

# PointNet:

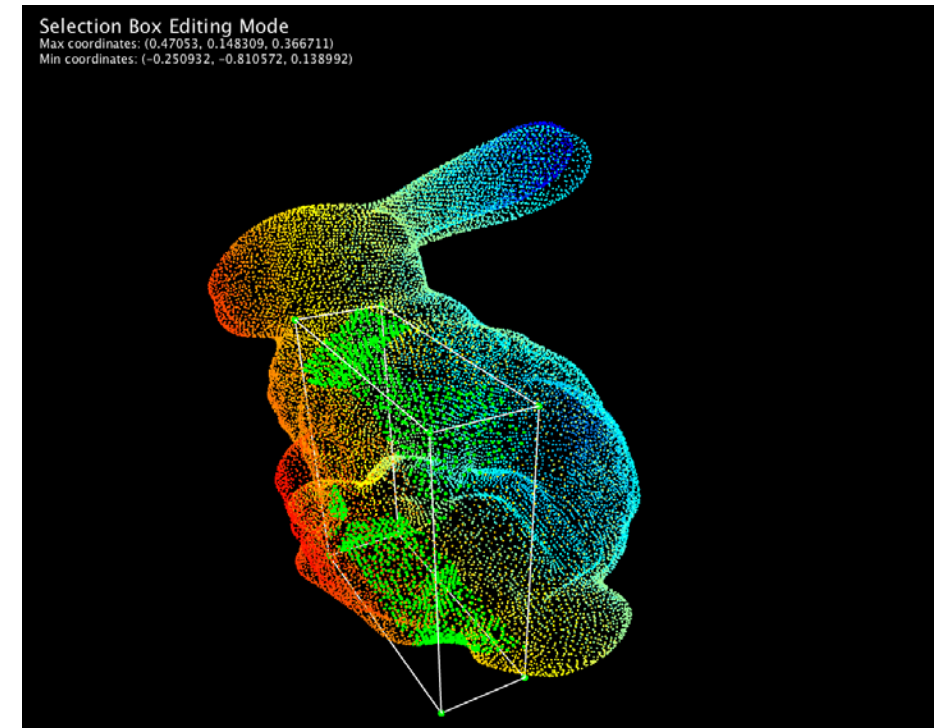
## Deep Learning on Point Sets for 3D Classification and Segmentation

CVPR 2017, Stanford

<http://stanford.edu/~rqi/pointnet/>

<https://github.com/fxia22/pointnet.pytorch> (not official)

# Point cloud



# Point cloud

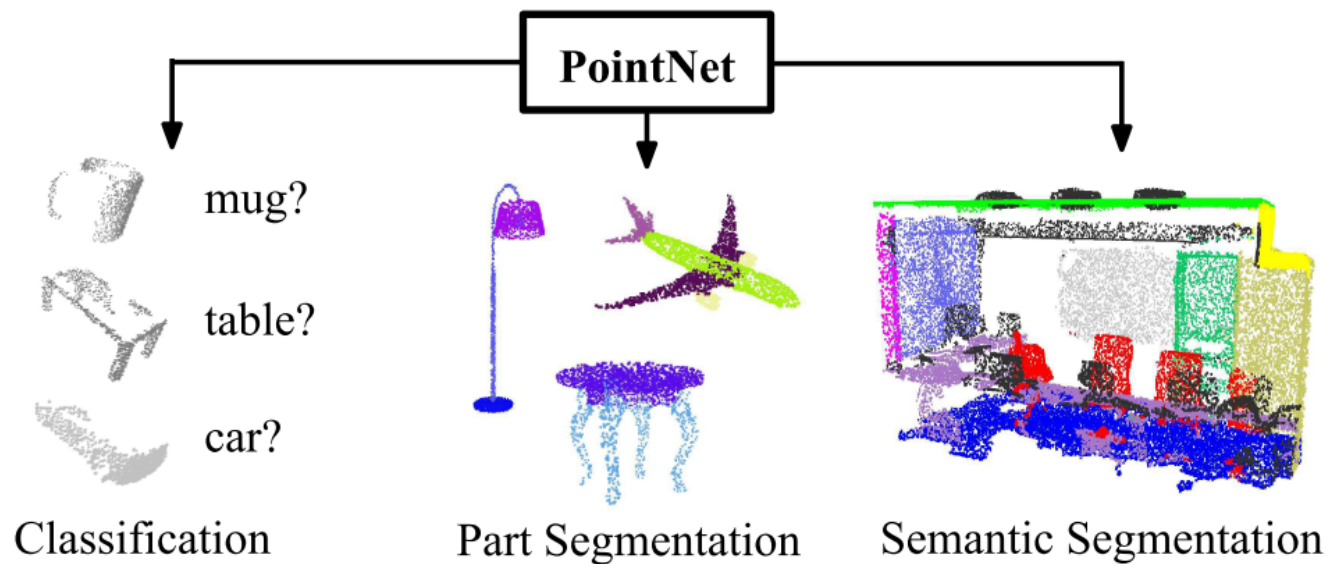


Figure 1. **Applications of PointNet.** We propose a novel deep net architecture that consumes raw point cloud (set of points) without voxelization or rendering. It is a unified architecture that learns both global and local point features, providing a simple, efficient and effective approach for a number of 3D recognition tasks.

# 3D featurization

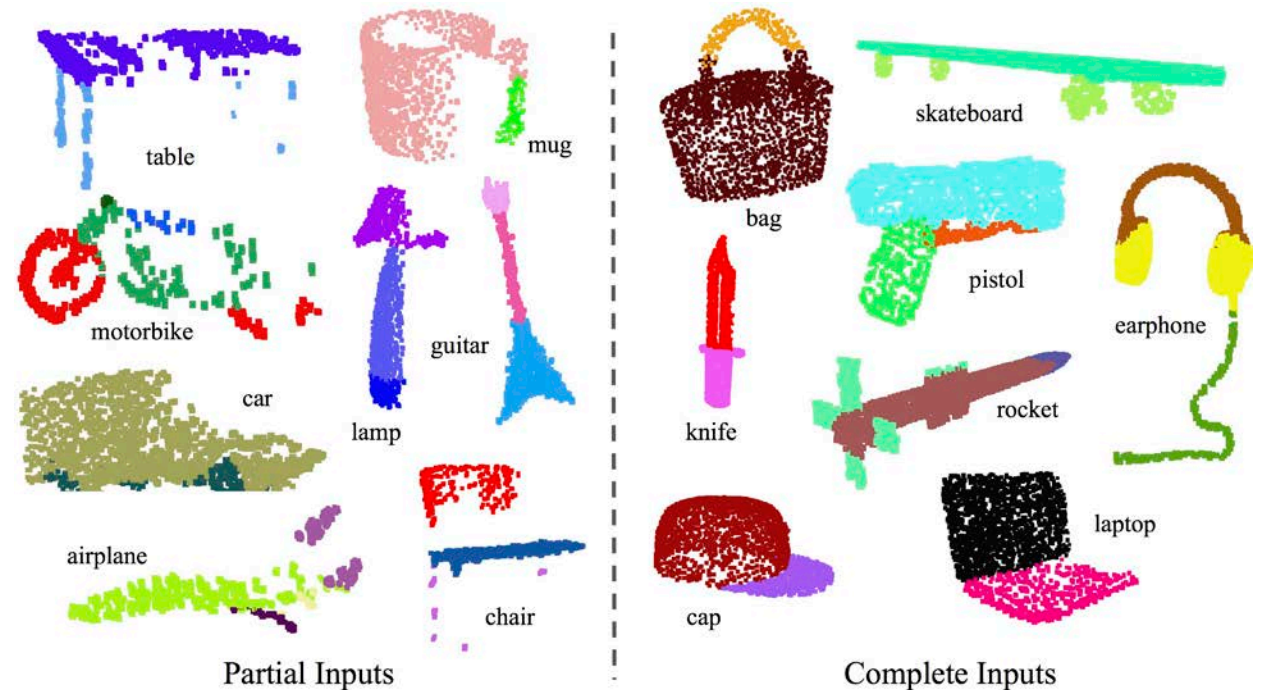
- Voxelization
  - Pixel → Voxel (마인크래프트 생각하면 됨)
  - 3D convolution: computation cost, sparsity, symmetry, ...
  - Voxelization 자체도 쉽지가 않다 (bleeding, ...)
- Rendering
  - Differentiable renderer / ray tracer
  - 3D object as multi-view 2D images
  - 얘기하기 시작하면 책 몇 권이 똑딱 나온다
- Point cloud
  - Canonical form

# Dataset: ShapeNet

- ShapeNet
  - <https://www.shapenet.org>
- PointNet provides preprocessed ShapeNet data
  - ShapeNetPart dataset
  - <https://github.com/charlesq34/pointnet>
  - [http://web.stanford.edu/~ericyi/project\\_page/part\\_annotation/index.html](http://web.stanford.edu/~ericyi/project_page/part_annotation/index.html)

# Dataset: ShapeNet (part)

- ShapeNetPart
  - Airplane, Bag, Cap, ...
  - 16 types of objects
- Each object has
  - Coordinates
  - Classes



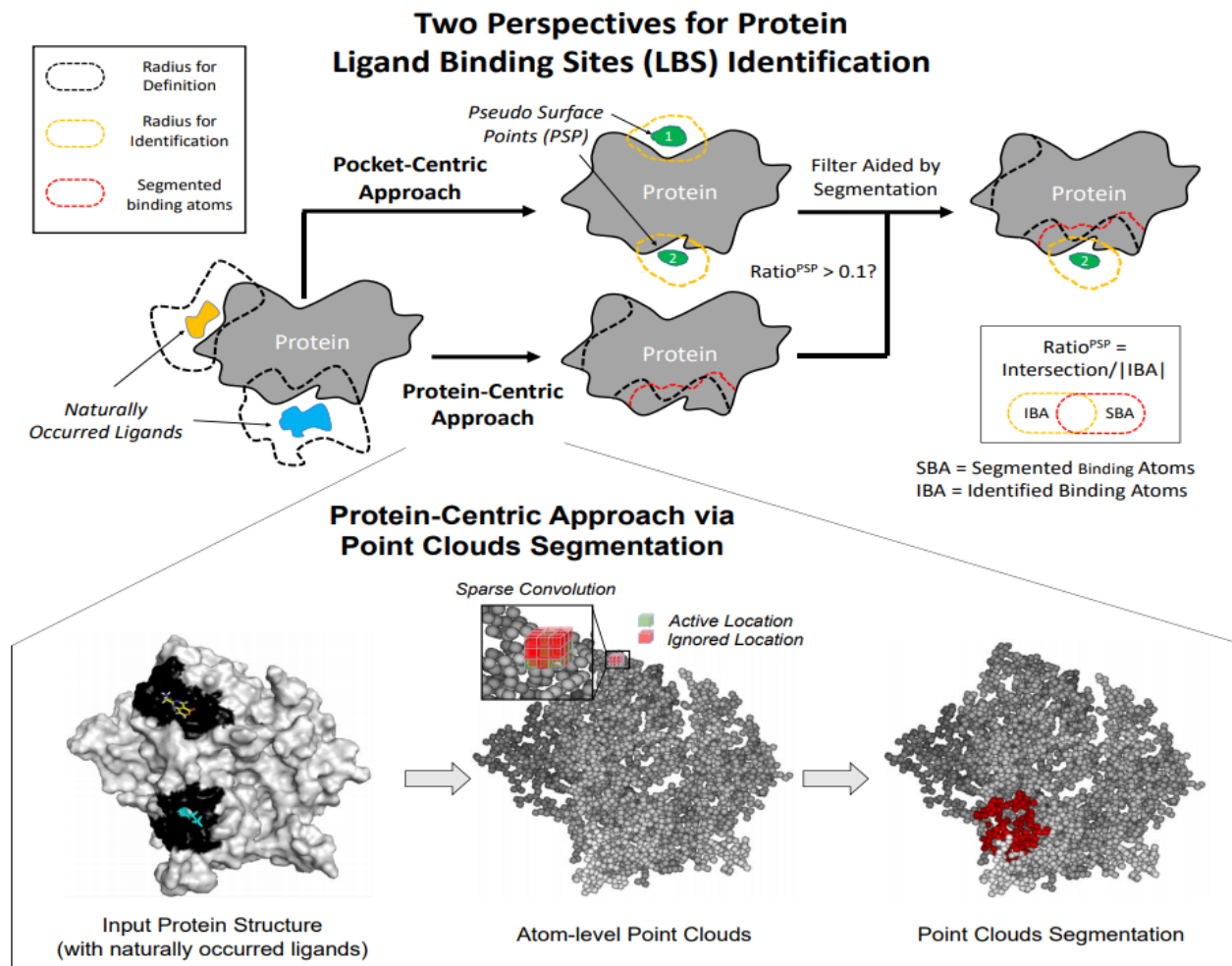
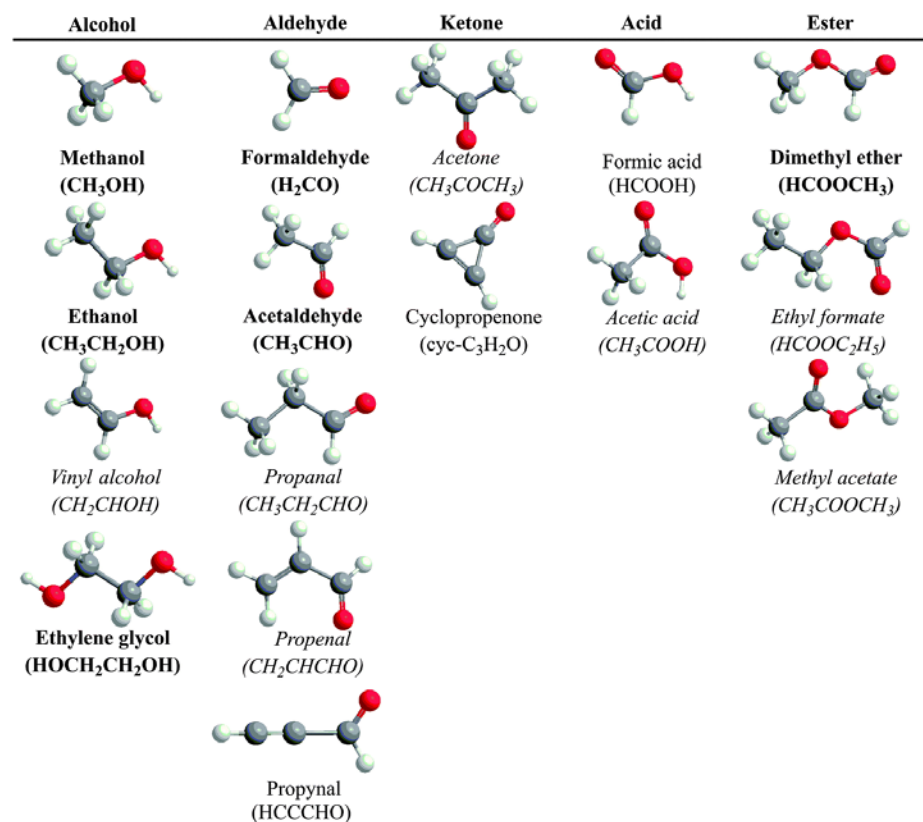
# Challenges

- Point cloud features
  - Rotation invariant
  - Translation invariant
  - Order invariant (unordered set)
- Further reading
  - Set transformer (ICML 2019)
  - 3D steerable CNNs (NIPS 2018)
  - Generalizing Convolutional Neural Networks for Equivariance to Lie Groups on Arbitrary Continuous Data (Arxiv)
  - Deep parametric continuous convolutional neural networks (CVPR 2018)
  - 아무튼 엄청 많음



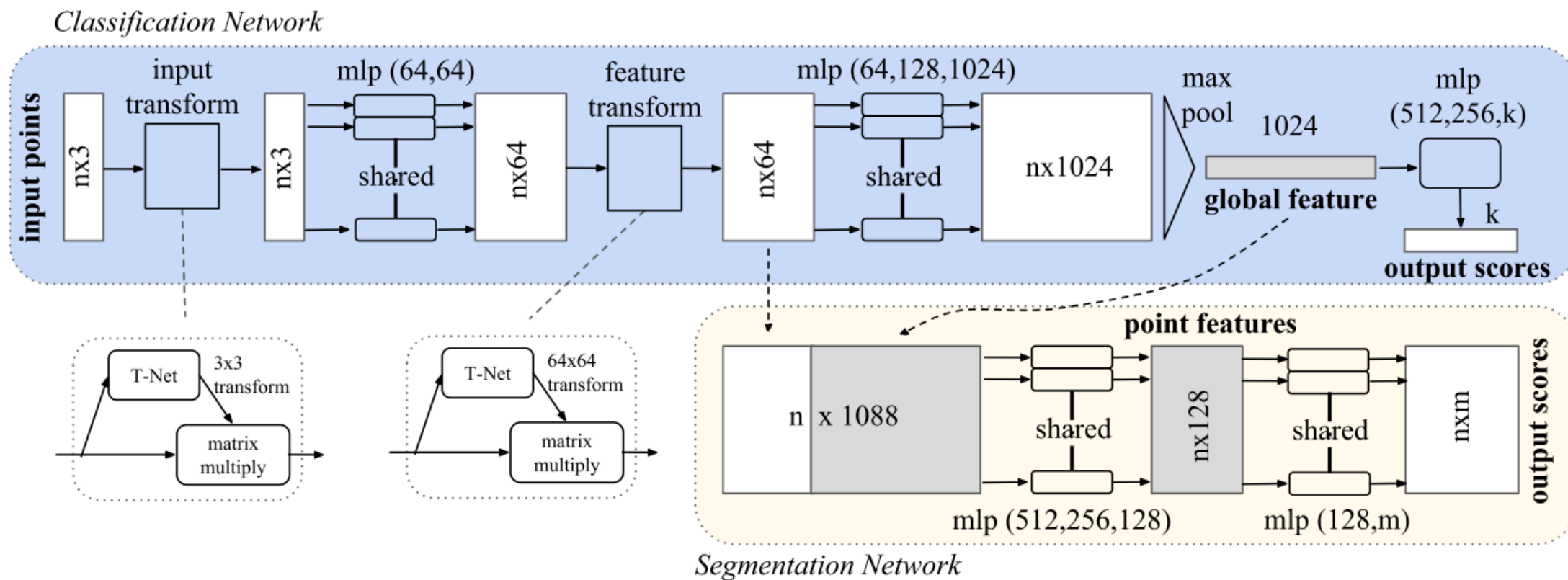
# Motivation

- My motivation





# Method



# Method: Implementation

- Main
  - $x$ :  $(b, 3, n)$
  - $t1 = \text{Tnet3d}(x)$ :  $(3, 3)$
  - $x = \text{bmm}(x, t1)$ :  $(b, 3, n) \rightarrow$  transpose 생략
  - $x = \text{ReLU}(\text{BatchNorm}(\text{Conv1d}(x)))$ :  $(b, 64, n)$
  - ...

# Method: Implementation

- Tnet3d
  - $x$ :  $(b, 3, n)$
  - $x = \text{ReLU}(\text{BatchNorm}(\text{Conv1d}(x)))$ :  $(b, 64, n)$
  - $x = \text{ReLU}(\text{BatchNorm}(\text{Conv1d}(x)))$ :  $(b, 128, n)$
  - $x = \text{ReLU}(\text{BatchNorm}(\text{Conv1d}(x)))$ :  $(b, 1024, n)$
  - $x = \max(x, \text{dim} = 2)$ :  $(b, 1024)$
  - $x = \text{MLP}(x)$ :  $(b, 9) \rightarrow 3$  layers
  - $x = x.\text{view}(b, 3, 3) + \text{Identity}$ :  $(b, 3, 3)$
  - return  $x$

# Method: Preprocessing

- Batch processing
  - Randomly select  $N$  points
- Normalization
  - Centering
  - Scaling (by dividing max distance)
- Data augmentation (optional)
  - Rotation (with fixed y-axis)
  - Jittering  $\sim N(0, 0.02)$

# Method: T-net

- T-net predict transformation matrix  $A$ 
  - Output: 3 by 3 (input), 64 by 64 (feature) matrix
- Orthogonal regularization
  - @ feature transformation matrix
  - $L_{reg} = \|I - AA^T\|_F^2$

Transform	accuracy
none	87.1
input (3x3)	87.9
feature (64x64)	86.9
feature (64x64) + reg.	87.4
both	<b>89.2</b>

Table 5. **Effects of input feature transforms.** Metric is overall classification accuracy on ModelNet40 test set.

# Result

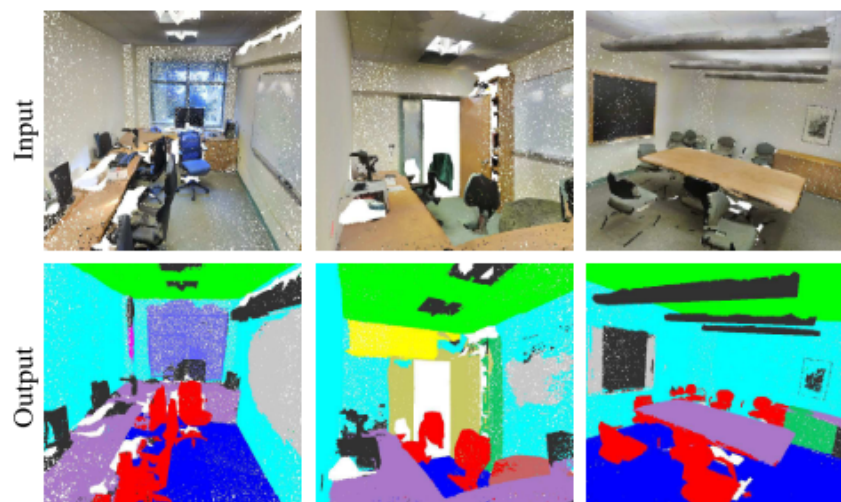


Figure 4. **Qualitative results for semantic segmentation.** Top row is input point cloud with color. Bottom row is output semantic segmentation result (on points) displayed in the same camera viewpoint as input.

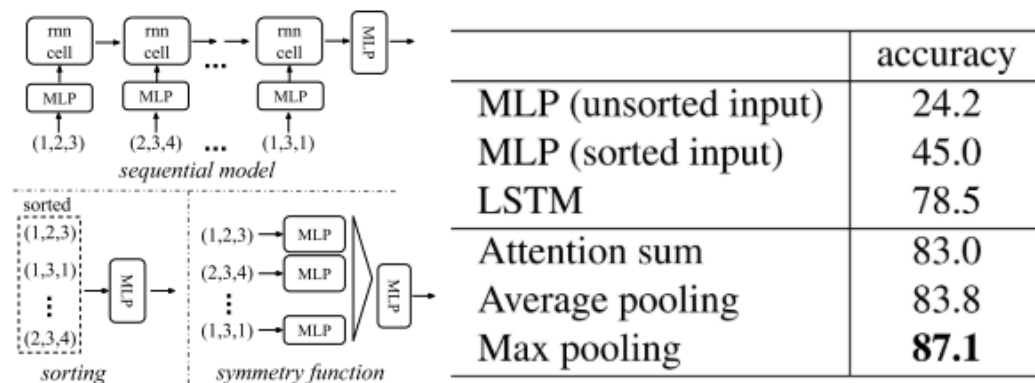


Figure 5. **Three approaches to achieve order invariance.** Multi-layer perceptron (MLP) applied on points consists of 5 hidden layers with neuron sizes 64,64,64,128,1024, all points share a single copy of MLP. The MLP close to the output consists of two layers with sizes 512,256.



# Result

- Baseline: Yi et al. (SIGGRAPH Asia 2016)
  - No deep learning (CRF based)

	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 [24]	-	63.2	-	-	-	73.5	-	-	-	74.4	-	-	-	-	-	-	74.8
Yi [26]	81.4	81.0	78.4	77.7	<b>75.7</b>	87.6	61.9	<b>92.0</b>	85.4	<b>82.5</b>	<b>95.7</b>	<b>70.6</b>	91.9	<b>85.9</b>	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	<b>83.7</b>	<b>83.4</b>	<b>78.7</b>	<b>82.5</b>	74.9	<b>89.6</b>	<b>73.0</b>	91.5	<b>85.9</b>	80.8	95.3	65.2	<b>93.0</b>	81.2	<b>57.9</b>	<b>72.8</b>	<b>80.6</b>

Table 2. **Segmentation results on ShapeNet part dataset.** Metric is mIoU(%) on points. We compare with two traditional methods [24] and [26] and a 3D fully convolutional network baseline proposed by us. Our PointNet method achieved the state-of-the-art in mIoU.

# Result

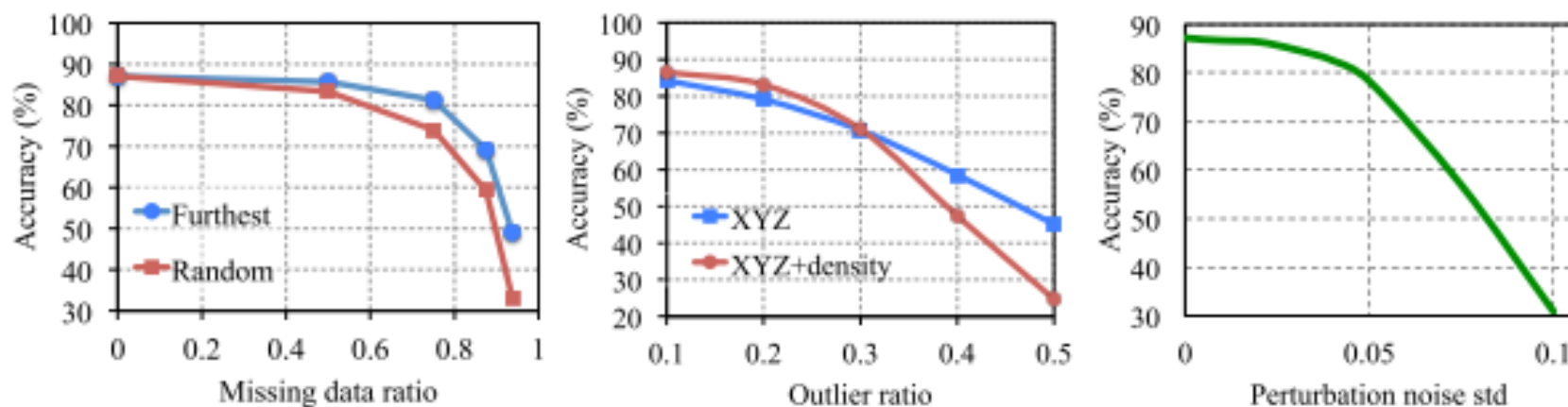


Figure 6. **PointNet robustness test.** The metric is overall classification accuracy on ModelNet40 test set. Left: Delete points. Furthest means the original 1024 points are sampled with furthest sampling. Middle: Insertion. Outliers uniformly scattered in the unit sphere. Right: Perturbation. Add Gaussian noise to each point independently.

# Result

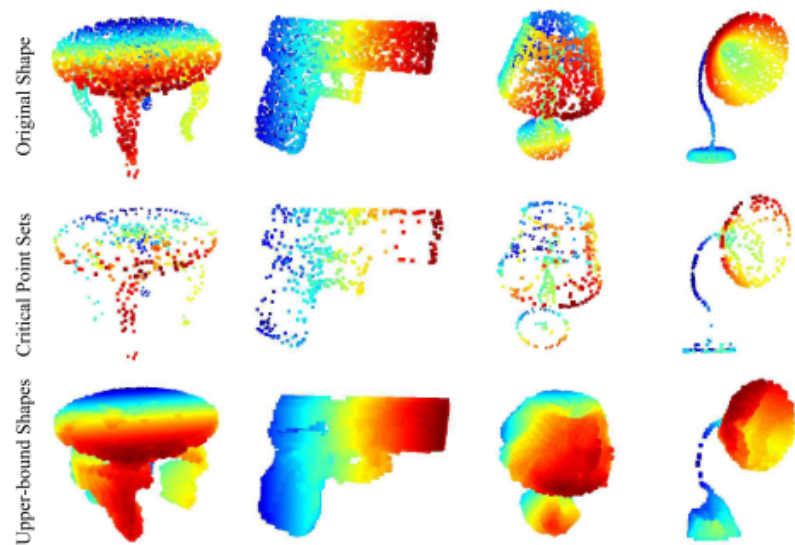


Figure 7. **Critical points and upper bound shape.** While critical points jointly determine the global shape feature for a given shape, any point cloud that falls between the critical points set and the upper bound shape gives exactly the same feature. We color-code all figures to show the depth information.

	#params	FLOPs/sample
PointNet (vanilla)	0.8M	148M
PointNet	3.5M	440M
Subvolume [16]	16.6M	3633M
MVCNN [20]	60.0M	62057M

Table 6. Time and space complexity of deep architectures for 3D data classification. PointNet (vanilla) is the classification PointNet without input and feature transformations. FLOP stands for floating-point operation. The “M” stands for million. Subvolume and MVCNN used pooling on input data from multiple rotations or views, without which they have much inferior performance.