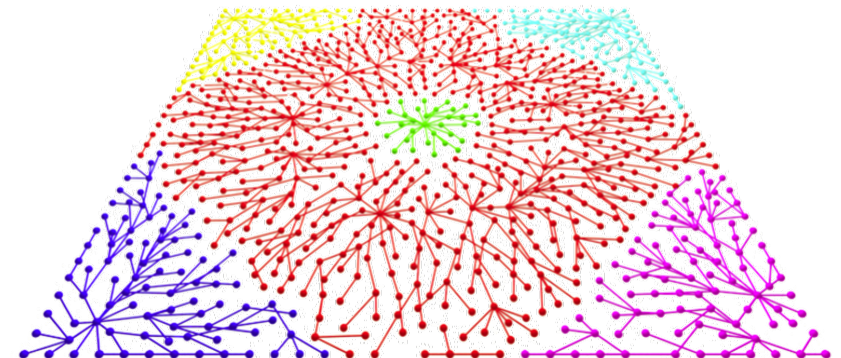


# CS233, CME251: Geometric and Topological Data Analysis

Leonidas Guibas  
Computer Science Department  
Stanford University



Lecture 14  
16 May 2022

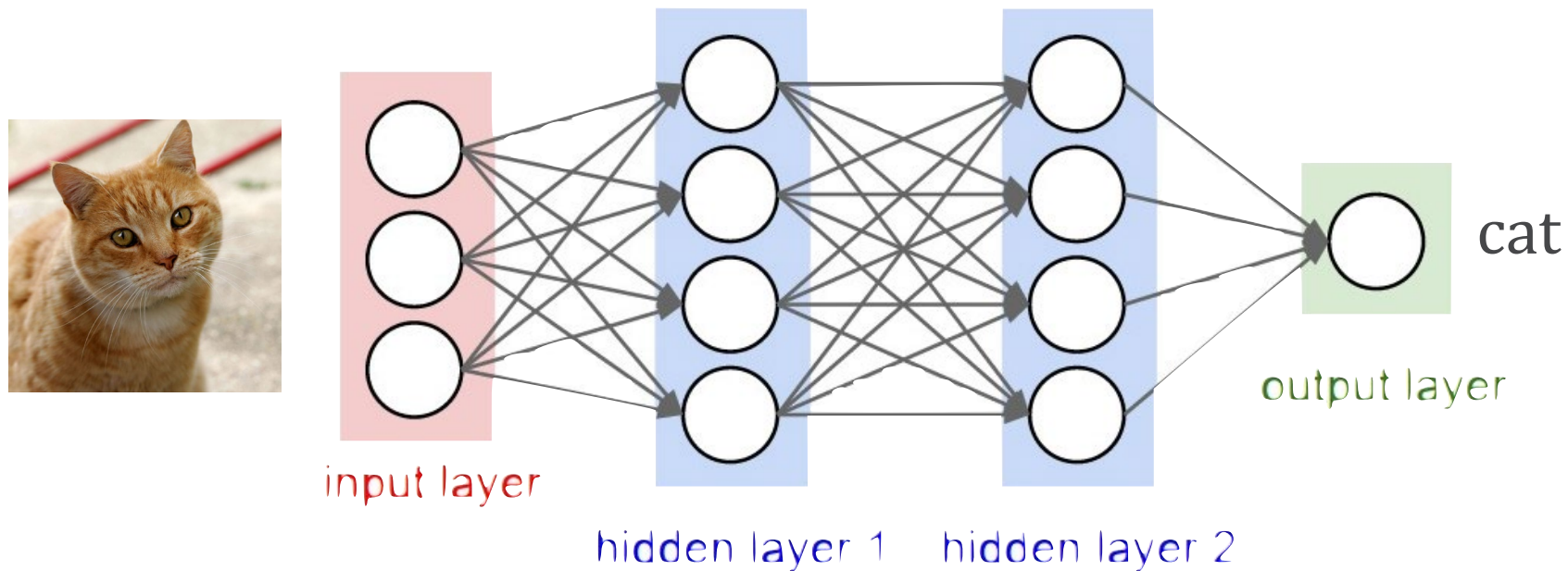


Last Time: Deep Nets, Multi-View and Volumetric Approaches to 3D, Graph/Mesh Neural Networks

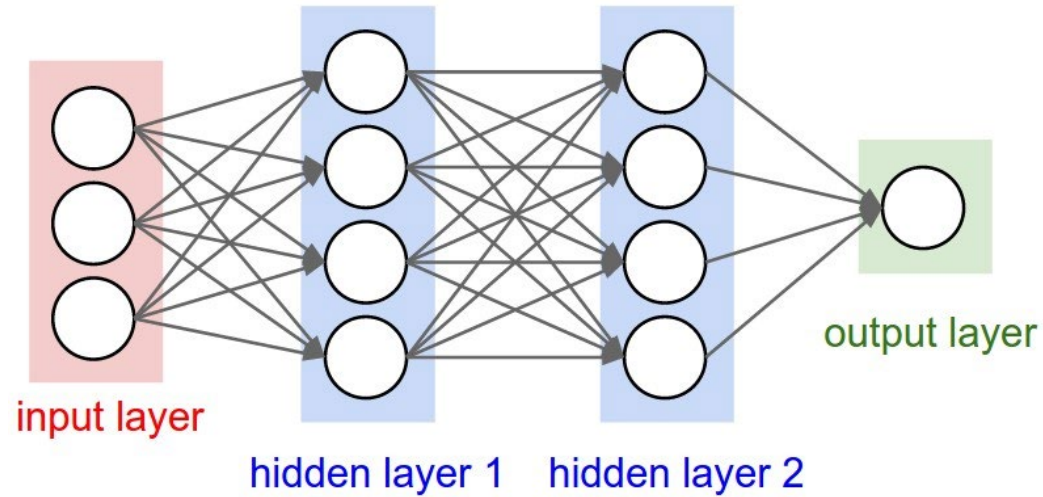
# Deep Learning

- Deep learning allows computational models that are composed of **multiple processing layers** to learn representations of data with **multiple levels of abstraction**.

*Deep Learning by Y. LeCun et al. Nature 2015*

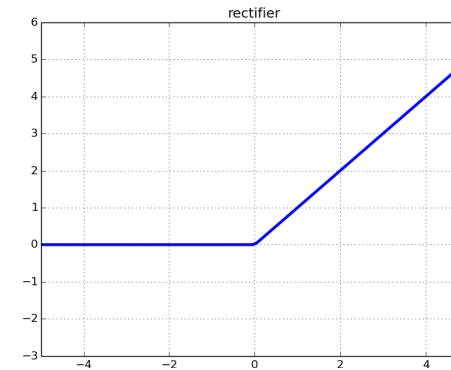


# Neural Networks Non-Linearities



$f$ : non-linear activation function

$$y = \text{ReLU}(x) = \max(\mathbf{0}, x)$$

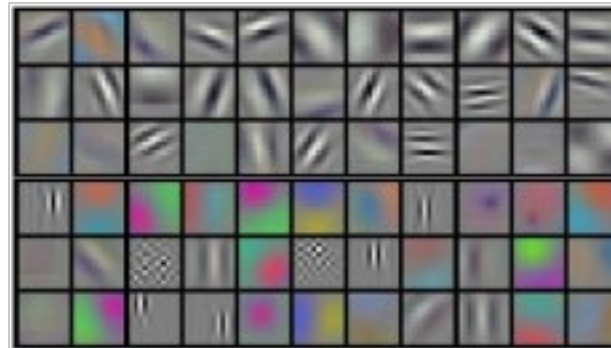
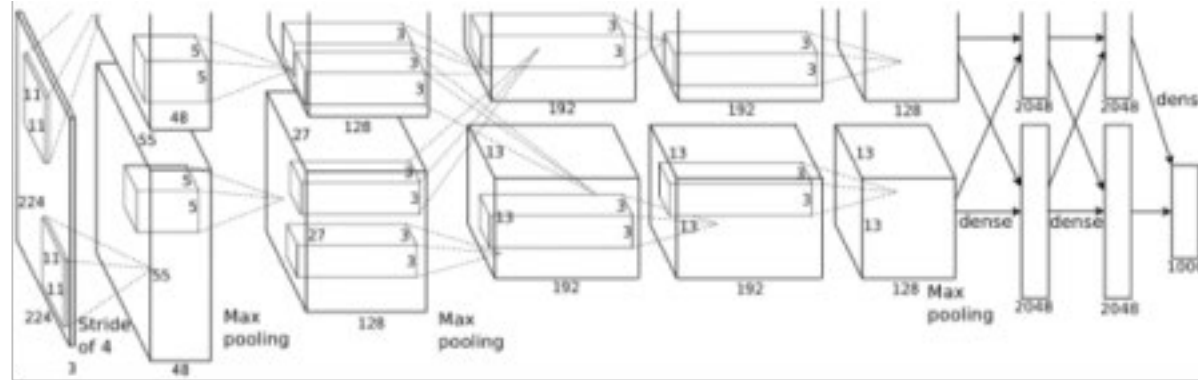
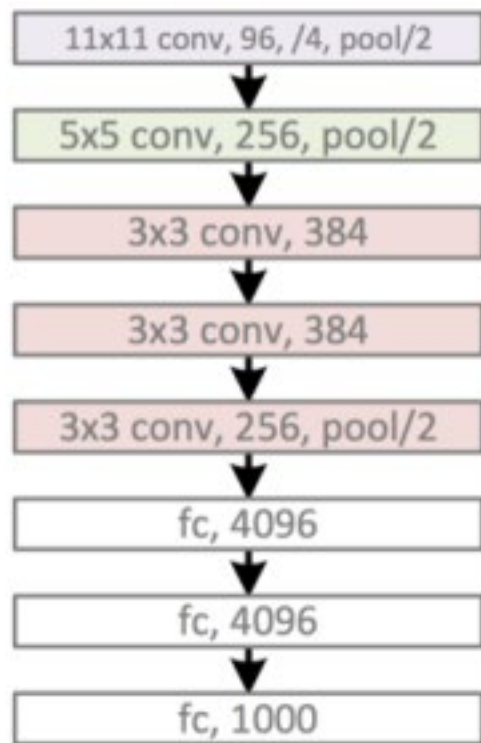


**Model:** Multi-Layer Perceptron (MLP)

$$y' = W_3 f(W_2 f(W_1 x + b_1) + b_2) + b_3$$

# Convolutional Neural Networks

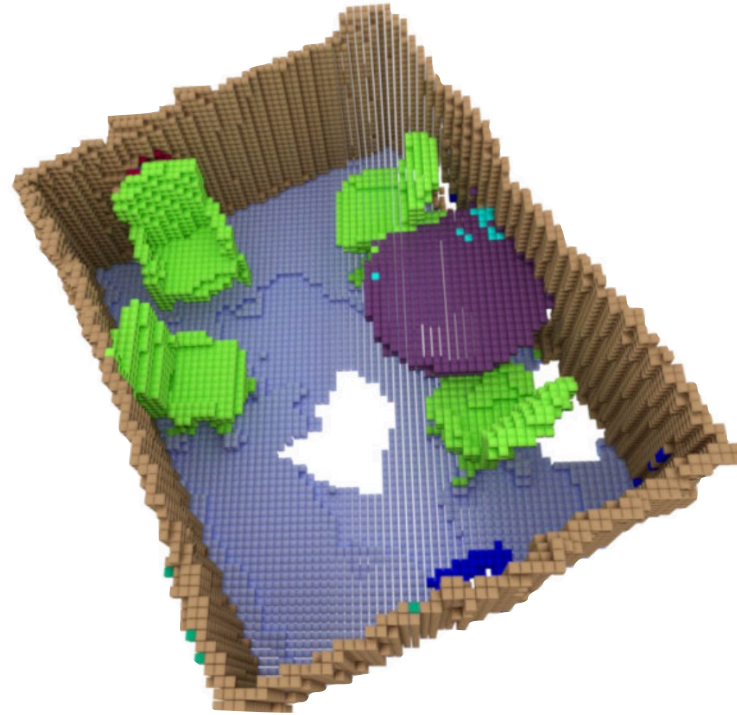
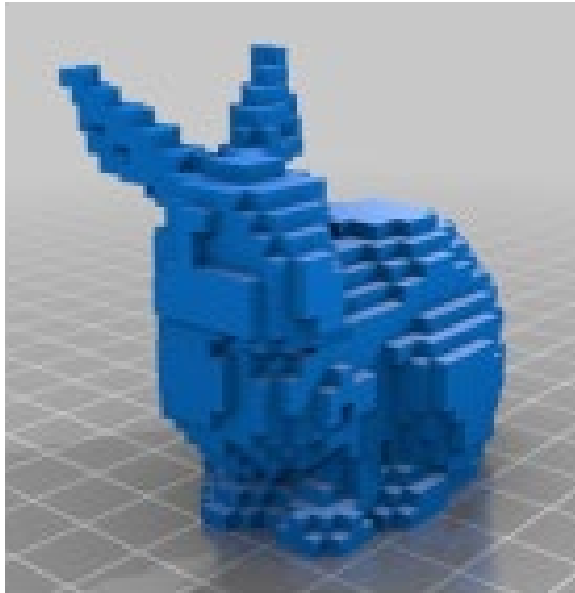
## AlexNet



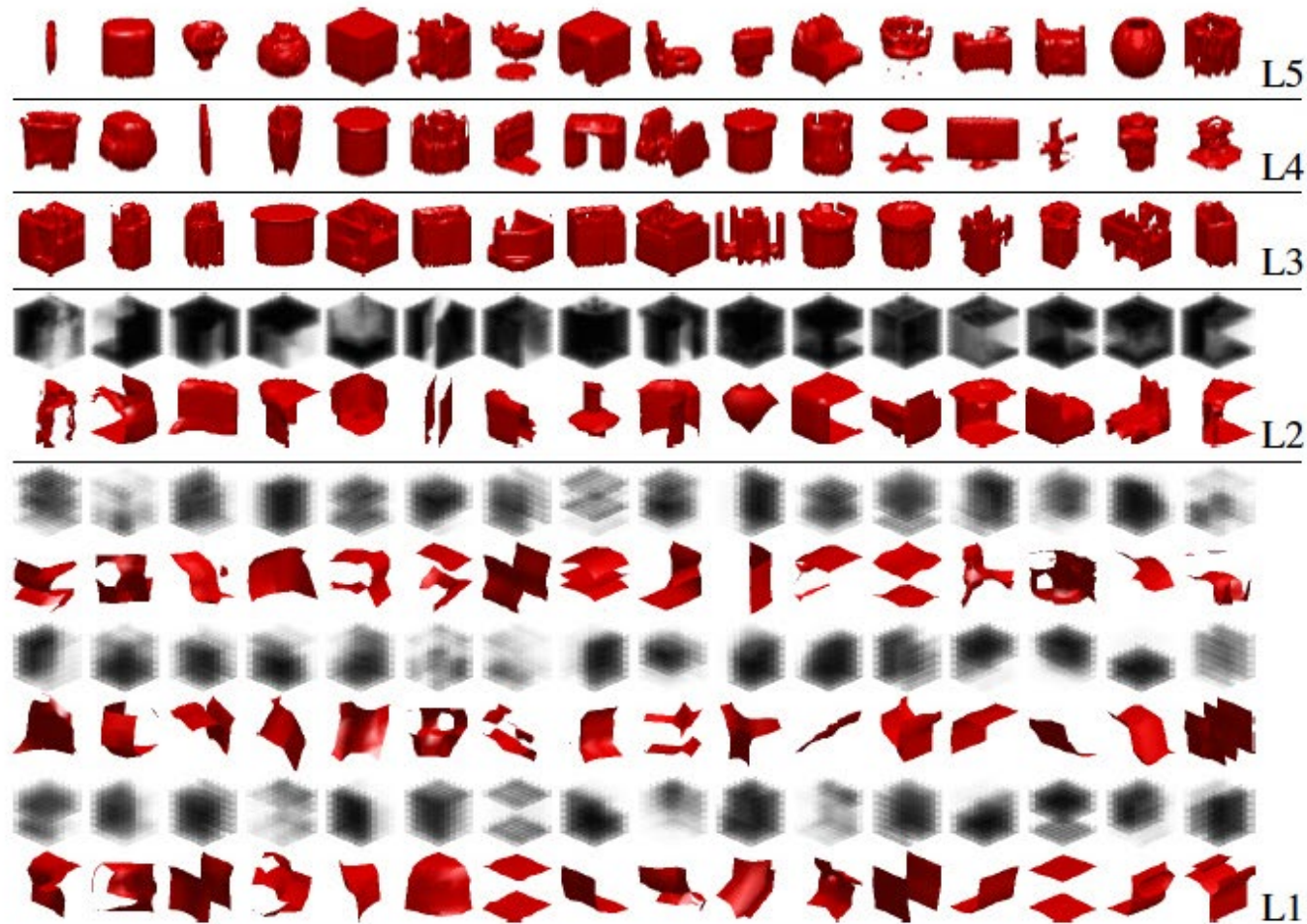
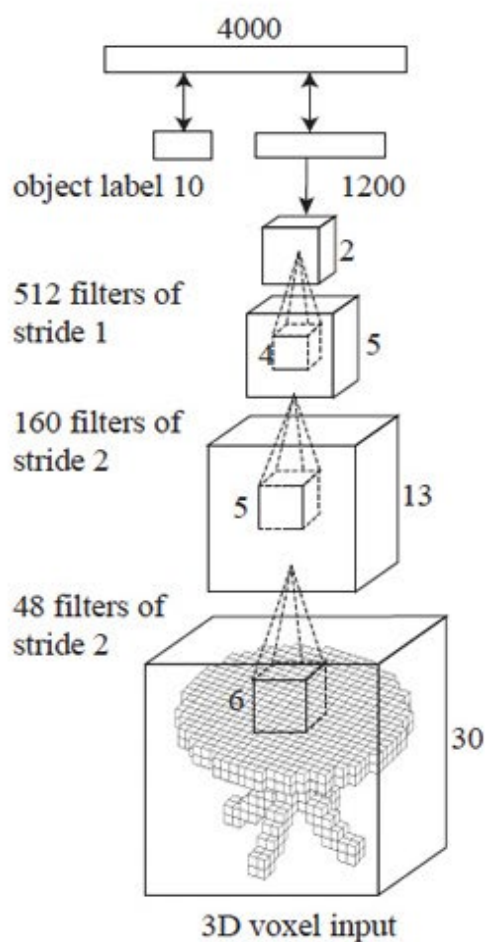
The first work that popularized Convolutional Networks in Computer Vision

# Volumetric Representation: 3D Geometry

Volumetric Representations achieve a similar discretization of the 3D space.

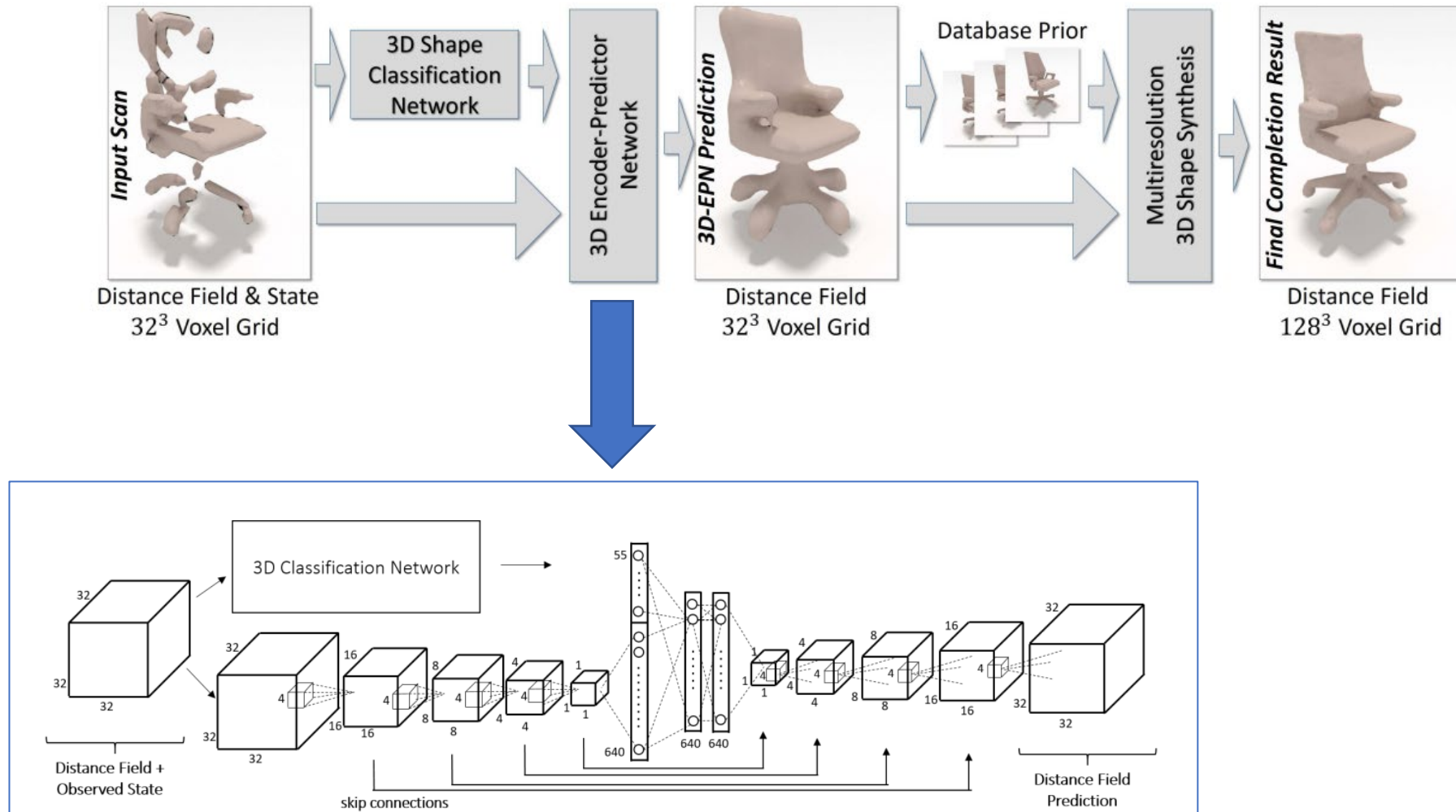


# Volumetric Deep Learning: Classification



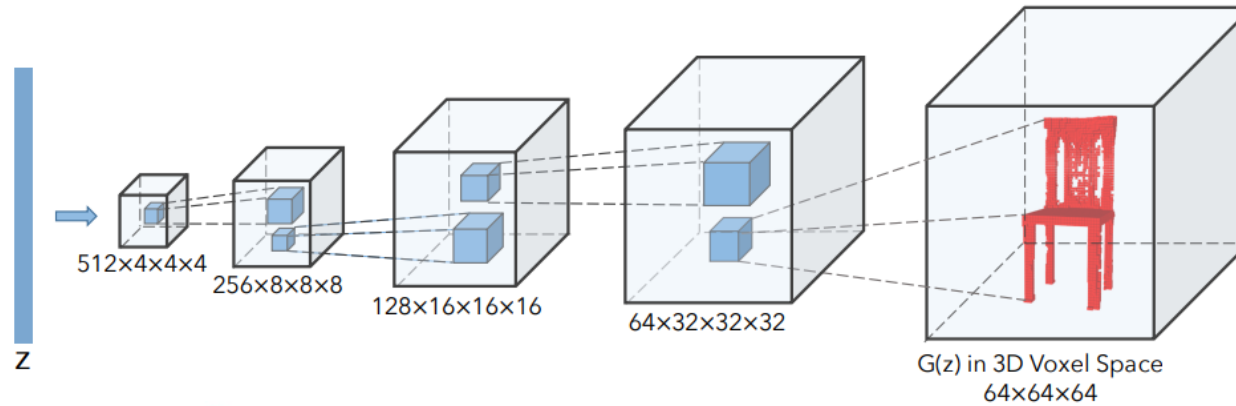
For each neuron, average the top 100 training examples with highest responses

# Volumetric Deep Learning: Completion

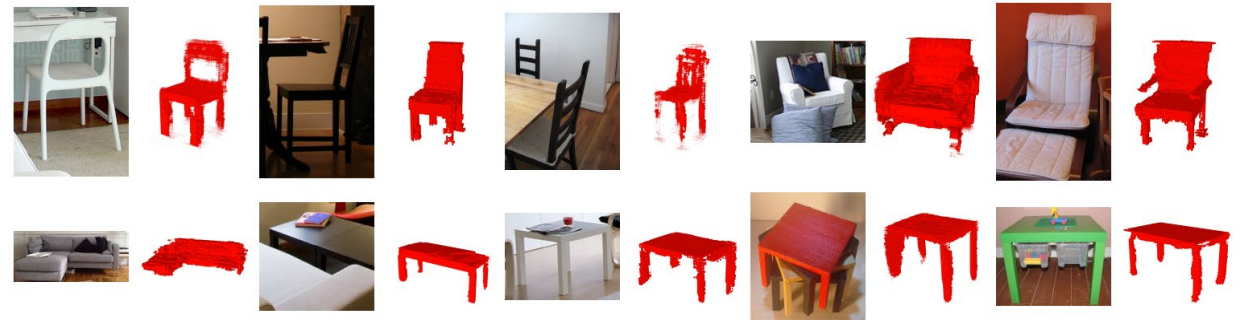


Shape Completion using 3D-Encoder-Predictor CNNs and Shape Synthesis, Angela Dai, Charles Ruizhongtai Qi, Matthias Nießner, CVPR 2017

# Volumetric Deep Learning: Generation



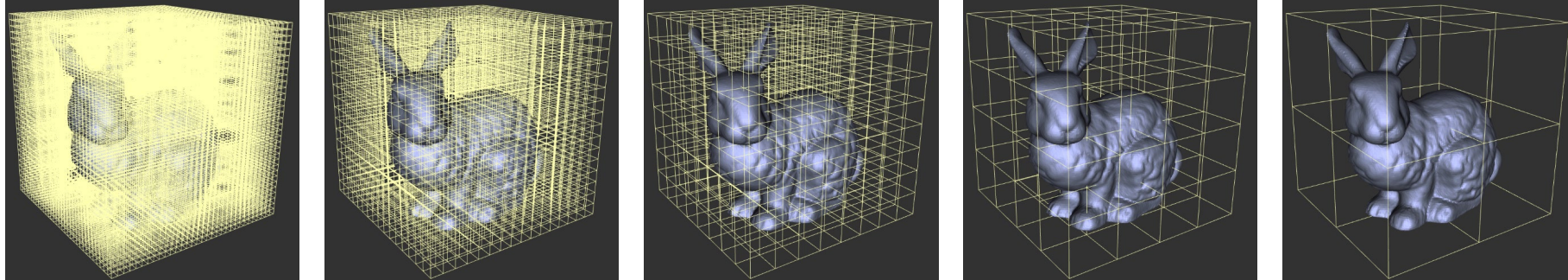
Random Samples



Singe-view 3D Reconstruction

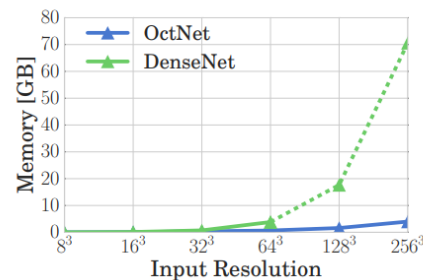
Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling, Jiajun Wu\* Chengkai Zhang\* Tianfan Xue, William T. Freeman, Joshua B. Tenenbaum, NeurIPS 2016

# Sparse Volumetric Representation

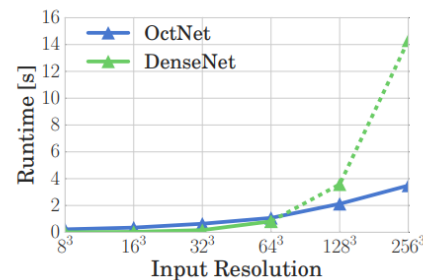


# Volumetric Deep Learning with Octrees

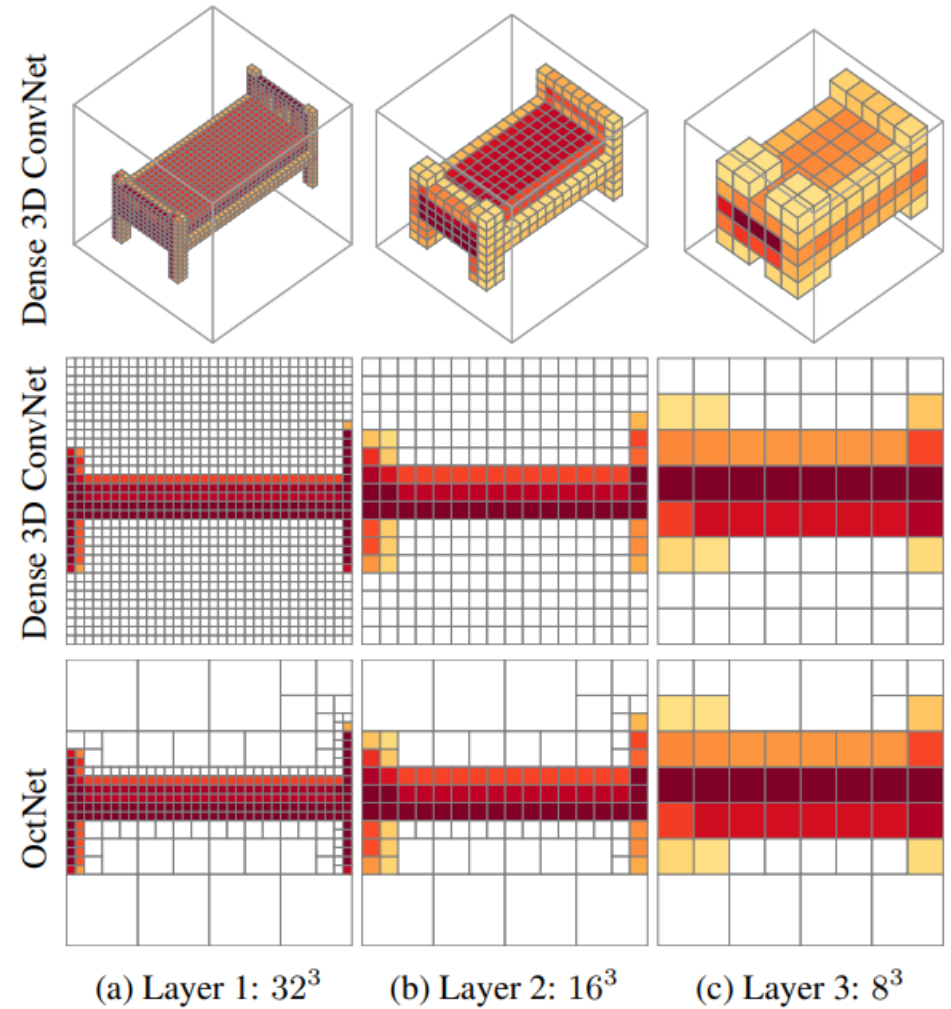
- **OctNet is a 3D convolutional network** that exploits this sparsity property, by hierarchically partitioning the 3D space into a set of unbalanced octrees.
- Instead of representing the entire high resolution 3D input with a single unbalanced octree, they propose to use multiple shallow octrees of maximal depth.
- Shallow octrees can be efficiently encoded using bit-strings
- Using these representation, they define various convolutional layers, pooling and unpooling layers,



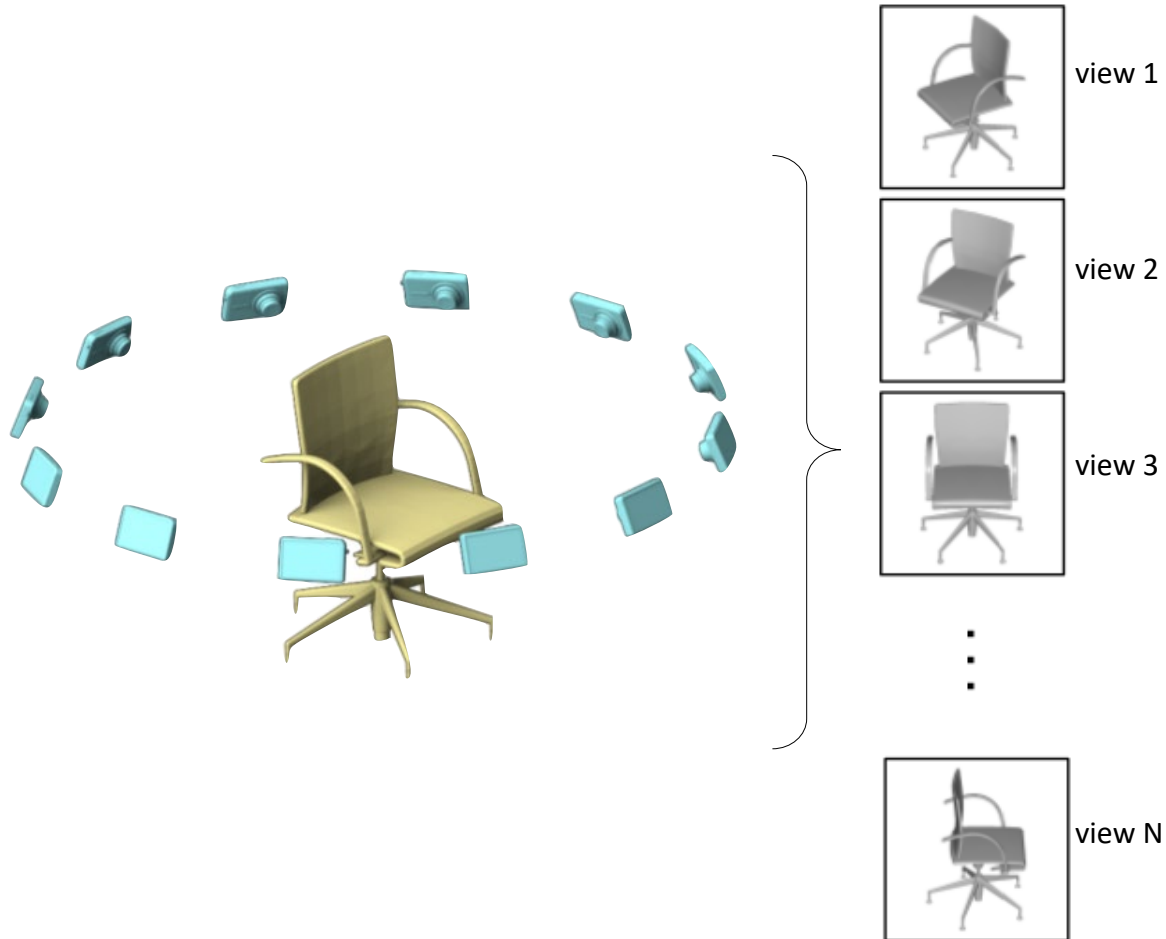
(a) Memory



(b) Runtime

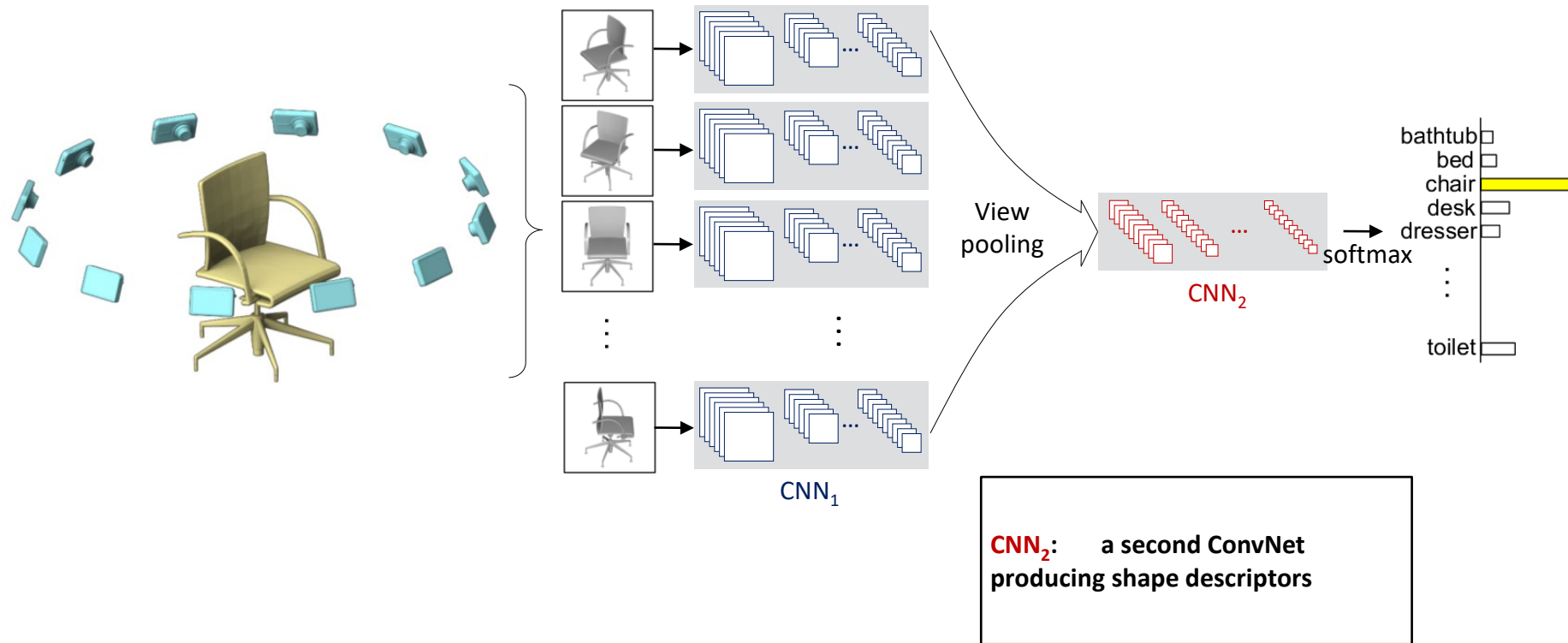


# Multi-view Representation



+ Powerful  
2D CNNs

# Multi-view CNNs: Classification



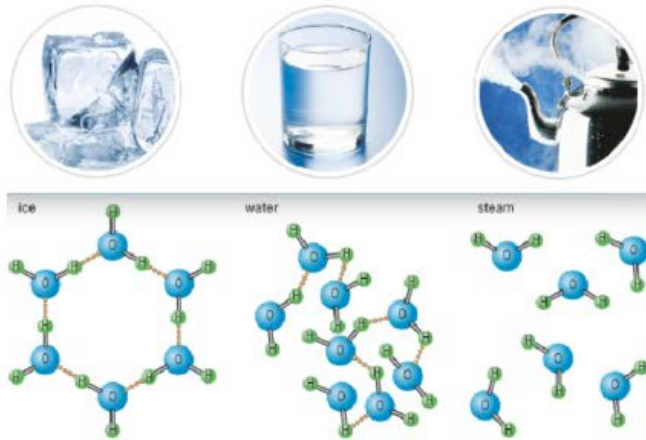
Hang Su, Subhransu Maji, Evangelos Kalogerakis, Erik Learned-

Miller, "Multi-view Convolutional Neural Networks for 3D

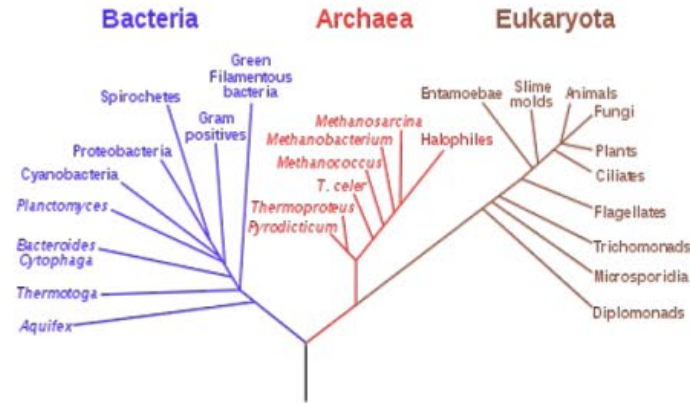
Shape Recognition", *Proceedings of ICCV 2015*

Image Credits: Hang Su

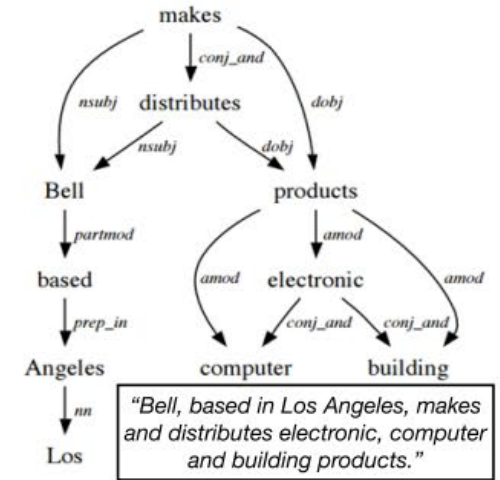
# Why do we need Graph Neural Networks?



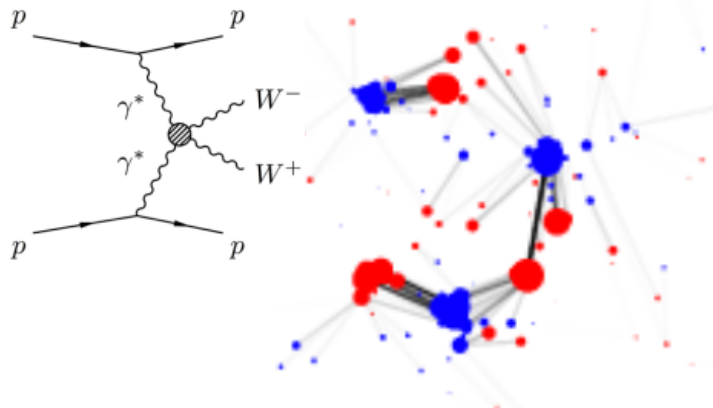
Molecules Predictions



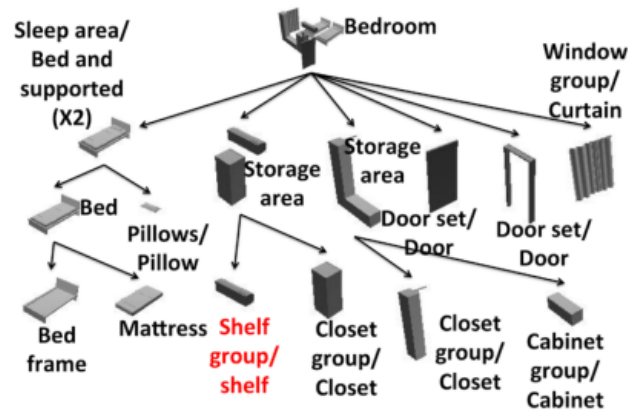
Biological Species Categorizations



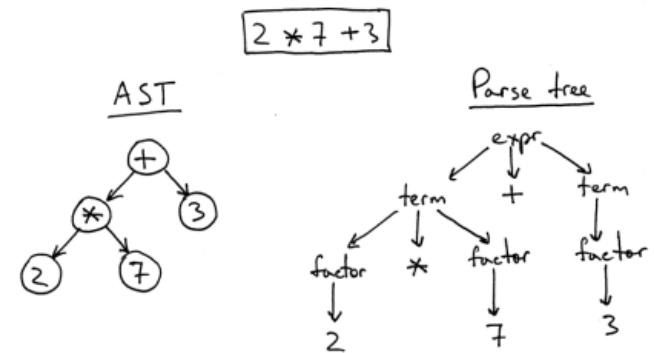
Natural Language Processing



Sub-atomic particles



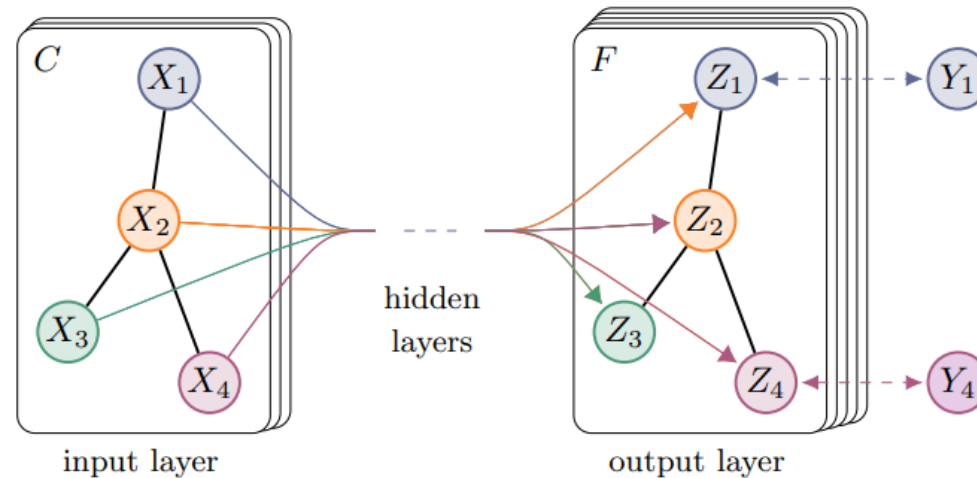
Everyday scenes



Code analysis

# Graph Convolutional Networks

Starting from a graph  $G = (V, E, X)$ , where  $V$  are the  $N$  nodes,  $E$  are the edges and  $X$  is the  $N \times D$  feature matrix for every node in the graph and a representation of its graph structure in a matrix form (such as the adjacency matrix) we want to learn a **node-level representation  $Z$** , defined as a  $N \times F$  feature matrix, where  $F$  are the output features per node.



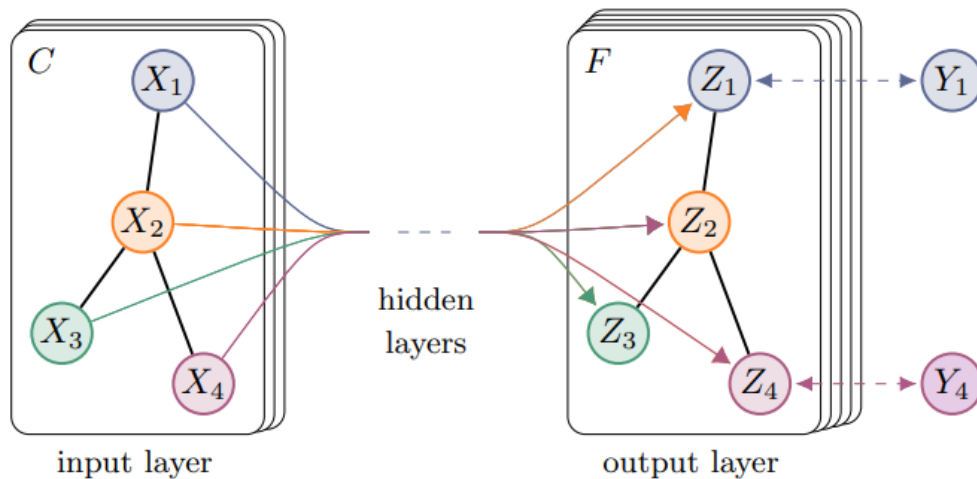
To this end we define a neural network layer as

$$H^{(l+1)} = f(H^{(l)}, A)$$

where  $H^0 = X$  and  $H^{(l)} = Z$  and  $L$  is the number of layers

# Graph Convolutional Networks

Starting from a graph  $G = (V, E, X)$ , where  $V$  are the  $N$  nodes,  $E$  are the edges and  $X$  is the  $N \times D$  feature matrix for every node in the graph and a representation of its graph structure in a matrix form (such as the adjacency matrix) we want to learn a **node-level representation  $Z$** , defined as a  $N \times F$  feature matrix, where  $F$  are the output features per node.



Let us assume the following very simple form of a layer

$$f(H^{(l)}, A) = \sigma(AH^{(l)}W^{(l)})$$

- **Multiplying with  $A$ , means to sum up all the feature vectors for all neighboring nodes** but not the node itself (unless there are self loops in the graph. **To enforce self-loops in the graph, we replace  $A$  with  $A+I$**

$$f(H^{(l)}, A) = \sigma((A + I)H^{(l)}W^{(l)})$$

- Instead of directly using  $A+I$ , we normalize it

$$f(H^{(l)}, A) = \sigma(D^{-\frac{1}{2}} \hat{A} D^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

where  $\hat{A} = A + I$

# Geodesic Convolutional Neural Networks

The key idea is to **employ a local system of geodesic polar coordinates** constructed at around a point  $x$ . **The radial coordinates** is constructed as  $p$ -level sets  $\{x' : d_x(x, x') = p\}$  of the geodesic distance (shortest path) for every  $p \in [0, p_0]$  where  $p_0$  is the radius of the geodesic disc



**Examples of Local Geodesic Patches**

We can define the **bijective map from the manifold into the local geodesic polar coordinates  $(\rho, \theta)$**  around point  $x$  as:

$$\Omega(x) : B_{p_0}(x) \rightarrow [0, p_0] \times [0, 2\pi]$$

We also define  $(D(x)f)(p, \theta) = (f \circ \Omega^{-1}(x))(p, \theta)$

where  $D(x)$  is called the **patch operator** that is used for **mapping the values of the function  $f$**  at a neighborhood of the point  $x$  into the local polar coordinates  $(\rho, \theta)$ .

# Patch Operator

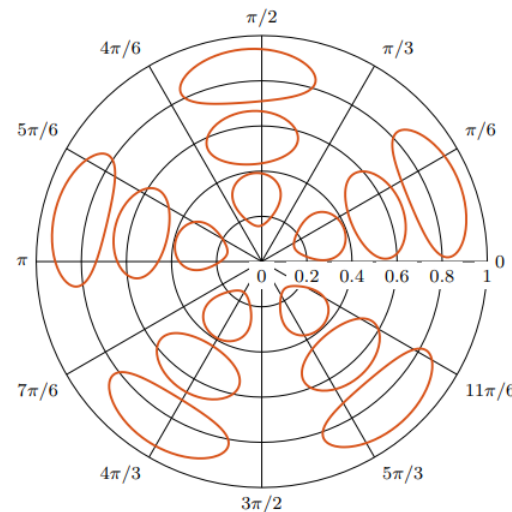
The **patch operator** maps the values of the function  $f$  at a neighborhood of the point  $x$  into the local polar coordinates  $(\rho, \theta)$ .

$$(D(x)f)(p, \theta) = \int_{\mathcal{X}} w_{p, \theta}(x, y) f(y) dy$$

$$w_{p, \theta} = \frac{v_p(x, x') v_\theta(x, x')}{\int_{\mathcal{X}} v_p(x, x') v_\theta(x, x') dx'}$$

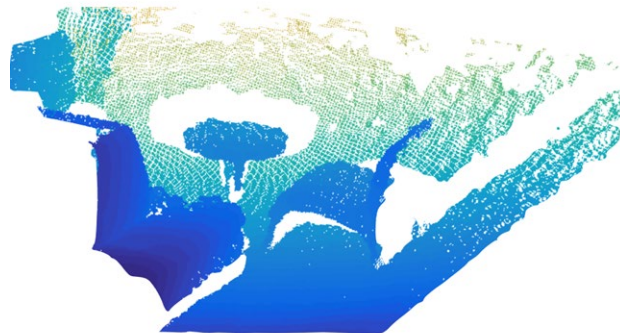


**Examples of Local Geodesic Patches**

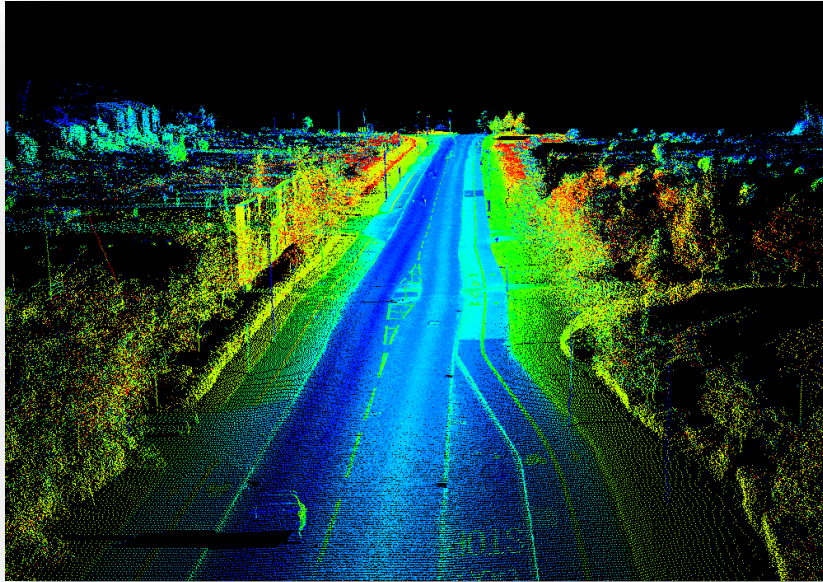


# Today: Deep Learning on Point Cloud Data

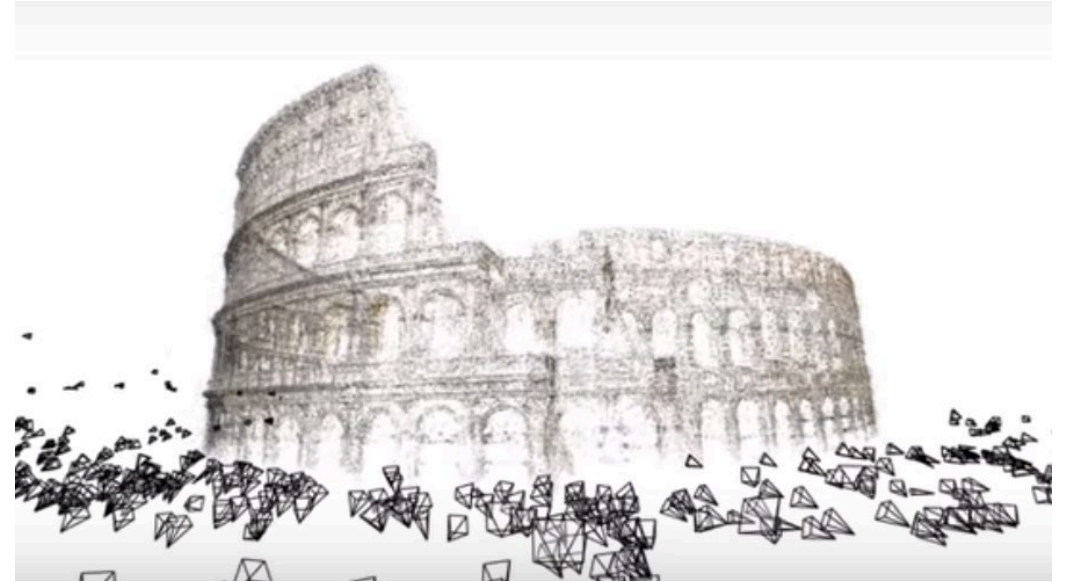
Non-regular, disconnected 3D data



# 3D Point Clouds from Many Sensors



Lidar point clouds (LizardTech)



Structure from motion (Microsoft)

Depth camera (Intel, Microsoft, Google)

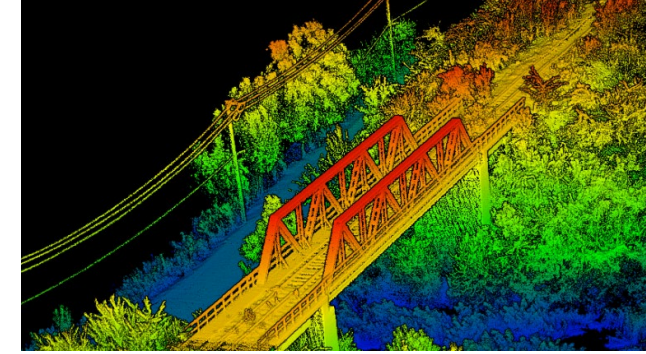


# 3D Point Cloud Data

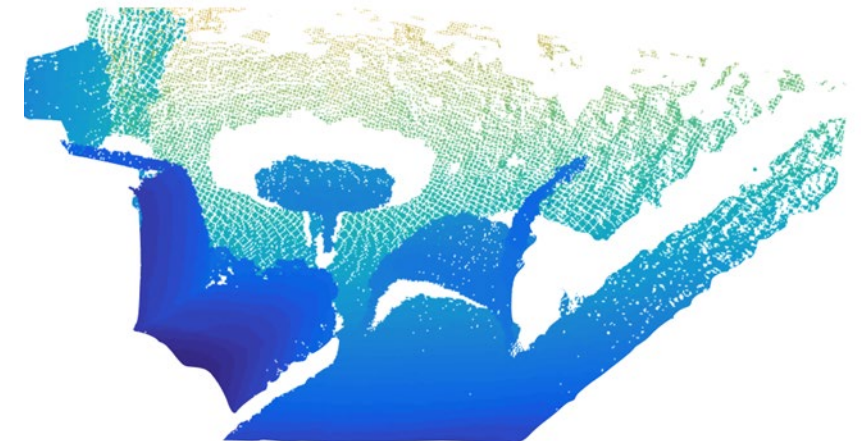
- Close to raw sensor data
- Representationally simple
- Irregular neighborhoods
- Variable density



LiDAR

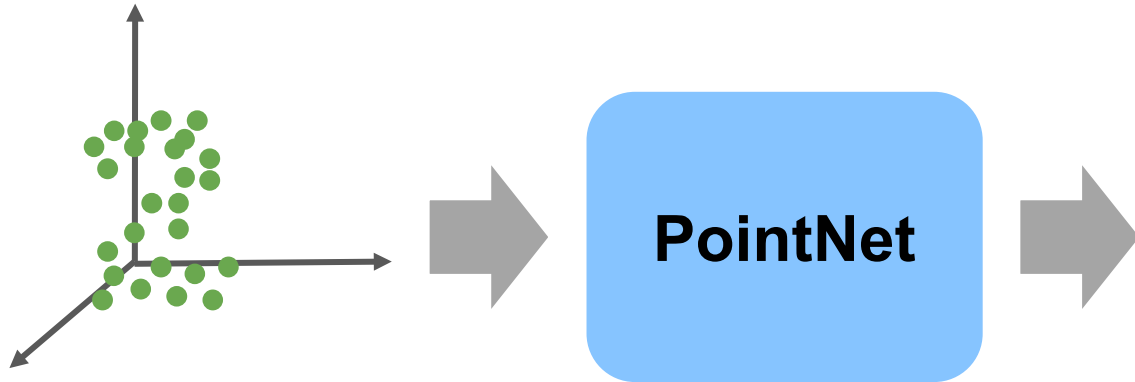


Depth Sensor



**Point Cloud**

# Deep Nets for PCs: PointNet and PointNet++



*Object Classification*

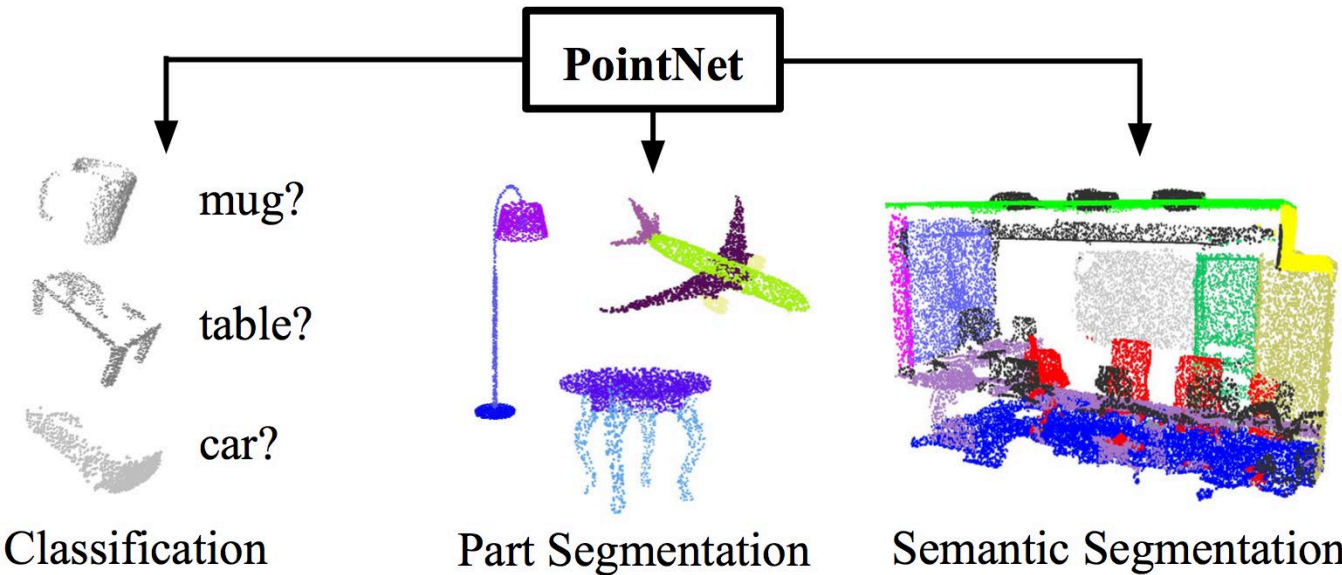
*Object Part Segmentation*

*Semantic Scene Parsing*

...

**End-to-end learning** for irregular point data

**Unified framework** for various tasks



Charles R. Qi, Hao Su, Kaichun Mo, Leonidas J. Guibas.  
*PointNet: Deep Learning on Point Sets for 3D Classification  
and Segmentation.* (CVPR'17)

# Invariances

*The model has to respect key desiderata for point clouds:*

## **Point Permutation Invariance**

Point cloud is a set of **unordered** points

## **Spatial Transformation Invariance**

Point cloud **rigid motions** should not alter classification results

## **Sampling Invariance**

Output a function of the underlying geometry and **not the sampling**

# Permutation Invariance: Symmetric Functions

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

## Examples:

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

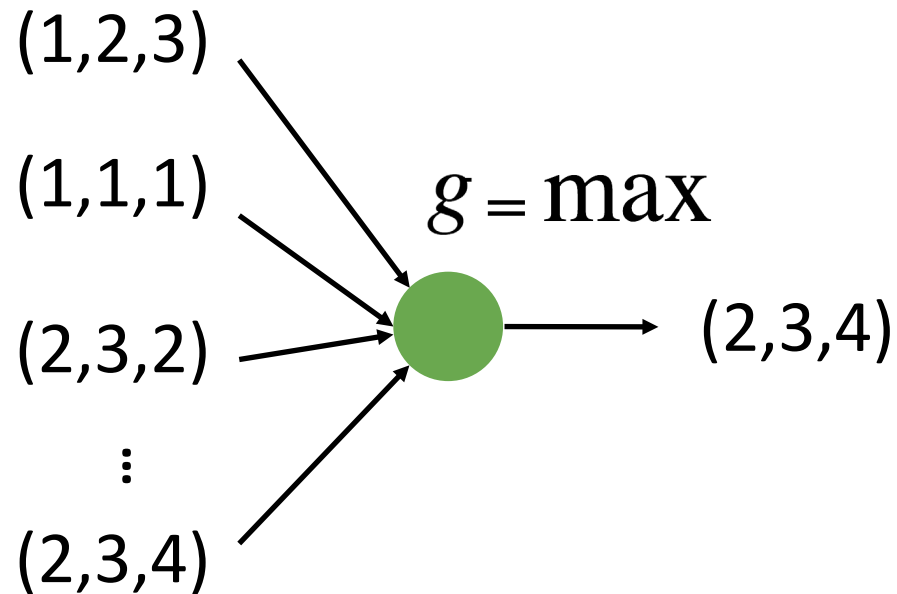
$$f(x_1, x_2, \dots, x_n) = x_1 + x_2 + \dots + x_n$$

...

**How can we construct a universal family of symmetric functions by neural networks?**

# Construct Symmetric Functions by Neural Networks

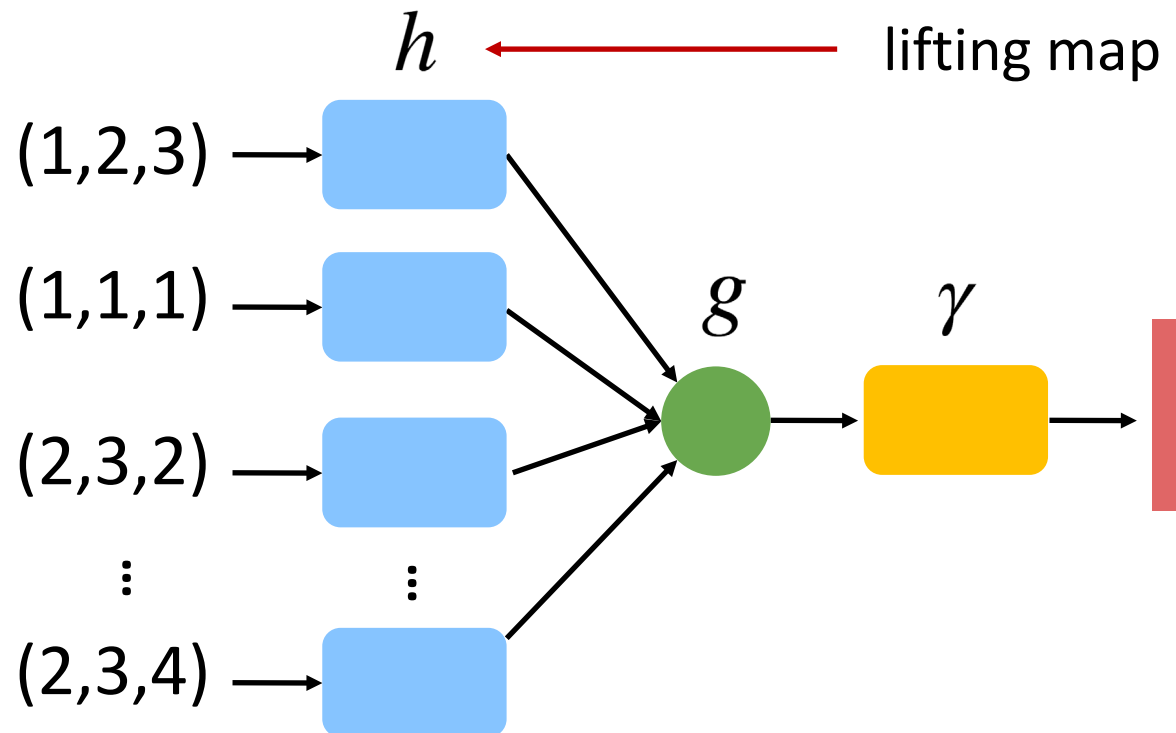
Simplest form: directly aggregate all points with a symmetric operator  $g$   
**Just discovers simple extreme/aggregate properties of the geometry.**



# Construct Symmetric Functions by Neural Networks

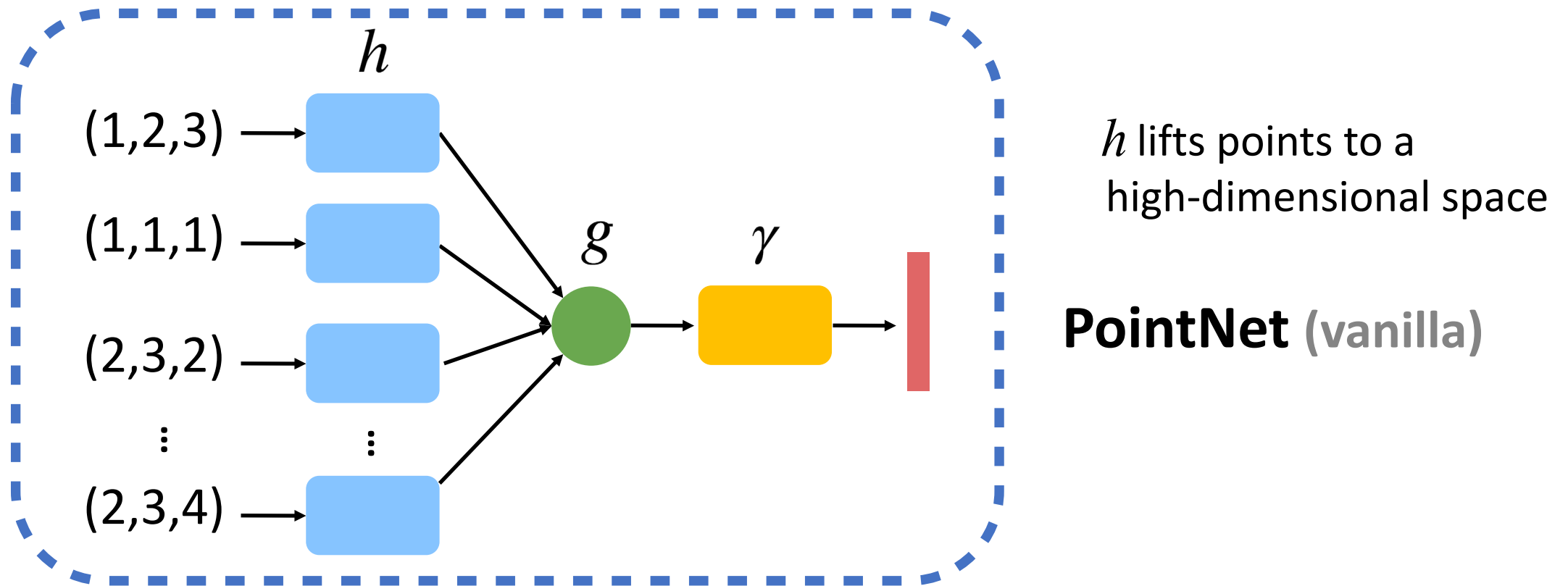
Embed points in a high-dim space before aggregation.

**Aggregation in the (redundant) high-dim space encodes more interesting properties of the geometry.**

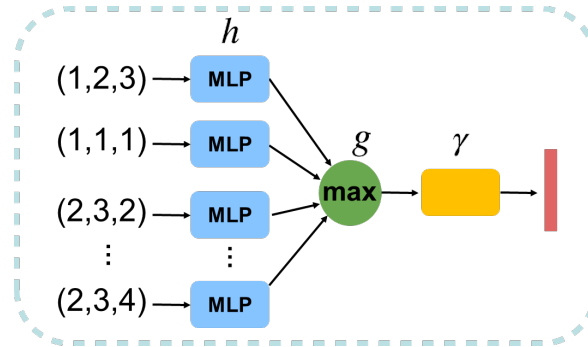


# Construct Symmetric Functions by Neural Networks

$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$  is symmetric if  $g$  is symmetric



# Symmetric Functions: Polynomials



$$2 \sum_{i \neq j} x_i x_j = \left( \sum_i x_i \right)^2 - \sum_i x_i^2$$

$$\sum_{i \neq j} (x_i - x_j)^2 = 3 \sum_i x_i^2 - \left( \sum_i x_i \right)^2$$

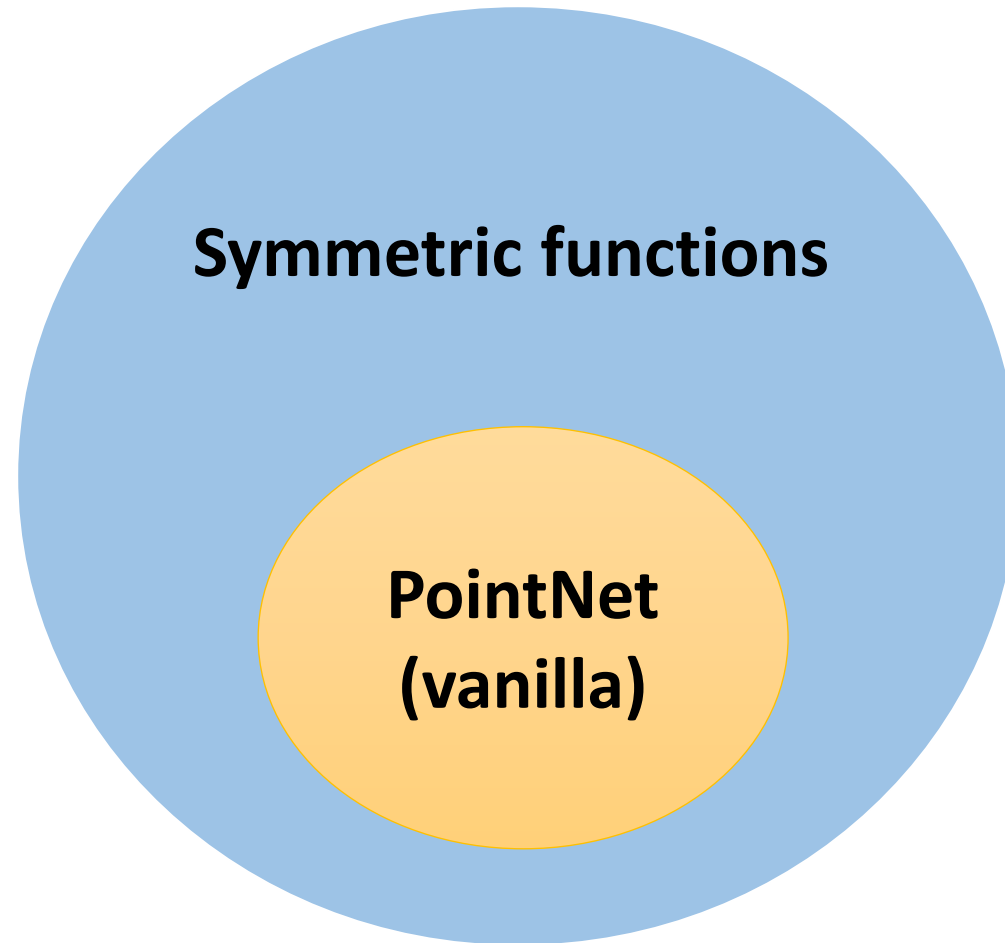
- In fact, **any** symmetric polynomial in the  $x_i$  can be expressed as a polynomial in sums of the form

$$\sum_i x_i^k$$

and can be computed by

$$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$$

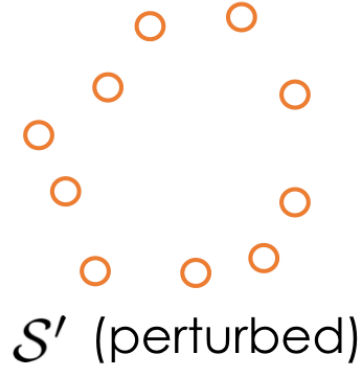
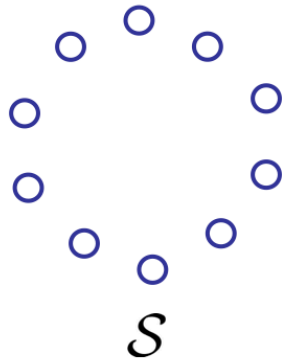
# What Symmetric Functions Can Be Constructed By PointNet?



# PointNet as a Universal Approximation to Set Functions

## Hausdorff continuous:

$f: 2^x \rightarrow \mathbb{R}$  is a continuous set function w.r.t Hausdorff distance



if  $d_{Hausdorff}(S, S') \approx 0$ , then  $f(S) \approx f(S')$

## Theorem

A Hausdorff continuous set function  $f: 2^x \rightarrow \mathbb{R}$  can be arbitrarily approximated by PointNet.

$$\left| f(S) - \gamma \left( \underset{x_i \in S}{\text{MAX}} \{h(x_i)\} \right) \right| < \epsilon$$

$$S \subseteq \mathbb{R}^d$$

**PointNet (vanilla)**

# Invariances

*The model has to respect key desiderata for point clouds:*

## **Point Permutation Invariance**

Point cloud is a set of **unordered** points

## **Spatial Transformation Invariance**

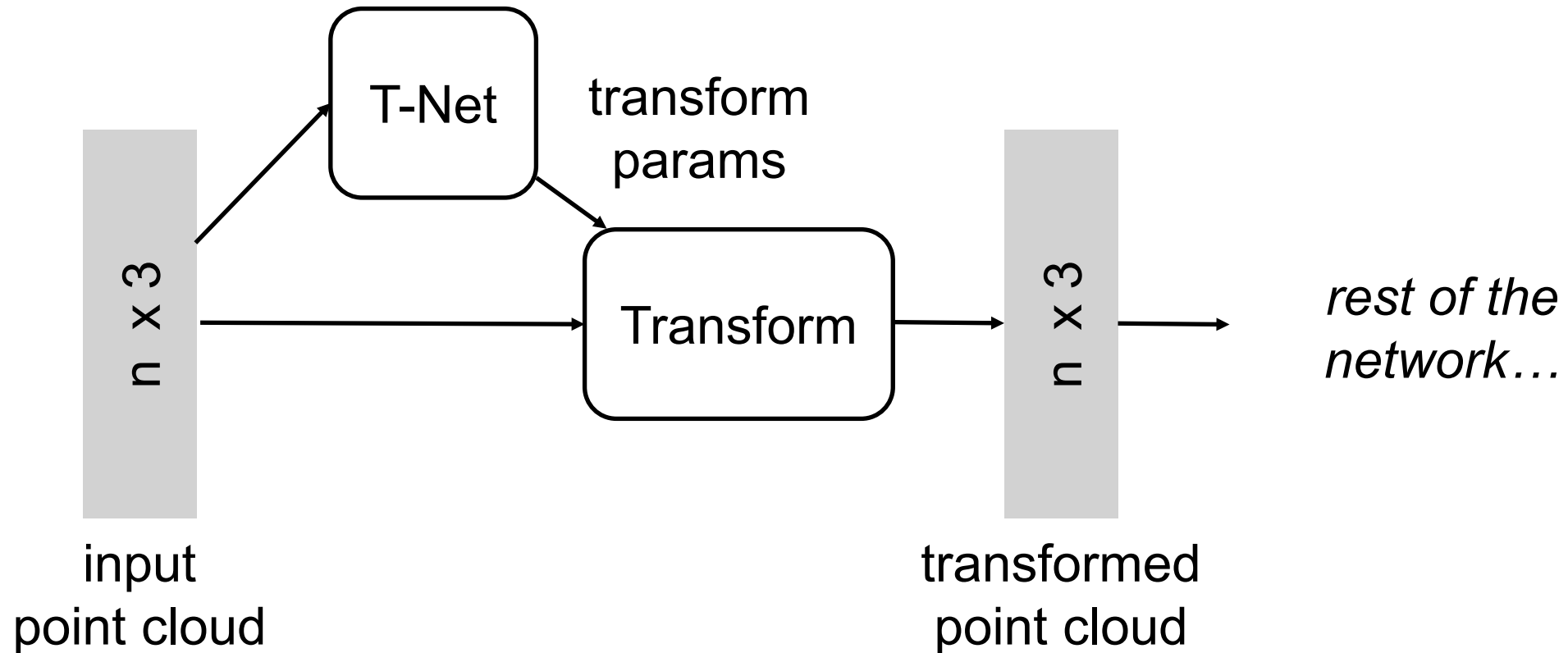
Point cloud **rigid motions** should not alter classification results

## **Sampling Invariance**

Output a function of the underlying geometry and **not the sampling**

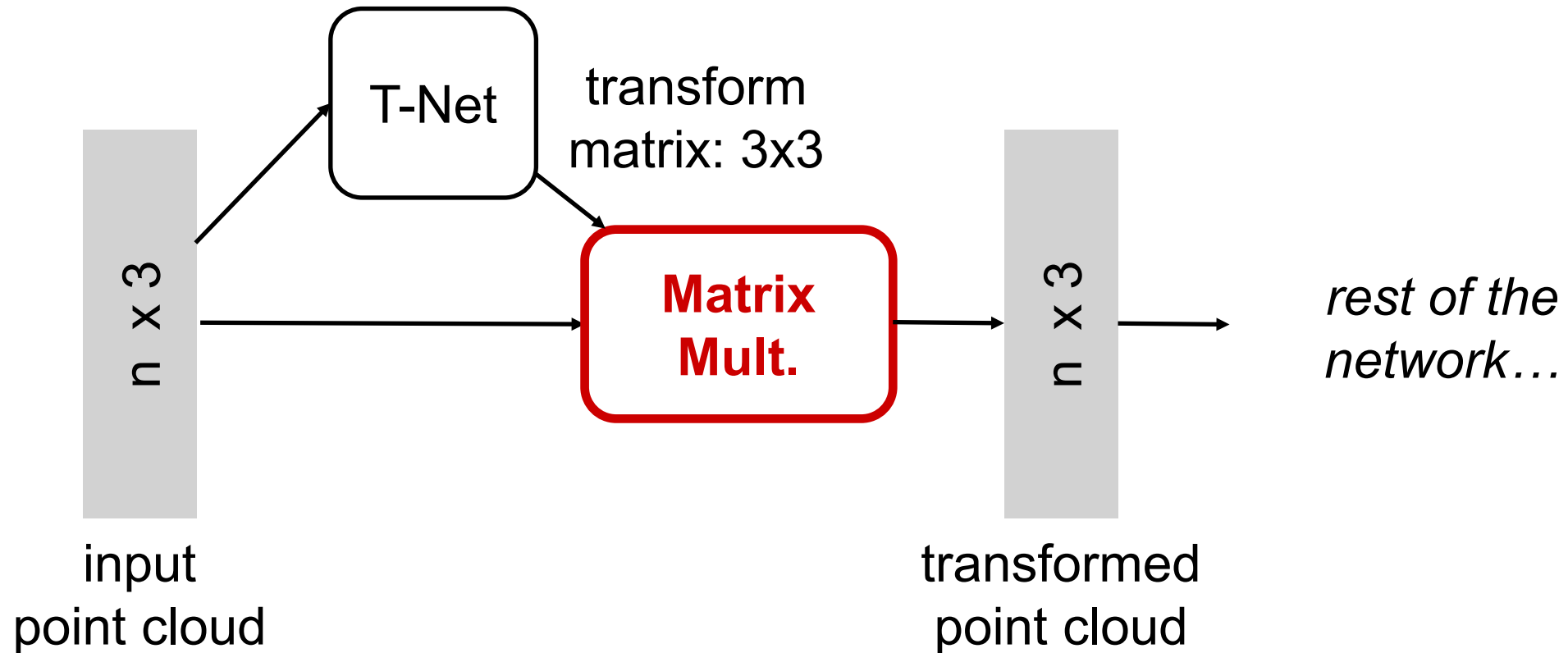
# Input Alignment by Transformer Network

Idea: Data dependent transformation for automatic alignment



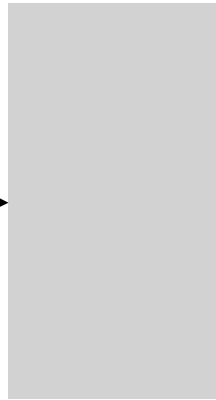
# Input Alignment by Transformer Network

Idea: Data dependent transformation for automatic alignment  
The transformation is just matrix multiplication!



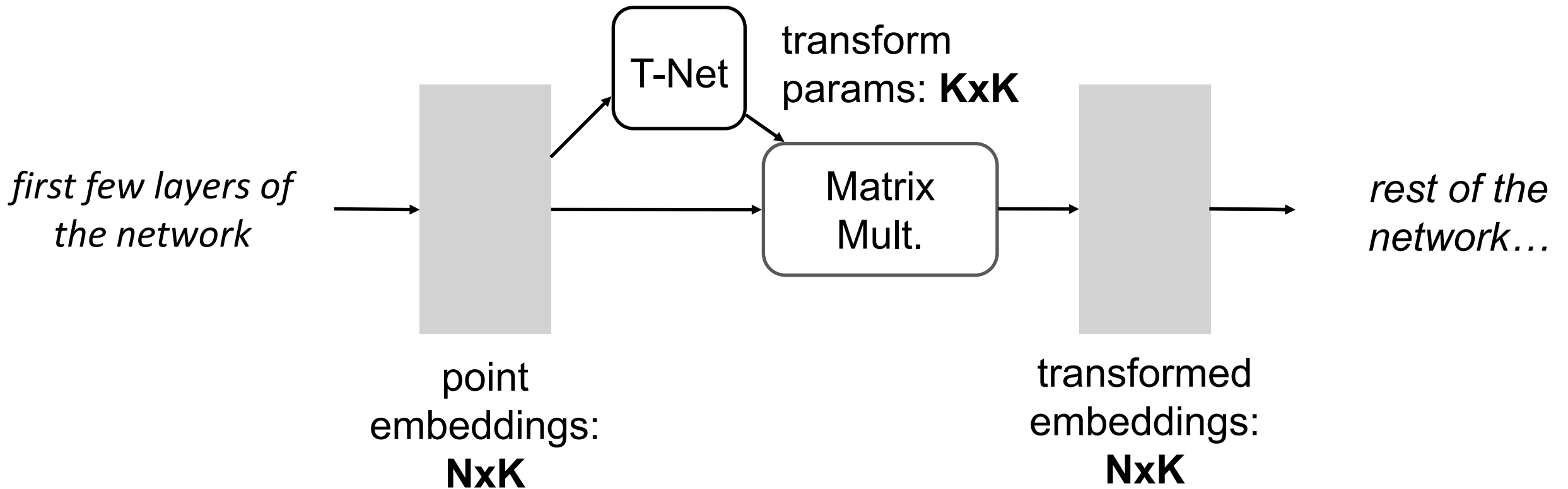
# Embedding Space Alignment

*first few layers of  
the network*

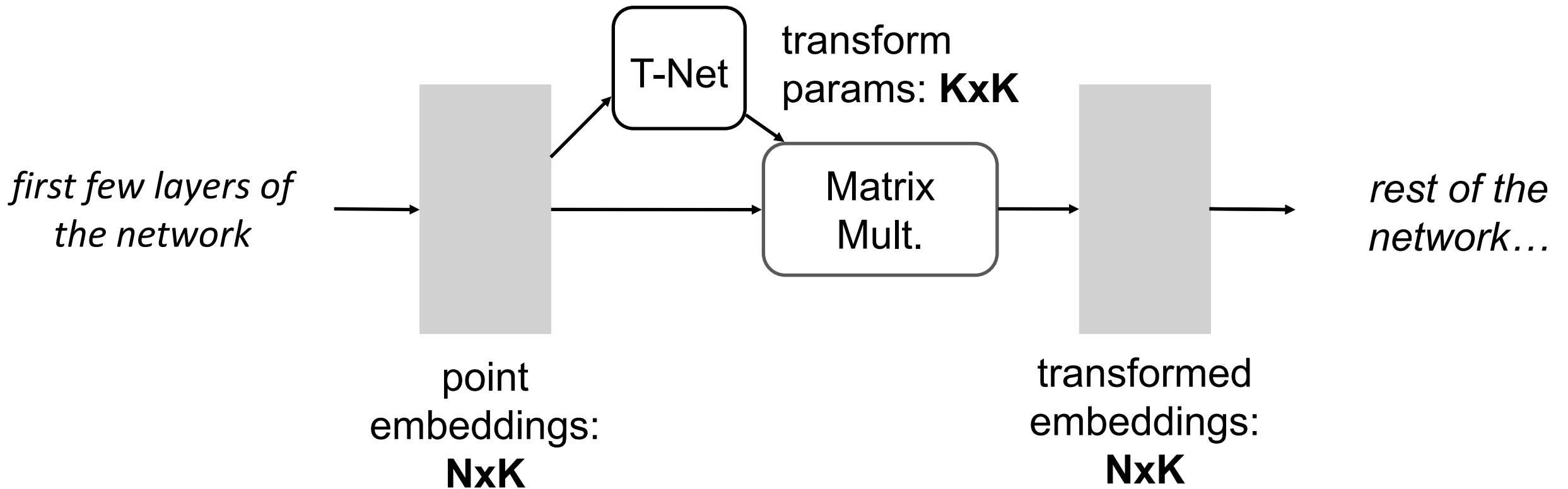


point  
embeddings:  
 **$N \times K$**

# Embedding Space Alignment



# Embedding Space Alignment



**Regularization loss:**

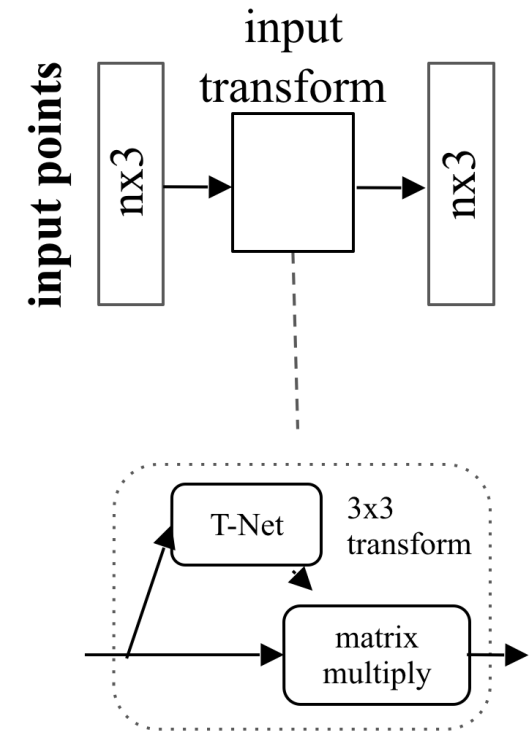
Transform matrix close to orthogonal:  $L_{reg} = \|I - AA^T\|_F^2$

# PointNet Classification Network

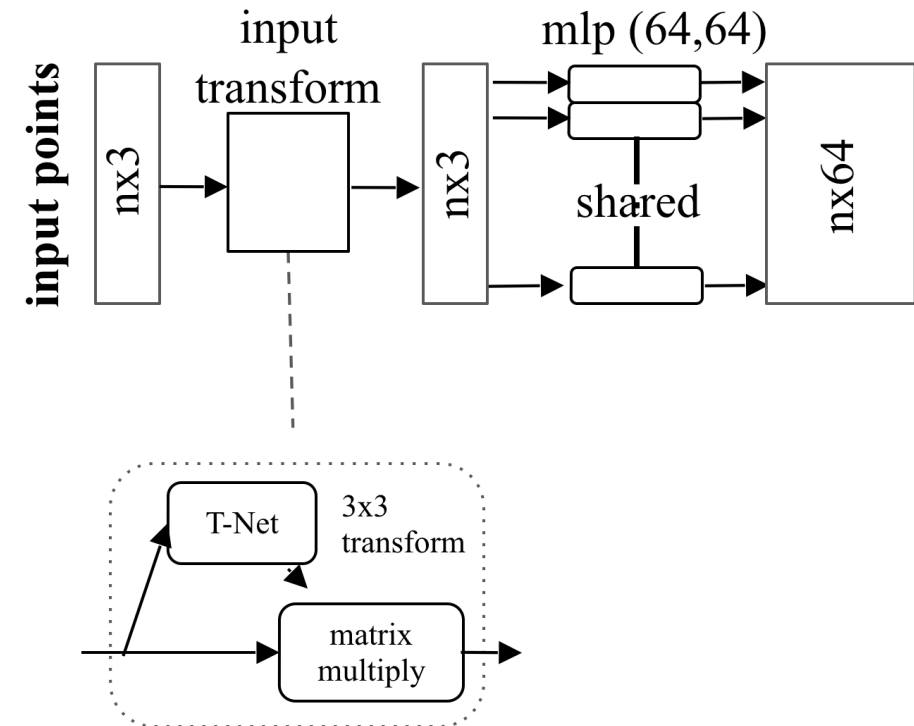
input points

$n \times 3$

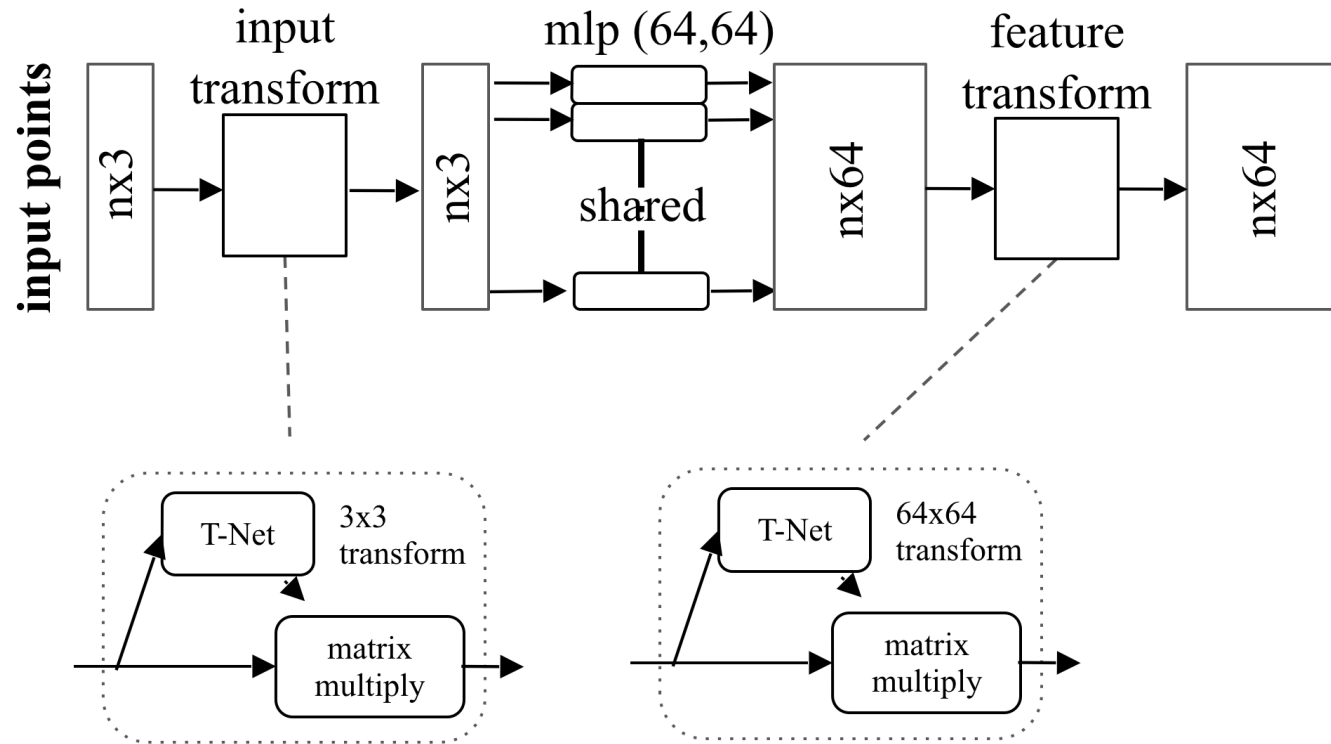
# PointNet Classification Network



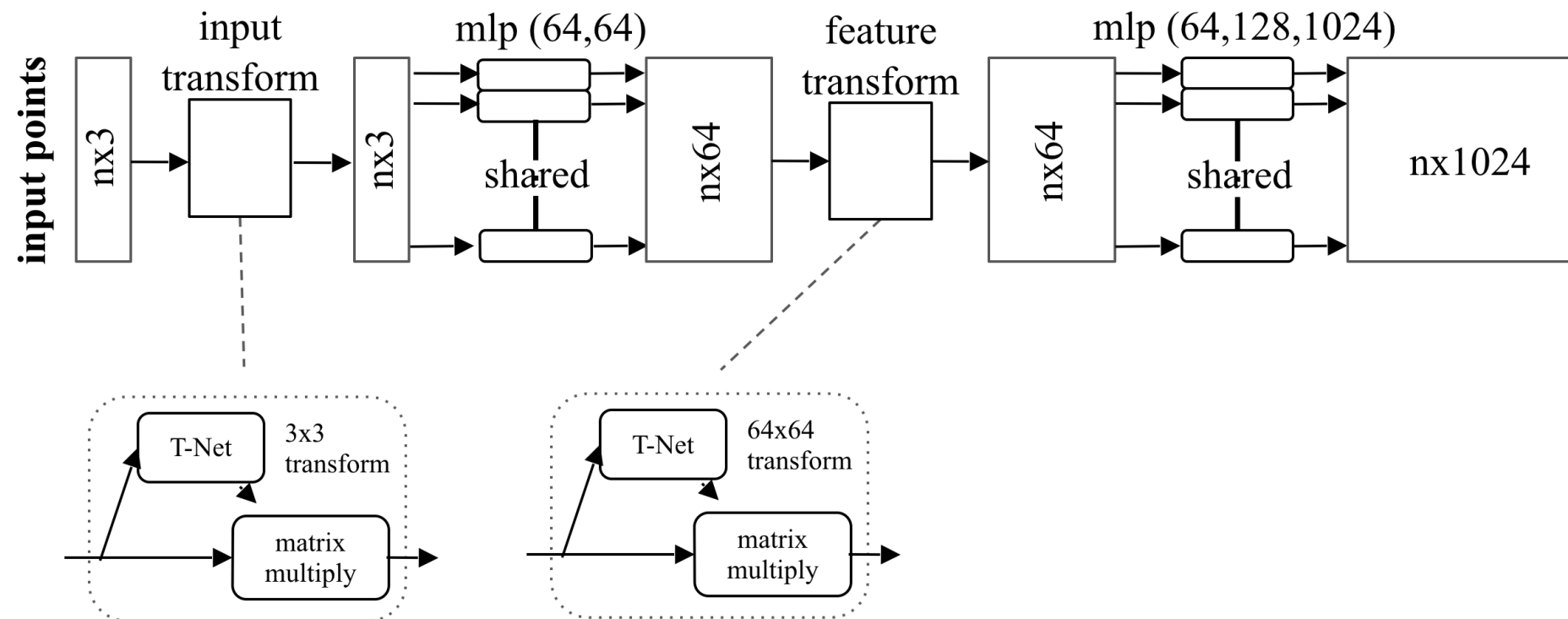
# PointNet Classification Network



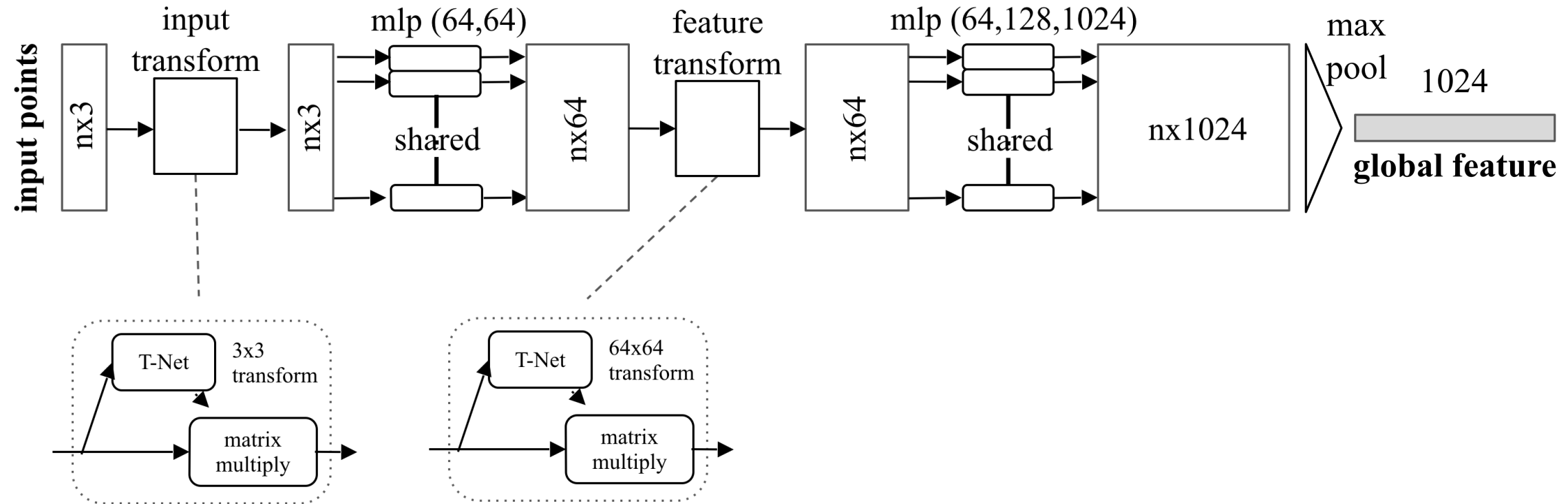
# PointNet Classification Network



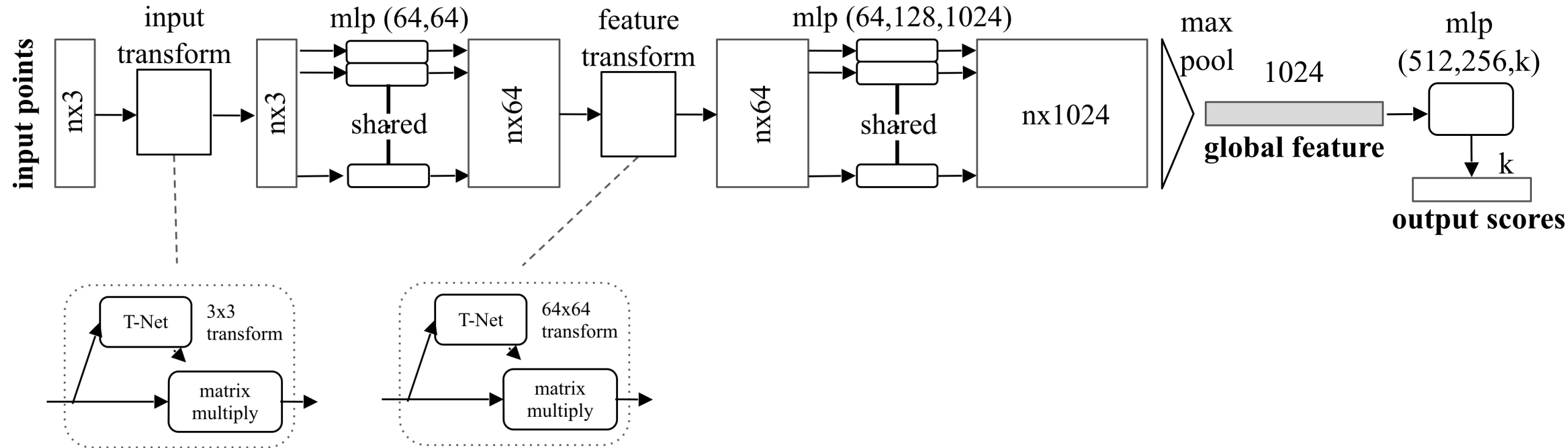
# PointNet Classification Network



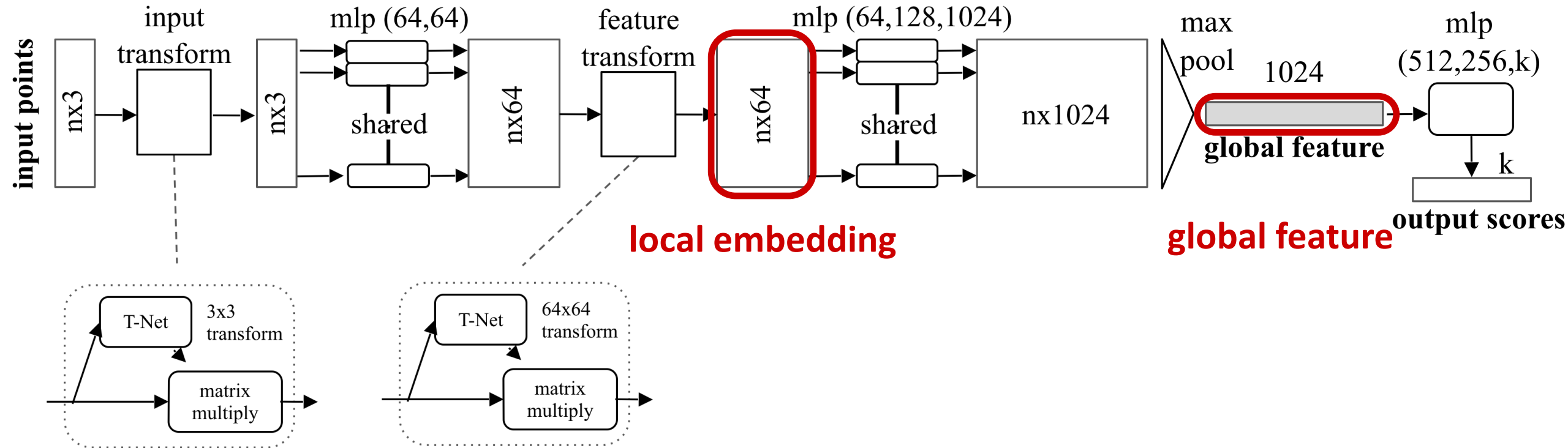
# PointNet Classification Network



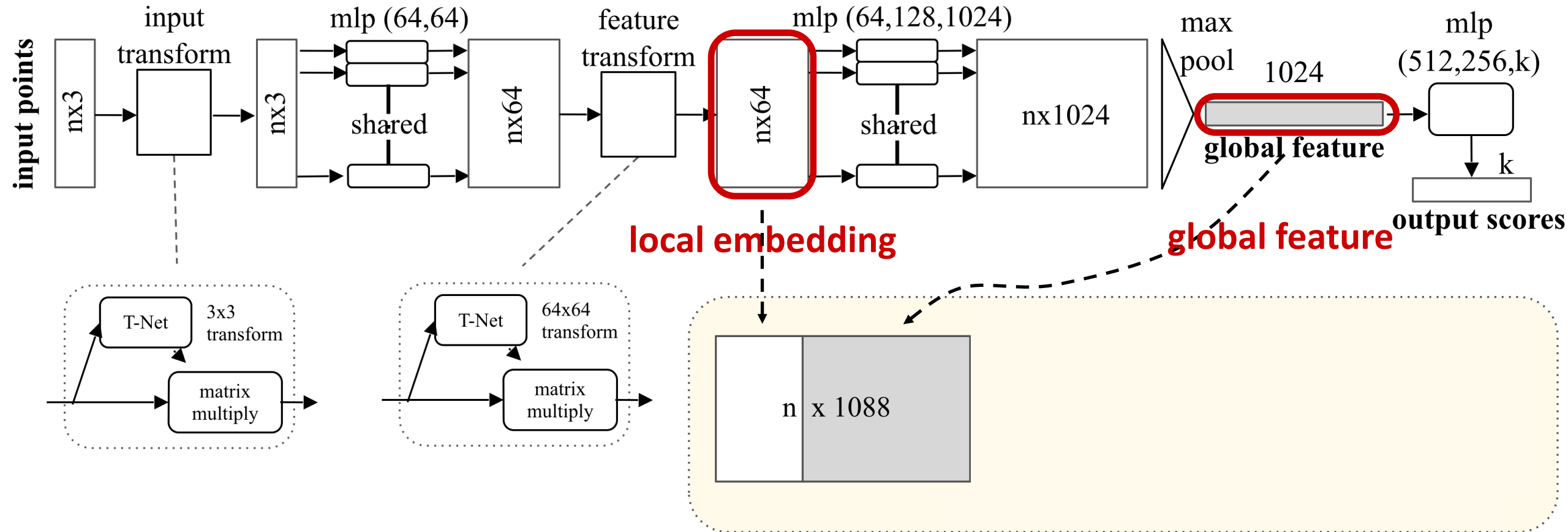
# PointNet Classification Network



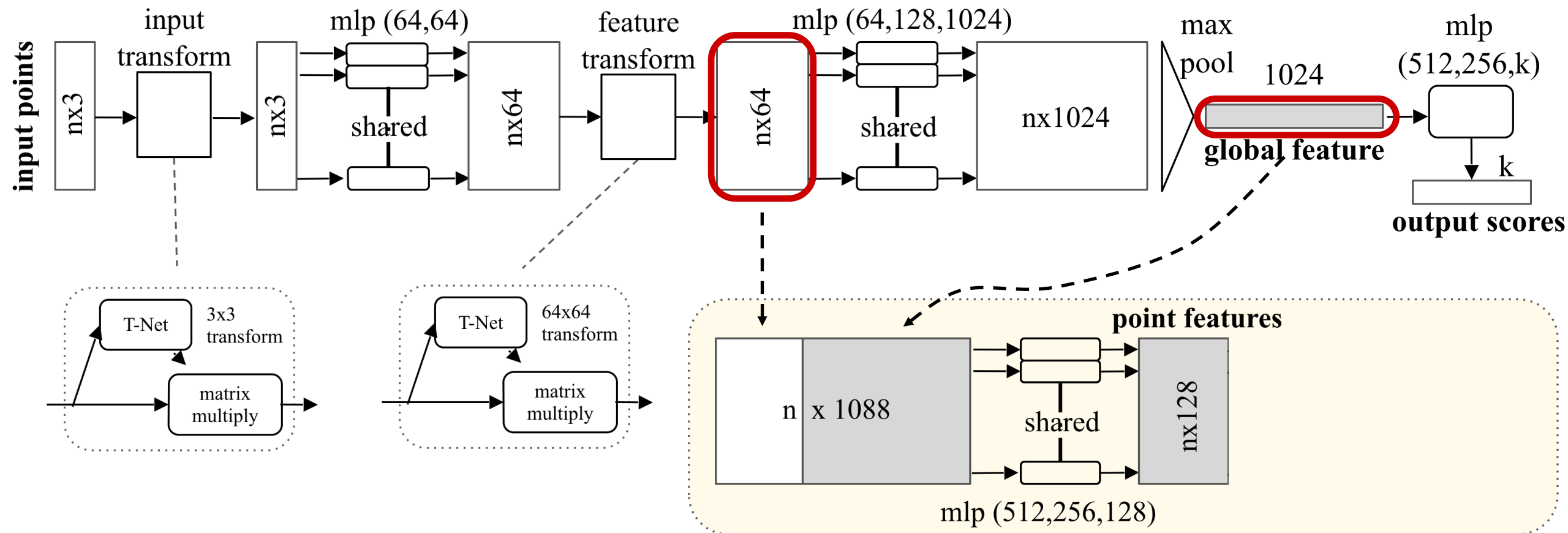
# Extension to PointNet Segmentation Network



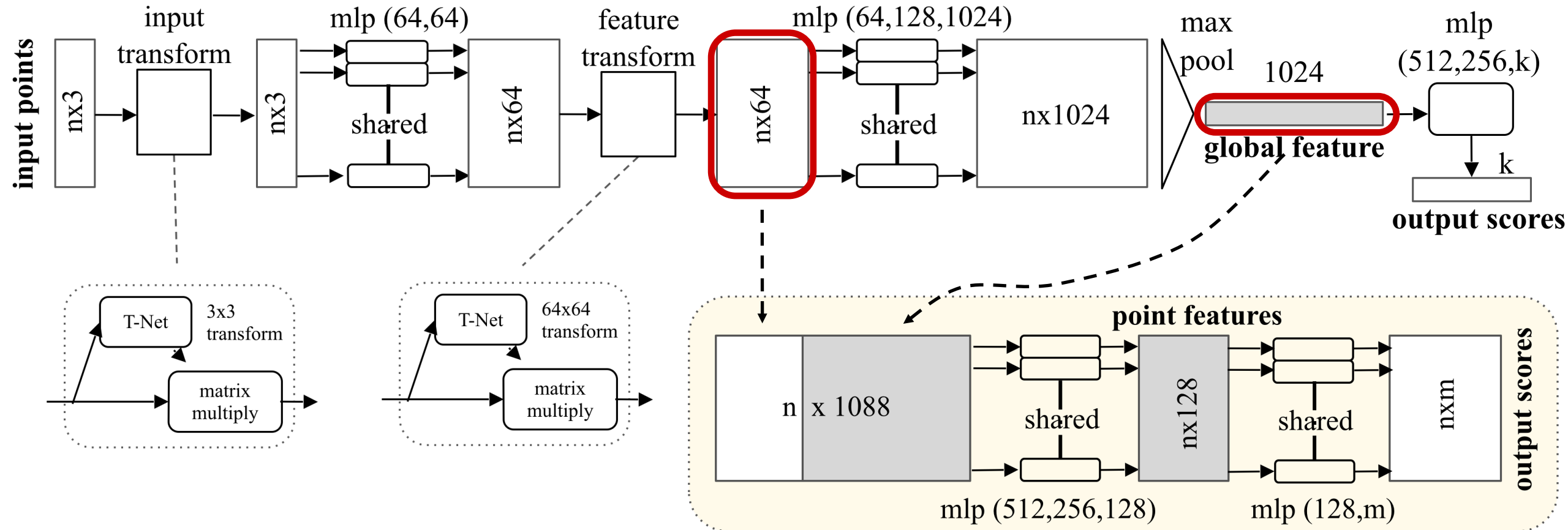
# Extension to PointNet Segmentation Network



# Extension to PointNet Segmentation Network



# Extension to PointNet Segmentation Network



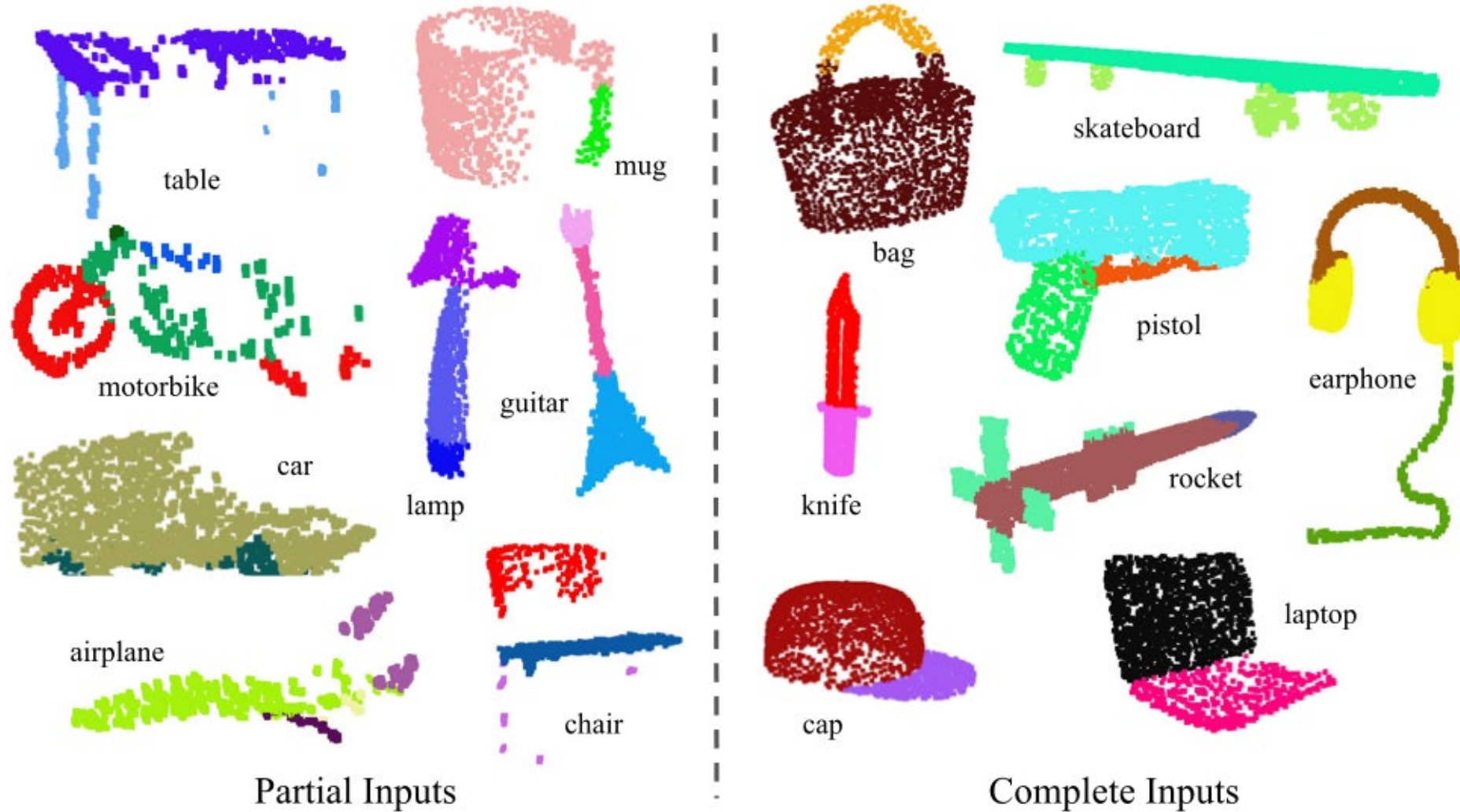
# Results on Object Classification

	input	#views	accuracy avg. class	accuracy overall
	mesh	-	68.2	
	volume	1	77.3	84.7
	volume	12	83.0	85.9
	volume	20	86.0	<b>89.2</b>
	image	10	75.5	-
	image	80	<b>90.1</b>	-
	point	-	72.6	77.4
	point	1	86.2	<b>89.2</b>

3D CNNs

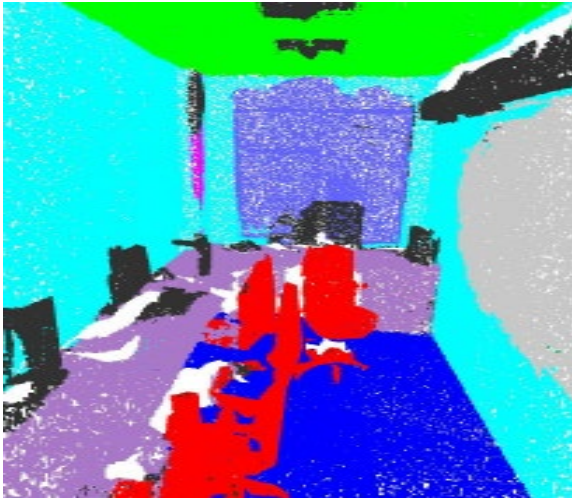
*dataset: ModelNet40; metric: 40-class classification accuracy (%)*

# Results on Object Part Segmentation

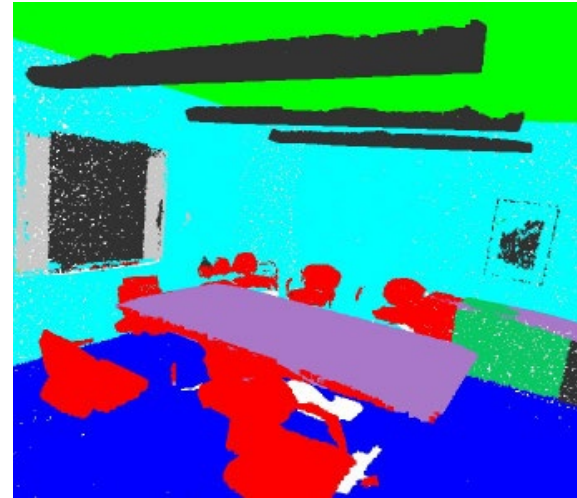


# Results on Semantic Scene Parsing

Input

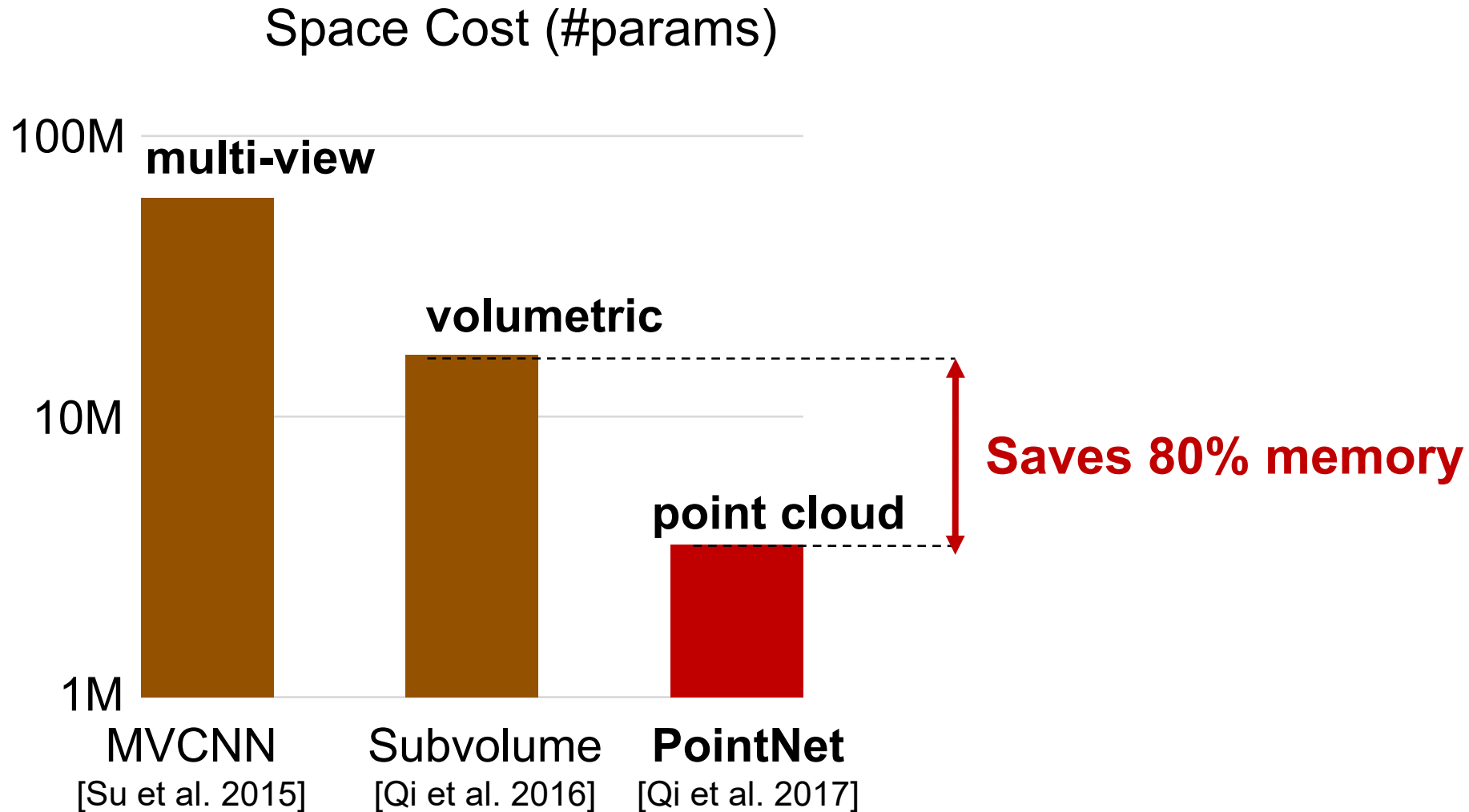


Output

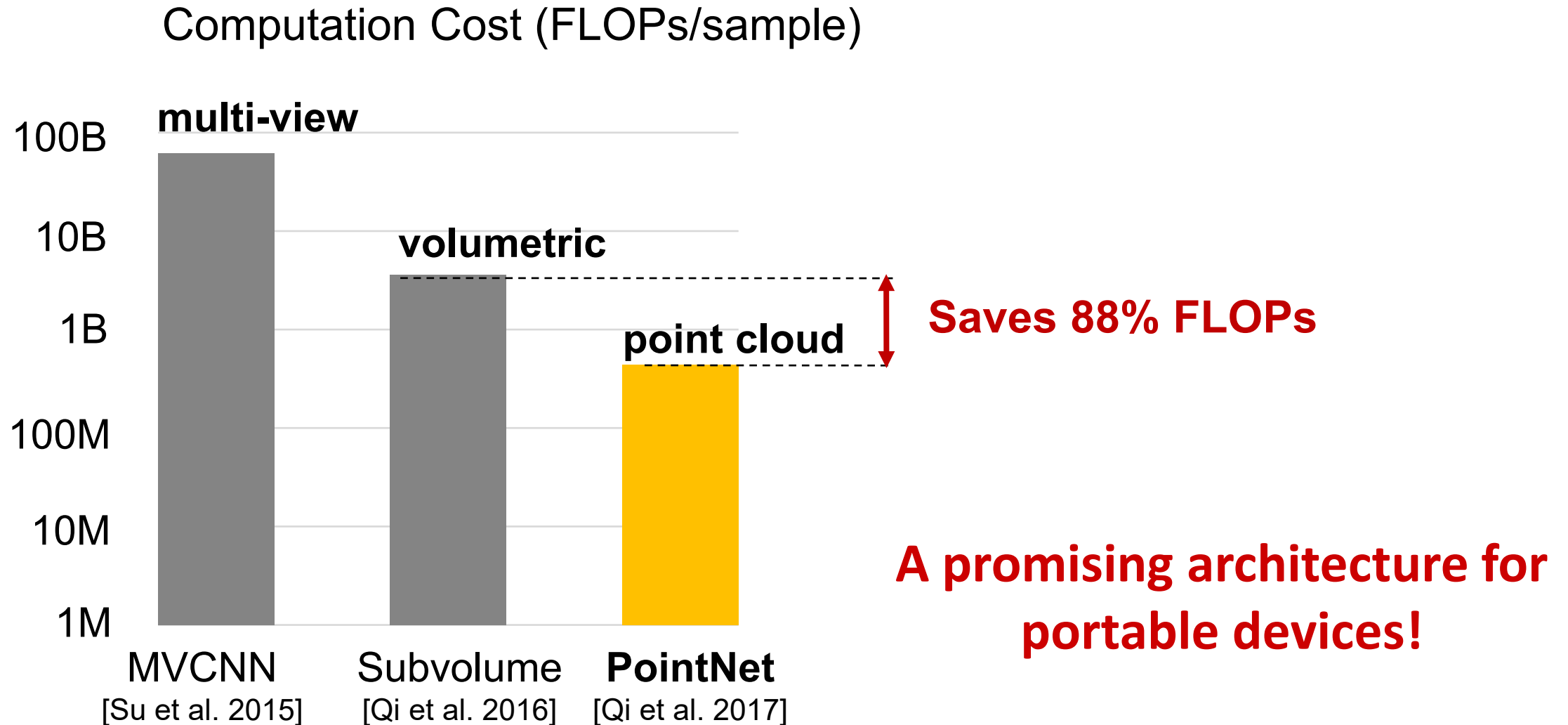


*dataset: Stanford 2D-3D-S (Matterport scans)*

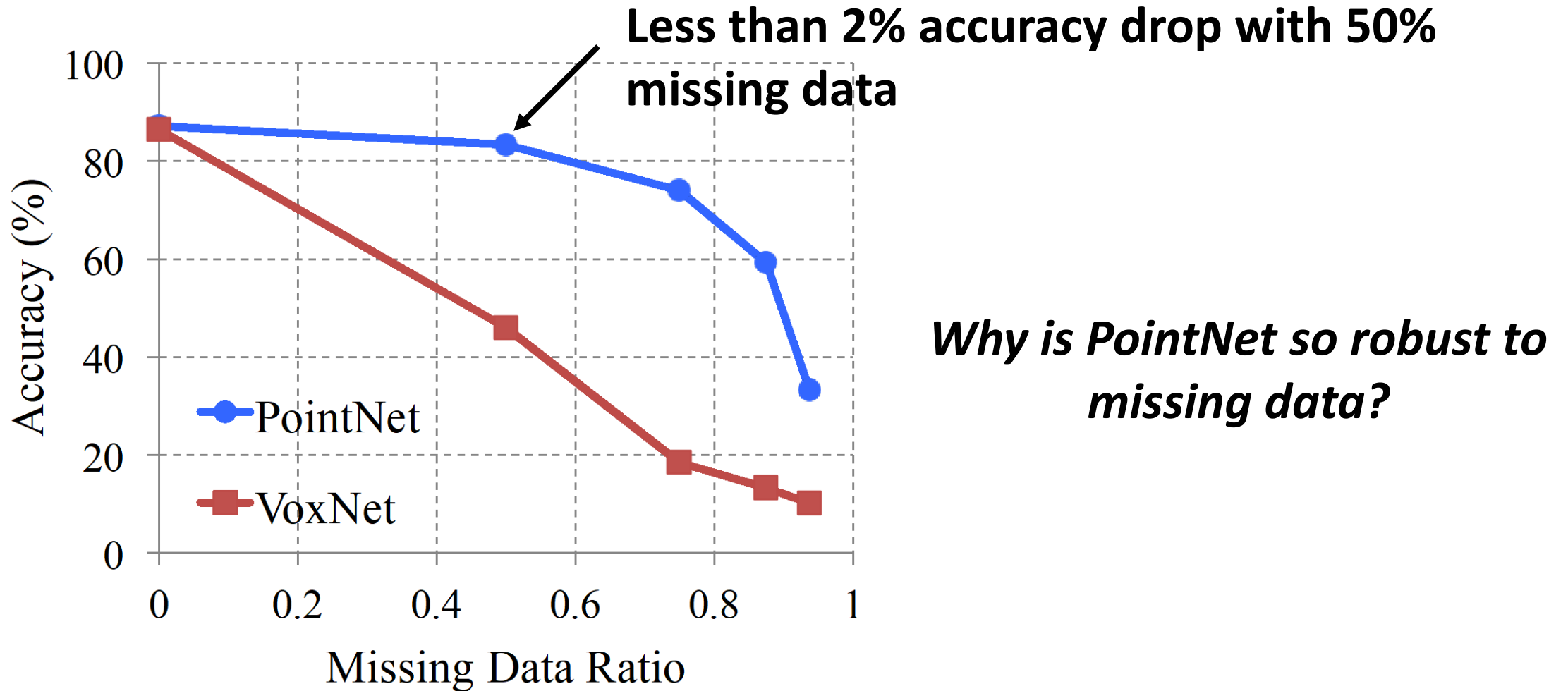
# PointNet is Light-Weight and Fast



# PointNet is Light-Weight and Fast



# PointNet is Robust to Data Corruption

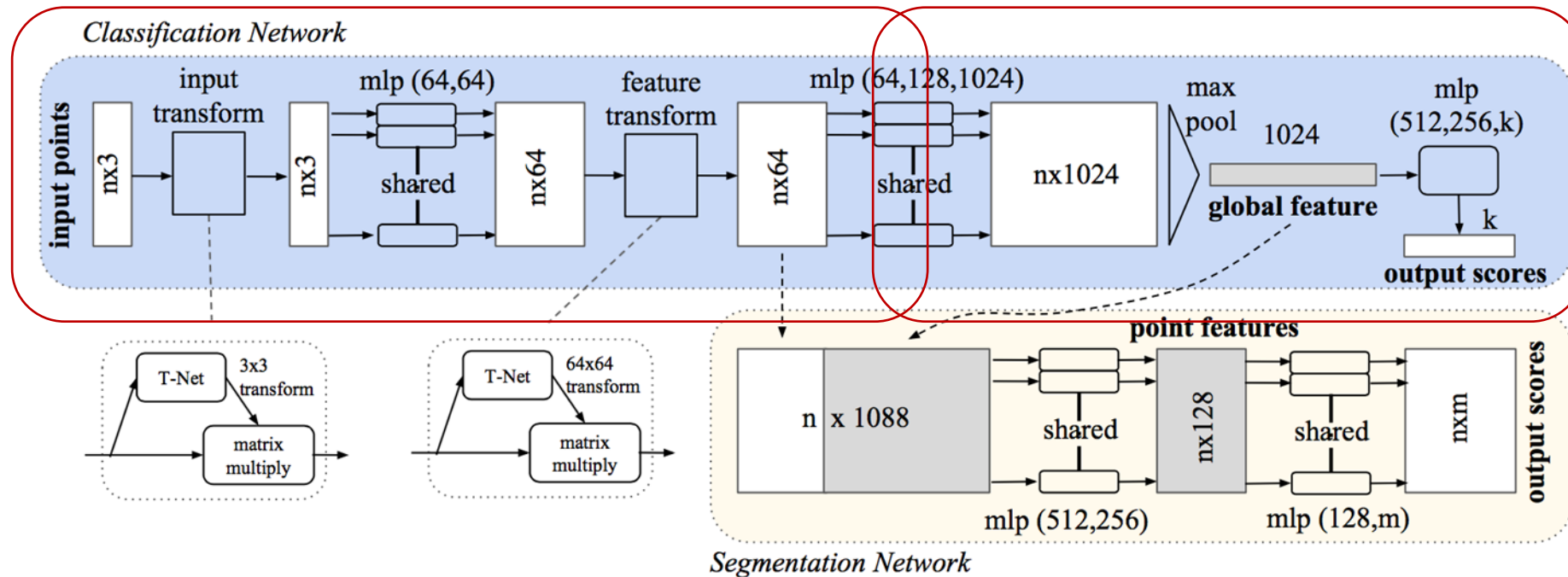


*dataset: ModelNet40; metric: 40-class classification accuracy (%)*

# Visualizing Global Point Cloud Features

Original Shape

# Learning Interesting Points

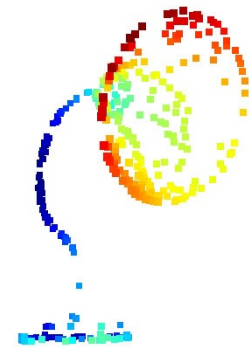
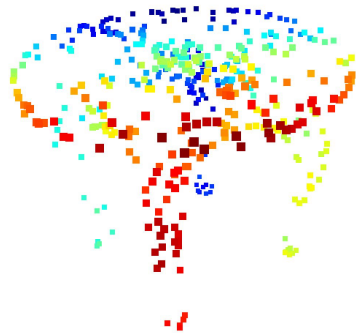


Pointnet learns optimization criteria, which in turn pick interesting points

# Visualizing Global Point Cloud Features

Original Shape

Critical Points

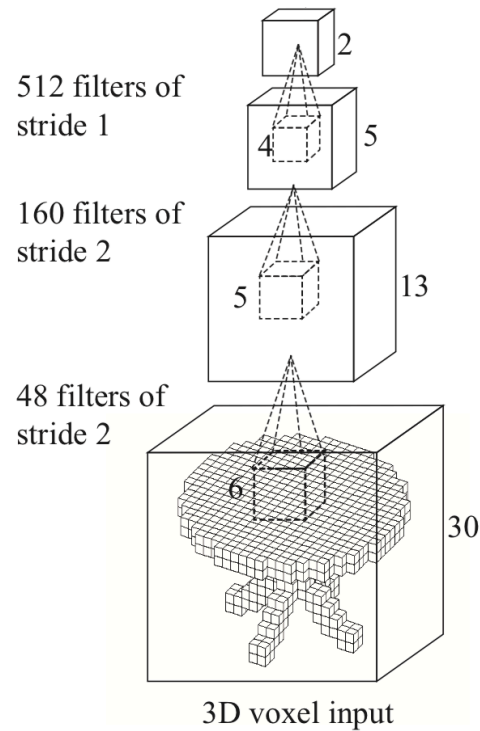


PointNet *learns to pick perceptually interesting points*  
A semantic *core-set* ...

# From PointNet to PointNet++

# Limitations of PointNet

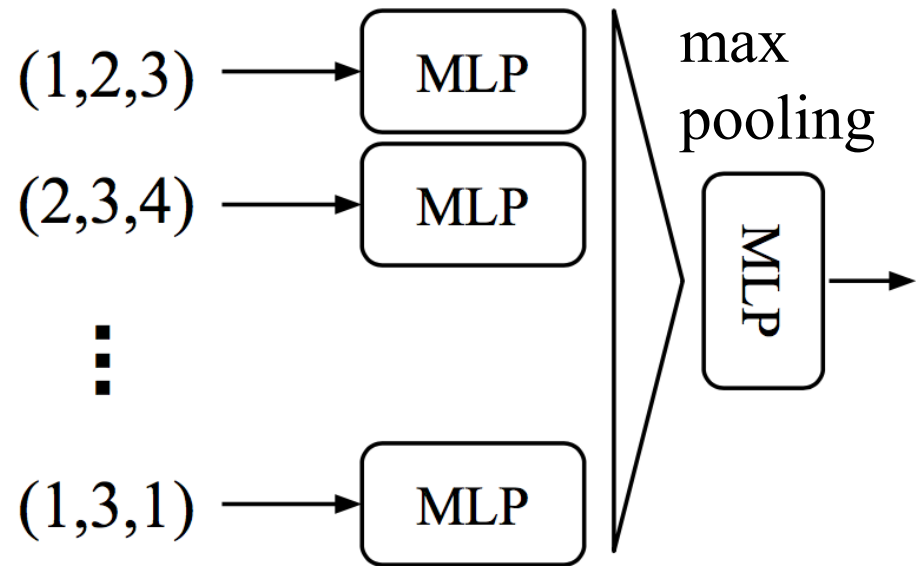
**Hierarchical feature learning**  
multiple levels of abstraction



3D CNN [Wu et al.2015]

**V.S.**

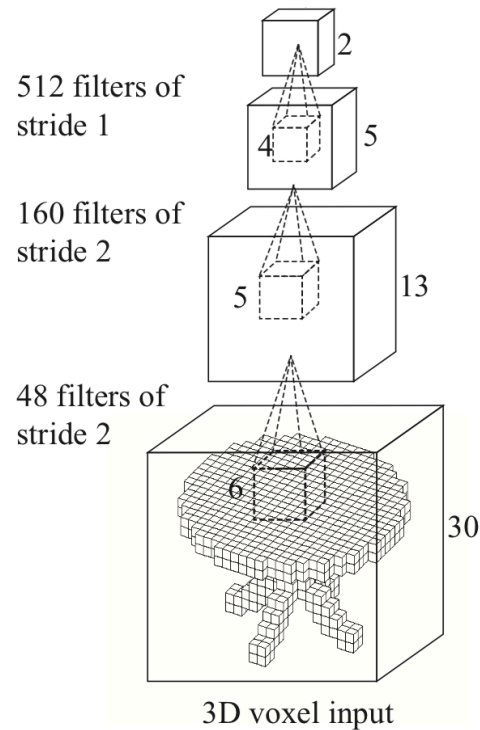
**Global feature learning**  
either **one point**, or **all points**



PointNet (vanilla) [Qi et al.2017]

# Limitations of PointNet

**Hierarchical feature learning**  
multiple levels of abstraction



3D CNN [Wu et al.2015]

**V.S.**

**Global feature learning**  
either **one** point or **all** points

**No local context**

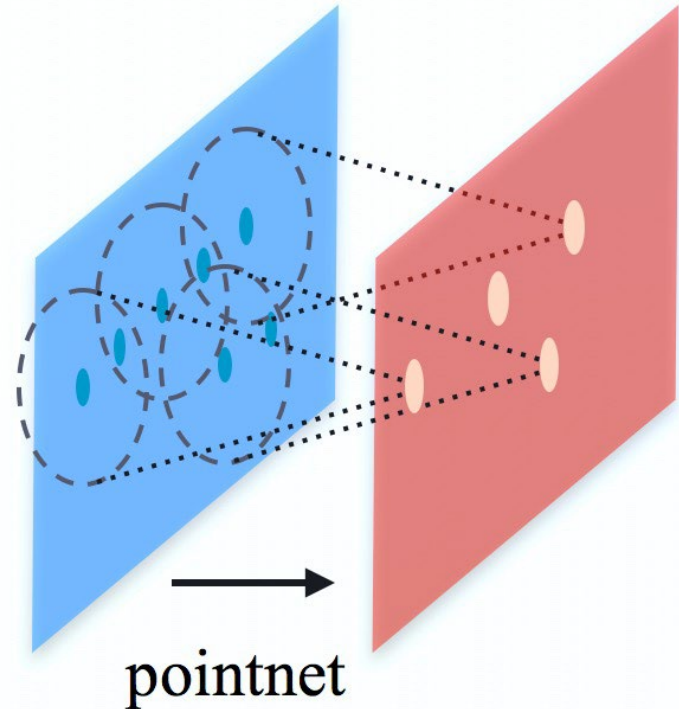
**Limited local invariance**

PointNet (vanilla) [Qi et al.2017]

# PointNet++

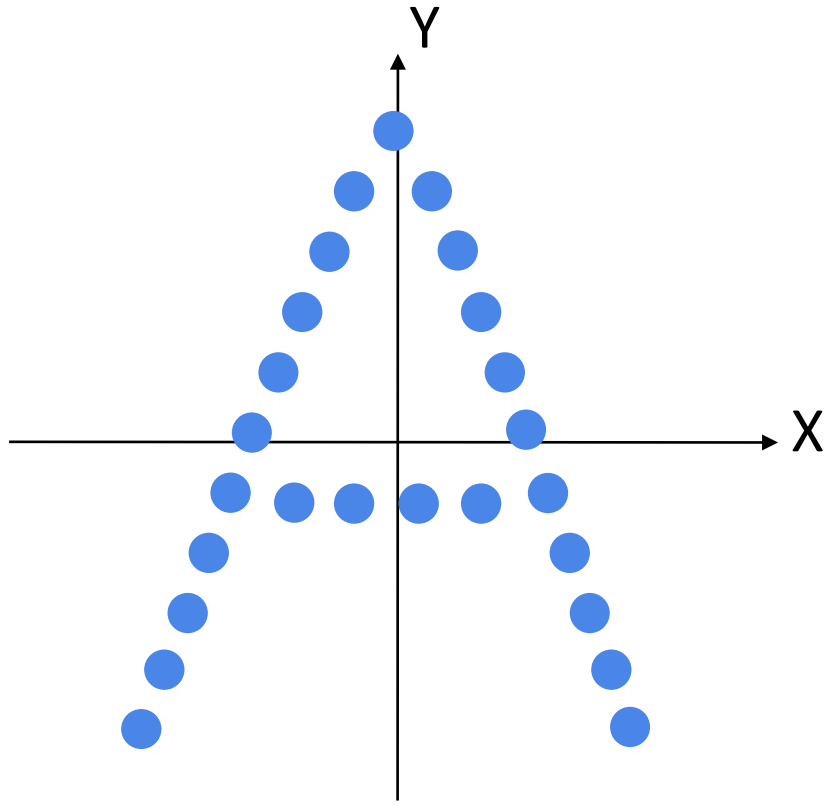
Basic idea: Recursively apply pointnet at local regions.

- ✓ Hierarchical feature learning
- ✓ Local translation invariance
- ✓ Permutation invariance



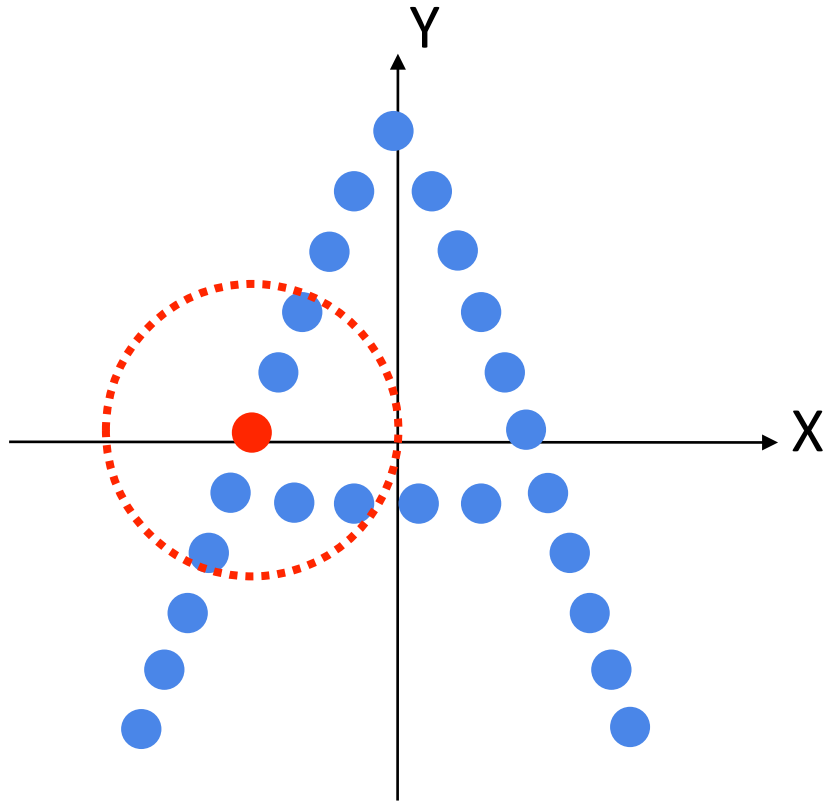
Charles R. Qi, Li Yi, Hao Su, Leonidas Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space (NIPS'17)

# Hierarchical Point Feature Learning



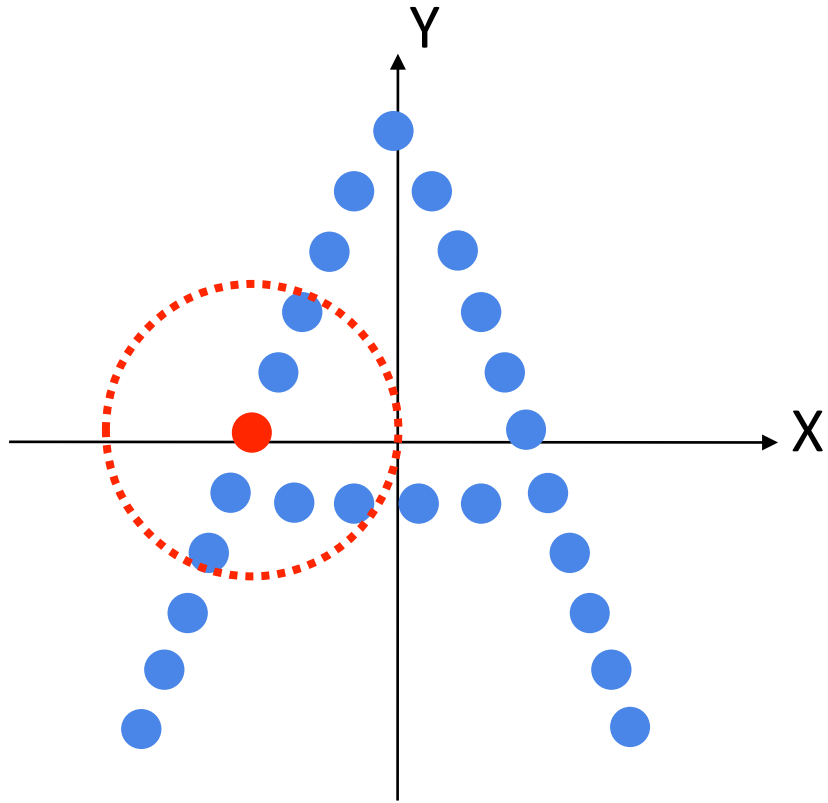
$N$  points in  $(X, Y)$

# Hierarchical Point Feature Learning

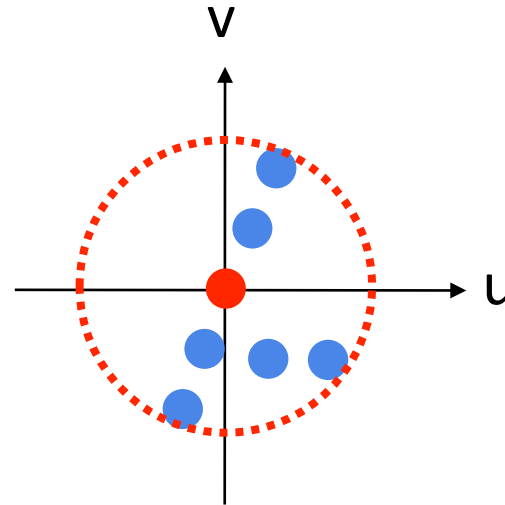


N points in (X,Y)

# Hierarchical Point Feature Learning

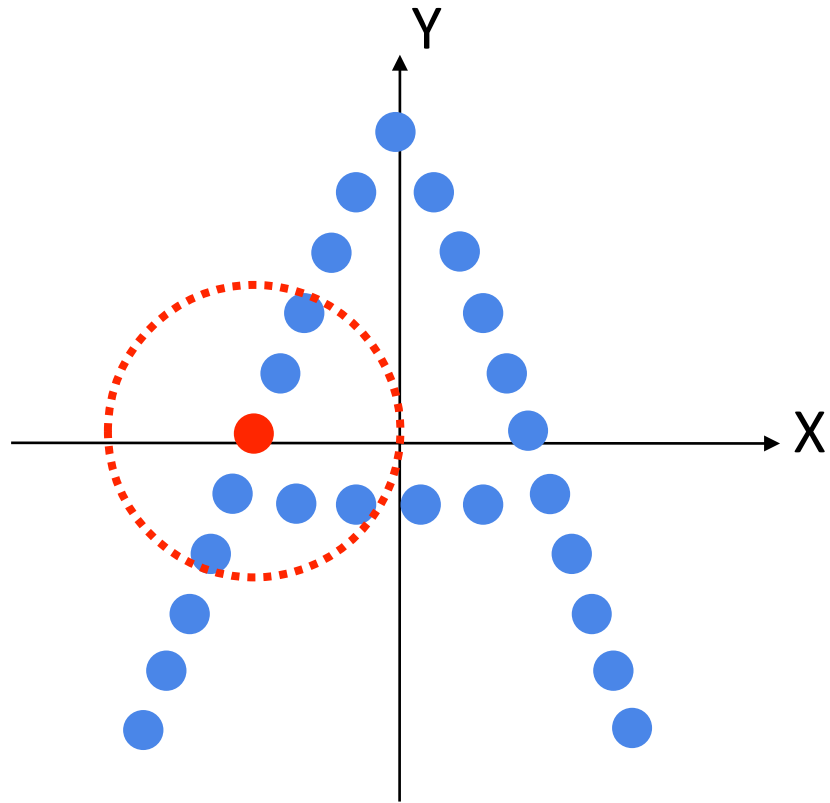


N points in  $(X, Y)$



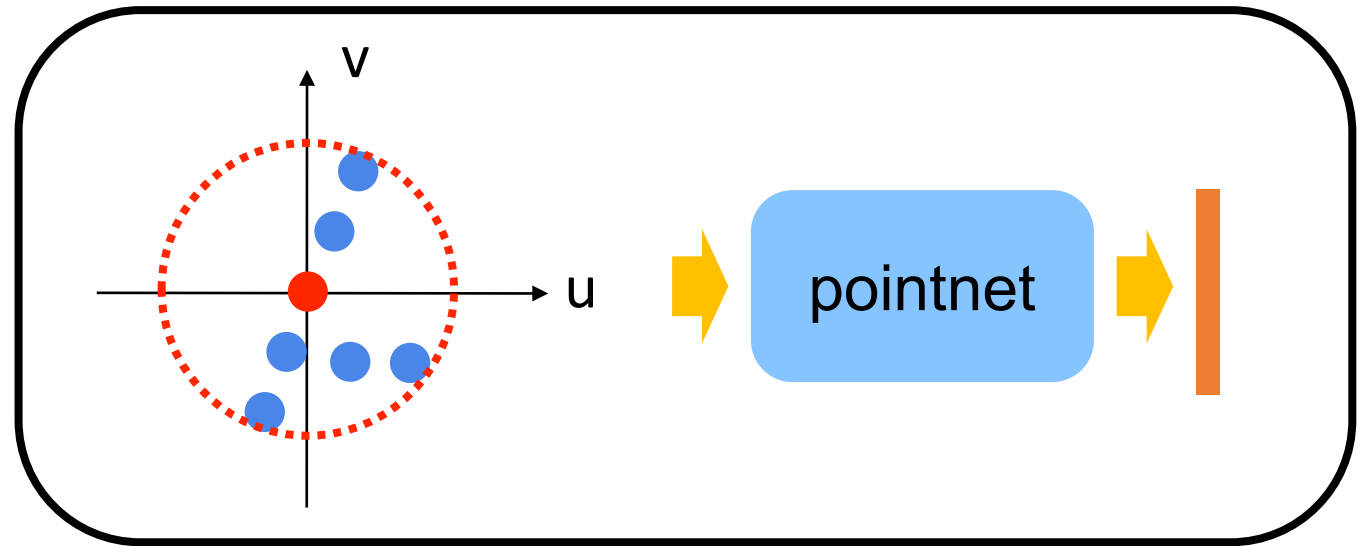
k points in local coordinates  $(u, v)$

# Hierarchical Point Feature Learning



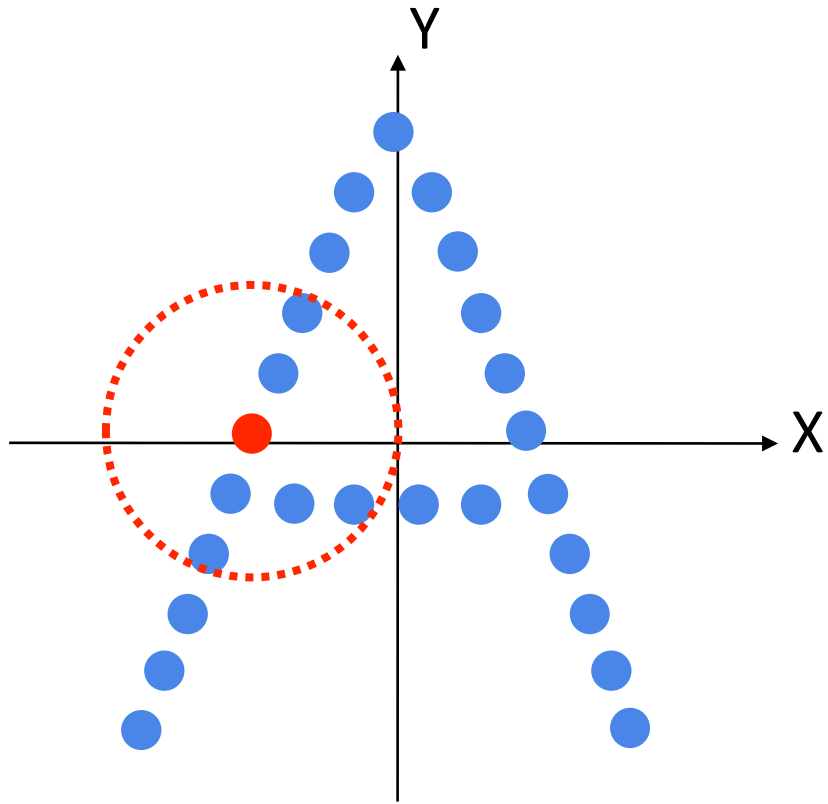
N points in (X,Y)

Apply pointnet at a local region

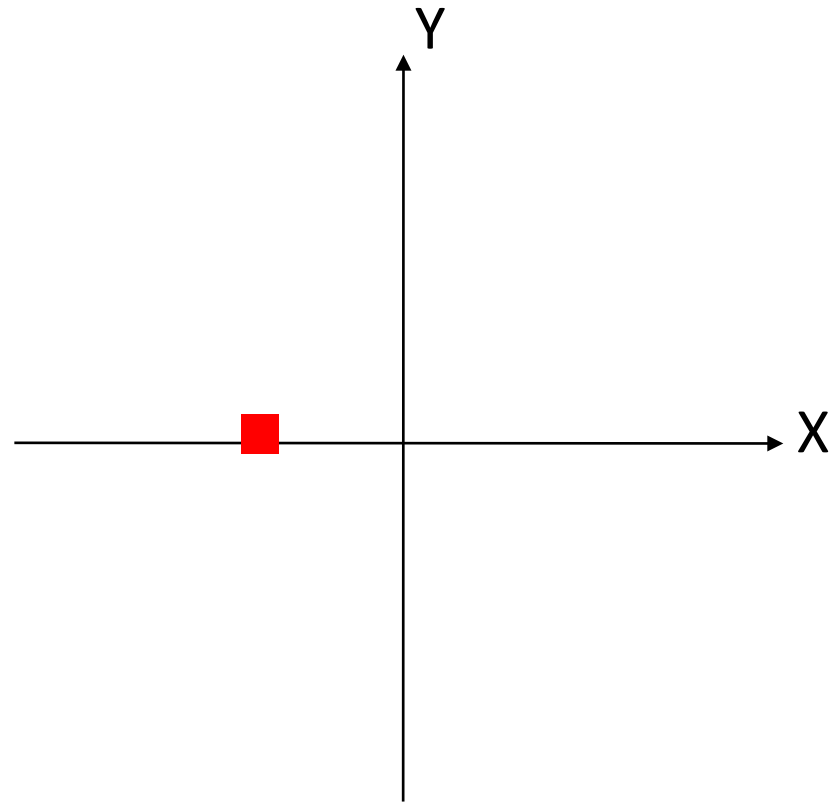


k points in local coordinates (u,v)

# Hierarchical Point Feature Learning



N points in  $(X,Y)$

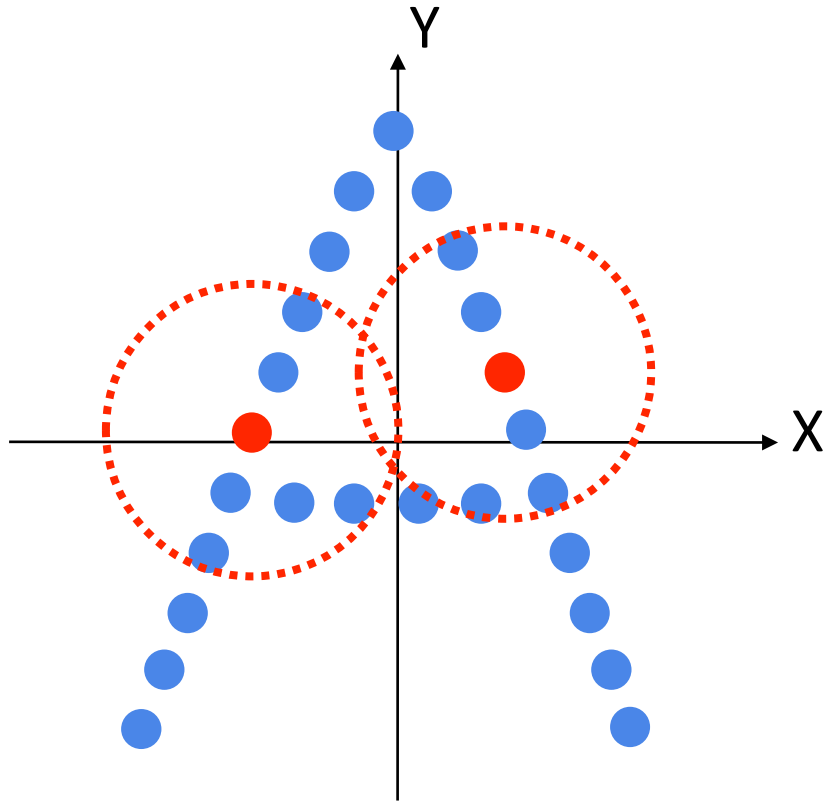


points in  $(X,Y, \mathbf{F})$

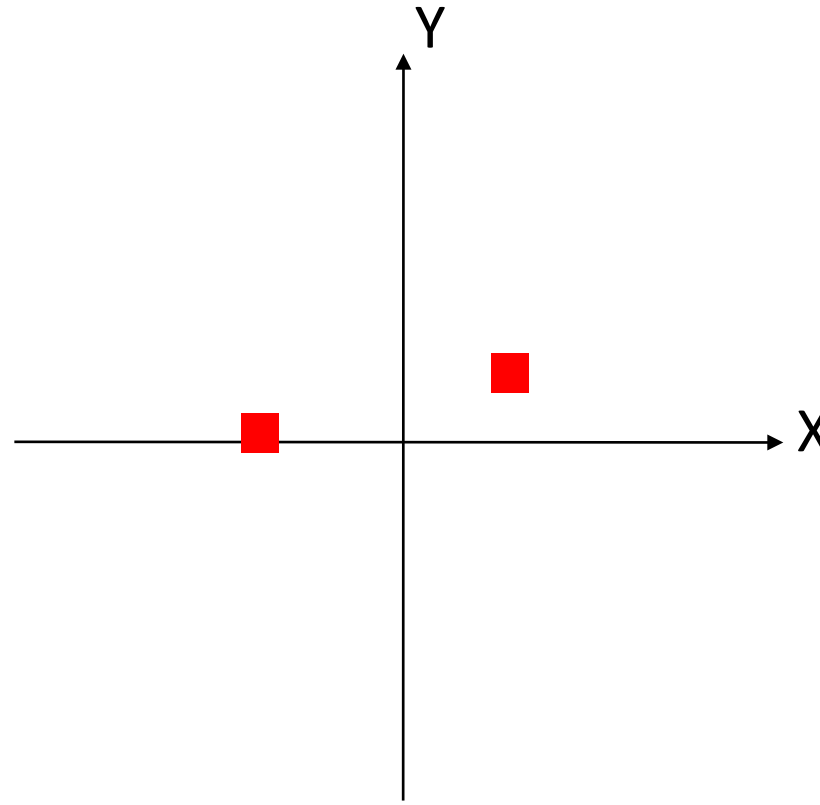
Euclidean space

**high-dim feature space**

# Hierarchical Point Feature Learning

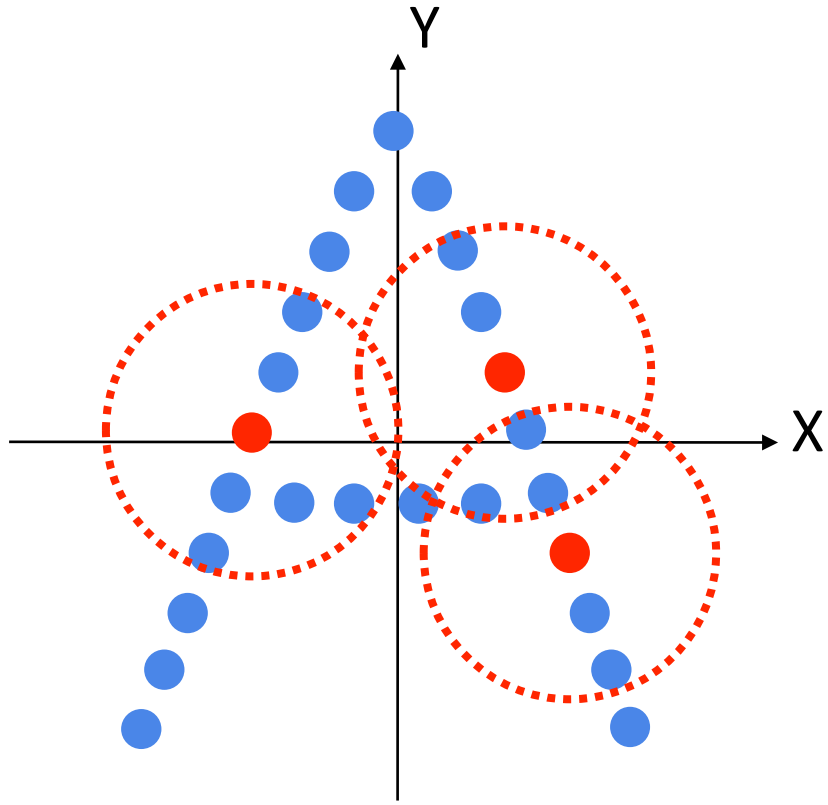


N points in  $(X, Y)$

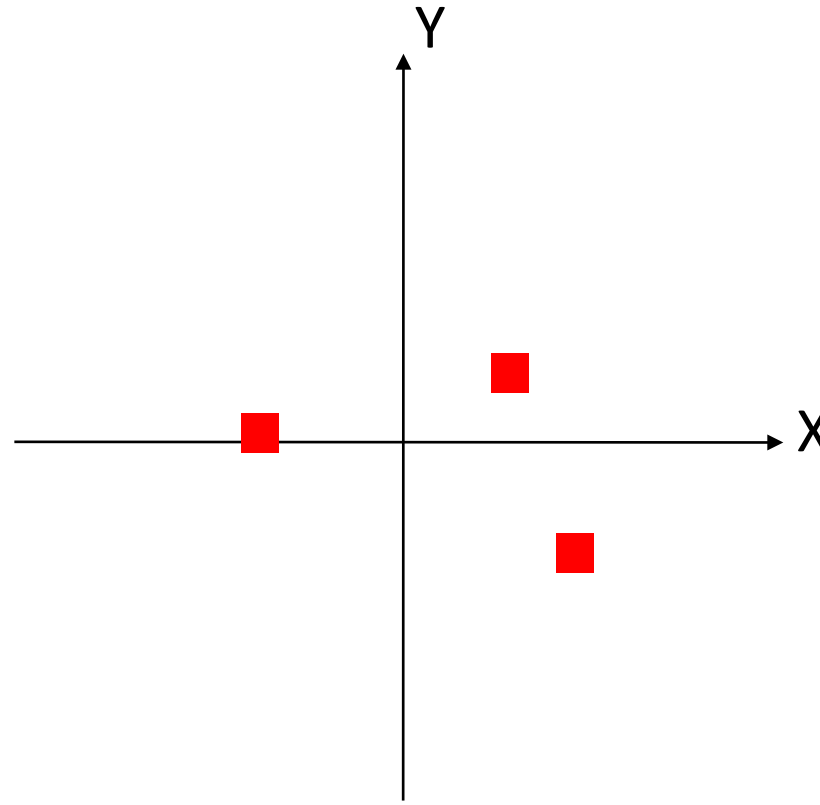


points in  $(X, Y, \mathbf{F})$

# Hierarchical Point Feature Learning

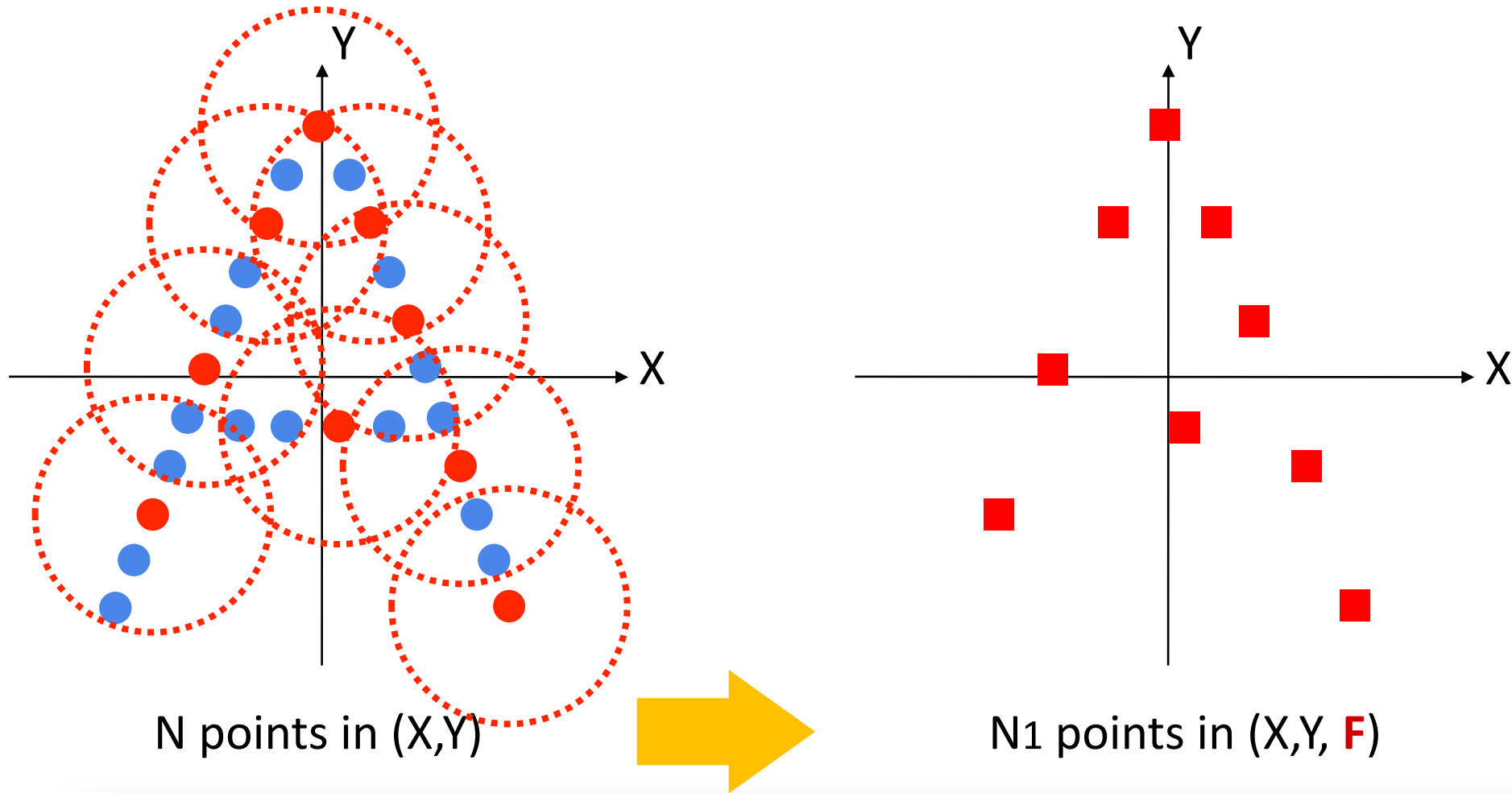


N points in (X,Y)



points in (X,Y, **F**)

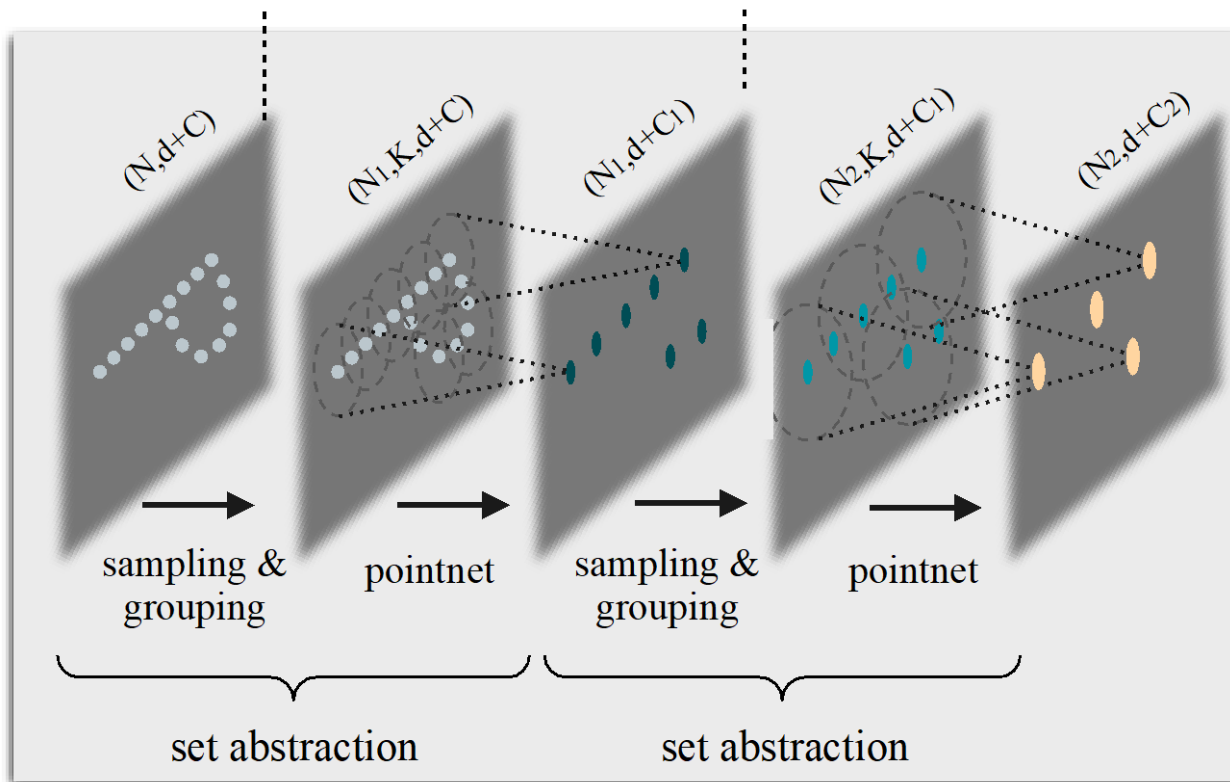
# Hierarchical Point Feature Learning



**Set Abstraction:** farthest point sampling + grouping + pointnet

# PointNet++ for Classification and Segmentation

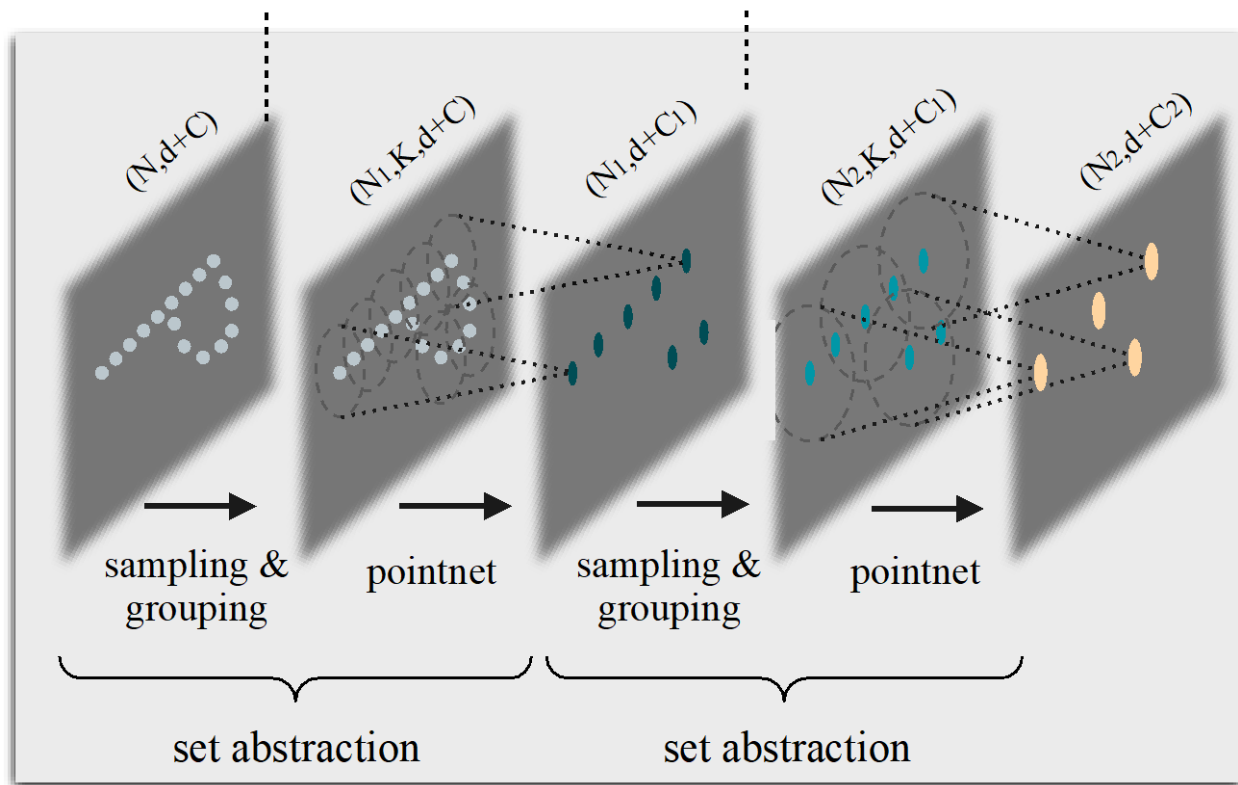
## *Hierarchical point set feature learning*



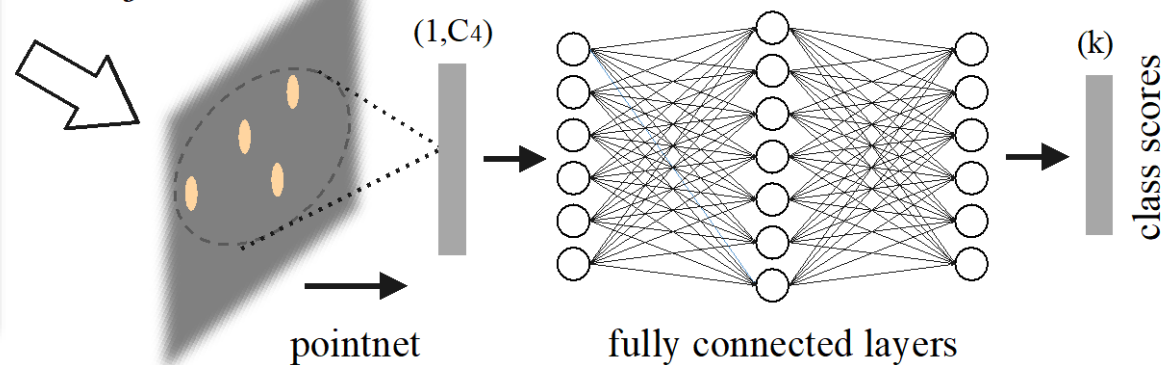
Caveat: Shouldn't feature dimensions from the lower layers affect connectivity at the higher layers?

# PointNet++ for Classification and Segmentation

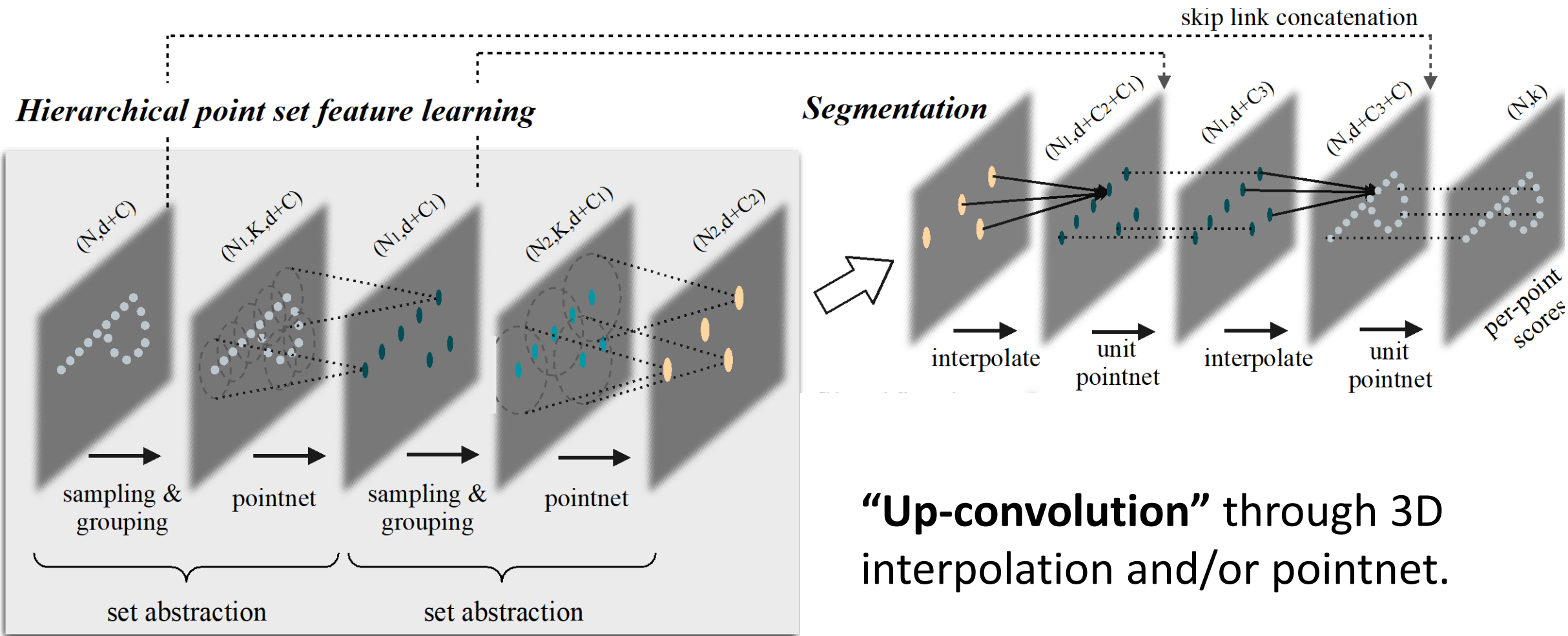
## *Hierarchical point set feature learning*



## *Classification*



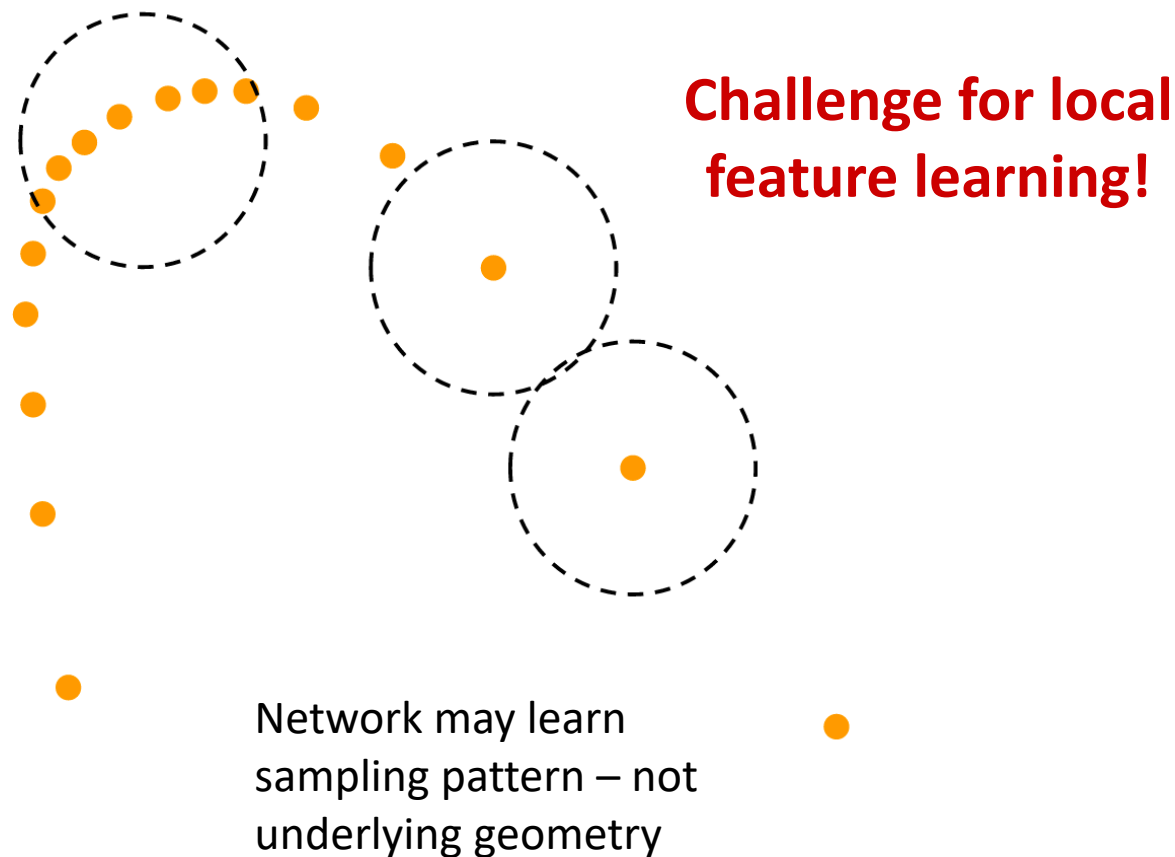
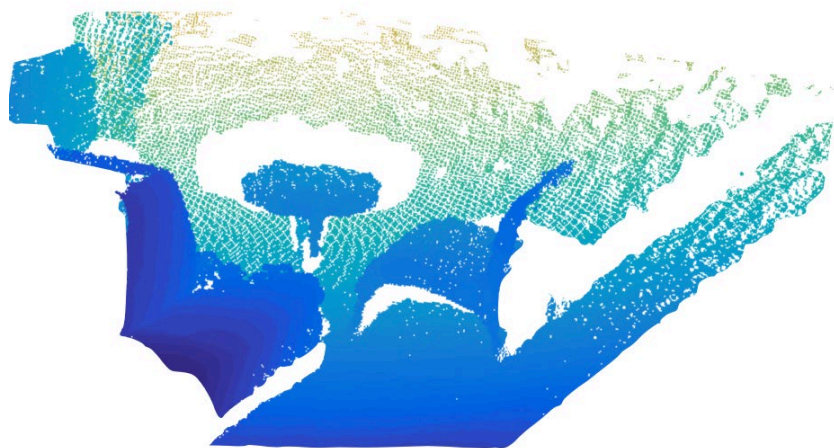
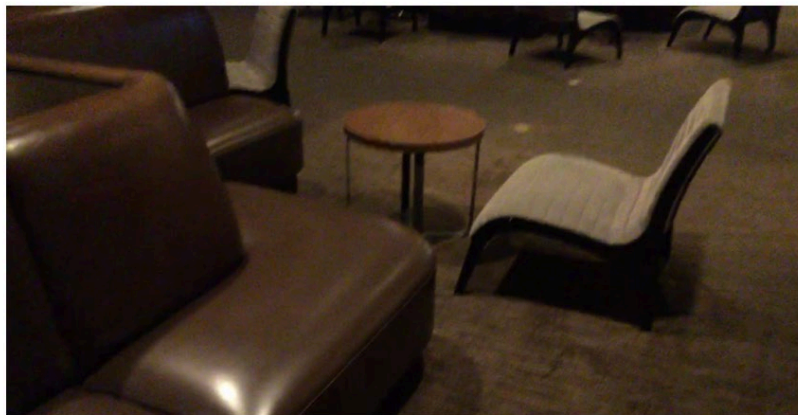
# PointNet++ for Classification and Segmentation



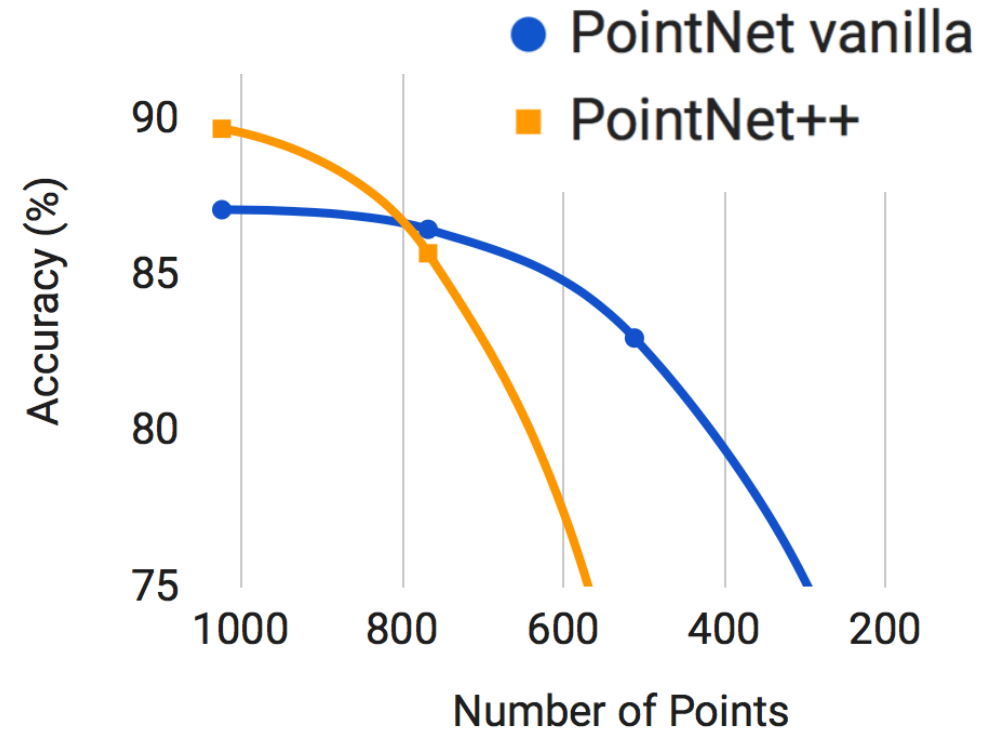
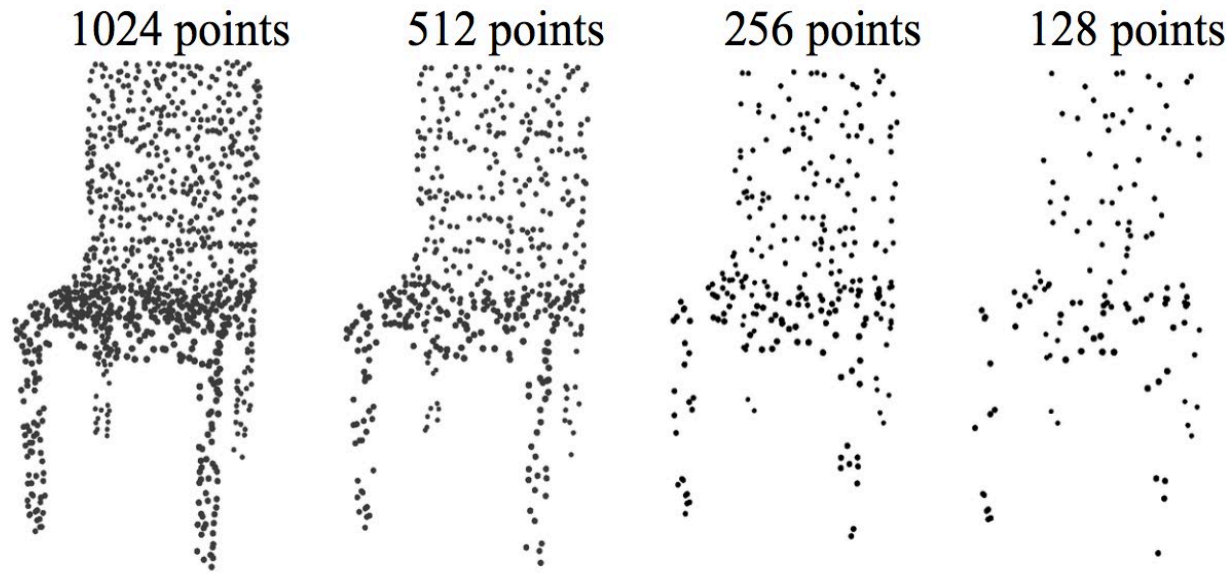
**“Up-convolution”** through 3D interpolation and/or pointnet.

# Non-uniform Sampling Density in Point Clouds

Density variation is a common issue in 3D point cloud processing  
- perspective effect, radial density variation, motion etc.

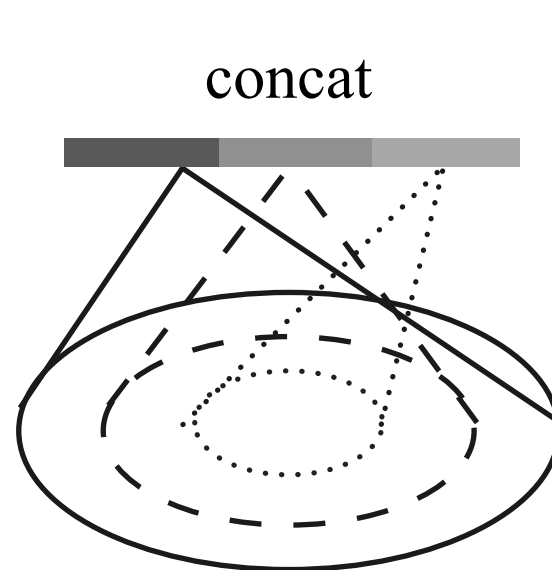
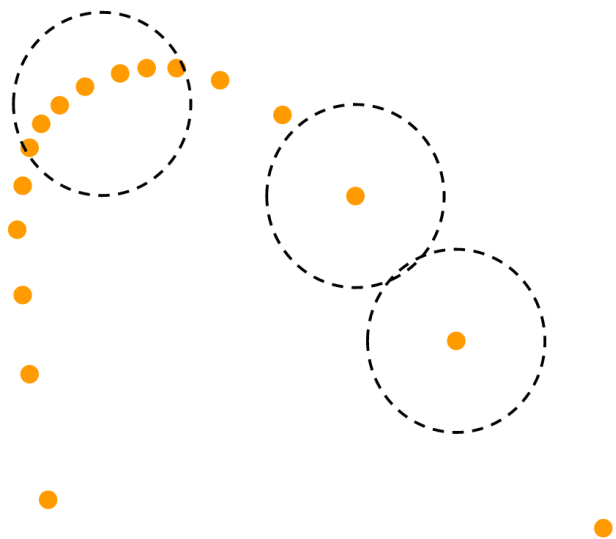


# Density Variation Affects Hierarchy



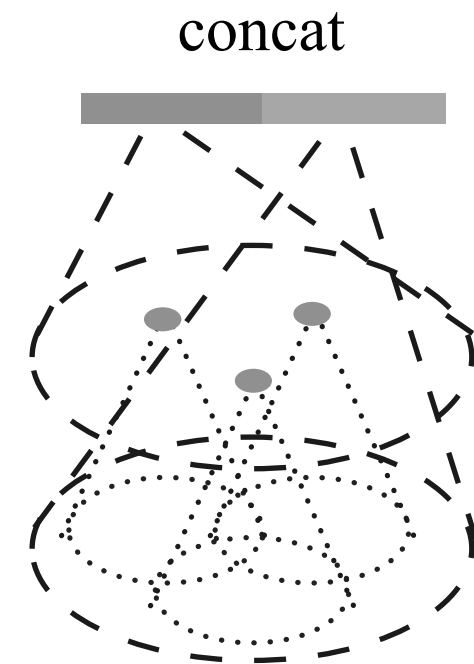
**Small kernels suffer from varying densities!**

# Robust Learning Under Varying Sampling Density



(a)

Multi-scale grouping (MSG)

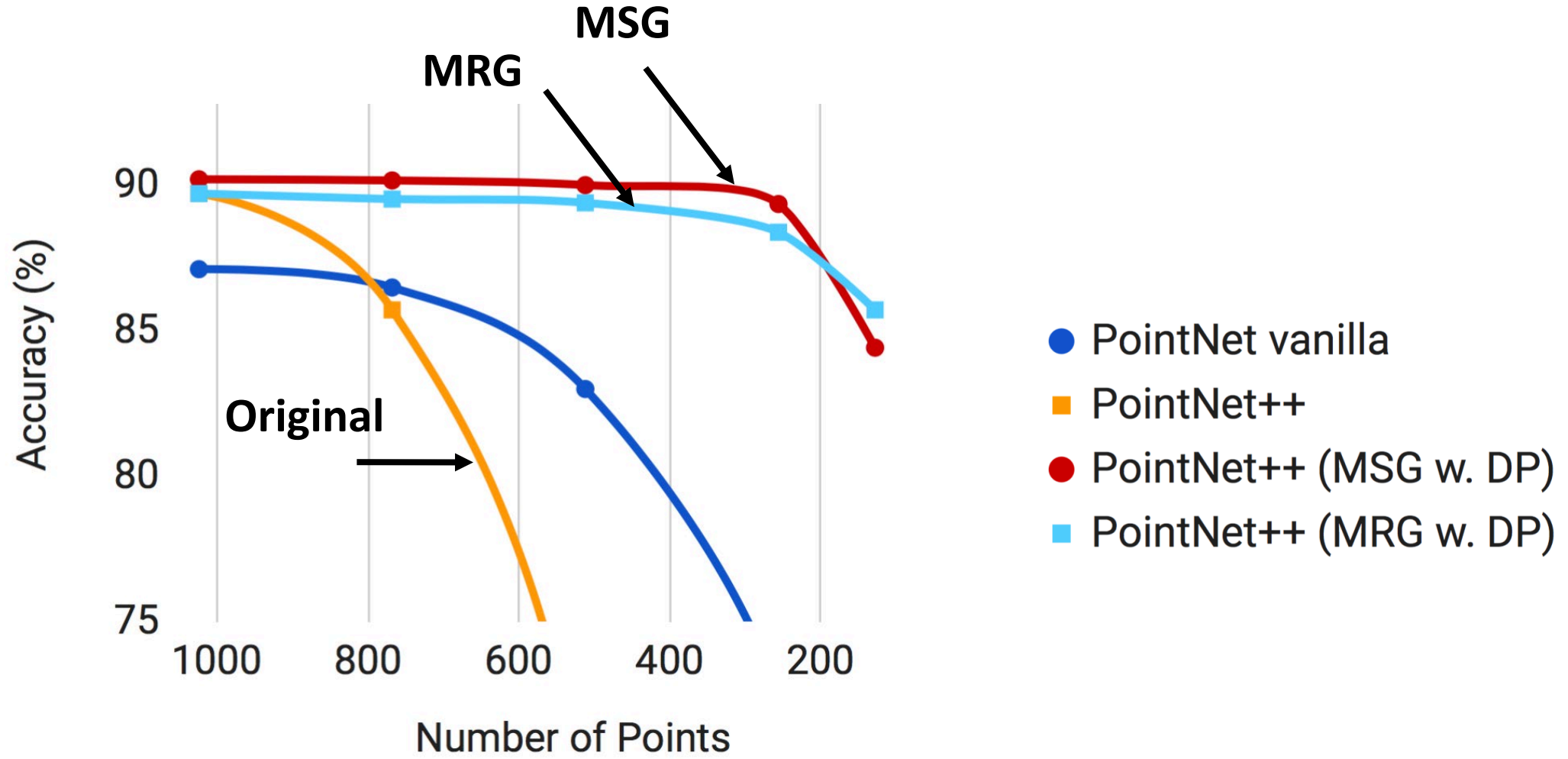


(b)

Multi-res grouping (MRG)

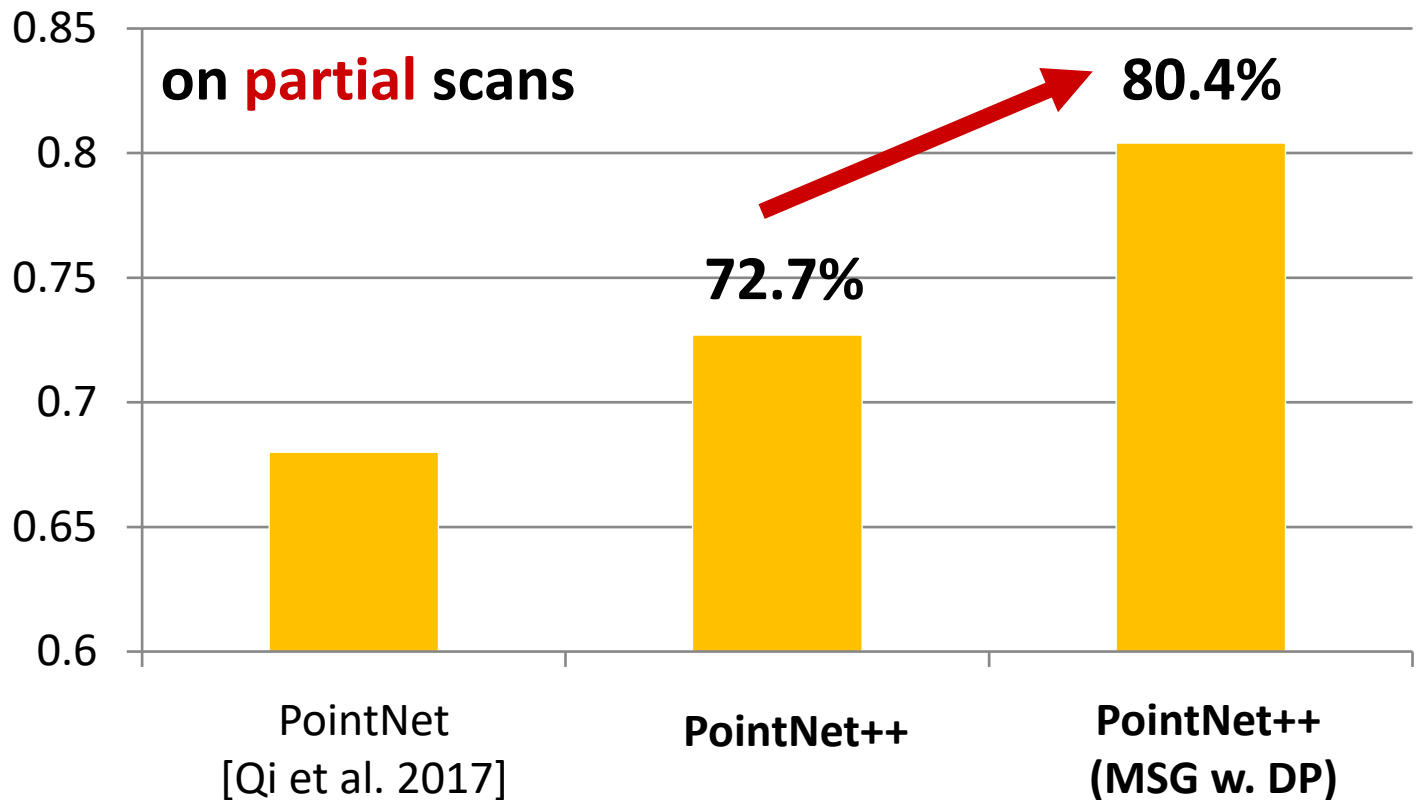
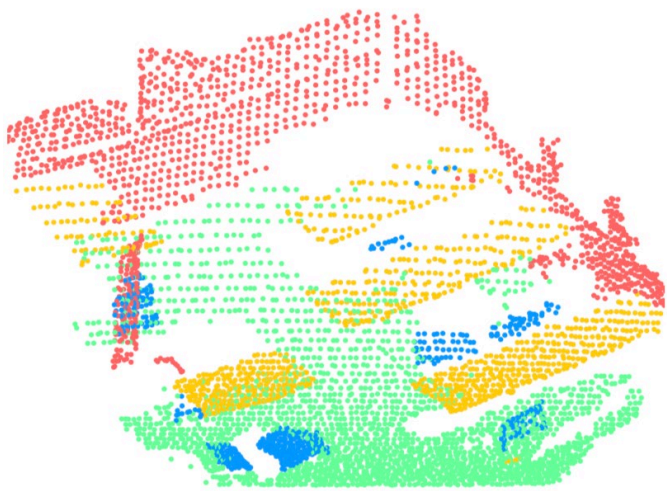
*During Training: input point dropout with random dropout ratio*

# Robust Learning Under Varying Sampling Density



# PointNet++ Results: Scene Parsing

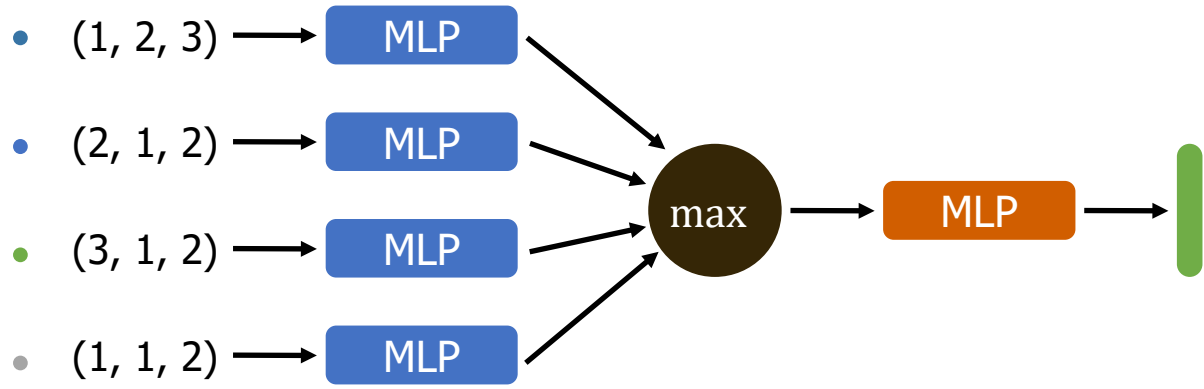
Robust layers for non-uniform densities (MSG) help a lot.



*dataset: ScanNet; metric: per-point semantic classification accuracy (%)*

# Graph Structures on Points: DGCNN

# Point Clouds and Graphs



PointNet family

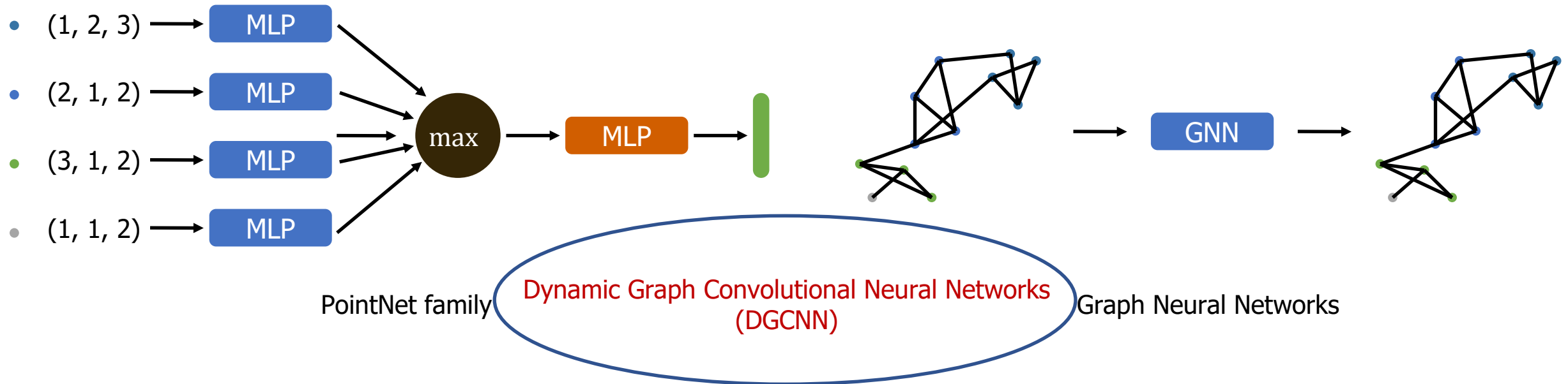


Graph Neural Networks

[Qi et al., CVPR 2017]

[Kipf et al., CVPR 2017]

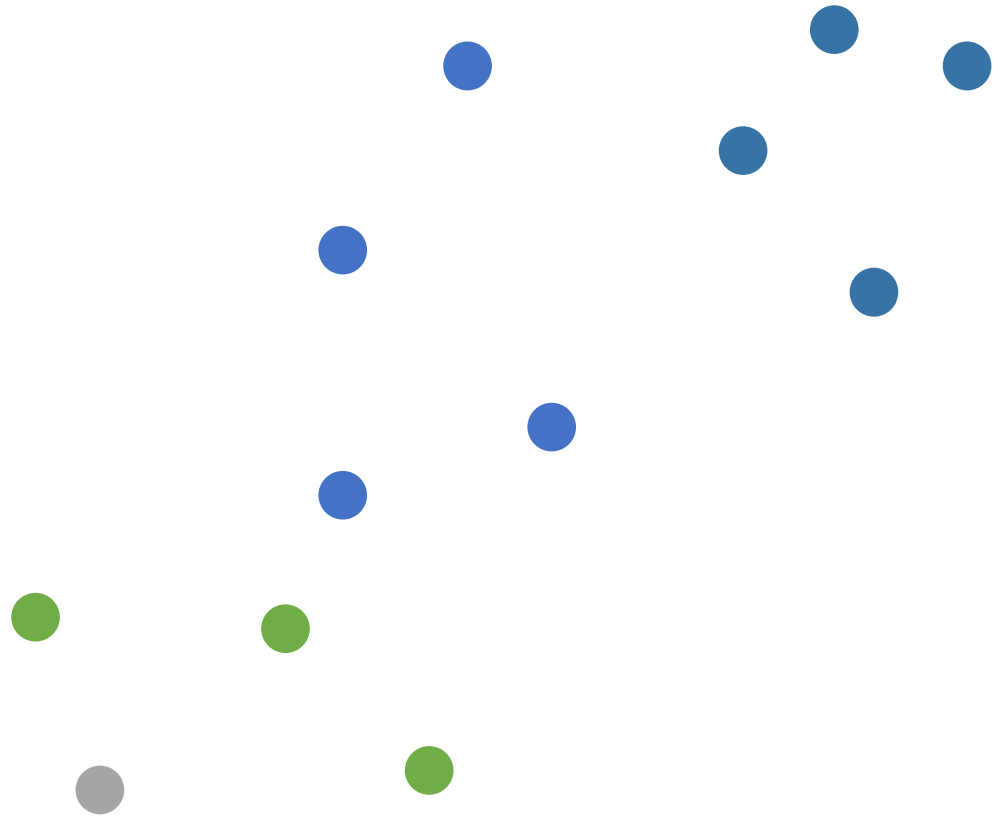
# Bridging the Two ... DGCNN



[Dynamic Graph CNN for Learning on Point Clouds

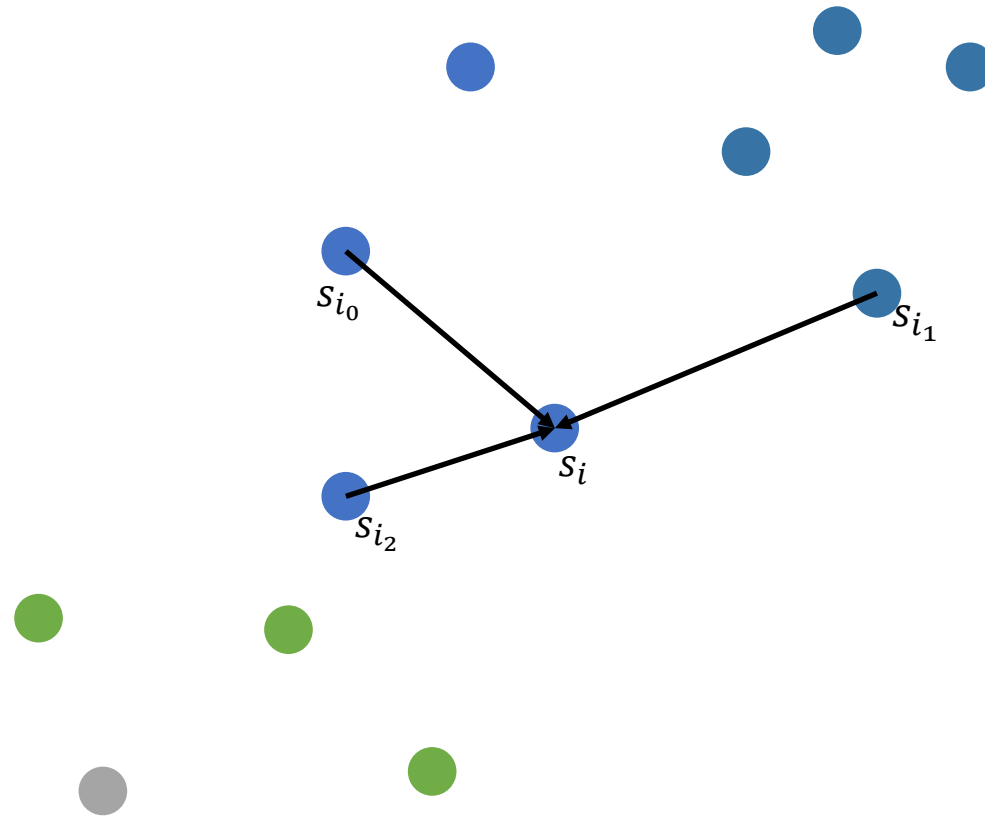
Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, Justin M. Solomon, TOG 2019]

# DGCNN



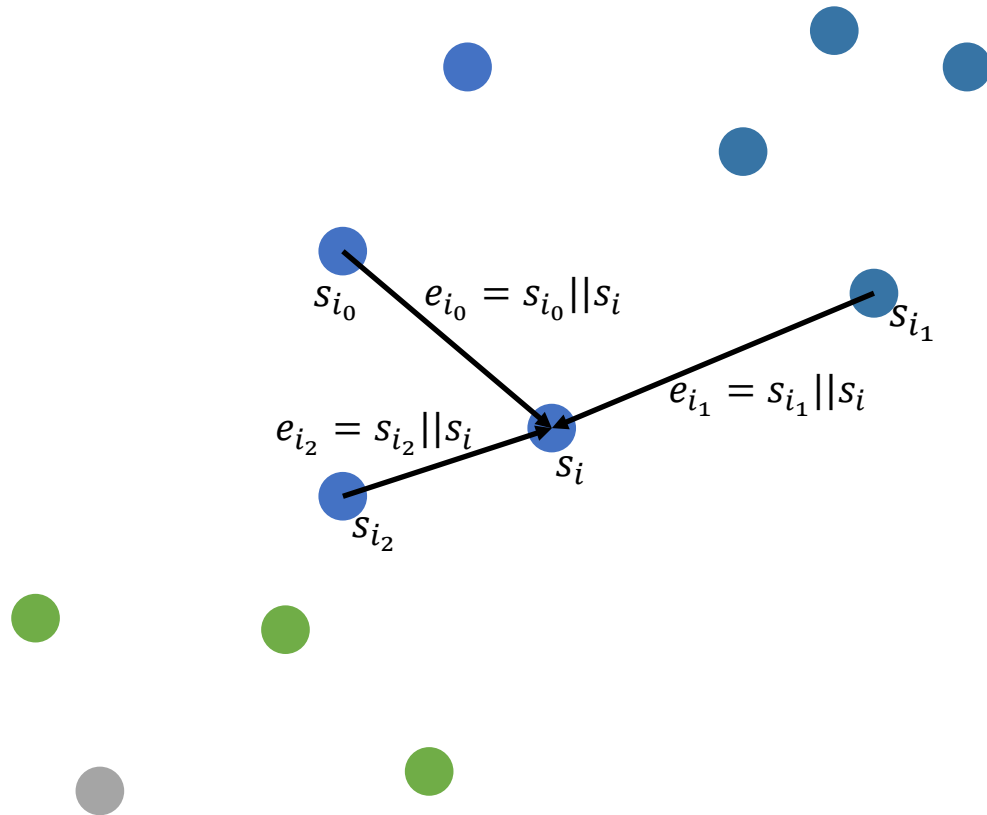
$$\mathcal{S} = \{(x_0, y_0, z_0), \dots, (x_i, y_i, z_i), \dots, (x_n, y_n, z_n)\}$$

# DGCNN



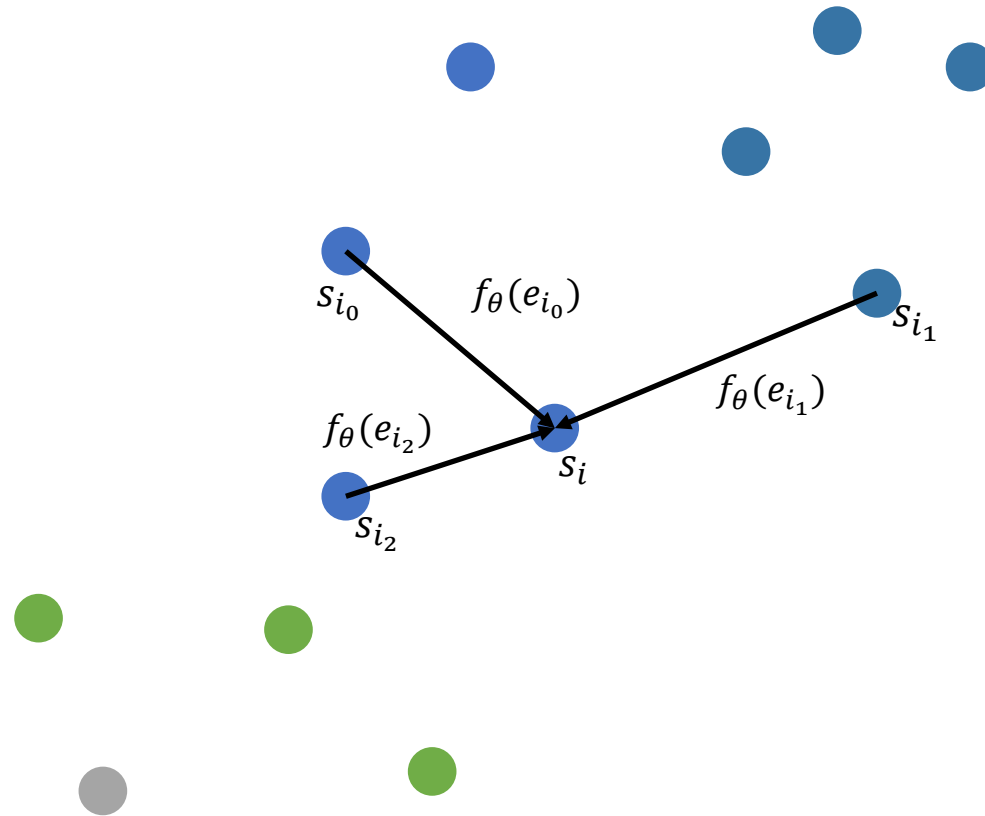
Each point (e.g.,  $s_i$ ) is connected to its  $k$  nearest neighbors

# DGCNN



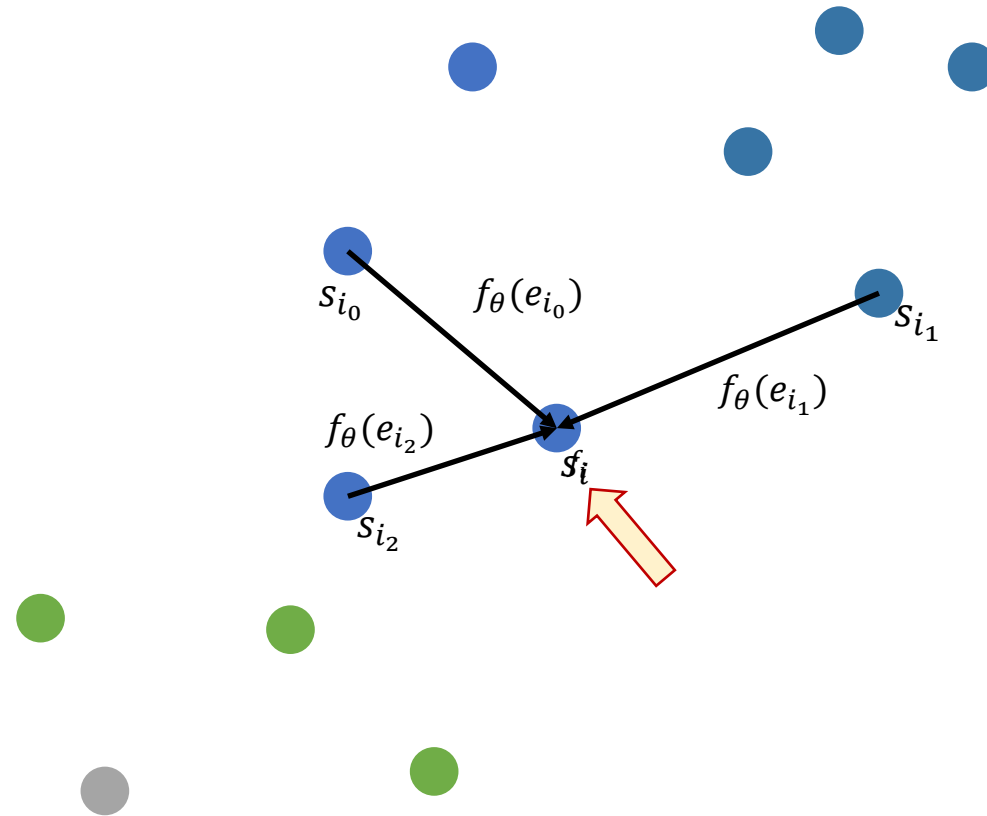
Edge features are defined by concatenating features of points

# DGCNN



A shared MLP further lifts edge features into high dimensional space

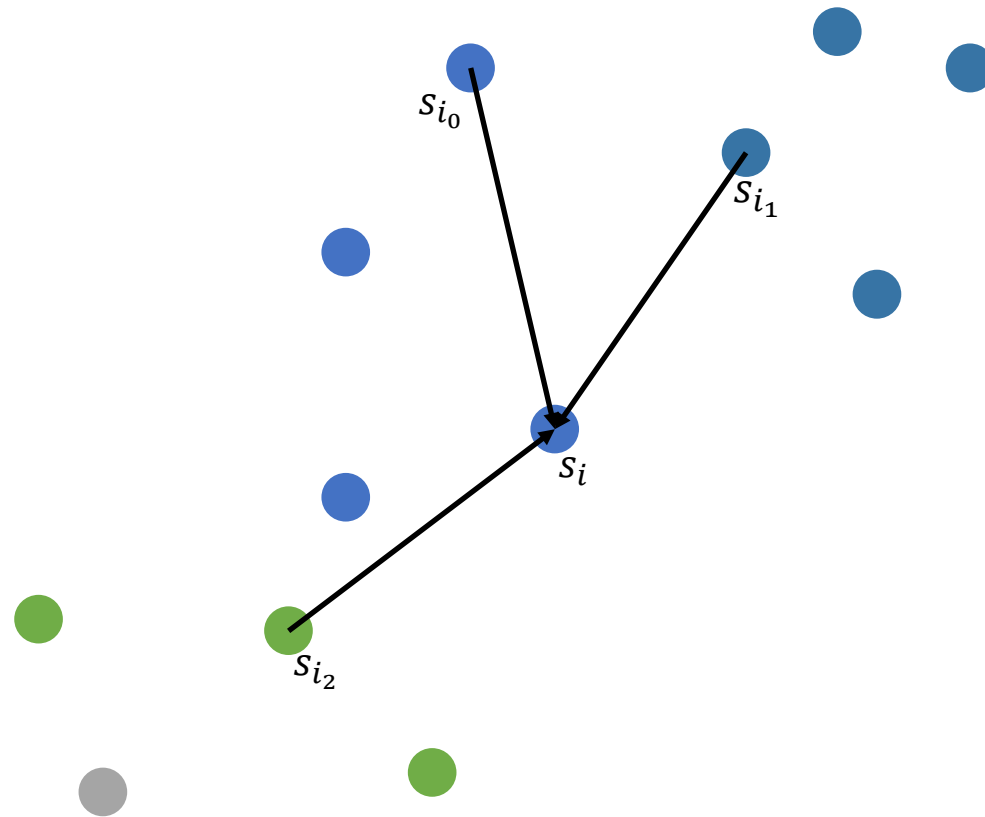
# DGCNN – the EdgeConv Operation



Features are pooled by a symmetric function

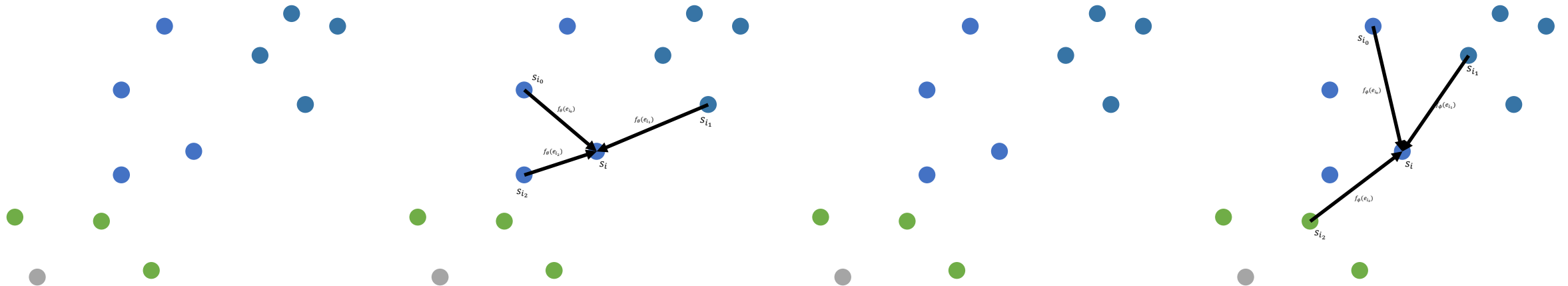
$$f_i = \max(f_{\theta}(e_{i_0}), f_{\theta}(e_{i_1}), f_{\theta}(e_{i_2}))$$

# DGCNN



Then, a new kNN graph is reconstructed based on features.

# DGCNN – Alternating Processing



DGCNN alternates feature learning (EdgeConvs) and graph reconstruction

# EdgeConv Operation

edge features:  $e_{ij} = h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j)$

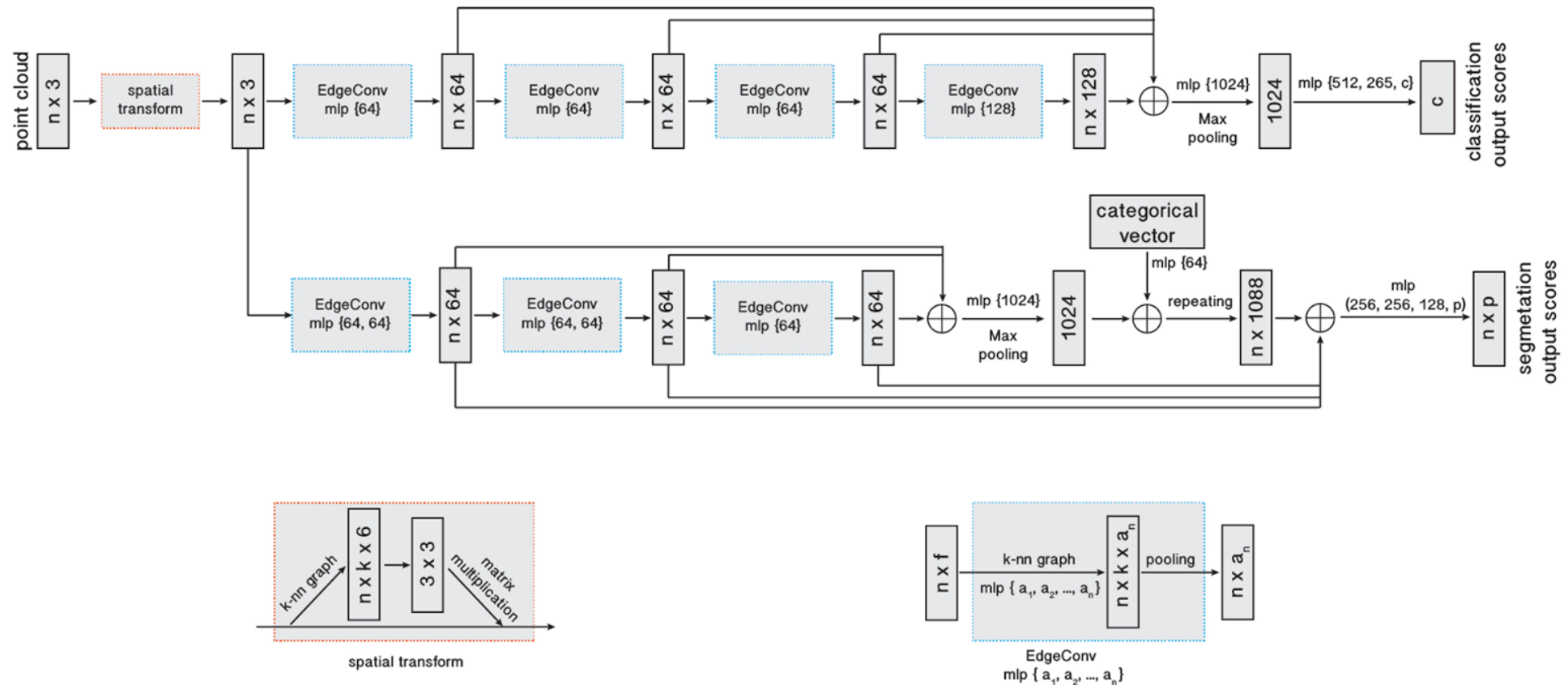
- general  $\mathbf{x}'_i = \square_{j:(i,j) \in \mathcal{E}} h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j).$  (1)
- weighted sum  $x'_{im} = \sum_{j:(i,j) \in \mathcal{E}} \theta_m \cdot \mathbf{x}_j.$  (2)
- global only  $h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j) = h_{\Theta}(\mathbf{x}_i),$  (3)
- local only  $h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j) = h_{\Theta}(\mathbf{x}_j - \mathbf{x}_i).$  (6)
- global + local  $h_{\Theta}(\mathbf{x}_i, \mathbf{x}_j) = \bar{h}_{\Theta}(\mathbf{x}_i, \mathbf{x}_j - \mathbf{x}_i).$  (7)

# Key EdgeConv Properties

- Easily implement and integrates into existing deep learning-based algorithm by switching the MLP component to EdgeConv
- EdgeConv is differentiable, which is an important property in ML and DL (convex optimization problem)
- Extracts local features without destroying the permutation invariance
- Dynamic Graph CNN => update the graph after each layer of network, i.e. recompute the k-nearest neighbors in the new feature space (and also recompute the edge features)

# DGCNN Architecture

## DGCNN



# DGCNN Results

	MEAN CLASS ACCURACY	OVERALL ACCURACY
3DShapeNets [Wu et al. 2015]	77.3	84.7
VoxNet [Maturana and Scherer 2015]	83.0	85.9
SubVolume [Qi et al. 2016]	86.0	89.2
VRN (Single View) [Brock et al. 2016]	88.98	-
VRN (Multiple Views) [Brock et al. 2016]	91.33	-
ECC [Simonovsky and Komodakis 2017]	83.2	87.4
PointNet [Qi et al. 2017b]	86.0	89.2
PointNet++ [Qi et al. 2017c]	-	90.7
KD-Net [Klokov and Lempitsky 2017]	-	90.6
PointCNN [Li et al. 2018a]	88.1	92.2
PCNN [Atzmon et al. 2018]	-	92.3
Ours (Baseline)	88.9	91.7
Ours	<b>90.2</b>	<b>92.9</b>
Ours (2048 Points)	<b>90.7</b>	<b>93.5</b>

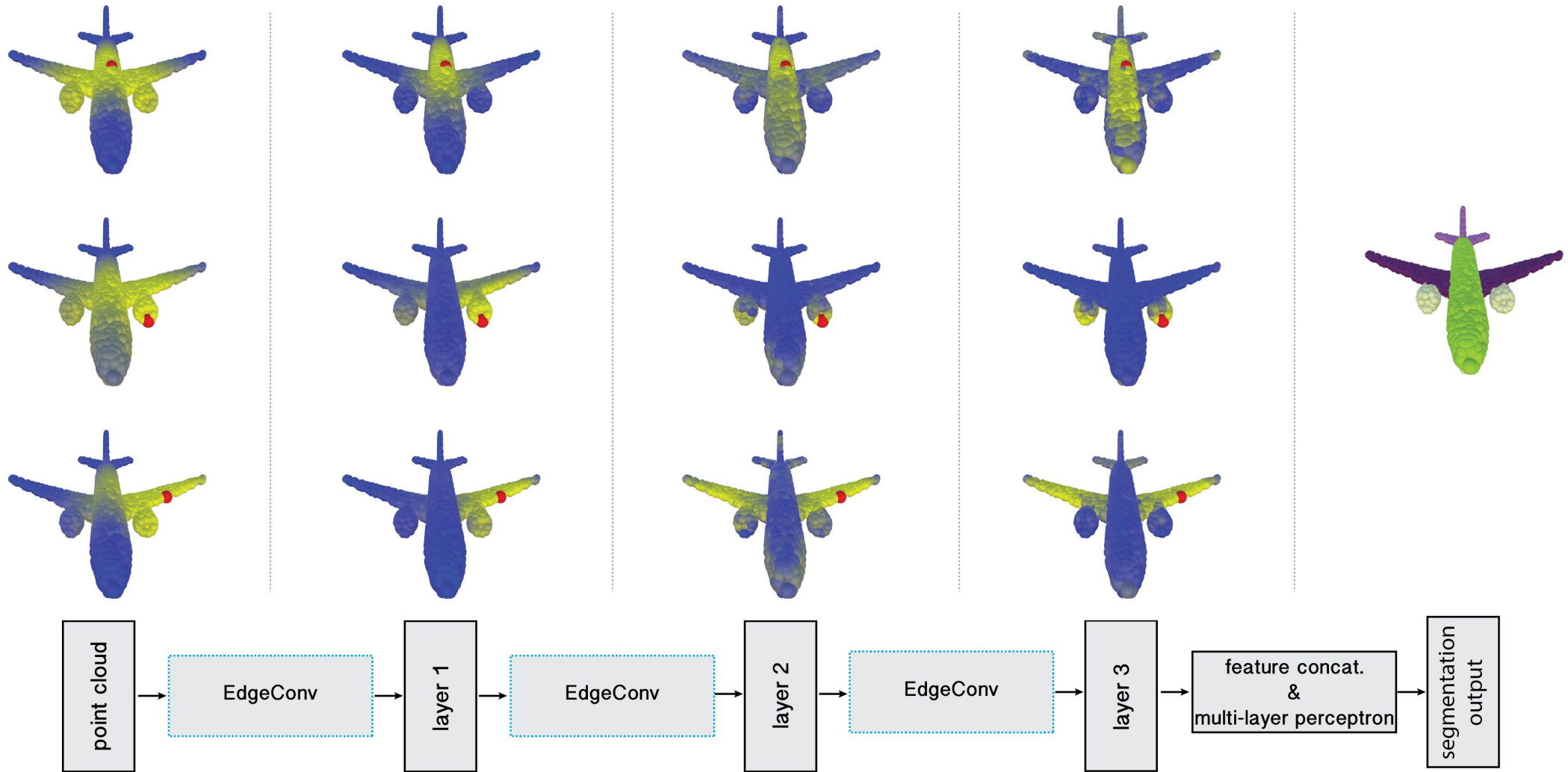
Table 2. Classification results on ModelNet40.

DGCNN achieves superior performance on shape classification tasks while maintaining simplicity!

	MODEL SIZE(MB)	TIME(MS)	ACCURACY(%)
PointNet (Baseline) [Qi et al. 2017b]	9.4	6.8	87.1
PointNet [Qi et al. 2017b]	40	16.6	89.2
PointNet++ [Qi et al. 2017c]	12	163.2	90.7
PCNN [Atzmon et al. 2018]	94	117.0	92.3
Ours (Baseline)	11	19.7	91.7
Ours	21	27.2	92.9

Table 3. Complexity, forward time, and accuracy of different models

# From Geometry to Semantics



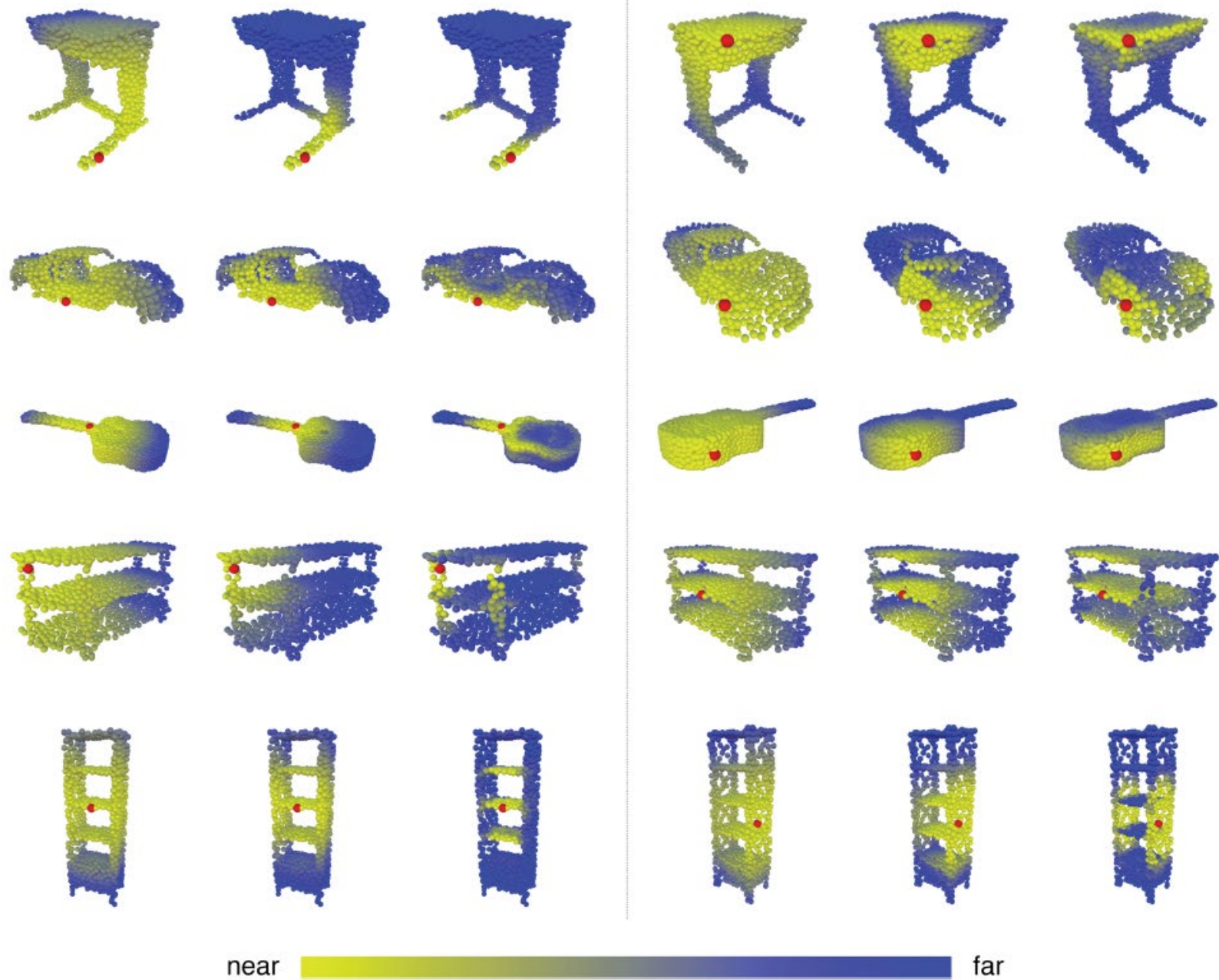
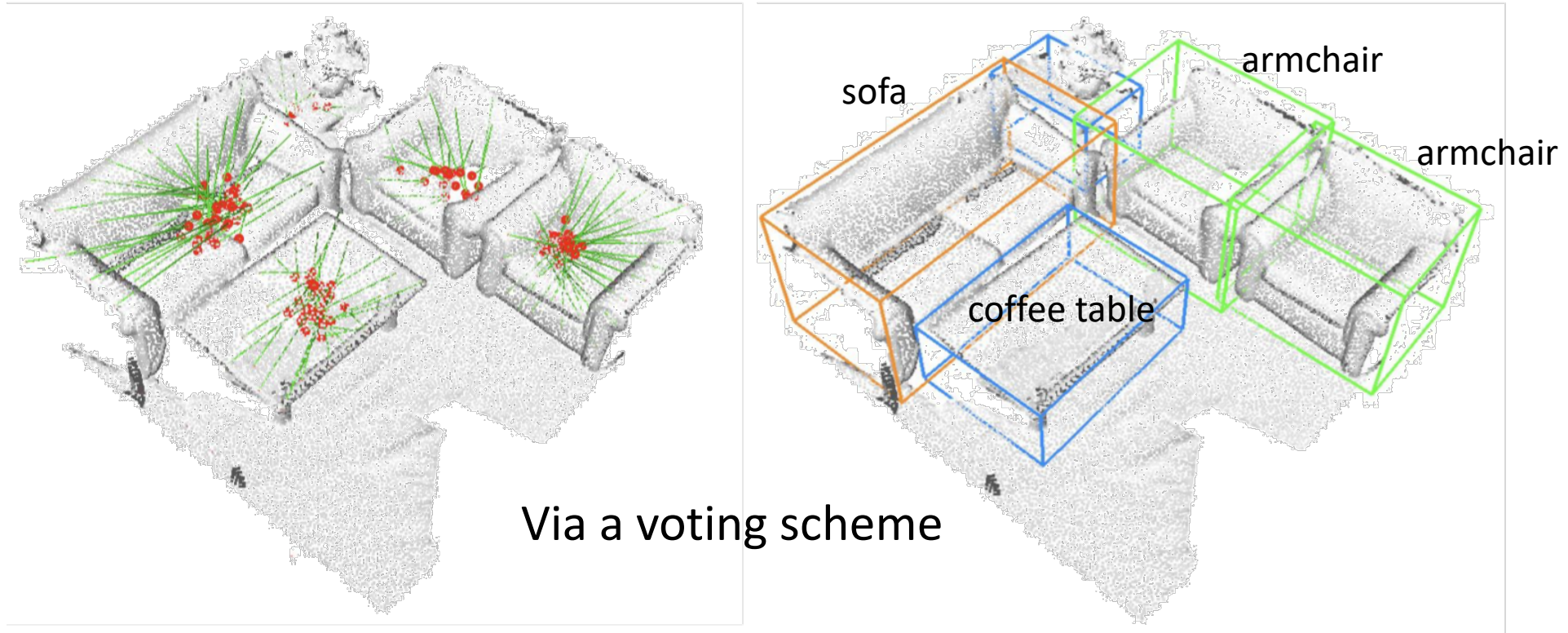


Fig. 4. Structure of the feature spaces produced at different stages of our shape classification neural network architecture, visualized as the distance between the red point to the rest of the points. For each set, Left: Euclidean distance in the input  $\mathbb{R}^3$  space; Middle: Distance after the point cloud transform stage, amounting to a global transformation of the shape; Right: Distance in the feature space of the last layer. Observe how in the feature space of deeper layers semantically similar structures such as shelves of a bookshelf or legs of a table are brought close together, although they are distant in the original space.

# Object Detection in Point Clouds, Indoors

# Point Cloud Object Amodal Bounding Box Detection



- Charles R. Qi, Or Litany, Kaiming He, Leonidas J. Guibas. *Deep Hough Voting for 3D Object Detection in Point Clouds*. ICCV 2019.
- Charles R. Qi, Xinlei Chen, Or Litany, Leonidas J. Guibas. *ImVoteNet: Boosting 3D Object Detection in Point Clouds with Image Votes*. CVPR 2020.

# Generalized Hough Transform

## GENERALIZING THE HOUGH TRANSFORM TO DETECT ARBITRARY SHAPES\*

D. H. BALLARD

Computer Science Department, University of Rochester, Rochester, NY 14627, U.S.A.

(Received 10 October 1979; in revised form 9 September 1980; received for  
publication 23 September 1980)

**Abstract**—The Hough transform is a method for detecting curves by  
a curve and parameters of that curve.

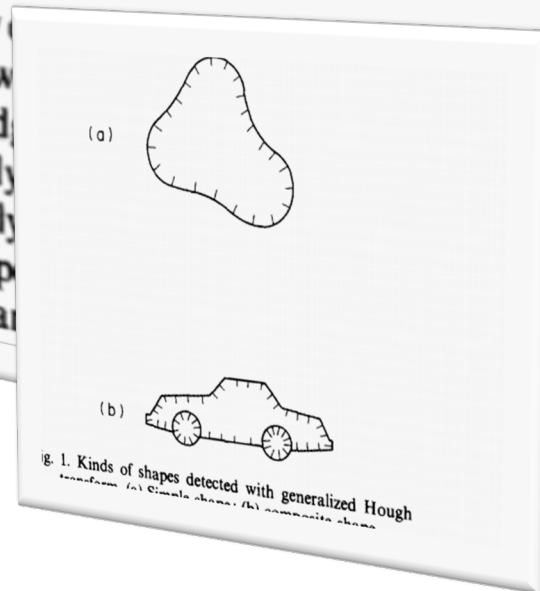
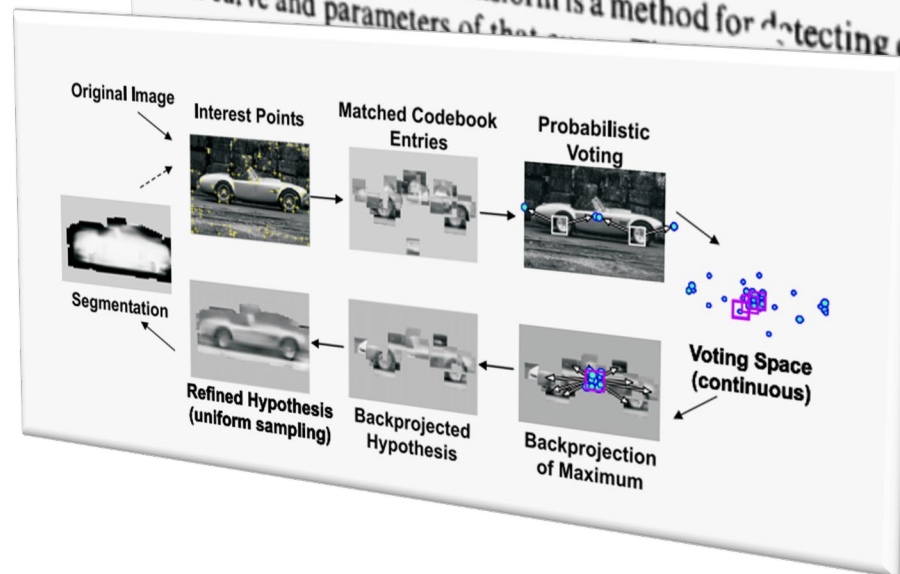
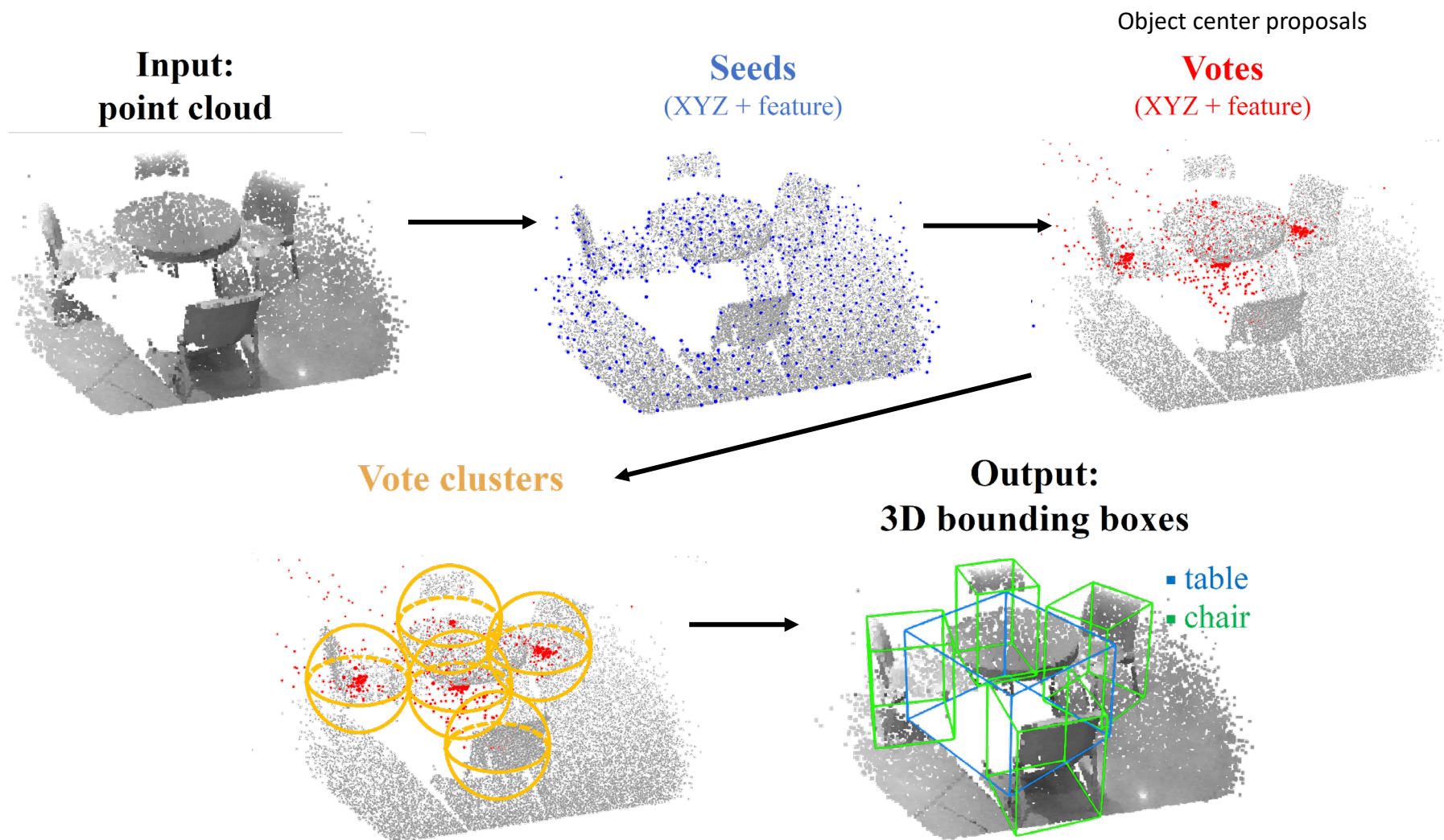
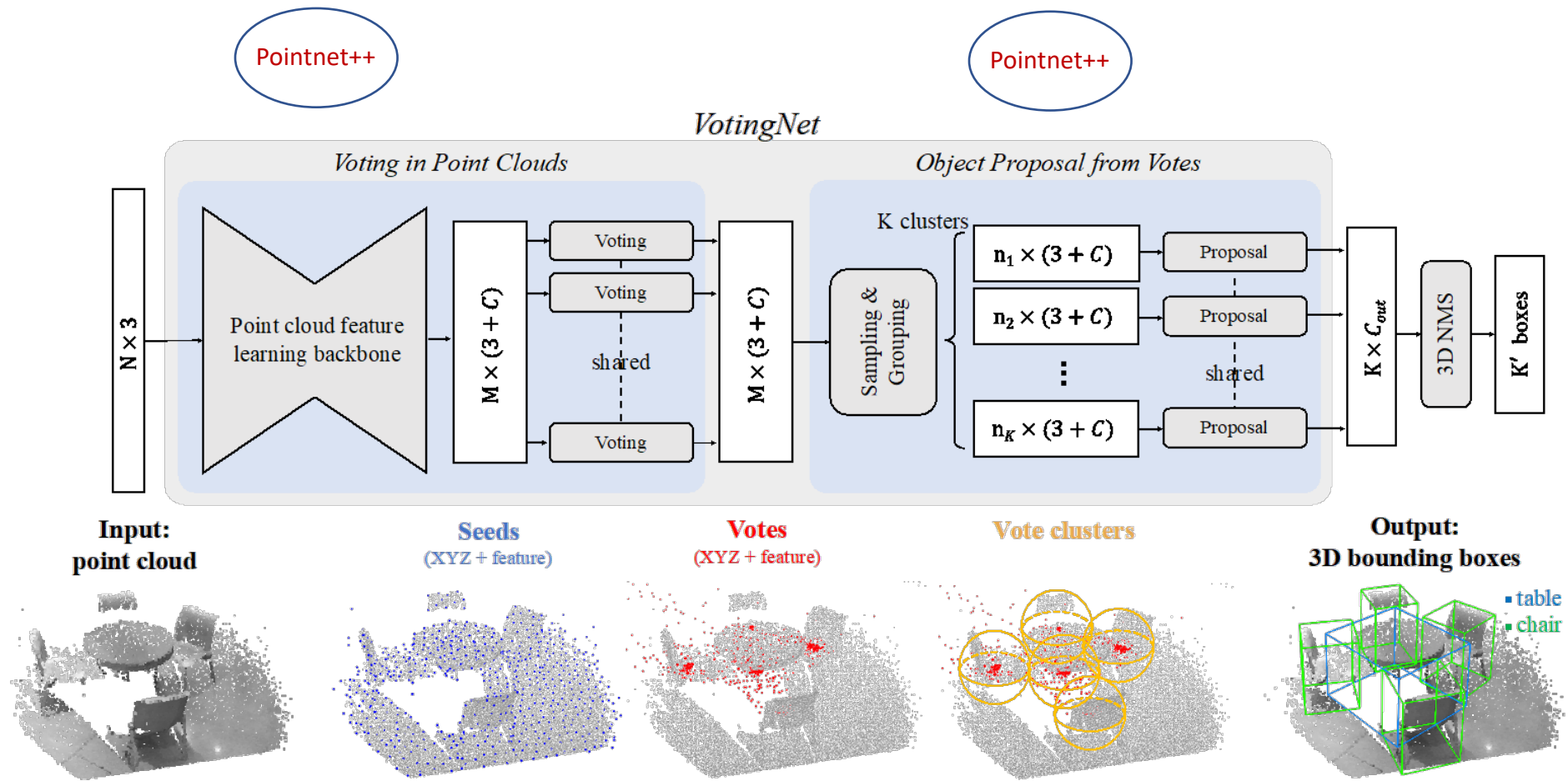


Fig. 1. Kinds of shapes detected with generalized Hough transform. (a) Simple shape; (b) composite shape.

# Deep Hough Voting – A Two-Stage Approach



# VoteNet – A Two-Stage Approach



A capsule/transformer network in disguise ...



# VoteNet Results on SUN RGB-D

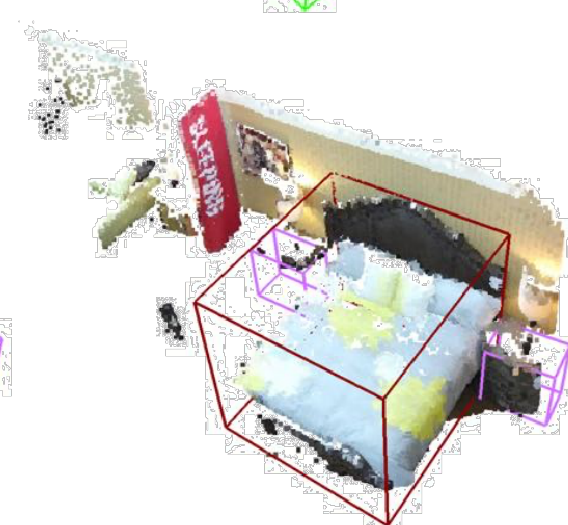
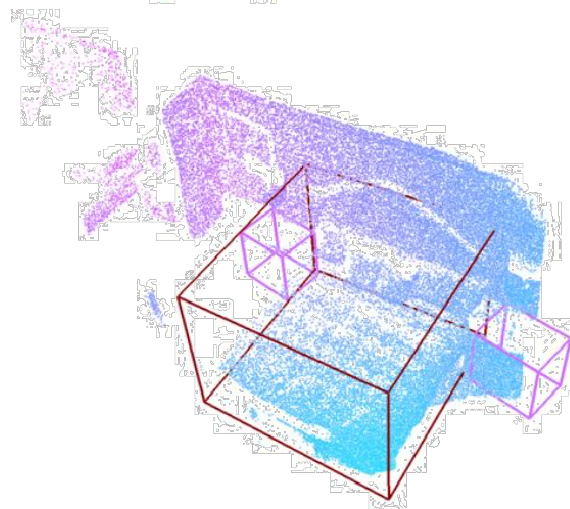
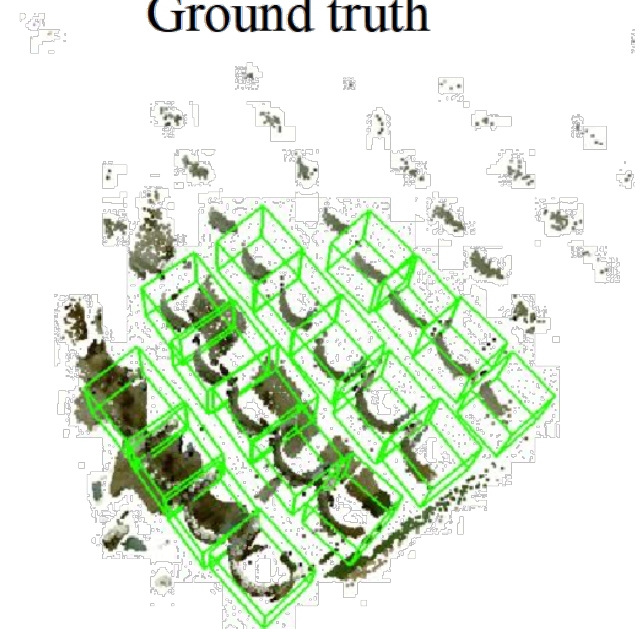
Image of the scene



VotingNet prediction



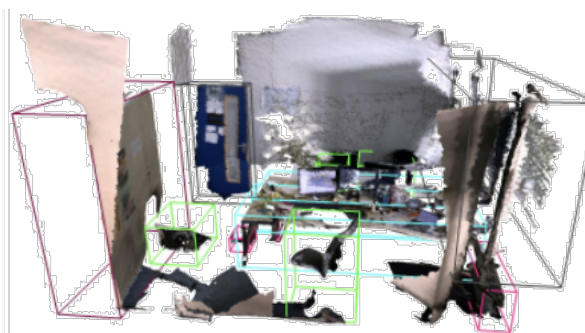
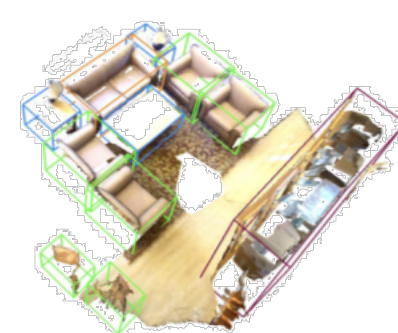
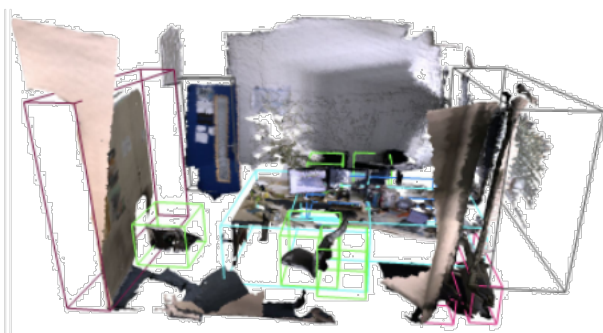
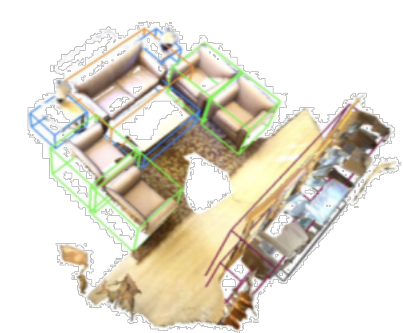
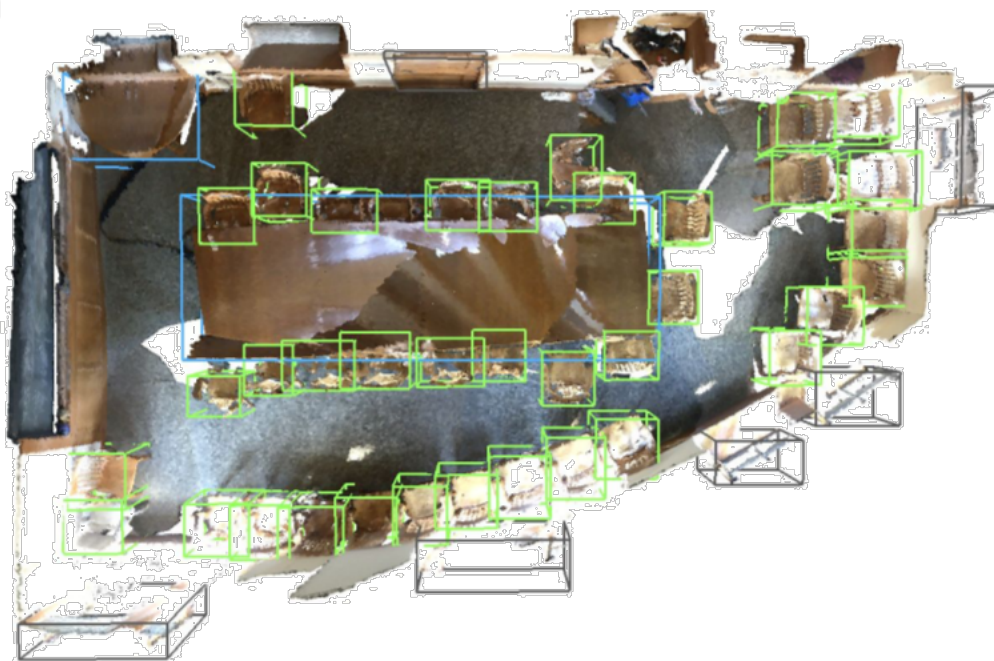
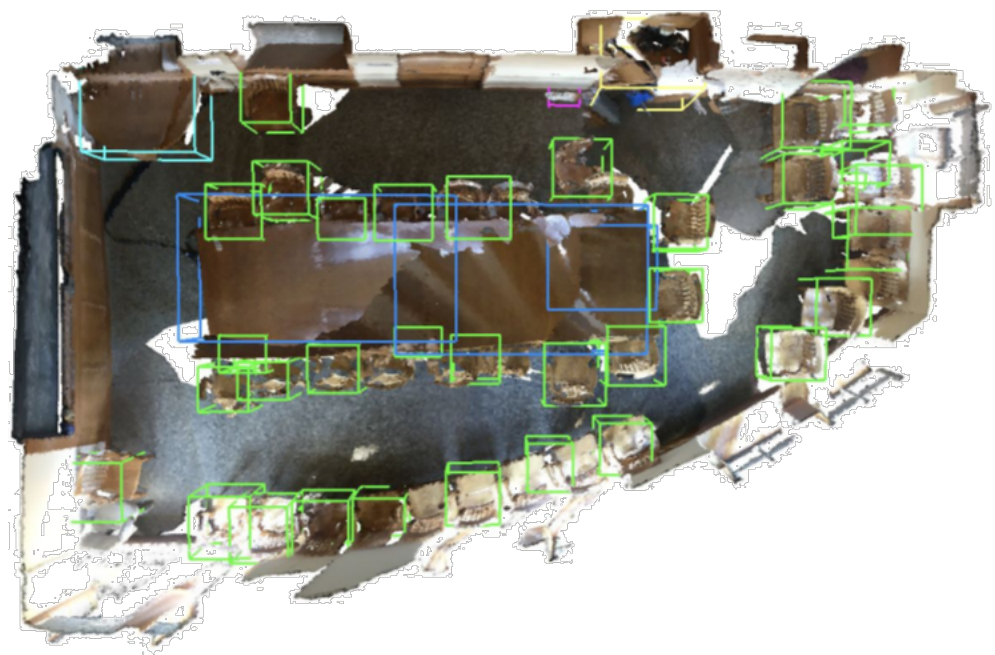
Ground truth



# VoteNet Results on ScanNet

VotingNet prediction

Ground truth



# VoteNet Quantitative Results

average precision with 3D IoU threshold 0.25

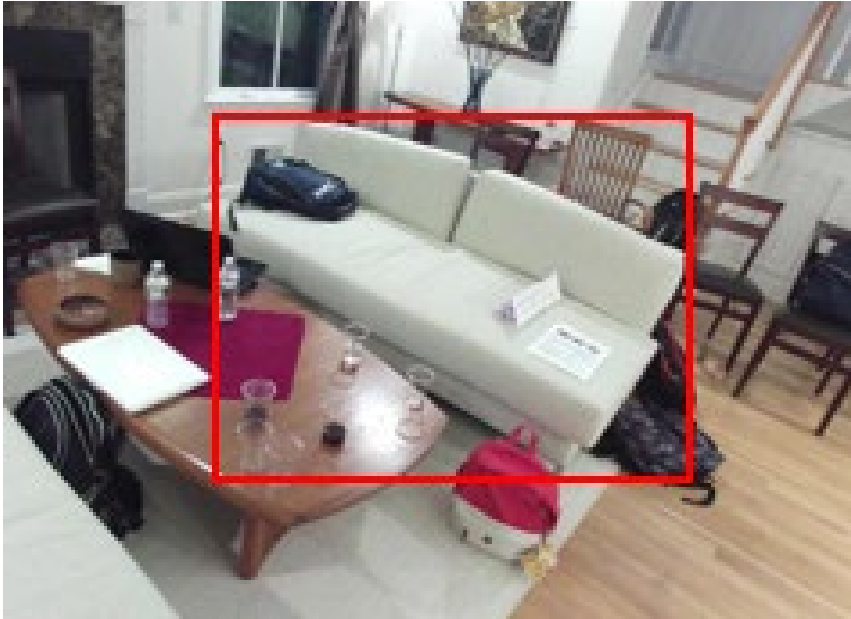
## SUN RGB-D

	Input	bathhtub	bed	bookshelf	chair	desk	dresser	nightstand	sofa	table	toilet	mAP
Deep sliding shapes	Geo + RGB	44.2	78.8	11.9	61.2	20.5	6.4	15.4	53.5	50.3	78.9	42.1
Clouds of oriented gradients	Geo + RGB	58.3	63.7	31.8	62.2	<b>45.2</b>	15.5	27.4	51.0	<b>51.3</b>	70.1	47.6
	Geo + RGB	43.5	64.5	31.4	48.3	27.9	25.9	41.9	50.4	37.0	80.4	45.1
Frustum pointnet	Geo + RGB	43.3	81.1	<b>33.3</b>	64.2	24.7	<b>32.0</b>	58.1	61.1	51.1	<b>90.9</b>	54.0
	Geo only	<b>74.4</b>	<b>83.0</b>	28.8	<b>75.3</b>	22.0	29.8	<b>62.2</b>	<b>64.0</b>	47.3	90.1	<b>57.7</b>

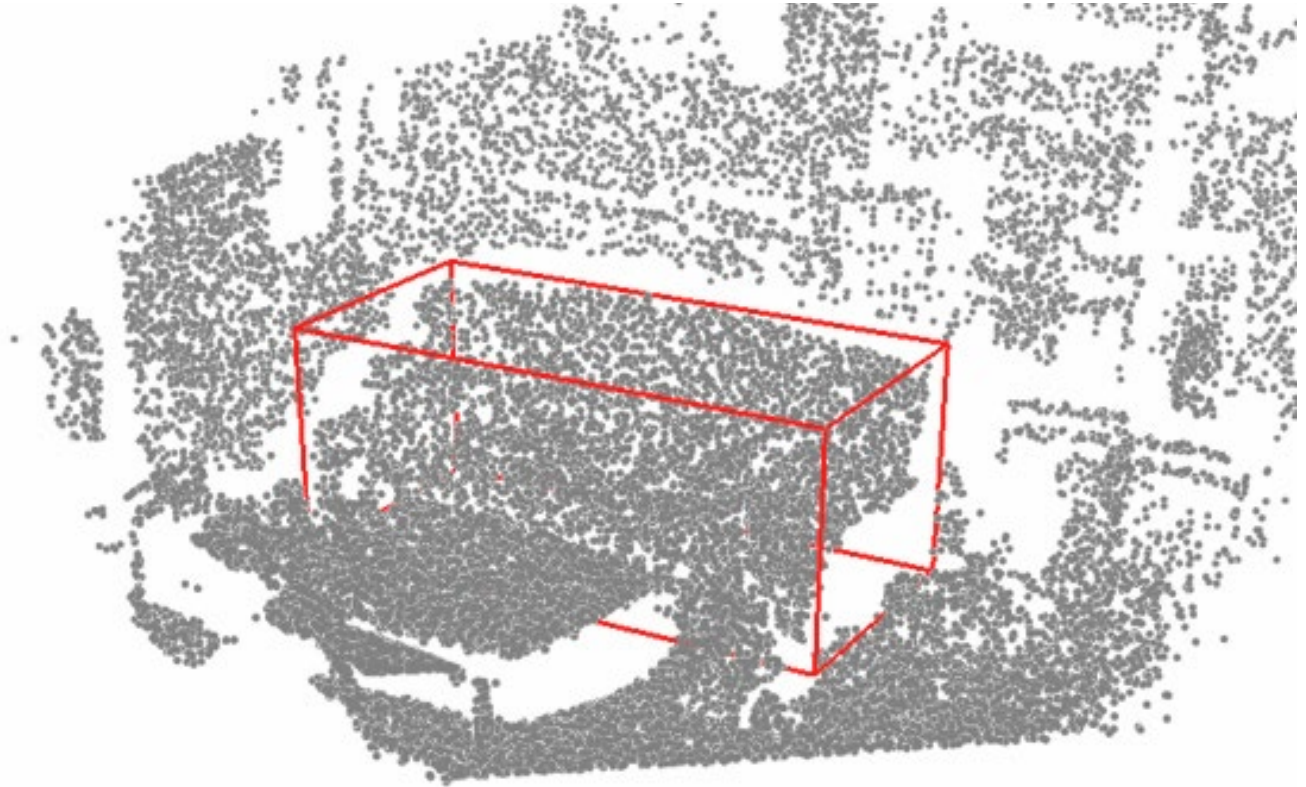
## ScanNetV2

	Input	mAP@0.25	mAP@0.5
DSS [42, 12]	Geo + RGB	15.2	6.8
MRCNN 2D-3D [11, 12]	Geo + RGB	17.3	10.5
F-PointNet [34, 12]	Geo + RGB	19.8	10.8
GSPN [54]	Geo + RGB	30.6	17.7
3D-SIS [12]	Geo + 1 view	35.1	18.7
3D-SIS [12]	Geo + 3 views	36.6	19.0
3D-SIS [12]	Geo + 5 views	40.2	22.5
3D-SIS [12]	Geo only	25.4	14.6
VoteNet (ours)	Geo only	<b>58.6</b>	<b>33.5</b>

# Modalities: Point Clouds Complement RGB Images



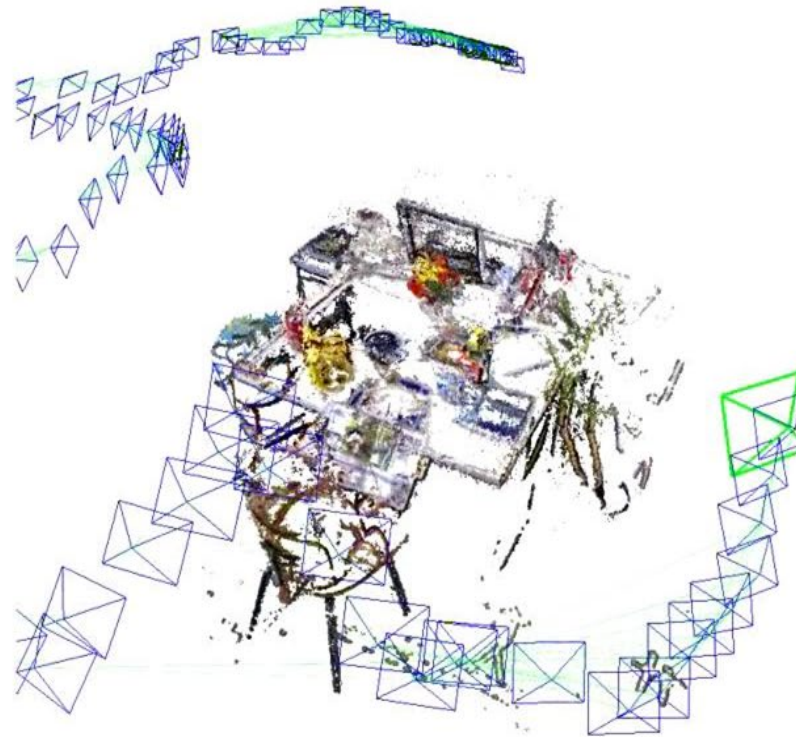
- + High resolution
- + Dense coverage
- - Subject to many imaging artifacts



- + Absolute depth and scale
- - Sparse, low rez

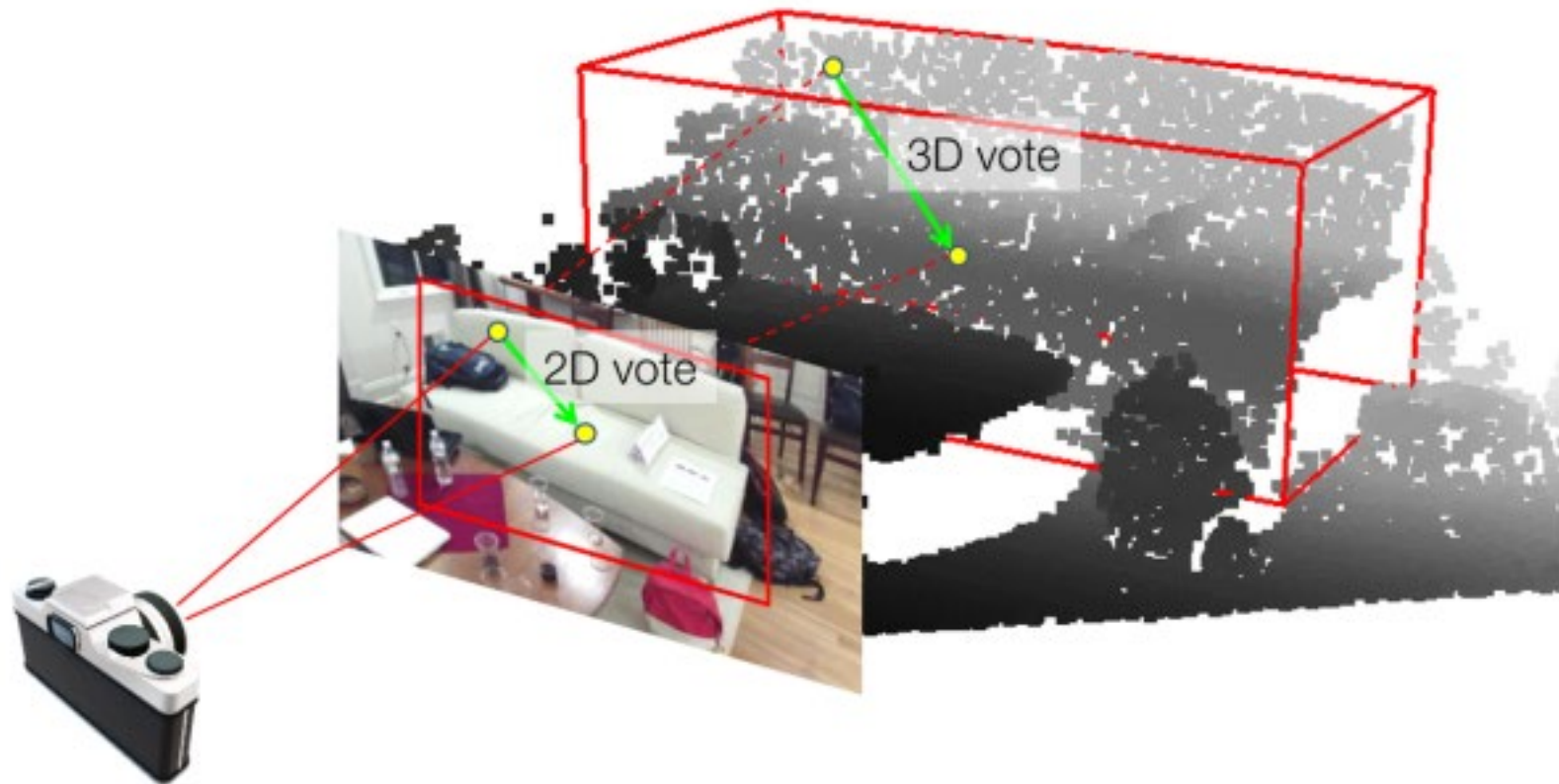
# 3D Detection with Sparse Points

**Application:** 3D detection from monocular video, using sparse SLAM keypoints.



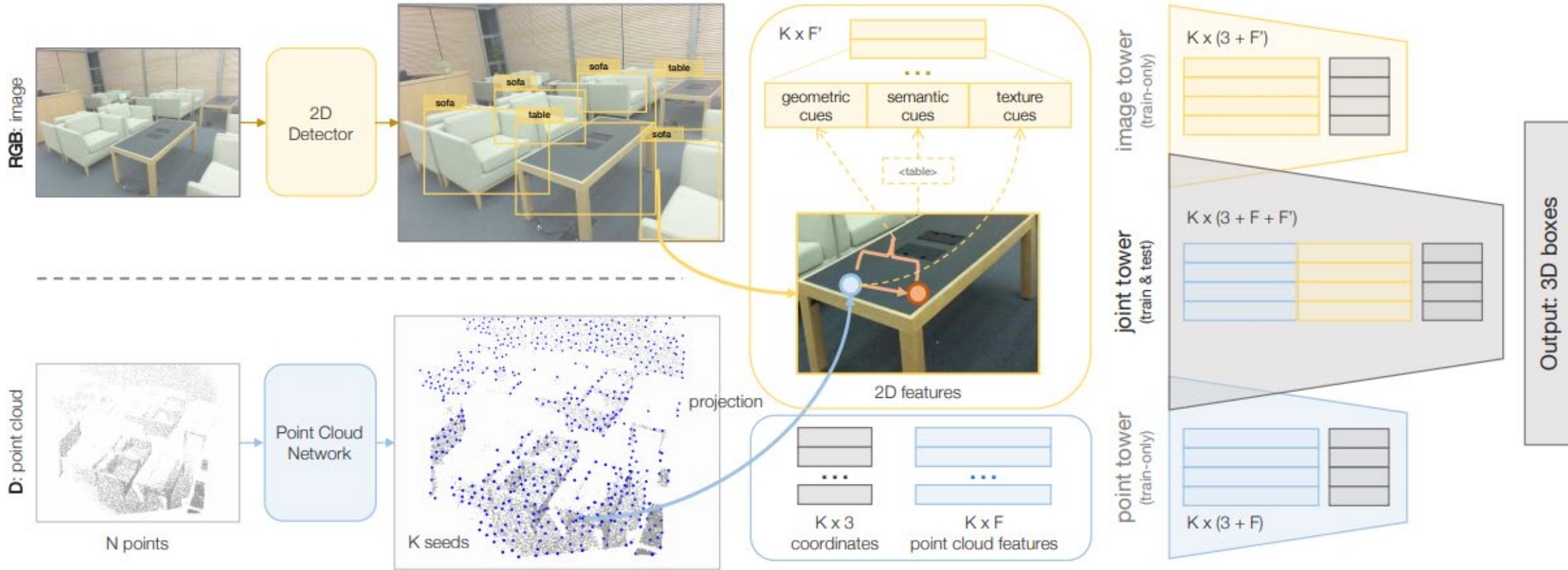
*Picture: ORB-SLAM results*

# Points and Images: ImVoteNet



How to best combine geometry and appearance info?

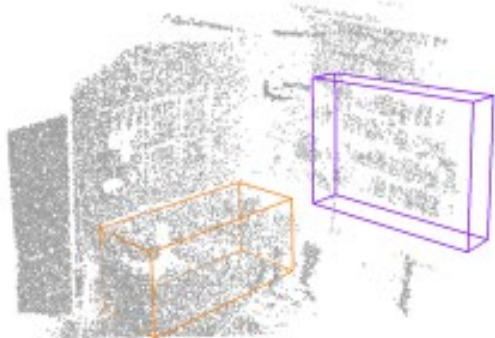
# ImVoteNet Architecture



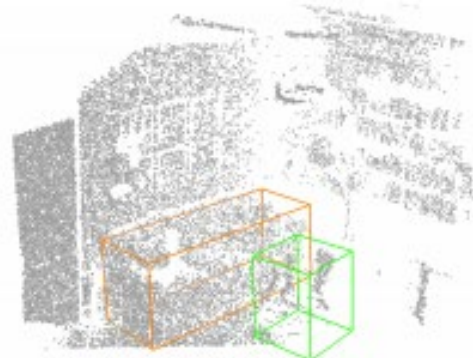
Ours 2D detection



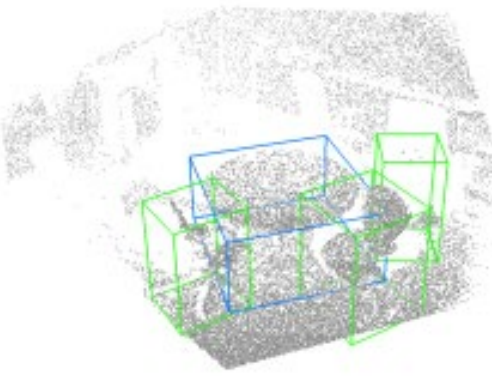
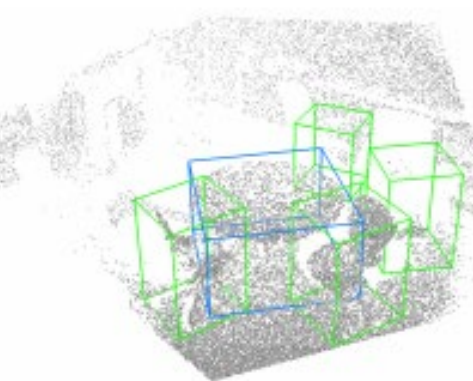
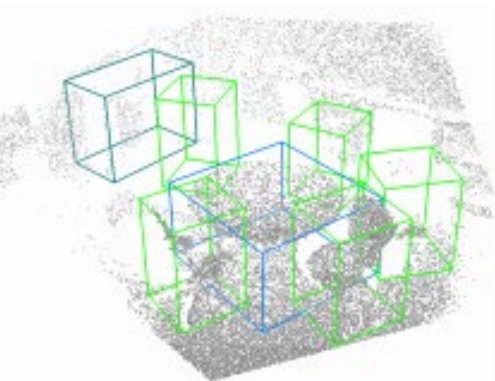
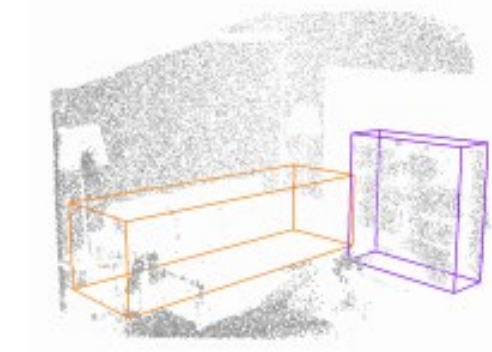
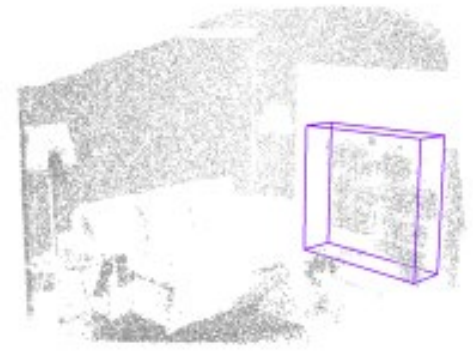
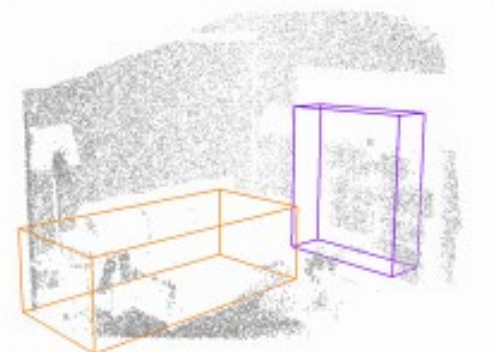
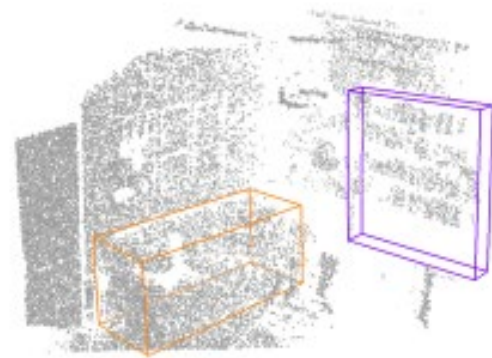
Ours 3D detection



VoteNet

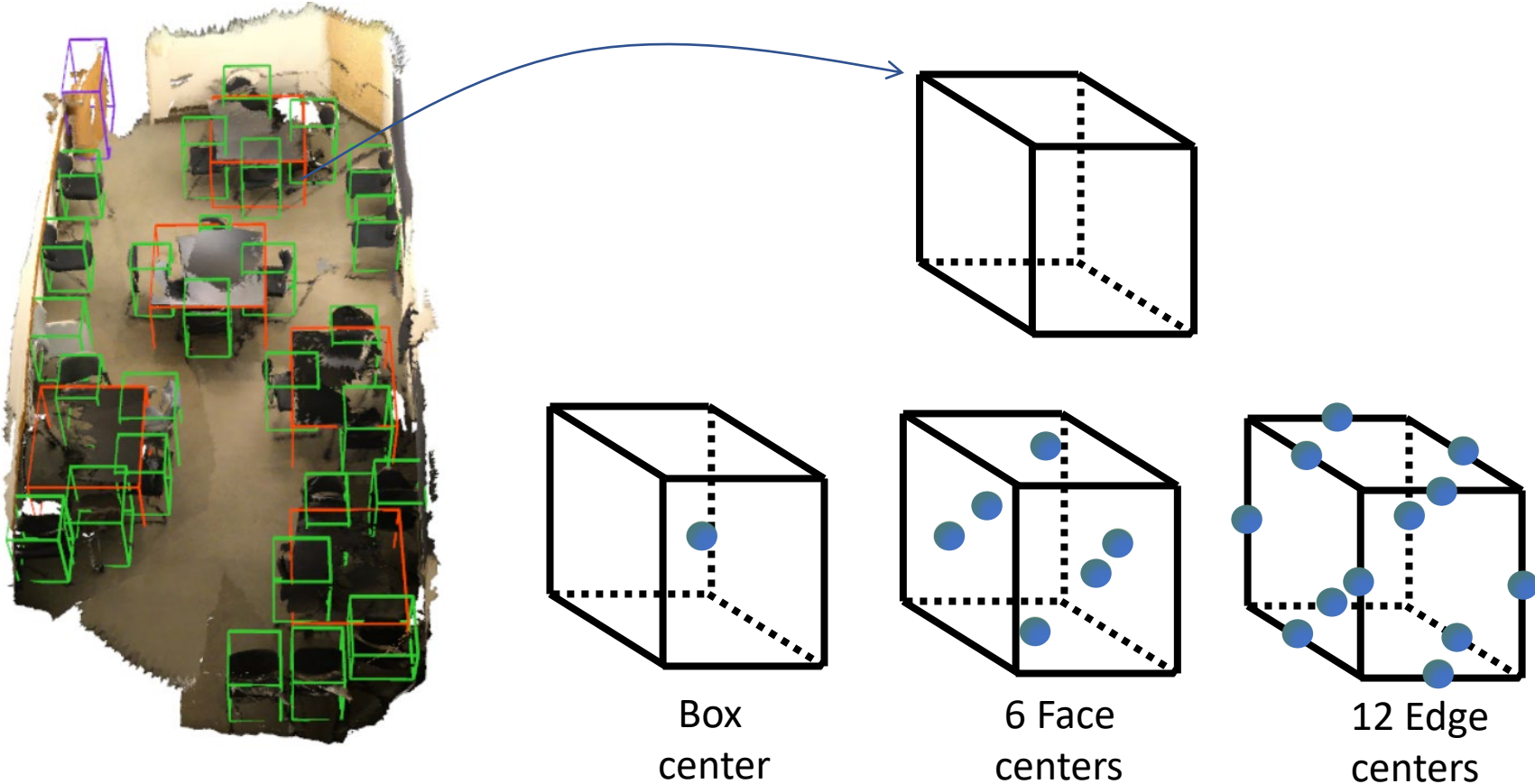


Ground truth



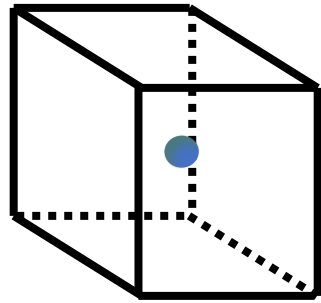
■ sofa ■ bookshelf ■ chair ■ table ■ desk

# Multiple 3D Geometry Representations

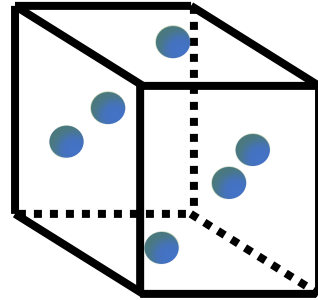


[Z. Zhang, B. Sun, H. Yang, Q. Huang.  
**H3DNet: 3D Object Detection Using Hybrid Geometric Primitives.**  
ECCV 2020]

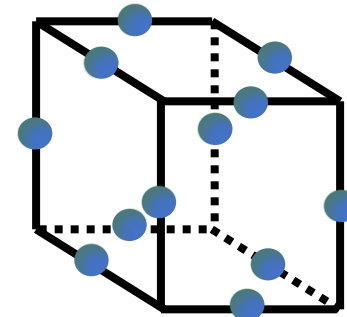
# Representations Best for Different Object Instances



Box  
center



6 Face  
centers

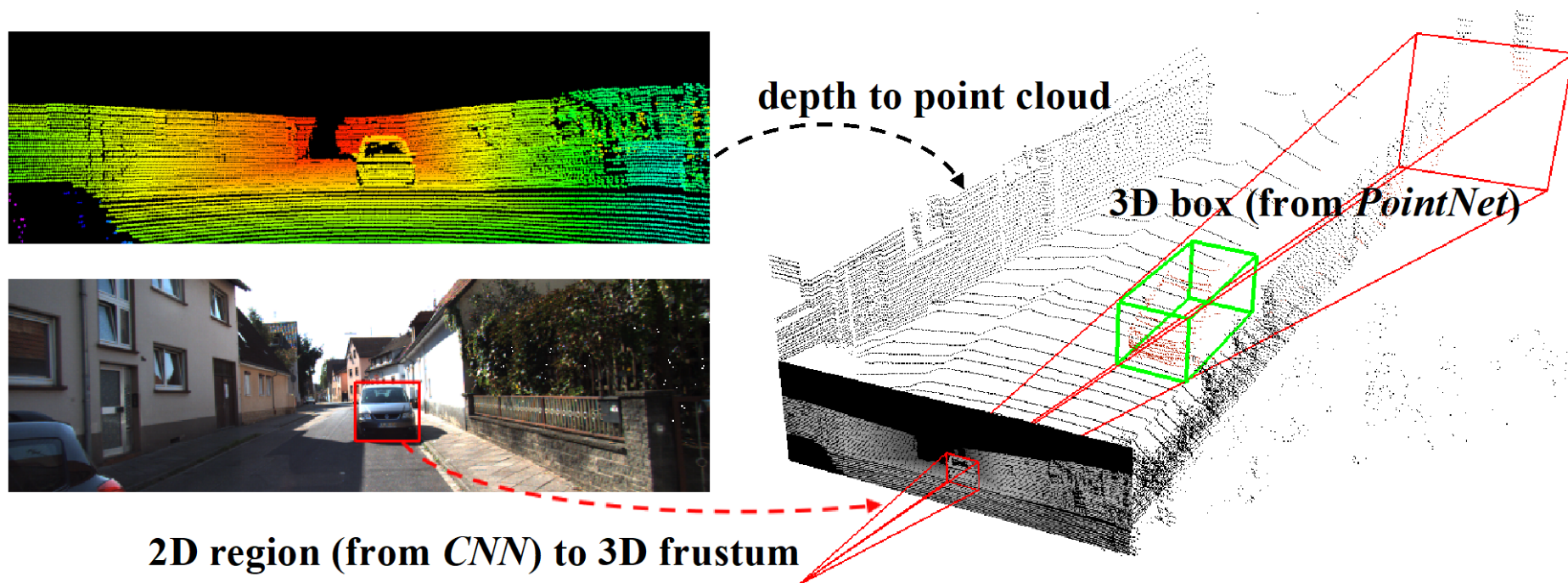


12 Edge  
centers



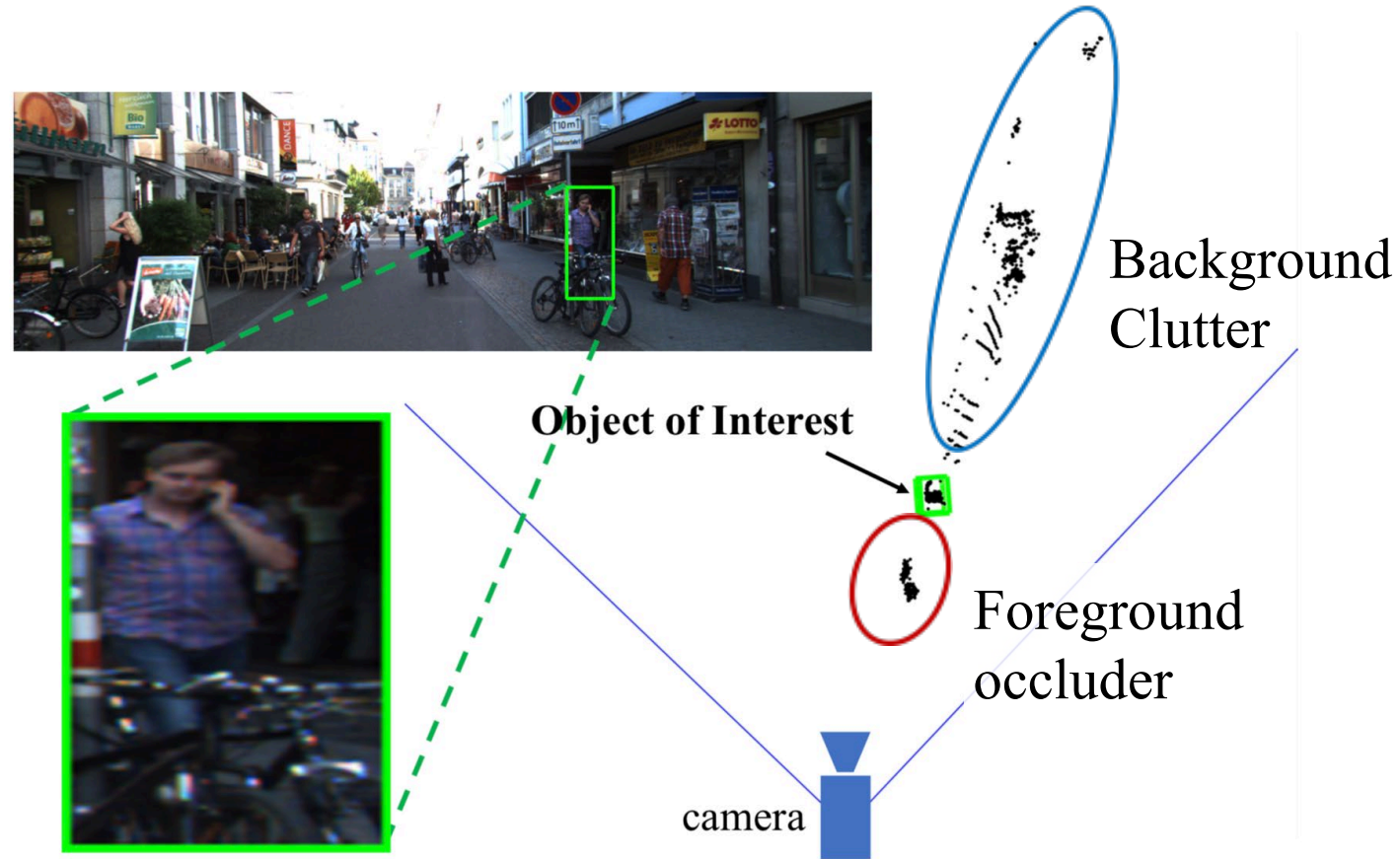
# Object Detection in Point Clouds, Outdoors

# Frustum PointNets for 3D Object Detection



- + **Leveraging mature 2D detectors** for region proposal. greatly reducing 3D search space.
- + Solving 3D detection problem with **3D data and 3D deep learning**.

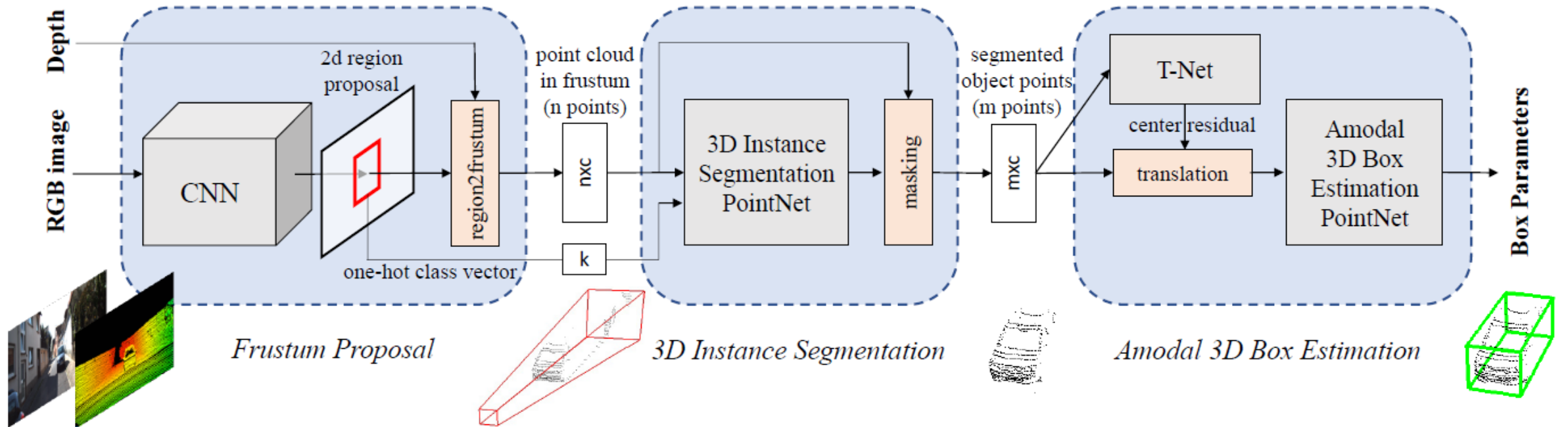
# Frustum-based 3D Object Detection: Challenges



- Occlusion and clutter is common in frustum point clouds
- Large range of point depths

# Frustum PointNets

Use **PointNets** for **data-driven** object detection in frustums.



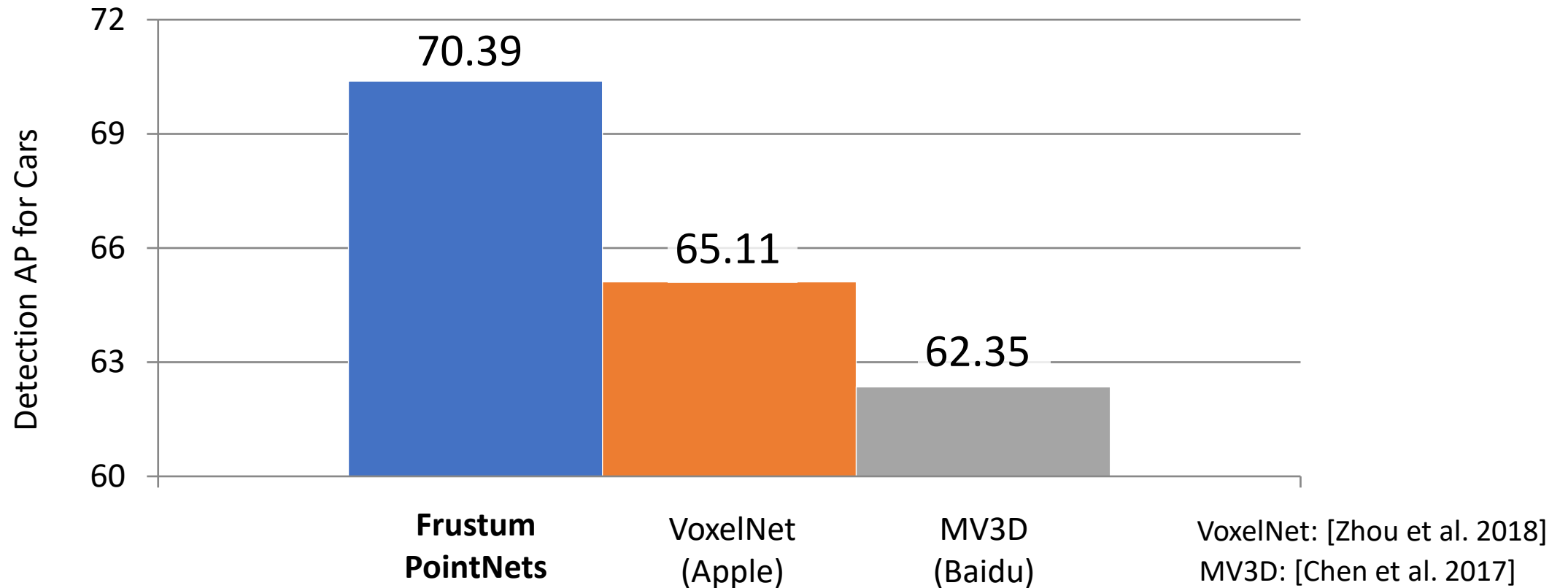
# Frustum PointNets: Key to Success

## *Respect and exploit 3D*

- **Use each modality (image, points) for what it's best at** — using 3D representation and 3D deep learning for the 3D problem.
- **Canonicalize the problem** — exploiting geometric transformations in point clouds.

# KITTI Results: Quantitative

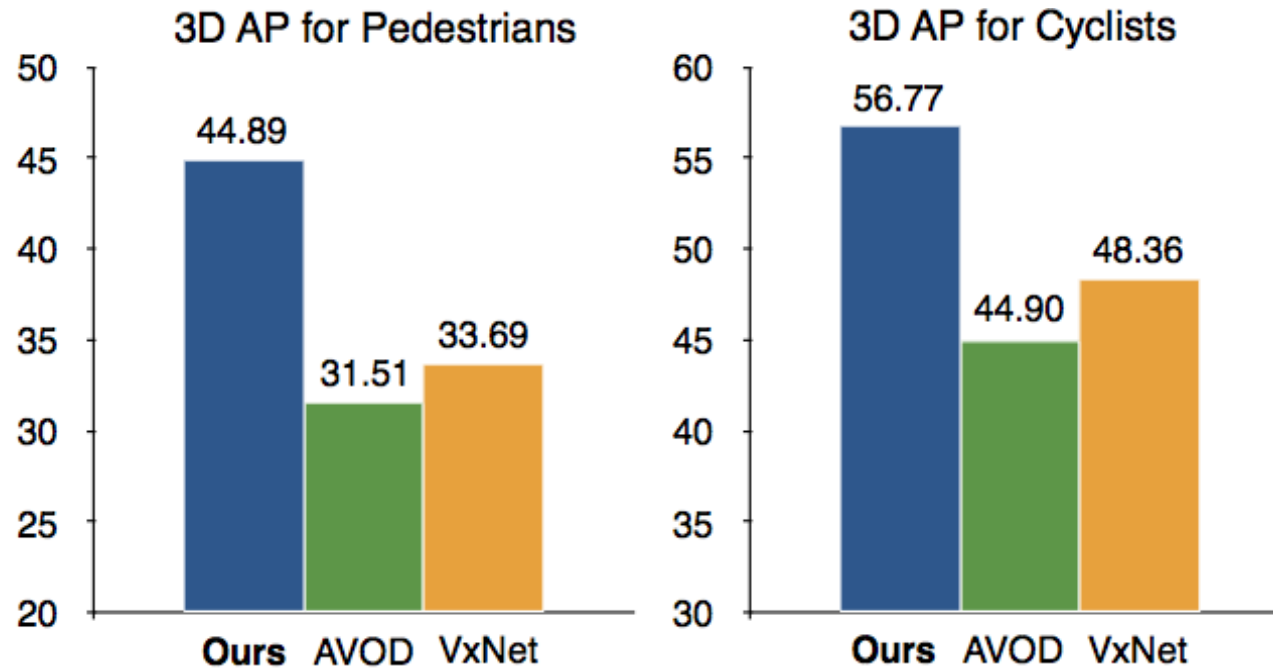
*Leading performance on KITTI benchmark*



# KITTI Results: Quantitative

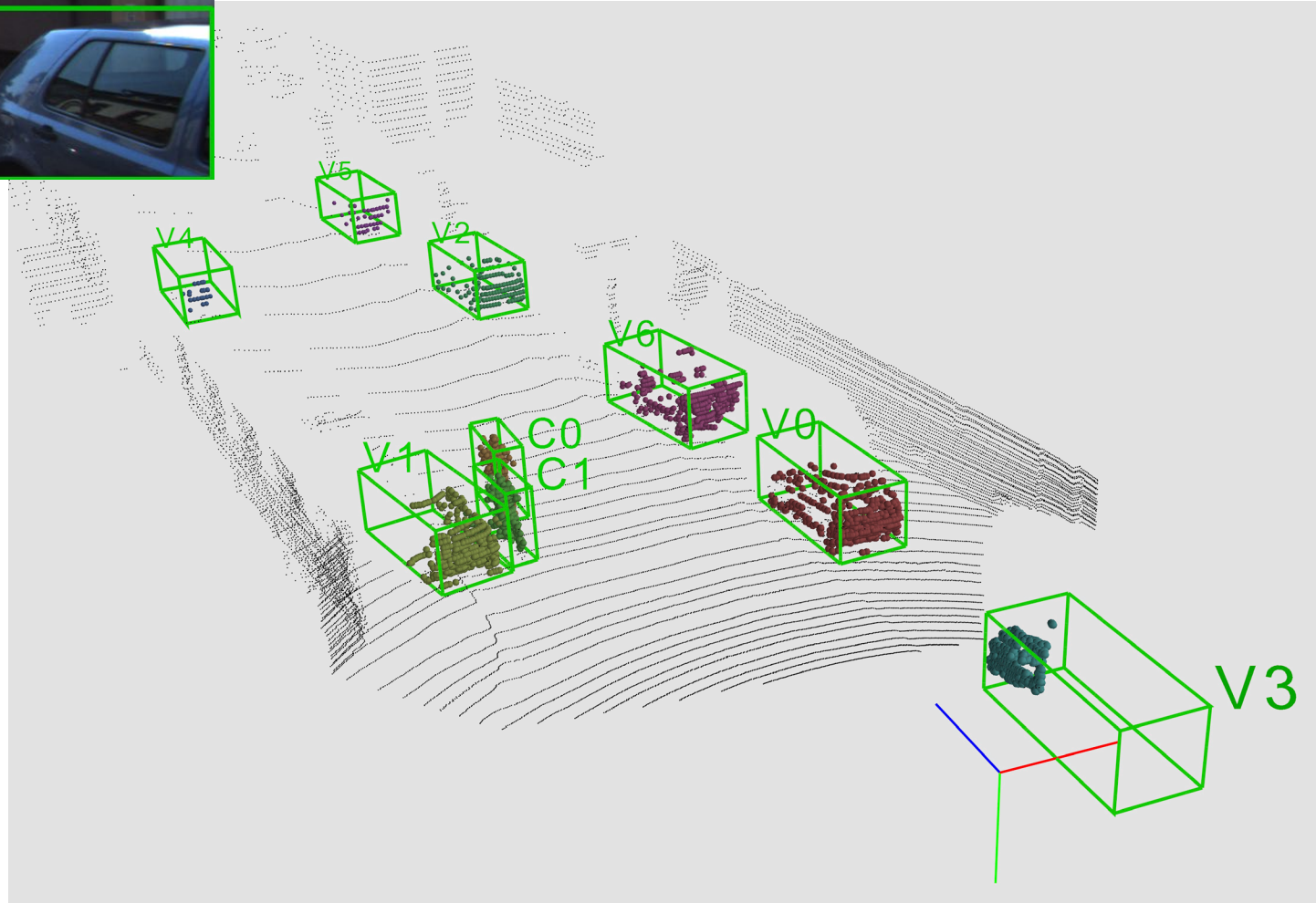
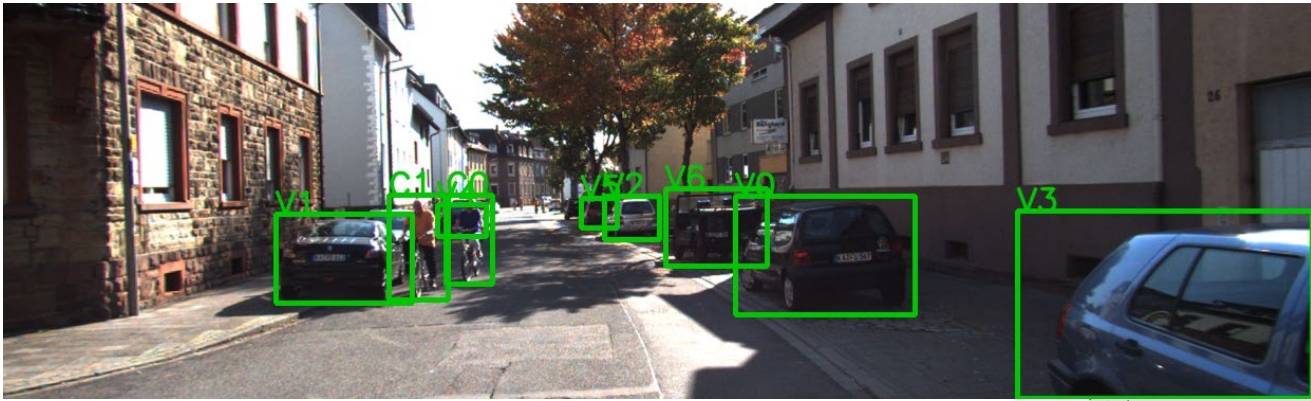
*Leading performance on KITTI benchmark*

Especially leading at smaller objects (pedestrians and cyclists)  
– hard to localize with 3D proposals only.



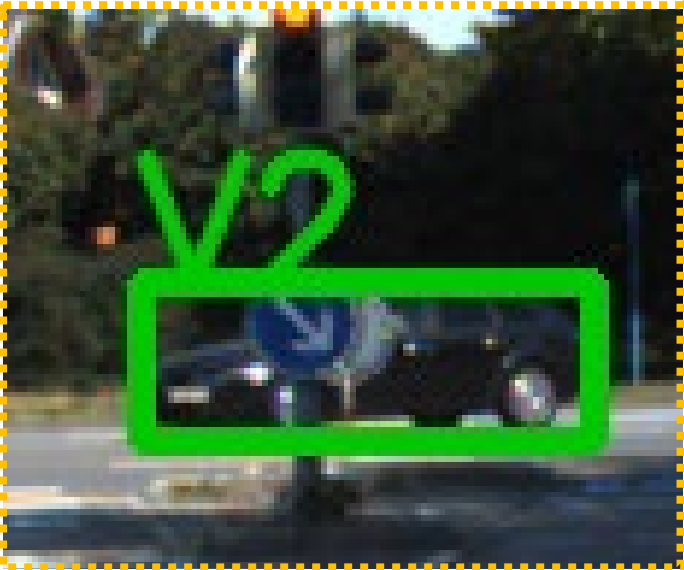
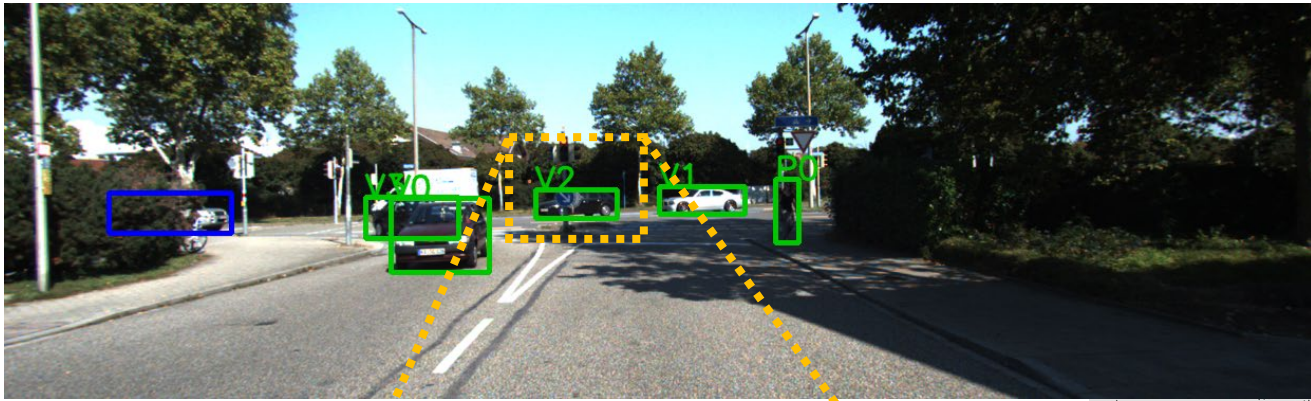
AVOD: [Ku et al. 2018]  
VxNet: [Zhou et al. 2017]

# KITTI Results: Qualitative

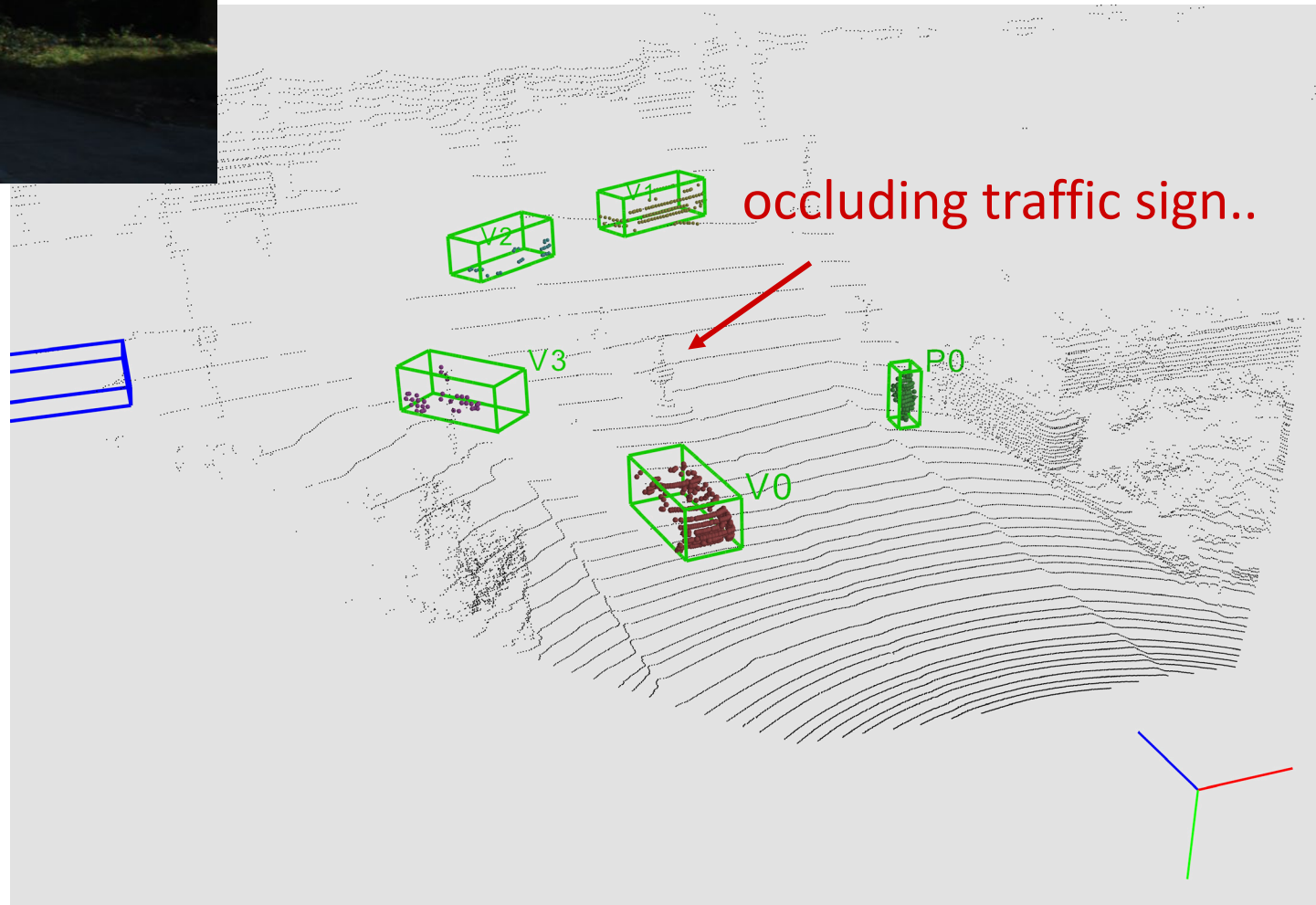


Remarkable box estimation accuracy even with a dozen of points or with very partial point clouds.

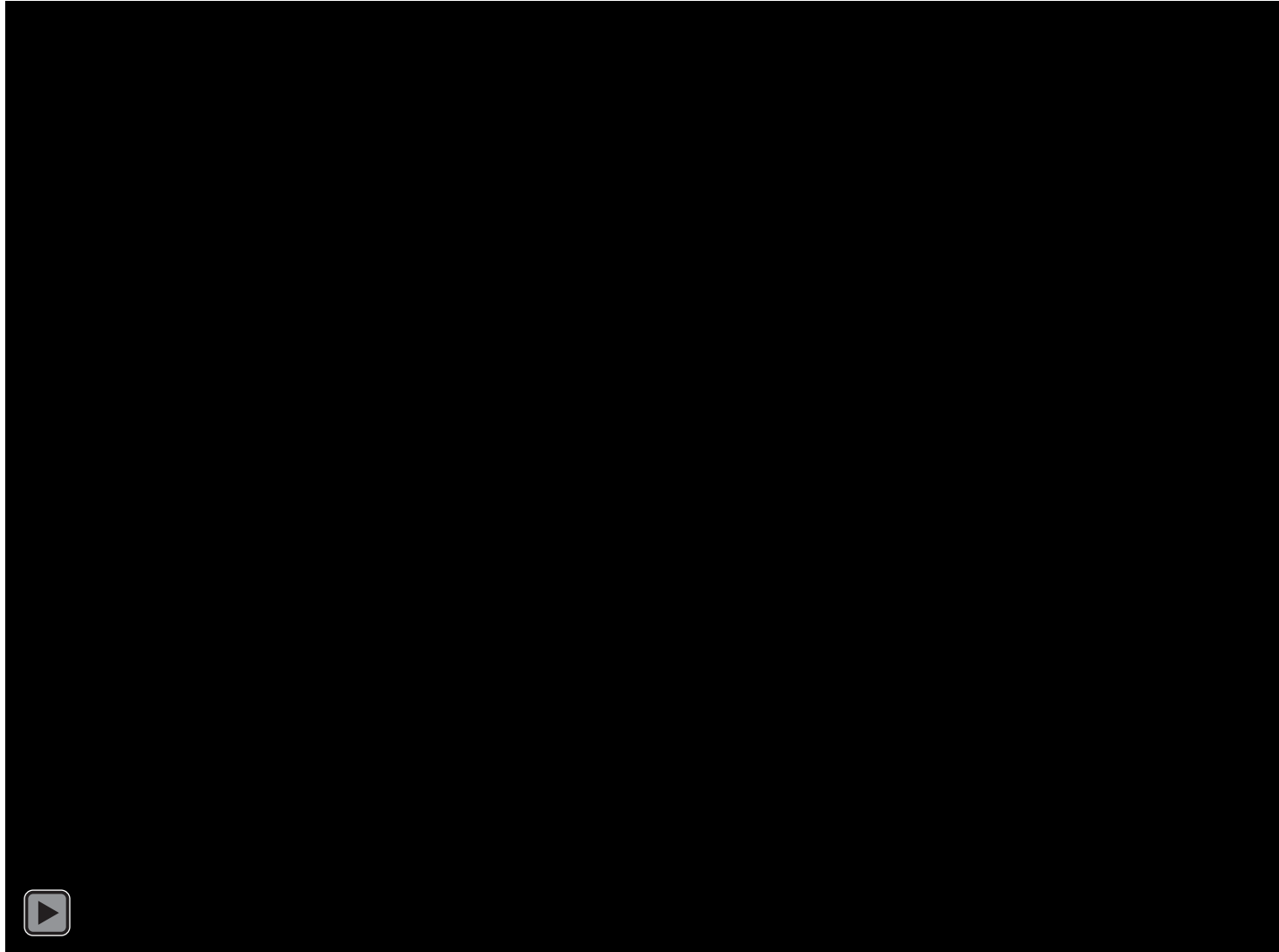
# KITTI Results: Qualitative



Correct segmentation in point clouds with heavy occlusion.



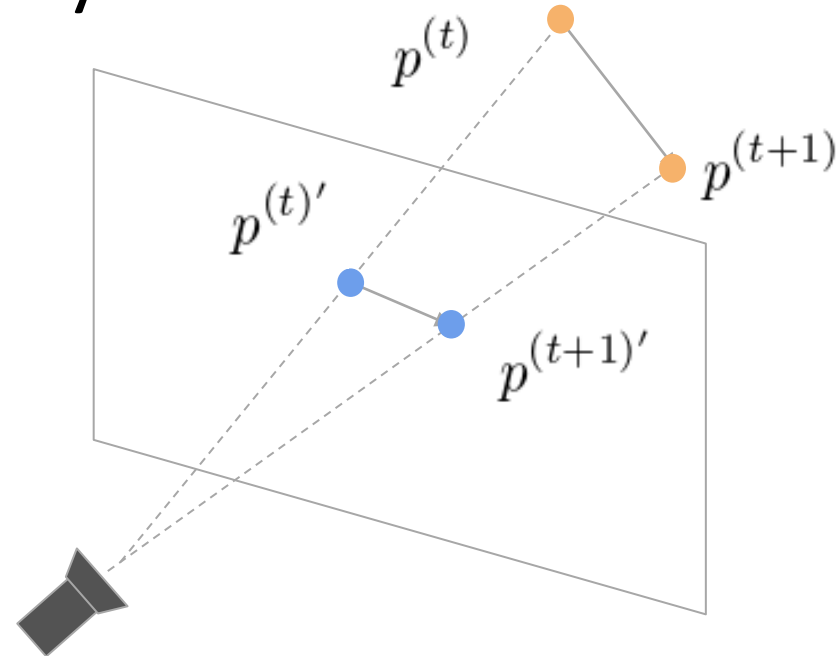
# KITTI Results: Example



# 3D Motion in Point Clouds

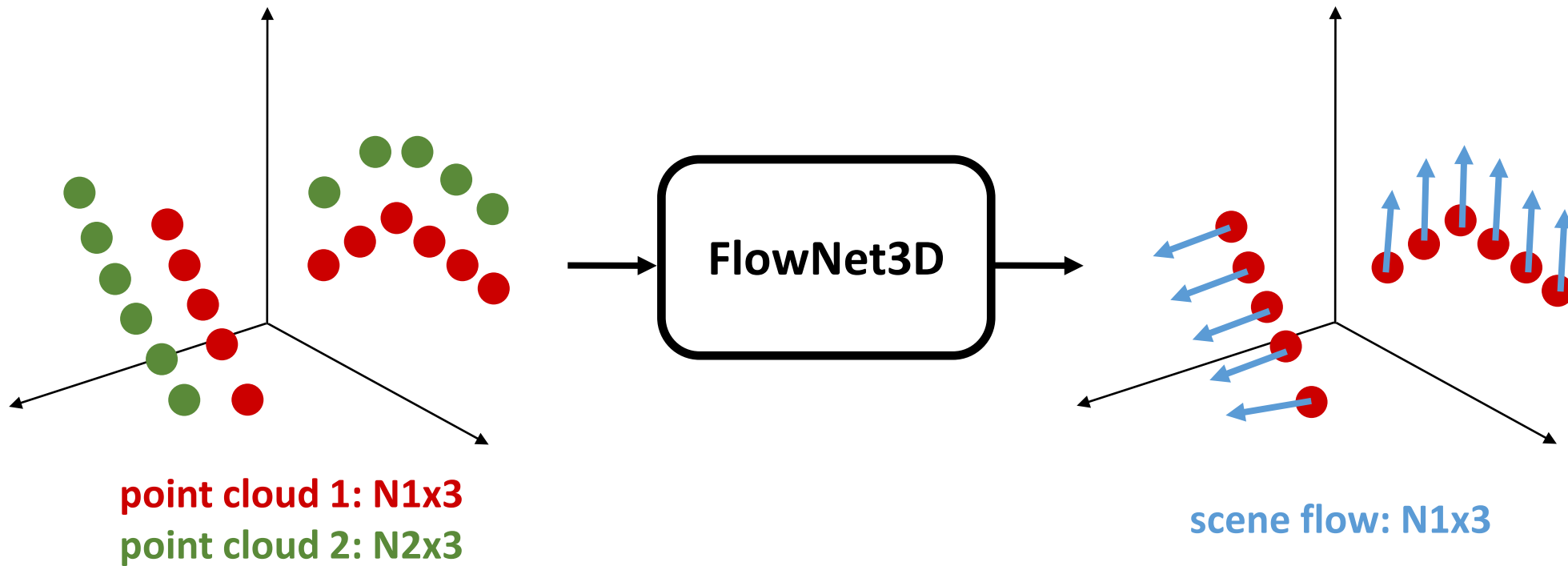
# Scene Flow [Vedula et al. 1999]

- Scene flow: 3D motion field of points
- Optical flow is its projection to 2D image plane.
- Low-level understanding of a dynamic environment



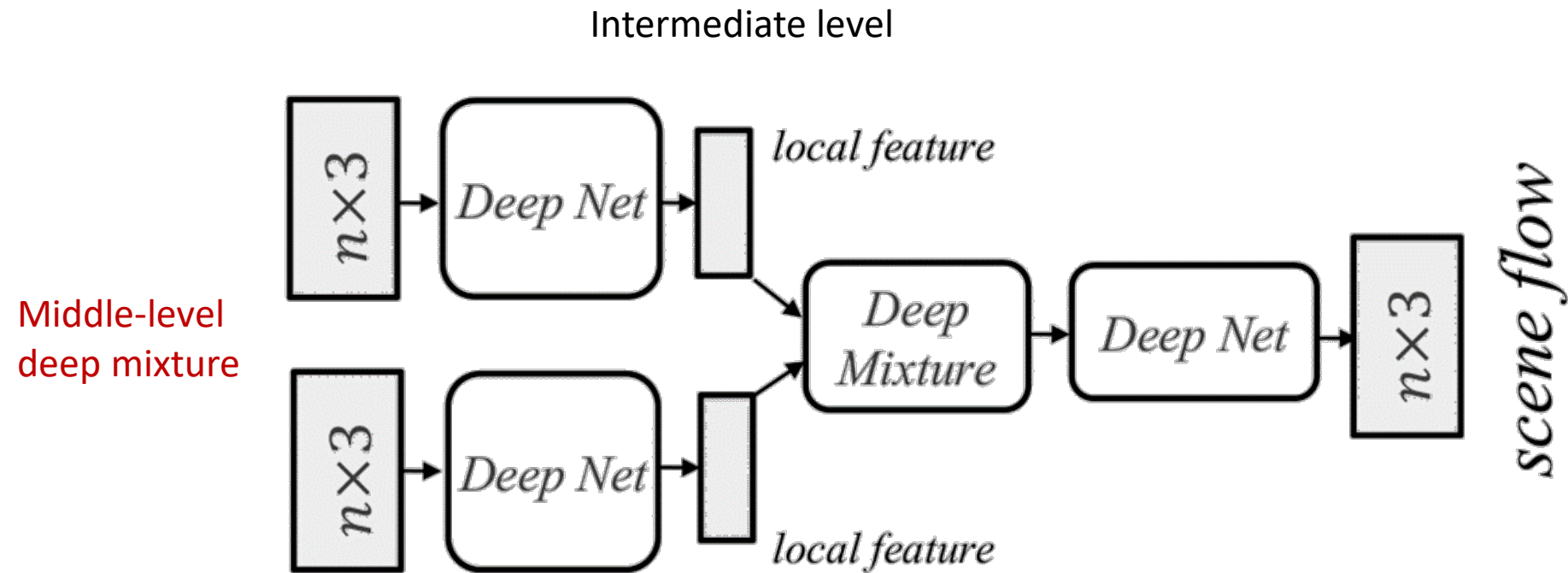
# Our Approach: FlowNet3D

- Directly learning scene flow in 3D point clouds, with 3D deep learning architectures.

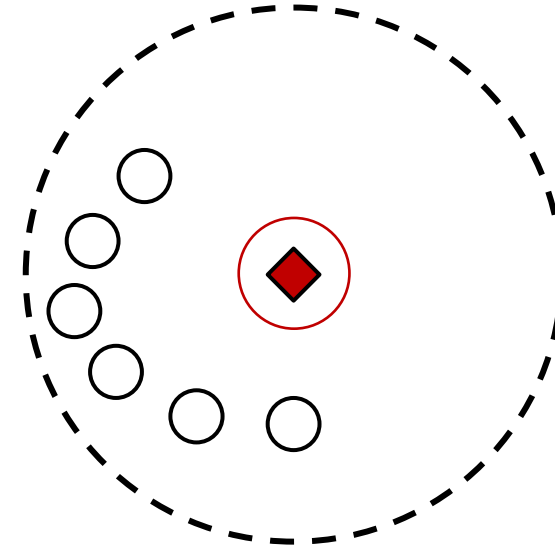
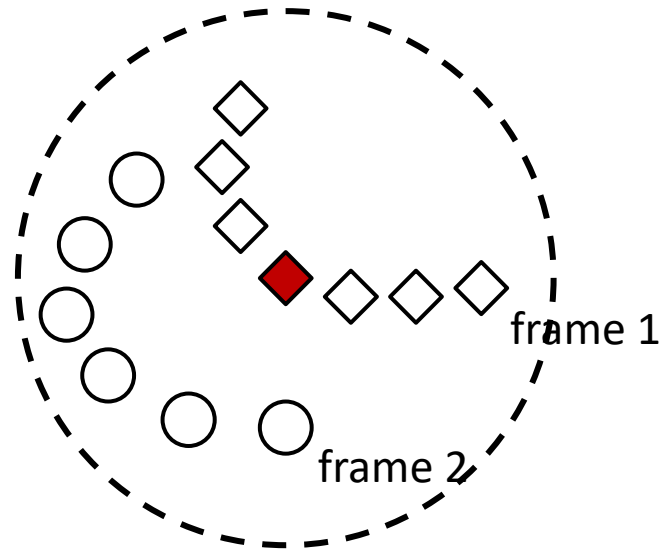


# Deep Net Architecture

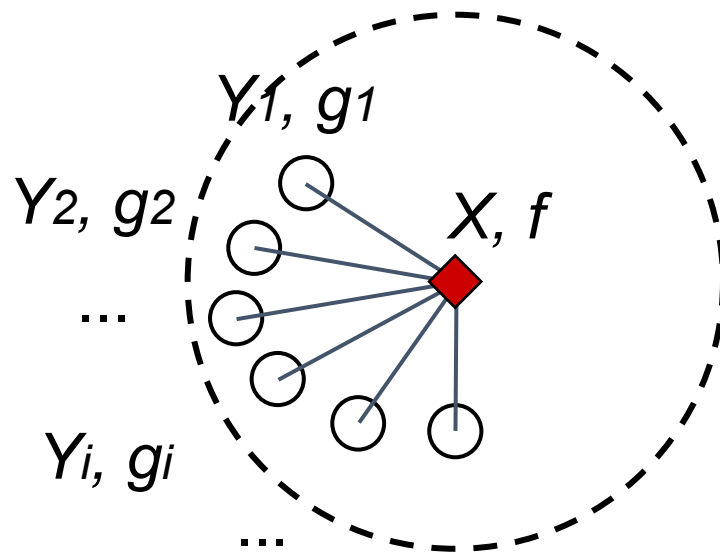
- How to learn point cloud features?
- Where in the network architecture to mix point features from consecutive frames?
- How to mix them?



# Middle-Level Mixing



# Point Attributes

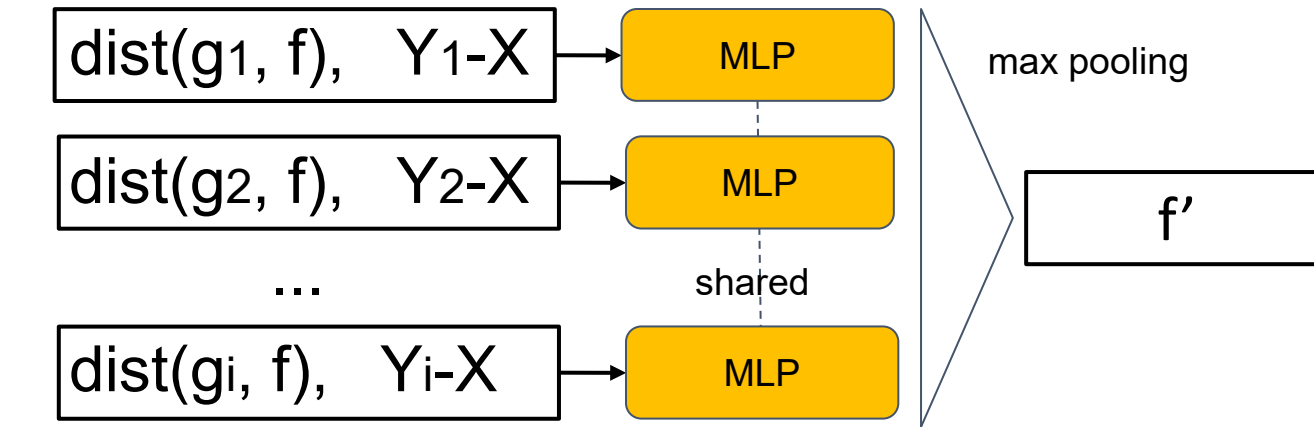


$\text{dist}(g_1, f), Y_1-X$   
 $\text{dist}(g_2, f), Y_2-X$   
 $\vdots$   
 $\text{dist}(g_i, f), Y_i-X$   
 $\vdots$

*Naive approach: concatenation*

$\text{dist}(g_1, f), Y_1-X$	$\text{dist}(g_2, f), Y_2-X$	$\dots$
------------------------------	------------------------------	---------

# A More Structured Approach



$\text{dist}(g_i, f)$

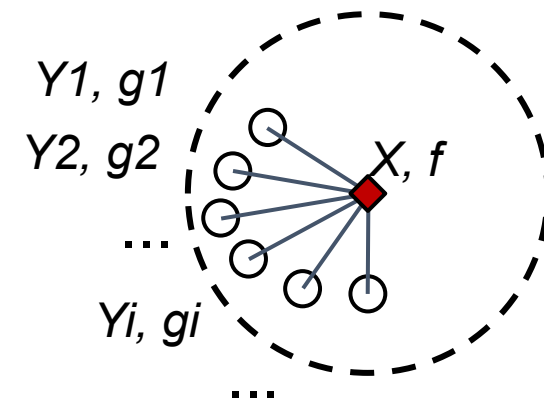
“Distance” functions:

Euclidean distance (scalar)

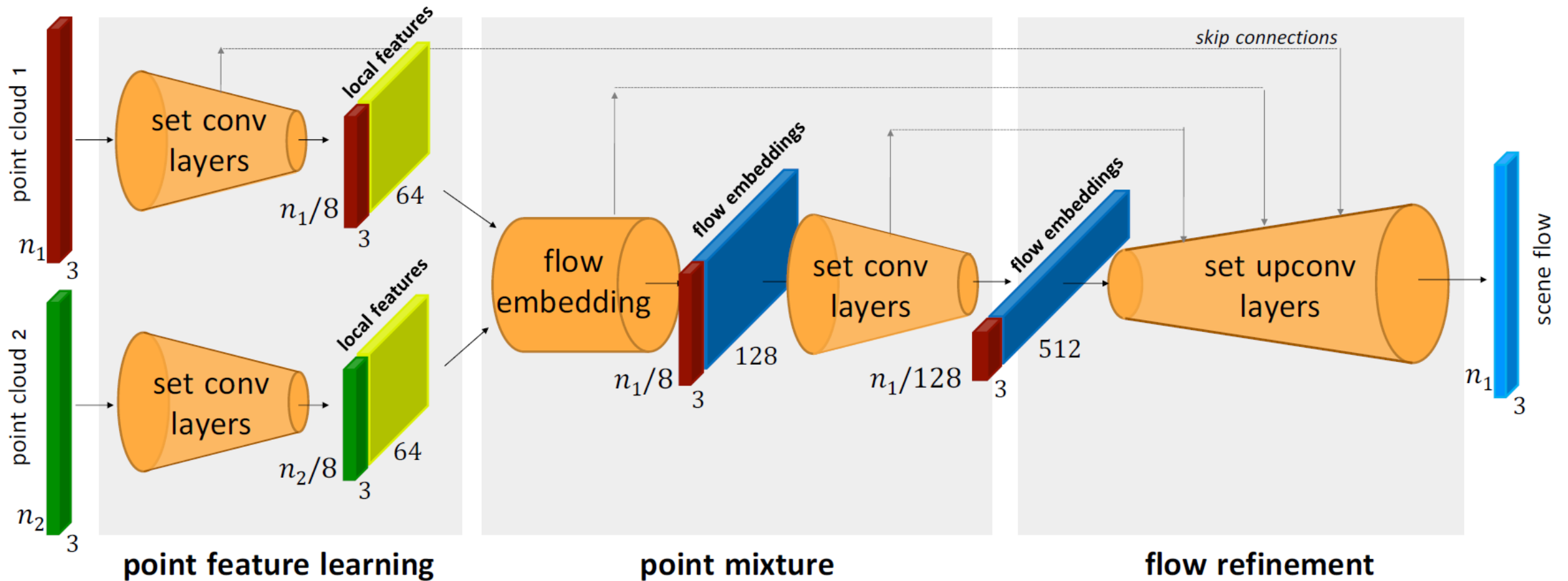
Cosine distance (scalar)

Element-wise product (vector)

Let the network learn the distance function ...



# FlowNet3D



set conv = set abstraction

Composed of many many mini-pointnet++ modules ...

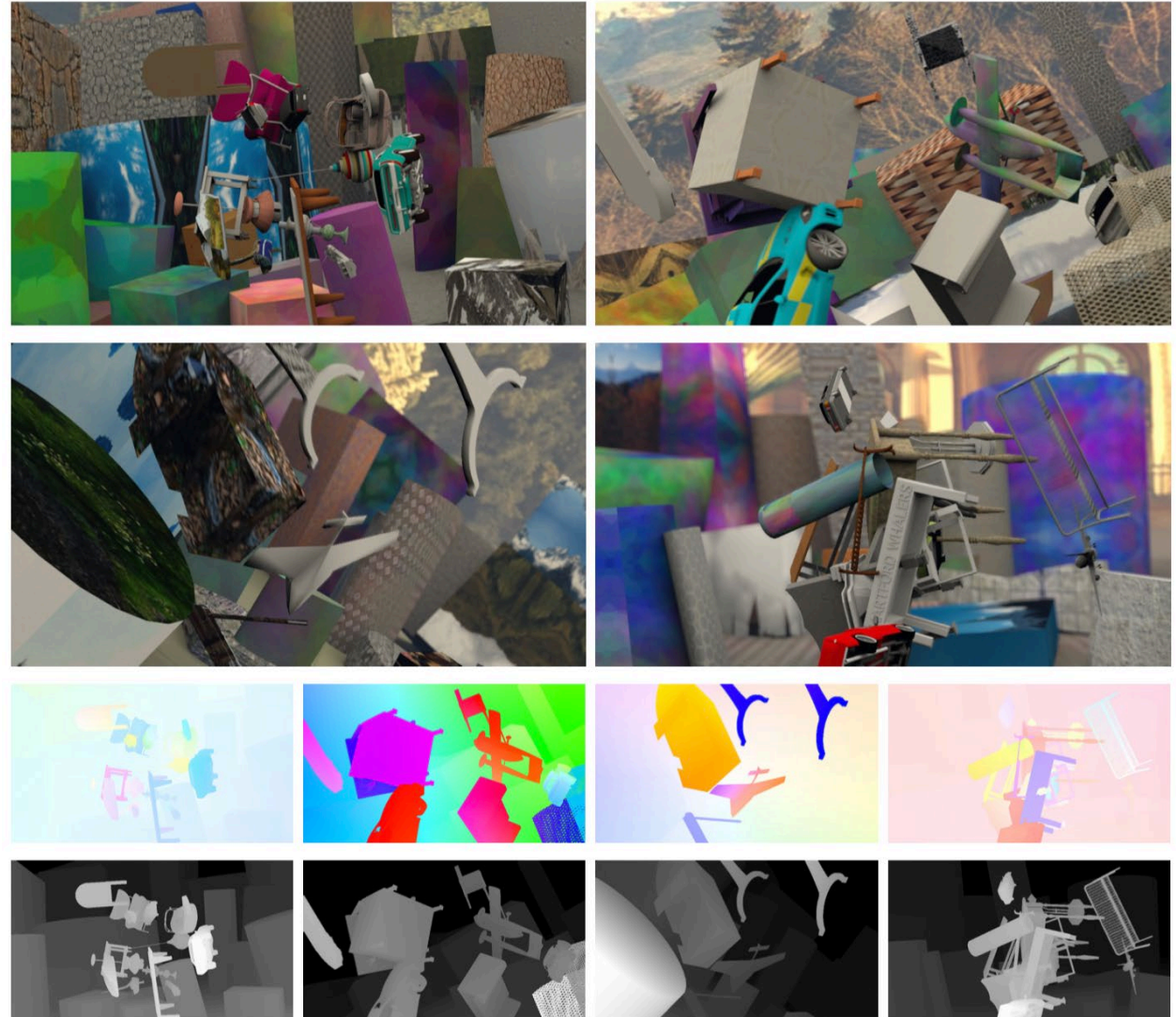
Pointnet++

# Training on Synthetic Data

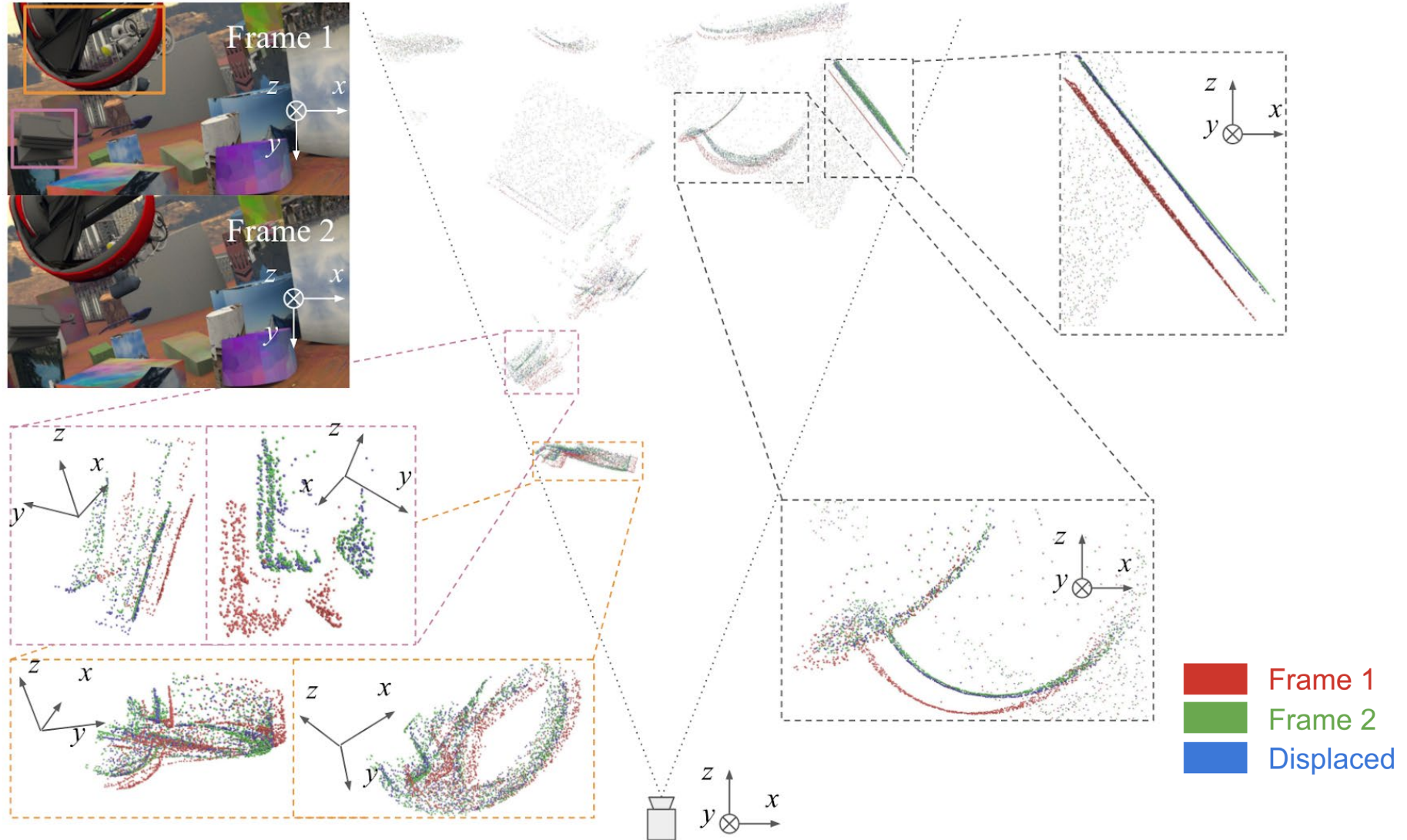
FlyingThings3D [Mayer et al. 2016]  
dataset from MPI

Random ShapeNet objects

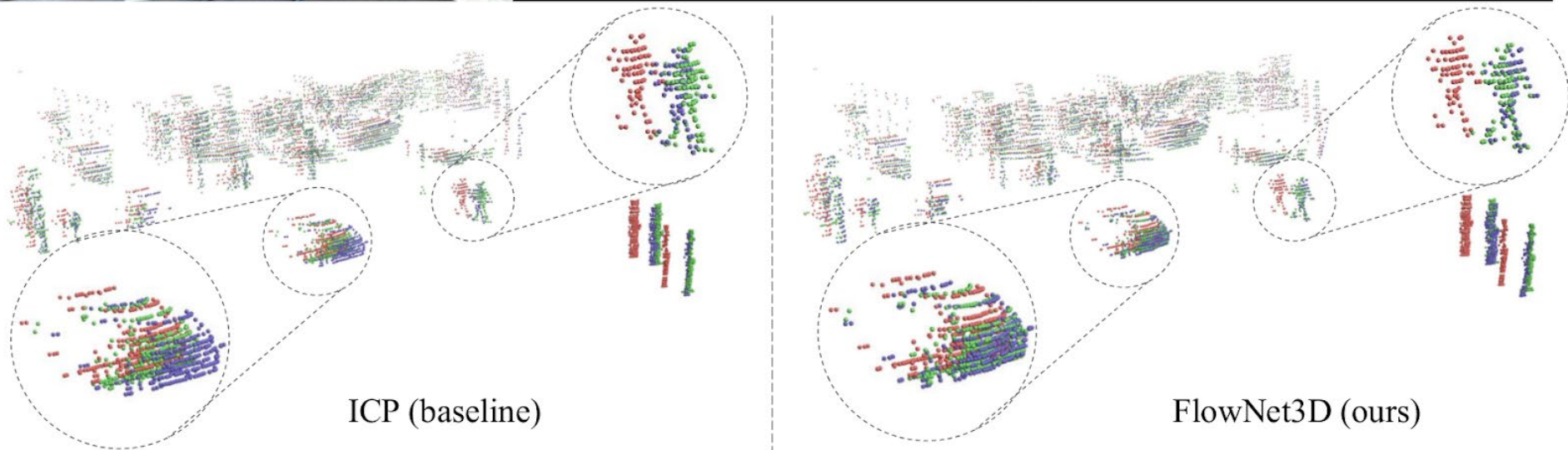
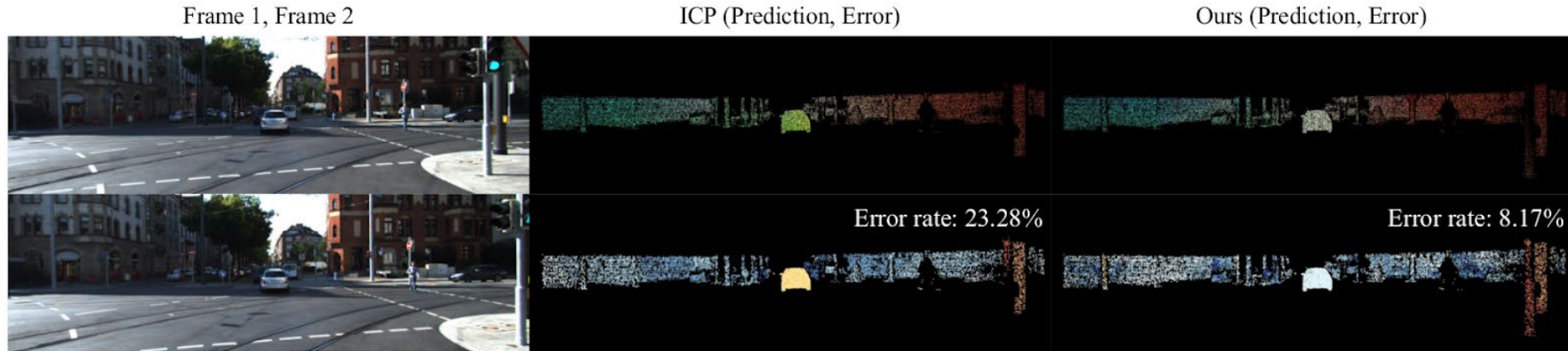
Very challenging dataset with  
strong occlusions and large motions.



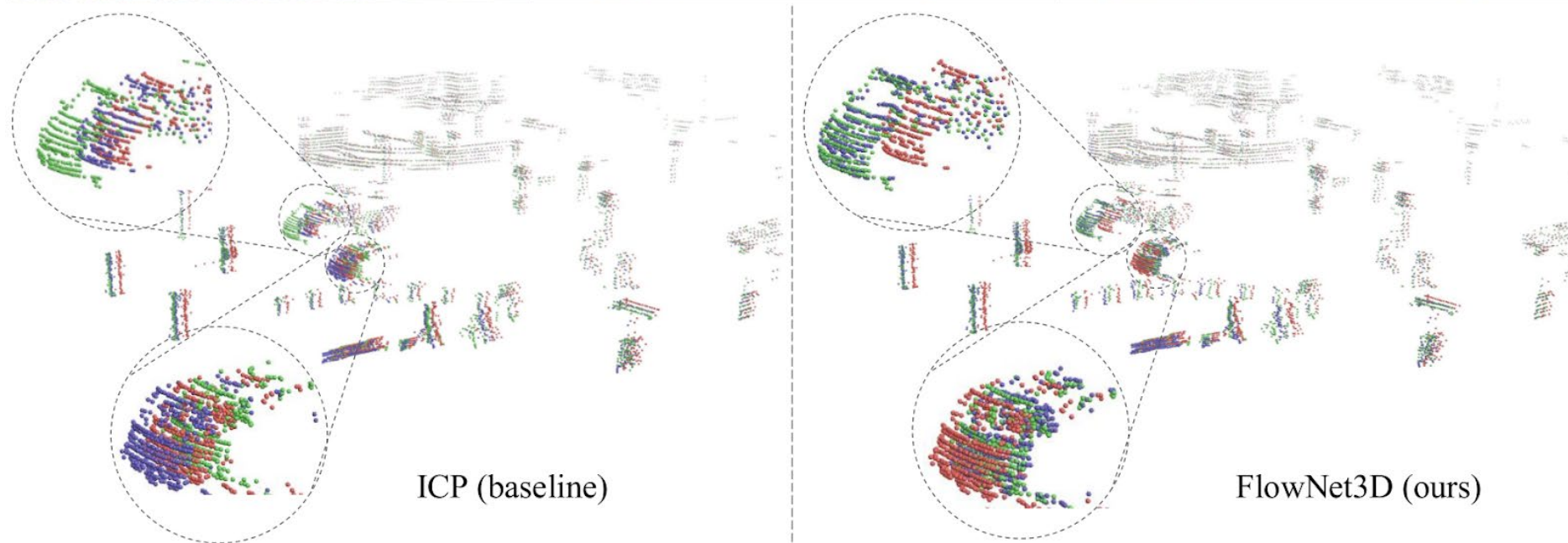
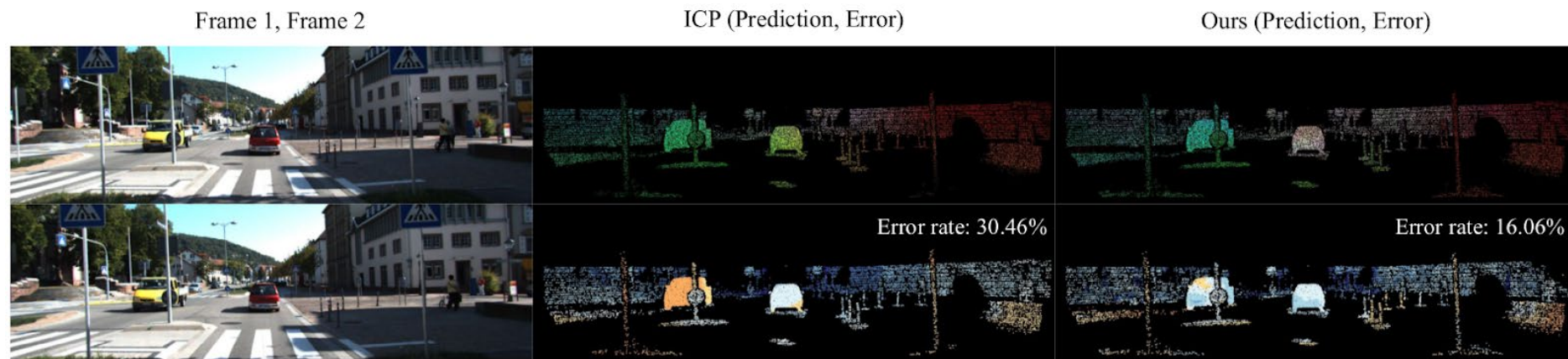
# FlyingThings3D Results



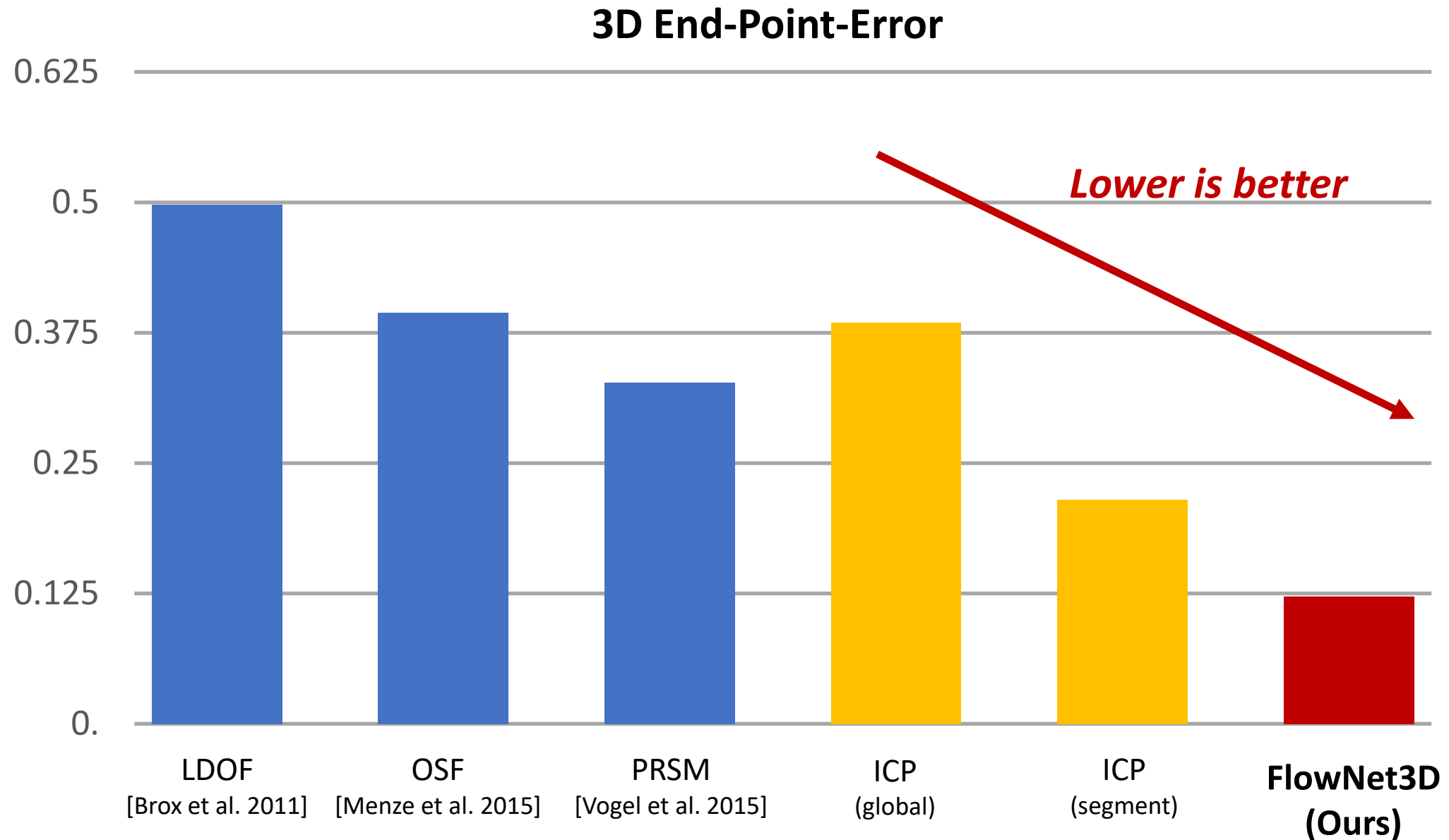
# KITTI Results



# KITTI Results

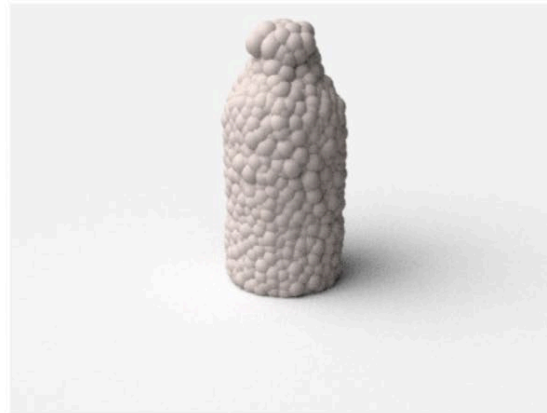
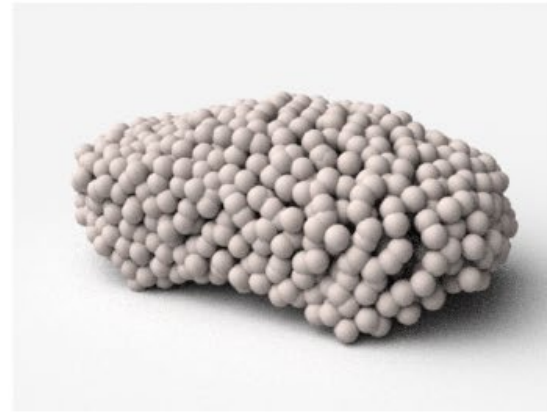
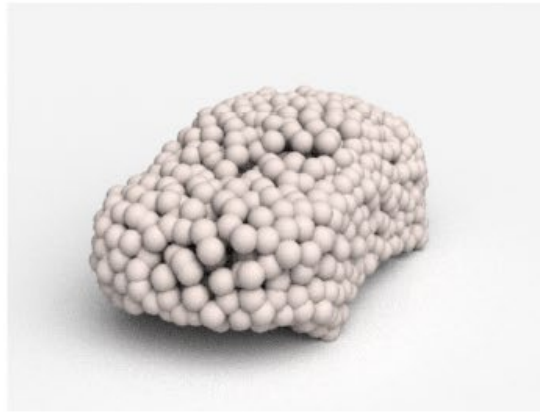


# Generalizing to KITTI: Quantitative



# Point-Set Generation

# Point Cloud Synthesis from a Single Image

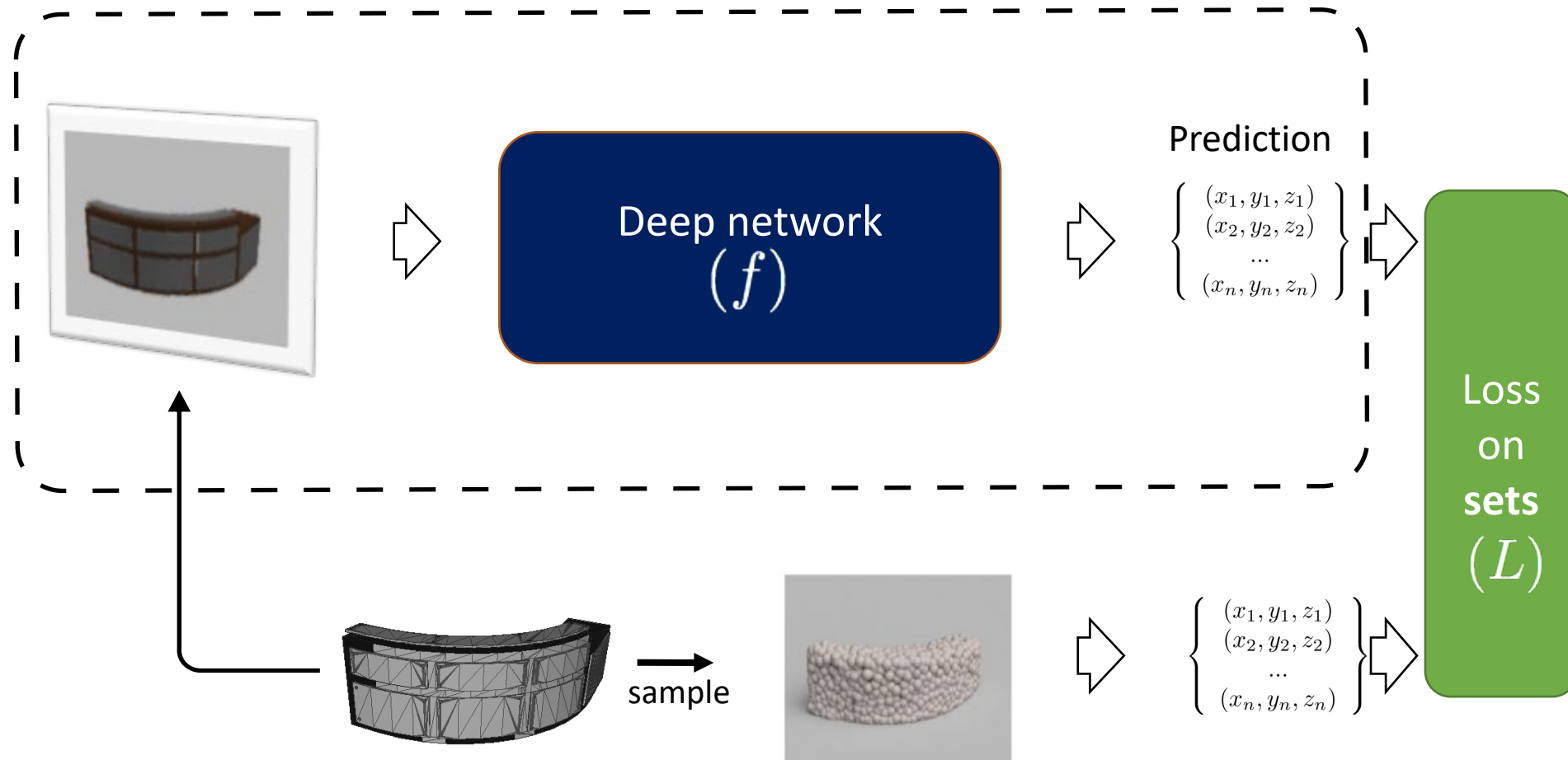


Input

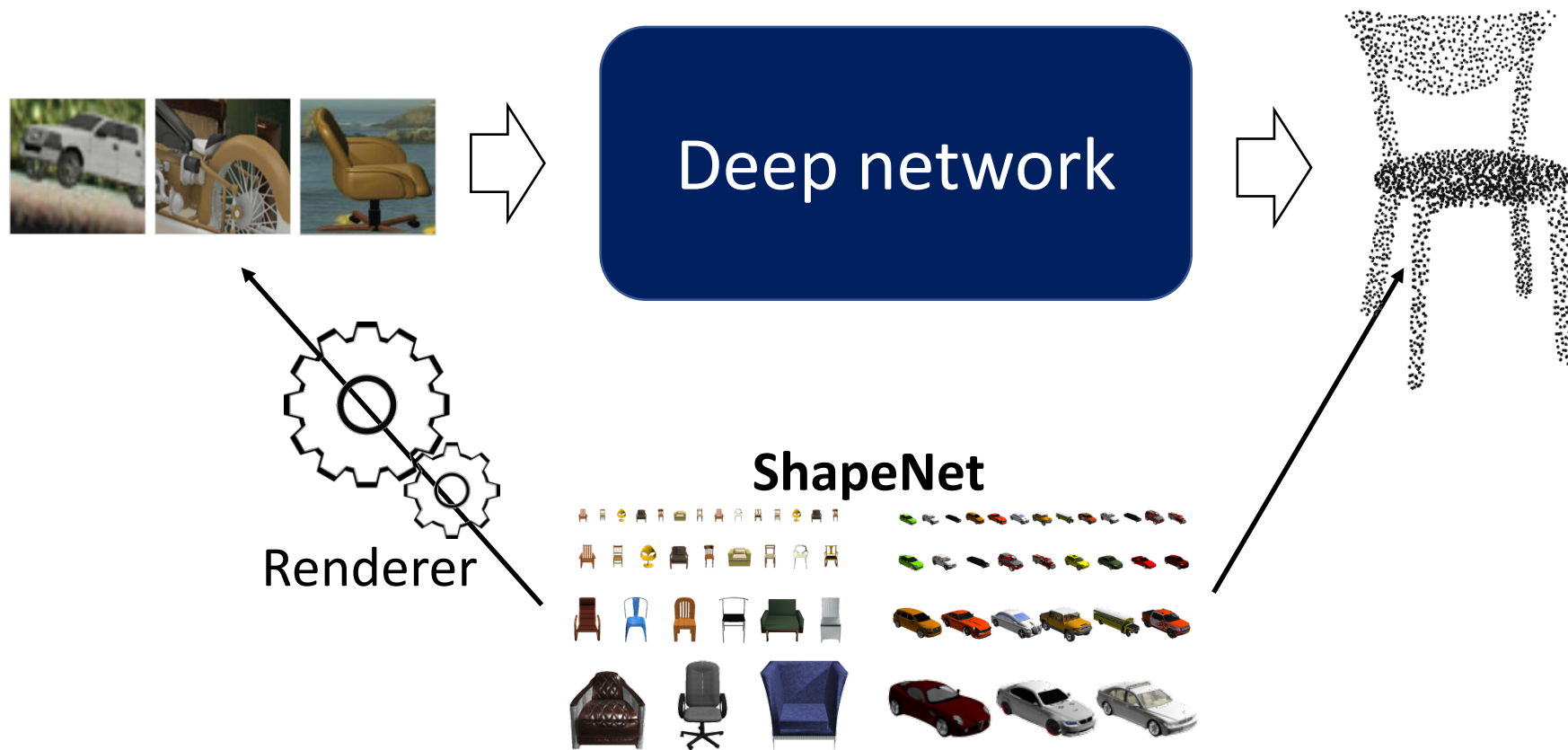
Reconstructed 3D point cloud

[H. Su, H. Fan, LG, 2017]

# End-to-End Learning

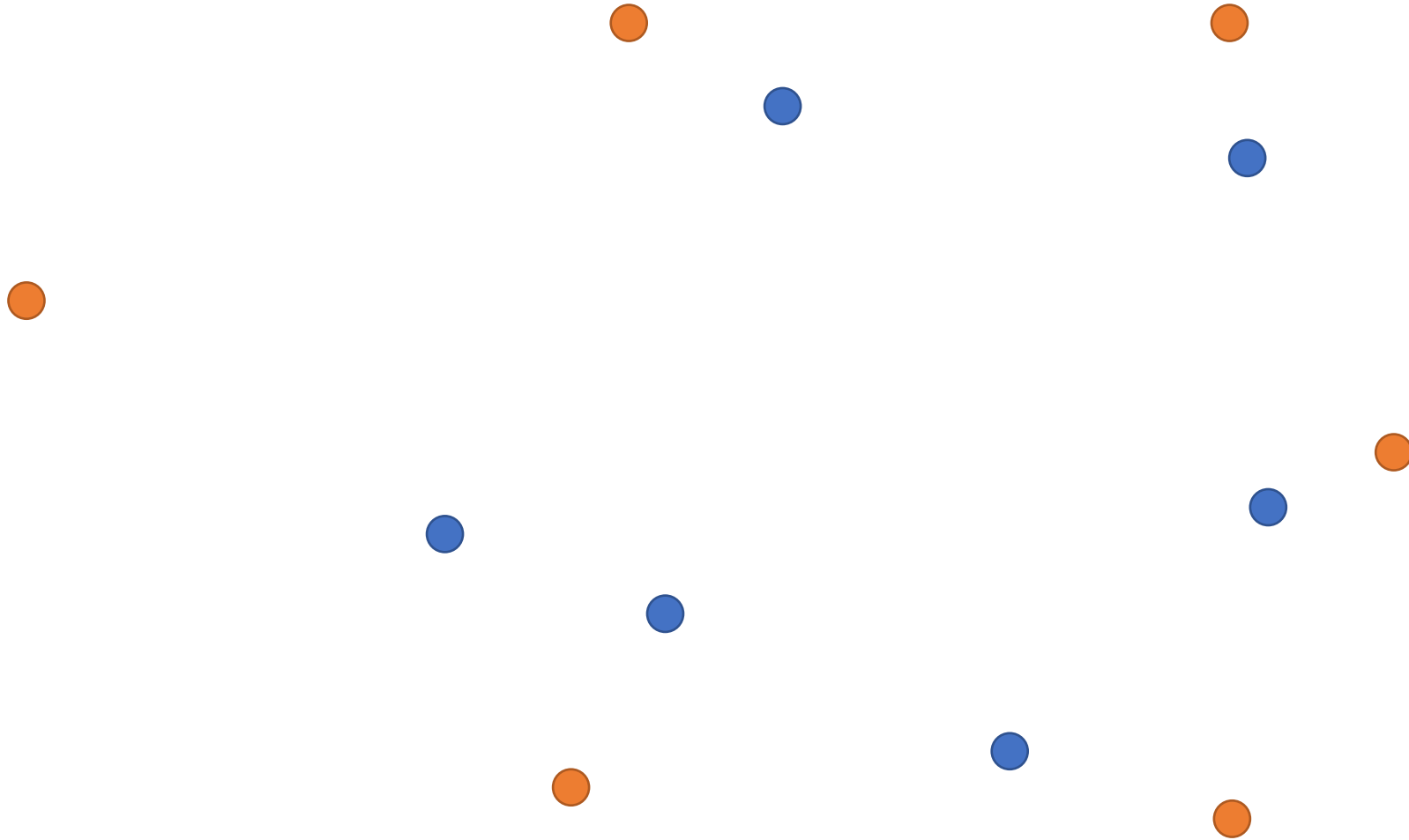


# Synthesize for Learning



# Distance Metrics Between Point Sets

Given two sets of points, measure their discrepancy

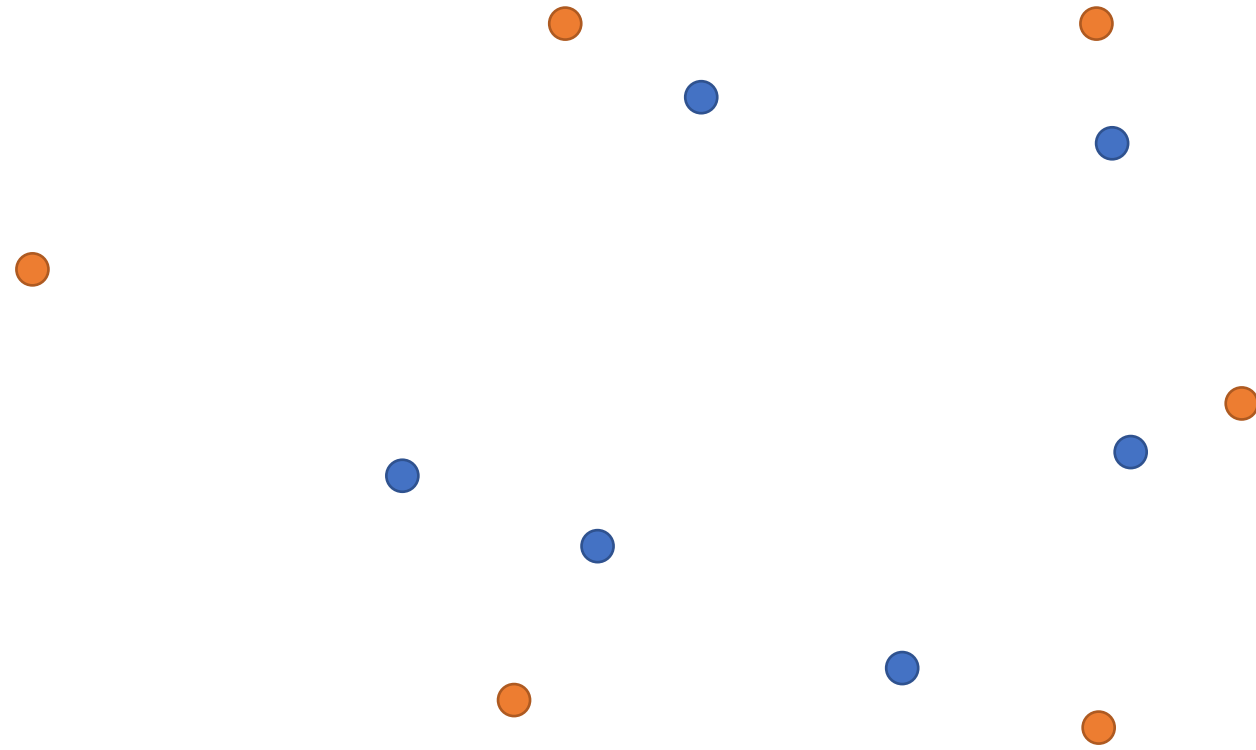


# Common Distance Metrics

Worst case: Hausdorff distance (HD)

Average case: Chamfer distance (CD)

Optimal case: Earth Mover's distance (EMD)

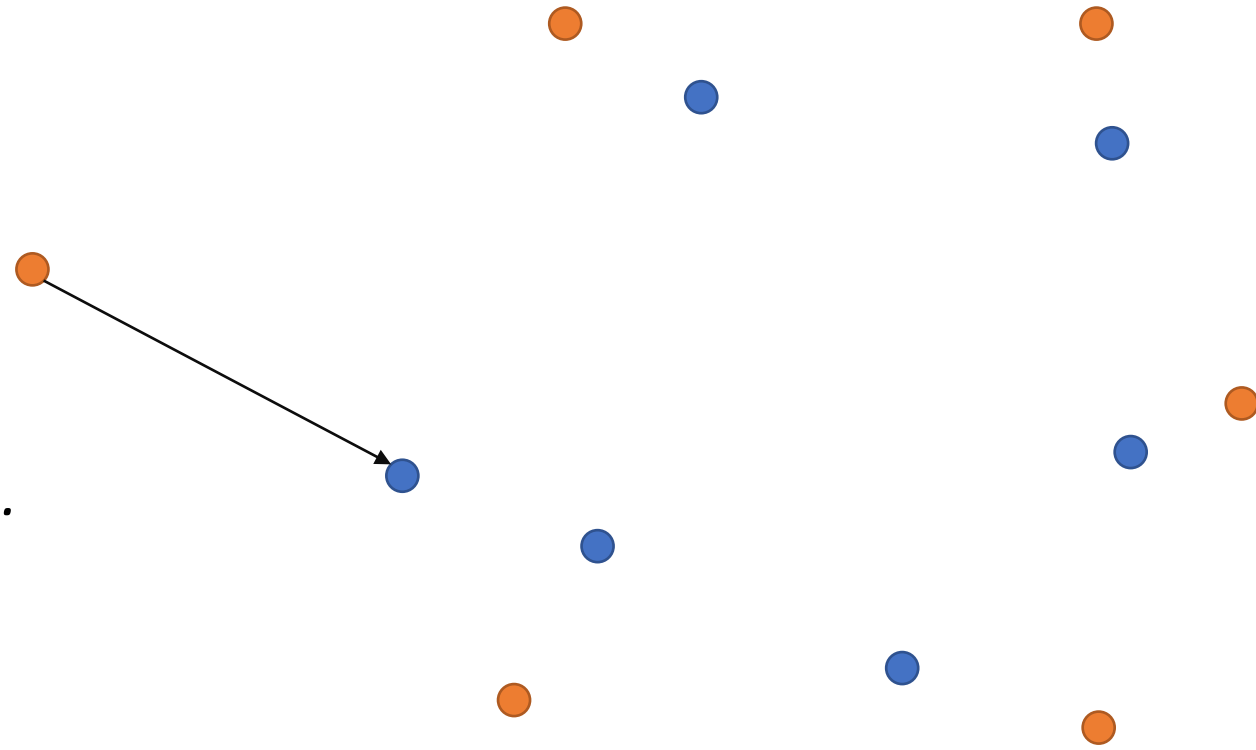


# Common Distance Metrics

$$d_{\text{HD}}(S_1, S_2) = \max\left\{ \max_{x_i \in S_1} \min_{y_j \in S_2} \|x_i - y_j\|, \max_{y_j \in S_2} \min_{x_i \in S_1} \|x_i - y_j\| \right\}$$

Worst case: Hausdorff distance (HD)

*A single distant pair determines the distance.  
In other words, **not robust to outliers!***



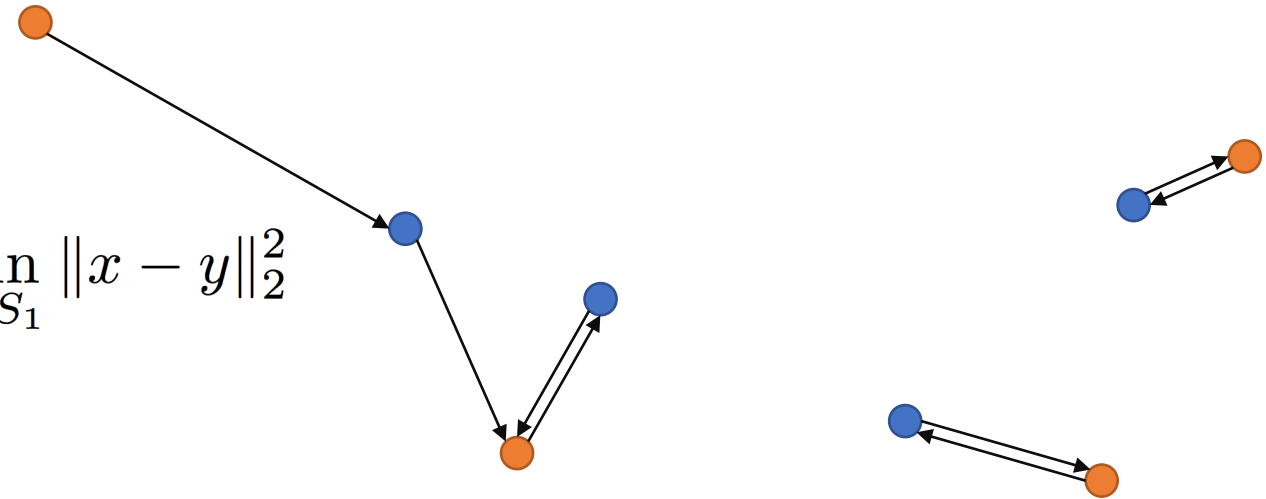
# Common Distance Metrics

Worst case: Hausdorff distance (HD)



Average case: Chamfer distance (CD)

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



*Average all the nearest neighbor distance by nearest neighbors*

# Common Distance Metrics

Worst case: Hausdorff distance (HD)

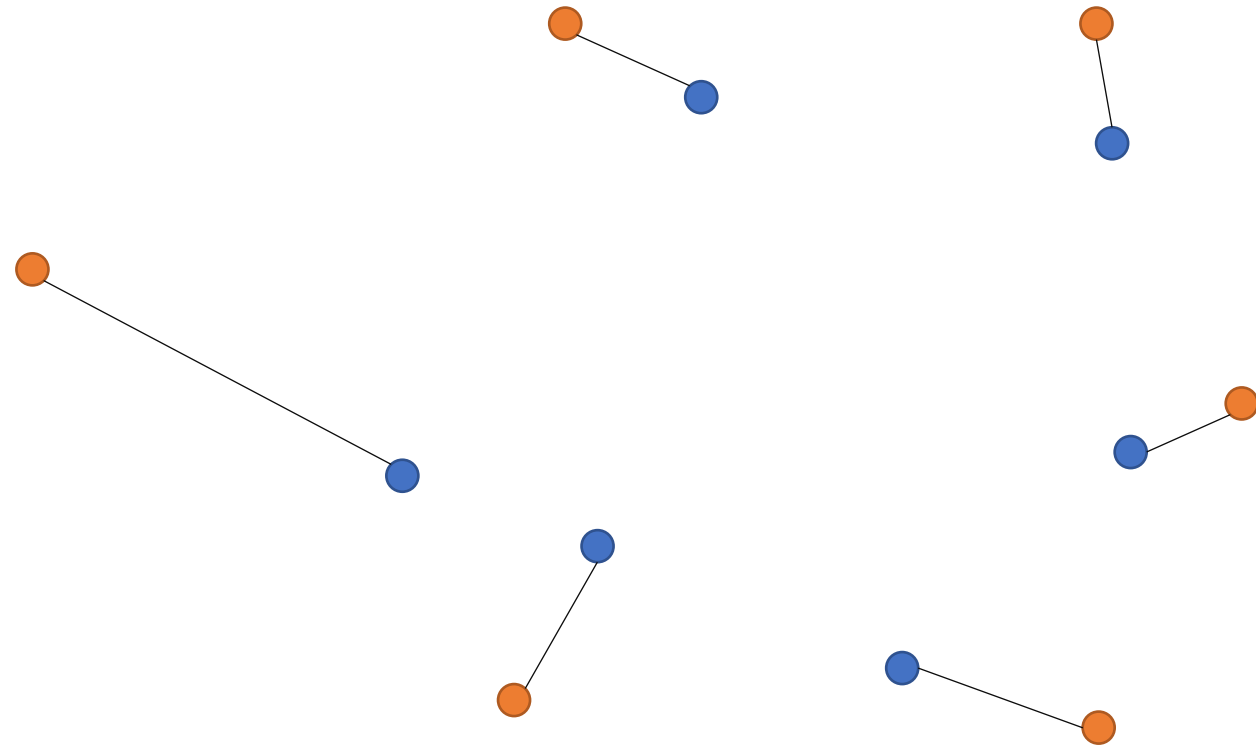
Average case: Chamfer distance (CD)

Optimal case: Earth Mover's distance (EMD)

$$d_{EMD}(S_1, S_2) = \min_{\phi: S_1 \rightarrow S_2} \sum_{x \in S_1} \|x - \phi(x)\|_2$$

where  $\phi : S_1 \rightarrow S_2$  is a bijection.

*Solves the optimal transportation (bipartite matching) problem!*



# Desired Properties of Distance Metrics

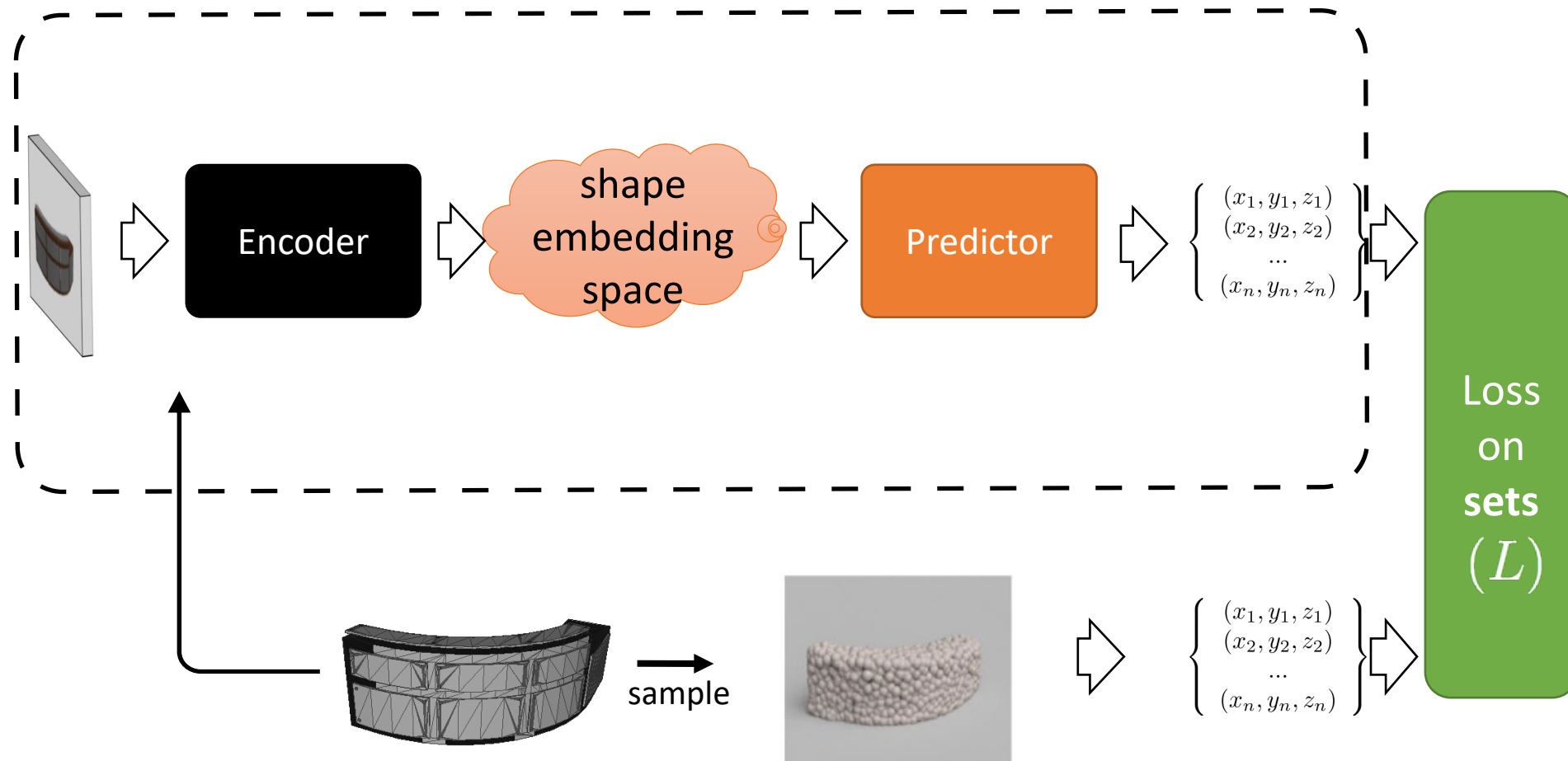
## Geometric requirement

- Induces a nice shape space
- In other words, a good metric should reflect the natural shape differences

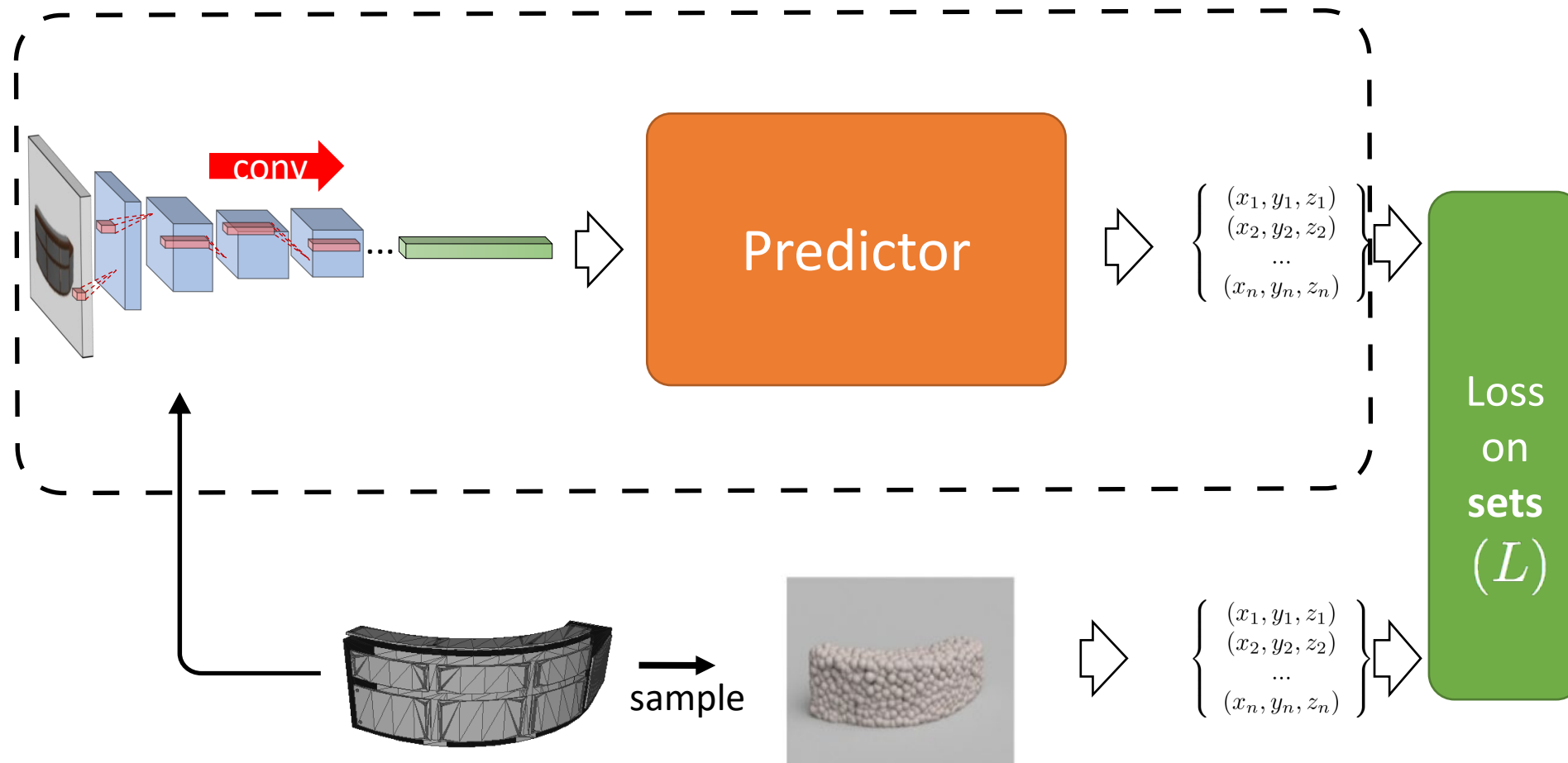
## Computational requirement

- Defines a loss that is numerically easy to compute and optimize

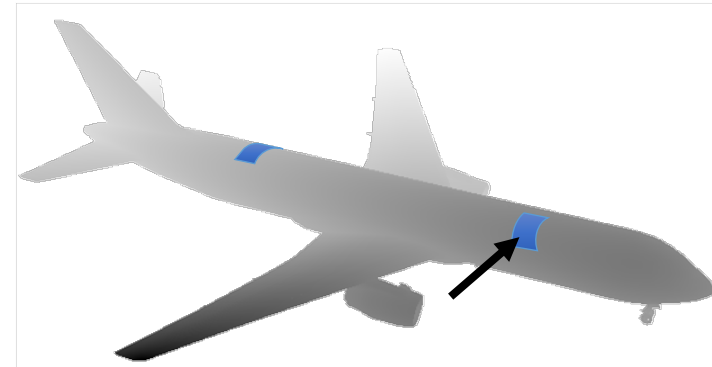
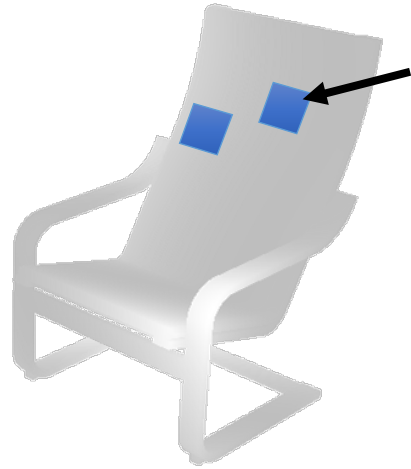
# End-to-End Learning



# End-to-End Learning



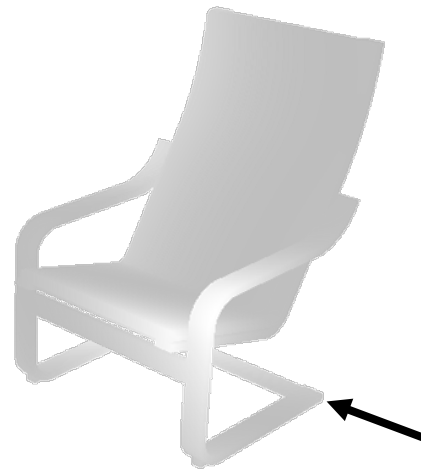
# Natural Statistics of Object Geometry



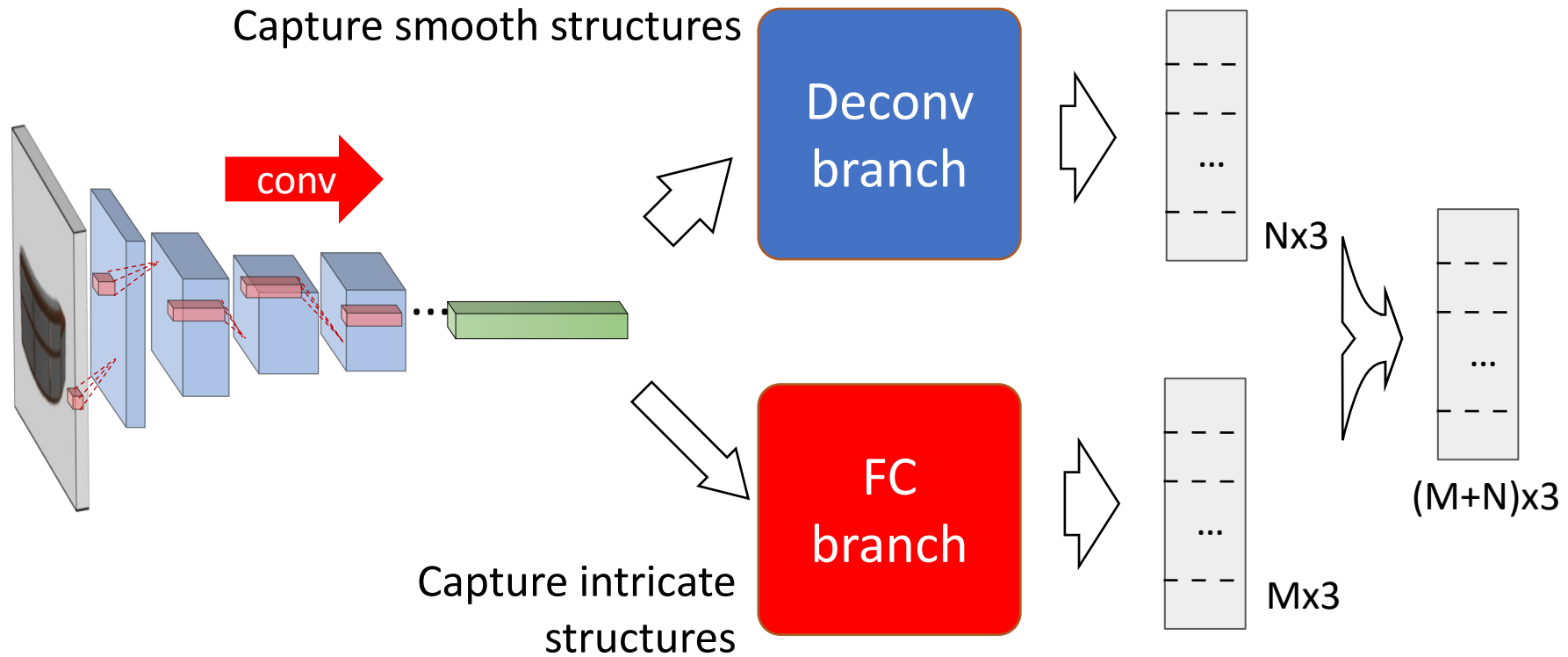
- Many local structures are common
  - e.g., planar patches, cylindrical patches
  - **strong local correlation** among point coordinates

# Natural Statistics of Object Geometry

- Many local structures are common/shared
  - e.g., planar patches, cylindrical patches
  - **strong local correlation** among point coordinates
- But also some intricate local structures
  - some points have **high variability** neighborhoods

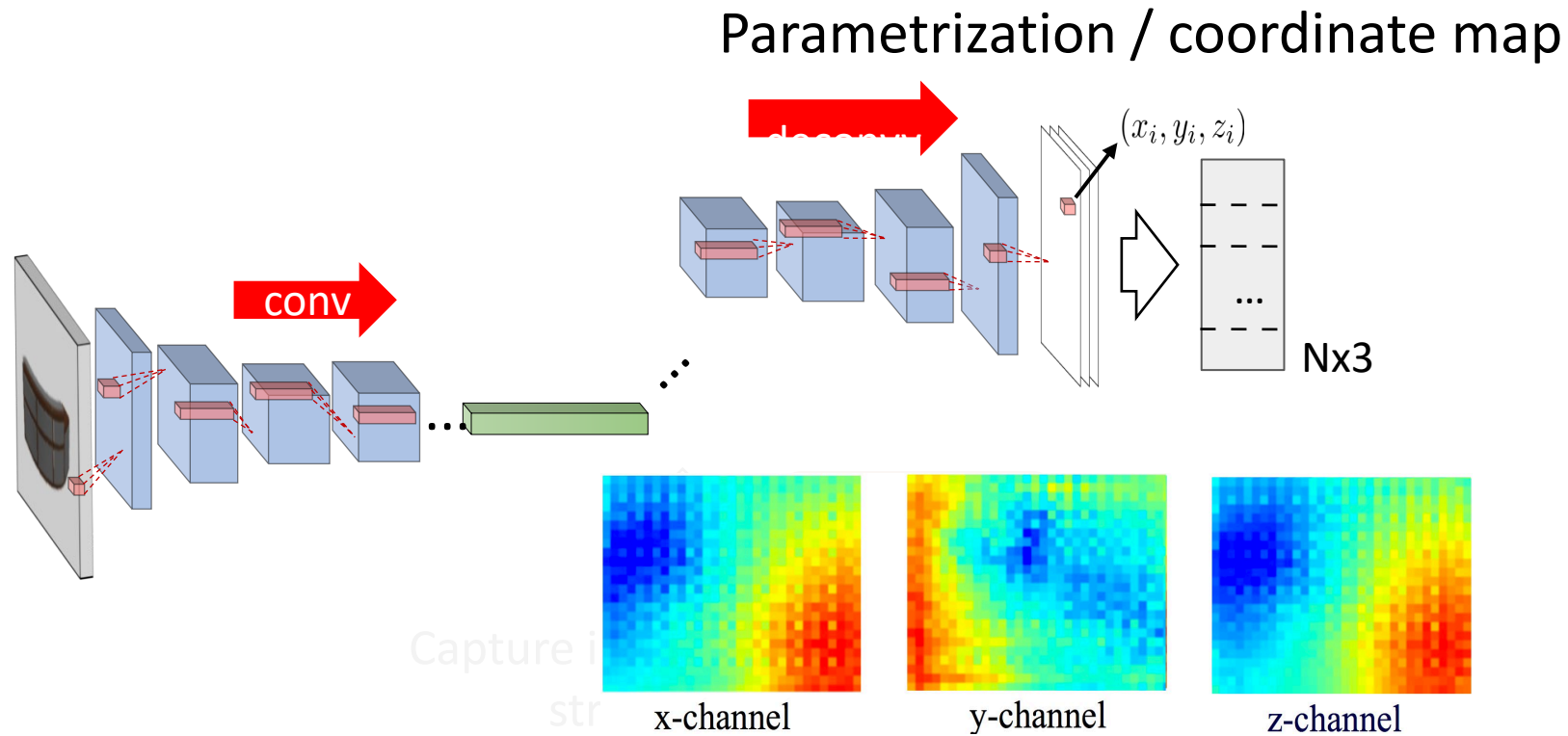


# Two-Branch Architecture



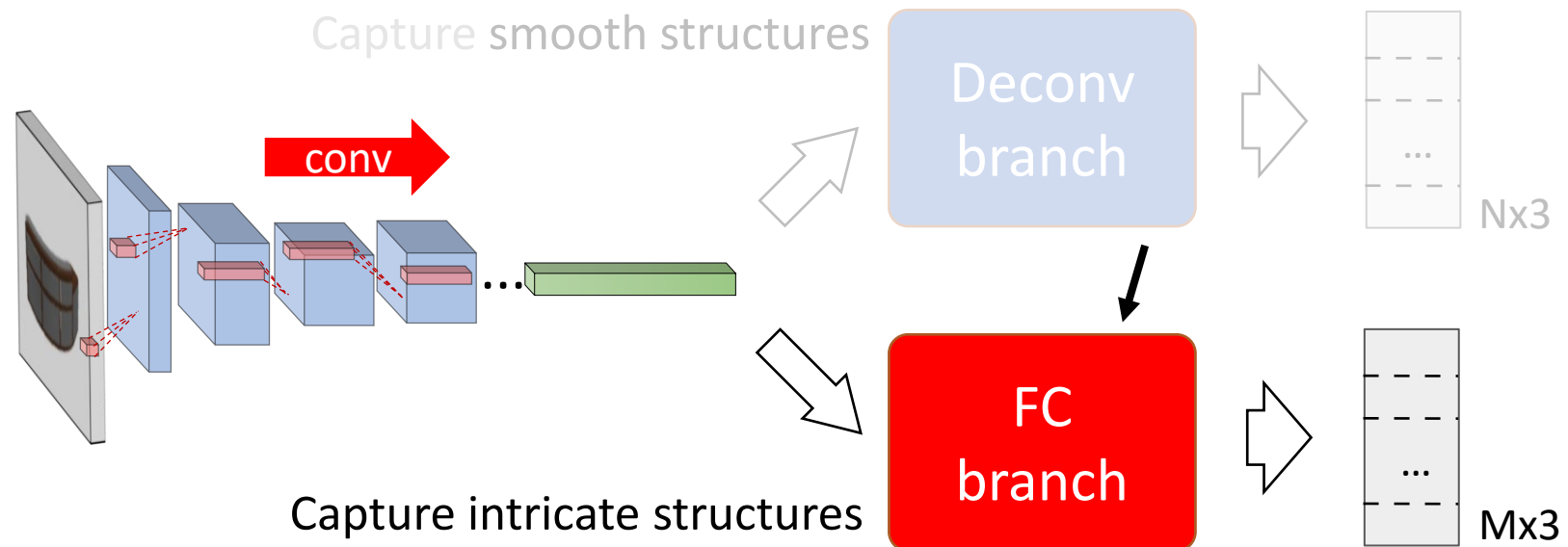
**Set union by array concatenation**

# Deconvolution Branch



- Deconvolution induces a smooth coordinate map
- Geometrically, learns a smooth parameterization

# Fully Connected Branch



# The Two Branches

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

**red:** fully-connected branch – **more intricate** structures



# Example Results

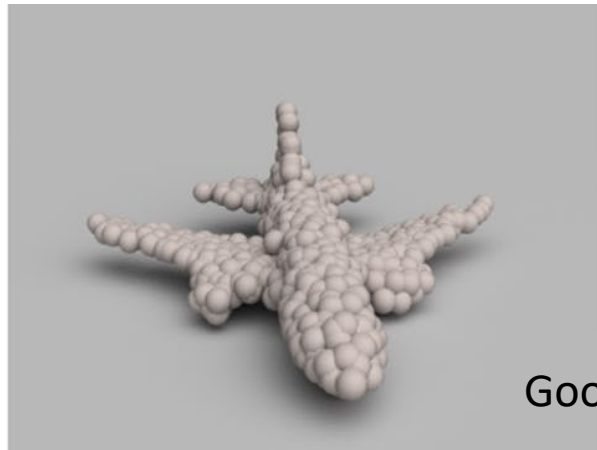
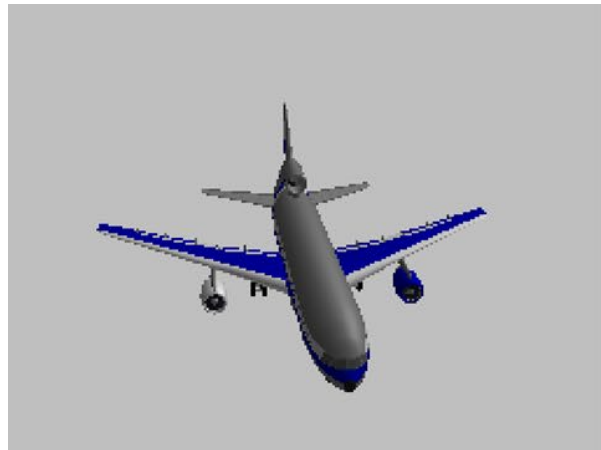


Same view

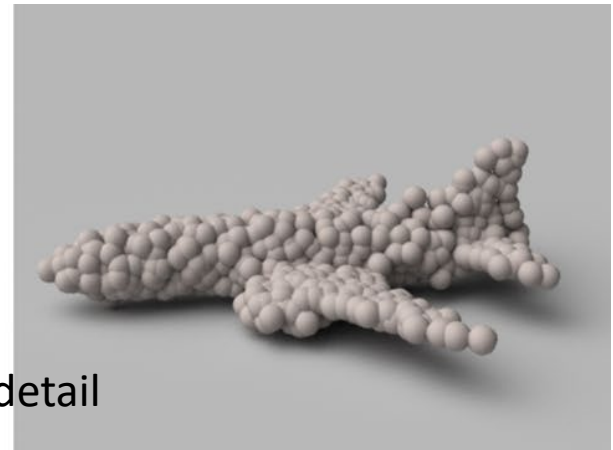


Good symmetry

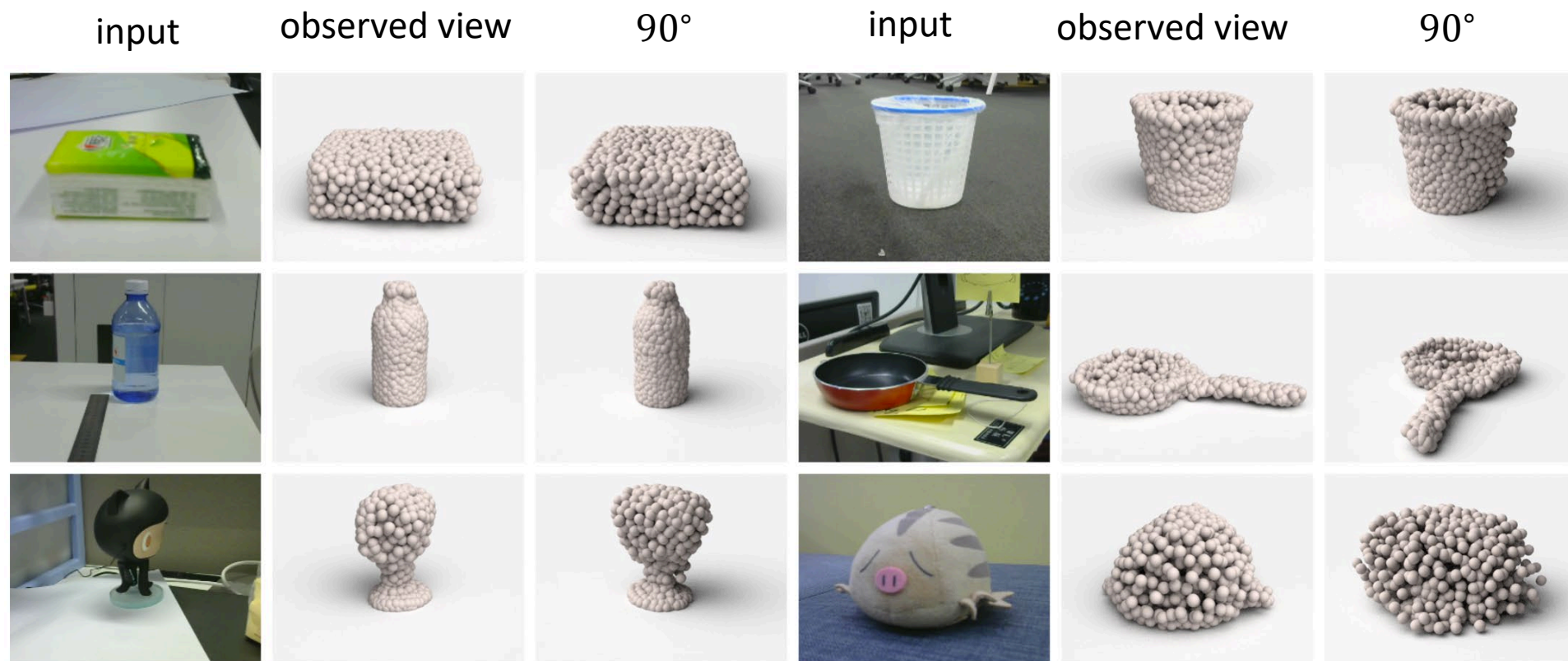
New view



Good detail



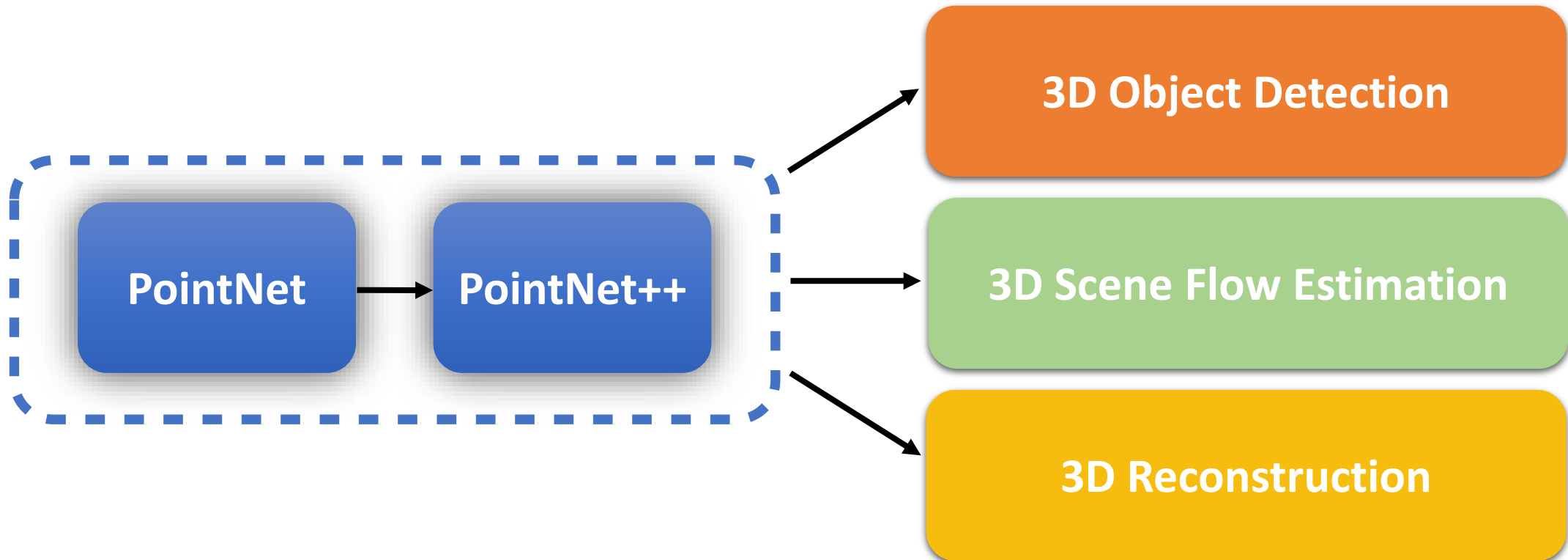
# From Real Images



Out of training categories

# Conclusions: Real-World 3D Understanding

- **Novel architectures for deep learning on point clouds** – PointNet and PointNet++, respecting invariances, light-weight and robust to data corruption, a unified framework for various tasks.
- **Successful applications in 3D scene understanding.**



# That's All

