

Day 10: Part 1

Speakers: Prof Christian Richardt, University of Bath

Title: Towards Reconstructing + Editing the Visual World

Space life  
Capture



## Roadmap

4 Papers  
will be  
discussed  
in the  
sessions

#	Project	Venue
1	MatryODShka	ECCV 2020 (oral)
2	Deep Video Portraits	SIGGRAPH 2018
3	HoloGAN	ICCV 2019
4	BlockGAN	NeurIPS 2020

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

4



## MatryODShka: Real-time 6DoF Video View Synthesis using Multi-Sphere Images

Benjamin Attal <sup>1,2</sup>



Selena Ling <sup>1</sup>



Aaron Gokaslan <sup>1</sup>



Christian Richardt <sup>3</sup>

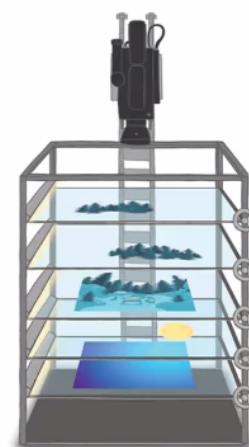


James Tompkin <sup>1</sup>



## Multi-plane and multi-sphere images

- Multi-plane images
  - Perspective view synthesis
  - RGB + alpha layers
  - Inferred from plane sweep volumes



Disney's multi-plane camera

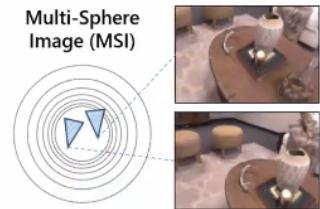
2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

6

## Multi-plane and multi-sphere images

- Multi-plane images
  - Perspective view synthesis
  - RGB + alpha layers
  - Inferred from plane sweep volumes



- Multi-sphere images
  - 360° view synthesis
  - RGB + alpha layers
  - Inferred from sphere sweep volumes



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

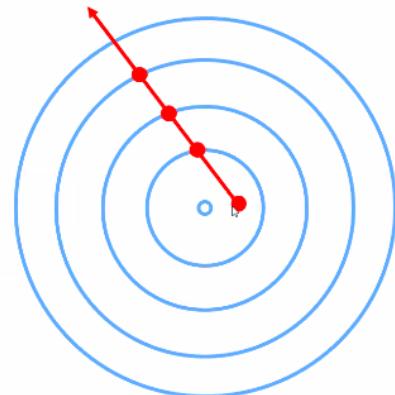
7

## Multi-sphere image rendering

1. Intersect ray with each layer of MSI
2. Over-composite colors  $\mathbf{c}_i$  and alphas  $\alpha_i$  of intersection points:

$$\mathbf{c} = \sum_{i=1}^N \mathbf{c}_i \cdot \alpha_i \cdot \prod_{j=1}^{i-1} (1 - \alpha_j)$$

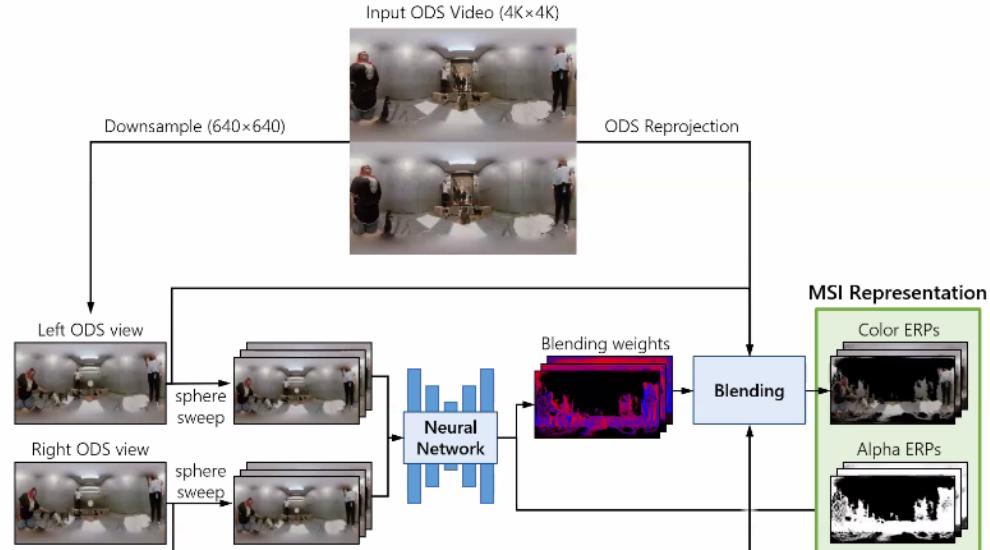
Net opacity of layer  $i$



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

8

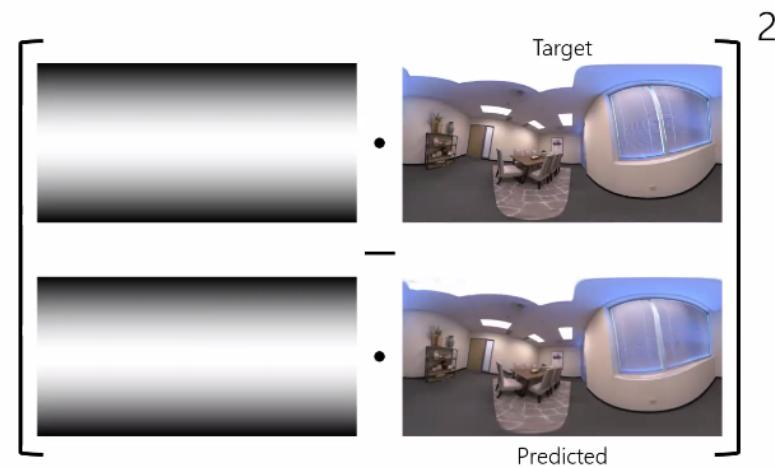


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

9

## Spherically adapted training loss

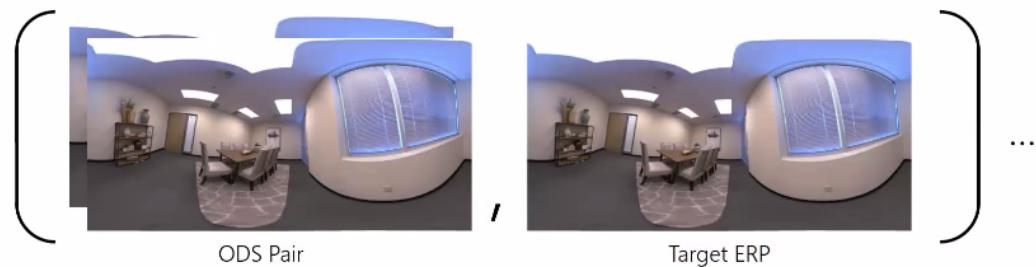


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

10

## Training data



Straub et al. (2019). The Replica dataset: A digital replica of indoor spaces

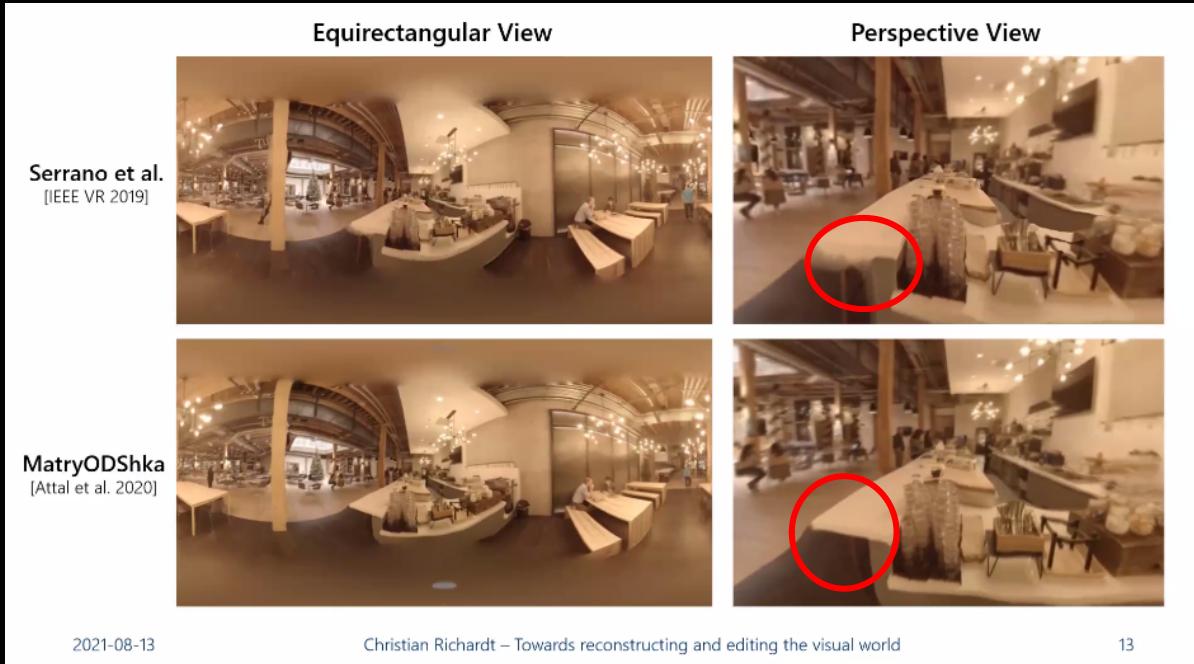
2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

11

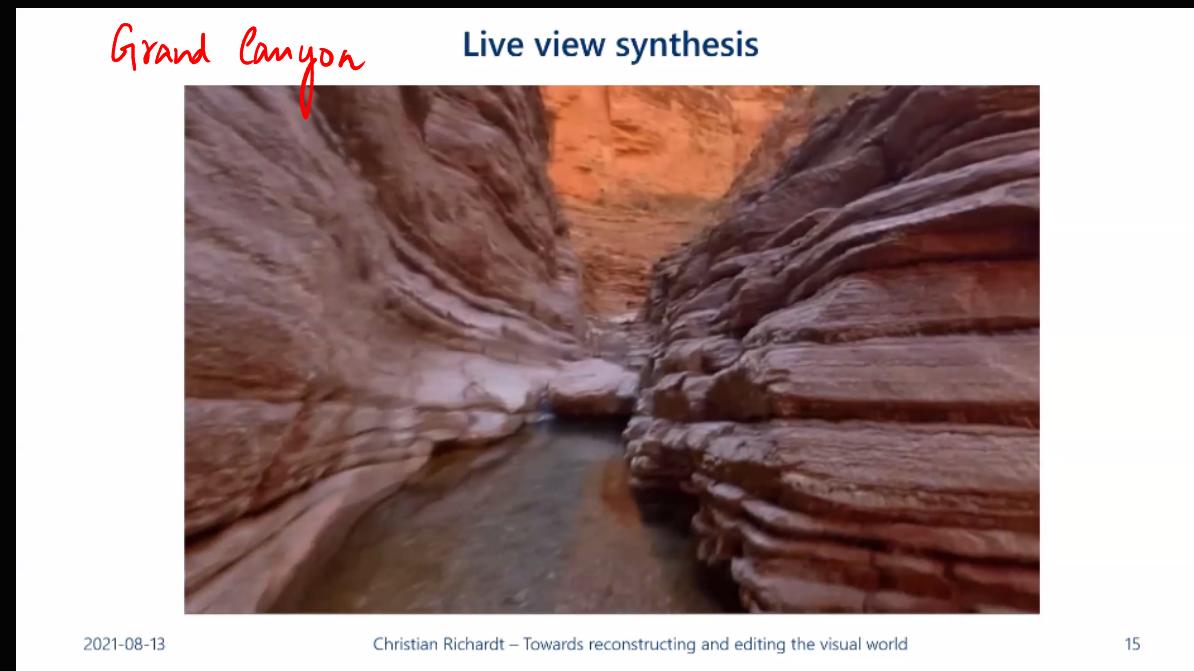
## Panoramic view synthesis



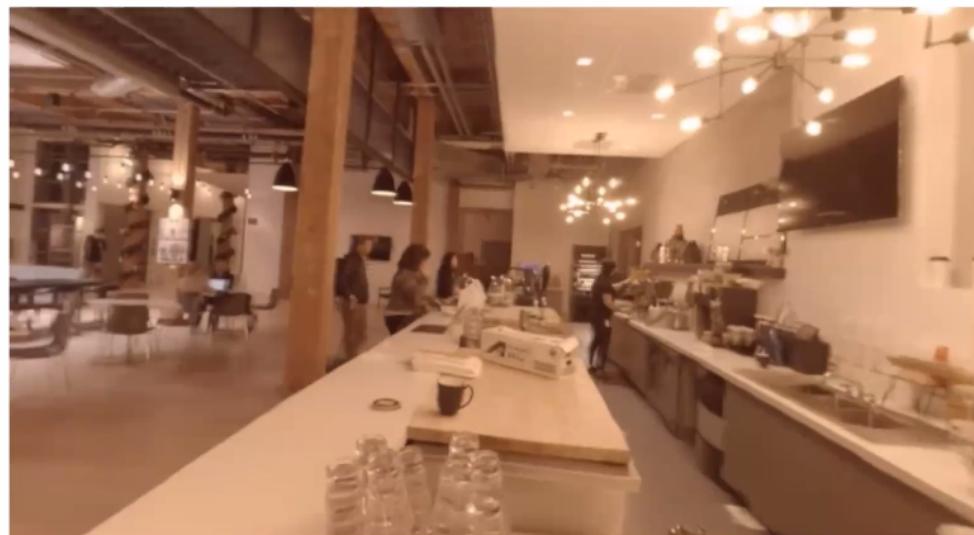


### Concurrent live view synthesis sphere images

<p><b>Immersive Light Field Video</b>  [Broxton et al., SIGGRAPH 2020]</p> <ul style="list-style-type: none"> <li>▪ 46 camera input</li> <li>▪ Offline (25 CPU hours per frame)</li> <li>▪ High-resolution inference</li> </ul>	<p><b>MatryODShka (Ours)</b>  [Attal et al., ECCV 2020]</p> <ul style="list-style-type: none"> <li>▪ ODS input</li> <li>▪ Online (30 Hz)</li> <li>▪ Low-resolution inference</li> </ul>	
		
2021-08-13	Christian Richardt – Towards reconstructing and editing the visual world	14



## Live view synthesis



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

16

# Deep Video Portraits

Hyeongwoo Kim<sup>1</sup> Pablo Garrido<sup>2</sup> Ayush Tewari<sup>1</sup> Weipeng Xu<sup>1</sup>

Justus Thies<sup>3</sup> Matthias Nießner<sup>3</sup> Patrick Pérez<sup>2</sup>

Christian Richardt<sup>4</sup> Michael Zollhöfer<sup>5</sup> Christian Theobalt<sup>1</sup>

<sup>1</sup> MPI Informatics



<sup>2</sup> Technicolor



<sup>3</sup> TU Munich



<sup>4</sup> University of Bath



<sup>5</sup> Stanford University



GENERATIONS / VANCOUVER  
SIGGRAPH 2018

## Contribution

Head Pose  
along with  
expression  
Captured



Original video

Pose editing

Expression editing

- Editing of head pose, rotation, face expression and eye gaze
- Combination of model-based face capture and CNN

Video courtesy of UK government (Open Government Licence)



18

## Related Work

### Model-based face capture and reenactment



Garrido et al., ToG 2016

Kemelmacher-Shlizerman et al., ECCV 2010  
 Shi et al., ToG 2014  
 Suwanjanakorn et al., ICCV 2015  
 Thies et al., CVPR 2016  
 Averbuch-Elor et al., ToG 2017  
 Thies et al., SIGGRAPH 2018

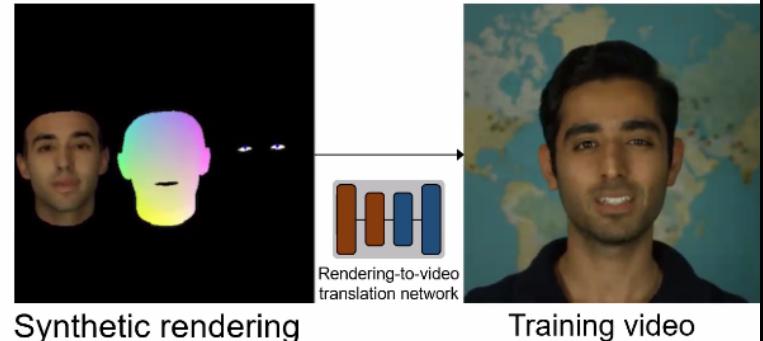
### CNN-based methods



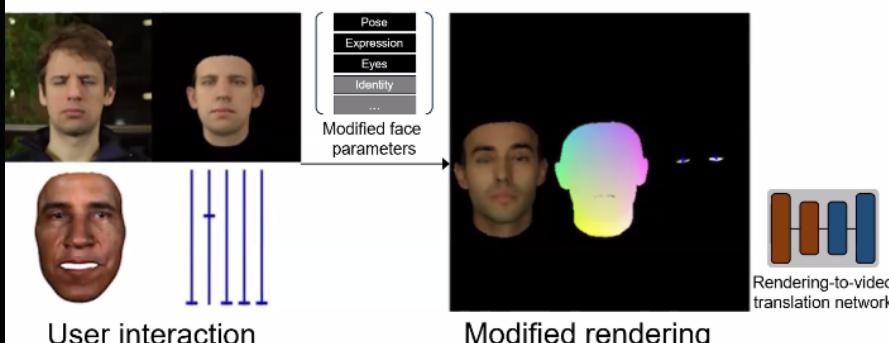
Karras et al., ICLR 2018

Goodfellow et al., NIPS 2014  
 Isola et al., CVPR 2017  
 Chen and Koltun, ICCV 2017  
 Tewari et al., ICCV 2017  
 Olszewski et al., ICCV 2018  
 Wang et al., CVPR 2018

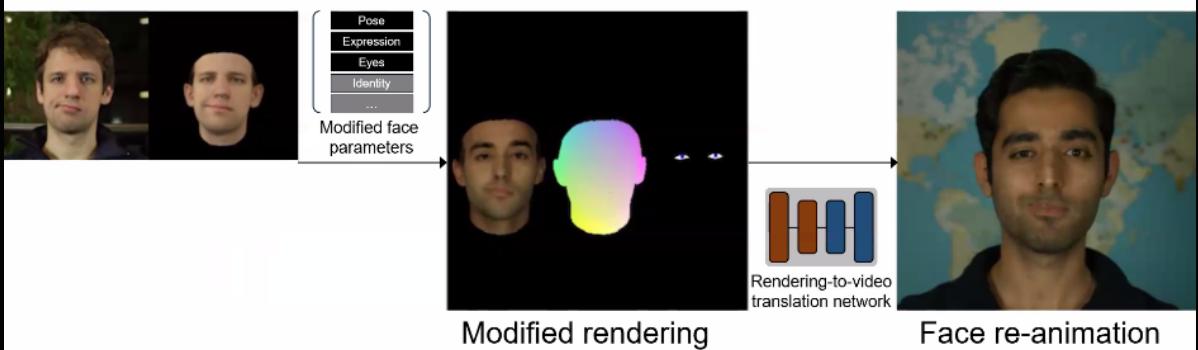
## Overview



## Overview

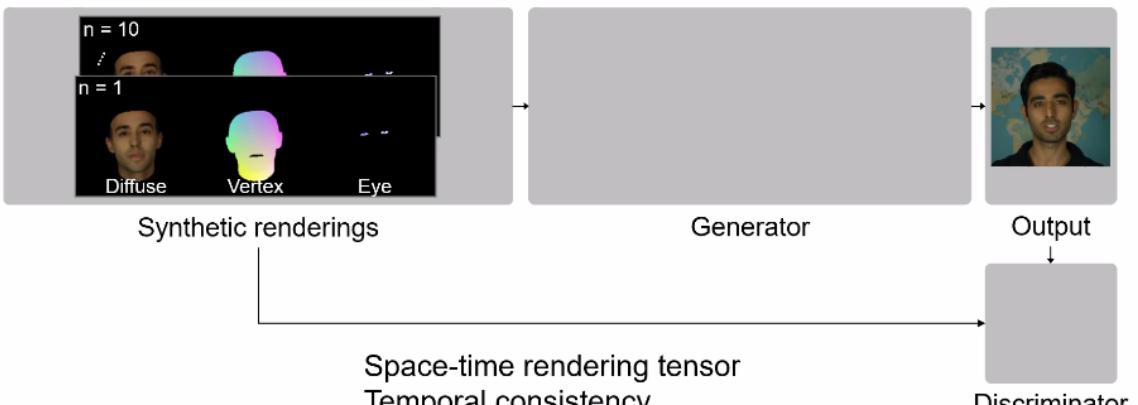


## Overview



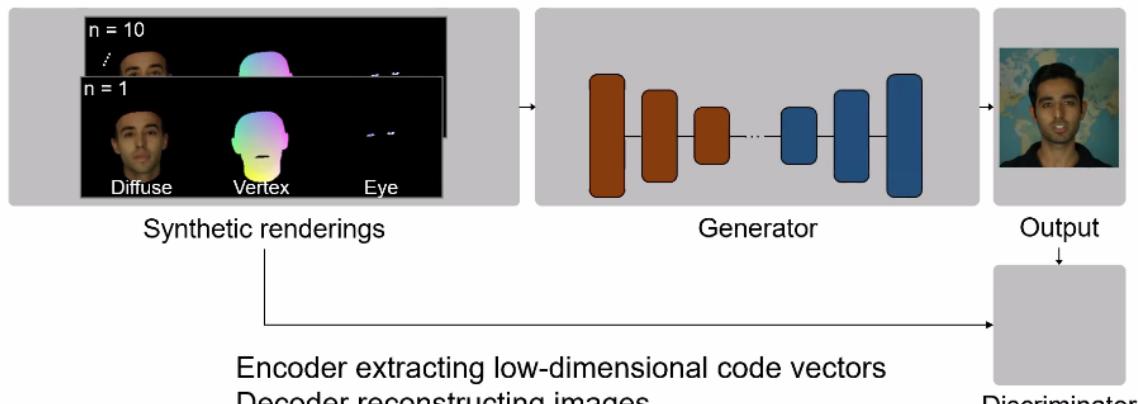
22

## Rendering-to-Video Translation Network



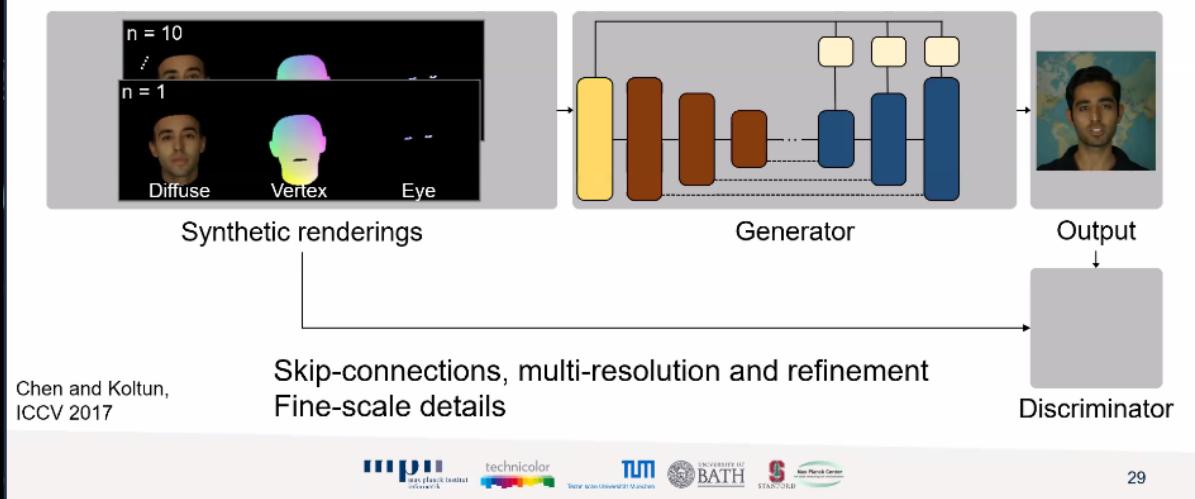
25

## Rendering-to-Video Translation Network

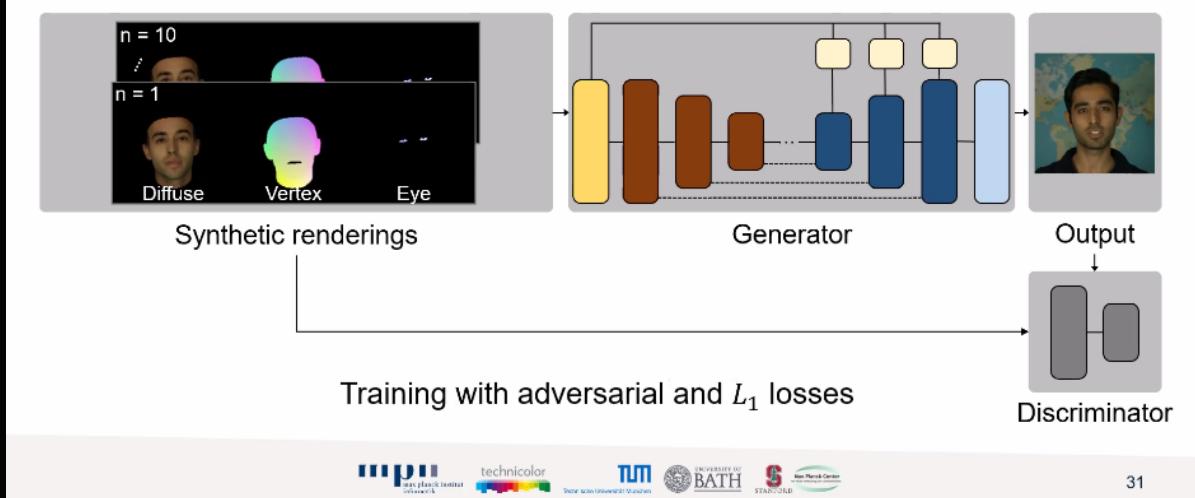


27

## Rendering-to-Video Translation Network

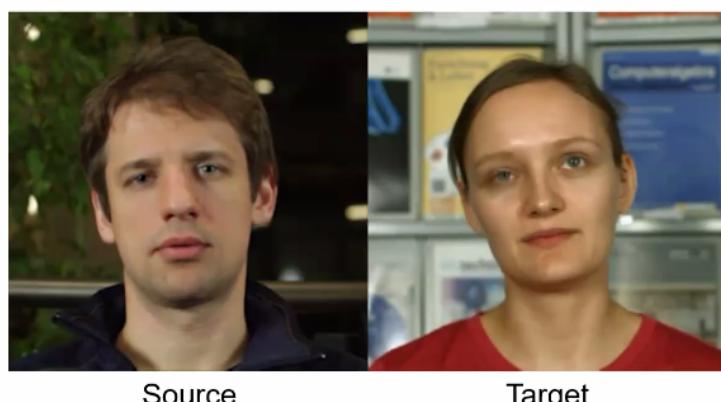


## Rendering-to-Video Translation Network



## Result: Facial Reenactment

Full reenactment of head pose, head rotation, face expression and eye gaze



## Result: Facial Reenactment

Full reenactment of head pose, head rotation, face expression and eye gaze



Source

Target

Face2Face  
(Thies et al., 2016)

## Result: Facial Reenactment



Source

Target

Result

## Result: Visual Dubbing

Visual discomfort due to the discrepancy between video and audio tracks



Dubbing actor video

Original video

## Result: Visual Dubbing

Modification of mouth motion to match audio tracks



Dubbing actor video

Dubbed video

Garrido et al., 2015

## Result: Interactive Editing



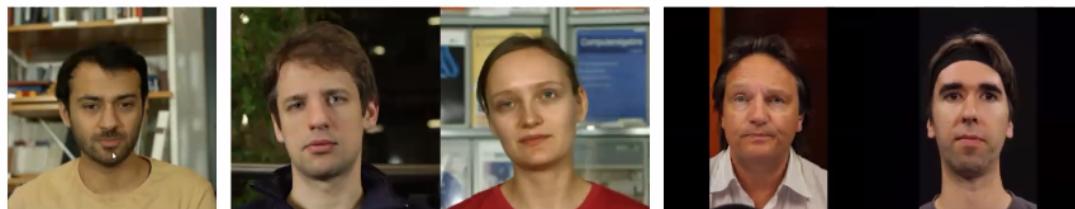
Pose

Expression

Shape

Approximately 9 fps

## Summary



Future work:

- Pushing toward higher quality and resolution
- Video authentication and forensics



# HoloGAN:

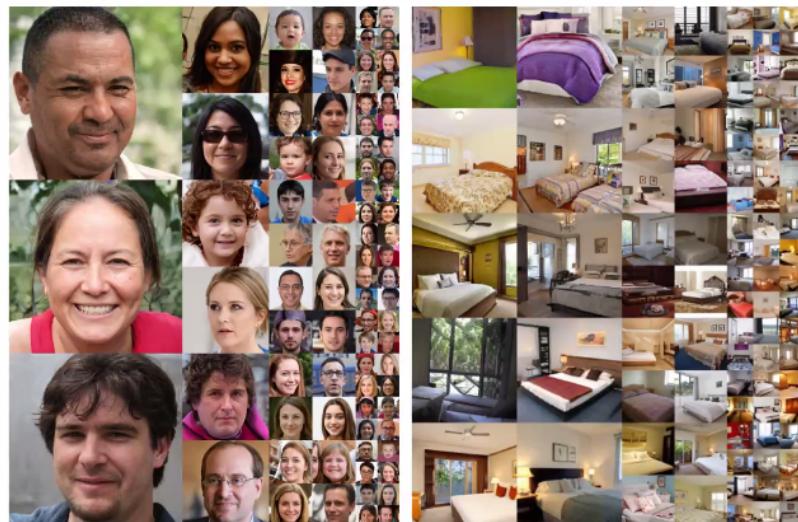
## Unsupervised learning of 3D representations from natural images

ICCV 2019

Thu Nguyen-Phuoc Chuan Li Lucas Theis Christian Richardt Yong-Liang Yang



### Generative adversarial networks



[Karras et al., StyleGAN, 2019]

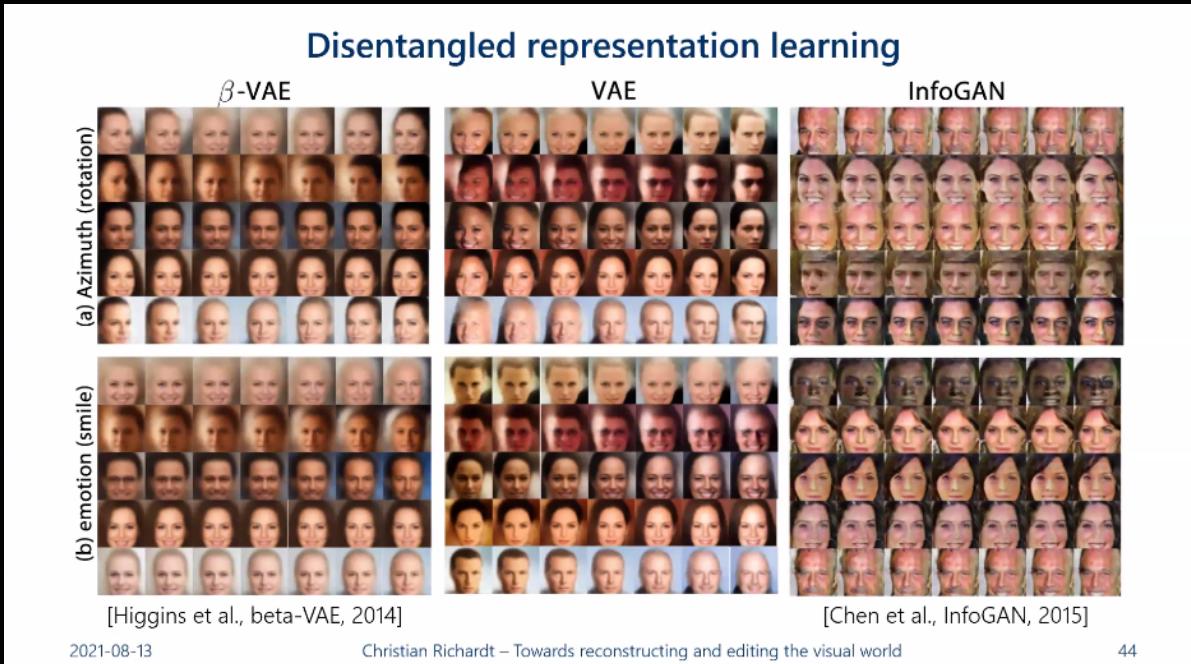
2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

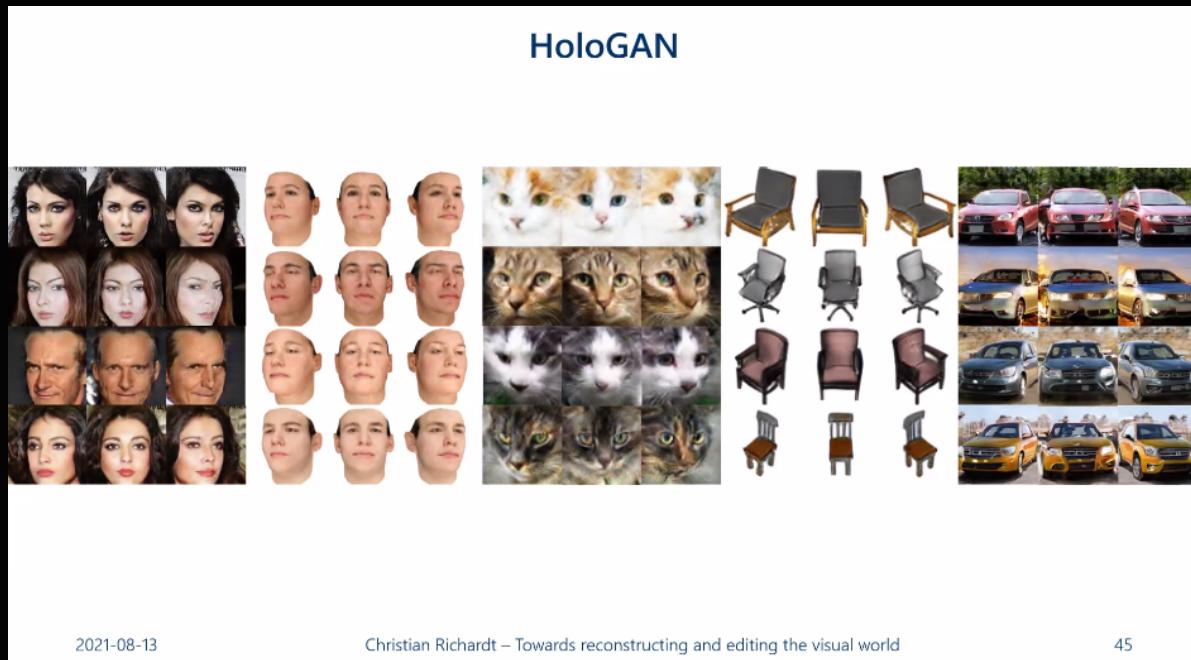
42

\*missed slide (to be updated)

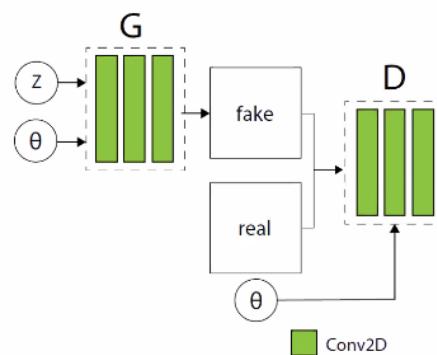
## Disentangled representation learning



## HoloGAN



## Conditional GANs

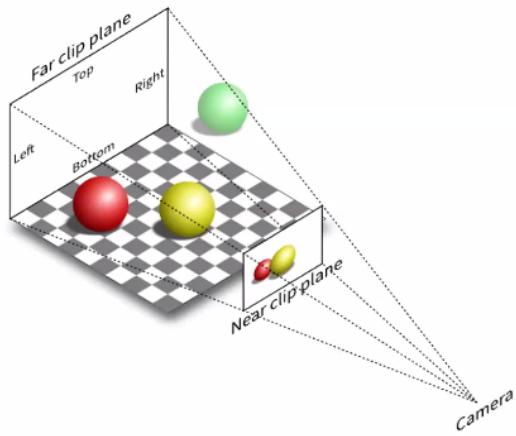


# Computer Graphics → Basics

## Back to the basics ...

Rendering comprises two main steps:

- Compute visibility
- Compute shading

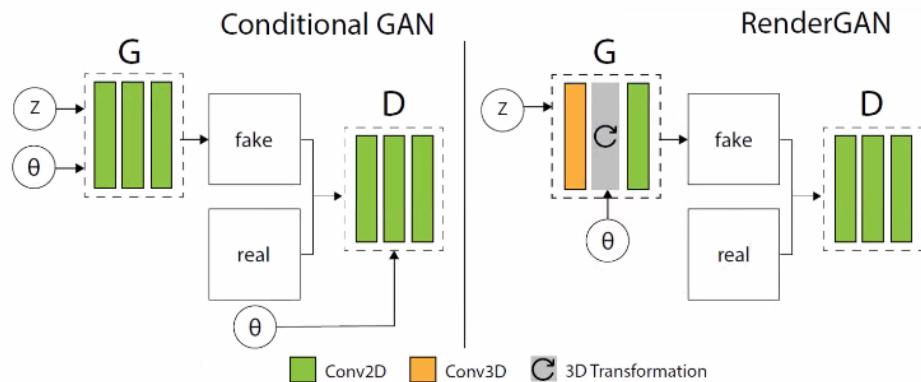


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

47

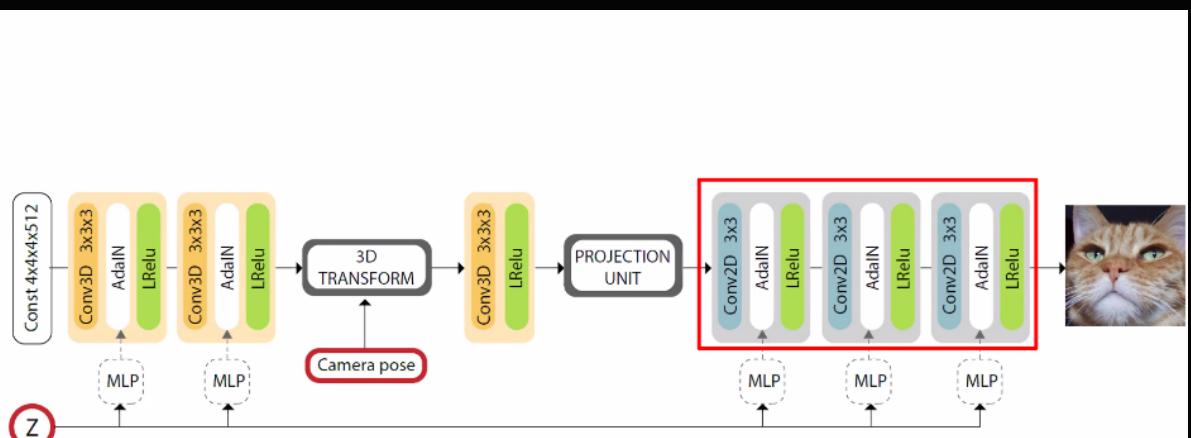
## GAN architecture comparison



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

48



$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(T(z)))))]$$

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

49

## Datasets



Basel



CelebA



Cats



Chairs



Cars

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

51

## Cats



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

52

## CelebA

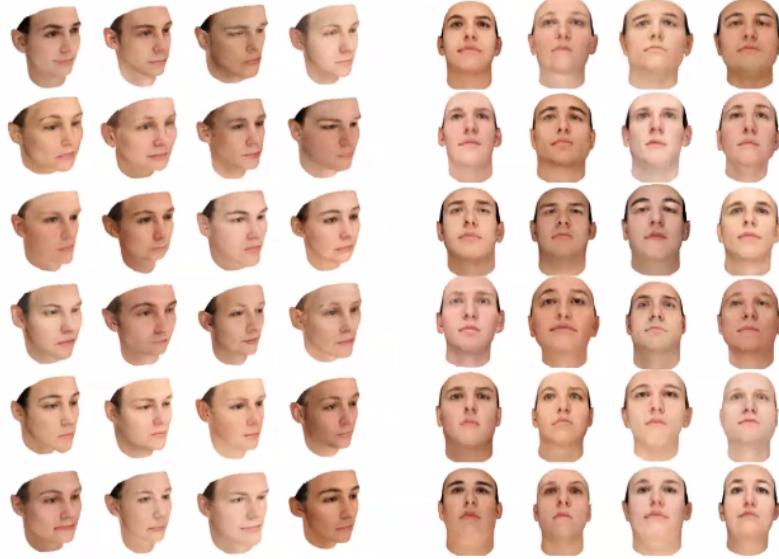


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

53

## Basel Face Model



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

54

## Chairs

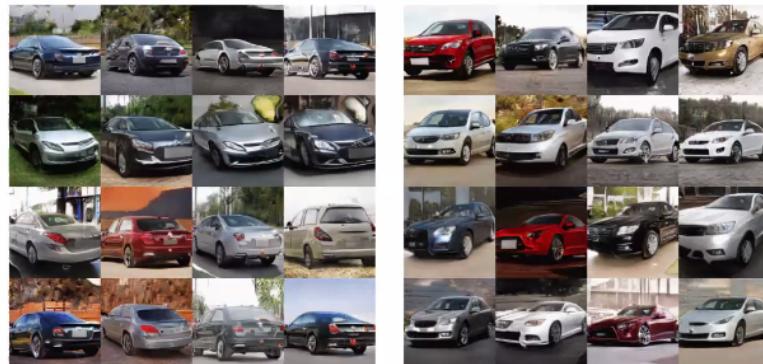


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

55

## Cars

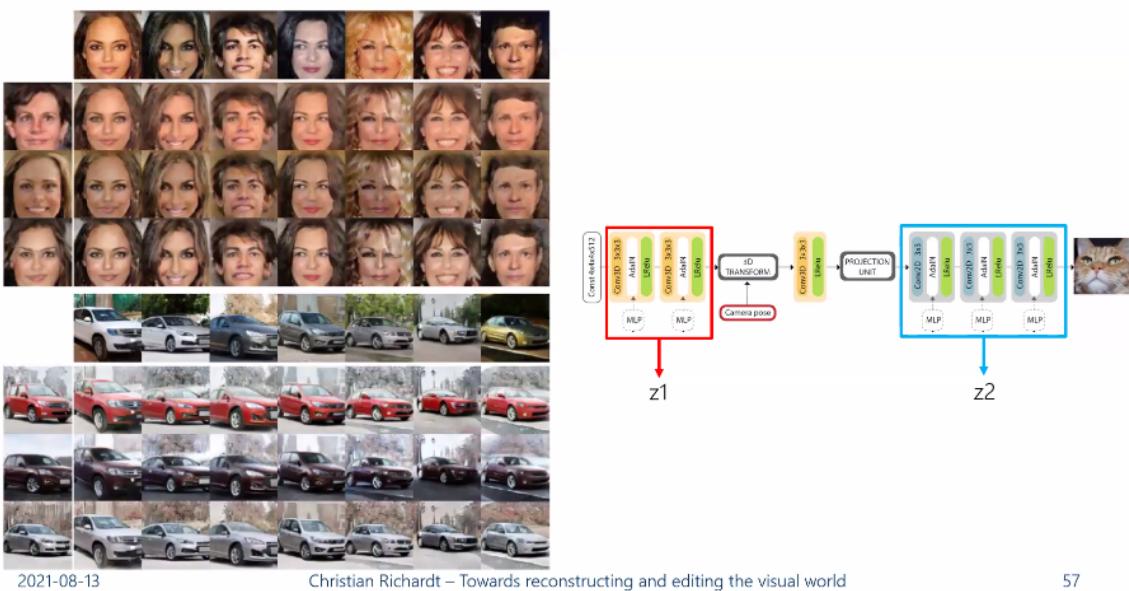


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

56

## Style mixing



## Separating shape + appearance



2021-08-13 Christian Richardt – Towards reconstructing and editing the visual world 59

## Separating shape + appearance



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

59

## The next frontier — complex scenes



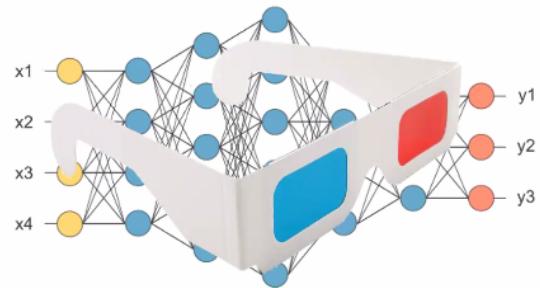
2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

61

## HoloGAN conclusion

- Adding inductive bias about the 3D world to a neural network
  - Better image quality
  - Better 3D understanding
  - Completely unsupervised



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

62



## BlockGAN

Learning 3D object-aware  
scene representations from unlabelled images

NeurIPS 2020

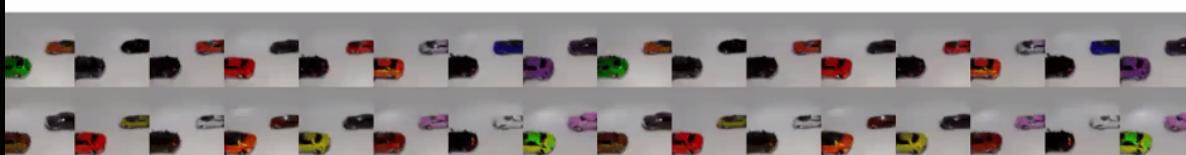
Thu Nguyen-Phuoc

Christian Richardt

Long Mai

Yong-Liang Yang

Niloy Mitra



## Current work



Scene Representation Networks  
[Sitzmann et al. 2019]



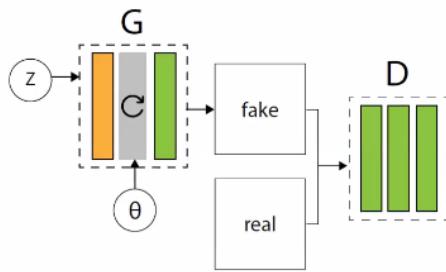
HoloGAN  
[Nguyen-Phuoc et al. 2019]

2021-08-13

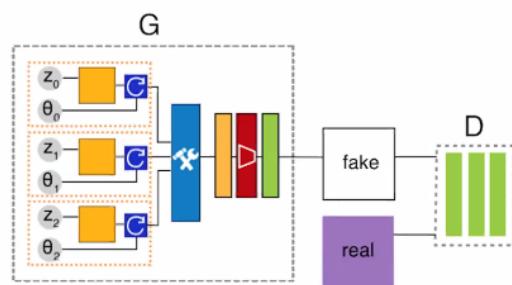
Christian Richardt – Towards reconstructing and editing the visual world

64

### HoloGAN



### BlockGAN

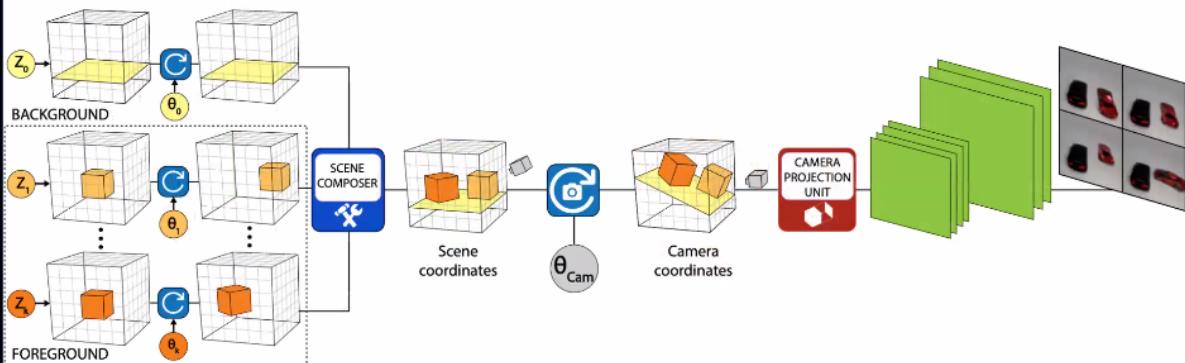


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

65

### BlockGAN architecture



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

66

## Datasets



Chair 1



Car 2



CLEVR 4



Cars (64×64)

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

67

## Rotation



\* all  
datasets  
rotated  
front &  
back

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

68

## Translation



\* Up &  
down  
translate -  
on

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

70

## Changing background appearance



light &  
shadow

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

71

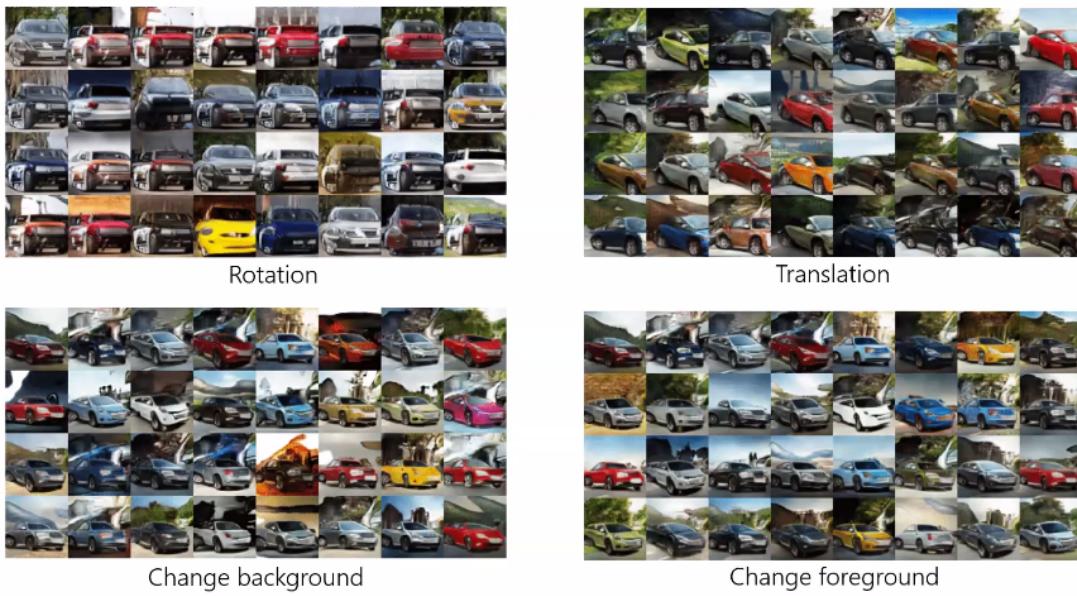
## Changing object #1 appearance



2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

72



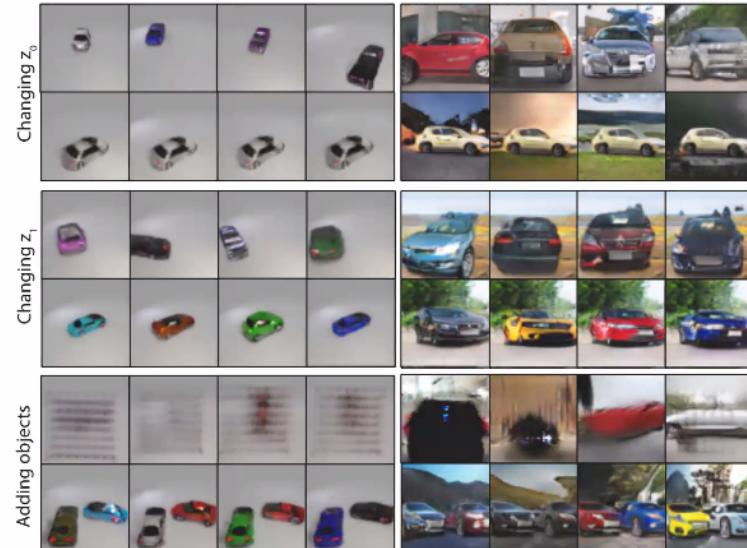
2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

73

## Generalisation

## Comparison with 2D methods

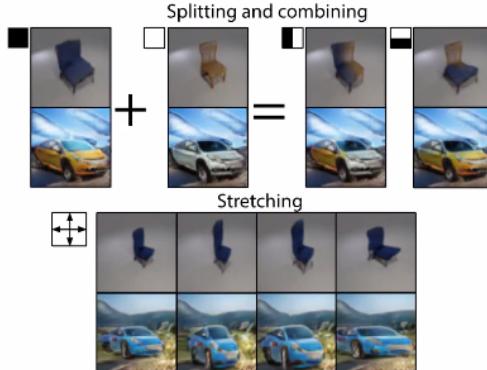


2021-08-13

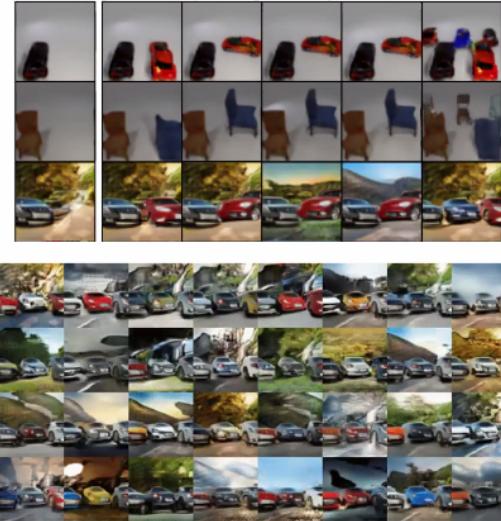
Christian Richardt – Towards reconstructing and editing the visual world

75

## Generalisation



Adding and manipulating new objects at test time

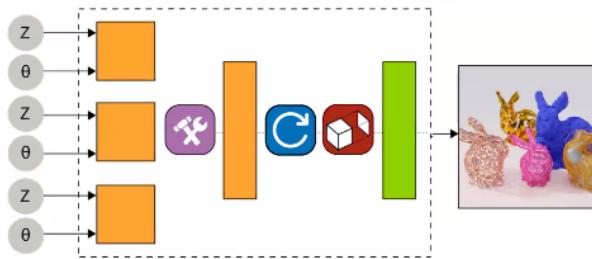


2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

76

## BlockGAN



- BlockGAN offers control over pose and appearance of individual object in the scene
- Deep voxel representations offer intuitive object spatial manipulation and composition

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

77

## Takeaways

- Traditional visual computing approaches build on explicit models
- Explicit models have limits:
  - sometimes difficult to model the real world, e.g. realistic faces
- Machine learning is reinvigorating visual computing:
  - can learn to model relationships from sufficient data (& expressive network)
  - no explicit model necessary, but often a black box ...
  - deep learning in particular has shown amazing results
- Benefits to applying domain knowledge to structure networks:
  - helps learning + we know what's going on!

2021-08-13

Christian Richardt – Towards reconstructing and editing the visual world

78

# Questions?

Christian Richardt

## Towards reconstructing + editing the visual world



CAMERA

Center for the Analysis of Motion,  
Entertainment Research and Applications



UNIVERSITY OF  
BATH

richardt.name  
 c\_richardt