# SCONE:

## Surface Coverage Optimization in uNknown Environments by Volumetric Integration

Antoine Guédon, Pascal Monasse, Vincent Lepetit

LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, France

## Introduction

## 3D Scene Rendering from images

Yu et al., Plenoxels: Radiance Fields without Neural Networks (CVPR 2022)



. . .

3D Reconstruction: NeRF-related models, SfM and MVS (COLMAP), etc.

## 3D Scene Rendering from images

## **Questions:**

- 1. How do we acquire these images?
- 2. Can we learn how to do it?

- → Path Planning
- → Next Best View (NBV)



## Next Best View (NBV)

- → Long standing problem in robotics
- → Depth sensors are commonly used (video stream, so even with RGB we consider we have access to partial point clouds)
- → Classic approaches: Usually, discretization into known and unknown voxels, then ray-casting entropy metrics. Very slow, not suited to real-time exploration.
- → Surface Coverage is a common metric to evaluate NBV

## Next Best View (NBV)

What about learning-based approaches?

- → Volumetric approaches: Rely on voxelization. Suited to path planning in scenes, but not scalable and far less precision.
- → Surface approaches: Work directly on the dense point cloud. Still rely on a single global encoding to regress NBV => Only suited to floating, small scale, single objects inside a sphere.

## Our approach: SCONE

## Input

## **1. Partial point clouds**

gathered by a LiDAR-class sensor or reconstructed from an RGB video stream (see DM)

2. Camera poses

**3. Bounding Box** (if needed, to delimit exploration area)



SCONE: Surface Coverage Optimization in Unknown Environments by Volumetric Integration (NeurIPS 2022)

- 1. Make a prediction about unseen geometry in the scene
- 2. Predict visibility of surface points you can't see right now

NBV = Camera with the most new visible surface points (highest coverage gain)

## Predict unseen geometry

## Naive approach

- 1. Predict unseen surface
- 2. Compute surface coverage for any camera from this surface.

In practice, does not work in unknown environments.

Predicting a surface requires high confidence in prediction.

## The need for a volumetric, probabilistic mapping



Our needs:

1. Virtually infinite resolution

**Deep Implicit Function** (very trendy in CV)

2. Generalization and Scalability

Focus on local features rather than a single global encoding









Supervision: MSE between predicted occupancy probabilities and GT occupancy map

## Estimate surface coverage gain

## Maximizing Surface Coverage with Volumetric Integration

Surface Coverage is a surface metric. Yet, we want to integrate on a volumetric occupancy probability field.

- → "Relax" the definition of visibility to a volumetric visibility field
- → Asymptotically, maximizing the volumetric integral is equivalent to maximizing the surface integral

**Theorem 1.** Under the previous regularity assumptions on the volume  $\chi$  of the scene and its surface  $\partial \chi$ , there exist  $\mu_0 > 0$  and M > 0 such that for all  $\mu < \mu_0$ , and any camera  $c \in C$ :

$$\left|\frac{1}{|\chi|_V} \int_{\chi} g_c^H(\mu; x) dx - \mu \frac{|\partial \chi|_S}{|\chi|_V} G_H(c)\right| \le M \mu^2 , \qquad (6)$$

where  $|\chi|_V$  is the volume of  $\chi$ .

## Predicting Visibility Gain

## How?

We want to avoid ray-casting operations.

Instead, we use directional, attention-based features between proxy points!

## Predicting Visibility Gain

- 1. Sample Proxy Points in the scene, based on their occupancy probability
- 2. For any new camera pose *c*, we **predict a Visibility Gain Score** for each *probabilistic* proxy point



## Predicting Visibility Gain: Input





Camera  $c_{t}$ 

Proxy Points with Occupancy Probability

## Predicting Visibility Gain: Input



**Spherical Harmonics** encoding Camera History

Camera  $c_t$ 

Proxy Points with **Occupancy Probability** 

## Predicting Visibility Gain: Output Camera $c_{t+1}$ **Spherical Harmonics** encoding Visibility Gain in all directions

Proxy Points with Visibility Gain in direction of new camera

## Predicting Coverage Gain



**Supervision**: Softmax then KL Divergence between predicted volumetric coverage gains and GT surface coverage gains, over a distribution of cameras

## Results

## Single Object Reconstruction: Quantitative Results

	Categories seen during training									
Method	1	Airplane	Cabinet	Car	Chair	Lamp	Sofa	Table	Vessel	Mean
Random		0.745	0.545	0.542	0.724	0.770	0.589	0.710	0.674	0.662
Proximity Count	8	0.800	0.596	0.591	0.772	0.803	0.629	0.753	0.706	0.706
Area Factor [25]		0.797	0.585	0.587	0.751	0.801	0.627	0.725	0.714	0.698
NBV-Net [16]		0.778	0.576	0.596	0.743	0.791	0.599	0.693	0.667	0.680
PC-NBV [36]		0.799	0.612	0.612	0.782	0.800	0.640	0.760	0.719	0.716
SCONE (Ours)		0.827	0.625	0.591	0.782	0.819	0.662	0.792	0.734	0.729
			Cate	gories no	ot seen du	ring train	ning			
Method	Bus	Bed	Bookshelf	Bench	Guitar	Motor	rbike	Skateboard	Pistol	Mean
Random	0.609	0.619	0.695	0.795	0.795	0.6	72	0.768	0.614	0.694
Proximity Count	0.646	0.645	0.749	0.829	0.854	0.7	05	0.828	0.660	0.740
Area Factor	0.629	0.631	0.742	0.827	0.852	0.7	18	0.799	0.660	0.732
NBV-Net	0.654	0.628	0.729	0.824	0.834	0.7	10	0.825	0.645	0.731
PC-NBV	0.667	0.662	0 740	0.845	0 849	0.7	28	0.840	0.672	0.750
	0.007	0.002	0.740	0.040	0.017	0.7	-0	0.010	0.012	0.700



#### Colosseum



#### Pisa Cathedral



#### Fushimi Castle



#### Christ the Redeemer



#### **Eiffel Tower**



#### Alhambra Palace



#### Manhattan Bridge



#### Pantheon (Rome, not Paris hehe)





#### **Dunnotar Castle**



#### Statue of Liberty





#### **Bannerman Castle**



#### Neuschwanstein Castle



#### Natural History Museum, London

	Method				
3D scene	Random Walk	SCONE-Entropy	SCONE		
unnottar Castle Ianhattan Bridge	$\begin{array}{c} 0.355 \pm 0.106 \\ 0.405 \pm 0.089 \end{array}$	$\begin{array}{c} 0.456 \pm 0.041 \\ 0.361 \pm 0.065 \end{array}$	$\begin{array}{c} \textbf{0.739} \pm 0.050 \\ \textbf{0.685} \pm 0.034 \end{array}$	ള്	
Ihambra Palace eaning Tower	$\begin{array}{c} 0.384 \pm 0.086 \\ 0.286 \pm 0.122 \end{array}$	$\begin{array}{c} 0.437 \pm 0.047 \\ 0.415 \pm 0.023 \end{array}$	$\begin{array}{c} \textbf{0.567} \pm 0.031 \\ \textbf{0.542} \pm 0.026 \end{array}$	overa	
Neuschwanstein Castle Colosseum	$\begin{array}{c} 0.403 \pm 0.032 \\ 0.308 \pm 0.061 \end{array}$	$\begin{array}{c} 0.538 \pm 0.040 \\ 0.512 \pm 0.024 \end{array}$	$\begin{array}{c} \textbf{0.653} \pm 0.025 \\ \textbf{0.571} \pm 0.024 \end{array}$	face c	
Eiffel Tower Fushimi Castle	$\begin{array}{c} 0.495 \pm 0.062 \\ 0.584 \pm 0.078 \end{array}$	$\begin{array}{c} 0.741 \pm 0.017 \\ 0.802 \pm 0.022 \end{array}$	$\begin{array}{c} \textbf{0.762} \pm 0.020 \\ \textbf{0.841} \pm 0.027 \end{array}$	Sur	
Pantheon Bannerman Castle	$\begin{array}{c} 0.175 \pm 0.065 \\ 0.321 \pm 0.121 \end{array}$	$\begin{array}{c} 0.351 \pm 0.020 \\ \textbf{0.667} \pm 0.023 \end{array}$	$\begin{array}{c} \textbf{0.396} \pm 0.036 \\ 0.642 \pm 0.047 \end{array}$	C	
Christ the Redeemer Statue of Liberty	$\begin{array}{c} 0.600 \pm 0.146 \\ 0.469 \pm 0.075 \\ 0.147 \pm 0.024 \end{array}$	$\begin{array}{c} 0.839 \pm 0.038 \\ 0.681 \pm 0.018 \\ 0.080 \pm 0.012 \end{array}$	$0.859 \pm 0.022 \\ 0.693 \pm 0.032 \\ 0.177 \pm 0.021$		
Mean	$\frac{0.147 \pm 0.024}{0.380}$	$0.080 \pm 0.010$ 0.529	$0.177 \pm 0.031$ 0.625		



# Thank you!