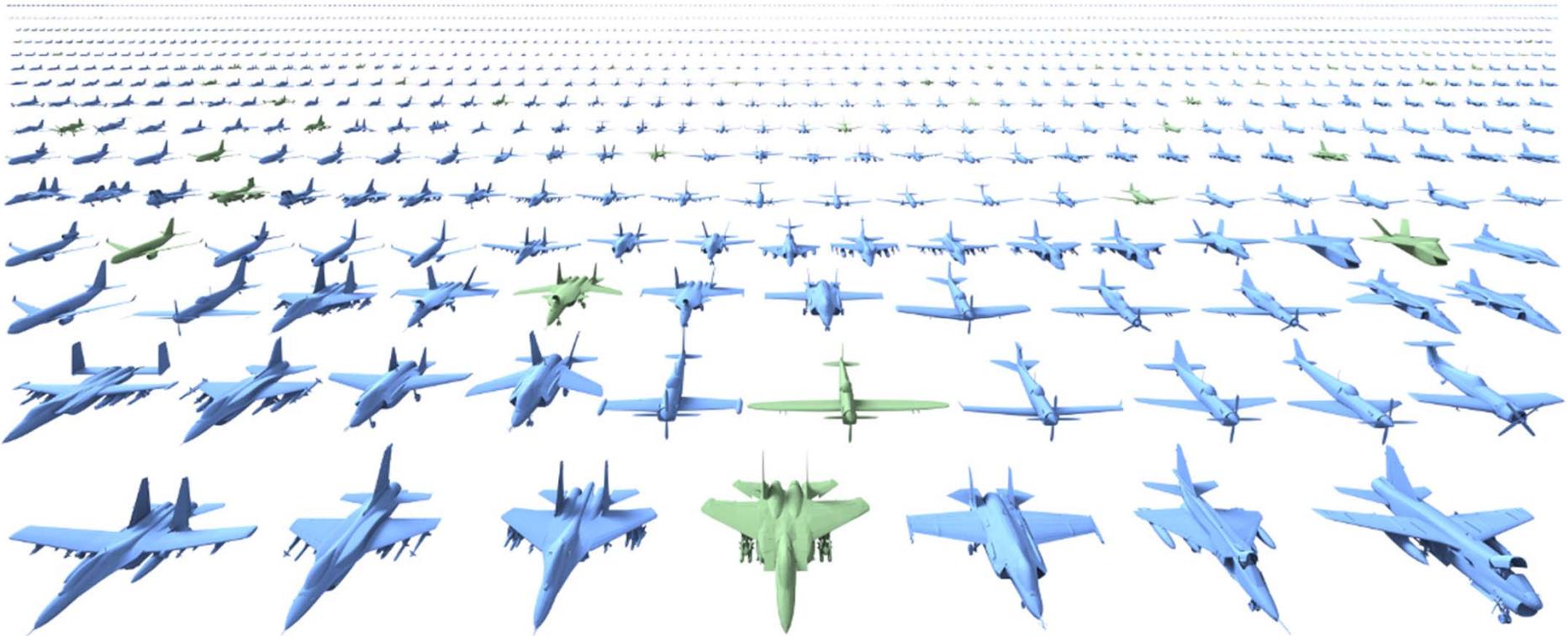


Generative models of 3D Shapes and Scenes



Evangelos Kalogerakis –
574/674

What representation should we generate?

- **Multiple views**



3D Shape
Reconstruction from
Sketches via Multi-view
Convolutional Networks,
Lun et al, 2017

- **Point cloud**



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

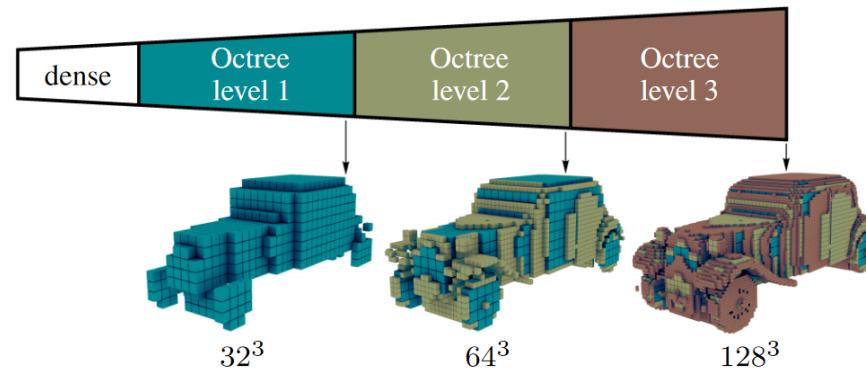
- **Voxels**



Learning a Probabilistic Latent
Space of Object Shapes via
3D Generative-Adversarial
Modeling, Wu et al 2017

What representation should we generate?

- **Octree cells**



Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs, Tatarchenko et al. 2017

- **Parametric Patches**



AtlasNet,
Groueix et al.
2018

- **Volumetric Primitives**



Learning Shape Abstractions by Assembling Volumetric Primitives Tulsiani et al. 2017

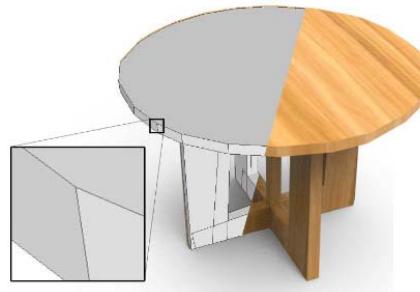
What representation should we generate?

- **Implicits**



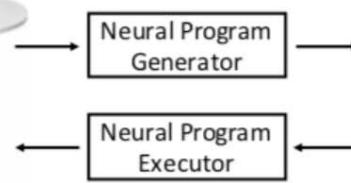
Convolutional Occupancy Networks
Peng al 2020

- **Meshes**



BSP-Net, Chen et al. 2020

- **Structure/Program**



ShapeAssembly, Jones et al 2020

```
Draw("Top", "Circle", position, geometry)
for(i < 2, "translation", a)
  for(j < 2, "translation", b)
    Draw("Leg", "Cub", position + i*a + j*b, geometry)
    for(i < 2, "translation", c)
      Draw("Layer", "Rec", position + i*c, geometry)
```

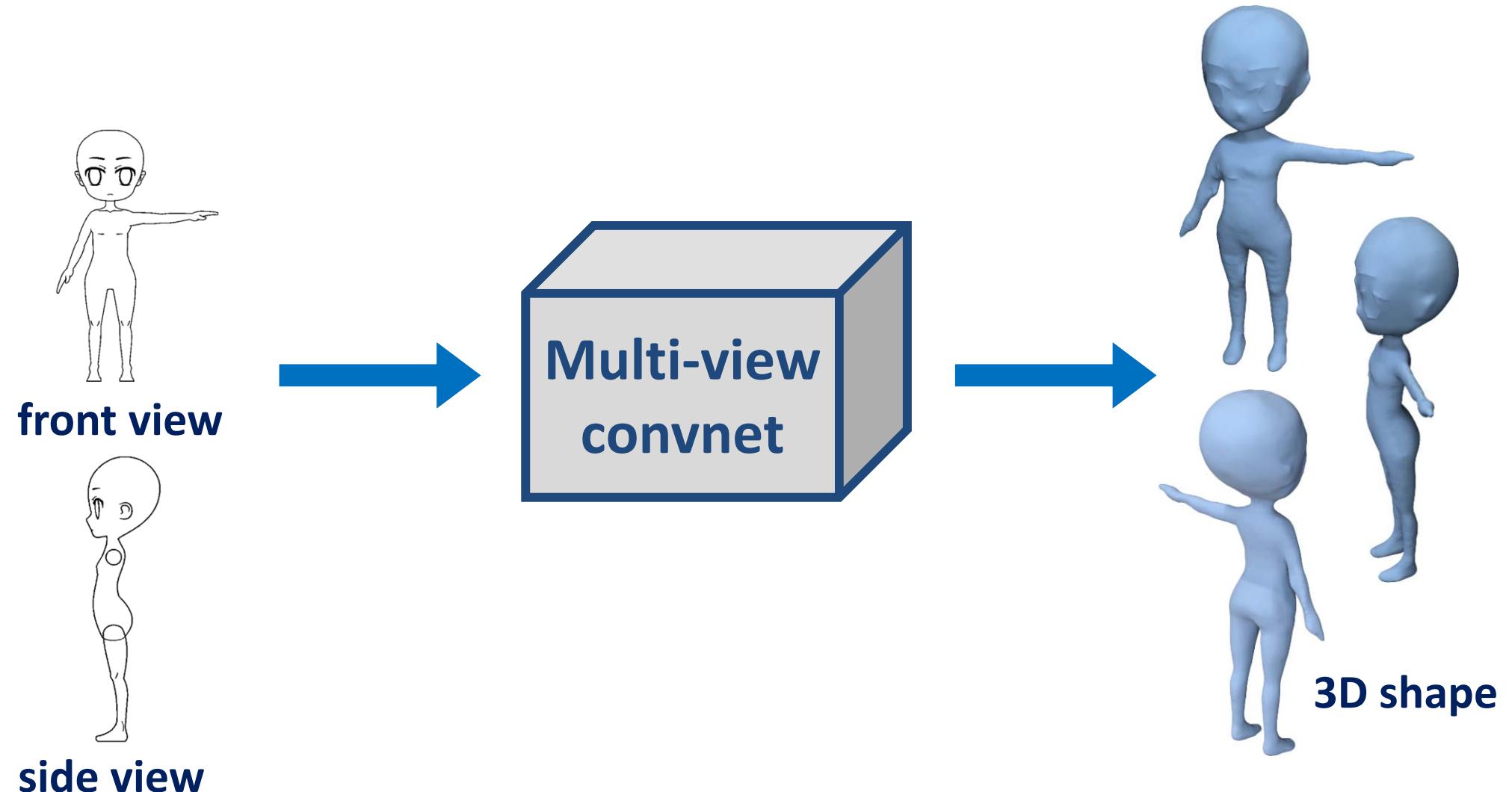
How to generate shapes/scenes?

- Encoder-Decoders
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

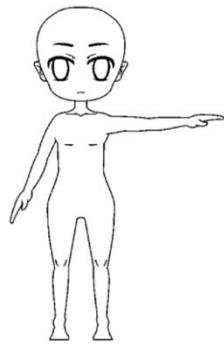
How to generate shapes/scenes?

- **Encoder-Decoders**
 - Case Study: Multi-view decoder
 - Case Study: Implicit decoder
 - Case Study: Patch decoder
 - Case Study: Mesh Decoder
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

Sketch-to-shape



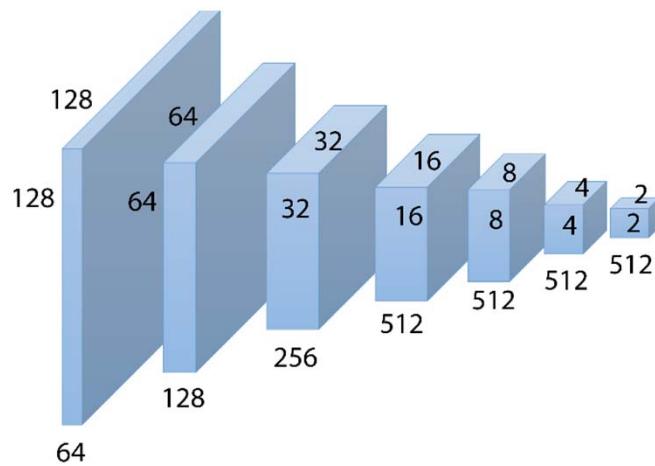
Encoder



front view



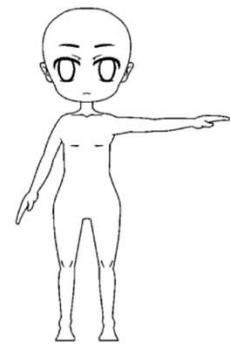
side view



Feature representations
capturing **increasingly larger**
context in the sketches

Decoder

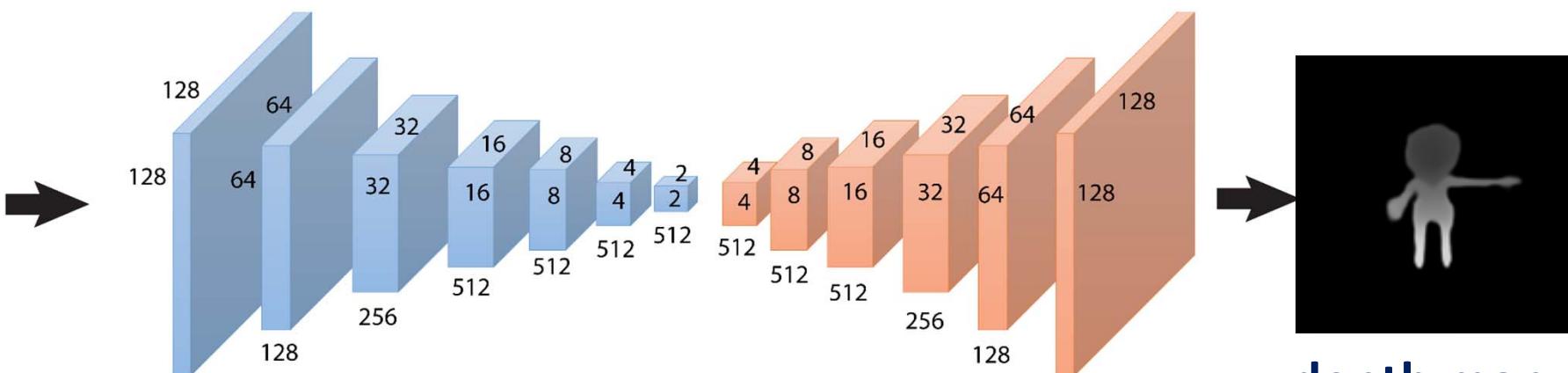
Infer depth and normal maps



front view



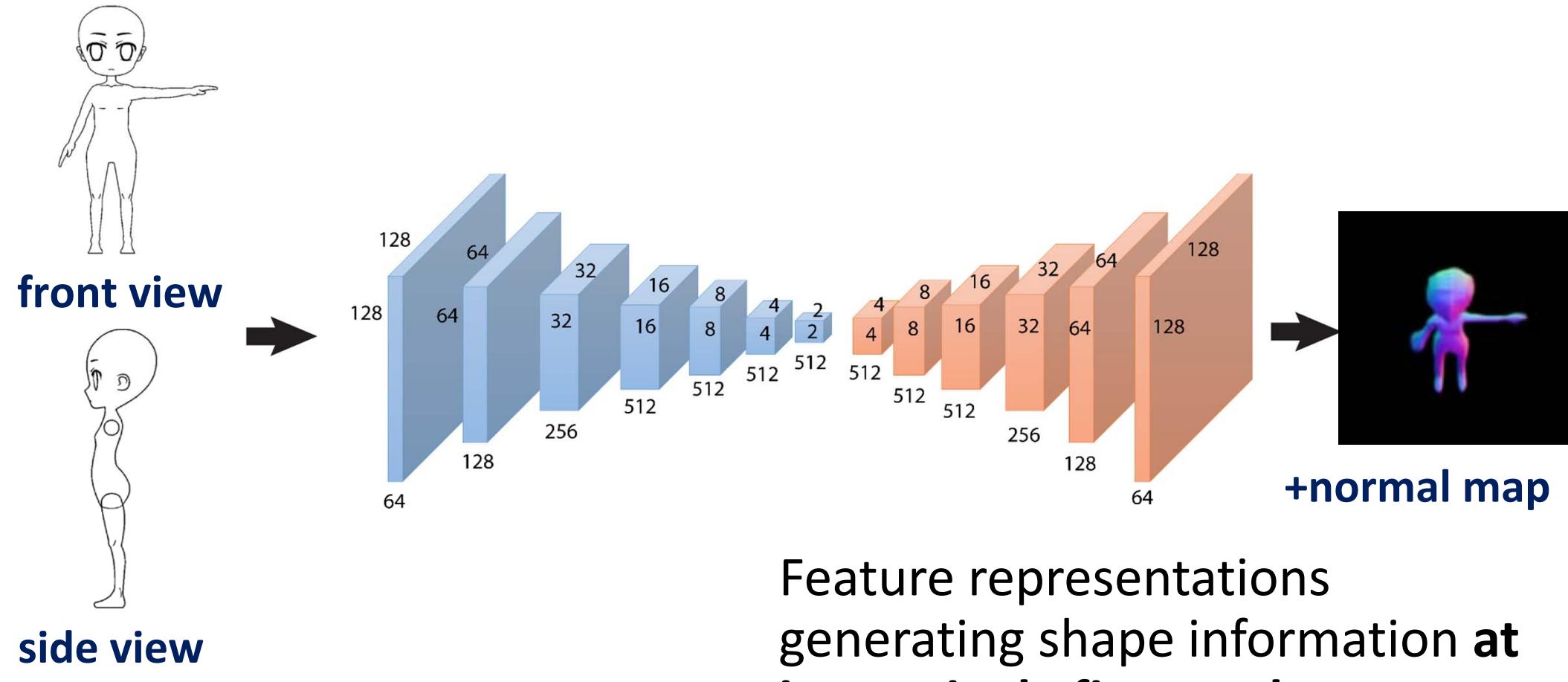
side view



Feature representations
generating shape information at
increasingly finer scales

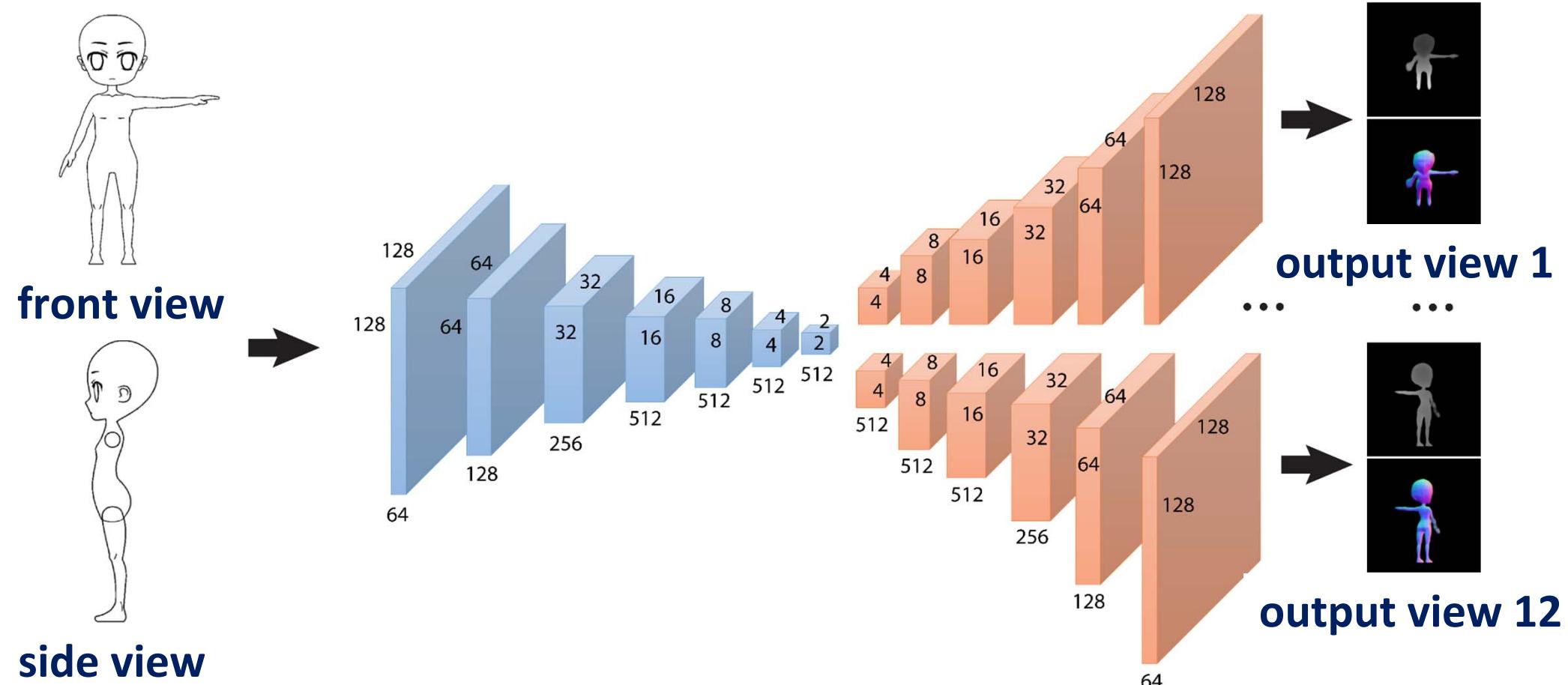
Decoder

Infer depth and normal maps



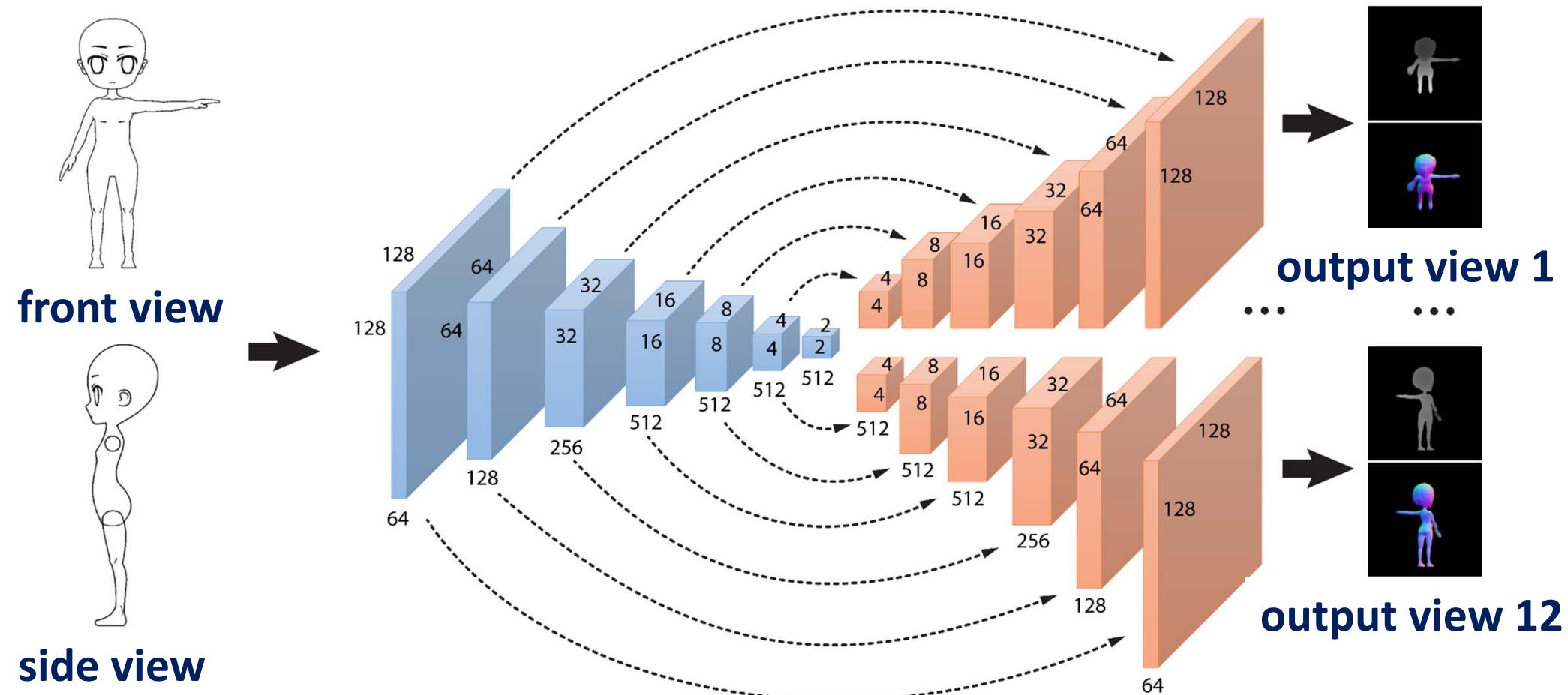
Multi-view decoder

Infer depth and normal maps for several views



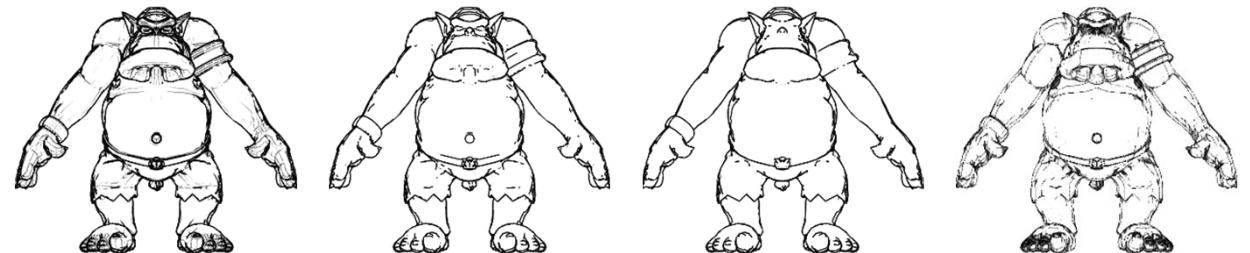
“U-Net”

Feature representations in the decoder depend on **previous layer & encoder's corresponding layer**



U-net: Ronneberger et al. 2015,
Isola et al. 2016

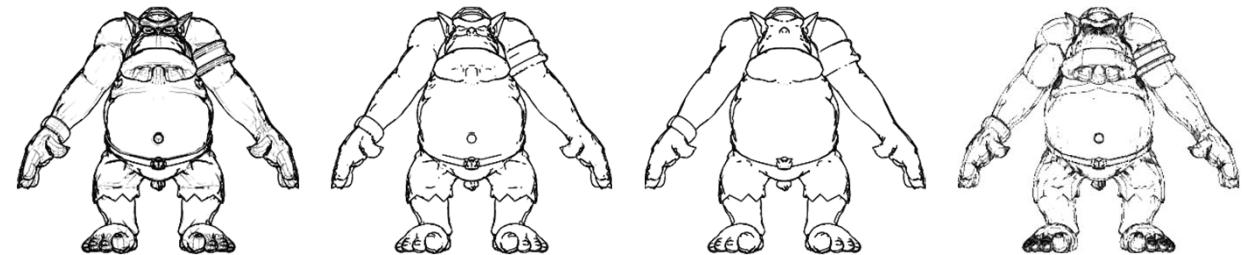
Training data



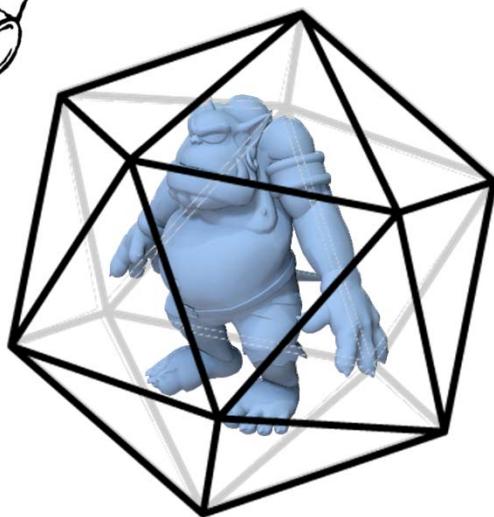
Synthetic line drawings



Training data



Synthetic line drawings



12 views



Training depth and normal maps

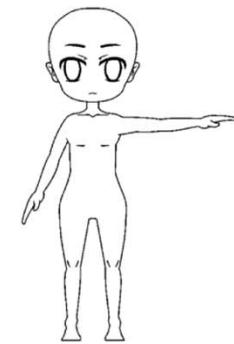
Initial training loss

Penalize per-pixel depth reconstruction error:

$$\sum_{pixels} |d_{pred} - d_{gt}|$$

& per-pixel normal reconstruction error:

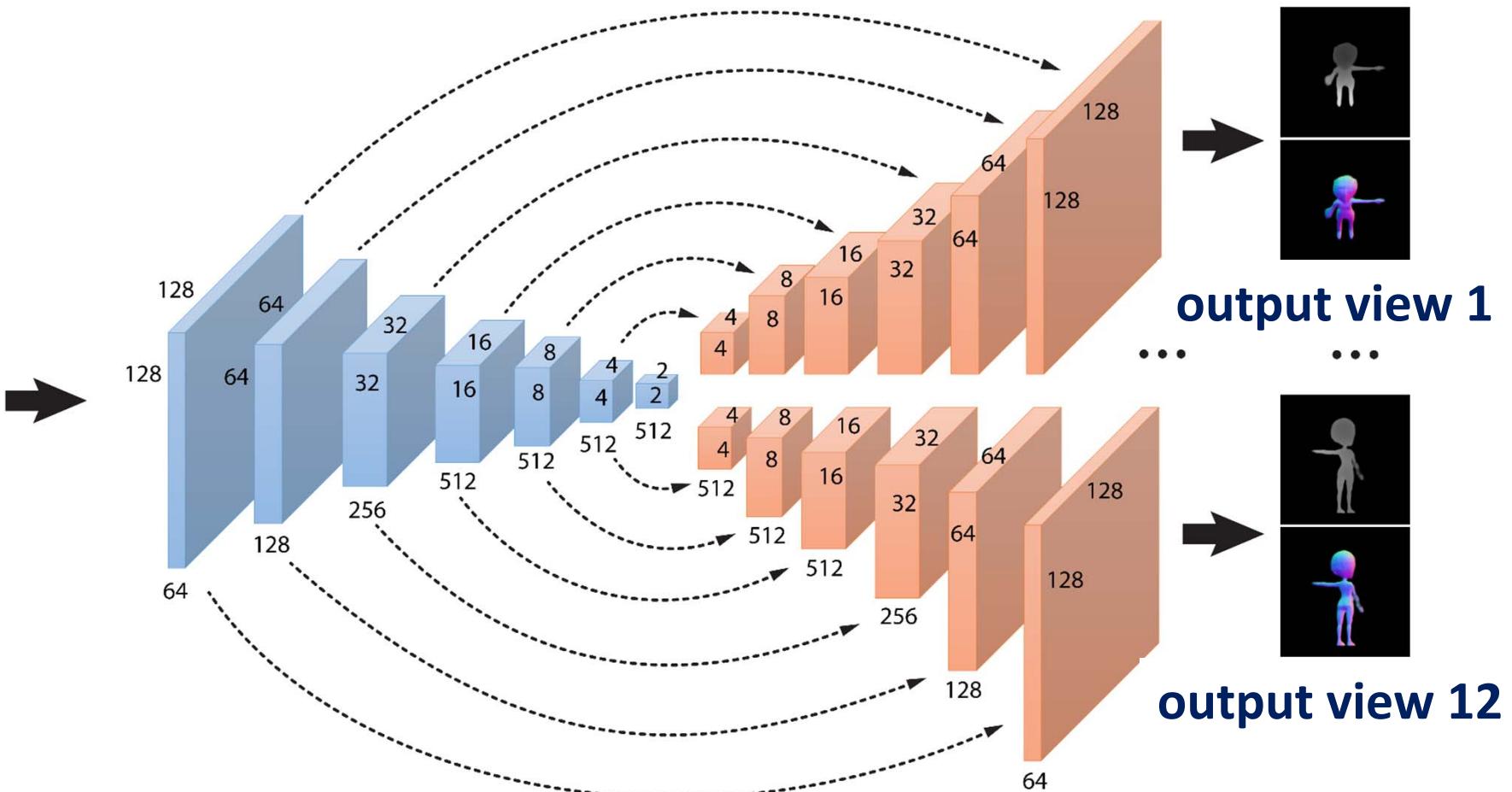
$$\sum_{pixels} (1 - \mathbf{n}_{pred} \cdot \mathbf{n}_{gt})$$



front view



side view



U-net: Ronneberger et al. 2015,
Isola et al. 2016

Training data



Character
10K models



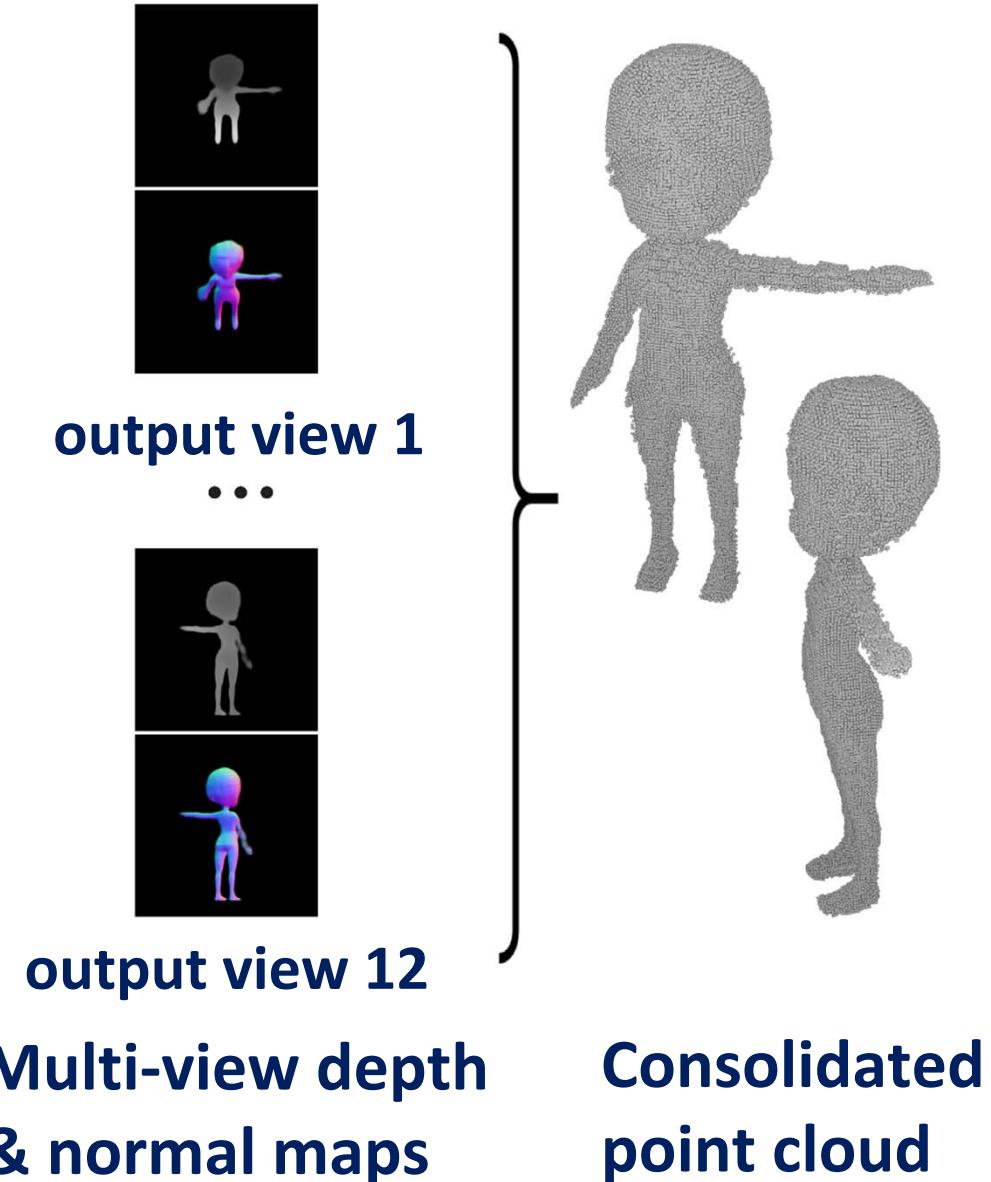
Chair
10K models



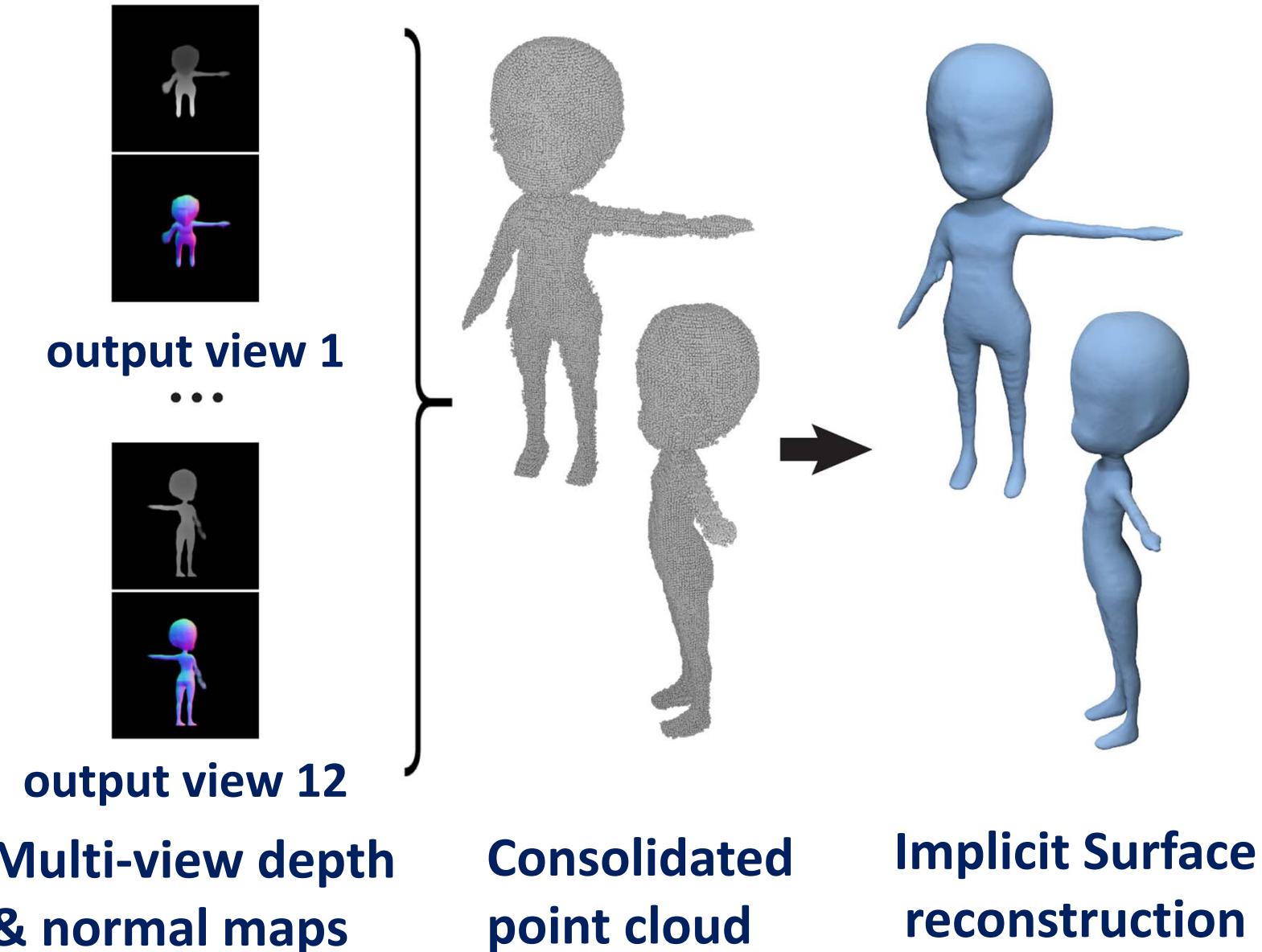
Airplane
3K models

**Models from “The Models Resource” &
3D Warehouse**

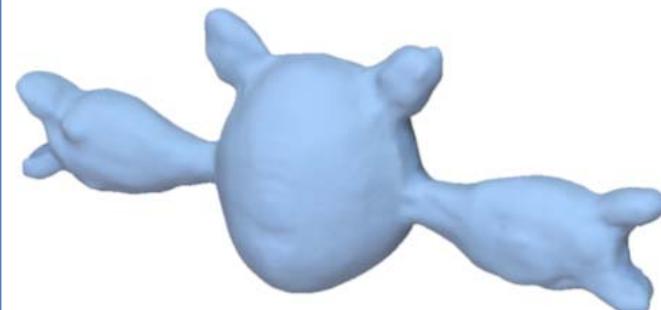
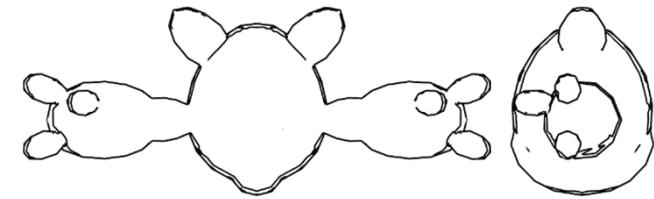
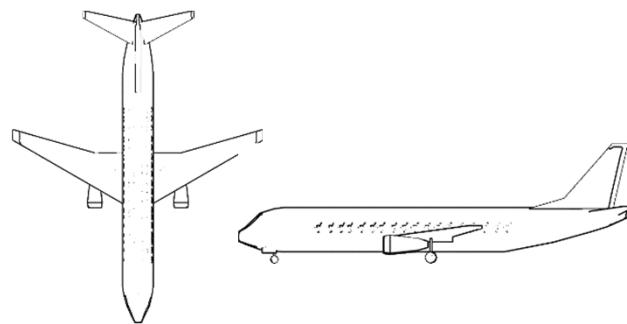
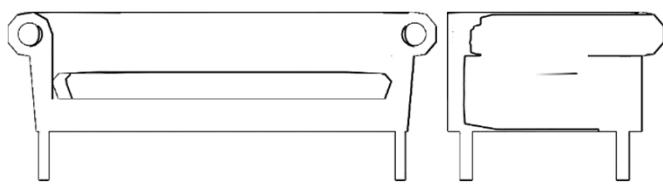
Multi-view depth & normal map fusion



Surface reconstruction



Qualitative results



How to generate shapes/scenes?

- **Encoder-Decoders**
 - Case Study: Multi-view decoder
 - **Case Study: Implicit decoder**
 - Case Study: Patch decoder
 - Case Study: Mesh Decoder
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

Image-to-shape

ResNet
encoder

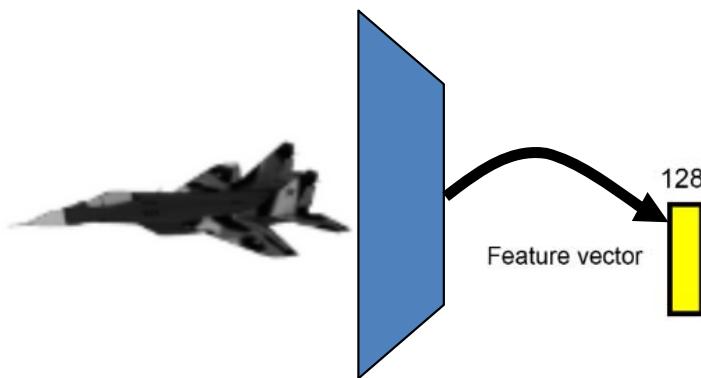
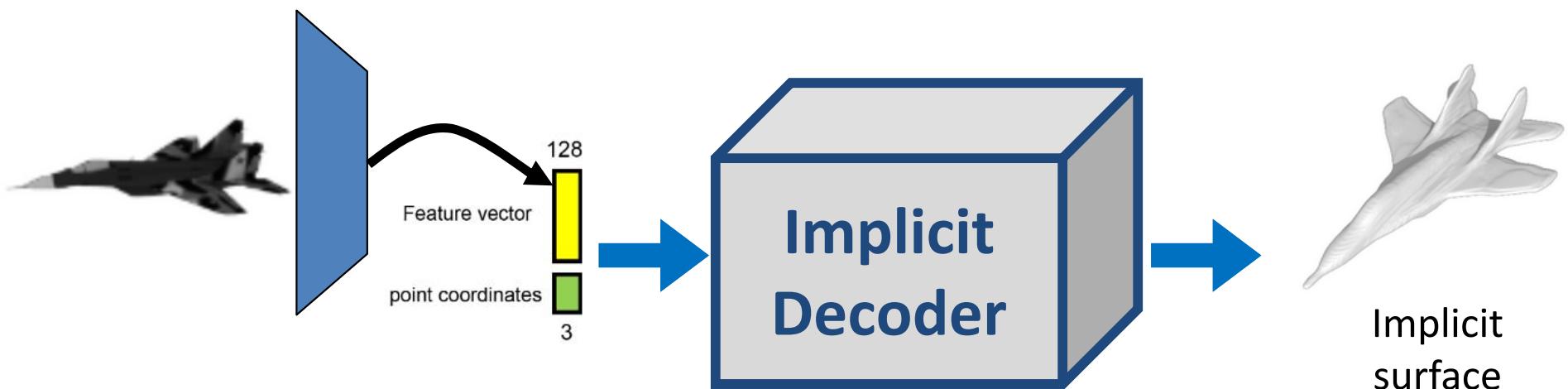


Image-to-shape

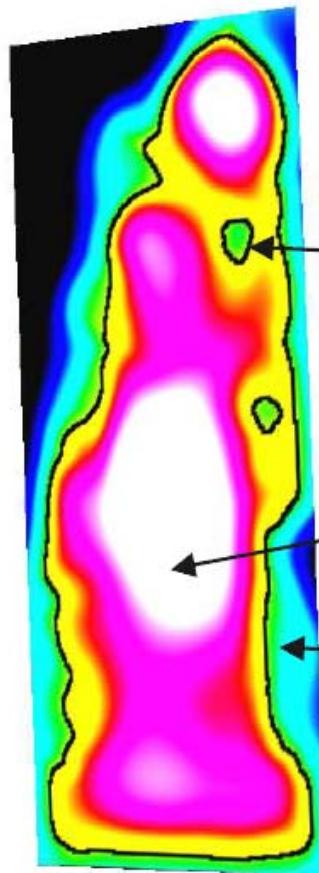
ResNet
encoder



Implicit surfaces

Surface is defined by **implicit** function $f(x,y,z) = 0$. In other words, implicit surfaces are isosurfaces through some scalar field in 3D.

3-D example



Implicit Surface:

$$f(x,y,z) = 0$$

Inside:

$$f(x,y,z) < 0$$

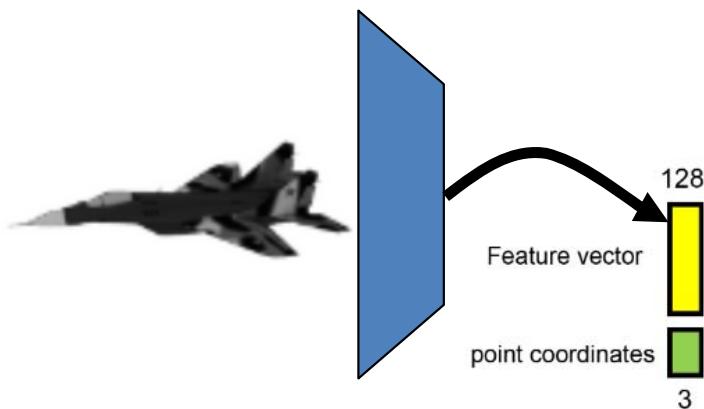
Outside:

$$f(x,y,z) > 0$$



Image-to-shape

ResNet
encoder



At test time, pass the
coordinates of each
point in a \mathbb{R}^3 grid

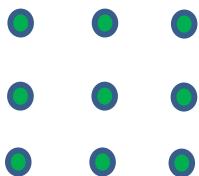


Image-to-shape

ResNet
encoder



Feature vector
point coordinates

128

3

Implicit
Decoder

2048

1024

512

256

128

1

Output SDF
per point

- Concatenate
- Copy and Concatenate
- FC, Leaky ReLU
- FC, Sigmoid

At test time, pass the
coordinates of each
point in a R^3 grid

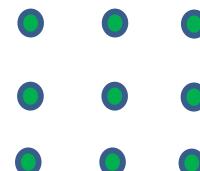
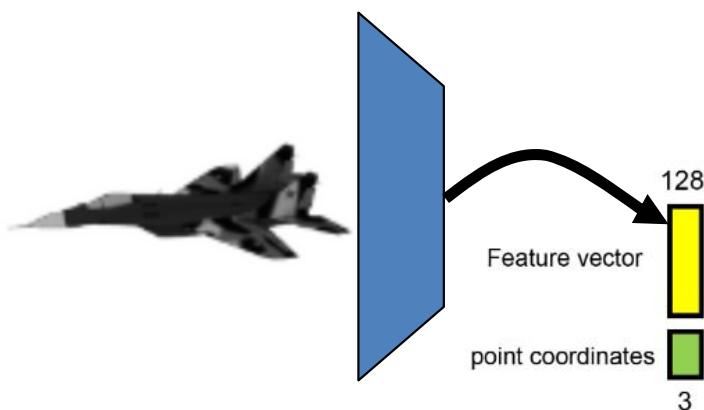
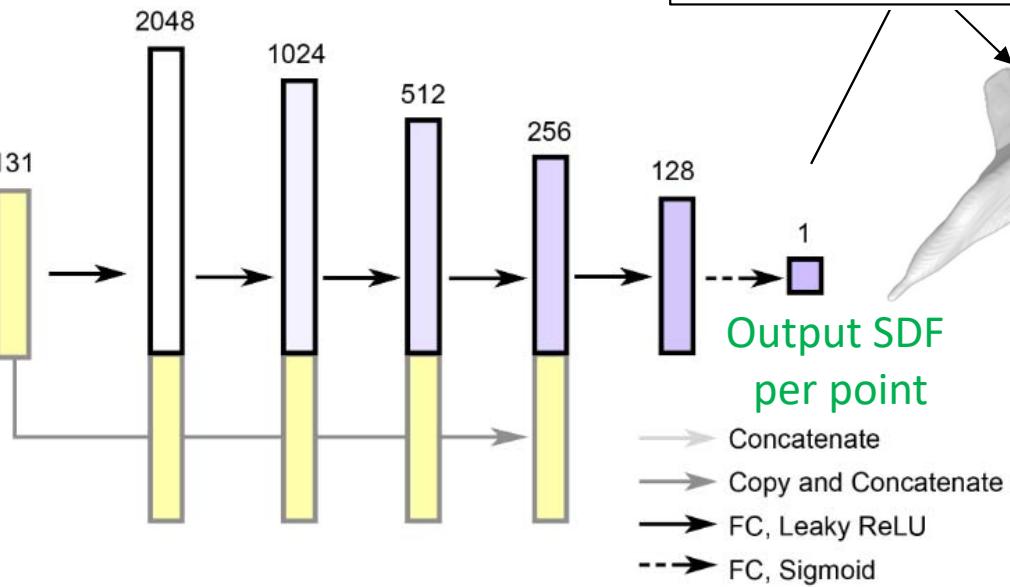


Image-to-shape

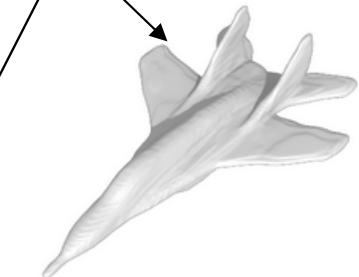
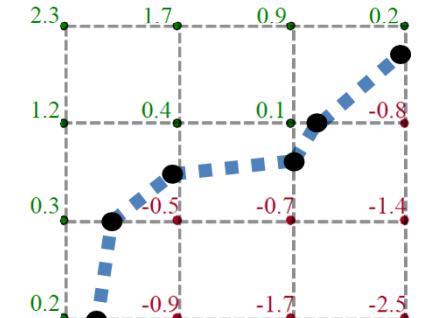
ResNet
encoder



Implicit
Decoder



Marching Cubes!



Training

Given a sample point \mathbf{p}_i' , the network output $f(\mathbf{p}_i')$, and its ground-truth SDF value s_i :

$$L(f(\mathbf{p}_i'), s_i) = | f(\mathbf{p}_i') - s_i |^2$$

One option is to use L2 loss (summed over sample points)

Training

Given a sample point \mathbf{p}_i' , the network output $f(\mathbf{p}_i')$, and its ground-truth SDF value s_i :

$$L(f(\mathbf{p}_i'), s_i) = | f(\mathbf{p}_i') - s_i |^2$$

One option is to use L2 loss (summed over sample points)

To create sample points, take each surface point and perturb it along its normal (e.g., use Gaussian for perturbing)

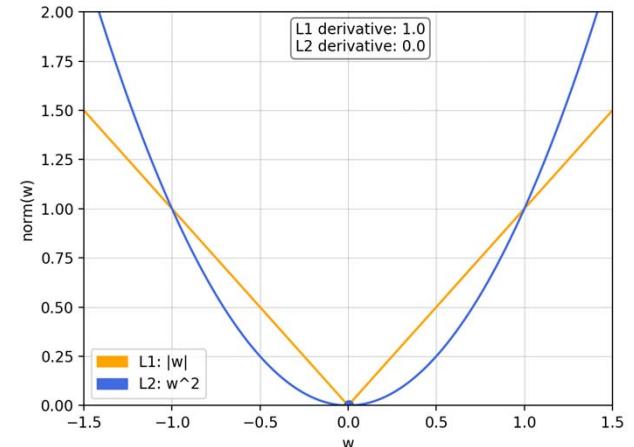
$$\mathbf{p}_i' = \mathbf{p}_i + \varepsilon \cdot \mathbf{n}_i$$

(where \mathbf{p}_i is a surface point, and \mathbf{n}_i is its normal)

Training

Another option is to use L1 loss:

$$L(f(\mathbf{p}_i'), s_i) = | f(\mathbf{p}_i') - s_i |$$



Better results are achieved with the use of **clamped L1 loss** (focus on predicting correct SDFs near the surface)

$$L(f(\mathbf{p}_i'), s) = | clamp(f(\mathbf{p}_i'), \sigma) - clamp(s, \sigma) |$$

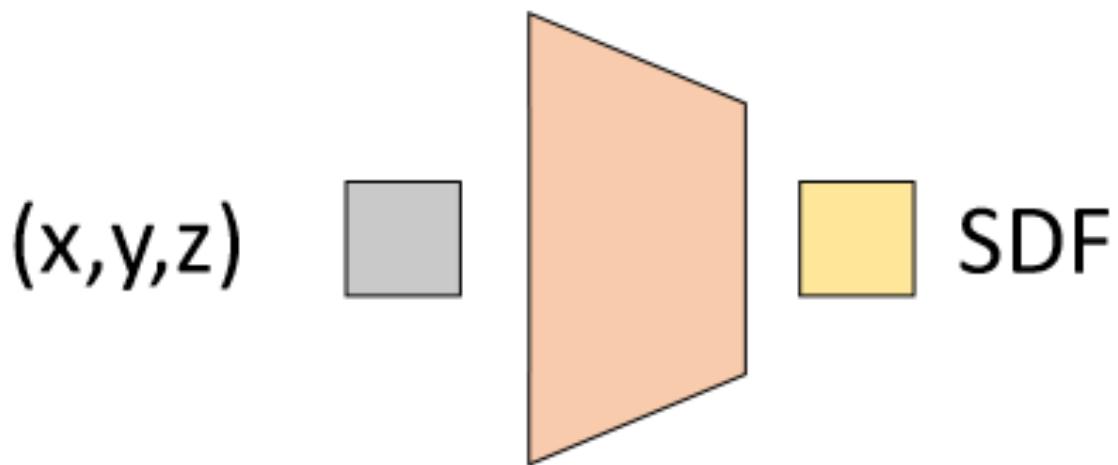
where $clamp(t, \delta) = \min(\delta, \max(-\delta, t))$

Results



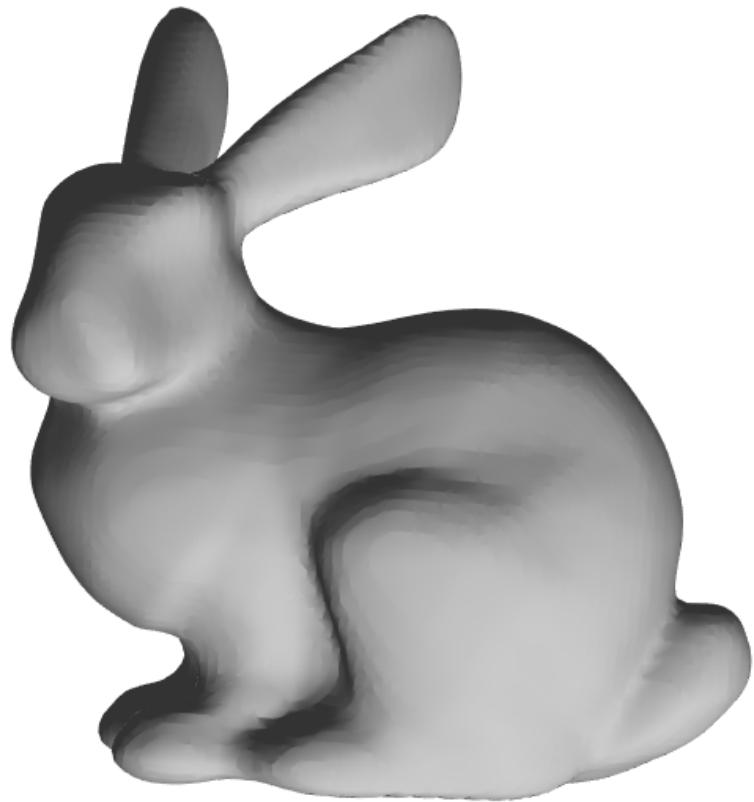
Assignment 4

A simpler version where you train a SDF neural network on a single point cloud (like in assignment 1)

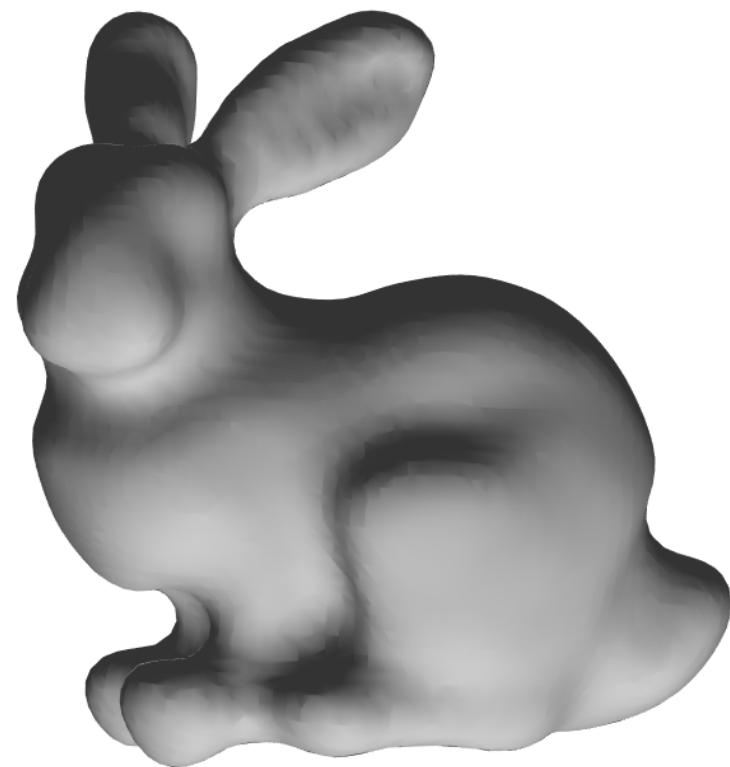


Single Shape DeepSDF

DeepSDF vs MLS

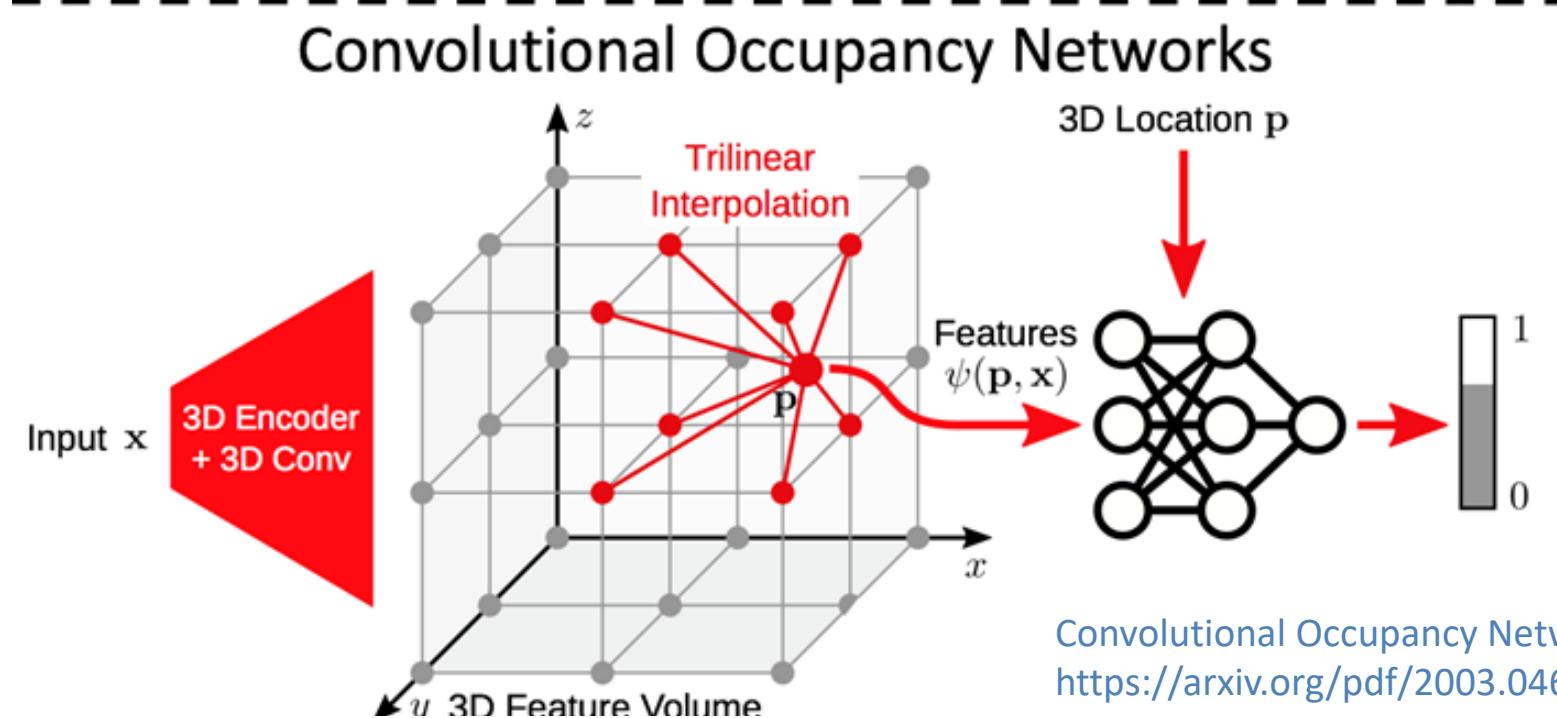
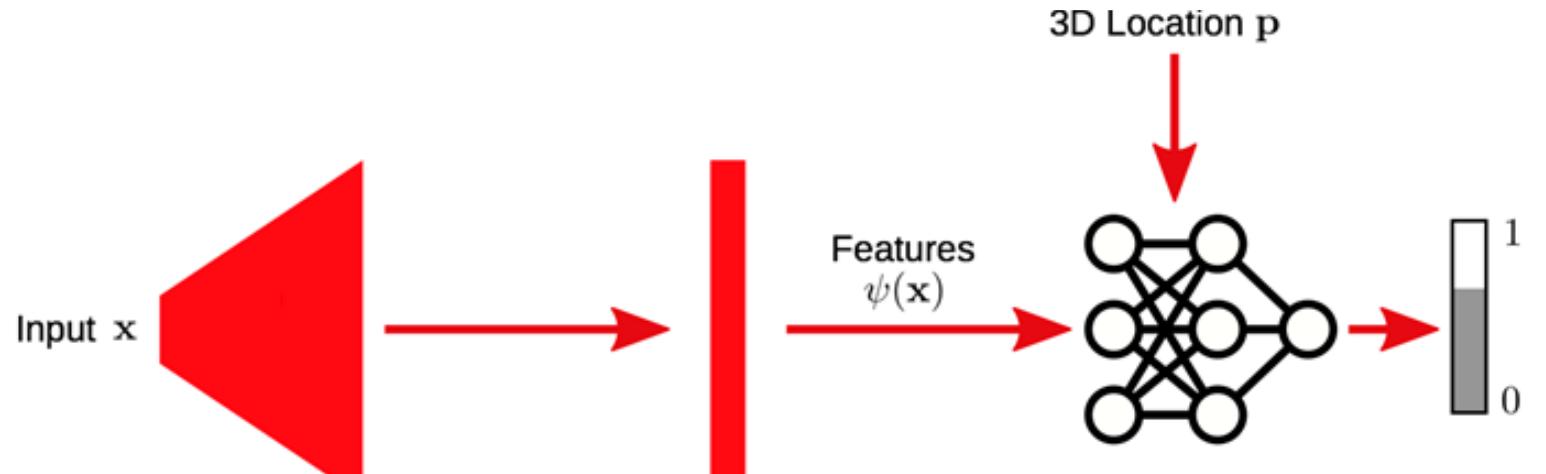


**Single-Shape
DeepSDF**

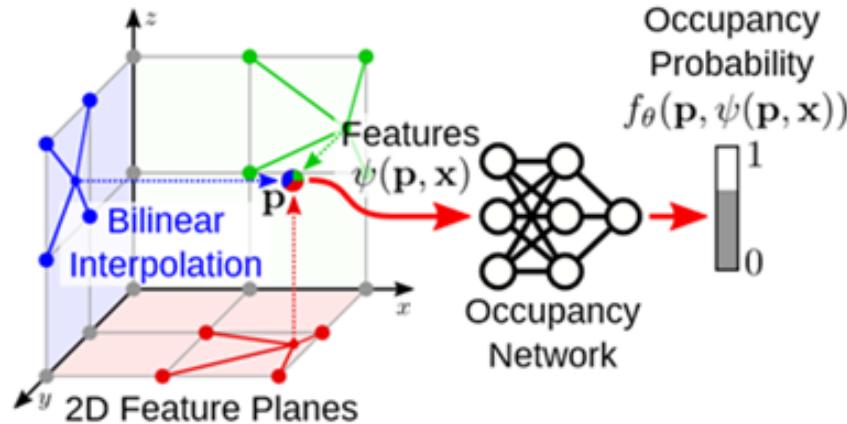


MLS

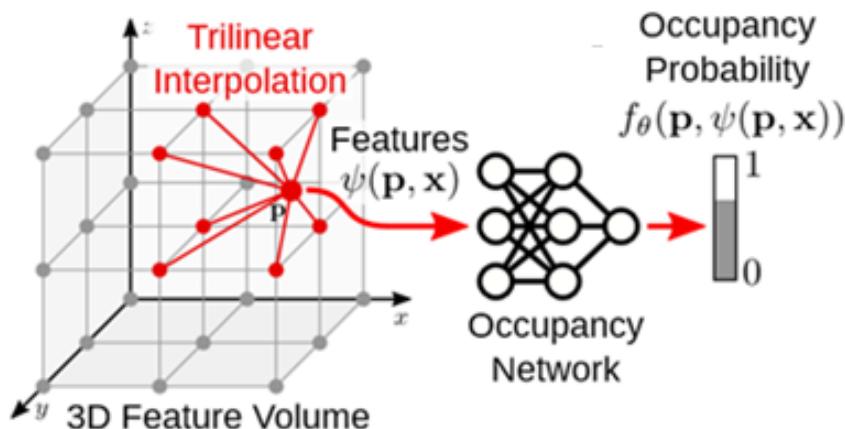
Instead of mapping a 1D **global** vector to per-point implicit, map 3D **local** features to per-point implicit.



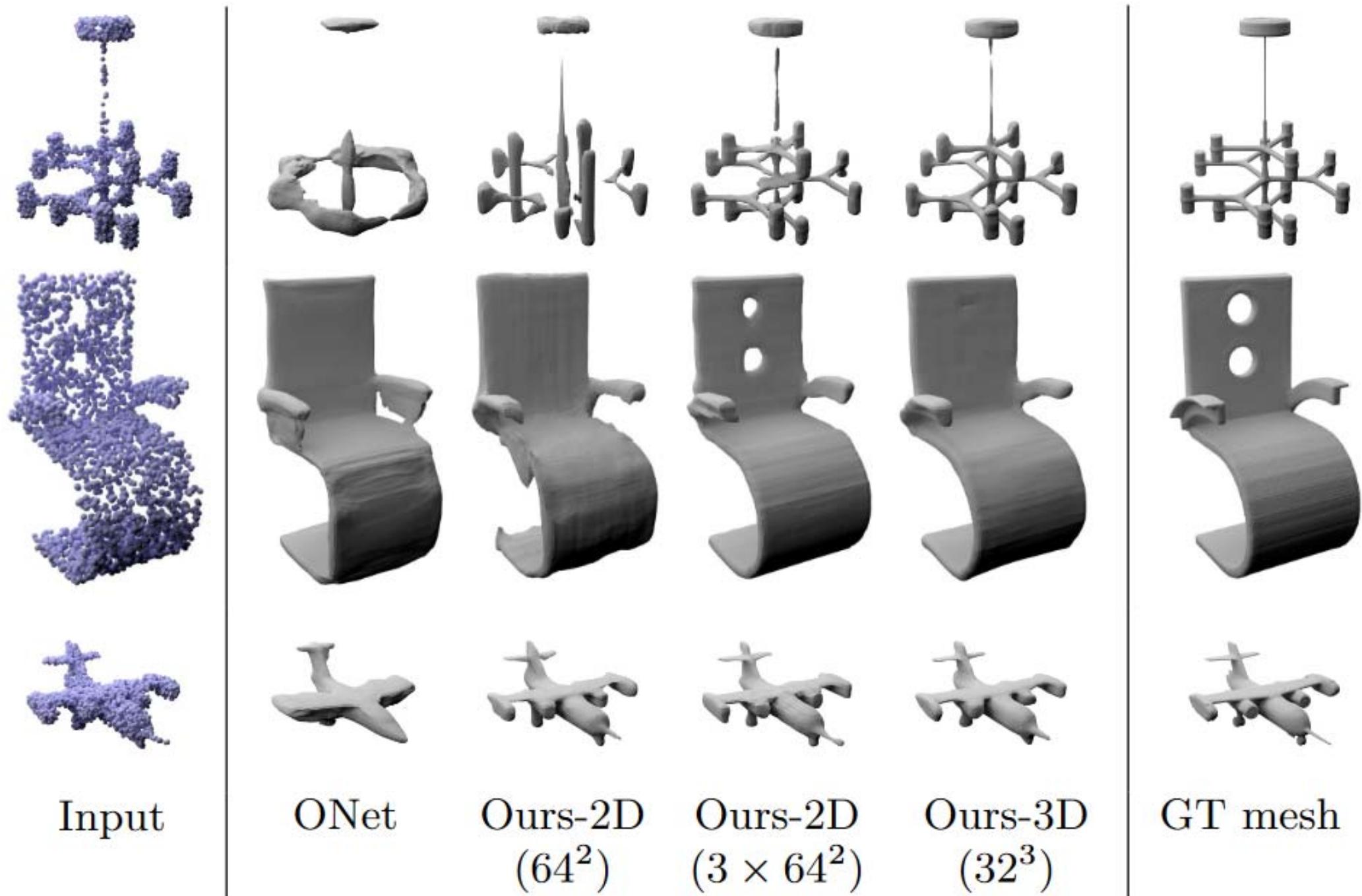
Instead of mapping a 1D **global** vector to per-point implicit, map “multi-plane” **local** features to per-point implicit.



Convolutional Multi-Plane Decoder



Convolutional Volume Decoder



How to generate shapes/scenes?

- **Encoder-Decoders**
 - Case Study: Multi-view decoder
 - Case Study: Implicit decoder
 - **Case Study: Patch decoder**
 - Case Study: Mesh Decoder
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

Image-to-patches

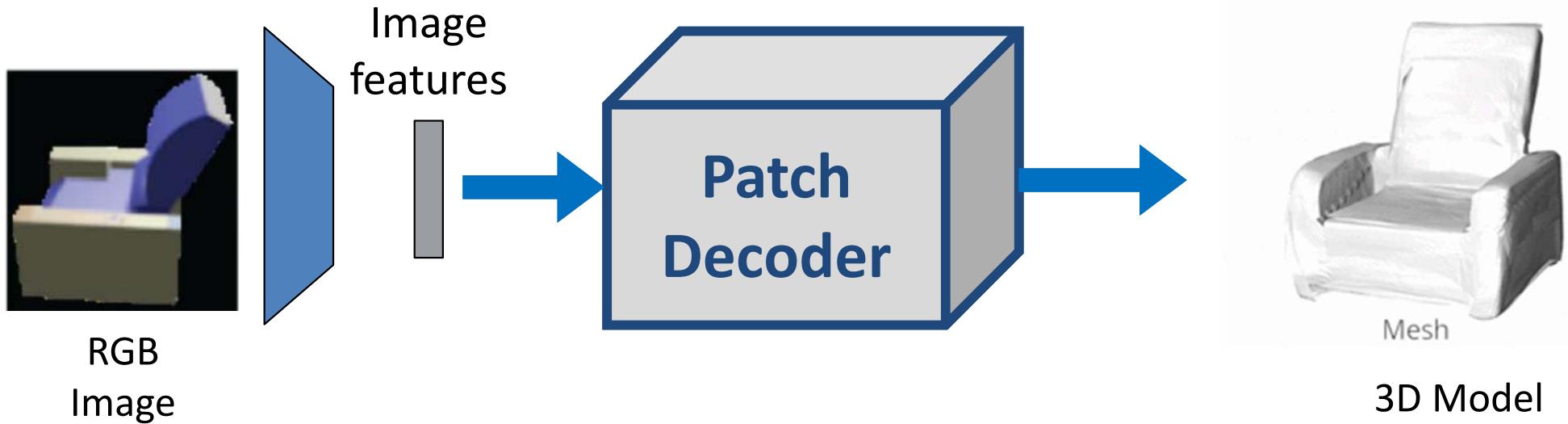


Image-to-patches

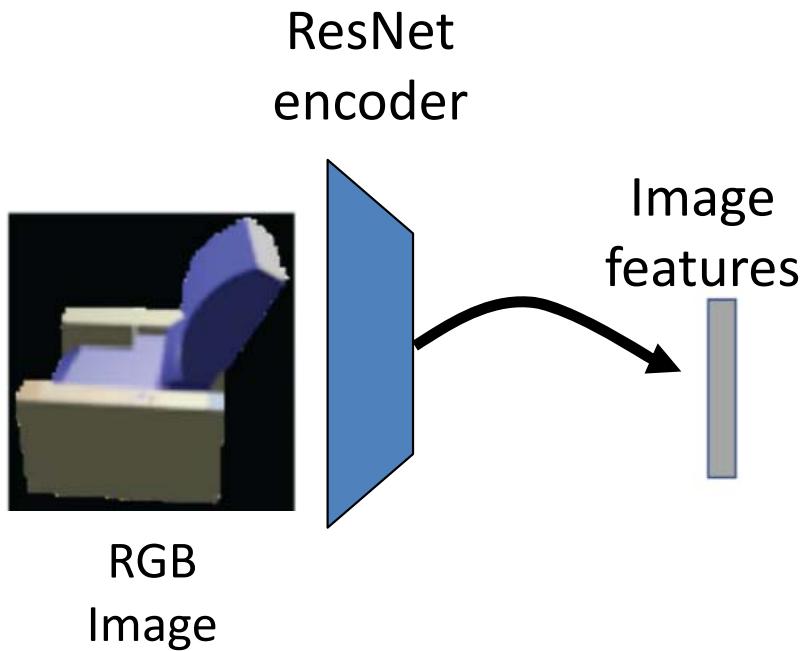


Image-to-patches

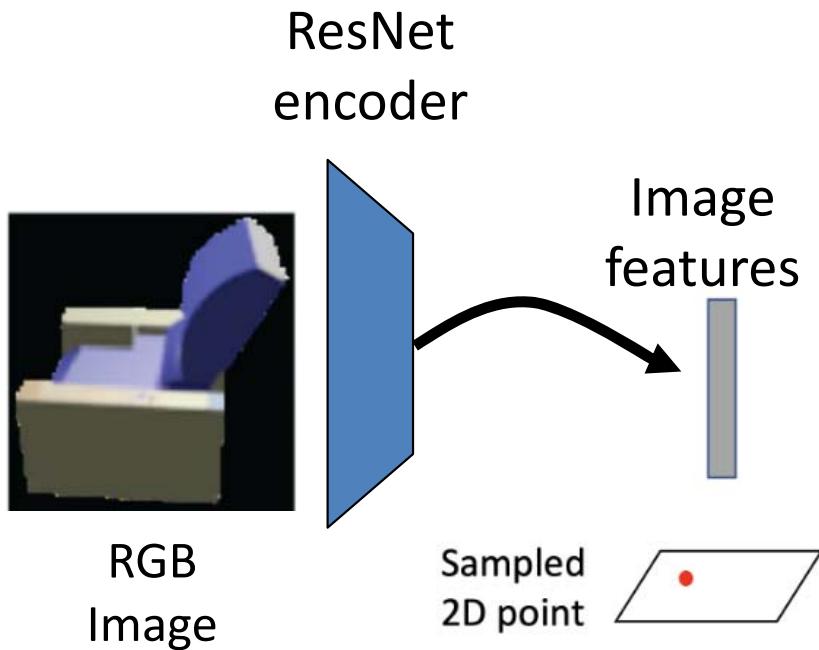
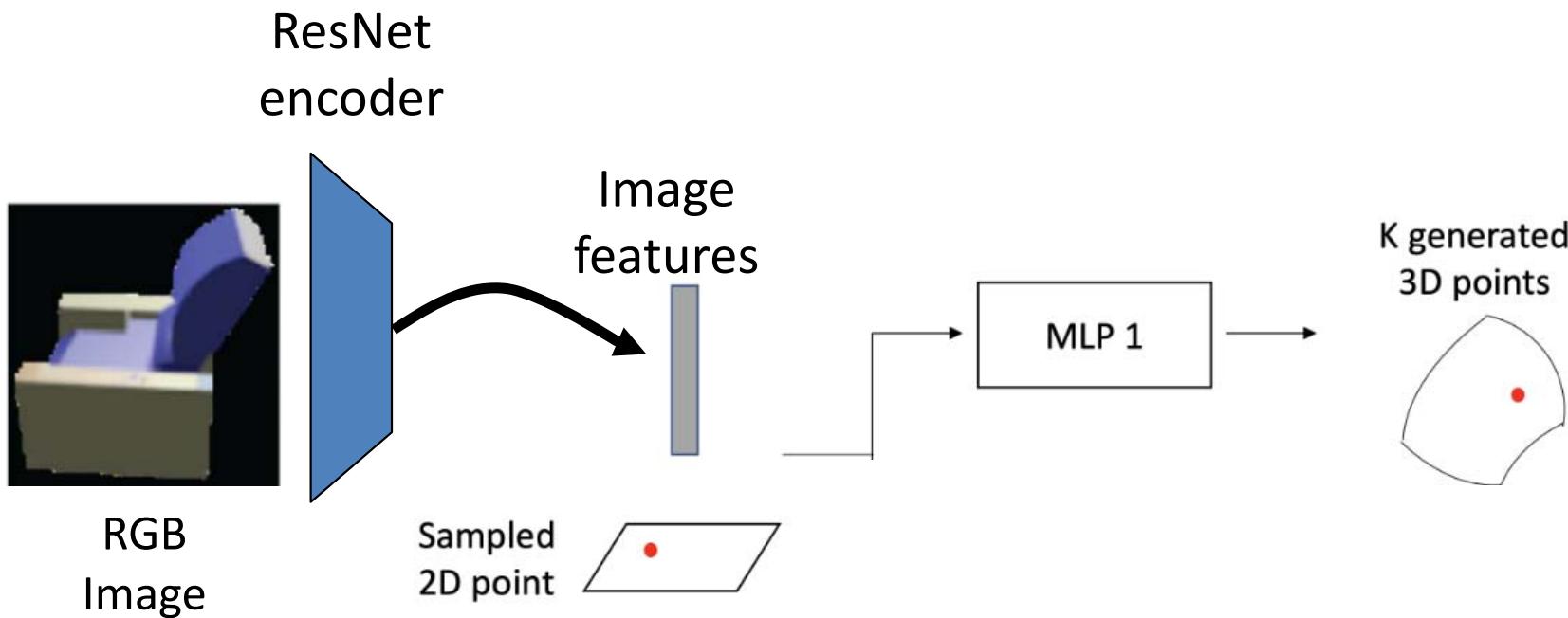
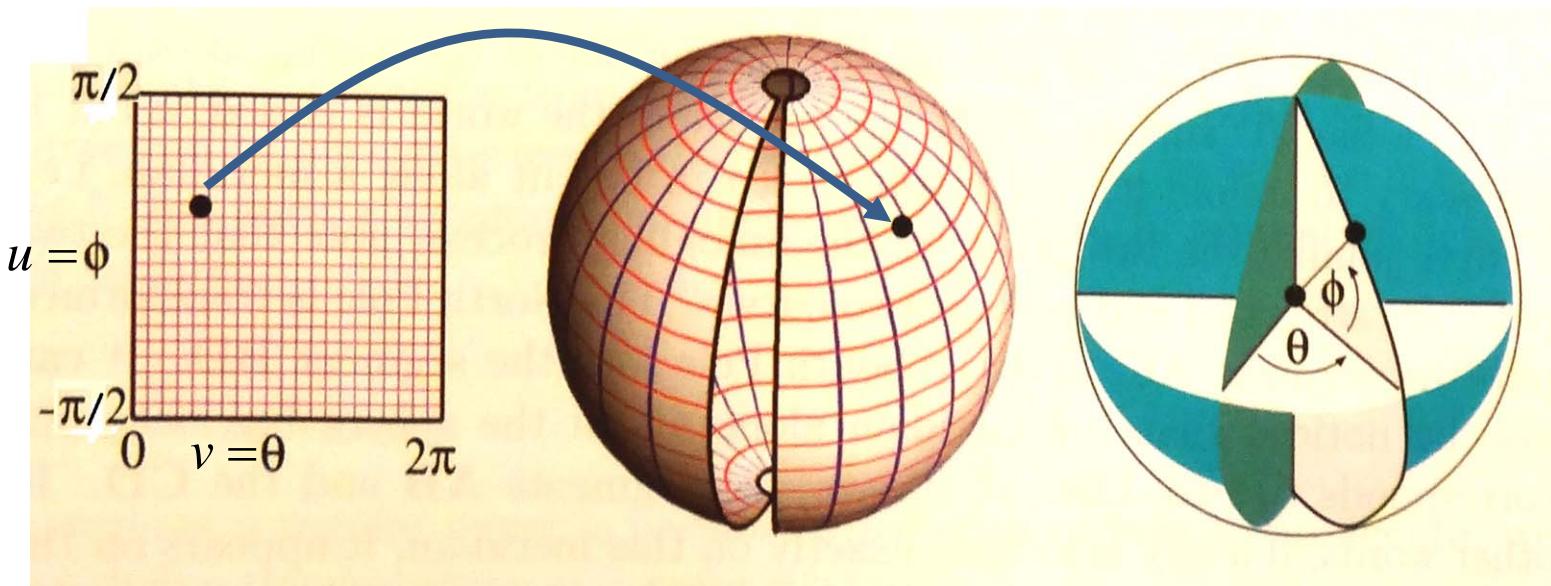


Image-to-patches



Why patches?

A (smooth) surface in 3D can be thought of as a map from a 2D set of parameters to 3D space



$$x(u, v) = R \cos u \sin v$$

$$y(u, v) = R \sin u \sin v$$

$$z(u, v) = R \cos v$$

$$u \in (0, 2\pi), v \in (-\pi/2, \pi/2)$$

Image-to-patches

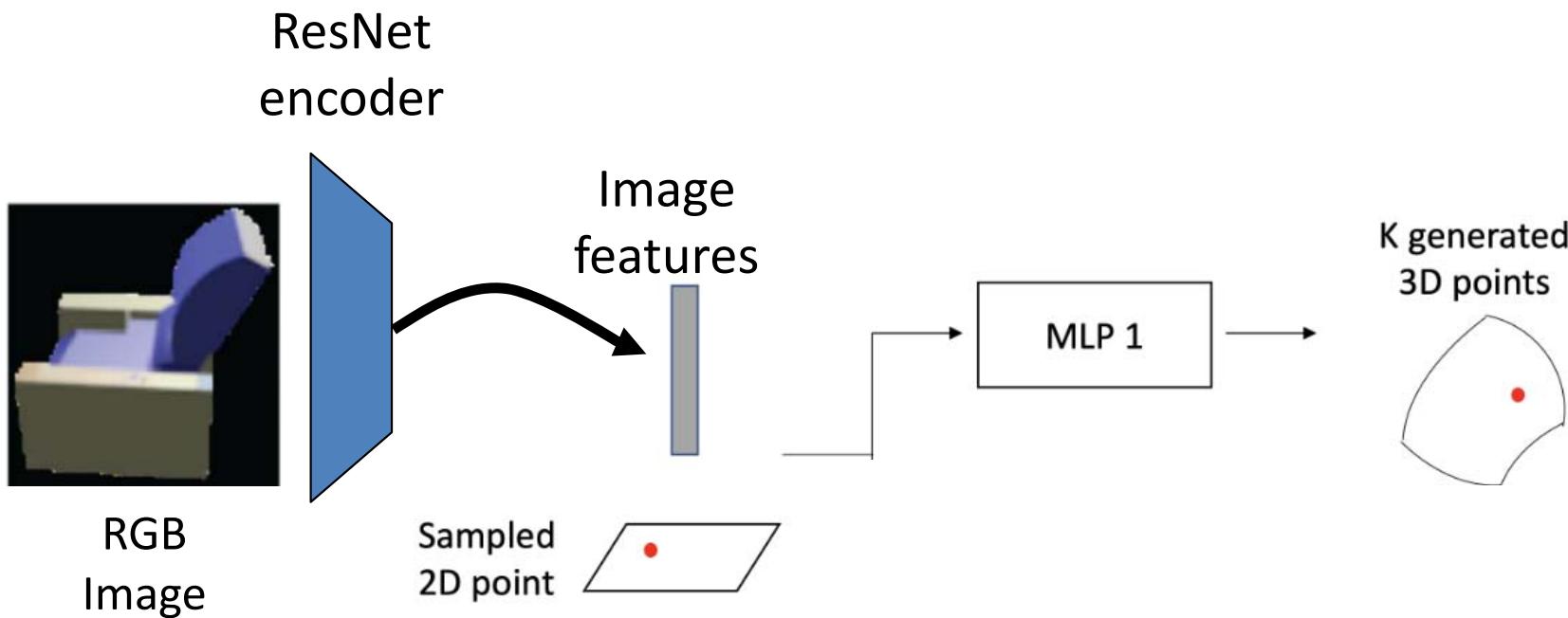


Image-to-patches

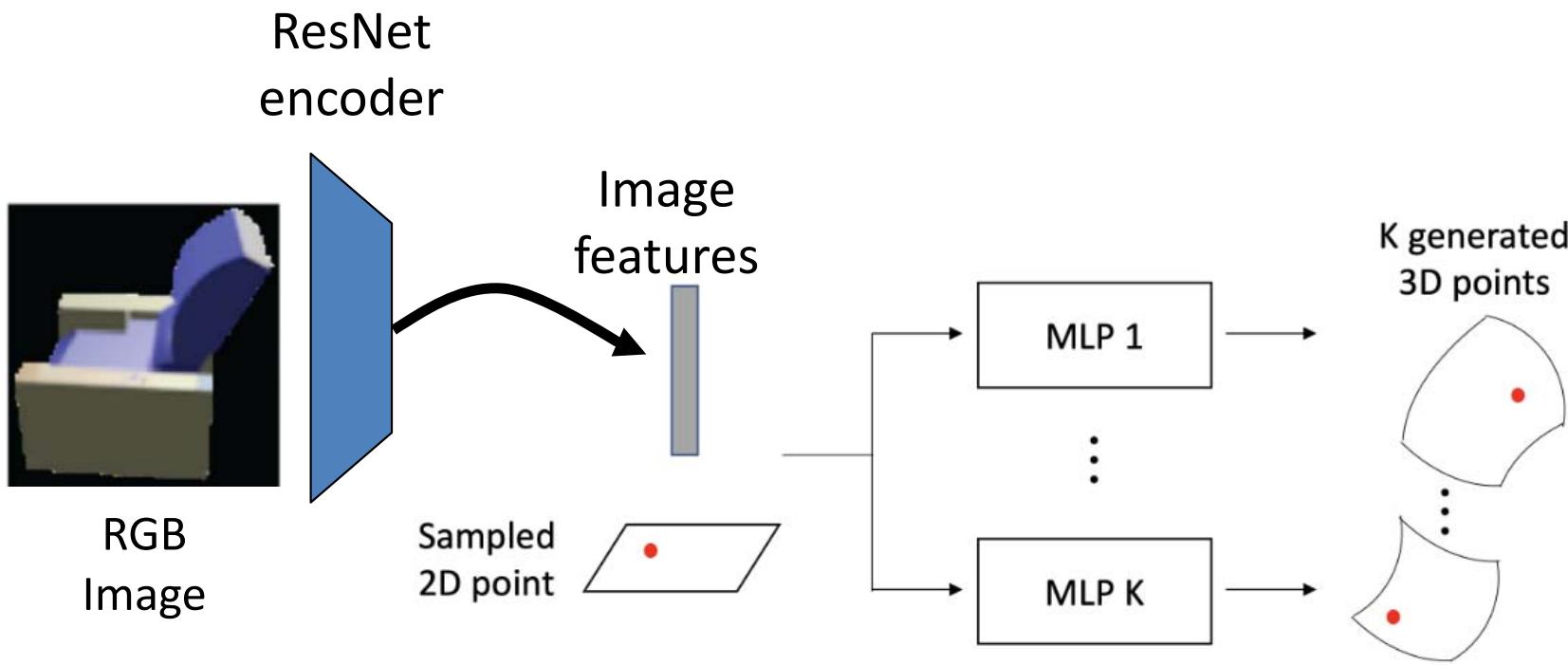


Image-to-patches

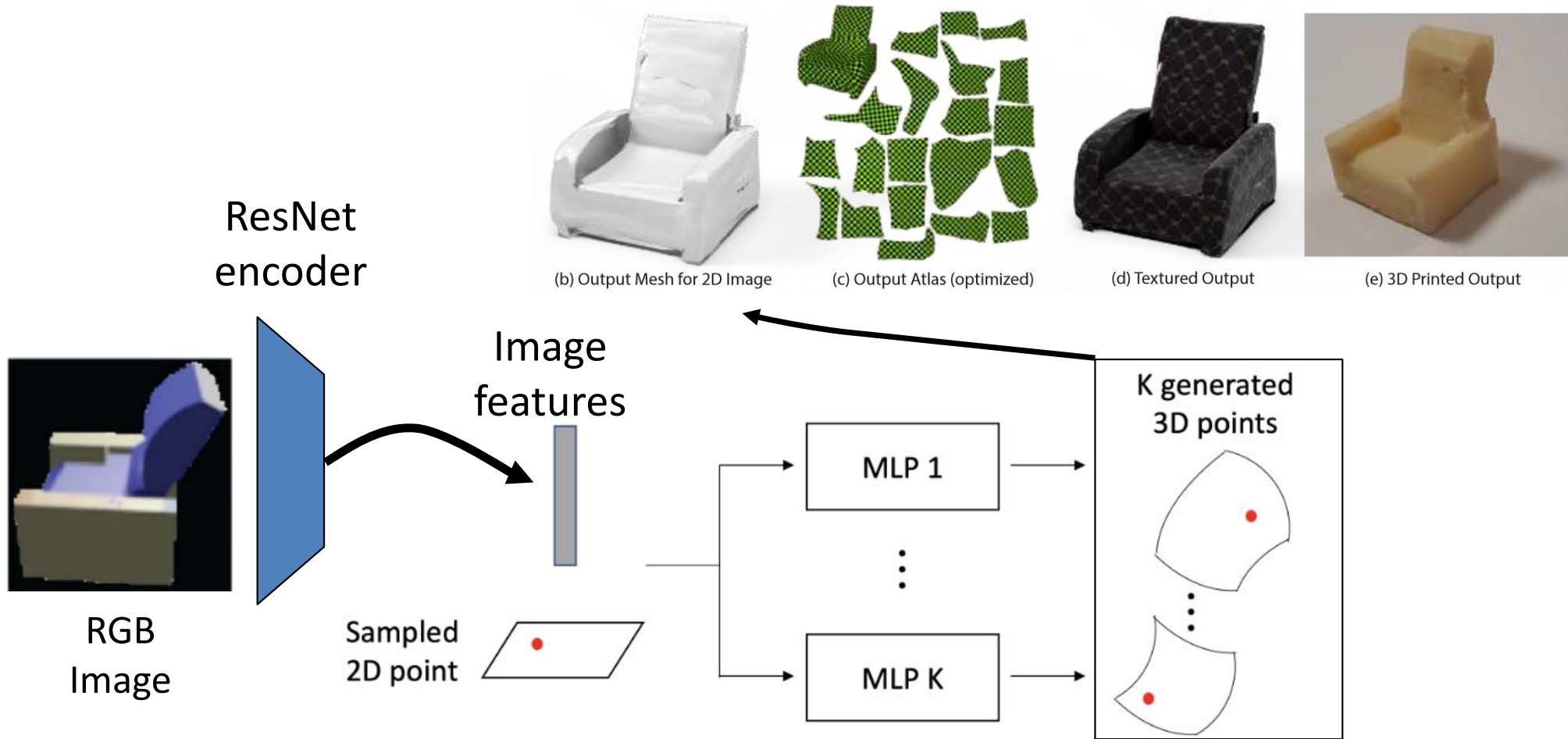
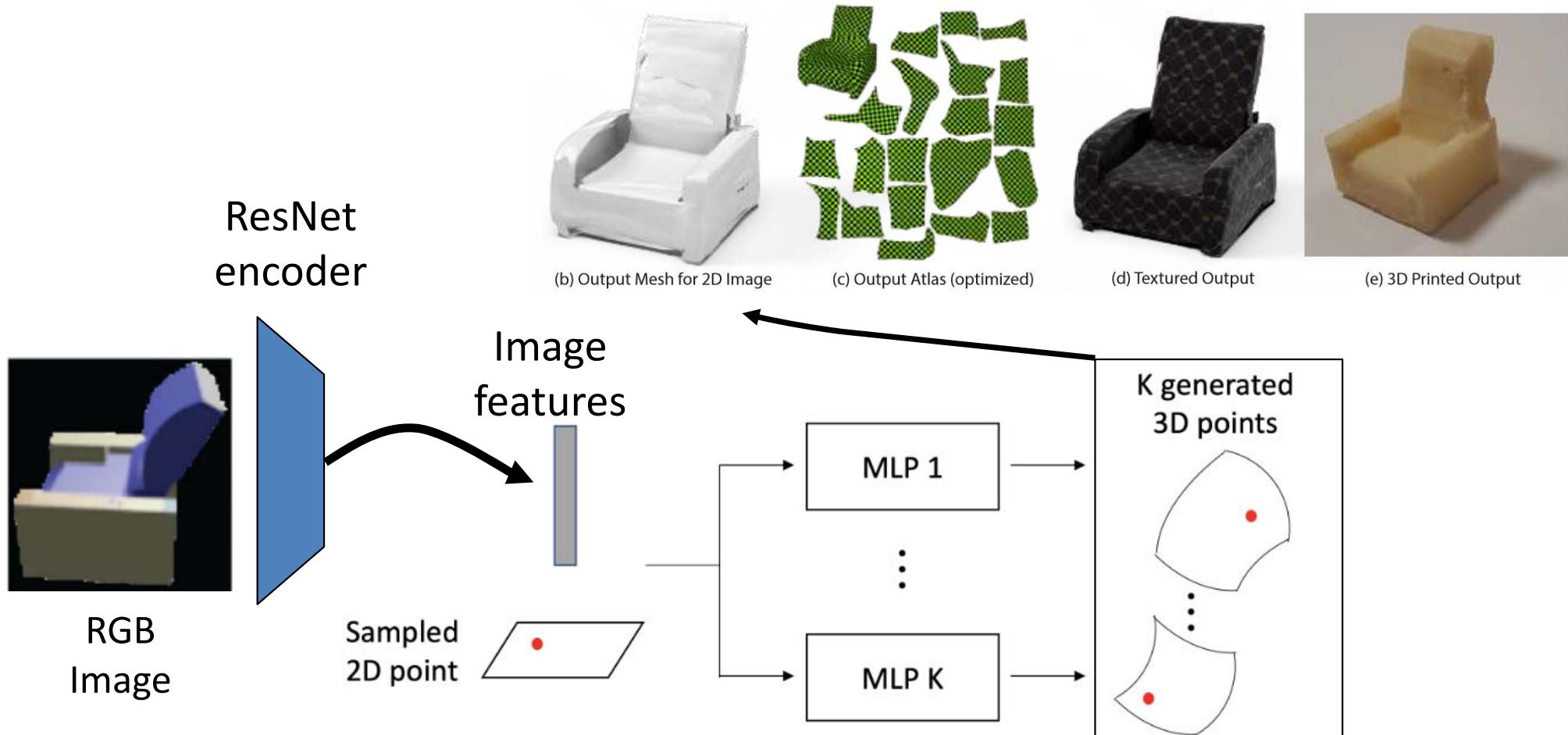


Image-to-patches



Given a ground-truth set T of points (samples from a training mesh) and the set G of generated points from all patches, use the **Chamfer distance** as loss:

$$d(T, G) = \sum_{t \in T} \min_{g \in G} \|t - g\| + \sum_{g \in G} \min_{t \in T} \|t - g\|$$

Results



RGB
Image input

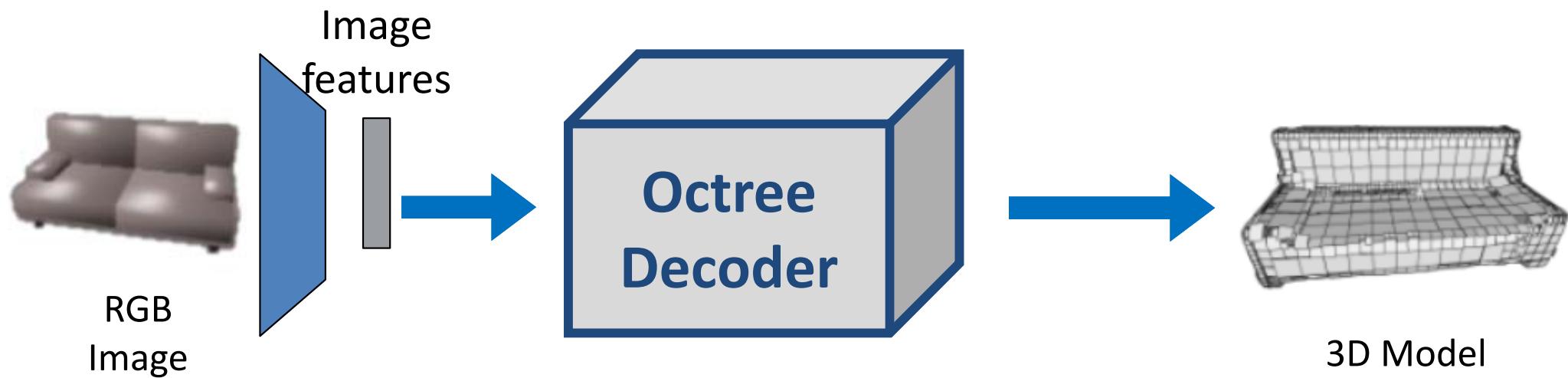
Voxel Network
(3DR2N2)

AtlasNet

How to generate shapes/scenes?

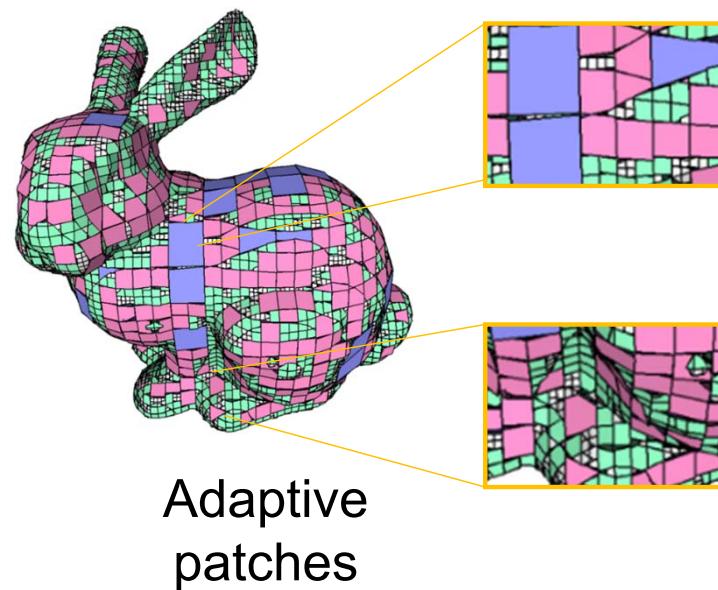
- **Encoder-Decoders**
 - Case Study: Multi-view decoder
 - Case Study: Implicit decoder
 - **Case Study: Patch decoder *within an octree***
 - Case Study: Mesh Decoder
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

Image-to-patches



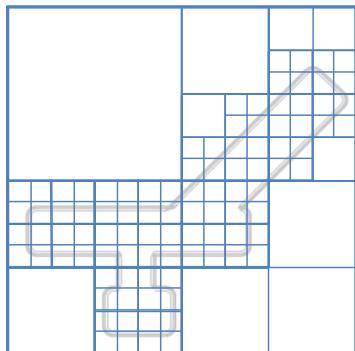
Adaptive O-CNN: key idea

Represent the input shape with adaptive planar patches



Adaptive O-CNN: key idea

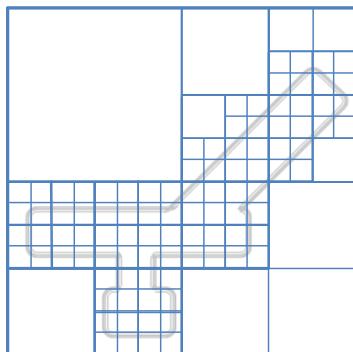
Patches are created for each cell in a **generated octree**



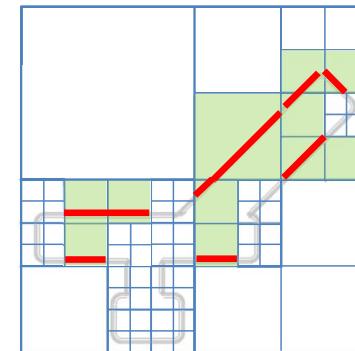
Octree:
subdivide if non-empty

Adaptive O-CNN: key idea

Patches are created for each cell in a **generated octree**

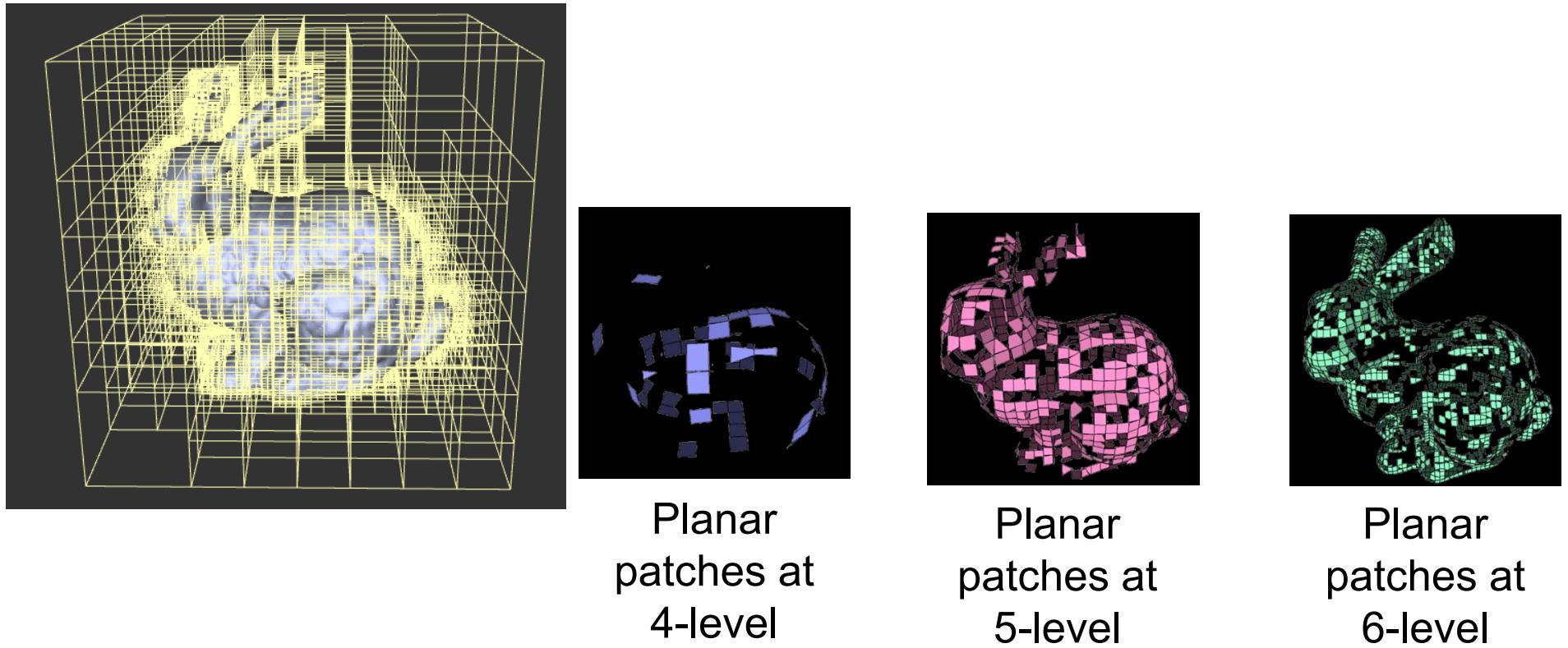


Octree:
subdivide if non-empty

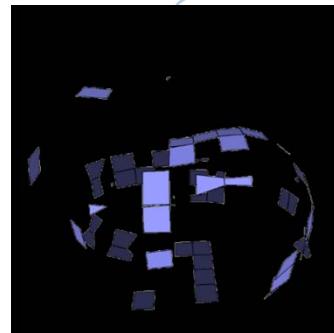
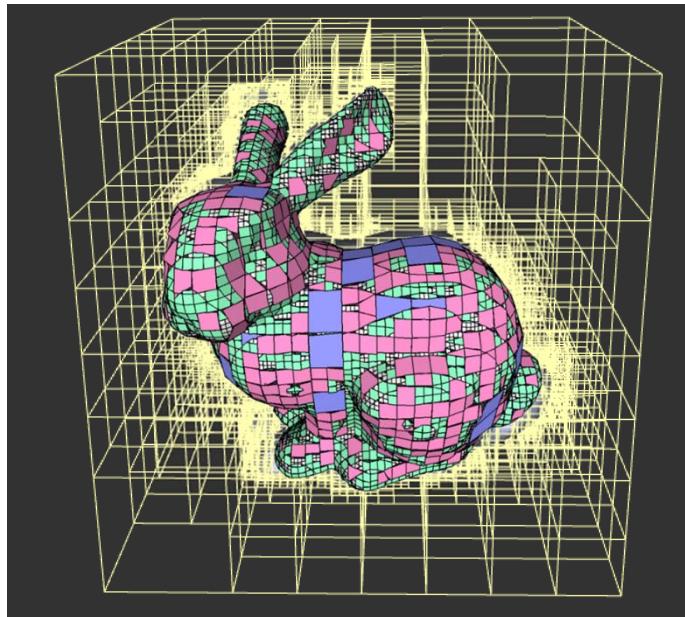


Adaptive octree:
subdivide if non-empty &&
surface-poorly-approximated

Patch-Based Adaptive Octree



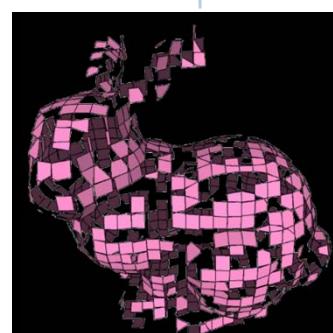
Patch-Based Adaptive Octree



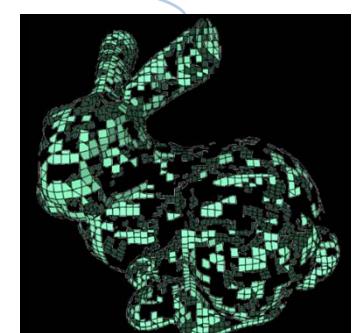
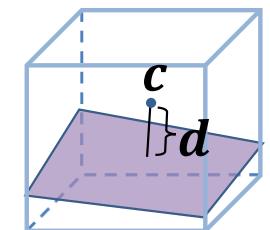
Planar
patches at
4-level

Planar patches

$$\mathbf{n} \cdot (\mathbf{x} - \mathbf{c}) + \mathbf{d} = 0$$

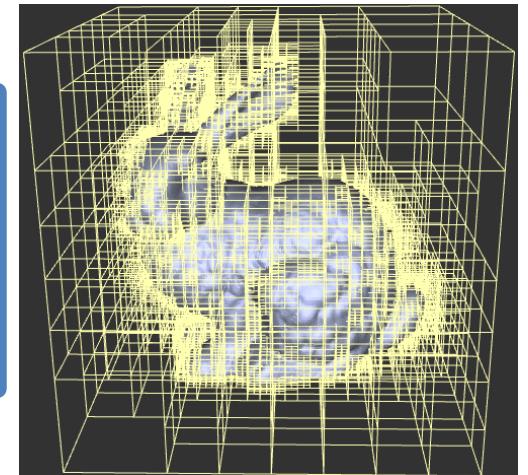
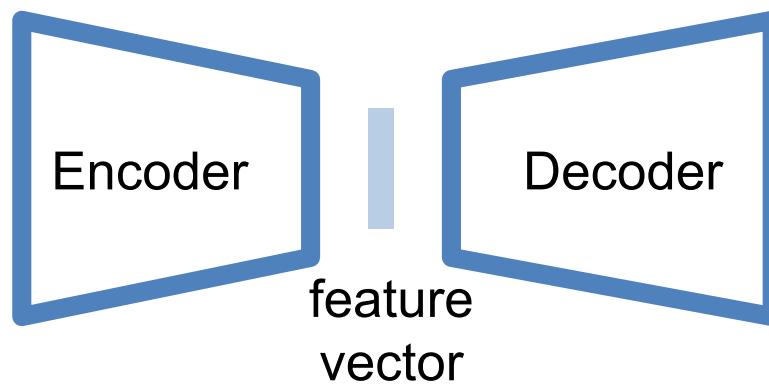


Planar
patches at
5-level



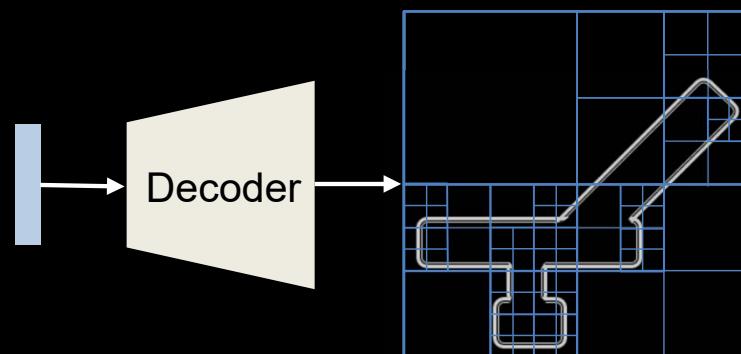
Planar
patches at
6-level

Encoder and Decoder

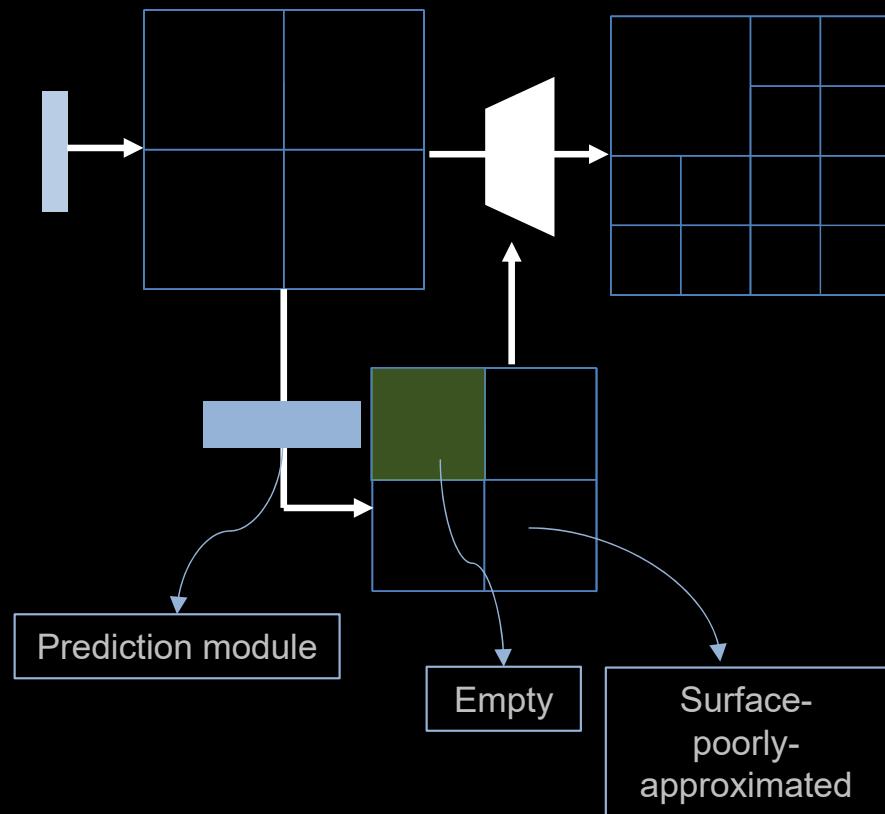


Adaptive O-CNN Decoder

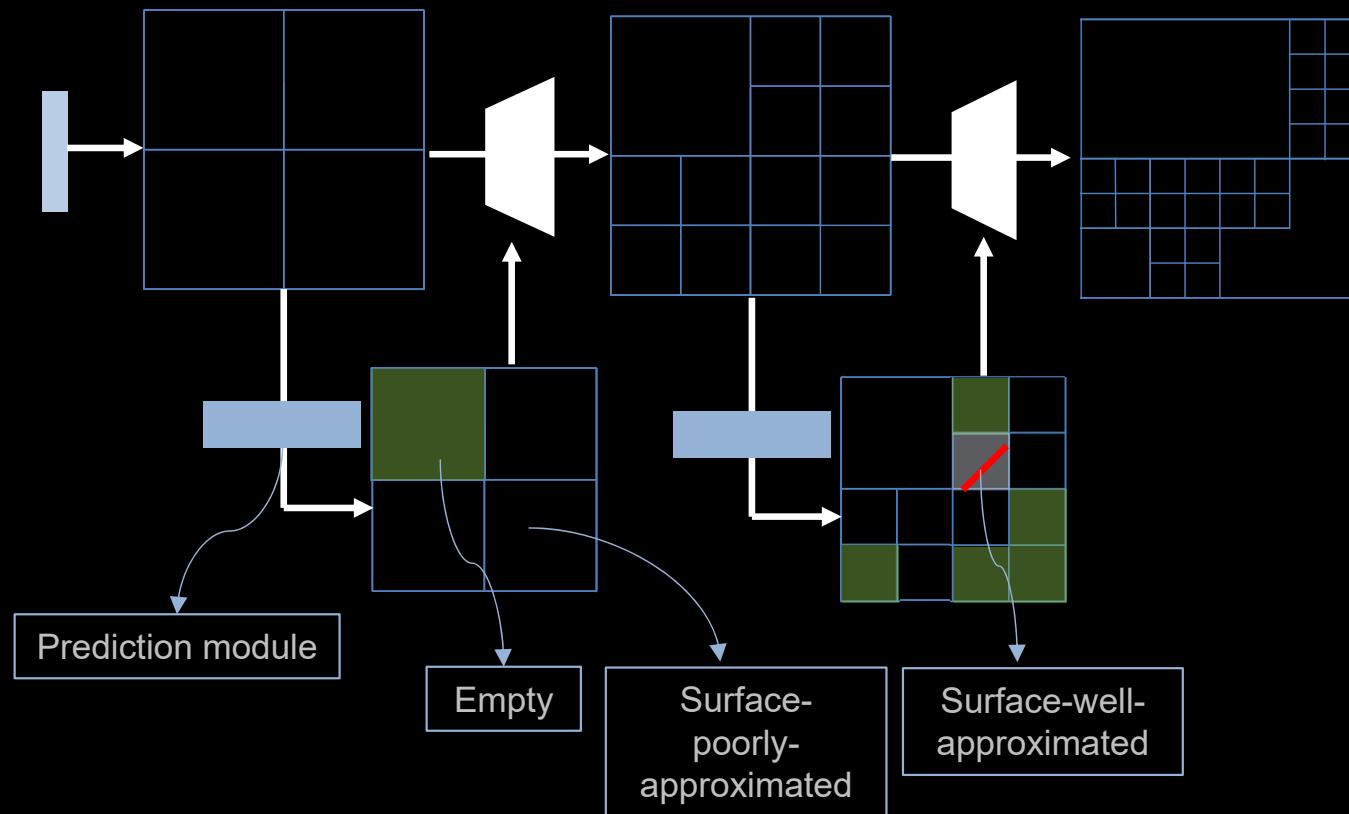
How to generate unknown octree structure?



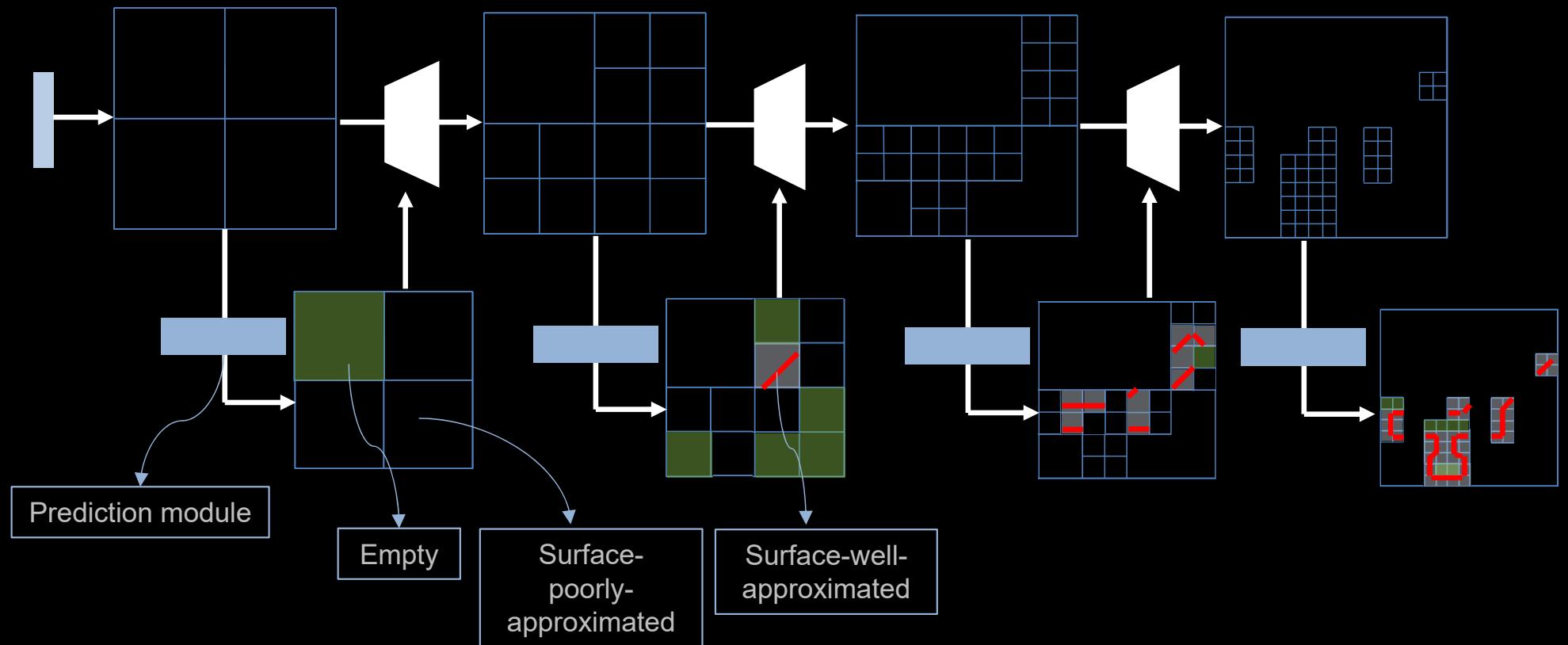
Adaptive O-CNN Decoder



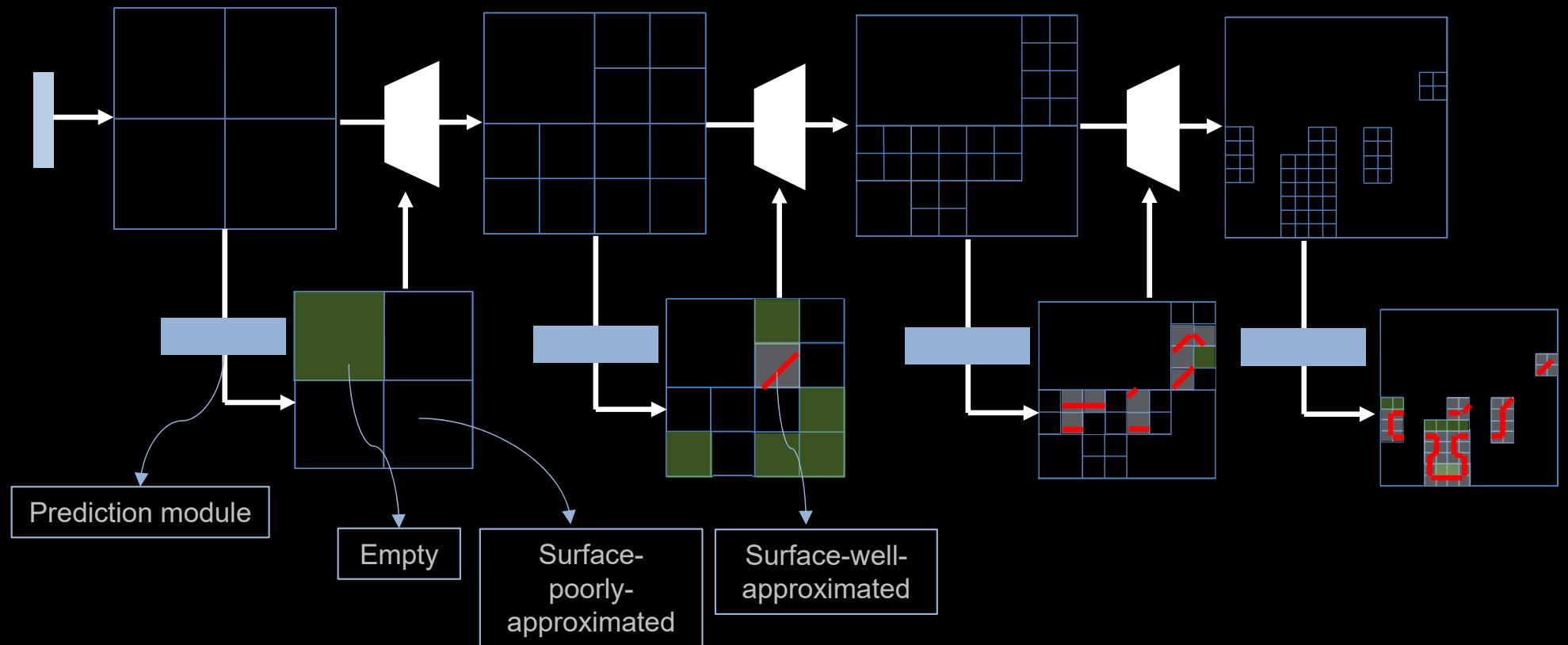
Adaptive O-CNN Decoder



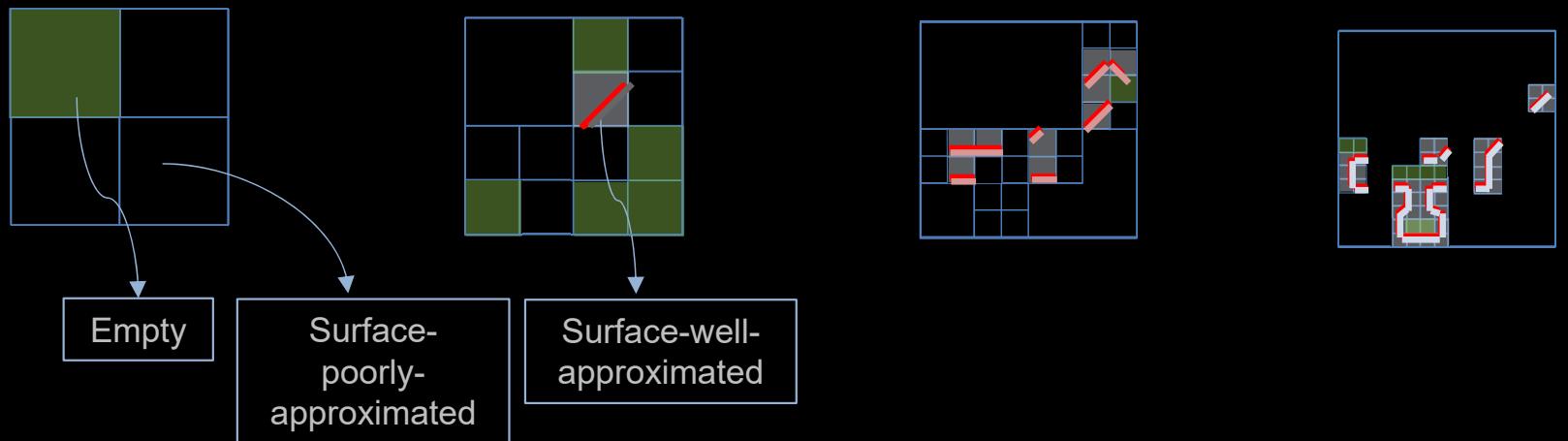
Adaptive O-CNN Decoder



Adaptive O-CNN Decoder

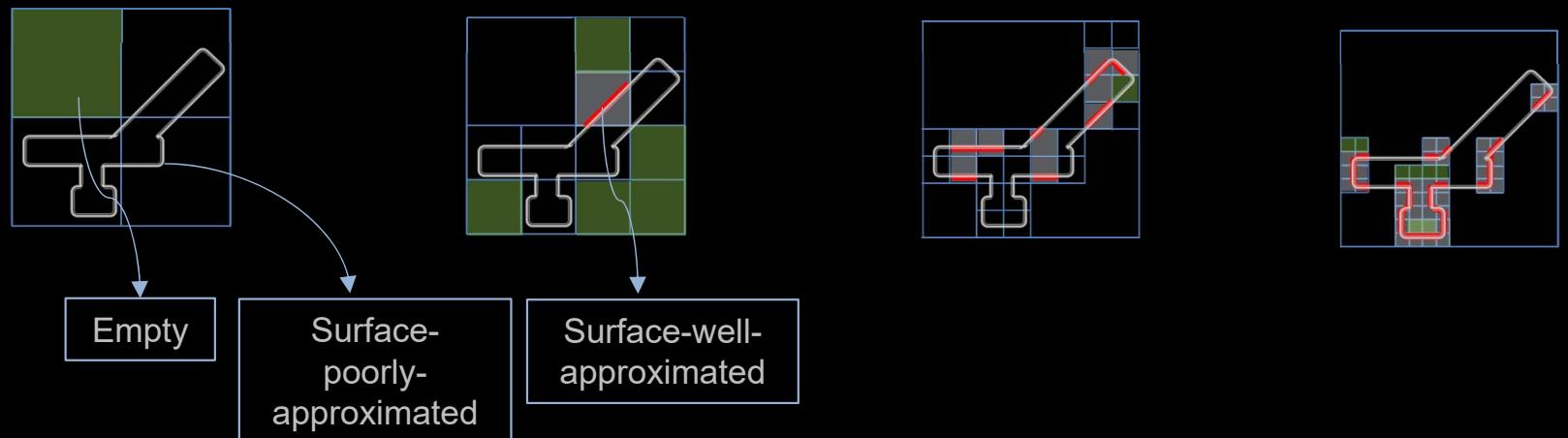


Adaptive O-CNN Decoder

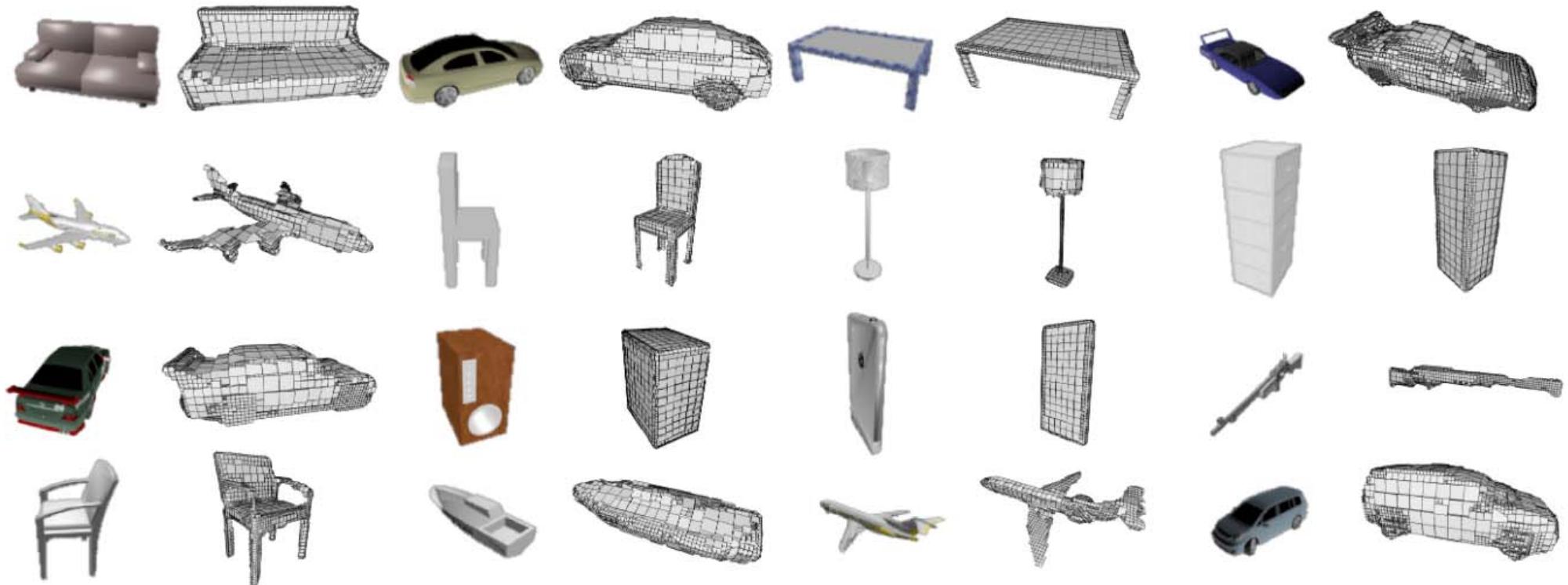


Adaptive O-CNN Decoder: Loss Function

Combine cross-entropy for classifying cells and regression for plane parameters



Results

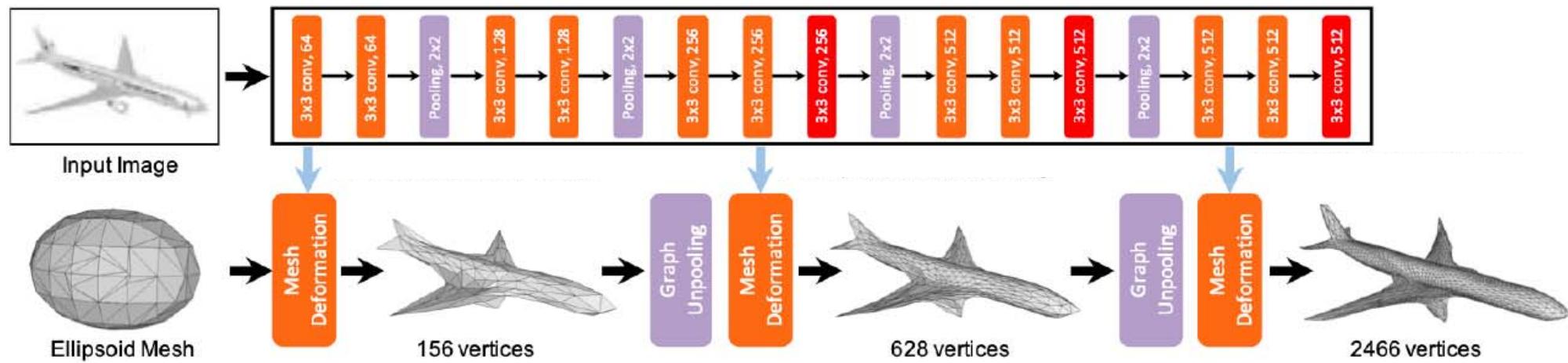


How to generate shapes/scenes?

- **Encoder-Decoders**
 - Case Study: Multi-view decoder
 - Case Study: Implicit decoder
 - Case Study: Patch decoder
 - **Case Study: Mesh Decoder**
- Generative Adversarial Networks
- Variational Autoencoders
- Autoregressive Models
- Diffusion Models

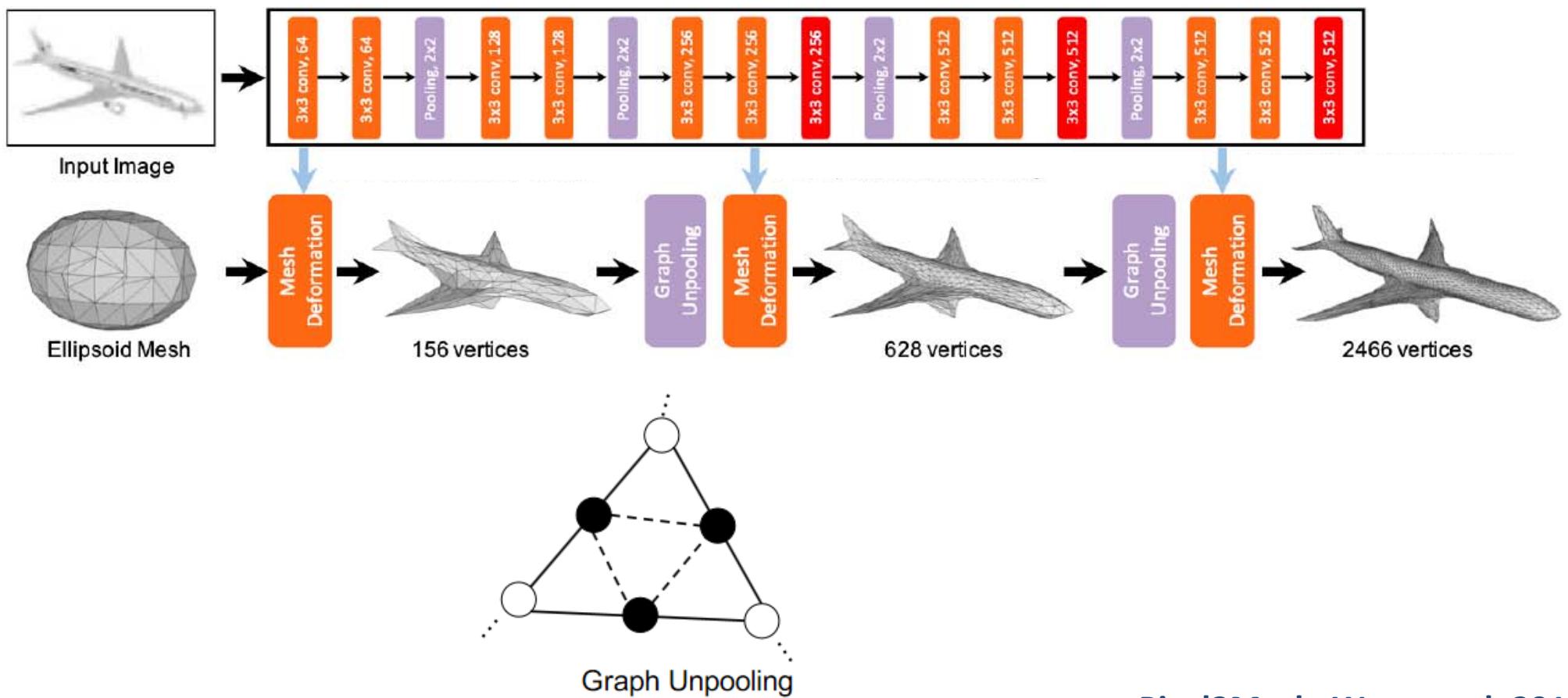
Pixel2Mesh

Learns to deform a template mesh guided by an input image.



Pixel2Mesh

Learns to deform a template mesh guided by an input image.
The decoder predicts vertex positions via graph conv+unpool.



Results



Pixel2Mesh, Wang et al. 2018

See also “Deep Mesh Reconstruction from Single RGB Images via Topology Modification Networks”, Pan et al., 2019 that tries to predict which triangles to remove to replicate holes