

Modelling Human Kinetics and Kinematics during Walking using Reinforcement Learning

Visak Kumar

Abstract—In this work, we develop an automated method to generate 3D human walking motion in simulation which is comparable to real-world human motion. At the core, our work leverages the ability of deep reinforcement learning methods to learn high-dimensional motor skills while being robust to variations in the environment dynamics. Our approach iterates between policy learning and parameter identification to match the real-world bio-mechanical human data. We present a thorough evaluation of the kinematics, kinetics and ground reaction forces generated by our learned virtual human agent. We also show that the method generalizes well across human-subjects with different kinematic structure and gait-characteristics.

I. INTRODUCTION

Assistive devices such as exoskeletons have garnered a lot of research interest because they can bring numerous benefits to injury/disability rehabilitation, elderly care, human-augmentation for people doing heavy lifting at work, exercise and military applications. However, there are a few challenges we need to overcome before assistive devices can become more prevalent in our society. Exoskeletons work by directly applying forces on the joints, hence safety during operation becomes paramount when designing control policies for such devices. Often, testing directly on a human-subject is not recommended.

Simulation has played a key role in offloading the burden of testing control policies in the real-world with human subjects. However, simulating the real-world accurately has been a long-standing challenge in robotics. Modelling human motion is especially difficult because of the high-dimensional nature of the state space and the natural variability that exists between different human-subjects.

Our focus is on modelling human motion during locomotion tasks such as walking and running. Humanoid locomotion is governed by highly non-linear and discontinuous dynamics and large number of unknown quantities such as joint damping, segment Center of Mass (COM), segment inertia, friction coefficient, muscle activation and muscle-tendon interaction. These are extremely challenging to measure. Traditional methods to identify simulation models, commonly known as system identification, are often insufficient when the system in question is high-dimensional. System identification works by collecting data on the real-system and then inferring model parameters based on the data collected. For a high-dimensional system, collecting data in the task-relevant state space can be quite challenging.

The difficulty in modelling human locomotion has led some researchers to focus on reduced-order models, such as Linear-inverted pendulum models (LIPM) (or even 3D LIPM) [1], to explain human locomotion. Although computationally more efficient, these models come at the cost of accuracy and hence are limited in their application.

Reinforcement learning (RL) and imitation learning have taken center stage in recent years in developing control algorithms to imitate human motion. In particular, Peng et al [2] and Yu et al [3] learn complex motor skills like walking, running, jumping and even performing a backflip. In addition to learning complex motion, policies trained using RL are also robust to variations in the underlying dynamics: variations caused by different dynamical parameters [4] as well as shape and size of the human-subject [?].

In this work, we leverage the following two advantages of RL, (1) Ability to learn complex motions in high-dimensional space and (2) robustness to variability in underlying dynamics, to learn control policies for human walking that are bio-mechanically accurate. We validate our approach by comparing the joint angles, joint moments and ground reaction forces generated by the policy to real-world data collected in human-subject experiments.

II. RELATED WORK

Since this work aims to leverage tools and data which are outcomes of research in two different fields - Deep Reinforcement Learning (DRL) and Biomechanics. The literature review is organized into two sections, each describes the current state of the art in each field and what aspects need improvement.

A. Deep Reinforcement Learning

DRL has seen remarkable success in recent years for learning complex tasks ranging from video games to controlling robots [5]–[7]. In particular, DRL has been successful in the field of developing locomotion controllers [2], [4], [8], [9]. However, these methods are seldom validated by comparing it to real human movement data generated by experiments. In Peng et al [2], a policy optimized using DRL algorithm was able to learn exceptional skills like walking, running, jumping and even doing a backflip with the help of motion capture data. But, the resulting policy, while visually pleasing, is not validated with real-world human data. Similarly, in Yu et al [9], several locomotion skills like walking and running at different speeds was learned from scratch using a curriculum learning approach. However, as was the issue with [2], this approach is not validated with real data. Additionally,

the above mentioned approach only adopts symmetry and low energy as metrics to enforce in the learned walking strategy, but some research in biomechanics points towards an existence of more complex relationships that gives rise to walking. For example, Wang et al [10] identified that the Center-of-mass velocity has a direct linear relationship to step-length. In Kumar et al [11], preliminary comparison of the walking motion learned using policy optimization and real-world human subject experiments were made. However, only joint kinematics and foot-step lengths were used as a metric for validation. We need a more thorough evaluation using joint kinematics, kinetics, ground reaction forces and foot step lengths across different individuals to be more confident about the simulation model.

B. Biomechanics

To study human motion during walking and recovery, biomechanics researchers often adopt an experimental approach. First, data is collected in the real-world, then control policies are synthesized using the real-world data. Winters et al [12] was among the first in the field to study human gait, and the data published in this work remains relevant to this day. Wang et al [10] and Hof et al [13] performed perturbation experiments and identified important relationships between COM velocity, step-lengths, center of pressure, stepping vs ankle strategy, etc.. We aim to leverage the finding of this research to validate some of the results. In Joshi et al [1], a balance recovery controller was derived using the results reported in [10], however, they use a 3D Linear Inverted pendulum model to approximate the human dynamics. A 3D LIPD does not capture the dynamics fully, for example, angular momentum about the center of mass. Most relevant to our work, Antoine et al [14], used a direct-collocation trajectory optimization to synthesize a walking controller for a 3D musculo-skeletal model in OpenSim (OpenSim gait2392 model). The gait generated by the controller closely matched experimental data. The proposed method, relies on understanding the basic principles that lead to walking, such as minimizing metabolic cost, muscle activations, etc. However, the proposed solution enforces left-right symmetry, which works for walking, but is not ideal for disturbance recovery. Hence its unclear how well this approach will perform when there is an external disturbance to the human.

III. METHOD

We propose a framework to automate the process of developing bio-mechanically accurate 3D Human walking policies in simulation. The difference between the dynamics of a simulated human agent and a real human subject is caused by multiple factors. Some quantities such as mass, height, etc.. can be easily measured. However other quantities such as dynamical parameters: joint damping, ground friction, lower-limb joint axis location and the accurate segment lengths can be challenging to estimate. Note that while the segment lengths can be easily measured on a real human, the bio mechanical data sets usually do not provide this information.

We present an iterative approach to develop accurate walking model in which, first, we use Deep Reinforcement Learning (DRL) to learn a walking policy using a nominal dynamical model. During training, domain randomization is used to ensure the policy is robust to dynamical and kinematic parameters. Second, once we have a policy to generate walking motion, we perform an optimization step to identify the optimal parameters that explain real-world walking motion, these two steps are repeated until convergence. To validate our simulated models, we compare the gait characteristics such as joint kinematics, kinetics and ground reaction forces generated by our policy to real-world data collected with 5 human participants.

We formulate this problem of learning human walking as a Markov Decision Processes (MDPs), $(\mathcal{S}, \mathcal{A}, \mathcal{T}, r, p_0, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{T} is the transition function, r is the reward function, p_0 is the initial state distribution and γ is a discount factor. We take the approach of model-free reinforcement learning to find a policy π , such that it maximizes the accumulated reward:

$$J(\pi) = \mathbb{E}_{\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_T} \sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t),$$

where $\mathbf{s}_0 \sim p_0$, $\mathbf{a}_t \sim \pi(\mathbf{s}_t)$ and $\mathbf{s}_{t+1} = \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t)$.

We denote the human walking policy as $\pi_h(\mathbf{a}_h | \mathbf{s}_h)$ where \mathbf{s}_h , \mathbf{a}_h represent the states and actions, respectively.

IV. DATA AND SIMULATION ENVIRONMENT

Fukuchi et al [15] published open-source biomechanics data of 38 human subjects walking overground and on treadmill. We use data of 10 subjects (for proof-of-concept) with different physical characteristics such as mass, height, leg-length, gender and speed of walking. For each subject, this dataset provides motion capture marker data, and individual gait characteristics like joint angles, joint moments and ground reaction forces.

We use DART physics [16] engine as our simulation environment in which the virtual human agent learns to walk. We first create 10 virtual agents with the same physical characteristics as the human subjects shown in table I. We choose hunt-crossley contact model to generate contact forces between two rigid bodies such as feet and ground. Hunt-crossley model is widely used in biomechanics research and is also one of the contact models used in OpenSim simulator [17].

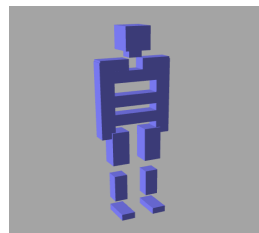


Fig. 1. We model a 31-Degree of Freedom(DoF) humanoid in PyDart.

Subject	Mass (kg)	Height (cm)	Age	Speed (m/s)
1 (M)	74	172.40	25	1.25
2 (M)	52.9	166.80	22	1.40
3 (F)	48.8	158	24	1.15
4 (M)	61.5	180.70	22	1.28
5 (F)	153	64	31	1.09
6 (M)	69.85	155	38	1.3
7 (F)	64.6	151.50	57	0.91
8 (M)	63.3	175	71	0.57
9 (F)	46.05	149.20	63	0.80
10 (M)	66.35	155.50	84	0.60

TABLE I
HUMAN SUBJECT PHYSICAL CHARACTERISTICS.

A. Human Walking Policy

We take a model-free reinforcement learning approach to develop a human locomotion policy $\pi_h(\mathbf{a}_h|\mathbf{s}_h)$. To achieve natural walking behaviors, we train a policy that imitates the human walking reference motion (motion capture data from the open-source data base) similar to Peng *et al.* [2]. The human 3D model (agent) consists of 25 actuated joints with a floating base as shown in Figure 1. This gives rise to a 71 dimensional state space $\mathbf{s}_h = [\mathbf{q}, \dot{\mathbf{q}}, \mathbf{v}_{com}, \boldsymbol{\omega}_{com}, \phi, \mu]$, including joint positions, joint velocities, linear and angular velocities of the center of mass (COM), and a phase variable(ϕ) that indicates the target frame in the motion clip. Vector μ (R^{14}) described as

$$\mu = [\beta, \sigma, f_l, s_l, t_l, f_j, h_j, s_j]$$

where β is joint damping, σ is friction coefficient between the foot and the ground. f_l, s_l, t_l and foot, shin and thigh segment lengths respectively and f_j, h_j, s_j are their corresponding joint axis location expressed in the robot root coordinate frame. These quantities are randomized during training so that the learned policy can be robust to variations in these quantities.

The action determines the target joint angles \mathbf{q}_t^{target} of the proportional-derivative (PD) controllers, deviating from the joint angles in the reference motion:

$$\mathbf{q}_t^{target} = \hat{\mathbf{q}}_t(\phi) + \mathbf{a}_t, \quad (1)$$

where $\hat{\mathbf{q}}_t(\phi)$ is the corresponding joint position in the reference motion at the given phase ϕ . Our reward function is designed to imitate the reference motion:

$$r_h(\mathbf{s}_h, \mathbf{a}_h) = w_q(\mathbf{q} - \hat{\mathbf{q}}(\phi)) + w_c(\mathbf{c} - \hat{\mathbf{c}}(\phi)) + w_e(\mathbf{e} - \hat{\mathbf{e}}(\phi)) - w_\tau \|\boldsymbol{\tau}\|^2 \quad (2)$$

where $\hat{\mathbf{q}}$, $\hat{\mathbf{c}}$, and $\hat{\mathbf{e}}$ are the desired joint positions, COM positions, and end-effector positions from the reference motion data, respectively. The reward function also penalizes the magnitude of torque $\boldsymbol{\tau}$. We use the same weight $w_q = 5.0$, $w_c = 2.0$, $w_e = 0.5$, and $w_\tau = 0.005$ for all experiments. We also use early termination of the rollouts, if the agent's pelvis drops below a certain height or if the base rotates about any axis beyond a threshold, we end the rollout and re-initialize the state.

We exert random forces to the agent during policy training. Each random force has a magnitude uniformly sampled

from $[0, 800]$ N and a direction uniformly sampled from $[-\pi/2, \pi/2]$, applied for 50 milliseconds on the agent's pelvis in parallel to the ground. The maximum force magnitude induces a velocity change of roughly 0.6m/sec. This magnitude of change in velocity is comparable to experiments found in literature such as [10], [18] and [13]. We also randomize the time when the force is applied within a gait cycle. Training in such a stochastic environment is crucial for reproducing the human motion. We represent a human policy as a multi-layered perceptron (MLP) neural network with two hidden layers of 128 neurons each. The formulated MDP is trained with Proximal Policy Optimization (PPO) [6].

B. Parameter Identification

As described in the previous section, the policy π_h is a function of vector μ . So, once the policy is trained we can search over the parameter μ to identify the values which best explains reference motion data. To do this, we perform an optimization step using CMA-ES algorithm [19]. We use the same reward function as described in equation 2. In this optimization process, at the beginning of each trajectory we set the value of μ , then generate the motion by executing the policy π_h , at the end of each generated trajectory the reward is computed using the equation 3. The goal is to maximize the sum of reward accumulated along the trajectory generated by the policy π_h .

$$\mu^* = \arg \max_{\mu} \sum_{t=0}^T (r_h^t(s_h^t, a_h^t)) \quad (3)$$

C. Error Threshold

We iterate between policy optimization and parameter identification until the error, defined in equation 4, is less than a predefined threshold κ .

$$\epsilon = \sum_{t=0}^T (||q - \hat{q}|| + ||\tau - \hat{\tau}|| + ||GRF - \hat{GRF}||) \quad (4)$$

Here, T is the time taken to complete one gait cycle, q is the joint angles generated by the policy and \hat{q} is the ground truth joint angles for that particular human subject. Similarly, τ is the joint moments and GRF is the ground reaction forces (vector of dimension R^3). $\hat{\tau}$ and \hat{GRF} are the ground truth joint moments and ground reaction forces measured in the real-world.

The procedure is outlined in algorithm 1.

V. RESULTS

We trained walking policies for 5 human subjects (first 5 in table I). We do the following evaluation to validate our approach,

- 1) Ablation study : How well does the joint angles, joint moments and ground reaction forces generated by our method compare to a baseline method which only uses RL? In other words, how much of a change does the second optimization step benefit the accuracy? We use Root-mean squared error as the metric for comparison.

Algorithm 1: Algorithm

```
1: for human subject  $i:1..n$  do
2:   Input: Dataset  $D_i$  - Motion capture trajectory of
      Human-subject  $i$ ,  $\kappa$  - error threshold
3:   Initialize CMA optimizer with generation size 8
4:   Initialize  $\pi_i(a|s, \mu)$ ,  $V_i(s)$ 
5:   while  $\epsilon > \kappa$  do
6:     Optimize policy  $\pi$  using reward function 2
7:     Optimize  $\mu^*$  using equation 3
8:     Compute error  $\epsilon$  using 4
9:   end while
10: end for
11: return  $\mu^*$  and  $\pi$ 
```

2) We highlight preliminary analysis on the generalizability of the approach by comparing the joint angles of 5 different human subjects.

A. Ablation Study

This ablation study is done only for one subject. In figure 2, we illustrate the benefit of our method, the additional optimization step has a clear effect on the RMSE error. Further, in figure 3, we also show the lower-limb joint angles before and after the CMA-ES optimization step. The ankle joint in particular matches the real-world data more closely. This highlights the ability of control policies trained using RL to be robust to variation in the dynamical parameters.



Fig. 2. RMSE error of our method compared to baseline approach which uses just RL to train the walking policy. We use 6 different initialization seeds for the policy, the bar plot indicates the RMSE error mean and standard deviation for the 6 policies for the same subject.

B. Joint moments, ground reaction forces and work loops

We also compare the joint moments and ground reaction forces generated by the policy to ground truth data by our method. The comparisons are illustrated in figures 5 and 6 respectively.

In addition to this, we also compare the torque loop at the hip joint and compare it to the bio-mechanical data reported in [12] (a long-standing gold standard). Torque loop is a plot of the torques generated at the hip joint in the y-axis and the joint angle in the x-axis during one-gait cycle. The

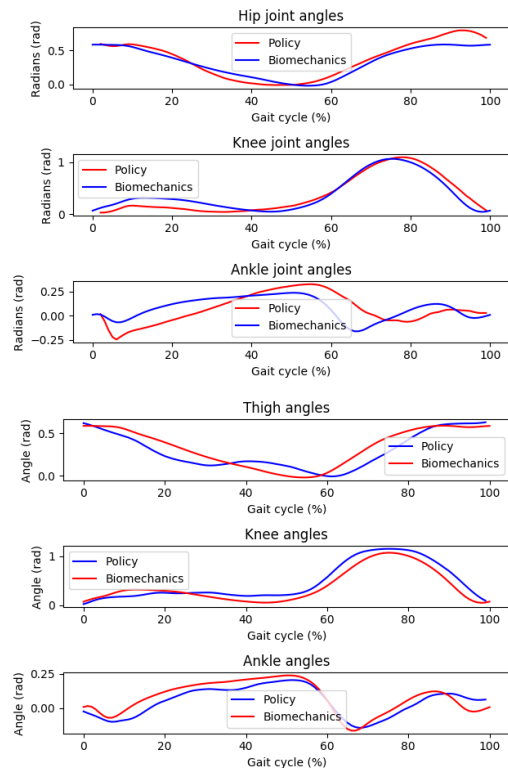


Fig. 3. Comparison of lower-limb joint angles during one gait cycle. **Top:** Joint angles after RL step, the ankle joints are noticeably different although the curve profile is similar. **Bottom:** After CMA-ES optimization step, the joint angles match better with real-world data.

arrows indicate the direction of movement and the points $[0, 50, 100]\%$ of the gait cycle match fairly well.

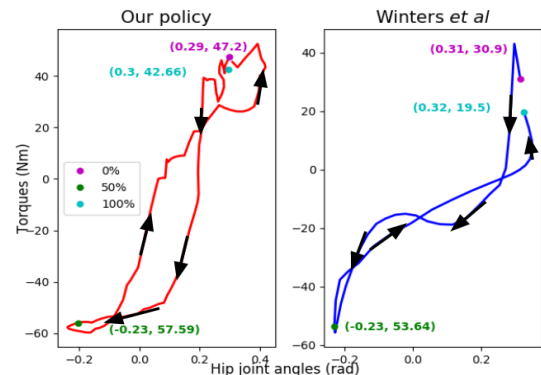


Fig. 4. Comparison of torque loops of a typical trajectory generated by our policy and human data reported by [12] at the hip of stance leg during a gait cycle. The green dots indicate the start and the black dots indicate 50% of the gait cycle. The arrows show the progression of the gait from 0% to 100%.

C. Results for five subjects

Our approach generalizes well to different human subjects. We apply our approach to 5 different human subjects walking at 5 different speeds. Figure 7 (**top**) illustrates the ground truth knee joint angle profiles for these 5 different human subjects. This shows the varied nature of the joint

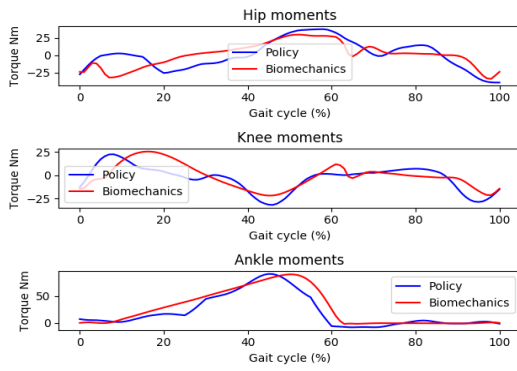


Fig. 5. Joint moments vs ground truth data.

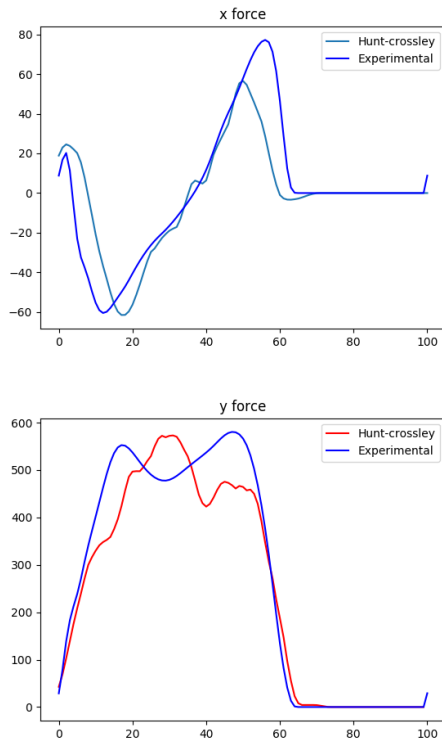


Fig. 6. Ground reaction forces **Top:** Tangential ground reaction force, its the force that propels a person forward while walking **Bottom:** Vertical ground reaction force, this is the normal force a person experiences while walking.

trajectories depending on speed and individual characteristics such as leg-length and mass. Our approach can capture these differences, as shown in figure 7 (**bottom**), the joint angles generated for the 5 different policies match well with the ground truth data.

VI. DISCUSSION AND CONCLUSION

In this work, we have shown that our proposed method can generate 3D human walking motion that matches well with real-world data. Our ablation study indicates there is an improvement over existing RL methods to learn walking motion. There are plenty of questions we can try to answer in future steps:

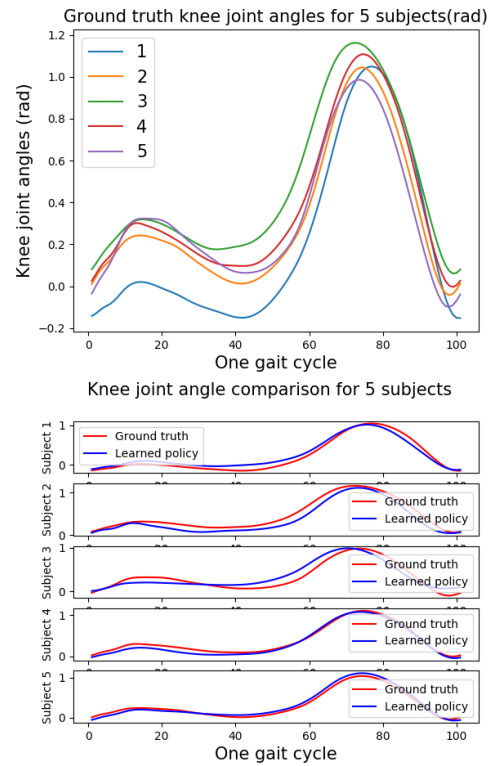


Fig. 7. Comparison of recovery performance when perturbation is applied at four different phases. **Top:** Comparison of stability region. **Bottom:** Comparison of COM velocity across five gait cycles. Perturbation is applied during the gait cycle 'p'. The increasing velocity after perturbation indicates that our policy is least effective at recovering when the perturbation occurs later in the swing phase.

- First, we would like to test our algorithm on all 42 subjects in the open source dataset [15].
- How well can we model a walking gait generated by a person with disability?
- How well can we model a walking gait generated by a person wearing an assistive device such as hip/ankle exoskeletons?
- Can we leverage this human model to learn interesting assistive strategies for an exoskeleton such as assistive walking to reduce metabolic cost of walking or to assist recovery when there is an external perturbation such as trips or slips?

REFERENCES

- [1] V. Joshi and M. Srinivasan, "A controller for walking derived from how humans recover from perturbations," 2019.
- [2] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne, "Deepmimic: Example-guided deep reinforcement learning of physics-based character skills," *ACM Transactions on Graphics (Proc. SIGGRAPH 2018)*, 2018.
- [3] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics*, vol. 37, no. 4, p. 1–12, Jul 2018. [Online]. Available: <http://dx.doi.org/10.1145/3197517.3201397>
- [4] W. Yu, V. C. V. Kumar, G. Turk, and C. K. Liu, "Sim-to-real transfer for biped locomotion," *CoRR*, vol. abs/1903.01390, 2019. [Online]. Available: <http://arxiv.org/abs/1903.01390>
- [5] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International Conference on Machine Learning*, 2015, pp. 1889–1897.

- [6] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [7] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," *CoRR*, vol. abs/1509.06461, 2015. [Online]. Available: <http://arxiv.org/abs/1509.06461>
- [8] X. B. Peng, G. Berseth, K. Yin, and M. Van De Panne, "Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, p. 41, 2017.
- [9] W. Yu, G. Turk, and C. K. Liu, "Learning symmetric and low-energy locomotion," *ACM Transactions on Graphics (Proc. SIGGRAPH 2018)*, vol. 37, no. 4, 2018.
- [10] Y. Wang and M. Srinivasan, "Stepping in the direction of the fall: The next foot placement can be predicted from current upper body state in steady-state walking," *Biology Letters*, vol. 10, no. 9, 2014.
- [11] V. C. V. Kumar, S. Ha, G. Sawicki, and C. K. Liu, "Learning a control policy for fall prevention on an assistive walking device," *International Conference of Robotics and Animation*, 2019.
- [12] D. A. Winter, *Biomechanics and motor control of human gait: normal, elderly and pathological*. Waterloo Biomechanics, 1991.
- [13] A. Hof, S. Vermerris, and W. Gjaltema, "Balance responses to lateral perturbations in human treadmill walking," *Journal of Experimental Biology*, vol. 213, no. 15, pp. 2655–2664, 2010.
- [14] F. Antoine, S. Gil, D. Christopher, G. Joris, J. Ilse, and F. Groote, "Rapid predictive simulations with complex musculoskeletal models suggest that diverse healthy and pathological human gaits can emerge from similar control strategies," 2019.
- [15] R. F. Claudine A Fukuchi and M. Duarte, "A public dataset of overground and treadmill walking kinematics and kinetics in healthy individuals," 2016. [Online]. Available: <https://doi.org/10.1080/00207540701450013>
- [16] S. Ha, "Pydart2," 2016. [Online]. Available: <https://github.com/sehoonha/pydart2>
- [17] OpenSim, "Hunt-Crossley model," 2011.
- [18] D. Martelli, V. Vashista, S. Micera, and S. K. Agrawal, "Direction-dependent adaptation of dynamic gait stability following waist-pull perturbations," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 24, no. 12, pp. 1304–1313, Dec 2016.
- [19] N. Hansen, "The CMA evolution strategy: A tutorial," *CoRR*, vol. abs/1604.00772, 2016. [Online]. Available: <http://arxiv.org/abs/1604.00772>