

# Real-world Visual Navigation in a Simulator: A New Benchmark

Henghui Bao\*, Kiran Lekkala\*, Laurent Itti  
University of Southern California  
Los Angeles, CA 90089, USA

## Abstract

In this paper, we explore advanced techniques in novel view rendering, particularly Gaussian Splatting, to create a simulator using a large-scale outdoor dataset. Our simulator, *Beogym*, is data-driven and built from data collected using a mobile robot. Our proposed pipeline processes the dataset to obtain an interconnected sequence of Gaussian splat files. These are then used by an engine to load appropriate splat files and render image frames during simulation. *Beogym* offers first-person view imagery, facilitating realistic training environments that could be used for enhancing and evaluating the learning capabilities of autonomous agents for visual navigation. It incorporates a sophisticated motion model and a sequence graph for seamless querying and loading of different sectors of the environment. The result closely resembles real-world navigation through smooth transitions across splat files.

## 1. Introduction

Learning to navigate without map is a task designed to enable agents to mimic human-like goal-oriented behaviors, relying solely on visual observations. Simulators are widely used in practice to seamlessly enable the agent to learn such behaviors. However, many recent works [1] attempt to fulfill only a few aspects out of the following; simulators that have photorealistic rendering, high performance, efficient utilization of compute resources and real-world transferability. Our method aims fulfill all the above requirements through advanced techniques in novel view rendering, such as Neural Radiance Fields (NeRF) [5] or Gaussian Splatting [3]. By interpolating these intermediate views, we seek to develop a real-time simulator that would not only facilitate more effective learning and navigation for robotic agents

\*Equal Contribution. Corresponding author: henghuib@usc.edu

This work was supported by the National Science Foundation (award 2318101), C-BRIC (one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA) and the Army Research Office (W911NF2020053). The authors affirm that the views expressed herein are solely their own, and do not represent the views of the United States government or any agency thereof.

but also enhance their applicability in real-world environments.

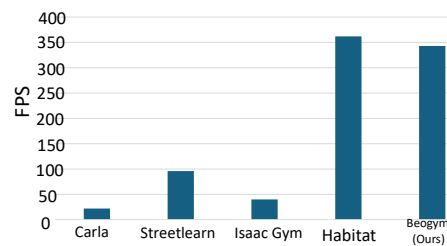


Figure 1. **Comparison with other simulators.** Rendering speed or Frames per Second (FPS) recorded for a single thread process with frame resolution of  $1280 \times 720$ , single episode  $\sim 200$  timesteps.

## 2. Proposed Simulator

We propose *Beogym*, a real-time simulator that allows an autonomous agent capable of navigating in an environment. After the agent executes an action/control signal  $\mathbf{u}_t$ , given an input image  $\mathbf{I}_t$  the simulator computes the pose  $\mathbf{x}_{t+1}$  at the next timestep using a motion model and renders an image, by querying from the splat file. One of the crucial aspects of a simulator apart from realistic quality is high performance in terms of rendering speed. Our method ensures that the visual feedback provided to the agent is both realistic and computationally efficient, facilitating more effective training and navigation. We compare the performance of our simulator with other SoTA simulators as shown in Figure 2.

### 2.1. Gaussian splat based rendering

The dataset that we use for our simulator consists of images and pointcloud data collected on a mobile robot that is collected across the USC campus [4]. In our methodology, each *sector*, defined as a splat file that is trained using a segmented portion of a trajectory. To optimize the efficiency and performance of reconstruction, we subdivided each session into several sectors, each containing 1250 im-

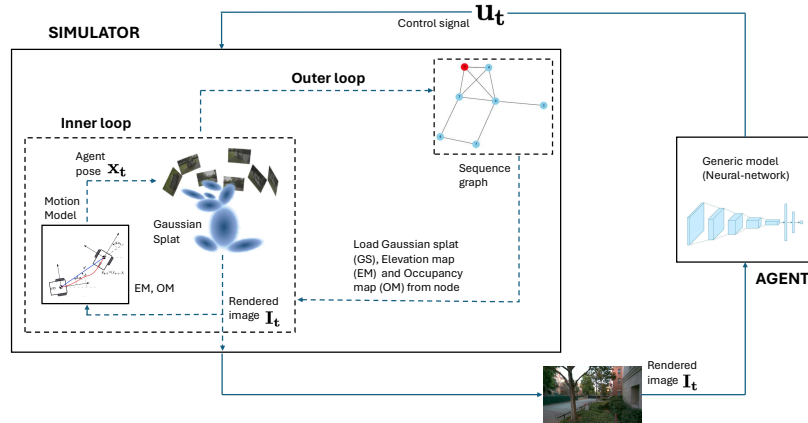


Figure 2. **Overview of our simulator.** The agent obtains a percept/image from the simulator and estimates a control signal  $u_t$ . The outer loop in the simulator determines whether the agent has passed the boundaries of the current *sector* and if a new splat file has to be loaded using the sequence graph. The inner loop corresponds to the motion model that computes the pose  $x_t$  that is then used to render the percept in the next timestep using the gaussian splat file.

056 ages captured from the five cameras. For every sector, we  
 057 utilize the collection of images and use COLMAP [7] to au-  
 058 tomatically obtain poses. This pose-annotated set of images  
 059 is passed as an input for training a splat file using Gaussian  
 060 Splatting [3]. The result would then be obtained as an ex-  
 061 plicit 3D representation of a specific scene or a *sector*. Sub-  
 062 sequently, we store these splat files, which are later used for  
 063 rendering purposes.

## 064 2.2. Sequence Graph for querying splat files

065 To facilitate navigation between sectors, we constructed a  
 066 sequence graph based on discrete image poses. In this  
 067 graph, each sector serves as a node corresponding to a splat  
 068 file and other related metadata. The connectivity between  
 069 sectors is represented by an edge that corresponds to a trans-  
 070 formation matrix  $T_a^b$  from a splat file  $a$  to the other  $b$ .

071 As the agent navigates within the simulation environ-  
 072 ment, a key element of realism is its interaction between  
 073 different sectors. When the agent approaches the boundary  
 074 of its current sector, the simulator recognizes this transition  
 075 and initiates the process of loading the appropriate splat file  
 076 for the new sector. This mechanism ensures a seamless vi-  
 077 sual experience as the agent moves through diverse parts of  
 078 the simulated environment.

## 079 2.3. Motion model

080 In the *Beogym* simulator, given a specific control signal, the  
 081 agent’s subsequent pose  $x_t$  at timestep  $t$  is governed by the  
 082 motion model. This model is crucial for simulating realis-  
 083 tic navigation behaviours, akin to those exhibited by actual  
 084 robots in real-world scenarios. One key challenge addressed  
 085 in our model is the maintenance of consistent height, raw  
 086 and yaw orientation relative to the ground by the agent. This

087 is crucial for ensuring realism and practical applicability,  
 088 especially environments where the terrain may vary, poten-  
 089 tially leading to the agent “floating” above the ground or  
 090 colliding with terrain features. To address this, we employ  
 091 a sophisticated method to compute the  $z$ -axis component or  
 092 the elevation of the agent’s pose accurately, using elevation  
 093 and occupancy maps obtained from the LiDAR data.

094 Furthermore, our simulation is designed to be adaptable  
 095 and responsive. An agent can be initialized at any location  
 096 within the simulation environment determined by the oc-  
 097 cupancy map. As the agent moves, guided by the motion  
 098 model, its pose is continually updated, and the rendering  
 099 process adapts accordingly, thus creating a seamless and  
 100 continuous simulation experience. This iterative process,  
 101 where the motion model predicts the next pose and the sim-  
 102 ulator renders the new view, forms the core inner-loop of  
 103 our simulation as shown in Figure 2.

104 As stated before, Within each sector, we employ  
 105 COLMAP [6] to obtain image poses that are then used for  
 106 training a splat file. However, it’s important to note that  
 107 these poses do not represent ground-truth poses obtained  
 108 from the LiDAR sensor, and merely are used for training a  
 109 Gaussian splat file. To utilize elevation maps derived from  
 110 ground-truth point clouds, we must transform the coordi-  
 111 nate system of the poses from COLMAP to those of the  
 112 elevation maps. We employ the *Kabsch* algorithm [2] to  
 113 compute this coordinate transformation. The ground truth  
 114 poses obtained from LiDAR in each sector are less dense  
 115 than image data and we pair these poses with images having  
 116 the closest timestamps. After forming all pairs, the *Kabsch*  
 117 algorithm is executed to derive optimal translation and ro-  
 118 tation matrices that minimize the root mean square (RMS)  
 119 deviation between the two sets of points.

120 **References**

- 121 [1] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio  
122 Lopez, and Vladlen Koltun. CARLA: An open urban driving  
123 simulator. In *Proceedings of the 1st Annual Conference on*  
124 *Robot Learning*, pages 1–16, 2017. 1
- 125 [2] Wolfgang Kabsch. A solution for the best rotation to relate  
126 two sets of vectors. *Acta Crystallographica Section A: Crys-*  
127 *tall Physics, Diffraction, Theoretical and General Crystallog-*  
128 *raphy*, 32(5):922–923, 1976. 2
- 129 [3] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and  
130 George Drettakis. 3d gaussian splatting for real-time radiance  
131 field rendering. *ACM Trans. Graph.*, 42(4):139:1–139:14,  
132 2023. 1, 2
- 133 [4] Laurent Itti Kiran Lekkala\*, Henghui Bao\*. Uscilab3d:  
134 A large-scale, long-term, semantically annotated outdoor  
135 dataset, 2024. 1
- 136 [5] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik,  
137 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:  
138 Representing scenes as neural radiance fields for view syn-  
139 thesis. In *Computer Vision - ECCV 2020 - 16th European*  
140 *Conference, Glasgow, UK, August 23-28, 2020, Proceedings,*  
141 *Part I*, pages 405–421. Springer, 2020. 1
- 142 [6] Johannes L. Schönberger and Jan-Michael Frahm. Structure-  
143 from-motion revisited. In *2016 IEEE Conference on Com-*  
144 *puter Vision and Pattern Recognition, CVPR 2016, Las Vegas,*  
145 *NV, USA, June 27-30, 2016*, pages 4104–4113. IEEE Com-  
146 puter Society, 2016. 2
- 147 [7] Johannes Lutz Schönberger and Jan-Michael Frahm.  
148 Structure-from-motion revisited. In *Conference on Computer*  
149 *Vision and Pattern Recognition (CVPR)*, 2016. 2