# 3D Scene Representation Learning

Martin Oswald

Computer Vision Group, University of Amsterdam

# Motivation: 3D reconstruction is hard!

# Motivation: 3D reconstruction is hard!

# Motivation: 3D reconstruction is hard!

# Video Generation: Sora



Prompt: This close-up shot of a chameleon showc

# Video Generation: Sora



Limited 3D consistency
due to lack of 3D modeling!

Prompt: Beautiful, snowy Tokyo city

# Vanishing Points!?

# Video Generation: Sora



Prompt: A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp.

# Dall-E3



Man looking through telescope



Woman on a surfboard

# "Sora is also a Physics Engine!



"Photorealistic closeup video of two pirate ships battling each other as they sail inside a cup of coffee"

# "Sora is also a Physics Engine!



Ski jumping man

# Scene Representations

**Classical**



Occupancy  Signed Distance

Spline/NURBS    Point Cloud    Surface Mesh    Tetrahedral Mesh    Voxel Grid    Voxel Octree    Voxel Hashing

**explicit** **(topology change=>hard)**          **implicit** **(topology change=>simple)**

**Learned**



$f_\theta(p) = \tau$

Voxel Grid    Point cloud    Mesh    (Neural) Classifier

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]

Object images & poses    DeepVoxels

Global
Optimization

Training
Testing

Novel Views

Rendering

[DeepVoxels CVPR 2019]

$F_\Theta$

Ray 1

Ray 2

[NeRF - Neural Radiance Fields, ECCV 2020]

Interpolation

3D Location **p**    Occupancy Probability $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Features $\psi(\mathbf{p}, \mathbf{x})$

**p**

Fully-Connected Network

3D Feature Volume

1

0

[Peng et al., Convolutional Occupancy Networks, ECCV 2020]

**Neural implicit**

# Scene Representations

Spline/NURBS    Point Cloud    Surface Mesh    Tetrahedral Mesh    Occupancy Voxel Grid    Signed Distance    Voxel Octree    Voxel Hashing

**Classical / Non-Neural**

3D Gaussian Splatting

Plenoxels / Plenoctrees

$$\text{minimize}_{\{\sigma, \bullet\}} \mathcal{L}_{recon} + \lambda\mathcal{L}_{TV}$$

c) Volumetric Rendering

**explicit**

**implicit**

**Neural**

Voxel Grid   Point cloud   Mesh   (Neural) Classifier

$f_\theta(p) = \tau$

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]

[DeepVoxels CVPR 2019]

[NeRF - Neural Radiance Fields, ECCV 2020]

3D Location **p**   Occupancy Probability $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Interpolation

Features $\psi(\mathbf{p}, \mathbf{x})$

Fully-Connected Network

3D Feature Volume

[Peng et al., Convolutional Occupancy Networks, ECCV 2020]

**Neural implicit**

# Scene Representations



[https://arxiv.org/pdf/1803.03352.pdf]

Spline/NURBS    Point Cloud    Surface Mesh    Tetrahedral Mesh    Occupancy    Signed Distance    Voxel Octree    Voxel Hashing

Voxel Grid

3D Gaussian Splatting

SfM Points → Initialization → 3D Gaussians → Camera → Projection → Differentiable Tile Rasterizer → Image → Adaptive Density Control

→ Operation Flow    → Gradient Flow

Plenoxels / Plenoctrees

Spherical Harmonics

Predicted Color

c) Volumetric Rendering

$$\text{minimize } \mathcal{L}_{recon} + \lambda \mathcal{L}_{TV}$$
$$\{\sigma, \bullet\}$$

Training Image

**explicit**

**implicit**

**Neural**

Voxel Grid    Point cloud    Mesh    (Neural) Classifier

$f_\theta(p) = \tau$

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]

Object images & poses → Global Optimization → DeepVoxels

Training / Testing → Rendering → Novel Views

[DeepVoxels CVPR 2019]

$F_\Theta$    Ray 2    Ray 1

[NeRF - Neural Radiance Fields, ECCV 2020]

Interpolation    3D Location $\mathbf{p}$    Occupancy Probability $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Features $\psi(\mathbf{p}, \mathbf{x})$

Fully-Connected Network

3D Feature Volume

[Peng et al., Convolutional Occupancy Networks, ECCV 2020]

**Neural implicit**

14

# Implicit Volumetric Representation

- *Voxel grid*: sample a volume containing the surface of interest uniformly

- Label each grid point as lying *inside* or *outside* the surface

**Signed distance function**

*SDF = 0*

*SDF > 0*

*SDF < 0*

**Occupancy function**

*OF = 0.5*

*OF = 0*

*OF = 1*

- The modeled surface is represented as an *isosurface* (e.g. SDF = 0 or OF = 0.5) of the labeling (implicit) function

- Advantages: simple handling of topological changes, watertight surfaces, no self-occlusions
  Disadvantages: Large memory requirement, bad scalability to large scenes (cubic growth)

# Represent Scenes with TSDFs



$F > \mu$

$F < -\mu$

[Newcombe & Lovegrove, Geometric Reconstruction Lecture]

# Real-time Mapping - KinectFusion

# Scene Representations

**Classical / Non-Neural**

Spline/NURBS    Point Cloud    Surface Mesh    Tetrahedral Mesh    Occupancy    Signed Distance    Voxel Octree    Voxel Hashing

Voxel Grid



3D Gaussian Splatting

**explicit**

Plenoxels / Plenoctrees

**implicit**

SfM Points   Initialization   Camera   Projection   3D Gaussians   Adaptive Density Control   Differentiable Tile Rasterizer   Image   → Operation Flow   → Gradient Flow

$$f_\theta(p) = \tau$$

Voxel Grid    Point cloud    Mesh    (Neural) Classifier

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]

[DeepVoxels CVPR 2019]

[NeRF - Neural Radiance Fields, ECCV 2020]

3D Location **p**   Occupancy Probability $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Interpolation

Features $\psi(\mathbf{p}, \mathbf{x})$

Fully-Connected Network

3D Feature Volume

[Peng et al., Convolutional Occupancy Networks, ECCV 2020]

**Neural**

**Neural implicit**

# Neural Implicit Scene Representations



Voxel Grid    Point cloud    Mesh    (Neural) Classifier

$f_\theta(p) = \tau$

$\mathbf{p}$

$(\mathbf{p} \in \mathbb{R}^3)$

**MLP**

SDF /
Occupancy
[Color, ...]

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]

# Neural Implicit Representations

## Occupancy Networks



Input $\mathbf{x}$ — 1D Encoder

3D Location $\mathbf{p}$

Features $\psi(\mathbf{x})$

1
0

## Convolutional Occupancy Networks



Input $\mathbf{x}$ — 3D Encoder + 3D Conv

$z$

Trilinear Interpolation

$\mathbf{p}$

$x$

$y$  3D Feature Volume

3D Location $\mathbf{p}$

Features $\psi(\mathbf{p}, \mathbf{x})$

1
0

# Scene Representations

Spline/NURBS  Point Cloud  Surface Mesh  Tetrahedral Mesh  Occupancy Voxel Grid  Signed Distance  Voxel Octree  Voxel Hashing

**Classical / Non-Neural**



3D Gaussian Splatting

SfM Points · Initialization · Camera · Projection · Differentiable Tile Rasterizer · Image · 3D Gaussians · Adaptive Density Control · Operation Flow · Gradient Flow



Plenoxels / Plenoctrees

$$\text{minimize}_{\{\sigma, \bullet\}} \mathcal{L}_{recon} + \lambda \mathcal{L}_{TV}$$

**explicit**      **implicit**

**Neural**



Voxel Grid  Point cloud  Mesh  (Neural) Classifier

$f_\theta(p) = \tau$

Learned / Deep Representations:
OccNet [https://arxiv.org/pdf/1812.03828.pdf]
DeepSDF [https://arxiv.org/pdf/1901.05103.pdf]
IM-Net [https://arxiv.org/pdf/1812.02822.pdf]



[DeepVoxels CVPR 2019]

[NeRF - Neural Radiance Fields, ECCV 2020]



3D Location **p**  Occupancy Probability $f_\theta(\mathbf{p}, \psi(\mathbf{p}, \mathbf{x}))$

Interpolation  Features $\psi(\mathbf{p}, \mathbf{x})$  Fully-Connected Network  3D Feature Volume

[Peng et al., Convolutional Occupancy Networks, ECCV 2020]

**Neural implicit**

# Scene Representations for 3D Reconstruction



[Metthew Brennan, "Photogrammetry / NeRF / Gaussian Splatting comparison", YouTube 2023]

# Structure-from-Motion



**Image Set**

**Image Association**

unknown cameras

**SfM**

**Sparse Model**

- **NeRF**
- **Gaussian Splatting**

**MVS**

known cameras

**Scene Graph**

**(Semi-) Dense Model**

[Large-scale 3D Modeling Tutorial, CVPR 2017]

# Structure-from-Motion (SfM)

**Rome** dataset

74,394 images

[Johannes L. Schönberger, Jan-Michael Frahm. **Structure-from-Motion Revisited**. CVPR, 2016; COLMAP]

# Neural Radiance Fields
## (NeRF)

Input Images          Optimize NeRF          Render new views

# Neural Implicit Representations

[3D point, viewing angle] → → [Color, Density]

≈ Occupancy

# Why view-dependent colors?

# Neural Radiance Fields (NeRFs)

# Neural Radiance Fields (NeRFs)

5D Input
Position + Direction

$(x,y,z,\theta,\phi) \rightarrow$ $F_\Theta$ $\rightarrow (RGB\sigma)$

Output
Color + Density

Ray 1
Ray 2

Volume Rendering

$\sigma$ Ray 1

$\sigma$ Ray 2
Ray Distance

Rendering Loss

$\left\| \quad - g.t. \right\|_2^2$

$\left\| \quad - g.t. \right\|_2^2$

(a)　(b)　(c)　(d)

# Neural Radiance Fields (NeRFs)

# Gaussian Splatting

# Point Splatting



Surface splatting with EWA
[Zwicker et al. 2001]

[Zwicker et al. 2001; Yifan Wang,  2019]

# Gaussian Splatting



[Matsuki et al., Gaussian Splatting SLAM, CVPR 2024]

# Gaussian Splatting

# Hybrid NeRF / GS: RadSplat



ZipNeRF (0.25FPS) | Ours (788FPS)

[Niemeyer et al., RadSplat: Radiance Field-Informed Gaussian Splatting for Robust Real-Time Rendering with 900+ FPS, Arxiv 2024]

# Neural Implicit Representations



Offline Methods

**ConvONet** [Peng et al., ECCV'20]          **NeRF** [Mildenhall et al., ECCV'20]

# Simultaneous Localization and Mapping



**Tracking & Mapping**

SLAM → Camera Pos

SLAM → Map

# Simultaneous Localization and Mapping



**Tracking & Mapping**

SLAM → Camera Pos

SLAM → Map

*Radiance Field*

"Loss-less"

Compression ≈

No storage

# Neural Implicit SLAM: iMAP

Imperial College London

# Dense SLAM with a Neural Implicit Scene Represenation

RGB-D Sequences

[Zhu et al., NICE-SLAM, CVPR 2022]

# Dense SLAM with a Neural Implicit Scene Represenation

[Zhu, Peng, Larsson, Xu, Bao, Cui, Oswald, Pollefeys, CVPR'22]



**Input Depth**

**Input RGB**

Hierarchical Feature Grid

Tri-linear Interpolation

$o_{\mathbf{p}}^0$
Coarse-level Occupancy   Coarse Level   $\phi_\theta^0(\mathbf{p})$

Volume

**Mapping**

Fine Level Occupancy
$o_{\mathbf{p}}$   $+$   $o_{\mathbf{p}}^1$   Mid Level   $\phi_\theta^1(\mathbf{p})$

$\Delta o_{\mathbf{p}}^1$   $\phi_\theta^1(\mathbf{p})$

Fine Level   $\phi_\theta^2(\mathbf{p})$

Color

**Tracking**

$\psi_\omega(\mathbf{p})$

$\mathbf{p}$

Ray –> Point Sampler

Camera Pose

[Zhu et al., NICE-SLAM, CVPR 2022]

# NICER-SLAM: RGB-only SLAM

# NICER-SLAM: RGB-only SLAM



GT     NICE-SLAM     Vox-Fusion

COLMAP     DROID-SLAM     NICER-SLAM

[Zhu, Peng, Larsson, Cui, Oswald, Geiger, Pollefeys, NICER-SLAM , Arxiv 2023]

# Gaussian-SLAM: Dense SLAM with Gaussian Splatting

Radience field sampling & feature aggregation

Set of Gaussians encodes geometry and color

# Gaussian-SLAM: Dense SLAM with Gaussian Splatting

ESLAM

Point-SLAM

Gaussian-SLAM

GT

# MAGiC-SLAM: Multi-Agent Gaussian SLAM

# Splat-SLAM

## Globally Optimized RGB-only SLAM with 3D Gaussians



Fixed Neural Network | Optimization Layer | Memory Buffer | Graph Structure

# ☀️ Results: Rendering on TUM-RGBD



GlORIE-SLAM      MonoGS      Ours      Ground Truth

# Results: Color & Depth Rendering on Replica

# Results: Rendering on ScanNet

| Method | Metric | 0000 | 0059 | 0106 | 0169 | 0181 | 0207 | Avg. |
|---|---|---|---|---|---|---|---|---|
| *RGB-D Input* | | | | | | | | |
| SplaTaM [24] | PSNR↑ | 19.33 | 19.27 | 17.73 | 21.97 | 16.76 | 19.80 | 19.14 |
| | SSIM ↑ | 0.66 | 0.79 | 0.69 | 0.78 | 0.68 | 0.70 | 0.72 |
| | LPIPS↓ | 0.44 | 0.29 | 0.38 | 0.28 | 0.42 | 0.34 | 0.36 |
| MonoGS [38] | PSNR↑ | 18.70 | 20.91 | 19.84 | 22.16 | 22.01 | 18.90 | 20.42 |
| | SSIM ↑ | 0.71 | 0.79 | 0.81 | 0.78 | 0.82 | 0.75 | 0.78 |
| | LPIPS↓ | 0.48 | 0.32 | 0.32 | 0.34 | 0.42 | 0.41 | 0.38 |
| Gaussian-SLAM [74] | PSNR↑ | 28.54 | 26.21 | 26.26 | 28.60 | 27.79 | 28.63 | 27.67 |
| | SSIM ↑ | **0.93** | **0.93** | **0.93** | **0.92** | **0.92** | **0.91** | **0.92** |
| | LPIPS↓ | 0.27 | 0.21 | 0.22 | 0.23 | 0.28 | 0.29 | 0.25 |
| *RGB Input* | | | | | | | | |
| GO-SLAM [79] | PSNR↑ | 15.74 | 13.15 | 14.58 | 14.49 | 15.72 | 15.37 | 14.84 |
| | SSIM ↑ | 0.42 | 0.32 | 0.46 | 0.42 | 0.53 | 0.39 | 0.42 |
| | LPIPS↓ | 0.61 | 0.60 | 0.59 | 0.57 | 0.62 | 0.60 | 0.60 |
| MonoGS [38] | PSNR↑ | 16.91 | 19.15 | 18.57 | 20.21 | 19.51 | 18.37 | 18.79 |
| | SSIM ↑ | 0.62 | 0.69 | 0.74 | 0.74 | 0.75 | 0.70 | 0.71 |
| | LPIPS↓ | 0.70 | 0.51 | 0.55 | 0.54 | 0.63 | 0.58 | 0.59 |
| GlORIE-SLAM* [75] | PSNR↑ | 23.42 | 20.66 | 20.41 | 25.23 | 21.28 | 23.68 | 22.45 |
| | SSIM ↑ | 0.87 | 0.87 | 0.83 | 0.84 | 0.91 | 0.76 | 0.85 |
| | LPIPS↓ | 0.26 | 0.31 | 0.31 | 0.21 | 0.44 | 0.29 | 0.30 |
| **Splat-SLAM (Ours)** | PSNR↑ | **28.68** | **27.69** | **27.70** | **31.14** | **31.15** | **30.49** | **29.48** |
| | SSIM ↑ | 0.83 | 0.87 | 0.86 | 0.87 | 0.84 | 0.84 | 0.85 |
| | LPIPS ↓ | **0.19** | **0.15** | **0.18** | **0.15** | **0.23** | **0.19** | **0.18** |

# Results: Reconstruction on Replica

| | GlORIE-SLAM | MonoGS | Ours | Ground Truth |

| Metrics | NeRF-SLAM [62] | DIM-SLAM [28] | GO-SLAM [79] | NICER-SLAM [81] | HI-SLAM [78] | MoD-SLAM* [80] | GlORIE-SLAM* [75] | Mono-GS[38] | Q-SLAM* [46] | Ours |
|---|---|---|---|---|---|---|---|---|---|---|
| Render Depth L1↓ | 4.49 | - | - | - | - | - | - | 27.24 | 2.76 | **2.41** |
| Accuracy ↓ | - | 4.03 | 3.81 | 3.65 | 3.62 | 2.48 | 2.96 | 30.61 | - | **2.43** |
| Completion ↓ | - | 4.20 | 4.79 | 4.16 | 4.59 | - | 3.95 | 12.19 | - | **3.64** |
| Comp. Rat. ↑ | - | 79.60 | 78.00 | 79.37 | 80.60 | - | 83.72 | 40.53 | - | **84.69** |

# Deblur Gaussian SLAM



Blurry Input

Reconstruction

# Deblur Gaussian SLAM

# Language and 3D



[LERF: Language Embedded Radiance Fields, ICCV 2023]

# Open-vocabulary Online SLAM

# Auto-Vocabulary Segmentation

**Fixed-Vocabulary Segmentation**

Image → Segmented Image ← Predefined Dataset Classes

**Known & Fixed Vocabulary**

**Open-Vocabulary Segmentation**

Image → Segmented Image ← User-Specified Classes

**Known & Open Vocabulary**

**Auto-Vocabulary Segmentation**

Image → Generated Classes → Segmented Image

**Unknown & Open Vocabulary**

Auto-Vocabulary

Fixed-Set Ground Truth

# 3D Auto-Vocabulary Segmentation for LiDAR



**Human Annotation**

**AutoVoc3D**

# Conclusion and Take-away

- 3D / 4D computer vision algorithms train faster and require less training data (vs. 2D)
- 3D modeling, but 2D supervision
- Scene understanding requires memory
- Photographic and deformable memory improves accuracy & enables new applications
- Self-supervised learning via re-rendering error minimization
- Scene representation is important (local updates, deformable, catastrophic forgetting)
- SLAM can be a useful stepping stone for continual scene understanding

# Future Directions

- Beyond semantics: multi-modal output open-vocabulary & foundation models

- 3D-Language maps and spatial language-based reasoning

- Learning and controlling forgetting (keeping track of task-relevant changes)

- Collaborative / distributed asynchronous learning with multiple agents

- Physics-based scene representations (metric units, weights, gravity, etc.)

- 3D generative multi-modal models

- Dynamic scenes and temporal representations