

Cooperative Multi-Robot Navigation in Dynamic Environment with Deep Reinforcement Learning

Seminar on Artificial Intelligence II

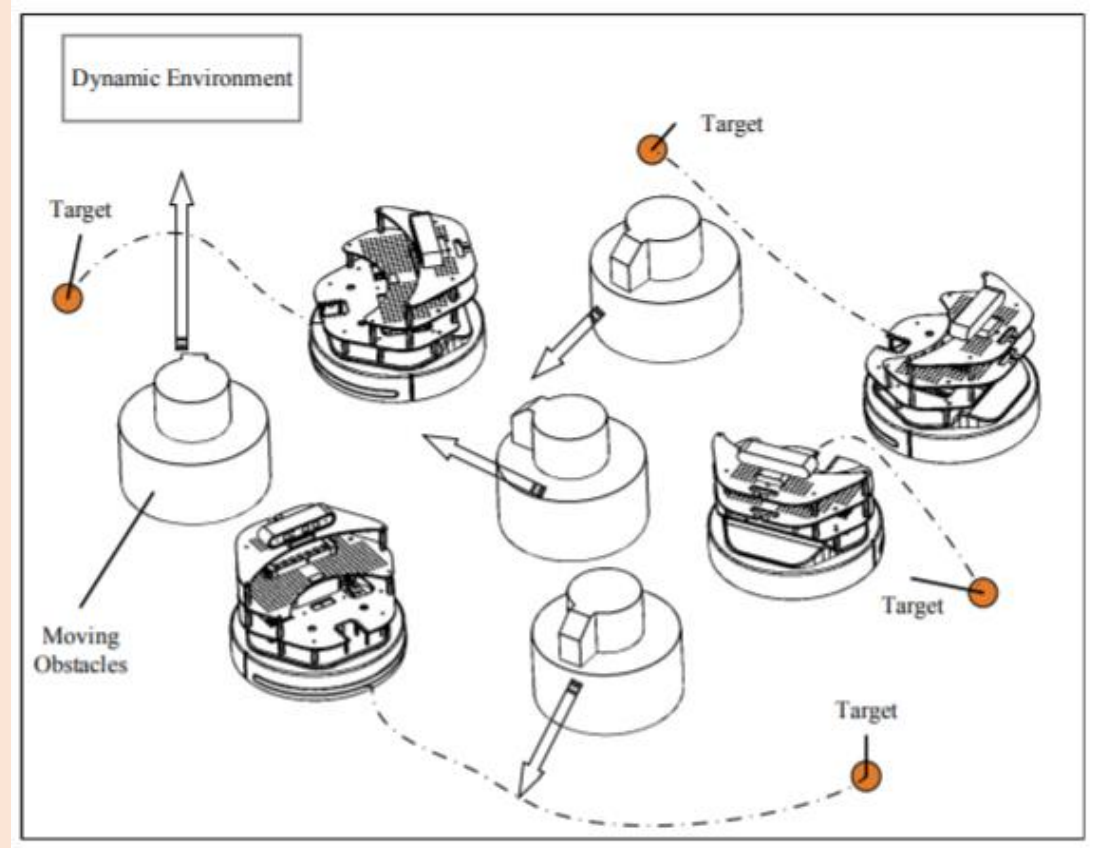
Gabriela Suchopárová, 16. 3. 2022

Outline

1. Introduction
2. Cooperation framework design
3. Deep reinforcement learning model
4. Target location allocation algorithm
5. How to solve the transfer to real world
6. Experiments

Multi-robot navigation problem

- N robots, M obstacles
 - Obstacles **move** as well
- Dynamic
- Partially observable
- **Multiple targets**
 - Robots can go to any target
 - Goal allocation
- Simulation and **Real World**
- Task: learn the **navigation policy**

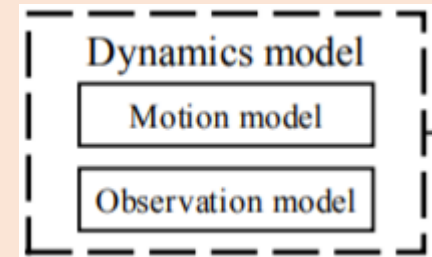
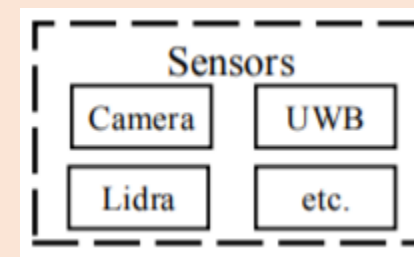
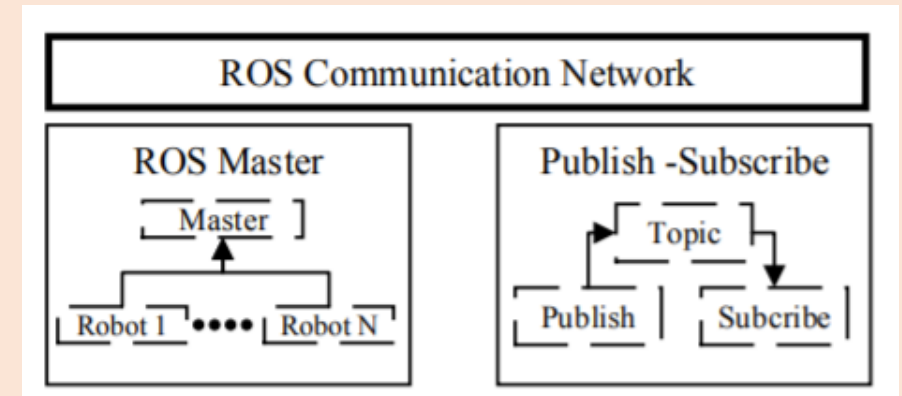


Challenges

- Efficient **target** location **allocation**
 - Allocate the goals fast
 - Reduce the total travel time
- Robot **cooperation**
 - How to combine the experience of all robots
- **Transfer** from **simulation** to **real world**
 - Noise in observations from sensors
 - Noise in motion when applying actions
 - The noise parameters differ across different scenarios (sensor type etc.)

System framework of multi-robot navigation

- They **share one neural network** and policy
 - The network is trained using input from all robots
- Agents communicate through ROS
 - They **share observations**
 - obstacle positions, their position,...
 - They **receive actions**
- Agents get information from **sensors**
- Agents have a **dynamics model**



Problem definition

- POMDP

- States

- **Robot** states: $\mathbf{s}_r^t = [p_{rx}^t, p_{ry}^t, \theta_r^t, v_{rx}^t, v_{ry}^t]$ (position, orientation, velocity)

- **Obstacle** states: $\mathbf{s}_o^t = [p_{ox}^t, p_{oy}^t, v_{ox}^t, v_{oy}^t, r_o^t]$ (position, velocity, radius)

- **Target** positions: $\mathbf{s}_g = [p_{gx}, p_{gy}]$

- Actions: $\mathbf{a}^t = [v_t^t, v_r^t] \sim \pi(\mathbf{a}^t | \mathbf{s}_r^t, \mathbf{s}_o^t, \mathbf{s}_g)$

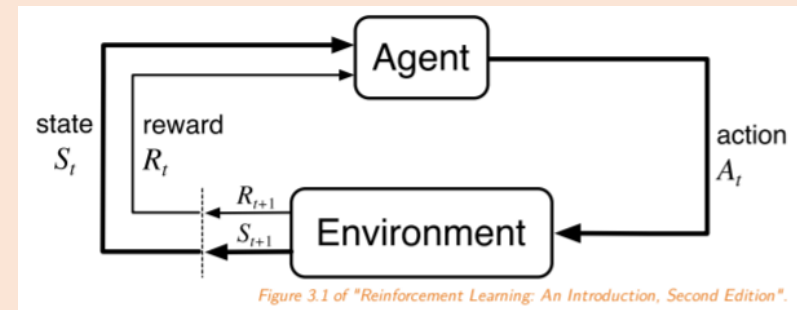
- Objective:

- **Minimize travel time** of all agents **to goals** while **avoiding collisions**

$$\begin{aligned} & \underset{\pi_\theta}{\operatorname{argmin}} \quad \mathbb{E} [T | \pi_\theta, \mathbf{s}_{r,1:N}, \mathbf{s}_{o,1:M}, \mathbf{s}_{g,1:N}] \\ & \text{s.t.} \quad \forall i, j \in [1, N], k \in [1, M] \\ & \quad \quad d_{rr,i,j} > 2r_r \\ & \quad \quad d_{ro,i,k} > r_r + r_o \\ & \quad \quad d_{g,i} < d_{\min} \end{aligned}$$

RL quick overview

- Estimate $v(s)$ or $q(s,a)$... average return from state s (starting with a)
 - $V(s) = \mathbb{E}[G_t | s_t = s] = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s_t = s]$
- Bellman equation
 - $V^\pi(s) = \sum_a \pi(s, a) \sum_{s'} P(s', r | s, a) (r + \gamma V^\pi(s'))$
- The equation is estimated, since we don't know the probability model
 - In POMDP, we even don't know the states
 - We also want to learn the policy π



1. Deep reinforcement learning framework

- **State** space: $\mathbf{s}^t = [\mathbf{s}_r^t, \tilde{\mathbf{s}}_r^t, \mathbf{s}_o^t, \mathbf{s}_g^t]$
 - Robot state, other robot states, obstacle state, *allocated* goal state
- **Action** space: $\mathbf{a}^t = [v_t^t, v_r^t]$
 - Velocities **clipped** to some range (hyperparameter)
 - The reason is limited obstacle detection speed
- **Reward** of robot i : $r_i^t = r_{g,i}^t + r_{c,i}^t$

$$r_{g,i}^t = \begin{cases} 5 & \text{if } d_g^t < 0.1 \\ 1 + d_g^t * 0.5 & \text{if } 0.1 < d_g^t < 0.4 \\ 10 * (d_g^{t-1} - d_g^t) & \text{otherwise} \end{cases}$$

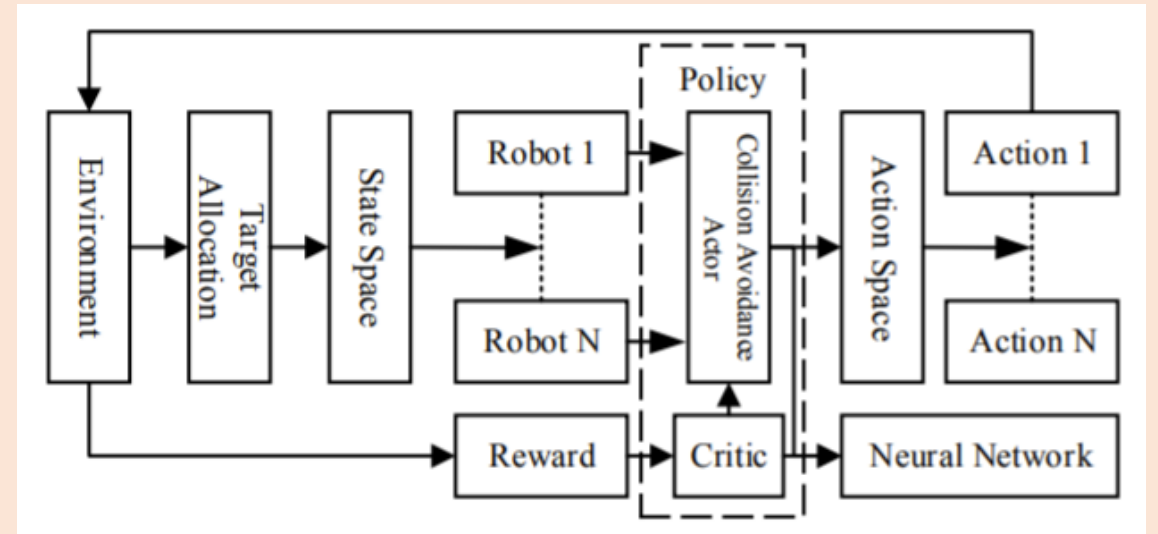
Goal reward

$$r_{c,i}^t = \begin{cases} -10 & \text{if } d_{rr}^t < 2 * r_r \\ -10 & \text{if } d_{ro}^t < r_r + r_o \\ 0 & \text{otherwise} \end{cases}$$

Collision avoidance reward

1. Deep reinforcement learning framework

- Model used – **PPO**
- Temporal-difference method
- Actor-critic – 2 dense networks
 - Actor – outputs **actions**
 - Critic – **estimates** the **value** function
- Predicted actions are **continuous**
 - Transitional and rotational velocities are modelled
 - Network outputs mean and standard deviation (Gaussian distribution)
 - The **action** is **sampled** from these distributions



2. Target location allocation

- *Why?* – While avoiding obstacles, the **closest target** may **change**
- For N robots and N targets, allocate targets such that the total **distance is minimal**
 - Allocation is a permutation of robots
- Presented algorithm – basically just **check all permutations** and select the best one
- Run only **after T_a steps**
 - T_a is a hyperparameter
 - Chosen according to no. of agents & computing power

3. Dynamics randomization

- Transfer from simulation to the real world
- Real sensors have noise that can be different for every device
- We want to **avoid retraining** the model on new data
- However, not retraining could lead to worse performance
- Solution: **add noise** during **training**
 - The noise stays the same during the episode
 - Uniformly selected from a predefined range
 - All noises are $N(0, \xi_i)$
 - At each step, we sample from it

- Noise in transitional velocity, ξ_1
- Noise in rotational velocity, ξ_2
- Noise in the position (coordinates) of robots, ξ_3
- Noise in position (coordinates) of obstacles, ξ_4
- Noise in the measurements of obstacles ξ_5
- Mass of robots m_r

Full training cycle

- Classic DRL setting
 - Simulate multiple episodes
- For every episode, **set dynamics**
- Train
 - Every few steps, **allocate goals**
 - Sample actions using the **actor**
 - Collect **next state, reward**
- After T steps, compute **losses**

Algorithm 2 Policy Training with PPO

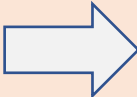

- 1: Initialize neural network π_θ
 - 2: **for** episode=1,2,... **do**
 - 3: Reset the environment with the initial state, \mathbf{s}_{init}
 - 4: Sample the dynamics parameters λ from a range γ uniformly, $\lambda \sim \gamma$
 - 5: **for** robot $i=1,2,\dots$ **do**
 - 6: Receive state \mathbf{s}_i^t , select the goal position $\mathbf{s}_{g,i}^t$
 - 7: Add noise, $\hat{\mathbf{s}}_i^t \sim \mathbf{s}_i^t + [\xi_3, \xi_4, \xi_5]$
 - 8: Sample action $\mathbf{a}_i^t \sim \pi_\theta(\mathbf{a}_i^t | \hat{\mathbf{s}}_i^t)$
 - 9: Add noise, $\hat{\mathbf{a}}_i^t \sim \mathbf{a}_i^t + [\xi_1, \xi_2]$
 - 10: Publish $\hat{\mathbf{a}}_i^t$ to robot i
 - 11: Collect state $\hat{\mathbf{s}}_i^t$, reward r_i^t and \mathbf{a}_i^t for T_i time steps
 - 12: Compute advantage estimates $\hat{A}_i^1, \dots, \hat{A}_i^{T_i}$
 - 13: **end for**
 - 14: Optimize surrogate loss $L^{CLIP}(\theta)$ wrt θ , with Adam optimizer and learning rate l_a for K epochs
 - 15: $\theta_{old} \leftarrow \theta$
 - 16: Optimize value loss $L^V(\psi)$ wrt ψ , with Adam optimizer and learning rate l_v for L epochs
 - 17: $\psi_{old} \leftarrow \psi$
 - 18: **end for**
-

Related work

1. DRL models

- Q-learning, DQN, A3C, DDPG, PPO
- The difference is what they estimate (state/action-value function) and how

2. Multi-agent learning

- MADDPG – shared critic, one actor per robot  High hardware cost, needs specific sensors
- SLCAP – similar approach, but PPO
- GA3C-CADRL – one policy for all agents, asynchronous  Asynchronous update inefficient for homogenous robots
- These methods need preallocated targets
- IDRL – allocates targets, but needs static obstacles

3. Real world transfer – either retrain, or complicated dynamics model

Experiments

- Simulation
 - Gym-Gazebo env, Turtlebot robots
 - Obstacles – turtlebots that periodically change velocity
- Three models
 - 1. Without target allocation and dynamics
 - 2. With target allocation
 - 3. Full model
- Real world scenario
 - UWB, Kinect v2 – coordinates, obstacle pos.

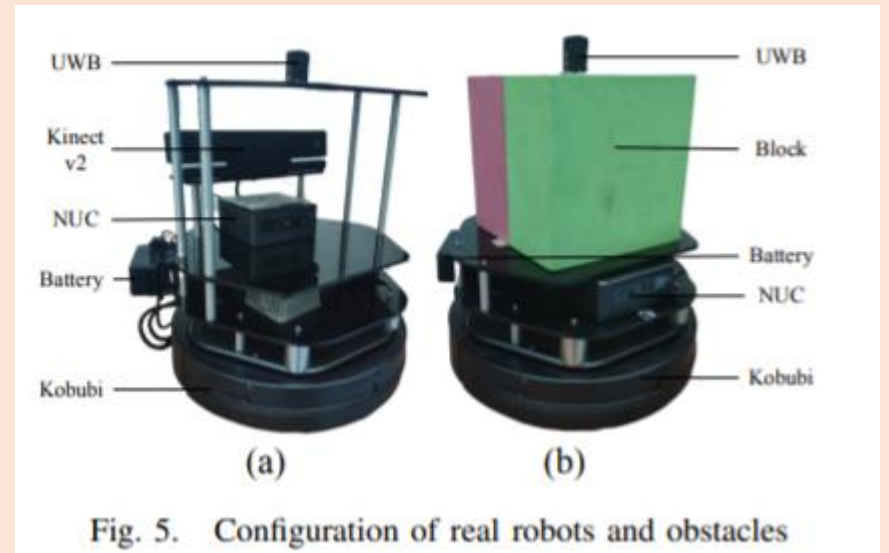


Fig. 5. Configuration of real robots and obstacles

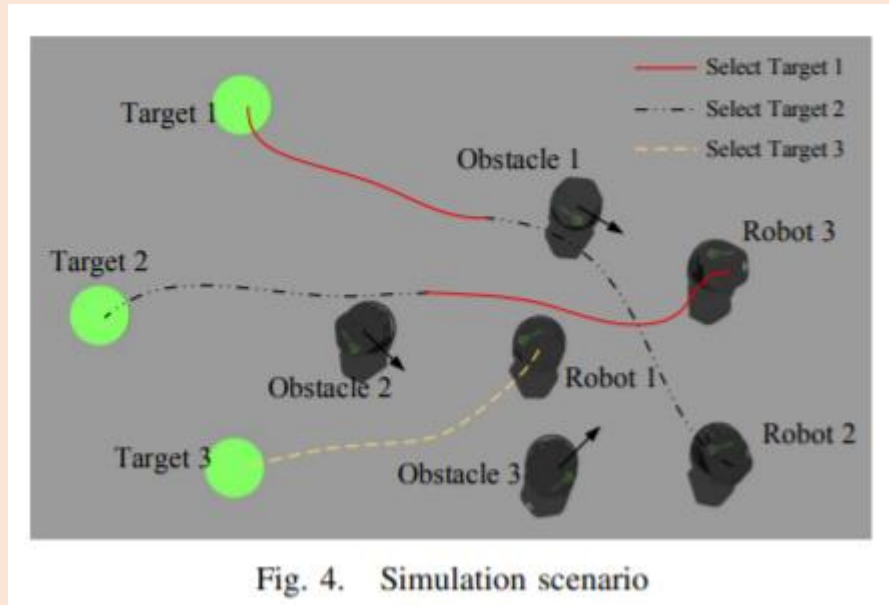
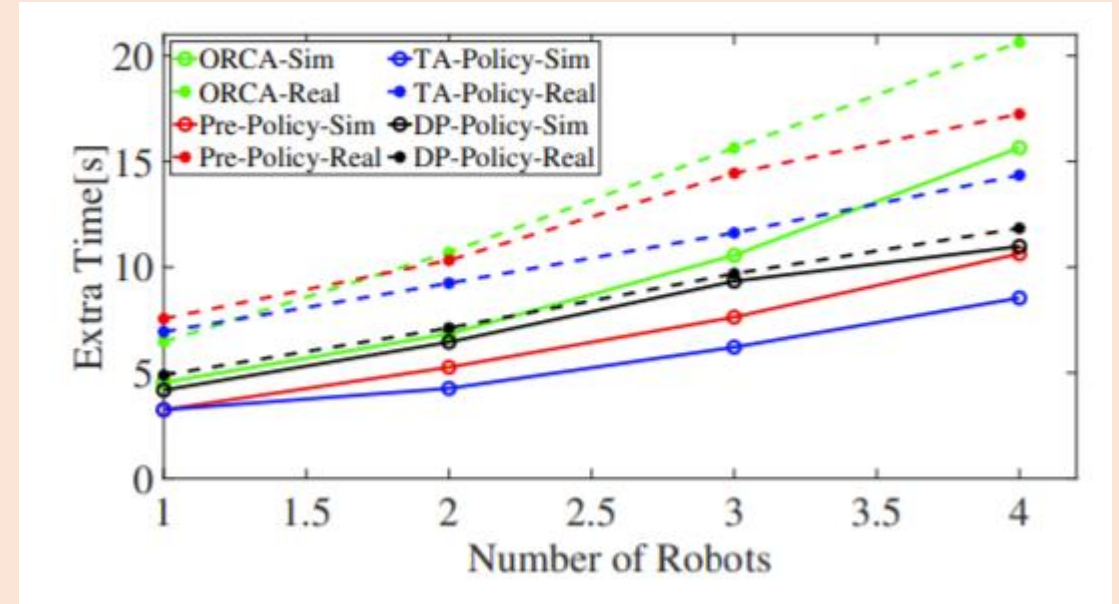
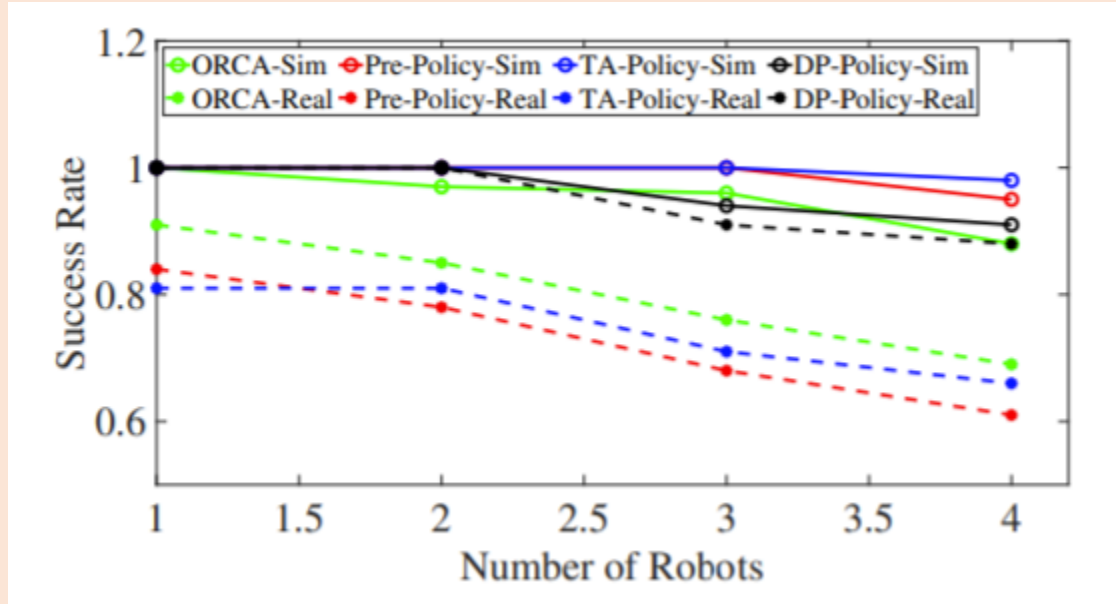


Fig. 4. Simulation scenario

Results

- Baseline – ORCA (computes “allowed” velocities as to not collide)
- Metrics
 - Success rate – # of bots that arrive to the goal
 - Extra time – (avg bot travel time - lower bound on travel)



Summary

- Introduced a **cooperative scheme** for multi-robot navigation
- Collision avoidance
- Closest **target** location **allocation**
- **Dynamics model** for real world transfer
- **Actor-critic** based neural network model
- Two experiments – simulation & real world

Thank you for your attention!

What I think about the paper

- The dynamics (noise) is nice
- Clearly written main results
- The model is not too complicated
- All steps could be easily used in other models and scenarios
- ~~The robots are cute~~
- The models with allocation/dynamics are trained longer
- Only one env/robot type, baseline model
- No confidence intervals in plots
- No code 😞
- The permutation allocation is lame

