



Lecture 9: 3D Generative Models

Li Yi

Apr 17, 2025

Recap

- Problem Definition
- Explicit Algorithms
 - Ball-Pivoting
- Implicit Algorithms
 - Estimate an implicit field function from data
 - Extract the zero iso-surface

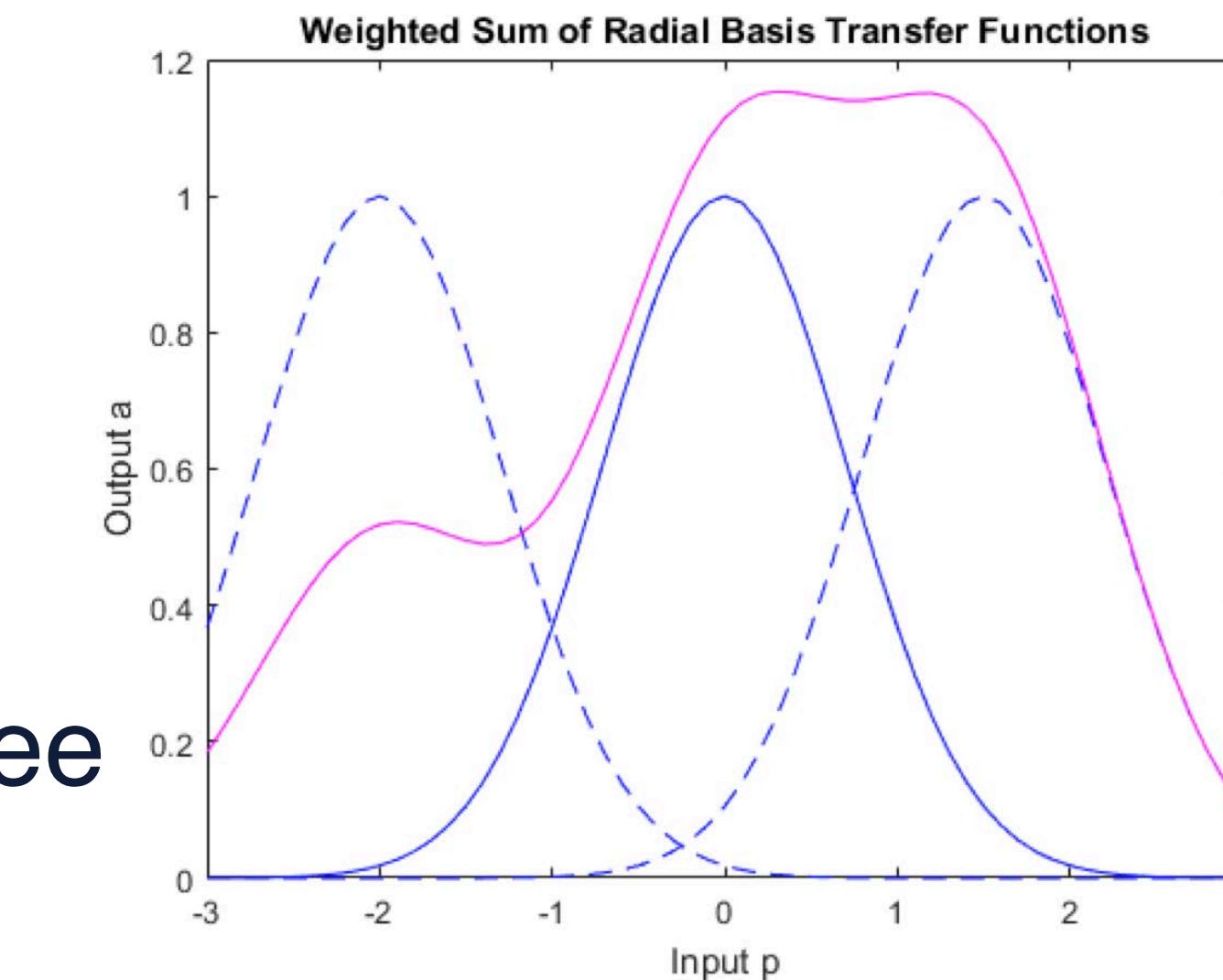
Radial Basis Functions

- Radial basis functions (RBF) $\phi_c(\mathbf{x})$: function value depends only on the distance from a center point c
- Use a weighted sum of radial basis functions to approximate the shape:

$$\phi_c(\mathbf{x}) = \phi(\|\mathbf{x} - \mathbf{c}\|)$$

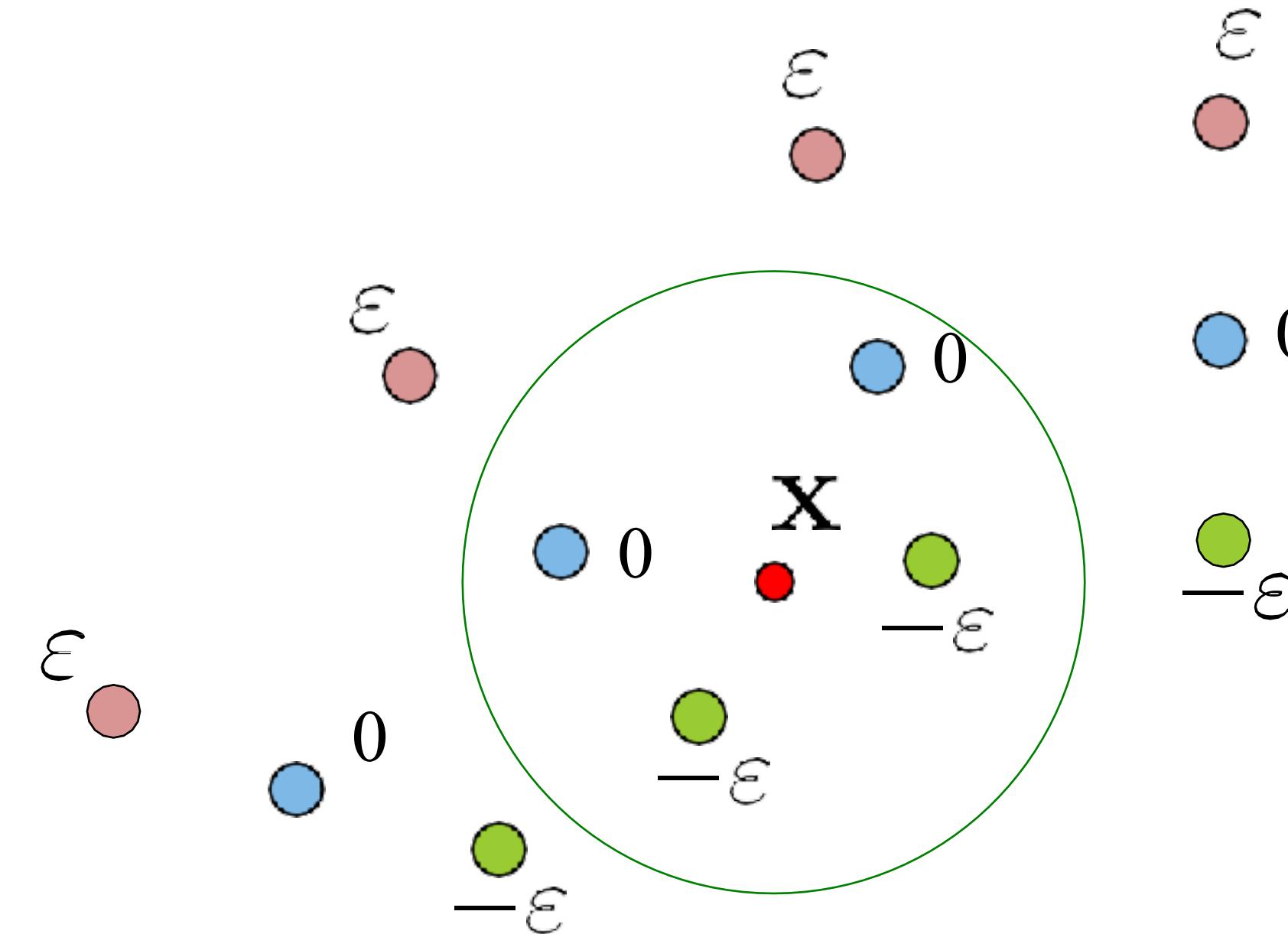
$$f(\mathbf{x}) = \sum_{i=1}^n \omega_i \phi\left(\|\mathbf{x} - \mathbf{x}_i\|\right) + p(x)$$

where p is a polynomial of low degree



Moving Least Squares (MLS)

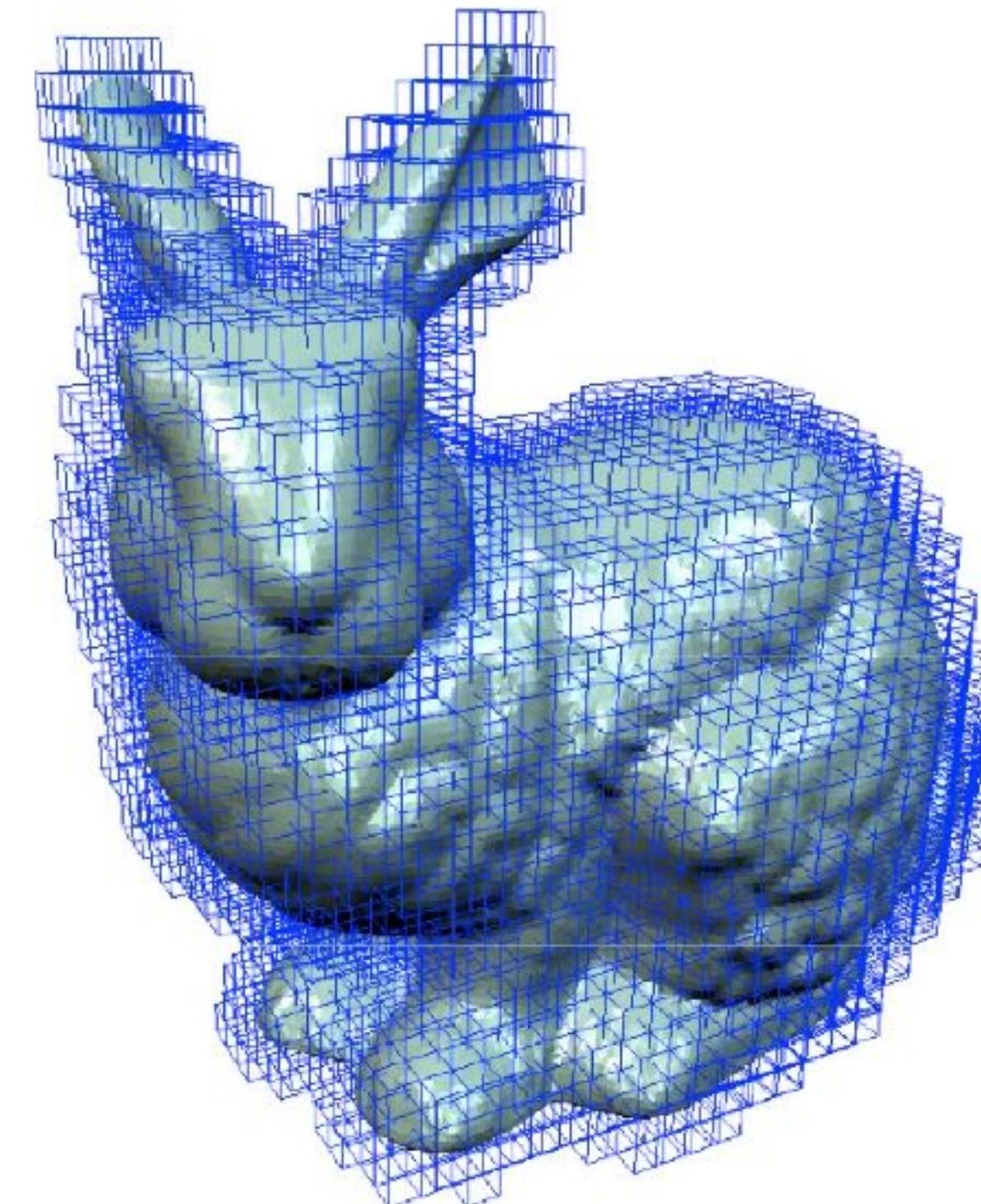
Do purely local approximation of the SDF
Weights change depending on where we are evaluating



Implicit Meshing Algorithm

- Two basic steps:
 1. Estimate an implicit field function from data
 - 2. Extract the zero iso-surface**

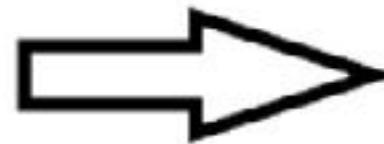
Input: a signed distance field
(Implicitly assumed knowing the
inside/outside of the shape,
often needs to be estimated with
normal information)



Classical Solution: 2D Marching Square

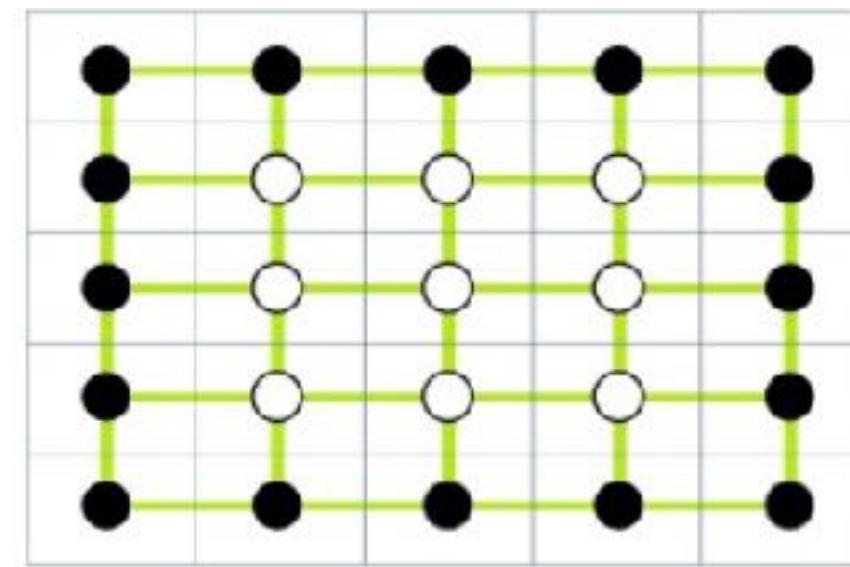
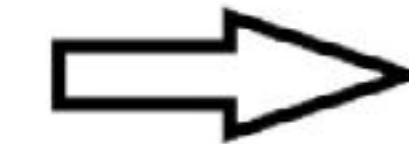
1	1	1	1	1
1	2	3	2	1
1	3	3	3	1
1	2	3	2	1
1	1	1	1	1

Threshold
with iso-value

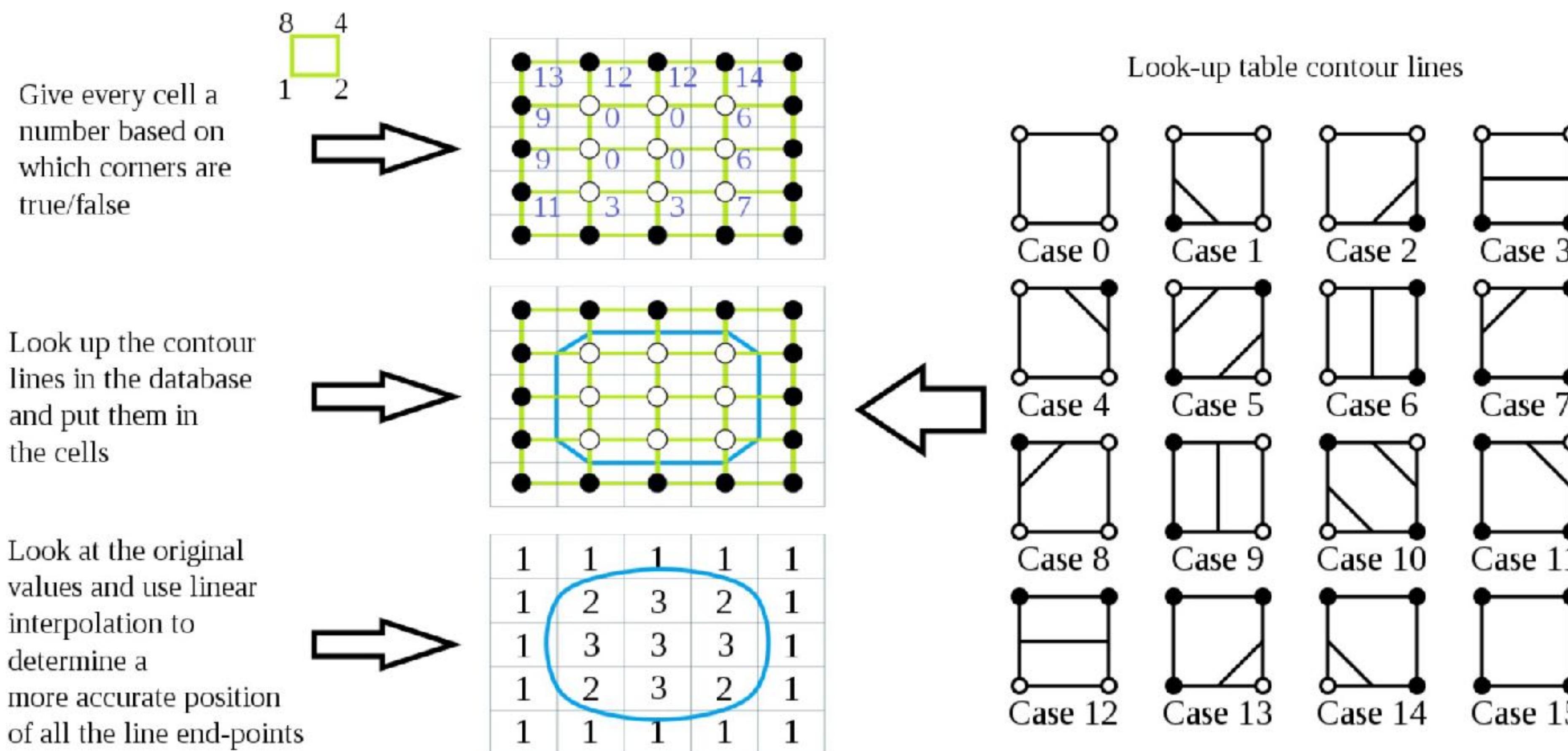


0	0	0	0	0
0	1	1	1	0
0	1	1	1	0
0	1	1	1	0
0	0	0	0	0

Binary image
to cells

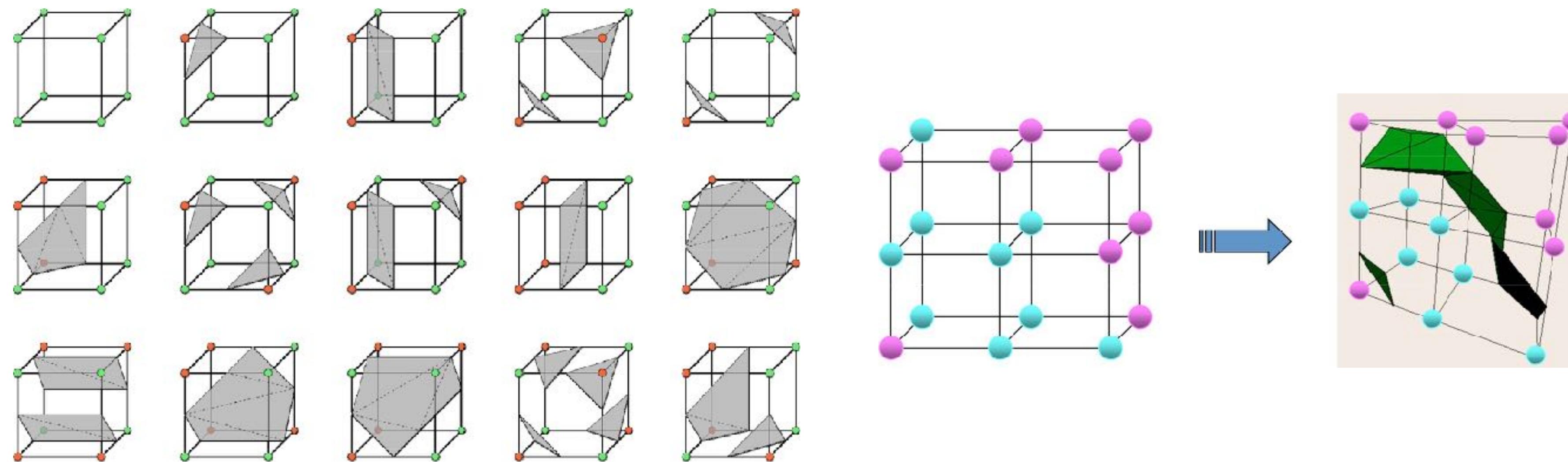


Classical Solution: 2D Marching Square



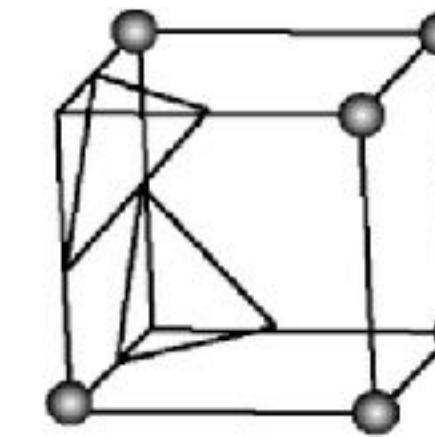
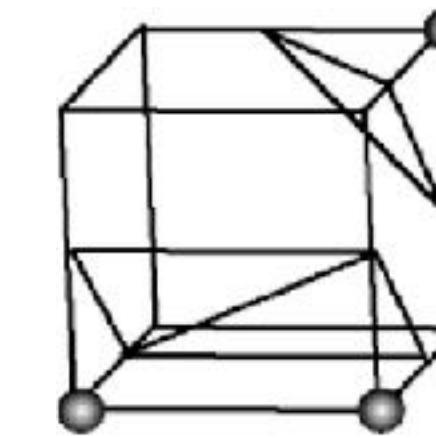
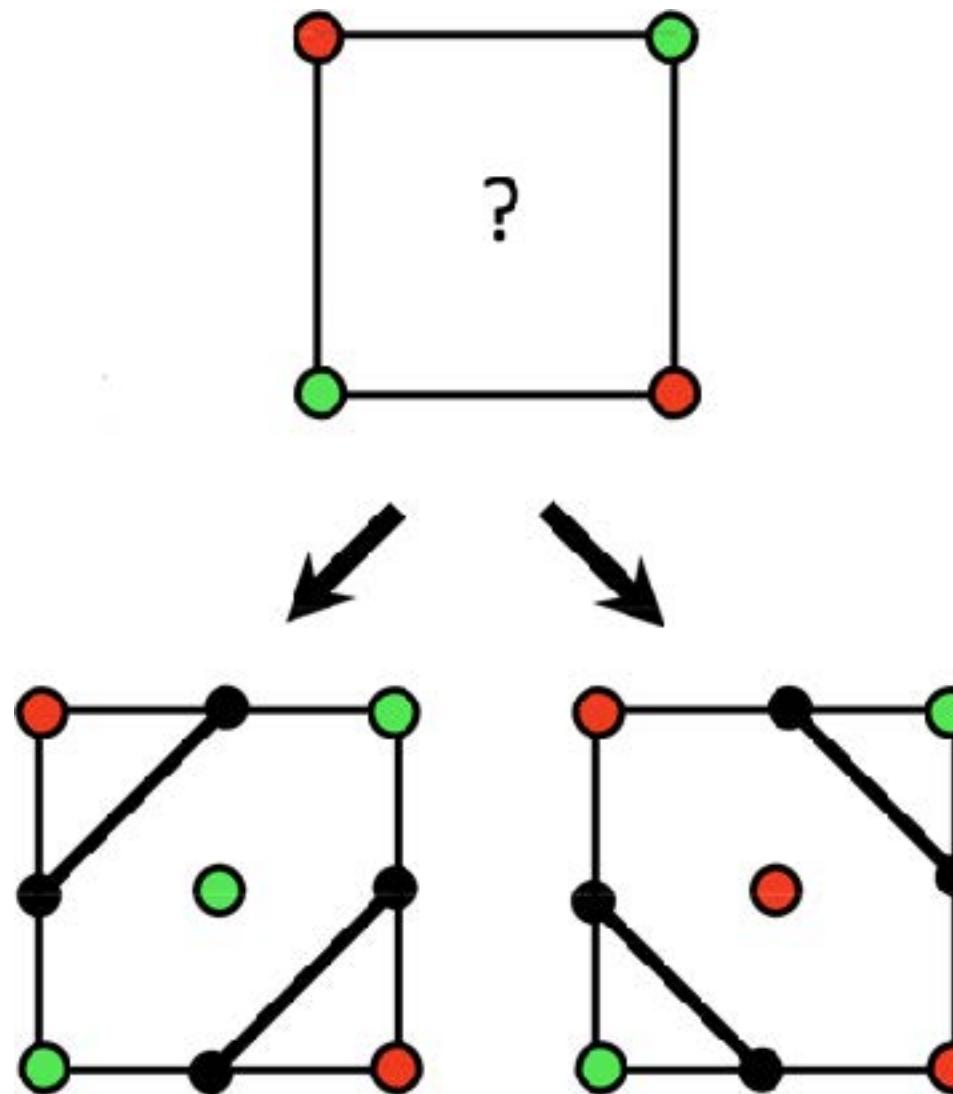
Classical Solution: 3D Marching Cube

- $2^8 = 256$ cases
- The first published version exploits rotation and inversion, and only considers 15 unique cases:



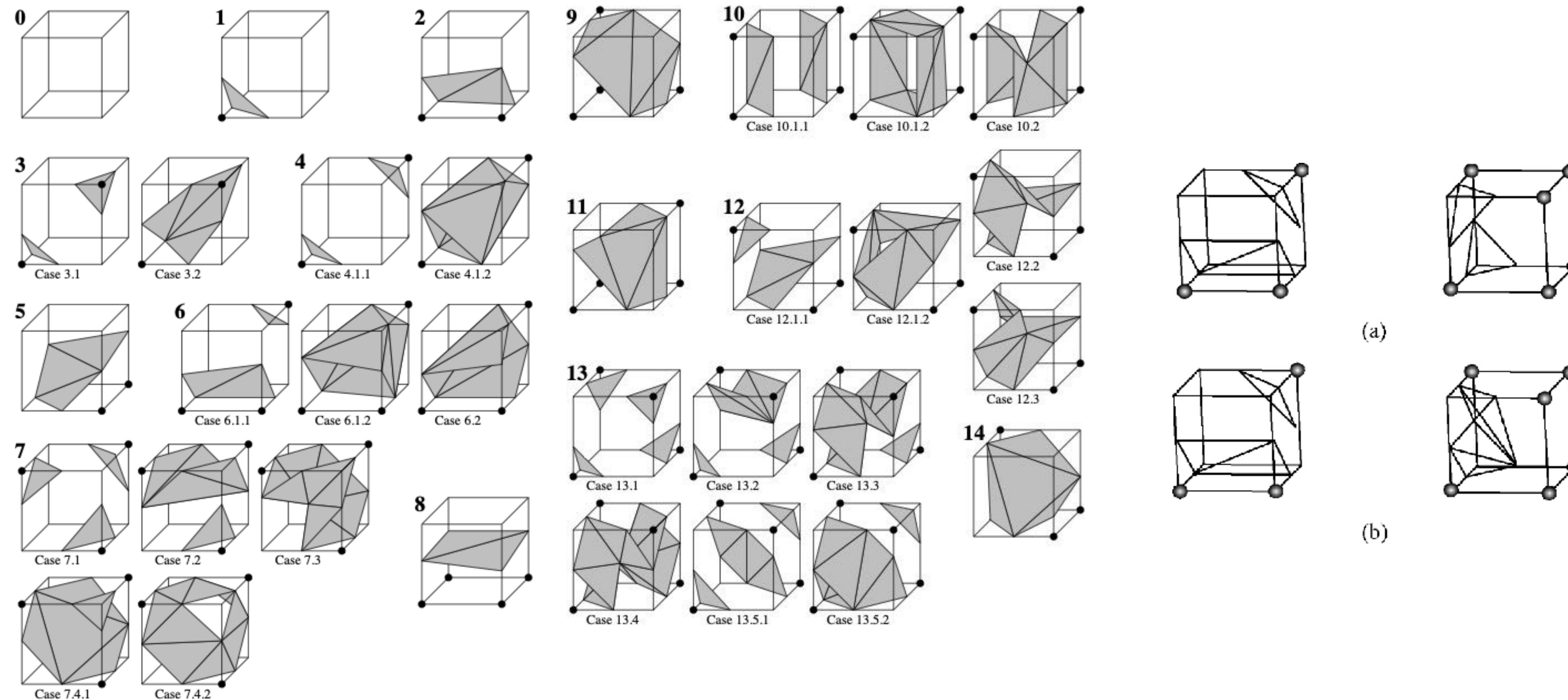
Ambiguity

- Ambiguity leads to holes:



Solution to Ambiguity

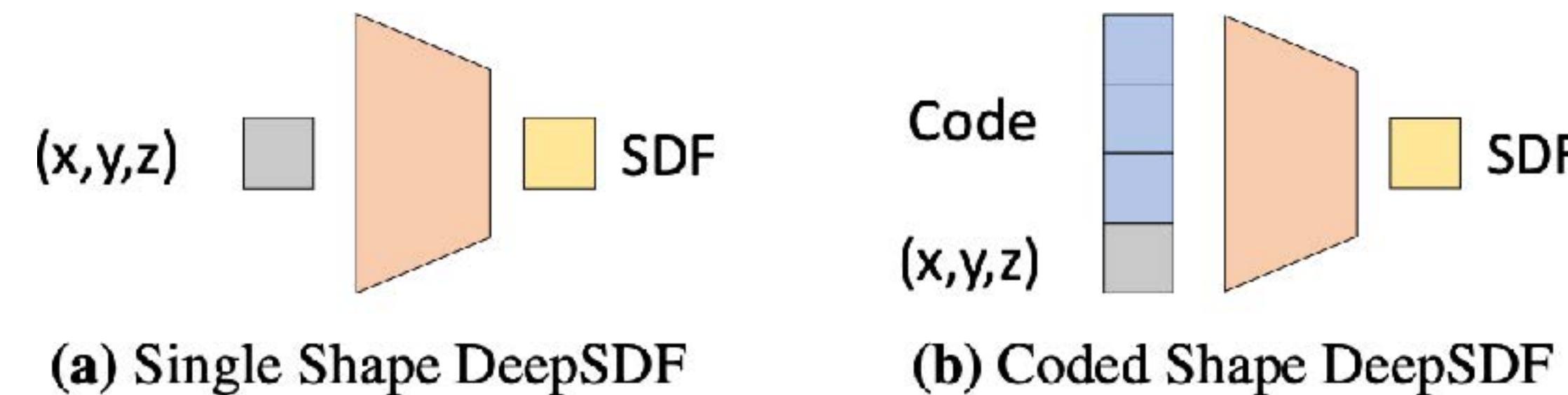
- Considering more cases in the look-up table by watching larger context:



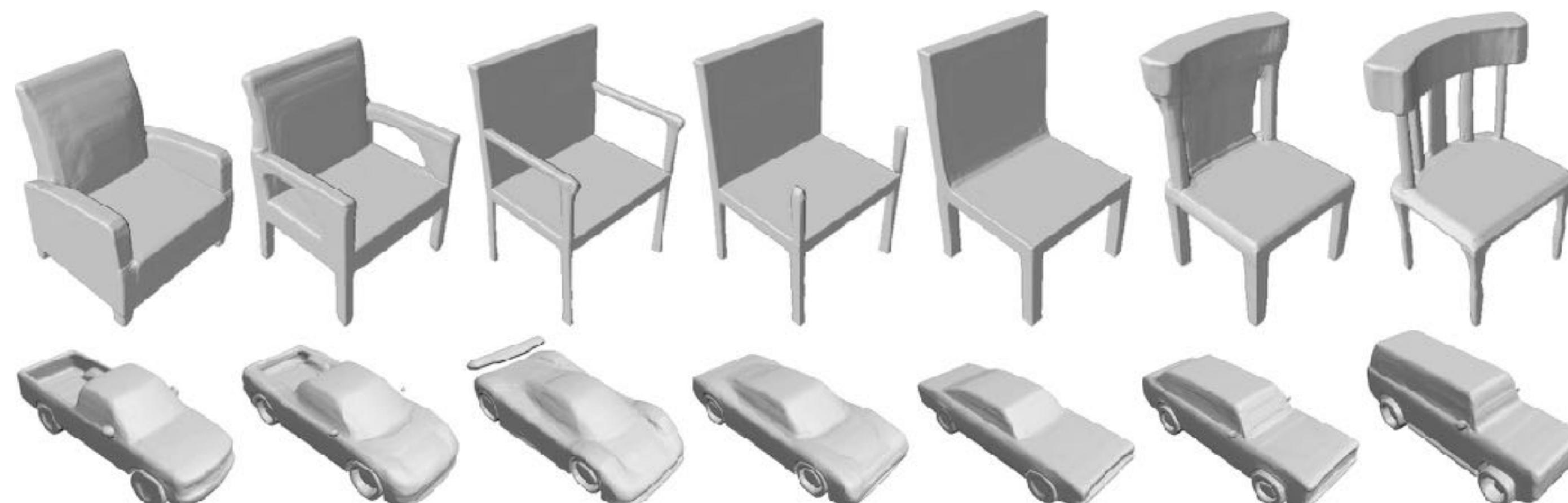
Comparison

	Explicit meshing (e.g., ball-pivoting)	Implicit meshing (e.g., RBF, MLS, Poission)
Sensitive to normals	No	Yes
Watertight manifold	No	Yes, in most cases
Complexity	Linear	<ul style="list-style-type: none">• Large-scale equations to estimate implicit function• Marching cubes• Dense voxelization

Use Neural Network to Approximate Implicit Field Function



- (a) use the network to overfit a single shape
- (b) use a latent code to represent a shape, so that the network can be used for multiple shapes



Learning-Based Marching Cube

Deep Marching Cubes: Learning Explicit Surface Representations

Yiyi Liao^{1,2} Simon Donné^{1,3} Andreas Geiger^{1,4}

¹Autonomous Vision Group, MPI for Intelligent Systems and University of Tübingen

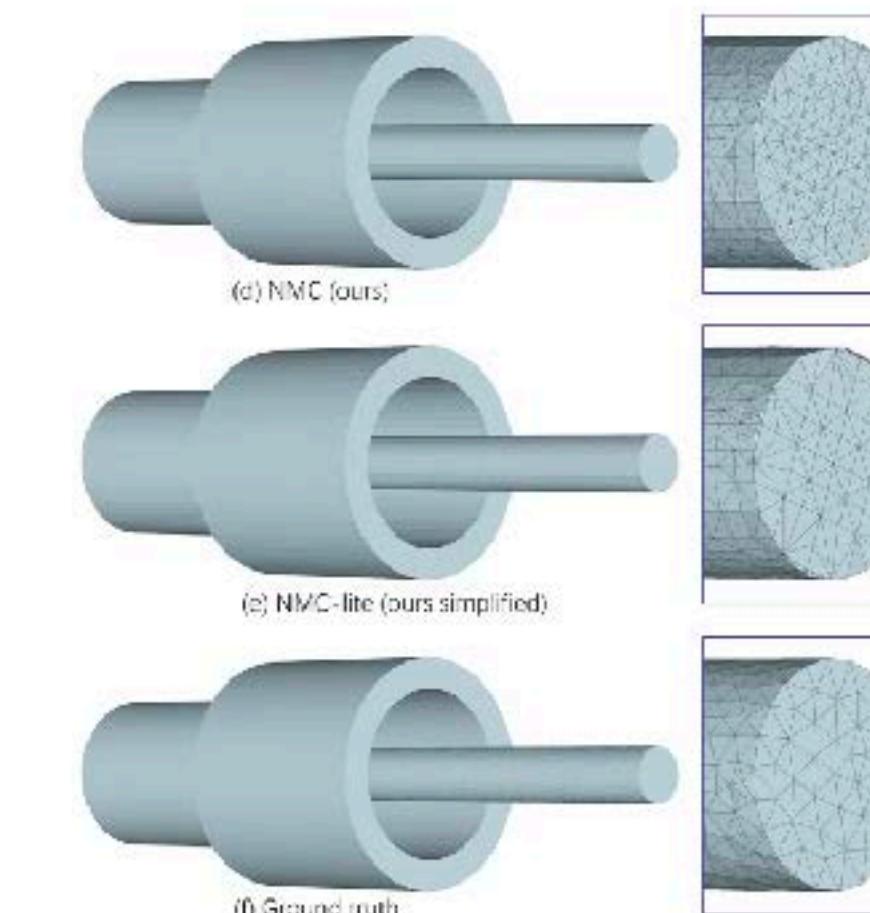
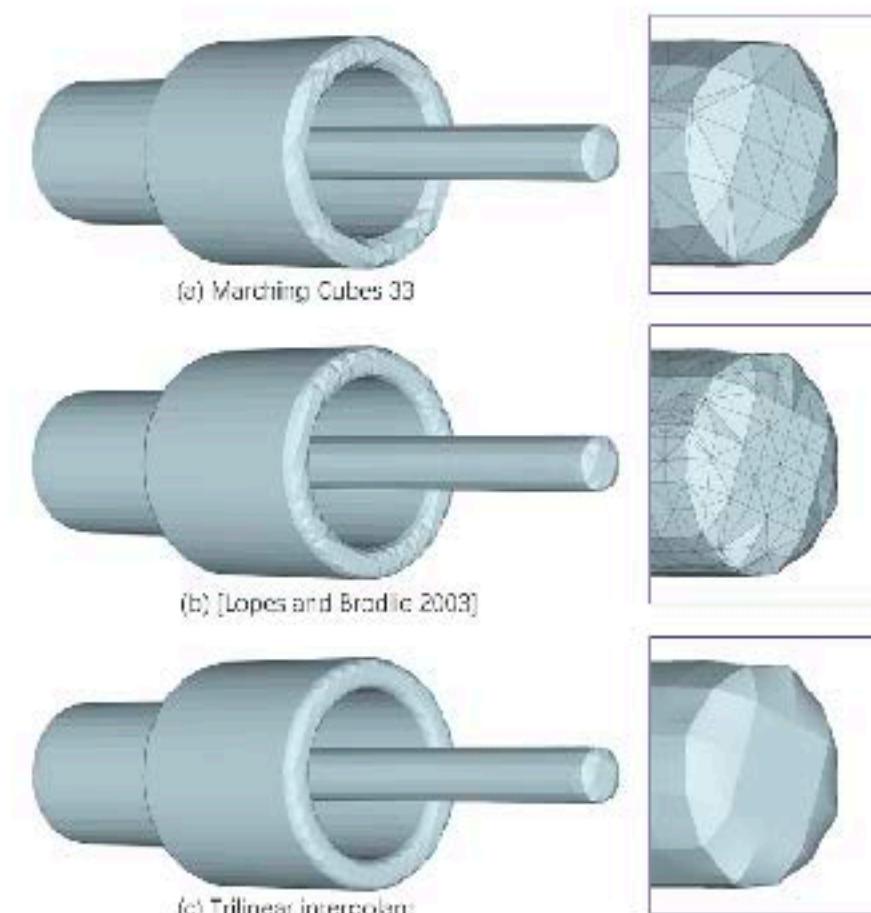
²Institute of Cyber-Systems and Control, Zhejiang University

³imec - IPI - Ghent University ⁴CVG Group, ETH Zürich

{yiyi.liao, simon.donne, andreas.geiger}@tue.mpg.de

Neural Marching Cubes

ZHIQIN CHEN, Simon Fraser University, Canada
HAO ZHANG, Simon Fraser University, Canada



Neural Dual Contouring

ZHIQIN CHEN, Simon Fraser University, Canada
ANDREA TAGLIASACCHI, Google Research, University of Toronto, Canada
THOMAS FUNKHOUSER, Google Research, USA
HAO ZHANG, Simon Fraser University, Canada

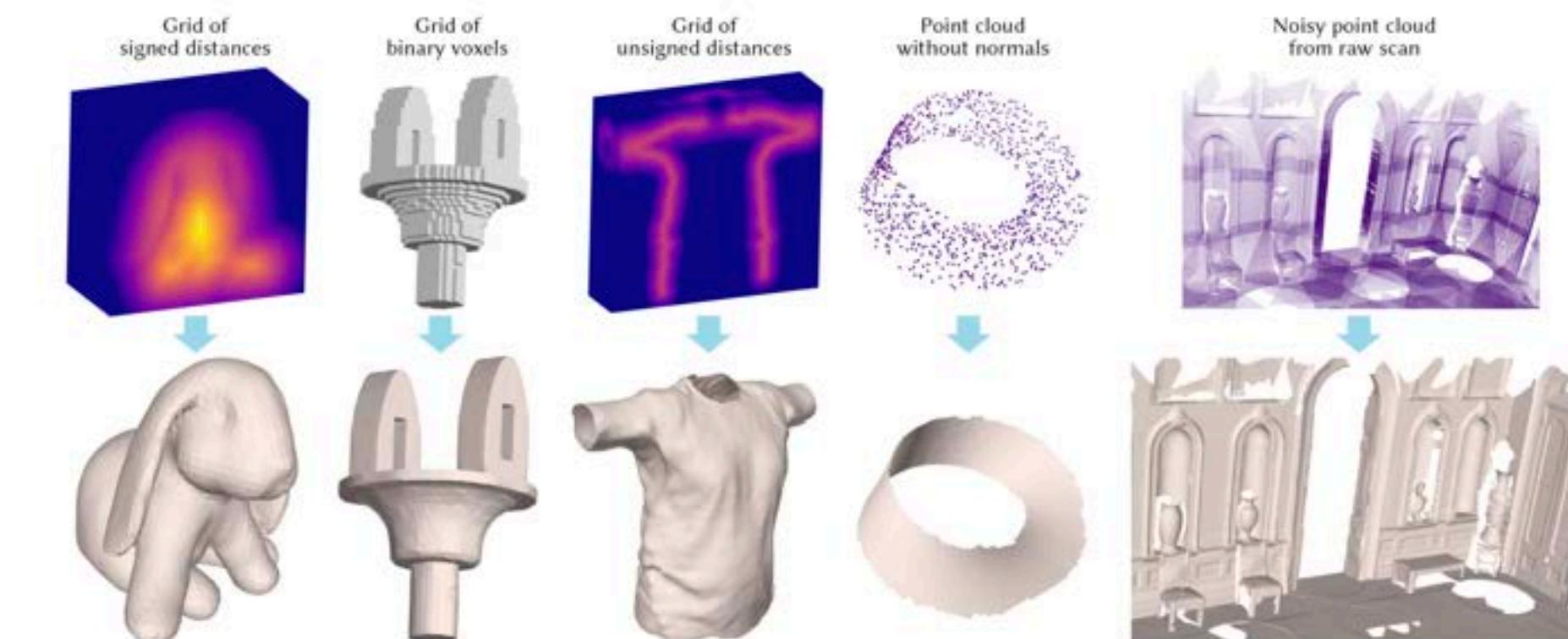


Fig. 1. Neural dual contouring (NDC) is a unified data-driven approach that learns to reconstruct meshes (bottom) from a variety of inputs (top): signed or unsigned distance fields, binary voxels, non-oriented point clouds, and noisy raw scans. Trained on CAD models, NDC generalizes to a broad range of shape types: CAD models with sharp edges, organic shapes, open surfaces for cloths, scans of indoor scenes, and even the non-orientable Möbius strip.

Consistent Orientation?

- Require the ground-truth ***signed*** implicit functions during training (e.g., signed distance, occupancy)
- However, 3D raw data, such as point cloud or triangle soup, are not necessarily consistently oriented
 - e.g., getting normal orientation for ShapeNet data is not easy

Sign Agnostic Learning of Shapes from Raw Data

- Unsigned distance is easy to obtain
 - Distance to the point cloud & triangle soup
- Learn ***signed*** distance from ***unsigned*** distance groundtruth
 - Require a special loss function

Sign Agnostic Learning

$$\text{loss}(\theta) = \mathbb{E}_{x \sim D_\chi} \tau(f(x; \theta), h_\chi(x))$$

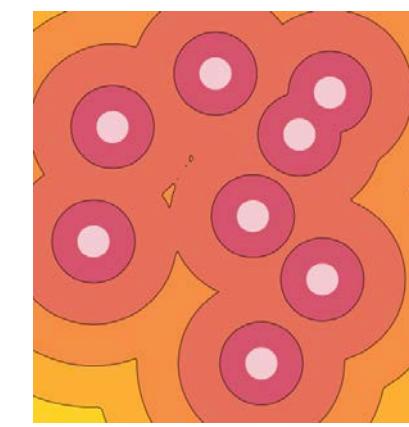
- $\chi \subset \mathbb{R}^3$: input raw data (e.g., a point cloud or a triangle soup)
- $f(x; \theta) : \mathbb{R}^3 \times \mathbb{R}^m \rightarrow \mathbb{R}$: learned signed function
- D_χ : distribution of the training samples defined by χ
- $h_\chi(x)$: some *unsigned* distance measure to χ
- $\tau : \mathbb{R} \times \mathbb{R}_+ \rightarrow \mathbb{R}$: a similarity function

Sign Agnostic Learning

$$\text{loss}(\theta) = \mathbb{E}_{x \sim D_\chi} \tau(f(x; \theta), h_\chi(x))$$

- $h_\chi(x)$: some *unsigned* distance measure to χ

$$h_2(z) = \min_{x \in \mathcal{X}} \|z - x\|_2$$



$$h_0(z) = \begin{cases} 0 & z \in \mathcal{X} \\ 1 & z \notin \mathcal{X} \end{cases}$$



- $\tau : \mathbb{R} \times \mathbb{R}_+ \rightarrow \mathbb{R}$: a similarity function

$$\tau_\ell(a, b) = |a - b|^\ell$$

Two Local Minima

$$\text{loss}(\theta) = \mathbb{E}_{x \sim D_\chi} \tau(f(x; \theta), h_\chi(x))$$

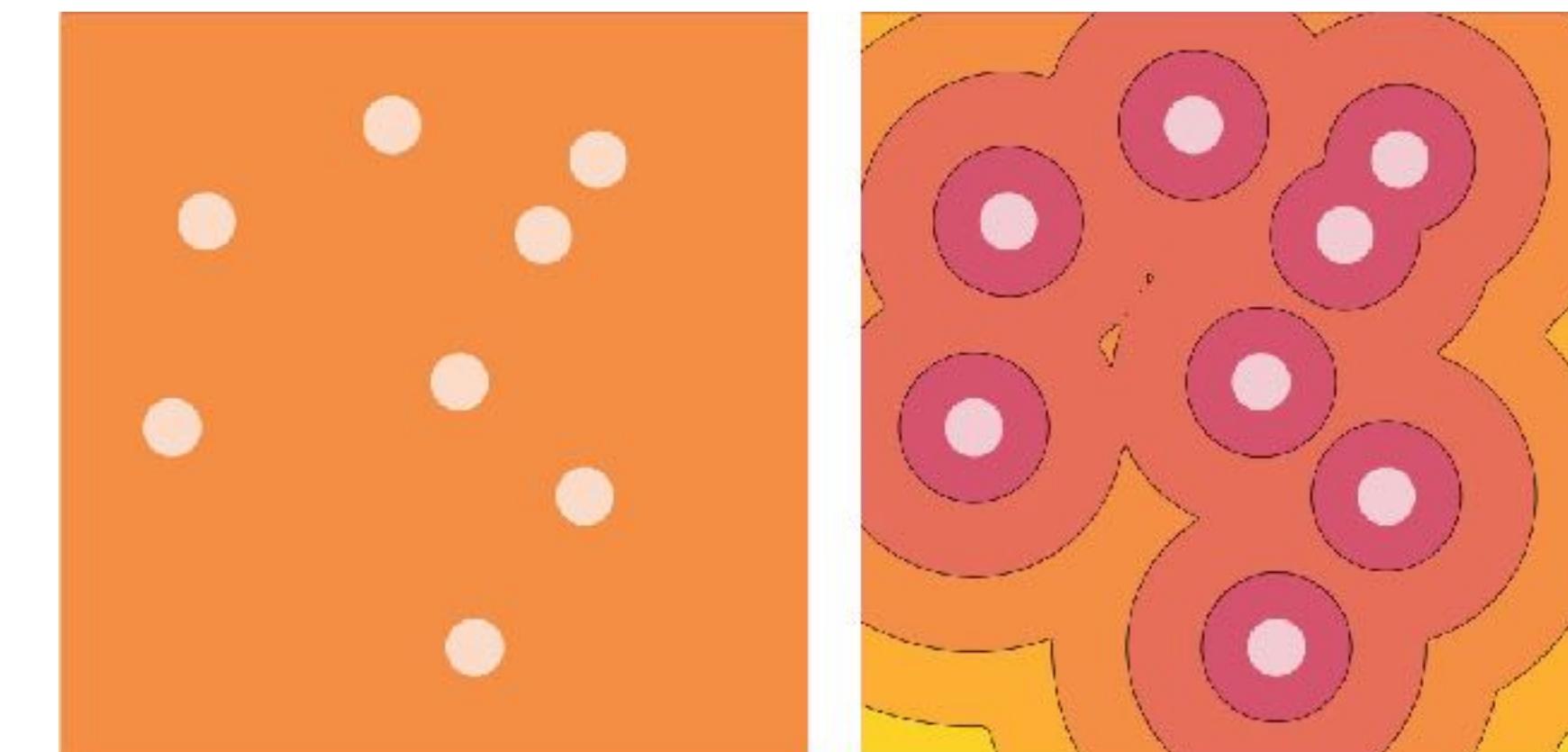
- f is an ***unsigned*** function that resembles $h_\chi(x)$
- f is a ***signed*** function and $|f|$ resembles $h_\chi(x)$
- We prefer the second case to use marching cube

Two Local Minima

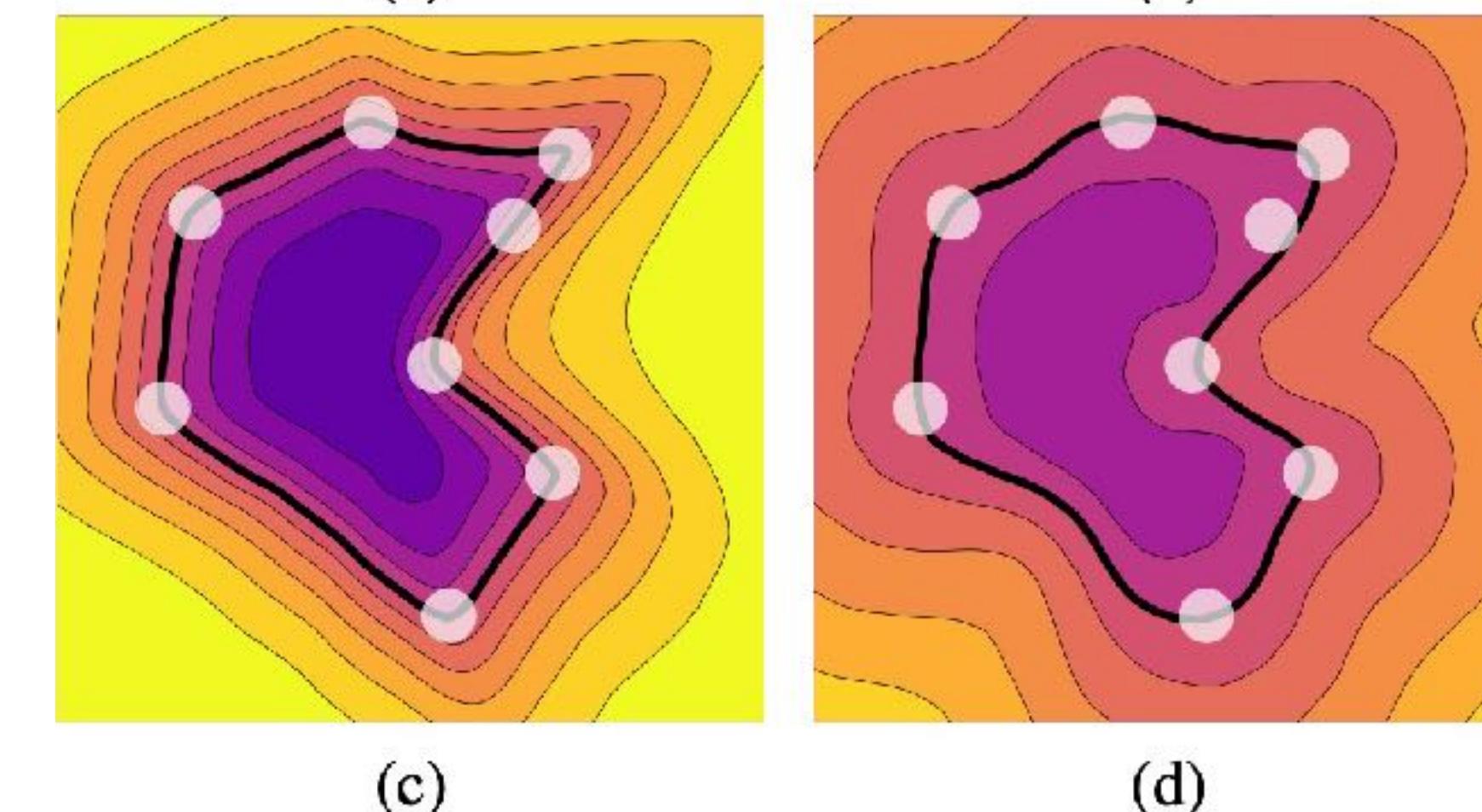
- Pick a special weight initialization θ^0 so that $f(x; \theta^0) \approx \varphi(\|x\| - r)$, where $\varphi(\|x\| - r)$ is the signed distance function to an r -radius sphere
 - $f > 0$ if $\|x\| > r$
 - $f < 0$ if $\|x\| < r$
- Under such an initialization, f is not easy to converge to the unsigned local minima.

2D Results

Unsigned
function as
supervision

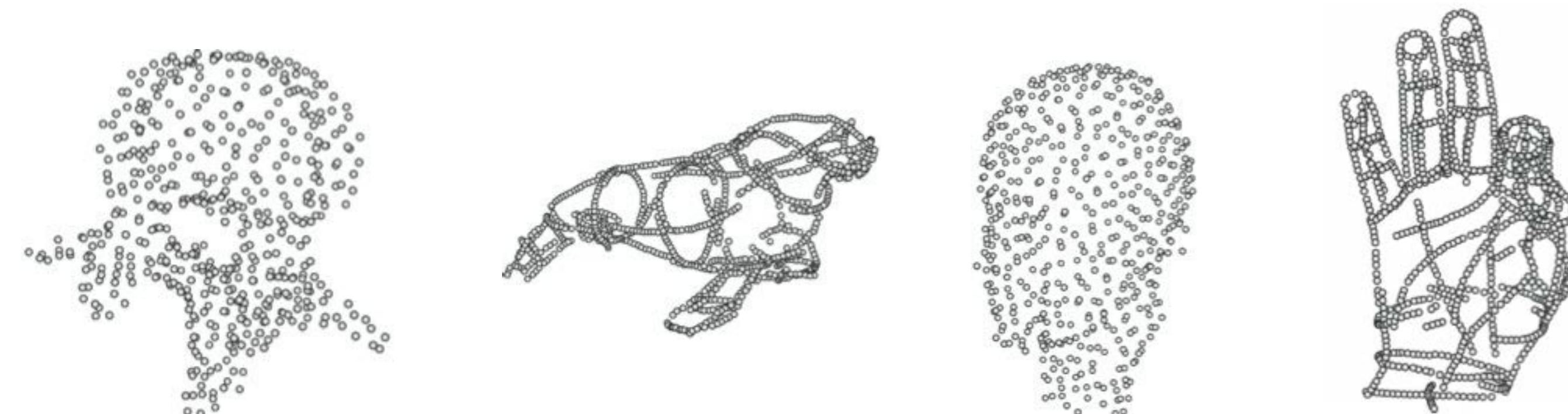


Learned
signed function



Surface Reconstruction Results from Raw Data

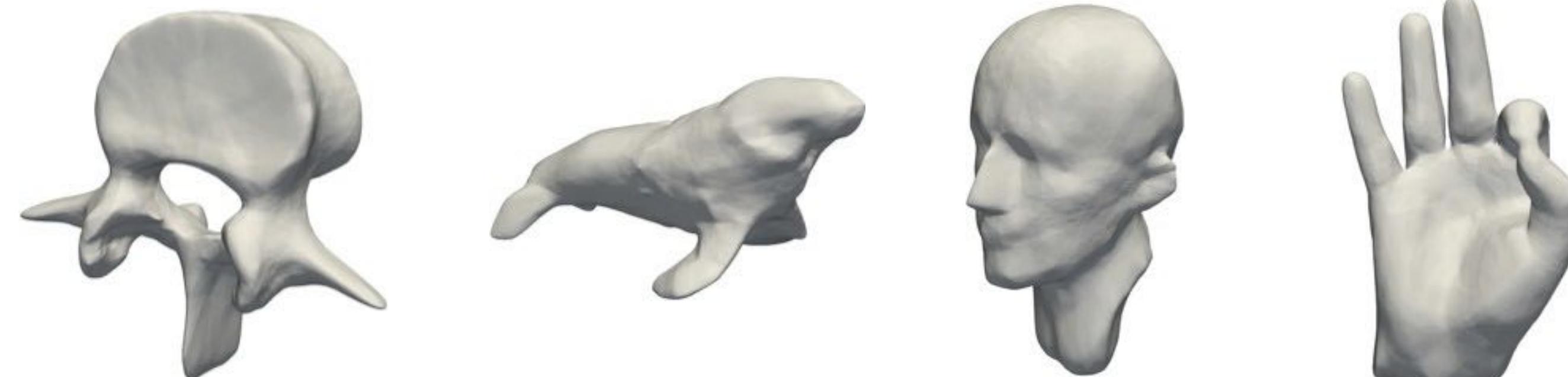
Raw
Point Cloud



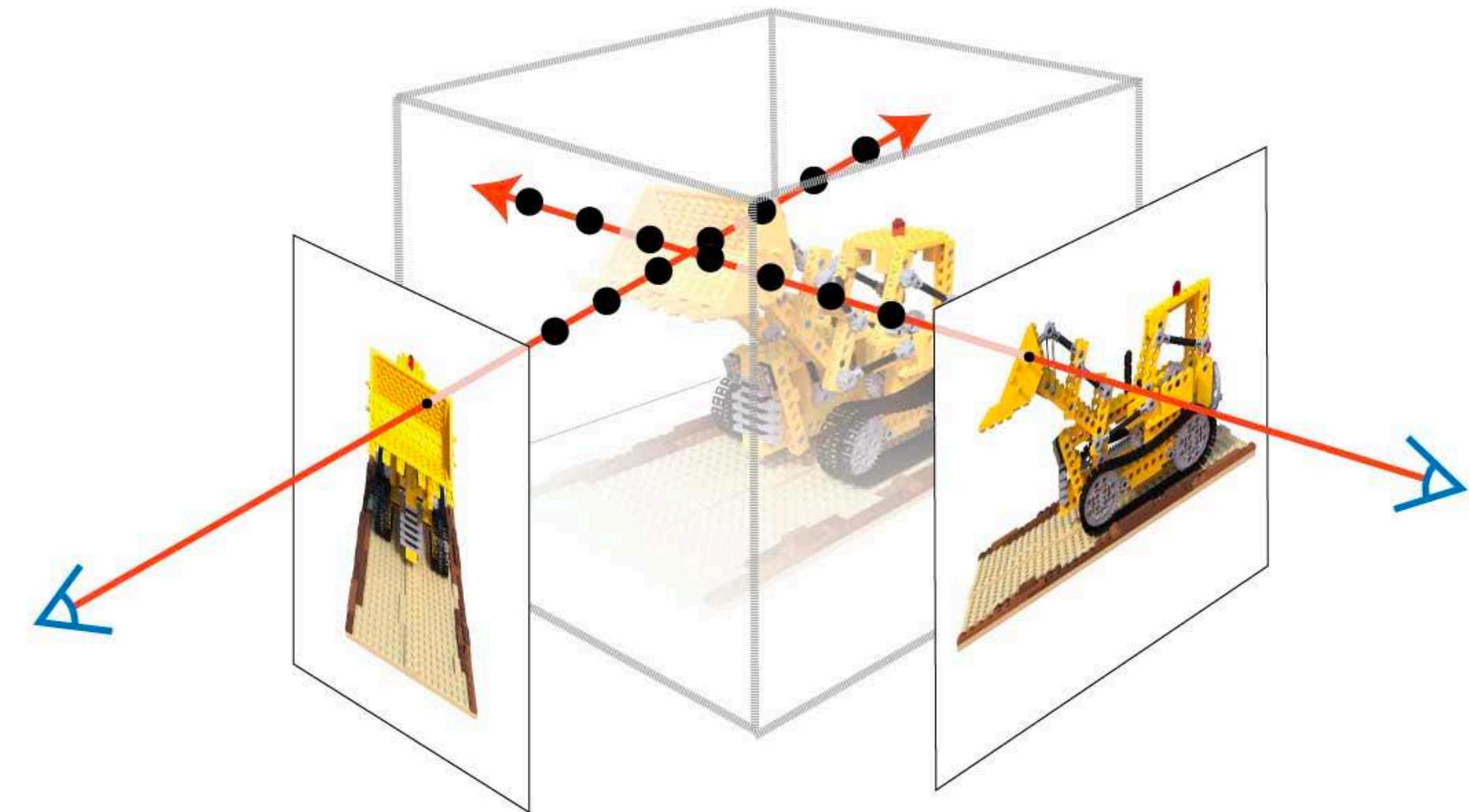
Ball-Pivoting
Algorithm



SAL



Recap: Single-View / Multi-View 3D



From Image(s) to 3D

Today's Focus



a DSLR photo of a squirrel



3D from scratch (Unconditional Generation)



Learn the space of (possibly textured) shapes

Should be able to **sample** from the learned shape space

3D from scratch (Unconditional Generation)



Why is this an important task?

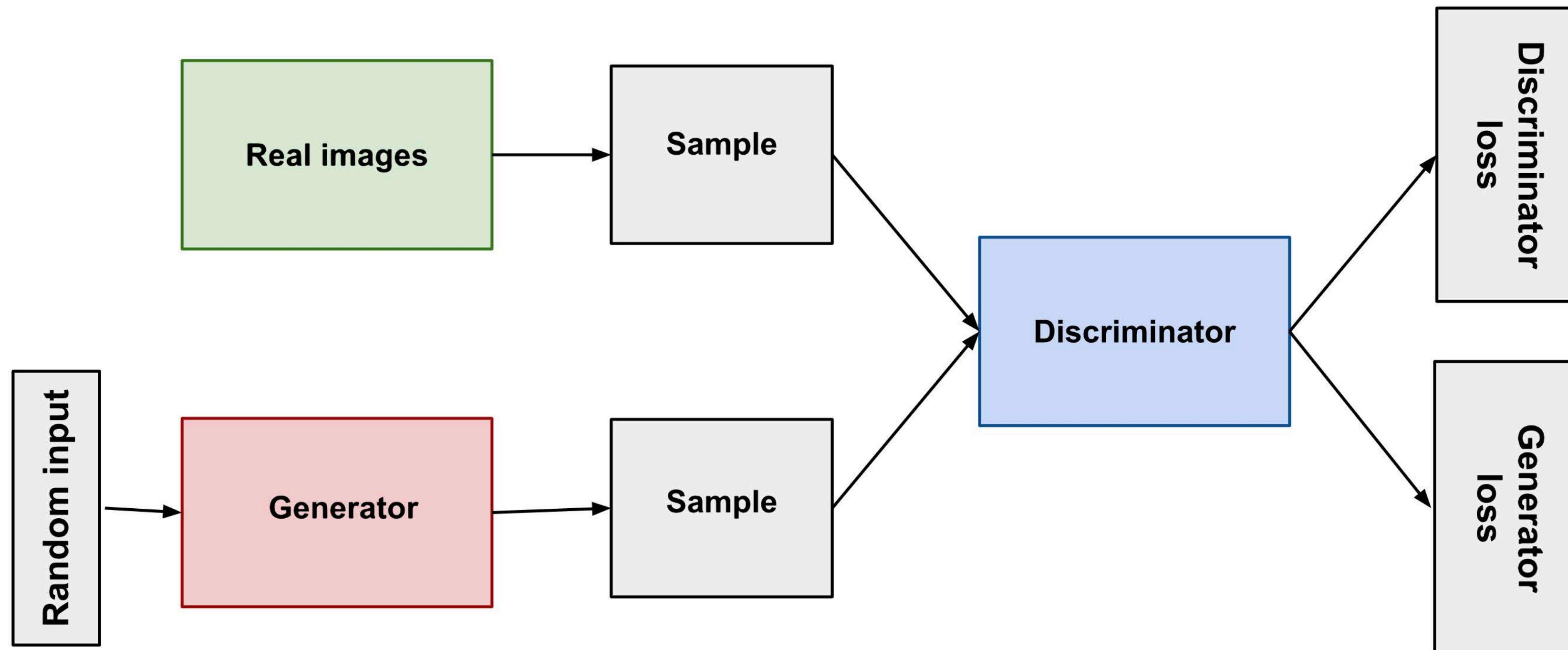
Can re-use the learned shape space for downstream conditional inference

Useful for artists exploring

Outline

- 3D GAN
- 3D Autoregressive Models
- 3D Diffusion Models
- 3D Generation without 3D Training Data
- Part-based 3D Generation

Background: Generative Adversarial Networks



$$L_{GAN}(G, D) = E_y[\log D(y) + E_{x,z}[\log(1 - D(G(x, z)))]$$

Learning Objective: Generate output that is indistinguishable from a ‘real’ example

StyleGAN2



StyleGAN3 (Ours)



Random latent walk using directions from StyleCLIP, GANSpace, and SeFa.



StyleGAN2



StyleGAN3 (Ours)

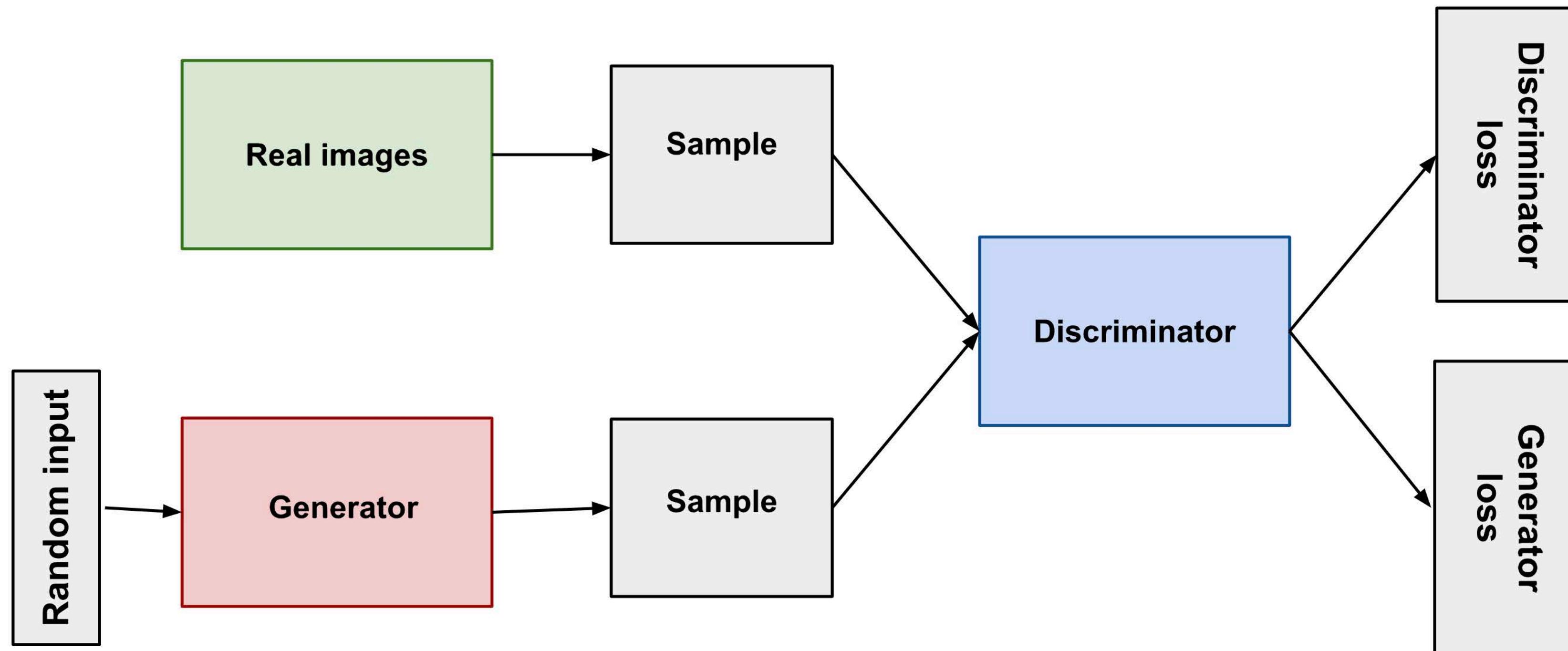


StyleGAN2



StyleGAN3 (Ours)

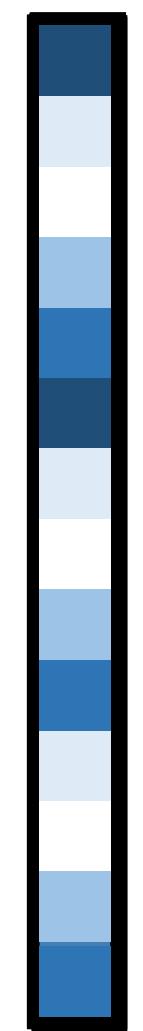
Background: Generative Adversarial Networks



$$L_{GAN}(G, D) = E_y[\log D(y) + E_{x,z}[\log(1 - D(G(x, z)))]$$

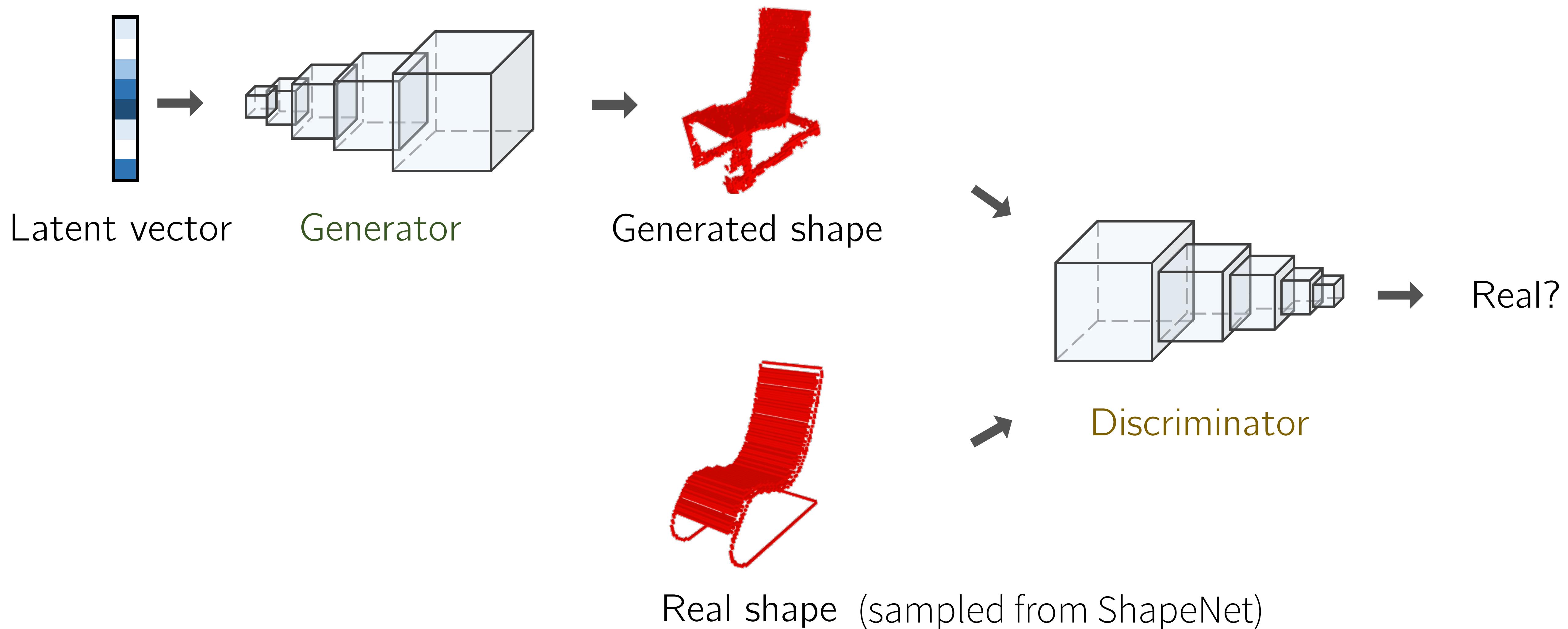
Can we similarly learn to generate 3D?

3D GANs

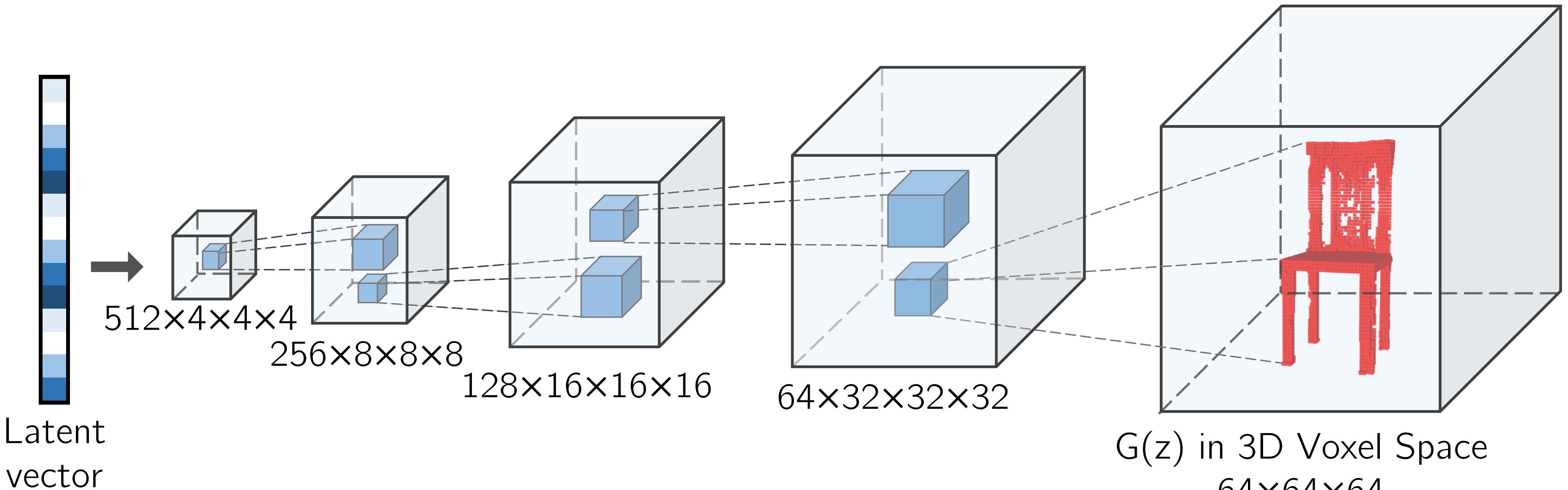


Latent vector

3D GANs



3D GAN: Generator



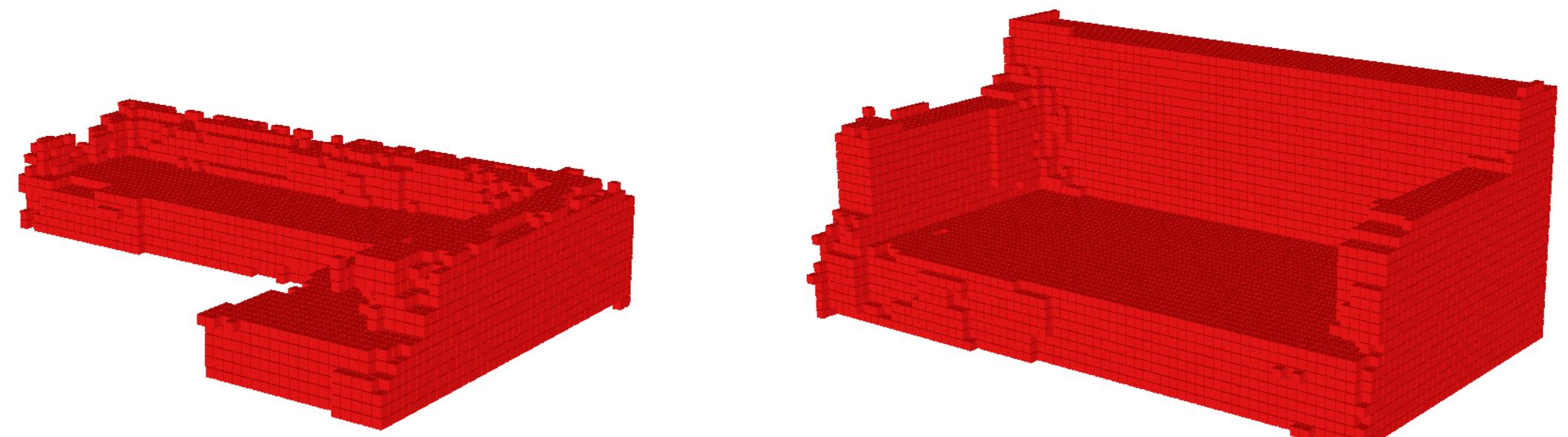
Incrementally increase resolution via convolutions and upsampling layers

3D GAN: Sampled Shapes

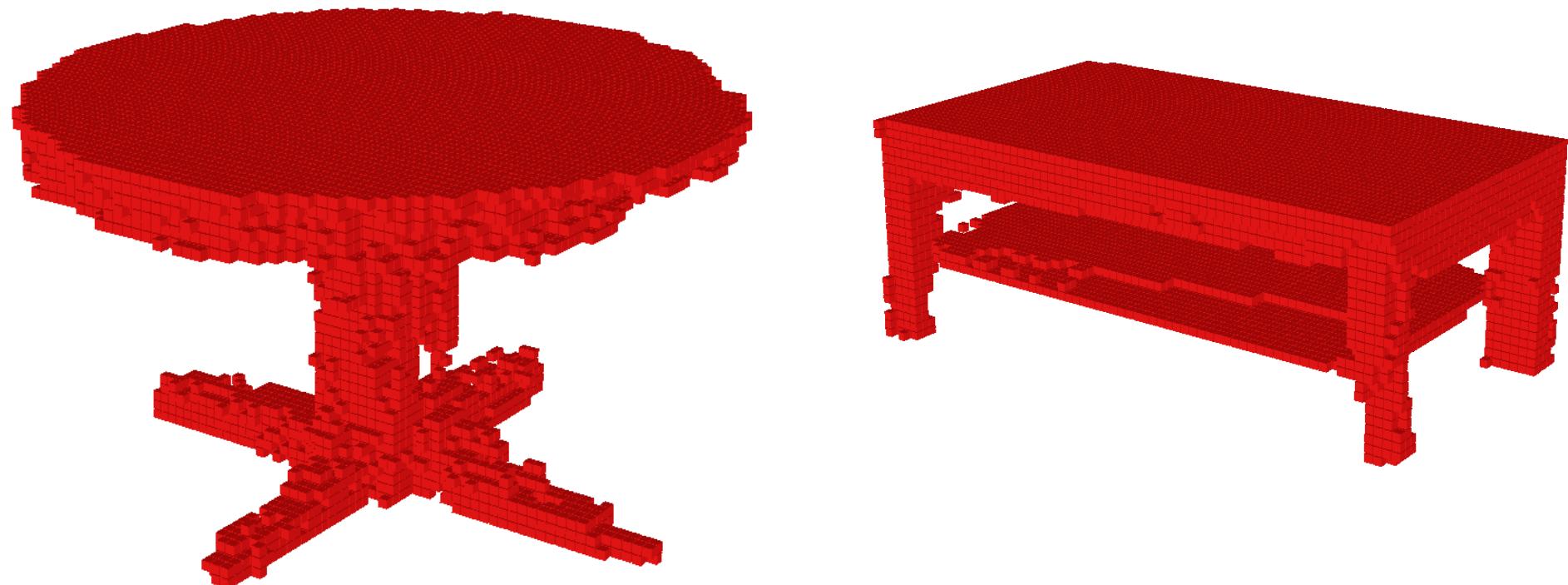
Chairs



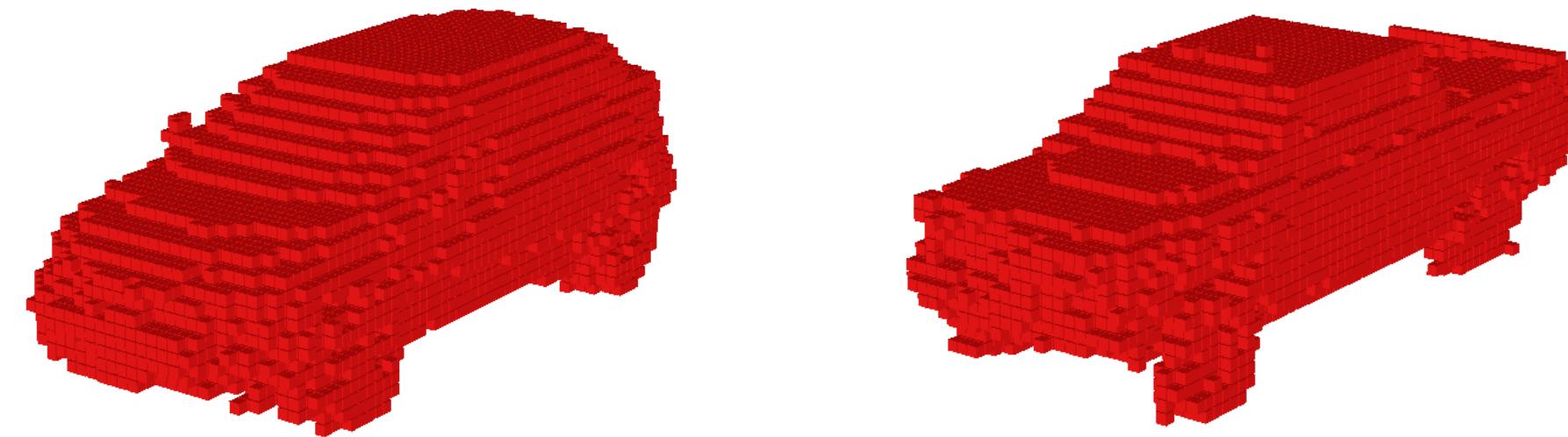
Sofas



Tables



Cars



3D GAN: Sampled Shapes

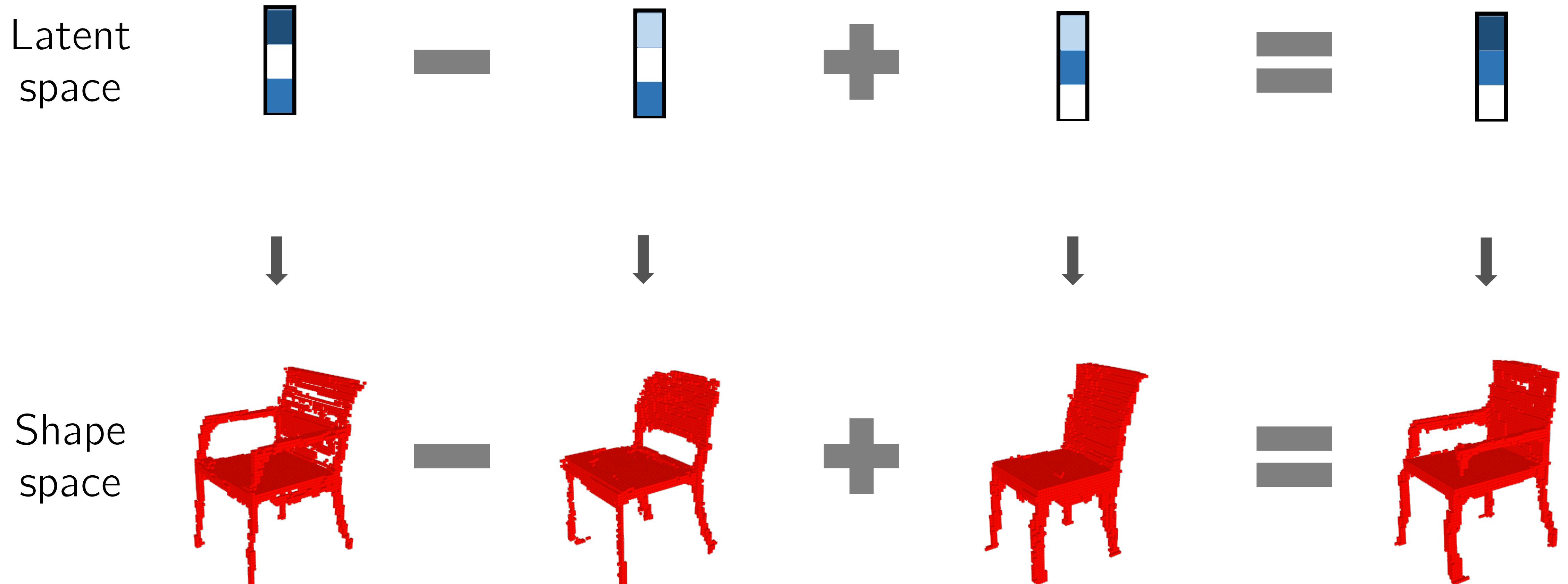


3D GAN: Vector Arithmetic

Interpolation in Latent Space

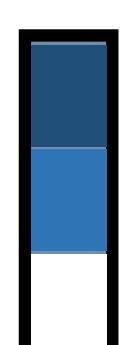
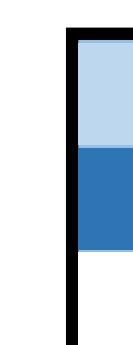
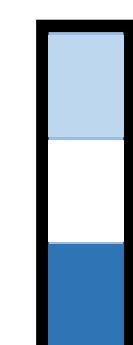


3D GAN: Vector Arithmetic

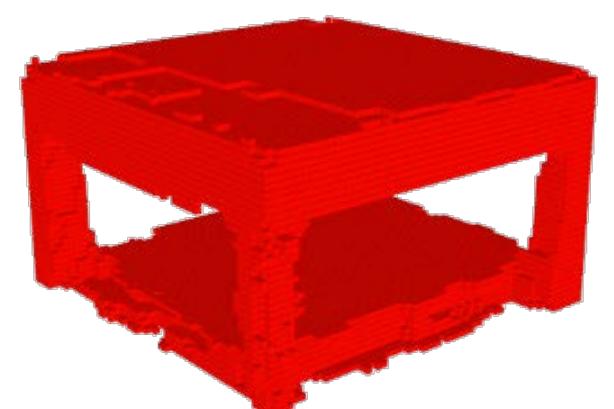
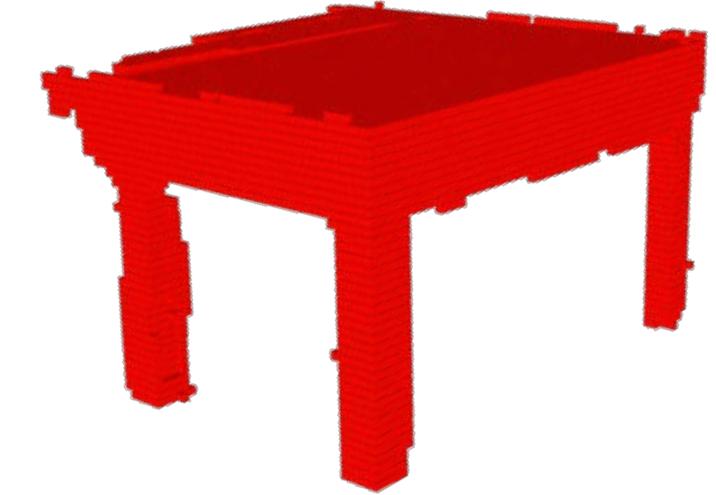
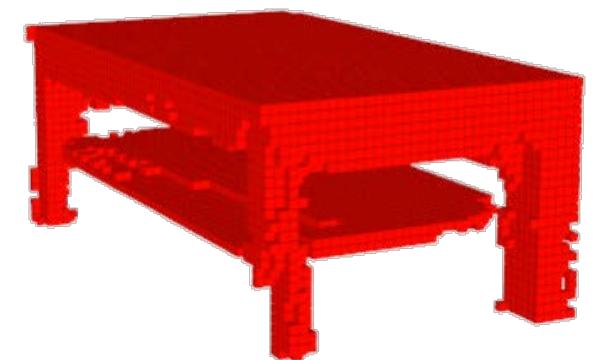


3D GAN: Vector Arithmetic

Latent
space

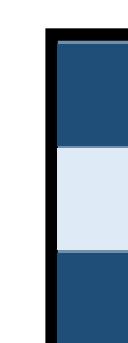
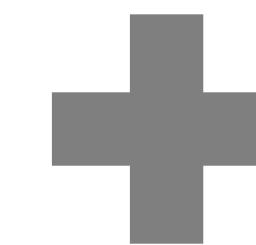
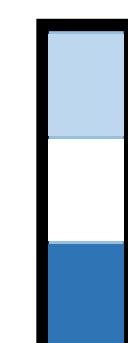
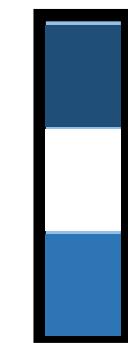


Shape
space

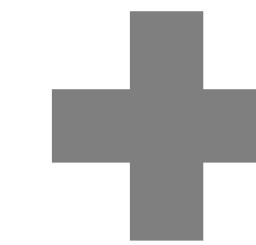
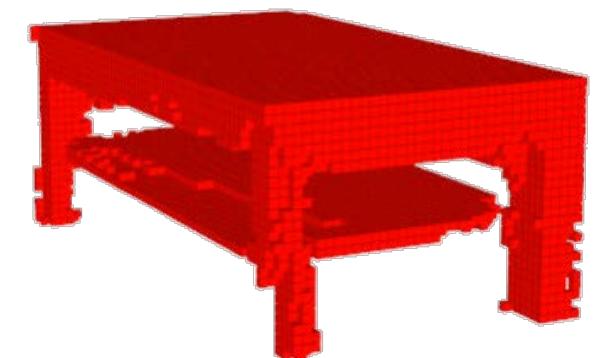


3D GAN: Vector Arithmetic

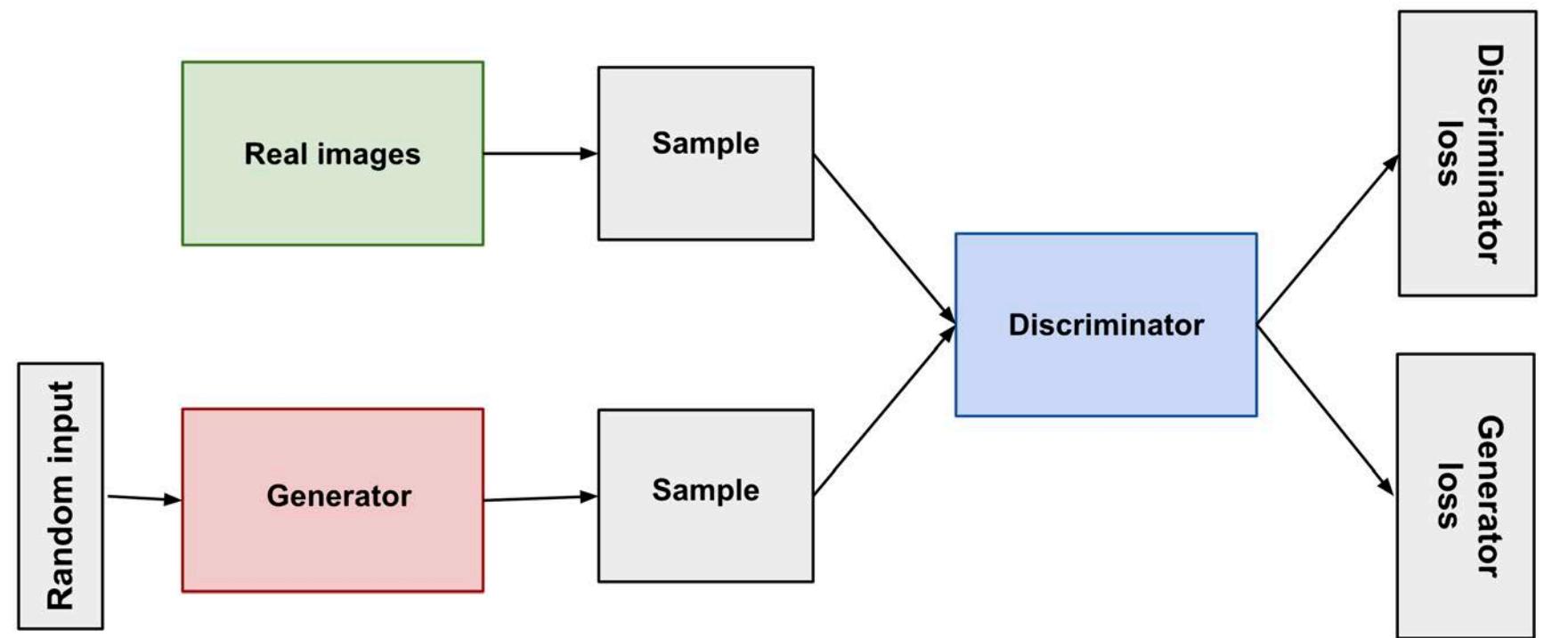
Latent
space



Shape
space

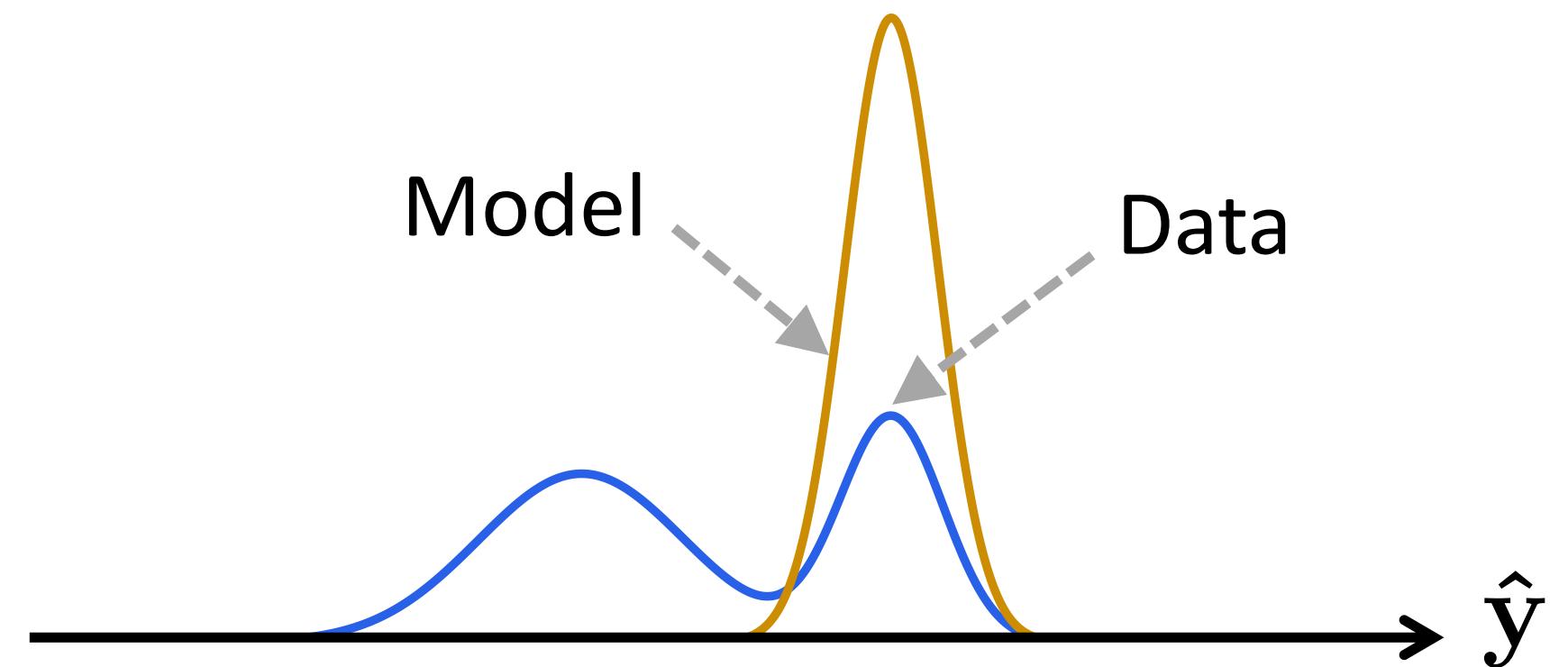


GAN Issues



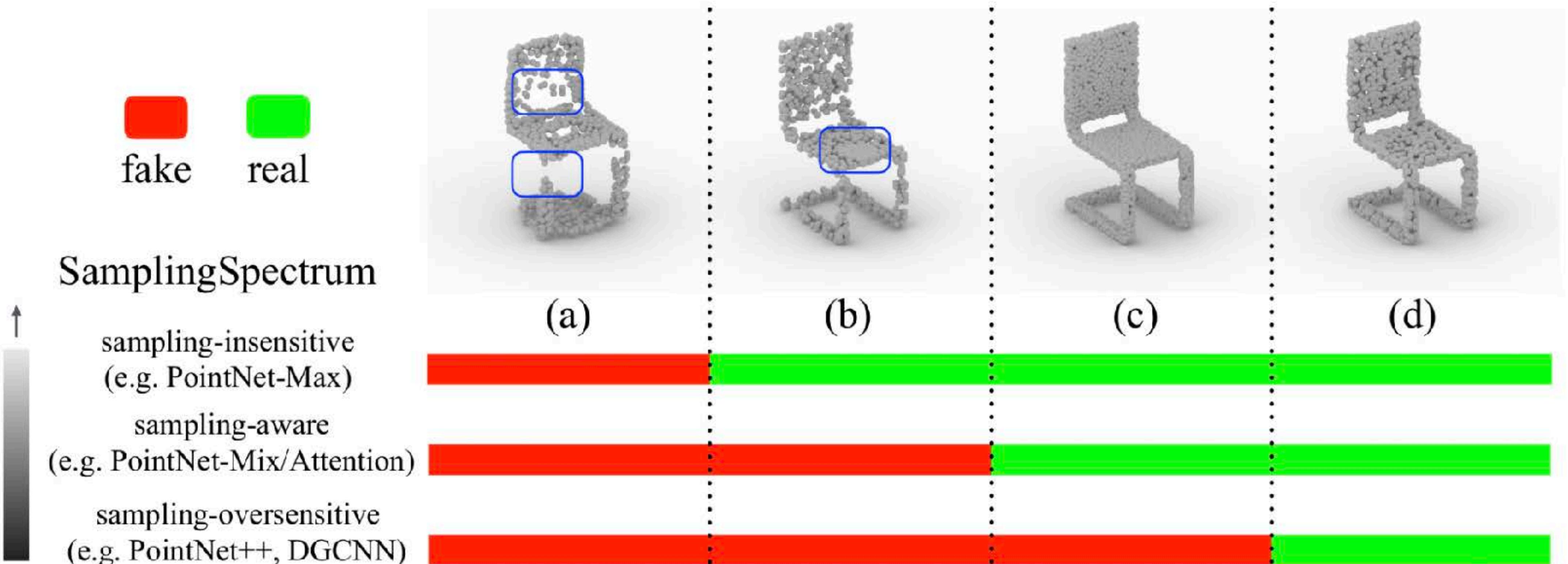
$$L_{GAN}(G, D) = E_y[\log D(y) + E_{x,z}[\log(1 - D(G(x, z)))]$$

Common issue: **mode collapse**



GANs don't maximize likelihood of data

GAN Issues

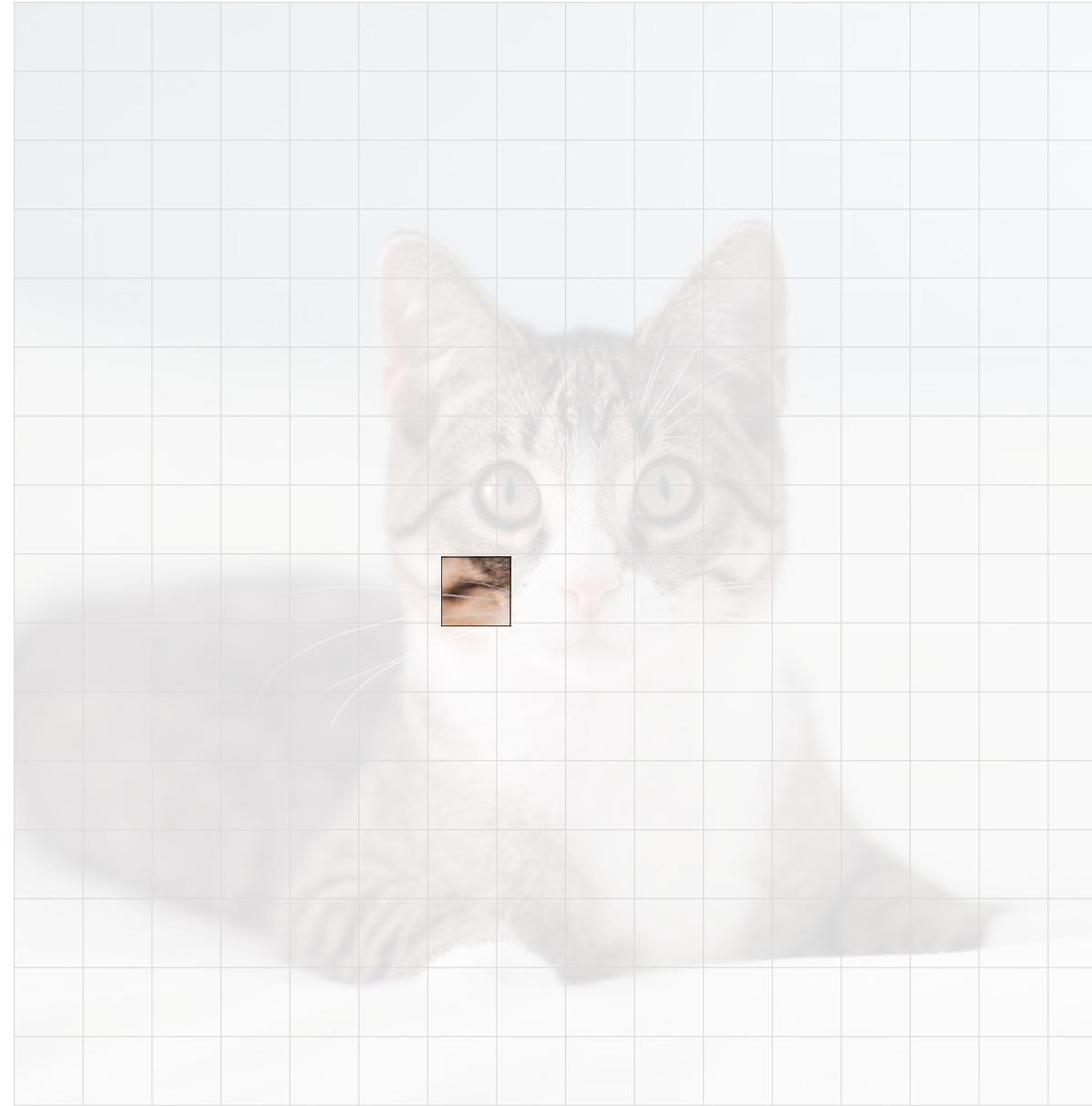


Outline

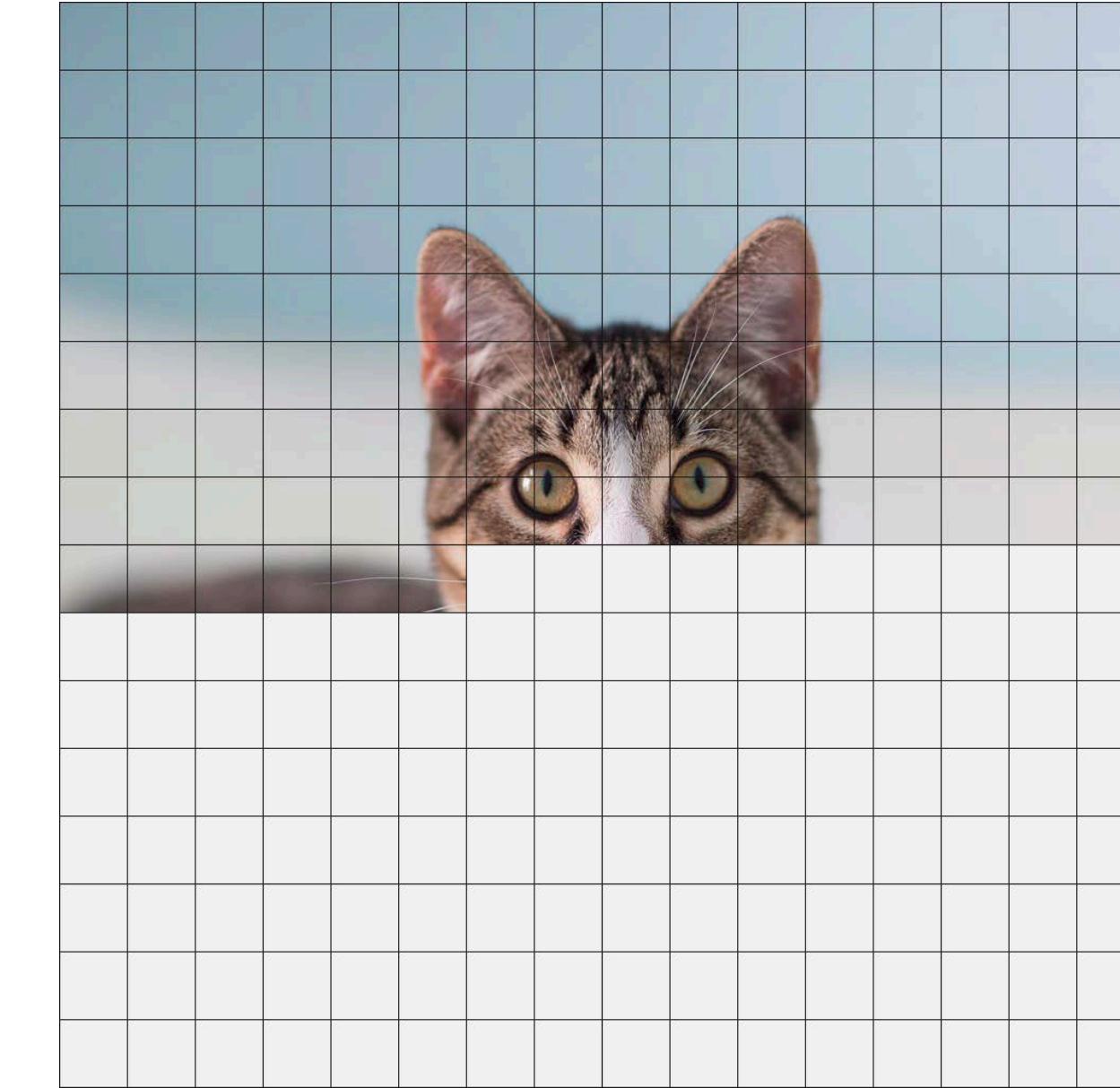
- 3D GAN
- 3D Autoregressive Models
- 3D Diffusion Models
- 3D Generation without 3D Training Data
- Part-based 3D Generation

Background: Autoregressive Models

$p($



|

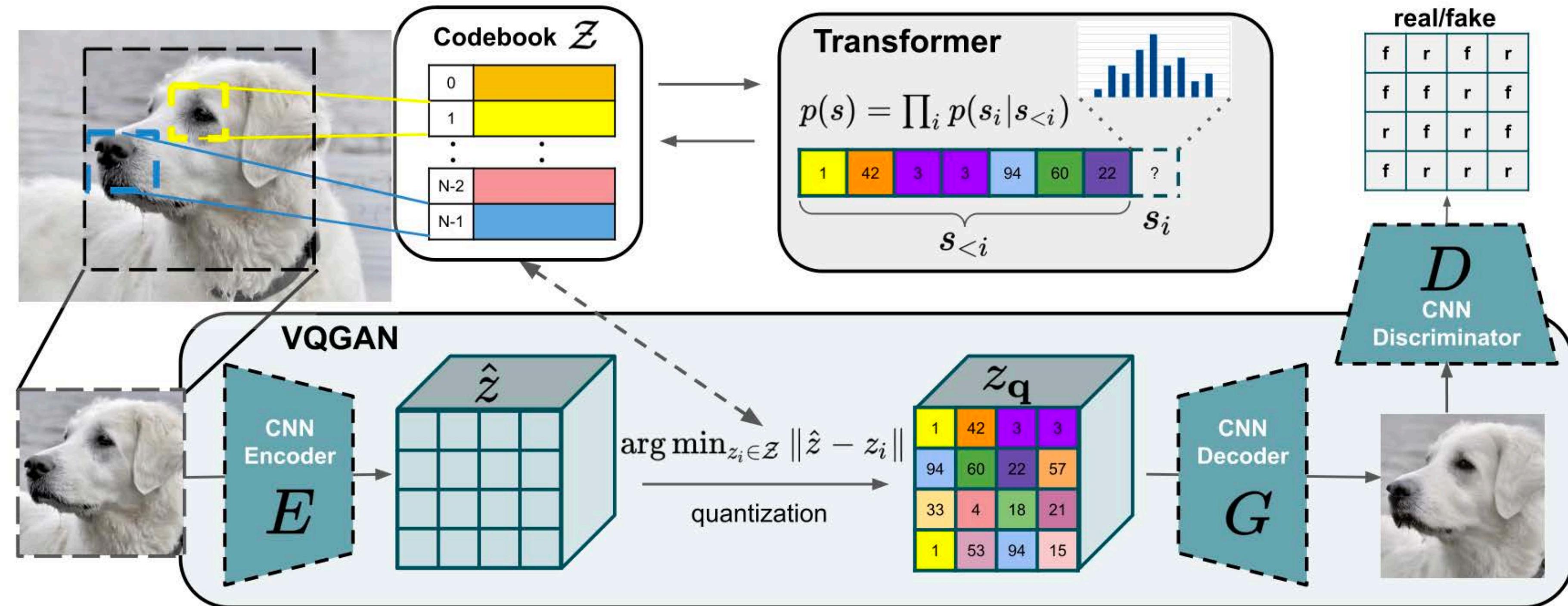


)

Training: $\max \mathbb{E}_{x \sim \mathcal{D}} \log p_{\theta}(x)$

Inference: Sample images, one pixel at a time

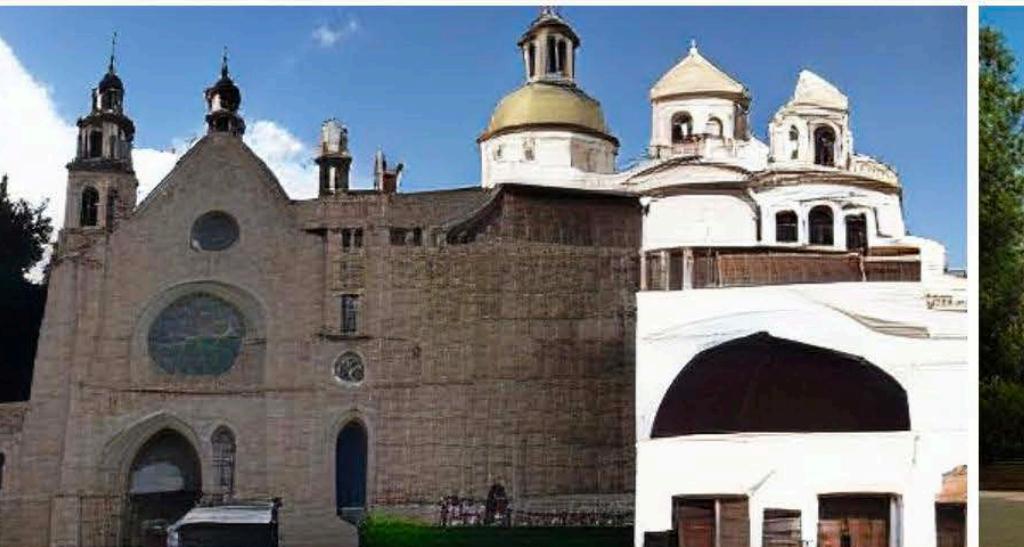
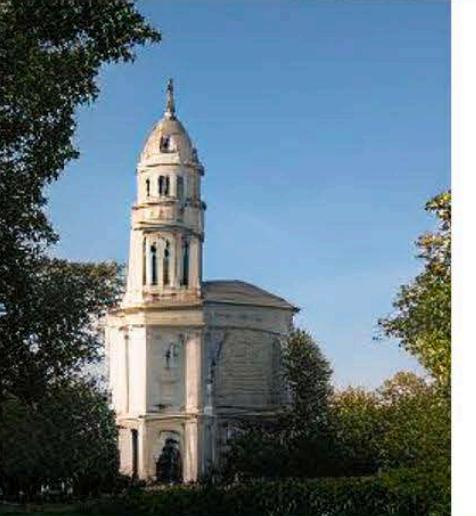
Background: Autoregressive Models



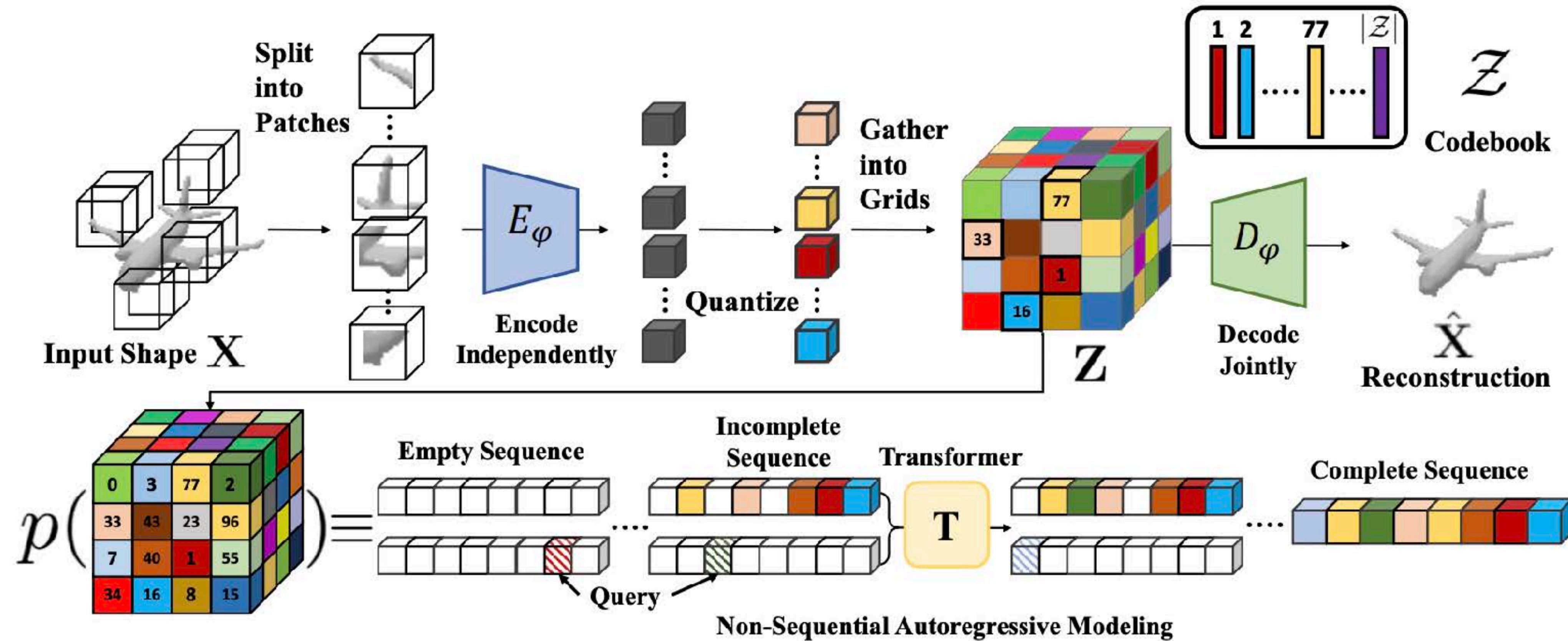
Step 1: Learn a ‘codebook’ of discrete patch representations

Step 2: Learn transformer-based Autoregressive models

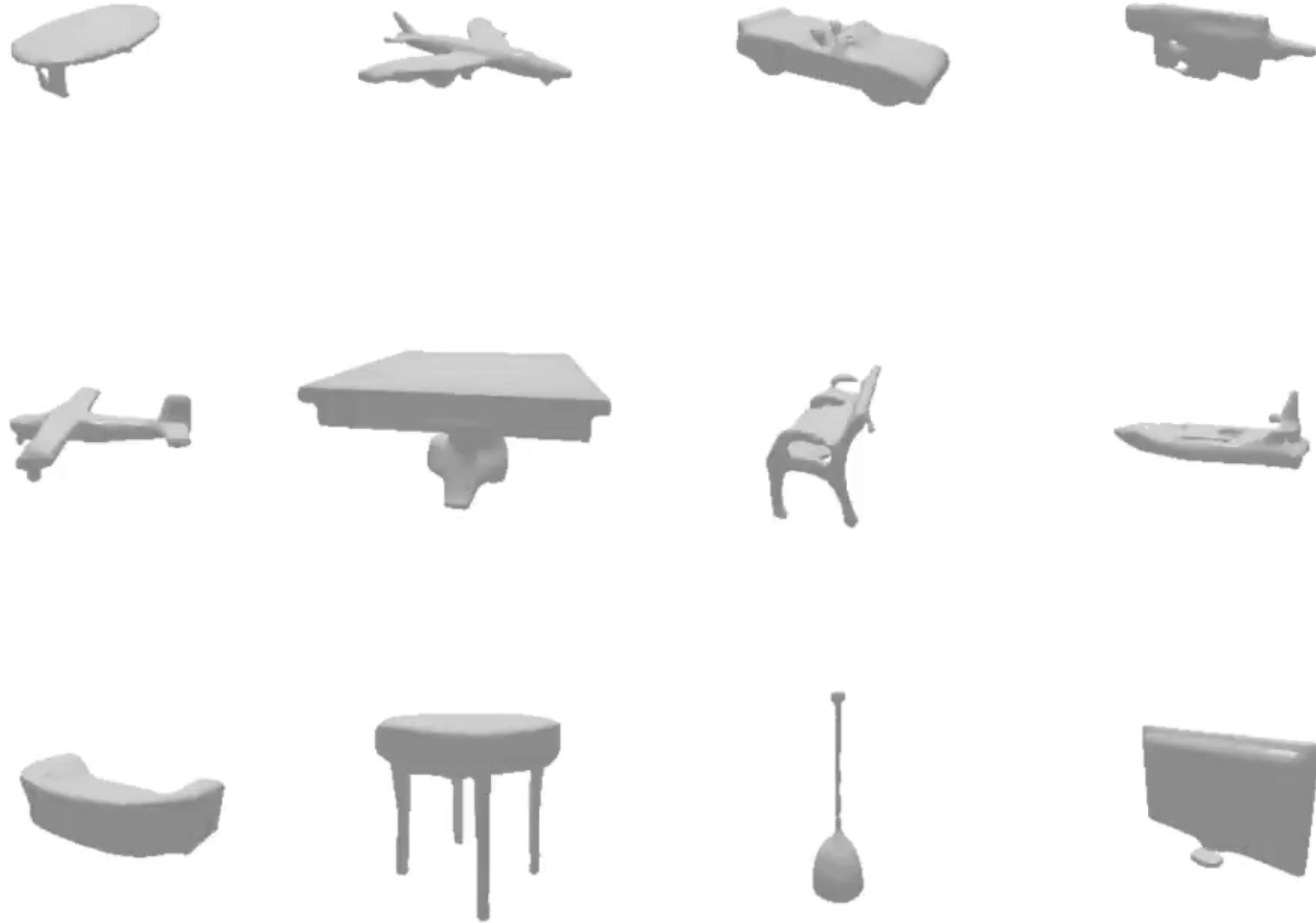
Background: Autoregressive Models



Autoregressive Generation of 3D Shapes



Autoregressive Generation of 3D Shapes



Unconditional generation
from model trained over
50 categories

Autoregressive Generation of 3D Shapes



Trivially allows shape completion (missing regions in red)

Autoregressive Generation of 3D Shapes

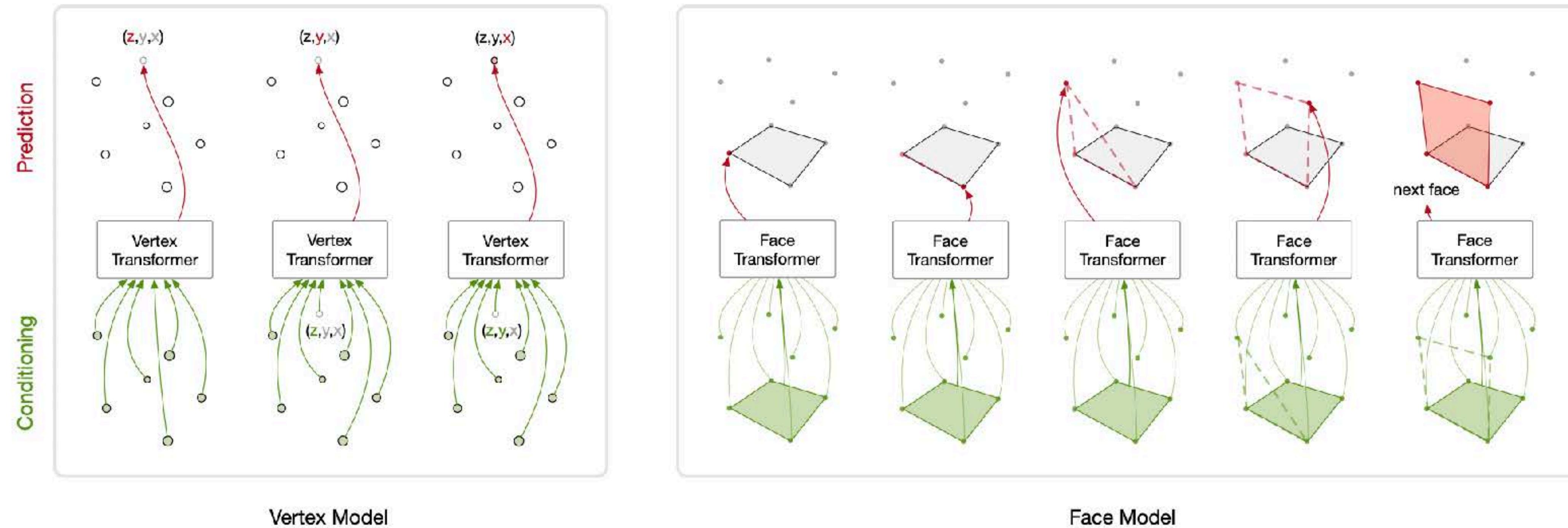
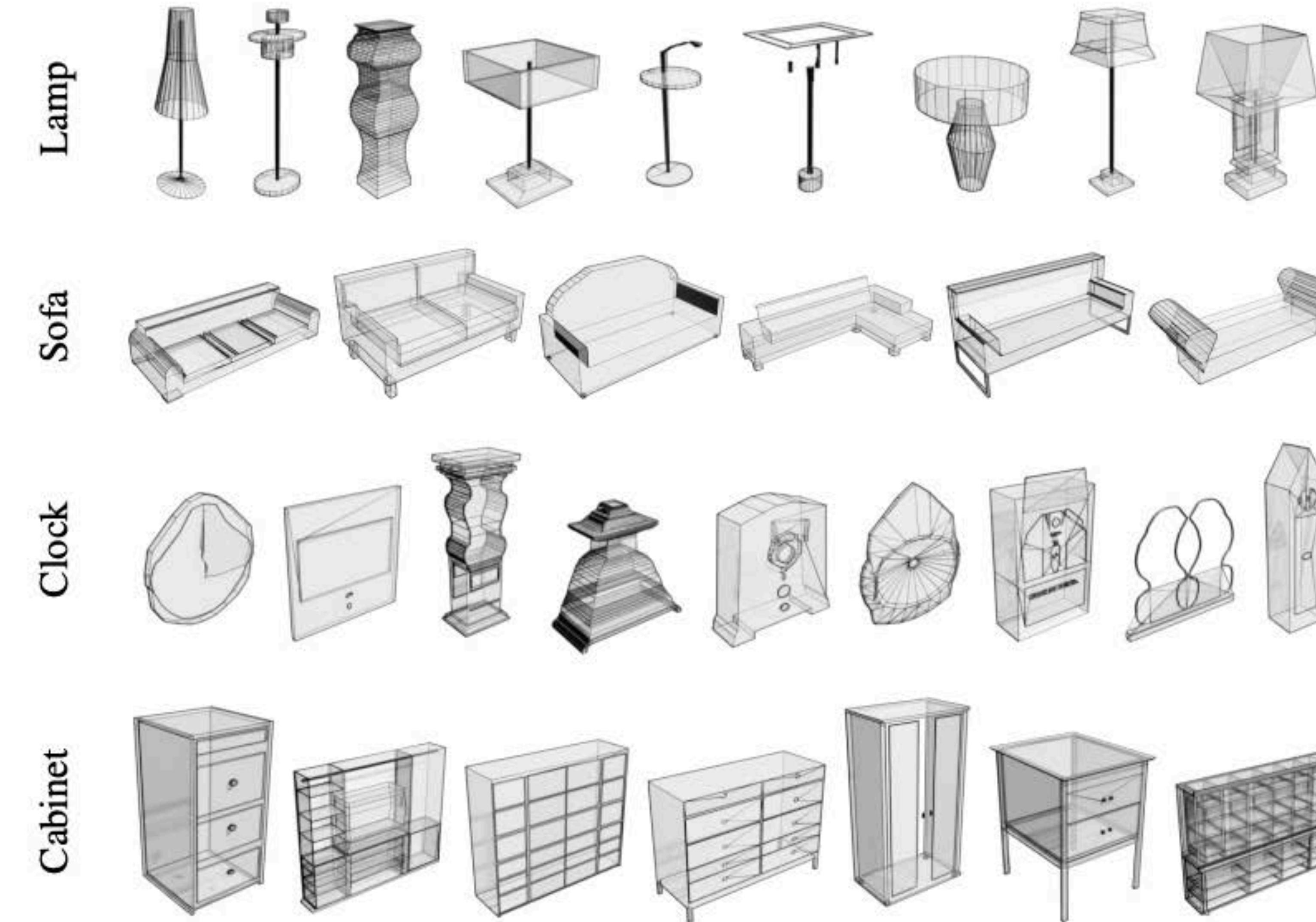


Figure 2. PolyGen first generates mesh vertices (left), and then generates mesh faces conditioned on those vertices (right). Vertices are generated sequentially from lowest to highest on the vertical axis. To generate the next vertex the current sequence of vertex coordinates is passed as context to a vertex Transformer, which outputs a predictive distribution for the next vertex coordinate. The face model takes as input a collection of vertices, and the current sequence of face indices, and outputs a distribution over vertex indices.

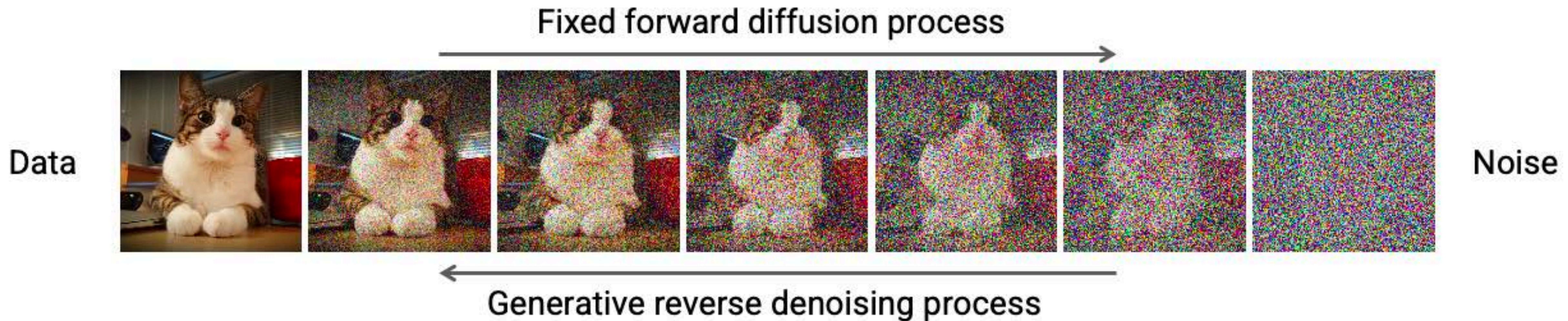
Autoregressive Generation of 3D Shapes



Outline

- 3D GAN
- 3D Autoregressive Models
- 3D Diffusion Models
- 3D Generation without 3D Training Data
- Part-based 3D Generation

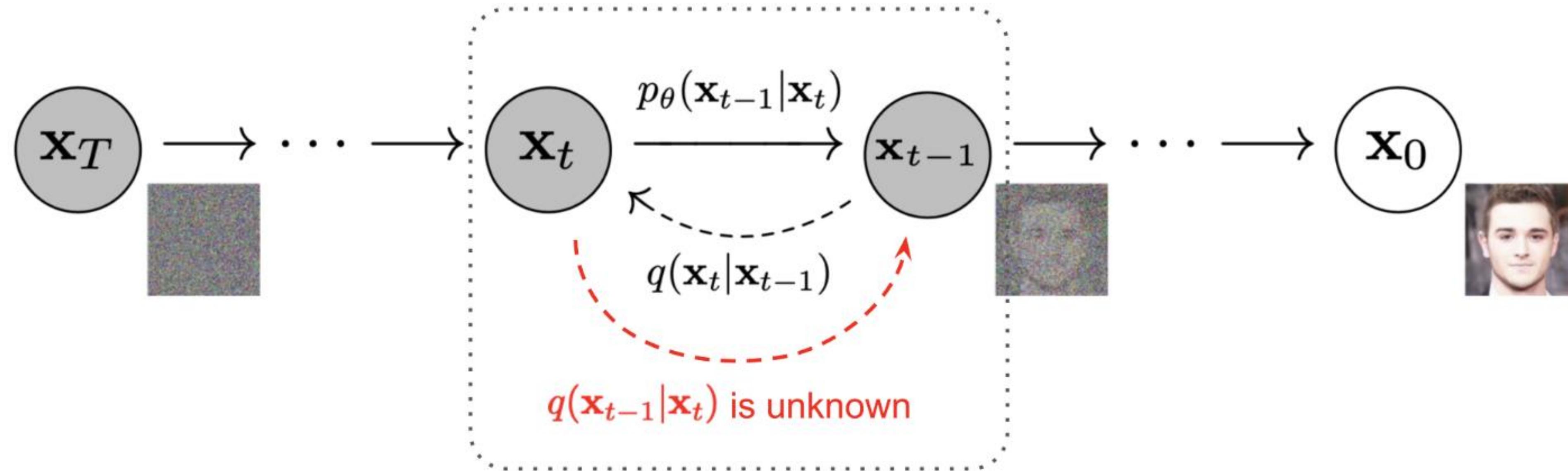
Background: Denoising Diffusion Models



synthesis by progressive denoising

Train a neural network to learn the reverse process

Background: Denoising Diffusion Models



Algorithm 1 Training

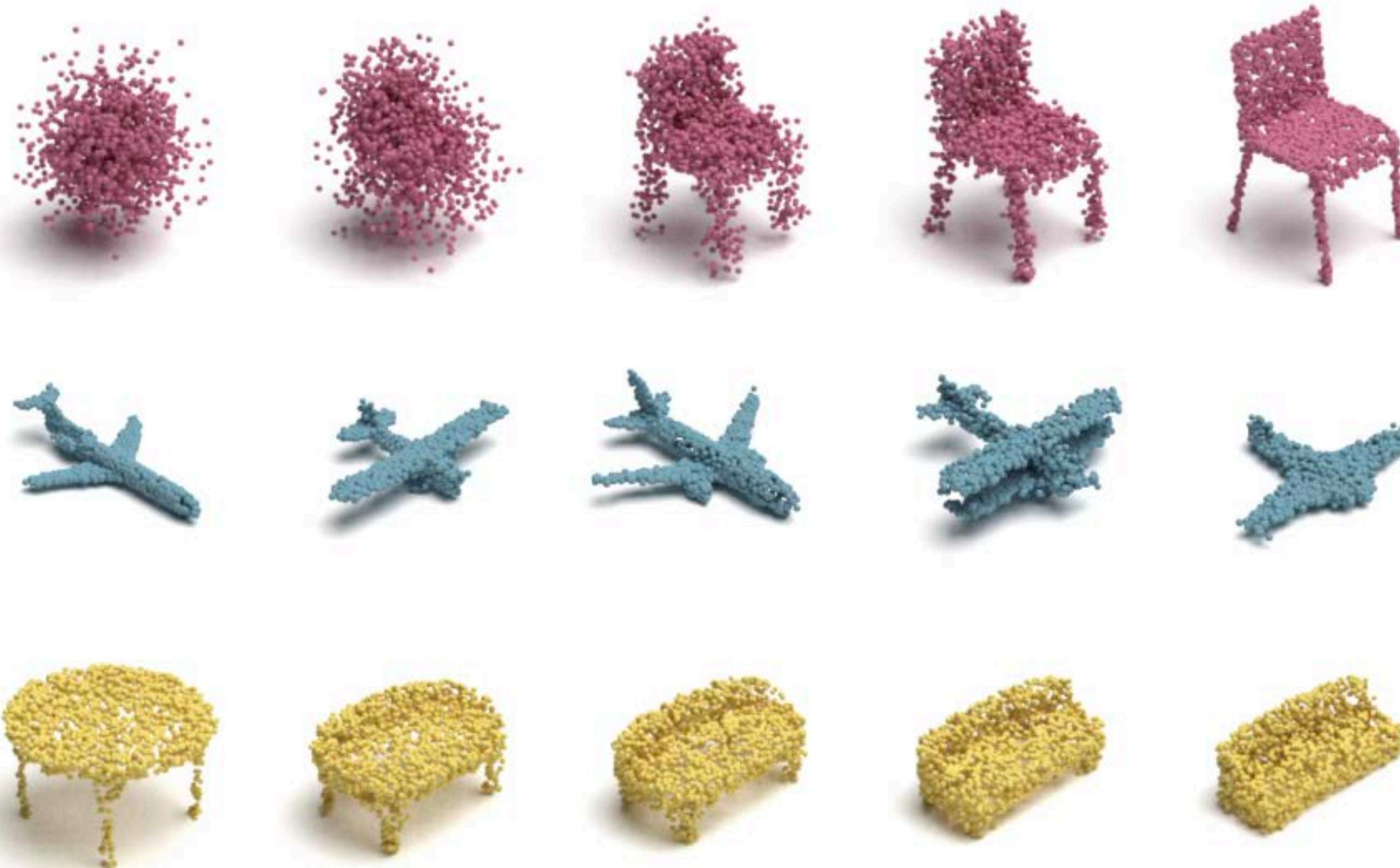
```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
        
$$\nabla_\theta \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$

6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Point Cloud Diffusion



Algorithm 1 Training (Simplified)

```
1: repeat
2:   Sample  $\mathbf{X}^{(0)} \sim q_{\text{data}}(\mathbf{X}^{(0)})$ 
3:   Sample  $\mathbf{z} \sim q_{\varphi}(\mathbf{z}|\mathbf{X}^{(0)})$ 
4:   Sample  $t \sim \text{Uniform}(\{1, \dots, T\})$ 
5:   Sample  $\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_N^{(t)} \sim q(\mathbf{x}^{(t)}|\mathbf{x}^{(0)})$ 
6:    $L_t \leftarrow \sum_{i=1}^N D_{\text{KL}} \left( q(\mathbf{x}_i^{(t-1)}|\mathbf{x}_i^{(t)}, \mathbf{x}_i^{(0)}) \middle\| p_{\theta}(\mathbf{x}_i^{(t-1)}|\mathbf{x}_i^{(t)}, \mathbf{z}) \right)$ 
7:    $L_z \leftarrow D_{\text{KL}}(q_{\varphi}(\mathbf{z}|\mathbf{X}^{(0)}) \| p(\mathbf{z}))$ 
8:   Compute  $\nabla_{\theta}(L_t + \frac{1}{T} L_z)$ . Then perform gradient descent.
9: until converged
```

Latent Diffusion for 3D Shape Generation

Diffusion-SDF: Conditional Generative Modeling of Signed Distance Functions

Gene Chou
Princeton University

Yuval Bahat
Princeton University

Felix Heide
Princeton University

SDFusion: Multimodal 3D Shape Completion, Reconstruction, and Generation

Yen-Chi Cheng¹ Hsin-Ying Lee² Sergey Tulyakov² Alexander Schwing^{1*} Liangyan Gui^{1*}
¹University of Illinois Urbana-Champaign ²Snap Research

Diffusion-SDF: Text-to-Shape via Voxelized Diffusion

Muheng Li¹, Yueqi Duan^{†,2}, Jie Zhou¹, Jiwen Lu¹

¹Department of Automation, Tsinghua University

²Department of Electronic Engineering, Tsinghua University

Several (concurrent) works exploring this, with essentially similar approach...

Background: Stable Diffusion

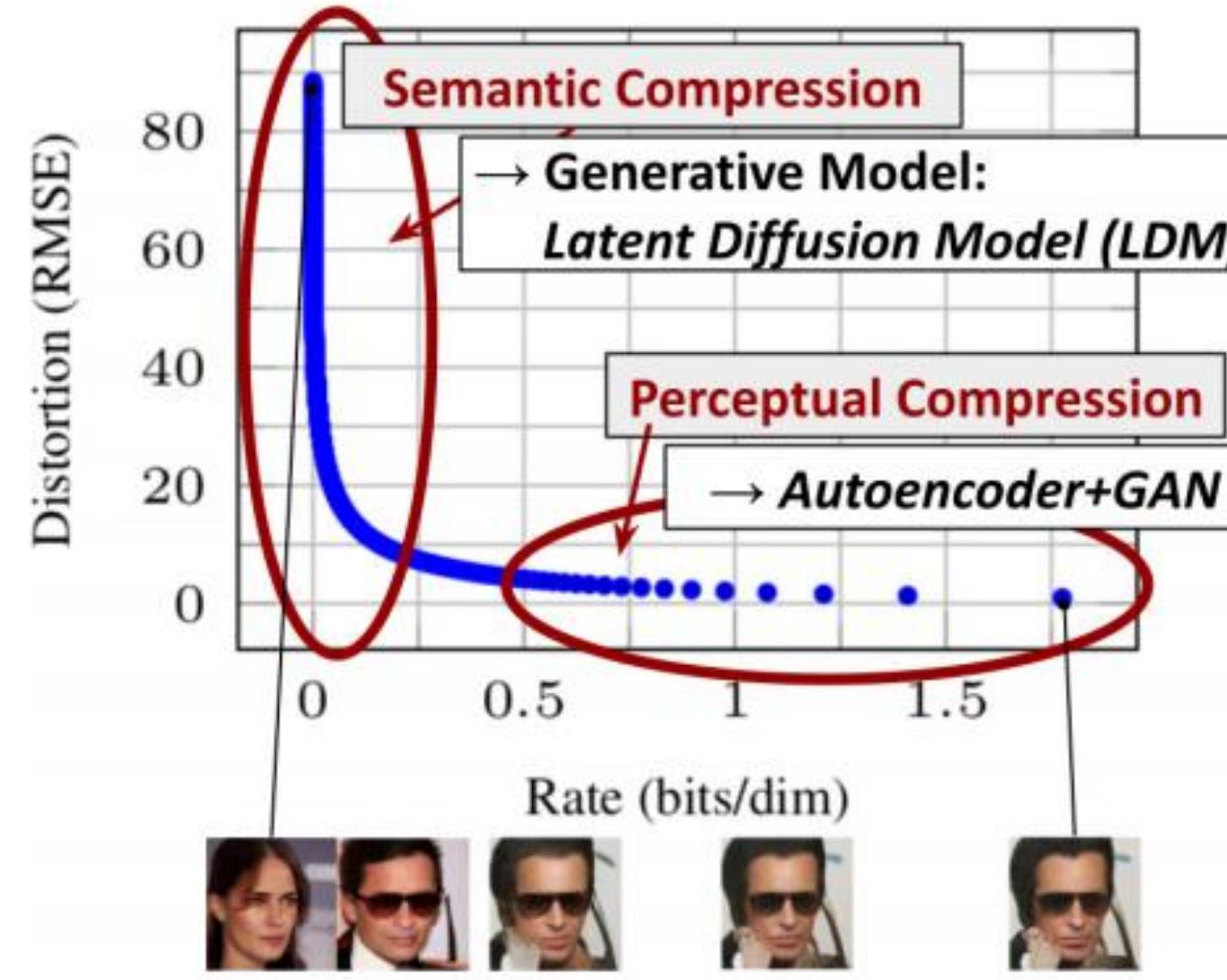
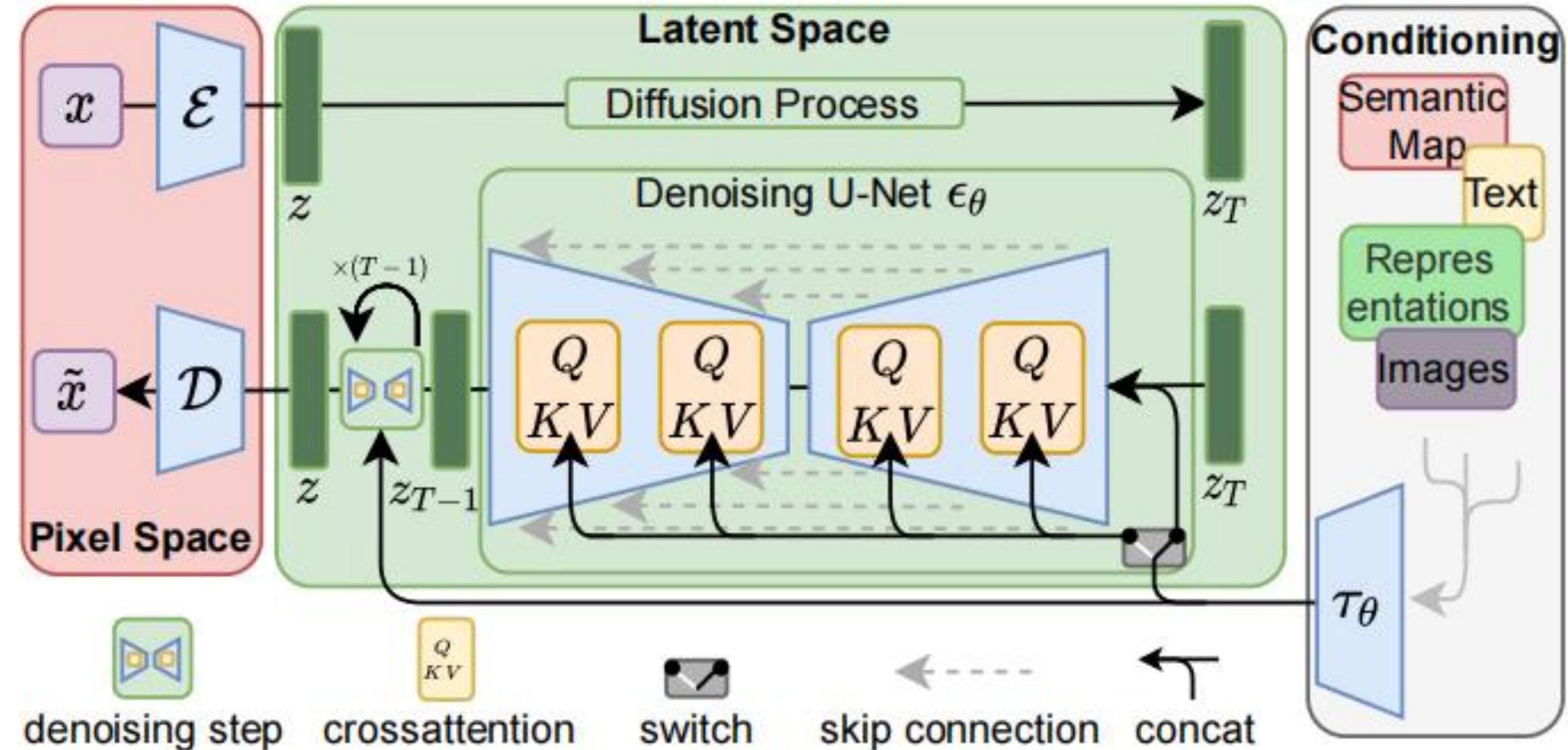
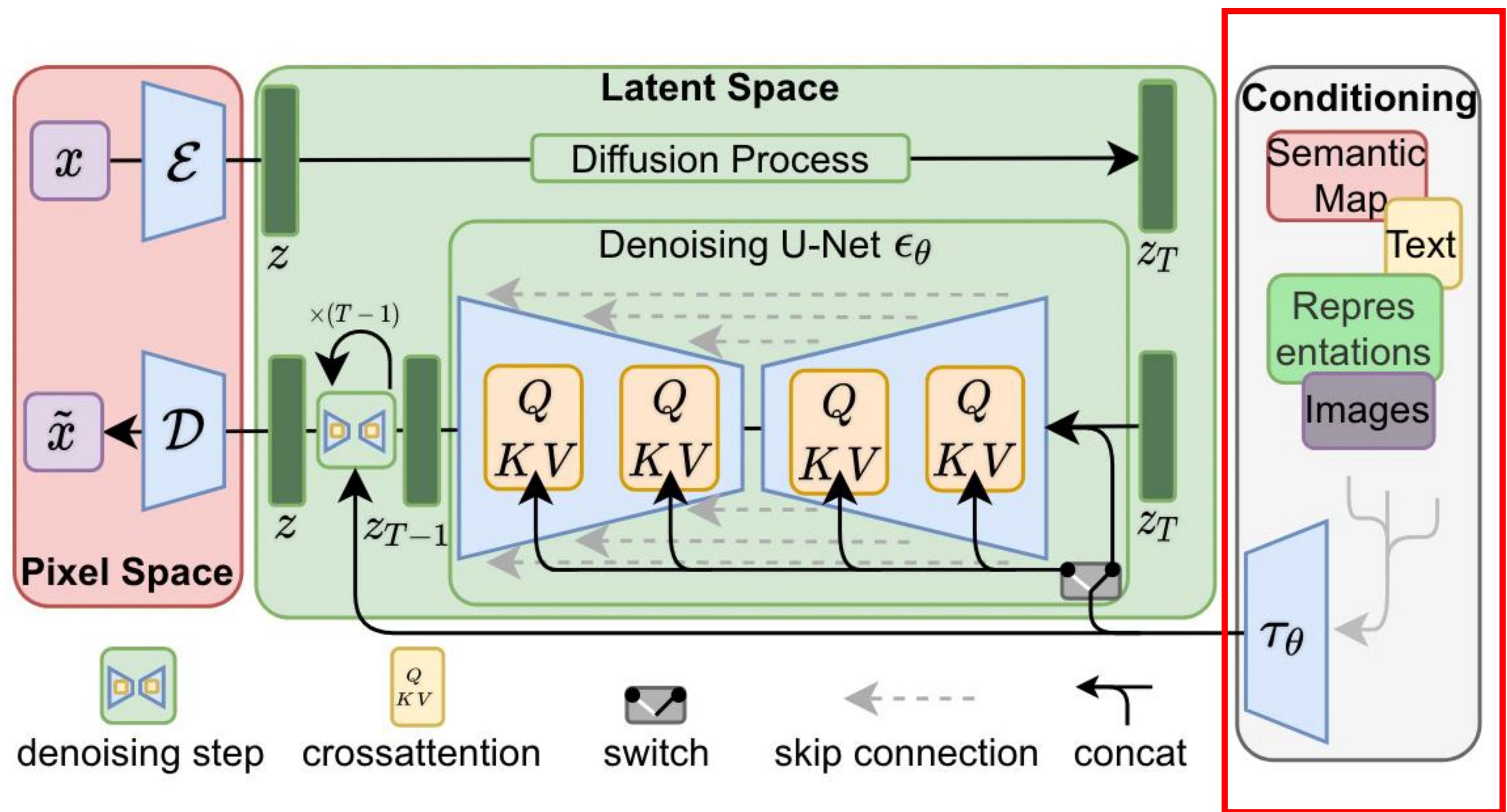


Figure 2. Illustrating perceptual and semantic compression: Most bits of a digital image correspond to imperceptible details. While DMs allow to suppress this semantically meaningless information by minimizing the responsible loss term, gradients (during training) and the neural network backbone (training and inference) still need to be evaluated on all pixels, leading to superfluous computations and unnecessarily expensive optimization and inference. We propose *latent diffusion models (LDMs)* as an effective generative model and a separate mild compression stage that only eliminates imperceptible details. Data and images from [30].

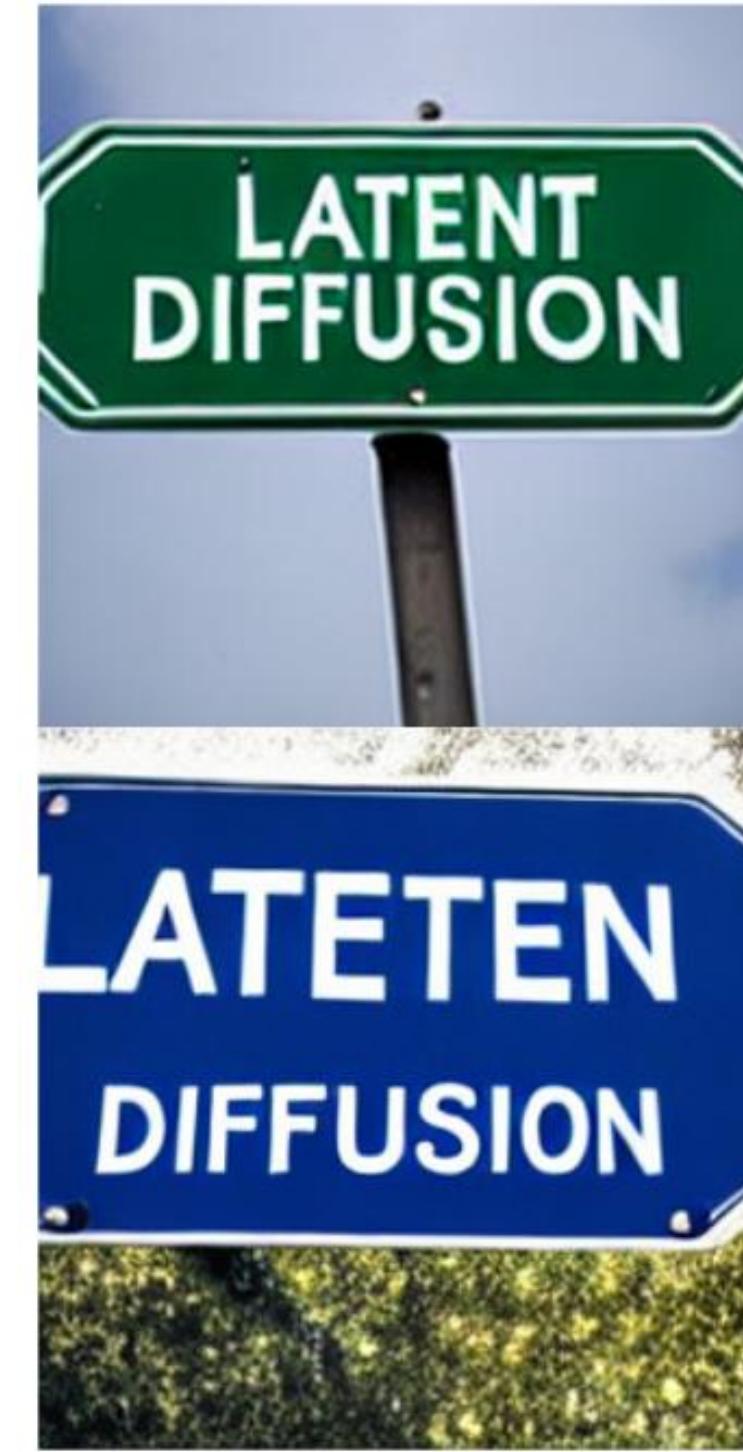


- VAE-encoder + diffusion model in latent space + VAE-decoder = Stable Generation

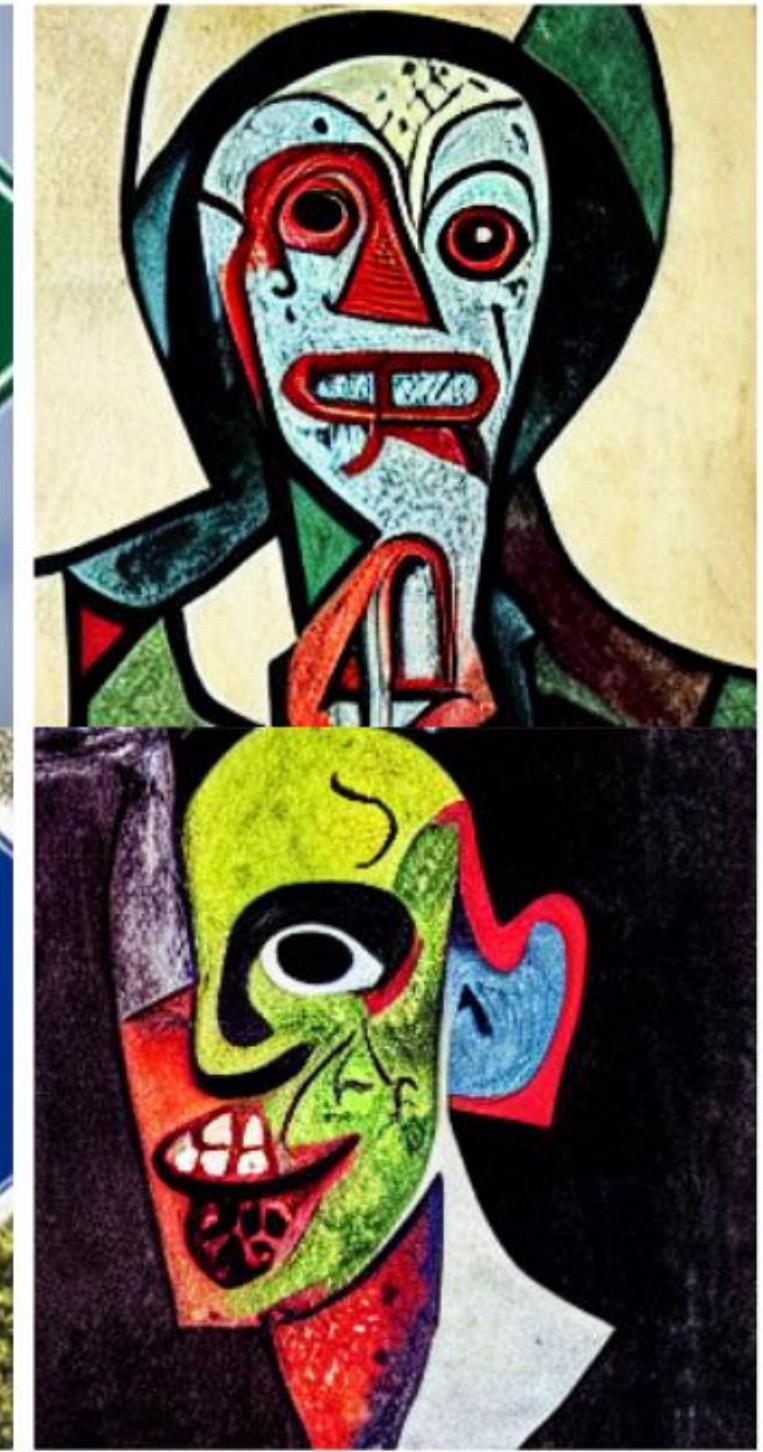
Background: Conditional Diffusion



'A street sign that reads
“Latent Diffusion”'



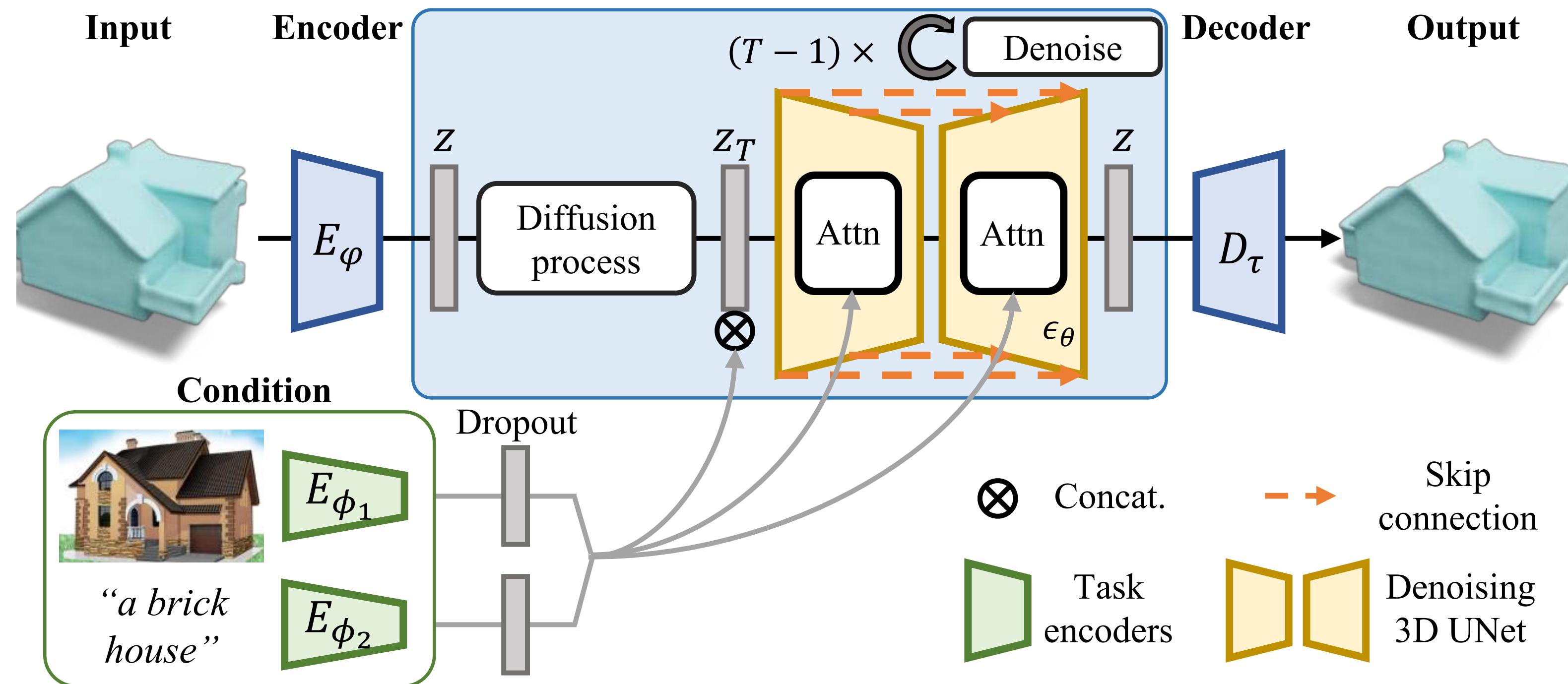
'A zombie in the
style of Picasso'



$$\mathbb{E}_{\mathbf{x}_0, \mathbf{y} \sim \tilde{p}(\mathbf{x}_0, \mathbf{y}), \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\left\| \boldsymbol{\varepsilon} - \boldsymbol{\epsilon}_\theta(\bar{\alpha}_t \mathbf{x}_0 + \bar{\beta}_t \boldsymbol{\varepsilon}, \mathbf{y}, t) \right\|^2 \right] \quad (22)$$

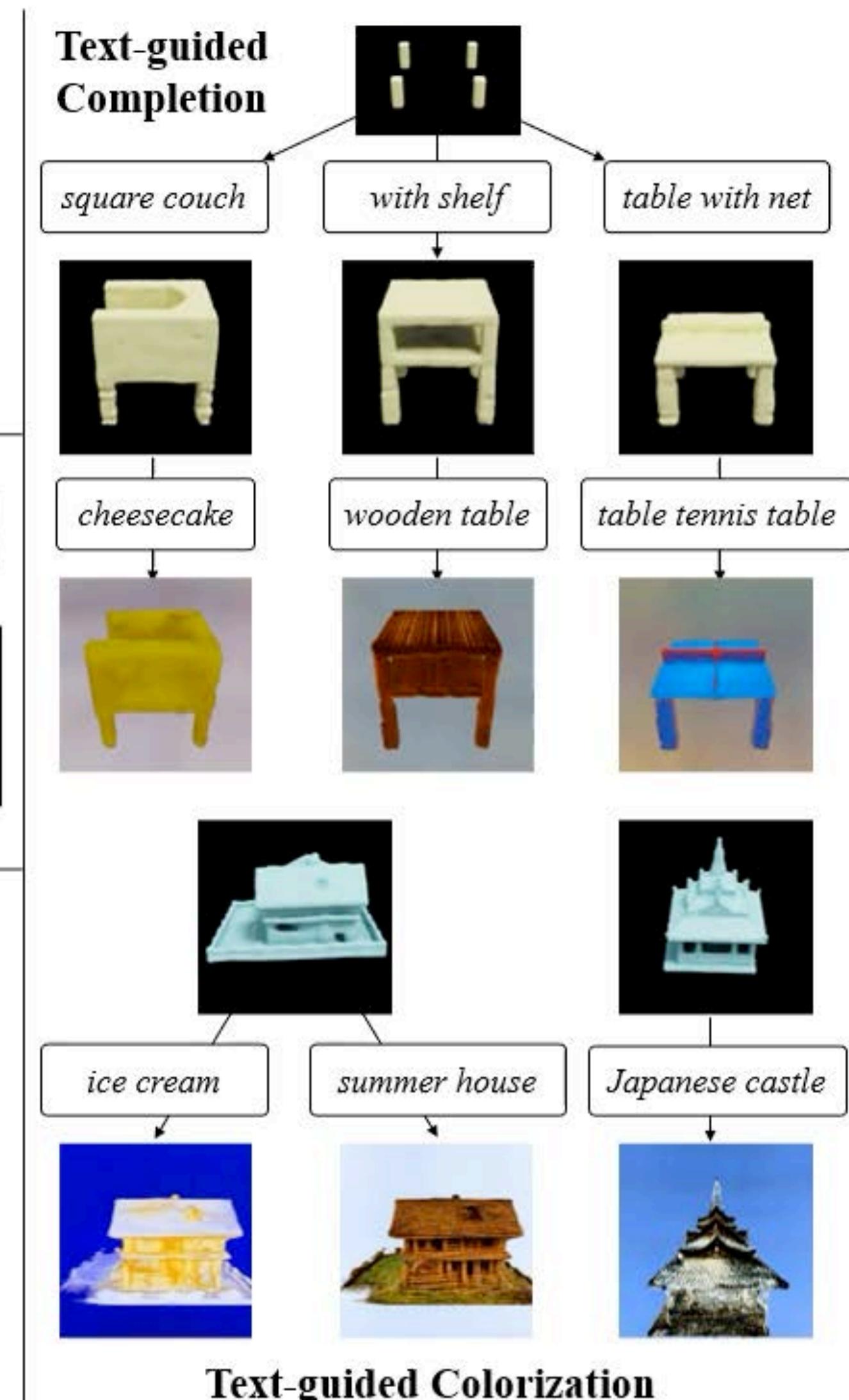
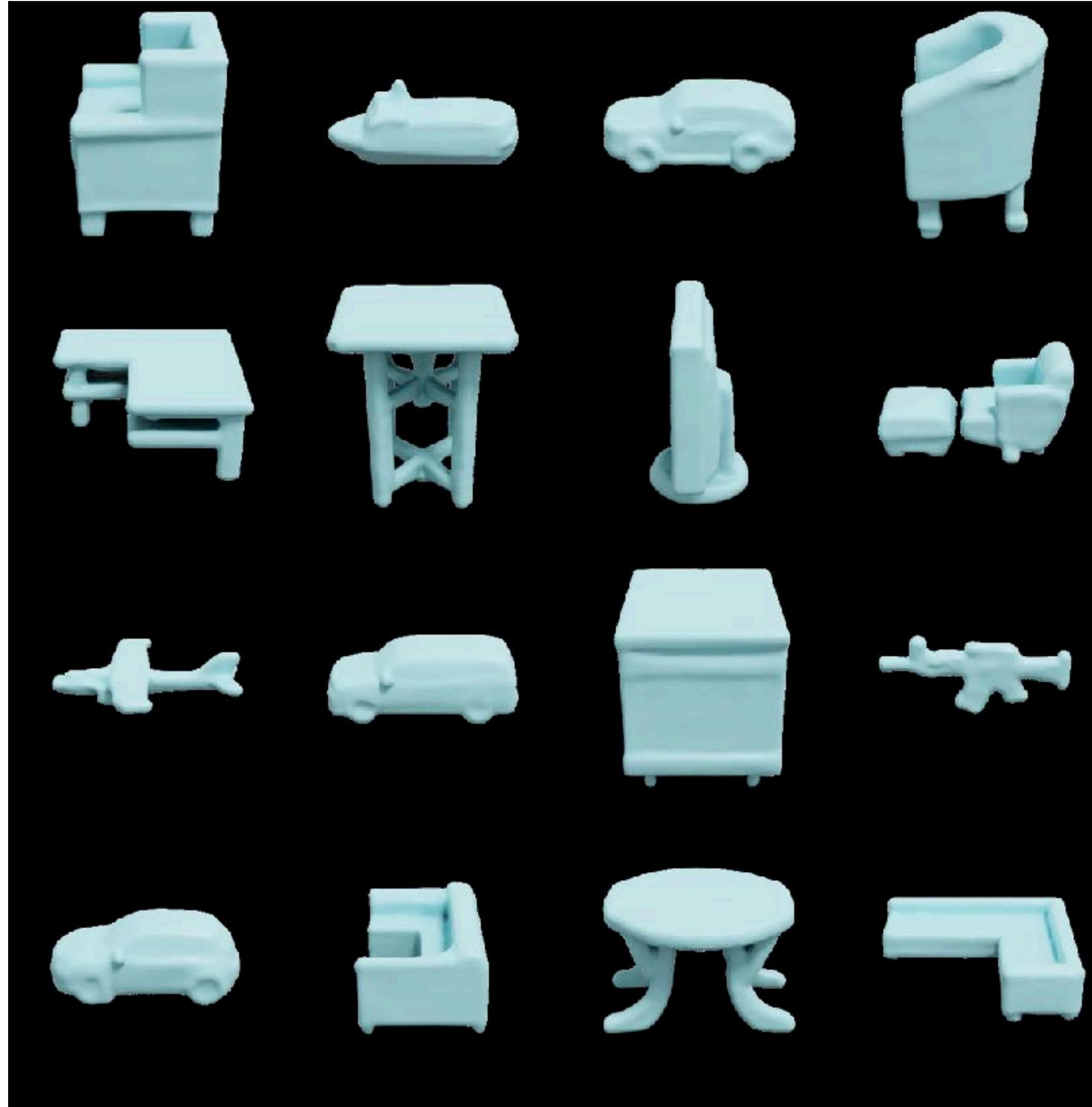
feed condition information \mathbf{y}
To de-noise network in training

Latent Diffusion for 3D Shape Generation



3D Diffusion model in latent space, with conditioning via attention layers

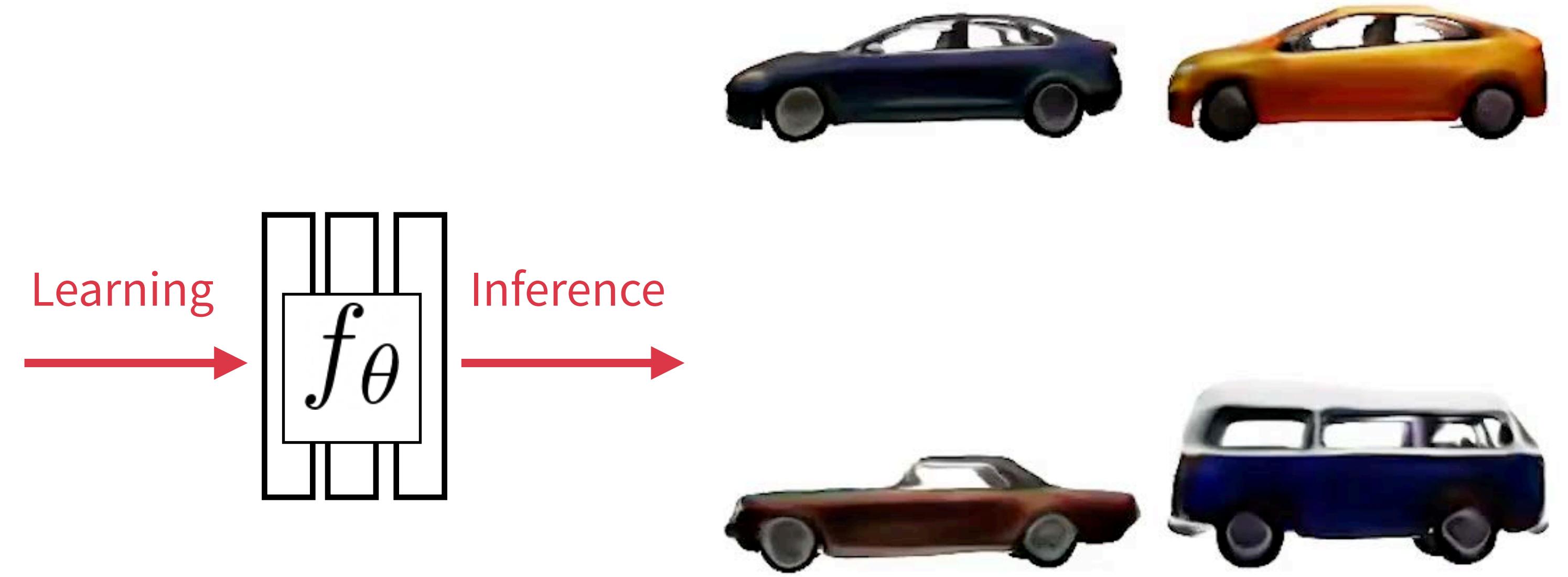
Latent Diffusion for 3D Shape Generation



Outline

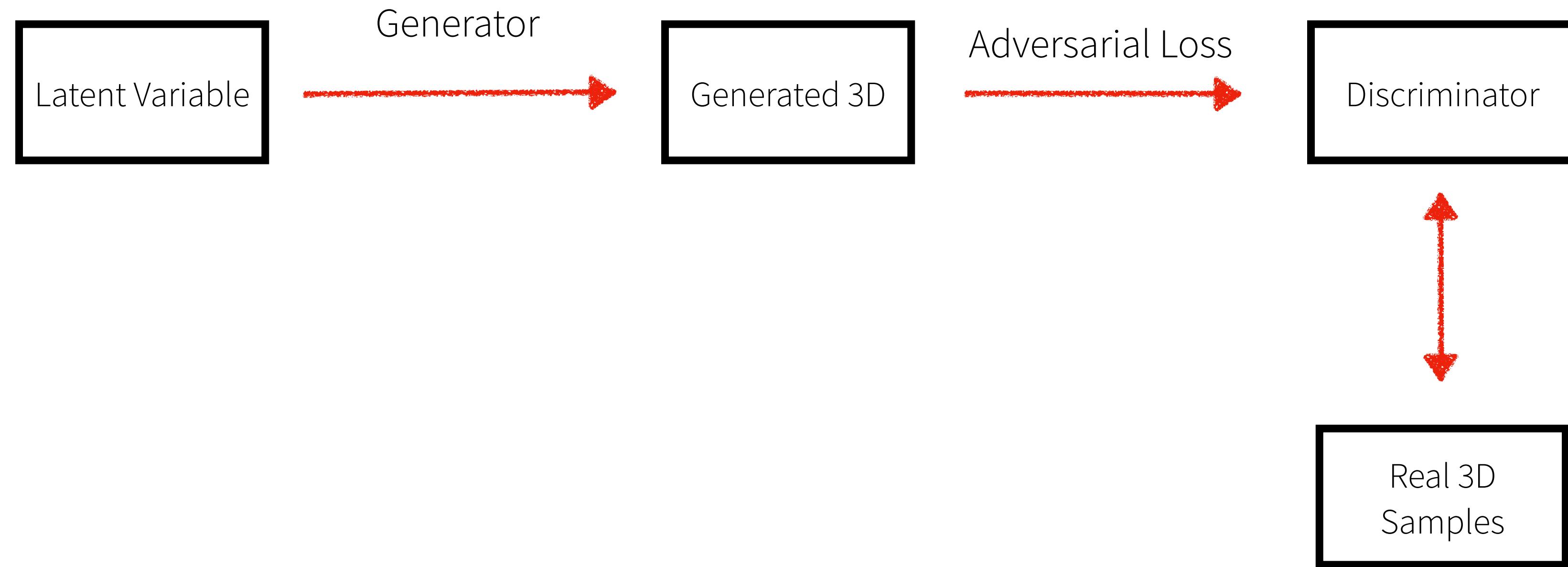
- 3D GAN
- 3D Autoregressive Models
- 3D Diffusion Models
- 3D Generation without 3D Training Data
- Part-based 3D Generation

Generating 3D Shapes without 3D training data



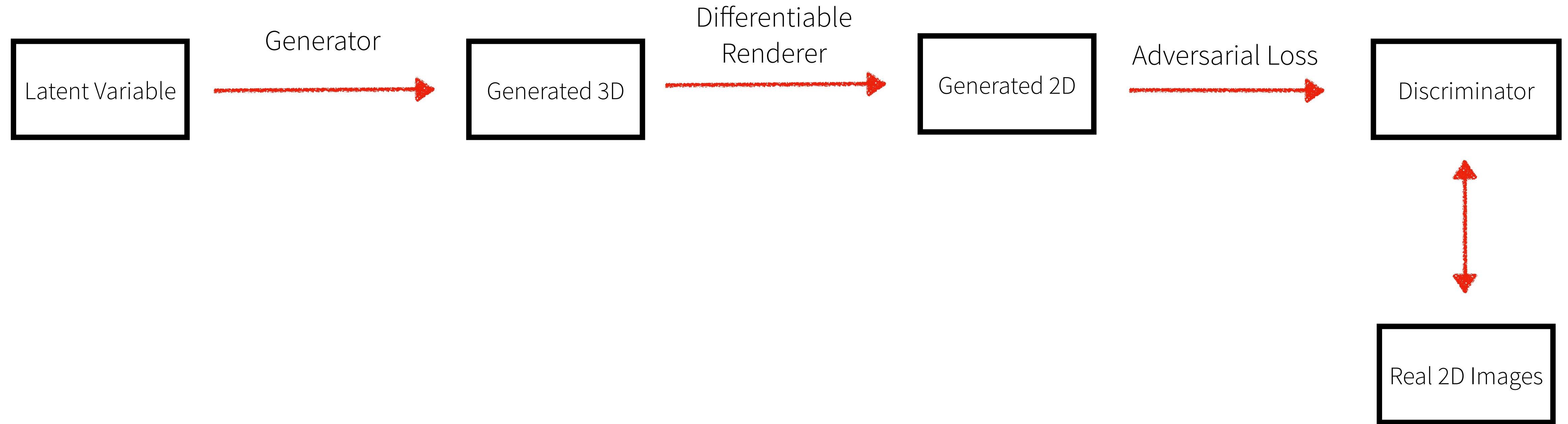
Can we learn a generative 3D model from 2D images?

Generating 3D Shapes without 3D training data



Key Idea: Generate 3D representations such that —
they are indistinguishable from real samples

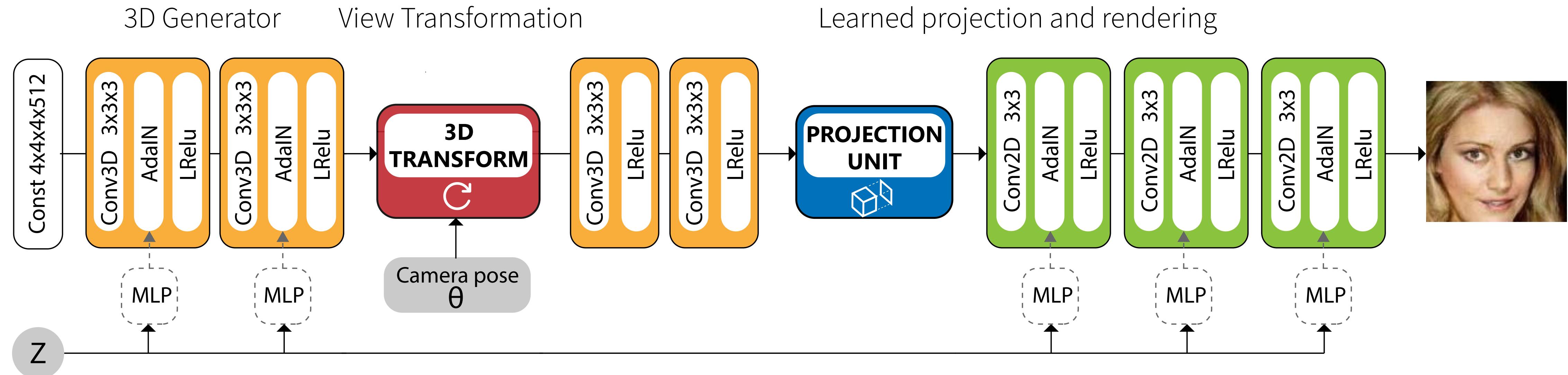
Generating 3D Shapes without 3D training data



Key Idea: Generate 3D representations such that —

their renderings are indistinguishable from real samples

Generating 3D using 2D Adversarial Nets



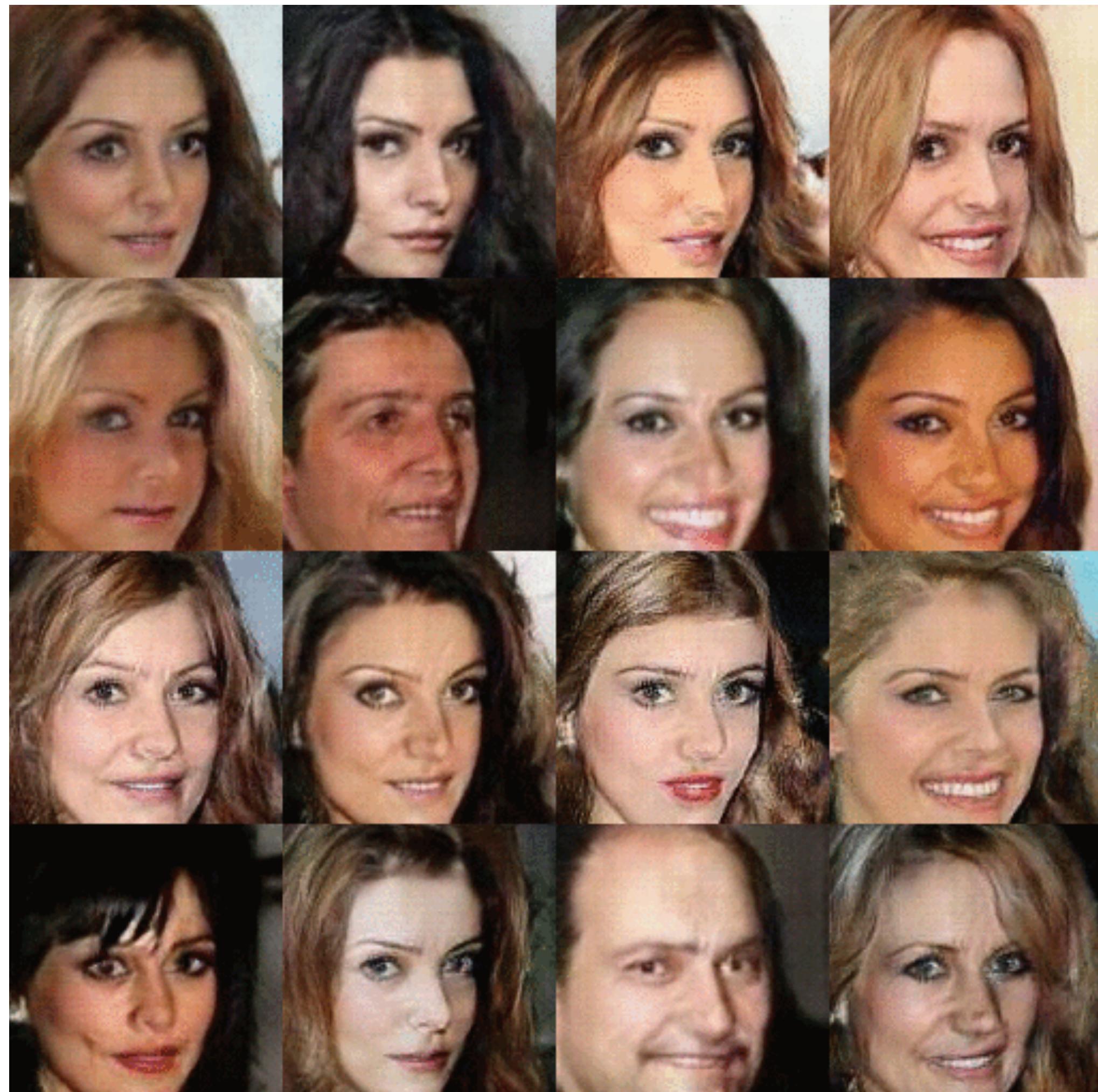
Key Idea: Generate 3D representations such that —

their renderings are indistinguishable from real samples

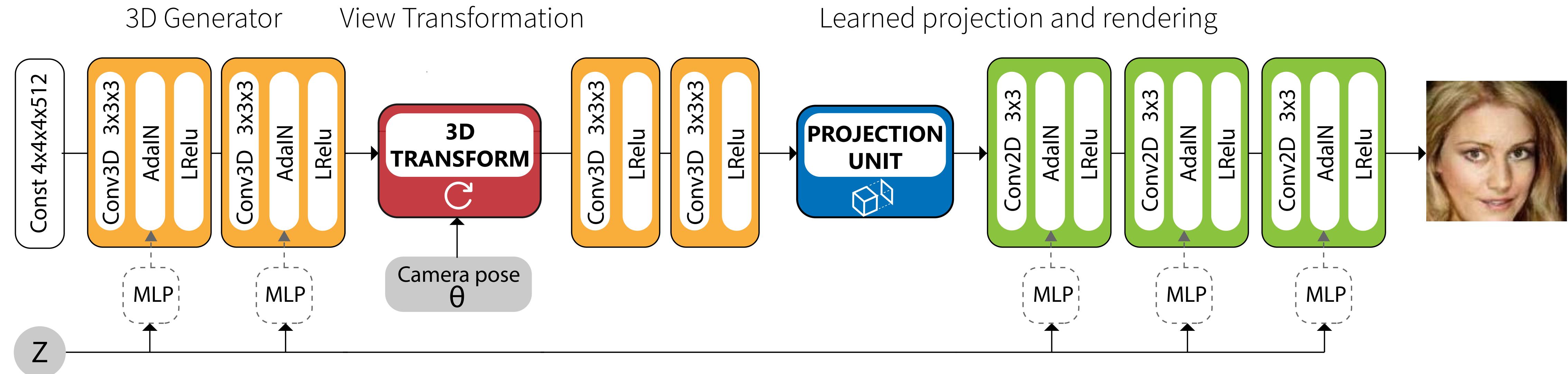
Generating 3D using 2D Adversarial Nets



Generating 3D using 2D Adversarial Nets



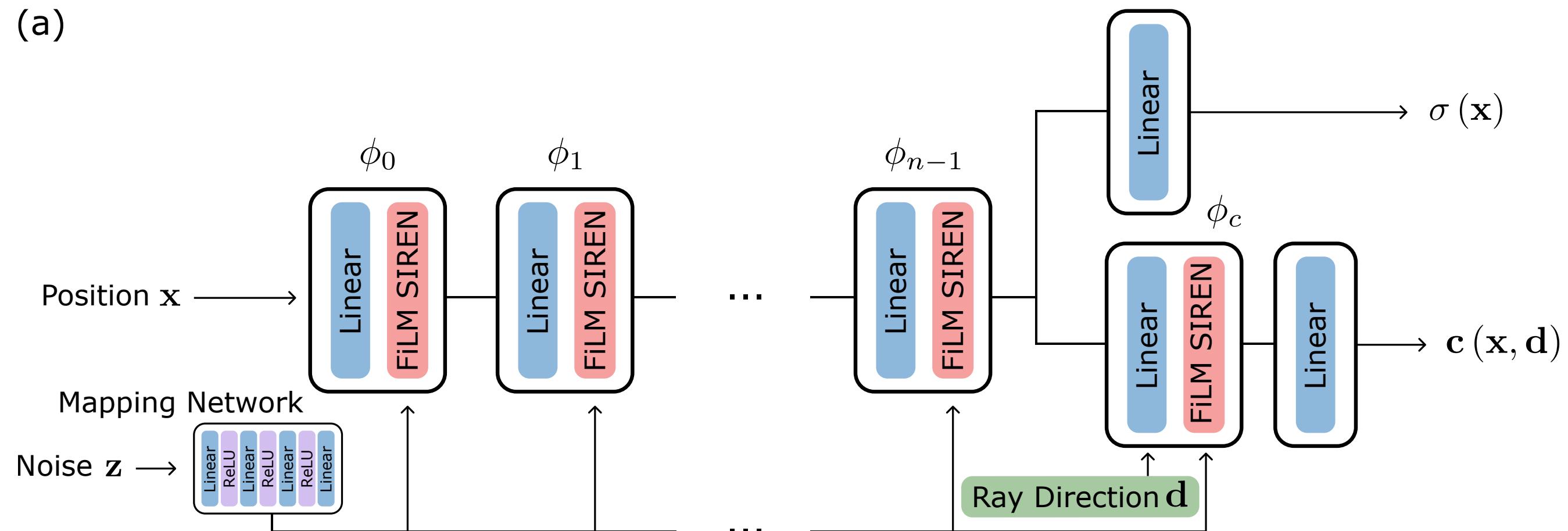
Generating 3D using 2D Adversarial Nets



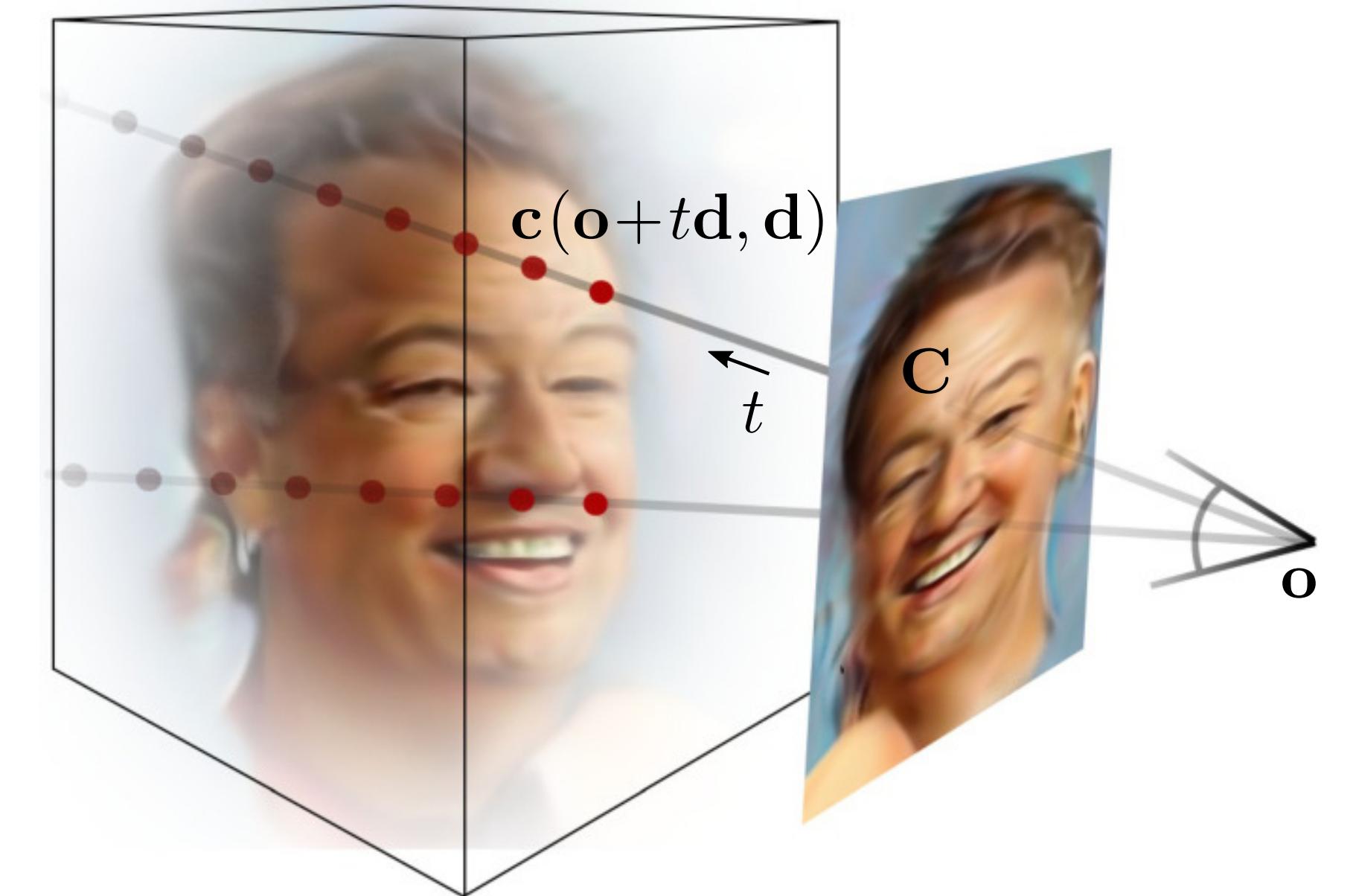
Some inconsistency across views

Some 3D inducting biases, but not enough (e.g. learned instead of volumetric rendering)

Generating 3D using 2D Adversarial Nets



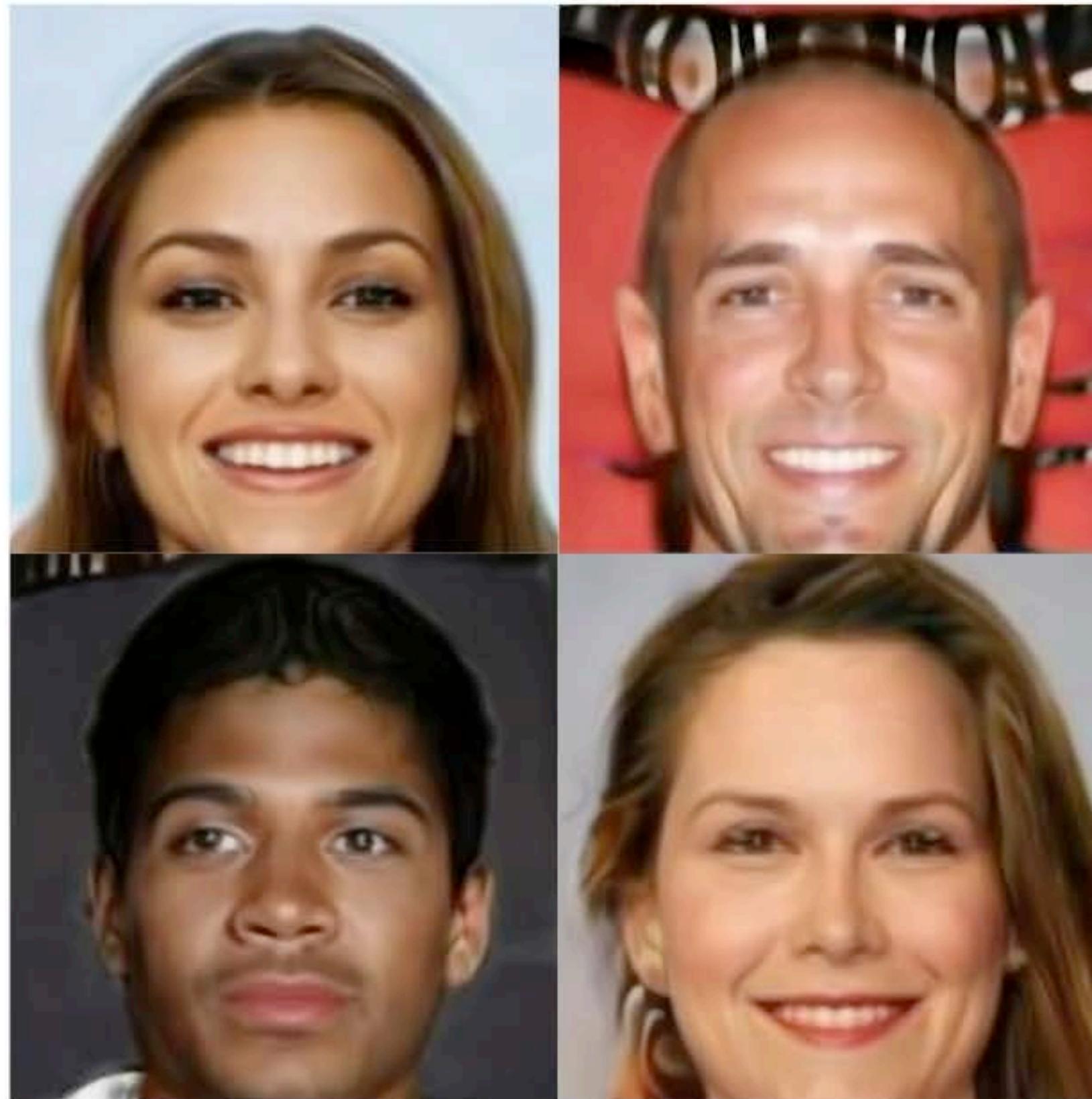
Generator: A conditional NeRF model



Renderer: Volume rendering

Generating 3D using 2D Adversarial Nets

CelebA



Cats



CARLA



Generating 3D using 2D Adversarial Nets

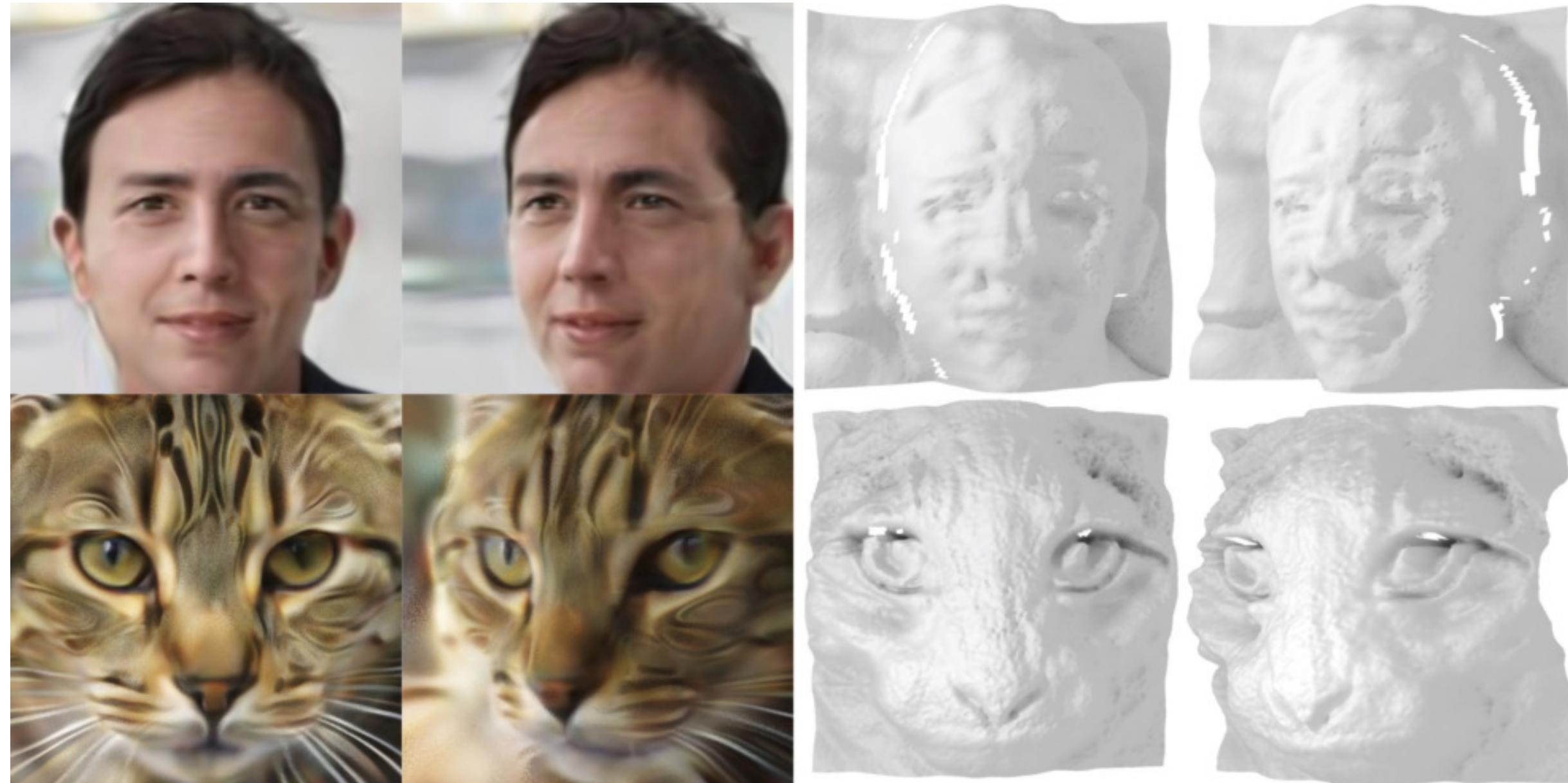


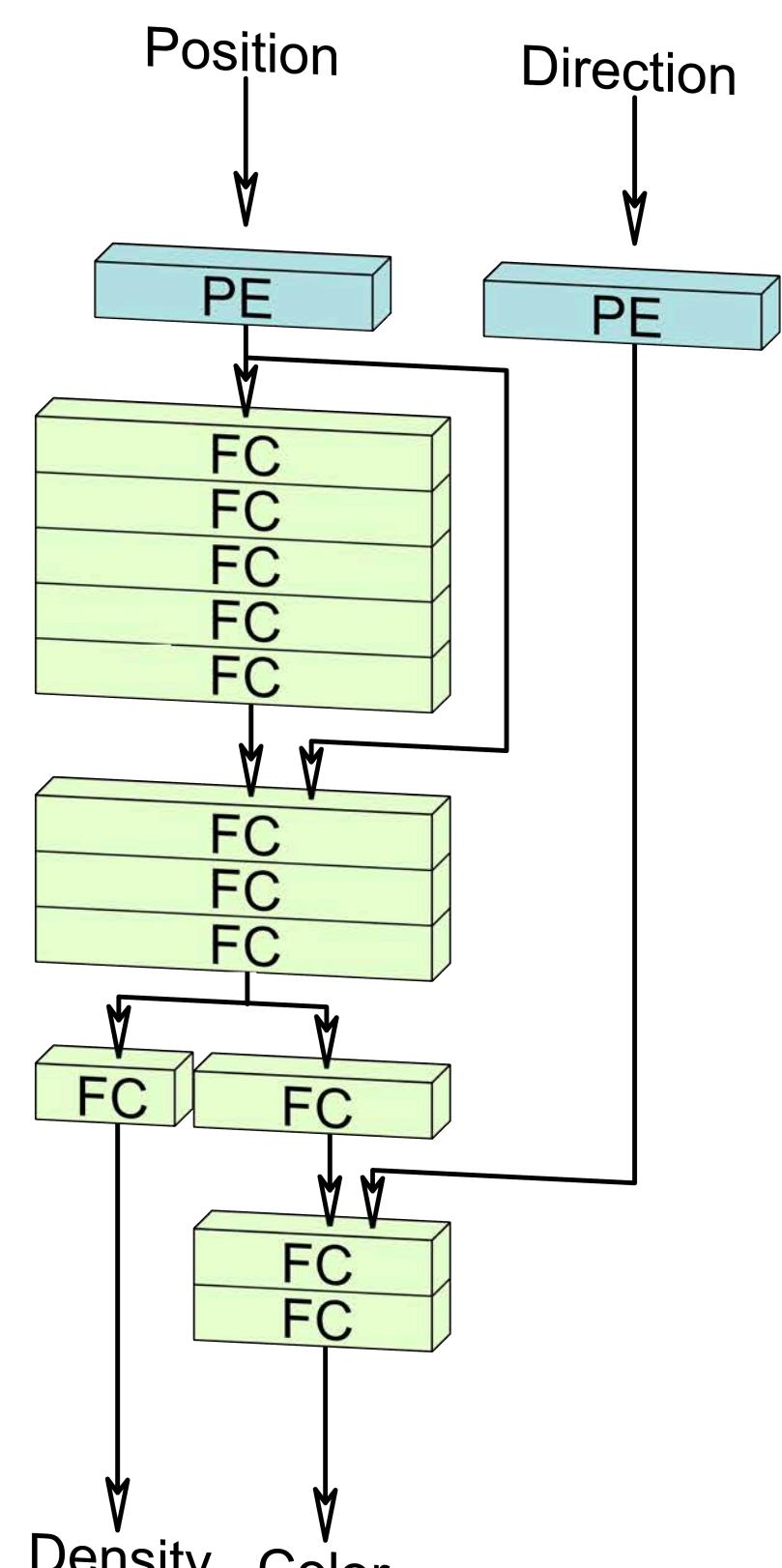
Figure 12: In a failure case reminiscent of the hollow-face illusion, our model sometimes generates objects with inverted sections.



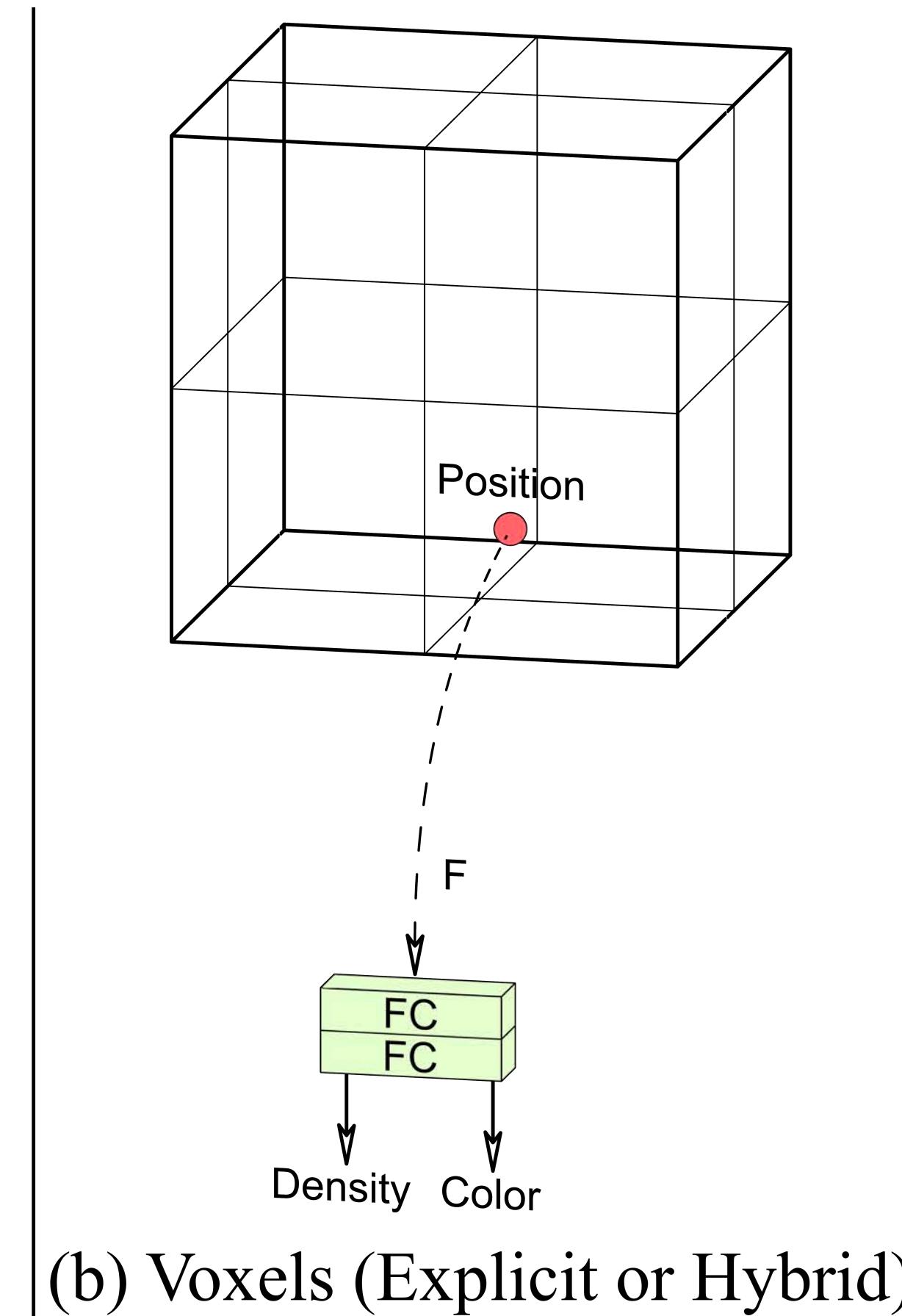
Also possible to extract some geometry

Although not always great

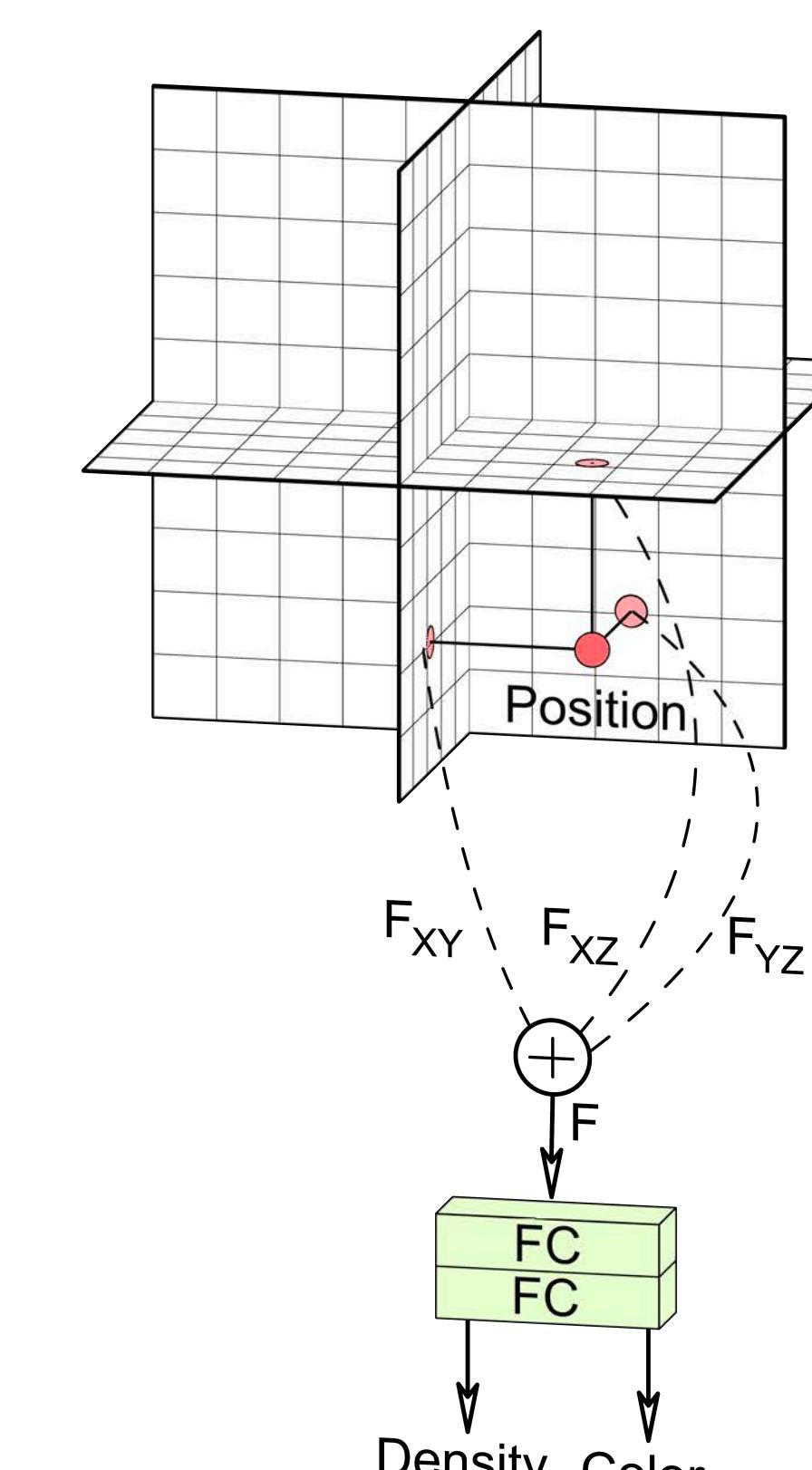
Generating 3D using 2D Adversarial Nets



(a) NeRF (Implicit)



(b) Voxels (Explicit or Hybrid)



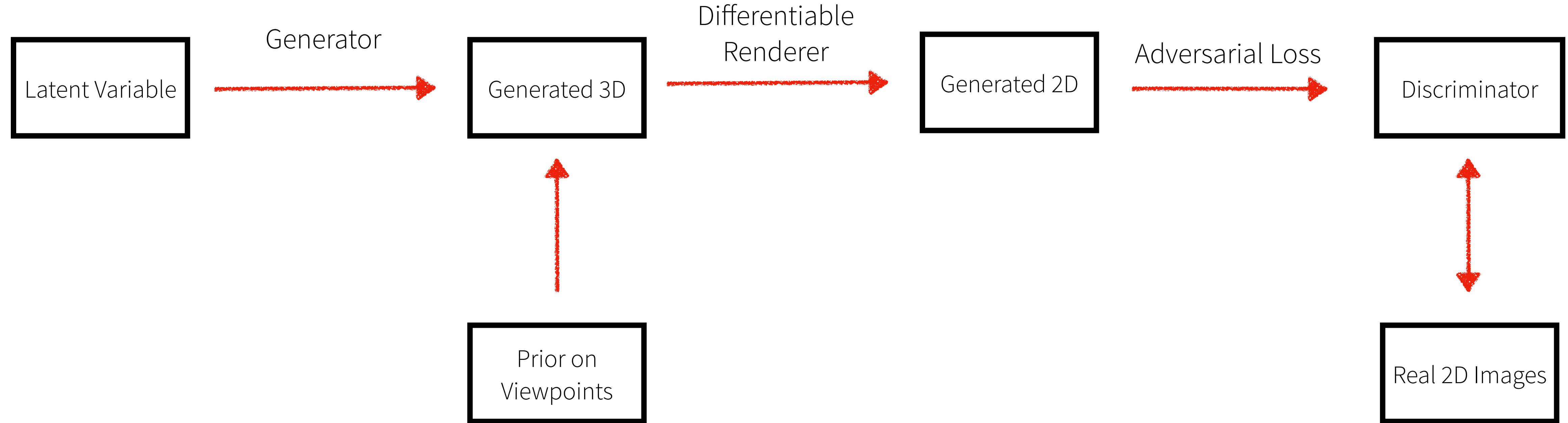
(c) Ours (Hybrid)

Replace NeRFs via hybrid representations

Generating 3D using 2D Adversarial Nets



Generating 3D Shapes without 3D training data



Crucial to leverage the 3D structure of the problem

Assumed known view distribution e.g. faces seen frontally

Typical success stories on ‘nice’ categories (lots of clean images, known views...)

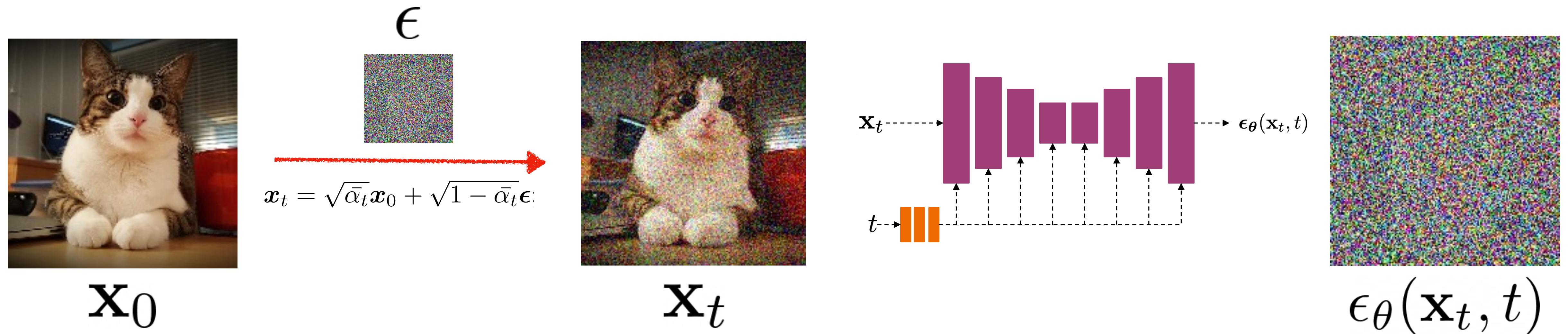
Optimizing Text-conditioned 3D



Let's see more combinations:
<https://dreamfusion3d.github.io/>

“A squirrel wearing a kimono and riding a motorbike”

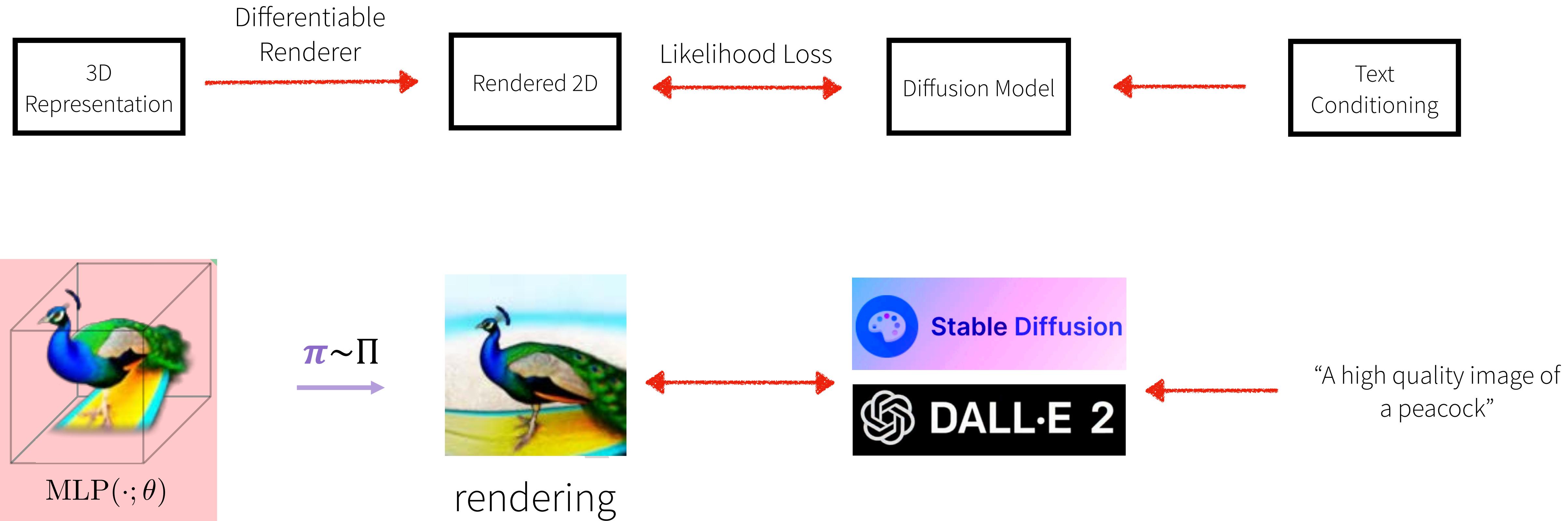
Background: Likelihood in Diffusion Models

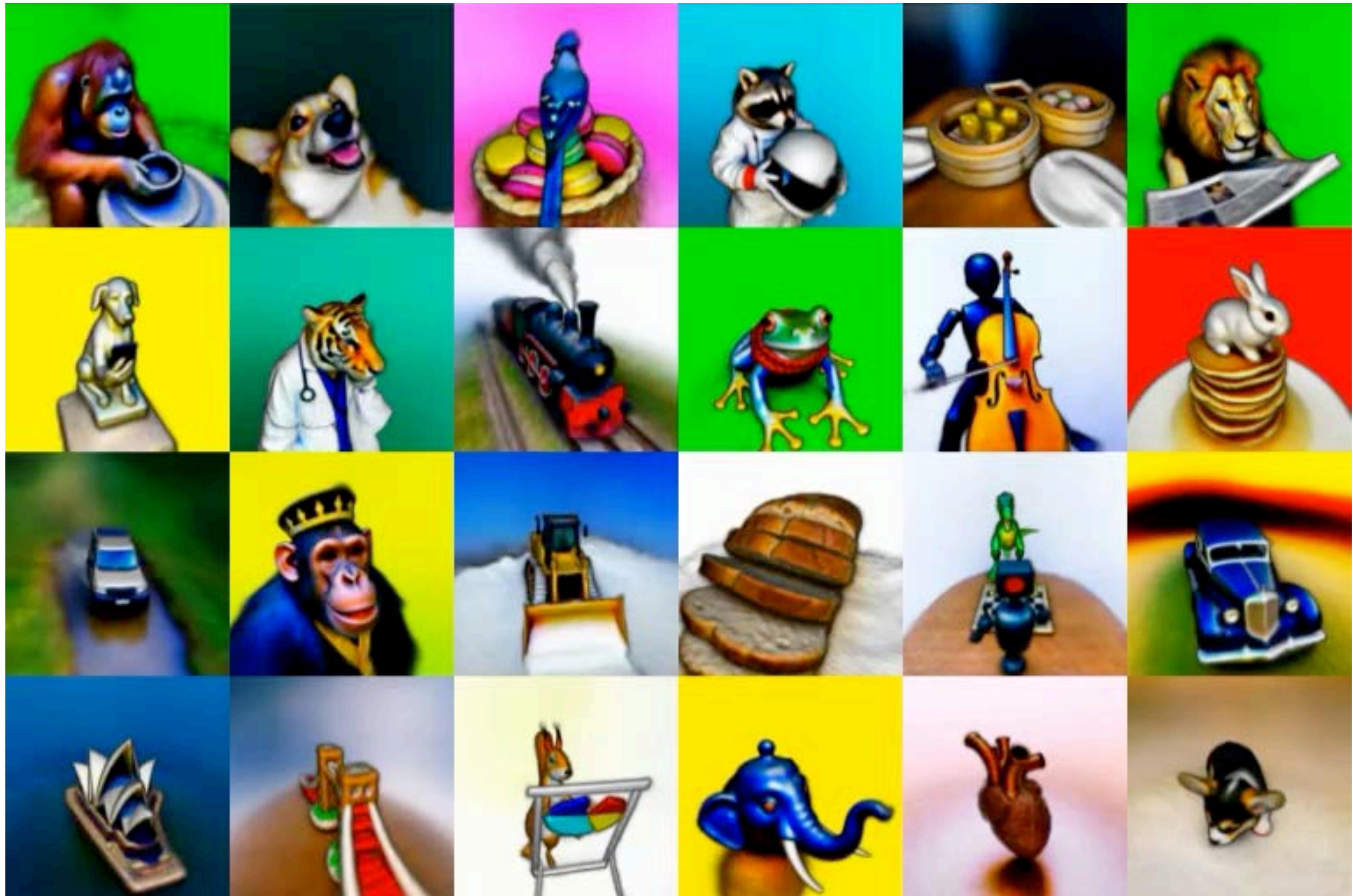


Training Objective: Noise prediction error $\|\epsilon_\theta(\mathbf{x}_t, t) - \epsilon\|^2$

Prediction error also be correlated to $\mathbf{p}(\mathbf{x})$ (images with high likelihood yield lower error)

Optimizing Text-conditioned 3D





Optimizing Text-conditioned 3D



Allows 3D generation for complex structures

Not a ‘generative’ model in a probabilistic sense — no distribution over 3D is inferred

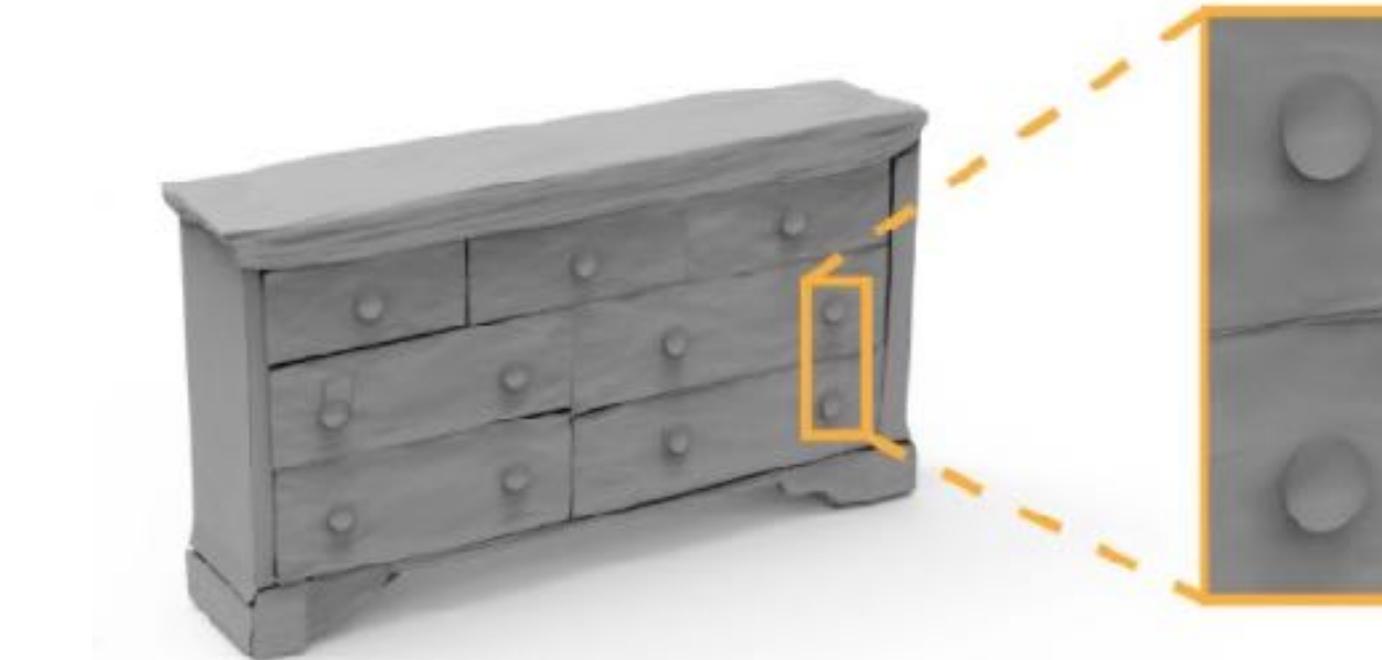
(relatively) compute intensive — per instance optimization

Outline

- 3D GAN
- 3D Autoregressive Models
- 3D Diffusion Models
- 3D Generation without 3D Training Data
- Part-based 3D Generation

Part-based 3D Shape Generation

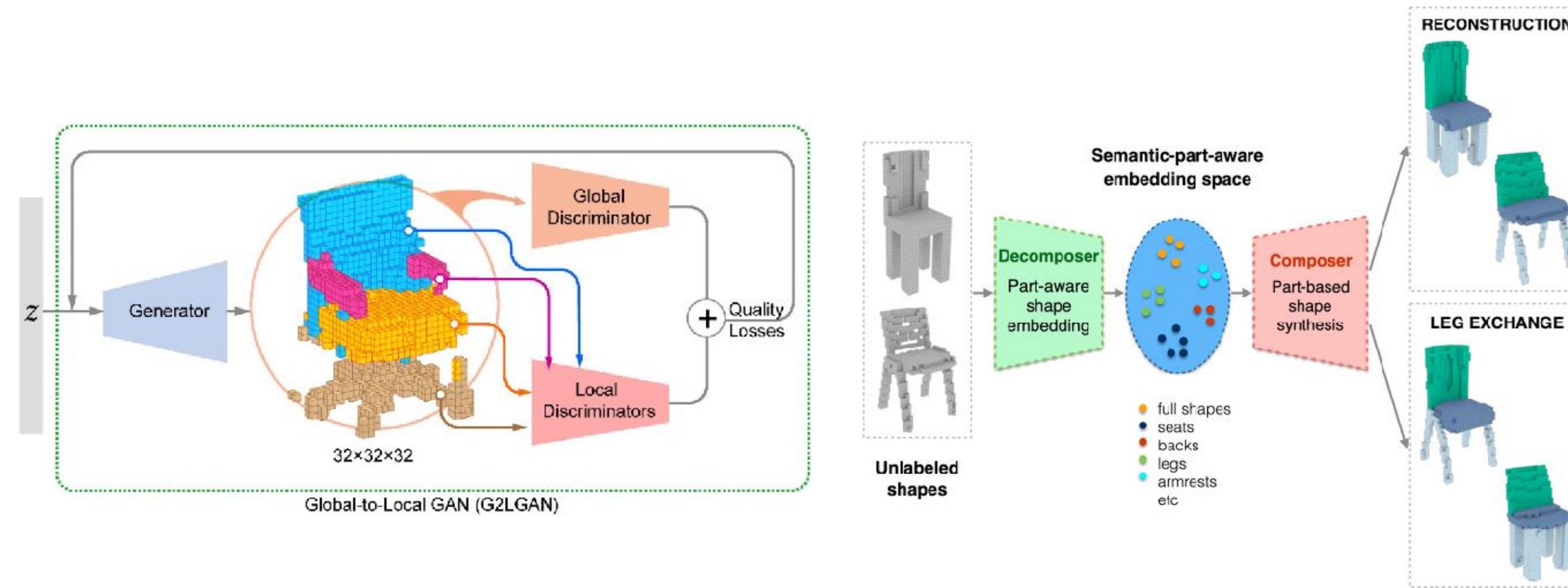
- Detailed Part Geometry
- Clear Shape Structures
- Sharp Part Boundaries
- With Semantic Part Labels



Yang and Mo et al., “**DSG-Net: Learning Disentangled Structure and Geometry for 3D Shape Generation**”, ACM ToG 2021

Semantic-level Synthesis and Assembly

Per-part Geometry Generation and Part Assembly

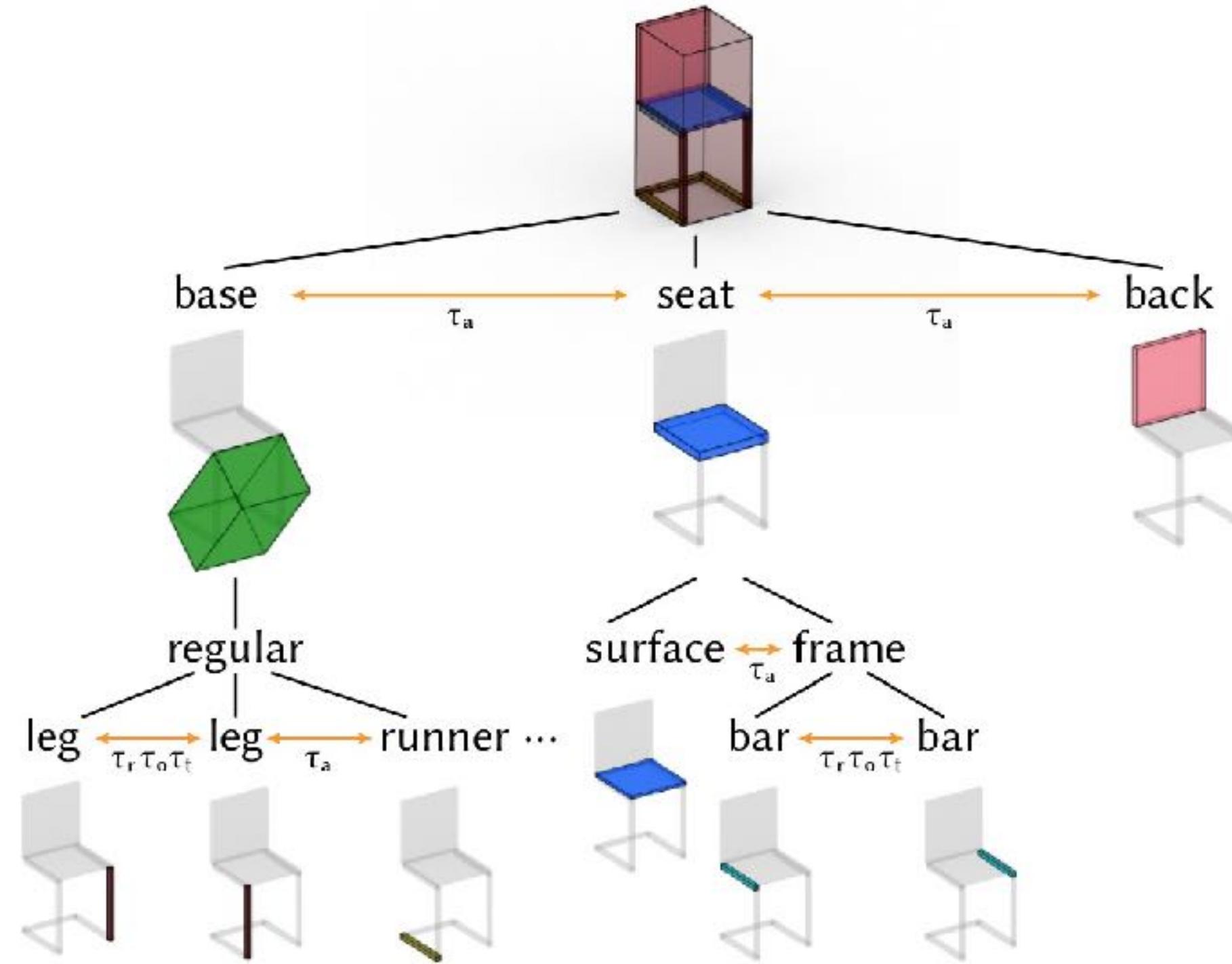


Wang and Schor et al., “**Global-to-Local Generative Model for 3D Shapes**”, Siggraph Asia 2018

Dubrovina et al., “**Composite Shape Modeling via Latent Space Factorization**”, ICCV 2019

Hierarchical Generation

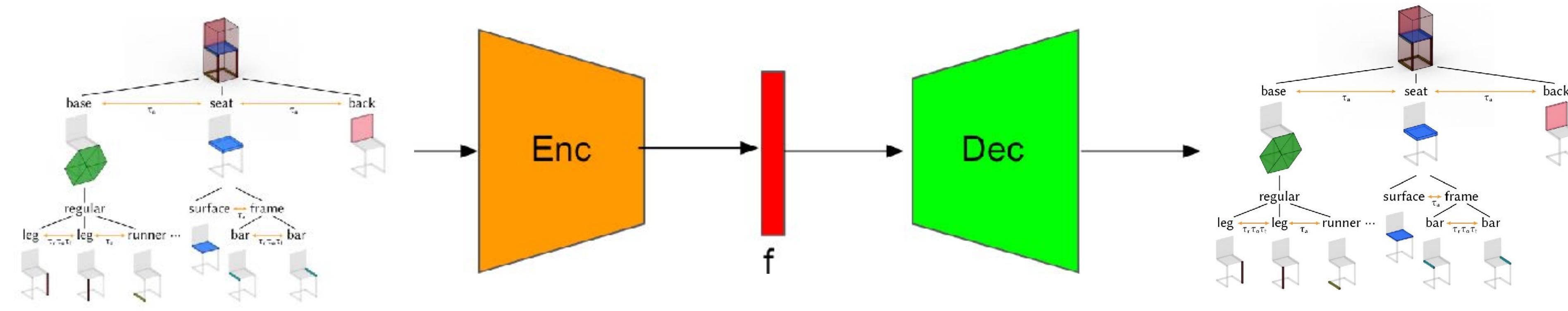
For fine-grained parts, we need coarse-to-fine generation



Mo et al., “**StructureNet: Hierarchical Graph Networks for 3D Shape Generation**”, *Siggraph Asia 2019*

Hierarchical Generation

We use Variational Autoencoder (VAE)



$$f \sim N(0, I)$$

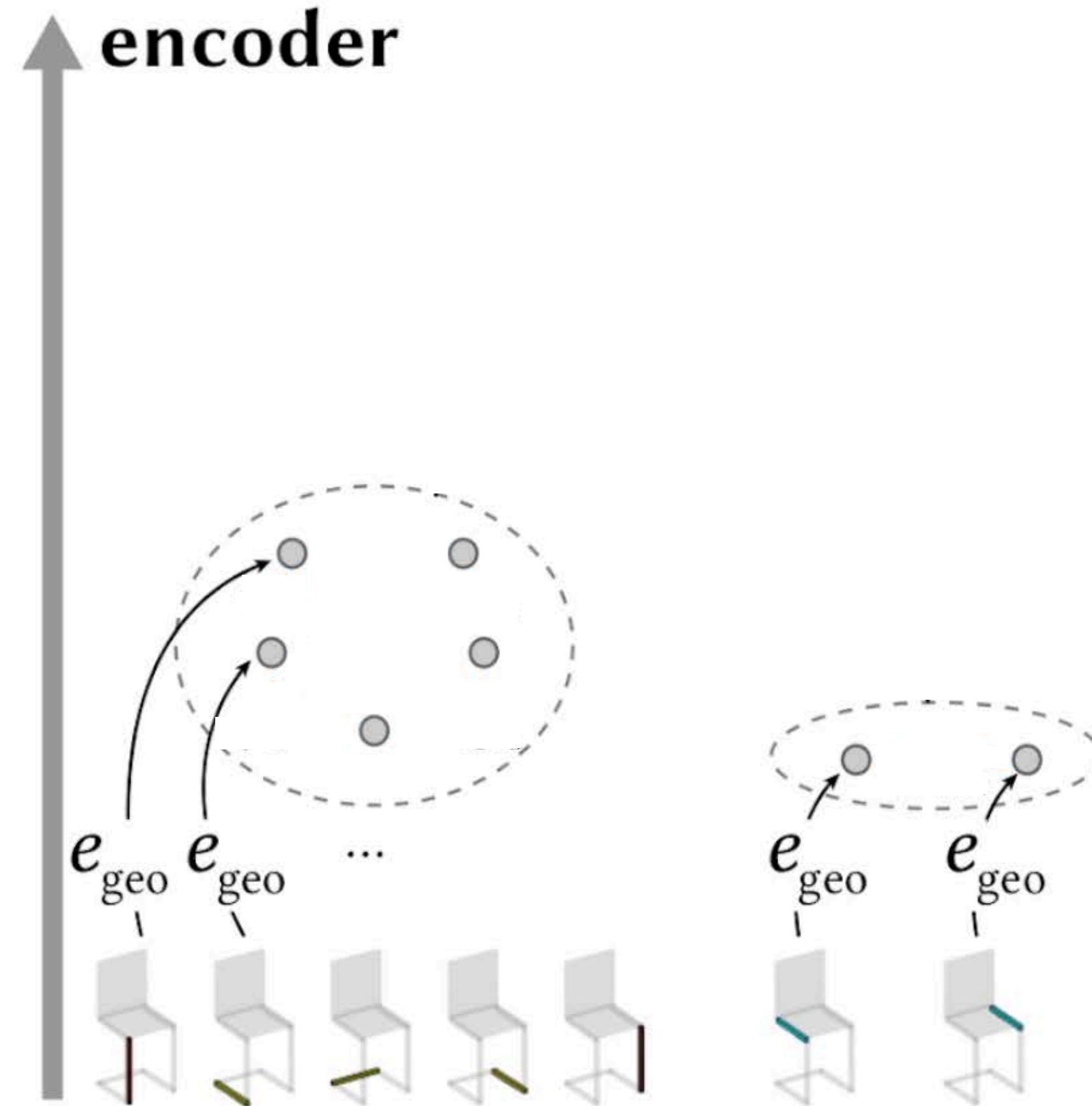
Mo et al., “**StructureNet: Hierarchical Graph Networks for 3D Shape Generation**”, *Siggraph Asia 2019*

Hierarchical Generation



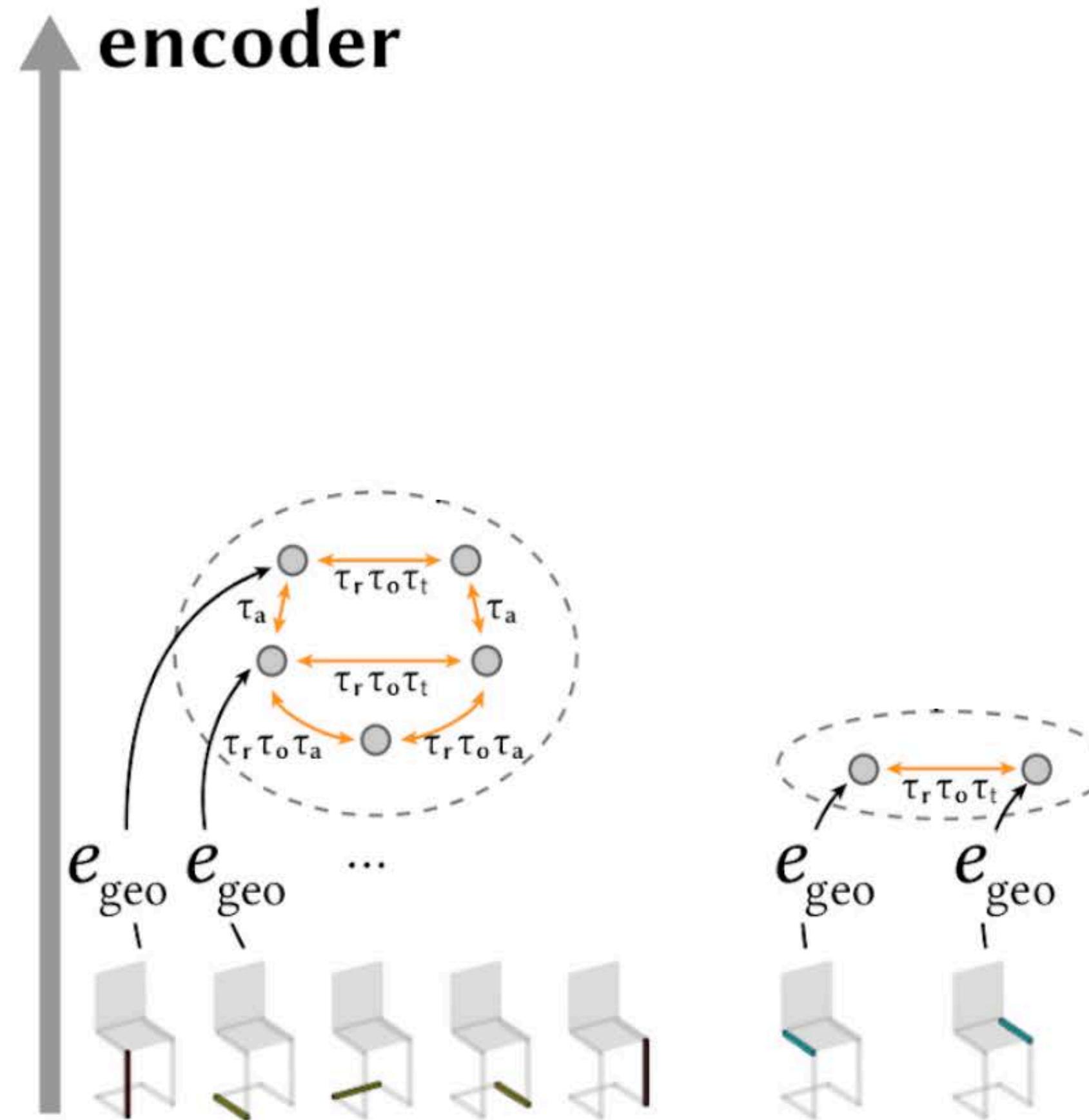
Mo et al., “**StructureNet: Hierarchical Graph Networks for 3D Shape Generation**”, *Siggraph Asia 2019*

Hierarchical Generation



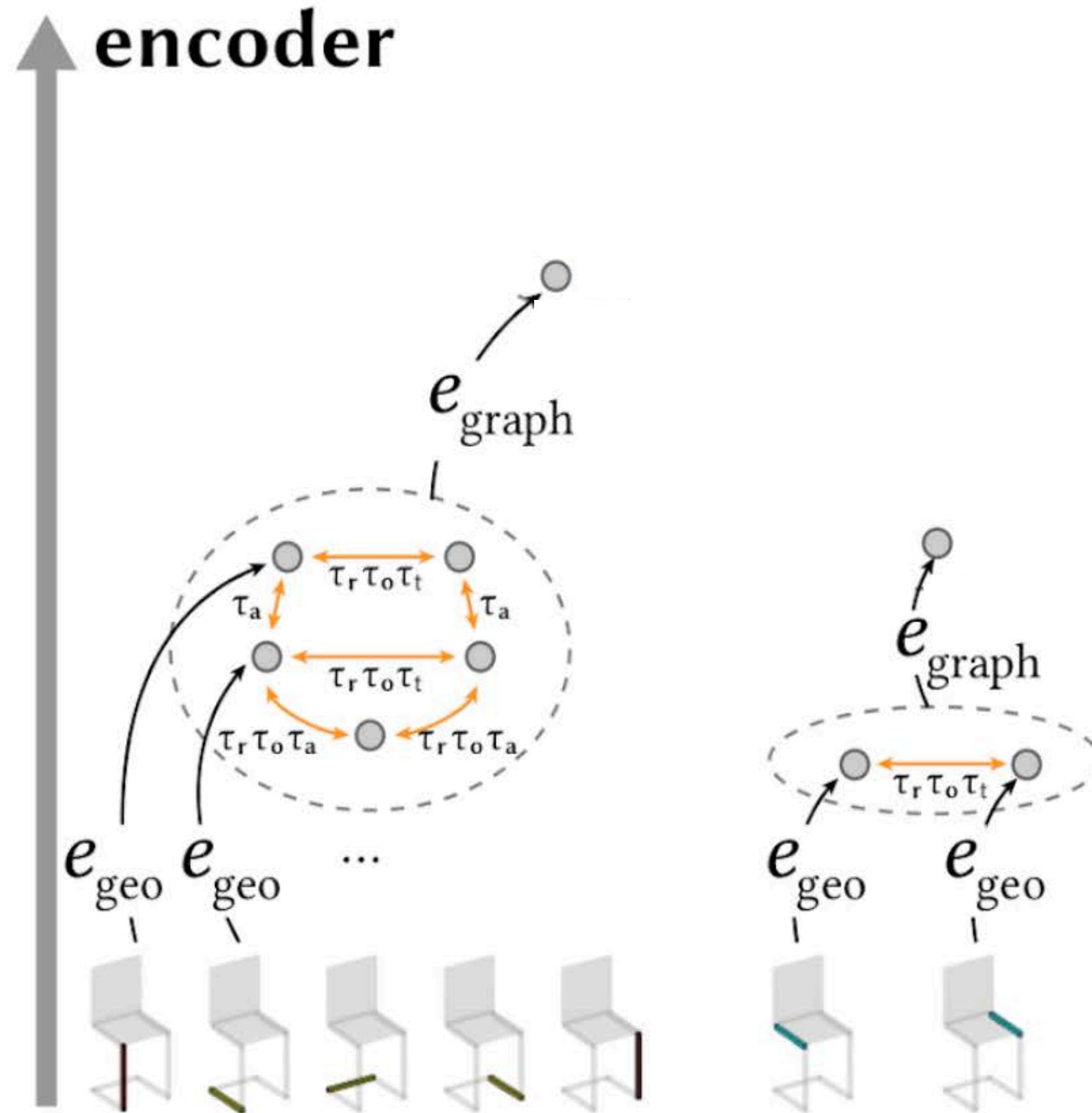
Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, *Siggraph Asia 2019*

Hierarchical Generation



Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

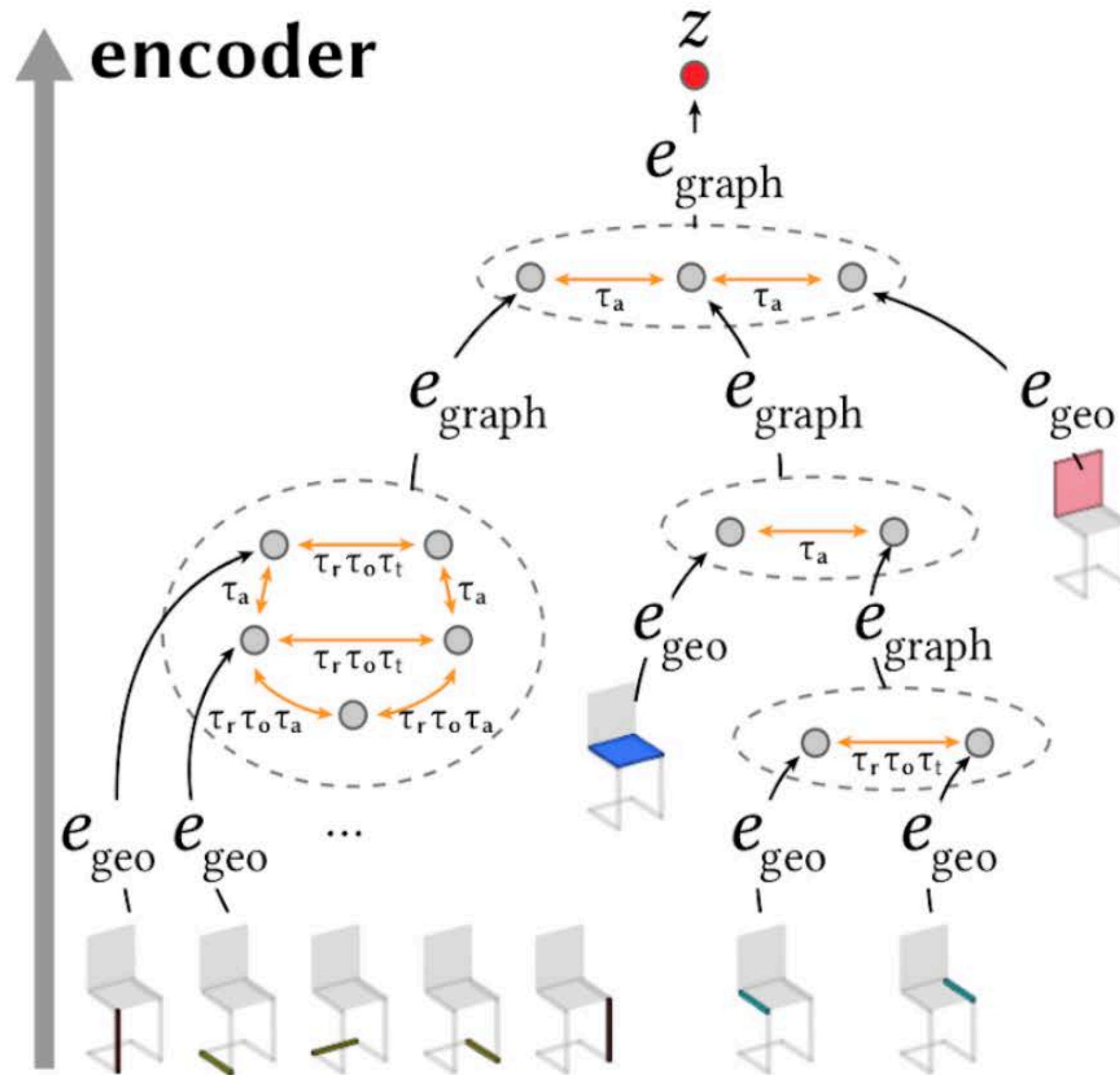
Hierarchical Generation



Graph Neural Network

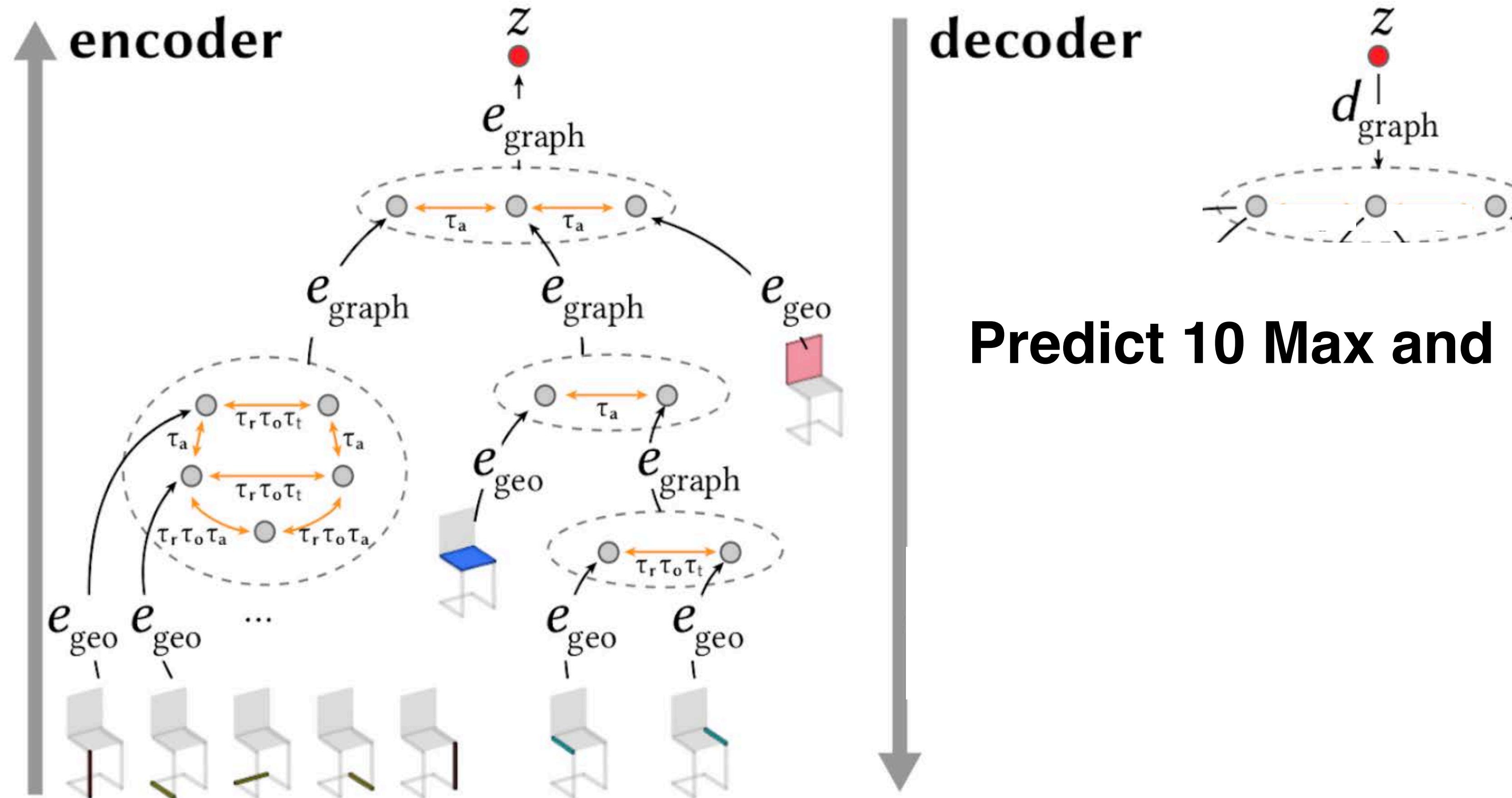
Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, *Siggraph Asia 2019*

Hierarchical Generation



Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

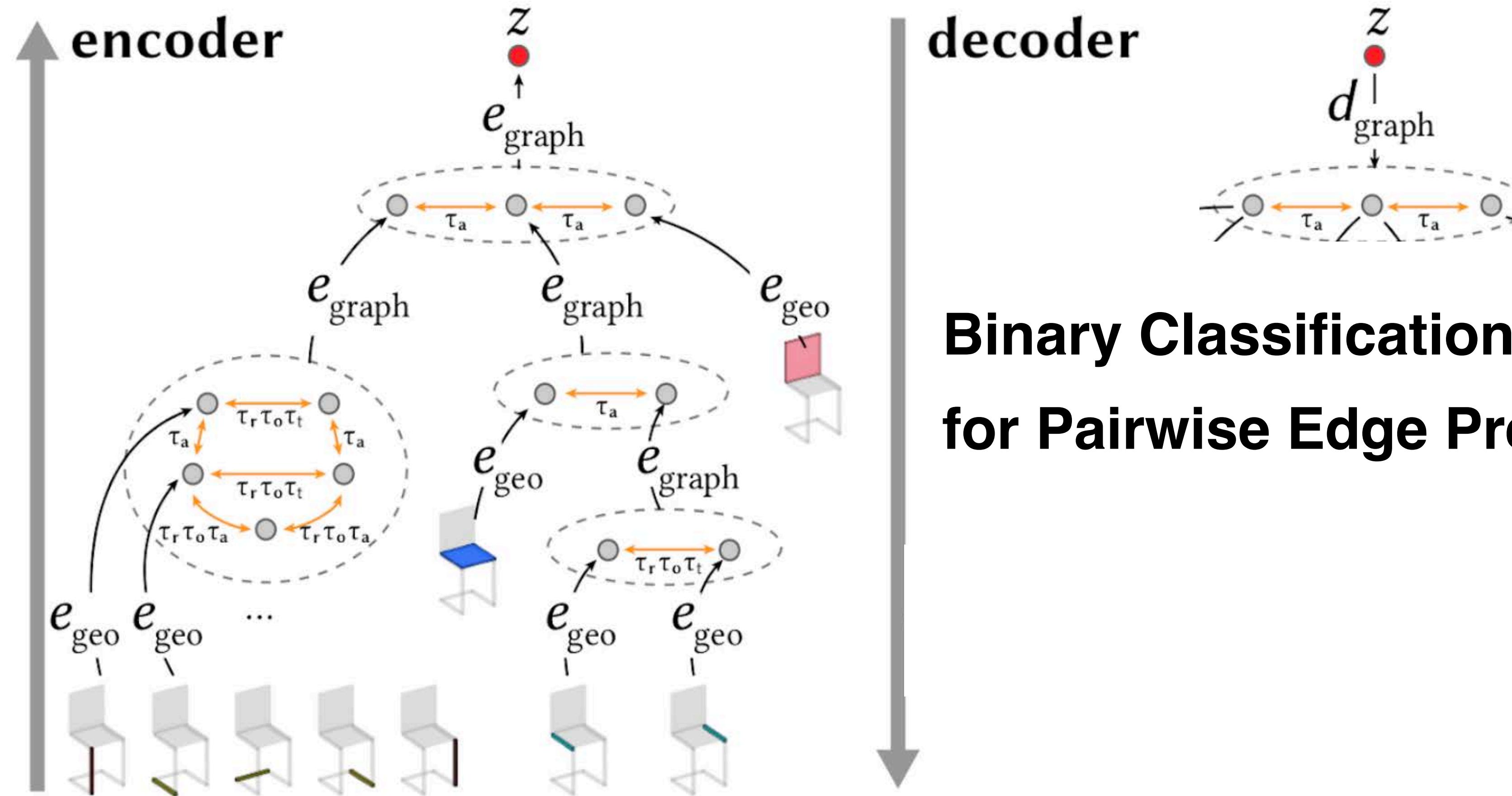
Hierarchical Generation



Predict 10 Max and Match

Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

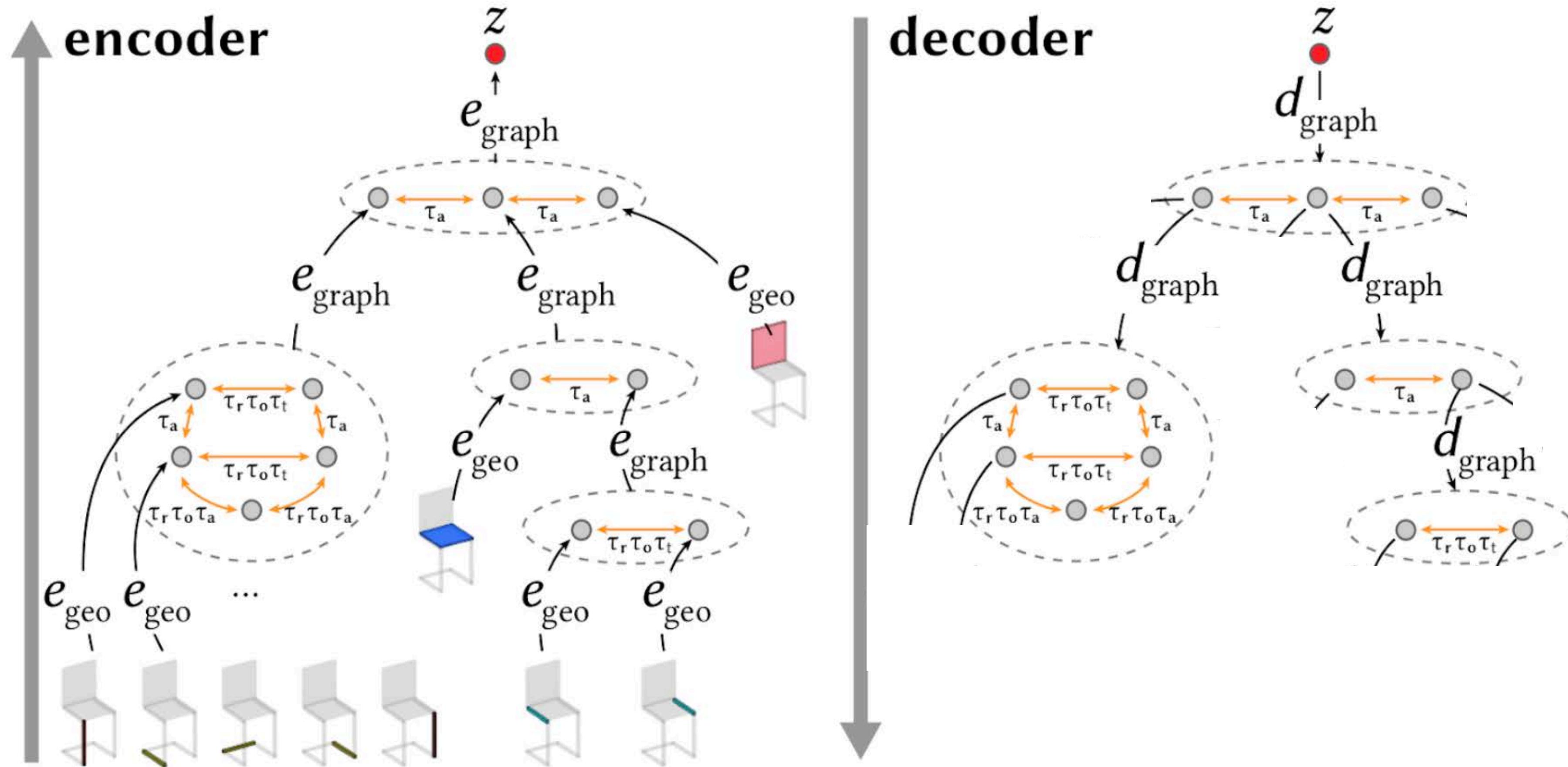
Hierarchical Generation



**Binary Classification
for Pairwise Edge Prediction**

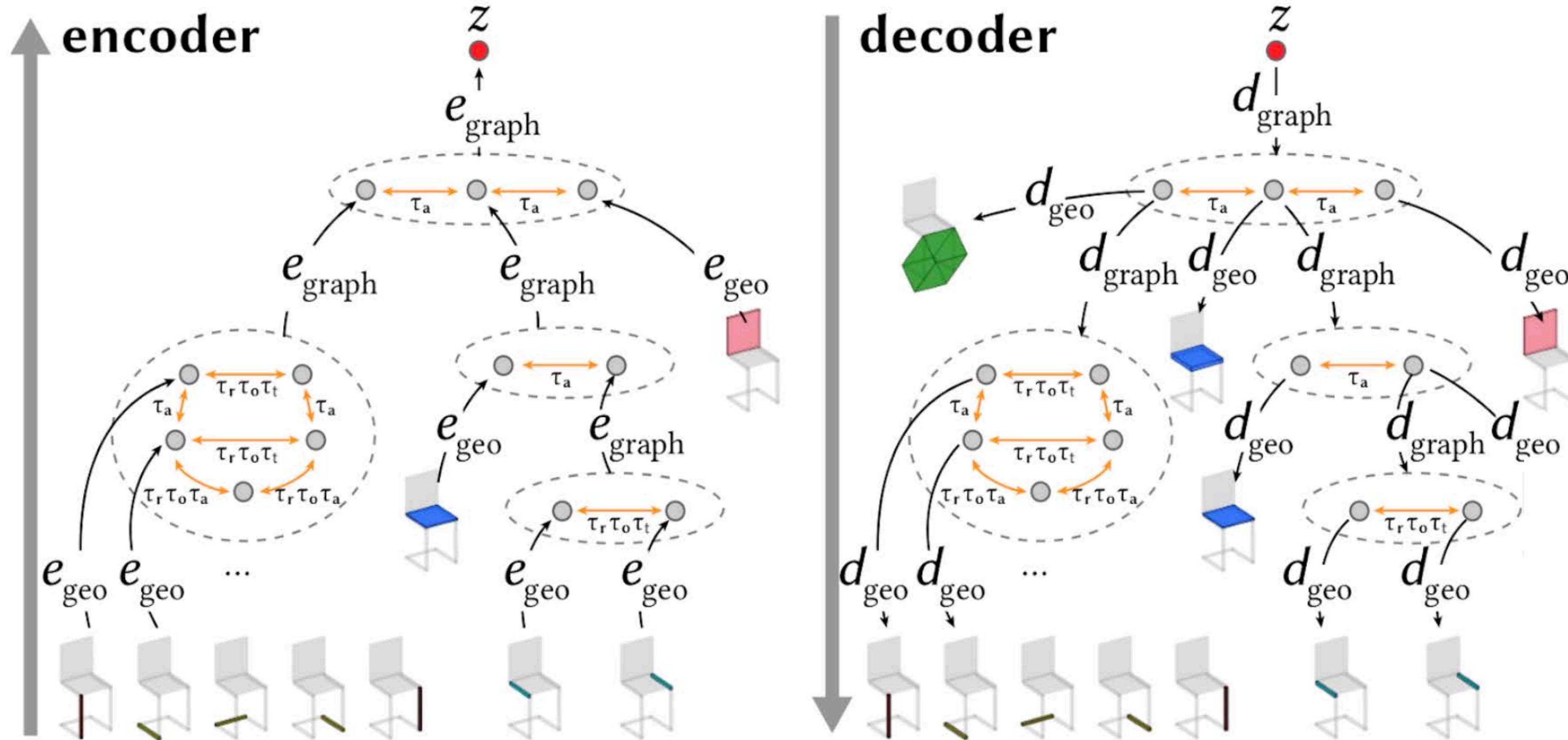
Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

Hierarchical Generation



Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

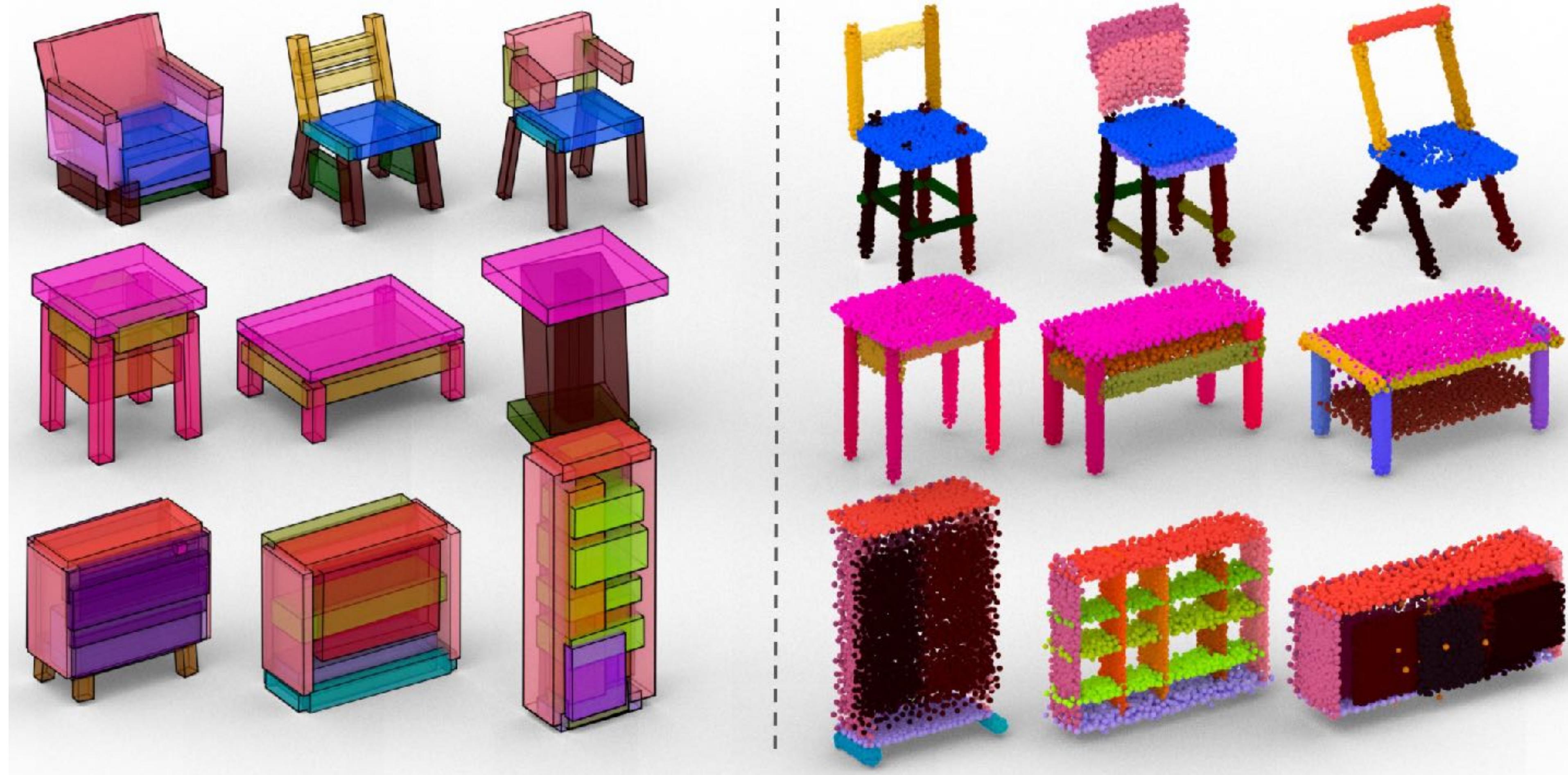
Hierarchical Generation



Recursive Neural Network

Mo et al., “StructureNet: Hierarchical Graph Networks for 3D Shape Generation”, Siggraph Asia 2019

Hierarchical Generation: Results



Mo et al., “**StructureNet: Hierarchical Graph Networks for 3D Shape Generation**”, *Siggraph Asia 2019*

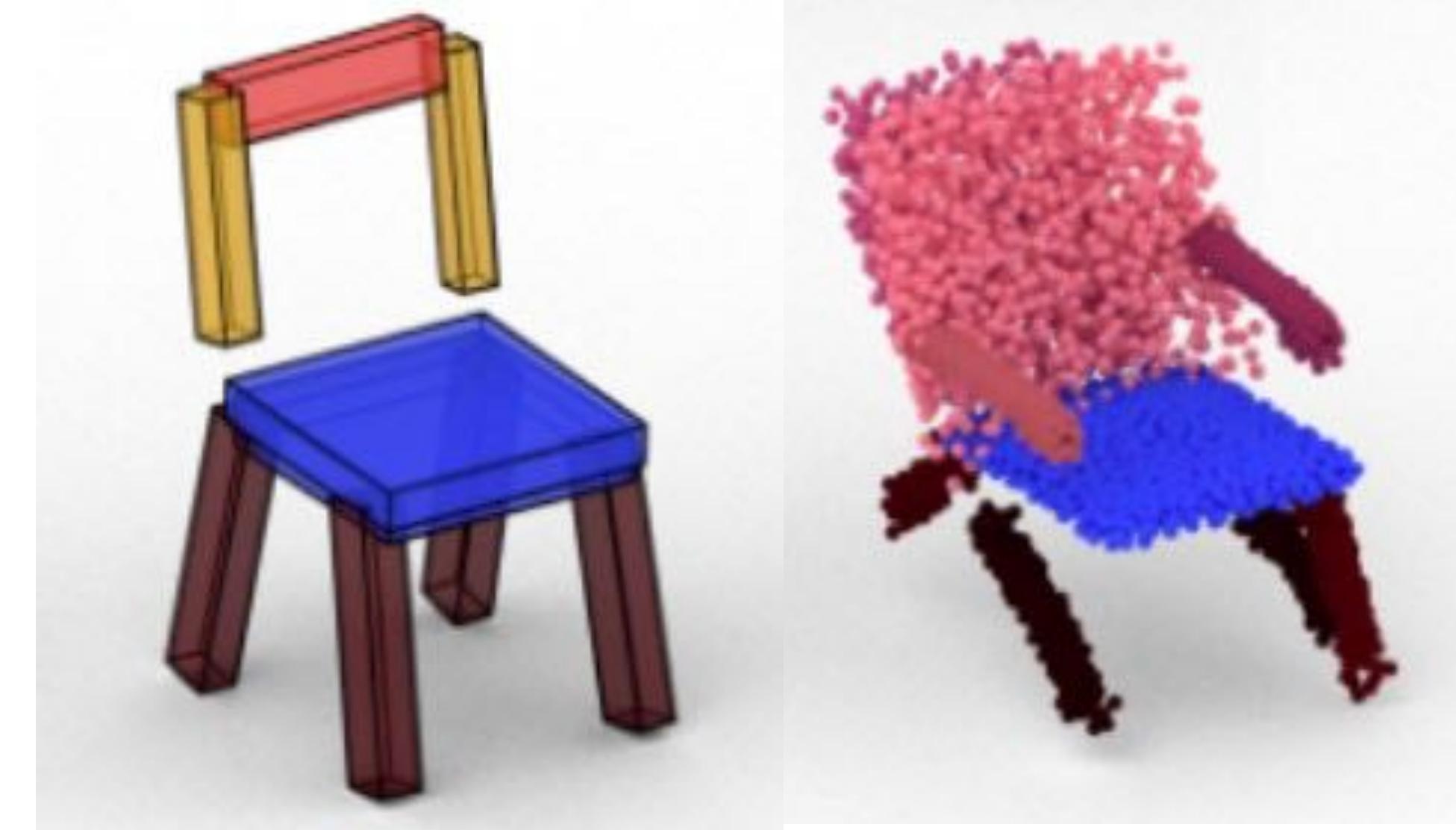
Hierarchical Generation: Results

Use mesh generation for leaf nodes



Part Connection Issue

There could be floating parts



Mo et al., “**StructureNet: Hierarchical Graph Networks for 3D Shape Generation**”, *Siggraph Asia 2019*

Part Connection Issue

Post-processing



*Optimize part poses
using the predicted part
adjacency relationships*

Yang and Mo et al., “DSG-Net: Learning Disentangled Structure and Geometry for 3D Shape Generation”, ACM ToG 2021

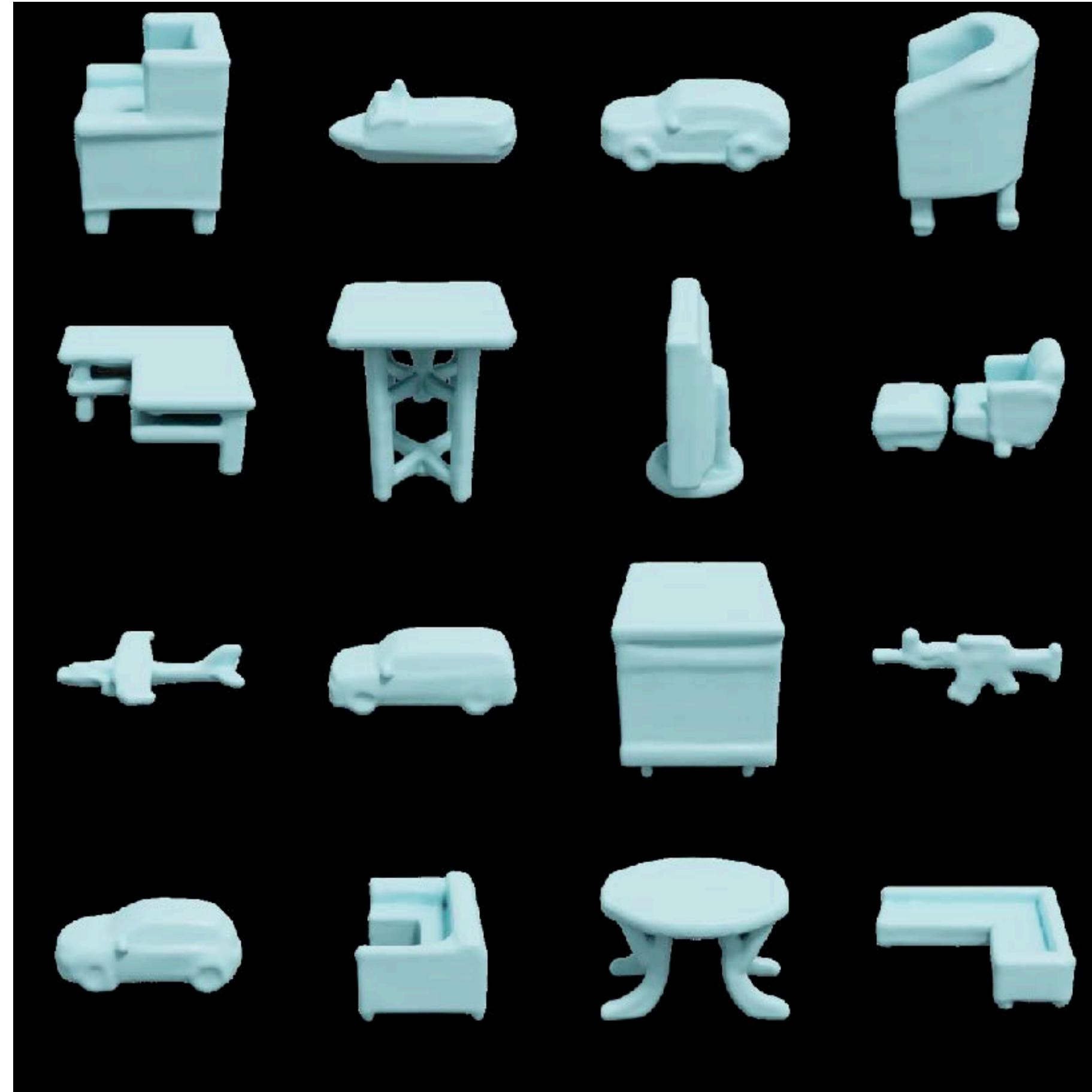
Joint Re-synthesis



*Erode the part connection
regions and re-generate
the local geometry*

Yin et al., “COALESCE: Component Assembly by Learning to Synthesize Connections”, 3DV 2020

Summary



Next Time

