

# A Policy training

## 1. Teacher policy

Privileged information  $d_t$

- mass  $m$
- COM  $p_b$
- strength  $[k_p, k_d]$
- terrain height  $H$

MLP encoder  $\phi$

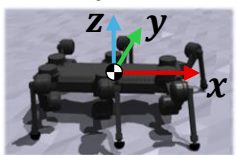
$z_t$

MLP  $\pi_b$

$a_t$

## 2. Student policy

Proprioceptive history  $o_{t-1} \dots o_{t-k}$



TCN encoder  $\sigma$

$\hat{z}_t$

MLP  $\pi_b$

$[\dot{x}_{com}, \varphi]$

## 3. Estimation policy

Observations  $o_t$

- commands  $r$
- orientation  $[\theta, \omega_b]$
- joint states  $[q_j, \dot{q}_j]$
- action history  $a_t^{last}$

MLP  $\pi_e$

$\hat{e}_t$

$[\varphi, \mu]$

Simulation feedback  $e_t$

# B Whole-body follower

## Hard constraints

1. Dynamic consistency
2. Kinematic limits
3. Torque limits

## Task

4. Joint space motion tracking

## Soft constraints

5. Contact motion constraints
6. F-T interaction constraints
7. Body space stabilization

Multi-Priority

Highest  
Lowest

# D Deployment architecture

A2 Student policy

$a_t$

A3 Estimation policy

$\hat{e}_t$

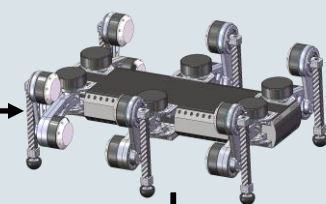
B Full-constrained Whole-body follower

$T$

C Hierarchical Optimization

$\tau_j$

$q_j, v_j$



# C Hierarchical optimization

## Problem formulation

$$\min. \quad \frac{1}{2} \xi_{p+1}^T \mathbf{H}_{p+1} \xi_{p+1} + \mathbf{c}_{p+1}^T \xi_{p+1}$$

$$\text{s.t.} \quad \tilde{\mathbf{D}}_{p+1} \xi_{p+1} \leq \tilde{\mathbf{f}}_{p+1}$$

optimize

## Optimal decision variables

$$x = [\ddot{q}_j^T, F_{\text{grf}}^T, \tau_j^T]^T$$

desired joint torque