

Sim-to-Real Learning of All Common Bipedal Gaits via Periodic Reward Composition

Jonah Siekmann*, Yesh Godse*, Alan Fern, Jonathan Hurst Collaborative Robotics and Intelligent Systems Institute Oregon State University

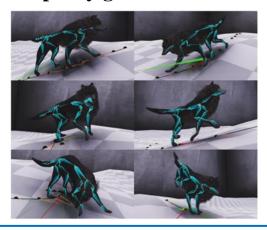
Article Presentation

Group 6: Xiaoyu Yan 395091 Zekun Wang 393910 Yuansheng Zhou 393916



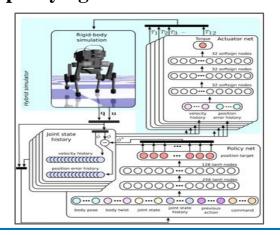
Previous Methods of Generating Gait Reward Functions

1. Use trajectory reference to specify gait rewards^[1]



- Unadaptability in varying environment
- Challenging to derive reference trajectories

2. Reference-free approaches for specifying reward functions^[2]



- Difficult to achieve specific gait characteristics
- Extending to other behaviors is hard



^[1] H. Zhang, S. Starke, T. Komura, and J. Saito, "Mode-adaptive neural networks for quadruped motion control," ACM Transactions on Graphics, vol. 37, no. 4, pp. 1–11, 2018, ISSN: 15577368. DOI: 10.1145/3197517.3201366

Main Idea and Contribution

The paper introduces a periodic reward composition framework to enable the bipedal robot Cassie to learn and perform a wide range of gaits (e.g., walking, hopping, running) through sim-to-real reinforcement learning, without relying on predefined trajectories.

Controbutions:

- 1. This work presents a principled framework for designing reward functions that can naturally capture all of the periodic bipedal locomotion gaits.
- 2. This work demonstrated this framework for sim-to-real RL of all common bipedal gaits, including walking, running, galloping, skipping, and hopping, without using a motion capture dataset or reference trajectories.



Key Formulas

Periodic Reward Composition

$$R(s,\phi) = \beta + \sum_{i} R_{i}(s,\phi)$$
$$R_{i}(s,\phi) = c_{i} \cdot I_{i}(\phi) \cdot q_{i}(s)$$

- c_i is a phase-specific coefficient controlling the weight of the reward for each phase.
- $I_i(\phi)$ is a phase indicator function that activates the reward depending on the current position in the gait cycle.
- $q_i(s)$ represents physical measurements (such as foot forces or velocities) relevant to the specific phase.

Phase Indicator Function

$$P(I_i(\phi) = x) = \begin{cases} P(A_i < \phi < B_i) & \text{if } x = 1\\ 1 - P(A_i < \phi < B_i) & \text{if } x = 0 \end{cases}$$

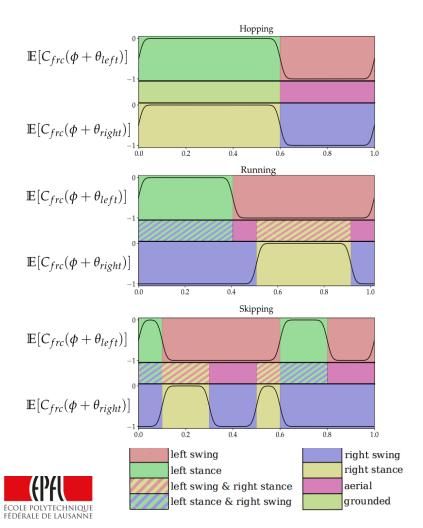
Here, A_i and B_i mark the start and end of the specific phase.

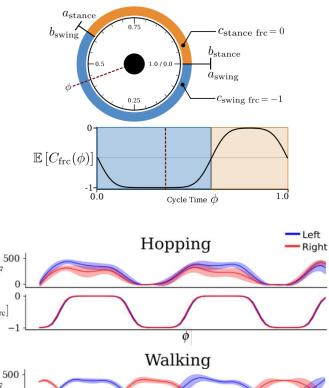
Bipedal Gait Reward Formula for Left and Right Feet

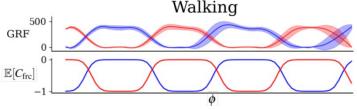
$$\mathbb{E}[R_{\text{bipedal}}(s, \phi)] = \mathbb{E}[C_{\text{frc}}(\phi + \theta_{\text{left}})] \cdot q_{\text{left frc}}(s) \\ + \mathbb{E}[C_{\text{frc}}(\phi + \theta_{\text{right}})] \cdot q_{\text{right frc}}(s) \\ + \mathbb{E}[C_{\text{spd}}(\phi + \theta_{\text{left}})] \cdot q_{\text{left spd}}(s) \\ + \mathbb{E}[C_{\text{spd}}(\phi + \theta_{\text{right}})] \cdot q_{\text{right spd}}(s)$$

- $C_{\rm frc}$ and $C_{\rm spd}$ represent phase reward functions, assessing foot forces and velocities respectively.
- θ_{left} and θ_{right} are the phase offsets for the left and right feet, ensuring alternating steps.
- $q_{\text{left frc}}(s)$ and $q_{\text{right frc}}(s)$ measure the foot forces for each foot, while $q_{\text{left spd}}(s)$ and $q_{\text{right spd}}(s)$ measure the velocities.









GRF

 $\mathbb{E}[C_{\mathrm{frc}}]$

Result

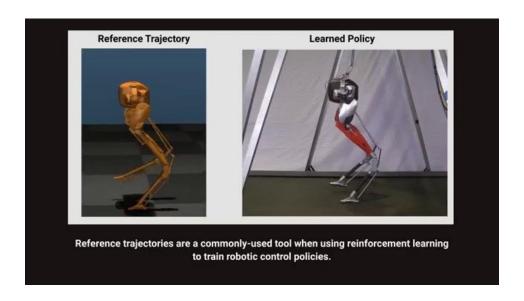
Single-Gait Policy Results: The framework successfully trained policies Single-Gait by holding specific parameters constant.

Multi-Gait Policy Results: Training policies for smooth transitions between gaits (like hopping to running) proved challenging. Adjusting both cycle offsets and phase ratios sometimes led to undesired behaviors, such as asymmetrical gaits or mixed behaviors.

Solution with Transition Penalties: "transition penalties" are specific rewards or costs added to distinguish different behaviors during training, ensuring smoother transitions.



Experiment Background and Objectives



Since the robot lacks external sensors, it relies entirely on its internal states for adjusting its gait, which is effectively 'blind'



Experiment





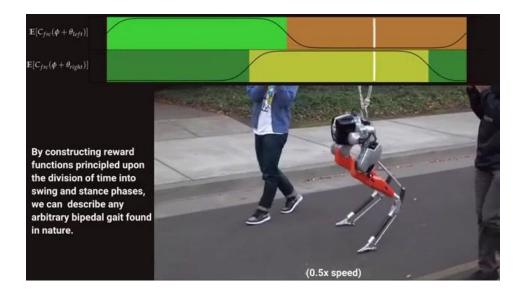




Tests involved the bipedal robot Cassie in environments like sidewalks, curbs, slopes, and stairs.



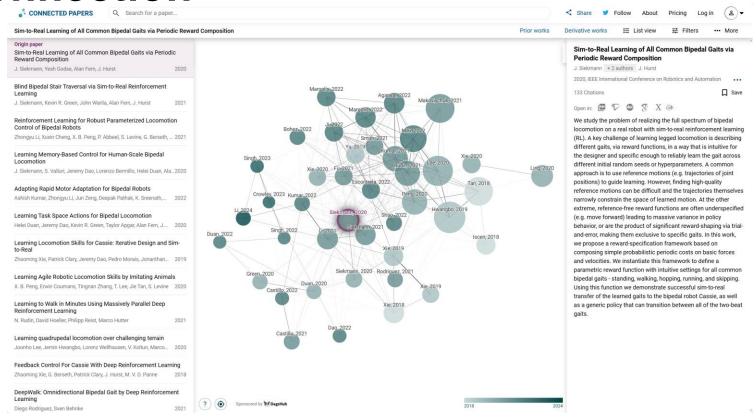
Experiment



Smooth transitions between gaits, such as moving from walking to running on a road



Connection





Connection





Article

Adaptive Gait Acquisition through Learning Dynamic Stimulus Instinct of Bipedal Robot

Yuanxi Zhang ¹, Xuechao Chen ^{1,2}, Fei Meng ^{1,2,*}, Zhangguo Yu ^{1,2}, Yidong Du ¹, Zishun Zhou ¹
and Iunvao Gao ^{1,2}

- School of Mechatronical Engineering, Beijing Institute of Technology, Beijing 100081, China; zhangyuanxi@bit.edu.cn (Y.Z.); chenxuechao@bit.edu.cn (X.C.); yuzg@bit.edu.cn (Z.Y.); duvidong@bit.edu.cn (Y.D.); zishun.zhou@outlook.com (Z.Z.); gaojunyao@bit.edu.cn (J.G.)
- 2 Key Laboratory of Biomimetic Robots and Systems, Ministry of Education, Beijing 100081, China
- Correspondence: mfly0208@bit.edu.cn

Abstract: Standard alternating leg motions serve as the foundation for simple bipedal gaits, and the effectiveness of the fixed stimulus signal has been proved in recent studies. However, in order to address perturbations and imbalances, robots require more dynamic gaits. In this paper, we introduce dynamic stimulus signals together with a bipedal locomotion policy into reinforcement learning (RL). Through the learned stimulus frequency policy, we induce the bipedal robot to obtain both three-dimensional (3D) locomotion and an adaptive gait under disturbance without relying on an explicit and model-based gait in both the training stage and deployment. In addition, a set of specialized reward functions focusing on reliable frequency reflections is used in our framework to ensure correspondence between locomotion features and the dynamic stimulus. Moreover, we demonstrate efficient sim-to-real transfer, making a bipedal robot called BITeno achieve robust locomotion and disturbance resistance, even in extreme situations of foot sliding in the real world. In detail, under a sudden change in torso velocity of -1.2 m/s in 0.65 s, the recovery time is within 1.5-2.0 s.

Keywords: reinforcement learning; bipedal robot; adaptive locomotion; period dynamic gait

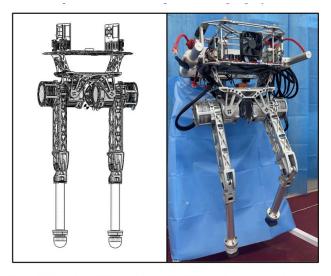


Figure 2. The design of BITeno platform. (**Left**) The mechanical design in simulation, were the feature of each link was assigned according to real materials. (**Right**) The physical robot with an electrical system onboard.



Pros & Cons

Pros

Versatile Gait Learning
Eliminates Dependency on Reference Trajectories
Adaptability to Environmental Disturbances

Cons

Limited Adaptability to High-Level Disturbances



Possible exam questions

- •How can reward functions be designed to achieve different gait behaviors without reference trajectories?
- •Explain the role of periodic reward composition in learning bipedal gaits for a robot.

Answer 1: Periodic reward composition allows for the specification of distinct gait characteristics by defining phases (swing and stance) with specific penalties or rewards. This framework guides the reinforcement learning process, enabling the robot to learn stable bipedal gaits like walking, running, hopping, and skipping by rewarding appropriate foot forces and velocities during each phase.

Answer 2: Without reference trajectories, reward functions can be designed using periodic reward composition. By assigning rewards or penalties to foot speed and force during different gait phases (such as penalizing foot force during the swing phase and penalizing foot speed during the stance phase), the framework encourages the robot to learn and achieve distinct gait characteristics without relying on specific motion trajectories.



Outdoor Experiments of Bipedal Robot Gait Learning

Aim to demonstrate the robot's stability and adaptability in real-world outdoor settings



