

Privileged Knowledge Distillation and Hierarchical Framework **in Reinforcement Learning**

CS586 – Robot Motion Planning and Applications

Speaker: Taegeun Yang

2025.04.09

Content

Recap: Reinforcement Learning

Privileged Knowledge Distillation

- Concept
- Teacher-Student Framework
- Regularized Online Adaptation (ROA)

Hierarchical Reinforcement Learning (HRL)

- Concept
- Related Works – *Manipulation, Locomotion, Navigation*

Our Work

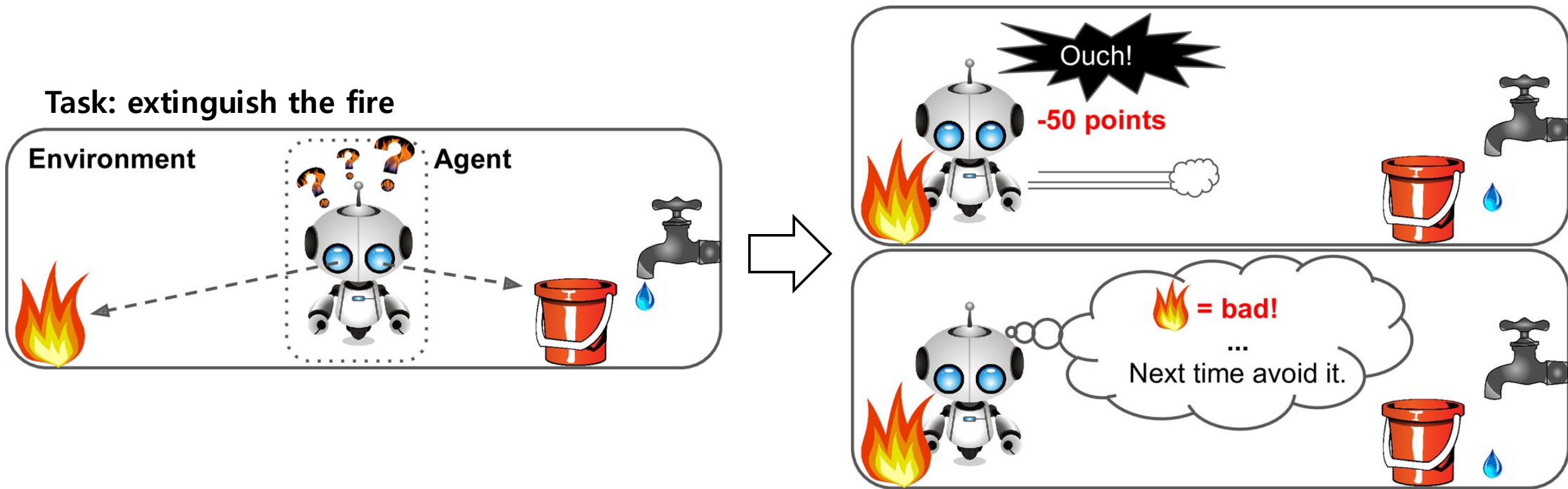
- Efficient Navigation Among Movable Obstacles using a Mobile Manipulator via Hierarchical Policy Learning

Reinforcement Learning (RL)

Reinforcement Learning (RL)

Agent learns to make decisions by interacting with an environment

- Trial and error interactions with an environment



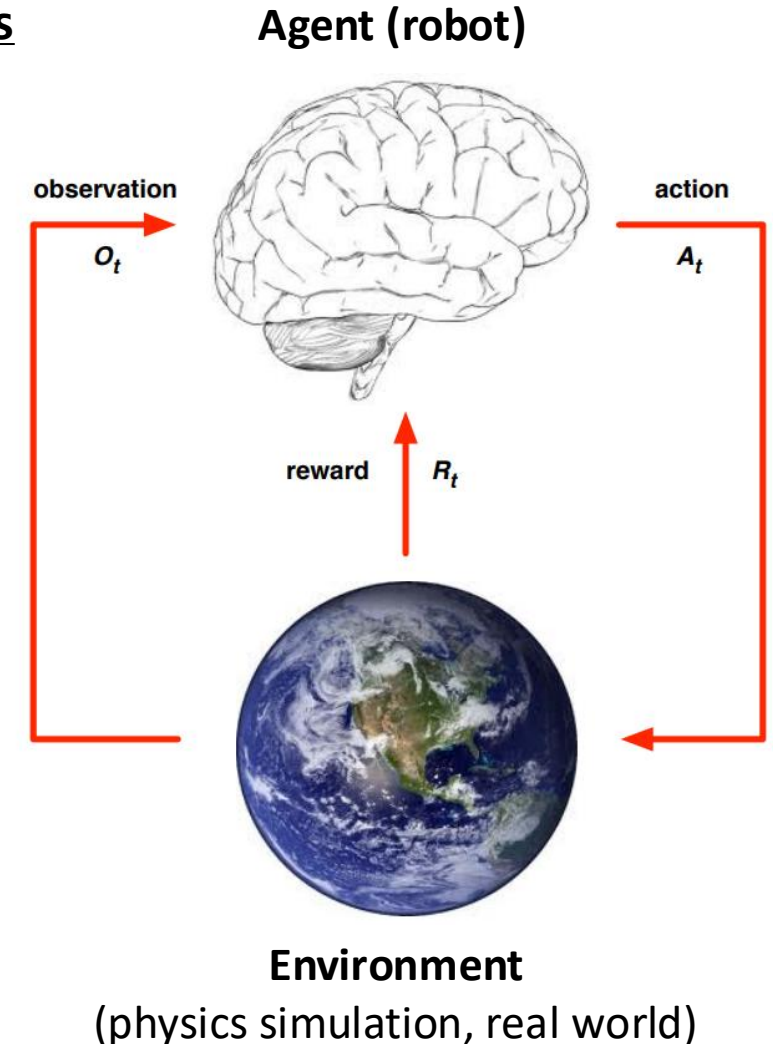
Reinforcement Learning (RL)

Agent learns to make decisions by interacting with an environments

- Trial and error interactions with an environment

At each time step t

- **Agent:**
 - Receives observation O_t (*e.g., sensing*)
 - Executes action A_t (*e.g., move forward*)
 - Receives reward R_t (*task related*)
- **Environment:**
 - Receives action A_t
 - Emits observation O_{t+1}
 - Emits reward R_{t+1}



Reinforcement Learning (RL)

Problem Formulation

- Markov Decision Problem (MDP)
 - $\langle S, A, R, T, \gamma \rangle$
 - S : state space
 - A : action space
 - R : reward function $R: S \times A \times S \rightarrow \mathbb{R}$
 - T : state transition function $T: S \times A \rightarrow S$
 - γ : discount factor $\gamma \in [0,1)$

Objective

- Maximize expected cumulative reward (*return*): $\sum_{t=0}^{\infty} \gamma^t r_t$
 - Policy $\pi_{\theta}: S \rightarrow A$
 - $\theta^* = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{a \sim \pi_{\theta}(s), s' \sim T(s,a)} [\sum_{t=0}^{\infty} \gamma^t R(s, a, s')]$

Reinforcement Learning (RL)

Example 1: Atari Example

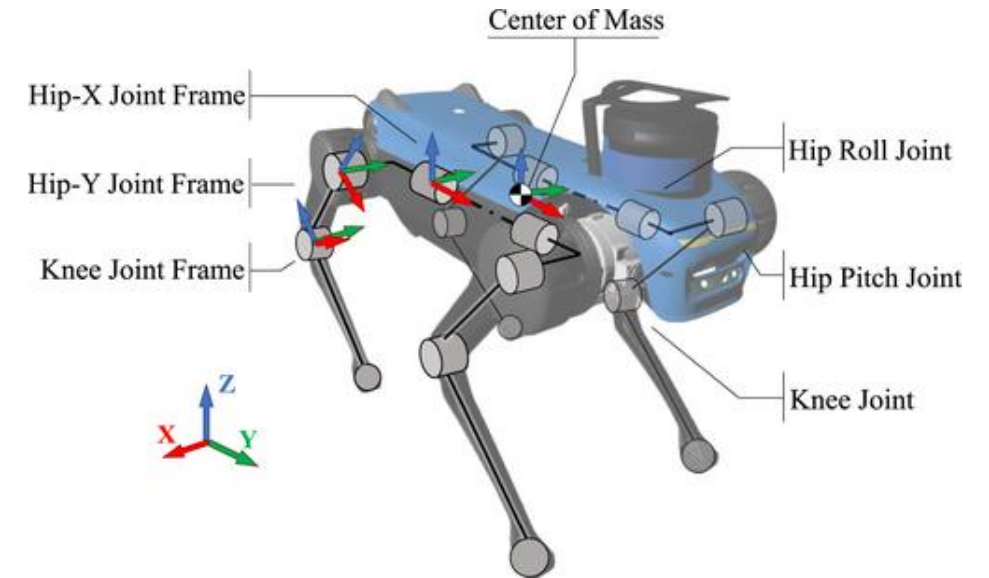
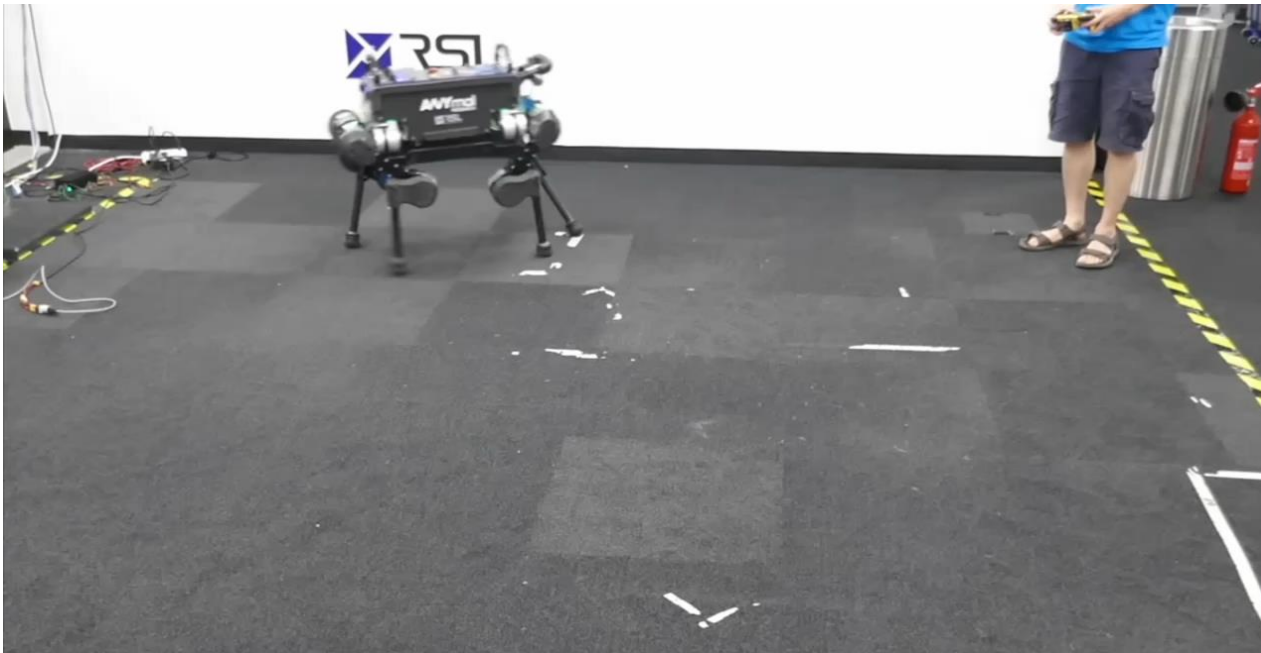
- S : state space – image
- A : action space – move left/right
- R : reward function
 - Breaking a block: +1
 - Missing a ball: -1
- T : state transition function
 - How the ball moves
 - Block breaks when hit by ball
- γ : discount factor $\gamma \in [0,1)$
- **Policy** determines where the paddle (agent) moves based on the image input



Reinforcement Learning (RL)

Example 2: Quadruped Robot Locomotion

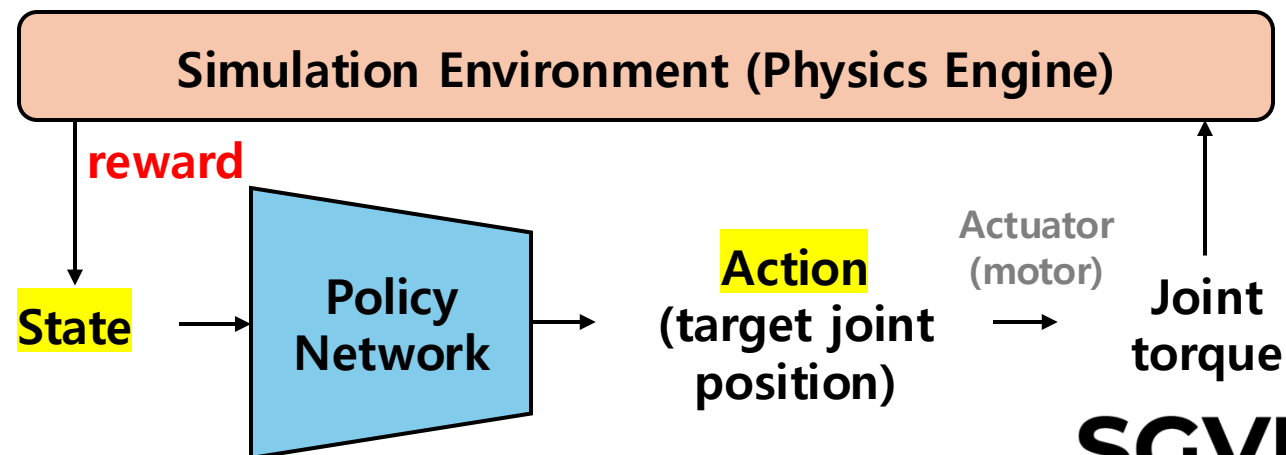
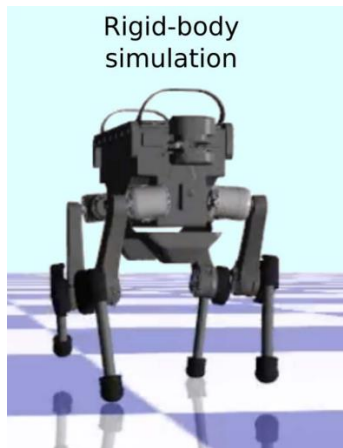
- Given movement command: forward velocity, lateral velocity, and angular velocity (yaw)
- Policy output: Action – target joint position (*12-dim*)



Reinforcement Learning (RL)

Example 2: Quadruped Robot Locomotion

- Given movement command: forward velocity, lateral velocity, and angular velocity (yaw)
- Policy *input* : **State** – current robot state, target movement command, etc.
- Policy *output* : **Action** – target joint position (12-dim)
- **Reward** – command tracking, stable locomotion (i.e., avoid falling), etc.
- **Transition** – Physics Engine (dynamics)



Privileged Knowledge Distillation

Privileged Knowledge Distillation

Privileged Knowledge

- Information that is available during training but not during deployment (testing)
- In the context of RL
 - Extra observations or state variables
 - E.g., ground friction, forces (*not be accessible in the real-world*)
- Provide more complete understanding of the environment

Knowledge Distillation

- Training technique - weaker policy learns by mimicking the behavior of stronger policy
student **teacher**
- Strong policy: trained with full state information (*w/ privileged information*)
- Weak policy: trained with partial state information (*w/o privileged information*)

Privileged Knowledge Distillation

Example 1: Teacher-Student Framework

- Train **teacher** policy with privileged information (RL)

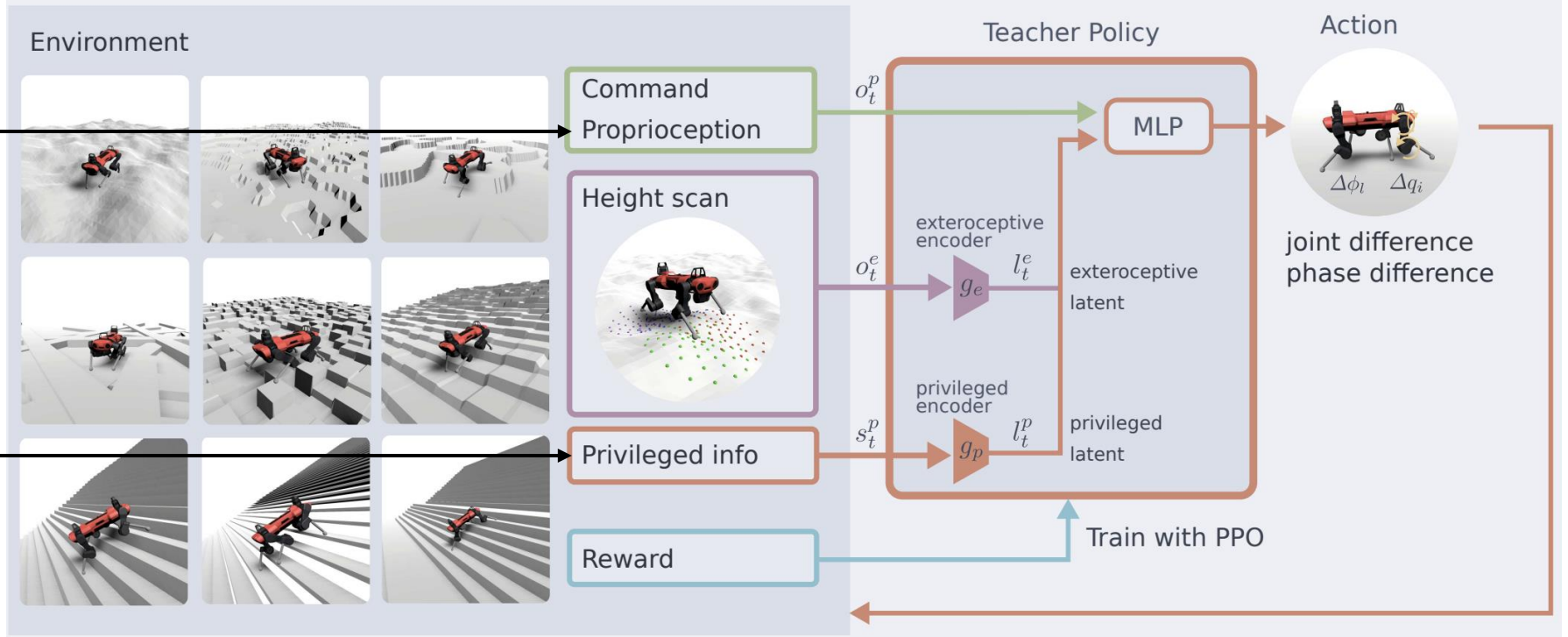
Robot's internal sensory data
- Obtainable

Body velocity
Body orientation
Joint position
Joint velocity
...

Contact states
Contact forces
Contact normal
Friction coef.
...

Unobtainable in real-world

1. Teacher policy training

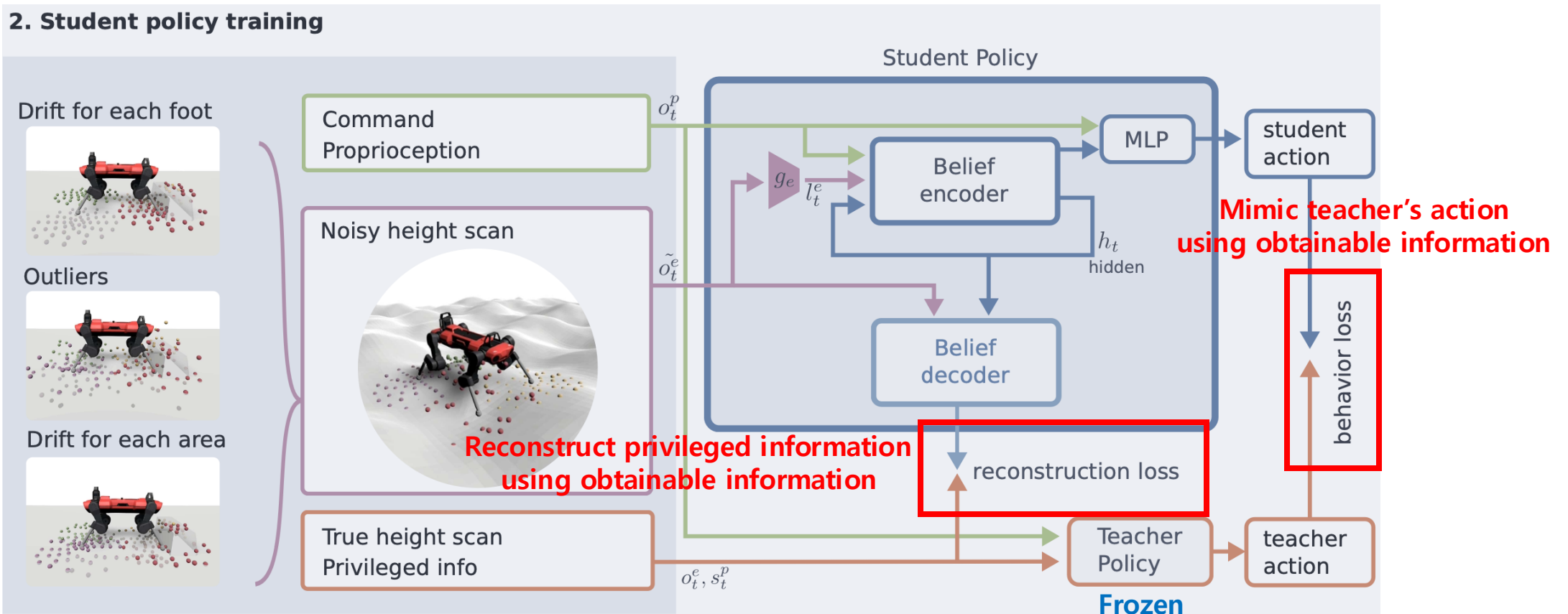


Privileged Knowledge Distillation

Example 1: Teacher-Student Framework

- Train **student** policy without privileged information (**Supervised Learning**)

2. Student policy training



Privileged Knowledge Distillation

Example 1: Teacher-Student Framework

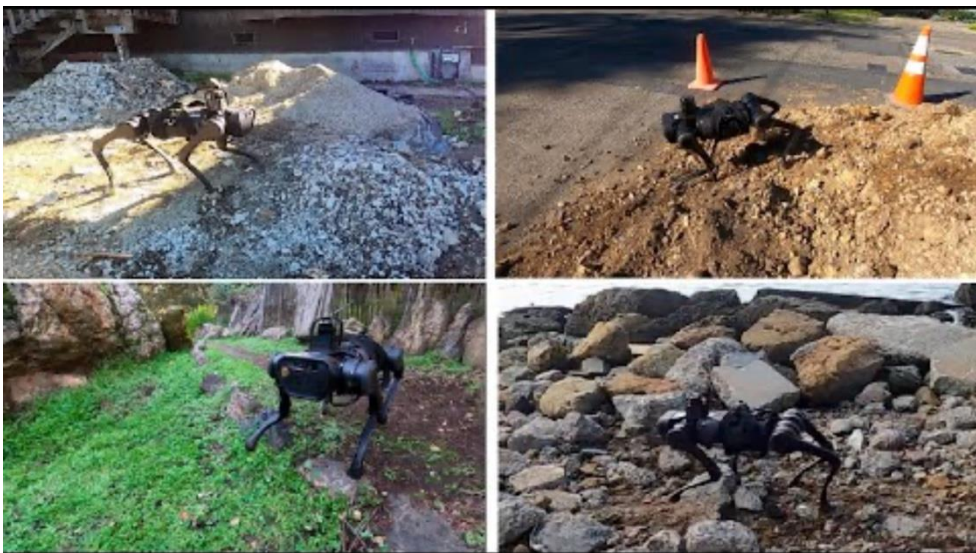
- Deploy using **student** policy (*real-world*)



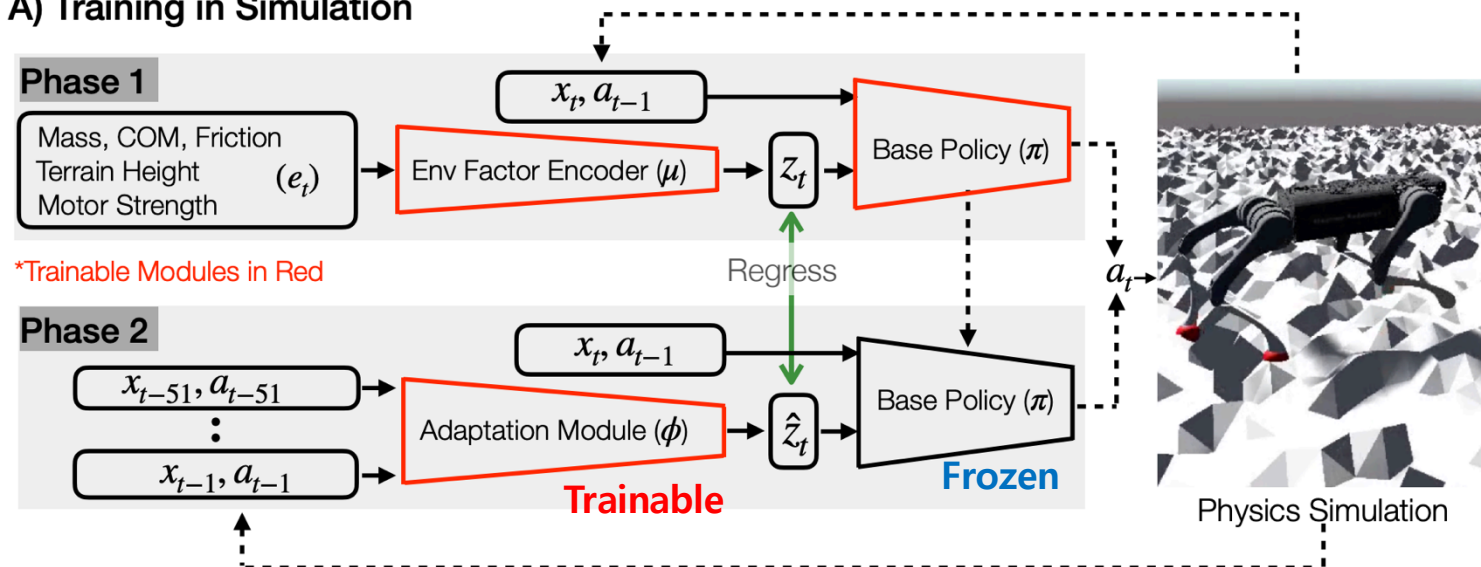
Privileged Knowledge Distillation

Example 2: Adaptation Module

- Estimate encoded privileged information using history of the robot's state
 - Perform robust and adaptive locomotion



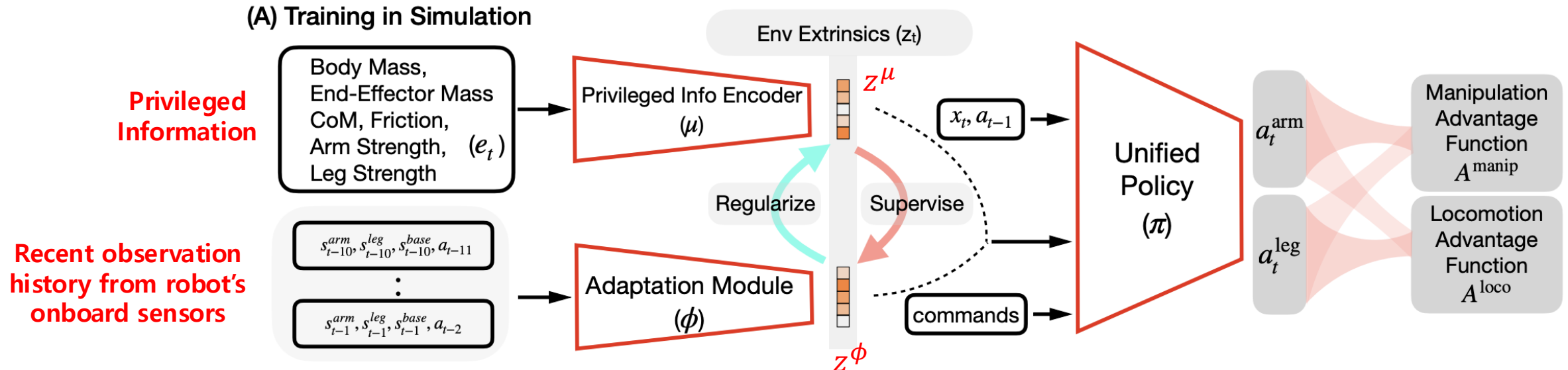
A) Training in Simulation



Privileged Knowledge Distillation

Example 3: Regularized Online Adaptation (ROA)

- 2-step training framework → teacher policy may not provide supervision that student can learn
 - Remove 2-step framework
 - Training Loss: $L(\theta_\pi, \theta_\mu, \theta_\phi) = \underbrace{-J(\theta_\pi, \theta_\mu)}_{\text{RL objective}} + \underbrace{\lambda \|z^\mu - \text{sg}[z^\phi]\|_2}_{\text{Increase from 0 to 1}} + \|\text{sg}[z^\mu] - z^\phi\|_2$ sg: stop gradient
- RL policy uses privileged information

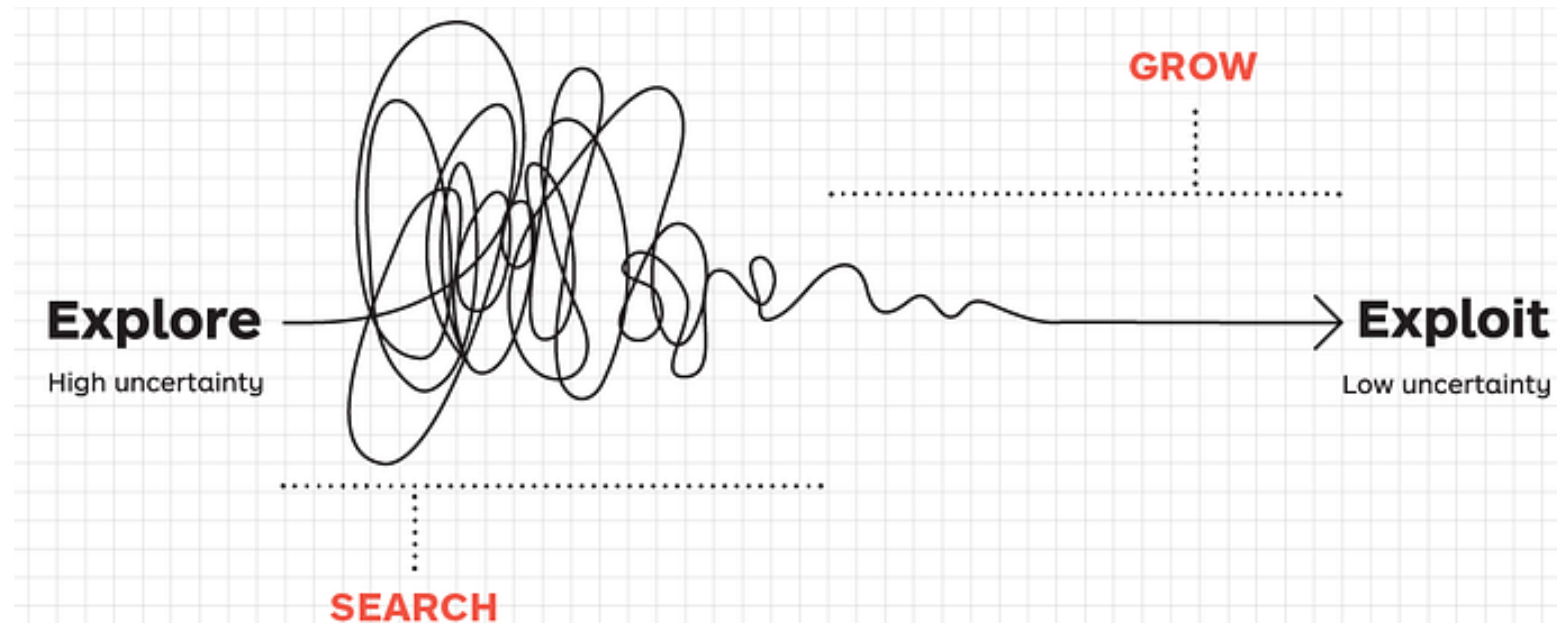


Hierarchical Reinforcement Learning (HRL)

Hierarchical Reinforcement Learning (HRL)

Exploration in Reinforcement Learning (RL)

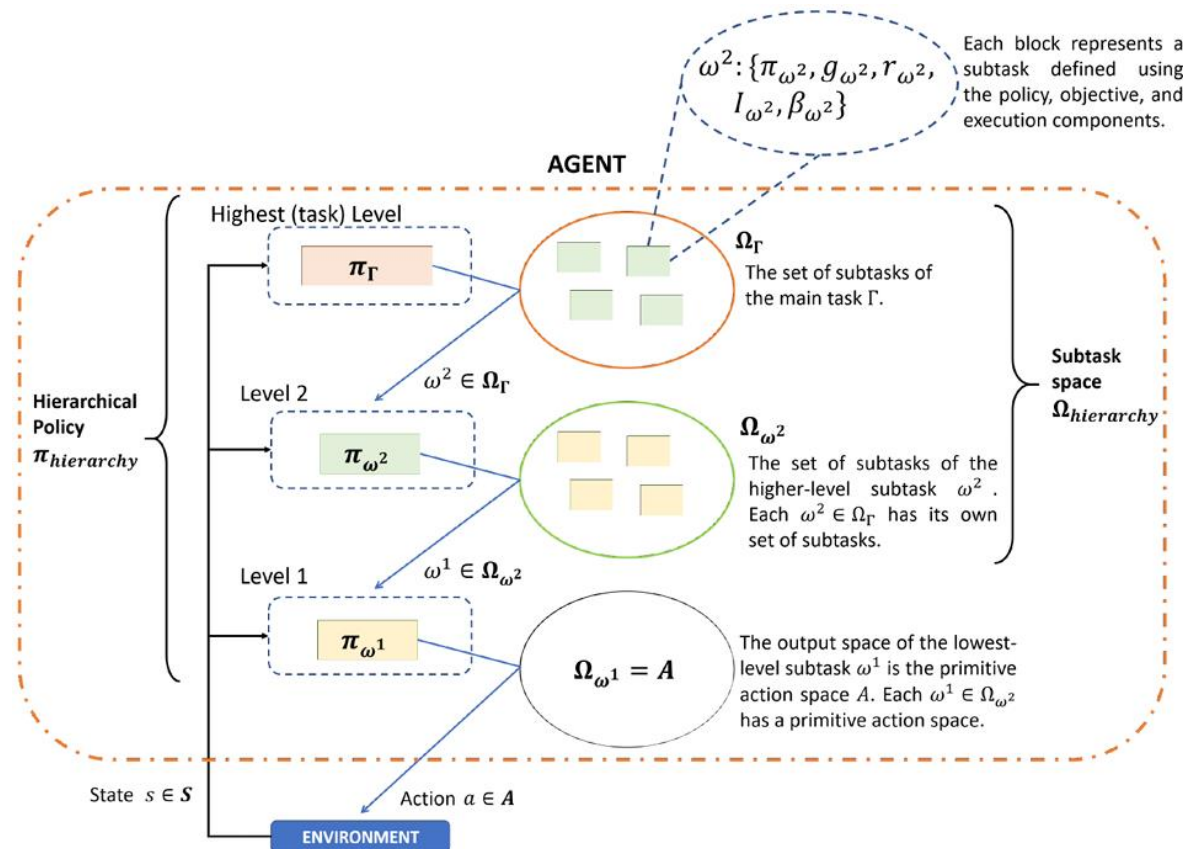
- Essential for learning optimal policies in uncertain environments
 - w/o exploration: agent may become stuck in local optima
- Become more challenging as the observation (*input*) & action (*output*) space expands



Hierarchical Reinforcement Learning (HRL)

Exploration in Reinforcement Learning (RL)

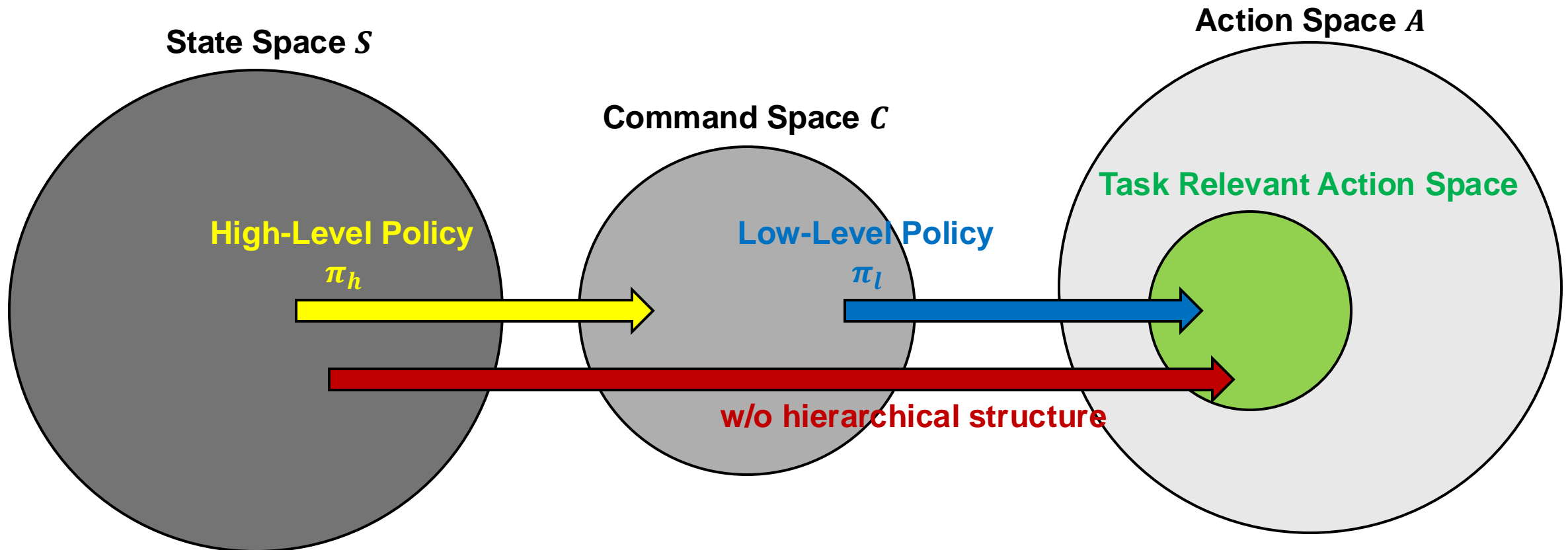
- Observation & Action space $\uparrow \rightarrow$ difficulty of effective exploration \uparrow
- Hierarchical structure
 - Break down a complex problem into multiple sub-problems



Hierarchical Reinforcement Learning (HRL)

Recent Hierarchical Framework in RL

- **High-Level Policy** : decide task-relevant command (*decision making*)
- **Low-Level Policy** : execute given command (*control*)

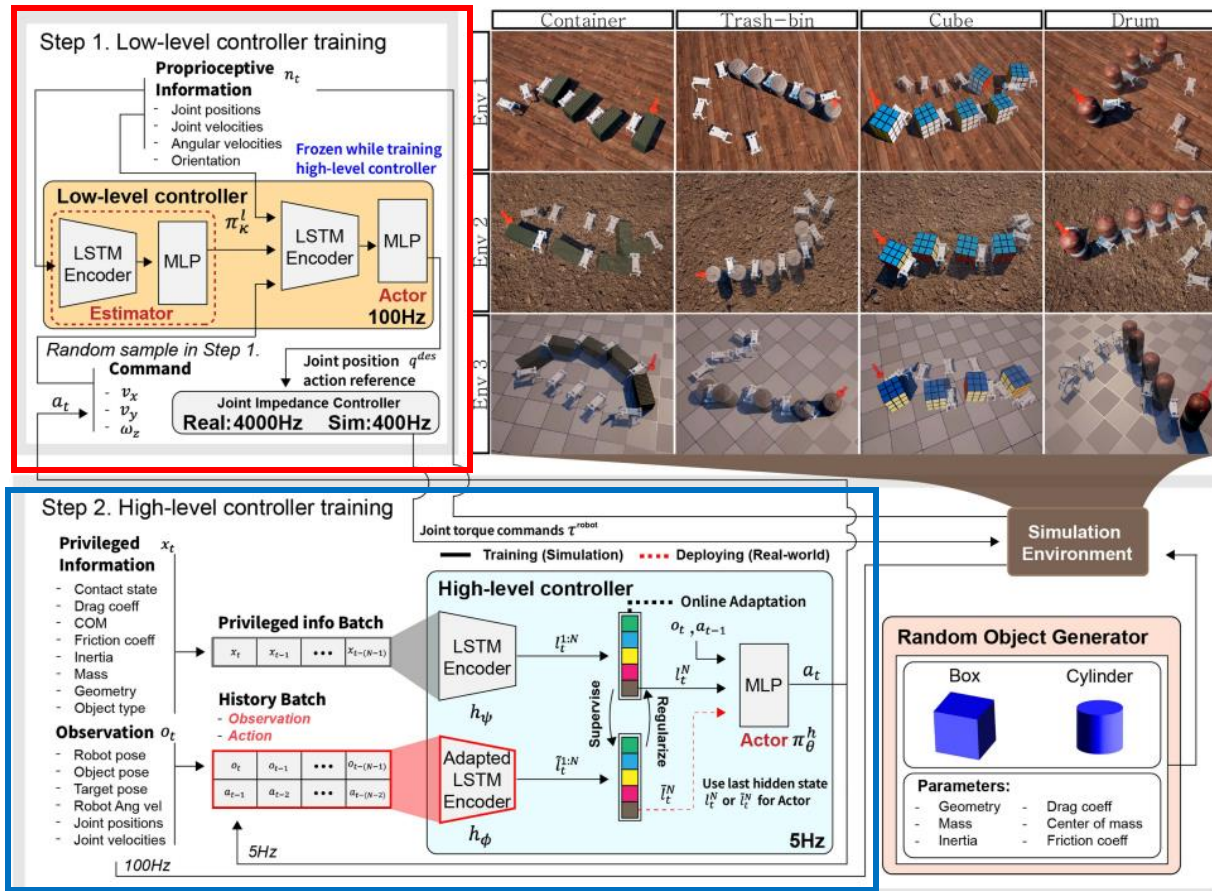


Hierarchical Reinforcement Learning (HRL)

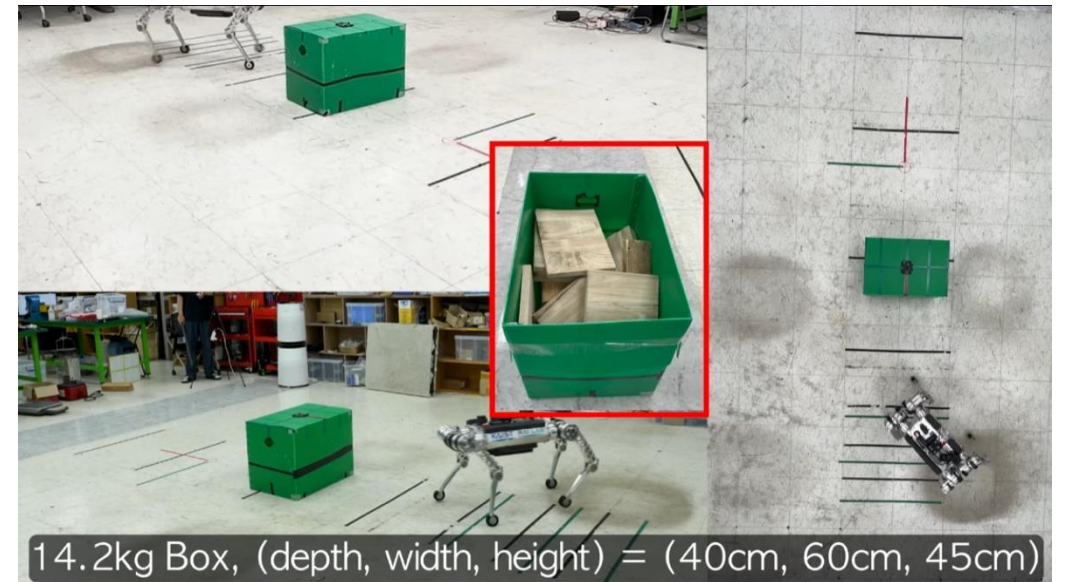
Example 1: Manipulation

- *High-Level Policy*: decide robot's velocity (*command*) to push the object to the target location

- Low-Level** • *Low-Level Policy*: execute given command (*velocity command* \rightarrow *target joint positions*)



* During High-Level training, Low-Level policy keeps frozen



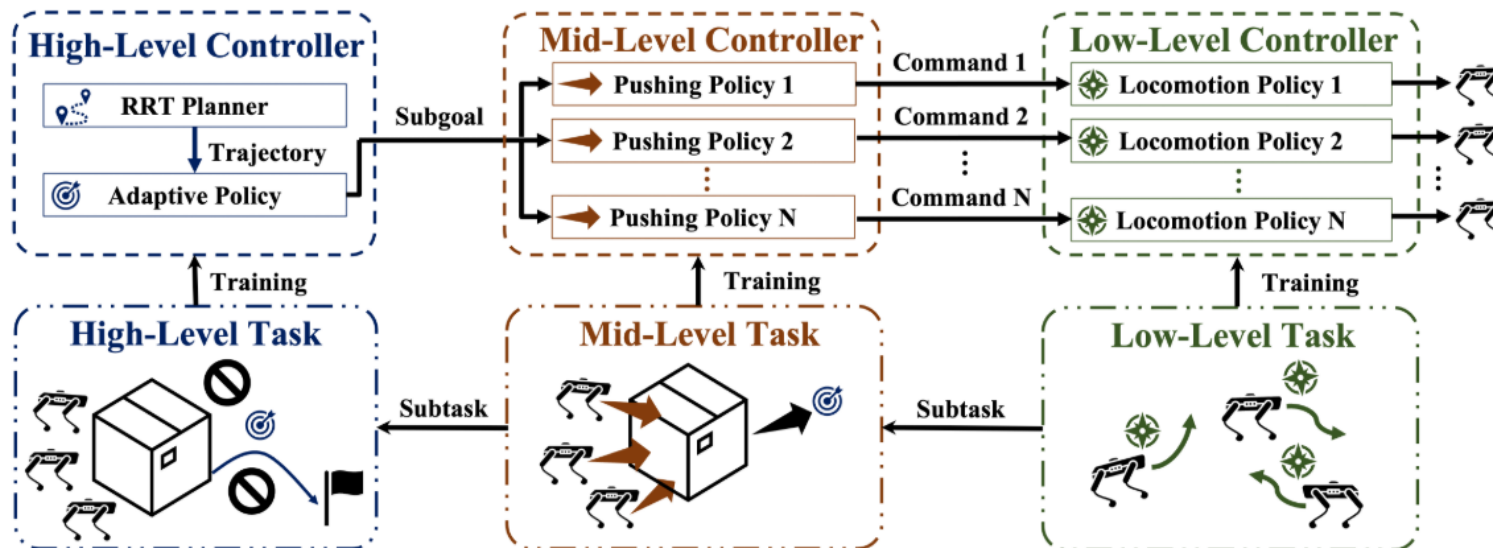
Hierarchical Reinforcement Learning (HRL)

Example 2: Manipulation

- *High-Level Policy*: generate subgoals for the object
- *Mid-Level Policy*: decide robots' velocity (*command*) to push the object to the subgoal
- *Low-Level Policy*: execute given velocity command

3-step training
process

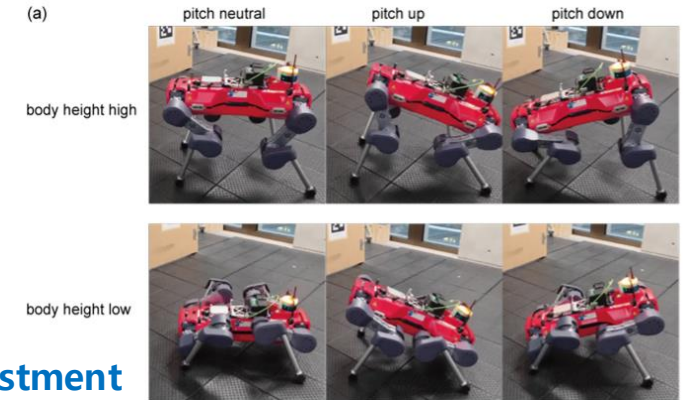
Methodology



Hierarchical Reinforcement Learning (HRL)

Example 3: Locomotion

- Follow given velocity command given by human in confined space
 (v_x, v_y, w_z)
- Requires body adjustments to avoid collision

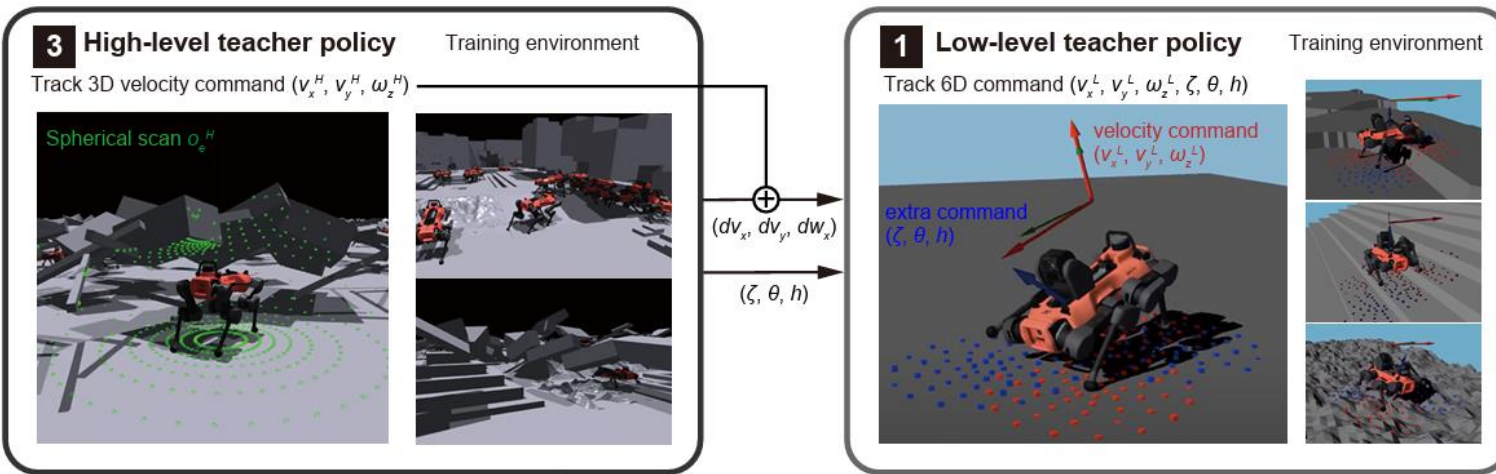


Hierarchical Reinforcement Learning (HRL)

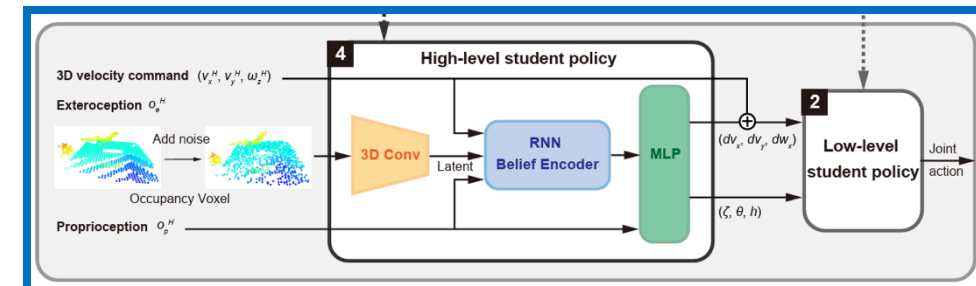
Example 3: Locomotion

Body velocity & Body adjustment motion

- *High-Level Policy*: generate 6-dim command \rightarrow avoid collision, track 3-dim velocity command
- *Low-Level Policy*: execute given 6-dim command (*command* \rightarrow *joint action*)



Hierarchical Structure



Knowledge Distillation

Hierarchical Reinforcement Learning (HRL)

Example 4: Navigation

- Autonomous navigation (*map* \rightarrow *planning* \rightarrow *follow*)
- Sensor data (*e.g., LiDAR, Camera*) \rightarrow Joint action (*avoid collision, path following, etc.*) : **Hard to learn**

A Workflow

i. Scan the environment



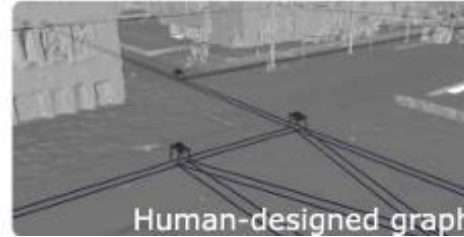
Laser scanner

ii. Process point cloud



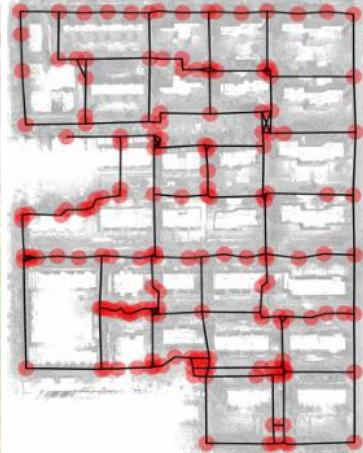
Collection time: ~90 min.

iii. Create navigation graph



Human-designed graph

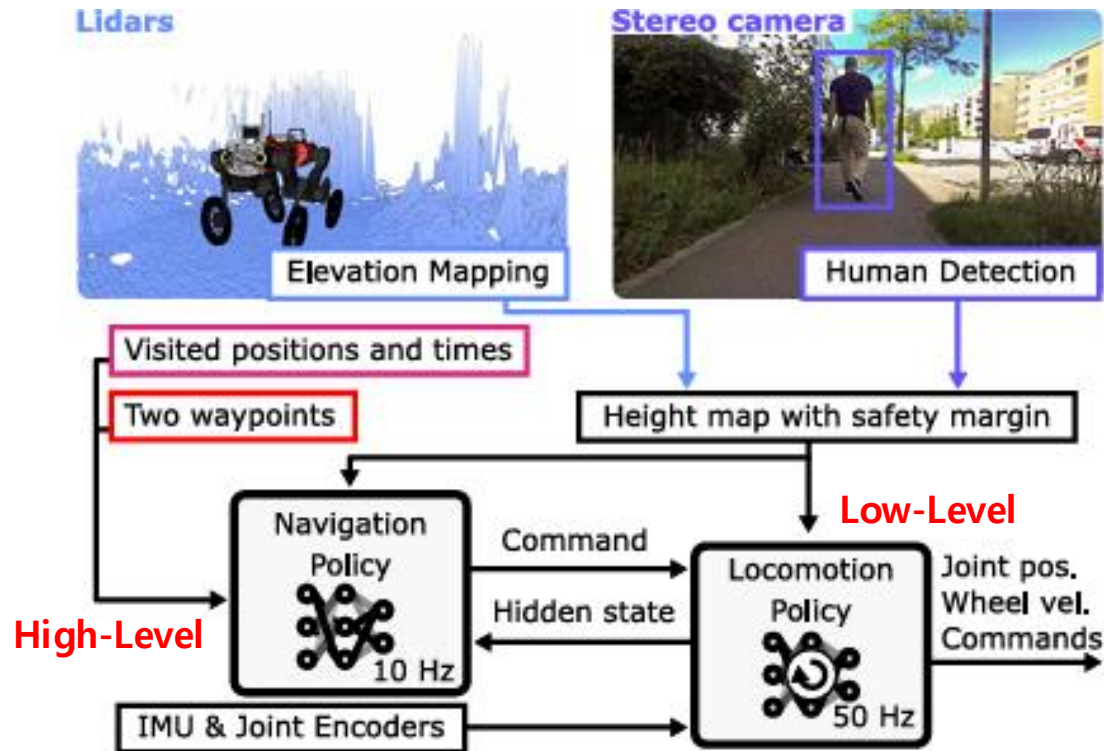
iv. Path planning & Follow



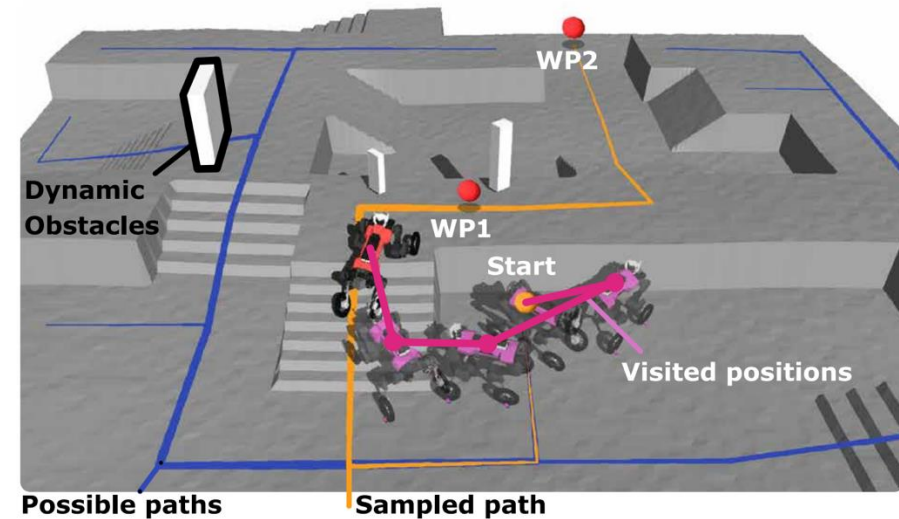
Hierarchical Reinforcement Learning (HRL)

Example 4: Navigation

- *High-Level Policy*: consider surrounding environment & path → generate velocity command
 (v_x, v_y, w_z)
- *Low-Level Policy*: execute given command (*velocity command* → *target joint position/velocity*)



C Training Environment



Thanks for your attention

Any question will be welcome

CS586 – Robot Motion Planning and Applications

Speaker: Taegeun Yang

2025.04.09