# Self-Supervised Monocular 4D Scene Reconstruction for Egocentric Videos

**Chengbo Yuan[1,2,3], Geng Chen[2,4†], Li Yi[1,2,3], Yang Gao[1,2,3‡]**

[1]Institute for Interdisciplinary Information Sciences, Tsinghua University
[2]Shanghai Artificial Intelligence Laboratory
[3]Shanghai Qi Zhi Institute  [4]UC San Diego

ycb24@mails.tsinghua.edu.cn  gec001@ucsd.edu
{ericyi, gaoyangiiis}@mail.tsinghua.edu.cn

## Abstract

*Egocentric videos provide valuable insights into human interactions with the physical world, which has sparked growing interest in the computer vision and robotics communities. A critical challenge in fully understanding the geometry and dynamics of egocentric videos is dense scene reconstruction. However, the lack of high-quality labeled datasets in this field has hindered the effectiveness of current supervised learning methods. In this work, we aim to address this issue by exploring an self-supervised dynamic scene reconstruction approach. We introduce **EgoMono4D**, a novel model that unifies the estimation of multiple variables necessary for <u>Ego</u>centric <u>Mono</u>cular <u>4D</u> reconstruction, including camera intrinsic, camera poses, and video depth, all within a fast feed-forward framework. Starting from pretrained single-frame depth and intrinsic estimation model, we extend it with camera poses estimation and align multi-frame results on large-scale unlabeled egocentric videos. We evaluate EgoMono4D in both in-domain and zero-shot generalization settings, achieving superior performance in dense pointclouds sequence reconstruction compared to all baselines. EgoMono4D represents the first attempt to apply self-supervised learning for pointclouds sequence reconstruction to the label-scarce egocentric field, enabling fast, dense, and generalizable reconstruction. The interactable visualization, code and trained models are released https://egomono4d.github.io/.*

## 1. Introduction

Egocentric videos, especially Hand Object Interaction (HOI) videos, capture a vast amount of knowledge about human interaction with the physical world, particularly in
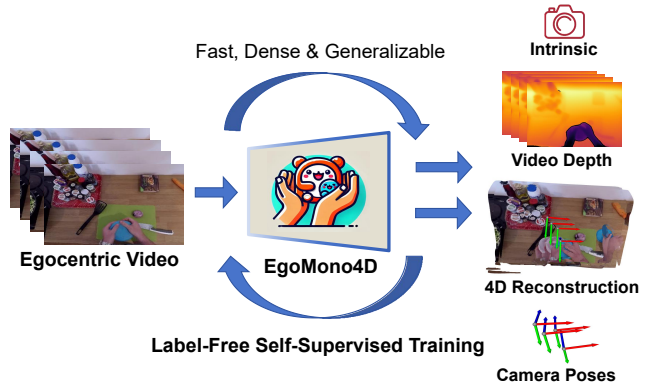


Figure 1. We propose EgoMono4D, a model that unifies the estimation of camera intrinsic, camera poses, and video depth for fast and dense 4D reconstruction of egocentric scenes. EgoMono4D is trained solely on large-scale unlabeled videos in an self-supervised learning framework.

tool usage. Due to this rich source of interaction knowledge, egocentric human videos have gained increasing interest from both the research community (e.g., computer vision [58, 100] and robotics [51, 96]) and industry (e.g., virtual reality [54]).

To better understand the geometry and dynamics in these activities, a crucial task is the dense 4D reconstruction from internet-scale egocentric video datasets [15, 23]. Here, dense 4D reconstruction is represented as **a pointclouds sequence which captures the 3D position of every pixel in each frame within a global coordinate system** [45, 84, 86, 99], preserving maximal geometry details. This task requires: (1) dense per-pixel reconstruction, (2) the ability to handle dynamic motions in egocentric videos. To facilitate large-scale usage [15, 90, 96], (3) strong generalization and fast processing speeds are also needed to adapt to large-scale unseen egocentric scenes.

However, current methods fail to meet these demands. Traditional methods [64, 77], such as Structure-from-

---

[1]† Work done during the internship at Shanghai AI Lab. ‡ The corresponding author.

Motion (SfM) or SLAM system combined with dense depth estimation [56, 95], face significant challenges in reconstructing dynamic scenes. These multi-step approaches also suffer from accumulated errors and inconsistencies between different modules [82, 86]. Test-time optimization techniques [32, 35, 46, 98], such as Gaussian Splatting [29], suffer from slow processing speeds, making them impractical for reconstructing large-scale egocentric datasets [23, 43]. Meanwhile, recent supervised learning approaches, such as DUSt3R [86] and MonSt3R [99], are limited by the scarcity of labeled egocentric videos, especially outside of controlled environments [15, 23].

To address these challenges, we explore an self-supervised learning approach for fast, dense, and generalizable reconstruction of highly dynamic egocentric videos. We propose **EgoMono4D**, a model trained without ground-truth labels that simultaneously predicts multiple variables necessary for dense 4D reconstruction, including camera intrinsic, camera poses, and video depth [42, 71].

Our key insight is to **extend pretrained single-frame scene reconstruction model to video version, and align multi-frame results with 4D constraints.** We begin by estimating per-frame depth and camera intrinsic to generate per-frame pointclouds predictions. Then, we align multi-frame results and derive camera extrinsics by minimizing the difference between (1) the 3D scene flow induced by the camera's motion through a static scene and (2) pre-computed 3D correspondences, similar to previous self-supervised methods [71–73, 94]. A confidence mask is used to exclude dynamic and unreliable areas during multi-frame alignment. To prevent model collapse and accelerate training convergence, we also regularize the model with predictions from state-of-the-art off-the-shelf models, such as Unidepth [56] for depth estimation and EgoHOS [100] for confidence mask prediction.

We evaluate EgoMono4D in both in-domain and zero-shot settings on unseen egocentric scenes. The model successfully recovers the 3D structure of scenes and the motion of dynamic parts, even in challenging synthetic surgery HOI videos [85]. We also provide a quantitative comparison with baseline methods that offer near-linear time complexity relative to the number of frames, focusing on two fundamental 4D tasks: dense pointclouds sequence reconstruction [71, 86] and long-term 3D scene flow recovery [3, 33, 96]. EgoMono4D demonstrates superior performance on evaluation metrics, outperforming all baseline methods.

In conclusion, our main contributions are as follows:

- We propose EgoMono4D, a model that unifies camera intrinsic, camera poses, and video depth estimation for fast, dense, and generalizable 4D reconstruction of egocentric videos.
- We introduce a novel self-supervised training method for egocentric scene reconstruction, training our model solely on large-scale, unlabeled monocular egocentric datasets, addressing the challenge of labeled data scarcity.
- EgoMono4D demonstrates promising performance in both in-domain and zero-shot unseen scenes, surpassing all baselines in pointclouds sequence reconstruction.

## 2. Related Works

### 2.1. (Self-Supervised) Monocular Depth Estimation

Monocular depth estimation has made significant progress in recent years [61]. Supervised learning models have shown strong generalization capabilities [4, 6, 25, 56, 88, 93, 95]. Our work builds on UniDepth [56], a state-of-the-art model that unifies camera intrinsic and depth estimation for single image. Another line of works focus on self-supervised training [5, 42, 69, 76, 94]. These methods train depth estimators purely on monocular videos using photometric error supervision [72, 73], often by leveraging camera labels or learning camera predictors. Despite recent progress, no methods have yet demonstrated strong generalization across both camera and depth. Our approach share similar intuition, which also unifies depth and camera prediction and trained solely on monocular video datasets [15] with photometric loss. However, our primary focus is on generalizable 4D reconstruction for egocentric scenes, which requires zero-shot prediction for both depth and camera parameters. To some extent, our work is the first attempt to extend self-supervised depth estimation to generalizable dense pointclouds sequence reconstruction.

### 2.2. Structure from Motion and SLAM Systems

Structure from Motion (SfM) [14, 64, 65, 103] and monocular visual SLAM systems [7, 49, 62, 77, 101] reconstruct 3D structures and estimate camera poses from image sequences. However, they struggle with dynamic scenes which pose ill condition for the epipolar constraint [86]. Moreover, most of them typically can not provide dense pixel-level reconstructions for all frames. Combining visual odometry [1, 8, 77, 87] with dense depth estimation [56, 95] helps constrain dynamic parts' geometry but can lead to accumulated and inconsistencies errors. Recently, FlowMap [71] offers a differentiable SfM for static scenes, optimizing depths, poses, and intrinsic simultaneously. We adopt its key ideas and extend it to a generalizable version for dynamic egocentric videos.

### 2.3. Dense 4D Reconstruction for Dynamic Scenes

Reconstructing 4D dynamic scenes remains a challenging problem in computer vision. Some approaches [12, 32, 35, 46, 75, 83, 97, 98, 102] first compute vision cues (e.g., camera pose, depth, optical flow) and then perform test-time optimization for each scene. While effective, these methods
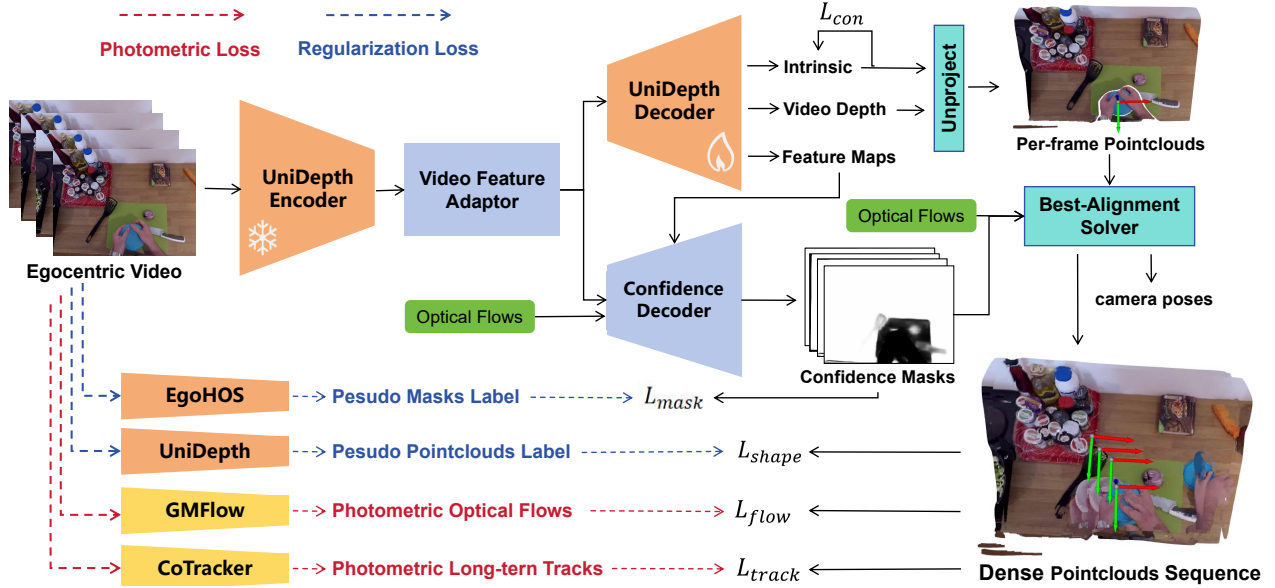
Figure 2. The overview of EgoMono4D and our self-supervised training framework. The model first simultaneously predicts camera intrinsic, video depth, and confidence maps (for camera pose estimation). Camera poses are then calculated by aligning unprojected pointclouds from different frames with confidence maps. The final dense pointclouds sequence reconstruction is assembled using all the predicted variables. We train our model purely on unlabeled egocentric video datasets, with both self-supervised photometric loss for depth alignment and regularization loss for training stablization.

are time-consuming and impractical for large-scale reconstructions [33].

Concurrent with our work, another line of research [17, 19, 26, 28, 38, 39, 45, 50, 81, 84, 86, 99] explores fast, feed-forward 4D reconstruction through supervised learning. Many of these methods are the extension of the pretrained end-to-end Multi-View Stereo (MVS) model DUSt3R [86]. For example, MonSt3R [99] and Stereo4D [26] fine-tune DUSt3R for dynamic scenes, while Spann3R [81] and CUT3R [84] incorporate memory buffers for temporal information merging to avoid post-optimization. Align3R [45] further enhances geometry estimation by merging depth priors. However, these methods rely on ground-truth labels for training, which are scarce for egocentric video datasets [34]. In this work, we aim to explore self-supervised methods [71, 72] for monocular egocentric scene reconstruction, leveraging large-scale unlabeled video datasets [15] instead.

### 2.4. Egocentric Video Understanding

Egocentric videos and datasets [15, 20, 22, 23, 34, 37, 43, 54, 85] capture human hand interactions with the environment, which are critical for robotics [2, 51, 90, 96], virtual reality [54], and intelligent agents [48, 89]. Tasks like detection [67], segmentation [100], intention recognition [37, 59], motion prediction [3, 60], and hand-object reconstruction [55, 57] are commonly used for egocentric videos. However, dense reconstruction of entire scenes remains a

challenge. EgoGaussian [98] solves this with time-intensive Gaussian Splatting [29], while our method provides a faster ($> 30$x speedup), feed-forward, and scalable solution.

## 3. Preliminary

### 3.1. Problem Definition

A monocular egocentric video can be represented as a sequence of frames $\{I_t \in R^{H \times W \times 3}\}$, where $t = 0, 1...T-1$ and $T$ is the total number of frames. Each frame is a RGB image with a resolution of $H \times W$. Given a video, our goal is to estimate a sequence of dense pointclouds [86], $\{\hat{S}_t \in R^{H \times W \times 3}\}$, which capture the 3D position of every pixel in each frame within a global coordinate system. To achieve this, we decompose the pointclouds sequence into: (1) video depth $\{\hat{D}_t \in R^{H \times W}\}$; (2) camera intrinsic $\hat{K} \in R^{3 \times 3}$ for unprojecting depth into per-frame pointclouds; and (3) camera poses $\{\hat{P}_t \in R^{4 \times 4}\}$ in camera-to-world format for projecting the per-frame pointclouds into global coordinates. Then the pointclouds sequence $\{\hat{S}_t\}$ could be calculated as:

$$\hat{X}_t = \hat{D}_t(i,j)\hat{K}^{-1}h(p(i,j)) \tag{1}$$

$$\hat{S}_t(i,j) = h^{-1}(\hat{P}_0^{-1}\hat{P}_t h(\hat{X}_t)) \tag{2}$$

where $(i, j)$ refers to the pixel position $p(i, j)$, $\hat{X}_t \in R^{H \times W \times 3}$ represents the per-frame pointclouds, and $h(\cdot)$

is the homogeneous operator that adds an extra dimension with a value of 1 to the coordinate system.

## 3.2. Camera Pose from Depth Alignment

Following FlowCam [70], we **reframe camera pose estimation as a depth alignment and confidence mask prediction problem.** This approach transforms the task of predicting sparse camera parameters into **a dense pixel-level prediction problem**, enhancing robustness and generalization [71]. Specifically, the model first predicts multi-frame depth $\hat{D}_t$ and camera intrinsic $\hat{K}$, and then unprojects the depths into per-frame pointcloudss $\hat{X}_t$. We use an off-the-shelf model [92] to compute the optical flow between adjacent frames, denoted as $\hat{u}_{i-1,i}$. The optimal camera pose transformation should best align the 3D point pairs induced by $\hat{u}_{i-1,i}$. Additionally, we predict a confidence mask $\hat{\mathcal{M}}_{i,i-1}$ for frame $i$ to exclude (1) dynamic regions, (2) occlusions, (3) scene edges, and (4) inaccuracies in optical flow during the alignment process.

Formally, let $\hat{X}_i^{\leftarrow i-1}$ denote the result of interpolating $\hat{X}_i$ using the points from $\hat{X}_{i-1}$ (this can be computed based on $\hat{u}_{i-1,i}$). The best-aligned camera pose transformation $\hat{P}_{i,i-1}$ can then be formulated as:

$$\hat{P}_{i,i-1} = \arg\min_{P \in SE(3)} ||\hat{\mathcal{M}}_{i,i-1}(\hat{X}_{i-1} - P\hat{X}_i^{\leftarrow i-1})|| \quad (3)$$

Then transformation $\hat{P}_{i,j}$ between arbitray frames $i$ and $j$ could be acquire by chain the nearby-frame results. Solving Equation 3 is known as the weighted procrustes-alignment problem [70], which can be solved in closed form using the singular value decomposition (SVD) [11], allowing a differentiable and learnable camera pose estimation process via gradient descent [71]. This means that we could get camera pose naturally by only focus on predicting video depth, camera intrinsic and confidence masks.

## 4. Methodology

### 4.1. Overview

We begin with the state-of-the-art pretrained single-frame scene reconstruction model, UniDepth [56], which predicts single-frame depth and intrinsic. Our goal is to extend it to a video-based model with label-free training. To achieve this, we need to (1) predict camera poses and (2) eliminate inconsistencies between multi-frame results.

To enable camera poses estimation, we adapt UniDepth from an image to a video estimator using adaptor blocks [9, 40]. We also introduce a new decoder to predict confidence masks, facilitating camera poses estimation as described in Section 3.2. To eliminate multi-frame inconsistencies, we employ self-supervised training losses based on 4D constraints to align multi-frame results, ensuring both temporal and spatial consistency in video predictions. The detailed approach is outlined in the following section.

### 4.2. Model Architecture

Our model aims to predict (1) video depth, (2) camera intrinsic, and (3) confidence masks. Camera poses are then derivated using multi-frame depth alignment following Section 3.2. Our model builds upon UniDepth [56], a universal estimator for predicting single-frame depth and camera intrinsic with an **encoder-decoder** architecture. More details about UniDepth can be found in Appendix B. We adopt its encoder and decoder for image encoding and depth and intrinsic prediction. To adapt to 4D video reconstruction, we introduce two modifications to the UniDepth backbone:

**From Image to Video Estimation**: To facilitate video prediction, we use adaptor blocks[9, 40] to extend the original image estimator to video version. Our model processes $N_w$ input images, extracting features from each frame individually through the encoder. The UniDepth encoder produces two types of features: (1) DINO [52] features $F_{dino} \in R^{T \times \frac{H}{s_h} \times \frac{W}{s_w} \times D_{dino}}$, where $s_h \times s_w$ represents the patch size in DINO and $D_{dino}$ is the feature dimension, and (2) global token features $F_{global} \in R^{T \times D_{global}}$. To fuse the features across time, we incorporate multiple adaptors [9]. Global token features are fused using a Transformer [79] on temporal dimension, while the patched DINO [52] features are fused using Unet3D [13] on both temporal and spatial dimension. The architecture for the depth and intrinsic decoders remains unchanged from the original UniDepth implementation.

**New Confidence Mask Decoder**: To enable camera poses derivation, we need to predict an extra confidence mask as mentioned in Section 3.2. We add a new confidence mask decoder adopted from [71], which is a 3-layer MLP with ReLU [24] activation. The decoder takes a concatenation of (1) fused shallow features from the video adaptors, (2) depth features and confidence maps from the UniDepth decoder [56], and (3) an interpolation of the above features, induced by optical flow [92] from neighboring frames (Secntion 3.2). A sigmoid function is applied to normalize the final confidence score within the range $[0, 1]$.

### 4.3. Self-supervised 4D Reconstruction Losses

Although the new architecture can predict camera intrinsics, poses, and video depth simultaneously, these variants remain inconsistent in both the temporal and spatial dimensions. To address these issues, we propose an self-supervised training method that optimizes and aligns these variants in an end-to-end manner. By leveraging several 4D geometric constraints, we design self-supervised training losses to enable label-free training. We categorize our losses into two types: (1) Photometric loss, which aligns depth, intrinsics, and extrinsics to ensure consistent 4D re-

construction, and (2) Regularization loss, which accelerates training convergence and helps prevent model collapse.

### 4.3.1. Photometric Loss from Flow and Track Prior

Similar to previous methods [45, 70, 72, 73, 94, 99], we use photometric loss to align multi-frame depth estimations. This also ensures consistency between the camera parameters, poses and depth predictions. Specifically, we first back-project depth and intrinsic into per-frame pointcloudss $\hat{X}_i$. Then, we align the multi-frame results by minimizing the difference between (1) the 3D scene flow and long-term tracking induced by the camera's movement through high-confidence areas and (2) pre-computed 3D correspondences (back-projected from the optical flow computed by GMFlow [92] and long-term tracking by CoTracker [27]).

Formally, suppose $i < j$, and $\hat{X}_j^{\leftarrow i}$ is the interpolation result of $\hat{X}_j$ based on the points of $\hat{X}_i$ (which can be computed using optical flow or tracking). The alignment minimizes the 3D reprojection error in high-confidence regions between frames $i$ and $j$ in a **scale-agnostic** manner, which can be expressed as:

$$\mathcal{L}_{flow/track} = \frac{||\tilde{M}_{i+1,i}\tilde{M}_{j,j-1}(\hat{X}_i - \hat{P}_{j,i}\hat{X}_j^{\leftarrow i})||}{F(\hat{X}_j)||\tilde{M}_{i+1,i}\tilde{M}_{j,j-1}||} \quad (4)$$

where $F(\cdot)$ computes the first principal component of the pointcloudss. We use $F(\hat{X}_i)$ as a proxy for the scale of $\hat{X}_t$ and place it in the denominator of $\mathcal{L}_{flow/track}$ to prevent the pointcloudss from collapsing to a single point. Compared to widely used 2D photometric loss [71, 72], $\mathcal{L}_{flow/track}$ encourages more intuitive 3D consistency in the predictions.

Note we use a pre-computed pseudo-confidence mask $\tilde{M}$ for the photometric loss instead of the predicted mask $\hat{M}$, since the pseudo-motion mask is able to approximate from pretrained segmentation model [100] in egocentric video. The predicted mask $\hat{M}$ is optimized by backpropagating $\mathcal{L}_{flow/track}$ through $\hat{P}_{j,i}$, as described in Section 3.2. Using $\tilde{M}$ improves model stability and robustness by promoting more correspondences and preventing $\hat{M}$ from shrinking into sparse predictions. The pseudo-mask $\tilde{M}$ is computed from two sources: (1) Pseudo-dynamic areas, using a hand and interacted objects mask from EgoHOS [100], which captures motion from hand-object interactions; (2) Pseudo-edges, derived from flying pixels [63] based on UniDepth [56] depth predictions. For the Epic-Kitchen [15] dataset, dynamic masks are also estimated using epipolar loss [44] when hands are outside the camera view.

### 4.3.2. Regularization Loss from Depth and HOI Prior

To stabilize the model training process and accelerate convergence, we also regularize the training with predictions from state-of-the-art off-the-shelf models, i.e.,



Figure 3. Visualization of dataset used for training and evaluation.

Unidepth [56] for depth estimation and EgoHOS [100] for confidence mask prediction.

**Shape Regularization Loss from Depth Prior** After predicting the depth $\hat{D}_t$ and camera intrinsic $\hat{K}$, we first recover per-frame pointclouds $\hat{X}_t$ using Equation (1). We then regularize the shape of $\hat{X}_t$ with the prediction $\tilde{X}_t$ from UniDepth [56]:

$$\mathcal{L}_{shape} = \frac{1}{H \times W} \min_{s,R,T} ||sR\hat{X}_t + T - \tilde{X}_t|| \quad (5)$$

Here, $(s, R, T)$ represents the scaled SE(3) transformation used for alignment, aiming to enable regularization on a relative scale [6]. The optimal transformation can be solved in closed form using SVD [11]. Note that $\mathcal{L}_{shape}$ helps constrain dynamic parts of scenes relative to static areas.

**Mask Regularization from HOI Prior** To speed up convergence, we also regularize the prediction of $\hat{M}$:

$$\mathcal{L}_{mask} = BCELoss(\hat{\mathcal{M}}, \tilde{\mathcal{M}}) \quad (6)$$

**Camera Consistency Self-Supervision** We additionally introduce a self-supervised loss $\mathcal{L}_{con}$ to enhance the consistency of camera intrinsic predictions. For two frame sequences $V_1$ and $V_2$ from the same video clip, we enforce the intrinsic predictions to be as similar as possible:

$$\mathcal{L}_{con} = ||\hat{K}_{V_1} - \hat{K}_{V_2}|| \quad (7)$$

### 4.3.3. Final Loss Function

Put it together, the final loss is:

$$\mathcal{L} = \alpha\mathcal{L}_{shape} + \beta\mathcal{L}_{flow} + \gamma\mathcal{L}_{track} + \lambda\mathcal{L}_{mask} + \mu\mathcal{L}_{con} \quad (8)$$

where we set $\alpha = 4$, $\beta = \gamma = 5$, $\lambda = 1$, $\mu = 0.005$ as loss weights by default to balance each supervision.

| | HOI4D | | | | H2O | | | | POV-Surgery[†] | | | | ARCTIC-HOI[†] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CD ↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ | CD ↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ | CD ↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ | CD ↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ |
| DS+UniDepth [77] | <u>6.7</u> | <u>23.2</u> | <u>53.4</u> | <u>79.9</u> | <u>5.1</u> | 36.5 | 75.0 | 93.2 | <u>39.1</u> | 7.7 | 20.0 | 38.7 | <u>2.9</u> | <u>22.2</u> | **60.2** | <u>84.7</u> |
| MFVR+UniDepth [1] | 8.9 | 15.4 | 38.2 | 68.0 | **4.7** | <u>47.8</u> | <u>81.5</u> | <u>94.2</u> | 223.3 | 3.6 | 9.6 | 19.1 | 5.8 | 12.6 | 33.8 | 63.2 |
| DUSt3R [86] | 8.6 | 24.0 | 53.2 | 76.8 | 8.8 | 23.1 | 56.0 | 82.7 | 55.4 | <u>9.7</u> | <u>24.7</u> | <u>45.1</u> | 3.2 | 20.8 | 53.7 | 83.1 |
| MonSt3R [99] | 7.6 | 21.4 | 50.8 | 77.9 | 14.7 | 15.1 | 42.4 | 70.0 | 94.5 | 6.5 | 17.8 | 34.8 | 5.4 | 9.9 | 31.8 | 65.0 |
| Align3R [45] | 7.1 | 22.3 | <u>53.4</u> | 78.9 | 11.5 | 20.0 | 45.4 | 72.0 | 41.1 | 7.9 | 21.0 | 40.4 | 4.5 | 14.2 | 41.5 | 75.8 |
| EgoMono4D (Ours) | **5.9** | **27.9** | **59.6** | **83.1** | <u>5.1</u> | **54.2** | **83.9** | **94.4** | **33.8** | **13.5** | **32.0** | **53.9** | **2.8** | **24.1** | <u>57.5</u> | **86.2** |

Table 1. The evaluation results for 4D pointclouds sequence reconstruction are presented, using 3D Chamfer Distance (CD, mm) and 3D Pointclouds F-score ($F_\delta$, %). † indicates zero-shot generalization for EgoMono4D. For ARCTIC-HOI, the evaluation focuses specifically on the reconstruction quality of the hand-object region. On average, EgoMono4D demonstrates a clear advantage across the metrics.

## 4.4. Inference Strategy

Finally, we describe the inference strategy for EgoMono4D. Due to GPU memory limitations, only a limited number of frames can be processed in a single feed-forward prediction. However, we can predict videos with infinite frames in a stream manner using a sliding window. The video is first split into $N_w$ frames sub-clips with $N_o$ overlapping frames between neighbors. Then we predict neighboring windows independently. Let $w_i$ represent the $i$-th sub-clip, $E_i$ the timestamp set of $w_i$, and $\hat{S}^{w_i}$, $\hat{S}^{w_{i+1}}$ the predictions for two neighboring sub-clips. The latter is then then aligned and concatenated to the former as follows:

$$(s^*, R^*, T^*) = \arg\min_{s,R,T} ||sR\hat{S}^{w_i}_{E^{ov}_{i+1}} + T - \hat{S}^{w_{i+1}}_{E^{ov}_{i+1}}|| \quad (9)$$

$$\hat{S}^{[w_i, w_{i+1}]} = [\hat{S}^{w_i}, s^* R^* \hat{S}^{w_{i+1}}_{E_{i+1} - E^{ov}_{i+1}} + T^*] \quad (10)$$

where $E^{ov}_{i+1}$ represents the overlapping timestamps between $w_i$ and $w_{i+1}$, and $[\cdot, \cdot]$ denotes the concatenation operator.

## 5. Experiments

### 5.1. Datasets

Figure 3 shows the datasets used for training and evaluation. For more details, refer to Appendix A. Our model is trained on egocentric videos from H2O[34], HOI4D[43], FPHA[22], EgoPAT3D[37], and Epic-Kitchen[15]. Each video is split into 20-frame sub-clips for batch training, totaling 11.2 million frames, with the majority (9.7M frames) from the unlabeled Epic-Kitchen dataset.

For evaluation, we use datasets with pointcloud sequence labels. In-domain evaluation is done using H2O[34] and HOI4D[43], with the datasets split by Scene ID to ensure no overlap between training and test sets. For zero-shot generalization, we use POV-Surgery[85] and ARCTIC[20] (only Mocap [78] Hand and Objects Interaction (HOI) labels). To avoid redundancy, we only use the first record from the first participant in each task. Videos are split into 40-frame sub-clips for batch evaluation. Note that ARCTIC provides only

hand and object labels, so we refer to it as ARCTIC-HOI to highlight its focus on foreground HOI reconstruction.

### 5.2. Implementation Details

We initialize our model with the pretrained UniDepthV2-L weights[56] and freeze the encoder. For input preprocessing, we resize the images to a resolution of $288 \times 384$. During training, each data point consists of 4 frames sampled from each sub-clip, with the interval between frames randomly selected from the range $[1, 4]$. We employ the Adam [30] optimizer with a learning rate of 5e-5 and a batch size of 16. The model is trained on 8 NVIDIA A800 GPUs for 350k iterations. For inference, we set the overlap to 1 frame. For stability, we set window size $N_w = 4$ by default, the same with training process.

### 5.3. Baseline

We compare our model to previous methods that (1) provide dense 4D reconstruction, (2) have nearly linear time complexity with respect to the number of frames, enabling large-scale reconstruction, and (3) demonstrate zero-shot generalization. **DS+UniDepth** (DROID-SLAM [77] + UniDepth [56]) combines depth estimation with learning-based RGBD visual odometry. **MFVR+UniDepth** (MapFreeVR [1] + UniDepth [56]) first estimates image depth, then computes camera poses with correspondence estimation (optical flow [92] in our setting) and depth-alignment. Additionally, we integrate HOI masks from EgoHOS[100] to filter out the HOI region for alignment. This baseline could be viewed as a modularized and no-training version of EgoMono4D. **DUSt3R** [86] supports end-to-end reconstruction. We use the "swin" mode from the original implementation to achieve $\mathcal{O}(T)$ inference. **MonSt3R** [99] is a finetuned version of DUSt3R on synthetic dynamic datasets, aimed at improving reconstruction performance in dynamic areas. Lastly, **Align3R** [45] further merge depth prior during finetuning to enhance the geometry estimation ability.

Since large-scale labeled 4D datasets (with both high-

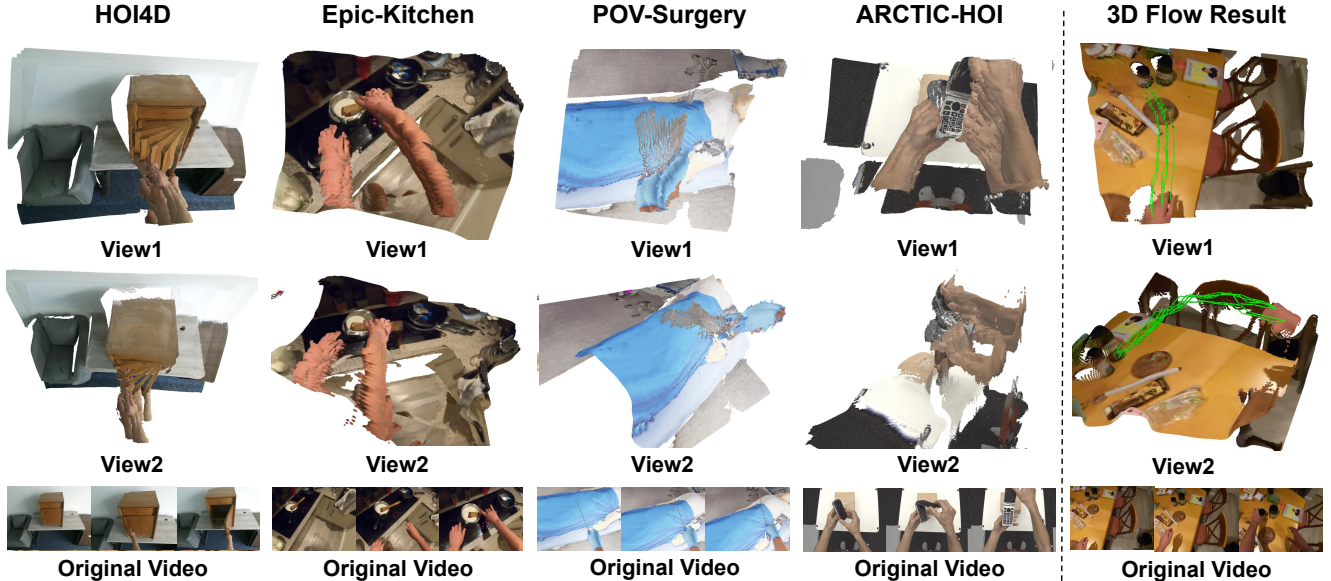| HOI4D | Epic-Kitchen | POV-Surgery | ARCTIC-HOI | 3D Flow Result |
| --- | --- | --- | --- | --- |
| View1 | View1 | View1 | View1 | View1 |
| View2 | View2 | View2 | View2 | View2 |
| Original Video | Original Video | Original Video | Original Video | Original Video |

Figure 4. The visualization of the dense pointcloud sequence reconstruction by EgoMono4D demonstrates its ability to effectively recover both the overall scene structure and dynamic motion elements to a significant extent. For additional visualizations, please refer to Appendix G. We also provide a **qualitative comparison** and visualization with baseline methods in Appendix H. **Video and interactable visualizations** can be found in project website.
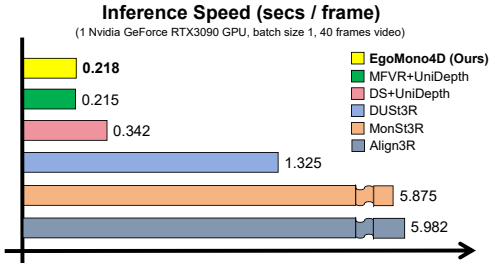


Figure 5. Inference speed comparison of different methods.

quality depth and camera labels) for in-the-wild egocentric videos are lacking, existing evaluation are limited to small, lab-based settings [20, 34, 43, 85]. Fine-tuning on these specialized datasets could lead to overfitting, resulting in an unfair comparison, since EgoMono4D is trained on large-scale, in-the-wild datasets and designed for generalization. To avoid this, and following previous self-supervised depth estimation works [72, 73, 94], we refrain from fine-tuning models on lab-based egocentric datasets. Moreover, our aim is to evaluate how self-supervised methods perform in label-scarce settings, such as egocentric videos. By excluding labels, we gain clearer insights into how these methods tackle the inherent challenges of such domains.

## 6. Result

Here we present the evaluation and ablation results of EgoMono4D. Detailed definitions of all metrics can be

found in Appendix C.

### 6.1. Dense Pointclouds Sequence Reconstruction

**Task and Metric** Dense pointclouds sequence reconstruction [71, 86, 99] requires the model to reconstruct the $H \times W$ pointclouds of each frame within a global coordinate system. Since the scale of the pointclouds is ambiguous for monocular reconstruction, we first apply an estimated globally scaled SE(3) transformation $(s, R, T)$ to align the prediction to the ground-truth. We then evaluate the accuracy of each frame's shape and average these results as the final evaluation score. Following [53, 56], we use the 3D Chamfer Distance (CD, mm) and the 3D Pointclouds F-score ($F_\delta$, %) to assess shape similarity. We use the notation $F_\delta$ to denote the F-score with a threshold of $\delta$ centimeters (cm) as the positive criterion.

**Results** The quantitative evaluation results are presented in Table 1, with visualizations shown in Figure 4. For additional visualizations, refer to Appendix G. We also provide a **qualitative comparison** and visualization with baseline methods in Appendix H. **Video and interactable visualization** can be found in project website .

EgomMono4D demonstrates superior performance across all evaluated methods. Notably, on the challenging POV-Surgery dataset [85], which contains complex surgical scenes with unrealistic textures and intricate actions, our model outperforms others by 10–20% in terms of F-score. DS+UniDepth also shows commendable performance

across several metrics. However, qualitative visualizations (Appendix H) reveal that it struggles to align the static portions of egocentric scenes.

The modular version of our model, MFVR+UniDepth, performs significantly worse across most metrics. This can be attributed to inconsistencies between the different modules. In contrast, our end-to-end self-supervised training (described in Section 4.3) ensures that the modules align through back-propagation of the 4D supervision loss, leading to improved 4D reconstruction accuracy.

Due to the limited availability of labeled egocentric data for training, DUSt3R, MonSt3R and Align3R exhibit suboptimal performance on egocentric videos. Furthermore, fine-tuning MonSt3R on small-scale dynamic datasets actually results in performance degradation compared to DUSt3R, primarily due to domain gap and overfitting. By introducing depth prior, Align3R alleviates this issue, but still fails to get precise monocular 4D estimation. These results highlight the challenges that supervised learning methods face in label-scarce scenarios. In contrast, our self-supervised approach performs effectively on unlabeled egocentric videos. It is important to note that this does not imply that supervised methods are inherently inferior to self-supervised ones. In domains with abundant labeled data, supervised methods may offer advantages, as demonstrated in depth estimation tasks [61].

## 6.2. Long-term 3D Scene Flow Recovery

We evaluate different models using long-term 3D scene flow [80], which captures both the structure and dynamics of egocentric scenes. Given a video and query points in the first frame, long-term 3D flow represents the future trajectory of each point in 3D space [33, 96]. EgoMono4D also outperforms all baseline in this task. Details and results are in Appendix D, with an prediction example in Figure 4.

## 6.3. Inference Speed

We evaluate the inference speed of all models (for 40 frames video). All measurements are conducted on a single NVIDIA GeForce 3090 GPU with a batch size of 1. Results are shown in Figure 5. Our model achieves the fastest speed, except for its modularized version (MFVR+UniDepth).

## 6.4. Ablation Study

We conduct an ablation study of EgoMono4D, training the model on 280K frames from the training set and testing its performance on the HOI4D dataset [43]. The variants tested are as follows: (1) w UniDepth-S: using small version of UniDepth [56]. (2) w/o V-Adaptor: remove video adaptor. (3) w $\alpha = 1$, changing the training weight. (4) w/ midas-loss: replacing $\mathcal{L}_{shape}$ with depth supervision at a relative scale [6]; (5) 2d-flow-loss: substituting $\mathcal{L}_{flow/track}$ with its 2D version in pixel space [71]; (6) w/o mask-loss: removing

|  | CD | $F_1$ | $F_{2.5}$ | $F_5$ |
|---|---|---|---|---|
| complete | 6.3 | 26.9 | 56.9 | 81.4 |
| w UniDepth-S | 6.5 ↓ | 23.8 ↓ | 54.5 ↓ | 79.2 ↓ |
| w/o V-Adaptor | 6.4 ↓ | 25.8 ↓ | 55.8 ↓ | 80.4 ↓ |
| w $\alpha = 1$ | 7.1 ↓ | 18.2 ↓ | 48.3 ↓ | 66.4 ↓ |
| w/ midas-loss | 6.3 · | 26.1 ↓ | 56.2 ↓ | 80.4 ↓ |
| w/ 2d-flow-loss | 6.2 ↑ | 26.7 ↓ | 57.0 ↑ | 81.0 ↓ |
| w/o mask-loss | 6.5 ↓ | 23.9 ↓ | 53.3 ↓ | 79.0 ↓ |
| w/o cc-loss | 6.3 · | 25.2 ↓ | 56.4 ↓ | 81.1 ↓ |

Table 2. Ablation study of the EgoMono4D model on the HOI4D dataset. ↓ indicates performance degradation relative to the complete model, while ↑ indicates improvement. The complete model outperforms other ablated variants on average.

| | POV-Surgery | | |
|---|---|---|---|
| | CD↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ |
| fps / 1 | 25.5 | 14.2 | 33.3 | 54.8 |
| fps / 2 | 25.4 | 13.8 | 32.8 | 54.9 |
| fps / 4 | 25.2 | 14.1 | 33.1 | 55.1 |
| fps / 12 | 28.5 | 13.0 | 31.0 | 53.5 |

Table 3. The POV-Surgery pointclouds sequence reconstruction results across different frames per second (fps) settings.

the confidence mask regularization loss $\mathcal{L}_{mask}$; and (4) w/o cc-loss: removing the self-supervised intrinsic loss $\mathcal{L}_{cc}$.

Results are presented in Table 2. The complete model outperforms all other variants on average, demonstrating the effectiveness of our design choices. From the results, we observe that mask and depth regularization plays a crucial role in stabilizing training. Additionally, the choice of the base model is important, and we anticipate that improvements in depth estimation will further benefit our self-supervised scene reconstruction approach in the future. In terms of shape regularization and photometric loss, applying constraints in 3D space yields moderately better results than applying them in 2D space on average.

We also conduct an ablation on the hyperparameters for model inference, including the window size $N_w$ and window overlap size $N_o$, which are detailed in Appendix E. For the window size $N_w$, maintaining consistency between training and inference is critical, likely because the video adaptor is specifically trained to fuse 4 frames. Regarding the overlap size $N_o$, the model performs comparably for $N_o = 1, 2, 3$. Therefore, we select $N_o = 1$ to maximize inference speed.

## 6.5. Impact on Video FPS

Since we use random frame intervals to sample data during training, our model is expected to be robust to variations in video fps to some extent. We test this on the pointclouds sequence reconstruction task using the zero-shot dataset POV-Surgery [85]. For a 40-frame sub-clip of POV-Surgery, we

select only frames 0, 12, 24, and 36 for evaluation. We define a $1/x$ fps video as a frame sequence sampled with an interval of $x$ (where 12 should be divisible by $x$). The evaluation results are shown in Table 3, and they demonstrate that our method is robust to changes in video fps within a certain range. However, when the fps becomes too low, performance degrades, likely because the optical flow module [92] we rely on may fail under such conditions.

## 7. Conclusion

We present EgoMono4D, an self-supervised model for 4D reconstruction of egocentric videos, trained solely on large-scale unlabeled data. By aligning video depth with 4D constraints, it achieves promising zero-shot results in dense, generalizable scene reconstruction. Additional visualization, details, discussion, limitations and future direction could be found in appendix.

## References

[1] Eduardo Arnold, Jamie Wynn, Sara Vicente, Guillermo Garcia-Hernando, Aron Monszpart, Victor Prisacariu, Daniyar Turmukhambetov, and Eric Brachmann. Map-free visual relocalization: Metric pose relative to a single image. In *European Conference on Computer Vision*, pages 690–708. Springer, 2022. 2, 6, 16, 17

[2] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13778–13790, 2023. 3

[3] Wentao Bao, Lele Chen, Libing Zeng, Zhong Li, Yi Xu, Junsong Yuan, and Yu Kong. Uncertainty-aware state space transformer for egocentric 3d hand trajectory forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13702–13711, 2023. 2, 3, 15, 16

[4] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 2

[5] Jiawang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised scale-consistent depth and ego-motion learning from monocular video. *Advances in neural information processing systems*, 32, 2019. 2

[6] Reiner Birkl, Diana Wofk, and Matthias Müller. Midas v3. 1–a model zoo for robust monocular relative depth estimation. *arXiv preprint arXiv:2307.14460*, 2023. 2, 5, 8

[7] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6): 1874–1890, 2021. 2

[8] Weirong Chen, Le Chen, Rui Wang, and Marc Pollefeys. Leap-vo: Long-term effective any point tracking for visual odometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19844–19853, 2024. 2

[9] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 4

[10] Suhwan Cho, Minhyeok Lee, Seunghoon Lee, Chaewon Park, Donghyeong Kim, and Sangyoun Lee. Treating motion as option to reduce motion dependency in unsupervised video object segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 5140–5149, 2023. 22

[11] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2514–2523, 2020. 4, 5

[12] Wen-Hsuan Chu, Lei Ke, and Katerina Fragkiadaki. Dreamscene4d: Dynamic multi-object scene generation from monocular videos. *arXiv preprint arXiv:2405.02280*, 2024. 2

[13] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part II 19*, pages 424–432. Springer, 2016. 4

[14] Zhaopeng Cui and Ping Tan. Global structure-from-motion by similarity averaging. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 864–872, 2015. 2

[15] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision*, pages 1–23, 2022. 1, 2, 3, 5, 6, 14

[16] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 14

[17] Bardienus Duisterhof, Lojze Zust, Philippe Weinzaepfel, Vincent Leroy, Yohann Cabon, and Jerome Revaud. Mast3r-sfm: a fully-integrated solution for unconstrained structure-from-motion. *arXiv preprint arXiv:2409.19152*, 2024. 3

[18] Ben Eisner, Harry Zhang, and David Held. Flowbot3d: Learning 3d articulation flow to manipulate articulated objects. *arXiv preprint arXiv:2205.04382*, 2022. 15

[19] Sven Elflein, Qunjie Zhou, Sérgio Agostinho, and Laura Leal-Taixé. Light3r-sfm: Towards feed-forward structure-from-motion. *arXiv preprint arXiv:2501.14914*, 2025. 3

[20] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J Black, and Otmar Hilliges. Arctic: A dataset for dexterous bimanual hand-object manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12954, 2023. 3, 6, 7, 14

[21] Clement Fuji Tsang, Maria Shugrina, Jean Francois Lafleche, Towaki Takikawa, Jiehan Wang, Charles Loop, Wenzheng Chen, Krishna Murthy Jatavallabhula, Edward Smith, Artem Rozantsev, Or Perel, Tianchang Shen, Jun Gao, Sanja Fidler, Gavriel State, Jason Gorski, Tommy Xiang, Jianing Li, Michael Li, and Rev Lebaredian. Kaolin: A pytorch library for accelerating 3d deep learning research. https://github.com/NVIDIAGameWorks/kaolin, 2022. 14

[22] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-person hand action benchmark with rgb-d videos and 3d hand pose annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 409–419, 2018. 3, 6, 14

[23] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 1, 2, 3

[24] Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns. *Advances in neural information processing systems*, 32, 2019. 4

[25] Wenbo Hu, Xiangjun Gao, Xiaoyu Li, Sijie Zhao, Xiaodong Cun, Yong Zhang, Long Quan, and Ying Shan. Depthcrafter: Generating consistent long depth sequences for open-world videos. *arXiv preprint arXiv:2409.02095*, 2024. 2

[26] Linyi Jin, Richard Tucker, Zhengqi Li, David Fouhey, Noah Snavely, and Aleksander Holynski. Stereo4d: Learning how things move in 3d from internet stereo videos. *arXiv preprint arXiv:2412.09621*, 2024. 3

[27] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Cotracker: It is better to track together. *arXiv preprint arXiv:2307.07635*, 2023. 5, 15

[28] Yoni Kasten, Wuyue Lu, and Haggai Maron. Fast encoder-based 3d from casual videos via point track processing. In *ECCV 2024 Workshop on Wild 3D: 3D Modeling, Reconstruction, and Generation in the Wild*. 3

[29] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2, 3

[30] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 22

[32] Johannes Kopf, Xuejian Rong, and Jia-Bin Huang. Robust consistent video depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1611–1621, 2021. 2

[33] Skanda Koppula, Ignacio Rocco, Yi Yang, Joe Heyward, João Carreira, Andrew Zisserman, Gabriel Brostow, and Carl Doersch. Tapvid-3d: A benchmark for tracking any point in 3d. *arXiv preprint arXiv:2407.05921*, 2024. 2, 3, 8, 15, 16

[34] Taein Kwon, Bugra Tekin, Jan Stühmer, Federica Bogo, and Marc Pollefeys. H2o: Two hands manipulating objects for first person interaction recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10138–10148, 2021. 3, 6, 7, 14

[35] Jiahui Lei, Yijia Weng, Adam Harley, Leonidas Guibas, and Kostas Daniilidis. Mosca: Dynamic gaussian fusion from casual videos via 4d motion scaffolds. *arXiv preprint arXiv:2405.17421*, 2024. 2

[36] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 21

[37] Yiming Li, Ziang Cao, Andrew Liang, Benjamin Liang, Luoyao Chen, Hang Zhao, and Chen Feng. Egocentric prediction of action target in 3d. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20971–20980. IEEE, 2022. 3, 6, 14

[38] Zhengqi Li, Richard Tucker, Forrester Cole, Qianqian Wang, Linyi Jin, Vickie Ye, Angjoo Kanazawa, Aleksander Holynski, and Noah Snavely. Megasam: Accurate, fast, and robust structure and motion from casual dynamic videos. *arXiv preprint arXiv:2412.04463*, 2024. 3

[39] Hanxue Liang, Jiawei Ren, Ashkan Mirzaei, Antonio Torralba, Ziwei Liu, Igor Gilitschenski, Sanja Fidler, Cengiz Oztireli, Huan Ling, Zan Gojcic, et al. Feed-forward bullet-time reconstruction of dynamic scenes from monocular videos. *arXiv preprint arXiv:2412.03526*, 2024. 3

[40] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022. 4

[41] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 22

[42] Xiaochen Liu, Tao Zhang, and Mingming Liu. Joint estimation of pose, depth, and optical flow with a competition–cooperation transformer network. *Neural Networks*, 171: 263–275, 2024. 2

[43] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21013–21022, 2022. 2, 3, 6, 7, 8, 14, 17

[44] Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng, Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf, and Jia-Bin Huang. Robust dynamic radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13–23, 2023. 5

[45] Jiahao Lu, Tianyu Huang, Peng Li, Zhiyang Dou, Cheng Lin, Zhiming Cui, Zhen Dong, Sai-Kit Yeung, Wenping Wang, and Yuan Liu. Align3r: Aligned monocular depth estimation for dynamic videos. *arXiv preprint arXiv:2412.03079*, 2024. 1, 3, 5, 6, 16, 17

[46] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023. 2

[47] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3061–3070, 2015. 15

[48] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *Advances in Neural Information Processing Systems*, 36, 2024. 3

[49] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 2

[50] Riku Murai, Eric Dexheimer, and Andrew J Davison. Mast3r-slam: Real-time dense slam with 3d reconstruction priors. *arXiv preprint arXiv:2412.12392*, 2024. 3

[51] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022. 1, 3

[52] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 4

[53] Evin Pınar Örnek, Shristi Mudgal, Johanna Wald, Yida Wang, Nassir Navab, and Federico Tombari. From 2d to 3d: Re-thinking benchmarking of monocular depth prediction. *arXiv preprint arXiv:2203.08122*, 2022. 7, 14

[54] Xiaqing Pan, Nicholas Charron, Yongqian Yang, Scott Peters, Thomas Whelan, Chen Kong, Omkar Parkhi, Richard Newcombe, and Yuheng Carl Ren. Aria digital twin: A new benchmark dataset for egocentric 3d machine perception. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20133–20143, 2023. 1, 3, 22

[55] Georgios Pavlakos, Dandan Shan, Ilija Radosavovic, Angjoo Kanazawa, David Fouhey, and Jitendra Malik. Reconstructing hands in 3d with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9826–9836, 2024. 3

[56] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. Unidepth: Universal monocular metric depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10106–10116, 2024. 2, 4, 5, 6, 7, 8, 14, 17, 18, 20

[57] Rolandos Alexandros Potamias, Jinglei Zhang, Jiankang Deng, and Stefanos Zafeiriou. Wilor: End-to-end 3d hand localization and reconstruction in-the-wild. *arXiv preprint arXiv:2409.12259*, 2024. 3

[58] Shraman Pramanick, Yale Song, Sayan Nag, Kevin Qinghong Lin, Hardik Shah, Mike Zheng Shou, Rama Chellappa, and Pengchuan Zhang. Egovlpv2: Egocentric video-language pre-training with fusion in the backbone. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5285–5297, 2023. 1

[59] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhuowen Tu, and Li Erran Li. Affordancellm: Grounding affordance from vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7587–7597, 2024. 3

[60] Yuzhe Qin, Hao Su, and Xiaolong Wang. From one hand to multiple hands: Imitation learning for dexterous manipulation from single-camera teleoperation. *IEEE Robotics and Automation Letters*, 7(4):10873–10881, 2022. 3

[61] Uchitha Rajapaksha, Ferdous Sohel, Hamid Laga, Dean Diepeveen, and Mohammed Bennamoun. Deep learning-based depth estimation methods from monocular image and videos: A comprehensive survey. *ACM Computing Surveys*, 2024. 2, 8

[62] Antoni Rosinol, John J Leonard, and Luca Carlone. Nerf-slam: Real-time dense monocular slam with neural radiance fields. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3437–3444. IEEE, 2023. 2

[63] Alexander Sabov and Jörg Krüger. Identification and correction of flying pixels in range camera data. In *Proceedings of the 24th Spring Conference on Computer Graphics*, pages 135–142, 2008. 5, 15

[64] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 1, 2

[65] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pages 501–518. Springer, 2016. 2

[66] Daniel Seita, Yufei Wang, Sarthak J Shetty, Edward Yao Li, Zackory Erickson, and David Held. Toolflownet: Robotic manipulation with tools via predicting tool flow from point clouds. In *Conference on Robot Learning*, pages 1038–1049. PMLR, 2023. 15

[67] Dandan Shan, Jiaqi Geng, Michelle Shu, and David F Fouhey. Understanding human hands in contact at internet scale. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9869–9878, 2020. 3

[68] Lin Shao, Parth Shah, Vikranth Dwaracherla, and Jeannette Bohg. Motion-based object segmentation based on dense rgb-d scene flow. *IEEE Robotics and Automation Letters*, 3 (4):3797–3804, 2018. 15

[69] Dongseok Shim and H Jin Kim. Swindepth: Unsupervised depth estimation using monocular sequences via swin trans-

former and densely cascaded network. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4983–4990. IEEE, 2023. 2

[70] Cameron Smith, Yilun Du, Ayush Tewari, and Vincent Sitzmann. Flowcam: training generalizable 3d radiance fields without camera poses via pixel-aligned scene flow. *arXiv preprint arXiv:2306.00180*, 2023. 4, 5

[71] Cameron Smith, David Charatan, Ayush Tewari, and Vincent Sitzmann. Flowmap: High-quality camera poses, intrinsics, and depth via gradient descent. *arXiv preprint arXiv:2404.15259*, 2024. 2, 3, 4, 5, 7, 8

[72] Jaime Spencer, Chris Russell, Simon Hadfield, and Richard Bowden. Kick back & relax: Learning to reconstruct the world by watching slowtv. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15768–15779, 2023. 2, 3, 5, 7

[73] Jaime Spencer, Chris Russell, Simon Hadfield, and Richard Bowden. Kick back & relax++: Scaling beyond ground-truth depth with slowtv & cribstv. *arXiv preprint arXiv:2403.01569*, 2024. 2, 5, 7

[74] Yukun Su, Jingliang Deng, Ruizhou Sun, Guosheng Lin, Hanjing Su, and Qingyao Wu. A unified transformer framework for group-based segmentation: Co-segmentation, co-saliency detection and video salient object detection. *IEEE Transactions on Multimedia*, 26:313–325, 2023. 22

[75] Jiakai Sun, Han Jiao, Guangyuan Li, Zhanjie Zhang, Lei Zhao, and Wei Xing. 3dgstream: On-the-fly training of 3d gaussians for efficient streaming of photo-realistic free-viewpoint videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20675–20685, 2024. 2

[76] Yihong Sun and Bharath Hariharan. Dynamo-depth: fixing unsupervised depth estimation for dynamical scenes. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[77] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems*, 34:16558–16569, 2021. 1, 2, 6, 16, 17, 20

[78] Vicon Vantage. Cutting edge, flagship camera with intelligent feedback and resolution. 6

[79] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 4, 14

[80] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins, and Takeo Kanade. Three-dimensional scene flow. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, pages 722–729. IEEE, 1999. 8, 15

[81] Hengyi Wang and Lourdes Agapito. 3d reconstruction with spatial memory. *arXiv preprint arXiv:2408.16061*, 2024. 3

[82] Jianyuan Wang, Nikita Karaev, Christian Rupprecht, and David Novotny. Vggsfm: Visual geometry grounded deep structure from motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21686–21697, 2024. 2, 21

[83] Qianqian Wang, Vickie Ye, Hang Gao, Jake Austin, Zhengqi Li, and Angjoo Kanazawa. Shape of motion: 4d reconstruction from a single video. *arXiv preprint arXiv:2407.13764*, 2024. 2, 16

[84] Qianqian Wang, Yifei Zhang, Aleksander Holynski, Alexei A Efros, and Angjoo Kanazawa. Continuous 3d perception model with persistent state. *arXiv preprint arXiv:2501.12387*, 2025. 1, 3

[85] Rui Wang, Sophokles Ktistakis, Siwei Zhang, Mirko Meboldt, and Quentin Lohmeyer. Pov-surgery: A dataset for egocentric hand and tool pose estimation during surgical activities. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 440–450. Springer, 2023. 2, 3, 6, 7, 8, 14, 17

[86] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20697–20709, 2024. 1, 2, 3, 6, 7, 16, 17, 21, 22

[87] Wenshan Wang, Yaoyu Hu, and Sebastian Scherer. Tartanvo: A generalizable learning-based vo. In *Conference on Robot Learning*, pages 1761–1772. PMLR, 2021. 2

[88] Yiran Wang, Min Shi, Jiaqi Li, Zihao Huang, Zhiguo Cao, Jianming Zhang, Ke Xian, and Guosheng Lin. Neural video depth stabilizer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9466–9476, 2023. 2

[89] Ying Wang, Dongdong Sun, Rui Chen, Yanlai Yang, and Mengye Ren. Egocentric video comprehension via large language model inner speech. In *3rd International Ego4D Workshop*, 2023. 3

[90] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023. 1, 3

[91] Junyu Xie, Charig Yang, Weidi Xie, and Andrew Zisserman. Moving object segmentation: All you need is sam (and flow). *arXiv preprint arXiv:2404.12389*, 2024. 22

[92] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8121–8130, 2022. 4, 5, 6, 9, 21

[93] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024. 2

[94] Zhenheng Yang, Peng Wang, Yang Wang, Wei Xu, and Ram Nevatia. Every pixel counts: Unsupervised geometry learning with holistic 3d motion understanding. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 2, 5, 7

[95] Wei Yin, Chi Zhang, Hao Chen, Zhipeng Cai, Gang Yu, Kaixuan Wang, Xiaozhi Chen, and Chunhua Shen. Metric3d: Towards zero-shot metric 3d prediction from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9043–9053, 2023. 2

[96] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. General flow as foundation affordance for scalable robot

learning. *arXiv preprint arXiv:2401.11439*, 2024. 1, 2, 3, 8, 15, 16

[97] Raza Yunus, Jan Eric Lenssen, Michael Niemeyer, Yiyi Liao, Christian Rupprecht, Christian Theobalt, Gerard Pons-Moll, Jia-Bin Huang, Vladislav Golyanik, and Eddy Ilg. Recent trends in 3d reconstruction of general non-rigid scenes. In *Computer Graphics Forum*, page e15062. Wiley Online Library, 2024. 2

[98] Daiwei Zhang, Gengyan Li, Jiajie Li, Mickaël Bressieux, Otmar Hilliges, Marc Pollefeys, Luc Van Gool, and Xi Wang. Egogaussian: Dynamic scene understanding from egocentric video with 3d gaussian splatting. *arXiv preprint arXiv:2406.19811*, 2024. 2, 3

[99] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *arXiv preprint arXiv:2410.03825*, 2024. 1, 2, 3, 5, 6, 7, 16, 17, 18, 20, 22

[100] Lingzhi Zhang, Shenghao Zhou, Simon Stent, and Jianbo Shi. Fine-grained egocentric hand-object segmentation: Dataset, model, and applications. In *European Conference on Computer Vision*, pages 127–145. Springer, 2022. 1, 2, 3, 5, 6

[101] Youmin Zhang, Fabio Tosi, Stefano Mattoccia, and Matteo Poggi. Go-slam: Global optimization for consistent 3d instant reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3727–3737, 2023. 2

[102] Zhoutong Zhang, Forrester Cole, Zhengqi Li, Michael Rubinstein, Noah Snavely, and William T Freeman. Structure and motion from casual videos. In *European Conference on Computer Vision*, pages 20–37. Springer, 2022. 2

[103] Wang Zhao, Shaohui Liu, Hengkai Guo, Wenping Wang, and Yong-Jin Liu. Particlesfm: Exploiting dense point trajectories for localizing moving cameras in the wild. In *European Conference on Computer Vision*, pages 523–542. Springer, 2022. 2

[104] Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetzstein, and Leonidas J Guibas. Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19855–19865, 2023. 22

# Appendix

This appendix file provides:

For **video and interactable visualization**, please refer to the project website.

## A. Details of Datasets

Figure 3 in the main paper visualizes samples from various datasets. Detailed information about the datasets is shown in Table 4. The data encompass different types of egocentric videos, with variations in camera motion (small/large), action complexity (simple/complex), and scene conditions (clean/cluttered).

For training, we use a combination of H2O[34] (40K frames), HOI4D[43] (540K frames), FPHA[22] (73K frames), EgoPAT3D[37] (823K frames), and Epic-Kitchen[15] (9.7M frames). The final training dataset contains a total of 11.2M frames, with Epic-Kitchen dominating the data (about 85%), providing large-scale diversity and a wide range of behavior modes. The other datasets contribute unique scene characteristics or behavior patterns. For instance, H2O[34] exclusively contains bi-manual operations with small camera motion in clean table scenes. Approximately 5% of the data is allocated for model validation.

For evaluation, we note a scarcity of egocentric datasets offering both high-quality depth and precise camera labels. We strongly encourage the computer vision and HOI communities to collect or synthesize larger-scale RGBD datasets with accurate pose annotations. We selected four datasets that meet the necessary criteria: H2O[34], HOI4D[43], ARCTIC[20], and POV-Surgery[85]. The sources for camera labels are provided in Table 4.

H2O[34] and HOI4D[43] feature simple scenes and actions and are used as test benchmarks for in-domain prediction performance. For zero-shot generalization evaluation, we use ARCTIC[20] and POV-Surgery[85]. ARCTIC[20] only provides labels for hands and objects, which is why we refer to it as ARCTIC-HOI, It is primarily used to assess the reconstruction quality of HOI. Both ARCTIC-HOI and POV-Surgery present significant challenges for EgoMono4D, as they exhibit a large domain gap from the training data. (1) For ARCTIC-HOI, the hand and object components occupy a much larger portion of the images compared to the training data. (2) For POV-Surgery, the dataset's unrealistic textures create a substantial visual domain gap, and the surgical scenes were not encountered during training.

## B. Demonstration of UniDepth Backbone

Our architecture builds upon the UniDepth backbone [56], an encoder-decoder architecture depth estimator. UniDepth decouples the tasks of depth and camera estimation by transforming the scene from Cartesian coordinates to a pseudo-spherical representation, which enables dense camera prediction in spherical space. It then incorporates scene scale information from the camera prediction into the depth estimation module using Laplace Spherical Harmonic Encoding (SHE) and a cross-attention mechanism [16, 79]. For more details, please check out the original paper Piccinelli et al. [56].

## C. Details of Evaluation Metrics

We provide the mathematical definitions of the metrics used to evaluate the performance of 3D point cloud sequence reconstruction and long-term 3D scene flow recovery.

### C.1. Metrics for Pointclouds Sequence

We follow Örnek et al. [53] in using 3D Chamfer Distance (CD, measured in millimeters) and the 3D Pointclouds F-score (F, measured as a percentage %) to evaluate shape similarity. For implementation, we leverage the Kaolin library [21]. Given the ambiguity in the scale of pointclouds, we first align the predicted point cloud sequence to the ground truth using an estimated best-aligned global scaled SE(3) transformation $(s, R, T)$.

**3D Chamfer Distance (CD, mm).** Given the predicted per-frame pointclouds $P \in R^{N \times 3}$ and the ground-truth $G \in R^{N \times 3}$, where $N$ is the number of points, the 3D Chamfer Distance (CD) is defined as:

$$CD(R, G) = \sum_{x \in G} \min_{y \in R} ||x - y|| + \sum_{y \in R} \min_{x \in G} ||x - y|| \quad (11)$$

**3D Pointclouds F-score (F).** The 3D F-score combines precision and recall to provide a balanced evaluation of the predicted surface quality. In the context of 3D pointclouds, precision measures how many points from the predicted surface are close to the ground truth surface, while recall measures how many ground truth points are captured by the predicted surface. Given a distance $\delta$ as the positive threshold, the precision $P_\delta(R, G)$, recall $R_\delta(R, G)$ and F-score

| Training | | | | | | |
|---|---|---|---|---|---|---|
| **Datasets** | **# of frames** | **Data Split** | **Camera Motion** | **Action** | **Scene** | **Note** |
| H2O | 40K | Original Split | Small | Simple | Clean | Bi-manual |
| HOI4D | 540K | Room ID | Medium | Simple | Clean | |
| FPHA | 73K | Task | Medium | Complex | Clutter | |
| EgoPAT3D | 823K | Scene ID | Large | Medium | Medium | Only pick & place |
| Epic-Kitchen | 9.7M | Scene ID | Large | Complex | Clutter | |
| Evaluation | | | | | | |
| **Datasets** | **# of frames** | **Label** | **Camera Motion** | **Action** | **Scene** | **Note** |
| H2O | 8K | Calibration | Small | Simple | Clean | Bi-manual |
| HOI4D | 12K | SfM | Medium | Simple | Clean | Contain noise |
| ARCTIC-HOI | 13K | Mocap | Medium | Complex | Clean | Only hand and object label |
| POV-Surgery | 26K | Synthesis | Large | Medium | Clutter | Unrealistic texture |

Table 4. Comparison of Different Datasets for Training and Evaluation

$F_\delta(R, G)$ are defined as:

$$P_\delta(R, G) = \frac{1}{|R|} \sum_{y \in R} [d_{y \to G} < \delta] \qquad (12)$$

$$R_\delta(R, G) = \frac{1}{|G|} \sum_{x \in G} [d_{x \to R} < \delta] \qquad (13)$$

$$F_\delta(R, G) = \frac{2P_\delta(R, G) R_\delta(R, G)}{P_\delta(R, G) + R_\delta(R, G)} \qquad (14)$$

### C.2. Metrics for Long-term 3D Scene Flow

Long-term 3D scene flow refers to predicting the future trajectories of multiple 3D query points in the pointclouds of the first frame [96]. Following previous works [3, 33, 96], we evaluate the precision of 3D flow recovery using three metrics: Average Displacement Error (ADE, measured in millimeters), Final Displacement Error (FDE, measured in millimeters), and Precision under Distance (P, measured as a percentage %). The 3D flow is generated by interpolating between the predicted and ground-truth pointclouds based on 2D tracking from CoTracker [27]. Before evaluation, we filter out trajectories affected by noise from flying pixels [63] using ground-truth depth information. To align the scale of scenes and the initial position of the 3D query points, we first perform a best-aligned scaled SE(3) transformation between the ground-truth and predicted pointclouds for the first frame.

**Average Displacement Error (ADE, mm).** Given the predicted 3D flow $F \in R^{T \times N \times 3}$ and ground-truth flow $G \in R^{T \times N \times 3}$, where $T$ represents the number of frames and $N$ represents the number of trajectories, ADE measures the average displacement across all timestamps.

$$ADE(F, G) = \frac{1}{N} \sum_{i=1}^{N} ||F_i - G_i|| \qquad (15)$$

**Final Displacement Error (FDE, mm).** FDE measures the displacement of the final timestamp.

$$FDE(F, G) = ||F_{T-1} - G_{T-1}|| \qquad (16)$$

**Precision under Distance (P, %).** The precision metric measures the average percentage of points with an error within $\delta$ centimeters (cm).

$$P_\delta = \frac{1}{N} \sum_{i=1}^{N} [||F_i - G_i|| < \delta] \qquad (17)$$

### D. Long-term 3D Scene Flow Recovery Results

**Task** Long-term 3D scene flow [80] captures both the structure and dynamics of egocentric scenes in a compact format, making it valuable for various applications such as perception [68], autonomous driving [47], and robot learning [18, 66, 96]. Given a video and a set of query points in the first frame, the 3D flow [33, 96] represents the future trajectory of each query point in 3D space. Since egocentric datasets lack explicit 3D flow labels, we first employ CoTracker [27], a high-precision pixel tracker, to generate 2D long-term tracking (with a 35×35 grid of query points). These 2D trajectories are then unprojected using ground-truth pointclouds sequence to create 3D trajectory labels. We use the same method to obtain predictions for models.

**Metric** To align the scale of scenes , we first perform a best-aligned scaled SE(3) transformation between the ground-truth and predicted pointclouds for the first frame. We adopt

| | HOI4D | | | | H2O | | | | POV-Surgery† | | | | ARCTIC-HOI† | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ADE↓ | FDE↓ | $P_5$ ↑ | $P_{10}$ ↑ | ADE↓ | FDE↓ | $P_5$ ↑ | $P_{10}$ ↑ | ADE↓ | FDE↓ | $P_5$ ↑ | $P_{10}$ ↑ | ADE↓ | FDE↓ | $P_5$ ↑ | $P_{10}$ ↑ |
| DS+UniDepth [77] | <u>64.0</u> | 75.1 | <u>54.0</u> | 82.5 | 46.4 | 48.6 | 67.6 | 95.0 | 168.9 | 199.7 | 13.2 | 39.8 | <u>53.9</u> | **72.1** | **60.1** | <u>86.8</u> |
| MFVR+UniDepth [1] | 88.8 | 94.8 | 23.1 | 69 | <u>40.9</u> | <u>42.8</u> | <u>75.8</u> | **97.0** | 504.5 | 767.6 | 3.2 | 13.6 | 102.0 | 154.6 | 33.8 | 62.0 |
| DUSt3R [86] | 79.2 | 75.8 | 47.7 | 75.1 | 71.7 | 71.5 | 54.3 | 75.9 | 412.8 | 407.2 | <u>14.8</u> | <u>44.9</u> | 214.3 | 181.6 | 52.2 | 82.3 |
| MonSt3R [99] | 70.2 | 71.0 | 50.7 | 80.6 | 96.3 | 97.1 | 27.6 | 67.3 | 201.8 | 208.9 | 8.4 | 31.8 | 83.9 | 106.3 | 28.4 | 69.9 |
| Align3R [45] | 66.3 | <u>67.8</u> | 53.4 | <u>83.7</u> | 75.3 | 75.5 | 39.2 | 79.7 | <u>157.6</u> | <u>163.7</u> | 14.5 | 43.1 | 64.8 | 83.3 | 43.8 | 83.5 |
| EgoMono4D (Ours) | **57.3** | **60.6** | **59.2** | **87.0** | **37.0** | **38.9** | **77.0** | <u>96.4</u> | **125.0** | **138.9** | **19.1** | **55.2** | **53.8** | 72.1 | <u>57.2</u> | **87.3** |

Table 5. The evaluation results for long-term 3D scene flow recovery are presented, with ADE (mm), FDE (mm), and Precision ($P_\delta$, %). † denotes zero-shot generalization for EgoMono4D. For ARCTIC-HOI, the evaluation focuses solely on hand-object recovery quality. Overall, EgoMono4D significantly outperforms the other baselines.



**HOI4D**    **Epic-Kitchen**    **POV-Surgery**    **ARCTIC**

View1    View1    View1    View1

View2    View2    View2    View2

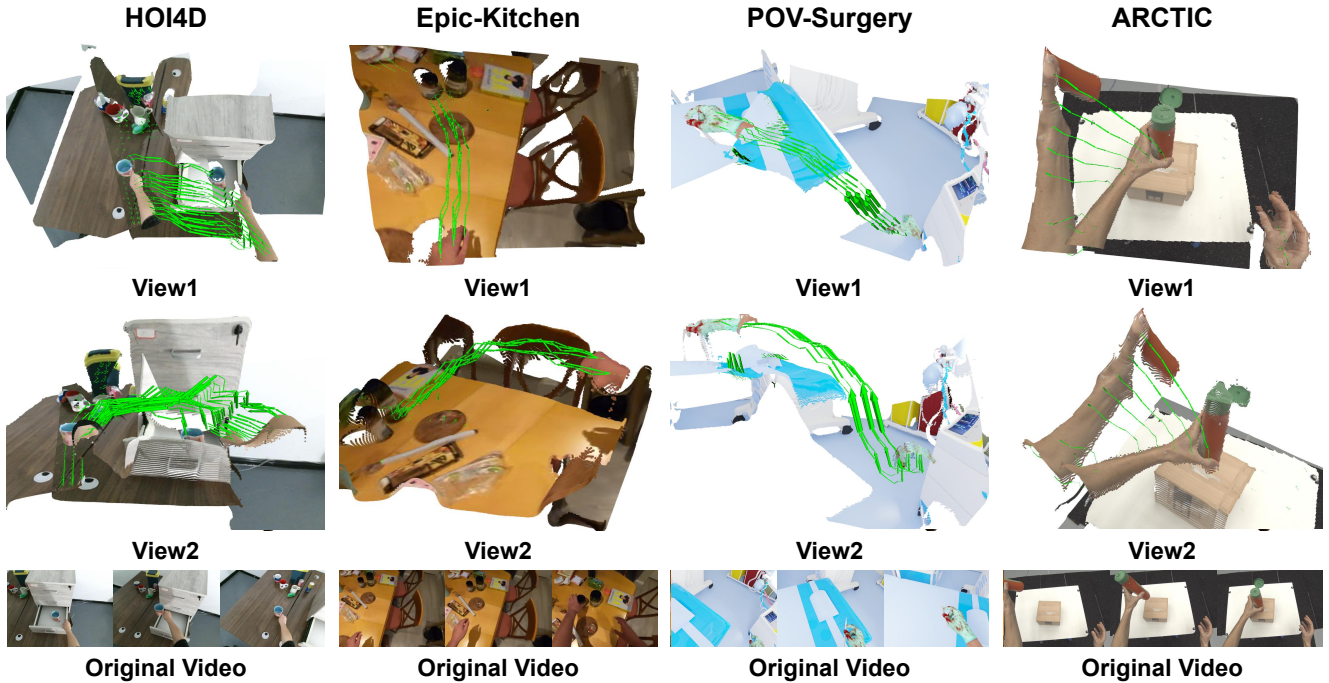Original Video    Original Video    Original Video    Original Video

Figure 6. The visualization of the long-term 3D scene flow recovery. For clarity, we display only the pointclouds from the first and last frame. The green arrows representing the estimated 3D flow. EgoMono4D successfully recovers the motion of dynamic parts while maintaining other regions static to some extent.

| (Only for HOI) | ADE↓ | FDE↓ | $P_5$↑ | $P_{10}$↑ |
|---|---|---|---|---|
| HOI4D | 79.3 / **76.6** | 84.6 / **78.4** | **48.3** / 46.2 | 77.8 / **81.0** |
| H2O | 43.6 / **40.5** | 45.8 / **41.4** | 68.6 / **71.2** | 95.7 / **98.2** |
| POV-Surgery† | 196.0 / **192.2** | 213.4 / **206.4** | **9.9** / **9.9** | 32.4 / **32.8** |

Table 6. Additional long-term 3D scene flow results on the HOI part. Values are presented **in the format 'DS+UniDepth / EgoMono4D (Ours)'**. Both models demonstrate comparable performance, with EgoMono4D showing a modest advantage.

metrics from related works [3, 33, 83, 96], including Average Displacement Error (ADE, mm), Final Displacement Error (FDE, mm), and Precision under Distance (P, %). We use the notation $P_\delta$ to represent precision with a $\delta$ centimeter (cm) threshold. Details are demonstrated in Appendix C.

**Result** Table 5 shows that our model outperforms all baselines for long-term 3D scene flow recovery on average. The visualizations are shown in Figure 6. DS+UniDepth performs well in hand-object motion estimation on the ARCTIC-HOI dataset. We also compared the HOI estimation of both models on three other datasets, with results in Table 6. Both models perform similarly, with EgoMono4D shows a modest advantage. However, DS+UniDepth's performance drops notably in full-scene estimation (e.g., POV-Surgery dataset) due to inconsistencies and accumulated errors between modules. Other models perform worse over-

| $N_w$ | $N_o$ | HOI4D | | | | POV-Surgery | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | CD↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ | CD↓ | $F_1$ ↑ | $F_{2.5}$ ↑ | $F_5$ ↑ |
| **4** | **1** | 5.9 | 27.9 | 59.6 | 83.1 | 33.8 | 13.5 | 32 | 53.9 |
| **2** | **1** | 17.9 | 11.5 | 29.4 | 51.6 | / | / | / | / |
| **8** | **1** | 6.4 | 26.0 | 55.6 | 80.0 | 35.2 | 12.6 | 30.8 | 53.1 |
| **12** | **1** | 6.7 | 25.1 | 53.5 | 78.4 | / | / | / | / |
| **4** | **2** | 5.9 | 27.9 | 59.7 | 83.1 | 34.6 | 13.4 | 32 | 53.9 |
| **4** | **3** | 5.9 | 27.9 | 59.7 | 83.1 | 35.0 | 13.4 | 31.9 | 53.8 |
| **8** | **4** | / | / | / | / | 36.3 | 16.6 | 30.5 | 52.6 |

Table 7. Comparison of pointclouds sequence reconstruction results across different window sizes ($N_w$) and overlapping sizes ($N_o$). Our model demonstrates robustness to variations in $N_o$, while maintaining consistency in $N_w$ between training and inference is essential.

| | POV-Surgery (Video Depth) | | |
|---|---|---|---|
| | AbsRel↓ | $\delta_{0.05}$ ↑ | $\delta_{0.1}$ ↑ |
| UniDepth [56] | **11.9** | 40.8 | **64.7** |
| DUSt3R [86] | 19.7 | 31.5 | 51.9 |
| MonSt3R [99] | 18.6 | 33.6 | 52.8 |
| Align3R [45] | 13.3 | **41.7** | 62.6 |
| EgoMono4D (Ours) | 12.6 | 41.1 | 63.9 |

| | POV-Surgery (Camera Poses) | | |
|---|---|---|---|
| | ATE↓ | RPE-T↓ | RPE-R↓ |
| DS+UniDepth [77] | 9.05 | 4.17 | 0.39 |
| MFVR+UniDepth [1] | 47.03 | 4.10 | 4.18 |
| MonSt3R [99] | 6.63 | 2.41 | 0.26 |
| Align3R [45] | **6.35** | **2.34** | **0.23** |
| EgoMono4D (Ours) | 11.54 | 4.01 | 0.43 |

Table 8. Result of POV-Surgery video depth and camera poses estimation. It shows that our model does not outperform others in estimating these geometric variables. Instead, our model focus on improving their consistency in 3D space, which leads to better pointclouds sequence reconstruction.

all. Figure 6 illustrates the estimated long-term 3D scene flow. It can be seen that EgoMono4D successfully reconstructs 3D dynamics across diverse scenes.

## E. Ablation on Inference Strategy

During inference, our model processes $N_w$ frames in a single feed-forward prediction. Theoretically, the window size $N_w$ can be any value greater than 1. For videos with more frames than $N_w$, the overlapping size $N_o$ between neighboring windows must also be determined. By default, we set $N_w = 4$ (consistent with training) and $N_o = 1$ (to optimize inference speed). We evaluate the impact of $N_w$ and $N_o$ on reconstruction performance using the pointclouds sequence

reconstruction task on HOI4D [43] and POV-Surgery [85], with results shown in Table 7.

For the window size $N_w$, maintaining consistency between training and inference is crucial, likely because the video adapter is trained specifically to fuse 4 frames. Regarding overlapping size $N_o$, the model exhibits comparable performance for $N_o = 1, 2, 3$. Therefore, we select $N_o = 1$ to maximize inference speed.

## F. Depth and Camera Results

We further evaluate the video depth and camera poses estimation performance on POV-Surgery [85], using the metrics from MonSt3R [99]. The results are presented in Table 8. Our model does not outperform other methods in estimating these independent geometric variables. The depth and camera prediction performance of EgoMono4D is only comparable with other baseline methods. Instead, it enhances the consistency of them in 3D space, leading to improved pointclouds sequence reconstruction results. UniDepth [56] achieves the best depth estimation, while Align3R [45] provides the most accurate camera poses.

## G. More Reconstruction Results Visualization

We provide more visualization of pointclouds sequence reconstruction results of EgoMono4D in Figure 7. More visualization of recovered long-term 3D scene flows are shown in Figure 8. Finally, we provide the visualization of intermediate geometric variables in Figure 10.

## H. Qualitative Comparison with Baseline

We qualitatively compare EgoMono4D with other baseline methods on the dense point cloud sequence reconstruction task. Visual comparisons are presented in Figure 9. EgoMono4D demonstrates superior performance compared to baseline methods in both static scene reconstruction and dynamic HOI motion recovery. DS+UniDepth [56, 77]
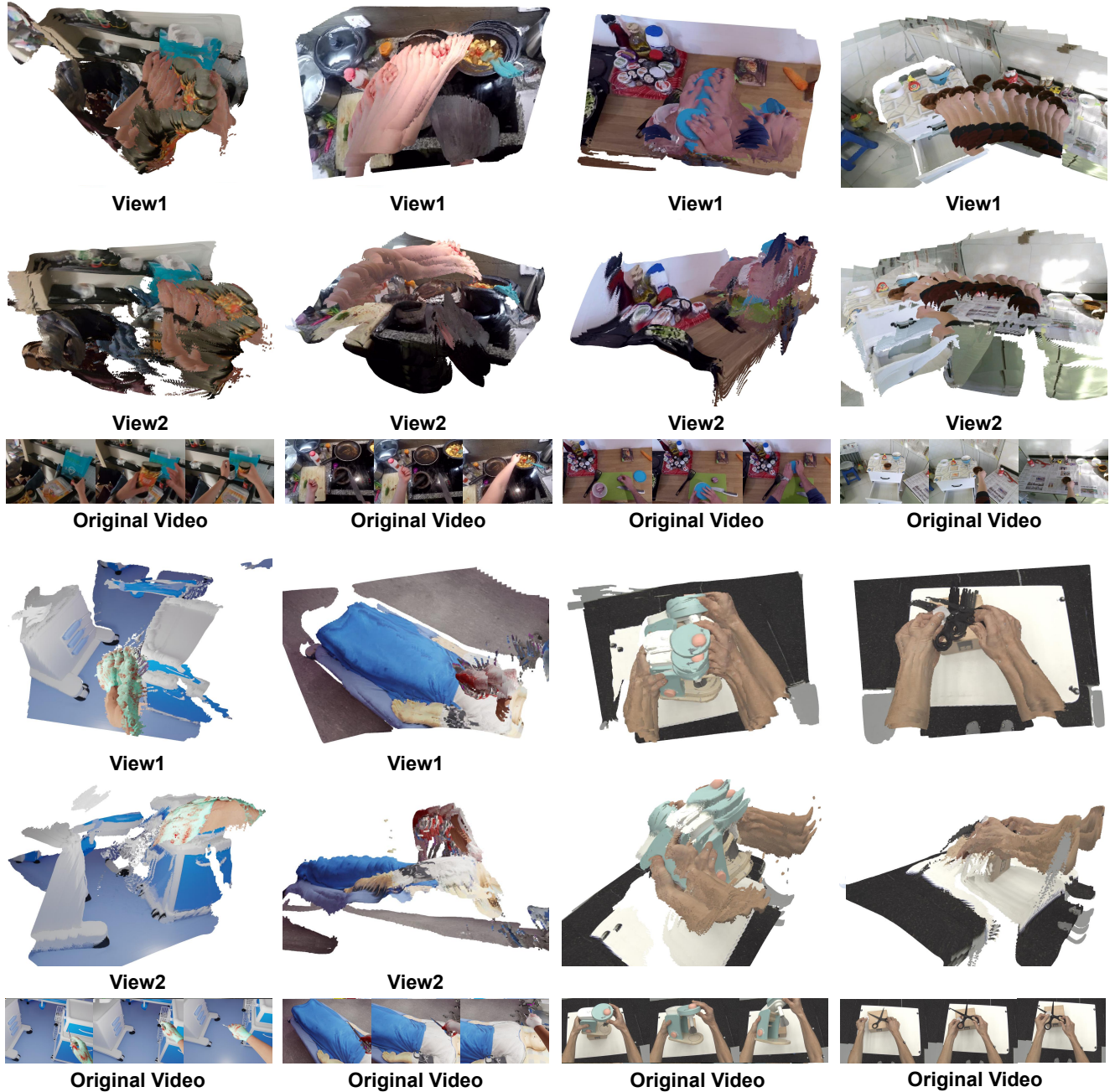
Figure 7. More visualization of pointclouds sequence reconstruction results from EgoMono4D.

struggles with static scene alignment, while MonSt3R [99] exhibits limitations in accurately reconstructing hand geometry and dynamics.

# I. Limitations and Future Directions

While EgoMono4D achieves impressive results in fast, dense, and generalizable dynamic HOI scene reconstruction, it still faces challenges with shape misalignment (see

Figure 5 in the main paper). This issue arises from inherent inconsistencies in UniDepth's [56] shape regularization and can be divided into two key problems. (1) Dynamic part (size) distortion: since we only regularize the shape of the dynamic part based on per-frame pointclouds predictions from UniDepth, the relative size of the dynamic part compared to the static part is determined by UniDepth's predictions. Any inaccuracy in this relative size may be carried over to EgoMono4D. (2) Static part misalignment: if the

View1   View1   View1

View2   View2   View2

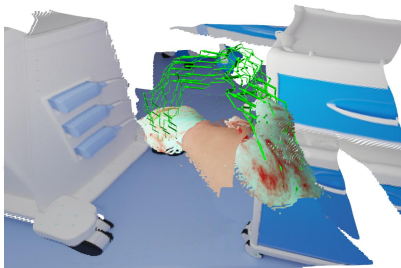Original Video  Original Video  Original Video
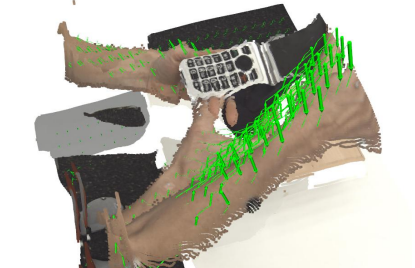
View1   View1   View1

View2   View2   View2

Original Video  Original Video  Original Video

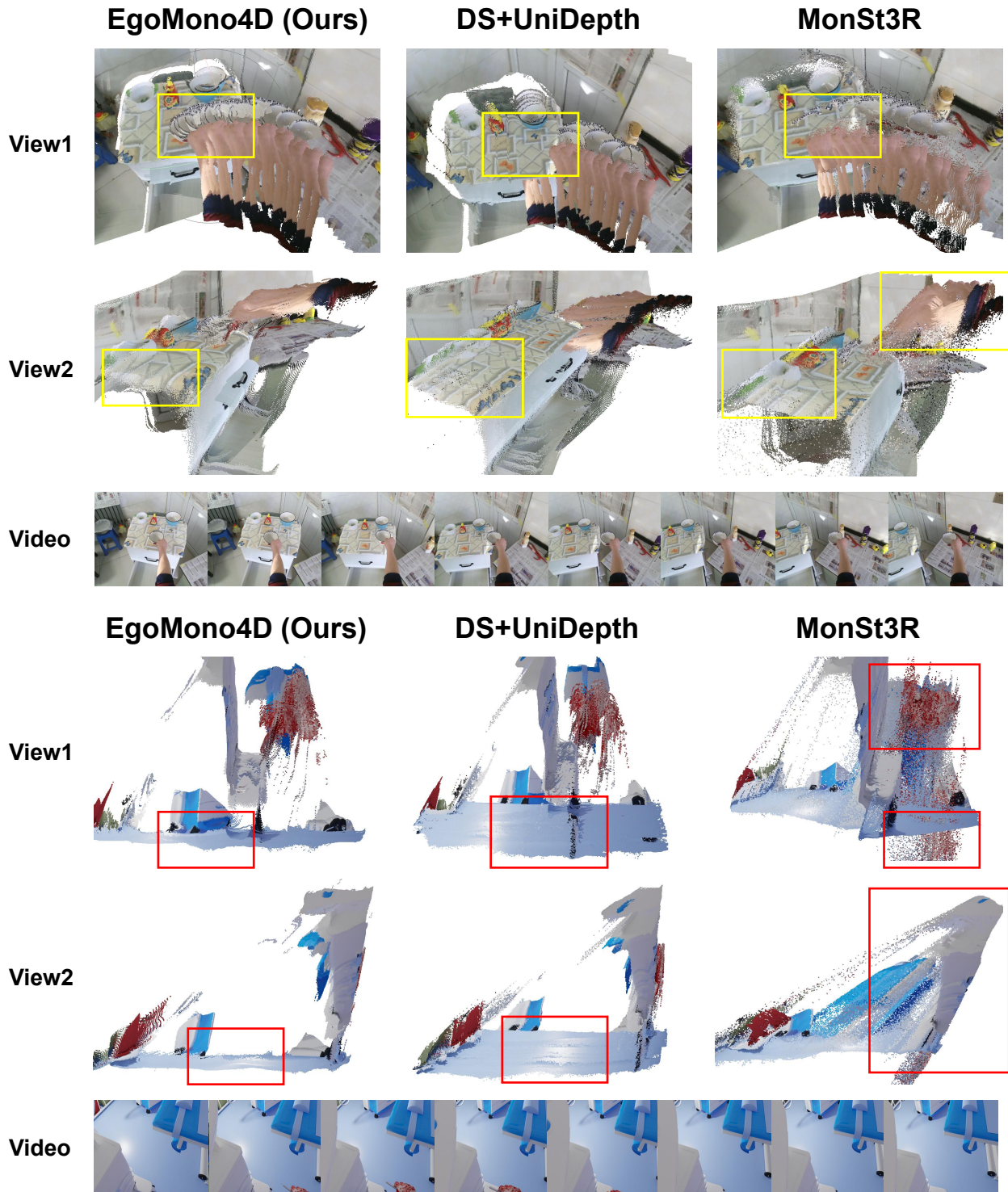Figure 8. More visualization of long-term 3D scene flow recovery results from EgoMono4D.

Figure 9. Qualitative comparison between EgoMono4D and other baseline methods. EgoMono4D demonstrates superior performance compared to baseline methods in both static scene reconstruction and dynamic HOI motion recovery. DS+UniDepth [56, 77] struggles with static scene alignment, while MonSt3R [99] exhibits limitations in accurately reconstructing hand geometry.
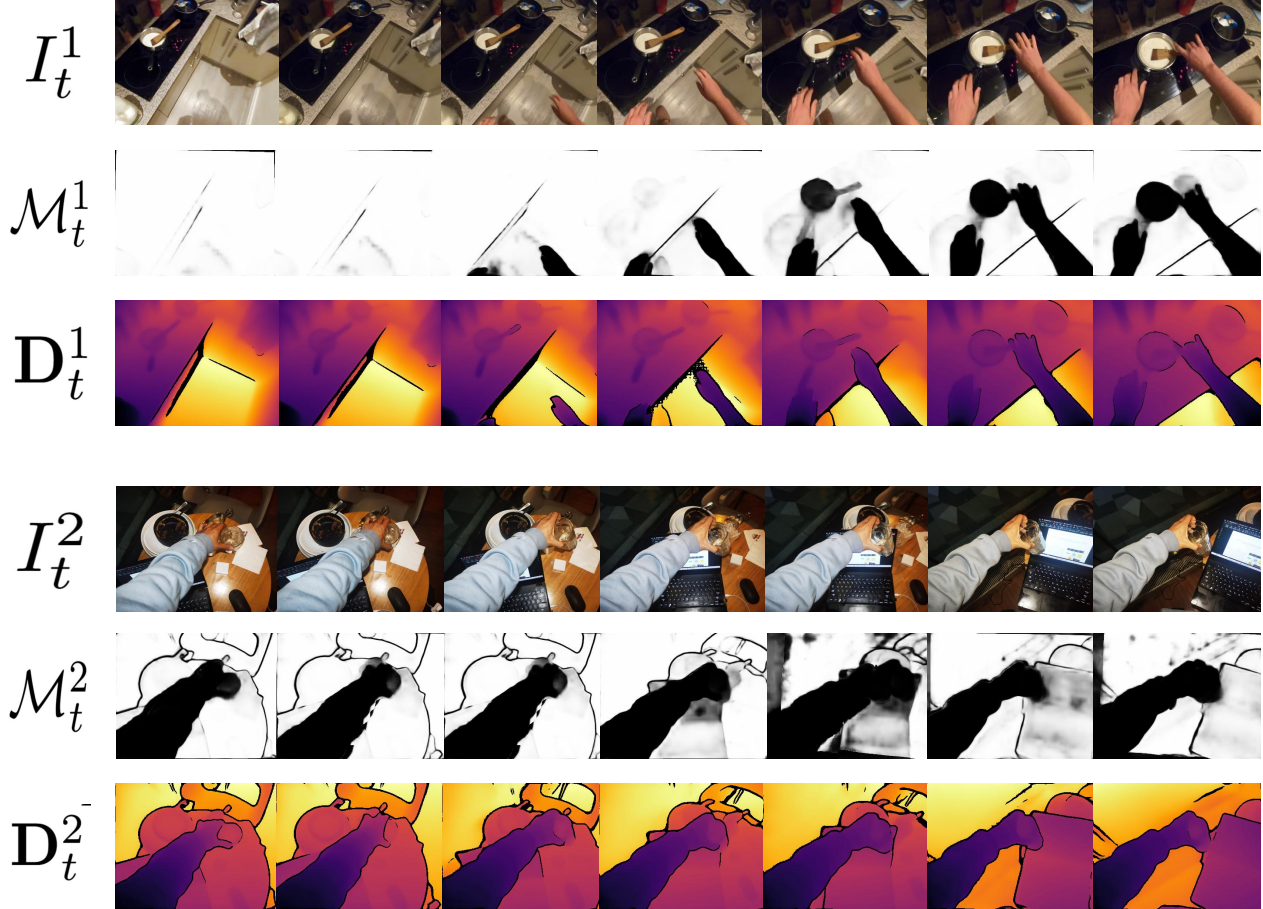
Figure 10. Visualization of original video $I_t$ and the predicted procrustes-alignment confidence maps $\mathcal{M}_t$ and video depth $D_t$ from EgoMono4D.
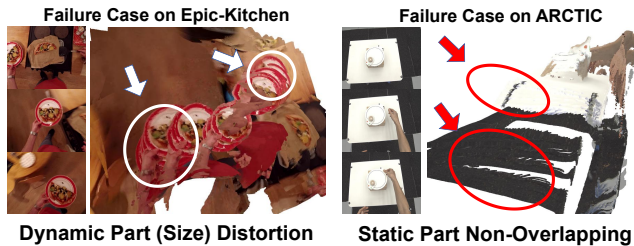


Figure 11. Two typical failure cases of EgoMono4D arise from inherent inconsistencies in UniDepth shape supervision.

original predicted shapes of the static areas from UniDepth differ significantly between frames, it becomes difficult to adjust and align them consistently.

Although precise posed datasets are limited, there are more datasets available with depth labels. Training on these datasets with ground-truth shape supervision, or using a combination of labeled and unlabeled datasets, could help address these issues. Improvements in monocular depth es-

timation and intrinsic parameter estimation may also alleviate these problems.

Another limitation of EgoMono4D is that it currently supports video reconstruction only for videos with an FPS above a certain threshold. This is due to its reliance on the off-the-shelf optical flow estimation module [92], which may perform poorly with sparse views. Integrating correspondence, matching, or optical flow prediction directly into the model to enable fully end-to-end training could address this limitation [36, 82]. Sparse-view data also needs to be incorporated during training [86] to solve this problem.

Additionally, our model currently only supports a resolution of 288×384. Training models at higher resolutions could enable more flexible applications. Our model also fails when the dynamic portion of the image is too large or contains too many moving objects beyond the ones being manipulated. Although we focus on egocentric HOI scenes, our training paradigm could be extended to more general cases. We plan to train this general scene model using motion mask priors derived from motion segmentation

[10, 91], salient video segmentation [74], and semantic segmentation [31, 41]. We also encourage the community to propose more synthetic datasets [54, 104] to explore supervised learning approaches [86, 99].