

# Legged Locomotion: Trajectory Optimization, Machine Learning, and Bio-Inspired Control

Guillaume Bellegarda



Skild AI, 2025



Kim et al., 2025

Is Legged Robots "Solved"?



Kungfu Kid

V6.0

Unitree G1

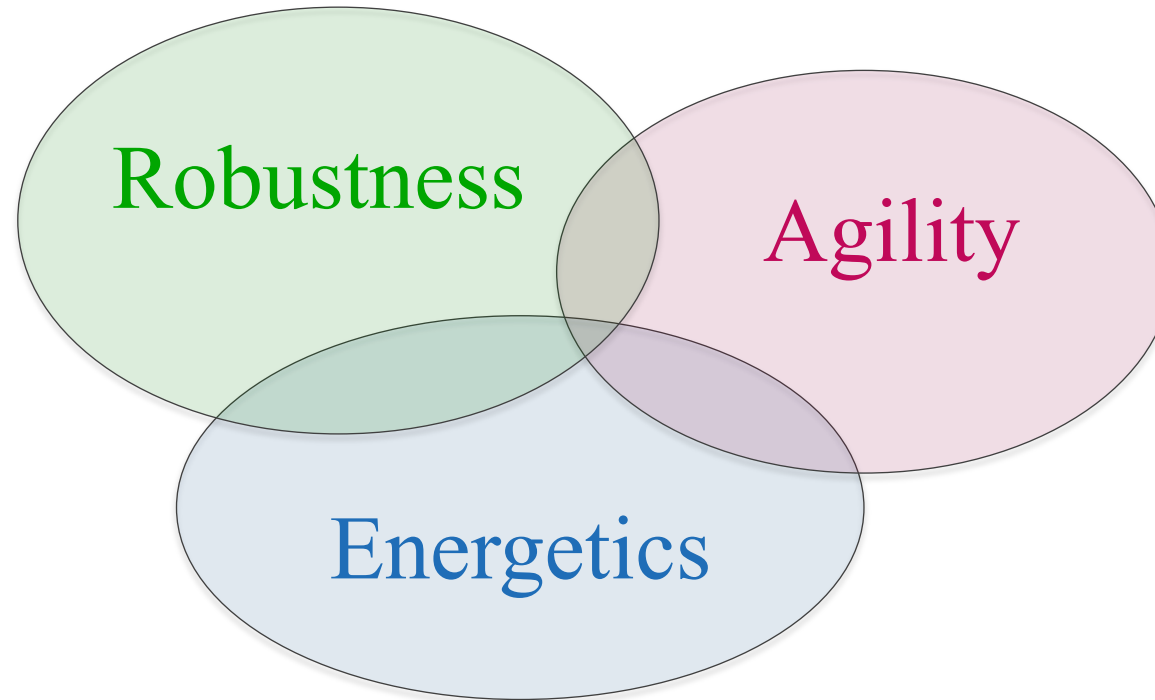
Unitree, 2025



1x speed

Yang et al., 2025

# Locomotion Goals



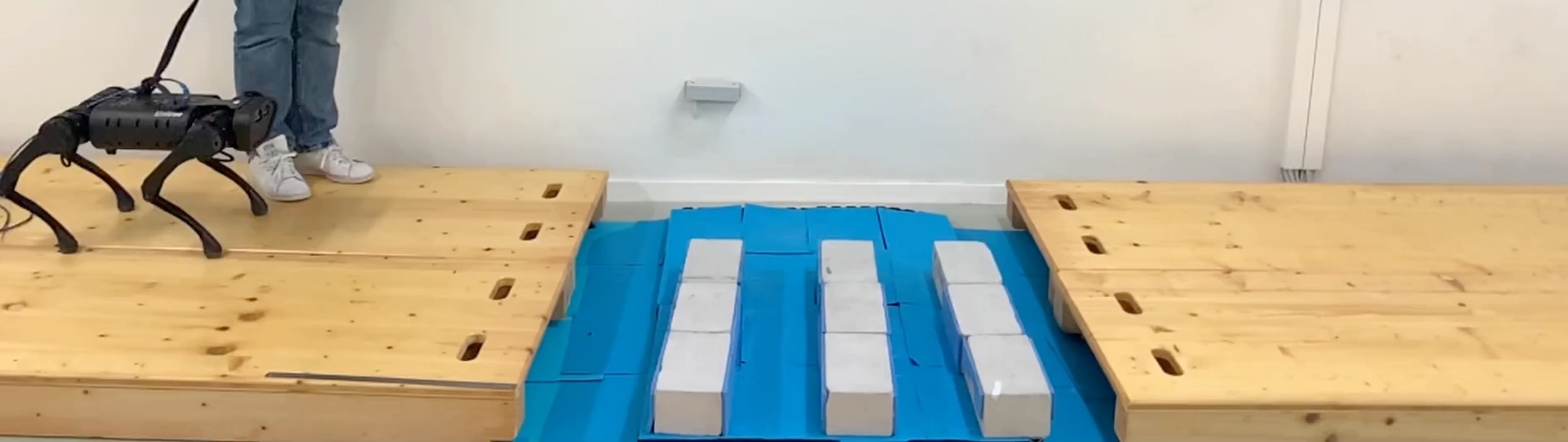
...in stochastic environments (variability)



Robustness  
unknown variability

Agility

Energetics

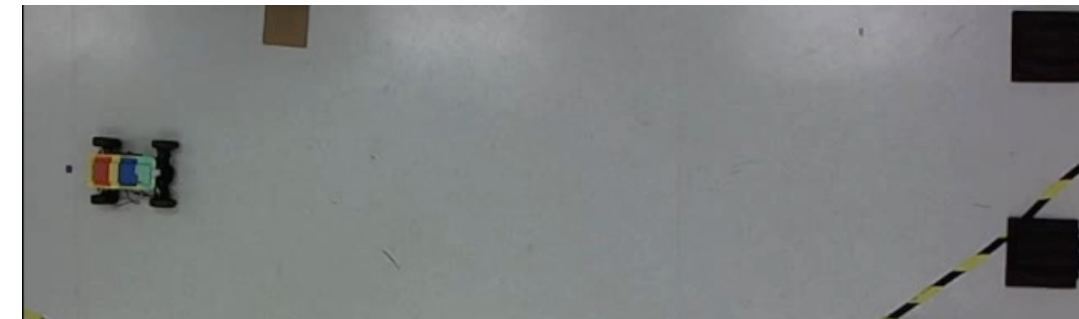


Robustness  
unknown variability



Agility  
known variability

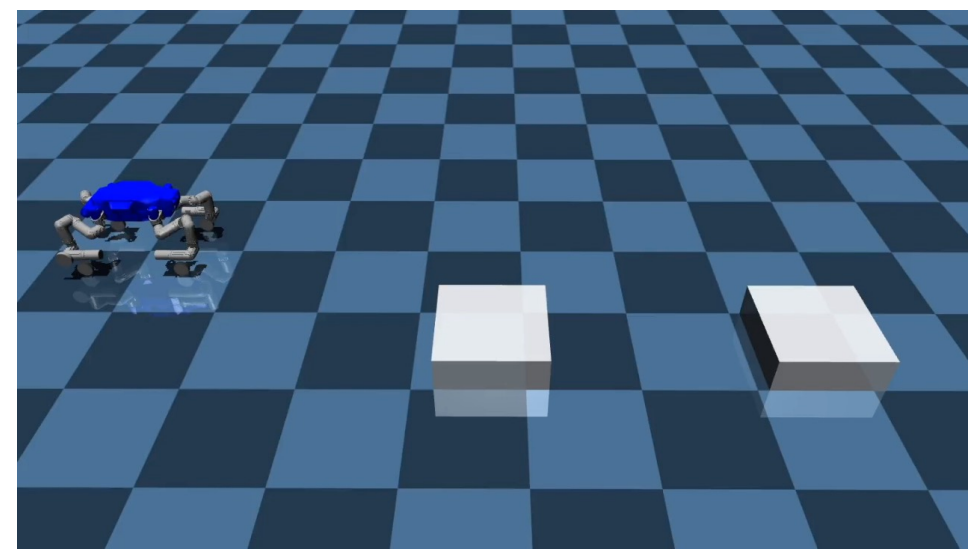
Energetics



Robustness  
unknown variability

Agility  
known variability

Energetics  
as efficient as practical



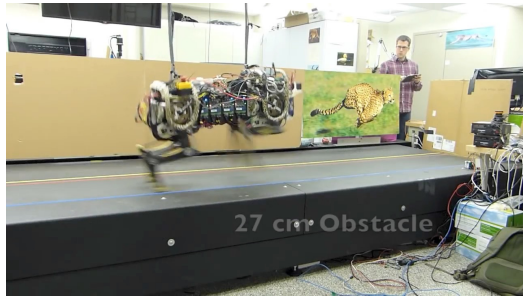
G. Bellegarda, K. Byl. "Trajectory Optimization for a Wheel-Legged System for Dynamic Maneuvers that Allow for Wheel Slip," CDC 2019.

G. Bellegarda, K. Byl. "Design and Evaluation of Skating Motions for a Dexterous Quadruped," ICRA 2018.

# Robotics Approaches

## Model-based methods

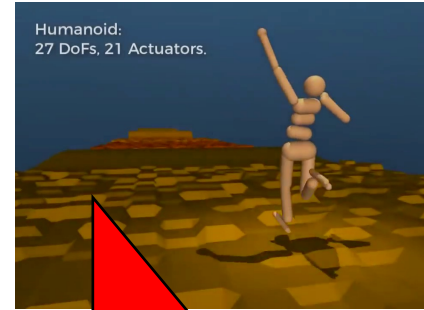
i.e. trajectory optimization, model-predictive control (MPC), Central Pattern Generators (CPGs)



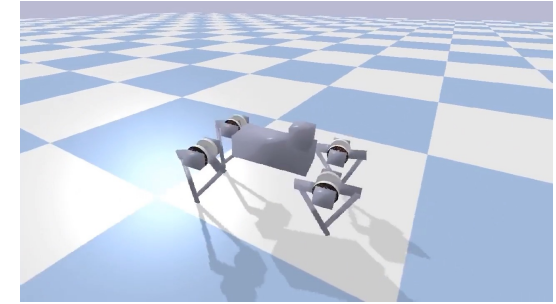
[Park et al. 2017]



[Boston Dynamics 2017]



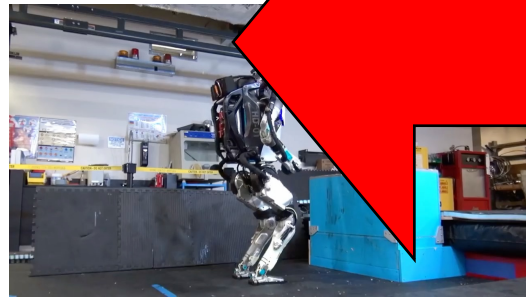
[Mettetal et al. 2017]



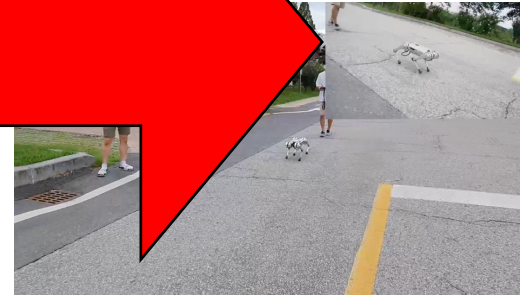
[Tan et al. 2018]



[Kim et al. 2019]



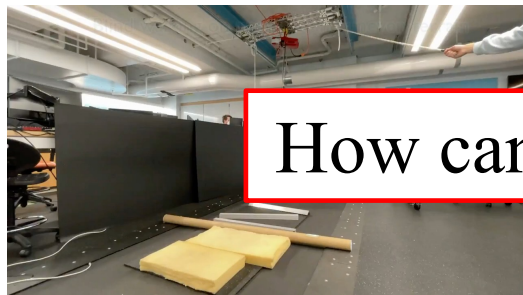
[Boston Dynamics 2018]



[Ji et al. 2022]



[Hwangbo et al. 2019]



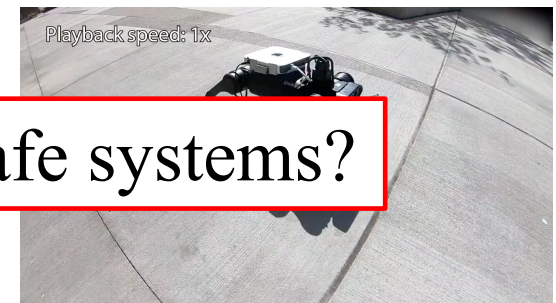
[Xiong et al. 2021]



[DRC 2015]



[Lee et al. 2020]



[Yang et al. 2021]

## Learning

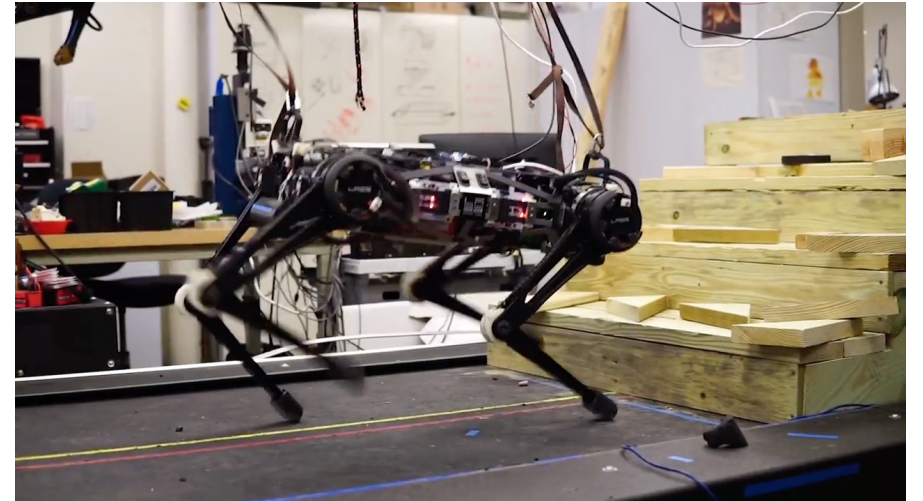
i.e. reinforcement learning, imitation learning

How can we build **generalizable, agile, robust, and safe systems?**

# How to create locomotion controllers?

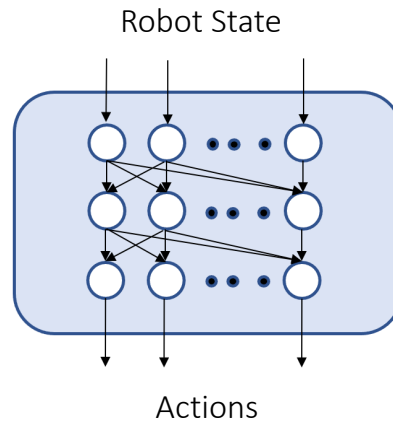
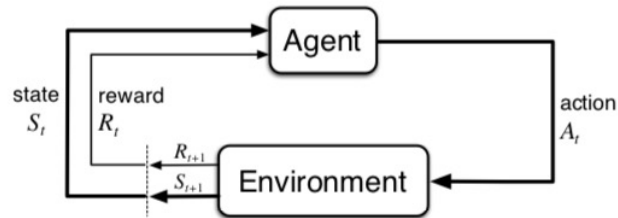
## 1. Trajectory optimization (model-predictive control)

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{u}} \quad & \sum_{i=0}^{k-1} \|\mathbf{x}_{i+1} - \mathbf{x}_{i+1, \text{ref}}\|_{\mathbf{Q}_i} + \|\mathbf{u}_i\|_{\mathbf{R}_i} \\ \text{subject to} \quad & \mathbf{x}_{i+1} = \mathbf{A}_i \mathbf{x}_i + \mathbf{B}_i \mathbf{u}_i, i = 0 \dots k-1 \\ & \underline{\mathbf{c}}_i \leq \mathbf{C}_i \mathbf{u}_i \leq \bar{\mathbf{c}}_i, i = 0 \dots k-1 \\ & \mathbf{D}_i \mathbf{u}_i = 0, i = 0 \dots k-1 \end{aligned}$$



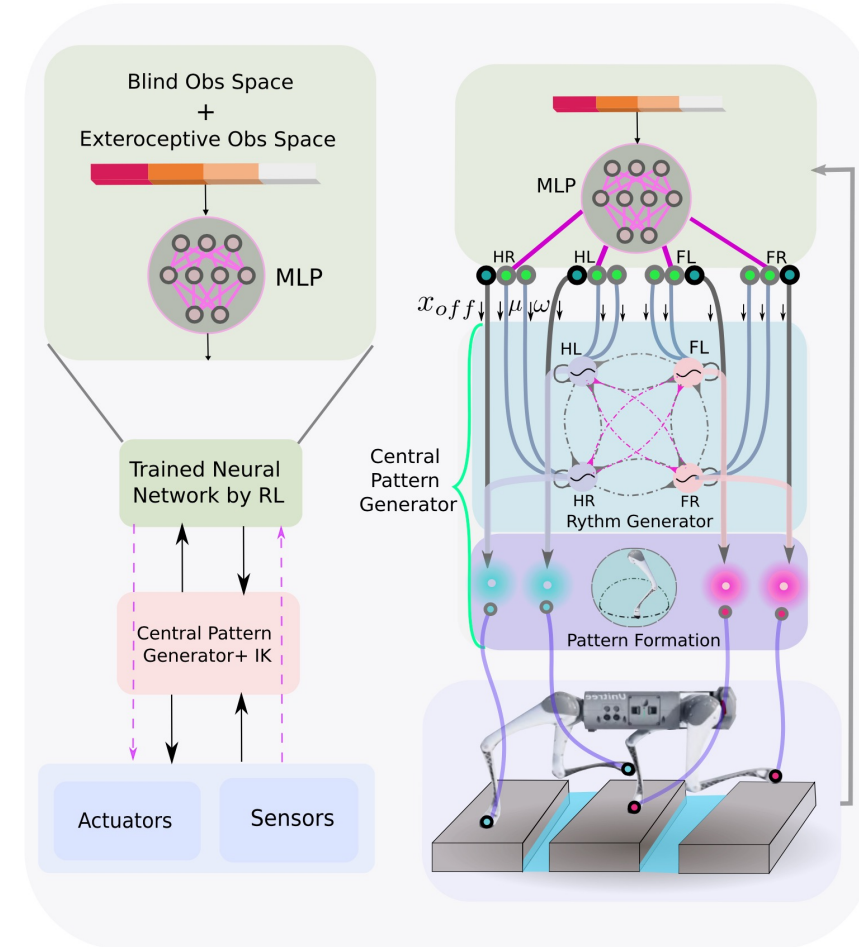
# How to create locomotion controllers?

1. Trajectory optimization (model-predictive control)
2. Deep reinforcement learning / imitation learning

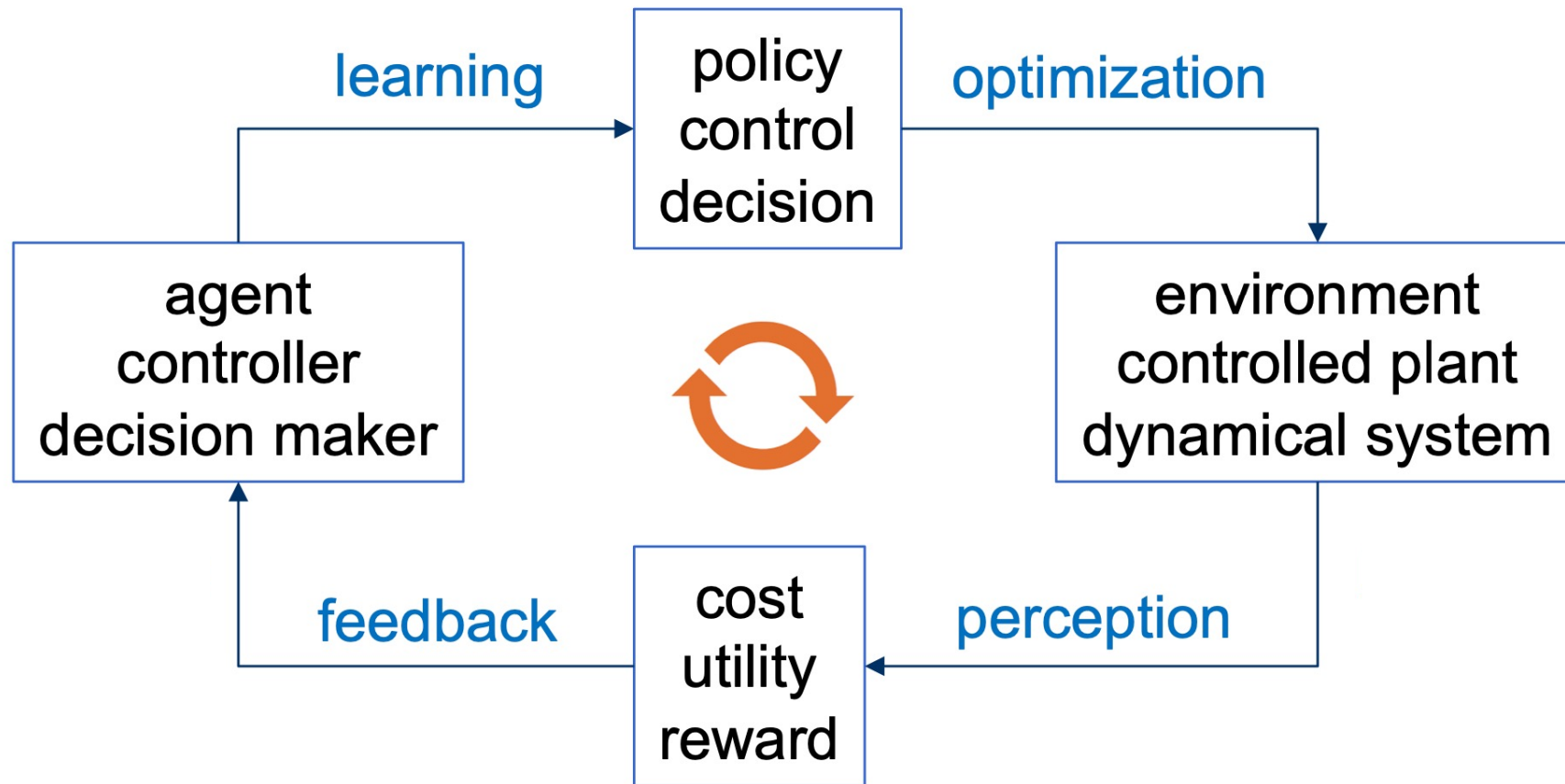


# How to create locomotion controllers?

1. Trajectory optimization (model-predictive control)
2. Deep reinforcement learning / imitation learning
3. Bio-inspired learning



# Common Setting: Closed-Loop Autonomous System



# Trajectory Optimization

minimize	$J(t_0, t_F, x(t_0), x(t_F)) + \int_{t_0}^{t_F} w(\tau, x(\tau), u(\tau)) d\tau$	
subject to	$\dot{x}(t) = f(t, x(t), u(t))$	System Dynamics
	$h(t, x(t), u(t)) \leq 0$	Path Constraint
	$g(t_0, t_F, x(t_0), x(t_F)) \leq 0$	Boundary Constraint
	$x_{low} \leq x(t) \leq x_{upp}$	Path Bound on State
	$u_{low} \leq u(t) \leq u_{upp}$	Path Bound on Control
	$t_{low} \leq t_0 < t_F \leq t_{upp}$	Bounds on Initial and Final Times
	$x_{0,low} \leq x(t_0) \leq x_{0,upp}$	Bound on Initial State
	$x_{F,low} \leq x(t_F) \leq x_{F,upp}$	Bound on Final State

# Direct Collocation

$$\begin{array}{ll}
 \underset{t_0, t_F, x(t), u(t)}{\text{minimize}} & J(t_0, t_F, x(t_0), x(t_F)) + \int_{t_0}^{t_F} w(\tau, x(\tau), u(\tau)) d\tau \\
 \text{subject to} & \dot{x}(t) = f(t, x(t), u(t)) \quad \text{System Dynamics} \\
 & h(t, x(t), u(t)) \leq 0 \quad \text{Path Constraint} \\
 & g(t_0, t_F, x(t_0), x(t_F)) \leq 0 \quad \text{Boundary Constraint} \\
 & x_{low} \leq x(t) \leq x_{upp} \quad \text{Path Bound on State} \\
 & u_{low} \leq u(t) \leq u_{upp} \quad \text{Path Bound on Control} \\
 & t_{low} \leq t_0 < t_F \leq t_{upp} \quad \text{Bounds on Initial and Final Times} \\
 & x_{0,low} \leq x(t_0) \leq x_{0,upp} \quad \text{Bound on Initial State} \\
 & x_{F,low} \leq x(t_F) \leq x_{F,upp} \quad \text{Bound on Final State}
 \end{array}$$



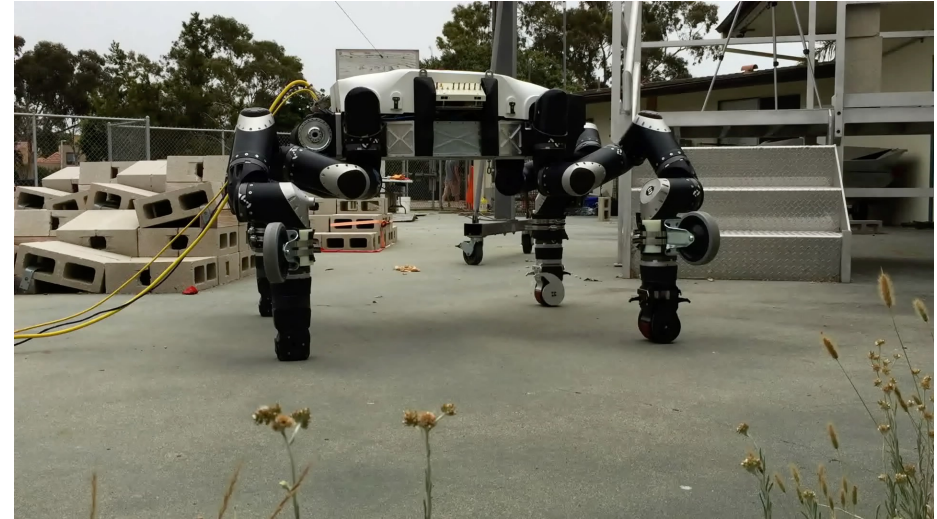
$$\begin{array}{ll}
 \underset{x_k, u_k; k=1 \dots N}{\text{minimize}} & J(x_N) + h \sum_{k=1}^N w(x_k, u_k) \\
 \text{subject to} & d(x_k, u_k, x_{k+1}, u_{k+1}) = 0, \quad k = 1, \dots, N - 1 \\
 & \phi(x_k, u_k) = 0, \quad k = 1, \dots, N \\
 & \psi(x_k, u_k) \geq 0, \quad k = 1, \dots, N
 \end{array}$$

Backward Euler integration:

$$d(x_k, u_k, x_{k+1}, u_{k+1}) := x_{k+1} - (x_k + hf(x_{k+1}, u_{k+1}))$$

# Applying Trajectory Optimization to JPL's RoboSimian

- Four identical limbs, each with 7 degrees of freedom (DOFs)
- Joint motor speeds peak at 1 (rad/s)
  - RoboSimian is stable, but very slow
- Passive single-wheel skate mounted at each forearm
  - Increase efficiency, speed
- 6 actuated DOFs available to set 6-DOF pose of each skate
- No sensing available at skate contact



# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
 subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q)$$

$$+ A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

End Location      Power

$$J = J_F(q_N, \dot{q}_N) + h \sum_{k=1}^N \sqrt{(\dot{q}_k^T U_k)^2 + \epsilon}$$

with

- $U = [0, 0, 0, 0, u_{x_i}, u_{y_i}, u_{z_i}, u_{\theta_i}]^T \in \mathbb{R}^{20}$  for  $i \in \{1, 2, 3, 4\}$
- $\epsilon = 10^{-6}$  is a regularization term to help smooth the cost function
- $h = \Delta t$  is the time interval between time steps
- $J_F(q_N, \dot{q}_N)$  is a distance measure to a set of goal coordinates, such as:

$$J_F = \alpha_x (x_g - x_N)^2 + \alpha_y (y_g - y_N)^2 + \alpha_\theta (\theta_g - \theta_N)^2$$

# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
 subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q)$$

$$+ A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

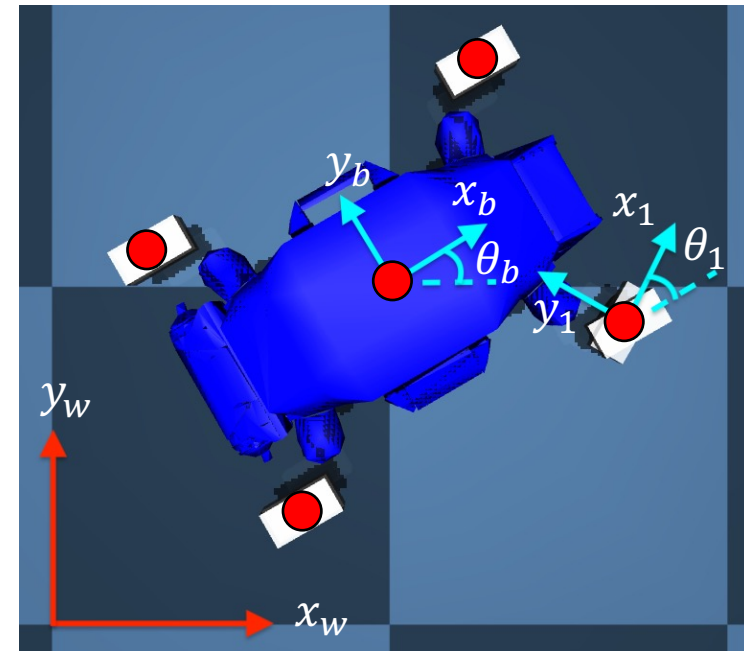
$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

- $q, \dot{q}$  are bounded by ranges for placing each skate with inverse kinematics, as well as by joint limits
- $u$  is bounded explicitly, as well as implicitly by  $\dot{q}$  ranges
- Known contact sequence, where  $F_n \geq 0$  for each skate in contact, and  $F_n = 0$  for skates not in contact



# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q) + A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

$$q_{k+1} = q_k + h\dot{q}_{k+1}$$

$$\dot{q}_{k+1} = \dot{q}_k + h\ddot{q}_{k+1}$$

with

$$\ddot{q}_{k+1} = M_{k+1}^{-1} (B_{k+1} u_{k+1} + F_{n_{k+1}} + J_{k+1}^T F_{fric_{k+1}} - C_{k+1} \dot{q}_{k+1} - G_{k+1})$$

# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
 subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q) + A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

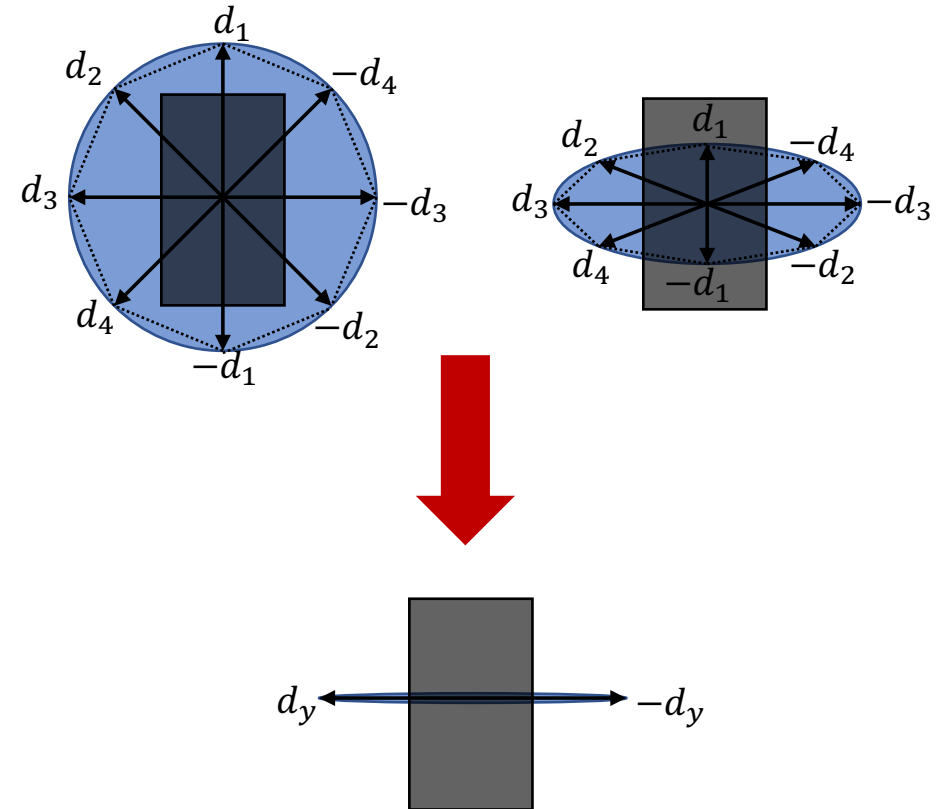
$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

Friction as a Linear Complementarity Problem (LCP)



**G. Bellegarda** and K. Byl. "Trajectory Optimization for a Wheel-Legged System for Dynamic Maneuvers that Allow for Wheel Slip," CDC 2019

**G. Bellegarda** and K. Byl. "Versatile Trajectory Optimization Using a LCP Wheel Model for Dynamic Vehicle Maneuvers," ICRA 2020

# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
 subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q) + A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

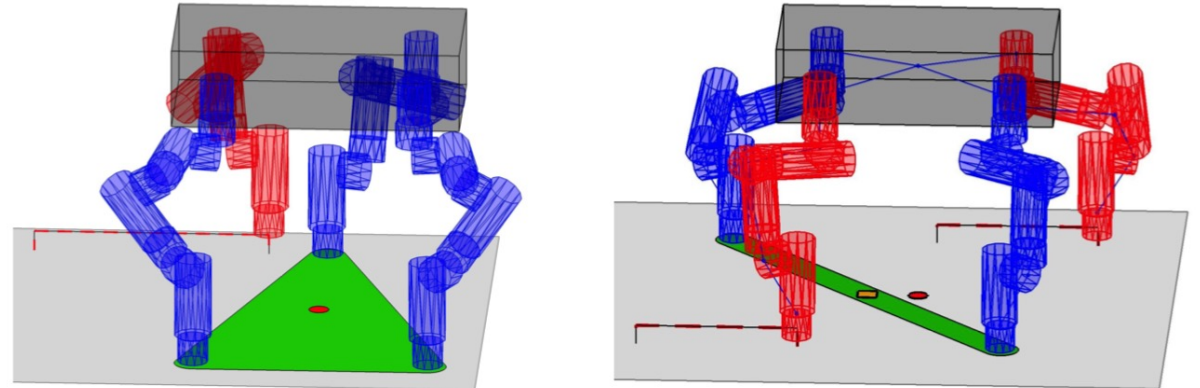
$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

Zero-Moment Point (ZMP): point on ground where the sum of all moments of the active forces is equal to zero



# Trajectory Optimization Framework

find  $q, \dot{q}, u, F_n, F_{fric}$  at discrete timesteps  $k = 1 \dots N$   
 subject to minimize cost  $J$

- State Constraints:

$$\phi(q, \dot{q}, u, F_n) = 0$$

$$\psi(q, \dot{q}, u, F_n) \geq 0$$

- Dynamics Constraints:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q)$$

$$+ A(q)^T \lambda = B(q)u + F_n + J(q)^T F_{fric}$$

- Friction Constraints (for each contact  $i$ ):

$$\gamma e + D^T v^{k+1} \geq 0, \quad \beta \geq 0$$

$$\mu F_{n_i} - e^T \beta \geq 0, \quad \gamma \geq 0$$

$$(\gamma e + D^T v^{k+1})^T \beta = 0$$

$$(\mu F_{n_i} - e^T \beta) \gamma = 0$$

$$F_{fric_i} = D \beta$$

- ZMP Constraints (RoboSimian Only)

- Gait Constraints (RoboSimian Only)

Static Trotting Gait

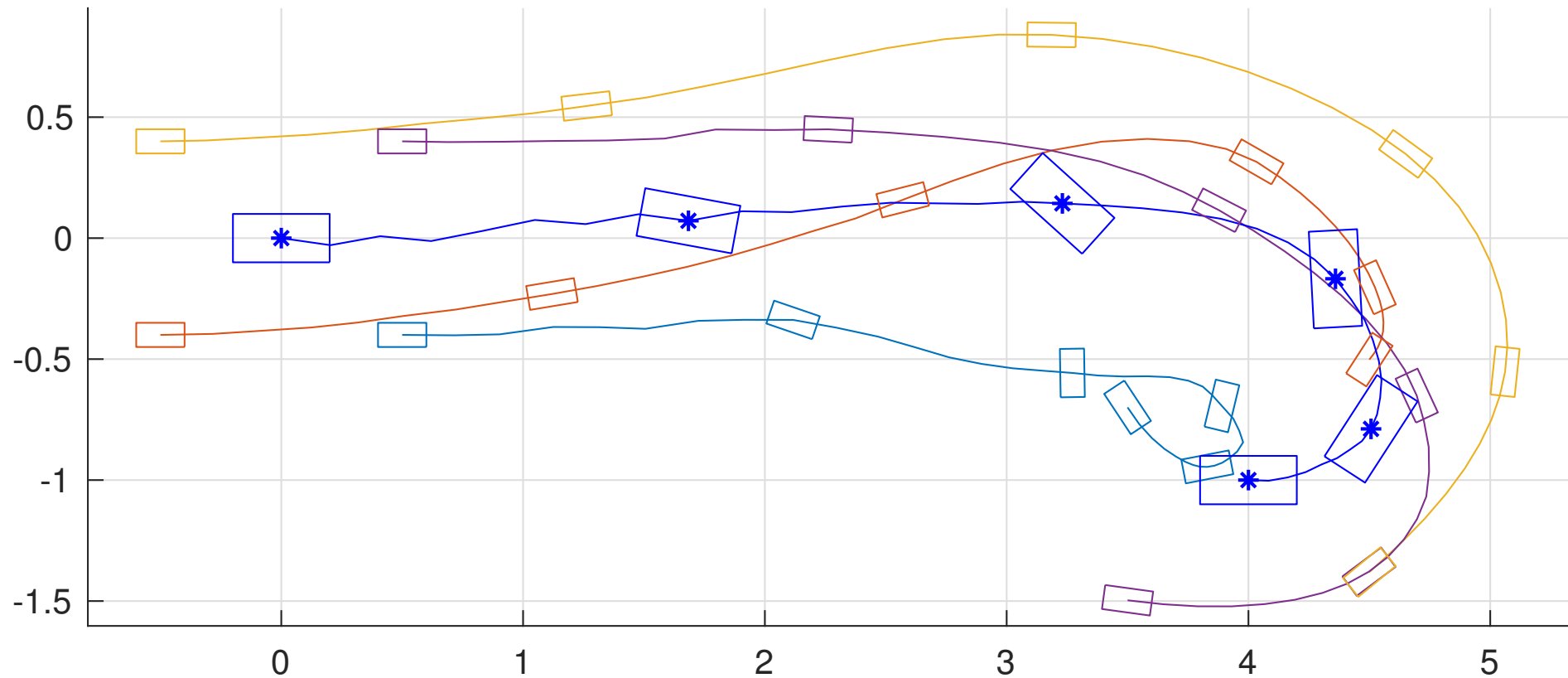
LH	Blue	Blue	White	Blue	White	Blue
LF	Blue	White	Blue	White	Blue	Blue
RF	Blue	Blue	White	Blue	White	Blue
RH	Blue	White	Blue	White	Blue	Blue

Static Walking Gait

LH	Blue	Blue	Blue	White	Blue	Blue	Blue	White
LF	Blue	Blue	Blue	Blue	White	Blue	Blue	Blue
RF	Blue	Blue	White	Blue	Blue	Blue	White	Blue
RH	Blue	White	Blue	Blue	Blue	White	Blue	Blue

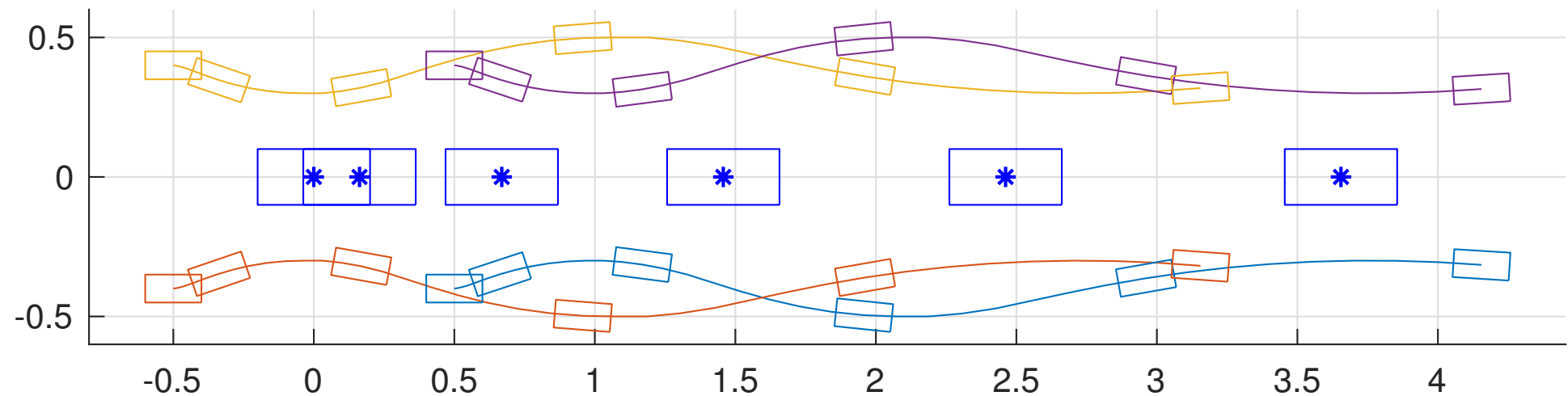
# Dynamic Parking Maneuver

- Goal: “park” at  $(x_b, y_b, \theta_b) = (4, -1, -\pi)$ , initialized at  $(0,0,0)$  with  $\dot{x}_b = 4$  (m/s)

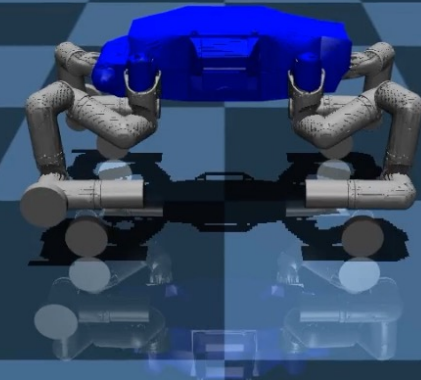
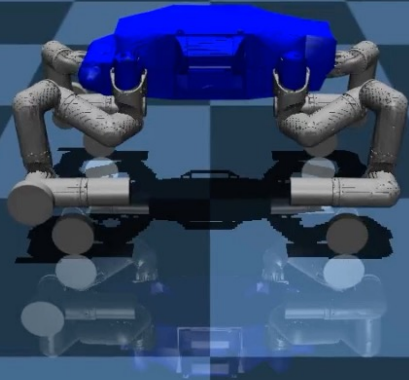


# Energy-Efficient Forward Locomotion

- Maximize final body position in the x-direction, while minimizing energy use



# Hybrid Rolling-Walking/Trotting



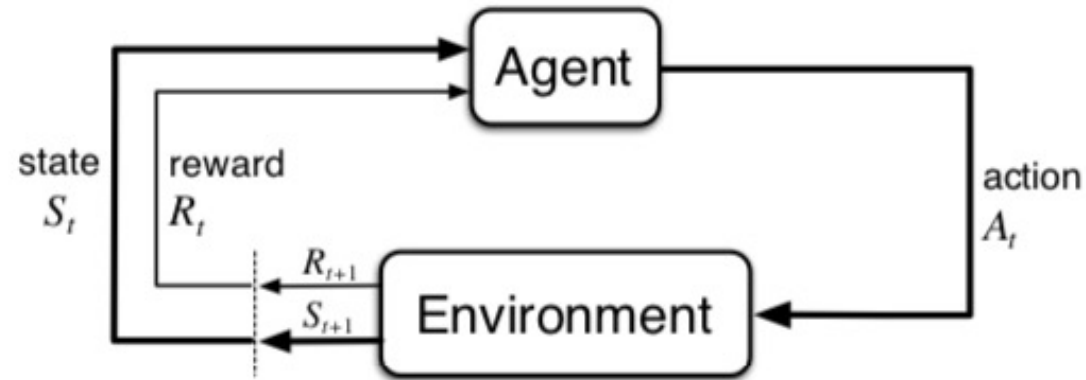
# TO/MPC Limitations

- Very complex systems → cannot solve in real-time!
  - Model simplifications
- Uncertainty?
  - Typically assumptions on friction, actuation, etc.
  - Environment may be different than planned!
- Biases?
  - Gait pattern, constraints, cost function
- Integrate vision, hierarchical control, etc.

# Reinforcement Learning as a Markov Decision Process (MDP)

An MDP is defined by:

- Set of states  $S$
- Set of actions  $A$
- Transition function  $P(s' | s, a)$
- Reward function  $R(s, a, s')$
- Start state  $s_0$
- Discount factor  $\gamma$
- Horizon  $H$



- Return over a trajectory  $\tau = (s_0, a_0, s_1, a_1, \dots)$

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$$

- Policy  $\pi(a_t | s_t)$  maps from states  $s_t$  to actions  $a_t$  (Goal: find policy maximizing above return)
- Value function:  $V^\pi(s) = \mathbb{E}_{\tau \sim \pi}[R(\tau) | s_0 = s]$
- Action-value function:  $Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[R(\tau) | s_0 = s, a_0 = a]$
- Advantage function:  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$

# Many Existing Tools for Reinforcement Learning

- RL algorithm implementations
  - rsl\_rl [https://github.com/leggedrobotics/rsl\\_rl](https://github.com/leggedrobotics/rsl_rl)
  - stable-baselines3 <https://github.com/DLR-RM/stable-baselines3>
  - ray[rllib] <https://github.com/ray-project/ray>
  - skrl <https://skrl.readthedocs.io/en/latest/>
  - rl\_games [https://github.com/Denys88/rl\\_games](https://github.com/Denys88/rl_games)
  - ... many others!
- Physics simulators
  - Isaac Gym <https://developer.nvidia.com/isaac-gym>
  - Isaac Lab <https://github.com/isaac-sim/IsaacLab>
  - pybullet <https://github.com/bulletphysics/bullet3>
  - MuJoCo <https://mujoco.org>
  - RaiSim <https://raisim.com>
  - ... and others!

# Reinforcement Learning Considerations

## Algorithm

- On/off policy
- Hyperparameters
- Network architecture
- Random seeds/trials

...implementation  
dependent!

## MDP Design Decisions

- Observation space
- Action space
- Reward function

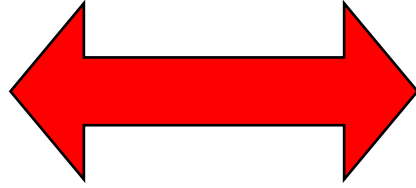
## Environment Parameters

- Simulator dynamics
- Control gains –  
joint/Cartesian
- Control/environment  
time step
- Noise, latency

# Learning in Simulation → Deploy in real world (Why?)

## Real World

- Slow
- Unsafe
- Expensive
- Human Supervision



## Simulation

- Fast
- Safe
- Cheap
- Scalable

**BUT: “Sim-to-real” gap → Simulation != Reality**

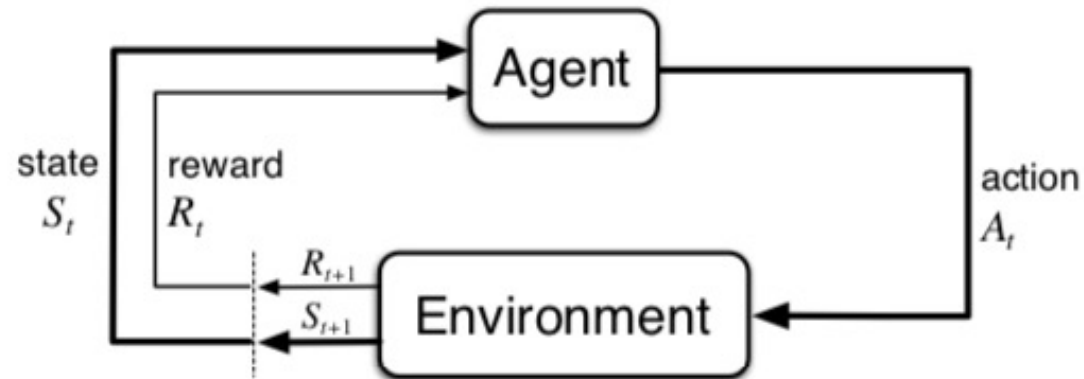
- Unmodeled dynamics
- Wrong simulation parameters
- Inaccurate contact models
- Latency
- Actuator dynamics
- Sensor noise
- Stochastic real world
- ...

# General RL Techniques for Locomotion

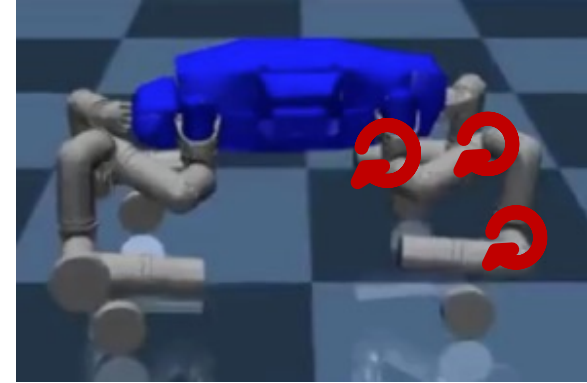
- System identification
  - Add real-world identified actuator dynamics to simulator
- Domain randomization
  - Randomize geometries, mass, inertia, friction, perturbations, etc.
- Curriculum learning
  - To learn a “difficult skill”, first start with an easier version of the task, then make it harder as the agent gets better
- Teacher-Student distillation
  - Train a teacher policy with privileged information (i.e. not known in real world)
  - Train a student policy to imitate the teacher action with supervised learning
- Reward shaping
  - Change or tune the reward function over time to prioritize different task-related components

# State/Action/Reward Space?

$s_t$  ? i.e.  
-body (z, r, p, y)  
-body velocities  
-joint states



$r_t$  ? i.e.  
-body linear velocity  
-energy penalty



$a_t$  ?  
-motor positions/torques

# Training in Joint Space (PPO)

Action Space:  $a_t = \tau_{1\dots N}$

Policy after training 1 million time steps in joint space,  
rewarding forward velocity



Humanoid:  
27 DoFs, 21 Actuators.



N. Heess et al. Emergence of Locomotion Behaviours in Rich Environments. *arXiv:1707.02286*, 2017

J. Schulman et al. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017

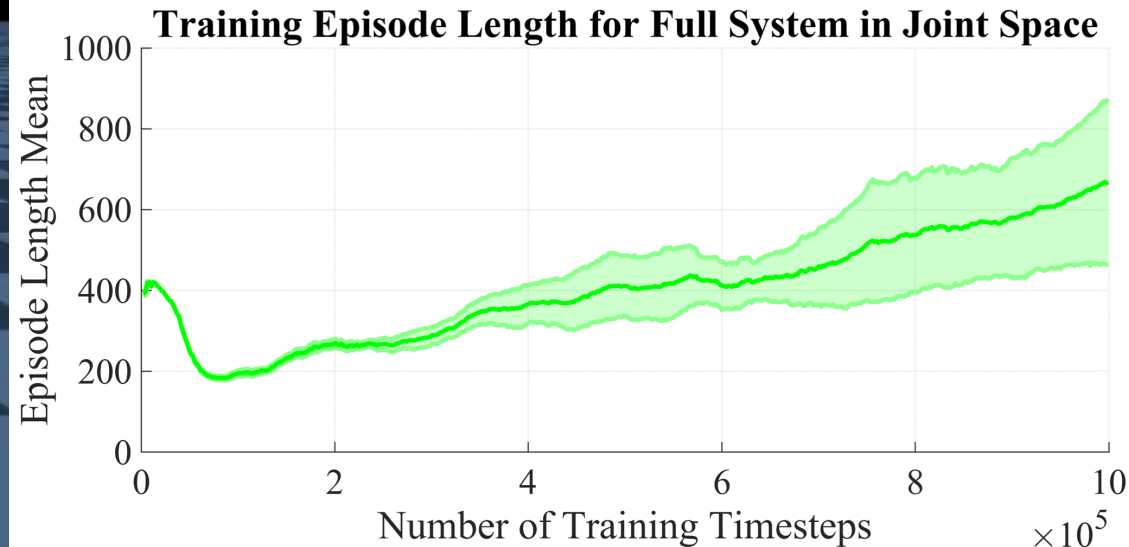
**G. Bellegarda** and K. Byl. “Training in Task Space to Speed Up and Guide Reinforcement Learning,” IROS 2019

# Training in Joint Space (PPO)

Action Space:  $a_t = \tau_{1\dots N}$

- Agent is essentially forced to learn forward and inverse kinematics!

Policy after training 1 million time steps in joint space,  
rewarding forward velocity



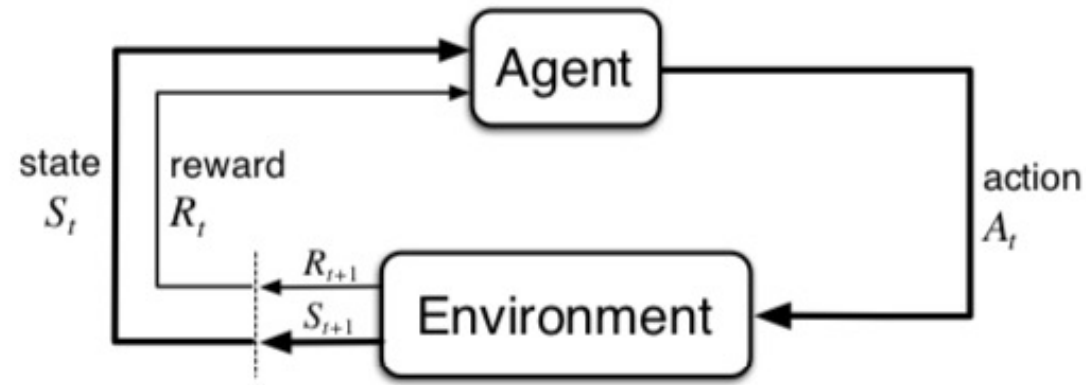
N. Heess et al. Emergence of Locomotion Behaviours in RichEnvironments. *arXiv:1707.02286*, 2017

J. Schulman et al. Proximal policy optimization algorithms. *arXiv:1707.06347*, 2017

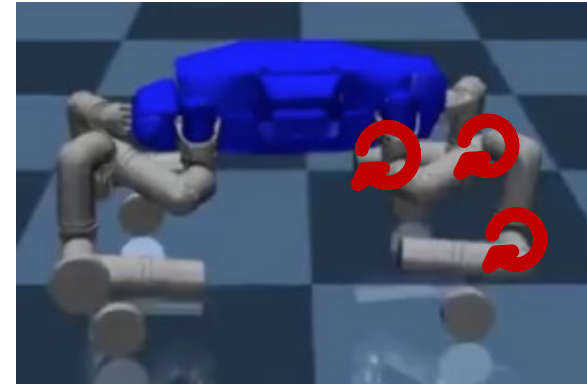
G. Bellegarda and K. Byl. "Training in Task Space to Speed Up and Guide Reinforcement Learning," IROS 2019

# State/Action/Reward Space?

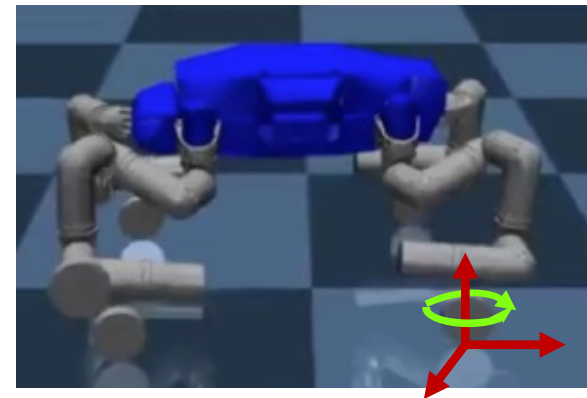
$s_t$  ? i.e.  
-body (z, r, p, y)  
-body velocities  
-joint states

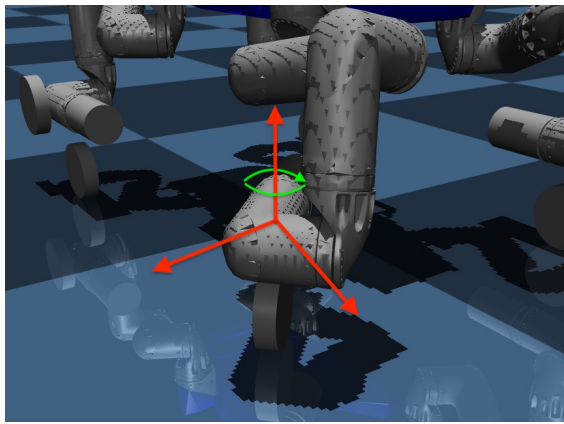


$r_t$  ? i.e.  
-body linear velocity  
-energy penalty



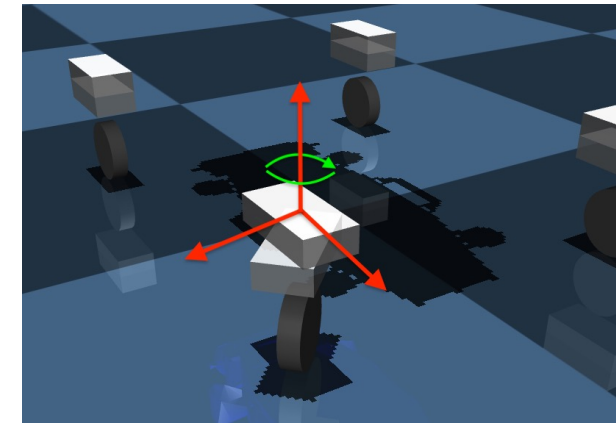
$a_t$  ?  
-motor positions/torques  
-Cartesian coordinates



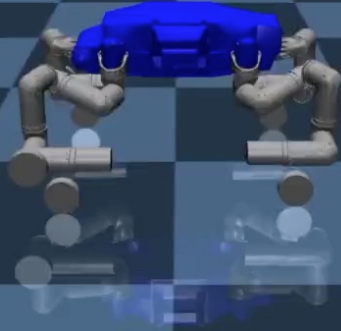


# Training in Task Space (PPO)

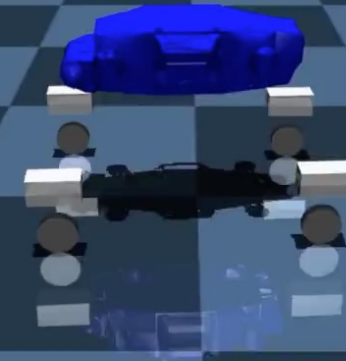
Action Space:  $a_t = [y_{ee_i}, \phi_{ee_i}]$  for  $i \in \{1,2,3,4\}$



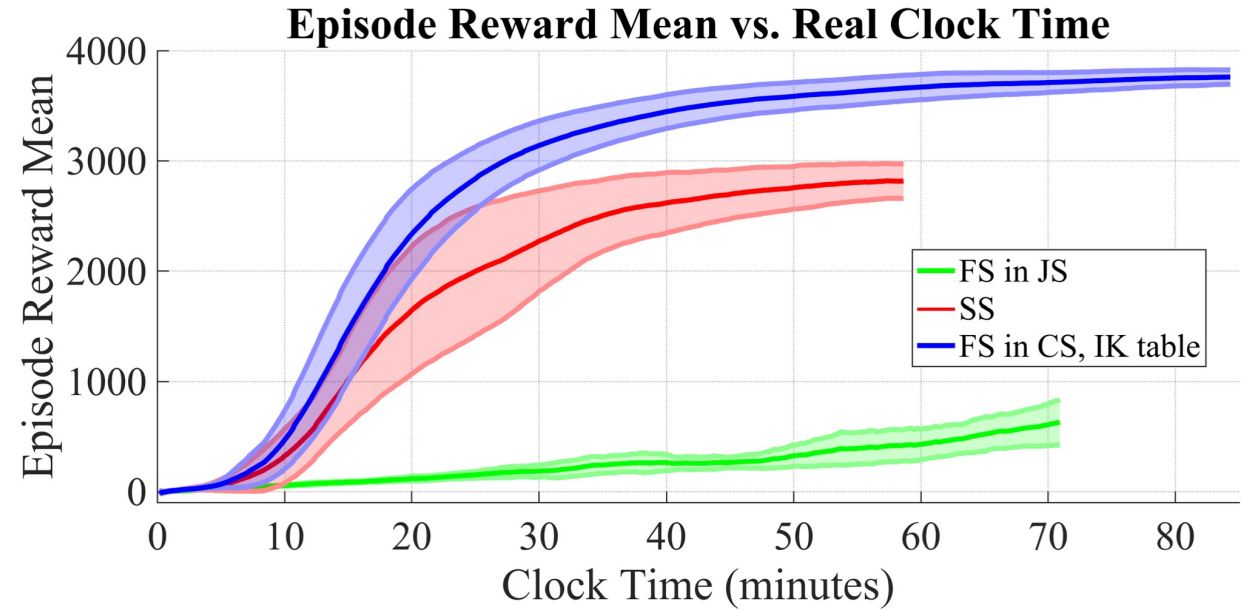
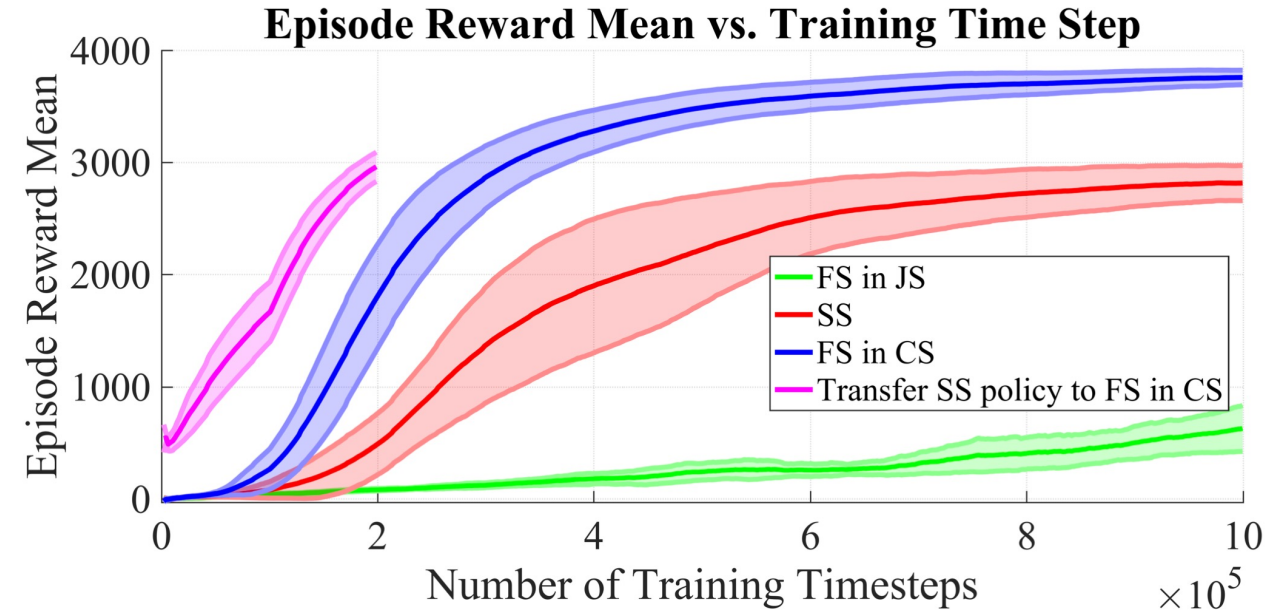
Policy after training 1 million time steps in Cartesian space,  
rewarding forward velocity



Policy after training 1 million time steps in Cartesian space,  
rewarding forward velocity



# Training in Joint Space vs. Task Space

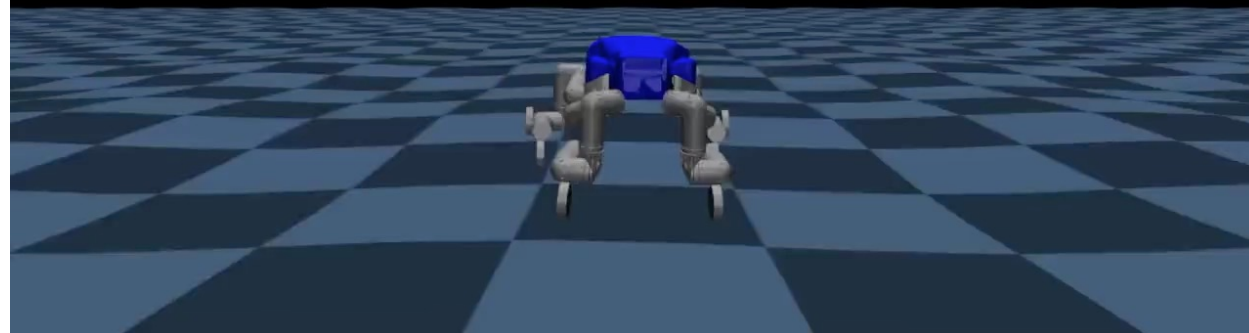


# Hand-designed Trajectory vs. Trained Policy over uneven terrain with varying coefficients of friction

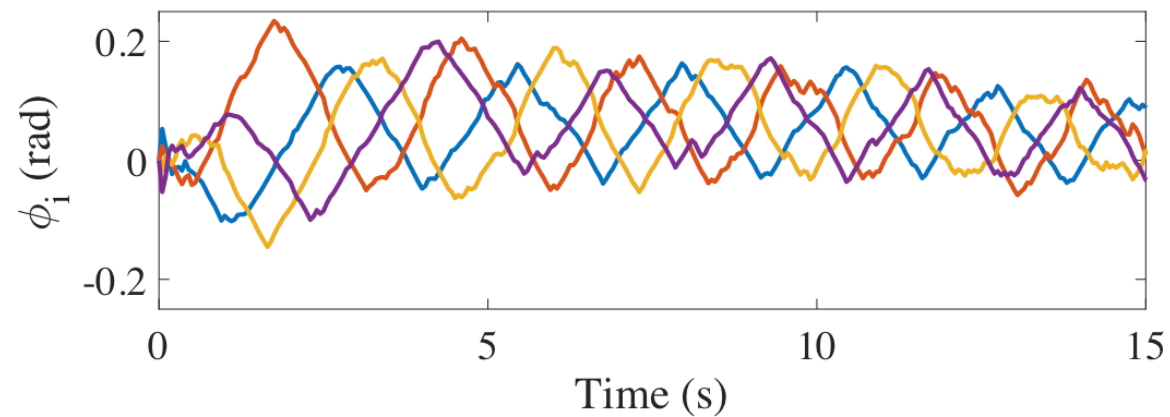
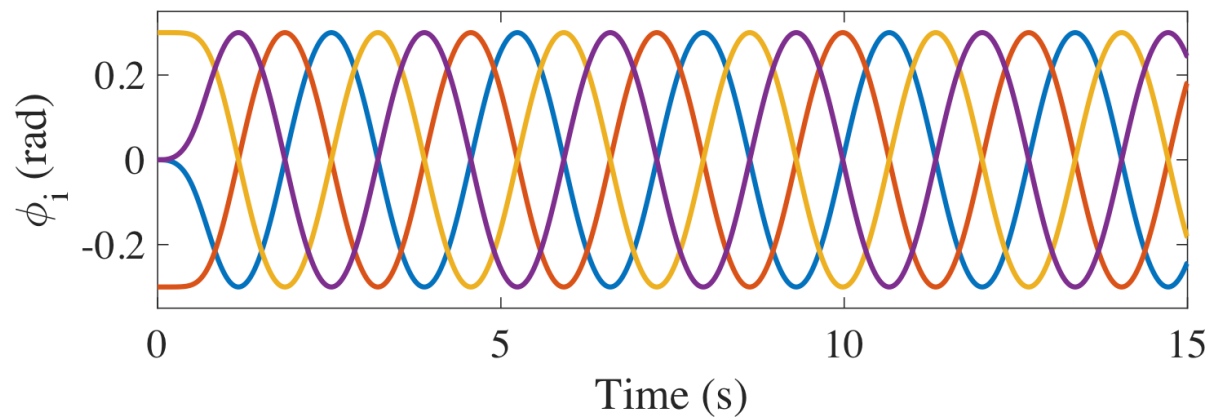
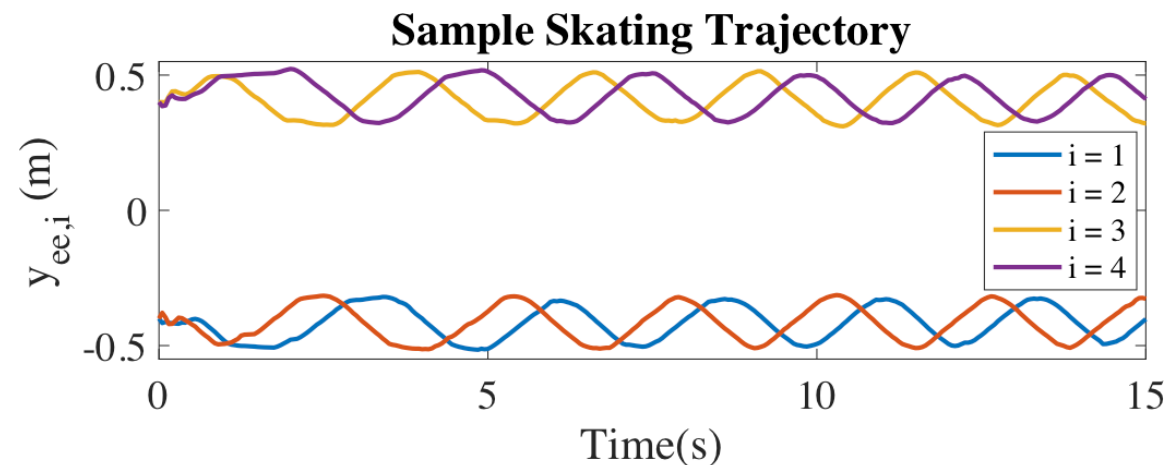
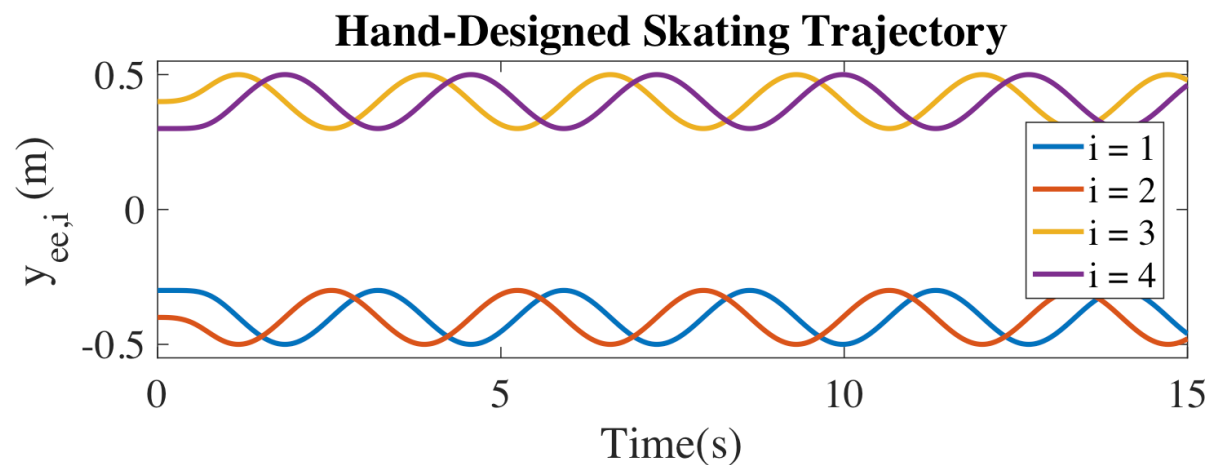
**Hand-designed trajectory to locomote Robosimian 5 meters forwards to (5,0)  
executed over uneven terrain with varying coefficients of friction**



**Policy trained 2 million time steps to locomote Robosimian to (5,0)  
executed over uneven terrain with varying coefficients of friction**

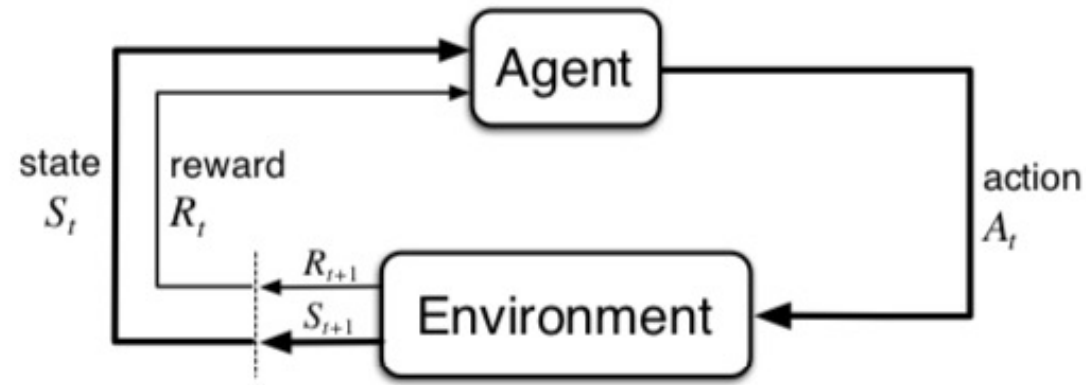


# Hand-designed Trajectory vs. Trained Policy

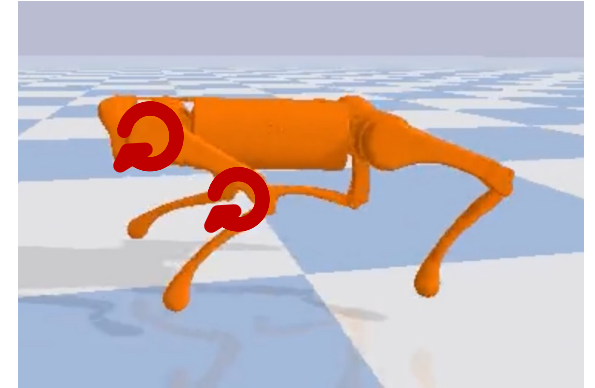


# State/Action/Reward Space: A1

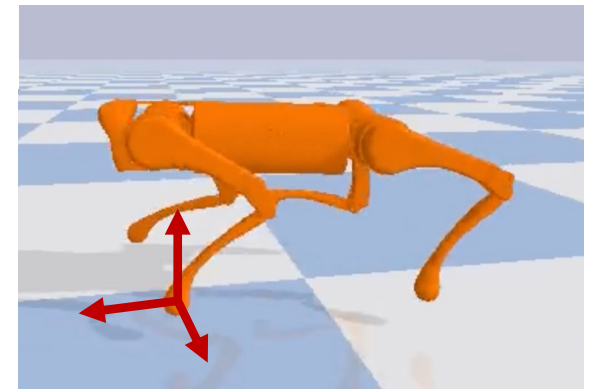
$s_t$  ? i.e.  
-body (z, r, p, y)  
-body velocities  
-joint states



$r_t$  ? i.e.  
-body linear velocity  
-energy penalty

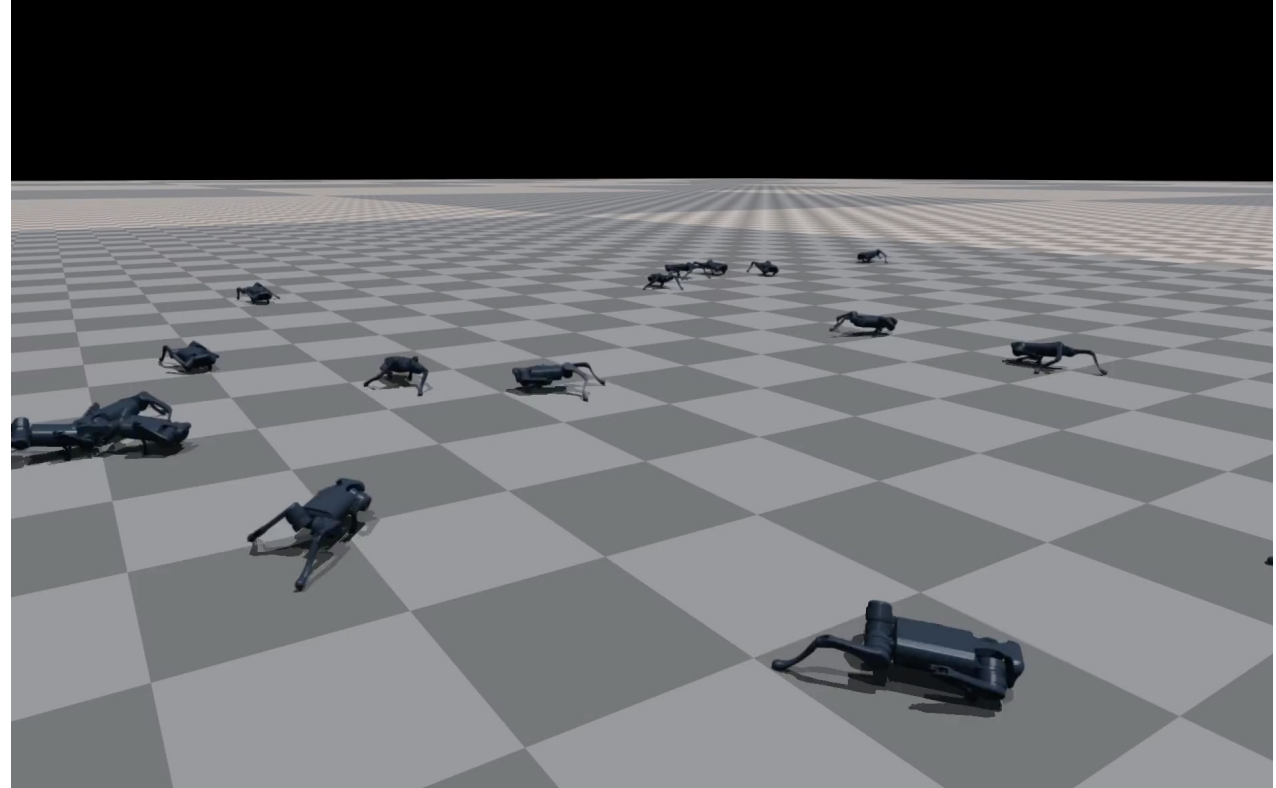
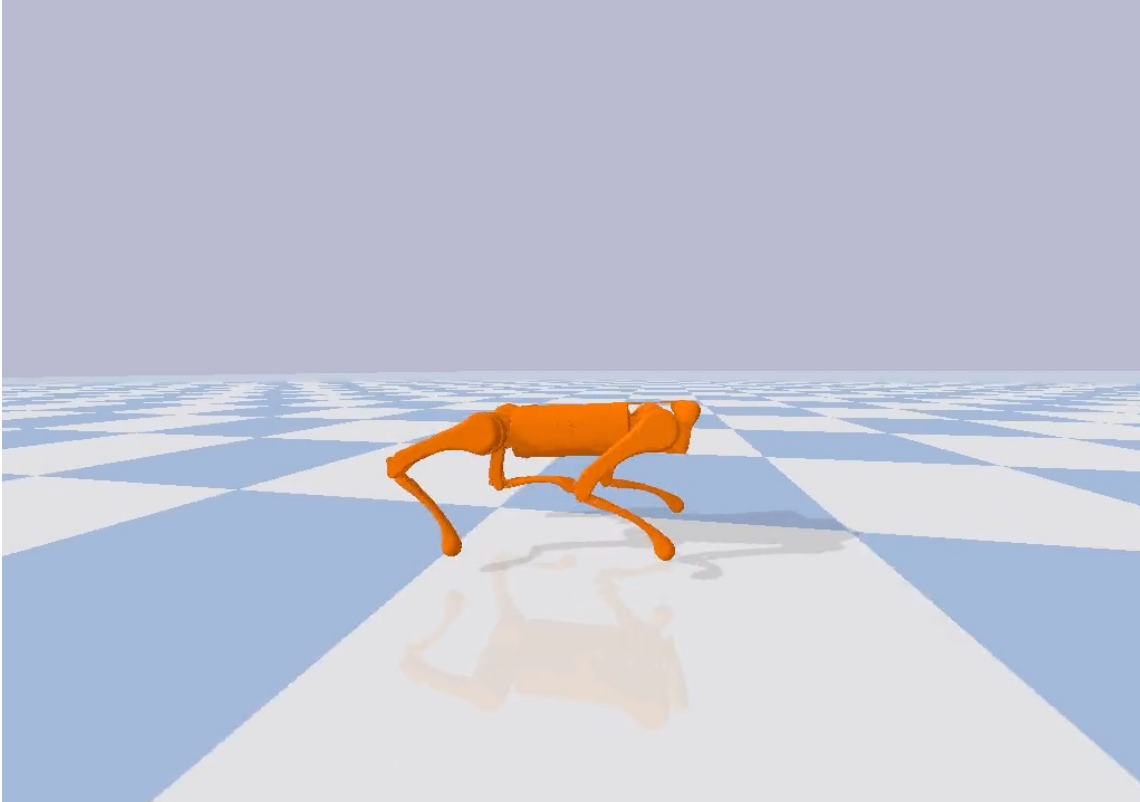


$a_t$  ?  
-motor positions/torques  
-Cartesian PD



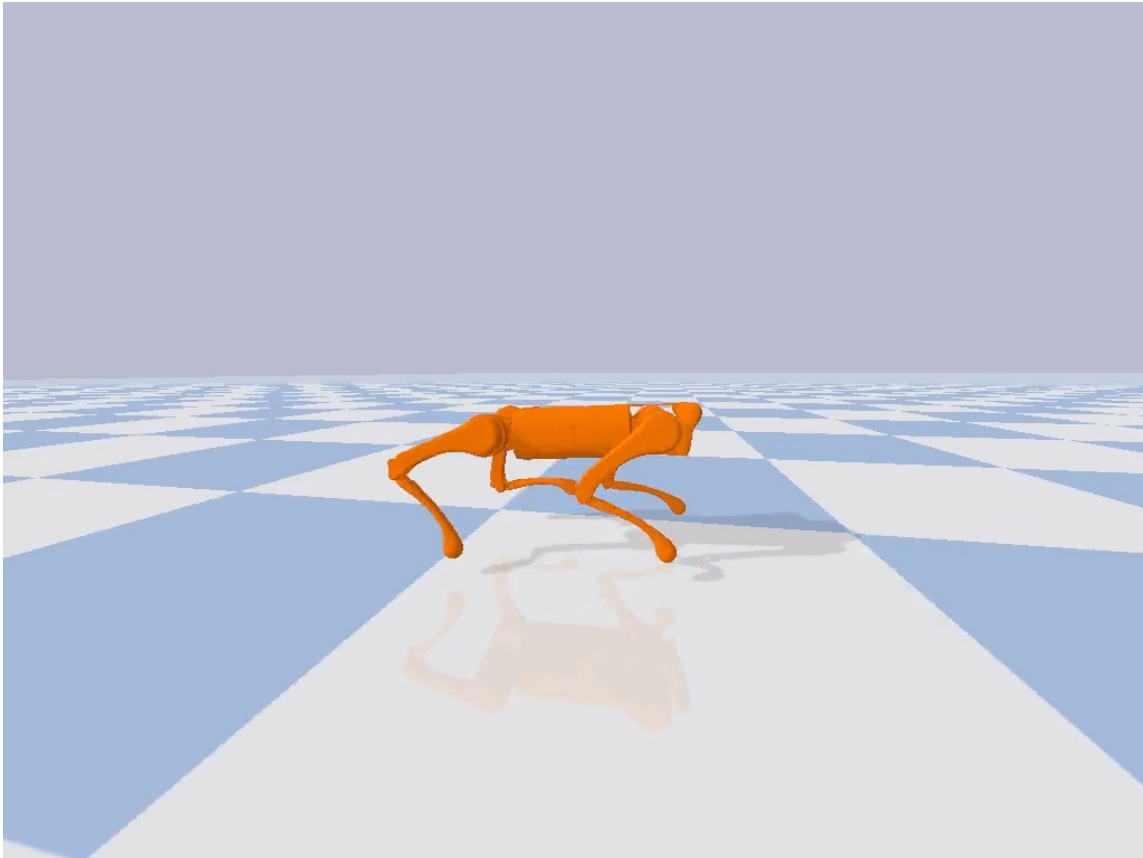
# Learning in Joint Space (PPO)

Action Space:  $a_t = q_{1...N}$

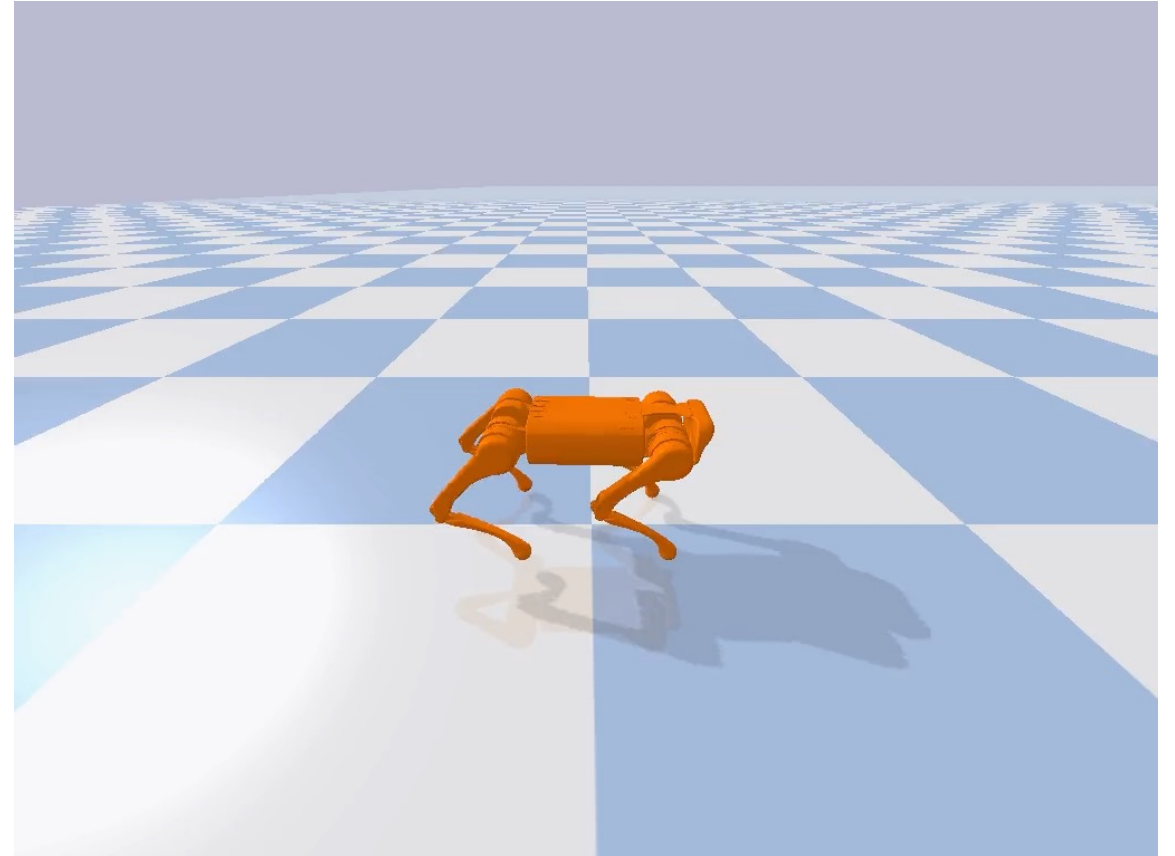


# Position Control vs. Cartesian PD Control (PPO)

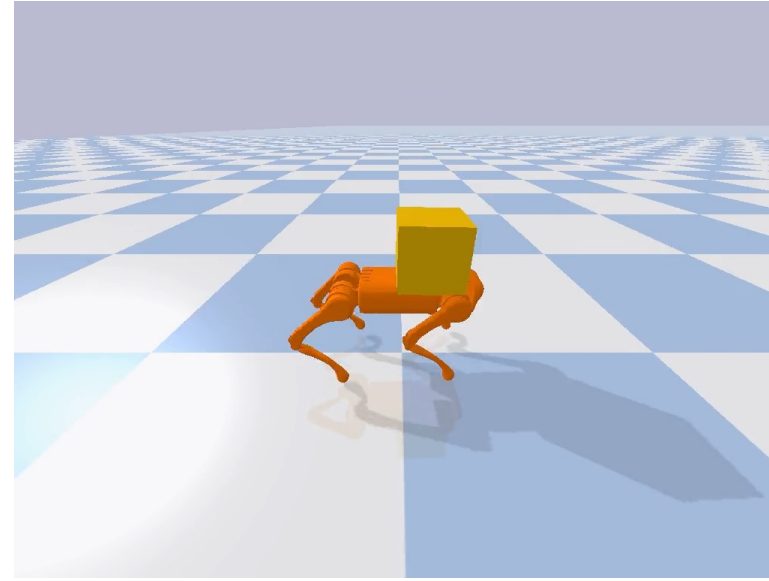
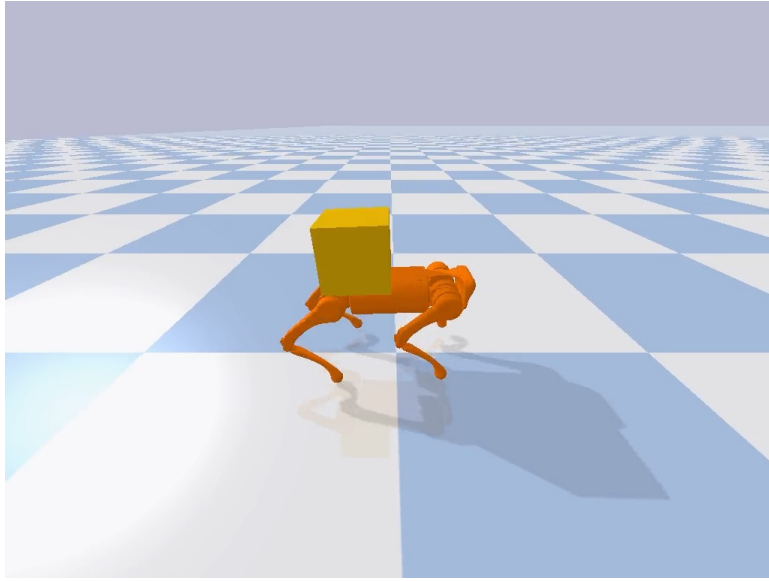
Action Space:  $a_t = q_{1\dots N}$



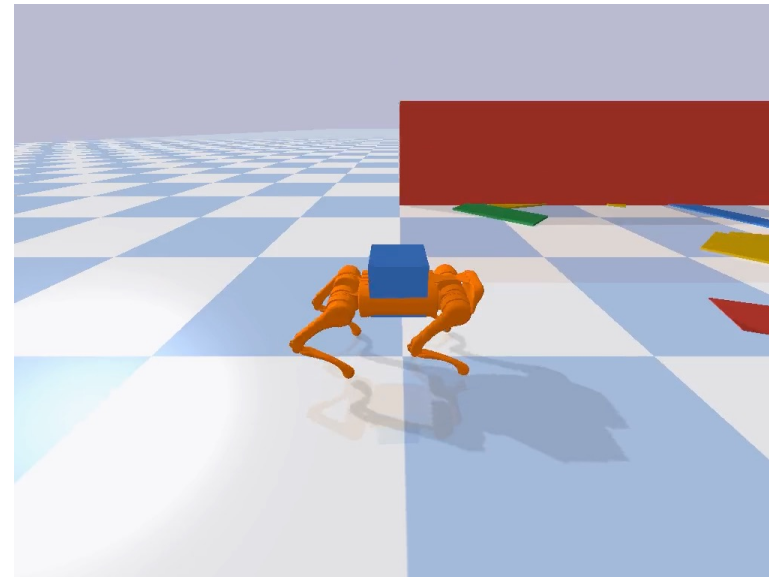
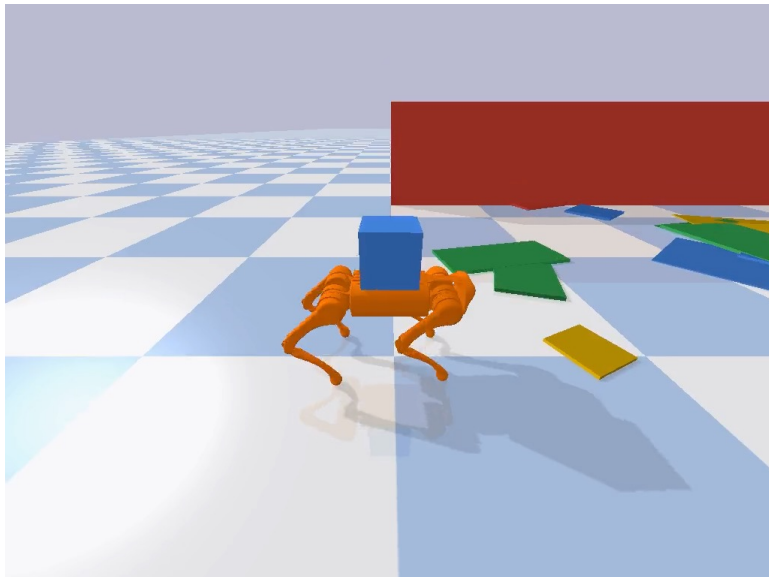
Action Space:  $a_t = [x_{ee_i}, y_{ee_i}, z_{ee_i}]$



Robust to disturbances (10kg load, and 20% body mass/inertia variability)

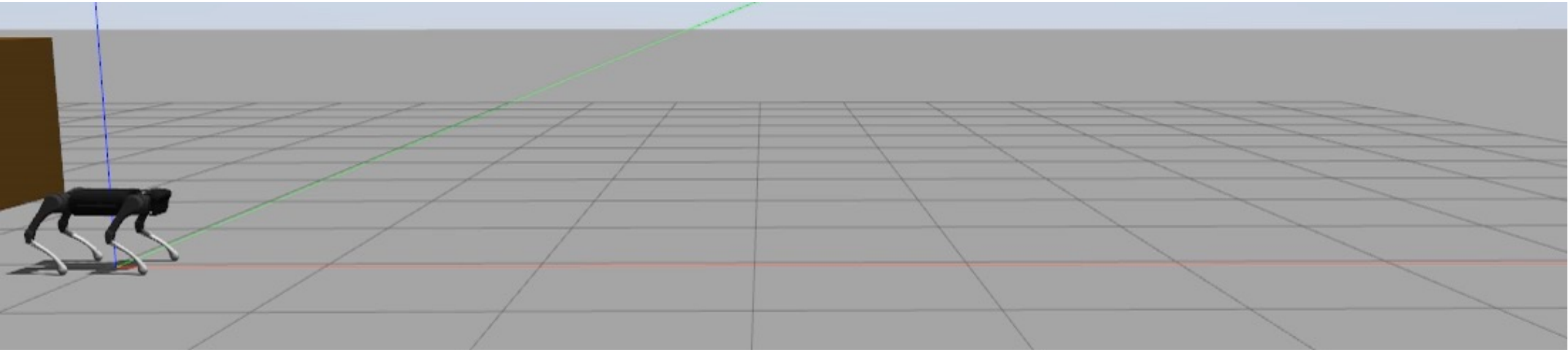


Robust to disturbances (6kg load, rough terrain up to 0.04m in height)

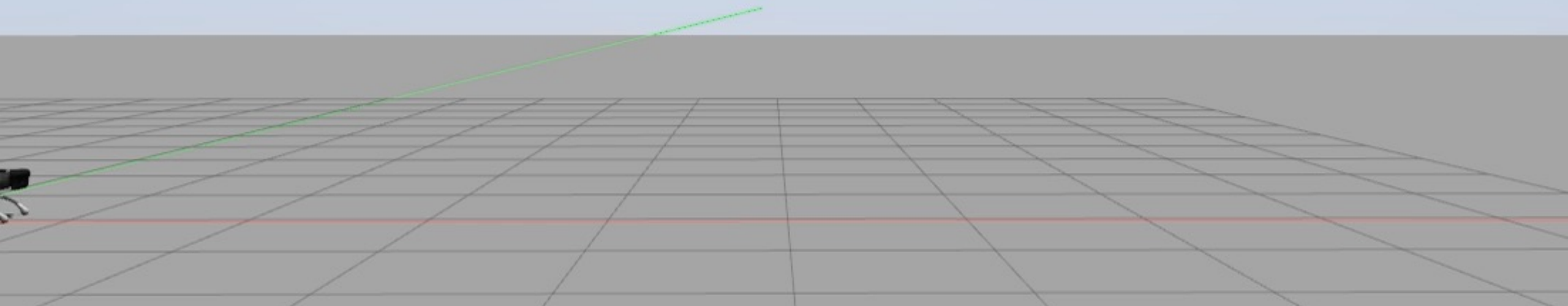


# Sim-to-sim Transfer

Running at 4 m/s



Running at 3.5 m/s with a 10kg load



# Sim-to-Real Additions

- Reward function terms:

Name	Formula	Weight
Velocity reward	$1 -  v_{base} - v_{des} $	0.1
Feet swing reward	$\sum_{i=0}^3 (t_i^k - t_i^{k-1} - 0.5)$	0.2
Energy penalty	$\int_t^{t+1}  \boldsymbol{\tau} \cdot \dot{\boldsymbol{q}}  dt$	-0.008
Orientation penalty	$\ \boldsymbol{\omega} - (0, 0, 0, 1)\ $	-0.1
Lateral drift penalty	$ y $	-0.1
Height penalty	$ z - 0.3 $	-0.1

- Dynamics randomization:
  - Terrain randomization, observation noise

Parameter	Lower Bound	Upper Bound
Mass (each body link)	80%	120%
Added mass	0 kg	5 kg
Added mass base offset	$[-0.15, -0.05, -0.05] m$	$[0.15, 0.05, 0.05] m$
Coefficient of friction	0.5	1

# Experimental Results

0kg load



3kg load



5kg load



5kg load



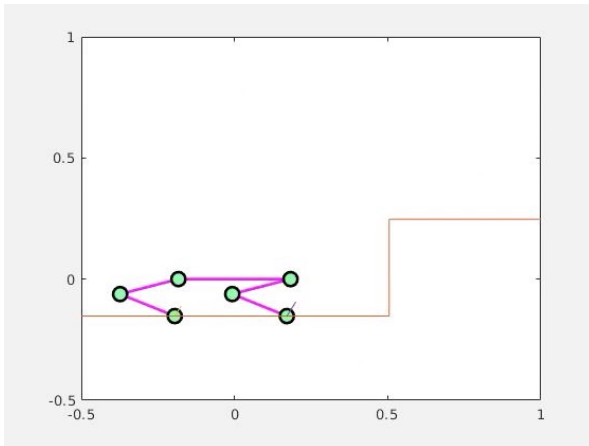
# Running Over Different Terrains



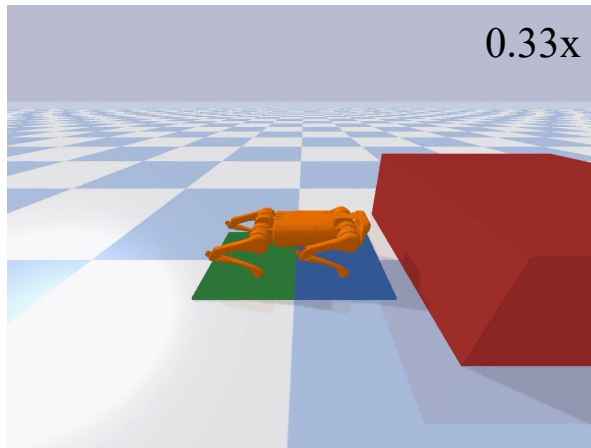
# How can we combine trajectory optimization and deep learning?

- Use learning to improve trajectory optimization

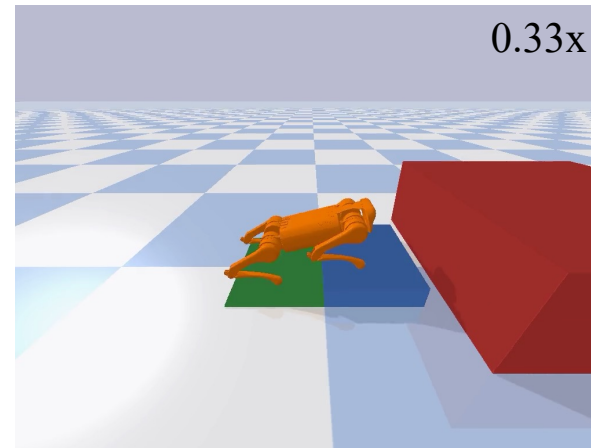
## Feedforward trajectory optimization



Nonlinear trajectory optimization generates jumping motions

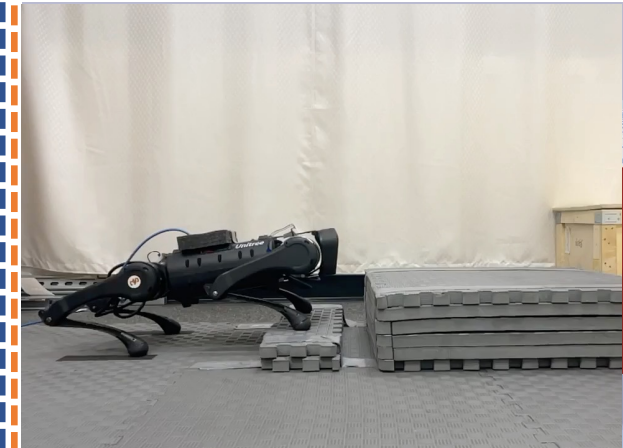


Feedforward execution



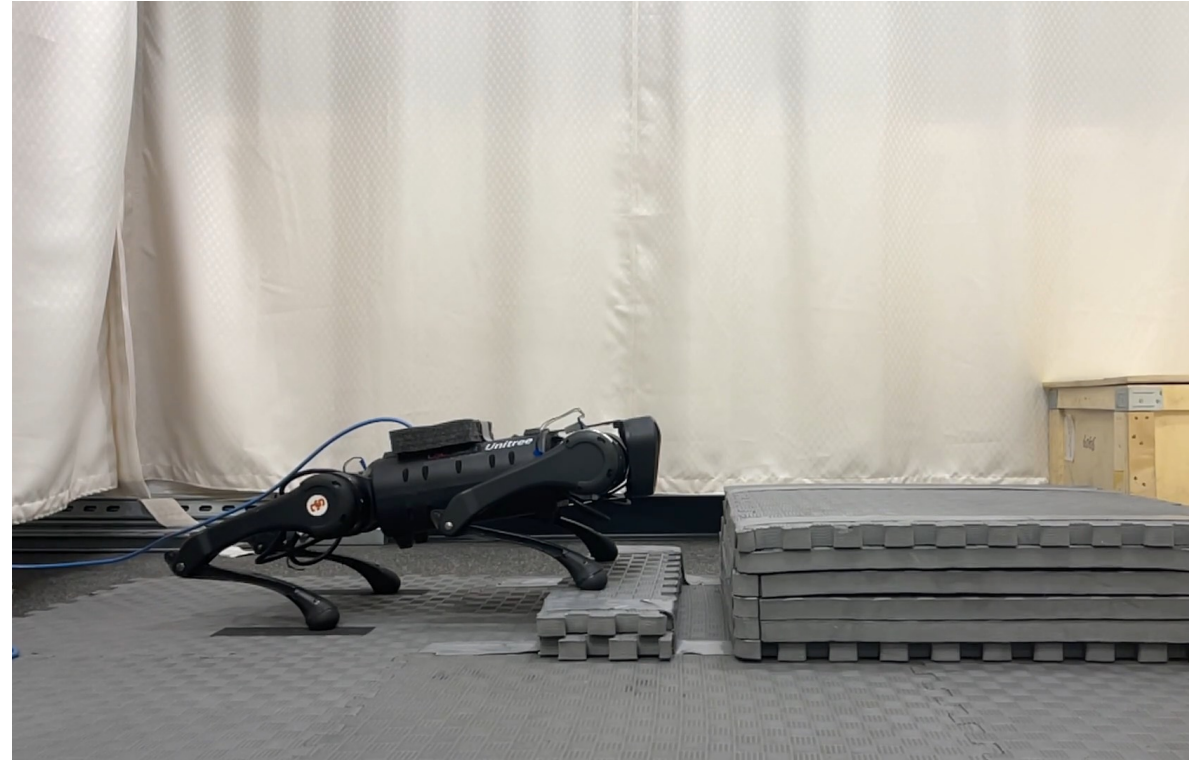
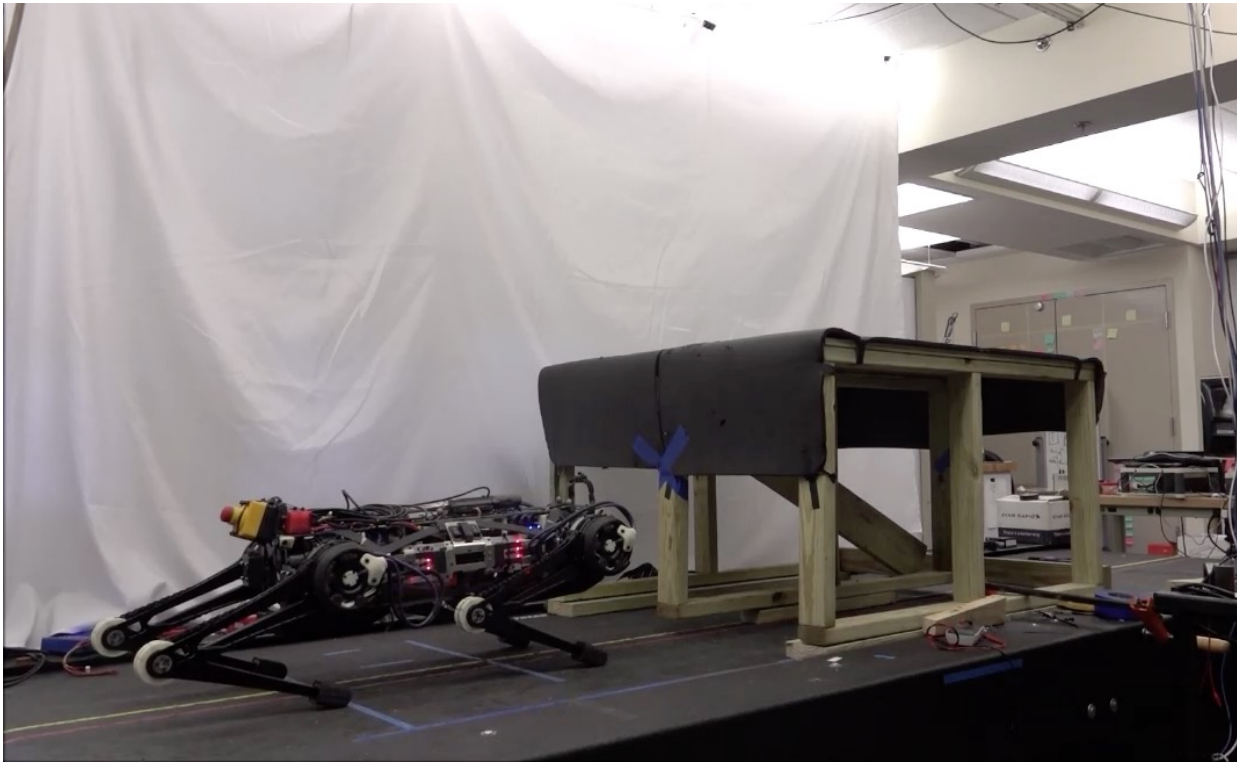
Feedforward execution **with disturbances** – doesn't work!

## TO + Learning



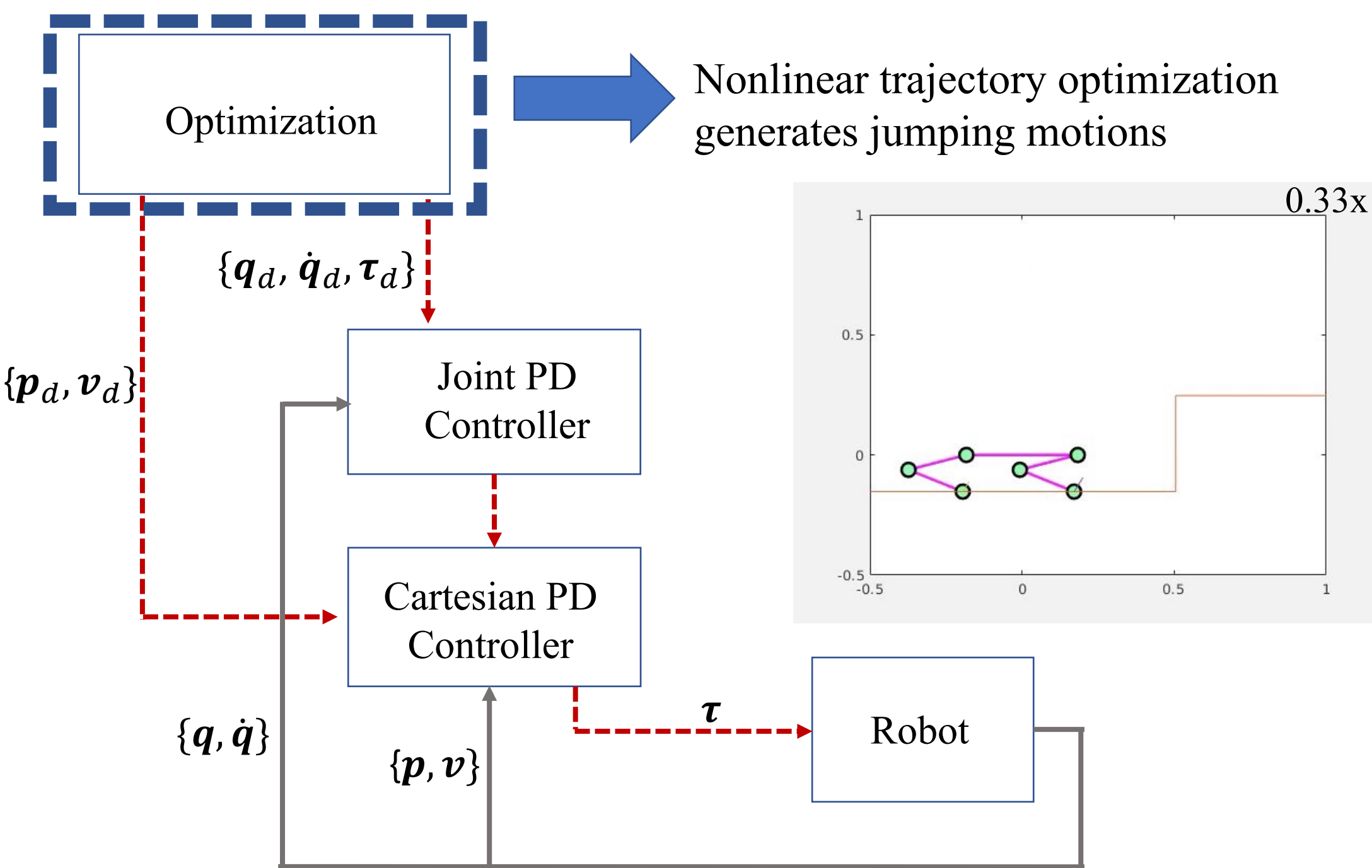
**Our method:**  
Trajectory Optimization + DRL

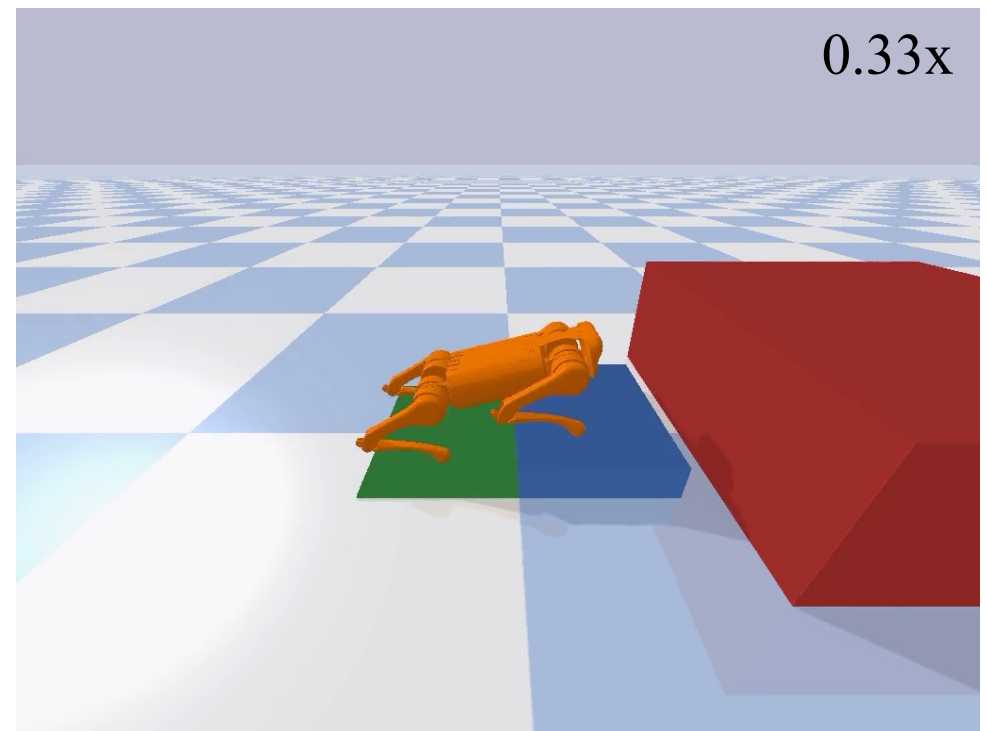
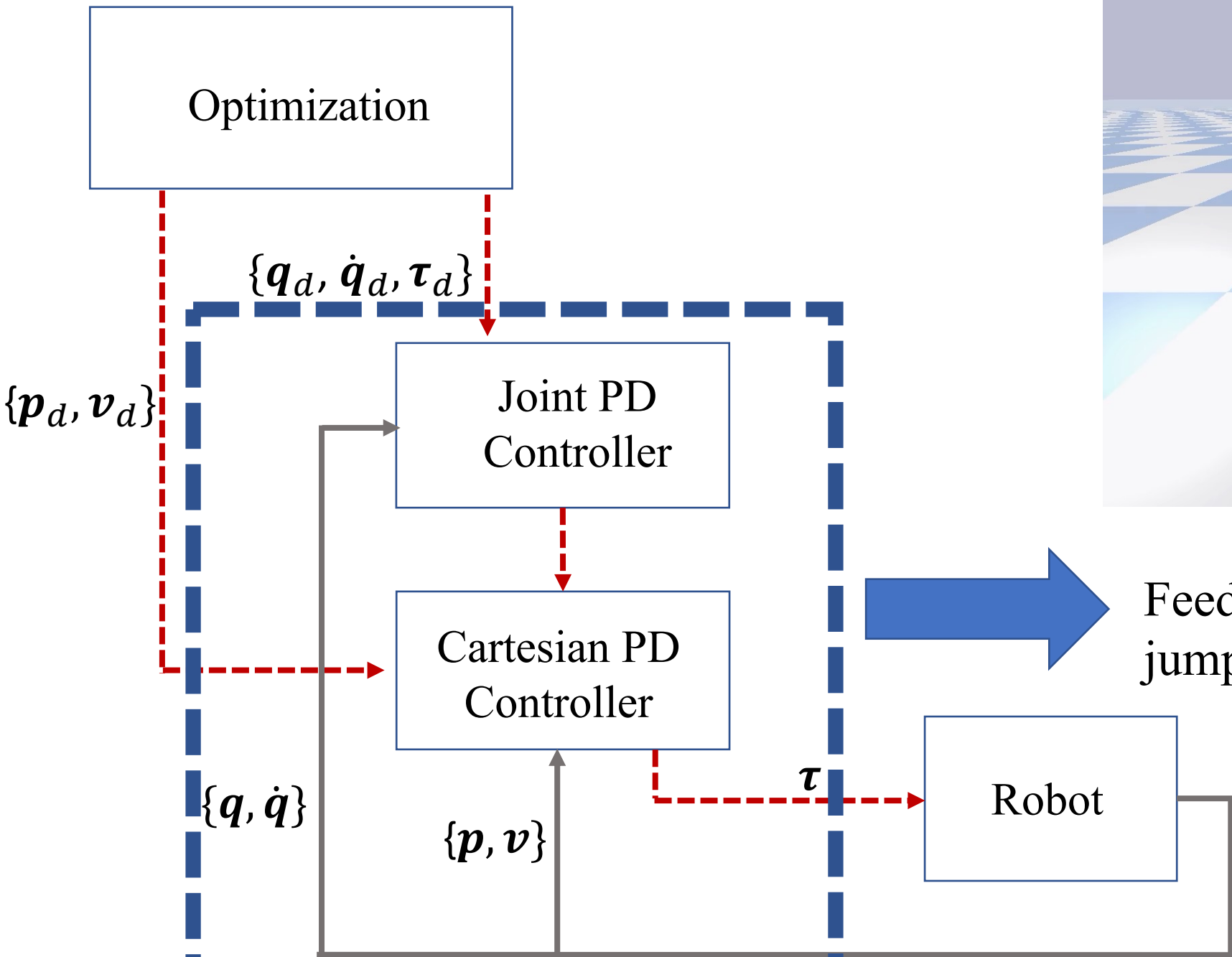
# Jumping Task (Robust?)



Q. Nguyen, M. Powell, B. Katz, J. D. Carlo, S. Kim. *Optimized Jumping on the MIT Cheetah 3 Robot*. ICRA 2019

G. Bellegarda, C. Nguyen, Q. Nguyen. "Robust Quadruped Jumping via Deep Reinforcement Learning," *Robotics and Autonomous Systems* 2024.



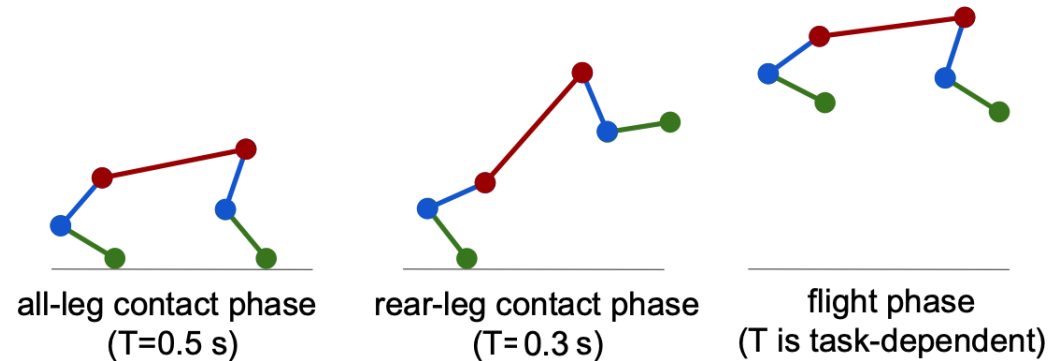


Feedforward controller tracks jumping motions...

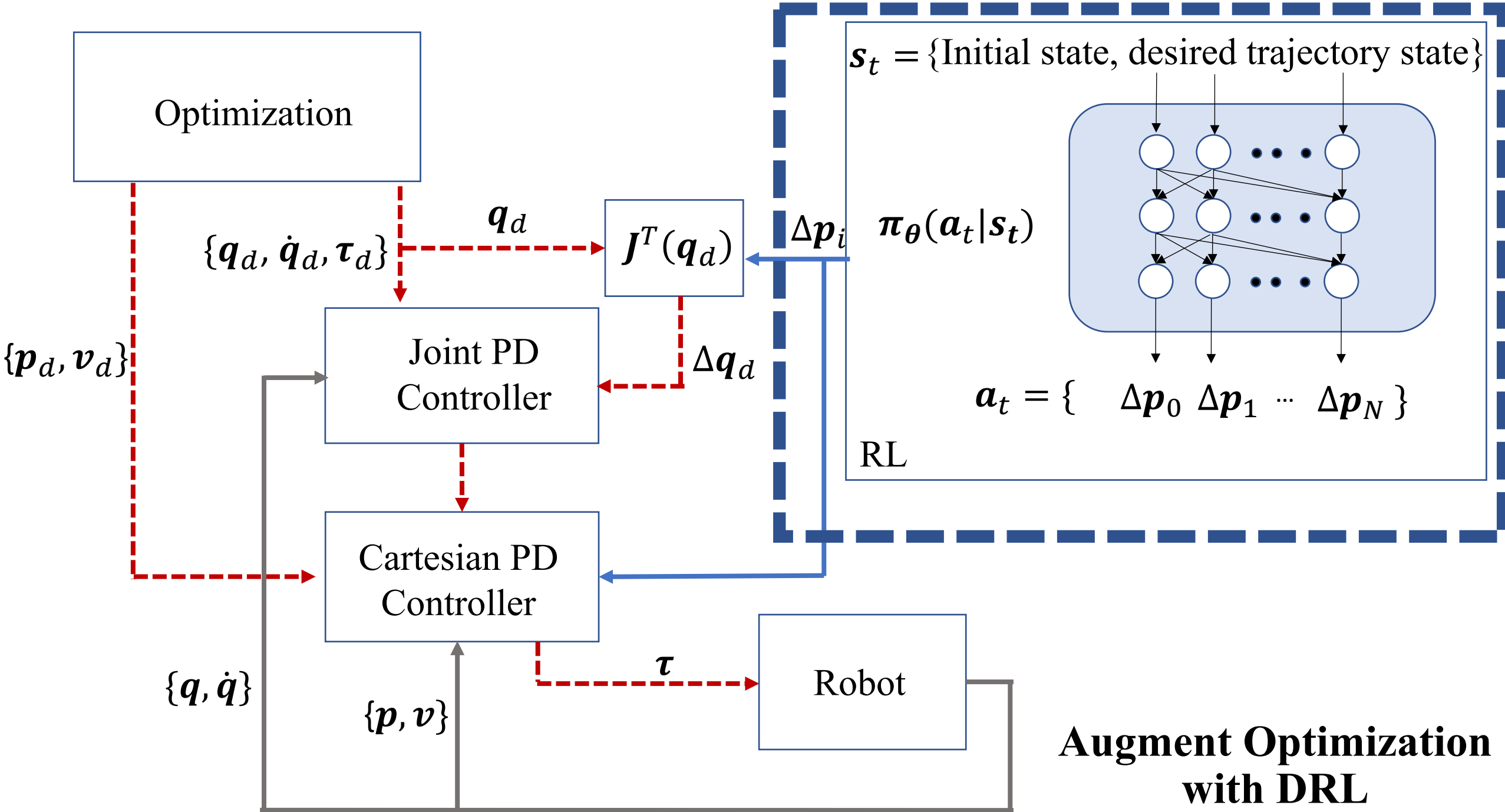
...but only works under ideal conditions!

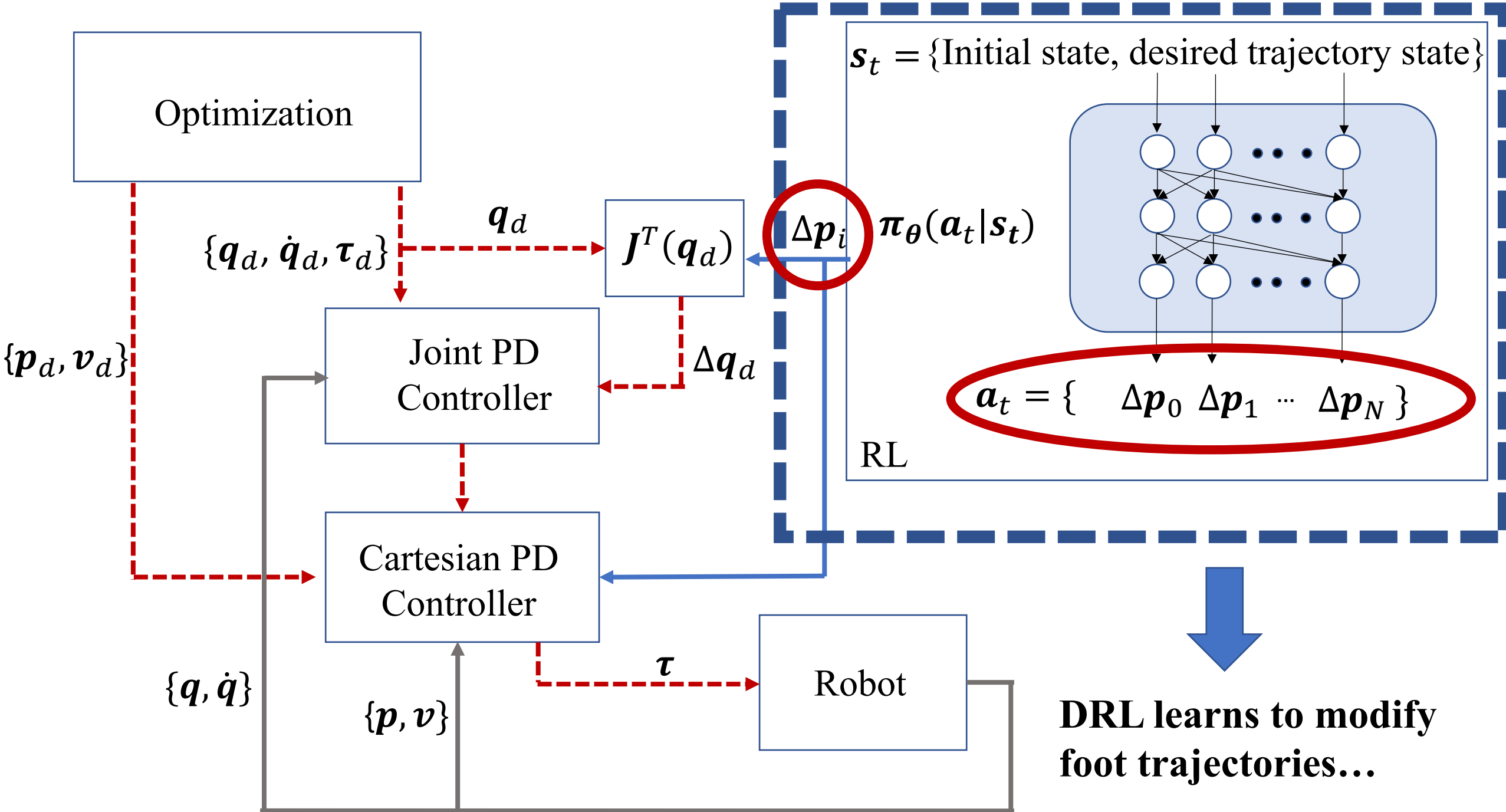
# How to improve jumping with learning?

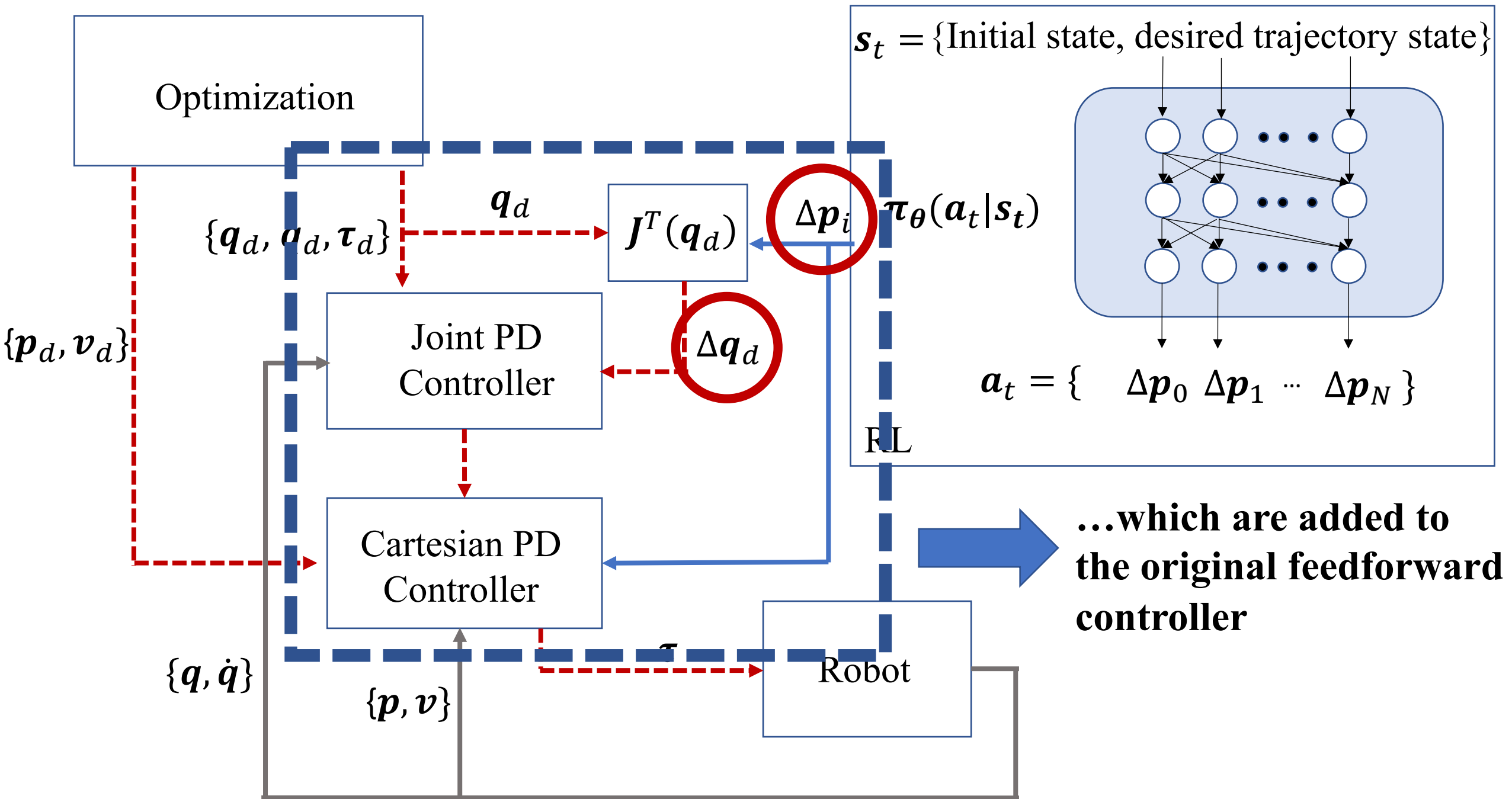
- Jumping involves very specific joint trajectories



- Difficult to learn from scratch..
  - Imitate trajectory?
  - Reward function tuning..
- Already have an approximate solution in ideal case
  - Learn to modify existing trajectories for robustness (residual learning)



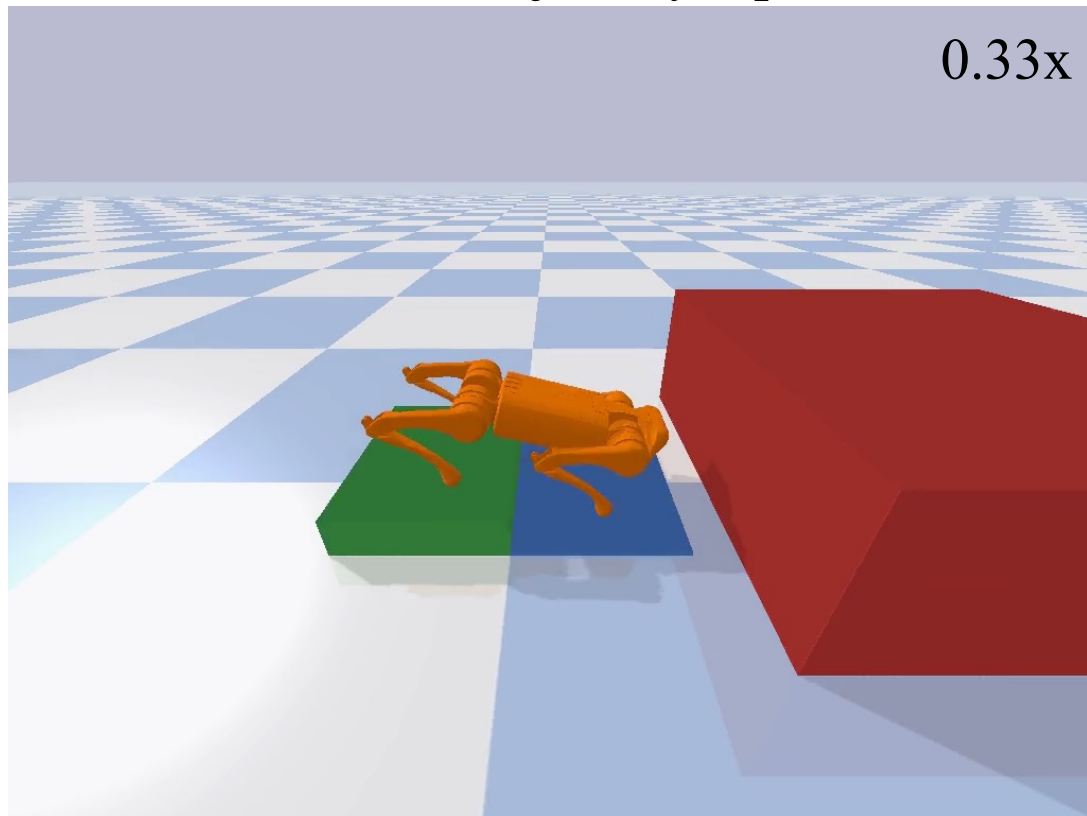




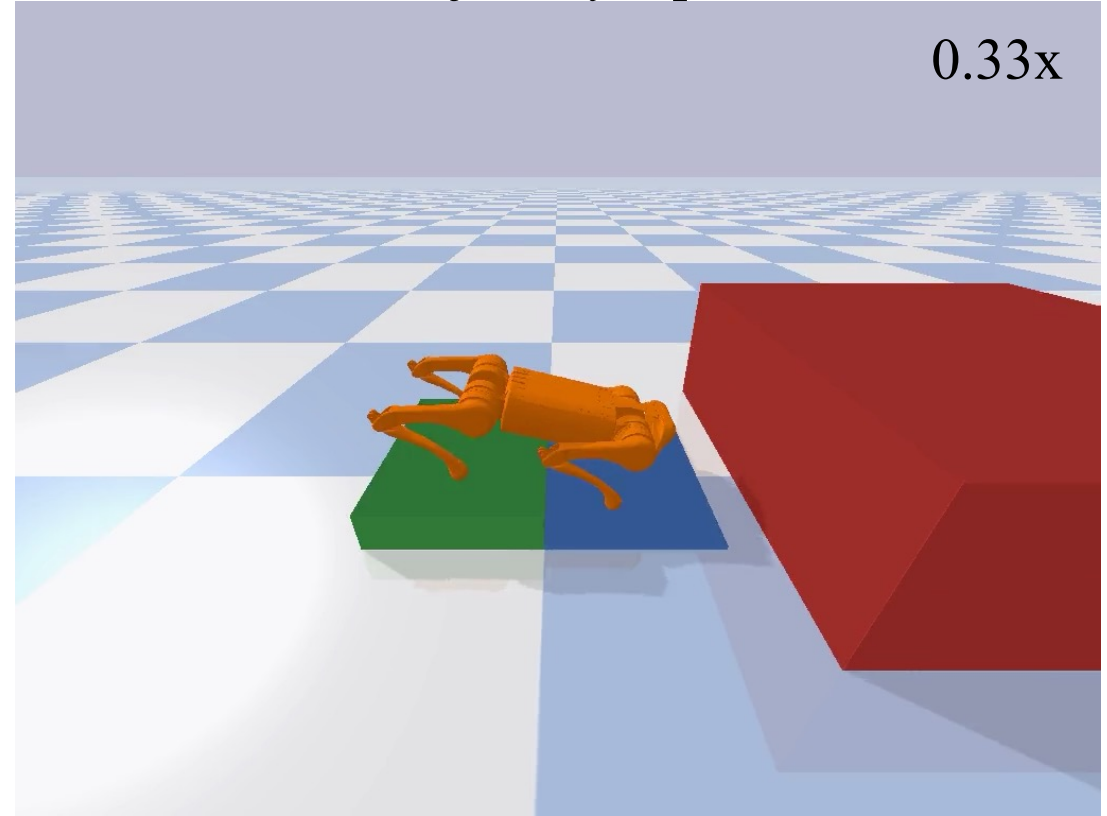
# Example Simulation Result

- Jumping height and distance: **(0.4, 0.7) m**
- Mass and inertia uncertainty: **5%**
- Ground uncertainty: **0.1 m** block under rear feet

**Baseline** - Pure Trajectory Optimization



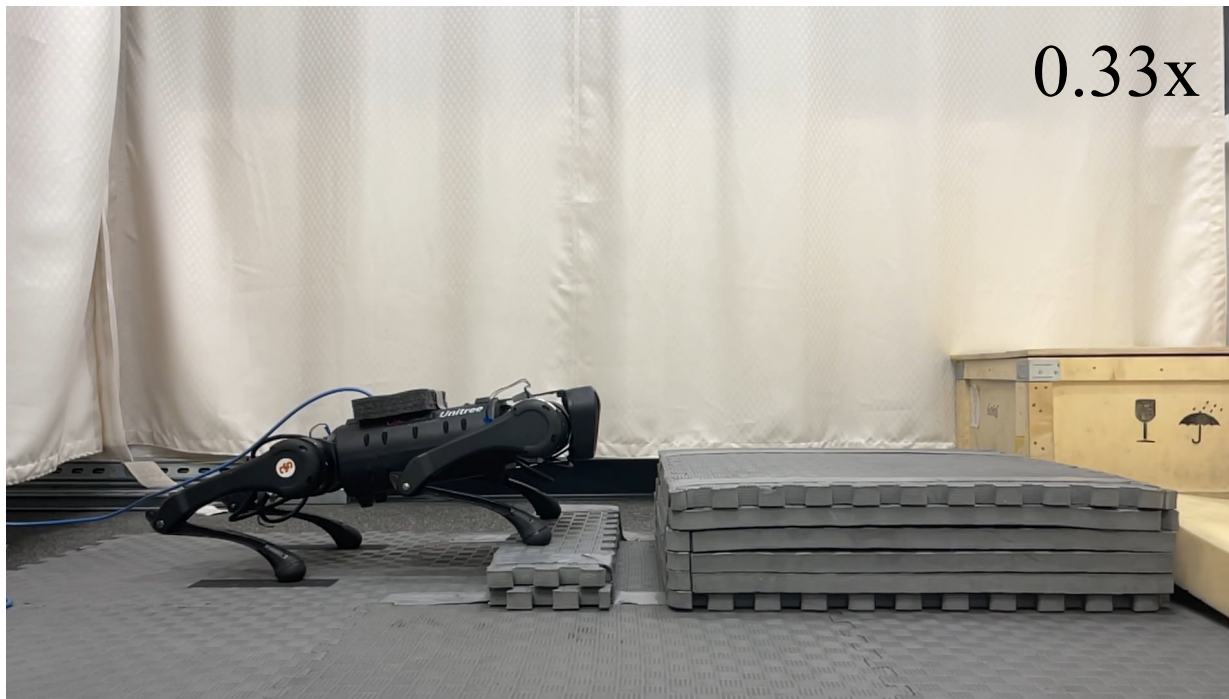
**Our method** – Trajectory Optimization + DRL



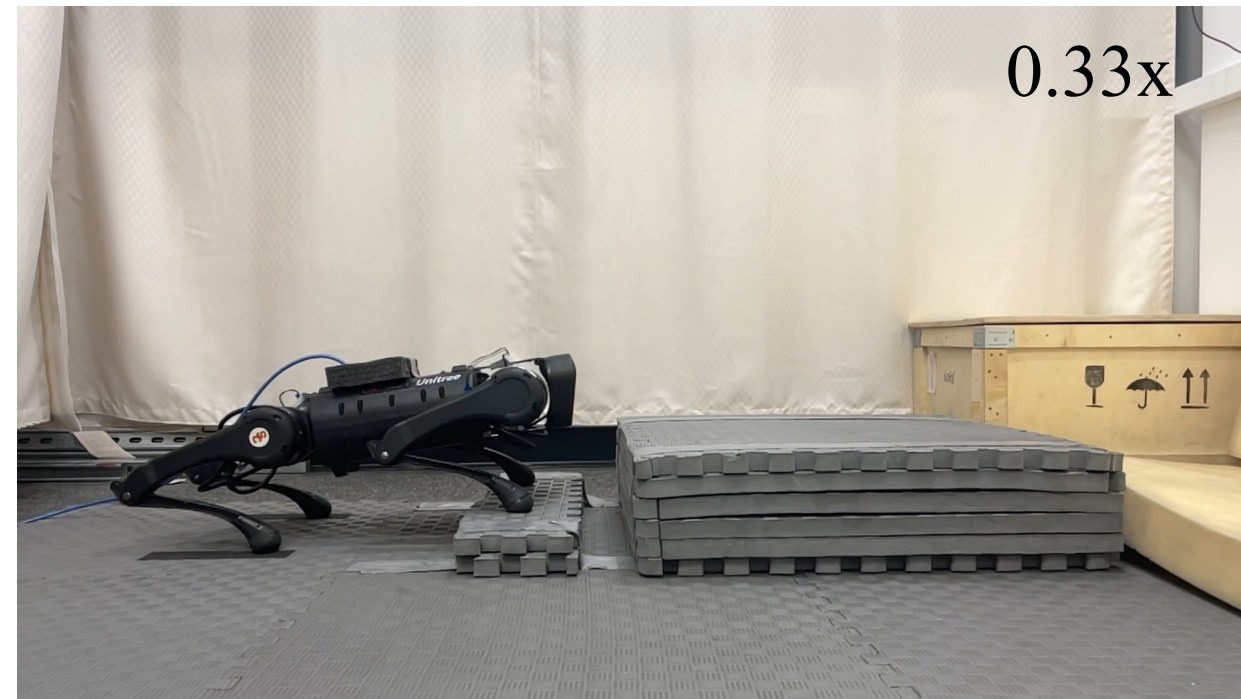
# Example Hardware Result

- Jumping height and distance: **(0.2, 0.6) m**
- Ground uncertainty: **0.06 m** block under front feet

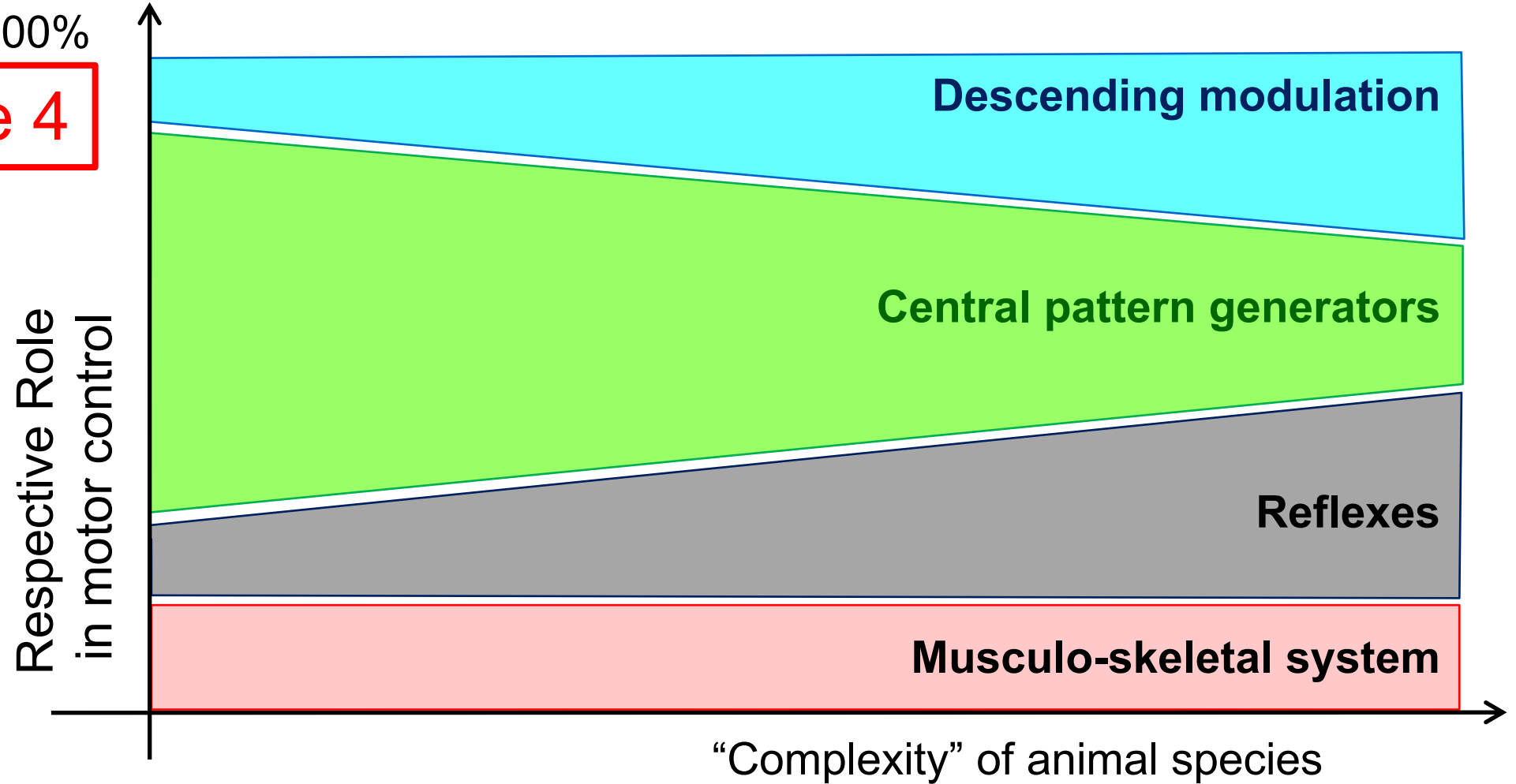
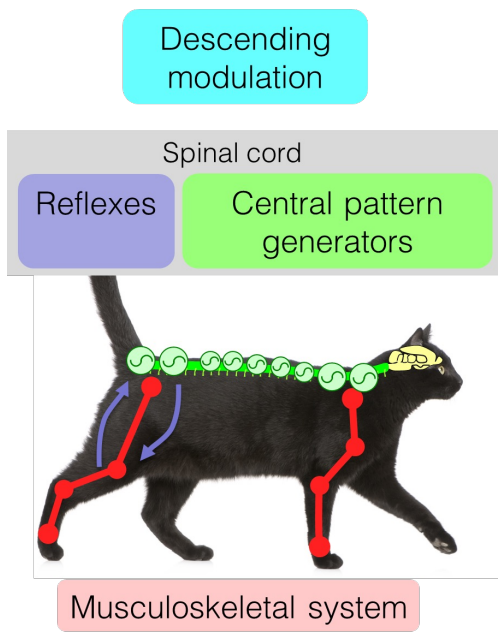
**Baseline** - Pure Trajectory Optimization



**Our method** – Trajectory Optimization + DRL

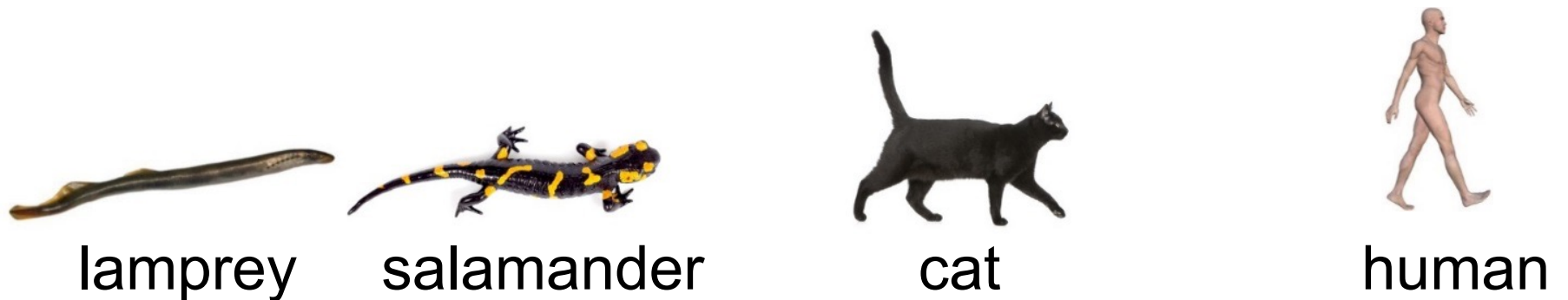


# From Lecture 4



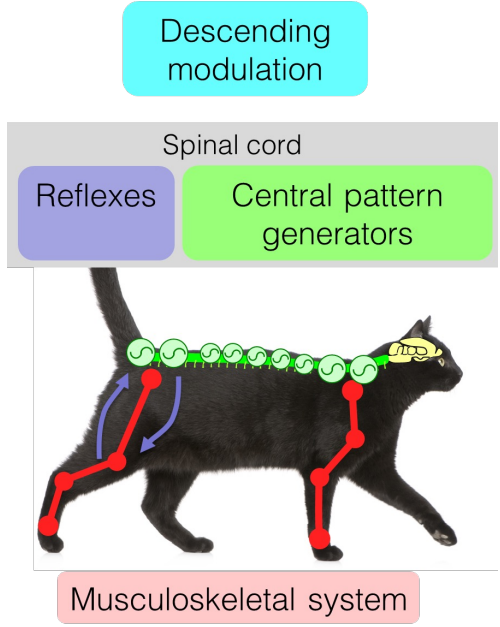
Ryczko, Simon, Ijspeert, Trends in Neuroscience, 2020

Minassian et al Neuroscientist 2017

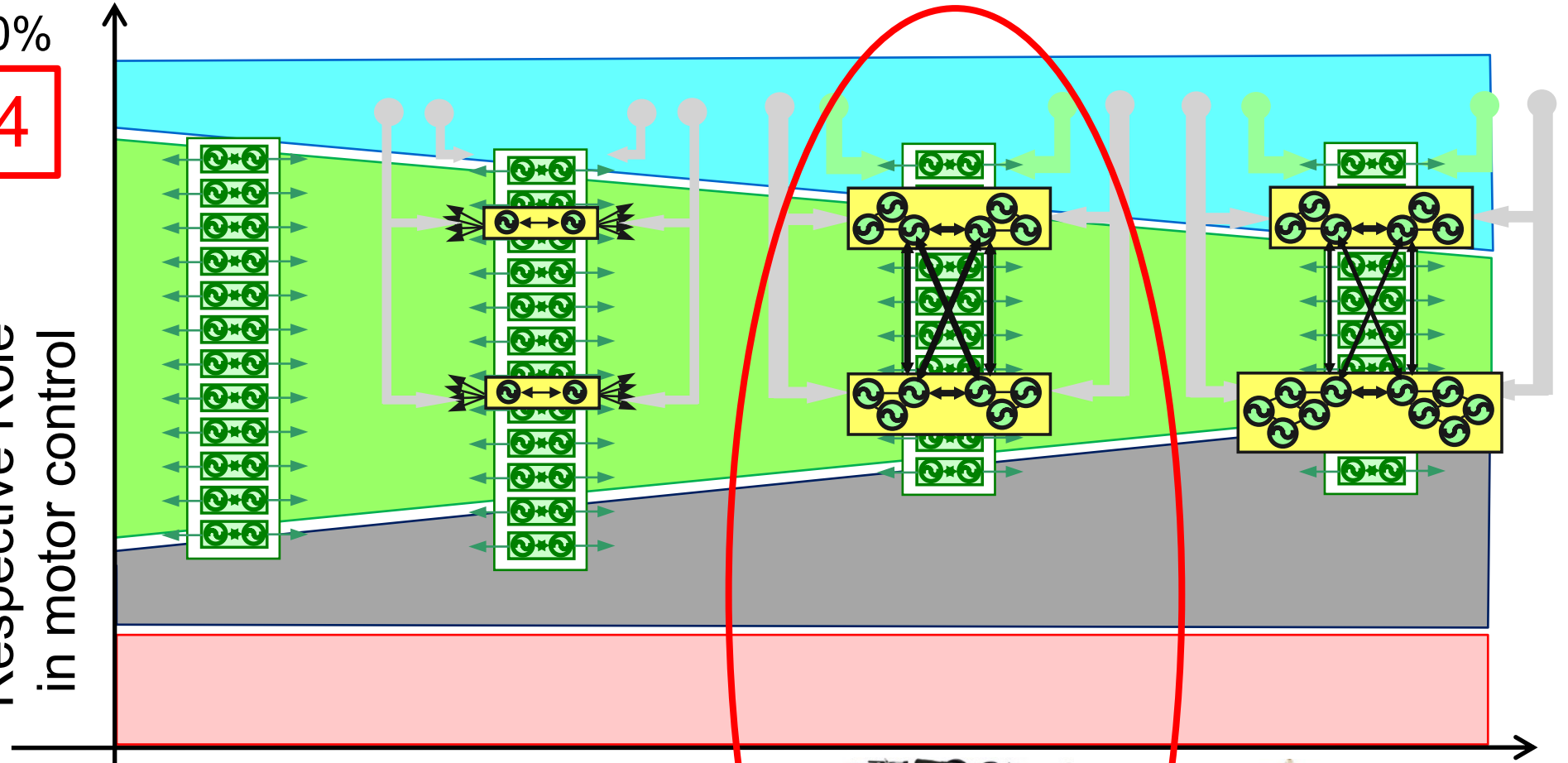


# From Lecture 4

100%



Respective Role  
in motor control



lamprey



salamander



"Comp animal sp



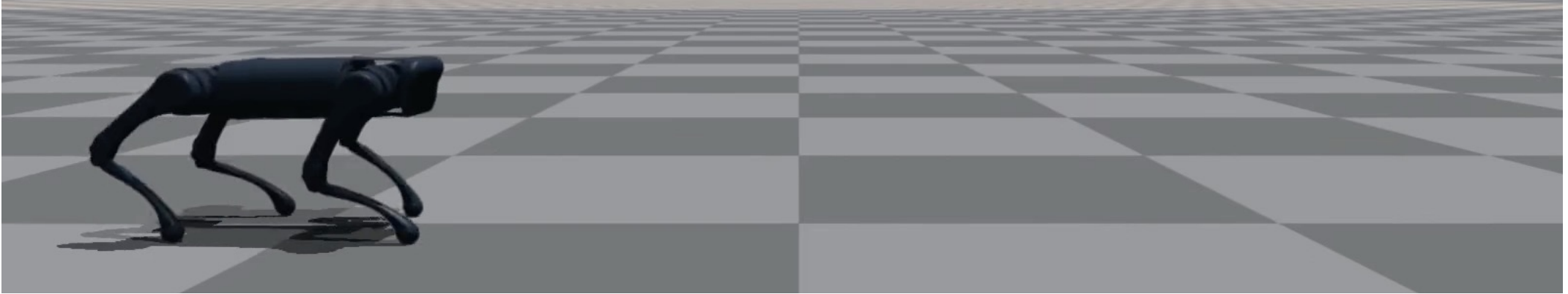
cat



human

# Traditional Central Pattern Generators for Quadruped Locomotion

Trot



Amplitude:

$$\ddot{r}_i = a \left( \frac{a}{4} (\mu_i - r_i) - \dot{r}_i \right)$$

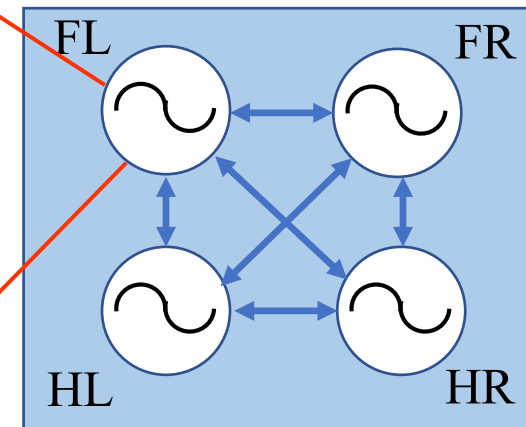
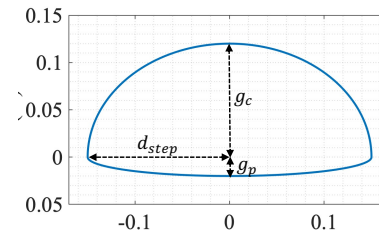
Phase:

$$\dot{\theta}_i = \omega_i + \sum_j r_j w_{ij} \sin(\theta_j - \theta_i - \phi_{ij}), \quad \omega_i = \begin{cases} \omega_{\text{swing}} & \text{if } 0 \leq \theta_i < \pi \\ \omega_{\text{stance}} & \text{if } \pi \leq \theta_i < 2\pi \end{cases}$$

Output:

$$x_{\text{foot},i} = -d_{\text{step}} r_i \cos(\theta_i)$$

$$z_{\text{foot},i} = \begin{cases} -h + g_c \sin(\theta_i) & \text{if } \sin(\theta_i) > 0 \\ -h + g_p \sin(\theta_i) & \text{otherwise} \end{cases}$$



# Traditional Central Pattern Generators for Quadruped Locomotion

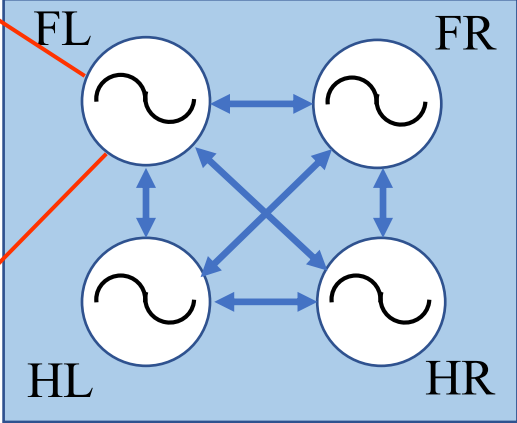
- Require careful parameter tuning (optimization-based or hand-tuned)
- Result in a single fixed gait (must re-tune for each desired gait and speed)
- Difficult to specify omni-directional motion (i.e. add control like VMC to turn)

Can we improve the general capabilities of CPGs with deep reinforcement learning?

Amplitude:  $\ddot{r}_i = a \left( \frac{a}{4} (\mu_i - r_i) - \dot{r}_i \right)$

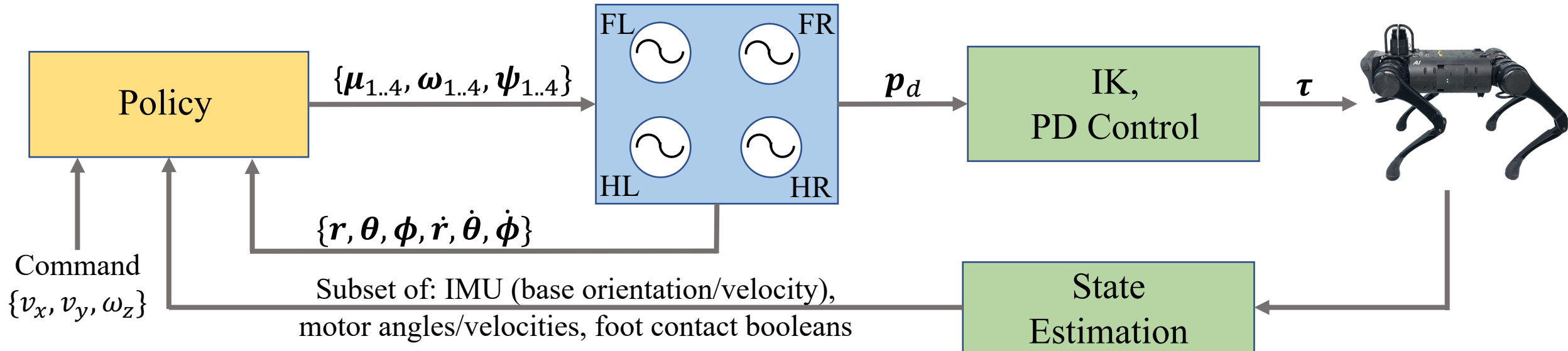
Phase:  $\dot{\theta}_i = \omega_i + \sum_j r_j w_{ij} \sin(\theta_j - \theta_i - \phi_{ij})$ ,  $\omega_i = \begin{cases} \omega_{\text{swing}} & \text{if } 0 \leq \theta_i < \pi \\ \omega_{\text{stance}} & \text{if } \pi \leq \theta_i < 2\pi \end{cases}$

Output:  $x_{\text{foot},i} = -d_{\text{step}} r_i \cos(\theta_i)$   
 $z_{\text{foot},i} = \begin{cases} -h + g_c \sin(\theta_i) & \text{if } \sin(\theta_i) > 0 \\ -h + g_p \sin(\theta_i) & \text{otherwise} \end{cases}$



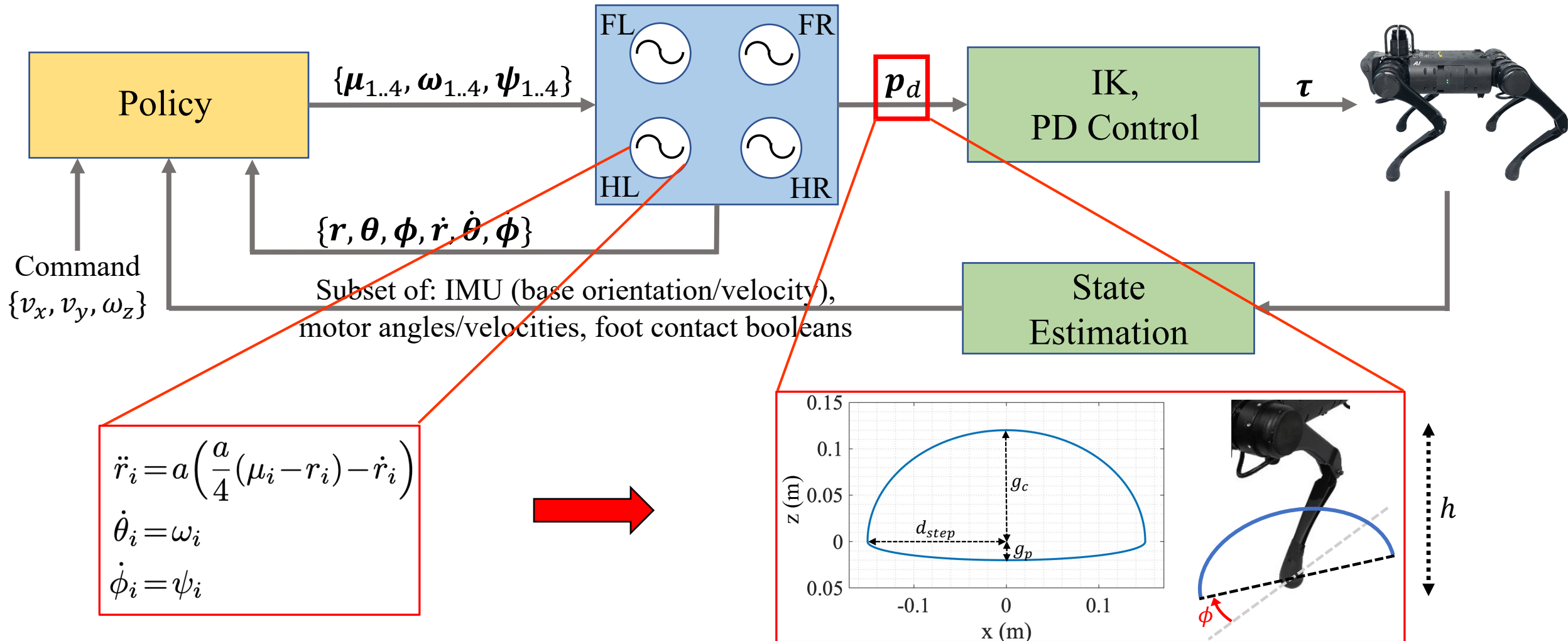
The diagram shows a central pattern generator (CPG) network for a quadruped. It consists of four nodes arranged in a 2x2 grid, labeled FL (Front Left), FR (Front Right), HL (Hind Left), and HR (Hind Right). Each node contains a sine wave symbol, representing the oscillatory output of the CPG. Bidirectional blue arrows connect every pair of adjacent nodes (FL-FR, HL-HR, FL-HL, FR-HR) and also cross-connections between FL-HR and FR-HL, indicating a fully interconnected network structure.

# CPG-RL: Learning Central Pattern Generators for Quadruped Locomotion



In this work, we propose to learn to modulate the neural oscillator intrinsic amplitude and frequency of each oscillator that together forms a Central Pattern Generator with deep reinforcement learning.

# CPG-RL: Learning Central Pattern Generators for Quadruped Locomotion

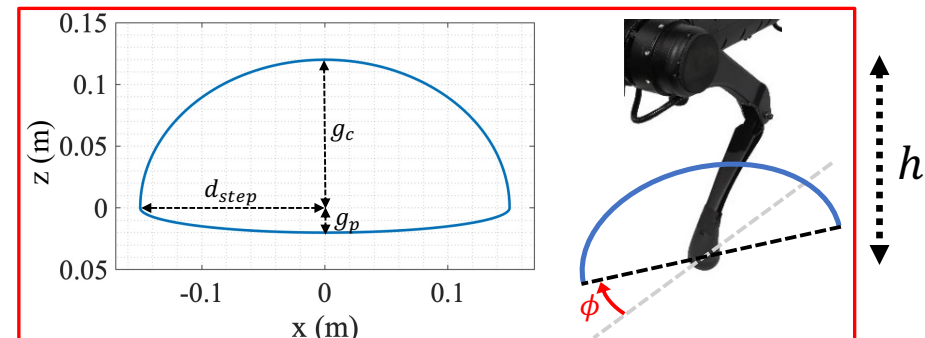
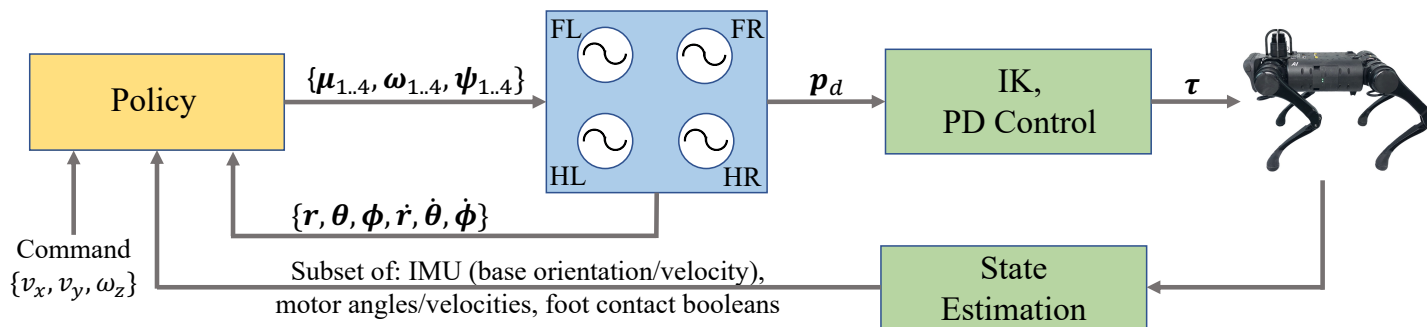
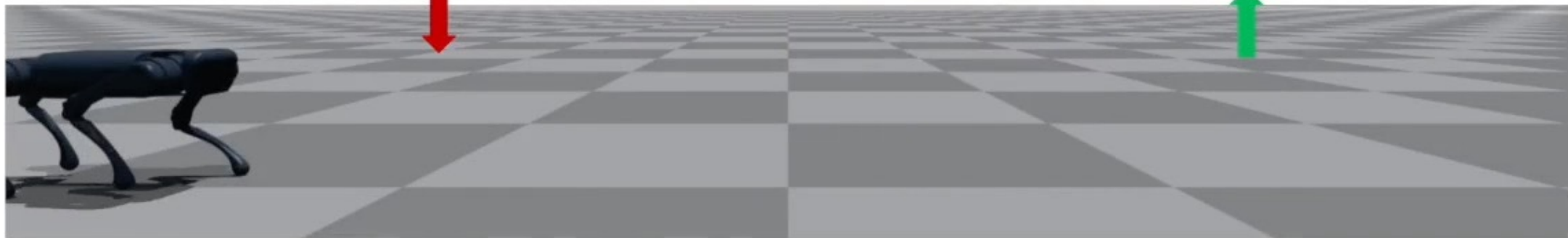


# Body Height and Swing Foot Height can be adjusted on the fly

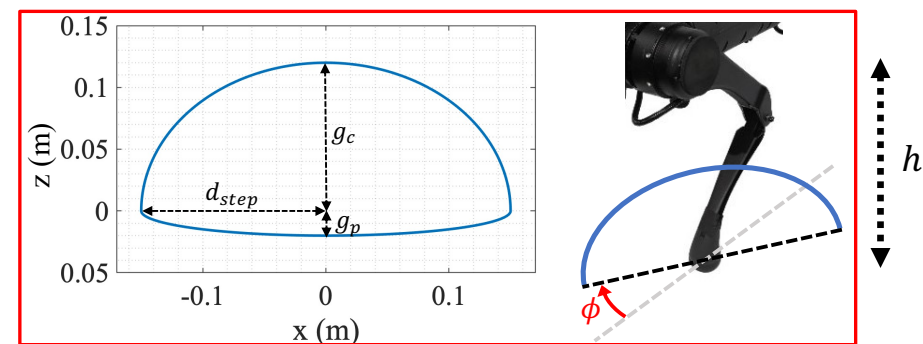
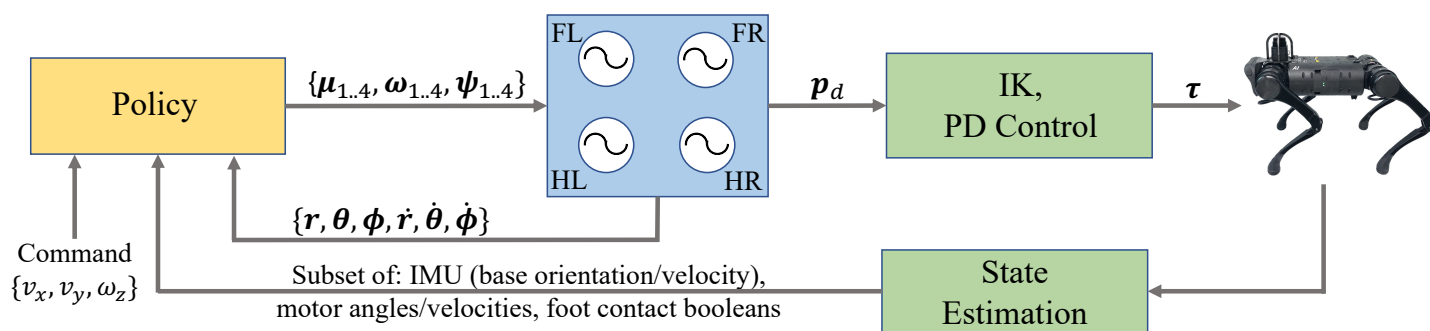
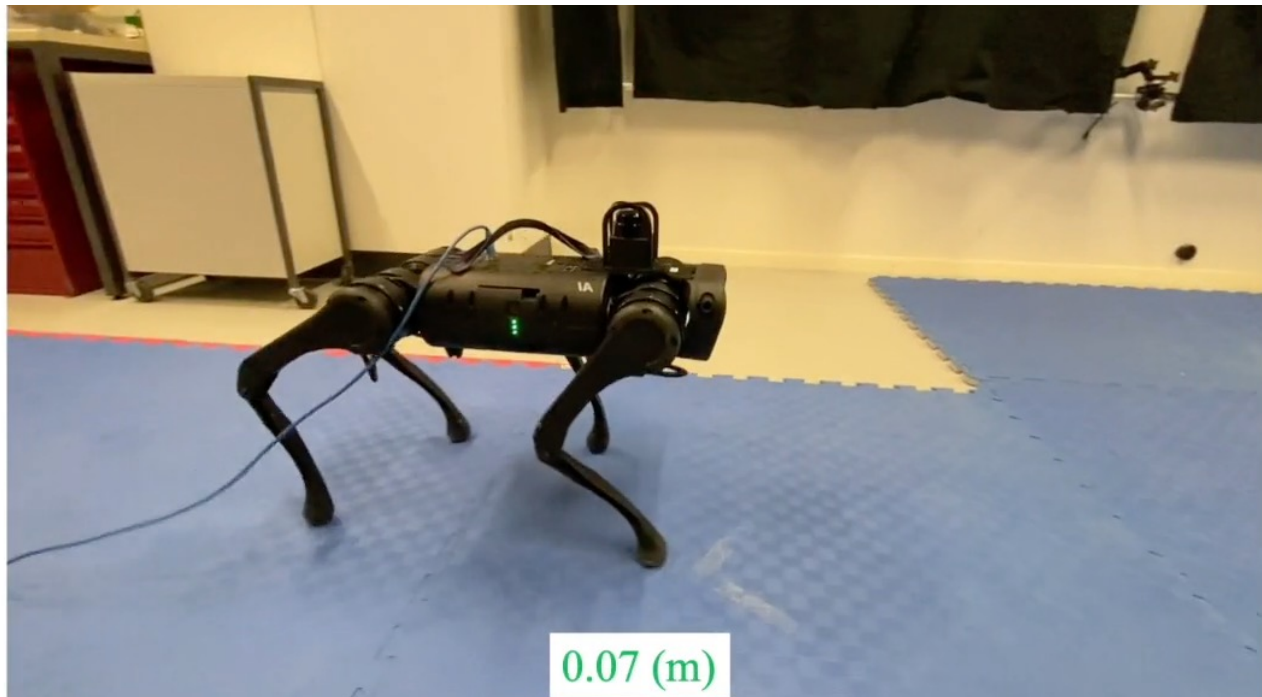
Body Height: 0.3 (m)

0.19 (m)

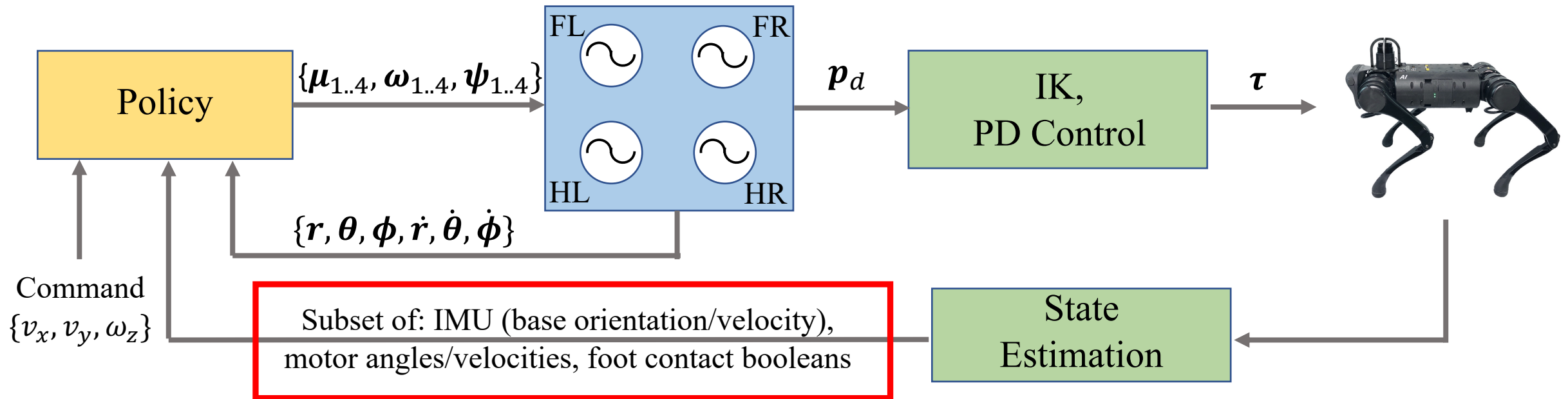
0.3 (m)



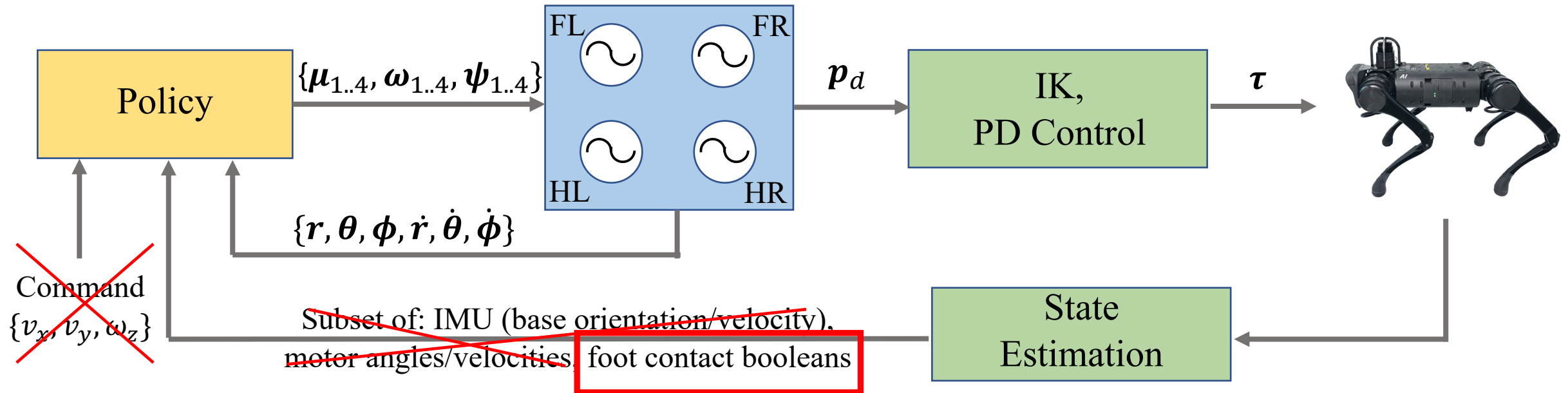
# Body Height and Swing Foot Height can be adjusted on the fly



# What sensory information should be in the observation space?



# What sensory information should be in the observation space?

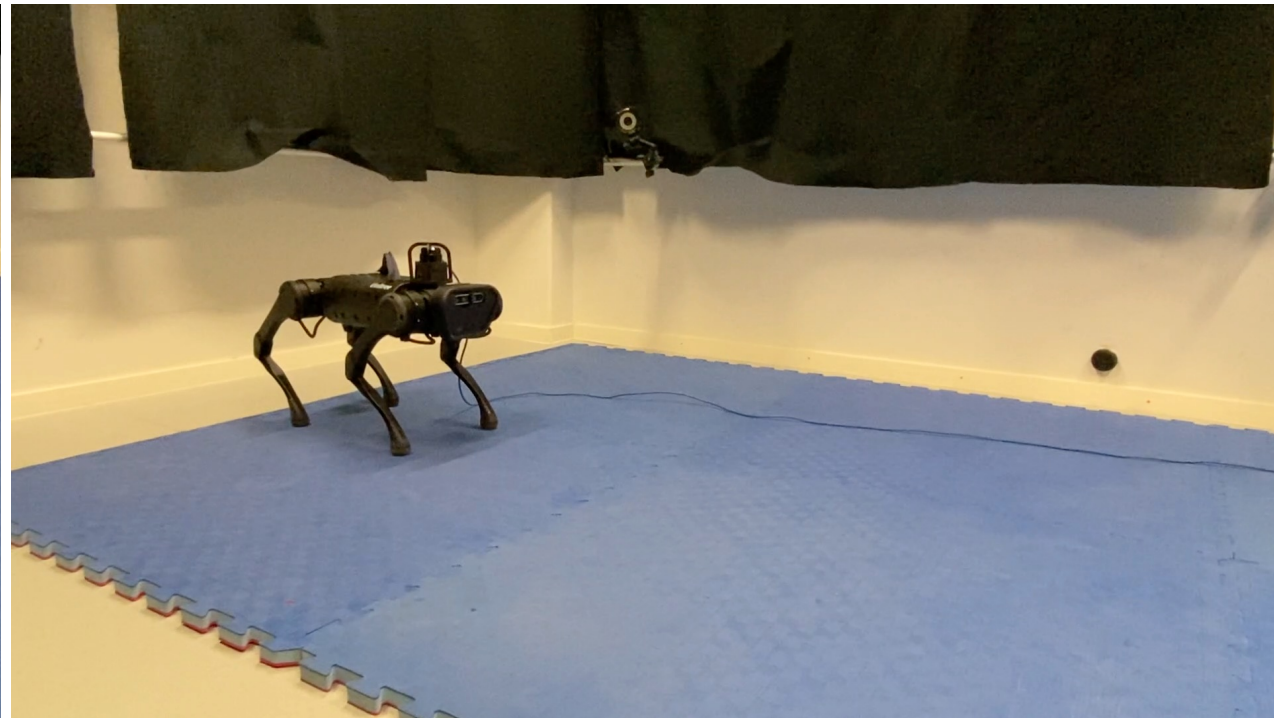
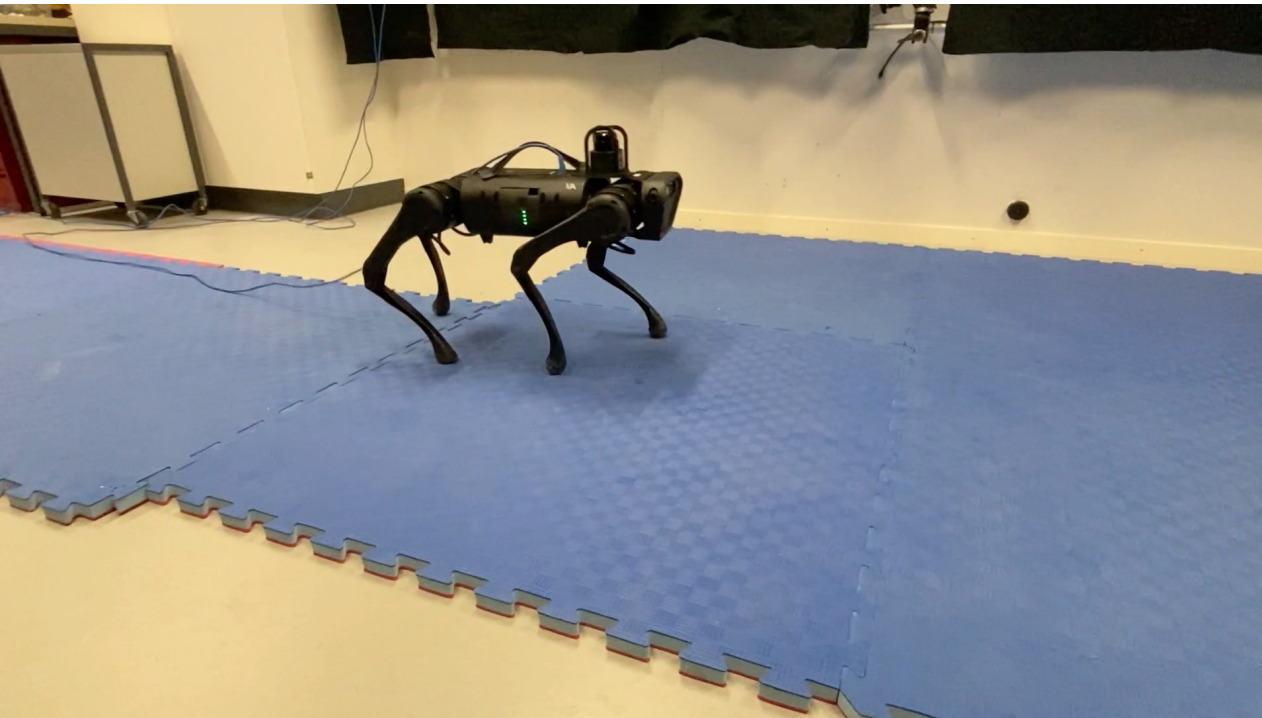


- Similar to the “Tegotae” feedback idea:

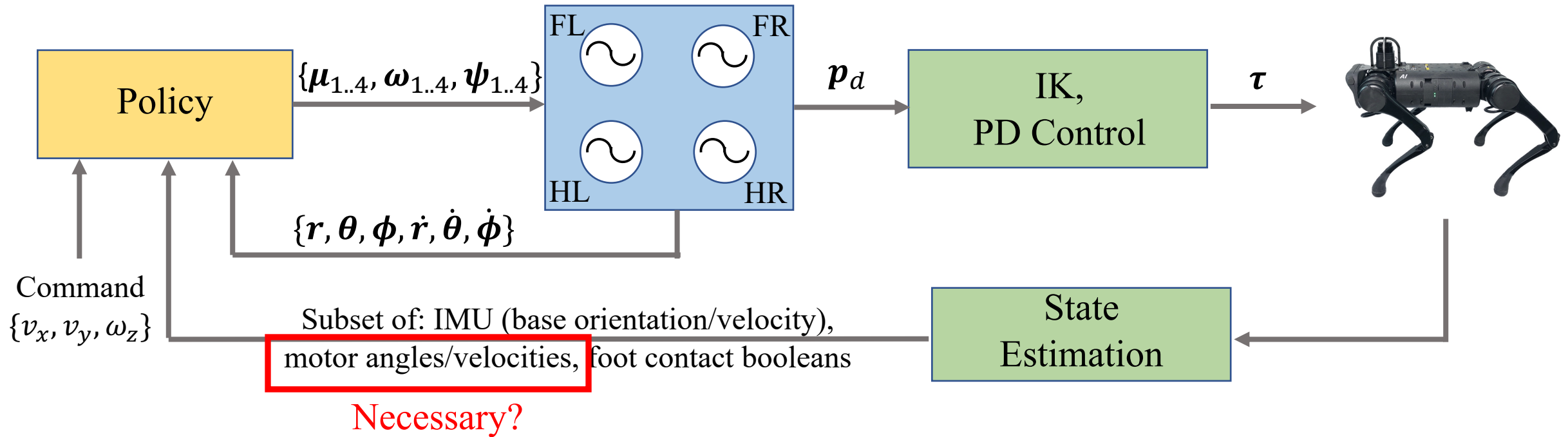
$$\dot{\theta}_i = \omega_i + \sum r_j w_{ij} \sin(\theta_j - \theta_i - \phi_{ij}) - \sigma N_i \cos \theta_i$$

# Observation space: {Foot contact booleans, CPG states}

- No domain randomization or noise during training



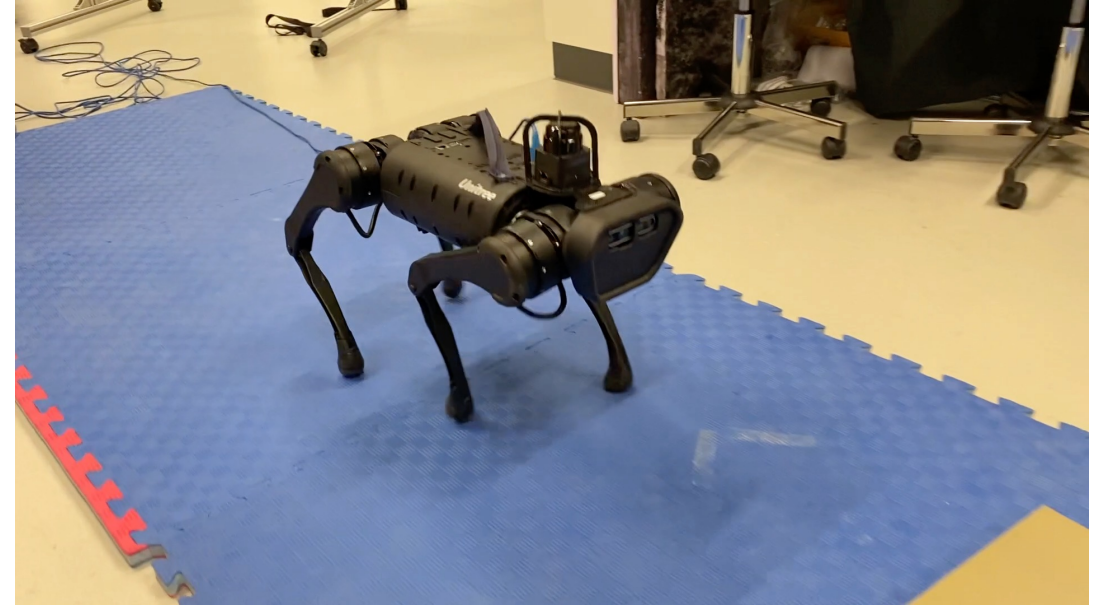
# What sensory information should be in the observation space?



Full Obs. Space,  $v_{b,x}^* = 0.3$  [m/s]



Medium Obs. Space,  $v_{b,x}^* = 0.8$  [m/s]  
(No joint states in obs. space)

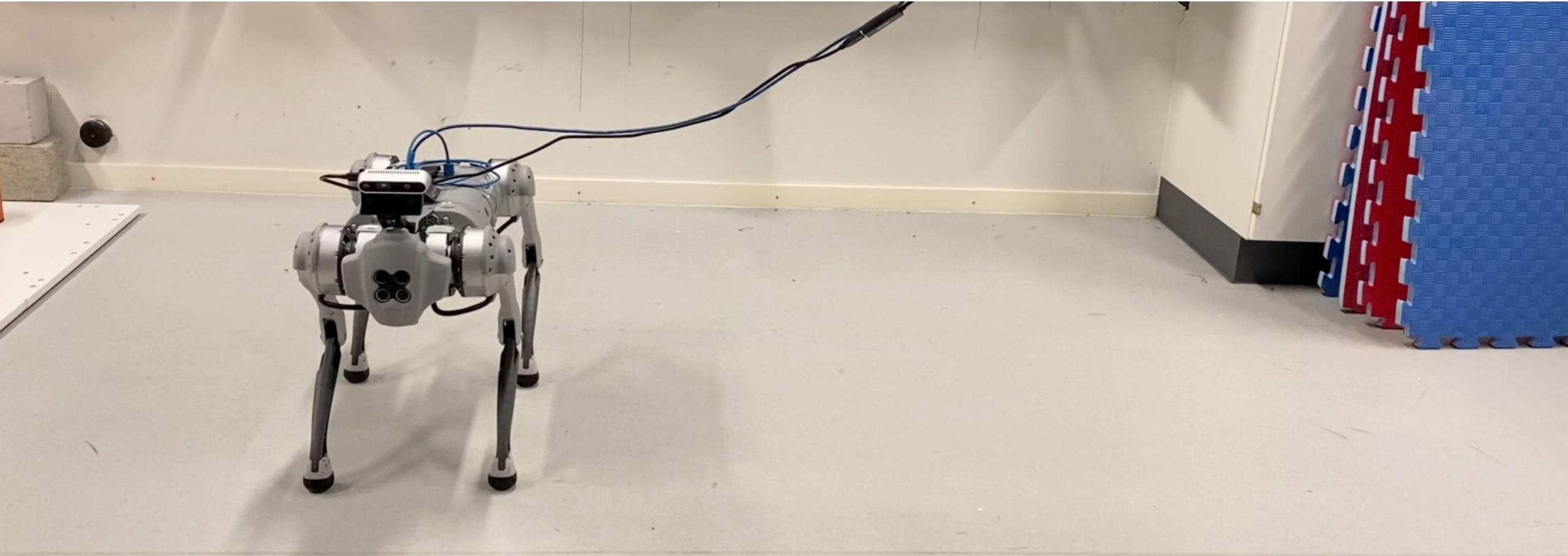


Dynamically Added Mass ( $5+5+2.5+1.25=13.75$  kg)  
Medium Observation Space,  $v_{b,x}^* = 0.8$  [m/s]

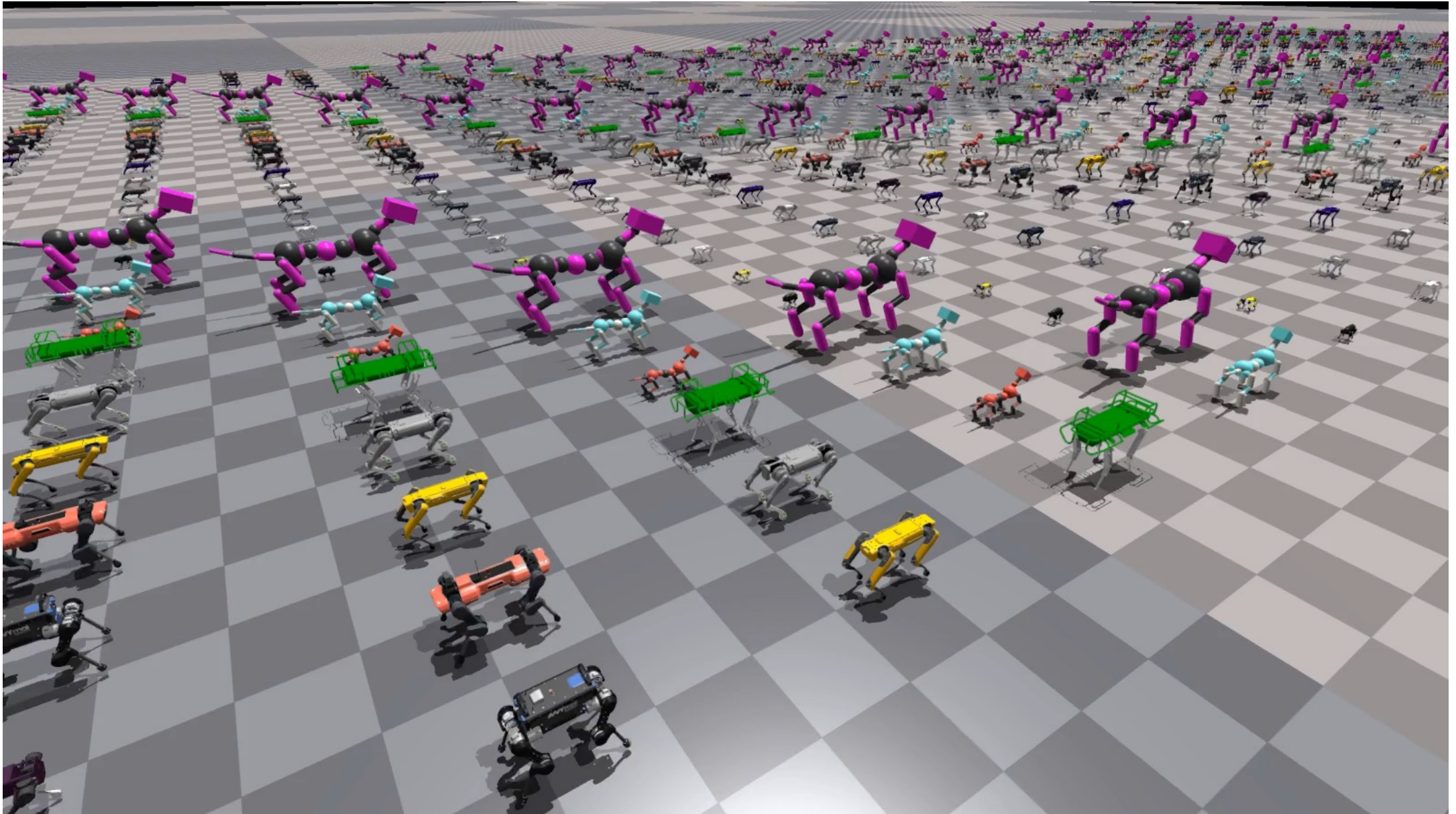


# Omnidirectional Locomotion

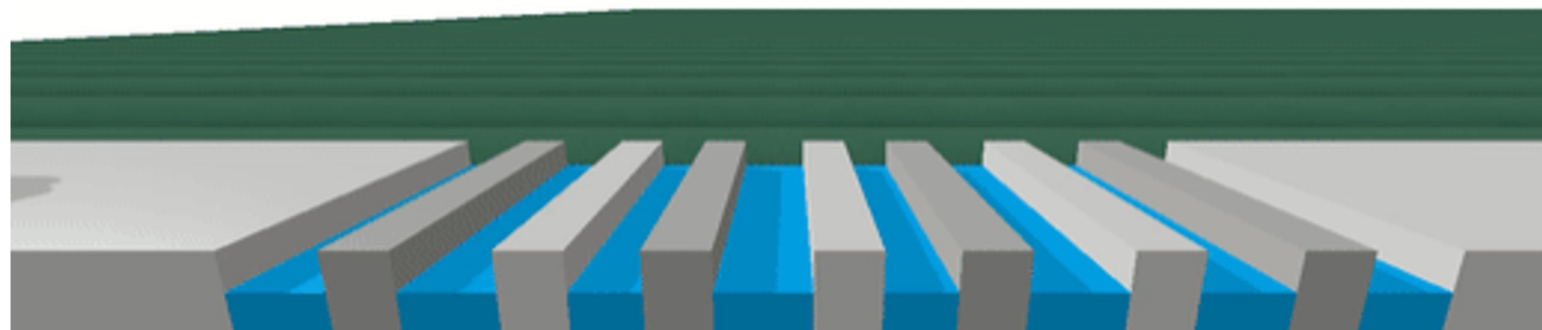
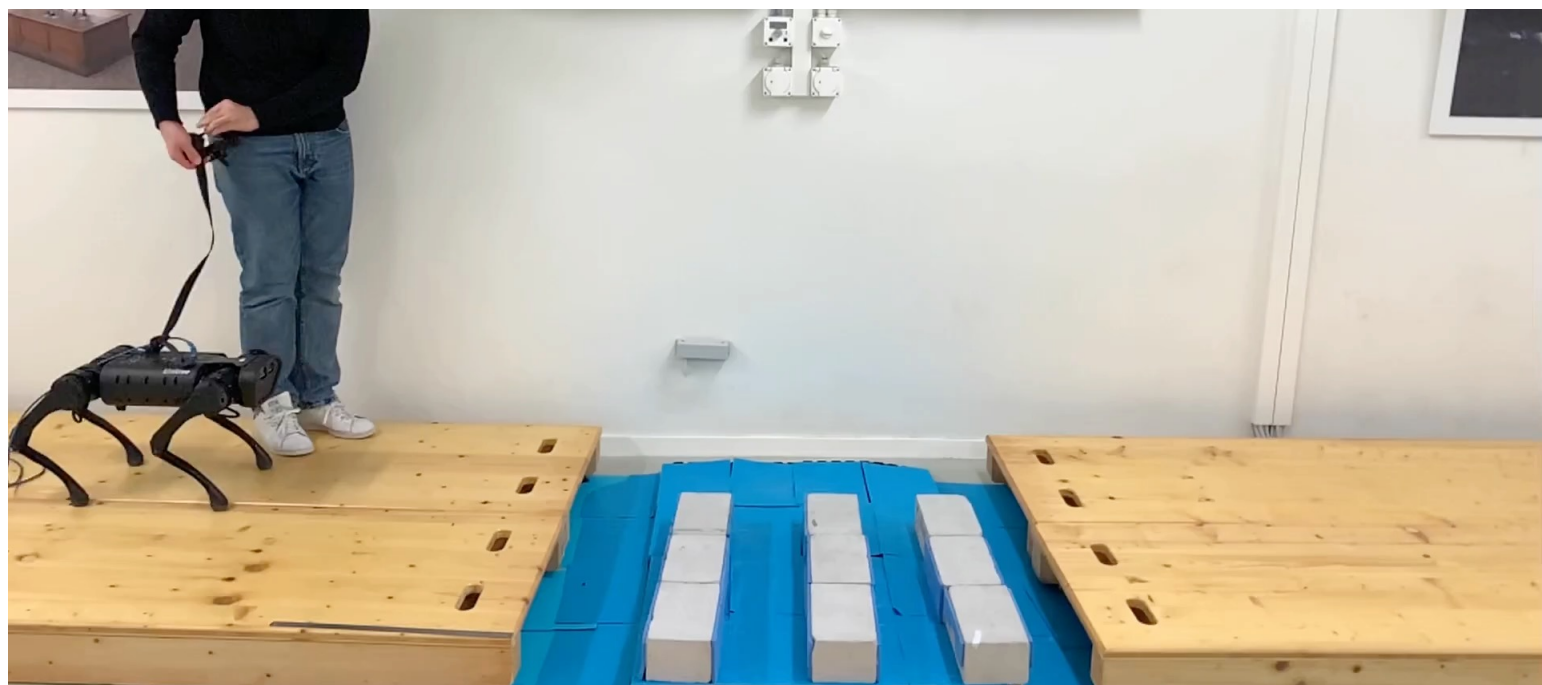
- Command:  $v_{b,y}^* = 0.4$  [m/s]



# ManyQuadrupeds: A single locomotion policy capable of controlling diverse robots

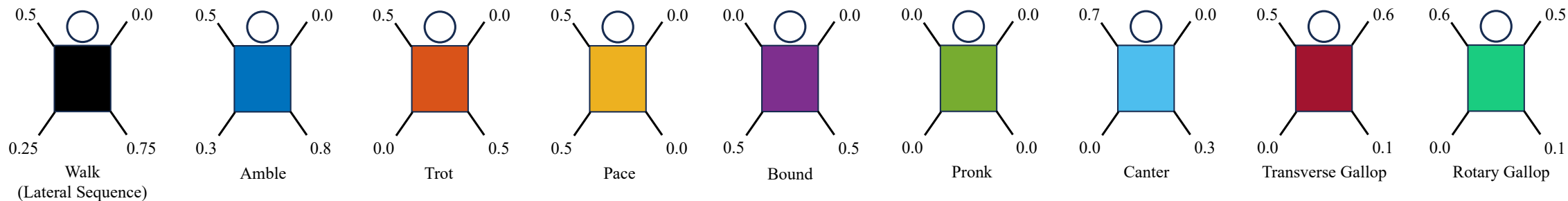
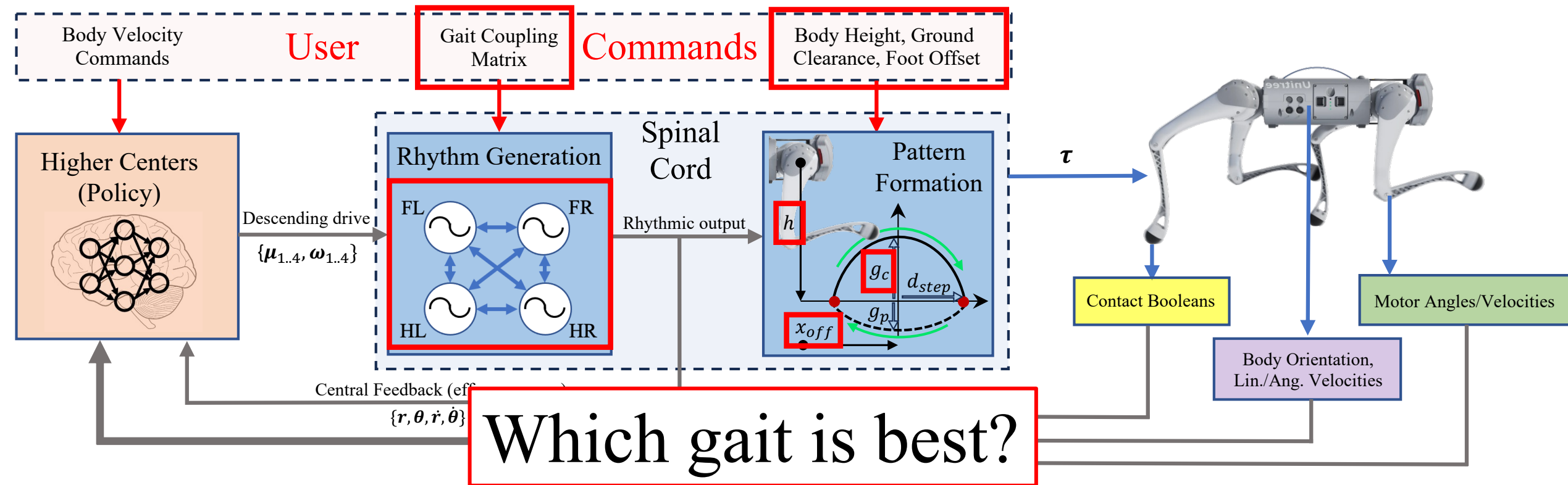


# Visual CPG-RL



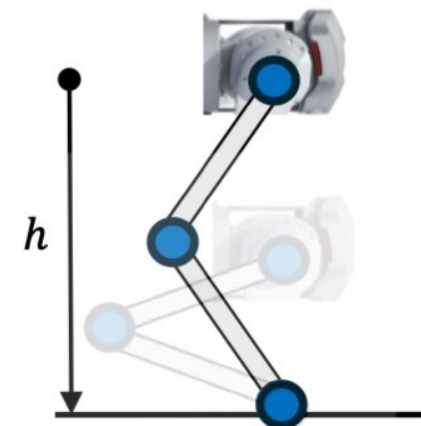
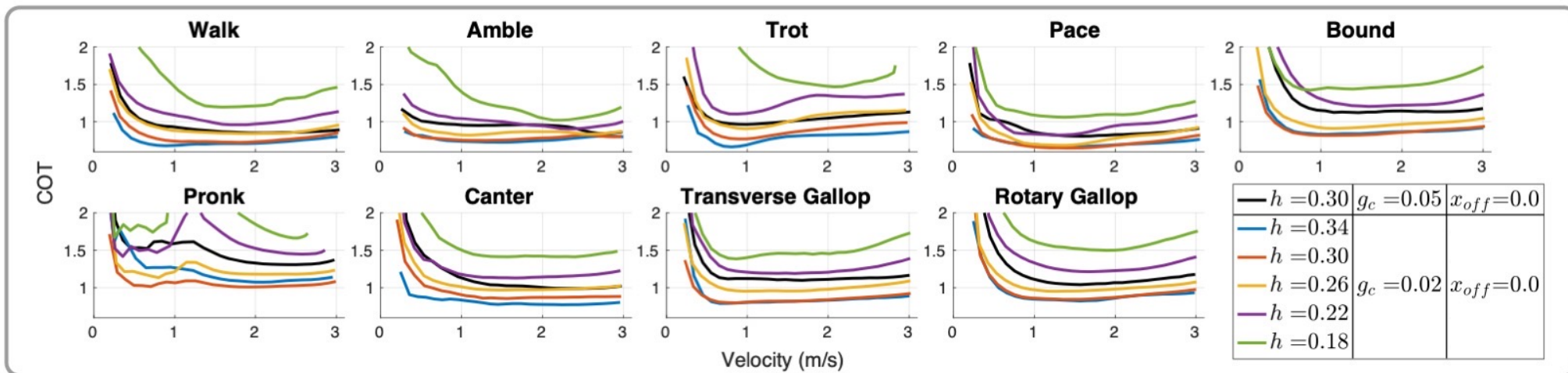
**G. Bellegarda**, M. Shafiee, A. Ijspeert. “Visual CPG-RL: Learning Central Pattern Generators for Visually-Guided Quadruped Locomotion,” ICRA 2024  
M. Shafiee, **G. Bellegarda**, A. Ijspeert. “Viability Leads to the Emergence of Gait Transitions in Learning Anticipatory Quadrupedal Locomotion Skills,” Nature Communications 2024  
M. Shafiee, **G. Bellegarda**, A. Ijspeert. “Learning Anticipatory Quadrupedal Locomotion Based on Interactions of a Central Pattern Generator and Supraspinal Drive,” ICRA 2023

# Learning All Gaits, Styles, and Transitions



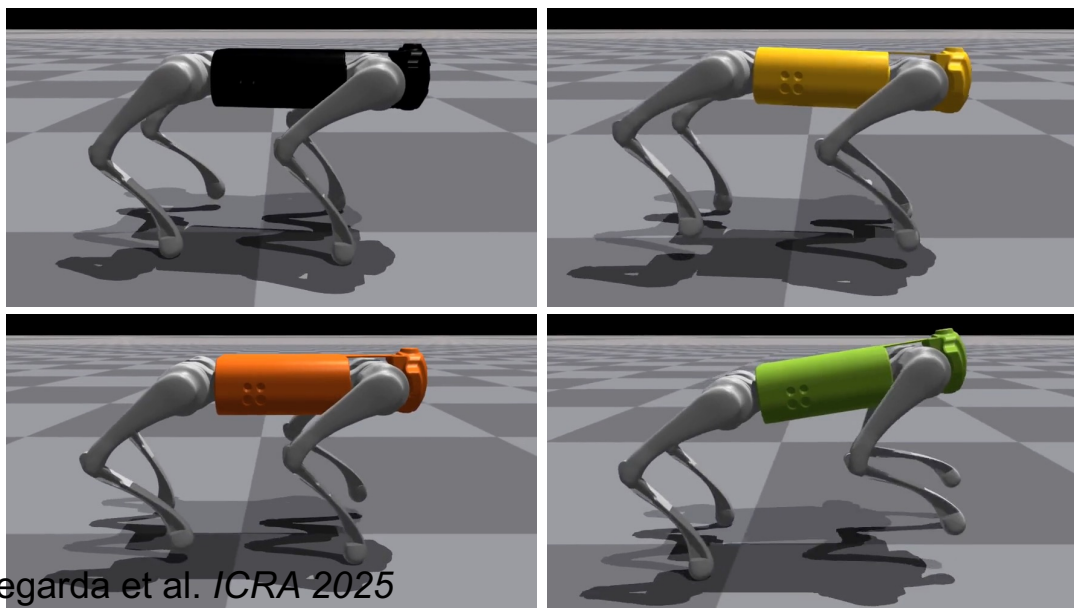
# We investigate the effects of different gait styles on the Cost of Transport (COT)

- **Changing body height ( $h$ ):** more upright postures are more energy-efficient

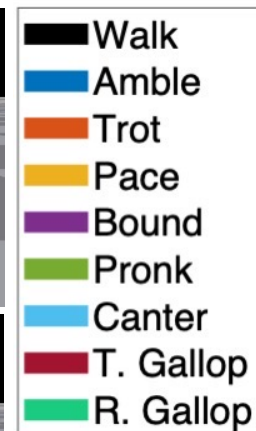
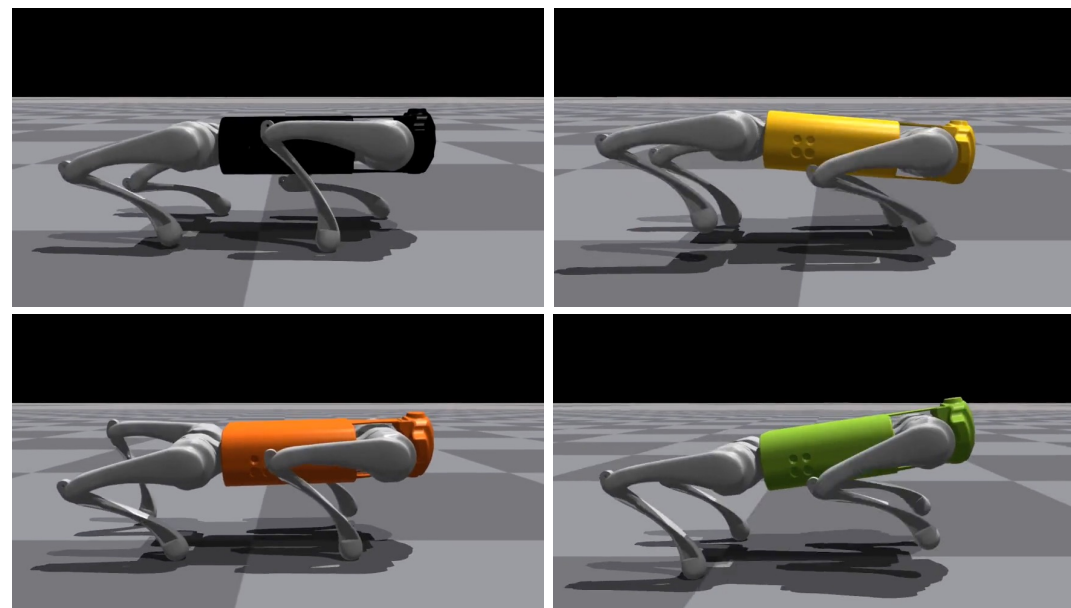


- **Examples:**

Higher Body Height (Lower COT)



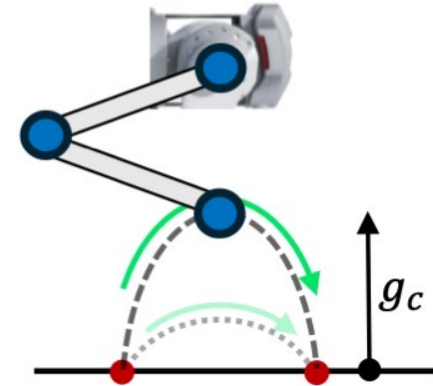
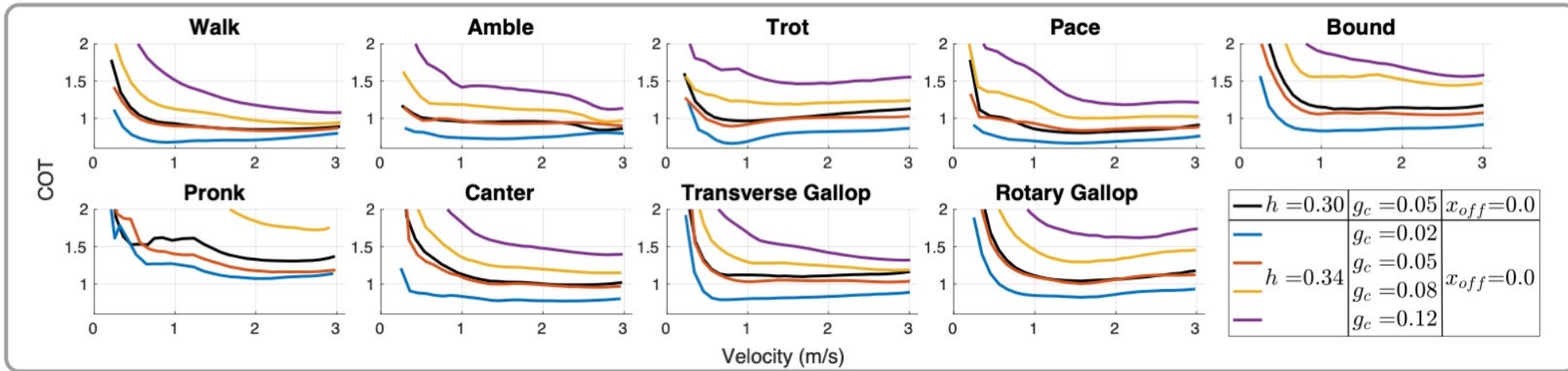
Lower Body Height (Higher COT)



0.5x

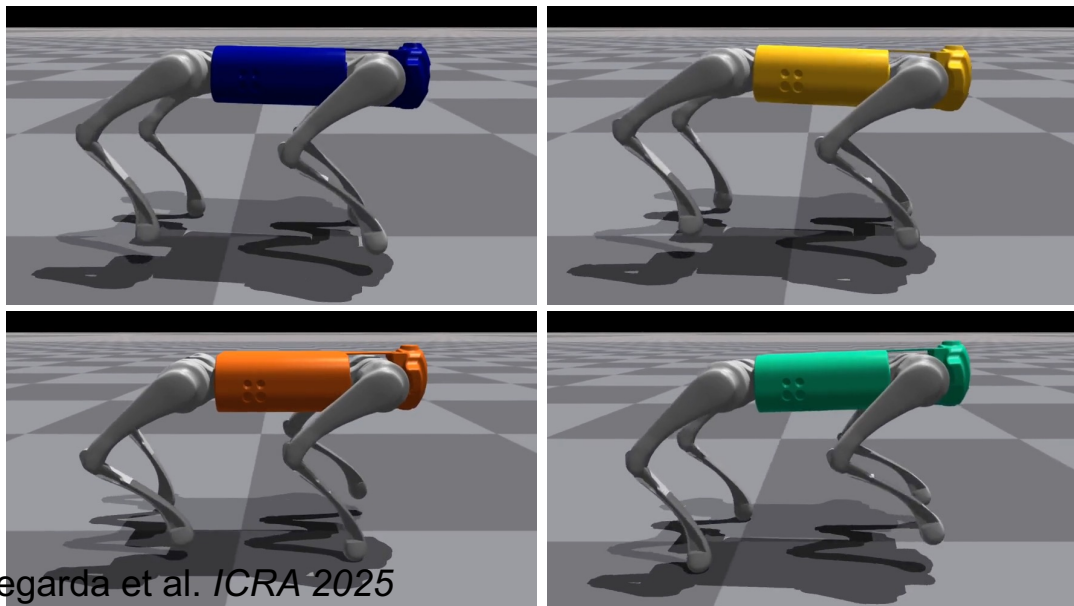
# We investigate the effects of different gait styles on the Cost of Transport (COT)

- **Changing swing foot height ( $g_c$ ):** lower swing foot height is more energy-efficient

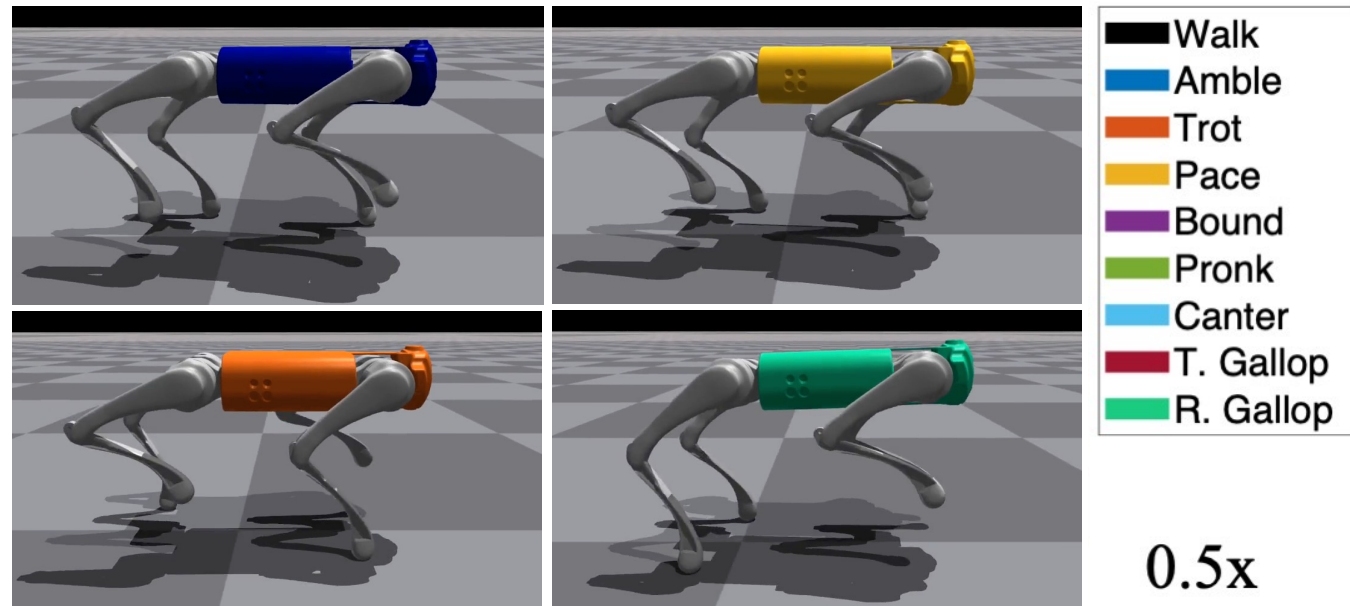


- **Examples:**

Lower Swing Foot Height (Lower COT)



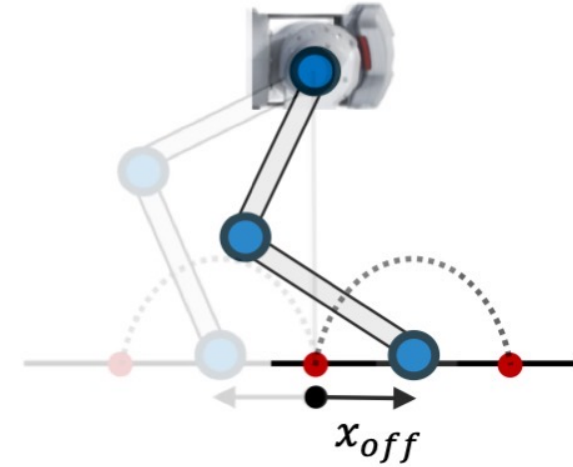
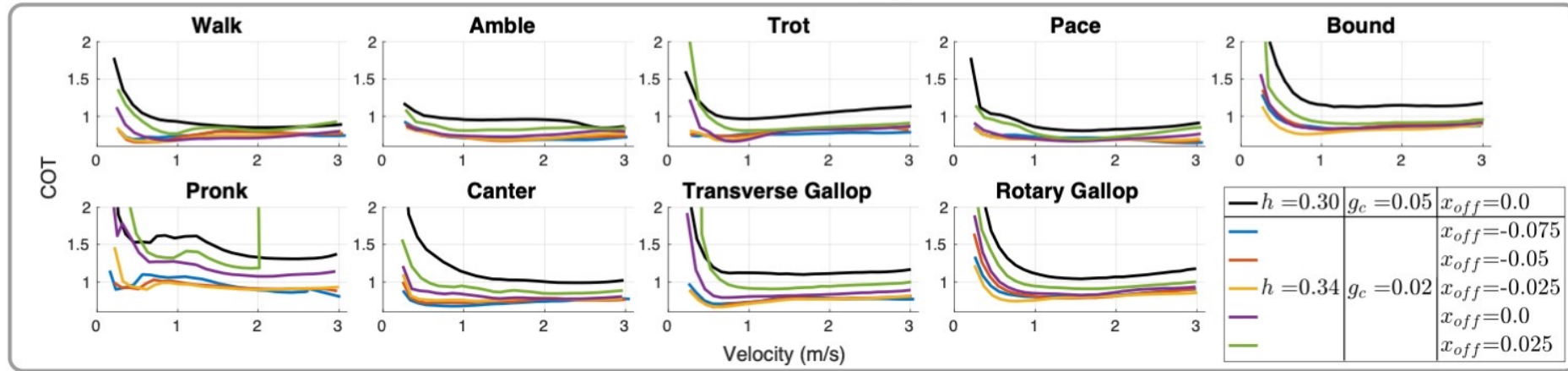
Higher Swing Foot Height (Higher COT)



0.5x

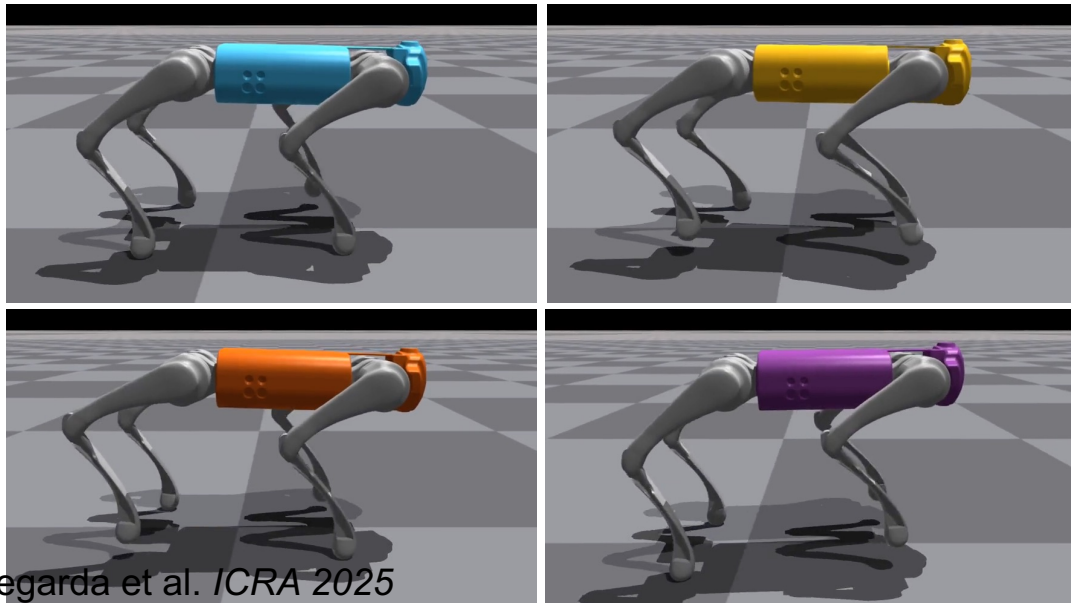
# We investigate the effects of different gait styles on the Cost of Transport (COT)

- **Changing  $x$  foot offset ( $x_{off}$ ):** variability based on velocity, but generally more energy-efficient to center the foot oscillations behind the nominal hip position

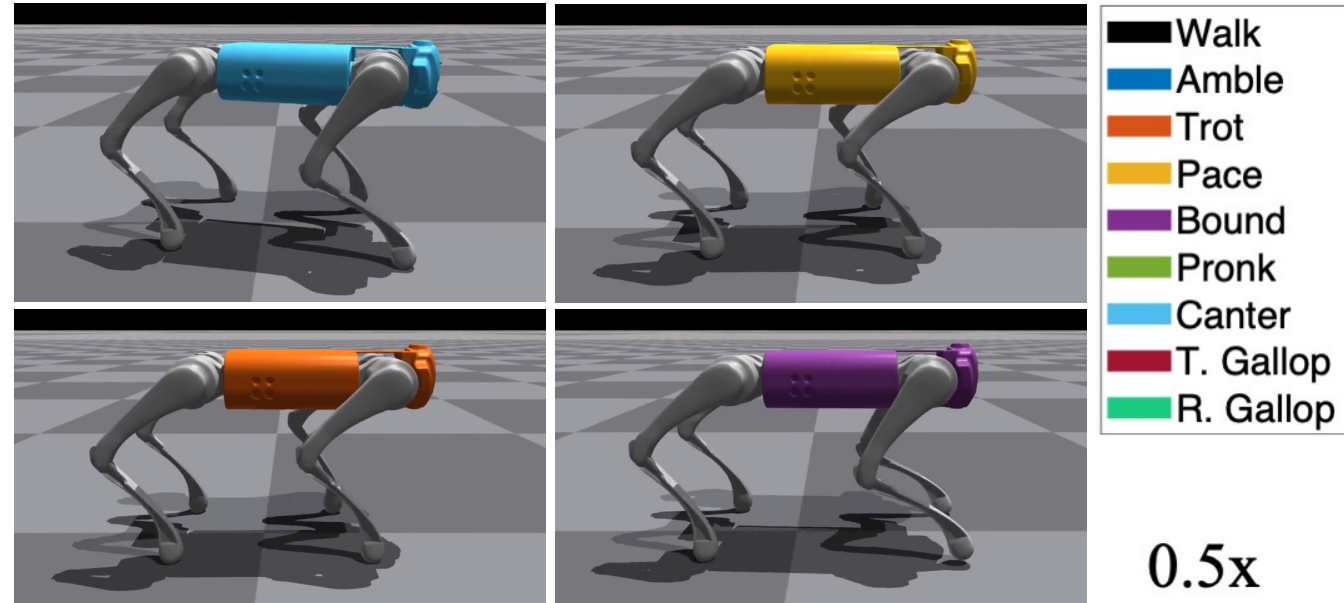


- **Examples:**

Foot offset behind hip (Lower COT)

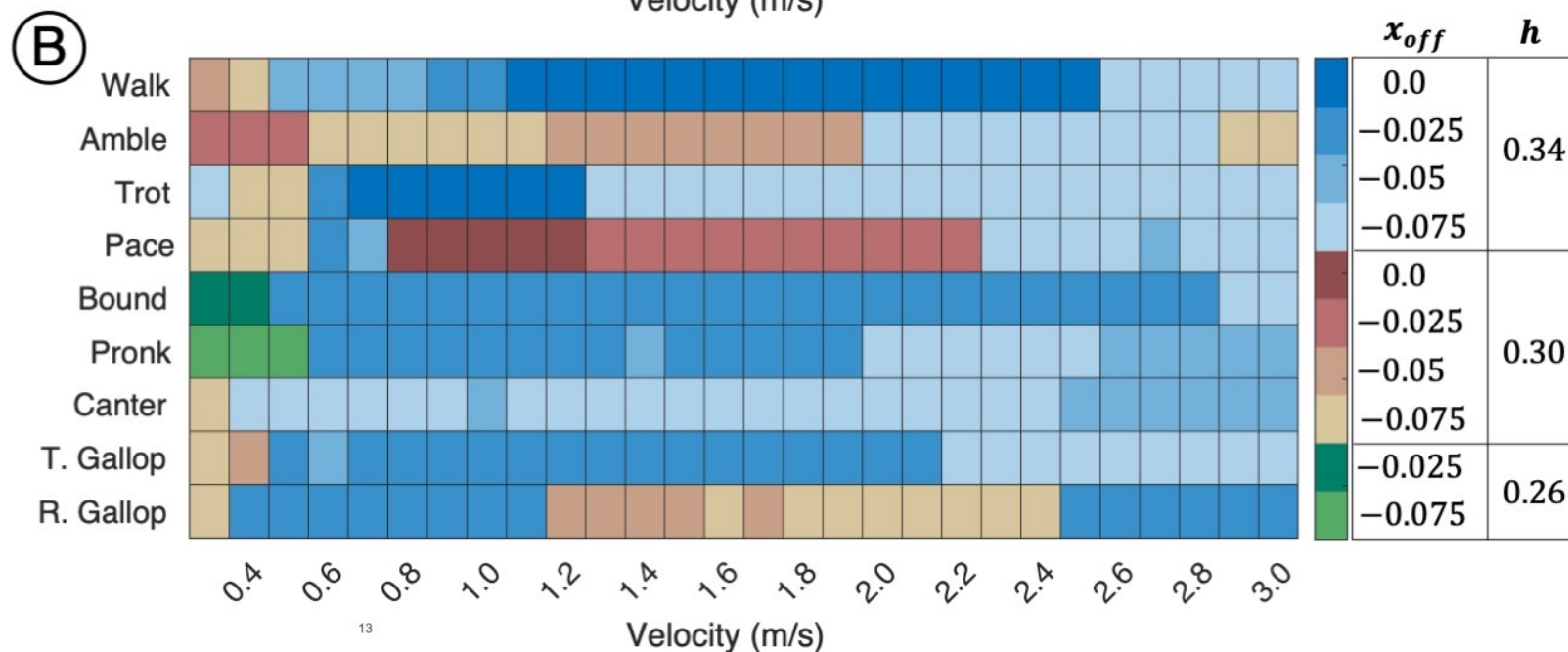
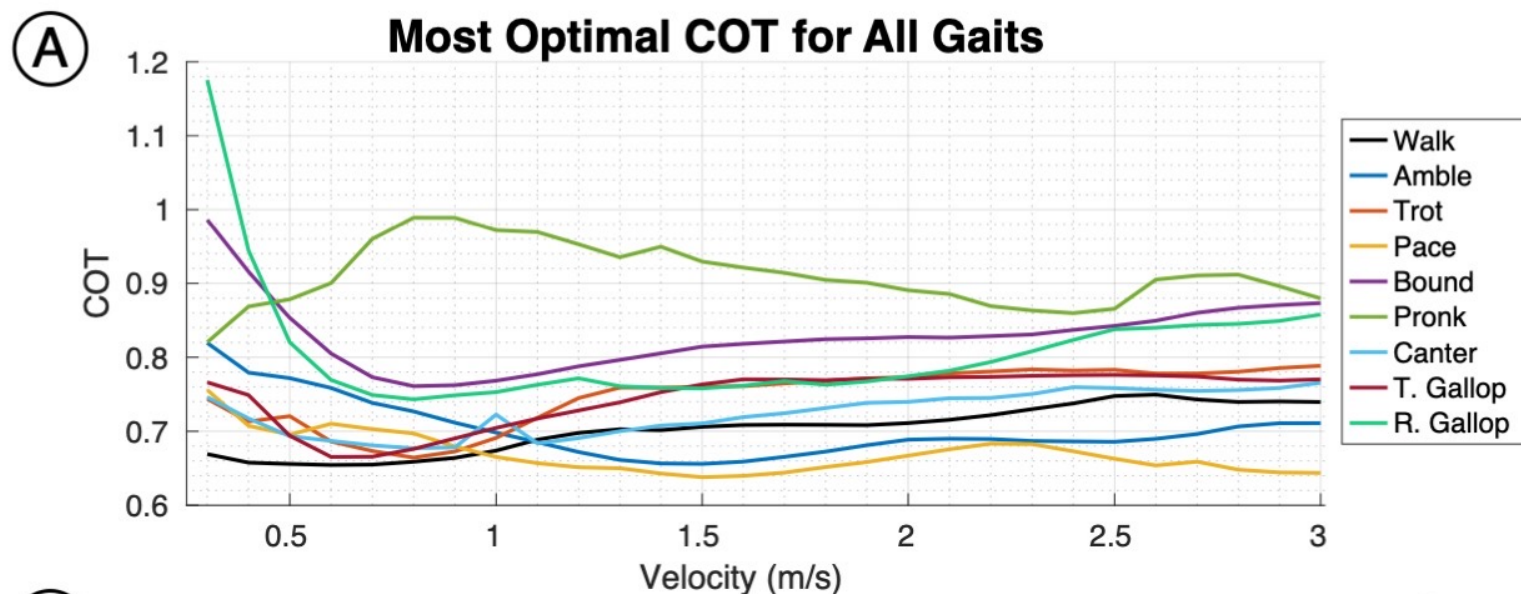


Foot offset in front of hip (Higher COT)

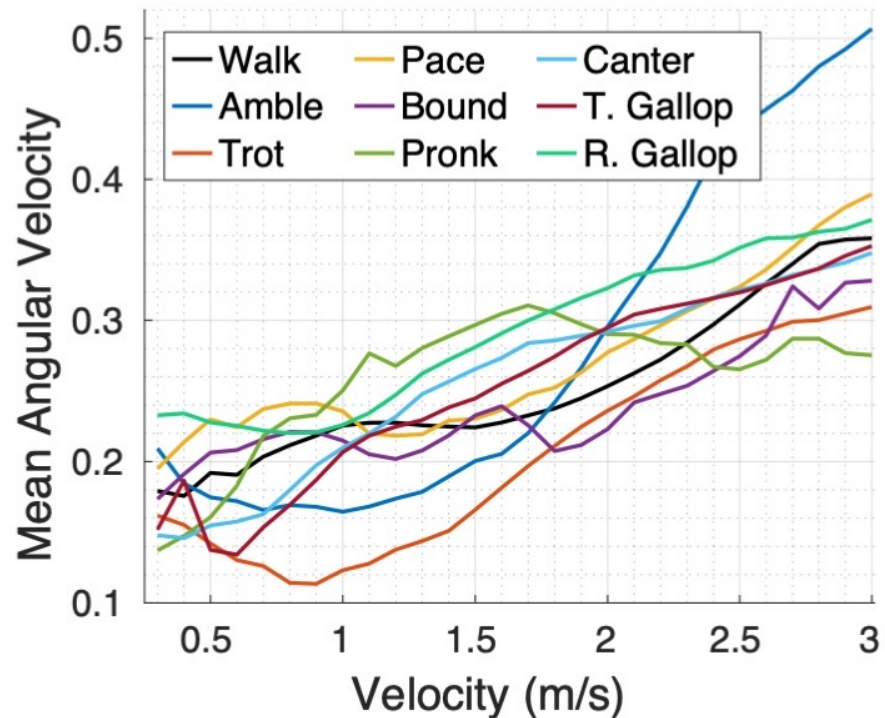


# Most Energy-Efficient Gaits

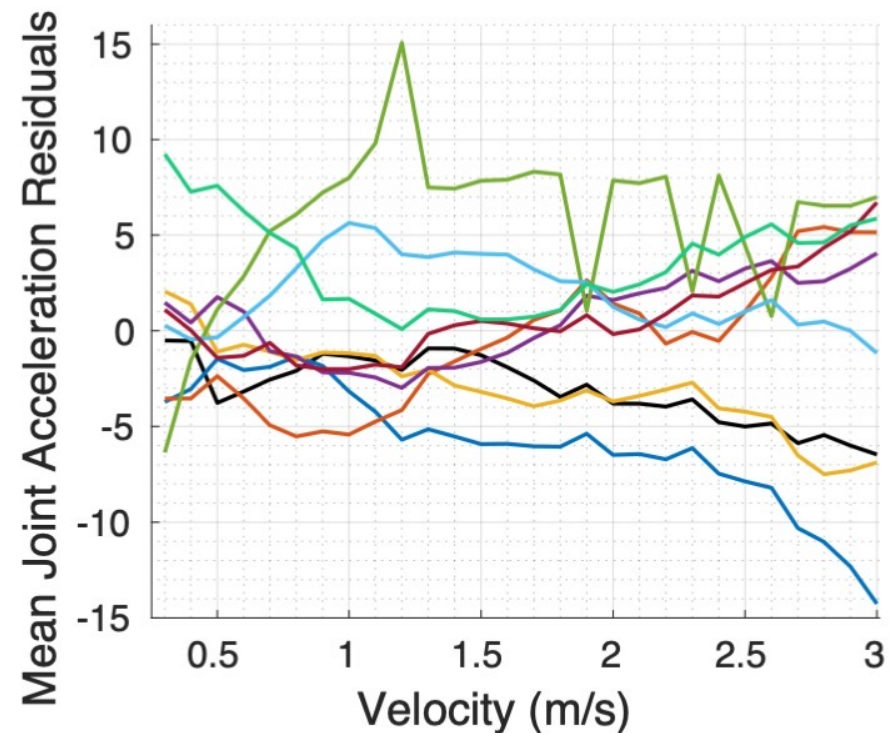
- **Walk** gait is most energy-efficient at low velocities
- **Pace** gait is most efficient at higher velocities
- Most gaits minimize the COT with:
  - **high** body height
  - **low** swing foot ground clearance
  - foot offsets **behind** the hip



## Considering Different Metrics Leads to Different ‘Optimal’ Gaits

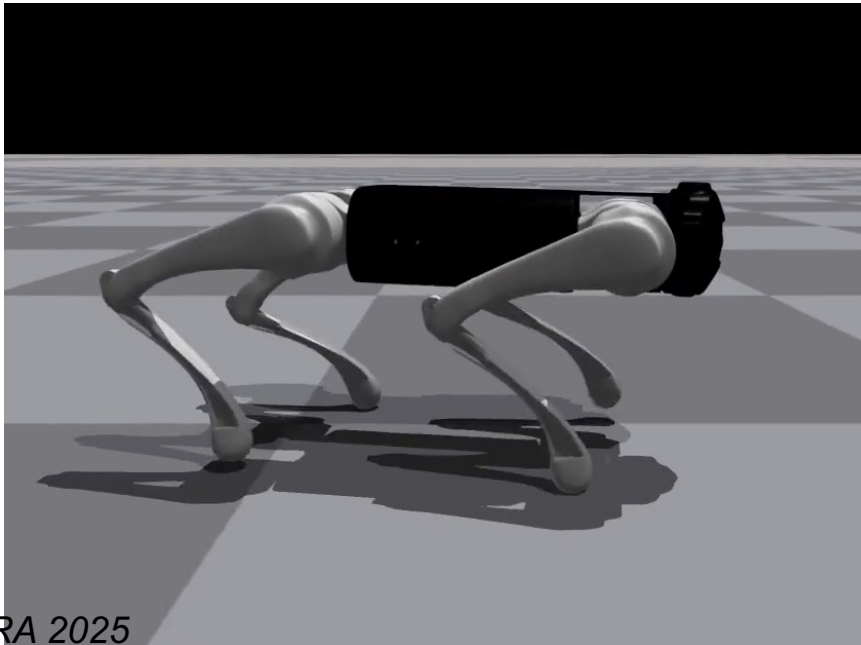
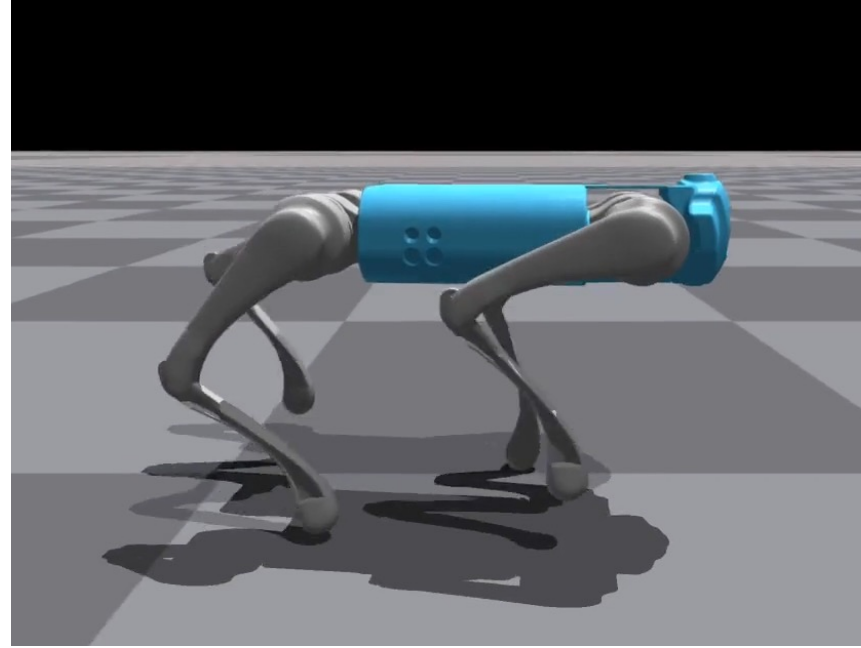
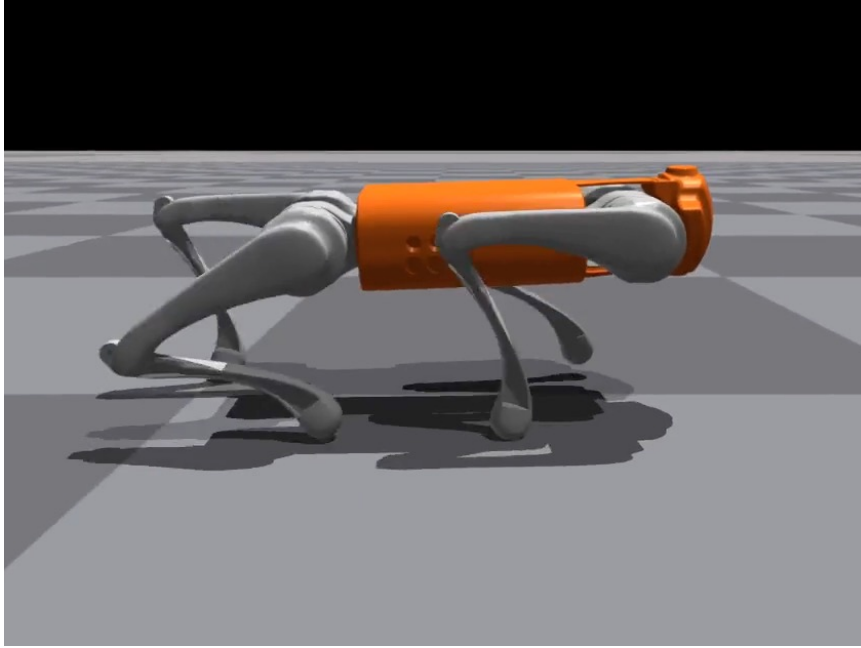


- **Trot** gait minimizes base angular velocities at low and middle speeds
- **Bound**, and then **pronk**, are ‘better’ at higher velocities



- While **amble** results in high mean base angular velocities, it minimizes mean joint acceleration relative to other gaits at the same velocity

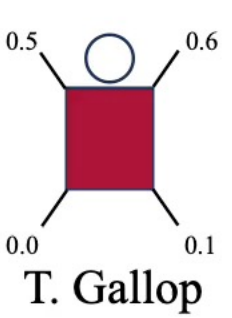
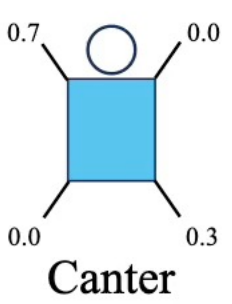
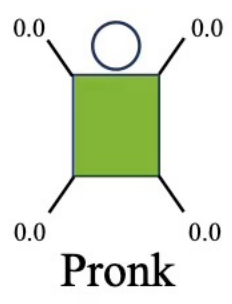
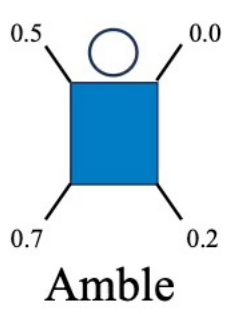
We can transition between any gait and style, at any velocity



- Walk
- Amble
- Trot
- Pace
- Bound
- Pronk
- Canter
- T. Gallop
- R. Gallop

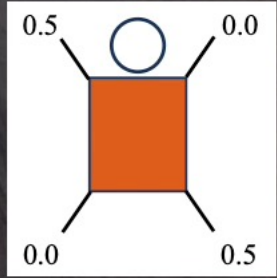
0.5x

0.5x

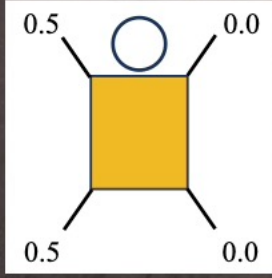


# Sample Gait Transitions

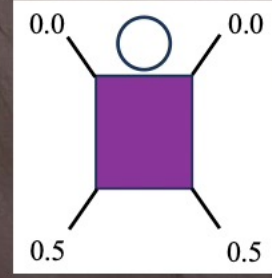
1.0x



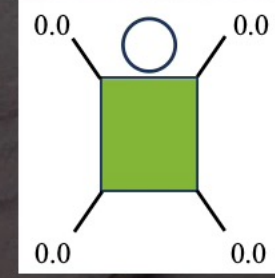
Trot → Pace



Pace → Bound



Bound → Pronk



0.5x

# Same Scenario: Robust to Disabling Leg(s)

1.0x



# Unseen Gaits During Training

1.0x

Trot → 3 Legs in Sync,  
Front Left out of phase

→ 3 Legs in Sync,  
Hind Right out of phase

→ Asymmetrical Trot



0.5x



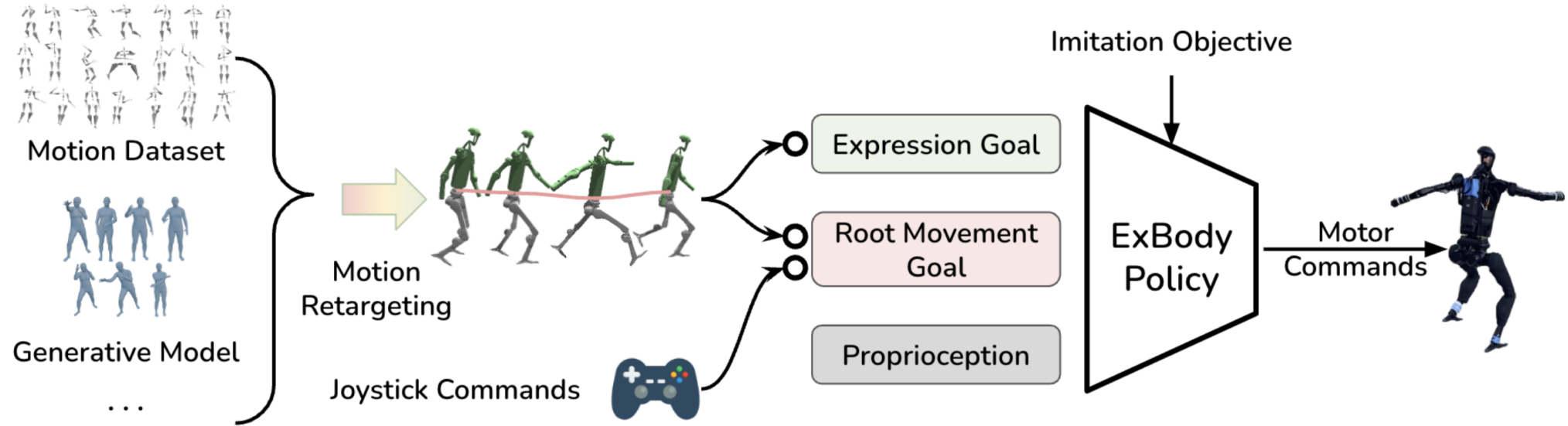
Bellegarda et al. *ICRA 2025*

The policy generalizes to new gaits specified with novel coupling matrices directly in hardware experiments, despite never experiencing such coupling or observations during training.

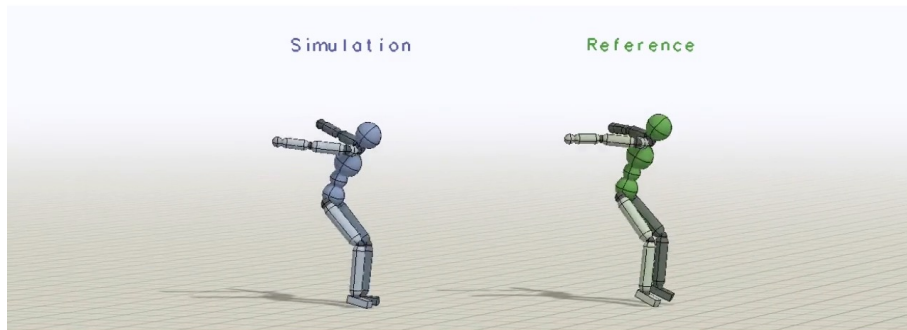
# CPG-RL Summary

- CPG-RL: integrated abstract oscillators into the deep reinforcement learning framework to learn quadruped locomotion
- Online modulation of body height and swing foot height
- Robust sim-to-real transfer (115% nominal mass load, uneven terrain)
- Minimal proprioceptive sensing (only contact Booleans)
- Training without any dynamics randomization or noise
- No reward function tuning!
- Can learn a single policy for Many Quadrupeds!
- Can learn all the gaits, styles, and gait transitions!
- Vision in-the-loop for navigation and gap crossing

# Current Humanoid Approaches – (Deep) imitation learning



## RSI + ET (Our Method)



With RSI and ET, the policy learns to perform the flip.



Peng et al. *DeepMimic: Example-Guided Deep Reinforcement Learning of Physics-Based Character Skills*. SIGGRAPH 2018

Cheng et al. *Expressive Whole-Body Control for Humanoid Robots*. RSS 2024

Chen et al. *GMT: General Motion Tracking for Humanoid Whole-Body Control*. Arxiv 2025

# Conclusion

- Trajectory optimization and deep reinforcement learning can both create locomotion controllers
- Both have pros and cons:
  - TO: model-based (explainable) and verifiable, but need a lot of online compute, and may not generalize to uncertainties
  - RL: can learn potentially anything, but “black-box”, and exploration / sim-to-real
- Abstractions and knowledge priors such as kinematics and Jacobians are very beneficial for learning
- The CPG is a powerful tool to combine with RL to simplify and robustify sim-to-real, and can be used to answer scientific questions
- At Inria we have several masters’ projects available, feel free to contact me

# Possible Exam Questions

- How are model-based methods and learning-based methods similar? How are they different?
- When might you consider using a model-based controller over a learning-based controller, and vice-versa?
- Assume you need to design a locomotion controller with trajectory optimization. What would your cost function include? What kind of constraints do you need?
- Assume you need to train a locomotion controller with deep reinforcement learning. How will you structure the Markov Decision Process (i.e. observation space, action space, reward function)?
- What are some benefits of integrating CPGs with reinforcement learning?
- What makes reinforcement learning challenging? What might you consider when applying reinforcement learning to a new problem?