

ABGICP: Photo-realistic Dense SLAM with Gaussian Splatting

<https://spla-tam.github.io>

Atticus Zhou, Atticus Zhou, Atticus Zhou, Atticus Zhou, Atticus Zhou, Atticus Zhou, Elias D. Striatum

July 14, 2024

ABSTRACT

We present a dense simultaneous localization and mapping (SLAM) method that uses 3D Gaussians as a scene representation. Our approach enables interactive-time reconstruction and photo-realistic rendering from real-world single-camera RGBD videos. To this end, we propose a novel effective strategy for seeding new Gaussians for newly explored areas and their effective online optimization that is independent of the scene size and thus scalable to larger scenes. This is achieved by organizing the scene into sub-maps which are independently optimized and do not need to be kept in memory. We further accomplish frame-to-model camera tracking by minimizing photometric and geometric losses between the input and rendered frames. The Gaussian representation allows for high-quality photo-realistic real-time rendering of real-world scenes. Evaluation on synthetic and real-world datasets demonstrates competitive or superior performance in mapping, tracking, and rendering compared to existing neural dense SLAM methods.

1 Introduction

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important. **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

2 Related Work

in the image, [1] typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important. **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** [2] The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

3 Method

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important. **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

3.1 Gaussian Splatting

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by [3] the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important. **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

$$\text{RMSE}_{eT} = \sqrt{\frac{1}{N} \sum_{i=1}^N (\|t_i - t'_i\|)^2}$$

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important . **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important . **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

4 Experiments

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important . **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

5 Conclusion

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units [4] used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention

something important . **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

in the image, typically given in pixels. f_x, f_y : Focal lengths of the camera in the x and y directions, respectively, typically given in pixels. d : Depth value at the pixel (u, v) [5] after scaling, given in meters. $I_d(u, v)$: Depth value from the *original depth image* at pixel (u, v) , typically given in the units used by the depth sensor. s : Scale factor to convert the depth values from the depth image units to meters. \mathbf{p} : 3D point in *camera coordinates*, represented as a vector. Keetha and colleagues mention something important . **Concatenate coordinates to form the 3D points with a homogeneous coordinate:** The final array, points, is reshaped to $(-1, 4)$ to flatten the point cloud into a two-dimensional array where each row represents a 3D point in homogeneous coordinates.

- [1] Y. Cai, W. Xu, and F. Zhang, “IkD-Tree: An Incremental K-D Tree for Robotic Applications.” Accessed: May 21, 2024. [Online]. Available: <http://arxiv.org/abs/2102.10808>
- [2] N. Keetha *et al.*, “SplaTAM: Splat, Track & Map 3D Gaussians for Dense RGB-D SLAM.” Accessed: Jun. 09, 2024. [Online]. Available: <http://arxiv.org/abs/2312.02126>
- [3] K. Koide, M. Yokozuka, S. Oishi, and A. Banno, “Voxelized gicp for fast and accurate 3d point cloud registration,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2021, pp. 11054–11059. doi: 10.1109/ICRA48506.2021.9560835.
- [4] S. Ha, J. Yeon, and H. Yu, “RGBD GS-ICP SLAM.” Accessed: May 23, 2024. [Online]. Available: <http://arxiv.org/abs/2403.12550>
- [5] M. Korn, M. Holzkothen, and J. Pauli, “Color supported generalized-ICP,” in *2014 International Conference on Computer Vision Theory and Applications (VISAPP)*, IEEE, 2014, pp. 592–599. Accessed: May 20, 2024. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/7295135/>