# Differentiable Dynamics for Articulated 3d Human Motion Reconstruction

Erik Gärtner[1,2]     Mykhaylo Andriluka[1]     Erwin Coumans[1]     Cristian Sminchisescu[1]

[1]**Google Research,** [2]**Lund University**

erik.gartner@math.lth.se

{mykhayloa,erwincoumans,sminchisescu}@google.com

## Abstract

*We introduce DiffPhy, a differentiable physics-based model for articulated 3d human motion reconstruction from video. Applications of physics-based reasoning in human motion analysis have so far been limited, both by the complexity of constructing adequate physical models of articulated human motion, and by the formidable challenges of performing stable and efficient inference with physics in the loop. We jointly address such modeling and inference challenges by proposing an approach that combines a physically plausible body representation with anatomical joint limits, a differentiable physics simulator, and optimization techniques that ensure good performance and robustness to suboptimal local optima. In contrast to several recent methods [40, 43, 56], our approach readily supports full-body contact including interactions with objects in the scene. Most importantly, our model connects end-to-end with images, thus supporting direct gradient-based physics optimization by means of image-based loss functions. We validate the model by demonstrating that it can accurately reconstruct physically plausible 3d human motion from monocular video, both on public benchmarks with available 3d ground-truth, and on videos from the internet.*

## 1. Introduction

We seek to contribute to the development of physics-based methodology as one of the building blocks in constructing accurate and robust 3d visual human sensing systems. Incorporating the laws of physics into the visual reasoning process is appealing as it promotes the plausibility of estimated motion and facilitates more efficient use of training examples [9]. We focus on articulated human motion as an epitome of a real-world prediction task that is both well studied and challenging. Existing state-of-the-art approaches demonstrate relatively high accuracy in terms of joint position estimation metrics [23, 24, 55, 63]. However,

predictions can sometimes be physically implausible, even for simple motions such as walking and running. For instance, estimates can include unreasonably abrupt transitions in world space, or artifacts such as foot skating or non-equilibrium states [40,43]. Many methods are typically trained on large motion capture datasets and encounter difficulties when tested on motions not well represented in those training sets. Arguably, imposing some form of physics-based generally valid prior on the articulated motion estimates should greatly improve the plausibility of results.

However, physics-based reasoning comes at the cost of substantial modeling and inference complexity. Typically, physics-based articulated estimation methods rely on rigid body dynamics (RBD) [10, 45], a formulation that introduces many auxiliary variables corresponding to forces acting at the body joints at each time step. Moreover, physical contact results in non-smooth effects where small changes to model parameters might result in substantially different motions. Therefore inferring physics variables given the inherent uncertainty in monocular video, and under contact discontinuities, becomes significantly difficult, algorithmically and computationally. Despite such challenges, a number of recent methods successfully apply physics-based constraints for articulated human motion estimation [2, 40, 43, 60]. One possibility to cope with modeling complexity, explored in recent work, is to simplify the physics and model contacts only between the body and the feet [40, 43, 56]. Others use auxiliary external forces applied at the body to compensate for modeling error [43, 60].

In this paper, we aim to broaden the methodology for physics-based articulated human motion estimation. Specifically, we demonstrate that we can successfully leverage recent progress in differentiable simulation [17, 19, 53] in order to incorporate physics-based constraints into the articulated 3d human motion reconstruction. Our approach, *DiffPhy*, relies on gradient-based optimization, connects end-to-end with images, and does not require simplifying assumptions on contacts or the introduction of external non-physical residual forces.
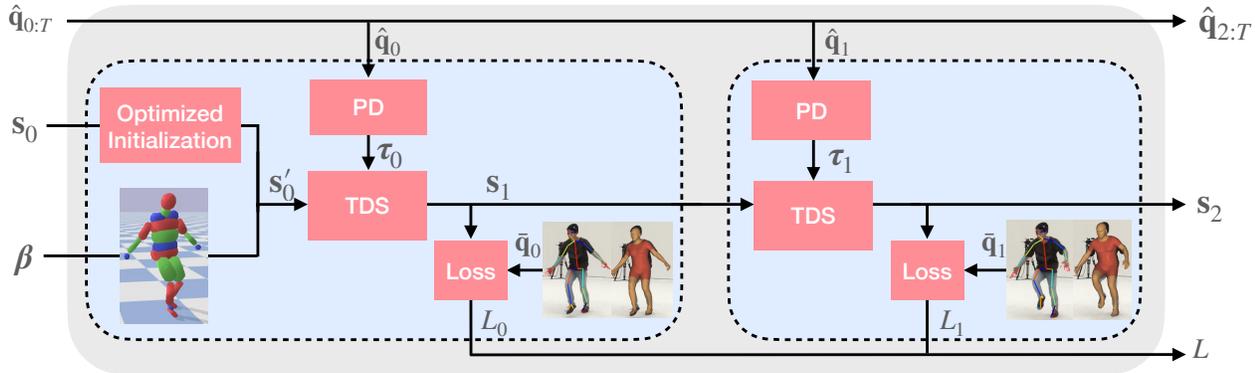
Figure 1. Overview of DiffPhy. Given kinematic estimates (described in §3.1) of a subject's body shape $\boldsymbol{\beta}$, the body's initial pose and velocity $\mathbf{s}_0$, and time-varying 3d poses $\bar{\mathbf{q}}_{0:T}$ with detected 2d keypoints, our model reconstructs the motion in physical simulation, by minimizing a differentiable loss $L$ (see §3.5). DiffPhy optimizes the control trajectory $\hat{\mathbf{q}}_{0:T}$ containing joint angle targets to PD controllers (*cf*. (4)). In turn, the PD controllers compute a torque vector $\boldsymbol{\tau}$, which actuates motors in the joints of the simulated body. DiffPhy integrates a full-featured differentiable simulator, TDS [17] (described in §3.2), that supports complex contacts. Each subject is represented by means of a personalised physical model (see §3.3). In addition, we optimize the initial state (see §3.6), which makes DiffPhy robust to low quality initial estimates. The outputs are 3d pose estimates that align with visual evidence and respect physical constraints.

## 2. Related work

**Kinematics-based 3d Human Pose Estimation.** The problem of monocular 3d pose estimation is usually addressed through end-to-end [30, 31, 62], or two-stage [8, 18] models where neural networks are used to predict 3d joint positions. This is an ill-posed problem due to depth ambiguities and occlusion. The networks are usually trained on vast pose datasets [21, 22, 29, 51] which usually supports good performance on poses previously observed during training. Several methods [24, 61, 63] directly regress the parameters of statistical body models [27, 57] (rather than 3d joint positions), including the subject's body shape as well as kinematic pose. The methods mentioned above take a purely visual inference approach to the problem and do not consider physics-based constraints. As observed by [40], this may cause artifacts such as jitter, ground-penetration, foot sliding, or unnatural leaning [43].

**Physics-based 3d Human Pose Estimation.** Recent work [15, 28, 40–43, 56, 60] aims to increase realism, by using physics to regularize reconstruction. This aims to enforce physical constraints such as proper contact and dynamic coherence. In [40] motion is reconstructed through optimization, but the method only accounts for collisions between the feet and the ground. Such simplifications are recurring in current approaches and limit the types of motions that can be reconstructed. In contrast, in this work, we use a full-featured physical simulator which supports contacts between all objects in the scene. PhysCap [43] is a real-time optimization-based approach, where feet contact is pre-detected based on a neural network. During the physics-based inference, contacts are considered fixed and thherefore cannot be corrected or improved. Moreover, following [59] the method uses non-physical "residual forces"

which improve 3d joint reconstruction metrics at the cost of altered physical plausibility. Since we aim to increase the physicality of reconstructed motions, we avoid using any residual forces. [60] follows on [34, 54] to learn a neural network that estimates torques to drive a model in the full-featured physical simulator MuJoCo [48]. However, MuJoCo is non-differentiable, hence the need to resort to expensive training using numerical gradients in a reinforcement learning setting. The method is trained for millions of steps using 3d ground-truth labels from a motion capture dataset, but the method's ability to generalize to in-the-wild is not demonstrated. Similarly to [40], the method assumes a known ground plane, whereas *DiffPhy* estimates it. [42] integrates a simplified physics approach, dubbed "physionical", into a neural network that estimates joint torques and ground-reaction forces. Similarly to [43] they detect foot contact using a neural network predictor rather than by means of physical simulation. Most recently, [56] introduced a method relying on a simplified physical formulation that makes it possible to refine 3D pose estimates well enough to train motion synthesis models based on that output. However, the method assumes a known ground plane, models only foot contact, and implements a simplified physical body scaled solely based on the estimated bone length rather than shape estimates. Finally, in our concurrent work [15], we perform physics-based human pose reconstruction of complex motions through trajectory optimization based on CMA-ES [16] in the non-differentiable simulator Bullet [7]. This general approach uses a mature and full-featured simulator which, while capable, is slow due to costly black-box optimization. The method does not optimize the initial state of the body (see §3.6) together with the joint control variables, being more vulnerable to unfavorable initialisation. In summary, this work takes the novel

| Method | Body | Cont. | DP | Trained | $\mathbf{T}_g$ | No RF |
|---|---|---|---|---|---|---|
| Rempe *et al.* [40] | Fixed | Feet | ✗ | Contacts | ✗ | ✓ |
| PhysCap [43] | Fixed | Feet | ✓ | Contacts | ✓ | ✗ |
| SimPoE [60] | Adapt | Full | ✗ | Yes | ✗ | ✗ |
| Shimada *et al.* [42] | Fixed | Feet | ✓ | Yes | ✓ | ✗ |
| Xie *et al.* [56] | Fixed | Feet | ✓ | No | ✗ | ✓ |
| Dynamics [15] | Adapt | Full | ✗ | Prior | ✓ | ✓ |
| DiffPhy | Adapt | Full | ✓ | No | ✓ | ✓ |

Table 1. Feature comparison against other physics-based methods. *Body* compares the type of physical body representation where "adapt" means individually constructed based on shape estimate, *Cont.* column compares what type of contacts are supported, *DP* whether the method uses a differentiable physical formulation, *Training* if the physical inference requires training, $\mathbf{T}_g$ compares if the ground plane is estimated (as opposed to assumed known), and *No RF* if the method avoid non-physical residual forces. Only our method does not require any additional training and uses a full-featured differentiable physics formulation.

approach of tightly integrating physics into the reconstruction process through a full-featured *differentiable* physics model. As a result, DiffPhy supports complex full-body contacts, connects pixels-to-physics using end-to-end differentiable losses, supports personalised body models, does not resort to residual forces, and is robust to poor initialization. See tab. 1 for an overview of physics-based methods.

It is worth mentioning that, aside from physical simulation, there exist many other approaches to grounding the human pose estimates using e.g., inertial estimates from IMUs [58], scene constraints [5,64], and motion priors [39].
**Differentiable Physics for Human Modeling.** Physical simulation is a mature area with several established simulation engines available [7, 25, 48]. These engines implement forward simulation but do not facilitate the computation of derivatives necessary for efficient gradient-based optimization. These simulators are well-suited for training with gradient-free methods such as reinforcement-learning or evolutionary algorithms and have been used for gradient-free optimization of human motion models [2,35,60]. More recently differentiable physics simulators have emerged [6,14,17,38,53]. Applying these to human motion reconstruction is difficult due to noisy gradients [19,33], and a non-convex objective function. We present a methodology using gradient-based local search with stochastic global optimization enabling the first use of a full-featured differentiable physics model [17] for human pose reconstruction from video. Furthermore, we show that our approach is magnitudes faster than a purely sampling-based approach.

## 3. Methodology

This section presents our approach to reconstructing 3d human shape and motion from video with differentiable physics in the loop. Given a monocular video of a human subject, we use a kinematic neural network to estimate 2d



Figure 2. Qualitative results on two in-the-wild sequences. Sports and dynamic activities are rarely found in motion capture datasets.

body keypoints, the body shape, and 3d body poses. Since estimating 3d pose from monocular video is ill-posed, due to e.g. depth-ambiguities and occlusion [44], the kinematic 3d reconstructions may suffer from self-penetration, inconsistent translation, jitters, floating above the ground, and non-physical leaning [40, 43]. We, therefore, reconstruct the motion in physical simulation, by jointly accounting for both visual evidence and the constraints of physical simulation (e.g. collisions, gravity, and Newton's laws of motion). See fig. 1 for an overview of our approach.

### 3.1. Kinematic Initialization

Given a sequence of monocular images $\{I_i\}$, we assume a pinhole camera with intrinsics $\mathbf{i} = [f_x, f_y, c_x, c_y]$ and constant camera extrinsics. We obtain the visual evidence used in our optimization objectives following the procedure introduced in [15]. This relies on HUND [61], a 3d pose estimator that produces per-frame 2d keypoints $\bar{\mathbf{x}}_i$ with confidence scores $\mathbf{c}_i$, 3d body poses $\boldsymbol{\theta}_i$, and 3d body shape $\boldsymbol{\beta}_i$, where $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ are the GHUM [57] posing and shape parameters, respectively.

Since HUND is a per-frame estimator, a temporally consistent shape is recovered by selecting the $N = 5$ highest-scoring frames according to keypoint confidences. For these frames, HUND image losses [61] are minimized using BFGS under the additional constraint of a constant shape, $\boldsymbol{\beta}$, across all frames. In addition, [15] introduces a final round of optimization where poses are updated under the time-consistent body shape and a temporal smoothness loss to reduce jittering.

Finally, as the ground plane location is not assumed to be known and HUND produces estimates in camera space $\mathbf{k}$, we estimate the global transform $\mathbf{T}_g \in \mathbb{R}^{3 \times 4}$ for the physical scene, with gravity along the y axis, as well as the ground plane at $y = 0$. This is achieved by minimizing

$$L_g(\mathbf{T}_g) = \sum_i^N \| \min(\delta, \mathbf{L}_y(\mathbf{T}_g[\mathbf{M}(\boldsymbol{\beta}, \boldsymbol{\theta}_i), 1])) \|^2, \quad (1)$$

where $\mathbf{L}_y$ is an operator that extracts the $k = 20$ smallest signed distances from the mesh vertices $\mathbf{M}(\boldsymbol{\beta}, \boldsymbol{\theta}_i)$ after the global transformation. This assumes the body is in ground plane contact for most of the sequence. To allow for frames where the subject is not in contact with the ground, we clip the maximum shortest distance to the ground to $\delta = 20$ cm.

## 3.2. Differentiable Physics Simulation Model

We implement our models in the framework of the "Tiny Differentiable Simulator" (TDS) [17]. This formulates rigid-body dynamics for articulated bodies in terms of reduced coordinates. Elements in the vector $\mathbf{q}$ represent the position of each joint, and elements in the vector $\dot{\mathbf{q}}$ represent joint space velocities, based on revolute and spherical joints. Given the state of the body $\vec{s}_t = (\mathbf{q}_t, \dot{\mathbf{q}}_t)$ at time $t$, as well as the vector of joint torques $\boldsymbol{\tau}_t$, and external forces $\mathbf{f}_t$, the computation shown in fig. 3 produces a new body state $\vec{s}_{t+\delta t}$ corresponding to the rigid multi-body dynamics with contacts. To that end, we first run forward kinematics to compute world space positions and velocities, as well as forward dynamics to compute unconstrained acceleration obtained without taking contacts into account. The forward dynamics computes the acceleration by solving the equation of motion for the kinematic tree given by

$$\boldsymbol{\tau}_t = \mathbf{H}(\mathbf{q}_t)\ddot{\mathbf{q}}_t + \mathbf{C}(\mathbf{q}_t, \dot{\mathbf{q}}_t, \mathbf{f}_t^x) \tag{2}$$

where $\mathbf{H}(\mathbf{q})$ is the joint-space inertia matrix, $\mathbf{C}$ the a joint space bias force and $\mathbf{f}^x$ is the vector of external forces. The forward dynamics is computed by propagation-based Articulated-Body Algorithm (ABA) [11] that traverses the kinematic chain of the body three times in order to compute quantities necessary to finally obtain the acceleration of each rigid component of the body[1]. The joint-space inertia matrix is computed using the Composite Rigid Body Algorithm (CRBA) [11].

Unconstrained accelerations $\ddot{\mathbf{q}}_{t+\delta t}^u$ are then used to compute unconstrained velocities, which in conjunction with the output of the forward kinematics $\mathbf{x}_{t+\delta t}$ are used to update the contact points between the articulated body and the environment. Contact points with positive (separating) distance are classified as inactive, while contact points with zero or negative distance are active. Active contacts generate a repulsive impulse that needs to be taken into account when computing the new body state. To that end, the forward dynamics computation is phrased as a linear complementarity problem (LCP) at the velocity level [46, 47]

$$\mathbf{J}_c \mathbf{H}^{-1} \mathbf{J}_c^\top \mathbf{p} + \mathbf{J}_c \dot{\mathbf{x}} = \mathbf{v} \tag{3}$$
$$\mathbf{v} = [\mathbf{v}_u, \mathbf{v}_b]$$
$$\text{s.t.} \quad \mathbf{v}_u^\top \mathbf{p}_u = 0 \quad \mathbf{v}_u \geq 0 \quad \mathbf{p}_u \geq 0 \quad \mathbf{v}_b = 0$$

where $\mathbf{J}_c$ is a contact Jacobian for the positions of contact points computed in the previous step, $\mathbf{p}$ is the vector of reaction impulses, and $\mathbf{v}$ is the vector of relative velocities. The indices $u$ and $b$ indicate the unilateral and bilateral portion of constraints, respectively. The LCP problem in (3) is then iteratively solved with a projected Gauss-Seidel method following the formulation in [47], by relying on a

---
[1]See tab. 7.1 in [11] for the Articulated-Body Algorithm.

per-contact LCP [20]. The final step of the computation is to obtain joint positions $\mathbf{q}_{t+\delta t}$ from joint velocities using semi-implicit Euler integration.
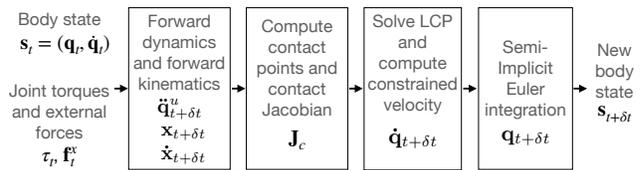


Figure 3. Overview of the simulation step of the physics model that updates the current state $S_t$ to a new state after time step $\delta_t$. For each computational block we include the output quantities used in the subsequent block.

## 3.3. Physical Human Body Modeling

In the physical simulation, we model the human body as rigid geometric primitives connected by joints. The model is comprised of 16 joints with a total of 48 degrees of freedom, joining together 26 capsules that represent the various body parts (*cf.* fig. 1). The shape and mass of the model is automatically adapted for various body shapes by relying on a statistical body model [57]. Given the 3d mesh $\mathbf{M}(\boldsymbol{\beta}, \emptyset)$ corresponding to a shape estimate $\boldsymbol{\beta}$ in rest pose, we infer the dimensions of the geometric primitives following the approach of [2]. The process is entirely automatic and yields individualized physical models for each subject. As a physical model requires mass, we first estimate the total mass of the body based on a human shape dataset [36] then distribute the weight according to an anatomical distribution [37]. Finally, the inertia of each primitive is computed based on its mass and dimensions.

DiffPhy reconstructs a motion in simulation by actuating torque motors in the joints of the body. Following prior work [1] we optimize over control targets to proportional-derivative (PD) controllers rather than over the torques directly. We define the body's angular joint positions as $\mathbf{q}_t$, and joint velocities as $\dot{\mathbf{q}}_t$, the associated 3d Cartesian coordinates of the joints as $\mathbf{x}_t$ for the time step $t$. Given a set of joint targets $\hat{\mathbf{q}}_{1:T} = \{\hat{\mathbf{q}}_1, \hat{\mathbf{q}}_2, \ldots, \hat{\mathbf{q}}_t\}$ the PD controllers infer the joint torques as

$$\boldsymbol{\tau}_t = k_p(\hat{\mathbf{q}}_t - \mathbf{q}_t) - k_d \dot{\mathbf{q}}_t, \tag{4}$$

where $k_p$ and $k_d$ are gain parameters of PD controllers. We may then specify a motion of length $T$ as the initial state $\mathbf{s}_0 = (\mathbf{q}_0, \dot{\mathbf{q}}_0)$, the world geometry $\mathbf{G}$ defining the position and orientation of the ground plane, and a target trajectory for the joints $\hat{\mathbf{q}}_{1:T}$. Given the loss presented in (5) we reconstruct the motion by minimizing $L = L(\mathbf{s}_0, \mathbf{G}, \hat{\mathbf{q}}_{1:T})$ with respect to $\hat{\mathbf{q}}_{1:T}$.

## 3.4. Gradient-Based Optimization

Given our loss function $L = L(\mathbf{S}_0, \mathbf{G}, \hat{\mathbf{q}}_{1:T})$ we can use any gradient-based optimization method to minimize the

| Method | # eval | MPJPE-G | MPJPE | MPJPE-PA | MPJPE-2d |
|---|---|---|---|---|---|
| CMA-ES | 80k | 206.7 | 125.7 | 77.4 | 16.9 |
| BFGS | 122 | 160.1 | 100.1 | 68.9 | 15.5 |
| Basin-BFGS | 509 | 144.9 | 84.6 | 61.1 | 12.6 |

Table 2. Comparison of optimization strategies on our Human3.6M validation set. BFGS and Basin-BFGS both use gradients, while CMA-ES is a gradient-free approach. Note that Basin-Hopping together with BFGS (Basin-BFGS) improves the performance of BFGS by combining it with stochastic global optimization. Using only a purely sampling-based approach (CMA-ES) requires magnitudes more function evaluations while still not finding better optima for our loss. BFGS was given a sufficiently large evaluation budget to converge.

loss with respect to $\hat{\mathbf{q}}_{1:T}$. Since the loss function is non-convex, convergence to suboptimal local minima is possible. Therefore, a global optimization combined with a local gradient-based search is expected to outperform a purely local method. One such method is the global stochastic optimization *Basin-Hopping* [52]. It uses a two-stage approach, which alternates between performing gradient-based local search and stochastic global search. Based on an initial candidate, it first performs a local search. It then randomly perturbs the local minimum, performs a local search again on the new candidate, and then either accepts or rejects the new solution based on the Metropolis criterion [32]. In our model, we use BFGS [13] for local optimization.

## 3.5. Optimization Objectives

Reconstructing a motion sequence amounts to finding the control trajectory $\hat{\mathbf{q}}_{1:T}$ that minimizes the reconstruction loss $L$ under the constraints of the simulation dynamics. In this work, we formulate $L$ as a weighted combination of loss functions

$$L = w_r L_r + w_j L_j + w_i L_i + w_l L_l, \quad (5)$$

with the weights $w_r = 10.0$, $w_j = 0.1$, $w_i = 0.01$, and $w_l = 0.01$. The root position loss $L_r$ measures errors between the 3d position of the simulated pelvis root joint $\mathbf{x}_t^{\text{root}}$ and the kinematically estimated position $\bar{\mathbf{x}}_t^{\text{root}}$

$$L_r(\hat{\mathbf{q}}_{1:T}) = \frac{1}{T} \sum_t^T \|\bar{\mathbf{x}}_t^{\text{root}} - \mathbf{x}_t^{\text{root}}\|^2 \quad (6)$$

at time $t$ where $T$ is the total length of the sequence. $L_j$ computes the rotational distance between the kinematic pose estimate and the simulated body's pose

$$L_j(\hat{\mathbf{q}}_{1:T}) = \frac{1}{TK} \sum_t^T \sum_k^K \arccos(|\mathbf{q}_t^k \cdot \bar{\mathbf{q}}_t^k|), \quad (7)$$

where $\bar{\mathbf{q}}_t^k$ and $\mathbf{q}_t^k$ are rotations expressed as quaternions for joint $k$ at time $t$ for kinematics and the simulated character

respectively. Note the difference between $\hat{\mathbf{q}}$ and $\mathbf{q}$, where the former are the PD control targets and the latter are the joint angles of the simulated model (*cf.* (4)). $L_i$ computes the 2d projection loss

$$L_i(\hat{\mathbf{q}}_{1:T}) = \frac{1}{TK} \sum_t^T \sum_k^K \mathbf{c}_t^k \|\bar{\mathbf{x}}_t^k - \Pi(\mathbf{x}_t^k, \mathbf{i})\|^2, \quad (8)$$

where $\Pi(\mathbf{x}_t^k, \mathbf{i})$ is the perspective operator projecting the simulated model's joint $\mathbf{x}_t^k$ onto the image with camera intrinsics $\mathbf{i}$ weighted by the keypoint detection confidence score $c_t^k$. Finally, $L_l$ is a regularizer that penalizes joints outside of human anatomical limits as present in the statistical body model [57]

$$L_l(\hat{\mathbf{q}}_{1:T}) = \frac{1}{TK} \sum_t^T \sum_k^K \| \max(z_{\text{lower}}^k - \mathbf{q}_t^k, 0)$$
$$+ \max(\mathbf{q}_t^k - z_{\text{upper}}^k, 0)\|^2, \quad (9)$$

where $z_{upper}^k$ and $z_{lower}^k$ are upper and lower bounds for joint $k$ respectively.

Note that in the above definitions, the positions of body joints angles $\mathbf{q}_t^k$ and 3d joint positions $\mathbf{x}_t^k$ are dependent on the control trajectory up until time $t$, as part of the physics formulation introduced in §3.2.

## 3.6. Optimized Initialization

We initialize the pose $\mathbf{q}_0$ in the first time step of the simulation to the kinematically estimated pose $\bar{\mathbf{q}}_0$ and estimate the velocity $\dot{\mathbf{q}}_0$ using finite differences between the first two kinematic poses $\{\bar{\mathbf{q}}_0, \bar{\mathbf{q}}_1\}$. However, if the initial kinematic pose estimate is poor, this might lead to a low quality starting pose from which the simulation cannot recover. Similarly, jitters in the kinematic poses may cause a significant error in the estimated initial velocity. We address these issues by including the initial pose and velocity as variables to optimize. We experimentally validate how such a relatively straightforward approach significantly impacts the results.

# 4. Experiments

**Datasets.** We quantitatively evaluate DiffPhy on the Human3.6M [21], and a subset of the AIST [49] pose datasets. The former contains a diverse set of motions from a motion capture laboratory, whereas the latter contains dance videos with triangulated 3d joints as pseudo-ground-truth. As only DiffPhy and SimPoe [60] supports full-body contacts (*cf.* tab. 1), PhysCap [43] proposed evaluating on a subset of the Human3.6M. This protocol eliminated all sequences requiring more than foot-floor contacts. Hence to allow for comparison, we use this subset in tab. 3, but note that our method is more general and supports contacts for all body parts. For ablations, we use 100 frames from 20 sequences
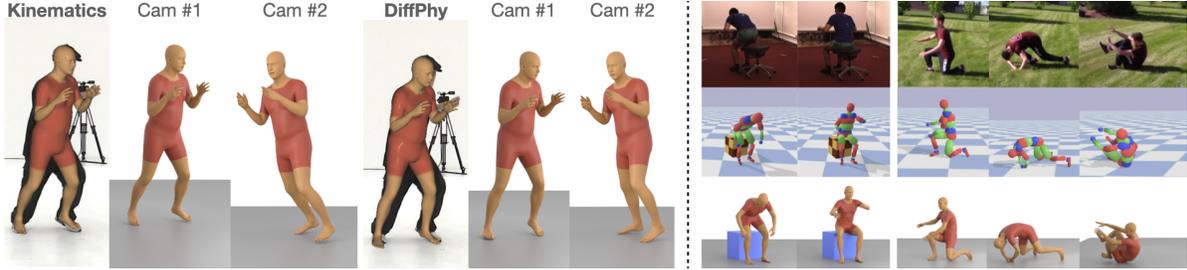
Figure 4. Qualitative examples on the AIST dataset (left) and of complex contacts (right). The AIST example shows that both kinematics and DiffPhy projects well into the image. However, when rendered from another viewpoint (cam #2) it becomes clear that kinematics exhibits unrealistic leaning while the physical constraints corrects the pose to keep the body in balance. See tiny.cc/diffphy for more.

| Dataset | Model | MPJPE-G | MPJPE | MPJPE-PA | MPJPE-2d | TV | Foot skate (%) |
|---------|-------|---------|-------|----------|----------|-----|---------------|
| Human3.6M | VIBE [24] | 207.7 | 68.6 | 43.6 | 16.4 | 0.32 | 27.4 |
| | PhysCap [43] | - | 97.4 | 65.1 | - | - | - |
| | SimPoE [60] | - | **56.7** | **41.6** | - | - | - |
| | Shimada *et al.* [42] | - | 76.5 | 58.2 | - | - | - |
| | Xie *et al.* [56] | - | 68.1 | - | - | - | - |
| | Kinematics | 145.3 | 83.0 | 55.4 | 13.4 | 0.34 | 47.5 |
| | DiffPhy | **139.1** | 81.7 | 55.6 | **13.1** | **0.20** | **7.4** |
| AIST | Kinematics | 155.7 | 107.4 | 66.9 | **10.4** | 0.52 | 50.9 |
| | DiffPhy | **150.2** | **105.5** | **66.0** | 12.1 | **0.44** | **19.6** |

Table 3. Quantitative evaluation on the Human3.6M and AIST datasets. Our full dynamic model improves over the kinematic estimates used as initialization with respect to standard joint position error metrics as well as reducing motion jitter and unnatural foot skating.

from a validation subset of Human3.6M. Finally, we quantitatively evaluate our method on real-world internet videos released under creative commons licenses. For additional details, refer to our supplementary material.

**Metrics.** We report the standard pose metrics such as mean per-joint position error in millimeter (MPJPE-G), mean Procrustes aligned joint error (MPJPE-PA), per-frame translation aligned error (MPJPE), and 2d mean per-joint error in pixels (MPJPE-2d). Note that many papers do not report global position errors since they consider only root-relative poses. We, however, are interested in measuring the pose error, including translational errors, since unnatural translation is a common (non-physical) reconstruction artifact. In addition, we also measure foot skating and the total variation in the joint acceleration per frame (TV). We measure foot skating as percentage of frames where a foot moves more than 2cm while in contact with the ground in two adjacent frames. Unlike [40], we do not assume foot contact annotations but instead heuristically detect foot contacts based on the distance between the foot mesh and the ground-plane. The total variation in acceleration is computed as $\frac{1}{T} \sum_{t \in T} \sum_{k \in K} |\ddot{x}_{t+1}^k - \ddot{x}_t^k|$, for the 3d acceleration $\ddot{x}_t^k$ of joint $k$ at time $t$ estimated using finite differences. Thus, high TV indicates motion jitter, and high foot skate implies motion that slides along the ground.

**Implementation Details.** We use the Tiny Differentiable Simulator [17] running at $1,000$ Hz with the gradients computed using the auto differentiation framework CppAD [3].

In addition, we use a Python implementation of Basin-Hopping and BFGS [50]. Since the length of the optimized trajectory may be great, we follow [1] and perform optimization in overlapping windows of length $N = 960$. The simulation steps take $\approx 5$s. For the large datasets in tab. 3, we compute the windows in parallel and stitch them together in order to speed up computation. We initialize the control targets $\hat{q}_{1:T}$ to 3d poses estimated by our kinematics. See our supplementary material for details.

### 4.1. Results

We compare DiffPhy against both state-of-the-art kinematic video models (VIBE [24]) and against physics-based methods. The results are summarized in tab. 3. Since VIBE predicts root-relative poses, we estimate the global translation (required to compute MPJPE-G) by minimizing 2d projection errors using a method similar to the one in [43]. For VIBE, we use the publicly available implementation. For the other methods, we give numbers presented by the authors. On both Human3.6M and AIST, our model improves with respect to the physical metrics (TV and foot skate) compared to the kinematic initialization. On Human3.6M, foot skating is only $7.4\%$ compared to $47.5\%$ for the kinematic initialization and $27.4\%$ for VIBE. On AIST, foot skating is reduced from $50.9\%$ to $19.6\%$. We believe that increased skating on AIST is due to actual skating motions performed as part of the hip-hop dances. On total variation, our model similarly improves over kinematics with $0.20$ and $0.44$ on Human3.6M and AIST, respectively. Fur-
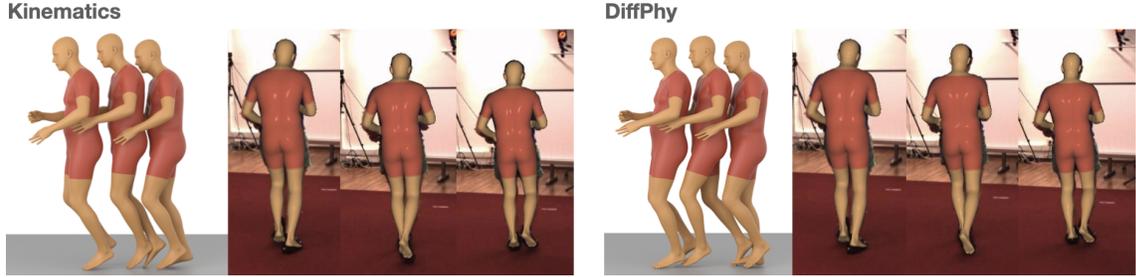
Figure 5. Qualitative examples on Human3.6M. DiffPhy infers plausible leg motion while kinematics skates unrealistically forward.

thermore, we note that our full model improves the global joint position error (MPJPE-G), a metric that measures pose and translation errors. On Human3.6M, DiffPhy has an error of 139.1 compared to 145.3 and 207.7 mm/joint for kinematics and VIBE, respectively. If we look at the error for foot joints only, we see an even larger improvement by including physics compared to kinematics alone (166.8 vs. 174.1 mm/joint). This result aligns with prior work [40], showing that physics improves foot position estimation. Furthermore, our method aligns well with image evidence when comparing 2d error, i.e., 13.1 px/joint vs. VIBE's 16.4 px/joint on Human3.6M. In terms of joint error including translation error (MPJPE), SimPoE [60], Xie *et al*. [56], Shimada *et al*. [42] outperform DiffPhy (56.7 vs. 68.1 vs. 76.5 vs. 81.7 mm/joint respectively), though in the case of SimPoE and Xie *et al*. this might stem from initializing from the already strong VIBE predictor (68.6 mm/joint). Furthermore, SimPoE is a neural network requiring extensive training using the 3d ground-truth from Human3.6M, whereas DiffPhy is a general method that requires no additional training (*cf*. tab. 1). Xie *et al*., PhysCap [43], Shimada *et al*. [42] on the other hand, focus only on feet-ground contacts while DiffPhy supports complex contacts. Unfortunately, this advantage cannot be demonstrated on subsets that exclude sequences with complex contacts.

Fig. 5 presents qualitative results where kinematics fails to estimate the positions of legs due to depth ambiguities. The reconstructed poses align well when projected into the image but are unrealistic since the model skates rather than walks forward. Since DiffPhy reconstructs the motion with physics in the loop, it must propel the model forward through bipedal locomotion, thus inferring feasible leg poses. Similarly, in fig. 4 kinematics estimates a pose that projects well into the image. However, when viewed from a side, it becomes clear that kinematics estimates a pose that leans unnaturally. Since DiffPhy is constrained by gravity, it must find a pose that is both physically plausible *and* aligns with 2d evidence. Fig. 4 also includes examples of object interactions and rolling motions requiring complex contacts. We manually modeled the chair as a box since DiffPhy does not estimate scene geometry. For the rolling motion, the kinematics were too noisy for DiffPhy to converge; hence, we manually corrected the worst kinematic

frames before running DiffPhy. Finally, fig. 2 shows two reconstruction examples for sequences in-the-wild. These videos exhibit poses and activities missing from standard laboratory-captured datasets.

**Ablation studies.** In tab. 6 we validate our choice of loss components in (5). We note that the 2d projection loss, as expected, plays an important role in aligning the reconstruction with the image evidence (17.1 vs. 12.6 px/joint). Furthermore, since 2d keypoints do not suffer from depth ambiguities, they are generally more reliable than 3d keypoints and thus serve as a strong signal. Therefore removing 2d evidence significantly increases MPJPE-G from 144.9 to 158.5 mm/joint. Removing the root position loss (6) has the largest impact on global position error (165.7 mm/joint) since without it, we do not provide DiffPhy with any supervision with respect to world positioning. This allows for suboptimal reconstructions that align well with the projected image (12.8 px/joint) but do not transition correctly in world space. Without the joint angle loss (7), DiffPhy is deprived of the per-frame 3d pose estimates, which, when predicted by neural networks such as HUND or VIBE that are trained on large pose datasets, provide useful guidance as long as their predictions do not contradict any physical constraints. Removing the joint angle limit regularizer (*cf*. (9)) demonstrates the usefulness of constraining the reconstructed motion to the space of anatomically valid poses even for everyday motions like those in Human3.6M. Finally, we validate the usefulness of optimizing the initial starting pose and velocity (see §3.6). Without it, the kinematic estimates for the initial frames must be accurate. If not, the simulation may start from an initial state from which DiffPhy may fail to recover, as seen by the largest MPJPE-PA in the ablation of 65.1 mm/joint.

Next, results in tab. 2 show that gradient-based methods are vastly more efficient for our physics loss compared to the commonly used gradient-free approach CMA-ES [16]. BFGS obtains a lower MPJPE-G error (160.1 vs 206.7 mm/joint), and requires a fraction of the computations (122 vs. 80k loss evaluations per windows). Next, We note that BFGS converges to suboptimal minima, but by combining BFGS with Basin-Hopping, we can reduce the errors further to 144.9 mm/joint. As Basin-Hopping can explore infinitely

| RF | MPJPE-G | MPJPE | MPJPE-PA | MPJPE-2d |
|---|---|---|---|---|
| 0 | 144.9 | 84.6 | 61.1 | 12.6 |
| 5 | 141.4 | 82.2 | 60.7 | 11.8 |
| 10 | 140.1 | 79.9 | 60.2 | 11.7 |
| 25 | 146.3 | 81.9 | 60.0 | 12.7 |
| 50 | 140.2 | 79.4 | 60.3 | 11.6 |
| 100 | 154.0 | 87.7 | 61.5 | 14.4 |

Table 4. Results on experiments on the effects of residual force. We note using a residual force decreases the error metrics, but we refrain from using it to avoid unexplained non-physical forces.

| Window | MPJPE-G | MPJPE | MPJPE-PA | MPJPE-2d |
|---|---|---|---|---|
| 240 | 390.1 | 224.1 | 96.6 | 40.3 |
| 480 | 165.6 | 97.2 | 63.8 | 13.2 |
| 720 | 148.9 | 87.2 | 61.8 | 12.6 |
| 960 | 144.9 | 84.6 | 61.1 | 12.6 |
| 1440 | 155.6 | 92.5 | 65.7 | 15.9 |

Table 5. Results on the effects of optimization window size. A balance needs to be found between a larger window size which allow for more visual evidence to be taken into account while a smaller reduces the dimensionality of the search space.

many basins, we set the limit to 5 basin steps, each with 50 BFGS iterations as a trade-off between accuracy and speed.

In tab. 5 we study the effect of the optimization window size. We find that a window of 960 simulation steps (containing 0.96s of video) is optimal for our setup. A larger window size increases the errors, most likely due to a larger search space combined with a larger gradient variance, as noted in [33]. On the other hand, smaller windows provide scarcer visual evidence and are sensitive to a few occluded frames, or to noisy estimates. Interestingly, a smaller window size performed better for experiments on ground-truth data (see supplementary material). This indicates that smaller apertures are better for noise-free inputs.

Several methods (*cf*. tab. 1) introduce "residual forces" acting on the root link of the physical body. This non-physical force allows the method to translate and rotate the body to align with visual evidence at the expense of physical realism. tab. 4 confirms that this indeed can be used to lower DiffPhy's joint errors (MPJPE-G from 144.9 to 140.2 mm/joint and 2d error from 12.6 to 11.6 px/joint when applying 50N for each of the six degrees of freedom). Interestingly, applying a too great residual force (100N) increased error, perhaps since it allows the model to circumvent some of the constraints of physical simulation. In this work, we avoid using residual forces, in order to keep all forces realistic, and avoid non-physical artifacts.

## 5. Discussion

In order to improve the realism of 3d human sensing, we have introduced *DiffPhy* – the first differentiable physics-based model for full-body articulated human motion estimation, that supports complex contacts, does not assume a

| Variant | MPJPE-G | MPJPE | MPJPE-PA | MPJPE-2d |
|---|---|---|---|---|
| Full model | 144.9 | 84.6 | 61.1 | 12.6 |
| No root | 165.7 | 84.8 | 60.7 | 12.8 |
| No 2d | 158.5 | 98.3 | 65.7 | 17.1 |
| No pose | 156.8 | 91.8 | 64.4 | 13.0 |
| No 3d loss | 216.6 | 122.4 | 76.3 | 12.6 |
| No limits | 146.8 | 86.5 | 62.2 | 13.0 |
| No opt. init. | 151.5 | 92.1 | 65.1 | 14.1 |

Table 6. Ablation of the model components introduced in §3. *No root* means without root position loss (6), *No 2d* without 2d keypoint loss (8), *No pose* without joint angle loss (7), *No 3d loss* without both root link position loss and joint angles losses, *No limit* without anatomical joint limits (9), and *No init. opt.* is without optimizing the initial state, *cf*. §3.6.

known ground plane, and avoids reliance on non-physical forces. This has the benefit of a human model with realistic physics interactions, that are constrained end-to-end by visual losses. Furthermore, such a model can provide a valuable non-learning-based component, which is always valid, complementing the statistical kinematic prediction and optimization techniques prevalent in the current state of the art. Visual 3d human motion reconstruction experiments on multiple datasets demonstrate that our methodology is competitive with other state of the art physics-based approaches.

**Limitations and Future Work.** An inherent limitation to physics-based approaches is the need to model objects in the scene. We hope to address this challenge in future work by integrating with 3d scene reconstruction techniques [4]. Ideally, we would be able to jointly optimize the control of the body and the world to match visual evidence. Another limitation is our current assumption of constant camera extrinsics. This limits our technique to videos captured using a static camera but can be easily relaxed. Finally, our reconstructions are limited to a single subject. Reconstructing multiple people interacting is interesting since these scenes are complex, and learning statistical models of interaction between humans is challenging [12]. A physics-based approach could help infer constraints and affordances.

**Ethical Considerations.** Our construction of physics-based models is motivated by the breadth of transformative 3d applications that would become possible, including fitness, personal well-being or special effects, or human-computer interaction, among others. In contrast, applications like visual surveillance and person identification would not be effectively supported, given that the model's output does not provide sufficient detail for these purposes. The same is true for the creation of potentially adversely-impacting deepfakes, as an appearance model or a joint audio-visual model are not included for photorealistic visual and voice synthesis. While our method is fundamentally applicable to a variety of human body types, we have not evaluated this aspect extensively and consider such a study an important objective for future work.

# References

[1] Mazen Al Borno, Martin de Lasa, and Aaron Hertzmann. Trajectory optimization for full-body movements with complex contacts. In *IEEE transactions on visualization and computer graphics*, volume 19, pages 1405–14, 08 2013. 4, 6

[2] Mazen Al Borno, Ludovic Righetti, Michael J. Black, Scott L. Delp, Eugene Fiume, and Javier Romero. Robust Physics-based Motion Retargeting with Realistic Body Shapes. In *Computer Graphics Forum*, 2018. 1, 3, 4

[3] B. Bell. Cppad: a package for c++ algorithmic differentiation, 2021. 6, 13

[4] Bharat Lal Bhatnagar, Xianghui Xie, Ilya Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2022. 8

[5] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *European Conference on Computer Vision*, pages 387–404. Springer, 2020. 3

[6] Justin Carpentier, Guilhem Saurel, Gabriele Buondonno, Joseph Mirabel, Florent Lamiraux, Olivier Stasse, and Nicolas Mansard. The pinocchio c++ library – a fast and flexible implementation of rigid body dynamics algorithms and their analytical derivatives. In *IEEE International Symposium on System Integrations (SII)*, 2019. 3

[7] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. http://pybullet.org, 2016–2019. 2, 3

[8] Rishabh Dabral, Anurag Mundhada, Uday Kusupati, Safeer Afaque, Abhishek Sharma, and Arjun Jain. Learning 3d human pose from structure and motion. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 668–683, 2018. 2

[9] Filipe de Avila Belbute-Peres, Kevin Smith, Kelsey Allen, Josh Tenenbaum, and J. Zico Kolter. End-to-end differentiable physics for learning and control. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. 1

[10] Roy Featherstone. *Rigid Body Dynamics Algorithms*. Springer-Verlag, Berlin, Heidelberg, 2007. 1

[11] Roy Featherstone. *Rigid Body Dynamics Algorithms*. Springer-Verlag, Berlin, Heidelberg, 2007. 4

[12] Mihai Fieraru, Mihai Zanfir, Teodor Szente, Eduard Bazavan, Vlad Olaru, and Cristian Sminchisescu. Remips: Physically consistent 3d reconstruction of multiple interacting people under weak supervision. *Advances in Neural Information Processing Systems*, 34, 2021. 8

[13] Roger Fletcher. *Practical Methods of Optimization*. John Wiley & Sons, New York, NY, USA, 1987. 5

[14] C. Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax - a differentiable physics engine for large scale rigid body simulation, 2021. 3

[15] Erik Gärtner, Mykhaylo Andriluka, Hongyi Xu, and Cristian Sminchisescu. Trajectory optimization for physics-based reconstruction of 3d human pose from monocular video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2, 3

[16] Nikolaus Hansen. *The CMA Evolution Strategy: A Comparing Review*, pages 75–102. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006. 2, 7

[17] Eric Heiden, David Millard, Erwin Coumans, Yizhou Sheng, and Gaurav S Sukhatme. NeuralSim: Augmenting differentiable simulators with neural networks. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2021. 1, 2, 3, 4, 6, 13

[18] Mir Rayat Imtiaz Hossain and James J Little. Exploiting temporal information for 3d human pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 68–84, 2018. 2

[19] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-Kelley, and Frédo Durand. Difftaichi: Differentiable programming for physical simulation. *arXiv preprint arXiv:1910.00935*, 2019. 1, 3

[20] Jemin Hwangbo, Joonho Lee, and Marco Hutter. Percontact iteration method for solving contact dynamics. *IEEE Robotics and Automation Letters*, 3(2):895–902, 2018. 4

[21] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, jul 2014. 2, 5, 12, 13

[22] H. Joo, T. Simon, and Y. Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8320–8329, 2018. 2

[23] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 1

[24] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. Vibe: Video inference for human body pose and shape estimation. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2020. 1, 2, 6

[25] Jeongseok Lee, Michael X. Grey, Sehoon Ha, Tobias Kunz, Sumit Jain, Yuting Ye, Siddhartha S. Srinivasa, Mike Stilman, and C. Karen Liu. Dart: Dynamic animation and robotics toolkit. *Journal of Open Source Software*, 3(22):500, 2018. 3

[26] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to dance with aist++: Music conditioned 3d dance generation, 2021. 12

[27] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 2

[28] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. *Advances in Neural Information Processing Systems*, 34, 2021. 2

[29] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision*, pages 5442–5451, Oct. 2019. 2

[30] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Srinath Sridhar, Gerard Pons-Moll, and Christian Theobalt. Single-shot multi-person 3d pose estimation from monocular rgb. In *2018 International Conference on 3D Vision (3DV)*, pages 120–130. IEEE, 2018. 2

[31] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. *ACM Transactions on Graphics (TOG)*, 36(4):1–14, 2017. 2

[32] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092, 1953. 5

[33] Luke Metz, C. Daniel Freeman, Samuel S. Schoenholz, and Tal Kachman. Gradients are not all you need, 2021. 3, 8, 13

[34] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graph.*, 37(4):143:1–143:14, July 2018. 2

[35] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. In *SIGGRAPH*, 2018. 3

[36] Leonid Pishchulin, Stefanie Wuhrer, Thomas Helten, Christian Theobalt, and Bernt Schiele. Building statistical shape spaces for 3d human modeling. *Pattern Recognition*, 2017. 4

[37] Stanley Plagenhoef, F Gaynor Evans, and Thomas Abdelnour. Anatomical data for analyzing human motion. *Research quarterly for exercise and sport*, 54(2):169–178, 1983. 4

[38] Yi-Ling Qiao, Junbang Liang, Vladlen Koltun, and Ming C Lin. Efficient differentiable simulation of articulated bodies. In *International Conference on Machine Learning*, pages 8661–8671. PMLR, 2021. 3, 13

[39] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. In *International Conference on Computer Vision (ICCV)*, 2021. 3

[40] Davis Rempe, Leonidas J. Guibas, Aaron Hertzmann, Bryan Russell, Ruben Villegas, and Jimei Yang. Contact and human dynamics from monocular video. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 3, 6, 7

[41] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick Pérez, and Christian Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics*, 40(4), aug 2021. 2

[42] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick P'erez, and Christian Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics (TOG)*, 40:1 – 15, 2021. 2, 3, 6, 7

[43] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. *ACM Transactions on Graphics*, 39(6), dec 2020. 1, 2, 3, 5, 6, 7

[44] C. Sminchisescu and B. Triggs. Kinematic jump processes for monocular 3d human tracking. In *CVPR*, 2003. 3

[45] Jakub Stepien. *Physics-Based Animation of Articulated Rigid Body Systems for Virtual Environments*. PhD thesis, 10 2013. 1

[46] David Stewart and J.C. (Jeff) Trinkle. An implicit time-stepping scheme for rigid body dynamics with coulomb friction. volume 1, pages 162–169, 01 2000. 4

[47] Jakub Stępień. *Physics-Based Animation of Articulated Rigid Body Systems for Virtual Environments*. PhD thesis, Silesian University of Technology, 2013. 4

[48] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012. 2, 3

[49] Shuhei Tsuchida, Satoru Fukayama, Masahiro Hamasaki, and Masataka Goto. Aist dance video database: Multi-genre, multi-dancer, and multi-camera database for dance information processing. In *Proceedings of the 20th International Society for Music Information Retrieval Conference, ISMIR 2019*, pages 501–510, Delft, Netherlands, Nov. 2019. 5, 12, 13

[50] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. 6

[51] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European Conference on Computer Vision (ECCV)*, sep 2018. 2

[52] David J. Wales and Jonathan P. K. Doye. Global optimization by basin-hopping and the lowest energy structures of lennard-jones clusters containing up to 110 atoms. *The Journal of Physical Chemistry A*, 101(28):5111–5116, 1997. 5

[53] Keenon Werling, Dalton Omens, Jeongseok Lee, Ionnis Exarchos, and C. Liu. Fast and feature-complete differentiable physics for articulated rigid bodies with contact. *ArXiv*, abs/2103.16021, 2021. 1, 3

[54] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. A scalable approach to control diverse behaviors for physically simulated characters. *ACM Trans. Graph.*, 39(4), 2020. 2

[55] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2019. 1

[56] Kevin Xie, Tingwu Wang, Umar Iqbal, Yunrong Guo, Sanja Fidler, and Florian Shkurti. Physics-based human motion estimation and synthesis from videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11532–11541, October 2021. 1, 2, 3, 6, 7, 12

[57] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. GHUM & GHUML: Generative 3d human shape and articulated pose models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6184–6193, 2020. 2, 3, 4, 5

[58] Xinyu Yi, Yuxiao Zhou, and Feng Xu. Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. *ACM Trans. Graph.*, 40(4), jul 2021. 3

[59] Ye Yuan and Kris Kitani. Residual force control for agile human behavior imitation and extended motion synthesis. In *Advances in Neural Information Processing Systems*, 2020. 2

[60] Ye Yuan, Shih-En Wei, Tomas Simon, Kris Kitani, and Jason Saragih. Simpoe: Simulated character control for 3d human pose estimation. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 3, 5, 6, 7

[61] Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Neural descent for visual 3d human pose and shape. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3

[62] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2148–2157, 2018. 2

[63] Mihai Zanfir, Andrei Zanfir, Eduard Gabriel Bazavan, William T. Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Thundr: Transformer-based 3d human reconstruction with markers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021. 1, 2

[64] Jason Y. Zhang, Sam Pepose, Hanbyul Joo, Deva Ramanan, Jitendra Malik, and Angjoo Kanazawa. Perceiving 3d human-object spatial arrangements from a single image in the wild. In *European Conference on Computer Vision (ECCV)*, 2020. 3

# Appendix

This supplement presents additional results (§A), a description of the datasets used (§B) together with a description of the usage of data with human subjects (§B.2), and additional details of the simulation setup (§C). Please refer to our video for qualitative results at tiny.cc/diffphy.

## A. Additional Results

Tab. 7 presents an ablation on window size performed using mocap data as initialization and reference trajectory rather than using the kinematic initialization. In this case, we note that a smaller window size of 480 outperforms the larger window size of 960 used in the main paper. We hypothesize that when the reference signal lacks noise, a smaller window is easier to optimize since the dimension of the problem is reduced. However, with noisy observations, a larger window is required for the method to be robust to missing or poor kinematic reconstructions.

| Window | MPJPE-G | MPJPE | MPJPE-PA |
|--------|---------|-------|----------|
| 240    | 112.8   | 75.9  | 40.1     |
| 480    | 39.4    | 33.4  | 21.9     |
| 720    | 46.1    | 42.1  | 29.4     |
| 960    | 77.8    | 68.4  | 44.9     |

Table 7. Ablation study of the optimization window size. Experiments were carried out on motion capture rather than the kinematic initialization as input. The experiment was performed on the same Human3.6M sequences as in the ablation in the main paper. Note that when using mocap rather than noisy observations, a smaller window size is better (480 vs. 960 in main paper).

## B. Datasets

We evaluate our method on the two established datasets Human3.6M [21] and AIST [49]. In addition, we evaluate our method on "real-world" internet videos.

**Human3.6M.** When comparing to the state-of-the-art methods, we evaluate on the Human3.6M Protocol P2 sequences while excluding the same sequences as by Xie et al. [56]. That leaves the sequences: *Directions, Discussions, Greeting, Posing, Purchases, Taking Photos, Waiting, Walking, Walking Dog and Walking Together*. We evaluate the motions using only camera *60457274*. Similar to [56], we down sample the Human3.6M data from 50 FPS to 25 FPS.

The ablation studies were performed on a smaller subset of four-second clips (frames 400-599) from a random camera, see tab. 8.

**AIST.** AIST provides dynamic dance motions not present in Human3.6M. We evaluate our method using the pseudo-ground-truth provided by [26]. We use the first four seconds

| Sequence | Subject | Camera Id | Frames |
|----------|---------|-----------|--------|
| Phoning | S11 | 55011271 | 400-599 |
| Posing_1 | S11 | 58860488 | 400-599 |
| Purchases | S11 | 60457274 | 400-599 |
| SittingDown_1 | S11 | 54138969 | 400-599 |
| Smoking_1 | S11 | 54138969 | 400-599 |
| TakingPhoto_1 | S11 | 54138969 | 400-599 |
| Waiting_1 | S11 | 58860488 | 400-599 |
| WalkDog | S11 | 58860488 | 400-599 |
| WalkTogether | S11 | 55011271 | 400-599 |
| Walking_1 | S11 | 55011271 | 400-599 |
| Greeting_1 | S9 | 54138969 | 400-599 |
| Phoning_1 | S9 | 54138969 | 400-599 |
| Purchases | S9 | 60457274 | 400-599 |
| SittingDown | S9 | 55011271 | 400-599 |
| Smoking | S9 | 60457274 | 400-599 |
| TakingPhoto | S9 | 60457274 | 400-599 |
| Waiting | S9 | 60457274 | 400-599 |
| WalkDog_1 | S9 | 54138969 | 400-599 |
| WalkTogether_1 | S9 | 55011271 | 400-599 |
| Walking | S9 | 58860488 | 400-599 |

Table 8. Human3.6M [21] sequences used for ablation studies. Note that we downsampled the sequences from 50 FPS to 25 FPS.

(120 frames) using a randomly selected camera from the sequences in tab. 9.

**Internet Videos.** Finally, we perform qualitative evaluation of our method on internet videos made public under creative common licences.

### B.1. Metrics

**Total variation.** We compute the total variation of the 3d joint acceleration as a measurement of the jitter in motion. This is given as

$$\frac{1}{T} \sum_{t \in T} \sum_{k \in K} |\ddot{x}_{t+1}^k - \ddot{x}_t^k|, \tag{10}$$

where $\ddot{x}_t^k$ is the 3d joint acceleration of joint $k$ at time $t$. We estimate the acceleration through finite differences.

**Foot skating.** We track unnatural foot skating artifacts by measuring the percentage of frames where either foot is "skating" along the ground. Our formulation doesn't rely on foot contact annotations but instead heuristic detect when foot contacts occur by measuring the distance between the foot mesh and the ground-plane. A contact is defined as $N = 10$ foot mesh vertices being within $d$ mm of the ground-plane. For kinematics we use $d = 5$ mm and for dynamics $d = 1$ mm to account for the capsule approximation being smaller than the foot mesh. We define skating as a foot moving $\geq 2$ cm between two frames while being in contact with the ground.

## B.2. Usage of data with human subjects

In this work, we employ two established pose benchmarks that are commonly used in the field of human pose estimation. Human3.6M [21] was recorded in a laboratory setting with the permission of the actors, and AIST [49] contains *"a shared database containing original street dance videos with copyright-cleared dance music. This is the first large-scale shared database focusing on street dances to promote academic research regarding Dance Information Processing"*[2]. As for the "in-the-wild" videos, these were released under creative common licenses granting express permission to *"copy and redistribute the material in any medium or format"* and *"remix, transform, and build upon the material for any purpose, even commercially"*. Finally, we do *not* intend to release these videos as part of a dataset. Instead we only use them to demonstrate our method on videos with poses and motion uncommon in laboratory captured datasets.

| Sequence | Frames |
|---|---|
| gBR_sBM_c06_d06_mBR4_ch06 | 1-120 |
| gBR_sBM_c07_d06_mBR4_ch02 | 1-120 |
| gBR_sBM_c08_d05_mBR1_ch01 | 1-120 |
| gBR_sFM_c03_d04_mBR0_ch01 | 1-120 |
| gJB_sBM_c02_d09_mJB3_ch10 | 1-120 |
| gKR_sBM_c09_d30_mKR5_ch05 | 1-120 |
| gLH_sBM_c04_d18_mLH5_ch07 | 1-120 |
| gLH_sBM_c07_d18_mLH4_ch03 | 1-120 |
| gLH_sBM_c09_d17_mLH1_ch02 | 1-120 |
| gLH_sFM_c03_d18_mLH0_ch15 | 1-120 |
| gLO_sBM_c05_d14_mLO4_ch07 | 1-120 |
| gLO_sBM_c07_d15_mLO4_ch09 | 1-120 |
| gLO_sFM_c02_d15_mLO4_ch21 | 1-120 |
| gMH_sBM_c01_d24_mMH3_ch02 | 1-120 |
| gMH_sBM_c05_d24_mMH4_ch07 | 1-120 |

Table 9. AIST [49] sequences used for evaluation.

## C. Differentiable Physics for Human Motion

Tiny Differentiable Simulator (TDS) [17] is a C++ simulator where the data type is templetized. In our experiments, we use the scalar from the automatic differentiation (AD) framework CppAD [3] to compute the simulation gradients. That is, we compute the gradients of the loss with respect to the input control variables at each time step:

$$\frac{\partial L}{\partial \hat{\mathbf{q}}_{1:T}} = \frac{\partial L}{\partial \mathbf{q}_{1:T}} \frac{\partial \mathbf{q}_{1:T}}{\partial \boldsymbol{\tau}_{1:T}} \frac{\partial \boldsymbol{\tau}_{1:T}}{\partial \hat{\mathbf{q}}_{1:T}}, \qquad (11)$$

where $L$ is objective function of the trajectory optimization, $\mathbf{q}_{1:T}$ are the simulated body's joint positions, and $\hat{\mathbf{q}}_{1:T}$

are the per-timestep control signal to the PD controllers in the body joints.

To speed up the optimization we implement our simulation as a fixed computational graph of the simulation rollout for a fixed number of steps and then repeatedly use it to compute the values of the gradients in (11). This greatly speeds up the optimization since the automatic differentiation framework doesn't need to setup the computational graph for each backward pass. To that end, we make the following adaptations to TDS to make it support a fixed graph.

**Differentiation and contact points.** Since at the time of graph construction it is not known in advance which contact points will be active for particular inputs we always include all contact points into the LCP formulation. This increases the graph size based on the number of contacts considered. The issue of large graph can be address by e.g. "checkpointing" the computation as described in [38].

**Dealing with exploding gradients.** As noted in [33], gradients from differentiable simulators may explode or vanishing when the window size is large. In this work, we experimentally found it possible to mitigate the issue by setting the LCP solver iterations to $K = 1$ without noticeable degradation of reconstruction quality.

**Implementation Details** In our experiments we run TDS with a step size of 1ms. This is partly due to the simpler PD controller, which requires smaller simulation steps to allow for stable control. We set the ground-plane friction to $0.8$ and the controller gains to $k_p = 200$ and $k_d = 5$. Evaluating our loss function and computing the gradients for a window of 960 simulation steps takes approximately $\approx 5$ seconds on a standard desktop computer with only feet contacts enabled. Enabling more contacts or simulating multiple objects increases memory and computation time.

---

[2]https://aistdancedb.ongaaccel.jp/