# GLAMR: Global Occlusion-Aware Human Mesh Recovery with Dynamic Cameras

Ye Yuan[1*]     Umar Iqbal[2]     Pavlo Molchanov[2]     Kris Kitani[1]     Jan Kautz[2]

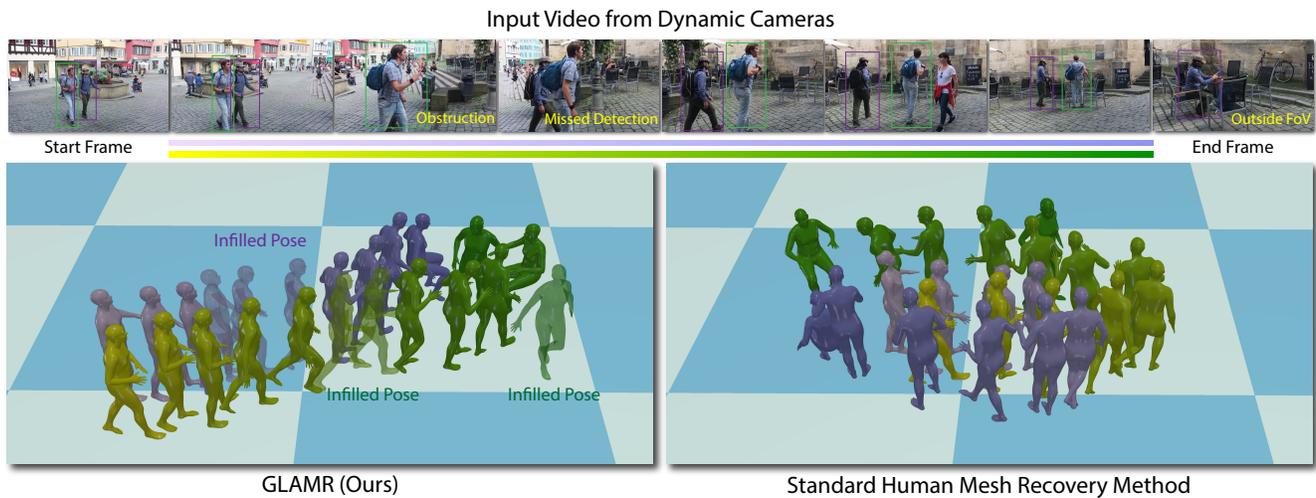[1]Carnegie Mellon University     [2]NVIDIA

https://www.ye-yuan.com/glamr

Figure 1. GLAMR (**Left**) recovers human meshes in consistent *global* coordinates and *infills missing poses* (transparent) due to various occlusions (obstruction, missed detection, outside field of view), while standard human mesh recovery methods (**Right**) fail to do so.

## Abstract

*We present an approach for 3D global human mesh recovery from monocular videos recorded with dynamic cameras. Our approach is robust to severe and long-term occlusions and tracks human bodies even when they go outside the camera's field of view. To achieve this, we first propose a deep generative motion infiller, which autoregressively infills the body motions of occluded humans based on visible motions. Additionally, in contrast to prior work, our approach reconstructs human meshes in consistent global coordinates even with dynamic cameras. Since the joint reconstruction of human motions and camera poses is underconstrained, we propose a global trajectory predictor that generates global human trajectories based on local body movements. Using the predicted trajectories as anchors, we present a global optimization framework that refines the predicted trajectories and optimizes the camera poses to match the video evidence such as 2D keypoints. Experiments on challenging indoor and in-the-wild datasets with dynamic cameras demonstrate that the proposed approach outperforms prior methods significantly in terms of motion infilling and global mesh recovery.*

## 1. Introduction

Recovering fine-grained 3D human meshes from monocular videos is essential for understanding human behaviors and interactions, which can be the cornerstone for numerous applications including virtual or augmented reality, assistive living, autonomous driving, *etc*. Many of these applications use dynamic cameras to capture human behaviors yet also require estimating human motions in global coordinates consistent with their surroundings. For instance, assistive robots and autonomous vehicles need a holistic understanding of human behaviors and interactions in the world to safely plan their actions even when they are moving. Therefore, our goal in this paper is to tackle the important task of recovering global human meshes from monocular videos captured by dynamic cameras.

However, this task is highly challenging for two main reasons. First, dynamic cameras make it difficult to estimate human motions in *consistent global coordinates*. Existing human mesh recovery methods estimate human meshes in

---

the camera coordinates [65, 115] or even in the root-relative coordinates [45, 66]. Hence, they can only recover global human meshes from dynamic cameras by using SLAM to estimate camera poses [58]. However, SLAM can often fail for in-the-wild videos due to moving and dynamic objects. It also has the problem of scale ambiguity, which often leads to camera poses that are inconsistent with the human motions. Second, videos captured by dynamic cameras often contain *severe and long-term occlusions* of humans, which can be caused by missed detection, complete obstruction by objects and other people, or the person going outside the camera's field of view (FoV). These occlusions pose serious challenges to standard human mesh recovery methods, which rely on detections or visible parts to estimate human meshes. Only a few works have attempted to tackle the occlusion problem in human mesh recovery [17, 36]. However, these methods can only address partial occlusions of a person and fail to handle severe occlusions when the person is completely invisible for an extended period of time.

To tackle the above challenges, we propose Global Occlusion-Aware Human Mesh Recovery (GLAMR), which can handle severe occlusions and estimate human meshes in consistent global coordinates – even for videos recorded with dynamic cameras. We start by using off-the-shelf methods (*e.g.*, KAMA [33] or SPEC [47]) to estimate the shape and pose sequences (motions) of visible people in the camera coordinates. These methods also rely on multi-object tracking and re-identification, which provide occlusion information, and the motion of occluded frames is not estimated. To tackle potentially severe occlusions, we propose a deep generative motion infiller that autoregressively infills the local body motions of occluded people based on visible motions. The motion infiller leverages human dynamics learned from a large motion database, AMASS [60]. Next, to obtain global motions, we propose a global trajectory predictor that can generate global human trajectories based on local body motions. It is motivated by the observation that the global root trajectory of a person is highly correlated with the local body movements. Finally, using the predicted trajectories as anchors to constrain the solution space, we further propose a global optimization framework that jointly optimizes the global motions and camera poses to match the video evidence such as 2D keypoints.

The contributions of this paper are as follows: **(1)** We propose the first approach to address long-term occlusions and estimate global 3D human pose and shape from videos captured by dynamic cameras; **(2)** We propose a novel generative Transformer-based motion infiller that autoregressively infills long-term missing motions, which considerably outperforms state-of-the-art motion infilling methods; **(3)** We propose a method to generate global human trajectories from local body motions and use the generated trajectories as anchors to constrain global motion and camera

optimization; **(4)** Extensive experiments on challenging indoor and in-the-wild datasets demonstrate that our approach outperforms prior state-of-the-art methods significantly in tackling occlusions and estimating global human meshes.

## 2. Related Work

**Camera-Relative Pose Estimation.** 3D human mesh recovery from RGB images or videos is an ill-posed problem due to the depth ambiguity. Most existing methods simplify the problem by estimating human poses relative to the pelvis (root) of the human body [1, 6, 8–10, 21, 37, 39, 40, 45, 48–52, 57, 66, 67, 70–72, 78, 81, 86, 89, 90, 97, 99, 105, 109, 112, 118]. These methods assume an orthographic camera projection model and neglect the absolute 3D translation of the person w.r.t. the camera. To address the lack of translation, recent methods start to estimate human meshes in the camera coordinates [33, 36, 53, 58, 74, 77, 84, 98, 106, 108, 110]. Several approaches recover the absolute translation of the person using an optimization framework [62–64, 80, 107]. A few methods exploit various scene constraints during the optimization process to improve depth prediction [95, 106]. Alternatively, recent approaches use physics-based constraints to ensure the physical plausibility of the estimated poses [12, 34, 84, 98, 104]. Iqbal *et al.* [32] exploit a limb-length constraint to recover the absolute translation of the person using a 2.5D representation. Some approaches approximate the depth of the person using the bounding box size [36, 65, 110]. HybrIK [53] and KAMA [33] employ inverse kinematics to estimate human meshes with absolute translations in the camera coordinates. Several methods directly predict the absolute depth of each person using a heatmap representation [16, 115]. Recently, SPEC [47] learns to predict the camera parameters (pitch, yaw, FoV) from the image, which are used for absolute pose regression in the camera coordinates. THUNDR [108] also adopts a similar strategy but uses known camera parameters. While these methods show impressive results, they cannot estimate global human motions from videos captured by dynamic cameras. In contrast, our approach can recover human meshes in consistent global coordinates for dynamic cameras and handle severe and long-term occlusions.

**Global Pose Estimation.** Most existing methods that estimate 3D poses in world coordinates rely on calibrated, synchronized, and static multi-view capture setups [5, 13, 15, 29, 38, 76, 77, 113, 114, 116]. Huang *et al.* [7] use uncalibrated cameras but still assume time synchronization and static camera setups. Hasler *et al.* [24] handle unsynchronized moving cameras but assume multi-view input and rely on audio stream for synchronization. More recently, Dong *et al.* [14] propose to recover 3D poses from unaligned internet videos of different actors performing the same activity from unknown cameras. However, they assume that multi-
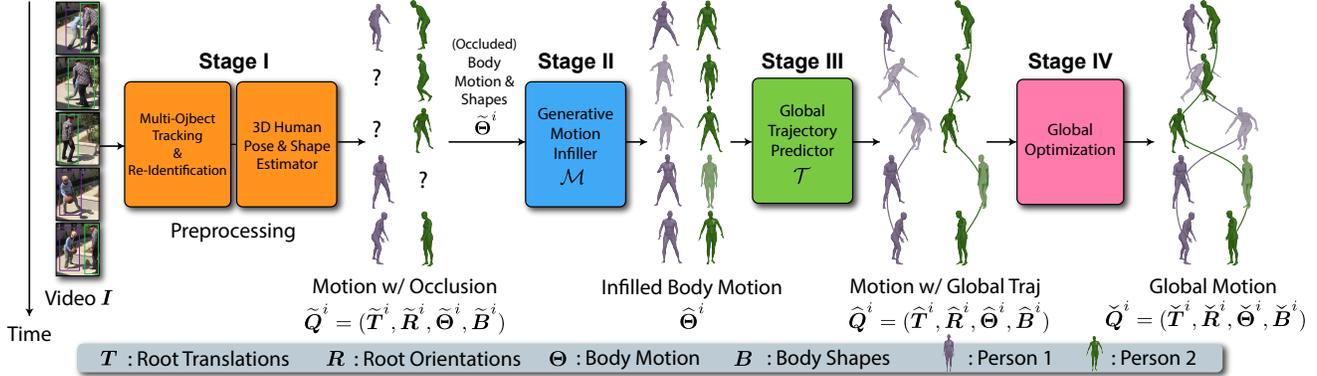
Figure 2. **Overview** of GLAMR. In **Stage I**, we preprocess the video with multi-object tracking, re-identification and human mesh recovery to extract each person's occluded motion $\widetilde{\boldsymbol{Q}}^i$ in the camera coordinates. In **Stage II**, we propose a generative motion infiller to infill the occluded body motion $\widetilde{\boldsymbol{\Theta}}^i$ to produce occlusion-free body motion $\widehat{\boldsymbol{\Theta}}^i$. In **Stage III**, we propose a global trajectory predictor that uses the infilled body motion $\widehat{\boldsymbol{\Theta}}^i$ to generate the global trajectory $(\widehat{\boldsymbol{T}}^i, \widehat{\boldsymbol{R}}^i)$ of each person and obtain their global motion $\widehat{\boldsymbol{Q}}^i$. In **Stage IV**, we jointly optimize the global trajectories of all people and the camera parameters to produce global motions $\check{\boldsymbol{Q}}^i$ consistent with the video.

ple viewpoints of the same pose are available in the videos. Different from these methods, our approach estimates human meshes in global coordinates from *monocular* videos recorded with dynamic cameras. Several methods rely on additional IMU sensors or pre-scanned environments to recover global human motions [22, 94], which is unpractical for large-scale adoption. Recently, Liu *et al.* [58] first obtain the camera poses and dense reconstruction of the scene from dynamic cameras using a SLAM algorithm, COLMAP [82]. The camera poses are used for camera-to-world transformation, while the reconstructed scene is used to encourage human-scene contacts. However, SLAM can often fail for the in-the-wild videos and is prone to error propagation. In contrast, our approach does not require SLAM but instead uses global trajectory prediction to constrain the joint reconstruction of human motions and camera poses. Additionally, our approach can also handle severe and long-term occlusions common in dynamic camera setups.

**Occlusion-Aware Pose Estimation.** Most existing human pose estimation methods assume the person is fully visible in the images and are not robust to strong occlusions. Only a few methods address the occlusion problem in pose estimation [17, 46, 78, 79, 112]. While these methods show impressive results under partial occlusions, they do not address severe and long-term occlusions when people are completely obstructed or outside the camera's FoV for a long time. In contrast, our approach leverages deep generative human motion models to tackle severe and long-term occlusions.

**Human Motion Modeling.** Extensive research has studied 3D human dynamics for various tasks including motion prediction and synthesis [2, 4, 18, 19, 25, 35, 56, 61, 73, 75, 93, 100–103]. Recent human pose estimation methods start to leverage learned human dynamics models to improve the accuracy of estimated motions [45, 78, 111]. Several motion

infilling approaches are also proposed to generate complete motions from partially observed motions [23, 28, 41, 42]. Additionally, recent work on motion capture shows that global human translations can be predicted from 3D local joint positions [83]. In contrast to prior work, our trajectory predictor does not require GT root orientations but can predict both global root translations and orientations. Furthermore, we also propose a novel generative autoregressive motion infiller that can use noisy poses as input instead of high-quality GT poses, and we demonstrate its effectiveness in tackling long-term occlusions in human pose estimation.

## 3. Method

The input to our framework is a video $\boldsymbol{I} = (\boldsymbol{I}_1, \ldots, \boldsymbol{I}_T)$ with $T$ frames, which is captured by a *dynamic camera*, *i.e.*, the camera poses can change every frame. Our goal is to estimate the global motion (pose sequence) $\{\boldsymbol{Q}^i\}_{i=1}^N$ of the $N$ people in the video in a *consistent global coordinate* system. The global motion $\boldsymbol{Q}^i = (\boldsymbol{T}^i, \boldsymbol{R}^i, \boldsymbol{\Theta}^i, \boldsymbol{B}^i)$ for person $i$ consists of the root translations $\boldsymbol{T}^i = (\boldsymbol{\tau}_{s_i}^i, \ldots, \boldsymbol{\tau}_{e_i}^i)$, root rotations $\boldsymbol{R}^i = (\boldsymbol{\gamma}_{s_i}^i, \ldots, \boldsymbol{\gamma}_{e_i}^i)$, as well as the body motion $\boldsymbol{\Theta}^i = (\boldsymbol{\theta}_{s_i}^i, \ldots, \boldsymbol{\theta}_{e_i}^i)$ and shapes $\boldsymbol{B}^i = (\boldsymbol{\beta}_{s_i}^i, \ldots, \boldsymbol{\beta}_{e_i}^i)$, where the motion spans from the the first frame $s_i$ to the last frame $e_i$, when the person $i$ is relevant in the video. In particular, each body pose $\boldsymbol{\theta}_t^i \in \mathbb{R}^{23 \times 3}$ and shape $\boldsymbol{\beta}_t^i \in \mathbb{R}^{10}$ corresponds to the pose parameters (excluding root rotation) and shape parameters of the SMPL model [59]. Using the root translation $\boldsymbol{\tau} \in \mathbb{R}^3$ and (axis-angle) rotation $\boldsymbol{\gamma} \in \mathbb{R}^3$, SMPL represents a human body mesh with a linear function $\mathcal{S}(\boldsymbol{\tau}, \boldsymbol{\gamma}, \boldsymbol{\theta}, \boldsymbol{\beta})$ that maps a global pose $\boldsymbol{q} = (\boldsymbol{\tau}, \boldsymbol{\gamma}, \boldsymbol{\theta}, \boldsymbol{\beta})$ to an articulated triangle mesh $\boldsymbol{\Phi} \in \mathbb{R}^{K \times 3}$ with $K = 6980$ vertices. We can therefore recover the global mesh sequence for each person from their global motion $\boldsymbol{Q}^i$ via SMPL.
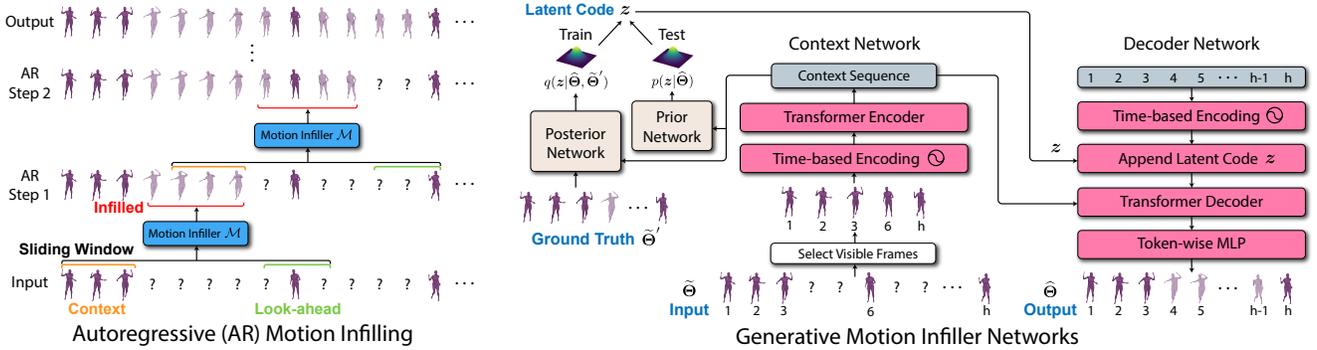
As outlined in Fig. 2, our framework consists of four

3

**Figure 3. Left:** We autoregressively infill the motion using a sliding window, where the first $h_c$ frames are already infilled to serve as context and the last $h_l$ frames are look-ahead to guide the ending motion. Frames between the context and look-ahead are infilled. **Right:** The CVAE-based motion infiller adopts a Transformer-based seq2seq architecture, where we encode only the visible frames of occluded body motion $\widetilde{\Theta}$ into a context sequence, which is used jointly with latent code $z$ by a decoder network to generate occlusion-free motion $\widehat{\Theta}$.

stages. In **Stage I**, we first use multi-object tracking (MOT) and re-identification algorithms to obtain the bounding box sequence of each person, which is input to a human mesh recovery method (*e.g.*, KAMA [33] or SPEC [47]) to extract the motion $\widetilde{Q}^i$ of each person (including translation) in the camera coordinates. The motion $\widetilde{Q}^i$ may be incomplete due to various occlusions (*e.g.*, obstruction, missed detection, going outside FoV), where bounding boxes from MOT are missing for some frames. In **Stage II** (Sec. 3.1), we propose a generative motion infiller to tackle the occlusions in the estimated body motion $\widetilde{\Theta}^i$ and produce occlusion-free body motion $\widehat{\Theta}^i$. In **Stage III** (Sec. 3.2), we propose a global trajectory predictor that uses the infilled body motion $\widehat{\Theta}^i$ to generate the global trajectory (root translations and rotations) of each person and obtain their global motion $\widehat{Q}^i$. In **Stage IV** (Sec. 3.3), we jointly optimize the global trajectories of all people and the camera parameters to produce global motions $\check{Q}^i$ consistent with the video evidence.

## 3.1. Generative Motion Infiller

The task of the generative motion infiller $\mathcal{M}$ is to infill the occluded body motion $\widetilde{\Theta}^i$ of each person to produce occlusion-free body motion $\widehat{\Theta}^i$. Here, we do not use the motion infiller $\mathcal{M}$ to infill other components in the estimated motion $\widehat{Q}^i$, *i.e.*, root trajectory $(\widetilde{T}^i, \widetilde{R}^i)$ and shapes $\widetilde{B}^i$. This is because it is difficult to infill the root trajectory $(\widetilde{T}^i, \widetilde{R}^i)$ using learned human dynamics, since it resides in the camera coordinates rather than a consistent coordinate system due to the dynamic camera. In Sec. 3.2, we will use the proposed global trajectory predictor to generate occlusion-free global trajectory $(\widehat{T}^i, \widehat{R}^i)$ from the infilled body motion $\widehat{\Theta}^i$. The trajectory $(\widetilde{T}^i, \widetilde{R}^i)$ from the pose estimator is not discarded and will be used in the global optimization (Sec. 3.3). For the shapes, we use linear interpo-

lation to produce occlusion-free shapes $\widehat{B}^i$ since a person's shape should stay close to a constant throughout the video.

Given a general occluded human body motion $\widetilde{\Theta} = (\widetilde{\theta}_1, \ldots, \widetilde{\theta}_h)$ of $h$ frames and its visibility mask $V = (V_1, \ldots, V_h)$ as input, the motion infiller $\mathcal{M}$ outputs a complete occlusion-free motion $\widehat{\Theta} = (\widehat{\theta}_1, \ldots, \widehat{\theta}_h)$. The visibility mask $V$ encodes the visibility of the occluded motion $\widetilde{\Theta}$, where $V_t = 1$ if the body pose $\widetilde{\theta}_t$ is visible in frame $t$ and $V_t = 0$ otherwise. Since the human pose for occluded frames can be highly uncertain and stochastic, we formulate the motion infiller $\mathcal{M}$ using the conditional variational autoencoder (CVAE) [44]:

$$\widehat{\Theta} = \mathcal{M}(\widetilde{\Theta}, V, z), \qquad (1)$$

where the motion infiller $\mathcal{M}$ corresponds to the CVAE decoder and $z$ is a Gaussian latent code. We can obtain different occlusion-free motions $\widehat{\Theta}$ by varying $z$.

**Autoregressive Motion Infilling.** To ensure that the motion infiller $\mathcal{M}$ can handle much longer test motions than the training motions, we propose an autoregressive motion infilling process at test time as illustrated in Fig. 3 (Left). The key idea is to use a sliding window of $h$ frames, where we assume the first $h_c$ frames of motion are already occlusion-free or infilled and serve as *context*, and we also use the last $h_l$ frames as *look-ahead*. The look-ahead is essential to the motion infiller since it may contain visible poses that can guide the ending motion and avoid generating discontinuous motions. Excluding the context and look-ahead frames, only the middle $h_o = h - h_c - h_l$ frames of motion are infilled. We iteratively infill the motion using the sliding window and advance the window by $h_o$ frames every step.

**Motion Infiller Network.** The overall network design of the CVAE-based motion infiller is outlined in Fig. 3 (Right). In particular, we employ a Transformer-based seq2seq architecture, which consists of three parts: (1) a *context network* that uses a Transformer encoder to encode the visi-

ble poses from the occluded motion $\widetilde{\boldsymbol{\Theta}}$ into a context sequence, which serves as the condition for other networks; (2) a *decoder network* that uses the latent code $\boldsymbol{z}$ and context sequence to generate occlusion-free motion $\widehat{\boldsymbol{\Theta}}$ via a Transformer decoder and a multilayer perceptron (MLP); (3) *prior and posterior networks* that generate the prior and posterior distributions for the latent code $\boldsymbol{z}$. In the networks, we adopt a time-based encoding that replaces the position in the original positional encoding [92] with the time index. Unlike prior CNN-based methods [28, 41], our Transformer-based motion infiller does not require padding missing frames, but instead restricts its attention to visible frames to achieve effective temporal modeling.

**Training.** We train the motion infiller $\mathcal{M}$ using a large motion capture dataset, AMASS [60]. To synthesize occluded motions $\widetilde{\boldsymbol{\Theta}}$, for any GT training motion $\widetilde{\boldsymbol{\Theta}}'$ of $h$ frames, we randomly occlude $H_{\mathrm{occ}}$ consecutive frames of motion where $H_{\mathrm{occ}}$ is uniformly sampled from $[H_{\mathrm{lb}}, H_{\mathrm{ub}}]$. Note that we do not occlude the first $h_{\mathrm{c}}$ frames which are reserved as context. We use the standard CVAE objective to train the motion infiller $\mathcal{M}$:

$$L_{\mathcal{M}} = \sum_{t=1}^{h} \|\widetilde{\boldsymbol{\theta}}_t - \widetilde{\boldsymbol{\theta}}'_t\|_2^2 + L_{\mathrm{KL}}^{\boldsymbol{z}}, \qquad (2)$$

where $L_{\mathrm{KL}}^{\boldsymbol{z}}$ is the KL divergence between the prior and posterior distributions of the CVAE latent code $\boldsymbol{z}$.

## 3.2. Global Trajectory Predictor

After we obtain occlusion-free body motion $\widehat{\boldsymbol{\Theta}}^i$ for each person using the motion infiller, a key problem still remains: the estimated trajectory $(\widetilde{\boldsymbol{T}}^i, \widetilde{\boldsymbol{R}}^i)$ of the person is still occluded and not in a consistent global coordinate system. To tackle this problem, we propose to learn a global trajectory predictor $\mathcal{T}$ that generates a person's occlusion-free global trajectory $(\widehat{\boldsymbol{T}}^i, \widehat{\boldsymbol{R}}^i)$ from the local body motion $\widehat{\boldsymbol{\Theta}}^i$.

Given a general occlusion-free body motion $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_m)$ as input, the trajectory predictor $\mathcal{T}$ outputs its corresponding global trajectory $(\boldsymbol{T}, \boldsymbol{R})$ including the root translations $\boldsymbol{T} = (\boldsymbol{\tau}_1, \ldots, \boldsymbol{\tau}_m)$ and rotations $\boldsymbol{R} = (\boldsymbol{\gamma}_1, \ldots, \boldsymbol{\gamma}_m)$. To address any potential ambiguity in the global trajectory, we also formulate the global trajectory predictor using the CVAE:

$$\boldsymbol{\Psi} = \mathcal{T}(\boldsymbol{\Theta}, \boldsymbol{v}), \qquad (3)$$

$$(\boldsymbol{T}, \boldsymbol{R}) = \texttt{EgoToGlobal}(\boldsymbol{\Psi}), \qquad (4)$$

where the global trajectory predictor $\mathcal{T}$ corresponds to the CVAE decoder and $\boldsymbol{v}$ is the latent code for the CVAE. In Eq. (3), the immediate output of the global trajectory predictor $\mathcal{T}$ is an egocentric trajectory $\boldsymbol{\Psi} = (\boldsymbol{\psi}_1, \ldots, \boldsymbol{\psi}_m)$, which by design can be converted to a global trajectory $(\boldsymbol{T}, \boldsymbol{R})$ using a conversion function $\texttt{EgoToGlobal}$.

**Egocentric Trajectory Representation.** The egocentric trajectory $\boldsymbol{\Psi}$ is just an alternative representation of the global trajectory $(\boldsymbol{T}, \boldsymbol{R})$. It converts the global trajectory into relative local differences and represents rotations and translations in the heading coordinates ($y$-axis aligned with the heading, *i.e.*, the person's facing direction). In this way, the egocentric trajectory representation is invariant of the absolute $xy$ translation and heading. It is more suitable for the prediction of long trajectories, since the network only needs to output the local trajectory change of every frame instead of the potentially large global trajectory offset.

The conversion from the global trajectory to the egocentric trajectory is given by another function: $\boldsymbol{\Psi} = \texttt{GlobalToEgo}(\boldsymbol{T}, \boldsymbol{R})$, which is the inverse of the function $\texttt{EgoToGlobal}$. In particular, the egocentric trajectory $\boldsymbol{\psi}_t = (\delta x_t, \delta y_t, z_t, \delta \phi_t, \boldsymbol{\eta}_t)$ at time $t$ is computed as:

$$(\delta x_t, \delta y_t) = \texttt{ToHeading}(\boldsymbol{\tau}_t^{xy} - \boldsymbol{\tau}_{t-1}^{xy}), \qquad (5)$$

$$z_t = \boldsymbol{\tau}_t^z, \quad \delta \phi_t = \boldsymbol{\gamma}_t^{\phi} - \boldsymbol{\gamma}_{t-1}^{\phi}, \qquad (6)$$

$$\boldsymbol{\eta}_t = \texttt{ToHeading}(\boldsymbol{\gamma}_t), \qquad (7)$$

where $\boldsymbol{\tau}_t^{xy}$ is the $xy$ component of the translation $\boldsymbol{\tau}_t$, $\boldsymbol{\tau}_t^z$ is the $z$ component (height) of $\boldsymbol{\tau}_t$, $\boldsymbol{\gamma}_t^{\phi}$ is the heading angle of the rotation $\boldsymbol{\gamma}_t$, $\texttt{ToHeading}$ is a function that converts translations or rotations to the heading coordinates defined by the heading $\boldsymbol{\gamma}_t^{\phi}$, and $\boldsymbol{\eta}_t$ is the local rotation. As an exception, $(\delta x_0, \delta y_0)$ and $\delta \phi_0$ are used to store the initial $xy$ translation $\boldsymbol{\tau}_0^{xy}$ and heading $\boldsymbol{\tau}_0^{\phi}$. These initial values are set to the GT during training and arbitrary values during inference (as the trajectory can start from any position and heading). The inverse process of Eq. (5)-(7) defines the inverse conversion $\texttt{EgoToGlobal}$ used in Eq. (4), which accumulates the egocentric trajectory to obtain the global trajectory. To correct potential drifts in the trajectory, in Sec. 3.3, we will optimize the global trajectory of each person to match the video evidence, which also solves the trajectory's starting point $(\delta x_0, \delta y_0, \delta \phi_0)$. More details about the egocentric trajectory are given in Appendix D.

**Network and Training.** The trajectory predictor adopts a similar network design as the motion infiller with one main difference: we use LSTMs for temporal modeling instead of Transformers since the output of each frame is the local trajectory change in our egocentric trajectory representation, which mainly depends on the body motion of nearby frames and does not require long-range temporal modeling. We will show in Sec. 4.2 that the egocentric trajectory and use of LSTMs instead of Transformers are crucial for accurate trajectory prediction. Please refer to Appendix D for the detailed network architectures. We use the standard CVAE objective to train the trajectory predictor $\mathcal{T}$:

$$L_{\mathcal{T}} = \sum_{t=1}^{m} \left( \|\boldsymbol{\tau}_t - \boldsymbol{\tau}'_t\|_2^2 + \|\boldsymbol{\gamma}_t \ominus \boldsymbol{\gamma}'_t\|_a^2 \right) + L_{\mathrm{KL}}^{\boldsymbol{v}}, \qquad (8)$$

5

where $\boldsymbol{\tau}'_t$ and $\boldsymbol{\gamma}'_t$ denote the GT translation and rotation, $\ominus$ computes the relative rotation, $\|\cdot\|_a$ computes the rotation angle, and $L^v_{\text{KL}}$ is the KL divergence between the prior and posterior distributions of the CVAE latent code $v$. We again use AMASS [60] to train the trajectory predictor $\mathcal{T}$.

### 3.3. Global Optimization

After using the generative motion infiller and global trajectory predictor, we have obtained an occlusion-free global motion $\widehat{\boldsymbol{Q}}^i = (\widehat{\boldsymbol{T}}^i, \widehat{\boldsymbol{R}}^i, \widehat{\boldsymbol{\Theta}}^i, \widehat{\boldsymbol{B}}^i)$ for each person in the video. However, the global trajectory predictor generates trajectories for each person independently, which may not be consistent with the video evidence. To tackle this problem, we propose a global optimization process that jointly optimizes the global trajectories of all people and the extrinsic camera parameters to match the video evidence such as 2D keypoints. The final output of the global optimization and our framework is $\check{\boldsymbol{Q}}^i = (\check{\boldsymbol{T}}^i, \check{\boldsymbol{R}}^i, \check{\boldsymbol{\Theta}}^i, \check{\boldsymbol{B}}^i)$ where $(\check{\boldsymbol{\Theta}}^i, \check{\boldsymbol{B}}^i) = (\widehat{\boldsymbol{\Theta}}^i, \widehat{\boldsymbol{B}}^i)$, *i.e.*, we directly use the occlusion-free body motion and shapes from the previous stages.

**Optimization Variables.** The first set of variables we optimize is the egocentric representation $\{\check{\boldsymbol{\Psi}}^i\}_{i=1}^N$ of the global trajectories $\{(\check{\boldsymbol{T}}^i, \check{\boldsymbol{R}}^i)\}_{i=1}^N$. We adopt the egocentric representation since it allows corrections of the translation and heading at one frame to propagate to all future frames. We will empirically demonstrate its effectiveness in Sec. 4.2.

The second set of optimization variables is the extrinsic camera parameters $\boldsymbol{C} = (\boldsymbol{C}_1, \dots, \boldsymbol{C}_T)$ where $\boldsymbol{C}_t \in \mathbb{R}^{4\times4}$ is the camera extrinsic matrix at frame $t$ of the video.

**Energy Function.** The energy function we aim to minimize is defined as

$$E(\{\check{\boldsymbol{\Psi}}^i\}_{i=1}^N, \boldsymbol{C}) = \lambda_{\text{2D}} E_{\text{2D}} + \lambda_{\text{traj}} E_{\text{traj}} \\ + \lambda_{\text{reg}} E_{\text{reg}} + \lambda_{\text{cam}} E_{\text{cam}} + \lambda_{\text{pen}} E_{\text{pen}}, \quad (9)$$

where we use five energy terms with their corresponding coefficients $\lambda_{\text{2D}}, \lambda_{\text{traj}}, \lambda_{\text{reg}}, \lambda_{\text{cam}}, \lambda_{\text{pen}}$.

The first term $E_{\text{2D}}$ measures the error between the 2D projection $\check{\boldsymbol{x}}^i_t$ of the optimized 3D keypoints $\widetilde{\boldsymbol{X}}^i_t \in \mathbb{R}^{J\times3}$ and the estimated 2D keypoints $\widetilde{\boldsymbol{x}}^i_t$ from a keypoint detector:

$$E_{\text{2D}} = \frac{1}{NTJ} \sum_{i=1}^N \sum_{t=1}^T V^i_t \|\check{\boldsymbol{x}}^i_t - \widetilde{\boldsymbol{x}}^i_t\|_F^2, \quad (10)$$

$$\check{\boldsymbol{x}}^i_t = \Pi\left(\widetilde{\boldsymbol{X}}^i_t, \boldsymbol{C}_t, \boldsymbol{K}\right), \quad \widetilde{\boldsymbol{X}}^i_t = \mathcal{J}(\check{\boldsymbol{\tau}}^i_t, \check{\boldsymbol{\gamma}}^i_t, \check{\boldsymbol{\theta}}^i_t, \check{\boldsymbol{\beta}}^i_t) \quad (11)$$

where $V^i_t$ is person $i$'s visibility at frame $t$, $\Pi$ is the camera projection with extrinsics $\boldsymbol{C}_t$ and approximated intrinsics $\boldsymbol{K}$, and $\widetilde{\boldsymbol{X}}^i_t$ is computed using the SMPL joint function $\mathcal{J}$ from the optimized global pose $\check{\boldsymbol{q}}^i_t = (\check{\boldsymbol{\tau}}^i_t, \check{\boldsymbol{\gamma}}^i_t, \check{\boldsymbol{\theta}}^i_t, \check{\boldsymbol{\beta}}^i_t) \in \check{\boldsymbol{Q}}^i$.

The second term $E_{\text{traj}}$ measures the difference between the optimized global trajectory $(\check{\boldsymbol{T}}^i, \check{\boldsymbol{R}}^i)$ viewed in the camera coordinates and the trajectory $(\widetilde{\boldsymbol{T}}^i, \widetilde{\boldsymbol{R}}^i)$ output by the pose estimator (*e.g.*, KAMA [33]) in Stage I:

$$E_{\text{traj}} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T V^i_t \left(\|\Gamma(\check{\boldsymbol{\gamma}}^i_t, \boldsymbol{C}_t) \ominus \widetilde{\boldsymbol{\gamma}}^i_t\|_a^2 \right. \\ \left. + w_t \|\Gamma(\check{\boldsymbol{\tau}}^i_t, \boldsymbol{C}_t) - \widetilde{\boldsymbol{\tau}}^i_t\|_2^2\right), \quad (12)$$

where the function $\Gamma(\cdot, \boldsymbol{C}_t)$ transforms the global rotation $\check{\boldsymbol{\gamma}}^i_t$ or translation $\check{\boldsymbol{\tau}}^i_t$ to the camera coordinates defined by $\boldsymbol{C}_t$, and $w_t$ is a weighting factor for the translation term.

The third term $E_{\text{reg}}$ regularizes the egocentric trajectory $\check{\boldsymbol{\Psi}}^i$ to stay close to the output $\widehat{\boldsymbol{\Psi}}^i$ of the trajectory predictor:

$$E_{\text{reg}} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left\|\boldsymbol{w}_\psi \circ \left(\check{\boldsymbol{\psi}}^i_t - \widehat{\boldsymbol{\psi}}^i_t\right)\right\|_2^2, \quad (13)$$

where $\circ$ denotes the element-wise product and $\boldsymbol{w}_\psi$ is a weighting vector for each element inside the egocentric trajectory. As an exception, we do not regularize each person's initial $xy$ position and heading $(\delta\check{x}^i_0, \delta\check{y}^i_0, \delta\check{\phi}^i_0) \subset \check{\boldsymbol{\psi}}^i_0$ as they need to be inferred from the video.

The fourth term $E_{\text{cam}}$ measures the smoothness of the camera parameters $\boldsymbol{C}$ and the uprightness of the camera:

$$E_{\text{cam}} = \frac{1}{T} \sum_{t=1}^T \langle \boldsymbol{C}^y_t, \boldsymbol{Y} \rangle \\ + \frac{1}{T-1} \sum_{t=1}^{T-1} \left\|\boldsymbol{C}^\gamma_{t+1} \ominus \boldsymbol{C}^\gamma_t\right\|_a^2 + \left\|\boldsymbol{C}^\tau_{t+1} - \boldsymbol{C}^\tau_t\right\|_2^2, \quad (14)$$

where $\langle\cdot,\cdot\rangle$ denotes the inner product, $\boldsymbol{C}^y_t$ is the $+y$ vector of the camera $\boldsymbol{C}_t$, and $\boldsymbol{Y}$ is the global up direction. $\boldsymbol{C}^\gamma_t$ and $\boldsymbol{C}^\tau_t$ denote the rotation and translation of the camera $\boldsymbol{C}_t$.

The final term $E_{\text{pen}}$ is an signed distance field (SDF)-based inter-person penetration loss adopted from [36].

## 4. Experiments

**Datasets.** We employ the following datasets in our experiments: (1) **AMASS** [60], which is a large human motion database with 11000+ human motions. We use AMASS to train and evaluate the motion infiller and trajectory predictor. (2) **3DPW** [94], which is an *in-the-wild* human motion dataset that uses videos and wearable IMU sensors to obtain GT poses, even when the person is occluded. We evaluate our approach using the test split of 3DPW. (3) **Dynamic Human3.6M** is a *new* benchmark for human pose estimation with dynamic cameras that we create from the Human3.6M dataset [31]. We simulate dynamic cameras and occlusions by cropping each frame with a small

view window that oscillates around the person (see Fig. 5). More details are provided in Appendix A and we will release the code for generating the dataset.

**Evaluation Metrics.** We use the following metrics for evaluation: (1) **G-MPJPE** and **G-PVE**, which extend the mean per joint position error (MPJPE) and per-vertex error (PVE) by computing the errors in the global coordinates. As errors in estimated global trajectories accumulate over time in our dynamic camera setting, we follow standard evaluations for open-loop reconstruction (*e.g.*, SLAM [87] and inertial odometry [27]) to compute errors using a sliding window (10 seconds) and align the root translation and rotation with the GT at the start of the window. (2) **PA-MPJPE**, which is the Procrustes-aligned MPJPE for evaluating estimated body poses. For invisible poses, since there can be many plausible poses beside the GT, we follow prior work [3, 102] to compute the best PA-MPJPE out of multiple samples for our probabilistic approach. (3) **Accel**, which computes the mean acceleration error of each joint and is commonly used to measure the jitter in estimated motions [45, 104]. (4) **FID**, which is an extension of the original Frechet Inception Distance that calculates the distribution distance between estimated motions and the GT. FID is a standard metric in motion generation literature to evaluate the quality of generated motions [30, 54, 55, 91]. Following prior work [55], we compute FID using the well-designed kinetic motion feature extractor in the fairmotion library [20].

**Implementation Details.** Thorough details about the entire framework are provided in Appendix A to E.

## 4.1. Evaluation of GLAMR

**Baselines.** Since no prior methods can estimate global motions from dynamic cameras and address long-term occlusions, we design various baselines by combining state-of-the-art human mesh recovery methods (KAMA [33] or SPEC [47]), motion infilling methods, and SLAM-based camera estimation (OpenSfM [68]). In particular, we use the estimated camera parameters to convert estimated motions from the camera coordinates to the global coordinates. For motion infilling, we use (1) linear interpolation, (2) last pose, *i.e.*, replicating the last visible pose, and (3) a state-of-the-art CNN-based motion infilling method, ConvAE [41].

The results on Dynamic Human3.6M and 3DPW are summarized in Table 1 and 2 respectively. We only report G-MPJPE and G-PVE on Dynamic Human3.6M since they require accurate GT trajectories, which 3DPW does not provide. It is evident that our approach, GLAMR, outperforms the baselines in almost all metrics. In particular, GLAMR achieves significantly lower G-MPJPE and G-PVE, which demonstrates its strong ability to reconstruct global human motions. Furthermore, GLAMR attains considerably lower FID and PA-MPJPE (with ten samples) for occluded (invis-

| Method | (All) G-MPJPE | (All) G-PVE | (Invisible) FID | (Invisible) PA-MPJPE | (Visible) PA-MPJPE | (All) Accel |
|---|---|---|---|---|---|---|
| KAMA [41] + Linear Interpolation | 1735.2 | 1744.1 | 30.2 | 74.8 | **47.4** | 8.0 |
| KAMA [41] + Last Pose | 1318.1 | 1330.3 | 36.7 | 88.8 | **47.4** | 12.3 |
| KAMA [41] + ConvAE [41] | 1737.8 | 1748.9 | 28.9 | 77.4 | 56.9 | 7.5 |
| SPEC [47] + Linear Interpolation | 2113.3 | 2119.5 | 29.7 | 78.7 | 55.7 | 14.2 |
| SPEC [47] + Last Pose | 1782.5 | 1790.9 | 36.2 | 92.6 | 55.7 | 18.8 |
| SPEC [47] + ConvAE [41] | 2113.3 | 2119.0 | 28.5 | 80.1 | 59.9 | 11.9 |
| Ours (GLAMR w/ SPEC) | 899.1 | 913.7 | **8.2** | 72.8 | 55.0 | 6.6 |
| Ours (GLAMR w/ KAMA) | **806.2** | **824.1** | 11.4 | **67.7** | 47.6 | **6.0** |

Table 1. Baseline comparison on Dynamic Human3.6M. We report results for visible, invisible (occluded), and all frames.

| Method | (Invisible) FID | (Invisible) PA-MPJPE | (Visible) PA-MPJPE | (All) Accel |
|---|---|---|---|---|
| KAMA [41] + Linear Interpolation | 30.7 | 87.5 | **50.8** | 24.2 |
| KAMA [41] + Last Pose | 40.3 | 96.3 | **50.8** | 25.4 |
| KAMA [41] + ConvAE [41] | 32.0 | 84.5 | 56.4 | 19.6 |
| SPEC [47] + Linear Interpolation | 33.6 | 85.6 | 53.3 | 33.1 |
| SPEC [47] + Last Pose | 39.5 | 92.4 | 53.3 | 34.2 |
| SPEC [47] + ConvAE [41] | 35.4 | 86.9 | 59.3 | 24.0 |
| Ours (GLAMR w/ SPEC) | 24.8 | 79.1 | 54.9 | 9.5 |
| Ours (GLAMR w/ KAMA) | **22.6** | **73.6** | 51.1 | **8.9** |

Table 2. Baseline comparison on 3DPW. G-MPJPE and G-PVE are not reported since 3DPW does not provide accurate GT global human trajectories. See also the caption of Table 1.

ible) poses. The lower FID means GLAMR can infill more humanlike motions, and the lower PA-MPJPE also shows GLAMR's probabilistic motion samples can cover the GT better. Finally, while GLAMR achieves almost the same PA-MPJPE for visible poses as the best method, it yields much smoother motions (smaller acceleration error). This is because our motion infiller leverages human dynamics learned from large motion dataset to produce motions.

**Qualitative Results.** Fig. 4 and 5 show qualitative comparisons of GLAMR against the strong baseline, KAMA + Linear Interpolation. Additionally, we provide abundant qualitative results in our supplementary video.

## 4.2. Evaluation of Key Components

**Benchmarking Motion Infiller.** We evaluate the proposed generative motion infiller on the test split of the AMASS dataset [60]. We compare against three motion infilling baselines: linear interpolation, replicating the last pose, and ConvAE [41]. As shown in Table 3, our generative motion infiller achieves significantly better PA-MPJPE for both the sampled motions (with five samples) and reconstructed motion for the infilled frames. Our approach also achieves considerably better FID, reducing the FID of ConvAE [41] by half, which indicates that the infilled motions by our approach are much closer to real human motions.

**Benchmarking Trajectory Predictor.** We also evaluate our global trajectory predictor against two variants on the test split of AMASS: (1) "Transformer", which replaces the LSTMs in the trajectory predictor with Transformers; (2) "Ours w/o Ego Trajectory", which directly outputs the 6-DoF global trajectory instead of the egocentric trajectory.
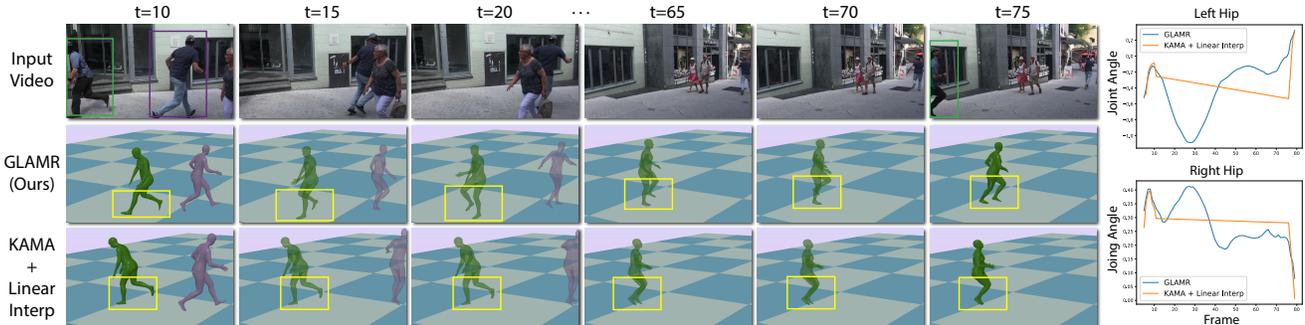
Figure 4. Qualitative comparison of GLAMR with a strong baseline on 3DPW. The infilled motion (transparent) by GLAMR is more natural especially for the legs, while the baseline has very slow leg motions due to interpolation in a large window (frame 10 to 75). On the **right**, we plot how the $x$-axis joint angles of left and right hips of the person (green) change over time for GLAMR and the baseline.
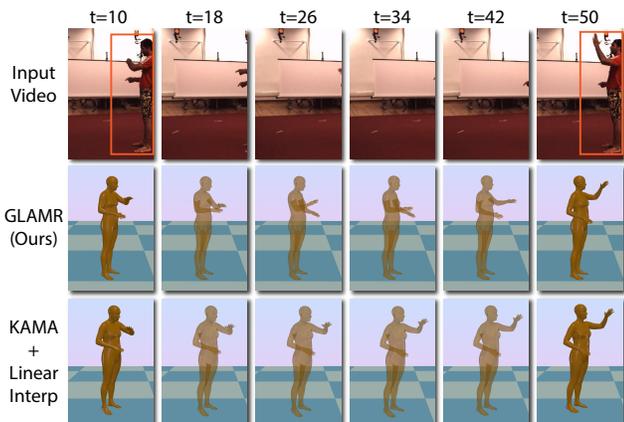


Figure 5. Qualitative comparison of GLAMR on Dynamic Human3.6M. GLAMR can generate natural hand motions for invisible frames instead of just doing linear interpolation.

| Method | (Sampled) PA-MPJPE | (Reconstructed) PA-MPJPE | (Sampled) FID |
|---|---|---|---|
| Linear Interpolation | 83.5 | 83.5 | 35.3 |
| Last Pose | 104.4 | 104.4 | 41.6 |
| ConvAE [41] | 72.8 | 72.8 | 31.4 |
| Ours | **61.4** | **36.1** | **16.7** |

Table 3. Benchmarking motion infiller on AMASS.

| Method | G-MPJPE | G-PVE | Accel |
|---|---|---|---|
| Transformer | 660.1 | 678.6 | 121.9 |
| Ours w/o Ego Trajectory | 763.0 | 780.6 | 8.7 |
| Ours | **466.9** | **472.5** | **5.8** |

Table 4. Benchmarking trajectory predictor on AMASS.

| Method | G-MPJPE | G-PVE | Accel |
|---|---|---|---|
| Ours w/o Trajectory Predictor | 1750.8 | 1761.4 | 12.6 |
| Ours w/o Opt Ego Trajectory | 877.3 | 895.0 | 15.5 |
| Ours (GLAMR) | **806.2** | **824.1** | **6.0** |

Table 5. Global optimization ablations on Dynamic Human3.6M.

Table 4 summarizes the results where the G-MPJPE and G-PVE are computed using the best of five trajectory samples. We can see that both variants lead to worse global trajectory prediction. We believe this is because the positional (time) encoding in Transformers may not generalize to longer motions in the test data, and directly predicting global trajectory offsets from local body motions is also difficult.

**Ablations for Global Optimization.** We further perform ablation studies on the effect of key components in our global optimization. Specifically, we design two variants: (1) "Ours w/o Trajectory Predictor", which does not use our trajectory predictor to generate the global human trajectories and uses camera parameters from OpenSfM [68] to obtain global trajectories instead; (2) "Ours w/o Opt Ego Trajectory", which does not employ the egocentric trajectory representation and directly optimizes the 6-DoF root trajectory instead. As shown in Table 5, both variants lead to significantly worse global trajectory reconstruction with large increases in G-MPJPE, G-PVE, and Accel. This demonstrates that both the global trajectory predictor and egocentric trajectory representation are vital in our approach.

## 5. Discussion and Limitations

In this paper, we proposed an approach to recover 3D human meshes in consistent global coordinates from videos captured by a dynamic camera. To achieve this, we first proposed a novel Transformer-based generative motion infiller to address severe occlusions that often come with dynamic cameras. To resolve ambiguity in the joint reconstruction of global human motions and camera poses, we proposed a new solution by predicting global human trajectories from their local body motions. Finally, we proposed a global optimization framework to refine the predicted trajectories and use them as anchors for camera optimization. Results on challenging datasets with dynamic cameras demonstrated the effectiveness of our approach, which marks a significant step towards global human mesh recovery in the wild.

Our approach is not without limitations. Currently, the generative motion infiller and global trajectory predictor operate for each person independently. Therefore, the generated motions and trajectories may not capture potentially

complex and nuanced interactions between occluded people such as hugging or dancing. Future work could address this limitation by employing new generative models that produce interaction-aware motions of multiple people.

## References

[1] Ijaz Akhter and Michael J. Black. Pose-conditioned joint angle limits for 3D human pose reconstruction. In *CVPR*, 2015. 2

[2] Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. Structured prediction helps 3d human motion modelling. In *ICCV*, 2019. 3

[3] Sadegh Aliakbarian, Fatemeh Sadat Saleh, Mathieu Salzmann, Lars Petersson, and Stephen Gould. A stochastic conditioning scheme for diverse human motion prediction. In *CVPR*, 2020. 7

[4] Emad Barsoum, John Kender, and Zicheng Liu. Hp-gan: Probabilistic 3d human motion prediction via gan. In *CVPR Workshops*, 2018. 3

[5] Vasileios Belagiannis, Sikandar Amin, Mykhaylo Andriluka, Bernt Schiele, Nassir Navab, and Slobodan Ilic. 3d pictorial structures for multiple human pose estimation. In *CVPR*, 2014. 2

[6] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *ECCV*, 2016. 2

[7] Tianshu Zhang Buzhen Huang, Yuan Shu and Yangang Wang. Dynamic multi-person mesh recovery from uncalibrated multi-view cameras. In *3DV*, 2021. 2

[8] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose. In *ECCV*, 2020. 2

[9] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Beyond static features for temporally consistent 3d human pose and shape from a video. In *CVPR*, 2021. 2

[10] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. Monocular expressive body regression through body-driven attention. In *ECCV*, 2020. 2

[11] MMTracking Contributors. MMTracking: OpenMMLab video perception toolbox and benchmark. https://github.com/open-mmlab/mmtracking, 2020. 13

[12] Rishabh Dabral, Soshi Shimada, Arjun Jain, Christian Theobalt, and Vladislav Golyanik. Gravity-aware monocular 3d human-object reconstruction. In *ICCV*, 2021. 2

[13] Junting Dong, Qi Fang, Wen Jiang, Yurou Yang, Hujun Bao, and Xiaowei Zhou. Fast and robust multi-person 3d pose estimation and tracking from multiple views. *TPAMI*, 2021. 2

[14] Junting Dong, Qing Shuai, Yuanqing Zhang, Xian Liu, Xiaowei Zhou, and Hujun Bao. Motion capture from internet videos. In *ECCV*, 2020. 2

[15] Zijian Dong, Jie Song, Xu Chen, Chen Guo, and Otmar Hilliges. Shape-aware multi-person pose estimation from multi-view images. In *ICCV*, 2021. 2

[16] Matteo Fabbri, Fabio Lanzi, Simone Calderara, Stefano Alletto, and Rita Cucchiara. Compressed volumetric heatmaps for multi-person 3d pose estimation. In *CVPR*, June 2020. 2

[17] Mihai Fieraru, Mihai Zanfir, Elisabeta Oneata, Alin-Ionut Popa, Vlad Olaru, and Cristian Sminchisescu. Three-dimensional reconstruction of human interactions. In *CVPR*, 2020. 2, 3

[18] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *ICCV*, 2015. 3

[19] Anand Gopalakrishnan, Ankur Mali, Dan Kifer, Lee Giles, and Alexander G Ororbia. A neural temporal model for human motion prediction. In *CVPR*, 2019. 3

[20] Deepak Gopinath and Jungdam Won. fairmotion - tools to load, process and visualize motion capture data. https://github.com/facebookresearch/fairmotion, 2020. 7

[21] Rıza Alp Guler and Iasonas Kokkinos. HoloPose: Holistic 3d human reconstruction in-the-wild. In *CVPR*, 2019. 2

[22] Vladimir Guzov, Aymen Mir, Torsten Sattler, and Gerard Pons-Moll. Human poseitioning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In *CVPR*, 2021. 3

[23] Félix G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)*, 39(4):60–1, 2020. 3

[24] Nils Hasler, Bodo Rosenhahn, Thorsten Thormahlen, Michael Wand, Jürgen Gall, and Hans-Peter Seidel. Markerless motion capture with unsynchronized moving cameras. In *CVPR*, 2009. 2

[25] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In *ICCV*, pages 11374–11384, 2021. 3

[26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 13

[27] Sachini Herath, Hang Yan, and Yasutaka Furukawa. Ronin: Robust neural inertial navigation in the wild: Benchmark, evaluations, & new methods. In *ICRA*, 2020. 7

[28] Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. Human motion prediction via spatio-temporal inpainting. In *CVPR*, 2019. 3, 5

[29] Buzhen Huang, Yuan Shu, Tianshu Zhang, and Yangang Wang. Dynamic multi-person mesh recovery from uncalibrated multi-view cameras. In *3DV*, 2021. 2

[30] Ruozi Huang, Huang Hu, Wei Wu, Kei Sawada, Mi Zhang, and Daxin Jiang. Dance revolution: Long-term dance generation with music via curriculum learning. In *ICLR*, 2021. 7

[31] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *TPAMI*, 36(7):1325–1339, 2014. 6, 13

[32] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning via multi-view images in the wild. In *CVPR*, 2020. 2

[33] Umar Iqbal, Kevin Xie, Yunrong Guo, Jan Kautz, and Pavlo Molchanov. Kama: 3d keypoint aware body mesh articulation. In *3DV*, 2021. 2, 4, 6, 7, 13, 15

[34] Mariko Isogawa, Ye Yuan, Matthew O'Toole, and Kris M Kitani. Optical non-line-of-sight physics-based 3d human pose estimation. In *CVPR*, 2020. 2

[35] Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *CVPR*, 2016. 3

[36] Wen Jiang, Nikos Kolotouros, Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Coherent reconstruction of multiple humans from a single image. In *CVPR*, 2020. 2, 6

[37] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human pose fitting towards in-the-wild 3d human pose estimation. In *3DV*, 2021. 2

[38] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *CVPR*, 2018. 2

[39] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 2, 13

[40] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 2

[41] Manuel Kaufmann, Emre Aksan, Jie Song, Fabrizio Pece, Remo Ziegler, and Otmar Hilliges. Convolutional autoencoders for human motion infilling. In *3DV*, 2020. 3, 5, 7, 8

[42] Tarasha Khurana, Achal Dave, and Deva Ramanan. Detecting invisible people. In *ICCV*, pages 3174–3184, 2021. 3

[43] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *arXiv preprint arXiv:1412.6980*, 2014. 14, 15

[44] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 4

[45] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 2, 3, 7

[46] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. PARE: Part attention regressor for 3D human body estimation. In *ICCV*, 2021. 3

[47] Muhammed Kocabas, Chun-Hao P. Huang, Joachim Tesch, Lea Müller, Otmar Hilliges, and Michael J. Black. SPEC: Seeing people in the wild with an estimated camera. In *ICCV*, 2021. 2, 4, 7, 13

[48] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 2

[49] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *CVPR*, 2019. 2

[50] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *ICCV*, 2021. 2

[51] Jogendra Nath Kundu, Mugalodi Rakesh, Varun Jampani, Rahul Mysore Venkatesh, and R. Venkatesh Babu1. Appearance consensus driven self-supervised human mesh recovery. In *ECCV*, 2020. 2

[52] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J. Black, and Peter V. Gehler. Unite the people: Closing the loop between 3D and 2D human representations. In *CVPR*, 2017. 2

[53] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *CVPR*, 2021. 2

[54] Jiaman Li, Yihang Yin, Hang Chu, Yi Zhou, Tingwu Wang, Sanja Fidler, and Hao Li. Learning to generate diverse dance motions with transformer. *arXiv preprint arXiv:2008.08171*, 2020. 7

[55] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *ICCV*, 2021. 7

[56] Zimo Li, Yi Zhou, Shuangjiu Xiao, Chong He, Zeng Huang, and Hao Li. Auto-conditioned recurrent networks for extended complex human motion synthesis. *arXiv preprint arXiv:1707.05363*, 2017. 3

[57] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *CVPR*, 2021. 2

[58] Miao Liu, Dexin Yang, Yan Zhang, Zhaopeng Cui, James M Rehg, and Siyu Tang. 4d human body capture from egocentric video via 3d scene grounding. In *3DV*, 2021. 2, 3

[59] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *SIGGRAPH Asia*, 34(6):248:1–248:16, 2015. 3

[60] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. In *ICCV*, 2019. 2, 5, 6, 7, 13

[61] Julieta Martinez, Michael J Black, and Javier Romero. On human motion prediction using recurrent neural networks. In *CVPR*, 2017. 3

[62] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3DV*, 2017. 2

[63] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Mohamed Elgharib, Pascal Fua, Hans-Peter Seidel, Helge Rhodin, Gerard Pons-Moll, and Christian Theobalt. XNect: Real-time multi-person 3D motion capture with a single RGB camera. In *SIGGRAPH*, 2020. 2

[64] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. VNect: Real-time 3D human pose estimation with a single RGB camera. In *SIGGRAPH*, 2017. 2

[65] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In *ICCV*, 2019. 2

[66] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image-to-lixel prediction network for accurate 3d human

pose and mesh estimation from a single rgb image. In *ECCV*, 2020. 2

[67] Lea Müller, Ahmed A. A. Osman, Siyu Tang, Chun-Hao P. Huang, and Michael J. Black. On self contact and human pose. In *CVPR*, 2021. 2

[68] Opensfm - a structure from motion library. https://github.com/mapillary/OpenSfM, 2021. 7, 8

[69] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 14, 15

[70] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, 2019. 2

[71] Georgios Pavlakos, Nikos Kolotouros, and Kostas Daniilidis. Texturepose: Supervising human mesh estimation with texture consistency. In *ICCV*, 2019. 2

[72] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3D human pose and shape from a single color image. In *CVPR*, 2018. 2

[73] Dario Pavllo, David Grangier, and Michael Auli. Quaternet: A quaternion-based recurrent model for human motion. In *BMVC*, 2018. 3

[74] Christian Payer, Thomas Neff, Horst Bischof, Martin Urschler, and Darko Stern. Simultaneous multi-person detection and single-person pose estimation with a single heatmap regression network. In *ICCV PoseTrack Workshop*, 2017. 2

[75] Mathis Petrovich, Michael J Black, and Gül Varol. Action-conditioned 3d human motion synthesis with transformer vae. *arXiv preprint arXiv:2104.05670*, 2021. 3

[76] Haibo Qiu, Chunyu Wang, Jingdong Wang, Naiyan Wang, and Wenjun Zeng. Cross view fusion for 3d human pose estimation. In *ICCV*, 2019. 2

[77] N. Dinesh Reddy, Laurent Guigues, Leonid Pischulini, Jayan Eledath, and Srinivasa Narasimhan. Tessetrack: End-to-end learnable multi-person articulated 3d pose tracking. In *CVPR*, 2021. 2

[78] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. In *ICCV*, 2021. 2, 3

[79] Chris Rockwell and David F. Fouhey. Full-body awareness from partial observations. In *ECCV*, 2020. 3

[80] Gregory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. LCR-Net: Localization-classification-regression for human pose. In *CVPR*, 2017. 2

[81] Yu Rong, Ziwei Liu, Cheng Li, Kaidi Cao, and Chen Change Loy. Delving deep into hybrid annotations for 3d human recovery in the wild. In *ICCV*, 2019. 2

[82] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 3

[83] Paul Schreiner, Maksym Perepichka, Hayden Lewis, Sune Darkner, Paul G Kry, Kenny Erleben, and Victor B Zordan. Global position prediction for interactive motion capture.

*Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 4(3):1–16, 2021. 3

[84] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. In *SIGGRAPH*, 2020. 2

[85] Leonid Sigal and Michael J Black. Humaneva: Synchronized video and motion capture dataset for evaluation of articulated human motion. *Brown Univertsity TR*, 120(2), 2006. 13

[86] Jie Song, Xu Chen, and Otmar Hilliges. Human body model fitting by learned gradient descent. In *ECCV*, 2020. 2

[87] Jürgen Sturm, Stéphane Magnenat, Nikolas Engelhard, François Pomerleau, Francis Colas, Daniel Cremers, Roland Siegwart, and Wolfram Burgard. Towards a benchmark for rgb-d slam evaluation. In *Rgb-d workshop on advanced reasoning with depth cameras at robotics: Science and systems conf.(rss)*, 2011. 7

[88] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, pages 5693–5703, 2019. 13

[89] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Michael J. Black, and Tao Mei. Monocular, one-stage, regression of multiple 3d people. In *ICCV*, 2021. 2

[90] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, Yili Fu, , and Tao Mei. Human mesh recovery from monocular images via a skeleton-disentangled representation. In *ICCV*, 2019. 2

[91] Guillermo Valle-Pérez, Gustav Eje Henter, Jonas Beskow, André Holzapfel, Pierre-Yves Oudeyer, and Simon Alexanderson. Transflower: probabilistic autoregressive dance generation with multimodal attention. *arXiv preprint arXiv:2106.13871*, 2021. 7

[92] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5, 14

[93] Ruben Villegas, Jimei Yang, Yuliang Zou, Sungryull Sohn, Xunyu Lin, and Honglak Lee. Learning to generate long-term future via hierarchical prediction. In *ICML*, 2017. 3

[94] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 3, 6, 13

[95] Zhenzhen Weng and Serena Yeung. Holistic 3d human and scene mesh estimation from single view images. In *CVPR*, 2021. 2

[96] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *ICIP*, pages 3645–3649. IEEE, 2017. 13

[97] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body and hands in the wild. In *CVPR*, 2019. 2

[98] Kevin Xie, Tingwu Wang, Umar Iqbal, Yunrong Guo, Sanja Fidler, and Florian Shkurti. Physics-based human motion estimation and synthesis from videos. In *ICCV*, 2021. 2

[99] Yuanlu Xu, Song-Chun Zhu, and Tony Tung. Denserac: Joint 3d pose and shape estimation by dense render-and-compare. In *ICCV*, 2019. 2

[100] Xinchen Yan, Akash Rastogi, Ruben Villegas, Kalyan Sunkavalli, Eli Shechtman, Sunil Hadap, Ersin Yumer, and Honglak Lee. Mt-vae: Learning motion transformations to generate multimodal human dynamics. In *ECCV*, 2018. 3

[101] Ye Yuan and Kris Kitani. Diverse trajectory forecasting with determinantal point processes. *arXiv preprint arXiv:1907.04967*, 2019. 3

[102] Ye Yuan and Kris Kitani. Dlow: Diversifying latent flows for diverse human motion prediction. In *ECCV*, 2020. 3, 7

[103] Ye Yuan and Kris Kitani. Residual force control for agile human behavior imitation and extended motion synthesis. In *NeurIPS*, 2020. 3

[104] Ye Yuan, Shih-En Wei, Tomas Simon, Kris Kitani, and Jason Saragih. Simpoe: Simulated character control for 3d human pose estimation. In *CVPR*, 2021. 2, 7

[105] Andrei Zanfir, Eduard Gabriel Bazavan, Hongyi Xu, William T. Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Weakly supervised 3d human pose and shape reconstruction with normalizing flows. In *ECCV*, 2020. 2

[106] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3d pose and shape estimation of multiple people in natural scenes the importance of multiple scene constraints. In *CVPR*, 2018. 2

[107] Andrei Zanfir, Elisabeta Marinoiu, Mihai Zanfir, Alin-Ionut Popa, and Cristian Sminchisescu. Deep network for the integrated 3d sensing of multiple people in natural images. In *NeurIPS*, 2018. 2

[108] Mihai Zanfir, Andrei Zanfir, Eduard Gabriel Bazavan, William T. Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Thundr: Transformer-based 3d human reconstruction with markers. In *ICCV*, 2021. 2

[109] Hongwen Zhang, Yating Tian, Xinchi Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. Pymaf: 3d human pose and shape regression with pyramidal mesh alignment feedback loop. In *ICCV*, 2021. 2

[110] Jianfeng Zhang, Dongdong Yu, Jun Hao Liew, Xuecheng Nie, and Jiashi Feng. Body meshes as points. In *CVPR*, 2021. 2

[111] Siwei Zhang, Yan Zhang, Federica Bogo, Marc Pollefeys, and Siyu Tang. Learning motion priors for 4d human body capture in 3d scenes. In *ICCV*, 2021. 3

[112] Tianshu Zhang, Buzhen Huang, and Yangang Wang. Object-occluded human shape and pose estimation from a single color image. In *CVPR*, 2020. 2, 3

[113] Yuxiang Zhang, Liang An, Tao Yu, Xiu Li, Kun Li, and Yebin Liu. 4D association graph for realtime multi-person motion capture using multiple video cameras. In *CVPR*, 2020. 2

[114] Yuxiang Zhang, Zhe Li, Liang An, Mengcheng Li, Tao Yu, and Yebin Liu. Lightweight multi-person total motion capture using sparse multi-view cameras. In *ICCV*, 2021. 2

[115] Jianan Zhen, Qi Fang, Jiaming Sun, Wentao Liu, Wei Jiang, Hujun Bao, and Xiaowei Zhou. SMAP: Single-shot multi-person absolute 3d pose estimation. In *ECCV*, 2020. 2

[116] Yang Zheng, Ruizhi Shao, Yuxiang Zhang, Tao Yu, Zerong Zheng, Qionghai Dai, and Yebin Liu. Deepmulticap: Performance capture of multiple characters using sparse multiview cameras. In *ICCV*, 2021. 2

[117] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 14

[118] Yuxiao Zhou, Marc Habermann, Ikhsanul Habibie, Ayush Tewari, Christian Theobalt, and Feng Xu. Monocular real-time full body capture with inter-part correlations. In *CVPR*, 2021. 2

## A. Details for the Datasets

**AMASS** [60] is a large human motion database with 11000+ human motions. We use AMASS to train and evaluate the motion infiller and trajectory predictor. Specifically, we use the Transitions, SSM, and HumanEva [85] subsets for testing and all other subsets for training.

**3DPW** [94] is an in-the-wild human motion dataset that consists of 60 videos recorded with dynamic cameras in diverse environments. The GT 3D poses are obtained using wearable IMU sensors. Since non-optical sensors are used to obtain GT data, the dataset also provides body pose information when the persons go outside the FoV of the camera. 3DPW also provides the global trajectories of people in the dataset. However, the global trajectories are quite inaccurate since they are estimated from IMU data. Therefore, we do not use 3DPW to evaluate global trajectory reconstruction in the paper. Since we do not use 3DPW for training, we use sequences from the entire 3DPW dataset for visualization. We use the official 3DPW test split to report quantitative results in the paper.

**Dynamic Human3.6M** is a new benchmark for global human pose estimation with dynamic cameras that we create from the Human3.6M dataset [31]. We simulate dynamic cameras and occlusions by cropping each frame with a view window of $300 \times 600$ that horizontally oscillates around the person's bounding box center with a period of 4.8 seconds and a magnitude of 200 pixels. In this way, we synthesize large camera motions and severe occlusions where the person is occluded for almost half of the time, which makes it very challenging for existing 3D human pose and shape estimation methods. Additionally, since Human3.6M provides accurate global human trajectories and human poses, we use Dynamic Human3.6M to evaluate global trajectory reconstruction and pose estimation for occluded frames. We follow the standard protocol [39] and use the official test split (subjects 9 and 11) for evaluation. Please refer to the supplementary video for an example sequence of the Dynamic Human3.6M dataset. We will release the code for generating Dynamic Human3.6M for users who have downloaded the original Human3.6M dataset [31].

## B. Implementation Details for Preprocessing

**3D Multi-Object Tracking and Re-identification.** We use DeepSORT [96] with ResNet-50 [26] in the MMTracking package [11] for 3D multi-object tracking (MOT) and re-identification. We use the GT tracks to evaluate our approach and the baselines, following the standard protocol for human pose estimation.

**Initial Human Pose and Shape Estimation.** As mentioned in the main paper, we use KAMA [33] or SPEC [47] to provide the initial human pose and shape estimation from the bounding boxes extracted by 3D MOT. We choose these two methods since both KAMA and SPEC estimate 3D human poses in the camera coordinates with absolute root translations, while many state-of-the-art human pose estimation methods do not provide the root translations. We also use HRNet [88] to extract 2D human keypoints from the video, which are used in the proposed global optimization framework.

## C. Implementation Details for Generative Motion Infiller

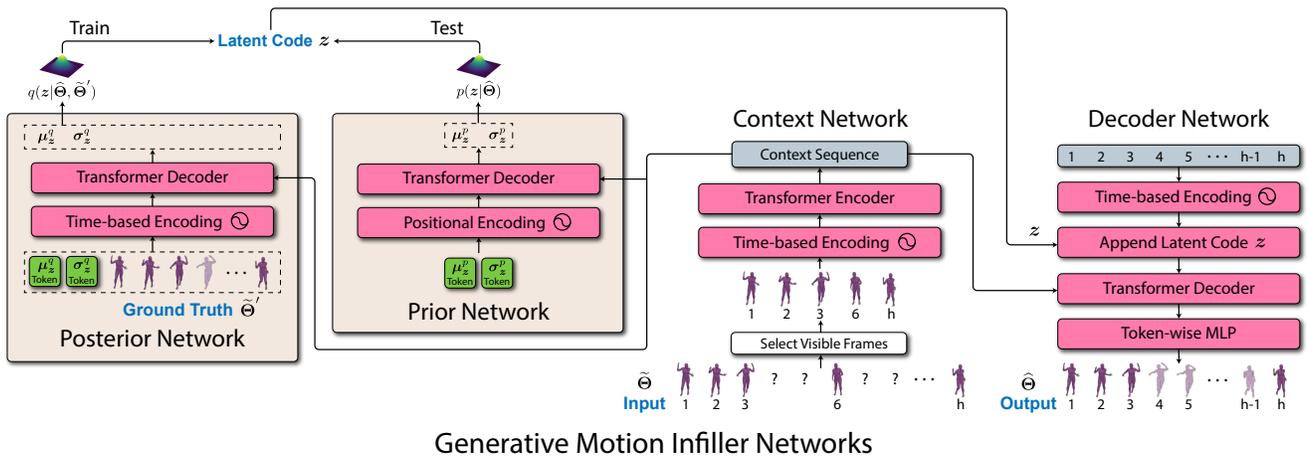

Generative Motion Infiller Networks

Figure 6. The detailed network architecture of the CVAE-based generative motion infiller. For all the Transformer modules, the dimensions for keys, queries, and values are set to 256, the number of transformer blocks is 2, the hidden dimensions of the feedforwards layers are 512, the dropout rate is 0.1, and 8 heads are used for the multi-head attention. Two hidden layers (512, 256) with ReLU activations are used for all the token-wise MLPs.

**Network Architecture.** The detailed network architecture of the CVAE-based generative motion infiller is outlined in Fig. 6. For all the Transformer [92] modules, the dimensions for keys, queries, and values are set to 256, the number of transformer blocks is 2, the hidden dimensions of the feedforwards layers are 512, the dropout rate is 0.1, and 8 heads are used for the multi-head attention. The time-based encoding takes the same sinusoidal form as the original positional encoding [92] but replaces the position with the time index. We use two hidden layers (512, 256) with ReLU activations for all the token-wise MLPs. In the prior network, two learnable tokens are used to form queries to produce the mean $\boldsymbol{\mu}_{\boldsymbol{z}}^p$ and standard deviation $\boldsymbol{\sigma}_{\boldsymbol{z}}^p$ of the prior distribution of the latent code $\boldsymbol{z}$. Similarly, in the posterior network, two learnable tokens are appended to the GT pose sequence $\tilde{\boldsymbol{\Theta}}'$ to output the mean $\boldsymbol{\mu}_{\boldsymbol{z}}^q$ and standard deviation $\boldsymbol{\sigma}_{\boldsymbol{z}}^q$ of the posterior distribution of the latent code $\boldsymbol{z}$.

**Hyperparameters and Training.** The dimension of the latent code $\boldsymbol{z}$ is 128. The sliding window size $h$ of the autoregressive motion infilling is 50. Both the number of context frames $h_{\mathrm{c}}$ and the number of look-ahead $h_1$ frames are 10. When synthesizing occluded motions, for any GT training motion of $h = 50$ frames, we randomly occlude $H_{\mathrm{occ}}$ consecutive frames of motion where $H_{\mathrm{occ}}$ is uniformly sampled from $[10, 40]$. Note that we do not occlude the first $h_{\mathrm{c}} = 10$ frames which are reserved as context. The KL divergence term in Eq. (2) uses a weighting factor of 0.001. We train the networks for 2000 epochs with a batch size of 1024 where each epoch uses a total of 10 million frames of motion. For optimization, we use the Adam optimizer [43] with a learning rate of 0.001 and clip the gradient if its norm is larger than 5. We use PyTorch [69] to implement and train the networks.

## D. Implementation Details for Global Trajectory Predictor

**Heading Coordinate and Egocentric Trajectory Representation.** The heading vector of a person points towards where the person is facing and is parallel to the ground. We obtain the heading vector by aligning the $z$-axis of the person's root coordinate with the world $z$-axis and use the resulting $y$-axis of the aligned root coordinate as the heading vector. This way of obtaining the heading is more stable than using the yaw of the Euler angle representation, which suffers from singularities and can be quite unstable. The heading coordinate is defined by first placing the world coordinate at the root position of the person and then rotating the world coordinate around the $z$-axis (vertical) to align the $y$-axis with the heading vector. By definition, representing and predicting human trajectories in the heading coordinate allows the predicted trajectory to be invariant of the person's absolute $xy$ translation and heading. In the egocentric trajectory representation $\boldsymbol{\psi}_t = (\delta x_t, \delta y_t, z_t, \delta \phi_t, \boldsymbol{\eta}_t)$, we use absolute height $z_t$ since the height of a person relative to the ground does not vary a lot and is highly correlated with the body motion of the person. For the local rotation $\boldsymbol{\eta}_t$, we adopt the 6D rotation representation [117] to avoid discontinuity.
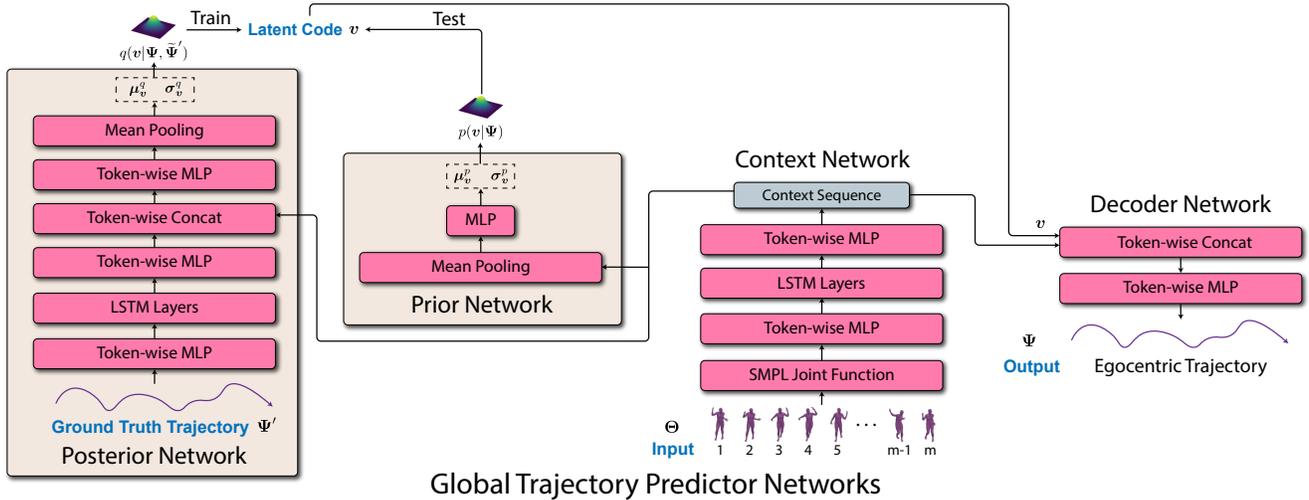


Figure 7. The network architecture of the CVAE-based global trajectory predictor. We use two bidirectional LSTM layers with hidden dimension 256 for all the LSTM blocks, and we use two hidden layers (512, 256) with ReLU activations for all the token-wise MLPs. Token-wise mean pooling is used in the prior and posterior networks to summary sequences into a single feature.

**Network Architecture.** The detailed network architecture of the CVAE-based global trajectory predictor is illustrated in Fig. 7. We use two bidirectional LSTM layers with hidden dimension 256 for all the LSTM blocks in the networks. We use two hidden layers (512, 256) with ReLU activations for all the token-wise MLPs. For the input poses, we first convert them to

3D joint positions using the SMPL joint function without global rotations and translations. This is because we find that using 3D joint positions leads to better performance than using joint rotations directly. In both the prior and posterior networks, token-wise mean pooling is used to produce a single feature from a sequence of tokens, which is then used to produce the parameters of the prior or posterior distribution of the latent code $v$.

**Hyperparameters and Training.** The dimension of the latent code $v$ is 128. The KL divergence term in Eq. (8) uses a weighting factor of 0.001. We train the networks for 2000 epochs with a batch size of 256 where each epoch uses a total of 2 million frames of motion. The training sequence length is 100 frames For optimization, we use the Adam optimizer [43] with a learning rate of 0.0001 and clip the gradient if its norm is larger than 5. We use PyTorch [69] to implement and train the networks.

## E. Implementation Details for Global Optimization

**Initialization.** We initialize the egocentric trajectories using the output from the global trajectory predictor. For the camera, we approximate the camera intrinsic parameters $K$ using the dimensions of the image where we assume the principal point is at the image center. Note that the camera intrinsics are kept fixed during the optimization process. For the camera extrinsic parameters $C$, we initialize them from the persons' global trajectories using the following equations:

$$C_t = \Omega \left( \frac{1}{\sum_{i=1}^{N} V_t^i} \sum_{i=1}^{N} V_t^i \cdot P_t^{i,\texttt{global}} P_t^{i,\texttt{cam}-1} \right) , \qquad (15)$$

where $V_t^i$ is the visibility of person $i$ at frame $t$, $P_t^{i,\texttt{global}} \in \mathbb{R}^{4\times4}$ is the person's transformation in the global coordinates based on the predicted global trajectory $(\widehat{T}^i, \widehat{R}^i)$, $P_t^{i,\texttt{cam}} \in \mathbb{R}^{4\times4}$ is the person's transformation in the camera coordinates based on the estimated trajectory $(\widetilde{T}^i, \widetilde{R}^i)$ by the pose estimator (e.g., KAMA [33]), $\Omega$ is a projection operator that projects the matrix into a valid transformation. If no person is visible at frame $t$, the camera extrinsics $C_t$ is initialized to the camera extrinsics of the most recent frame with visible people. Eq. (15) is the least squares solutions of the following (transposed) linear systems:

$$P_t^{i,\texttt{global}} = C_t P_t^{i,\texttt{cam}}, \qquad \forall i, V_t^i = 1 . \qquad (16)$$

**Hyperparameters and Optimization.** The optimization loss coefficients $(\lambda_{\texttt{2D}}, \lambda_{\texttt{traj}}, \lambda_{\texttt{reg}}, \lambda_{\texttt{cam}}, \lambda_{\texttt{pen}})$ in Eq. (9) are set to (1, 100000, 100, 10000, 100000) for 3DPW and (1, 100000, 100, 10000, 0) for Human3.6M. We do not use the inter-person penetration loss for Human3.6M since it only has one person in each video. The weighting factor $w_t$ for the translation term in Eq. (12) is set to 0 since the translation estimated by the pose estimator can be quite noisy. The trajectory regularization weighting factor $w_\psi$ in Eq. (13) is set to (3,10,10000,5,10000) for each element in the egocentric trajectory $\psi_t = (\delta x_t, \delta y_t, z_t, \delta \phi_t, \eta_t)$, where we use large weights to penalize changes in height $z_t$ and local rotation $\eta_t$. The global optimization is also implemented in PyTorch [69], where we use the Adam optimizer [43] with a learning rate of 0.001 to optimize the global trajectories and camera extrinsics.

## F. Evaluation of Global Optimization on 3DPW

We also perform experiments on 3DPW with and without our global optimization framework to study the importance of global optimization when there are multiple people in the video. Although 3DPW does not provide accurate GT human trajectories in the global coordinates, the relative translations and rotations between people in 3DPW are quite accurate. Therefore, we compute the relative translation and rotation errors between pairs of humans as an alternative way to evaluate global reconstruction quality. As shown in Table 6, using global optimization can greatly reduce the relative translation and rotation errors between humans, which means our global optimization framework can greatly help to reconstruct the spatial relationships of humans in the video.

| Method | Relative Translation Error | Relative Rotation Error |
|---|---|---|
| Ours w/o Global Optimization | 1.92 | 1.07 |
| Ours (GLAMR) | **0.66** | **0.30** |

Table 6. Evaluation of our global optimization framework on 3DPW. We evaluate the relative translation error (in meters) and relative rotation error (in angles) between pairs of humans.