Neural Implicit Representations for 3D Vision and Beyond

> Prof. Dr.-Ing. Andreas Geiger Autonomous Vision Group





1 Implicit Neural Representations

2 Differentiable Volumetric Rendering

3 Neural Radiance Fields

4 Generative Radiance Fields

5 Further Applications

1 Implicit Neural Representations

Traditional 3D Reconstruction Pipeline



3D Datasets and Repositories



[Newcombe et al., 2011]



[Wu et al., 2015]



[Choi et al., 2011]



[Chang et al., 2015]



[Dai et al., 2017]



[Chang et al., 2017]

Can we learn 3D Reconstruction from data?







Input Images

Neural Network

3D Reconstruction

What is a good **output** representation?

3D Representations



Occupancy Networks

Key Idea:

- ► Do not represent 3D shape explicitly
- Instead, consider surface implicitly as decision boundary of a non-linear classifier:



Remarks:

• The function f_{θ} models an **occupancy field**





Occupancy Networks

Key Idea:

- ► Do not represent 3D shape explicitly
- Instead, consider surface implicitly as decision boundary of a non-linear classifier:



Remarks:

- The function f_{θ} models an **occupancy field**
- ► Also possible: signed distance field [Park et al., 2019]



Network Architecture



Training Objective

Occupancy Network:

$$\mathcal{L}(\theta, \psi) = \sum_{j=1}^{K} \mathsf{BCE}(f_{\theta}(p_{ij}, z_i), o_{ij}) + KL \left[q_{\psi}(z | (p_{ij}, o_{ij})_{j=1:K}) \| p_0(z) \right]$$

- K: Randomly sampled 3D points (K = 2048)
- ► BCE: Cross-entropy loss

Training Objective

Variational Occupancy Encoder:

$$\mathcal{L}(\theta, \psi) = \sum_{j=1}^{K} \mathsf{BCE}(f_{\theta}(p_{ij}, z_i), o_{ij}) + KL \left[q_{\psi}(z | (p_{ij}, o_{ij})_{j=1:K}) \| p_0(z) \right]$$

- K: Randomly sampled 3D points (K = 2048)
- ► BCE: Cross-entropy loss
- ► q_{ψ} : Encoder

Occupancy Networks



Multiresolution IsoSurface Extraction (MISE):

- Build octree by incrementally querying the occupancy network
- Extract triangular mesh using marching cubes algorithm (1-3 seconds in total)

Results



Representing View-Dependent Appearance



Representing Motion



- Extending Occupancy Networks to 4D is hard (curse of dimensionality)
- Represent **shape** at t = 0 using a 3D Occupancy Network
- ► Represent **motion** by temporally and spatially continuous vector field
- ► Relationship between 3D trajectory s and velocity v given by (differentiable) ODE:

$$\frac{\partial \mathbf{s}(t)}{\partial t} = \mathbf{v}(\mathbf{s}(t), t)$$

Limitations



Structure of implicit neural representations:

- ► Global latent code ⇒ no local information, overly smooth geometry
- ► Fully connected architecture ⇒ does not exploit translation equivariance

Limitations

Implicit models work well for simple objects but poorly on complex scenes:



Convolutional Occupancy Networks



2 Differentiable Volumetric Rendering

DVR: Differentiable Volumetric Rendering



Forward Pass (Rendering)

DVR: Differentiable Volumetric Rendering

Forward Pass:

- $\blacktriangleright\,$ For all pixels ${\bf u}$
- Find surface point p̂ along ray w via ray marching and root finding (secant method)
- Evaluate texture field $\mathbf{t}_{\theta}(\hat{\mathbf{p}})$ at $\hat{\mathbf{p}}$
- ► Insert color $\mathbf{t}_{\theta}(\hat{\mathbf{p}})$ at pixel \mathbf{u}



Backward Pass (Differentiation)

DVR: Differentiable Volumetric Rendering

Backward Pass:

- ► Image Observation I
- \blacktriangleright Loss $\mathcal{L}(\mathbf{\hat{I}},\mathbf{I}) = \sum_{\mathbf{u}} \|\mathbf{\hat{I}_u} \mathbf{I_u}\|$
- ► Gradient of loss function:

$$\begin{array}{lll} \frac{\partial \mathcal{L}}{\partial \theta} & = & \sum_{\mathbf{u}} \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{I}}_{\mathbf{u}}} \cdot \frac{\partial \hat{\mathbf{I}}_{\mathbf{u}}}{\partial \theta} \\ \frac{\partial \hat{\mathbf{I}}_{\mathbf{u}}}{\partial \theta} & = & \frac{\partial \mathbf{t}_{\theta}(\hat{\mathbf{p}})}{\partial \theta} + \frac{\partial \mathbf{t}_{\theta}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \cdot \frac{\partial \hat{\mathbf{p}}}{\partial \theta} \end{array}$$

• Differentiation of $f_{\theta}(\hat{\mathbf{p}}) = \tau$ yields:

$$\frac{\partial \hat{\mathbf{p}}}{\partial \theta} = -\mathbf{w} \left(\frac{\partial f_{\theta}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \cdot \mathbf{w} \right)^{-1} \frac{\partial f_{\theta}(\hat{\mathbf{p}})}{\partial \theta}$$



⇒ Analytic solution and no need for storing intermediate results

Results



24

3 Neural Radiance Fields

NeRF: Representing Scenes as Neural Radiance Fields



- ► Task: Given a set of images of a scene, render image from novel viewpoint
- ► Vanilla ReLU MLP that maps from location/view direction to color/density
- Conditioning on view direction allows for modeling view-dependent effects

Volume Rendering



Volume rendering works very similar to traditional ray tracing in graphics

- Shoot ray through scene, sample points, apply alpha composition to render pixel
- ► Fourier features improve visual fidelity (i.e., texture details) significantly

NeRF Results





Reiser, Peng, Liao and Geiger: KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. ICCV, 2021.

4 Generative Radiance Fields



- Generative model for radiance fields
- Train from unstructured and unposed 2D image collections
- ► Challenges: Can't discriminate in 3D, but volumetric rendering is slow



- A radiance field g_{θ} maps a 3D point \mathbf{x}_{r}^{i} and viewing direction \mathbf{d}_{r} to color/density
- \blacktriangleright By sampling N points along the ray, we can **render** a pixel's color \mathbf{c}_r



- GRAF conditions the radiance field g_{θ} on additional latent codes z
- \blacktriangleright Here, \mathbf{z}_s is a shape latent code and \mathbf{z}_a is an appearance latent code



▶ The generator repeats this process for *R* rays, sparsely sampled on a 2D grid

• The camera intrinsics K, extrinsics ξ and 2D grid ν are drawn randomly



- ► This generates **image patches** of size 32 × 32 pixels
- ► The sampling pattern changes the **location** and **stride** (scale)



- Similarly to the generator, we can extract real patches of the same size
- Here I denotes a real image sampled from the data distribution $p_{\mathcal{D}}$



- We can now compare both patches using a 2D discriminator
- ► GRAF implements the 2D discriminator as a simple 4 layer ConvNet





Schwarz, Liao, Niemeyer, Geiger: GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. NeurIPS, 2020.

GIRAFFE: Compositional Generative Neural Feature Fields



5 Further Applications

COIN: Compression with Implicit Neural representations



Parameters of trained neural network act as compressed representation of data

Dupont, Golinski, Alizadeh, Teh and Doucet: COIN: COmpression with Implicit Neural representations. Arxiv, 2021.

SMD-Nets: Stereo Mixture Density Networks



► MLP decoder enables to query stereo disparity maps at arbitrary resolution

Tosi, Liao, Schmitt and Geiger: SMD-Nets: Stereo Mixture Density Networks. CVPR, 2021.

PointRend: Image Segmentation As Rendering



MLP decoder allows for accurate instance segmentation in Mask R-CNN

Kirillov, Wu, He and Girshick: PointRend: Image Segmentation As Rendering. CVPR, 2020.

iLabel: Interactive Neural Scene Labelling



Predicting geometry and semantics jointly enables efficient few-shot annotation

LENS: Localization enhanced by NeRF synthesis



NVS improves downstream tasks like localization via data augmentation

Neural Reflectance Fields for Appearance Acquisition



$$\sum_{j} \tau_{c}(\mathbf{x}_{j}) \tau_{l}(\mathbf{x}_{j}) (1 - \exp(-\sigma(\mathbf{x}_{j}) \Delta t_{j})) [f_{r}(\mathbf{x}_{j})] L_{l}(\mathbf{x}_{j})$$
Reflectance Model
$$f_{r}(\mathbf{x}_{j}, \boldsymbol{\omega}_{o}, \boldsymbol{\omega}_{i}, \boldsymbol{n}(\mathbf{x}_{j}), \boldsymbol{R}(\mathbf{x}_{j}))$$

$$x_{j} = (x, y, z)$$
Neural Reflectance Field
$$R(\mathbf{x}_{j})$$

Beyond light fields: modeling material properties using coordinate-based models

MeshfreeFlowNet: Deep continuous space-time super-resolution



Space-time super-resolution of turbulent flows (Rayleigh-Bénard convection)

SNARF: Differentiable Forward Skinning for Animating Neural Shapes





Learned Canonical 3D Shape and Skinning Weights

Learning articulated human shape representations (geometry and skinning)

Chen, Zheng, Black, Hilliges and Geiger: SNARF: Differentiable Forward Skinning for Animating Non-Rigid Neural Implicit Shapes. ICCV, 2021.

DiffSDFSim: Differentiable Rigid-Body Dynamics With Implicit Shapes



► Learning differentiable rigid-body dynamics (mass, friction) with implicit shapes

Neural Descriptor Fields



► Few-shot generalization of manipulation tasks by predicting spatial descriptors

Simeonov, Du, Tagliasacchi, Tenenbaum, Rodriguez, Agrawal and Sitzmann: Neural Descriptor Fields: SE(3)-Equivariant Object Representations for Manipulation. Arxiv, 2021. 45

NEAT: Neural Attention Fields for End-to-End Autonomous Driving



Continuous state representation and waypoint prediction for self-driving

Summary

Neural Networks as Continuous Representations

Occupancy Networks

 $\begin{array}{l} [\text{Mescheder et al. 2019}] \\ (x,y,z) \rightarrow \text{occupancy} \end{array}$



Scene Representation Networks

[Sitzmann et al. 2019] $(x, y, z) \rightarrow$ latent vec. (color, dist)



DeepSDF [Park et al. 2019]

 $(x, y, z) \rightarrow \text{distance}$



Differentiable Volumetric Rendering

[Niemeyer et al. 2020] $(x, y, z) \rightarrow \text{color, occ.}$



Neural Radiance Fields

 $\begin{aligned} & [\mathsf{Mildenhall et al. 2020}] \\ & (x,y,z,\theta,\phi) \to \mathrm{color, \ density} \end{aligned}$



 $\begin{array}{l} \mbox{Generative Radiance Fields} \\ \mbox{[Schwarz et al. 2020]} \\ \mbox{(}x,y,z,\theta,\phi,{\bf z}) \rightarrow {\rm color,\ density} \end{array}$



Summary

Coordinate-Based Networks:

- ► Effective output representation for shape, appearance, material, motion
- ► No discretization, model arbitrary topology
- ► Can be learned from images via differentiable rendering
- ► Many applications: reconstruction, motion, view synthesis, robotics

However:

- Geometry must be extracted in post-processing step (1 sec for ONet)
- Extension to 4D not straightforward (curse of dimensionality)
- ► Fully connected architecture and global condition lead to oversmooth results
- ► Promising: Local features (ConvONet), Better input encoding (NeRF)

Thank you!

http://autonomousvision.github.io

