



dynamic objects, and disentangling ego-motion from object motion through the inherent attention mechanisms of the brain [58]. Yet, the learning process rarely relies on explicit dynamic labels.

We observe that DUST3R implicitly learned a similar mechanism, and based on this, we introduce Easi3R, a training-free method to achieve dynamic object segmentation, dense point map reconstruction, and robust camera pose estimation from dynamic videos, as shown in Figure 1. DUST3R uses attention layers at its core, taking two image features as input and producing pixel-aligned point maps as output. These attention layers are trained to directly predict pointmaps in the reference view coordinate space, implicitly matching the image features between the input views [4] and estimating the rigid view transformation in the feature space. In practice, performance drops significantly when processing pairs with object dynamics [73], as shown in Figure 2. By analyzing the attention maps in the transformer layers, we find that regions with less texture, under-observed, and dynamic objects can yield low attention values. Therefore, we propose a simple yet effective decomposition strategy to isolate the above components, which enables long-horizon dynamic object detection and segmentation. With this segmentation, we perform a second inference pass by applying a re-weighting [17] in the cross-attention layers, enabling robust dynamic 4D reconstruction and camera motion recovery without fine-tuning on a dynamic dataset, all at minimal additional cost to DUST3R.

Despite its simplicity, we demonstrate that our inference-time scaling approach for 4D reconstruction is remarkably robust and accurate on in-the-wild casual dynamic videos. We evaluate our Easi3R adaptation on the DUST3R and MonST3R backbones in three task categories: camera pose estimation, dynamic object segmentation, and pointcloud reconstruction in dynamic scenes. Easi3R performs surprisingly well across a wide range of datasets, even surpassing concurrent methods (e.g., CUT3R [63], MonST3R [73], and DAS3R [68]) that are trained on dynamic datasets.

## 2. Related Work

**SfM and SLAM.** Structure-from-Motion (SfM) [2, 41, 42, 48, 51, 52] and Simultaneous Localization and Mapping (SLAM) [9, 13, 32, 34] have long been the foundation for 3D structure and camera pose estimation. These methods are done by associating 2D correspondences [5, 10, 28, 32, 47] or minimizing photometric errors [12, 13], followed by bundle adjustment (BA) [3, 6, 55, 57, 59, 62] to refine structure and motion estimates. Although highly effective with dense input, these approaches often struggle with limited camera parallax or ill-posed conditions, leading to performance degeneracy. To overcome these limitations, DUST3R [64] introduced a learning-based approach that di-

rectly predicts two pointmaps from an image pair in the coordinate space of the first view. This approach inherently matches image features and rigid body view transformation. By leveraging a Transformer-based architecture [11] and direct point supervision on large-scale 3D datasets, DUST3R establishes a robust Multi-View Stereo (MVS) foundation model. However, DUST3R and the follow-up methods [27, 33, 56, 61] assume primarily static scenes, which can lead to significant performance degradation when dealing with videos with dynamic objects.

**Pose-free Dynamic Scene Reconstruction.** Modifications to SLAM for dynamic scenes involve robust pose estimation to mitigate moving object interference, dynamic map management for updating changing environments, including techniques like semantic segmentation [72], optical flows [75], enhance SLAM’s resilience in dynamic scenarios. Another line of work focuses on estimating stable video depth by incorporating geometric constraints [29] and generative priors [18, 49]. These methods enhance monocular depth accuracy but lack global point cloud lifting due to missing camera intrinsics and poses. For joint pose and depth estimation, optimization-based methods such as CasualSAM [74] fine-tune a depth network [45] at test time using pre-computed optical flow [66]. Robust-CVD [22] refines pre-computed depth [45] and camera pose by leveraging masked optical flow [16, 66] to improve stability in occluded and moving regions. Concurrently, MegaSaM [25] further enhances pose and depth accuracy by integrating DROID-SLAM [57], optical flow [66], and depth initializations from [40, 71], achieving state-of-the-art results. Alternatively, point-map-based approaches like MonST3R [73] extend DUST3R to dynamic scenes by fine-tuning with dynamic datasets and incorporating optical flow [66] to infer dynamic object segmentation. DAS3R trains a DPT [44] on top of MonST3R, enabling feedforward segmentation estimation. CUT3R [63] fine-tunes MAST3R [24] on both static and dynamic datasets, achieving feedforward reconstruction but without predicting dynamic object segmentation, thereby entangling the static scene with dynamic objects. Although effective, these methods require costly training on diverse motion patterns to generalize well.

In contrast, we take an opposite path, exploring a training-free and plug-in-play adaptation that enhances the generalization of DUST3R variants for dynamic scene reconstruction. Our method requires no fine-tuning and comes at almost no additional cost, offering a scalable and efficient alternative for handling real-world dynamic videos.

**Motion Segmentation.** Motion segmentation aims to predict dynamic object masks from video inputs. Classical approaches generally rely on optical flow estimation [26, 31, 67, 70] and point tracking [7, 21, 35, 50, 69] to distinguish moving objects from the background. Being trained solely on 2D data, they often struggle with occlusions and