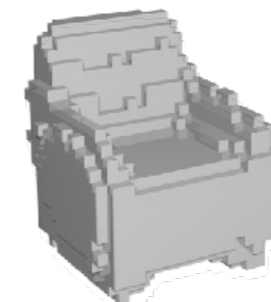
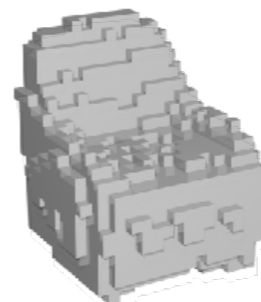
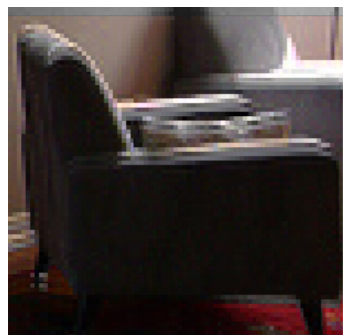
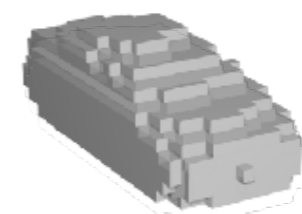
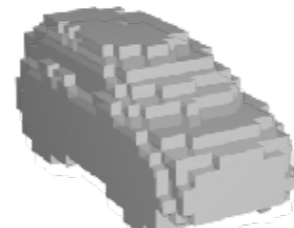
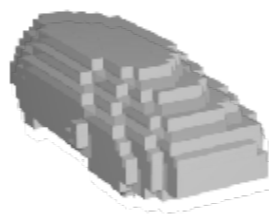
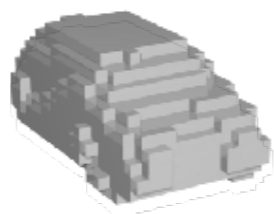
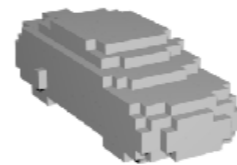
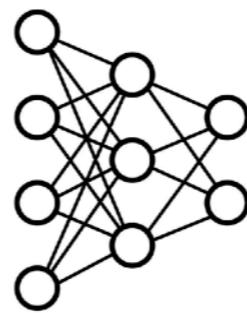
DRC  
(Pascal)

SNet 3D

DRC  
(Joint)

'Ground-  
Truth'

# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

DRC  
(Pascal)

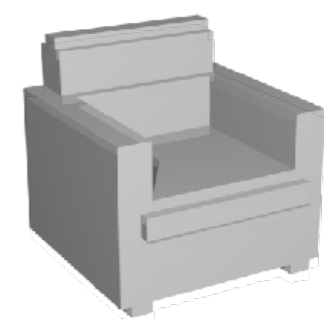
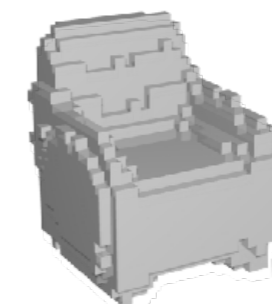
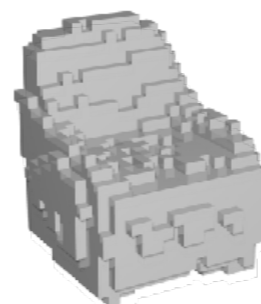
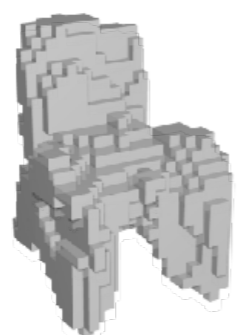
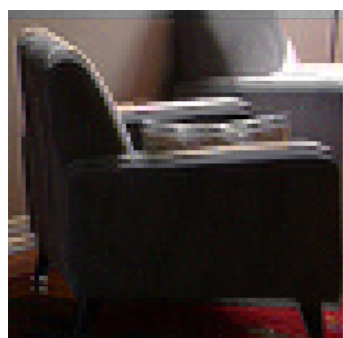
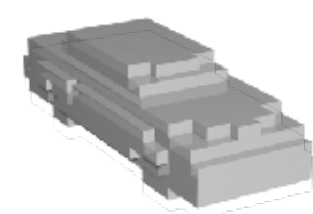
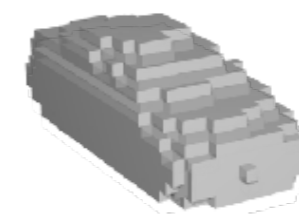
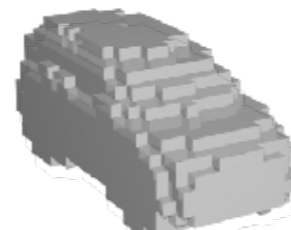
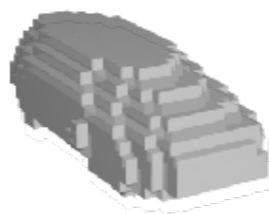
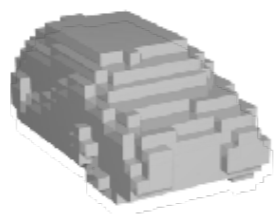
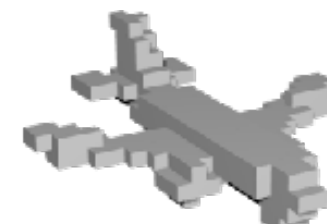
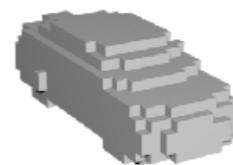
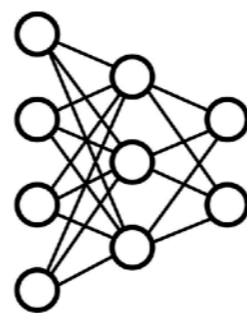
SNet 3D

DRC  
(Joint)

'Ground-  
Truth'



# Experiments - PASCAL VOC



Input

CSDM  
(Kar et. al.)

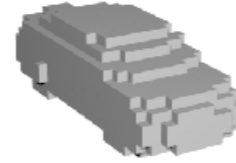
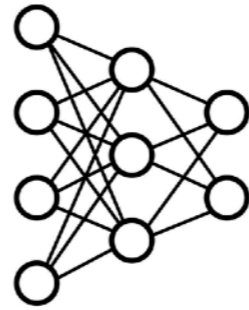
DRC  
(Pascal)

SNet 3D

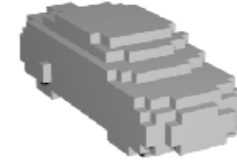
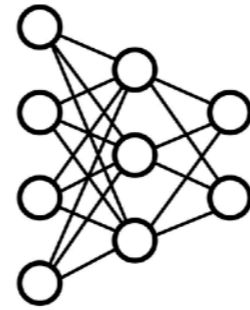
DRC  
(Joint)

'Ground-  
Truth'

# Experiments - PASCAL VOC



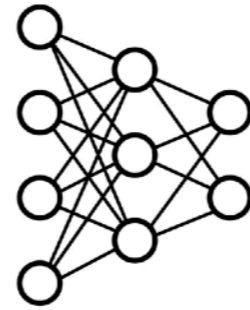
# Experiments - PASCAL VOC



Input



# Experiments - PASCAL VOC

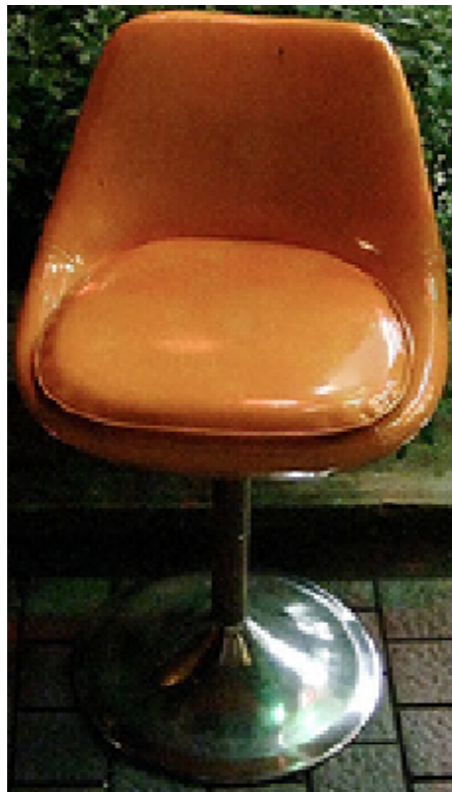
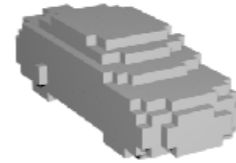
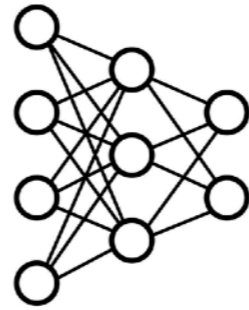


Input



Prediction

# Experiments - PASCAL VOC

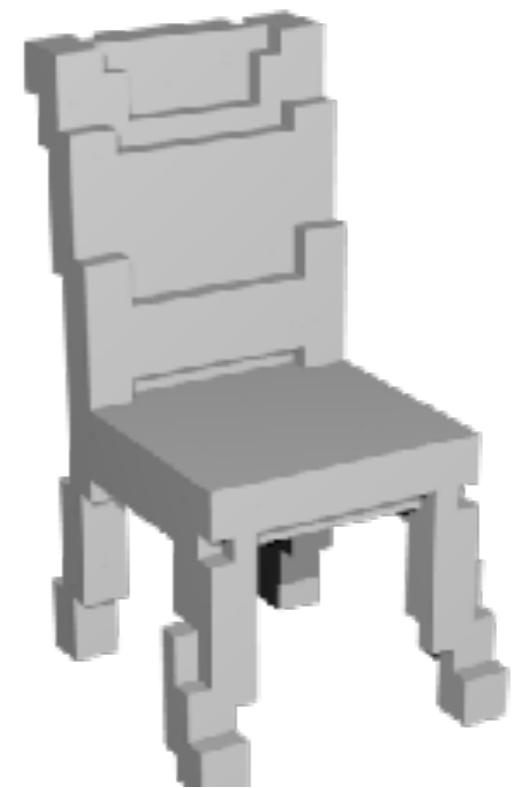
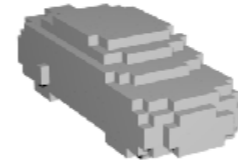
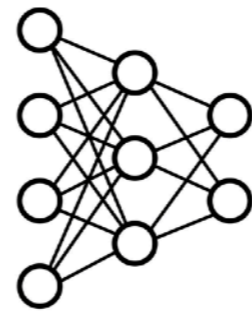


Input

Prediction

'Ground-truth'

# Experiments - PASCAL VOC



Input

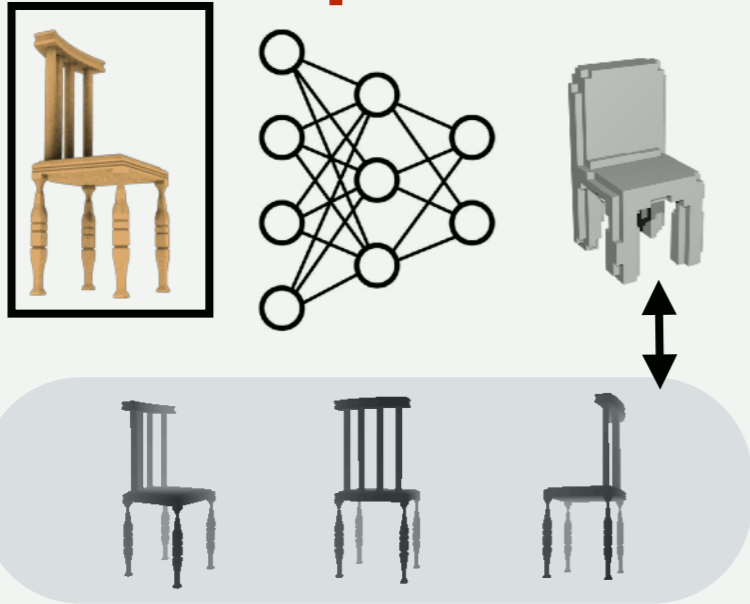
Prediction

'Ground-truth'

Collecting 'ground-truth' 3D is hard !

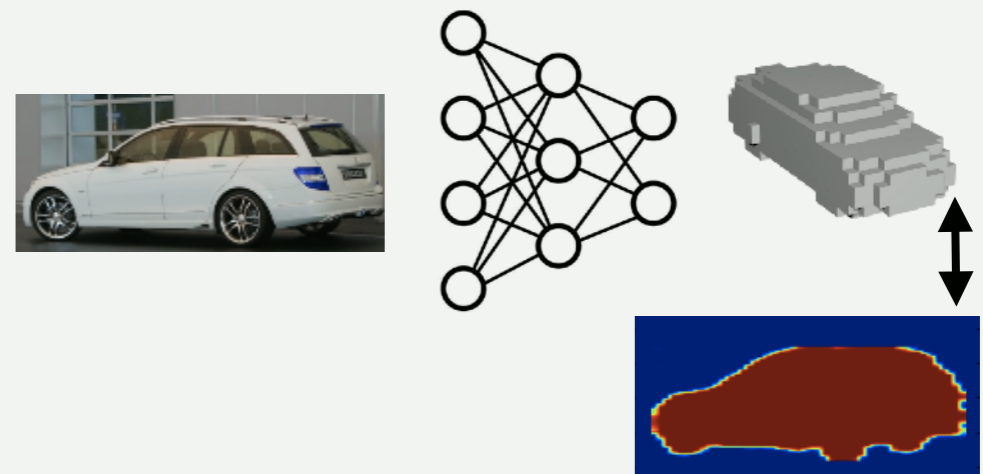
# Learning Single-view Reconstruction

## ShapeNet



Supervision : Pose + Depth/Mask

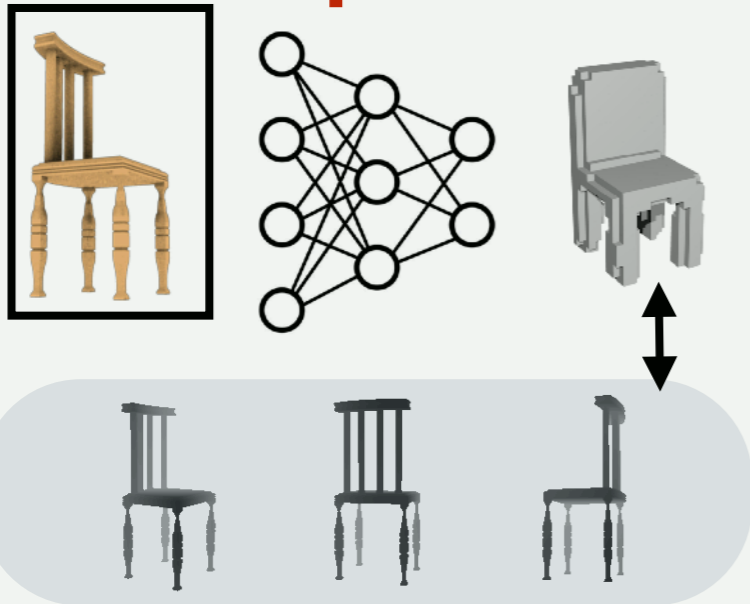
## PASCAL VOC



Supervision : Pose + Mask

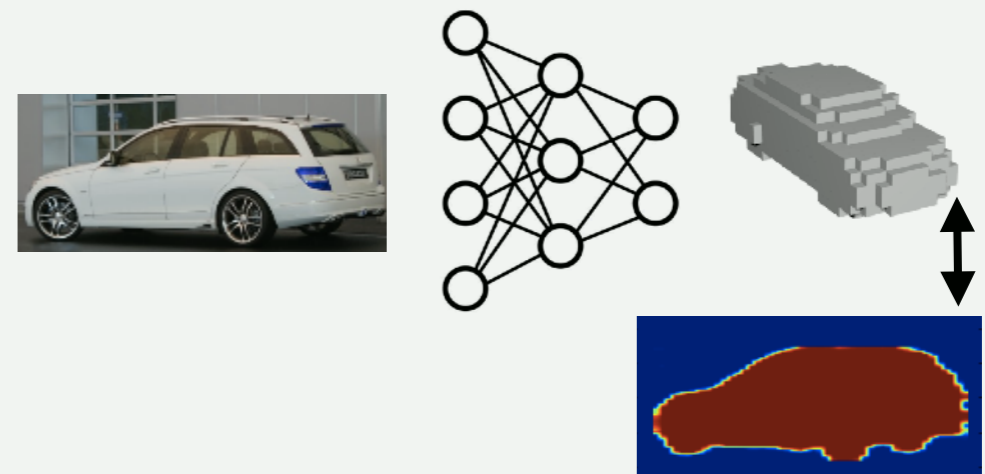
# Learning Single-view Reconstruction

## ShapeNet



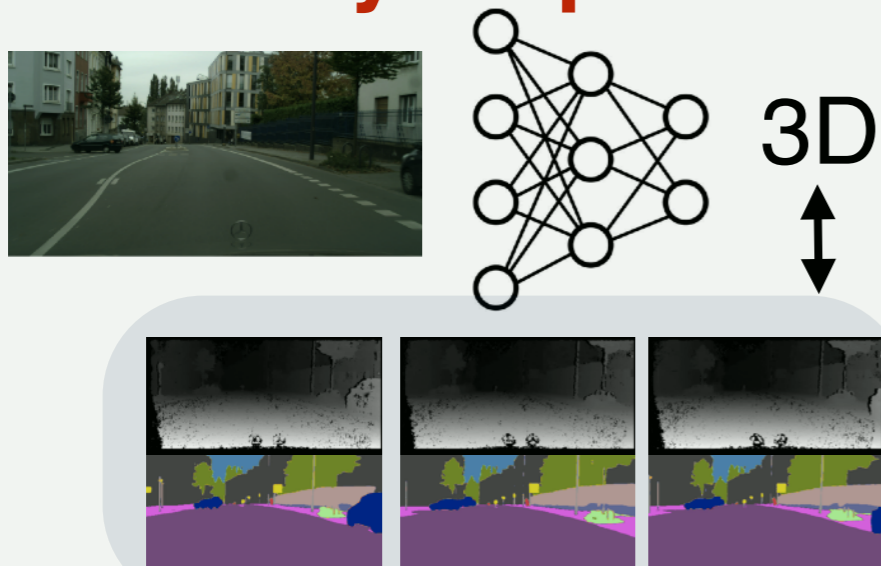
Supervision : Pose + Depth/Mask

## PASCAL VOC



Supervision : Pose + Mask

## Cityscapes

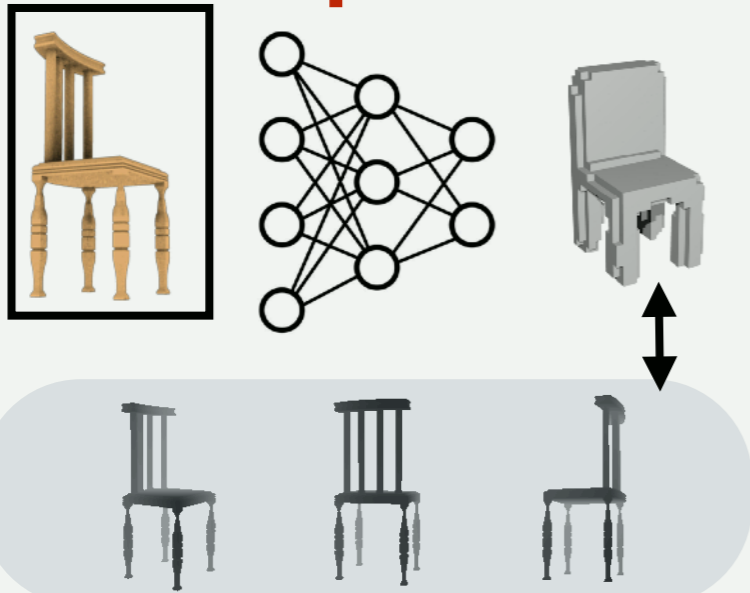


Supervision : Ego-motion,  
Depth, Semantics



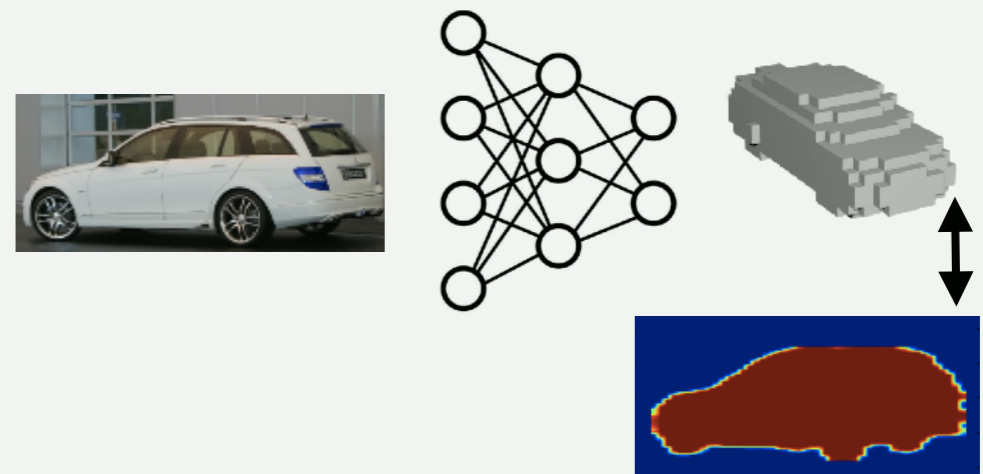
# Learning Single-view Reconstruction

## ShapeNet



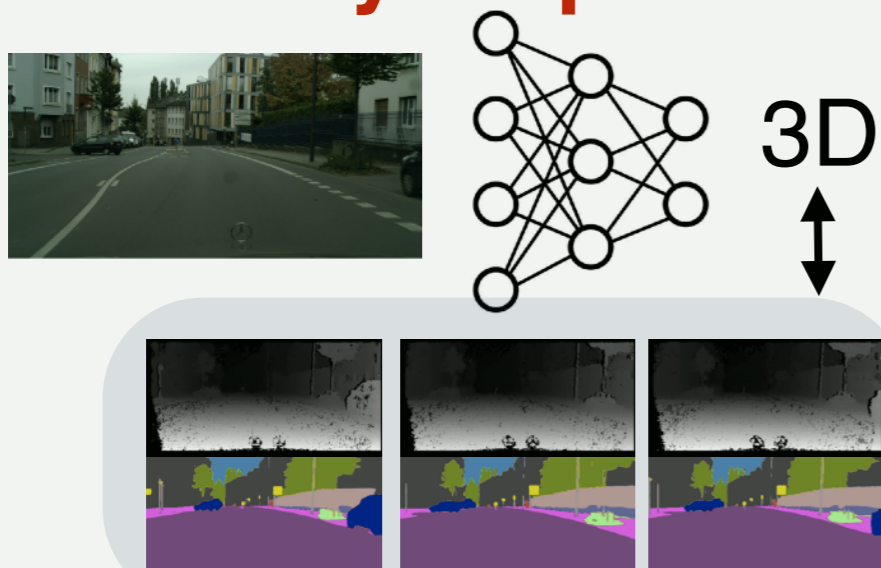
Supervision : Pose + Depth/Mask

## PASCAL VOC



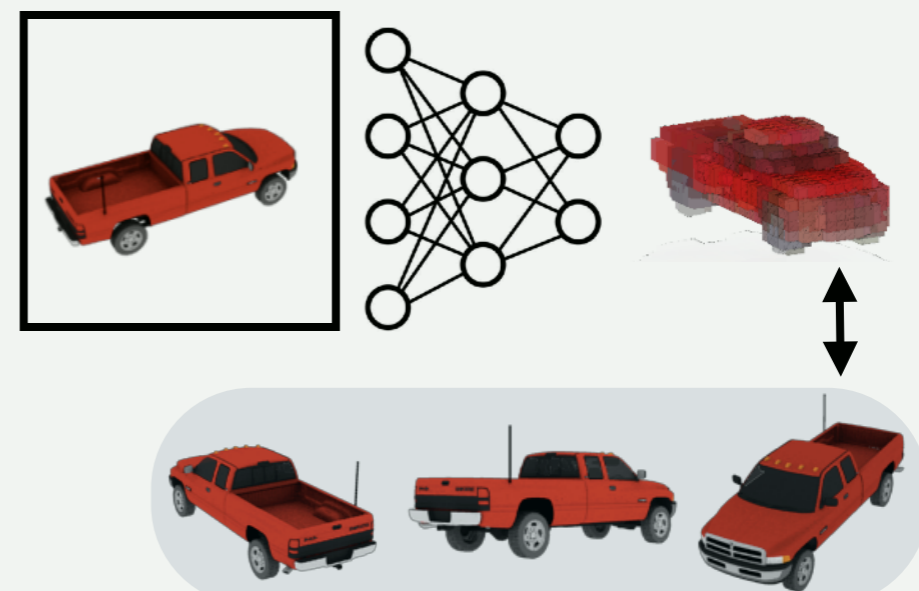
Supervision : Pose + Mask

## Cityscapes



Supervision : Ego-motion,  
Depth, Semantics

## ShapeNet (color supervised)

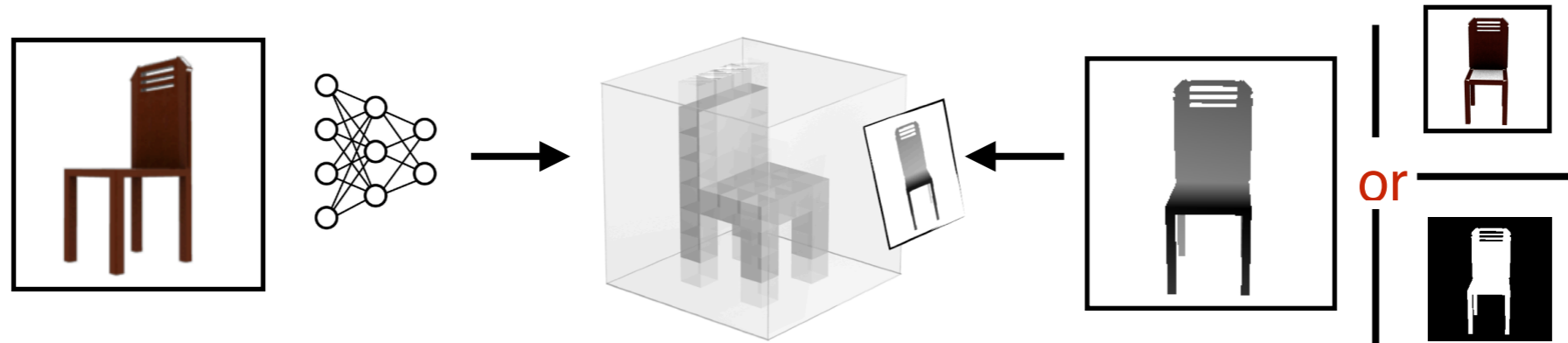


Supervision : Pose + RGB

# Conclusion

# Conclusion

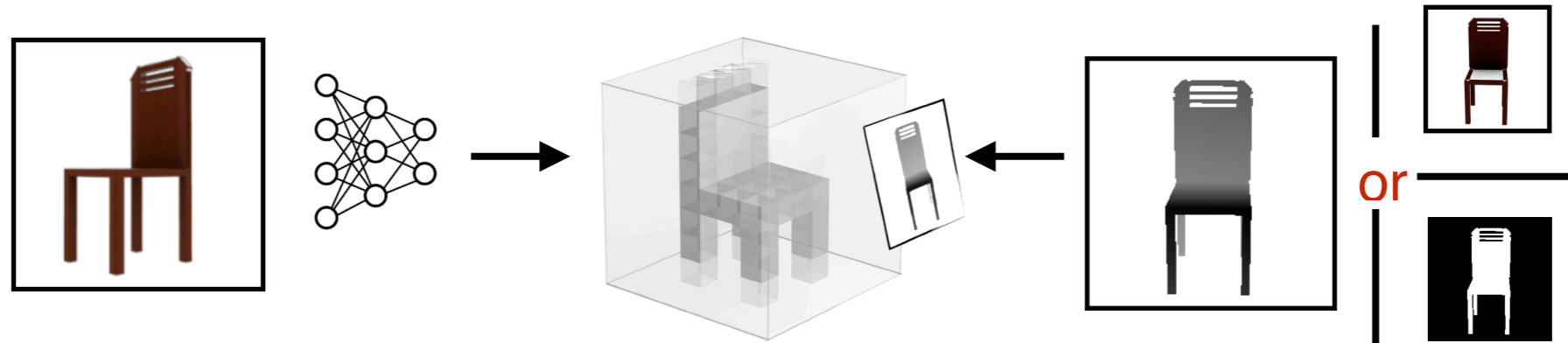
- Learning 3D via Geometric Consistency



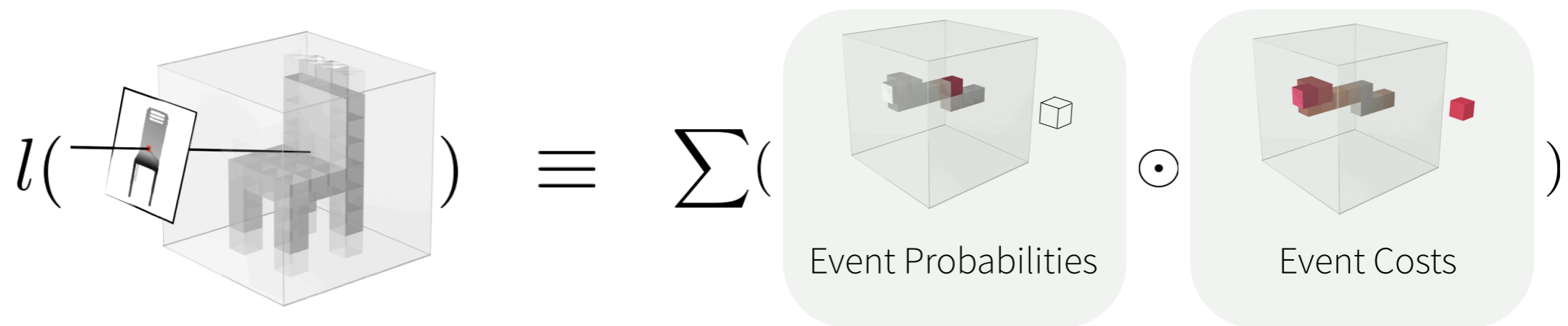


# Conclusion

- Learning 3D via Geometric Consistency

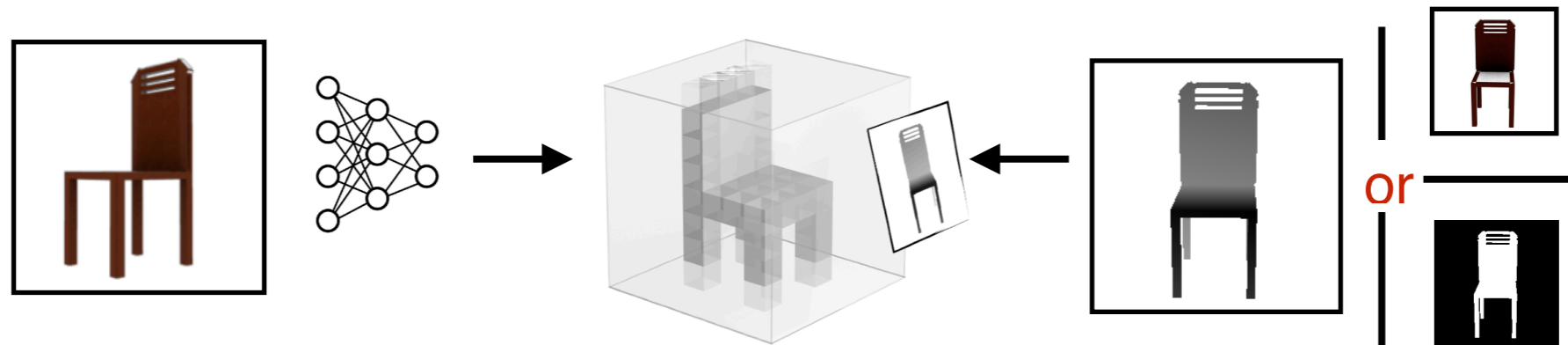


- Differentiable Ray Consistency Formulation

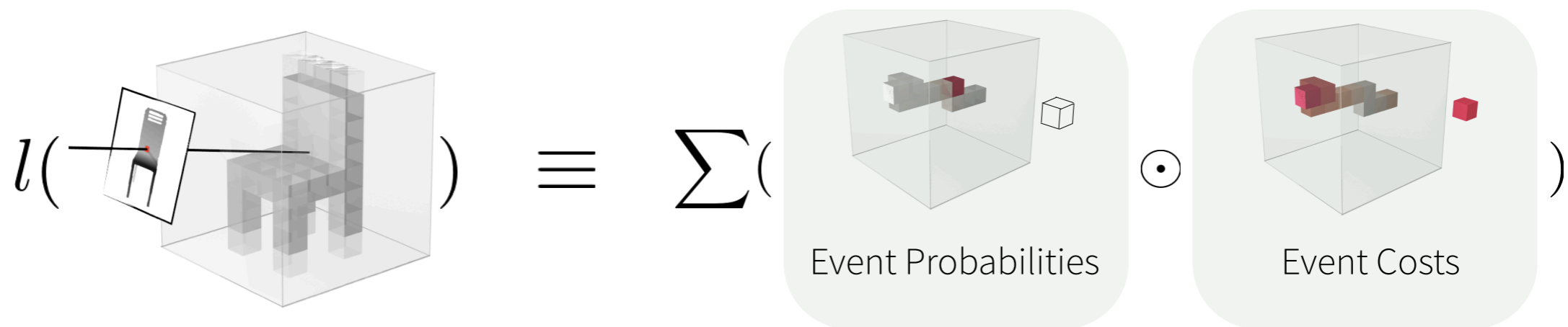


# Conclusion

- Learning 3D via Geometric Consistency



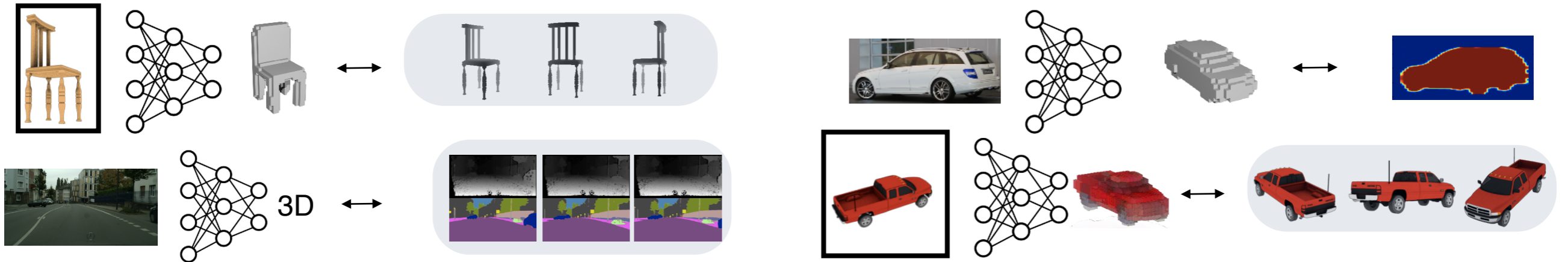
- Differentiable Ray Consistency Formulation



- Applications across scenarios



# Thank You



[Code : https://github.com/shubhtuls/drc](https://github.com/shubhtuls/drc)