

Seeing the unseen

Data-driven 3D Understanding from Single Images

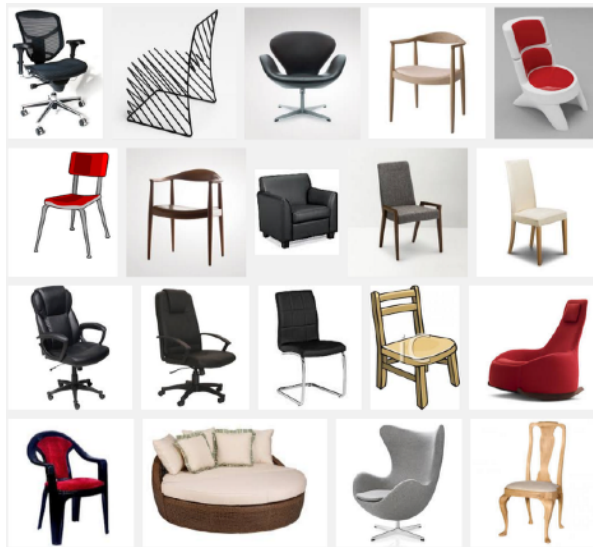
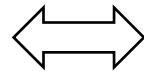


Image world



Shape world

Hao Su

Stanford
University

3D perception from a single image

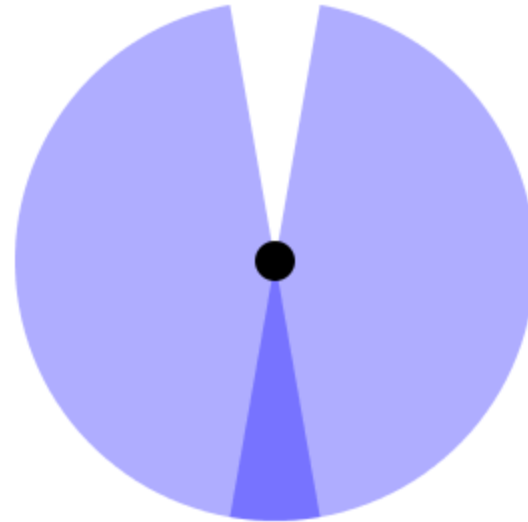


Monocular vision

a typical prey



Pigeon

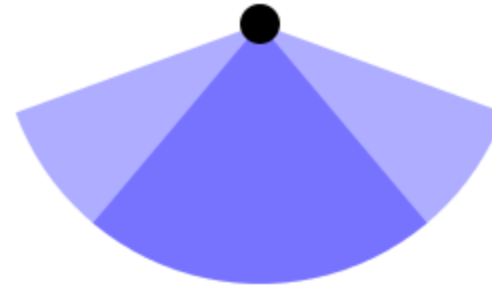


■ Binocular vision

a typical predator



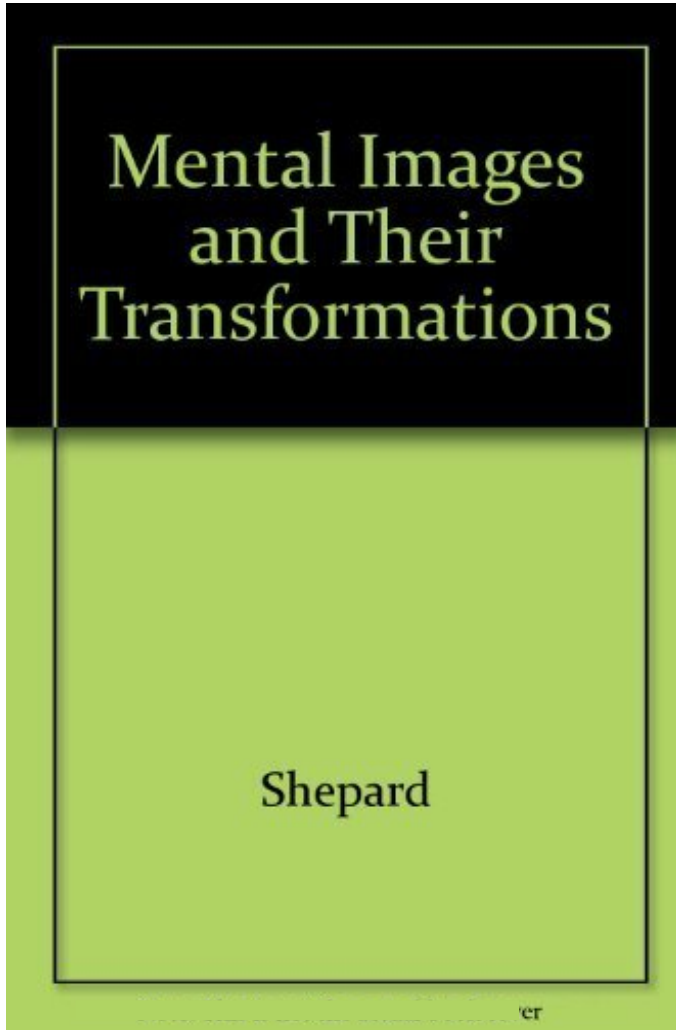
Owl



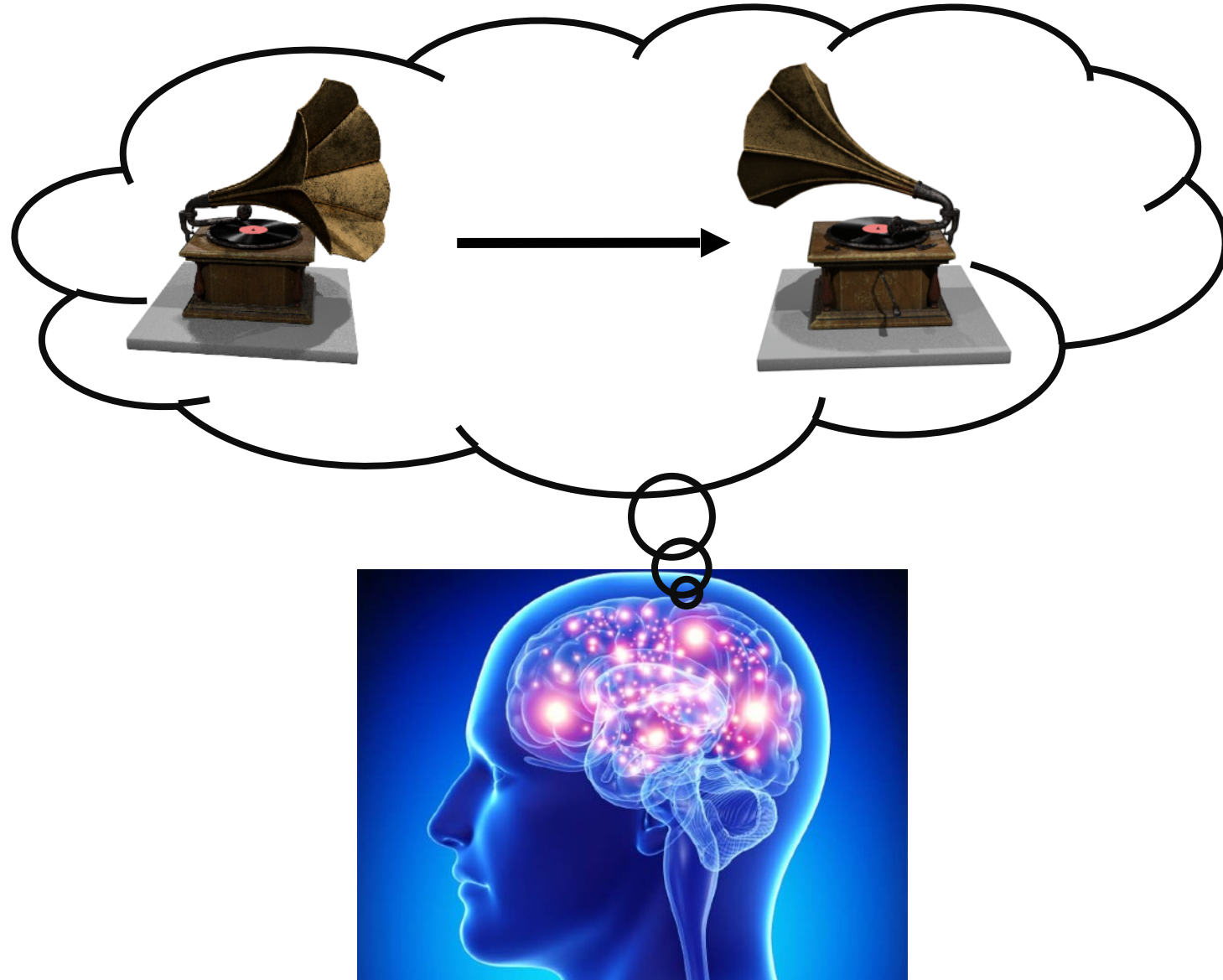
■ Monocular vision

Cited from https://en.wikipedia.org/wiki/Binocular_vision

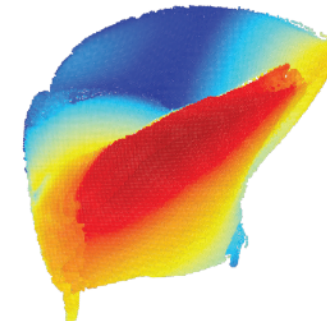
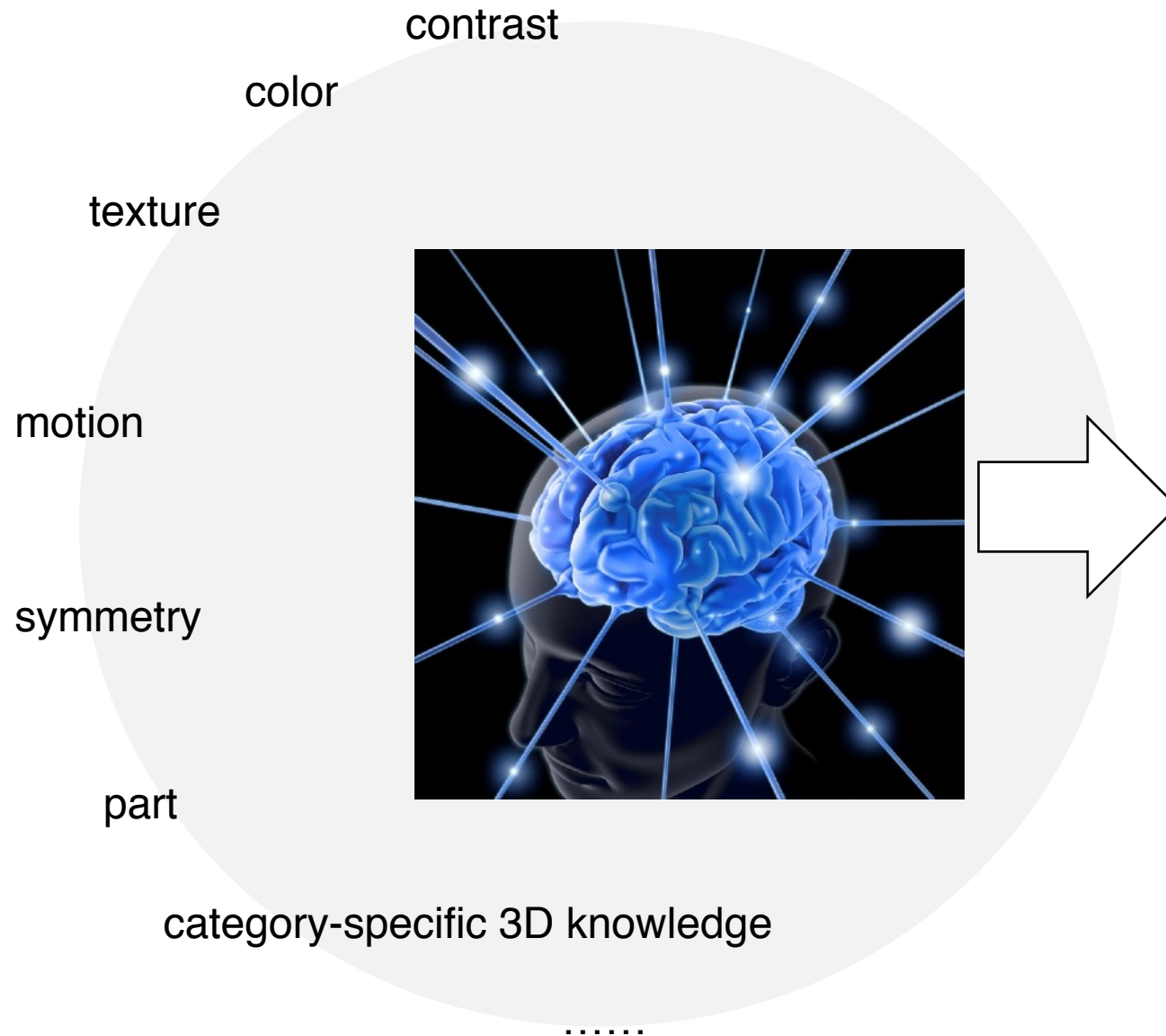
A psychological evidence – mental rotation



by Roger N. Shepard, National Science Medal Laureate
and Lynn Cooper, Professor at Columbia University

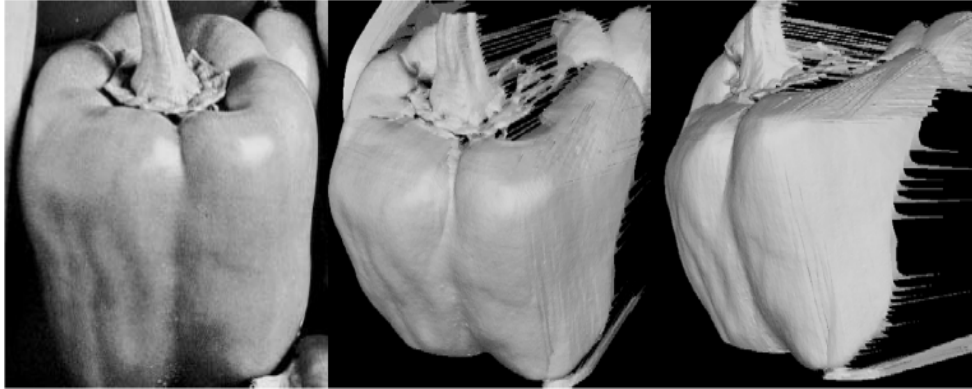


Visual cues are complicated

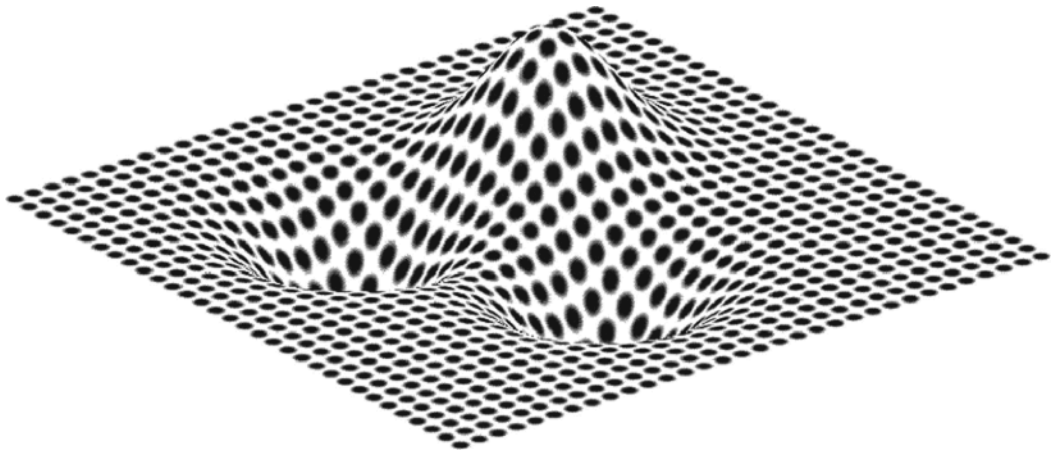


Status review of monocular vision algorithms

- Shape from X (texture, shading, ...)



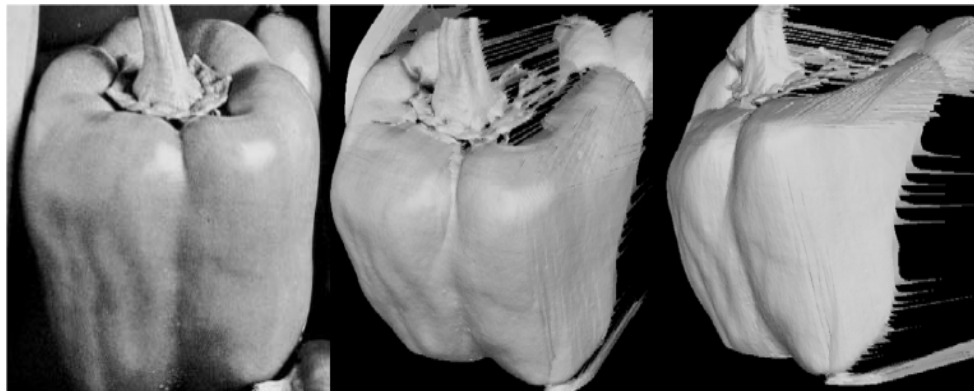
[Horn, 1989]



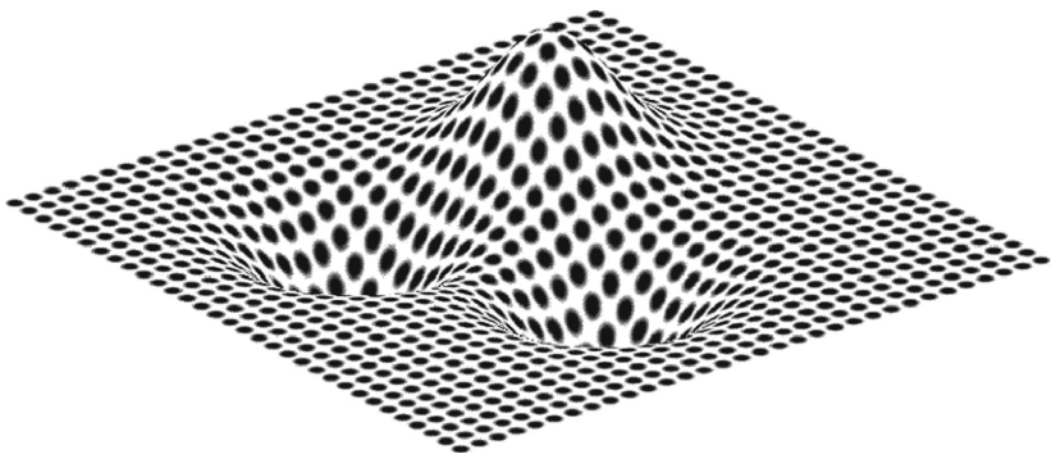
[Kender, 1979]

Status review of monocular vision algorithms

- Shape from X (texture, shading, ...)

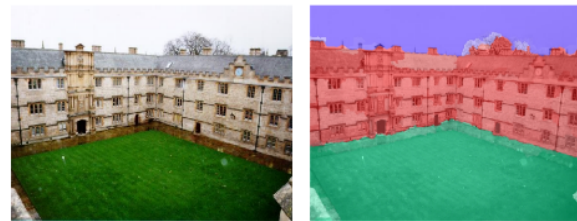


[Horn, 1989]



[Kender, 1979]

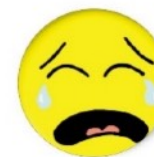
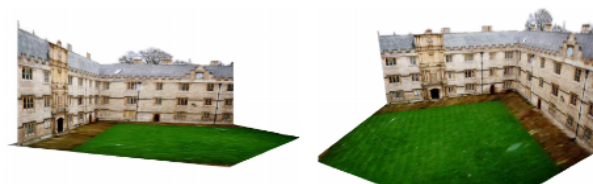
- Learning-based (from small data)



Hoiem et al, ICCV'05
Saxena et al, NIPS'05
...



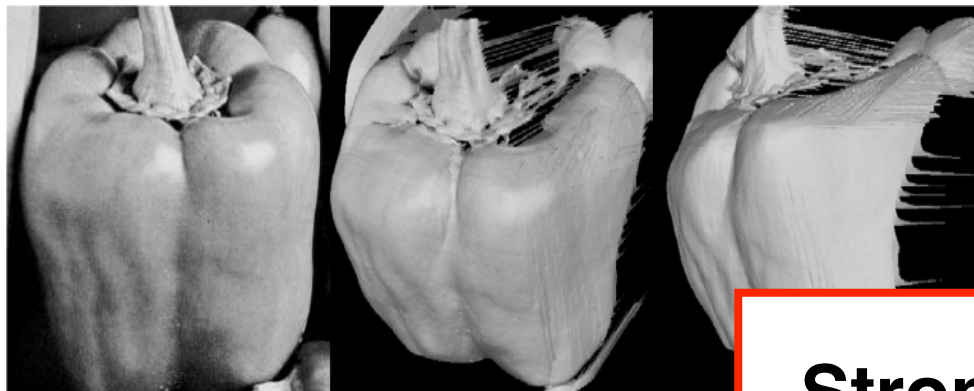
- large planes



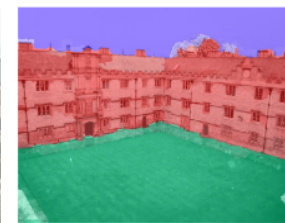
- fine structure
- topological variation
- ...

Status review of monocular vision algorithms

- Shape from X (texture, shading, ...)
- Learning-based (from small data)



[Horn, 1989]

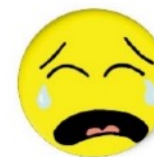


Hoiem et al, ICCV'05
Saxena et al, NIPS'05
...

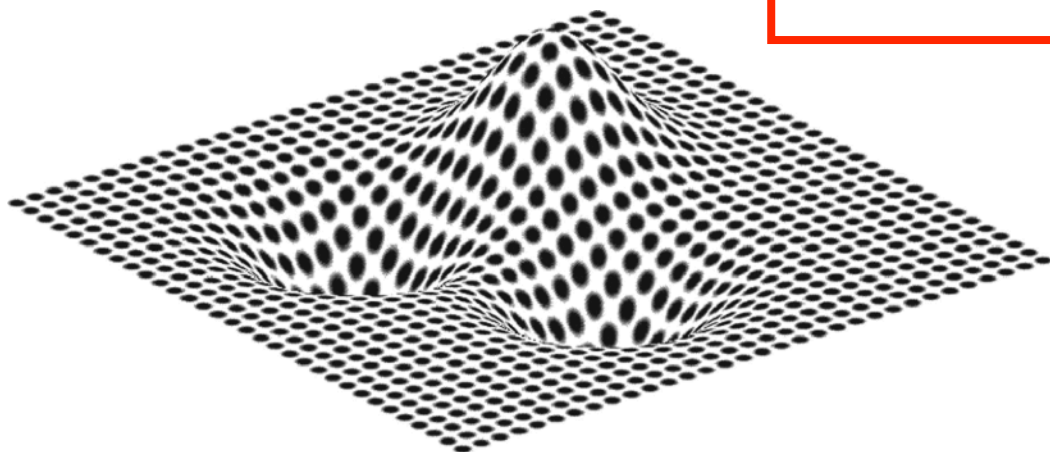


- large planes

Strong assumption
Not robust



- fine structure
- topological variation
- ...

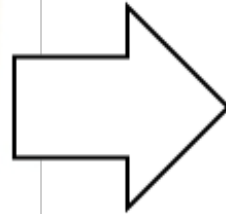


[Kender, 1979]

Data-driven 2D-3D lifting

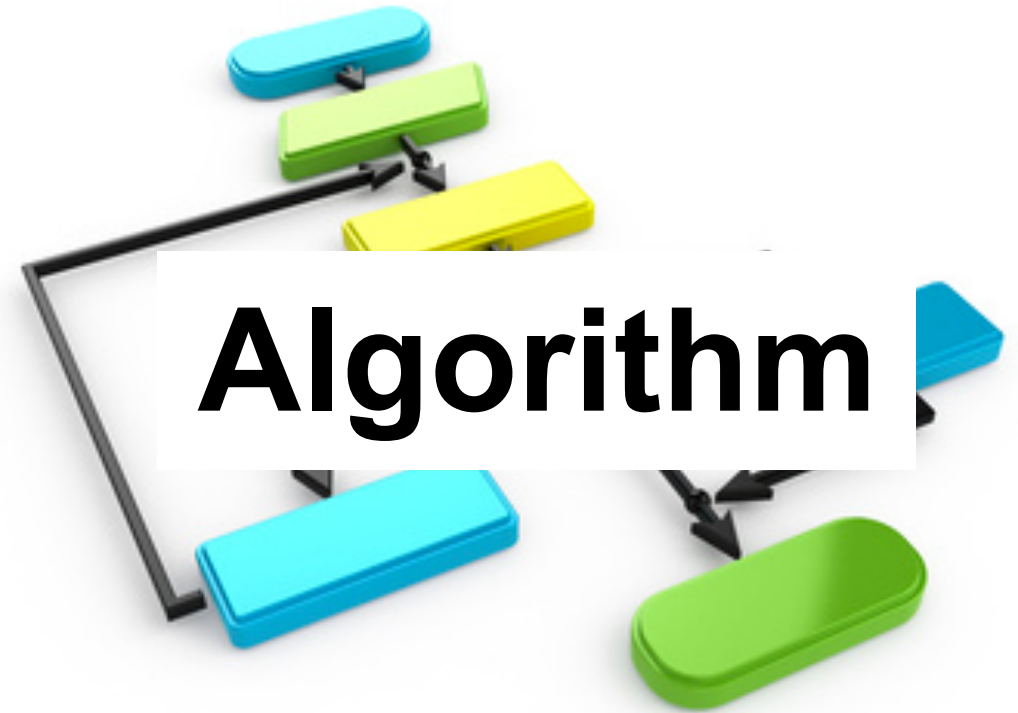


Many 3D objects



A priori knowledge of
the 3D world

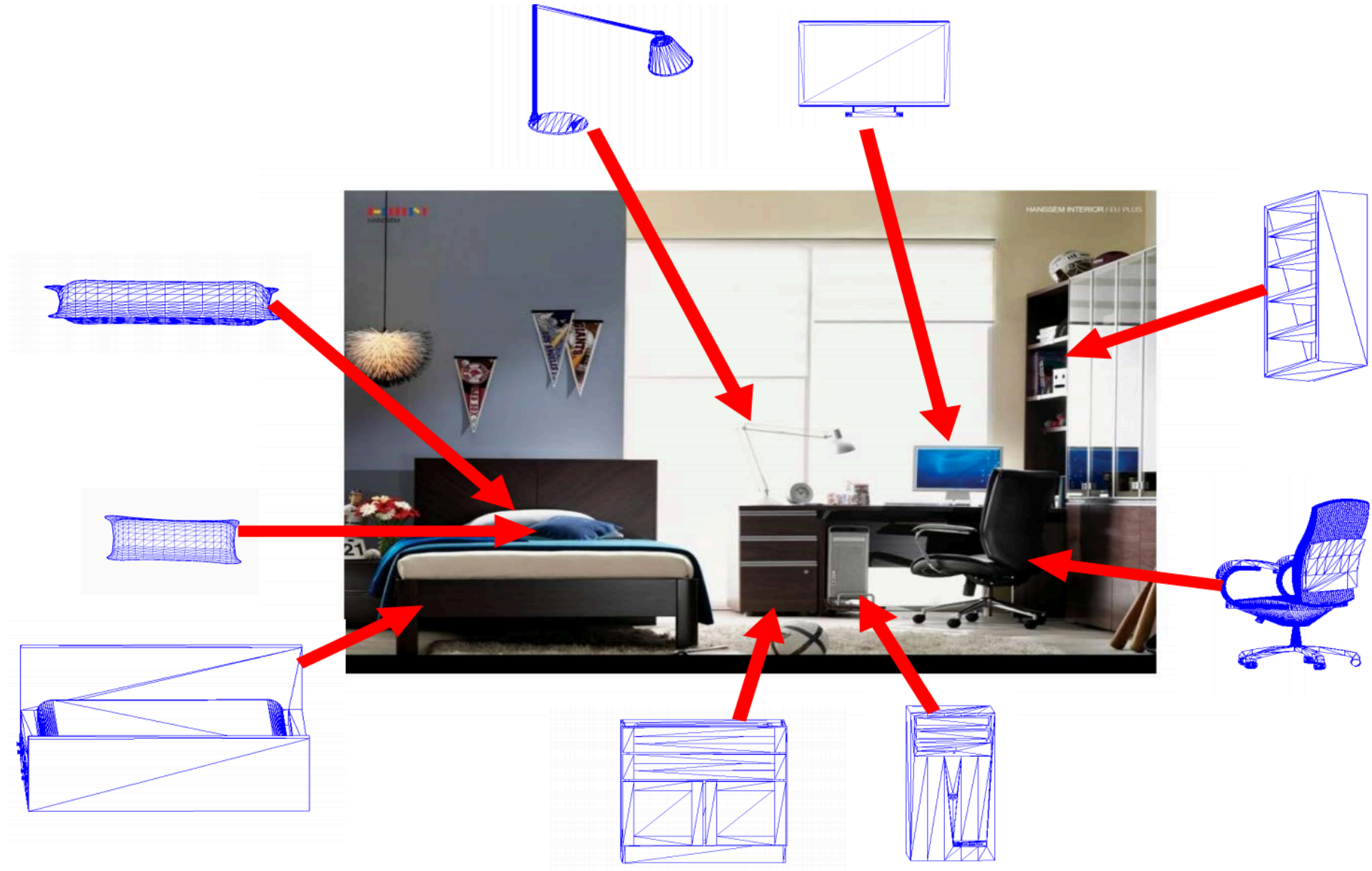
Two pillars for learning-based 3D understanding



Outline

- Motivation
- **Data**
- Algorithms for 2D-3D lifting
- Algorithm for 3D representation learning

Need many 2D images with 3D labels

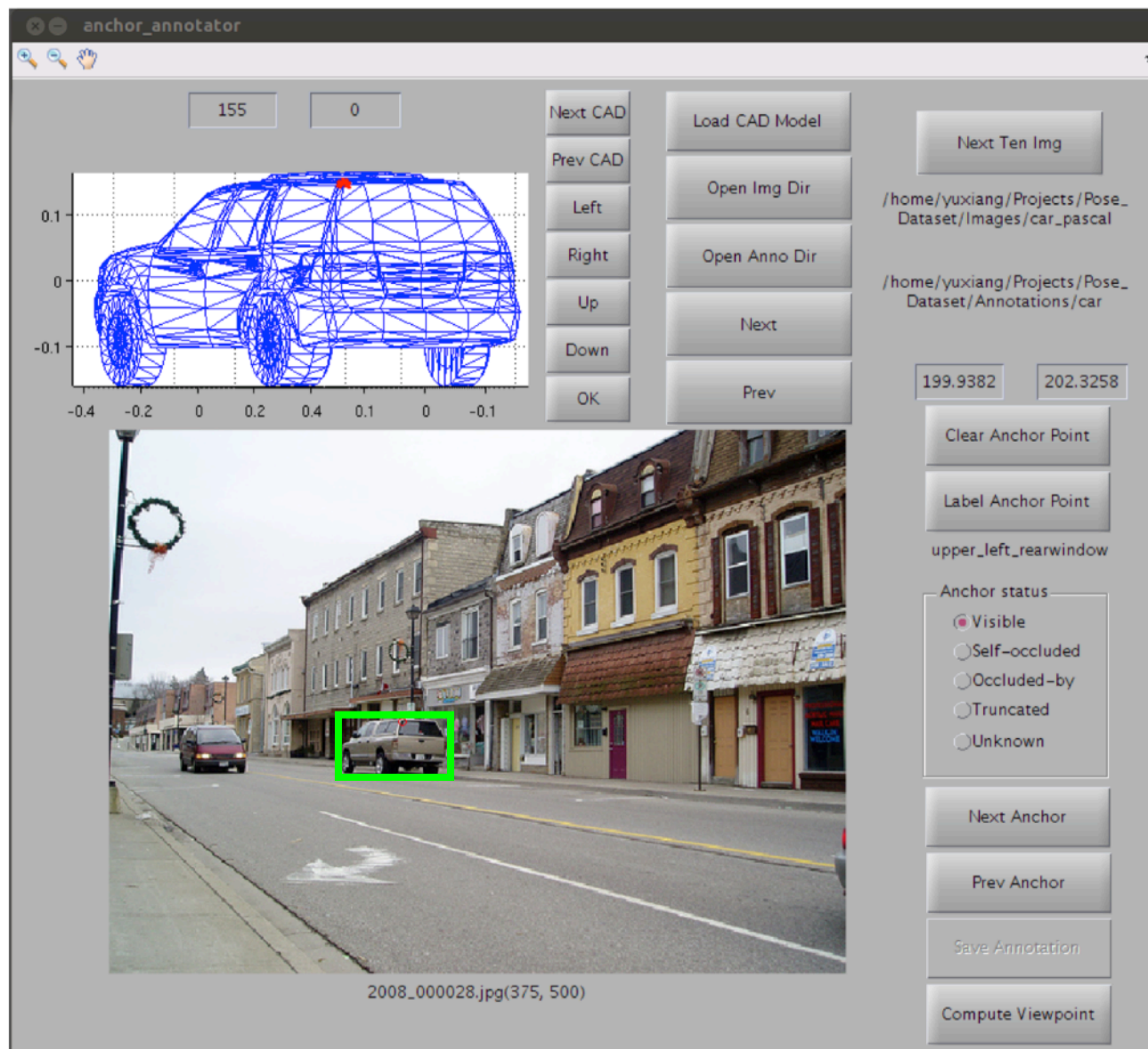


Accurate 3D label annotation is expensive

Example: viewpoint estimation

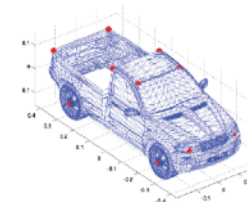
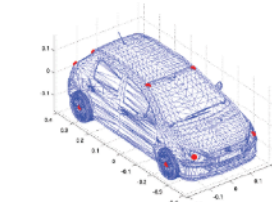
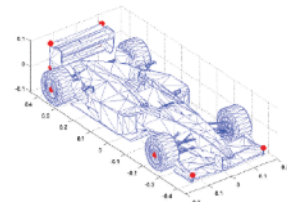
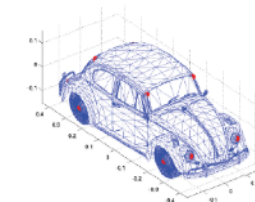
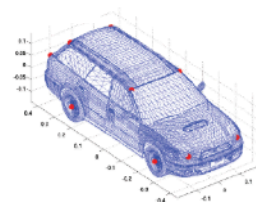
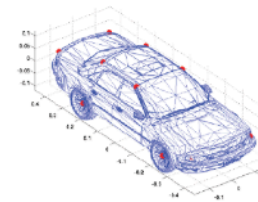
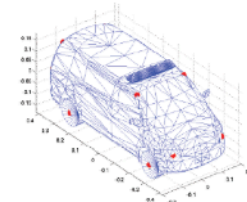
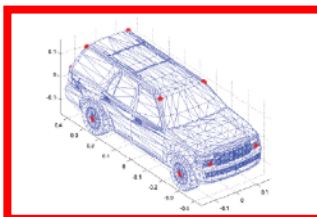
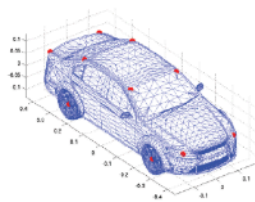
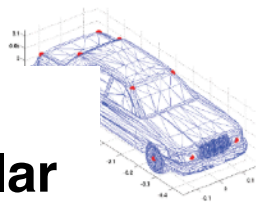


PASCAL3D+ dataset [Xiang et al.]



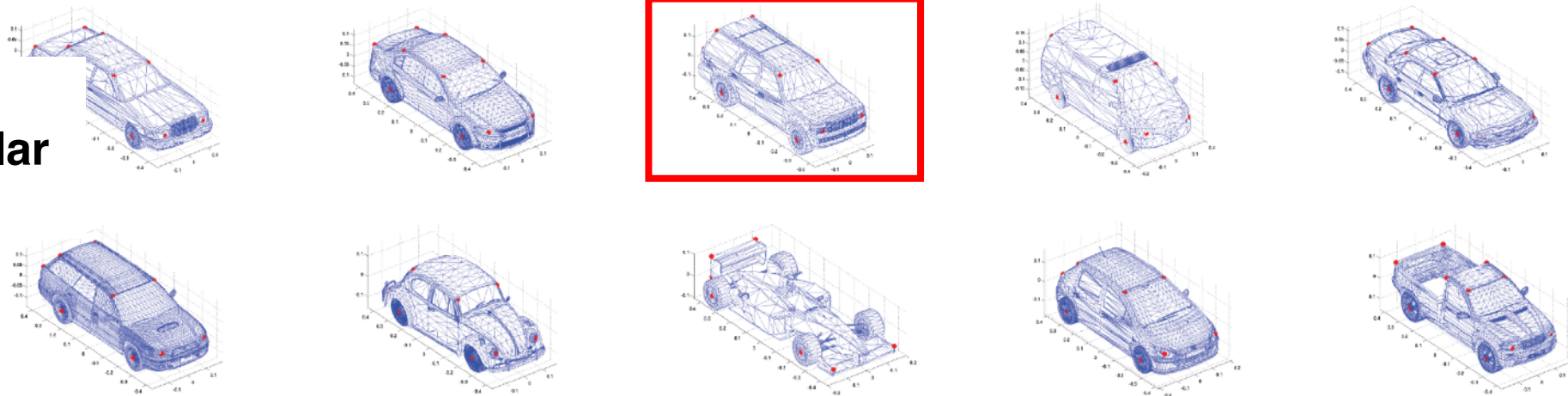
Accurate 3D label annotation is expensive

Step1:
Choose similar
model

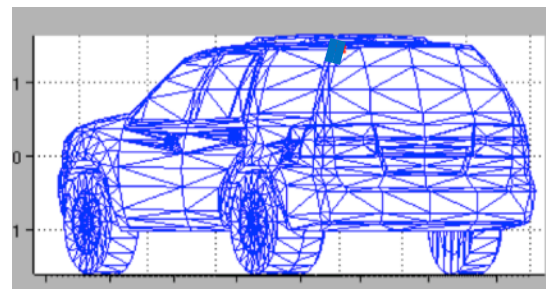


Accurate 3D label annotation is expensive

Step1:
Choose similar
model

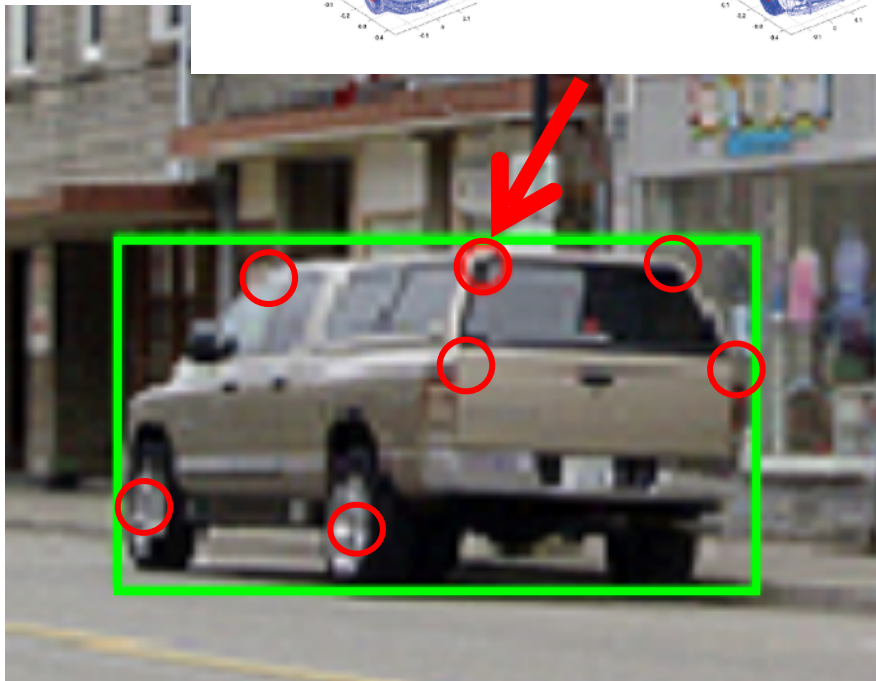
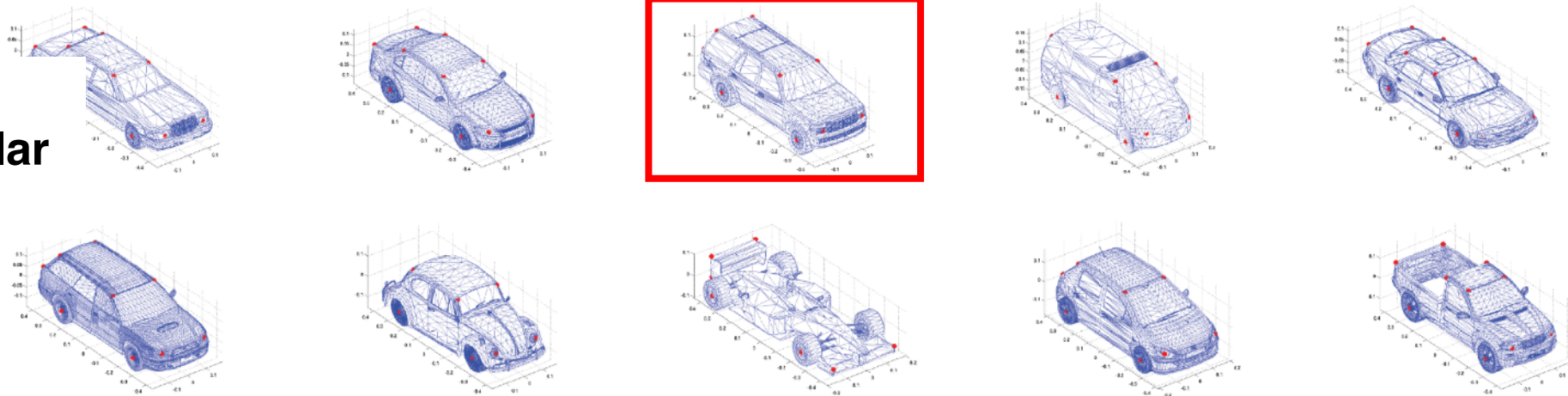


Step2:
Coarse Viewpoint
Labeling

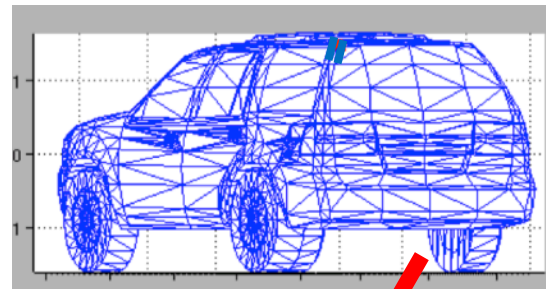


Accurate 3D label annotation is expensive

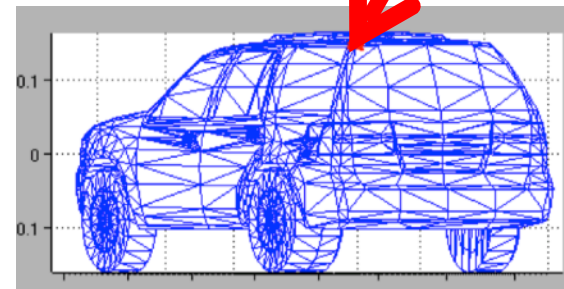
Step1:
Choose similar
model



Step2:
Coarse Viewpoint
Labeling



Step3:
Label keypoints
For alignment



Annotation takes
~1 min per object

High-cost Label Acquisition

Mismatch!

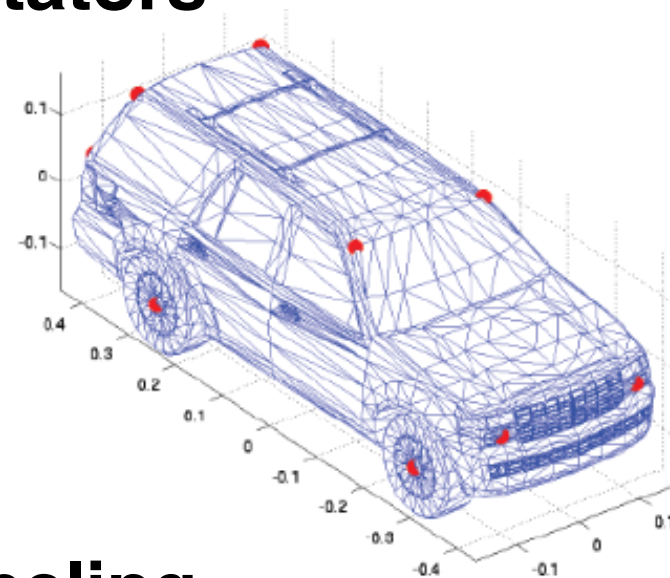
High-capacity Model

Needs millions of images to train

How to get MORE images with ACCURATE labels?



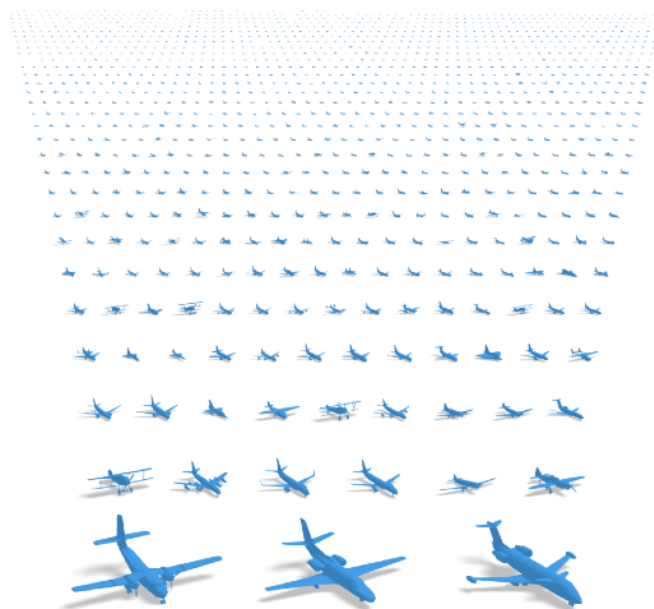
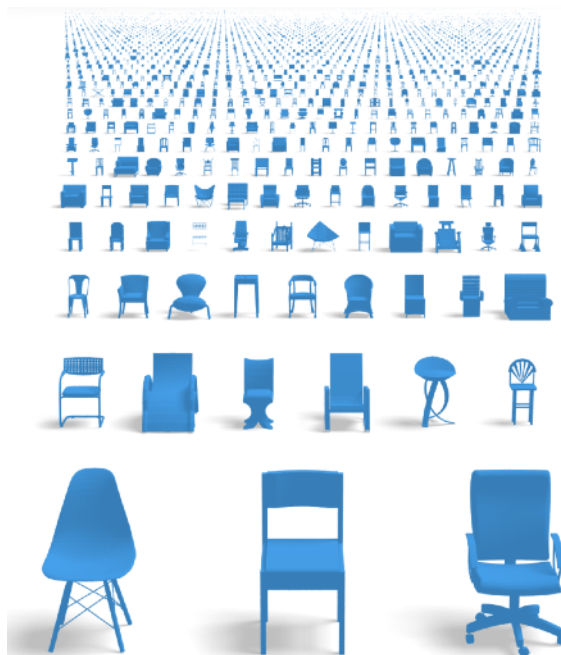
Manual labeling
by annotators



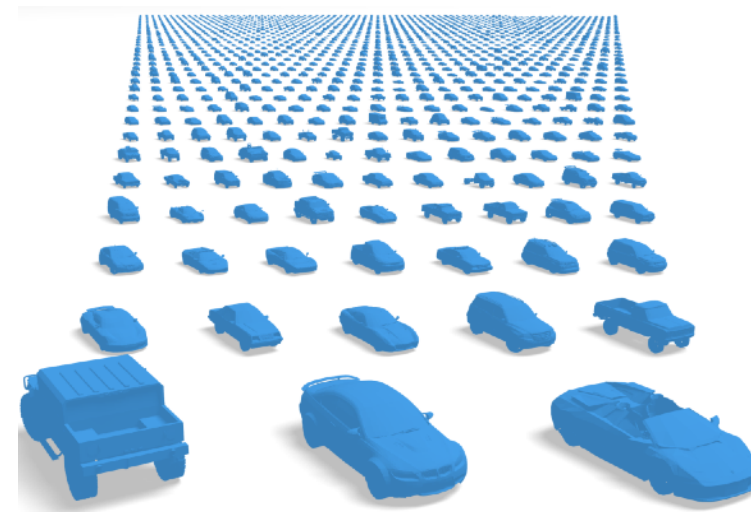
Auto labeling
through rendering

Good News: ShapeNet

<http://shapenet.cs.stanford.edu>

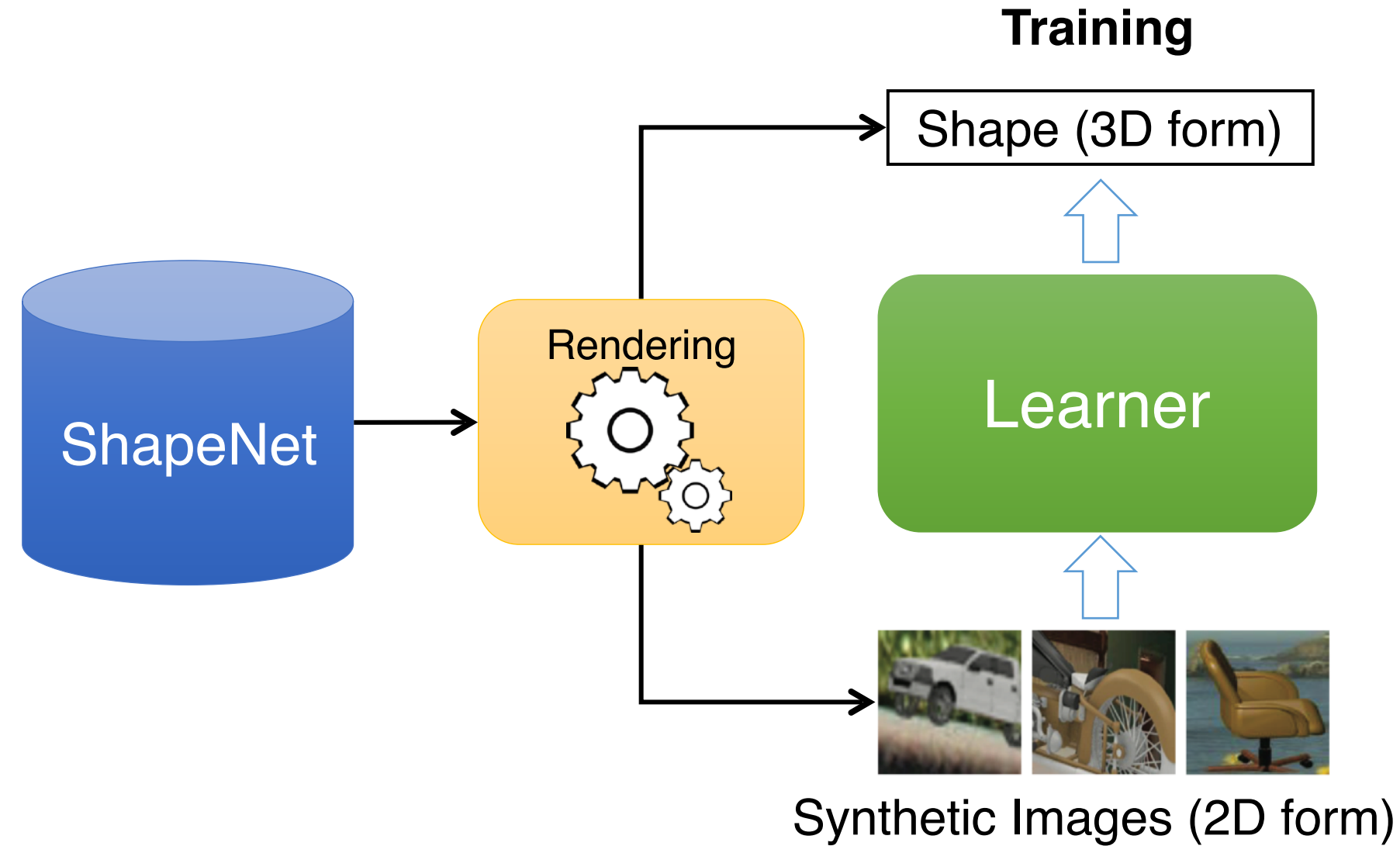


...

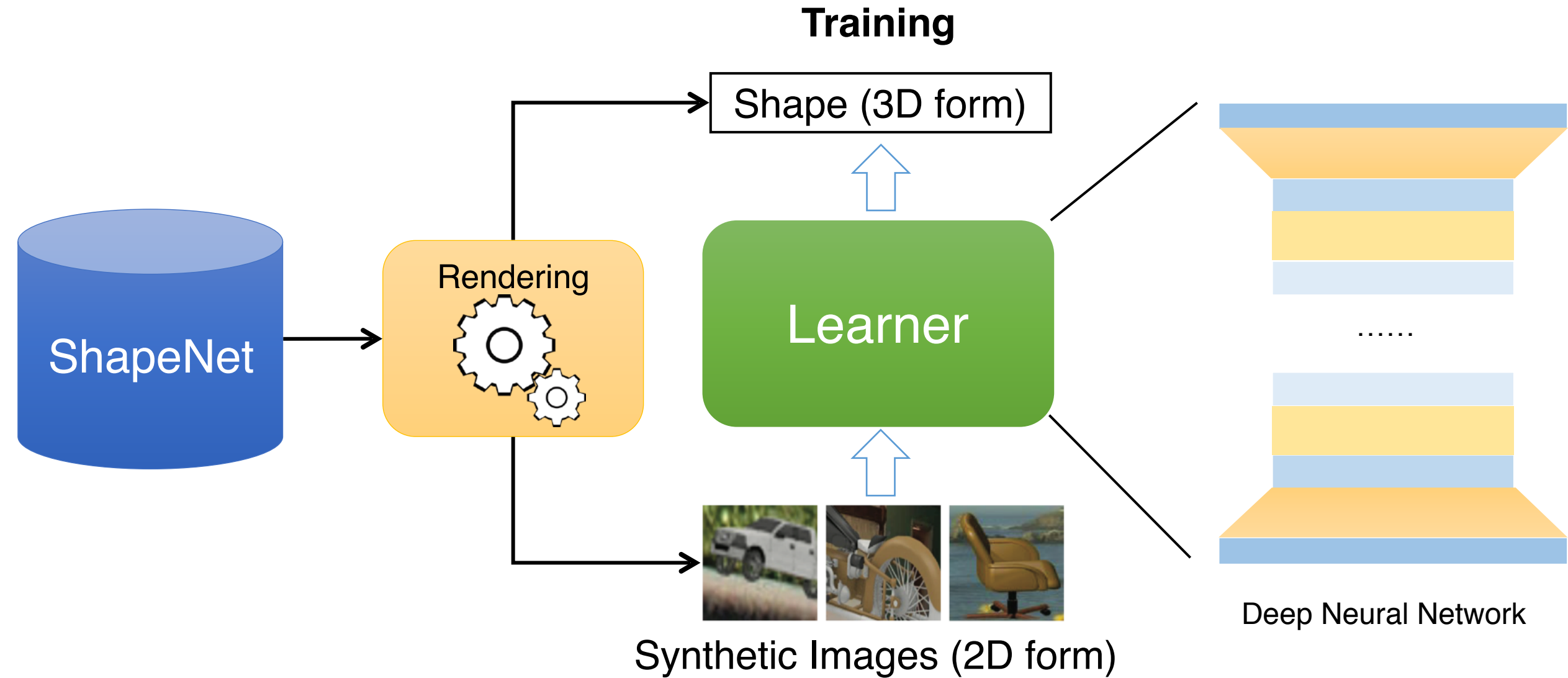


millions of shapes with **geometric** and **physical** annotations
(*ongoing efforts*)

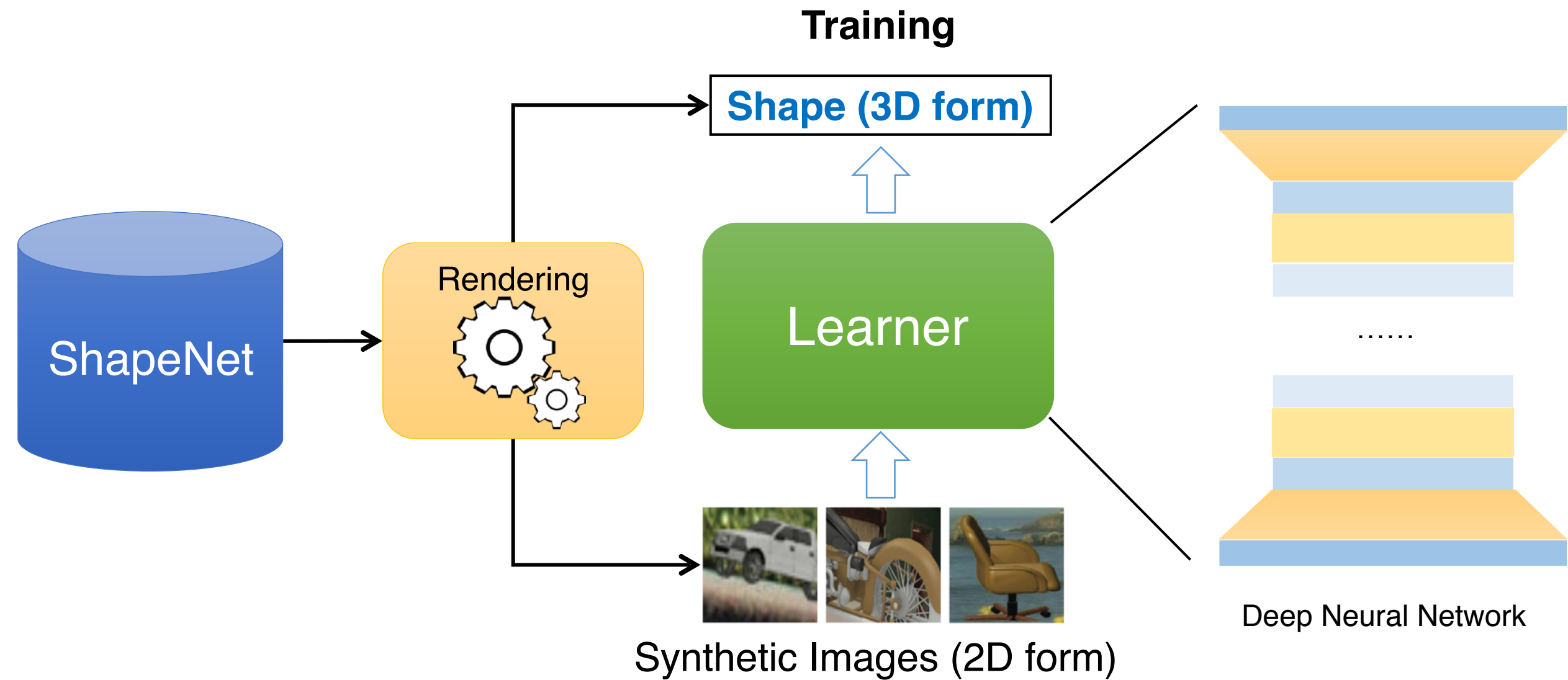
Synthesize for learning



Synthesize for learning



How shall we represent 3D shapes?



Many 3D representations are available

Candidates:

multi-view images

depth map

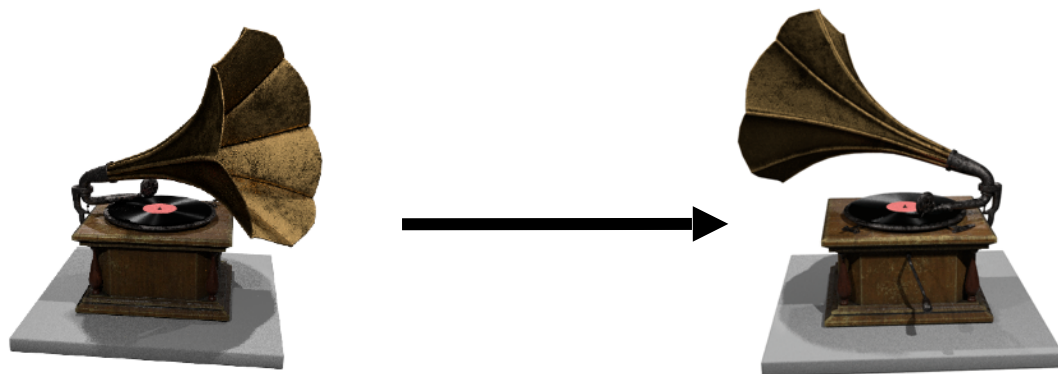
volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



Novel view image synthesis

[Su et al., ICCV15]

[Dosovitskiy et al., ECCV16]

Candidates:

multi-view images

depth map

volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



Candidates:

multi-view images

depth map

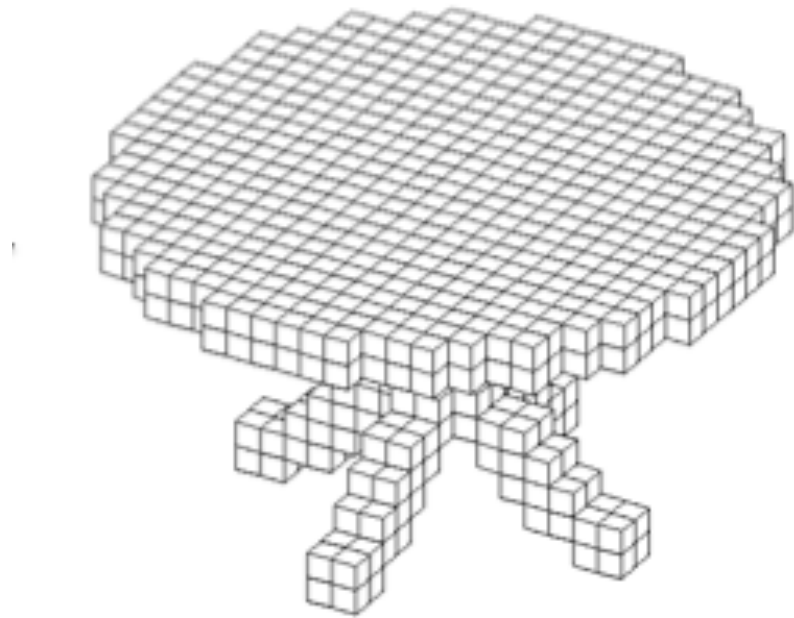
volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



Candidates:

multi-view images

depth map

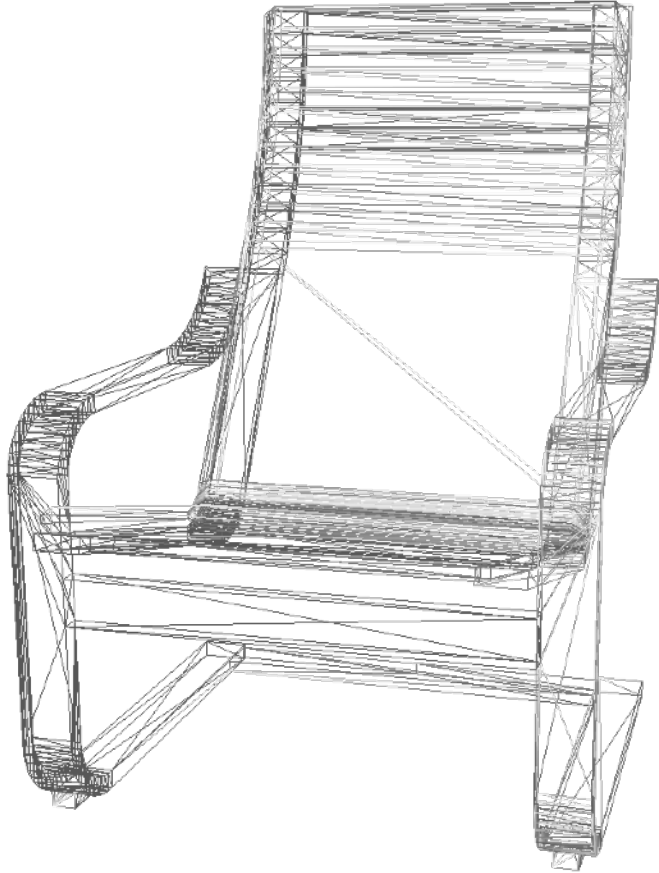
volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



Candidates:

multi-view images

depth map

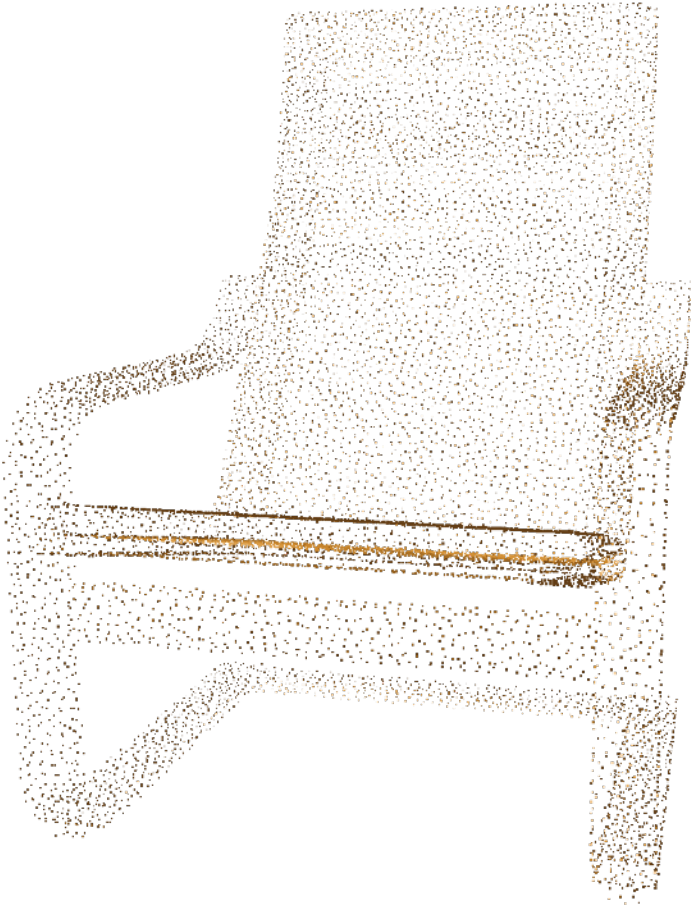
volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



Candidates:

multi-view images

depth map

volumetric

polygonal mesh

point cloud

primitive-based CAD models

3D representation



a chair assembled by cuboids

Candidates:

multi-view images

depth map

volumetric

polygonal mesh

point cloud

primitive-based CAD models

Two groups of representations

**Rasterized form
(regular grids)**

**Geometric form
(irregular)**

Candidates:

multi-view images

depth map

volumetric

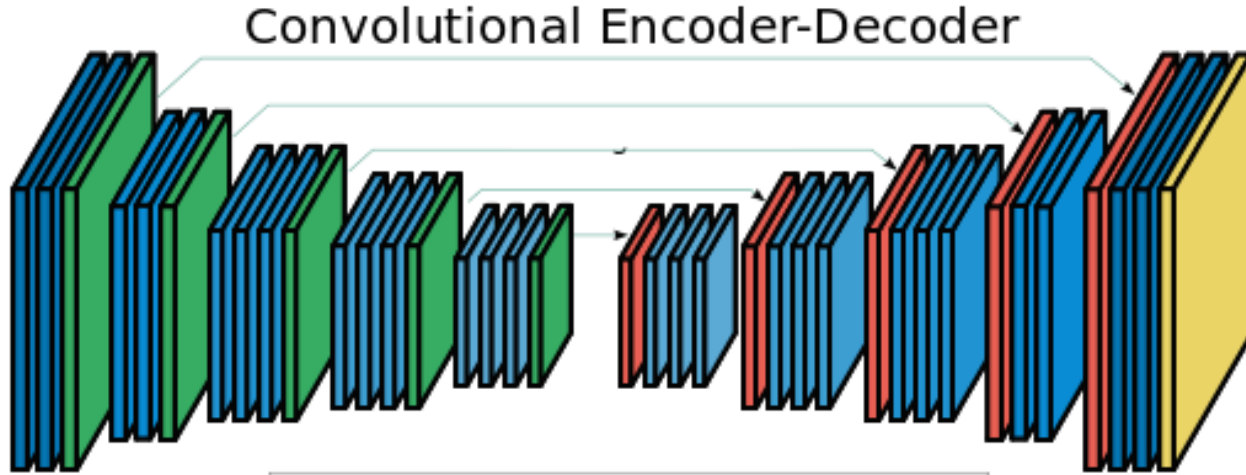
polygonal mesh

point cloud

primitive-based CAD models

Extant 3D DNNs work on grid-like representations

Candidates:



multi-view images

depth map

volumetric

polygonal mesh

point cloud

primitive-based CAD models

Ideally, a 3D representation should be

Friendly to learning

- easily formulated as the output of a neural network
- fast forward-/backward- propagation
- etc.

Ideally, a 3D representation should be

Friendly to learning

- easily formulated as the output of a neural network
- fast forward-/backward- propagation
- etc.

Flexible

- can precisely model a great variety of shapes
- etc.

Ideally, a 3D representation should be

Friendly to learning

- easily formulated as the output of a neural network
- fast forward-/backward- propagation
- etc.

Flexible


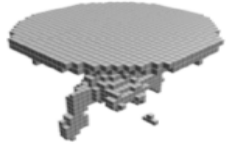

- can precisely model a great variety of shapes
- etc.

Geometrically manipulable


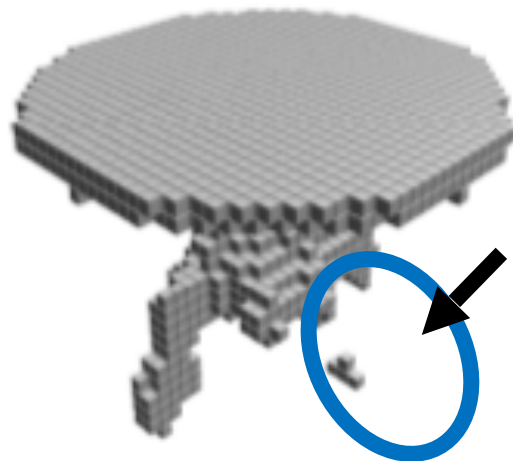





- geometrically deformable, interpolable and extrapolable for networks
- convenient to impose structural constraints
- etc.

Others

The problem of grid representations

	Affability to learning	Flexibility	Geometric manipulability
Multi-view images 	✓	✓	✗
Volumetric occupancy 	✗ Expensive to compute: $O(N^3)$	✓	✗
Depth map 	✓	✗ Cannot model "back side"	✓

Typical artifacts of volumetric reconstruction

	Affability		Geometric manipulability
Multi-view images 	 <p>Missing thin structures due to improper shape space structure</p> <p>hard for the network to rotate / deform / interpolate</p>		
Volumetric occupancy 			
Depth map 			

How about learning to predict geometric forms?

Candidates:

**Rasterized form
(regular grids)**

multi-view images

depth map

volumetric

**Geometric form
(irregular)**

polygonal mesh



point cloud



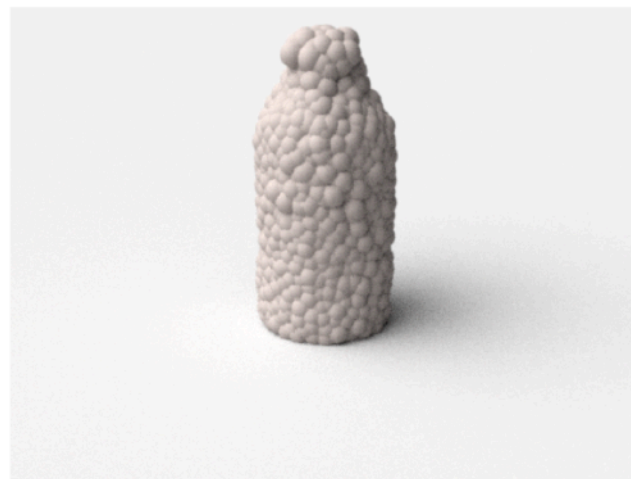
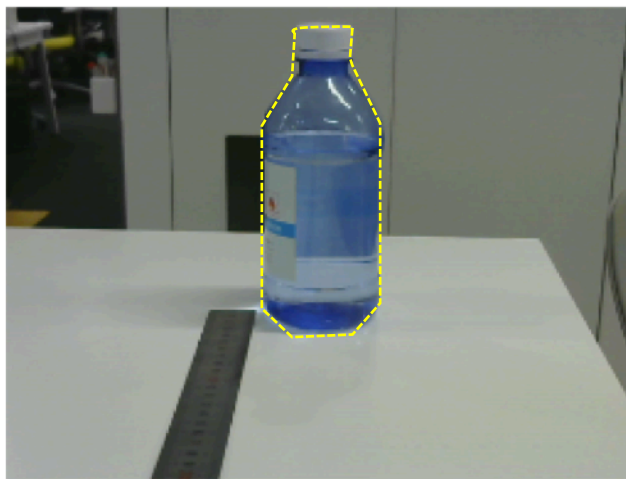
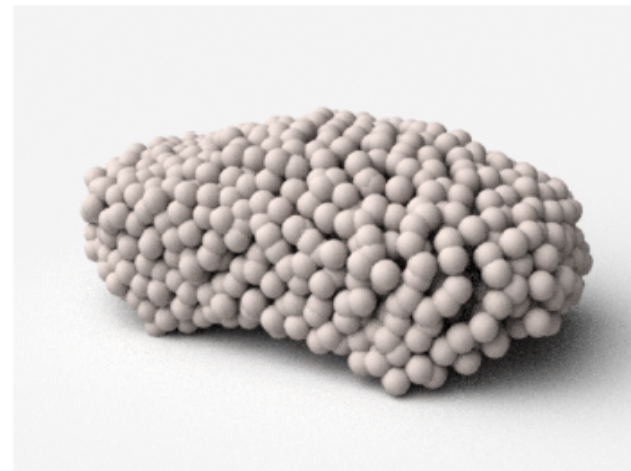
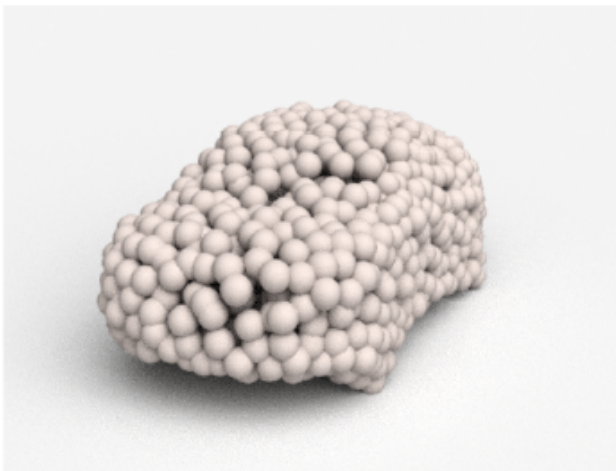
primitive-based CAD models

Outline

- Motivation
- Data
- **Algorithms for 2D-3D lifting**
- Shape abstraction by volumetric primitives

Our result: 3D reconstruction from real Images

CVPR 2017, A Point Set Generation Network for 3D Object Reconstruction from a Single Image

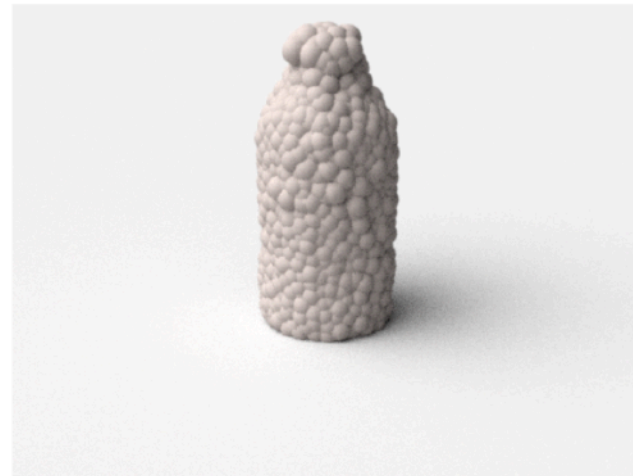
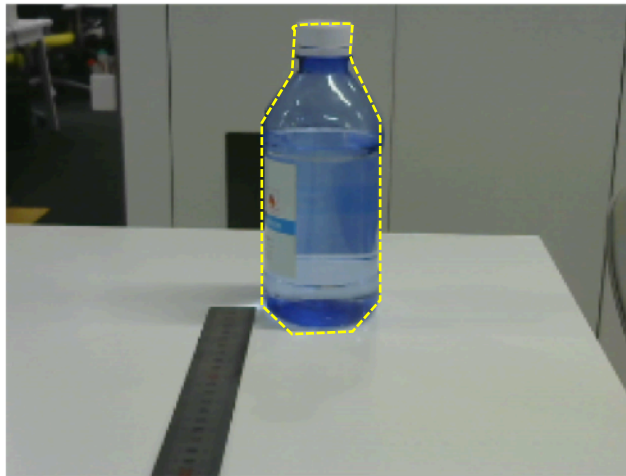
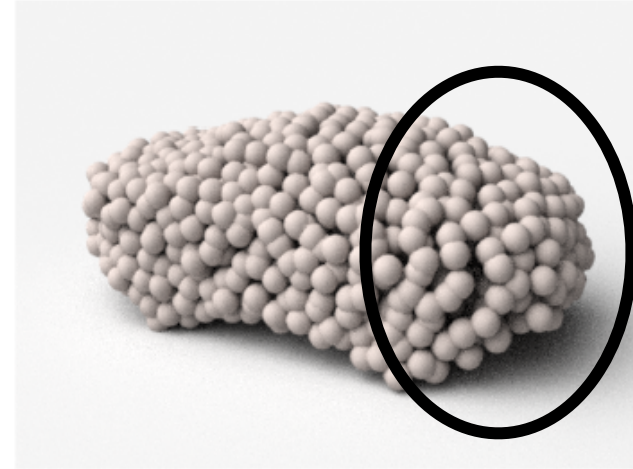
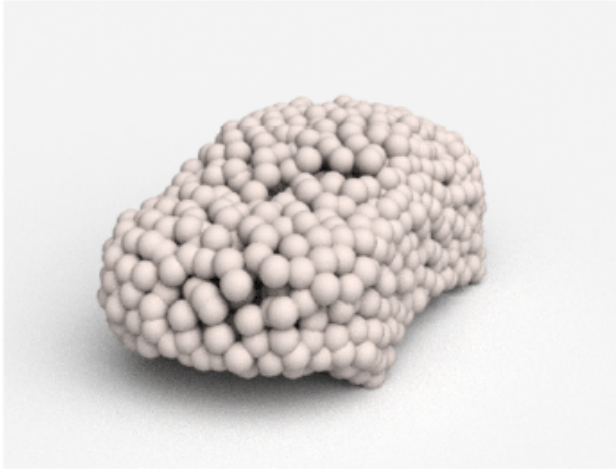


Input

Reconstructed 3D point cloud

Our result: 3D reconstruction from real Images

CVPR 2017, A Point Set Generation Network for 3D Object Reconstruction from a Single Image



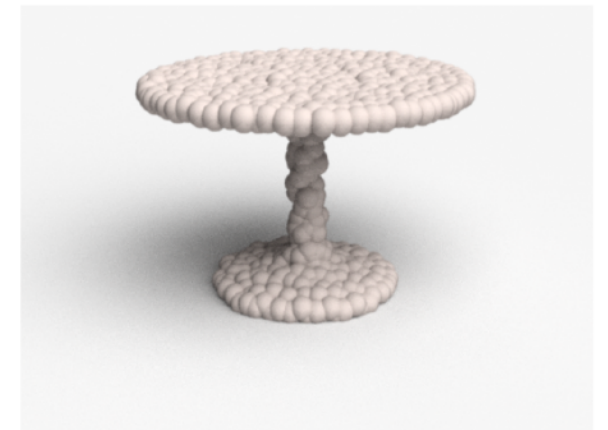
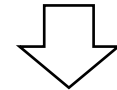
Input

Reconstructed 3D point cloud

3D point clouds

✓ Flexible

- a few thousands of points can precisely model a great variety of shapes



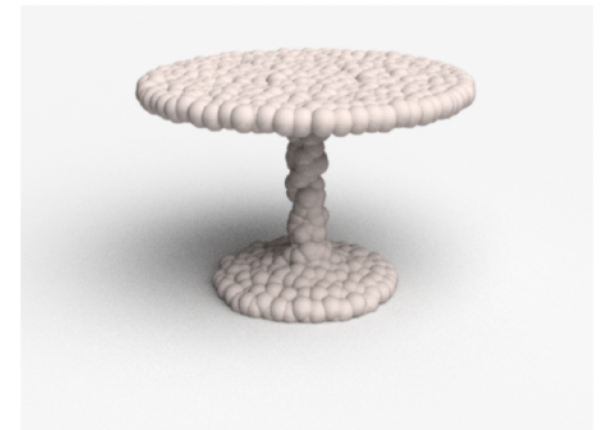
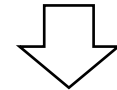
3D point clouds

✓ Flexible

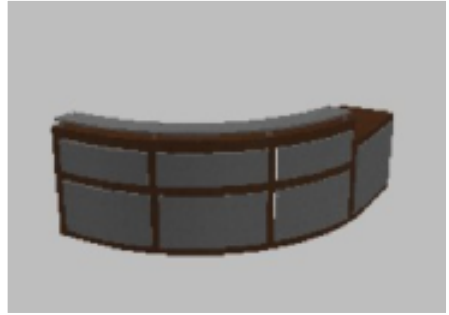
- a few thousands of points can precisely model a great variety of shapes

✓ Geometrically manipulable

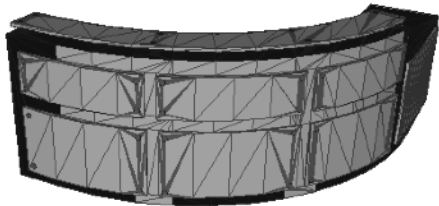
- deformable
- interpolable, extrapolable
- convenient to impose structural constraints



Pipeline



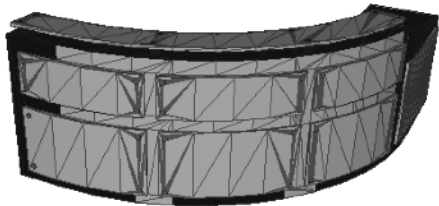
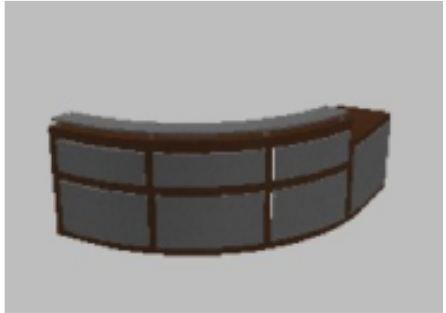
render



Pipeline

2K object categories
200K shapes
~10M image/point set pairs

render



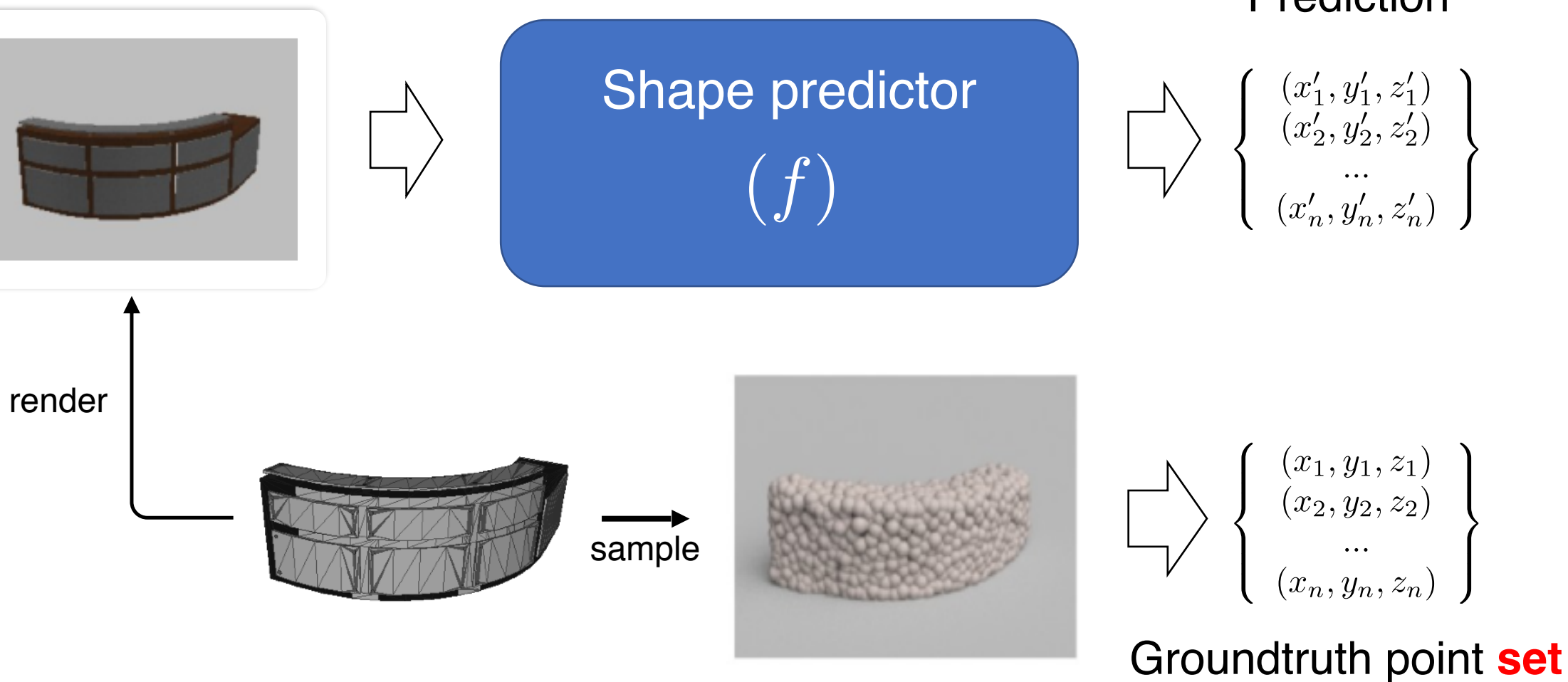
sample



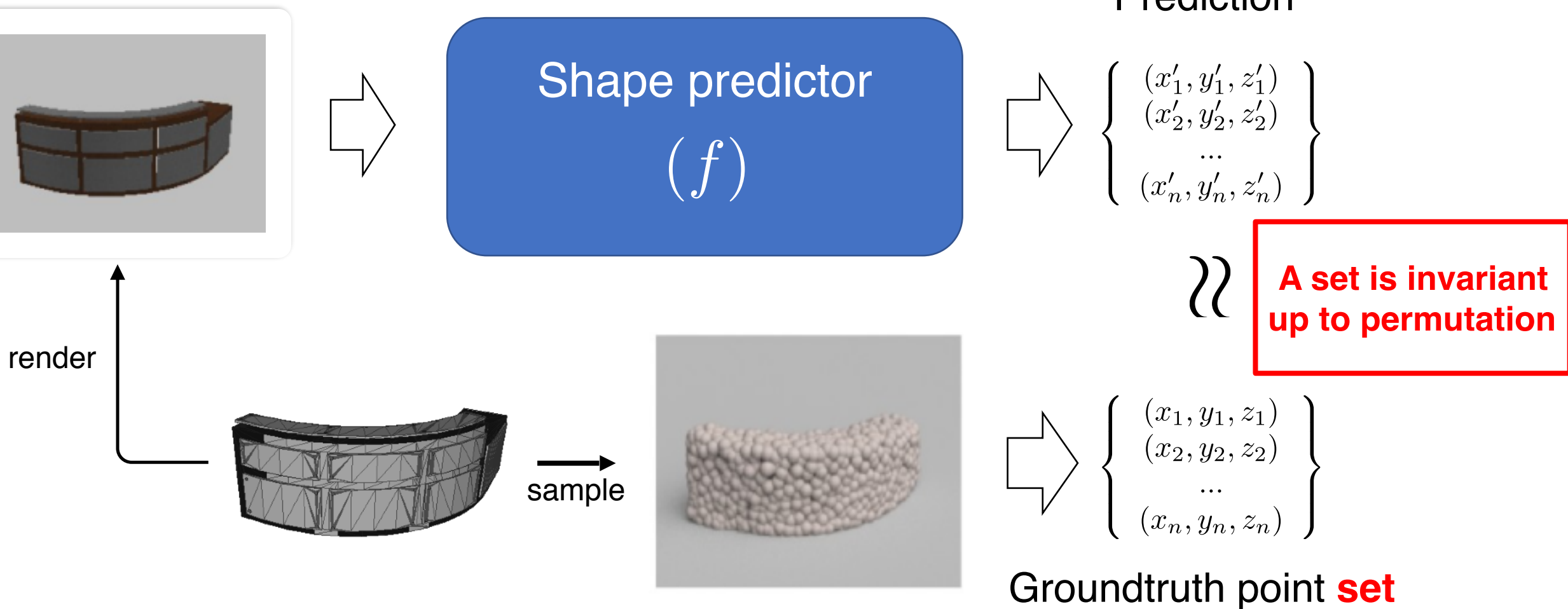
$$\left\{ \begin{array}{c} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ \dots \\ (x_n, y_n, z_n) \end{array} \right\}$$

Groundtruth point **set**

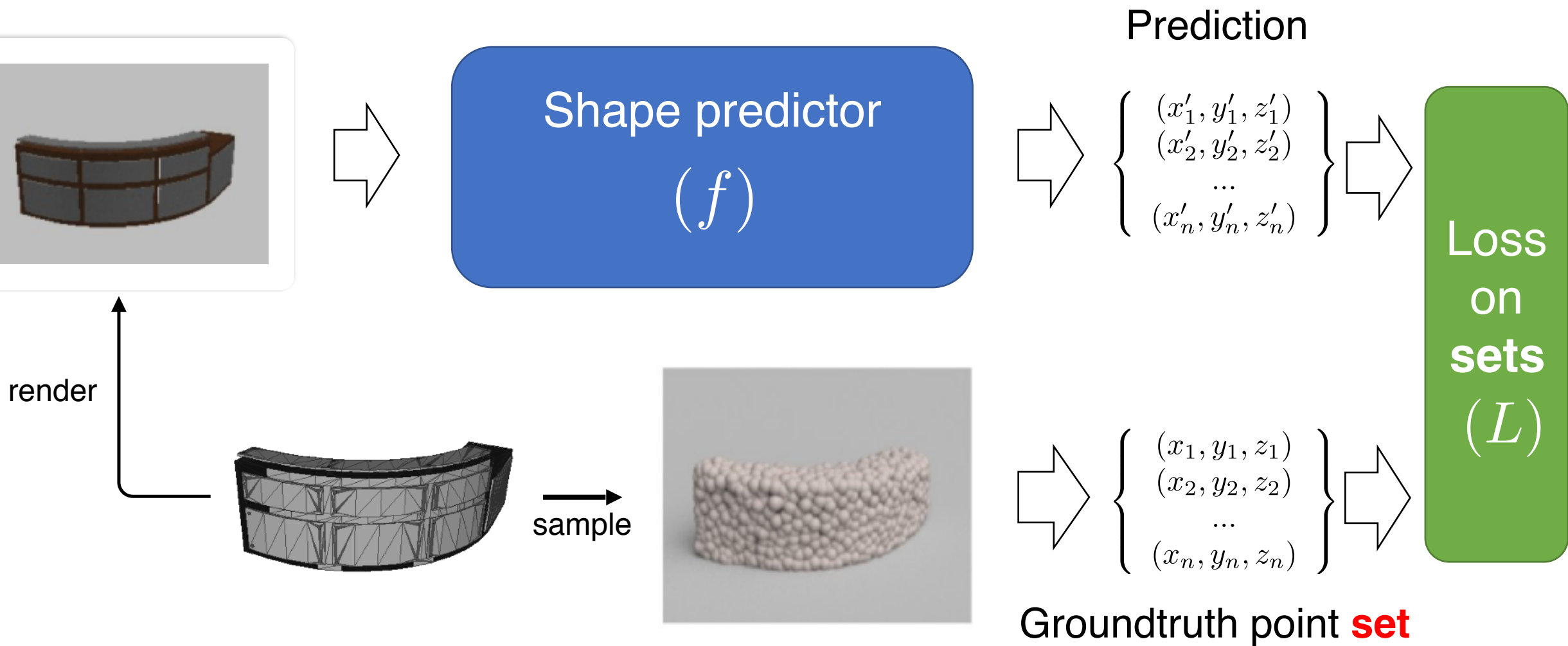
Pipeline



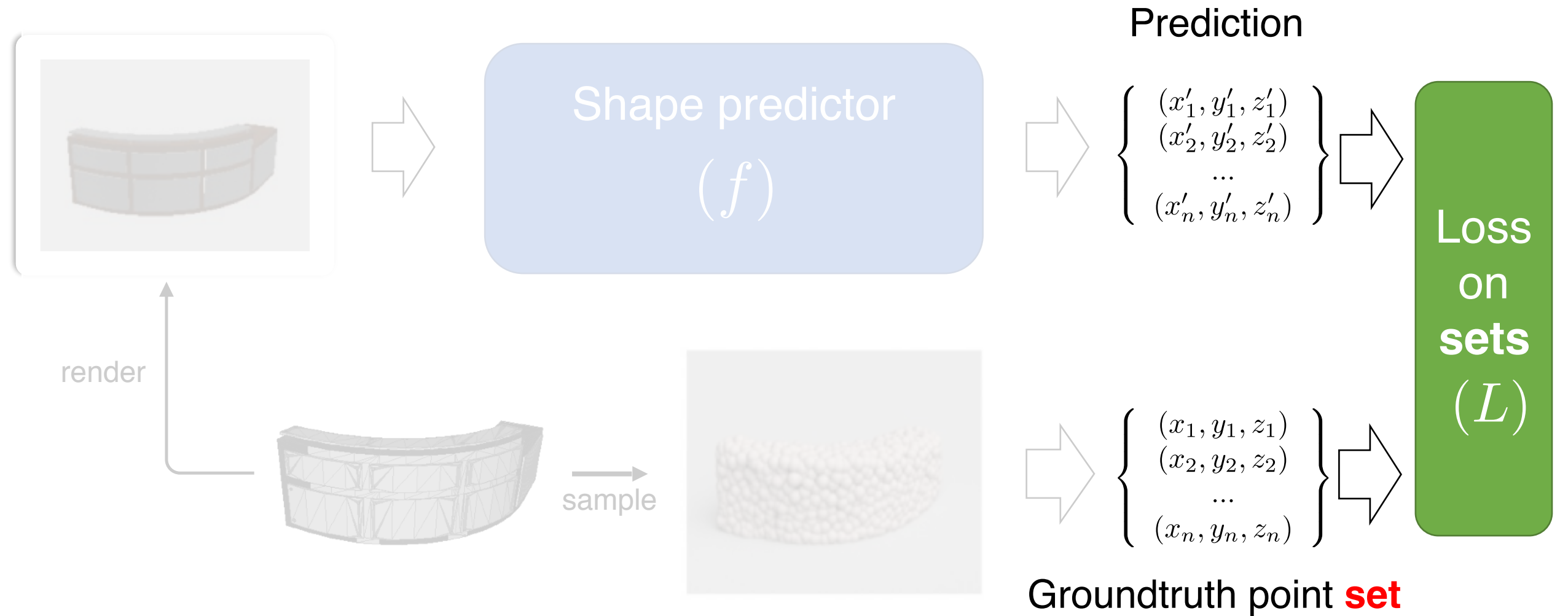
Pipeline



Pipeline

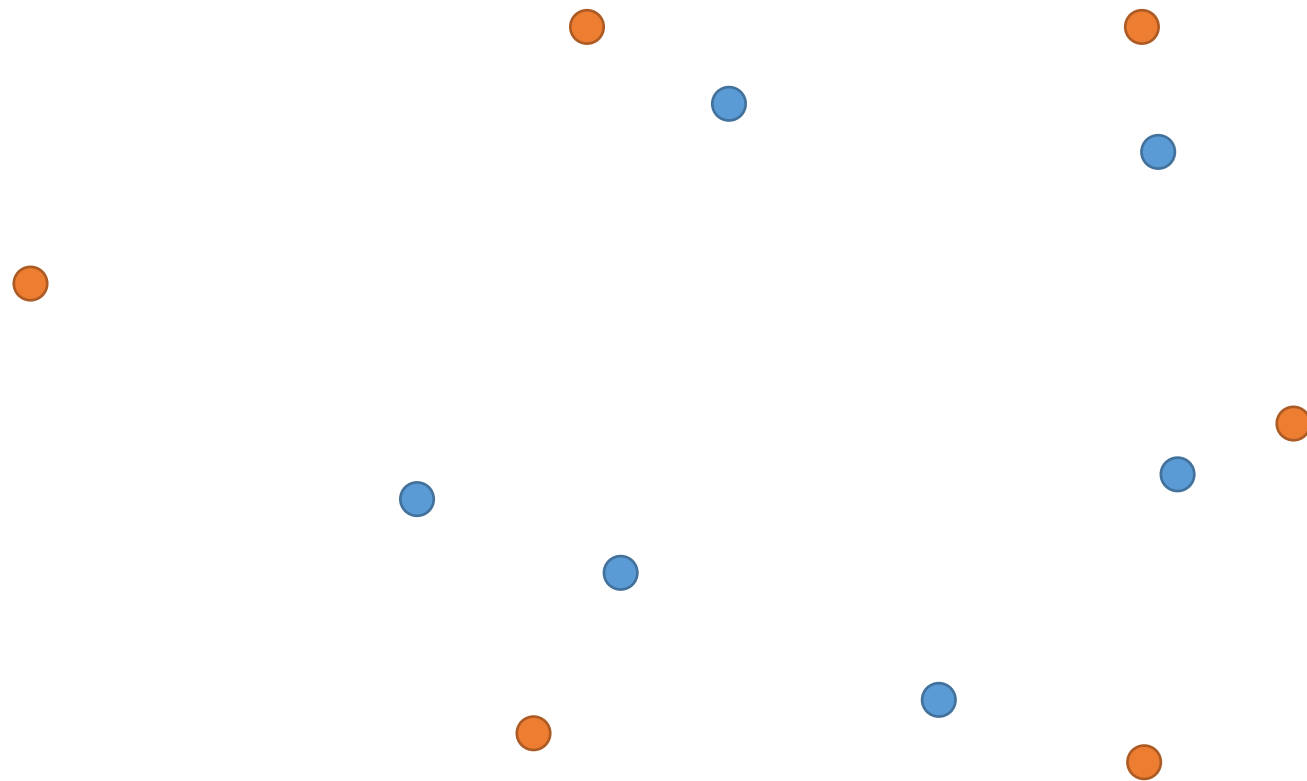


Pipeline



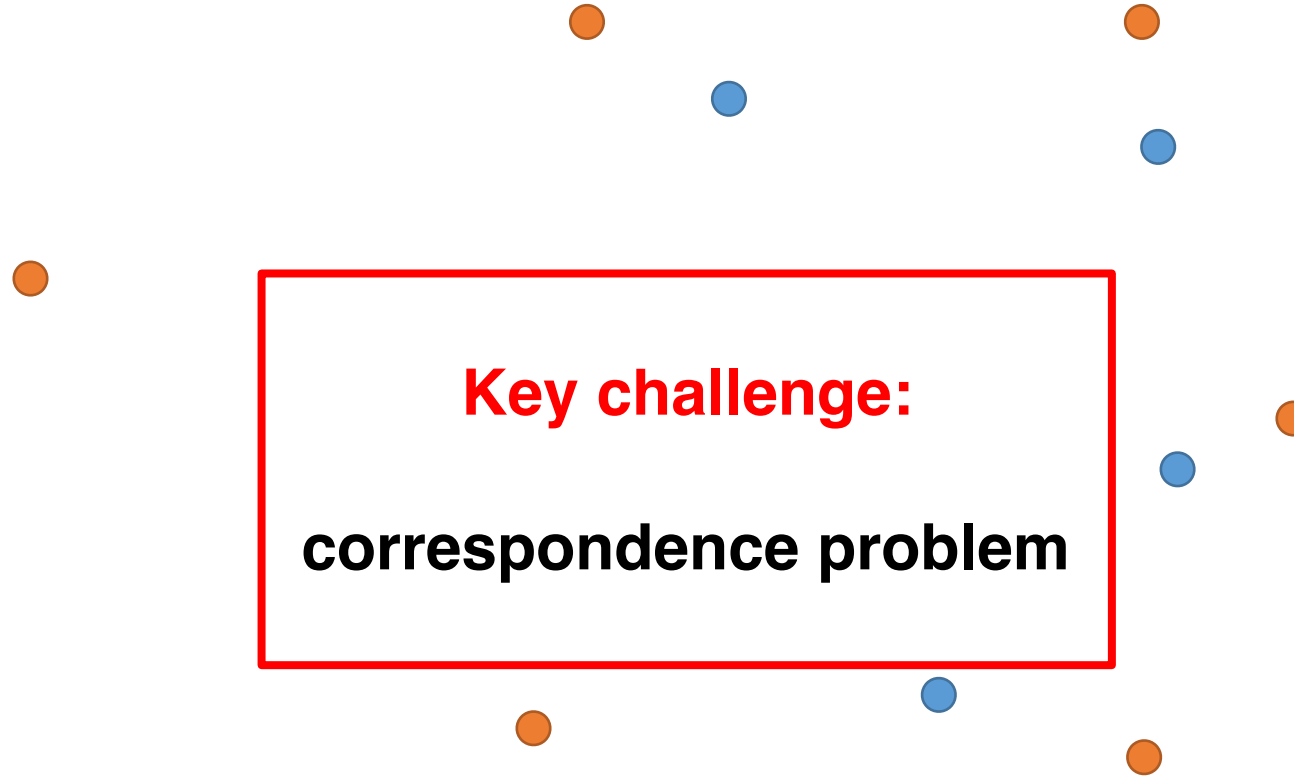
Set comparison

Given two sets of points, measure their discrepancy



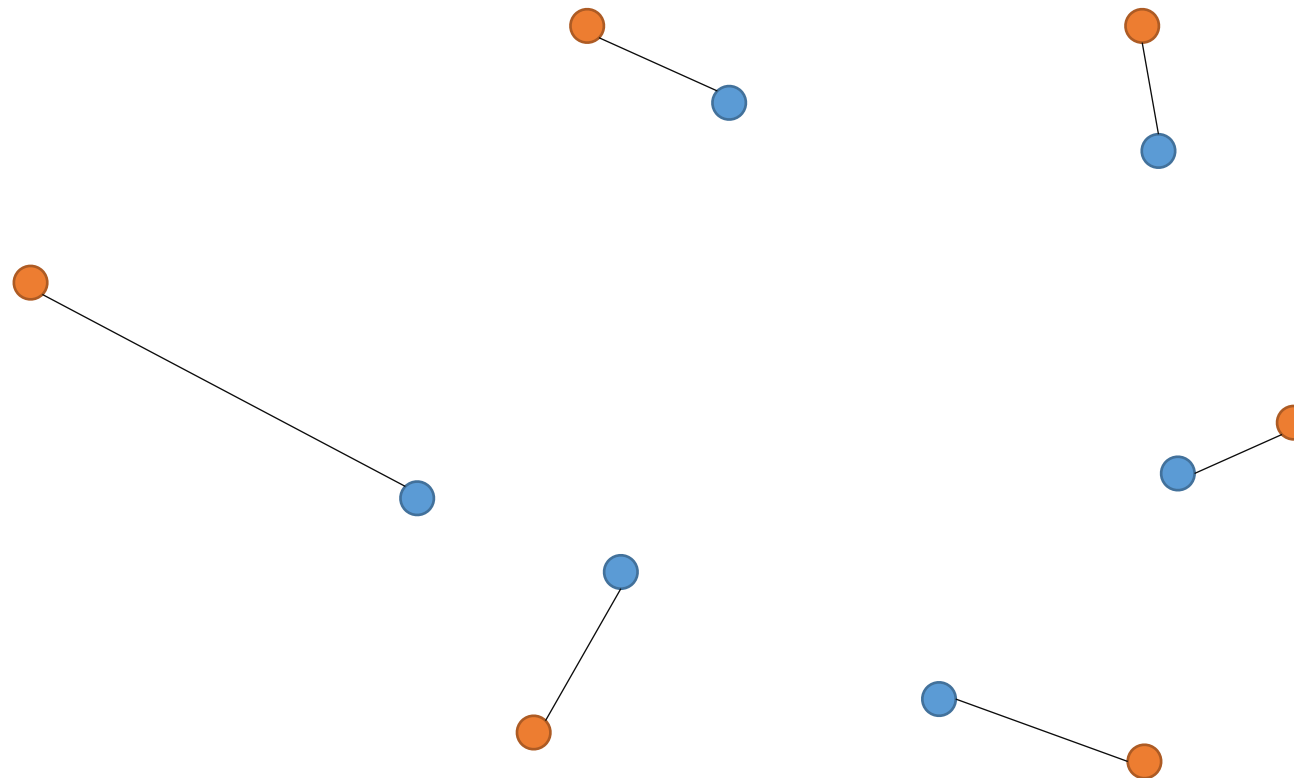
Set comparison

Given two sets of points, measure their discrepancy



Correspondence (I): optimal assignment

Given two sets of points, measure their discrepancy

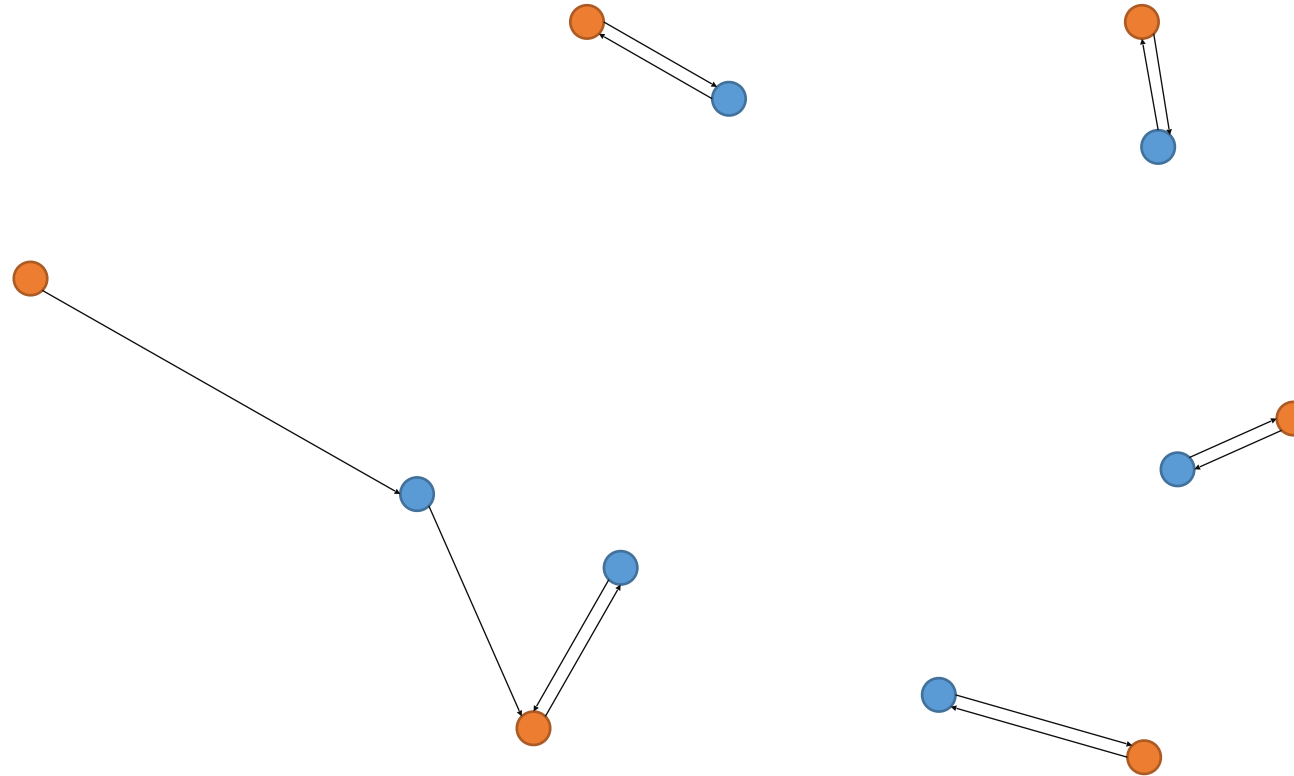


a.k.a Earth Mover's distance (EMD)

$$d_{EMD}(S_1, S_2) = \min_{\phi: S_1 \rightarrow S_2} \sum_{x \in S_1} \|x - \phi(x)\|_2 \quad \text{where } \phi : S_1 \rightarrow S_2 \text{ is a bijection.}$$

Correspondence (II): closest point

Given two sets of points, measure their discrepancy



a.k.a Chamfer distance (CD)

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$

Required properties of distance metrics

Geometric requirement

Computational requirement

Required properties of distance metrics

Geometric requirement

- Reflects natural shape differences
- Induce a nice space for *shape interpolations*

Computational requirement

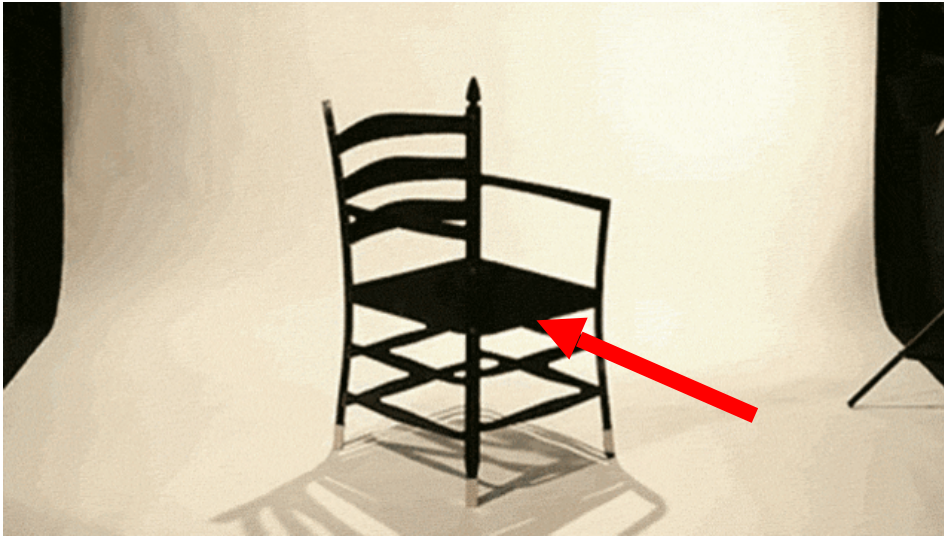
How distance metric affects learning?

A fundamental issue: inherent ambiguity in 2D-3D dimension lifting



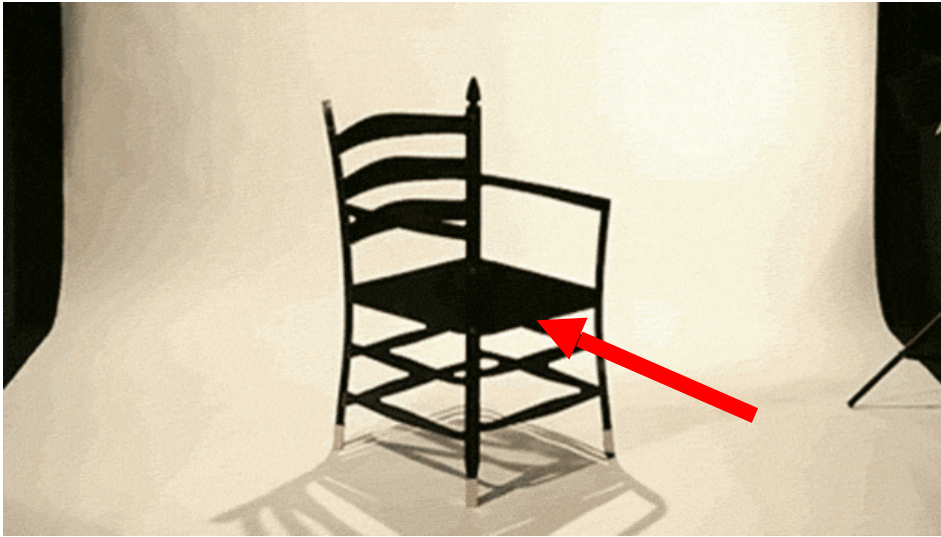
How distance metric affects learning?

A fundamental issue: inherent ambiguity in 2D-3D dimension lifting



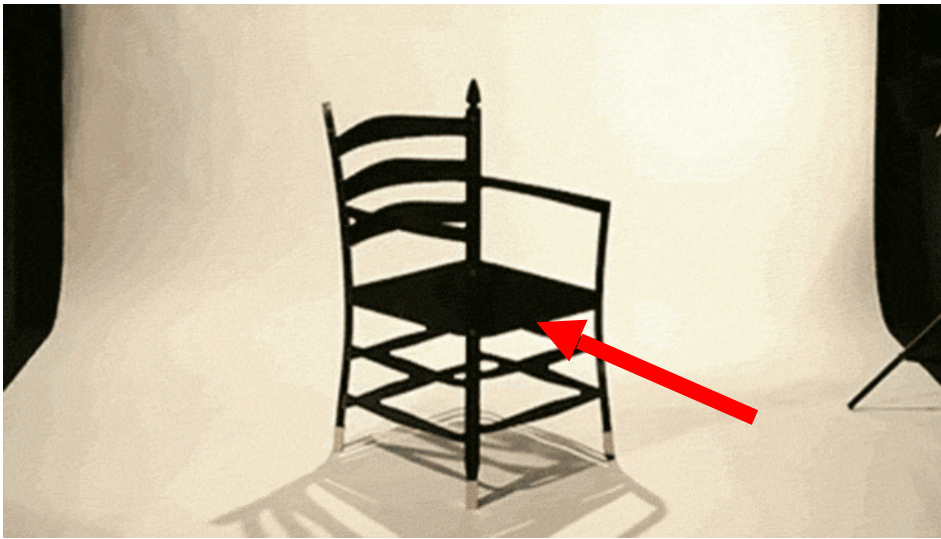
How distance metric affects learning?

A fundamental issue: inherent ambiguity in 2D-3D dimension lifting



How distance metric affects learning?

A fundamental issue: inherent ambiguity in 2D-3D dimension lifting



- By loss minimization, the network tends to predict a “**mean shape**” that **averages out** uncertainty

Distance metrics affect mean shapes

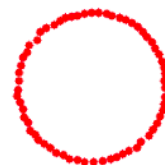
The mean shape carries characteristics of the distance metric

$$\bar{x} = \operatorname{argmin}_x \mathbb{E}_{s \sim \mathcal{S}} [d(x, s)]$$

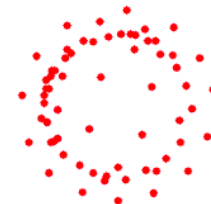
continuous
hidden variable
(radius)



Input



EMD mean



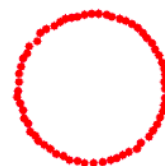
Chamfer mean

Mean shapes from distance metrics

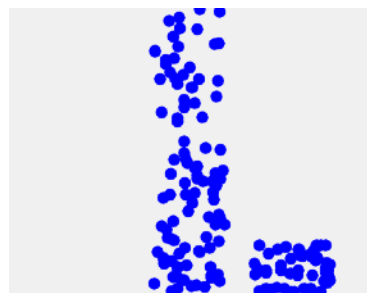
The mean shape carries characteristics of the distance metric

$$\bar{x} = \operatorname{argmin}_x \mathbb{E}_{s \sim \mathcal{S}} [d(x, s)]$$

continuous
hidden variable
(radius)



discrete
hidden variable
(add-on location)



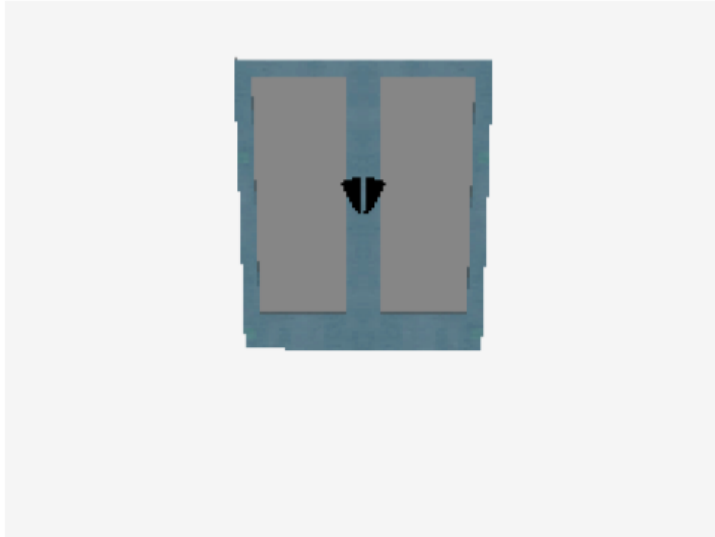
Input

EMD mean

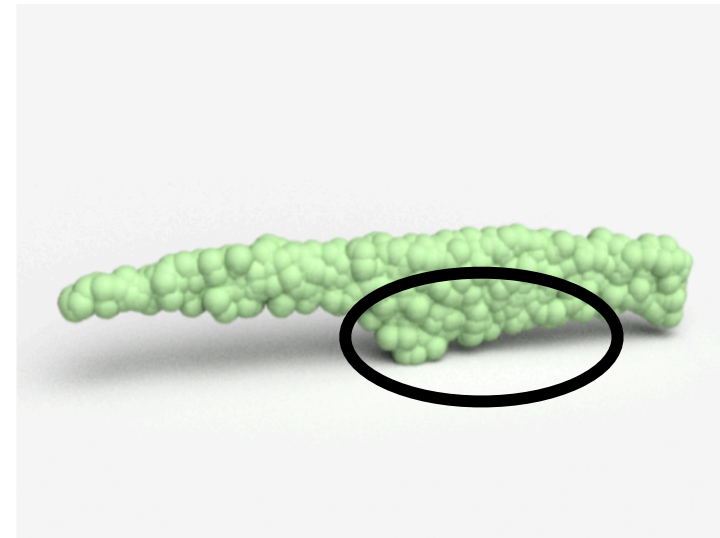
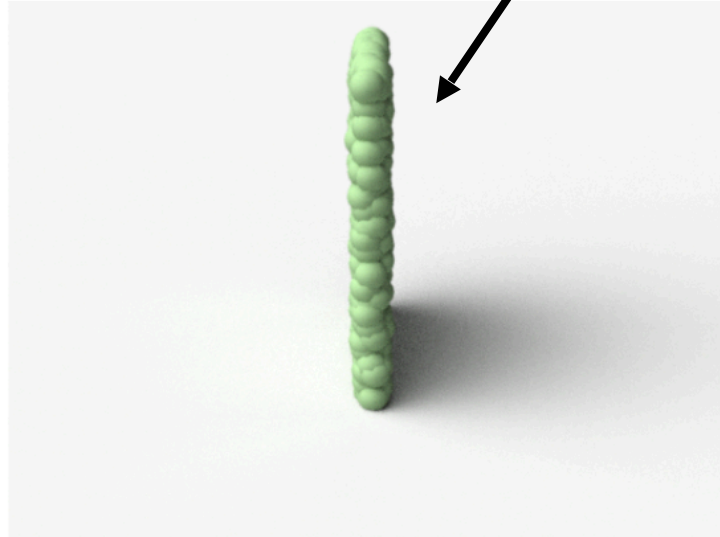
Chamfer mean

Comparison of predictions by EMD versus CD

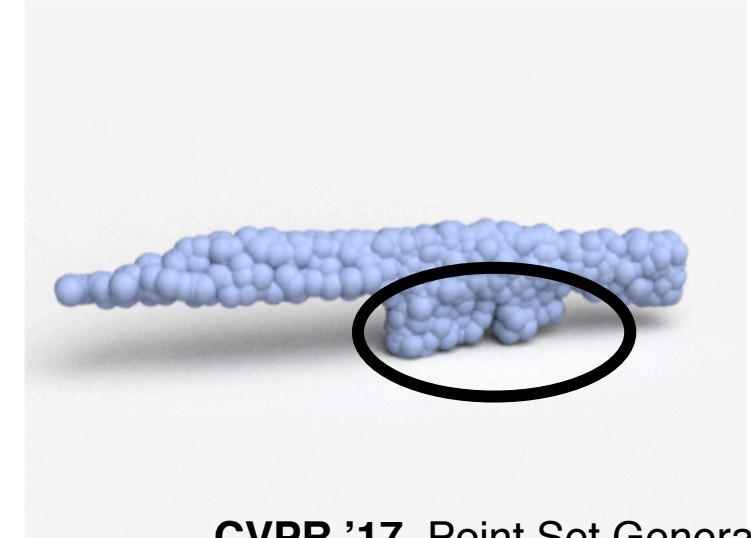
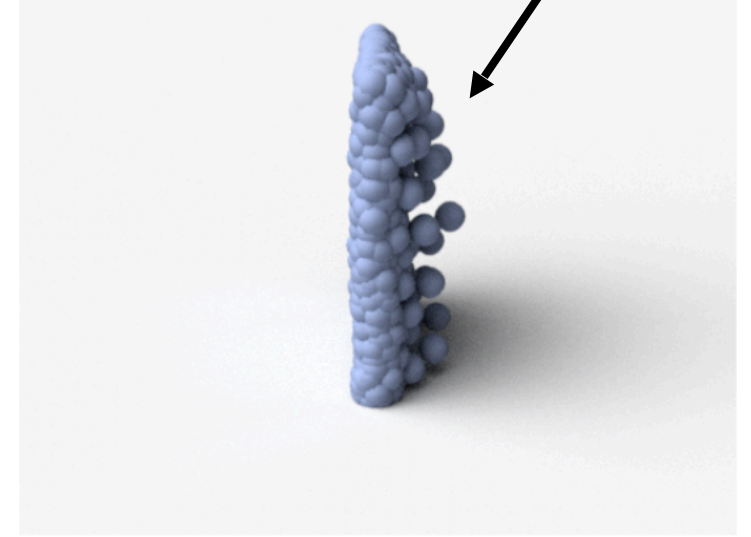
Input



EMD



Chamfer



Required properties of distance metrics

Geometric requirement

- Reflects natural shape differences
- Induce a nice space for shape interpolations

Computational requirement

- Defines a loss function that is numerically easy to optimize

Computational requirement of metrics

To be used as a loss function, the metric has to be

- **Differentiable** with respect to point locations
- **Efficient** to compute

Computational requirement of metrics

- **Differentiable** with respect to point location

Chamfer distance

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$



Earth Mover's distance

$$d_{EMD}(S_1, S_2) = \min_{\phi: S_1 \rightarrow S_2} \sum_{x \in S_1} \|x - \phi(x)\|_2 \quad \text{where } \phi: S_1 \rightarrow S_2 \text{ is a bijection.}$$



- Simple function of coordinates
- In general positions, the correspondence is unique
- **With infinitesimal movement, the correspondence does not change**

Conclusion: differentiable almost everywhere

Computational requirement of metrics

- **Differentiable** with respect to point location

- For many **algorithms** (sorting, shortest path, network flow, ...),
- an infinitesimal change to model parameters (almost) does not change solution structure,

leads to **differentiable a.e.!**

Computational requirement of metrics

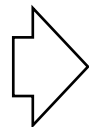
- **Efficient** to compute

Chamfer distance: trivially parallelizable on CUDA

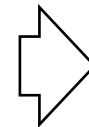
Earth Mover's distance (optimal assignment):

- We implement a **distributed** approximation algorithm on CUDA
- Based upon [Bertsekas, 1985], $(1 + \epsilon)$ -approximation

Pipeline

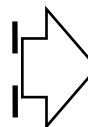


Deep network
(f)

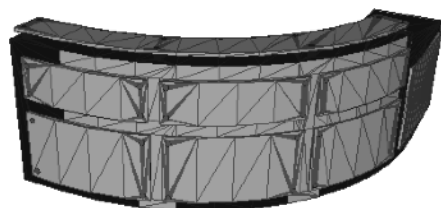


Prediction

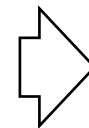
$$\begin{Bmatrix} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ \dots \\ (x_n, y_n, z_n) \end{Bmatrix}$$



Loss
on
sets
(L)



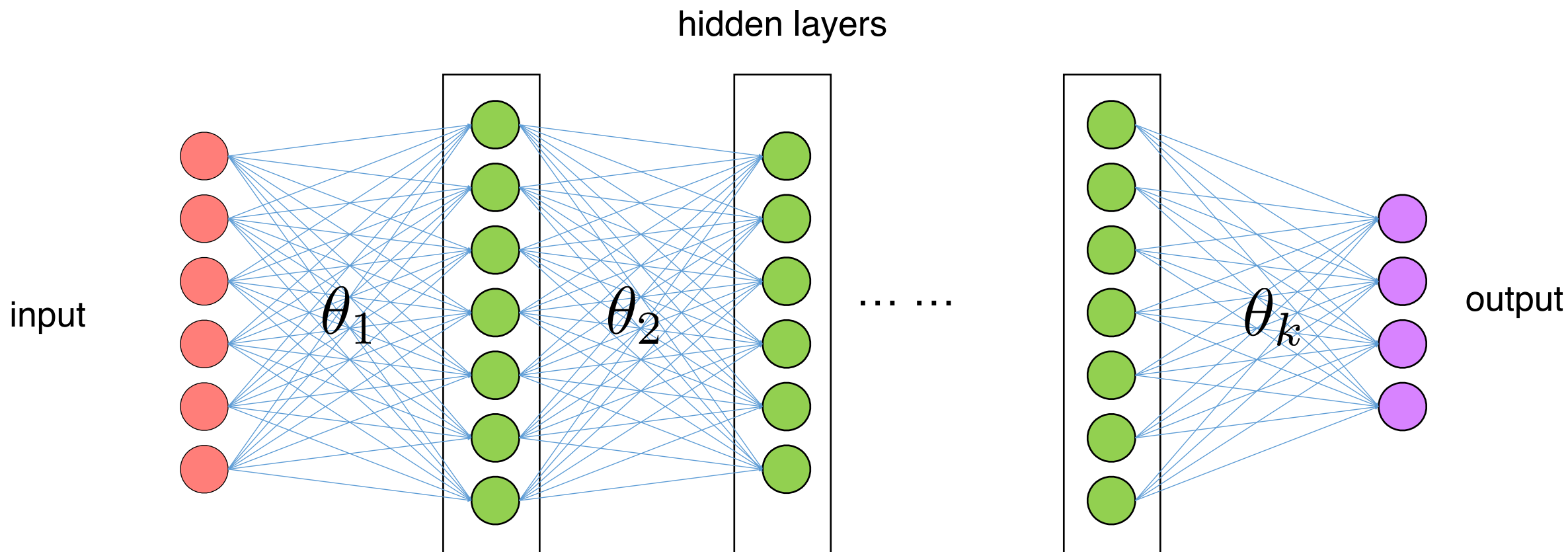
sample



$$\begin{Bmatrix} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ \dots \\ (x_n, y_n, z_n) \end{Bmatrix}$$



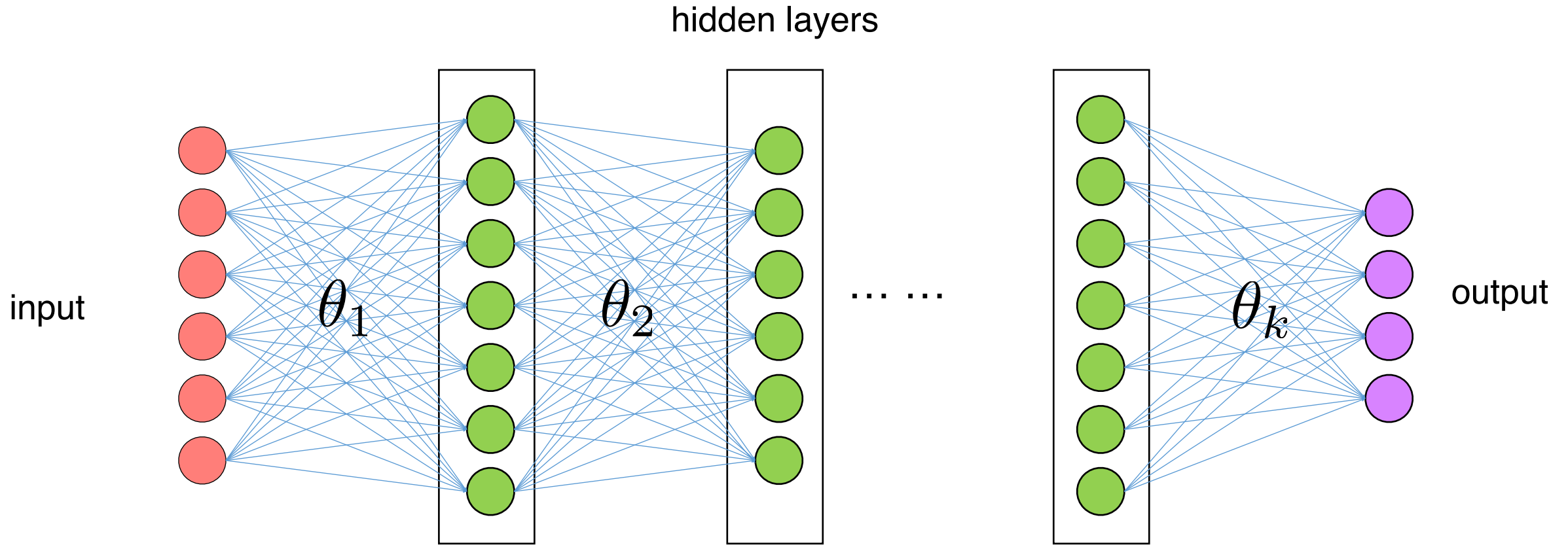
Deep neural network



Universal function approximator

- A cascade of layers

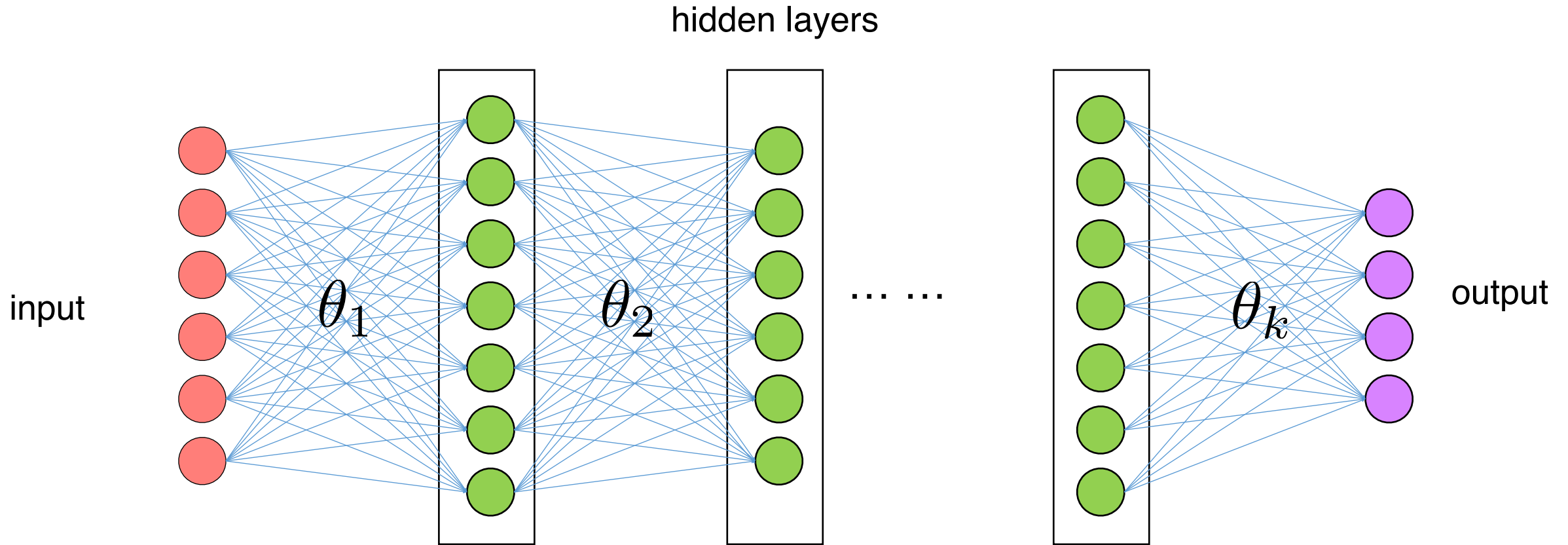
Deep neural network



Universal function approximator

- A cascade of layers
- Each layer conducts a simple transformation (parameterized)

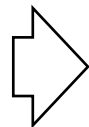
Deep neural network



Universal function approximator

- A cascade of layers
- Each layer conducts a simple transformation (parameterized)
- Millions of parameters, has to be fitted by many data

Pipeline

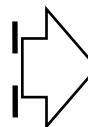


Deep network
(f)

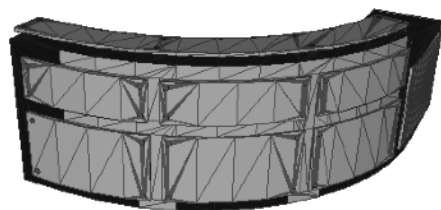


Prediction

$$\begin{Bmatrix} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ \dots \\ (x_n, y_n, z_n) \end{Bmatrix}$$



Loss
on
sets
(L)



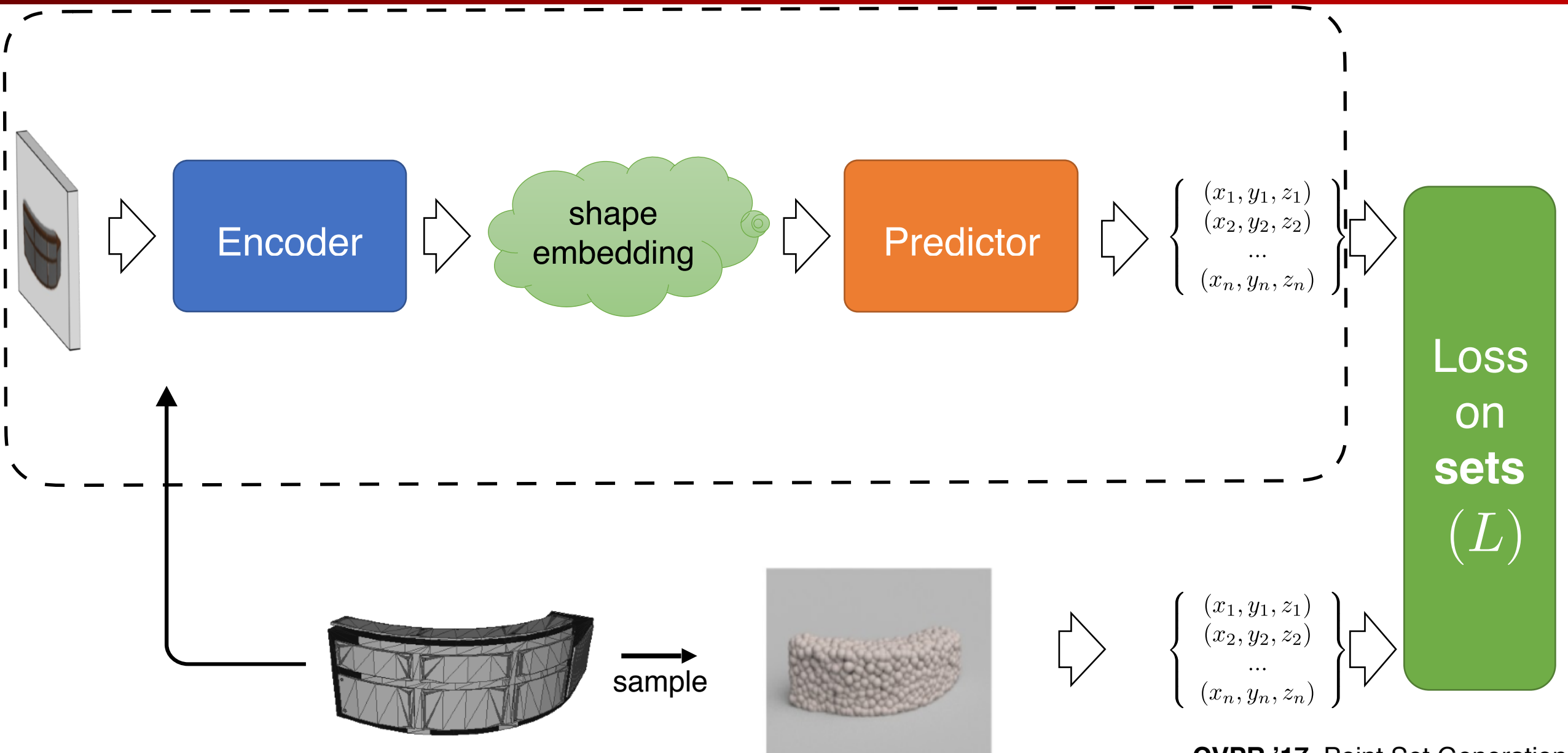
sample



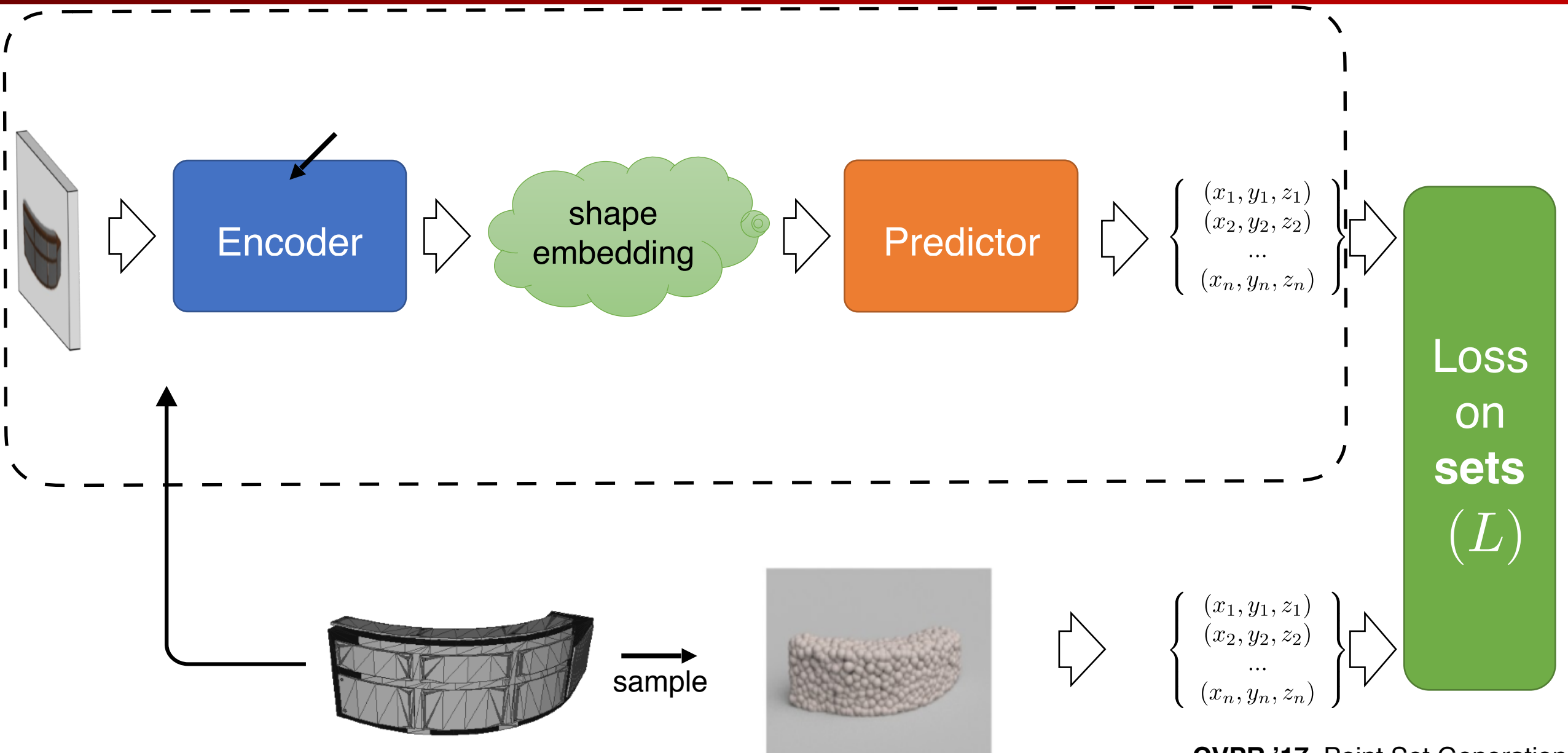
$$\begin{Bmatrix} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ \dots \\ (x_n, y_n, z_n) \end{Bmatrix}$$



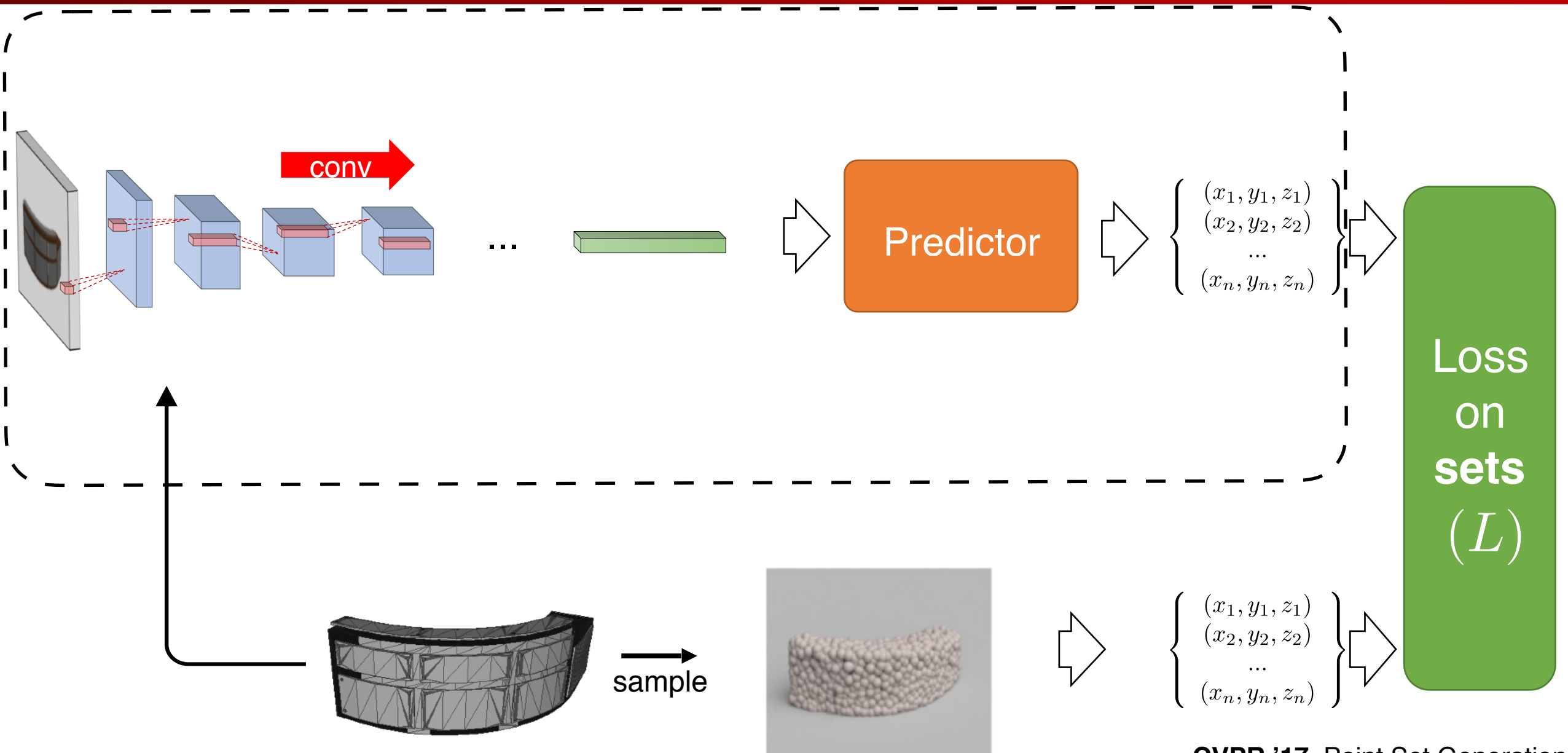
Pipeline



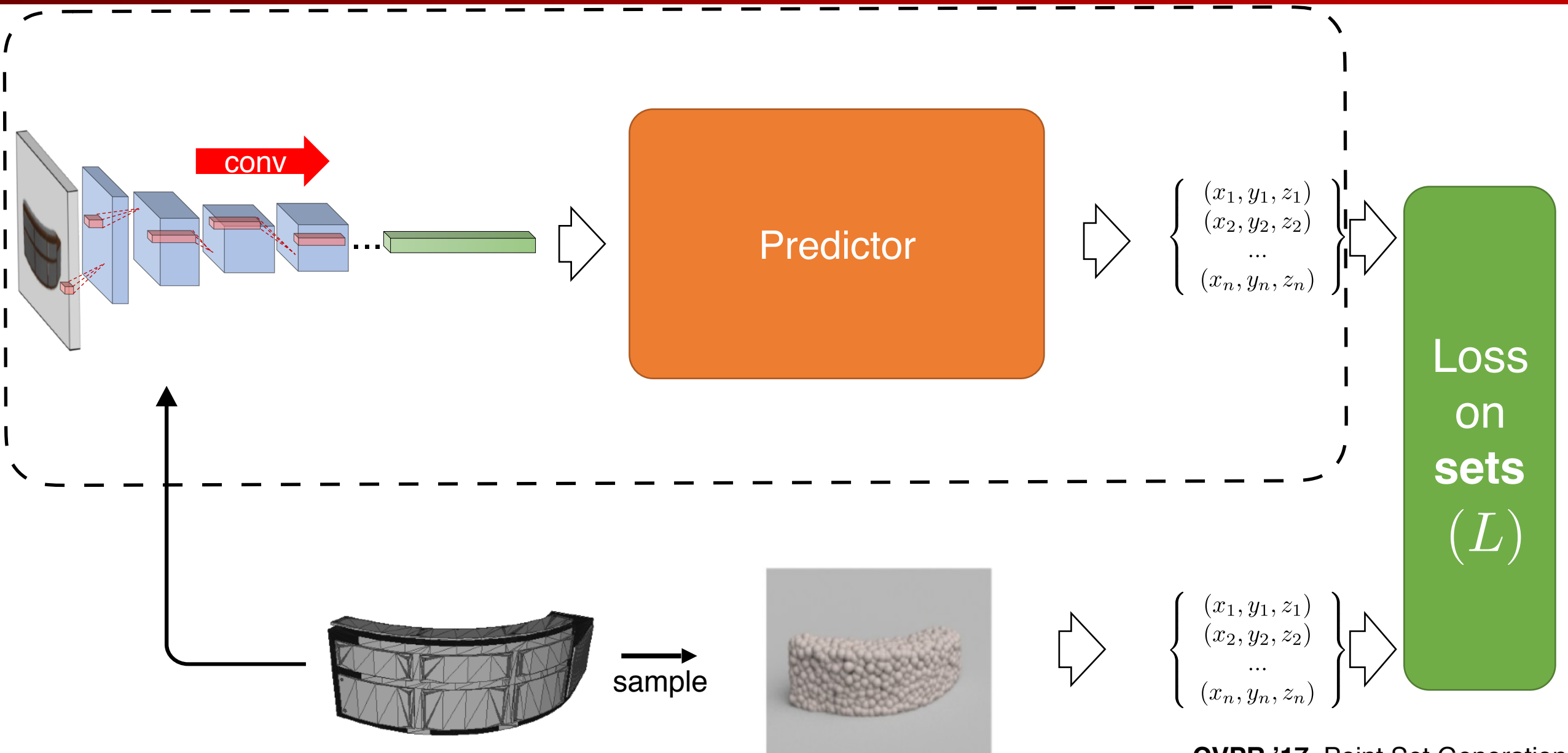
Pipeline



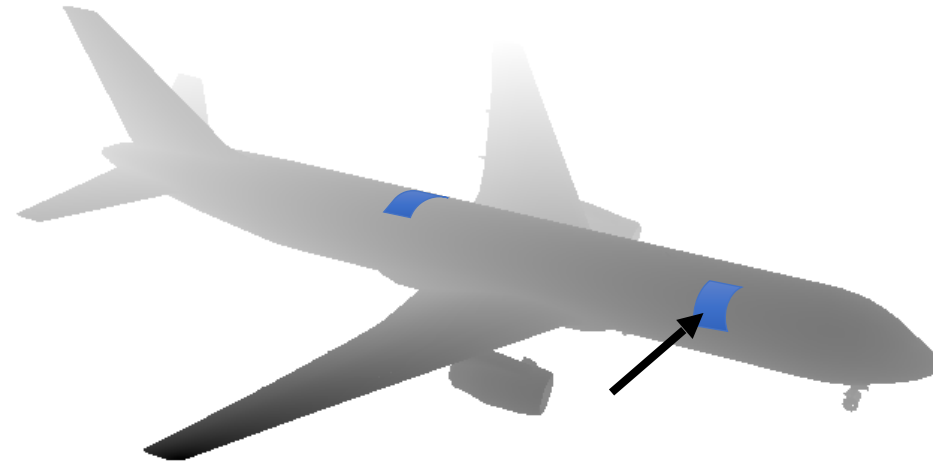
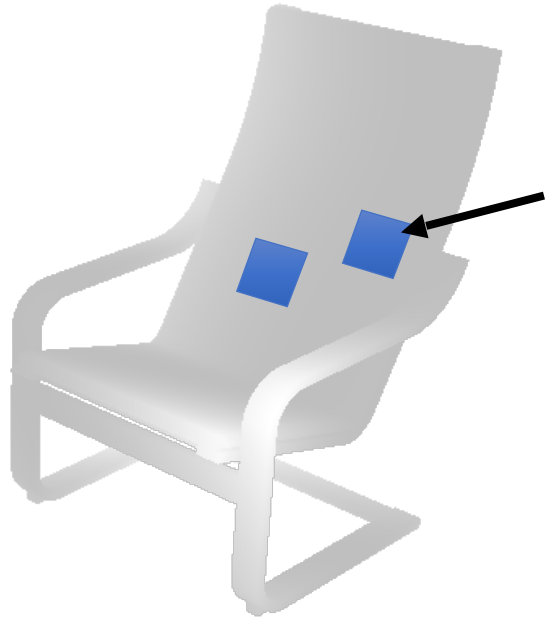
Pipeline



Pipeline

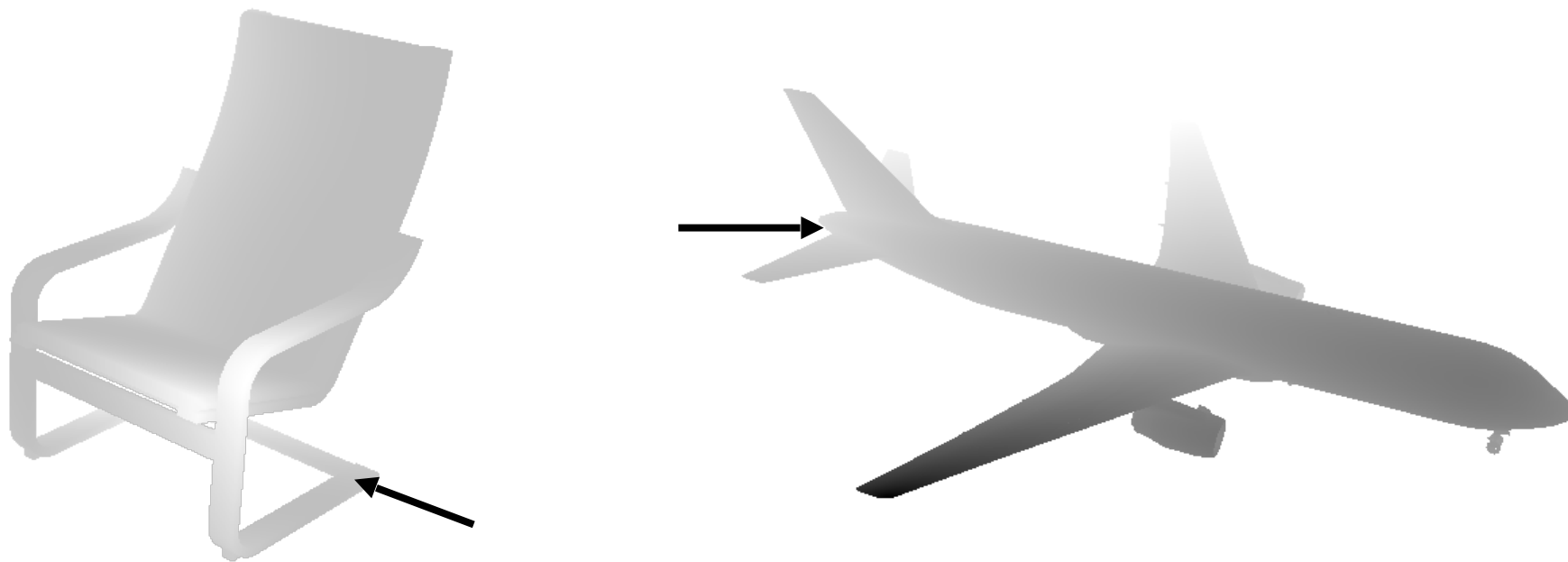


Natural statistics of geometry



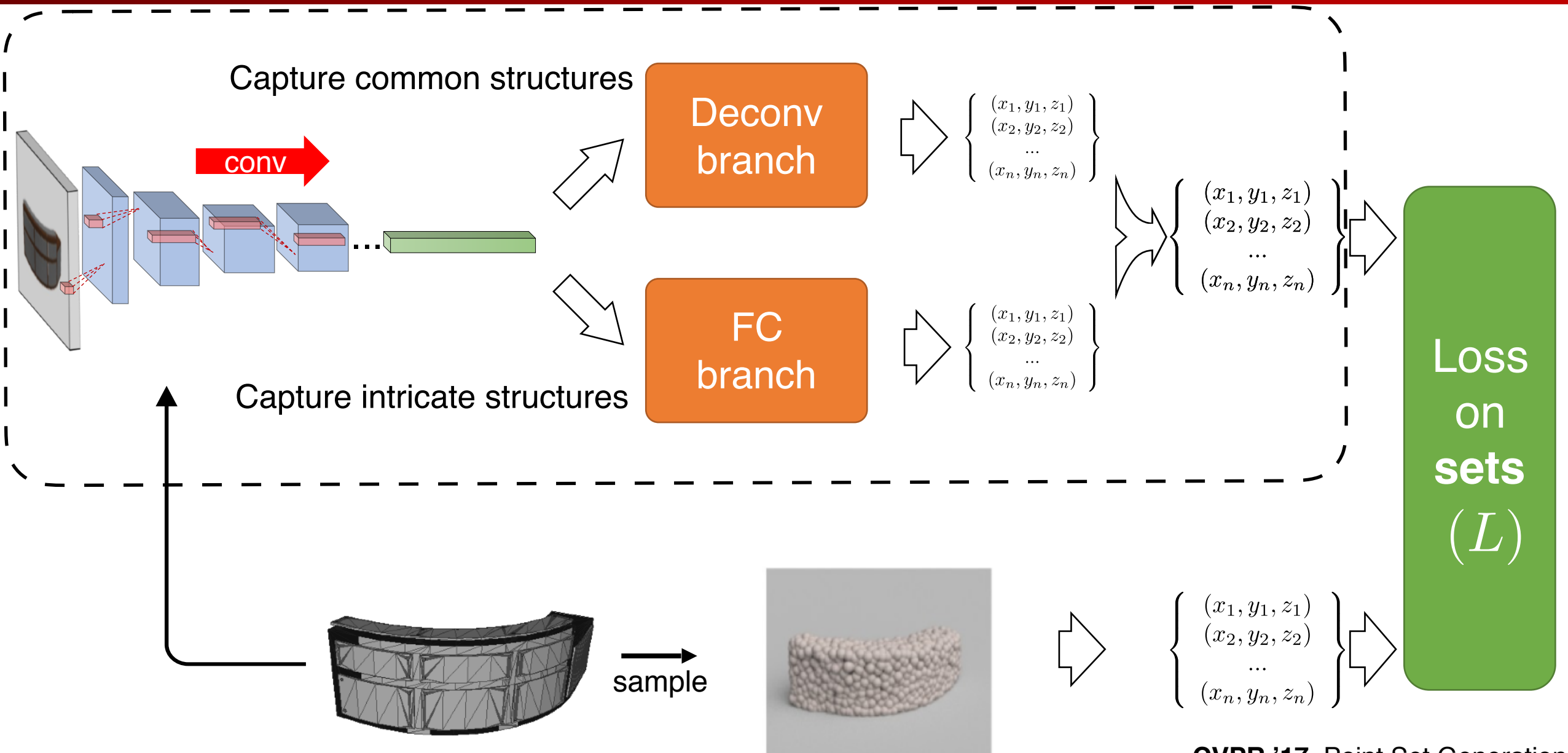
- Many local structures are common
 - e.g., planar patches, cylindrical patches
 - **strong local correlation** among point coordinates

Natural statistics of geometry

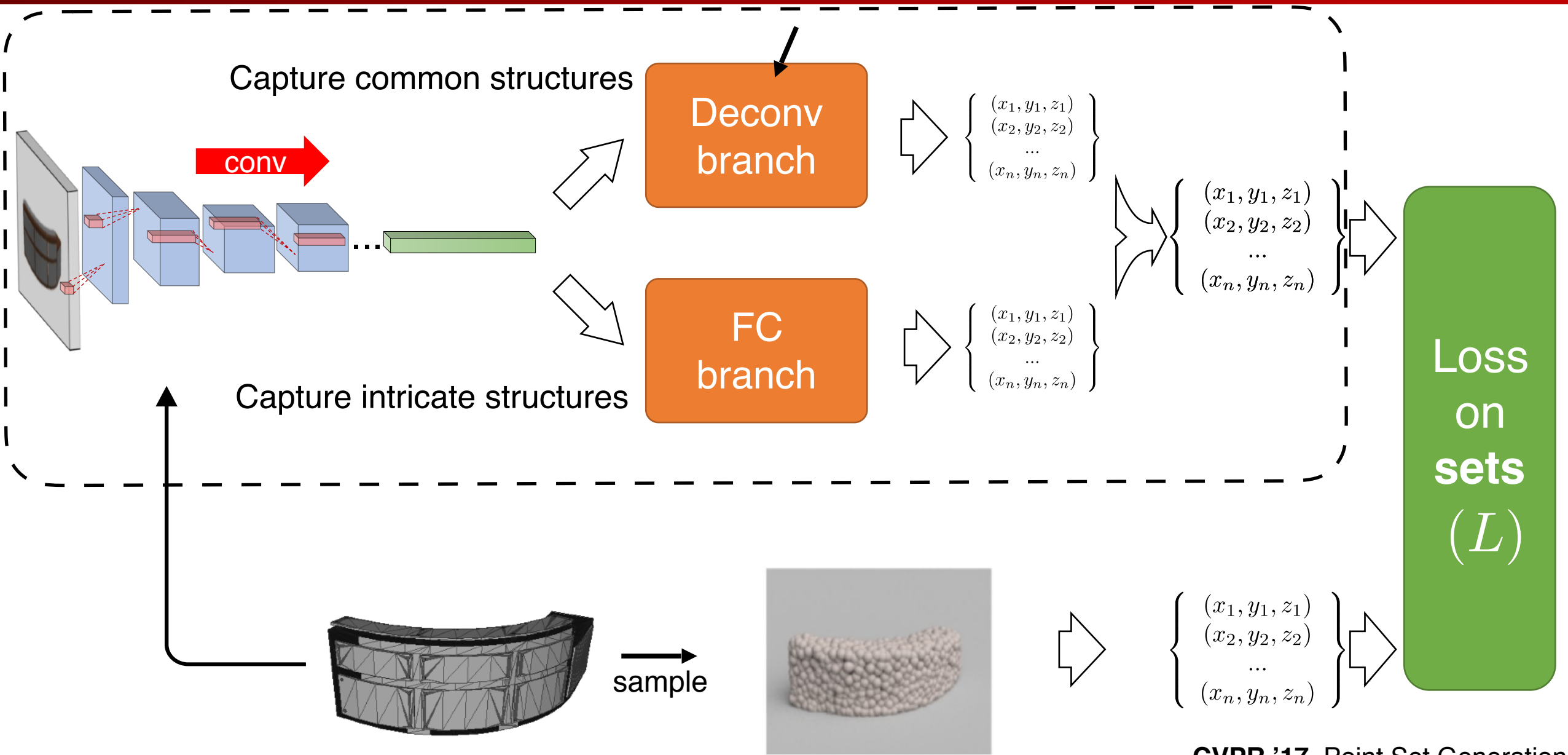


- Many local structures are common
 - e.g., planar patches, cylindrical patches
 - **strong local correlation** among point coordinates
- Also some intricate structures
 - points have **high local variation**

Pipeline

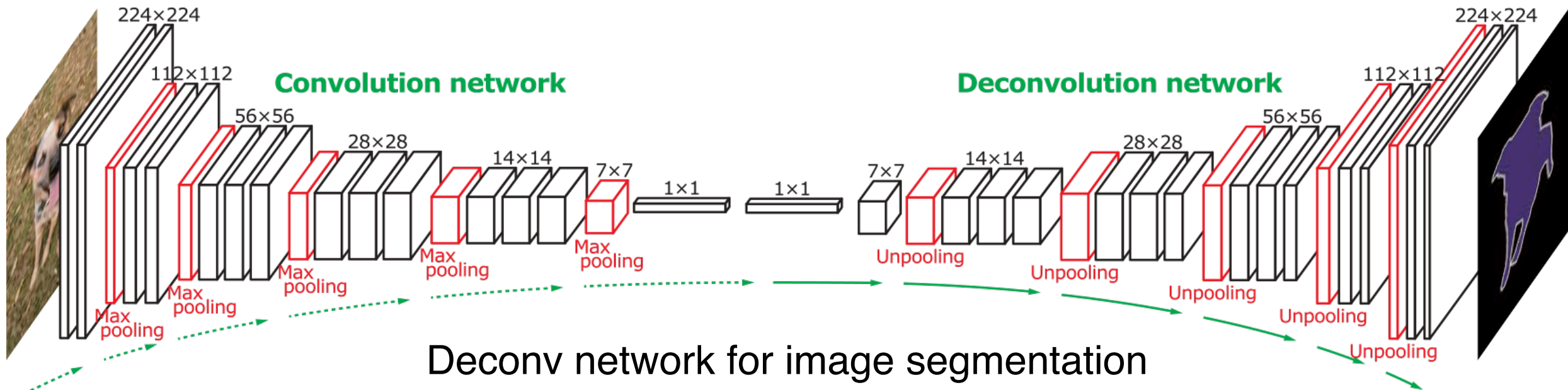


Pipeline



Review: deconv network

- Output n D arrays, e.g., 2D segmentation map
- **Common local patterns** are **learned from data**
- Predict **locally correlated** data well
- Weight sharing reduces the number of params



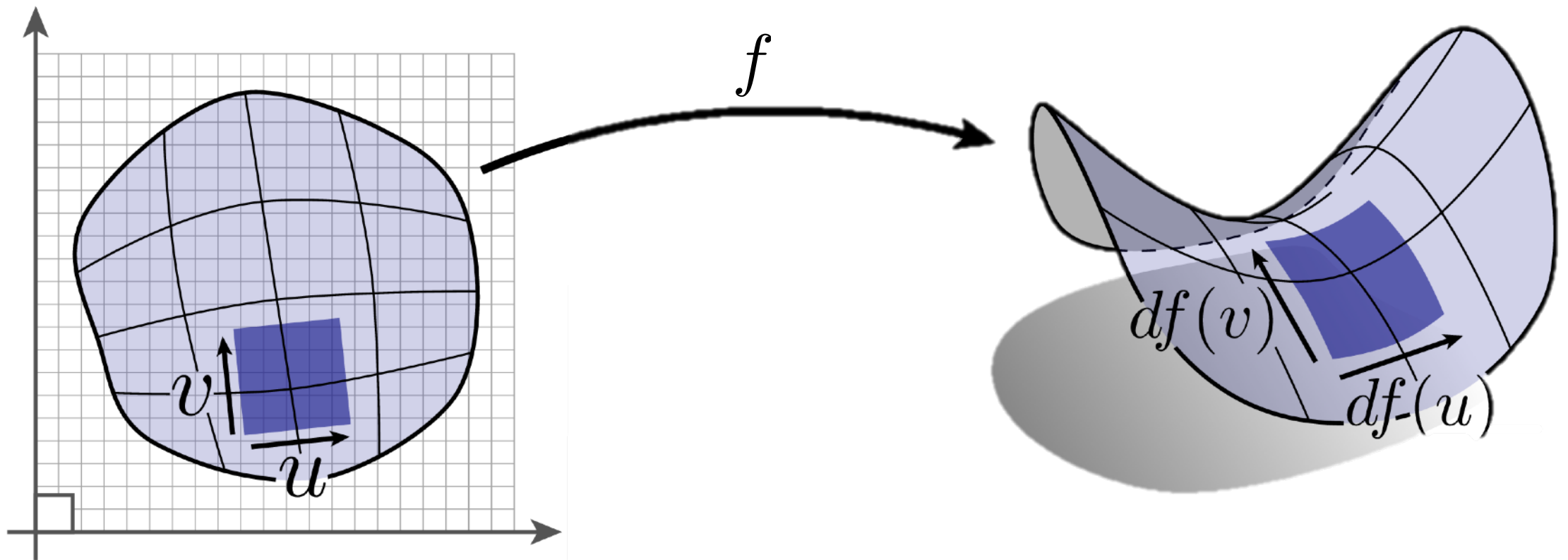
Review: deconv network

- Output n D arrays, e.g., 2D segmentation map
- **Common local patterns** are **learned from data**
- Predict **locally correlated** data well
- Weight sharing reduces the number of params



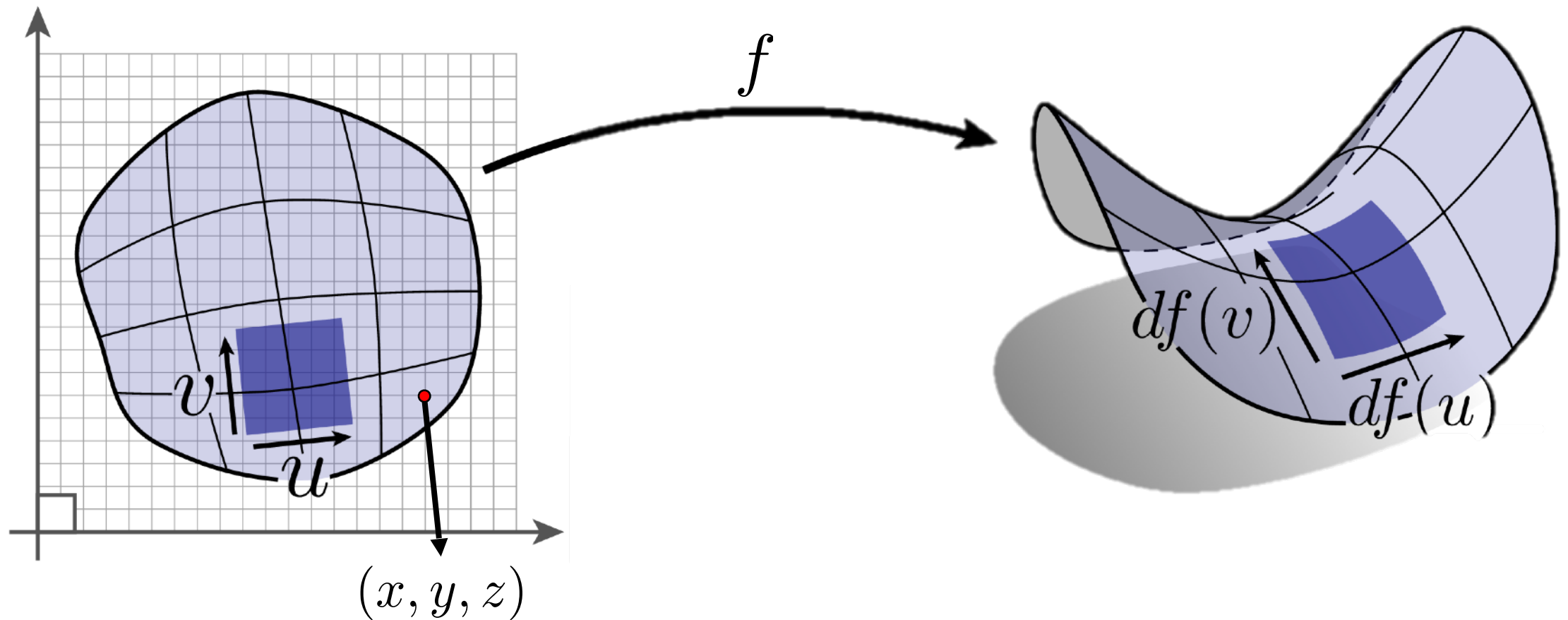
Prediction of curved 2D surfaces in 3D

- Surface parametrization ($2D \leftrightarrow 3D$ mapping)



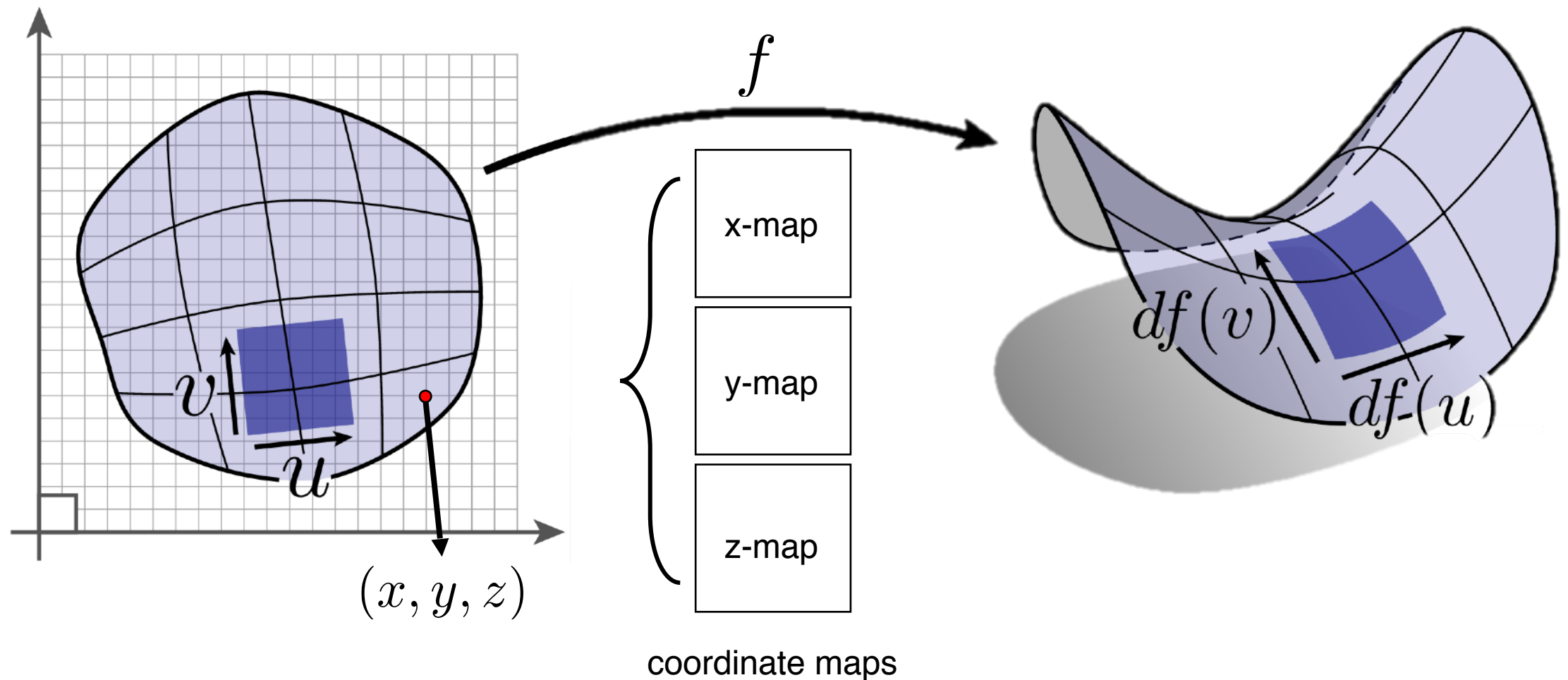
Prediction of curved 2D surfaces in 3D

- Surface parametrization ($2D \leftrightarrow 3D$ mapping)

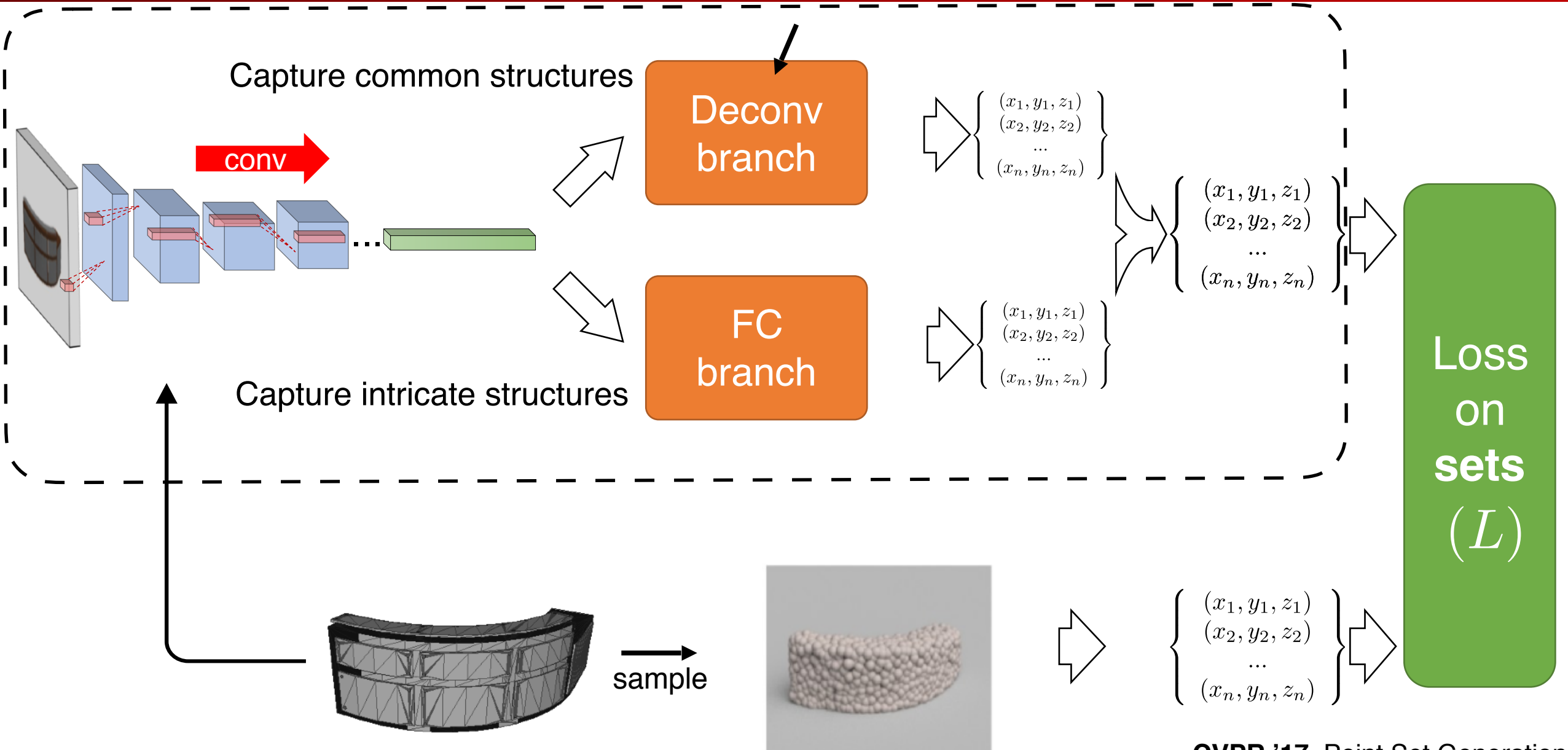


Prediction of curved 2D surfaces in 3D

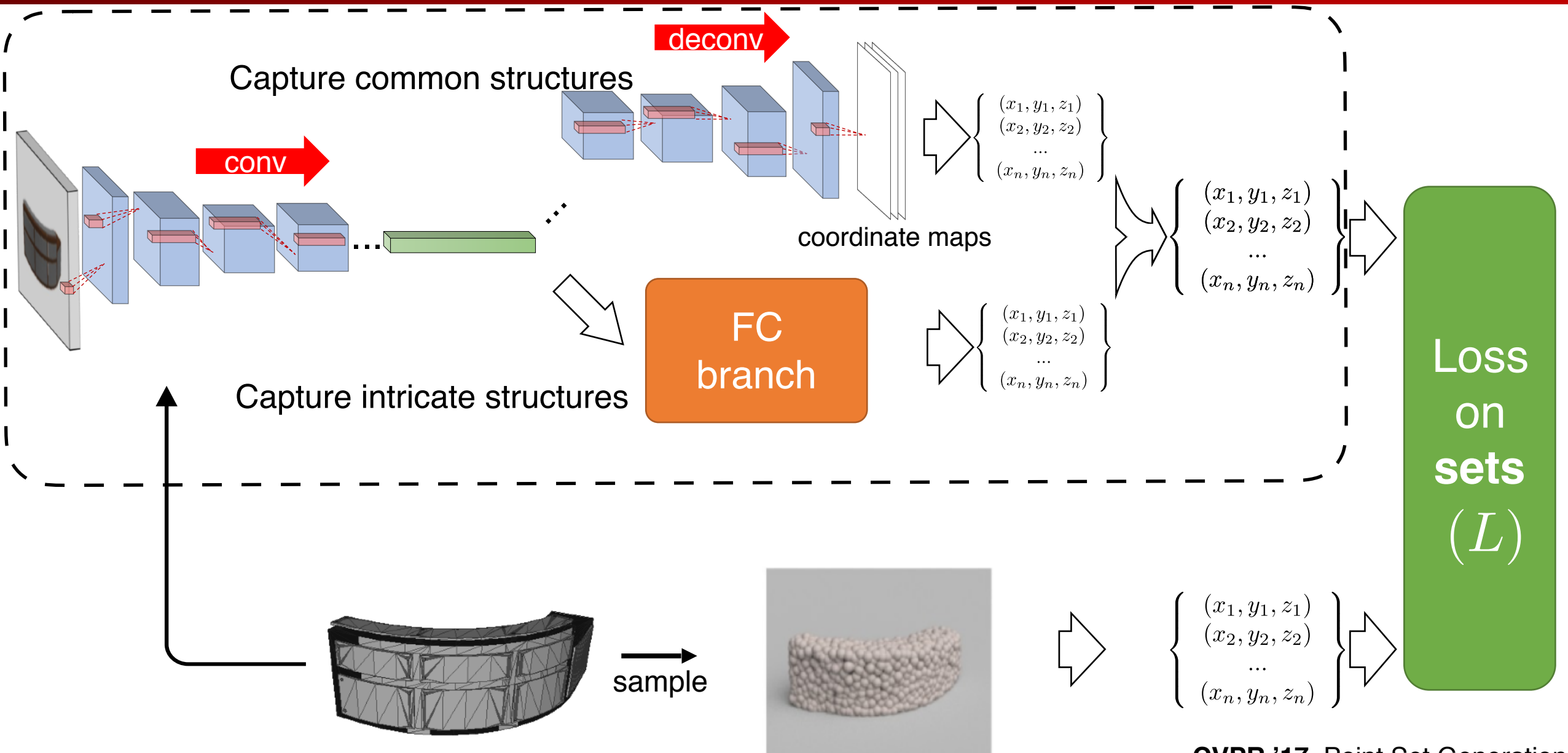
- Surface parametrization ($2D \leftrightarrow 3D$ mapping)



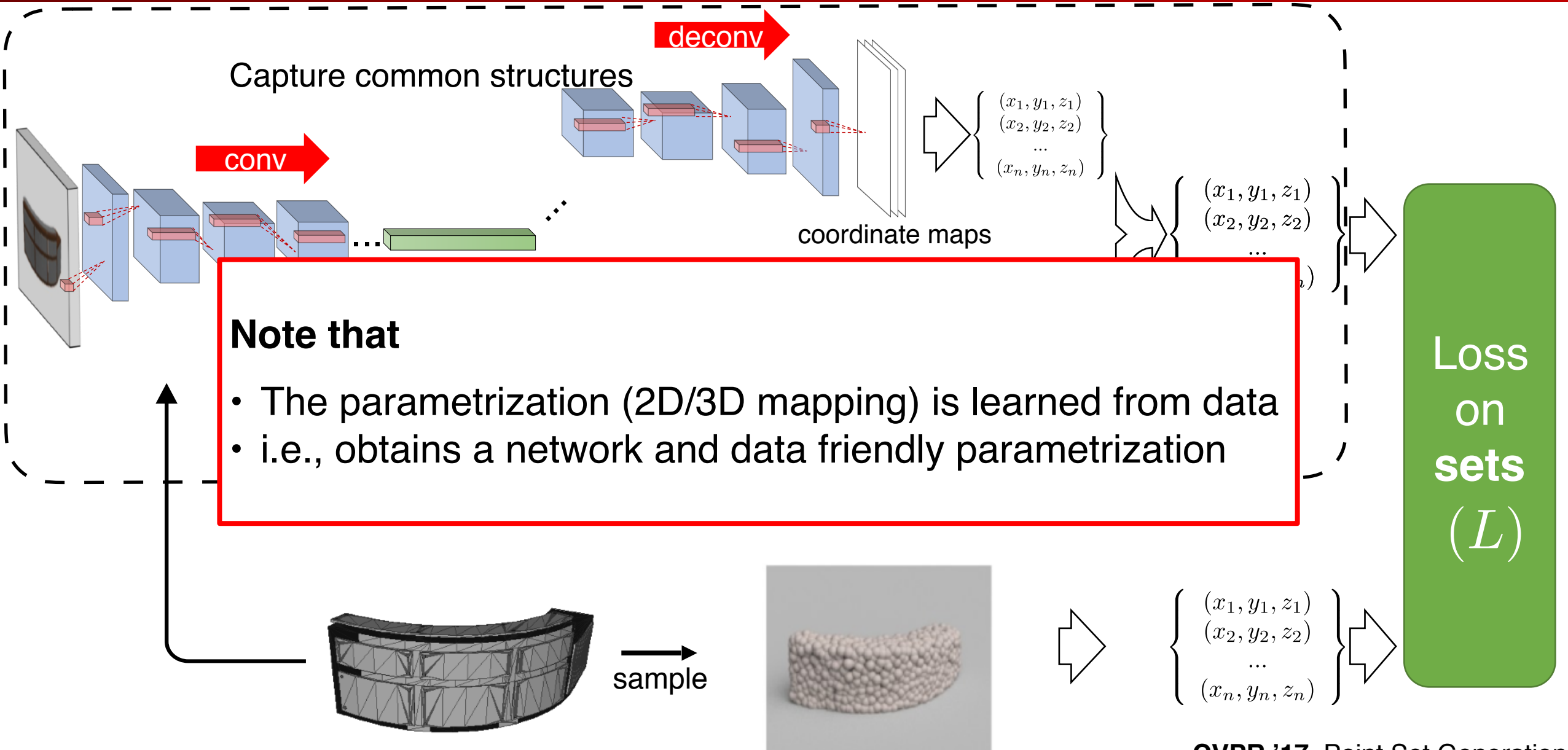
Parametrization prediction by deconv network



Parametrization prediction by deconv network



Parametrization prediction by deconv network

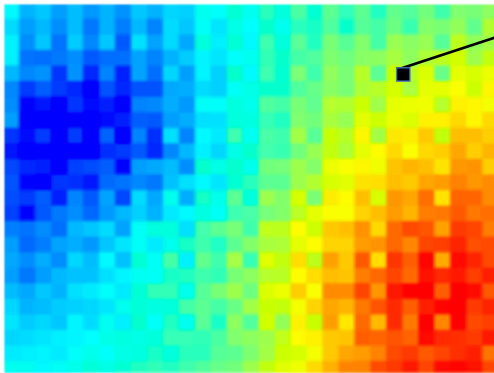
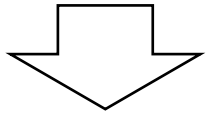


Visualization of the learned parameterization

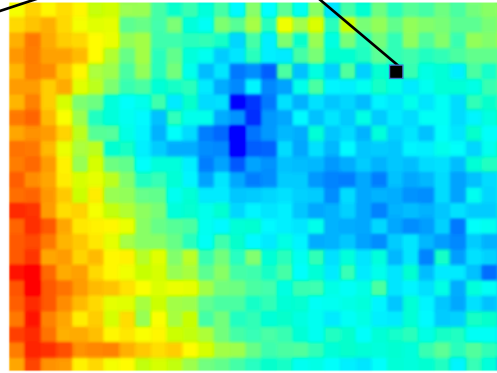


Observation:

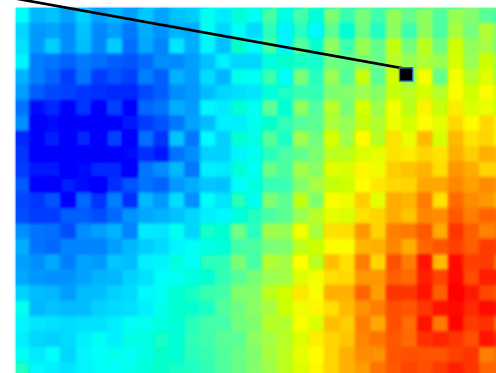
- Learns a **smooth** parametrization
- Because deconv net tends to predict data with local correlation



map of x coord



map of y coord



map of z coord

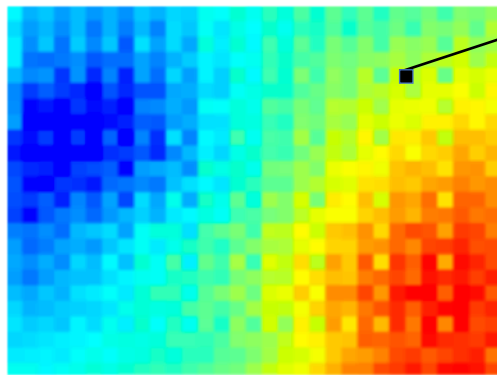
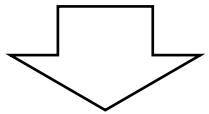
(x_k, y_k, z_k)

Visualization of the learned parameterization

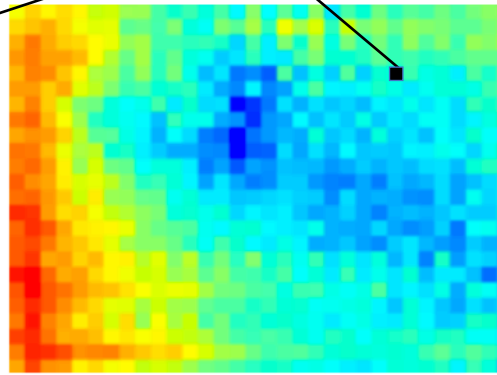


Observation:

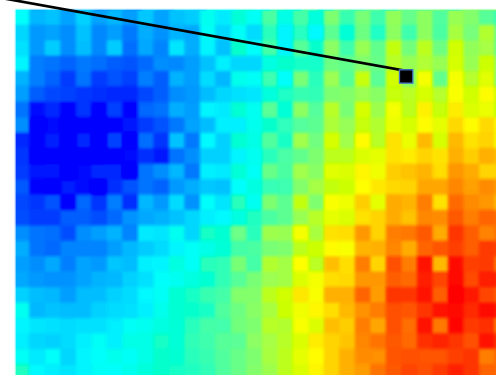
- Learns a **smooth** parametrization
- Because deconv net tends to predict data with local correlation
- Corresponds to **smooth surfaces!**



map of x coord

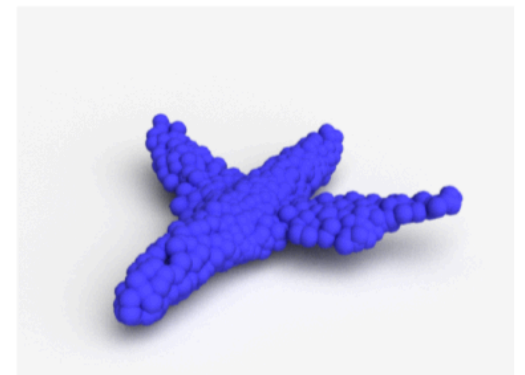
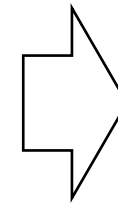


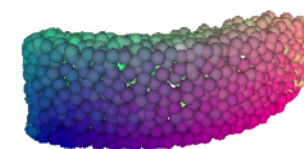
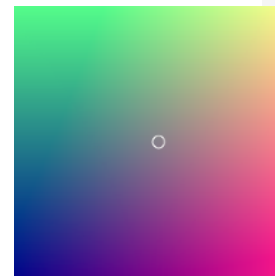
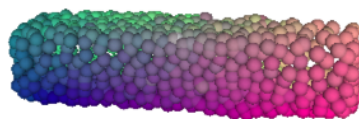
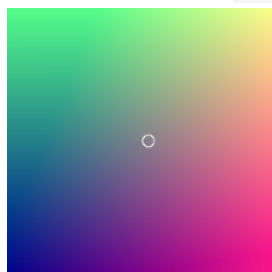
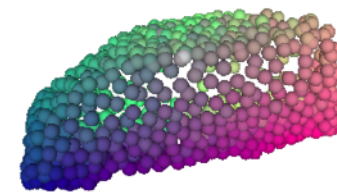
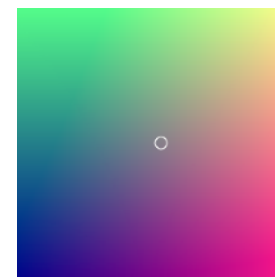
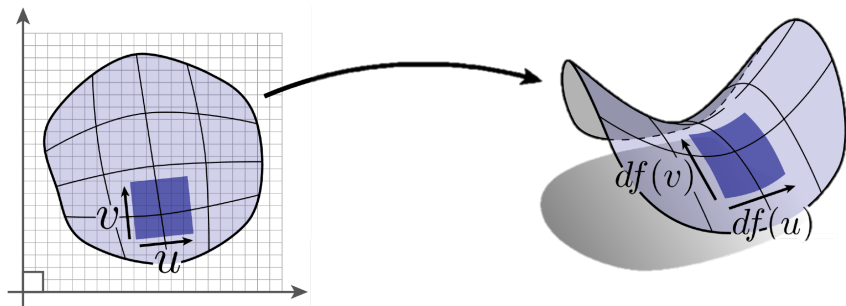
map of y coord



map of z coord

(x_k, y_k, z_k)



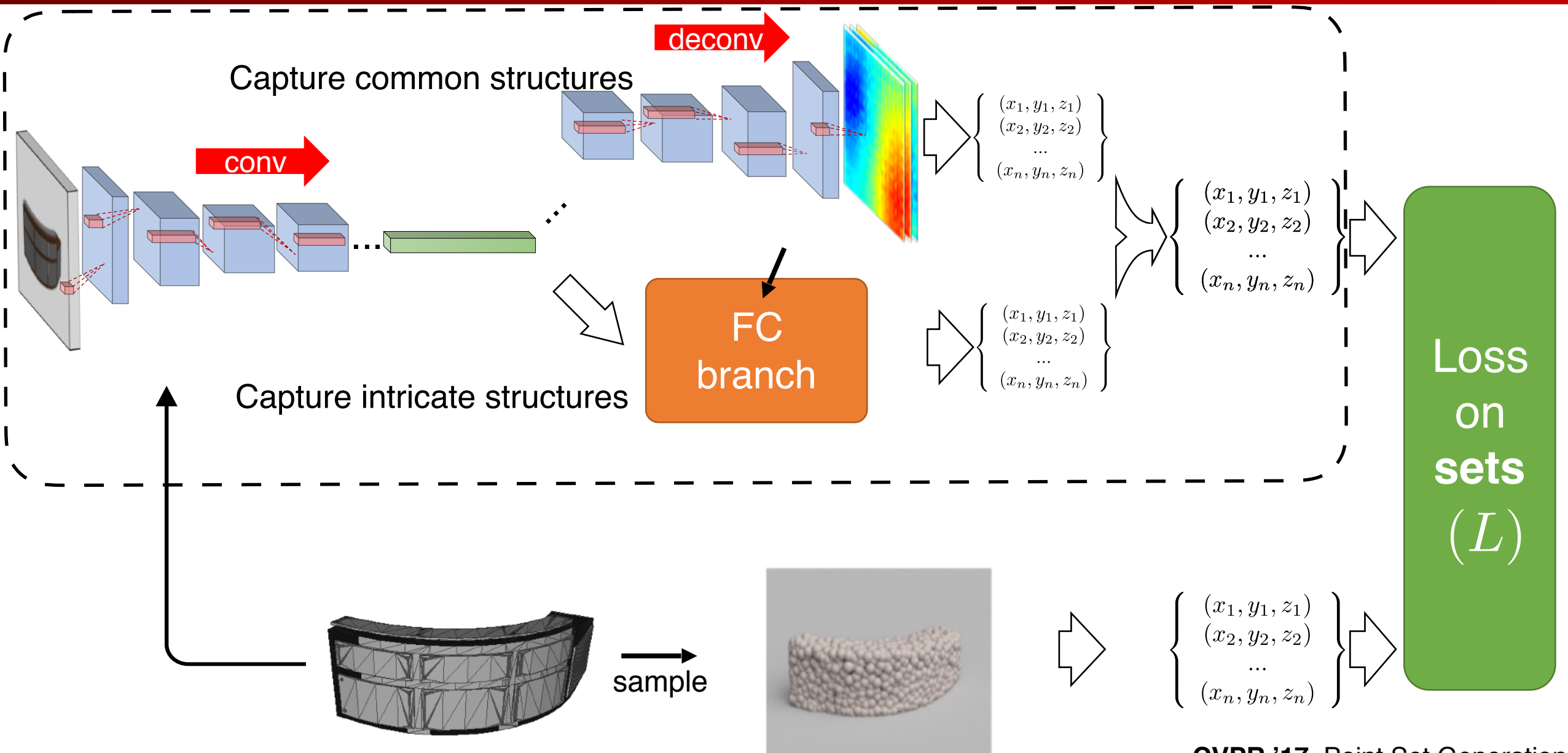


Natural statistics of geometry

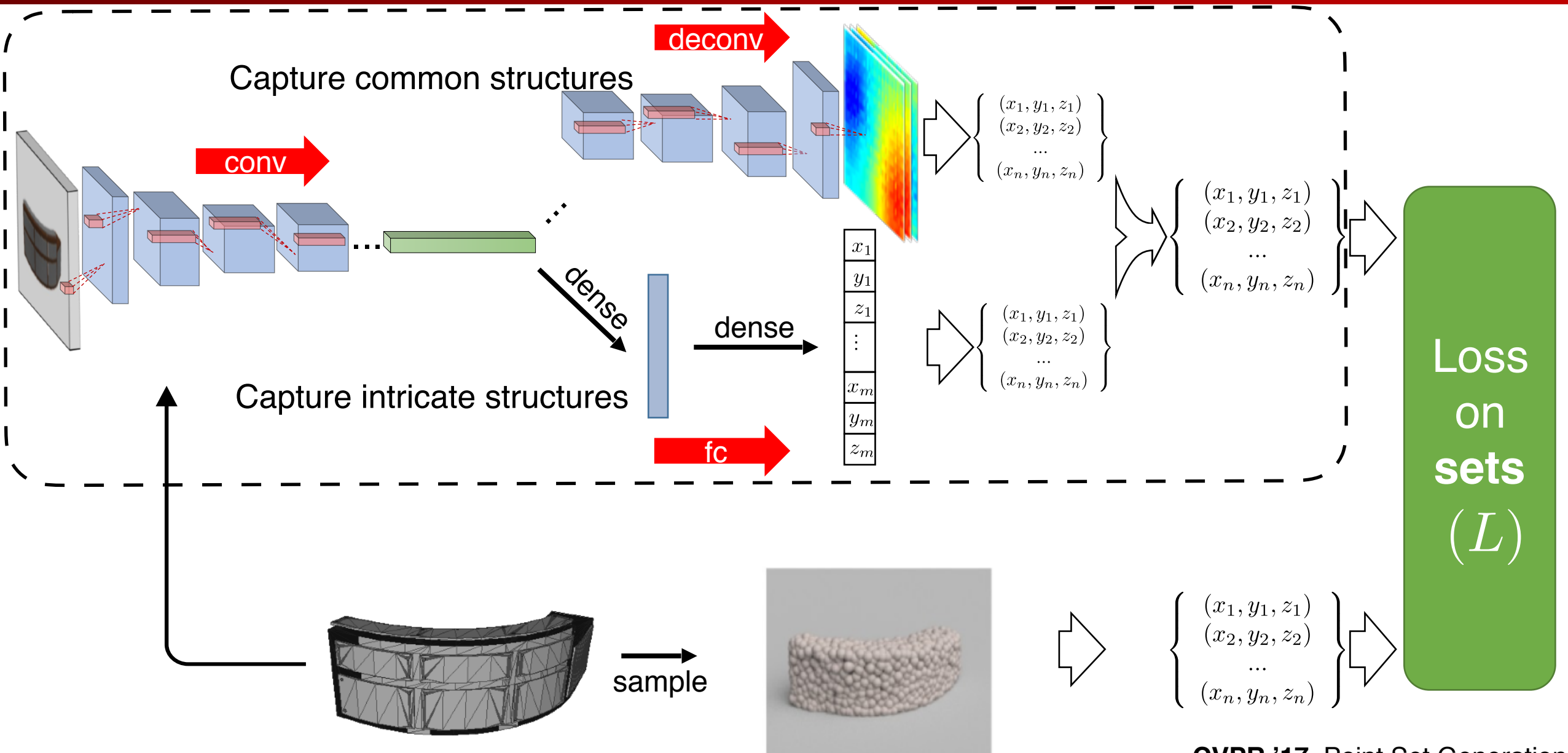


- Many local structures are common
 - e.g., planar patches, cylindrical patches
 - **strong local correlation** among point coordinates
- Also some intricate structures
 - points have **high local variation**

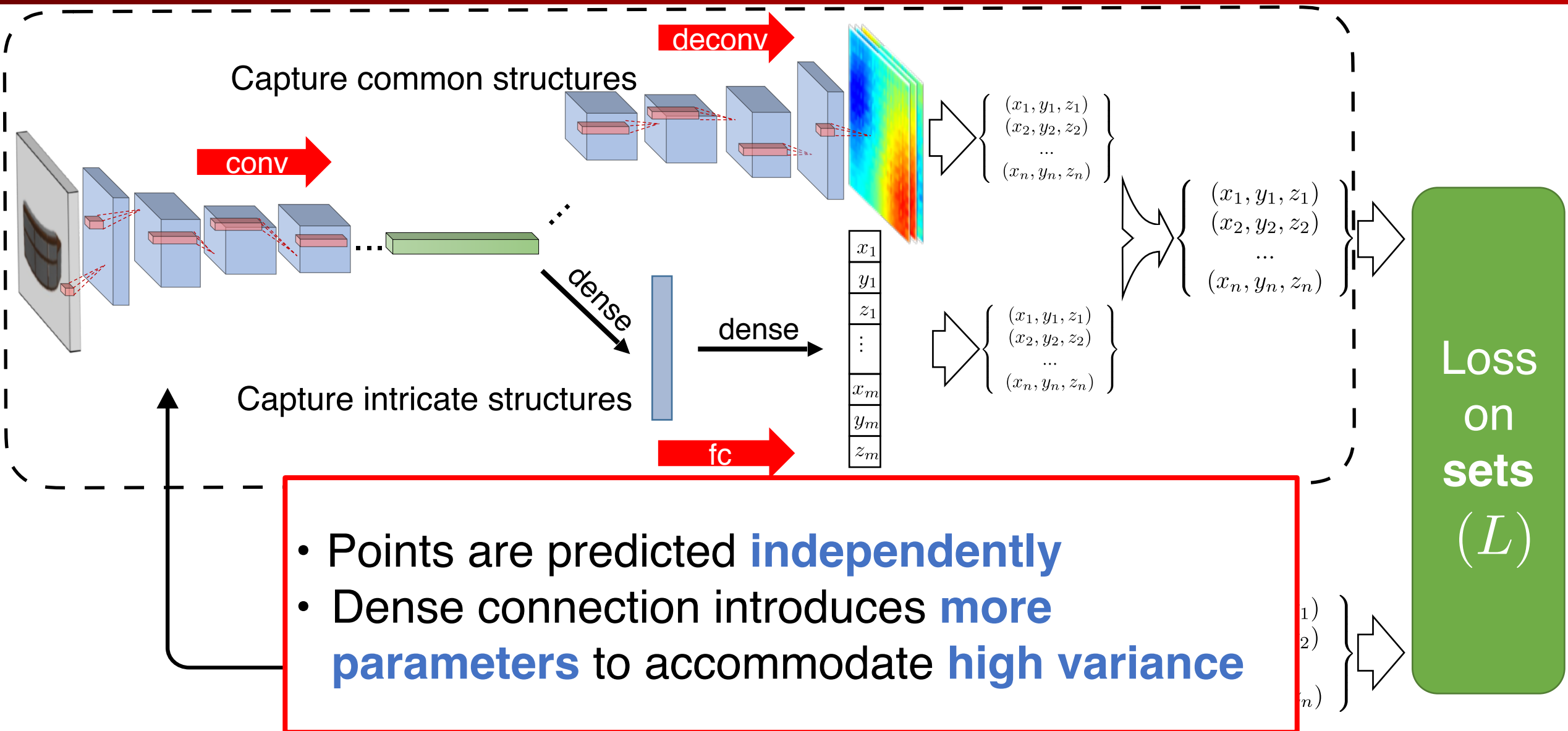
Pipeline



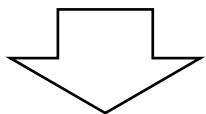
Pipeline



Pipeline

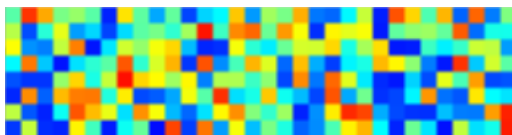


Visualization of the effect of FC branch

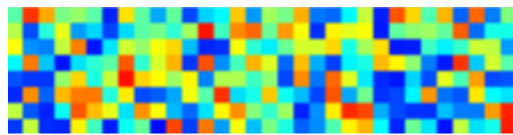


Observation:

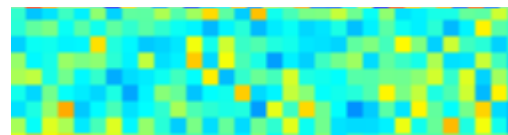
- The arrangement of predicted points are uncorrelated



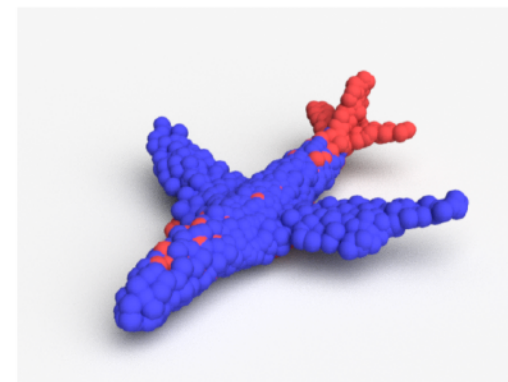
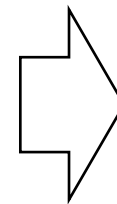
x-coord



y-coord

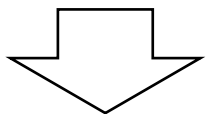


z-coord



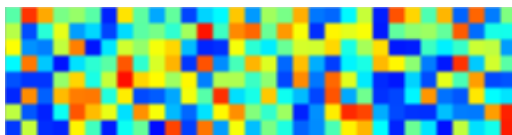
red

Visualization of the effect of FC branch

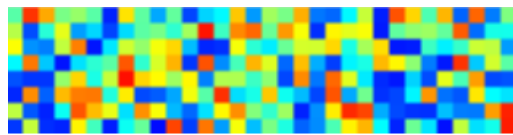


Observation:

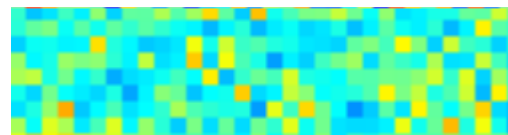
- The arrangement of predicted points are uncorrelated
- Located at **fine** structures



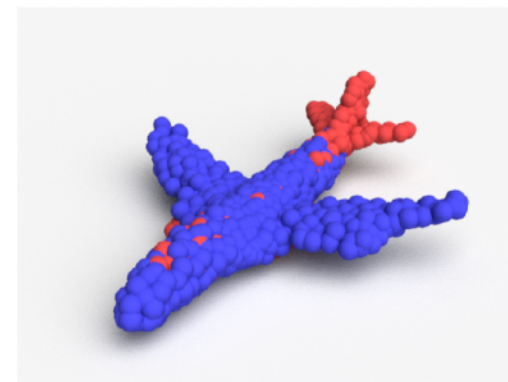
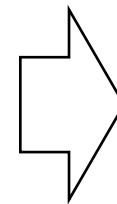
x-coord



y-coord

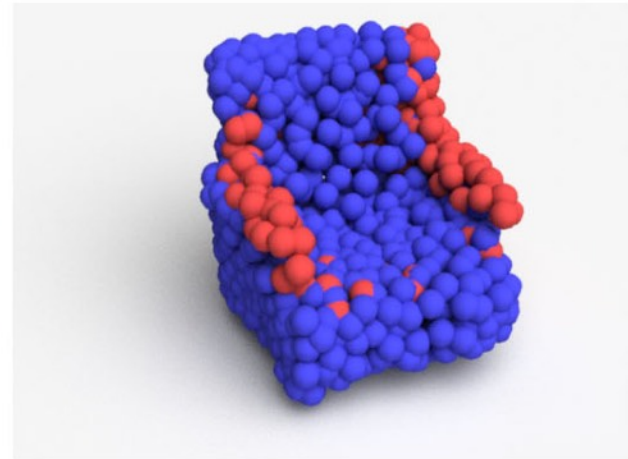
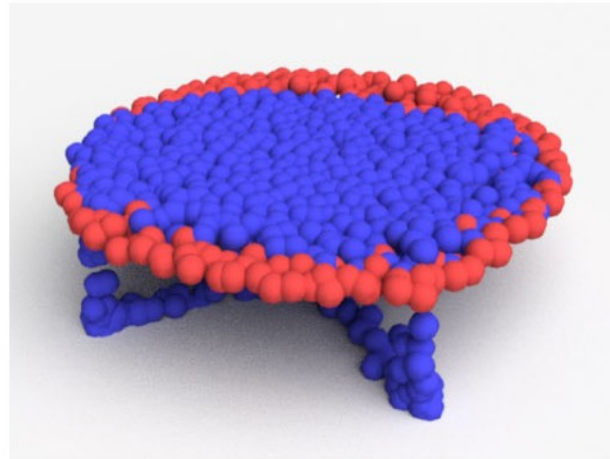
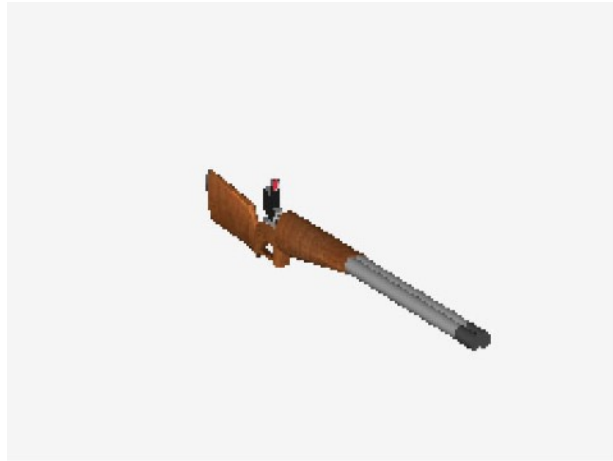


z-coord



red

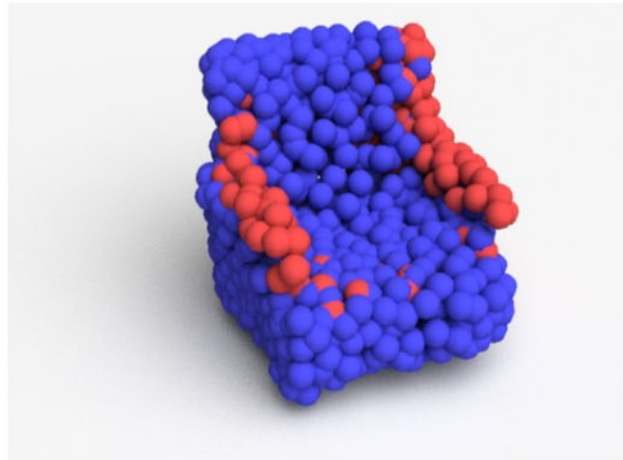
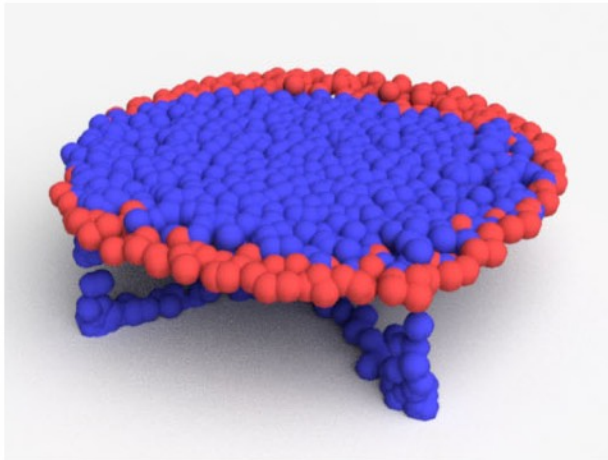
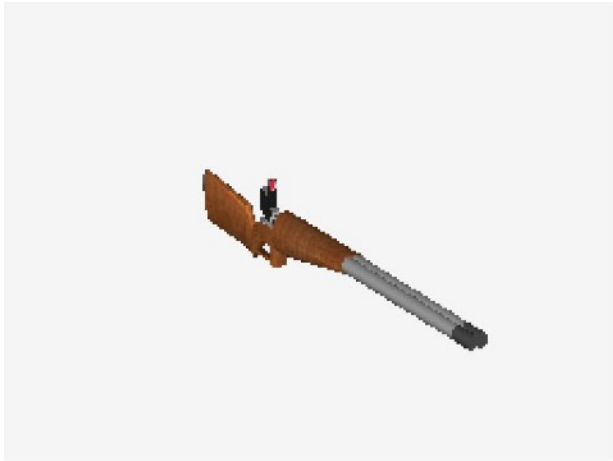
Q: Which color corresponds to the deconv branch? FC branch?



Q: Which color corresponds to the deconv branch? FC branch?

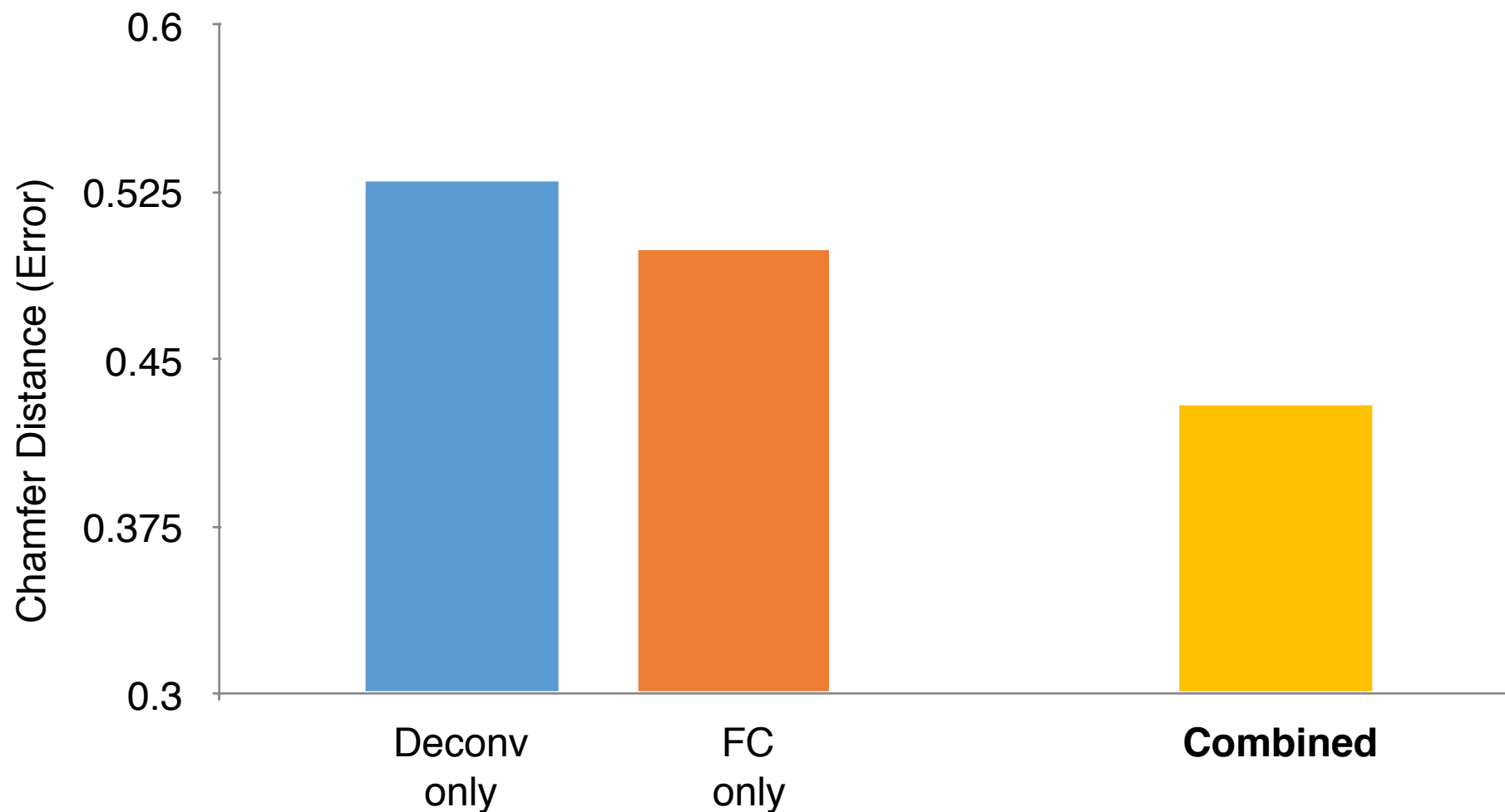
blue: deconv branch – **large, smooth** structures

red: FC branch – **intricate** structures



Effect of combining two branches

Train/tested on 2K object categories

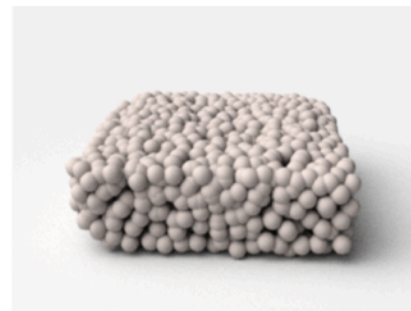


Real-world results

input

observed view

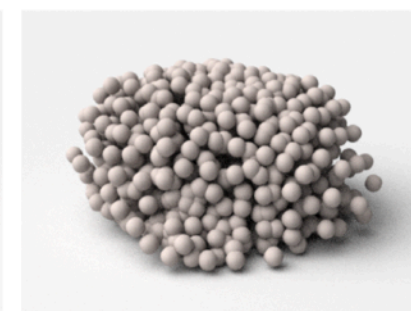
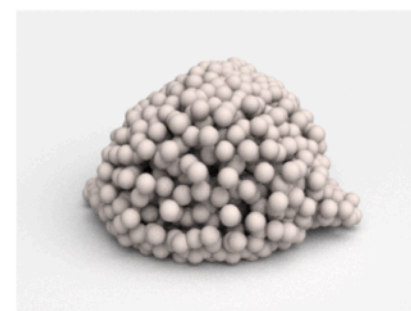
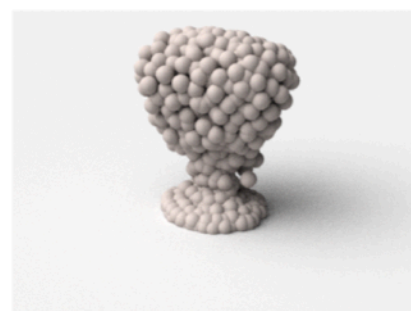
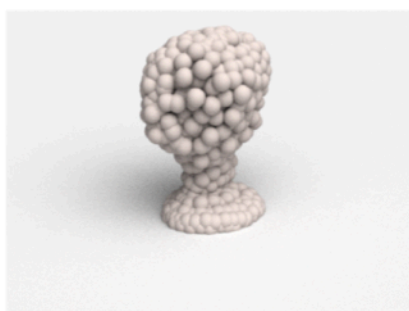
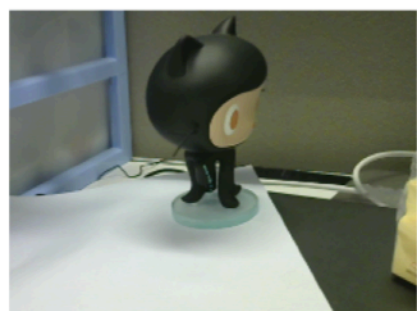
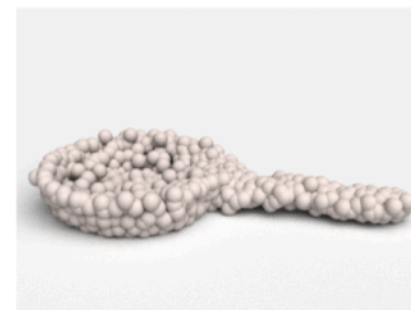
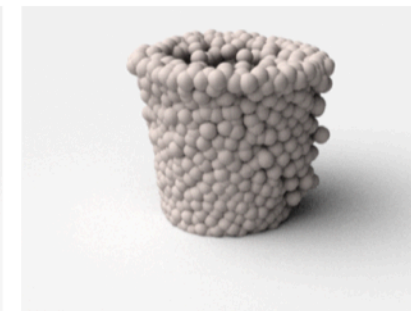
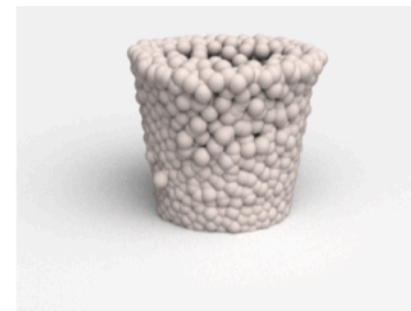
90°



input

observed view

90°



Generalization to unseen categories

input

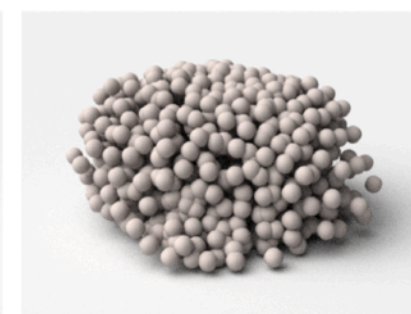
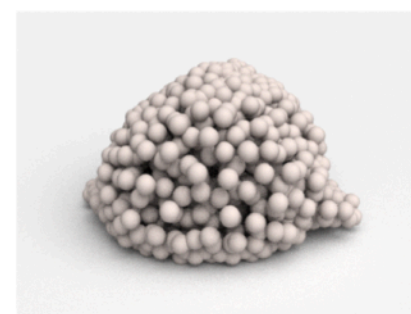
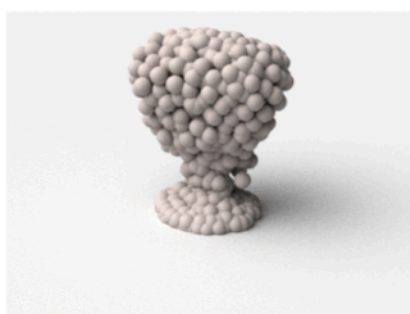
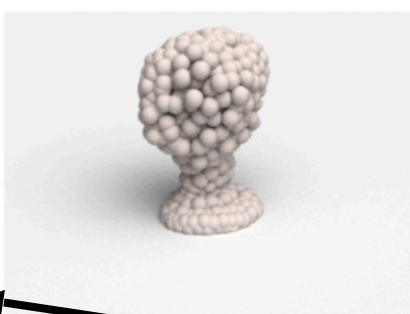
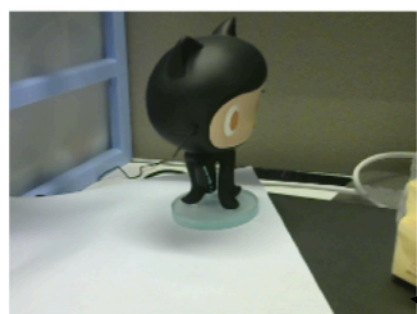
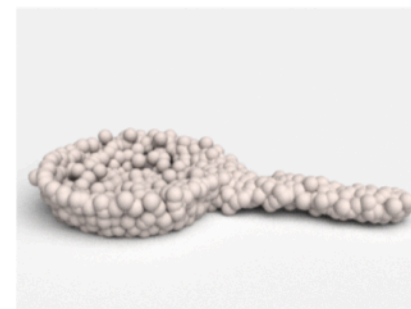
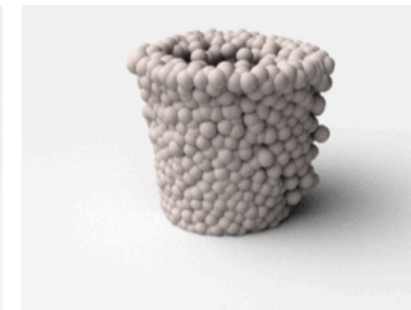
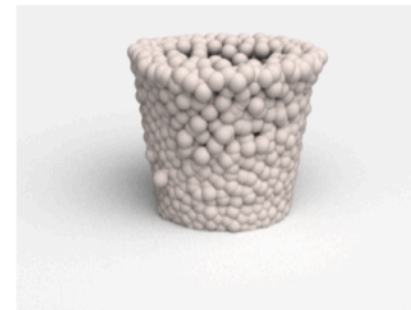
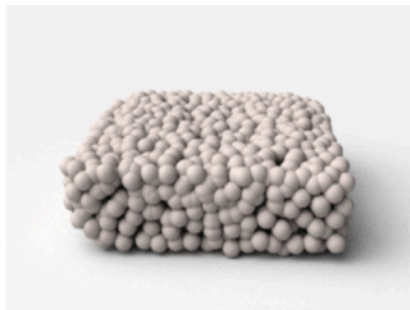
observed view

90°

input

observed view

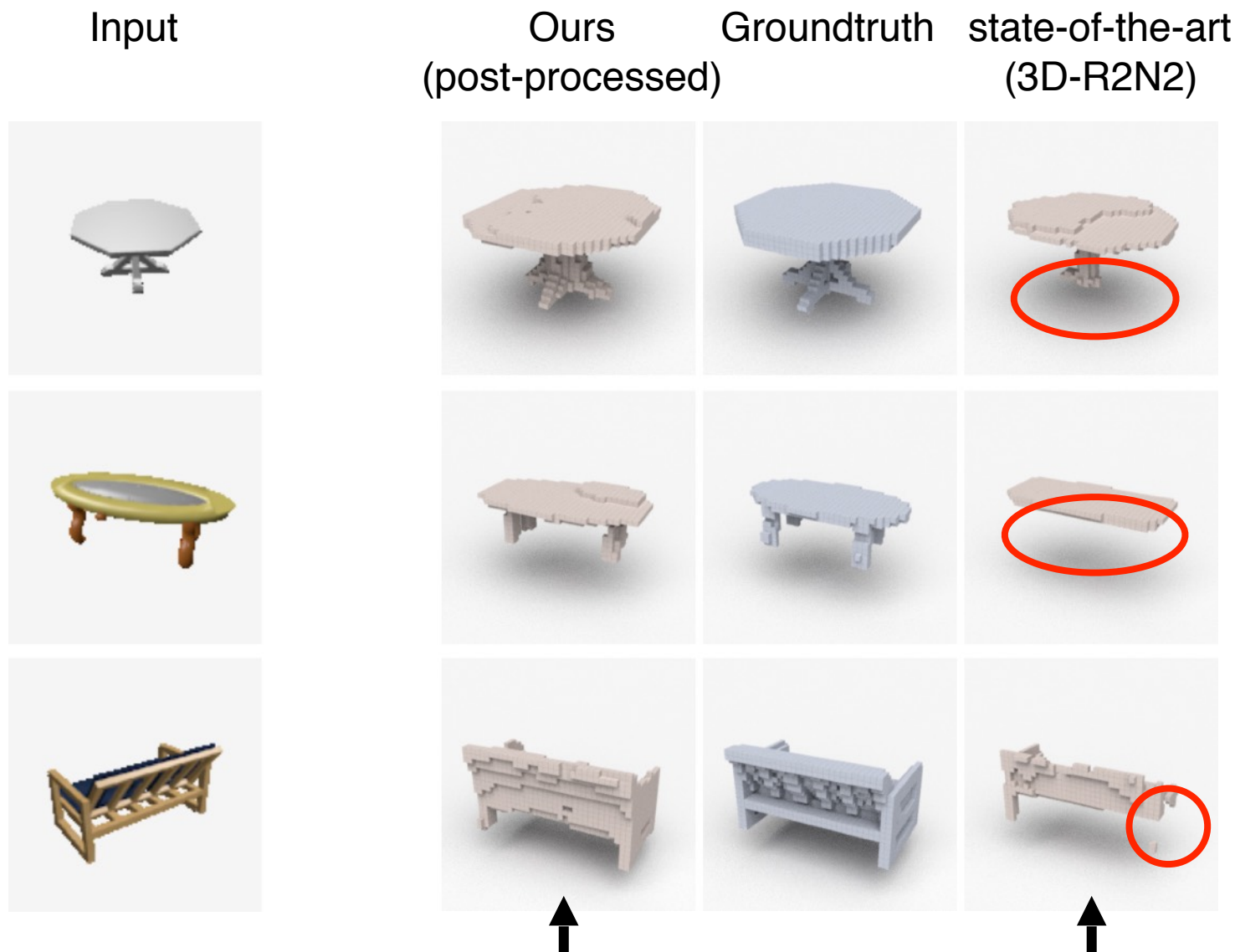
90°



Out of training categories

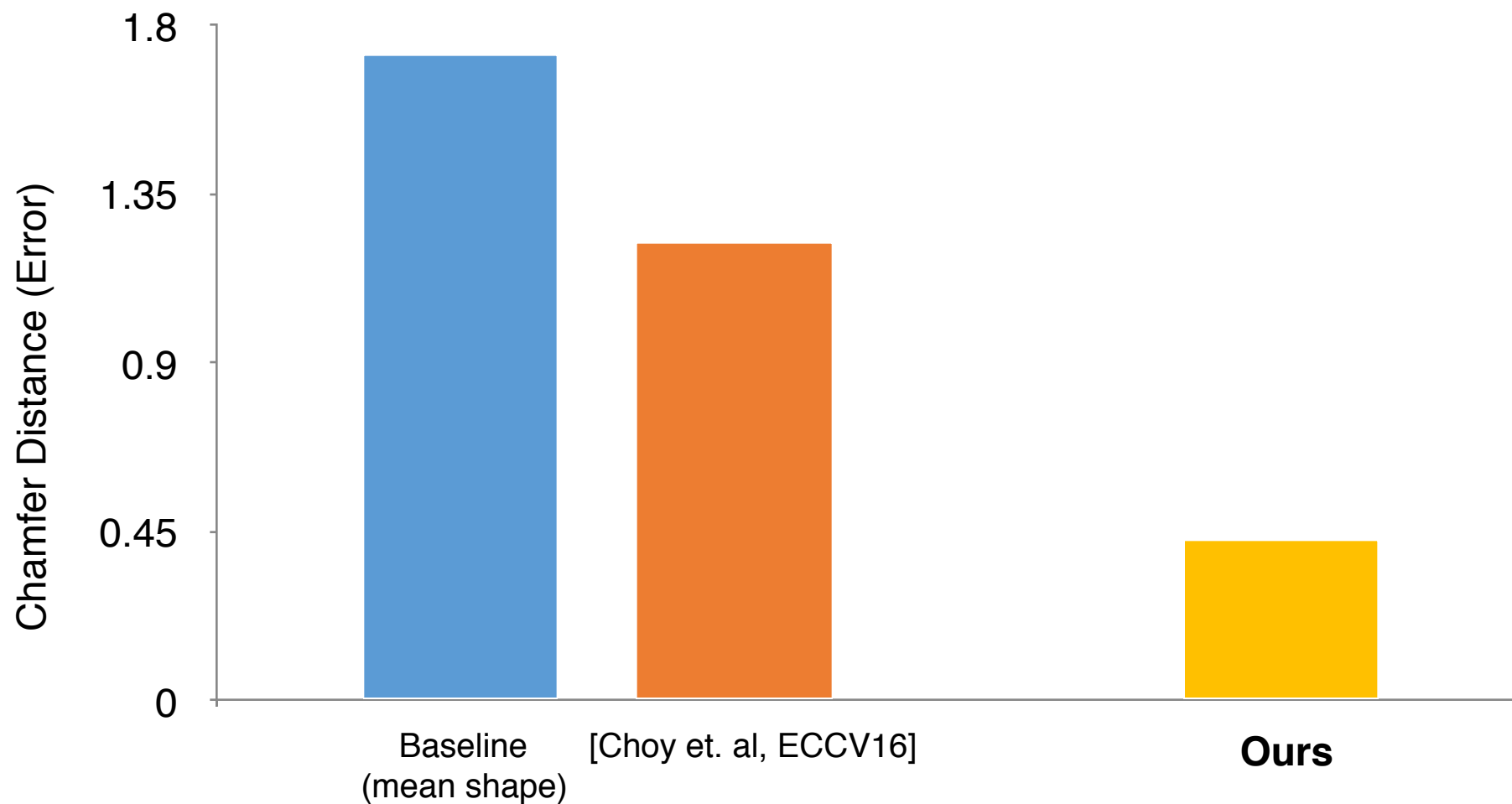
Comparison to state-of-the-art

- Better global structure
- Better details

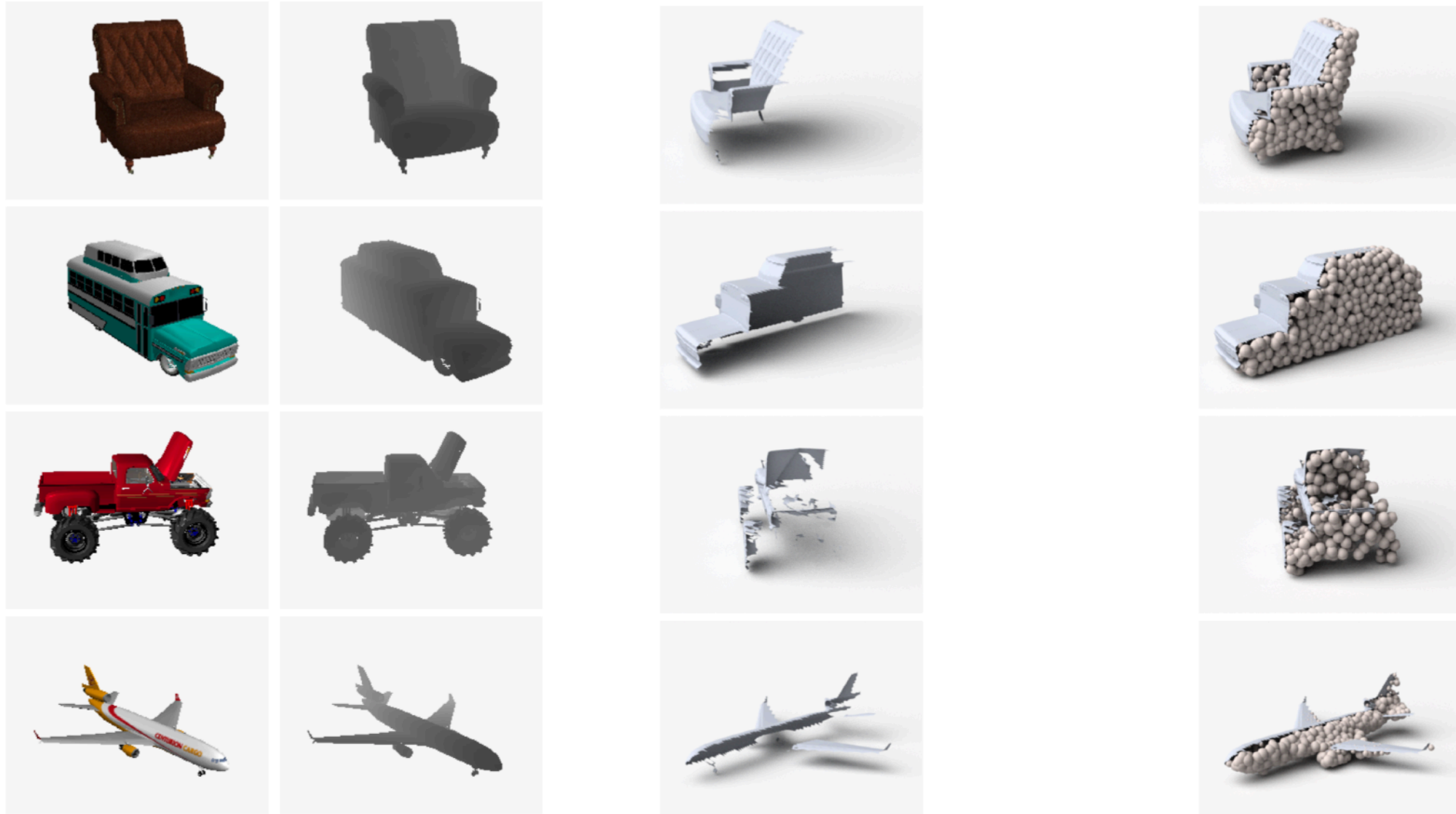


Comparison to state-of-the-art

Trained/tested on 2K object categories



Extension: shape completion for RGBD data



RGBD map (input)

90° view of input

output: completed point cloud

Open problems

A better metric that takes the best of Chamfer and EMD?

How to add further structure constraints?

How to extend the pipeline to scene level?

How generalizable the method is?

In principle, what is the generalizability of a geometry estimator? To what extent is 3D perception ability innate or learned?

Outline

- Motivation
- Data
- Algorithms for 2D-3D lifting
- **Shape abstraction by volumetric primitives**

How about learning to predict geometric forms?

Candidates:

**Rasterized form
(regular grids)**

multi-view images

depth map

volumetric

**Geometric form
(irregular)**

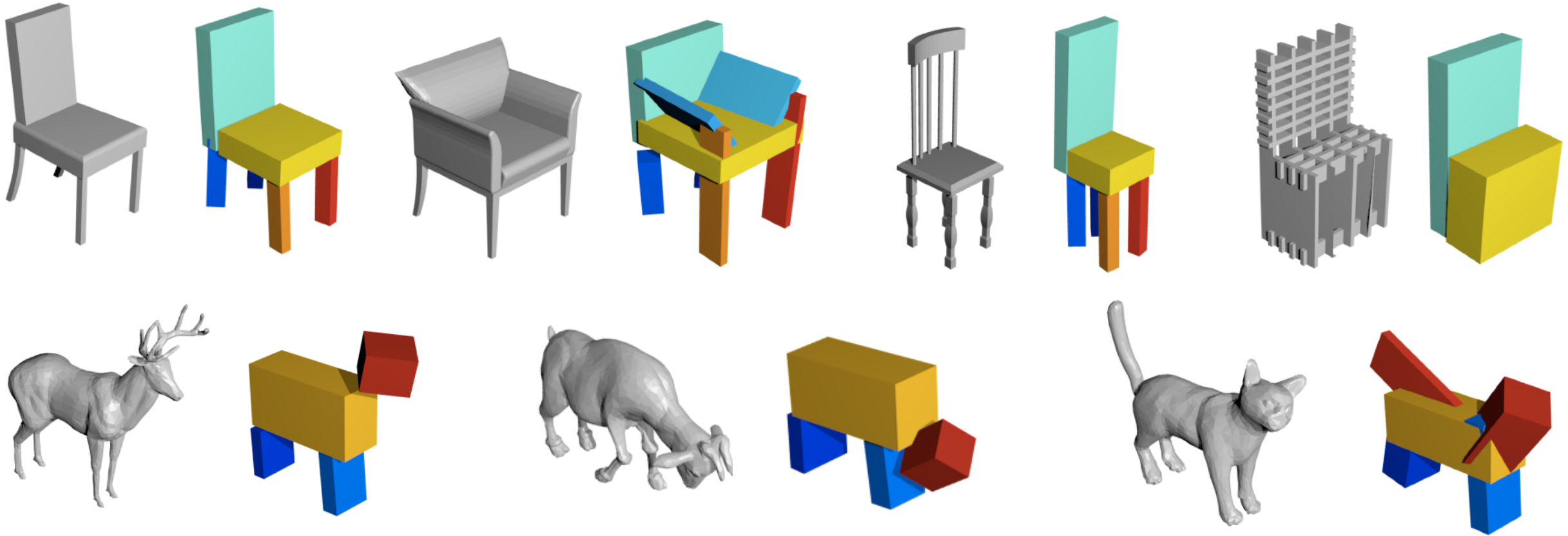
polygonal mesh

point cloud



primitive-based CAD models

Primitive-based assembly



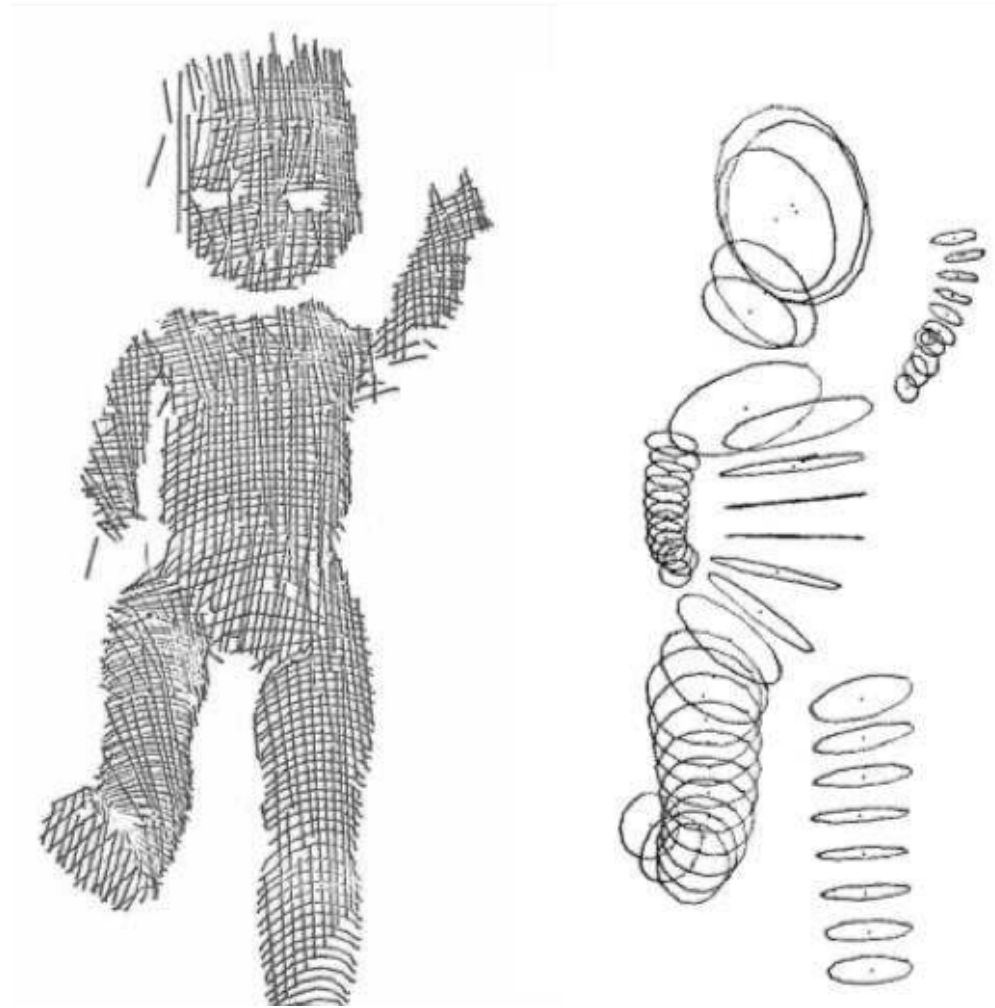
We **learn** to predict a corresponding shape composed by primitives. It allows us to predict **consistent** compositions across objects.

Unsupervised parsing



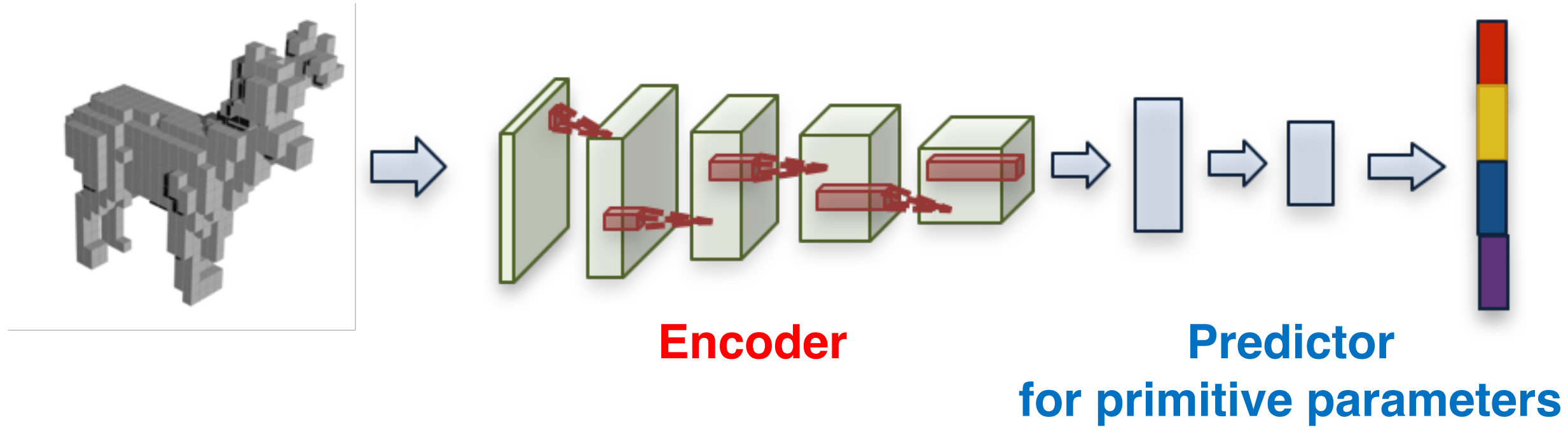
Each point is colored according to the assigned primitive

A historical overview



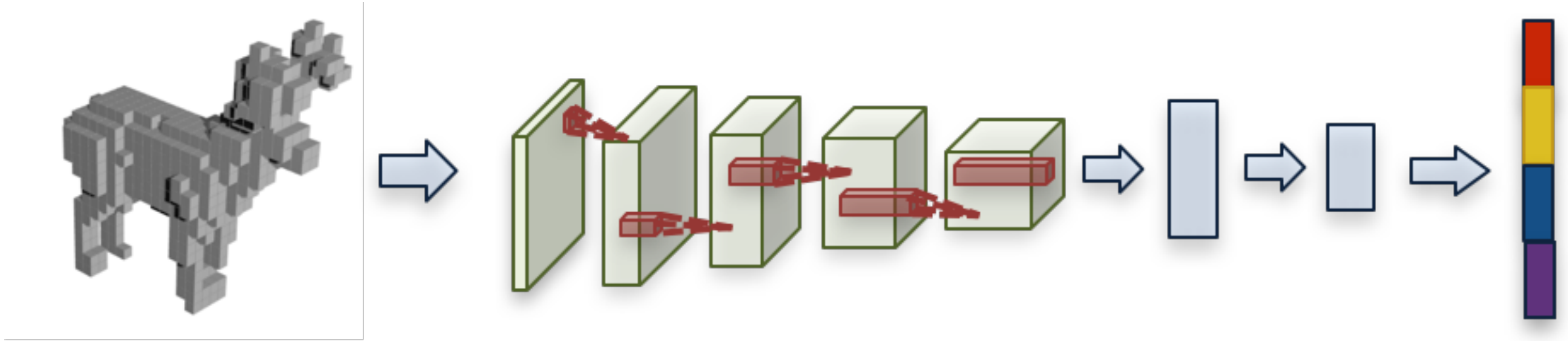
Generalized Cylinders, Binfond (1971)

Approach



We predict primitive parameters: size, rotation, translation of M cuboids.

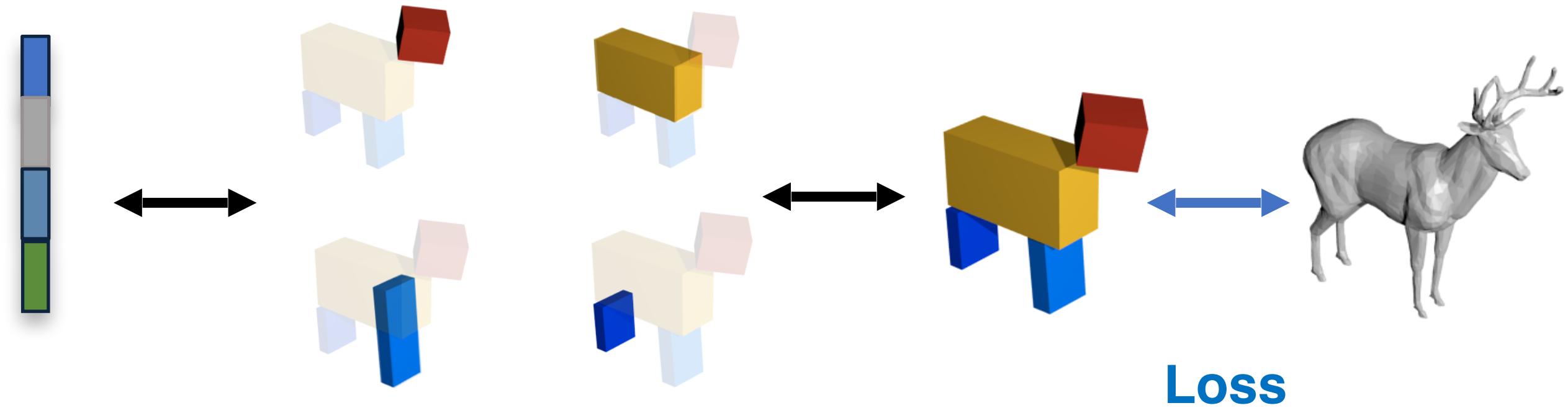
Approach



We predict primitive parameters: size, rotation, translation of M cuboids.

Variable number of parts? We predict “primitive existence probability”

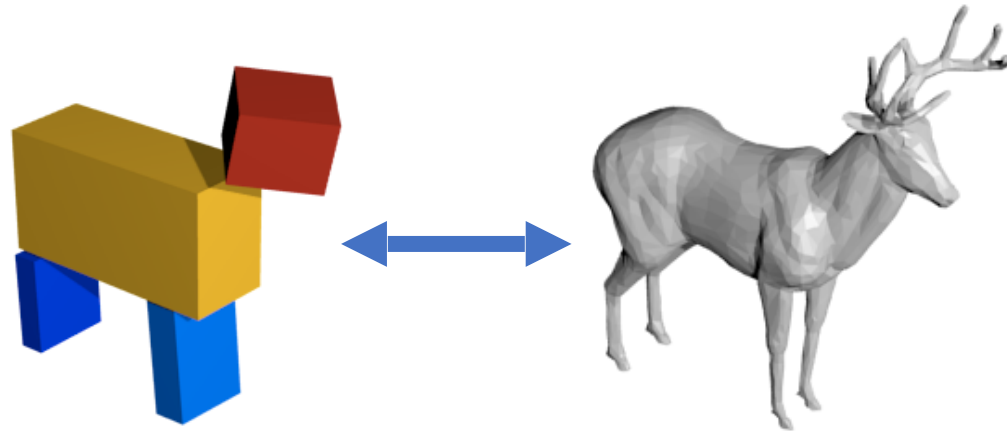
Loss function



Loss function construction

Basic idea: **Chamfer distance!**

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$



Loss function construction

Sample points on the groundtruth mesh and predicted assembly

$$\Delta(\text{deer}, \text{assembly}) + \Delta(\text{deer}, \text{assembly}) + \Delta(\text{deer}, \text{assembly}) \dots + \Delta(\text{deer}, \text{assembly})$$

Each point is a **linear function** of mesh/primitive vertex coordinates

Differentiable!

Loss function construction

Sample points on the groundtruth mesh and predicted assembly

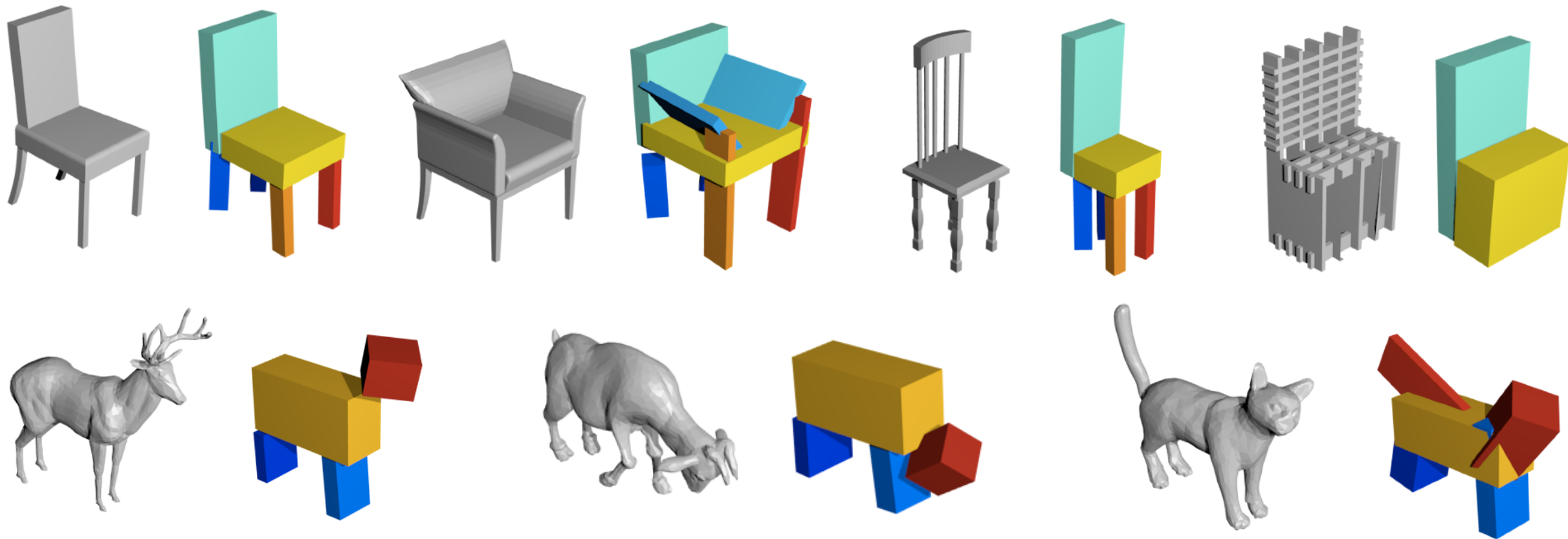
$$\Delta(\text{deer}, \text{assembly}) + \Delta(\text{deer}, \text{assembly}) + \Delta(\text{deer}, \text{assembly}) \dots + \Delta(\text{deer}, \text{assembly})$$

Each point is a **linear function** of mesh/primitive vertex coordinates

Differentiable!

Speed up the computation leveraging parameterization of primitives

Consistent primitive configurations



Primitive locations are **consistent** due to the **smoothness** of primitive prediction network

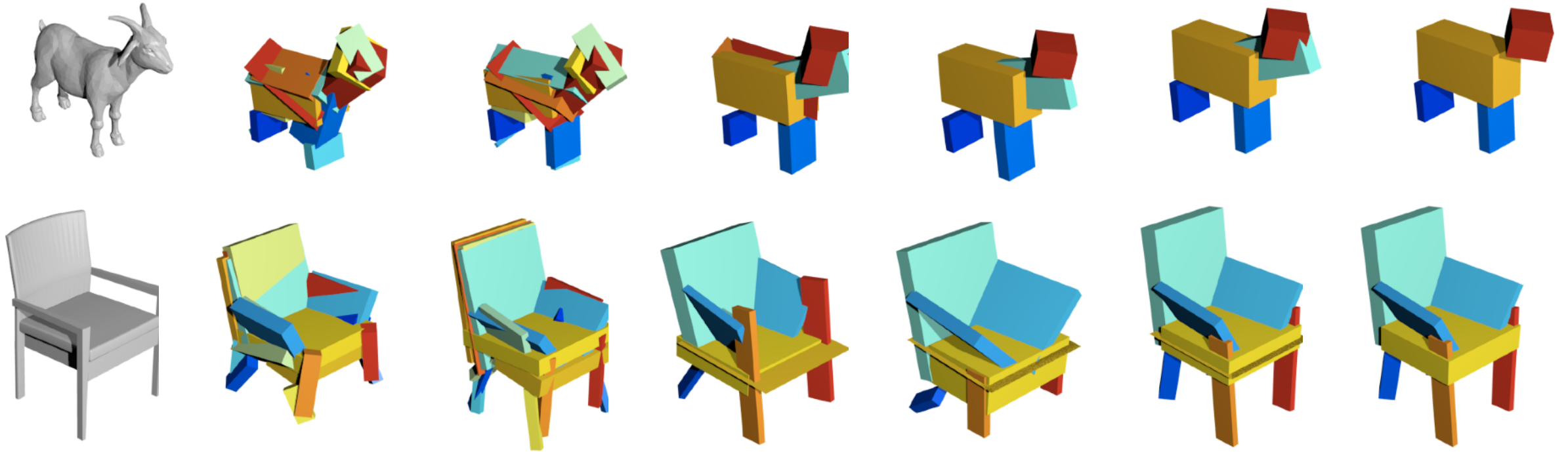
Unsupervised parsing



Method	[31] (initial)	[31] (refined)	Ours
Accuracy	78.6	84.8	89.0

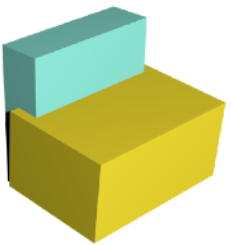
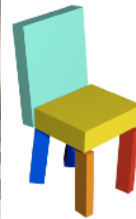
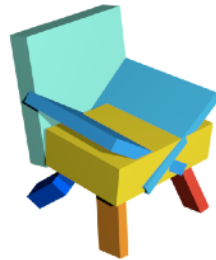
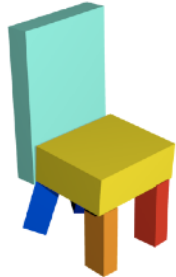
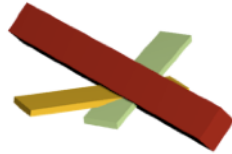
Mean accuracy (face area) on Shape COSEG chairs.

Analysis



Shapes become more parsimonious as training progresses (due to our parsimony reward)

Image-based modeling



Open problems

How to introduce other primitives types?

Towards image based modeling, how to add more operations to edit those primitives?

- e.g., Deform? Extrude? Loop cut?

How to use it for design purposes? For example, with certain structural and functional constraints.

Ultimately, we expect to automate the modeling process from images, as artists do.

Resources

- For both works, paper and source codes are available.
- Haoqiang Fan*, Hao Su*, Leonidas Guibas, *A Point Set Generation Network for 3D Object Reconstruction from a Single Image*, CVPR2017 (oral)
- Shubham Tulsiani, Hao Su, Leonidas Guibas, Alexei Efros, Jitendra Malik, *Learning Shape Abstractions by Assembling Volumetric Primitives*, CVPR2017 (oral)

To sum up

We explore to generate geometric representations by neural networks

Point cloud based reconstruction has better quality than state-of-the-art volumetric shape generators

Primitive-based CAD models can be generated to abstract polygonal meshes in an unsupervised manner

Keys:

- network structure leveraging geometric natural statistics
- loss function design