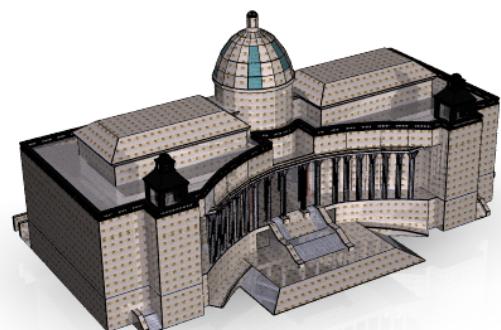
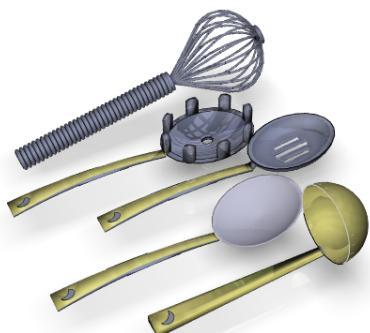
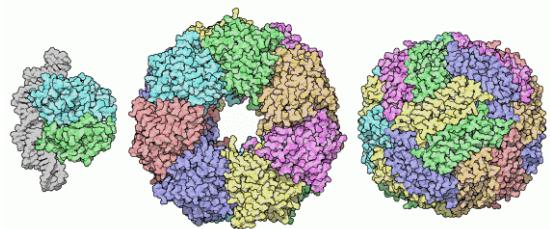
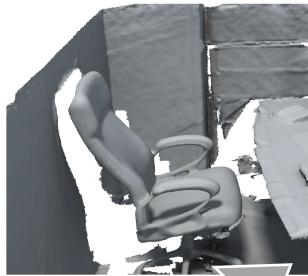


Tutorial on 3D Deep Learning

Hao Su
Jiayuan Gu
Minghua Liu



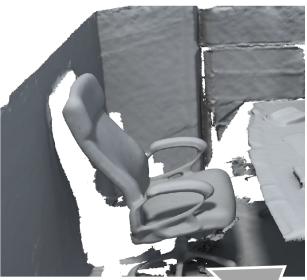
Broad Applications of 3D data



Robotics



Broad Applications of 3D data



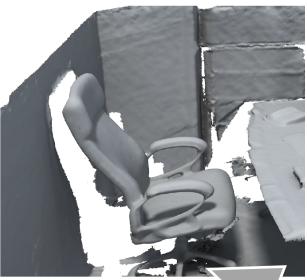
Robotics



**Augmented
Reality**



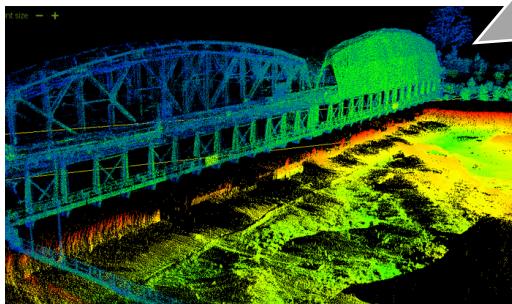
Broad Applications of 3D data



Robotics



Augmented Reality



Autonomous driving



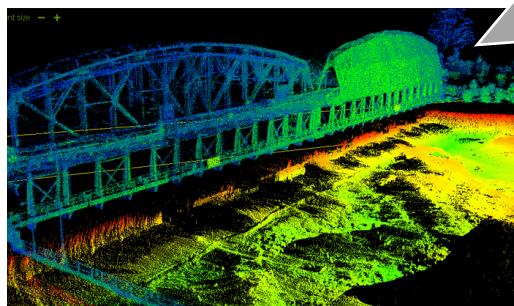
Broad Applications of 3D data



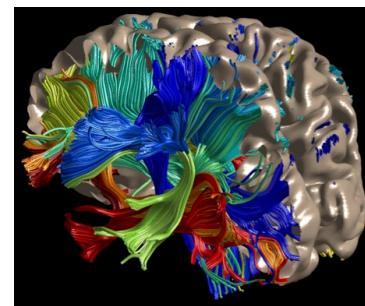
Robotics



Augmented Reality



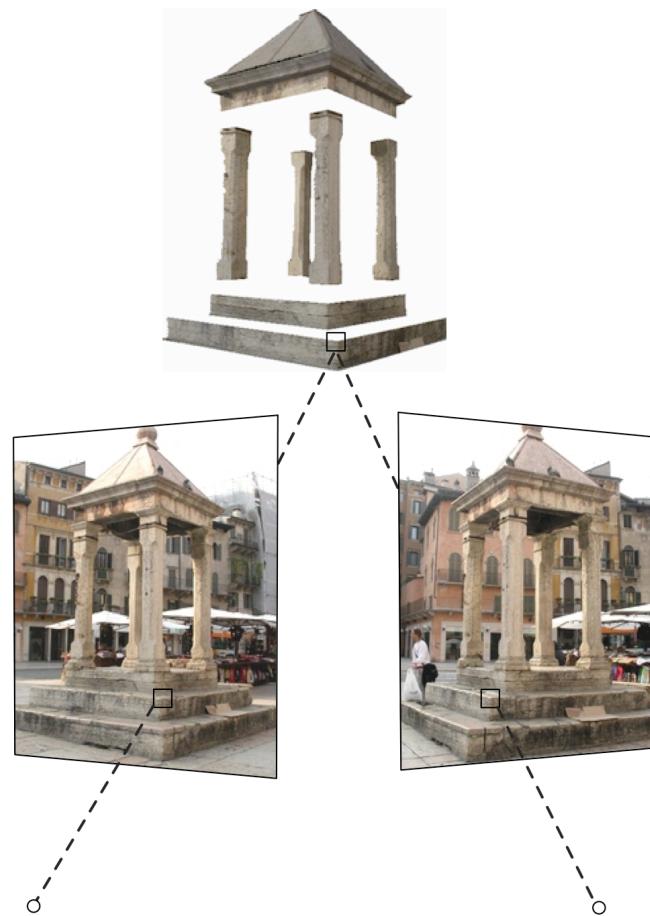
Autonomous driving



Medical Image Processing

Traditional 3D Vision

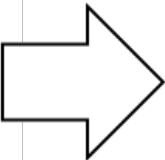
Multi-view Geometry: Physics based



3D Learning: Knowledge Based

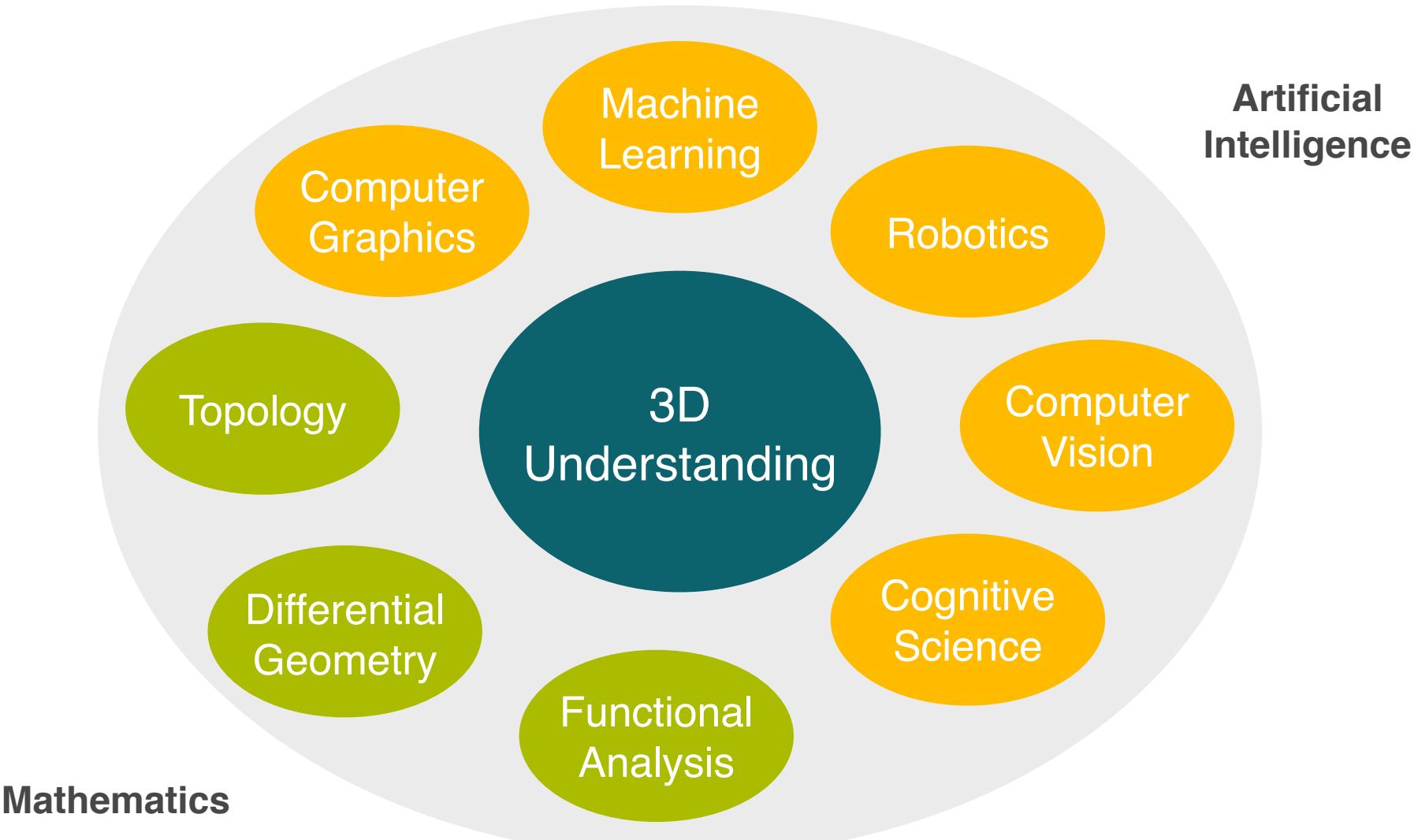


Acquire Knowledge of 3D World by Learning

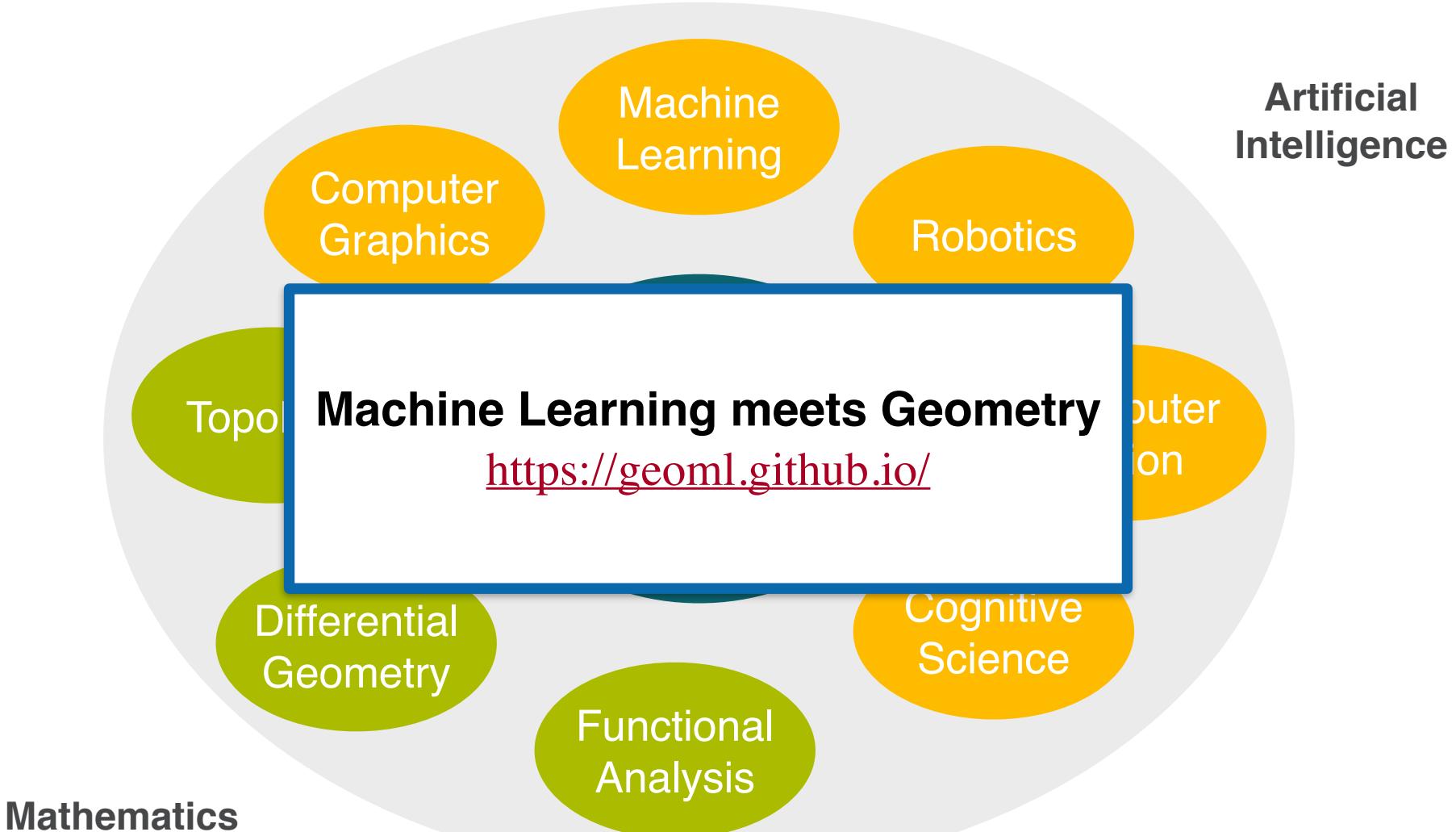


A priori knowledge of
the 3D world

A New Rising Field



A New Rising Field



Instructor Team



Hao Su, Asst Professor



Jiayuan Gu, Ph.D. Student



Minghua Liu, Ph.D. Student

Schedule

- Time: 12:30PM-2:20PM
- Part I: 3D Data, by Hao Su
- Part II: Classification, by Hao Su
- Part II: Segmentation & Detection, by Jiayuan Gu
- Part III: 3D Data Synthesis, by Minghua Liu

Topics

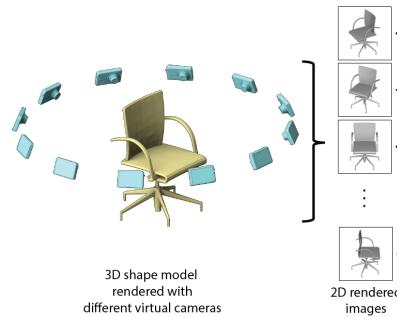
- **3D Data**
- Classification
- Segmentation and Detection
- Reconstruction

The Representation Challenge of 3D Deep Learning

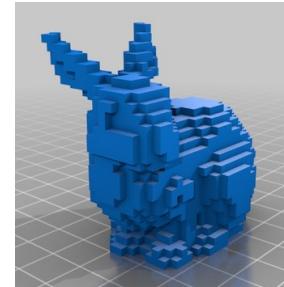
Rasterized form
(regular grids)

Geometric form
(irregular)

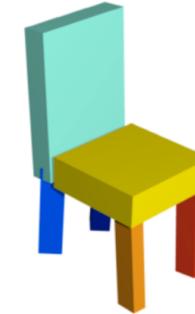
The Representation Challenge of 3D Deep Learning



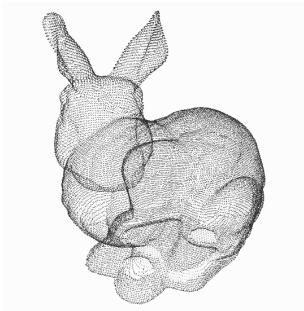
Multi-view



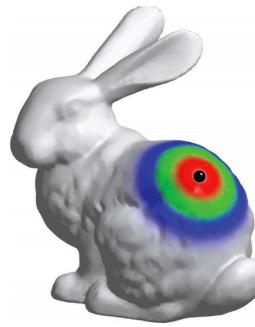
Volumetric



Part Assembly



Point Cloud



Mesh (Graph CNN)

$$F(x) = 0$$

Implicit Shape

Datasets for 3D Objects

Large-scale Synthetic Objects: ShapeNet



3DScan: Consumer-grade 3D scanning (click to open)

ModelNet: absorbed by ShapeNet

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

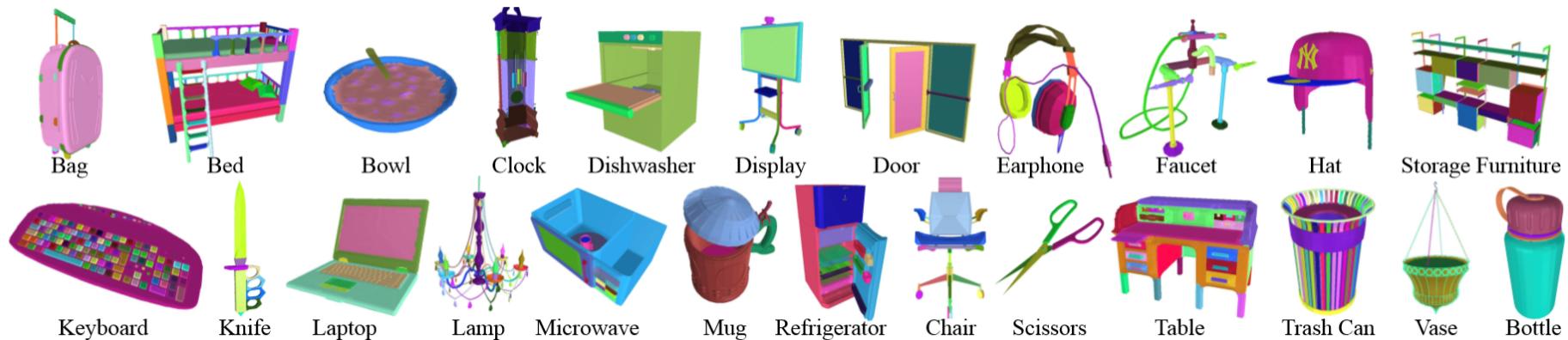
Wu et al., "3D ShapeNets: A deep representation for volumetric shapes", *CVPR 2015*

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

Datasets for 3D Object Parts

Fine-grained Part: PartNet (ShapeNetPart2019)

- Fine-grained (towards mobility)
- Instance-level
- Hierarchical



Datasets for Indoor 3D Scenes

Large-scale Synthetic Scenes: SceneNet

- 3D meshes
- 5M Photorealistic Images

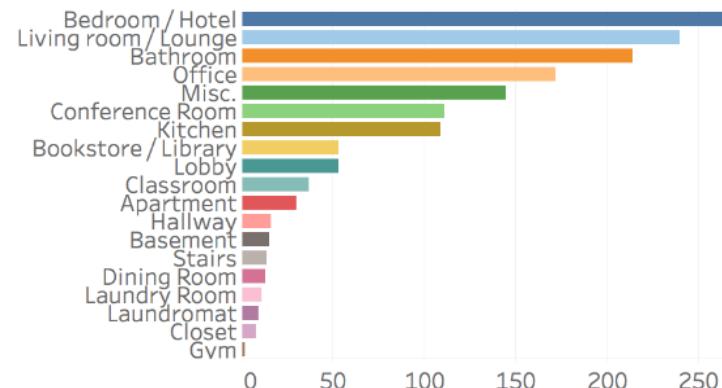


Ankur et al., "Understanding RealWorld Indoor Scenes with Synthetic Data", *CVPR 2016*
McCormac et al., "SceneNet RGB-D: Can 5M Synthetic Images Beat Generic ImageNet Pre-training on Indoor Segmentation?", *ICCV 2017*

Datasets for Indoor 3D Scenes

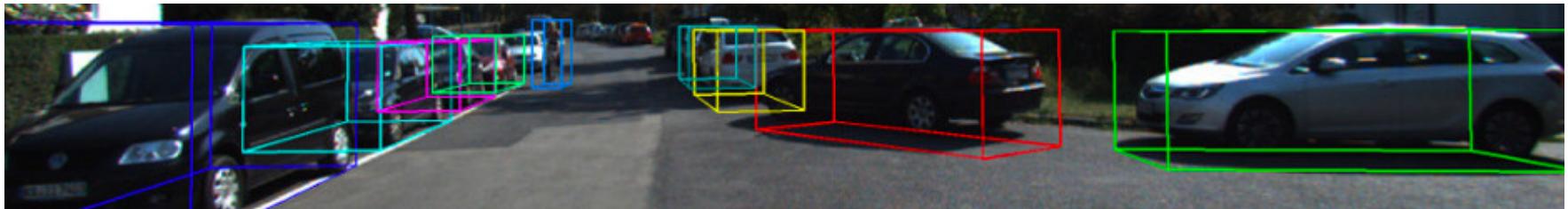
Large-scale Scanned Real Scenes: ScanNet

- 2.5 M Views in 1500 RGBD scans
- 3D camera poses
- surface reconstructions
- Instance-level semantic segmentations



Datasets for Outdoor 3D Scenes

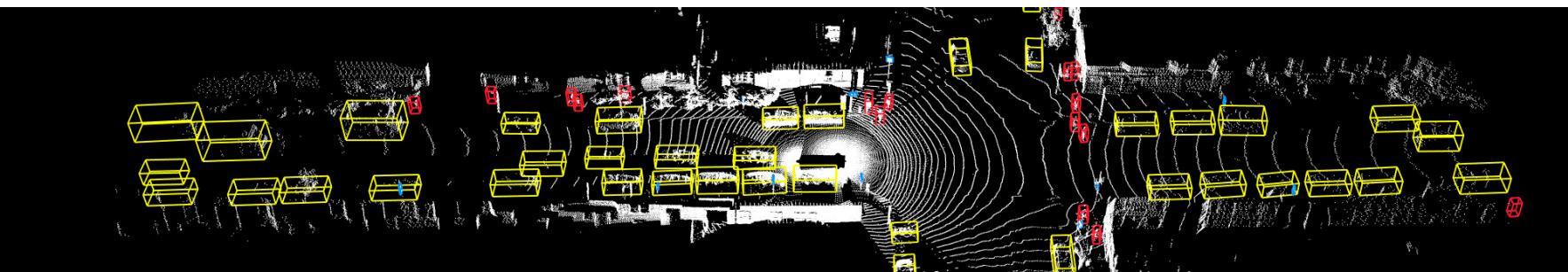
KITTI: LiDAR data, labeled by 3D b.boxes



Semantic KITTI: LiDAR data, labeled per point



Waymo Open Dataset: LiDAR data, labeled by 3D b.boxes



Topics

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

Task: 3D Classification



This is a chair!

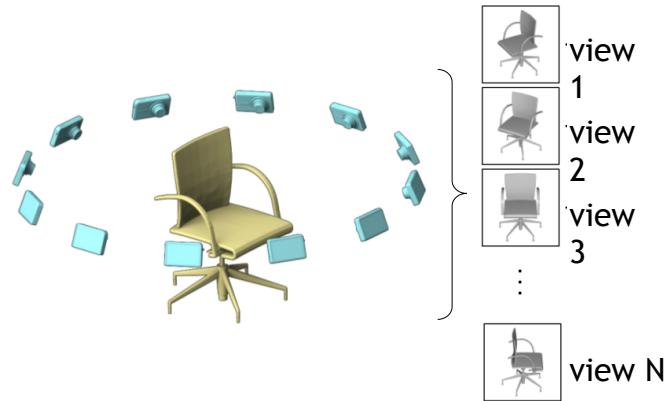
Covered methods: Volumetric CNN, OctNet, O-CNN, SparseConvNet, PointNet, PointNet++, RS CNN, DGCNN, Point ConvNet, KPConv, Monte Carlo Point Convolution, PConv, Multi-View CNN, Spectral CNN, Synchronized Spectral CNN, Spherical CNN

Multi-View CNN

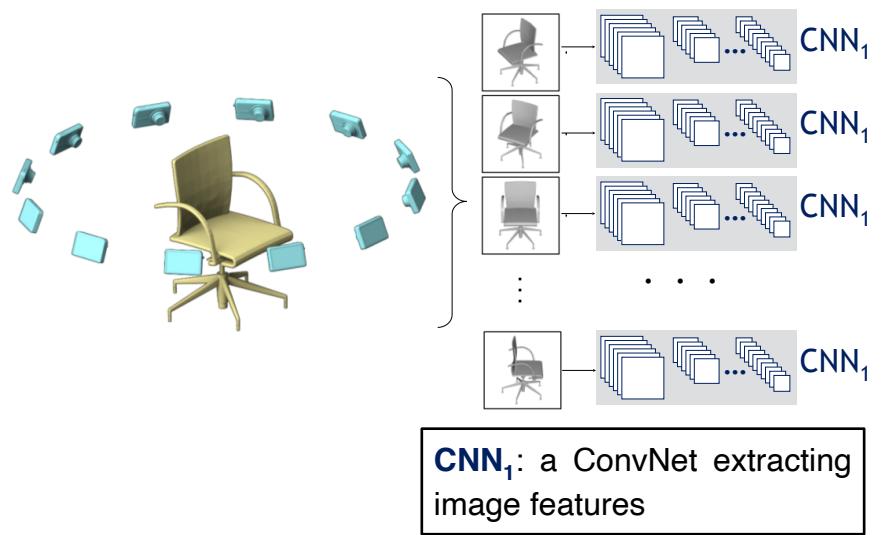
Given an Input Shape



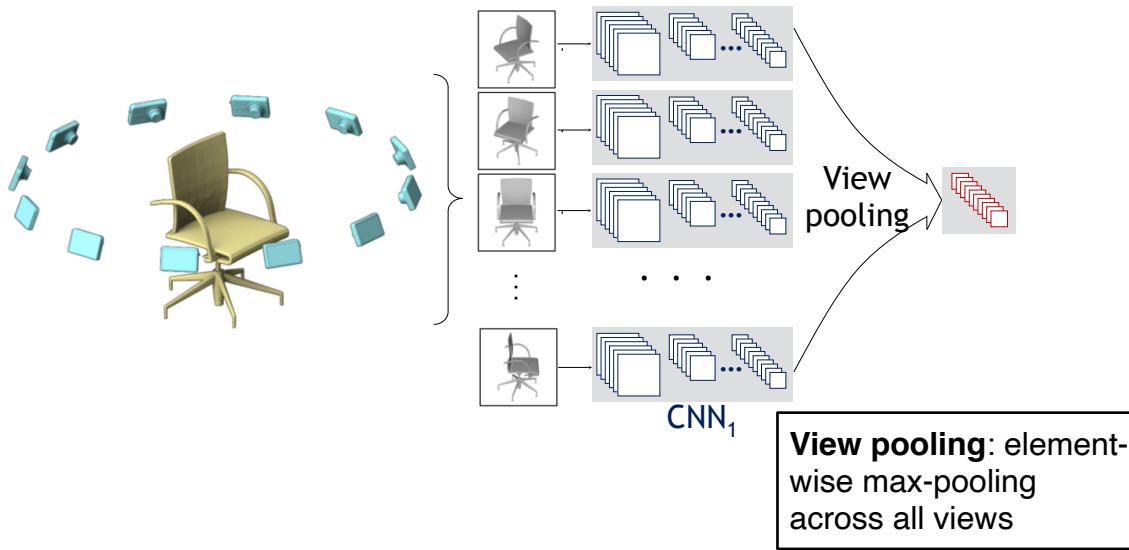
Render with Multiple Virtual Cameras



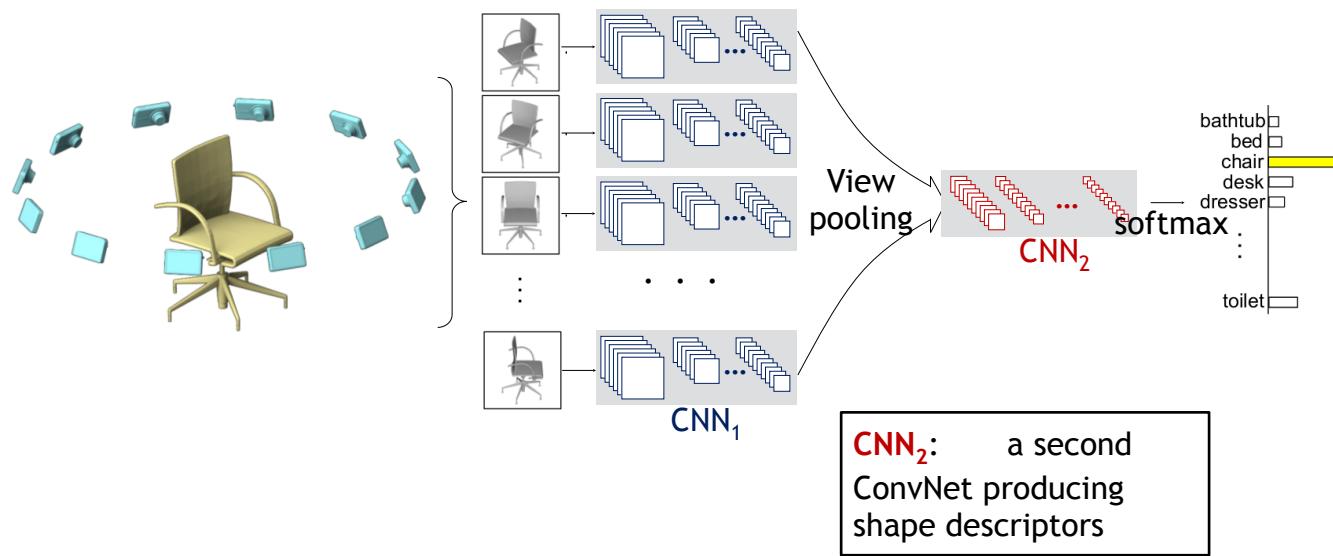
The Rendered Images are Passed through CNN_1 for Image Features



All Image Features are Combined by View Pooling



... and then Passed through CNN_2 and to Generate Final Predictions



Experiments – Classification & Retrieval

Method	Classification (Accuracy)	Retrieval (mAP)
Non-deep {	SPH [16]	68.2%
	LFD [5]	75.5%
	3D ShapeNets [37]	77.3%
	FV, 12 views	84.8%
	CNN, 12 views	88.6%
	MVCNN, 12 views	89.9%
	MVCNN+metric, 12 views	89.5%
	MVCNN, 80 views	90.1%
	MVCNN+metric, 80 views	90.1%
		79.5%

On ModelNet40

[credit: Hang Su]

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

Volumetric CNN

Can we use CNNs without 2D-3D projection?

Straight-forward idea: 3D native convolution

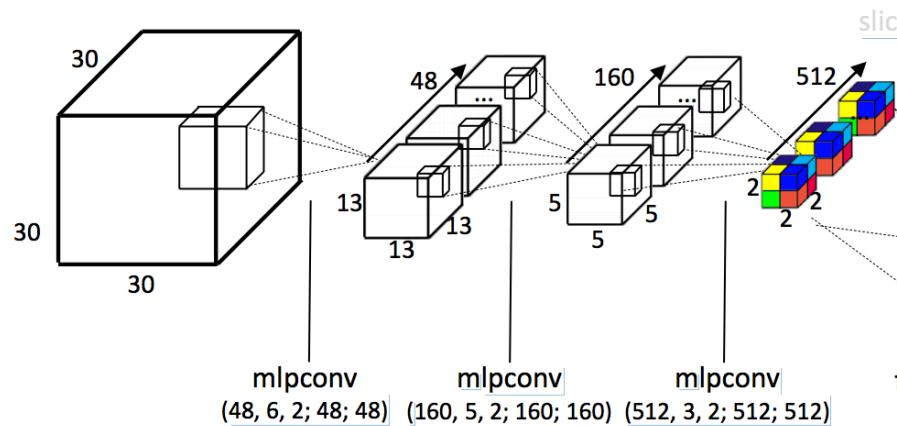
Voxelization

Represent the occupancy of regular 3D grids

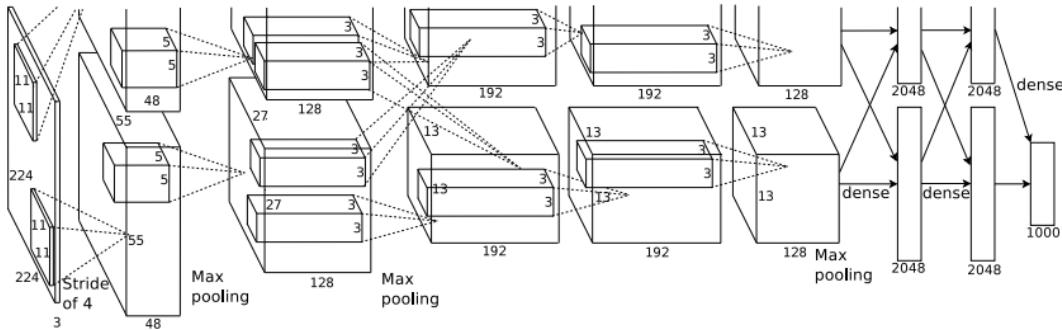


3D CNN on Volumetric Data

3D convolution uses 4D kernels



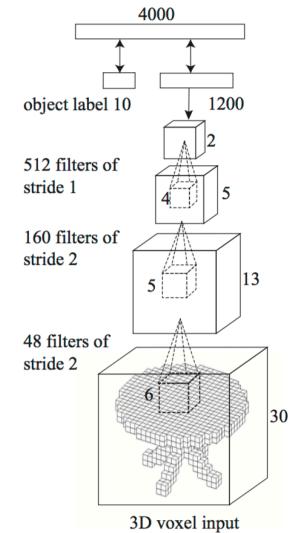
Complexity Issue



AlexNet, 2012

Input resolution: 224x224

$$224 \times 224 = 50176$$

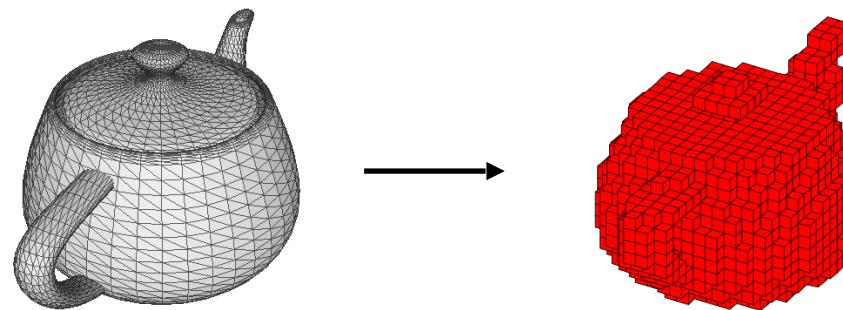


3DShapeNets,
2015

Input resolution: 30x30x30

$$224 \times 224 = 27000$$

Complexity Issue



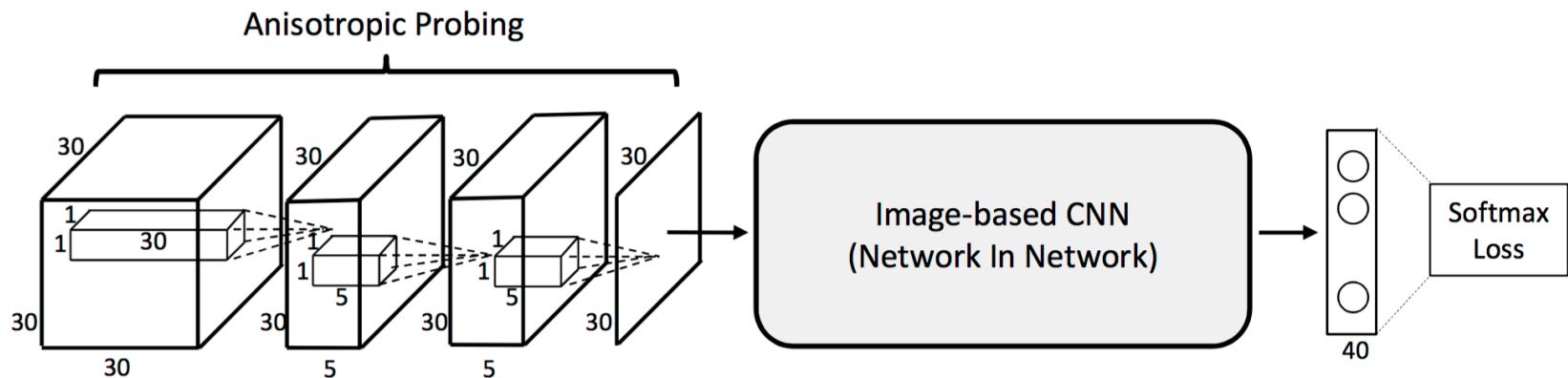
Polygon Mesh

Occupancy Grid
 $30 \times 30 \times 30$

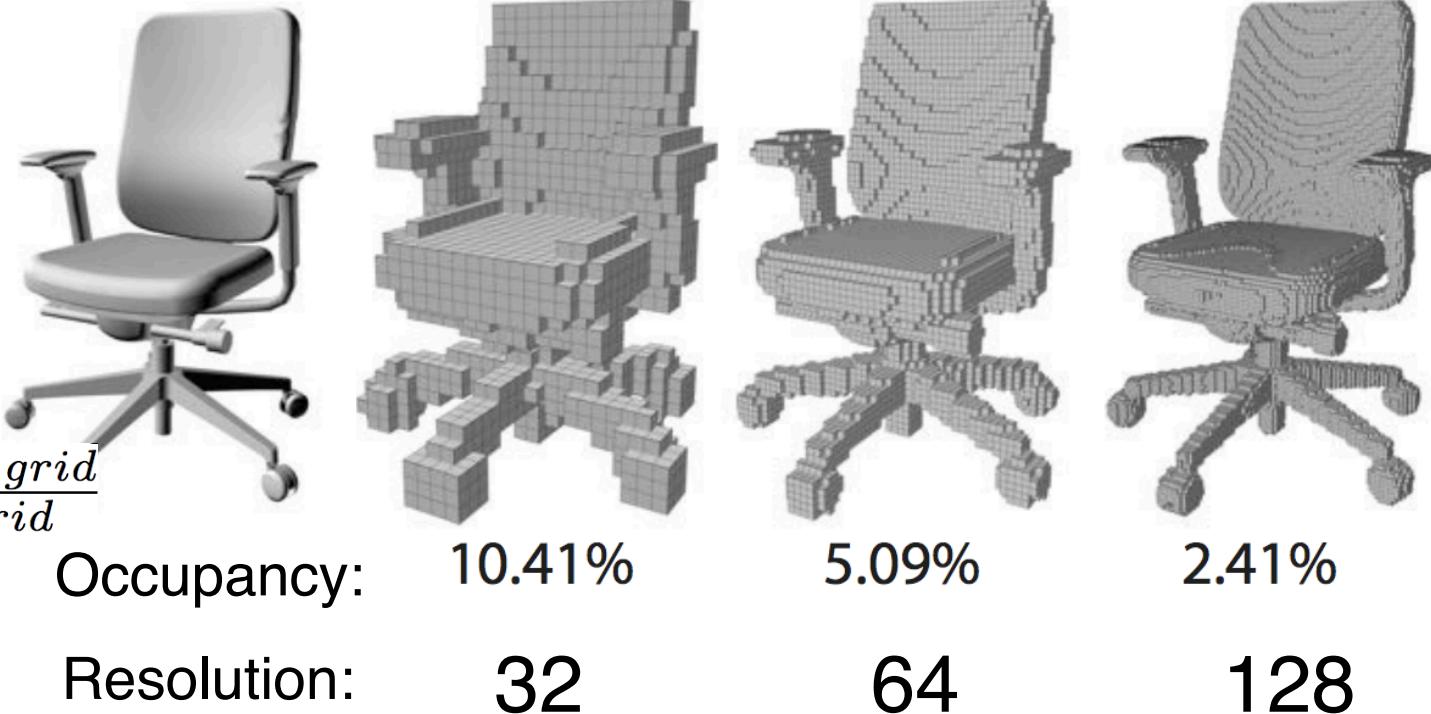
Information loss in voxelization

Idea 1: Learn to Project

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

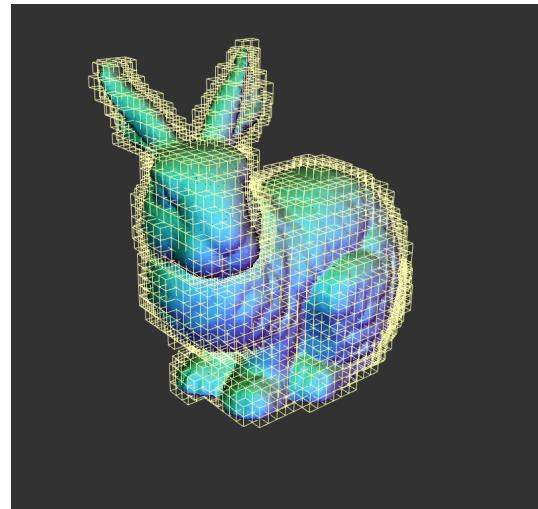
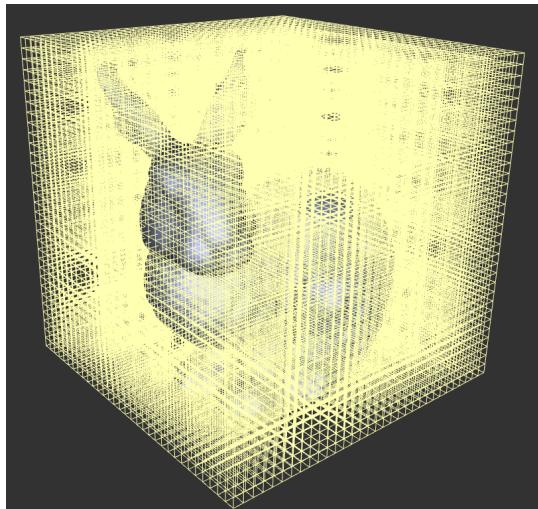


More Principled: Leverage Sparsity of 3D Surface Data



Store only the Occupied Grids

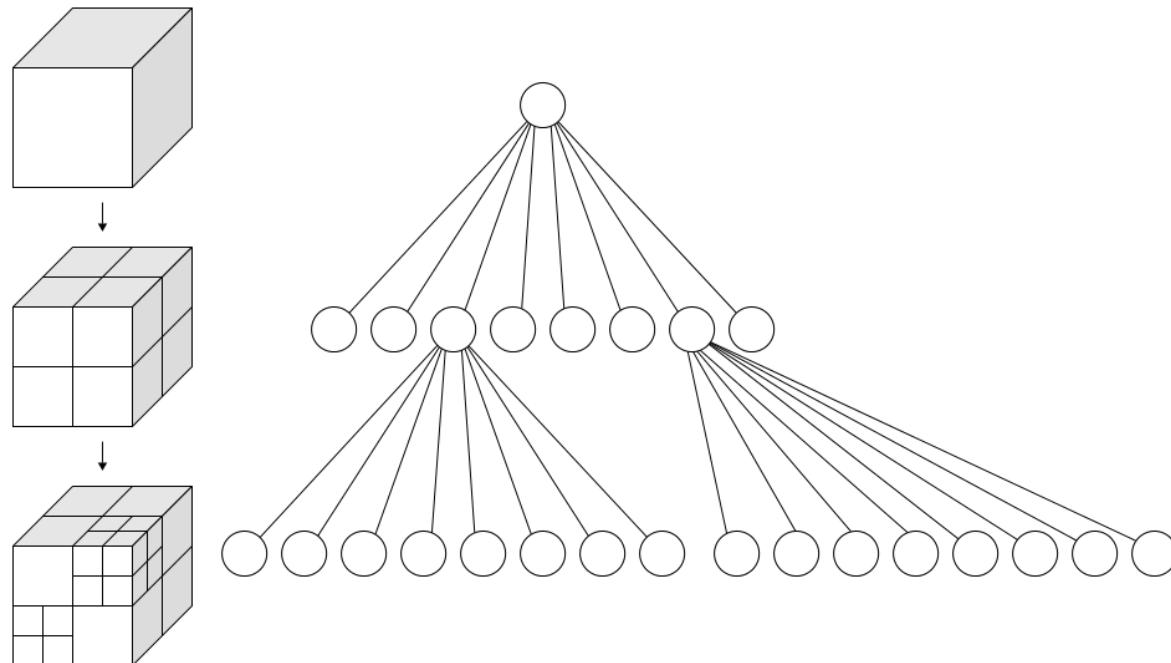
- Store the sparse surface signals
- Constrain the computation near the surface



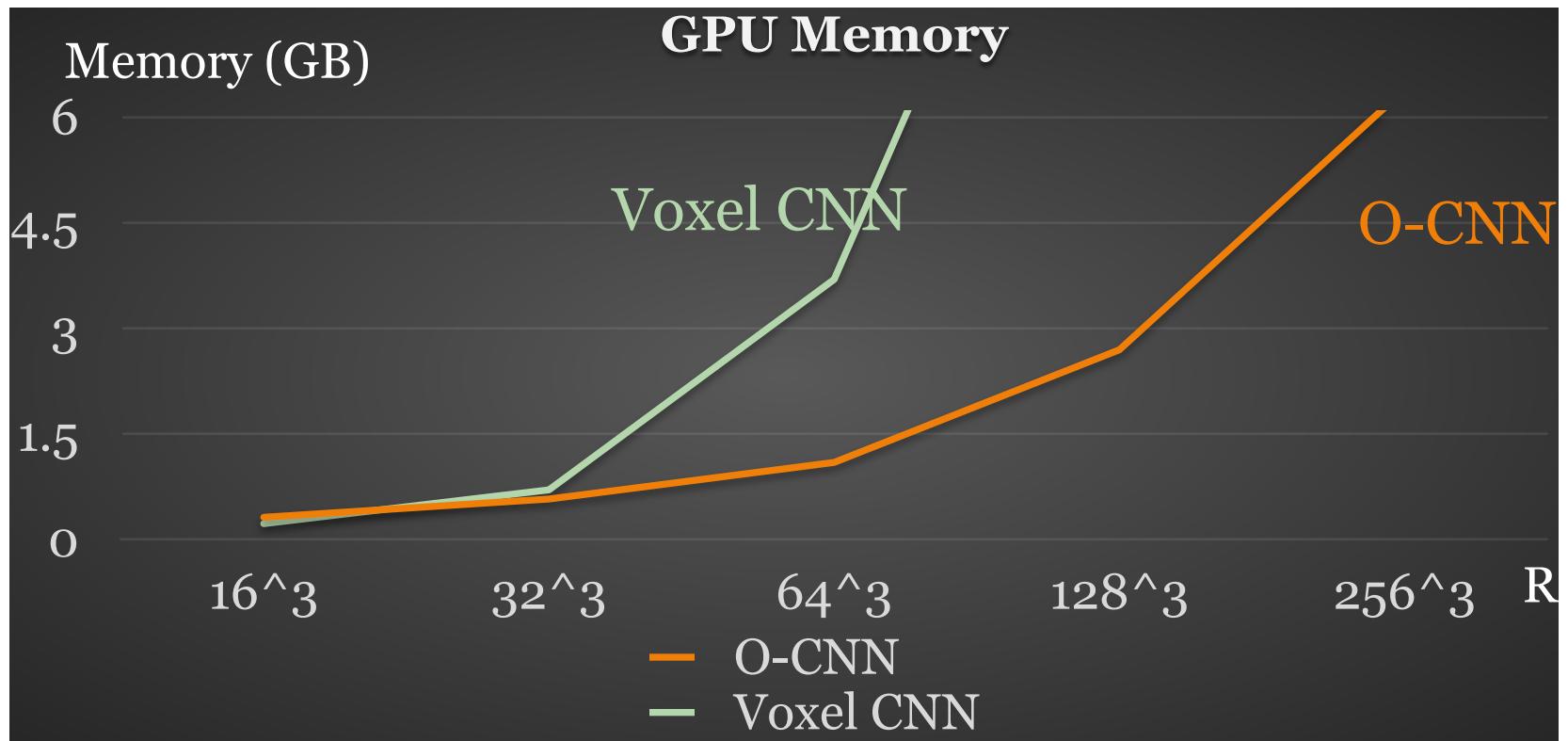
Octree: Recursively Partition the Space

Each **internal node** has exactly eight **children**

Neighborhood searching: Hash table



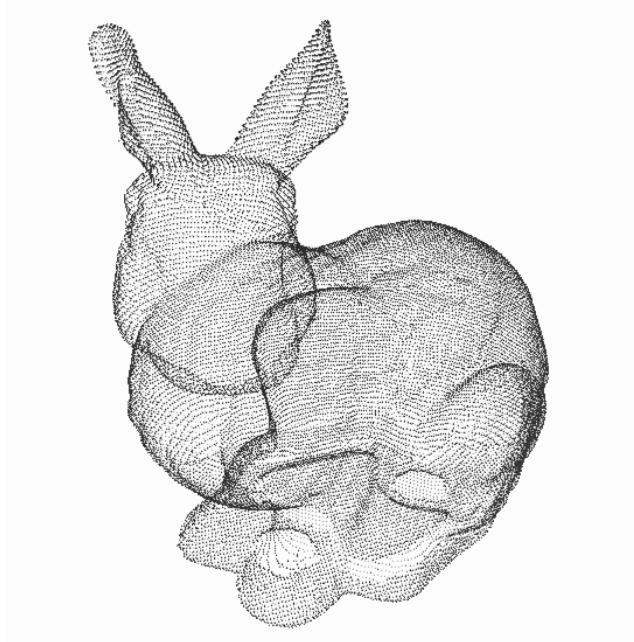
Memory Efficiency



Implementation

- SparseConvNet
 - [https://github.com/facebookresearch/
SparseConvNet](https://github.com/facebookresearch/SparseConvNet)
 - Uses ResNet architecture
 - State-of-the-art for 3D analysis
 - Takes time to train

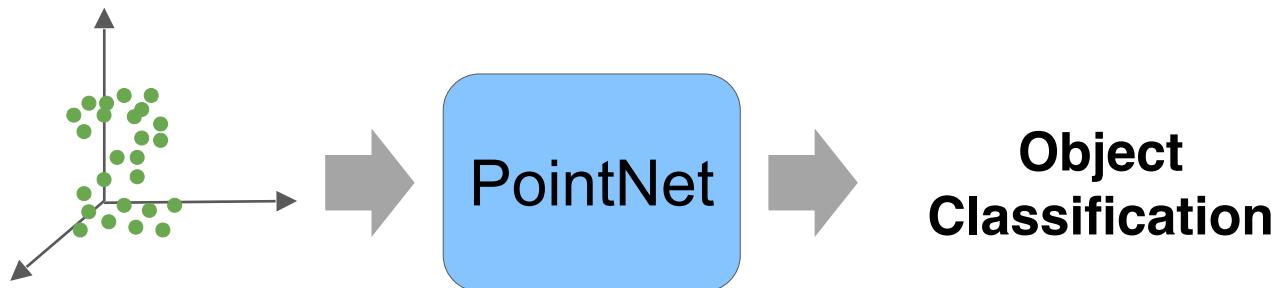
Point Networks



Point cloud
(The most common 3D sensor data)

Directly Process Point Cloud Data

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

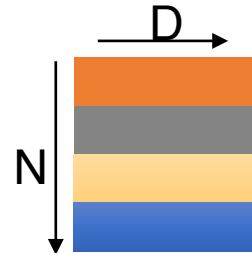


Qi, Charles R., et al. "Pointnet: Deep learning on point sets for 3d classification and segmentation", CVPR
2017

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

Permutation invariance

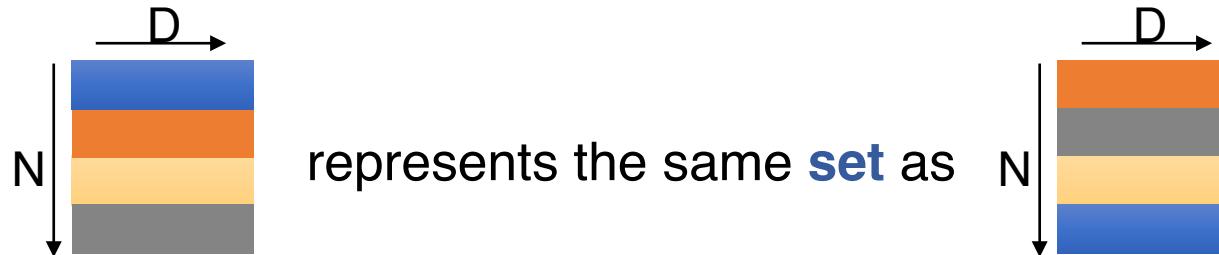
Point cloud: N **orderless** points, each represented by a D dim coordinate



2D array representation

Permutation invariance

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



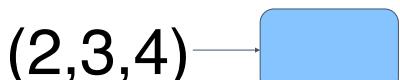
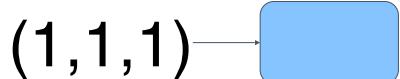
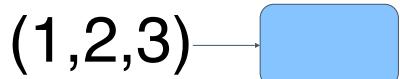
2D array representation

Construct a Symmetric Function

Observe:

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

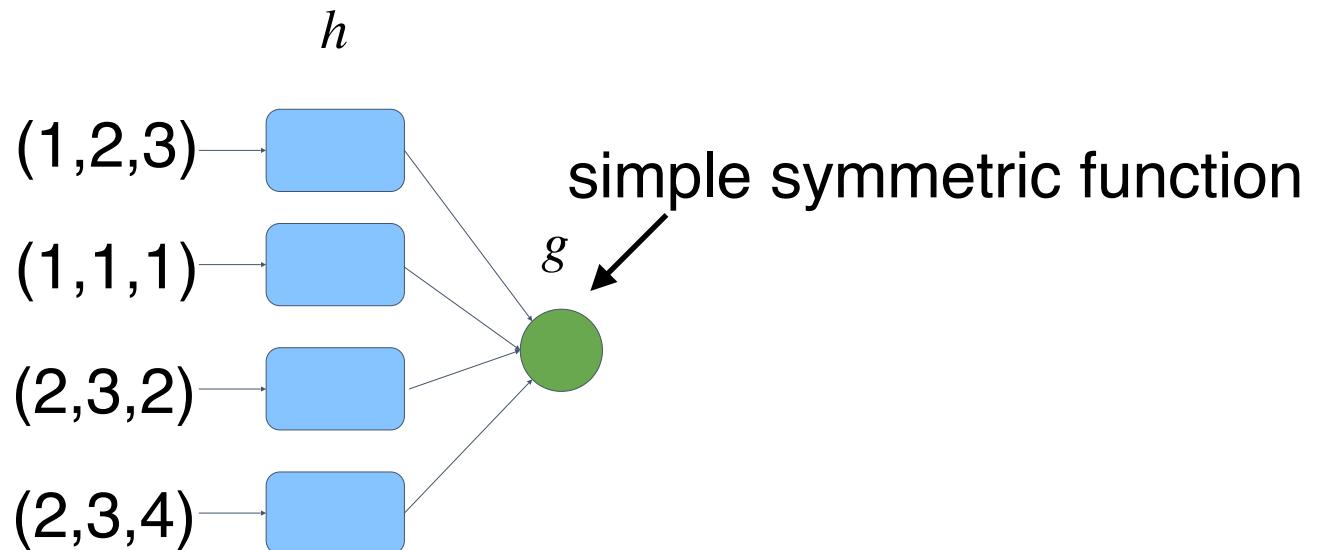
h



Construct a Symmetric Function

Observe:

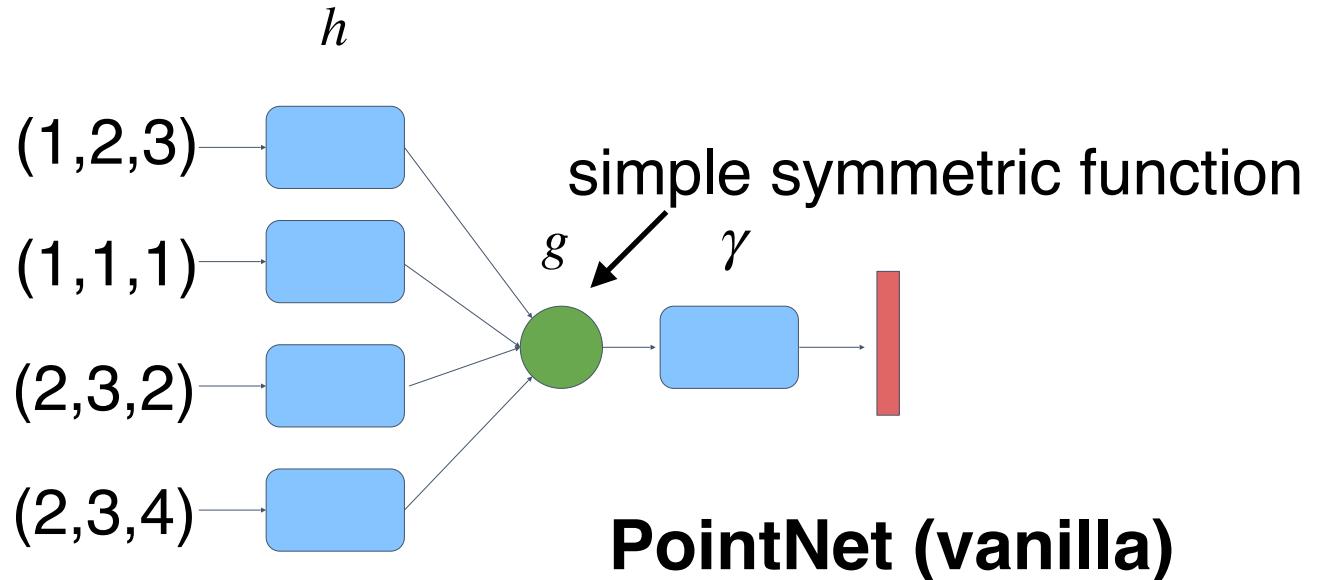
$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$ is symmetric if g is symmetric



Construct a Symmetric Function

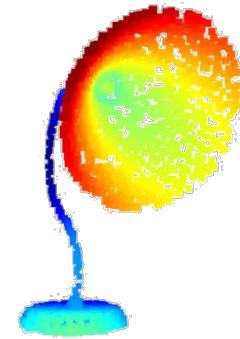
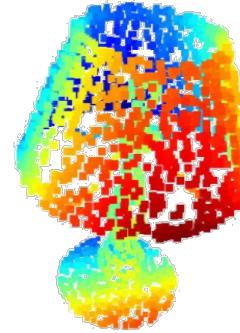
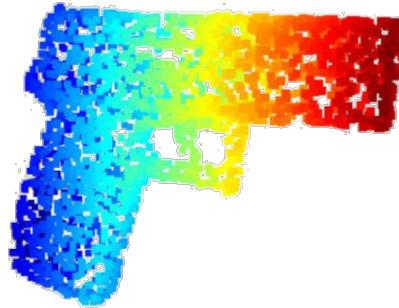
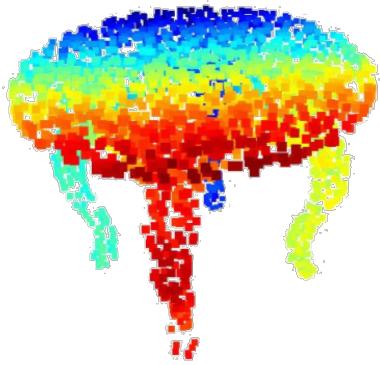
Observe:

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

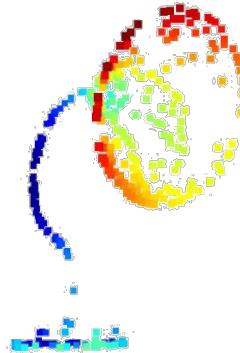
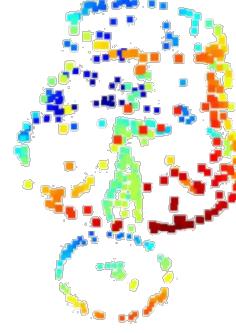
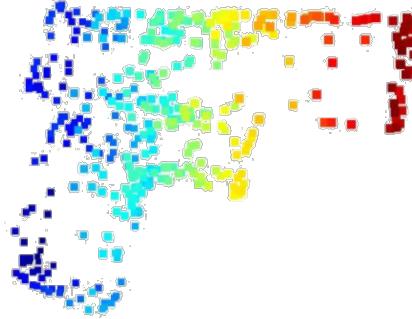
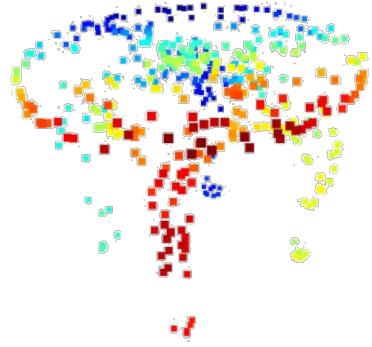


Visualize What is Learned by Reconstruction

Original Shape



Critical Point Sets



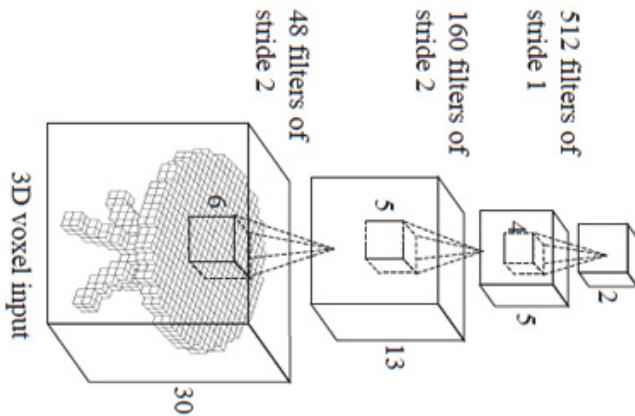
Salient points are discovered!

Limitations of PointNet

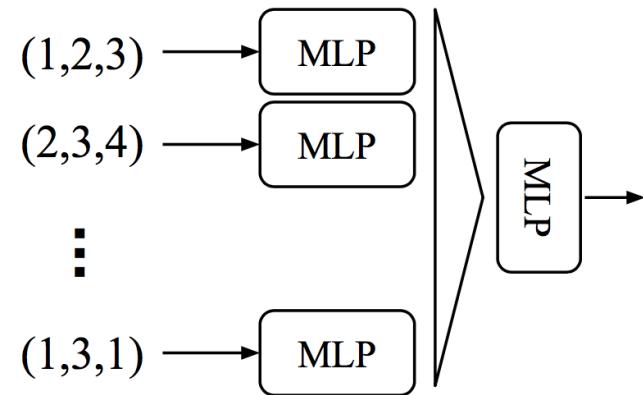
Hierarchical feature learning
Multiple levels of abstraction

v.s.

Global feature learning
Either one point or all points



3D CNN (Wu et al.)



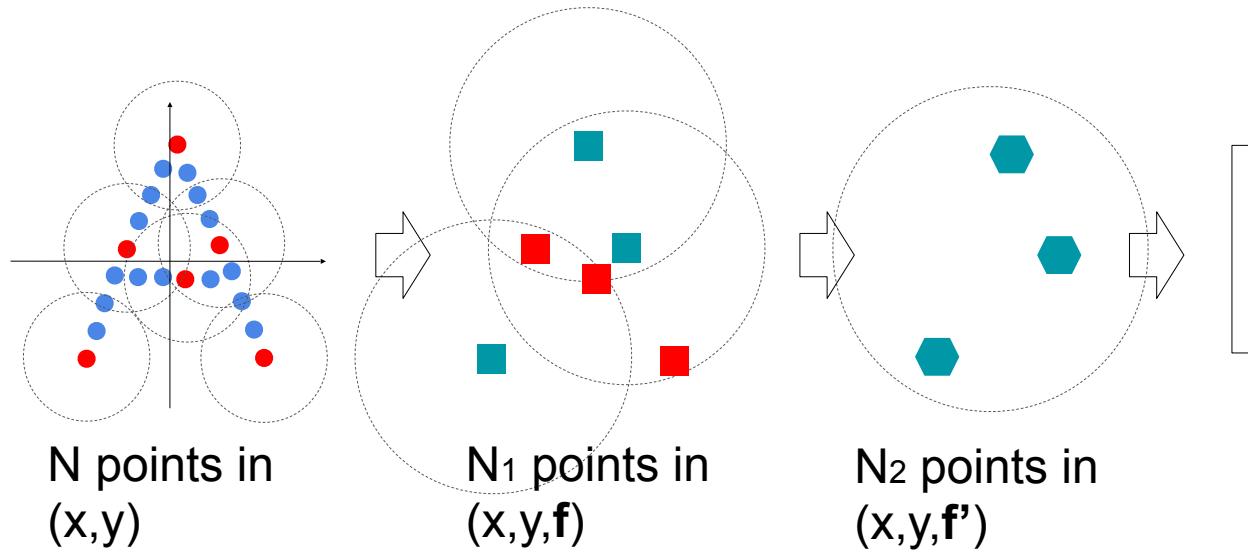
PointNet (vanilla) (Qi et al.)

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

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 v2.0: Multi-Scale PointNet

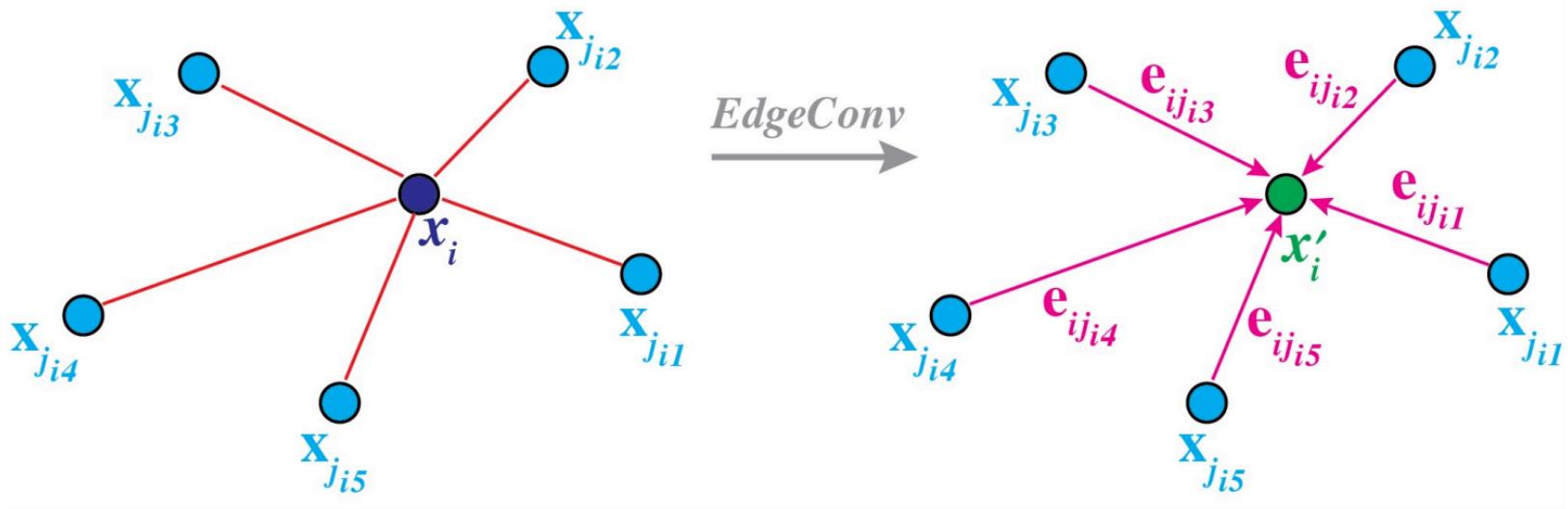


Repeat

- 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



Wang et al., “Dynamic Graph CNN for Learning on Point Clouds”,
Transactions on Graphics, 2019

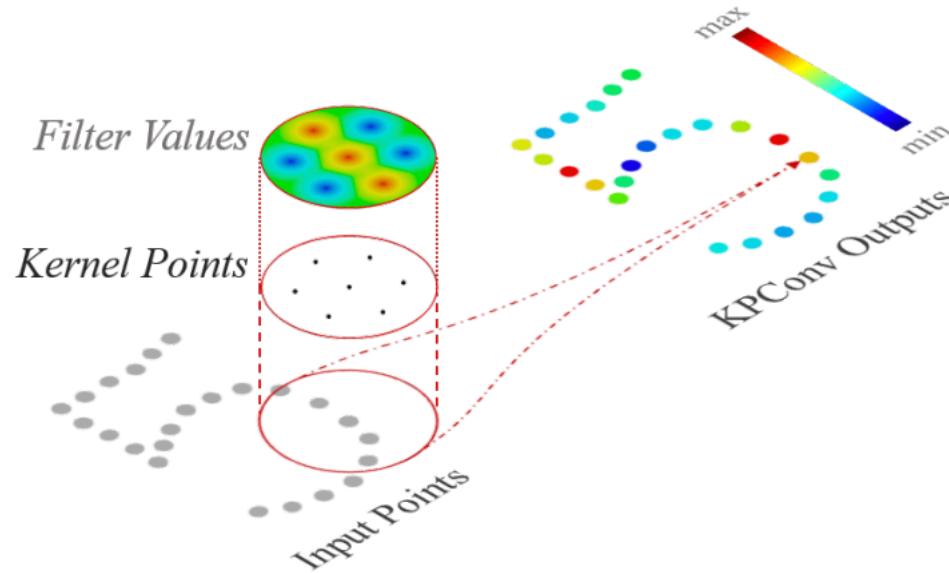
Liu et al., “Relation-Shape Convolutional Neural Network for Point
Cloud Analysis”, *CVPR 2019*

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

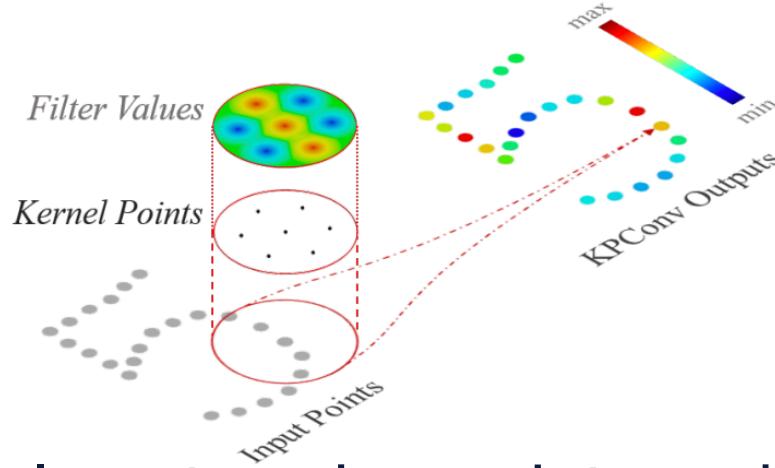
Mathematically Proper Conv. Discretization

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



Interpolated Kernel for Convolution

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



- Learn kernel value at anchor points and interpolate to build continuous kernel

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

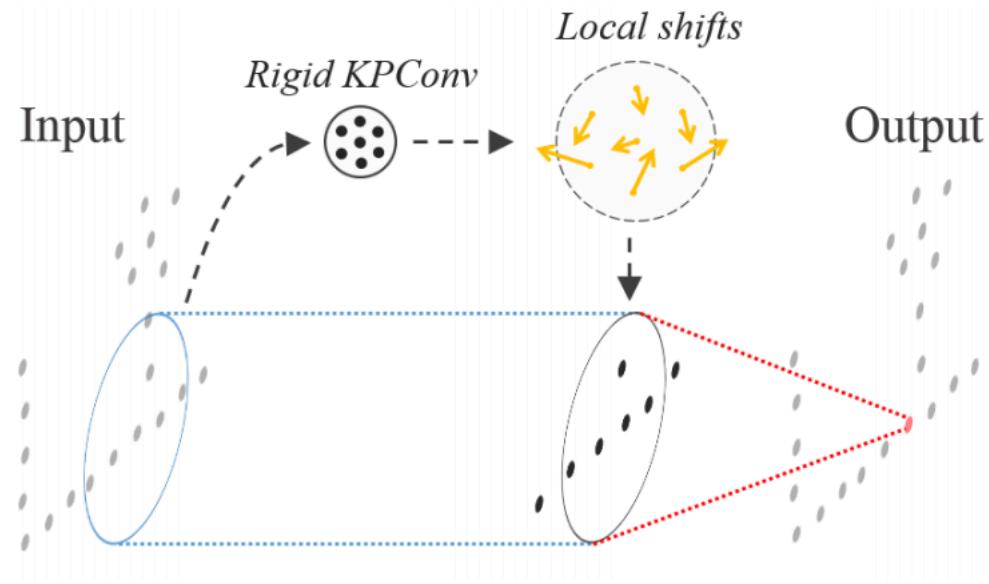
Φ : RBF kernel

Atzmon et al., “Point Convolutional Neural Networks by Extension Operators”, *Trans. on Graphics*, 2018

Thomas et al., “KPConv: Flexible and Deformable Convolution for Point Clouds”, *ICCV* 2019

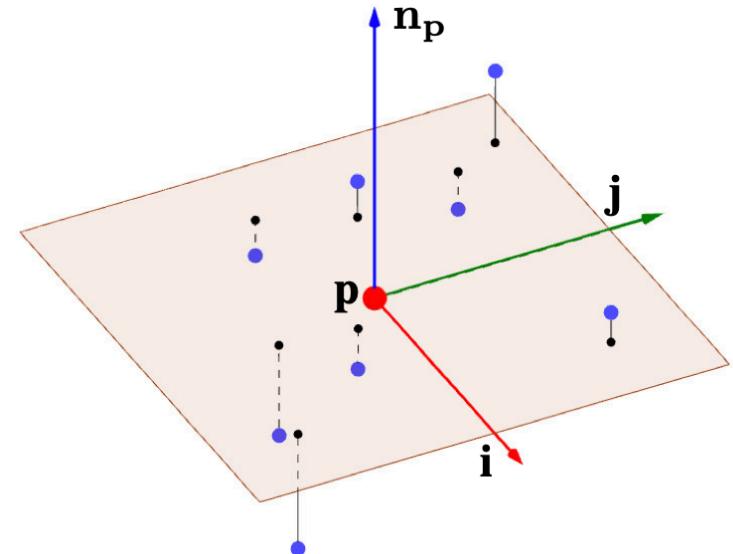
Deformable Kernel for Deformable Objects

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

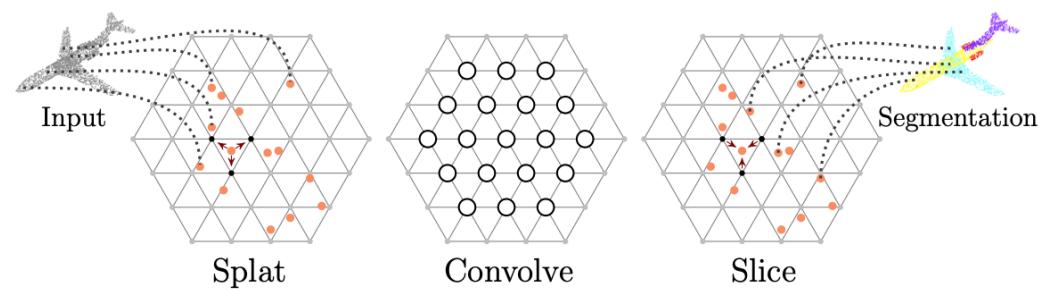


Other Ways to Address 3D Conv

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



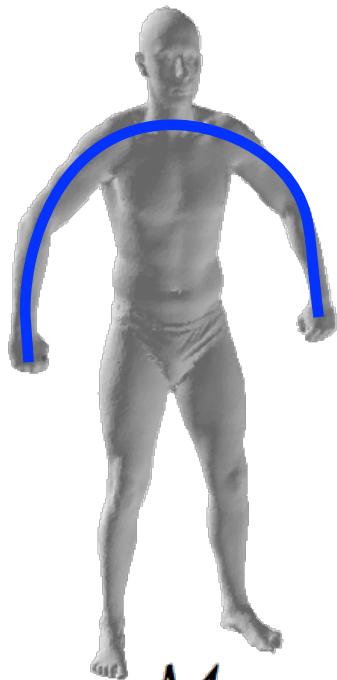
Tatarchenko, Maxim, et al. "Tangent convolutions for dense prediction in 3d", CVPR 2018.



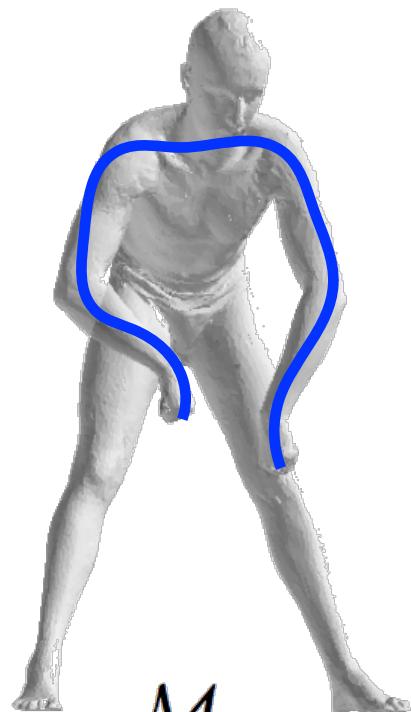
Su, Hang, et al. "Splatnet: Sparse lattice networks for point cloud processing", CVPR 2018

Spectral Convolution

Recognition with Isometric Invariance?



geodesic = intrinsic M_1



isometry = length-preserving transform M_2



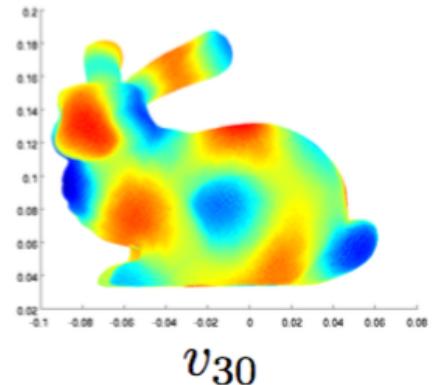
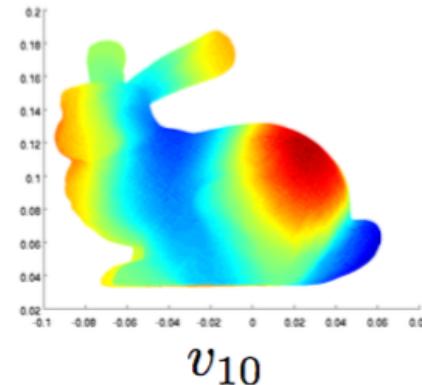
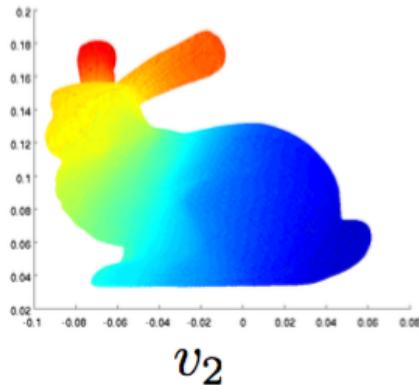
Spectral CNN

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

Challenge: Obtain Fourier Basis

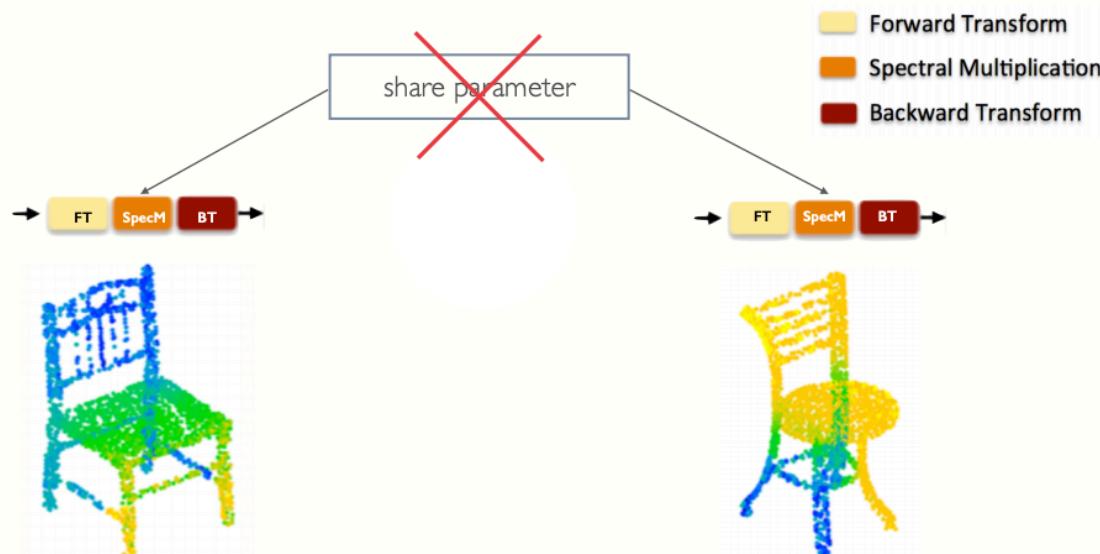
- Derived by eigenfunctions of self-adjoint operators, e.g. Laplacian-Bertrami or Dirac operator

“Fourier basis” of the graph: V : Eigenvectors of Δ



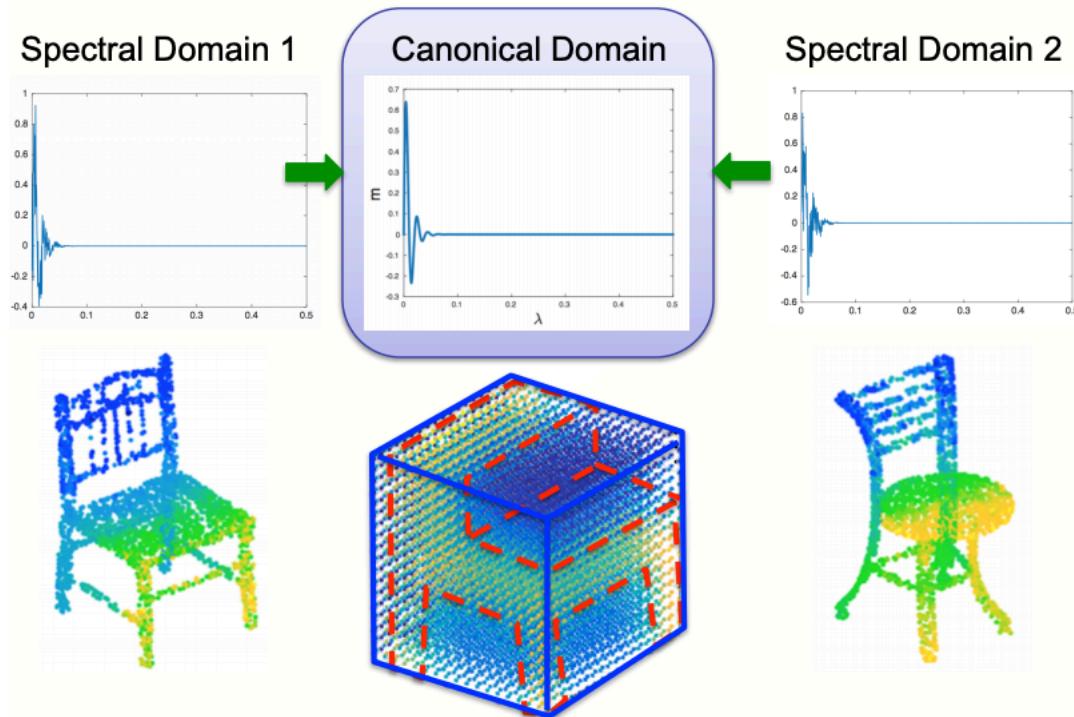
Fundamental Challenge of Spectral CNN

- 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



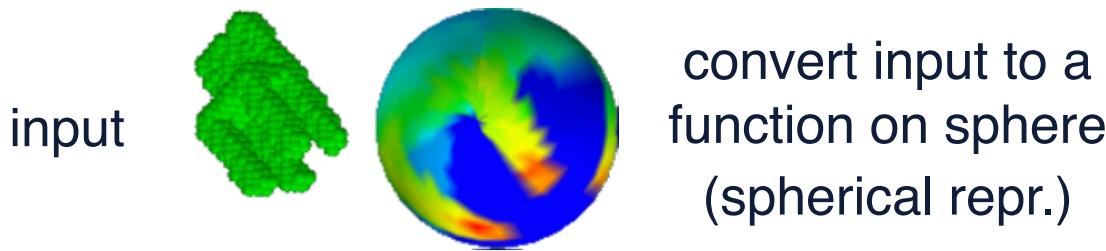
Domain Synchronization

- Rescue: Synchronize spectral domains by functional maps



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

Cohen et al., “Spherical CNN”, *ICLR 2018*

Esteves et al., “Learning $SO(3)$ Equivariant Representations with Spherical CNNs”, *ECCV 2018*

Learned Filters

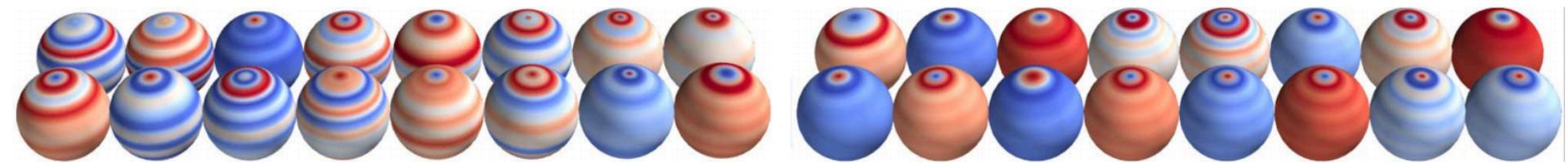


Fig. 4: Filters learned in the first layer. The filters are zonal. *Left*: 16 nonlocalized filters. *Right*: 16 localized filters. Nonlocalized filters are parameterized by all spectral coefficients (16, in the example). Even though locality is not enforced, some filters learn to respond locally. Localized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

A Special Case: Spherical CNN

- Rotation invariance guaranteed

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.

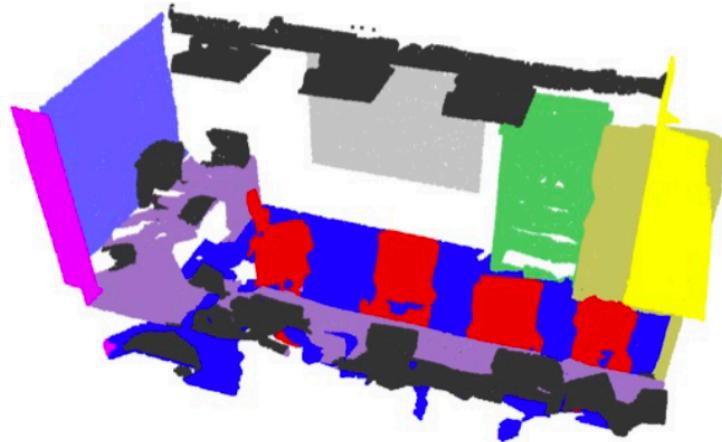
Method	z/z	SO3/SO3	z/SO3	params	inp. size
PointNet [7]	89.2	83.6	14.7	3.5M	2048 x 3
PointNet++ [38]	89.3	85.0	28.6	1.7M	1024 x 3
VoxNet [29]	83.0	73.0	-	0.9M	30^3
SubVolSup [8]	88.5	82.7	36.6	17M	30^3
SubVolSup MO [8]	89.5	85.0	45.5	17M	20×30^3
MVCNN 12x [9]	89.5	77.6	70.1	99M	12×224^2
MVCNN 80x [9]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

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

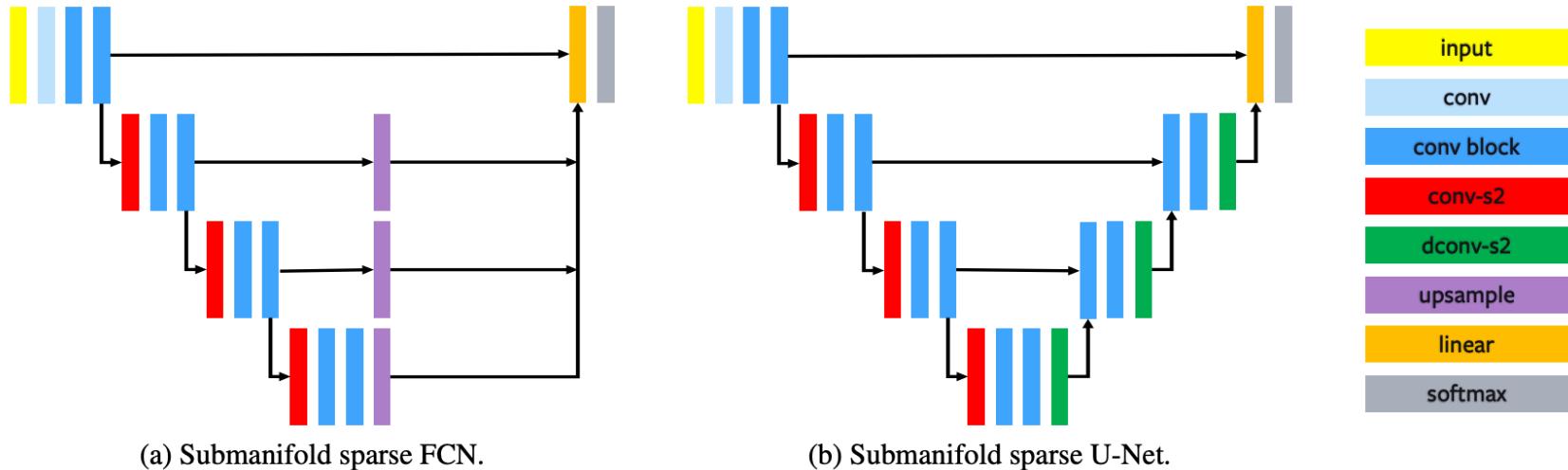
Topics

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

Task: 3D Semantic Segmentation



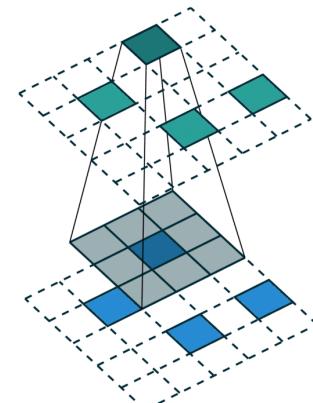
Encoder-Decoder



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

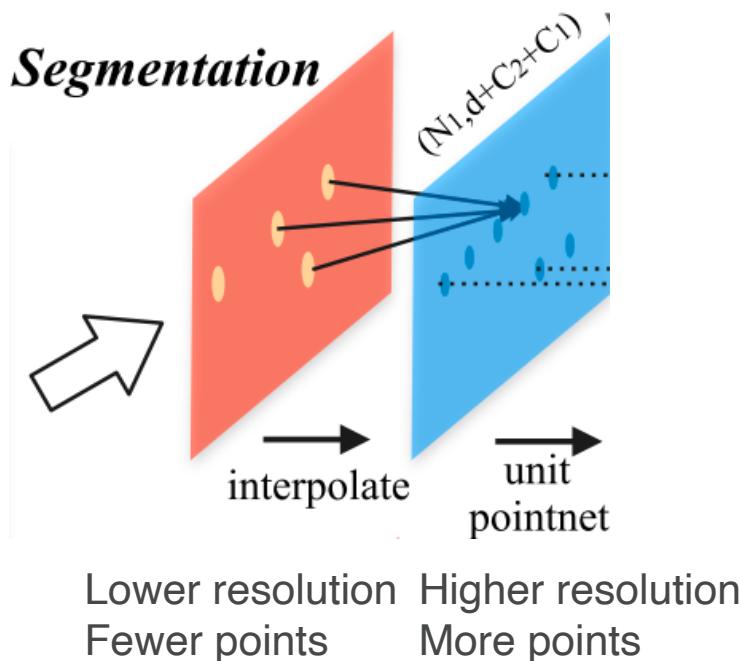
Choy, Christopher, et al. "4D Spatio-Temporal Convnets: Minkowski Convolutional Neural Networks." CVPR 2019

Sparse Convolution



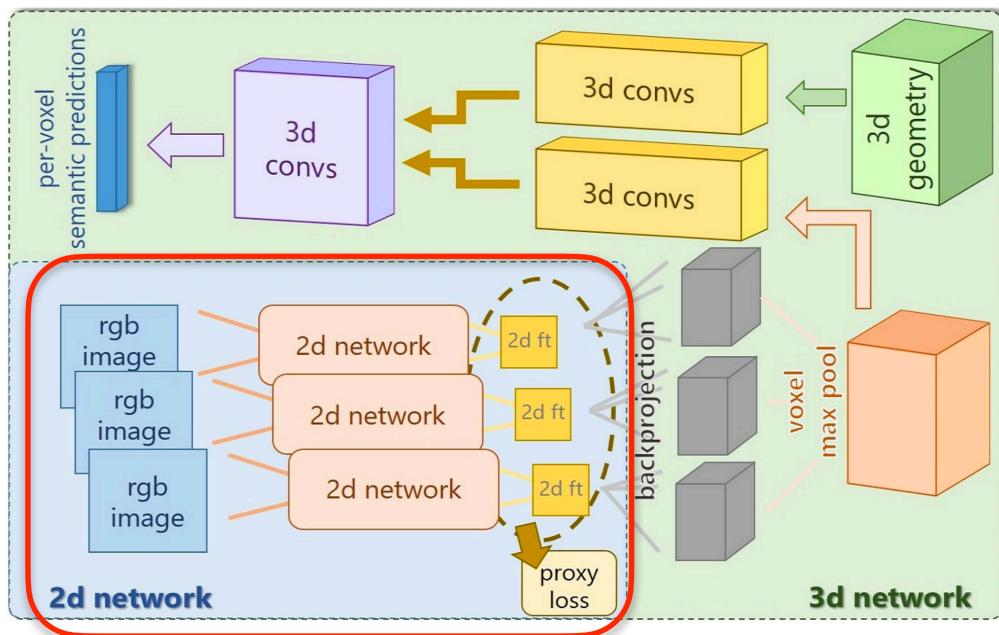
Sparse Deconvolution

Encoder-Decoder



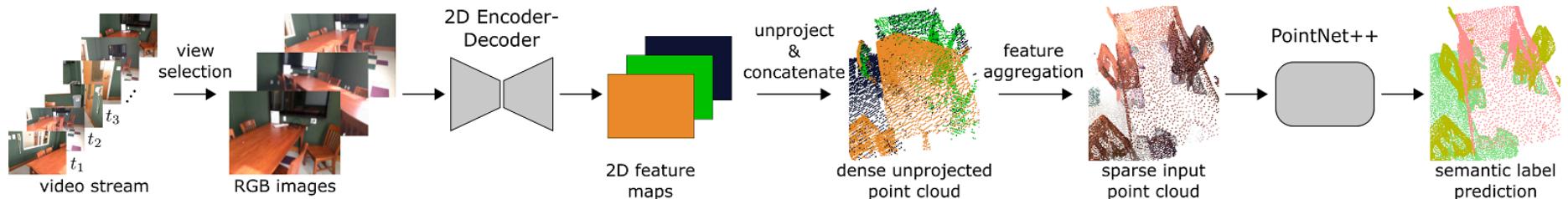
Upsample point cloud features by interpolating from 3 nearest points in the point cloud of lower resolution.

Multimodal



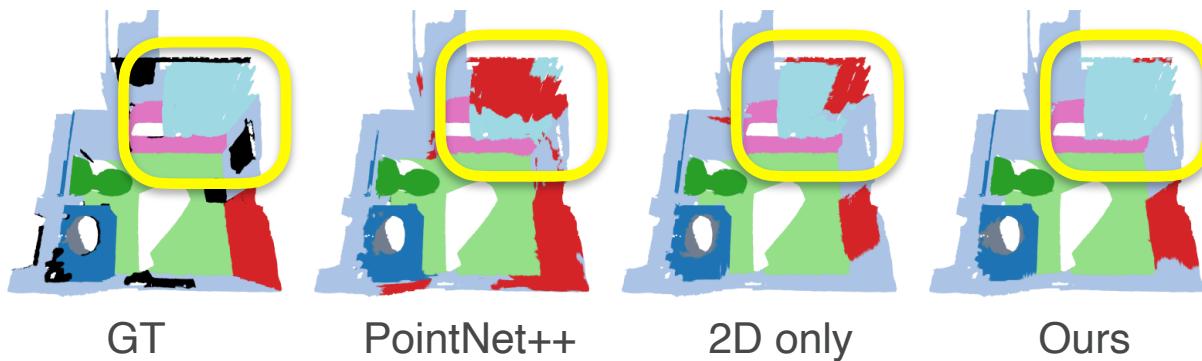
- Backproject 2D features to 3D voxels
- Apply voxel-wise max-pooling across multiple views
- Fuse 2D and 3D features at the intermediate level

Multimodal



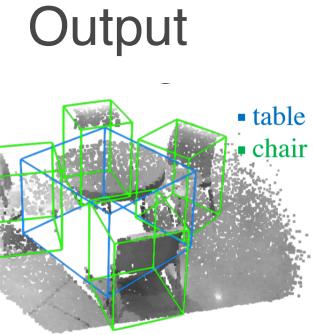
- Backproject 2D features to 3D points
- Apply PointNet to aggregate multi-view features
- Fuse 2D and 3D features as input to 3D network

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



Object Segmentation



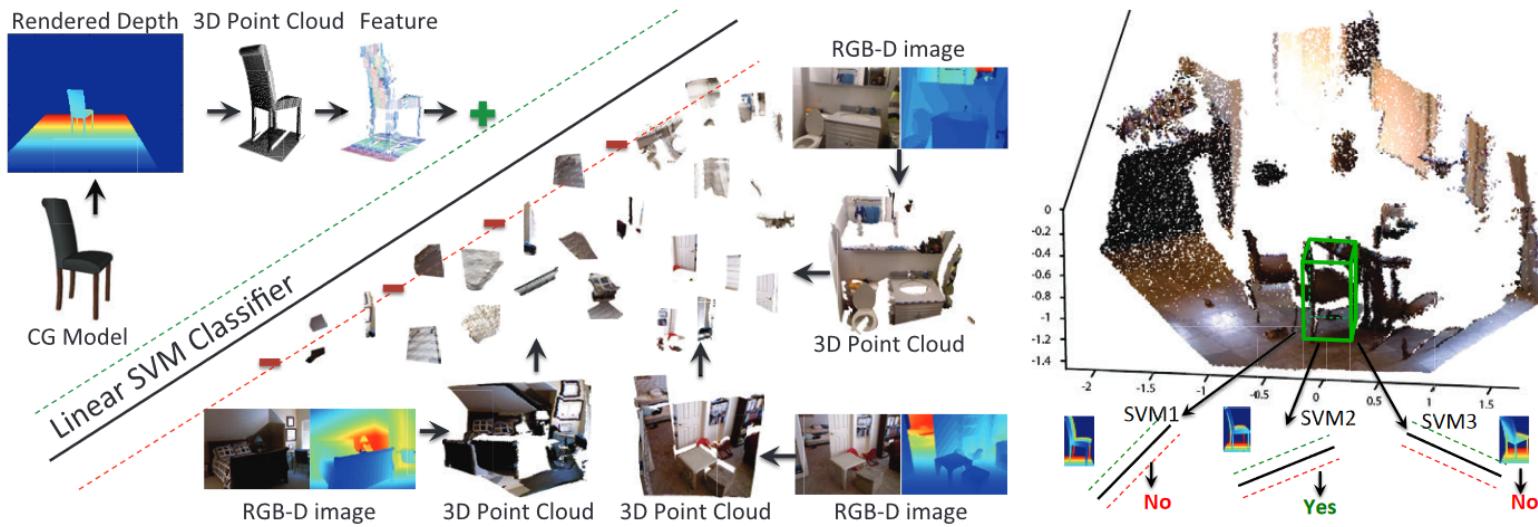
Part Segmentation



Top-down Methods

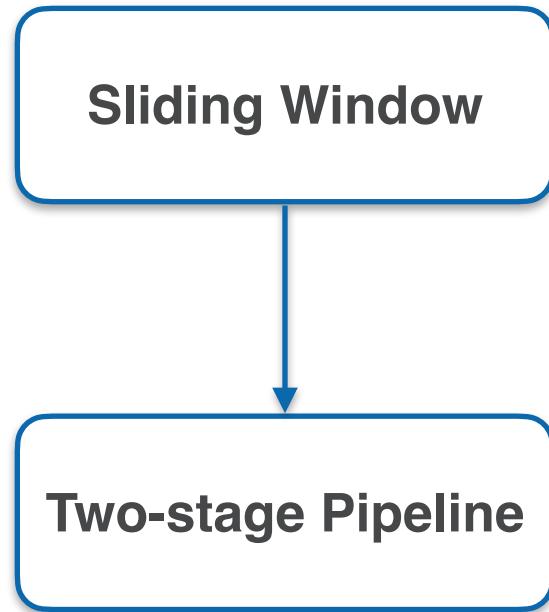
- Does this point cloud contain an object?
- How to select point clouds to classify objectness?

Sliding Shapes



Sliding window to walk over
the entire space

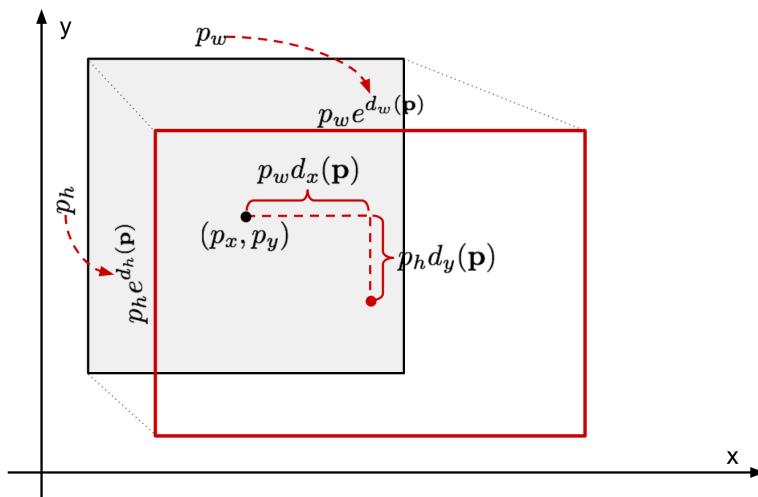
Expensive !



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

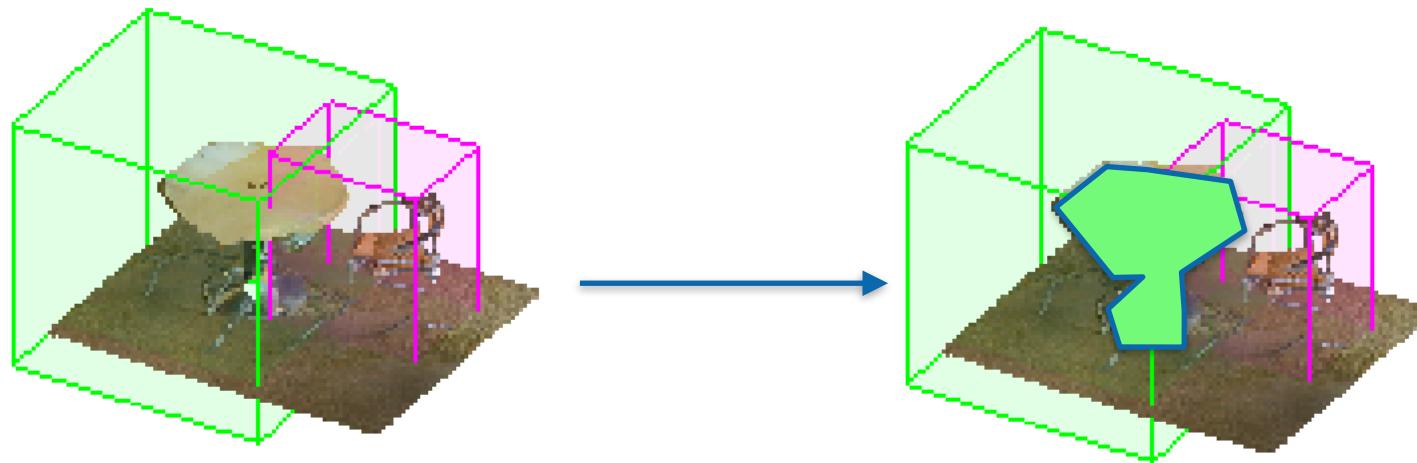
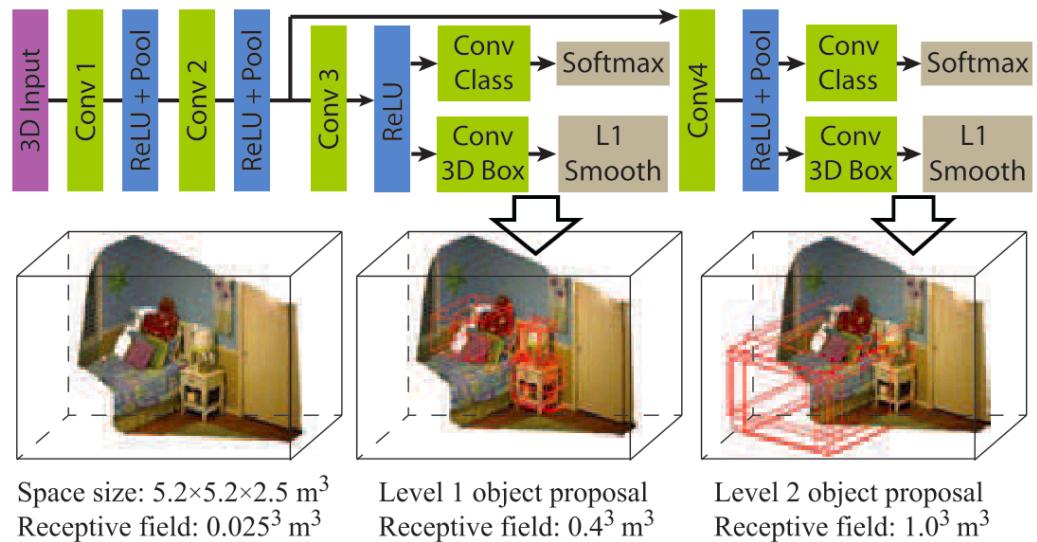


Figure from “Learning Object Bounding Boxes for 3D Instance Segmentation on Point Clouds”

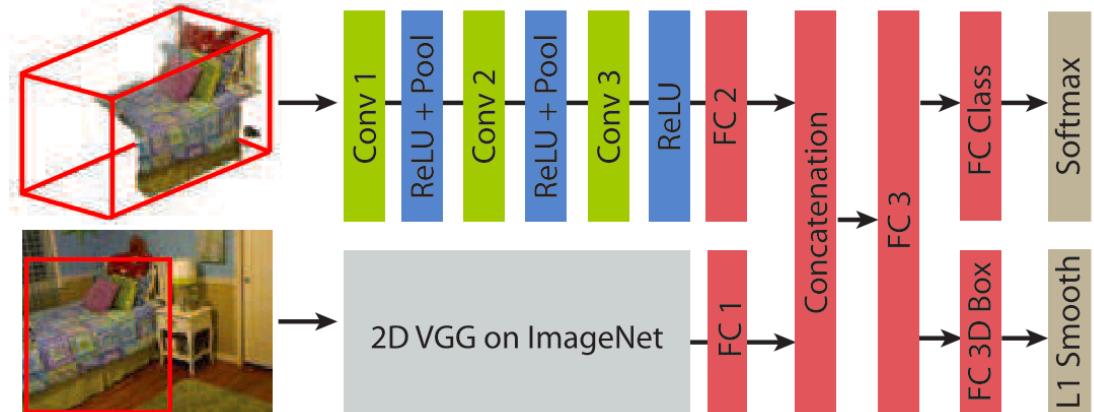
Volumetric R-CNN

Stage1: 3D Region
Proposal Network

TSDF



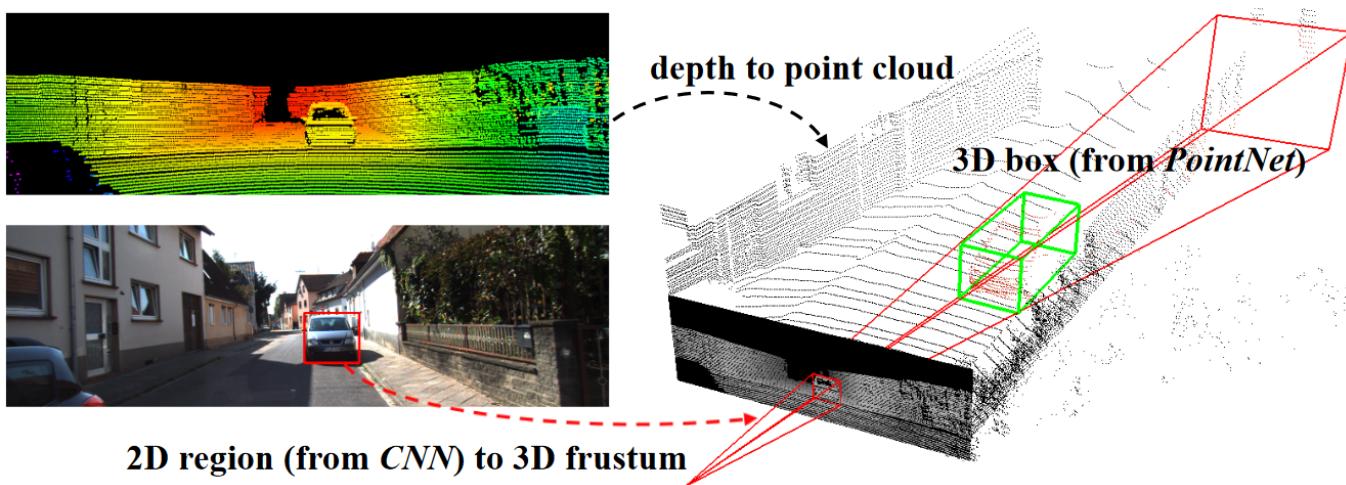
Stage2: Joint Object
Recognition Network



Song et al., “Deep Sliding Shapes for Amodal 3D
Object Detection in RGB-D Images”, CVPR 2016

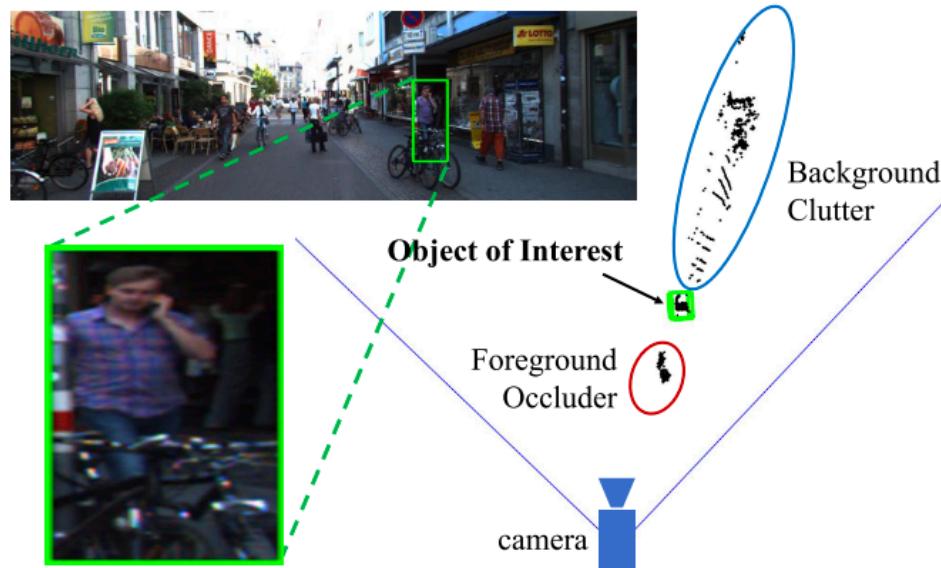
View-based Proposal

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



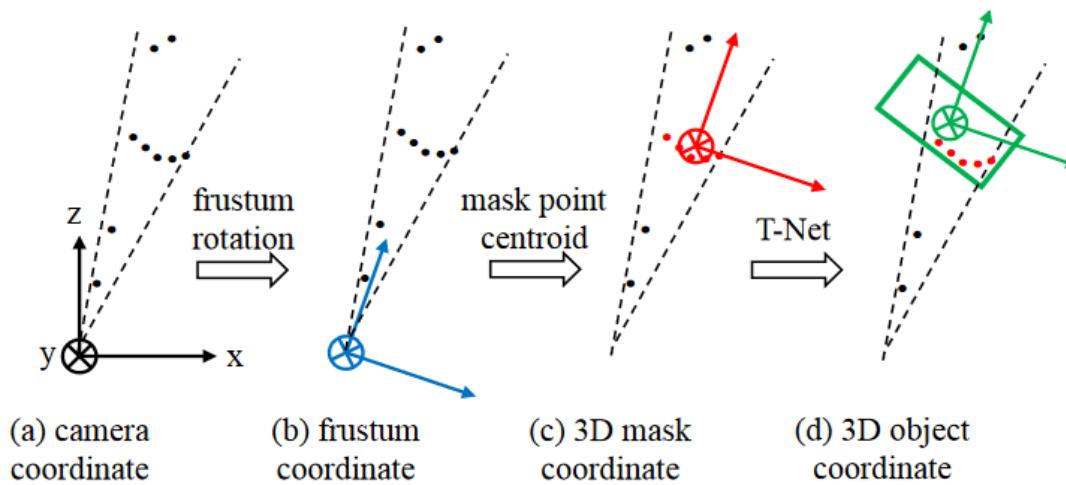
Stage 2: FG/BG Segmentation

Instance segmentation to remove other foreground instances and background clutter



Stage 2: Coordinate Normalization

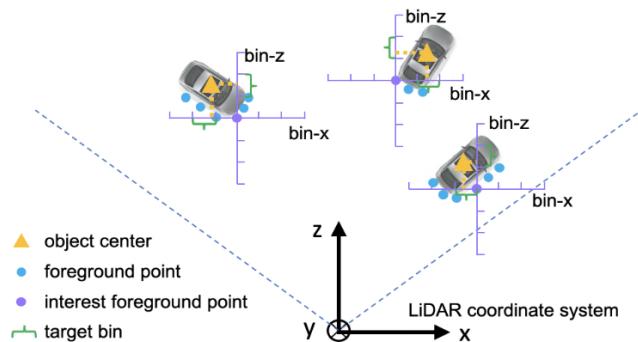
Handle perspective variation in frustum point cloud by a series of coordinates normalization



Point-based Proposal

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

Stage-2: Refine proposals in the canonical coordinates



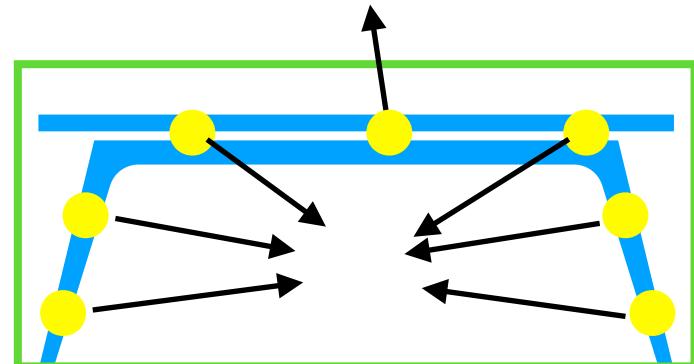
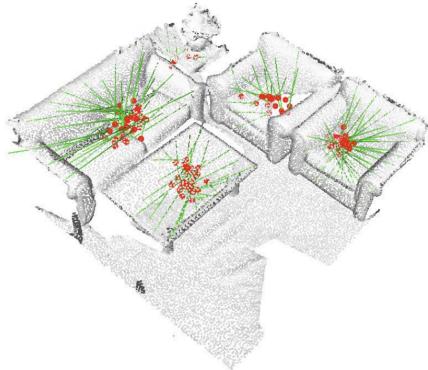
Bin-based Box Representation

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

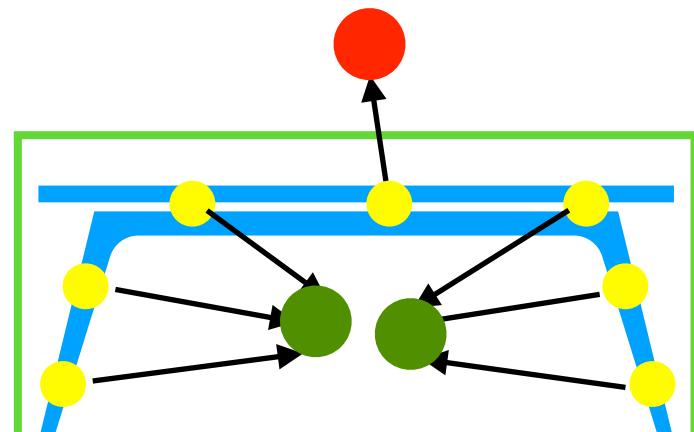
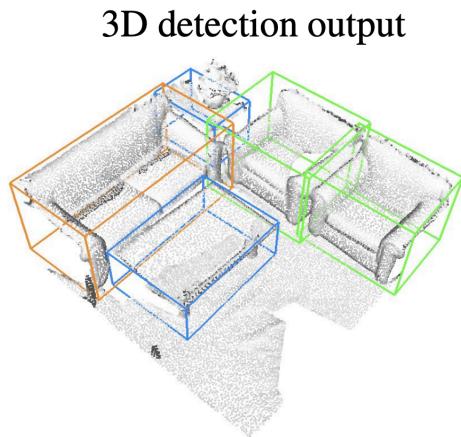
Voting from input point cloud



Proposal from Voting

Challenge: 3D object centroid can be far from any surface point, thus hard to regress accurately

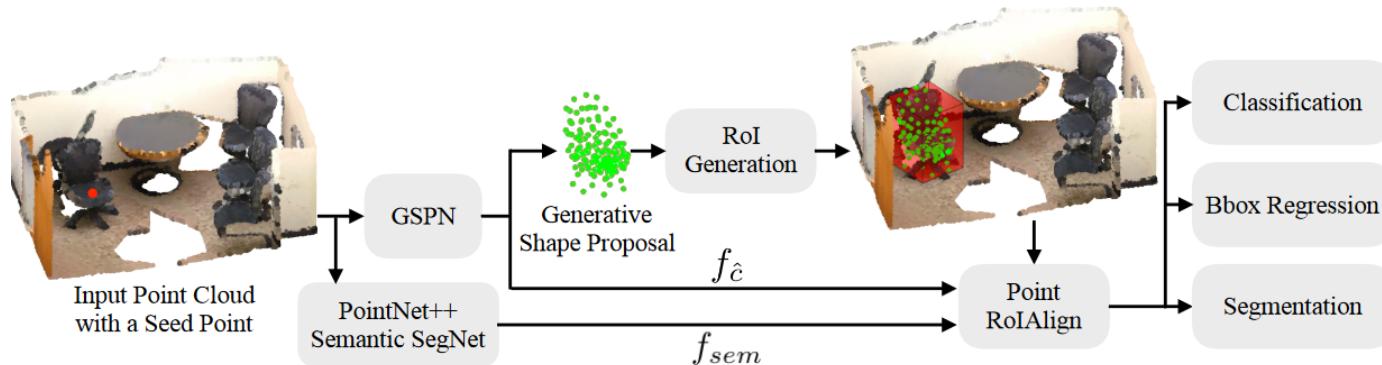
- Votes cluster near object centers
(Farthest point sample votes to cluster)



● Positive
● Negative

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



Bottom-up Methods

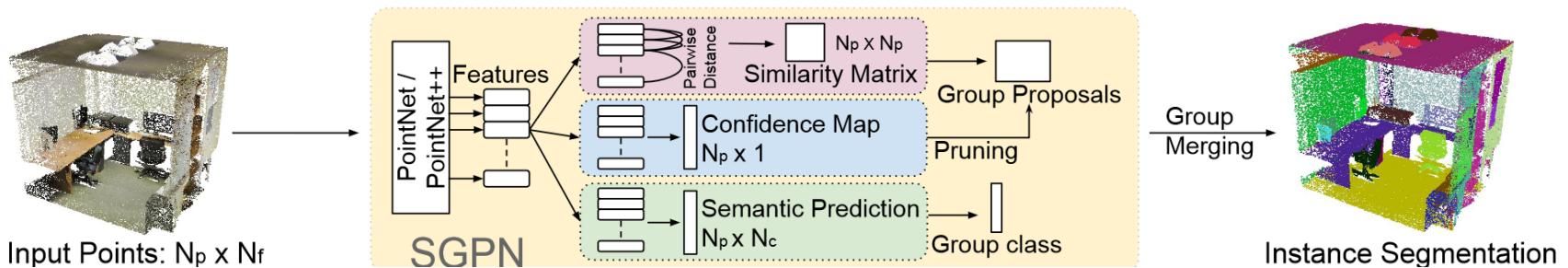
- Do these points belong to the same instance?
- How to measure the similarity between two points?

Associative Embedding

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

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

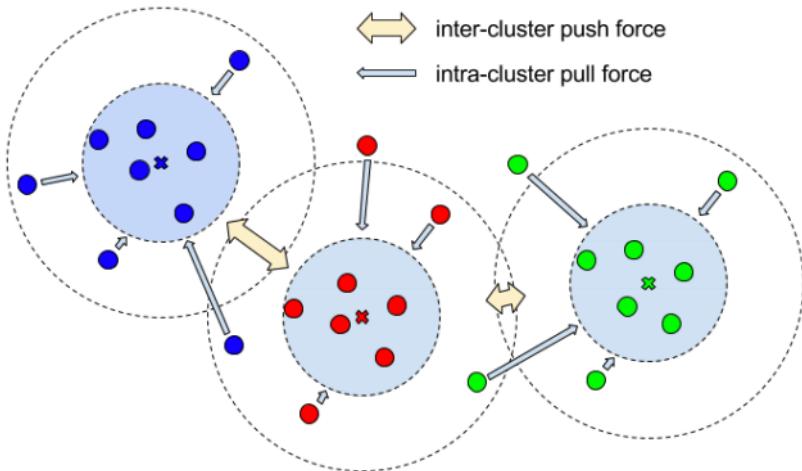
- Clustering gives proposals



JSIS3D

- A discriminative function to present the embedding loss

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

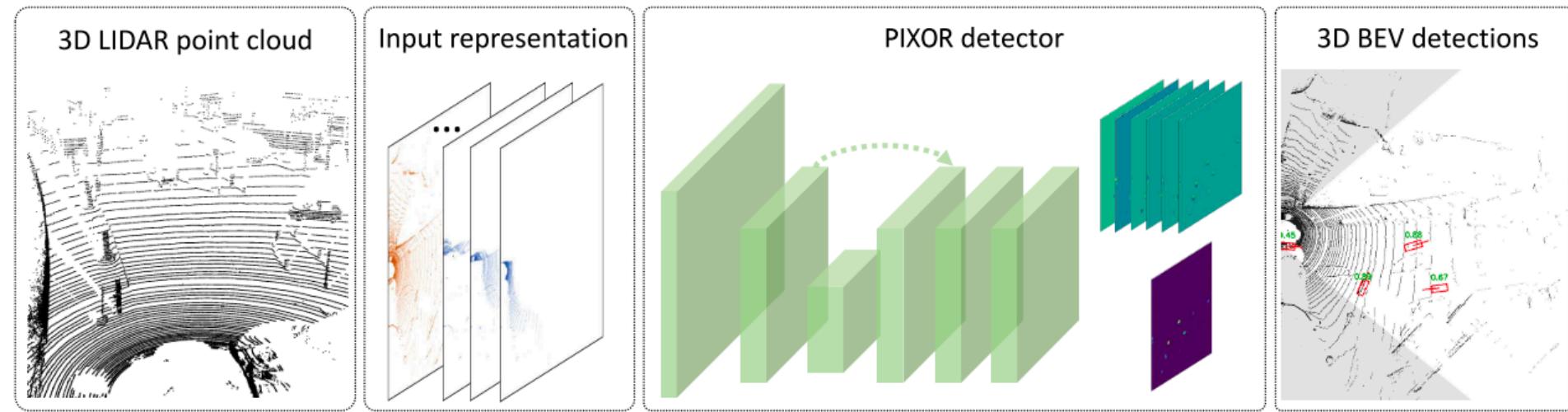


$$\mathcal{L}_{pull} = \frac{1}{K} \sum_{k=1}^K \frac{1}{N_k} \sum_{j=1}^{N_k} [\|\boldsymbol{\mu}_k - \mathbf{e}_j\|_2 - \delta_v]_+^2$$

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

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

BEV (bird-eye view)

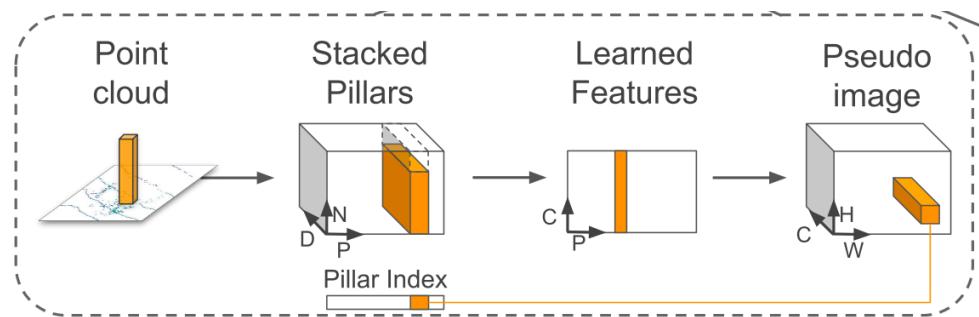


$H \times W \times D$ voxels

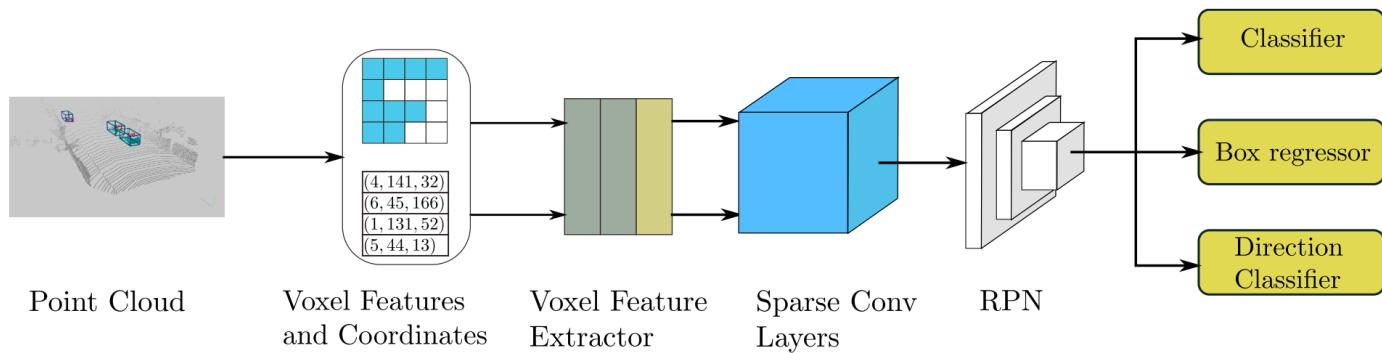
$H \times W$ image
with D channels

Feature Encoder

- Replace occupancy of each pixel with point cloud features of the pillar



- Sparse convolution to extract features

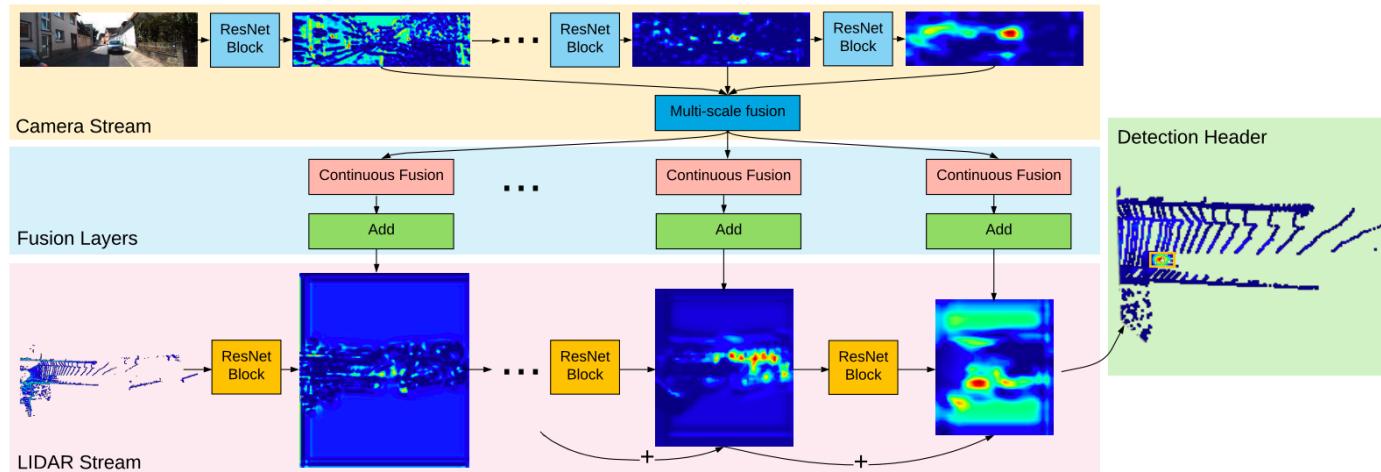


Yan, Yan, et al. “Second: Sparsely Embedded Convolutional Detection.” Sensors 2018

Lang, Alex H., et al. “Pointpillars: Fast Encoders for Object Detection from Point Clouds.” CVPR 2019

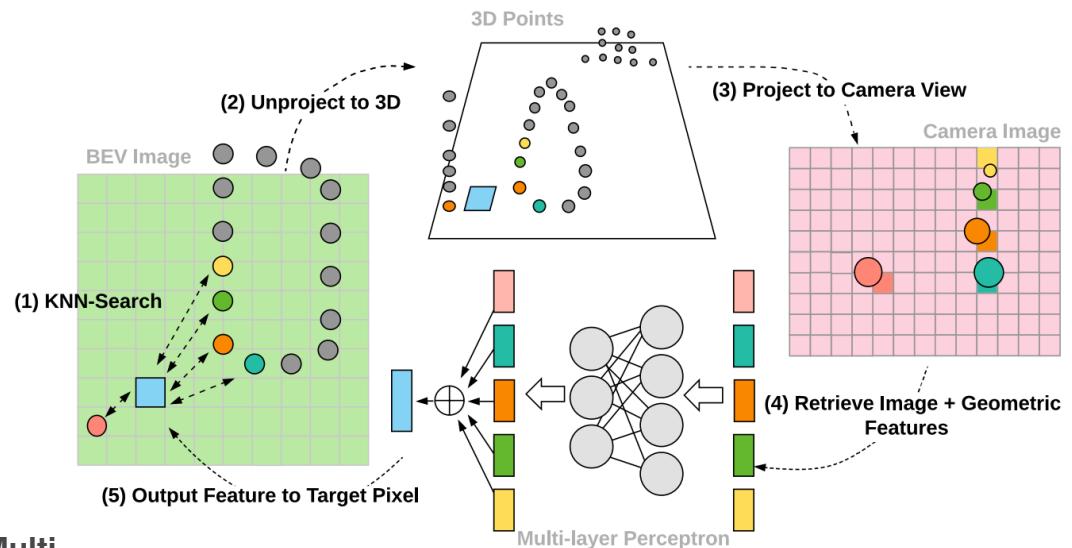
ContFuse

Image feature network: 2D image



BEV network: 3D voxel (HxWxC)

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



Topics

- 3D Data
- Classification
- Segmentation and Detection
 - Few-shot/Zero-shot Learning
- Reconstruction

Why Few-shot/Zero-shot Learning by 3D?

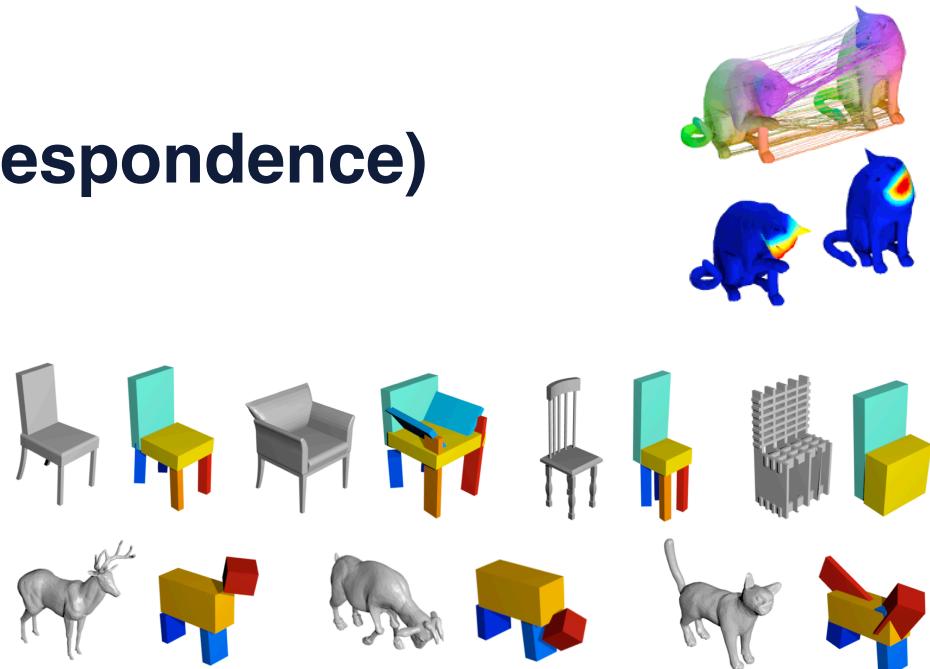
- Can be a better platform than images
- Shapes are **pure**, with no contamination by projection distortion, illumination, ...

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

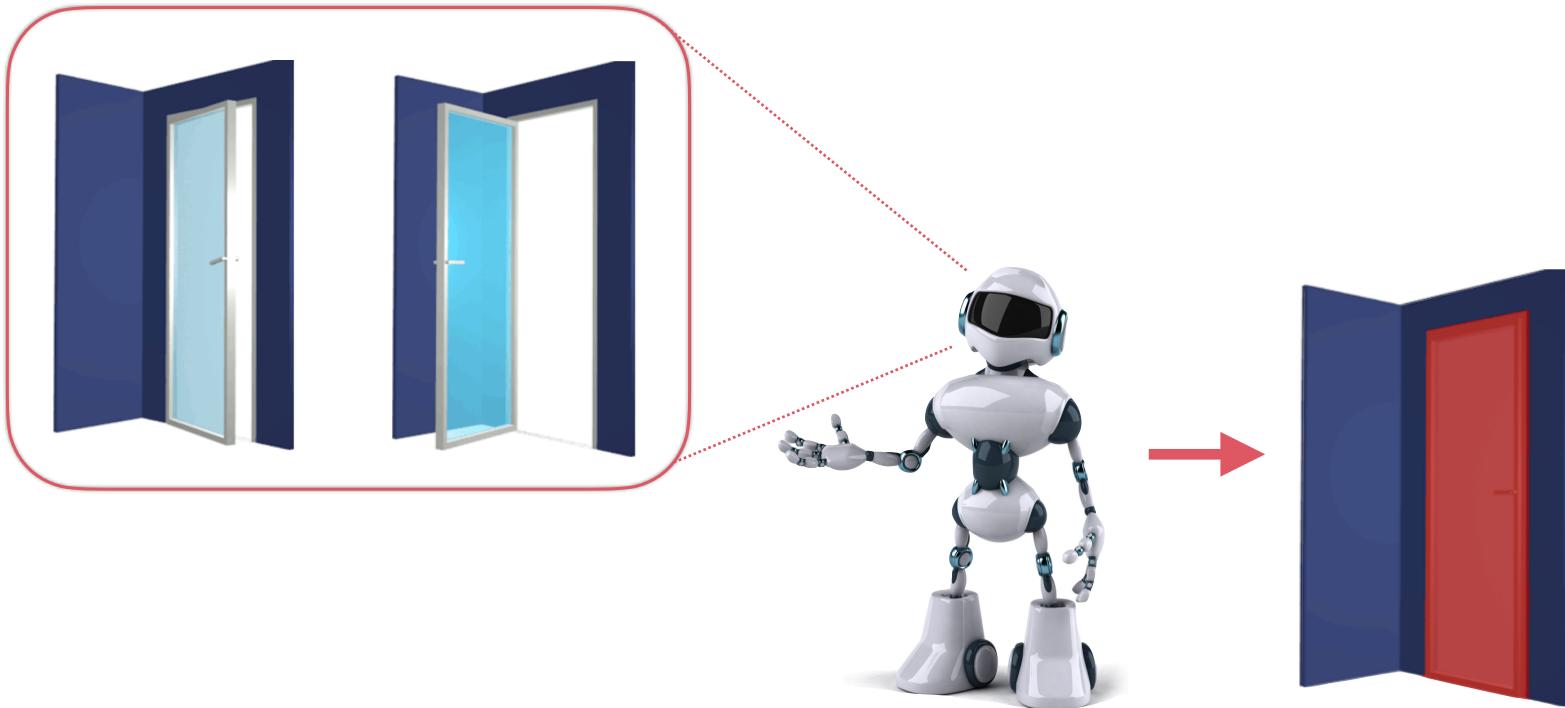
Why Few-shot/Zero-shot Learning by 3D?

Algorithmically, 3D shapes are:

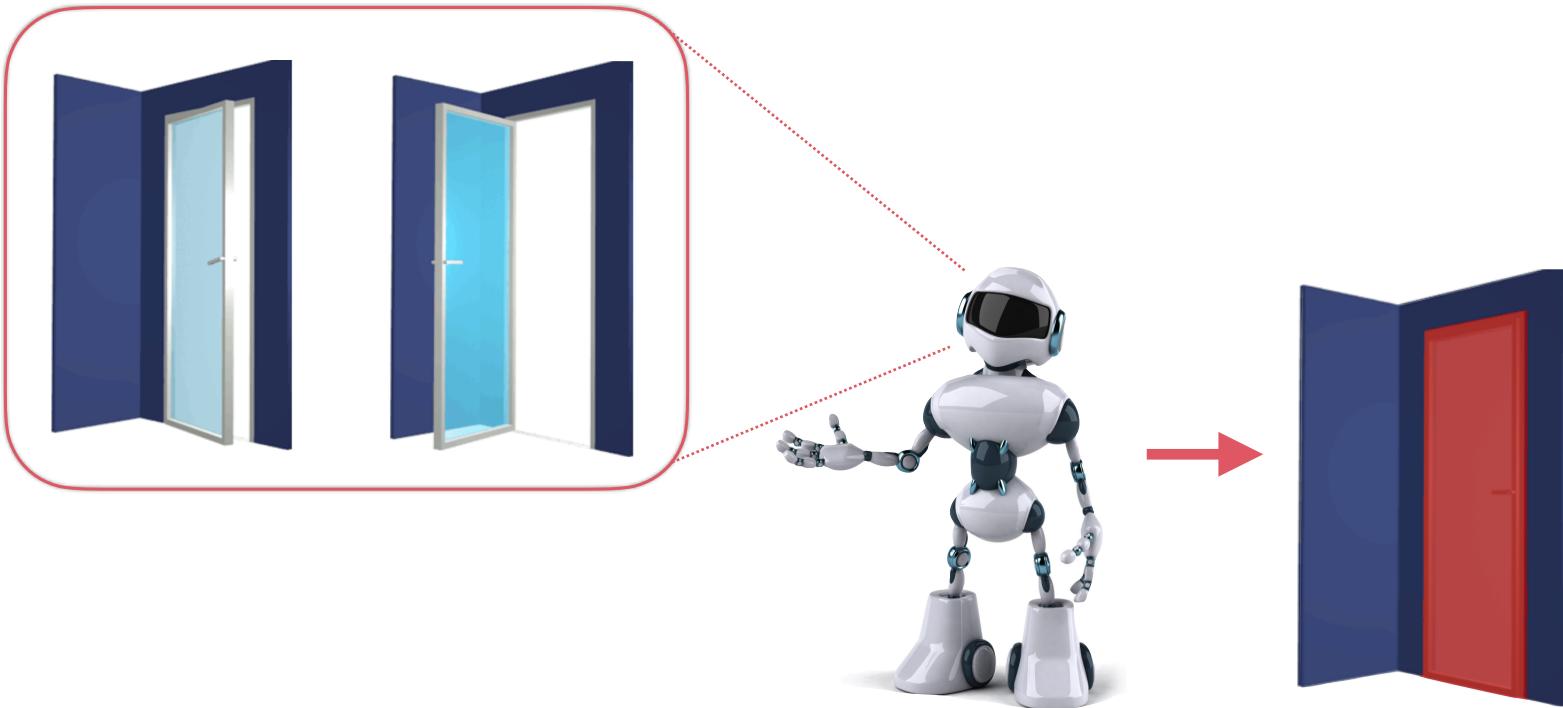
- easier to be **related (correspondence)**
- easier to be **compared**
- easier to **abstracted**



Task: Few-shot Structure Induction

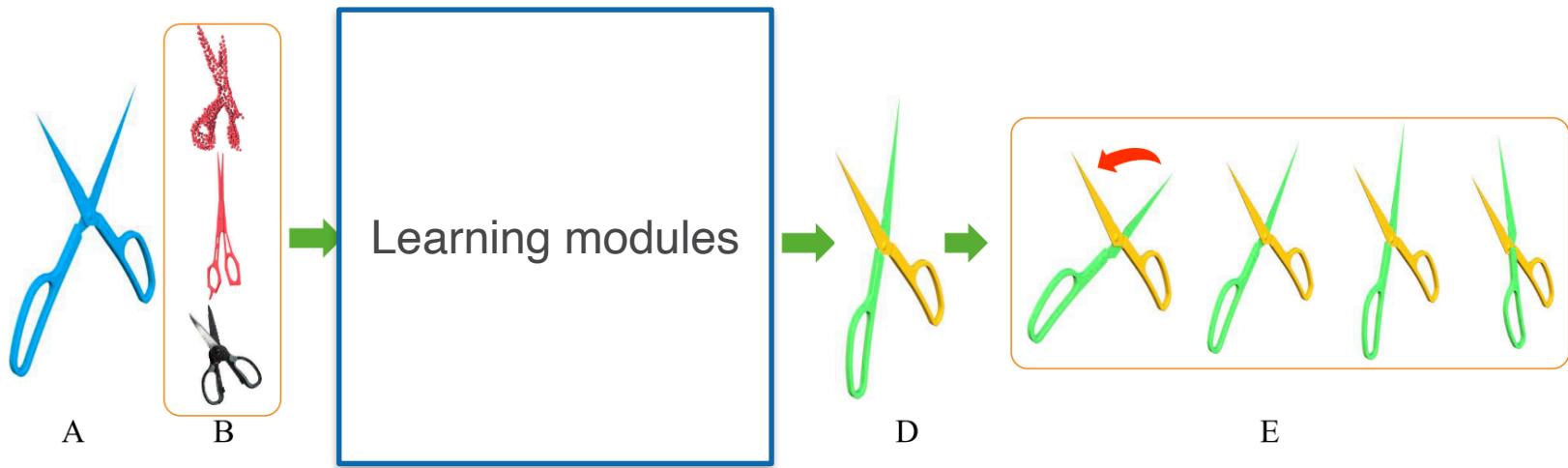


Emergence of Structure by Persistence

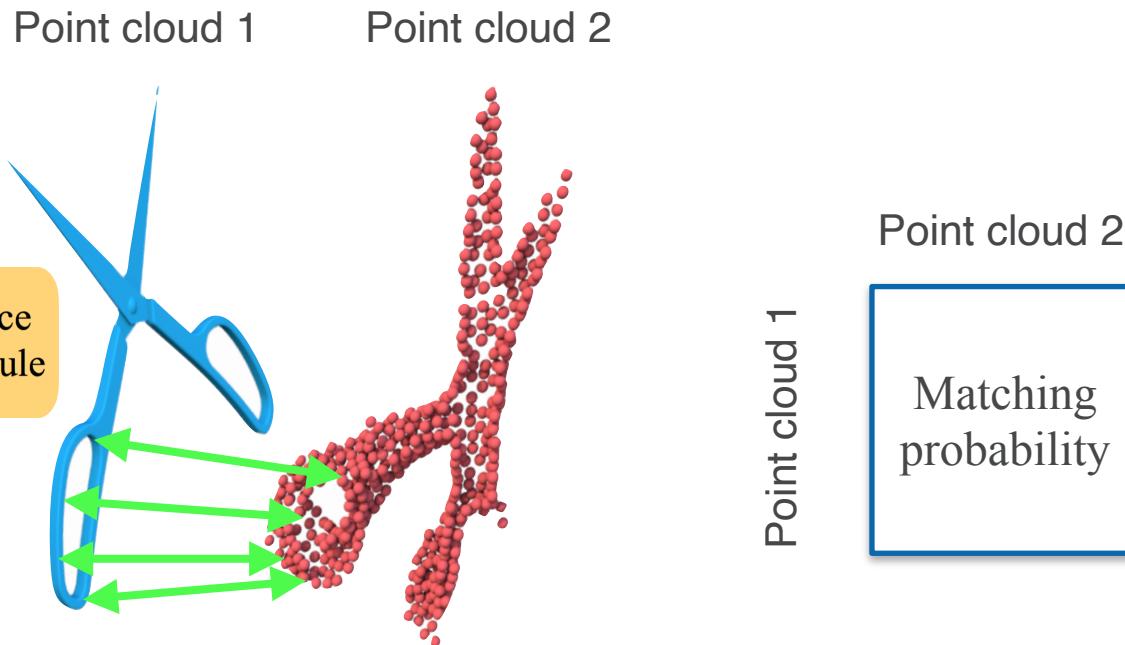


Capture re-occurring units!

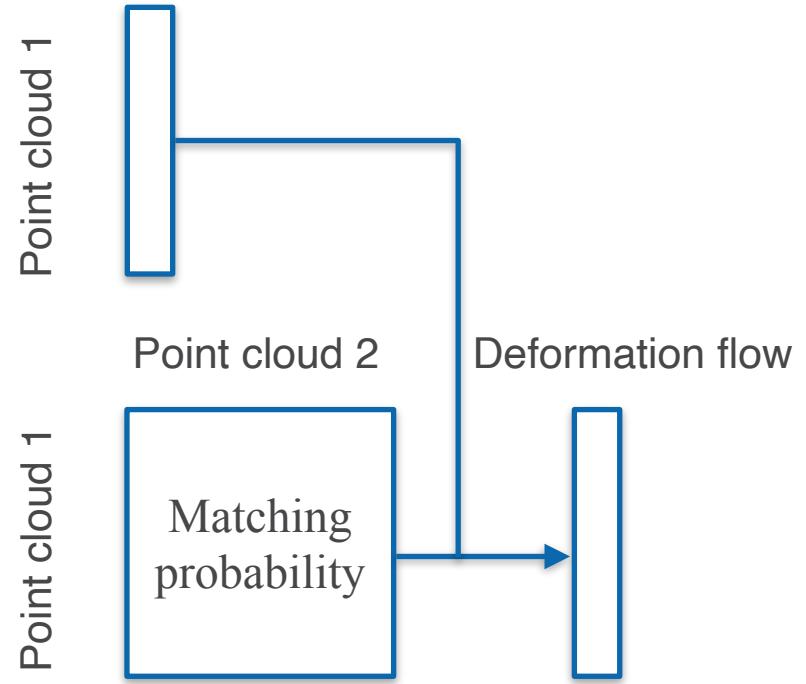
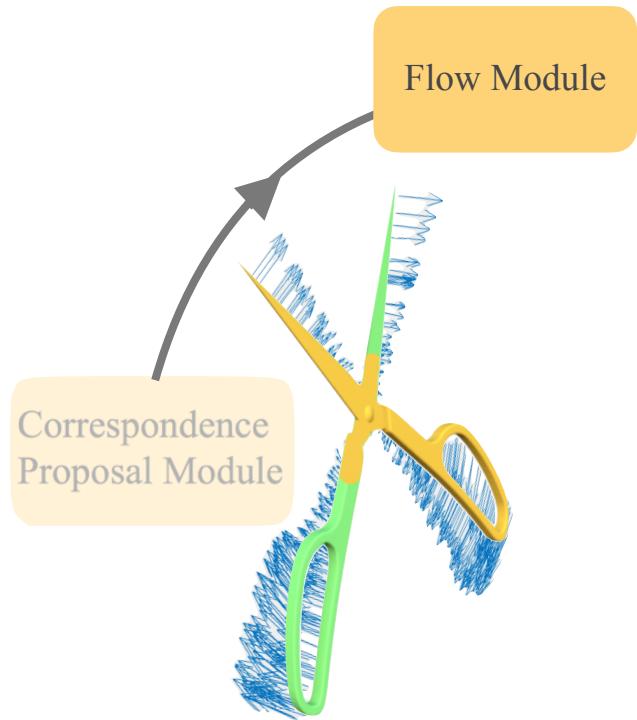
Part Induction by Relating Shapes



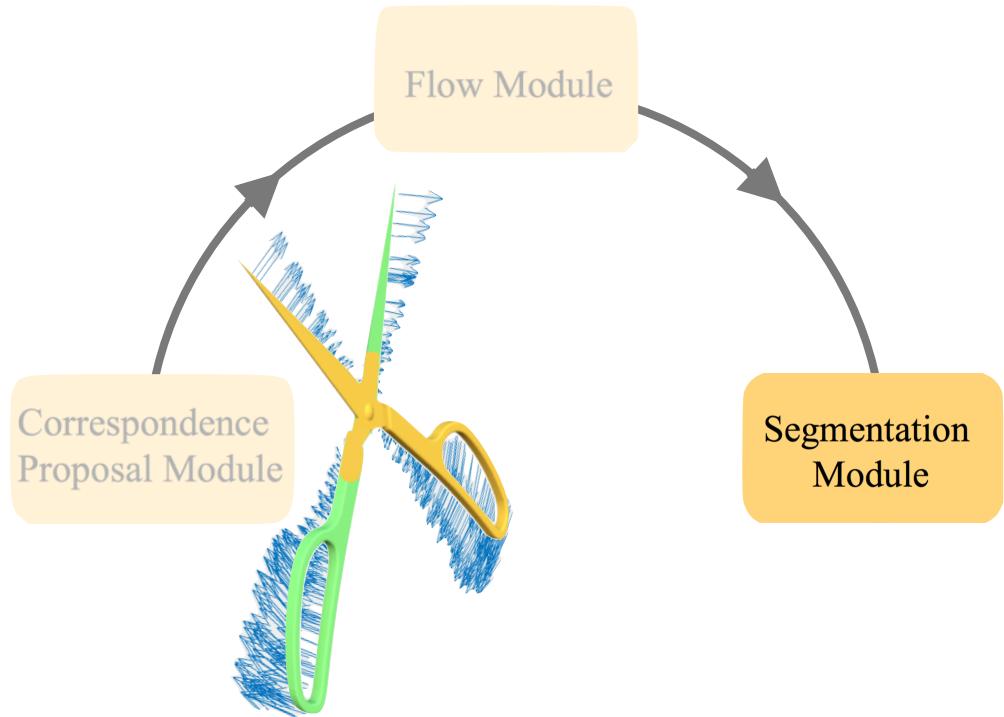
Part Induction by Relating Shapes



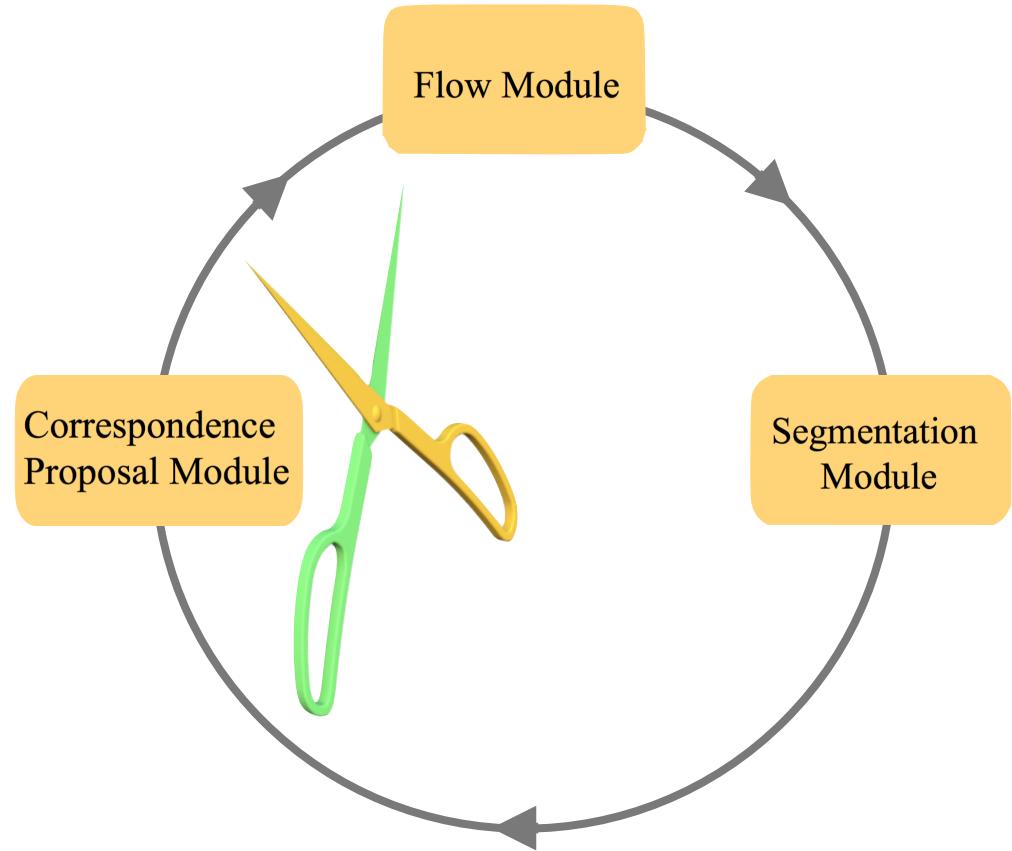
Part Induction by Relating Shapes



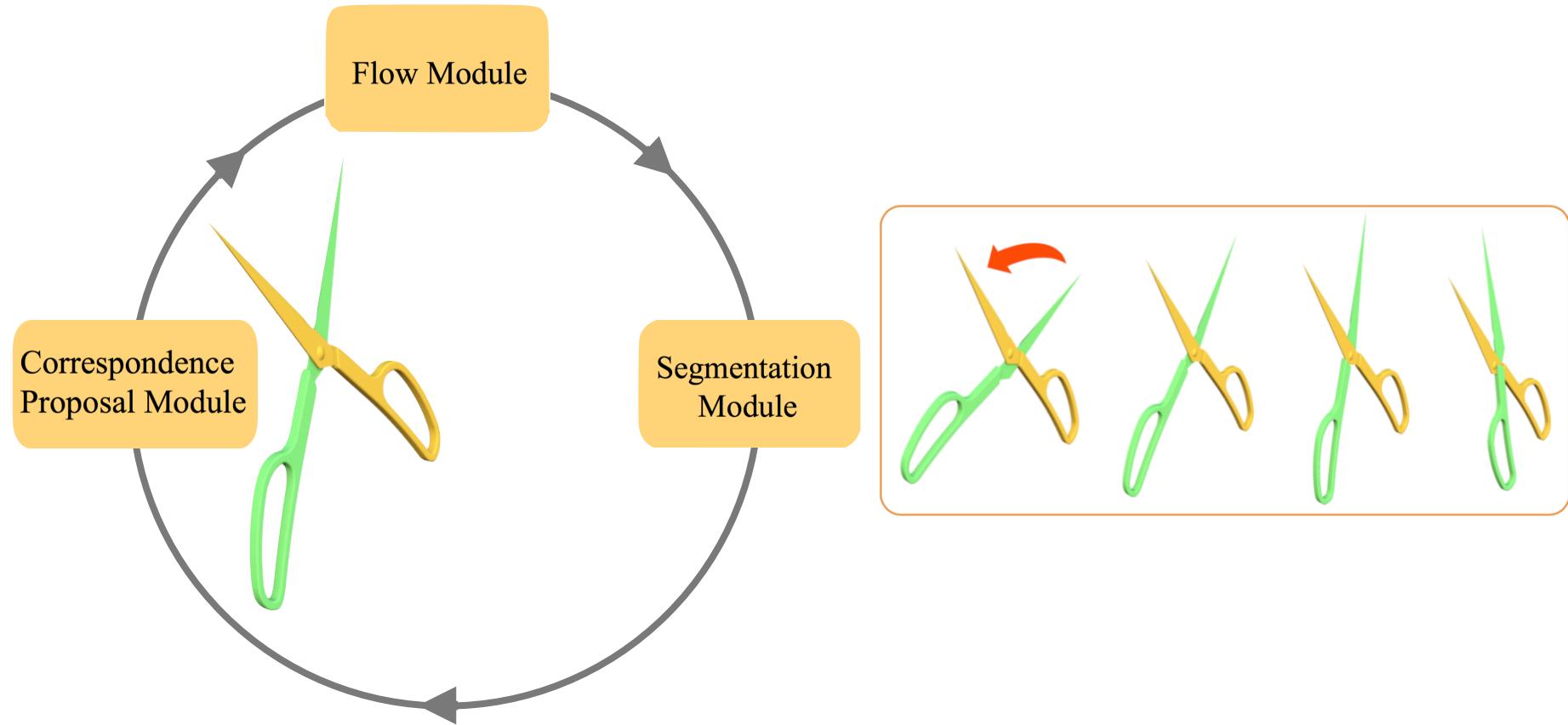
Part Induction by Relating Shapes



Part Induction by Relating Shapes



Mobility Induction

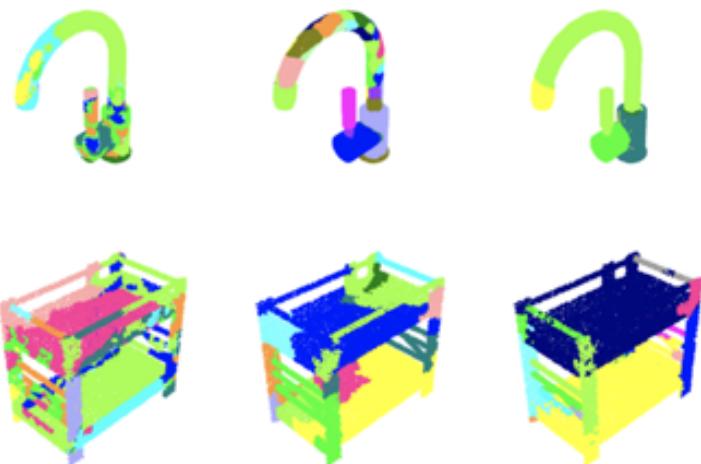


Task: Zero-shot Part Discovery

Train set



Test set

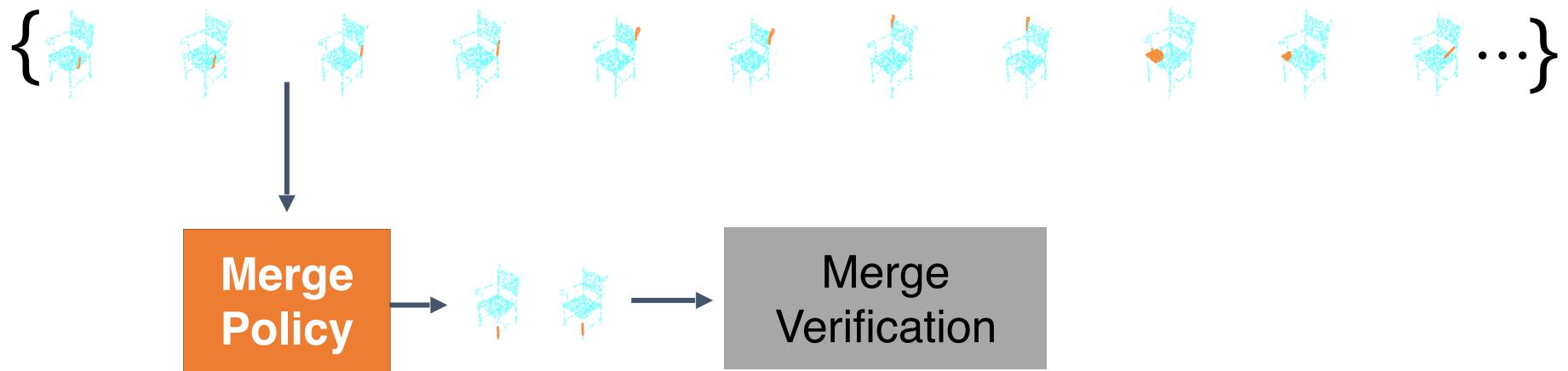


SOTA of
Deep Learning Classics
(PartNet) SOTA of
(WCSeg)

Ours

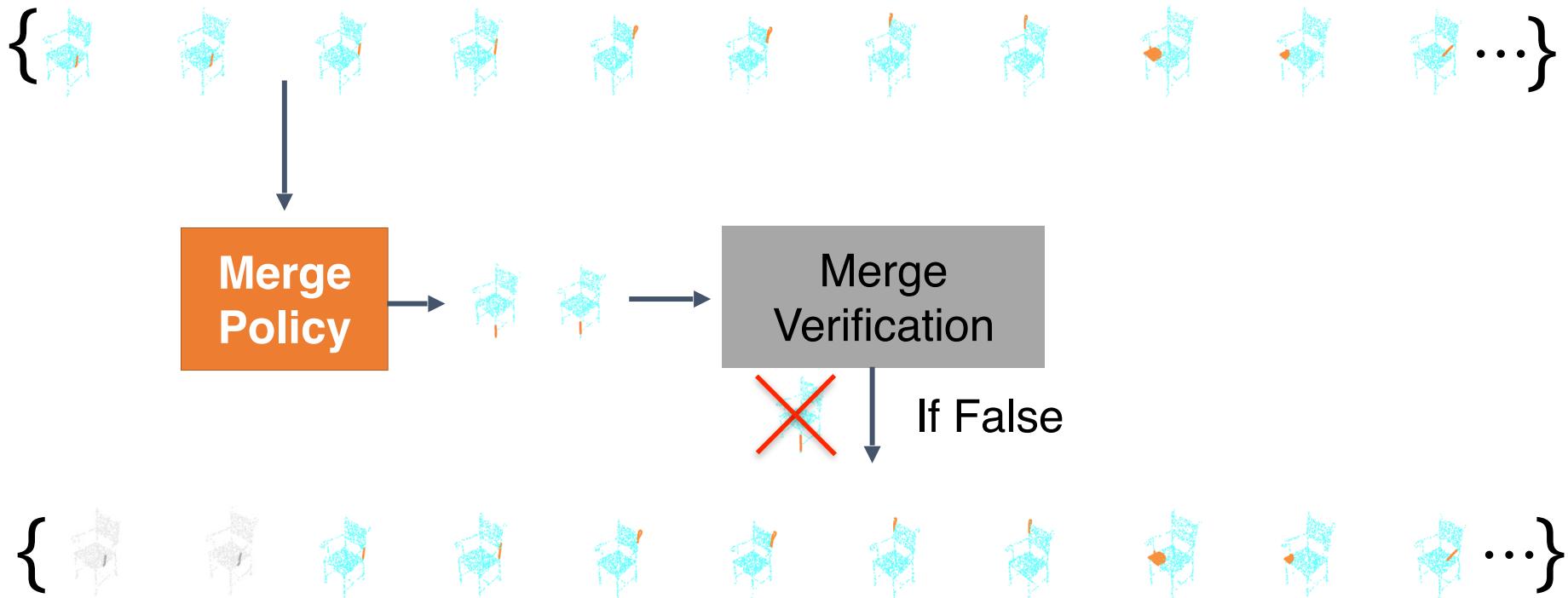
Learning to Group

Sub-Part Pool



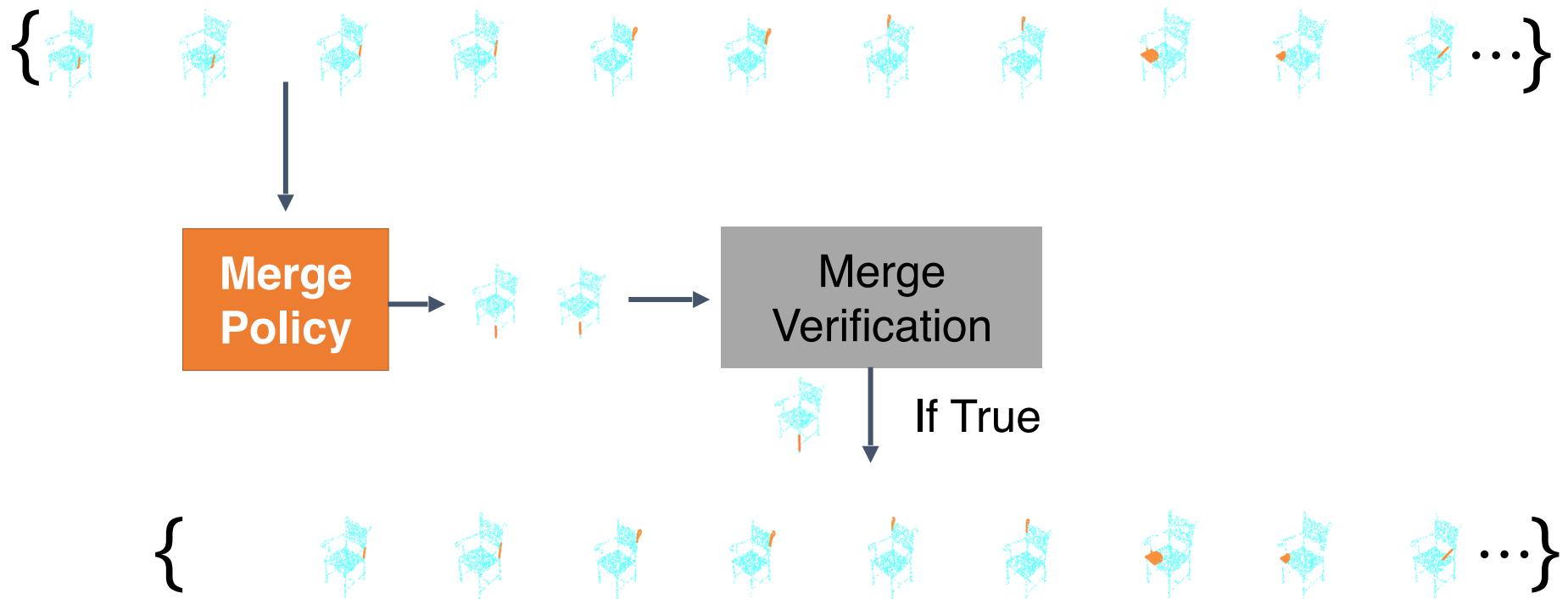
Learning to Group

Sub-Part Pool



Learning to Group

Sub-Part Pool



Learning to Group

Sub-Part Pool



Topics

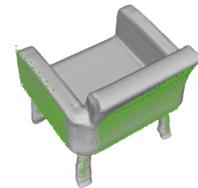
- 3D Data
- Classification
- Segmentation and Detection
- Reconstruction
 - Generative Model
 - Multi-View Stereo

Task

Conditional generation

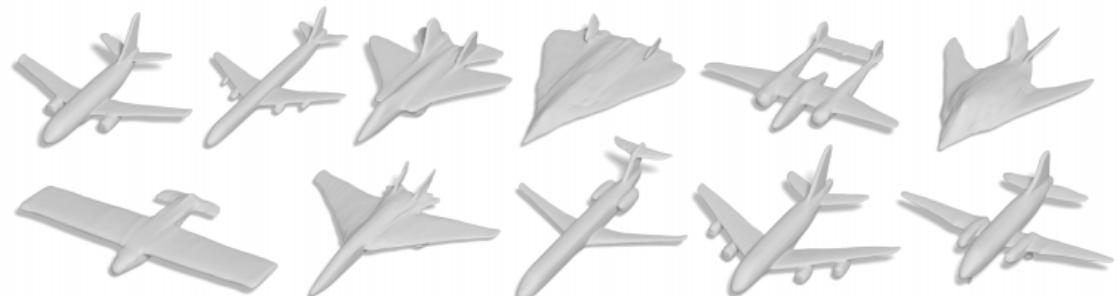


Single-image
3D reconstruction



Shape Completion

Free generation



Gaussian Noise

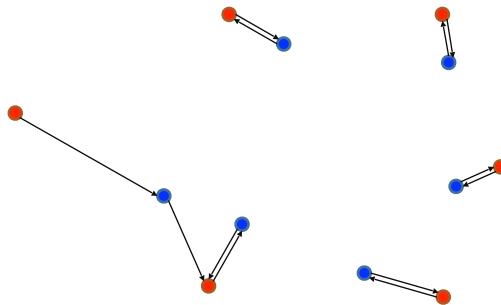
Metric

First of all,

how to evaluate the generated shapes?

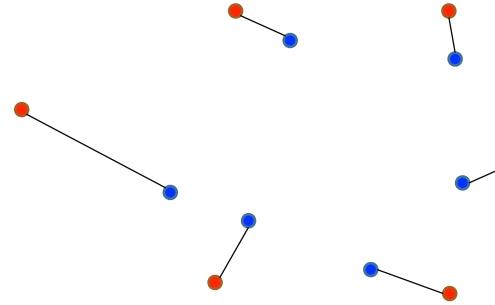
Metric for Point Clouds

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.

F-score

precision and recall are calculated by checking the percentage of points in one point cloud that can find a neighbor from the other point cloud within a threshold. F-score is then calculated as the harmonic mean.

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

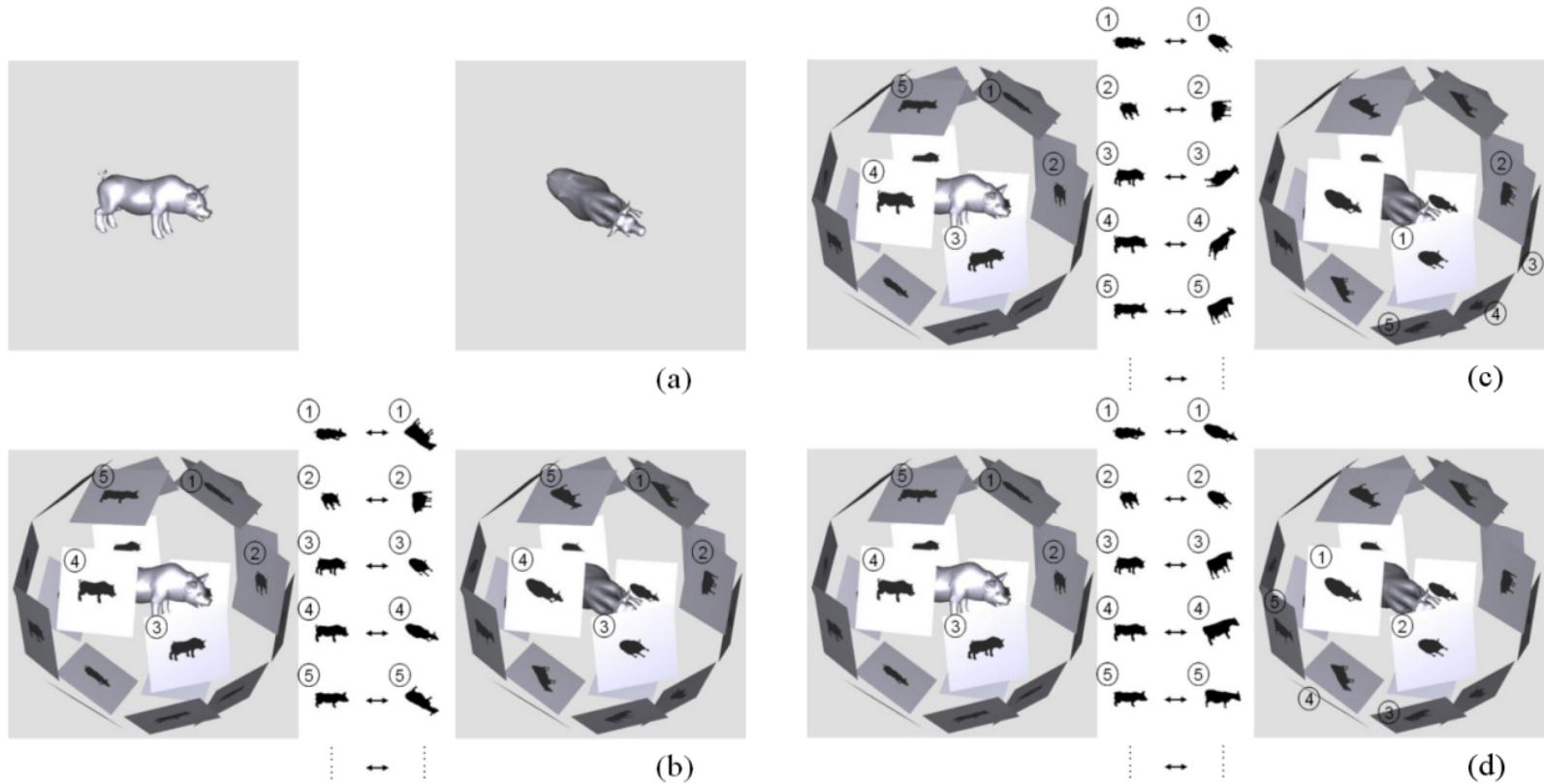
Wang et al., "Pixel2mesh: Generating 3d mesh models from single rgb images", ECCV 2018

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

Metric by Projection

Light Field Descriptor (LFD)

- Extract features from orthogonal projections

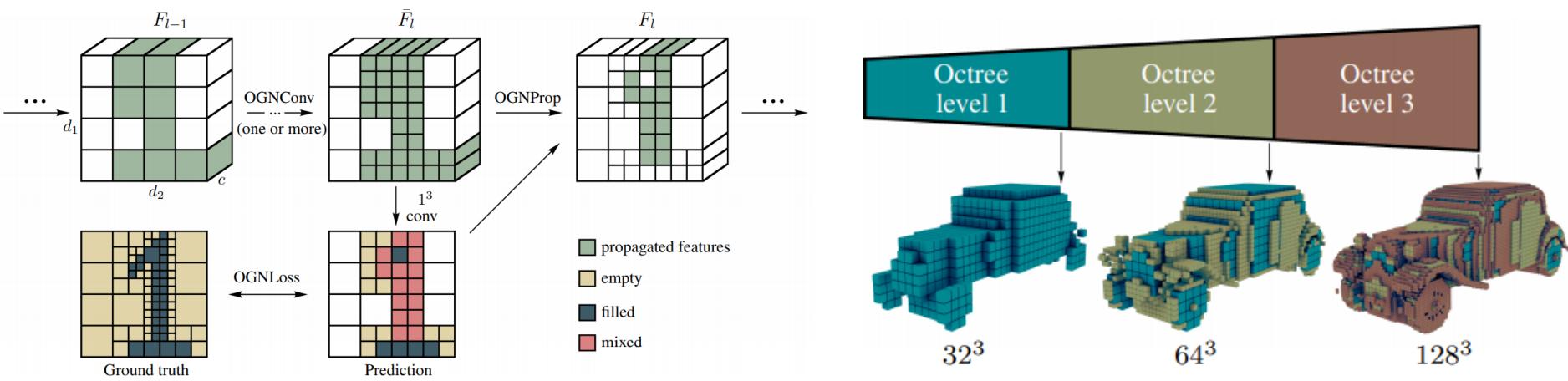


Algorithm for Conditional Generation

From Single Image to Volume

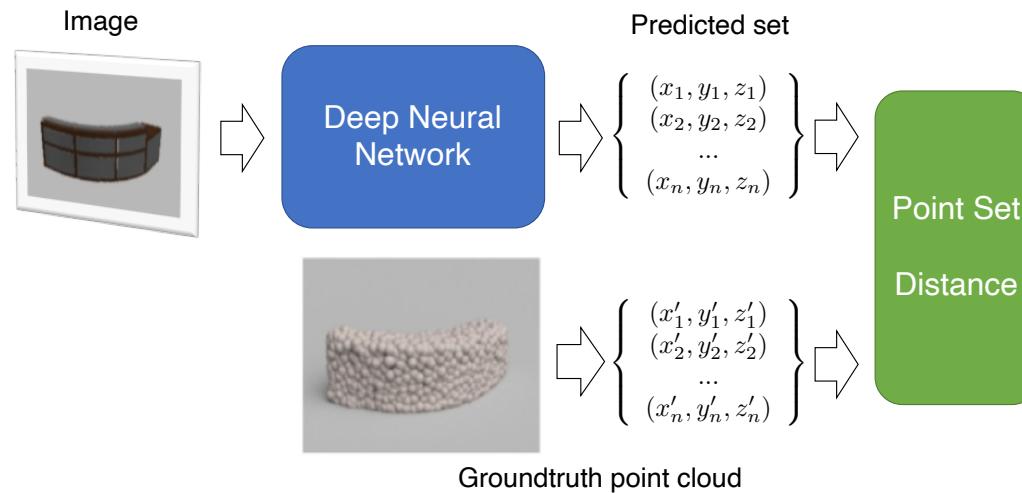
Avoid $\mathcal{O}(n^3)$ reconstruction

- Octree representation of shapes
- Generate the octree layer by layer



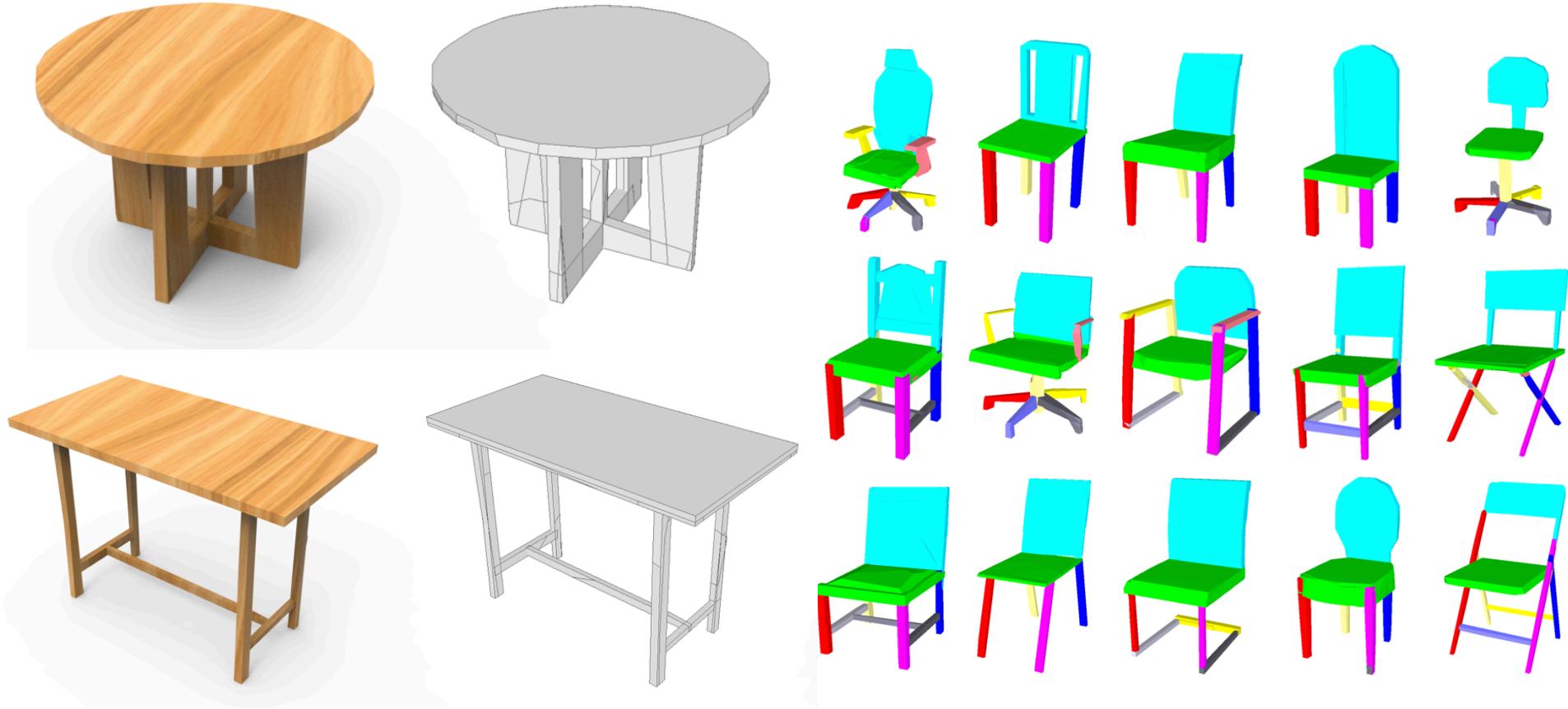
From Single Image to Point Cloud

- It is possible to generate a **set** (permutation invariant)



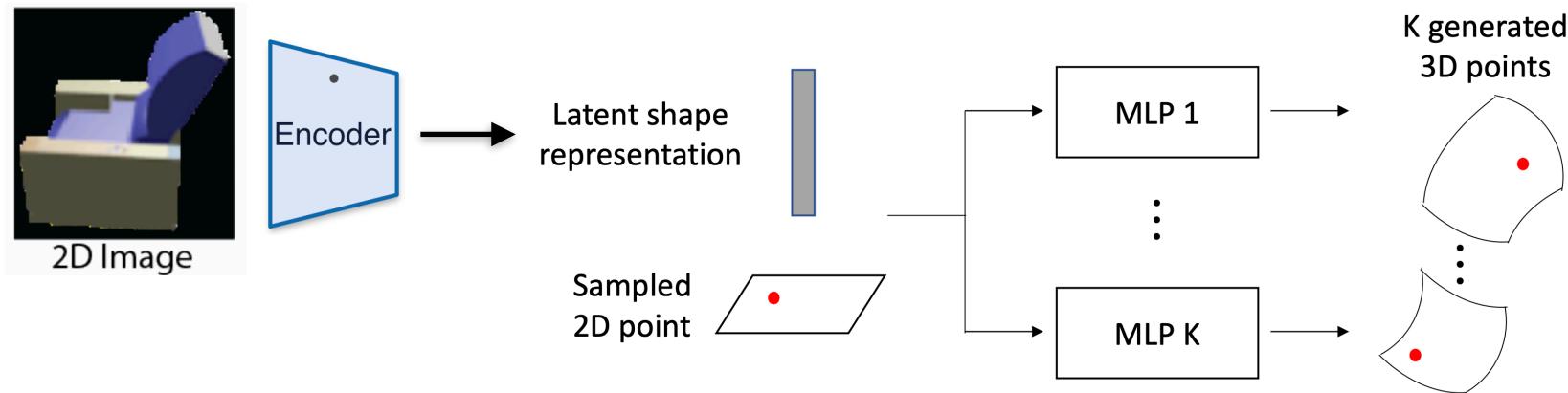
From Image to Shape

- Planes -> convex polytopes -> shape



From Image to Surface

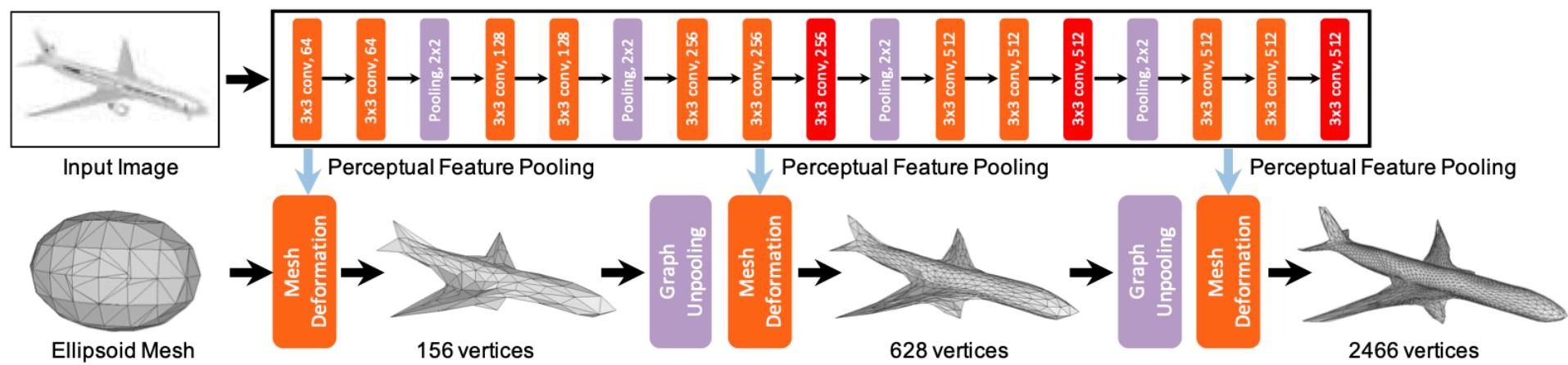
- Learn to warp a plane to surface



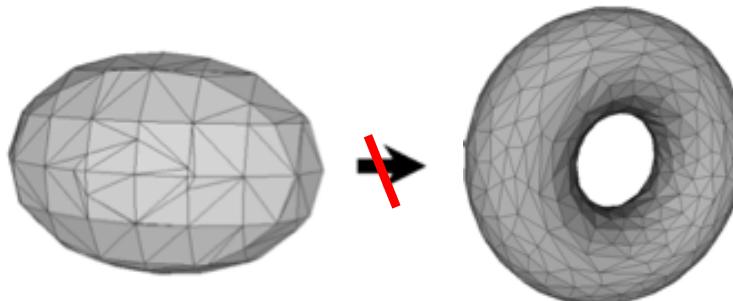
Groueix et al., "AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation", CVPR 2018
Yang et al., "Foldingnet: Point cloud auto-encoder via deep grid deformation", CVPR 2018

From Image to Surface

- Polygon meshes are irregular
- Learn to deform a template mesh



cannot change the topology of the template mesh

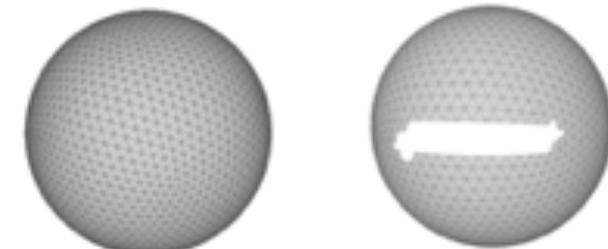


From Image to Surface

part-level deformation

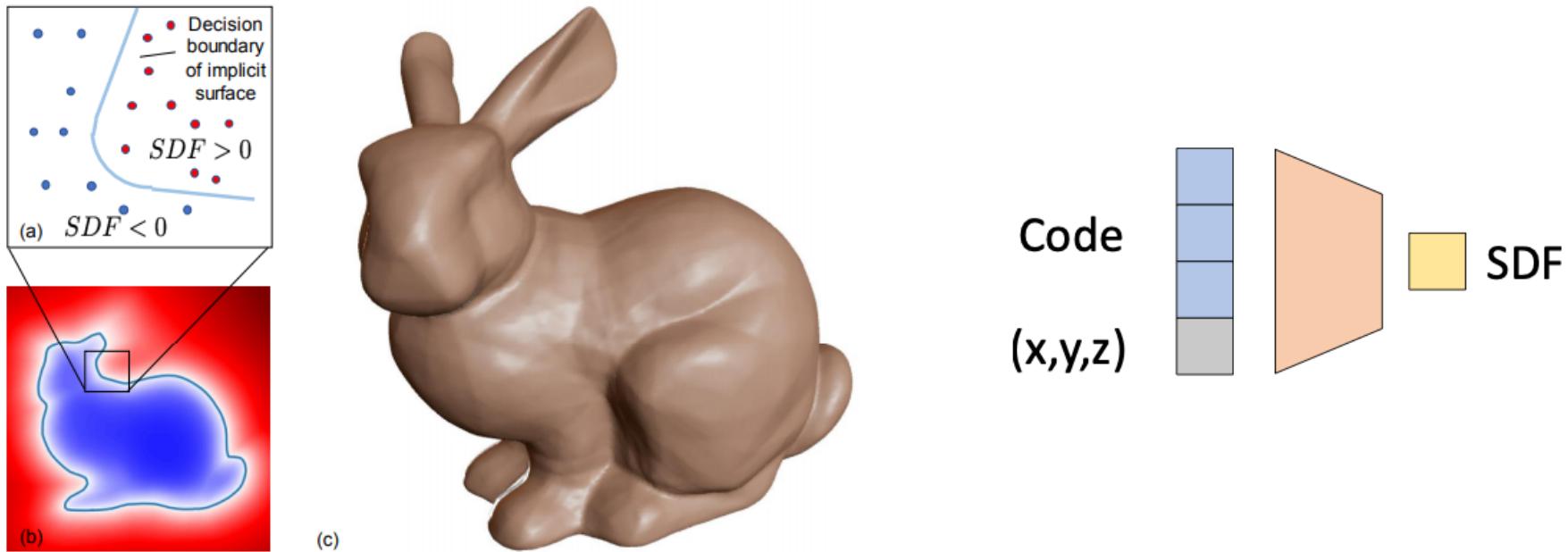


modify the topology of
the template mesh



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

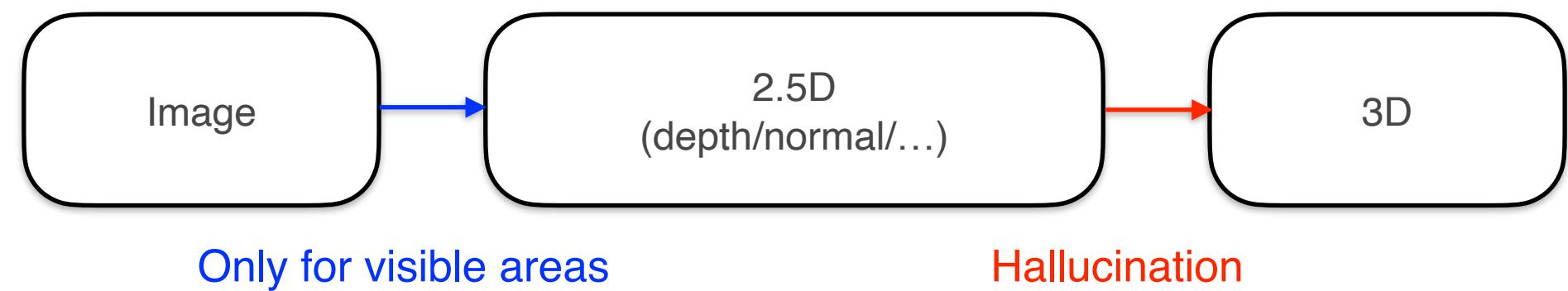
Other two similar paper on implicit representation:

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

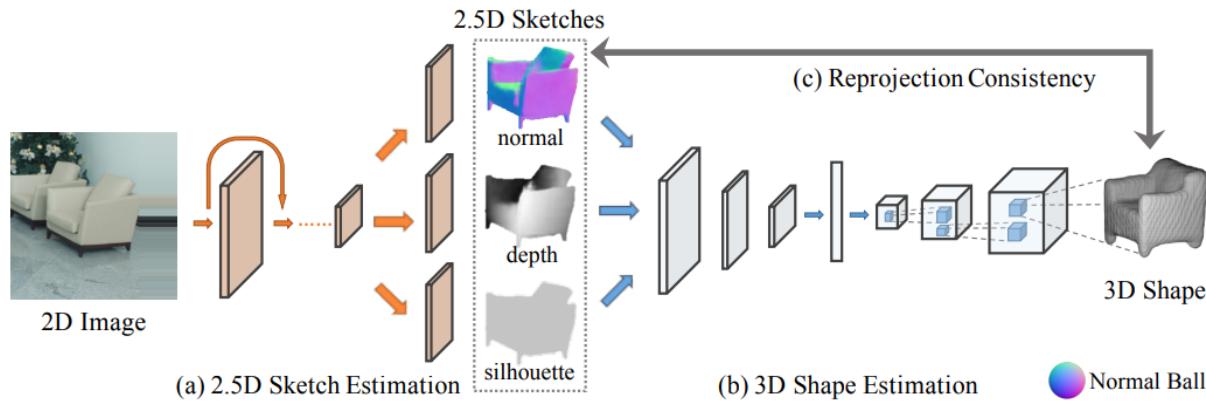
Chen et al., “Learning Implicit Fields for Generative Shape Modeling”, CVPR 2019

- In general,
 - First map the input to a shape embedding
 - Then reconstruct by decoding
- Limitation
 - Output is not explicitly grounded on the input
 - Structures of 3D objects are not explicitly leveraged
 - Cannot generalize to unseen objects

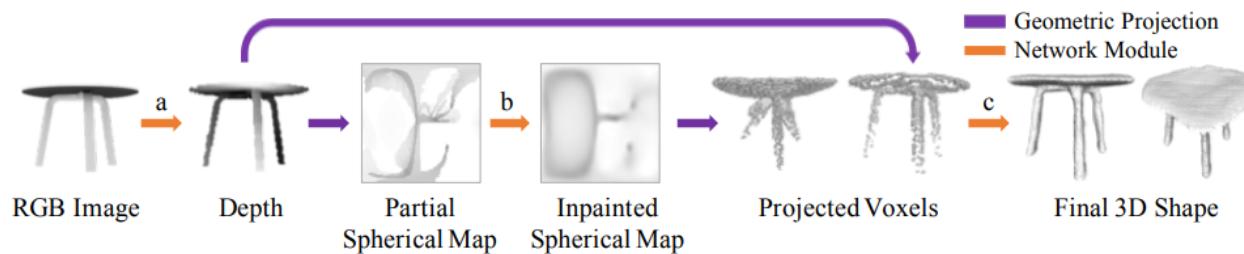
Visually Grounded Prediction: 2.5D to Bridge



Visually Grounded Prediction: 2.5D to Bridge



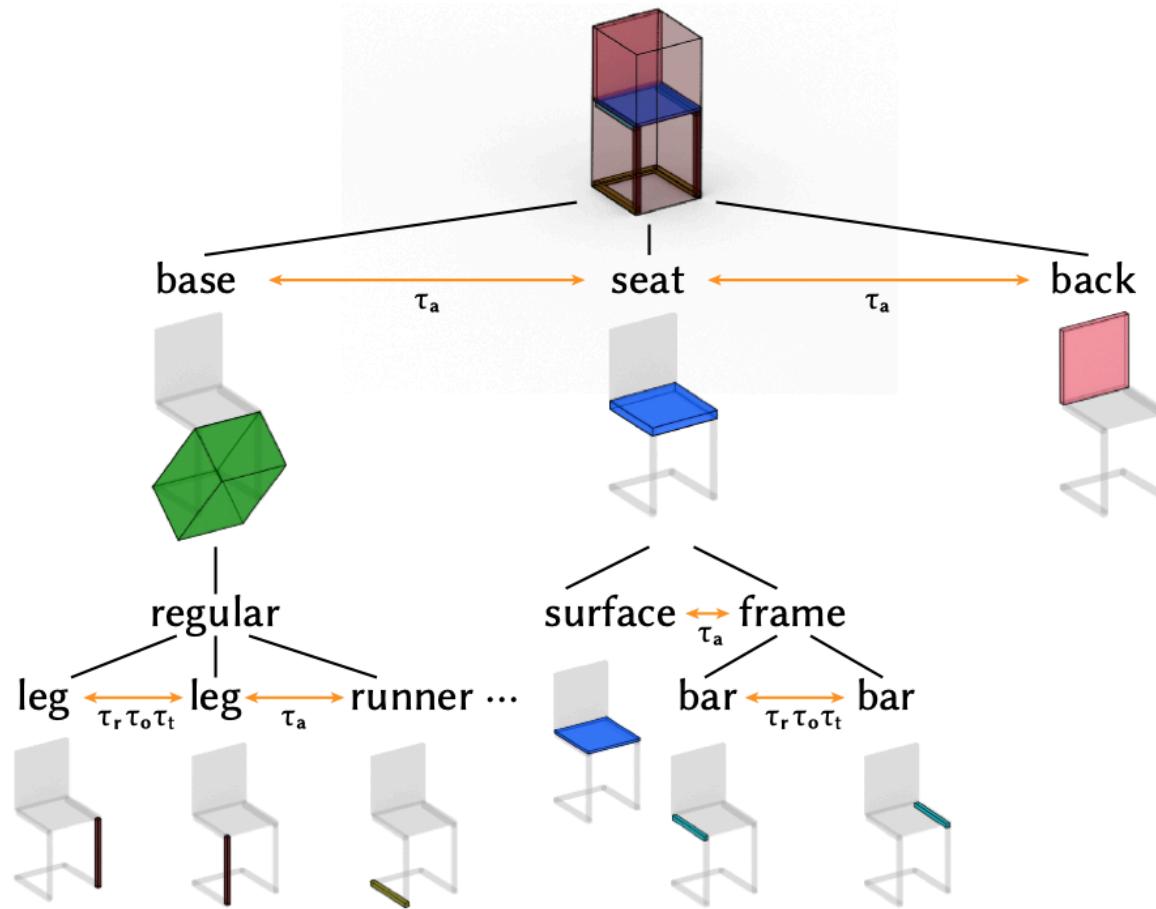
Wu et al., “MarrNet: 3D Shape Reconstruction via 2.5D Sketches”, NeurIPS 2017



Zhang et al., “Learning to Reconstruct Shapes from Unseen Classes”, NeurIPS 2018

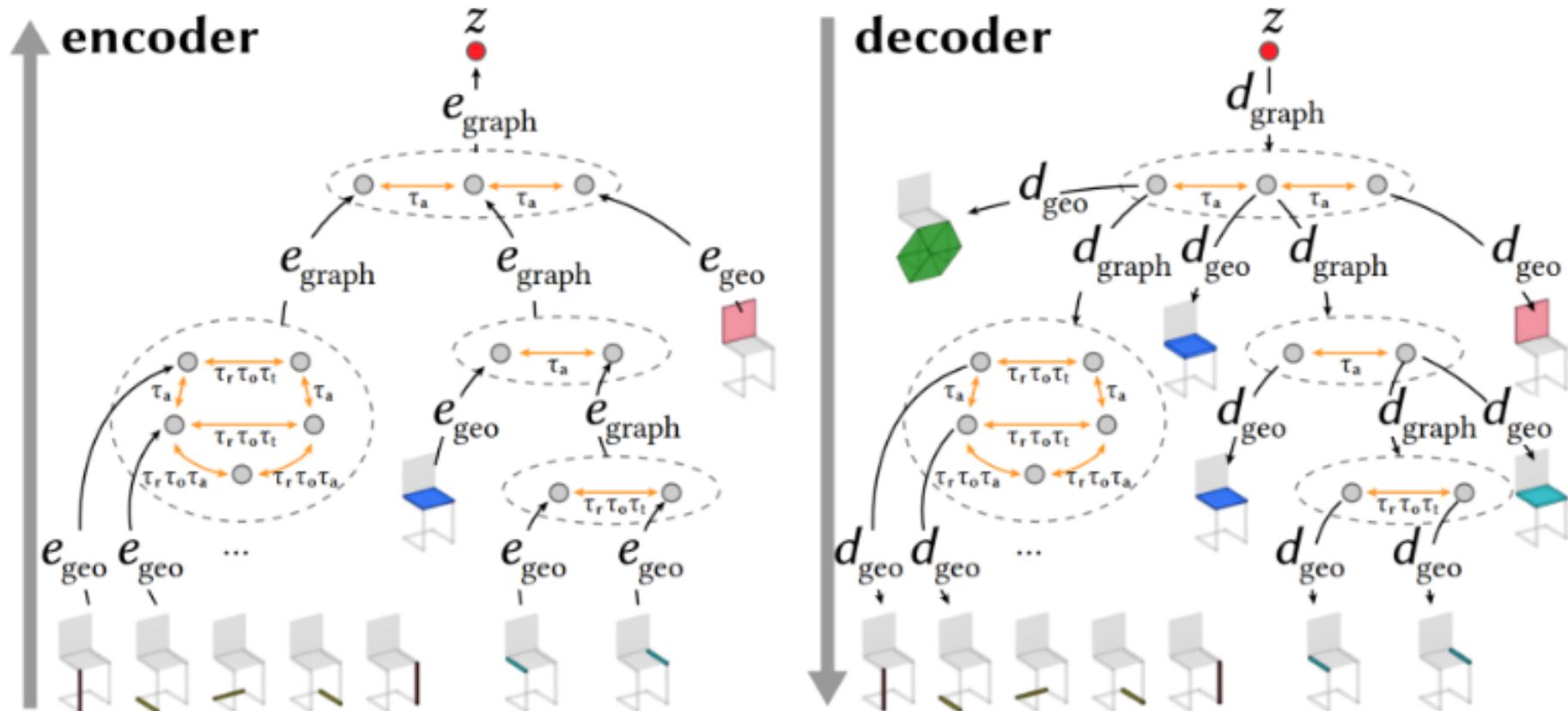
Structured Prediction: Part-based

Hierarchical Graph



Structured Prediction: Part-based

Recursive Network for Hierarchical Graph AE



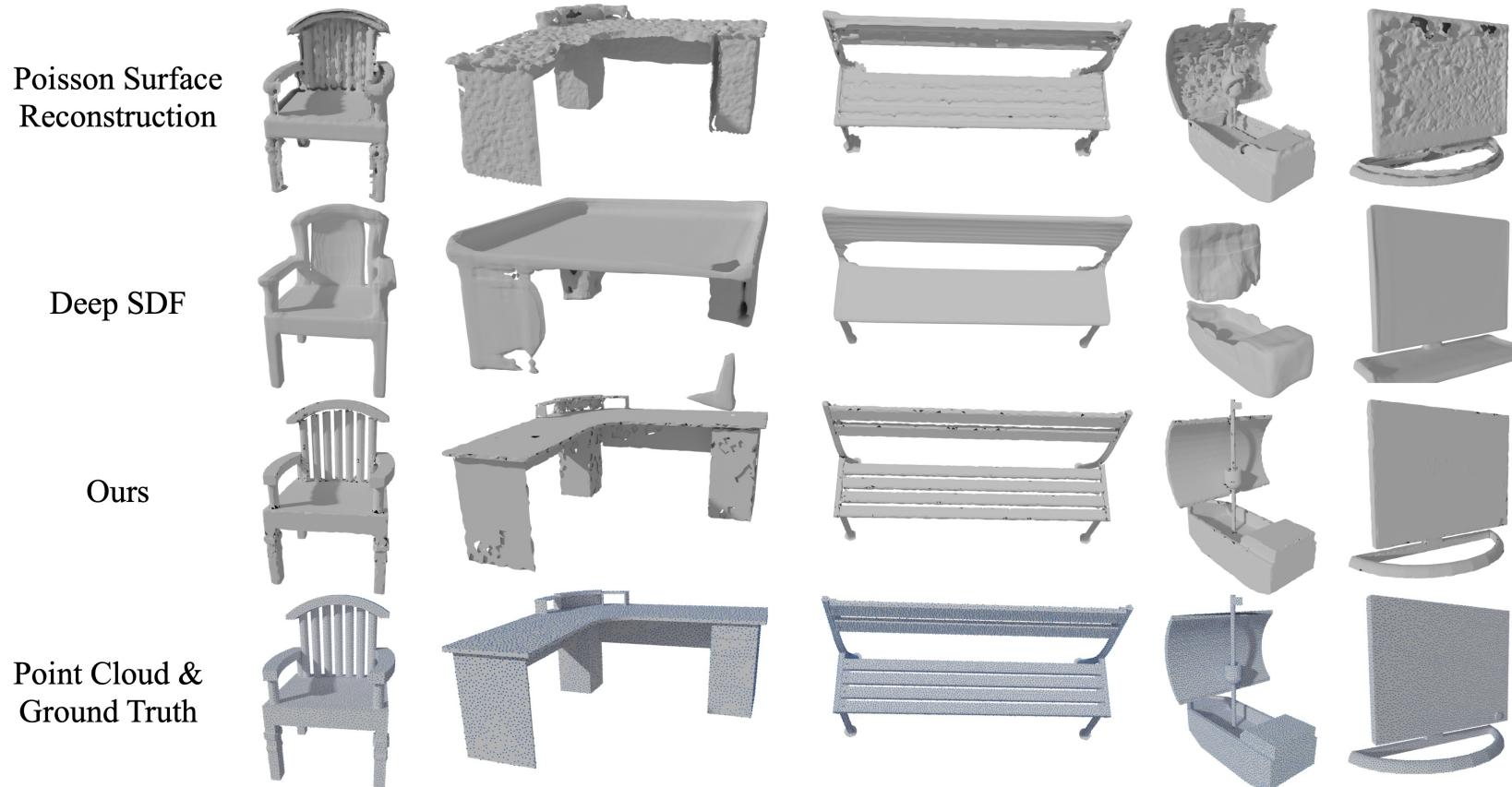
Structured Prediction: Part-based



Mo et al., “**StructureNet, a hierarchical graph network for learning PartNet shape generation**”, *Siggraph Asia 2019*

From Point Cloud to Mesh

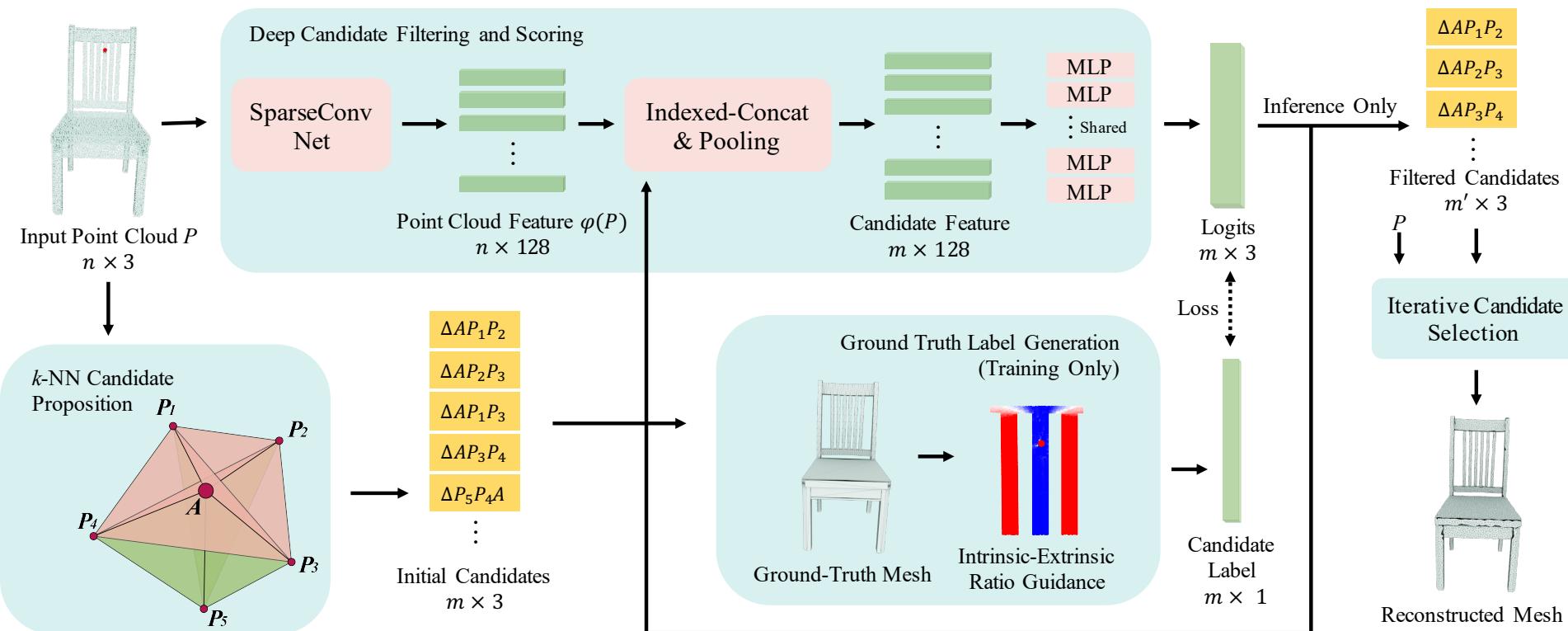
- Traditional methods: cannot handle ambiguous structures (e.g., thin structures) when the resolution of point clouds is limited
- Existing learning-based methods: cannot generate fine-grained details or generalize to unseen objects



From Point Cloud to Mesh

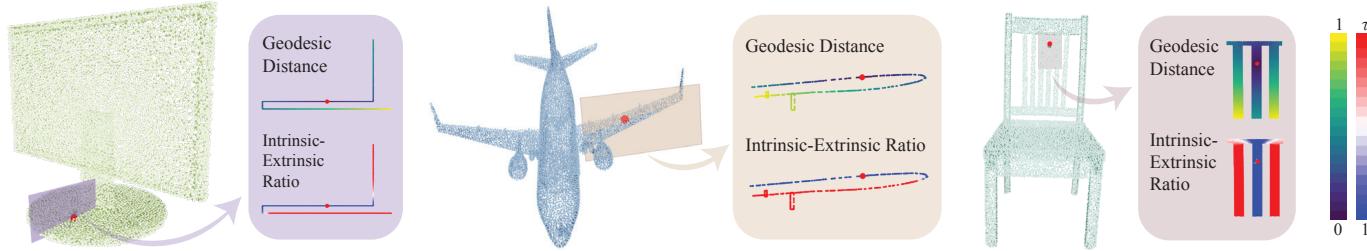
Fully utilize grounded input point clouds

1. Propose candidate triangles
2. Classify the candidates (filter out the incorrect triangles)
3. Merge the remaining candidates in a greedy way

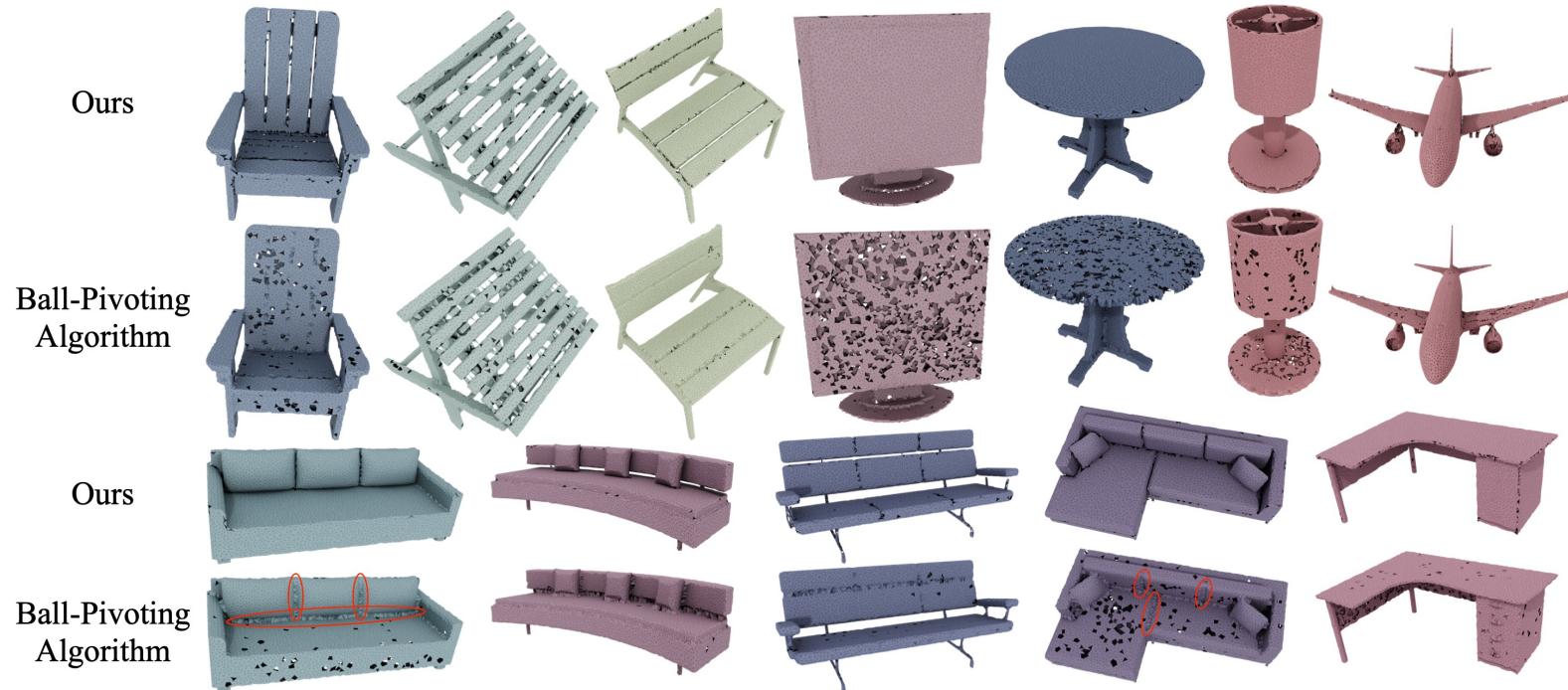


From Point Cloud to Mesh

Utilize “Intrinsic-Extrinsic Ratio” to label the candidate triangles



Leverage local priors and thus possess strong generalizability



Algorithm for Free Generation (GAN)

Challenges

Similar challenges as GAN for images:

- Good by human eye v.s. Good by objective metric

Metrics

Comparison between two sets of point clouds A and B:

- A: generated point clouds
- B: a held-out test set

Geometry Quality of Generated Shape

- e.g., minimum matching ($B \rightarrow A$) distance for CD/EMD

Coverage

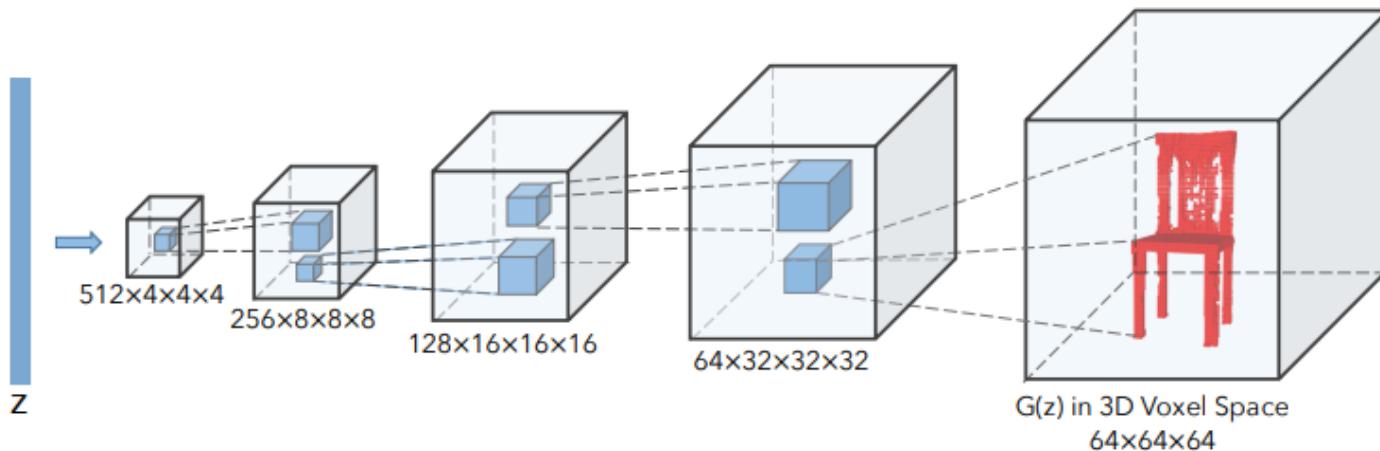
- The fraction of the shapes in B that were matched to shapes in A

Perceptually Correct

- Feature distribution distance (e.g. Frechet Point Cloud Distance)

$$\text{FPD}(\mathbb{P}, \mathbb{Q}) = \|\mathbf{m}_{\mathbb{P}} - \mathbf{m}_{\mathbb{Q}}\|_2^2 + \text{Tr}(\Sigma_{\mathbb{P}} + \Sigma_{\mathbb{Q}} - 2(\Sigma_{\mathbb{P}} \Sigma_{\mathbb{Q}})^{\frac{1}{2}})$$

Volumetric Generation



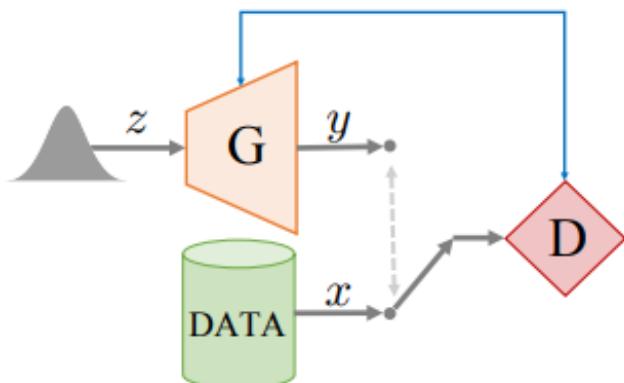
$$\log D(x) + \log(1 - D(G(z)))$$



Wu et al., "Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling", NeurIPS 2016

Point Cloud Generation

- FC layer as generator
- PointNet as discriminator
- WGAN



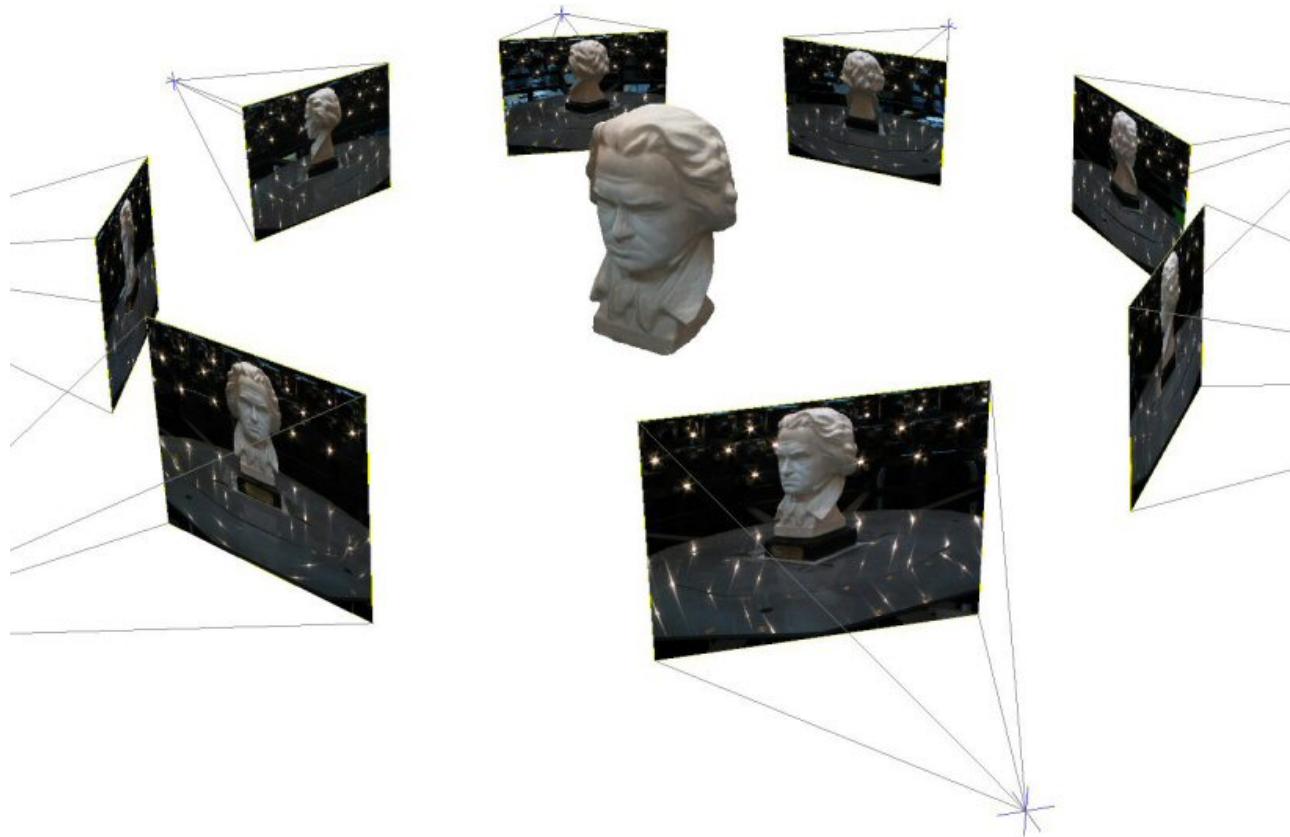
Many Issues

- 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

Topics

- 3D Data
- Classification
- Segmentation and Detection
- Reconstruction
 - Generative Model
 - Multi-View Stereo

Task: Reconstruction

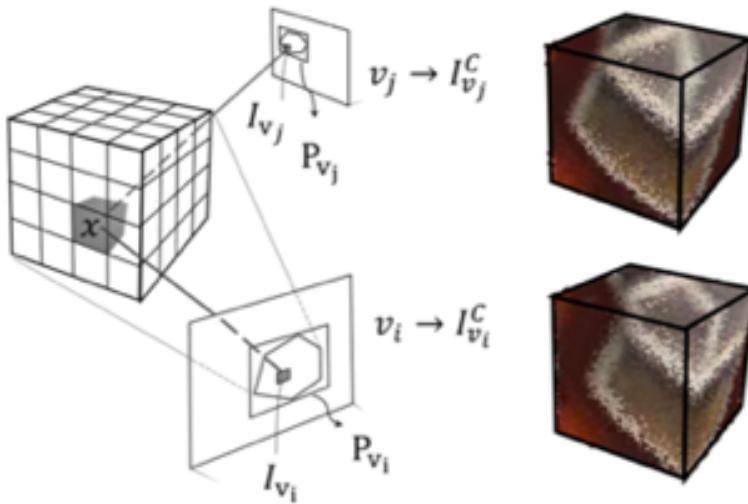


[image: oswald]

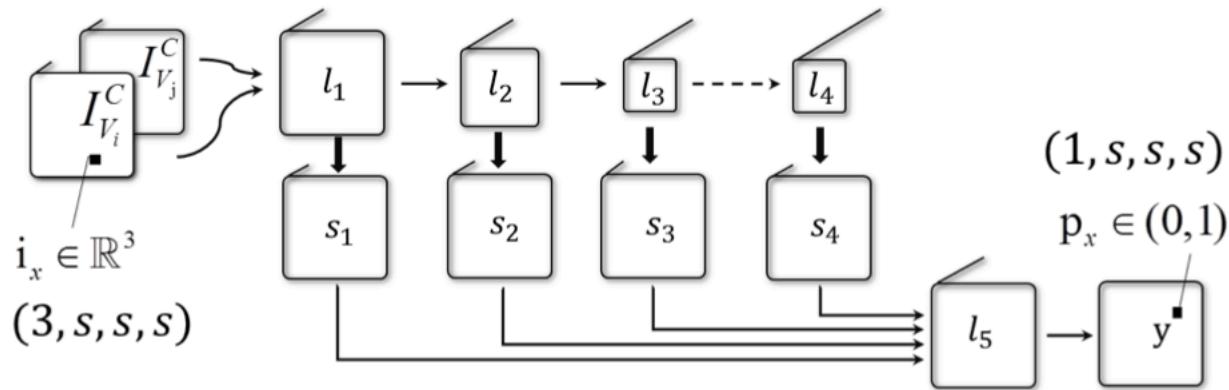
Covered methods: SurfaceNet, LSM, GC-Net, MVSNet, R-MVSNet, PointMVSNet, BA-Net

Surface Reconstruction as Voxel Occupancy Prediction

project along viewing rays to build colored voxel cubes.



Predict the surface confidence for each voxel:



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

binarization as post-processing

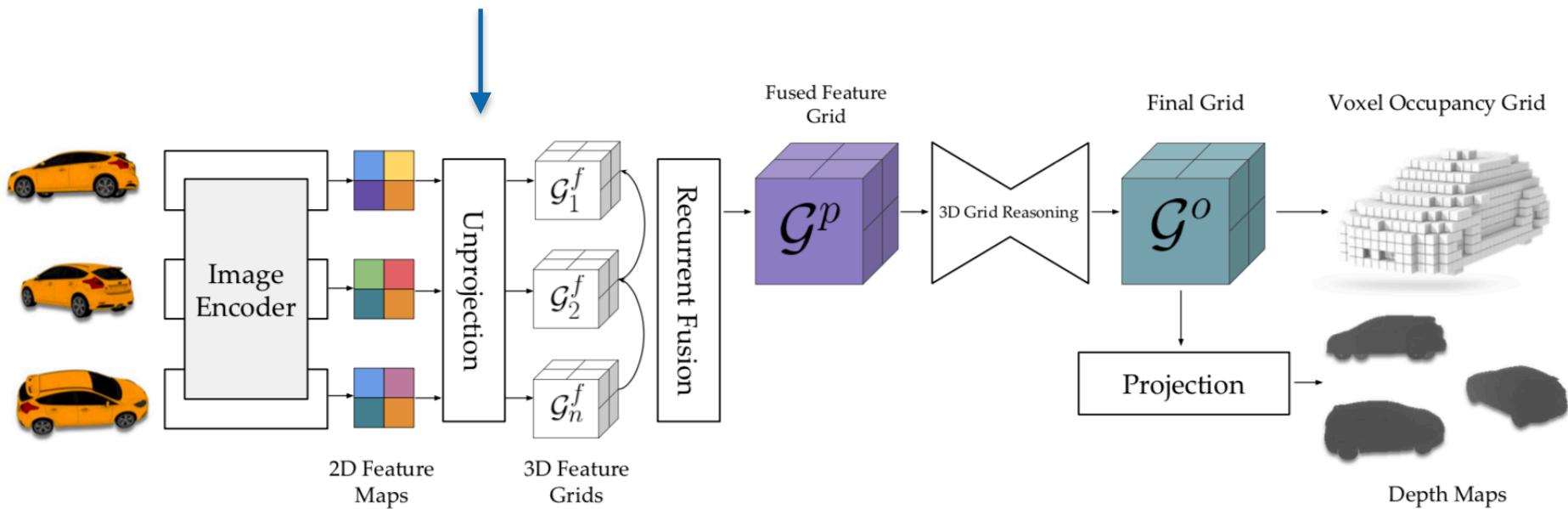
Limitations:

- Pre-computed voxel cubes can only take RGB colors at coarse resolution
- Voxel binarization introduces quantization errors.

Learning-Based Stereopsis

End-to-end learning of deep features for each image pixel

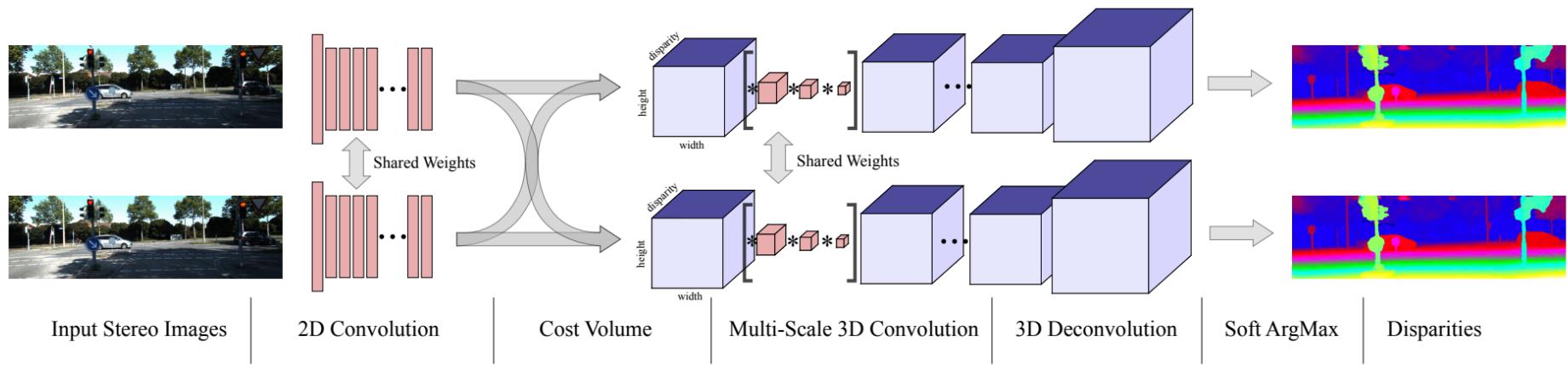
Unproject image features into 3D grids





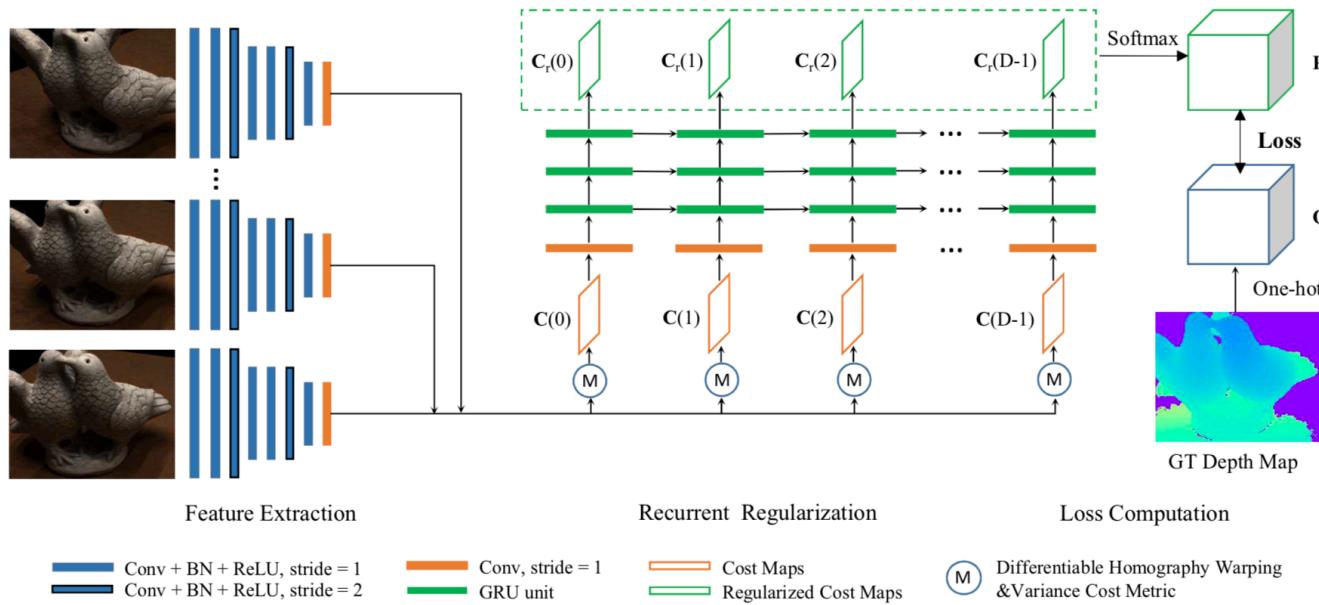
Still very coarse resolution (32x32x32) due to volumetric representation.

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



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

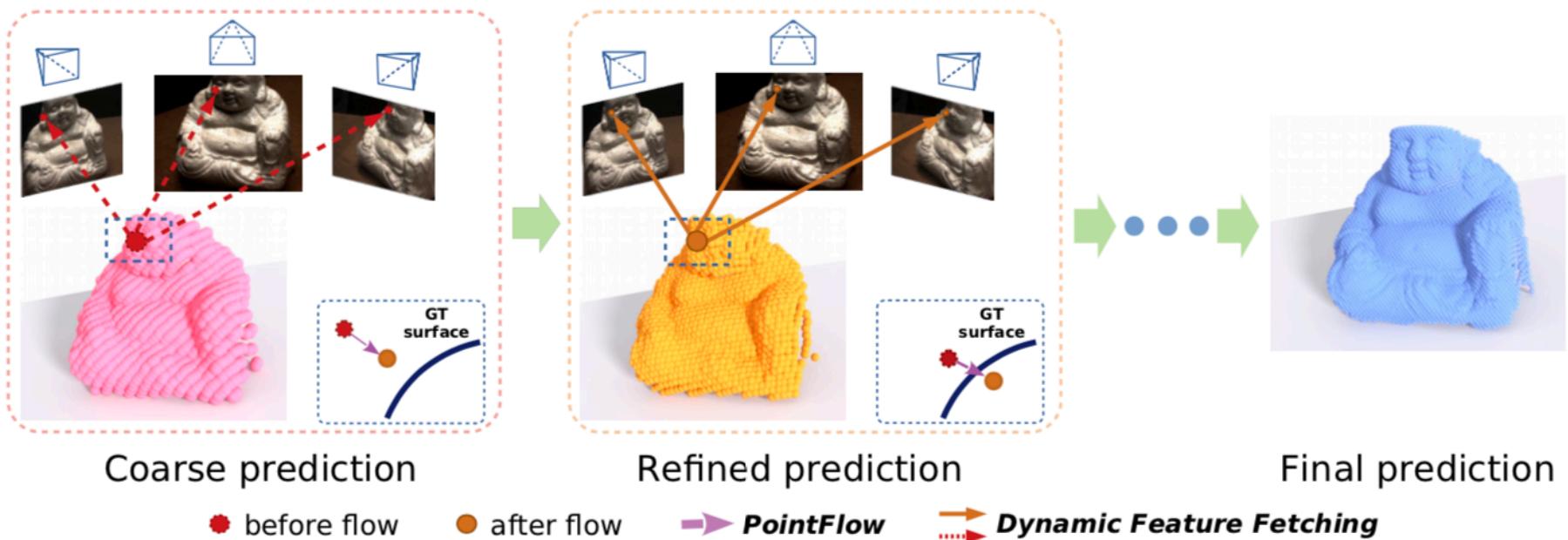
Idea 1: Slide-by-Slide Processing of Cost Volume by Recurrent Neural Network



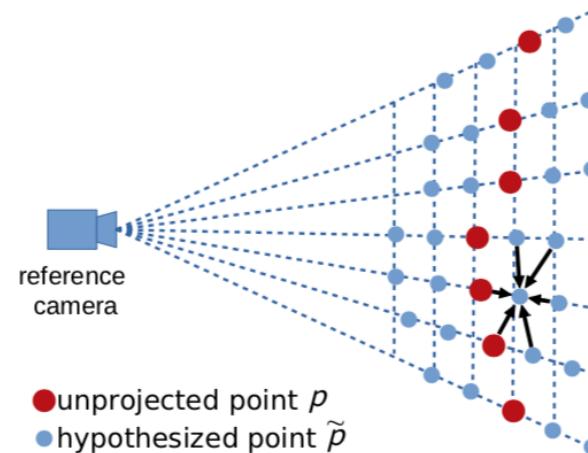
The cost volume is sequentially regularized along the depth direction.

Idea 2: Point-based MVS

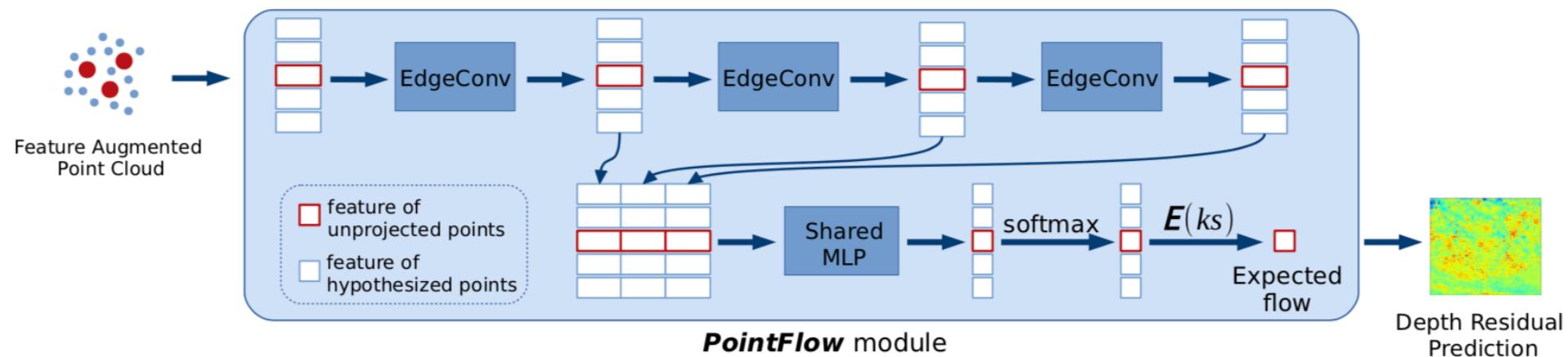
- Point-based representation for computational efficiency.
- Iteratively update the location of points and spawn more points.
- More flexible and accurate.



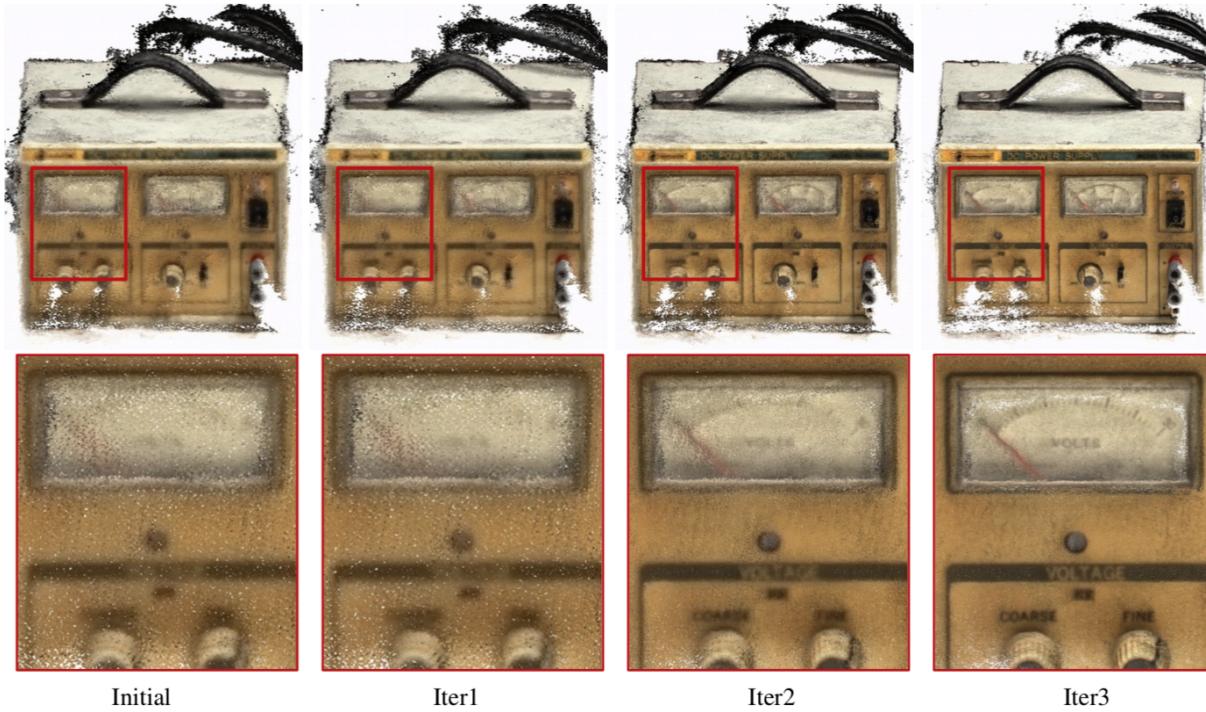
hypothesized point generation
along camera direction:



- PointFlow: iteratively pull all the points to the right position.



Iterative refinement:



Results on DTU benchmark

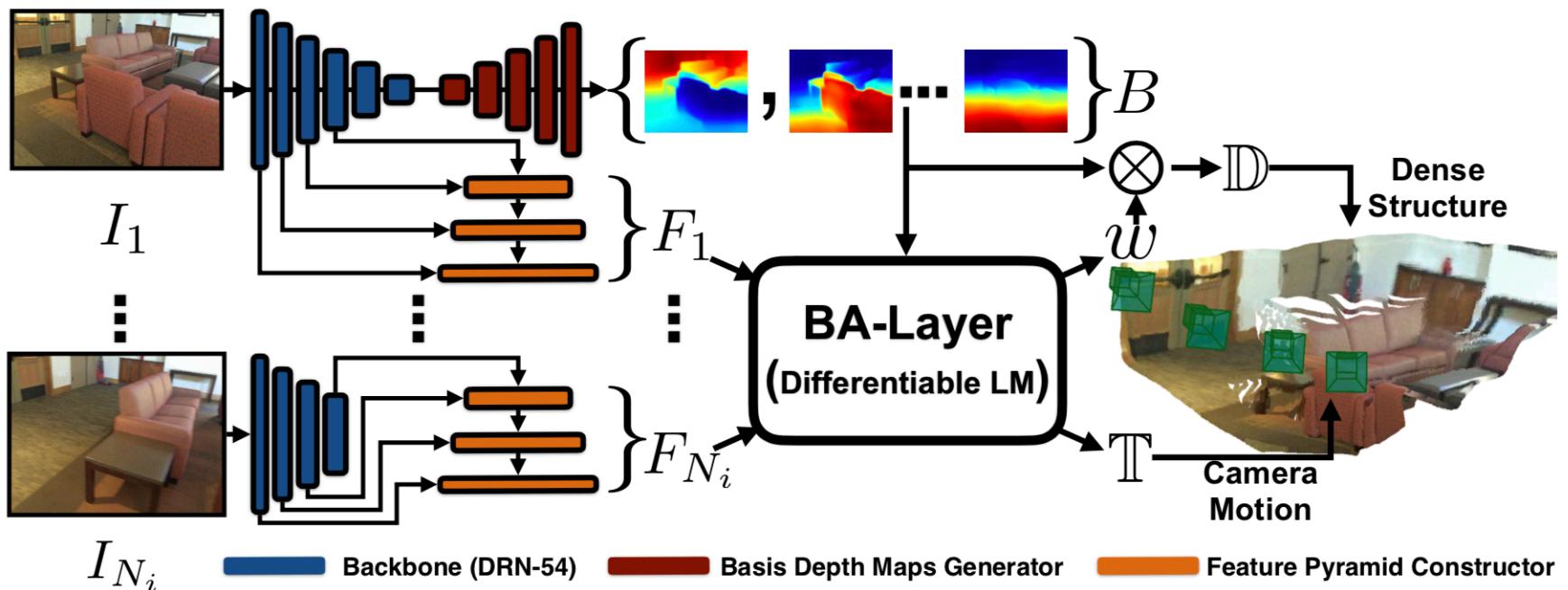
Iter.	Acc. (mm)	Comp. (mm)	Overall (mm)	0.5mm <i>f-score</i>	Depth Map Res.	Depth Interval (mm)	GPU Mem. (MB)	Runtime (s)
-	0.693	0.758	0.726	47.95	160×120	5.30	7219	0.34
1	0.674	0.750	0.712	48.63	160×120	5.30	7221	0.61
2	0.448	0.487	0.468	76.08	320×240	4.00	7235	1.14
3	0.361	0.421	0.391	84.27	640×480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

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

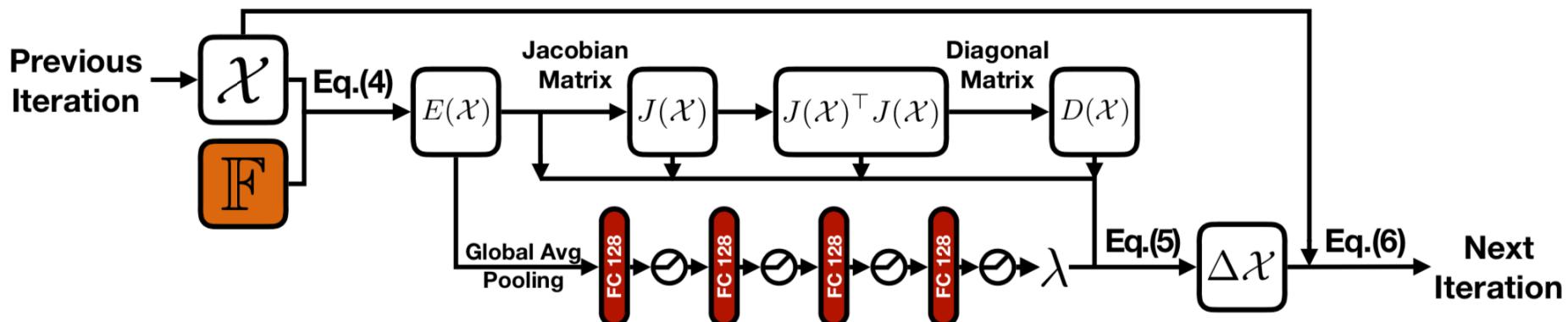
End-to-end pipeline for SfM with differentiable bundle adjustment.



Differentiable LM algorithm:

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

BA-layer:



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

Slides posted on
<http://ai.ucsd.edu/~haosu/>

(Homepage of Prof. Hao Su)