

Intent-aware Multi-agent Reinforcement Learning

Siyuan Qi and Song-Chun Zhu

2018 IEEE International Conference on Robotics and Automation(ICRA)

Soutrik Bandyopadhyay

2020EEZ8451



Term Paper Presentation for ELV700

Apr 10, 2023



Table of Contents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

① Introduction

② Intent-aware Multi-agent Reinforcement Learning

③ Library Implementation

④ Simulation Results

⑤ Conclusions



Introduction

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

- In the course we saw how to predict “**intent**”
- After figuring out the intent, how can we go about considering complex tasks ?
 - ① Cooperative tasks
 - ② Safety critical tasks
 - ③ Competitive tasks



On Utility

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

We need a way to encode the inter-relationship between different agents

- How good a particular goal is for a given agent
- Expected long term reward for given agent pursuing the given goal



Example

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

You and your arch-nemesis are selecting universities to study in

Arch-Nemesis

Uni 1
My utility = 2
His utility = 3

Uni 2
My utility = 5
His utility = 7

Uni 3
My utility = 7
His utility = 1

Me

If both of you choose the same uni , a penalty of 10 is doled out



Example

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

Arch-Nemesis

Uni 1
My utility = 2
His utility = 3

Uni 2
My utility = 5
His utility = 7

Uni 3
My utility = 7
His utility = 1

Me

My Action	My belief of other agent		
	1	2	3
1	-10	2	2
2	5	-10	5
3	7	7	-10



Intent aware multi-agent RL

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

- Intrinsic value for the different goals unknown
- Can we learn the intrinsic value from reward data considering the intent of other agents



Table of Contents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.
Results

Conclusions

References

① Introduction

② Intent-aware Multi-agent Reinforcement Learning

③ Library Implementation

④ Simulation Results

⑤ Conclusions



Problem Statement

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

We consider the multi-agent problem with

- n agents
- m goals
- Each agent has an inner-loop controller that brings the agent to the desired goal location
- Objective is to choose the desired goal location based on the intent of the other agents.



Mathematical Notation

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

At each time step t

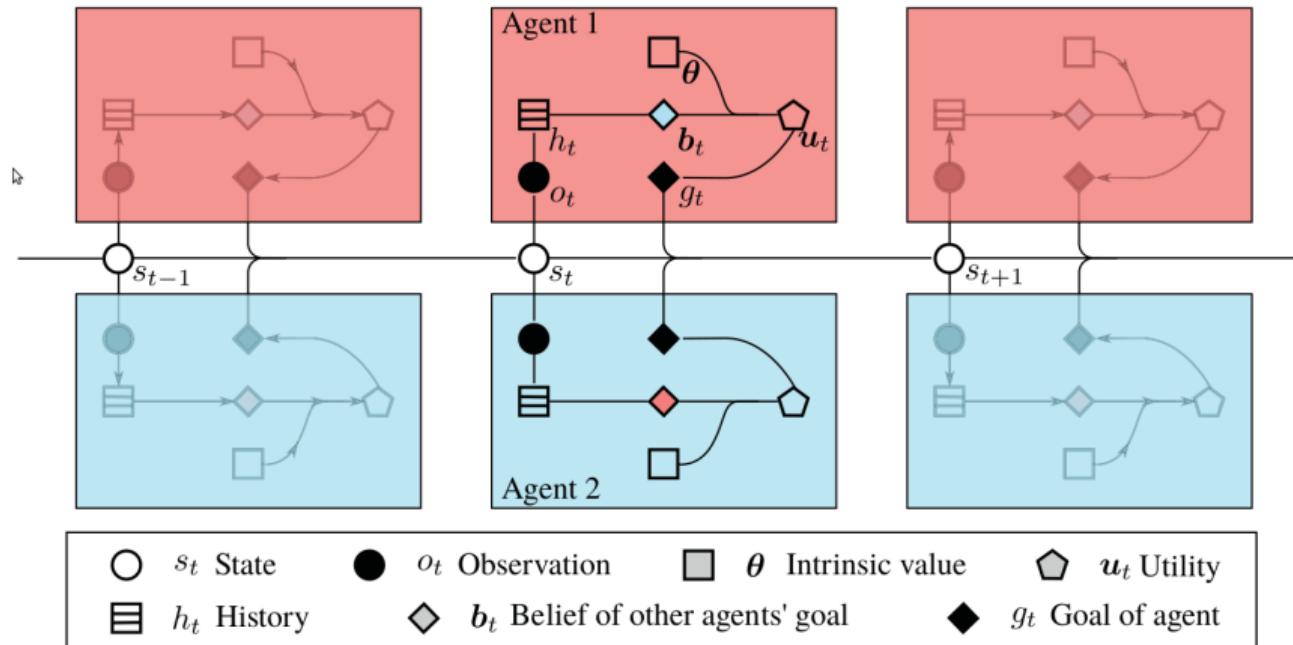
- Agent makes an observation o_t from overall state s_t (unknown)
- Past history of observations h_t
- Based on history agent infers intent of other agents (b_t , probability distribution over agents)
- Compute the desired goal location (g) for each agent based on the intent of other agents.



Pictorial Representation

Intent
aware RL

Intent-aware RL





Utility of the agents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

$$g_i = \operatorname{argmax}_g u(g|b_{-i}; h, \theta) \quad (1)$$

- u denotes the utility function
- g_i denotes the chosen goal for the current agent
- b_{-i} denotes belief of other agents
- h denotes history of observations
- θ denotes the intrinsic value



Example

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

Arch-Nemesis

Uni 1
My utility = 2
His utility = 3

Uni 2
My utility = 5
His utility = 7

Uni 3
My utility = 7
His utility = 1

Me

My Action	My belief of other agent		
	1	2	3
1	-10	2	2
2	5	-10	5
3	7	7	-10



Generalizability of the framework

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

- Safety problem (Avoid the same goal)
- Cooperation for the same goal
- Cooperation for two different goals



Q-function

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

$$u(g_{ik}|b_{-i}; h, \theta) = \sum_j \sum_l \theta_{ik,jl} p(g_{jl}|h) \quad (2)$$

where $\theta = (\theta_{ik,jl}) \in \mathbb{R}^{m \times n \times m}$ is the intrinsic utility (Agent i going to goal k and Agent j going to goal l)

The Q function is defined as

$$Q(h, g_{ik}) = u(g_{ik}|b_{-i}; h, \theta) = \langle \theta, \phi_{ik}(h) \rangle \quad (3)$$

where $\phi_{jl}(h) = p(g_{jl}|h)$ and $p(g_{ik}|h) = 1$ and $p(g_{iq}|h) = 0$ for $q \neq k$



How to compute $p(g|h)$?

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

$$p(g|h) = \frac{p(g)p(h|g)}{p(h)} \quad (4)$$
$$\propto p(g)p(h|g)$$

where $p(h|g)$ is given by

$$p(h|g) = \frac{1}{z} \exp(-\beta d(\Gamma_o, \Gamma_p)) \quad (5)$$

- Γ_o is observed trajectory from history h
- Γ_p is predicted trajectory assuming agent is heading for g

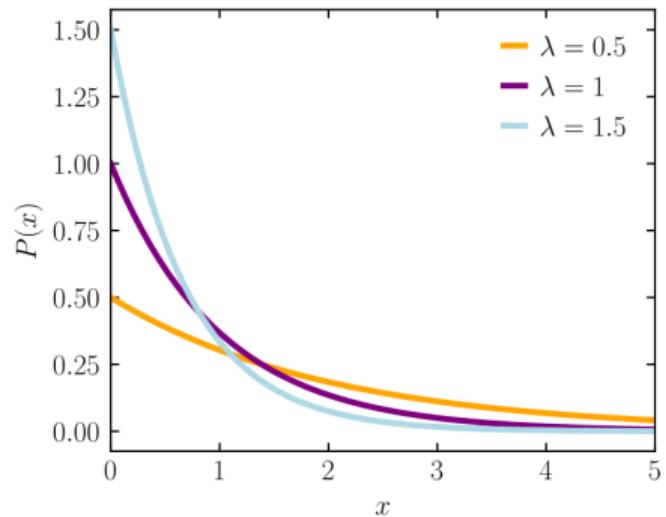


Figure: Boltzmann Distribution



Dynamic Time Warping algorithm

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

How to compute distance between vectors of unequal lengths ?¹

Dynamic Time Warping algorithm

Let Γ_o have length n and Γ_p have length m

$$d(\Gamma_o[:i], \Gamma_p[:j]) = d(\Gamma_o[i], \Gamma_p[j]) + \min(d(\Gamma_o[:i-1], \Gamma_p[:j]), \\ d(\Gamma_o[:i], \Gamma_p[:j-1]), \\ d(\Gamma_o[:i-1], \Gamma_p[:j-1]))$$

corresponding to insertion, deletion, match.

¹ “Dynamic time warping,” *Information Retrieval for Music and Motion*, pp. 69–84, 2007.
DOI: 10.1007/978-3-540-74048-3_4. Slide 17/35



Example

Intent aware RL

S. Bandyopadhyay

Introduction

Intent-aware RL

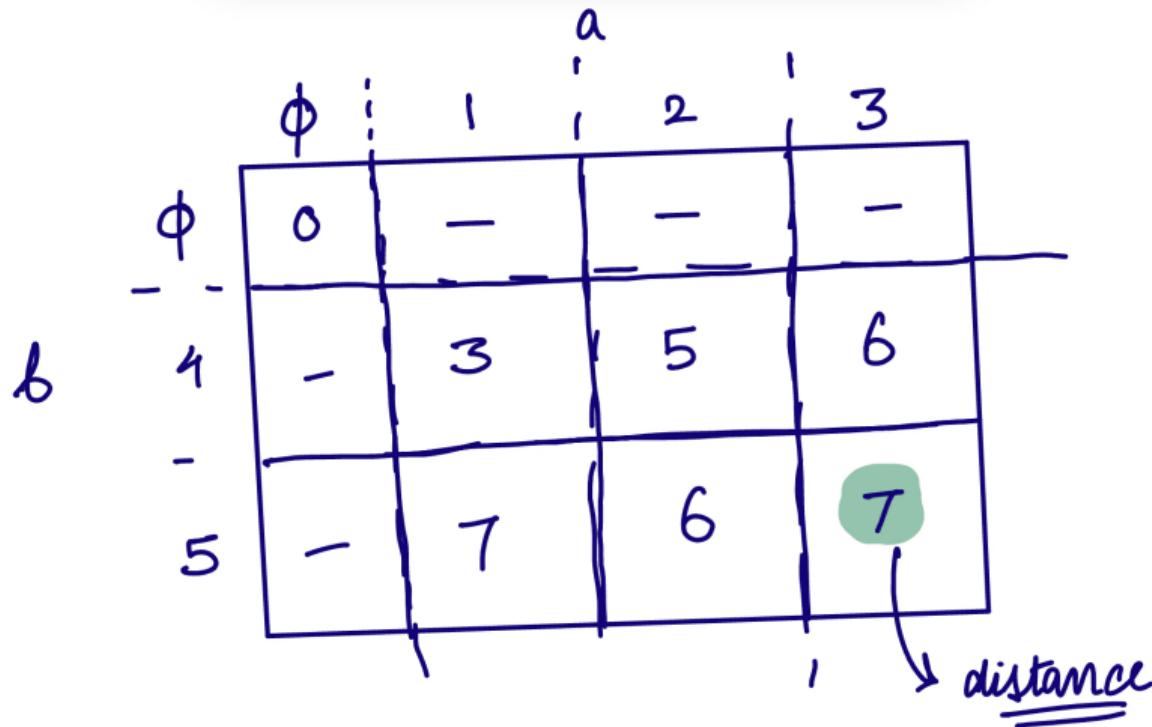
Implement.

Results

Conclusions

References

Computing distance between $a = [1, 2, 3]$ and $b = [4, 5]$





Learning the unknown θ

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

Temporal Difference(TD) error

$$\delta = r_{t+1} - \bar{r}_t + Q(h_{t+1}, g_{t+1}; \theta) - Q(h_t, g_t; \theta) \quad (6)$$

where

$$\theta_{t+1} = \theta_t + \alpha \delta \nabla Q(h_t, g_t; \theta_t) \quad (7)$$

$$\bar{r}_{t+1} = \bar{r}_t + \beta \delta \quad (8)$$

Algorithm 1: Intent-aware learning in continuing tasks

Input : Linear approximation function
 $Q(h, g; \theta) = \langle \theta, \phi(h) \rangle$

Parameters: Learning rate $\alpha, \beta > 0$
Goal update frequency f

- 1 Initialize θ, \bar{r} arbitrarily (e.g., $\theta = 0, \bar{r} = 0$)
- 2 Initialize observation o and goal g
- 3 Initialize history h as an empty stack
- 4 **while** *True* **do**
- 5 Plan for g , observe r, o
- 6 Push o into h : $h' \leftarrow [h, o]$
- 7 **if** $r == 0$ *and* $\text{random}() < \frac{1}{f}$ **then**
 // Keep goal with randomness
- 8 **continue**
- 9 **else**
- 10 Compute $\phi(h, o)$ by inferring intent of other agents
- 11 Choose goal g' as a function of $Q(h, \cdot; \theta)$ (e.g., ϵ -greedy)
- 12 **if** $r != 0$ **then**
- 13 $\delta \leftarrow r - \bar{r} + Q(h', g'; \theta) - Q(h, g; \theta)$
- 14 $\bar{r} \leftarrow \bar{r} + \beta\delta$
- 15 $\theta \leftarrow \theta + \alpha\delta\nabla Q(h, g; \theta)$
- 16 **end**
- 17 $o \leftarrow o', g \leftarrow g'$
- 18 **end**
- 19 **end**



Table of Contents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

① Introduction

② Intent-aware Multi-agent Reinforcement Learning

③ Library Implementation

④ Simulation Results

⑤ Conclusions



My implementation

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

https://github.com/SoutrikBandyopadhyay/intent_aware_rl

A library I developed for implementing the Q learning framework mentioned in this paper.



A brief technical aside

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

- Project written in Rust (programming language) from scratch
- Why Rust?
 - ① Excellent performance (close to C)
 - ② Memory Safe (reduces a significant number of memory related vulnerabilities)
 - ③ Friendly syntax



Code snippets

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

```
1 2 3 4 5 6 7 lib.rs - Doom Emacs | 12:12 AM / 67% | VOL 90% | us | RAM 49% | CPU 2% | Fully charged
16
15 //Adapted from https://towardsdatascience.com/dynamic-time-warping-3933f25fcdd
14 pub fn dtw_distance<T: Distance>(s: &Vec<T>, t: &Vec<T>) -> f64 {
13     //Create a 2D matrix to store the DTW values
12     let n = s.len();
11     let m = t.len();
10
9     let mut dtw: Vec<Vec<f64>> = vec![vec![f64::INFINITY; m + 1]; n + 1];
8     dtw[0][0] = 0.0;
7     for i in 1..n + 1 {
6         for j in 1..m + 1 {
5             let cost = s[i - 1].distance(&t[j - 1]);
4             let temp1 = dtw[i - 1][j];
3             let temp2 = dtw[i][j - 1];
2             let temp3 = dtw[i - 1][j - 1];
1             let last_min = f64::min(f64::min(temp1, temp2), temp3);
46             dtw[i][j] = cost + last_min;
1         }
2     }
3     return dtw[n][m];
4 }
5
| ● 2.5k intent_aware_rl/src/lib.rs 46:0 26% Rustic 1.67.1 master ✓

```



Code snippets (contd.)

Intent
aware RL

Implement.



Table of Contents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.
Results

Conclusions

References

① Introduction

② Intent-aware Multi-agent Reinforcement Learning

③ Library Implementation

④ Simulation Results

⑤ Conclusions



Simulation Setup

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

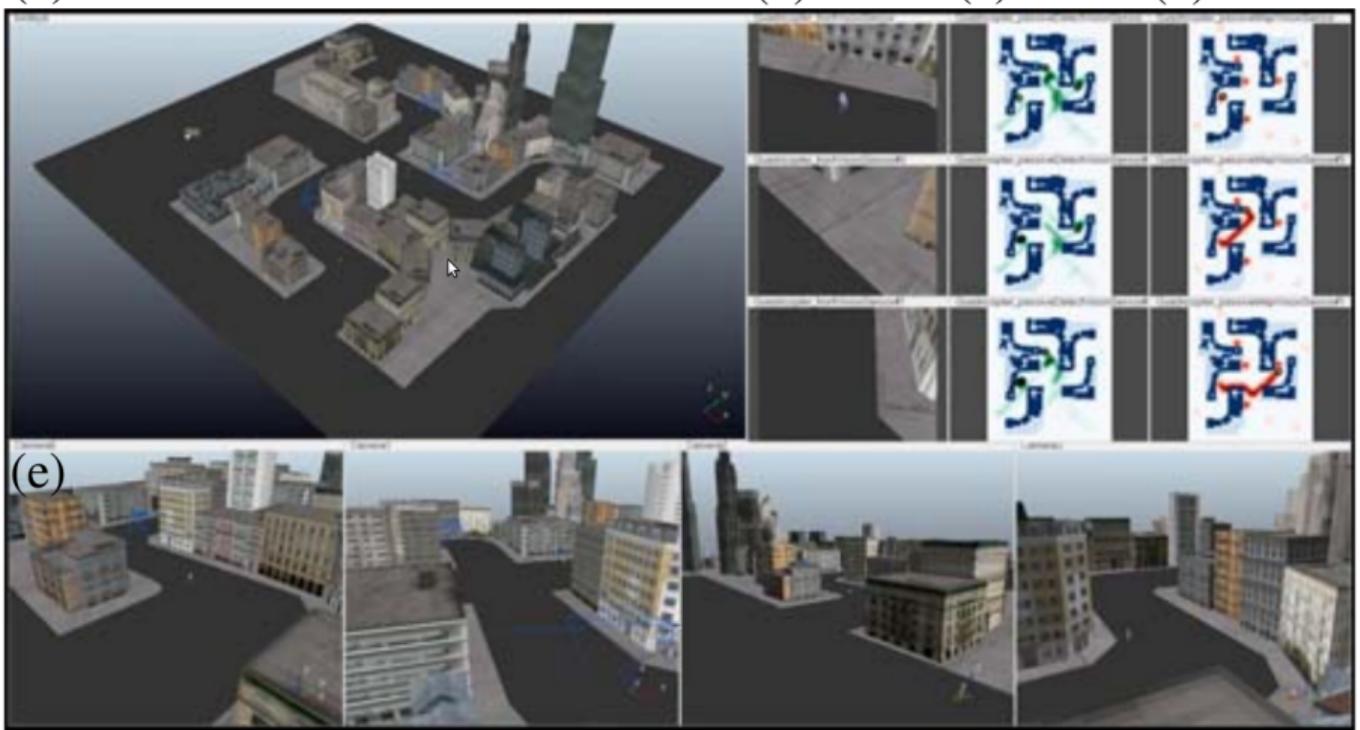
References

Multi-robot aerial surveillance

- There are 5 buildings, and a person is trying to sneak into the buildings.
- There are 3 aerial robots that are tasked to not let the human walk into a building undetected

Intent must be identified

- ① Human - For suspicious activity
- ② Other robots - For collaboration



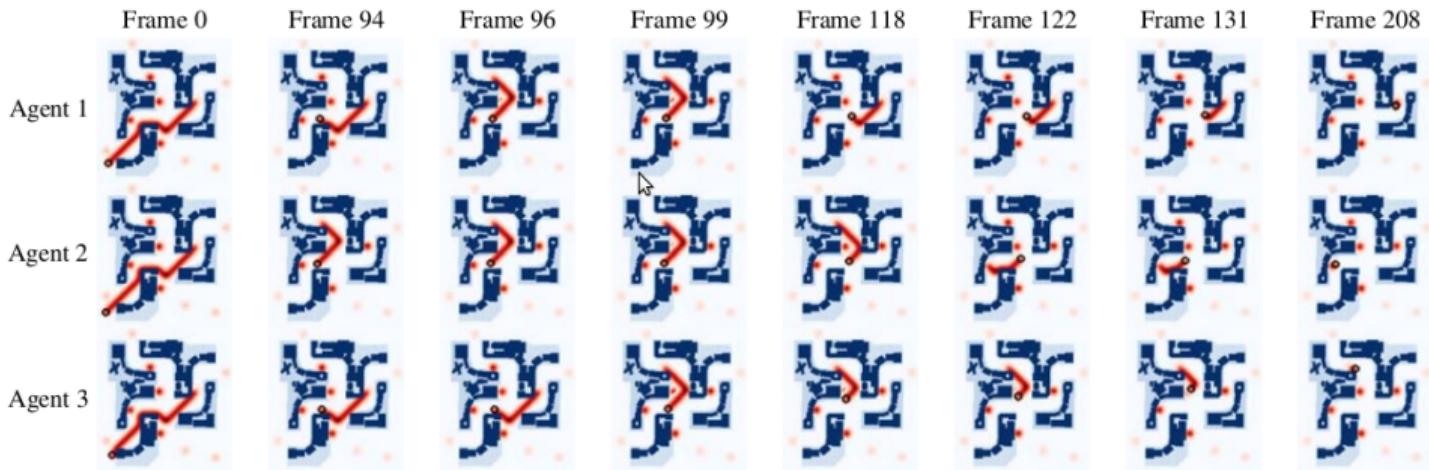


Fig. 8: Cooperative behavior of the learned model. Each row demonstrates the values of different goals (the higher the value the higher intensity of red), the chosen goal, and the planned path for a robot at different time frames. The values of monitoring buildings are higher than monitoring entrances or crossroad. At the beginning, all robot agents intended for the farthest building since they thought it is less likely for the other robots to monitor. During the process, the agents realized that the other agents might have chosen the same goal as themselves and chose another goal accordingly (e.g., agent 2 at frame 94). The agents chose different goals iteratively and finally reached an equilibrium that three robots were monitoring three different buildings. This dynamic process reflects how humans naturally interact with each other in a complex environment.



Comparision

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

		Testing Accuracy(%)				
		Training Scene				Testing Scene
Training Iterations		100	200	300	500	1000
random		29.7	29.7	29.7	29.7	29.7
greedy		31.4	31.4	31.4	31.4	31.4
RNN-POMDP		17.3	22.9	20.8	21.0	23.6
ours		13.8	19.1	60.3	62.4	59.3



Table of Contents

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

① Introduction

② Intent-aware Multi-agent Reinforcement Learning

③ Library Implementation

④ Simulation Results

⑤ Conclusions



Concluding Remarks

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.

Results

Conclusions

References

- Proposed an intent-aware multi-agent RL algorithm
- Generalizable framework
- Q function is linear (guaranteed convergence)



Possible Extensions to the work

Intent
aware RL

S. Bandy-
opadhyay

Introduction

Intent-
aware RL

Implement.
Results

Conclusions

References

- Sophisticated intent estimators
- Sophisticated high-level planning algorithms

References I

- [1] “Dynamic time warping,” *Information Retrieval for Music and Motion*, pp. 69–84, 2007. doi: [10.1007/978-3-540-74048-3_4](https://doi.org/10.1007/978-3-540-74048-3_4).

Thank you :)