

---

# LiDAR-aid Inertial Poser: Large-scale Human Motion Capture by Sparse Inertial and LiDAR Sensors

---

Chengfeng Zhao<sup>†</sup> Yiming Ren<sup>†</sup> Yannan He Peishan Cong Han Liang  
Jingyi Yu Lan Xu Yuexin Ma  
School of Information Science and Technology  
ShanghaiTech University

## Abstract

We propose a multi-sensor fusion method for capturing challenging 3D human motions with accurate consecutive local poses and global trajectories in large-scale scenarios, only using a single LiDAR and 4 IMUs. Specifically, to fully utilize the global geometry information captured by LiDAR and local dynamic motions captured by IMUs, we design a two-stage pose estimator in a coarse-to-fine manner, where point clouds provide the coarse body shape and IMU measurements optimize the local actions. Furthermore, considering the translation deviation caused by the view-dependent partial point cloud, we propose a pose-guided translation corrector. It predicts the offset between captured points and the real root locations, which makes the consecutive movements and trajectories more precise and natural. Extensive quantitative and qualitative experiments demonstrate the capability of our approach for compelling motion capture in large-scale scenarios, which outperforms other methods by an obvious margin. We will release our code and captured dataset to stimulate future research.<sup>1</sup>

## 1 Introduction

Human motion capture (mocap) has developed rapidly in recent years, which plays an important role in various applications, such as VR/AR, gaming, sports analysis, and movie production. However, accurate capture of challenging human motions in large-scale scenarios hasn't been settled well, which is critical for the reconstruction, simulation, and generation of sporting mega-events, stage performances, interactions of crowds, etc.

By far, optical-based solutions take the majority of human mocap. The high-end marker-based solutions [53, 54, 39] require outside-in multi-camera setup or dense optical markers, and thus confine the performers to a constrained captured area, making large-scale capture impractical. Recently, learning-based monocular methods [21, 22, 23, 67, 13, 14, 26, 43, 24] enable robust motion capture under light-weight setting. Although alleviating expensive facilities and fixed captured region, they remain vulnerable to occlusions, lack of textures and severe changes of environment lighting, etc.. Moreover, their inherent lack of depth measurement makes them unstable to accurately track global trajectories of humans, especially under large-scale setting.

In contrast, motion capture using inertia information recorded by Inertial Measurement Units (IMUs) is occlusion-unaware and environment-independent. The high-end solutions [58, 38] require a large amount of body-worn IMUs (from 8 to 17), making them unsuitable to capture human motions with everyday apparel. Recent data-driven advances [18, 64, 63] enable real-time motion capture

---

<sup>1</sup>†: Equal Contribution

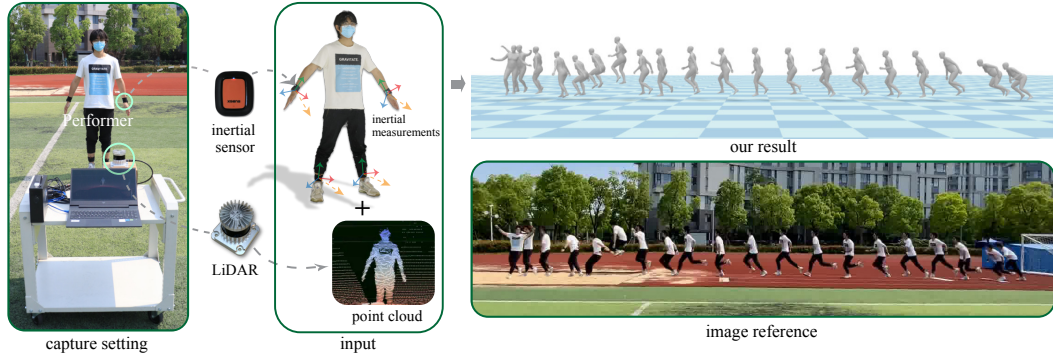


Figure 1: We propose a multi-modal motion capture approach, LIP, that estimates challenging human motions with accurate local pose and global translation in large-scale scenarios using single LiDAR and four inertial measurement units.

with sparse IMUs. They can obtain global location by accumulating the IMU observation and even partially alleviate the drifting with foot-ground contact or physical constraints. Yet, such purely inertial methods still inherently suffer from temporally cumulative error of global localization, especially for capturing large-scale motions with a long duration. Only recently, LiDAR-based mocap is gaining increasing attention due to the tremendous progress of LiDAR in large-scale perception [71, 70, 7]. Notably, LiDARCap [28] leverages a graph-based convolutional network to predict daily human poses from the point clouds captured by a LiDAR sensor within range of 30 m. However, it’s fragile to capture challenging human motions with strong self-occlusions due to the sparsity of point clouds input, especially for large-scale capture.

To tackle the above challenges, we propose a novel approach called LiDAR-aid Inertial Poser (LIP), to capture challenging human motions in large-scale scenarios. In stark contrast with previous mocap system, our LIP adopt a novel and light-weight hardware setup using only a monocular LiDAR and 4 IMUs. These two kinds of sensors are complementary to each other, providing both global geometry and local dynamics information of the captured performer. Meanwhile, they are naturally privacy-preserving and insensitive to the lighting, which is appropriate to be generalized to novel scenes. As shown in Figure 1, our system faithfully reconstructs both local skeletal poses and global trajectories of the performer under a novel sensor-fusion setting.

Generating robust mocap results using such novel multi-modal inputs is non-trivial. First, point cloud only capture the global visual information, while IMU measurements encode the local actions physically. Furthermore, the partially captured point clouds of human body will significantly change under various poses and views of the performer, which negatively affects the accurate localization and naturally consecutive motion capture. To this end, in LIP, we adopt two consistent technical modules to accurately estimate human poses and translations, respectively. For the former one, we introduce a multi-modal pose estimation module, which jointly utilize the point cloud and motion inertia observations in a coarse-to-fine manner. Specifically, we first extract pose and motion patterns from point cloud by regressing 24 joints of the parametric human model SMPL [31] and the 6D rotation of the root joint through a convolutional temporal encoder (PointNet-GRU). Next, a hierarchical pose estimation process with fused IMU measurements is appended to refine the coarse skeleton by a joint-map estimator and solve the inverse kinematics(IK) problem for the joint rotations by a body-pose estimator. To precisely capture the 3D global motion trajectories, we introduce the second module called Pose-guided Translation Correction, so as to learn the deviation between the partially captured point cloud and the real location of the root joint. Specifically, considering the fact that diverse poses can affect the deviation, we uses another isomorphic PointNet-GRU structure to model the shape and pose characteristics from point cloud, cooperated with the estimated joint rotations and refined skeleton. We further utilize the fused pose feature to eliminate the inherent translation deviation and random noise from the captured point cloud. We demonstrate the capability of our LIP on a variety of real and synthetic datasets with LiDAR and IMU measurements, covering various large-scale scenes and challenging human motions. To summarize, our main contributions include:

- We propose the first LiDAR-IMU hybrid approach for 3D human motion capture in large-scale scenarios and achieves state-of-the-art performance.

- We propose an effective two-stage coarse-to-fine fusion method to fully utilize the complementary features of multi-modal input for accurate pose estimation.
- We propose an approach to eliminate the translation deviation in a pose-guided manner, achieving accurate global trajectories and natural consecutive actions.

## 2 Related Work

**Optical Motion Capture.** Marker-based motion capture studios [53, 54, 39] enable capturing high-quality professional motions, which have achieved success and are widely used in the industry. However, these systems are costly and cumbersome, and performers usually need to wear the marker suits, which means unavailable for daily usage. To overcome these problems, the exploration of markerless mocap [5, 8, 51, 60, 16, 47, 48, 50, 19, 59, 66, 33] has made great progress. Previous multi-view algorithms [1, 6, 9, 44, 45, 41, 49] demonstrate robust motion capture even in the wild. Although the cost and intrusiveness is drastically reduced, synchronizing and calibrating multi-camera systems is still tedious. Thus various monocular mocap approaches have been proposed, which estimate 3D human pose and shape by optimizing [17, 27, 4, 25] or directly regressing [21, 22, 23, 67]. To overcome various flaws of monocular setup, template-based approaches [62, 12, 61, 13, 14], probabilistic approaches [26, 43] and semantic-modeling approaches [24] are proposed. However, the inherent depth ambiguity is still unsolved. Some approaches [46, 3, 57, 65, 11] using the commodity depth cameras enable alleviating this, but these active IR-based cameras are unsuitable for outdoor capture and the capture volume is limited. Recently, Li et al. [28] employs a consumer-level LiDAR which provides large-scale depth information, to enable large-scale 3D human motion capture, but it still suffers from severe self-occlusion.

**Inertial Motion Capture.** In contrast to optical approaches based on cameras, purely inertial approaches using IMUs are free from occlusion and restricted recording environment and volume. However, commercial purely inertial solutions such as Xsens MVN [58] rely on large amounts of IMUs. Performers are usually required to wear tight suits with densely bounded IMUs, which is intrusive, tedious and inconvenient, and prompt the community forward to the sparse setup. A pioneering work, SIP [56], presents an optimization-based offline method using only 6 IMUs and achieves promising results. Inspired by it, recent data-driven approaches [18, 64, 63] with the same setup achieve great improvements in accuracy and efficiency, which enable real-time pose and translation estimation. Nevertheless, these purely inertial approaches still suffer from substantial drifts while performing high-speed or large-scale capture.

**Hybrid Motion Capture.** As optical and inertial mocap solutions suffer from occlusions and drifts, respectively. Approaches that fuse these two types of sensors to benefit from complementarity, so as to achieve more robust mocap, have attracted much attention. Preceding methods propose to combine IMUs with RGB cameras, which can be achieved by either optimization [42, 20, 37, 35, 36] or regression [10, 52, 55, 68]. Recently, Liang et al. [30] presents a learning-and-optimization method fusing a single camera with only 4 IMUs, which demonstrates robust challenging motion capture. Besides, some works fuse IMUs with depth cameras [15, 69] or optical markers [2], which achieve promising results. Nevertheless, these methods still suffer from limited capture volume, which limits the usage for large-scale motion capture. In this work, we propose a novel hierarchical framework fusing a single LiDAR and only 4 IMUs, which alleviates both occlusions and drifts and achieves accurate large-scale 3d challenging motion capture.

## 3 Overview and Preliminaries

Our goal is to capture challenging 3D human motions in large-scale scenario with consistently accurate local pose and global trajectory estimation using single LiDAR and 4 IMUs. Figure 2 illustrates an overview of the entire pipeline, which consists of two cooperative modules to infer pose and translation, respectively. For pose inference module, we propose a two-stage network working in a coarse-to-fine manner to distill rough human body skeleton and global orientation from raw point cloud first, and then regress human body pose from joint positions refined by IMU measurements (Section 4.1). For translation inference module, we design a pose-guided approach to learn the latent domain gap between LiDAR measurements and real global movements, through which the inherent translation deviation and random noise from point cloud can be eliminated (Section 4.2).

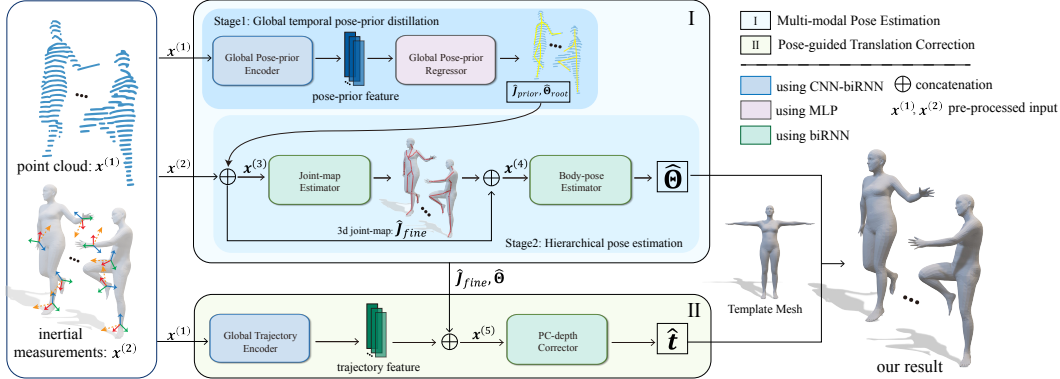


Figure 2: Overview of our pipeline. It consists of two cooperative modules: *Multi-modal Pose Estimation* and *Pose-guided Translation Correction* to estimate skeletal pose and global translation from sparse IMU and LiDAR inputs.

### 3.1 Data Preparation

To learn the local pose and global translation estimator working for large-scale scenario, a huge LiDAR-IMU hybrid mocap dataset is required, which is unfortunately unavailable at present to the best of our knowledge. Due to the inconvenience and constraint on performers to act arbitrary motions, it is difficult to collect large amounts of data with dense IMU sensors. Thus, we simulate plenty of synthetic LiDAR-IMU measurements with ground-truth SMPL pose parameters on DIP-IMU, LiDARHuman26M, AIST++ and a subset of AMASS, including ACCAD, BMLMovi, CMU, and TotalCapture. For simulation details and rationality analysis, please refer to B. Furthermore, in order to demonstrate that our network has great generalization capability to handle completely real LiDAR-IMU hybrid signals, we record LIP-Dataset(LIPD) by Xsens Dot IMU sensors[58] and OUSTER-1 LiDAR[40] in diverse wild scenes with many challenging poses, like dancing and sports. All datasets combined consist of 267 subjects, 4238 motion genres, and over 25 hours of data. As for the protocol of dataset splitting, we take DIP-IMU, TotalCapture, LIPD and the testing set of LiDARHuman26M only for evaluation. Table 1 gives an overview of all the datasets.

Table 1: Overview of datasets we use. The data mode of point cloud and IMU measurements for each dataset is reported. Moreover, "Capture distance" shows the maximum distance between performer and capture device; "Activity range" stands for the average 3D space size in which the performer moves. Both of them reflect the scale attribute of each dataset.

Dataset	Point cloud	IMU recording	Capture distance(m)	Activity range(m <sup>3</sup> )
DIP-IMU [18]	Synthetic	Real	N/A	N/A
AMASS [34]	Synthetic	Synthetic	3.42	0.27
AIST++ [29]	Synthetic	Synthetic	4.23	0.67
LiDARHuman26M [28]	Real	Synthetic	28.05	142.95
LIPD	Real	Real	21.07	366.34

## 4 Method

In this section, we give detailed explanations for the design of our network. Before that, we clarify our system input pre-processing and the most frequently used math symbols in the following.

**System Input Pre-processing:** Since raw point cloud sequence with variable  $N_p(t)$  points at different time frame  $t$  is temporally inconsistent, we normalize the point cloud in every single frame  $\tilde{\mathbf{x}}^{(1)}(t) \in \mathbb{R}^{N_p(t) \times 3}$  to  $\mathbf{x}^{(1)}(t) \in \mathbb{R}^{N_{fps} \times 3}$  by subtracting its arithmetic mean and upsampling to fixed  $N_{fps}$  points using farthest point sampling algorithm(FPS). In our implementation, we set  $N_{fps} = 256$ . For IMU measurements, we transform sensors' inertia in inertial coordinate frame to bones' inertia in LiDAR coordinate frame, and formulate single-frame inertial input as  $\mathbf{x}^{(2)}(t) = [\mathbf{R}_{lw}, \mathbf{R}_{rw}, \mathbf{R}_{la}, \mathbf{R}_{ra}, \mathbf{a}_{lw}, \mathbf{a}_{rw}, \mathbf{a}_{la}, \mathbf{a}_{ra}] \in \mathbb{R}^{48}$ , where  $\mathbf{R}$  denotes the rotation in flattened



rotation matrix form while  $\mathbf{a}$  indicates the free acceleration, and  $lw, rw, la, ra$  mean left wrist, right wrist, left ankle, and right ankle, respectively.

**Definition of Motion Representations:** We define  $N_j, N_s$  as the amount of body joints and IMU sensors;  $\hat{\mathbf{J}}(t), \mathbf{J}^{GT}(t), \tilde{\mathbf{J}}^{GT}(t) \in \mathbb{R}^{3N_j}$  as predicted root-relative joint positions, ground-truth root-relative joint positions, and ground-truth absolute joint positions at time  $t$ ;  $\hat{\Theta}(t), \Theta^{GT}(t) \in \mathbb{R}^{6N_j}$  as predicted joint rotations and ground-truth joint rotations in 6D rotation representation at time frame  $t$ . Note that all the formulation of loss functions defined below omit a common factor  $\frac{1}{T}$  where  $T$  is the total time length of training motion sequences.

#### 4.1 Multi-modal Pose Estimation

Purely inertial or LiDAR-based methods more or less suffer from insufficient observations. On one hand, without global visual cue, the reconstruction from local inertia to 3D joint position is ambiguous. On the other hand, the lack of direct measurements of local motion disables accurate pose inference and robustness to occlusions. Therefore, we propose to estimate body pose with multi-modal input, which includes both global geometry information and local dynamic motion inertia. However, since point cloud is in spatial form while IMU measurement is physical quantity, it is irrational to concatenate  $\mathbf{x}^{(1)}(t)$  and  $\mathbf{x}^{(2)}(t)$  as network input directly. Instead, we formulate this module as a two-stage network working in a coarse-to-fine manner so that the point cloud can combine with motion inertia efficiently. *Global Temporal Pose-prior Distillation* is designed to extract hidden geometric feature and motion pattern from normalized point cloud and express it explicitly. After that, *Hierarchical Pose Estimation* fuses IMU measurements in and regresses joint rotations from refined 3D joint positions.

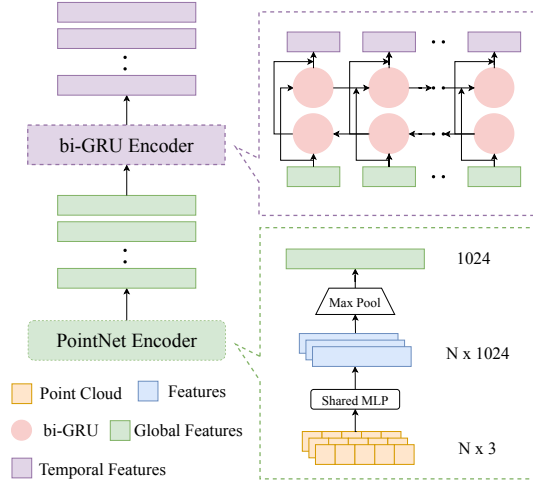


Figure 3: Detailed architecture of *Encoder Block*, consisting of PointNet and bidirectional GRU.

**Global Temporal Pose-prior Distillation:** Considering that the raw point cloud can directly represent the coarse human shape information and sparse IMUs can not directly estimate the root joint orientation, we distill the human-skeleton-joint positions and the root joint orientation from the point cloud in the first stage and then use the four IMU sensors in the bone to refine the result and regress the shape parameters. Specifically, We propose a global temporal pose-prior distillation to regress the 24 joints of the SMPL mesh and the 6D rotation representation of the root joint, which is composed of global feature extractor PointNet and temporal encoder two-way GRU(bi-GRU), as Figure 3 shows. PointNet is used as encoder to extract human skeleton geometric information from the raw point cloud as a feature vector  $v(t) \in \mathbb{R}^k$  from each frame  $F(t)$ , where  $k = 1024$ . We feed the frame-wise features  $v(t)$  into the two-way GRU(bi-GRU) to generate the hidden variables  $h(t)$  to extract temporal information. Then, we use the MLP decoder to predict the joint positions  $\hat{\mathbf{J}}_{prior}(t)$  and the orientation of root joint  $\hat{\Theta}_{root}(t)$ , which are the part of the input in the second stage. We extract 24 joints on the SMPL mesh  $\mathbf{J}^{GT}(t)$  and select 6D rotation representation of the root joint from SMPL pose parameters  $\Theta_{root}^{GT}(t)$  as the ground truth. The losses of the above two supervision information can be formulated as

$$\mathcal{L}_{joint-prior} = \sum_t \|\hat{\mathbf{J}}_{prior}(t) - \mathbf{J}^{GT}(t)\|_2^2, \quad (1)$$

$$\mathcal{L}_{ori-prior} = \sum_t \|\hat{\Theta}_{root}(t) - \Theta_{root}^{GT}(t)\|_2^2, \quad (2)$$

$$\mathcal{L}_{prior} = \lambda_1 \mathcal{L}_{joint-prior} + \lambda_2 \mathcal{L}_{ori-prior}, \quad (3)$$

where  $\lambda_1$  and  $\lambda_2$  are hyper-parameters.

**Hierarchical Pose Estimation:** Although the explicit pose-prior distilled from the point cloud is available for a neural Inverse Kinematics (IK) solver, it is too coarse to give accurate 3D joint positions and robust inference for challenging motions due to the lack of locally precise motion observations. Therefore, we hierarchically organize this stage as two estimators, which refine the root-relative 3D joint positions with the aid of IMU measurements and then regress the joint rotations. These two estimators have the same architecture, composed of three GRUs connected sequentially with a skip connection, which can make use of temporal motion information and maintain a healthy back propagation for our deep network during training. The former one, *Joint-map Estimator*, takes  $\mathbf{x}^{(3)}(t) = [\mathbf{x}^{(2)}(t), \hat{\mathbf{J}}_{prior}(t), \hat{\Theta}_{root}(t)] \in \mathbb{R}^{12N_s+3N_j+6}$  as the input and outputs a refined root-relative 3D joint-map with more accurate positions  $\hat{\mathbf{J}}_{fine}(t)$ , supervised by loss function

$$\mathcal{L}_{fineJ} = \sum_t \|\hat{\mathbf{J}}_{fine}(t) - \mathbf{J}^{GT}(t)\|_2^2. \quad (4)$$

After that,  $\mathbf{x}^{(4)}(t) = [\mathbf{x}^{(3)}(t), \hat{\mathbf{J}}_{fine}(t)] \in \mathbb{R}^{12N_i+3N_j+3N_j+6}$  is fed into *Body-pose Estimator* which learns to solve an IK-like problem to regress joint rotations  $\hat{\Theta}(t)$ , supervised by the loss function

$$\mathcal{L}_{ik} = \sum_t \|\hat{\Theta}(t) - \Theta^{GT}(t)\|_2^2. \quad (5)$$

Because of the limited expression capability of joint rotations, we introduce loss function 6 for *Body-pose Estimator* by Forward Kinematics (FK) to better reconstruct the 3D joint positions.

$$\mathcal{L}_{fk} = \sum_t \|\text{FK}(\hat{\Theta}(t)) - \mathbf{J}^{GT}(t)\|_2^2. \quad (6)$$

The complete loss function of this module is formulated as

$$\mathcal{L}_{pose} = \lambda_3 \mathcal{L}_{fineJ} + \lambda_4 \mathcal{L}_{ik} + \lambda_5 \mathcal{L}_{fk}. \quad (7)$$

## 4.2 Pose-guided Translation Correction

Global translation estimation, especially in the large-scale scenario, is challenging for purely inertial mocap methods since no localization information can be provided by IMUs and IMUs also suffer from the drifting problem. Even in the state-of-the-art work (TransPose [64]), the capture of global movements should be performed under strong assumptions such as level ground and sufficient foot-ground contacts. LiDAR, however, brings significant benefits to this task by measuring distances directly. LiDARCap [28] utilizes the average of points as global translation, however, there exists a deviation between scanned point cloud and real movement positions, because LiDAR can only collect partial points on the performer in the perspective view. To solve this problem, we combine *encoder* and *estimator* block to model the motion pattern from the point cloud and learn this deviation, rather than regress the global translation directly. We treat the translation discrepancy as a per-pose variable, which is only correlated with the pose performed. For example, standing still and bending result in similar global translation positions, while different averages of points. We construct the relationships between poses and translation deviations. The ground-truth translation discrepancy  $\mathbf{D}^{GT}(t) \in \mathbb{R}^3$  can be calculated as follows:

$$\mathbf{D}^{GT}(t) = \tilde{\mathbf{J}}_{root}^{GT}(t) - \text{avg}(\tilde{\mathbf{x}}^{(1)}(t)), \quad (8)$$

where  $\tilde{\mathbf{J}}_{root}^{GT}(t)$  is the ground-truth absolute position of the root joint and the operator  $\text{avg}(\cdot)$  calculates the approximate center of given point cloud by averaging the positions of all the points. Finally, we use loss function 9 to train this module:

$$\mathcal{L}_{trans} = \sum_t \|\hat{\mathbf{D}}(t) - \mathbf{D}^{GT}(t)\|_2^2, \quad (9)$$

where  $\hat{\mathbf{D}}(t)$  denotes the predicted translation discrepancy result.

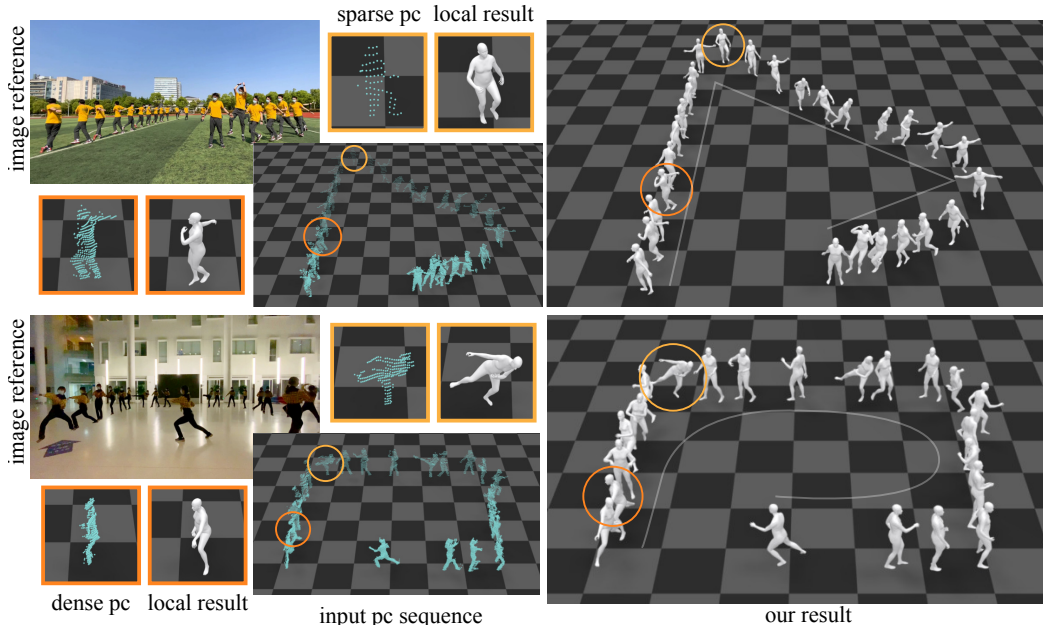


Figure 4: Our results. We provide image reference, point cloud input, and mocap results in 3D scene in this figure. Some key frames are circled and zoomed in to show detailed local pose.

### 4.3 Implementation Details

We implement our network using PyTorch 1.8.1 with CUDA 11.1. The training of the two modules is separate, which uses batch size of 32, collated time split of 32, learning rate of  $10^{-4}$  and decay rate of  $10^{-4}$  with AdamW optimizer [32]. The weights of loss functions are:  $\lambda_1 = 1$ ,  $\lambda_2 = 1$ ,  $\lambda_3 = 0.2$ ,  $\lambda_4 = 0.7$ ,  $\lambda_5 = 0.7$ . All the training and evaluation are conducted on a server with an Intel(R) Xeon(R) E5-2678 CPU and an NVIDIA RTX3090 graphics card. Furthermore, we use free acceleration input, which escapes from the influence of gravity in inertial coordinate frame.

## 5 Experiments

In this section, we compare our method with State-Of-The-Art(SOTA) methods qualitatively and quantitatively to show the superiority of the proposed LIP in Section 5.1. In addition, in Section 5.2, sufficient ablation studies are conducted to further demonstrate the rationality of our architecture design and input configuration. For evaluation metrics, we use 1) MPJPE: mean per root-relative joint position error in millimeters; 2) Mesh Err: mean per SMPL mesh vertex position error in millimeters; 3) Ang Err: mean per global joint rotation error in degrees to evaluate local pose, and 4) CD: chamfer distance between vertices of SMPL mesh result and raw point cloud in centimeters to evaluate global translation. For all these metrics, the lower is better.

### 5.1 Comparison

We conduct comparisons to illustrate that our proposed LIP method enables more accurate capture in challenging motions and more precise translations in large-scale spaces. We select current SOTA Transpose [64] and LiDARCap [28] for comparative analysis. We use the model provided by Transpose and reproduced LiDARCap network (no released model) for comparison. The superior results with different evaluation metrics are illustrated in Table 2. As shown in Figure 5. Benefiting from LiDAR-IMU hybrid modal input and our coarse-to-fine design, our method outperforms others by an obvious margin. To take advantage of the 3D distinguish-ability of the LiDAR point clouds, we design the pose-guided translation correction. We use our LIPD dataset which consists of abundant challenge motions with large-scale spaces to compare the estimated translation. Our method performs much better than Transpose and LiDARCap as shown in Figure 5.

Table 2: Quantitative comparisons between LIP and related methods on our evaluation dataset. Note that LiDARHuman26M dataset only contains data with rate of 10fps, which is not applicable(N/A) for TransPose. Also, the CD metric is not used for DIP-IMU dataset since no translation is provided.

	TotalCapture [52]				DIP-IMU [18]			LiDARHuman26M [28]			
	MPJPE	Mesh Err	Ang Err	CD	MPJPE	Mesh Err	Ang Err	MPJPE	Mesh Err	Ang Err	CD
LiDARCap[28]	44.7	56.9	9.1	6.4	43.2	56.1	12.8	79.0	96.6	19.7	8.1
TransPose[64]	55.7	63.6	12.3	102.7	49.0	58.3	7.6	N/A	N/A	N/A	N/A
LIP	<b>28.6</b>	<b>37.2</b>	<b>5.8</b>	<b>0.7</b>	<b>29.7</b>	<b>38.3</b>	<b>8.3</b>	<b>59.7</b>	<b>73.4</b>	<b>6.14</b>	<b>3.5</b>

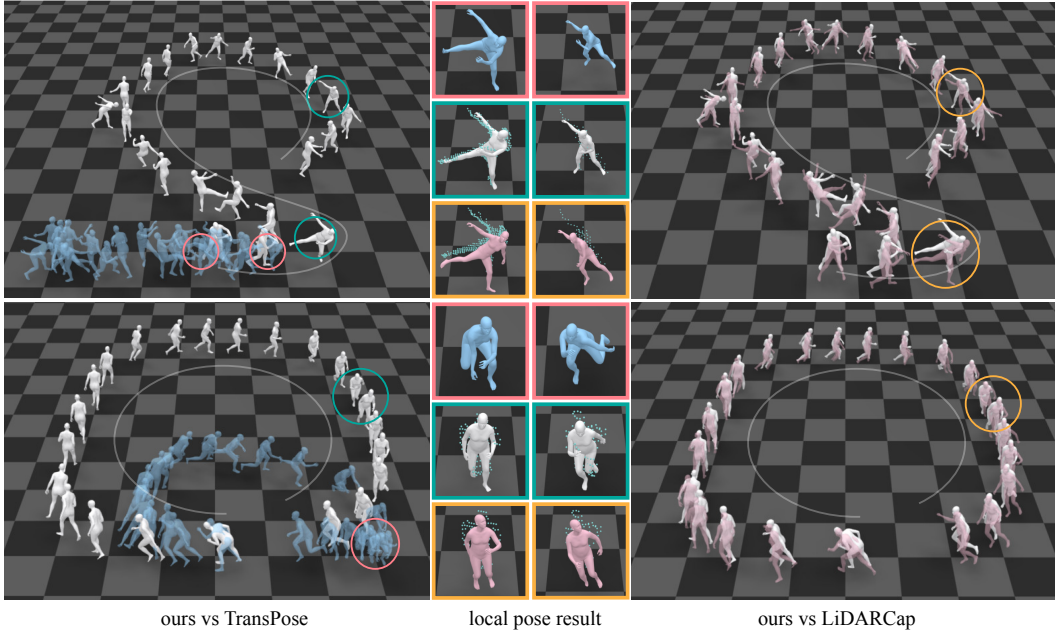


Figure 5: Qualitative results of LIP(white) compared with TransPose(blue) and LiDARCap(pink). The trajectory is sketched by a grey curve. Detailed comparisons of some key frames are circled and zoomed in to show the local pose and alignment with point cloud. For intuitive visualization, we use the arithmetic mean of raw point cloud as the global translation for LiDARCap in which translation inference is not included.

## 5.2 Evaluation

We validate the effectiveness of our network design and the reasonability of the configuration of our multi-modal hybrid input by ablation studies. Qualitative evaluations are illustrated in Figure 6.

**Ablation Study on Network Architecture:** We demonstrate the superiority of our two-stage coarse-to-fine structure and evaluate the following two variants of our network: 1) the temporal features obtained from the raw point cloud are directly combined with IMU data as the input of the second stage, without the coarse joints and root orientation estimation; 2) the model directly regresses the pose parameters from coarse joint estimation in stage one, without any refinement by IMUs. We compare the above variants of our pipeline on the testing set mentioned in Section 3.1 in Table 3. It can be observed that regressing the coarse skeleton and root orientation as the intermediate task helps the pose estimation, which is more intuitive to guide the subsequent generation of finer joint positions instead of using high-dimensional features. The second refinement stage further benefits more accurate results in the estimation of joints and vertices. To verify the effectiveness of Pose-guided Translation Correction module, we conduct comparison with other location estimation methods and ablation study. As Table 3 shows, the design of deviation regression and the usage of pose features benefit the performance a lot.

**Ablation Study on Input Configuration:** To verify that single LiDAR with 4 IMUs is a necessary and compact configuration for high-quality mocap, we experiment various combinations of this two modal inputs with the same network architecture and training strategy. Table 7 gives quantitative comparisons which demonstrate that LiDAR-IMU hybrid input overall outperforms single modality

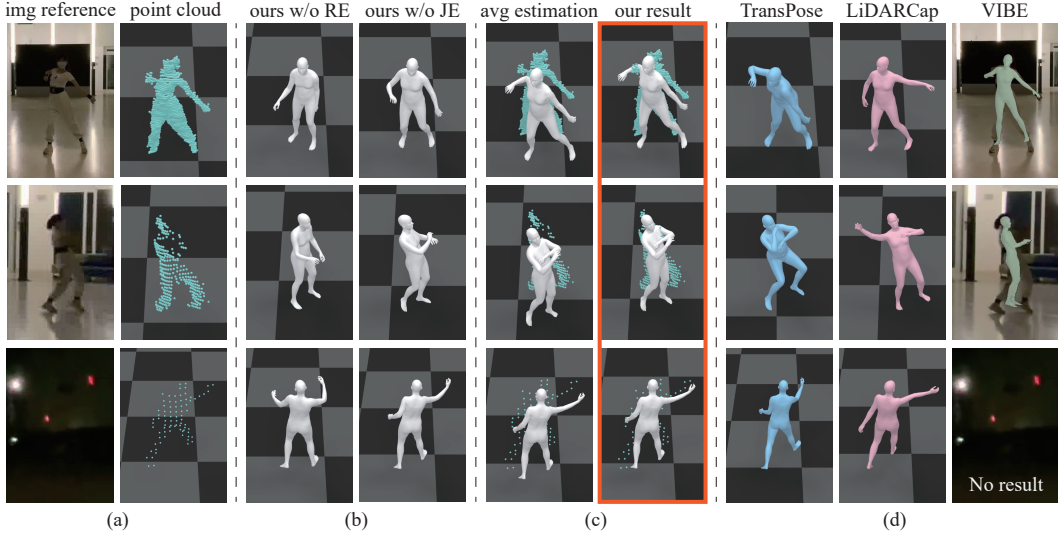


Figure 6: Qualitative evaluations of our design. We provide (a) image reference and point cloud input; (b) local pose evaluation; (c) global translation evaluation; (d) results of other methods.

input. Specifically, direct measurements of 3D scene provided by LiDAR reduce the ambiguities of regression from sparsely local inertia to 3D joint positions, which so that decrease the MPJPE and mesh error by more than 30 and 40 millimeters and angular error by 2 degrees over purely inertial input in average. Meanwhile, with the aid of local inertia, inaccurate estimation of joint positions and rotations on point cloud can be corrected in detail. Especially for LiDARHuman26M dataset, of which the point cloud data is real and imperfect with random noise, LIP improves the performance of LiDAR-only input by 10.5 millimeters lower MPJPE, 16 millimeters lower mesh error and 1.8 degrees lower angular error. Besides, Table 7 illustrates that 4-IMU-aid configuration brings considerable improvements with only two more sensors compared with 2-IMU-aid version, and keeps acceptable performance gap to even 12-IMU-aid setting. Therefore, 4 IMU sensors is an appropriate choice to provide sufficient inertial aid while maintain a light-weight capture setting, considering the complexity of capture system and convenience for performers.

Table 3: Ablation study on network architecture. "RE" means *Global pose-prior Regressor* and "JE" indicates *Joint-map Estimator*.

Pose estimation	MPJPE	Mesh Err	Ang Err
Ours w/o RE	58.6	70.7	11.5
Ours w/o JE	38.1	47.7	9.9
Ours	<b>31.2</b>	<b>40.0</b>	<b>9.0</b>
Translation estimation			CD
Average estimation			7.2
Ours w/o PG			2.1
Ours			<b>1.9</b>

Table 4: Ablation study on input configuration.

	MPJPE	Mesh Err	Ang Err
LiDAR Only	32.9	43.9	14.7
4 IMUs Only	87.4	98.6	13.2
5 IMUs Only	74.5	85.2	11.4
LiDAR+2 IMUs	31.8	41.6	10.8
LiDAR+4 IMUs(Ours)	<b>31.2</b>	<b>40.0</b>	<b>9.0</b>
LiDAR+5 IMUs	30.3	39.3	8.8
LiDAR+12 IMUs	30.2	39.0	5.5

### 5.3 Limitation and Future Work

Due to the low frame rate of the utilized LiDAR (10fps), it is difficult to capture extremely fast motions, resulting in unsmooth fast motion sequences. In future, we will improve the prediction modules using LSTM or Transformer for interpolating global poses to capture high-frequency motion details. Moreover, we will enhance the capability of our method on handling domain gaps between different LiDAR sensors with different point distributions or between complete and incomplete point clouds caused by occlusions. In addition, limited by the simulated training data using the SMPL model, the generalization capability of our method requires improvement. Thus, we will enrich the diversity of our training dataset including people with different clothing, multi-person interactions, human-object interactions, etc.

## 6 Conclusion

In this paper, we propose a new solution, LiDAR-aid Inertial Poser, for human motion capture with high-quality estimation of poses and translations in large-scale scenarios, which is occlusion-free and environment-independent. Based on our multi-modal captured data, we design an effective fusion strategy for taking advantage of the global location and geometry features of LiDAR point cloud and local fine-grained motion information from IMU measurements. We also eliminate the translation deviation in a pose-guided manner and gain accurate global trajectories and natural consecutive actions. Extensive comparisons illustrate the superiority of our method, while ablation studies verify the reasonability of our input configuration and effectiveness of our network modules.

## References

- [1] Sikander Amin, Mykhaylo Andriluka, Marcus Rohrbach, and Bernt Schiele. Multi-view pictorial structures for 3D human pose estimation. In *British Machine Vision Conference (BMVC)*, 2009.
- [2] Sheldon Andrews, Ivan Huerta, Taku Komura, Leonid Sigal, and Kenny Mitchell. Real-time physics-based motion capture with sparse sensors. pages 1–10, 12 2016.
- [3] Andreas Baak, Meinard Müller, Gaurav Bharaj, Hans-Peter Seidel, and Christian Theobalt. A data-driven approach for real-time full body pose reconstruction from a depth camera. In *International Conference on Computer Vision (ICCV)*, 2011.
- [4] 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 Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *Computer Vision – ECCV 2016*, pages 561–578, Cham, 2016. Springer International Publishing.
- [5] Chris Bregler and Jitendra Malik. Tracking people with twists and exponential maps. In *Computer Vision and Pattern Recognition (CVPR)*, 1998.
- [6] Magnus Burenius, Josephine Sullivan, and Stefan Carlsson. 3D pictorial structures for multiple view articulated pose estimation. In *Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [7] Peishan Cong, Xinge Zhu, Feng Qiao, Yiming Ren, Xidong Peng, Yuenan Hou, Lan Xu, Ruigang Yang, Dinesh Manocha, and Yuexin Ma. Stcrowd: A multimodal dataset for pedestrian perception in crowded scenes. *arXiv preprint arXiv:2204.01026*, 2022.
- [8] Edilson De Aguiar, Carsten Stoll, Christian Theobalt, Naveed Ahmed, Hans-Peter Seidel, and Sebastian Thrun. Performance capture from sparse multi-view video. pages 1–10, 2008.
- [9] Ahmed Elhayek, Edilson de Aguiar, Arjun Jain, Jonathan Tompson, Leonid Pishchulin, Mykhaylo Andriluka, Chris Bregler, Bernt Schiele, and Christian Theobalt. Efficient ConvNet-based marker-less motion capture in general scenes with a low number of cameras. In *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [10] Andrew Gilbert, Matthew Trumble, Charles Malleon, Adrian Hilton, and John Collomosse. Fusing visual and inertial sensors with semantics for 3d human pose estimation. *International Journal of Computer Vision*, 127(4):381–397, 2019.
- [11] Kaiwen Guo, Jonathan Taylor, Sean Fanello, Andrea Tagliasacchi, Mingsong Dou, Philip Davidson, Adarsh Kowdle, and Shahram Izadi. Twinfusion: High framerate non-rigid fusion through fast correspondence tracking. In *International Conference on 3D Vision (3DV)*, pages 596–605, 2018.
- [12] Marc Habermann, Weipeng Xu, Michael Zollhöfer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions on Graphics (TOG)*, 38(2):14:1–14:17, 2019.
- [13] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [14] Yannan He, Anqi Pang, Xin Chen, Han Liang, Minye Wu, Yuexin Ma, and Lan Xu. Challengcap: Monocular 3d capture of challenging human performances using multi-modal references. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11400–11411, 2021.
- [15] Thomas Helten, Meinard Muller, Hans-Peter Seidel, and Christian Theobalt. Real-time body tracking with one depth camera and inertial sensors. In *Proceedings of the IEEE international conference on computer vision*, pages 1105–1112, 2013.
- [16] Michael B. Holte, Cuong Tran, Mohan M. Trivedi, and Thomas B. Moeslund. Human pose estimation and activity recognition from multi-view videos: Comparative explorations of recent developments. *Journal of Selected Topics in Signal Processing*, 6(5):538–552, 2012.



- [17] Y. Huang, F. Bogo, C. Lassner, A. Kanazawa, P. V. Gehler, J. Romero, I. Akhter, and M. J. Black. Towards accurate marker-less human shape and pose estimation over time. In *2017 International Conference on 3D Vision (3DV)*, pages 421–430, 2017.
- [18] Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J Black, Otmar Hilliges, and Gerard Pons-Moll. Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Transactions on Graphics (TOG)*, 37(6):1–15, 2018.
- [19] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *International Conference on Computer Vision (ICCV)*, 2015.
- [20] Tomoya Kaichi, Tsubasa Maruyama, Mitsunori Tada, and Hideo Saito. Resolving position ambiguity of imu-based human pose with a single rgb camera. *Sensors*, 20(19):5453, 2020.
- [21] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [22] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [23] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [24] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. Pare: Part attention regressor for 3d human body estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11127–11137, October 2021.
- [25] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [26] Nikos Kolotouros, Georgios Pavlakos, Dinesh Jayaraman, and Kostas Daniilidis. Probabilistic modeling for human mesh recovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11605–11614, 2021.
- [27] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6050–6059, 2017.
- [28] Jialian Li, Jingyi Zhang, Zhiyong Wang, Siqi Shen, Chenglu Wen, Yuexin Ma, Lan Xu, Jingyi Yu, and Cheng Wang. Lidarcap: Long-range marker-less 3d human motion capture with lidar point clouds. *arXiv preprint arXiv:2203.14698*, 2022.
- [29] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++, 2021.
- [30] Han Liang, Yannan He, Chengfeng Zhao, Mutian Li, Jingya Wang, Jingyi Yu, and Lan Xu. Hybridcap: Inertia-aid monocular capture of challenging human motions. *arXiv preprint arXiv:2203.09287*, 2022.
- [31] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: A skinned multi-person linear model. *ACM Trans. Graph.*, 34(6):248:1–248:16, October 2015.
- [32] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [33] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *Advances in Neural Information Processing Systems*, 34, 2021.
- [34] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- [35] Charles Malleson, John Collomosse, and Adrian Hilton. Real-time multi-person motion capture from multi-view video and imus. *International Journal of Computer Vision*, pages 1–18, 2019.
- [36] Charles Malleson, John Collomosse, and Adrian Hilton. Real-time multi-person motion capture from multi-view video and imus. *International Journal of Computer Vision*, 128(6):1594–1611, 2020.
- [37] Charles Malleson, Andrew Gilbert, Matthew Trumble, John Collomosse, Adrian Hilton, and Marco Volino. Real-time full-body motion capture from video and imus. In *2017 International Conference on 3D Vision (3DV)*, pages 449–457. IEEE, 2017.
- [38] Noitom Motion Capture Systems. <https://www.noitom.com/>, 2015.
- [39] OptiTrack Motion Capture Systems. <https://www.optitrack.com/>, 2009.
- [40] OUSTER High Performance Digital Lidar Solutions. <https://ouster.com/>, 2021.



- [41] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. Harvesting multiple views for marker-less 3d human pose annotations. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [42] Gerard Pons-Moll, Andreas Baak, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bodo Rosenhahn. Multisensor-fusion for 3d full-body human motion capture. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 663–670. IEEE, 2010.
- [43] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11488–11499, October 2021.
- [44] Helge Rhodin, Nadia Robertini, Christian Richardt, Hans-Peter Seidel, and Christian Theobalt. A versatile scene model with differentiable visibility applied to generative pose estimation. In *International Conference on Computer Vision (ICCV)*, 2015.
- [45] Nadia Robertini, Dan Casas, Helge Rhodin, Hans-Peter Seidel, and Christian Theobalt. Model-based outdoor performance capture. In *International Conference on 3D Vision (3DV)*, 2016.
- [46] Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, and Andrew Blake. Real-time human pose recognition in parts from single depth images. In *Computer Vision and Pattern Recognition (CVPR)*, 2011.
- [47] Leonid Sigal, Alexandru O. Bălan, and Michael J. Black. HumanEva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision (IJCV)*, 2010.
- [48] Leonid Sigal, Michael Isard, Horst Haussecker, and Michael J. Black. Loose-limbed people: Estimating 3D human pose and motion using non-parametric belief propagation. *International Journal of Computer Vision (IJCV)*, 98(1):15–48, 2012.
- [49] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multiview bootstrapping. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [50] Carsten Stoll, Nils Hasler, Juergen Gall, Hans-Peter Seidel, and Christian Theobalt. Fast articulated motion tracking using a sums of Gaussians body model. In *International Conference on Computer Vision (ICCV)*, 2011.
- [51] Christian Theobalt, Edilson de Aguiar, Carsten Stoll, Hans-Peter Seidel, and Sebastian Thrun. Performance capture from multi-view video. In *Image and Geometry Processing for 3-D Cinematography*, pages 127–149. Springer, 2010.
- [52] Matthew Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John Collomosse. Total capture: 3d human pose estimation fusing video and inertial sensors. In *Proceedings of 28th British Machine Vision Conference*, pages 1–13. University of Surrey, 2017.
- [53] Vicon Motion Capture Systems. <https://www.vicon.com/>, 2010.
- [54] Daniel Vlasic, Rolf Adelsberger, Giovanni Vannucci, John Barnwell, Markus Gross, Wojciech Matusik, and Jovan Popović. Practical motion capture in everyday surroundings. *ACM transactions on graphics (TOG)*, 26(3):35–es, 2007.
- [55] Timo Von Marcard, Gerard Pons-Moll, and Bodo Rosenhahn. Human pose estimation from video and imus. *IEEE transactions on pattern analysis and machine intelligence*, 38(8):1533–1547, 2016.
- [56] Timo Von Marcard, Bodo Rosenhahn, Michael J Black, and Gerard Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. In *Computer Graphics Forum*, volume 36, pages 349–360. Wiley Online Library, 2017.
- [57] Xiaolin Wei, Peizhao Zhang, and Jinxiang Chai. Accurate realtime full-body motion capture using a single depth camera. *SIGGRAPH Asia*, 31(6):188:1–12, 2012.
- [58] Xsens Technologies B.V. <https://www.xsens.com/>, 2011.
- [59] L. Xu, Z. Su, L. Han, T. Yu, Y. Liu, and L. FANG. Unstructuredfusion: Realtime 4d geometry and texture reconstruction using commercialrgbd cameras. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2019.
- [60] Lan Xu, Yebin Liu, Wei Cheng, Kaiwen Guo, Guyue Zhou, Qionghai Dai, and Lu Fang. Flycap: Markerless motion capture using multiple autonomous flying cameras. *IEEE Transactions on Visualization and Computer Graphics*, 24(8):2284–2297, Aug 2018.
- [61] Lan Xu, Weipeng Xu, Vladislav Golyanik, Marc Habermann, Lu Fang, and Christian Theobalt. Eventcap: Monocular 3d capture of high-speed human motions using an event camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

- [62] Weipeng Xu, Avishek Chatterjee, Michael Zollhöfer, Helge Rhodin, Dushyant Mehta, Hans-Peter Seidel, and Christian Theobalt. Monoperfcap: Human performance capture from monocular video. *ACM Transactions on Graphics (TOG)*, 37(2):27:1–27:15, 2018.
- [63] Xinyu Yi, Yuxiao Zhou, Marc Habermann, Soshi Shimada, Vladislav Golyanik, Christian Theobalt, and Feng Xu. Physical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial sensors. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022.
- [64] Xinyu Yi, Yuxiao Zhou, and Feng Xu. Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. 40(4), July 2021.
- [65] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.
- [66] Ye Yuan and Kris Kitani. Residual force control for agile human behavior imitation and extended motion synthesis. *Advances in Neural Information Processing Systems*, 33:21763–21774, 2020.
- [67] Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Neural descent for visual 3d human pose and shape. *arXiv preprint arXiv:2008.06910*, 2020.
- [68] Zhe Zhang, Chunyu Wang, Wenhui Qin, and Wenjun Zeng. Fusing wearable imus with multi-view images for human pose estimation: A geometric approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2200–2209, 2020.
- [69] Zerong Zheng, Tao Yu, Hao Li, Kaiwen Guo, Qionghai Dai, Lu Fang, and Yebin Liu. Hybridfusion: Real-time performance capture using a single depth sensor and sparse imus. In *European Conference on Computer Vision (ECCV)*, Sept 2018.
- [70] Xinge Zhu, Yuexin Ma, Tai Wang, Yan Xu, Jianping Shi, and Dahua Lin. Ssn: Shape signature networks for multi-class object detection from point clouds. In *European Conference on Computer Vision*, pages 581–597. Springer, 2020.
- [71] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Wei Li, Yuexin Ma, Hongsheng Li, Ruigang Yang, and Dahua Lin. Cylindrical and asymmetrical 3d convolution networks for lidar-based perception. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

## A Details of Data Collection

Recall that we collect our LIPD dataset to evaluate our LIP approach, which covers various challenging large-scale scenarios with real point-cloud and IMU measurements. In this section, we provide more details about our multi-sensor calibration (Sec. A.1) and our LIPD dataset (Sec. A.2). Our calibration scheme and the whole LIPD dataset will be made publicly available for future research about multi-sensor human motion capture.

### A.1 Multi-sensor System Calibration

**Preliminary.** Note that we attach the four IMU sensors  $s_i$  to the four limbs  $b_i$ , ( $i = 1, 2, 3, 4$ ) of the performer. Let  $\mathcal{F}_I, \mathcal{F}_L, \mathcal{F}_M, \mathcal{F}_{s_i}$  and  $\mathcal{F}_{b_i}$  denote the coordinate systems for the inertia, LiDAR, SMPL model and the  $i$ -th sensor and bone, respectively. Then, let  $\mathbf{R}_o^{f(p)} \in \mathbb{R}^{3 \times 3}$  and  $\mathbf{a}_o^f \in \mathbb{R}^3$  further denote the orientation in terms of rotation matrix and the free acceleration of object  $o$  in the coordinate system  $f$  when applied the skeletal pose  $p$ . Consistently, we use  $\mathbf{R}_{fg}$  to represent the rotation matrix between two coordinate systems  $f$  and  $g$ .

**Calibration Scheme.** Here, we introduce a efficient two-frame calibration scheme to obtain both the inertia-LiDAR calibration  $\mathbf{R}_{I2L}$  from  $\mathcal{F}_I$  to  $\mathcal{F}_L$ , as well as the sensor-bond calibration  $\mathbf{R}_{s_i 2b_i}$  from  $\mathcal{F}_{s_i}$  to  $\mathcal{F}_{b_i}$ . These calibration results are essential for unifying the original measurements from the sensors relative to  $\mathcal{F}_I$  into the attached bones' inertia relative to  $\mathcal{F}_L$ . During calibration, the performer keeps still in both A-pose and T-pose for several seconds, so as to obtain the known rotation pairs  $(\mathbf{R}_{b_i}^{M(A)}, \mathbf{R}_{s_i}^{I(A)})$  and  $(\mathbf{R}_{b_i}^{M(T)}, \mathbf{R}_{s_i}^{I(T)})$ . Then, to calculate  $\mathbf{R}_{s_i 2b_i}$ , we formulate its equivalence for both A-pose and T-pose during the whole motion sequence as follows:

$$\mathbf{R}_{s_i 2b_i}^{(A)} = \mathbf{R}_{s_i 2b_i}^{(T)}, \quad (10)$$

which is equivalent to

$$(\mathbf{R}_{b_i}^{M(A)})^{-1}(\mathbf{R}_{M2I})^{-1}\mathbf{R}_{s_i}^{I(A)} = (\mathbf{R}_{b_i}^{M(T)})^{-1}(\mathbf{R}_{M2I})^{-1}\mathbf{R}_{s_i}^{I(T)}, \quad (11)$$

where  $\mathbf{R}_{M2I}$  is an unknown rotation offset from  $\mathcal{F}_M$  to  $\mathcal{F}_I$ . Therefore, we optimize  $\mathbf{R}_{M2I}$  by solving a non-linear least squares problem iteratively, and so that  $\mathbf{R}_{s_i 2b_i}$  can be determined according to any side of equation 11.

Besides, after optimizing  $\mathbf{R}_{M2I}$ , we obtain the calibration result  $\mathbf{R}_{I2L}$  from  $\mathcal{F}_I$  to  $\mathcal{F}_L$  as follows:

$$\mathbf{R}_{I2L} = \mathbf{R}_{M2L}(\mathbf{R}_{M2I})^{-1}. \quad (12)$$

Here,  $\mathbf{R}_{M2L}$  is a predefined rotation offset from  $\mathcal{F}_M$  to  $\mathcal{F}_L$ , since we assume that the performer facing the LiDAR also stands horizontally and vertically aligned during calibration.

After the above calibration, we transform the original IMU measurements  $\mathbf{R}_{s_i}^I, \mathbf{a}_{s_i}^I$  into  $\mathbf{R}_{b_i}^L, \mathbf{a}_{b_i}^L$ :

$$\begin{aligned} \mathbf{R}_{b_i}^L &= \mathbf{R}_{I2L} \mathbf{R}_{s_i}^I (\mathbf{R}_{s_i 2b_i})^{-1}, \\ \mathbf{a}_{b_i}^L &= \mathbf{R}_{I2L} \mathbf{a}_{s_i}^I, \end{aligned} \quad (13)$$

which are fed into our hierarchical pose estimation network.

### A.2 Design of LIPD Collection

As listed in Table. 5, our LIPD includes 10 motion sequences of 6 performers (4 males and 2 females) covering various motion genres in large-scale scenarios, with the multi-modal LiDAR and Inertial observations. Specifically, we collect various sports sequences like ‘‘football’’, ‘‘taekwondo’’ or ‘‘long-jump’’, which obtain challenging, fast and large-scene movements. We further capture those challenging sequences with complicated terrain or poor lighting and the one ‘‘pas de deux’’ with multiple performers. These sequences in LIPD serve as a strong supplement to existing dataset, which points out the limitation of purely inertial or RGB-based solutions for capturing large-scene motions. By introducing the dataset and our LIP algorithm on top of it, we would like to push the technical boundary of human motion capture into the realm of large-scale and multi-modal setting.

Table 5: Metadata of LIPD. We report specific motion genres, gender of the performer and characteristics of each motion sequence in our dataset.

motion genre	gender of performer	characteristics
football	male	
basketball	male	
taekwondo	female	
running	male	challenging and fast; large-scale movements
boxing	male	
dancing	female	
long-jump	male	
climbing	male	extremely complicated terrain
wotagei	male	extremely poor lighting
pas de deux	males	multi-person

## B Data Synthesis

In order to ensure the diversity of training data, we synthesize point cloud data B.1 or IMU measurements B.2 on public datasets[18, 28, 29, 34], which do not contain both modalities as input. In this section, we will explain our detailed implementations of data synthesis.

### B.1 Point Cloud Synthesis

LiDAR works in a time-of-flight way with simple principles of physics, which can be easily simulated with a small gap to real data. We generate simulated LiDAR point cloud by emitting regular lights from the LiDAR center according to its horizontal resolution and vertical resolution. The light can be reflected back when encountering obstacles, generating a point at the intersection on the surface of obstacles. We conduct the simulation according to the parameters of Ouster(OS1-128), which is also the device used in collecting LIPD. Its horizontal resolution is 2048 and its vertical resolution is 128 lines. Each emission direction is described by unit vector in spherical coordinate system  $d = [\cos \varphi \sin \theta, \cos \varphi \cos \theta, \sin \varphi]$ , where  $\varphi$  represents the angle between the emission direction and the plane XY,  $\theta$  indicates the azimuth, and  $c = [0, 0, 2]$  is the LiDAR center. The intersection point  $p = [p_x, p_y, p_z]$  is calculated by

$$p = c + d \frac{n^T(q - c)}{n^T d}, \quad (14)$$

where  $n$  represents the normal vector of corresponding mesh and  $q$  denotes any vertex point of the mesh.

For the process of calculating the intersection of LiDAR and mesh surface, there are mainly three steps. First is to calculate the intersection of light and triangular patch, and the second is to judge whether the intersection of light and plane is inside the triangle. Due to occlusions, one light should only have the intersection with the first touched mesh of object. Finally, we filter the intersections to only keep the ones first occurred in the LiDAR view.

### B.2 IMU Measurement Synthesis

We simulate the magnetometer and accelerometer of IMU to synthesize rotation and acceleration data. To do so, we virtually attach four imaginary IMUs to specific vertices of the SMPL mesh corresponding to left wrist, right wrist, left ankle and right ankle with perfectly aligned sensor-bone coordinate frames. Then, the rotation of each IMU sensor can be calculated through operating forward kinematics(FK) on ground-truth SMPL pose parameters, and the acceleration can be approximated by formulation 15

$$\mathbf{a}_{s_i}(t) = \frac{\mathbf{p}_{s_i}(t - s) + \mathbf{p}_{s_i}(t + s) - 2\mathbf{p}_{s_i}(t)}{(s\Delta t)^2}, \quad i = 1, 2, 3, 4, \quad (15)$$

where  $\mathbf{p}_{s_i}(t) \in \mathbb{R}^3$  denotes the absolute position of imaginary IMU sensor  $s_i$  at time frame  $t$  and  $\Delta t$  means the time interval between two consecutive frames, which is  $0.1s$  in our implementation. Meanwhile, we refer to TransPose to set a smooth factor  $s = 4$  for filtering the jitters.

Since the synthesized acceleration data is actually free acceleration, we transform real IMU acceleration  $\tilde{\mathbf{a}}_{s_i}(t) \in \mathbb{R}^3$  recorded in LIPD, which is relative to sensor coordinate frame with gravity  $\mathbf{g} = [0, 0, 9.81]$ , to free acceleration form by formulation 16

$$\mathbf{a}_{s_i}(t) = \mathbf{R}_{s_i}(t)\tilde{\mathbf{a}}_{s_i}(t) - \mathbf{g}, \quad (16)$$

where  $\mathbf{R}_{s_i}(t) \in \mathbb{R}^{3 \times 3}$  is the raw IMU rotation recording at time frame  $t$ .

## C More Results of LIP

In this section, we provide more results of LIP on diverse challenging sequences, which illustrates the effectiveness of our approach for accurate skeletal pose and global translation estimation in more general scenarios.

### C.1 Extreme Environment-lighting

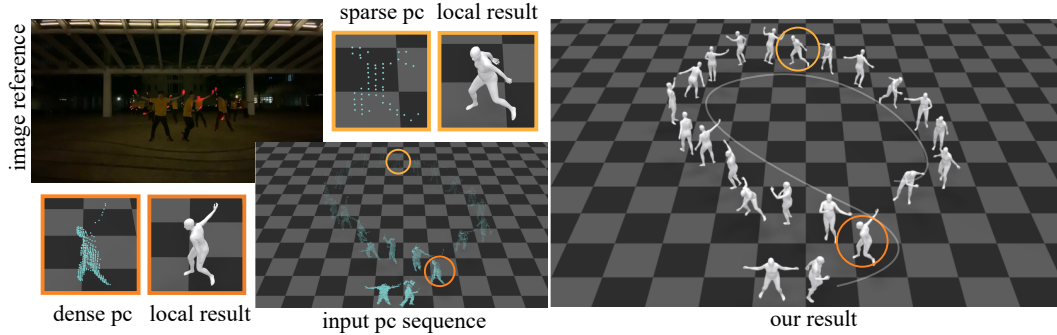


Figure 7: Capture result of our method on a night activity motion sequence “wotagei” (cyalume dance). Benefiting from lighting-independent sensors, our method can robustly capture the global motions of the performer even with extremely poor lighting.

In Figure 7 we demonstrate our capture results where the actor is performing wotagei in the night. Such scenario of large capturing range and poor lighting conditions is extremely challenging for traditional RGB-based mocap methods. While it is even already difficult for human eyes to spot the actor in the reference images, our approach yields high-quality results, thanks to the lighting-independent and 3D-aware property of our multi-sensor input. Thus, our LIP brings huge potential to capture those motions with flick or dark lighting, such as stage performance or night activities.

### C.2 Complicated Terrain

In Figure 8, we provide our capture results for handling various terrain where the actor is moving around a long stairs. Previous purely RGB or IMU-based mocap methods will suffer from global localization artifacts. Differently, thanks to the global 3D-aware property of our multi-modal input, our LIP approach can reconstruct the global skeletal motions of the performer in a wide 3D space. It demonstrates the huge potential of LIP to be further combined with those static reconstruction approaches for human-scene modeling.

### C.3 Multiple Persons

We also explore the potential of LIP in the application of multi-person mocap task. Figure 9 shows that LIP can naturally extend to multi-person capture, and works well when few human interactions and occlusions happen in the motion sequence. Whereas, LIP cannot make accurate inference when the point cloud of human body is not complete due to self- or external-occlusion, which can be potentially addressed by data augmentation for partial point cloud data and taking advantage of RNN

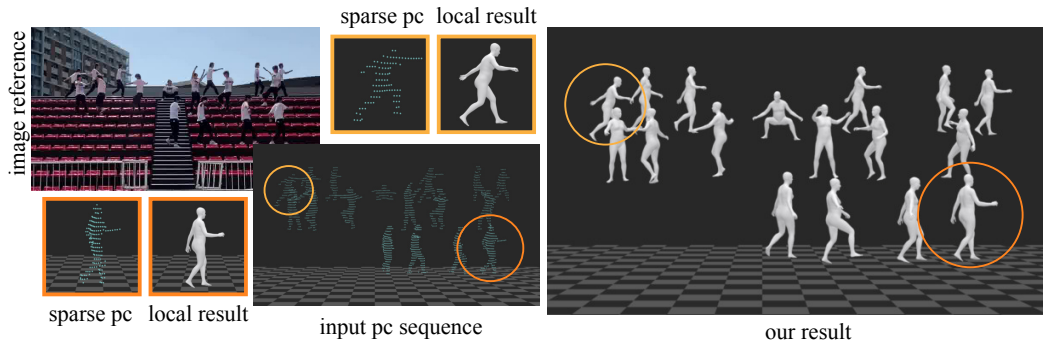


Figure 8: Capture result of our method on the “climbing” sequence. Thanks to the 3D-aware property of our approach, we faithfully recover the performer’s motions when he is climbing across the stairs.

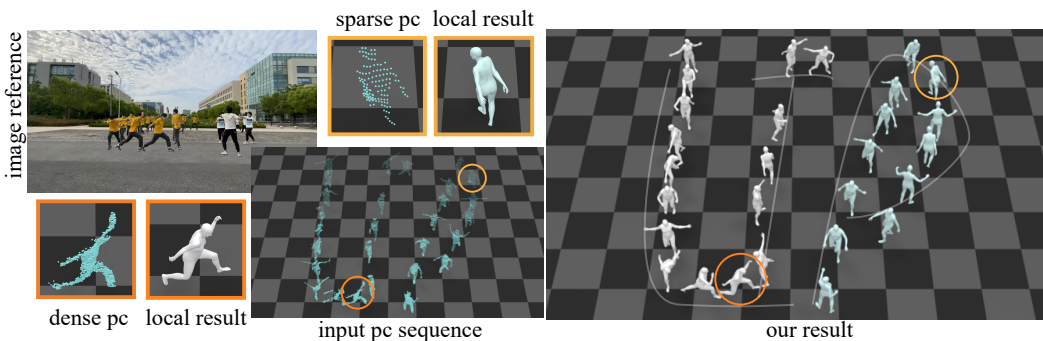


Figure 9: Capture result of our method on a multi-person motion sequence “pas de deux”. Two performers act and dance freely in a large area and finally “draw” the abbreviation “LIP” using their global motion trajectories.

or Transformer to enhance the prediction capability of the network. It deserves further efforts to provide more effective and robust solutions to handle more challenging situations in the future.

## D More Experiments

In this section, we first provide more qualitative results of comparison, and then show more detailed ablation study of our network modules and diverse input configurations, and finally offers detailed analysis for the generalization capability of our method.

### D.1 More Qualitative Comparisons

We show the comparison result in the climbing status in Figure 10. Benefiting from global localization information provided by LiDAR and precise inertia measurements, our method can track complicated 3D trajectory not only in horizontal but also in vertical very well and also maintain accurate local pose estimation. Figure 12 shows the visualization of comparison results of local pose estimation on the testing data of synthetic datasets. Our method is obviously superior to others with results more close to the ground truth.

### D.2 Detailed Ablation Study of Network Design

In the main body of the submission, we report the weighted mean of each evaluation metric, while in Table 6 we reports detailed evaluation results of our network design on different datasets, both for *Multi-modal Pose Estimation* and *Pose-guided Translation Correction*.

The subtable on the top demonstrates that *Global Pose-prior Regressor* is imperative for multi-modal input fusion, without which the implicit motion pattern representation, point cloud feature, cannot combine with explicit motion inertia effectively. Furthermore, *Joint-map Estimator* is proved

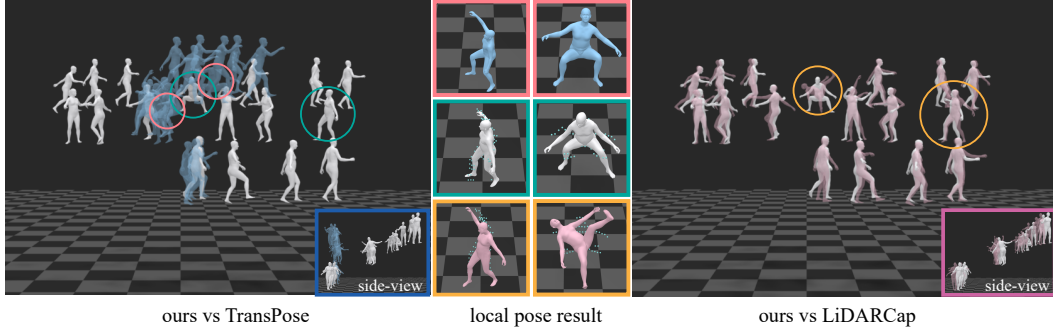


Figure 10: Qualitative comparison of capture performance on the “climbing” sequence with extremely complicated terrain.

Table 6: Ablation study of network design. We evaluate the importance of *Global Pose-prior Regressor* and *Joint-map Estimator* block for local pose estimation. Meanwhile, the benefits of pose-guiding information(“PG”) in translation estimation is evaluated together with average translation estimation and direct regression, using only CD metric.

Local pose	TotalCapture			DIP-IMU			LiDARHuman26M		
	MPJPE	Mesh Err	Ang Err	MPJPE	Mesh Err	Ang Err	MPJPE	Mesh Err	Ang Err
Ours w/o RE	55.6	68.1	9.96	58.6	70.8	11.75	67.9	77.0	12.73
Ours w/o JE	37.2	47.8	8.56	36.5	45.8	10.03	60.7	72.2	11.42
Ours	<b>28.6</b>	<b>37.3</b>	<b>6.81</b>	<b>29.7</b>	<b>38.3</b>	<b>9.35</b>	<b>58.2</b>	<b>68.9</b>	<b>11.26</b>

Global translation	TotalCapture	LiDARHuman26M	LIPD
Average estimation	6.4	8.1	7.9
ICP estimation	2.0	4.7	4.1
Ours w/o PG	0.9	3.7	3.4
Ours	<b>0.7</b>	<b>3.5</b>	<b>3.2</b>

important for a hierarchical pose estimation pipeline, which refines the human body joint positions with IMU measurements so that the pose solved from that is more accurate.

The subtable at the bottom indicates that *Pose-guided Translation Correction* outperforms other naive translation estimation methods, such as taking ICP result or arithmetic mean of raw point cloud as the translation estimation, with a significant margin. Moreover, results show that the pose-guiding information can further benefit the accuracy of trajectory tracking.

### D.3 Detailed Ablation Study of Input Configuration

Table 7: Ablation study for input configuration. We evaluate our choice of single LiDAR with 4 IMUs by experimenting various combinations of input modalities on our evaluation dataset.

	TotalCapture			DIP-IMU			Lidarhuman26M		
	MPJPE	Mesh Err	Ang Err	MPJPE	Mesh Err	Ang Err	MPJPE	Mesh Err	Ang Err
LiDAR Only	30.4	41.2	12.10	30.9	41.5	14.87	65.5	81.5	19.42
4 IMUs Only	73.6	90.3	10.12	90.3	100.1	13.81	89.7	103.7	13.60
5 IMUs only	56.3	70.2	8.66	79.0	88.8	11.96	70.1	82.6	12.29
LiDAR+2 IMUs	29.1	38.4	8.49	30.1	39.7	11.04	61.7	75.1	14.01
LiDAR+4 IMUs(Ours)	<b>28.6</b>	<b>37.3</b>	<b>6.81</b>	<b>29.7</b>	<b>38.3</b>	<b>9.35</b>	<b>58.2</b>	<b>68.9</b>	<b>11.26</b>
LiDAR+5 IMUs	27.7	36.8	6.73	29.0	37.8	9.17	55.4	66.0	10.47
LiDAR+12 IMUs	21.8	27.3	3.63	30.7	40.1	5.76	49.0	59.3	7.45

In Table 7, the performance of varied combinations of multi-modal input on each evaluation dataset are reported. Consistently, our configuration(LiDAR+4 IMUs) outperforms any single-modal system and achieves comparable accuracy with dense sensors configuration(LiDAR+12 IMUs), demonstrating that our configurations balance the high-quality performance and light-weight device.



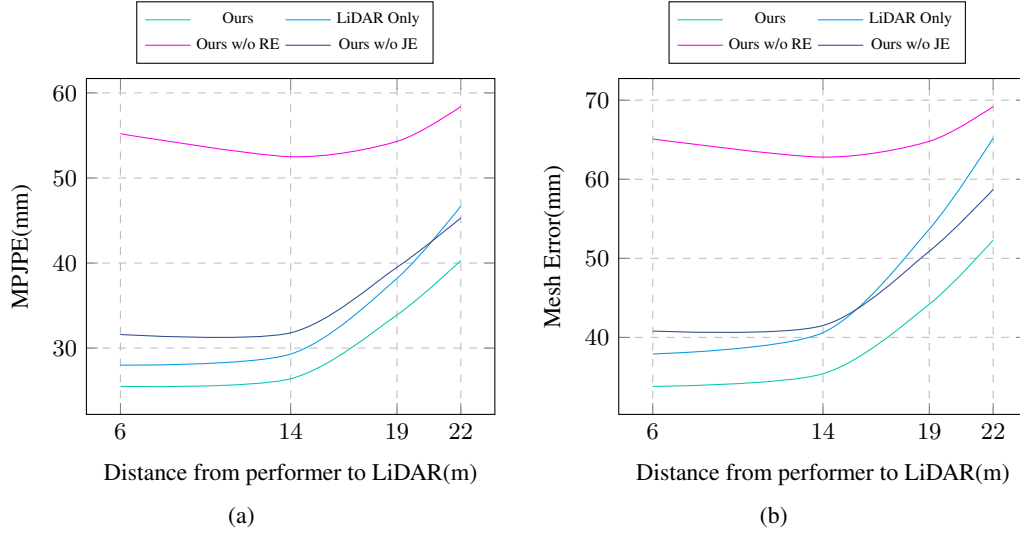


Figure 11: Experiment for generalization capability. We evaluate LIP and LiDAR-based methods in different distance under the same motion. The metrics of MPJPE and mesh error can directly reflect the comparison

#### D.4 Generalization Capability

To evaluate the generalization capability of LIP to point clouds captured in various ranges with different sparsity, we simulate a bunch of point clouds from TotalCapture dataset at multiple virtual capture distances from 6 to 25 meters. Figure 11 illustrates that our LIP performs consistently well in a wide range about 20 meters. Meanwhile, various curves in Figure 11 further demonstrate that both our multimodal input configuration and network design contribute to the generalization capability of our method.



Figure 12: Qualitative comparison of capture performance on synthetic data. Our method outperforms state-of-the-art works by more accurate local pose estimation on three different datasets. Note that LiDARHuman26M dataset is not applicable for TransPose since the synthetic IMU measurement from LiDARHuman26M can only be 10fps, while frame rate of 60fps is required by TransPose.