# Learning Velocity-based Humanoid Locomotion: Massively Parallel Learning with Brax and MJX*

William Thibault[1], William Melek[1], and Katja Mombaur[1,2]

[1]University of Waterloo, Waterloo, ON N2L 3G1, Canada
[2]Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany

June 30, 2024

### Abstract

Humanoid locomotion is a key skill to bring humanoids out of the lab and into the real-world. Many motion generation methods for locomotion have been proposed including reinforcement learning (RL). RL locomotion policies offer great versatility and generalizability along with the ability to experience new knowledge to improve over time. This work presents a velocity-based RL locomotion policy for the REEM-C robot. The policy uses a periodic reward formulation and is implemented in Brax/MJX for fast training. Simulation results for the policy are demonstrated with future experimental results in progress.

## 1 Introduction

Recent interest in humanoid robots as general purpose robots has lead to a significant increase in humanoid robotics research and development in both industry and academia. A main reason for the interest is because humanoid robots have the ability to perform repetitive and dull tasks in human environments. A core skill necessary for many tasks, like moving boxes around a warehouse, is robust locomotion. Locomotion planning and control algorithms vary greatly from linear inverted pendulum walking (LIPM) [4] to online whole-body MPC walking [3]. Reinforcement learning (RL) has also been a method of choice recently for robotic motion generation given its ability to adapt to different environments or conditions and generalize well to many scenarios. Additionally, RL is a data-driven approach allowing it to improve as it is exposed to new experience. New RL frameworks that optimize GPU usage have enabled faster train times allowing for robust policy learning with domain randomization in mere hours [2,7]. A method for generating a controlled, periodic bipedal gait on Cassie was developed by Siekmann et al. [8] through the use of a periodic reward composition where the rewards were provided using an indicator function for the robot matching the gait phase, otherwise it provides penalties. Singh et al. [9] proposed a method for learning planned 3D footsteps in simulation for larger, heavier robots like HRP-5P leveraging a similar periodic reward composition. This work was followed by a current controlled approach for using the trained policy on HRP-5P for planar walking with domain randomization to bridge the sim to real gap [10]. In this work, a humanoid locomotion RL policy for 2D velocity-based locomotion is proposed for the REEM-C robot. This work uses the periodic reward composition previously mentioned while leveraging faster train times than the previous work through the use of the recent Brax and MJX RL training framework [1]. Simulation results for the walking policy are presented in this work along with plans for experimental testing on the real REEM-C robot.

## 2 Reinforcement Learning Problem

### 2.1 REEM-C Humanoid Robot

The RL problem developed in this work is focused on joystick-style velocity locomotion for the full-size humanoid robot REEM-C. This humanoid robot has 30 degrees of freedom with 6 degree of freedom legs, 7 degree of freedom

---

*This work has been accepted at the CLAWAR 2024 conference in Kaiserslautern, Germany

Table 1: Summary of observation space

| State | Dimensions | Scaling |
|---|---|---|
| Yaw Rate | 1 | 0.25 |
| Projected Gravity of Base | 3 | - |
| Command | 3 | [2.0,2.0,0.25] |
| Joint Angles | 12 | - |
| Joint Velocities | 12 | - |
| Last Action | 12 | - |
| Periodic Clock Signal | 2 | - |

arms, 2 degrees of freedom for the torso and 2 degrees of freedom for the head. In the wrists and feet of the robot are 6D force-torque sensors. To simplify the learning problem, the upper body joints are held fixed and only the position controlled leg joints, with complete dynamics, are controllable in the MJX simulator.

## 2.2 Markov Decision Process

The Markov Decision Process for this work has a 12 dimensional action space composed of the 12 leg joint positions. The observation space includes the joint state of robot, the periodic clock signal as defined in Singh et al. [9] for a gait cycle with a double support time of 0.35 s and a single support time of 0.75 s, a command velocity of $(x, y, \psi)$. The input velocity has a range of $(x, y, \psi) = ([-0.3, 1.0]m/s, [-0.3, 0.3]m/s, [-0.5, 0.5]rad/s)$, based on reasonable limits and existing joystick LIPM walking for REEM-C. Table 1 summarizes the observation space. The rewards include the following rewards or penalties:

1. Center of mass linear velocity tracking reward:
   $r_1 = 3 * e^{-((\dot{x}_{input} - \dot{x}_{base})^2 + (\dot{y}_{input} - \dot{y}_{base})^2/\sigma)}$, where $\sigma$ is a scaling factor

2. Center of mass angular velocity tracking reward:
   $r_2 = 3 * e^{-((\dot{\psi}_{input} - \dot{\psi}_{base})^2/\sigma)}$, where $\sigma$ is a scaling factor

3. Linear z velocity penalty:
   $r_3 = -2 * (\dot{z}_{base})^2$

4. Torso roll and pitch rate penalty:
   $r_4 = -0.1 * (\dot{\phi})^2 + (\dot{\theta})^2$

5. Non-zero torso roll and pitch penalty:
   $r_5 = -10 * (\phi)^2 + (\theta)^2$

6. Torque penalization:
   $r_6 = -0.0002 * \sqrt{\Sigma \tau^2} + \Sigma|\tau|$

7. Action rate penalization:
   $r_7 = -0.01 * \Sigma(A_t - A_{t-1})^2$

8. Stand still no motion penalty:
   $r_8 = -0.5 * \Sigma|q - q_{default}|$, where this penalty is only included if the normalized velocity command is less than 0.1.

9. Early termination penalty:
   $r_9 = -1$, when the base height drops below 0.7 m, joint angles are exceeded or the robot is falling.

10. Gait foot contact stance reward or penalty:
    $r_{10} = 0.02 * (I_{stance}(F_{z,right}) + I_{stance}(F_{z,left}))$, where the indicator function $I_{stance}$ is defined by the gait phase as shown in Figure 1.

11. Gait foot flight velocity reward or penalty:
    $r_{11} = 0.2 * I_{flight}(\Sigma(v_{right})^2) + I_{flight}(\Sigma(v_{left})^2)$, where the indicator function $I_{flight}$ is defined by the gait phase as shown in Figure 1.

12. High foot contact penalty:
   $r_{12} = -0.01 * (F_{z,right} + F_{z,left})$, for contact forces exceeding the maximum contact force of 1500 N.

13. Gait phase reward:
   $r_{13} = 1$, when the gait phase changes to promote taking steps.

14. Velocity rate penalty:
   $r_{14} = -0.001 * \Sigma(\dot{q}_n - \dot{q}_{n-1})^2 + \Sigma(\dot{q}_n - 2.0 * \dot{q}_{n-1} + \dot{q}_{n-2})^2$

15. Pose bias penalty:
   $r_{15} = -0.4 * \Sigma(q - q_{default})^2$, where $q_{default}$ is the starting pose.

16. Head to base projection penalty:
   $r_{16} = -2 * ((x_{base} - x_{head})^2 + (y_{base} - y_{head})^2)$

17. Base height penalty:
   $r_{17} = -100 * (z_{base} - 0.8)^2$, where 0.8 m is the reference base height.

The reward functions and weights were largely inspired by Singh et al. [9], Radosavovic et al. [6] and the open-sourced Brax/MJX implementation of the work by Caluwaerts et al [2].
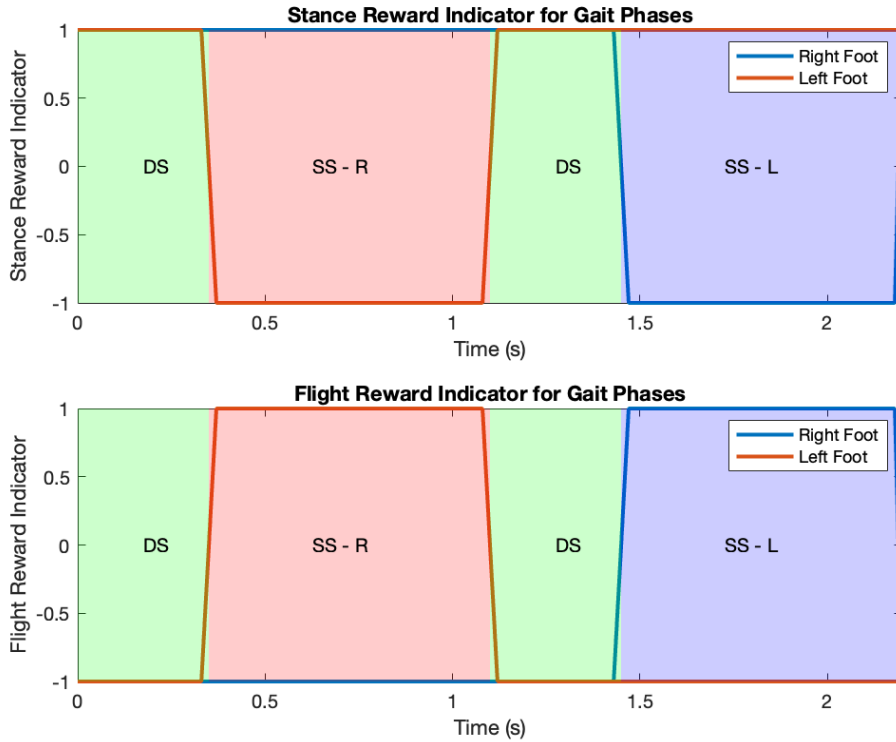


Figure 1: Stance and flight reward indicators depending on different gait phases

# 3 Results

The RL algorithms of Brax using the MJX simulator, a version of MuJoCo that supports the XLA compiler for GPU and TPU based simulation using JAX, is used to train the locomotion policy. This training pipeline was selected due to the highly parallel approach that can exploit the GPU for fast physics simulation steps and in turn shorter training times. The training is run for 200,000,000 training steps with 8192 parallel environments using PPO. The networks for the PPO algorithm had 4 hidden layers of 128 neurons. This training was performed on a PC with an AMD 7950x CPU, 64 GB RAM, and an RTX 4090 24 GB VRAM and completed training in

approximately 56 minutes. Note that to achieve fast and high performance training an optimized model of the robot was used, containing only lower body collisions and recommended MJX performance tuning techniques [5]. A sample motion sequence of the walking policy can be see in Figure 2. With these promising simulation results,
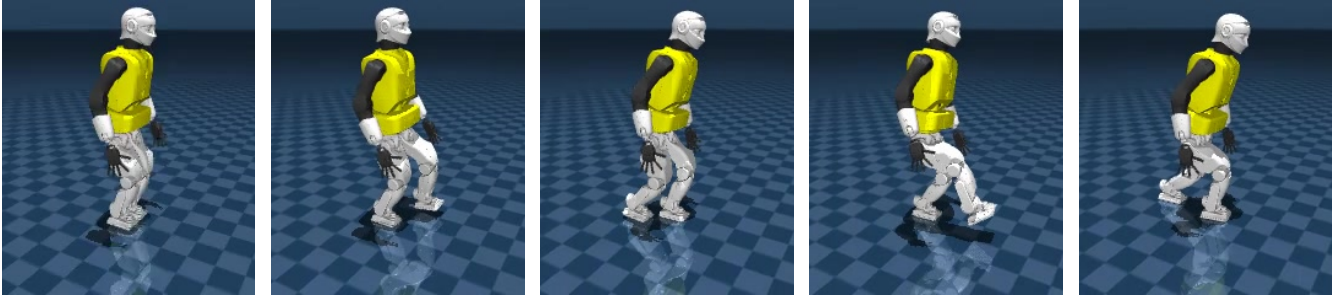


Figure 2: REEM-C walking forward at 1.0 m/s

plans for real hardware experiments are in progress. Mainly, the sim to real gap is being closed through better actuator modelling for the simulation and a task space inverse dynamics approach [11] for tracking the joint angles, center of mass position and foot contacts.

## 4    Discussion

This work presents initial results of a velocity-based RL locomotion policy for the REEM-C humanoid robot. The RL problem formulation is based on a periodic reward formulation to enforce a periodic stepping gait and is implemented in Brax/MJX for fast training performance, the first implementation for a real bipedal robot to the best knowledge of the authors. Preliminary results produced a stable, periodic walking policy at different command velocities and future hardware experiments will demonstrate the robustness and versatility.

## References

[1]  Brax. Accessed: 2024-04-30., https://github.com/google/brax

[2]  Caluwaerts, K., Iscen, A., Kew, J.C., Yu, W., Zhang, T., Freeman, D., Lee, K.H., Lee, L., Saliceti, S., Zhuang, V., Batchelor, N., Bohez, S., Casarini, F., Chen, J.E., Cortes, O., Coumans, E., Dostmohamed, A., Dulac-Arnold, G., Escontrela, A., Frey, E., Hafner, R., Jain, D., Jyenis, B., Kuang, Y., Lee, E., Luu, L., Nachum, O., Oslund, K., Powell, J., Reyes, D., Romano, F., Sadeghi, F., Sloat, R., Tabanpour, B., Zheng, D., Neunert, M., Hadsell, R., Heess, N., Nori, F., Seto, J., Parada, C., Sindhwani, V., Vanhoucke, V., Tan, J.: Barkour: Benchmarking animal-level agility with quadruped robots (arXiv:2305.14654) (May 2023), http://arxiv.org/abs/2305.14654, arXiv:2305.14654 [cs]

[3]  Dantec, E., Naveau, M., Fernbach, P., Villa, N., Saurel, G., Stasse, O., Taix, M., Mansard, N.: Whole-body model predictive control for biped locomotion on a torque-controlled humanoid robot. In: 2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids). p. 638–644. IEEE, Ginowan, Japan (Nov 2022), https://ieeexplore.ieee.org/document/10000129/

[4]  Kajita, S., Kanehiro, F., Kaneko, K., Fujiwara, K., Harada, K., Yokoi, K., Hirukawa, H.: Biped walking pattern generation by using preview control of zero-moment point. In: 2003 IEEE International Conference on Robotics and Automation (Cat. No.03CH37422). p. 1620–1626. IEEE, Taipei, Taiwan (2003), http://ieeexplore.ieee.org/document/1241826/

[5]  MuJoCo XLA MJX. Accessed: 2024-06-14., https://mujoco.readthedocs.io/en/stable/mjx.html

[6]  Radosavovic, I., Xiao, T., Zhang, B., Darrell, T., Malik, J., Sreenath, K.: Learning humanoid locomotion with transformers (arXiv:2303.03381) (Mar 2023), http://arxiv.org/abs/2303.03381, arXiv:2303.03381 [cs]

[7]  Rudin, N., Hoeller, D., Reist, P., Hutter, M.: Learning to walk in minutes using massively parallel deep reinforcement learning. In: Proceedings of the 5th Conference on Robot Learning. p. 91–100. PMLR (Jan 2022), https://proceedings.mlr.press/v164/rudin22a.html

[8] Siekmann, J., Godse, Y., Fern, A., Hurst, J.: Sim-to-real learning of all common bipedal gaits via periodic reward composition (arXiv:2011.01387) (Mar 2021), http://arxiv.org/abs/2011.01387, arXiv:2011.01387 [cs]

[9] Singh, R.P., Benallegue, M., Morisawa, M., Cisneros, R., Kanehiro, F.: Learning bipedal walking on planned footsteps for humanoid robots. In: 2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids). p. 686–693. IEEE, Ginowan, Japan (Nov 2022), https://ieeexplore.ieee.org/document/10000067/

[10] Singh, R.P., Xie, Z., Gergondet, P., Kanehiro, F.: Learning bipedal walking for humanoids with current feedback. IEEE Access 11, 82013–82023 (2023), arXiv:2303.03724 [cs]

[11] Efficient task space inverse dynamics (TSID) based on pinocchio. Accessed: 2024-04-30., https://github.com/stack-of-tasks/tsid