

Adaptive Spatio-Temporal Organization in Groups of Robots

Torbjørn S. Dahl, Maja J. Matarić and Gaurav S. Sukhatme

*Robotics Research Laboratory
Computer Science Department, University of Southern California
tdahl|mataric|gaurav@usc.edu*

Abstract

This paper presents experiments, in simulation, with a group of robots that improve their performance on a straightforward transportation task by using reinforcement learning to associate input states with a set of abstract behaviors. We show that the improvement in performance is a result of the group adapting its spatio-temporal organization to the given environment. Spatio-temporal adaptation is a general form of adaptation in that it can improve performance over a range of different tasks and environments. Hence it increases the general applicability and autonomy of robotic systems. Lastly, we present two communication strategies that improve this ability to adapt by generally improving learning rates for cooperative robots in highly dynamic domains.

1 Introduction

Adaptive capabilities in groups of robots are attractive because they can increase the applicability and efficiency of these systems. Learning provides a way of producing behaviors that are difficult or impossible to program because of unfamiliar or complex problem spaces. As an example of such problems, we study the adaptation of the spatio-temporal organization of a group of robots with the aim to improve the group's performance on a straightforward transportation task.

Behavior-Based Systems have provided examples of how abstract behaviors and communication can provide tractable problem spaces for learning in multi-robot systems [7]. Previous work in reinforcement learning (RL) [13] has shown the possibility of using communication as a general way of improving learning rates for cooperative agents. To demonstrate the applicability of the RL methods however, it is necessary to show that similar results can be achieved in real multi-robot systems with real world data and scalable communication.

We present experiments, in simulation, with a group of five robots that use reinforcement learning on ab-

stract behaviors to adapt their spatio-temporal organization. This form of adaptation is general enough to be used to improve group performance over a range of different tasks and environments and as such, it exemplifies the kind of capability that increases the general applicability and autonomy of robotic systems. The policies produced by the adaptive controllers during our experiments show that the improved performance is indeed a result of spatio-temporal reorganization. The learned policies also show that the robots learn to use data relating to the position and orientation of the other robots to decide whether to move ahead or stop.

Lastly, we present two communication strategies developed to improve learning rates in groups of cooperative robots. The strategies are particularly well suited for use in highly dynamic domains such as multi-robot systems. The experimental data provided show that these communication strategies produce statistically significant increases in convergence rates for both individual and group performance.

2 The Transportation Problem

In order to demonstrate how the learning and communication affect the group's performance, we studied a straightforward transportation problem where the robots had to make their way back and forth between a source and a sink. We did not yet require physical objects to be transported. A traversal corresponded to one crossing from the source to the sink or vice versa. Transportation is similar to foraging and collection [1, 4, 9], but simpler as we assume that the positions of both the source and the sink are known to the robots, thus excluding issues of searching.

We measure individual performance on the transportation task by the time taken for one traversal, and group performance by the period between the arrival at either the source or the sink by any of the robots. We call this period the *target time*. By providing a simple problem in an open environment, we focus the learning on the complexity that comes from the dynamics of robot interaction rather than

on the complexity that comes from the task or from the environment.

3 Experimental Setup

We instantiated the transportation problem by simulating five Pioneer 2DX robots with PTZ cameras and SICK laser range-finders traversing a six-by-eight meter rectangular area. The source and sink were unique beacons made from highly reflective material, recognizable by laser. Each of the robots also wore color markings recognizable using ActiveMedia’s Color-Tracking Software (ACTS). The markings emphasized the orientation of the robot with red on the left side, green on the right side, and yellow at the rear. The prototype markings as they appear on a real Pioneer are shown in Figure 1.



Figure 1: Pioneer Robot with Color Markings

We performed the experiments in simulation on the *Player/Stage*[3] software platform¹. From experience, controllers written for the *Stage* simulator work with little or no modification on the real Pioneers.

The experiments presented below used only local communication in order to be scalable. Local communication was implemented using a wireless communication network and a Global Positioning Device (GPS). Every message was tagged with the location of the sender, and robots would only receive messages from senders that were within a two-meter radius. The system only used GPS to simulate local communication. The controllers used odometry for position estimation when the landmark beacons were not visible. Relying strictly on local sensors increases the applicability of our solution and ensures the scalability of the control algorithms.

3.1 Reward Structure

During the experiments the robots received rewards related to their individual efficiency every time they

¹*Player* is a server and protocol that connects robots, sensors and control programs across a network. *Stage* simulates a population of Player devices, allowing off-line development of control algorithms. *Player* and *Stage* were developed jointly at the USC Robotics Research Labs and HRL Labs and are freely available under the GNU Public License from <http://playerstage.sourceforge.net>.

arrived at the source or the sink. The reward, r , received after a traversal of time, t , is given in Equation 1. The number 250 reflects the shortest possible traversal time, giving a reward between 1.0 and 2.0.

$$r = 1 + \frac{250}{t} \quad (1)$$

4 Controller Architecture

We implemented a controller with three abstract movement behaviors: *stopping*, *wall-following* and *direct target approach*. These three behaviors made use of the robots’ lower level capabilities, such as obstacle avoidance and position estimation. They also kept track of what target to visit next, the source or the sink. The mapping of input states to these behaviors was left to the adaptive control algorithm, as described below. On top of the adaptive control structure, we implemented two different communication behaviors that could be expressed in parallel with the movement behaviors. The communication behaviors were fixed and not available to the adaptive control algorithm. In Figure 2 we show the activation and parallel expression of the main parts of our control architecture. The horizontal arrows indicate parallel behavior expression, while the vertical arrows describe behavior activation.

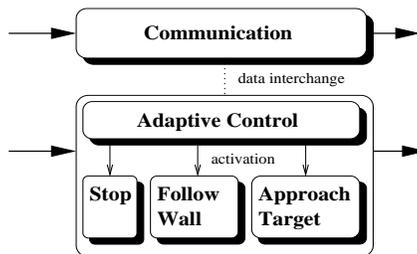


Figure 2: Behavior Activation and Expression

The Underlying Behaviors. The pre-programmed behaviors *stopping*, *wall-following* and *direct target approach* generally steered the robots toward their targets. The *direct target approach* behavior servoed the robot directly toward the current target. The *wall following* behavior followed the boundaries of the area while moving toward the target. Finally, the *stopping* behavior froze the robot in its current position and orientation. Together, these three behaviors allowed the robots to organize in space by learning e.g., when to concentrate on the shortest path between the source and the sink, when to disperse around the boundaries of the environment, and when to stop and wait.

4.1 The Adaptive Control Algorithm

The adaptive element of the controller decided the behavioral policy of each individual robot controller by relating each of the possible input states to one of the three available movement behaviors. It used reinforcement learning to improve the robot’s performance according to the rewards received.

Random application of the underlying behaviors provided a baseline for performance evaluation. Hence, the conceptual goal of the learning algorithm was to improve on that performance rather than to develop a transport capability from scratch. The base performance was clearly suboptimal, with robots getting in each other’s way and stopping regularly. This left ample opportunity for the robots to learn spatial and temporal structures that would reduce interference and optimize performance.

The State-Action Space. The robots had three abstract behaviors available to them: *stopping*, *wall-following* and *direct target approach*. These behaviors were all applicable in all states.

The input state was represented as five bits. The first bit indicated the presence of a target, the source, or the sink. The second bit was on for relatively short (1-30 second) periods of time at random intervals. The additional states produced by this bit allowed the robot to express certain behaviors for a short period at random times independent of the rest of the input state. This fragmentation of the input state space was an unbiased way of increasing the sophistication of learnable policies. For example, instead of always following the walls whenever a target is not visible, the robot can learn to approach the target directly in this state, but with short random periods of wall following. The three final bits represented the presence of color-blobs corresponding to the three colored markings on the robots, indicating the presence and orientation of another robot. When several robots were present, these bits only reported on the visible markings of the closest one, distinguished by the height of the visible color-blobs.

The five input bits and three behaviors made up a state-action space of size $2^5 * 3 = 96$. In the future we intend to enlarge this space to allow more interesting organizational structures to emerge. However, even with this reduced space we are able to demonstrate both spatial and temporal adaptation as well as use of data directly describing other robots.

The Learning Algorithm. It is difficult or impossible to create useful models of the complex dynamics of multi-robot interaction. In our scenario, the lack of a world model limits the information that can be communicated between robots about features

of the environment and past experiences.

To update the adaptive part of the control policy we used *naive* $Q(\lambda)$, a model-free, off-policy, temporal difference (TD) learning algorithm with eligibility traces. We used an accumulating rather than a replacing trace. Though the *naive* $Q(\lambda)$ algorithm is not proved to converge to optimality, related algorithms have shown superior performance in empirical tests [12]. The value estimate for a state, s , and an action, a , resulting in a reward, r , and a subsequent state, s' , is, according to *naive* $Q(\lambda)$, identical to the estimate for standard TD-learning, as shown in Equations 2 and 3 with learning rate, α , return discount rate, γ , and TD-error, δ .

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta \quad (2)$$

$$\delta = r + \gamma \max_b Q(s', b) - Q(s, a) \quad (3)$$

The estimated values for a previous state, s_n and action, a_n , recorded in the eligibility trace n actions ago, is given in Equation 4 with trace decay rate, λ .

$$Q(s_n, a_n) \leftarrow Q(s_n, a_n) + \alpha (\gamma \lambda)^n \delta \quad (4)$$

Action Selection. Actions were chosen by the adaptive controller using the Boltzmann Softmax algorithm. With this algorithm the probability of an action being selected is related to its estimated value in the current state. As the time-sensitive parameter, τ , an analog of temperature, decreases toward 0, so the relative probability of the highest estimated value being chosen increases. The probability of choosing action, a , in state, s , according to the Boltzmann Softmax algorithm is given by Equation 5.

$$P(a|s) = \frac{e^{\frac{Q(a,s)}{\tau}}}{\sum_b e^{\frac{Q(b,s)}{\tau}}} \quad (5)$$

5 Experimental Validation of Spatio-Temporal Adaption

To evaluate the adaptive controller, we ran five trials of five robots each performing two-hundred traversals, a total of 1000 traversals per trial. The initial individual performance, measured by traversal time, had an average, \overline{T}_i , of 1112.5 seconds and a standard deviation, s_i , of 1023.7 seconds. These values and the corresponding values after 50, 100, 150 and 200 traversals, t_i , are given in columns 2 and 3 of Table 1. The group performance, measured in average *target time*, the time between the arrival at the source or the sink by any robot, had an initial average, \overline{T}_g , of 245.2 seconds and a standard deviation, s_g , of 165.3

seconds. Columns 5 and 6 of Table 1 give these values and the corresponding values after 250, 500, 750 and 1000 target arrivals, t_g .

t_i	\bar{T}_i	s_i	t_g	\bar{T}_g	s_g
0	1112.5	1023.7	0	245.2	165.3
25	843.9	628.0	125	159.8	102.7
50	759.7	559.8	250	161.4	130.4
75	639.6	445.9	375	133.6	84.4
100	668.3	504.5	500	139.5	114.4
150	537.7	343.0	750	123.1	86.7
200	524.8	296.8	1000	100.0	72.8

Table 1: Improvement from Adaptation

The results from these experiments prove that the robot group improved its performances by producing traversal and target times that are lower than the initial random behavior application by a statistically significant amount.

5.1 The Learned Policies

The Q-tables produced during the experiments show that the robots on average learned policies dominated by the *direct target approach* behavior, but with significant portions of the *wall-following* behavior. These policies created a common circular path, minimizing interference.

We also found evidence that the robots adapted their temporal organization according to the presence and orientation of other robots. Of the input states where yellow markings were visible, i.e. states with a one in bit four, we calculated the number of times the *stopping* behavior had the highest q-value. We called this number the *stop-rate* for *yellow states*. The stop-rate for the yellow states was significantly lower than the stop-rate for states in which robots were visible, but no yellow markings. The stop-rate was also significantly lower for yellow states than for the corresponding *red states* and *green states*. This indicates that the robots in general learned to stop only when they saw a robot but could not see its yellow markings. Intuitively, visible yellow or rear markings indicate that the other robot is facing, and likely moving, away. Hence, it will not be likely to cause an obstruction.

6 Communication Strategies for Improved Adaptability

In order to improve the adaptive capabilities demonstrated above we developed two communication strategies for cooperative robot learning. The *expressed behavior evaluation* strategy was inspired by

calls of encouragement as found in animals and humans. Such calls express the value a critic puts on a subject’s behaviors. The *imitation* strategy is inspired by imitation learning, short periods of interaction between robots where any number of students repeat the actions of a teacher. Mataric [8, 6] has also previously used communication and imitation in the context of reinforcement learning for groups of robots.

6.1 Expressed Behavior Evaluation

This communication strategy presupposes the ability to accurately estimate the current input state and active behavior of another robot. This is in general very difficult for real world robots. However, by having the robots continuously communicate these two pieces of information locally, we get the same effect for a tolerable increase in communication. To allow individual rather than a broadcast feedback, the robots also communicate a unique identifier.

Being able to give feedback implies having an internal estimate of the value of an observed state-action pair. As the robots are homogeneous, this estimate already exists in a robot’s own Q-table. Communicating early estimates of Q-values provides little information, so we introduced a limitation on what state-action combinations a robot would give feedback on. A feedback message was only sent when a robot had tried all the actions in the given state and one of the action values was significantly higher than all the others by some given margin, ϕ . The feedback message contained only the unique identifier of the robot that was evaluated.

When a robot received a feedback message, it treated it as an independent estimate of the current state and active behavior doing a full backup over its current trace, independently of other rewards received or backups performed in the current state. The feedback-inspired backup uses the exceptional behavior margin ϕ as the TD-error, replacing δ with ϕ in Equations 2 and 4 in Section 4.1.

6.2 Imitation

The second strategy was inspired by imitation learning. During imitation a student imitates the actions of a teacher, hence sharing the state-action sequence of a trial. Rather than physically doing the imitation, this strategy communicated the state-action sequence of the best trial a robot had experienced so far. The best trace was communicated at random intervals averaging five seconds, providing roughly one imitation experience per robot per trial.

Upon reception of an imitation message, the robots went through the received trace as if they were experiencing the relevant states and taking the related

t_i	\overline{T}_i	s_i	t_g	\overline{T}_g	s_g
0	1027.9	727.3	0	209.2	142.5
25	647.8	405.9	125	155.3	106.3
50	589.7	350.2	250	121.3	87.9
75	592.2	295.9	375	104.2	77.1
100	621.6	441.8	500	124.4	86.5

Table 2: Performance for Expr. Behavior Eval.

t_i	\overline{T}_i	s_i	t_g	\overline{T}_g	s_g
0	937.776	614.9	0	237.3	153.0
25	609.9	316.9	125	122.7	94.0
50	548.9	321.6	250	129.1	89.3
75	567.0	378.7	375	106.5	86.8
100	565.0	376.3	500	116.5	94.9

Table 3: Performances for Imitation

actions themselves, using Equations 2, 3 and 4, but without actually taking any actions.

7 Experimental Evaluation of the Communication Strategies

In order to evaluate the contribution the communication made to the group learning, we performed five trials of 200 traversals using each of the communication strategies and compared the learning rates to those found for the communication-free learning experiment presented in Section 5. The individual and group performance data are presented in Tables 2 and 3, using the same format as in Table 1.

The improvements of the group performance in the three different experiments are plotted as graphs in Figure 3, showing that both communication strategies improve on the communication-free or individual learning. Analysis also shows that on a 95% confidence level, the traversal times for the individual controllers using communication are significantly better than those of the communication-free controller after 25, 50 and 75 traversals. After 100 traversals, however, the communication-free controller has caught up, as can be seen in Tables 2 and 3. The same result is apparent for group performance after 125, 250 and 375 arrivals, but only on a 90% confidence level. Only after 125 and 375 arrivals is the difference significant on a 95% confidence level.

The stopping behavior, which was dependent on the markings of other robots, became more pronounced when communication took place than it was in the communication-free case.

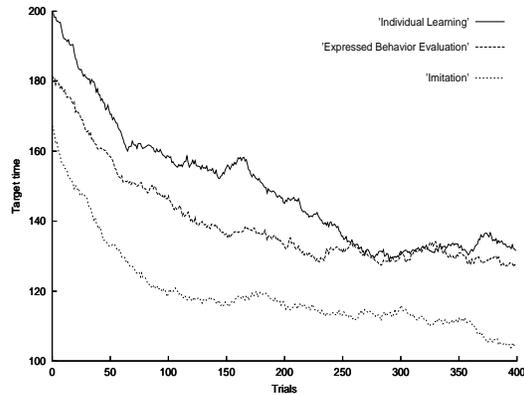


Figure 3: Convergence Rates

8 Related Work

Stone and Veloso [11] use robust, low-bandwidth communication to modify group structure in the domain of simulated soccer. The modifications, however, are purely switches between pre-set formations and do not include any adaptation. The same is the case for Ünsal and Bay’s [14] work on spatial organization of swarm systems, where the group adapts to its environment, but the individual agents are non-adaptive. Itsuki, Hitoshi and Kazuo [5] also describe experiments in simulated soccer. Using artificial neural networks, they demonstrate two agents learning whether to pass or shoot in different but related situations also involving two hostile agents.

Balch and Arkin [2] demonstrate that improvement of group behavior is possible using cooperative communication in multi-agent systems, but explicitly restrict themselves to reactive systems. Balch [1] also investigates the effects of different reward structures on the diversity produced in a group of foraging, BB robots, that adapt using RL. Mataric [7, 6, 8] has previously shown the feasibility of group learning in behavior space and demonstrated uses of communication and imitation for cooperative learning.

Sen, Sekaran and Hale [10] have previously shown that group behavior can improve through learning even in the case when the participating agents are oblivious of each other’s existence. Tan [13] has shown three communication strategies that improve the learning rate in groups of cooperative agents. These strategies are broadcast based and imply updating shared data structures. As such they are unlikely to scale up to larger groups. Like Sen et al., Tan’s work is also based on an abstract world, modeling only simple physics.

Our work goes beyond previous work on reinforcement learning by doing experiments in simulations that are a significantly more realistic than the abstract world models of Sen et al. and Tan and the

simulated soccer environment used by Stone and Veloso, and Itsuki et al. We also add to previous work on cooperative mechanism by presenting a new strategy for cooperation, *expressed behavior evaluation*. This strategy goes further than previously demonstrated strategies in minimizing communication for cooperative sharing of experiences. The adaptive capabilities we present here are more general than Matarić's previous uses of learning in groups, and as such apply to a wider range of problems and go further in increasing general robot adaptability and autonomy.

9 Conclusion

Our experiments demonstrate the possibility of adaptive spatio-temporal organization in groups of real robots. This work takes one step toward providing robots with more general adaptive capabilities, hence making them more applicable and autonomous.

Our experiments also show that our cooperative communication strategies can contribute significantly to improving a group's adaptability. Our *imitation* strategy is very similar to Tan's *episode exchange* strategy, and our experiments show that strategies of this kind are useful in groups of real robots as well as in the grid-world example Tan uses.

10 Future Work

In order to evaluate our experimental results we plan to run tests to see whether the performance improvement observed in the simulation is also observable on groups of real robots, using up to 9 Pioneers available in our laboratory. To show the generality of the adaptive capabilities we will also run experiments using different environments and different tasks, initially expanding the current framework to look at the transportation of physical objects.

As a next step we aim to enlarge the problem space presented to the adaptive controller in order to produce more advanced organization patterns such as lanes, territories, or bucket brigading. The work presented here is an initial step to developing these more advanced capabilities. We also plan to investigate the use of social structures such as hierarchies to facilitate group adaptation.

Acknowledgments

This work is supported in part by DOE RIM grant DE-FG03-01ER45905 for "Multi-Robot Learning in Tightly-Coupled, Inherently Cooperative Tasks" and in part by ONR DURIP grant 00014-00-1-0638.

References

- [1] T. R. Balch. The impact of diversity on performance in multi-robot foraging. In Oren Etzioni, Jörg P. Müller, and Jeffrey M. Bradshaw, editors, *The proceedings of the Third International Conference on Autonomous Agents*, pages 92–99, Seattle, WA, May 1-5 1999. ACM Press.
- [2] T. R. Balch and R. C. Arkin. Communication in Reactive Multiagent Robotic Systems. *Autonomous Robots*, (1):1–25, 1994.
- [3] B. P. Gerkey, R. T. Vaughan, K. Støy, A. Howard, G. S. Sukhatme, and M. J. Matarić. Most Valuable Player: A Robot Device Server for Distributed Control. In *Proceedings of the 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'01)*, pages 1226–1231, Wailea, Hawaii, October 29 - November 3 2001.
- [4] D. Goldberg and M. J. Matarić. Design and Evaluation of Robust Behavior-Based Controllers for Distributed Multi-Robot Collection Tasks. In T. Balch and L. E. Parker, editors, *Robot Teams: From Diversity to Polymorphism*, pages 315–244. A K Peters Ltd, 2001.
- [5] N. Itsuki, M. Hitoshi, and H. Kazuo. Learning Cooperative Behavior in Multi-agent Environment - Case Study. In *Proceedings of PRICAI'96*, pages 570–579, Cairns, Australia, August 26 - 30 1996. Springer.
- [6] M. J. Matarić. Learning social behavior. *Robotics and Autonomous Systems*, 20:191–204, 1997.
- [7] M. J. Matarić. Reinforcement Learning in the Multi-Robot Domain. *Autonomous Robots*, 4(1):73–83, 1997.
- [8] M. J. Matarić. Reducing locality through communication in distributed multi-agent learning. *Journal of Experimental and Theoretical Artificial Intelligence, Special issue on Learning in DAI Systems*, 10(3):357–369, 1998.
- [9] E. Østergaard, G. S. Sukhatme, and M. J. Matarić. Emergent Bucket Brigading - A simple mechanism for improving performance in multi-robot constrained-space foraging tasks. In *Proceedings of the 5th International Conference on Autonomous Agents*, pages 29–30, Montreal, Canada, May 28 - June 1 2001. ACM Press.
- [10] S. Sen, M. Sekaran, and J. Hale. Learning to coordinate without sharing information. In *Proceedings of the Twelfth National Conference of Artificial Intelligence (AAAI'94)*, pages 426–431, Seattle, WA, August 1-4 1994.
- [11] P. Stone and M. Veloso. Task Decomposition, Dynamic Role Assignment, and Low-Bandwidth Communication for Real-Time Strategic Teamwork. *Artificial Intelligence*, 110(2):241–273, 1999.
- [12] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. Adaptive Computation and Machine Learning. MIT Press, 1998.
- [13] M. Tan. Multi-Agent Reinforcement Learning: Independent vs. Cooperative Agents. In *Proceedings of the 10th International Conference on Machine Learning (ICML'93)*, pages 330–337, Amherst, MA, June 27-29 1993. Morgan Kaufmann.
- [14] C. Ünsal and J. S. Bay. Spatial Self-Organization in Large Populations of Mobile Robots. In *IEEE International Symposium on Intelligent Control*, pages 249–254, Columbus, Ohio, August 16-18 1994.