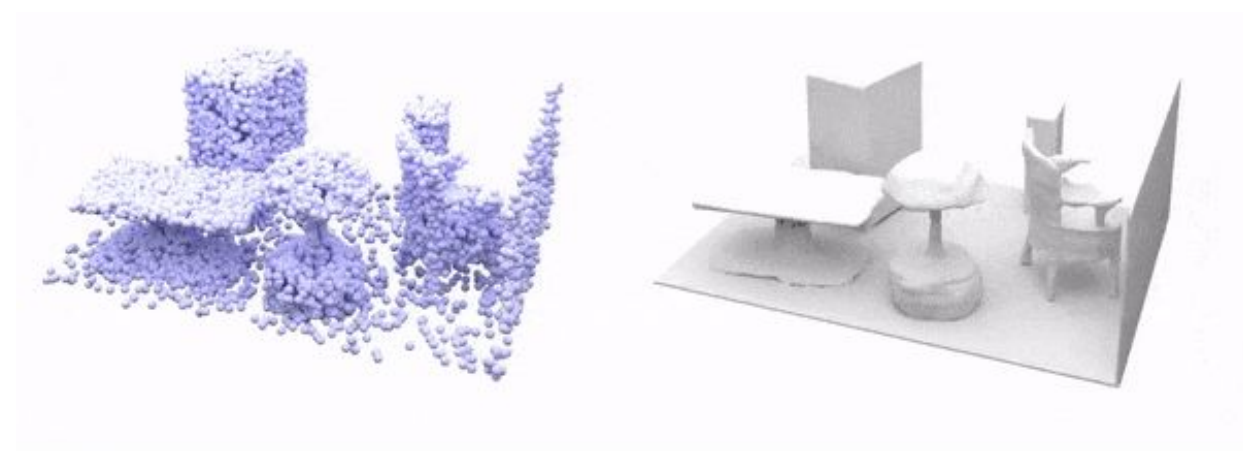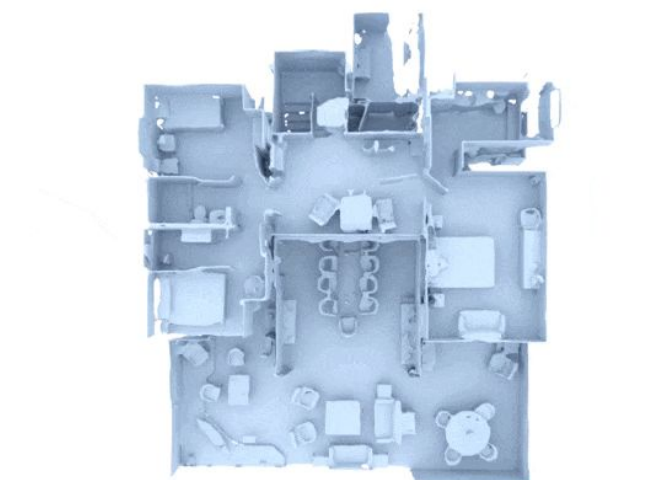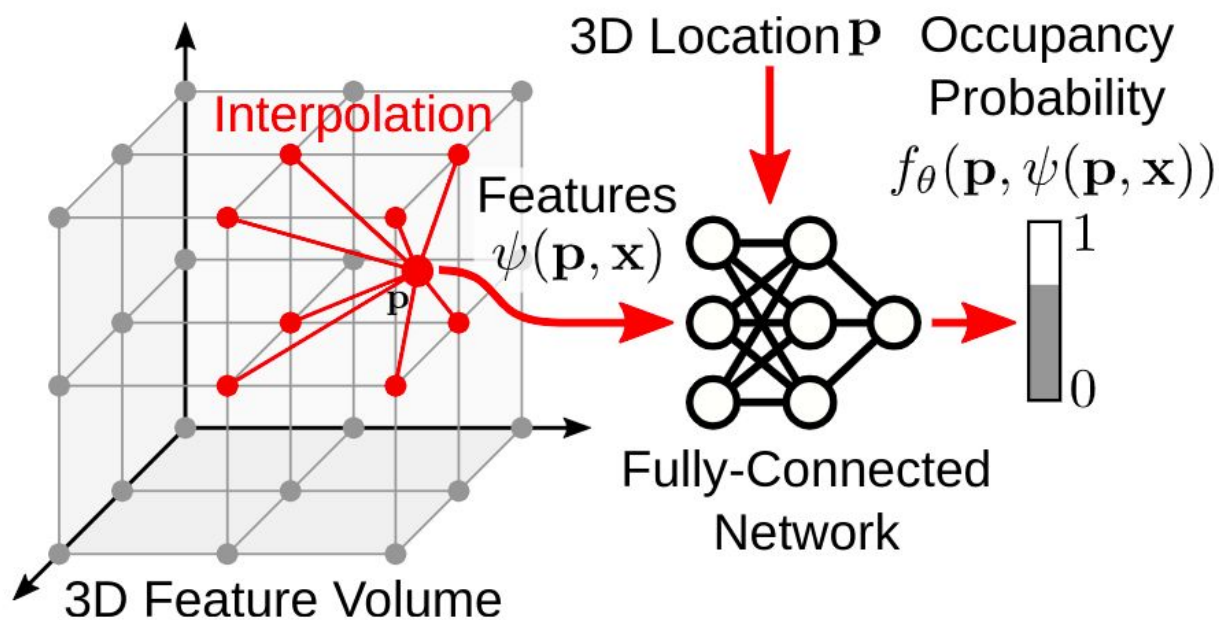# Seminar Presentation: Recent Advances in 3D Computer Vision
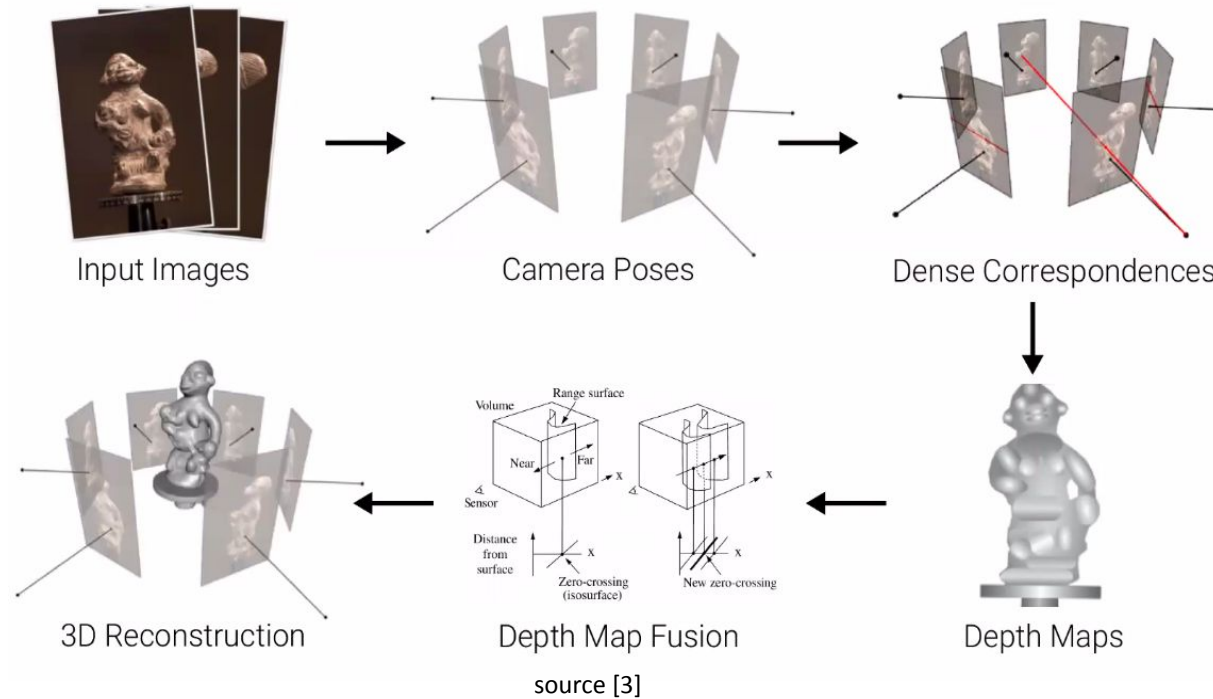
Panagiotis Petropoulakis, and Björn Häfner

Munich, October 2022

# Outline

1. Introduction of the problem
2. Approach
3. Results
4. Personal comments
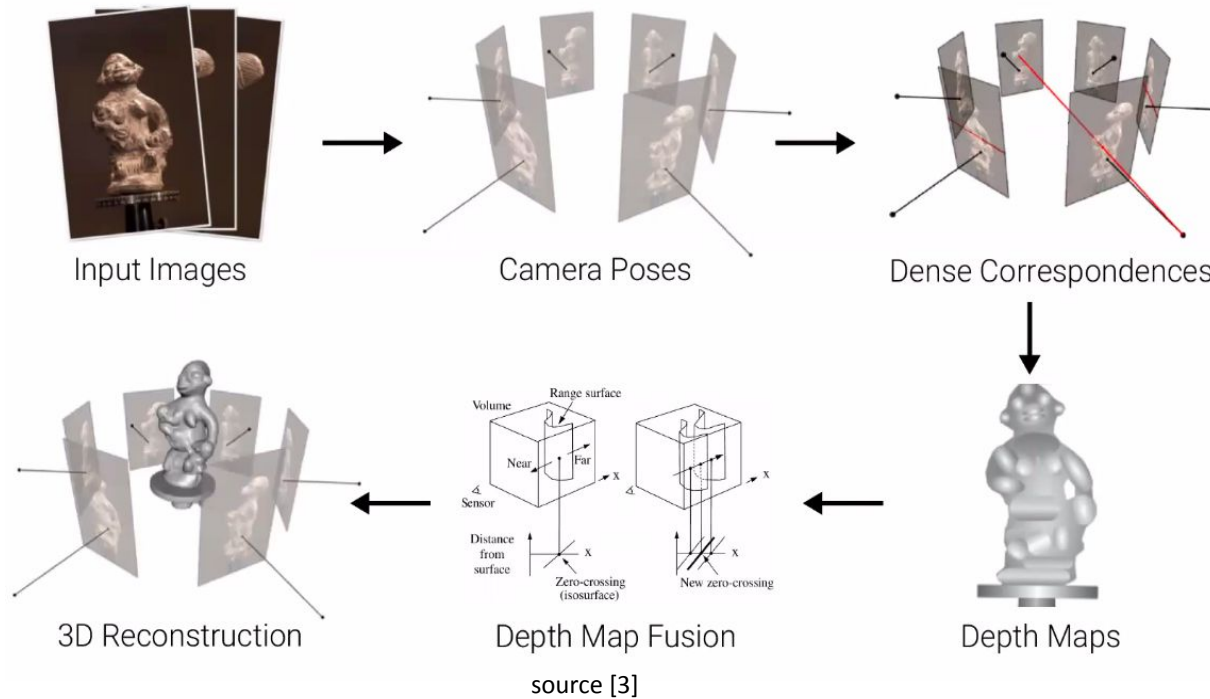5. Summary

# 1. Introduction of the problem

- **Traditional 3D reconstruction**



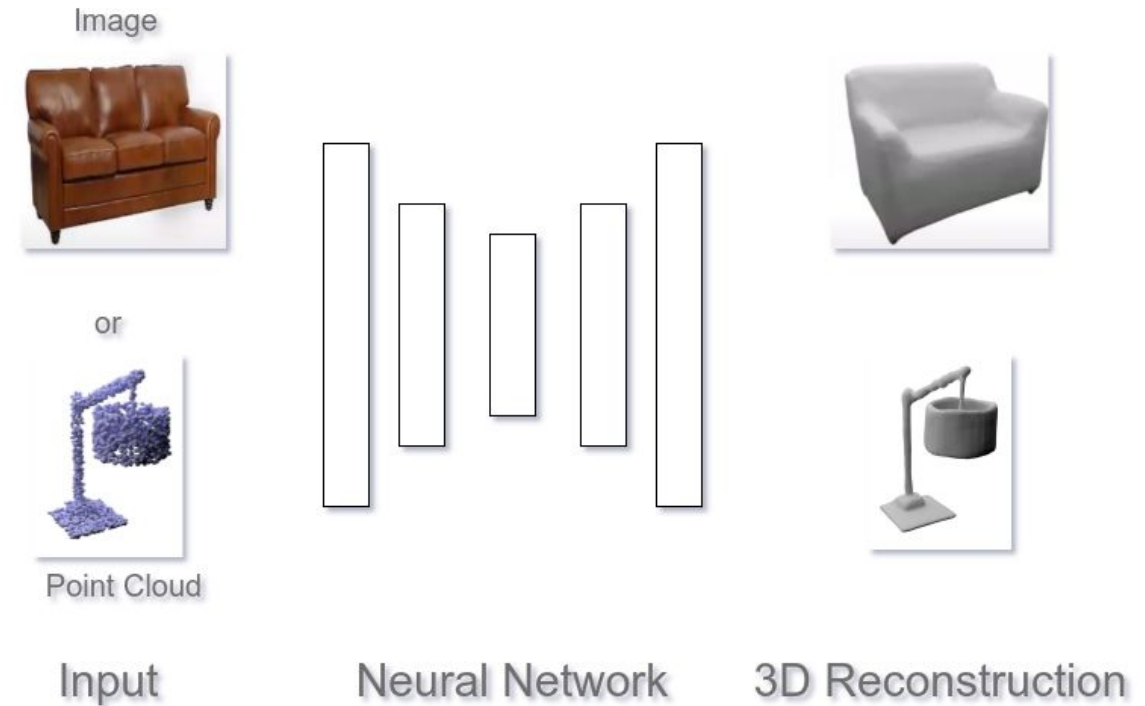source [3]

-> Multiple images are required as an input at test time

# 1. Introduction of the problem

- **Traditional 3D reconstruction**



Input Images → Camera Poses → Dense Correspondences

3D Reconstruction ← Depth Map Fusion ← Depth Maps

source [3]

-> Multiple images are required as an input at test time

- **Learning-based 3D reconstruction**



Image

or

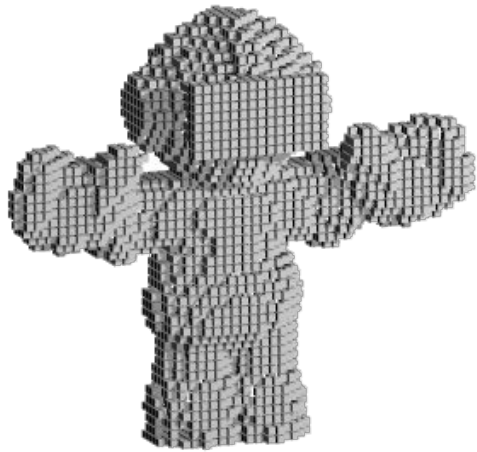Point Cloud

Input → Neural Network → 3D Reconstruction

-> Learn the 3D shape

-> 3D reconstruction from a single input

# 1. Introduction of the problem

**Common output representation of Learning-based 3D Reconstruction methods**

Voxels:
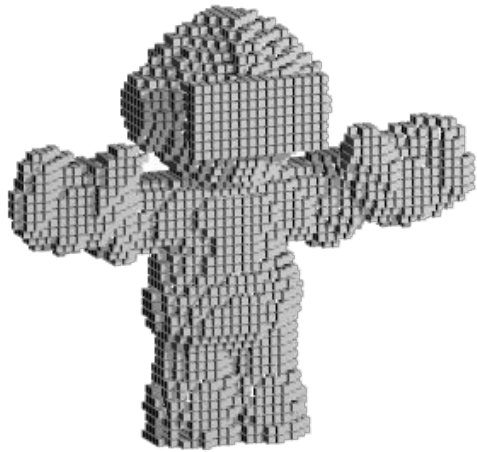


source [4]

-> Discretize into a grid

- High memory consumption

# 1. Introduction of the problem

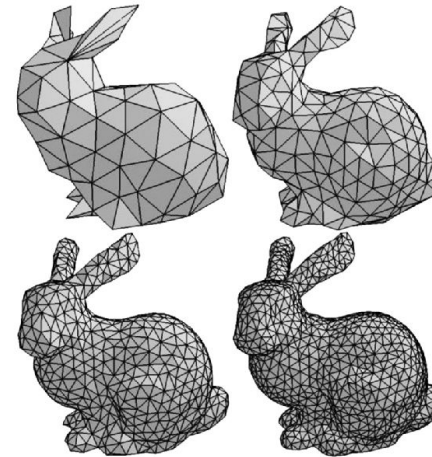**Common output representation of Learning-based 3D Reconstruction methods**

Voxels:



source [4]

-> Discretize into a grid
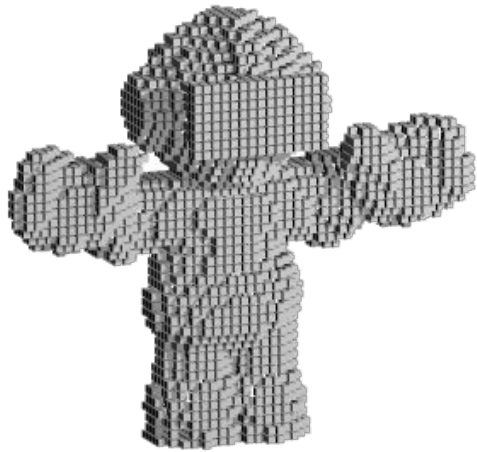
- High memory consumption

Meshes:



source [5]

-> Regress into vertices & faces

- Non-watertight reconstructions
- Often require deforming
  a template

# 1. Introduction of the problem

**Common output representation of Learning-based 3D Reconstruction methods**
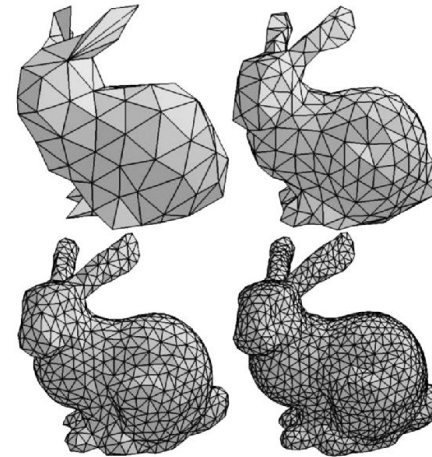
<u>Voxels:</u>

-> Discretize into a grid

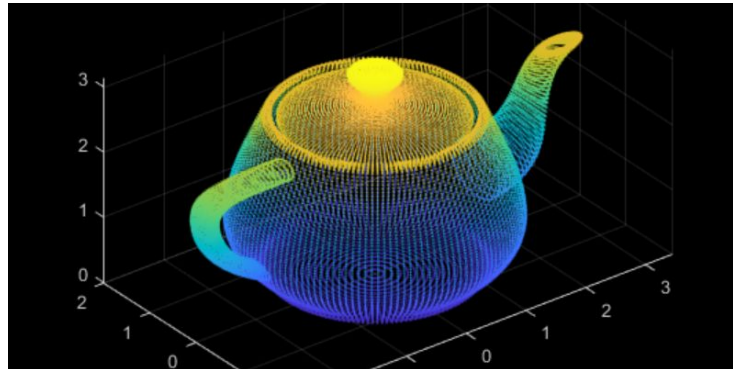  - High memory consumption



source [4]

<u>Meshes:</u>



-> Regress into vertices & faces

  - Non-watertight reconstructions

  - Often require deforming
    a template

source [5]

<u>Point Clouds:</u>



-> Predict the coordinates of 3D points

  - Limited number of points

  - Topological relations are lost

source [6]

# 1. Introduction of the problem

## Neural Implicit Representation

- No discretization of the 3D space
- No topological restrictions
- Independent of the camera viewpoint



$$f_\theta(p) = \tau$$

-> Represent the 3D shape implicitly

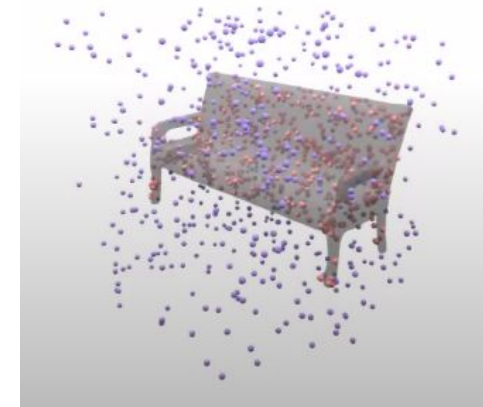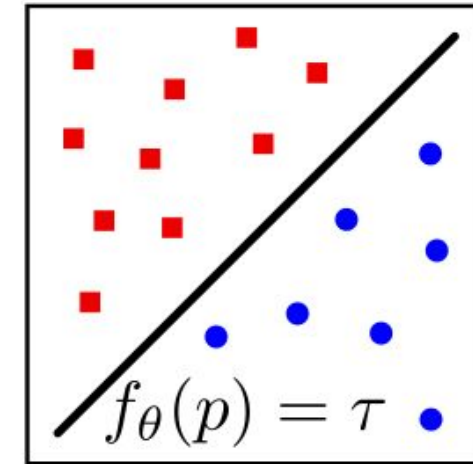-> Surface <=> Decision boundary of a non-linear classifier



source [2, 3]

# 1. Introduction of the problem

Neural Implicit Representation

- No discretization of the 3D space
- No topological restrictions
- Independent of the camera viewpoint

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

3D point       Condition       Occupancy Probability

-> Represent the 3D shape implicitly

-> Surface <=> Decision boundary of a non-linear classifier



$$f_\theta(p) = \tau$$

source [2, 3]

# 1. Introduction of the problem

**Problems with previous works**

-> Occupancy Networks

# 1. Introduction of the problem

**Problems with previous works**

-> Occupancy Networks



- Local details are not preserved
  - Overly smooth reconstruction

- No Translation Equivariance

- Mainly for simple objects

# 2. Approach

## 2D Method

Encoder

# 2. Approach

## 2D Method

Encoder



Input x  Point Cloud

PointNet Encoder

2D Feature Plane

1. Refine features

    -> 2D PointNet

    +Preserves local information

# 2. Approach

## 2D Method

Encoder



Input **x** — Point Cloud — PointNet Encoder — 2D Feature Plane

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
   +Not depend on a global frame

# 2. Approach

## 2D Method

### Encoder

Input **x**    Point Cloud



PointNet Encoder

2D Feature Plane

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
   +Not depend on a global frame

### Decoder



2D U-Net

Features
$\psi(p, x)$

**p**

# 2. Approach

**2D Method**

Encoder

Decoder

Feature space from the encoder

Input x  Point Cloud

PointNet Encoder

2D Feature Plane

2D U-Net

Features $\psi(p, x)$

p

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
   +Not depend on a global frame

# 2. Approach

## 2D Method

### Encoder

Input x    Point Cloud

PointNet Encoder

2D Feature Plane

### Decoder

Feature space from the encoder

New feature space with global information
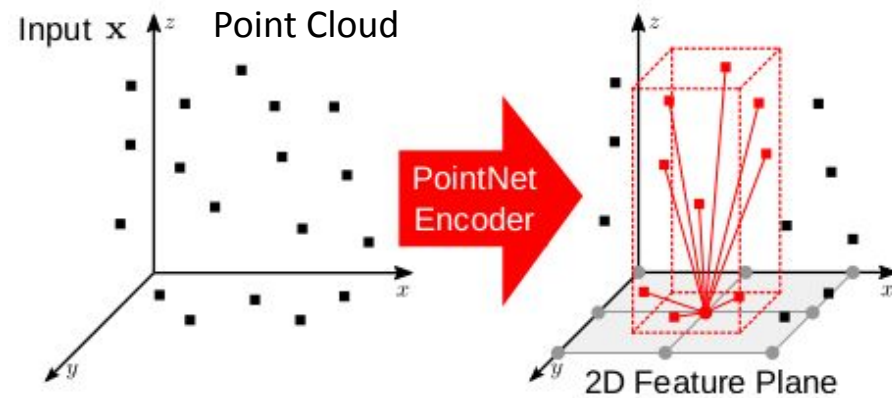
2D U-Net

Features
$\psi(p, x)$

p

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
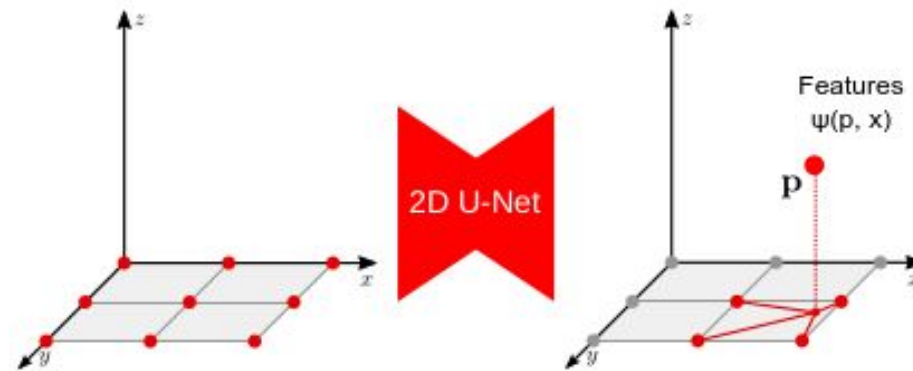   +Not depend on a global frame

# 2. Approach

## 2D Method

### Encoder

Input x    Point Cloud



PointNet
Encoder

2D Feature Plane

### Decoder

Feature space from the encoder

New feature space with global information



2D U-Net

Features
ψ(p, x)

p

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
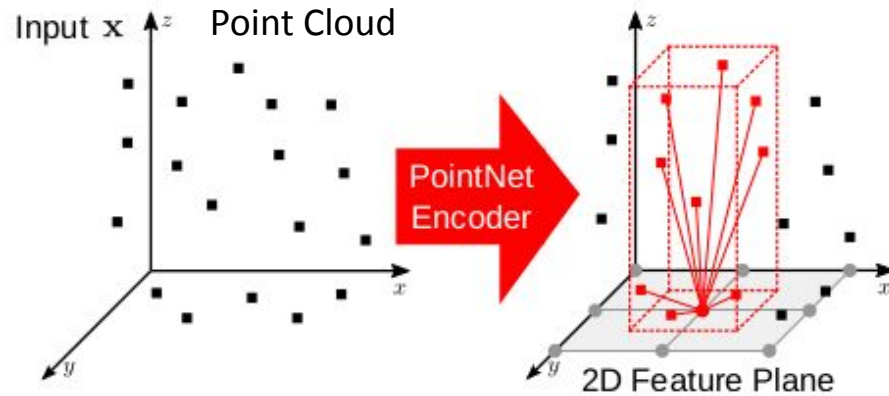   +Preserves local information
   +Not depend on a global frame

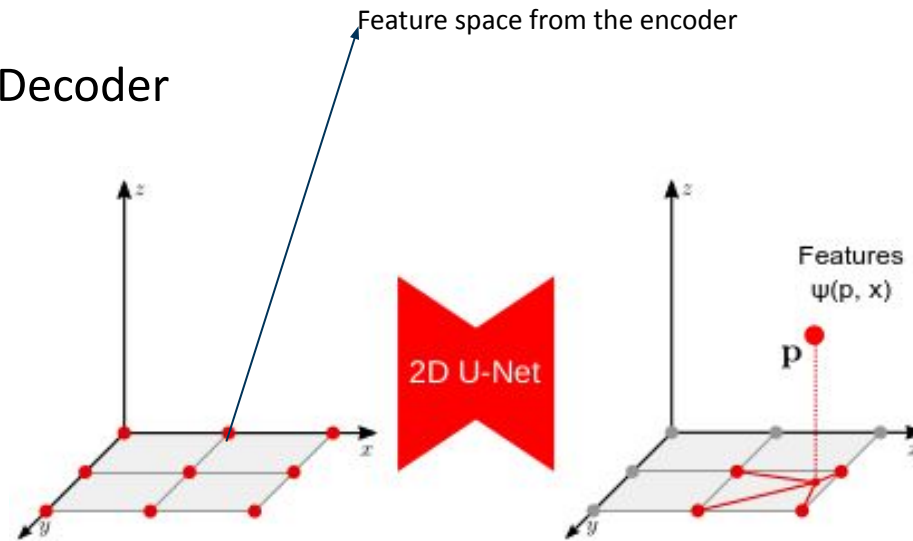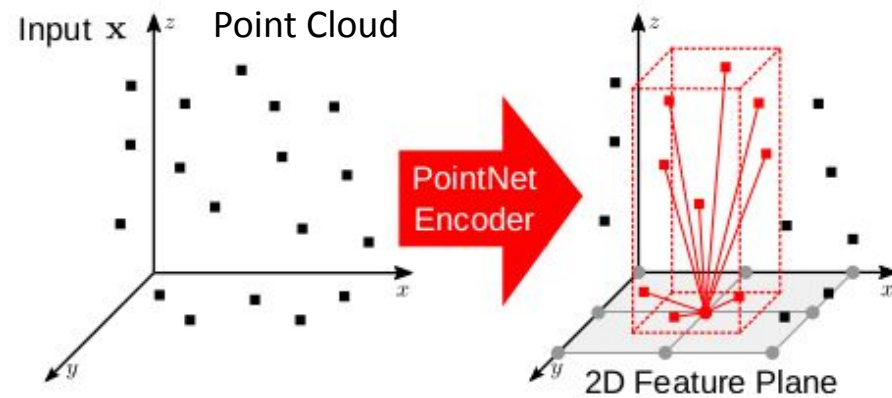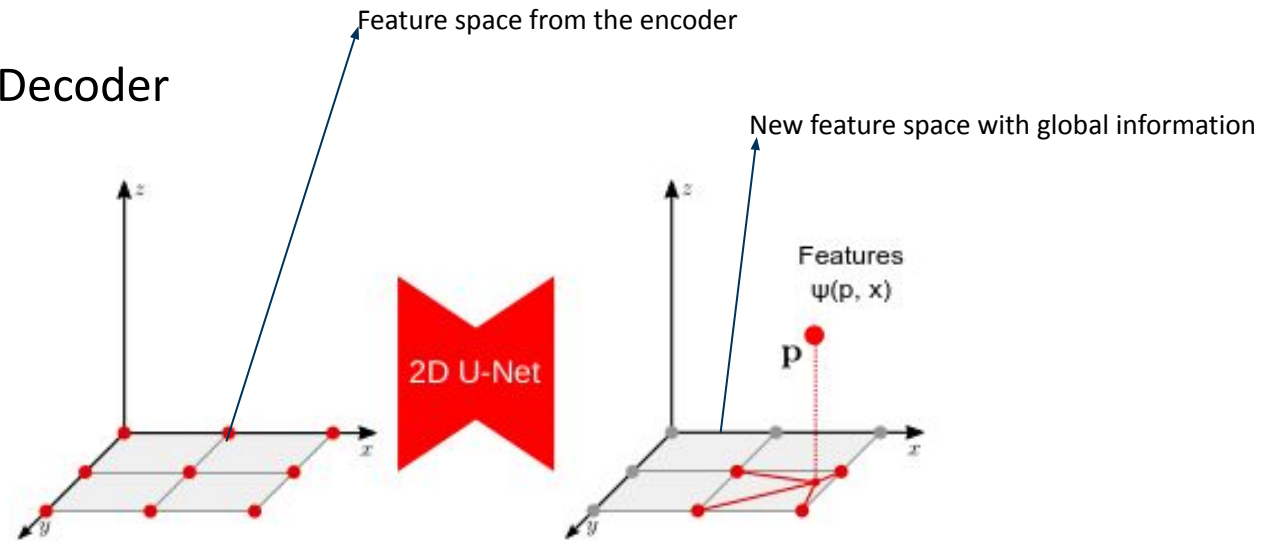1. Process the Feature Plane (space)
   -> 2D U-Net
   +Integrate global information
   +Translation equivariance

# 2. Approach

## 2D Method

Encoder



Input **x**   Point Cloud

PointNet Encoder

2D Feature Plane

Decoder

or Query Point

Occupancy Prediction

of point p

3D Location **p**

2D U-Net

Features $\psi(\mathbf{p}, \mathbf{x})$

**p**

Occupancy Probability

$f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

1

0

Occupancy Network

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
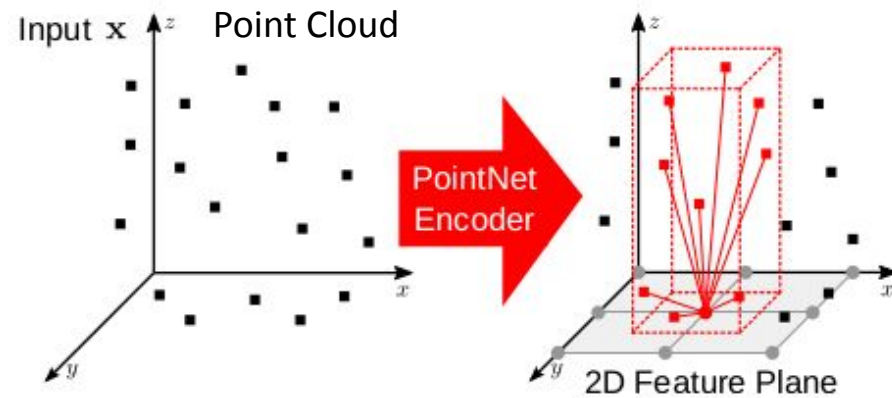   +Not depend on a global frame

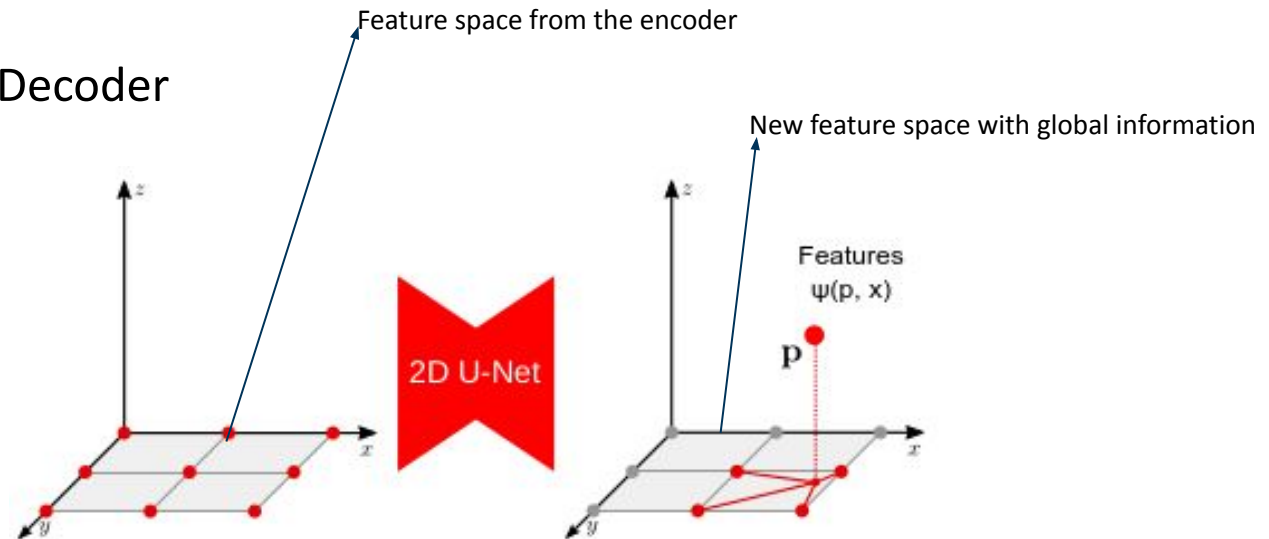1. Process the Feature Plane (space)
   -> 2D U-Net
   +Integrate global information
   +Translation equivariance

# 2. Approach

**2D Method**

Encoder

Input **x** $z$    Point Cloud



PointNet Encoder

2D Feature Plane

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
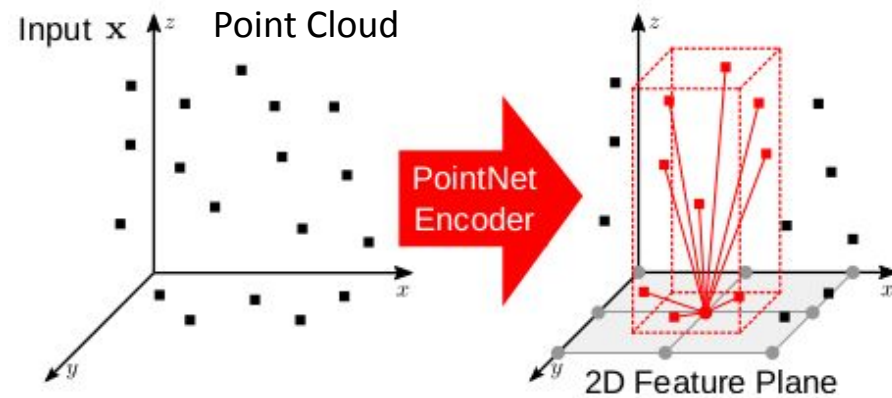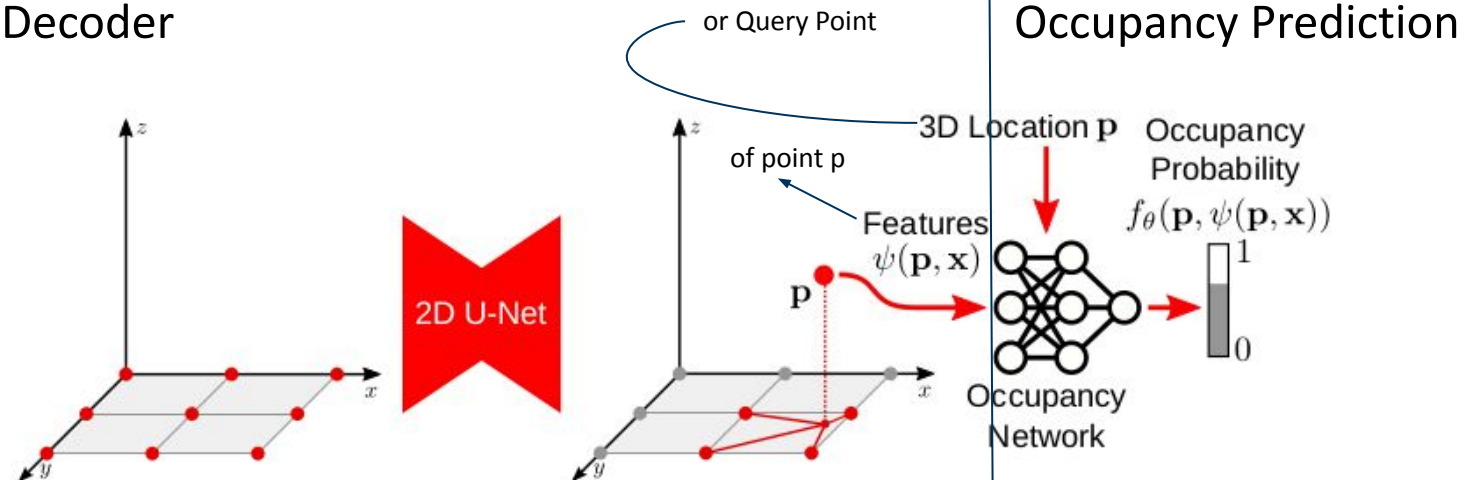   +Preserves local information
   +Not depend on a global frame

Decoder



2D U-Net

or Query Point

of point p

3D Location **p**

Features $\psi(\mathbf{p}, \mathbf{x})$

**p**

Occupancy Prediction

Occupancy Probability
$f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Occupancy Network

1. Process the Feature Plane (space)
   -> 2D U-Net
   +Integrate global information
   +Translation equivariance

1. Query a 3D point
   -> Use interpolation
   -> Predict the Occupancy Prob.

# 2. Approach

## 2D Method

Encoder

Input **x**   Point Cloud



PointNet
Encoder

2D Feature Plane

1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
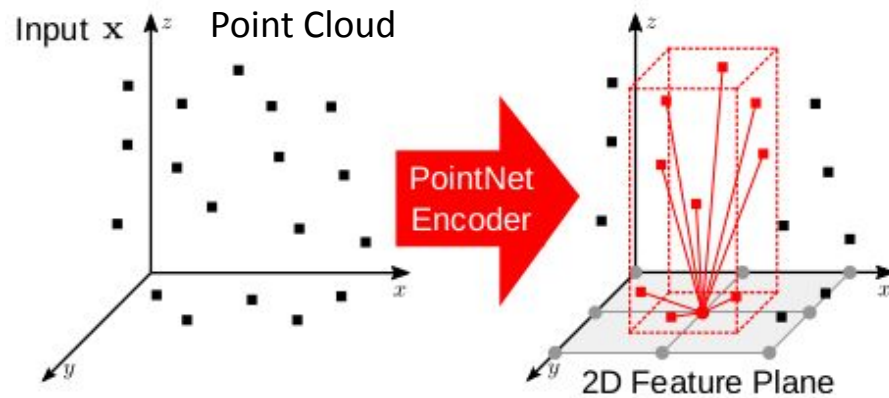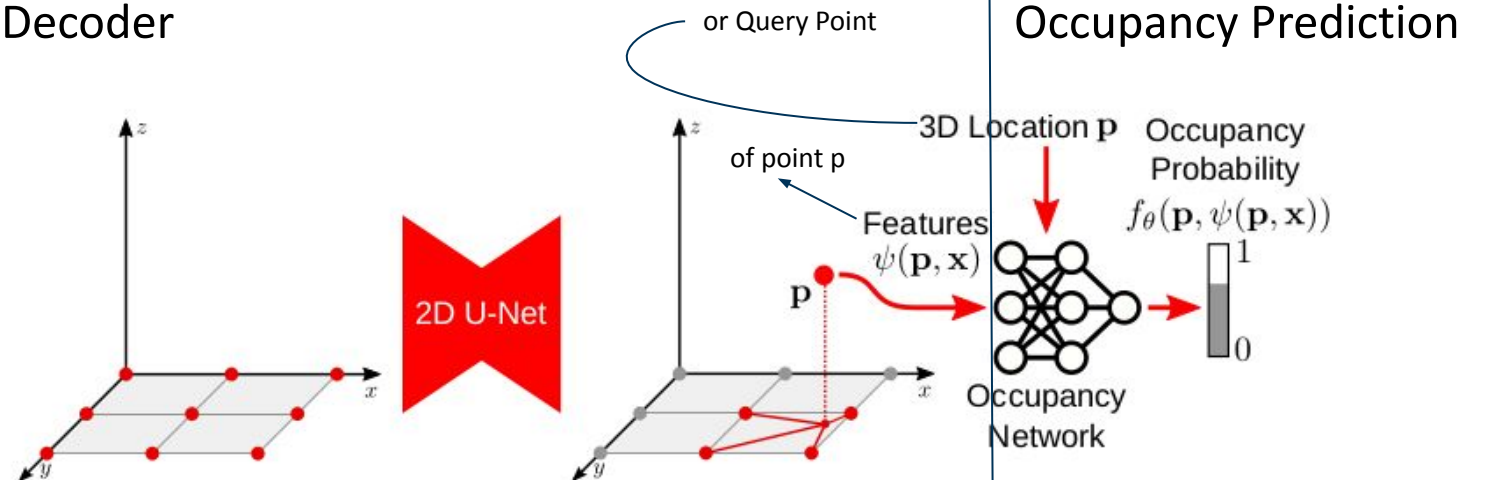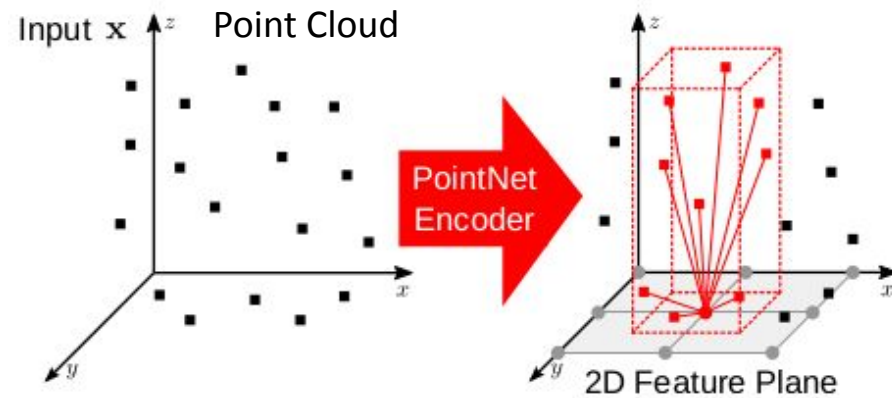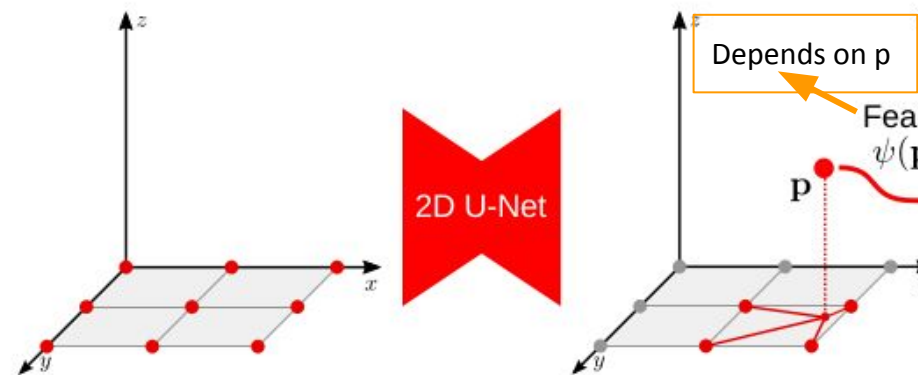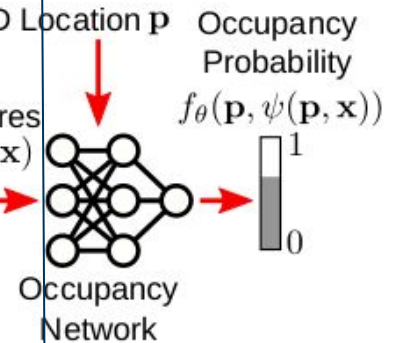   +Preserves local information
   +Not depend on a global frame

Decoder



2D U-Net

1. Process the Feature Plane (space)
   -> 2D U-Net
   +Integrate global information
   +Translation equivariance

Occupancy Prediction

Depends on p
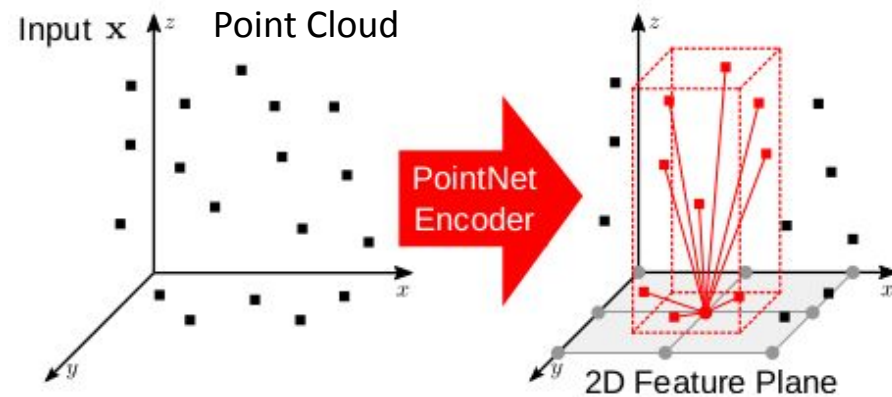
3D Location **p**   Occupancy Probability

$f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Features
$\psi(\mathbf{p}, \mathbf{x})$

**p**

Occupancy Network

1. Query a 3D point
   -> Use interpolation
   -> Predict the Occupancy Prob.

# 2. Approach

**2D Method**

Multiple Canonical Planes

Encoder

Decoder

Occupancy Prediction



Input $\mathbf{x}$    Point Cloud

PointNet Encoder

2D Feature Plane

3D Location $\mathbf{p}$   Occupancy Probability

$f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Features $\psi(\mathbf{p}, \mathbf{x})$

2D U-Net

Bilinear Interpolation

2D Feature Planes

Occupancy Network
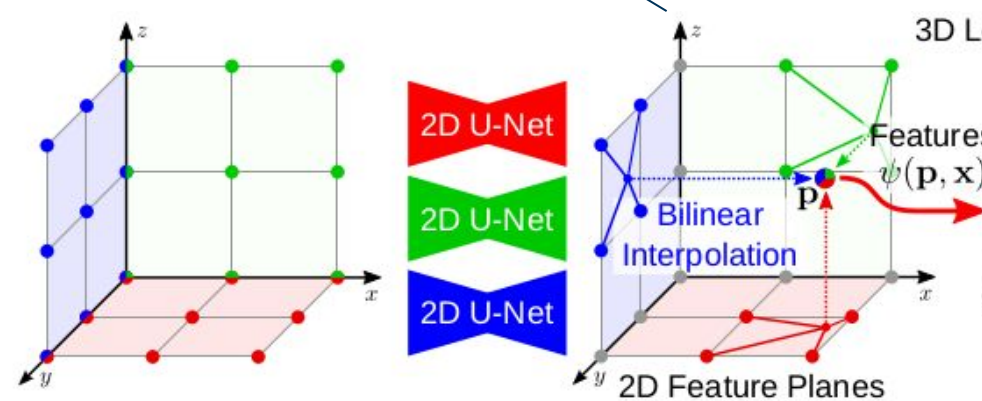
1. Refine features
   -> 2D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
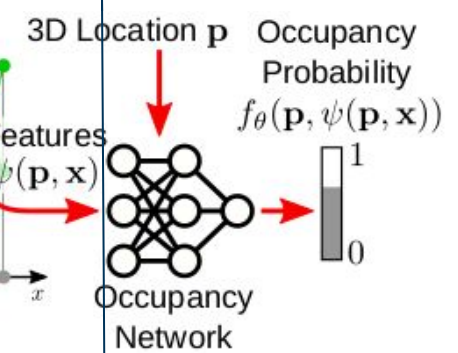   +Preserves local information
   +Not depend on a global frame

1. Process the Feature Plane (space)
   -> 2D U-Net
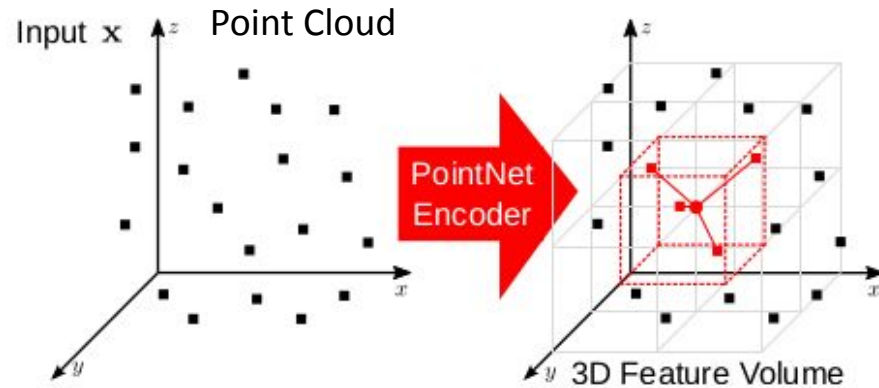   +Integrate global information
   +Translation equivariance

1. Query a 3D point
   -> Use interpolation
   -> Predict the Occupancy Prob.

# 2. Approach

**3D Method - Volumetric Repr.**

Encoder



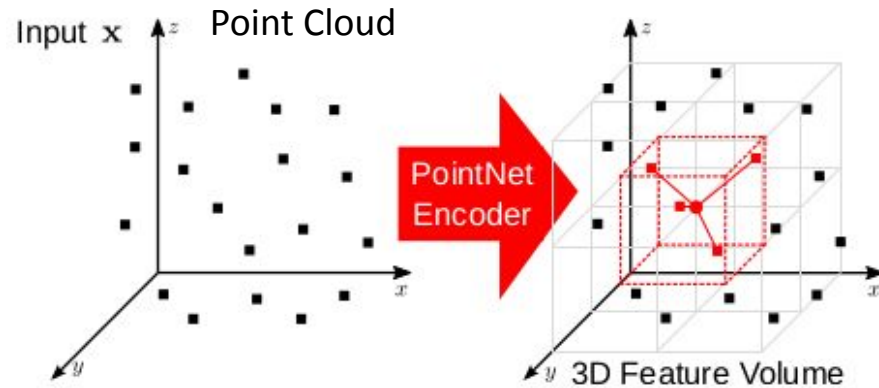Input x    Point Cloud

PointNet
Encoder

3D Feature Volume

1. Refine features
   -> 3D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
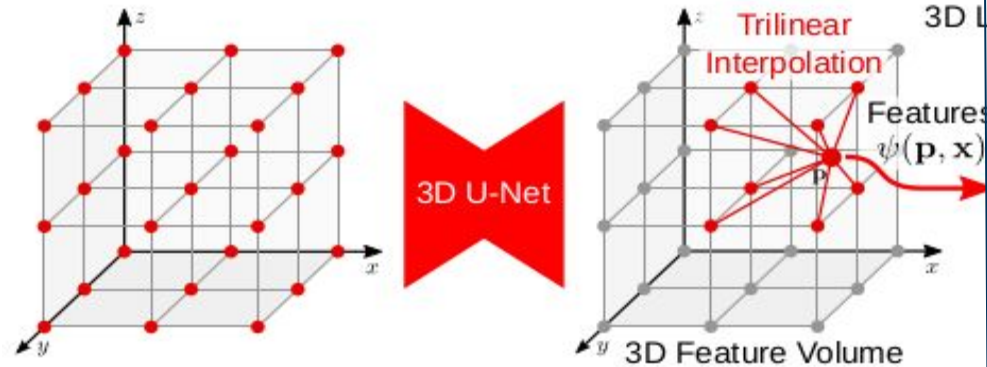   +Not depend on a global frame

# 2. Approach

## 3D Method - Volumetric Repr.

Encoder



Input x  Point Cloud  PointNet Encoder  3D Feature Volume

Decoder



Trilinear Interpolation  3D U-Net  3D Feature Volume  Features $\psi(\mathbf{p}, \mathbf{x})$

Occupancy Prediction

3D Location $\mathbf{p}$  Occupancy Probability  $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$
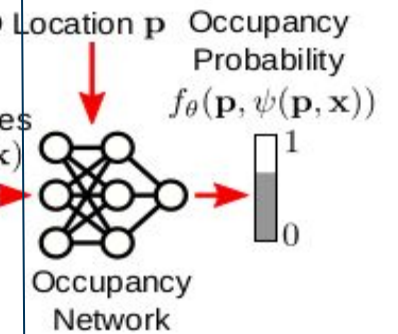
Occupancy Network

1. Refine features
   -> 3D PointNet
   +Preserves local information

2. Project to canonical plane
   -> Aggregate local neighbors
   +Preserves local information
   +Not depend on a global frame

1. Process the Feature Plane (space)
   -> 3D U-Net
   +Integrate global information
   +Translation equivariance

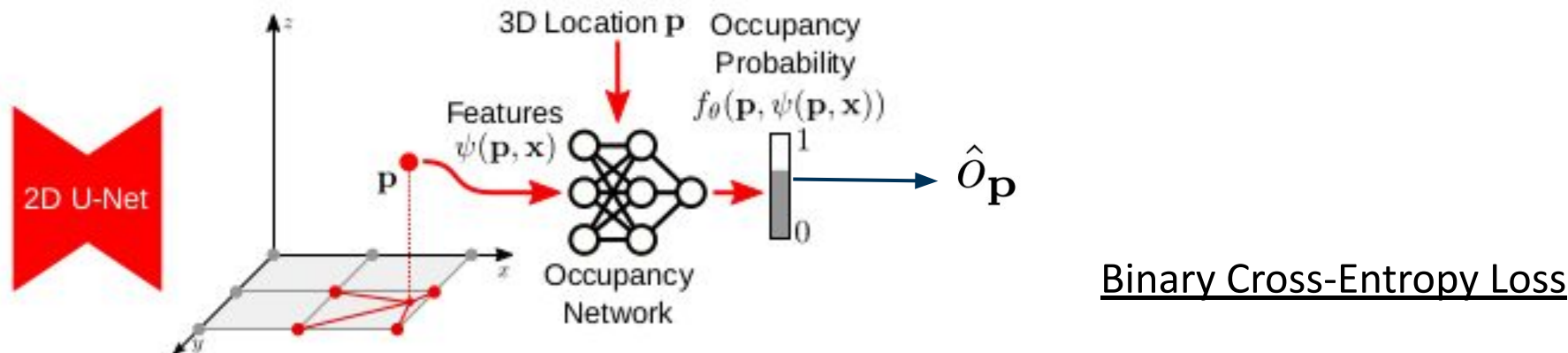1. Query a 3D point
   -> Use interpolation
   -> Predict the
   Occupancy Prob.

# 2. Approach

**Training**

- Train the Occupancy network
  - Sample query points p from 3D objects using the train set



Binary Cross-Entropy Loss

$$\mathcal{L}(\hat{o}_\mathbf{p}, o_\mathbf{p}) = -[o_\mathbf{p} \cdot \log(\hat{o}_\mathbf{p}) + (1 - o_\mathbf{p}) \cdot \log(1 - \hat{o}_\mathbf{p})]$$

True occupancy prob.          Predicted occupancy prob.

# 2. Approach

**Training**

- Train the Occupancy network
  - Sample query points p from 3D objects using the train set
  - The Encoder is pre-trained / task-specific: classf. & segm.
    - feature space is ready to use (ψ)



Binary Cross-Entropy Loss

$$\mathcal{L}(\hat{o}_{\mathbf{p}}, o_{\mathbf{p}}) = -[o_{\mathbf{p}} \cdot \log(\hat{o}_{\mathbf{p}}) + (1 - o_{\mathbf{p}}) \cdot \log(1 - \hat{o}_{\mathbf{p}})]$$
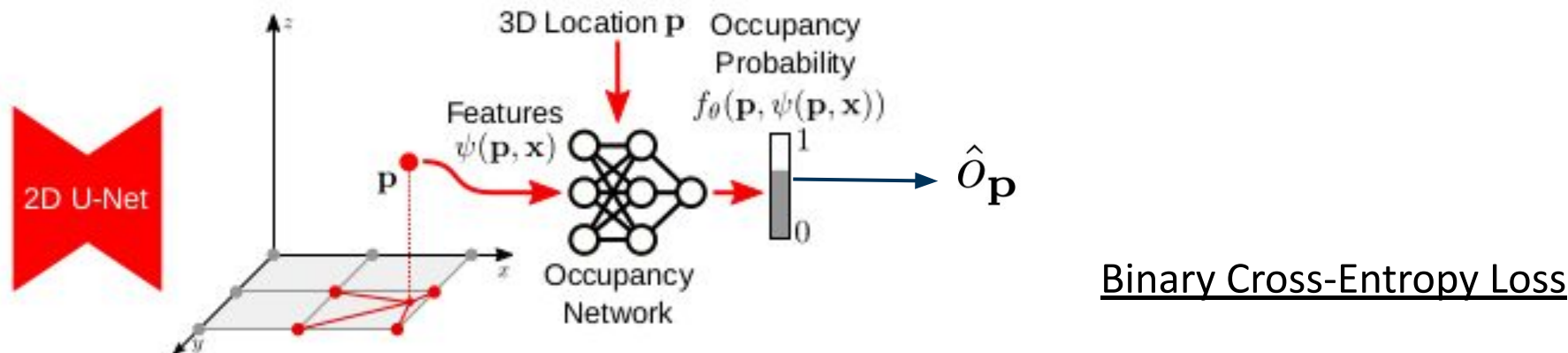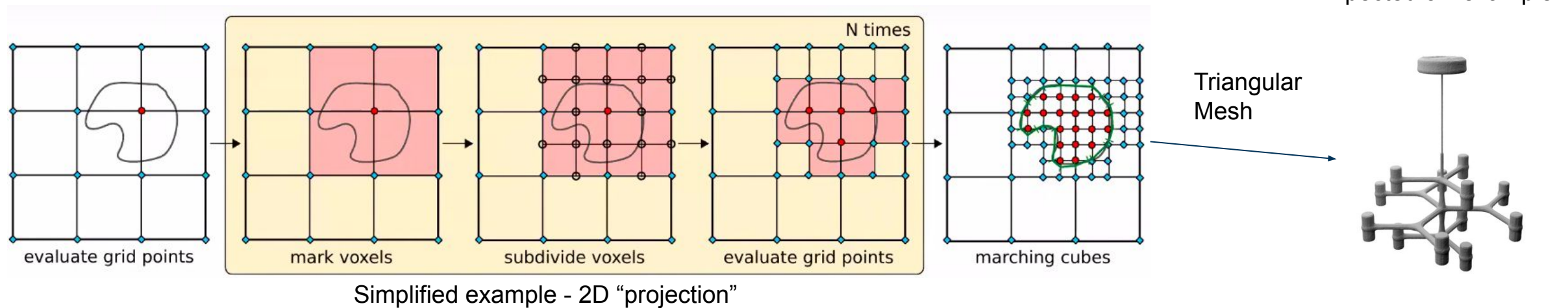
True occupancy prob.          Predicted occupancy prob.

# 2. Approach

**Rendering - Generate a 3D Mesh**

Multiresolution IsoSurface Extraction (MISE)

Expected 3D example



N times

evaluate grid points | mark voxels | subdivide voxels | evaluate grid points | marching cubes

Triangular Mesh

Simplified example - 2D "projection"

1. Partition the 3D space
   - Build octree incrementally

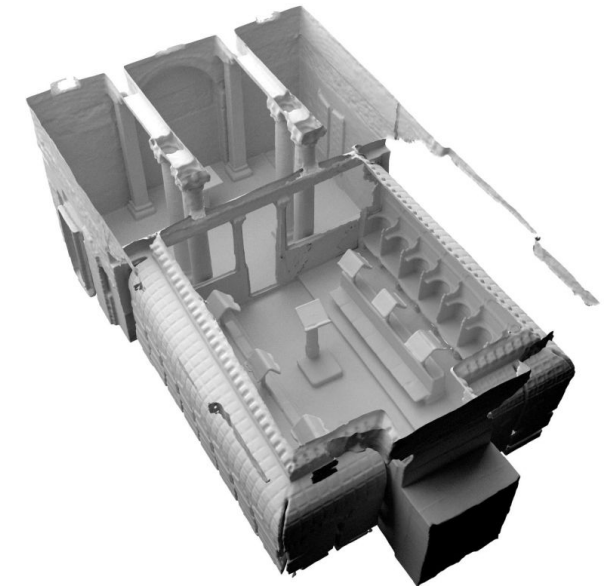2. Query the occupancy network
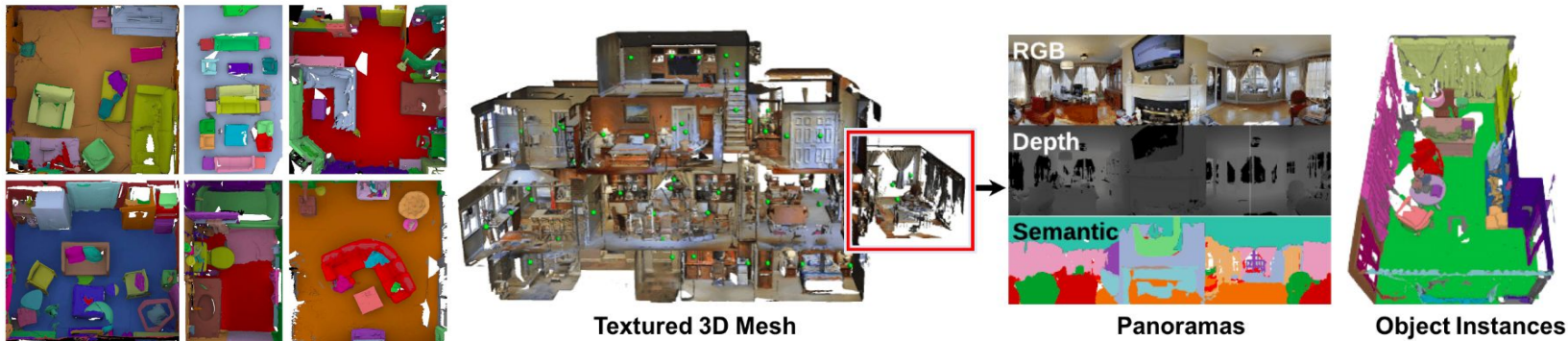
# 3. Results
## Datasets - 4 in total

ShapeNet

Synthetic Indoor Scene Dataset

Matterport3D

ScanNet v2

Textured 3D Mesh

RGB

Depth

Semantic

Panoramas

Object Instances

# 3. Results

**Metrics**

-> Volumetric IoU

$$\text{IoU}(\mathcal{M}_{\text{pred}}, \mathcal{M}_{\text{GT}}) \equiv \frac{|\mathcal{M}_{\text{pred}} \cap \mathcal{M}_{\text{GT}}|}{|\mathcal{M}_{\text{pred}} \cup \mathcal{M}_{\text{GT}}|}$$

-> Chamfer Distance

$$\text{Accuracy}(\mathcal{M}_{\text{pred}}|\mathcal{M}_{\text{GT}}) \equiv \frac{1}{|\partial\mathcal{M}_{\text{pred}}|} \int_{\partial\mathcal{M}_{\text{pred}}} \min_{\mathbf{q} \in \partial\mathcal{M}_{\text{GT}}} \|\mathbf{p} - \mathbf{q}\| d\mathbf{p}$$

$$\text{Chamfer-}L_1(\mathcal{M}_{\text{pred}}, \mathcal{M}_{\text{GT}}) = \qquad \text{Completeness}(\mathcal{M}_{\text{pred}}|\mathcal{M}_{\text{GT}}) \equiv \frac{1}{|\partial\mathcal{M}_{\text{GT}}|} \int_{\partial\mathcal{M}_{\text{GT}}} \min_{\mathbf{p} \in \partial\mathcal{M}_{\text{pred}}} \|\mathbf{p} - \mathbf{q}\| d\mathbf{q}$$

$$\frac{1}{2}(\text{Accuracy}(\mathcal{M}_{\text{pred}}|\mathcal{M}_{\text{GT}}) + \text{Completeness}(\mathcal{M}_{\text{pred}}|\mathcal{M}_{\text{GT}}))$$

-> Normal Consistency

$$\text{Normal-Con.}(\mathcal{M}_{\text{pred}}, \mathcal{M}_{\text{GT}}) \equiv \frac{1}{2|\partial\mathcal{M}_{\text{pred}}|} \int_{\partial\mathcal{M}_{\text{pred}}} |\langle n(\mathbf{p}), n(\text{proj}_2(\mathbf{p})) \rangle| d\mathbf{p}$$

$$+ \frac{1}{2|\partial\mathcal{M}_{\text{GT}}|} \int_{\partial\mathcal{M}_{\text{GT}}} |\langle n(\text{proj}_1(\mathbf{q})), n(\mathbf{q}) \rangle| d\mathbf{q}$$
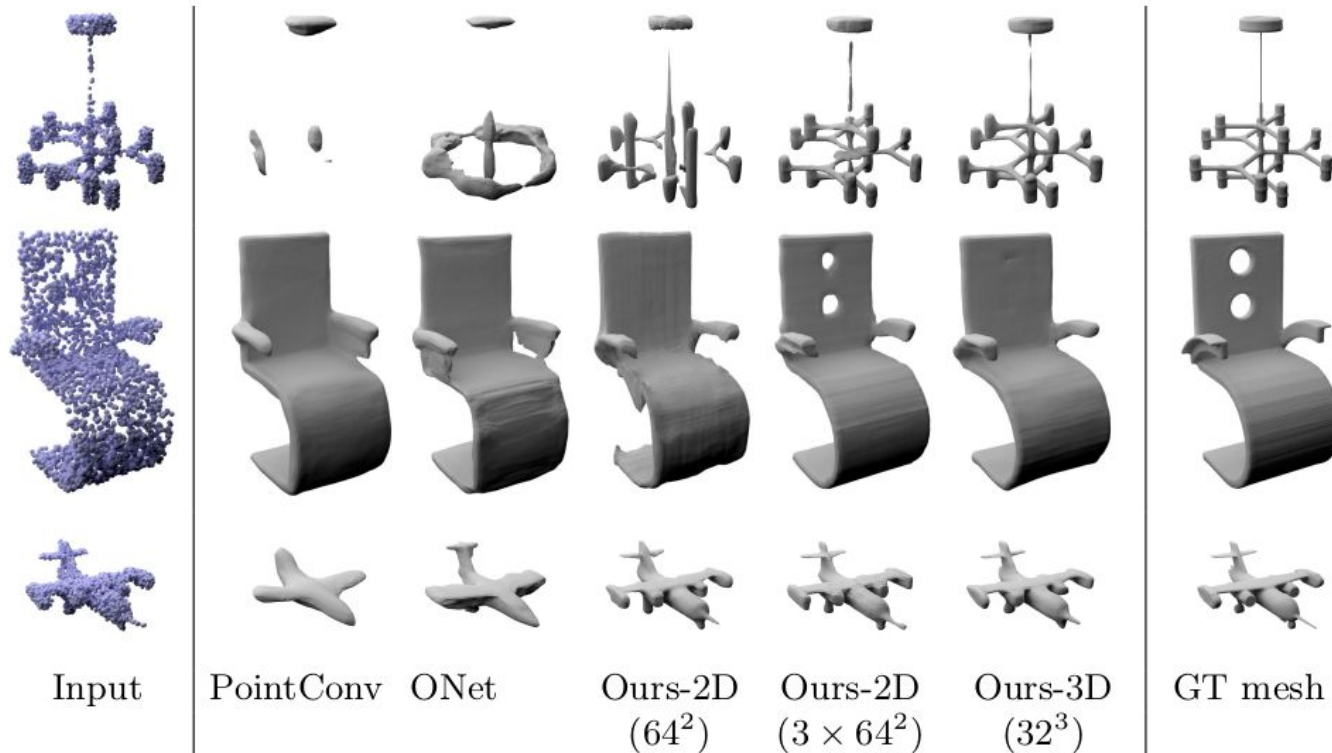
-> F-Score

$$\text{F-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$
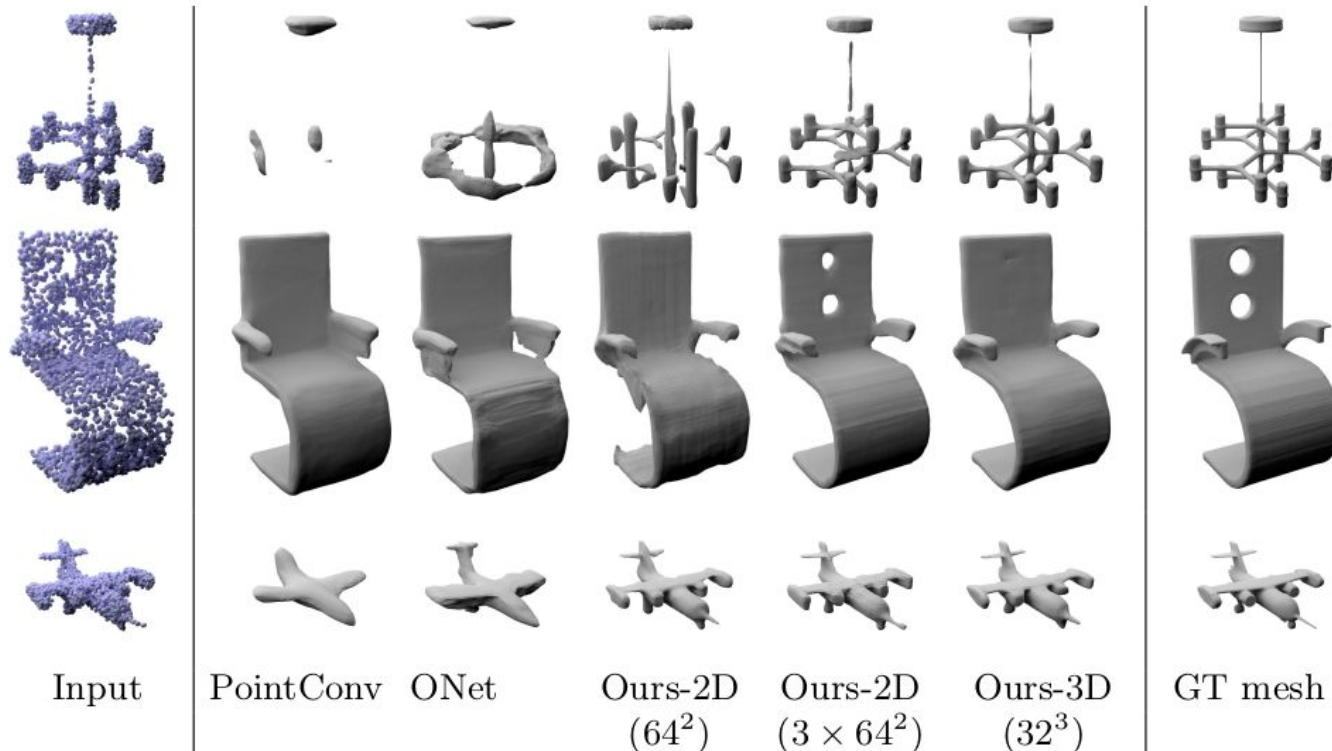
# 3. Results

**Object-Level 3D Reconstruction**

Reconstruction from Point Clouds

# 3. Results

**Object-Level 3D Reconstruction**

Reconstruction from Point Clouds



Input | PointConv | ONet | Ours-2D $(64^2)$ | Ours-2D $(3 \times 64^2)$ | Ours-3D $(32^3)$ | GT mesh
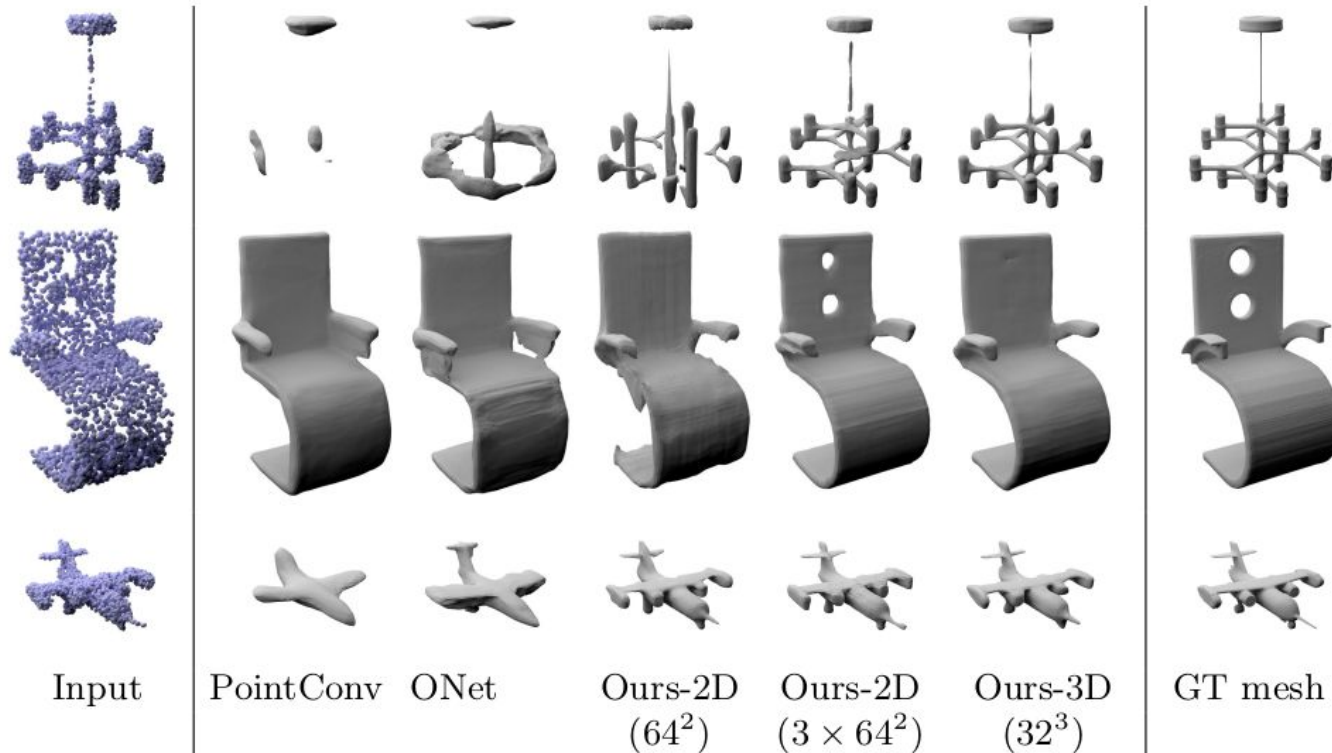
- Baseline: PointConv
  - PointNet++ encoder
    - Remove canonical planes
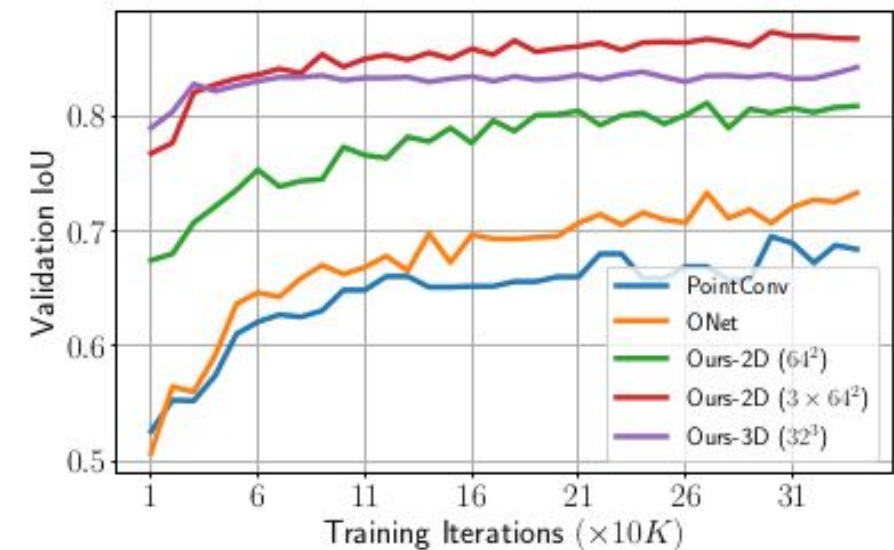  - Instead of the 2D decoder and interpolation
    - Gaussian regression

# 3. Results

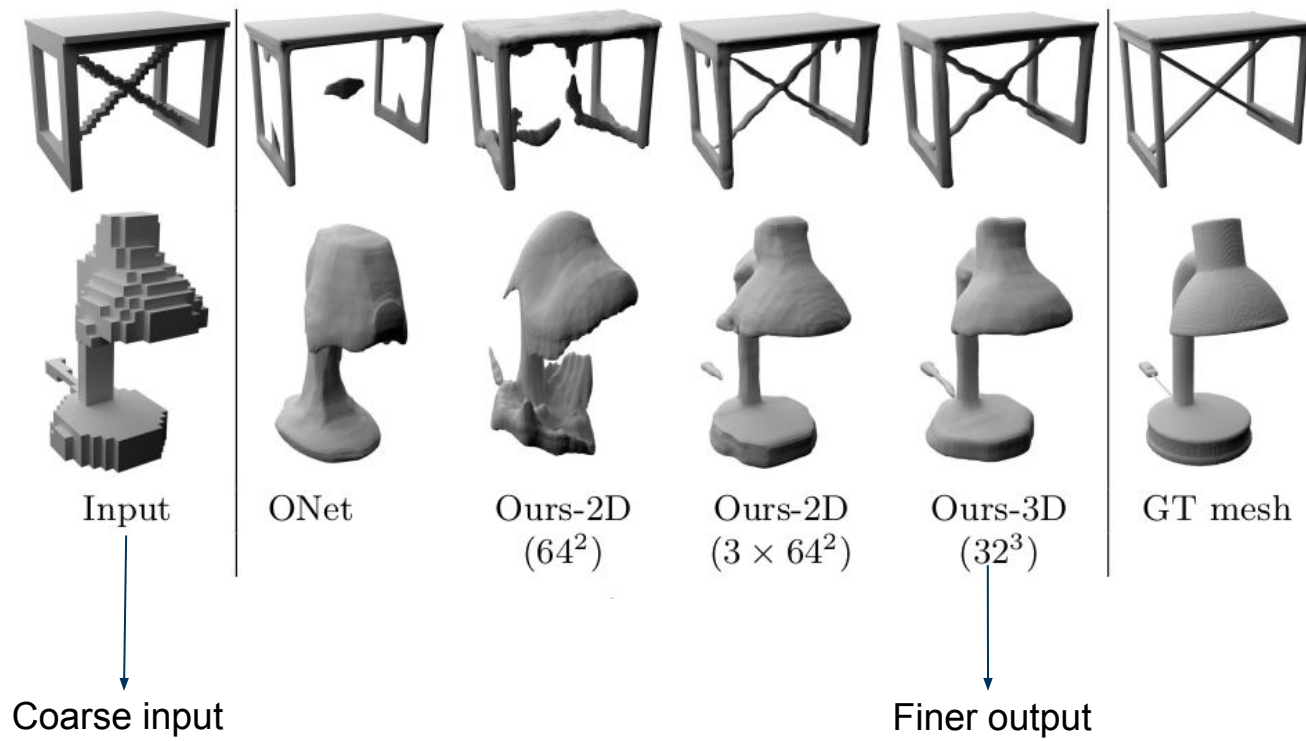**Object-Level 3D Reconstruction**

Reconstruction from Point Clouds



- Convolutional Occupancy Networks
  - Reconstruction of complex shapes
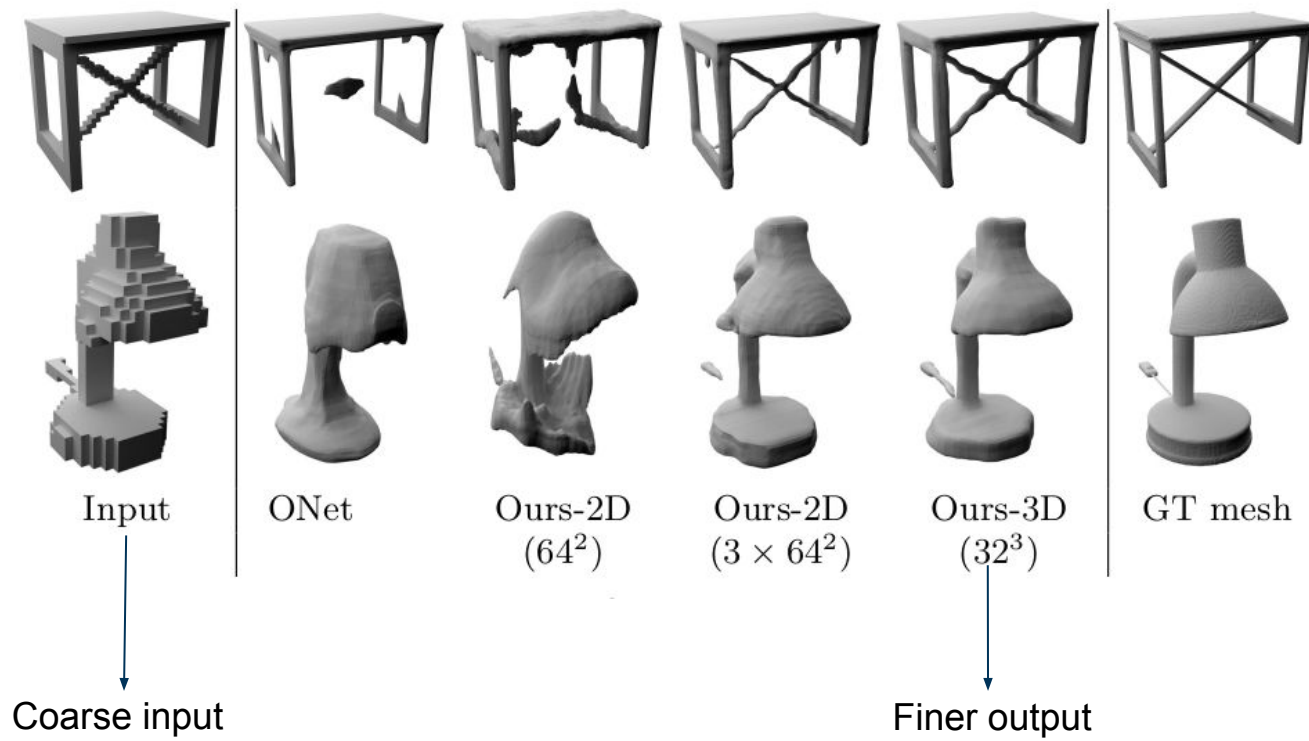  - Faster convergence

# 3. Results

**Object-Level 3D Reconstruction**

Voxel Super-Resolution



Coarse input

Finer output

# 3. Results

**Object-Level 3D Reconstruction**

Voxel Super-Resolution



Input | ONet | Ours-2D $(64^2)$ | Ours-2D $(3 \times 64^2)$ | Ours-3D $(32^3)$ | GT mesh

Coarse input

Finer output

- Convolutional Occupancy Networks
  - Recover high-resolution details
  - Three planes perform similar to the volumetric encoder while consuming 37% of the GPU
  - The single-plane approach is not powerful

# 3. Results

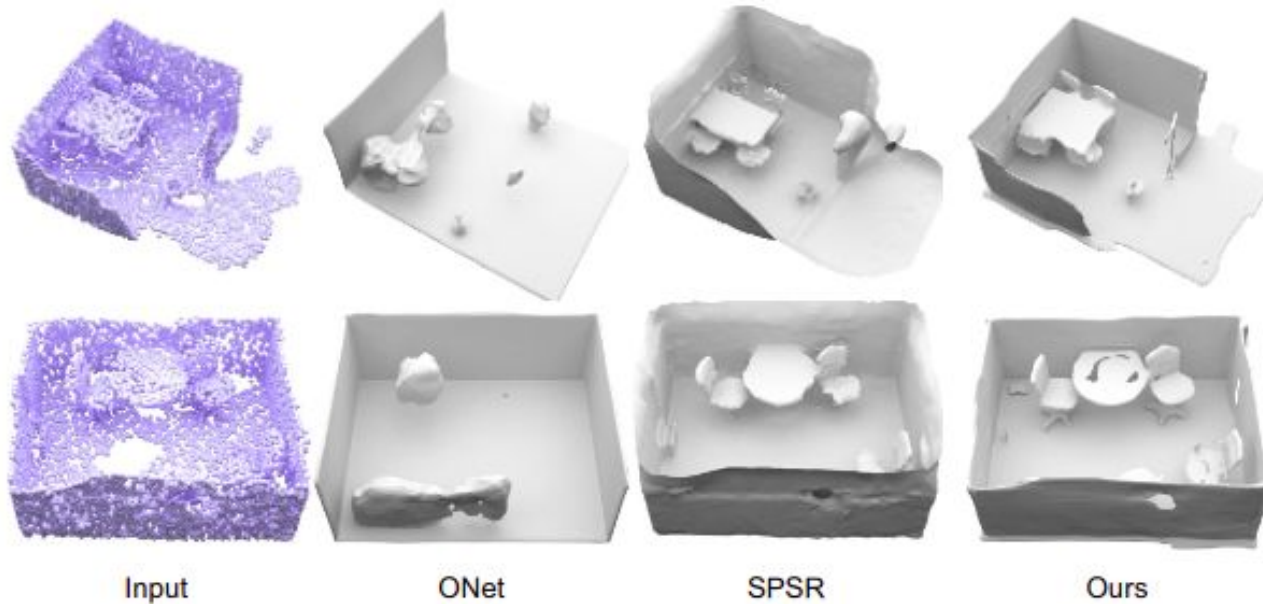**Scene-Level Reconstruction**

Synthetic dataset evaluation



Input    ONet    SPSR    Ours    GT

- Occupancy Networks
  - Can not scale to bigger scenes

- SPSR
  - Requires the normals of the points
  - Noisy results

# 3. Results

**Scene-Level Reconstruction**
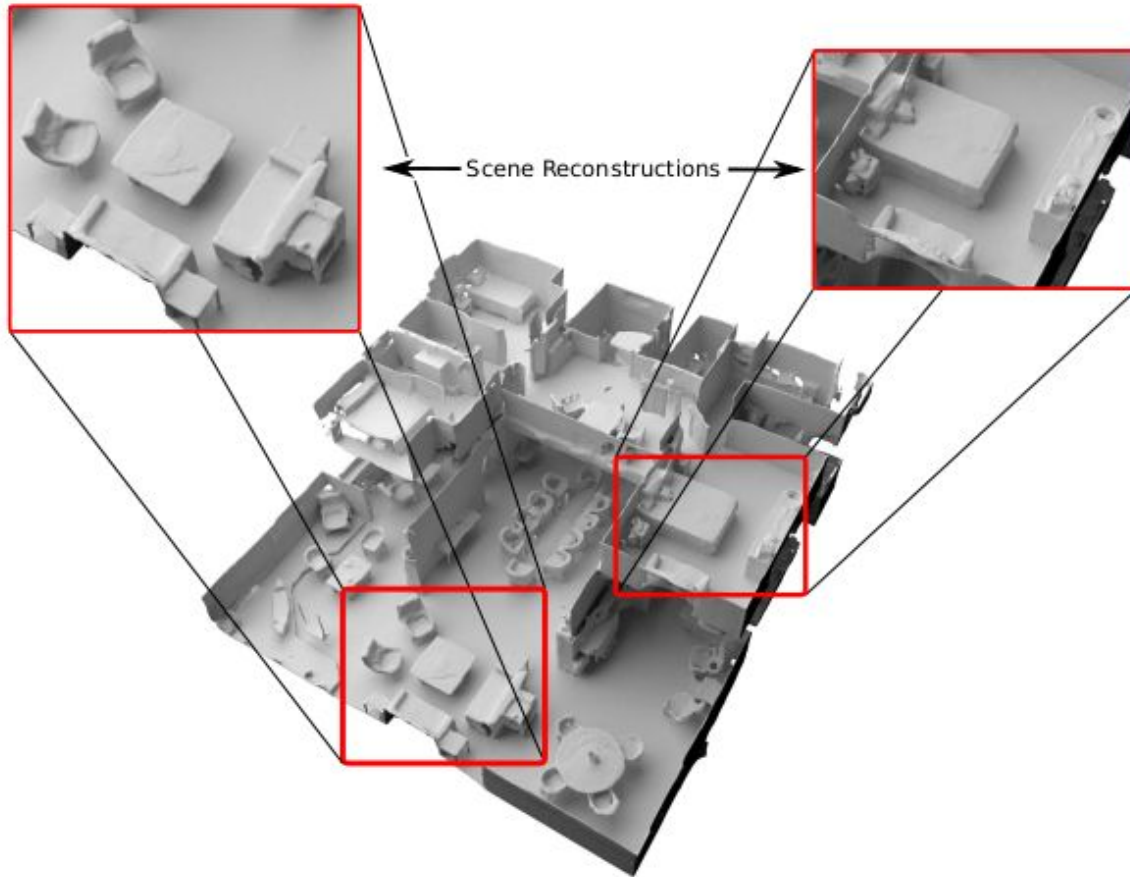
Trained on synthetic and transfer to ScanNet v2



Input        ONet        SPSR        Ours

-> All previous methods mostly fail on this task

# 3. Results

**Large-Scale Reconstruction**
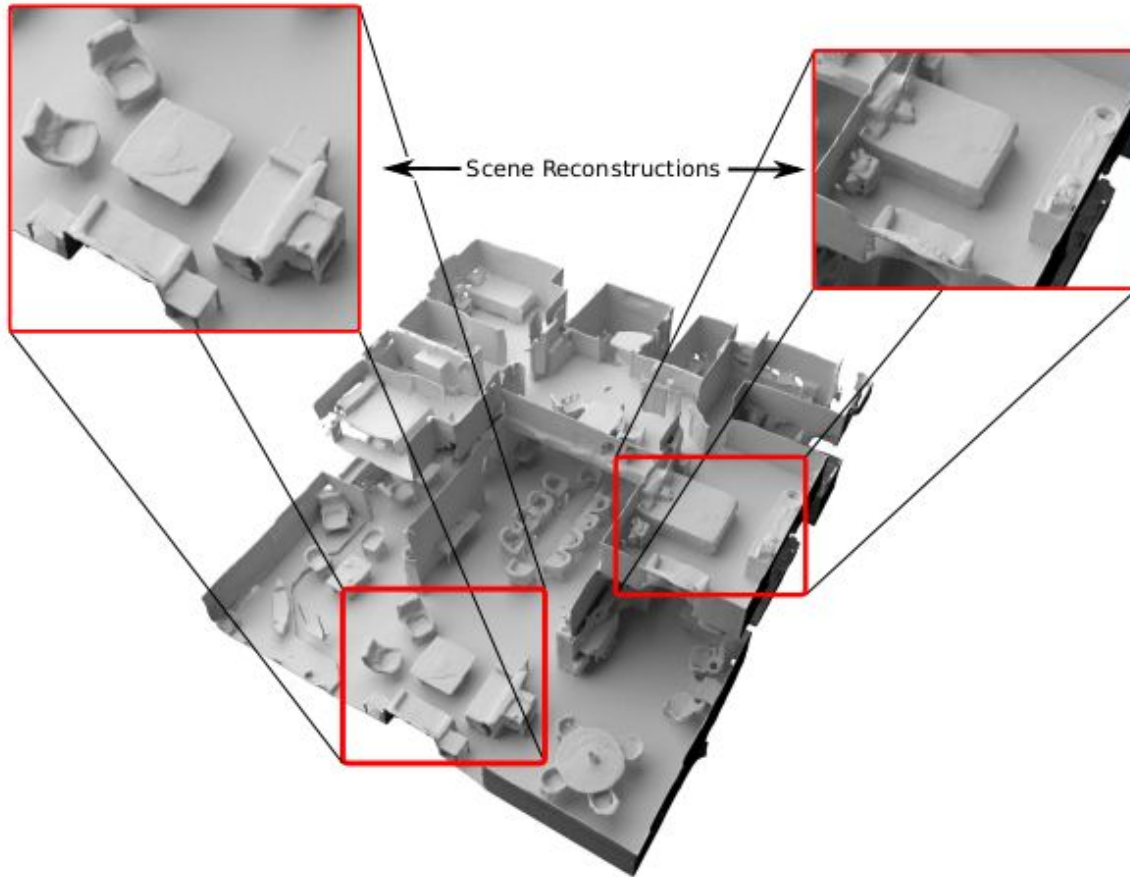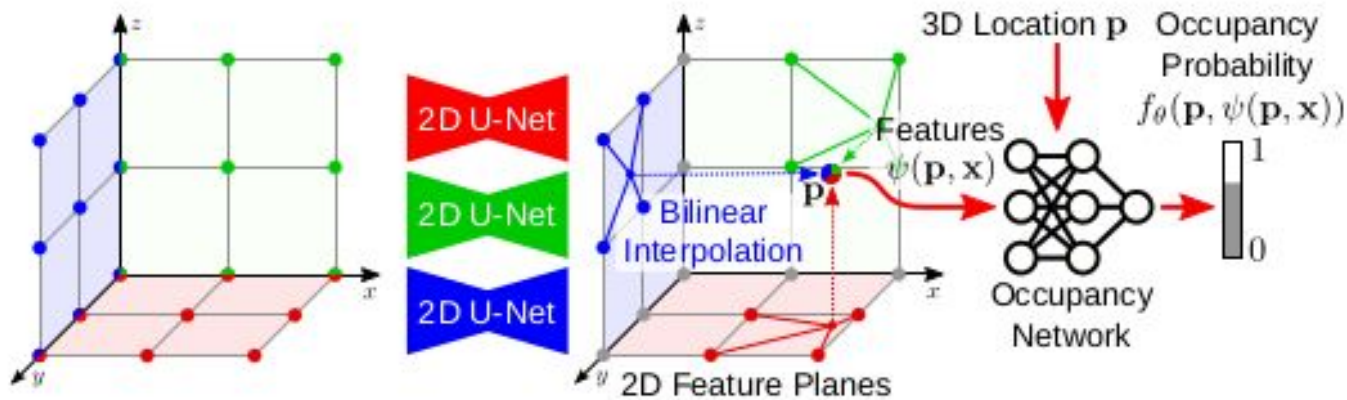
Trained on synthetic and transfer to Matterport3D



-> Trained on synthetic crops

-> During inference, use sliding window

-> 3D CNN performed the best

# 3. Results

**Large-Scale Reconstruction**

Trained on synthetic and transfer to Matterport3D



Scene Reconstructions

-> Trained on synthetic crops

-> During inference, use sliding window

-> 3D CNN performed the best

-> The authors do not explain how to merge the patches
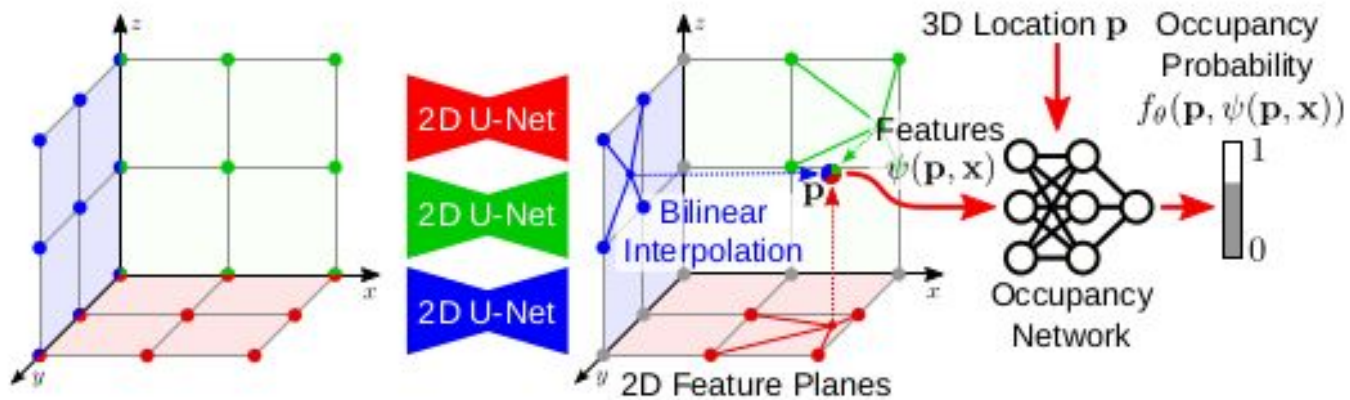  - What happens with the artifacts of the overlapping windows?
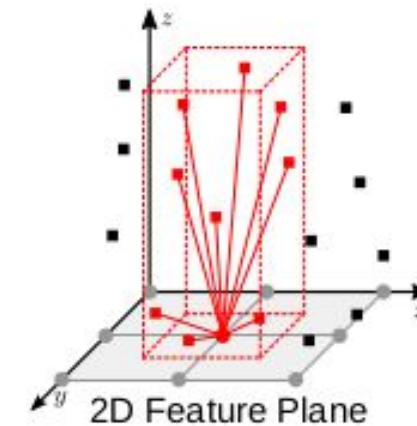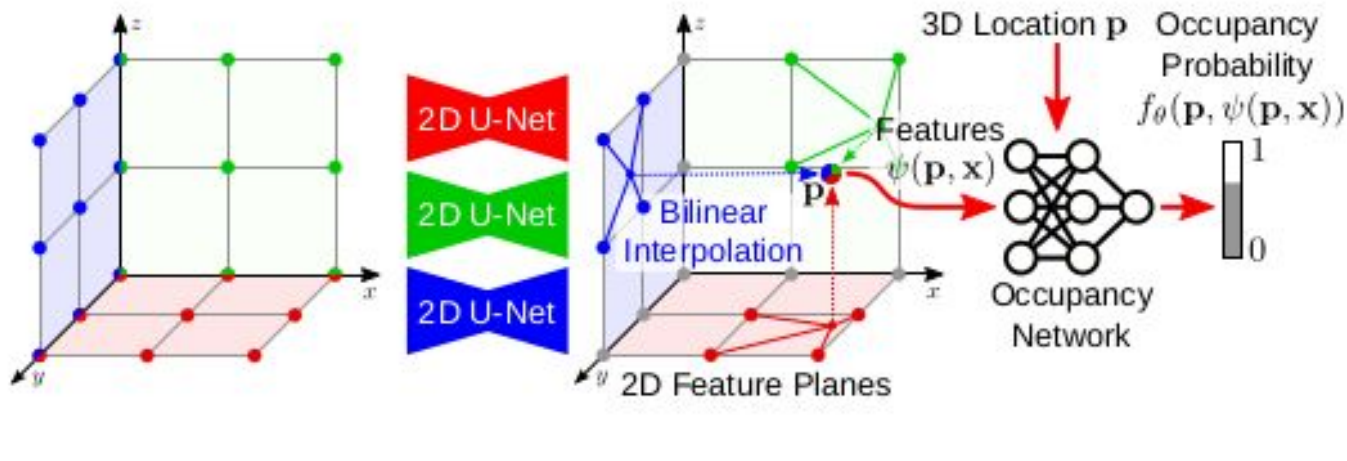
# 4. Personal comments

1. Shared 2D U-Nets?

# 4. Personal comments

1. Shared 2D U-Nets?     2. Use shallow Neural Net instead of sum?
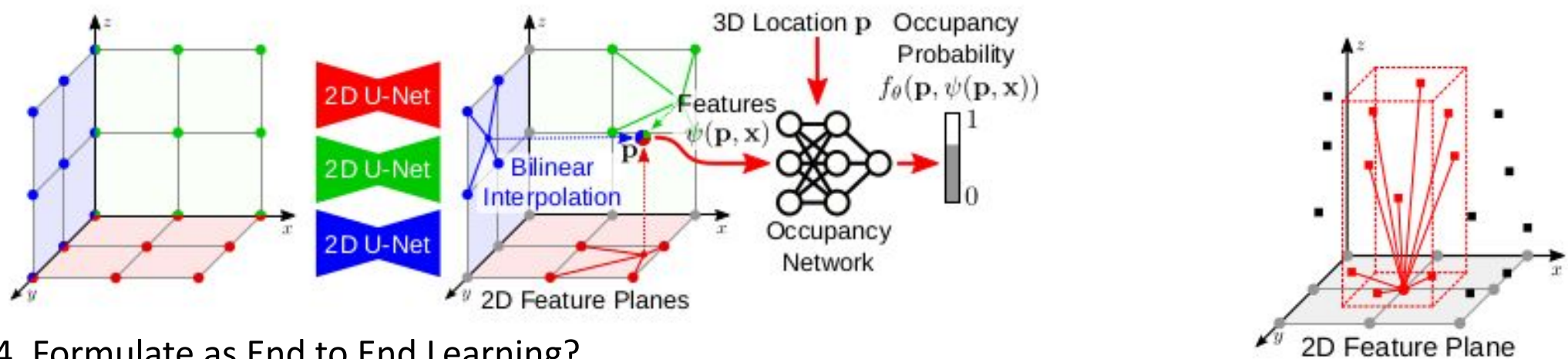
# 4. Personal comments

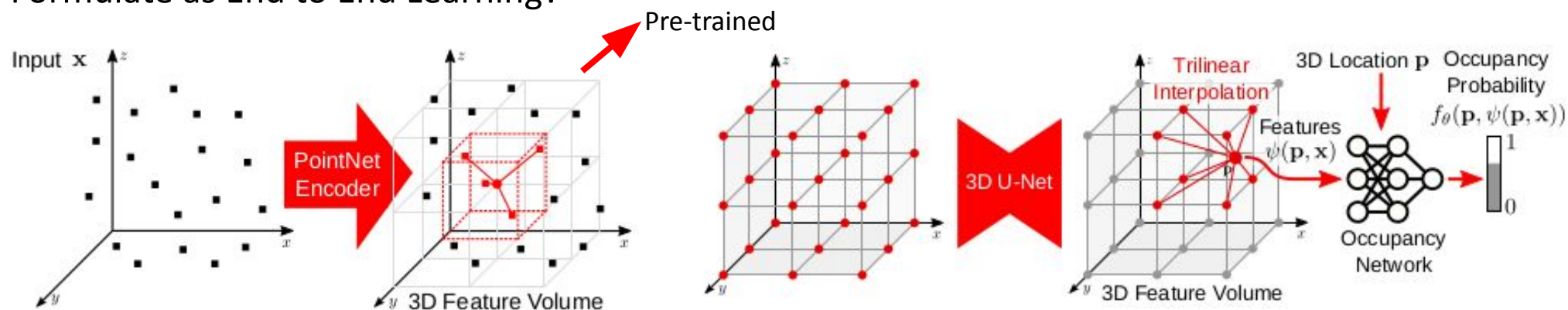1. Shared 2D U-Nets?    2. Use shallow Neural Net instead of sum?    3. Average or max pooling aggregation?

# 4. Personal comments

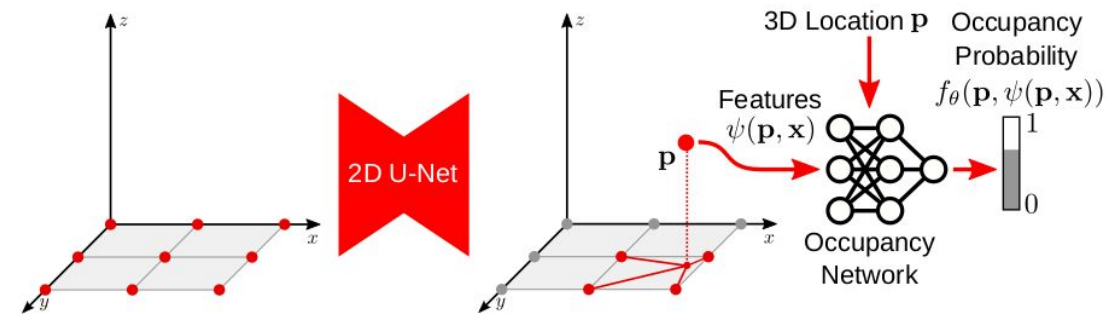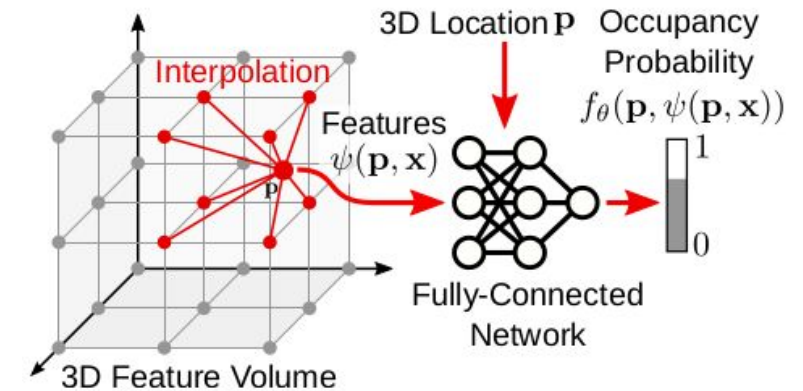1. Shared 2D U-Nets?    2. Use shallow Neural Net instead of sum?    3. Average or max pooling aggregation?


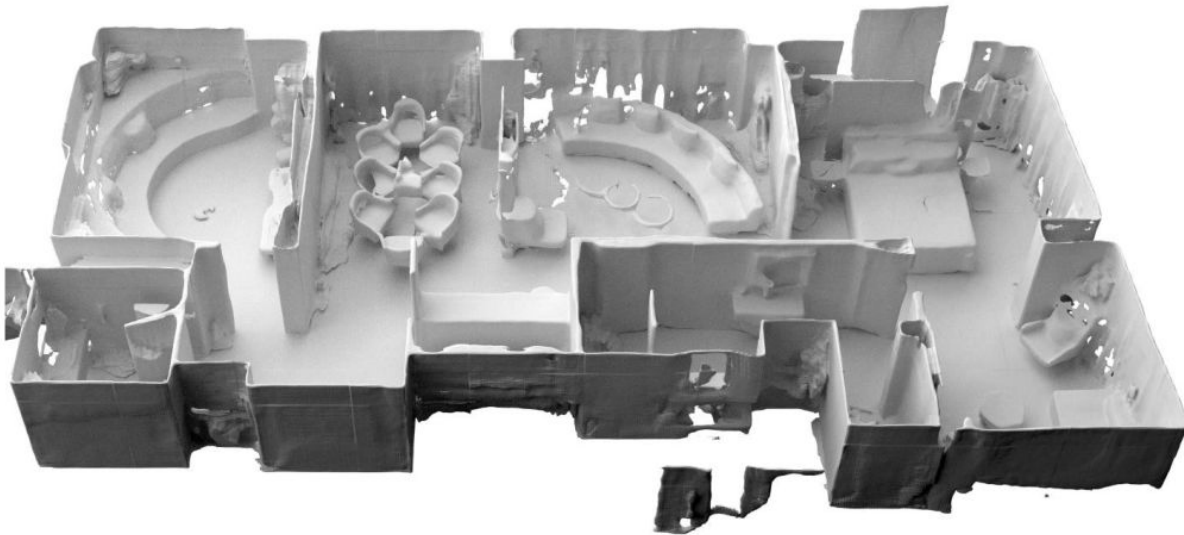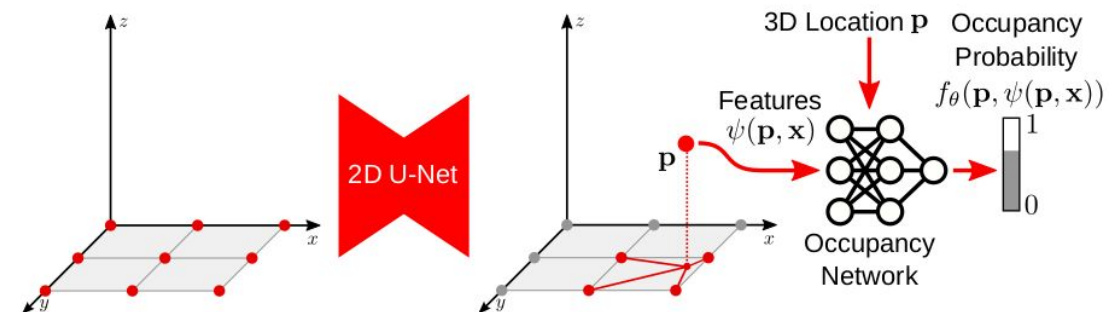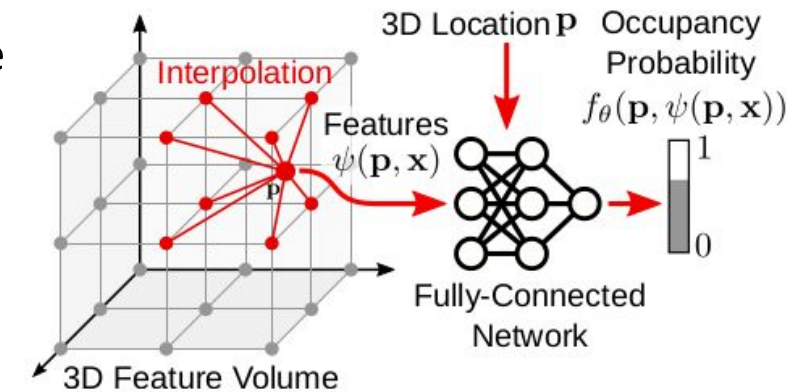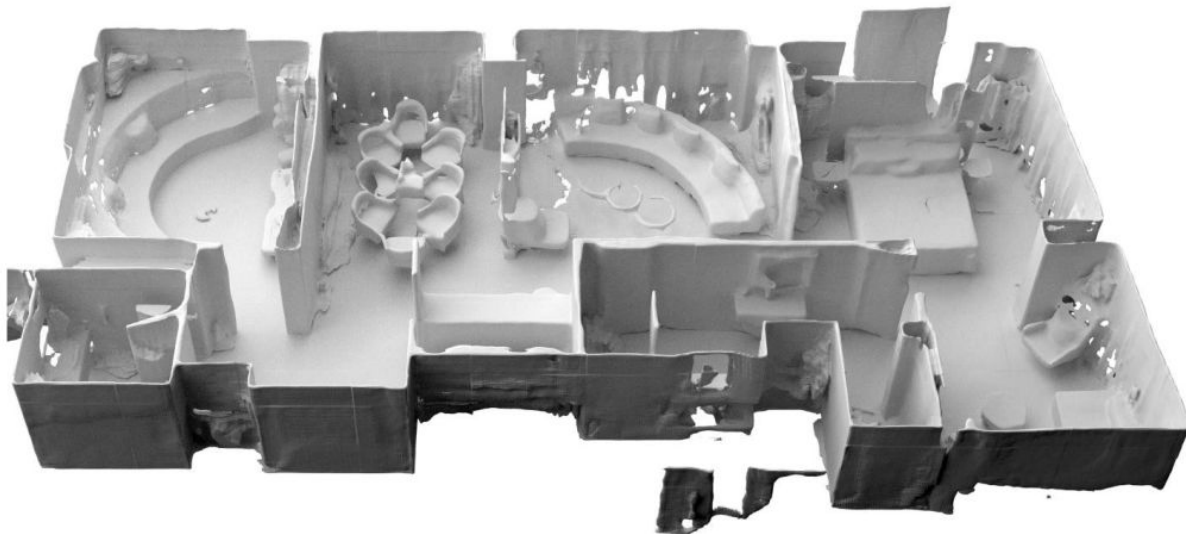
4. Formulate as End to End Learning?

# 5. Summary

- Conv. Occupancy Networks can be transferred to noisy real large-scale scenes
- Incorporate global and local information
- Faster training

# 5. Summary

- Conv. Occupancy Networks can be transferred to noisy real large-scale scenes
- Incorporate global and local information
- Faster training
- **But**
  - Translation equivariant w.r.t to multiples translations of the voxel size
  - No rotation equivariant
  - Reality gap is still present

# Thank you for the attention!
# Questions?

# References

[1] S. Peng, M. Niemeyer, L. Mescheder, M. Pollefeys, and A. Geiger, "Convolutional Occupancy Networks Computer", in ECCV, 2020

[2] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space", in: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2019

[3] https://www.cvlibs.net/talks/talk_oxford_2020_04_27.pdf

[4] https://www.gamersnexus.net/gg/762-voxels-vs-vertexes-in-games

[5]https://www.researchgate.net/figure/3D-mesh-triangles-with-different-resolution-3D-Modelling-for-programmers-Available-at_fig2_322096576

[6] https://www.revopoint3d.com/point-cloud-and-3d-image/

[7] A. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "Shapenet: An information-rich 3d model repository", in arXiv.org 1512.03012, 2015

[8] A. Dai, A. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Niessner, "Scannet: Richly-annotated 3d reconstructions of indoor scenes", in: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2017

[9] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3D: Learning from RGB-D data in indoor environments", in Proc. of the International Conf. on 3D Vision (3DV), 2017