# Supplementary Material for MultiEgo

## 1 MultiEgo Dataset Instruction

The dataset contains 5 scenes: talking, statement, concert, sword, and presentation. Each scene provide video, camera intrinsic, camera poses, timestamp, and a sparse point cloud of the first frame scene.

The file construction is as follows:

```
scene
|-cam1
| |-<scene>-cam1.mp4
| |-intrinsic.txt
| |-camera_poses.txt
| |-sampletime.txt
|-cam2
|-cam3
|-cam4
|-cam5
|-sparse
|-camera.bin
|-images.bin
|-points3D.bin
|-points3D.ply
```

where `<scene>-camx.mp4` is the egocentric video of the performer x in the scene. If frame extraction is performed on all videos, it is recommended to reserve 25 GB of storage space.

`intrinsic.txt` is the intrinsic matrix of the camera x, in the format as:

$$\begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \tag{1}$$

`camera_poses.txt` is the camera poses matrix of the frames in the `<scene>-camx.mp4`. The camera poses are represented as camera-to-world transformations in the world coordinate system. The pose in the format as:

$$\begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix} \tag{2}$$

`sampletime.txt` is the capture time of the acquisition system. The data in `sampletime.txt` is in the unit of nano-second.

The `sparse` directory contains COLMAP [2] binary files for all images, including intrinsic camera parameters (`camera.bin`) and world-to-camera extrinsic transformations (`images.bin`).

The `images.bin` file names follow the naming convention `camx_frame_00000.png`. Additionally, we provide sparse 3D point clouds reconstructed from the first frame's images and supplementary images, stored in `points3D.bin` and `points3D.ply`.

## 2 Data Processing Details

In the following part, we will explain the details of data annotation process. We assume that after a data acquisition, the $i$-th AR glasses acquires a sequence of image frames $X_i$, and a sequence of gyroscopic pose frames $G_i$.

### 2.1 Monocular Pose Tracking

As described in Section 3.2 in the paper, each image frame and gyroscopic pose frame has its own timestamp, with image frames captured at 30Hz and gyroscopic pose frames at 50Hz. To align these data streams, we perform Spherical Linear Interpolation (SLERP) on the gyroscopic pose frames to obtain rotation data $\hat{G}_i$ corresponding to the exact capture times of the image frames. Specifically, let $q_0$ and $q_1$ denote the quaternions at times $t_0$ and $t_1$, respectively. The interpolated quaternion $q$ at time $t \in (t_0, t_1)$ is given by:

$$q = q_0(q_0^{-1}q_1)^{\frac{t-t_0}{t_1-t_0}} \tag{3}$$

Then we employ several different image-based camera pose estimation methods to obtain multiple camera trajectories, in this paper we use Anycam [4], Mega-SAM [5], CUT3R [3], MonST3R [5] and PySLAM [1]. We let $P_{i,j}$ denote the $j$-th trajectory of $i$-th image frame sequence $X_i$, where the translation part is $t_{i,j}$ and the rotation part is $r_{i,j}$. It's notable that we . Subsequently, we fuse all the trajectories based on the rotation data $q$ obtained by SLERP. Specifically, we calculate the importance $m_j$ of $j$-th method based on the $L_1$ norms of the difference between $\hat{G}_i$ and $r_{i,j}$:

$$m_j = \frac{1}{\sum_i^I |r_{i,j}^{-1}\hat{G}_i|/I} \tag{4}$$

where $I$ is the number of AR glasses. We obtain the weight $w_j$ of the $j$-th method based on $m_j$:

$$w_j = \frac{m_j}{\sum_j^n m_n} \tag{5}$$

After the calculation above, a normalized monocular camera trajectory $\bar{P}_i$ of the $i$-th AR glasses is given by:

$$\bar{P}_i = (\sum_j^J w_j \cdot \frac{t_{i,j}}{\|t_{i,j,max}\|}, \frac{\sum_j^J w_j \cdot r_{i,j}}{\|\sum_j^J w_j \cdot r_{i,j}\|}) \tag{6}$$

where $\sum_j^J w_j \cdot \frac{t_{i,j}}{\|t_{i,j,max}\|}$ denotes the translation part, and $\frac{\sum_j^J w_j \cdot r_{i,j}}{\|\sum_j^J w_j \cdot r_{i,j}\|}$ denotes the rotation part. We abbreviate them as $\bar{t}_i$ and $\bar{r}_i$, respectively.

### 2.2 Multi-camera Pose Synthesis

Before data acquisition, we capture supplementary image sequence $X_s$ of the first frame static scene. We process the supplementary image sequence $X_s$ and the first frame of all the image frame sequence $X_i$ by SfM pipeline of COLMAP [2] to reconstruction a static scene. In this scene, we obtain the absolute pose of different AR glasses at first frame $P_{i,0}$. Then we add the images in $X_i$ which have the

max translation value, into the static scene to obtain the absolute pose of these images. We denote the displacement value between the first frame pose and the corresponding max translation pose as $\Delta t_{i,max}$. To scaling the normalized monocular trajectory $\bar{P}_i$ to the size of the static scene, we calculate a scale factor $s_i$:

$$s_i = \frac{\|\Delta t_{i,max}\|}{\|\bar{t}_{i,max}\|} \tag{7}$$

Then, based on normalized monocular pose $\bar{P}_i$ and scale factor $s_i$, the absolute pose sequence of $i$-th view $P_i$ is given by:

$$P_i = (t_{i,0} + s_i \cdot \bar{t}_i \cdot r_{i,0}, \quad r_{i,0} \cdot \bar{r}_i) \tag{8}$$

where $t_{i,0}$ and $r_{i,0}$ denotes the translation and rotation of first frame pose, $t_{i,0} + s_i \cdot \bar{t}_i \cdot r_{i,0}$ and $r_{i,0} \cdot \bar{r}_i$ represent the translation and rotation of the final absolute pose.

## 3 Consent Forms

Consent forms of performers are shown in figure 1.



**Figure 1: Consent Forms of Performers**

## References

[1] Luigi Freda. 2025. pySLAM: An open-source, modular, and extensible framework for SLAM. *arXiv preprint arXiv:2502.11955* (2025).

[2] Johannes Lutz Schönberger and Jan-Michael Frahm. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.

[3] Qianqian Wang*, Yifei Zhang*, Aleksander Holynski, Alexei A. Efros, and Angjoo Kanazawa. 2025. Continuous 3D Perception Model with Persistent State. In *CVPR*.

[4] Felix Wimbauer, Weirong Chen, Dominik Muhle, Christian Rupprecht, and Daniel Cremers. 2025. AnyCam: Learning to Recover Camera Poses and Intrinsics from Casual Videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

[5] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. 2024. MonST3R: A Simple Approach for Estimating Geometry in the Presence of Motion. *arXiv preprint arxiv:2410.03825* (2024).