

Physical Inertial Poser (PIP): Physics-aware Real-time Human Motion Tracking from Sparse Inertial Sensors

Xinyu Yi¹

Yuxiao Zhou¹

Marc Habermann²

Soshi Shimada²

Vladislav Golyanik²

Christian Theobalt²

Feng Xu¹

¹School of Software and BNRist, Tsinghua University

²Max Planck Institute for Informatics, Saarland Informatics Campus

Abstract

Motion capture from sparse inertial sensors has shown great potential compared to image-based approaches since occlusions do not lead to a reduced tracking quality and the recording space is not restricted to be within the viewing frustum of the camera. However, capturing the motion and global position only from a sparse set of inertial sensors is inherently ambiguous and challenging. In consequence, recent state-of-the-art methods can barely handle very long period motions, and unrealistic artifacts are common due to the unawareness of physical constraints. To this end, we present the first method which combines a neural kinematics estimator and a physics-aware motion optimizer to track body motions with only 6 inertial sensors. The kinematics module first regresses the motion status as a reference, and then the physics module refines the motion to satisfy the physical constraints. Experiments demonstrate a clear improvement over the state of the art in terms of capture accuracy, temporal stability, and physical correctness.

1. Introduction

Capturing the motion of real humans is a long-standing and challenging problem with many applications in computer vision and graphics, movie production, gaming, AR, and VR. However, due to its articulated structure, capturing the highly complex and potentially fast movements of the human body is challenging and many works have been proposed in the past [4, 69, 72, 73, 77, 92].

One category of approaches are image-based where the actor motion is recovered by analyzing the image data,

This work was supported by Beijing Natural Science Foundation (JQ19015), the NSFC (No.61727808, 62021002), the National Key R&D Program of China 2018YFA0704000. This work was supported by THUICS, Tsinghua University and BLBCI, Beijing Municipal Education Commission. This work was partially supported by the ERC consolidator grant 4DReply (770784). We thank Notiom [40] for the extensive support on inertial sensors, and Liuqing Yang, Liangdi Ma, Siyuan Teng, Wenbin Lin for the help on live demos. Feng Xu is the corresponding author.

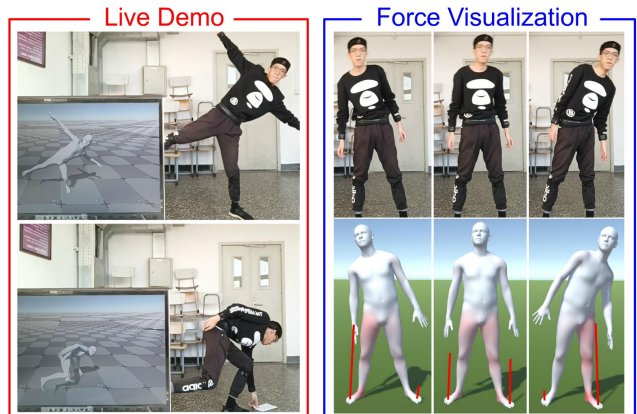


Figure 1. PIP captures physically correct human motion, joint torques, and ground reaction forces solely from a sparse set of six IMUs. Importantly, PIP runs at 60 frames per second with only 16ms latency, which enables real-time applications.

which can be either multi-view imagery [4, 6, 61, 86], depth images [75–77, 88], or a single RGB stream [5, 8, 15, 25, 28, 38, 84, 92]. From the setting, it becomes clear that occlusions (either object-actor or self occlusions) can lead to a significantly reduced tracking quality. Besides, these methods are sensitive to the lighting and the appearance of the actor as distinct features need to be extracted from images. Moreover, many methods assume a static camera, resulting in a limited space where the subject can be captured. These drawbacks limit the usability of optical motion capture.

Recently, researchers start to explore alternative sensing devices such as inertial measurement units (IMUs). Production-ready solutions [40, 71] can track the body motion accurately solely from inertial sensors. However, they rely on special suits with densely placed sensors (usually 17 IMUs), which are difficult to wear. Besides, the large number of IMUs can hinder the actor’s movement. Having a sparser set of IMUs on the body is clearly advantageous and more flexible. However, recent sparse methods [20, 66, 73] struggle with physical correctness and cannot disambiguate poses with similar sensor measurements such as sitting and standing; they are non-causal, *i.e.*, need future information,

which introduces large delays; their accuracy is still limited while temporal artifacts such as jitter become visible.

To this end, we propose Physical Inertial Poser (PIP), a new real-time method for motion capture as well as joint torque and ground reaction force estimation using only six IMUs (see Fig. 1). In contrast to previous works [20, 73] that require future information, our method only requires the information already available at any given time, which means no additional delay is introduced. Our algorithm has two stages: 1) learning-based motion estimation and 2) physics-based motion optimization, which leverage both human kinematics and dynamics in motion capture.

In the estimation stage, we regress human pose, joint velocities, and foot-ground contact probabilities from the inertia inputs using recurrent neural networks (RNNs). We estimate leaf-to-full joint positions as intermediate tasks to improve the tracking accuracy as proposed by TransPose [73]. To resolve the pose ambiguity arising from the sparse IMU placement, we further propose a *learning-based* RNN state initialization strategy, which helps the networks better learn the change of body pose from input inertia measurements. This results in a significant accuracy improvement especially for ambiguous motions such as sitting still.

In the optimization stage, we recover the physically correct motion, joint torques, and ground reaction forces from the kinematic estimations, leveraging a torque-controlled floating-base simulated character model. Different from previous works that independently control the rotation of each degree of freedom of the character using proportional-derivative (PD) rules [21, 54, 55, 80], we propose a novel *dual PD controller* to incorporate the global holistic control of the character’s pose. This is achieved by applying PD rules on both joint positions and rotations. The proposed technique significantly improves the translation accuracy and physical plausibility of the motion.

In summary, our main contributions are:

- The first physics-aware real-time approach that estimates human motion, joint torques, and ground reaction forces with only six IMUs, which we call PIP.
- A learning-based RNN state initialization scheme, which helps to better disambiguate human motion regression from sparse IMU measurements (Sec. 3.1).
- A dual PD controller, which achieves the combined control of local and global pose to improve the motion tracking accuracy and physical plausibility (Sec. 3.2).

Our experiments demonstrate that PIP significantly outperforms previous sparse IMU-based methods in terms of tracking accuracy, physical plausibility, and disambiguation of challenging poses.

2. Related Work

Human motion capture (mocap) has a long research history. Many works have been devoted to this topic, which

can be mainly categorized into optical, inertial, and hybrid approaches. Since our method only requires IMU measurements as input, we do not discuss purely image-based approaches [15, 16, 24, 27, 39, 48, 58]. Here, we focus on hybrid and inertial mocap solutions, and the previous efforts on the physical plausibility of human motion.

Optical-inertial Hybrid Motion Capture. As image-based mocap solutions suffer from occlusions, fusing images with IMUs, which aims at achieving more robust motion tracking, has recently attracted much attention. This can be achieved by either energy-based optimization [22, 35–37, 43, 65] which optimizes human pose to fit both image features and inertia measurements, or feature-based estimation [13, 62] which regresses human pose from the combined features derived from images and IMUs. Zhang et al. [87] propose to exploit IMUs in the 2D pose estimation by fusing the image features of each pair of joints linked by the IMUs. Some works fuse IMUs with depth images [18, 23, 90] or optical markers [1] to perform human motion/performance capture. Nevertheless, these methods are still substantially limited under low light conditions and heavy occlusions, and require the actor to move within the viewing frustum of the camera. Our method requires no visual input, and thus is free from these limitations.

Motion Capture from Inertial Sensors. Inertial mocap approaches do not suffer from occlusions or restricted moving space. Commercial solutions [40, 71] and the extended work [14] rely on 17 IMUs to perform motion capture. They usually require the actor to wear a tight suit with densely bounded IMUs, which is inconvenient, intrusive, and obstructive. It is clear that having a reduced set of IMUs on the body is preferable. However, motion capture from sparse inertial sensors is very ambiguous and challenging. Some works [31, 59, 64] leverage ultrasonic sensors for additional position information to resolve some ambiguities, but the use of distance sensors limits the recording range. Early purely-inertial works [47, 57, 60] use sparse accelerometers to reconstruct human pose by database search. Schwarz et al. [50] use sparse orientation measurements to perform person-specific pose estimation. To improve the accuracy, recent works [12, 20, 44, 66, 73] leverage both acceleration and orientation measurements. Marcard et al. [66] present an offline method for human motion capture from only 6 IMUs, which achieves promising accuracy. Huang et al. [20] propose the first deep learning method, which uses a bidirectional recurrent neural network (biRNN) to estimate the human pose from 6 IMUs in real-time. However, their method does not allow to locate the person in the 3D space, *i.e.*, the root translation is not estimated. The current state-of-the-art method, TransPose [73], introduces the first real-time pose and translation estimation framework, which achieves an accurate capture quality while also only using 6 IMUs. However, all of these works have a non-negligible

delay due to the inherent need of future information, cannot stably capture ambiguous poses, and has many non-physical artifacts such as jitter and foot-sliding. In contrast, our method does not rely on any future information and is free from the delay while even achieving higher accuracy. In addition, we are the first to combine physics-based motion optimization with sparse inertial motion capture, and we show that such a carefully orchestrated design significantly improves the physical correctness of the motion.

Physical Plausibility of Motion. To ensure the physical plausibility of motion, many works address the awareness of physics in their approach. One category of works only impose physical constraints (*e.g.*, foot contacts [7, 53, 93], temporal consistency [16, 39, 66], and collision [81]) without considering human dynamics such as forces and masses. Some works leverage an explicitly reconstructed scene to constrain the motion [14, 17, 85]. However, due to the articulated structure of humans, it is considerably difficult to track the complicated body movements by imposing such naive constraints. Another category of works leverage physics-based human models and estimate forces to control the motion, targeting a more accurate modeling of real-world human movements. Some works [30, 46, 55, 67, 70, 83] use optimization-based methods to solve the optimal forces and human motion, which satisfy the physical constraints and laws such as the equation of motion [9]. Zell et al. [82] propose a weakly-supervised learning framework for dynamics estimation from human motion. Shimada et al. [54] present a fully-differentiable framework for learning-based motion and force estimation from videos. Reinforcement learning is also used in physics-based character control [2, 21, 32, 41, 42, 74, 78, 80], which can utilize advanced non-differentiable physics simulators. Among these works, our physics module is most similar to the work of Shimada et al. [55]. The major differences are the input to the respective method and the control of the physical character. Our method assumes sparse inertia measurements of the moving body as input while theirs [55] leverages images of the actor. Moreover, their approach [55] uses a proportional-derivative (PD) controller to control the rotation of each joint of the character independently. In contrast, our method uses a novel *dual PD controller* to introduce the global control of the character, aiming at better accuracy. In other words, our method is the first that leverages explicit physics-based optimization into sparse IMU-based motion capture.

3. Method

Our task is to track human motion in real-time using 6 IMUs. The input of our method is the sequential measurements of accelerations and orientations of the 6 IMUs mounted on the left/right forearms, left/right lower legs, head, and pelvis (Fig. 2). The output of our method is the subject’s motion in terms of joint angles and global trans-

lation, together with physical properties including ground reaction forces and joint torques. The method incorporates two modules: 1) *the kinematics module*: a neural kinematics estimator, which infers the human motion from the IMU measurements, followed by 2) *the dynamics module*: a physics-aware motion optimizer, which refines the human motion and outputs the physical properties.

3.1. Neural Kinematics Estimator

The task of the kinematics module is to estimate the current motion status (specified later). We use the same kinematic tree as in SMPL [33], which contains $J = 24$ joints. We refer to the wrists, ankles, and head as the *leaf joints*, and the pelvis as the *root joint*. For all the J joints, their 3D positions are denoted as $\mathbf{p} \in \mathbb{R}^{3J}$; their linear velocities are denoted as $\mathbf{v} \in \mathbb{R}^{3J}$; and their rotations are denoted as $\varphi \in \mathbb{R}^{6J}$ in the 6D representation [91]. Since the IMUs do not provide any positional measurement, all these estimations are in the local coordinate frame (relative to the root joint). Similar to TransPose [73], we perform a T-Pose calibration at the beginning, and then at each time step we stack the IMU measurements, *i.e.*, the aligned accelerations and rotation matrices, into a single input vector $\mathbf{x} \in \mathbb{R}^{72}$. In the following, we first give an overview of our network structure (Sec. 3.1.1). Then, we dive into our novel learning-based initialization for the RNN hidden states during training (Sec. 3.1.2). The new initialization method helps the network learn to capture the state-change signals, which is crucial to resolve the pose ambiguity in our task.

3.1.1 Motion Estimation Network

The structure of the network follows the one of TransPose [73], which first estimates leaf joint properties and then full body status in a multi-stage style. Different from their method [73], we choose to use RNN instead of biRNN as the basic network structure (elaborated in Sec. 3.1.2), and we estimate the velocity of *all joints* instead of *only the root joint*, since we find that using full-joint velocities in combination with the physics part allows better character control.

Specifically, as shown in Fig. 2, we first use an RNN P_L to regress leaf joint positions $\mathbf{p}_{\text{leaf}} \in \mathbb{R}^{15}$ from the IMU measurements \mathbf{x} . Then, the concatenated vector $[\mathbf{p}_{\text{leaf}} \ \mathbf{x}]$ is fed into the second RNN P_A , which estimates all joint positions $\mathbf{p} \in \mathbb{R}^{3J}$. Next, we feed the vector $[\mathbf{p} \ \mathbf{x}]$ into three RNNs R_A , V_A , and C_F to estimate the joint rotations $\varphi \in \mathbb{R}^{6J}$ (the root orientation is directly measured by the IMU placed on the pelvis), linear velocities $\mathbf{v} \in \mathbb{R}^{3J}$, and foot-ground contact probabilities $\mathbf{c} \in \mathbb{R}^2$. Finally, φ , \mathbf{v} , and \mathbf{c} , which we call *motion status*, are fed into the subsequent dynamics module. During training, we use an L2 loss for P_L , P_A , and R_A , a binary cross-entropy loss for C_F , and the cumulative loss proposed in [73] for V_A .

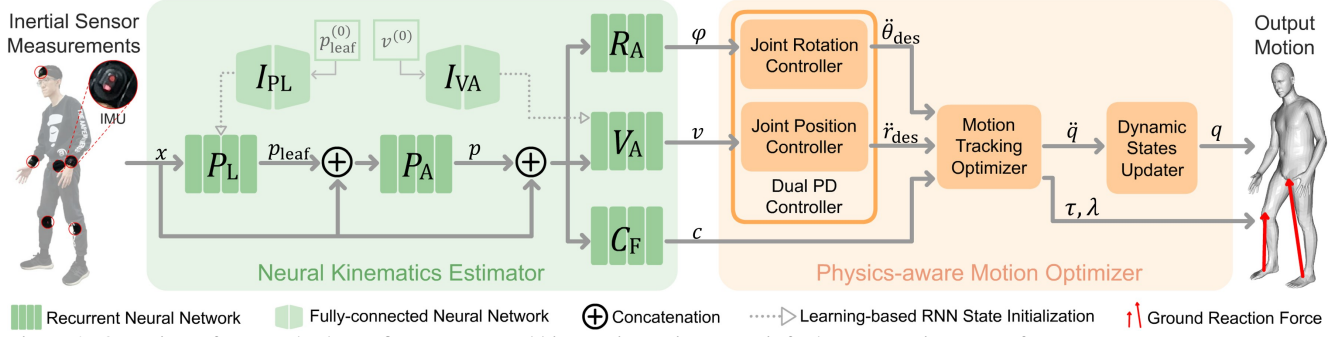


Figure 2. Overview of our method. We first use a neural kinematics estimator to infer human motion status from sparse IMU measurements. Then, we use a physics-aware motion optimizer to obtain physically correct human motion, joint torques, and ground reaction forces.

3.1.2 RNN with Learning-based Initialization

Full-body motion tracking from sparse inertial sensors is severely under-constrained and ambiguous. For example, due to the sparsity of the IMUs, it is impossible to distinguish *standing still* and *sitting still* since the IMU measurements are identical: orientations are the same and accelerations are zero. To cope with this ambiguity, leveraging the temporal information by capturing and memorizing the state-change signals in historical frames is a must. Previous state-of-the-art works [20, 73] leverage bidirectional recurrent neural networks (biRNN) [49] to learn such temporal information. However, the design of biRNNs only allows a fixed frame window in the real-time setting, which prevents the access to state changes that happened outside this temporal window, *e.g.*, when the subject remains seated for a longer time. In consequence, these methods fail to capture such ambiguous poses correctly. To overcome this limitation, we use RNNs to retain complete historical information and capture the crucial state-change signals.

To capture the state-change signals, not only does the architecture have to be updated, but also a new training strategy is necessary. Traditionally, an RNN is trained in a mini-batch manner and always starts with a zero initialization for hidden states in each batch. However, in our setting, a constant initial state is incorrect (the subject may start from sitting, standing, lying, *etc.*); and when the initial state is wrong, the model can never learn how to change its hidden state according to the signals afterward due to the mismatch at the beginning. To address this problem, we propose a learning-based RNN initialization strategy. Specifically, we have a separate fully-connected neural network (FCN), which regresses the initial state of the RNN from body pose information. The FCN and RNN are trained jointly: for each mini-batch, the ground-truth pose at the beginning is fed into the FCN, then the output of the FCN is assigned to the hidden state of the RNN, then the RNN is trained as usual. As the RNN implementation is not modified, the proposed strategy is highly effective and compatible with highly optimized black-box RNN libraries. Dur-

ing inference, we assume the initial pose of the subject is known, which can be obtained from the calibration step. Notice that the FCN only initializes the RNN for the first frame. The FCNs we use for initialization during training are shown in Fig. 2 as I_{PL} and I_{VA} , which take the beginning leaf joint positions $p_{\text{leaf}}^{(0)}$ and joint velocities $v^{(0)}$ as input, respectively. This initialization is only applied to P_L and V_A , which suffer most from the ambiguity.

3.2. Physics-aware Motion Optimizer

The output of the kinematics module may still contain artifacts like jitter and ground penetration. We therefore introduce the dynamics module to explicitly apply the physical constraints as similar to [55]. The input to this module is the motion status φ , v , and c estimated by the kinematics module, which serve as the *reference* in the physics-based optimization. The task of the dynamics module is to obtain the motion, internal joint torques, and ground reaction forces that align with the reference but also satisfy physical constraints. Specifically, based on the physics model (Sec. 3.2.1), we first use a novel *dual PD controller* (Sec. 3.2.2) to compute the desired acceleration for the simulated character which can fully reproduce the reference motion, and then use a motion optimizer (Sec. 3.2.3) to solve for the acceleration and force that the character can actually produce within the physical constraints. Finally, we update the character status and compute the final output motion (Sec. 3.2.4).

3.2.1 Physics Model

We use a torque-controlled floating-base simulated character [89] as our physics model and follow the same mass distribution as in [55]. We initialize the subject’s global position at the origin. We refer to the joint positions in the global coordinate frame as $r \in \mathbb{R}^{3J}$ and the translation as $r_{\text{root}} \in \mathbb{R}^3$. The time derivative \dot{r} and \ddot{r} refer to the linear velocity and acceleration in the global frame. We refer to the local joint rotations (*i.e.*, pose) in Euler angles as $\theta \in \mathbb{R}^{3J}$, and its time derivative $\dot{\theta}$ and $\ddot{\theta}$ are the angular ve-

locity and acceleration. The configuration of the character is described by its pose and translation, which we denote as $\mathbf{q} = [\mathbf{r}_{\text{root}} \ \boldsymbol{\theta}] \in \mathbb{R}^N$ where $N = 3 + 3J$ is the degree of freedom (DoF). The time derivative $\dot{\mathbf{q}}$ and $\ddot{\mathbf{q}}$ are the generalized velocity and acceleration. The character is controlled by the vector of force $\boldsymbol{\tau} \in \mathbb{R}^N$ where each dimension refers to the force on the corresponding DoF. In the real world, the character is actuated only by the torques at the non-root joints, while no force is applied to the root joint. However, to compensate for the dynamics mismatch between our physics model and real humans, we allow a small residual force at the root joint as prior works [29, 54, 55, 79, 80] do. In our notation, the first six entries $\boldsymbol{\tau}_{:6}$ correspond to the residual force at the root joint, and $\boldsymbol{\tau}_6$ are the actuated joint torques. The generalized acceleration $\ddot{\mathbf{q}}$ and the force $\boldsymbol{\tau}$ follow the equation of motion [9]:

$$\boldsymbol{\tau} + \mathbf{J}_c(\mathbf{q})^T \boldsymbol{\lambda} = \mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{h}(\mathbf{q}, \dot{\mathbf{q}}), \quad (1)$$

where $\mathbf{M} \in \mathbb{R}^{N \times N}$ is the inertia matrix; $\mathbf{h} \in \mathbb{R}^N$ is the non-linear effect term that accounts for gravity, Coriolis, and centripetal forces; $\boldsymbol{\lambda} \in \mathbb{R}^{3n_c}$ is the external contact forces applied at n_c character-ground contact points; $\mathbf{J}_c \in \mathbb{R}^{3n_c \times N}$ is the contact point Jacobian, which maps the generalized velocity $\dot{\mathbf{q}}$ to the contact point velocities:

$$\dot{\mathbf{r}}_c = \mathbf{J}_c \dot{\mathbf{q}}. \quad (2)$$

Readers are referred to [9] for more details. In our model, we assume all external forces (except for gravity) are the support and frictional forces exerted at the contact points by the ground, which we call *ground reaction force (GRF)*.

3.2.2 Dual Proportional-differential Controller

To control the character, previous methods [21, 54, 55, 80] use a single PD controller to compute either angular accelerations or joint torques to reproduce the reference motion. However, since the configuration of the character is parameterized in *local* Euler angles, such methods only focus on the independent control of local joint rotations, which may result in an undesirable global pose. Simply applying PD control on *global* joint rotations will make the optimization problem non-quadratic, introducing a large computation cost. We find that imposing a PD controller on joint positions will constrain the global pose, while still keeping the problem quadratic. To this end, we propose a *dual PD controller*, which contains 1) a rotation controller controlling the local pose in joint rotational space and 2) an additional position controller controlling the global pose in joint positional space. Below, we elaborate the two controllers.

Joint Rotation Controller. This controller computes the desired joint angular acceleration $\ddot{\boldsymbol{\theta}}_{\text{des}}$ from the estimated reference joint rotations $\boldsymbol{\varphi}$ using:

$$\ddot{\boldsymbol{\theta}}_{\text{des}} = k_{p_\theta}(\mathbf{E}(\boldsymbol{\varphi}) - \boldsymbol{\theta}) - k_{d_\theta} \dot{\boldsymbol{\theta}}, \quad (3)$$

where $\boldsymbol{\theta}$ and $\dot{\boldsymbol{\theta}}$ are the current joint angles and angular velocities; $\mathbf{E}(\cdot)$ transforms the reference pose to local Euler angles; $k_{p_\theta} = 2400$ and $k_{d_\theta} = 60$ are the gain parameters.

Joint Position Controller. This controller computes the desired linear joint acceleration $\ddot{\mathbf{r}}_{\text{des}}$. Different from the rotation controller, we do not have reference joint positions since we have no direct distance measurements. Thus, we compute them from the current joint positions \mathbf{r} and the estimated joint velocity \mathbf{v} as:

$$\mathbf{r}_{\text{ref}} = \mathbf{r} + \mathbf{T}(\mathbf{v})\Delta t, \quad (4)$$

where $\mathbf{T}(\cdot)$ maps the joint velocity from local frame to global frame and Δt is the simulation time step. Then, the joint position controller is defined as:

$$\ddot{\mathbf{r}}_{\text{des}} = k_{p_r}(\mathbf{r}_{\text{ref}} - \mathbf{r}) - k_{d_r} \dot{\mathbf{r}}, \quad (5)$$

where $k_{p_r} = 3600$ and $k_{d_r} = 60$ are the gain parameters; $\dot{\mathbf{r}}$ is the current joint velocity, which can be computed by:

$$\dot{\mathbf{r}} = \mathbf{J} \dot{\mathbf{q}}, \quad (6)$$

where $\mathbf{J} \in \mathbb{R}^{3J \times N}$ is the joint Jacobian.

3.2.3 Motion Tracking Optimizer

The motion tracking optimizer solves a quadratic programming problem and estimates acceleration $\ddot{\mathbf{q}}$, joint torques $\boldsymbol{\tau}$, and GRF $\boldsymbol{\lambda}$. The optimization problem can be written as:

$$\begin{aligned} \arg \min_{\ddot{\mathbf{q}}, \boldsymbol{\lambda}, \boldsymbol{\tau}} \quad & \mathcal{E}_{\text{PD}} + \mathcal{E}_{\text{reg}} \\ \text{s.t.} \quad & \boldsymbol{\tau} + \mathbf{J}_c^T \boldsymbol{\lambda} = \mathbf{M} \ddot{\mathbf{q}} + \mathbf{h} \quad (\text{equation of motion}) \\ & \boldsymbol{\lambda} \in \mathcal{F} \quad (\text{friction cone}) \\ & \dot{\mathbf{r}}_j(\ddot{\mathbf{q}}) \in \mathcal{C} \quad (\text{no sliding}). \end{aligned} \quad (7)$$

\mathbf{M} is computed from \mathbf{q} using the composite rigid body algorithm [9]. \mathbf{h} is computed from \mathbf{q} and $\dot{\mathbf{q}}$ using the recursive Newton-Euler algorithm [9]. The energy function and the three constraints will be elaborated in the following.

Contact Point Determination. To apply the three constraints in Eq. 7, we first need to acquire all the contact points between the body and the ground. We determine whether joint j contacts the ground by its vertical distance to the ground d_j , and for the foot joint we additionally leverage the predicted contact probability c for better accuracy since feet touch the ground more often. Specifically, a foot joint f is considered in contact if 1) $d_f < 0.5\text{cm}$ or 2) $d_f < 3\text{cm}$ and the ground contact probability $c_f > 0.5$; a non-foot joint n is considered in contact only if $d_n < 0.5\text{cm}$. We then draw an $L \times L$ square centered at each contact joint and take its 4 vertices as the contact points. This is based on our finding that assuming facet-contacts instead of point-contacts produces more stable results. We take $L = 20\text{cm}$

which is roughly the size of a foot. The number of the contact joints is denoted as n_j , then the number of the contact points is $n_c = 4n_j$.

Dual PD Controller Term \mathcal{E}_{PD} . To reproduce the kinematic estimation, the character should generate the angular and linear joint accelerations given by the dual PD controller. Thus, in Eq. 7, \mathcal{E}_{PD} consists of two components \mathcal{E}_θ and \mathcal{E}_r , which control the angular and linear accelerations:

$$\begin{aligned}\mathcal{E}_{PD} &= k_\theta \mathcal{E}_\theta + k_r \mathcal{E}_r, \\ \mathcal{E}_\theta &= \|\ddot{\mathbf{q}}_3 - \ddot{\boldsymbol{\theta}}_{\text{des}}\|^2, \mathcal{E}_r = \|\mathbf{J}\ddot{\mathbf{q}} + \dot{\mathbf{J}}\dot{\mathbf{q}} - \ddot{\mathbf{r}}_{\text{des}}\|^2,\end{aligned}\quad (8)$$

where k_θ and k_r are the weight terms both set to 1.

Regularization Term \mathcal{E}_{reg} . Our regularization term \mathcal{E}_{reg} in Eq. 7 contains three energy terms: 1) \mathcal{E}_λ penalizes violations of the Signorini's conditions of contacts [56]; 2) \mathcal{E}_{res} constrains the magnitude of the root residual force; and 3) \mathcal{E}_τ confines the norms of the joint torques:

$$\begin{aligned}\mathcal{E}_{\text{reg}} &= k_\lambda \mathcal{E}_\lambda + k_{\text{res}} \mathcal{E}_{\text{res}} + k_\tau \mathcal{E}_\tau, \\ \mathcal{E}_\lambda &= \sum_{c=1}^{n_c} d_c \|\boldsymbol{\lambda}_c\|^2, \mathcal{E}_{\text{res}} = \|\boldsymbol{\tau}_{:6}\|^2, \mathcal{E}_\tau = \|\boldsymbol{\tau}_6\|^2,\end{aligned}\quad (9)$$

where d_c is the vertical height of the contact point c ; $\boldsymbol{\lambda}_c$ is the GRF at point c ; k_λ , k_{res} , and k_τ are the corresponding weights, which are set to 10, 0.1, and 0.01, respectively.

Friction Cone and Sliding Constraints. These two constraints in Eq. 7 are only applied to the contacts. We assume the GRF at the contact points should be inside the friction cone¹ and the contact joints do not slide. Specifically, we denote the force/velocity along the y (vertical) axis of the global frame as \cdot^y and the same for the x, z (horizontal) axis. The friction cone constraint can be linearized as:

$$\begin{aligned}\mathcal{F}_c &= \{\boldsymbol{\lambda}_c \in \mathbb{R}^3 | \boldsymbol{\lambda}_c^y \geq 0, |\boldsymbol{\lambda}_c^x| \leq \mu \boldsymbol{\lambda}_c^y, |\boldsymbol{\lambda}_c^z| \leq \mu \boldsymbol{\lambda}_c^y\}, \\ \mathcal{F} &= \{[\boldsymbol{\lambda}_1 \cdots \boldsymbol{\lambda}_{n_c}] \in \mathbb{R}^{3n_c} | \boldsymbol{\lambda}_c \in \mathcal{F}_c, c = 1, 2, \dots, n_c\},\end{aligned}\quad (10)$$

which means the vertical force from the ground must be upward, and the horizontal forces should not be larger than the maximum frictional force. We empirically set the friction coefficient $\mu = 0.6$. For the sliding constraint, we have:

$$\begin{aligned}\mathcal{C}_j &= \{\dot{\mathbf{r}}_j \in \mathbb{R}^3 | \dot{\mathbf{r}}_j^y \geq 0, |\dot{\mathbf{r}}_j^x| \leq \sigma, |\dot{\mathbf{r}}_j^z| \leq \sigma\}, \\ \mathcal{C} &= \{[\dot{\mathbf{r}}_1 \cdots \dot{\mathbf{r}}_{n_j}] \in \mathbb{R}^{3n_j} | \dot{\mathbf{r}}_j \in \mathcal{C}_j, j = 1, 2, \dots, n_j\},\end{aligned}\quad (11)$$

which confines the sliding velocity of every contact joint smaller than $\sigma = 0.01$ while preventing ground penetration. The contact joint velocity $\dot{\mathbf{r}}_j(\ddot{\mathbf{q}})$ in Eq. 7 is computed by:

$$\dot{\mathbf{r}}_j(\ddot{\mathbf{q}}) = \mathbf{J}_j(\dot{\mathbf{q}} + \ddot{\mathbf{q}}\Delta t), \quad (12)$$

where $\mathbf{J}_j \in \mathbb{R}^{3n_j \times N}$ is the contact joint Jacobian.

¹Friction cone: the set of all forces that can be transmitted through a Coulomb friction contact. See [3].

3.2.4 Dynamic States Updater

We use a finite difference method for dynamic state updates:

$$\begin{aligned}\mathbf{q}^{(t+1)} &= \mathbf{q}^{(t)} + \dot{\mathbf{q}}^{(t)} \Delta t, \\ \dot{\mathbf{q}}^{(t+1)} &= \dot{\mathbf{q}}^{(t)} + \ddot{\mathbf{q}}^{(t)} \Delta t,\end{aligned}\quad (13)$$

where $\ddot{\mathbf{q}}^{(t)}$ is the estimated acceleration from the optimizer and $\mathbf{q}^{(t+1)}$ is the updated pose and translation. Since our system runs at 60 fps, Δt is set to 1/60 second.

4. Experiments

In this section, we first compare our approach with previous works (Sec. 4.1). Then, we perform an ablative study of the key components (Sec. 4.2). Finally, we show the potential applications of our methods (Sec. 4.3).

Datasets. The training and evaluation involve the AMASS dataset [34], the DIP-IMU dataset [20], and the TotalCapture dataset [62]. Following [73], we first train the model on AMASS using synthesized IMUs and then fine-tune it on the train split of DIP-IMU. The evaluations are performed on TotalCapture and the test split of DIP-IMU. The acceleration measurement in TotalCapture is constantly biased and we re-calibrated it (detailed in the supplemental document). All the reported numbers are online results.

Metrics. We use the following metrics to evaluate our method. 1) *SIP Error* measures the mean orientation error of the upper arms and legs in the global space in degrees. 2) *Mesh Error* measures the mean vertex distance between the reconstructed and ground-truth meshes with both root position and orientation aligned in cm. 3) *Jitter* measures the mean jerk (time derivative of acceleration) of all body joints in the global space in km/s^3 , which reflects the smoothness of the motion [11]. 4) *Zero-Moment Point (ZMP) distance* measures the mean distance from the fictitious ZMP [68] position to the Base of Support² of the character in m. It represents the intensity of the perturbation moment which lets the character fall, and should be zero for a real human in dynamic equilibrium [68]. Previous works [45, 51] leverage Center of Pressure (CoP) to quantify the equilibrium, which is related to the ZMP distance. A discussion about ZMP and CoP can be found in the supplemental document. Among these metrics, 1) and 2) measure pose accuracy, 3) and 4) measure physical plausibility. We further evaluate *cumulative translation error* which means the global position error *w.r.t* the real travelled distance; and *latency* which measures the time from receiving the inertia measurements to outputting the pose and translation for the corresponding frame in ms, using a laptop with an Intel(R) Core(TM) i7-10750H CPU and an NVIDIA RTX2080 Super graphics card. For all these metrics, the lower, the better.

²Base of Support (BoS): the area around the outside edge of the body sections in contact with the ground. Also called the support polygon.

Method	DIP-IMU				
	SIP Err	Mesh Err	Jitter	ZMP Dist	Latency
DIP [20]	17.10	8.96	-	-	117
TransPose [73]	16.68	7.09	1.46	1.67	94
PIP	15.02	5.95	0.24	0.12	16

Method	TotalCapture				
	SIP Err	Mesh Err	Jitter	ZMP Dist	Latency
DIP [20]	18.62	11.22	-	-	117
TransPose [73]	16.58	7.42	1.87	1.40	94
PIP	12.93	6.51	0.20	0.23	16

Table 1. Comparison with the state-of-the-art methods on DIP-IMU [20] and TotalCapture [62]. Metrics and units are detailed at the beginning of Sec. 4. PIP achieves a reduction of 15% of the pose error, 87% of the jitter, and 89% of the motion imbalance with 83% lower latency compared with the SOTA TransPose [73].

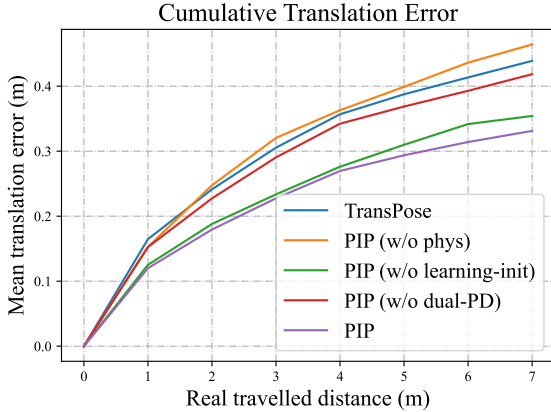


Figure 3. Comparison of the translation estimation on TotalCapture [62]. Our method has the lowest cumulative error because of the learning-based RNN initialization and the dual PD controller.

4.1. Comparisons

Quantitative. We compare our method to state-of-the-art methods DIP [20] and TransPose [73] which also target motion capture from sparse IMUs. Note that DIP does *not* estimate global translation. The results are shown in Tab. 1 and Fig. 3. Our method not only significantly outperforms previous works on capture accuracy and physical plausibility, but also largely reduces the delay. The pose accuracy improvement is attributed to the RNN-based kinematics estimator which makes use of complete historical information and better captures state-change signals. The improvement of the motion smoothness, equilibrium, and translation accuracy is attributed to the physics-aware motion optimizer with the novel dual PD controller. Thus, the proposed combination of learning-based kinematics and optimization-based physics leads to the overall best result.

Qualitative. In Fig. 4, we show the mesh error distribution of DIP [20], TransPose [73], and our method on TotalCapture. We take 4 examples at 1) 10%, 2) mode, 3) median, and 4) 95%. In the first two cases, our method estimates arm and leg orientations better than the previous

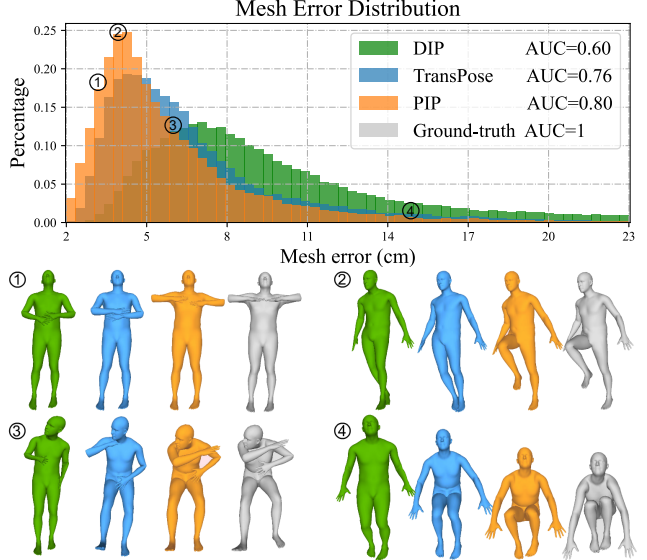


Figure 4. Qualitative results on TotalCapture [62] dataset. We show the mesh error distribution and the Area Under Curve (AUC) value of different methods, and select 4 examples for visualization. Colors of different methods are shown in the legend. Our approach is visually the most accurate over all the methods.

Method	DIP-IMU		TotalCapture	
	SIP Error	Jitter	SIP Error	Jitter
w/o learning-init	15.12	0.27	13.70	0.23
w/o dual-PD	15.04	0.28	12.93	0.32
w/o physics module	15.04	0.48	12.84	0.51
Ours	15.02	0.24	12.93	0.20

Table 2. Ablation study on the learning-based RNN initialization, the dual PD controller, and the physics-based optimizer. It demonstrates the help of our key components on pose accuracy (shown in SIP Error) and physical plausibility (shown in Jitter).

works. In the third case, we reconstruct the full-body pose faithfully while others nearly fail. In the last challenging example, although the estimated upper legs slightly defer from the ground truth, our result still looks similar and outperforms others. Again, we attribute this superiority to the RNN-based estimator and the physics-based optimizer. The ambiguity in these cases comes from the fact that the subject can perform very different poses while keeping the forearm/lower leg orientation unchanged, and the key to resolving the ambiguity is the temporal information of the state-change signals. Compared with previous works, we make better use of such information due to our learning-based RNN initialization. In consequence, our networks regress more accurate pose and velocities, which, in combination with the dual PD controller, further improve the results.

4.2. Evaluations

Physical Properties. In Fig. 5, we demonstrate the estimated physical properties. The left figure shows the GRF of two feet when the subject is walking. We can see that

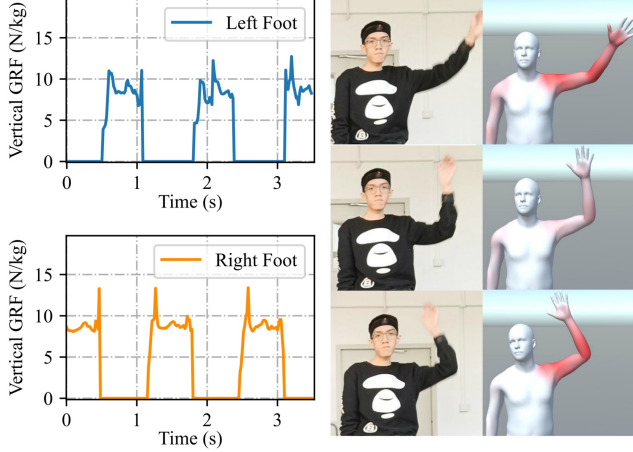


Figure 5. *Left*: evaluation on the estimated contact forces over time in a walking motion. *Right*: visualization of the arm joint torques for a waving-hand motion (red for large torques).

the two feet support the body in turns and the GRF approximately equals gravity, which is reasonable according to [52, 82]. The right figure is a series of frames where the subject waves his hand and we visualize the internal torques of the arm (red for large torques). When the arm starts and stops moving, the torque is larger due to the acceleration.

Learning-based RNN Initialization. To examine the effect of the learning-based initialization, we train the same networks (P_L and V_A) without learning-based initialization and perform the evaluation. As shown in Tab. 2 and Fig. 3, this variant becomes less accurate and stable, which is reflected in the larger errors across all metrics. The learned initialization is most effective when the motion is highly ambiguous (e.g., sitting), but this advantage is numerically averaged out by the common (non-ambiguous) motion in the test dataset. To this end, we pick a long-sitting sequence in the DIP-IMU test dataset and plot the upper leg orientation error over time in Fig. 6. During the sequence, the subject started from standing, then sat down immediately, and kept sitting to the end. The curves show that the method without learning-based initialization starts to fail as time goes by while our model stably tracks the sequence. By examining a few selected frames from the sequence, we can see that the zero-initialized version (in blue) is correct at the beginning but goes wrong after a long period because it loses the historical state-change information from standing to sitting. As a result, its prediction stands up again. In contrast, ours (in orange) always gives the correct estimation due to the good memorizing of such information.

Physics and Dual PD Controller. We evaluate 1) removing the physics-based optimizer, i.e., only using the kinematics module and integrating root velocities to obtain the translation; and 2) removing the joint position controller, i.e., replacing the dual PD controller with a single PD controller that only watches q as in [55]. As shown in Fig. 3,

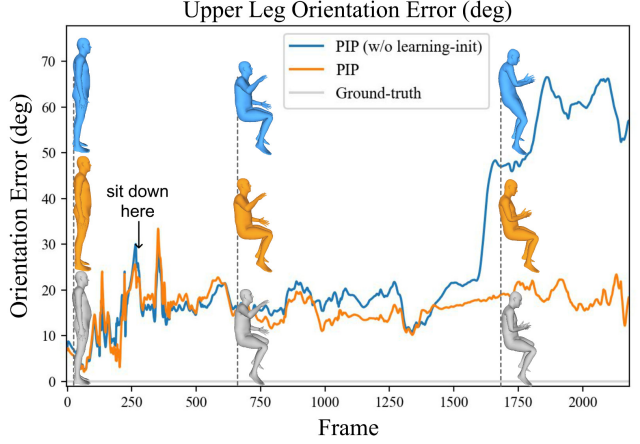


Figure 6. Ablation study on the learning-based RNN initialization skill using a long-sitting sequence. We plot the upper leg orientation error over time and pick three frames for visualization.

without the physics-based optimizer or the dual PD controller, the translation accuracy deteriorates significantly. Besides, the method without physics-based optimization will *always* estimate a character floating in the air or sinking into the ground due to the accumulated error of velocities, and the foot-sliding artifacts are also severe. These facts demonstrate the necessity of physics awareness and our dual PD controller in mitigating translation error accumulation. Quantitative results in Tab. 2 also show a significant reduction of the motion jitter in our full method. However, with the physics-based optimization, the SIP error for TotalCapture is slightly larger (0.1°). This is because the physics module is mainly helpful for estimating translation and improving the physical correctness of the motion.

4.3. Applications

Our method enables several applications such as real-time animation of a virtual character and motion re-targeting. Also note that reducing the latency from 94ms to 16ms is critical to enable applications such as gaming. Please see our supplementary materials for more results.

5. Conclusion and Limitations

In this work, we present the first real-time physics-aware approach that estimates human motion, joint torques, and ground reaction forces from solely 6 IMUs. Combining the kinematics and the physics modules leads to higher accuracy and realism, as shown in our experiments. We also demonstrate exciting applications like live motion capture.

However, we simplify the real world too much, e.g., assuming a flat ground, which makes our method incapable of capturing humans walking upstairs. Besides, the current approach is based on the assumption of a known body shape. For different body shapes, we only need to adjust the bone lengths and the mass distribution of the physics model.

References

- [1] Sheldon Andrews, Ivan Huerta, Taku Komura, Leonid Sigal, and Kenny Mitchell. Real-time physics-based motion capture with sparse sensors. In *Proceedings of the 13th European Conference on Visual Media Production (CVMP 2016)*. Association for Computing Machinery, 2016. 2
- [2] Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. Drecon: Data-driven responsive control of physics-based characters. *ACM Trans. Graph.*, 38, nov 2019. 3
- [3] Stéphane Caron. Friction cones - robot locomotion. Website. <https://scaron.info/robot-locomotion/friction-cones.html>. 6
- [4] Long Chen, Haizhou Ai, Rui Chen, Zijie Zhuang, and Shuang Liu. Cross-view tracking for multi-human 3d pose estimation at over 100 fps. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1
- [5] Rishabh Dabral, Soshi Shimada, Arjun Jain, Christian Theobalt, and Vladislav Golyanik. Gravity-aware monocular 3d human-object reconstruction. In *International Conference on Computer Vision (ICCV)*, 2021. 1
- [6] Junting Dong, Wen Jiang, Qixing Huang, Hujun Bao, and Xiaowei Zhou. Fast and robust multi-person 3d pose estimation from multiple views. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1
- [7] Xiaoxiao Du, Ram Vasudevan, and Matthew Johnson-Roberson. Bio-lstm: A biomechanically inspired recurrent neural network for 3-d pedestrian pose and gait prediction. *IEEE Robotics and Automation Letters*, 4, 2019. 3
- [8] Sai Kumar Dwivedi, Nikos Athanasiou, Muhammed Kocabas, and Michael J. Black. Learning to regress bodies from images using differentiable semantic rendering. In *Proc. International Conference on Computer Vision (ICCV)*, oct 2021. 1
- [9] Roy Featherstone. *Rigid Body Dynamics Algorithms*. Springer US, 01 2008. 3, 5
- [10] Martin Felis. Rbdl: an efficient rigid-body dynamics library using recursive algorithms. *Autonomous Robots*, 41, 02 2017. 13
- [11] Tamar Flash and Neville Hogan. The coordination of arm movements: An experimentally confirmed mathematical model. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 5, 08 1985. 6, 14
- [12] Jack H. Geissinger and Alan T. Asbeck. Motion inference using sparse inertial sensors, self-supervised learning, and a new dataset of unscripted human motion. *Sensors*, 20, 2020. 2
- [13] Andrew Gilbert, Matthew Trumble, Charles Malleison, Adrian Hilton, and John Collomosse. Fusing visual and inertial sensors with semantics for 3d human pose estimation. *Int. J. Comput. Vision*, 127, apr 2019. 2
- [14] Vladimir Guzov, Aymen Mir, Torsten Sattler, and Gerard Pons-Moll. Human positioning system (hps): 3d human pose estimation and self-localization in large scenes from body-mounted sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, jun 2021. 2, 3
- [15] Marc Habermann, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, jun 2020. 1, 2
- [16] Marc Habermann, Weipeng Xu, Michael Zollhöfer, Gerard Pons-Moll, and Christian Theobalt. LiveCap: Real-time human performance capture from monocular video. *ACM Trans. Graph.*, 38, mar 2019. 2, 3
- [17] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3
- [18] Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Christian Theobalt. Real-time body tracking with one depth camera and inertial sensors. In *2013 IEEE International Conference on Computer Vision*, 2013. 2
- [19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9, 12 1997. 12
- [20] 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, (Proc. SIGGRAPH Asia)*, 37, nov 2018. 1, 2, 4, 6, 7, 12, 13
- [21] Mariko Isogawa, Ye Yuan, Matthew O’Toole, and Kris Kitani. Optical non-line-of-sight physics-based 3d human pose estimation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3, 5
- [22] Tomoya Kaichi, Tsubasa Maruyama, Mitsunori Tada, and Hideo Saito. Resolving position ambiguity of imu-based human pose with a single rgb camera. *Sensors*, 20, 2020. 2
- [23] Christoph Kalkbrenner, Steffen Hacker, Maria-Elena Algorri, and Ronald Blechschmidt. Motion capturing with inertial measurement units and kinect - tracking of limb movement using optical and orientation information. In *Proceedings of the International Conference on Biomedical Electronics and Devices - BIODEVICES, (BIOSTEC 2014)*, 03 2014. 2
- [24] 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. 2
- [25] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1
- [26] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 12 2014. 13
- [27] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. Vibe: Video inference for human body pose and shape estimation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2
- [28] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. PARE: Part attention regressor for 3D human body estimation. In *Proceedings International Conference on Computer Vision (ICCV)*. IEEE, oct 2021. 1
- [29] Sergey Levine and Jovan Popović. Physically plausible simulation for character animation. In *Proceedings of the ACM*

SIGGRAPH/Eurographics Symposium on Computer Animation. Eurographics Association, 2012. 5

- [30] Zongmian Li, Jiri Sedlar, Justin Carpentier, Ivan Laptev, Nicolas Mansard, and Josef Sivic. Estimating 3d motion and forces of person-object interactions from monocular video. In *Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [31] Huajun Liu, Xiaolin Wei, Jinxiang Chai, Inwoo Ha, and Tae-hyun Rhee. Realtime human motion control with a small number of inertial sensors. In *Symposium on Interactive 3D Graphics and Games*. Association for Computing Machinery, 2011. 2
- [32] Libin Liu and Jessica Hodgins. Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning. *ACM Trans. Graph.*, 37, jul 2018. 3
- [33] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34, oct 2015. 3
- [34] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2019. 6, 12, 13
- [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*, 128, 06 2020. 2
- [36] 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)*, 2017. 2
- [37] 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, 2016. 2
- [38] 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 *ACM Transactions on Graphics*, volume 39, July 2020. 1
- [39] 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 *ACM Transactions on Graphics*, volume 36, July 2017. 2, 3
- [40] Noitom. Noitom motion capture system. Website. <https://www.noitom.com/>. 1, 2, 13
- [41] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graph.*, 37, jul 2018. 3
- [42] Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, and Sergey Levine. Sfv: Reinforcement learning of physical skills from videos. *ACM Trans. Graph.*, 37, nov 2018. 3
- [43] 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*, 2010. 2
- [44] Patrik Puchert and Timo Ropinski. Human pose estimation from sparse inertial measurements through recurrent graph convolution. *CoRR*, abs/2107.11214, 2021. 2
- [45] Bharadwaj Ravichandran. Biopose-3d and pressnet-kl: A path to understanding human pose stability from video. 2020. 6, 14
- [46] Davis Rempe, Leonidas J. Guibas, Aaron Hertzmann, Bryan Russell, Ruben Villegas, and Jimei Yang. Contact and human dynamics from monocular video. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 3
- [47] Qaiser Riaz, Guanhong Tao, Björn Krüger, and Andreas Weber. Motion reconstruction using very few accelerometers and ground contacts. *Graph. Models*, 79, may 2015. 2
- [48] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2020. 2
- [49] Mike Schuster and Kuldip K Paliwal. Bidirectional recurrent neural networks. *IEEE transactions on Signal Processing*, 45, 1997. 4
- [50] Loren Arthur Schwarz, Diana Mateus, and Nassir Navab. Discriminative Human Full-Body Pose Estimation from Wearable Inertial Sensor Data. In *Modelling the Physiological Human. 3DPH 2009.*, Nov 2009. 2
- [51] Jesse Scott, Bharadwaj Ravichandran, Christopher Funk, Robert T. Collins, and Yanxi Liu. From image to stability: Learning dynamics from human pose. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII*, 2020. 6, 14
- [52] Erfan Shahabpoor and Aleksandar Pavic. Measurement of walking ground reactions in real-life environments: A systematic review of techniques and technologies. *Sensors*, 17(9), 2017. 8, 14
- [53] Mingyi Shi, Kfir Aberman, Andreas Aristidou, Taku Komura, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Motionet: 3d human motion reconstruction from monocular video with skeleton consistency. *ACM Transactions on Graphics (TOG)*, 2020. 3
- [54] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick Pérez, and Christian Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics*, 40, aug 2021. 2, 3, 5, 14
- [55] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: physically plausible monocular 3d motion capture in real time. *ACM Transactions on Graphics*, 39, dec 2020. 2, 3, 4, 5, 8, 14
- [56] A Signorini. Questioni di elastostatica linearizzata e semilinearizzata. *Rend. Mat. Appl.*, XVIII, 1959. 6
- [57] Ronit Slyper and Jessica Hodgins. Action capture with accelerometers. *ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 01 2008. 2

- [58] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Black Michael J., and Tao Mei. Monocular, one-stage, regression of multiple 3d people. In *ICCV*, October 2021. 2
- [59] Luke Wicent F. Sy, Nigel H. Lovell, and Stephen J. Redmond. Estimating lower limb kinematics using a lie group constrained extended kalman filter with a reduced wearable imu count and distance measurements. *Sensors*, 20, 2020. 2
- [60] Jochen Tautges, Arno Zinke, Björn Krüger, Jan Baumann, Andreas Weber, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bernhard Eberhardt. Motion reconstruction using sparse accelerometer data. *ACM Transactions on Graphics*, 30, 05 2011. 2
- [61] Denis Tome, Matteo Toso, Lourdes Agapito, and Chris Russell. Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture. In *2018 International Conference on 3D Vision (3DV)*, 2018. 1
- [62] Matthew Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John Collomosse. Total capture: 3d human pose estimation fusing video and inertial sensors. In *2017 British Machine Vision Conference (BMVC)*, 09 2017. 2, 6, 7, 12, 13
- [63] Vicon. Award winning motion capture systems. Website. <https://www.vicon.com/>. 13
- [64] Daniel Vlasic, Rolf Adelsberger, Giovanni Vannucci, John Barnwell, Markus Gross, Wojciech Matusik, and Jovan Popović. Practical motion capture in everyday surroundings. *ACM Trans. Graph.*, 26, jul 2007. 2
- [65] 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 *European Conference on Computer Vision (ECCV)*, sep 2018. 2
- [66] Timo von Marcard, Bodo Rosenhahn, Michael Black, and Gerard Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. *Computer Graphics Forum 36(2), Proceedings of the 38th Annual Conference of the European Association for Computer Graphics (Eurographics)*, 2017. 1, 2, 3
- [67] Marek Vondrak, Leonid Sigal, Jessica Hodgins, and Odest Jenkins. Video-based 3d motion capture through biped control. *ACM Trans. Graph.*, 31, jul 2012. 3
- [68] Miomir Vukobratovic and Branislav Borovac. Zero-moment point - thirty five years of its life. *I. J. Humanoid Robotics*, 1, 03 2004. 6, 14
- [69] Jian Wang, Lingjie Liu, Weipeng Xu, Kripasindhu Sarkar, and Christian Theobalt. Estimating egocentric 3d human pose in global space. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021. 1
- [70] Xiaolin Wei and Jinxiang Chai. Videomocap: Modeling physically realistic human motion from monocular video sequences. *ACM Trans. Graph.*, 29, jul 2010. 3
- [71] Xsens. Xsens 3d motion tracking. Website. <https://www.xsens.com/>. 1, 2
- [72] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, jun 2020. 1
- [73] Xinyu Yi, Yuxiao Zhou, and Feng Xu. Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. *ACM Transactions on Graphics*, 40, 08 2021. 1, 2, 3, 4, 6, 7, 12, 13
- [74] Ri Yu, Hwangpil Park, and Jehee Lee. Human dynamics from monocular video with dynamic camera movements. *ACM Trans. Graph.*, 40, 2021. 3
- [75] Tao Yu, Kaiwen Guo, Feng Xu, Yuan Dong, Zhaoqi Su, Jianhui Zhao, Jianguo Li, Qionghai Dai, and Yebin Liu. Bodyfusion: Real-time capture of human motion and surface geometry using a single depth camera. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2017. 1
- [76] Tao Yu, Jianhui Zhao, Zhang Zerong, Kaiwen Guo, Dai Quionhai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performance with inner body shape from a depth sensor. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, july 2019. 1
- [77] Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR2021)*, June 2021. 1
- [78] Ye Yuan and Kris Kitani. Ego-pose estimation and forecasting as real-time pd control. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3
- [79] Ye Yuan and Kris Kitani. Residual force control for agile human behavior imitation and extended motion synthesis. In *Advances in Neural Information Processing Systems*, 2020. 5
- [80] Ye Yuan, Shih-En Wei, Tomas Simon, Kris Kitani, and Jason Saragih. Simpoe: Simulated character control for 3d human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3, 5
- [81] 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 *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [82] Petrisa Zell, Bodo Rosenhahn, and Bastian Wandt. Weakly-supervised learning of human dynamics. In *ECCV*, 07 2020. 3, 8, 14
- [83] Petrisa Zell, Bastian Wandt, and Bodo Rosenhahn. Joint 3d human motion capture and physical analysis from monocular videos. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017. 3
- [84] Hongwen Zhang, Yating Tian, Xinchu Zhou, Wanli Ouyang, Yebin Liu, Limin Wang, and Zhenan Sun. Pymaf: 3d human pose and shape regression with pyramidal mesh alignment feedback loop. In *Proceedings of the IEEE International Conference on Computer Vision*, 2021. 1
- [85] Siwei Zhang, Yan Zhang, Federica Bogo, Pollefeys Marc, and Siyu Tang. Learning motion priors for 4d human body capture in 3d scenes. In *International Conference on Computer Vision (ICCV)*, Oct. 2021. 3

- [86] 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 *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1
- [87] Zhe Zhang, Chunyu Wang, Wenhui Qin, and Wenjun Zeng. Fusing wearable imus with multi-view images for human pose estimation: A geometric approach. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [88] 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 *IEEE Conference on Computer Vision (ICCV 2021)*, 2021. 1
- [89] Yu Zheng and Katsu Yamane. Human motion tracking control with strict contact force constraints for floating-base humanoid robots. In *2013 13th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, 2013. 4
- [90] 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. 2
- [91] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3, 12
- [92] Yuxiao Zhou, Marc Habermann, Ikhsanul Habibie, Ayush Tewari, Christian Theobalt, and Feng Xu. Monocular real-time full body capture with inter-part correlations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021. 1
- [93] Yuliang Zou, Jimei Yang, Duygu Ceylan, Jianming Zhang, Federico Perazzi, and Jia-Bin Huang. Reducing footskate in human motion reconstruction with ground contact constraints. In *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020. 3

A. Implementation Details

Network Structure. We schematically visualize the network structures in our kinematics module in Fig. 7. The recurrent neural network (RNN) P_L , P_A , R_A , V_A , and C_F share the same structure. Each network includes a linear input layer with a ReLU activation, two Long Short-term Memory (LSTM) [19] layers with the width of 256, and a linear output layer. A 40% dropout is applied to prevent over-fitting. The RNN C_F is finally activated by a Sigmoid function to obtain probability values. The initial states of P_L and V_A are regressed from the starting leaf joint positions $\mathbf{p}_{\text{leaf}}^{(0)}$ and joint velocities $\mathbf{v}^{(0)}$ using the fully-connected network (FCN) I_{PL} and I_{VA} , respectively. Each FCN consists of 3 fully-connected (FC) layers with the width of 256, 512, and 1024 using the ReLU activation.

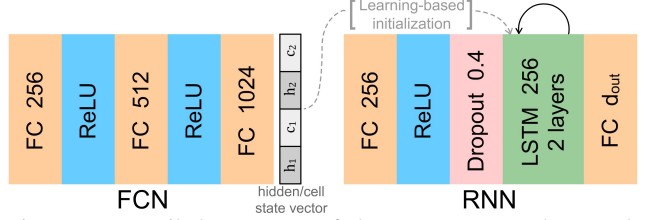


Figure 7. Detailed structures of the recurrent neural network (RNN) and the fully-connected network (FCN) in our kinematics module. "FC" represents a fully-connected layer. The output dimension and other hyper-parameters are marked in each block.

The output of the FCN is used to initialize the hidden/cell states of the two LSTM layers of the RNN at the beginning.

Rotation Representation. The inertia input vector \mathbf{x} consists of accelerations and *rotation matrices*, which are obtained after the calibration. The output of R_A is the non-root joint rotations w.r.t the root parameterized in the *6D representation* [91]. Combining the estimated non-root joint rotations with the root orientation measured by the IMU placed on the pelvis, we obtain the vector φ . The character pose in the physics module is described by local joint rotations (i.e., each joint relative to its parent) in *Euler angles*, which is denoted as θ . The configuration vector $\mathbf{q} = [\mathbf{r}_{\text{root}} \ \theta]$ is then composed of the root translation and the pose in Euler angles.

Datasets. Following [73], we use the AMASS [34] dataset and the train split of the DIP-IMU [20] dataset for the network training, and use the TotalCapture [62] dataset and the test split of the DIP-IMU dataset for evaluation. For AMASS, we synthesize the IMU measurements and foot-ground contact labels as proposed by Yi et al. [73], and synthesize the ground-truth joint velocities using:

$$\mathbf{v}^{\text{GT}}(t) = (\mathbf{R}_{\text{root}}^{\text{GT}}(t))^{-1}(\mathbf{r}^{\text{GT}}(t) - \mathbf{r}^{\text{GT}}(t-1))/\Delta t, \quad (14)$$

where $\mathbf{R}_{\text{root}}^{\text{GT}}(t) \in \mathbb{R}^{3 \times 3}$ is the ground-truth root orientation at frame t ; $\mathbf{r}^{\text{GT}} \in \mathbb{R}^{3J}$ is the ground-truth joint global positions; Δt is the frame interval. We also re-calibrate the acceleration measurements in TotalCapture, as we find that they are constantly biased (see Fig. 8). Specifically, to remove the bias, we synthesize the accelerations for TotalCapture using the method of Yi et al. [73] and align the mean acceleration measurement for each sequence to the mean synthetic values by adding or subtracting a constant.

Gain Parameters for PD Controllers. The gain parameters k_{p_θ} , k_{d_θ} , k_{p_r} , and k_{d_r} of the dual PD controller introduced in Sec. 3.2.2 are derived as follows. Take the joint rotation controller (controlling θ) as an example. As we use first-order approximations in the dynamic states updater (Sec. 3.2.4), we apply first-order Taylor expansion on θ and $\dot{\theta}$, and rearrange the equation, which writes:

$$\ddot{\theta}(t) = \frac{1}{\Delta t^2}(\theta(t+2\Delta t) - \theta(t+\Delta t)) - \frac{1}{\Delta t}\dot{\theta}(t), \quad (15)$$

Method		DIP-IMU							
		SIP Error	Ang Error	Pos Error	Mesh Error	Rel Jitter	Abs Jitter	ZMP Dist	Latency
Offline	DIP [20]	16.36	14.41	6.98	8.56	2.34	-	-	-
	TransPose [73]	13.97	7.62	4.90	5.83	0.13	0.85	0.59	-
Online	DIP [20]	17.10	15.16	7.33	8.96	3.01	-	-	117
	TransPose [73]	16.68	8.85	5.95	7.09	0.61	1.46	1.67	94
	PIP (Ours)	15.02	8.73	5.04	5.95	0.23	0.24	0.12	16

Method		TotalCapture							
		SIP Error	Ang Error	Pos Error	Mesh Error	Rel Jitter	Abs Jitter	ZMP Dist	Latency
Offline	DIP [20]	18.47	17.54	9.47	11.19	2.91	-	-	-
	TransPose [73]	14.71	12.19	5.44	6.22	0.16	0.91	0.76	-
Online	DIP [20]	18.62	17.22	9.42	11.22	3.62	-	-	117
	TransPose [73]	16.58	12.89	6.55	7.42	0.95	1.87	1.40	94
	PIP (Ours)	12.93	12.04	5.61	6.51	0.20	0.20	0.23	16

Table 3. Comparison with the state-of-the-art methods on more metrics. PIP outperforms previous online methods on all metrics with much less latency, while also achieves comparable capture accuracy but higher physical correctness when compared with previous offline methods. This demonstrates the superiority of our system which runs in real-time with extremely small latency.

where $\Delta t = 1/60$ is the time interval between frames. By associating this equation with Eq. 3 and Eq. 5 in the main paper, the proportional gain k_{p_θ} and k_{p_r} should be 3600, and the derivative gain k_{d_θ} and k_{d_r} should be 60. For the joint rotation controller, setting the proportional gain k_{p_θ} to a lower value gives smoother angular accelerations. Thus, we set k_{p_θ} to 2400 in our experiments.

Other Details. We use a laptop with an Intel(R) Core(TM) i7-10750H CPU and an NVIDIA RTX2080 Super graphics card to run the experiments and the live demos. We use PyTorch 1.8.1 with CUDA 11.1 to implement our kinematics estimator, and leverage the Rigid Body Dynamics Library [10] to implement our physics-based optimizer. The live demo is implemented using Unity3D. We use Noitom Perception Neuron Laboratory [40] IMU sensors in our demo. Both training and evaluation assume 60 fps sensor input. The training data is additionally clipped into short sequences in 200-frame lengths for more effective learning. Specifically, we separately train each RNN in the kinematics module using the synthetic AMASS [34] dataset with a batch size of 256 using the Adam [26] optimizer, and fine-tune P_L (together with I_{PL}), P_A , and R_A on the train split of the DIP-IMU dataset, following [73]. We do not train V_A and C_F on DIP-IMU as it does not contain global movements.

B. Comparisons on More Metrics

In this section, we show the comparison results with the previous state-of-the-art methods [20, 73] on more metrics. In addition to the metrics used in the main paper, we also evaluate 1) *Angular Error*: the mean rotation error of all body joints in the global space in degrees; 2) *Positional Error*: the mean position error of all body joints in the global space with the root position and orientation aligned

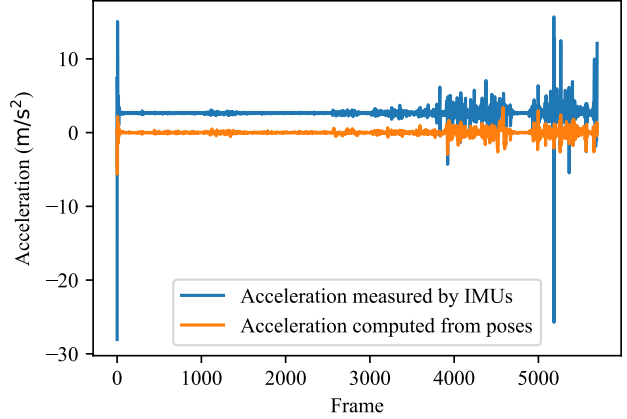


Figure 8. The acceleration measurements in the TotalCapture [62] dataset is constantly biased. We visualize the accelerations (x -axis component) measured by IMUs in blue and the one computed from the subject motions based on Vicon [63] by a finite-difference method in orange. We can see an obvious constant bias in the IMU acceleration measurements (blue) based on the fact that real accelerations should be approximately zero-centered.

in cm; 3) *Relative Jitter*: the jitter calculated in the local (root-relative) frame in km/s^3 , where the root translation is not considered. Notice that due to the length limit of the main paper, we only showed the mesh error as it incorporates both angular and positional error, and the SIP error as it is directly related to motion ambiguities in the main text. Here, we report the results on more metrics for a fair comparison. We also evaluate previous offline methods for references, which need to pre-record the inertia measurements during the whole motion and estimate the motion with the help of the complete inertia sequence. The results on TotalCapture [62] and the test split of DIP-IMU [20] dataset are shown in Tab. 3. We outperform previous on-

line methods on all metrics with largely reduced latency, which demonstrates the accuracy and effectiveness of our approach. Moreover, compared with the offline methods, PIP achieves comparable motion accuracy (reflected in the first 5 metrics) but higher physical plausibility (reflected in Absolute Jitter and ZMP Distance). We attribute this to the physics-based motion optimizer proposed in the main paper. Most importantly, our system runs *in real-time*, while the offline approaches require the access to the complete inertia sequence. Thus, our approach significantly closes the gap between online and offline methods, and enables a wide variety of real-time applications such as gaming.

C. Discussions and Future Works

Quantitative Evaluations of Physics. A direct quantitative evaluation of physics (*e.g.*, joint torques and ground reaction forces) would be advantageous. However, to the best of our knowledge, there is no public dataset containing both IMU measurements and ground-truth forces (either joint torques or ground reaction forces). We believe that creating such a dataset requires research on its own, and would have great value for the community. For now, we can only provide qualitative visualization of torques/GRFs in our supplemental video and Fig. 5, which is intuitively plausible and in line with the references [52, 82]. Besides, as the output motion is *entirely* driven by the estimated forces, the quantitative evaluation of the motion can also implicitly demonstrate the quality of our force estimation. Furthermore, we use jitter (jerk) and ZMP distance as indirect quantitative evaluations of the physics estimation, which reflect the naturalness [11] and equilibrium [68] of the motion, respectively. Since we do not adopt any explicit penalty on these two metrics, nor do we use any temporal filter or balancing technique on the motion, the better results on these two metrics actually suggest the improved physical correctness achieved by our motion optimizer.

Regarding the ground contact evaluation, previous works [54, 55] use mean penetration error to evaluate the non-physical foot penetration. As we explicitly model the contacts as hard constraints, both sliding and ground penetration are *strictly* avoided with any contacting part of the body. Thus, these errors would be zero.

Zero Moment Point vs. Center of Pressure. Previous works [45, 51] use Center of Pressure (CoP) accuracy to quantify the force estimation, which is related to our Zero Moment Point (ZMP) distance. Here we point out the difference between these two notations and the reason why we choose to use ZMP distance. The pressure between the human body and the ground can be represented by a force exerted at the CoP. If such a force can balance all active forces acting on the human body during the motion, the human body is in dynamic equilibrium, and ZMP coincides with CoP (*i.e.*, within the support polygon). However, when the

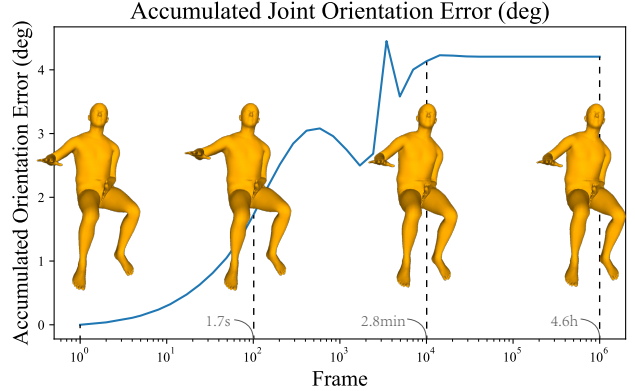


Figure 9. Pose drifts in a perfectly-still sitting pose. We evaluate PIP on 4.6-hour artificial inertia measurements with zero accelerations and fixed orientations of a sitting pose. We plot the accumulated orientation error of all body joints over time and pick four frames for visualization. Our system stably estimates a sitting pose during the entire sequence with a total drift of 4.2 degrees.

force acting on the CoP cannot balance other forces, the human will fall down about the foot edge, and the ZMP (more precisely, the fictitious ZMP) will be outside the support polygon, whose distance to the polygon is proportional to the intensity of the unbalanced force. In such cases, CoP is on the border of the support polygon as the ground reaction forces cannot escape the polygon. Thus, the reason to use ZMP distance in our physics evaluation becomes clear: since the estimated motion cannot be perfectly physically correct and contains unbalanced movements, the ZMP distance can better reflect the disequilibrium in the captured motion. On the other hand, evaluating CoP accuracy needs a more sophisticated modeling of human feet (rather than a simplified square facet contact) and ground-truth pressure annotations, which we leave as a future work. For more detailed introductions of ZMP, readers are referred to [68].

Drifts in Long-term Tracking. As a purely inertial sensor based approach, PIP inevitably suffers from drifts in long-term tracking. As measured in Fig. 3, the translation drift of our system depends on how far the subject moves, and is about 4.6% in our experiments. Regarding the subject's pose, we do not see an evident drift in our experiments. This may be because the subject is always moving, and the orientation and acceleration measurements effectively confine the possible human pose. Therefore, it is interesting to examine the *pose drift in still poses*, especially for the ambiguous ones like sitting. However, as the IMUs always have small noises and humans cannot keep perfectly still for a long time, it is difficult to quantify the pose drifts in real settings. Thus, we conduct a toy experiment where we artificially set all acceleration measurements to zero and orientations unchanged at the point after the sit-down motion in Fig. 6, *i.e.*, to simulate a perfectly-still sitting pose. As shown in Fig. 9, our system can keep estimating sitting

poses stably with a total drift of 4.2 degrees for all body joints at 1 million (4.6 hours) frames. This demonstrates the robustness of our system in long-term tracking, which is ensured by the RNNs and the learning-based RNN initialization scheme. We also conduct a live experiment where our method can track long-period sitting for half an hour stably and is not getting worse as time goes by. Please refer to our supplementary video for more results.