

# **Surfaces, Objects, Procedures: Integrating Learning and Graphics for 3D Scene Understanding**

Jiajun Wu

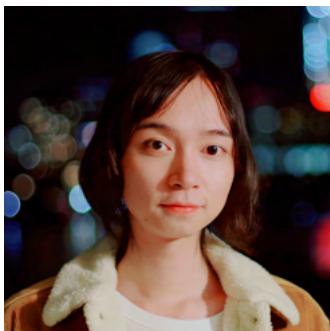




Xiuming Zhang



Zhoutong Zhang



Jiayuan Mao



Yikai Li



Bill Freeman



Josh Tenenbaum



Noah Snavely



Jun-Yan Zhu



# 3D Scene Understanding



How many cars are there in the scene?

What's the color of the closest car?

How would the car look like from the front?

What's going to happen?

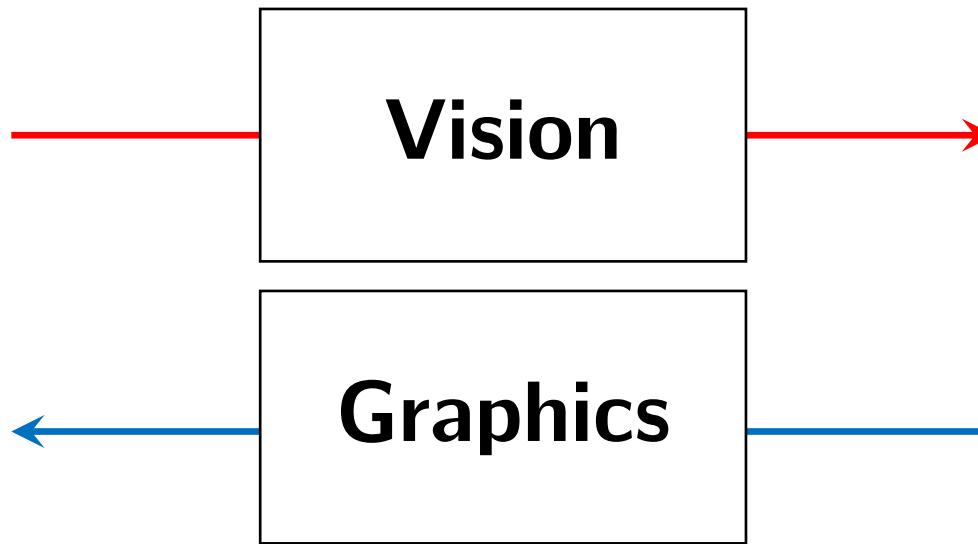
What if it is rainy?



Also applies to **novel** objects



2D Image

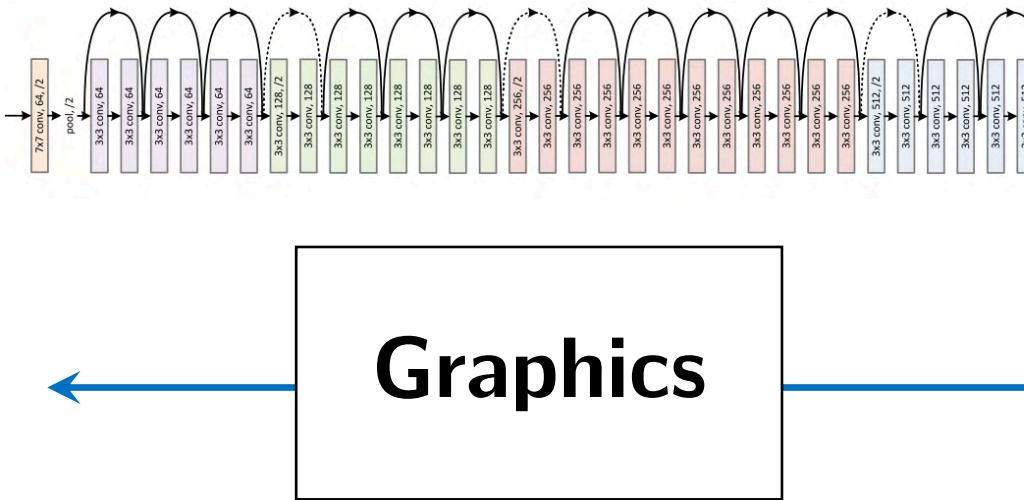


3D Shape

# Deep Learning

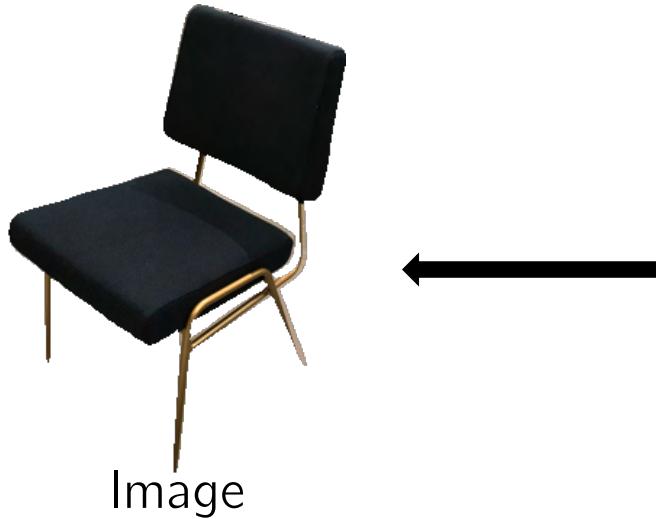


2D Image



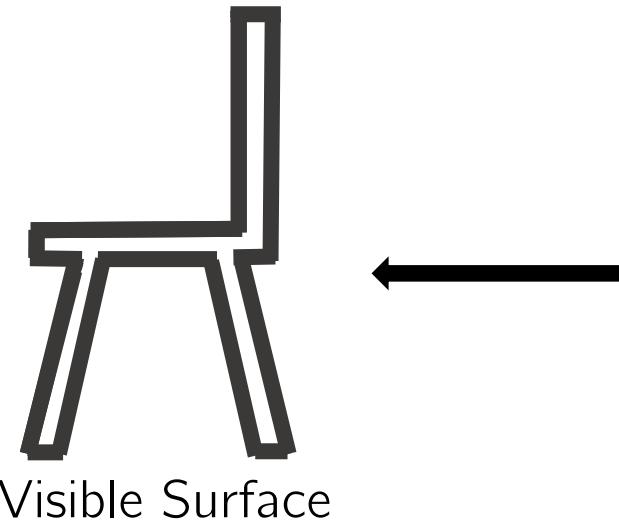
3D Shape

## Forward: Image Formation



Image

/



Visible Surface

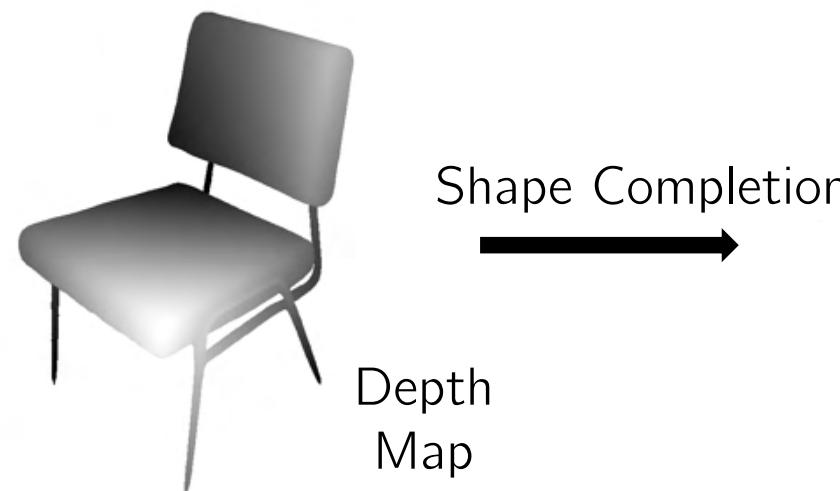


3D Shape

## Inverse: Shape Estimation



Depth Estimation



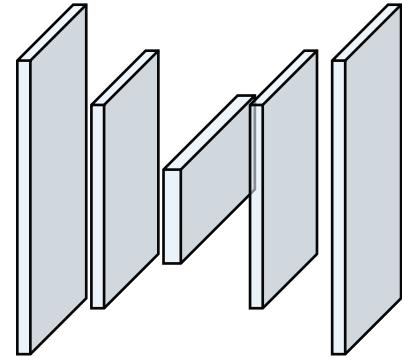
Shape Completion

Depth Map

# Inverting the Graphics Engine



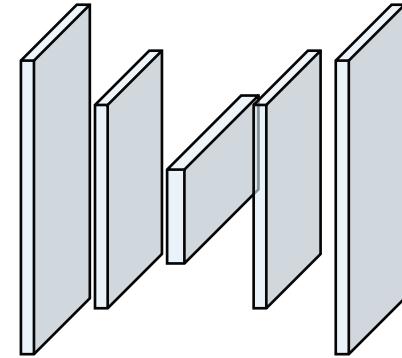
Image



Depth Estimation



Depth

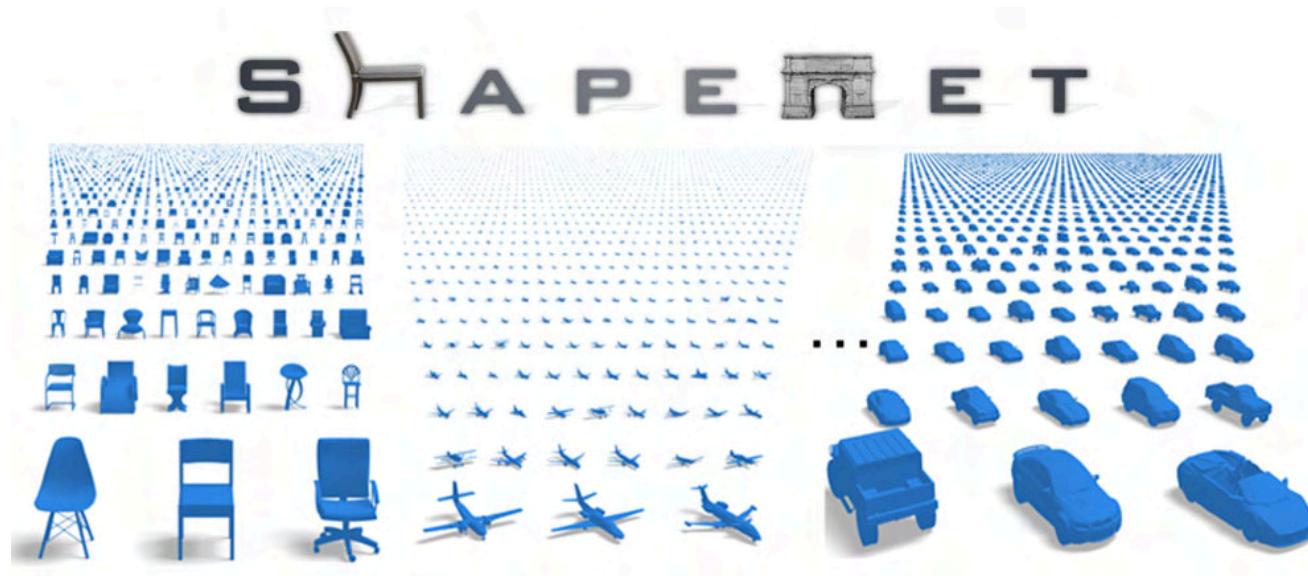


Shape Completion

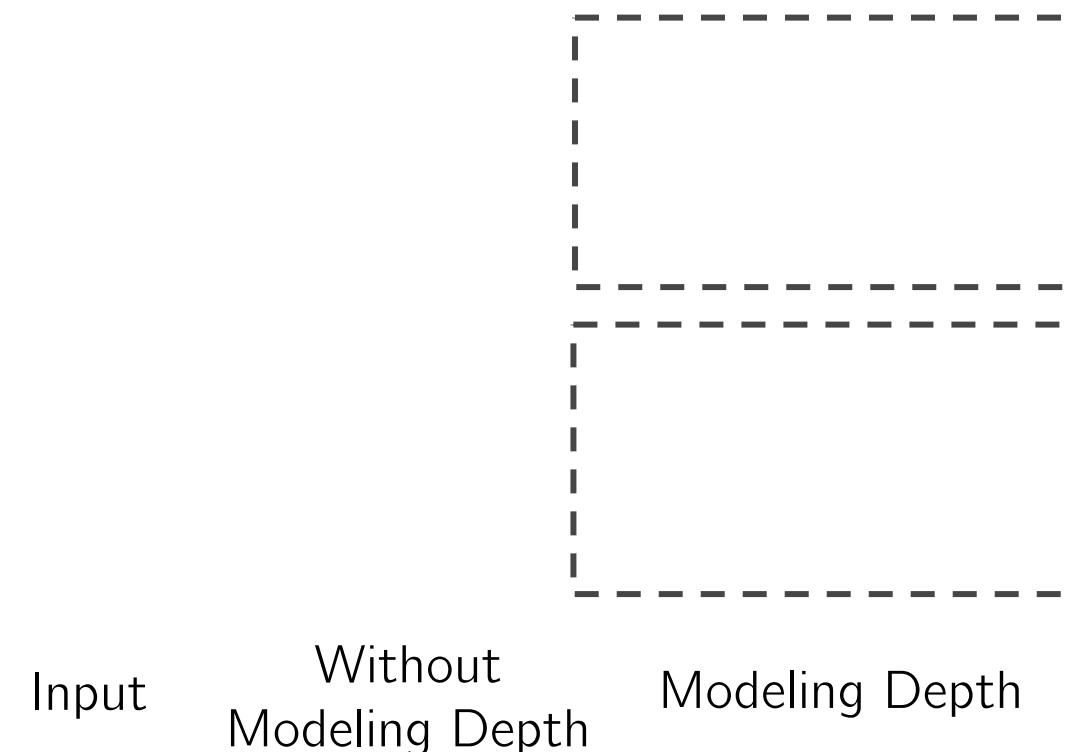
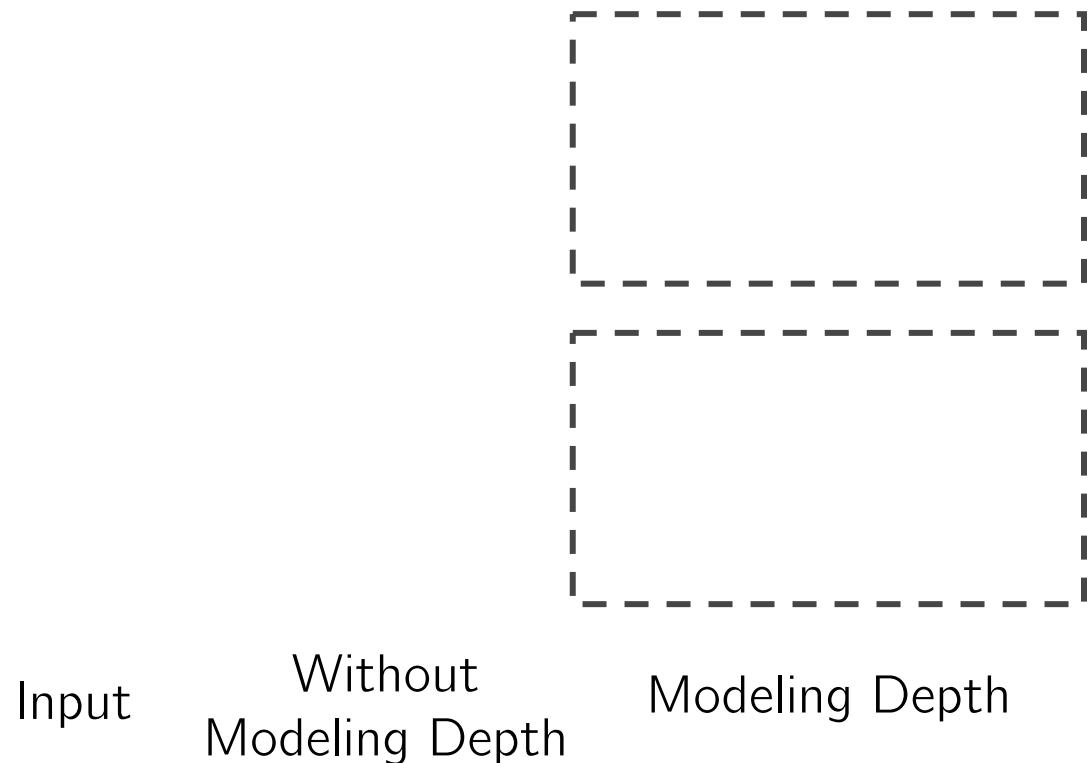


Shape

Trained on:  
ShapeNet [Chang et al., 2016]



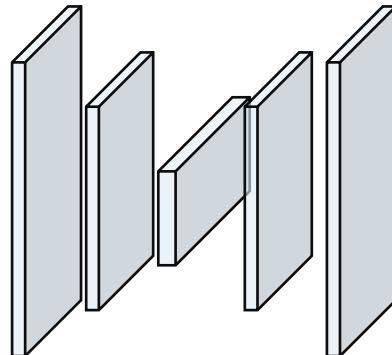
with Yifan Wang



# Generalization to Unseen Classes



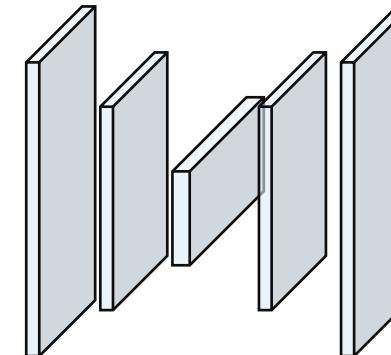
Image



Depth Estimation



Depth



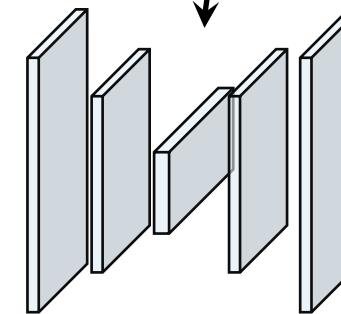
Shape Completion

- 3D-R2N2 [ECCV'16]
- DRC [CVPR'17]
- PSGN [CVPR'17]
- OGN [ICCV'17]
- MarrNet [NeurIPS'17]
- AtlasNet [CVPR'18]
- ShapeHD [ECCV'18]
- Multi-View [ECCV'18]

If we take models trained on chairs, planes, cars



Training Images



Test Image (Table)



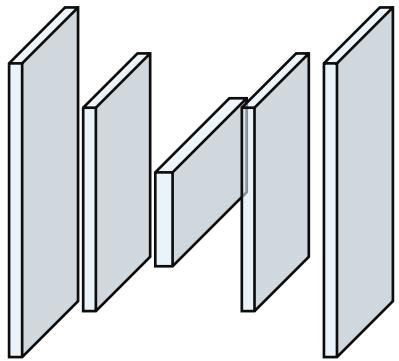
Network

Direct Prediction

# Inverting the Graphics Engine



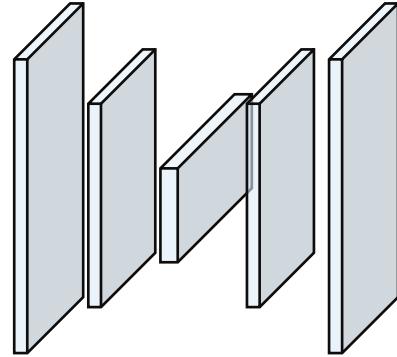
Image



Depth Estimation



Depth



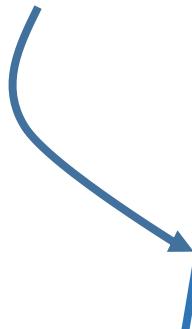
Shape Completion



Shape

Projecting depth into 3D is a deterministic,  
fully differentiable process.

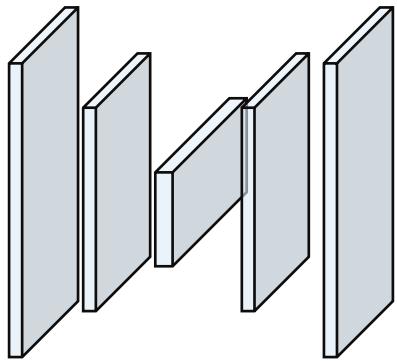
Shape completion network is over-parameterized:  
learning a deterministic mapping.



# Inverting the Graphics Engine



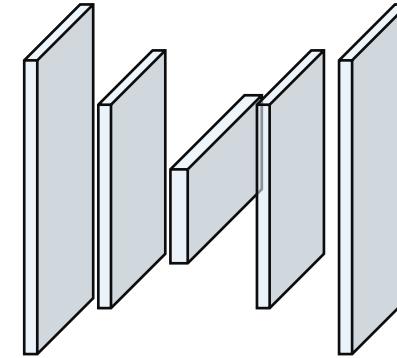
Image



Depth Estimation



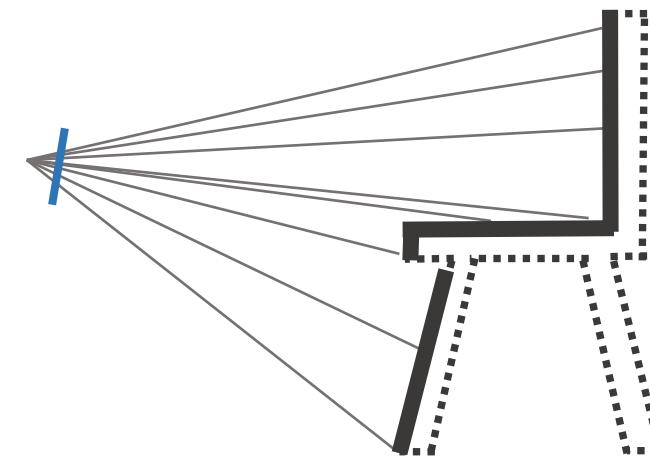
Depth



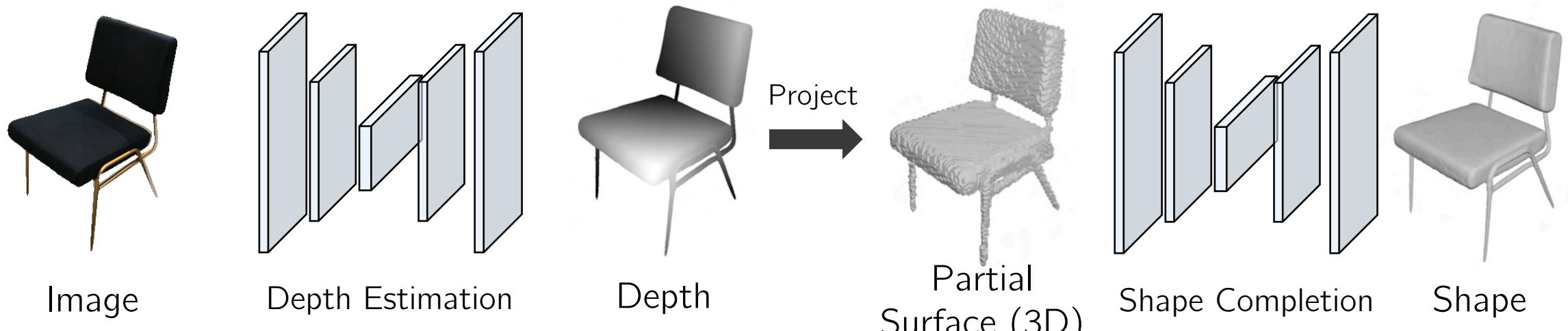
Shape Completion



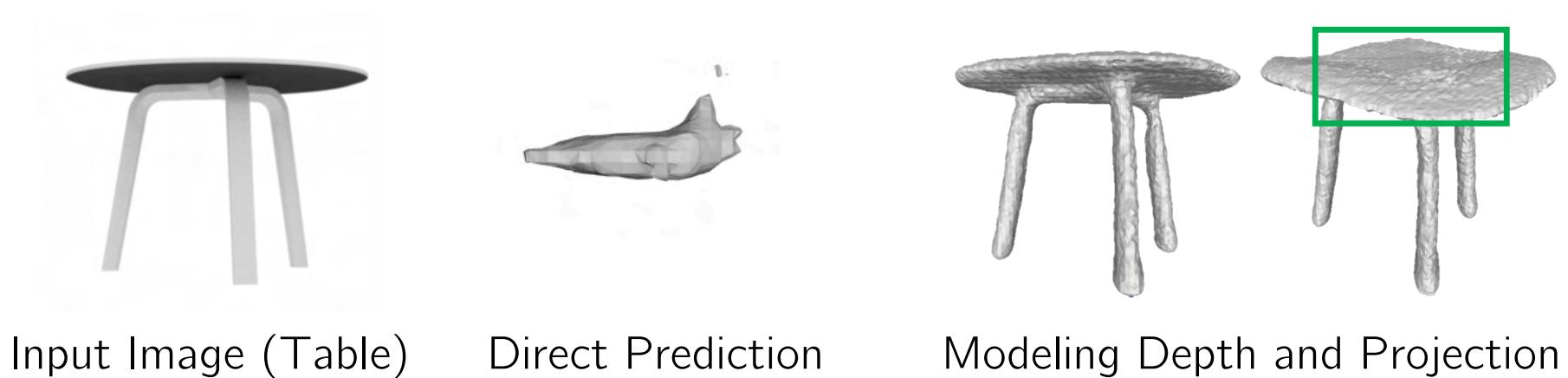
Shape



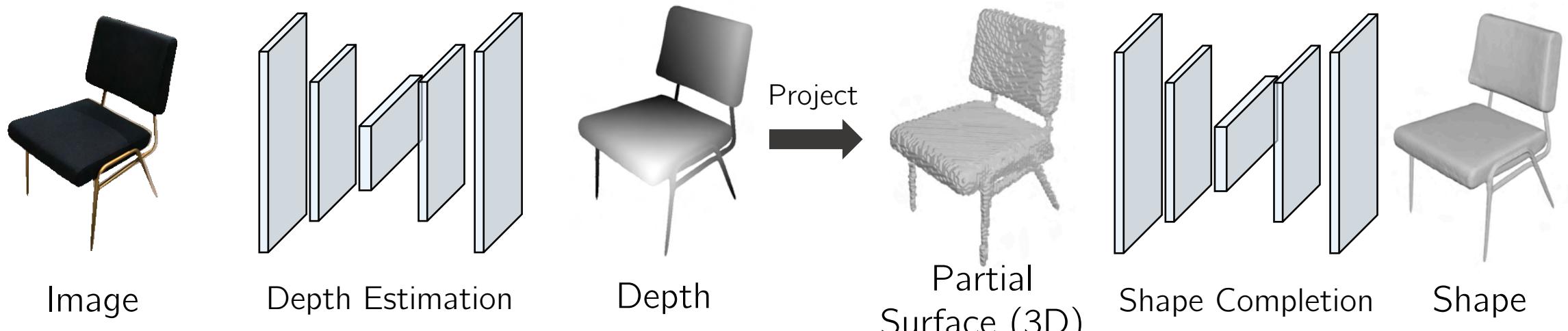
# Inverting the Graphics Engine



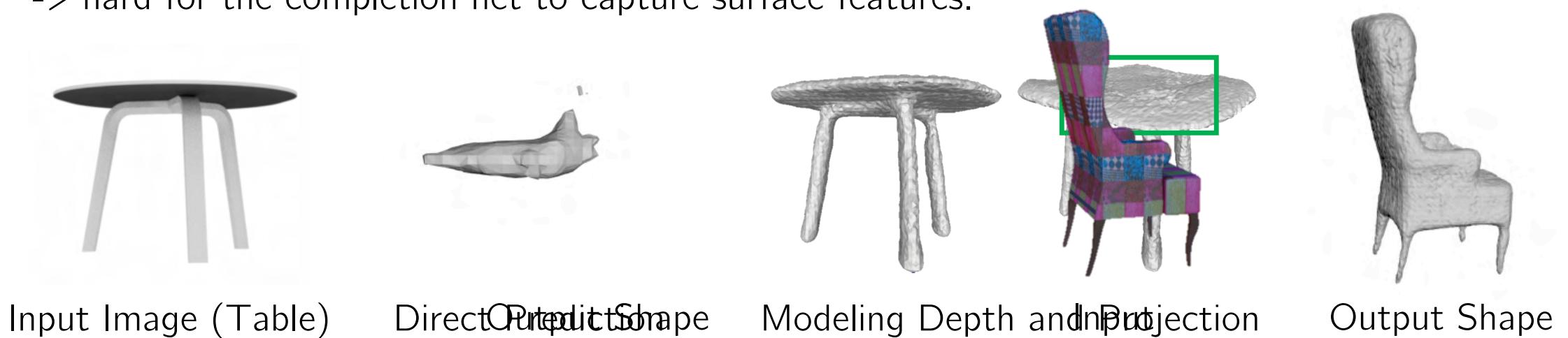
If we take models trained on chairs, planes, cars



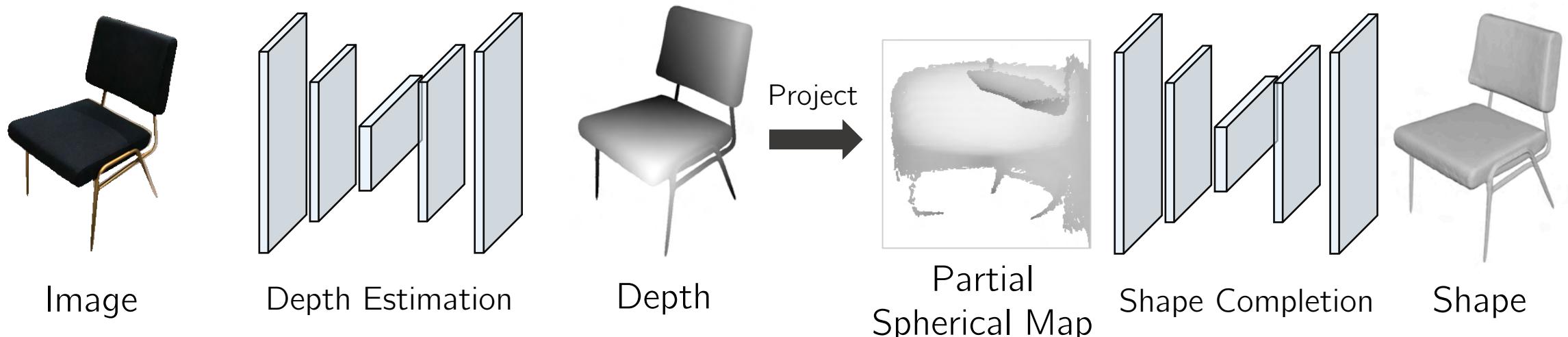
# Inverting the Graphics Engine



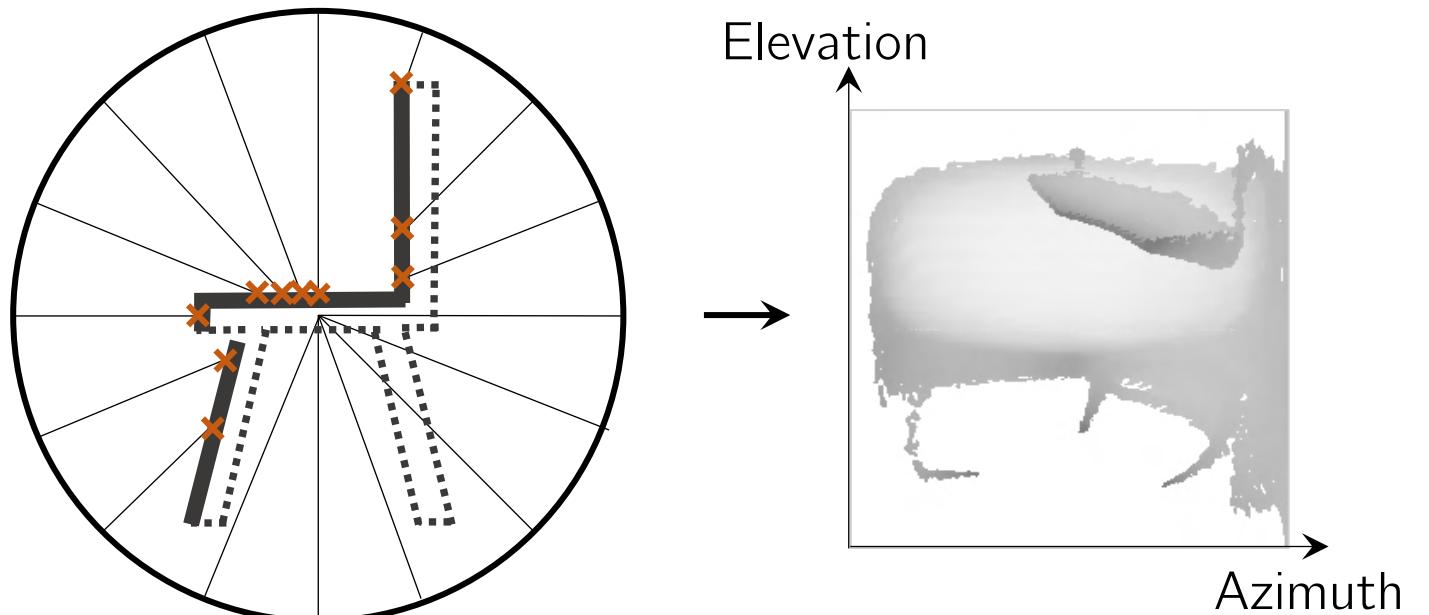
Partial surface in 3D is very sparse  
-> hard for the completion net to capture surface features.



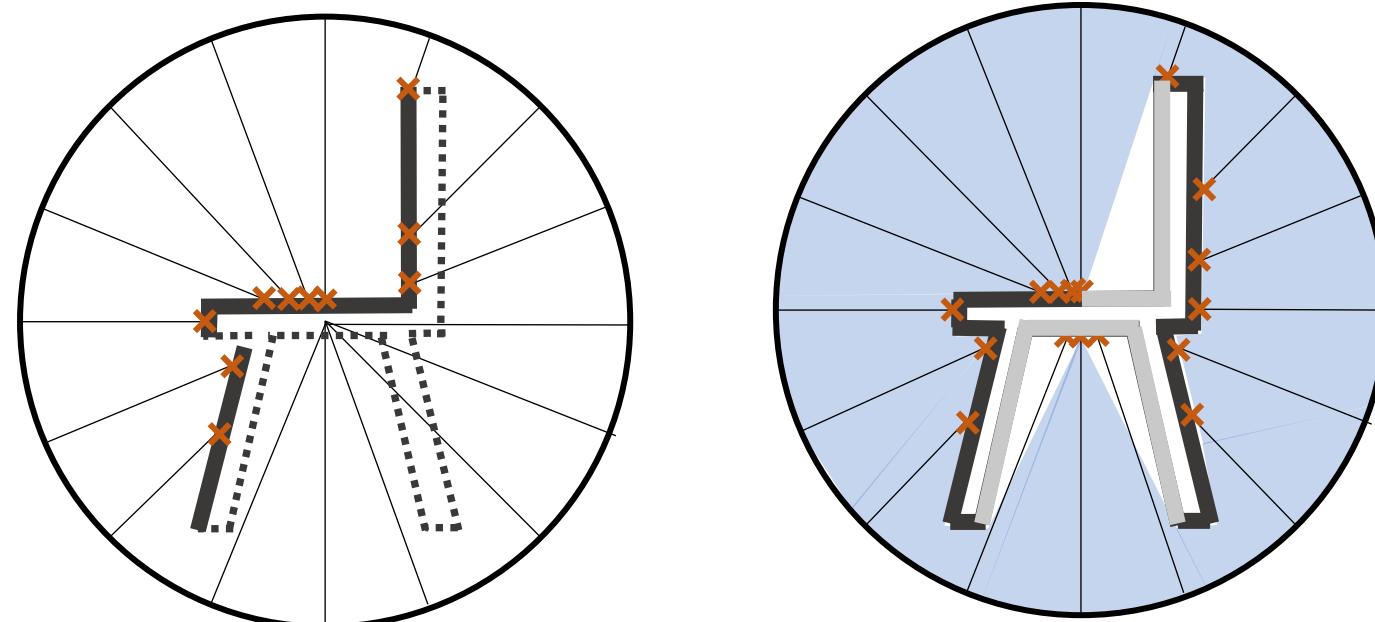
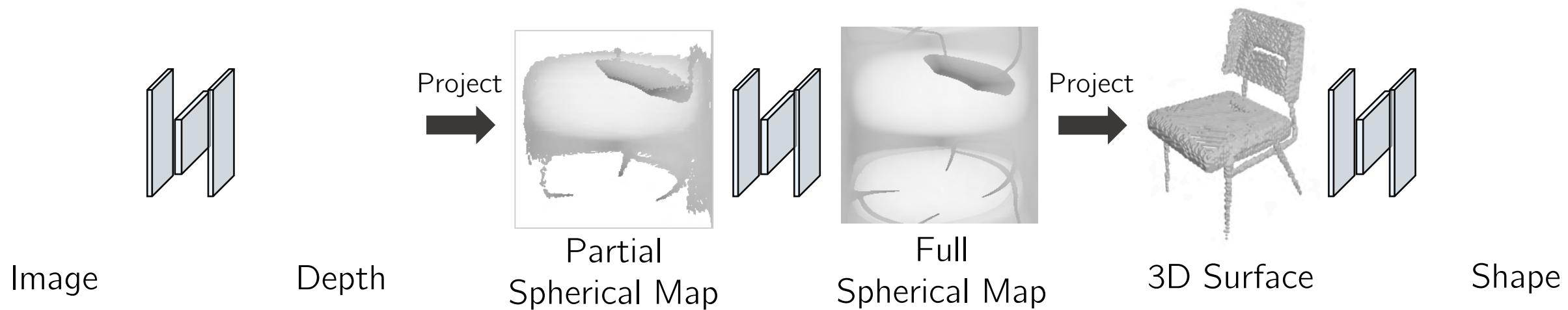
# Inverting the Graphics Engine



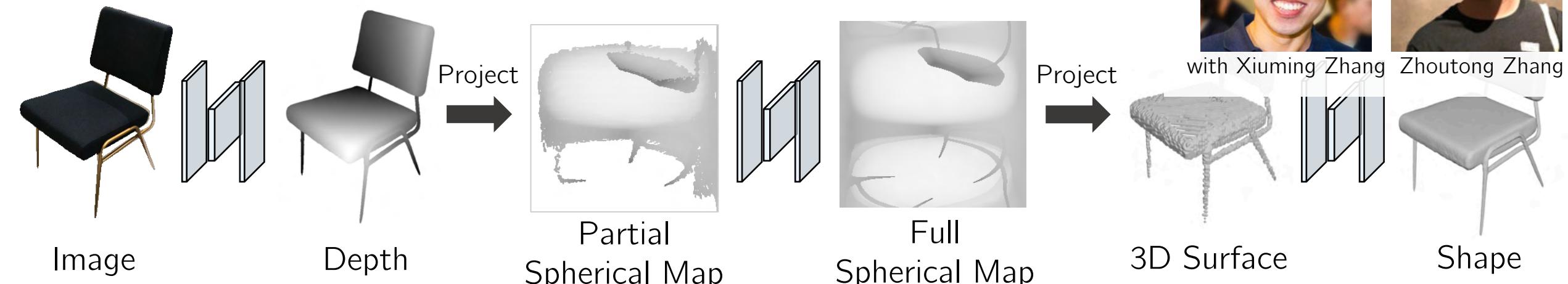
Spherical map as a surrogate representation for surfaces in 3D



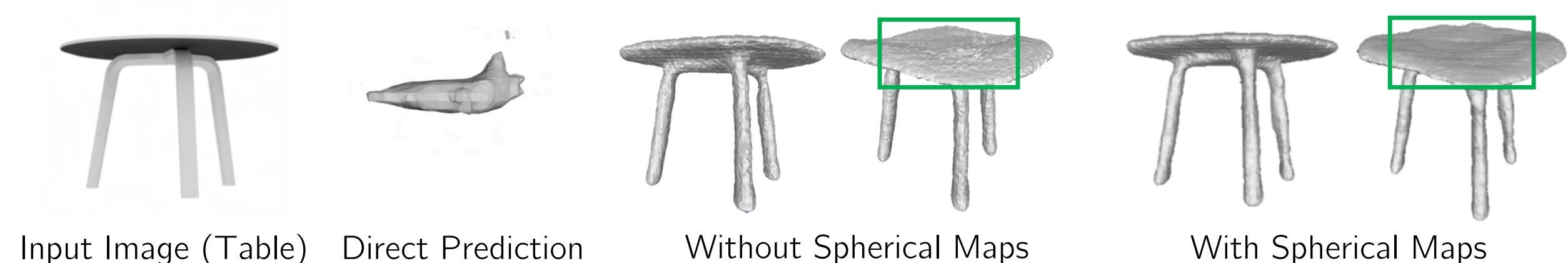
# Inverting the Graphics Engine



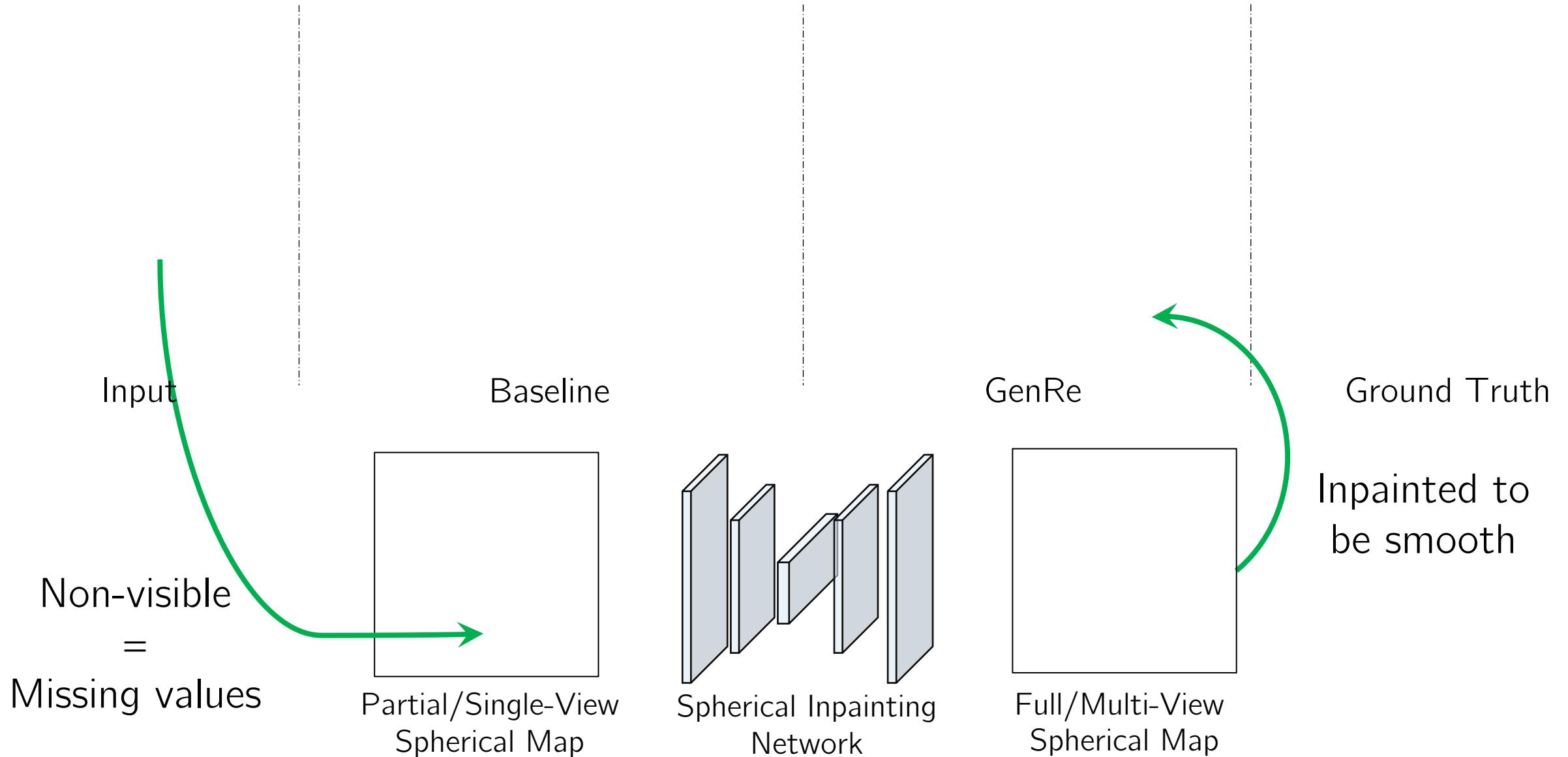
# Invertible Graphistr Engine (GenRe)



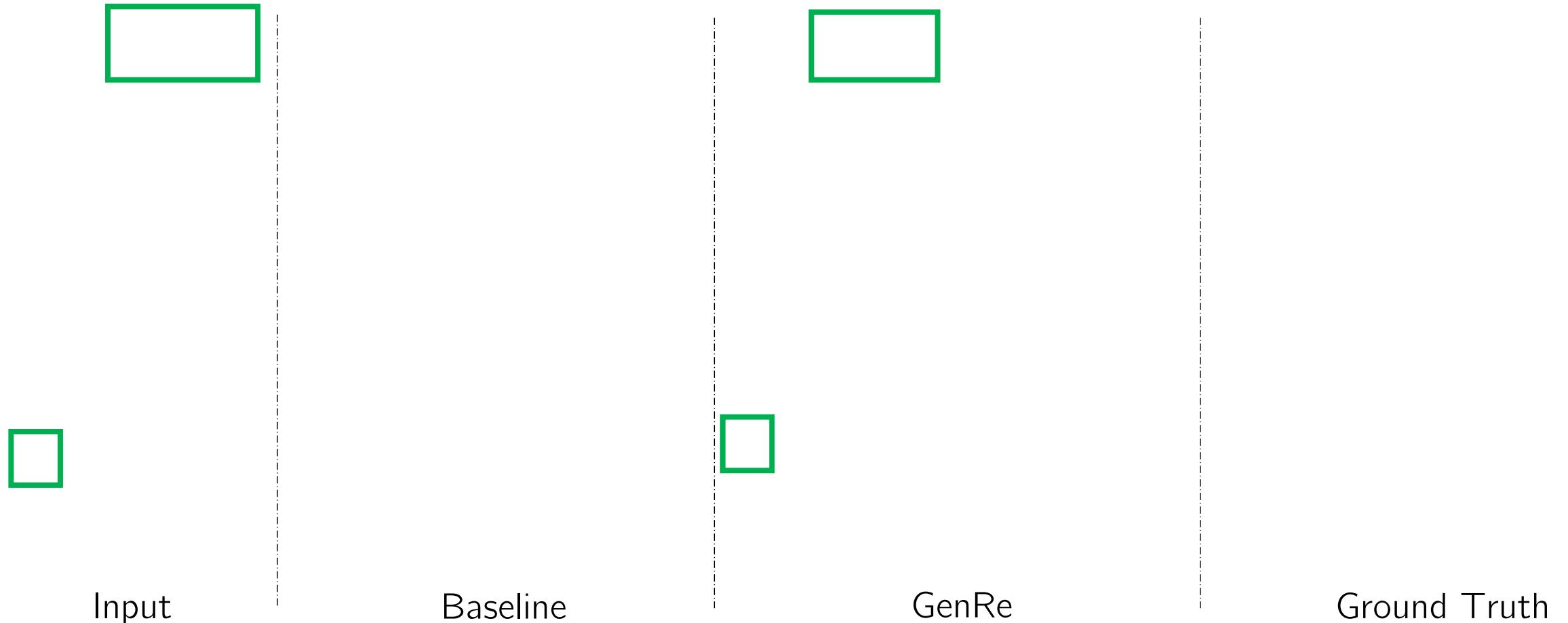
If we take models trained on chairs, planes, cars



# Results: Generalizing to Unseen Classes



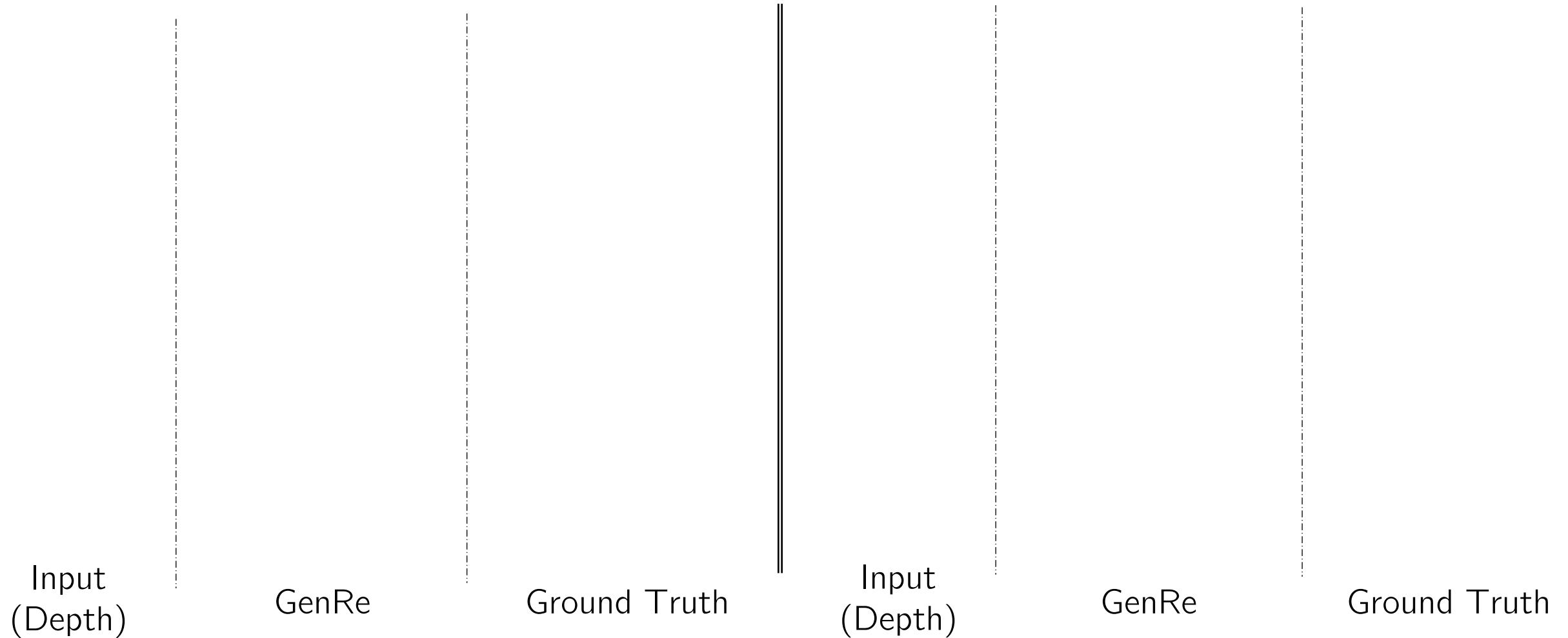
# Results: Generalizing to Unseen Classes



**Training:** cars, chairs, airplanes

**Testing:** beds, tables, benches

# Results: Generalizing to Non-Rigid Classes

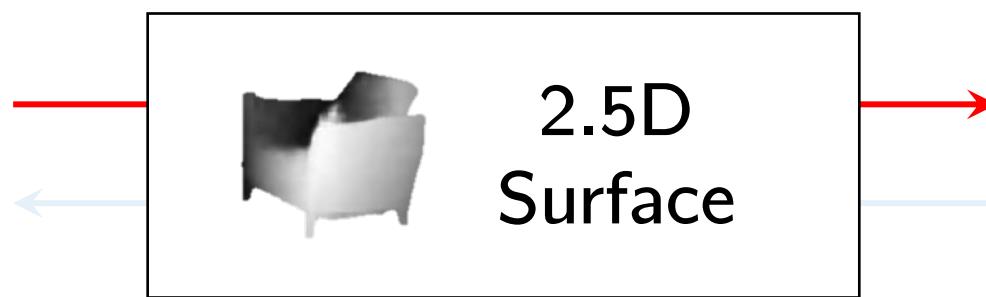


Training: cars, chairs, airplanes

Testing: humans, horses



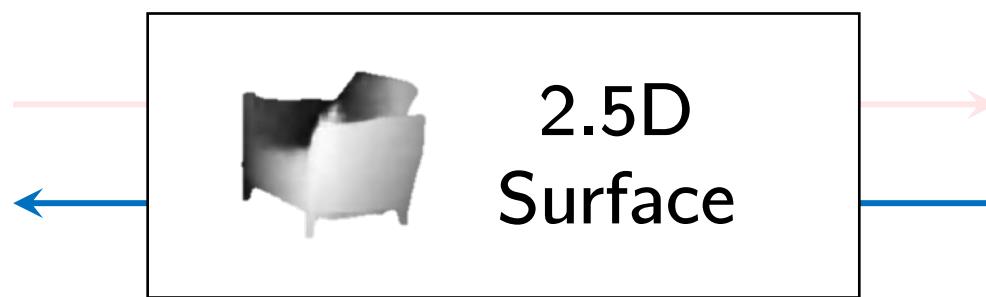
2D Image



3D Shape

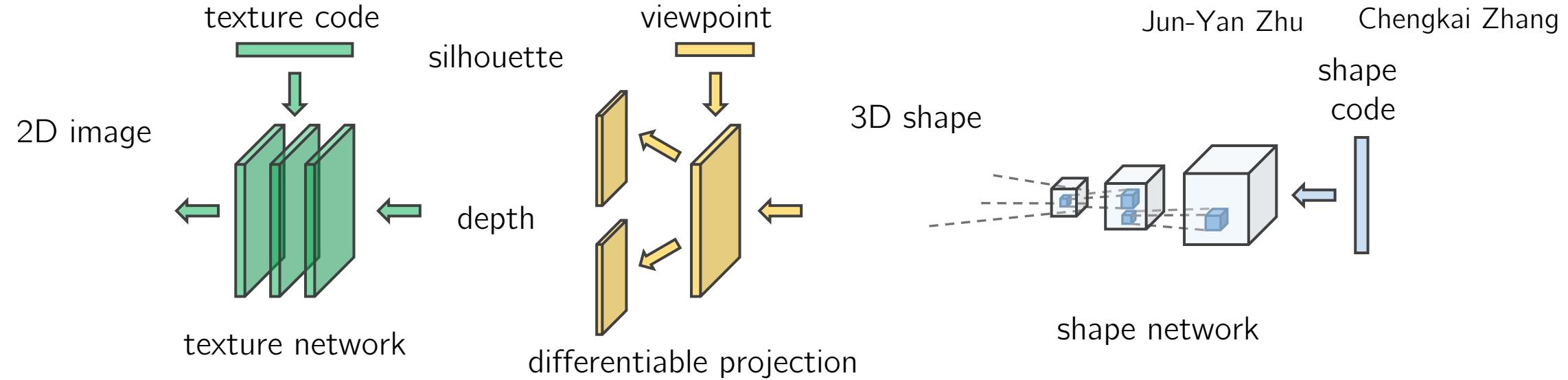


2D Image



3D Shape

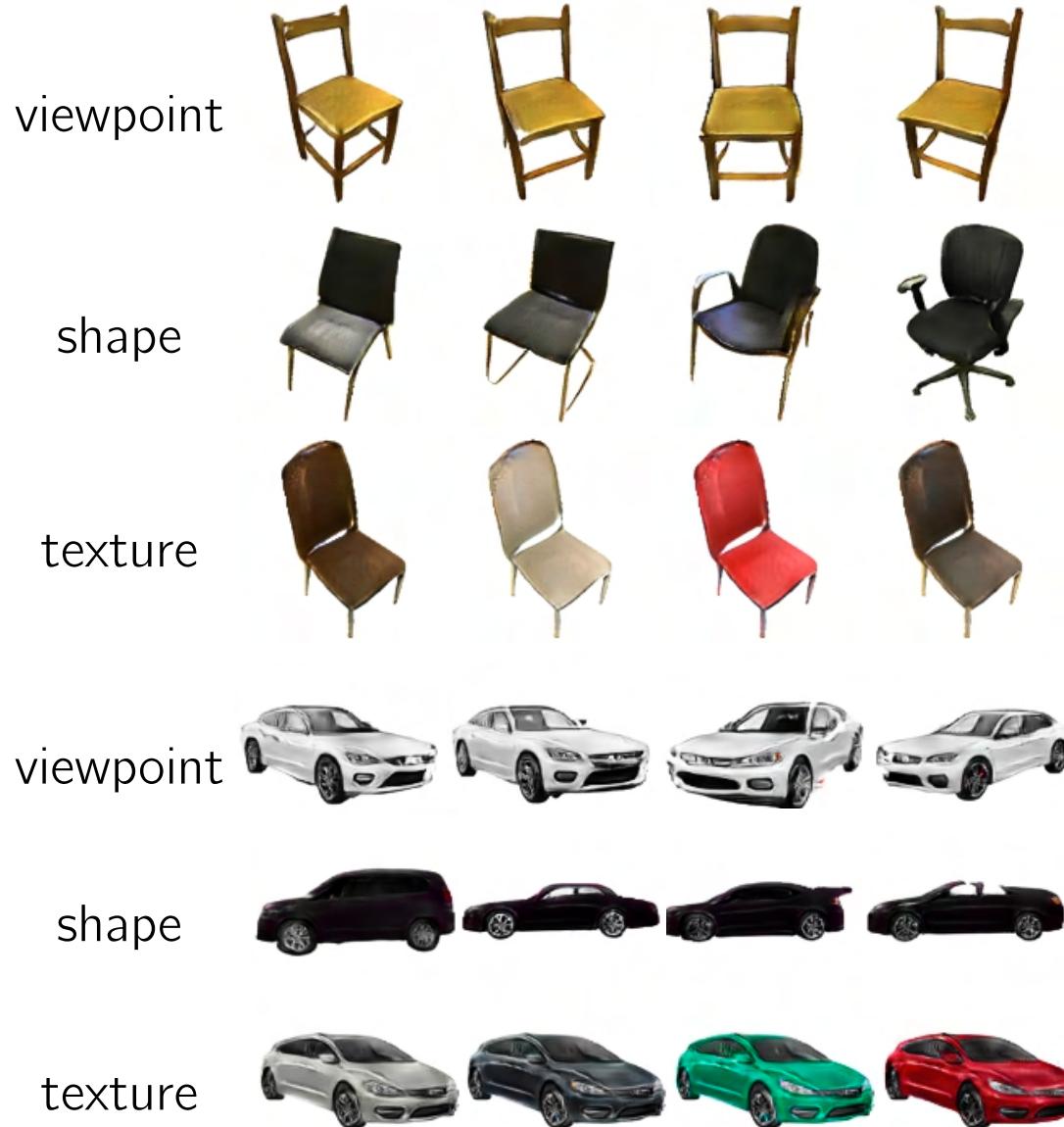
# Shape and Texture Synthesis



2D

3D

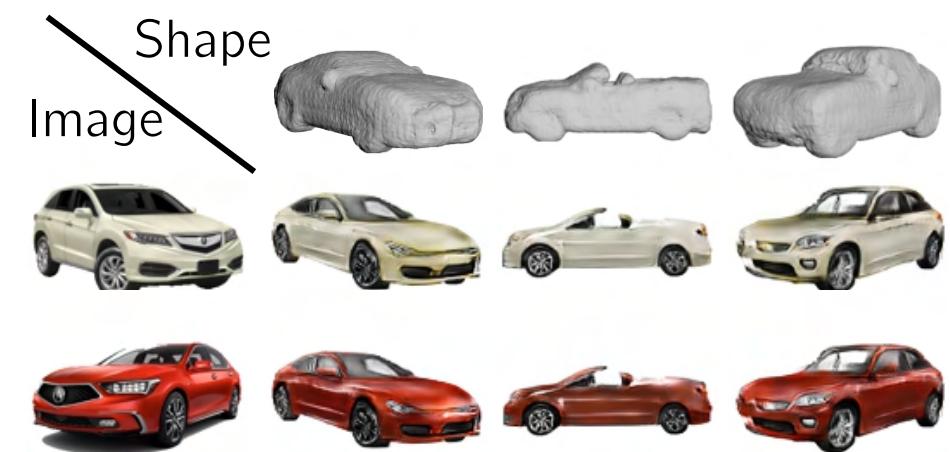
Editing viewpoint, shape, and texture



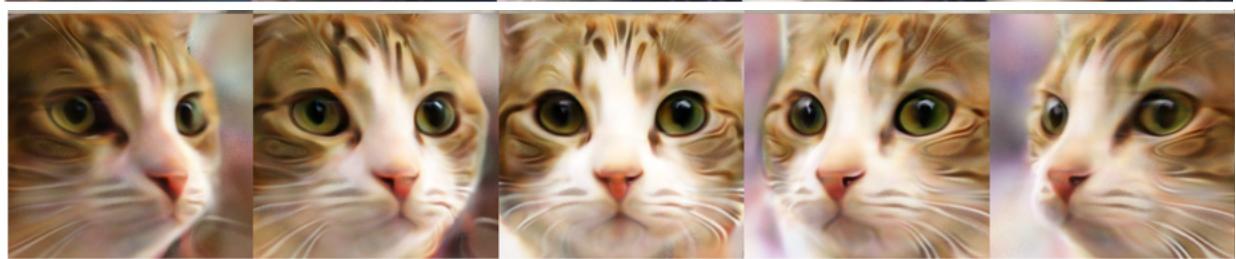
Interpolation in the latent space



Transferring shape and texture

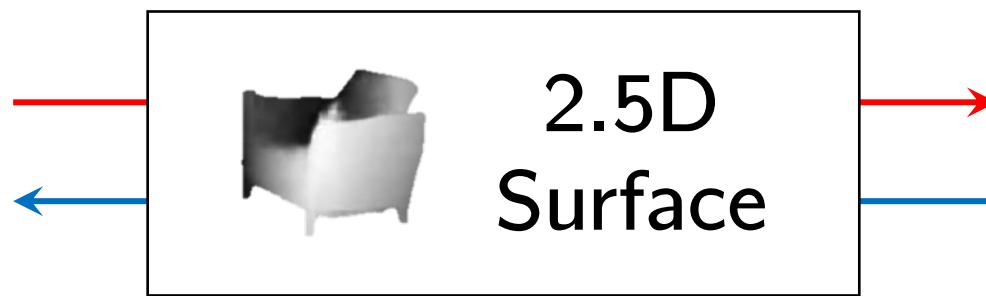


# $\pi$ -GAN: Implicit Representations





2D Image



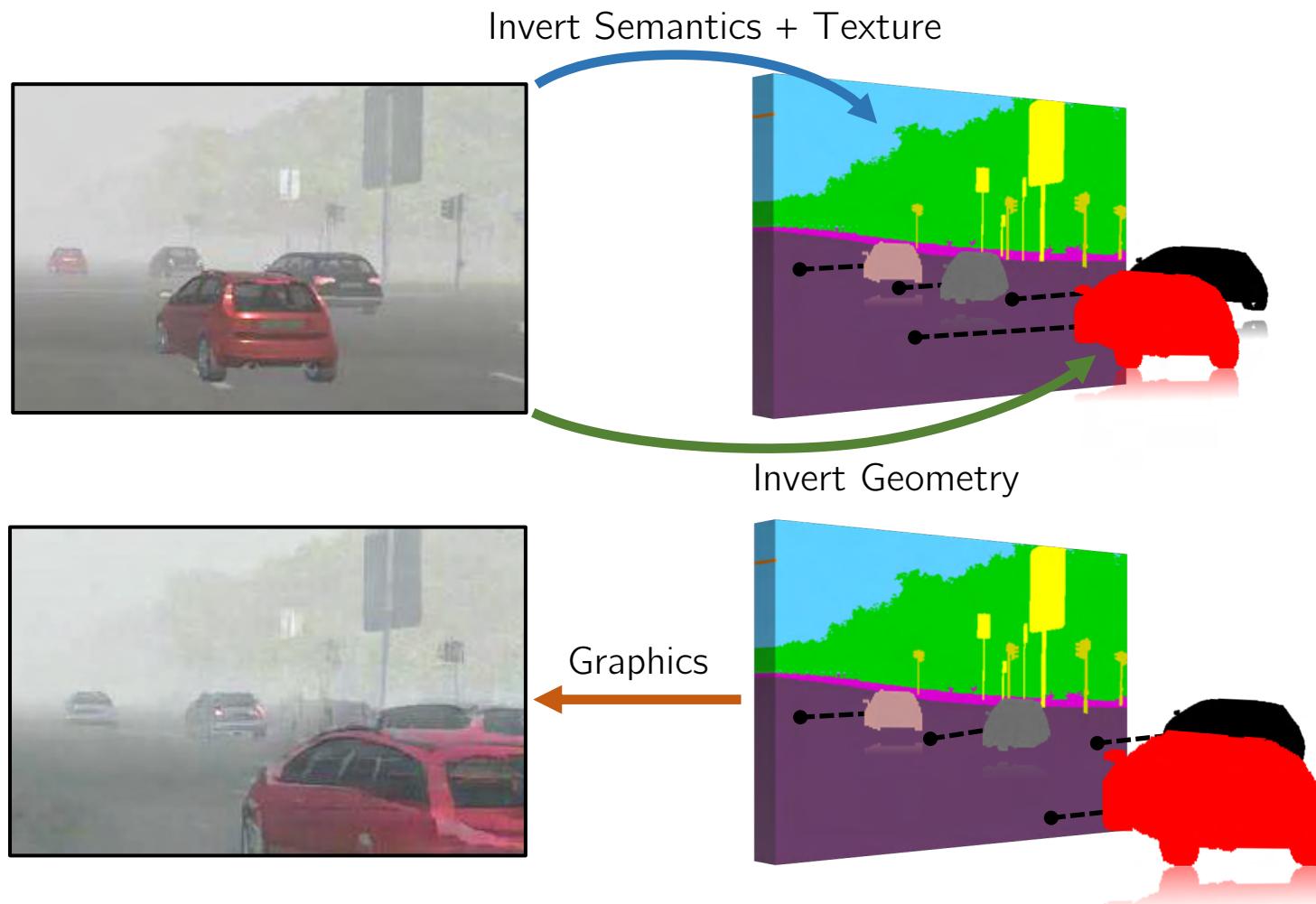
3D Shape

# Extension to Scenes



Shunyu Yao

Harry Hsu



```
<seg sky code=[0.53, -0.60, -0.28...]>  
<seg tree code=[0.61, -0.64, -0.22...]>  
...
```

```
<obj type=car center forward black>  
<obj type=car center forward red>
```

```
<seg sky code=[0.5, -0.60, -0.28...]>  
<seg tree code=[0.61, -0.64, -0.22...]>  
...
```

```
<obj type=car right forward black>  
<obj type=car right forward red>
```

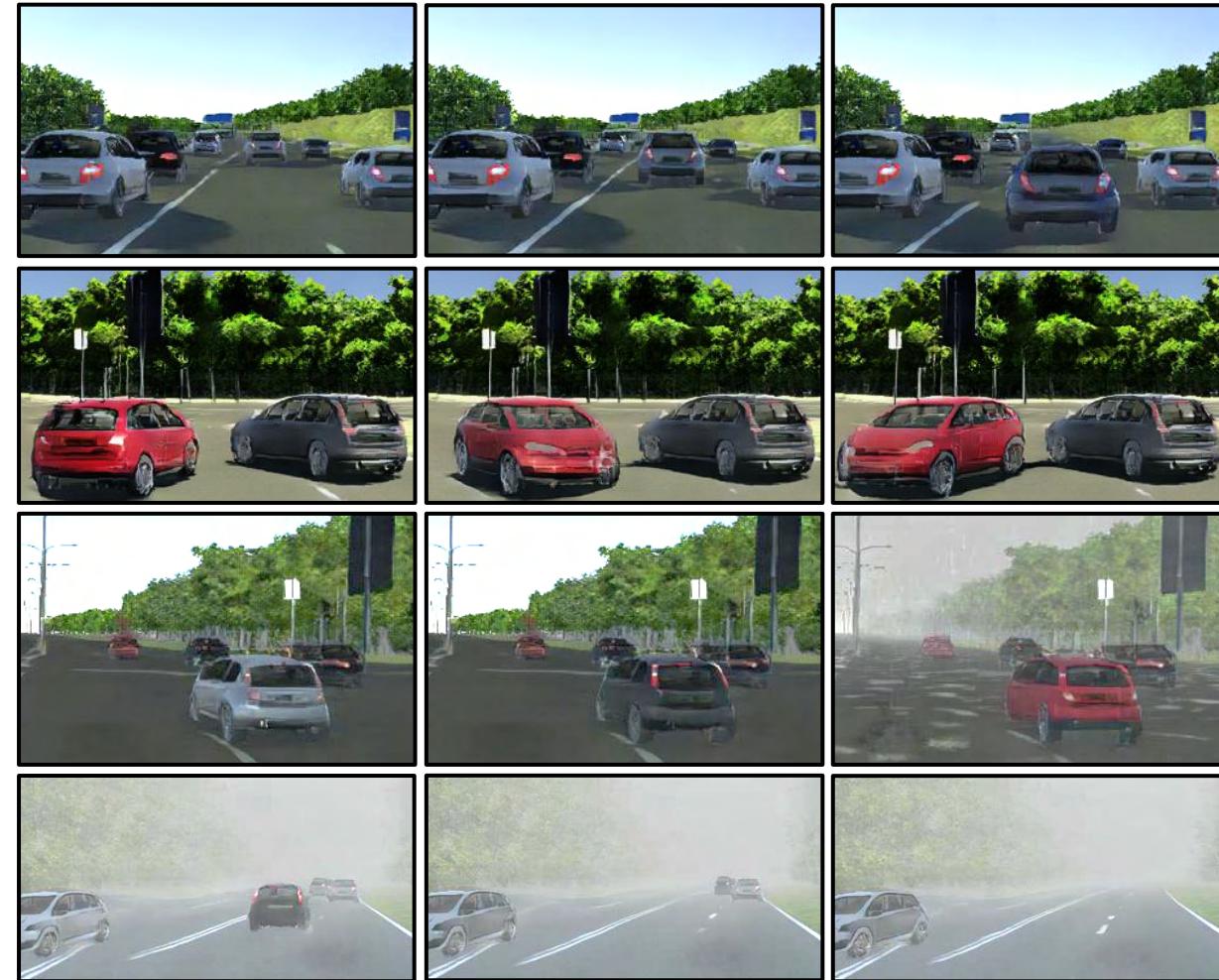
Manipulate

# Image Editing on Virtual KITTI

Original images



Edited images



# Image Editing on CityScapes (Real Images)

Original images

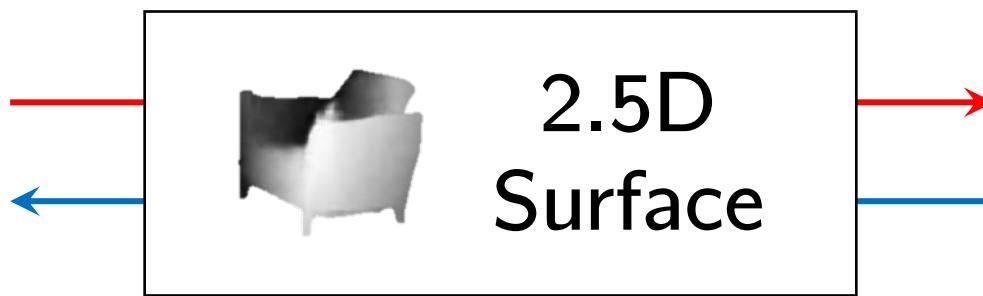


Edited images





2D Image



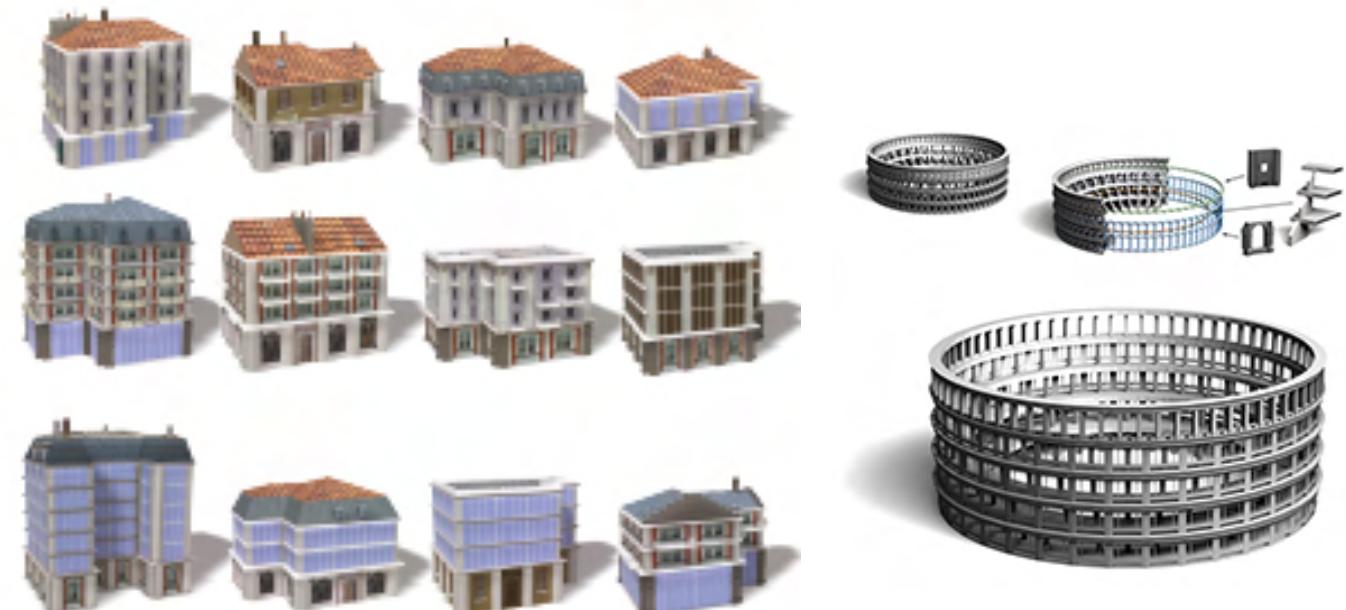
3D Shape

# L-Systems

# PCFG, Symmetry

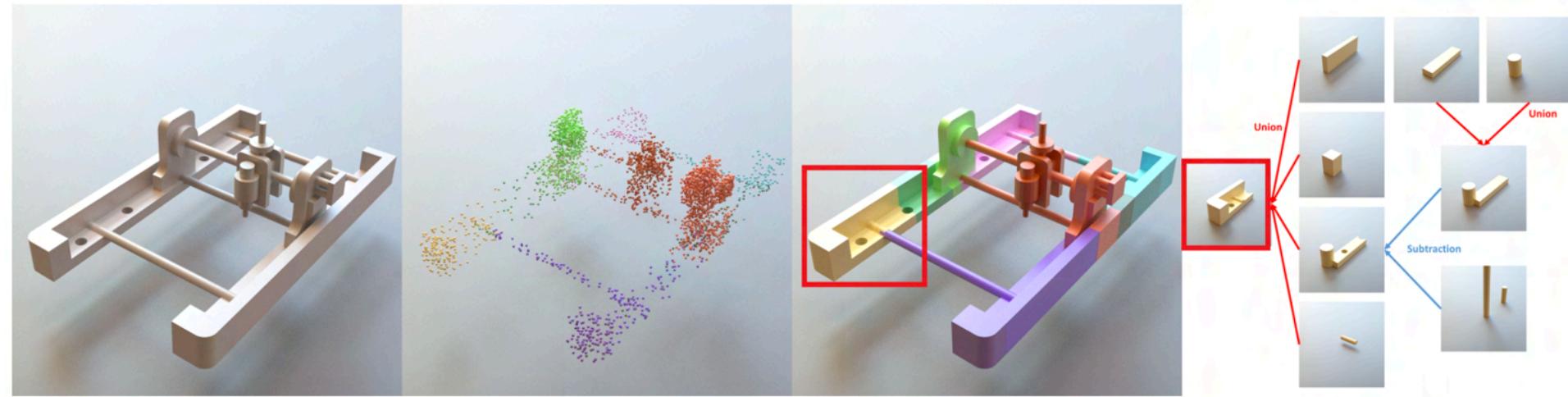
start: X

rules:  $(X \rightarrow F+[[X]-X]-F[-FX]+X)$   
 $(F \rightarrow FF)$



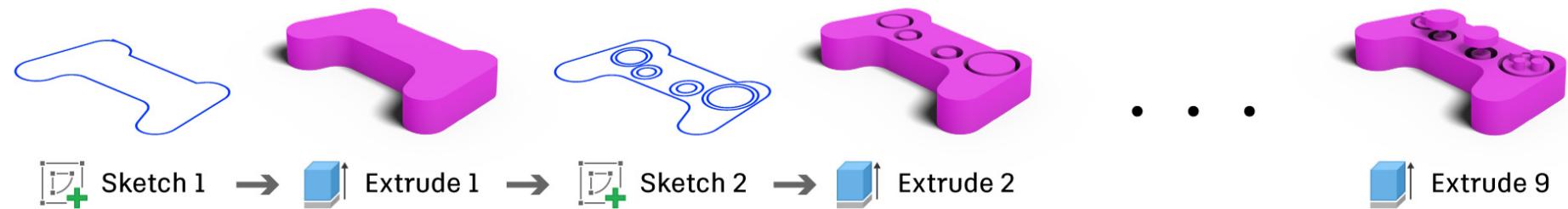
## InverseCSG

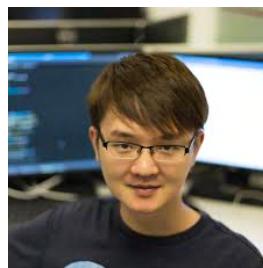
Du et al. SIGAsia 2018



## Fusion 360 Gallery

Willis et al. arXiv 2020





Yonglong Tian

# From 3D Reconstruction to Abstraction



Input Image  
(Table)



Direct  
Prediction



GenRe Without  
Spherical Maps



GenRe



GenRe +  
Shape Programs

Shapes often have abstract, program-like structure.



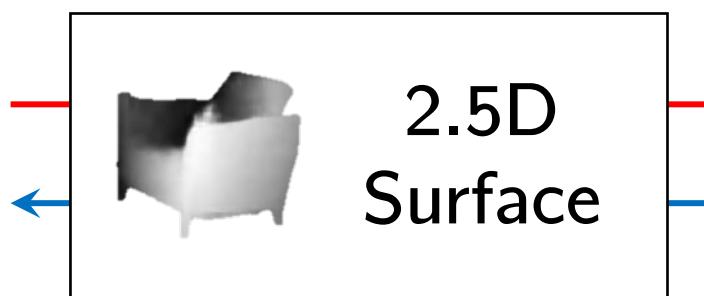
```
Draw("Top", "Circle")
For(i < 2, "trans")
  For(j < 2, "trans")
    Draw("Cuboid")
  For(i < 2, "trans")
    Draw("Rect")
```

Shape domain-specific language (Shape DSL)

Program	→ Statement; Program
Statement	→ Draw(Semantics, Shape, Position_Parms, Geometry_Parms)
Statement	→ For(For_Parms); Program; EndFor
Semantics	→ semantics 1   semantics 2   semantics 3   ...
Shape	→ Cuboid   Cylinder   Rectangle   Circle   Line   ...
Position_Parms	→ (x, y, z)
Geometry_Parms	→ (g <sub>1</sub> , g <sub>2</sub> , g <sub>3</sub> , g <sub>4</sub> , ...)
For_Parms	→ Translation_Parms   Rotation_Parms
Translation_Parms	→ (times i, orientation u)
Rotation_Parms	→ (times i, angle θ, axis a)



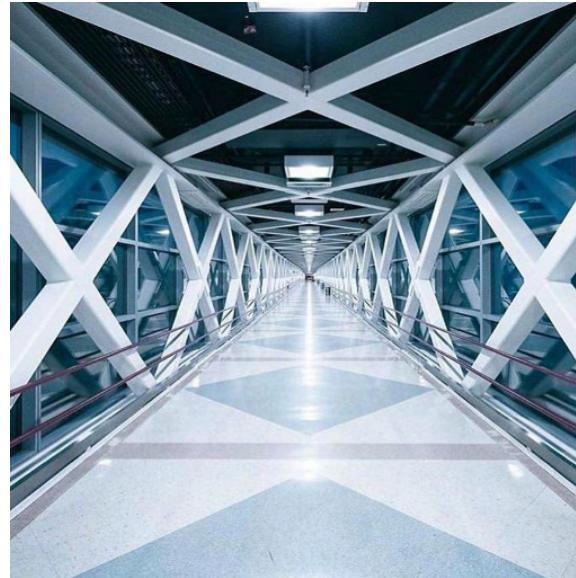
2D Image



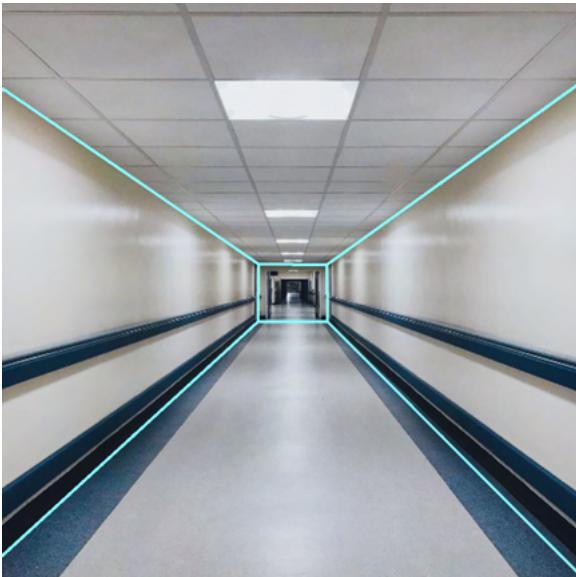
3D Shape

`draw(Back, Cuboid)`  
`for i in range(0, 1):`  
`draw(Side, Cuboid, i)`  
`draw(Bottom, Cuboid)`

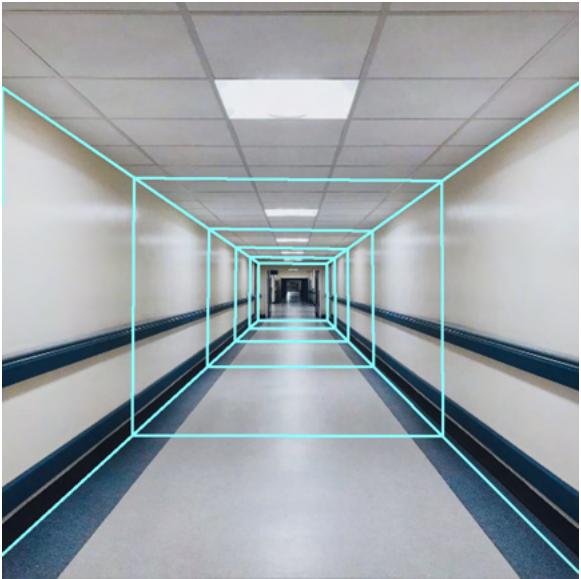
# Rich Structure in 3D Scenes



# Planes/Surfaces, Symmetry

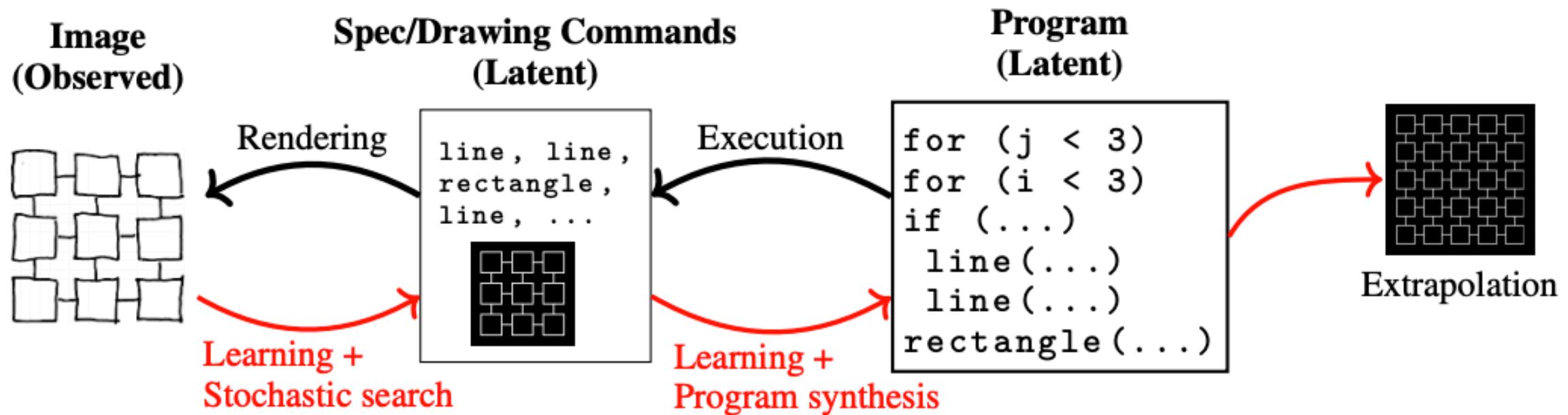


# Repetition



How can we leverage such structural information for 3D scene understanding and editing?

# Program Synthesis for Visual Data



(sphere, small, metal,  
green, x=2, y=2, z=0) ...

(a) Input Image

(b) Object & Group Detection

```
for i in range(0, 5):
    for j in range(0, 3 - i):
        sphere(pos = (1 + i, 1 + j, 0),
               color = 6 - j)
for i in range(0, 4):
    cylinder(pos = (3, 3, i),
              color = 7)
cylinder(pos = (4, 2, 0),
          color = 3)
```

(c) Program

Input Image

Patches

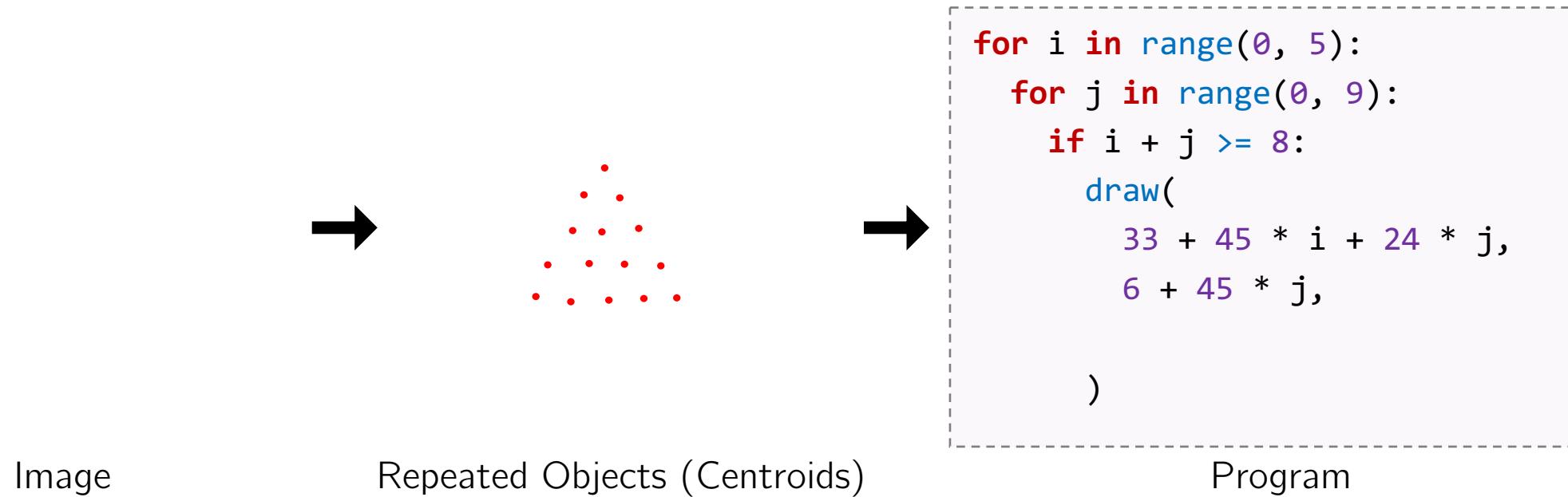
Edited Image

Input Images

Patches

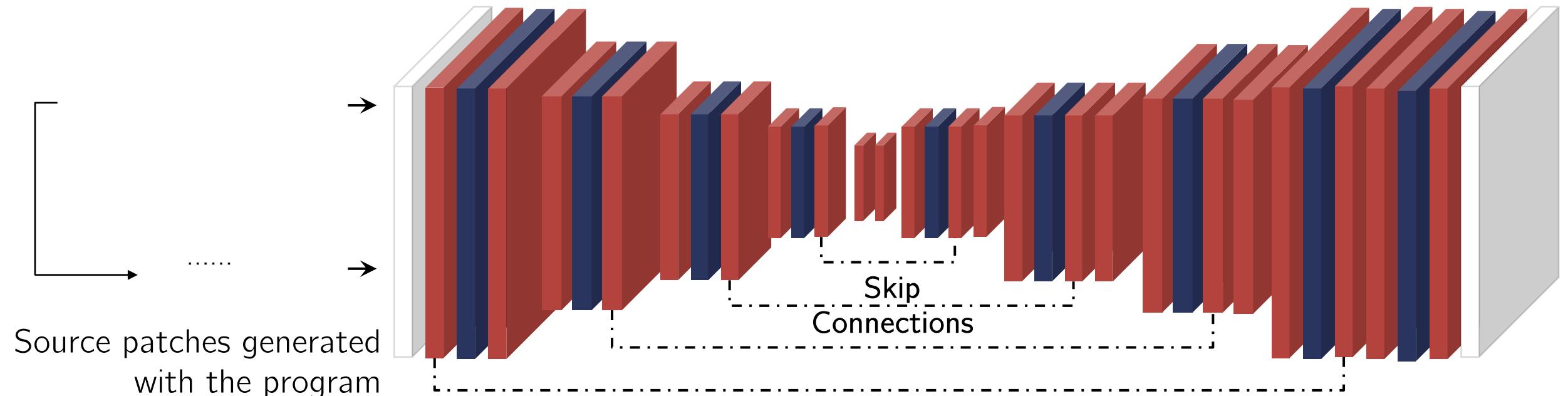
Edited Image

# Generalizing to Natural Images



Key Idea: Internal learning (single-image learning)

[Internal Statistics] Zontak and Irani. CVPR 2011  
[Repeated Pattern Detection] Lettry et al. WACV 2017



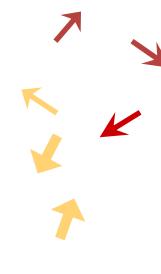
Input (Corrupted) Output (Inpainted)

A. Inpainting

Input (Partial) Output (Extrapolated)

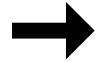
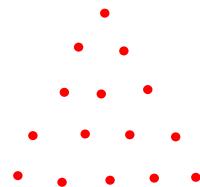
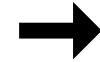
B. Extrapolation

Input (Regular) Output (Irregular)



C. Regularity Editing

Image



Repeated Objects (Centroids)

Corrupted

Inpainted

Partial

Extrapolated

Original

Irregular

A. Inpainting

B. Extrapolation

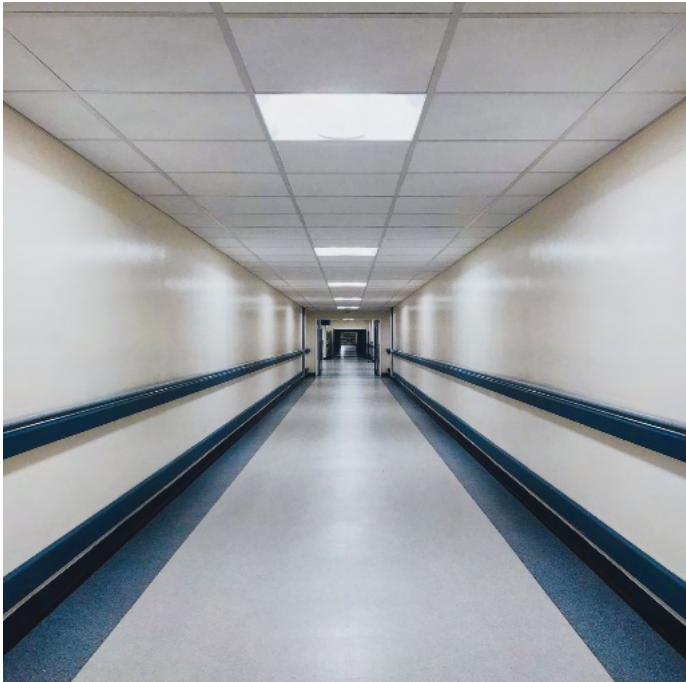
C. Regularity Editing

```
! "# i $% range(0, 5):  
! "# j $% range(0, 9):  
$! i + j >= 8:  
draw(  
    33 + 45 * i + 24 * j,  
    6 + 45 * j,  
)
```

Program



# Generalizing to Multiple Planes



Input Image

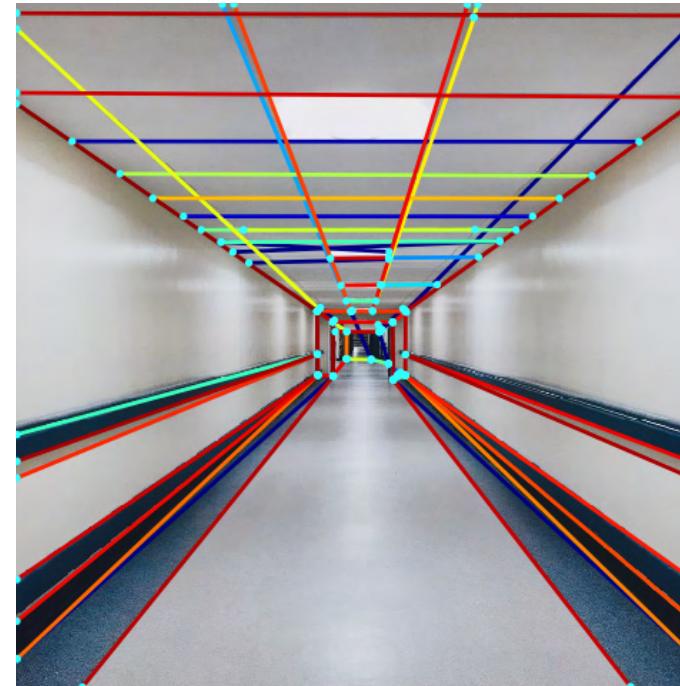
Program

```
: (; C&?(#&)*+/DE-
6/78&%(/9
: (; 78&%( )<"=>9+/%" #?&8>9-
! "# $ $% #&%' ()*+ , - . /
0#&1)//////////+
2 3 4, 5/$-
6/78&%(/@
: (; 78&%( )<"=>@+/%" #?&8>@-
! "# $ $% #&%' ()*+ , - . /
0#&1)//////////+
2 3 4, 5/$-
6/78&%(/A/BBB
```

# Visual Cue Extraction



Vanishing Point  
(NeurVPS [Zhou et al. 2019])



Wireframes  
(L-CNN [Zhou et al. 2019])

Bottom-up visual cues for program synthesis:

Lines and shapes

for line drawings

Repeated objects

for single-plane images

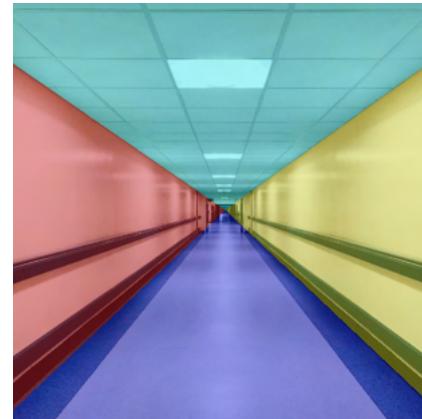
Vanishing point and wireframes

for multi-plane images

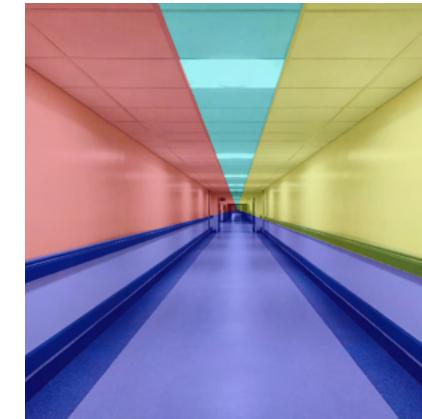
# Candidate Plane Partition Generation



Candidate 1

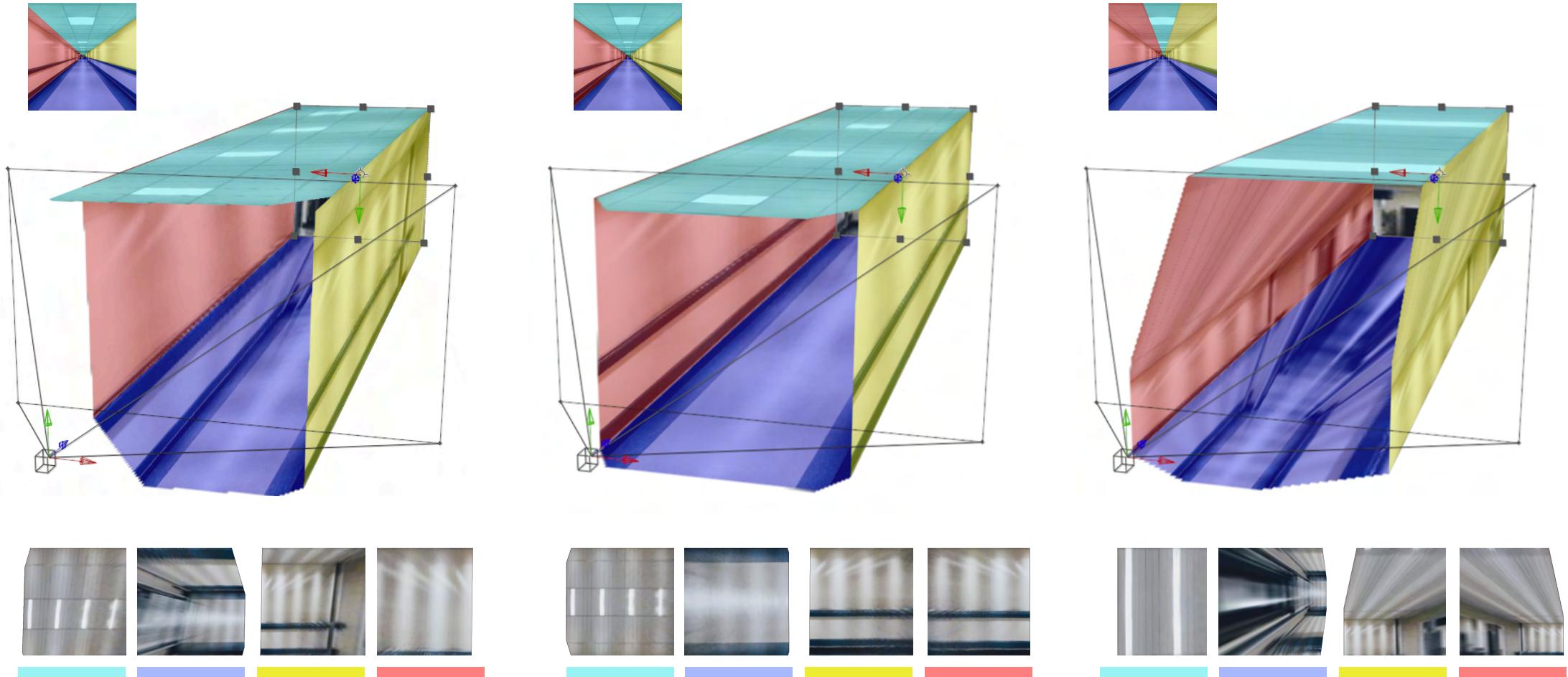


Candidate 2



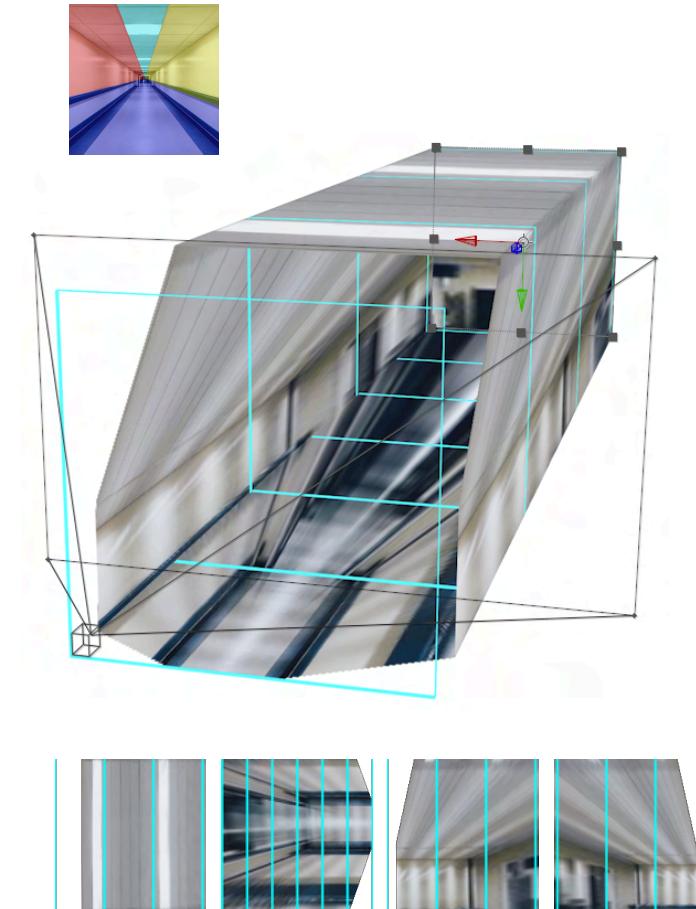
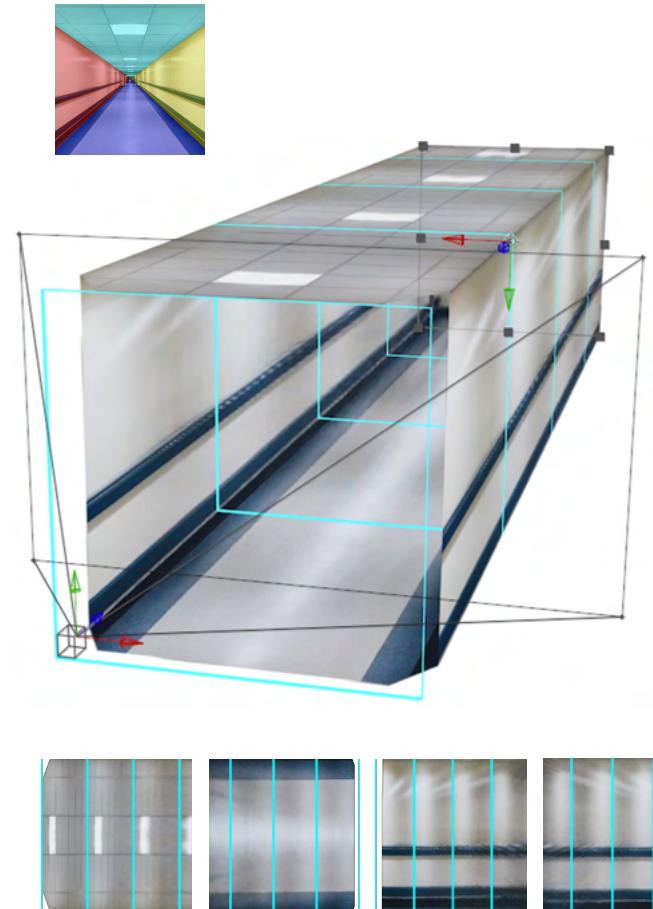
Candidate 3

# Plane Rectification



# Visual Program Synthesis

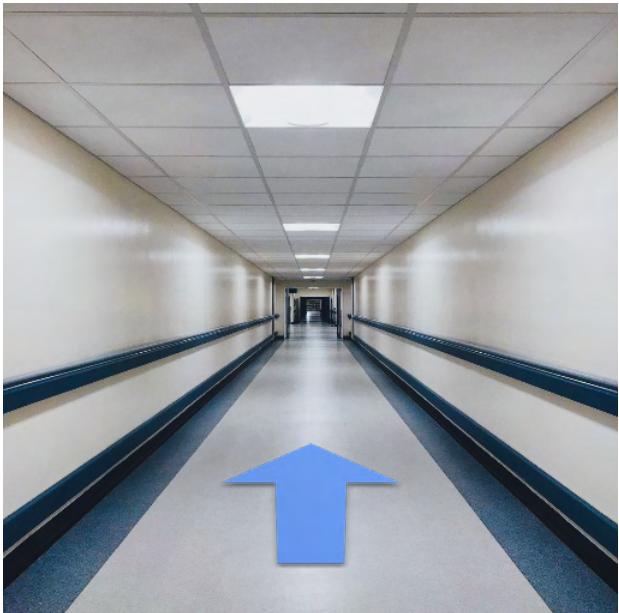
Plane partition and program synthesis help each other.



# Visual Program Synthesis



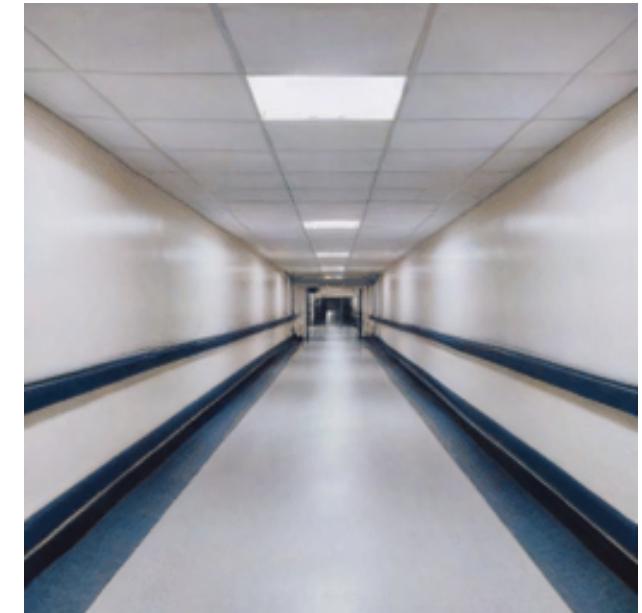
# View Synthesis



Input Image

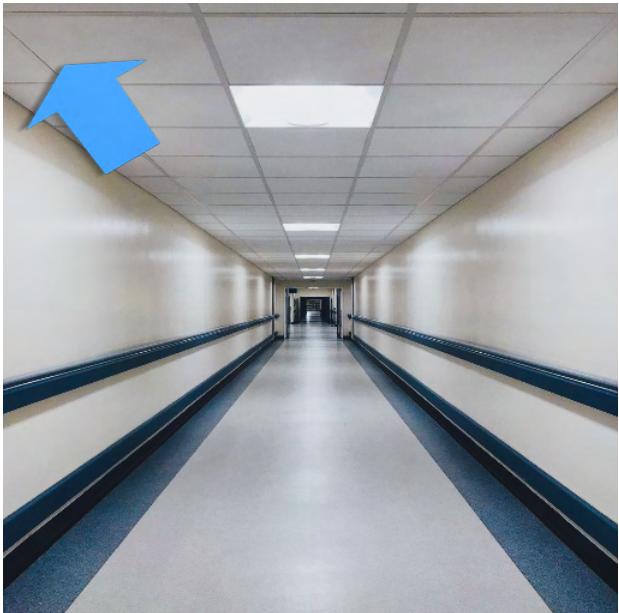


Ours

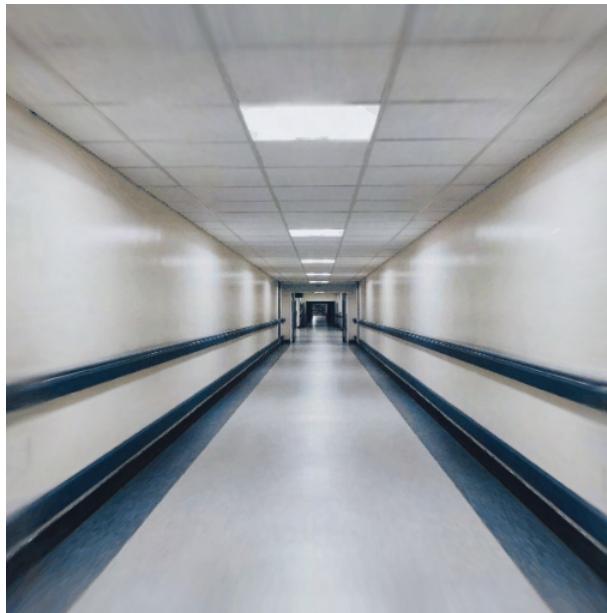


SynSin [Wiles et al. 2020]

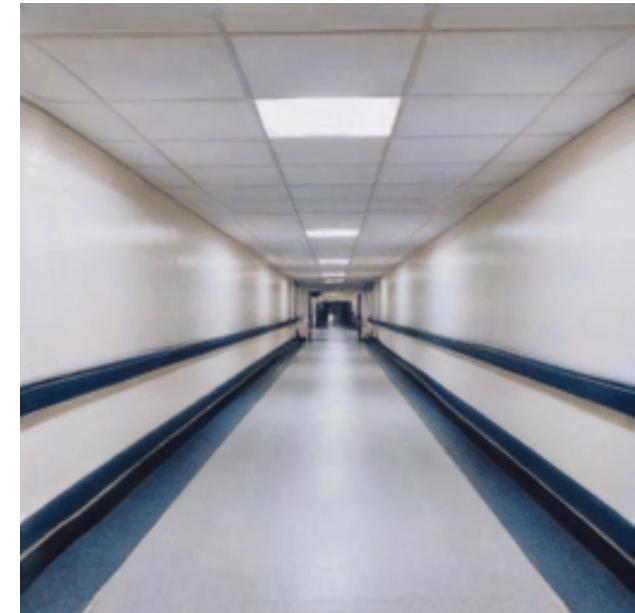
# View Synthesis



Input Image

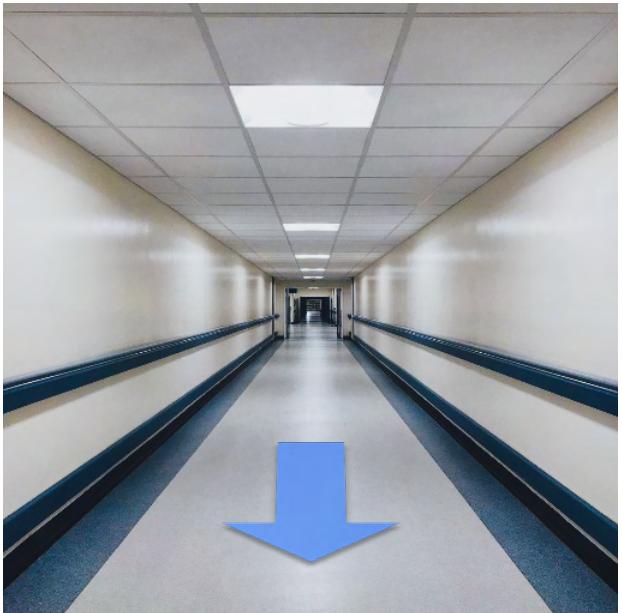


Ours

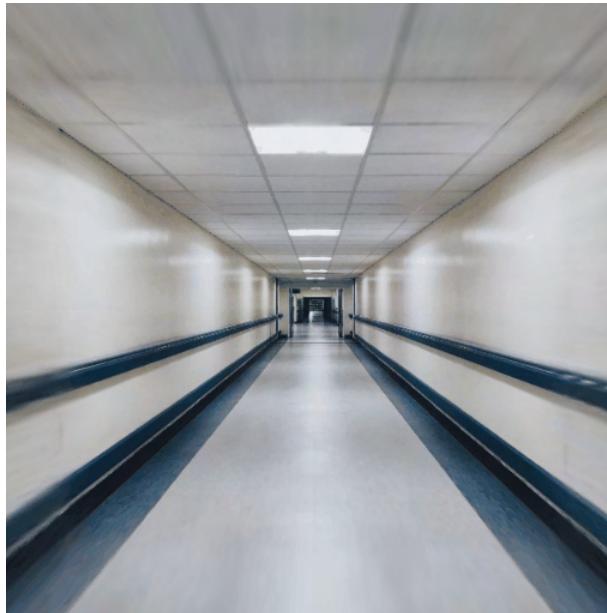


SynSin [Wiles et al. 2020]

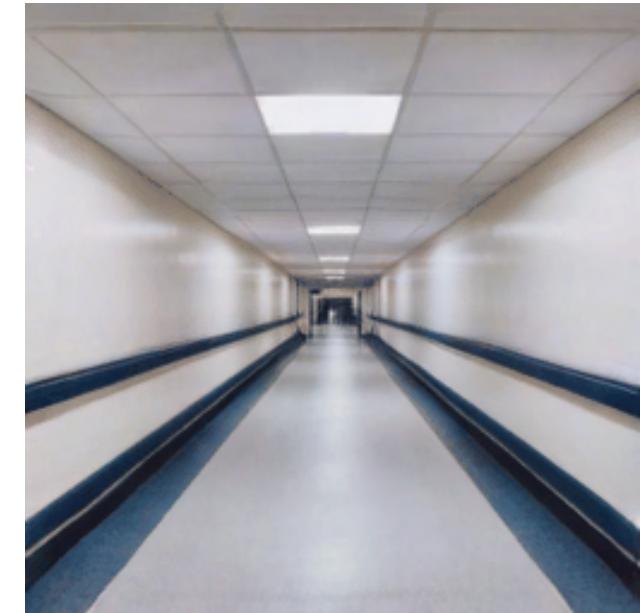
# View Synthesis



Input Image

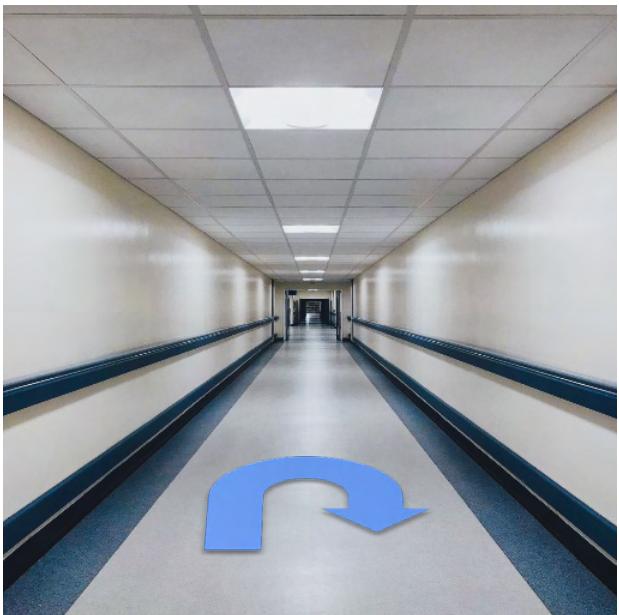


Ours

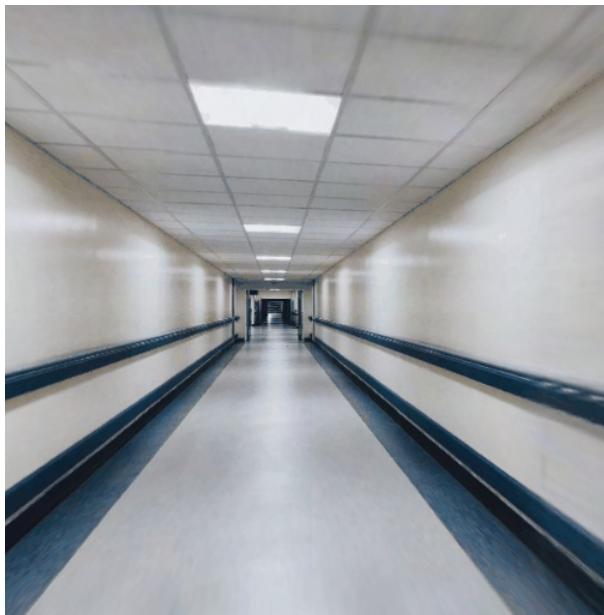


SynSin [Wiles et al. 2020]

# View Synthesis



Input Image



Ours



SynSin [Wiles et al. 2020]

# Image Extrapolation

Jiayuan Mao

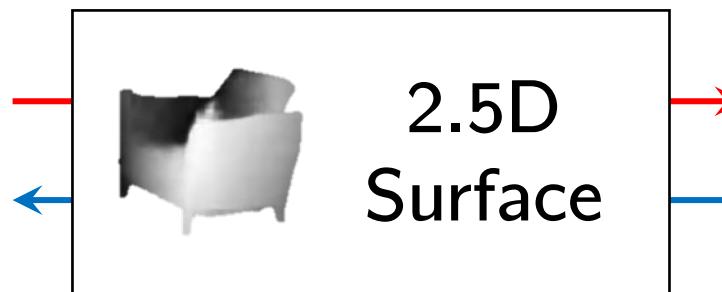
Yikai Li

Input Image

Content-Aware Scale Kaspar et al. [2015] InGAN [2019] Huang et al. [2014] Ours



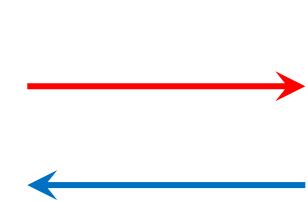
2D Image



2.5D  
Surface



3D Shape



```
draw(Back, Cuboid)  
  
for i in range(0, 1):  
    draw(Side, Cuboid, i)  
  
draw(Bottom, Cuboid)
```

Programs

### Key innovations:

- Neural nets for recognition +
- Domain knowledge for generalization: surfaces, objects, procedures

### Takeaways:

- Think about the causal structure behind visual data.
- Model what we can, learn what we can't.