

ChunkFusion

A Learning-based RGB-D 3D Reconstruction Framework via Chunk-wise Integration

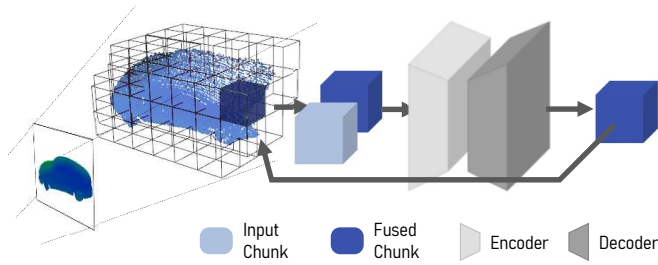
Chaozheng Guo, Lin Zhang, Ying Shen, Yicong Zhou
Tongji University, University of Macau



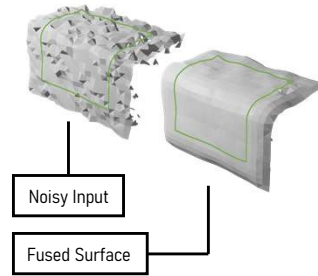
Introduction

In this paper, we devote our efforts to try to fill in the research gap in online RGB-D 3D reconstruction by proposing a scalable and robust RGB-D 3D reconstruction framework, namely ChunkFusion.

In ChunkFusion, sparse voxel management is exploited to improve the scalability of online reconstruction. Besides, a chunk-wise TSDF (truncated signed distance function) fusion network is designed to perform a robust integration of the noisy depth measurements on the sparsely allocated voxel chunks.



The overall pipeline of ChunkFusion. Based on the point cloud projected from the scanned depth map, the corresponding chunks will be updated by fusing the newly measured TSDF via a two-stage 3D convolutional network.



Example of chunk-wise surface fusing result.

Methods

For a frame scanned by an RGB-D camera with a known pose, ChunkFusion first allocates chunks according to the distribution of the point cloud. Then the newly allocated chunks storing the depth information will be fused individually with the historical chunks by a two-stage fusion network. Subsequently, the standard iso-surface mesh extraction will be conducted on the fused chunks. As a result, the global consistent 3D model of the scanned object can be yielded.

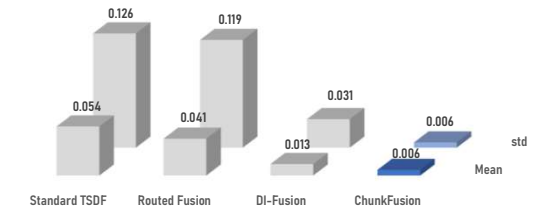
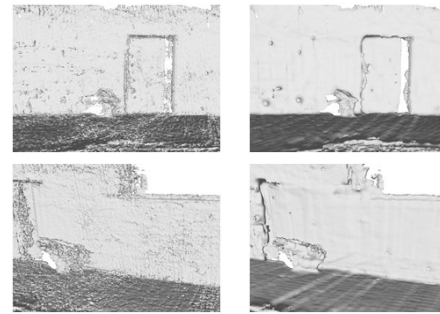
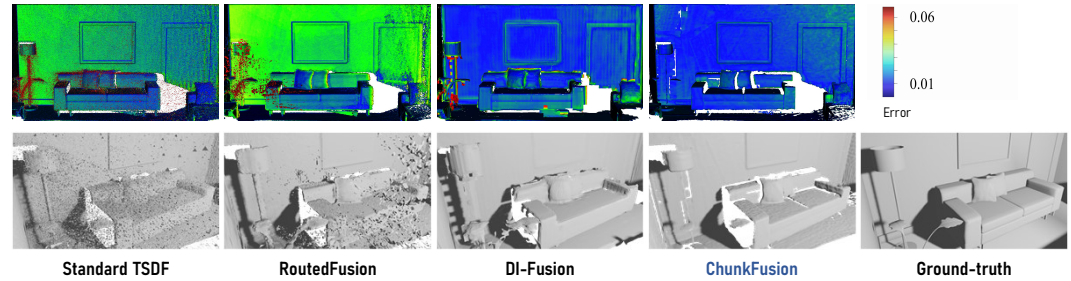
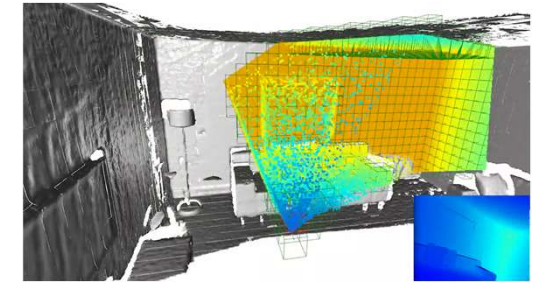
The loss function of the fusion network is defined as: $\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_{sign} + \mathcal{L}_{grad}$, where $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_{sign}, \mathcal{L}_{grad}$ are the L1 distance, the L2 distance, the gradient loss and the sign loss respectively. The gradient loss is introduced to restrict the 3D gradient of TSDF values. We also exploit the binary cross entropy to guarantee the correctness of the sign of TSDF values. The gradient loss and sign loss are defined as:

$$\mathcal{L}_{sign} = BCE(\text{sign}(\mathcal{F}(\hat{v}_i^t, v_i^{t-1})), \text{sign}(v_i^*)) \quad \mathcal{L}_{grad} = \sum_{j=x,y,z} \frac{1}{k^3} \|\nabla_j(\mathcal{F}(\hat{v}_i^t, v_i^{t-1})) - \nabla_j(v_i^*)\|_1$$

where $\mathcal{F}(\cdot)$ represents the TSDF fusion network, $\nabla_i(\cdot)$ returns the 3D Sobel gradient along the axis $x/y/z$, and v_i^* is the corresponding ground-truth TSDF value.

Results

Experiments are conducted on both synthetic dataset and self-collected real-world data. ChunkFusion can achieve lower error than state-of-the-art competitors. Defects caused by noise and outliers are shown more clearly in the figure below, in which a large number of fragments corrupt the reconstructed results. ChunkFusion manages to suppress the influence of these outliers and restore the smooth and accurate surface of the scene.



Mean surface error and standard deviation results on the ICL-NUIM dataset. ChunkFusion achieves better surface accuracy.

Conclusion

- ChunkFusion manages to employ the scalable voxel hashing scheme for the learning-based TSDF integration. Such a novel strategy eliminates the restriction of previous learning-based schemes and makes it possible to adapt the learning-based TSDF fusion to scenes with various scales.
- A two-stage fusion network is designed to perform the TSDF integration in an end-to-end manner. It is demonstrated that our fusion network can accurately restore the actual surfaces from noisy depth maps, yielding satisfactory reconstruction results both qualitatively and quantitatively.
- The proposed method fully exploits the sparsity of voxel representation by utilizing the chunk-wise fusion strategy and the sparsity-aware 3D convolutional network. Such an implementation scheme can further improve the computational efficiency and the surface quality of reconstruction.

Contact



GitHub Page

Email
gchaozheng@163.com