### The Potential of AI for 3D Scene Digitization



Recent Learning-Based Innovations on Handling Challenging Scenarios

**Michael Weinmann** Delft University of Technology, Netherlands



The Potential of Al for 3D Scene Digitization





The Potential of Al for 3D Scene Digitization

- Replication/simulation of human intelligence in machines
  - "Ability to learn and perform suitable techniques to solve problems and achieve goals appropriate to the context in an uncertain, ever-varying world" – Manning, 2020



https://www.javatpoint.com/application-of-ai

From (visual) navigation (e.g. Google Maps/Earth) ...



... to scene understanding ...

- Based on scene geometry
- Based on spatial layout or arrangement of objects





https://www.thedubrovniktimes.com/news/dubrovnik/item/1765-amazing-video-from-dubrovnikto-king-s-landing?fb comment id=1072838662825020 1123759387732947

How does lighting influence object appearance?

Images partially taken from presentermedia.com

#### https://www.youtube.com/watch?v=kzNVkc4gB6U

Ideally multi-modal immersive experience

Beyond standard displays 

"Holodeck" experience



Telepresence: "Subjective experience of being in an environment that may differ from the user's actual local physical surrounding"







**TUDelft** M. Weinmann, The Potential of AI for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

### How to (digitally) represent scenes?

How to represent and capture 3D scenes?

### What do we see?





**TUDelft** M. Weinmann, The Potential of AI for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

### What do we see?

### Color snapshot

- Intensity of light
  - Seen from a single view point
  - At a single time
  - As a function of wavelength







### What do we see?



The plenoptic function

 Reconstruction of scene appearance under every possible view, at every moment, from every position, at every wavelength

 $\rightarrow$  It completely captures our visual reality!



Slide credit: A. Efros

# Model geometry or just capture images?

#### Geometry-image continuum ...



# Model geometry or just capture images?

Reduce complexity?

- Assumption: known surfaces
- → (traditional)
  computer graphics



Slide by Rick Szeliski and Michael Cohen



Background: 2D image synthesis (via rendering)

- Underlying principle: ray tracing (standard technique in graphics)
- Idea:
  - Measure light that arrives in the camera image
  - "Shooting" rays from a viewpoint through an image grid into the 3D scene
  - Intersection with geometric structures (e.g., ray-triangle intersection for meshes)
  - Color assignment (per pixel) based on evaluation of reflectance model



Image credits: C. Dachsbacher

# How should we represent complex 3D scenes?

But how can we get information about geometry/reflectance?

TUDelft M. Weinmann, The Potential of AI for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

### **Conventional 3D Scanning**

Separate 3D scanning with classical techniques:

- Structured light systems (active)
- Laser scanners (active)
- Multi-view stereo (passive)



M. Weinmann, C. Schwartz, R. Ruiters, and R. Klein. *A Multi-Camera, Multi-Projector Super-Resolution Framework for Structured Light*. 3DIMPVT, 2011



distance **d** 



McCann, 3D Reconstruction from Multiple Images



### **Conventional 3D Scanning**



#### Example:



### **Conventional 3D Scanning**



**3x speed** 

#### Example:

 Visual SLAM (online/real-time reconstruction)

VR-based telepresence/teleoperation [TVCG 19, 15MAR 19, 3DV'19, CVPRW'19, IROS'19, ICCVW'23]

### **Conventional 3D Scanning**



#### Example:

 Visual SLAM (online/real-time reconstruction)





VR-based telepresence/teleoperation [TVCG'19, ISMAR'19, 3DV'19, CVPRW'19, IROS'19, ICCVW'23]

## Classical 3D Scanning Techniques

#### Separate 3D scanning with classical techniques











**TUDelft** 



https://www.youtube.com/



https://www.youtube.com/watch?v=sQmlxPVtOGk

M. Weinmann, The Potential of Al for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

# Conventional 3D Scanning Techniques

#### Challenges:

Dynamic range



# Conventional 3D Scanning Techniques

### Challenges:

- Dynamic range
- Device resolution



## Conventional 3D Scanning Techniques

#### Challenges:

- Dynamic range
- Device resolution
- Optically complicated materials





Overcome challenges of classical techniques based on:

HDR Scanning



M. Weinmann, R. Ruiters, A. Osep, C. Schwartz, and R. Klein, A Multi-Camera, Multi-Projector Super-Resolution Framework for Structured Light, 3DIMPVT 2011



Overcome challenges of classical techniques based on:

Super-Resolution Framework for Structured Light, 3DIMPVT 2011

- HDR Scanning
- Superresolution

(4,4)(4,3)(2,3)(2,2)(1,5) (4,2)(4,1) projector 1 L50px 150px 150px projector 2 150p 0p superresolution 150px 150px M. Weinmann, R. Ruiters, A. Osep, C. Schwartz, and R. Klein, A Multi-Camera, Multi-Projector 3D/





Overcome challenges of classical techniques based on:

- HDR Scanning
- Superresolution
- Combination of different techniques to compensate for their individual limitations



M. Weinmann, R. Ruiters, A. Osep, C. Schwartz, and R. Klein, Fusing Structured Light Consistency and Helmholtz Normals for 3D Reconstruction, BMVC 2012



Overcome challenges of classical techniques based on:

- HDR Scanning
- Superresolution
- Combination of different techniques to compensate for their individual limitations
- Techniques tailored to complicated objects/surfaces



observations ightarrow surface consistency





M. Weinmann, A. Osep, R. Ruiters and Reinhard Klein. Multi-View Normal Field Integration for 3D Reconstruction of Mirroring Objects, ICCV 2013



# How should we represent complex 3D scenes?

Now that we have geometry, how do we get reflectance?

### Why do we need reflectance?



#### Important visual cue regarding "how do things feel"



Photo



Texture

BTF

### How to model light exchange at the surface?

Describe the material appearance decoupled from the environment, lighting and observer characteristics





How to capture appearance characteristics?

- Measurement of reflectance samples
  - Irregular
  - Occlusion
  - Not measured





Separate geometry and reflectance reconstructions



C. Schwartz, M. Weinmann, R. Ruiters, and R. Klein, Integrated High-Quality Acquisition of Geometry and Appearance for Cultural Heritage, VAST 2011

Separate geometry and reflectance reconstructions

**TU**Delft







C. Schwartz, M. Weinmann, R. Ruiters, and R. Klein, Integrated High-Quality Acquisition of Geometry and Appearance for Cultural Heritage, VAST 2011

### High-quality results, but ...



#### Trade-off: Acquisition technology vs. expressiveness of models:



setup complexity
# How should we represent complex 3D scenes?

Are we done?

Unfortunately not ...









# How should we represent complex 3D scenes?



## Al-based Scene Capture/Modeling



what we have to provide

### Traditional learning-based methods:

Supervised learning

Problem:





## Al-based Scene Capture/Modeling

#### Traditional learning-based methods:

Supervised learning

Problem:



#### Challenges:

- Different objects/scenes
- Different materials
- Different view conditions
- Different lighting conditions

→ Large datasets!





## Al-based Scene Capture/Modeling

How to get adequate training data?



# How to get adequate training data?



Approaches:

- Manual collection of labeled training data
  - Time-consuming (too many configurations)
  - Costly



Easy to label



(https://europe.naverlabs.com/research/computer-vision/proxy-virtual-worlds-vkitti-2/)

## Al-based Scene Capture and Modeling

AI-based scene capture and modeling? How?

Given adequate datasets, we can ...

... learn to recognize materials (e.g. to steer capture)



M. Weinmann, J. Gall, R. Klein. Material classification based on training data synthesized using a BTF database, ECCV, 2014

… learn feature correspondences



P. Truong, M. Danelljan, L. Van Gool, R. Timofte. Learning Accurate Dense Correspondences and When to Trust Them, CVPR, 2021

Given adequate datasets, we can ...

 ... learn to estimate geometry from a single image (also used for SLAM ...)

#### Training data: RGB + Depth pairs



NYU Depth V2 dataset

#### Predictions:



 $(R, G, B) \longrightarrow$ 

https://paperswithcode.com/task/monocular-depth-estimation

#### Loss components:

• Depth / depth relations

 $\rightarrow d$ 

- Surface normal
- Gradient information

Various applications:

Turning pictures into 3D experiences
(→ geometry estimation from a single image + warping + inpainting)



Joseph Redmon, CSE455: Computer Vision, University of Washington

Given adequate datasets, we can ...

- ... learn to estimate reflectance from a single image
  - Combine AI ...
  - ... and CG



V Deschaintre, M Aittala, F Durand, G Drettakis, A Bousseau. Single-image svbrdf capture with a rendering-aware deep network, ACM Transactions on Graphics, 2018

surface normals

diffuse component

Given adequate datasets, we can ...

In the synthesize novel views from 2 input views



Given adequate datasets, we can ...

In the synthesize novel views from 2 input views



#### Given adequate datasets, we can ...

Input views from 2 input views



Fig. 6: Self-supervised soft masking. Predicted image (left) and learned mask activations (right). The masks are normalized (dark = high activation). The PSV layer depth is up to scale. Note how the learned masks correlate with depth in the scene.

#### Given adequate datasets, we can ...

Input views from 2 input views



Input 1

Input 2

Extrapolated view



#### Given adequate datasets, we can ...



## Al-based Scene Capture and Modeling

Seems complicated... Can we avoid pre-training models?

### What can we represent based on AI?

Self-supervision for ...

- Reflectance estimation from a single image
  - Combine AI ...
  - ... and CG



Aittala et al., Reflectance Modeling by Neural Texture Synthesis





#### Self-supervision for ...

Geometry reconstruction from multiple views





P.-H. Huang, K. Matzen, J. Kopf, N. Ahuja, J.-B. Huang. DeepMVS: Learning Multi-View Stereopsis, CVPR, 2018



## Can we use AI to encode/represent whole scenes?

### Self-supervision for ...

■ Novel view synthesis (→ no need for additional datasets!)

- Idea:
  - Train a model which overfits to one object/scene!
  - Leverage inverse rendering setup
- Assumption:
  - Good results = expressive/accurate model





#### Neural Radiance Fields [Mildenhall et al. 2020]



Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., Ng, R., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV, 2020



Neural Radiance Fields [Mildenhall et al. 2020]:

Blurry results ...



NeRF (Naive)

Neural Radiance Fields [Mildenhall et al. 2020]:

How to get MLPs to represent high-frequency functions?

Ground truth image

Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., Ng, R., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV, 2020

#### age Standard fully-connected net







Neural Radiance Fields [Mildenhall et al. 2020]:

How to get MLPs to represent high-frequency functions? 

Ground truth image

With "positional encoding"

Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., Ng, R., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV, 2020





 $(x, y) \longrightarrow$ 



 $\blacktriangleright$  (R,G,B)

Neural Radiance Fields [Mildenhall et al. 2020]:

 Improvement for high-frequency embedding of input

NeRF (with positional encoding)







 $(\mathbf{x}) \longrightarrow (\mathbf{c})$ 





### Results

32

Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., Ng, R.. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV, 2020

**TUDelft** M. Weinmann, The Potential of AI for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

64

Extension to unconstrained photo collections ...







65



#### ... close to real-time ...



one of the input images



live training

Müller, T., Evans, A., Schied, C. and Keller, A., Instant Neural Graphics Primitives with a Multiresolution Hash Encoding, SIGGRAPH, 2022



#### ... close to real-time ...



ongoing work on real-time NeRF



AI Artists with Instant NeRF | NVIDIA





#### Results:

- Dense reconstruction (ideally surfaces ...)
- Preservation of details
- Inspection under novel views



Images taken from Alex Trevithick and Bo Yang. GRF: Learning a General Radiance Field for 3D Scene Representation and Rendering

#### Current trend: Back to explicit 3D representations

■ Keep optimizer, but replace neural network → 3D Gaussian Splatting





https://towardsdatascience.com/a-comprehensiveoverview-of-gaussian-splatting-e7d570081362









B. Kerbl, G. Kopanas, T. Leimkühler, G. Drettakis.3D Gaussian Splatting for Real-Time Radiance Field Rendering, SIGGRAPH 2023





B. Kerbl, G. Kopanas, T. Leimkühler, G. Drettakis.3D Gaussian Splatting for Real-Time Radiance Field Rendering, SIGGRAPH 2023





B. Kerbl, G. Kopanas, T. Leimkühler, G. Drettakis.3D Gaussian Splatting for Real-Time Radiance Field Rendering, SIGGRAPH 2023
### Al-based Scene Represent acc

accurate novel view synthesis



Current trend: Back to explicit 3D

Keep optimizer, but replace neuro

lots of views (with camera parameters) required





[1, ] huge memory footprint

https://www.magnopus.com/blog/the-rise-of-3d-gaussian-splatting

B. Kerbl, G. Kopanas, T. Leimkühler, G. Drettakis. 3D Gaussian Splatting for Real-Time Radiance Field Rendering, SIGGRAPH 2023

# Al-based Scene Capture





# **Remaining Limitations**





TUDelft M. Weinmann, The Potential of AI for 3D Scene Digitization - Recent Learning-Based Innovations on Handling Challenging Scenarios

# **AI-based Scene Capture**



How to address the aforementioned limitations?

### Robustness to Distractors

Observation:

NeRFs suffer from artifacts by distraction and occlusion



B. Buschmann, A. Dogaru, E. Eisemann, M. Weinmann, B. Egger, RANRAC: Robust Neural Scene Representations via Random Ray Consensus, ECCV 2024





### Idea:

- Outlier filtering via RANSAC-like scheme
   Random Ray Consensus (iterative subset evaluation)
  - Sampling of observations from the images with camera poses
  - Fitting of hypothesis/model
  - Use obtained model to render images for unused input views
  - Hypothesis validation ( $\rightarrow$  inlier/outlier classification)
  - Selection of model with most inliers and re-estimation of model based on this consensus set

B. Buschmann, A. Dogaru, E. Eisemann, M. Weinmann, B. Egger, RANRAC: Robust Neural Scene Representations via Random Ray Consensus, ECCV 2024



### Idea:

Outlier filtering via Random Ray Consensus (iterative subset evaluation)
 NeRF
 RobustNeRF
 RANRAC
 Ground Truth\*



RANRAC: Robust Neural Scene Representations via Random Ray Consensus, ECCV 2024

### **Robustness to Distractors**





B. Buschmann, A. Dogaru, E. Eisemann, M. Weinmann, B. Egger, RANRAC: Robust Neural Scene Representations via Random Ray Consensus, ECCV 2024

80

### **Sparse-view Scenarios**





D. Haitz, M. Hermann, A. S. Roth, M. Weinmann, M. Weinmann. The Potential of Neural Radiance Fields and 3D Gaussian Splatting for 3D Reconstruction from Aerial Imagery, ISPRS Annals, 2024

### Sparse-view Scenarios

Observation:

GS (20 views)

 NeRFs and Gaussian splatting produces a lot of noise for configurations with only a few views

198.88

B. Buschmann, E. Eisemann, M. Weinmann. ongoing work on robust sparse-view Gaussian splatting







#### Potential solution for few-shot capture

 Integration of neural inference of 2D-3D mappings, 3D Gaussian Splatting and camera refinement



B. Buschmann, E. Eisemann, M. Weinmann. ongoing work on robust sparse-view Gaussian splatting Improved practical relevance! (robustness, few-shot, less priors, ...)



#### Potential solution for few-shot capture

- Scene priors → feature space (beyond RGB)
  - Projection of sampling points on the ray emitted from the target view into the reference views
  - Obtain the features in the reference views by interpolation





### Potential solution for few-shot capture

■ Scene priors → feature space (beyond RGB)



#### Feature extraction module:

- extracts deep features of the reference view image
- aggregates the features under different reference views as priors



### Potential solution for few-shot capture

■ Scene priors → feature space (beyond RGB)



NeRF takes the scene priors along with position and orientation encodings to generate color and opacity



### Potential solution for few-shot capture

Scene priors + camera refinement



Idea: Leverage inconsistencies in observations for mirroring surfaces

Allows detection and handling of mirrors



L. V. Holland, M. Weinmann, P. Stotko and R. Klein, NeRFs are Mirror Detectors:



Idea: Leverage inconsistencies in observations for mirroring surfaces

Allows detection and handling of mirrors



L. V. Holland, M. Weinmann, P. Stotko and R. Klein, NeRFs are Mirror Detectors:



Idea: Leverage inconsistencies in observations for mirroring surfaces

Allows detection and handling of mirrors



L. V. Holland, M. Weinmann, P. Stotko and R. Klein, NeRFs are Mirror Detectors:





91

Idea: Leverage inconsistencies in observations for mirroring

Allows detection and handling of mirrors



Idea: Leverage inconsistencies in observations for mirroring surfaces

Allows detection and handling of mirrors



Ours







Best Baseline Mip-NeRF 360 [3] Ground Truth



MS-NeRF [112] PSNR: 33.55 dB

Ours G PSNR: 35.27 dB

L. V. Holland, M. Weinmann, P. Stotko and R. Klein, NeRFs are Mirror Detectors:







### Idea:

Scene representation with semantics
 e.g. Gaussian Splats + Segment Anything



Ground truth (Multiview input frames) Rendered Gaussian Splats Rendered Segmentation maps

Distilled feature fields (Object-Ids for scene editing)



#### Idea:

Representations per sub-band combined with material-wise shading model



S. N. Sinha, J.Kühn, H. Grat, M. Weinmann. SpectralSplatsViewer: An Interactive Web-Based Tool for Visualizing Cross-Spectral Gaussian Splats, Web3D, 2024

S. N. Sinha, H. Graf, M. Weinmann. SpectralGaussians: A relightable spectral Gaussian splatting framework to generate photorealistic semantic Gaussian splats in different spectra, ongoing work



Extension to Multi-Spectral Scene Representation



S. N. Sinha, H. Graf, M. Weinmann. SpectralGaussians: A relightable spectral Gaussian splatting framework to generate photorealistic semantic Gaussian splats in different spectra, ongoing work



**Results**:



S. N. Sinha, J.Kühn, H. Graf, M. Weinmann. SpectralSplatsViewer: An Interactive Web-Based Tool for Visualizing Cross-Spectral Gaussian Splats, Web3D, 2024

S. N. Sinha, H. Graf, M. Weinmann. SpectralGaussians: A relightable spectral Gaussian splatting framework to generate photorealistic semantic Gaussian splats in different spectra, ongoing work

Method	Spectral NeRF Synthetic Dataset[27]						Average
	kitchen	Living room	Digger	Spaceship	Vintage car	Cartoon knight	Average
	PSNR ↑						
NeRF[7]	34.583	33.172	30.658	30.126	33.478	34.485	32.400
Mip-NeRF[8]	-	-	33.301	31.495	33.883	35.102	33.945
Aug-NeRF[37]	34.480	32.540	31.538	30.929	33.639	33.908	32.677
SpectralNeRF[27]	35.115	33.665	33.378	31.951	34.480	34.915	33.610
Ours	37.035	37.989	40.218	41.233	42.636	36.723	38.456
	SSIM ↑						
NeRF[7]	0.8943	0.9929	0.9187	0.9358	0.7958	0.9273	0.9123
Mip-NeRF[8]	-	-	0.9290	0.9475	0.8166	0.9572	0.9126
Aug-NeRF[37]	0.9026	0.9649	0.9248	0.9402	0.8002	0.9287	0.9163
SpectralNeRF[27]	0.9117	0.9931	0.9357	0.9482	0.8169	0.9573	0.9349
Ours	0.9747	0.9733	0.9923	0.9951	0.9893	0.9572	0.9801
	LPIPS ↓						
NeRF[7]	0.1650	0.0578	0.0413	0.0275	0.1319	0.1545	0.0722
Mip-NeRF[8]	-	-	0.0435	0.0535	0.1747	0.1526	0.1061
Aug-NeRF[37]	0.1603	0.0706	0.0341	0.0389	0.1536	0.1705	0.0973
SpectralNeRF[27]	0.1637	0.0479	0.0259	0.0250	0.1499	0.1510	0.0733
Ours	0.0739	0.0525	0.0109	0.0084	0.0527	0.0741	0.0438



arXiv preprint arXiv:1912.06354, 2019

activity maps and action recognition

J. Tanke et al., Bonn Activity Maps: Dataset Description,

### **Further Trends**

### Tighter integration of context into inference/capture

... behavioral patterns





hybrid (static+dynamic) scene representation

L. V. Holland, P. Stotko, S. Krumpen, R. Klein, M. Weinmann. Efficient 3D Reconstruction, Streaming and Visualization of Static and Dynamic Scene Parts for Multi-client Live-telepresence in Large-scale Environments, ICCV Workshops 2023 L. Bruckschen, K. Bungert, M. Wolter, S. Krumpen, M. Weinmann, R. Klein, M. Bennewitz. Where can i help? Human-aware placement of service robots, International Conference on Robot and Human Interactive Communication, 2020







### **Further Trends**



### Combination with simulations



S. Krumpen, M. Weinmann, R. Klein. Interactive Appearance Manipulation of Fiber-based Materials, International Conference on Computer Graphics Theory and Applications 2017



#### Physics-informed machine learning

N. Wandel, M. Weinmann, R. Klein. Learning Incompressible Fluid Dynamics from Scratch-Towards Fast, Differentiable Fluid Models that Generalize, ICLR 2021

N. Wandel, M. Weinmann, R. Klein. Teaching the incompressible Navier–Stokes equations to fast neural surrogate models in three dimensions, Physics of Fluids, 2021

N. Wandel, M. Weinmann, M. Neidlin, R. Klein. Spline-PINN: Approaching PDEs without data using fast, physics-informed Hermite-spline CNNs, AAAI 2022

Editable representations

Semantic scene stylization



S. N. Sinha, H. Graf, M. Weinmann. SemanticSplatStylization: Semantic scene stylization based on 3D Gaussian splatting and class-based style transfer, GCH 2024

### Further Trends







#### Robust extensions of NeRFs and Gaussian Splatting



# Conclusions

- AI ...
  - ... offers breakthrough potential for numerous applications
  - ... leads to paradigm shift in capture technology
  - ... will make us re-think conventional processes
  - ... will help to tackle more challenging scenarios

### But ...

- ... cannot solve everything
- ... may be outperformed in certain conditions
- ... may benefit from coupling to traditional principles
- ... we should not forget to use natural intelligence







# Acknowledgements





+ numerous more collaborators

Saptarshi Neil

Sinha



Klein

Martin Weinmann



Matthias B. Hullin



Ralf Sarlette



Julian Iseringhausen





Andre

Rochow

Leif Van Holland



Jan Uwe Müller



Schwarz



Sven Behnke



Holger

Graf

Roland Ruiters

Nils

Wandel



Christopher Schwarz



Patrick Stotko

Elena

Trunz



Stefan Krumpen





Lukas Bode



Eisemann









Sebastian

Merzbach

102



Benno





# Thank you for your attention!



#### Questions? More Information?



### **Michael Weinmann** TU Delft, Netherlands

m.weinmann@tudelft.nl