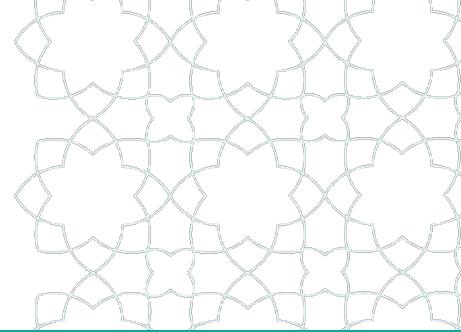




جامعة الملك عبد الله  
لعلوم والتكنولوجيا  
King Abdullah University of  
Science and Technology

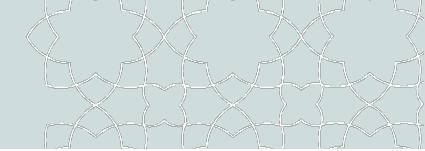
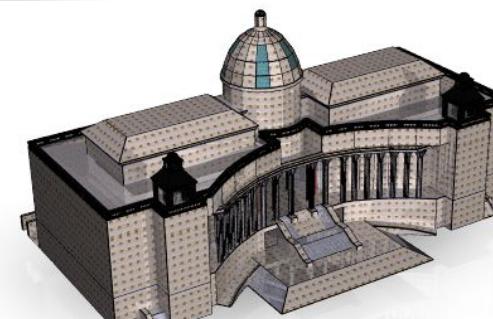
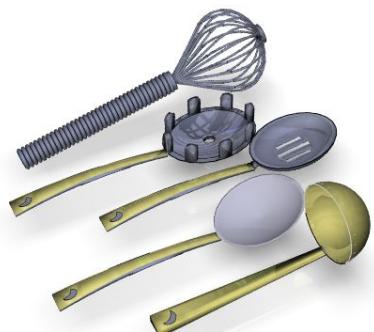
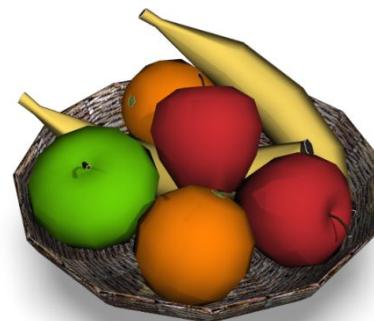
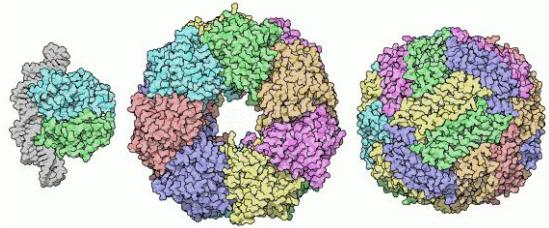


# 3D Deep Learning

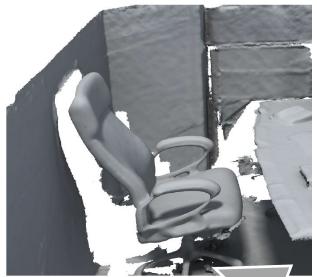
PROF. BERNARD GHANEM



# We live in a 3D world



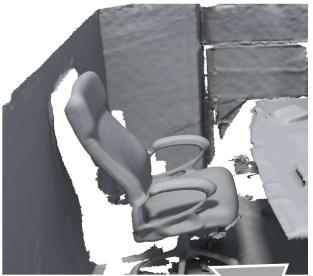
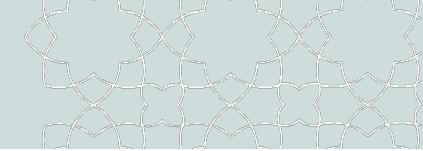
# Broad Application of 3D data



**Robotics**



# Broad Application of 3D data

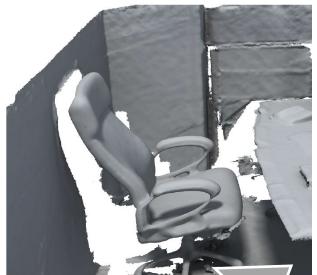
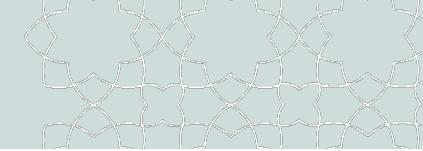


**Robotics**

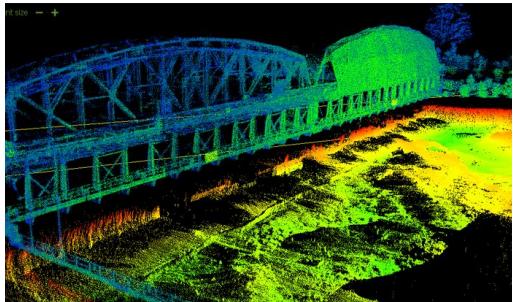


**Augmented  
Reality**

# Broad Application of 3D data



**Robotics**

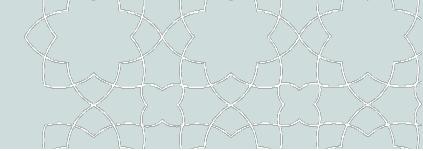


**Autonomous  
driving**



**Augmented  
Reality**

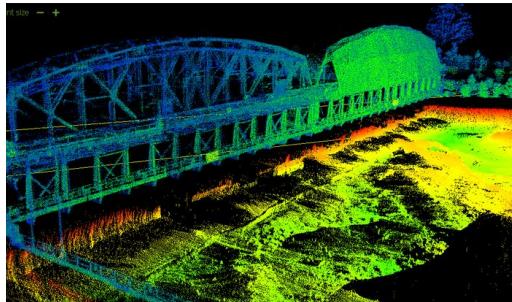
# Broad Application of 3D data



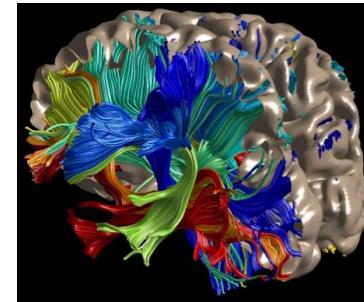
**Robotics**



**Augmented  
Reality**

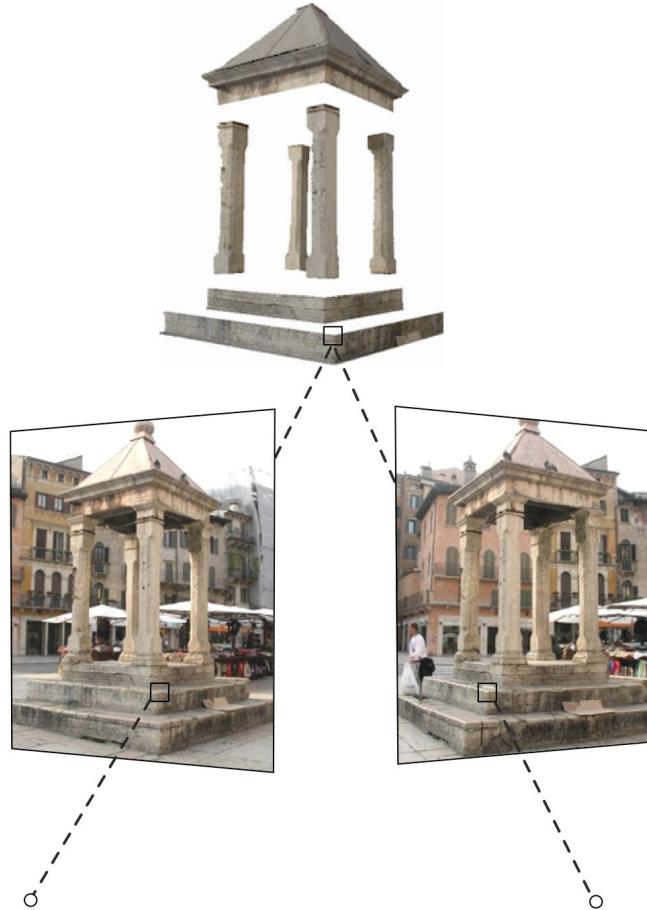
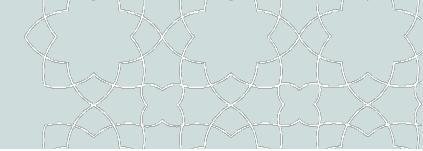


**Autonomous  
driving**



**Medical Image  
Processing**

# Traditional 3D Vision: Physics based

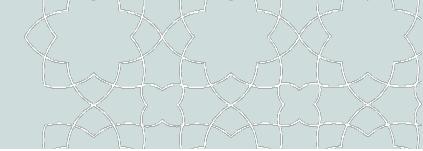


Multi-view  
Geometry

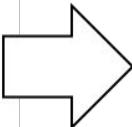
# 3D Learning: Knowledge Based



# 3D Learning: Knowledge Based

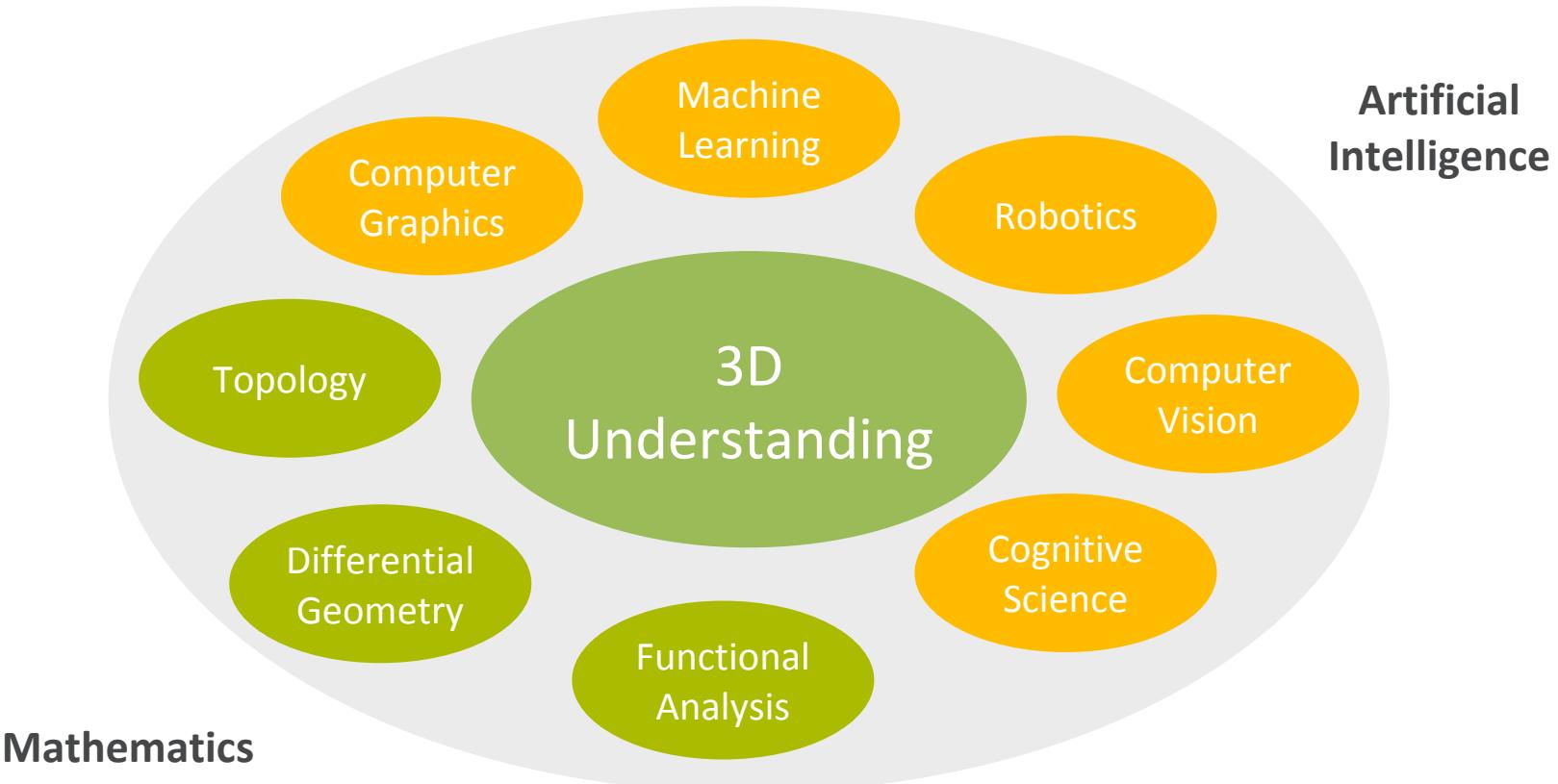


## Acquire Knowledge of 3D World by **Learning**



A priori knowledge of  
the 3D world

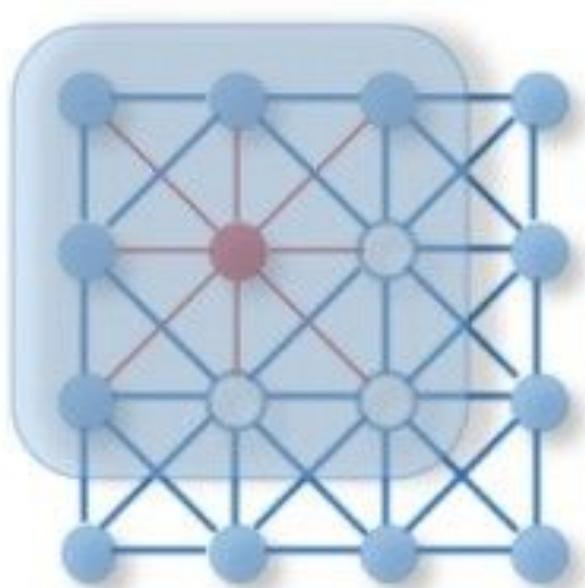
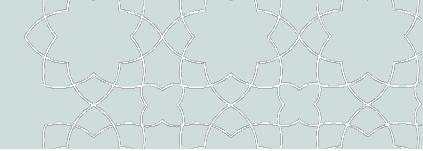
# A New Rising Field



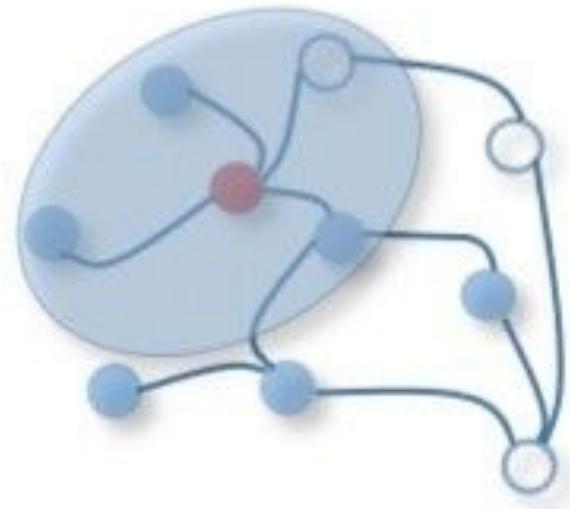
# Outline

- 3D Data
- Classification
- Segmentation and Detection
- 3D Data synthesis

# The 3D Representation Challenge

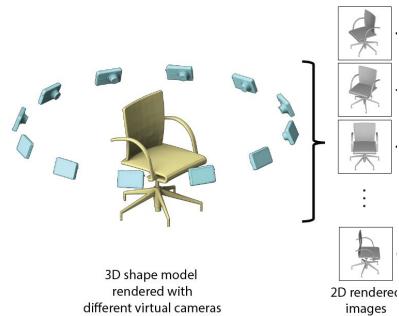


**Rasterized form  
(regular grids)**

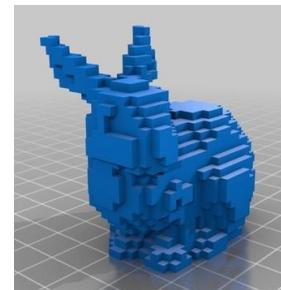


**Geometric form  
(irregular)**

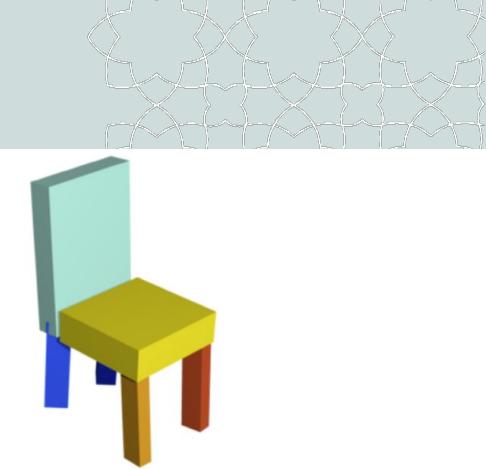
# The 3D Representation Challenge



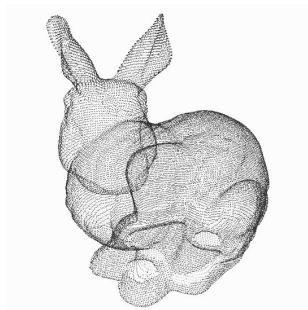
Multi-view



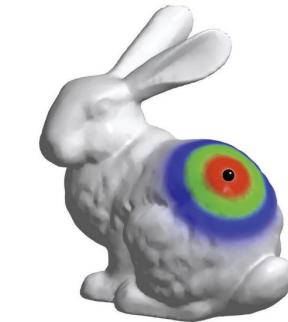
Volumetric



Part Assembly



Point Cloud

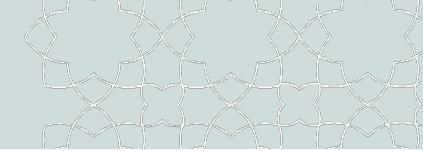


Mesh (Graph CNN)

$$F(x) = 0$$

Implicit Shapes

# Datasets for 3D Objects



- **ModelNet:** Classification - absorbed by ShapeNet
- **ShapeNet:** Large-scale Synthetic Objects
- **3DScan:** Consumer-grade 3D Scanning

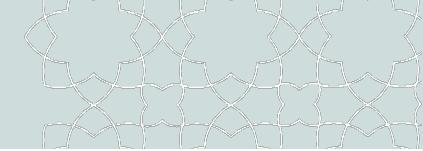


Chang et al., "ShapeNet: An Information-Rich 3D Model Repository" , *arXiv 2015*

Wu et al., "3D ShapeNets: A deep representation for volumetric shapes", *CVPR 2015*

Choi et al., "A Large Dataset of Object Scans", *arXiv 2016*

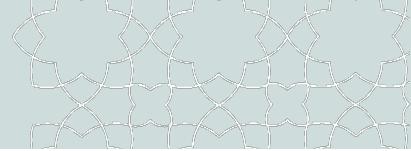
# Datasets for 3D Object Parts



- **ShapeNetPart2019:** Fine-grained Part
  - Fine-grained (towards mobility)
  - Instance-level / Hierarchical



# Datasets for Indoor 3D Scenes



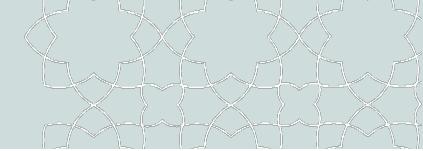
- **SceneNet**: Large-scale Synthetic Scenes
  - 3D meshes / 5M Photorealistic Images



Ankur et al., "Understanding RealWorld Indoor Scenes with Synthetic Data", CVPR 2016

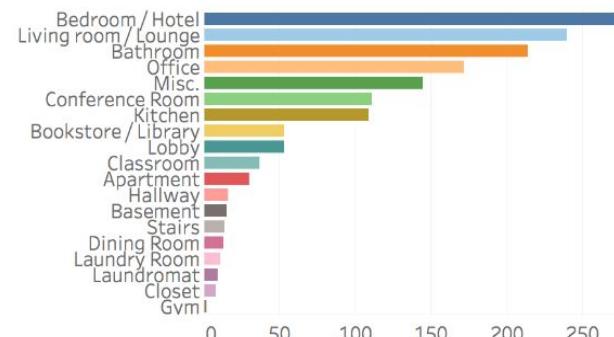
McCormac et al., "SceneNet RGB-D: Can 5M Synthetic Images Beat Generic ImageNet Pre-training on Indoor Segmentation?", ICCV 2017

# Datasets for Indoor 3D Scenes

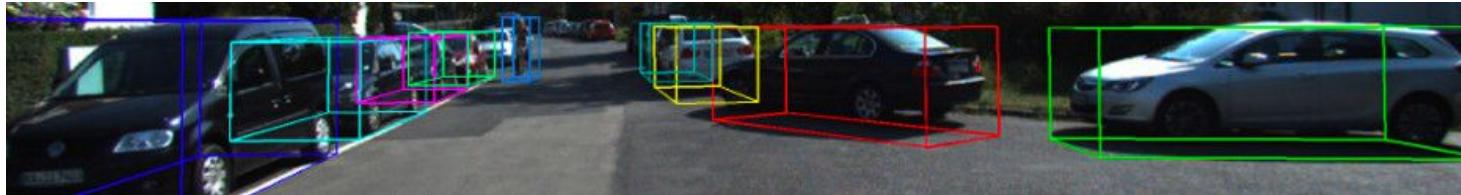
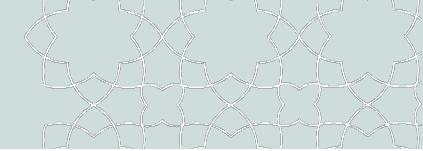


- **ScanNet: Large-scale Scanned Real Scenes**

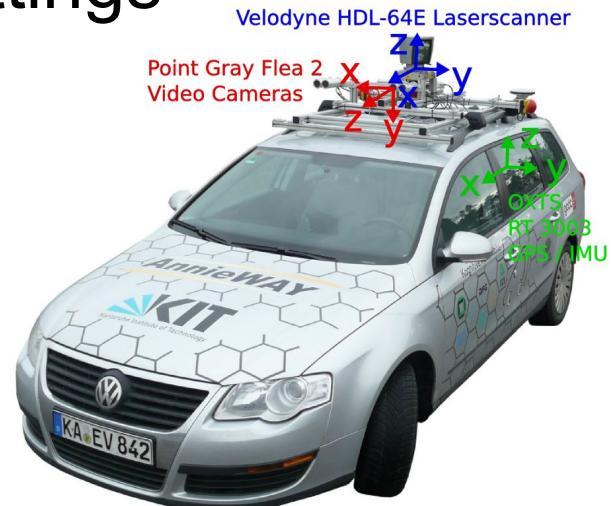
- 2.5 M Views
- 1500 RGBD scans
- 3D camera poses
- Surface reconstruction
- Instance-level semantic segmentation



# Datasets for Outdoor 3D Scenes



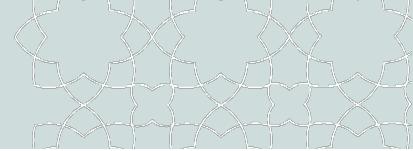
- **KITTI: Autonomous navigation settings**
  - Camera, LiDAR, GPS, IMU
  - 14000 frames training/testing
  - Stereo, Flow, Depth, Odometry, ...
  - Detection, Tracking, Road, Semantic, ...



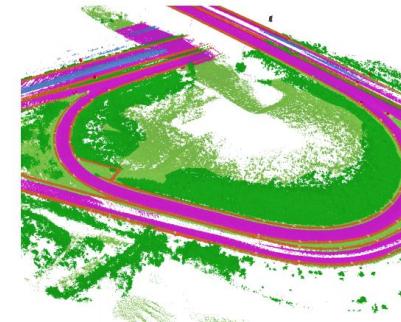
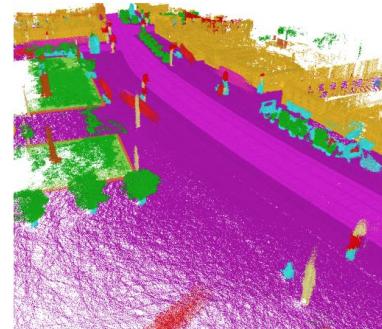
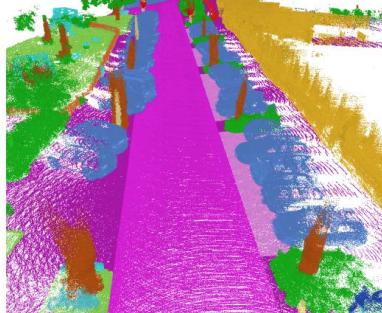
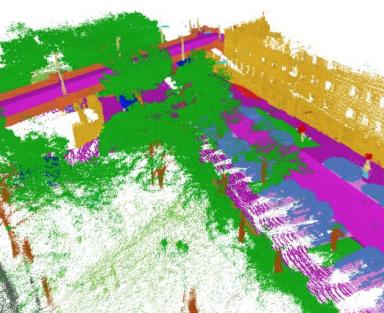
Geiger et al., "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite", CVPR 2012

Geiger et al., "Vision meets Robotics: The KITTI Dataset", IJRR 2013

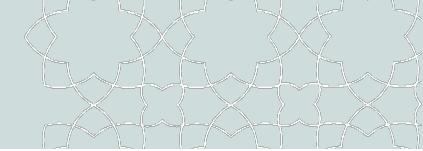
# Datasets for Outdoor 3D Scenes



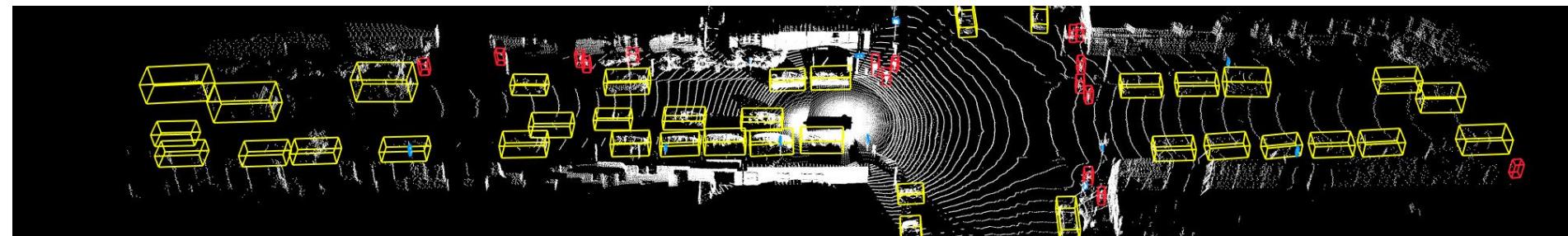
- **Semantic KITTI:** LiDAR data, labeled per point
  - Semantic Segmentation
  - Panoptic Segmentation
  - Semantic Scene Completion



# Datasets for Outdoor 3D Scenes



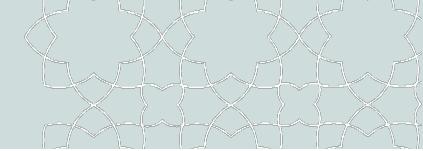
- **Waymo Open Dataset:** Autonomous Navigation at Large Scale
  - Camera / LiDAR data
  - 230.000 frames (1150 temporal sequences, 12M BB)
  - 2D/3D Object Detection/Tracking
  - Domain Adaptation (700 extra sequences from different locations)



# Outline

- 3D Data
- **Classification**
- Segmentation and Detection
- 3D Data synthesis

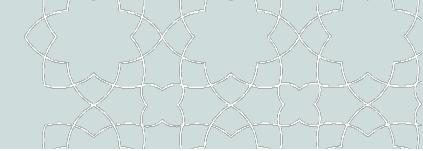
# Task: 3D Classification



This is a chair!

- **Covered methods:** Volumetric CNN, OctNet, O-CNN, SparseConvNet, PointNet, PointNet++, RS CNN, DGCNN, Point ConvNet, KPConv, Monte Carlo Point Convolution, PConv, Multi-View CNN, Spectral CNN, Synchronized Spectral CNN, Spherical CNN

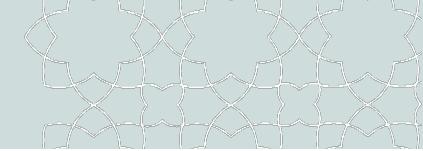
# Multi-View CNN



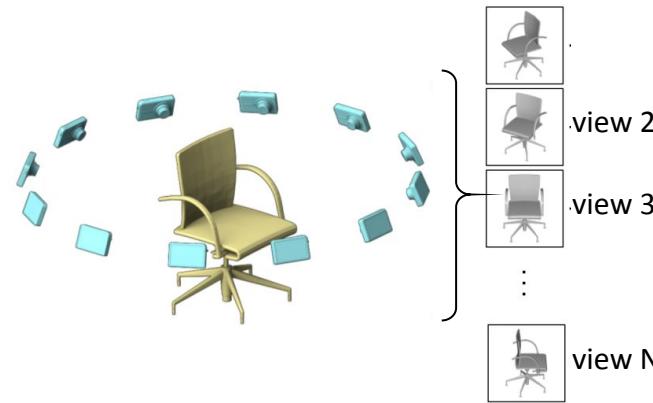
- Given an Input Shape



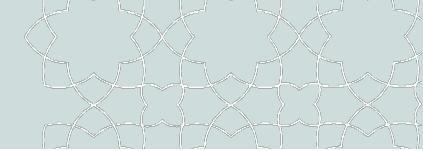
# Multi-View CNN



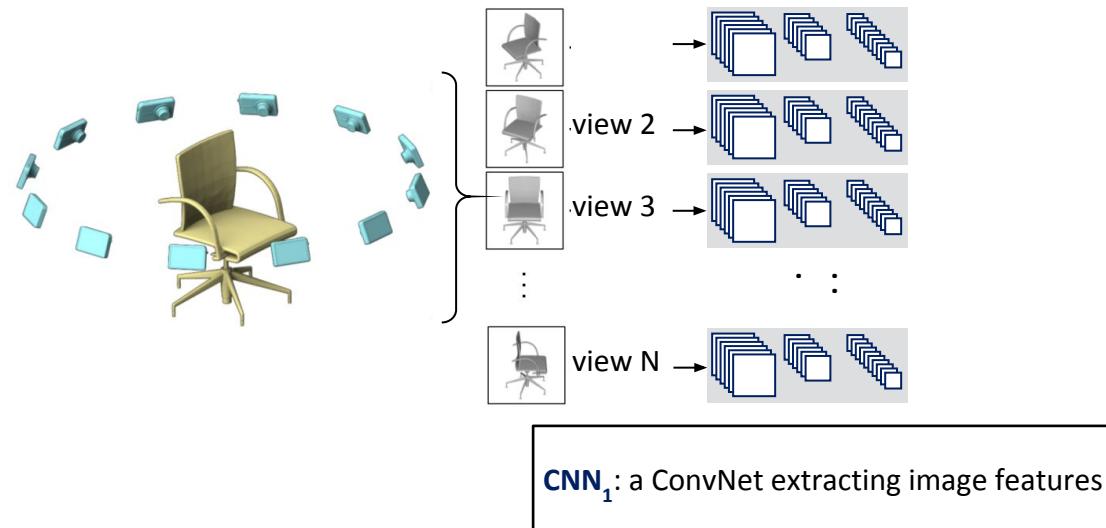
- Render with Multiple Virtual Cameras



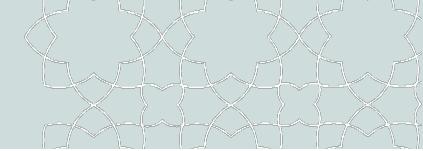
# Multi-View CNN



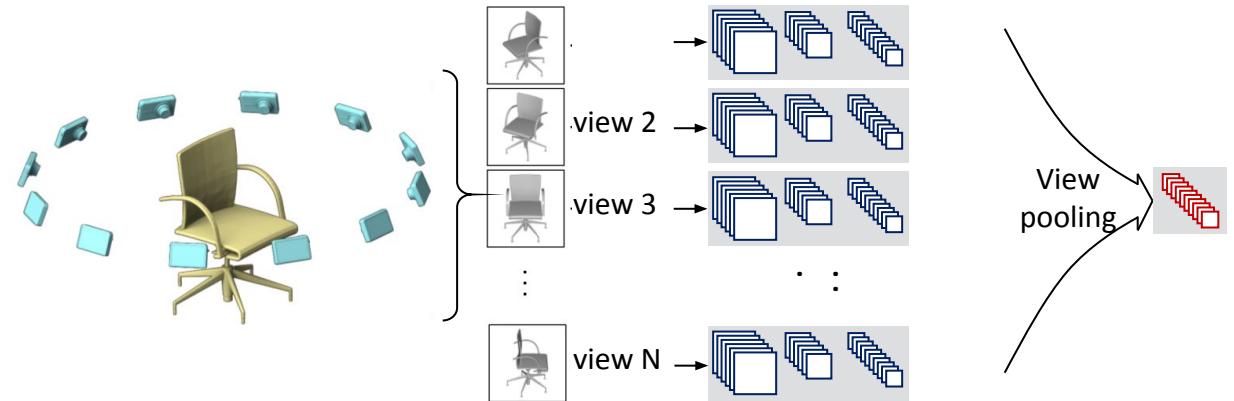
- The Rendered Images are Passed through  $\text{CNN}_1$  for Image Features



# Multi-View CNN

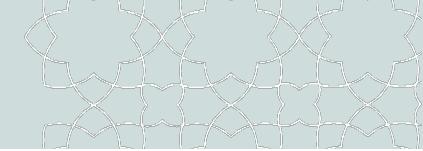


- All Image Features are Combined by View Pooling

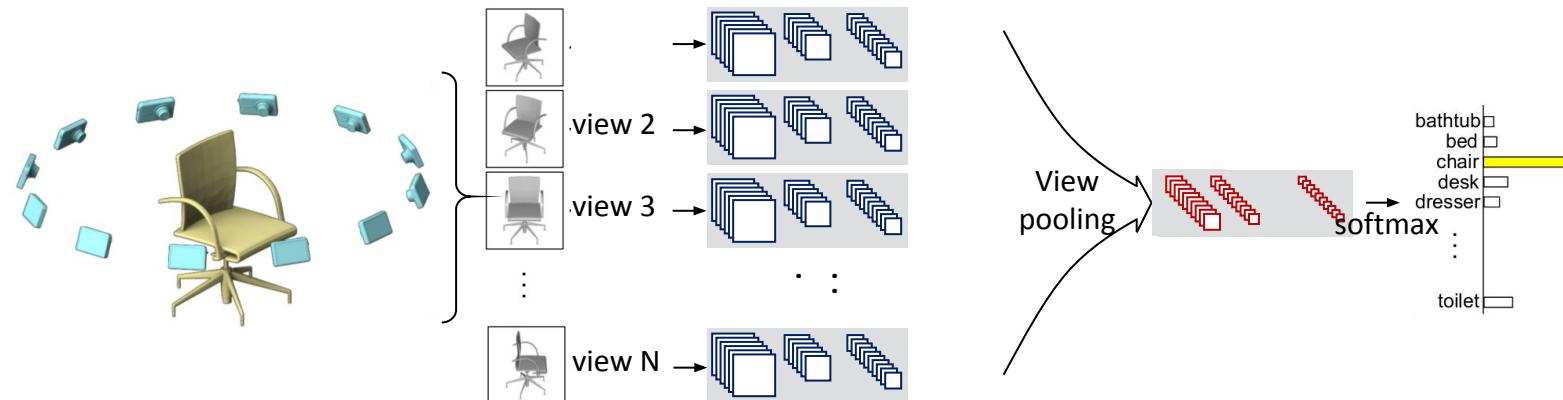


**View pooling:**  
element-wise max-pooling  
across all views

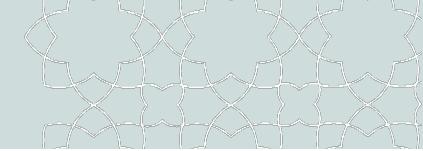
# Multi-View CNN



- ... and then Passed through  $\text{CNN}_2$  and to Generate Final Predictions



**CNN<sub>2</sub>:** a second ConvNet producing shape descriptors



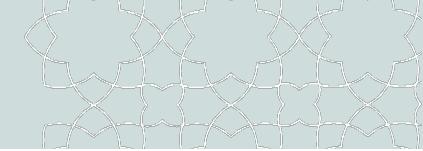
- Experiments – Classification & Retrieval

Method	Classification (Accuracy)	Retrieval (mAP)
Non-deep { SPH [16]	68.2%	33.3%
LFD [5]	75.5%	40.9%
3D ShapeNets [37]	77.3%	49.2%
FV, 12 views	84.8%	43.9%
CNN, 12 views	88.6%	62.8%
MVCNN, 12 views	<b>89.9%</b>	70.1%
MVCNN+metric, 12 views	89.5%	<b>80.2%</b>
MVCNN, 80 views	90.1%	70.4%
MVCNN+metric, 80 views	<b>90.1%</b>	<b>79.5%</b>

On ModelNet40

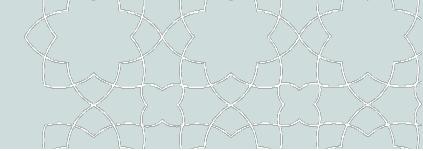
[credit: Hang Su]

# Multi-View CNN



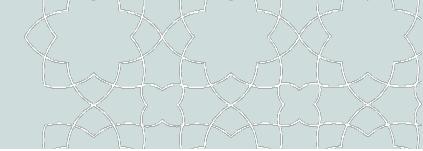
- Indeed gives good performance
- Can leverage vast literature of image classification
- Can use pre-trained features
- Needs projection
- What if the input is noisy and/or incomplete? e.g., point cloud

# Volumetric CNN



- Can we use CNNs without 2D-3D projection?
- Straightforward idea: 3D native convolution

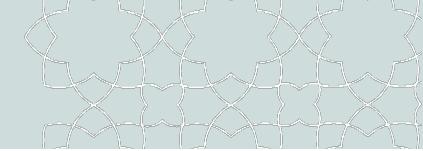
# Volumetric CNN: Voxelization



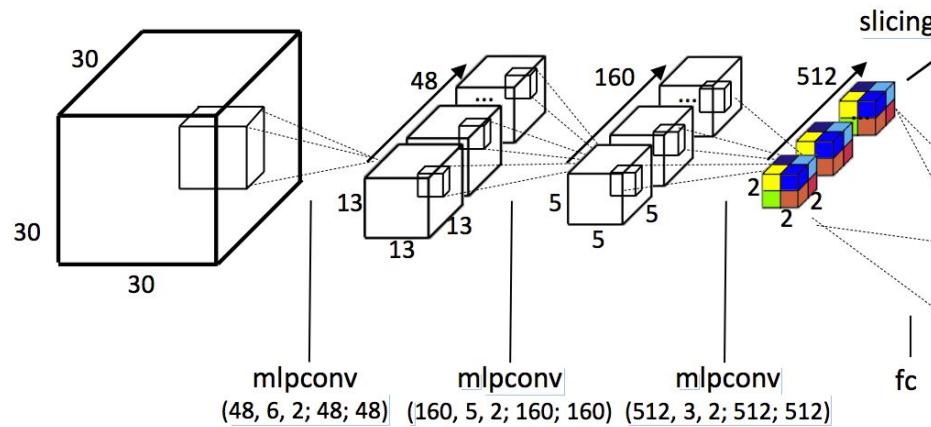
- Represent the occupancy of regular 3D grids



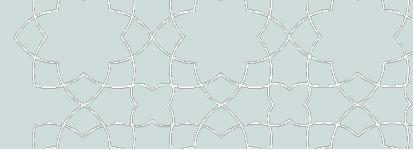
# Volumetric CNN: Voxelization



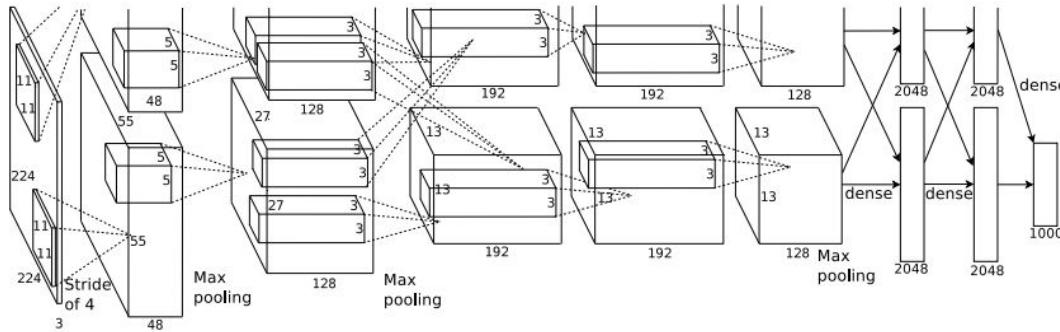
- 3D convolution uses 4D kernels



# Volumetric CNN: Voxelization



- Complexity Issue

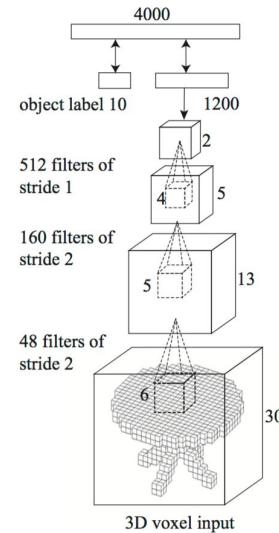


AlexNet, 2012

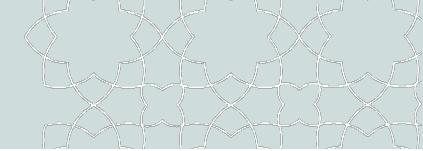
Input resolution:  
224x224=50176

3DShapeNets, 2015

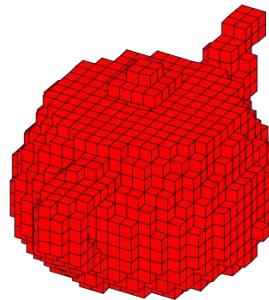
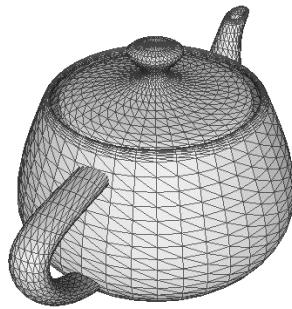
Input resolution:  
30x30x30 =27000



# Volumetric CNN: Voxelization



- Complexity Issue
- Information loss in voxelization



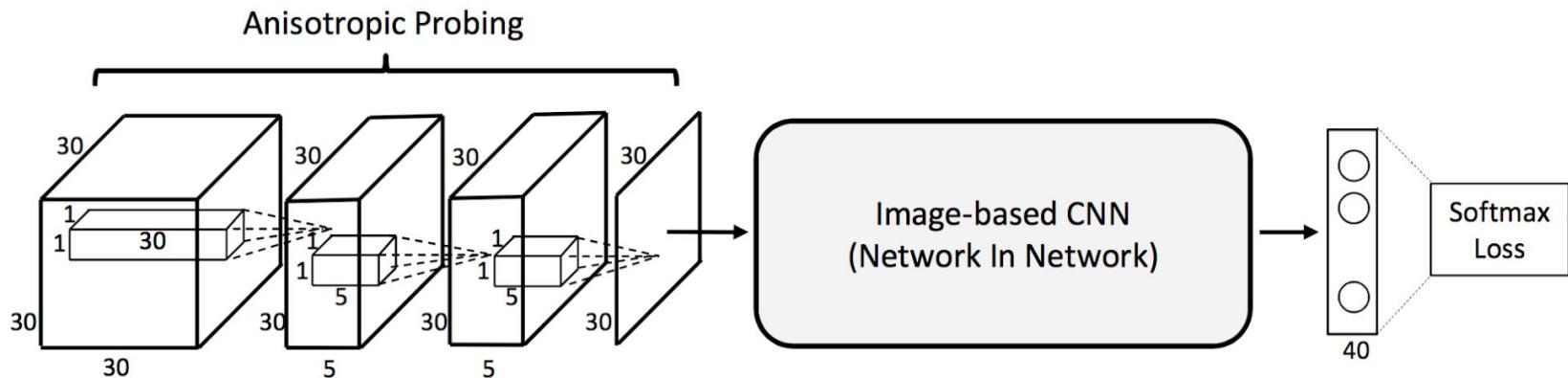
Polygon Mesh

Occupancy Grid  
 $30 \times 30 \times 30$

# Volumetric CNN: Voxelization

- Idea 1: Learn to Project

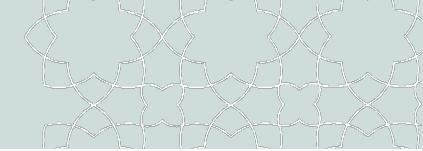
Idea: “X-ray” rendering + Image (2D) CNNs  
***very low #param, very low computation***



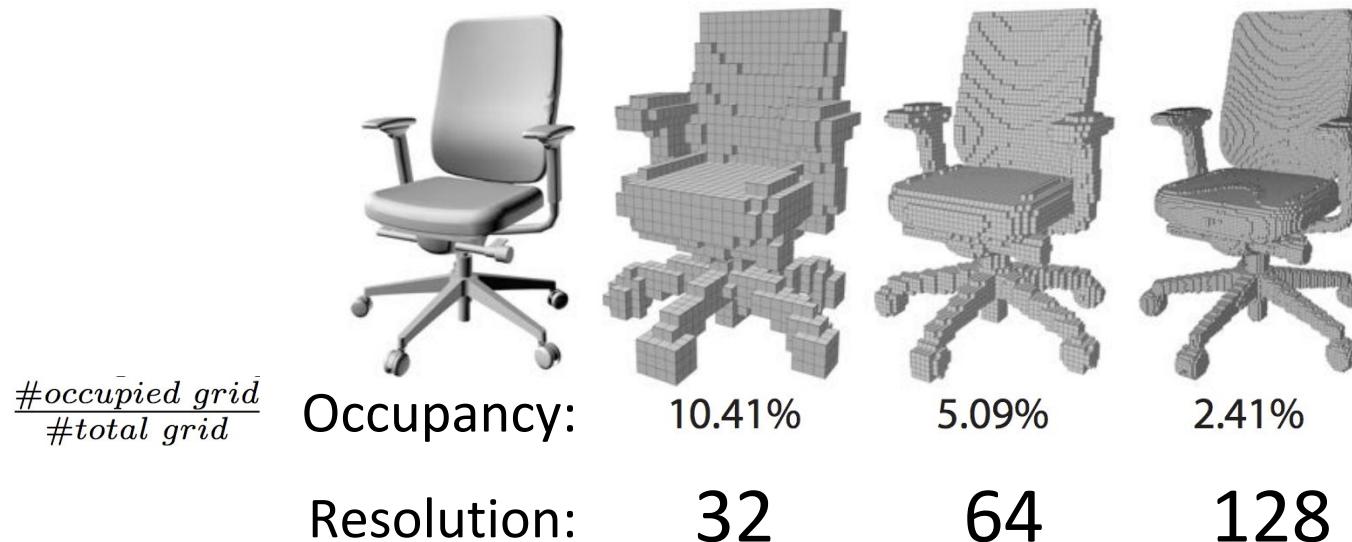
Su et al., “Volumetric and Multi-View CNNs for Object Classification on 3D Data”. CVPR 2016

Many other works in autonomous driving that uses bird's eye view for object detection.

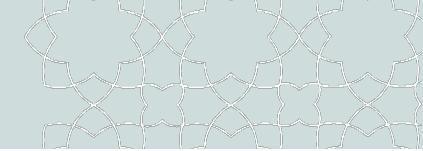
# Volumetric CNN: Voxelization



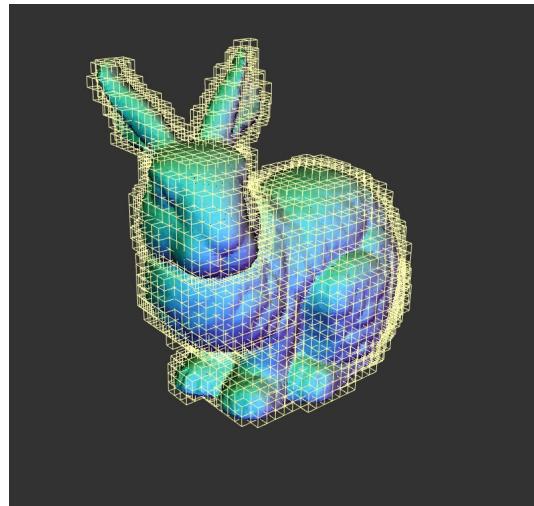
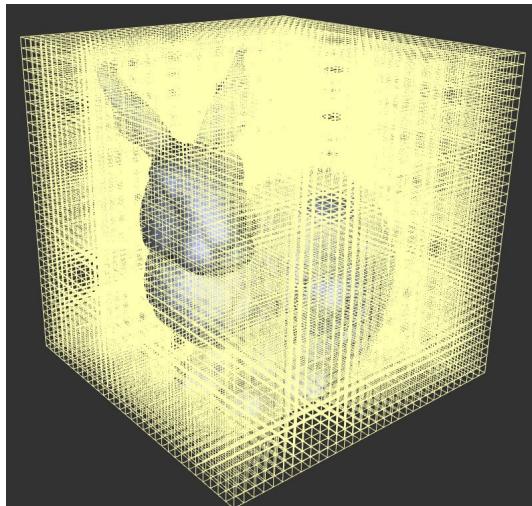
- More Principled: Leverage Sparsity of 3D Surface



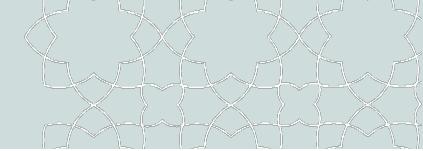
# Volumetric CNN: Voxelization



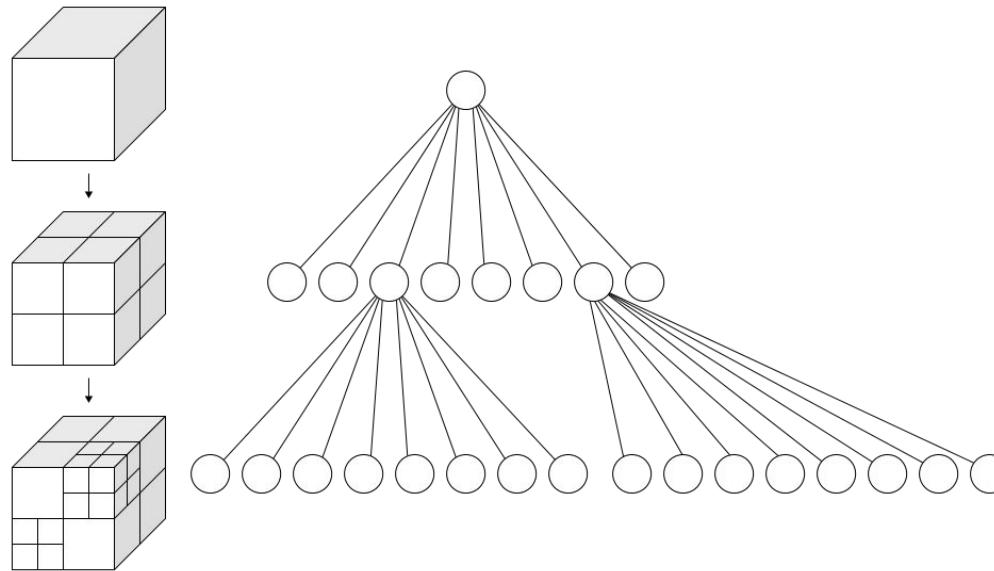
- Store only the Occupied Grids
  - Store the sparse surface signals
  - Constrain the computation near the surface



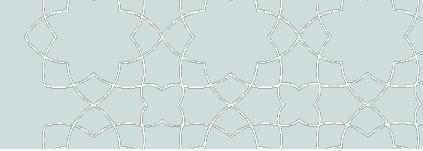
# Volumetric CNN: Voxelization



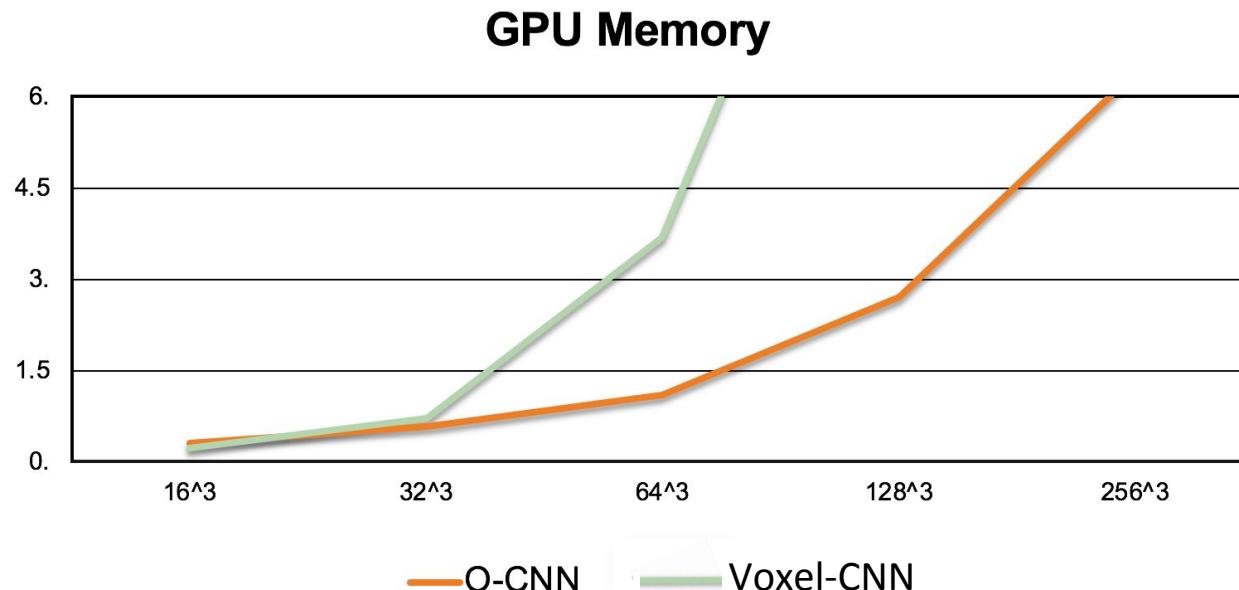
- Octree: Recursively Partition the Space
  - Each internal node has exactly eight children
  - Neighborhood searching: Hash table



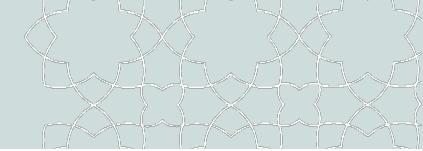
# Volumetric CNN: Voxelization



- Octree: Recursively Partition the Space
  - Memory Efficiency

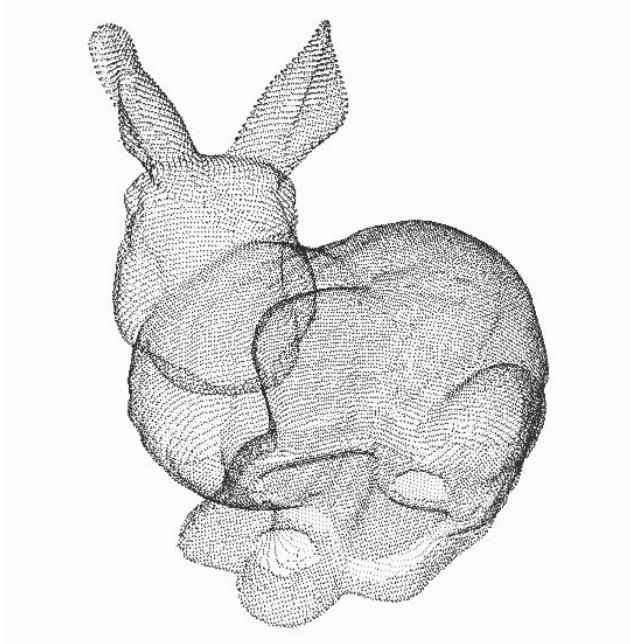
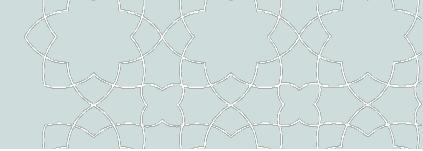


# Volumetric CNN: Voxelization



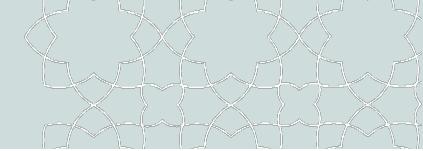
- Implementation
  - SparseConvNet
  - <https://github.com/facebookresearch/SparseConvNet>
  - Uses ResNet architecture
  - State-of-the-art for 3D analysis
  - Takes time to train

# Point Networks

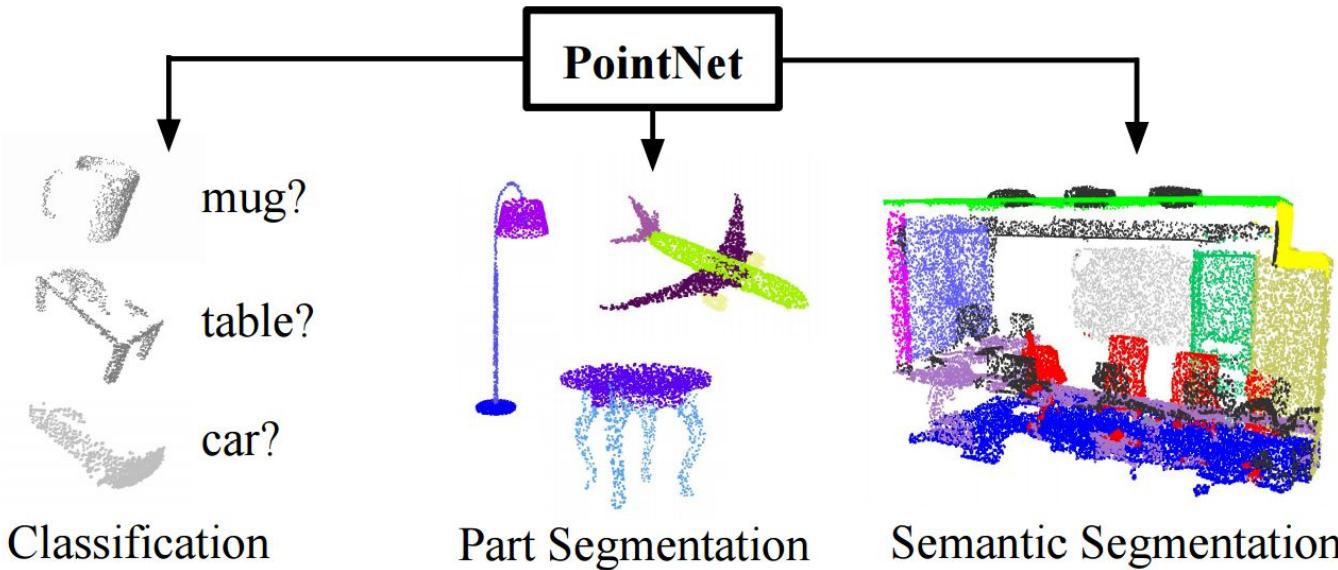


**Point cloud**  
(The most common 3D sensor data)

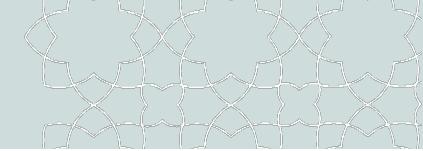
# Point Networks



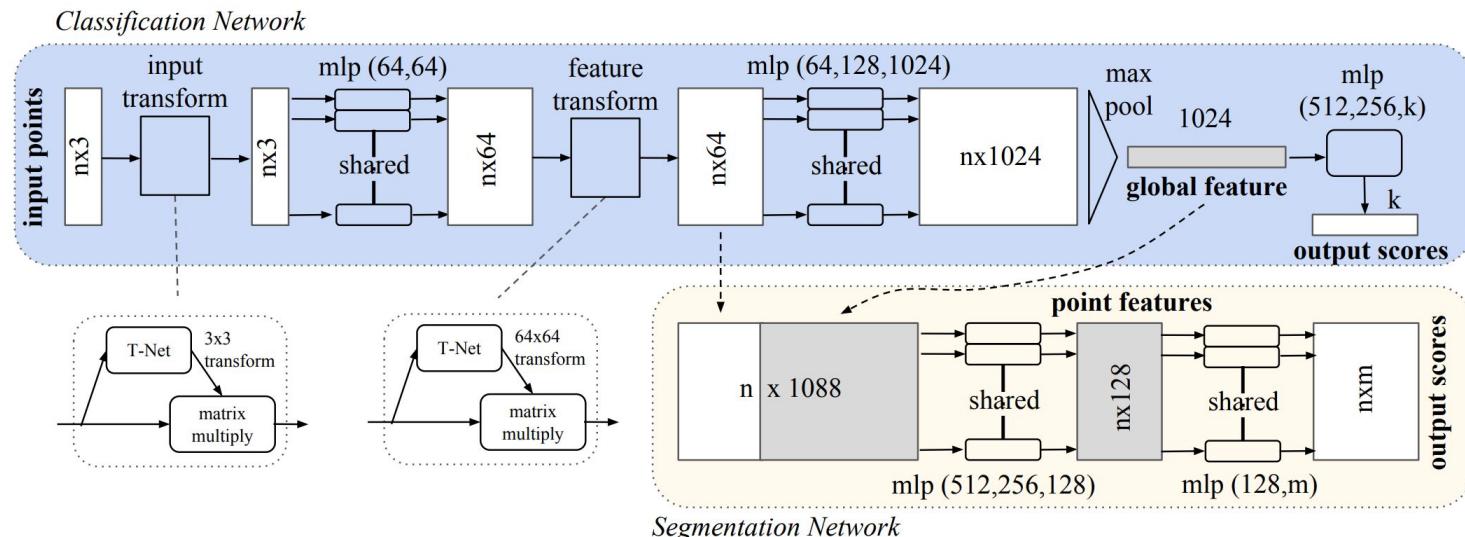
- Directly Process Point Cloud Data
  - End-to-end learning for **unstructured, unordered** point data



# Point Networks



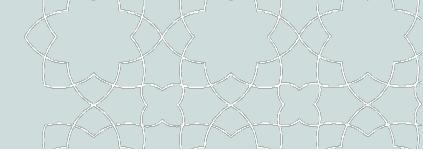
- Directly Process Point Cloud Data
  - **Segmentation:** MLP for each points
  - **Classification:** MaxPooling between point features



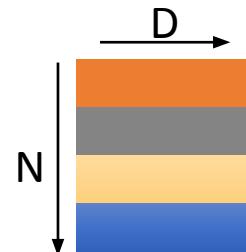
Qi, Charles R., et al. "Pointnet: Deep learning on point sets for 3d classification and segmentation", CVPR 2017

Zaheer, Manzil, et al. "Deep sets", NeurIPS 2017

# Point Networks

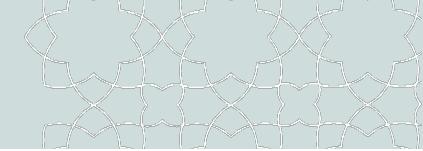


- Permutation invariance:
  - Point cloud:  $N$  **orderless** points, each represented by a  $D$  dim coordinate

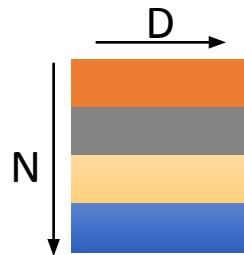


2D array representation

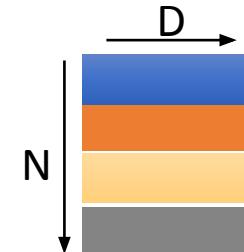
# Point Networks



- Permutation invariance:
  - Point cloud:  $N$  **orderless** points, each represented by a  $D$  dim coordinate

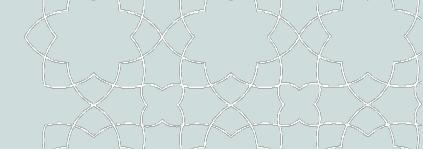


represents the same **set** as



2D array representation

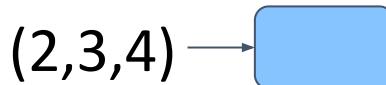
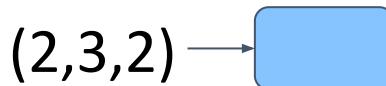
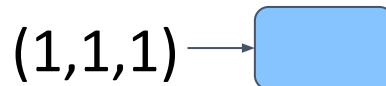
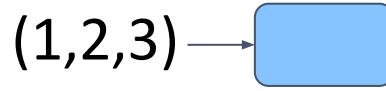
# Point Networks



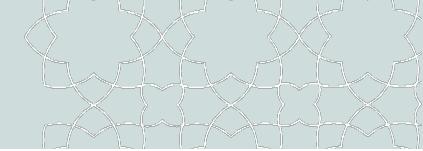
- Construct a Symmetric Function

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

$h$

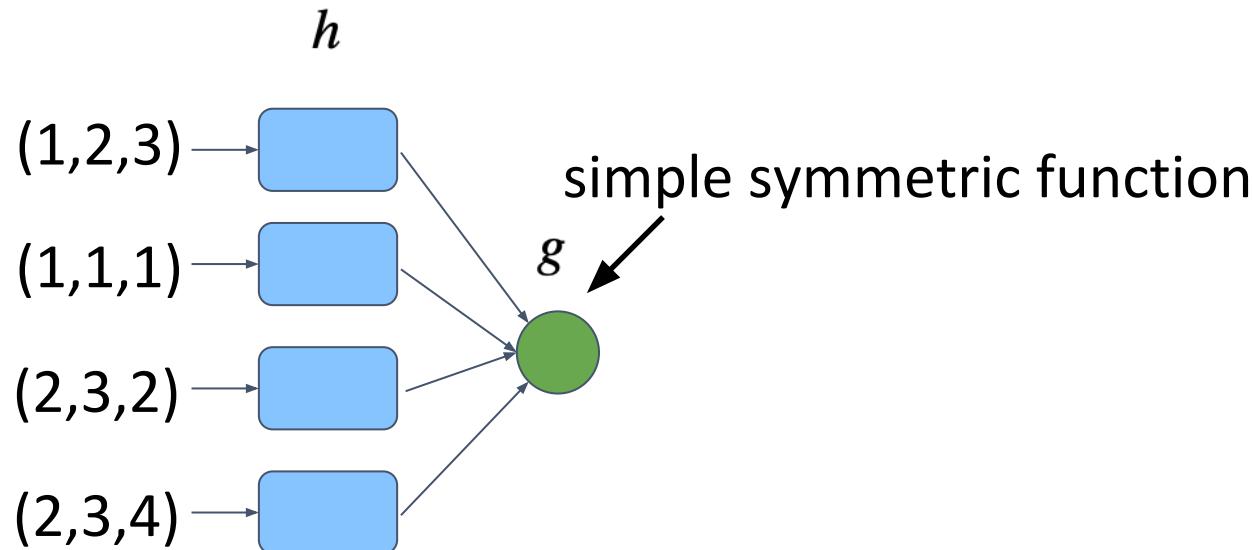


# Point Networks

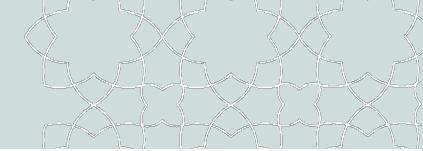


- Construct a Symmetric Function

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

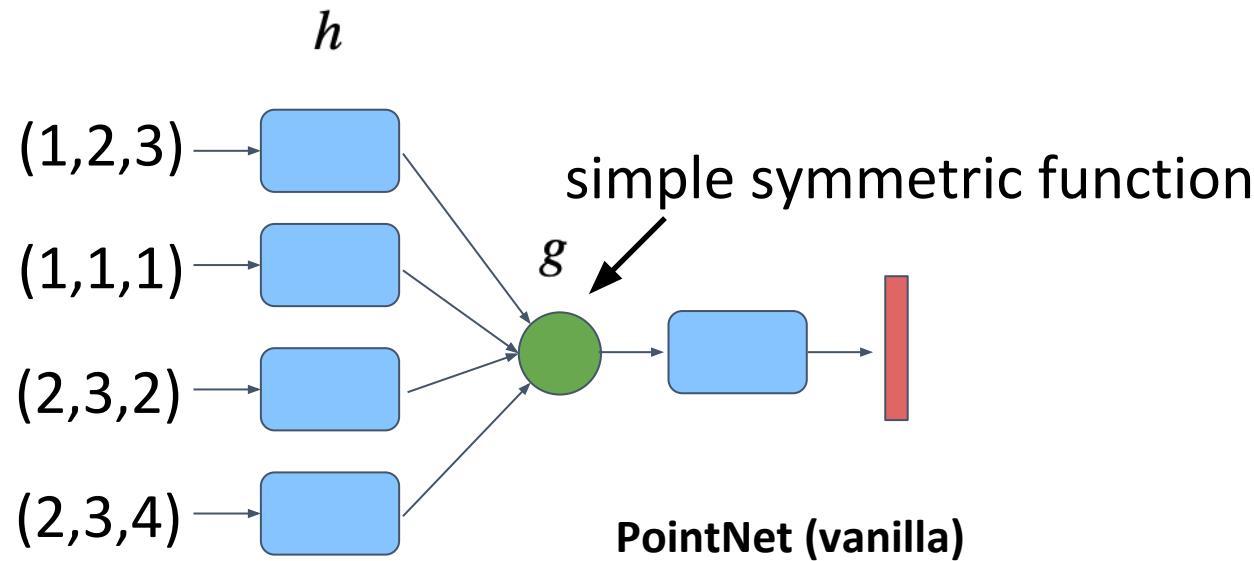


# Point Networks



- Construct a Symmetric Function

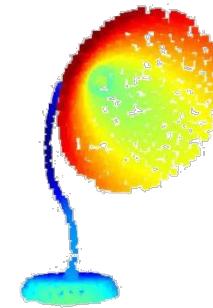
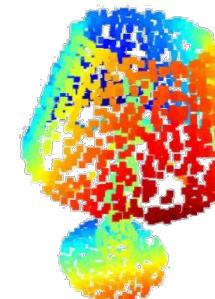
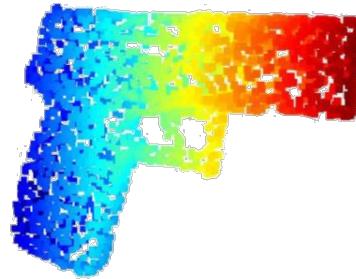
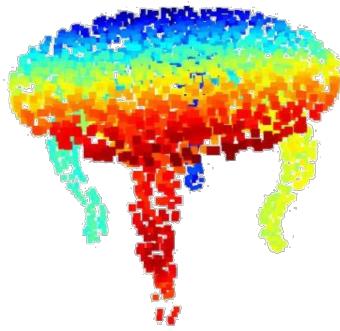
$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$  is symmetric if  $g$  is symmetric



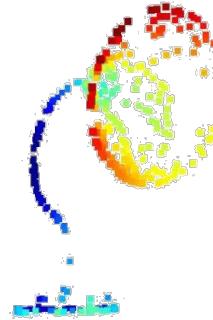
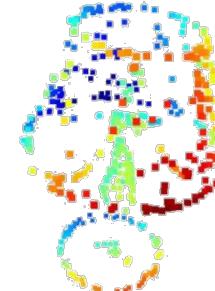
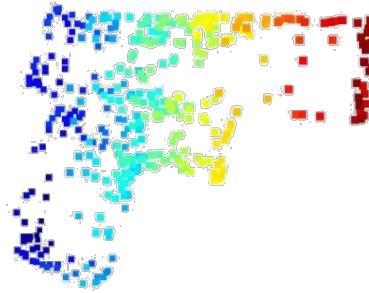
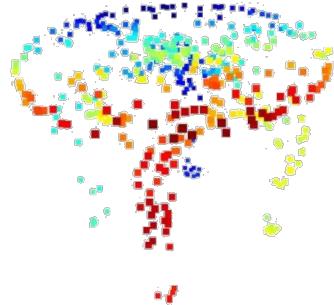
# Point Networks

- Salient points are discovered!

Original Shape



Critical Point Sets

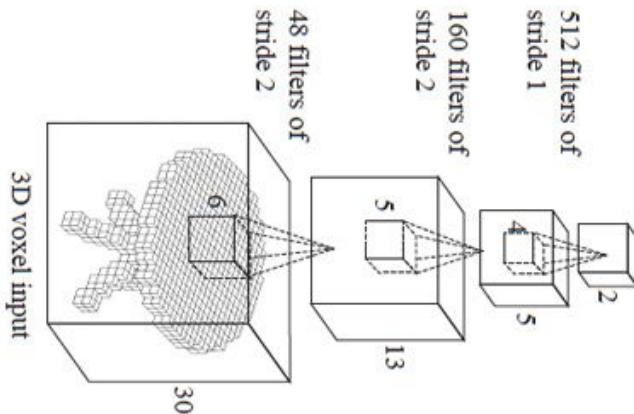


# Point Networks: Limitations

Hierarchical feature learning  
Multiple levels of abstraction

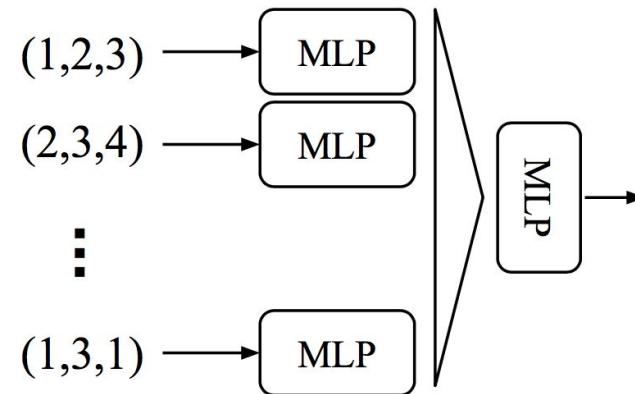
v.s.

Global feature learning  
Either one point or all points



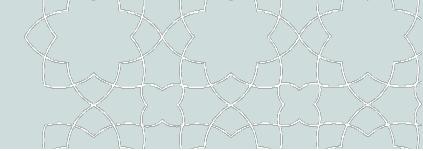
**3D CNN (Wu et al.)**

No local context for each point!  
Global feature depends on absolute coordinate.  
Hard to generalize to unseen scene configurations!



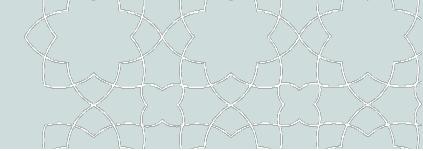
**PointNet (vanilla) (Qi et al.)**

# Point Networks

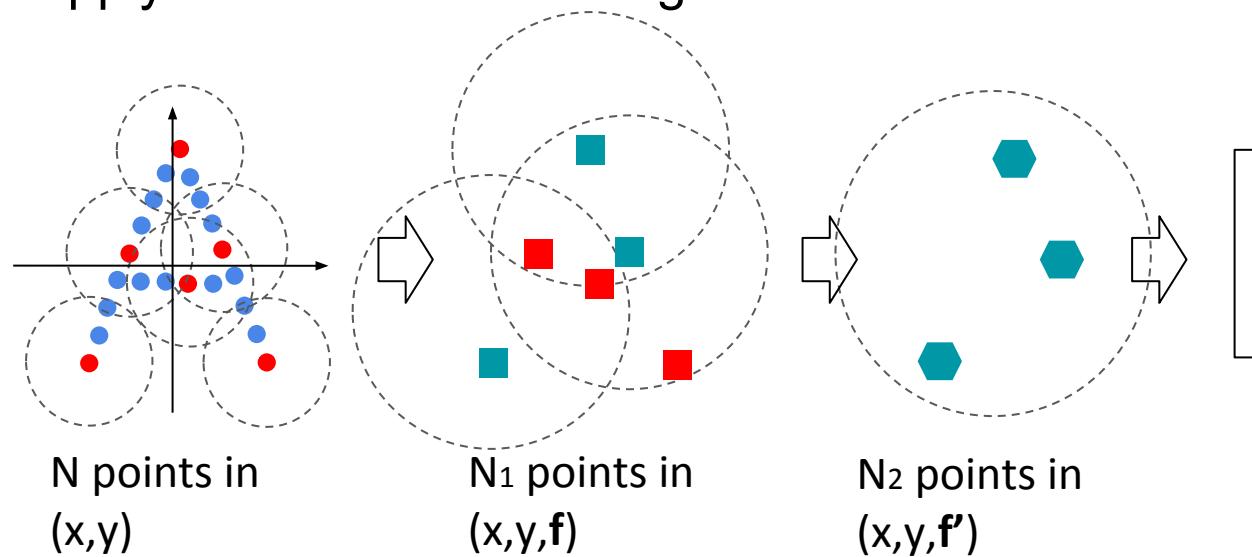


- Points in Metric Space:
  - Learn “kernels” in 3D space and conduct convolution
  - Kernels have compact spatial support
  - For convolution, we need to find neighboring points
  - Possible strategies for range query
    - Ball query (results in more stable features)
    - k-NN query (faster)

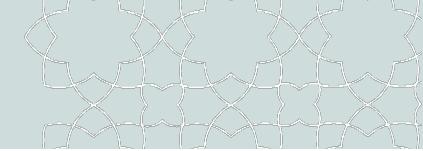
# PointNet++: Multi-Scale PointNet



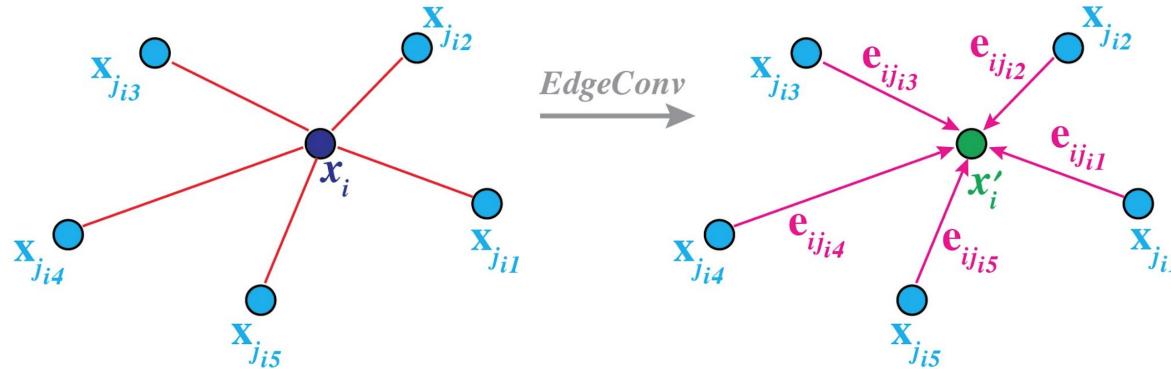
- Repeat
  - Sample anchor points
  - Find neighborhood of anchor points
  - Apply PointNet in each neighborhood to mimic convolution



# Point Convolution As Graph Convolution

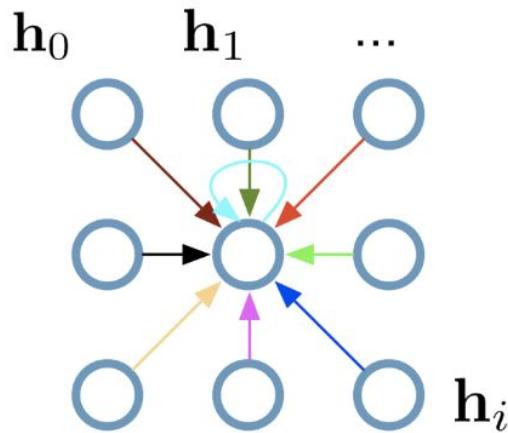


- Points -> Nodes
- Neighborhood -> Edges
- Graph CNN for point cloud processing



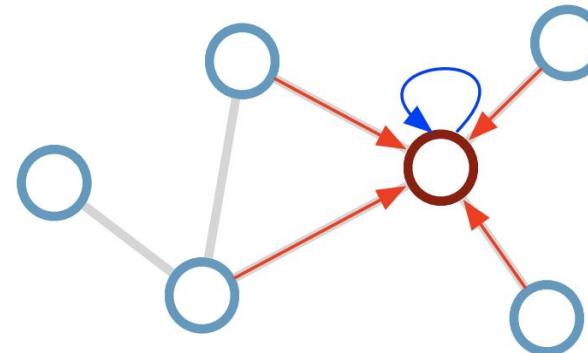
# Point Convolution As Graph Convolution

Convolutional Neural Network (CNN)



$$\mathbf{h}_4^{(l+1)} = \sigma \left( \mathbf{W}_0^{(l)} \mathbf{h}_0^{(l)} + \mathbf{W}_1^{(l)} \mathbf{h}_1^{(l)} + \cdots + \mathbf{W}_8^{(l)} \mathbf{h}_8^{(l)} \right)$$

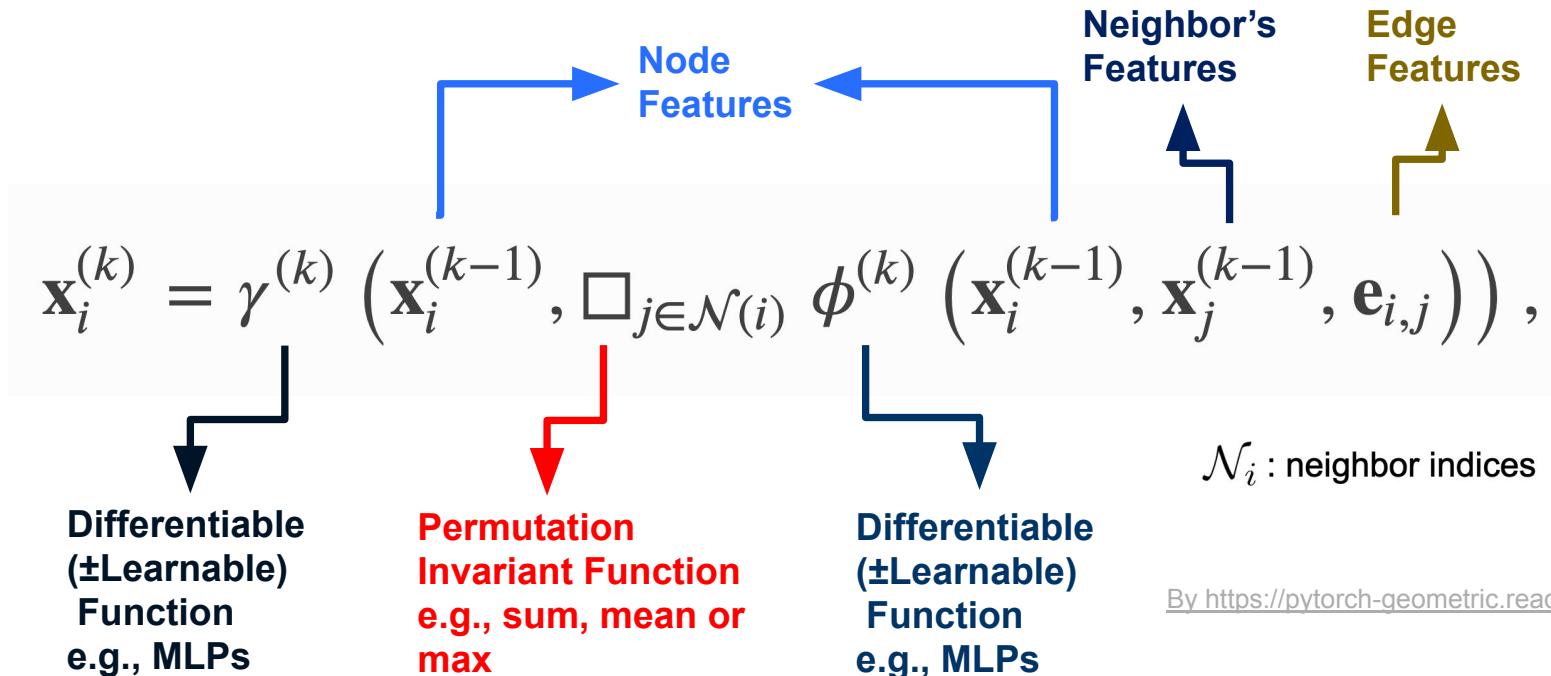
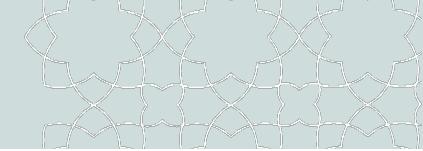
Graph Convolutional Network (GCN)



$$\mathbf{h}_i^{(l+1)} = \sigma \left( \mathbf{h}_i^{(l)} \mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} \mathbf{h}_j^{(l)} \mathbf{W}_1^{(l)} \right)$$

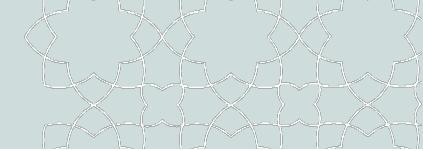
$\mathcal{N}_i$ : neighbor indices     $c_{ij}$ : norm. constant  
(fixed/trainable)

# Point Convolution As Graph Convolution

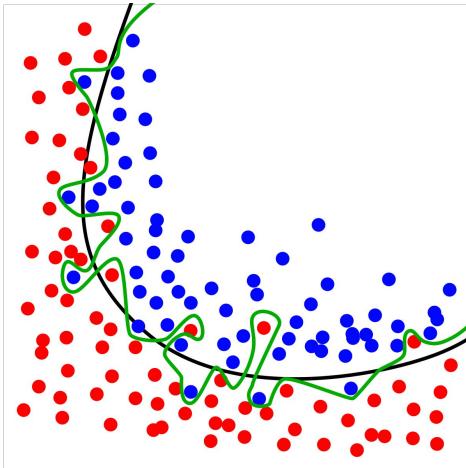


By <https://pytorch-geometric.readthedocs.io>

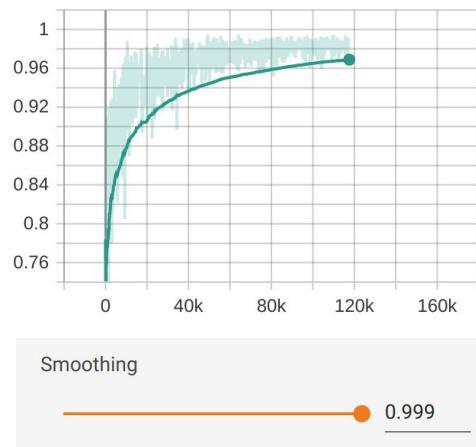
# Point Convolution As Graph Convolution



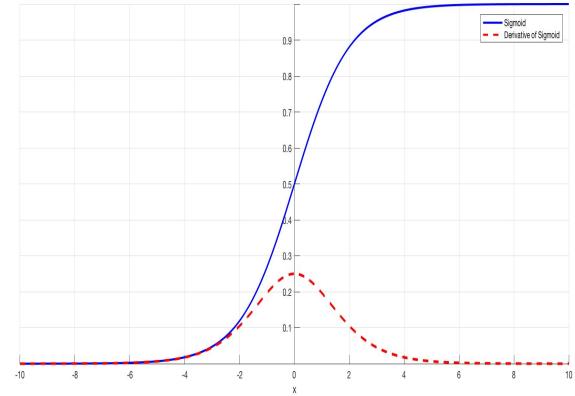
- GCN: Main Limitations



Over-fitting



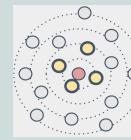
Over-smoothing



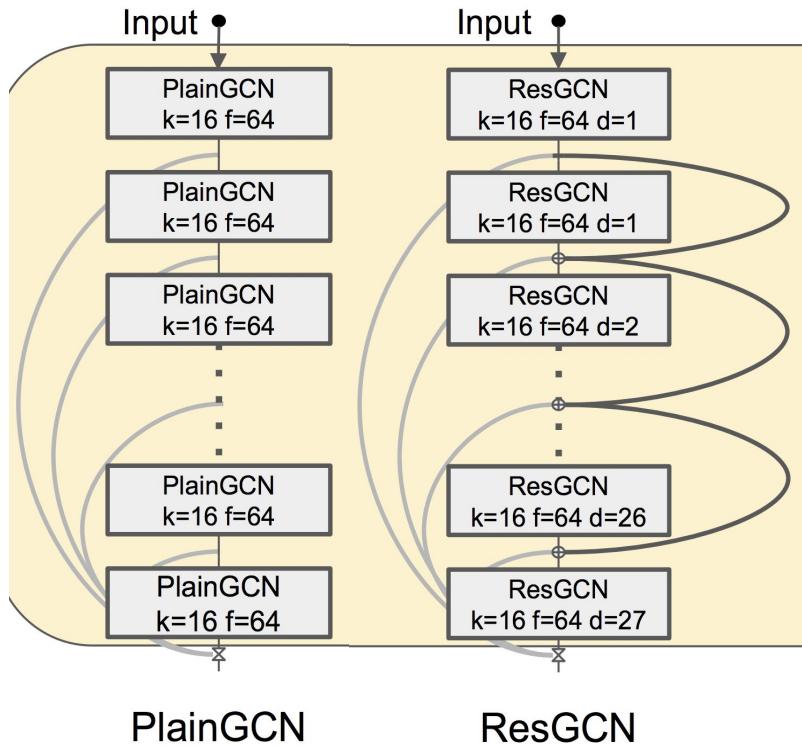
Vanishing Gradient

Figures from <https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>

# deepGCN: Residual Connections

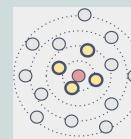


DeepGCNs.org

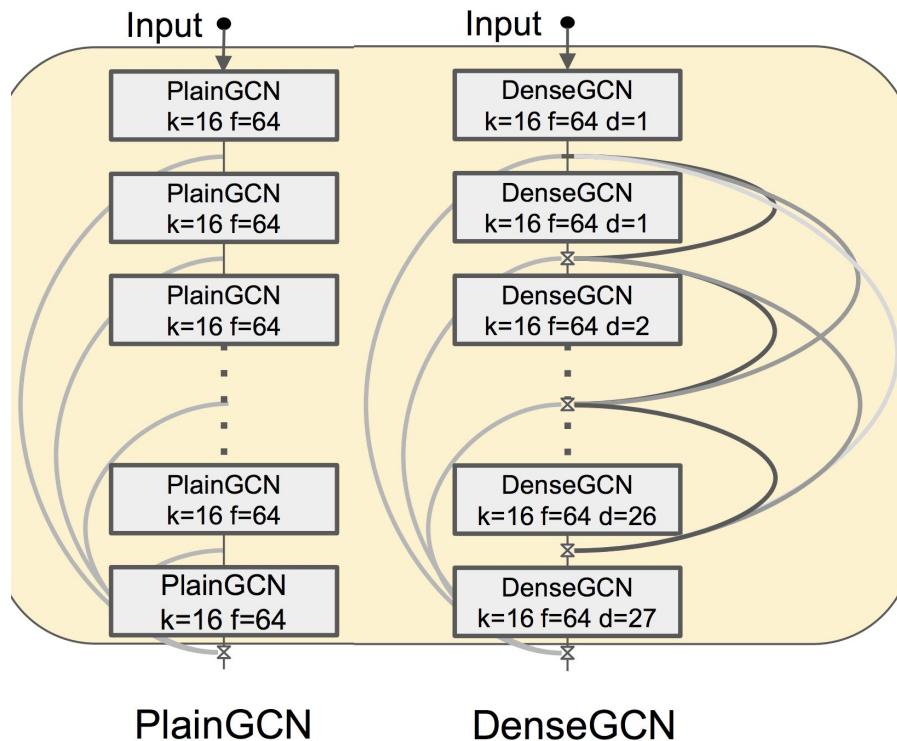


$$\begin{aligned}\mathcal{G}_{l+1} &= \mathcal{H}(\mathcal{G}_l, \mathcal{W}_l) \\ &= \mathcal{F}(\mathcal{G}_l, \mathcal{W}_l) + \mathcal{G}_l.\end{aligned}$$

# deepGCN: Dense Connections

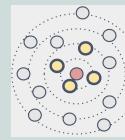


DeepGCNs.org

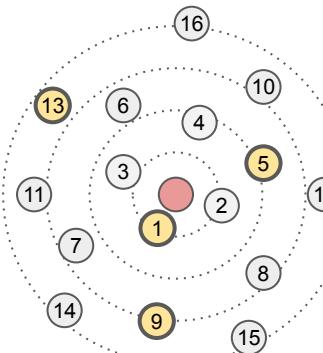
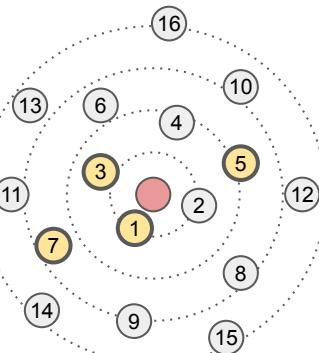
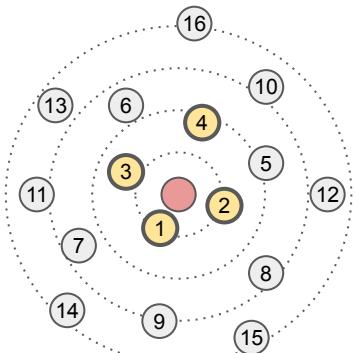
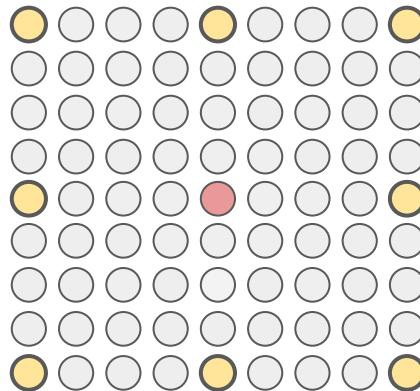
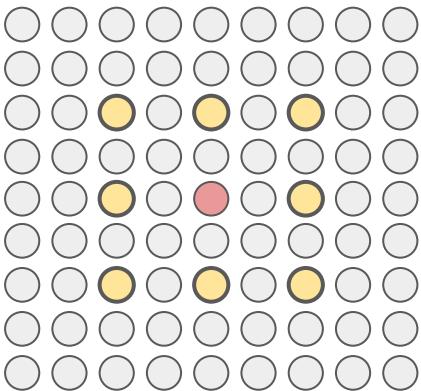
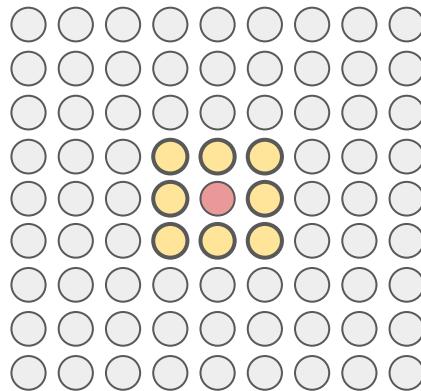


$$\begin{aligned}\mathcal{G}_{l+1} &= \mathcal{H}(\mathcal{G}_l, \mathcal{W}_l) \\ &= \mathcal{T}(\mathcal{F}(\mathcal{G}_l, \mathcal{W}_l), \mathcal{G}_l) \\ &= \mathcal{T}(\mathcal{F}(\mathcal{G}_l, \mathcal{W}_l), \dots, \mathcal{F}(\mathcal{G}_0, \mathcal{W}_0), \mathcal{G}_0).\end{aligned}$$

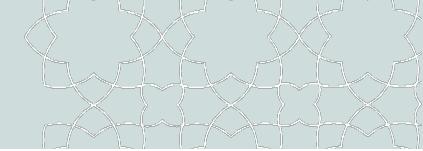
# deepGCN: Dilated GCN



DeepGCNs.org



# Standard GCNs are not Geometry-Aware



- Note that points are **sampled** from surfaces
- Ideally, features describe the geometry of underlying surface.  
Should be sample invariant
- But GCNs lack design to address sample invariance
- Remind us “density estimation” from a population
- Rescue: Estimate the continuous kernel and point density for continuous convolution

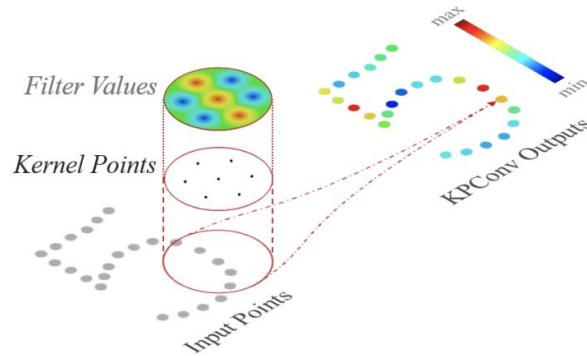
# Mathematically Proper Conv. Discretization

- Continuous conv:
- Empirical conv:
- Learn kernel value at anchor points and interpolate to build continuous kernel

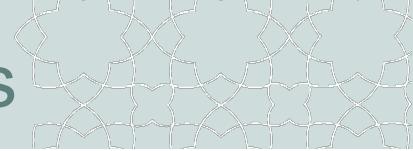
$$\kappa_{jm}(z) = \sum_l k_{ljm} \Phi(|z - y_l|)$$

$\Phi$ : RBF kernel

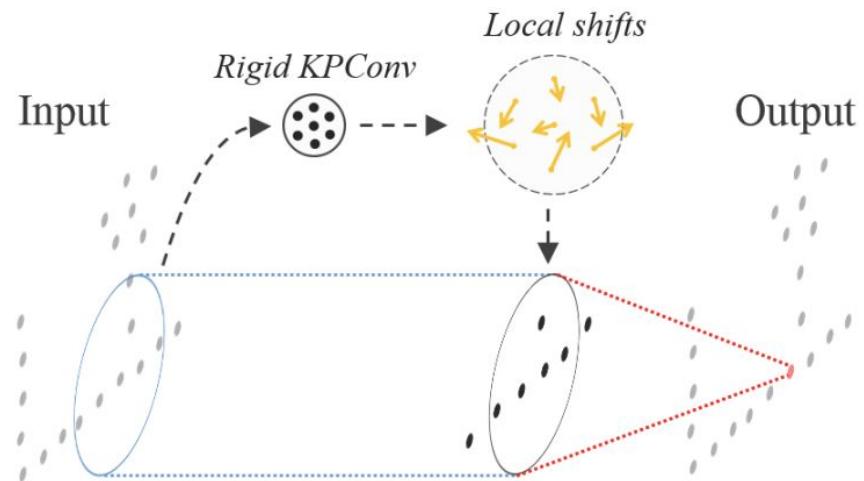
$$(\mathcal{F} * g)(x) = \int g(y - x)f(y)dy$$
$$(\mathcal{F} * g)(x) = \sum_{x_i \in \mathcal{N}_x} g(x_i - x)f_i$$



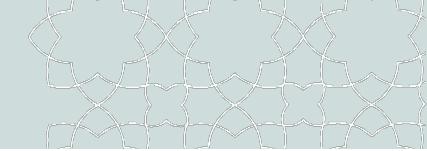
# Deformable Kernel for Deformable Objects



- Deformable point-based kernel
- The 3D version of 2D deformable convolution

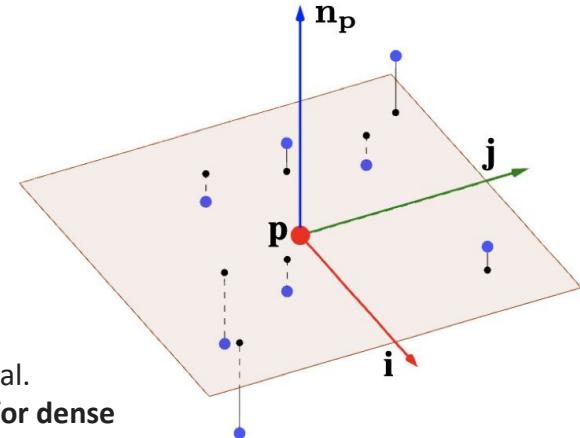


# Other Ways to Address 3D Conv

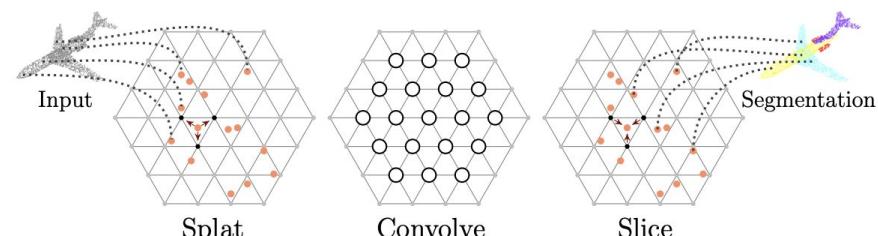


- **Tangent convolution:**
  - Project & interpolate local points to the tangent plane (PCA for orientation)
- **Lattice**
  - high-dimensional space

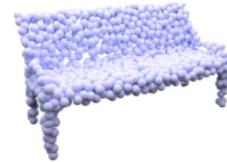
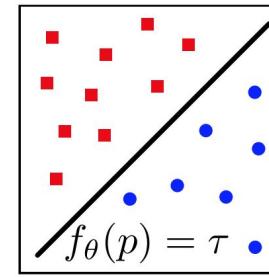
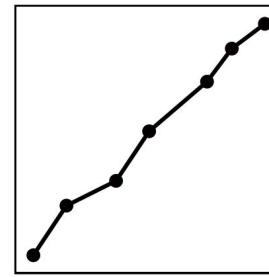
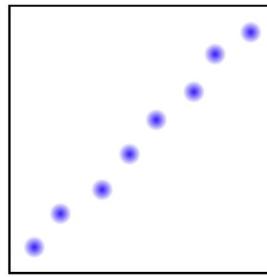
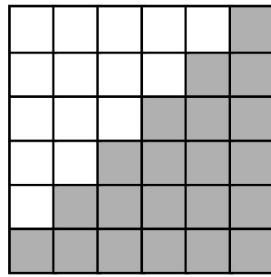
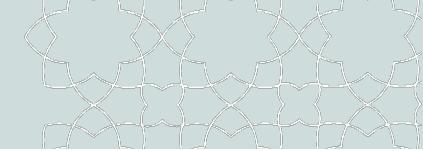
Tatarchenko, Maxim, et al.  
"Tangent convolutions for dense  
prediction in 3d", CVPR 2018.



Su, Hang, et al. "Splatnet: Sparse lattice networks for point cloud processing", CVPR 2018



# Implicit Representation?



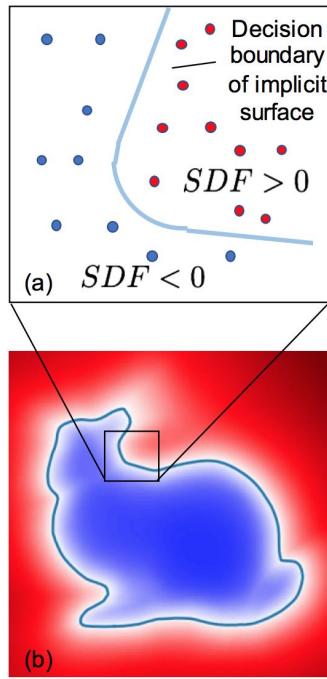
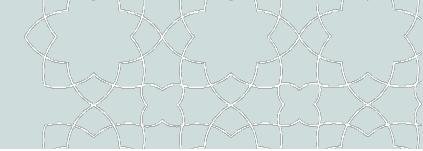
$$F(x) = 0$$

Mescheder et al., “Occupancy Networks: Learning 3D Reconstruction in Function Space”, CVPR 2019

Park et al., “DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation”, CVPR 2019

Chen et al., “Learning Implicit Fields for Generative Shape Modeling”, CVPR 2019

# Implicit Representation?



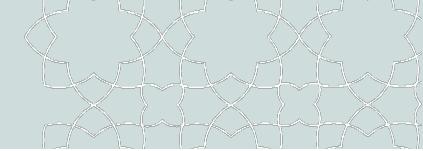
$$F(x) = 0$$

Mescheder et al., "Occupancy Networks: Learning 3D Reconstruction in Function Space", CVPR 2019

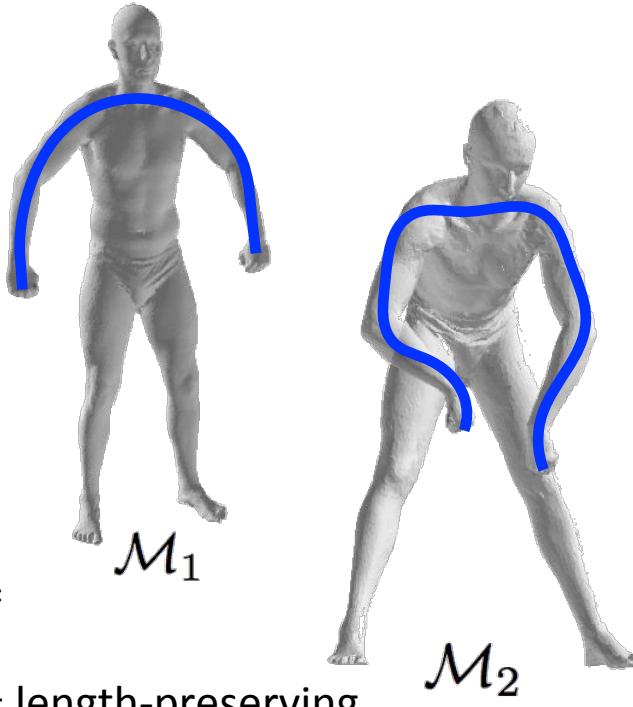
Park et al., "DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation", CVPR 2019

Chen et al., "Learning Implicit Fields for Generative Shape Modeling", CVPR 2019

# Spectral Convolution



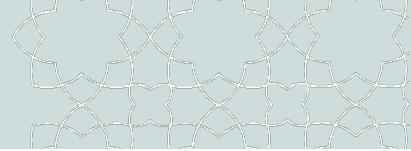
- Recognition with Isometric Invariance?



geodesic =  
intrinsic  
isometry = length-preserving  
 $f$

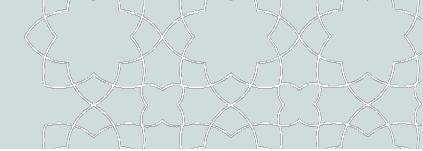


# Spectral Convolution



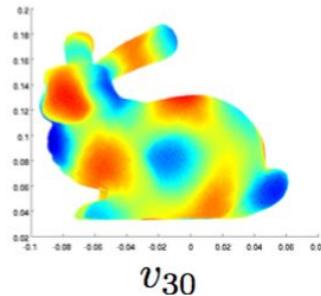
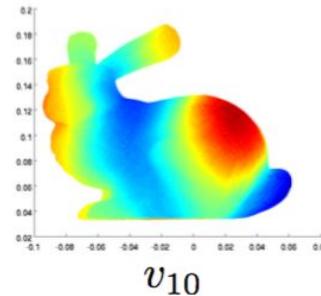
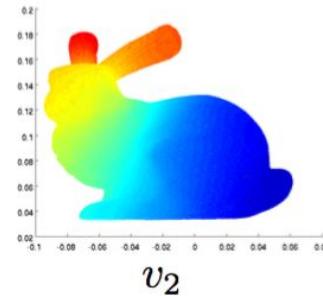
- Convolution done in the spectral domain
- Kernels are also built in spectral domain
- Activation done in the spatial domain

# Spectral Convolution

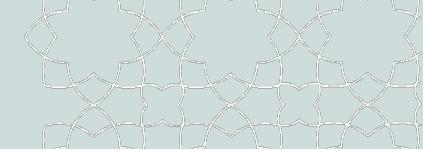


- Challenge: Obtain Fourier Basis
  - Derived by eigenfunctions of self-adjoint operators, e.g. Laplacian-Bertrami or Dirac operator

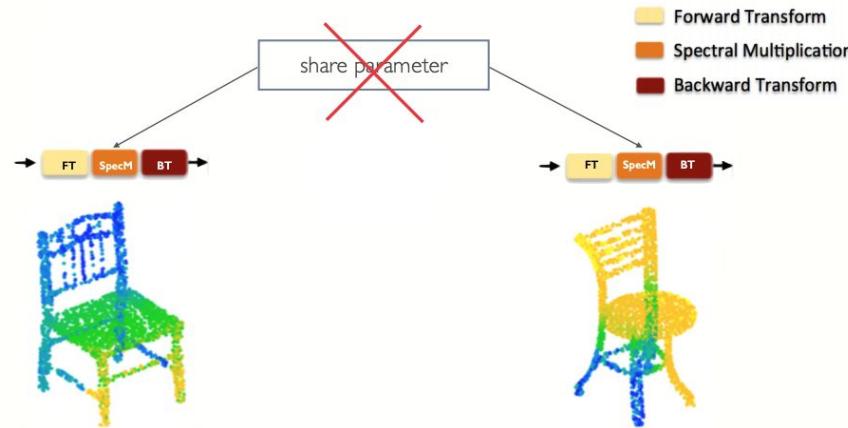
“Fourier basis” of the graph:  $V$  : Eigenvectors of  $\Delta$



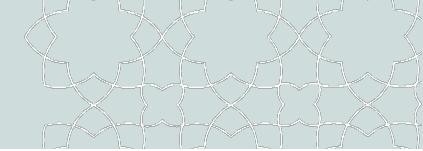
# Spectral Convolution



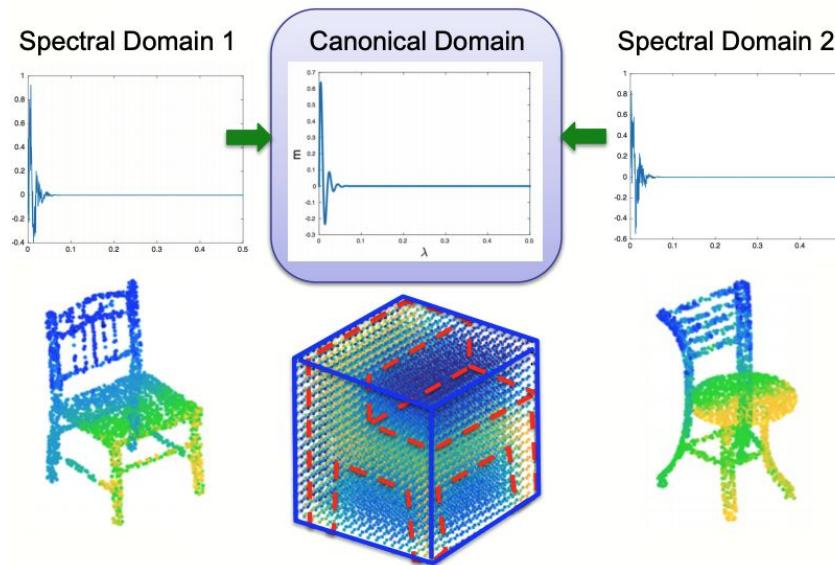
- Fundamental Challenge of Spectral CNN
  - If the shapes to compare are not isometric, their spectral domains are not aligned
  - Function bases are derived by Laplacian operator, which is geometry dependent



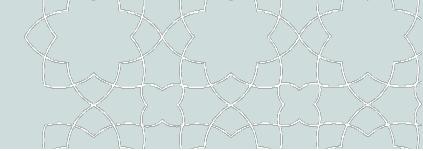
# Spectral Convolution



- Domain Synchronization
  - Rescue: Synchronize spectral domains by functional maps

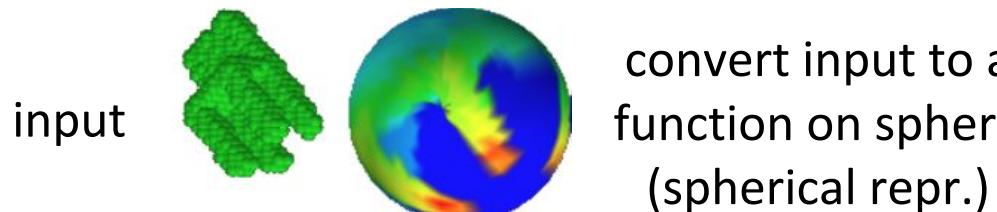


# Spectral Convolution



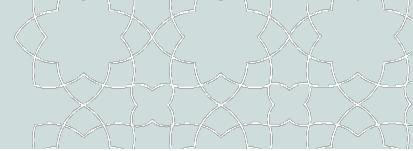
- A Special Case: Spherical CNN

- If the surface is always a SPHERE, no worry about the functional space alignment anymore
- Generate a spherical representation



- Do Spherical CNN
  - Has numerical tricks exploiting the symmetry of sphere

# Spectral Convolution



- Learned Filters

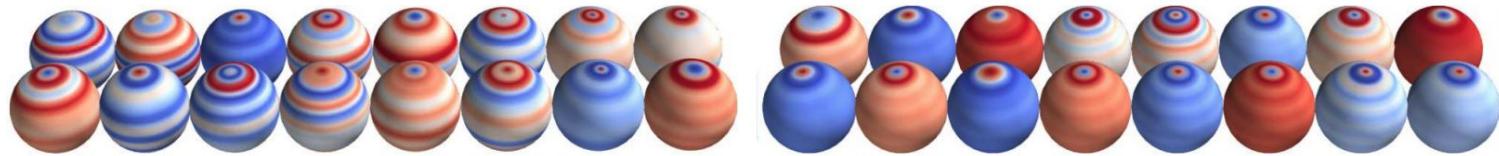
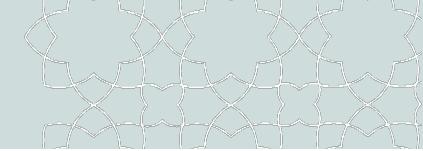


Fig. 4: Filters learned in the first layer. The filters are zonal. *Left*: 16 nonlocalized filters. *Right*: 16 localized filters. Nonlocalized filters are parameterized by all spectral coefficients (16, in the example). Even though locality is not enforced, some filters learn to respond locally. Localized filters are parameterized by a few points of the spectrum (4, in the example), the rest of the spectrum is obtained by interpolation.

# Spectral Convolution



- Rotation invariance guaranteed
- Can be used to improve the rot. invariance of MVCNN, as well

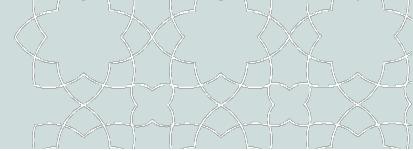
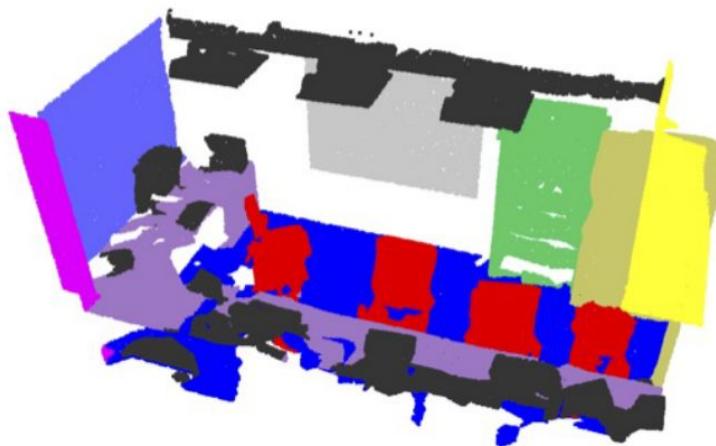
Table 1: ModelNet40 classification accuracy per instance. Spherical CNNs are robust to arbitrary rotations, even when not seen during training, while also having one order of magnitude fewer parameters and faster training.

Method	z/z	SO3/SO3	z/SO3	params	inp. size
PointNet [7]	89.2	83.6	14.7	3.5M	2048 x 3
PointNet++ [38]	89.3	85.0	28.6	1.7M	1024 x 3
VoxNet [29]	83.0	73.0	-	0.9M	$30^3$
SubVolSup [8]	88.5	82.7	36.6	17M	$30^3$
SubVolSup MO [8]	89.5	85.0	45.5	17M	$20 \times 30^3$
MVCNN 12x [9]	89.5	77.6	70.1	99M	$12 \times 224^2$
MVCNN 80x [9]	<b>90.2</b>	86.0	- <sup>2</sup>	99M	$80 \times 224^2$
RotationNet 20x [30]	<b>92.4</b>	80.0	20.2	58.9M	$20 \times 224^2$
Ours	88.9	<b>86.9</b>	<b>78.6</b>	<b>0.5M</b>	<b><math>2 \times 64^2</math></b>

# Outline

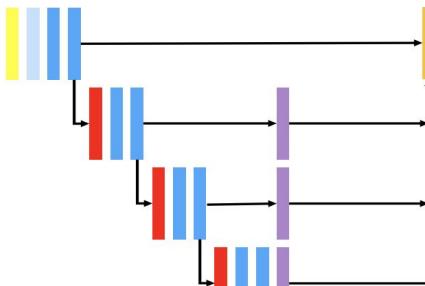
- 3D Data
- Classification
- **Segmentation and Detection**
- 3D Data synthesis

# Task: 3D Semantic Segmentation

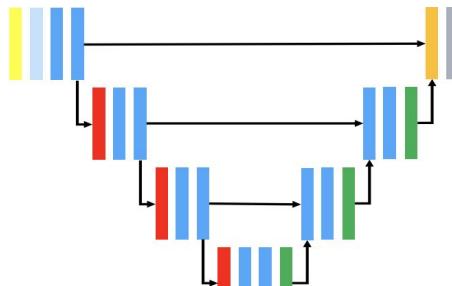


# Task: 3D Semantic Segmentation

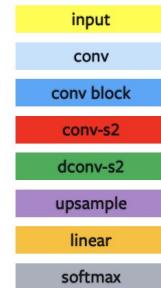
- Encoder – Decoder



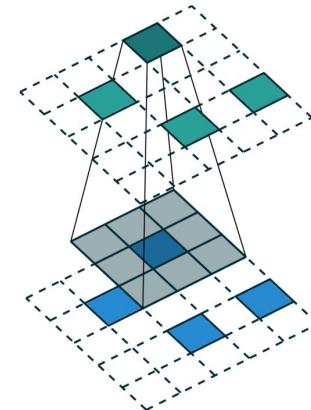
(a) Submanifold sparse FCN.



(b) Submanifold sparse U-Net.



Sparse Convolution

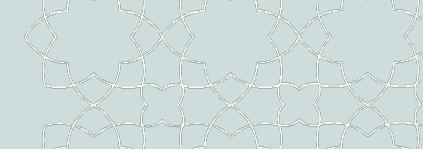


Sparse Deconvolution

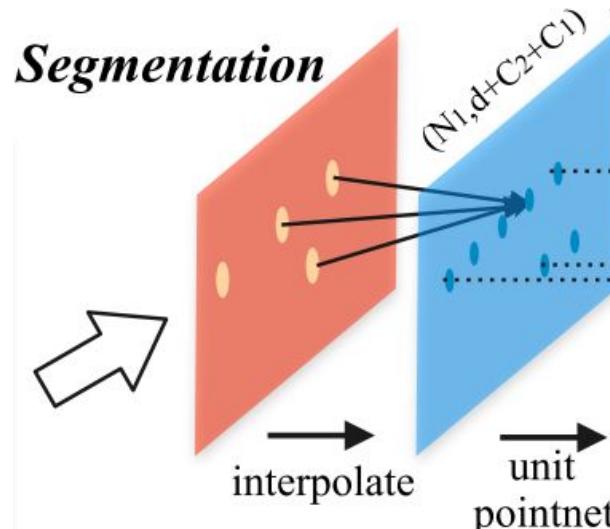
Graham, Benjamin, Martin Engelke, and Laurens van der Maaten. "3d semantic segmentation with submanifold sparse convolutional networks." CVPR 2018.

Choy, Christopher, et al. "4D Spatio-Temporal Convnets: Minkowski Convolutional Neural Networks." CVPR 2019

# Task: 3D Semantic Segmentation



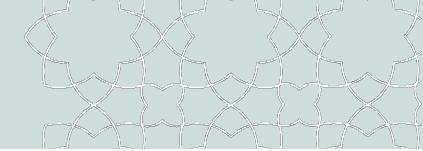
- Decoder



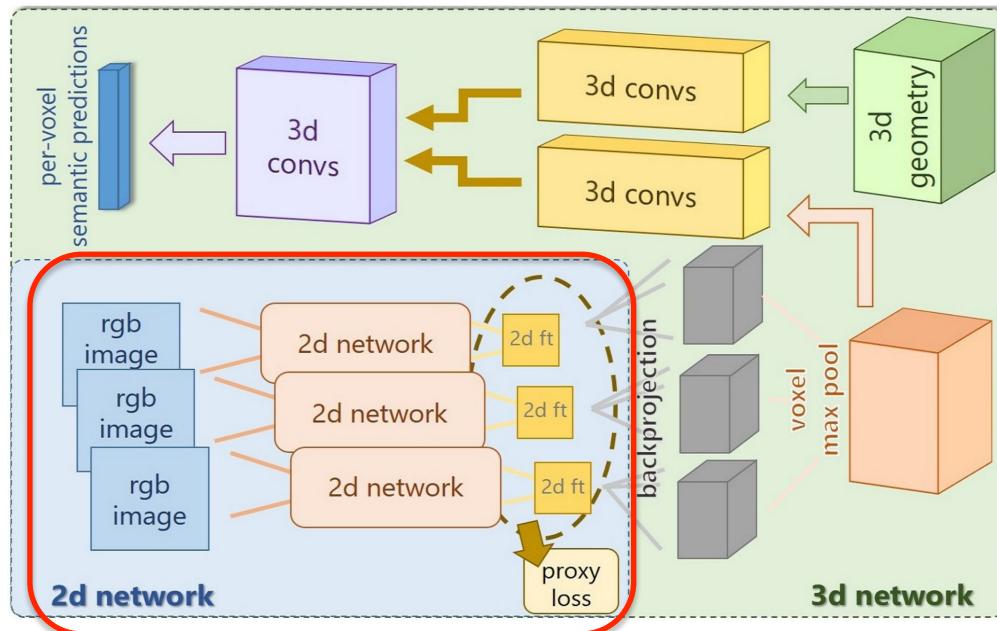
**Upsample** point cloud features by interpolating from 3 nearest points in the point cloud of lower resolution.

Lower resolution      Higher resolution  
Fewer points      More points

# Task: 3D Semantic Segmentation



- Multimodal

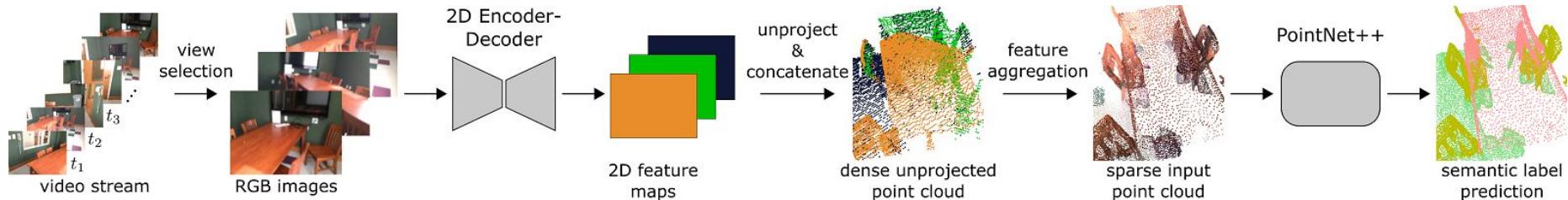


Backproject 2D features to 3D voxels

Apply voxel-wise max-pooling across multiple views

Fuse 2D and 3D features at the intermediate level

# Task: 3D Semantic Segmentation

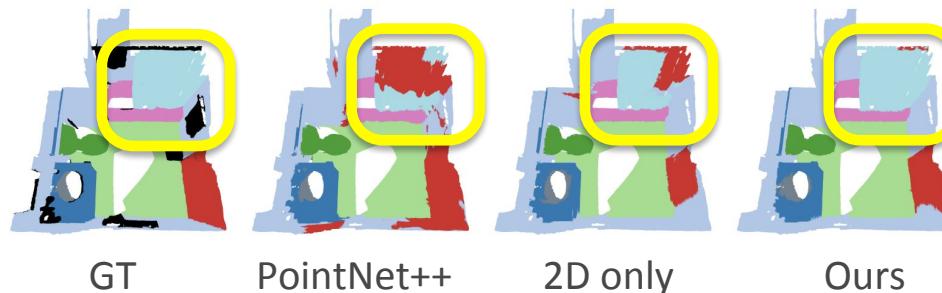


Backproject 2D features to 3D points

Apply PointNet to aggregate multi-view features

Fuse 2D and 3D features as input to 3D network

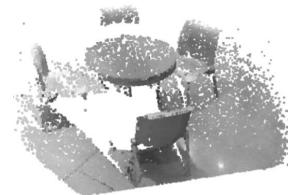
ignore ■ wall ■ floor ■ cabinet ■ bed ■ chair ■ sofa ■ table ■ door ■ window ■ bookshelf ■ picture ■ counter ■ desk ■ curtain ■ refrigerator ■ shower curtain ■ toilet ■ sink ■ bathtub ■ otherfurniture ■



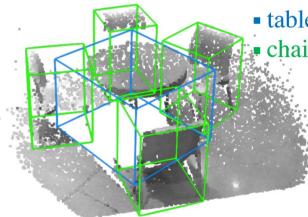
# Task: Instance-level Understanding

Object Detection

Input



Output



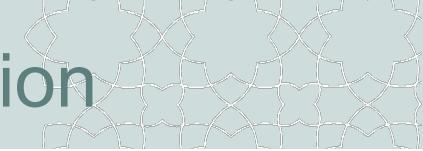
Object Segmentation



Part Segmentation



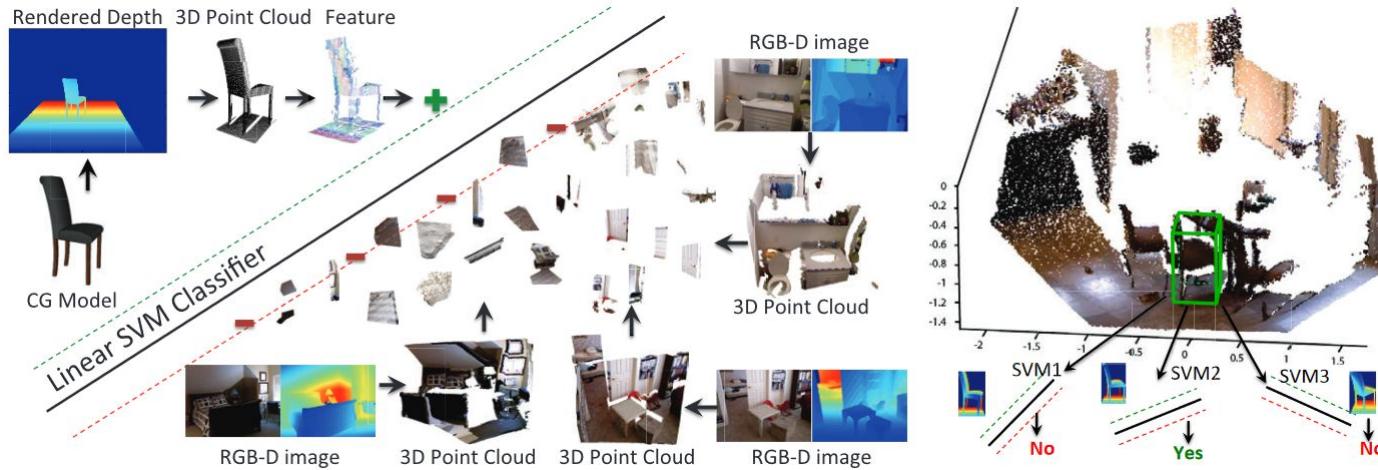
# Task: 3D Detection & Instance Segmentation



- Top-down Methods
  - Does this point cloud contain an object?
  - How to select point clouds to classify objectness?

# Task: 3D Detection & Instance Segmentation

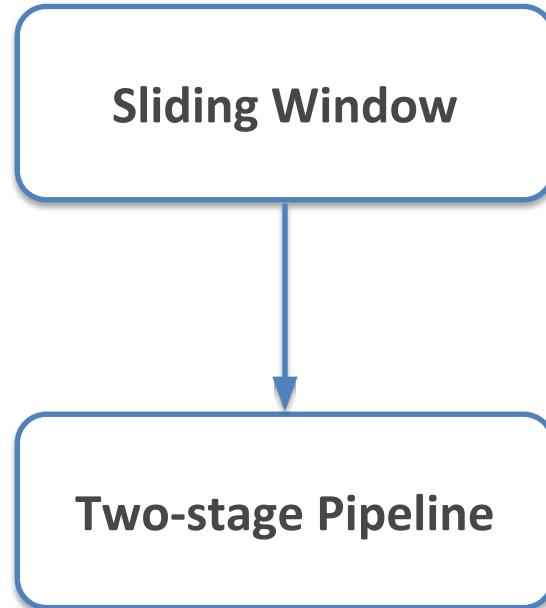
- Sliding Shapes



Sliding window to walk over the entire space

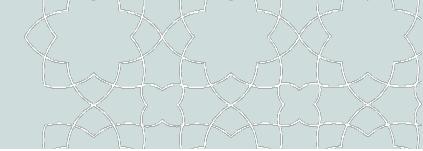
**Expensive !**

# Task: 3D Detection & Instance Segmentation

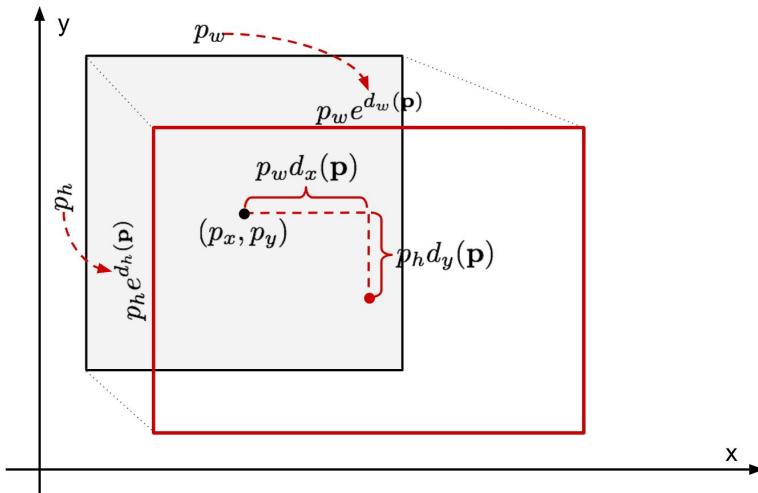


First stage: Proposal  
Second stage: Refinement

# Localization



- 3D bounding box regression
- Predefined bounding boxes (anchors) at each sliding window
- Re-parameterize as relative offsets



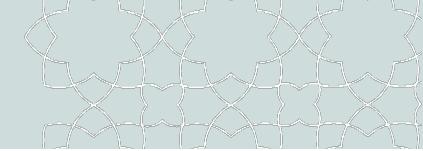
$$\Delta_x = \frac{\mu - \mu_{anchor}}{\phi_{anchor}}$$

Center shift

$$\Delta_w = \ln\left(\frac{\phi}{\phi_{anchor}}\right)$$

Size scale

# From Box to Instance Segmentation



- Box includes background points or other instances
- foreground/background segmentation

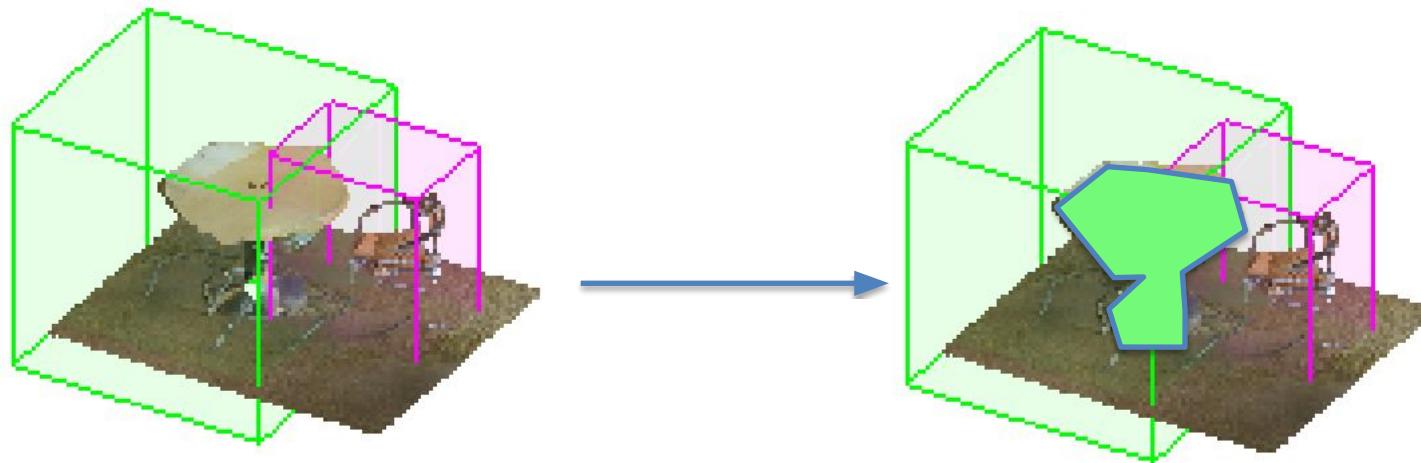


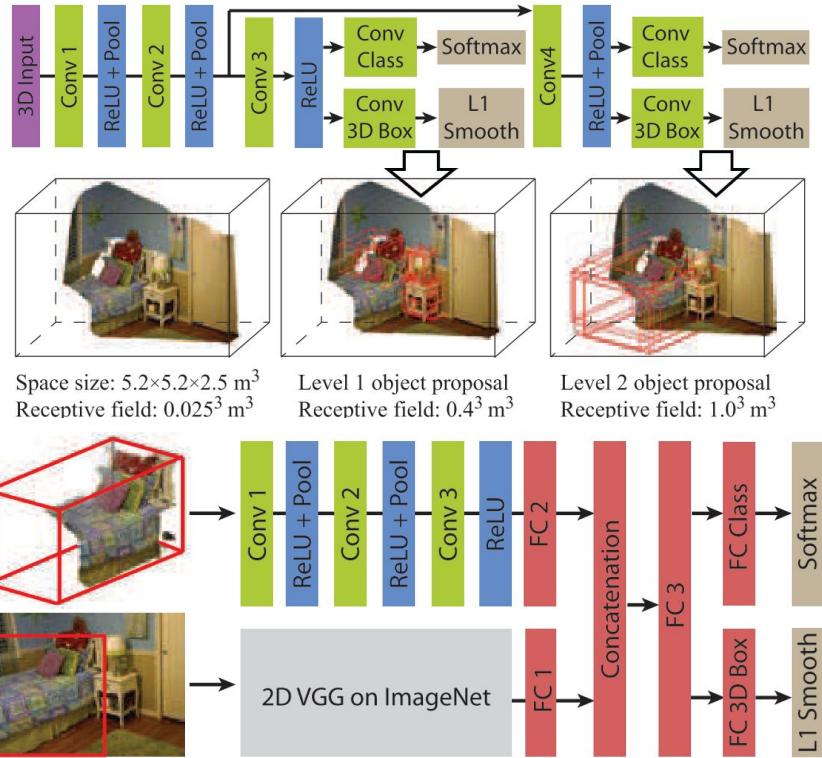
Figure from “Learning Object Bounding Boxes for 3D Instance Segmentation on Point Clouds”

# Volumetric R-CNN

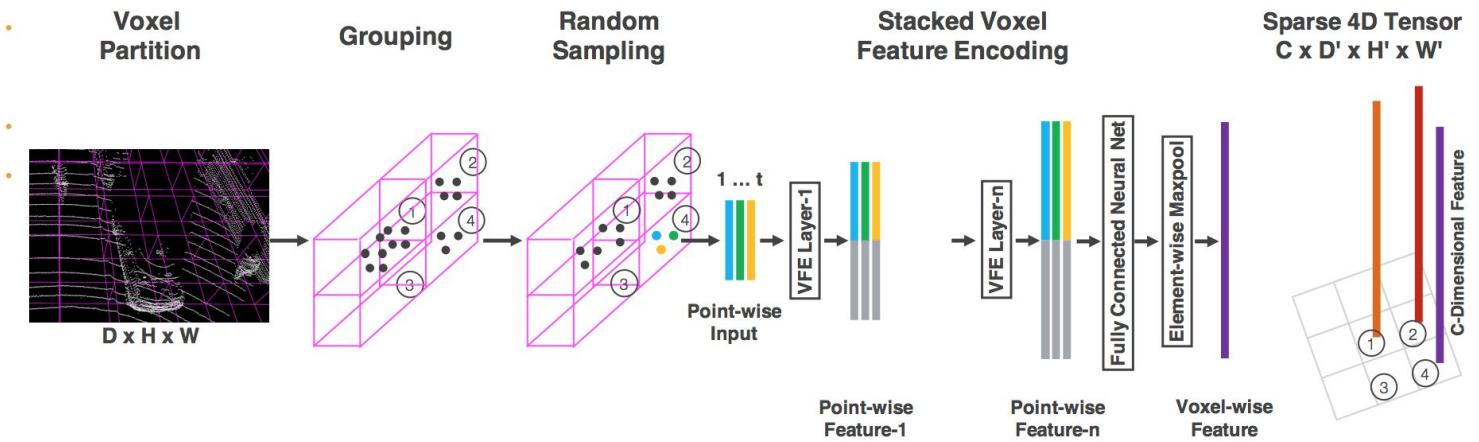
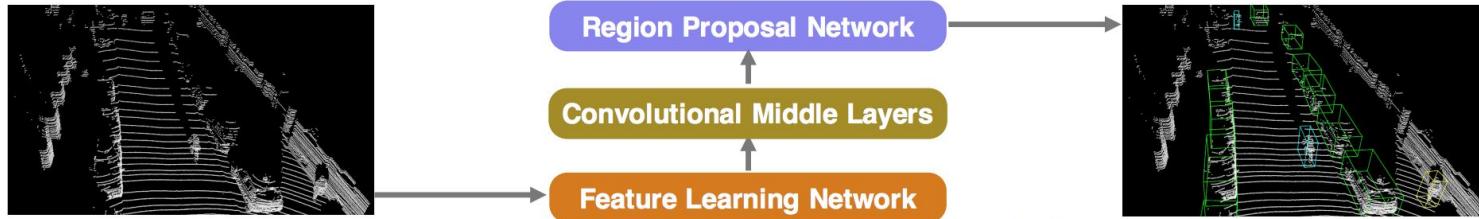
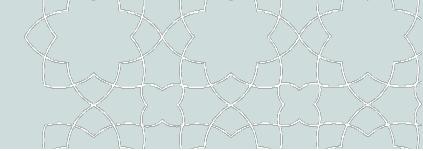
**Stage1:** 3D Region Proposal Network

TSDF

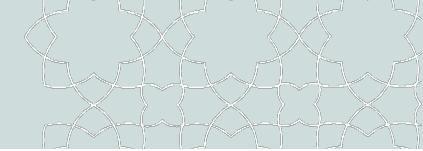
**Stage2:** Joint Object Recognition Network



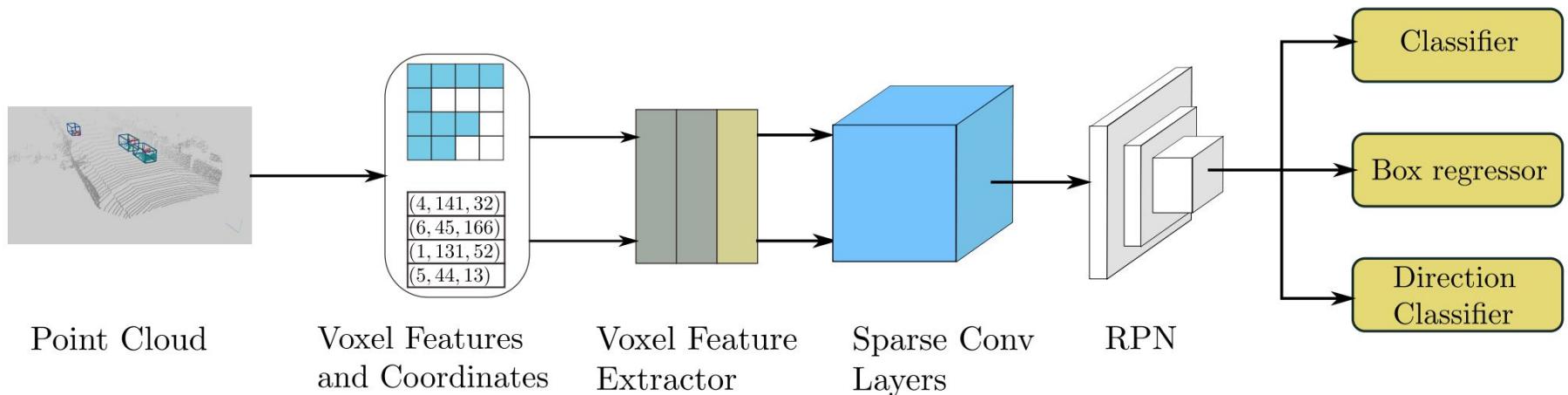
# VoxelNet: 3D Convolution



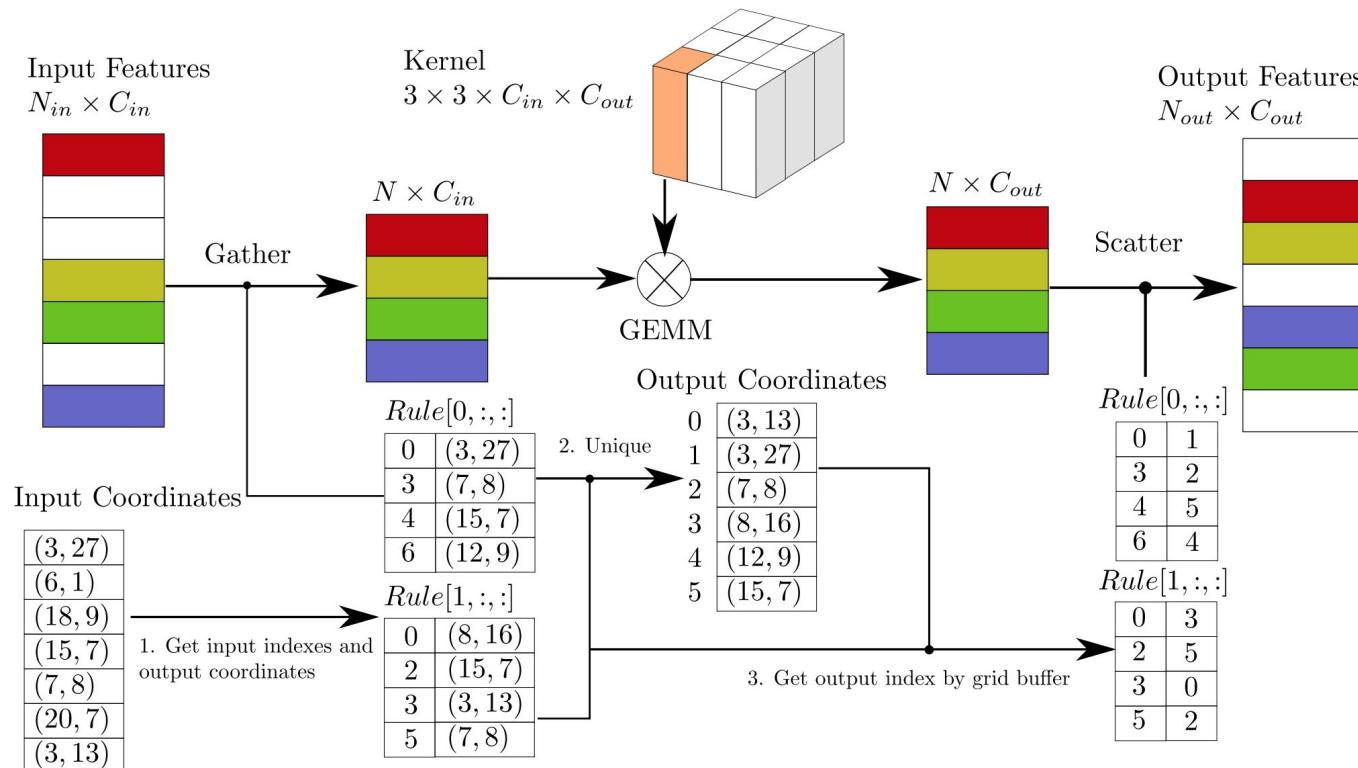
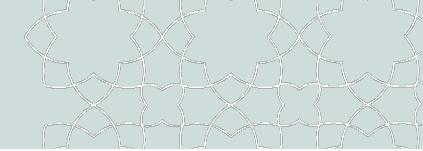
# SECOND: Sparse 3D Convolution



- Sparse 3D Convolution:
  - Reduce complexity
  - Avoid computation on empty voxels

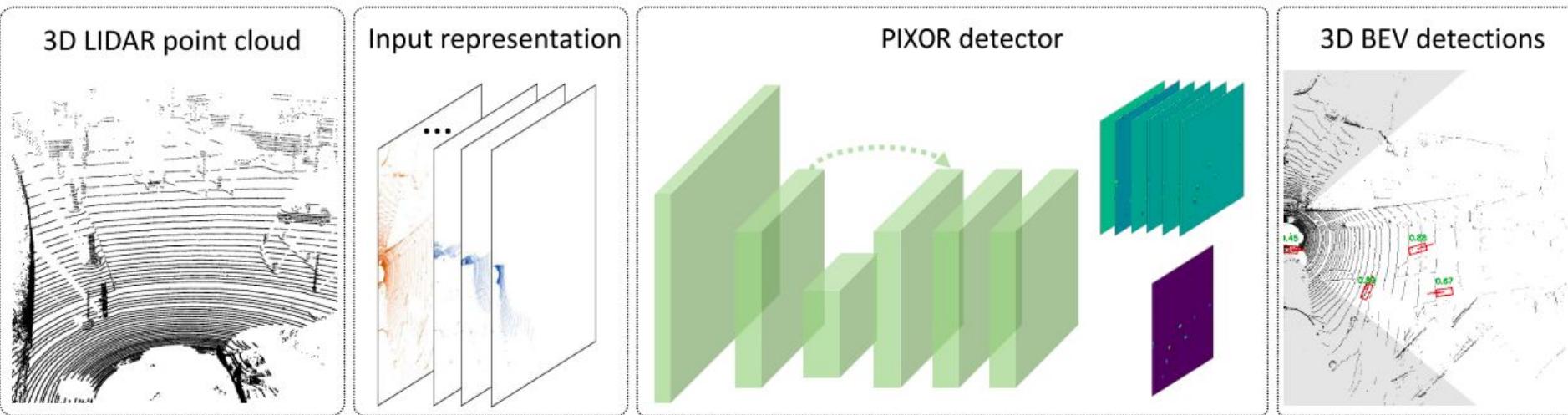


# SECOND: 3D Sparse Convolution



# BEV (bird-eye view)

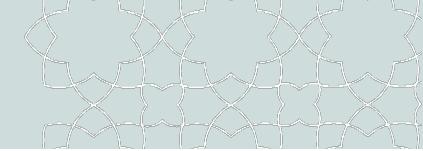
- Top View Representation (2D Projection)



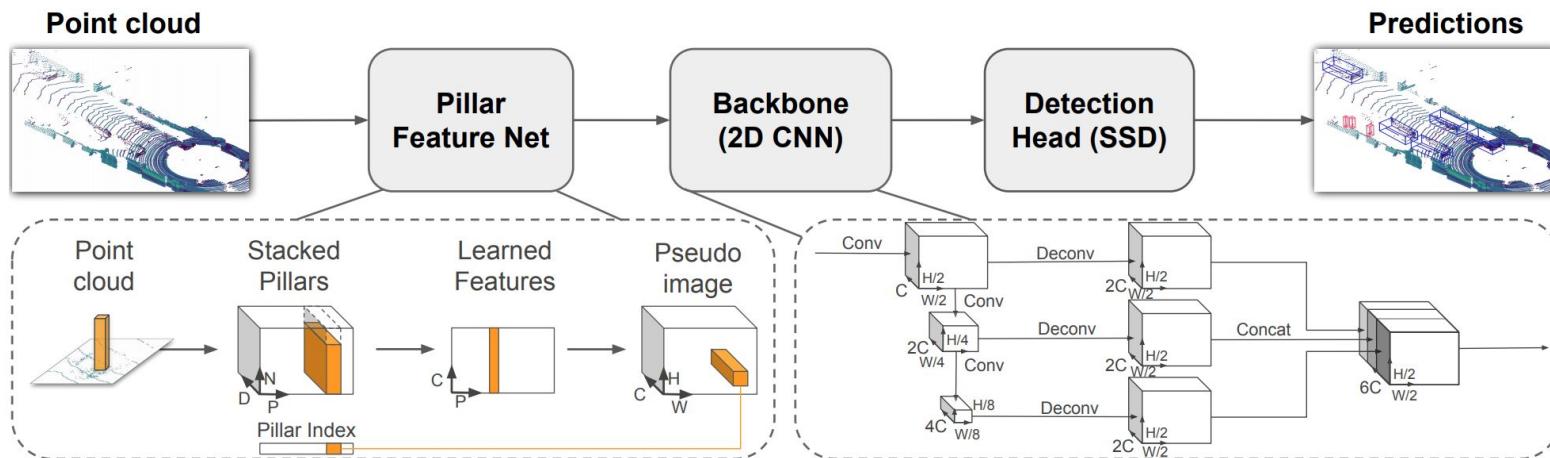
$H \times W \times D$  voxels

$H \times W$  image with  
D channels

# PointPillars: Bird Eye View Convolution

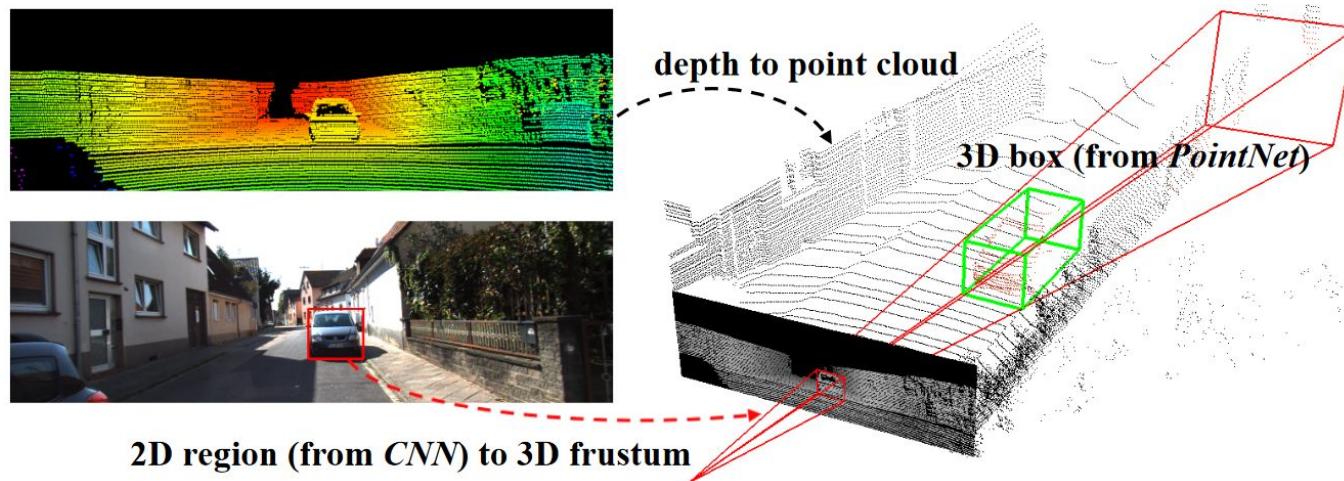


- 2D Convolution from top view:
  - Replace occupancy of each pixel with point cloud features of the pillar



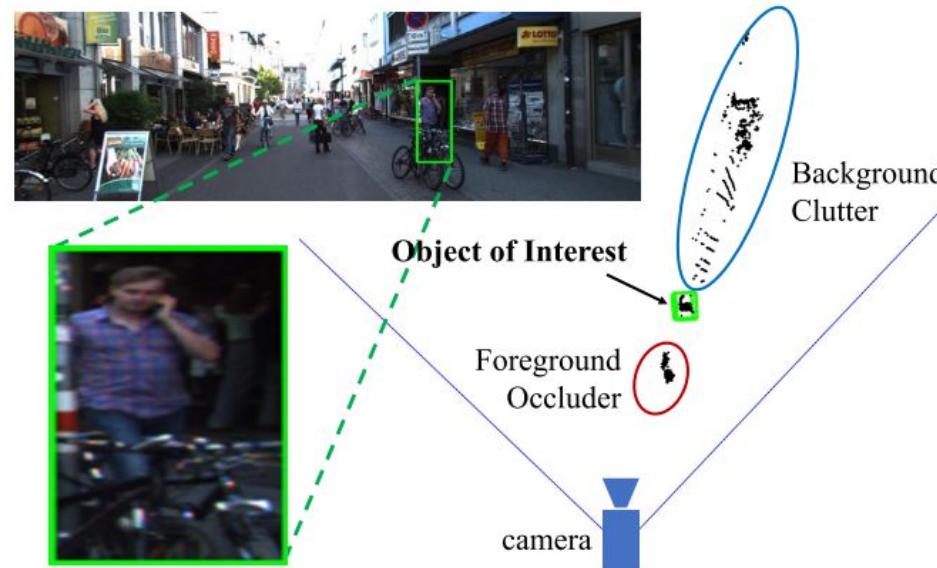
# Frustum PointNets: View-based Proposal

- Generate object proposals from a view  
(e.g., using SSD)



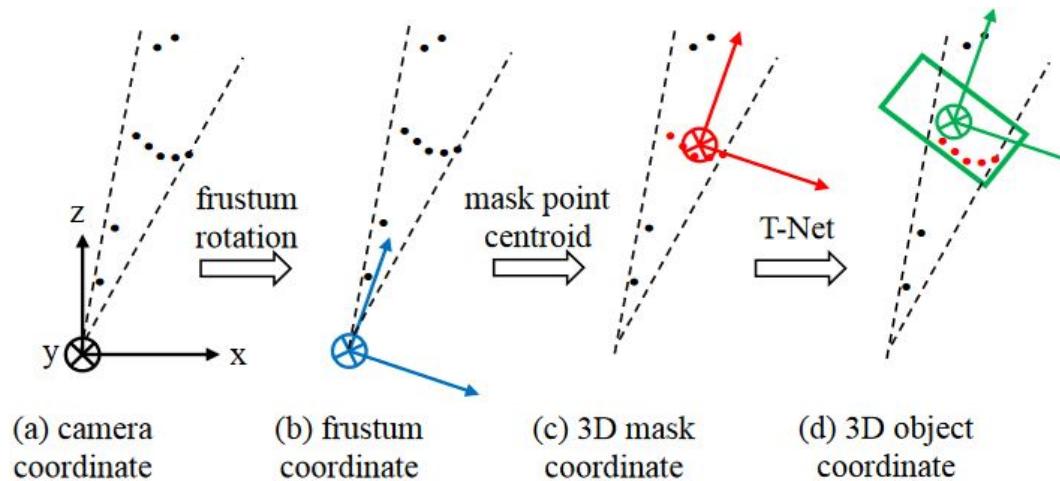
# Frustum PointNets: FG/BG Segmentation

- Stage 2: FG/BG Segmentation
  - Instance segmentation to remove other foreground instances and background clutter

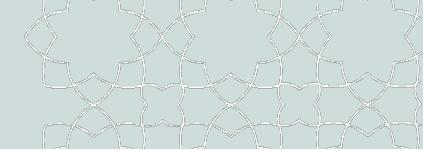


# Frustum PointNets: Coordinate Normalization

- Instance segmentation to remove other foreground instances and background clutter

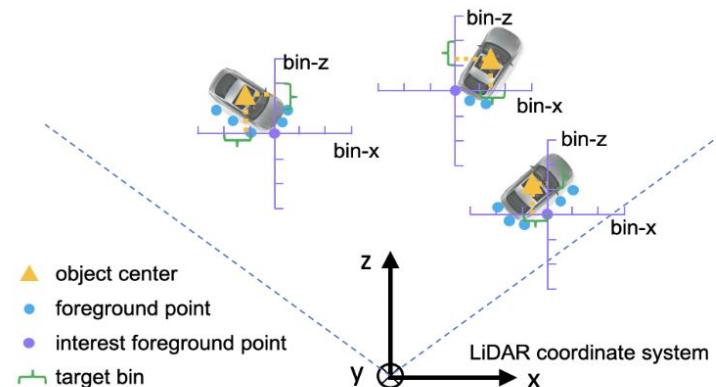


# PointRCNN: Region Proposal Network

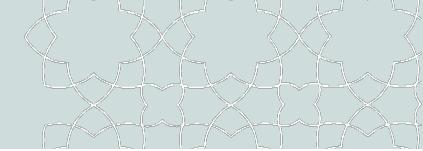


**Stage-1:** Foreground/Background segmentation and generate 3D proposals for each foreground point

**Stage-2:** Refine proposals in the canonical coordinates

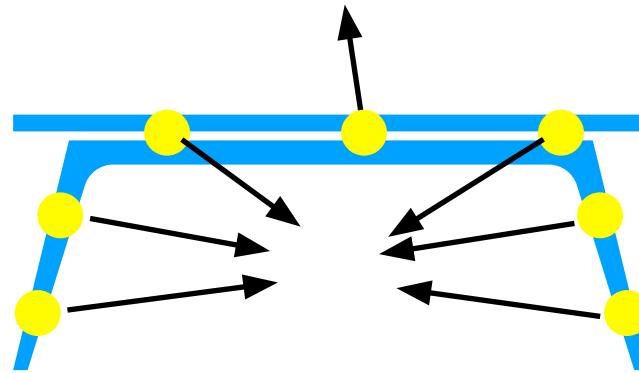
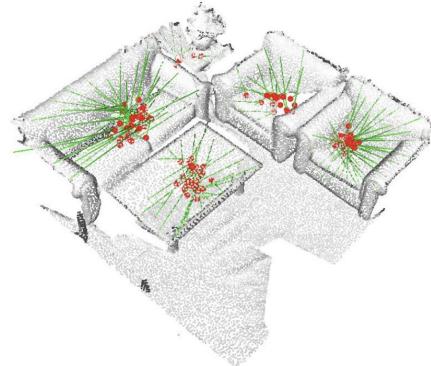


# VoteNet: Proposal from Voting

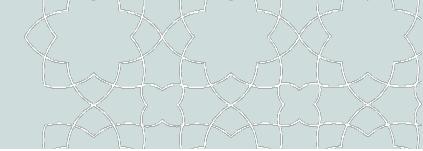


- **Challenge:** 3D object centroid can be far from any surface point, thus hard to regress accurately
  - Sample a set of seed points and generate votes, targeting at object centers

Voting from input point cloud

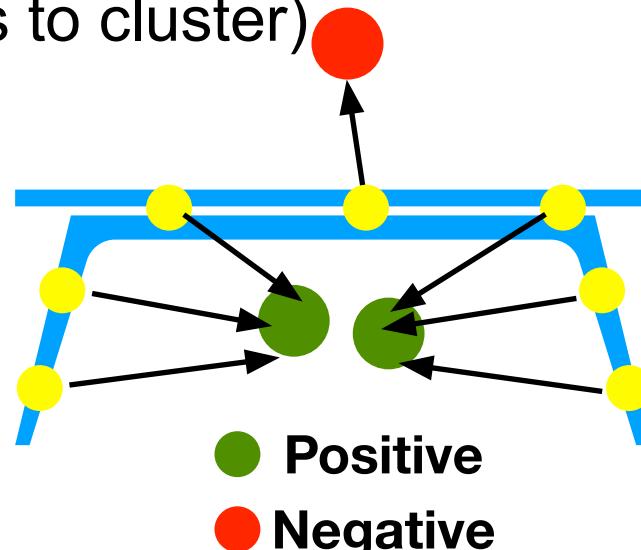
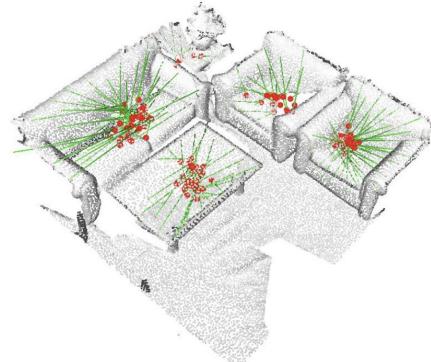


# VoteNet: Proposal from Voting



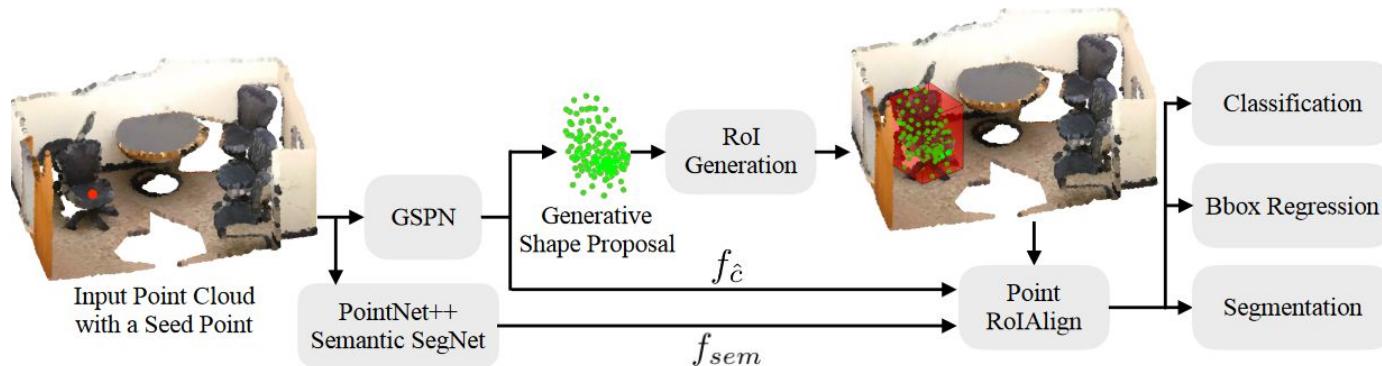
- **Challenge:** 3D object centroid can be far from any surface point, thus hard to regress accurately
  - Votes cluster near object centers  
(Farthest point sample votes to cluster)

Voting from input point cloud

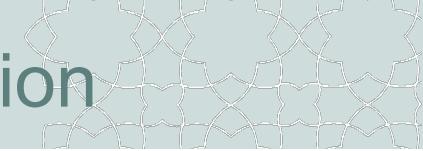


# GSPN : Proposal from Generative Network

- Randomly sample seeds points
- Take point cloud and a seed point as input, use conditional VAE to generate a point cloud as proposal
- Convert the proposal to an ROI box
- R-PointNet (mask RCNN) to segment the object

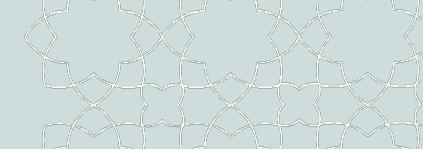


# Task: 3D Detection & Instance Segmentation



- Bottom-Up Methods
  - Do these points belong to the same instance?
  - How to measure the similarity between two points?

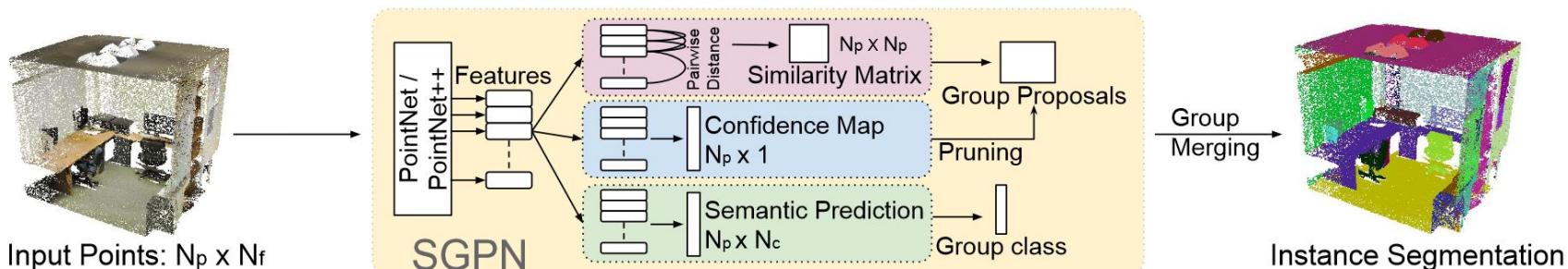
# Associative Embedding



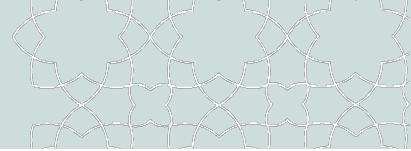
- Learn a per-point embedding, so that points from the same instance have similar embeddings

$$l(i, j) = \begin{cases} \|F_{SIM_i} - F_{SIM_j}\|_2 & C_{ij} = 1 \\ \alpha \max(0, K_1 - \|F_{SIM_i} - F_{SIM_j}\|_2) & C_{ij} = 2 \\ \max(0, K_2 - \|F_{SIM_i} - F_{SIM_j}\|_2) & C_{ij} = 3 \end{cases}$$

- Clustering gives proposals

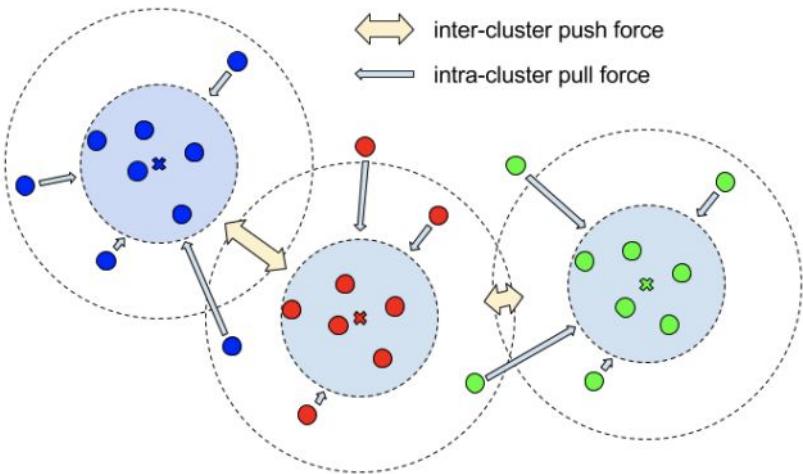


# JSIS3D: “push/pull” forces



- A discriminative function to present the embedding loss

$$\mathcal{L}_{embedding} = \alpha \cdot \mathcal{L}_{pull} + \beta \cdot \mathcal{L}_{push} + \gamma \cdot \mathcal{L}_{reg}$$



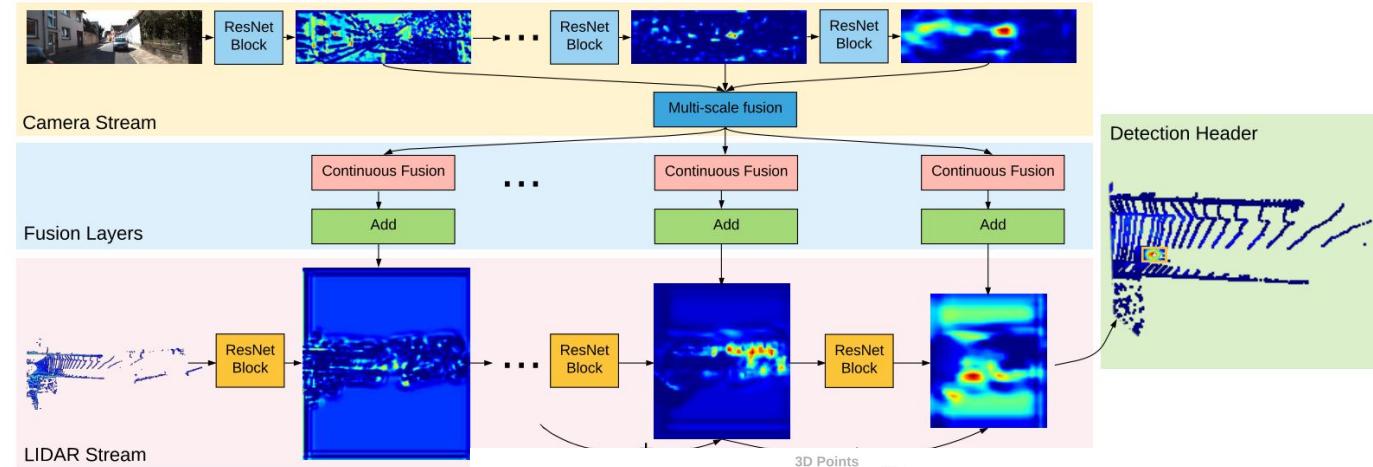
$$\mathcal{L}_{pull} = \frac{1}{K} \sum_{k=1}^K \frac{1}{N_k} \sum_{j=1}^{N_k} [\|\boldsymbol{\mu}_k - \mathbf{e}_j\|_2 - \delta_v]_+^2$$

$$\mathcal{L}_{push} = \frac{1}{K(K-1)} \sum_{k=1}^K \sum_{m=1, m \neq k}^K [2\delta_d - \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_m\|_2]_+^2$$

$$\mathcal{L}_{reg} = \frac{1}{K} \sum_{k=1}^K \|\boldsymbol{\mu}_k\|_2$$

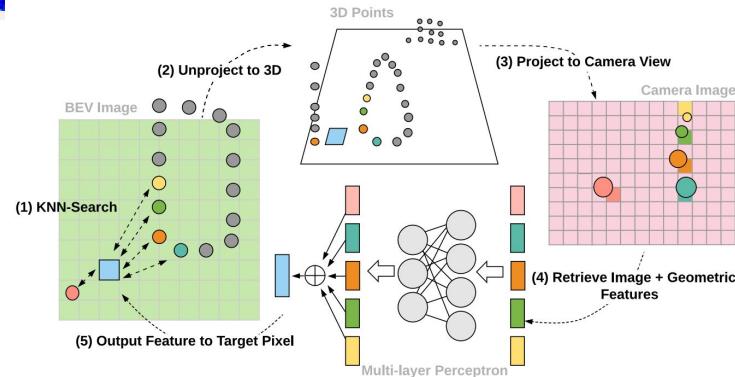
# Multi-Modal Continuous Fusion

Image feature network: 2D image

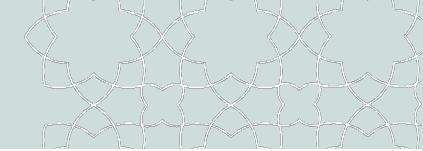


BEV network: 3D voxel (HxWxC)

3D points as the intermediate media to transform from the camera to BEV



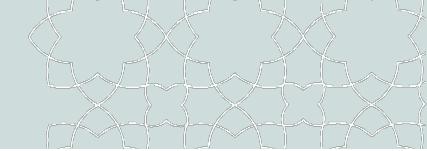
# Task: Few-shot/Zero-shot Learning



- Can be a better platform than images
- Shapes are **pure**, with no contamination by projection distortion, illumination, ...

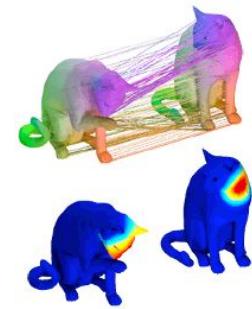
If one image is more than a thousand words,  
then one shape is more than one thousand  
pictures

# Task: Few-shot/Zero-shot Learning

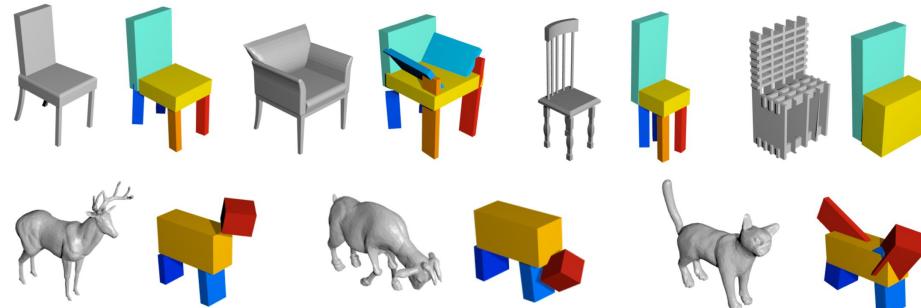


Algorithmically, 3D shapes are:

- easier to be **related (correspondence)**

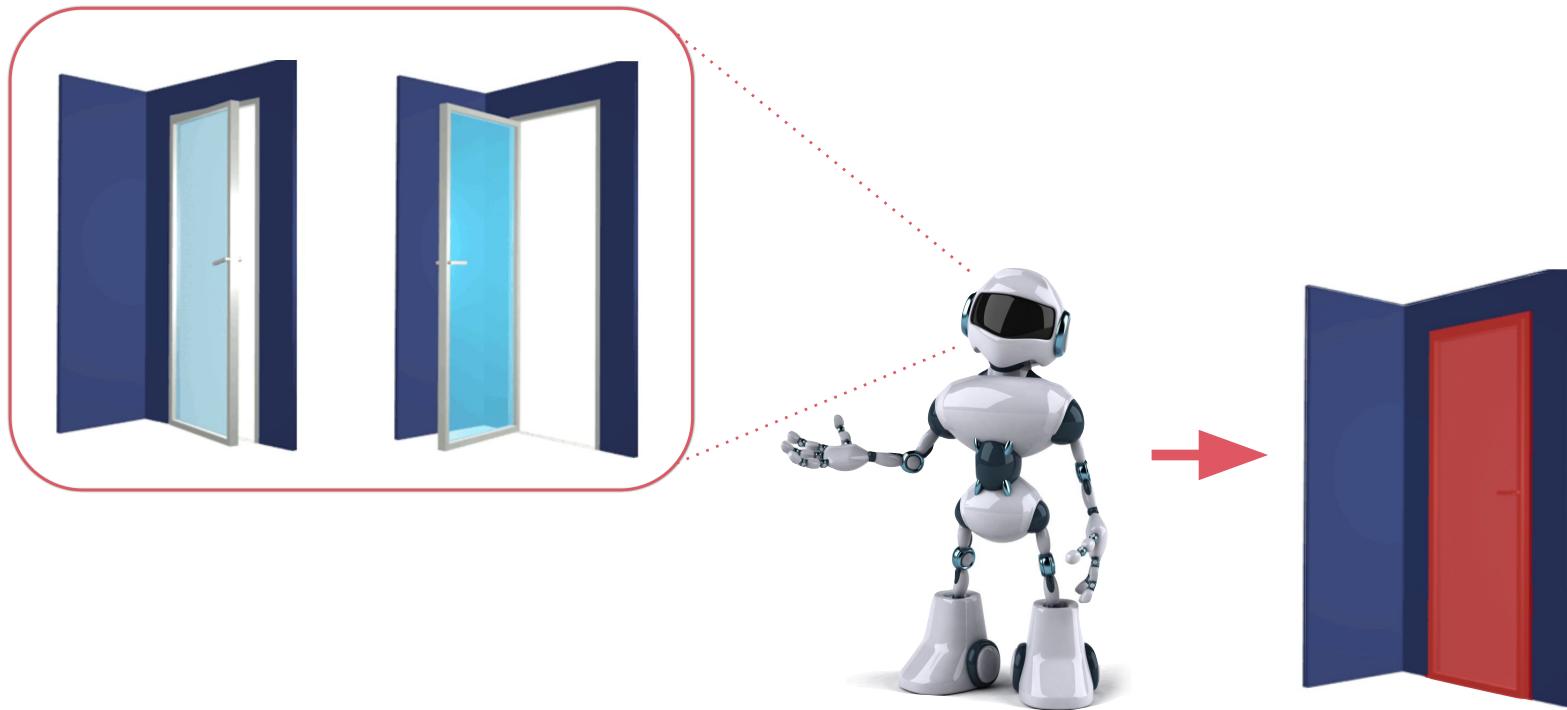
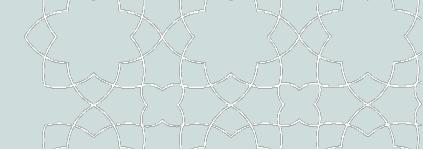


- easier to be **compared**

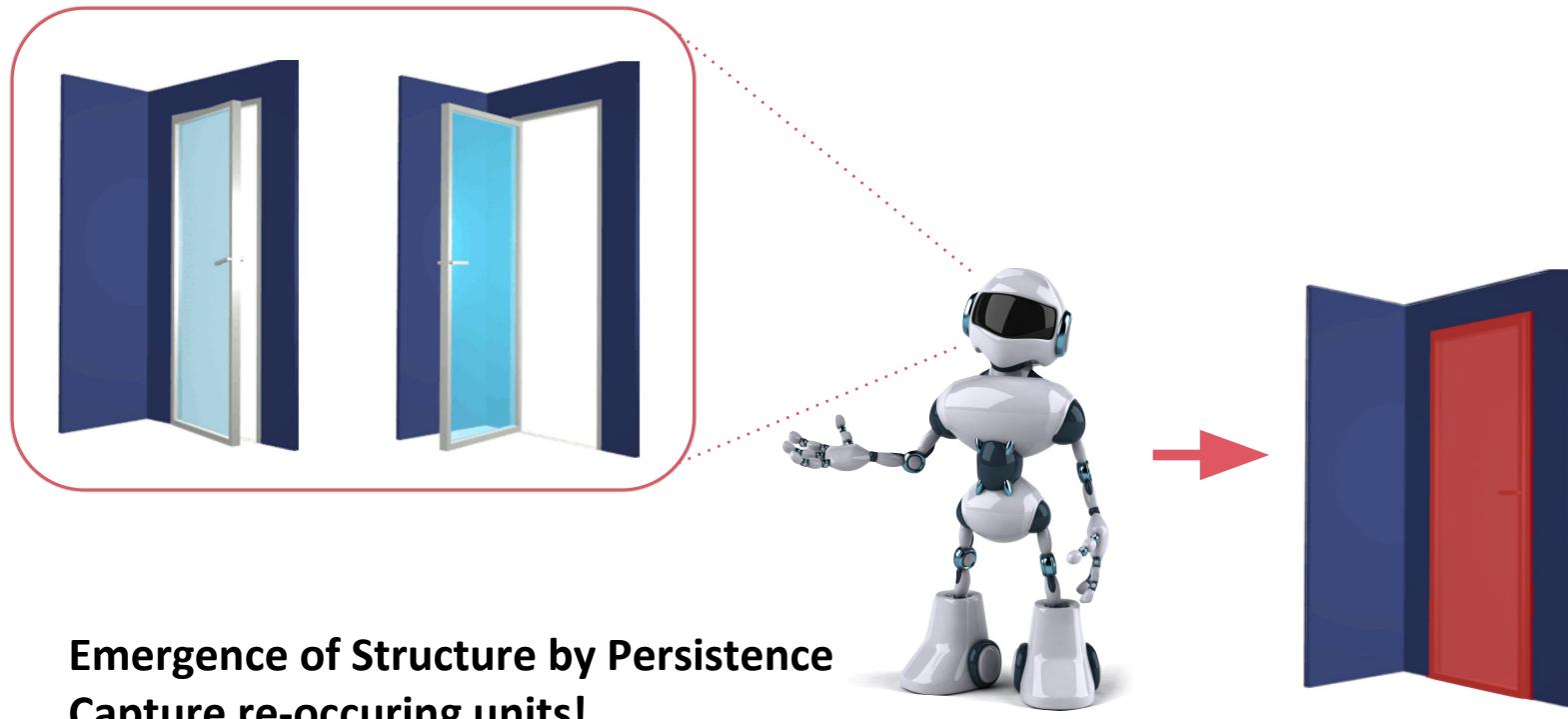
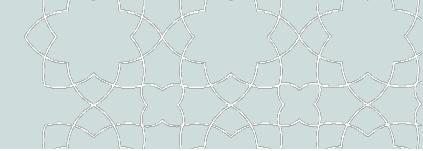


- easier to **abstracted**

# Task: Few-shot Structure Induction

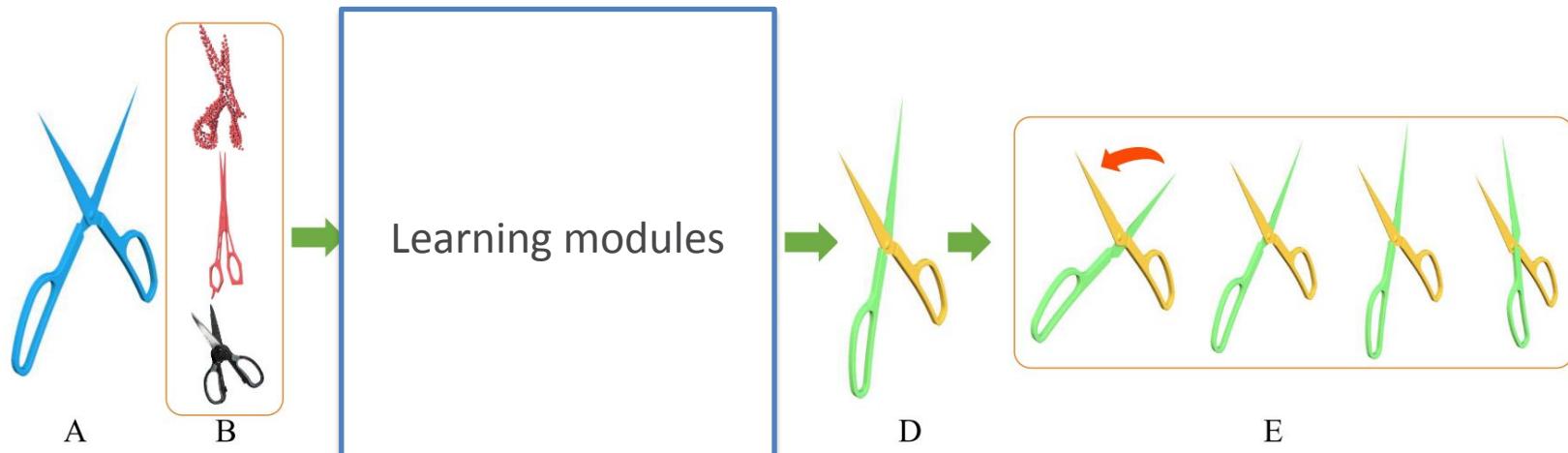
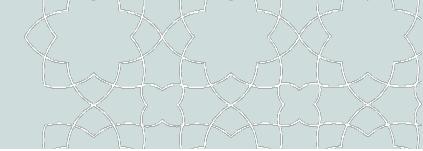


# Task: Few-shot Structure Induction

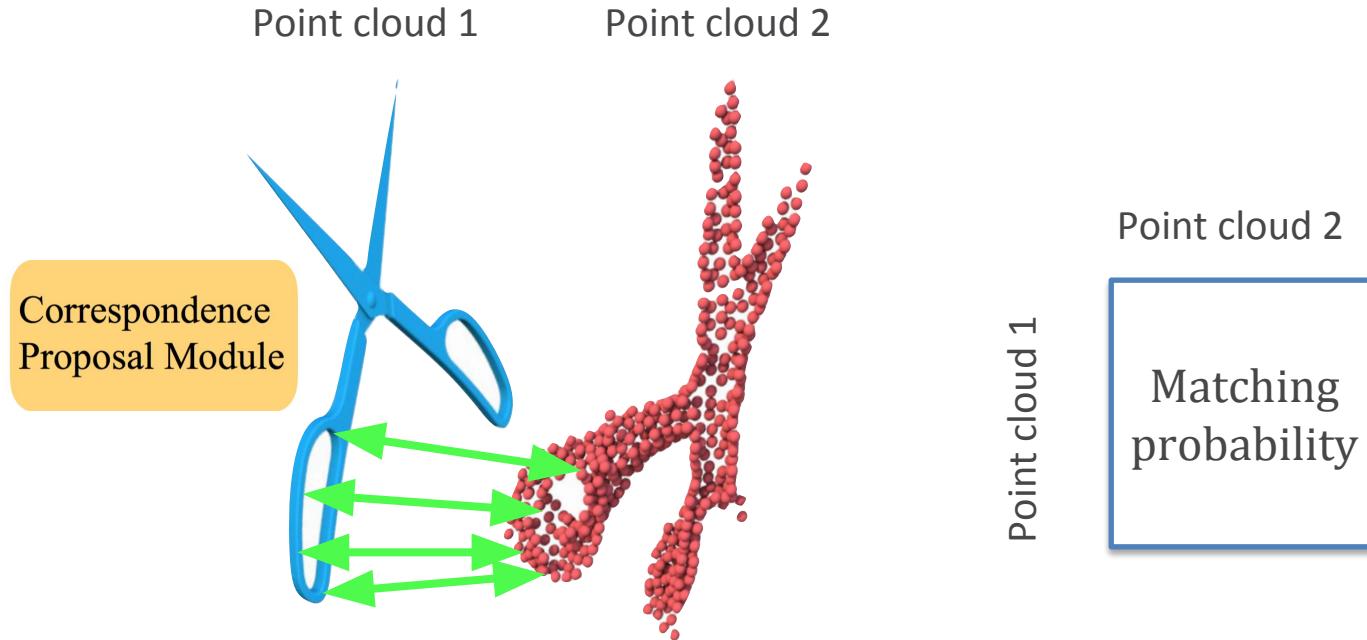
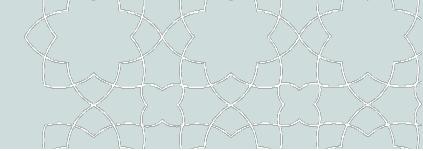


**Emergence of Structure by Persistence  
Capture re-occurring units!**

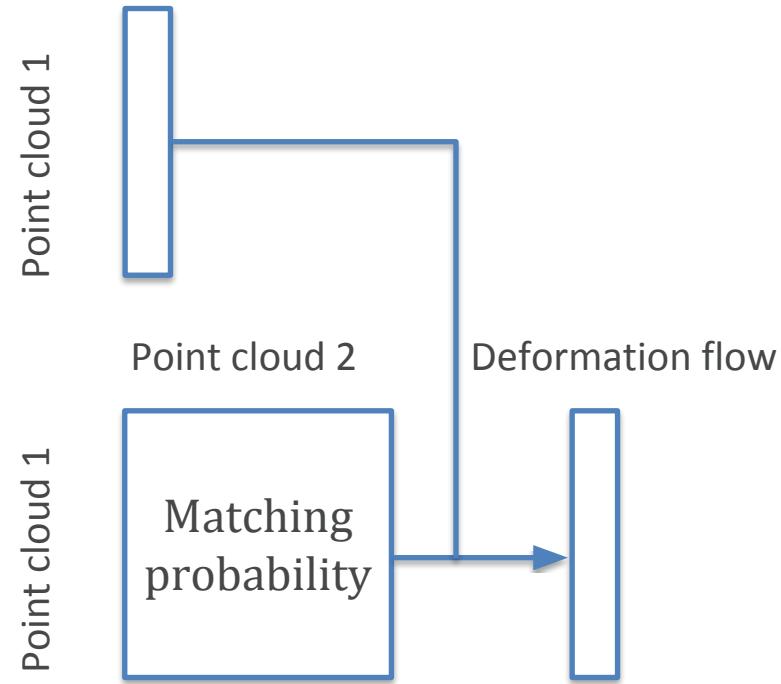
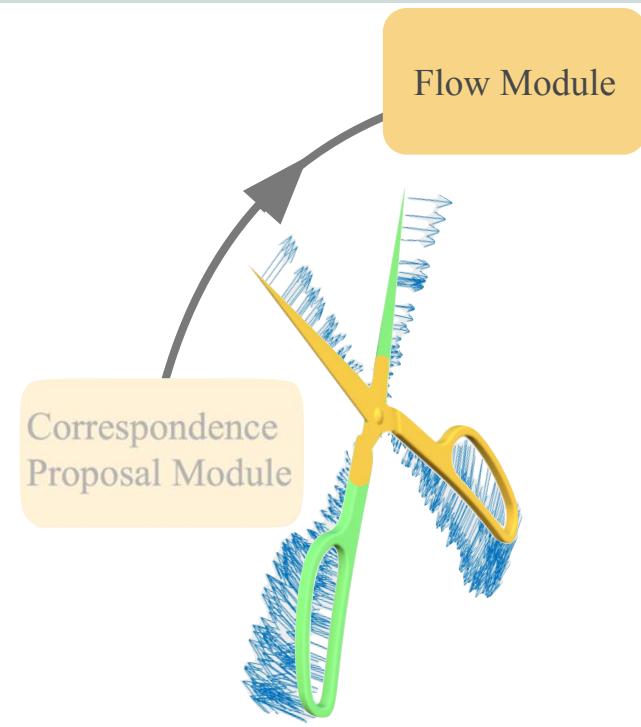
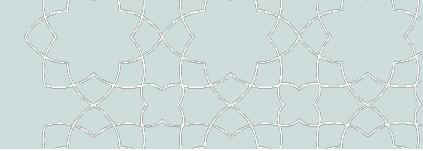
# Part Induction by Relating Shapes



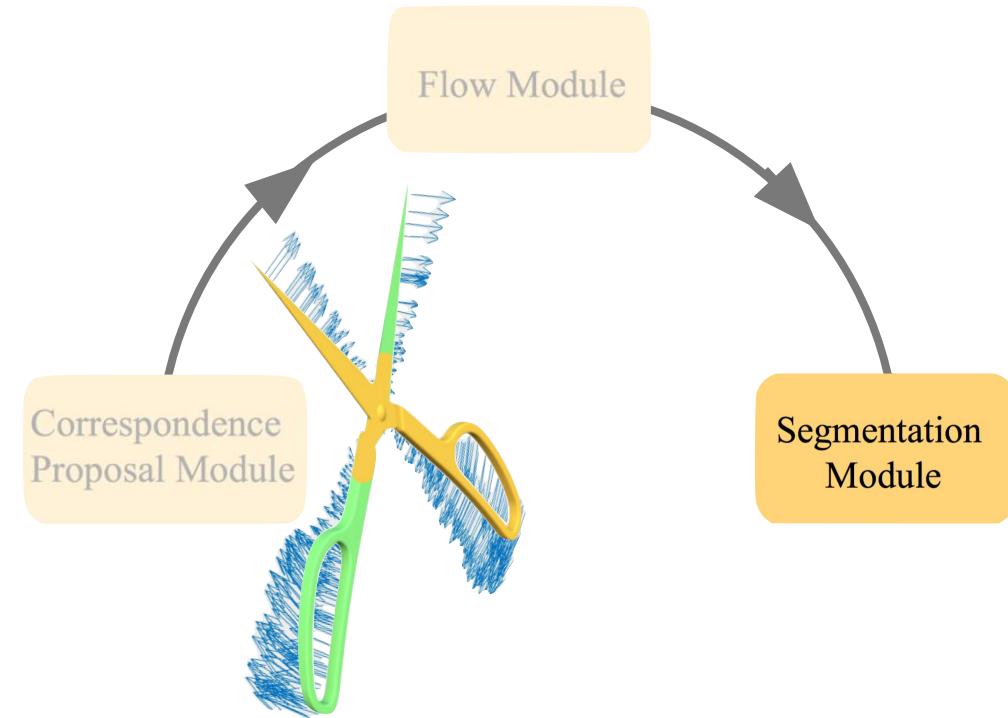
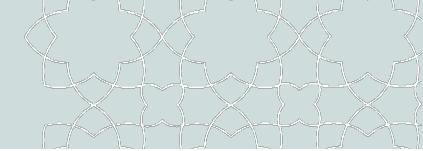
# Part Induction by Relating Shapes



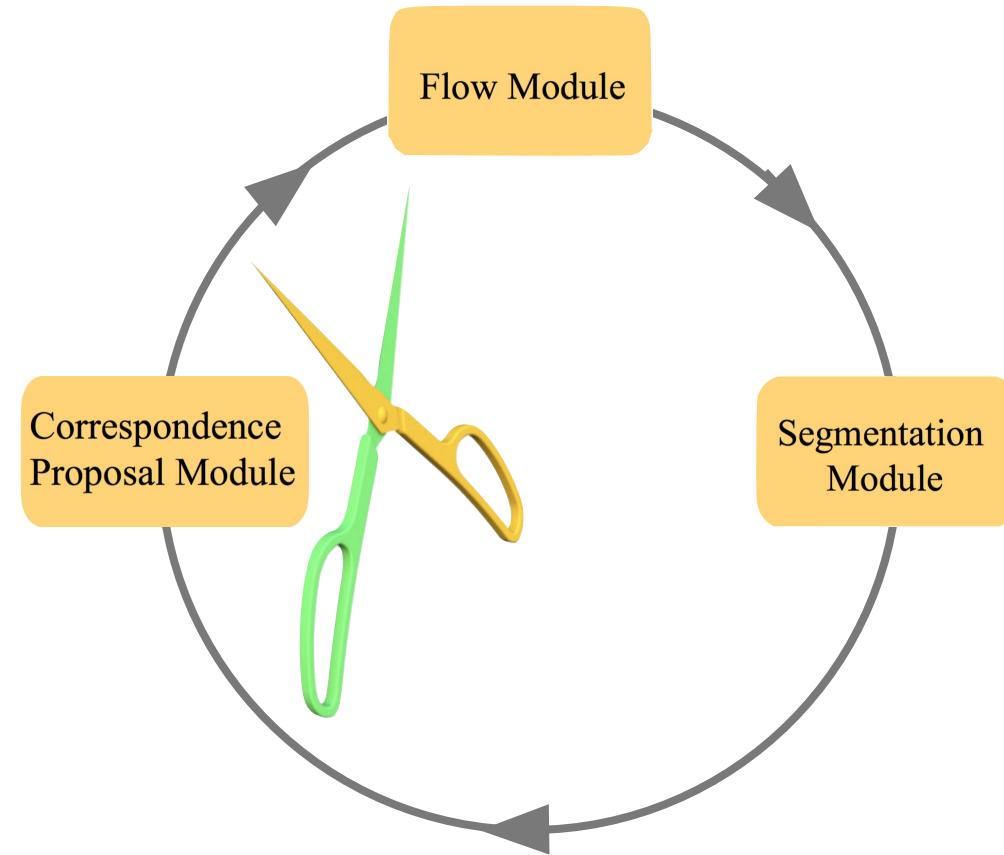
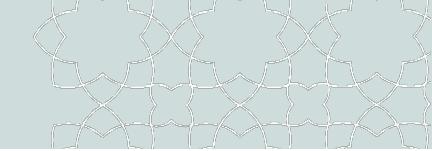
# Part Induction by Relating Shapes



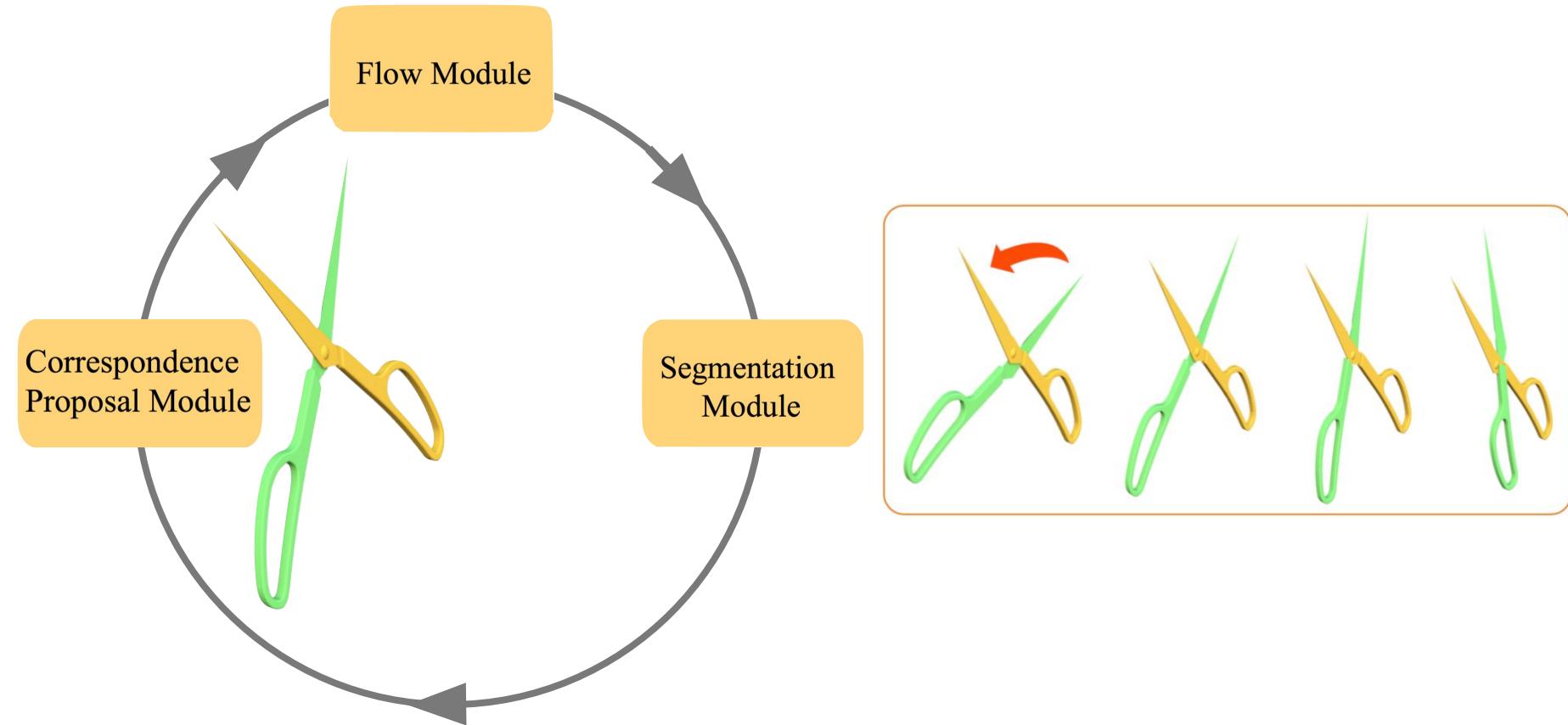
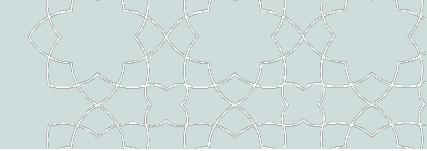
# Part Induction by Relating Shapes



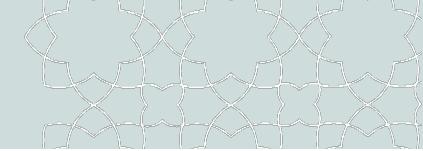
# Part Induction by Relating Shapes



# Part Induction by Relating Shapes



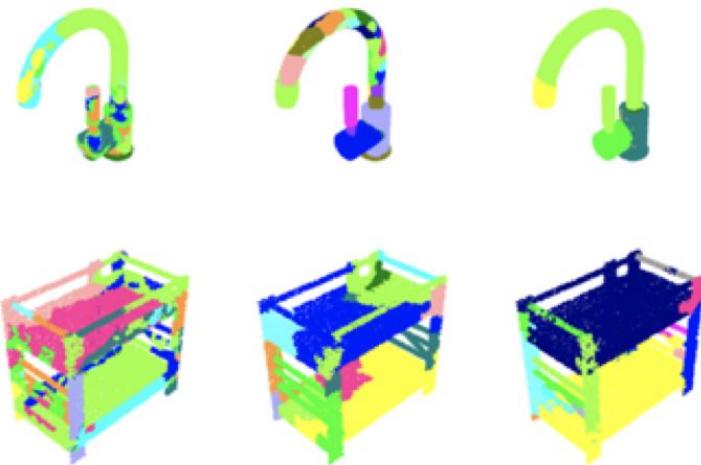
# Task: Zero-shot Part Discovery



Train set



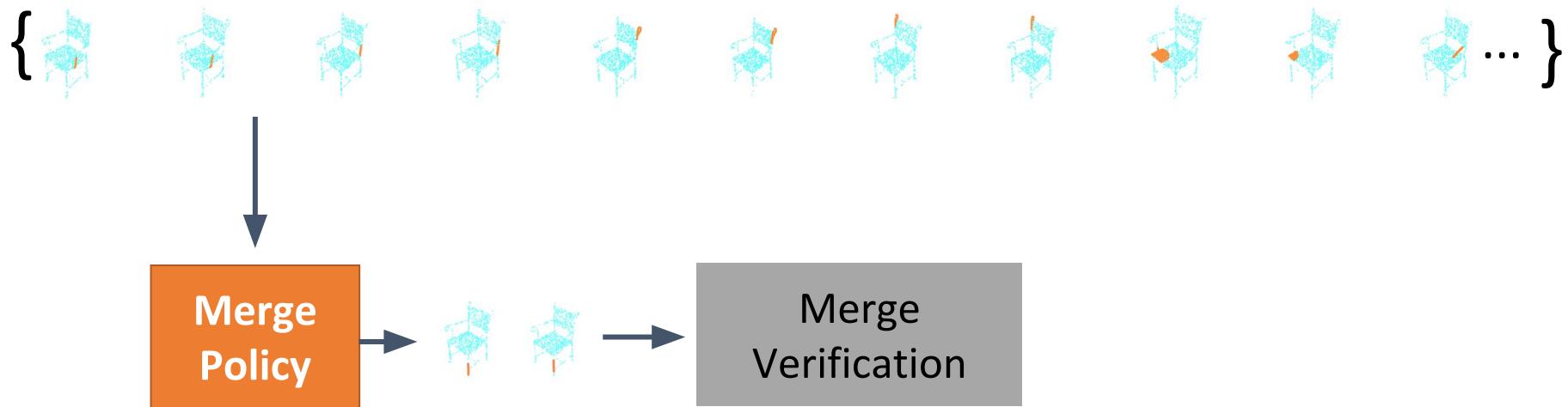
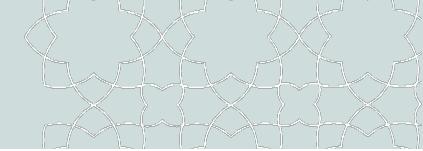
Test set



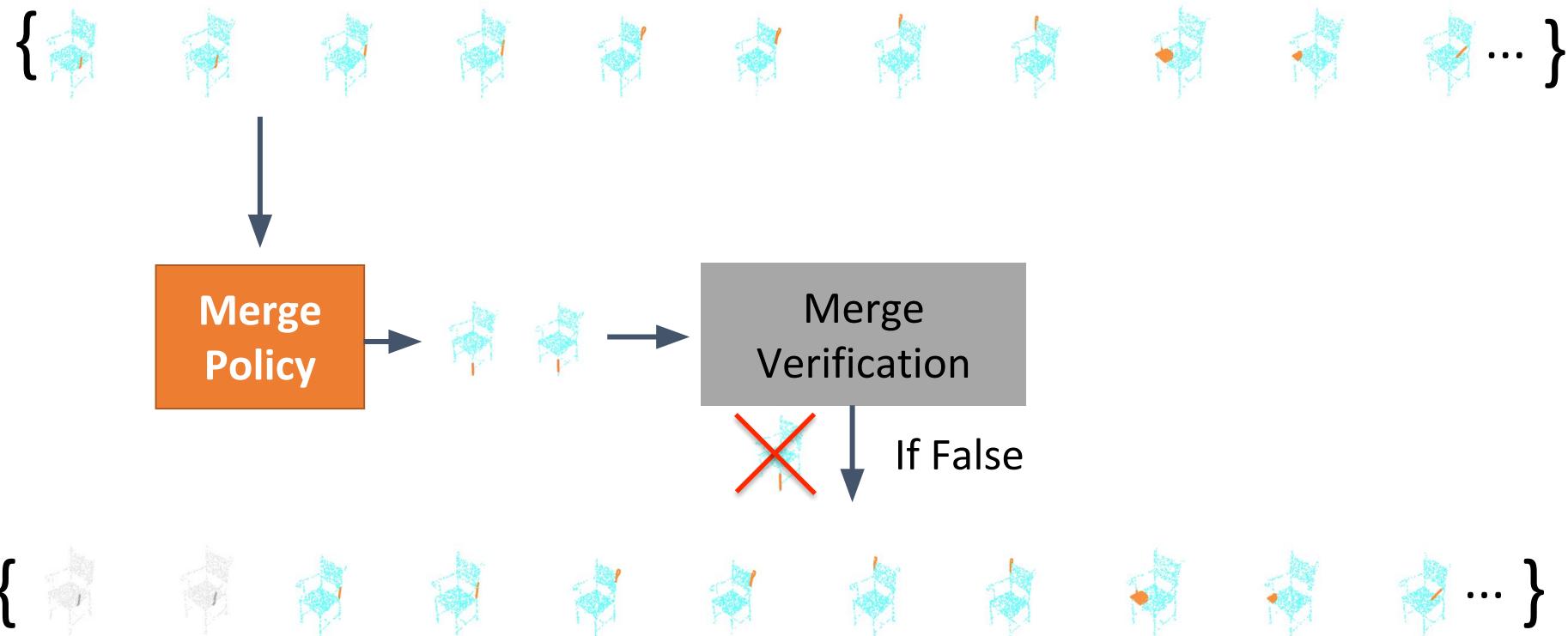
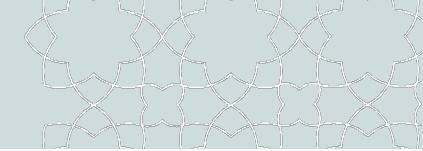
SOTA of  
Deep Learning  
(PartNet)      SOTA of  
Classics  
(WCSeg)

**Ours**

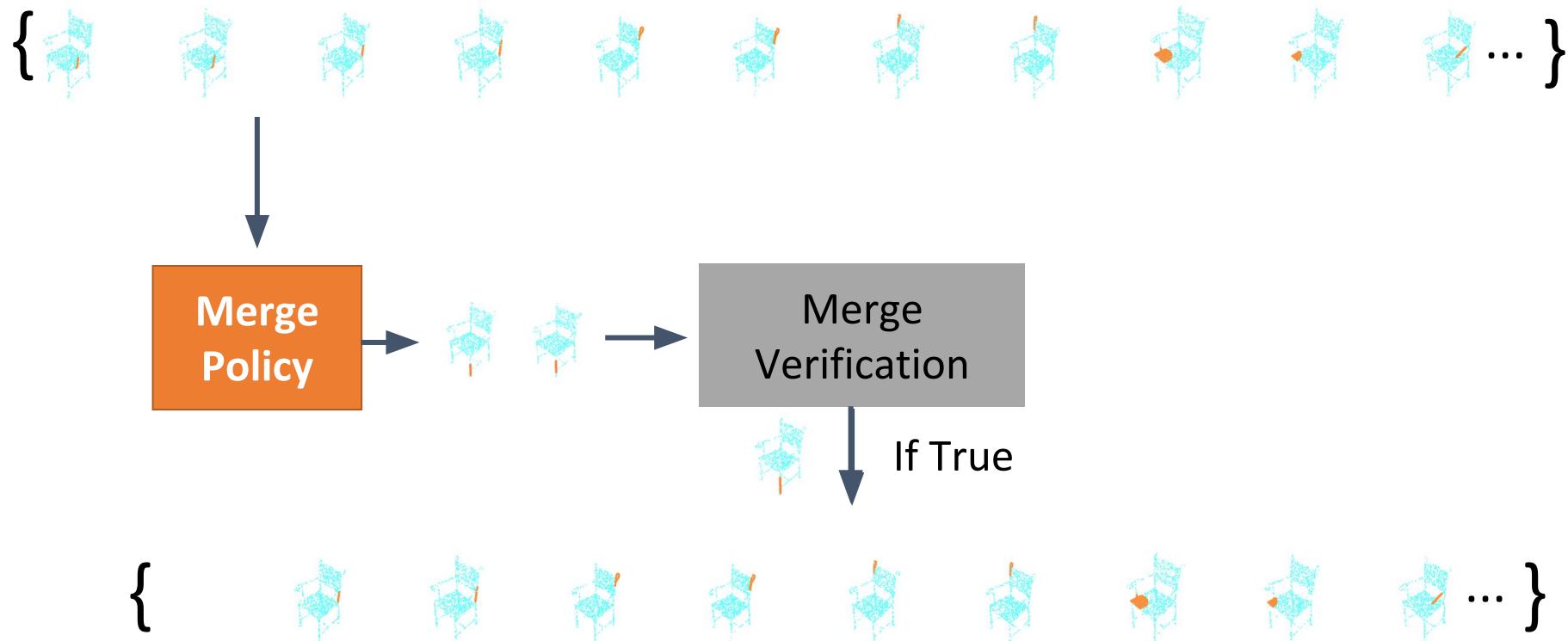
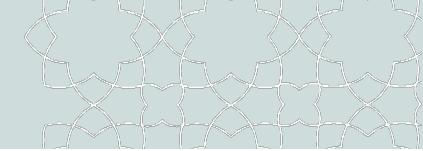
# Learning to Group: Sub-Part Pool



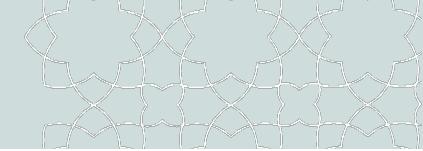
# Learning to Group: Sub-Part Pool



# Learning to Group: Sub-Part Pool



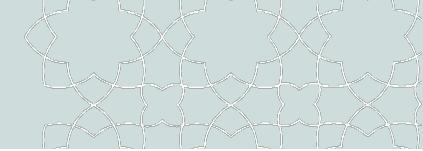
# Learning to Group: Sub-Part Pool



# Outline

- 3D Data
- Classification
- Segmentation and Detection
- **3D Data synthesis**

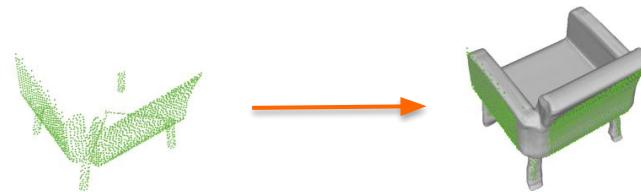
# Task: 3D Data Synthesis



- Conditional generation



Single-image  
3D reconstruction

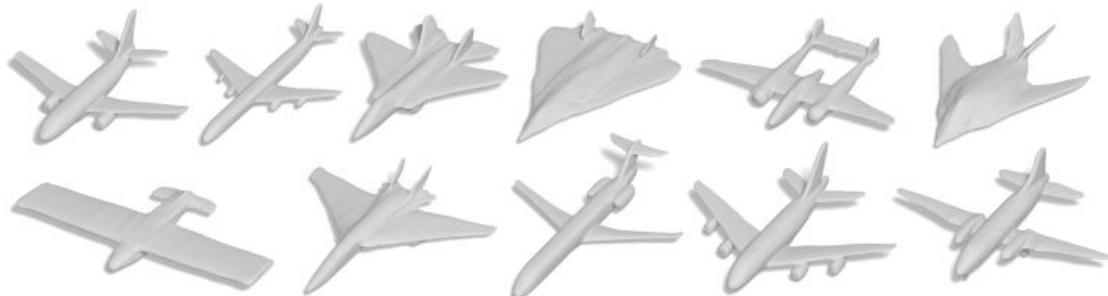


Shape Completion

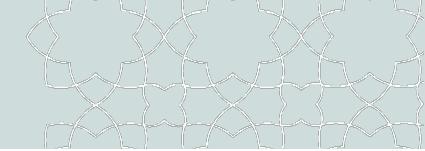
- Free generation



Gaussian Noise

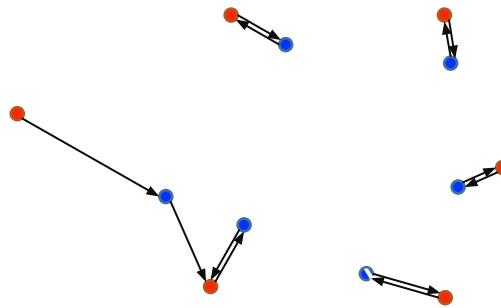


# Task: 3D Data Synthesis



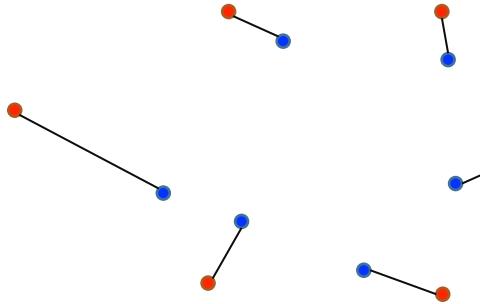
- Metrics to Compare Point Clouds

## Chamfer Distance



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

## Earth Mover's Distance



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

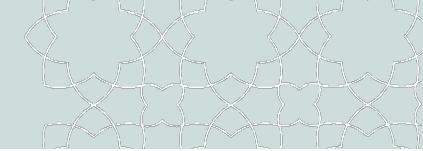
where  $\phi : S_1 \rightarrow S_2$  is a bijection.

Fan et al., "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

Wang et al., "Pixel2mesh: Generating 3d mesh models from single rgb images", ECCV 2018

Mescheder et al., "Occupancy Networks: Learning 3D Reconstruction in Function Space", CVPR 2019

# Task: 3D Data Synthesis



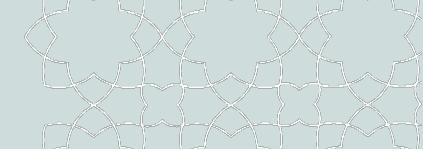
- Metrics to Compare Point Clouds
  - **F-score**: harmonic mean of precision/recall.  
Percentage of points in one point cloud that can find a neighbor from the other point cloud within a threshold.
  - **Normal Consistency**: dot product of the normals of each point and its nearest neighbor

Fan et al., “A Point Set Generation Network for 3D Object Reconstruction from a Single Image”, CVPR 2017

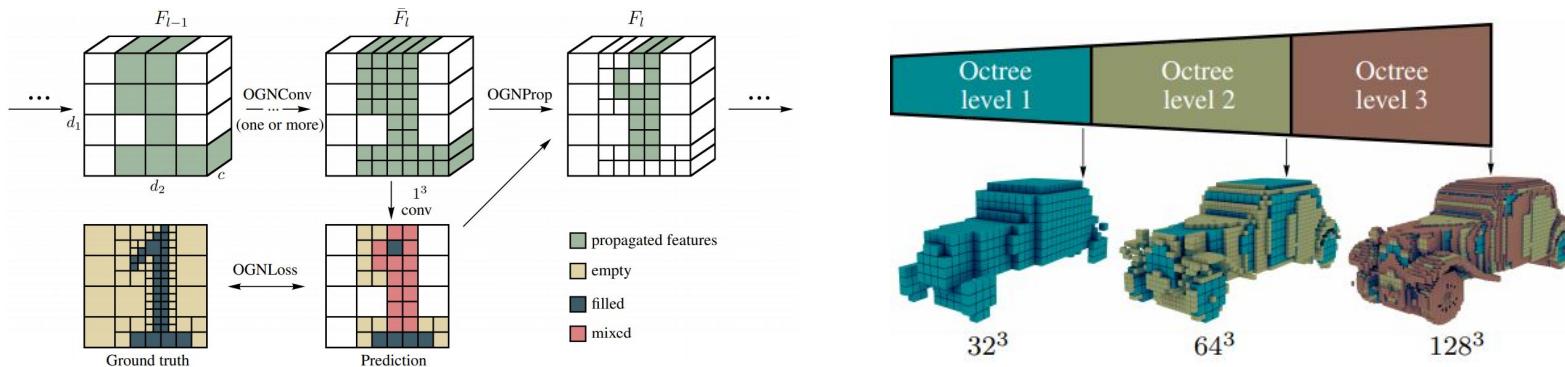
Wang et al., “Pixel2mesh: Generating 3d mesh models from single rgb images”, ECCV 2018

Mescheder et al., “Occupancy Networks: Learning 3D Reconstruction in Function Space”, CVPR 2019

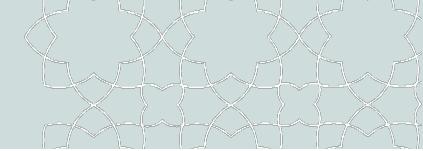
# Task: Conditional Generation



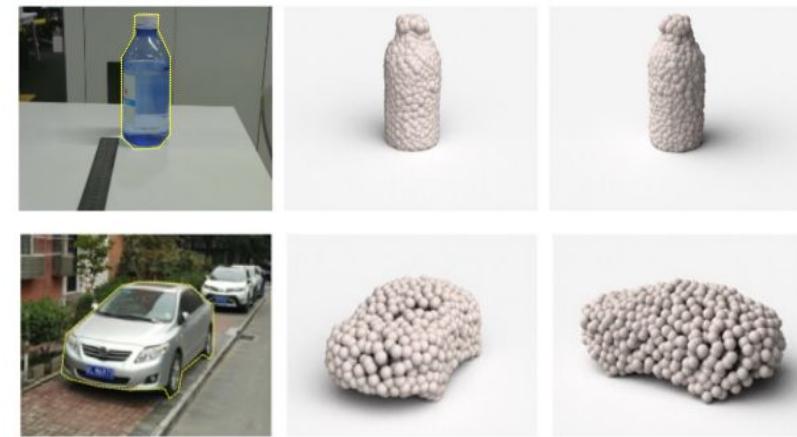
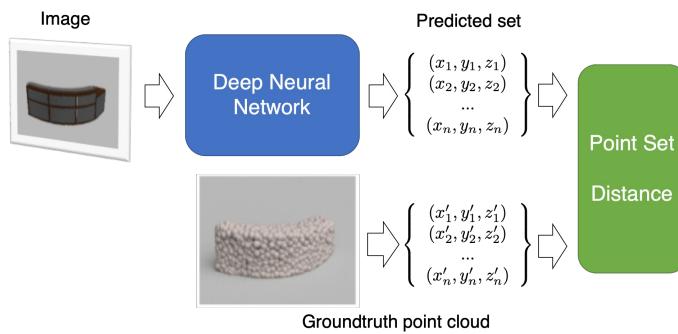
- **From Single Image to Volume:**
  - Octree representation of shapes
  - Avoid  $\mathcal{O}(n^3)$  reconstruction
  - Generate the octree layer by layer



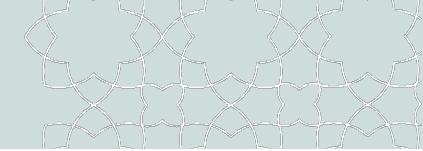
# Task: Conditional Generation



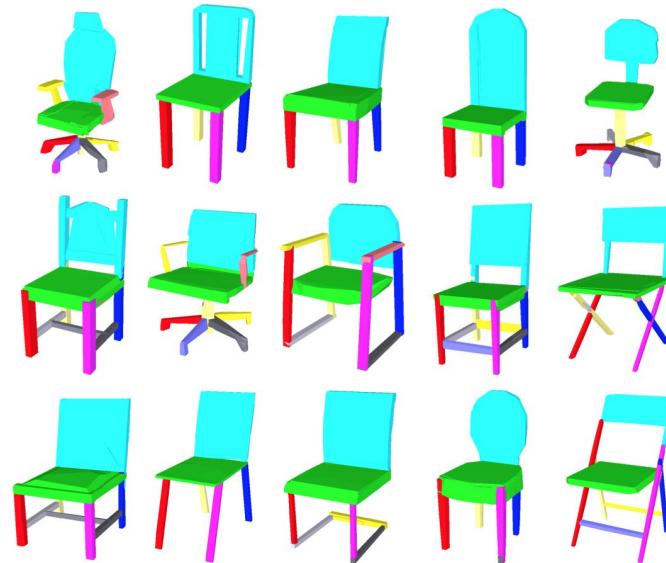
- **From Single Image to Point Cloud:**
  - Point set generation



# Task: Conditional Generation

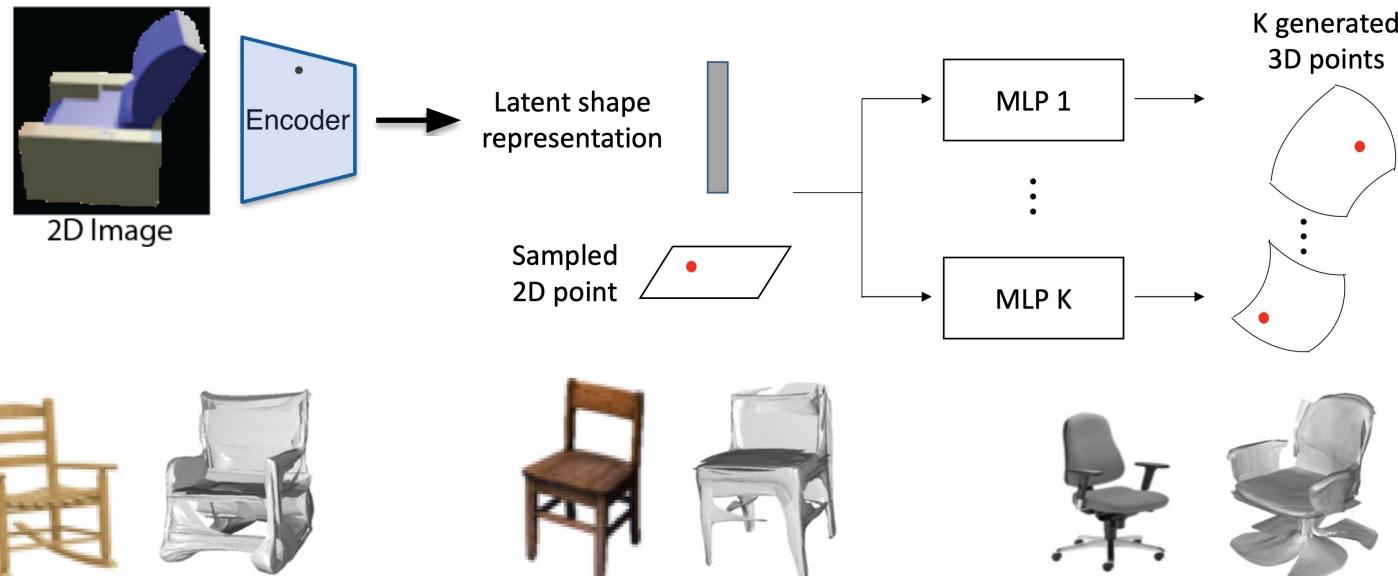


- **From Image to Shape:**
  - Planes  $\rightarrow$  convex polytopes  $\rightarrow$  shape



# Task: Conditional Generation

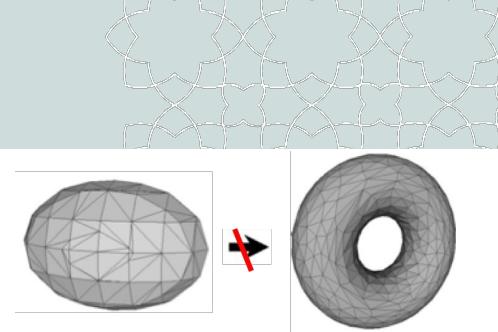
- From Image to Surface:
  - Learn to warp a plane to surface



Groueix et al., "AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation", CVPR 2018

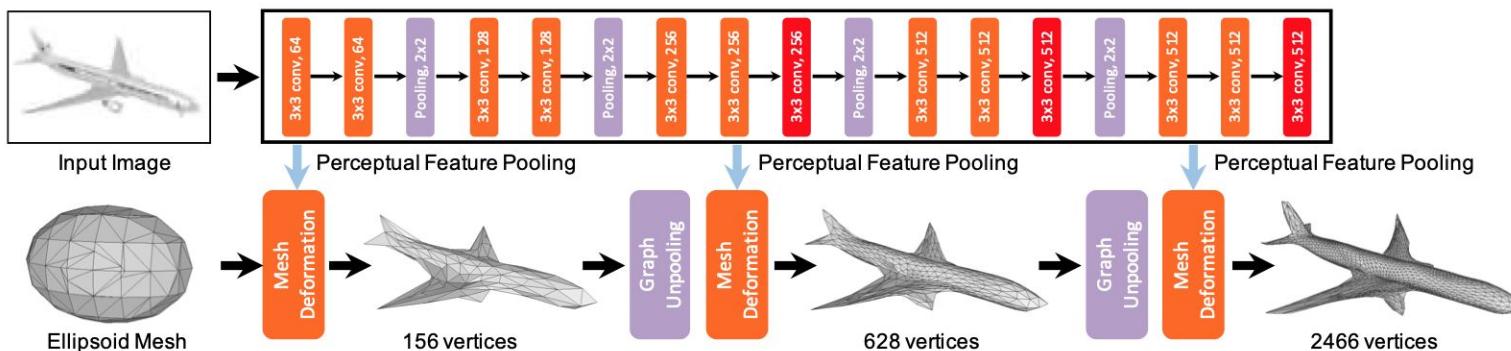
Yang et al., "Foldingnet: Point cloud auto-encoder via deep grid deformation", CVPR 2018

# Task: Conditional Generation



- **From Image to Surface:**
  - Polygon meshes are irregular
  - Learn to deform a template mesh

cannot change  
topology of the  
template mesh



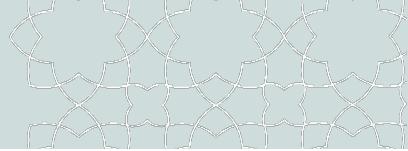
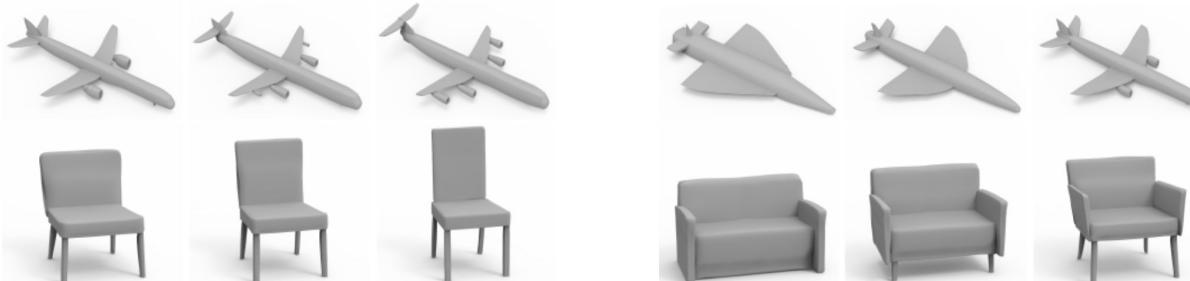
Groueix et al., “AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation”, CVPR 2018

Yang et al., “Foldingnet: Point cloud auto-encoder via deep grid deformation”, CVPR 2018

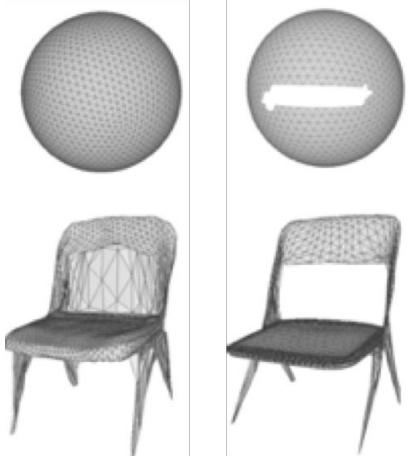
# Task: Conditional Generation

- **From Image to Surface:**

- modify the topology of the template mesh
- part-level deformation

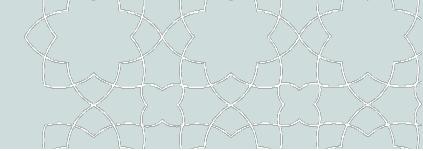


Pan et al., “Deep Mesh Reconstruction from Single RGB Images via Topology Modification Networks.”, ICCV 2019

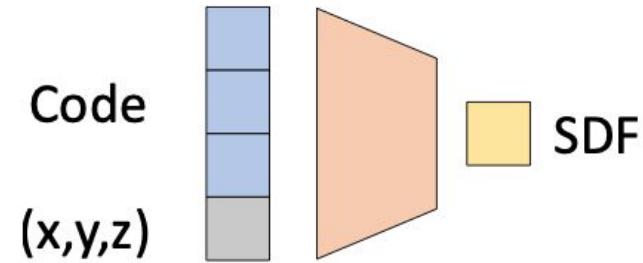
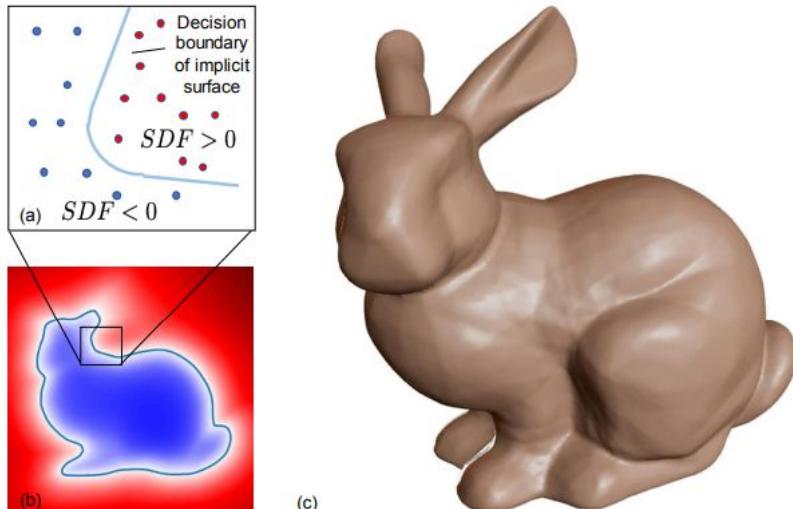


Gao et al., “SDM-NET: Deep generative network for structured deformable mesh.”, Siggraph Asia 2019

# Implicit Surface Reconstruction



- Implicit field function  $F(x)$  (e.g., signed distance)
- Extract the iso-surface  $F(x) = 0$

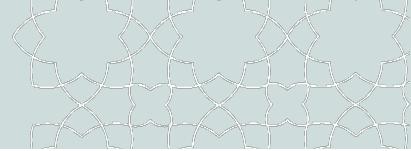


Park et al., "DeepSDF: Learning Continuous Signed Distance Functions for Shape Representation", CVPR 2019

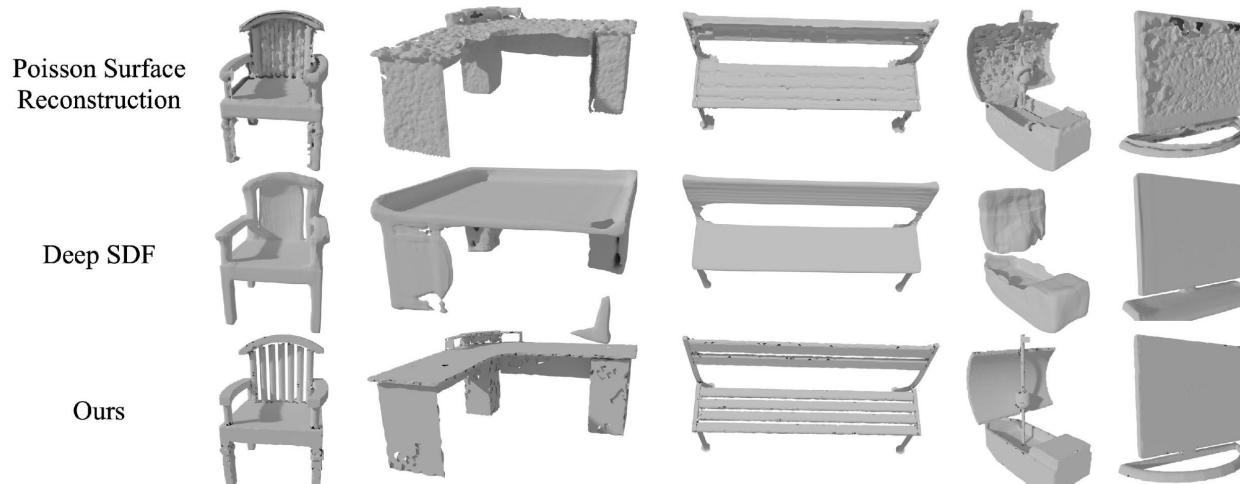
Mescheder et al., "Occupancy Networks: Learning 3D Reconstruction in Function Space", CVPR 2019

Chen et al., "Learning Implicit Fields for Generative Shape Modeling", CVPR 2019

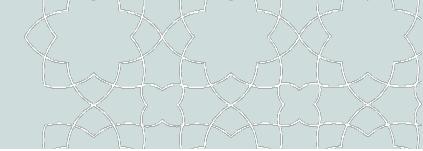
# From Point Cloud to Mesh



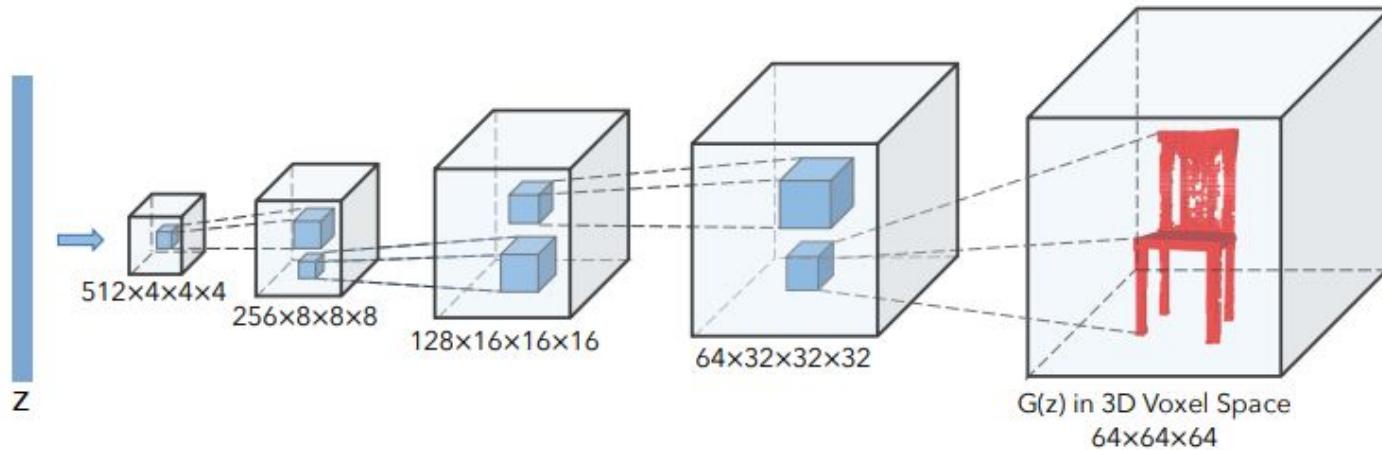
- **Traditional methods:** cannot handle ambiguous structures (e.g., thin structures) when the resolution of point clouds is limited
- **Learning-based methods:** cannot generate fine-grained details or generalize to unseen objects



# Task: Free Generation

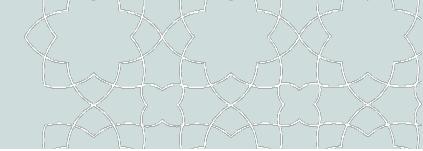


- Volumetric Generation



$$L_{\text{3D-GAN}} = \log D(x) + \log(1 - D(G(z)))$$

# Task: Free Generation

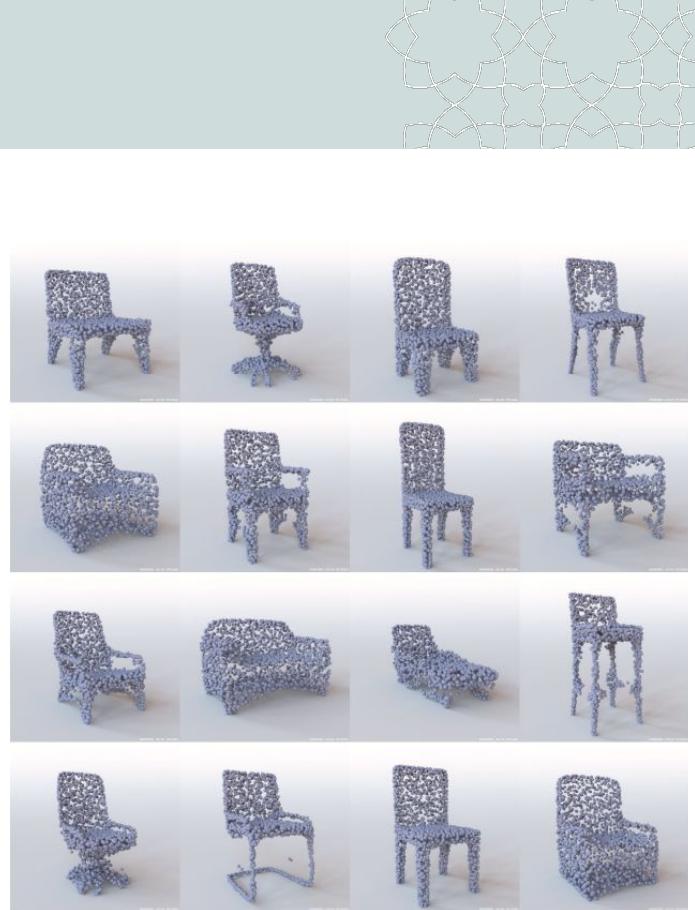
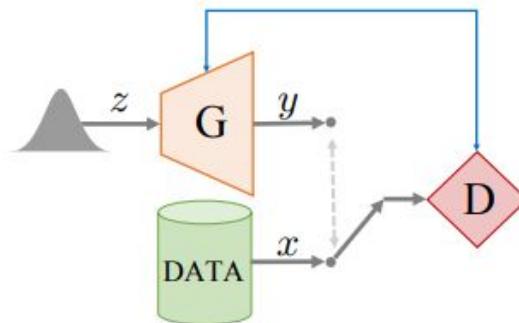


- Volumetric Generation



# Task: Free Generation

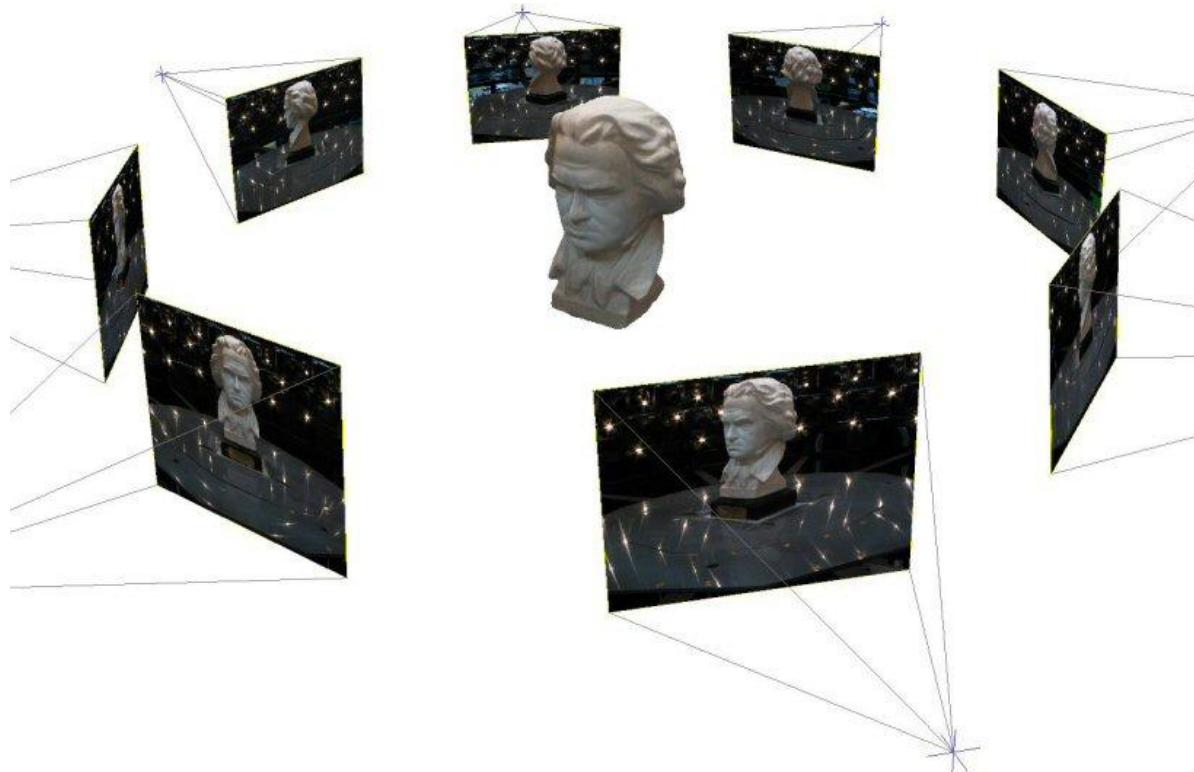
- Point Cloud Generation:
  - FC layer as generator
  - PointNet as discriminator
  - WGAN



# Task: Free Generation: Current challenges

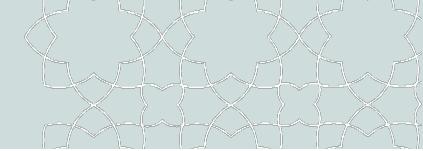
- Still cannot generate high quality local details
- Still hard to generate complex structures
- For point clouds, using strong classifiers (than PointNet) as discriminator is highly tricky

# Task: Multi View Stereo Reconstruction

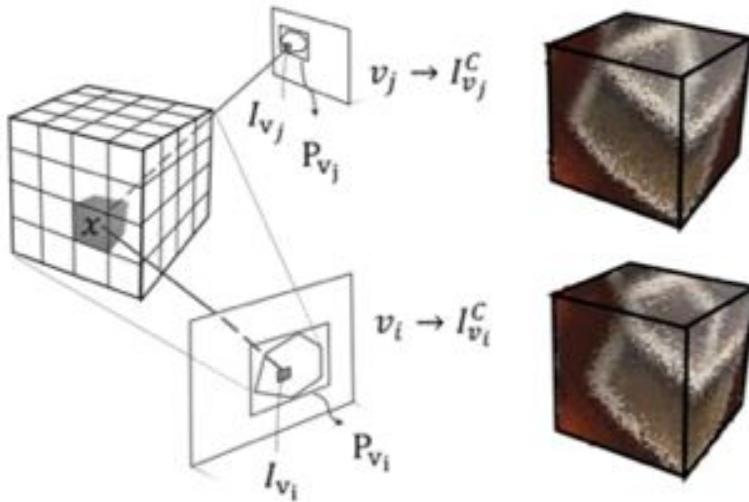


**Covered methods:** SurfaceNet, LSM, GC-Net, MVSNet, R-MVSNet, PointMVSNet, BA-Net

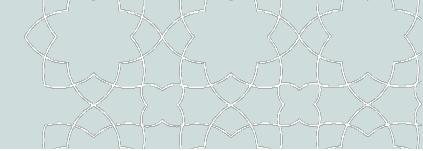
# Task: Multi View Stereo Reconstruction



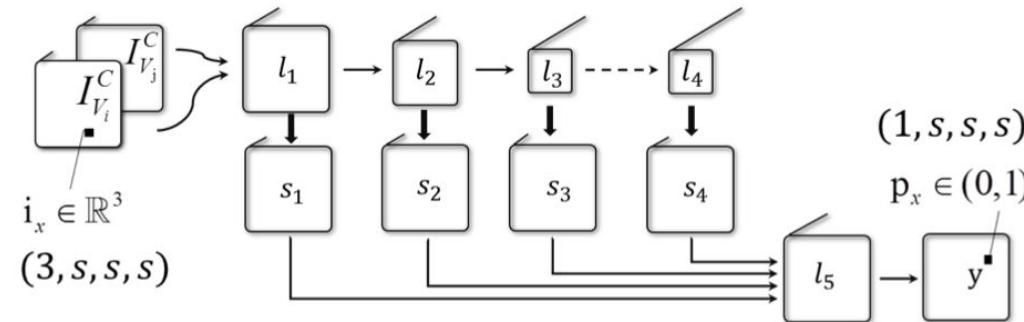
- Surface Reconstruction as Voxel Occupancy Prediction
  - project along viewing rays to build colored voxel cubes.



# Task: Multi View Stereo Reconstruction

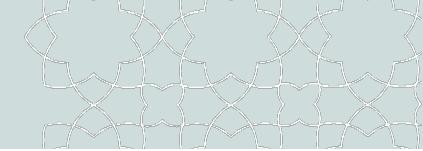


- Surface Reconstruction as Voxel Occupancy Prediction
  - Predict the surface confidence for each voxel



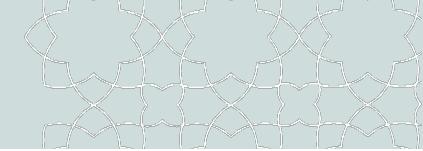
$$\begin{aligned} L(I_{v_i}^C, I_{v_j}^C, \hat{S}^C) = & \\ - \sum_{x \in C} \{ & \alpha \hat{s}_x \log p_x + (1 - \alpha)(1 - \hat{s}_x) \log(1 - p_x) \} \end{aligned}$$

# Task: Multi View Stereo Reconstruction

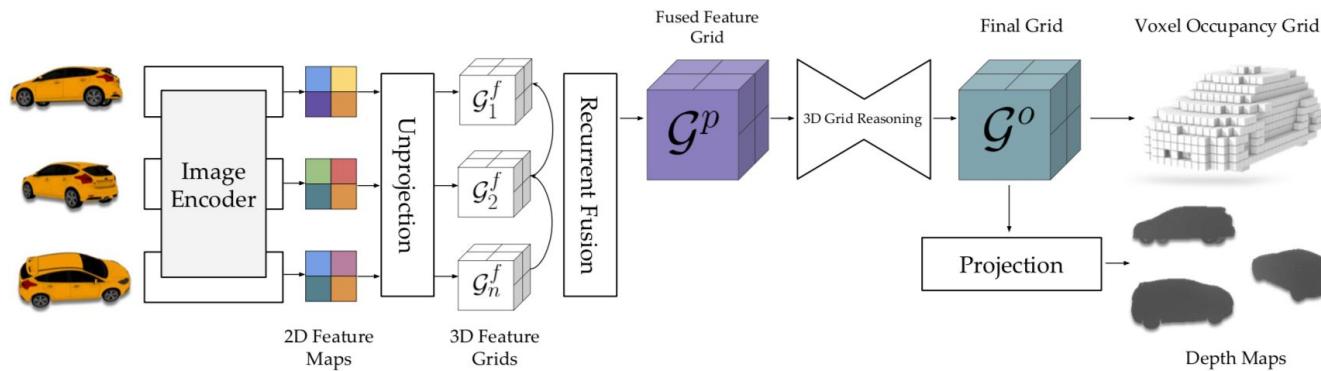


- Surface Reconstruction as Voxel Occupancy Prediction
  - Limitations:
    - Pre-computed voxel cubes can only take RGB colors at coarse resolution
    - Voxel binarization introduces quantization errors.

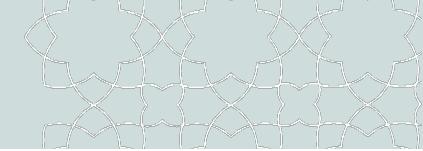
# Task: Multi View Stereo Reconstruction



- Learning-Based Stereopsis
  - End-to-end learning of deep features for each image pixel
  - Unproject image features into 3D grids



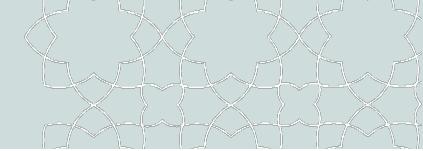
# Task: Multi View Stereo Reconstruction



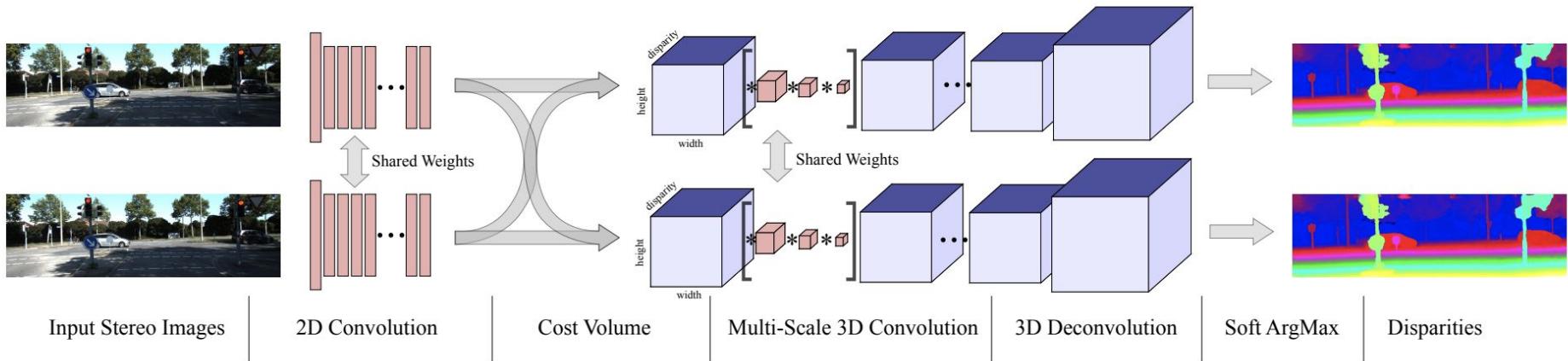
- Learning-Based Stereopsis
  - Still very coarse resolution ( $32 \times 32 \times 32$ ) due to volumetric representation.



# Task: Multi View Stereo Reconstruction

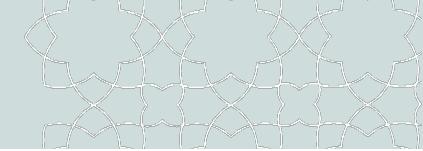


- View-aligned cost-volume construction
- Differentiable soft-argmin to achieve sub-pixel accuracy

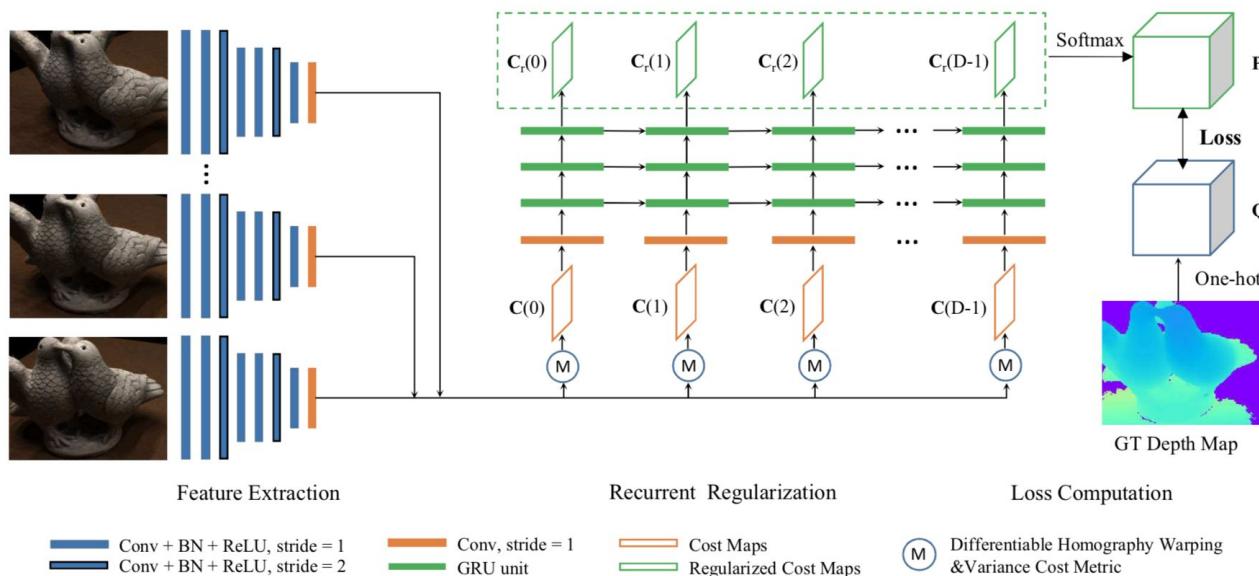


$$\text{soft argmin} := \sum_{d=0}^{D_{\max}} d \times \sigma(-c_d)$$

# Task: Multi View Stereo Reconstruction

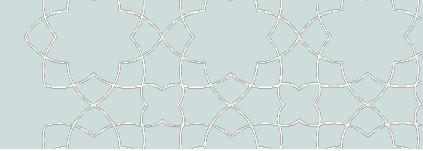


- **Idea 1:** Slide-by-Slide Processing of Cost Volume by Recurrent Neural Network

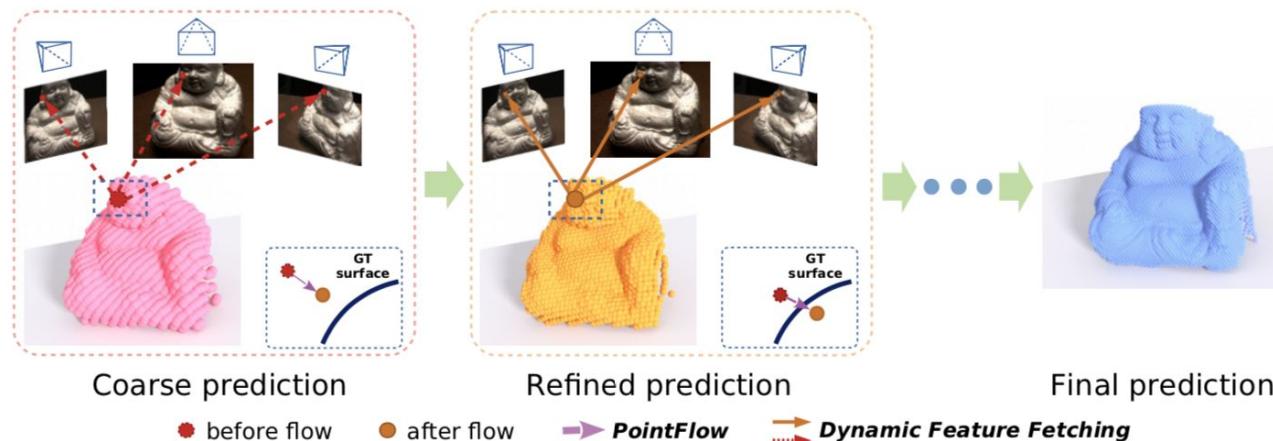


The cost volume is sequentially regularized along the depth direction.

# Task: Multi View Stereo Reconstruction



- **Idea 2: Point-based MVS**
  - Point-based representation for computational efficiency.
  - Iteratively update the location of points and spawn more points.
  - More flexible and accurate.

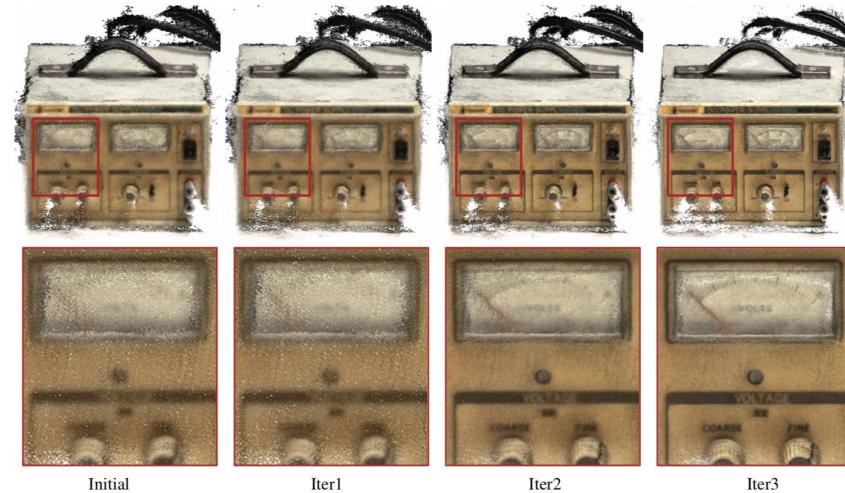


# Task: Multi View Stereo Reconstruction

- **Idea 2: Point-based MVS**

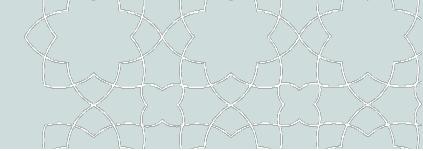
- Iterative refinement

- Results on DTU benchmark



Iter.	Acc. (mm)	Comp. (mm)	Overall (mm)	0.5mm <i>f-score</i>	Depth Map Res.	Depth Interval (mm)	GPU Mem. (MB)	Runtime (s)
-	0.693	0.758	0.726	47.95	160×120	5.30	<b>7219</b>	<b>0.34</b>
1	0.674	0.750	0.712	48.63	160×120	5.30	7221	0.61
2	0.448	0.487	0.468	76.08	320×240	4.00	7235	1.14
3	<b>0.361</b>	<b>0.421</b>	<b>0.391</b>	<b>84.27</b>	<b>640×480</b>	<b>0.80</b>	8731	3.35
MVSNet[29]	0.456	0.646	0.551	71.60	288×216	2.65	10805	1.05

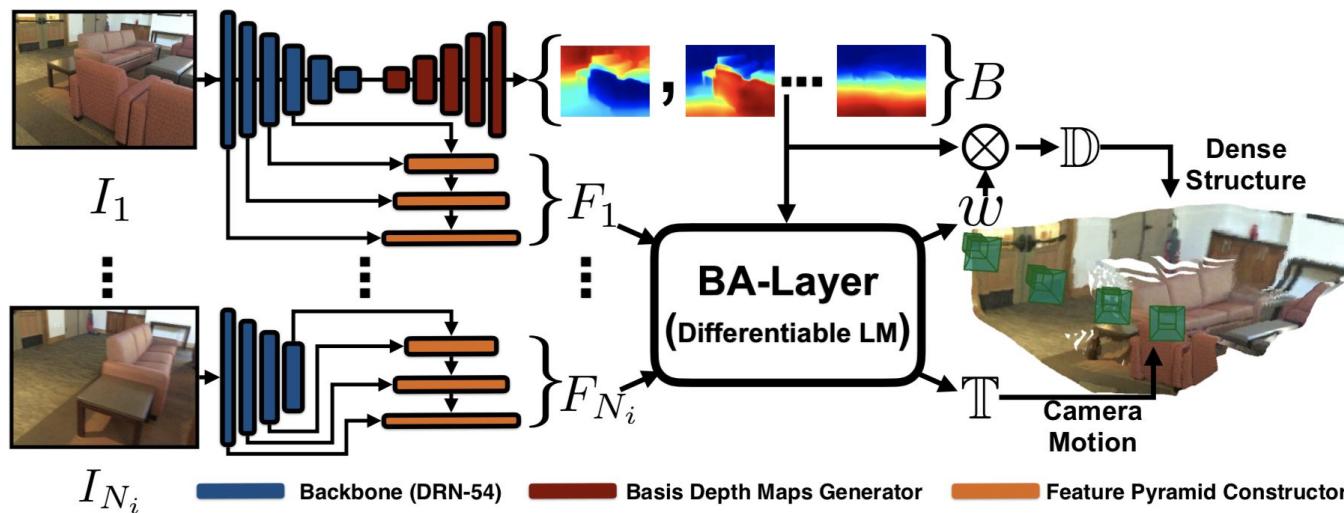
# Task: Learning for SfM



- Above learning-based MVS methods all **assume camera poses are available**
- What if not?
  - Classic 3D: Bundle Adjustment
- **Learning-based bundle adjustment**

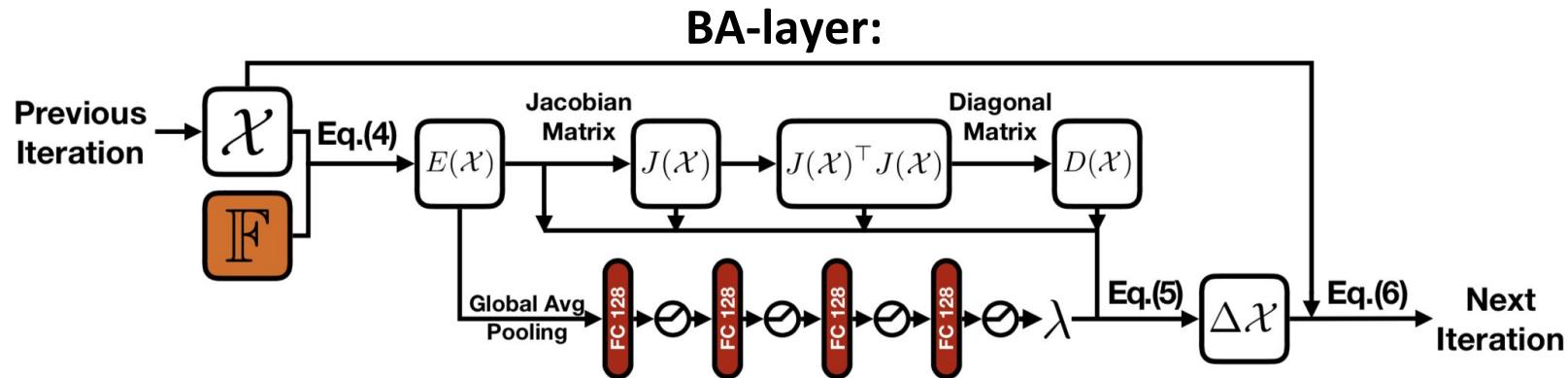
# BA-Net: Dense Bundle Adjustment Network

- End-to-end pipeline for SfM with differentiable bundle adjustment.

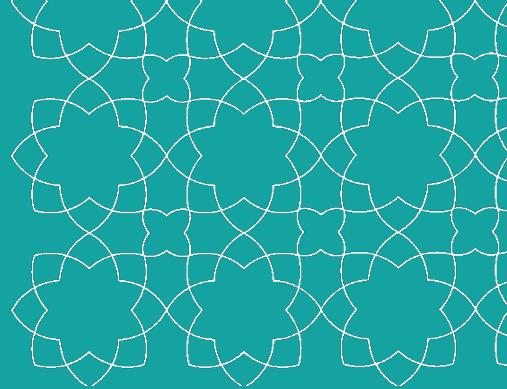


# BA-Net: Dense Bundle Adjustment Network

- Differentiable LM algorithm:
  - Iterative update as rollout of network layers
  - Use network to predict the damping factor lambda.



$$\Delta \mathcal{X} = (J(\mathcal{X})^\top J(\mathcal{X}) + \lambda D(\mathcal{X}))^{-1} J(\mathcal{X})^\top E(\mathcal{X}).$$



Prof. Bernard Ghanem

# 3D DEEP LEARNING