Contact and Human Dynamics from Monocular Video

Davis Rempe^{1,2}, Leonidas J. Guibas¹, Aaron Hertzmann², Bryan Russell², Ruben Villegas², and Jimei Yang²

> ¹ Stanford University ² Adobe Research geometry.stanford.edu/projects/human-dynamics-eccv-2020

Abstract. Existing deep models predict 2D and 3D kinematic poses from video that are approximately accurate, but contain visible errors that violate physical constraints, such as feet penetrating the ground and bodies leaning at extreme angles. In this paper, we present a physicsbased method for inferring 3D human motion from video sequences that takes initial 2D and 3D pose estimates as input. We first estimate ground contact timings with a novel prediction network which is trained without hand-labeled data. A physics-based trajectory optimization then solves for a physically-plausible motion, based on the inputs. We show this process produces motions that are significantly more realistic than those from purely kinematic methods, substantially improving quantitative measures of both kinematic and dynamic plausibility. We demonstrate our method on character animation and pose estimation tasks on dynamic motions of dancing and sports with complex contact patterns.

1 Introduction

Recent methods for human pose estimation from monocular video [1,17,30,43] estimate accurate overall body pose with small absolute differences from the true poses in body-frame 3D coordinates. However, the recovered motions in world-frame are visually and physically implausible in many ways, including feet that float slightly or penetrate the ground, implausible forward or backward body lean, and motion errors like jittery, vibrating poses. These errors would prevent many subsequent uses of the motions. For example, inference of actions, intentions, and emotion often depends on subtleties of pose, contact and acceleration, as does computer animation; human perception is highly sensitive to physical inaccuracies [14,34]. Adding more training data would not solve these problems, because existing methods do not account for physical plausibility.

Physics-based trajectory optimization presents an appealing solution to these issues, particularly for dynamic motions like walking or dancing. Physics imposes important constraints that are hard to express in pose space but easy in terms of dynamics. For example, feet in static contact do not move, the body moves smoothly overall relative to contacts, and joint torques are not large. However,



Fig. 1. Our contact prediction and physics-based optimization corrects numerous physically implausible artifacts common in 3D human motion estimations from, e.g., Monocular Total Capture (MTC) [43] such as foot floating (top row), foot penetrations (middle), and unnatural leaning (bottom).

full-body dynamics is notoriously difficult to optimize [36], in part because contact is discontinuous, and the number of possible contact events grows exponentially in time. As a result, combined optimization of contact and dynamics is enormously sensitive to local minima.

This paper introduces a new strategy for extracting dynamically valid fullbody motions from monocular video (Figure 1), combining learned pose estimation with physical reasoning through trajectory optimization. As input, we use the results of kinematic pose estimation techniques [4,43], which produce accurate overall poses but inaccurate contacts and dynamics. Our method leverages a reduced-dimensional body model with centroidal dynamics and contact constraints [7,42] to produce a physically-valid motion that closely matches these inputs. We first infer foot contacts from 2D poses in the input video which are then used in a physics-based trajectory optimization to estimate 6D center-of-mass motion, feet positions, and contact forces. We show that a contact prediction network can be accurately trained on synthetic data. This allows us to separate initial contact estimation from motion optimization, making the optimization more tractable. As a result, our method is able to handle highly dynamic motions without sacrificing physical accuracy. We focus on single-person dynamic motions from dance, walking, and sports. Our approach substantially improves the realism of inferred motions over stateof-the-art methods, and estimates numerous physical properties that could be useful for further inference of scene properties and action recognition. We primarily demonstrate our method on character animation by retargeting captured motion from video to a virtual character. We evaluate our approach using numerous kinematics and dynamics metrics designed to measure the physical plausiblity of the estimated motion. The proposed method takes an important step to incorporating physical constraints into human motion estimation from video, and shows the potential to reconstruct realistic, dynamic sequences.

2 Related Work

We build on several threads of work in computer vision, computer animation, and robotics, each with a long history [9]. Recent vision results are detailed here.

Recent progress in pose estimation can accurately detect 2D human keypoints [4,12,27] and infer 3D pose [1,17,30] from a single image. Several recent methods extract 3D human motions from monocular videos by exploring various forms of temporal cues [18,26,44,43]. While these methods focus on explaining human motion in pixel space, they do not account for physical plausibility. Several recent works interpret interactions between people and their environment in order to make inferences about each [6,11,45]; each of these works uses only static kinematic constraints. Zou et al. [46] infer contact constraints to optimize 3D motion from video. We show how dynamics can improve inference of human-scene interactions, leading to more physically plausible motion capture.

Some works have proposed physics constraints to address the issues of kinematic tracking. Brubaker et al. [3] propose a physics-based tracker based on a reduced-dimensional walking model. Wei and Chai [41] track body motion from video, assuming keyframe and contact constraints are provided. Similar to our own work, Brubaker and Fleet [2] perform trajectory optimization for full-body motion. To jointly optimize contact and dynamics, they use a continuous approximation to contact. However, soft contact models introduce new difficulties, including inaccurate transitions and sensitivity to stiffness parameters, while still suffering from local minima issues. Moreover, their reduced-dimensional model includes only center-of-mass positional motion, which does not handle rotational motion well. In contrast, we obtain accurate contact initialization in a preprocessing step to simplify optimization, and we model rotational inertia.

Li et al. [23] estimate dynamic properties from videos. We share the same overall pipeline of estimating pose and contacts, followed by trajectory optimization. Whereas they focus on the dynamics of human-object interactions, we focus on videos where the human motion itself is much more dynamic, with complex variation in pose and foot contact; we do not consider human-object interaction. They use a simpler data term, and perform trajectory optimization in full-body dynamics unlike our reduced representation. Their classifier training requires hand-labeled data, unlike our automatic dataset creation method.



Fig. 2. Method overview. Given an input video, our method starts with initial estimates from existing 2D and 3D pose methods [4,43]. The lower-body 2D joints are used to infer foot contacts (orange box). Our optimization framework contains two parts (blue boxes). Inferred contacts and initial poses are used in a kinematic optimization that refines the 3D full-body motion and fits the ground. These are given to a reduceddimensional physics-based trajectory optimization that applies dynamics.

Prior methods learn character animation controllers from video. Vondrak et al. [38] train a state-machine controller using image silhouette features. Peng et al. [32] train a controller to perform skills by following kinematically-estimated poses from input video sequences. They demonstrate impressive results on a variety of skills. They do not attempt accurate reconstruction of motion or contact, nor do they evaluate for these tasks, rather they focus on control learning.

Our optimization is related to physics-based methods in computer animation, e.g., [8,16,20,24,25,33,40]. Two unique features of our optimization are the use of low-dimensional dynamics optimization that includes 6D center-of-mass motion and contact constraints, thereby capturing important rotational and footstep quantities without requiring full-body optimization, and the use of a classifier to determine contacts before optimization.

3 Physics-Based Motion Estimation

This section describes our approach, which is summarized in Figure 2. The core of our method is a physics-based trajectory optimization that enforces dynamics on the input motion (Section 3.1). Foot contact timings are estimated in a preprocess (Section 3.2), along with other inputs to the optimization (Section 3.3). Similar to previous work [23,43], in order to recover full-body motion we assume there is no camera motion and that the full body is visible.

3.1 Physics-Based Trajectory Optimization

The core of our framework is an optimization which enforces dynamics on an initial motion estimate given as input (see Section 3.3). The goal is to improve the plausibility of the motion by applying physical reasoning through the objective and constraints. We aim to avoid common perceptual errors, e.g., jittery, unnatural motion with feet skating and ground penetration, by generating a smooth trajectory with physically-valid momentum and static feet during contact.

The optimization is performed on a reduced-dimensional body model that captures overall motion, rotation, and contacts, but avoids the difficulty of optimizing all joints. Modeling rotation is necessary for important effects like arm swing and counter-oscillations [13,20,25], and the reduced-dimensional *centroidal* dynamics model can produce plausible trajectories for humanoid robots [5,7,28]. Our method is based on a recent robot motion planning algorithm from Winkler et al. [42] that leverages a simplified version of centroidal dynamics, which treats the robot as a rigid body with a fixed mass and moment of inertia. Their method finds a feasible trajectory by optimizing the position and rotation of the centerof-mass (COM) along with feet positions, contact forces, and contact durations as described in detail below. We modify this algorithm to suit our computer vision task: we use a temporally varying inertia tensor which allows for changes in mass distribution (swinging arms) and enables estimating the dynamic motions of interest, we add energy terms to match the input kinematic motion and foot contacts, and we add new kinematics constraints for our humanoid skeleton.

Inputs. The method takes initial estimates of: COM position $\bar{\mathbf{r}}(t) \in \mathbb{R}^3$ and orientation $\bar{\theta}(t) \in \mathbb{R}^3$ trajectories, body-frame inertia tensor trajectory $\mathbf{I}_b(t) \in \mathbb{R}^{3\times 3}$, and trajectories of the foot joint positions $\bar{\mathbf{p}}_{1:4}(t) \in \mathbb{R}^3$. There are four foot joints: left toe base, left heel, right toe base, and right heel, indexed as $i \in \{1, 2, 3, 4\}$. These inputs are at discrete timesteps, but we write them here as functions for clarity. The 3D ground plane height h_{floor} and upward normal is provided. Additionally, for each foot joint at each time, a binary label is provided indicating whether the foot is in contact with the ground. These labels determine initial estimates of contact durations for each foot joint $\bar{T}_{i,1}, \bar{T}_{i,2}, \ldots, \bar{T}_{i,n_i}$ as described below. The distance from toe to heel ℓ_{foot} and maximum distance from toe to hip ℓ_{leg} are also provided. All quantities are computed from video input as described in Sections 3.2 and 3.3, and are used to both initialize the optimization variables and as targets in the objective function.

Optimization Variables. The optimization variables are the COM position and Euler angle orientation $\mathbf{r}(t), \boldsymbol{\theta}(t) \in \mathbb{R}^3$, foot joint positions $\mathbf{p}_i(t) \in \mathbb{R}^3$ and contact forces $\mathbf{f}_i(t) \in \mathbb{R}^3$. These variables are continuous functions of time, represented by piece-wise cubic polynomials with continuity constraints. We also optimize contact timings. The contacts for each foot joint are independently parameterized by a sequence of phases that alternate between contact and flight. The optimizer cannot change the type of each phase (contact or flight), but it can modify their durations $T_{i,1}, T_{i,2}, \ldots, T_{i,n_i} \in \mathbb{R}$ where n_i is the number of total contact phases for the *i*th foot joint.

Objective. Our complete formulation is shown in Figure 3. E_{data} and E_{dur} seek to keep the motion and contacts as close as possible to the initial inputs, which

min	$\sum_{t=0}^{T} \left(E_{data}(t) + E_{vel}(t) + E_{acc}(t) \right) + E_{dur}$	
s.t.	$m\ddot{\mathbf{r}}(t) = \sum\nolimits_{i=1}^{4} \mathbf{f}_{i}(t) + m\mathbf{g}$	(dynamics)
	$\mathbf{I}_{w}(t)\dot{\boldsymbol{\omega}}(t) + \boldsymbol{\omega}(t) \times \mathbf{I}_{w}(t)\boldsymbol{\omega}(t) = \sum\nolimits_{i=1}^{4} \mathbf{f}_{i}(t) \times$	$(\mathbf{r}(t) - \mathbf{p}_i(t))$
	$\dot{\mathbf{r}}(0) = \dot{\overline{\mathbf{r}}}(0), \dot{\mathbf{r}}(T) = \dot{\overline{\mathbf{r}}}(T)$	(velocity boundaries)
	$ \mathbf{p}_1(t) - \mathbf{p}_2(t) = \mathbf{p}_3(t) - \mathbf{p}_4(t) = \ell_{foot}$	(foot kinematics)
	for every foot joint i :	
	$ \mathbf{p}_i(t) - \mathbf{p}_{hip,i}(t) \le \ell_{leg}$	(leg kinematics)
	$\sum_{j=1}^{n_i} T_{i,j} = T$	(contact durations)
	for foot joint i in contact at time t :	
	$\dot{\mathbf{p}}_i(t) = 0$	(no slip)
	$p_i^z(t) = h_{floor}(\mathbf{p}_i^{xy})$	(on floor)
	$0 \le \mathbf{f}_i(t)^T \hat{\mathbf{n}} \le f_{max}$	(pushing/max force)
	$ \mathbf{f}_i(t)^T \hat{\mathbf{t}}_{1,2} < \mu \mathbf{f}_i(t)^T \hat{\mathbf{n}}$	(friction pyramid)
	for foot joint i in flight at time t :	
	$p_i^z(t) \ge h_{floor}(\mathbf{p}_i^{xy})$	(above floor)
	$\mathbf{f}_i(t) = 0$	(no force in air)

Fig. 3. Physics-based trajectory optimization formulation. Please see text for details.

are derived from video, at discrete steps over the entire duration T:

$$E_{data}(t) = w_r ||\mathbf{r}(t) - \overline{\mathbf{r}}(t)||^2 + w_\theta ||\boldsymbol{\theta}(t) - \overline{\boldsymbol{\theta}}(t)||^2 + w_p \sum_{i=1}^4 ||\mathbf{p}_i(t) - \overline{\mathbf{p}}_i(t)||^2$$
(1)

$$E_{dur} = w_d \sum_{i=1}^{4} \sum_{j=1}^{n_i} (T_{i,j} - \bar{T}_{i,j})^2$$
(2)

We weigh these terms with $w_d = 0.1$, $w_r = 0.4$, $w_{\theta} = 1.7$, $w_p = 0.3$.

The remaining objective terms are regularizers that prefer small velocities and accelerations resulting in a smoother optimal trajectory:

$$E_{vel}(t) = \gamma_r ||\dot{\mathbf{r}}(t)||^2 + \gamma_\theta ||\dot{\boldsymbol{\theta}}(t)||^2 + \gamma_p \sum_{i=1}^4 ||\dot{\mathbf{p}}_i(t)||^2$$
(3)

$$E_{acc}(t) = \beta_r ||\ddot{\mathbf{r}}(t)||^2 + \beta_\theta ||\ddot{\boldsymbol{\theta}}(t)||^2 + \beta_p \sum_{i=1}^4 ||\ddot{\mathbf{p}}_i(t)||^2$$
(4)

with $\gamma_r = \gamma_{\theta} = 10^{-3}$, $\gamma_p = 0.1$ and $\beta_r = \beta_{\theta} = \beta_p = 10^{-4}$.

Constraints. The first set of constraints strictly enforce valid rigid body mechanics, including linear and angular momentum. This enforces important properties

of motion, for example, during flight the COM must follow a parabolic arc according to Newton's Second Law. During contact, the body motion acceleration is limited by the possible contact forces e.g., one cannot walk at a 45° lean.

At each timestep, we use the world-frame inertia tensor $\mathbf{I}_{w}(t)$ computed from the input $\mathbf{I}_{b}(t)$ and the current orientation $\boldsymbol{\theta}(t)$. This assumes that the final output poses will not be dramatically different from those of the input: a reasonable assumption since our optimization does not operate on upper-body joints and changes in feet positioning are typically small (though perceptually important). We found that using a constant inertia tensor (as in Winkler et al. [42]) made convergence difficult to achieve. The gravity vector is $\mathbf{g} = -9.8\hat{\mathbf{n}}$, where $\hat{\mathbf{n}}$ is the ground normal. The angular velocity $\boldsymbol{\omega}$ is a function of the rotations $\boldsymbol{\theta}$ [42].

The contact forces are constrained to ensure that they push away from the floor but are not greater than $f_{max} = 1000$ N in the normal direction. With 4 feet joints, this allows 4000 N of normal contact force: about the magnitude that a 100 kg (220 lb) person would produce for extremely dynamic dancing motion [19]. We assume no feet slipping during contact, so forces must also remain in a friction pyramid defined by friction coefficient $\mu = 0.5$ and floor plane tangents $\hat{\mathbf{t}}_1, \hat{\mathbf{t}}_2$. Lastly, forces should be zero at any foot joint not in contact.

Foot contact is enforced through constraints. When a foot joint is in contact, it should be stationary (no-slip) and at floor height h_{floor} . When not in contact, feet should always be on or above the ground. This avoids feet skating and penetration with the ground.

In order to make the optimized motion valid for a humanoid skeleton, the toe and heel of each foot should maintain a constant distance of ℓ_{foot} . Finally, no foot joint should be farther from its corresponding hip than the length of the leg ℓ_{leg} . The hip position $\mathbf{p}_{hip,i}(t)$ is computed from the COM orientation at that time based on the hip offset in the skeleton detailed in Section 3.3.

Optimization Algorithm. We optimize with IPOPT [39], a nonlinear interior point optimizer, using analytical derivatives. We perform the optimization in stages: we first use fixed contact phases and no dynamics constraints to fit the polynomial representation for COM and feet position variables as close as possible to the input motion. Next, we add in dynamics constraints to find a physically valid motion, and finally we allow contact phase durations to be optimized to further refine the motion if possible.

Following the optimization, we compute a full-body motion from the physicallyvalid COM and foot joint positions using Inverse Kinematics (IK) on a desired skeleton \mathbf{S}_{tqt} (see supplement for details).

3.2 Learning to Estimate Contacts

Before performing our physics-based optimization, we need to infer when the subject's feet are in contact with the ground, given an input video. These contacts are a target for the physics optimization objective and their accuracy is crucial to its success. To do so, we train a network that, for each video frame, classifies whether the toe and heel of each foot are in contact with the ground.

The main challenge is to construct a suitable dataset and feature representation. There is currently no publicly-available dataset of videos with labeled foot contacts and a wide variety of dynamic motions. Manually labeling a large, varied dataset would be difficult and costly. Instead, we generate synthetic data using motion capture (mocap) sequences. We automatically label contacts in the mocap and then use 2D joint position features from OpenPose [4] as input to our model, rather than image features from the raw rendered video frames. This allows us to train on synthetic data but then apply the model to real inputs.

Dataset. To construct our dataset, we obtained 65 mocap sequences for the 13 most human-like characters from www.mixamo.com, ranging from dynamic dancing motions to idling. Our set contains a diverse range of mocap sequences, retargeted to a variety of animated characters. At each time of each motion sequence, four possible contacts are automatically labeled by a heuristic: a toe or heel joint is considered to be in contact when (i) it has moved less than 2 cm from the previous time, and (ii) it is within 5 cm from the known ground plane. Although more sophisticated labeling [15,21] could be used, we found this approach sufficiently accurate to learn a model for the videos we evaluated on.

We render these motions (see Figure 5(c)) on their rigged characters with motion blur, randomized camera viewpoint, lighting, and floor texture. For each sequence, we render two views, resulting in over 100k frames of video with labeled contacts and 2D and 3D poses. Finally, we run a 2D pose estimation algorithm, OpenPose [4], to obtain the 2D skeleton which our model uses as input.

Model and Training. The classification problem is to map from 2D pose in each frame to the four contact labels for the feet joints. As we demonstrate in Section 4.1, simple heuristics based on 2D velocity do not accurately label contacts due to the ambiguities of 3D projection and noise.

For a given time t, our labeling neural network takes as input the 2D poses over a temporal window of duration w centered on the target frame at t. The 2D joint positions over the window are normalized to place the root position of the target frame at (0,0), resulting in relative position and velocity. We set w = 9video frames and use the 13 lower-body joint positions as shown in Figure 4. Additionally, the OpenPose confidence c for each joint position is included as input. Hence, the input to the network is a vector of (x, y, c) values of dimension 3 * 13 * 9 = 351. The model outputs four contact labels (left/right toe, left/right heel) for a window of 5 frames centered around the target. At test time, we use majority voting at overlapping predictions to smooth labels across time.

We use a five-layer multilayer perceptron (MLP) (sizes 1024, 512, 128, 32, 20) with ReLU non-linearities [29]. We train the network entirely on our synthetic dataset split 80/10/10 for train/validation/test based on motions per character, i.e., no motion will be in both train and test on the same character, but a training motion may appear in the test set retargeted to a different character. Although 3D motions may be similar in train and test, the resulting 2D motions (the network input) will be very different after projecting to differing camera viewpoints. The network is trained using a standard binary cross-entropy loss.

3.3 Kinematic Initialization

Along with contact labels, our physics-based optimization requires as input a ground plane and initial trajectories for the COM, feet, and inertia tensor. In order to obtain these, we compute an initial 3D full-body motion from video. Since this stage uses standard elements, e.g., [10], we summarize the algorithm here, and provide full details in the supplement.

First, Monocular Total Capture [43] (MTC) is applied to the input video to obtain an initial noisy 3D pose estimate for each frame. Although MTC accounts for motion through a texture-based refinement step, the output still contains a number of artifacts (Figure 1) that make it unsuitable for direct use in our physics optimization. Instead, we initialize a skeleton \mathbf{S}_{src} containing 28 body joints from the MTC input poses, and then use a kinematic optimization to solve for an optimal root translation and joint angles over time, along with parameters of the ground plane. The objective for this optimization contains terms to smooth the motion, ensure feet are stationary and on the ground when in contact, and to stay close to both the 2D OpenPose and 3D MTC pose inputs.

We first optimize so that the feet are stationary, but not at a consistent height. Next, we use a robust regression to find the ground plane which best fits the foot joint contact positions. Finally, we continue the optimization to ensure all feet are on this ground plane when in contact.

The full-body output motion of the kinematic optimization is used to extract inputs for the physics optimization. Using a predefined body mass (73 kg for all experiments) and distribution [22], we compute the COM and inertia tensor trajectories. We use the orientation about the root joint as the COM orientation, and the feet joint positions are used directly.

4 Results

Here we present extensive qualitative and quantitative evaluations of our contact estimation and motion optimization.

4.1 Contact Estimation

We evaluate our learned contact estimation method and compare to baselines on the synthetic test set (78 videos) and 9 real videos with manually-labeled foot contacts. The real videos contain dynamic dancing motions and include 700 labeled frames in total. In Table 1, we report classification accuracy for our method and numerous baselines.

We compare to using a velocity heuristic on foot joints, as described in Section 3.2, for both the 2D OpenPose and 3D MTC estimations. We also compare to using different subsets of joint positions. Our MLP using all lower-body joints is substantially more accurate on both synthetic and real videos than all baselines. Using upper-body joints down to the knees yields surprisingly good results.

In order to test the benefit of contact estimation, we compared our full optimization pipeline on the synthetic test set using network-predicted contacts

 Table 1. Classification accuracy of estimating foot contacts from video. Left: comparison to various baselines, Right: ablations using subsets of joints as input features.

Baseline Method	Synthetic Accuracy	Real Accuracy	MLP Input Joints	Synthetic Accuracy	Real Accuracy
Random	0.507	0.480	Upper down to hips	0.919	0.692
Always Contact	0.677	0.647	Upper down to knees	0.935	0.865
2D Velocity	0.853	0.867	Lower up to ankles	0.933	0.923
3D Velocity	0.818	0.875	Lower up to hips	0.941	0.935



Fig. 4. Foot contact estimation on a video using our learned model compared to a 2D velocity heuristic. All visualized joints are used as input to the network which outputs four contact labels (left toes, left heel, right toes, right heel). Red joints are labeled as contacting. Key differences are shown with orange boxes.

versus contacts predicted using a velocity heuristic on the 3D joints from MTC input. Optimization using network-predicted contacts converged for 94.9% of the test set videos, compared to 69.2% for the velocity heuristic. This illustrates how contact prediction is crucial to the success of motion optimization.

Qualitative results of our contact estimation method are shown in Figure 4. Our method is compared to the 2D velocity baseline which has difficulty for planted feet when detections are noisy, and often labels contacts for joints that are stationary but off the ground (e.g. heels).

4.2 Qualitative Motion Evaluation

Our method provides key qualitative improvements over prior kinematic approaches. We urge the reader to **view the supplementary video** in order to fully appreciate the generated motions. For qualitative evaluation, we demonstrate animation from video by retargeting captured motion to a computeranimated character. Given a target skeleton \mathbf{S}_{tgt} for a character, we insert an IK retargeting step following the kinematic optimization as shown in Figure 2 (see supplement for details), allowing us to perform the usual physics-based op-



Fig. 5. Qualitative results on synthetic and real data. a) results on a synthetic test video with a ground truth alternate view. Two nearby frames are shown for the input video and the alternate view. We fix penetration, floating and leaning prevalent in our method's input from MTC. b) dynamic exercise video (top) and the output full-body motion (middle) and optimized COM trajectory and contact forces (bottom).

timization on this new skeleton. We use the same IK procedure to compare to MTC results directly targeted to the character.

Figure 1 shows that our proposed method fixes artifacts such as foot floating (top row), foot penetrations (middle), and unnatural leaning (bottom). Figure 5(a) shows frames comparing the MTC input to our final result on a synthetic video for which we have a ground truth alternate view. For this example only, we use the true ground plane as input to our method for a fair comparison (see Section 4.3). From the input view, our method fixes feet floating and penetration. From the first frame of the alternate view, we see that the MTC pose is in fact extremely unstable, leaning backward while balancing on its heels; our method has placed the contacting feet in a stable position to support the pose, better matching the true motion.

Figure 5(b) shows additional qualitative results on a real video. We faithfully reconstruct dynamic motion with complex contact patterns in a physically accurate way. The bottom row shows the outputs of the physics-based optimization stage of our method at multiple frames: the COM trajectory and contact forces at the heel and toe of each foot.

Quantitative Motion Evaluation 4.3

Quantitative evaluation of high-quality motion estimation presents a significant challenge. Recent pose estimation work evaluates average positional errors of joints in the local body frame up to various global alignment methods [31]. However, those pose errors can be misleading: a motion can be pose-wise close to ground truth on average, but produce extremely implausible dynamics, including vibrating positions and extreme body lean. These errors can be perceptually objectionable when remapping the motion onto an animated character, and prevent the use of inferred dynamics for downstream vision tasks.

Therefore, we propose to use a set of metrics inspired by the biomechanics literature [2,13,16], namely, to evaluate *plausibility* of physical quantities based on known properties of human motion.

We use two baselines: MTC, which is the state-of-the-art for pose estimation, and our kinematic-only initialization (Section 3.3), which transforms the MTC input to align with the estimated contacts from Section 3.2. We run each method on the synthetic test set of 78 videos. For these quantitative evaluations only, we use



Fig. 6. Contact forces from our physicsbased optimization for a walking and dancing motion. The net contact forces around 1000 N are 140% of the assumed body weight (73 kg), a reasonable estimate compared to prior force plate data [2].

the ground truth floor plane as input to our method to ensure a fair comparison. Note that our method does not *need* the ground truth floor, but using it ensures a proper evaluation of our primary contributions rather than that of the floor fitting procedure, which is highly dependent on the quality of MTC input (see supplement for quantitative results using the estimated floor).

Dynamics Metrics. To evaluate dynamic plausibility, we estimate net ground reaction forces (GRF), defined as $\mathbf{f}_{GRF}(t) = \sum_i \mathbf{f}_i(t)$. For our full pipeline, we use the physics-based optimized GRFs which we compare to implied forces from the kinematic-only initialization and MTC input. In order to infer the GRFs implied by the kinematic optimization and MTC, we estimate the COM trajectory of the motion using the same mass and distribution as for our physics-based optimization (73 kg). We then approximate the acceleration at each time step and solve for the implied GRFs for all time steps (both in contact and flight).

We assess plausibility using GRFs measured in force plate studies, e.g., [2,13,35]. For walking, GRFs typically reach 80% of body weight; for a dance jump, GRFs can reach up to about 400% of body weight [19]. Since we do not know body weights of our subjects, we use a conservative range of 50kg–80kg for evaluation. Figure 6 shows the optimized GRFs produced by our method for a walking and swing dancing motion. The peak GRFs produced by our method match the data: for the walking motion, 115–184% of body weight, and 127–204% for dancing. In contrast, the kinematic-only GRFs are 319–510% (walking) and 765–1223% (dancing); these are implausibly high, a consequence of noisy and unrealistic joint accelerations.

We also measure GRF plausibility across the whole test set (Table 2(left)). GRF values are measured as a percentage of the GRF exerted by an idle 73 kg person. On average, our estimate is within 1% of the idle force, while the kinematic motion implies GRFs as if the person were 24.4% heavier. Similarly, the peak force of the kinematic motion is equivalent to the subject carrying an extra 830 kg of weight, compared to only 174 kg after physics optimization. The Max GRF for MTC is even less plausible, as the COM motion is jittery

Table 2. Physical plausibility evaluation on synthetic test set. $Mean/Max \ GRF$ are contact forces as a proportion of body weight; see text for discussion of plausible values. *Ballistic GRF* are unexplained forces during flight; smaller values are better. Foot position metrics measure the percentage of frames containing typical foot contact errors per joint; smaller values are better.

Dynamics (Contact forces)				Kinematics (Foot positions)			
Method	Mean GRF	Max GRF	Ballistic GRF	Floating	Penetration	Skate	
MTC [43]	143.0%	9055.3%	115.6%	58.7%	21.1%	16.8%	
Kinematics (ours)	124.4%	1237.5%	255.2%	2.3%	2.8%	1.6%	
Physics (ours)	99.0%	$\mathbf{338.6\%}$	0.0%	8.2%	0.3%	3.6%	

before smoothing during kinematic and dynamics optimization. *Ballistic GRF* measures the median GRF on the COM when no feet joints should be in contact according to ground truth labels. The GRF should be exactly 0%, meaning there are no contact forces and only gravity acts on the COM; the kinematic method obtains results of 255%, as if the subject were wearing a powerful jet pack.

Kinematics Metrics. We consider three kinematic measures of plausibility (Table 2(right)). These metrics evaluate accuracy of foot contact measurements. Specifically, given ground truth labels of foot contact we compute instances of foot *Floating*, *Penetration*, and *Skate* for heel and toe joints. *Floating* is the fraction of foot joints more than 3 cm off the ground when they should be in contact. *Penetration* is the fraction penetrating the ground more than 3 cm at any time. *Skate* is the fraction moving more than 2 cm when in contact.

After our kinematics initialization, the scores on these metrics are best (lower is better for all metrics) and degrade slightly after adding physics. This is due to the IK step which produces full-body motion following the physics-based optimization. Both the kinematic and physics optimization results substantially outperform MTC, which is rarely at a consistent foot height.

Positional Metrics. For completeness, we evaluate the 3D pose output of our method on variations of standard positional metrics. Results are shown in Table 3. In addition to our synthetic test set, we evaluate on all walking sequences from the training split of HumanEva-I [37] using the known ground plane as input. We measure the mean **global** per-joint position error (mm) for ankle and toe joints (*Feet* in Table 3) and over all joints (*Body*). We also report the error after aligning the root joint of only the first frame of each sequence to the ground truth skeleton (*Body-Align 1*), essentially removing any spurious constant offset from the predicted trajectory. Note that this differs from the common practice of aligning the roots at every frame, since this would negate the effect of our trajectory optimization and thus does not provide an informative performance measure. The errors between all methods are comparable, showing at most a difference of 5 cm which is very small considering global joint position. Though

Table 3. Pose evaluation on synthetic and HumanEva-I walking datasets. We measure mean global per-joint 3D position error (no alignment) for feet and full-body joints. For full-body joints, we also report errors after root alignment on only the first frame of each sequence. We remain competitive while providing key physical improvements.

	Synthetic Data				HumanEva-I Walking			
Method	Feet	Body	Body-Align 1	Feet	Body	Body-Align 1		
MTC [43]	581.095	560.090	277.215	511.59	532.286	402.749		
Kinematics (ours)	573.097	562.356	281.044	496.671	525.332	407.869		
Physics (ours)	571.804	573.803	323.232	508.744	499.771	421.931		

the goal of our method is to improve physical plausibility, it does not negatively affect the pose on these standard measures.

5 Discussion

Contributions. The method described in this paper estimates physically-valid motions from initial kinematic pose estimates. As we show, this produces motions that are visually and physically much more plausible than the state-of-the-art methods. We show results on retargeting to characters, but it could also be used for further vision tasks that would benefit from dynamical properties of motion.

Estimating accurate human motion entails numerous challenges, and we have focused on one crucial sub-problem. There are several other important unknowns in this space, such as motion for partially-occluded individuals, and ground plane position. Each of these problems and the limitations discussed below are an enormous challenge in their own right and are therefore reserved for future work. However, we believe that the ideas in this work could contribute to solving these problems and open multiple avenues for future exploration.

Limitations. We make a number of assumptions to keep the problem manageable, all of which can be relaxed in future work: we assume that feet are unoccluded, there is a single ground plane, the subject is not interacting with other objects, and we do not handle contact from other body parts like knees or hands. These assumptions are permissible for the character animation from video mocap application, but should be considered in a general motion estimation approach. Our optimization is expensive. For a 2 second (60 frame) video clip, the physical optimization usually takes from 30 minutes to 1 hour. This runtime is due primarily to the adapted implementation from prior work [42] being ill-suited for the increased size and complexity of human motion optimization. We expect a specialized solver and optimized implementation to speed up execution.

Acknowledgments. This work was in part supported by NSF grant IIS-1763268, grants from the Samsung GRO program and the Stanford SAIL Toyota Research Center, and a gift from Adobe Corporation. We thank the following YouTube channels for permitting us to use their videos: Dance FreaX (Fig. 4), Dancercise Studio (Fig. 1&2), Fencer's Edge (Fig. 5), and MihranTV(Fig. 1).

References

- Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) Computer Vision – ECCV 2016. pp. 561–578. Springer International Publishing (2016)
- Brubaker, M.A., Sigal, L., Fleet, D.J.: Estimating contact dynamics. In: The IEEE International Conference on Computer Vision (ICCV). pp. 2389–2396. IEEE (2009)
- Brubaker, M.A., Fleet, D.J., Hertzmann, A.: Physics-based Person Tracking using the Anthropomorphic Walker. International Journal of Computer Vision 87(1), 140–155 (2010)
- 4. Cao, Z., Hidalgo, G., Simon, T., Wei, S.E., Sheikh, Y.: Openpose: Realtime multiperson 2d pose estimation using part affinity fields. pp. 1–1. IEEE (2019)
- 5. Carpentier, J., Mansard, N.: Multicontact locomotion of legged robots. IEEE Transactions on Robotics **34**(6), 1441–1460 (2018)
- Chen, Y., Huang, S., Yuan, T., Qi, S., Zhu, Y., Zhu, S.C.: Holistic++ scene understanding: Single-view 3d holistic scene parsing and human pose estimation with human-object interaction and physical commonsense. In: The IEEE International Conference on Computer Vision (ICCV). pp. 8648–8657. IEEE (2019)
- Dai, H., Valenzuela, A., Tedrake, R.: Whole-body motion planning with centroidal dynamics and full kinematics. In: IEEE-RAS International Conference on Humanoid Robots. pp. 295–302. IEEE (2014)
- Fang, A.C., Pollard, N.S.: Efficient synthesis of physically valid human motion. ACM Trans. Graph. 22(3), 417–426 (Jul 2003)
- Forsyth, D.A., Arikan, O., Ikemoto, L., O'Brien, J., Ramanan, D.: Computational studies of human motion: Part 1, tracking and motion synthesis. Foundations and Trends in Computer Graphics and Vision 1(2–3), 77–254 (2006)
- Gleicher, M.: Retargetting motion to new characters. In: Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques. pp. 33–42. SIGGRAPH '98, ACM, New York, NY, USA (1998)
- Hassan, M., Choutas, V., Tzionas, D., Black, M.J.: Resolving 3d human pose ambiguities with 3d scene constraints. In: The IEEE International Conference on Computer Vision (ICCV). pp. 2282–2292. IEEE (October 2019)
- 12. He, K., Gkioxari, G., Dollar, P., Girshick, R.: Mask r-cnn. In: The IEEE International Conference on Computer Vision (ICCV). pp. 2961–2969. IEEE (Oct 2017)
- Herr, H., Popovic, M.: Angular momentum in human walking. Journal of Experimental Biology 211(4), 467–481 (2008)
- Hoyet, L., McDonnell, R., O'Sullivan, C.: Push it real: Perceiving causality in virtual interactions. ACM Trans. Graph. 31(4), 90:1–90:9 (Jul 2012)
- Ikemoto, L., Arikan, O., Forsyth, D.: Knowing when to put your foot down. In: Proceedings of the 2006 Symposium on Interactive 3D Graphics and Games. pp. 49–53. I3D '06, ACM, New York, NY, USA (2006)
- Jiang, Y., Van Wouwe, T., De Groote, F., Liu, C.K.: Synthesis of biologically realistic human motion using joint torque actuation. ACM Trans. Graph. 38(4), 72:1–72:12 (Jul 2019)
- Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7122–7131. IEEE (June 2018)
- Kanazawa, A., Zhang, J.Y., Felsen, P., Malik, J.: Learning 3d human dynamics from video. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5614–5623. IEEE (June 2019)

- 16 D. Rempe et al.
- Kulig, K., Fietzer, A.L., Jr., J.M.P.: Ground reaction forces and knee mechanics in the weight acceptance phase of a dance leap take-off and landing. Journal of Sports Sciences 29(2), 125–131 (2011)
- de Lasa, M., Mordatch, I., Hertzmann, A.: Feature-based locomotion controllers. In: ACM SIGGRAPH 2010 Papers. pp. 131:1–131:10. SIGGRAPH '10, ACM, New York, NY, USA (2010)
- Le Callennec, B., Boulic, R.: Robust kinematic constraint detection for motion data. In: ACM SIGGRAPH/Eurographics Symposium on Computer Animation. pp. 281–290. SCA '06, Eurographics Association, Aire-la-Ville, Switzerland, Switzerland (2006)
- 22. de Leva, P.: Adjustments to zatsiorsky-seluyanov's segment inertia parameters. Journal of Biomechanics **29**(9), 1223 – 1230 (1996)
- Li, Z., Sedlar, J., Carpentier, J., Laptev, I., Mansard, N., Sivic, J.: Estimating 3d motion and forces of person-object interactions from monocular video. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 8640– 8649. IEEE (June 2019)
- Liu, C.K., Hertzmann, A., Popović, Z.: Learning physics-based motion style with nonlinear inverse optimization. ACM Trans. Graph. 24(3), 1071–1081 (Jul 2005)
- Macchietto, A., Zordan, V., Shelton, C.R.: Momentum control for balance. In: ACM SIGGRAPH 2009 Papers. pp. 80:1–80:8. SIGGRAPH '09, ACM, New York, NY, USA (2009)
- Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.P., Xu, W., Casas, D., Theobalt, C.: Vnect: Real-time 3d human pose estimation with a single rgb camera. ACM Trans. Graph. 36(4) (Jul 2017)
- Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) Computer Vision ECCV 2016. pp. 483–499. Springer International Publishing (2016)
- Orin, D.E., Goswami, A., Lee, S.H.: Centroidal dynamics of a humanoid robot. Autonomous Robots 35(2-3), 161–176 (2013)
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alche Buc, F., Fox, E., Garnett, R. (eds.) Advances in Neural Information Processing Systems 32. pp. 8026–8037. Curran Associates, Inc. (2019)
- Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3d hands, face, and body from a single image. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10975–10985. IEEE (June 2019)
- Pavllo, D., Feichtenhofer, C., Grangier, D., Auli, M.: 3d human pose estimation in video with temporal convolutions and semi-supervised training. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 7753–7762. IEEE (June 2019)
- Peng, X.B., Kanazawa, A., Malik, J., Abbeel, P., Levine, S.: Sfv: Reinforcement learning of physical skills from videos. ACM Trans. Graph. 37(6), 178:1–178:14 (Dec 2018)
- Popović, Z., Witkin, A.: Physically based motion transformation. In: Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques. pp. 11–20. SIGGRAPH '99, ACM Press/Addison-Wesley Publishing Co., New York, NY, USA (1999)

- Reitsma, P.S.A., Pollard, N.S.: Perceptual metrics for character animation: Sensitivity to errors in ballistic motion. ACM Trans. Graph. 22(3), 537–542 (Jul 2003)
- 35. Robertson, D.G.E., Caldwell, G.E., Hamill, J., Kamen, G., Whittlesey, S.N.: Research Methods in Biomechanics. Human Kinetics (2004)
- Safonova, A., Hodgins, J.K., Pollard, N.S.: Synthesizing physically realistic human motion in low-dimensional, behavior-specific spaces. In: ACM SIGGRAPH 2004 Papers. pp. 514–521. SIGGRAPH '04, ACM, New York, NY, USA (2004)
- 37. Sigal, L., Balan, A., Black, M.: Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. International Journal of Computer Vision 87, 4–27 (03 2010)
- Vondrak, M., Sigal, L., Hodgins, J., Jenkins, O.: Video-based 3d motion capture through biped control. ACM Trans. Graph. 31(4), 27:1–27:12 (Jul 2012)
- Wächter, A., Biegler, L.T.: On the implementation of an interior-point filter linesearch algorithm for large-scale nonlinear programming. Mathematical Programming 106(1), 25–57 (2006)
- Wang, J.M., Hamner, S.R., Delp, S.L., Koltun, V.: Optimizing locomotion controllers using biologically-based actuators and objectives. ACM Trans. Graph. 31(4) (Jul 2012)
- Wei, X., Chai, J.: Videomocap: Modeling physically realistic human motion from monocular video sequences. In: ACM SIGGRAPH 2010 Papers. pp. 42:1–42:10. SIGGRAPH '10, ACM, New York, NY, USA (2010)
- 42. Winkler, A.W., Bellicoso, D.C., Hutter, M., Buchli, J.: Gait and trajectory optimization for legged systems through phase-based end-effector parameterization. IEEE Robotics and Automation Letters (RA-L) 3, 1560–1567 (July 2018)
- 43. Xiang, D., Joo, H., Sheikh, Y.: Monocular total capture: Posing face, body, and hands in the wild. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10965–10974. IEEE (June 2019)
- 44. Xu, W., Chatterjee, A., Zollhöfer, M., Rhodin, H., Mehta, D., Seidel, H.P., Theobalt, C.: Monoperfcap: Human performance capture from monocular video. ACM Trans. Graph. 37(2), 27:1–27:15 (May 2018)
- Zanfir, A., Marinoiu, E., Sminchisescu, C.: Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2148–2157. IEEE (June 2018)
- Zou, Y., Yang, J., Ceylan, D., Zhang, J., Perazzi, F., Huang, J.B.: Reducing footskate in human motion reconstruction with ground contact constraints. In: The IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 459– 468. IEEE (March 2020)