# CS468: 3D Deep Learning on Point Cloud Data
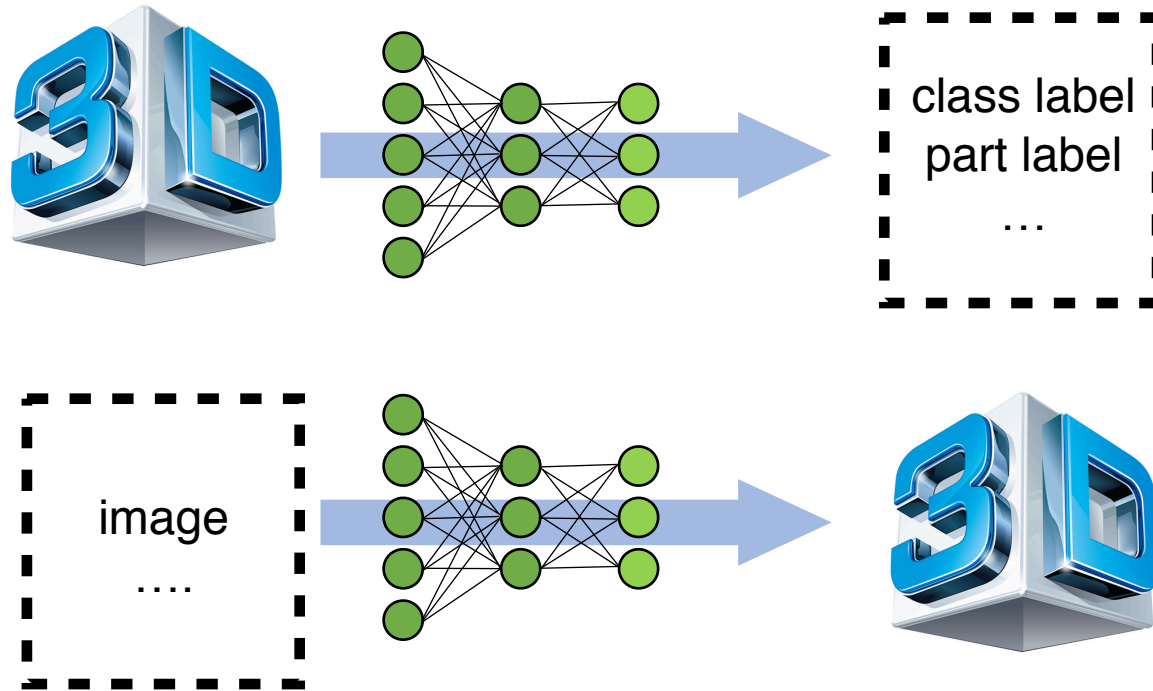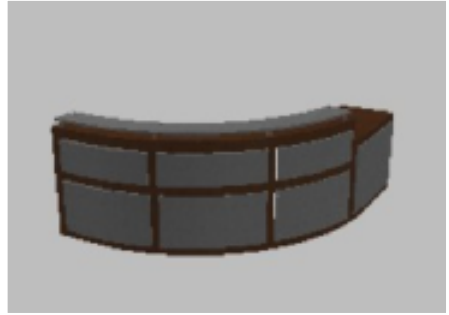


Hao Su

**Stanford** University

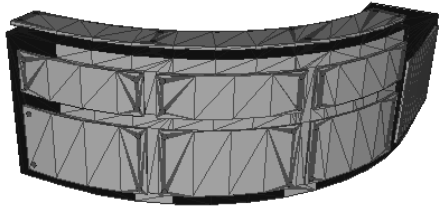May 10, 2017

# Agenda

- **Point cloud generation**

- Point cloud analysis
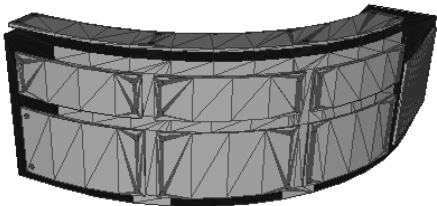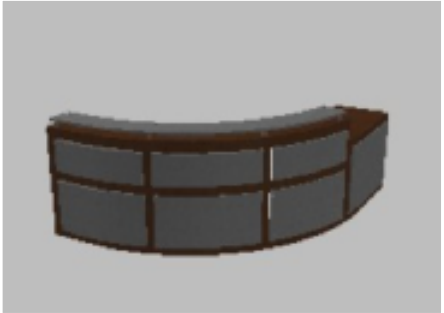
# Pipeline



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



Prediction

$$\left\{\begin{array}{c}(x'_1, y'_1, z'_1)\\ (x'_2, y'_2, z'_2)\\ ...\\ (x'_n, y'_n, z'_n)\end{array}\right\}$$

Shape predictor

$(f)$

render

sample

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

Groundtruth point **set**

# Pipeline



Prediction

$$\left\{ \begin{array}{c} (x'_1, y'_1, z'_1) \\ (x'_2, y'_2, z'_2) \\ ... \\ (x'_n, y'_n, z'_n) \end{array} \right\}$$

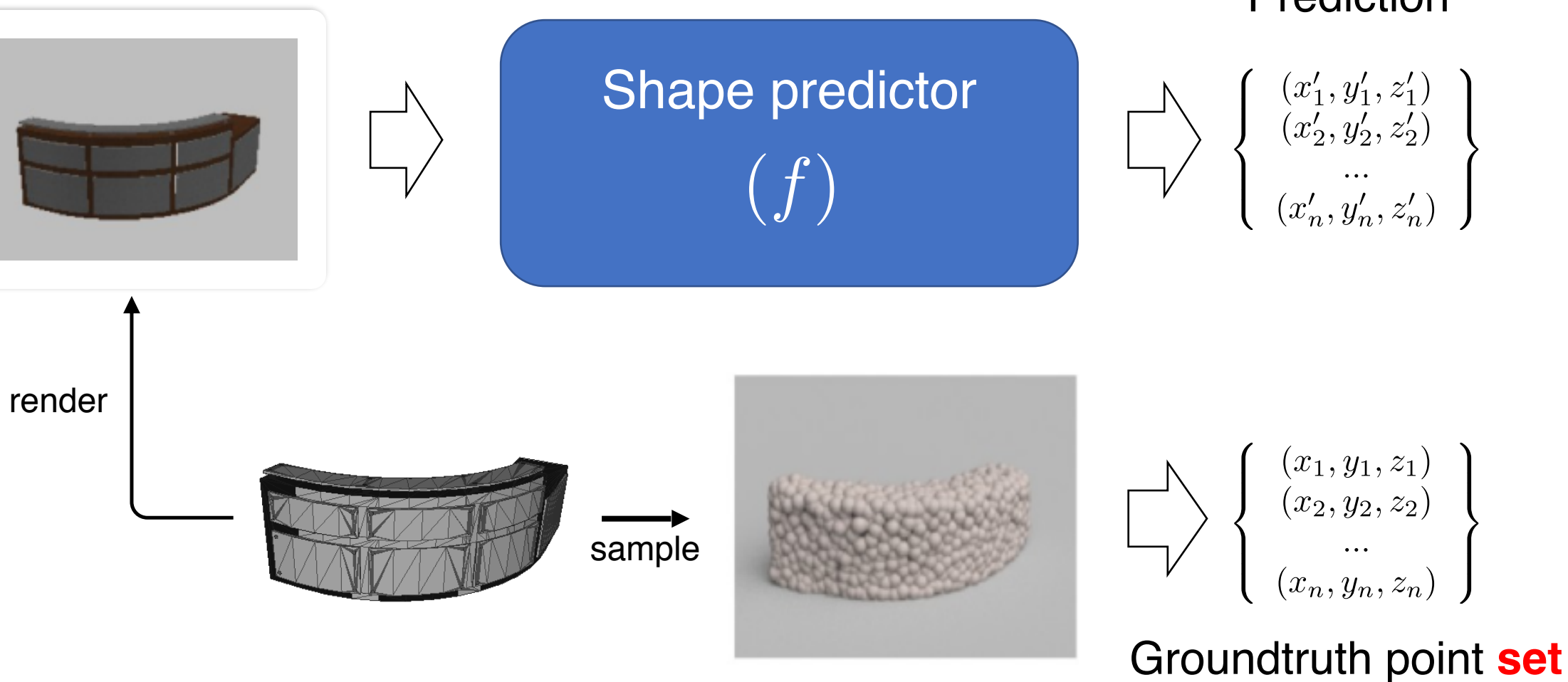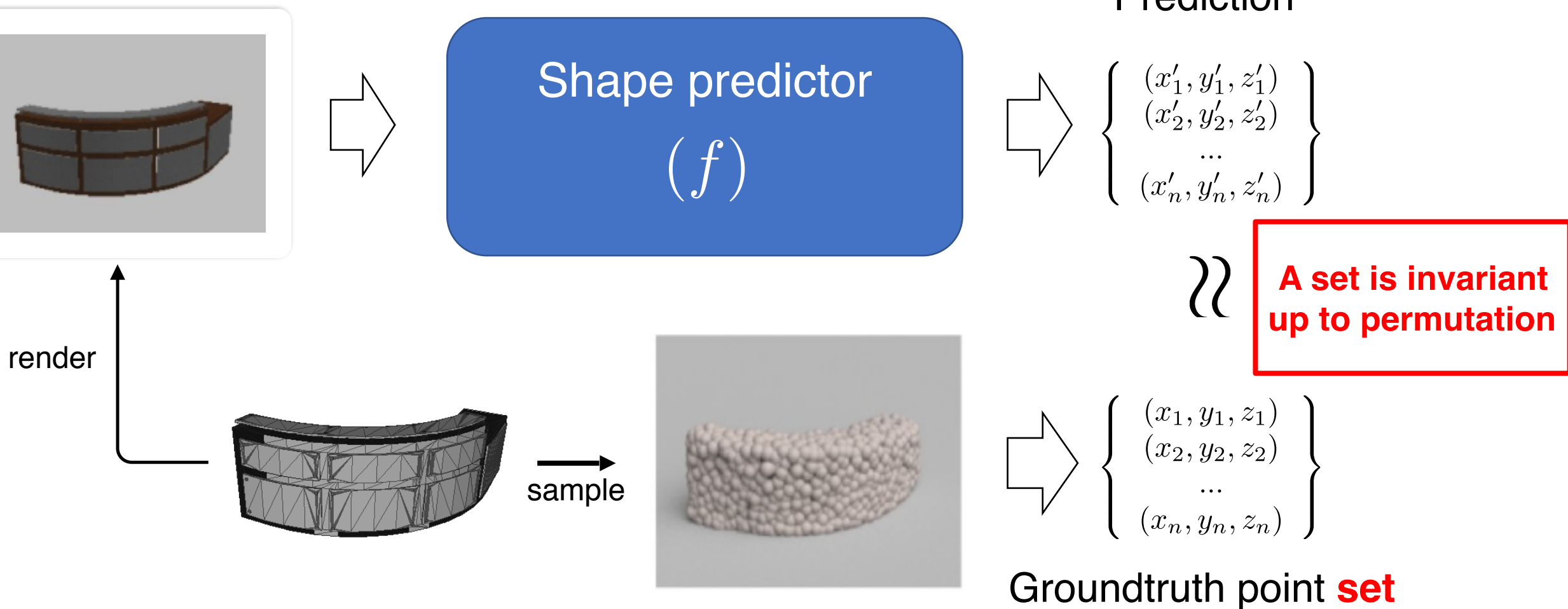**A set is invariant up to permutation**

render

sample

$$\left\{ \begin{array}{c} (x_1, y_1, z_1) \\ (x_2, y_2, z_2) \\ ... \\ (x_n, y_n, z_n) \end{array} \right\}$$
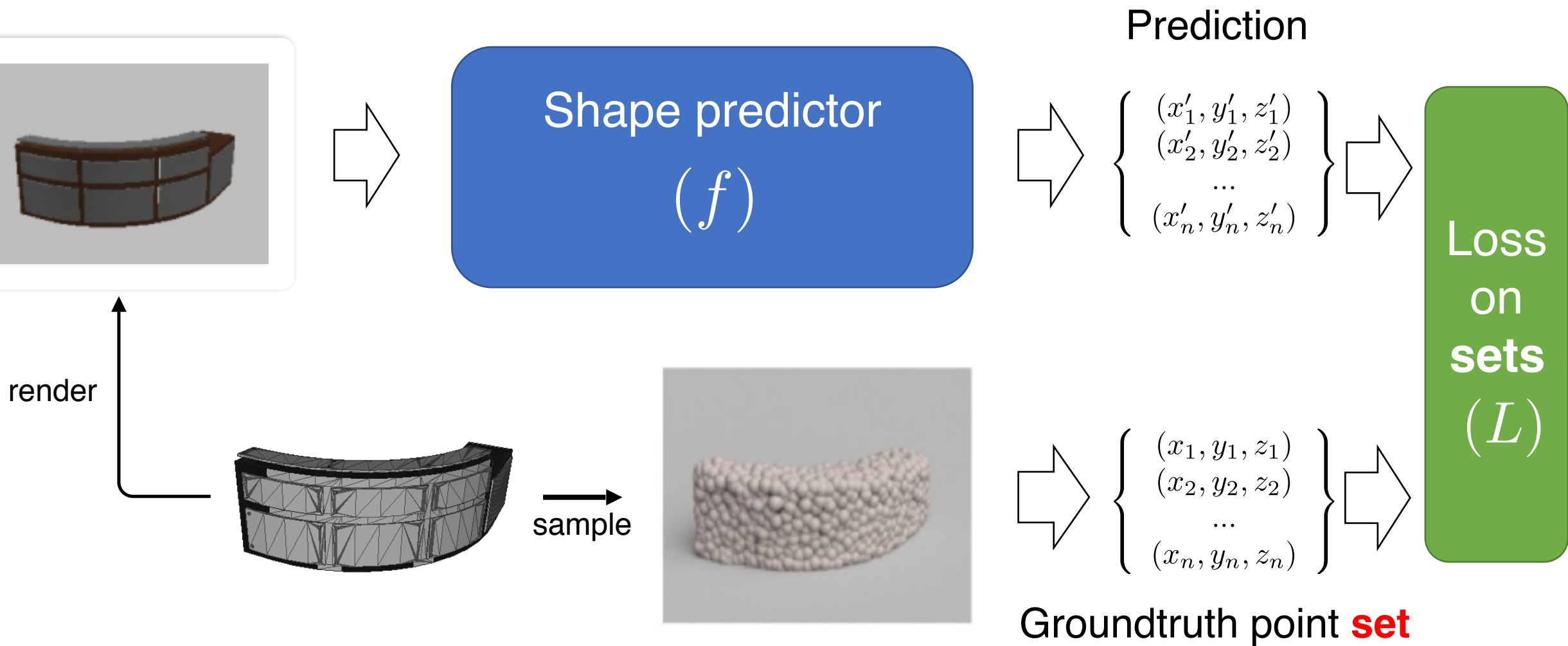
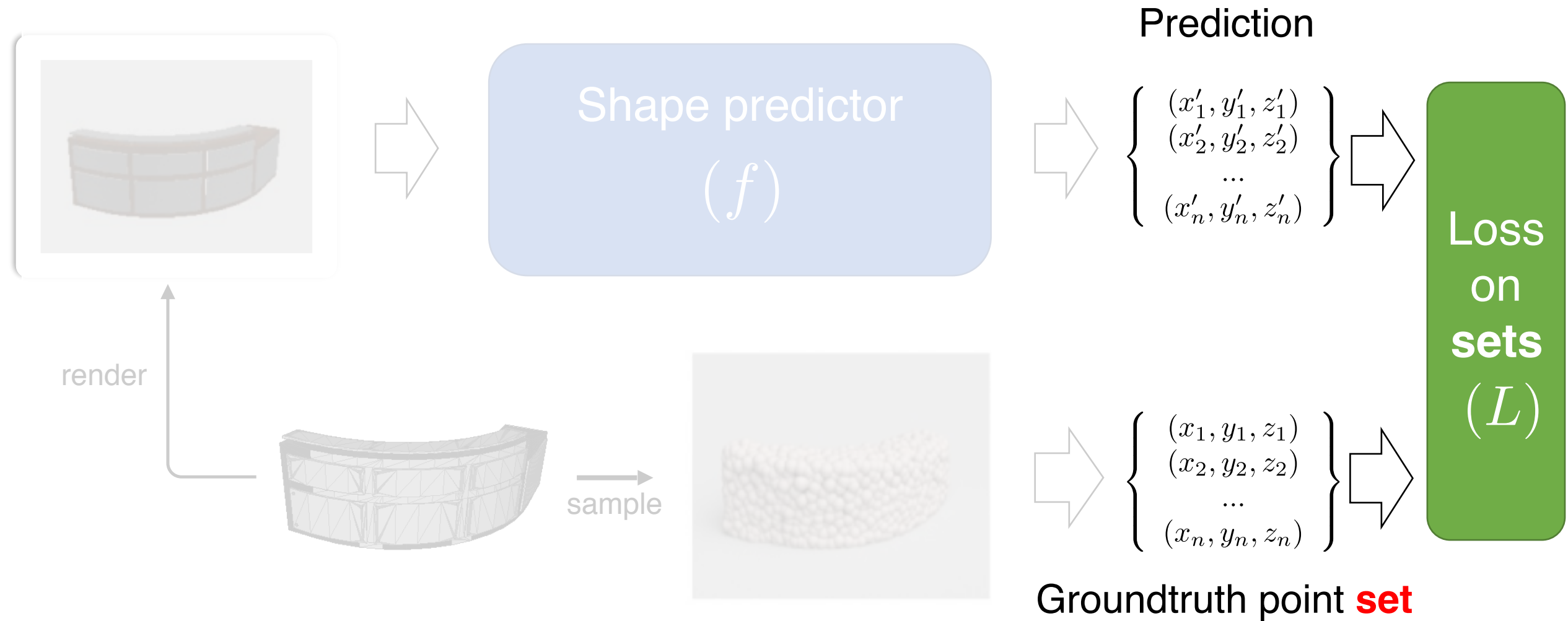Groundtruth point **set**

# Pipeline



Prediction

$$\left\{ \begin{array}{c} (x'_1, y'_1, z'_1) \\ (x'_2, y'_2, z'_2) \\ ... \\ (x'_n, y'_n, z'_n) \end{array} \right\}$$

Shape predictor $(f)$

render

sample

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

Groundtruth point **set**

Loss on **sets** $(L)$

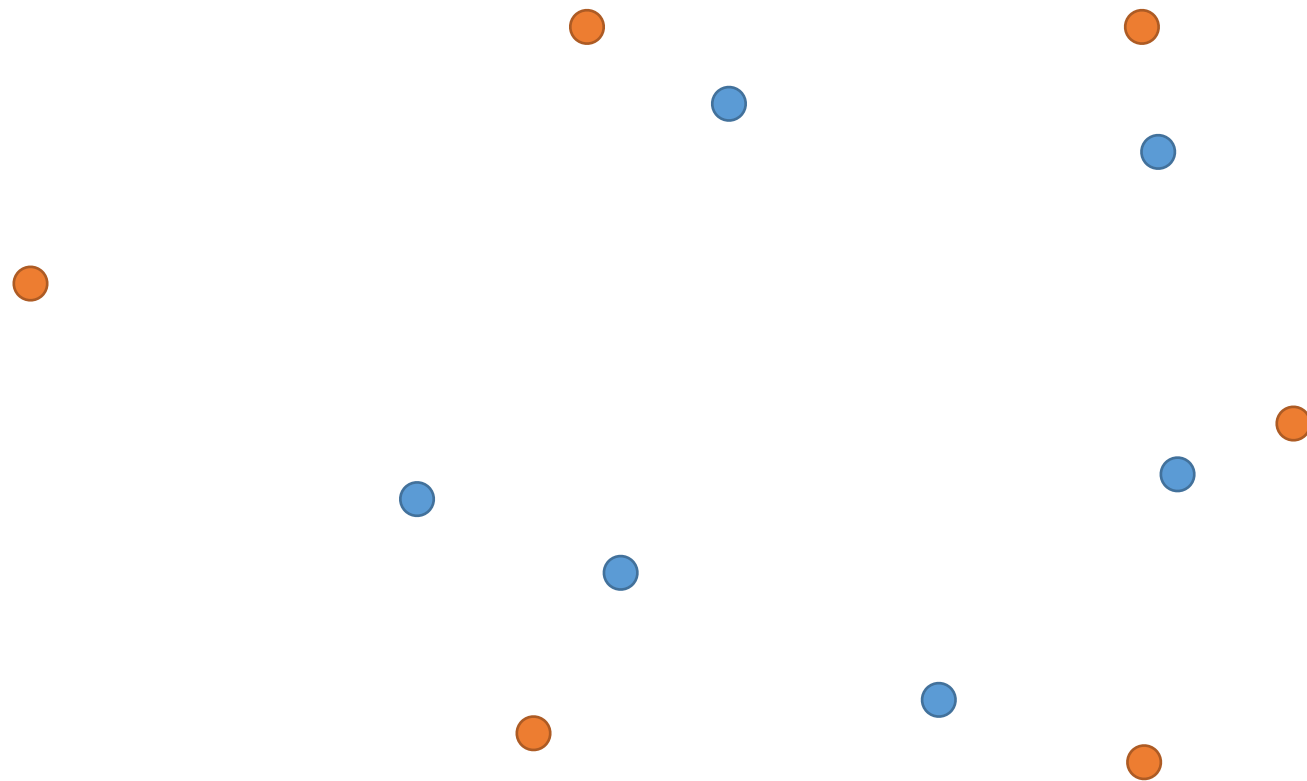**CVPR '17,** Point Set Generation

Prediction

Shape predictor

$(f)$

$$\left\{ \begin{array}{c} (x'_1, y'_1, z'_1) \\ (x'_2, y'_2, z'_2) \\ ... \\ (x'_n, y'_n, z'_n) \end{array} \right\}$$

render

sample

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

Groundtruth point **set**

Loss on **sets** $(L)$

**CVPR '17,** Point Set Generation

# Set comparison

Given two sets of points, measure their discrepancy

# Set comparison

Given two sets of points, measure their discrepancy

**Key challenge:**

**correspondence problem**

# Correspondence (I): optimal assignment
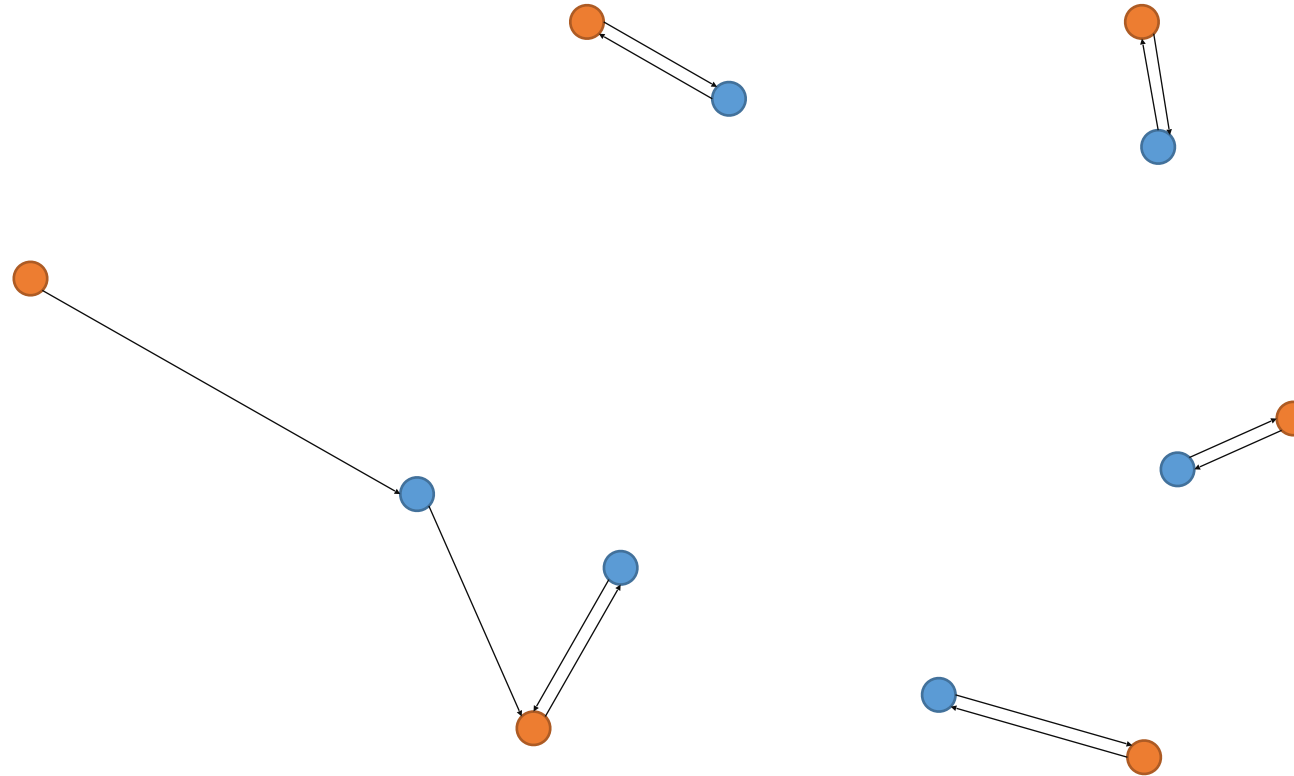
Given two sets of points, measure their discrepancy



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

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

**CVPR '17,** Point Set Generation

# 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

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

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

The mean shape carries characteristics of the distance metric

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

continuous
hidden variable
(radius)



Input                    EMD mean               Chamfer mean

# Mean shapes from distance metrics

The mean shape carries characteristics of the distance metric

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

continuous
hidden variable
(radius)

discrete
hidden variable
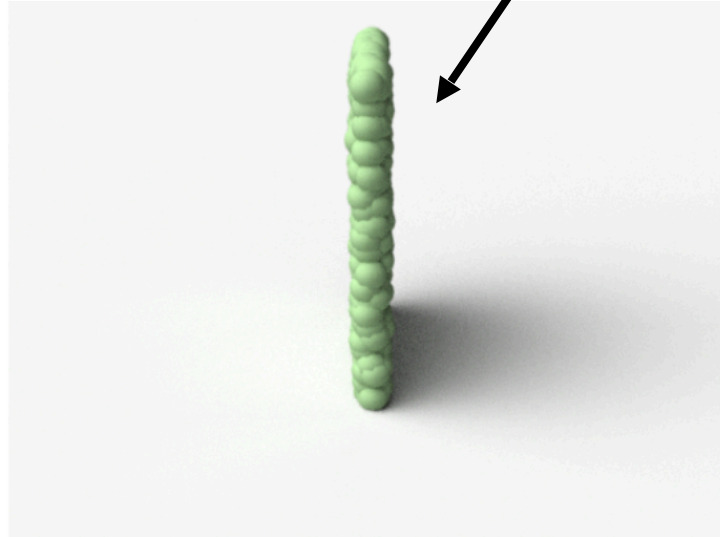(add-on location)



Input

EMD mean

Chamfer mean

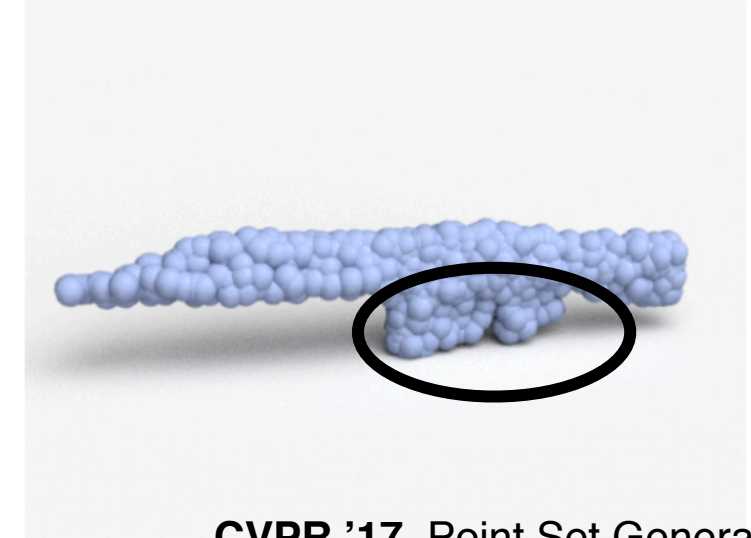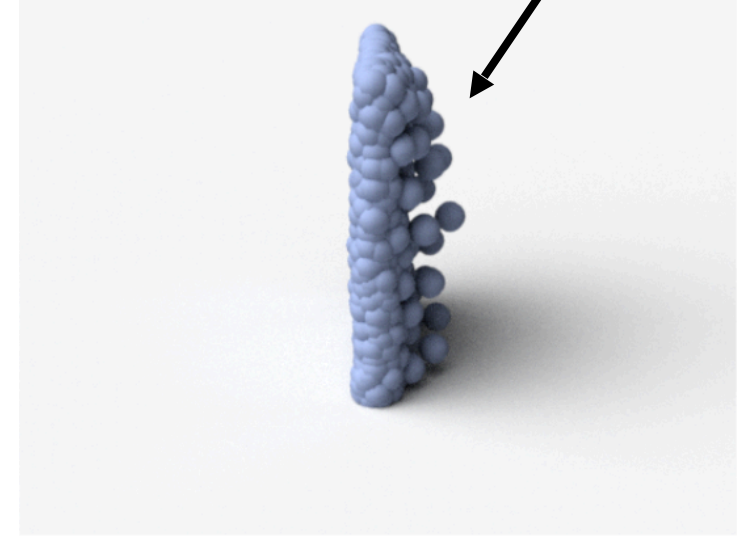# Comparison of predictions by EMD versus CD



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

- **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

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

  ✓

  - 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

- **Efficient** to compute

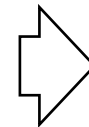  Chamfer distance: trivially parallelizable on CUDA

  Earth Mover's distance (optimal assignment):

  - We implement a **distributed** approximation algorithm on CUDA

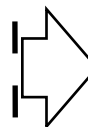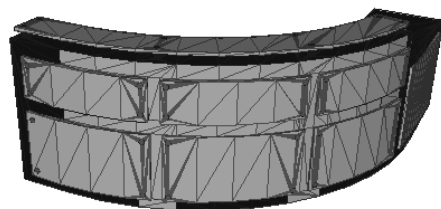  - Based upon [Bertsekas, 1985], $(1 + \epsilon)$-approximation
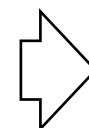
Deep network

$(f)$

Prediction

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

# Pipeline



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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

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

Predictor

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

conv

Predictor

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

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

- 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**

Capture common structures

conv

Deconv branch

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

Capture intricate structures

FC branch

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

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

CVPR '17, Point Set Generation

# 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



Deconv network for image segmentation

*Credit: FCNN, Long et al.*

# 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



**How to predict curved surfaces in 3D?**

Deconv network for image segmentation

*Credit: FCNN, Long et al.*

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



$f$

$df(v)$

$df\text{-}(u)$

$v$

$u$

*Credit: Discrete Differential Geometry, Crane et al.*

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



$f$

$df(v)$

$df\text{-}(u)$

$v$

$u$

$(x, y, z)$

*Credit: Discrete Differential Geometry, Crane et al.*

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



$f$

x-map

y-map

z-map

coordinate maps

$(x, y, z)$

$v$

$u$

$df(v)$

$df\text{-}(u)$

*Credit: Discrete Differential Geometry, Crane et al.*

# Parametrization prediction by deconv network

Capture common structures

conv

Deconv branch

FC branch

Capture intricate structures

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

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

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

# Parametrization prediction by deconv network

deconv

Capture common structures

conv

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

coordinate maps

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

FC branch

Capture intricate structures

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

**Note that**

- The parametrization (2D/3D mapping) is learned from data
- i.e., obtains a network and data friendly parametrization

coordinate maps

sample

Loss
on
**sets**
$(L)$

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

Capture common structures

**CVPR '17,** Point Set Generation

Observation:

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

$(x_k, y_k, z_k)$

map of x coord     map of y coord     map of z coord

Observation:

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

$(x_k, y_k, z_k)$

map of x coord          map of y coord          map of z coord

- 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**

Capture common structures

deconv

conv

FC branch

Capture intricate structures

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

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

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

Capture common structures

deconv

conv

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

dense

dense

$\begin{matrix} x_1 \\ y_1 \\ z_1 \\ \vdots \\ x_m \\ y_m \\ z_m \end{matrix}$

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

Capture intricate structures

fc

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

Loss on **sets** $(L)$

sample

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

**CVPR '17,** Point Set Generation

Capture common structures

deconv

conv

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

dense

Capture intricate structures

dense

$$\begin{bmatrix} x_1 \\ y_1 \\ z_1 \\ \vdots \\ x_m \\ y_m \\ z_m \end{bmatrix}$$

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

fc

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

Loss
on
**sets**
$(L)$

- Points are predicted **independently**
- Dense connection introduces **more parameters** to accommodate **high variance**

**CVPR '17,** Point Set Generation

# Visualization of the effect of FC branch



Observation:

- The arrangement of predicted points are uncorrelated

x-coord     y-coord     z-coord     **red**

Observation:

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

x-coord          y-coord          z-coord          **red**

**CVPR '17,** Point Set Generation

**CVPR '17,** Point Set Generation

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

**red**: FC branch – **intricate** structures

# Comparison to state-of-the-art



Input    Ours    Ours (post-processed)    Groundtruth    state-of-the-art (3D-R2N2)

- Better global structure
- Better details

**CVPR '17,** Point Set Generation

# Comparison to state-of-the-art

Trained/tested on 2K object categories



**CVPR '17,** Point Set Generation

RGBD map (input)        $90°$ view of input        output: completed point cloud

**CVPR '17,** Point Set Generation

# 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**

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

Each point is colored according to the assigned primitive

*Generalized Cylinders,* Binford (1971)

**Encoder**

**Predictor
for primitive parameters**

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

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

Variable number of parts? We predict "primitive existence probability"

**Loss**

Basic idea: **Chamfer distance!**

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

Sample points on the groundtruth mesh and predicted 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



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

## Differentiable!

Speed up the computation leveraging parameterization of primitives

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

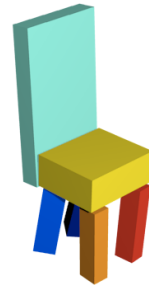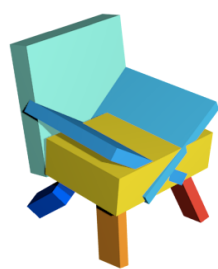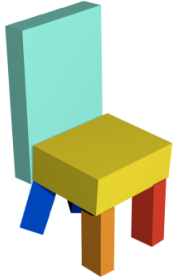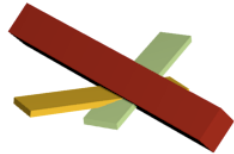| Method | [31] (initial) | [31] (refined) | Ours |
|--------|---------------|----------------|------|
| Accuracy | 78.6 | 84.8 | 89.0 |

Mean accuracy (face area) on Shape COSEG chairs.

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

# Agenda

- Point cloud generation

- **Point cloud analysis**

# Applications of Point Set Learning

- **Robot Perception**

What and where are the objects in a LiDAR scanned scene?



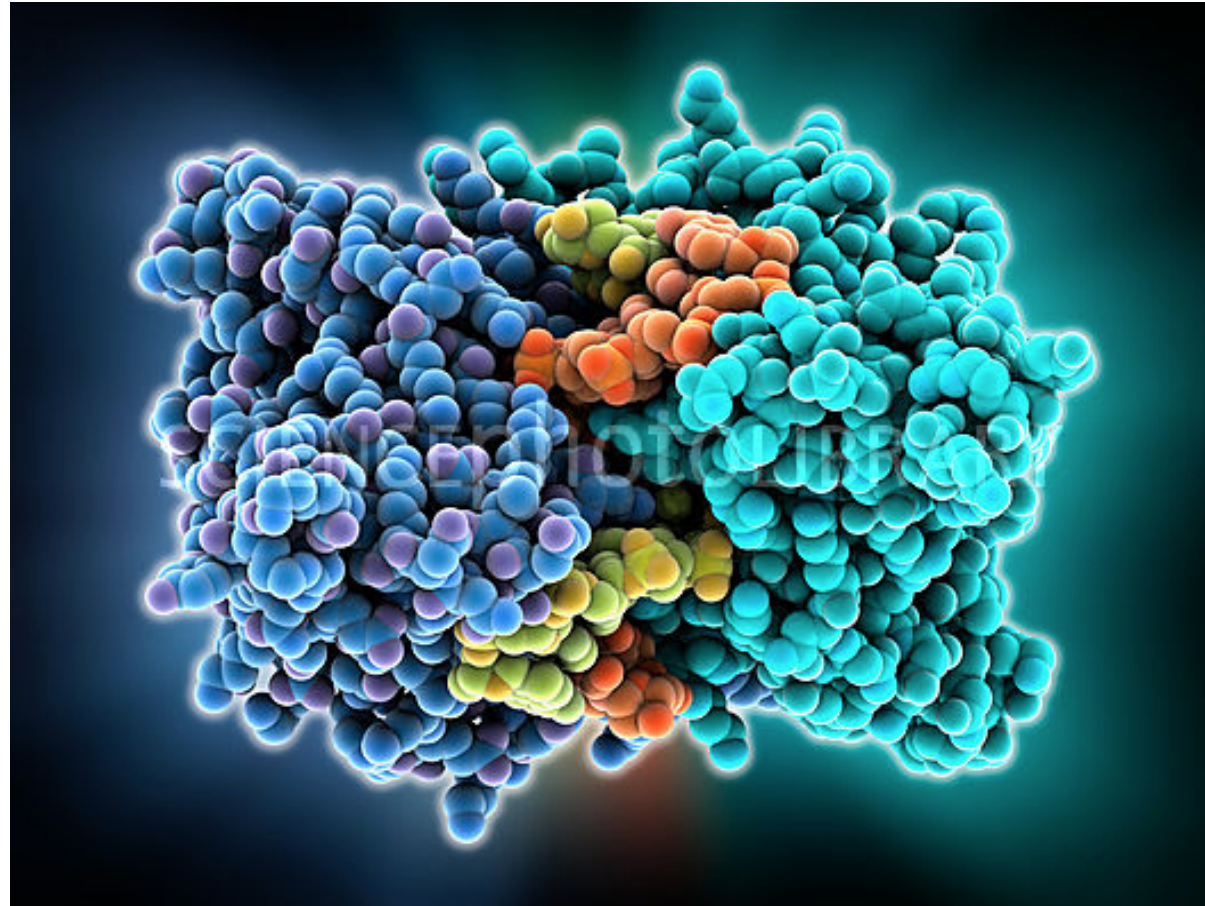*https://3dprint.com/116569/self-driving-cars-privacy/*

# Applications of Point Set Learning
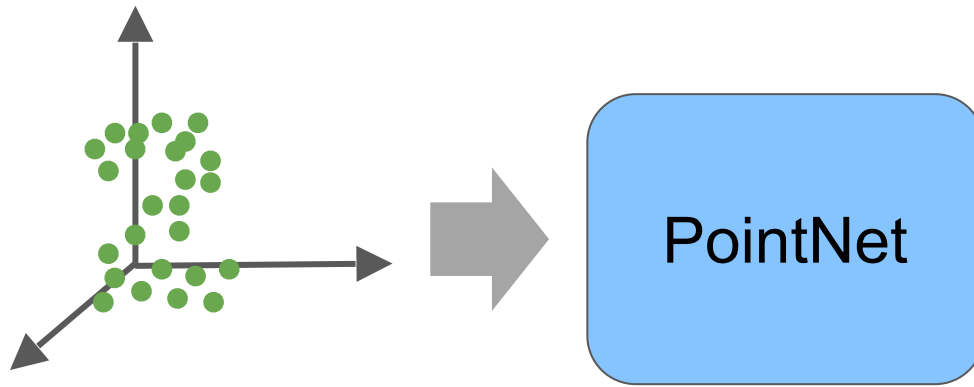
- **Molecular Biology**

Can we infer an enzyme's category (reactions they catalyze) from its structure?



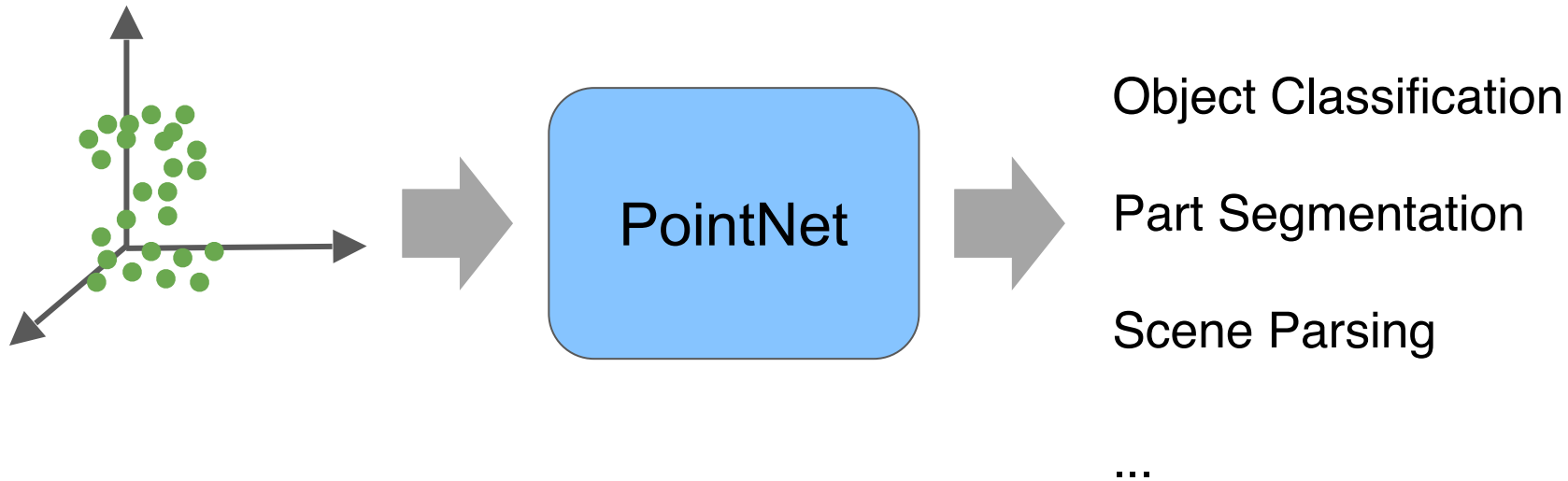*EcoRV restriction enzyme molecule, LAGUNA DESIGN/SCIENCE PHOTO LIBRARY*

End-to-end learning for **unstructured, unordered** point data
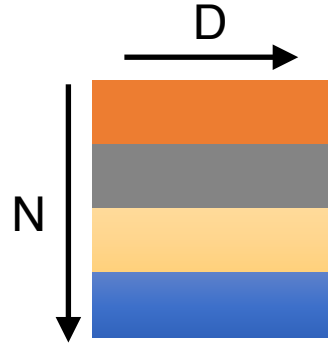
# Directly process point cloud data

End-to-end learning for **unstructured, unordered** point data

**Unified** framework for various tasks



PointNet

Object Classification

Part Segmentation

Scene Parsing

...

Point cloud: N **orderless** points, each represented by a D dim coordinate



2D array representation

Point cloud: N **orderless** points, each represented by a D dim coordinate



2D array representation

**Permutation invariance**

**Transformation invariance**

Point cloud: N **orderless** points, each represented by a D dim coordinate



represents the same **set** as

2D array representation

**Permutation invariance**

$$f(x_1, x_2, \ldots, x_n) \equiv f(x_{\pi_1}, x_{\pi_2}, \ldots, x_{\pi_n}), \quad x_i \in \mathbb{R}^D$$

**Examples:**

$$f(x_1, x_2, \ldots, x_n) = \max\{x_1, x_2, \ldots, x_n\}$$

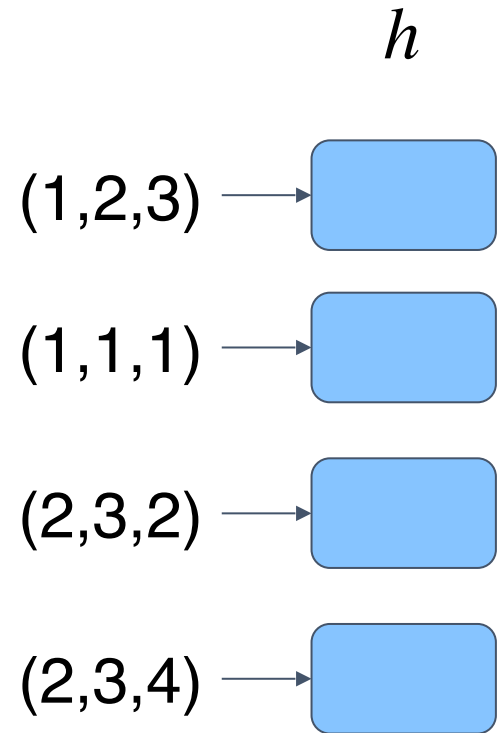$$f(x_1, x_2, \ldots, x_n) = x_1 + x_2 + \ldots + x_n$$
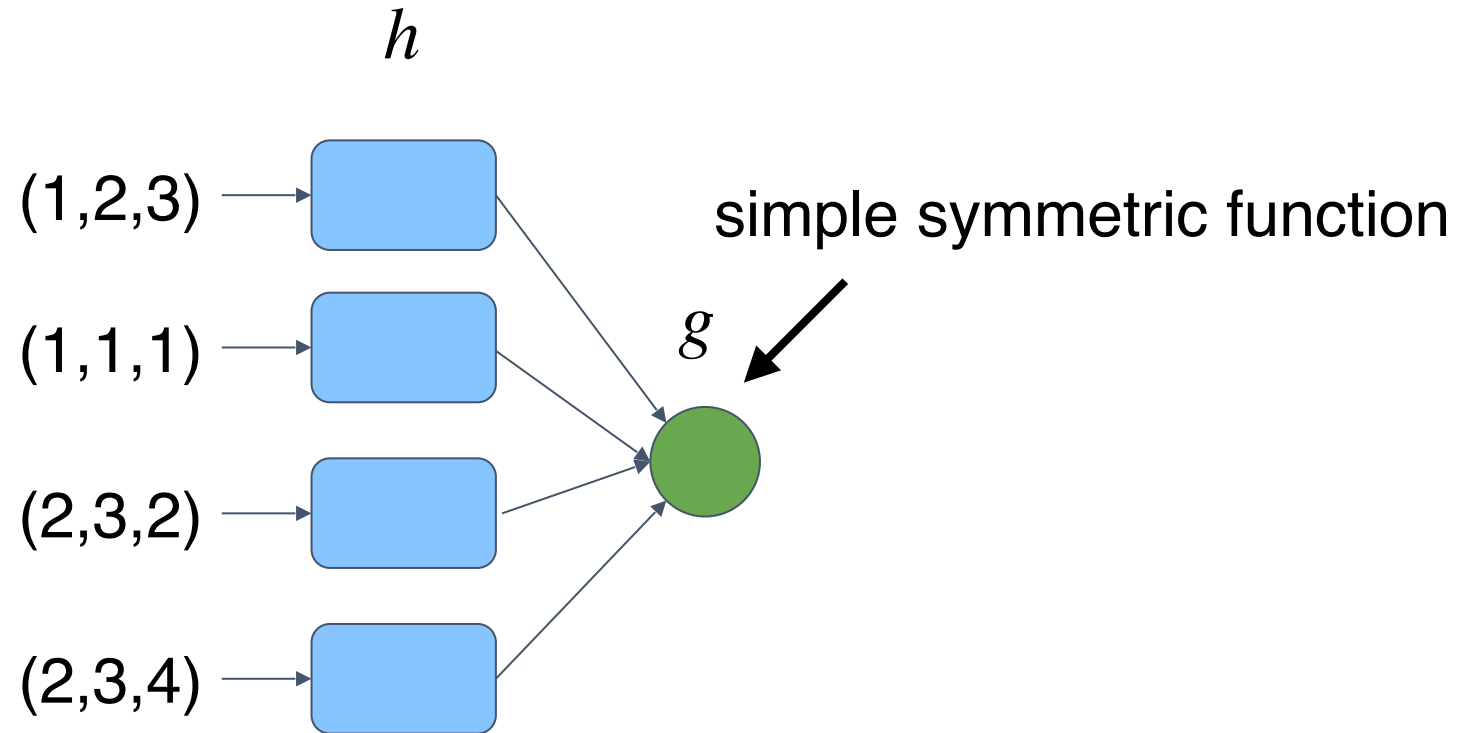
$$\ldots$$

# Construct symmetric function family

**Observe:** $f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$ is symmetric if $g$ is symmetric

# Construct symmetric function family

**Observe:** $f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$ is symmetric if $g$ is symmetric

$$h$$

(1,2,3) →
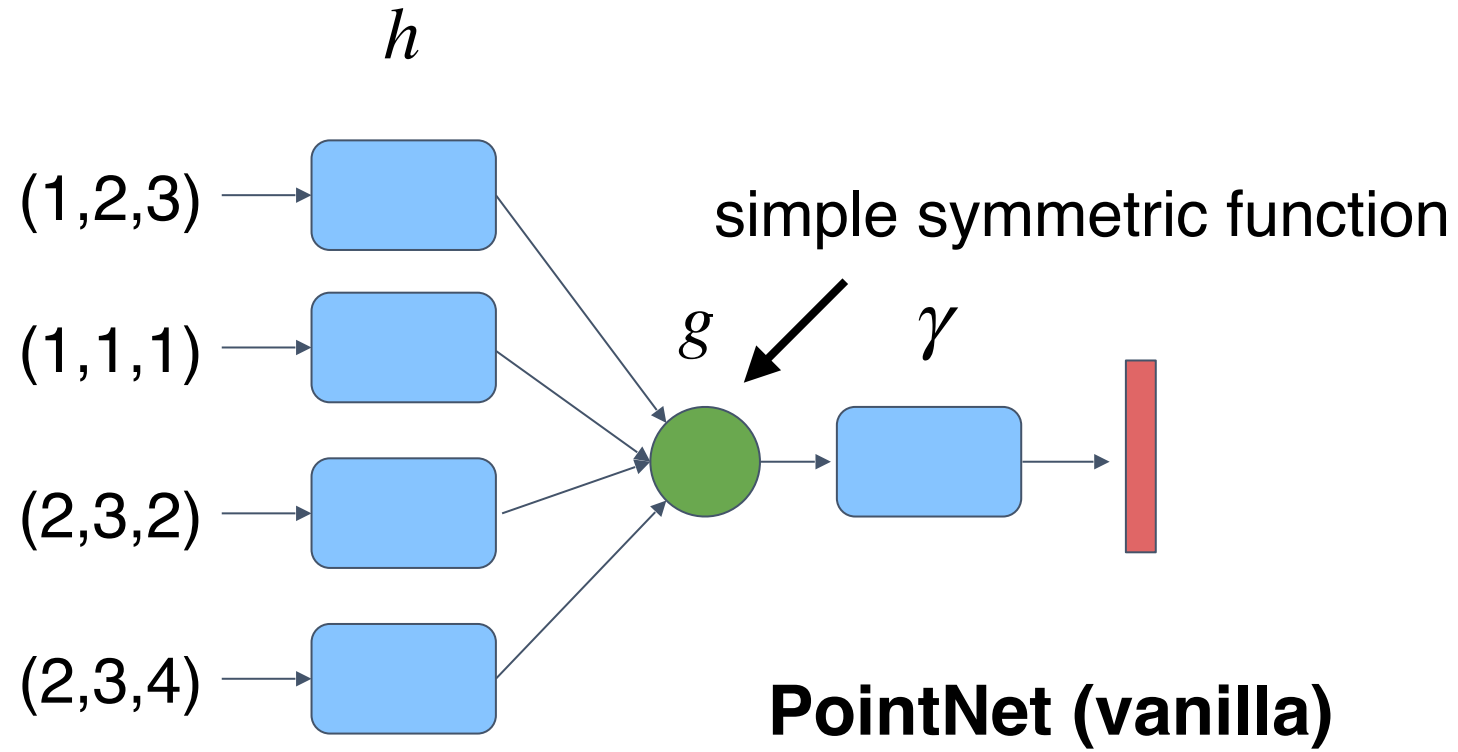(1,1,1) →
(2,3,2) →
(2,3,4) →

# Construct symmetric function family

**Observe:** $f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$ is symmetric if $g$ is symmetric
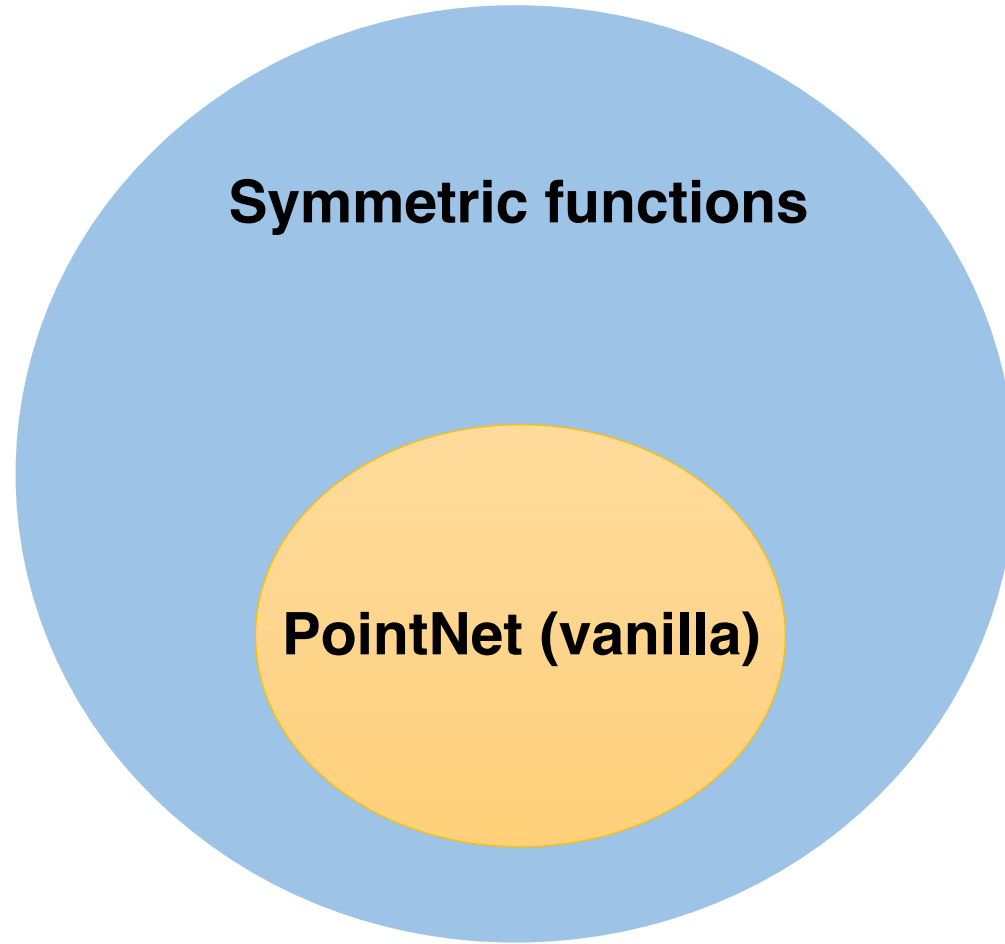
# Construct symmetric function family

**Observe:** $f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$ is symmetric if $g$ is symmetric



$h$

(1,2,3)

(1,1,1)

(2,3,2)

(2,3,4)

simple symmetric function

$g$ $\gamma$

**PointNet (vanilla)**

**Symmetric functions**

**PointNet (vanilla)**

## Theorem:

A Hausdorff continuous symmetric function $f : 2^{\mathcal{X}} \to \mathbb{R}$ can be arbitrarily approximated by PointNet.
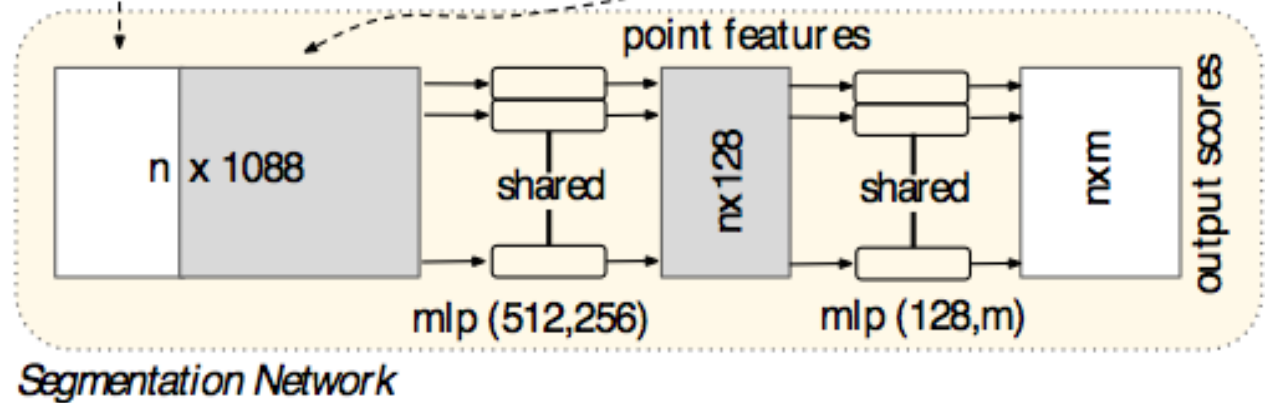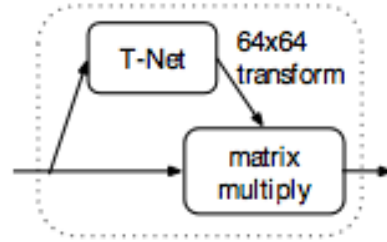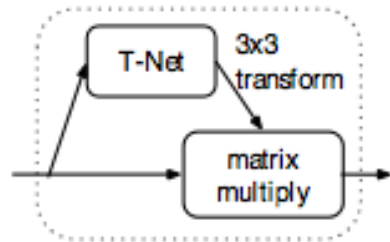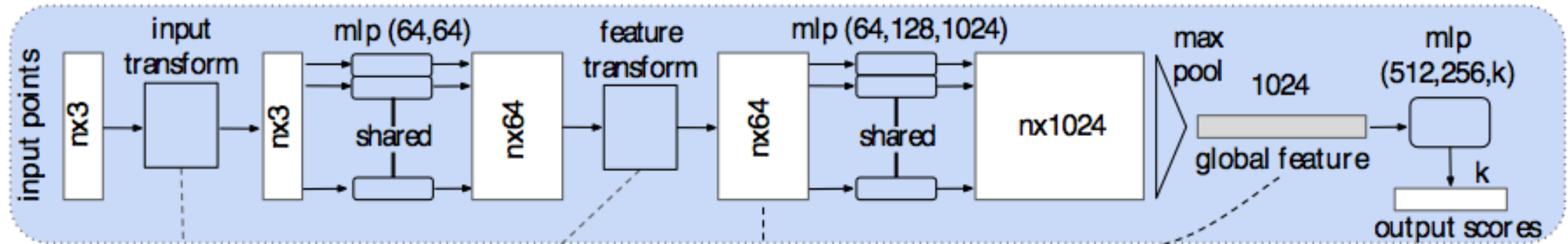
$$\left| f(S) - \gamma \left( \underset{x_i \in S}{\mathrm{MAX}} \{h(x_i)\} \right) \right| < \epsilon$$

$$S \subseteq \mathbb{R}^d,$$
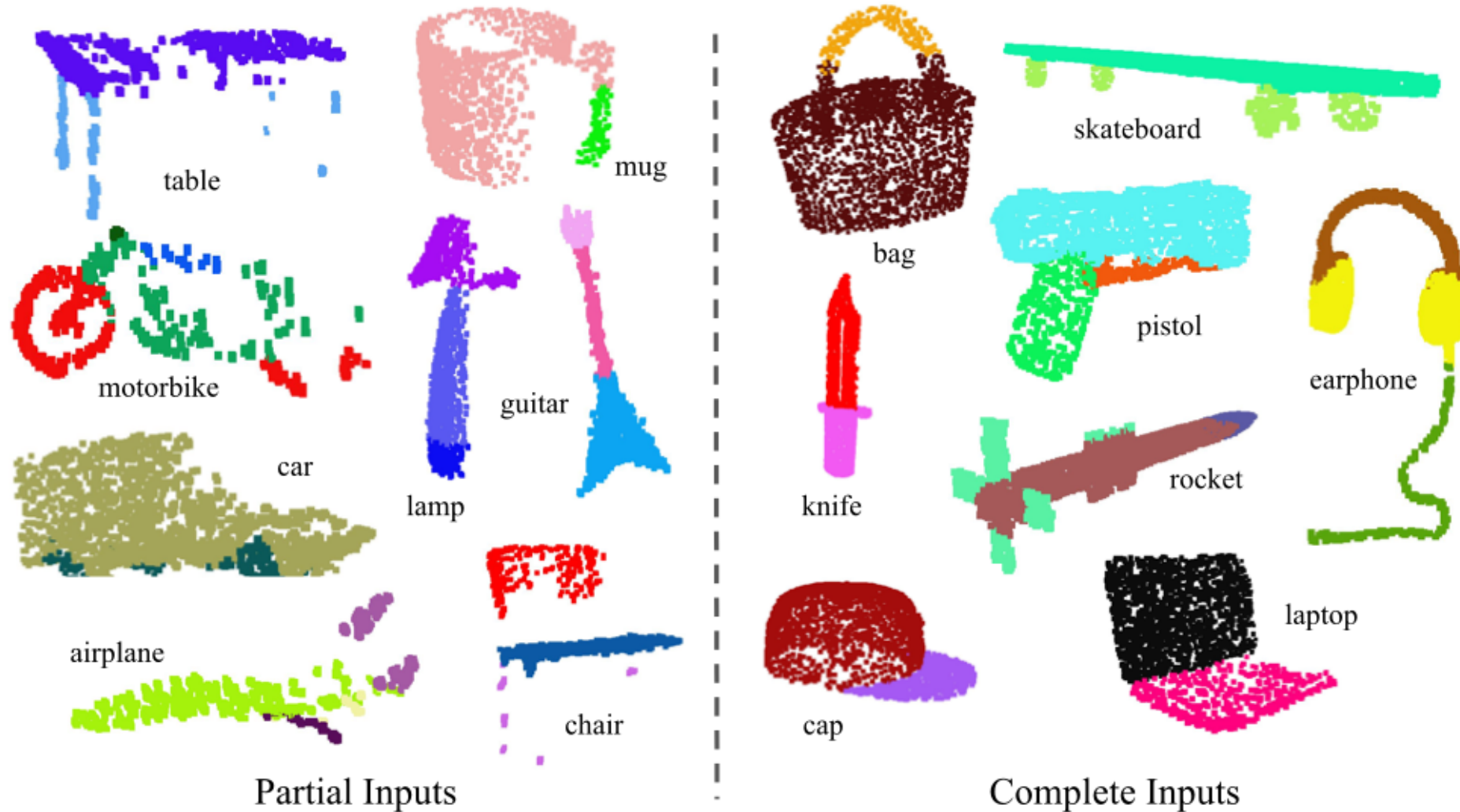
**PointNet (vanilla)**

# PointNet Architecture

# Results on Object Classification

Object Classification Accuracy on ModelNet40

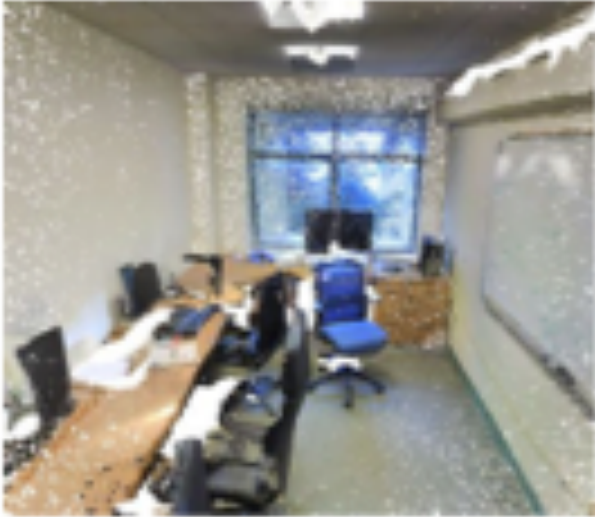|  | input | #views | accuracy avg. class | accuracy overall |
|---|---|---|---|---|
| SPH [12] | mesh | - | 68.2 | - |
| 3DShapeNets [29] | volume | 1 | 77.3 | 84.7 |
| VoxNet [18] | volume | 12 | 83.0 | 85.9 |
| Subvolume [19] | volume | 20 | 86.0 | **89.2** |
| LFD [29] | image | 10 | 75.5 | - |
| MVCNN [24] | image | 80 | **90.1** | - |
| Ours baseline | point | - | 72.6 | 77.4 |
| Ours PointNet | point | 1 | 86.2 | **89.2** |

Partial Inputs

Complete Inputs

# Results on Object Part Segmentation

Part Segmentation mIoU on ShapeNet Part Dataset

| | mean | aero | bag | cap | car | chair | ear phone | guitar | knife | lamp | laptop | motor | mug | pistol | rocket | skate board | table |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # shapes | | 2690 | 76 | 55 | 898 | 3758 | 69 | 787 | 392 | 1547 | 451 | 202 | 184 | 283 | 66 | 152 | 5271 |
| Wu [28] | - | 63.2 | - | - | - | 73.5 | - | - | - | 74.4 | - | - | - | - | - | - | 74.8 |
| Yi [30] | 81.4 | 81.0 | 78.4 | 77.7 | **75.7** | 87.6 | 61.9 | **92.0** | 85.4 | **82.5** | **95.7** | **70.6** | 91.9 | **85.9** | 53.1 | 69.8 | 75.3 |
| 3DCNN | 79.4 | 75.1 | 72.8 | 73.3 | 70.0 | 87.2 | 63.5 | 88.4 | 79.6 | 74.4 | 93.9 | 58.7 | 91.8 | 76.4 | 51.2 | 65.3 | 77.1 |
| Ours | **83.7** | **83.4** | **78.7** | **82.5** | 74.9 | **89.6** | **73.0** | 91.5 | **85.9** | 80.8 | 95.3 | 65.2 | **93.0** | 81.2 | **57.9** | **72.8** | **80.6** |

# Results on Semantic Scene Segmentation

# Results on Semantic Scene Parsing

Semantic Segmentation (point based)
on Stanford Semantic Parsing dataset

|               | mean IoU | overall accuracy |
|---------------|----------|------------------|
| Ours baseline | 20.12    | 53.19            |
| Ours PointNet | **47.71** | **78.62**       |

3D Object Detection (bounding box based)

|                   | table     | chair     | sofa     | board     | mean      |
|-------------------|-----------|-----------|----------|-----------|-----------|
| # instance        | 455       | 1363      | 55       | 137       |           |
| Armeni et al. [2] | 46.02     | 16.15     | **6.78** | 3.91      | 18.22     |
| Ours              | **46.67** | **33.80** | 4.76     | **11.72** | **24.24** |

# Robustness to Data Corruption

# Robustness to Data Corruption

# Visualizing Point Functions

Compact View:

1x3

1x1024

FCs

Expanded View:

1x3

1x1024

FC · 64 · FC · 64 · FC · 64 · FC · 128 · FC
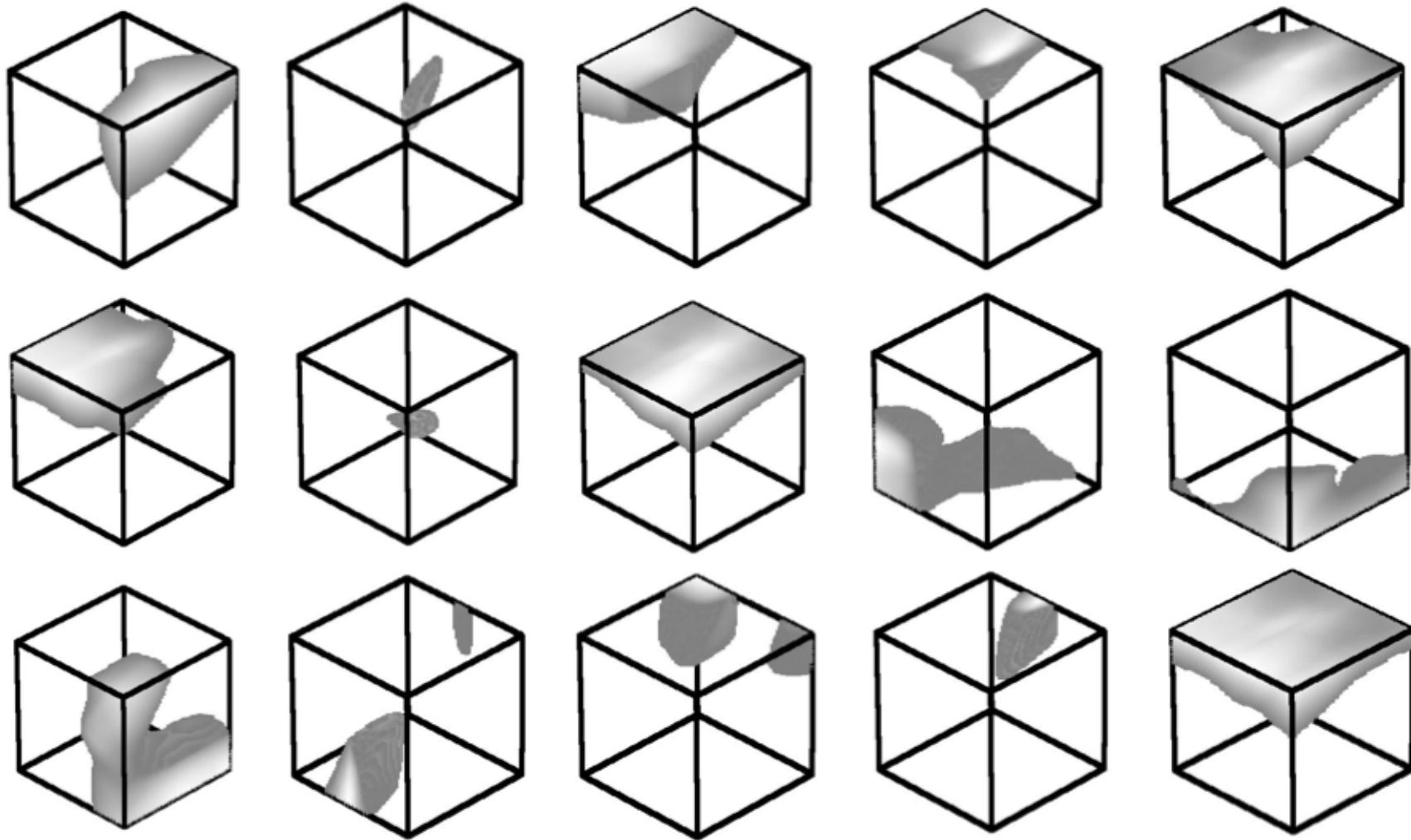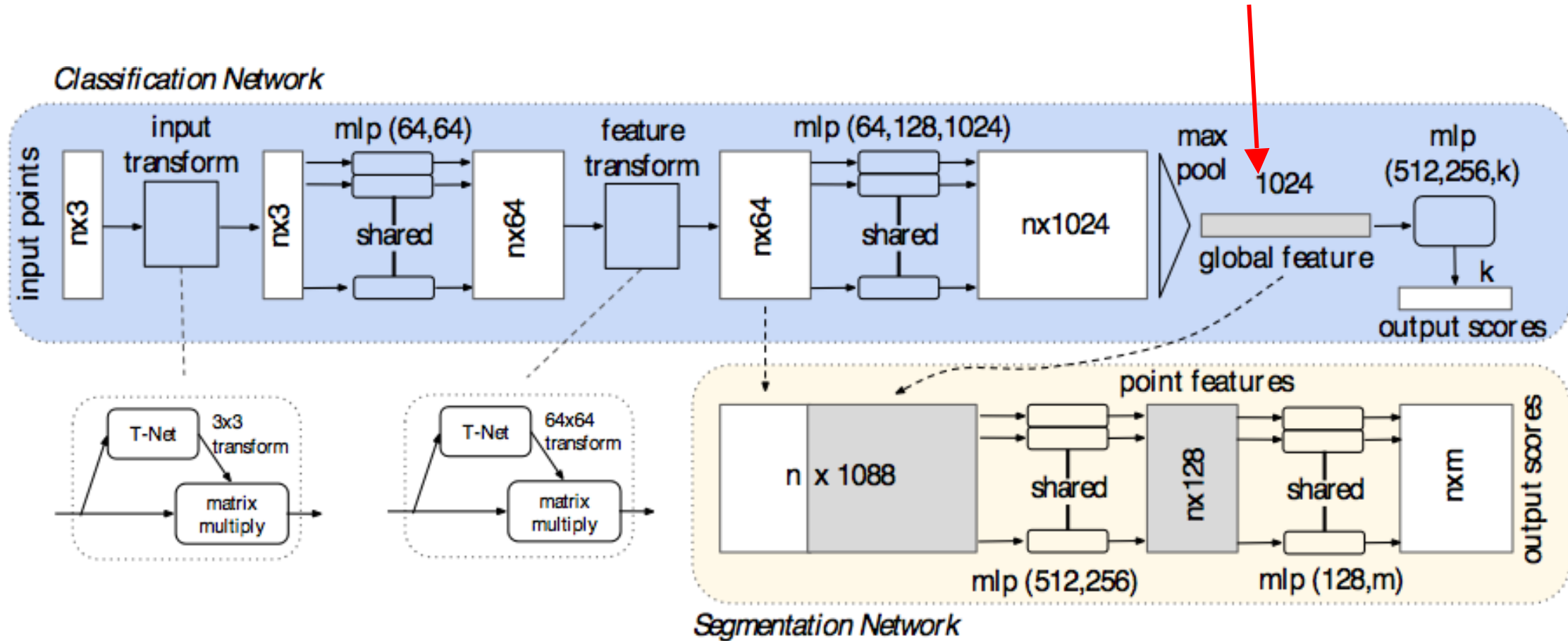
**Which input point will activate neuron j?**

Find the top-K points in a dense volumetric grid that activates neuron j.
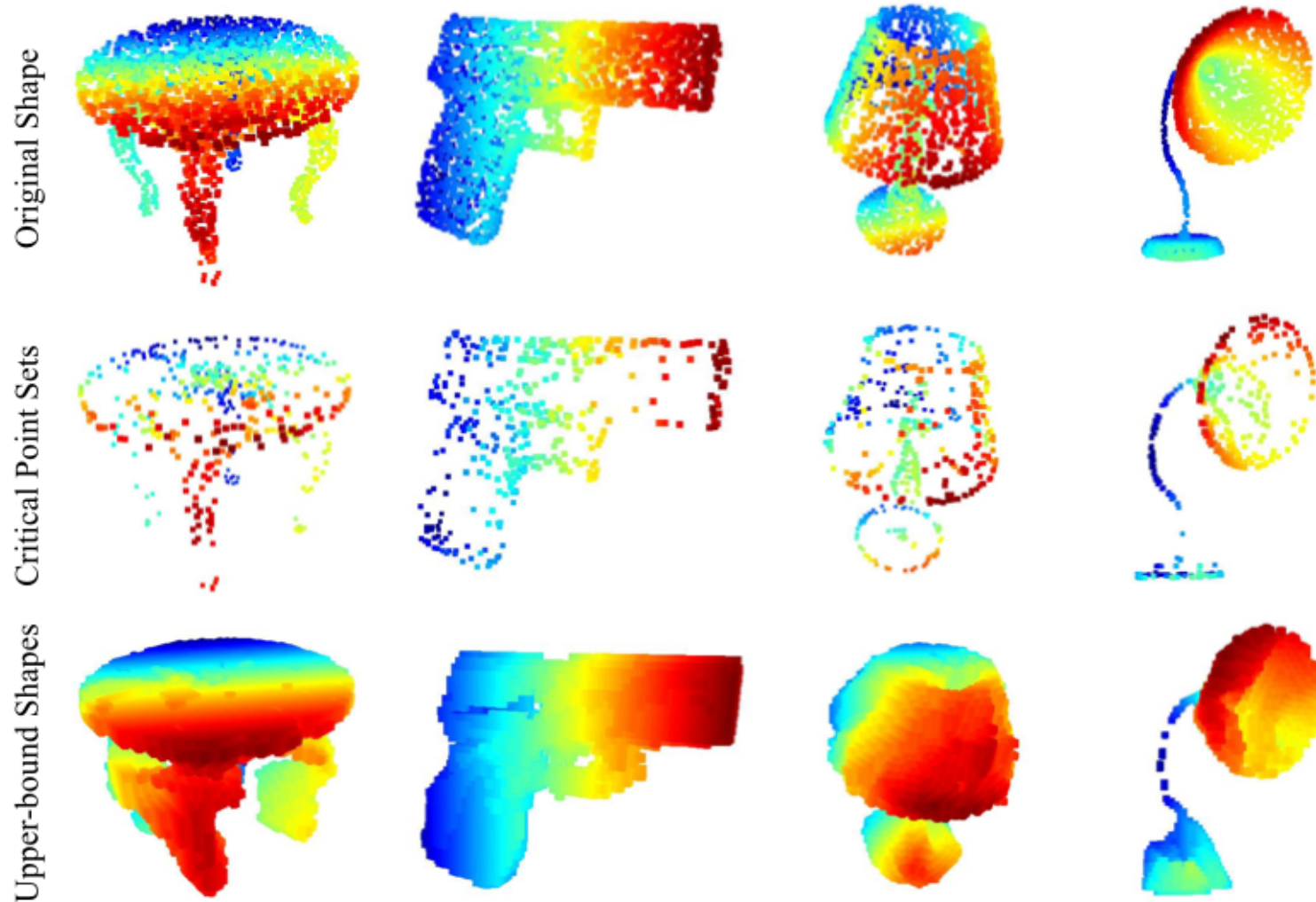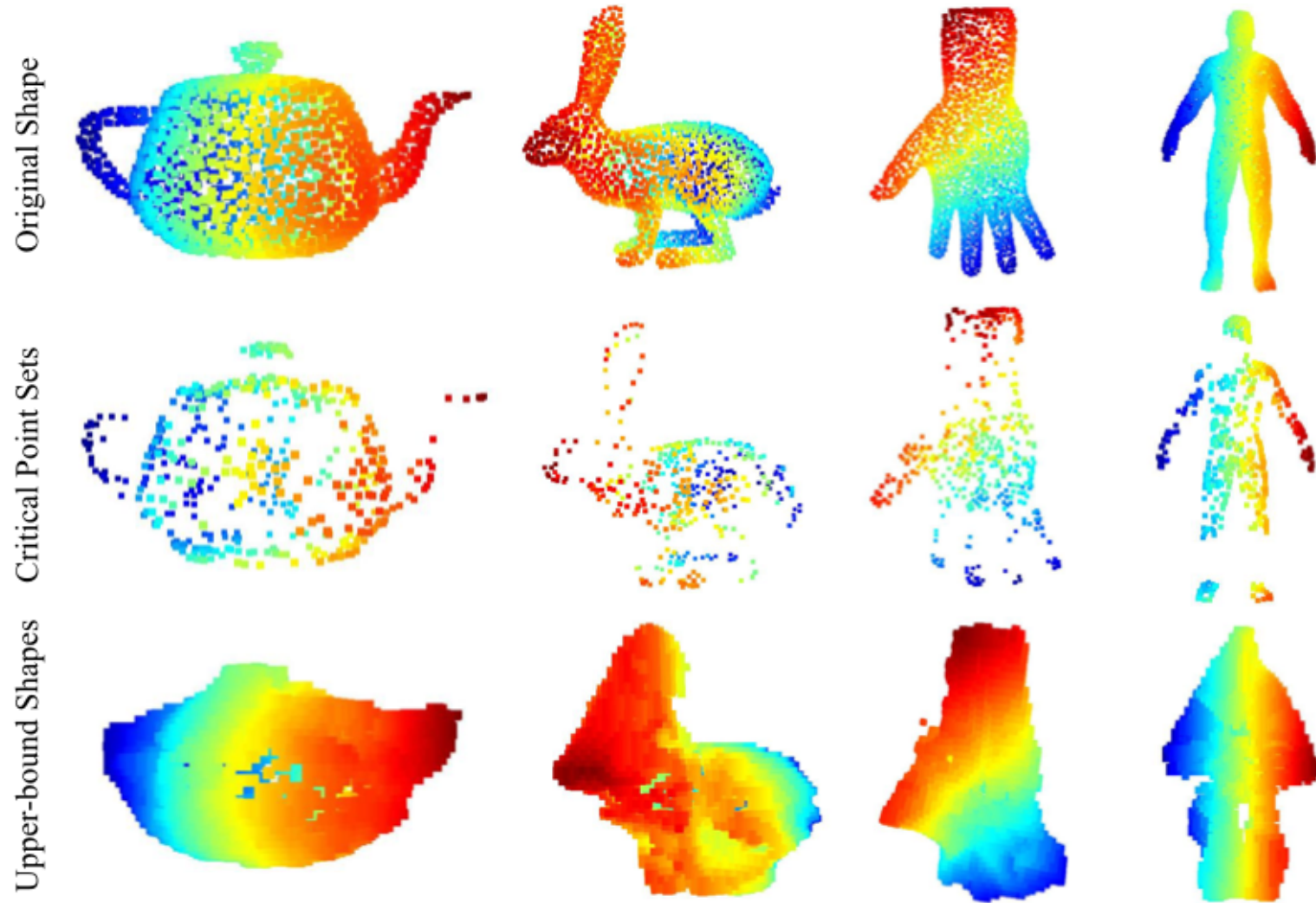
# Visualizing Global Point Cloud Features



*What's captured and left out here?*
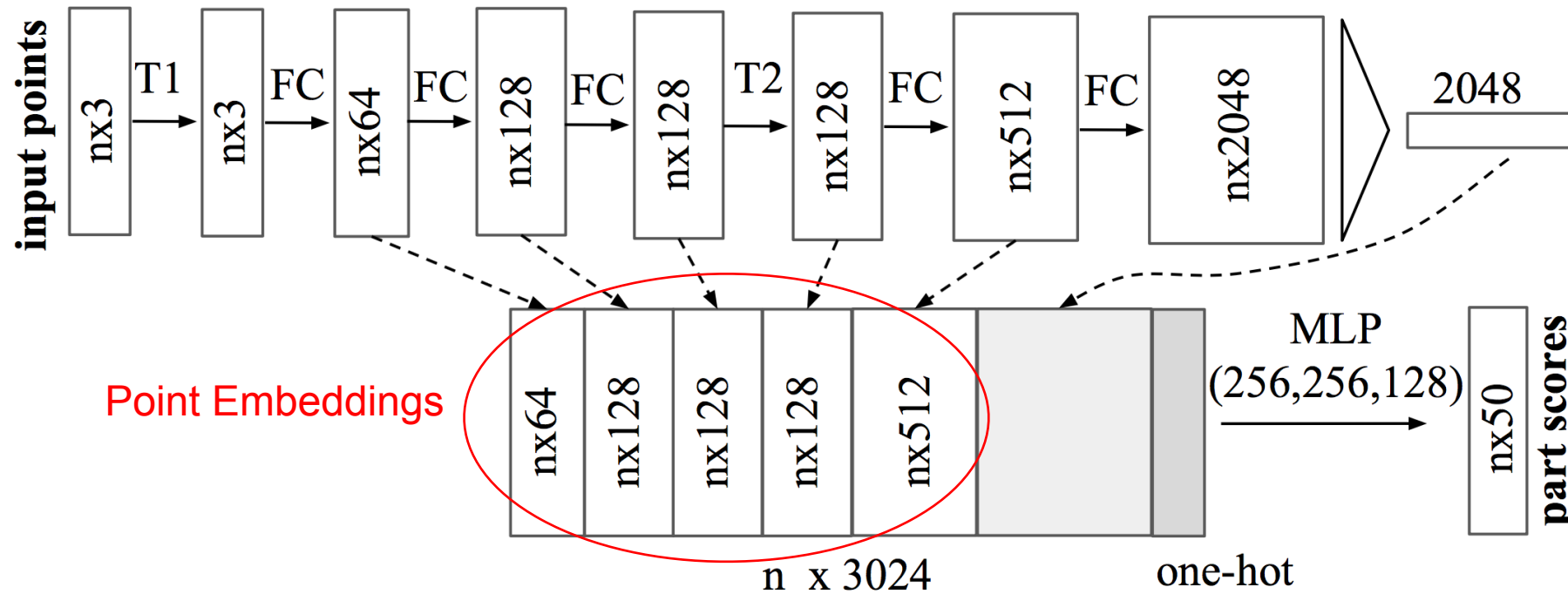
# Visualizing Global Point Cloud Features

*Segmentation Network*

*Segmentation Network*

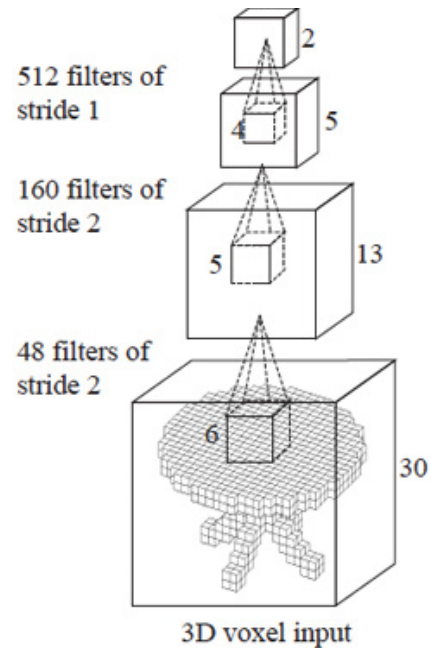*Segmentation Network*



• No local context for each point!

- Hierarchical Feature Learning

- Increasing receptive field



512 filters of stride 1

160 filters of stride 2

48 filters of stride 2

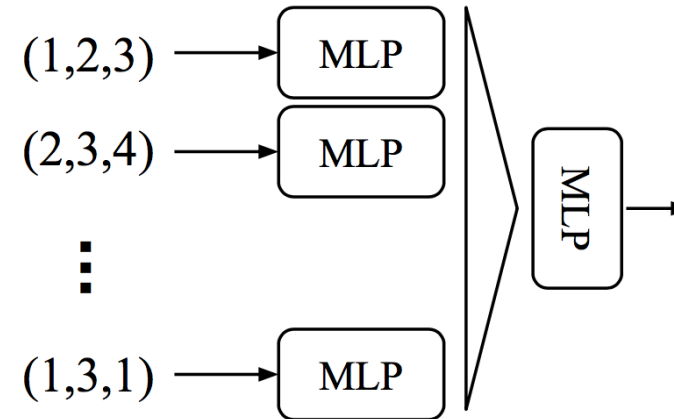3D voxel input

3D CNN (Wu et al.)

# Limitations of PointNet v1.0

- **Hierarchical Feature Learning**

- **Increasing receptive field**

**Global Feature Learning**
**Receptive field:**
**one point OR all points**



512 filters of stride 1

160 filters of stride 2

48 filters of stride 2

3D voxel input

V.S.

$(1,2,3) \longrightarrow$ MLP

$(2,3,4) \longrightarrow$ MLP

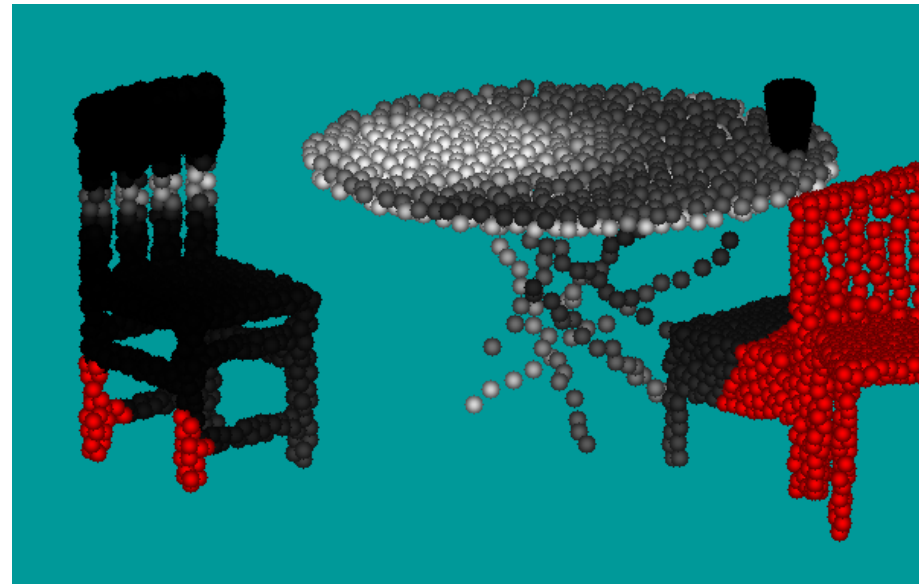$(1,3,1) \longrightarrow$ MLP

MLP

3D CNN (Wu et al.)

PointNet (vanilla) (Qi et al.)
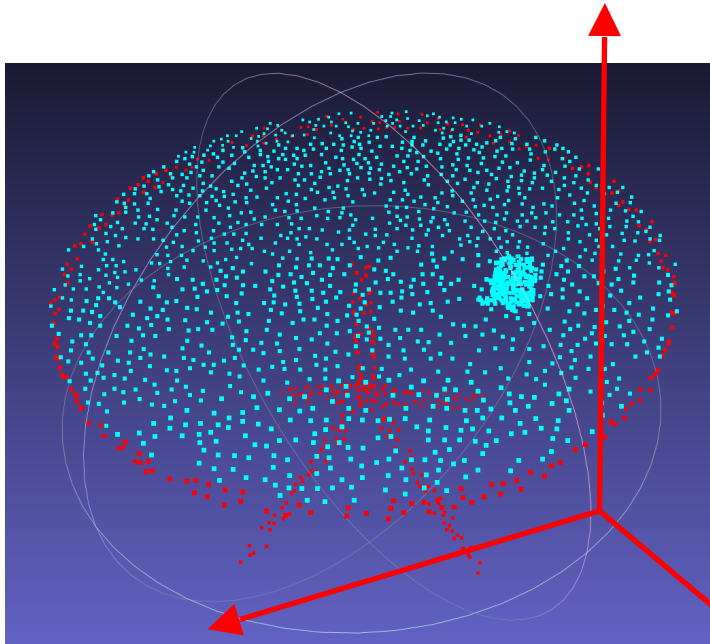
# Limitations of PointNet v1.0

Artifacts in segmentation tasks:



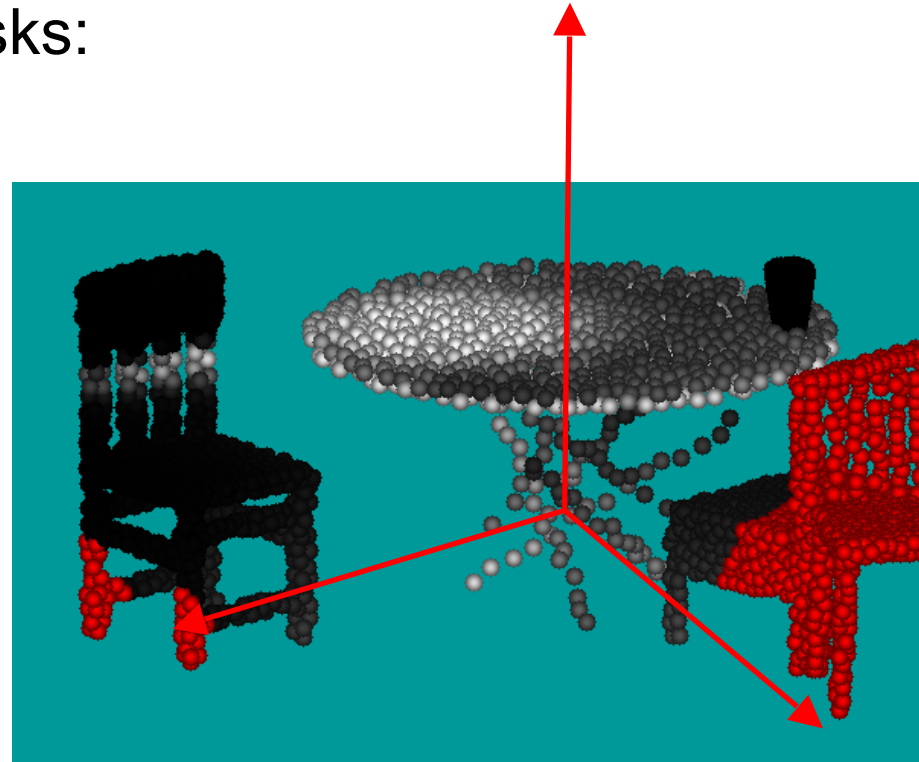Semantic segmentation in randomly translated table-cup scene.



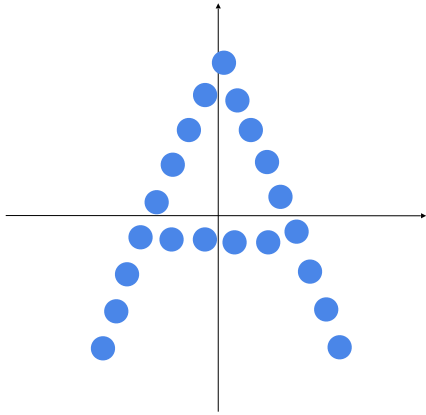Instance segmentation in table-chair-cup scene

Artifacts in segmentation tasks:



Semantic segmentation in randomly translated table-cup scene.

Instance segmentation in table-chair-cup scene

- Global feature depends on absolute XYZ!
- Hard to generalize to unseen point configurations

# Question

- How to learn local context feature for points?

- Use PointNet in local regions, aggregate local region features by PointNet again..
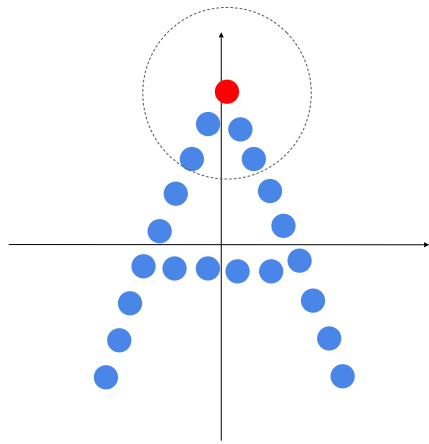
- 

- Hierarchical feature learning!

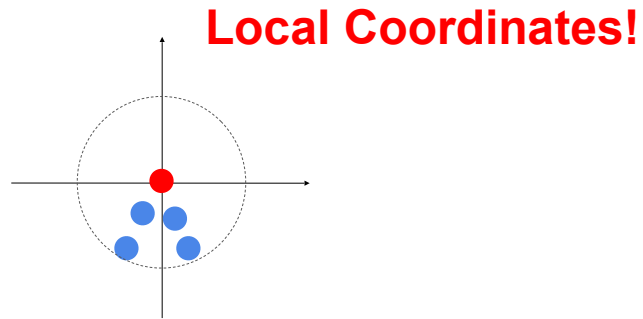# Multi-Scale PointNet for Hierarchical Feature Learning
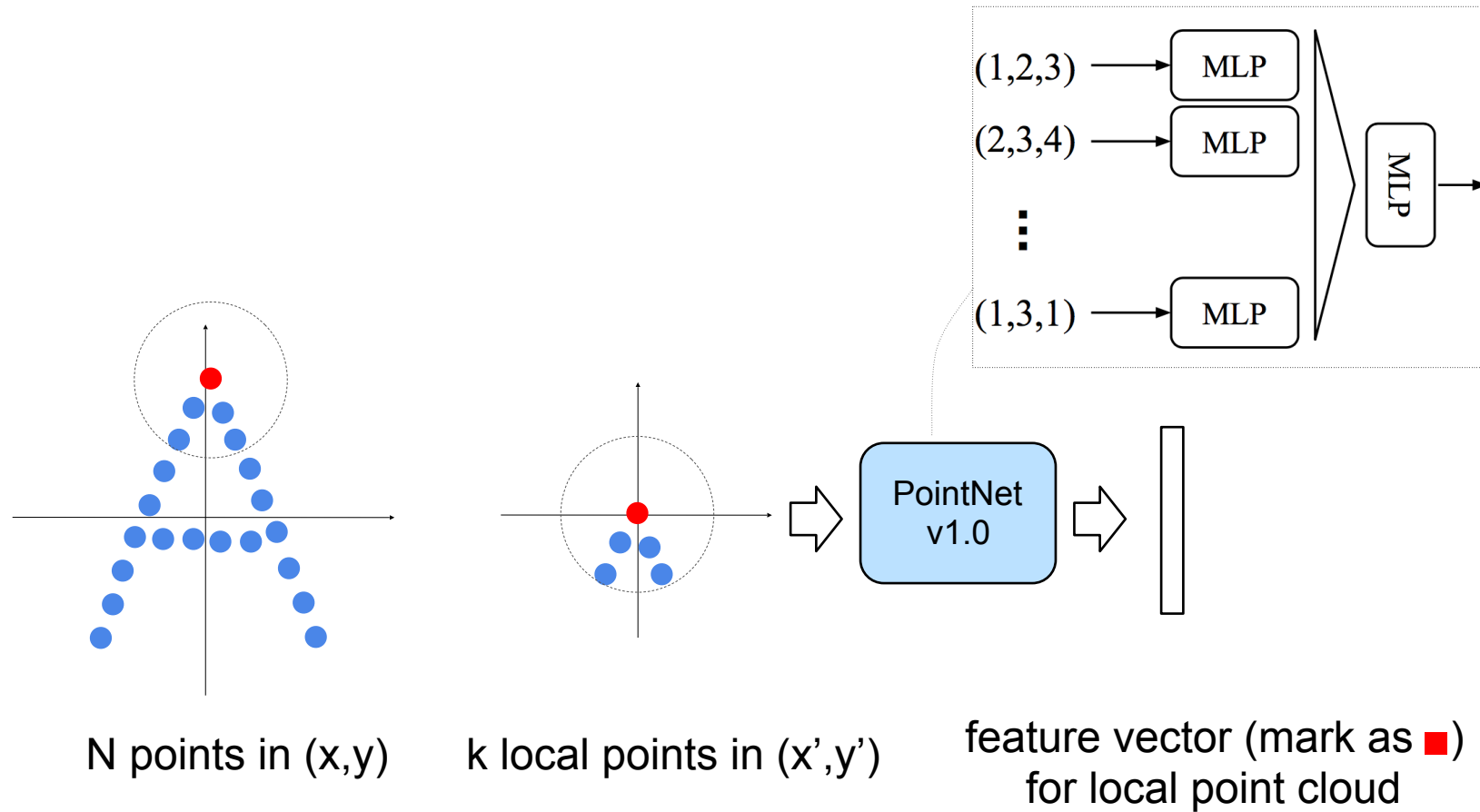
N points in (x,y)

# PointNet v2.0: Multi-Scale PointNet



Local Coordinates!

N points in (x,y)          k local points in (x',y')

# PointNet v2.0: Multi-Scale PointNet



(1,2,3) ⟶ MLP

(2,3,4) ⟶ MLP

⋮

(1,3,1) ⟶ MLP

MLP ⟶

PointNet v1.0

N points in (x,y)      k local points in (x',y')      feature vector (mark as ■) for local point cloud

**PointNet Module/Layer:** Farthest Point Sampling + Grouping + PointNet v1.0



N points in (x,y)          $N_1$ points in (x,y,**f**)

# PointNet v2.0: Multi-Scale PointNet



N points in (x,y) $\Rightarrow$ N₁ points in (x,y,**f**) $\Rightarrow$ N₂ points in (x,y,**f'**)

N points in (x,y)          N₁ points in (x,y,**f**)          N₂ points in (x,y,**f'**)

# PointNet v2.0: Multi-Scale PointNet



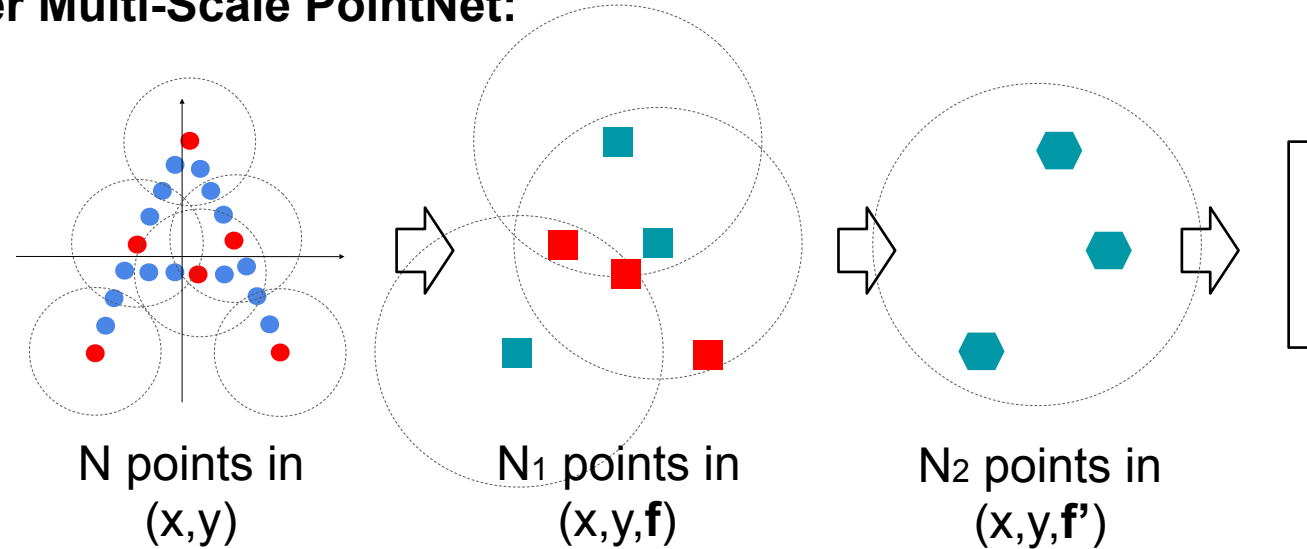N points in (x,y)        $N_1$ points in (x,y,**f**)        $N_2$ points in (x,y,**f'**)

1. Larger receptive field in higher layers ✓
2. Less points in higher layers (more scalable) ✓
3. Weight sharing ✓
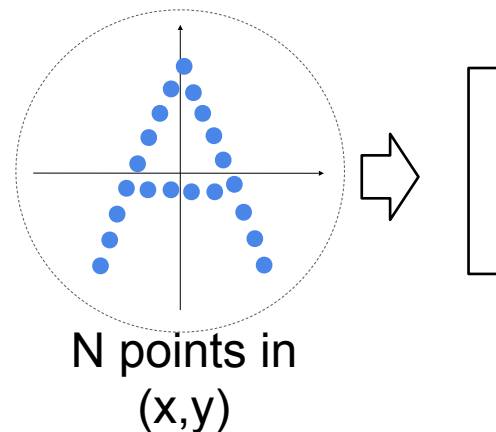4. Translation invariance (local coordinates in local regions) ✓

# Discussions on Multi-Scale PointNet

# Multi-Scale PointNet v.s. PointNet v1.0
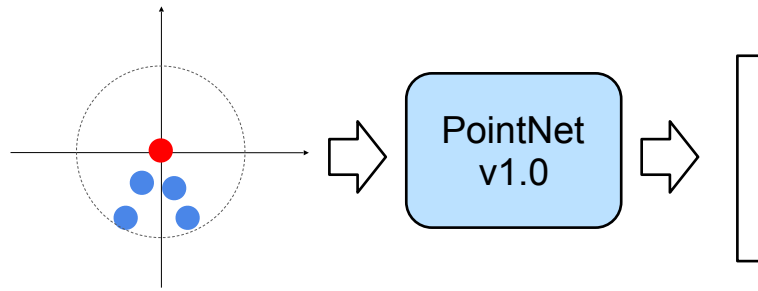
**Three-layer Multi-Scale PointNet:**



N points in
(x,y)

$N_1$ points in
(x,y,**f**)

$N_2$ points in
(x,y,**f'**)

**One-layer Multi-Scale PointNet <=> PointNet v1.0**



N points in
(x,y)

# PointNet Layer v.s. Convolution Layer



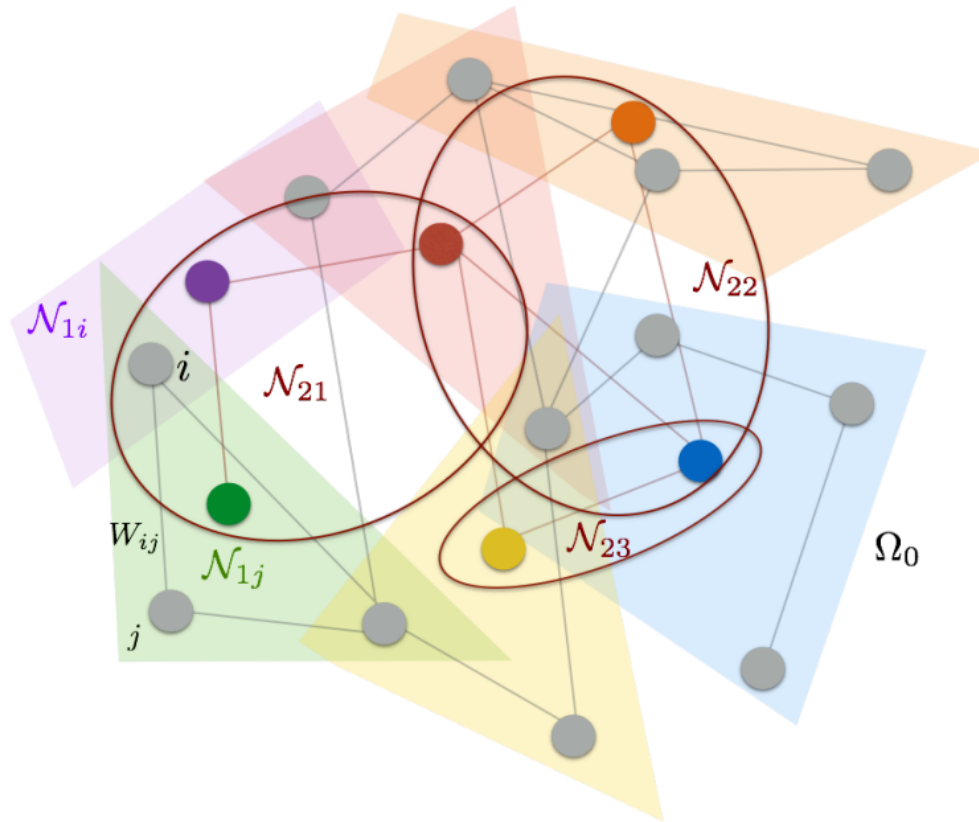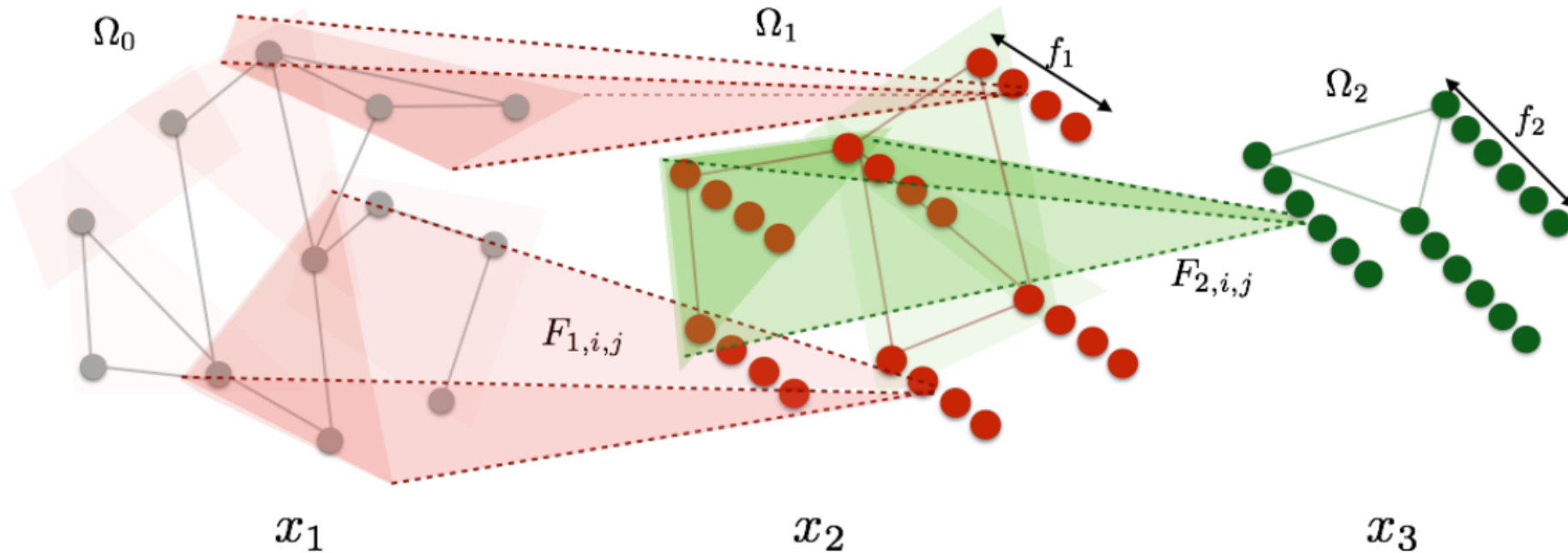|  | PointNet Layer | Convolution Layer |
|---|---|---|
| Input: | Point set | Dense array |
| Operation: | MLP + max pooling | Multiply and add |
| Neighbor-hood: | Distance query | Array index |

# Multi-Scale PointNet v.s. Graph CNN

- Unexpectedly strong relation with Graph CNN:



*Joan Bruna et al. Spectral Networks and Deep Locally Connected Networks on Graphs. ICLR 2014*
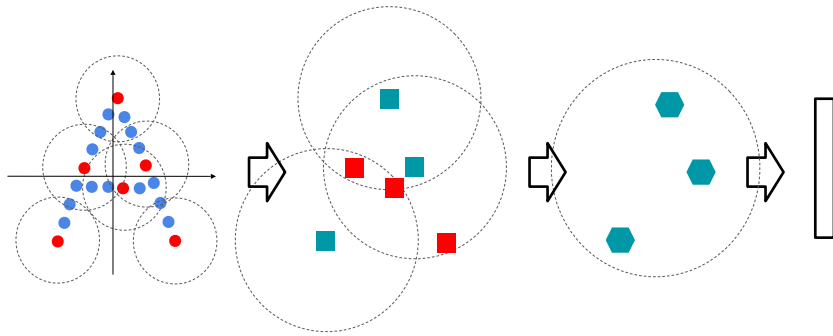
- Local feature extraction, graph coarsening, repeat..



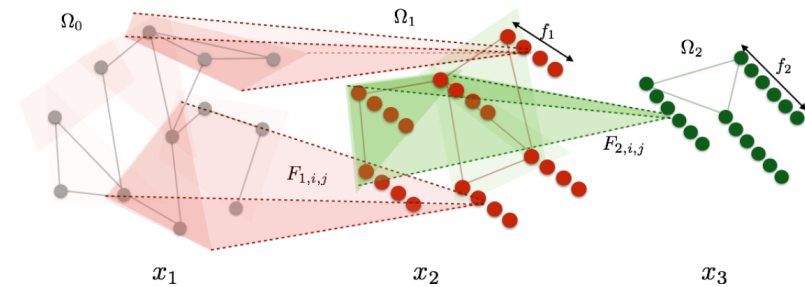*Joan Bruna et al. Spectral Networks and Deep Locally Connected Networks on Graphs. ICLR 2014*

- In Graph CNN's perspective:

- Multi-Scale PointNet defines

  1. Graph connectivity through Euclidean distance

  2. Graph coarsening by farthest point sampling

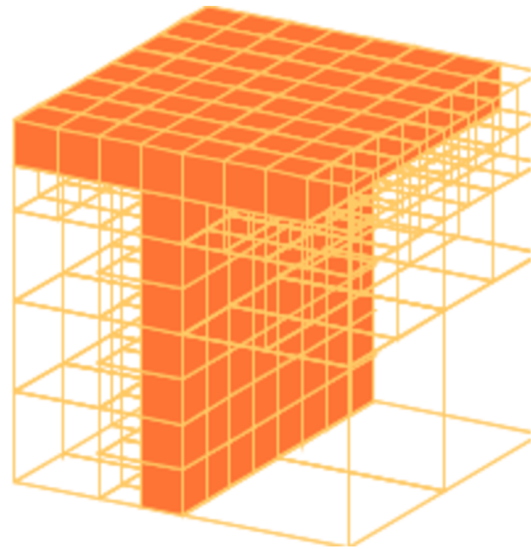  3. Local feature extraction with PointNet (v1.0)



Multi-scale PointNet

Graph CNN

OctNet in Graph CNN's perspective:

1. Both connectivity and graph coarsening are defined by the Octree.
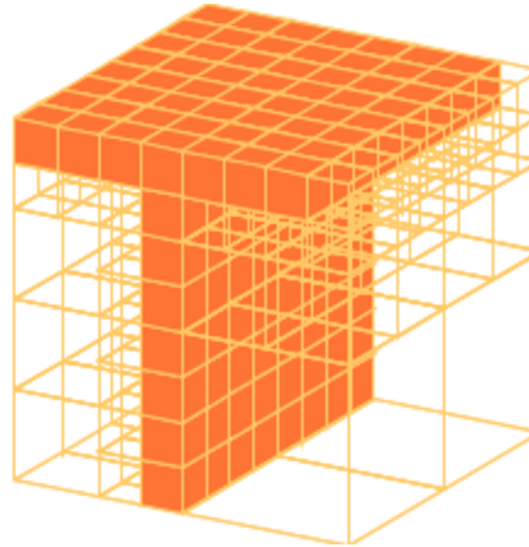2. Local feature extraction by convolution layer.



*OctNet: Learning Deep 3D Representations at High Resolutions*
*Gernot Riegler, Ali Osman Ulusoy and Andreas Geiger*

OctNet in Graph CNN's perspective:

<span style="color:red">In Multi-Scale PointNet</span>

1. Both connectivity and graph coarsening are defined by the Octree. <span style="color:red">By ground distance</span>
2. Local feature extraction by convolution layer. <span style="color:red">By PointNet (v1.0)</span>



*OctNet: Learning Deep 3D Representations at High Resolutions*
*Gernot Riegler, Ali Osman Ulusoy and Andreas Geiger*