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

Hansol Kang, Combined Course, 9th semester
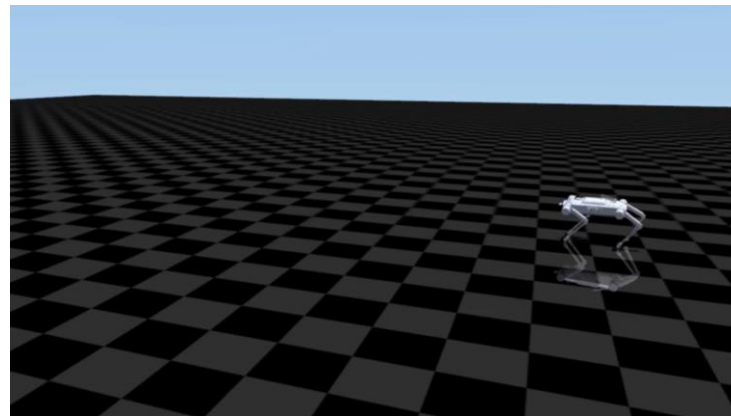
SK UNIVERSITY

Robotics Innovatory

http://mecha.skku.ac.kr

# Contents

# Introduction

▶ Quadruped robot control via RL



Walk on flat ground (same as training)

ANYmal runs faster than ever before.

Learning to Walk via Deep Reinforcement Learning. Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, Sergey Levine. Robotics: Science and Systems (RSS). 2019.
Learning agile and dynamic motor skills for legged robots, Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen Koltun, Marco Hutter, Science Robotics, 2019

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



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
{siekmanj, godsey, afern, jhurst}@oregonstate.edu

Abstract— We study the problem of realizing the full spectrum of bipedal locomotion on a real robot with sim-to-real reinforcement learning (RL). A key challenge of learning legged locomotion is describing different gaits, via reward functions, in a way that is intuitive for the designer and specific enough to reliably learn the gait across different initial random seeds or hyperparameters. A common approach is to use reference motions (e.g. trajectories of joint positions) to guide learning. However, finding high-quality reference motions can be difficult and the trajectories themselves narrowly constrain the space of learned motion. At the other extreme, reference-free reward functions are often underspecified (e.g. move forward) leading to massive variance in policy behavior, or are the product of significant reward-shaping via trial-and-error, making them exclusive to specific gaits. In this work, we propose a reward-specification framework based on composing simple probabilistic periodic costs on basic forces and velocities. We instantiate this framework to define a parametric reward function with intuitive settings for all common bipedal gaits - standing, walking, hopping, running, and skipping. Using this function we demonstrate successful sim-to-real transfer of the learned gaits to the bipedal robot Cassie, as well as a generic policy that can transition between all of the two-beat gaits.

## I. INTRODUCTION

Using reinforcement learning (RL) to learn all of the common bipedal gaits found in nature for a real robot is an unsolved problem. A key challenge of learning a specific locomotion gait via RL is to communicate the gait behavior through the reward function. In general, a specific gait can be viewed as a dynamic process that has a characteristic periodic structure, but is also able to flexibly adapt to moderate environment disturbances. This suggests two considerations when designing a gait reward function. First, the reward must be specific enough to produce the desired gait characteristic when optimized. Second, to account for the fact that there is uncertainty about the exact details of a gait in the context of specific terrain and dynamic conditions, the reward should not be overly constraining.

The common use of reference trajectories to specify gait-specific rewards, e.g., [1]–[5], partly addresses the first consideration above, but mostly ignores the second. In particular, a reference trajectory only captures a small part of the variation needed to realize a gait characteristic under varying conditions. Thus, attempting to adhere to such a trajectory can prevent learning a characteristic gait that is more robust and/or efficient, not to mention that deriving

feasible reference trajectories for a particular desired gait can be very challenging in the first place.

Reference-free approaches to specifying reward functions for locomotion are often highly underspecified, for example, those used in the OpenAI Gym [6] locomotion benchmarks. With this starting point, achieving a specific gait characteristic requires iterations of heuristic reward-function adjustments, based on observed RL performance, until arriving at a desired behavior. This approach can be tedious when it works and is unreliable as a general framework. Other reference-free approaches structure the reward around a specific type of locomotion behavior [7] without being easily extended to other behaviors.

The first contribution of our work is to present a principled framework for designing reward functions that can naturally capture all of the periodic bipedal locomotion gaits. We are motivated by the fact that all common bipedal gaits can be defined by periodic *swing phases* (foot swinging in the air) and *stance phases* (foot planted on the ground) for each foot [8]. A fundamental distinction between swing and stance
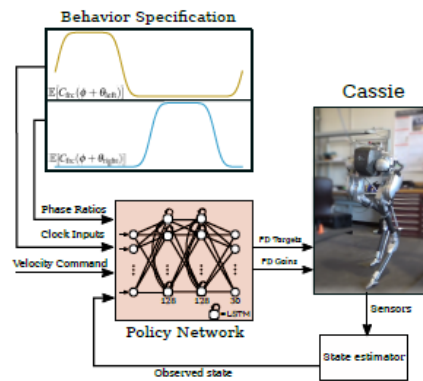
Fig. 1. In this work, we present a reward design framework which makes it easy to learn policies which can stand, walk, run, gallop, hop, and skip on hardware. We condition the reward function based on a number of gait parameters, and also provide these parameters to the LSTM policy, which outputs PD joint position targets and PD gains to the robot.

* denotes equal contribution, order determined by coin toss

▶ 2020.11. arXiv

▶ 2021 ICRA Submitted paper

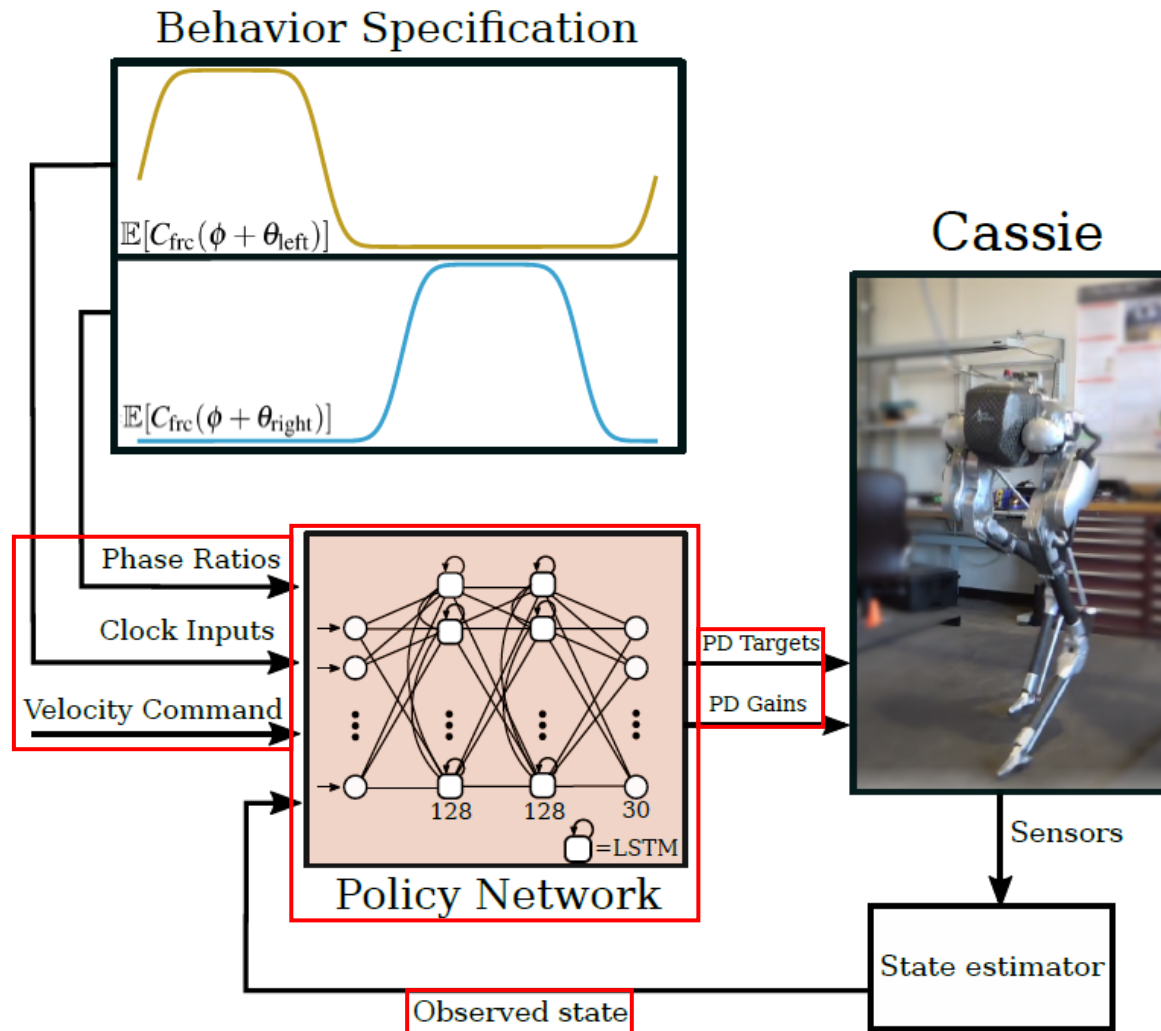▶ Jonah Siekmann*, Yesh Godse*, Alan Fern, Jonathan Hurst

▶ Oregon State University

▶ Propose a reward function dependent on a periodic signal, and through this, the agent is possible to learn to walk in various gait according to the user's command, and natural gait transition is possible.

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



Input : Phase Ratio($r$), Clock Inputs($Clock$),
         Velocity Command($C$), Observation($O$)

Long Short-Term Memory
2 Hidden Layers(128 Hidden nodes each)

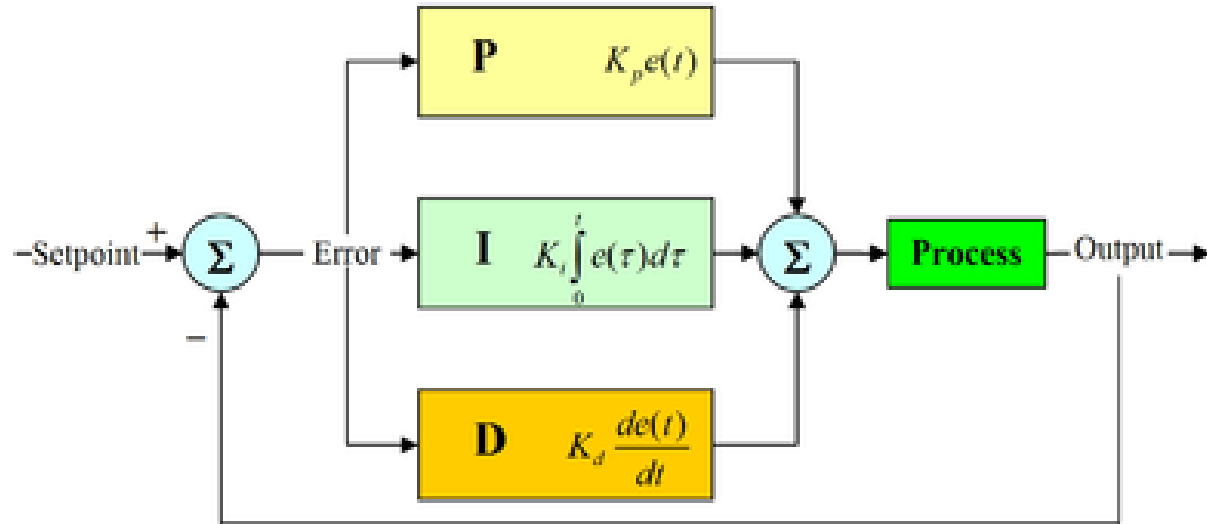Output : Target Joint Position(10)
          P, D Gain(10 × 2)

Algorithm : PPO
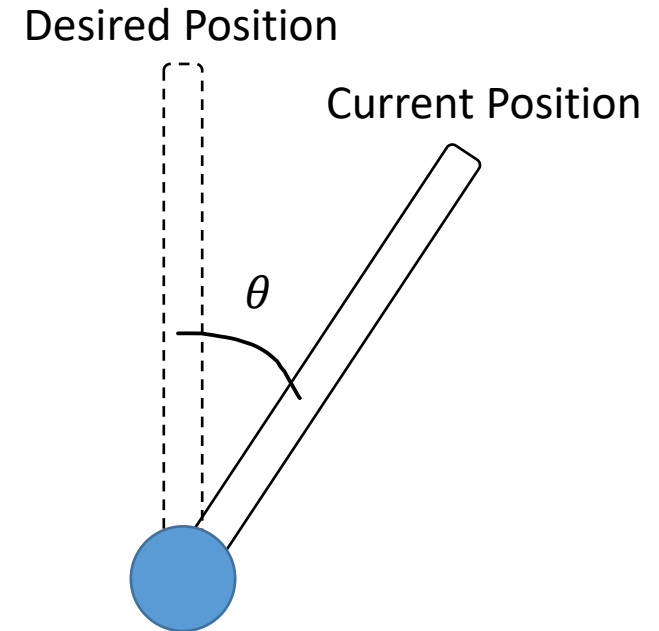
# Preliminaries for robot control



$$Manipulated\ Variable = K_p e(t) + K_i \int_0^t e(t)dt + K_d \frac{de}{dt}$$

▶ Sim-to-Real Leaning of All Common Bipedal Gaits via <mark>Periodic Reward Composition</mark>

$$R(\boldsymbol{s}, \phi) = \beta + \sum_i R_i(\boldsymbol{s}, \phi)$$

Normalized time period $\phi$

$$R_i(\boldsymbol{s}, \phi) = c_i \cdot I_i(\phi) \cdot q_i(\boldsymbol{s})$$

Phase coefficient $c_i$
Phase indicator $I_i(\phi)$
Phase reward measurement $q_i(\boldsymbol{s})$

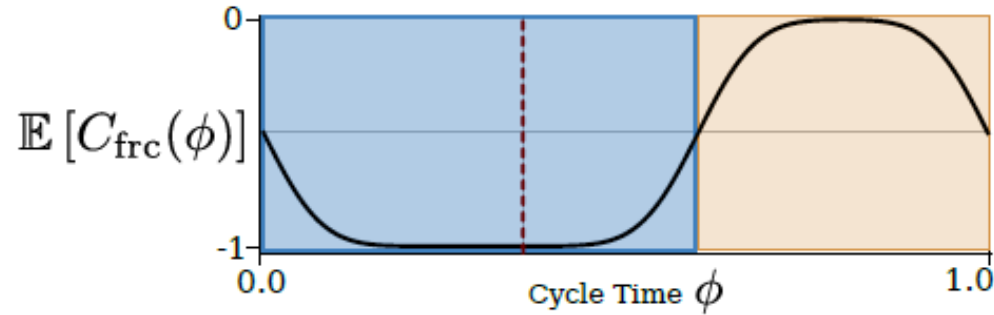$$\mathbb{E}[R(\boldsymbol{s}, \phi)] = \sum_i^n c_i \cdot \mathbb{E}[I_i(\phi)] \cdot q_i(\boldsymbol{s}) + \beta$$

$$I_i(\phi) = \begin{cases} 0, & Stance \\ 1, & Swing \end{cases}$$

$$\mathbb{E}[C_{frc}(\phi)] = c_{swing\ frc} \cdot \mathbb{E}[I_{swing\ frc}(\phi)] + c_{stance\ frc} \cdot \mathbb{E}[I_{stance\ frc}(\phi)]$$

$$c_{stance\ frc} = 0, \quad c_{swing\ frc} = -1$$

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



- Unipedal

$$\mathbb{E}[R_{unipedal}(\boldsymbol{s}, \phi)] = \mathbb{E}[C_{frc}(\phi)] \cdot q_{frc}(\boldsymbol{s}) + \mathbb{E}[C_{spd}(\phi)] \cdot q_{spd}(\boldsymbol{s})$$
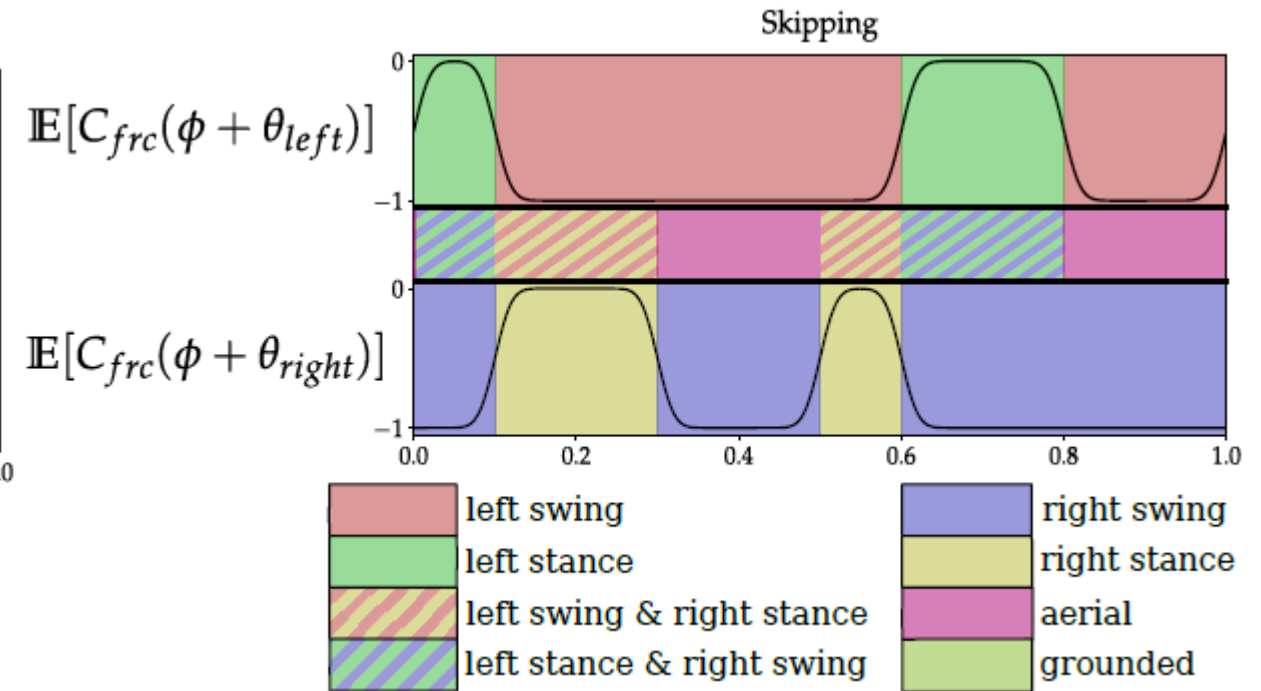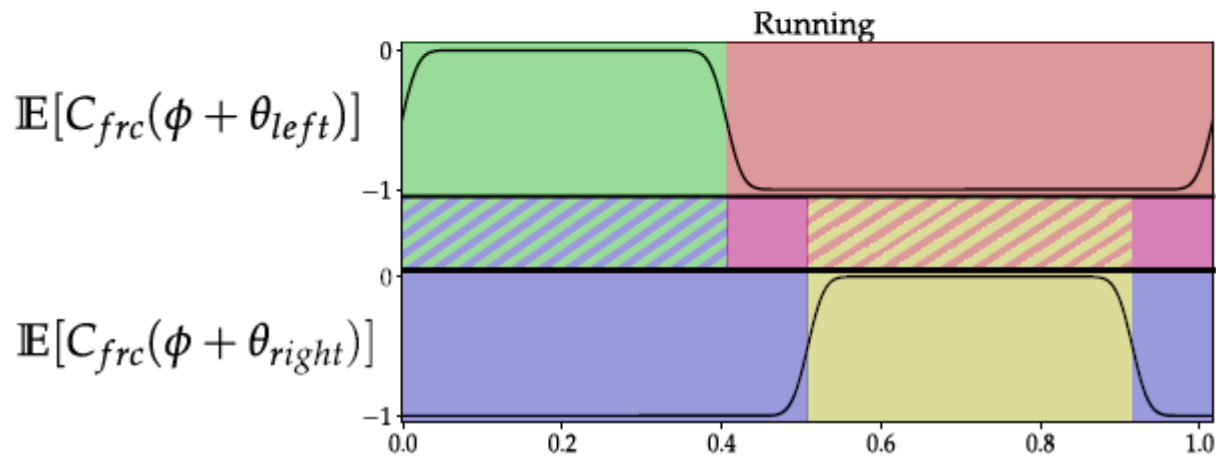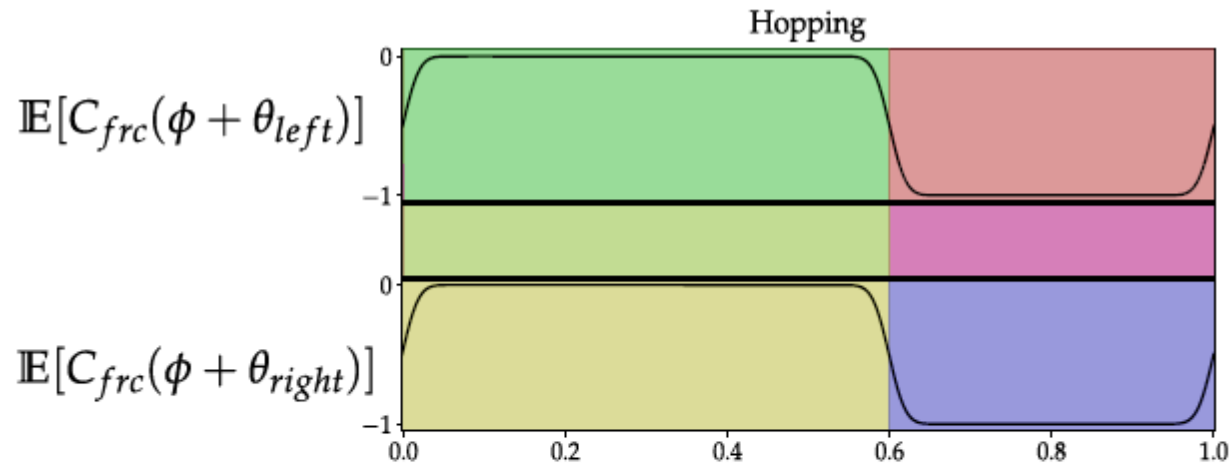
- Biped

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

⋮

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

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

- Total Reward Function

$$R(\boldsymbol{s}, \phi) = \beta + \sum_i R_i(\boldsymbol{s}, \phi)$$

$$\mathbb{E}[R(\boldsymbol{s}, \phi)] = \mathbb{E}[R_{biped}(\boldsymbol{s}, \phi)] + R_{cmd}(\boldsymbol{s}) + R_{smooth}(\boldsymbol{s}) + \beta$$

$$R_{cmd}(\boldsymbol{s}) = (-1) \cdot q_{\dot{x}}(\boldsymbol{s}) + (-1) \cdot q_{\dot{y}}(\boldsymbol{s}) + (-1) \cdot q_{orientation}(\boldsymbol{s})$$

$$R_{smooth}(\boldsymbol{s}) = (-1) \cdot q_{action\_diff}(\boldsymbol{s}) + (-1) \cdot q_{torque}(\boldsymbol{s}) \\ + (-1) \cdot q_{pelvis\_acc}(\boldsymbol{s})$$

- Dynamic Randomization

| Parameter | Unit | Range |
|---|---|---|
| Joint damping | Nms/rad | $[0.3, 4.0] \times$ default values |
| Joint mass | kg | $[0.5, 1.5] \times$ default values |
| Ground Friction | – | $[0.35, 1.1]$ |
| Ground Slope | rad | $[-0.03, 0.03]$ |
| Joint Encoder Offset | rad | $[-0.05, 0.05]$ |

TABLE I. The ranges for randomization of several dynamics parameters during training. We use a uniform distribution over the given ranges for all listed parameters.
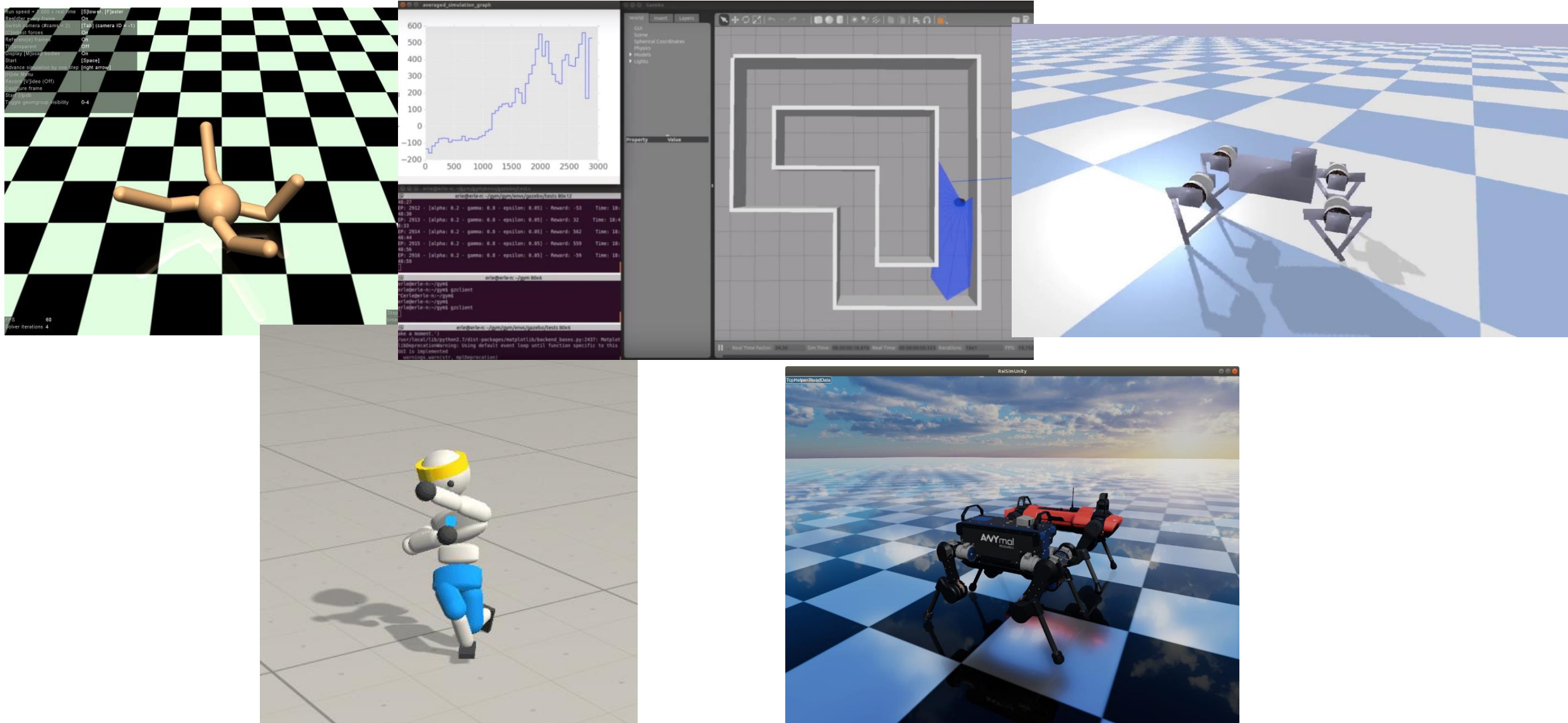
Sim-to-Real Leaning of All Common Bipedal Gaits via Periodic Reward Composition

In this work, we introduce a framework which makes it possible to construct reward functions for learning locomotion behaviors without the use of reference trajectories.

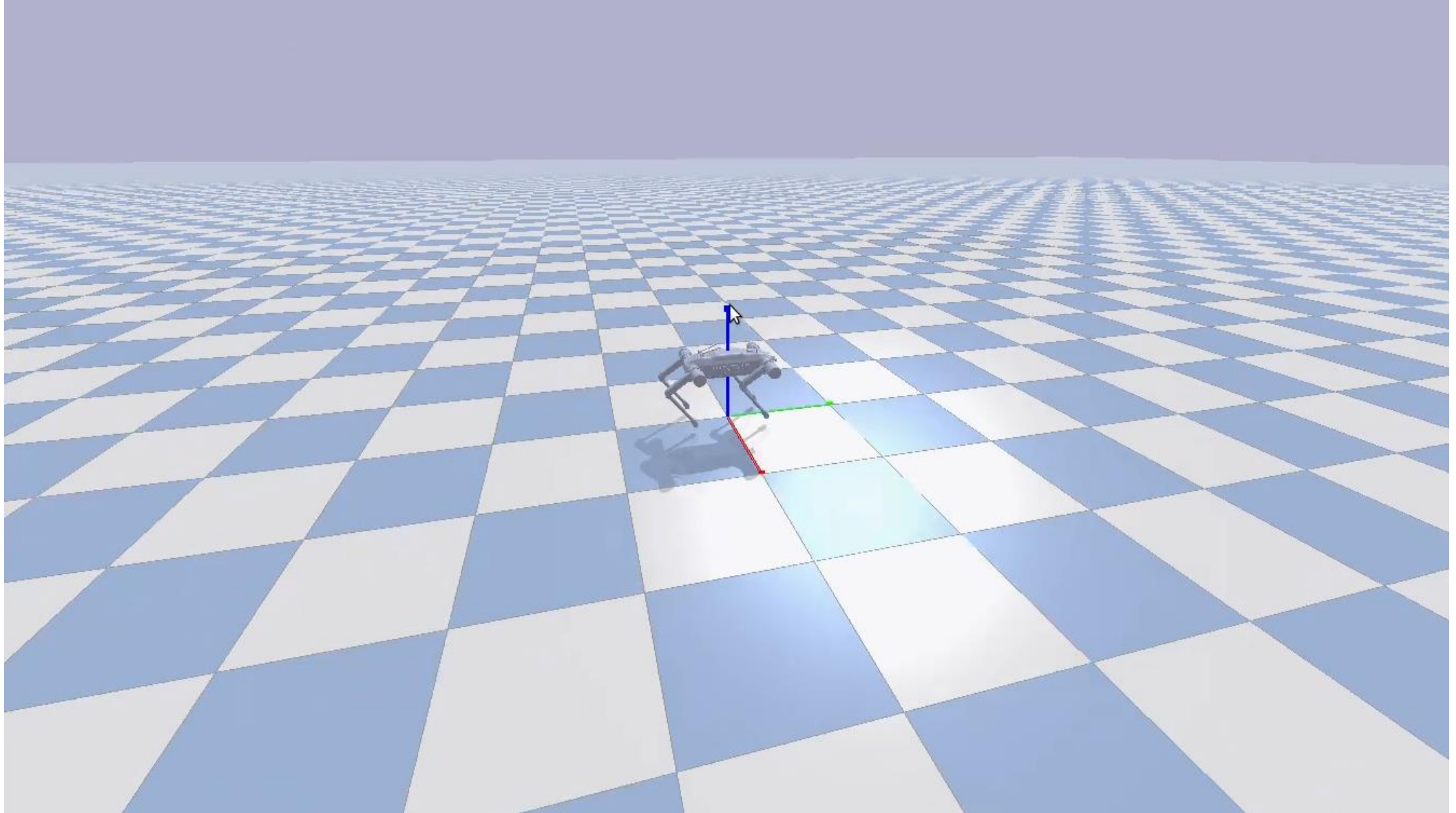# Simulator

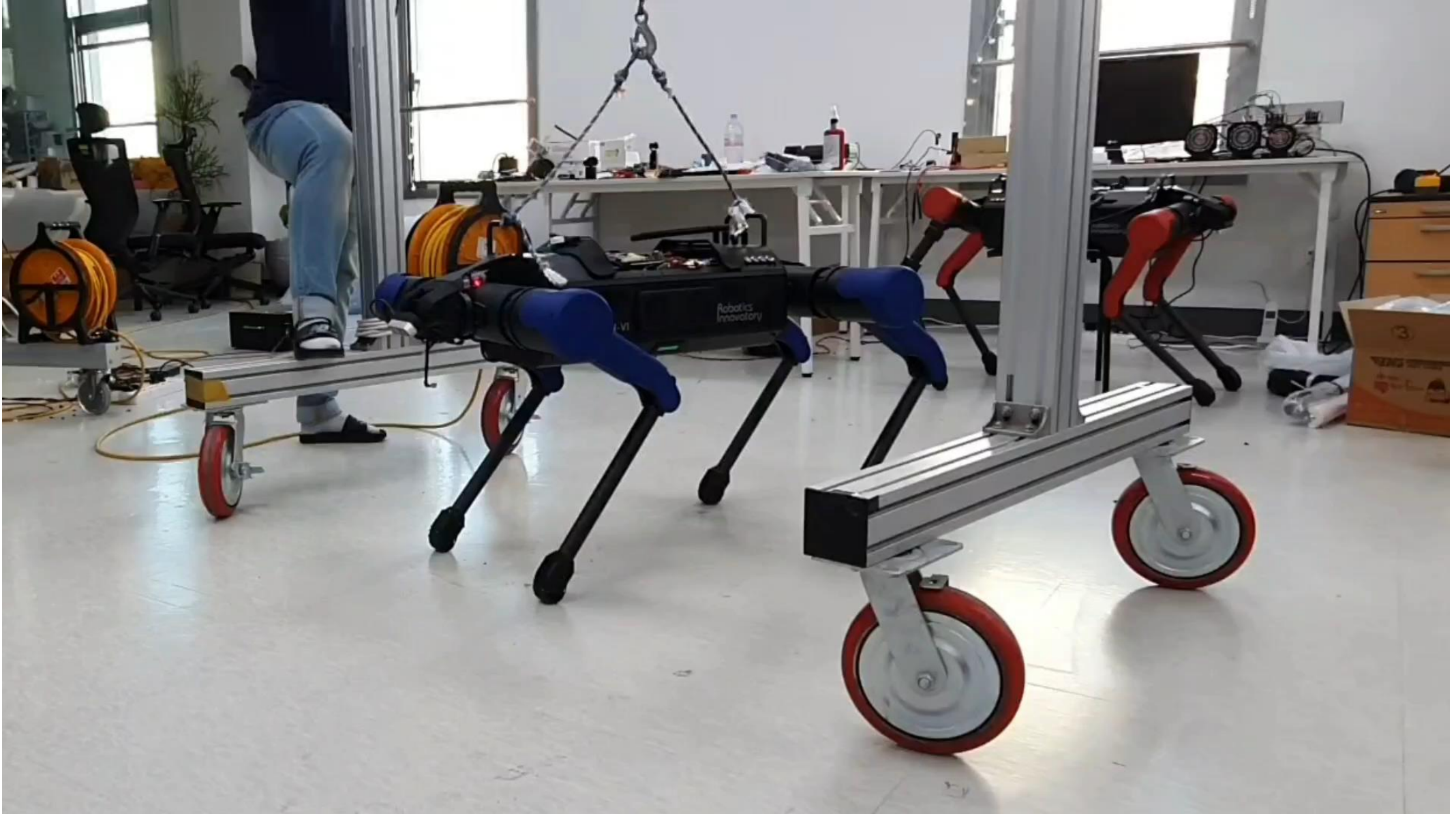▶ MuJoCo? Gazebo? Pybullet? Unity? RaiSim?

# Various Gaits of Quadruped

▶ 12-DOF AiDIN-VI + Task Space + Periodic Reward

▶ Sim-to-Real Test

# Q & A

## Thank you for your attention