

# Learning to generate 3D shapes

**Subhransu Maji**

College of Information and Computer Sciences

University of Massachusetts, Amherst

<http://people.cs.umass.edu/smaji>

August 10, 2018 @ Caltech



# Creating 3D shapes is not easy

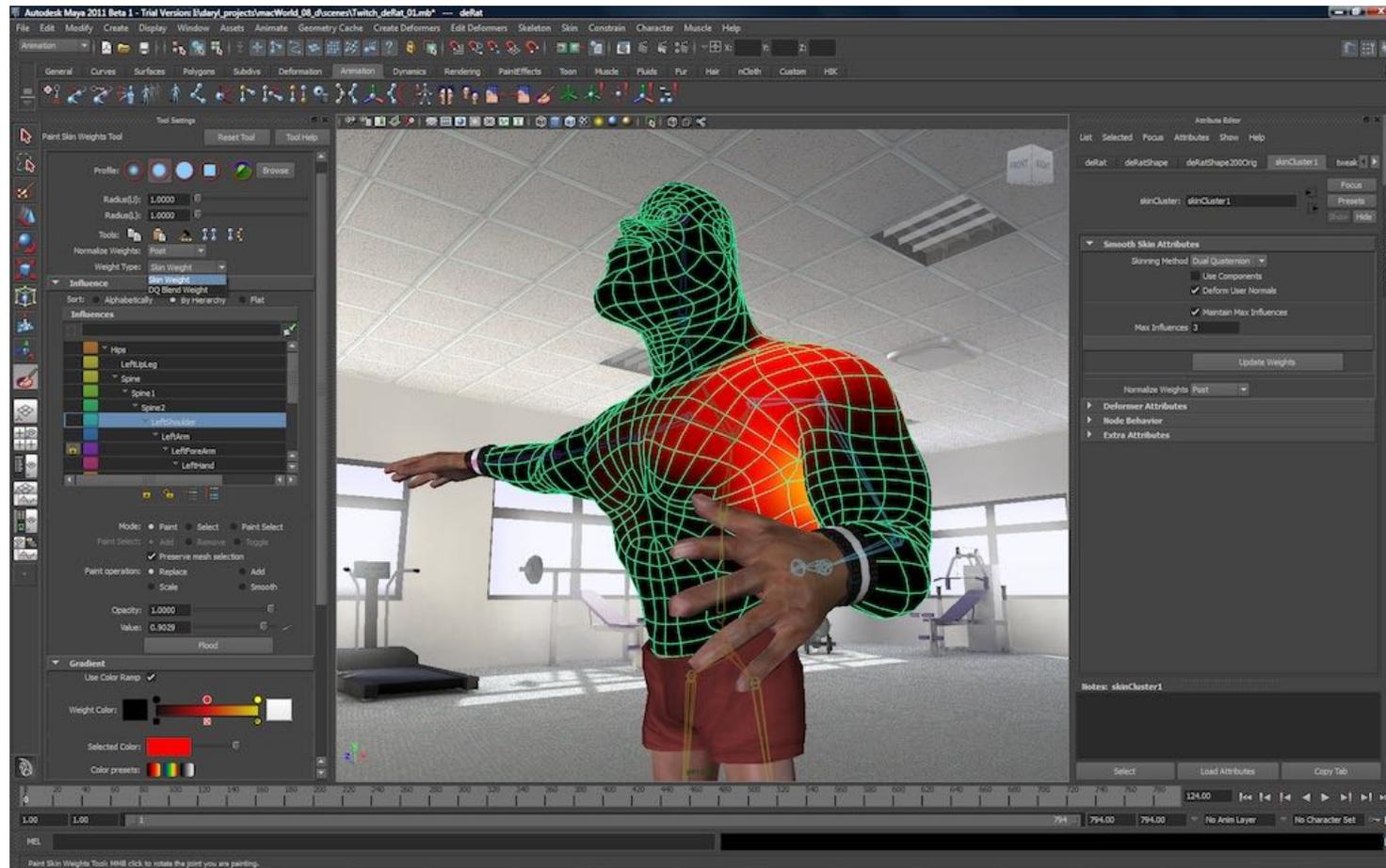
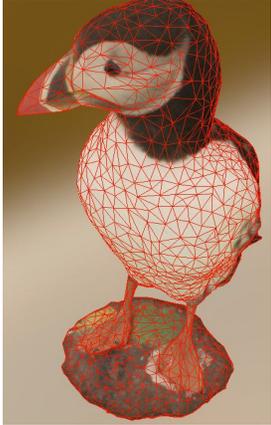
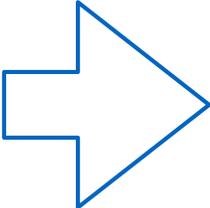


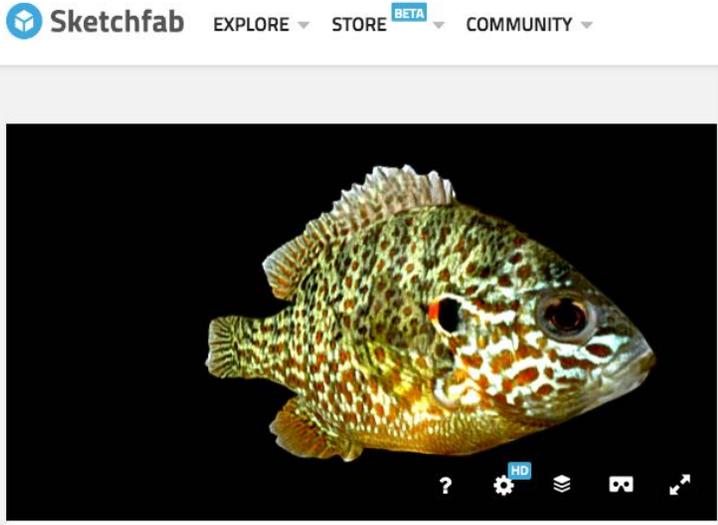
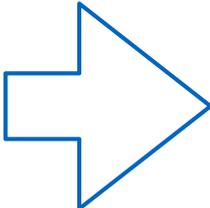
Image from Autodesk 3D Maya

# Inferring 3D shapes from images

What shapes are puffins?



What shapes are pumpkinseed fish?

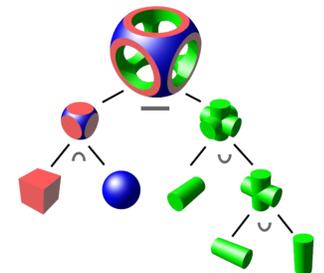
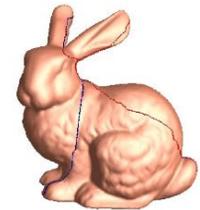
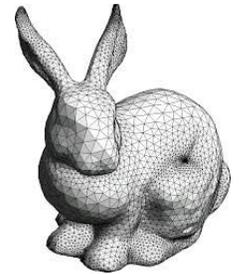


# Creating 3D shapes is not easy

- Many techniques for recognizing 3D data, but relatively few techniques for generating them

- Representations for generation?

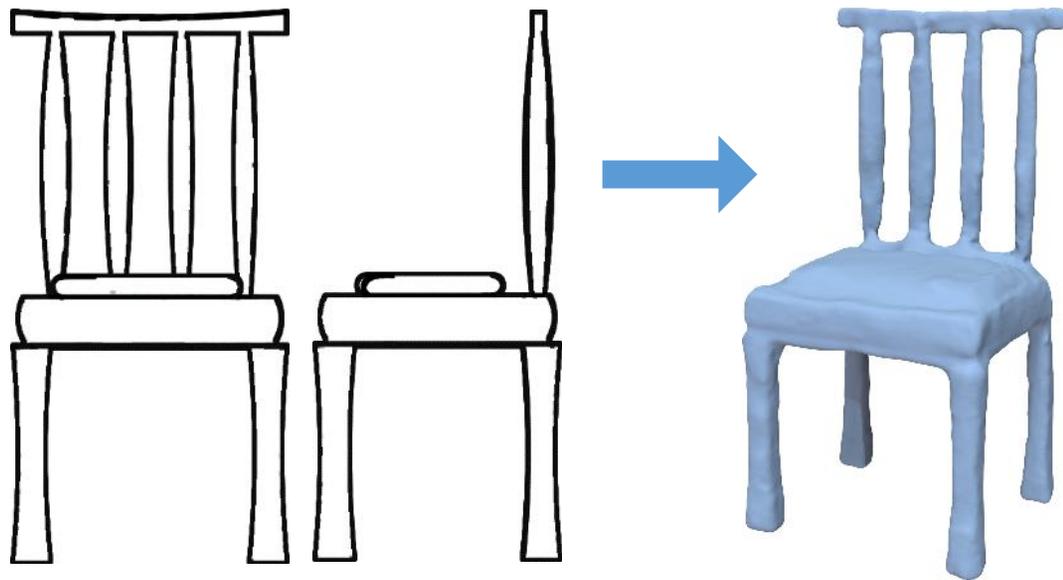
- Voxels
- Multiview
- Geometry images
- Shape basis
- Set-based (points, triangles, etc.)
- Procedural, e.g., constructive solid geometry



# Talk overview

- Generative models for 3D shapes and applications
  - Multiview [3DV'17]
  - Multiresolution tree networks [ECCV'18]
  - Constructive solid geometry [CVPR'18]
- Learning 3D shapes with weak supervision [3DV'17]

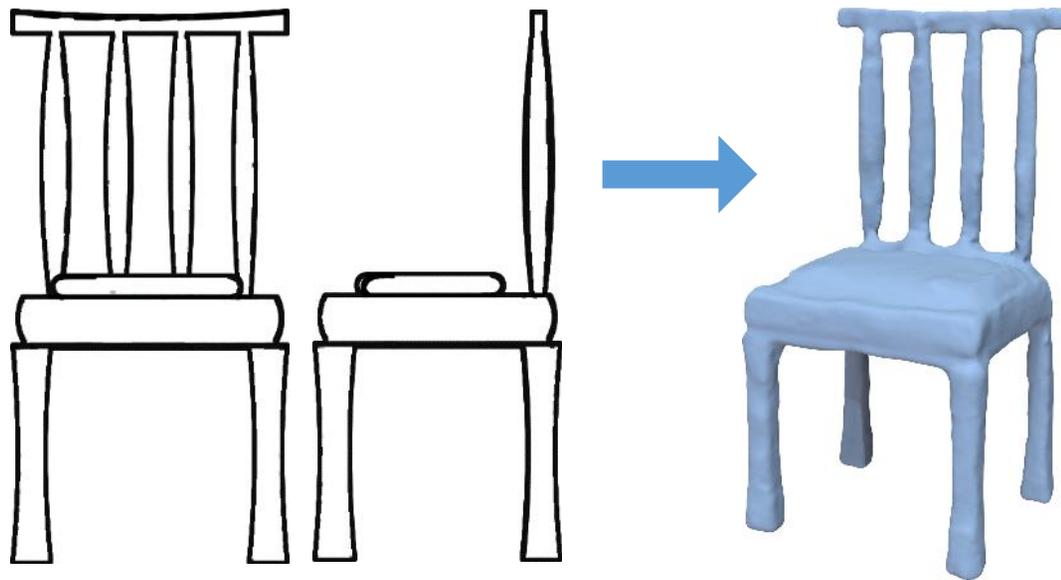
# 3D Shape Reconstruction from Sketches via Multi-view Convolutional Networks



Zhaoliang Lun  
Matheus Gadelha  
Evangelos Kalogerakis  
Subhransu Maji  
Rui Wang

3DV 2017

# 3D Shape Reconstruction from Sketches via Multi-view Convolutional Networks



Zhaoliang Lun  
Matheus Gadelha  
Evangelos Kalogerakis  
Subhransu Maji  
Rui Wang

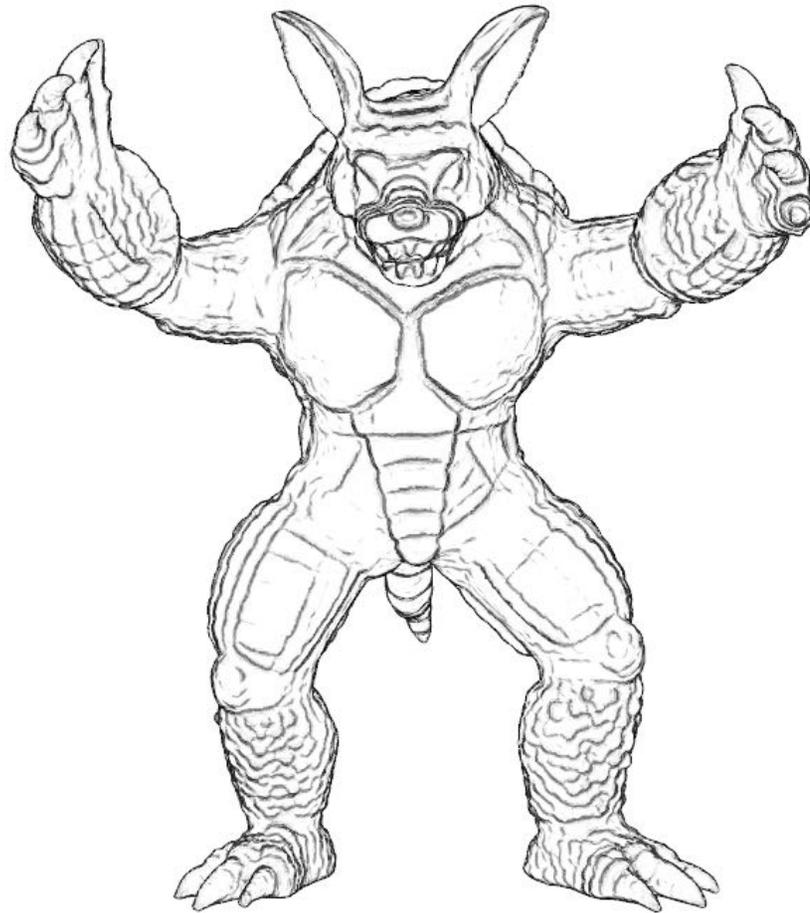
3DV 2017

**What Does the Occluding Contour Tell Us about Solid Shape?**

[Jan J Koenderink](#)

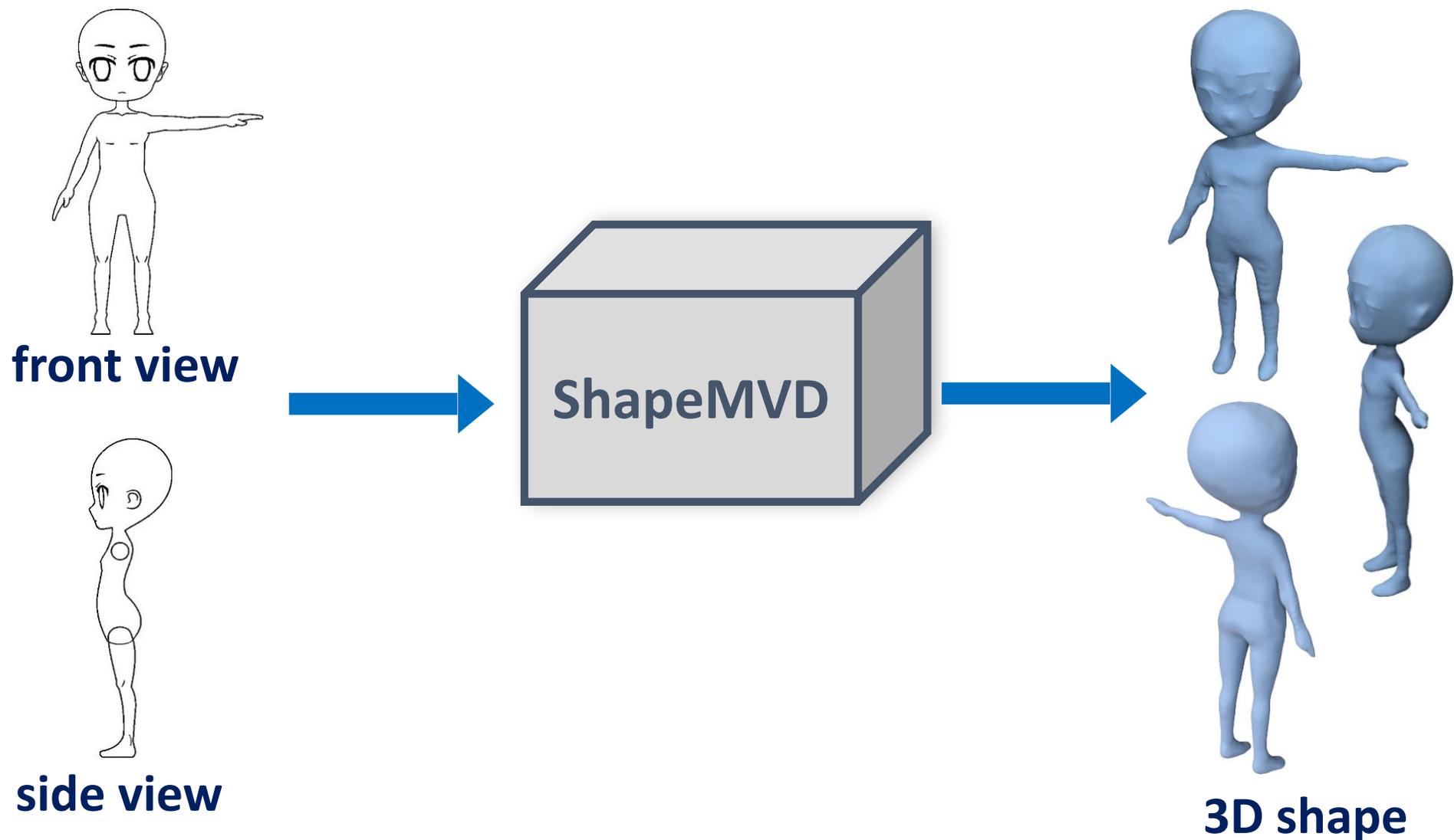
First Published June 1, 1984 | Research Article

Why line drawings? Simple & intuitive medium to convey shape!



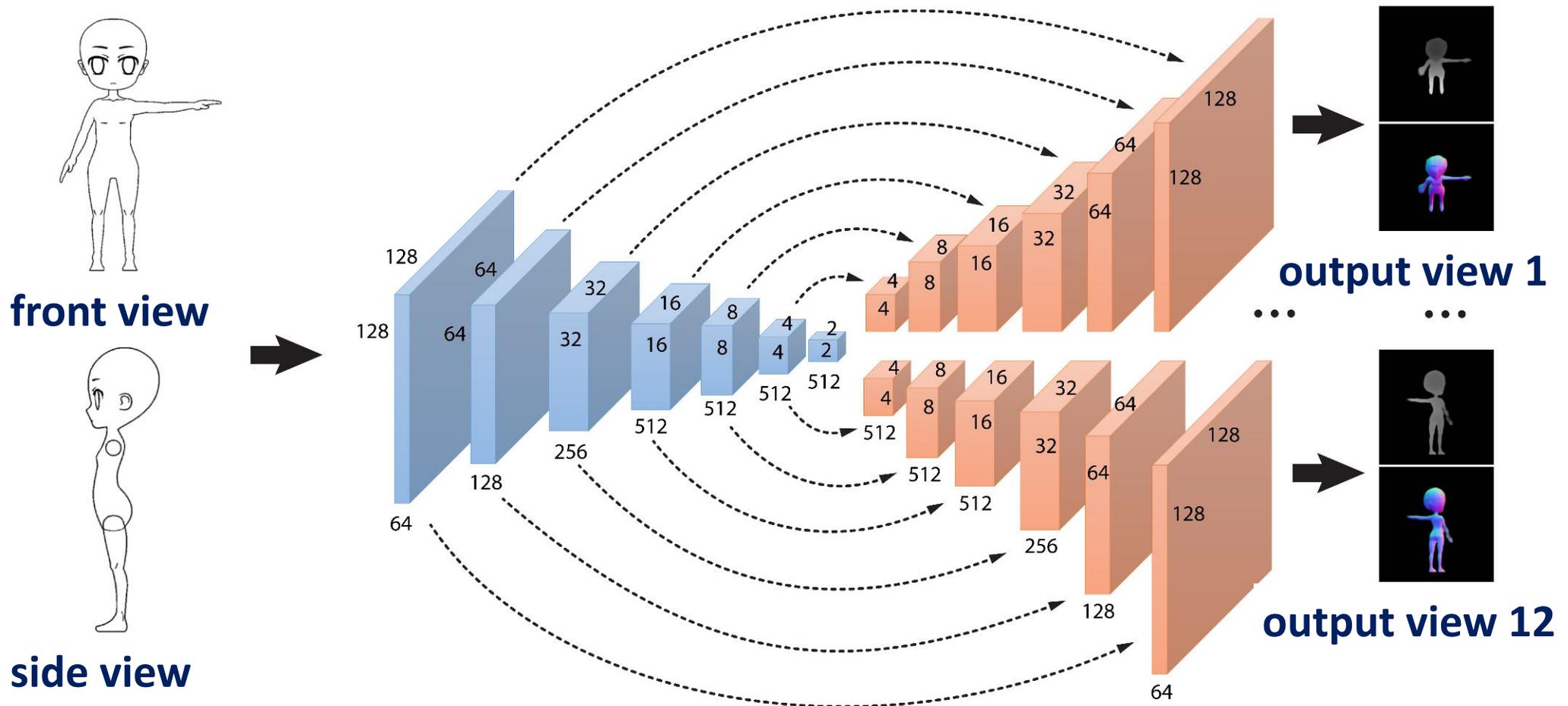
**Image from Suggestive Contour  
Gallery, DeCarlo et al. 2003**

Goal: 2D line drawings in, 3D shapes out!



# Deep net architecture: U-net structure

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



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

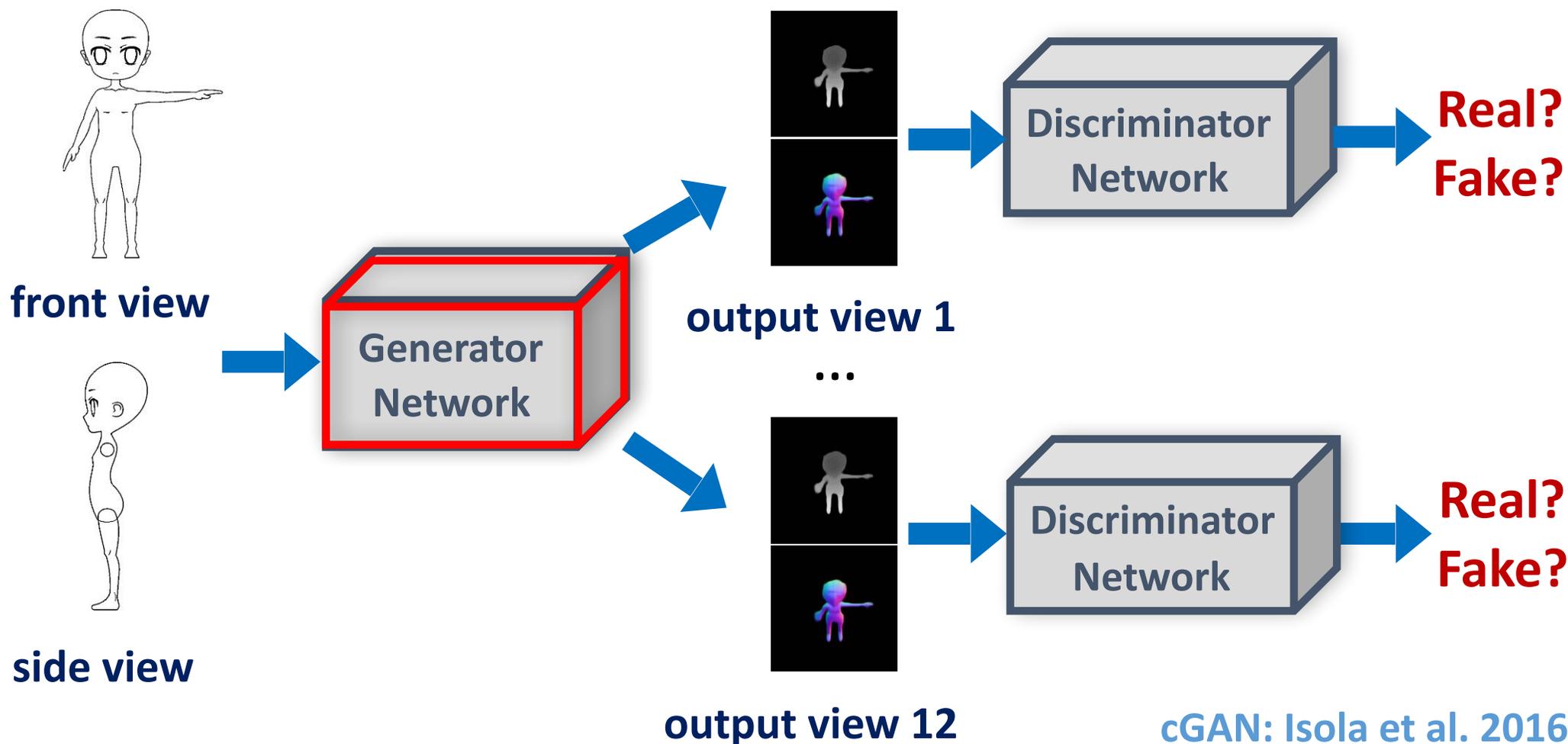
# Training: full loss

Penalize **per-pixel depth reconstruction error:**

& **per-pixel normal reconstruction error:**

& **“unreal” outputs:**

$$\sum_{pixels} |d_{pred} - d_{gt}|$$
$$\sum_{pixels} (1 - n_{pred} \cdot n_{gt})$$
$$-\log P(real)$$



# Training data



**Character**  
**10K models**



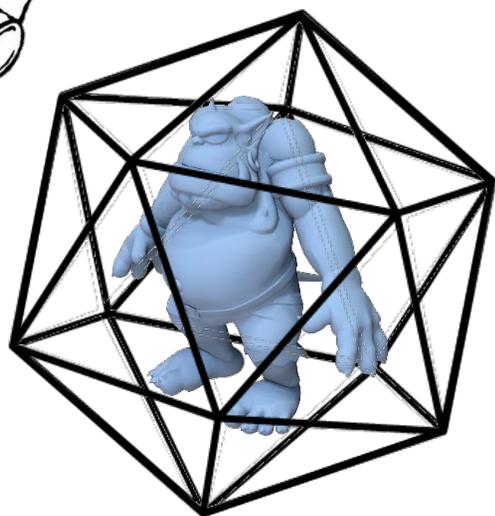
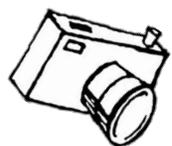
**Chair**  
**10K models**



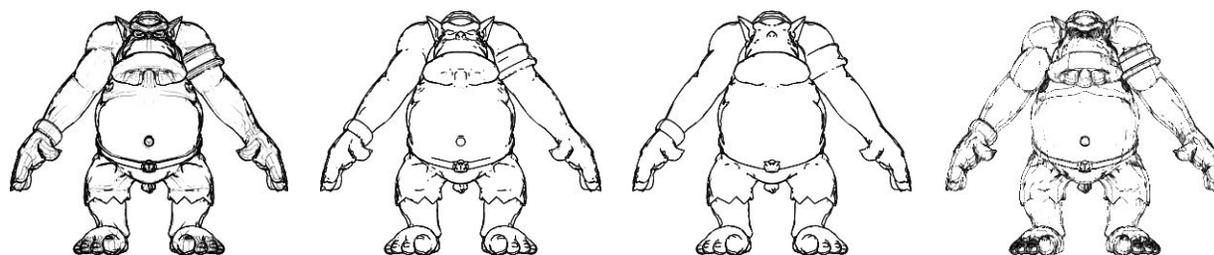
**Airplane**  
**3K models**

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

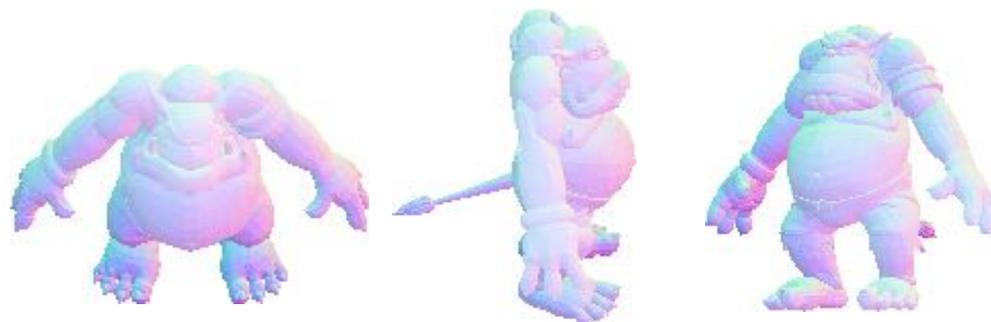
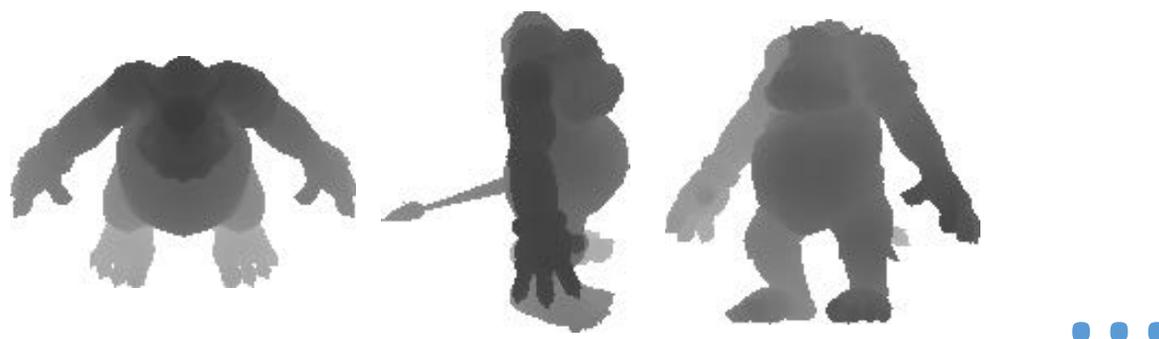
# Training data



12 views



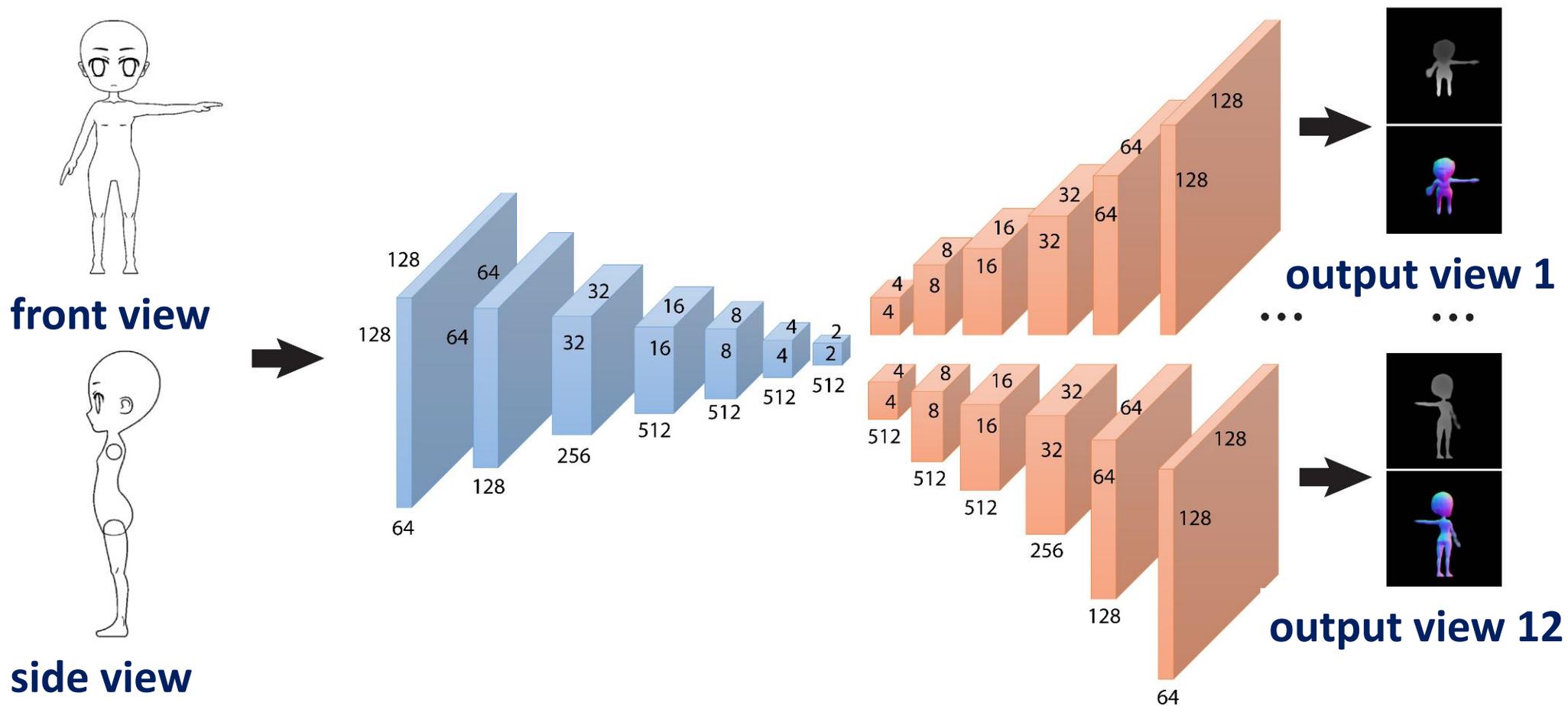
Synthetic line drawings



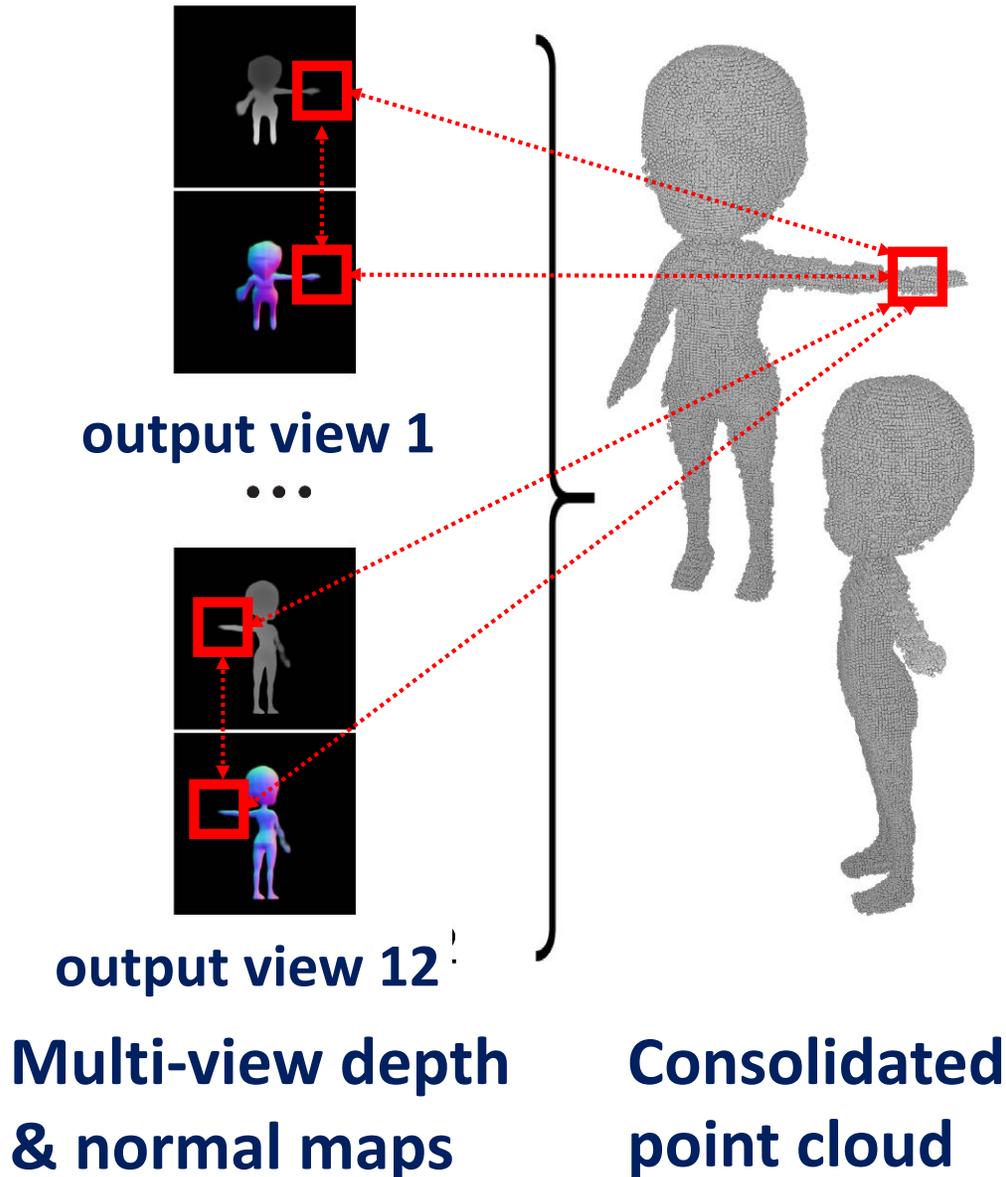
Training depth and normal maps

# Test time

Predict multi-view depth and normal maps!



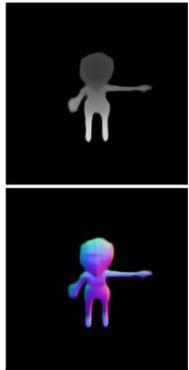
# Multi-view depth & normal map fusion



Optimization problem

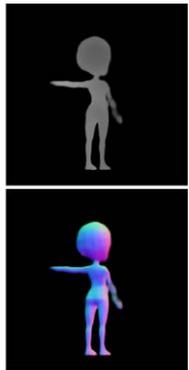
- Depth derivatives should be consistent with normals
- Corresponding depths and normals across different views should agree

# Surface reconstruction



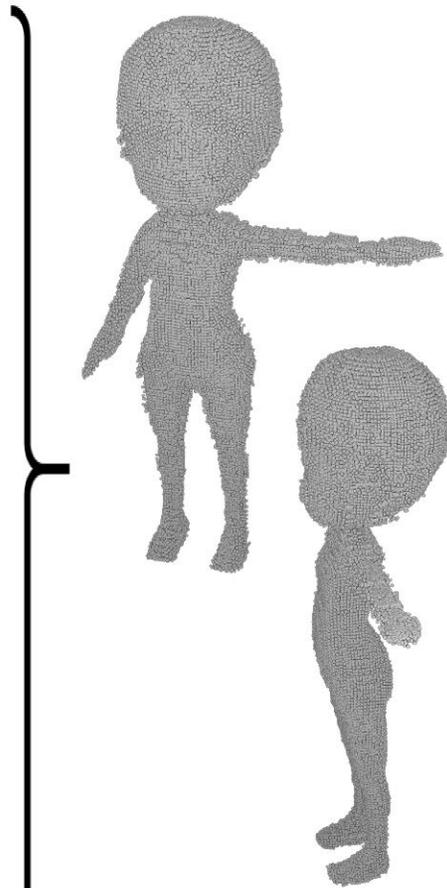
output view 1

...

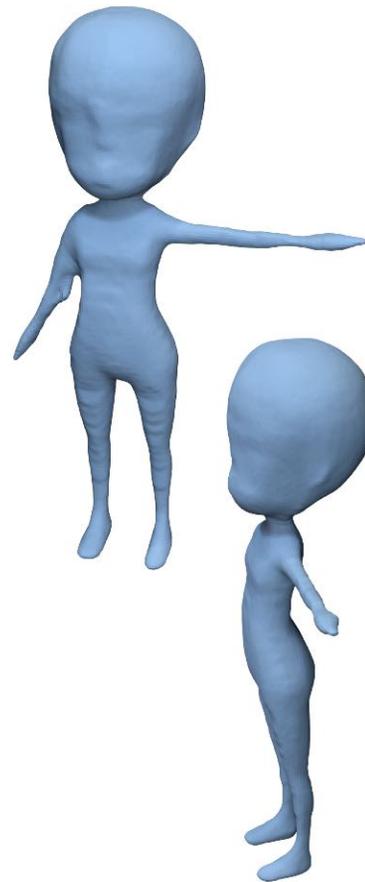


output view 12

**Multi-view depth  
& normal maps**

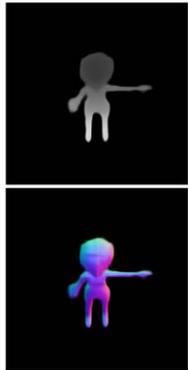


**Consolidated  
point cloud**



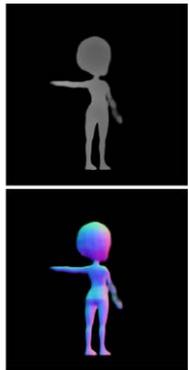
**Surface  
reconstruction**  
[Kazhdan et al. 2013]

# Surface deformation



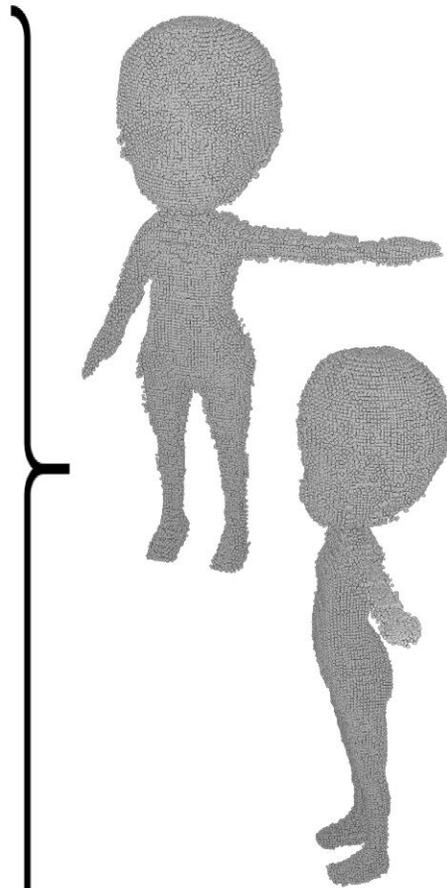
output view 1

...

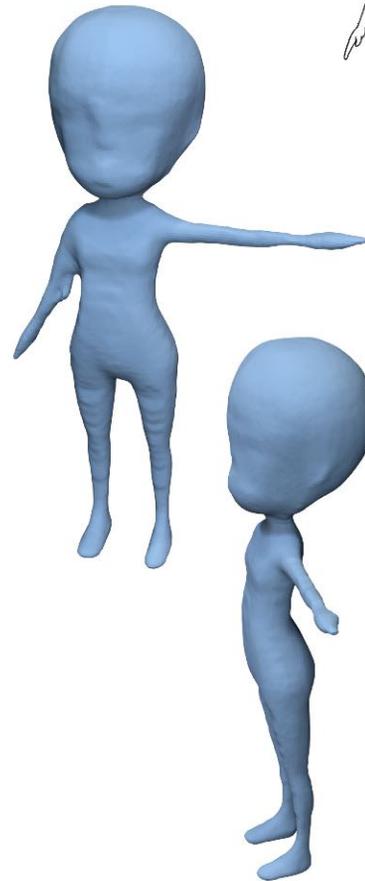


output view 12

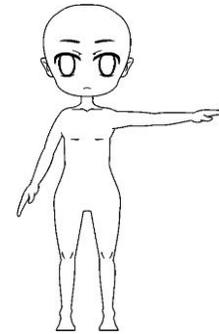
**Multi-view depth  
& normal maps**



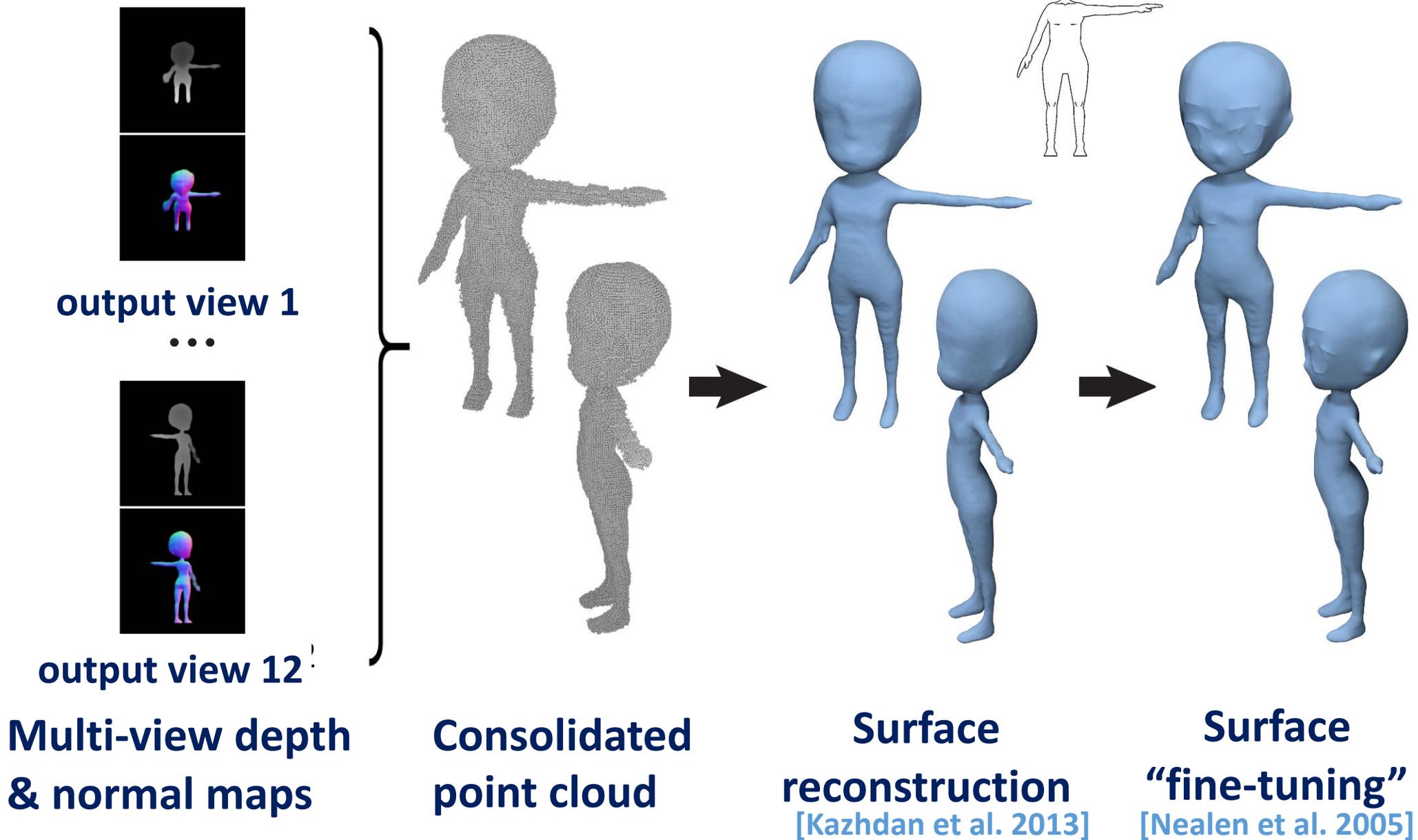
**Consolidated  
point cloud**



**Surface  
reconstruction**  
[Kazhdan et al. 2013]



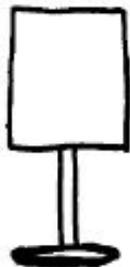
# Surface deformation



# Experiments

# Qualitative Results

**reference  
shape**

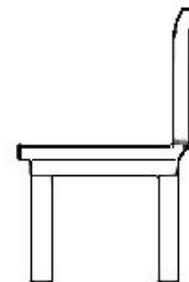
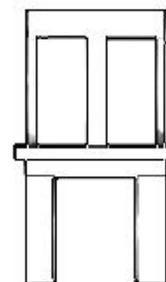


**nearest  
retrieval**

**our  
result**

**volumetric  
net**

**reference  
shape**



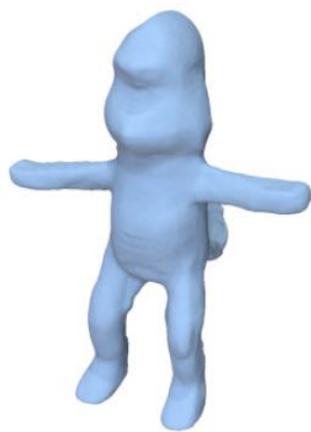
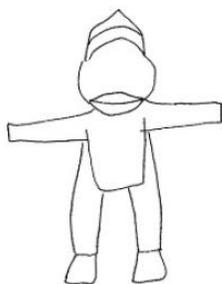
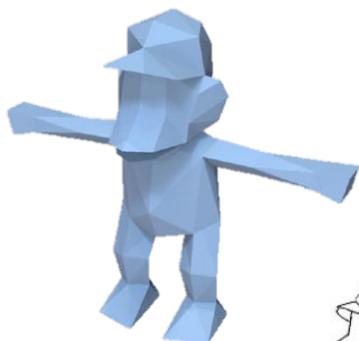
**nearest  
retrieval**

**our  
result**

**volumetric  
net**

# Qualitative Results

reference  
shape

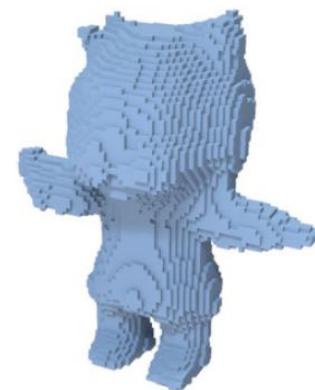


nearest  
retrieval

our  
result

volumetric  
net

reference  
shape



nearest  
retrieval

our  
result

volumetric  
net

# Quantitative Results

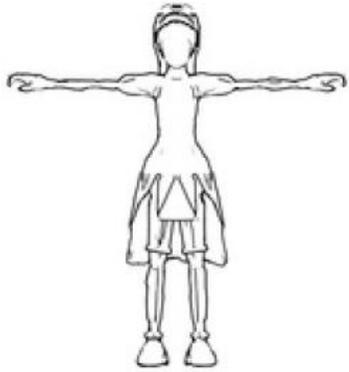
## Character (human drawing)

Metric	Our method	Volumetric	NN
Hausdorff distance	<b>0.120</b>	0.638	0.242
Chamfer distance	<b>0.023</b>	0.052	0.045
normal distance	<b>34.27</b>	56.97	47.94
volumetric distance	<b>0.309</b>	0.497	0.550

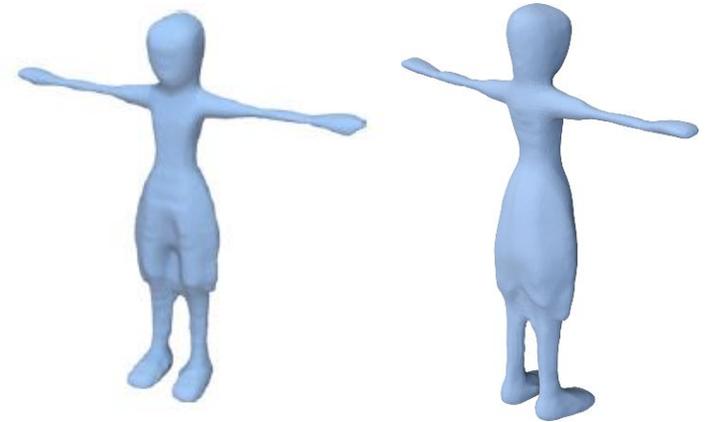
## Man-made (human drawing)

Hausdorff distance	<b>0.171</b>	0.211	0.228
Chamfer distance	<b>0.028</b>	0.032	0.038
normal distance	<b>34.19</b>	48.81	43.75
volumetric distance	<b>0.439</b>	0.530	0.560

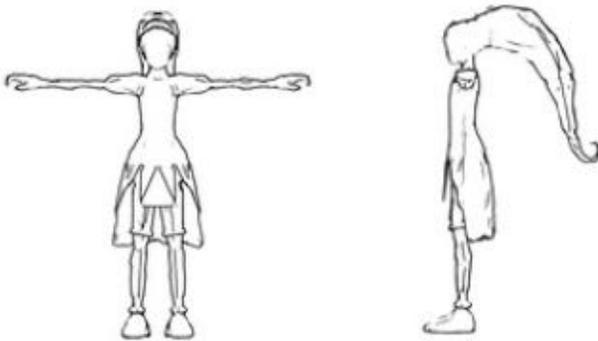
# Single vs two input line drawings



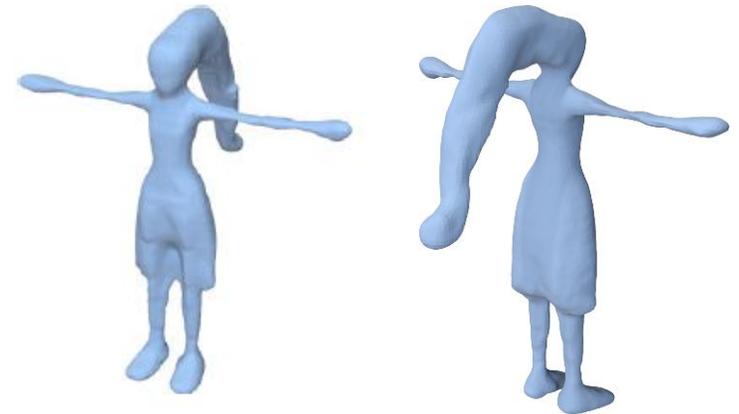
**Single sketch**



**Resulting shape**

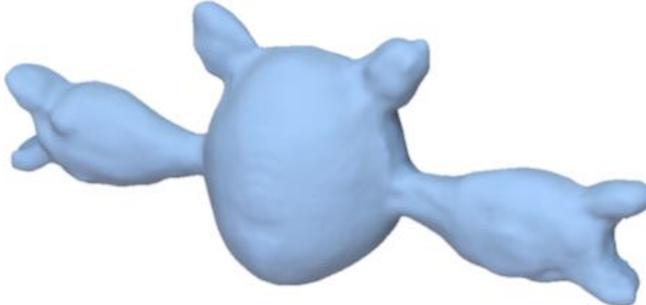
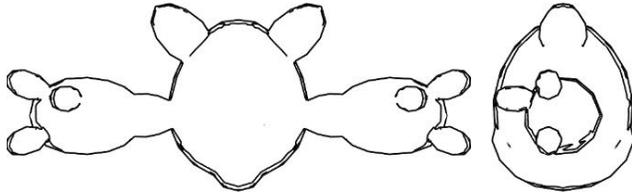
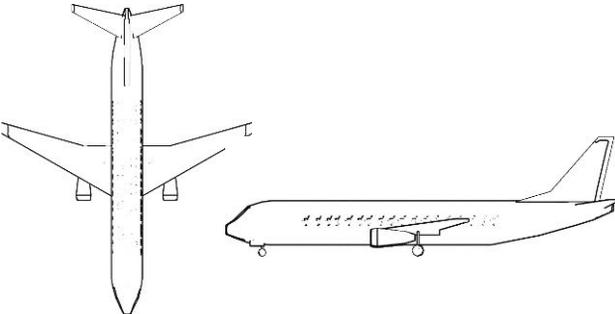
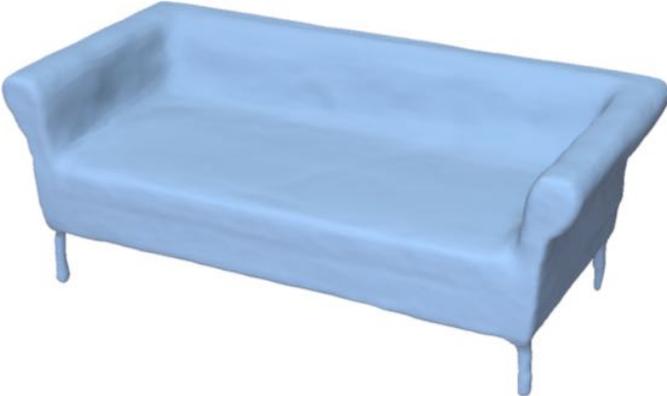
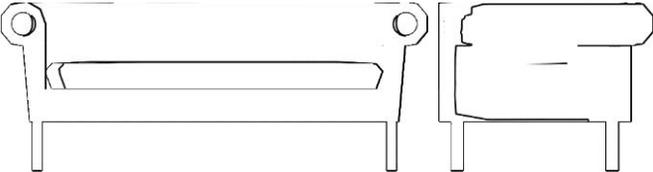


**Two sketches**

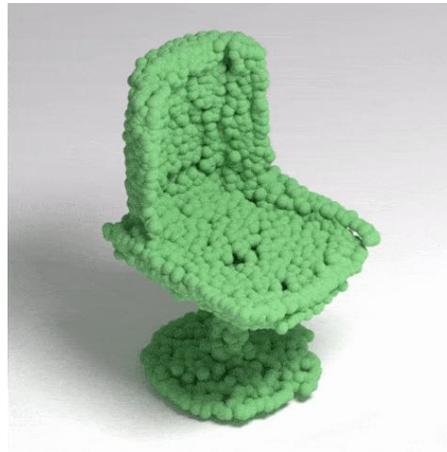
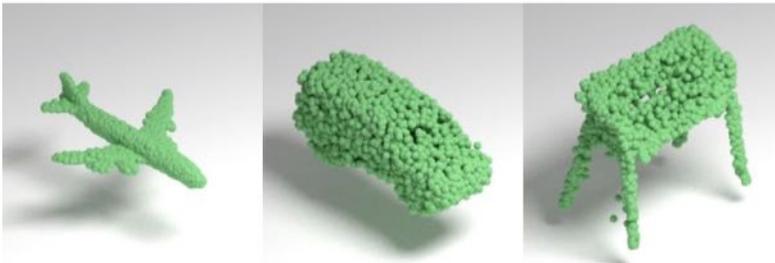


**Resulting shape**

# More results



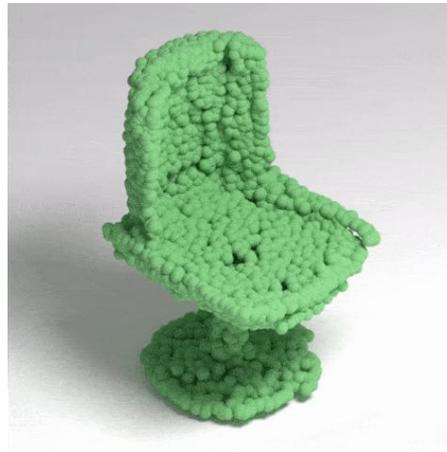
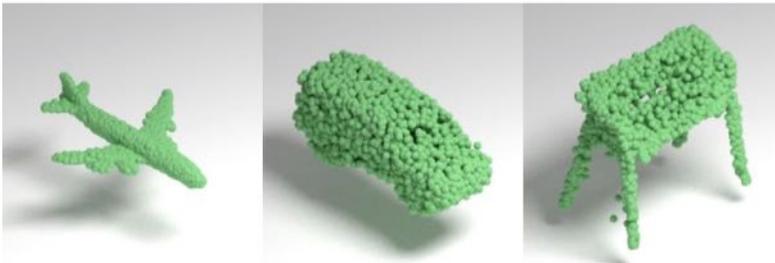
# Multiresolution Tree Networks for 3D Point Cloud Processing



Matheus Gadelha  
Subhransu Maji  
Rui Wang

ECCV 18

# Multiresolution Tree Networks for 3D Point Cloud Processing

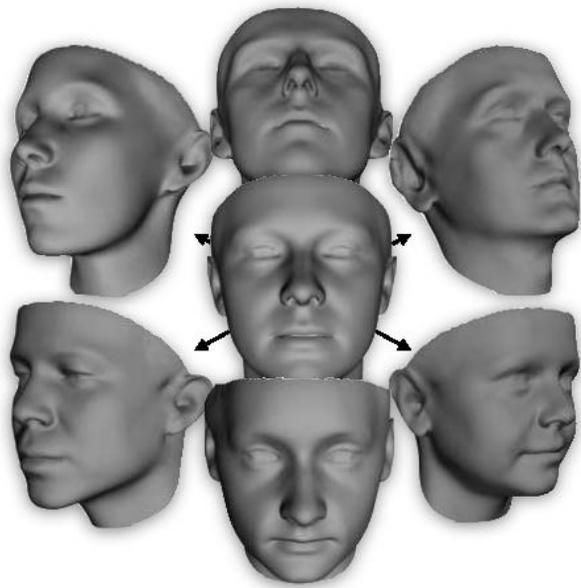


Matheus Gadelha  
Subhransu Maji  
Rui Wang

ECCV 18

# Point-cloud decoders

- Global shape basis or fully-connected decoders



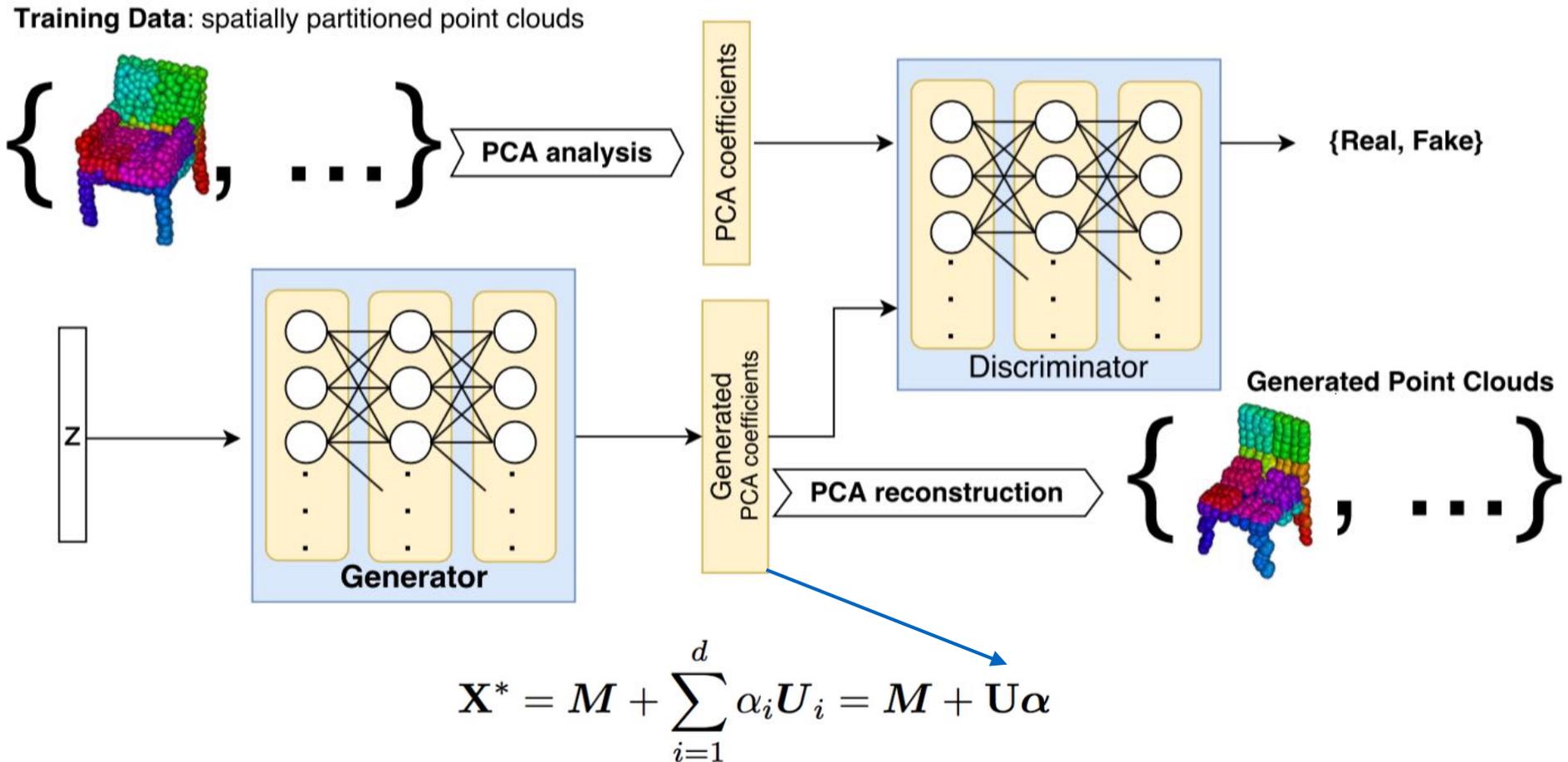
$$\mathbf{X}^* = \mathbf{M} + \sum_{i=1}^d \alpha_i \mathbf{U}_i = \mathbf{M} + \mathbf{U}\boldsymbol{\alpha}$$

Requires perfect correspondence

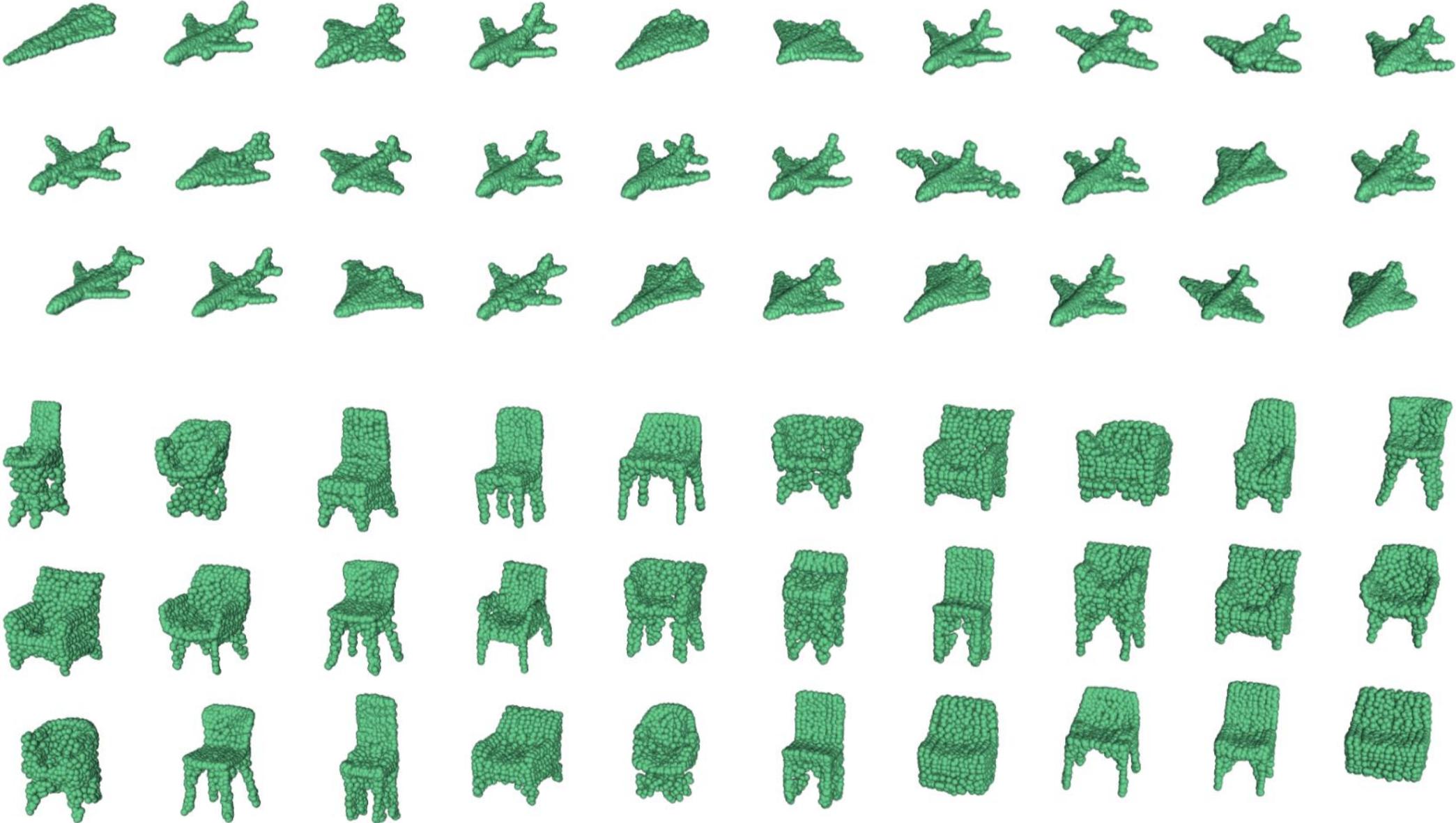
Morphable models

[Figure from Booth et al., 16]

# How important is correspondence?

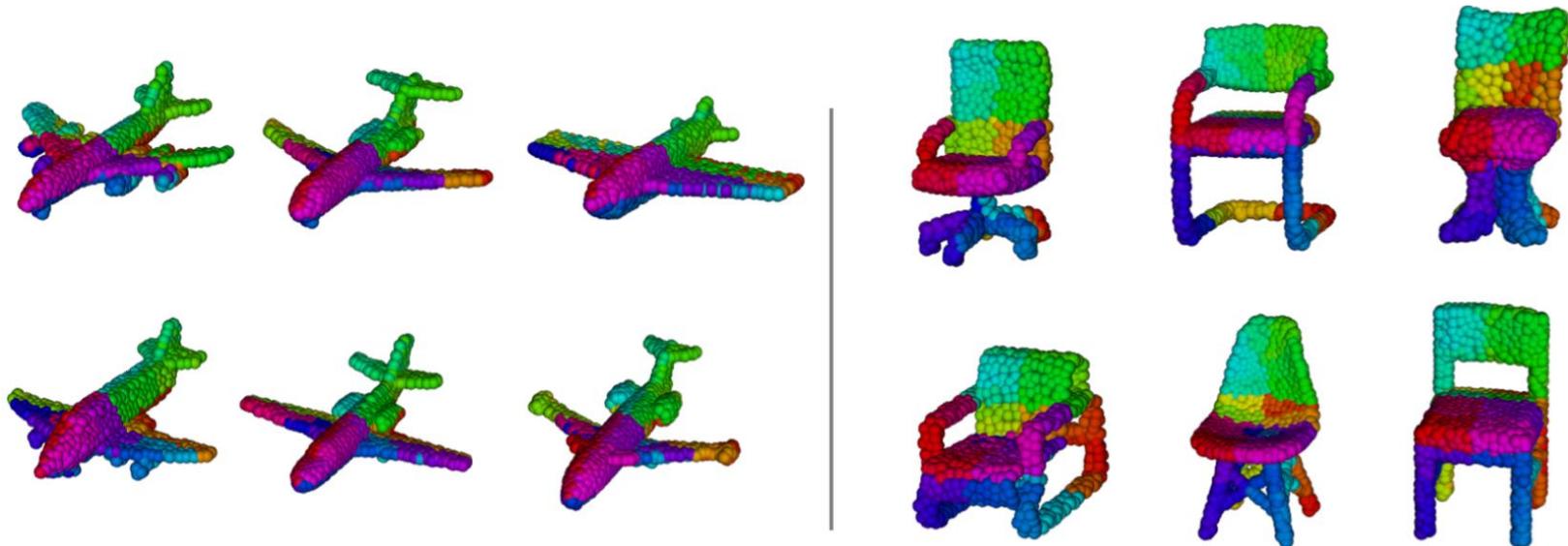


# How important is correspondence?



# Multiresolution tree networks

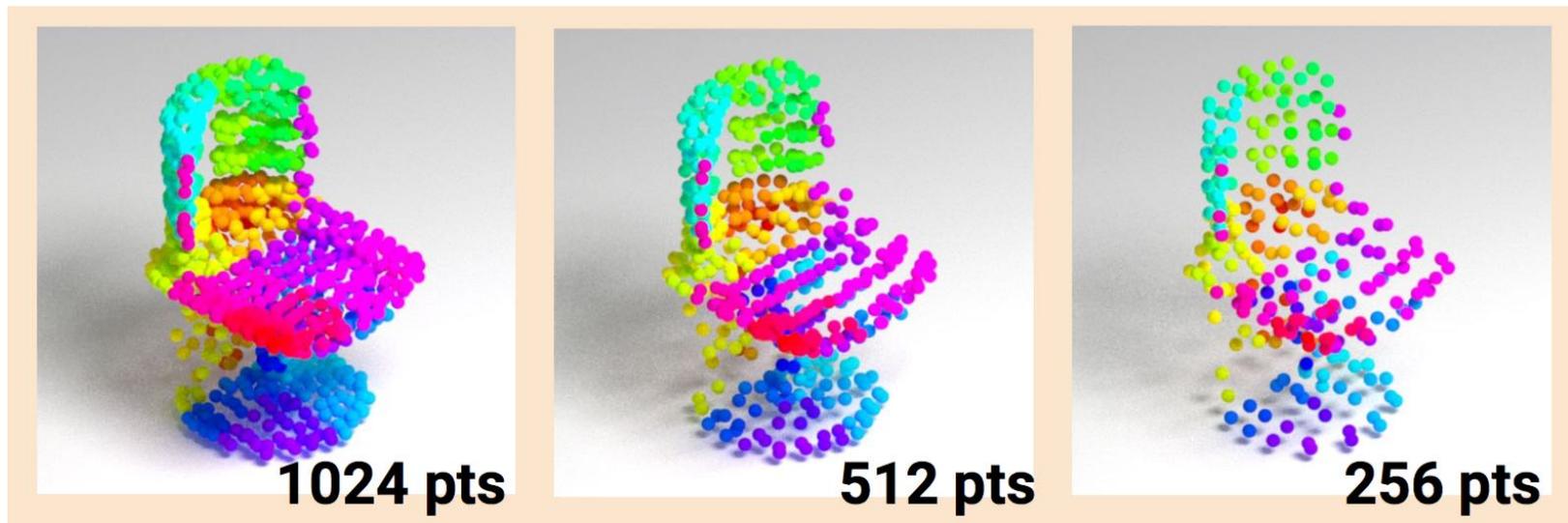
- Addresses the lack of
  - convolutional structure, and
  - coarse to fine reasoning
- Basic idea: linearize 3D points and use 1D convolutions



Points colored with kdtree sort index

# Multiresolution tree networks

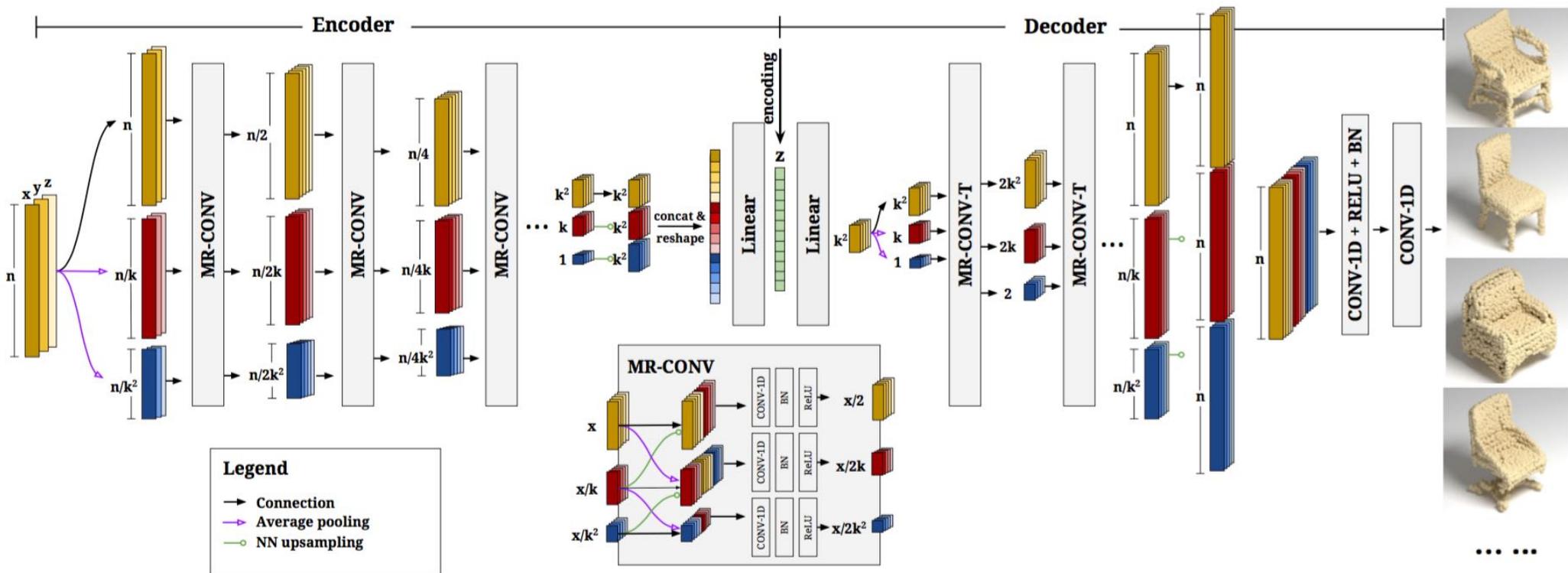
- Addresses the lack of
  - convolutional structure, and
  - coarse to fine reasoning
- Basic idea: linearize 3D points and use 1D convolutions



Implicit multiresolution structure

# Multiresolution tree networks

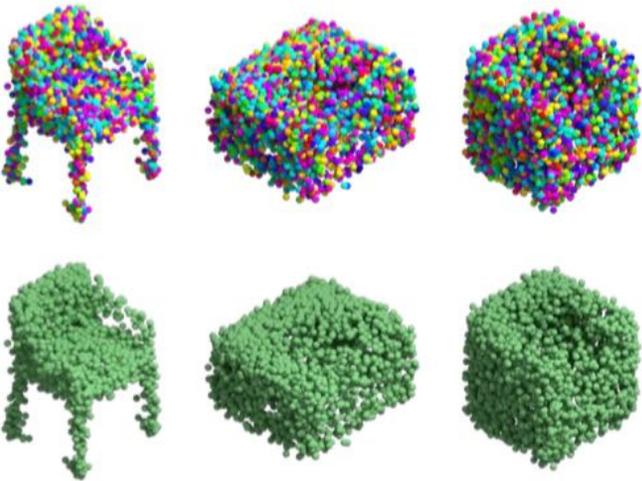
## Architecture for encoding and decoding



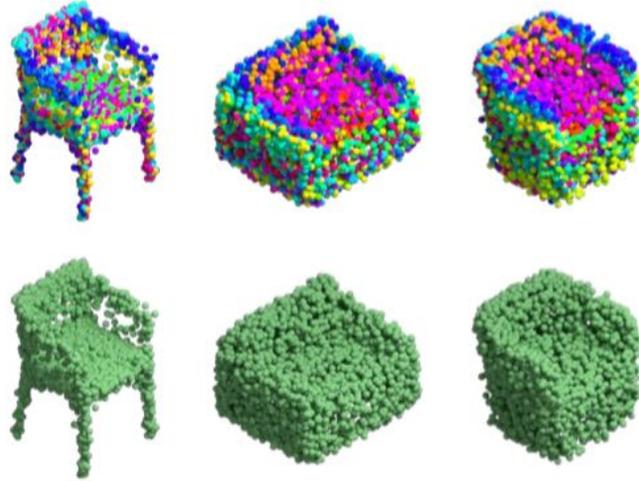
## Multiresolution convolution block

# Does multiresolution analysis help?

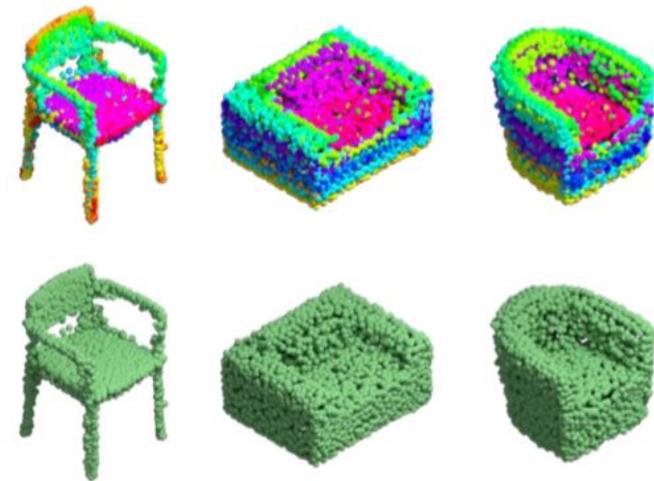
fully-connected



single-resolution

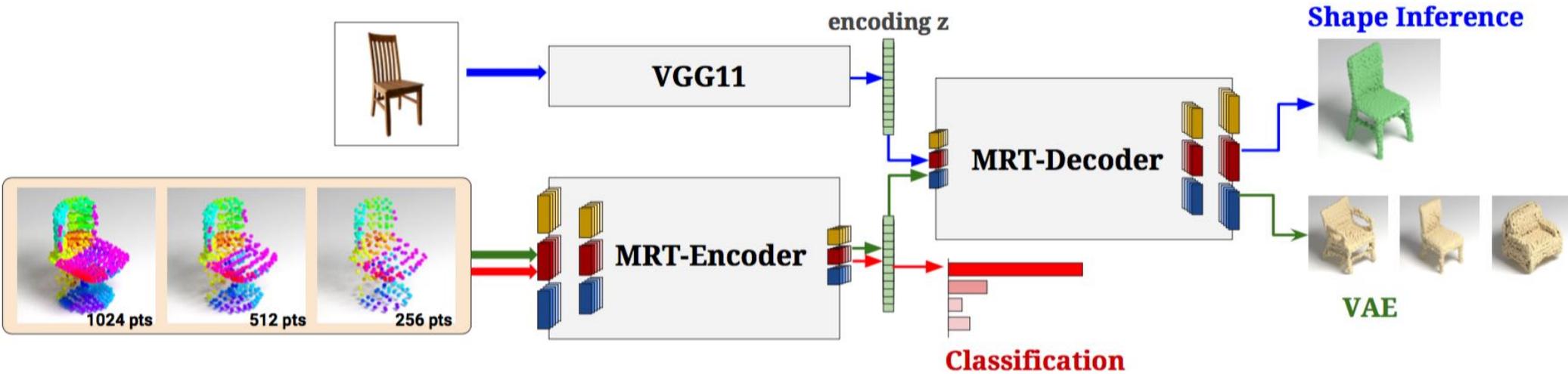


multi-resolution

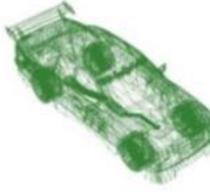
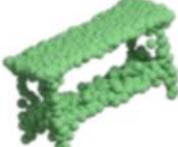
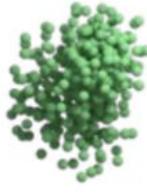
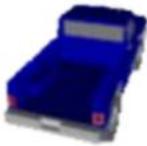
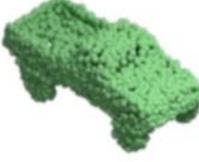
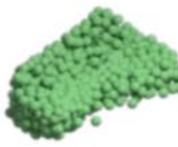


Color indicates the position in the list

# Other shape tasks with MRTNet



# Single image shape reconstruction

Input	G.T.	MRTNet	Fan [12]	Choy [9]
				
				
				
				
				
				
				
				

# Quantitative evaluation: ShapeNet dataset

**Chamfer distance:** pred → GT / GT → pred

voxel-based

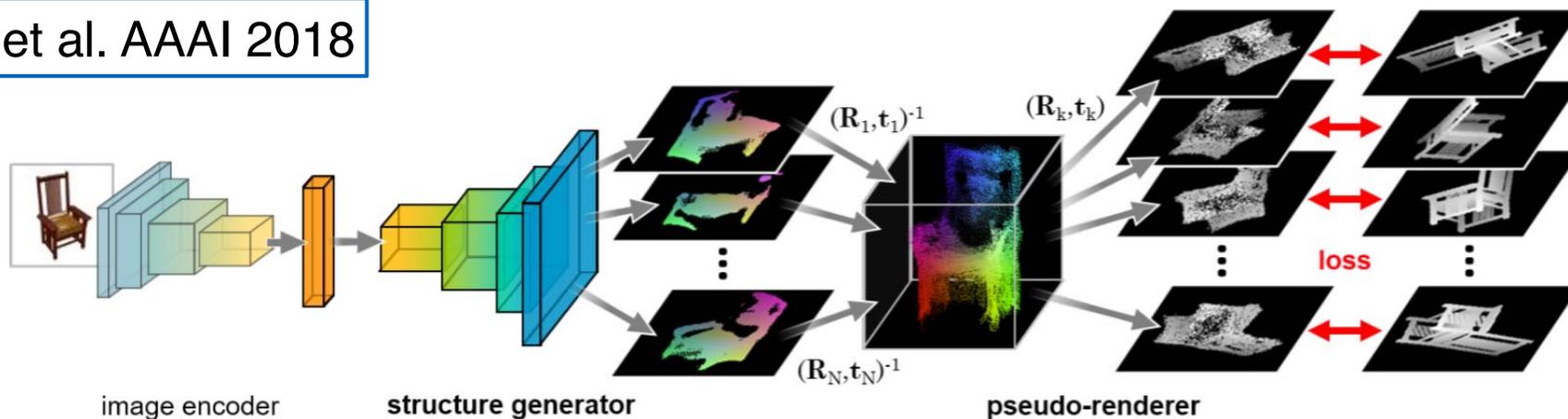
fully-conn.

multiview

Category	3D-R2N2 [9]			Fan et al. [12]	Lin et al. [26]	MRTNet
	1 view	3 views	5 views	(1 view)	(1 view)	(1 view)
mean	3.345 / 4.102	2.702 / 3.465	2.588 / 3.342	1.982 / 2.146	1.846 / 1.701	<b>1.559 / 1.529</b>

$$Ch(\mathbf{x}, \mathbf{y}) = \frac{1}{|\mathbf{x}|} \sum_{x \in \mathbf{x}} \min_{y \in \mathbf{y}} \|x - y\|_2 + \frac{1}{|\mathbf{y}|} \sum_{y \in \mathbf{y}} \min_{x \in \mathbf{x}} \|x - y\|_2$$

Lin et al. AAAI 2018



# Quantitative evaluation: ShapeNet dataset

**Chamfer distance:** pred  $\rightarrow$  GT / GT  $\rightarrow$  pred

voxel-based

fully-conn.

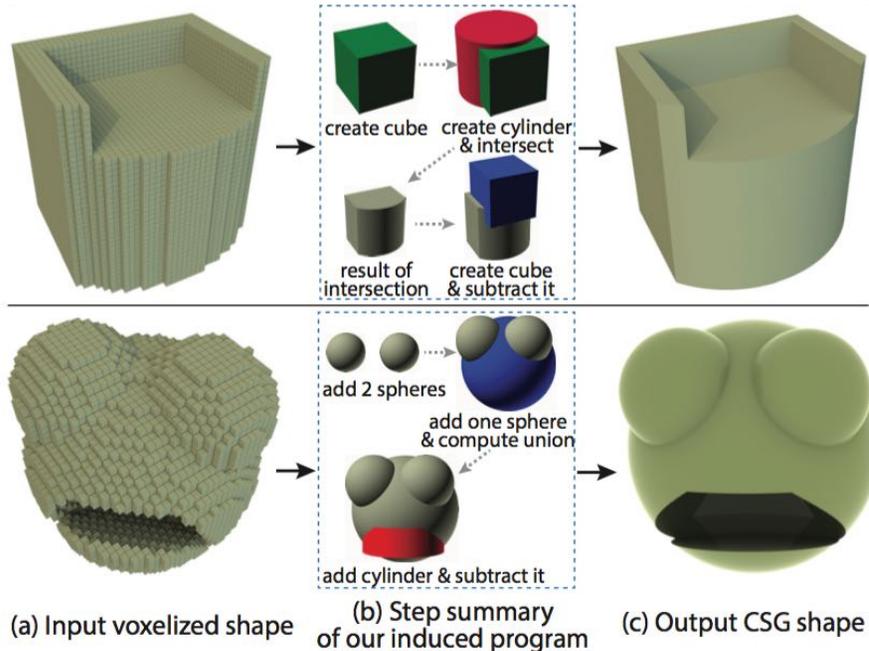
multiview

Category	3D-R2N2 [9]			Fan et al. [12]	Lin et al. [26]	MRTNet
	1 view	3 views	5 views	(1 view)	(1 view)	(1 view)
airplane	3.207 / 2.879	2.521 / 2.468	2.399 / 2.391	1.301 / 1.488	1.294 / 1.541	<b>0.976 / 0.920</b>
bench	3.350 / 3.697	2.465 / 2.746	2.323 / 2.603	1.814 / 1.983	1.757 / 1.487	<b>1.438 / 1.326</b>
cabinet	1.636 / 2.817	1.445 / 2.626	<b>1.420</b> / 2.619	2.463 / 2.444	1.814 / <b>1.072</b>	1.774 / 1.602
car	1.808 / 3.238	1.685 / 3.151	1.664 / 3.146	1.800 / 2.053	1.446 / <b>1.061</b>	<b>1.395</b> / 1.303
chair	2.759 / 4.207	1.960 / 3.238	1.854 / 3.080	1.887 / 2.355	1.886 / 2.041	<b>1.650 / 1.603</b>
display	3.235 / 4.283	2.262 / 3.151	2.088 / 2.953	1.919 / 2.334	2.142 / <b>1.440</b>	<b>1.815</b> / 1.901
lamp	8.400 / 9.722	6.001 / 7.755	5.698 / 7.331	2.347 / 2.212	2.635 / 4.459	<b>1.944 / 2.089</b>
speaker	2.652 / 4.335	2.577 / 4.302	2.487 / 4.203	3.215 / 2.788	2.371 / <b>1.706</b>	<b>2.165</b> / 2.121
rifle	4.798 / 2.996	4.307 / 2.546	4.193 / 2.447	1.316 / 1.358	1.289 / 1.510	<b>1.029 / 1.028</b>
sofa	2.725 / 3.628	2.371 / 3.252	2.306 / 3.196	2.592 / 2.784	1.917 / <b>1.423</b>	<b>1.768</b> / 1.756
table	3.118 / 4.208	2.268 / 3.277	2.128 / 3.134	1.874 / 2.229	1.689 / 1.620	<b>1.570 / 1.405</b>
telephone	2.202 / 3.314	1.969 / 2.834	1.874 / 2.734	1.516 / 1.989	1.939 / <b>1.198</b>	<b>1.346</b> / 1.332
watercraft	3.592 / 4.007	3.299 / 3.698	3.210 / 3.614	1.715 / 1.877	1.813 / 1.550	<b>1.394 / 1.490</b>
<b>mean</b>	3.345 / 4.102	2.702 / 3.465	2.588 / 3.342	1.982 / 2.146	1.846 / 1.701	<b>1.559 / 1.529</b>

# MRTNet summary

- A generic architecture for
  - Point cloud classification (**91.7%** on ModelNet40)
  - Semantic segmentation (see results in the paper)
  - Generation
- Project page: <http://mgadelha.me/mrt/index.html>

# CSGNet: Neural Shape Parser for Constructive Solid Geometry



Gopal Sharma

Rishabh Goyal

Difan Liu

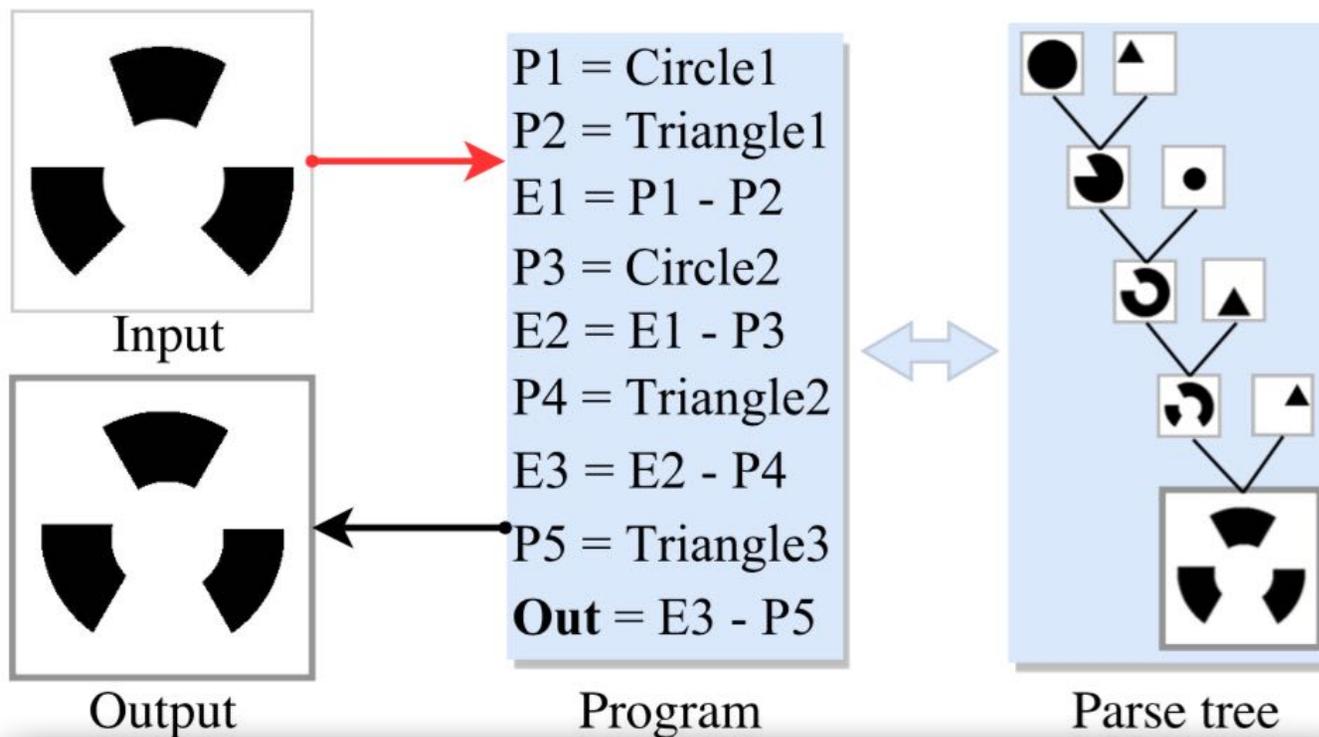
Evangelos Kalogerakis

Subhansu Maji

CVPR 18

interpretable and editable

# Constructive 2D geometry



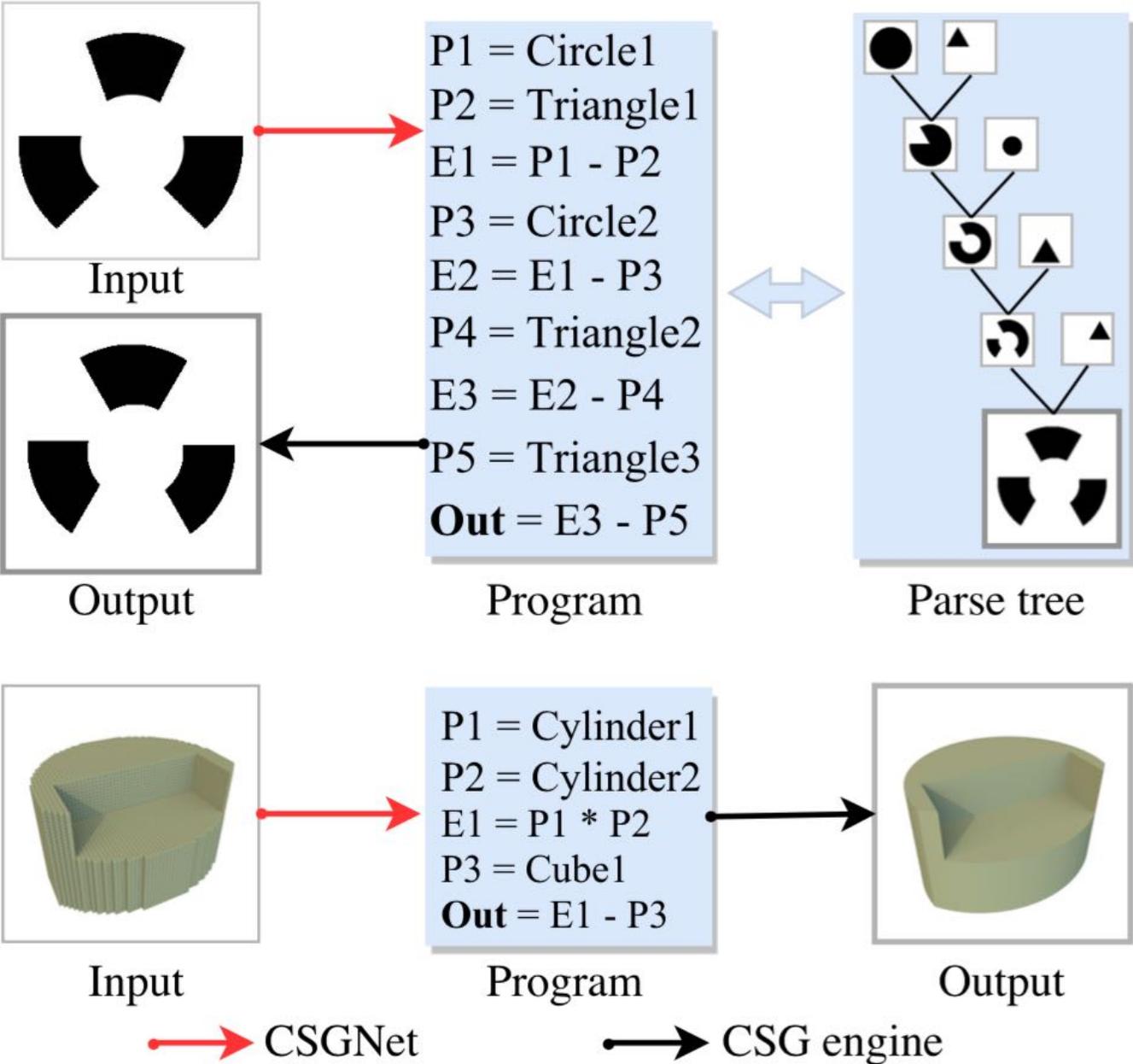
$$S \rightarrow E$$

$$E \rightarrow E E T \mid P$$

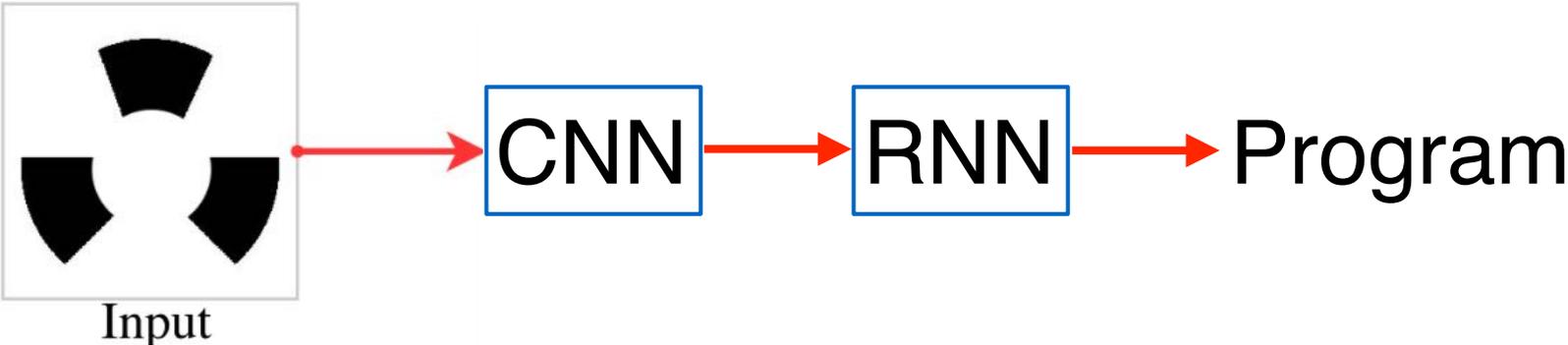
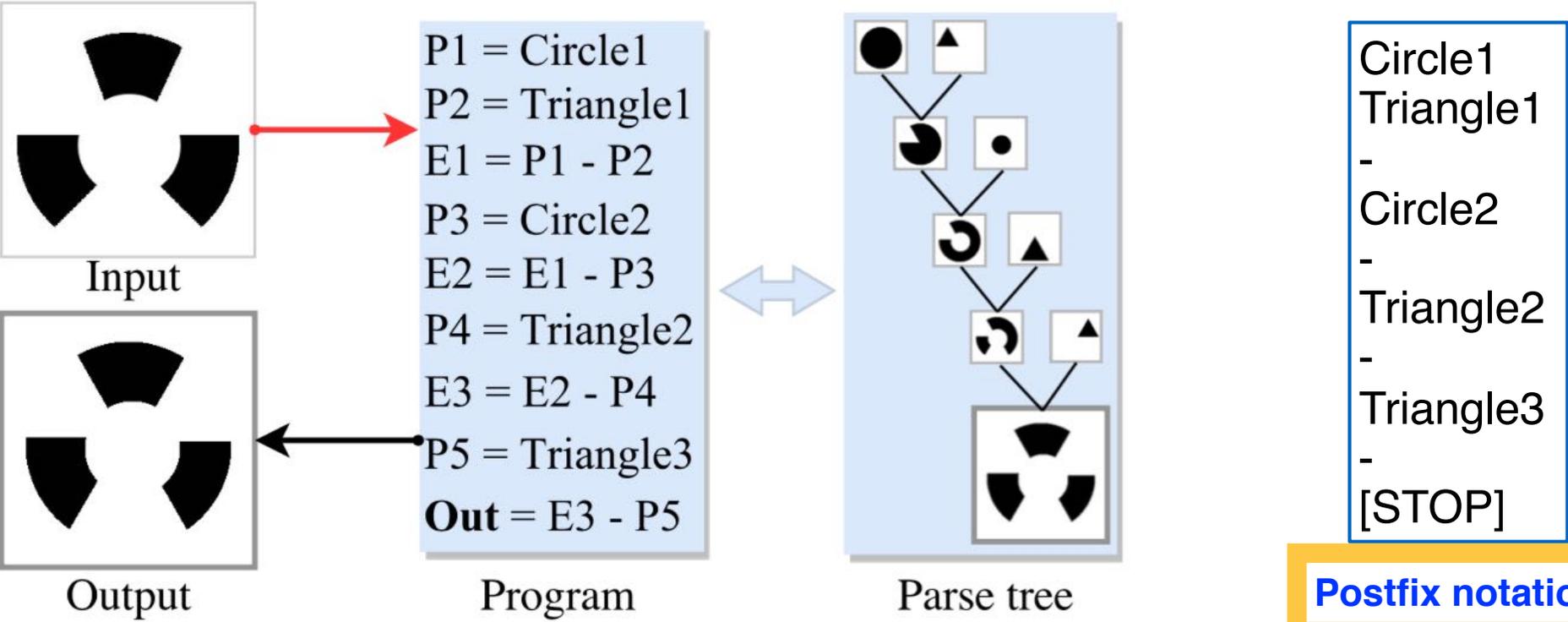
$$T \rightarrow OP_1 \mid OP_2 \mid \dots \mid OP_m$$

$$P \rightarrow SHAPE_1 \mid SHAPE_2 \mid \dots \mid SHAPE_n$$

# Constructive solid geometry

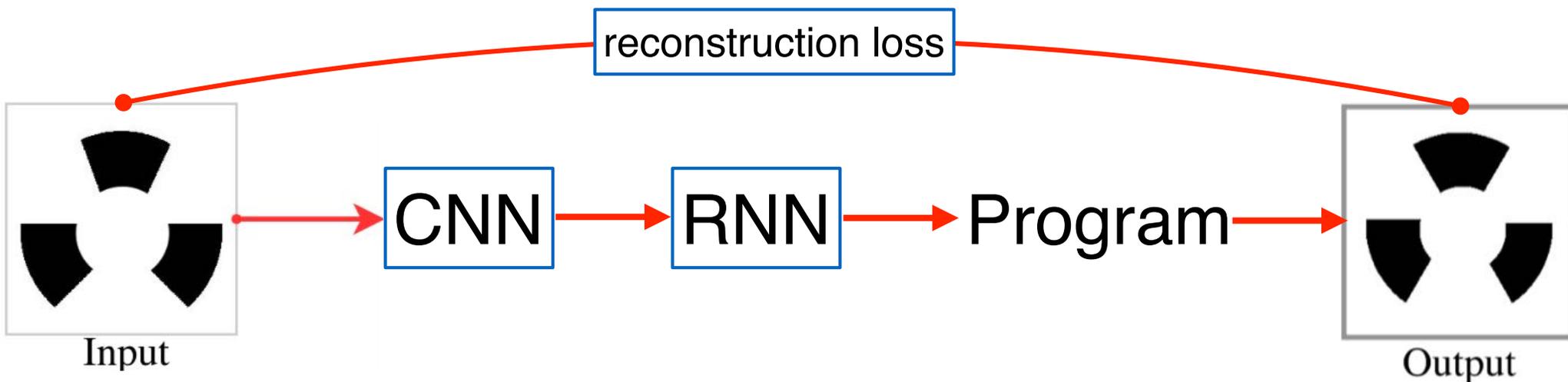


# Constructive solid geometry

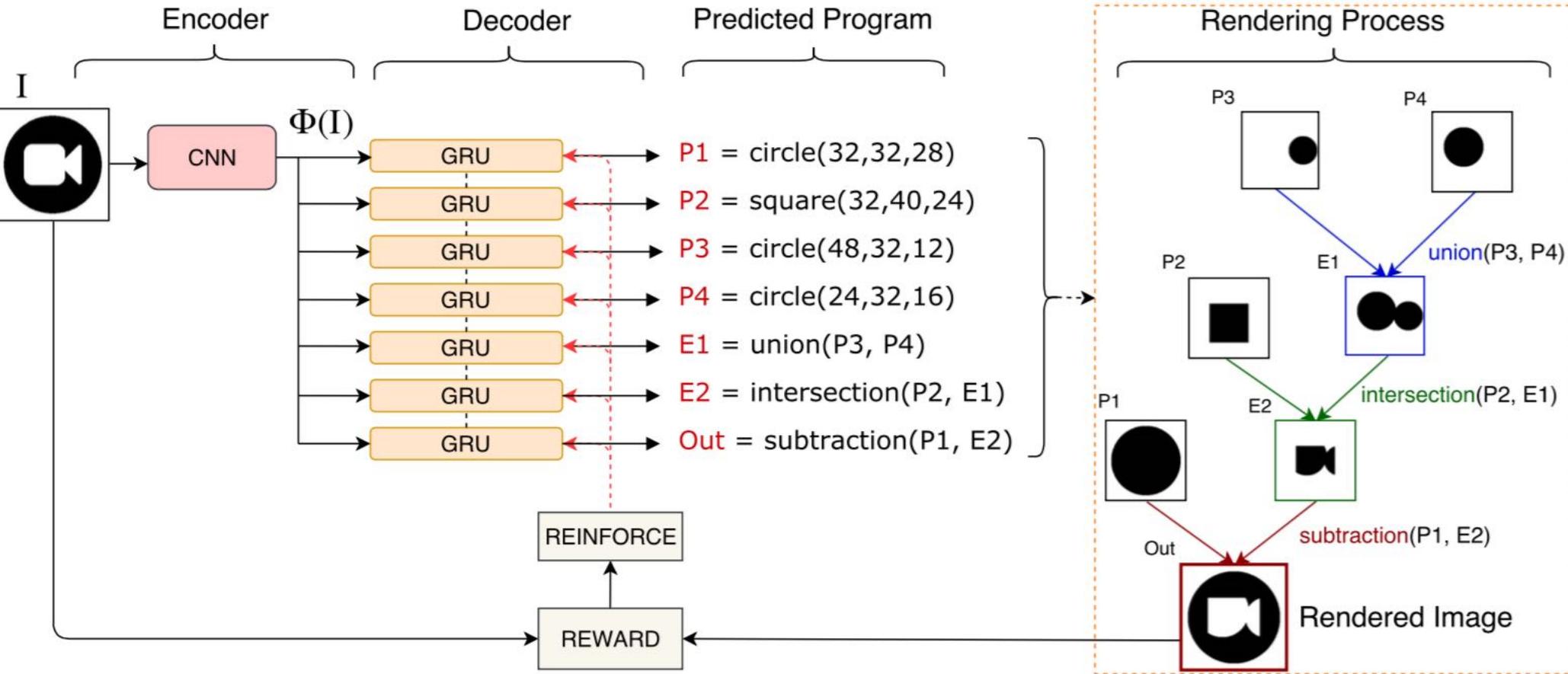


# Learning

- **Supervised setting:** learn to predict programs directly
- **Unsupervised setting:** No ground-truth programs.
  - Learn parameters to minimize a reconstruction error through policy gradients [REINFORCE, Williams 1992]

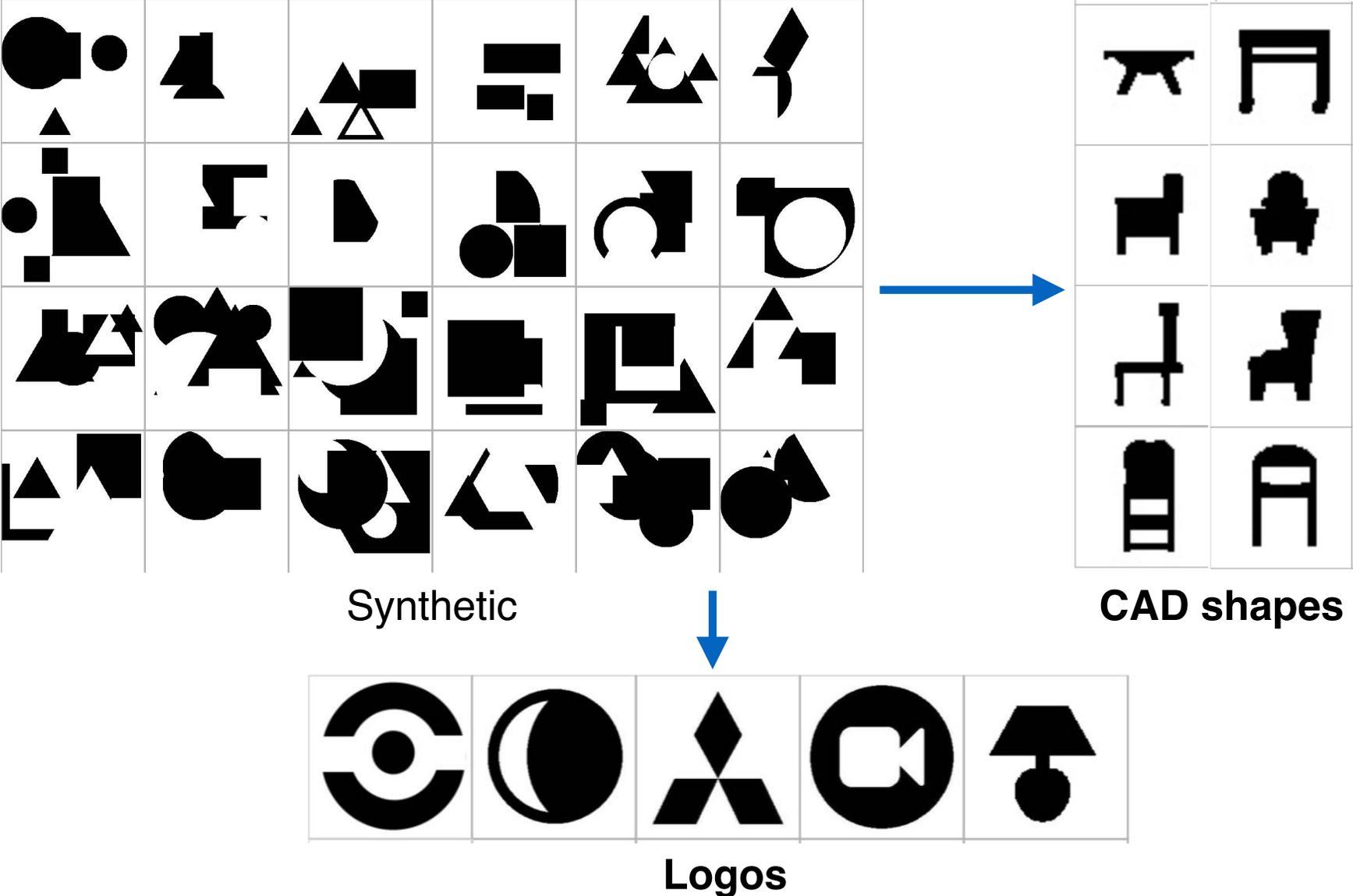


# CSGNet: 2D/3D $\Leftrightarrow$ programs



# CSGNet: 2D/3D $\Leftrightarrow$ programs

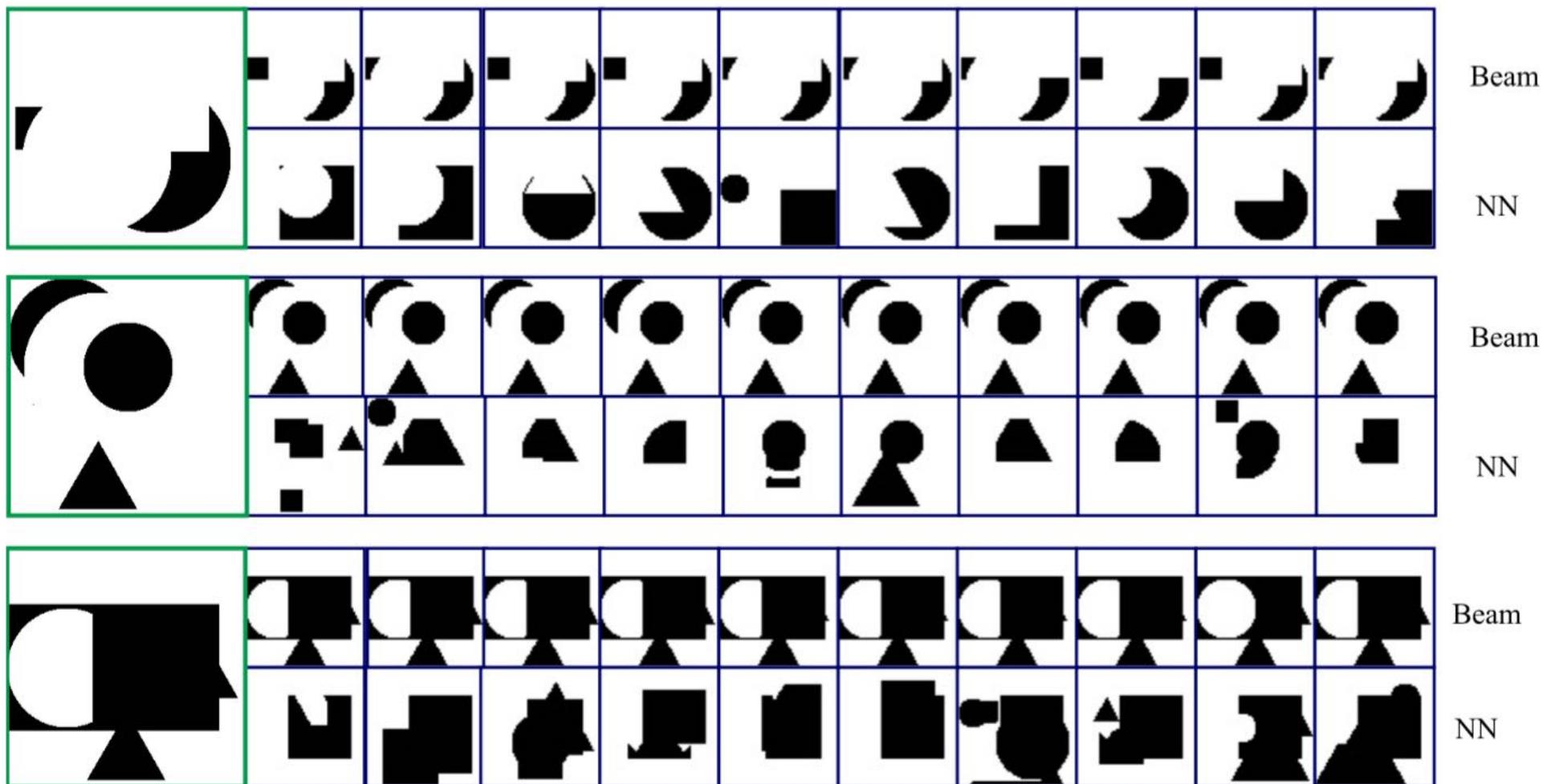
Train on synthetic data and adapt to new domains using policy gradients



# CSGNet: 2D/3D $\Leftrightarrow$ programs

How well does the nearest neighbor perform?

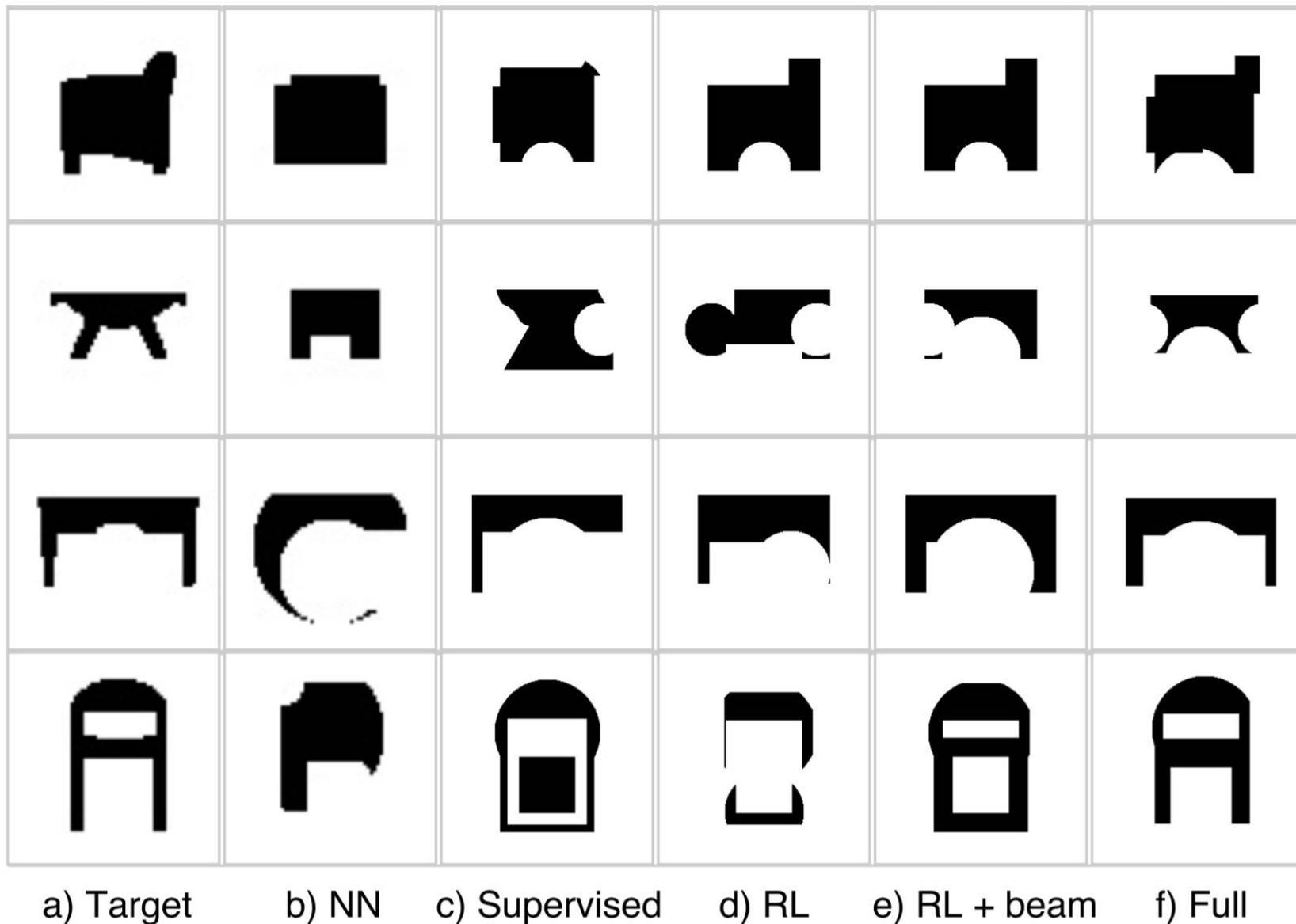
Chamfer distance **1.88** (NN), **1.36** (CSGNet) with 675K training examples





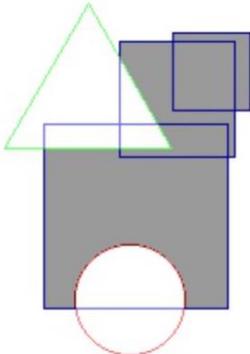
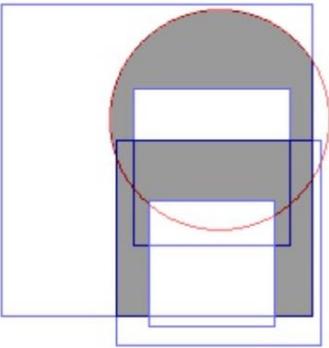
# CSGNet: 2D/3D $\Leftrightarrow$ programs

CAD shapes dataset: Chamfer distance **1.94** (NN), **0.51** (CSGNet)



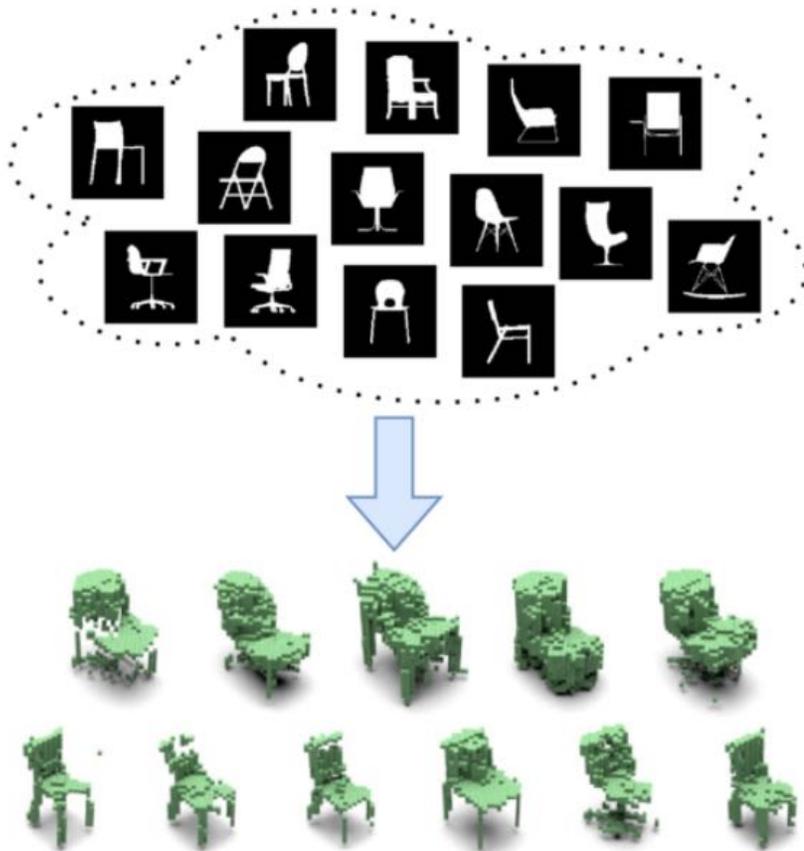
# CSGNet: 2D/3D $\Leftrightarrow$ programs

Input

			<pre>s(34,40,18) c(33,21,9) union s(40,31,25) intersect c(33,52,9) subtract s(29,35,12) union</pre>
			<pre>s(34,41,21) s(33,44,13) subtract c(34,23,16) s(33,30,16) subtract union s(25,29,32) intersect</pre>

- More results in the paper: reward shaping, comparison to Faster R-CNN for primitive detection, results on 3D, etc.
- Preprint available: <https://arxiv.org/abs/1712.08290>

# Learning 3D Shape Representations with Weak Supervision



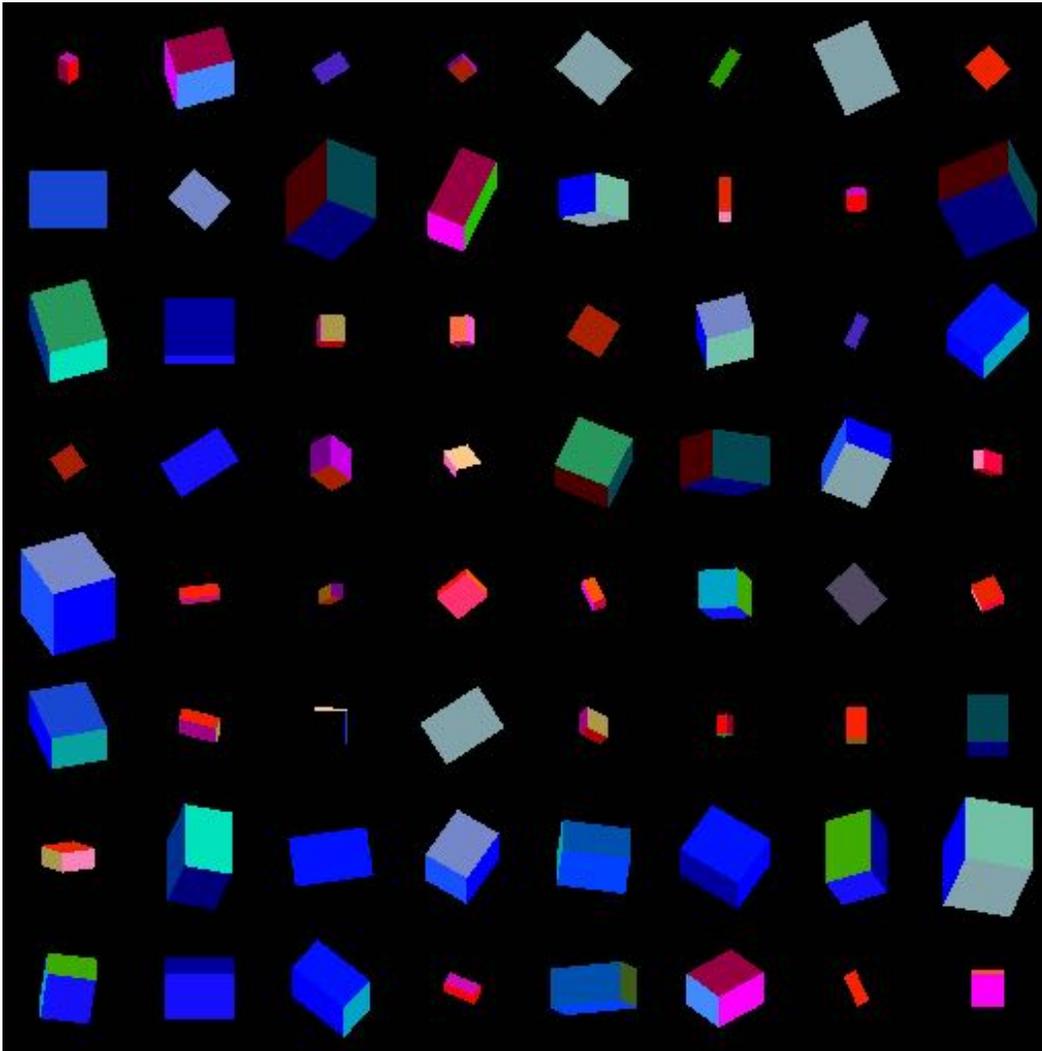
Matheus Gadelha  
Subhransu Maji  
Rui Wang

3DV 2017

# Related Work

- 3D shape from collection of images
  - Visual hull — same instance, known viewpoints
  - Photometric stereo — same instance, known lighting, simple reflectance
  - Structure from motion — same instance (or 3D)
  - Non-rigid structure from motion — known shape family (e.g., faces)
  - **Our work** — unknown shape family, unknown viewpoints
- 3D shape from single image
  - Optimization-based approaches;
  - Recognition-based approaches;

# A motivating example

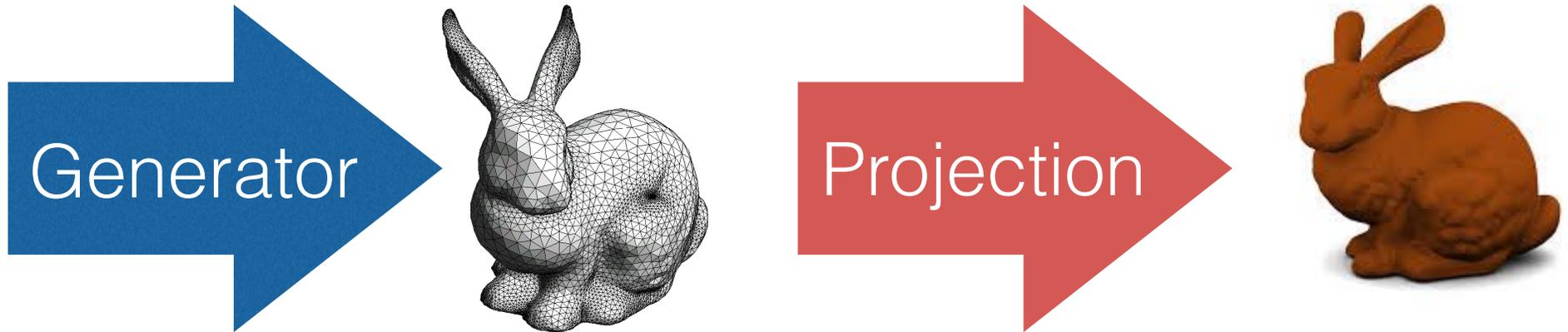


Small cubes are reddish  
Big cubes are bluish

**Hypothesis: It is easier to generate these images by reasoning in 3D**

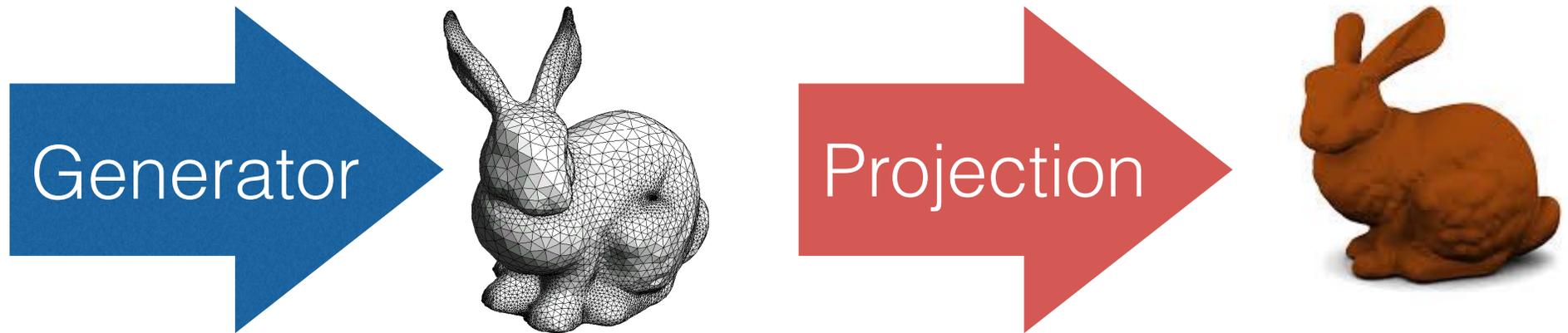
# Approach

- Our goal is to learn a 3D shape generator whose projections match the provided set of the views



# Approach

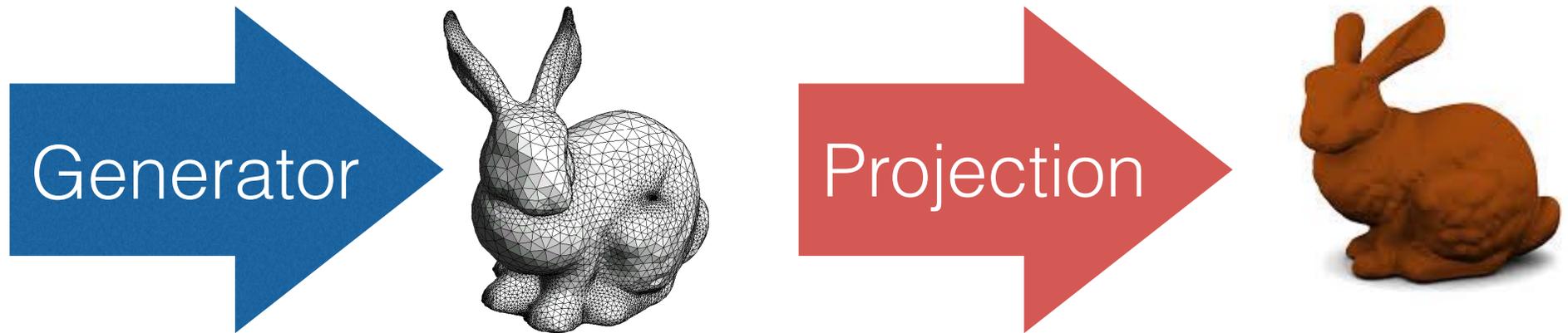
- Our goal is to learn a 3D shape generator whose projections match the provided set of the views



How do we match distributions?

# Approach

- Our goal is to learn a 3D shape generator whose projections match the provided set of the views

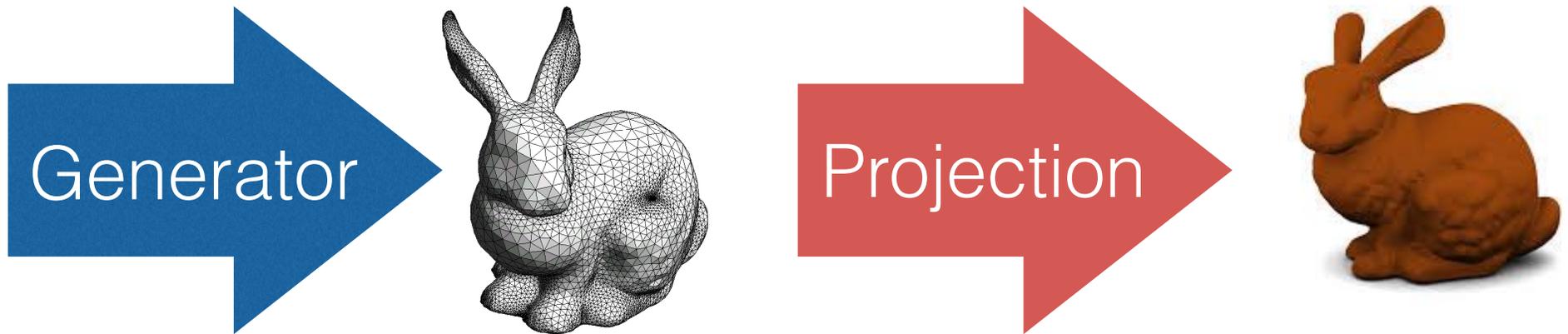


How do we match distributions?

$$\min_G D_{\text{KL}}(\overset{\text{generated}}{G} || \overset{\text{true}}{D}) = \min_{z \sim G} \mathbb{E} \left[ \log \frac{G(z)}{D(z)} \right]$$

# Approach

- Our goal is to learn a 3D shape generator whose projections match the provided set of the views



How do we match distributions?

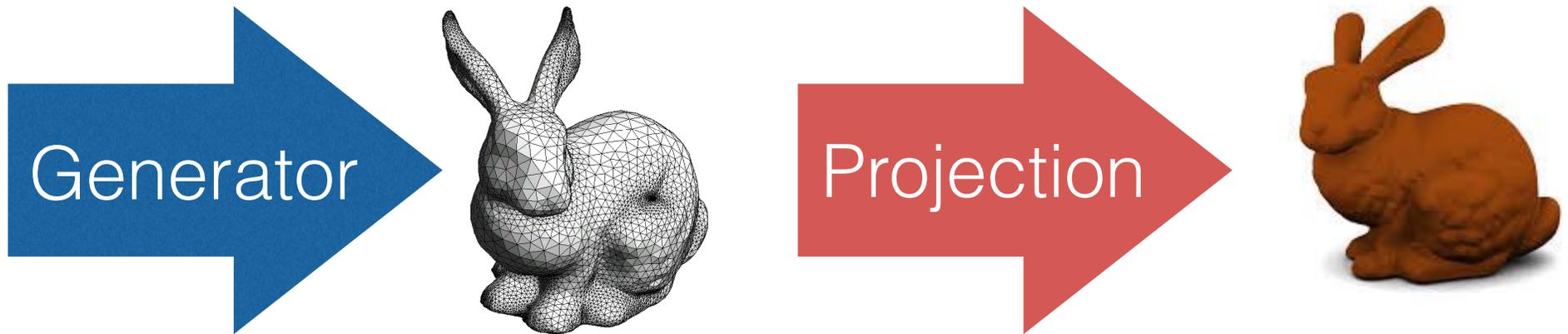
$$\min_G D_{\text{KL}}(G || D) = \min_{z \sim G} \mathbb{E} \left[ \log \frac{G(z)}{D(z)} \right]$$

generated      true

estimate using logistic regression

# Approach

- Our goal is to learn a 3D shape generator whose projections match the provided set of the views



How do we match distributions?

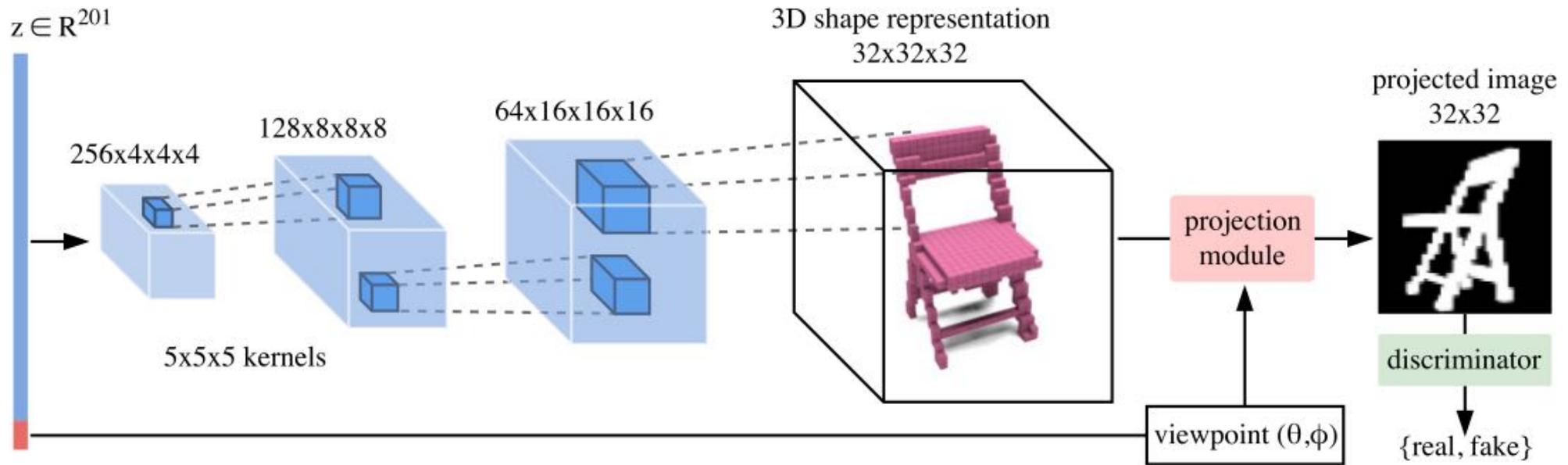
$$\min_G D_{\text{KL}}(\overset{\text{generated}}{G} || \overset{\text{true}}{D}) = \min_{z \sim G} \mathbb{E} \left[ \log \frac{G(z)}{D(z)} \right] \quad \begin{array}{l} \text{estimate using} \\ \text{logistic regression} \end{array}$$

$$\min_G \max_d \mathbb{E}_{x \sim D} [\log d(x)] + \mathbb{E}_{z \sim G} [\log(1 - d(z))]$$

Generative adversarial networks [Goodfellow et al.]

# PrGAN

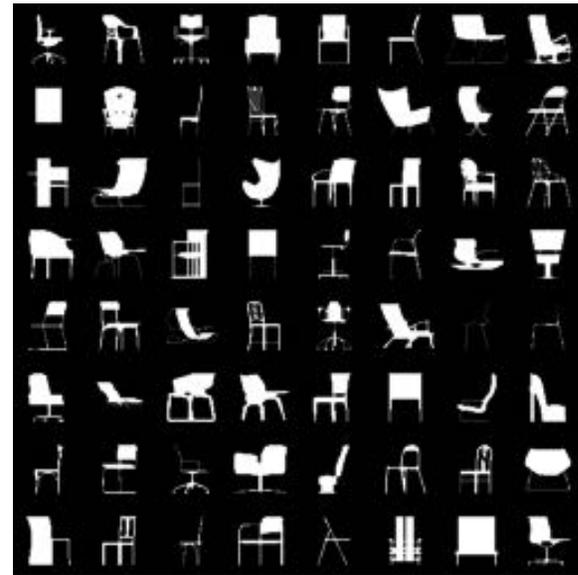
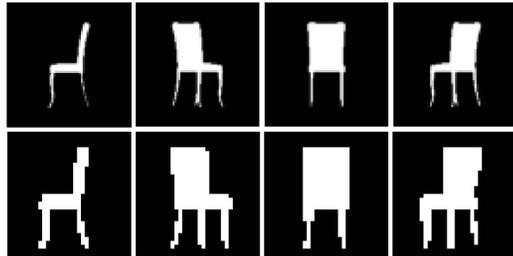
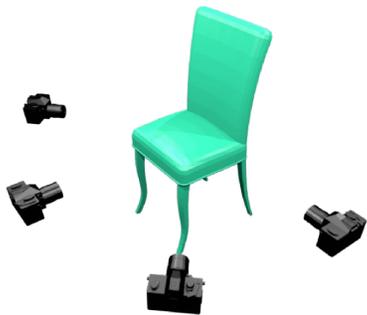
Generator maps  $z$  to a voxel occupancy grid and a viewpoint



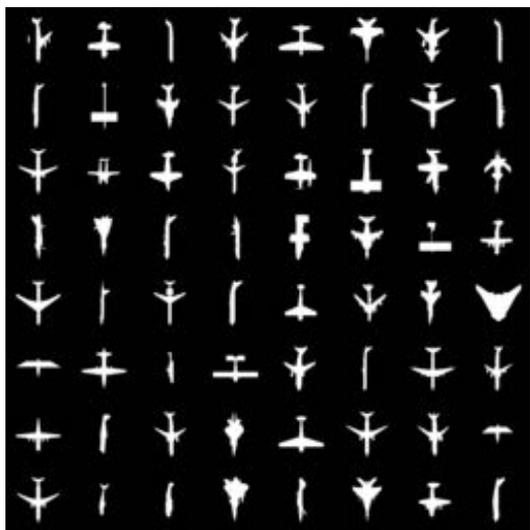
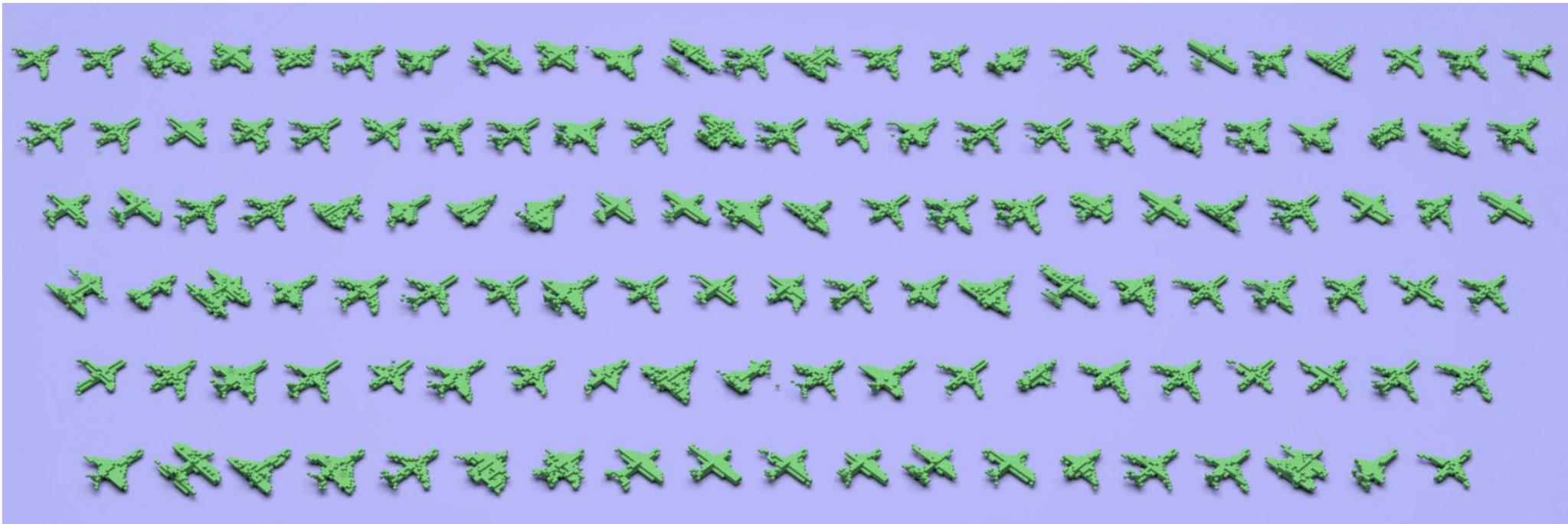
Projection using line integration along the view direction

$$I(\mathbf{x}) = 1 - \exp\left(-\int_0^\infty V(\mathbf{x} + \mathbf{r})dr\right)$$

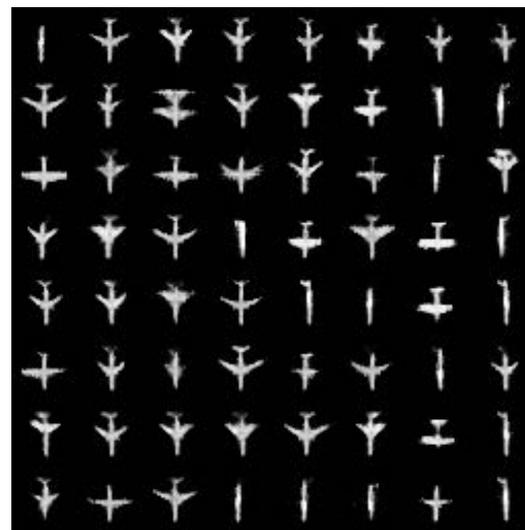
# Dataset generation



# Airplanes

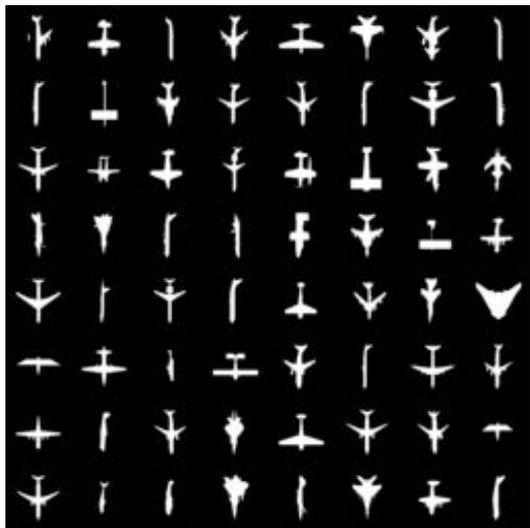
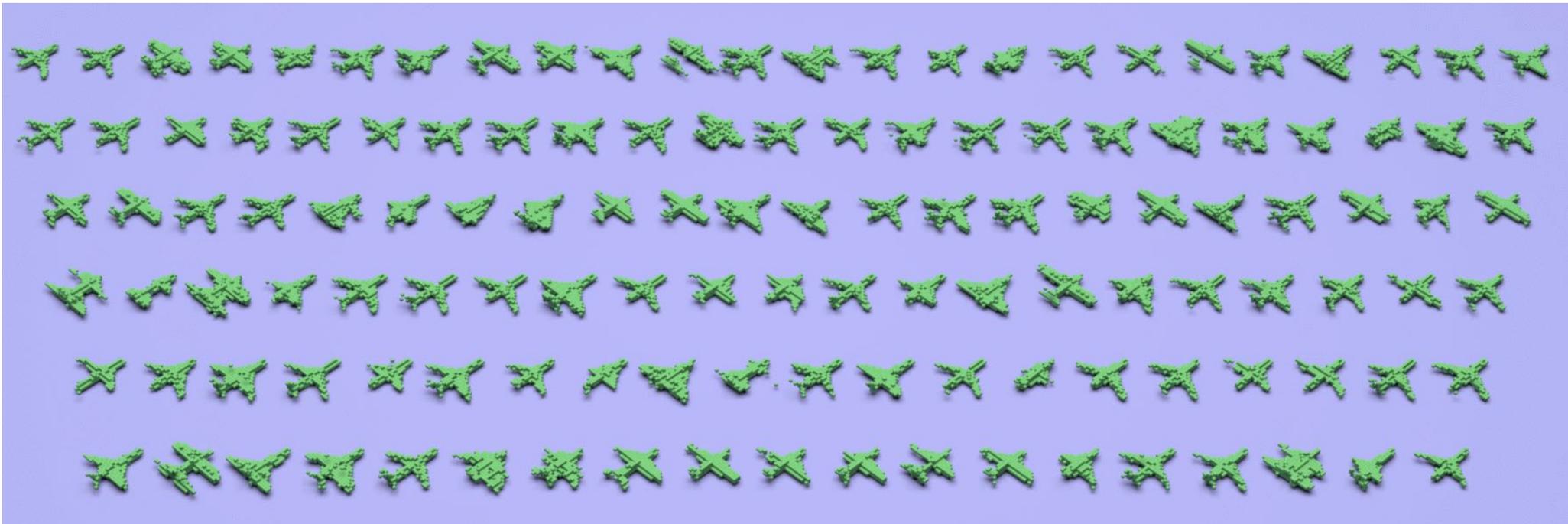


input

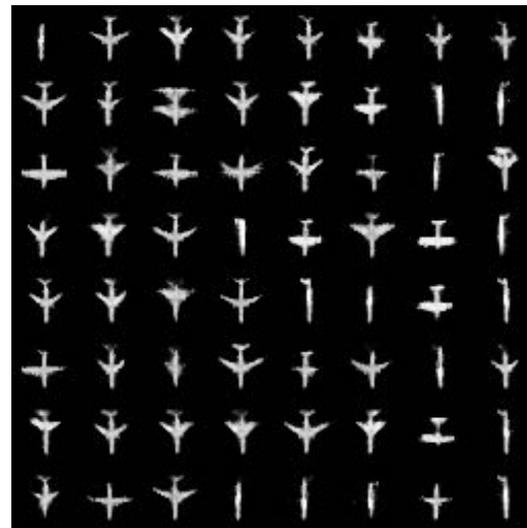


generated

# Airplanes



input

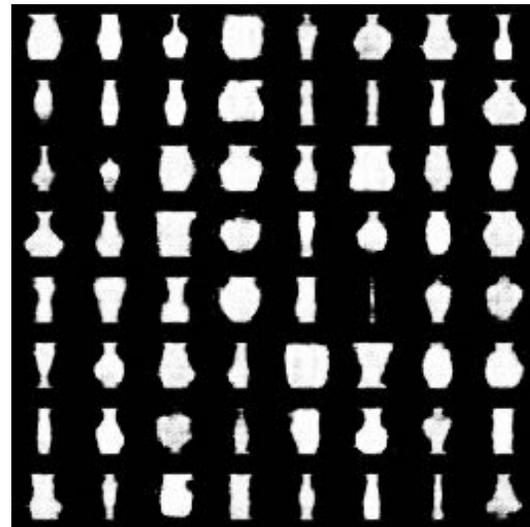


generated

# Vases

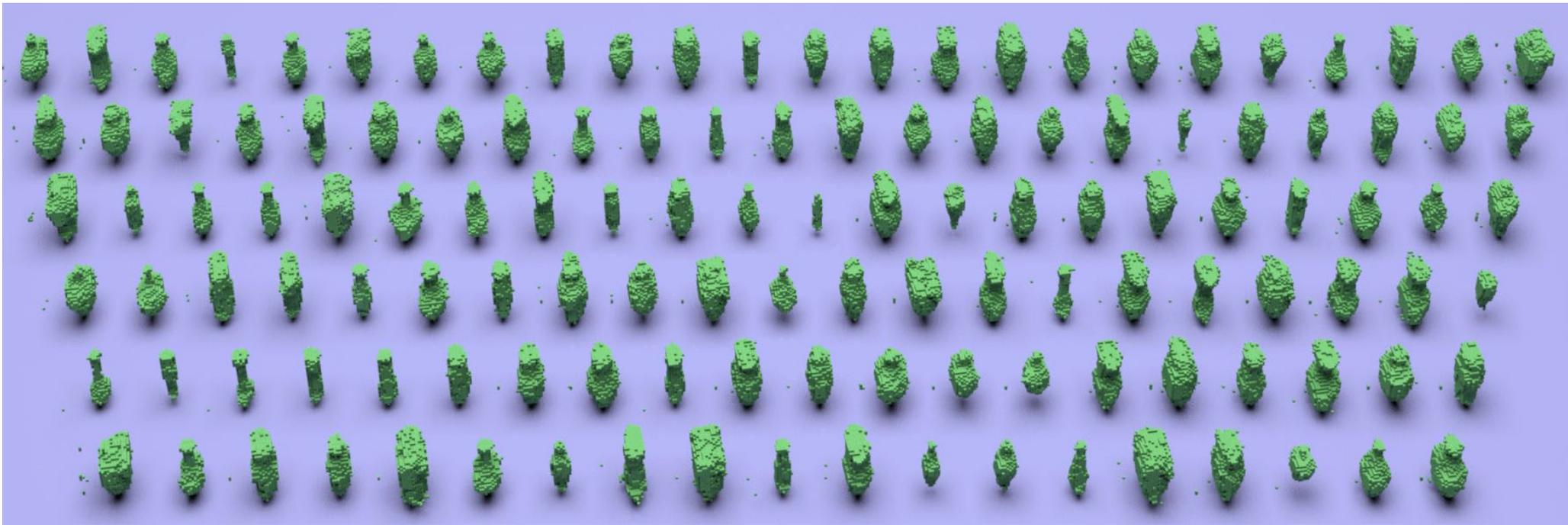


input

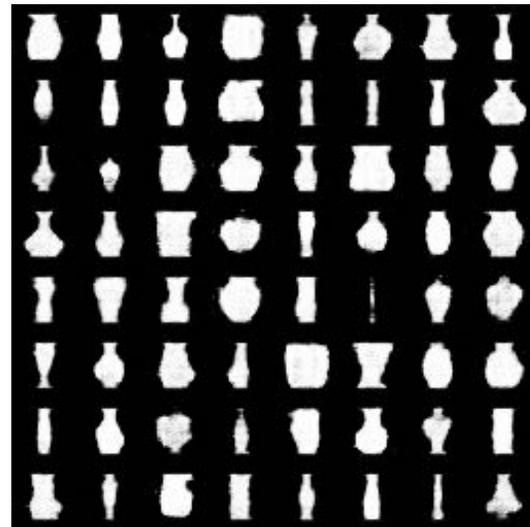


generated

# Vases



input

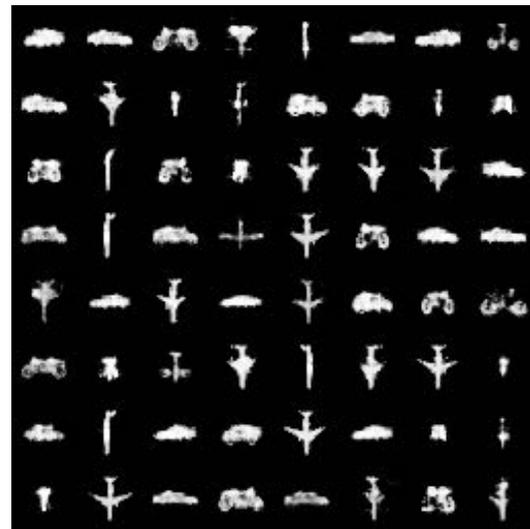


generated

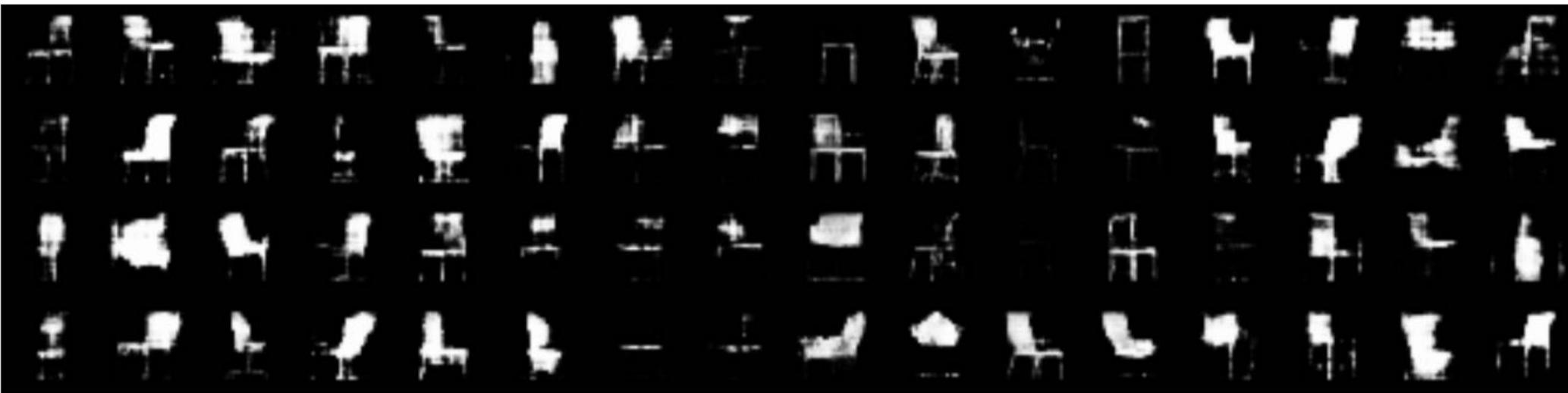
# Mixed categories



input



generated

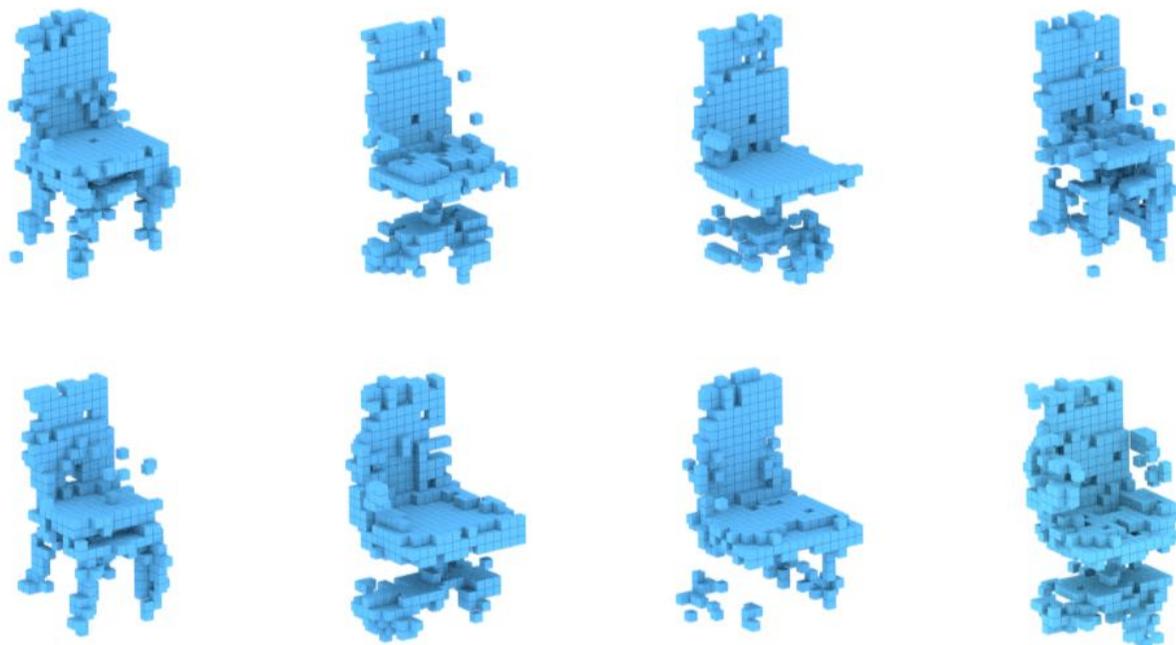


(a) Results from 2D-GAN.



(a) Results from PrGAN.

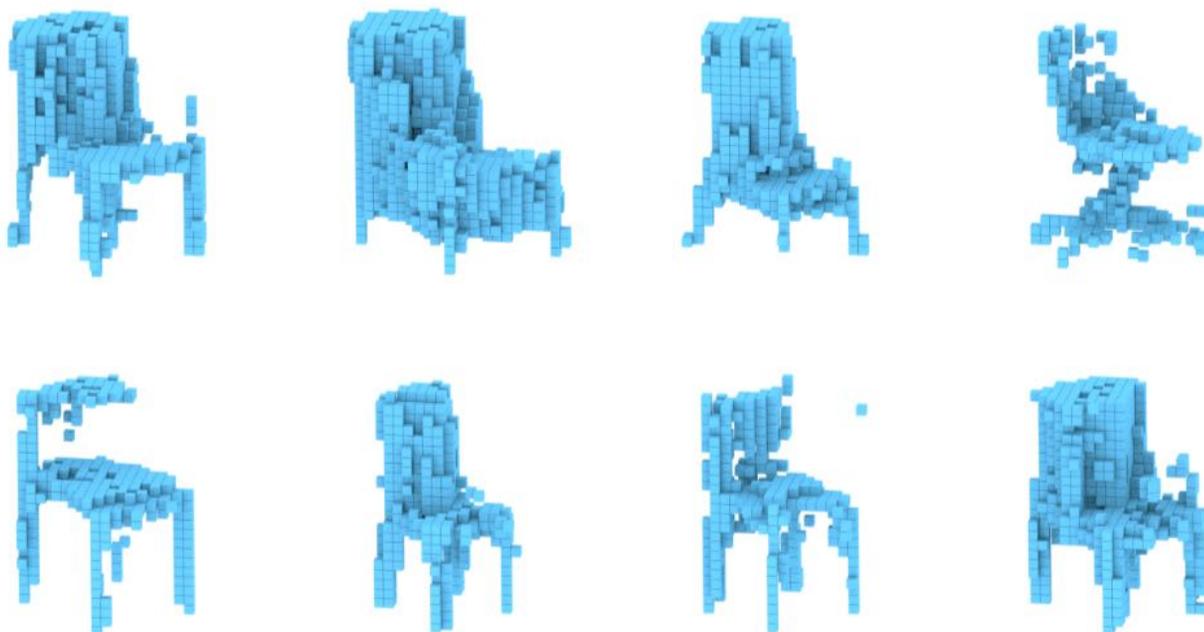
**MMD metric: 2D-GAN 90.1, PrGAN 88.3**



Trained with 3D supervision

(a) Results from 3D-GAN.

**MMD metric**  
3D-GAN 347.5  
PrGAN 442.9



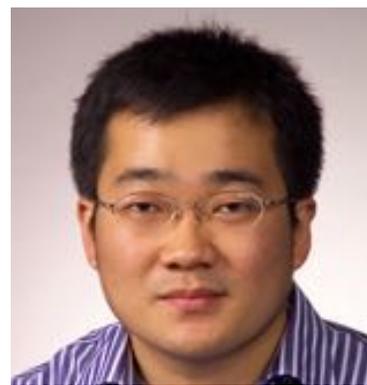
(a) Results from PrGAN.

# Projection GAN

- The model is able to recover the coarse 3D structure
- But should use side information when available
  - Viewpoint
  - Landmarks / part labels
  - Pose estimates
- Iterative: bootstrap 3D to estimate pose & viewpoint

# Thank you!

- Collaborators: Matheus Gadelha, Zhaoliang Lun, Gopal Sharma, Rui Wang, Evangelos Kalogerakis



- Funding from NSF, NVIDIA, Facebook
- <https://people.cs.umass.edu/smaji/projects.html>