

IMDY: HUMAN INVERSE DYNAMICS FROM IMITATED OBSERVATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Inverse dynamics (ID), which aims at reproducing the driven torques from human kinematic observations, has been a critical tool for human motion analysis. However, it is hindered from wider application to general motion due to its limited scalability. Conventional optimization-based ID requires expensive laboratory setups, restricting its availability. To alleviate this problem, we propose to exploit the recently progressive human motion imitation algorithms to learn human inverse dynamics in a *data-driven* manner. The key insight is that the human ID knowledge is implicitly possessed by motion imitators, though not directly applicable. In light of this, we devise an efficient data collection pipeline with state-of-the-art motion imitation algorithms and physics simulators, resulting in a large-scale human inverse dynamics benchmark as **Imitated Dynamics (ImDy)**. ImDy contains over **150 hours** of motion with joint torque and full-body ground reaction force data. With ImDy, we train a data-driven human inverse dynamics solver **ImDyS(olver)** in a fully supervised manner, which conducts ID and ground reaction force estimation simultaneously. Experiments on ImDy and real-world data demonstrate the impressive competency of ImDyS in human inverse dynamics and ground reaction force estimation. Moreover, the potential of ImDy(-S) as a fundamental motion analysis tool is exhibited with downstream applications. Our data and code would be made publicly available.

1 INTRODUCTION

The rapid progress in human motion capture based on computer vision has made an enormous amount of human motion data available to the research community (Mahmood et al., 2019; Mandery et al., 2016). The accumulation of human motion manages to push motion understanding forward in various tasks, including behavior understanding (Punnakkal et al., 2021; Shahroudy et al., 2016) and character animation (Guo et al., 2022; Tevet et al., 2023). However, given the vision-based nature, most current efforts focus only on visible kinematics information. The invisible factors, especially the dynamic factors, which could carry deeper insights into the underlying production mechanism of human motion, are typically overlooked, such as *driven torques* and *ground reaction forces*. This limits the current motion understanding algorithms from wider applications to domains where physical constraints must be seriously considered, such as **robotics** (Figueredo et al., 2020; Teramae et al., 2017), **healthcare** (Yao et al., 2018), and **sports training** (Caruntu & Moreno, 2019). To alleviate this, we focus on identifying the driven torques and ground reaction forces for human motion from pure kinematics MoCap data, known as human inverse dynamics (ID).

Human inverse dynamics, as a basic step toward physical motion modeling, has been extensively discussed by the biomechanics community for applications like gait analysis. A fundamental obstacle is that it could not be measured non-intrusively. Therefore, computationally expensive optimization-based methods are widely adopted and mature software is developed (Delp et al., 2007; Damsgaard et al., 2006; Werling et al., 2021). However, accurately measured ground reaction forces are required to ensure a determinate solution, which could be expensive and applicable only in restricted laboratory settings. Also, the optimization process could be sensitive to small disturbances in either motion capture noises or subject variances. These make it hard to scale up for wider applications to general motion. Given the success achieved by data-driven methods in CV and NLP, deep-learning-based methods are proposed (Zell & Rosenhahn, 2015; Zell et al., 2017; Lv et al., 2016), aiming at scalable human inverse dynamics with only kinematic observations as inputs.

Unfortunately, *data acquisition* becomes a major bottleneck since laboratory setups are still required for ground-truth acquisition.

Given this, we project our sights on the recent progress of Imitation Learning (IL) (Luo et al., 2021; 2023), which replicates recorded human motion through fully simulated humanoids with physical control signals, namely, joint torques. A key insight is that with the goal of kinematics phenomenon imitation, IL might also *implicitly* imitate the dynamics production mechanism, known as ID. However, IL is not directly applicable to ID. Despite the visual resemblances between the recorded and simulated motion, kinematic errors still exist. These errors could be neglected for kinematic analyses, however, for dynamic analysis, they could be amplified drastically (Uchida & Seth, 2022). Moreover, existing successful IL algorithms are typically based on joint-actuated SMPL (Loper et al., 2015)

avatars, whose physical properties and topology differ from real humans. To this end, extracting ID knowledge from IL becomes critical. Here, we adopt the state-of-the-art motion IL algorithm (Luo et al., 2023) and physics simulator (Makoviychuk et al., 2021) to imitate recorded motions, extracting the observed kinematic states, joint torques, and the ground reaction forces, resulting in a large-scale human inverse dynamics database named **Imitated Dynamics** (ImDy) with more than 150-hour human motion. There are two major merits of ImDy. First, it is *scalable*. Multiple samples could be concurrently collected in the simulator without expensive laboratory setups, extending the border of ID data acquisition. As shown in Fig. 1, we could even pair some rather complex motions with ID data, which is hard to achieve in laboratories. Second, it is *holistic*. Beyond the ground reaction force and ID typically recorded in laboratories for previous efforts (Zell et al., 2020; Mourot et al., 2022; Han et al., 2023), the physics simulator enables us to access the GRFs and joint torques of all human body segments, as shown in Fig. 1.

With the accumulated data, we could address the human inverse dynamics in a fully supervised manner. Given the observed kinematics states that describe a motion transition in a certain period, we train a data-driven solver as ImDyS(olver) to estimate the ground reaction forces and the internal dynamics to drive the transition. We also devise losses to regulate ImDyS with forward dynamics awareness and motion plausibility constraints.

We demonstrate the efficacy of ImDyS through a wide span of experiments. First, we evaluate our method on ImDy for a basic performance illustration with simulated ImDy. Then ImDyS is evaluated on GroundLink (Han et al., 2023), which contains real-world ground reaction force. Furthermore, we demonstrate the efficacy of ImDy on the recent real-world human dynamics dataset AddBiomechanics (Werling et al., 2024).

Our contribution could be summarized as: (1) We propose a novel pipeline for human inverse dynamics data collection, introducing a large-scale benchmark as ImDy. (2) Based on ImDy, a data-driven ID solver is instantiated as ImDyS. (3) Extensive experiments are conducted with analyses of the proposed data-driven methodology, demonstrating the feasibility of ImDyS.

2 BACKGROUND

Conventional Inverse Dynamics. Inverse dynamics, known as inferring forces/moments from kinematic observations, have been discussed for long in the biomechanics community. In this literature, it is formulated as an optimization problem: given a representative model of a subject, the joint kinematics over time w.r.t. the subject model, and the external forces, find the driving torques that produce the motion (Uchida & Delp, 2021). The Newtonian dynamic equations are involved as

$$M(q)\ddot{q} + C(q, \dot{q}) + G(q) = J\lambda + \tau, \quad (1)$$

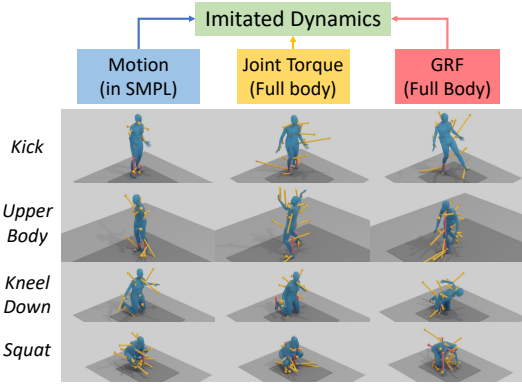


Figure 1: ImDy pairs diverse SMPL motion data with dynamics including full-body torques and ground reaction forces (GRF) like the right knee GRF for kneeling, which could be hard to achieve under conventional laboratory setups.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

where $M(q)$ is the generalized human inertia matrix w.r.t. generalized coordinate q , $C(q, \dot{q})$ is the Coriolis and centrifugal forces, $G(q)$ represents gravity, J is the Jacobian matrix mapping external forces λ to the generalized coordinates. Thus, the driven torques τ could be obtained by minimizing the difference between the left and right terms of Eq. 1. Mature software based on this has been developed like OpenSim (Delp et al., 2007), AnyBody (Damsgaard et al., 2006), and Nimble (Werling et al., 2021). In addition, many efforts are made for clinical motion analysis (Fukuchi et al., 2018; Schreiber & Moissenet, 2019). However, these efforts are not as extensively recognized by the computer vision and computer graphics community as expected due to the scalability issue. Despite the elegant formulation, the efficacy of optimization-based heavily relies on the quality of external force λ (like GRF) measurement, whose cost could be non-trivial. Therefore, most of them focused on limited motion in laboratory settings. Some resort to wearable devices (Latella et al., 2016; 2019) to partially mitigate the limitation. In addition, fitting the raw captured kinematic observations to a specific human model for joint kinematics could be time-consuming and unstable, even with recent progress on it (Keller et al., 2023; Werling et al., 2023).

Learning-based Inverse Dynamics. With the progress in deep learning, there have been efforts to adopt neural networks to address the human ID problem. Many efforts focus on lower-body-only (Johnson & Ballard, 2014; Xiong et al., 2019) or upper-body-only (Manukian et al., 2023) inverse dynamics. More recently, Lv et al. (2016) collected over 1 hour of motion with an optical MoCap system, four force plates, and a pair of pressure insoles. The ground truth was obtained through optimization and a Gaussian mixture framework was devised. Zell & Rosenhahn (2015); Zell et al. (2017); Zell & Rosenhahn (2017) introduced a predictive dynamics-based human modeling for the acquisition of ground truth. Hundreds of motions were collected and different data-driven techniques were adopted for joint torque regression. Zell et al. (2020) proposed a weakly supervised method based only on motion for gait analysis. These efforts were constrained by costly data acquisition in real-world scenarios, resulting in limited data scale. Very recently, Werling et al. (2024) aggregated multiple existing biomechanics datasets, considerably boosting the data scale. However, most of the collected sequences contained only regular exercise motion with limited diversity. Some efforts focused on ground reaction forces such as (Rempe et al., 2020; Scott et al., 2020), UnderPressure (Mourot et al., 2022), and GroundLink (Han et al., 2023). Some recent works incorporated inverse dynamics into vision-based markerless MoCap systems. Shimada et al. (2021) and Li et al. (2022) simultaneously captured motion and joint torques with customized fully differentiated pipelines. A series of works (Yi et al., 2022; Gartner et al., 2022; Gärtner et al., 2022; Huang et al., 2022; Wang et al., 2023) imitated the captured motion in physical simulators with PD controllers and obtained the torques. However, an inherent problem is the amplification effect from kinematic errors to dynamic errors. As measured by Uchida & Seth (2022), only a 2-cm uncertainty of marker placement in a marker-based MoCap system could result in a peak ankle plantarflexion moment of $26.6 N \cdot m$. Considering the precision of current markerless MoCap algorithms, the accuracy of the accompanied inverse dynamics could be questionable. Also, among all these efforts for learning-based inverse dynamics, only a few (Zell et al., 2017; Zell & Rosenhahn, 2017; Zell et al., 2020) were quantitatively evaluated with limited locomotion data. A scalable benchmark for learning-based inverse dynamics is still not available.

Motion Imitation. IL for human motion replicates recorded human motion sequences with physically controlled simulated characters, which could be inherently close to ID. Most early efforts focus on specified usages with limited generalizability (Bergamin et al., 2019; Peng et al., 2021; Won et al., 2021; 2022; Peng et al., 2022). With residual force control (Yuan & Kitani, 2019), which imposed supernatural forces at the root joint of the humanoid, Luo et al. (2021) generalized to 97% sequences in AMASS (Mahmood et al., 2019). Luo et al. (2023) eliminated the supernatural root force and achieved a 98.9% success rate on AMASS with fall-state recovery. The progress in human motion IL makes it possible to collect human-like motions with full dynamics, shedding new light on the scalable human ID data collection.

3 CONSTRUCTING IMITATED DYNAMICS

ImDy aims to exploit the inherent closeness of inverse dynamics and imitation learning. Generally, the inverse dynamics (ID) and imitation algorithms (IL) could be abstracted as

$$\tau = ID(s_o^t, s_o^{t+1}), \tau = IL(s_o^t, s_i^{t+1}), \quad (2)$$

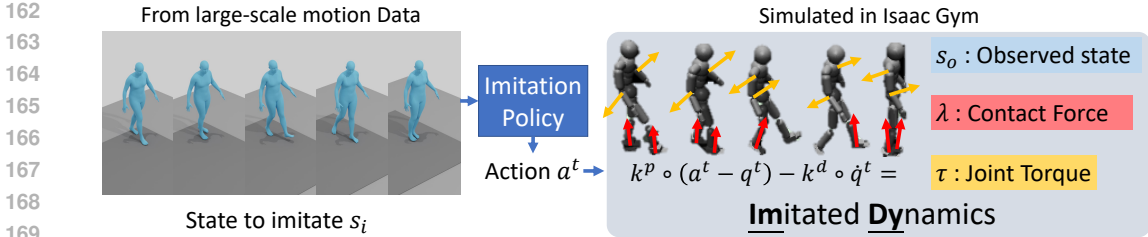


Figure 2: ImDy construction. We first train a motion imitation policy following Luo et al. (2023). Then, the policy is adopted to imitate arbitrary motions, with the imitated states recorded as ImDy.

with driven torque τ , timestamp t , observed kinematic states s_o , and the state to imitate s_i . Both ID and IL learn the dynamic production mechanism of human motion. However, IL algorithms are not directly applicable to ID due to the non-equivalence between s_o^{t+1} and s_i^{t+1} . The errors in kinematics could be magnified in dynamics (Uchida & Seth, 2022). This also makes ID algorithms that are deeply coupled with markerless MoCap less reliable. However, it is possible to extract knowledge from IL for ID. In this section, we introduce a simple but effective ID data collection pipeline with IL algorithms. First, the adopted IL algorithm (Luo et al., 2023) is briefly covered in Sec. 3.1. Then, the data collection pipeline is introduced in Sec. 3.2. An overview is given in Fig. 2.

3.1 IMITATION LEARNING BASICS

A motion imitator $\pi(a^t | s_o^t, s_i^t)$ is trained following Luo et al. (2023) to solve the Markov Decision Process $\mathcal{M} = \langle \mathcal{T}, \mathcal{S}, \mathcal{A}, \mathcal{R}, \gamma \rangle$. The transition dynamics \mathcal{T} and states \mathcal{S} are governed by the physics simulator. For each timestamp t , the policy π produces action $a^t \in \mathcal{A}$ and the reward \mathcal{R} , based on state $s \in \mathcal{S}$. The training goal is maximizing the reward expectation $\mathbf{E}(\sum_{t=1}^T \gamma^{t-1} r^t)$.

Transition. IsaacGym (Makoviychuk et al., 2021) is adopted for simulation. A 24-joint humanoid with SMPL (Loper et al., 2015) kinematics and physical properties following Luo et al. (2021; 2023) is adopted with variable shape parameter $\beta \in \mathbb{R}^{10}$. Thus, a human pose at timestamp t could be defined as $q^t = \{\theta^t, p^t\}$, where $\theta^t \in \mathbb{R}^{J \times 6}$ is the joint rotation in the 6d representation (Zhou et al., 2019) and $p^t \in \mathbb{R}^{J \times 3}$ is the 3D joint position.

State. At timestamp t , s^t contains the observed s_o^t and s_i^{t+1} to imitate. s_o^t is defined in simulator as $s_o^t = (q_t, \dot{q}_t, \beta)$ with 3D body pose q_t , velocity \dot{q}_t , and body shape β . s_i^{t+1} is defined similarly except that it is the reference motion with finite-differentiated velocities.

Action. All joints but the pelvis are actuated with proportional derivative (PD) controllers, with a^t as the PD target. The torque applied could be calculated as

$$\tau^t = k^p \circ (a^t - q^t) - k^d \circ \dot{q}^t. \quad (3)$$

Reward. The reward is composed of four terms: motion imitation reward for minimizing the difference between the imitated states and the expected states, fail-state recovery reward (Luo et al., 2023), AMP reward (Peng et al., 2021), and energy reward to reduce jittering.

Training. Following PHC, three primitive policies are progressively trained with hard negative mining, two for pure motion imitation, and one for fail-state recovery. Then, a composer learns to combine the primitives dynamically. PPO (Schulman et al., 2017) is adopted to train the policies.

3.2 IMITATED DATA ACQUISITION

With the imitator π , we pursue to extract its inherent ID knowledge. As in Eq. 2, though the imitator-produced τ is not accurate for $s_o^t \rightarrow s_i^{t+1}$ since s_i^{t+1} is not guaranteed to reach, τ is accurate for $s_o^t \rightarrow s_o^{t+1}$. Thus, the idea could be as simple as using π to imitate arbitrary motions in the simulator, then collecting all the **observed** states s_o , the applied torques τ , and the full-body GRF λ .

We adopt AMASS (Mahmood et al., 2019) and KIT (Krebs et al., 2021) as two major data sources. Sequences involving humans interacting with objects other than the ground are excluded, resulting in over 50 hours of motion. All the sequences are re-sampled to 30FPS, with the z-axis as the

Table 1: ImDy compared to related human dynamics datasets. Zell et al. (2020) recorded full-body data but simplified the upper body with a single torso segment. All previous efforts contain only GRF for feet (indicated with *), while we include full body GRF.

Dataset	#Subj	Duration (h)	Dynamics	Body Repr.	Style
Zell et al. (2020)	22	0.07	GRF* & Torques	Partial Skeleton*	Real
Scott et al. (2020)	10	7.6	vertical GRF*	Skeleton	Real
UnderPressure (Mourot et al., 2022)	10	5.5	vertical GRF*	Skeleton	Real
GroundLink (Han et al., 2023)	7	1.5	GRF*	SMPL	Real
AddBiomechanics (Werling et al., 2024)	273	57.6	GRF* & Torques	Rajagopal et al. (2016)	Real
ImDy	435	152.3	GRF & Torques	SMPL	Simulated

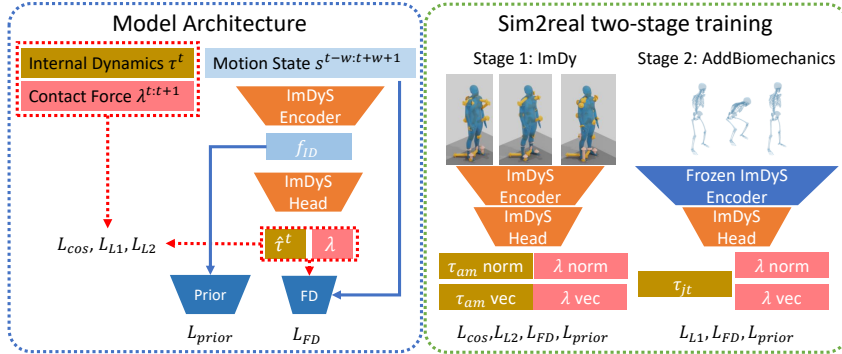


Figure 3: ImDyS overview. Taking a motion transition, ImDyS predicts the internal dynamics and ground reaction forces. Moreover, a prior discriminator is trained with the feature from ImDyS. A two-stage sim2real training curriculum is further designed.

gravity axis. Then, the sequences are imitated three times by the two primitive policies and the multiplicative policy with a simulation frequency of 60Hz, resulting in over 150 hours of human motion data with dynamics. States including q, \dot{q}, β are recorded in synchronous with the torque τ , all restored in the format of SMPL (Loper et al., 2015) if possible. Moreover, GRFs for the whole body are also recorded, resulting in ImDy, a large-scale human motion dynamics dataset.

Detailed statistics of ImDy are demonstrated in Tab. 1. There are three major advantages. First, a considerably larger data scale is $100\times$ compared to previous efforts with full-body dynamics data, covering a wide span of human motion, which could be hard to acquire in laboratory setups. Second, thanks to the advanced simulator (Makoviychuk et al., 2021), we could include ground reaction forces for the whole body instead of the two feet only like in previous efforts. Finally, we represent humans with SMPL (Loper et al., 2015), increasing availability.

4 LEARNING IMDYS

With the collected ImDy, we could address the human inverse dynamics in a full-supervised manner with a data-driven solver ImDyS. In Sec. 4.1, we first introduce the formulation of data-driven inverse dynamics. Then, the proposed data-driven solver is introduced in Sec. 4.2. The overall pipeline of ImDyS is illustrated in Fig. 3.

4.1 FORMULATION

Recall the abstraction of ID in Eq. 2, which we rewrite as

$$(\tau^t, \lambda^{t:t+1}) = \text{ImDyS}(s^{t-w:t+w+1}). \quad (4)$$

Given the kinematics states from timestamp $t - w$ to $t + w + 1$, ImDyS is required to estimate the internal dynamics τ^t for the transition from s^t to s^{t+1} and the ground reaction forces λ that the subject bears in timestamp t and $t + 1$.

Motion States s could be represented by either SMPL parameters, joint angles, joint coordinates, or marker coordinates. However, due to the topology divergence, the conversion among SMPL parameters, joint angles, and joint coordinates is non-trivial with limited performances. To guarantee that ImDyS could be seamlessly adopted to both ImDy and real-world biomechanics data, we adopt

marker coordinates as motion state representation for ImDyS. The state $s^t = (m^t, \dot{m}^t)$ is composed of marker coordinates m^t and finite-differentiated velocities \dot{m}^t at timestamp t , which are easy to obtain for both ImDy and AddBiomechanics (Werling et al., 2024). Two temporal windows before and after the transition with a length of w are included for contextual information. Notice that human physical properties like height and weight could also be implicitly represented by the markers. The states are canonicalized w.r.t. the heading direction of s^t .

Internal Dynamics τ . For ImDy, the imposed angular momentum τ_{am} is adopted for dynamics representation. Notice that in Sec. 3.2, the original sequences are in 30FPS, while the simulation runs at 60FPS. This means for each motion transition (s^t, s^{t+1}) , two torques were applied sequentially, each for $\frac{1}{60}s$. Predicting both torques is a plausible design choice. However, the second torque is based on the un-recorded mid-state between s^t, s^{t+1} . Predicting it involves the forward dynamics from s^t to the mid-state, with increased complexity. To this end, instead of predicting instantaneous torques, we switch to predicting the imposed angular momentum $\tau_{am} \in \mathbb{R}^{(J-1) \times 3}$, the time-accumulation effect of torque, for each motion transition. Thus, the modeling could stay consistent with proper complexity, only needing to sum the two torques up for s^t, s^{t+1} and then multiply it with the delta time. For AddBiomechanics, joint torque τ_{jt} is adopted for dynamics representation.

Ground Reaction Forces λ . Different from previous efforts (Mourot et al., 2022; Han et al., 2023; Werling et al., 2024) with foot GRFs only, we predict full-body GRF $\lambda \in \mathbb{R}^{J \times 3}$ as in Fig. 1.

4.2 DATA-DRIVEN IMDYS

Model architecture. With the enormous data scale of ImDy, we would like to keep ImDyS simple. An encoder-head structure is adopted. $s^{t-w:t+w+1} \in \mathbb{R}^{M \times (2w+2) \times 6}$ is first flattened as $\tilde{s}^{t-w:t+w+1} \in \mathbb{R}^{M \times (12w+12)}$ with window size w and M markers. Then, a transformer encoder converts \tilde{s} into ID feature $f_{ID} \in \mathbb{R}^d$, where d is the feature dimension. For prediction, we decompose τ_{am} and λ into magnitudes $|\tau_{am}^t|, |\lambda^{t:t+1}|$ and direction vectors $\vec{\tau}_{am}^t, \vec{\lambda}^{t:t+1}$ and predict each of them with a linear head. τ_{jt}^t is predicted with another linear head. The final predictions are $\hat{\tau}_{am}^t = |\tau_{am}^t| \vec{\tau}_{am}^t, \hat{\lambda}^{t:t+1} = |\lambda^{t:t+1}| \vec{\lambda}^{t:t+1}$ and τ_{jt}^t .

Loss terms. L1 loss, cosine loss, and L2 loss are adopted to optimized the predicted magnitudes $|\tau_{am}^t|, |\lambda^{t:t+1}|$, direction vectors $\vec{\tau}_{am}^t, \vec{\lambda}^{t:t+1}$, and joint torques τ_{jt} as L_{mag}, L_{cos}, L_{L2} respectively. Besides, a forward dynamics (FD) loss L_{fd} is proposed with an auxiliary FD model to inform the learning with the ID-FD cycle. The FD model takes $s^{t-w:t}, \tau^t = (\tau_{am}^t, \tau_{jt}^t), \lambda^t$ as input, predicts the next-frame joint angles. The FD loss is thus computed with cycle consistency as

$$L_{FD} = |s^{t+1} - FD(s^{t-w:t}, \hat{\tau}^t, \hat{\lambda}^t)|. \quad (5)$$

Finally, we devise a loss term similar to Peng et al. (2021), which encourages the ImDy feature f_{ID} to model physically plausible motion transitions. A linear discriminator takes f_{ID} and outputs a logit indicating whether the motion transition is plausible. To train the discriminator, besides the positive samples from ImDy and AMASS (Mahmood et al., 2019), we propose two negative sample generation strategies. First, $s^{t-w:t+w+1}$ is randomly permuted along the temporal axis. Second, random Gaussian noises are added on $s^{t-w:t+w+1}$. Binary cross-entropy loss is adopted as L_{cls} .

Sim2Real training curriculum is devised in a simple two-stage manner. In the first stage, ImDyS is trained on ImDy, with the overall loss as $\mathcal{L}_{s1} = \alpha_1 L_{mag} + \alpha_2 L_{cos} + \alpha_3 L_{FD} + \alpha_4 L_{cls}$. In the second stage, we freeze the encoder and train the linear head for joint torques τ_{jt} . The loss is calculated as $\mathcal{L}_{s2} = \alpha_3 L_{FD} + \alpha_4 L_{cls} + \alpha_5 L_{L2}$. Results show that ImDy pre-trained encoder converges fast on AddBiomechanics, indicating that it holds useful knowledge on real-world human dynamics.

5 EXPERIMENTS

5.1 IMPLEMENTATION DETAILS

PHC (Luo et al., 2023) adopted the position-control mode implemented by IsaacGym (Makoviychuk et al., 2021), where the imposed torque is calculated differently from the naive PD controller and inaccessible. Therefore, we re-trained the PHC on AMASS (Mahmood et al., 2019) with the effort-control mode, and a naive PD controller was adopted. Training the PHC took approximately 10

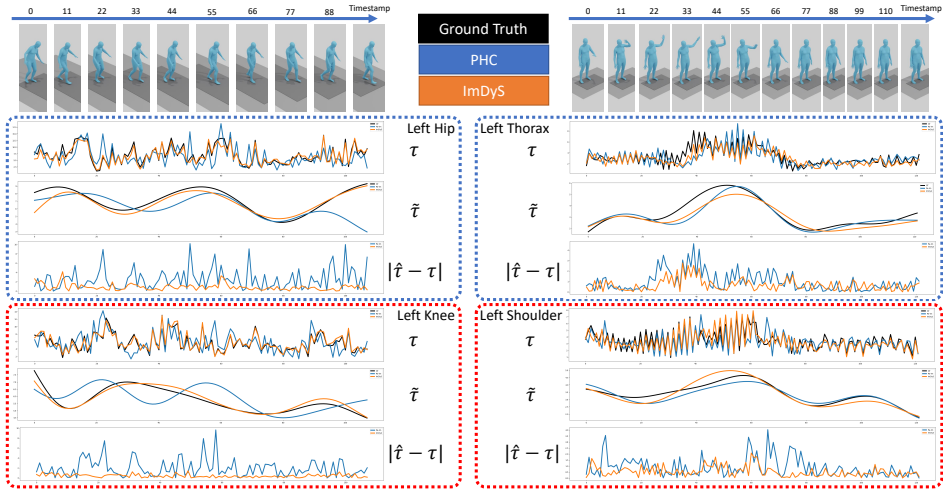


Figure 4: Qualitative results on ImDy. $\tilde{\cdot}$ indicates a low-pass filter at 14Hz is applied. A typical gait sample and an arm-waving sample are visualized.

days, with a success rate on AMASS of 91.3%. The window size w is set as 2 to keep a short-term motion modeling, which is proven helpful in Sec. 5.3. The encoder of ImDyS is a three-layer transformer with a dimension of 64, ReLU activation, and LayerNorm. The loss weights are set as $\alpha_1 = \alpha_3 = 0.01$, $\alpha_2 = \alpha_4 = \alpha_5 = 1$ to maintain all terms at similar numerical scales for training stability. ImDyS, the prior discriminator, and the FD model are all trained using the AdamW optimizer with a batch size of 2,400 for 140 epochs on ImDy for the first stage. For the second stage, ImDyS is further tuned on AddBiomechanics for only 10 epochs with the same hyperparameters. When generating negative samples for the prior discriminator, the two strategies are randomly adopted with a positive-negative ratio of 1:1. We split ImDy into a training set of 27,501 sequences and a test set of 3,055 sequences. All the data collection processes and experiments are conducted on a single NVIDIA RTX3090 GPU.

5.2 EVALUATION ON IMDY

Metric. We calculate the mPJE (mean Per Joint Error) for τ and λ as

$$mPJE_{\tau} = \frac{1}{J} \sum_{j=1}^J |\tau_j - \hat{\tau}_j|_2, \quad mPJE_{\lambda} = \frac{1}{J} \sum_{j=1}^J |\lambda_j - \hat{\lambda}_j|_2, \quad (6)$$

where J is the number of joints. The result is further normalized by body weight to align different subjects, with units of $N \cdot m \cdot s/kg$ and N/kg . Specifically, the mPJE for the GRF on both feet $mPJE_{\lambda_{l,f}}$, $mPJE_{\lambda_{r,f}}$ is also reported.

Baseline. Few efforts except IL algorithms are feasible as baselines. To this end, we introduce PHC as a baseline, where the sequences in ImDy are re-imitated by the re-trained PHC. The imposed angular momentums and the GRF obtained via the re-imitation process are adopted as the baseline predictions. With this baseline, we demonstrate the amplification effect from the kinematics error to the dynamics error, thus validating the performance of directly adopting IL for ID.

Results. Quantitative results are shown in Tab. 2. PHC produces an mPJPE of 56.13 mm, which is admirable for kinematics but results in high dynamics errors. ImDyS demonstrates considerably better performance. We further visualize two qualitative samples in Fig. 4. Since the raw data could be jittering, we also filter the predictions with a low-pass filter at 14Hz, denoted as $\tilde{\tau}$, $\tilde{\lambda}$, which helps reveal the general trend of the predictions. For the gait sample at the left, the imposed angular momentum τ at the left hip and the left knee are plotted, along with the GRF λ at the left toe. We also plot the error between the predicted values and GT values. ImDyS manages to faithfully reconstruct τ for the left knee and hip with minor errors. Meanwhile, PHC typically produces higher errors due to phase mismatch. As shown, it tends to lag behind the input motion. For GRF, ImDyS also produces reasonable predictions. Besides a typical gait analysis sample, we also demonstrate the performance of ImDyS with an arm-waving motion. The τ at directly related body segments

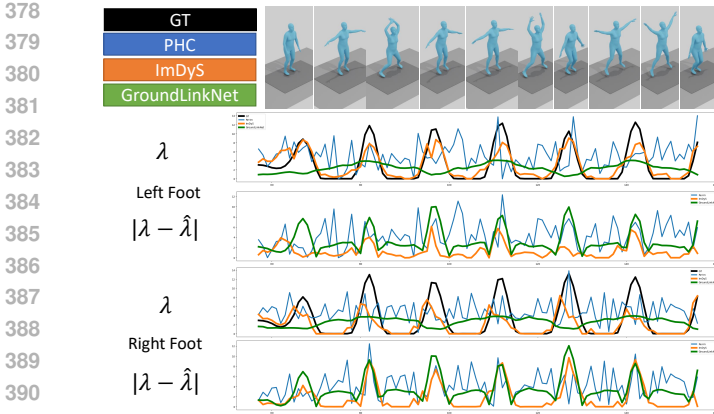


Figure 5: Qualitative results on GroundLink including PHC, GroundLinkNet, and ImDyS. The GRF λ for both feet are shown. Surprisingly, ImDyS provides better consistency with the ground truth.

including the left thorax and shoulder is visualized. ImDyS reproduces the dynamic status with better alignment to GT compared to PHC. Generally, ImDyS produces reasonable ID predictions. A potential issue is the jittering prediction, which is a consequence of the jittering observations in ImDy. However, we show that ImDyS could handle real-world smooth observations well even when trained only on jittering ImDyS. More demonstrations are available in the supplementary video.

5.3 EVALUATION ON GROUNDLINK

Metrics. GroundLink (Han et al., 2023) provides 1.5-hour motion from 7 subjects with GRF. We adopt subject 7 for evaluation. $mPJE_\lambda$ at both feet normalized by body weight is reported.

Baselines. PHC is evaluated similarly to Sec. 5.2. We also report the performance of GroundLinkNet (Han et al., 2023). PHC and ImDyS are not exposed to GroundLink during training, resulting in a **zero-shot** evaluation for ImDyS and the PHC baseline. Also, GroundLinkNet operates on 250FPS motion, while ImDyS and the PHC re-imitation baseline only operate on 30FPS motion. Finally, GroundLinkNet predicts GRF for both feet, while ImDyS and PHC could decouple feet into ankles and toes, and predict GRF separately for each part. We add up the ankle GRF and the toe GRF as the foot GRF. All predictions are re-sampled to 30FPS.

Results. Quantitative results are illustrated in Tab. 3. Surprisingly, both ImDyS and PHC manage to outperform the specifically trained GroundLinkNet. We attribute this to the enormous scale of AMASS and ImDy, which is much larger than GroundLink. Moreover, even though ImDyS is trained on simulated ImDy only, it generalizes to real-world data with competitive performance. We visualize the results in Fig. 5, 7. The PHC re-imitation baseline produces jittering predictions similar to Fig. 4. GroundLinkNet, though specifically trained on GroundLink, fails to capture the rapid GRF changes in this jumping jack motion, resulting in a relatively flat output. In contrast, ImDyS surprisingly presents good consistency with GT, and even faithfully reproduces the intense peak GRFs for the left foot for the jumping jack. Besides, the prediction is not as jittering as in Fig. 4, indicating ImDyS could handle real-world smooth data well.

5.4 EVALUATION ON ADDBIOMECHANICS

Metrics. AddBiomechanics (Werling et al., 2024) is recently proposed with over 50 hours of human dynamics data from 273 subjects. We adopt the armless part of this dataset. We follow the train/test split in Addbiomechanics and report mPJE for the joint torque normalized by body weight.

Results. A baseline model trained only on AddBiomechanics for 150 epochs with the same architecture as ImDyS is reported to showcase the generalization from ImDy to real-world dynamics. All data are re-sampled to 30 FPS. Quantitative results are illustrated in Tab. 4. ImDyS outperforms the baseline with faster convergence, indicating the efficacy of ImDyS in pre-training and mitigating

Table 2: Quantitative results on ImDy. mPJE is normalized by the body weight.

Methods	PHC	ImDyS
$mPJE_\tau$ (Nm/s/kg) ↓	0.095	0.021
$mPJE_\lambda$ (N/kg) ↓	0.409	0.289
$mPJE_{\lambda_{lf}}$ (N/kg) ↓	3.034	1.843
$mPJE_{\lambda_{rf}}$ (N/kg) ↓	2.866	1.842

Table 3: Ground reaction force prediction results on GroundLink.

Methods	PHC	GroundLinkNet	ImDyS
$mPJE_{\lambda_{lf}}$ (N/kg) ↓	2.362	5.423	0.986
$mPJE_{\lambda_{rf}}$ (N/kg) ↓	2.636	2.891	1.149

Table 4: Quantitative results on Ad-Biomechanics.

Methods	Baseline (150 epochs)	ImDyS (10 epochs)
$mPJE_\tau$ (Nm/kg) ↓	0.1699	0.1626 _{↓4.30%}
$mPJE_\lambda$ (Nm/kg) ↓	1.0876	1.0633 _{↓2.18%}

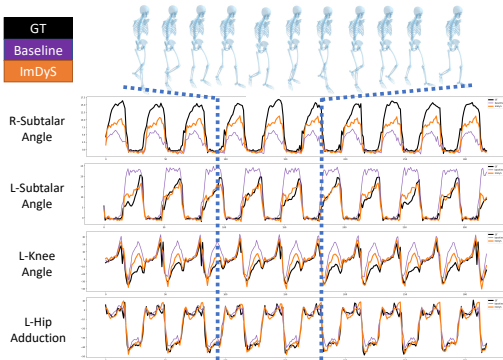


Figure 6: Joint torque predictions on AdBiomechanics.

the sim2real gap. Qualitative results are shown in Fig. 6, 8, where ImDyS shows better alignment with GT and more precise magnitude predictions. More analyses on the relationship between performance, data distribution and quality are in the appendix.

5.5 ABLATION STUDIES

Different Motion Representations are evaluated on ImDy in Tab. 5. Though SMPL and joint-based representations perform better, we adopt marker-based representation for its generality.

Different Loss Terms are evaluated in Tab. 5. L_{FD} is proven to contribute more than L_{cls} .

Different Window Sizes w are evaluated on AddBiomechanics in Tab. 6. ImDyS achieves the best balance between rich contexts and conciseness with $w = 2$.

6 DISCUSSION

Given the fully simulated nature of ImDy, a reasonable question is the sim2real problem. ImDy could be unnaturally jittering as in Fig. 4. Also, the physical properties of the simulated humanoid differ from those of real humans. Empirically, experiments show that ImDyS generalizes well to real-world data, partially mitigating this gap. The reason could be threefold. First, the jitters are unnatural but still physically plausible given that ImDy faithfully preserves consistent information for the simulated physics phenomena. Second, the small window size of ImDyS prevents it from relying on long-term contexts, where jitters are more salient. Finally, the enormous scale of ImDy is helpful for generalization. To further mitigate the sim2real gap with ImDy is a meaningful goal to pursue. Besides, ImDyS is designed as a first-step baseline to demonstrate the efficacy of ImDy. Introducing more sophisticated designs to regulate the behavior of ImDyS would be preferable. Moreover, ImDy only considers GRF, while other external forces are not involved. Also, interaction with other entities is absent. Exploration of these would be interesting for future works.

7 CONCLUSION

Leveraging the inherent resemblance between inverse dynamics and imitation learning, we proposed a novel human dynamics dataset ImDy, which contained over 150 hours of human motion paired with full-body driven torques and GRFs from well-developed simulator and imitation algorithms. Based on ImDy, a data-driven human inverse dynamics solver ImDyS is devised to reconstruct the driven angular momentum and contact forces from kinematic observations. ImDyS demonstrated impressive performance on both simulated and real-world data. As a first step toward scalable and easily accessible human inverse dynamics, we hope ImDy can shed new light on the data-driven physical analysis of human motion.

Table 5: Ablation study on ImDy.

Methods	$mPJE_\tau$ (Nm/s/kg)↓	$mPJE_\lambda$ (N/kg)↓	$mPJE_{\lambda_{lf}}$ (N/kg)↓	$mPJE_{\lambda_{rf}}$ (N/kg)↓
ImDyS	0.021	0.289	1.843	1.842
ImDyS-SMPL	0.011	0.272	1.746	1.764
ImDyS-Joint	0.014	0.273	1.755	1.775
w/o L_{FD}	0.023	0.302	1.962	1.976
w/o L_{cls}	0.022	0.294	1.884	1.891

Table 6: Ablation study on AddBiomechanics.

Methods	ImDyS $w=2$	ImDyS $w=1$	ImDyS $w=3$
$mPJE_\tau$ (Nm/kg) ↓	0.1626	0.1690	0.1720
$mPJE_\lambda$ (Nm/kg) ↓	1.0633	1.0990	1.1030

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

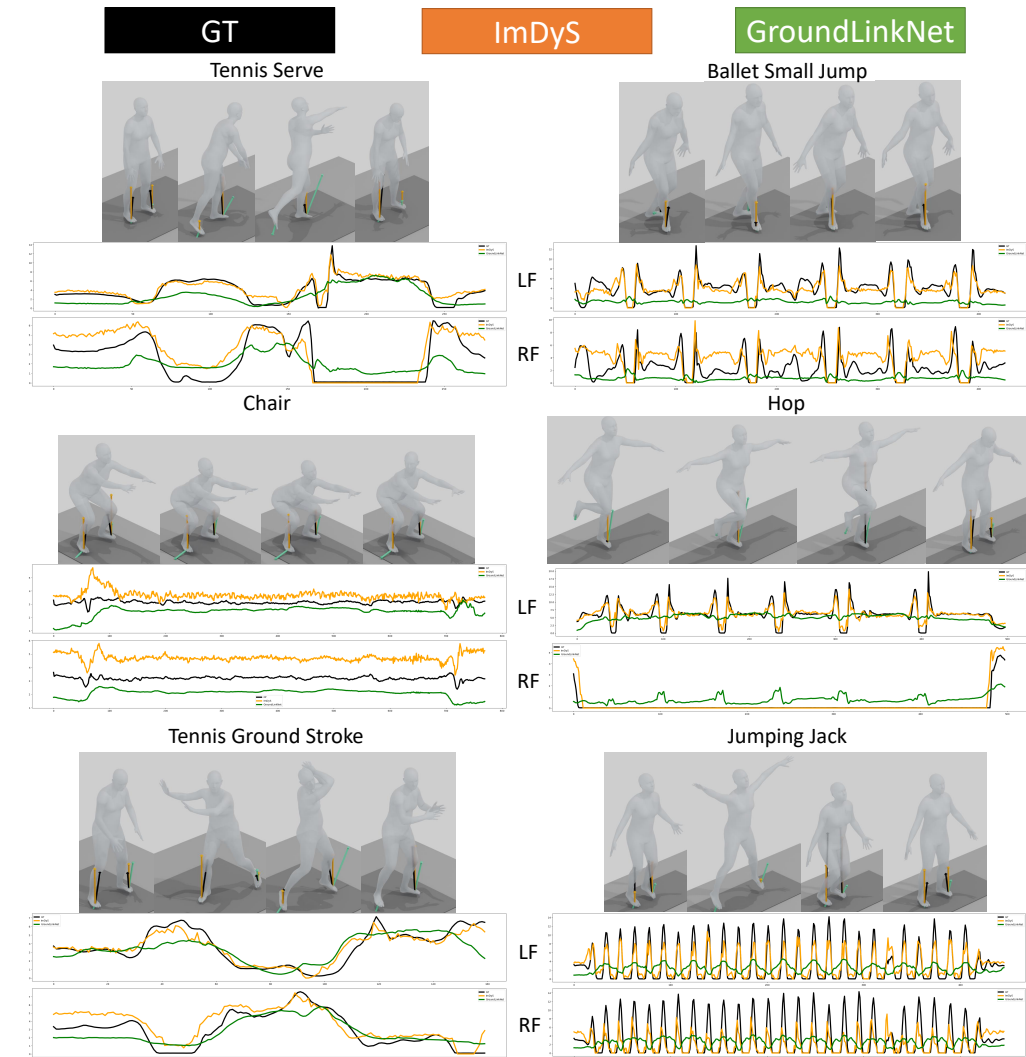
516

517

518

519

520



521 Figure 7: Extensive visualization on GroundLink. For various motions, ImDyS shows superior
 522 alignment with GT compared to specifically trained GroundLinkNet, showcasing the efficacy of
 523 ImDy. Especially, the intense peaks are also reproduced by ImDyS.

524

525

526

527

528

529

530

531

532

533

534

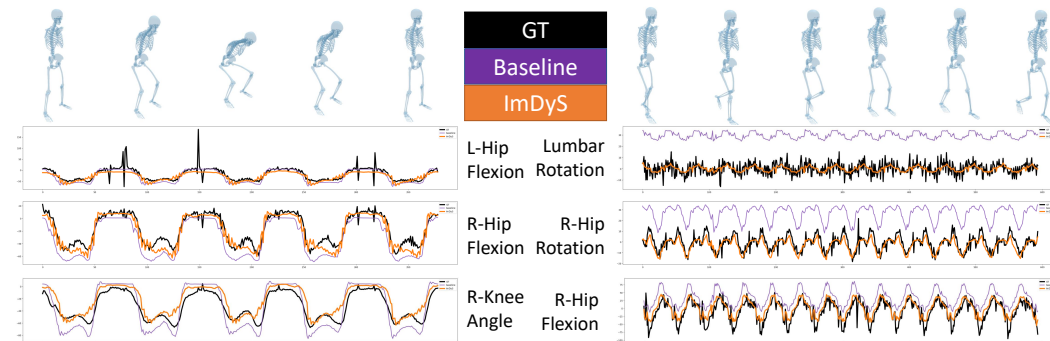
535

536

537

538

539



537 Figure 8: Extensive visualization on AddBiomechanics. ImDyS demonstrates superior performance
 538 to the baseline, indicating ImDy’s generalization ability.

REFERENCES

- 540
541
542 Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. Drecon: data-driven
543 responsive control of physics-based characters. *ACM Transactions On Graphics (TOG)*, 38(6):
544 1–11, 2019.
- 545
546 Dumitru I Caruntu and Ricardo Moreno. Human knee inverse dynamics model of vertical jump
547 exercise. *Journal of Computational and Nonlinear Dynamics*, 14(10):101005, 2019.
- 548
549 Michael Damsgaard, John Rasmussen, Søren Tørholm Christensen, Egidijus Surma, and Mark
550 De Zee. Analysis of musculoskeletal systems in the anybody modeling system. *Simulation Mod-
551 elling Practice and Theory*, 14(8):1100–1111, 2006.
- 552
553 Scott L Delp, Frank C Anderson, Allison S Arnold, Peter Loan, Ayman Habib, Chand T John,
554 Eran Guendelman, and Darryl G Thelen. Opensim: open-source software to create and analyze
555 dynamic simulations of movement. *IEEE transactions on biomedical engineering*, 54(11):1940–
556 1950, 2007.
- 557
558 Luis FC Figueredo, Rafael Castro Aguiar, Lipeng Chen, Samit Chakrabarty, Mehmet R Dogar, and
559 Anthony G Cohn. Human comfortability: Integrating ergonomics and muscular-informed metrics
560 for manipulability analysis during human-robot collaboration. *IEEE Robotics and Automation
561 Letters*, 6(2):351–358, 2020.
- 562
563 Claudiane A Fukuchi, Reginaldo K Fukuchi, and Marcos Duarte. A public dataset of overground
564 and treadmill walking kinematics and kinetics in healthy individuals. *PeerJ*, 6:e4640, 2018.
- 565
566 E. Gartner, M. Andriluka, E. Coumans, and C. Sminchisescu. Differentiable dynamics for artic-
567 ulated 3d human motion reconstruction. In *2022 IEEE/CVF Conference on Computer Vi-
568 sion and Pattern Recognition (CVPR)*, pp. 13180–13190, Los Alamitos, CA, USA, jun 2022.
569 IEEE Computer Society. doi: 10.1109/CVPR52688.2022.01284. URL [https://doi.
570 ieecomputersociety.org/10.1109/CVPR52688.2022.01284](https://doi.ieeeecomputersociety.org/10.1109/CVPR52688.2022.01284).
- 571
572 Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating
573 diverse and natural 3d human motions from text. In *2022 IEEE/CVF Conference on Computer
574 Vision and Pattern Recognition (CVPR)*, pp. 5142–5151, Los Alamitos, CA, USA, 2022. IEEE
575 Computer Society. doi: 10.1109/CVPR52688.2022.00509.
- 576
577 Erik Gärtner, Mykhaylo Andriluka, Hongyi Xu, and Cristian Sminchisescu. Trajectory optimization
578 for physics-based reconstruction of 3d human pose from monocular video. In *2022 IEEE/CVF
579 Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13096–13105, Los Alami-
580 tos, CA, USA, 2022. IEEE Computer Society. doi: 10.1109/CVPR52688.2022.01276.
- 581
582 Xingjian Han, Ben Senderling, Stanley To, Deepak Kumar, Emily Whiting, and Jun Saito.
583 Groundlink: A dataset unifying human body movement and ground reaction dynamics. In
584 *SIGGRAPH Asia 2023 Conference Papers*, SA '23, New York, NY, USA, 2023. Association
585 for Computing Machinery. ISBN 9798400703157. doi: 10.1145/3610548.3618247. URL
586 <https://doi.org/10.1145/3610548.3618247>.
- 587
588 B. Huang, L. Pan, Y. Yang, J. Ju, and Y. Wang. Neural mocon: Neural motion control for physically
589 plausible human motion capture. In *2022 IEEE/CVF Conference on Computer Vision and Pattern
590 Recognition (CVPR)*, pp. 6407–6416, Los Alamitos, CA, USA, jun 2022. IEEE Computer Society.
591 doi: 10.1109/CVPR52688.2022.00631. URL [https://doi.ieeeecomputersociety.
592 org/10.1109/CVPR52688.2022.00631](https://doi.ieeeecomputersociety.org/10.1109/CVPR52688.2022.00631).
- 593
594 Leif Johnson and Dana H. Ballard. Efficient codes for inverse dynamics during walking. In Carla E.
595 Brodley and Peter Stone (eds.), *Proceedings of the Twenty-Eighth AAAI Conference on Artificial
596 Intelligence*, pp. 343–349, Québec City, Québec, Canada, 2014. AAAI Press. doi: 10.1609/AAAI.
597 V28I1.8747. URL <https://doi.org/10.1609/aaai.v28i1.8747>.
- 598
599 Marilyn Keller, Keenon Werling, Soyong Shin, Scott Delp, Sergi Pujades, C Karen Liu, and
600 Michael J Black. From skin to skeleton: Towards biomechanically accurate 3d digital humans.
601 *ACM Transactions on Graphics (TOG)*, 42(6):1–12, 2023.

- 594 Franziska Krebs, Andre Meixner, Isabel Patzer, and Tamim Asfour. The kit bimanual manipulation
595 dataset. In *IEEE/RAS International Conference on Humanoid Robots (Humanoids)*, pp. 499–506,
596 2021.
- 597 Claudia Latella, Naveen Kuppaswamy, Francesco Romano, Silvio Traversaro, and Francesco Nori.
598 Whole-body human inverse dynamics with distributed micro-accelerometers, gyros and force
599 sensing. *Sensors*, 16(5), 2016. ISSN 1424-8220. doi: 10.3390/s16050727. URL <https://www.mdpi.com/1424-8220/16/5/727>.
- 600
601
- 602 Claudia Latella, Silvio Traversaro, Diego Ferigo, Yeshasvi Tirupachuri, Lorenzo Rapetti, Fran-
603 cisco Javier Andrade Chavez, Francesco Nori, and Daniele Pucci. Simultaneous floating-base
604 estimation of human kinematics and joint torques. *Sensors*, 19(12), 2019. ISSN 1424-8220. doi:
605 10.3390/s19122794. URL <https://www.mdpi.com/1424-8220/19/12/2794>.
- 606 Jiefeng Li, Siyuan Bian, Chao Xu, Gang Liu, Gang Yu, and Cewu Lu. D & d: Learning human
607 dynamics from dynamic camera. In *Computer Vision – ECCV 2022: 17th European Conference,
608 Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part V*, pp. 479–496, Berlin, Heidelberg,
609 2022. Springer-Verlag. ISBN 978-3-031-20064-9. doi: 10.1007/978-3-031-20065-6_28. URL
610 https://doi.org/10.1007/978-3-031-20065-6_28.
- 611 Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl:
612 a skinned multi-person linear model. In *ACM Transactions on Graphics*, volume 34, New York,
613 NY, USA, oct 2015. Association for Computing Machinery. doi: 10.1145/2816795.2818013.
614 URL <https://doi.org/10.1145/2816795.2818013>.
- 615 Z. Luo, J. Cao, A. Winkler, K. Kitani, and W. Xu. Perpetual humanoid control for real-
616 time simulated avatars. In *2023 IEEE/CVF International Conference on Computer Vision
617 (ICCV)*, pp. 10861–10870, Los Alamitos, CA, USA, oct 2023. IEEE Computer Society. doi:
618 10.1109/ICCV51070.2023.01000. URL [https://doi.ieeecomputersociety.org/
619 10.1109/ICCV51070.2023.01000](https://doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.01000).
- 620 Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy
621 for egocentric pose estimation. *Advances in Neural Information Processing Systems*, 34:25019–
622 25032, 2021.
- 623 Xiaolei Lv, Jinxiang Chai, and Shihong Xia. Data-driven inverse dynamics for human motion. *ACM
624 Transactions on Graphics (TOG)*, 35(6):1–12, 2016.
- 625 N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. Black. Amass: Archive of motion
626 capture as surface shapes. In *2019 IEEE/CVF International Conference on Computer Vision
627 (ICCV)*, pp. 5441–5450, Los Alamitos, CA, USA, nov 2019. IEEE Computer Society. doi: 10.
628 1109/ICCV.2019.00554. URL [https://doi.ieeecomputersociety.org/10.1109/
629 ICCV.2019.00554](https://doi.ieeecomputersociety.org/10.1109/ICCV.2019.00554).
- 630 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
631 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State.
632 Isaac gym: High performance GPU based physics simulation for robot learning. In
633 Joaquin Vanschoren and Sai-Kit Yeung (eds.), *Proceedings of the Neural Information Pro-
634 cessing Systems Track on Datasets and Benchmarks, Virtual*, 2021. Curran Associates
635 Inc. URL [https://datasets-benchmarks-proceedings.neurips.cc/paper/
636 2021/hash/28dd2c7955ce926456240b2ff0100bde-Abstract-round2.html](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/28dd2c7955ce926456240b2ff0100bde-Abstract-round2.html).
- 637 Christian Mandery, Ömer Terlemez, Martin Do, Nikolaus Vahrenkamp, and Tamim Asfour. Uni-
638 fying representations and large-scale whole-body motion databases for studying human motion.
639 *IEEE Transactions on Robotics*, 32(4):796–809, 2016.
- 640 Mykhailo Manukian, Serhii Bahdasariants, and Sergiy Yakovenko. Artificial physics engine for
641 real-time inverse dynamics of arm and hand movement. *Plos one*, 18(12):e0295750, 2023.
- 642 Lucas Mourot, Ludovic Hoyet, François Le Clerc, and Pierre Hellier. Underpressure: Deep learning
643 for foot contact detection, ground reaction force estimation and footskate cleanup. *Computer
644 Graphics Forum*, 41(8):195–206, 2022. doi: <https://doi.org/10.1111/cgf.14635>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/cgf.14635>.

- 648 Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial
649 motion priors for stylized physics-based character control. *ACM Transactions on Graphics (ToG)*,
650 40(4):1–20, 2021.
- 651 Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale
652 reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions On*
653 *Graphics (TOG)*, 41(4):1–17, 2022.
- 654 A. R. Punnakal, A. Chandrasekaran, N. Athanasiou, A. Quiros-Ramirez, and M. J. Black. Babel:
655 Bodies, action and behavior with english labels. In *2021 IEEE/CVF Conference on*
656 *Computer Vision and Pattern Recognition (CVPR)*, pp. 722–731, Los Alamitos, CA, USA,
657 jun 2021. IEEE Computer Society. doi: 10.1109/CVPR46437.2021.00078. URL <https://doi.ieeecomputersociety.org/10.1109/CVPR46437.2021.00078>.
- 658 Apoorva Rajagopal, Christopher L Dembia, Matthew S DeMers, Denny D Delp, Jennifer L Hicks,
659 and Scott L Delp. Full-body musculoskeletal model for muscle-driven simulation of human gait.
660 *IEEE transactions on biomedical engineering*, 63(10):2068–2079, 2016.
- 661 Davis Rempe, Leonidas J. Guibas, Aaron Hertzmann, Bryan Russell, Ruben Villegas, and Jimei
662 Yang. Contact and human dynamics from monocular video. In *Computer Vision – ECCV*
663 *2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V*, pp.
664 71–87, Berlin, Heidelberg, 2020. Springer-Verlag. ISBN 978-3-030-58557-0. doi: 10.1007/
665 978-3-030-58558-7_5. URL https://doi.org/10.1007/978-3-030-58558-7_5.
- 666 Céline Schreiber and Florent Moissenet. A multimodal dataset of human gait at different walking
667 speeds established on injury-free adult participants. *Scientific data*, 6(1):111, 2019.
- 668 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
669 optimization algorithms. preprint on webpage at arXiv:1707.06347, 2017.
- 670 Jesse Scott, Bharadwaj Ravichandran, Christopher Funk, Robert T. Collins, and Yanxi Liu. From
671 image to stability: Learning dynamics from human pose. In *Computer Vision – ECCV 2020:*
672 *16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII*, pp.
673 536–554, Berlin, Heidelberg, 2020. Springer-Verlag. ISBN 978-3-030-58591-4. doi: 10.1007/
674 978-3-030-58592-1_32. URL https://doi.org/10.1007/978-3-030-58592-1_32.
- 675 A. Shahroudy, J. Liu, T. Ng, and G. Wang. Ntu rgb+d: A large scale dataset for 3d human
676 activity analysis. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*
677 *(CVPR)*, pp. 1010–1019, Los Alamitos, CA, USA, jun 2016. IEEE Computer Society.
678 doi: 10.1109/CVPR.2016.115. URL <https://doi.ieeecomputersociety.org/10.1109/CVPR.2016.115>.
- 679 Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick Pérez, and Christian Theobalt. Neural
680 monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics*
681 *(ToG)*, 40(4):1–15, 2021.
- 682 Tatsuya Teramae, Tomoyuki Noda, and Jun Morimoto. Emg-based model predictive control for
683 physical human–robot interaction: Application for assist-as-needed control. *IEEE Robotics and*
684 *Automation Letters*, 3(1):210–217, 2017.
- 685 Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit Haim Bermano.
686 Human motion diffusion model. In *The Eleventh International Conference on Learning Representations*
687 *ICLR 2023, Kigali, Rwanda, 2023*. OpenReview.net. URL <https://openreview.net/pdf?id=SJ1kSyO2jwu>.
- 688 Thomas K Uchida and Scott L Delp. *Biomechanics of movement: the science of sports, robotics,*
689 *and rehabilitation*. Mit Press, 2021.
- 690 Thomas K Uchida and Ajay Seth. Conclusion or illusion: Quantifying uncertainty in inverse anal-
691 yses from marker-based motion capture due to errors in marker registration and model scaling.
692 *Frontiers in Bioengineering and Biotechnology*, 10:874725, 2022.

- 702 J. Wang, Y. Yuan, Z. Luo, K. Xie, D. Lin, U. Iqbal, S. Fidler, and S. Khamis. Learning
703 human dynamics in autonomous driving scenarios. In *2023 IEEE/CVF International Con-*
704 *ference on Computer Vision (ICCV)*, pp. 20739–20749, Los Alamitos, CA, USA, oct 2023.
705 IEEE Computer Society. doi: 10.1109/ICCV51070.2023.01901. URL [https://doi.](https://doi.ieeeecomputersociety.org/10.1109/ICCV51070.2023.01901)
706 [ieeeecomputersociety.org/10.1109/ICCV51070.2023.01901](https://doi.ieeeecomputersociety.org/10.1109/ICCV51070.2023.01901).
- 707
708 Keenon Werling, Dalton Omens, Jeongseok Lee, Ioannis Exarchos, and C. Karen Liu. Fast and
709 Feature-Complete Differentiable Physics Engine for Articulated Rigid Bodies with Contact Con-
710 straints. In *Robotics: Science and Systems XVII, Virtual Event, July 12-16, 2021*, Virtual, July
711 2021. RSS Foundation. doi: 10.15607/RSS.2021.XVII.034. URL [https://doi.org/10.](https://doi.org/10.15607/RSS.2021.XVII.034)
712 [15607/RSS.2021.XVII.034](https://doi.org/10.15607/RSS.2021.XVII.034).
- 713 Keenon Werling, Nicholas A Bianco, Michael Raitor, Jon Stingel, Jennifer L Hicks, Steven H
714 Collins, Scott L Delp, and C Karen Liu. Addbiomechanics: Automating model scaling, inverse
715 kinematics, and inverse dynamics from human motion data through sequential optimization. *Plos*
716 *one*, 18(11):e0295152, 2023.
- 717 Keenon Werling, Janelle Kaneda, Alan Tan, Rishi Agarwal, Six Skov, Tom Van Wouwe, Scott
718 Uhlich, Nicholas Bianco, Carmichael Ong, Antoine Falisse, et al. Addbiomechanics dataset:
719 Capturing the physics of human motion at scale. *arXiv preprint arXiv:2406.18537*, 2024.
- 720
721 Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Control strategies for physically simulated
722 characters performing two-player competitive sports. *ACM Transactions on Graphics (TOG)*, 40
723 (4):1–11, 2021.
- 724 Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Physics-based character controllers using
725 conditional vaes. *ACM Transactions on Graphics (TOG)*, 41(4):1–12, 2022.
- 726
727 Baoping Xiong, Nianyin Zeng, Han Li, Yuan Yang, Yurong Li, Meilan Huang, Wuxiang Shi, Min
728 Du, and Yudong Zhang. Intelligent prediction of human lower extremity joint moment: an arti-
729 ficial neural network approach. *Ieee Access*, 7:29973–29980, 2019.
- 730
731 Shaowei Yao, Yu Zhuang, Zhijun Li, and Rong Song. Adaptive admittance control for an ankle
732 exoskeleton using an emg-driven musculoskeletal model. *Frontiers in neurorobotics*, 12:16, 2018.
- 733
734 X. Yi, Y. Zhou, M. Habermann, S. Shimada, V. Golyanik, C. Theobalt, and F. Xu. Phys-
735 ical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial
736 sensors. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
737 *(CVPR)*, pp. 13157–13168, Los Alamitos, CA, USA, jun 2022. IEEE Computer Society. doi:
738 [10.1109/CVPR52688.2022.01282](https://doi.ieeeecomputersociety.org/10.1109/CVPR52688.2022.01282). URL [https://doi.ieeeecomputersociety.org/](https://doi.ieeeecomputersociety.org/10.1109/CVPR52688.2022.01282)
739 [10.1109/CVPR52688.2022.01282](https://doi.ieeeecomputersociety.org/10.1109/CVPR52688.2022.01282).
- 740
741 Y. Yuan and K. Kitani. Ego-pose estimation and forecasting as real-time pd control. In *2019*
742 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10081–10091, Los Alami-
743 totes, CA, USA, nov 2019. IEEE Computer Society. doi: 10.1109/ICCV.2019.01018. URL
744 <https://doi.ieeeecomputersociety.org/10.1109/ICCV.2019.01018>.
- 745
746 P. Zell and B. Rosenhahn. Learning-based inverse dynamics of human motion. In *2017 IEEE*
747 *International Conference on Computer Vision Workshop (ICCVW)*, pp. 842–850, Los Alamitos,
748 CA, USA, oct 2017. IEEE Computer Society. doi: 10.1109/ICCVW.2017.104. URL [https://doi.](https://doi.ieeeecomputersociety.org/10.1109/ICCVW.2017.104)
749 [ieeeecomputersociety.org/10.1109/ICCVW.2017.104](https://doi.ieeeecomputersociety.org/10.1109/ICCVW.2017.104).
- 750
751 P. Zell, B. Wandt, and B. Rosenhahn. Joint 3d human motion capture and physical analysis from
752 monocular videos. In *2017 IEEE Conference on Computer Vision and Pattern Recognition*
753 *Workshops (CVPRW)*, pp. 17–26, Los Alamitos, CA, USA, jul 2017. IEEE Computer Society.
754 doi: 10.1109/CVPRW.2017.9. URL [https://doi.ieeeecomputersociety.org/10.](https://doi.ieeeecomputersociety.org/10.1109/CVPRW.2017.9)
755 [1109/CVPRW.2017.9](https://doi.ieeeecomputersociety.org/10.1109/CVPRW.2017.9).
- 756
757 Petri Zell and Bodo Rosenhahn. A physics-based statistical model for human gait analysis. In
758 Juergen Gall, Peter Gehler, and Bastian Leibe (eds.), *Pattern Recognition*, pp. 169–180, Cham,
759 2015. Springer International Publishing. ISBN 978-3-319-24947-6.

756 PetriSSa Zell, Bodo Rosenhahn, and Bastian Wandt. Weakly-supervised learning of human dynam-
757 ics. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *Computer*
758 *Vision – ECCV 2020*, pp. 68–84, Cham, 2020. Springer International Publishing. ISBN 978-3-
759 030-58574-7.

760 Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations
761 in neural networks. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recog-*
762 *niton (CVPR)*, pp. 5738–5746, Los Alamitos, CA, USA, jun 2019. IEEE Computer Society.
763 doi: 10.1109/CVPR.2019.00589. URL [https://doi.ieeecomputersociety.org/
764 10.1109/CVPR.2019.00589](https://doi.ieeecomputersociety.org/10.1109/CVPR.2019.00589).
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

APPENDIX

A APPLICATIONS OF IMDyS

In this section, we demonstrate some downstream applications of ImDyS.

Human Work Analysis. With the predicted τ , we could calculate the work conducted at each joint. Visualizations are in Fig. 9, reasonably revealing the energy flow during human motion.

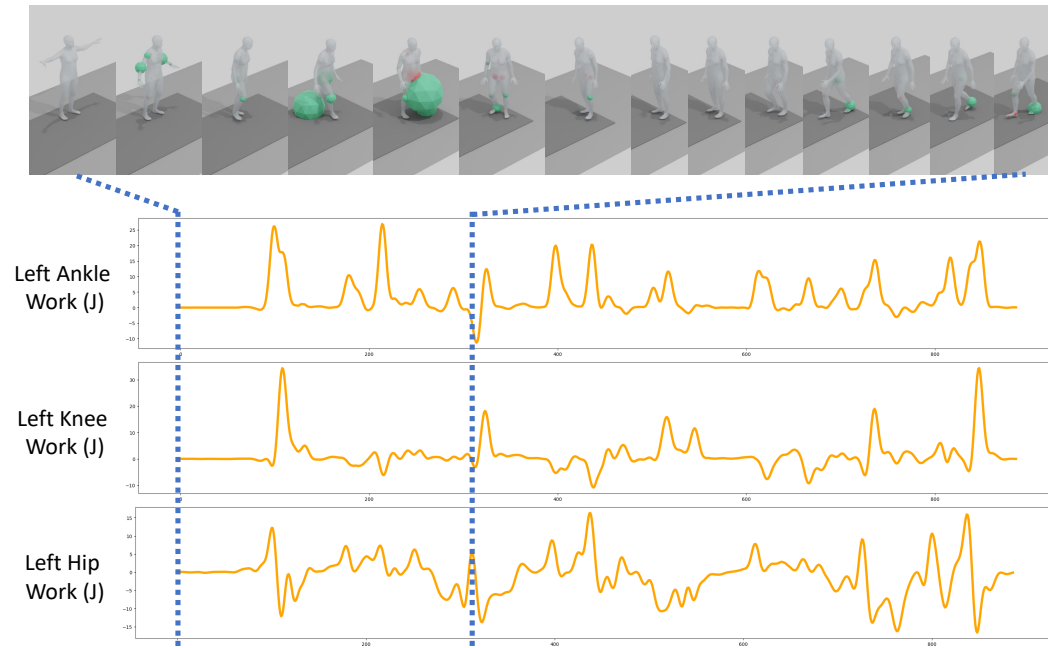


Figure 9: Human work visualization with ImDyS prediction. Green indicates positive work and red indicates negative work.

Motion Assessment. Another interesting application of ImDyS is based on the discriminator introduced in Sec. 4.2. Besides facilitating ImDyS learning, it could also assess whether a motion transition is physically plausible as in Fig. 10. Specifically, we adopt ImDyS to assess the motion generated from MDM Tevet et al. (2023). ImDyS reasonably tells when the motion starts to deviate from realism.

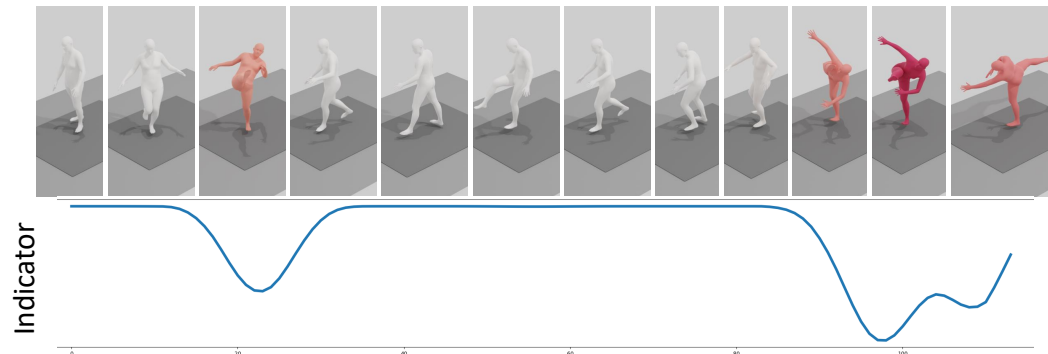
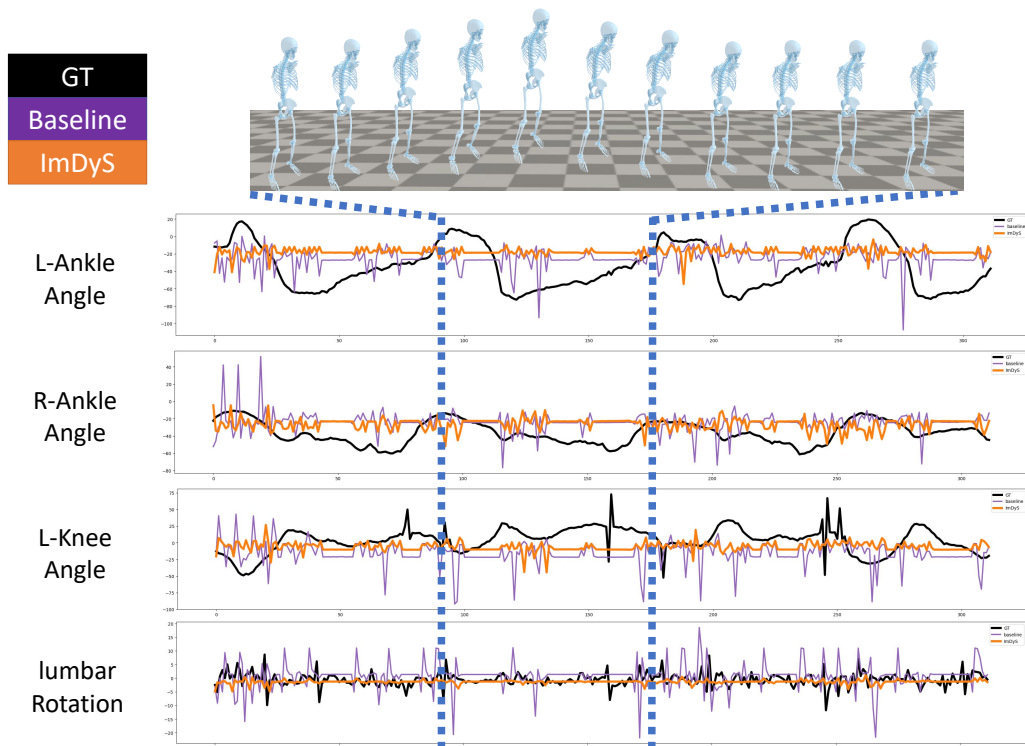


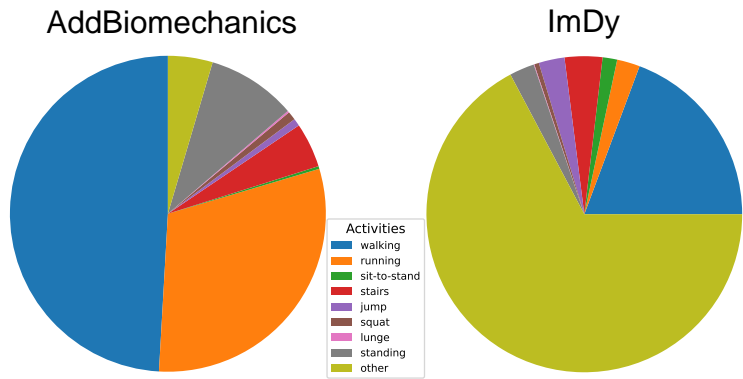
Figure 10: Motion assessment visualization. The motion artifacts are annotated with red with a low indicator value from ImDyS. As shown, ImDyS manages to identify implausible transitions in a kicking motion generated by MDM.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888



889 Figure 11: Visualization of a failed joint torque prediction case on AddBiomechanics. For the “jumping” motion, the baseline and ImDyS both perform sub-optimally. Neither of them correctly predicts the joint torques.

893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917



906 Figure 12: Data distribution of AddBiomechanics and ImDy. Among all activities, walking and running account for over 75% of AddBiomechanics. In comparison, according to the annotations from BABEL (Punnakkal et al., 2021), ImDyS is less imbalanced with better diversity.

910 B ADDBIOMECHANICS RESULTS ANALYSIS

912 We visualize a failure case on the AddBiomechanics dataset in Fig. 11. As shown, neither the baseline nor ImDyS manages to faithfully predict the joint torques for the jumping motion. In the following, we discuss the reasons for the failure.

916 **Data distribution.** Fig. 12 shows the data distribution of AddBiomechanics (Werling et al., 2024). As shown, over 75% of the data are either walking, running, or standing, which are extremely limited. Though ImDyS is empowered with the diverse ImDy as shown in Fig. 12, it still requires

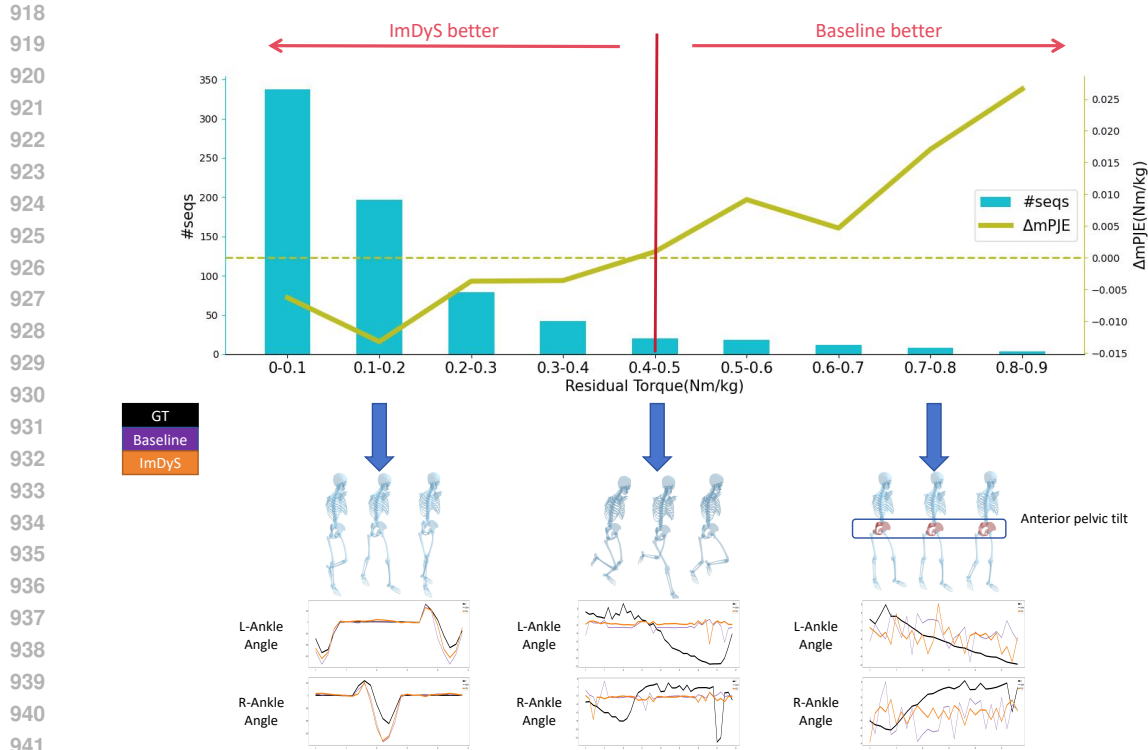


Figure 13: Relationship between data quality and model performance differences. Higher residual torque indicates lower data quality with lower reliability of the optimized GT torques. $\Delta mPJE$ is the difference between the mPJE of ImDyS and the baseline. #seqs is the number of sequences. With the residual torque increasing, the baseline provides lower mPJE than ImDyS, indicating the baseline overfits low-quality data. Instead, ImDyS, with the knowledge inherited from ImDy, shows less overfitting for these cases.

data to learn the mapping between simulated torques and real torques for non-gait data. These result in ImDyS' poor performance when processing non-gait data. Further mitigating the limited data issue for non-gait motions would be a meaningful goal to pursue.

Data quality. Besides the distribution, the quality is also limited in AddBiomechanics. As shown in Fig. 11, joint torques for some joints (like the lumbar) suffer from unstable optimization with jittering results. According to Werling et al. (2024), 21.2% of AddBiomechanics are classified with clinical-grade high quality (residual torque $< 0.1 \cdot \text{body weight} \cdot \text{height}$). There exists a 1.6829 Nm/kg average root residual torque of the optimized GTs in AddBiomechanics, which is considerably higher than the mPJE of ImDyS (0.1626 Nm/kg). We further analyze the relationship between the data quality and the model performances. We adopt residual torques as an indicator of the data quality and calculate $\Delta mPJE = mPJE_{ImDyS} - mPJE_{Baseline}$ of sequences with different residual torques. Notice that higher residual torque indicates lower data quality with lower reliability of the optimized GT torques. Results are shown in Fig. 13. As shown, some samples could suffer from bad kinematics fitting (like the unnatural anterior pelvic tilt in Fig. 13), resulting in less reliable GT optimized joint torques. An interesting phenomenon is that the lower the residual torques are, the better ImDyS performs, which means ImDyS performs better for high-quality samples. This indicates the baseline might overfit low-quality data with high residual torques. Instead, ImDyS, with the knowledge inherited from the large-scale diverse ImDy, manages to resist the negative influences from low-quality samples. We also show how the mPJE of ImDyS changes with data quality in Fig. 14. As shown, the performance of ImDyS degenerates synchronously with data quality.

Per-Joint Performance Analysis. It is also noticeable in Fig. 11 that the gap between GT and prediction differs for different joints. To this end, we further analyze the per-joint performance of ImDyS. The per-joint mPJE of ImDyS and the per-joint mPJE of ImDyS in each frame for samples

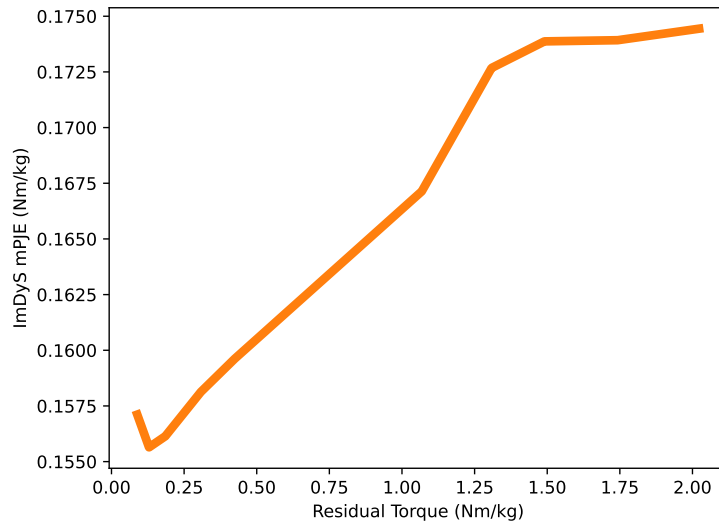


Figure 14: Relationship between data quality and ImDyS performance. Higher residual torque indicates lower data quality with lower reliability of the optimized GT torques. The performance of ImDyS degenerates synchronously with data quality.

Table 7: Per-Joint mPJE of ImDyS for samples with clinical-grade quality.

Joint Name	Right mPJE $_{\tau} \downarrow$		Left mPJE $_{\tau} \downarrow$	
	ImDyS	Baseline	ImDyS	Baseline
Hip Flexion	0.267	0.270	0.273	0.274
Hip Adduction	0.196	0.212	0.198	0.199
Hip Rotation	0.087	0.109	0.083	0.098
Knee	0.193	0.188	0.195	0.209
Ankle	0.197	0.195	0.202	0.199
Subtalar	0.069	0.071	0.079	0.084
MTP	0.0005	0.0006	0.0007	0.0008
Lumbar Extension	0.328	0.339	-	-
Lumbar Bending	0.255	0.255	-	-
Lumbar Rotation	0.113	0.113	-	-

with clinical-grade quality (residual torque < 0.1 body weight * height) is demonstrated in Tab. 7. ImDyS manages to improve the performance on most joints compared to the baseline without ImDy, especially for the hips. An interesting phenomenon is that ImDyS performs slightly better on the right half of the body.

C SIM2REAL ANALYSIS

We further analyze the Sim2Real effect of ImDy(S) via Fig. 15. An interesting question is the performance of ImDyS without any fine-tuning on AddBiomechanics. Though this could be inapplicable for most joints due to the human model definition discrepancy between Rajagopal’s model in AddBiomchanics and SMPL in ImDy, the knee joints in the two models could roughly correspond to each other. Therefore, we visualize the knee torque magnitudes of ImDyS and ImDyS w/o Sim2Real finetuning on AddBiomechanics in Fig. 15. Even without fine-tuning, ImDyS could reproduce the trends of knee torque magnitudes. However, artifacts could also be observed in two aspects. First, ImDyS w/o Sim2Real tends to produce much larger torques. Second, ImDyS w/o Sim2Real could be over-active compared to real humans and ImDyS like in the red circles. The reason could be two-fold. First, the simulation parameters used by ImDy, like mass and inertia, are different from real humans. Second, though the knee joints could roughly correspond, the knee in SMPL has more

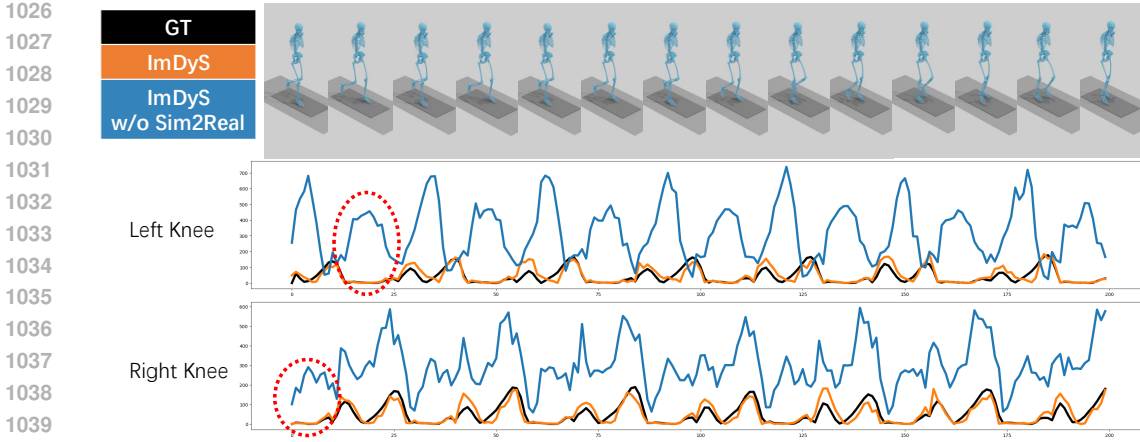


Figure 15: Knee torque magnitude visualization of ImDyS and ImDyS w/o Sim2Real fine-tuning on AddBiomechanics. ImDyS w/o Sim2Real produces larger magnitudes and over-active torques w/o Sim2Real fine-tuning as circled in red.

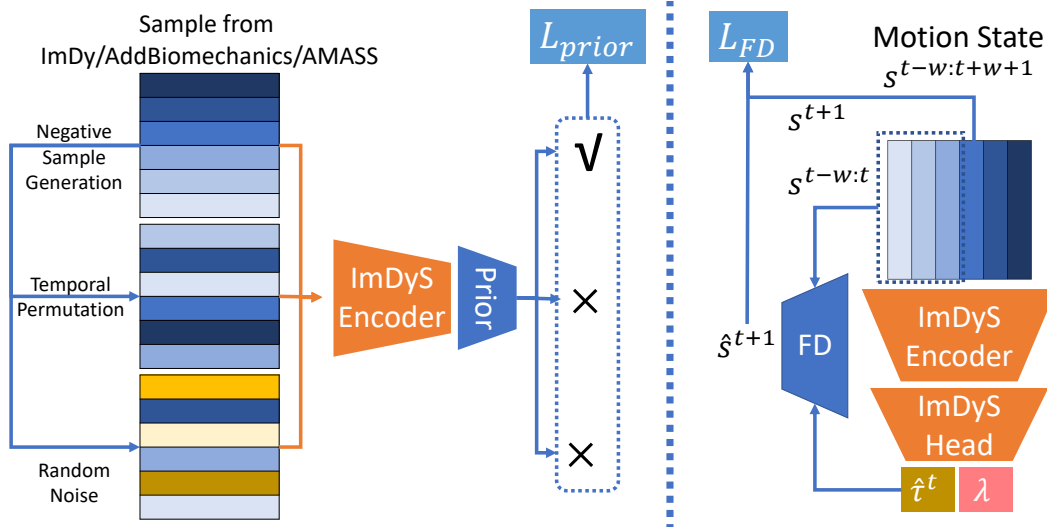


Figure 16: Details of L_{prior} and L_{FD} .

DoFs than Rajagopal’s model, which might require larger torques to produce similar motions. With the simple Sim2Real fine-tuning of ImDyS, the issues could be alleviated. Further exploration for better Sim2Real performance would be meaningful future work.

D DETAILS ON DATA FLOW

Details of the adopted L_{prior} and L_{FD} are illustrated in Fig. 16.

For L_{prior} , the input motion state is treated as the positive case, and we generate corresponding negative cases by either temporal permutation or adding random noises. The samples are fed to the encoder, and the prior discriminator predicts whether the sample is positive.

For L_{FD} , we first feed ImDyS with motion state $s^{t-w:t+w+1}$, obtaining τ, λ . Then, $\tau, \lambda, s^{t-w:t}$ are fed into the FD model, outputting \hat{s}^{t+1} . The FD loss is computed as $L_{FD} = |s^{t+1} - \hat{s}^{t+1}|$.

Table 8: Extended results on the FD model on Addbiomechanics.

Methods	Baseline	ImDyS	Nimble
RMSE	0.0302	0.0194	0.0186

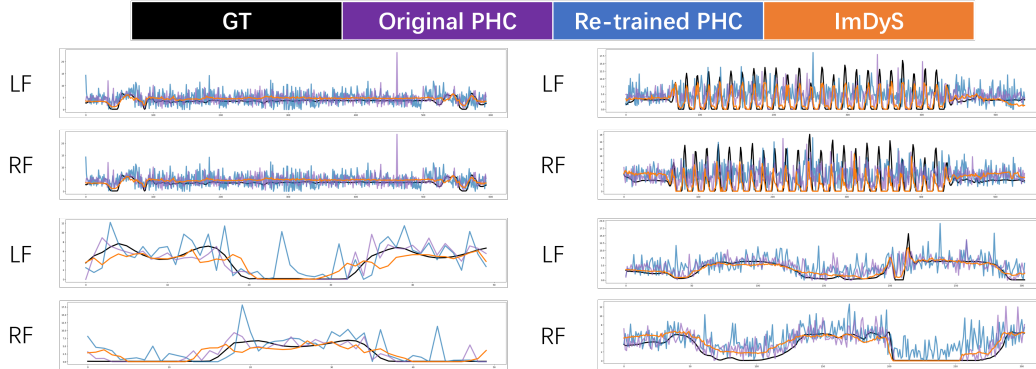


Figure 17: Original PHC on GroundLink.

E ANALYSIS ON FD MODEL

We report the marker RMSE of the FD model on the AddBiomechanics test set as Tab. 8. ImDyS noticeably outperforms the baseline trained on AddBiomechanics only, indicating the importance of ImDy pre-training. Moreover, ImDyS is competitive even compared to the differentiable simulator Nimble.

F COMPARISON WITH ORIGINAL PHC

Due to the inaccessible torques, we did not include the original PHC as a baseline for ImDy. However, it is noticeable that the original PHC can also conduct GRF prediction. To this end, we also evaluate the original PHC on GroundLink. It provides a left-foot mPJE of 1.559 and a right-foot mPJE of 3.518, which are comparable to re-trained PHC and worse than our proposed ImDyS. Some visualizations are included in Fig. 17. Even without the naive PD controller, the original PHC could suffer from jittering predictions, which could result from the non-perfect contact simulation. In contrast, ImDyS could produce smoother predictions with higher precision.