# **IBD-SLAM:** Learning Image-Based Depth Fusion for Generalizable SLAM

Minghao Yin<sup>1</sup> Shangzhe Wu<sup>2</sup> Kai Han<sup>1</sup>

<sup>1</sup>Visual AI Lab, The University of Hong Kong <sup>2</sup>Visual Geometry Group, University of Oxford

### **Background & Contribution**

**AiL** 

We address the problem of visual SLAM. We aim to develop a generalizable visual SLAM system that does not require network optimization during the mapping process.





Key contributions:

- We design IBD-SLAM, which can generalize to novel scenes without the need to retrain the model for scene-specific representation.
- We propose to utilize *xyz*-map for high-quality depth fusion.
- Our method runs  $10 \times$  faster than the previous state-of-the-art methods during the mapping stage.

#### Motivation

- Existing SLAM methods require high time consumption and increased memory usage as the scene scales up.
- Existing SLAM systems based on implicit representations require per-scene optimization for mapping, which limits the generalizability to unseen scenes.



### Method

• IBD-SLAM predicts the target view RGB image & *xyz*-map by fusing multi-view inputs from previous frames.



Depth maps suffer from multi-view inconsistency in camera coordinates. We propose to utilize *xyz*-map in world coordinates for fusing.

Tracking: Matched points on the reference and rendered images should share the same *xvz* values. By minimizing their xyz differences, we optimize the target camera pose. Mapping: The final colors and *xyz* values are obtained as a weighted sum of their correspondences in reference views.

#### ≻Ablation study

|                   | Depth L1 $\downarrow$ | Acc. $\downarrow$ | Comp. $\downarrow$ | (   |  |
|-------------------|-----------------------|-------------------|--------------------|-----|--|
| w/o novel         | 1.80                  | 2.66              | 2.91               |     |  |
| Shared-net        | 2.35                  | 3.02 3.8          |                    |     |  |
| Ours              | 1.53                  | 1.83              | 2.02               |     |  |
|                   | (a) Ablat             | ion study o       | of model des       | igr |  |
|                   | Depth L1              | ↓ Acc.↓           | Comp.↓             |     |  |
| iMAP              | 4.39                  | 4.77              | 5.02               |     |  |
| iMAP <sup>‡</sup> | 5.17                  | 6.19              | 6.87               |     |  |
| NICE-SLAM         | 2.49                  | 2.42              | 2.62               |     |  |
| NICE-SLAM         | <sup>‡</sup> 3.13     | 3.05              | 3.16               |     |  |
| Ourst             | 1.72                  | 2.05              | 2.30               |     |  |
| Ouro              |                       |                   |                    |     |  |

## **Qualitative results**

≻Reconstruction: Replica



► Reconstruction: Scannet





xvz values in 3D space xvz-mar





# **Quantitative Results**

▶ Reconstruction results & Time consumption

|           | Depth L1 $\downarrow$  | Acc. $\downarrow$      | Comp.↓                | Comp Ratio                     |           | Track↓[ms x it] | Map↓[ms x it] | #param . |
|-----------|------------------------|------------------------|-----------------------|--------------------------------|-----------|-----------------|---------------|----------|
| Orb-SLAM2 | $4.49/3.35^{\dagger}$  | $3.97/3.36^{\dagger}$  | $4.05/3.60^{\dagger}$ | $82.4/86.3^{\dagger}$          | NICE-SLAM | 7.8x10          | 82.5x60       | 17.4M    |
| NICE-SLAM | $13.55/2.49^{\dagger}$ | $2.87/2.42^{\dagger}$  | $3.13/2.65^{\dagger}$ | $87.1/90.3^{\dagger}$          | ESLAM     | 6.9x8           | 18.4x15       | 9.29M    |
| ESLAM     | 2.30/1.29 <sup>†</sup> | $2.82/2.34^{\dagger}$  | $2.97/2.14^{\dagger}$ | 89.5/ <b>94.7</b> <sup>†</sup> | Co-SLAM   | 5.8x10          | 9.8x10        | 0.26M    |
| Co-SLAM   | $2.59/1.60^{\dagger}$  | $2.66/2.21^{\dagger}$  | $3.21/2.36^{\dagger}$ | $88.9/92.7^{\dagger}$          | Ours      | 5.4x20          | 12.3x1        | 0.04M    |
| Ours      | $2.41/1.53^{\dagger}$  | 2.25/1.83 <sup>†</sup> | $2.93/2.02^{\dagger}$ | 90.9/93.8 <sup>†</sup>         |           |                 |               |          |

≻Tracking results on TUM-RGBD and Scannet datasets

|               | fr1/desk | fr2/xyz | fr3/office | Avg |               | 0000  | 0059   | 0106  | 0169  | Avg   |
|---------------|----------|---------|------------|-----|---------------|-------|--------|-------|-------|-------|
| MAP[35]       | 7.2      | 2.1     | 9.0        | 6.1 | iMAP[35]      | 55.95 | 32.06  | 17.50 | 70.51 | 44.00 |
| DI-Fusion[15] | 4.4      | 2.1     | 15.6       | 7.4 | DI-Fusion[15] | 66.99 | 128.00 | 18.50 | 75.80 | 72.32 |
| NICE-SLAM[53] | 2.7      | 1.8     | 3.0        | 2.5 | NICE-SLAM[53] | 8.64  | 12.25  | 8.09  | 10.28 | 9.89  |
| Ours          | 1.8      | 1.8     | 2.7        | 2.2 | Ours          | 7.96  | 9.19   | 7.13  | 7.98  | 9.44  |
|               |          |         |            |     |               |       |        |       |       |       |





| Ļ    | Comp. Ratio † |
|------|---------------|
|      | 91.2          |
|      | 87.7          |
|      | 93.8          |
| desi | ign           |
| p. ↓ | Comp. Ratio ↑ |
| )2   | 75.5          |
|      | <i>c</i> • •  |

In table (b), † denotes results without regularization losses. And ‡ denotes Poisson surface reconstruction results.

|     | Comp Ratio 1 |             | Depth L1 $\downarrow$ | Acc. $\downarrow$ | Comp. $\downarrow$ | Comp. Ratio ↑ |
|-----|--------------|-------------|-----------------------|-------------------|--------------------|---------------|
| · + | Comp. Rano   | w/o Caus    | 3.02                  | 2.81              | 3 23               | 88.2          |
| 2   | 75.5         | w/o Calanth | 2.09                  | 2.23              | 2.45               | 92.0          |
| 7   | 61.4         | w/o Creat   | 2.17                  | 2.35              | 2.64               | 91.6          |
| 2   | 90.3         | w/o Lnormal | 1.65                  | 1.98              | 2.20               | 92.7          |
| 6   | 87.2         | w/o Common  | 1.62                  | 1.92              | 2.13               | 93.3          |
| 0   | 92.2         | w/o C.com   | 1.57                  | 1.86              | 2.20               | 93.1          |
| 2   | 93.8         | Ours        | 1.53                  | 1.83              | 2.02               | 93.8          |
|     |              |             |                       |                   |                    |               |

(c) Ablation study of pretraining loss functions



➤Tracking: Replica

#### ▶ Reconstruction results with different fusion methods

xyz-map based fusion