

Transformer Inertial Poser: Attention-based Real-time Human Motion Reconstruction from Sparse IMUs

YIFENG JIANG, Stanford University, USA

YUTING YE, Reality Labs Research, Meta, USA

DEEPAK GOPINATH, AI Research, Meta, USA

JUNGDAM WON, AI Research, Meta, USA

ALEXANDER W. WINKLER, Reality Labs Research, Meta, USA

C. KAREN LIU, Stanford University, USA

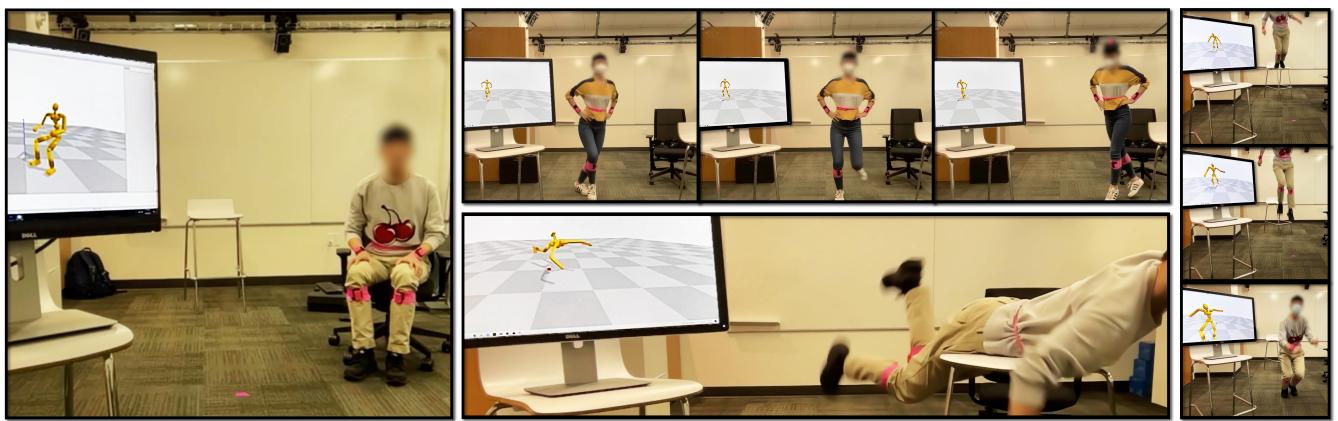


Fig. 1. We develop an attention-based deep learning method to reconstruct full-body motion from six IMU sensors in real-time. In addition to common locomotion, our method can produce stable stationary motion without drifting, such as sitting still, and dynamic motions, such as kicking and dancing, performed on arbitrary terrains.

Real-time human motion reconstruction from a sparse set of wearable IMUs provides an non-intrusive and economic approach to motion capture. Without the ability to acquire absolute position information using IMUs, many prior works took data-driven approaches that utilize large human motion datasets to tackle the under-determined nature of the problem. Still, challenges such as temporal consistency, global translation estimation, and diverse coverage of motion or terrain types remain. Inspired by recent success of Transformer models in sequence modeling, we propose an attention-based deep learning method to reconstruct full-body motion from six IMU sensors in real-time. Together with a physics-based learning objective to predict "stationary body points", our method achieves new state-of-the-art results both quantitatively and qualitatively, while being simple to implement and

smaller in size. We evaluate our method extensively on synthesized and real IMU data, and with real-time live demos.

CCS Concepts: • Computing methodologies → Motion capture.

Additional Key Words and Phrases: Attention in Neural Networks, Wearable Devices, Inertial Measurement Units, Kinematic Constraints, Human Motion

ACM Reference Format:

Yifeng Jiang, Yuting Ye, Deepak Gopinath, Jungdam Won, Alexander W. Winkler, and C. Karen Liu. 2022. Transformer Inertial Poser: Attention-based Real-time Human Motion Reconstruction from Sparse IMUs. *ACM Trans. Graph.* 1, 1 (March 2022), 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Authors' addresses: Yifeng Jiang, Stanford University, USA; Yuting Ye, Reality Labs Research, Meta, USA; Deepak Gopinath, AI Research, Meta, USA; Jungdam Won, AI Research, Meta, USA; Alexander W. Winkler, Reality Labs Research, Meta, USA; C. Karen Liu, Stanford University, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

0730-0301/2022/3-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Real-time reconstruction of 3D human motion is crucial for applications in various domains, such as biomechanics and sports analysis, motion-based video games, and virtual presence in VR/AR systems. While marker-based optical motion capture systems [vic [n.d.]] remain an ideal option for research labs and professional studios due to the superior accuracy, more and more applications demand a portable, less costly, and minimally-invasive mocap system that reconstructs human movements in real-time and can be used anywhere by everyone.

Among many proposed sensing modalities, such as RGB cameras [Cao et al. 2019; Güler et al. 2018; Kanazawa et al. 2019], depth

cameras [Taylor et al. 2012; Wei et al. 2012], or wearable electromagnetic sensors [Kaufmann et al. 2021], inertial measurement sensors (IMUs) [Huang et al. 2018; von Marcard et al. 2017] stand out for many unique advantages they provide. IMU-based mocap is egocentric and untethered, applicable to both indoor and outdoor activities; is unsusceptible to occlusion; and is less sensitive to privacy issues. The state-of-the-art [Yi et al. 2021] has shown that by learning multiple RNN models from large-scale human motion databases, such as AMASS [Mahmood et al. 2019], it is possible to estimate full 3D skeletal poses from only six IMUs with good accuracy.

This paper continues to push the state-of-the-art in real-time human motion reconstruction using six IMU sensors. Inspired by recent success of GPT-style models for natural language generation [Radford et al. 2018], we propose to learn a small GPT model which takes as input a recent history of IMU readings and predicts a current full-body pose. Comparing to the state-of-the-art, our model improves human motion reconstruction both qualitatively and qualitatively across a large variety of motion types. In particular, previous RNN-based models do not explicitly make use of its own previous predictions, and thus struggle to distinguish between motion trajectories with similar orientations and persistent zero-acceleration at sensor locations, such as sitting versus standing [Huang et al. 2018]. Because GPT models are conditioned on the entire history of past predictions when making the next prediction, our method can create consistent long sequence, such as sitting motion without switching between sitting and standing poses.

Another fundamental challenge in IMU-based systems is the reconstruction of the root motion in which small errors accumulate rapidly due to the lack of feedback correction. Previous methods, e.g. [Yi et al. 2021], rely on heuristics to detect or learn the current supporting foot and enforce the foot to be stationary. This approach tends to generate "foot locking" artifacts and fails to handle non-locomotion or locomotion over non-flat terrains (e.g. slope, stairs). We introduce a physics-based learning objective to predict the "stationary body points" (i.e. a Cartesian point on the character that has zero velocity) from the previous time step. We then use the predicted stationary points to correct the predicted pose. Learned station body points are particularly helpful when there is a shift in distribution between the test and training sets due to noisy or corrupted IMU sensors, or unseen motion.

Our algorithm achieves new state-of-the-art results on both synthesized and real IMU datasets. In practice, our model is easy to implement and smaller in size (therefore faster to inference), compared to existing methods. Besides the AMASS dataset, we evaluate our method on DanceDB, DIP, TotalCapture datasets across a wide range of challenging scenarios. We also show live demos of our system on free-form motions.

2 Related Work

Human motion reconstruction from various sensor inputs has been studied for a long time especially in Computer Graphics and Computer Vision communities. We review mainly the prior work that is most relevant to our work, which use any part of the inputs that constitute IMU sensors. We also review motion generation models based on Transformer [Vaswani et al. 2017] because it constitutes the core of our reconstruction model.

A typical IMU sensor includes an accelerometer measuring 3-axis linear acceleration, a gyroscope measuring 3-axis angular velocity, a magnetometer identifying the vector towards Earth's magnetic North. Sensor fusion algorithms based on Kalman filter or its extended version are used to provide measures of orientation and heading [Bachmann et al. 2001; Del Rosario et al. 2018; Foxlin 1996; Vitali et al. 2021]. Although many commercial solutions can provide stable orientations, the estimation of absolute positions are still inaccurate and noisy due to the inherent nature of IMU sensors, where orientations are the values integrated once with the gyroscope inputs whereas positions are the value integrated twice with the accelerometer inputs. Because both accurate orientation and position information are important in human motion reconstruction, it has traditionally been considered that IMUs work best when combined with other sensors (i.e. sensor fusion). One of the most popular fusion options is using vision-based sensors such as RGB or RGB-D cameras. In a high level view, we can regard the approaches as adding extra constraints by IMUs to the motions predicted from the vision inputs, where the constraints are used in off-line optimization [Helten et al. 2013; Pons-Moll et al. 2011, 2010; von Marcard et al. 2018, 2016; Zheng et al. 2018], online per-frame optimization [Charles Malleson 2020; Malleson et al. 2017; Zhang et al. 2020], or learning deep neural networks [Gilbert et al. 2019; Trumble et al. 2017]. There have also been other sensor fusions with IMUs such as optical markers [Andrews et al. 2016] or ultrasonic [Liu et al. 2011; Vlasic et al. 2007].

As IMUs sensors are getting smaller and cheaper, they received increasing attentions from both industry and research communities as a standalone body tracking solution. Popular commercial products such as [xse [n.d.]] and [rok [n.d.]] can generate high-quality human motions ready to be used in real-time game engines. However, they are still not accessible to everyday users because they require a sophisticated setup with at least 17 IMU sensors all over the body. Researcher therefore proposed body tracking systems with a small number of IMUs, usually utilizing existing statistical body models or high quality optical motion capture data as prior information. Marcard et al. [2017] developed an off-line system (SIP) with only 6 IMUs, which optimizes the parameters of the SMPL body model [Loper et al. 2015] so that it fits to the sparse sensor input. Huang et al. [2018] learned a deep neural network model (DIP) from a large amount of motion capture data to directly map the IMU input to a body pose. Their model is based on bidirectional recurrent neural networks (BRNN), so the system can run in an online manner while considering both the past and future sensor inputs with a negligible delay. An ensemble of BRNNs was further adopted by Nagaraj et al. [2020] to improve upon the results. However, the two real-time solutions mostly focus on reconstructing the local body motion without global transformation. The current state-of-the-art system, *TransPose* [Yi et al. 2021], also trained a deep neural network model where the progressive up-scaling of joint position estimation showed more accurate pose estimation. They can additionally generate accurate global root motions by using a confidence-based fusion of a supporting-foot heuristics and a small learned deep network. We instead propose a much simpler method that is faster, easy-to-implement, and most importantly, producing better reconstruction results for both the body and the root.

The blooming AR/VR industry draws attention to full body motion tracking and synthesis from IMUs on headsets and controllers. These devices can usually produce 6D poses by sensor fusion. However, no sensors are available on the lower body. Utilizing deep reinforcement learning, Luo et al. [2021] is able to generate physically valid locomotion from only the 6D egocentric handset pose. Cha et al. [2021] complements pose estimation from headset cameras with IMUs only when the hands are out-of-view. Choutas et al. [2021] and Dittadi et al. [2021] experimented with deep generative models conditioned on headset and controllers poses to regress parameters of the SMPL model. Most similarly to our work is LoBSTR [Yang et al. 2021], where they include an IMU on the waist in addition to IMUs on the headset and controllers. Using a recurrent network, they can synthesize both sitting and running motions from only 4 sensors. We opt to use 6 IMUs with lower body information for accurate motion reconstruction rather than synthesis, but these sparser setups are fruitful future directions.

Other types of lightweight wearable sensors have also been explored to reconstruct human motion. Comparing to vision-based sensors, they have similar advantages as IMUs, as being robust to occlusions and adversarial lighting. Earlier work finds nearest neighbors of input signals in a motion database [Chai and Hodgins 2005] to output continuous motion from 5 accelerometers [Slyper and Hodgins 2008; Tautges et al. 2011]. Similar approaches were later applied to as few as 3 accelerometers placed on both wrists and lower trunk [Riaz et al. 2015]. Electromagnetic-field (EM) sensing recently becomes another viable solution, in which an EM sensor measures its position and orientation relative to a magnetic field emitter. Kaufmann et al. [2021] demonstrated a wearable pose estimation system with 12 wireless EM sensors and an emitter, all worn on the body. They used learned gradient descent to fit SMPL parameters to these 6D sensor poses.

Since the inception of Transformer, attention-based models have been applied to many problems involving sequence data and become state-of-the-art, such as in language translation [Brown et al. 2020] and audio generation [Dhariwal et al. 2020]. It is natural to also apply Transformer models in synthesizing motion sequences. Aksan et al. [2021] developed a generative model using dual attention mechanism to capture spatial and temporal correlations, which predicts future full-body motions given a short history. Petrovich et al. [2021] used a Transformer and a variational autoencoder conditioned on action labels, such as *walking* or *jumping*, to generate full body motions. Valle-Pérez et al. [2021] instead combined a Transformer with normalizing flows to synthesize dancing motion from music, building on a similar prior work [Li et al. 2021]. For motion reconstruction, Kim et al. [2021] experimented with a Transformer encoder-decoder model with sparse synthetic input features, and found it more effective than recurrent networks. In our case, we face the additional challenge of handling noise from real IMU sensors.

3 Transformer Inertial Poser (TIP)

We introduce a real-time human motion reconstruction technique from six IMU sensors placed on the user’s legs, wrists, head, and pelvis. Our data-driven approach trains a neural network model to output a current estimation of full-body joint angles q_t and root velocity v_t from a real-time stream of orientation readings R and

acceleration readings A provided by the IMU sensors (Figure 2). Additionally, our model predicts stationary body points c_t which improve the accuracy of the reconstructed motion. Recent predictions q_t and c_t are fed back to the model so the model can explicitly condition next predictions on its own past predictions. During test time, we run contact filtering on predicted root velocity v_t and smoothing on predicted q_t as post-processing steps.

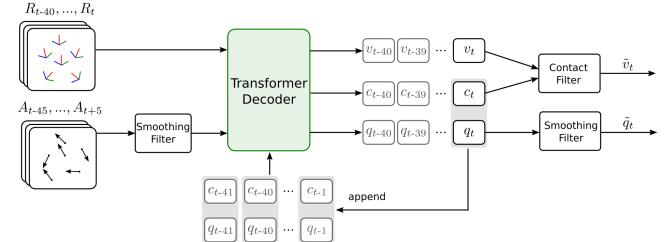


Fig. 2. Overview of our pose estimation algorithm.

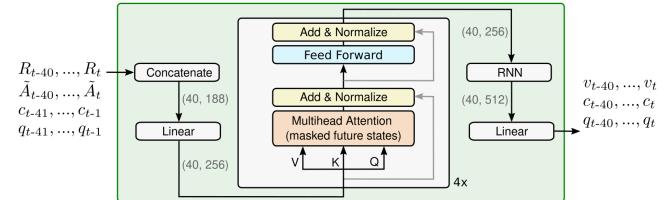


Fig. 3. Transformer Decoder architecture.

3.1 The Transformer Model

At the core of our algorithm is a Transformer Decoder model, which has shown excellent performance for sequence prediction tasks, comparing to alternatives, such as LSTM [Hochreiter and Schmidhuber 1997] or temporal convolution [Bai et al. 2018]. Following previous learning-based methods, our model input includes the IMU orientations $R \in \mathbb{R}^{54 \times 40}$, represented as flattened rotation matrix (length 9) with window size of 40. The input also includes the raw IMU acceleration readings $A \in \mathbb{R}^{18 \times 50}$, smoothed by a moving average filter to $\tilde{A} \in \mathbb{R}^{18 \times 40}$ (Section 3.1). The model is trained to predict the human joint pose $q_t \in \mathbb{R}^{57}$, represented by the axis-angle of 19 major joints defined in the SMPL [Loper et al. 2015] human model. The model also predicts the root linear velocity, $v_t \in \mathbb{R}^3$, which can be integrated to recover the root translation. The root orientation is given directly by one of the IMUs placed on the pelvis.

Our method adapts a Transformer Decoder commonly seen in the GPT natural language model [Radford et al. 2018]. In summary, Transformer models are fully-connected residual networks with temporal attention mechanism (please see [Vaswani et al. 2017] for details). A Transformer Decoder takes in the sequence of past-to-present predictions and generates the most plausible next prediction. Unlike free-form language generation [Radford et al. 2018], we have additional constraints from the IMU sensors. As such, we also

feed the model the sequence of IMU readings in parallel to its past predictions (Figure 3).

Instead of learning only to predict the next frame t at run-time, for efficiency during training time the model is asked to predict in parallel the whole sequence from $t - 40$ to t . To prevent the model from learning simply to shift the input, a causal mask [Vaswani et al. 2017] is added to hide attention information such that the prediction of a frame $t - i$ is only allowed to use information from $t - 40$ to $t - i - 1$. Because the model is auto-regressive (taking its own output as input) at run-time, c_t and q_t could deviate from those used during training. We therefore add a 80% dropout [Srivastava et al. 2014] to the input c_t and q_t to prevent the model from overly relying on its past predictions. In practice, we found that excluding v_t from history is important to prevent auto-regressive drifting of the root translation during test time, possibly as v_t tends to be close to constant in a time window which the model during training could easily exploit and overfit.

Optical-based motion capture datasets are abundant in quantity and diverse in motion types, but they do not have the corresponding IMU data for supervised learning. Following previous work [Huang et al. 2018], we place virtual IMU sensors on virtual characters driven by captured motions to synthesize IMU orientation and acceleration readings. We created synthetic IMU data using the AMASS [Mahmood et al. 2019] motion dataset, which is 20 times larger than the real IMU training data available, which for us is the DIP dataset [Huang et al. 2018] consisting real IMU measurements on 10 different actors, paired with ground-truth SMPL poses without root translation.

However, synthetic and real IMU data exhibit vastly different noise profiles. Acceleration data in the real dataset are noisy, but not in the same way as the noise in the synthetic dataset, which is caused by double differentiation of mocap data (Figure 4 Top). On the other hand, orientation data are usually less noisy because they are processed by the Kalman filter [Kalman et al. 1960]. Previous work [Huang et al. 2018] recognized this distribution mismatch problem and proposed to first train the model exclusively on the synthetic data and then fine-tune it on a smaller real dataset. This two-step solution leads to a more complex training procedure that requires careful tuning to avoid overfitting the real dataset.

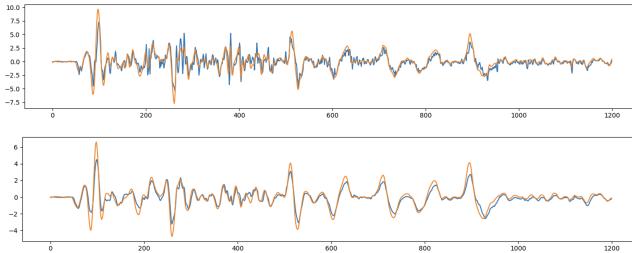


Fig. 4. Example of synthesized (orange) and real (blue) acceleration data, before (top) and after (bottom) moving average filtering.

We found that simply running an average filter on both synthetic and real acceleration data (with window length of 11) would bring

the two data sources sufficiently close to each other (Fig. 4 Bottom). We then train the model only once on the combined dataset without having to fine-tune. As shown in the result section, our model can perform well on both synthetic and real holdout sets.

In practice, filtering causes latency during real-time inference, as computing moving average requires future IMU readings. We use five times steps (83ms) of future readings, the same requirement as the state-of-the-art [Yi et al. 2021], though they require future readings as a part of the model input while we merely use them for filtering.

3.2 Predicting Stationary Body Points (SBP)

Predicting global (root) motion from IMUs is inherently difficult because, unlike optical-based motion capture, IMUs do not sense global position directly and the integration of noisy acceleration information often results in large drifting over time. Previous works designed special loss functions [Yi et al. 2021] and/or rely on heuristics to detect and enforce the current supporting foot [Rempe et al. 2021; Shimada et al. 2020; Yi et al. 2021]. All of them have demonstrated improvements, but the heuristics sometimes result in visual artifacts, such as a foot unnaturally "locked" on the ground, especially when applied to non-locomotion. These heuristics also do not handle non-flat terrains.

We propose a physics-inspired learning objective: predicting the *stationary body points (SBP)* c_t in human movements. As a most common example, when a person is in contact with a static environment, the contacting areas may have zero velocity due to friction forces, but the body links in contact could still be moving (Figure 5 illustrating heel-to-toe rolling contact during walking). Since our articulated rigid-body human model is non-deformable, we approximate the contacting area using a single point with zero velocity in the world frame, while allowing the rest of points in that body frame to move. The neural network is trained to predict the onsets and locations of such static points, $c_t = [b_t, r_t]$, where $b_t \in \{0, 1\}$ represents whether there exists an active stationary point, and if so, $r_t \in \mathbb{R}^3$ represents its local coordinate in the body frame.

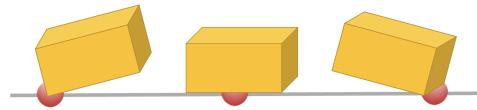


Fig. 5. During heel-to-toe contact in locomotion, contact patch is stationary but the foot link velocity is not zero.

To create labels for training, we need to detect the SBPs in the AMASS training dataset. Applying optimization approaches like least-squares, often run into numerical issues, such as under-determinacy or nonexistence of the solutions, due to noisy mocap data, as pointed out by [Le Callennec and Boulic 2006]. We instead take a sampling-based approach which samples N ($N \approx 1000$) points in each body frame that can potentially be in contact and evaluate all candidate points by:

$$l(\mathbf{r}) = \|\boldsymbol{\omega} \times \mathbf{R}_B \mathbf{r} + \mathbf{v}\| + 0.3 \|\mathbf{r} - \mathbf{r}_{t-1}\|, \quad (1)$$

where R_B is the body's global orientation, v and ω are the body's linear and angular velocity in the global frame respectively, and r_{t-1} is the solution in the previous frame, if existing. For each animation frame t in the AMASS dataset, we evaluate Equation 1 for every point candidate r_i in every body link k and choose the point $r^* = \operatorname{argmin}_i l(r_i)$ that has the smallest value. If $l(r^*)$ is less than a manually chosen threshold (0.15), we label $c_t(k) = [1, r^*]$. Otherwise, $c_t(k) = [0, (0, 0, 0)]$. Evaluation of Equation 1 for all candidates can be done in parallel efficiently using matrix operations.

Our current implementation assumes potential body links in contact to be only the feet, therefore $c_t \in \mathbb{R}^8$. This same methodology can be trivially extended to other body parts which we leave for future work. For example, we can include the pelvis link for drift mitigation during sitting, and the hand links for break-dancing motions as well. As the DIP real IMU data do not have root motion, we use a pre-trained model to label pseudo ground-truth SBPs for the DIP motions.

Using SBPs at run-time ("Contact Filtering"). We found that the learned transformer decoder works well on the AMASS dataset, but the performance on real IMU dataset can be further improved if we also use predicted c_t to correct the predicted root velocity v_t at run-time. If there is only one body link (i.e. one foot) that has an active stationary point, we adjust the root velocity by $\tilde{v}_t = v_t - v_{c(t)}$ so that the corrected \tilde{v}_t will make the velocity of predicted body point, $v_{c(t)}$, zero. If there are more than one link (in our implementation, both feet) that have an active stationary point, we randomly choose one to perform the same procedure. This simple correction method works for all the motion types in the AMASS dataset, but we expect that future work is required for more challenging cases.

3.3 Implementation Details

We use the AMASS dataset to generate synthetic training data following the smoothing procedure in Sec. 3.1. It consists of over a dozen different motion capture datasets performing a variety of activities. In addition, we include 8 out of 10 subjects' data from the DIP dataset. We use bullet [Coulmans and Bai 2016] for calculating forward kinematics during data synthesis, SBP label generation, contact filtering, and final visualizations.

We adopted a similar data calibration and normalization scheme as in TransPose [Yi et al. 2021], with the only difference that our model takes in pelvis orientation and acceleration directly in the global frame rather than character frame. We found this detail unimportant since the large training dataset covers all character headings relatively well. Note that our model requires an initial full-body pose given in the first step of prediction. In practice this is always the case since the sensors need to be calibrated with a T pose before each use, as they are allowed to be slightly differently worn.

We use standard loss functions for the model outputs, i.e., mean-squared error for joint rotations (represented as first two columns of the rotation matrix for unique and numerically stable ground-truth labels), mean-squared error for v_t and Cartesian elements of c_t , and binary cross-entropy for onset elements of c_t , (i.e. b_t). Since our model during training time predicts a whole trajectory window, we experimented with a jerk loss penalizing deviation of neighboring frames, but it did not produce visible improvements. This might be due to the fact that during test time we still only use the last

prediction at each step. Instead, we pass our output through an exponential moving average filter as post processing.

Our model is trained in PyTorch [Paszke et al. 2019] using the Adam optimizer [Kingma and Ba 2014], with a batch size of 256 and a learning rate of 0.0001 multiplied with a cosine schedule [Loshchilov and Hutter 2016]. We perform training for 1500 epochs, which takes around 8 hours with a GeForce GTX 2080Ti GPU. Once trained, our model is small enough to run at 60 fps on a 2080Ti machine. Our model contains a total number of 3,663,479 parameters, comparing to 4,798,771 in TransPose and 10,801,934 in DIP. Our source code will be released upon publication.

4 Results

Our experiments in this section demonstrate motion reconstruction results of our method on a variety of activities, as well as quantitative improvements on common metrics. They are best seen in the supplementary video. We also evaluate two key design choices in ablations, namely history feedback and stationary body point prediction. We conclude with a live demo and discussions on failure cases.

4.1 Evaluation

We describe the datasets and metrics used for evaluation, and show how our model performs in comparison to state-of-the-art methods.

Datasets We evaluate our model on both synthetic data and real data that cover a wide range of challenging scenarios.

- **Difficult categories from AMASS (synthetic training):** We randomly select 250 sequences from five motion categories in the AMASS dataset used for training: *parkour, dances and jumps, rolling, uneven terrain, losing balance*. They represent rare, long-tail actions in the dataset with non-cyclic movements and complex contacts or root dynamics.
- **DanceDB (synthetic heldout):** A large dataset of contemporary hip-hop dances unique to any other training dataset. Note that DanceDB is part of AMASS but we intentionally hold it out from our training data.
- **DIPEval (real heldout):** Data from two held-out subjects in the DIP dataset.
- **TotalCapture (real heldout):** We held out real IMU measurements from the TotalCapture dataset [Trumble et al. 2017] for evaluation, but still use its ground truth and synthesized IMU readings as part of the AMASS training set, following the same practice as previous works.

Metrics We define the following metrics that are commonly used to evaluate motion reconstruction quality. They are first computed for each frame of a motion, then averaged over all frames in each evaluation dataset.

- **Mean Joint Angle Error (in degrees):** Joint angle (represented in axis-angles) difference between reconstruction and ground-truth, averaged over all joints.
- **Mean Root-Relative Joint Position Error (in meters):** Joint position difference between the reconstruction and ground-truth by aligning at the root, averaged over all joints.
- **Root Error 2s/5s/10s (in meters):** Root translation error measured in l_2 norm during a continuous period of 2s/5s/10s.

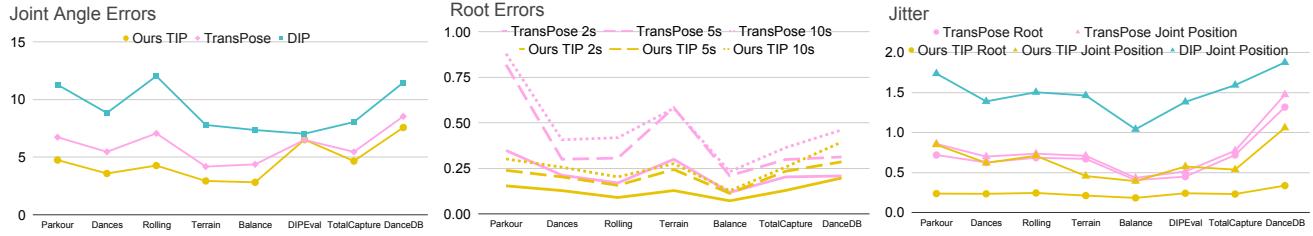


Fig. 6. Plots of metrics in Table 1 comparing our DIP model, TransPose model and DIP model.

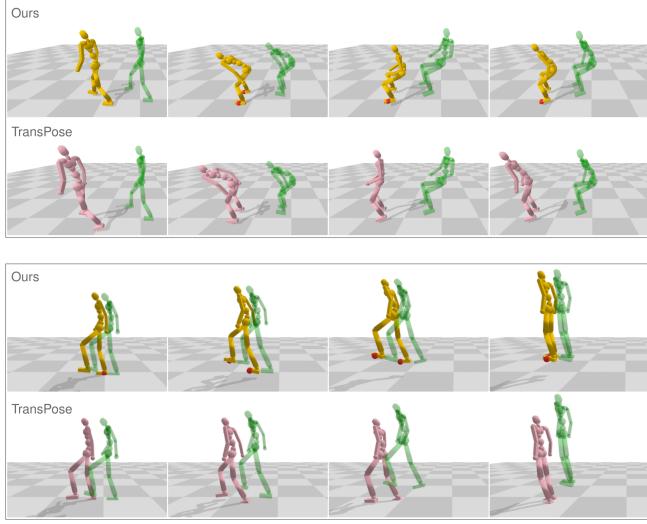


Fig. 7. Motion reconstruction for sitting on a chair (top) and climbing steps (bottom). Our character is shown in yellow, TransPose in purple and Ground-Truth motion is shown in green. The red spheres are predicted SBPs.

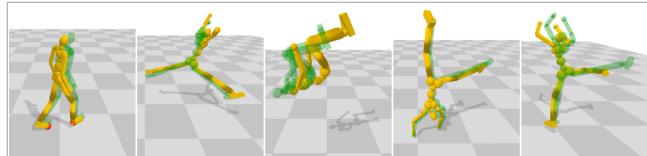


Fig. 8. Our model (yellow) can track a variety of motions, including dynamic motions and hand contacts, such as jumping and break dancing.

- **Mean Joint Position Jitter (in m/s^2):** Joint position jitter computed using the same formula as in TransPose, averaged over all joints.
- **Root Jitter (in m/s^2):** Root position jitter computed using the same formula as above.

Experiments We present quantitative metrics on the evaluation datasets between our model, TransPose, and DIP in Table 1. We used the best performing models published by the authors in this comparison. Overall, our small and simple model achieves better accuracy and lower jitter in almost all evaluations.

Our auto-regressive decoder generates consistent motions when the input is ambiguous. For example, IMU readings are similar between sitting or standing still, since all accelerations are close to

zero and the limb sensors are similarly oriented. Both TransPose and DIP have a hard time distinguishing these two actions, and produce unstable motions that constantly switch between them. On the other hand, we can generate a stable sitting posture based on the prediction history feedback (Fig. 7 Top).

We reduced root position errors over both short and long duration for all categories, oftentimes by more than 30% (Fig. 6 Middle). We found that enforcing SBPs at runtime has a noticeable improvement for real data. Our root motions are also much smoother (Fig. 6 Right), thanks to the smoothing filter on root prediction, at the expense of slightly larger joint errors caused by a small filter-window delay.

Since we made no assumptions about the environment or contact patterns, our model greatly improves the root motion accuracy when walking on uneven terrains (Fig 7 Bottom), or when performing dynamic actions such as jumping or break dancing (Fig 8), where we observed long duration without active SBPs to enable a fluid motion. In addition, we found that the predicted SBPs coincide with environment contacts and follows the heel-toe rolling pattern.

4.2 Ablations

We perform two ablations to validate our key design choices: (1) without history feedback; (2) without SBP prediction. Table 2 summarizes the results on the TotalCapture dataset.

As expected, when we removed the history feedback from model input, we observed the same unstable prediction for ambiguous input as existing work. We can no longer generate a stable sitting posture, but instead a jittery motion that constantly switches between sitting and standing (ref. Video). The lack of history feedback from SBP predictions also has a big negative impact to root accuracy: when SBPs are inconsistent across frames, they are not nearly as useful in correcting for drifts. On the other hand, we observed a slight improvement in joint angle accuracy on average without history feedback. This can be partly explained by the discrepancy in history input between training and testing, where we use ground truth for training but the model output for testing. Another potential reason is that history feedback could also lead to persistent mistakes. If the model makes a wrong prediction due to ambiguous input, it would be difficult to recover later.

The model without predicting SBPs performs significantly worse especially on long-term root trajectory accuracy. This is expected because we can no longer adjust the root motion to prevent drifts. We observed a greater effect on real datasets than on synthetic data, probably because v_t predictions are less robust on noises, while the SBPs are more resistant to input noise.

Table 1. Comparison of model performance on evaluation datasets. Bold numbers indicate the best performing entries.

Our TIP Model								
	Parkour	Dances	Rolling	Terrains	Balance	DIPEval	TotalCapture	DanceDB
joint angle errors (degree)	4.74376	3.57834	4.26563	2.9291	2.81225	6.51343	4.65692	7.56777
joint position errors (meter)	0.06167	0.04078	0.04858	0.03173	0.03139	0.06687	0.05271	0.08059
root errors in 2s (meter)	0.15387	0.1277	0.08986	0.128	0.07187		0.12865	0.19691
root errors in 5s (meter)	0.23898	0.20274	0.15678	0.24366	0.11263		0.23402	0.28554
root errors in 10s (meter)	0.30146	0.25598	0.20274	0.27567	0.12658		0.25683	0.39387
joint position jitter (m/s^2)	0.84803	0.62009	0.70928	0.4557	0.3925	0.57487	0.53599	1.0587
root jitter (m/s^2)	0.23667	0.23384	0.24435	0.21186	0.18294	0.2408	0.2311	0.33655

TransPose Model								
	Parkour	Dances	Rolling	Terrains	Balance	DIPEval	TotalCapture	DanceDB
joint angle errors (degree)	6.72435	5.45769	7.05947	4.17613	4.36833	6.49189	5.45348	8.53172
joint position errors (meter)	0.06784	0.05044	0.07002	0.03388	0.03684	0.06331	0.05419	0.08039
root errors in 2s (meter)	0.34825	0.21207	0.16982	0.2995	0.11767		0.20273	0.20812
root errors in 5s (meter)	0.81832	0.29996	0.30606	0.58249	0.20944		0.29827	0.31241
root errors in 10s (meter)	0.87868	0.40665	0.41824	0.57941	0.23171		0.36359	0.46058
joint position jitter (m/s^2)	0.85816	0.69872	0.73561	0.70862	0.43393	0.51272	0.77223	1.47254
root jitter (m/s^2)	0.71842	0.62605	0.68413	0.67024	0.40534	0.44845	0.71907	1.31713

DIP Model								
	Parkour	Dances	Rolling	Terrains	Balance	DIPEval	TotalCapture	DanceDB
joint angle errors (degree)	11.28364	8.83262	12.03162	7.78829	7.34856	7.02777	8.04311	11.46666
joint position errors (meter)	0.15187	0.09099	0.15029	0.07274	0.06719	0.07331	0.08756	0.12303
joint position jitter (m/s^2)	1.73595	1.38855	1.5027	1.46214	1.04003	1.38213	1.59157	1.87524

Table 2. Comparison of ablation models on TotalCapture evaluation dataset.

	No History	No SBPs	TIP (w/ both)
joint angle errors	4.23954	4.77988	4.65692
joint position errors	0.04934	0.05429	0.05271
root errors in 2s	0.15833	0.16987	0.12865
root errors in 5s	0.28575	0.32543	0.23402
root errors in 10s	0.41114	0.5051	0.25683
joint position jitter	0.68268	0.48611	0.53599
root jitter	0.26743	0.01809	0.2311

4.3 Live Demo

We test our system live with 6 Xsens IMU sensors. Our video visualizes live performance side-by-side with real-time reconstructions, with a slight latency caused by our pre-processing filter. We cover a variety of motion tasks in our demo, both commons ones such as locomotion and whole-body manipulation, and more challenging highly-dynamical ones such as jumping from a high place, "swimming" on a stool, rolling on the floor, or swirl kicks. We tested our system on one male and one female subjects, and observed degraded performance on the female subject which we did not expect. We will discuss this in more detail in Sec. 5.

4.4 Failure Cases

Our system still has a few drawbacks. Even though SBPs can help correct drifts in the root motion, they are susceptible to drift themselves because they are relative to the previous frame and not to a static global reference. The effect is more pronounced in vertical

motion, which could be more distracting than errors in horizontal directions. For example, in a jumping jack sequence, the character appears to slowly sink to the floor on every downward strike. Similarly, we also observed larger root translation errors in fast dancing motions since the predicted SBPs are not always static. If we have access to environment geometry information, we may be able to combine SBPs and contact detection to mitigate these cases.

5 Discussions and Conclusion

This paper presents a new data-driven method for human motion reconstruction from six wearable IMUs. By utilizing recent advances in sequence modeling with attention-based neural networks, combined with a physics-inspired task of learning stationary body points to mitigate out-of-distribution collapses of neural networks, we make one principled step forward in tackling this real-world, under-constrained problem where noises and unmodeled dynamics constantly impact systems' performance and real-world usability.

Though we have shown clear improvement on existing challenges of temporal consistency due to ambiguity, dynamic motion coverage, and terrain coverage, there is still room for improvements. The most visible problem we observe is the model's bias to body types. Our synthesized data were generated from a virtual character with a fixed height at 1.6m. We observed that the algorithm generalizes better to taller users than shorter ones. We hypothesize that this phenomenon is due to the magnitude of acceleration, as the model might be more easily confused by smaller signals from a shorter user. Similarly, existing real IMU datasets might have a bias in human

shapes. Some personalized training and fine-tuning of the model may eventually be necessary for reconstructing more accurate and detailed motion for each individual user.

Another potential future direction is to estimate the geometry of the environment from human motion. Terrain height estimation remains challenging since they are very sensitive to different body types and noises (for example, our algorithm could predict a locomotion on a slightly bumpy ground while the user really is moving on flat surface). In practice, we found that the SBP corrections are not as helpful in the vertical direction than they are in the horizontal directions, and will address this issue in future work. Finally, the success of highly expressive Transformer models depend on the availability of large-scale data. Compared with commonly used training datasets for natural language modeling, human motion databases are much smaller in size and variation, causing less-ideal coverage for rare motions. Learning a model that can generalize to rarely seen motions remain a challenge for our method.

References

- [n.d.]. Rokoko <https://www.rokoko.com/>. Last visited: 01/26/2022.
- [n.d.]. Vicon Motion Systems <https://www.vicon.com/>. Last visited: 01/26/2022.
- [n.d.]. Xsens <https://www.xsens.com/>. Last visited: 01/26/2022.
- Emre Aksan, Manuel Kaufmann, Peng Cao, and Otmar Hilliges. 2021. A Spatio-temporal Transformer for 3D Human Motion Prediction. *International Conference on 3D Vision (3DV)* (2021).
- Sheldon Andrews, Ivan Huerta, Taku Komura, Leonid Sigal, and Kenny Mitchell. 2016. Real-Time Physics-Based Motion Capture with Sparse Sensors. In *Proceedings of the 13th European Conference on Visual Media Production (CVMP 2016) (CVMP 2016)*. Article 5.
- Eric R. Bachmann, Robert B. McGhee, Xiaoping Yun, and Michael J. Zyda. 2001. Inertial and Magnetic Posture Tracking for Inserting Humans into Networked Virtual Environments. In *Proceedings of the ACM Symposium on Virtual Reality Software and Technology (VRST '01)*. 9–16.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. 2018. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271* (2018).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. (2020). arXiv:2005.14165 [cs.CL]
- Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. 2019. OpenPose: Real-time Multi-Person 2D Pose Estimation using Part Affinity Fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019).
- Young-Woon Cha, Husan Shaik, Qian Zhang, Fan Feng, Andrei State, Adrian Ilie, and Henry Fuchs. 2021. Mobile. Ego-centric Human Body Motion Reconstruction Using Only Eyeglasses-mounted Cameras and a Few Body-worn Inertial Sensors. In *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*.
- Jinxiang Chai and Jessica K. Hodgins. 2005. Performance Animation from Low-Dimensional Control Signals. *ACM Trans. Graph.* 24, 3 (2005).
- Adrian Hilton Charles Malleson, John Collomosse. 2020. Real-Time Multi-person Motion Capture from Multi-view Video and IMUs. *International Journal of Computer Vision* 128 (06 2020).
- Vasileios Choutas, Federica Bogo, Jingjing Shen, and Julien Valentin. 2021. Learning to Fit Morphable Models. *CoRR* abs/2111.14824 (2021). arXiv:2111.14824 <https://arxiv.org/abs/2111.14824>
- Erwin Coumans and Yunfei Bai. 2016. Pybullet, a python module for physics simulation for games, robotics and machine learning. (2016).
- Michael B. Del Rosario, Heba Khamis, Phillip Ngo, Nigel H. Lovell, and Stephen J. Redmond. 2018. Computationally Efficient Adaptive Error-State Kalman Filter for Attitude Estimation. *IEEE Sensors Journal* 18, 22 (2018), 9332–9342.
- Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. 2020. Jukebox: A Generative Model for Music. arXiv:2005.00341 [eess.AS]
- Andrea Dittadi, Sebastian Dziadzio, Darren Cosker, Ben Lundell, Tom Cashman, and Jamie Shotton. 2021. Full-Body Motion From a Single Head-Mounted Device: Generating SMPL Poses From Partial Observations. In *International Conference on Computer Vision 2021*.
- E. Foxlin. 1996. Inertial head-tracker sensor fusion by a complementary separate-bias Kalman filter. In *Proceedings of the IEEE 1996 Virtual Reality Annual International Symposium*. 185–194.
- Andrew Gilbert, Matthew Trumble, Charles Malleson, Adrian Hilton, and John Collomosse. 2019. Fusing Visual and Inertial Sensors with Semantics for 3D Human Pose Estimation. *International Journal of Computer Vision* 127 (04 2019), 1–17.
- Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. 2018. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7297–7306.
- Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Christian Theobalt. 2013. Real-Time Body Tracking with One Depth Camera and Inertial Sensors. In *2013 IEEE International Conference on Computer Vision*. 1105–1112.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J. Black, Otmar Hilliges, and Gerard Pons-Moll. 2018. Deep Inertial Poser: Learning to Reconstruct Human Pose from Sparse Inertial Measurements in Real Time. *ACM TOG* 37, 6 (12 2018).
- Rudolph Emil Kalman et al. 1960. A new approach to linear filtering and prediction problems [J]. *Journal of basic Engineering* 82, 1 (1960), 35–45.
- Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. 2019. Learning 3D Human Dynamics from Video. In *Computer Vision and Pattern Recognition (CVPR)*.
- Manuel Kaufmann, Yi Zhao, Chengcheng Tang, Lingling Tao, Christopher Twigg, Jie Song, Robert Wang, and Otmar Hilliges. 2021. EM-POSE: 3D Human Pose Estimation from Sparse Electromagnetic Trackers. In *International Conference on Computer Vision (ICCV)*.
- Seong Uk Kim, Hanyoung Jang, Hyeyounseung Im, and Jongmin Kim. 2021. Human motion reconstruction using deep transformer networks. *Pattern Recognition Letters* 150 (2021), 162–169.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Benoit Le Calennec and Ronan Boulic. 2006. Robust kinematic constraint detection for motion data. In *Proceedings of the 2006 ACM SIGGRAPH/Eurographics symposium on Computer animation*. 281–290.
- Rui long Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. 2021. AI Choreographer: Music Conditioned 3D Dance Generation with AIST++.
- Huajun Liu, Xiaolin Wei, Jinxiang Chai, Inwoo Ha, and Taehyun Rhee. 2011. Realtime Human Motion Control with a Small Number of Inertial Sensors. In *Symposium on Interactive 3D Graphics and Games (IS3D '11)*. 133–140.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM TOG* 34, 6 (Oct. 2015), 248:1–248:16.
- Ilya Loschilov and Frank Hutter. 2016. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983* (2016).
- Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. 2021. Dynamics-Regulated Kinematic Policy for Egocentric Pose Estimation. In *NeurIPS*.
- Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. 2019. AMASS: Archive of Motion Capture as Surface Shapes. In *ICCV*. 5442–5451.
- Charles Malleson, Marco Volino, Andrew Gilbert, Matthew Trumble, John Collomosse, and Adrian Hilton. 2017. Real-time Full-Body Motion Capture from Video and IMUs. In *Int. Conf. 3D Vis.*.
- Deepak Nagaraj, Erik Schake, Patrick Leiner, and Dirk Werth. 2020. An RNN-Ensemble Approach for Real Time Human Pose Estimation from Sparse IMUs. In *Proceedings of the 3rd International Conference on Applications of Intelligent Systems (Las Palmas de Gran Canaria, Spain) (APPIIS 2020)*. Article 32, 6 pages.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. PyTorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- Mathis Petrovich, Michael J. Black, and Guil Varol. 2021. Action-Conditioned 3D Human Motion Synthesis with Transformer VAE. In *International Conference on Computer Vision (ICCV)*. 10985–10995.
- Gerard Pons-Moll, Andreas Baak, Juergen Gall, Laura Leal-Taixé, Meinard Müller, Hans-Peter Seidel, and Bodo Rosenhahn. 2011. Outdoor human motion capture using inverse kinematics and von mises-fisher sampling. In *2011 International Conference on Computer Vision*. 1243–1250.
- Gerard Pons-Moll, Andreas Baak, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bodo Rosenhahn. 2010. Multisensor-fusion for 3D full-body human motion capture. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. 663–670.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. (2018).
- Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. 2021. HuMoR: 3D Human Motion Model for Robust Pose Estimation. In *ICCV*.

- Qaiser Riaz, Guanhong Tao, Björn Krüger, and Andreas Weber. 2015. Motion Reconstruction Using Very Few Accelerometers and Ground Contacts. *Graph. Models* 79, C (may 2015), 23–38.
- Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. 2020. PhysCap: Physically Plausible Monocular 3D Motion Capture in Real Time. *ACM TOG* 39, 6 (12 2020).
- Ronit Slyper and Jessica K. Hodgins. 2008. Action Capture with Accelerometers. In *Proceedings of the 2008 ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA '08)*. 193–199.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The Journal of Machine Learning Research* 15, 1 (2014), 1929–1958. <http://dl.acm.org/citation.cfm?id=2627435.2670313>
- Jochen Tautges, Arno Zinke, Björn Krüger, Jan Baumann, Andreas Weber, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bernd Eberhardt. 2011. Motion Reconstruction Using Sparse Accelerometer Data. *ACM Trans. Graph.* 30, 3 (2011).
- Jonathan Taylor, Jamie Shotton, Toby Sharp, and Andrew Fitzgibbon. 2012. The Vitruvian manifold: Inferring dense correspondences for one-shot human pose estimation. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*. 103–110.
- Matt Trumble, Andrew Gilbert, Charles Malleson, Adrian Hilton, and John Collomosse. 2017. Total Capture: 3D Human Pose Estimation Fusing Video and Inertial Sensors. In *BMVC*.
- Guillermo Valle-Pérez, Gustav Eje Henter, Jonas Beskow, Andre Holzapfel, Pierre-Yves Oudeyer, and Simon Alexanderson. 2021. Transflower: Probabilistic Autoregressive Dance Generation with Multimodal Attention. *ACM Trans. Graph.* 40, 6, Article 195 (dec 2021), 14 pages. <https://doi.org/10.1145/3478513.3480570>
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, Vol. 30.
- Rachel V. Vitali, Ryan S. McGinnis, and Noel C. Perkins. 2021. Robust Error-State Kalman Filter for Estimating IMU Orientation. *IEEE Sensors Journal* 21, 3 (2021), 3561–3569.
- Daniel Vlasic, Rolf Adelsberger, Giovanni Vannucci, John Barnwell, Markus Gross, Wojciech Matusik, and Jovan Popović. 2007. Practical Motion Capture in Everyday Surroundings. *ACM Trans. Graph.* 26, 3 (2007).
- Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. 2018. Recovering Accurate 3D Human Pose in The Wild Using IMUs and a Moving Camera. In *European Conference on Computer Vision (ECCV)*.
- Timo von Marcard, Gerard Pons-Moll, and Bodo Rosenhahn. 2016. Human Pose Estimation from Video and IMUs. *Transactions on Pattern Analysis and Machine Intelligence (PAMI)* (jan 2016).
- Timo von Marcard, Bodo Rosenhahn, Michael Black, and Gerard Pons-Moll. 2017. 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), 349–360.
- Xiaolin Wei, Peizhao Zhang, and Jinxiang Chai. 2012. Accurate Realtime Full-Body Motion Capture Using a Single Depth Camera. *ACM Trans. Graph.* 31, 6, Article 188 (nov 2012).
- Dongseok Yang, Doyeon Kim, and Sung-Hee Lee. 2021. LoBSTR: Real-time Lower-body Pose Prediction from Sparse Upper-body Tracking Signals. *Computer Graphics Forum* (2021). <https://doi.org/10.1111/cgf.142631>
- Xinyu Yi, Yuxiao Zhou, and Feng Xu. 2021. TransPose: Real-time 3D Human Translation and Pose Estimation with Six Inertial Sensors. *ACM TOG* 40, 4 (8 2021).
- Zhe Zhang, Chunyu Wang, Wenhui Qin, and Wenjun Zeng. 2020. Fusing Wearable IMUs with Multi-View Images for Human Pose Estimation: A Geometric Approach. In *CVPR*.
- Zerong Zheng, Tao Yu, Hao Li, Kaiwen Guo, Qionghai Dai, Lu Fang, and Yebin Liu. 2018. HybridFusion: Real-Time Performance Capture Using a Single Depth Sensor and Sparse IMUs. In *European Conference on Computer Vision (ECCV)*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (Eds.). 389–406.