

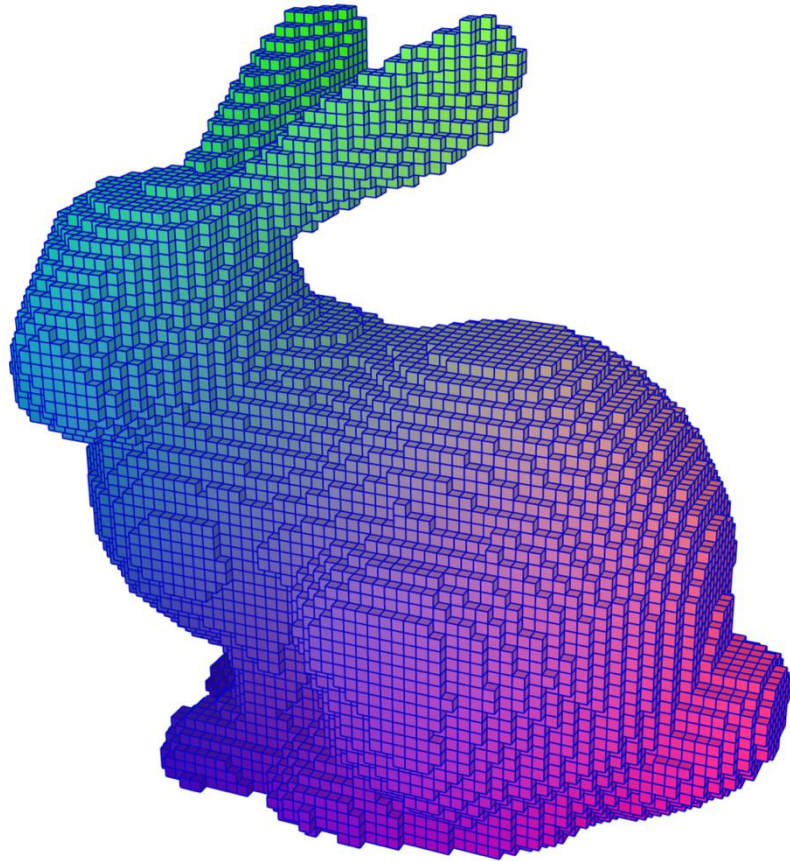
Deep 3D Vision

Xiaolong Wang

This Class – 3D representations

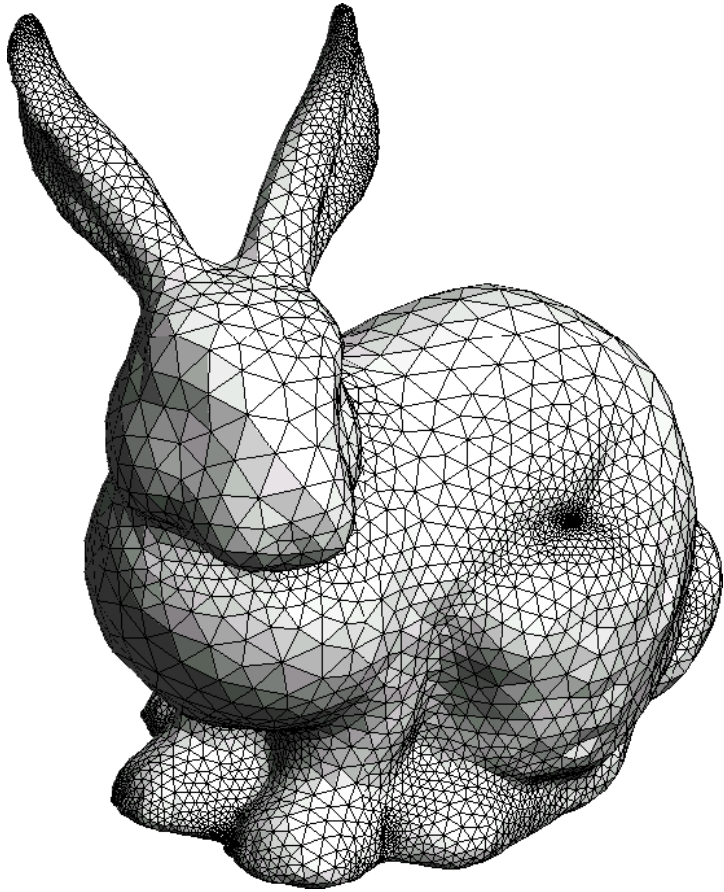
- Volumetric representation
- Mesh
- Point cloud
- Implicit functions

How to represent objects in 3D?



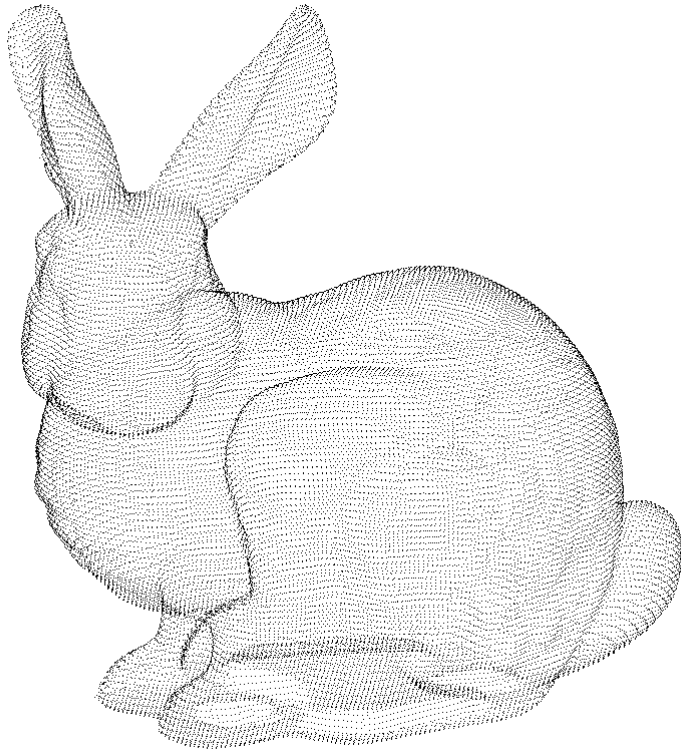
- Volumetric
- Mesh
- Point cloud
- Implicit function

How to represent objects in 3D?



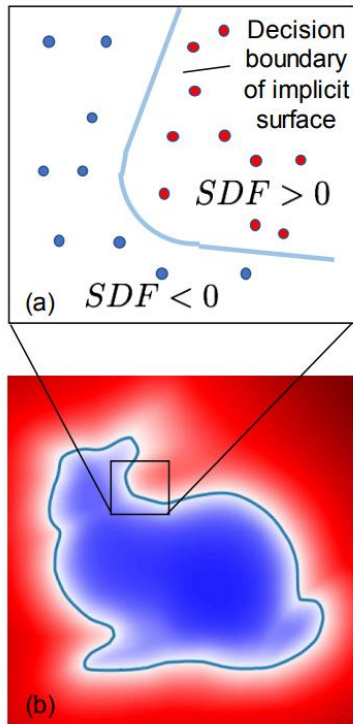
- Volumetric
- **Mesh**
- Point cloud
- Implicit function

How to represent objects in 3D?



- Volumetric
- Mesh
- **Point cloud**
- Implicit function

How to represent objects in 3D?



- Volumetric
- Mesh
- Point cloud
- **Implicit function**

1. Voxel: Discretizing into grids



Voxels

Pros:

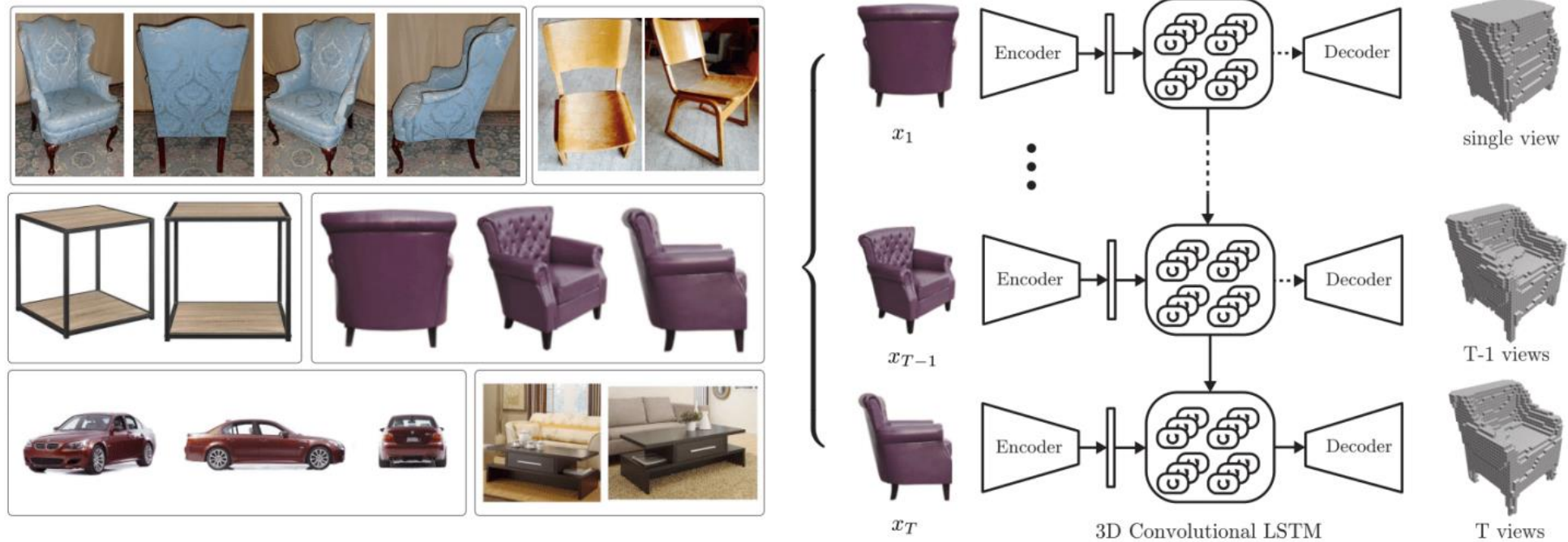
- *Easy to process with Neural Networks*

Cons:

- *Cubic memory*
- *Limited resolution*

3D-R2N2

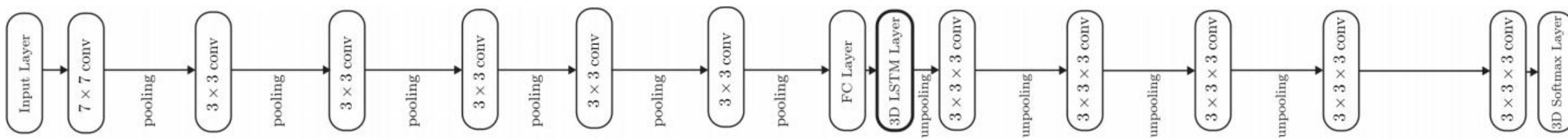
Given single/multi-view images, output 3D voxel occupancies



(a) Images of objects we wish to reconstruct (b) Overview of the network

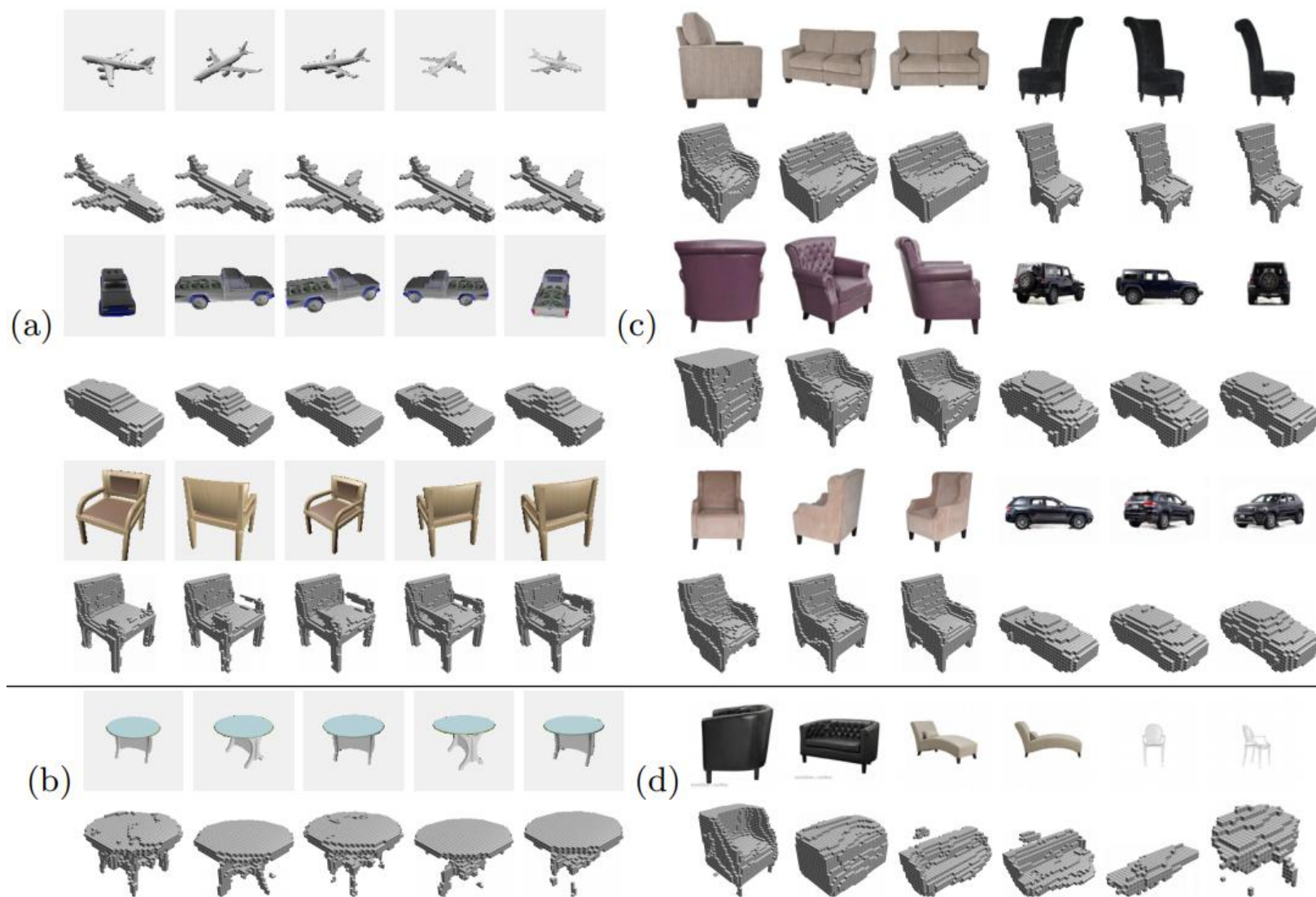
3D-R2N2

Overall architecture



3D-R2N2

Do not need to be trained
per category.



2.Mesh: Discretizing into vertices



Mesh

Pros:

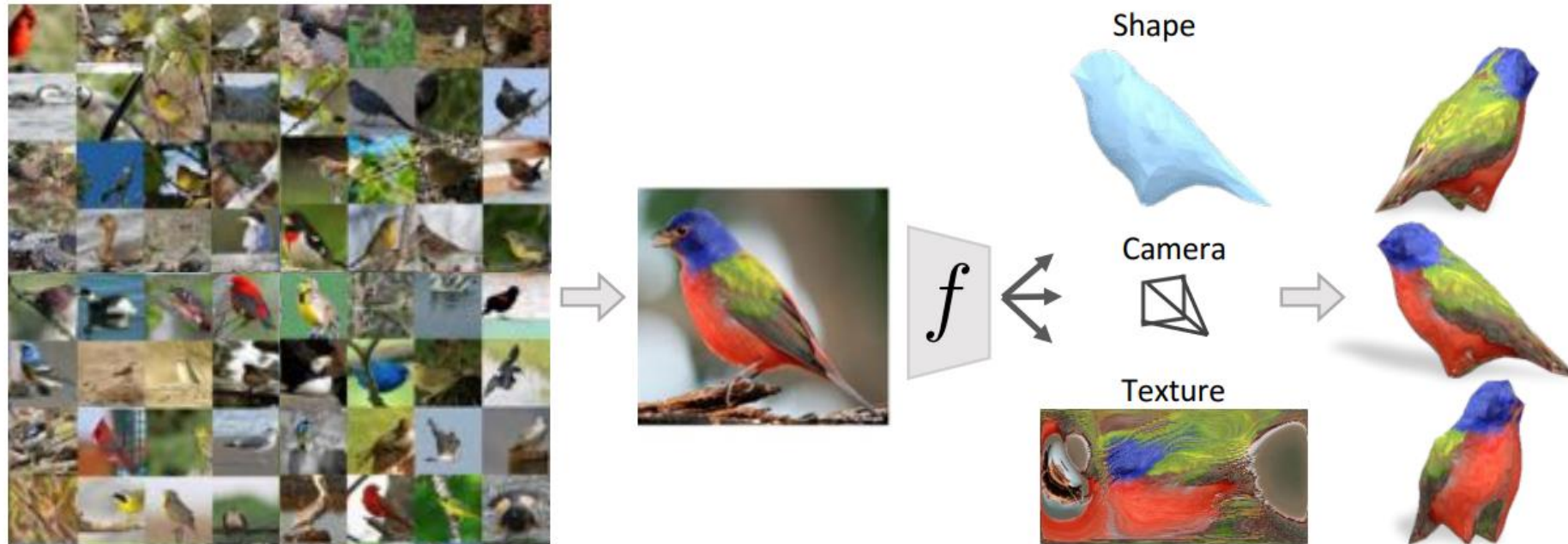
- *Easy to process with Neural Networks*

Cons:

- *Usually require class-specific template*
- *Limited by the number of vertices*

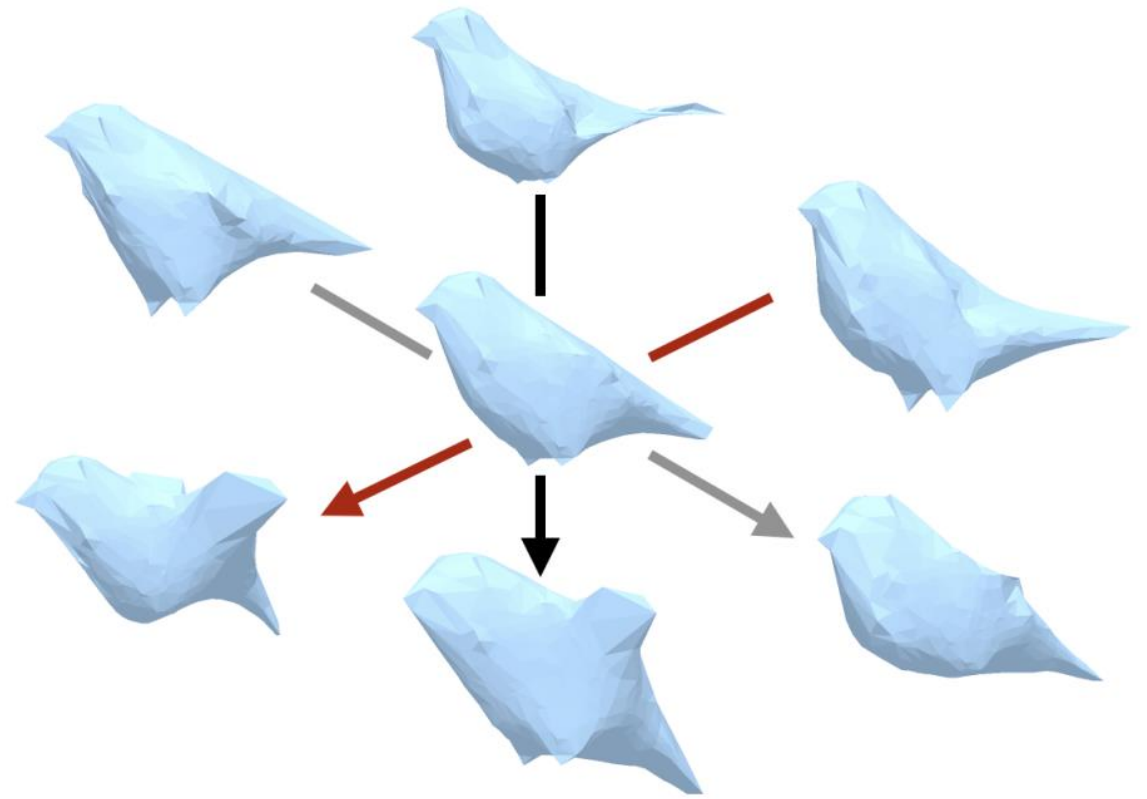
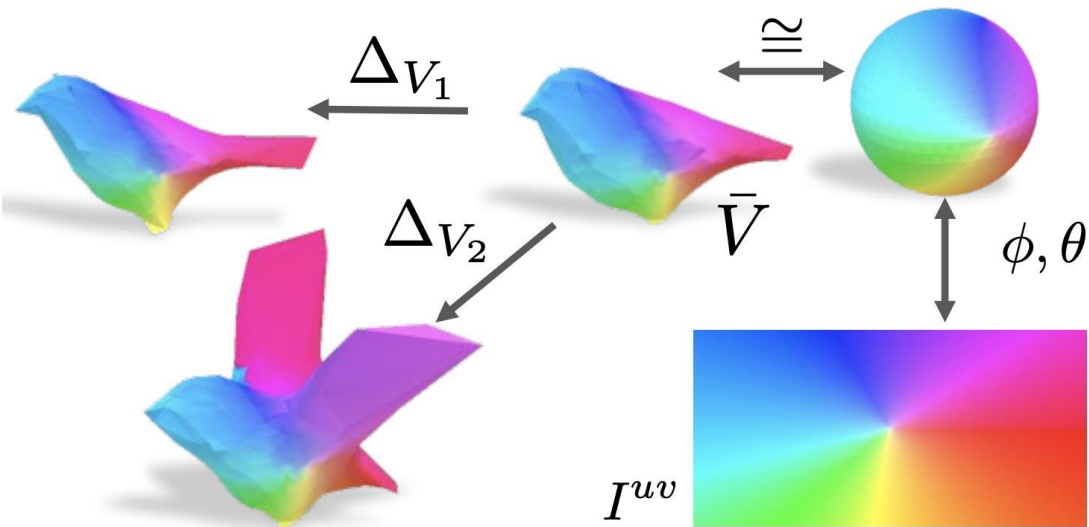
Learn mesh deformation

1. *Given a collection of images from the same category*
2. *Jointly predict 3D mesh, texture and camera parameters*



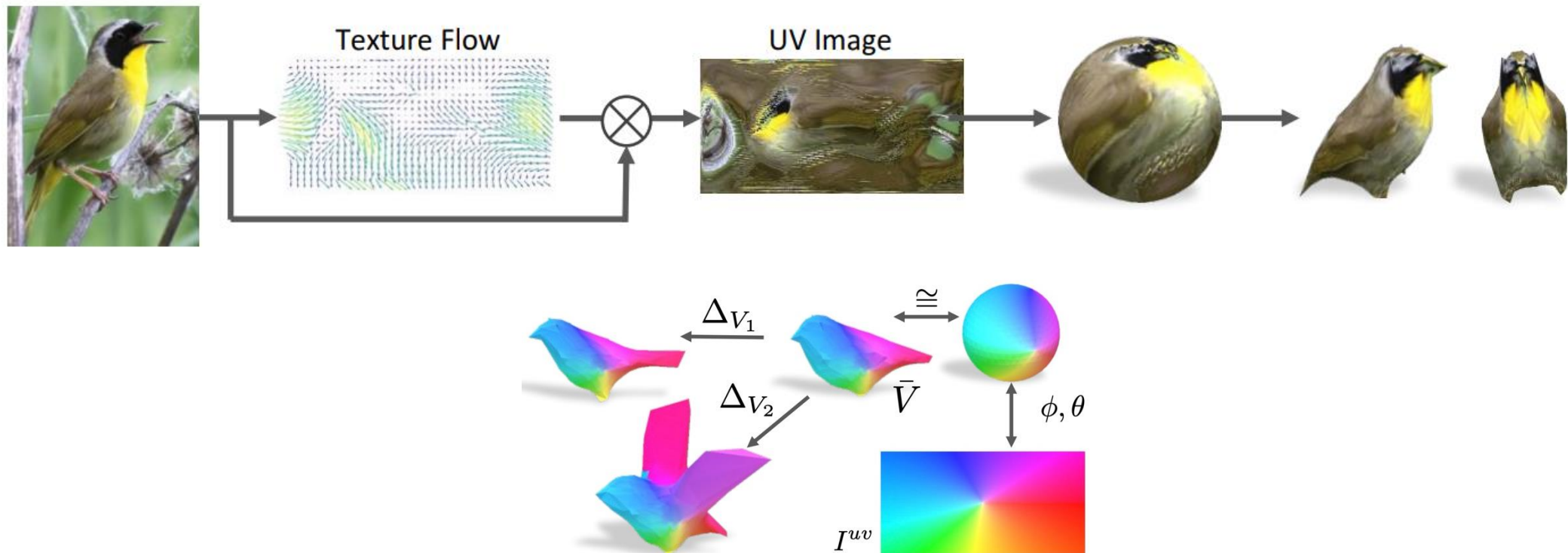
Learn mesh deformation

Key 1: Learned deformation from mean shape

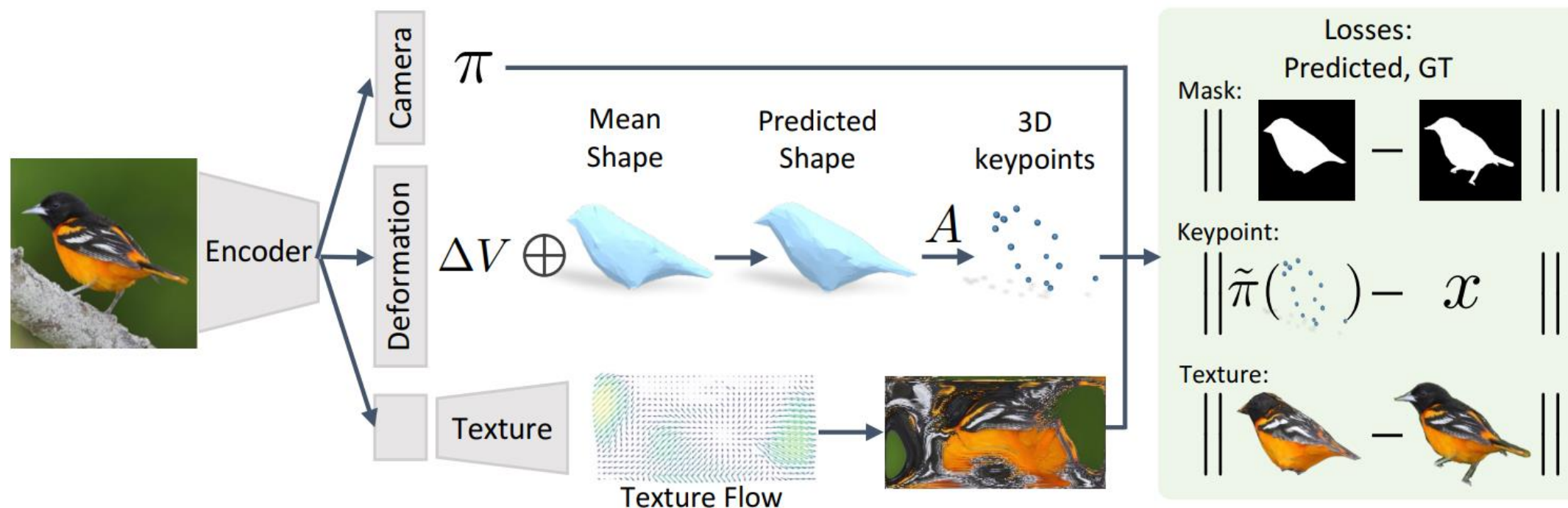


Learn mesh deformation

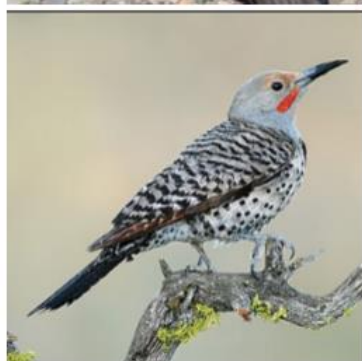
Key 2: Predict texture UV map



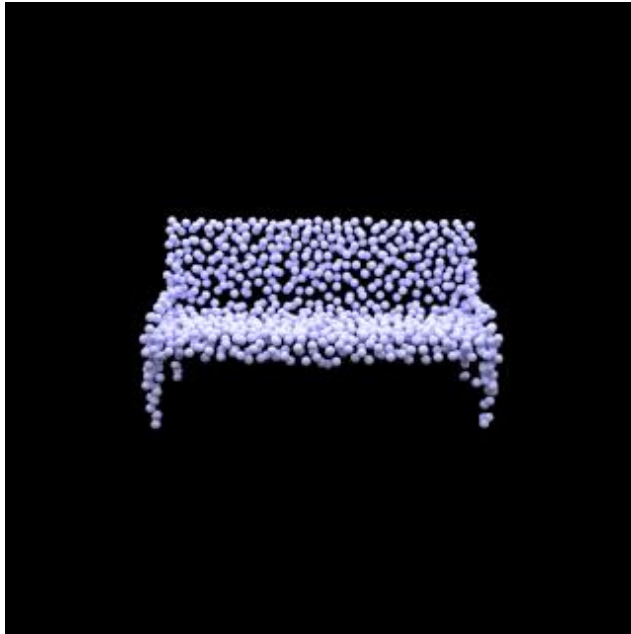
Learn mesh deformation



Learn mesh deformation



3.Point cloud: Discretizing into points



Point cloud

Pros:

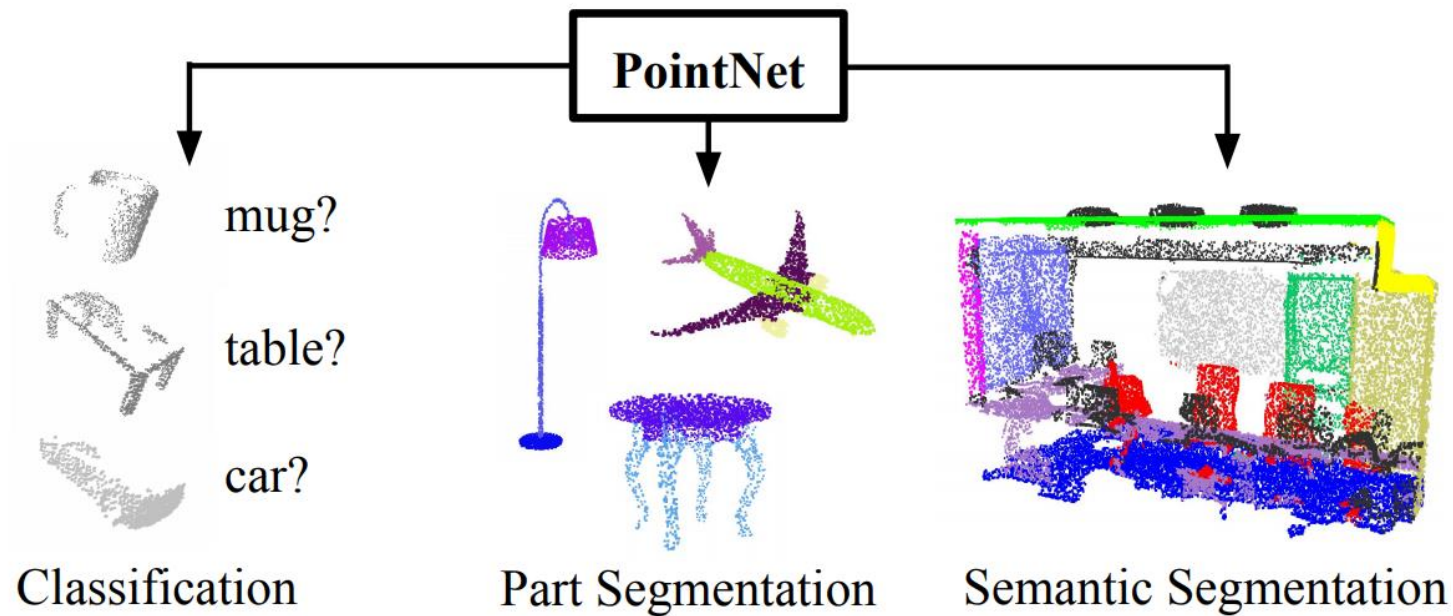
- *Flexible and memory efficient*

Cons:

- *Not model topology / connection*
- *Limited by the number of input points*

PointNet

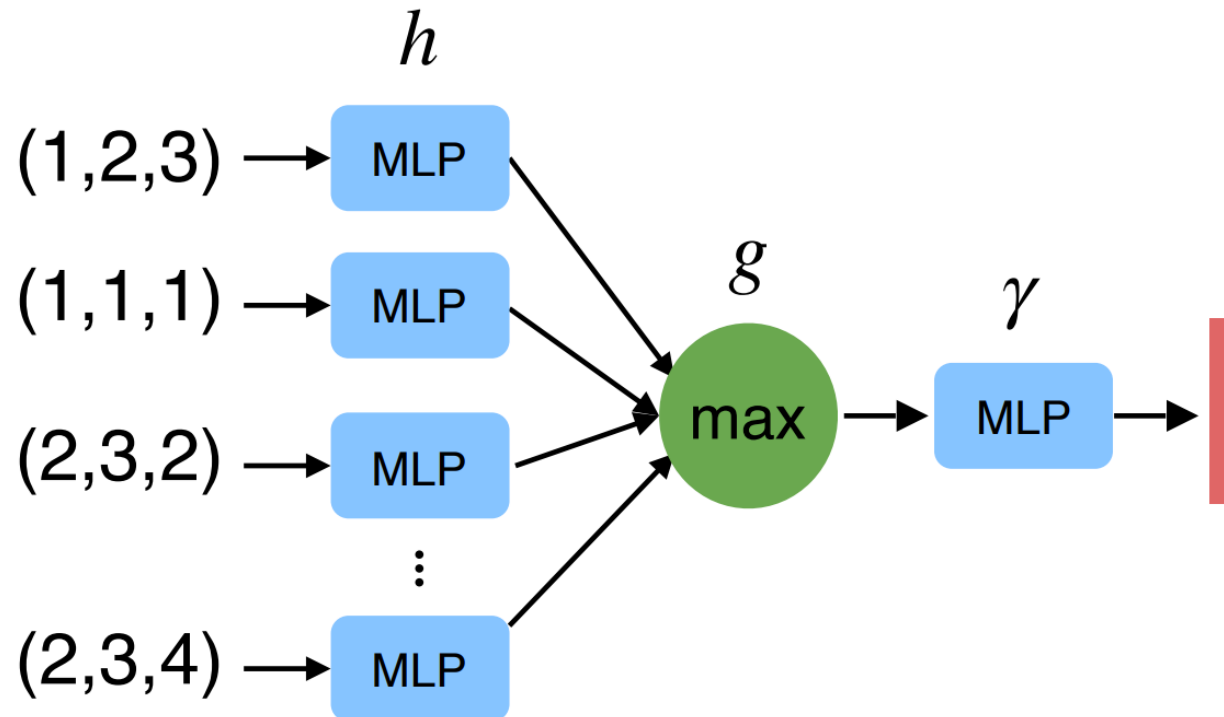
End-to-end learning given unordered and scattered point cloud, extracted features can be applied to multiple down-streaming tasks



PointNet

Key 1: the model should be invariant to input point permutations

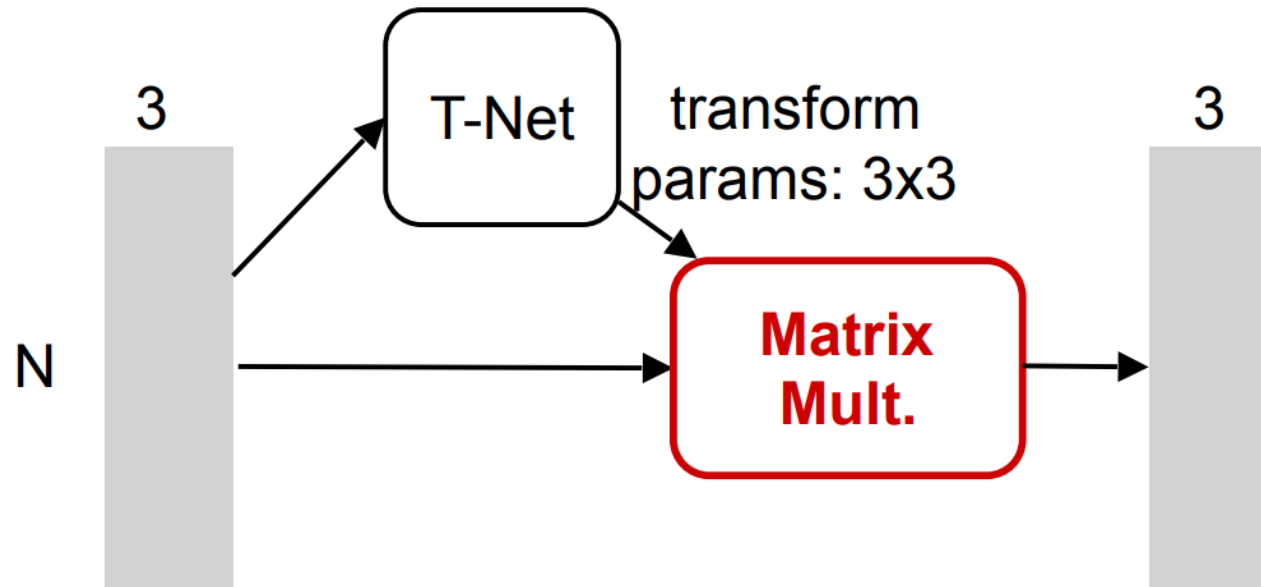
Proposed solution: use symmetric functions, e.g., max pooling



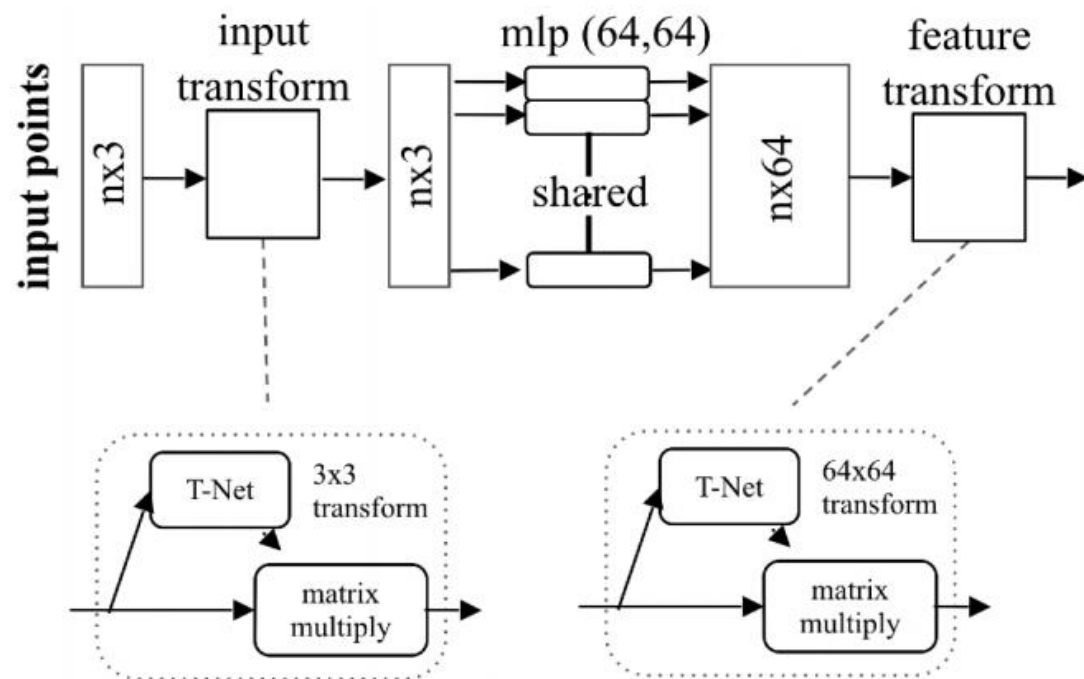
PointNet

Key 2: invariance under geometric transformations

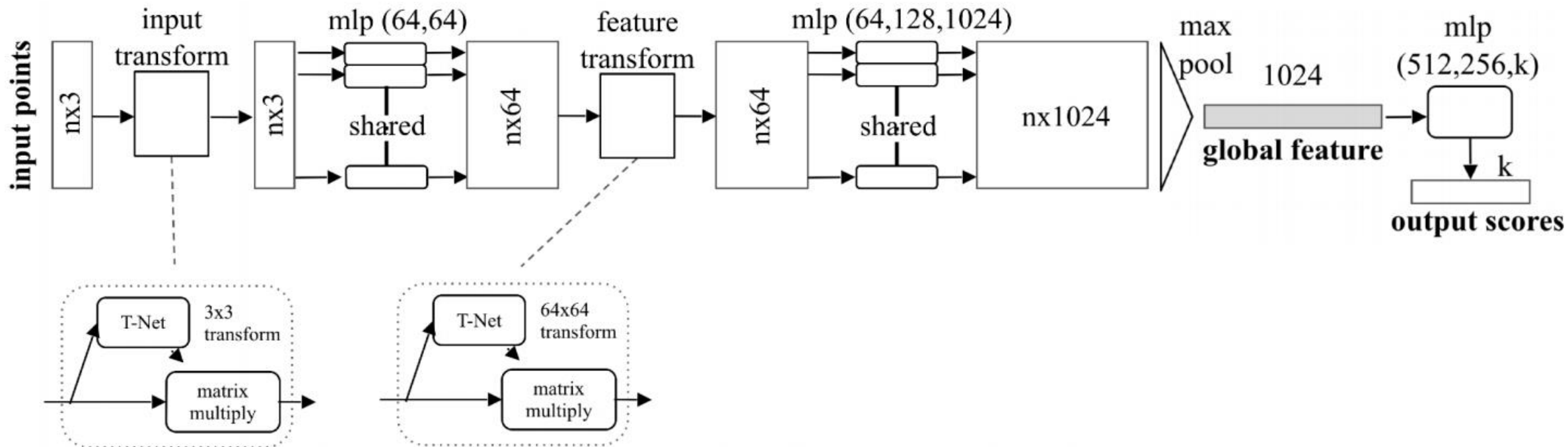
Proposed solution: affine transformation



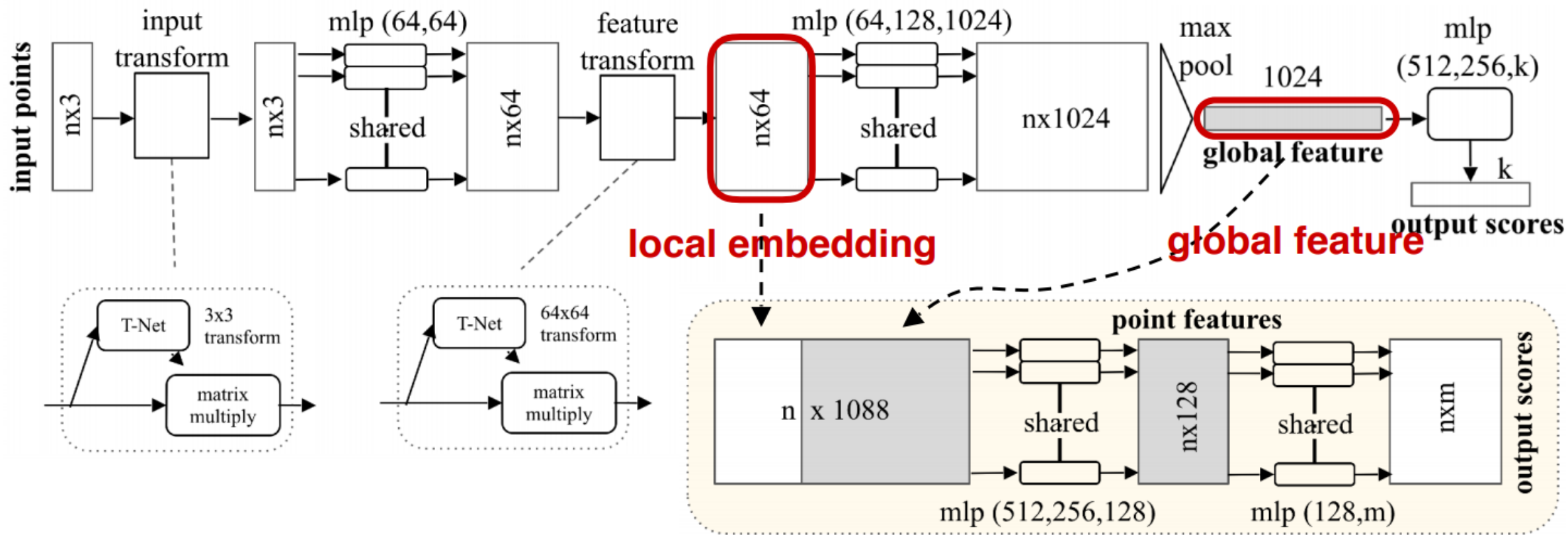
PointNet



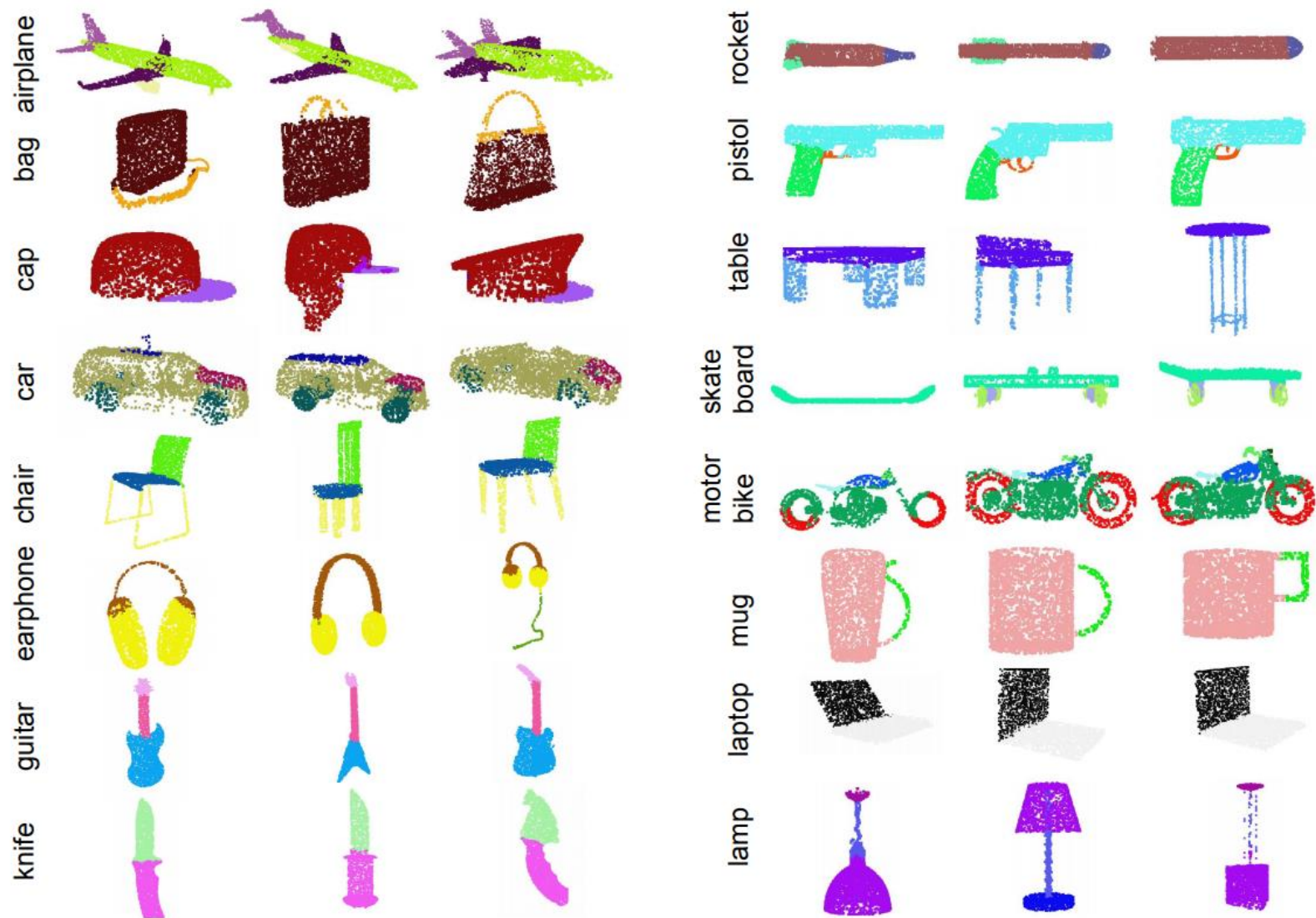
PointNet



PointNet



PointNet



4. Implicit function: continuous representation



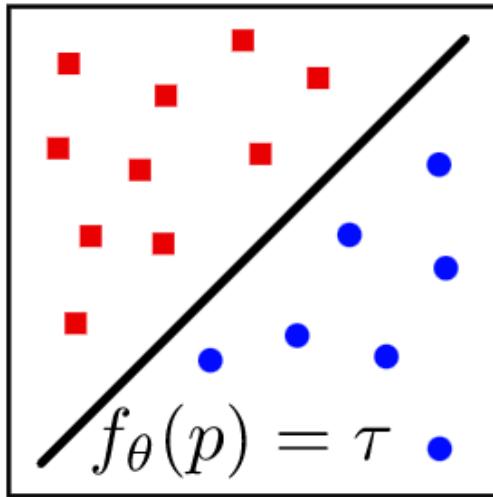
Pros:

- *No discretization*
- *Arbitrary resolution and memory efficient*
- *Arbitrary topology*
- *Not restricted to specific class*

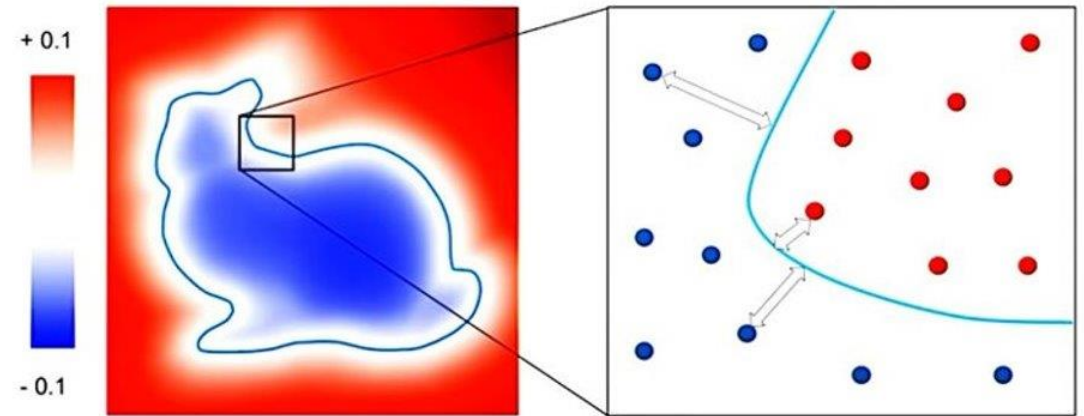
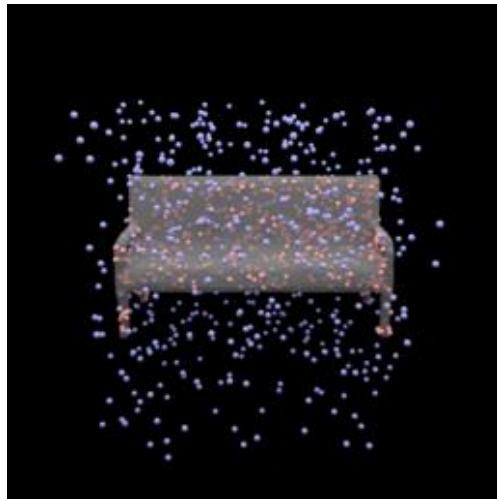
Cons:

- *Post-processing to extract meshes*

Occupancy VS Signed Distance

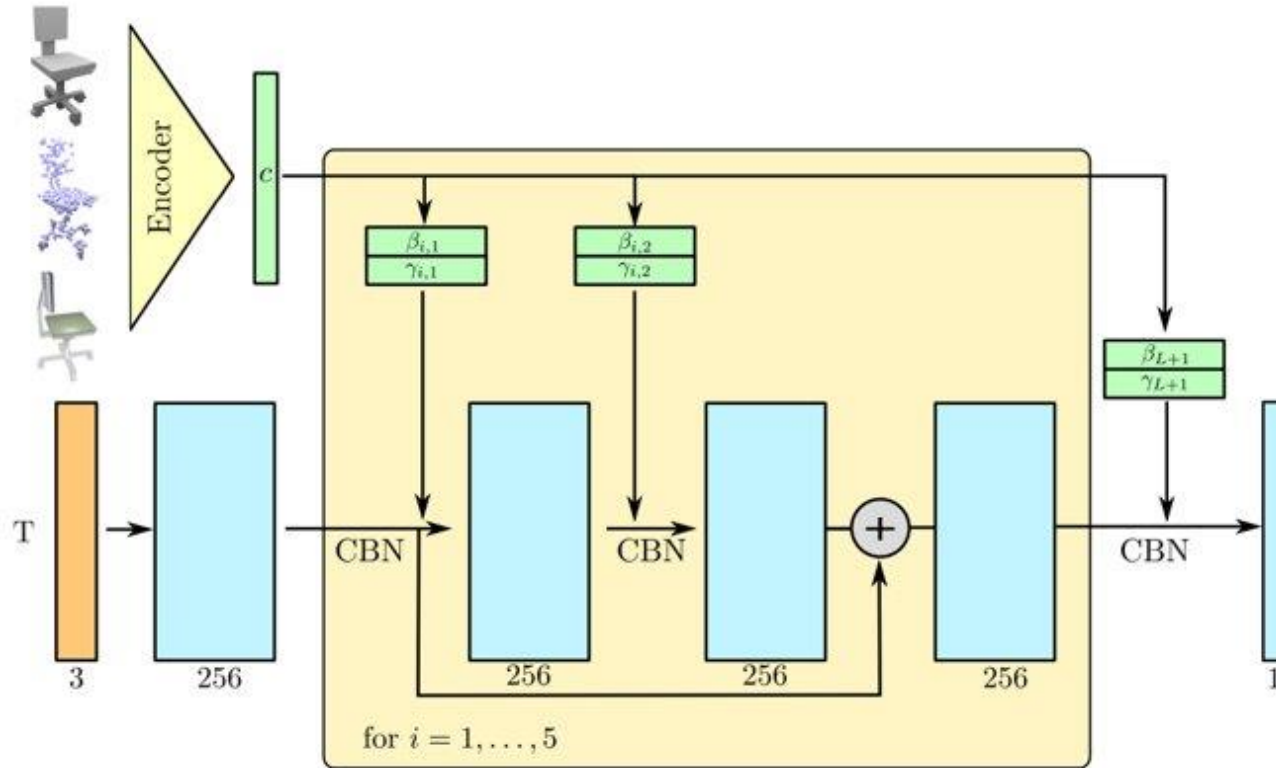


Occupancy Network



Deep Signed Distance Function

Occupancy Network



- *Input sampled 3D points and encoded features; output 0/1*

$$f_{\theta} : \mathbb{R}^3 \times \mathcal{X} \rightarrow [0, 1]$$

- *Fully connected layers + conditional batch normalization (CBN)*
- *Handling various input data*
 - *Image: Resnet*
 - *Point cloud: PointNet*
 - *Voxel: 3D CNN*

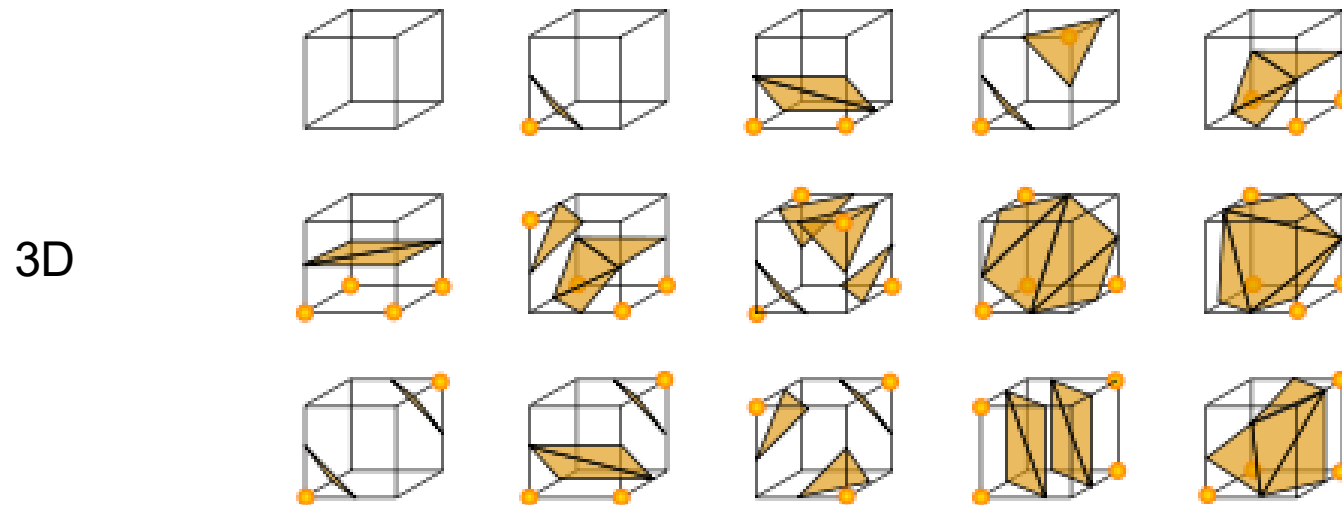
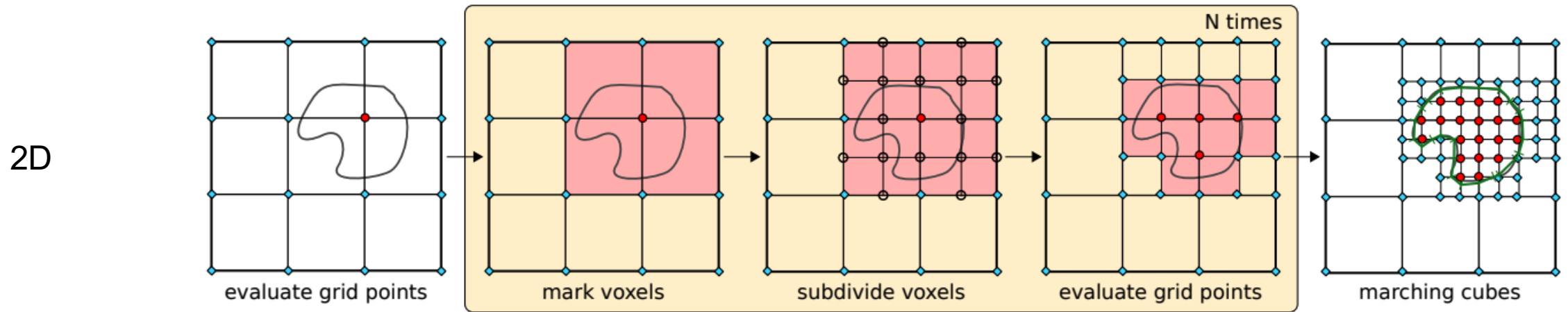
Occupancy Network

Occupancy Classification $\mathcal{L}_i(\theta)$

$$\mathcal{L}_i(\theta) = \sum_{j=1}^K BCE(f_{\theta}(p_{ij}, x_i), o_{ij})$$

- p_{ij} : location of the j^{th} point in the i^{th} sample.
- x_i : the i^{th} sample the f_{θ} conditioned.
- o_{ij} : occupancy of the j^{th} point in the i^{th} sample.

Marching cubes to extract meshes



Highlighted points: outside points

Occupancy Network

Discretize Vs continuous



16^3



32^3



64^3

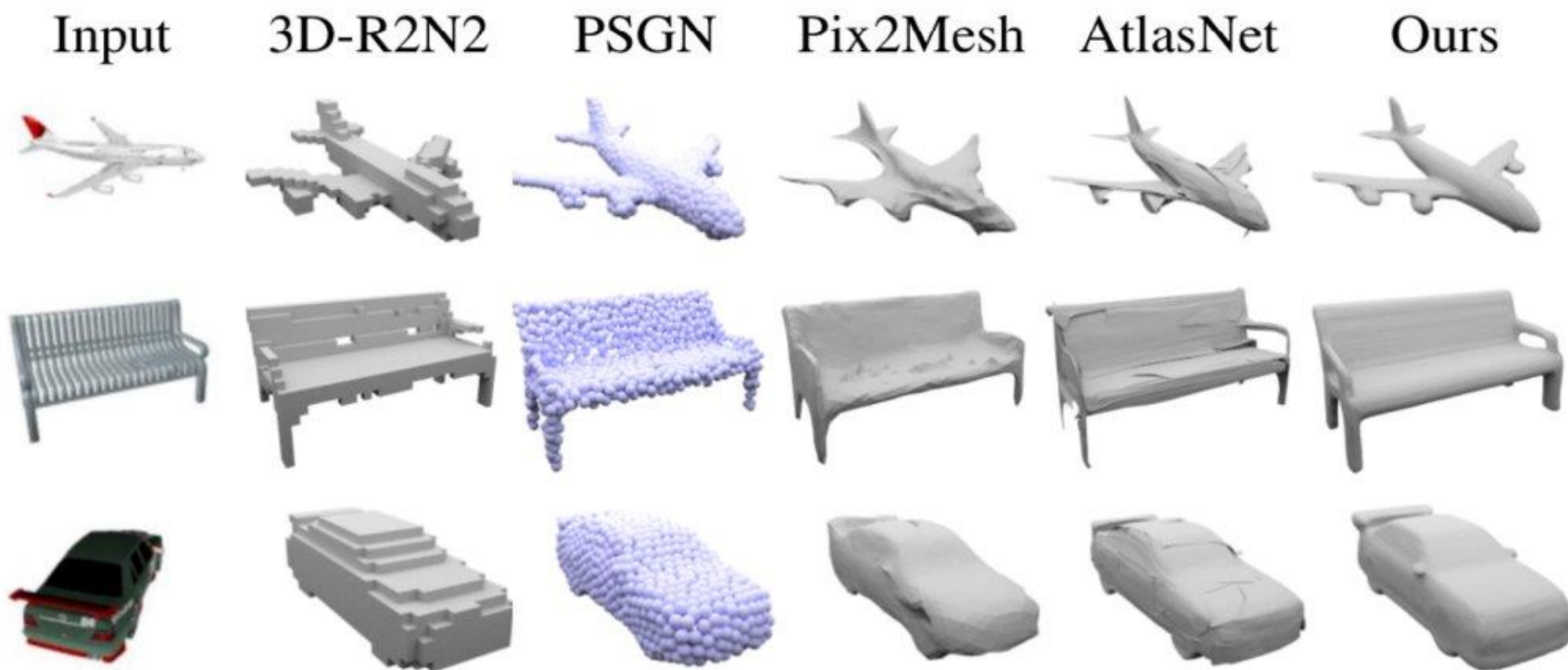


128^3



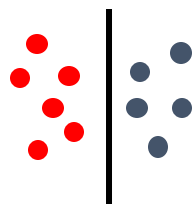
ours

Occupancy Network



DeepSDF

Occupancy

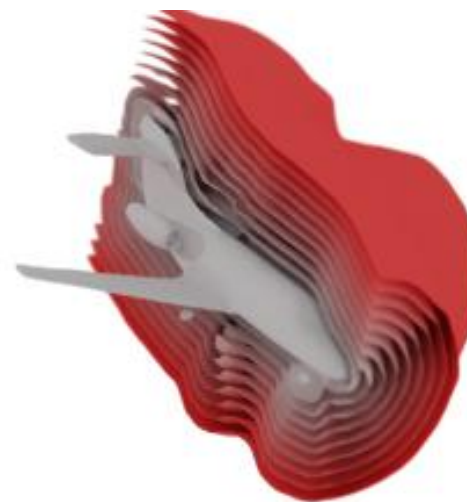
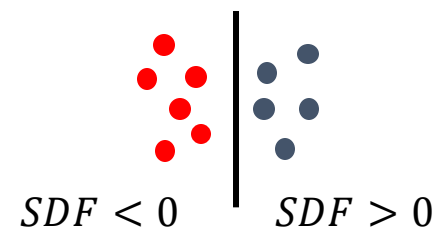


$$SDF = 0$$



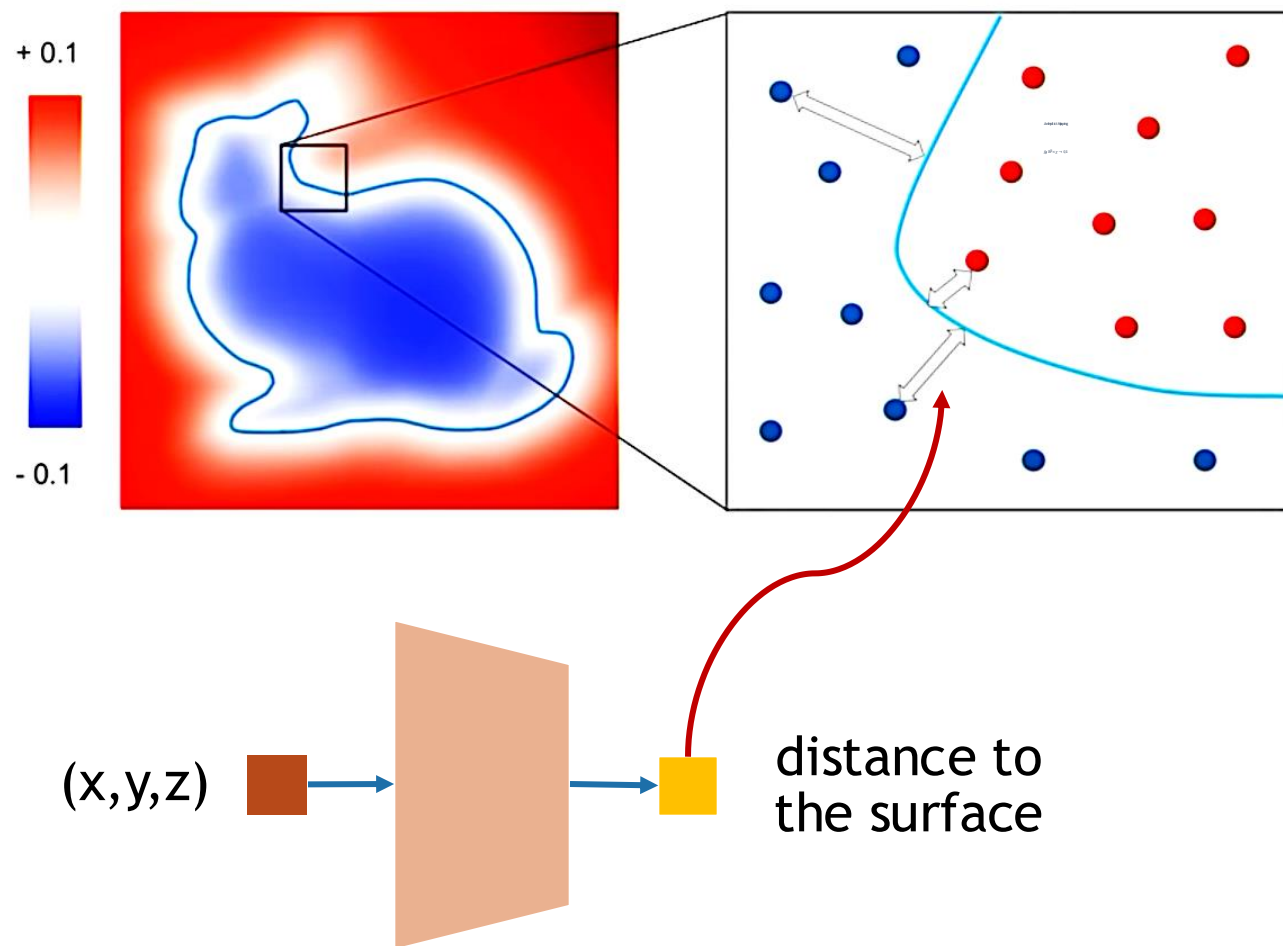
0 level-set

Sign-distance-function



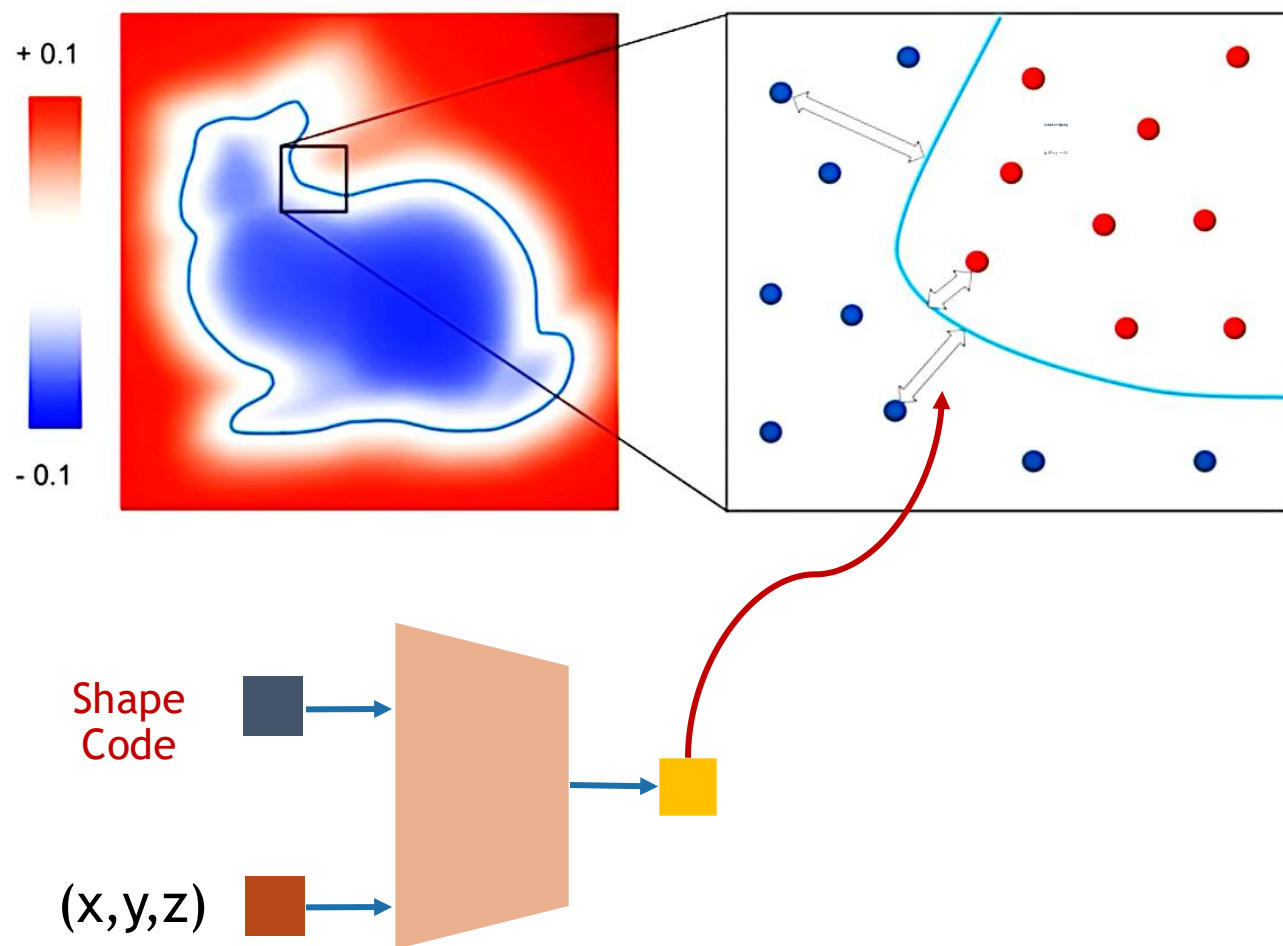
full level-set

DeepSDF



$$f_{\theta}(x) \approx SDF(x), \forall x \in \Omega$$

DeepSDF

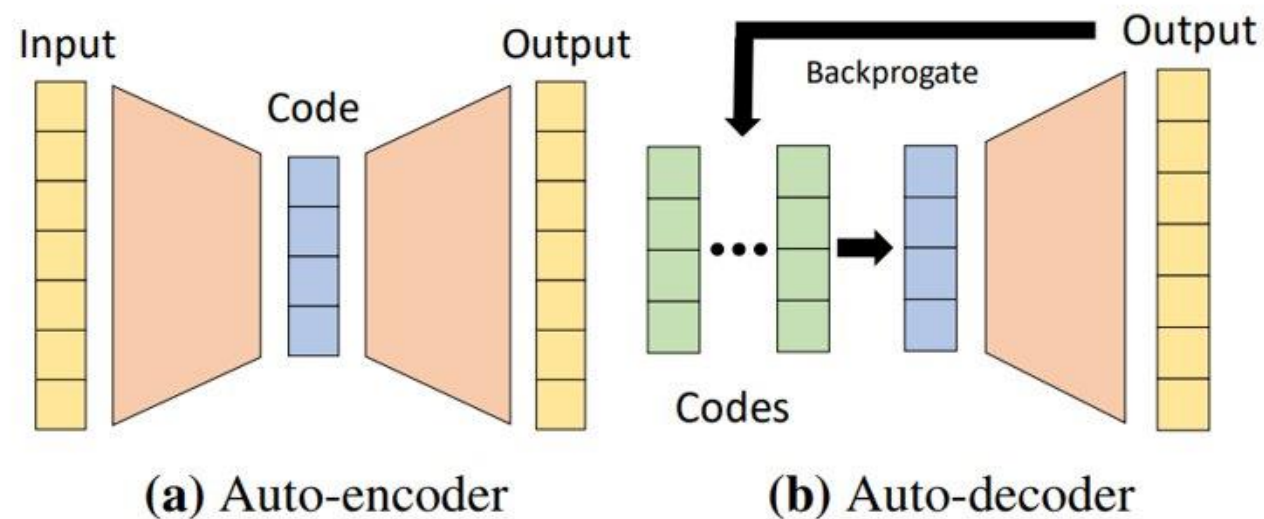


$$f_{\theta}(z_i, x) \approx SDF^i(x)$$

Function

- 8 fully-connected layers
- weight-norm

DeepSDF



- *Auto-decoder architecture*

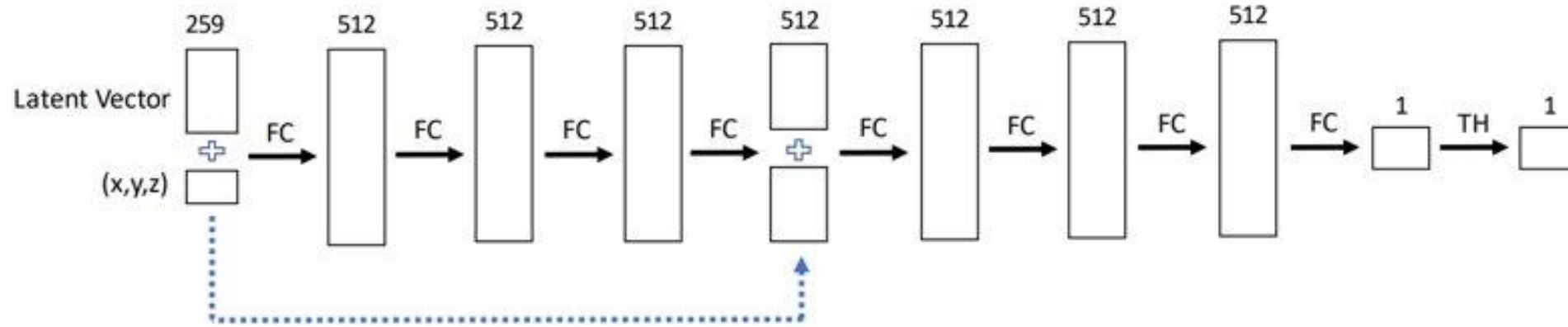
$$f_{\theta}(z_i, x) \approx SDF^i(x)$$

Shape code

Sampled point

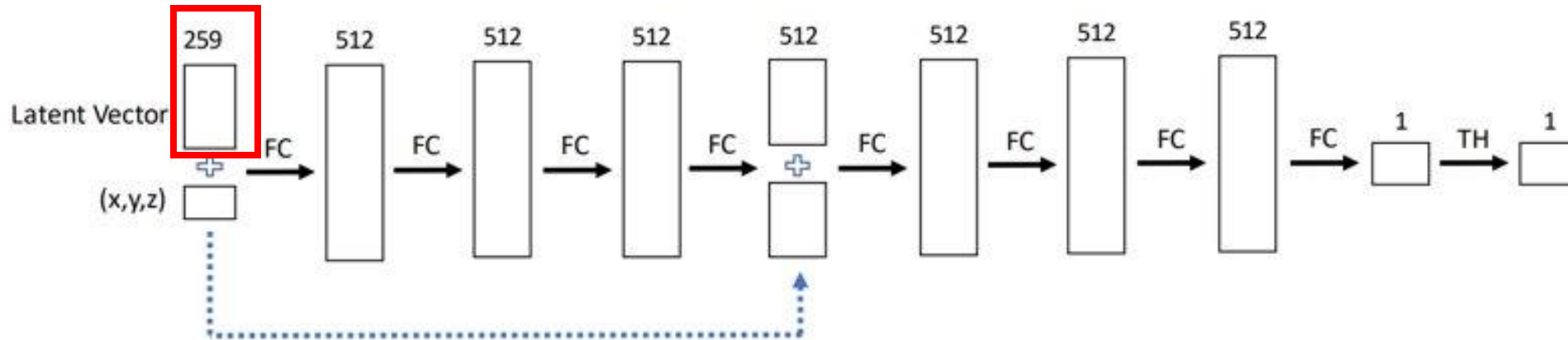
- *Training:*
 - *Jointly learning of shape codes and model parameters*
 - *Shape codes are randomly initialized*

DeepSDF



- *Fully connected layers + Weight normalization*
- *Tanh() at the output*
- *Signed Distance clipping*

DeepSDF



- *Optimization-based inference*
 - *Fix model parameters*
 - *Infer the shape code of the observation*

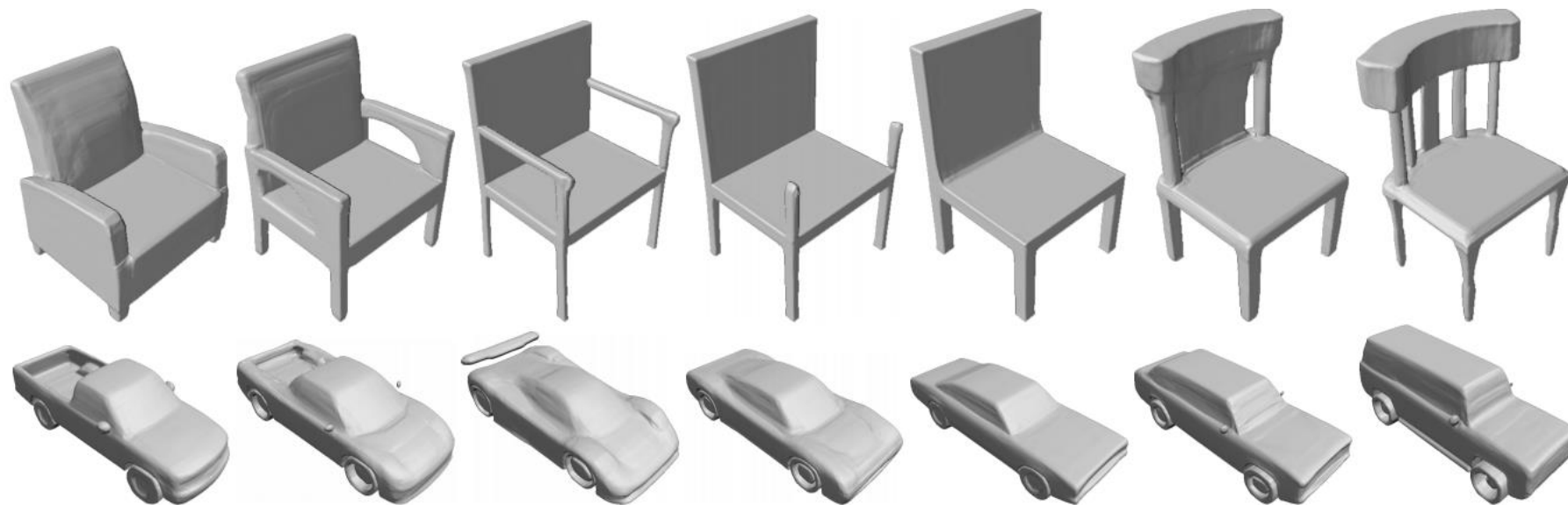
$$\arg \min_{\theta, \{z_i\}_{i=1}^N} \sum_{i=1}^N \left(\sum_{j=1}^K \mathcal{L}(f_{\theta}(z_i, \mathbf{x}_j), s_j) + \frac{1}{\sigma^2} \|z_i\|_2^2 \right)$$

Signed Distance Loss

Shape code regularization

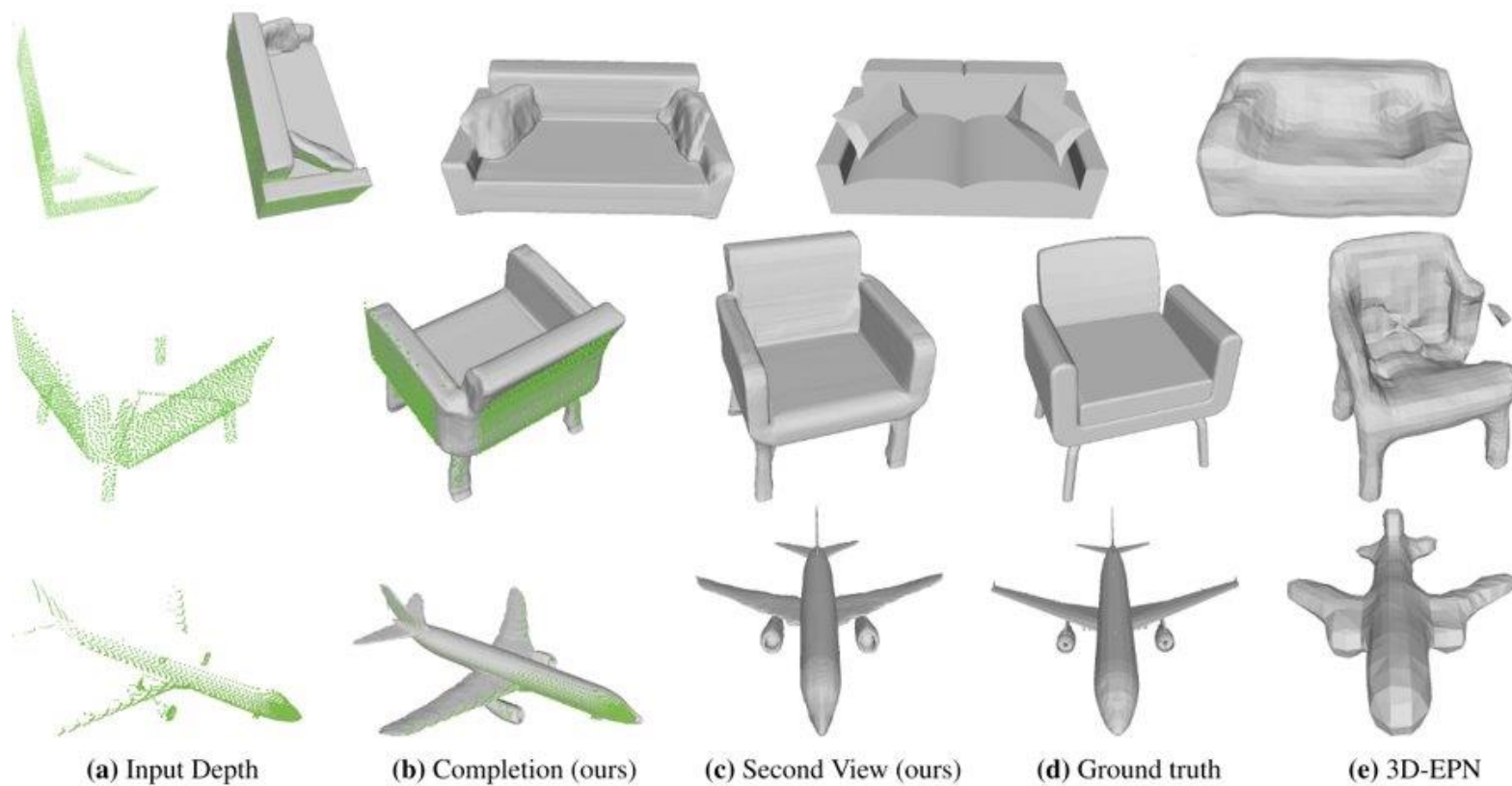
DeepSDF

Generate shapes via shape code interpolation



DeepSDF

Partial point cloud completion



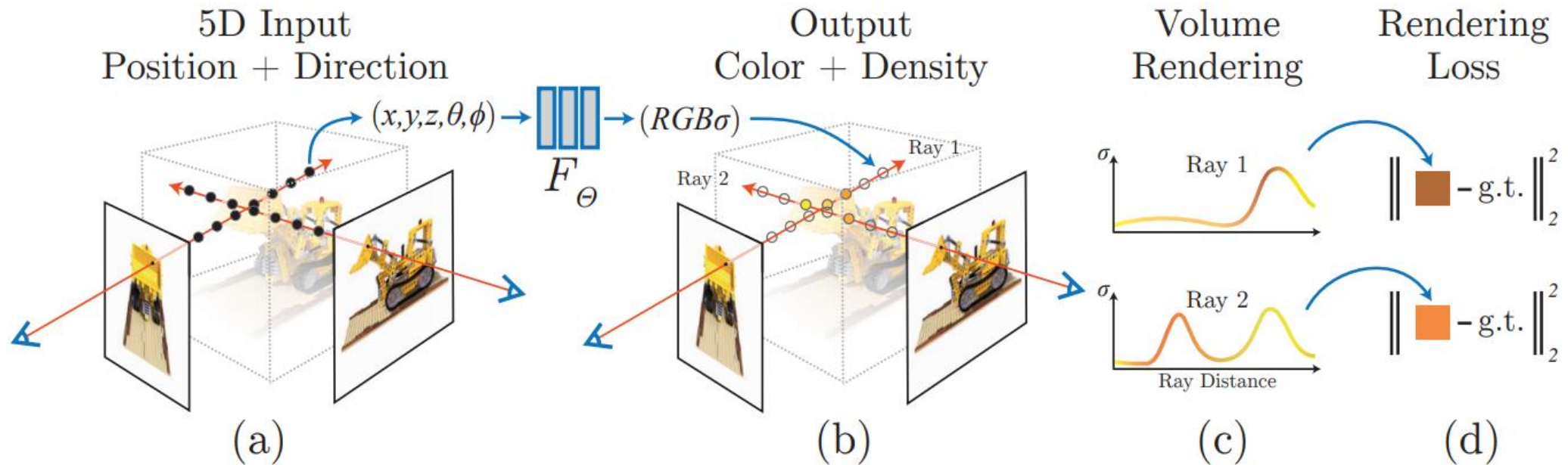
NeRF

Applied to image rendering



NeRF

1. Given sampled 3D points + camera directions, output RGB values and density
2. Image can be generated via volume rendering



Summary

- Volumetric representation
- Mesh
- Point cloud
- Implicit functions