

A machine learning method for improving task allocation in distributed multi-robot transportation

Torbjørn S. Dahl

Norwegian Defence Research Establishment

Torbjorn-Semb.Dahl@ffi.no

Maja J Matarić and Gaurav S. Sukhatme

Center for Robotics and Embedded Systems

University of Southern California

mataric|gaurav@usc.edu

Machine learning is a means of automatically generating solutions that perform better than those that are hand-coded by human programmers. We present a general behavior-based algorithm that uses reinforcement learning to improve the spatio-temporal organization of a homogeneous group of robots. In this algorithm each robot applies the learning at the level of individual behavior selection. We demonstrate how the interactions within the group affect the individual learning in a way that produces group-level effects, such as lane-formation and specialization, and improves the group's performance. We also present a model of multi-robot task allocation as resource distribution through *vacancy chains*, a distribution method common in human and animal societies, and an algorithm for multi-robot task allocation based on that model. The model explains and predicts the task allocation achieved by our algorithm and highlights its limitations. We present experimental results that validate our model and show that our algorithm outperforms pre-programmed solutions. Last, we present an extension of our algorithm that makes it applicable to heterogeneous groups of robots

by making it sensitive to differences in individual robot performance levels.

1.1 Introduction

Machine learning (ML) [24] is a means of automatically generating solutions that perform better than those that are hand-coded by human programmers. Such improvement is possible in problem domains where optimal solutions are difficult to identify, i.e., when there are no models available that can accurately relate a system’s dynamics to its performance. One such domain is the control of multi-robot systems.

Mobile robots are notoriously difficult to control in a robust, reliable, and repeatable fashion. The challenges stem from uncertainty inherent in physically embodied systems, including in sensors, effectors, and interactions between the system components and the environment. The *behavior-based* (BB) control paradigm [22, 1] provides a means of structuring robot controllers into collections of task-achieving modules or *behaviors*, such as exploration and obstacle avoidance. The modules operate in parallel and interact within the system and also through their effects on the environment. When properly designed, the resulting controller produces robust, repeatable, and reliable overall behavior for the robot.

The BB paradigm takes inspiration from biology and the metaphors used to describe animal behavior. Behaviors also provide a basis for employing ML to improve the robot’s performance [23]. ML can be employed at the level of behavior coordination, or at the level of the behaviors themselves. The behaviors typically have a set of *preconditions* used to decide when they are suitable for execution. These preconditions and behaviors together form a tractable search space for improvement through ML.

In Section 1.4 we present an algorithm that uses ML at the level of behavior selection, applied to the multi-robot control domain. Specifically, we show how a group of robots improved its performance on a cooperative transportation task by using individual reinforcement learning to adjust individual behavior selection policies from feedback related to individual performance.

In Section 1.5 we present a BB solution to multi-robot task allocation (MRTA) where robots represent tasks internally as behaviors and use ML to automatically improve their individual task selection policies. Through interaction the robots specialize on different tasks and, from these specializations, global task allocation emerges. We also present a model of MRTA called task allocation through *vacancy chains* (TAVC). A vacancy chain [7] is a resource distribution mechanism common in animal and human societies. The TAVC model represents opportunities to improve a group’s performance by reallocating tasks, as chains of task-related vacancies to be distributed between the robots. The TAVC model captures the effects of group dynamics in such systems and explains and predicts the task allocations produced by distributed, communication-free MRTA. We demonstrate a TAVC algorithm, in realistic simulations, using a ho-

mogeneous group of robots in a dynamic environment where task values change and where robots are susceptible to failure. Our results validate the TAVC model of MRTA and show that our TAVC algorithm improves on pre-programmed solutions to this problem.

In Section 1.6 we show that our TAVC algorithm also works for heterogeneous groups of robots using a Boltzmann-softmax action selection function. Our TAVC algorithm is a valuable alternative to existing MRTA algorithms as these are typically not sensitive to the complex effects of group dynamics, such as interference and synergy.

1.2 Motivation

As robot technology becomes increasingly ubiquitous, so will multi-robot systems. Control and coordination of groups of mobile robots is an active field of research; multi-robot task allocation (MRTA) is a key aspect of this problem. MRTA encompasses issues of skill utilization, cooperation, communication, and scheduling. One of the central problems in dealing with these issues is modeling the underlying group dynamics and their effects on the group’s performance. We have developed TAVC, a model of MRTA inspired by vacancy chains, a robust and efficient resource allocation mechanism common in biological systems. The advantage of TAVC over current MRTA models is that it is explicit about the effects of group dynamics. Based on the TAVC model, we have developed a novel MRTA algorithm that is sensitive to group dynamics. To achieve this sensitivity, the TAVC algorithm uses distributed RL to make individual task utility estimates. During learning these estimated utilities diverge according to the effects of the prevailing group dynamics. The divergence leads to a general improvement in the group’s performance.

1.2.1 Group dynamics

By the term *group dynamics*, we refer to the interaction among the members of a group of robots. Group dynamics can have positive or negative effects on the group’s performance. Interference is a typical negative effect in confined spaces, but positive or *synergistic* effects are also possible when the actions of one robot unintentionally facilitates the actions of another robot, or when robots intentionally cooperate. Existing multi-robot control algorithms and related theory such as *scheduling theory* [6] assume *task independence*. Task are independent when there are no significant interference-related or synergistic effects to be considered. However, in multi-robot systems in general, and in MRTA in particular, such effects are often critical to group performance. When this is the case, understanding and modeling those effects is important for constructing and predicting the quality of the control algorithms. There are major challenges related to modeling group dynamics in general. Each member of a group of robots is a physically embedded system with inherent uncertainty. Interactions among a group of physically embodied systems increase uncertainty to a level where

practical accurate models are extremely difficult to construct. Previous work on modeling group dynamics can be divided into *microscopic* (*simulation-based* or otherwise) and *macroscopic* [19]. Microscopic models explicitly model each agent. Simulation-based models simulate the actions taken by each so that properties of the system can be recorded as the agents interact. Macroscopic models describe system properties in terms of abstract features such as the number and general distribution of agents. We briefly overview representative relevant examples from both areas next.

Microscopic Models Game theory [25] provides an abstract model of interacting agents as *players* in formalized n -person stochastic *games*. Game theory represents players in terms of their *strategies* for playing a given game and their related *payoff matrices* for possible moves. Using these factors, game theory predicts optimality and stability features of games. A game theoretic approach to modeling group dynamics would require a specification of the payoff matrices involved. However, these cannot be known for any but the most trivial multi-robot systems. Littman [20] presented the Minimax-Q algorithm as a means of finding payoff matrices. The algorithm is guaranteed to find the equilibrium of any stochastic game. The exponential growth, n^m , of a payoff matrix for n tasks and m robots also means that it is not feasible to estimate the payoffs for each state/action combination. Bowling *et al.* [5] presented experimental data demonstrating how a reinforcement-learning algorithm based on policy gradient ascent and the WoLF (Win or Learn Fast) principle for adjusting the learning rate can overcome this complexity.

Simulation-Based Models To learn statistical models of interaction in a space of abstract robot behaviors, Goldberg and Mataric [16] developed Augmented Markov Models (AMMs). AMMs are transition probability matrices with additional temporal information. They used these models to maximize reward in multi-robot foraging/mine collection tasks. Balch *et al.* [3] used live insect colonies to construct three-dimensional histograms of insect presence over a discretized area, with the long-term goal of combining spatio-temporal models with Hidden Markov Models for behavior recognition [18] in order to recognize colony behaviors. Yan and Mataric [30] modeled group behaviors from spatial data describing both human and robot activity. Their models were based on proxemics and considered data related to individuals, pairs, clusters, and the environment. They also used three-dimensional histograms to identify and describe activity distributions produced by underlying behaviors. Seth [27] presented a simulation-based model using genetic algorithms to evolve foraging behaviors for multiple agents in spatially explicit environments. The evolved systems were able to reproduce interference functions previously described in field studies of insects, but not previously reproduced by microscopic models.

Macroscopic Models Mataric [21] proposed *group density* as a macroscopic model of interference. Group density is defined as the ratio between the agents'

footprints and the available *interaction space*. An agent’s footprint is defined by its sphere of influence, including geometry, motion constraints, and sensor range. When the number and size of the agents are known, a *mean free path* can be computed and used to estimate the *number of expected collisions* for agents executing random walks. Lerman *et al.* [19] studied three different types of models of cooperation and interference in groups of robots: sensor-based simulations, microscopic numerical models, and macroscopic numerical models. The three different models produced corresponding results but the macroscopic model had the advantage of being very fast and independent of the number of robots modeled.

Current macroscopic models are not sophisticated enough for optimizing the control of multi-robot systems with regards to specific problems. Simulation-based models are problem-specific but the time needed to for model-construction makes them unsuitable for real-time control of a multi-robot system. Currently there are no models of group dynamics with the speed, generality, and predictive capabilities needed to specify the effects of interaction on group performance. Simulation-based models are, in general, too slow and macroscopic models make too many simplifying assumptions to be of predictive use.

1.3 Cooperative transportation

Our experiments are related to the problem of *cooperative transportation*, where group dynamics often have a critical impact on performance. In transportation problems, a group of robots traverses a given environment in order to transport items between sources and sinks. We call the time taken to traverse the environment once from a sink via a source *and* back to a sink again the *traversal time*. We call the time between two arrival events at either source or sink the *target time*. To perform optimally the robots must maximize the number of traversals over time.

The basic transportation problem is one of the sub-problems of foraging [2, 17]. Foraging is reduced to a problem of transportation when the locations of the sources and sinks are known. When there are multiple sources and sinks, the transportation problem becomes the MRTA problem where transportation tasks must be allocated to robots in a way that optimizes overall performance. We refer to a problem as a *prioritized transportation problem* when sources and sinks have (sets of) different values. To optimize performance on prioritized transportation problems, measured in total value of delivered goods over time, a group of robots must strike the correct balance between the positive contribution of high target values and the negative contribution of interference resulting in increased delivery times.

1.4 Improving performance on cooperative tasks

Taking advantage of the reduction in state/action space made possible by the BB approach, we implemented a general algorithm, based on distributed reinforcement learning [28], that adapts the spatio-temporal properties of each robot’s behavior in a way that improves the group’s performance on cooperative tasks [9]. The robots learn from performance-related feedback. In homogeneous groups this feedback mainly reflects effects of the prevailing group dynamics and allows the robots to adjust their behavior so as to minimize the negative effects and maximize the positive ones.

In addition to the problem of accurately modeling the effects of group dynamics, multi-robot systems also face the inherent problem of scaling. Algorithms that make decisions based on global system state need complex communication structure to synchronize and update state information and coordinate individual robots. This makes it difficult to scale such centralized algorithms to work efficiently for systems with many robots.

This section presents our distributed algorithm for spatio-temporal organization. The algorithm works purely through *stigmergy*, i.e., implicit communication through the effects of actions on the environment, and hence avoids the scale-related difficulties of communication-dependent algorithms. We have demonstrated the effectiveness of our algorithm and analyzed the qualitative effects it had on individual action selection policies. Our results show that the algorithm changes the spatio-temporal organization of a group of robots according to the effects of the group dynamics. The most interesting result was a clear tendency toward policy differentiation, i.e., specialization. This effect was the motivation for adapting the algorithm as presented in this section for use on MRTA problems, where existing solutions have typically been insensitive to the effects of group dynamics.

1.4.1 The learning algorithm

Our algorithm for spatio-temporal organization is built on two main principles: *individual learning* and *emergence*. Each robot individually estimates the utility of the available actions and chooses from these in a way that allows it both to keep its estimates current through continuous experimentation and to exploit the actions with high estimated utilities. We used Q-learning for utility estimation.

Our algorithm depends on a reward structure that reflects the effects of the current group dynamics. Q-learning is not sensitive to the frequency of rewards. Hence, the estimated values of actions do not necessarily correspond to the action’s contribution to performance over time. In order to use Q-learning to optimize performance over time it is necessary to make the temporal aspect of performance explicit in the reward function. Such a reward function, using the last task processing time, t , and task value, w_i , is given in Equation 1.1.

$$r = w_i/t \tag{1.1}$$

Because we want the system to remain adaptive to changes in the environment, the experimentation rate ϵ does not decrease over time, as is common.

Using a reward function that is sensitive to the effects of the group dynamics allows the robots to learn to take actions that minimize these effect with respect to their individual performance. We assume that those individual changes in action selection policies are enough to produce spatio-temporal effects on a group level, such as turn-taking and lanes. The emergence of such phenomena would indicate that the changes in the individual action selection policies can be interpreted as spatio-temporal organization at a group level.

In Section 1.4.2 we demonstrate how, from individual Q-learning and the effects of the group dynamics, consistent and repeatable group-level effects emerge. In particular we show the effects of the emergent organization on the group's performance.

1.4.2 Experimental validation of spatio-temporal organization

We hypothesized that it was be possible for recognizable spatio-temporal features of group behavior such as turn-taking, lanes, and territories to emerge in a group of robots using our adaptive algorithm. Each robot was given a set of behaviors that performed the same general task, but did this in ways with different spatio-temporal qualities. We assumed that the spatio-temporal features listed above had genuine positive effects on group performance. If so, these features would represent easily identifiable optima in the problem space of spatially and temporally constrained multi-robot control. As such, these features may emerge from the interaction between robots that use Q-learning to individually estimate each behavior's utility in different locally perceived world states.

To test this hypothesis we choose to the problem of cooperative transportation, a problem that could be given clear spatial and temporal constraints. We also designed a set of transportation-related behaviors with different spatio-temporal qualities.

Simulated Environment For the experiments presented in this section, we simulated five Pioneer 2DX robots with PTZ cameras and SICK laser range-finders in a 6 by 8-meter rectangular area. There was one source and one sink, each placed in the center of one of the short sides of the rectangle. The robots wore color markings recognizable using ActiveMedia's Color-Tracking Software (ACTS), allowing them to perceive each other's orientation. The markings, inspired by ship lanterns which serve the same purpose at sea, were 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 1.1.

The source and sink were simulated with unique bar-codes (as shown in Figure 1.2) made from highly reflective material readable by the SICK laser. Both the recognition of the color markings and the reading of laser bar-codes have been validated on real robot tests.



Figure 1.1: Prototype Pioneer 2DX robot with color markings

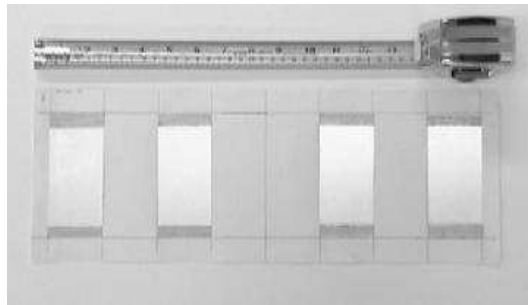


Figure 1.2: Prototype bar-code readable by the SICK laser

We performed our experiments in simulation, using the Player/Stage [15] software platform. From experience, controllers written for the Stage simulator work with little or no modification on the real Pioneers. A graphical rendering of the simulated environment is presented in Figure 1.3.

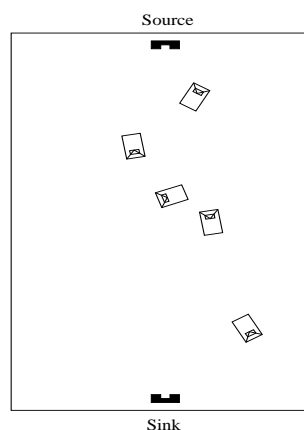


Figure 1.3: The simulated environment used for experiments on adaptive spatio-temporal organization

Controller architecture To support flexibility in the spatio-temporal features of the transportation, we implemented three transport-related behaviors:

- **direct target approach:** servoed the robot directly toward the current target.
- **wall following:** followed the boundaries of the simulated environment, eventually arriving at a source or a sink.
- **stopping:** froze the robot in its current position and orientation.

These three behaviors made use of lower level behaviors such as *obstacle avoidance*. The mapping of the transport-related behaviors to input states was done by the learning algorithm presented in Section 1.4.1.

Learning parameters The input state presented to the learning algorithm consisted of five bits. The first bit indicated the presence of a source or a sink. The second bit was on for random intervals of time (1-30 seconds), allowing the robot to express certain behaviors for short periods independent of the rest of the input state. For example, instead of always following the walls whenever a target is not visible, the robot could learn to approach the target directly *most of the time* in this state, interspersed with short 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 reported only 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 $2^5 * 3 = 96$. Actions were chosen by the adaptive controller using a Boltzmann-softmax action selection function.

Results We ran a total of five trials. During every trial the five robots each performed around two hundred traversals, resulting in a total of 1000 traversals per trial. The initial individual performance had a mean traversal time, $\hat{\mu}_i$, 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, n , are given in columns 2 and 3 of Table 1.1. The group performance had an initial mean target time, $\hat{\mu}_g$, of 245.2 seconds and a standard deviation, s_g , of 165.3 seconds. Columns 5 and 6 of Table 1.1 show these values and the corresponding values after 250, 500, 750 and 1000 target arrivals, m .

n	$\hat{\mu}_i$	s_i	m	$\hat{\mu}_g$	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 1.1: Converging individual and group performance

The converging of the mean target times over a window of 10 trials is shown graphically in Figure 1.4.

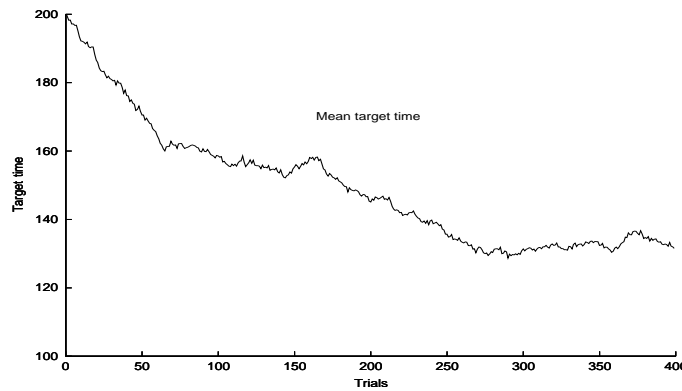


Figure 1.4: Converging group performance, mean target time, $\hat{\mu}_g$

All mean traversal and target times are lower than the ones produced by the initial random behavior application, at $n = m = 0$, by a statistically significant amount.

Analysis By analyzing the Q-tables produced during this experiment we found that the robots learned policies dominated by the *direct target approach* behavior, but with significant portions of the *wall-following* behavior. These policies made the robots follow a common circular path from the source to the sink and back, while keeping to the left side to minimize the interference with robots coming from the opposite direction. In most of the experiments, one or two of the robots would learn specialized policies containing disproportionately small amounts of the *wall-following* behavior. When observing the behavior of the specialized robots, it appeared that they effectively cut across the empty center of the common circular path, benefiting from the general dispersion of the rest of the group. The policies show that a structuring of the group’s spatial organization is partly responsible for the improvement in performance.

We also found evidence that the robots adapted their temporal organization according to the presence and orientation of other robots. From each robot’s Q-table, using only the input states where yellow markings were visible, we calculated the number of times the *stopping* behavior had the highest Q-value. We called this number the *yellow stop-rate* of a robot. The yellow stop-rates were significantly lower than the stop-rate for states in which other robots were visible without displaying any yellow markings, i.e., states where the other robots were seen from the front rather than from the back. The yellow stop-rate was also significantly lower than the corresponding *red* and *green* stop rates. This indicates that the robots in general learned to stop only when they saw a robot without visible yellow markings. Intuitively, visible yellow or rear markings indicate that the other robot is facing, and likely moving, away, and hence not likely to cause interference. Simply put, the robots learned to stop for oncoming traffic only.

1.4.3 Conclusions

Our analysis of the learned behaviors showed that they had natural interpretations at the group level. The tendency to keep to the left can be seen as the formation of two lanes. The observed utilization of the center space, by one or two of the robots, can be seen as specialization. Our analysis, together with the demonstrated improvements in performance, verify our hypothesis that individual learning is sufficient to produce recognizable forms of group behavior in our domain.

1.5 Task-allocation through specialization

In Section 1.4 we demonstrated a distributed learning algorithm that consistently improved a group’s performance on a cooperative transportation problem by

minimizing the negative effects of the group dynamics. One of the ways in which the performance was improved was the specialization of the members in an initially homogeneous group. The specialization was a result of the algorithm’s sensitivity to the effects of the group dynamics. In this section we show how the algorithm can be used for handling MRTA problems.

MRTA is a problem domain that presents a third major challenge, in addition to those we reviewed in Sections 1.1 and 1.4. This challenge concerns the fact that if the effects of the group dynamics are considered explicitly, MRTA is an \mathcal{NP} -complete problem [10]. This inherent complexity of MRTA implies that it is not feasible to find optimal solutions as the size of the system, i.e., the number of tasks and robots, grows. For such high-complexity problems, solutions must be found through heuristic algorithms that rely on appropriate assumptions about the problem to construct solutions of the necessary quality. Machine learning methods are commonly used to produce the necessary heuristics for such algorithms.

In this section we present the TAVC model of MRTA and the associated TAVC algorithm [11]. Our TAVC model explicitly represents the effects of group dynamics and allowed us to develop the TAVC algorithm, whose main contribution is its sensitivity to those effects. The TAVC algorithm uses Q-learning to make individual estimates of task utilities; in it each robot follows a *greedy* rule for choosing its next task based on its estimated utility. The algorithm goes beyond traditional greedy algorithms [8] in that the utility estimates are sensitive to the effects of the group dynamics produced by the current task allocation.

1.5.1 Multi-robot task allocation

Multi-robot task allocation (MRTA), is the problem of allocating a set of tasks to the robots in a group in a way that optimizes some group-level performance metric. This problem is receiving growing attention. Recently, Gerkey and Mataric [14, 13] presented a task- and architecture-independent framework in which to study and compare existing empirical solutions. The solution methods include L-ALLIANCE [26], a system for MRTA based on local task allocation decisions only. Robots using L-ALLIANCE individually estimate their own progress and the progress made by other robots, and take over or acquiesce tasks according to those estimates. However, the progress estimation functions are problem-specific and there is no explicit consideration of the effects of interaction.

The M+ algorithm [4] uses a centralized task allocation mechanism based on negotiation through the Contract Net Protocol. The algorithm relies on pre-specified capabilities and does not consider effects of interaction. Also, a distributed negotiation system carries with it a high communication overhead that introduces problems when scaling the system up to handle larger numbers of robots.

Similarly, MURDOCH [12] used a negotiation mechanism based on the Contract Net Protocol for task allocation. However, MURDOCH leaves open a *suitability* parameter to be filled in by the individual robots when bidding for

tasks. Though it is was done in their work, the suitability value could reflect the effects of the group dynamics.

The broadcast of local eligibility algorithm [29] defines a similar parameter called *eligibility*. The eligibility value is calculated by the robots individually and used in an inter-robot coordination mechanism based on port-arbitrated behaviors (PABs). PABs are a generalization, for multi-robot systems, of the traditional behavior coordination mechanisms. As in MURDOCH, the locally calculated eligibility could reflect the effects of the group dynamics but was not used in that manner.

Among previous MRTA algorithms, none deal explicitly with the effects of group dynamics. Our TAVC algorithm addresses that deficiency.

1.5.2 The vacancy chain model

The inspiration for the TAVC algorithm is the vacancy chain resource distribution process [7] common in human and animal societies. The typical example of resource distribution through a vacancy chain is a bureaucracy, where the retirement of a senior employee creates a vacancy to be filled by a less senior employee. The promotion, in turn, creates a second vacancy to be filled, and so on. The initial retirement results in a chain of vacancies linked by the promotions. The resources that are distributed in this example are the positions, and the consumers are the employees. Chase [7] has proposed that major human consumer goods such as cars, boats, and airplanes, also move through vacancy chains and that vacancy chains are common in other species as well.

Inspired by this process we developed a model of MRTA that describes how task values and task processing times combine to produce task allocations among greedy consumers [10]. Our model has the important property of breaking the group performance down into individual robot contributions, a property that allows us to use distributed control algorithms.

According to our TAVC model, any number of robots can be assigned to tasks from a given *class*. When a robot, j , is assigned to a task from a class, i , currently being *serviced* by a set of robots, J , we say that *service-slot*, (i, J, j) is being filled. A particular set of robots, J , servicing the same task, i , will have a task processing frequency, $c_{i,J}$, dependent on the degree to which the robots are able to work concurrently without interference. The difference in the group's task processing frequency before and after robot j joined, together with the task value, w_i , defines robot j 's contribution. We call this contribution, which can be negative, the *slot-value*, $s_{i,J,j}$. The formal definition of the slot-value is given in Equation 1.2.

$$s_{i,J,j} = w_i(c_{i,J \cup \{j\}} - c_{i,J}) \quad (1.2)$$

Assigning a robot to a task can lead to a decrease in the task processing frequency of the group of robots servicing tasks from that class. The slot-value then becomes negative. When all the available service-slots for a class of tasks have negative slot-values, we say the task is *saturated*.

When the service-slots are allocated optimally, a failure in a robot servicing a high-value task, or the creation of a new class of high-value tasks, will result in empty high-value service-slots that must be re-allocated for the system to perform optimally. If the value of a vacant slot is greater than the value of the occupied service-slots, the vacant slot will have to be filled in order to restore optimal performance. In vacancy chain terminology, a vacant, high-value service-slot is a resource to be distributed among the robots.

We make two simplifying assumptions: the robots are identical and the slot-values decrease monotonically with the number of robots servicing a class of tasks. When these assumptions are true, the task allocation algorithm can be distributed by letting each robot individually optimize the value of the service-slot it occupies. We will relax the first assumption later, in Section 1.6. Even with these simplifying assumptions, the optimal task allocations cannot be predicted without a model of the group dynamics, and such models are, as discussed in Section 1.2.1, difficult to construct.

1.5.3 The vacancy chain algorithm

In the TAVC algorithm, a task corresponds to an *action* in RL terms. Each robot uses a Q-table to keep individual estimates of task utilities and chooses its next task using an ϵ -greedy selection function.

The task processing is implemented in terms of pre-programmed high-level behaviors, one per each of the available tasks. The state-space describes what class of tasks the robot last serviced. This state-space allows the robots to learn policies that are not dedicated to one class of tasks; the learned policies can switch between classes of tasks in order to construct optimal task sequences. The action space corresponds to the available classes of tasks.

If a robot consistently occupies a service-slot that is suboptimal due to a high level of interference, the high mean traversal time for that class of tasks will reduce the estimated utility for that class of tasks until it falls below the estimated utility of the class of tasks that contains of the optimal service-slot. This change in estimated utility will attract the robot to the optimal slot. The robots keep a constant level of exploration, allowing them to identify and migrate to new vacancies.

All the robots in our demonstration used the same adaptive, behavior-based controller [22]. However, our vacancy chain task allocation algorithm is independent of the underlying architecture, being defined purely in terms of distributed reinforcement learning of task utilities.

1.5.4 Experimental validation of TAVC

In this section we present a set of experiments that had two aims. First, the experiments were designed to demonstrate the validity of the TAVC model, i.e., show that a simulated multi-robot system would behave as predicted by the model. Second, they were designed to demonstrate that the TAVC algorithm im-

proved system performance beyond what can be expected from pre-programmed solutions.

The ‘initial task distribution’ experiment was designed to demonstrate that the interactions between the robots could produce specializations where a number of robots would choose a task with sub-optimal value as a result of the effects of the group dynamics on individual utility estimates.

The ‘vacancy from robot breakdown’ experiment was designed to demonstrate that it is possible for the TAVC algorithm to recover from individual robot failures. The TAVC model predicts that the failure of a robot servicing tasks from a high-value class of tasks creates a vacant service-slot that is distributed among the robots currently servicing lower-value service-slots. This experiment was designed to show how the individual task utility estimates of the robots change in a way that leads to exactly this kind of migration between classes of tasks by a single robot. This experiment also demonstrates that our TAVC algorithm can perform better than a pre-programmed system.

The ‘breakdown without vacancy’ experiment was designed as a control experiment, where the failure of a robot leads to a situation where the TAVC model predicts no change. This experiment was designed to show that the TAVC model’s prediction was correct by showing that the task allocation in the simulated system did not change.

The ‘changing sub-task values’ was also a control experiment. Again a vacancy was created, but this time not from a robot failure, but from a change in sub-task values. This experiment was designed to demonstrate that the simulated system produces the patterns predicted by the TAVC model regardless of the underlying reason for a change in the predicted task allocation.

The four experiments above were also designed to show whether the TAVC algorithm could improve the performance of the system to a level significantly above that of a system where task were randomly allocated to robots. Random allocation is a pessimistic approximation of a pre-programmed system, because it might be possible to make universally applicable assumptions, such as ‘minimize changes between classes of tasks’ or ‘distribute robots evenly among the classes of tasks’, that would improve the performance of general MRTA algorithms above the performance of random allocation. It is, however, not unreasonable to use a pessimistic approximation since, for most heuristics, one can construct pathological cases where they would perform significantly below the random allocation solution. In such cases, our TAVC algorithm would outperform both the pre-programmed and the randomized solutions.

The simulated environment The experiments presented in this section were done in simulation using six robots identical to those described in Section 1.4.2. The simulated environment was an 8 by 12-meter rectangular area with two unique laser bar-codes along one of the longer sides of the rectangle, indicating sources, and two along the opposing side indicating sinks. A depiction of the simulated environment is presented in Figure 1.5; common paths between the sources and sinks are indicated by dashed arrows.

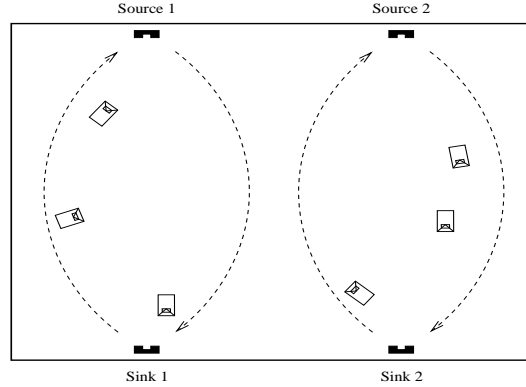


Figure 1.5: The simulated environment used for experiments on TAVC

Controller architecture All of the task-related behaviors used by the TAVC algorithm were implemented as a collection of sub-behaviors responsible for different aspects of a safe and efficient task execution.

These sub-behaviors were:

- **obstacle avoidance:** avoided obstacles detected (by laser range finder) in the desired path.
- **visible target approach:** approached the desired target when it was visible (to the laser range finder).
- **target location approach:** approached the location of the desired target when the target itself was not visible.
- **wall following:** followed the walls to the first available target when the desired target was not visible and localization (based on odometry) was deemed to be inaccurate.

The localization was deemed to be inaccurate whenever the desired target was not visible, but should have been so according to the robot's estimated position and orientation. On encountering a target, the localization estimate was reset and again deemed to be accurate.

Learning parameters The Q-tables were initialized with random values between -0.1 and 0.1 , the learning rate, α , was set to 0.1 , and the discount factor γ was set to 0.95 . For action selection we used a greedy- ϵ strategy [28], where ϵ was set to 0.1 .

We call the circuit with the highest related reward the *high-value* circuit and, correspondingly, the circuit with the lowest related reward is called the *low-value* circuit. In order to demonstrate task distribution through vacancy chains, we estimated circuit values that would produce an initial task allocation

where three robots were allocated tasks from the high-value circuit and three were allocated tasks from the low-value circuit.

In order to produce this task allocation, a situation needed to be set up in which it was less attractive to be one of four robots servicing the high-value circuit than to be one of three servicing the low-value circuit. This constraint on the reward function is presented formally in Equation 1.3.

$$\forall(x, y). \bar{r}_{x,4} < \bar{r}_{y,3} \quad (1.3)$$

We also wanted to demonstrate the filling of a vacancy on the high-value circuit. In order for a vacancy in the high-value circuit to be filled, it had to be the case that it was more attractive to be the third robot on that circuit than to be the third robot on the low-value circuit. This is expressed formally in Equation 1.4, where p denotes the preferred circuit.

$$\forall(x \neq p). \bar{r}_{x,3} < \bar{r}_{p,3} \quad (1.4)$$

We empirically estimated the relevant mean traversal times. To satisfy the given constraints we chose the circuit values, as defined in Equation 1.1, to be $w_1 = 2200$ and $w_2 = 2000$. The optimal allocation pattern emerges as each robot’s utility estimates converge as a result of individual experience only. The fact that the robots do not know the task values allows new allocation patterns to emerge when external factors, such as the task values or the group size, change.

1.5.5 Results

For each experiment we defined a convergence period and a stable period according to the stability of the system performance. Our student-t tests for statistical significance were all done at a 90% confidence level.

Initial task distribution We performed 20 individual trials of 5 hours, each averaging 3000 traversals or 500 traversals per robot. The convergence period was 1.25 hours.

To show the structure that emerged we consider the last target visited by each robot, which yields seven possible system states. We refer to each state using the notation $h : l$, where h is the number of robots whose last target was on the high-value circuit. Correspondingly, l is the number of robots whose last target was on the low-value circuit. The rows labeled *TAVC* in Table 1.2 show the mean, $\hat{\mu}$, and standard deviation, s , of the time the system spent in each of the states when controlled by the TAVC algorithm. The values are percentages of the total stable period. The rows labeled *R* describe the same values for a set of 20 control trials using a group of robots where tasks were allocated at random.

The row labeled *T* in Table 1.2 lists the number of different ways to choose a sample of size n from a population of size m , as a percentage of all possible samples, according to Equation 1.5. The time distribution produced by the six

State		0:6	1:5	2:4	3:3	4:2	5:1	6:0
TAVC	$\hat{\mu}$	0.1	2.8	19.3	44.5	27.6	5.3	0.4
	s	0.1	1.7	5.4	6.0	7.0	2.6	0.4
R	$\hat{\mu}$	1.0	7.4	22.3	33.7	25.3	9.3	1.0
	s	0.6	1.6	1.5	2.2	1.6	1.4	0.6
T		1.6	9.4	23.4	31.2	23.4	9.4	1.6

Table 1.2: State/time distributions with six robots

random controllers is closely aligned with this theoretical estimate of random distribution, though the differences are statistically significant.

$$T = \frac{m!}{n!(m-n)!2^m} \quad (1.5)$$

The two time distributions given in Table 1.2 are presented as histograms in Figure 1.6 with the standard deviation indicated by the error bars for each state.

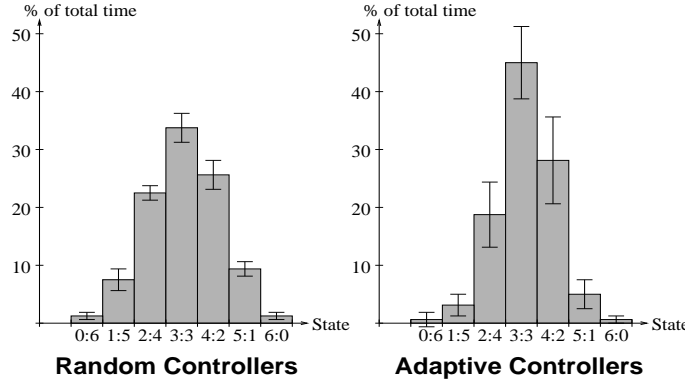


Figure 1.6: State/time distribution with six robots

The increase in the amount of time spent in state 3 : 3 is statistically significant. The time the adaptive group spends in state 3 : 3 is also significantly higher than the time spent in any of the other states.

Figure 1.7 presents the mean performance of a group of robots controlled by the TAVC algorithm over both the convergence period and the stable period. The TAVC performance is indicated by the thick, solid line. The performance of random allocation is indicated by the dashed line. The values used for the performance plots presented here are 10.0 and 1.0 for the high-value and low-value circuits, respectively.

The performance of a group of robots controlled by the TAVC algorithm is significantly higher than the performance of random allocation.

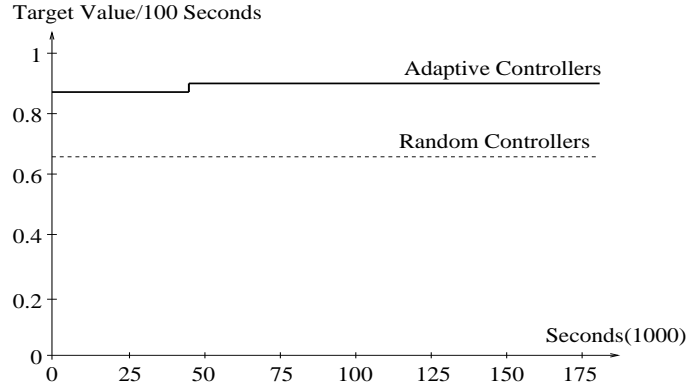


Figure 1.7: Performance with six robots

Vacancy from robot breakdown This experiment started with the system in the converged state of the ‘initial task distribution’ experiment. We choose at random one of the robots dedicated to the high-value circuit for removal. According to the MRTA model, that circuit then had a vacancy. We performed 20 experiments of 7.5 hours. The convergence period was 2.5 hours. The values of the time distribution in the stable period are given in Table 1.3 and a graphical presentation is provided in Figure 1.8.

State		0:5	1:4	2:3	3:2	4:1	5:0
TAVC	$\hat{\mu}$	0.3	6.5	35.7	46.2	10.4	0.6
	s	0.3	3.8	9.0	8.6	4.7	0.6
R	$\hat{\mu}$	2.2	13.9	31.0	32.9	16.1	2.8
	s	0.8	2.6	2.2	2.3	2.6	0.4
T		3.1	15.6	31.3	31.3	15.6	3.1

Table 1.3: State/time distribution after breakdown creating a vacancy

The converged controllers kept the system in state 3 : 2 for a significantly larger amount of time than a group of five random controllers. The group adapted its structure from one that promotes the 3 : 3 state to one that promotes the 3 : 2 state. This change implies that a robot from the low-value circuit has filled the vacancy we created in the high-value circuit.

The performance data presented in Figure 1.9 show that on the removal of a robot from the high-value circuit, the performance drops sharply. After the re-convergence period, however, the performance rises again to a level that is significantly higher than the performance of five random controllers. The new performance level is also significantly higher than the mean performance, over 20 trials, of a group of robots controlled by a static TA algorithm optimized for six robots. The mean performance of the static group is indicated by the thin solid line.

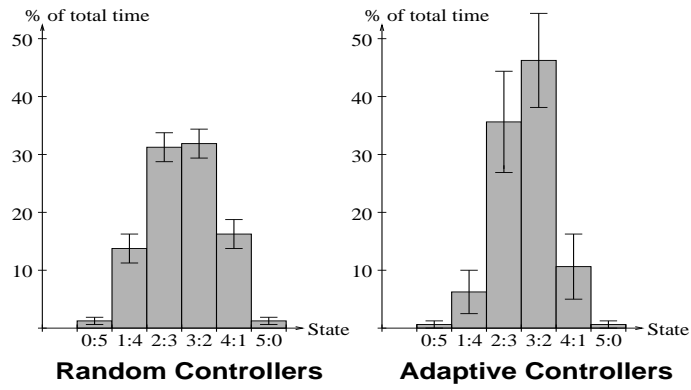


Figure 1.8: State/time distribution after breakdown creating a vacancy

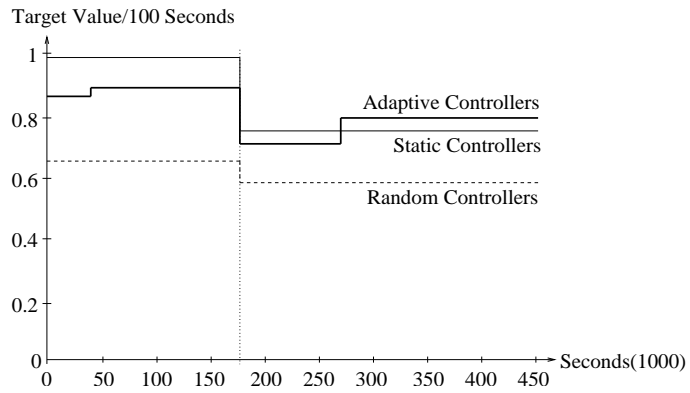


Figure 1.9: Performance during breakdown creating a vacancy

Breakdown without vacancy This experiment had an initial state identical to the ‘vacancy from robot breakdown’ experiment. This time, however, we removed a robot from the low-value circuit. We performed 20 trials. The convergence time was 2.5 hours. The time distribution during the stable period of this experiment, presented in Table 1.4, was not significantly different from the distribution produced during the experiment where a vacancy was created by the removal of a robot.

State	0:5	1:4	2:3	3:2	4:1	5:0
$\hat{\mu}$	0.3	6.7	34.6	47.1	10.5	0.7
s	0.3	3.7	9.2	9.4	3.8	0.4

Table 1.4: State/time distribution after breakdown not creating a vacancy

As shown in Figure 1.10, performance fell significantly when the robot was removed, but remained significantly higher than the performance of random allocation among five robots. There was no significant difference in the performance during the stable period of this experiment and the stable period during the experiment where a vacancy was created. Also, there was no significant difference in performance between the convergence period and the stable period. This consistency in performance likely reflects the fact that the group structure remained unchanged.

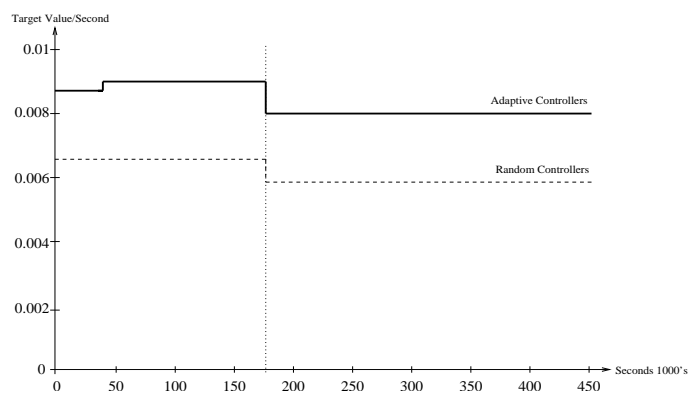


Figure 1.10: Performance during breakdown not creating a vacancy

This result demonstrates that our algorithm produces the group structure required for VC distribution, independent of which robot should fail.

Changing sub-task values The start state of this experiment was the converged state of the ‘vacancy from robot breakdown’ experiment. We switched the values of the tasks so that the high-value circuit became the low-value circuit and vice versa. According to the TAVC model, the new setup had a vacancy on

the new high-value circuit. We ran 20 individual trials using this setup. The convergence time was 50 minutes.

The time distribution for the stable period, given in Table 1.5, had significantly higher values for states 3 : 2 and 4 : 1 than the time distribution produced by the experiment where a vacancy was created by the removal of a robot. Compared to the time distribution produced by a group of five robots where tasks were allocated at random, the distribution produced by the TAVC algorithm had a significantly higher value for state 3 : 2 and significantly lower values for *all* other states.

State	0:5	1:4	2:3	3:2	4:1	5:0
$\hat{\mu}$	0.2	5.4	30.3	48.7	14.4	0.9
s	0.5	3.9	5.2	6.1	1.6	0.3

Table 1.5: State/time distribution after changing sub-task values

When the task values were changed the performance, presented in Figure 1.11, fell significantly. After the re-convergence period, however, it was back up to a level that was not significantly different from the initial level.

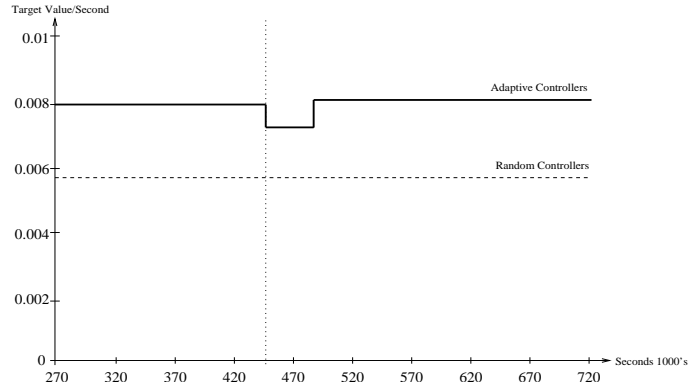


Figure 1.11: Performance during changing sub-task values

This experiment showed that the TAVC algorithm is not sensitive to how a vacancy is created, whether by robot failure or by a change in task values.

1.5.6 Conclusions

In the ‘initial task distribution’ experiment we showed that the performance of the TAVC algorithm was significantly above the performance level of a system where tasks were allocated at random.

In the ‘vacancy from robot breakdown’ experiment we showed that the performance of the TAVC algorithm after re-convergence was higher than a mis-configured pre-programmed solution. Assuming that a model of the group dy-

namics is not available in general, a pre-programmed solution must be based on some heuristics that will occasionally lead to sub-optimal task allocations. Our experiment shows that the TAVC algorithm has the potential to improve on pre-programmed solutions.

Together, the experiments validate our TAVC model by reproducing the task allocations predicted by the model. They also demonstrate that the TAVC algorithm outperforms random allocation and can outperform pre-programmed solutions under the assumption that pre-programmed solutions must occasionally make sub-optimal task allocations.

1.6 Performance-related specialization

The analysis of the TAVC experiments presented in Section 1.5.4 indicated that persistence with a class of tasks was a determining factor in the estimated utility that class. By persisting with a class of tasks a robot contributes to the level of interference around the related source and sink. This leads to a general fall in the estimated utility of that class of tasks for all the robots in the group. When the fall leads to a situation where one of the robots involved estimates another class of tasks to have a higher utility, that robot will migrate to the class of tasks. We refer to this effect of the group dynamics as *forced migration through persistence*.

A severe limitation of our TAVC algorithm as presented in Section 1.5 is that it is limited to homogeneous groups of robots. To remove this limitation, we modified it to be sensitive to the individual performance levels of the robots in a *heterogeneous* group of robots.

We hypothesized that this could be achieved by using the perceived ability of robots to force migration through persistence. If performance could be related to persistence in a way that made robots with high performance levels more persistent, then the observed effect of forced migration through persistence would allocate robots with high performance levels to high-value classes of tasks and robots with low performance levels to low-value classes of tasks. This would produce a completely communication-free algorithm for MRTA that is sensitive to individual robot performance levels as a result of being sensitive to the general effects of the group dynamics.

In this section we present action selection through a Boltzmann-softmax function as a way of relating a robot's performance level to its persistence. We present experimental evidence showing that the TAVC algorithm, when using a Boltzmann-softmax action selection function, is sensitive to the performance level of the individual robots to a degree which produces a significant difference in system performance over a system that allocates tasks to robots at random.

1.6.1 Action Selection

The TAVC algorithm used a greedy- ϵ action selection function. The problem with using this function is that all the tasks, apart from the one with the

highest estimated utility, have equal probability, ϵ , of being explored. With a Boltzmann-softmax action selection function however, the probability of trying a task is correlated with the relative estimated utility of that task. The probability of selecting action, a , according to a Boltzmann-softmax function is given in Equation 1.6.

$$P(a) = \frac{e^{Q_t(a)/\tau}}{\sum_{b=1}^n e^{Q_t(b)/\tau}} \quad (1.6)$$

A robot that on average has task processing time of \bar{p} will have a difference in estimated task utility that correlates with the expression $\frac{w_h - w_l}{\bar{p}}$ where w_h and w_l are the values of high- and low-value tasks respectively. A robot with a lower mean task processing time will have a correspondingly higher difference in estimated utility of high- and low-value tasks. This greater difference in estimated utility implies a greater persistence on high-value tasks in terms of a higher probability of servicing these tasks. This higher persistence can hypothetically lead to fast robots forcing the slow robots to migrate from high-value tasks to low-value tasks.

The experimental results presented in this section demonstrate that the ability of the Boltzmann-softmax action selection function to differentiate between suboptimal actions can indeed be made to work on an inter-robot level, as a mechanism for allocating high-value tasks to high-performance robots without explicit communication.

1.6.2 Experimental Validation

We designed an experiment to test our hypothesis that a Boltzmann-softmax action selection rule could relate a robot’s individual performance level to its persistence and that a difference in persistence would lead to high-value tasks being allocated to robots with high performance levels and low-value tasks to robots with low performance levels.

The simulated environment was identical to the one used in the experiments presented in Section 1.5.4.

Control architecture We made one change to the system architecture presented in Section 1.5.4. The original group of six robots was made heterogeneous by dividing it into two sub-groups of three. One sub-group was made to operate at a default speed of 300 mm/sec while the other was made to operate at a default speed of 200 mm/sec. These speeds were chosen as equidistant points from the default operating speed of 250 mm/sec used in the homogeneous robot experiments. The equidistance was intended to preserve the allocation of three robots to the high-value circuit and three to the low-value circuit.

The general difference in default speed of operation encompasses more specific differences in robot morphology and task competence. Hence, the results presented here generalize to heterogeneous multi-robot systems where the differ-

ences between the participating robots can be expressed in terms of differences in task processing speed.

Learning parameters The Q-tables were initialized with random values between -0.1 and 0.1 , the learning rate, α , was set to 0.1 , and the discount factor γ was set to 0.95 .

For the Boltzmann-softmax function, the temperature parameter, τ , was set empirically to 0.005 . Because we wanted the system to remain adaptive to changes in the environment we did not decrease τ over time, as is common. The reward structure was identical to the one presented in Section 1.5.4.

1.6.3 Results

We defined a convergence period of 15 hours based on the stability of the system performance. To examine the performance we consider which of the robots visited the high-value circuit last. We used three fast and three slow robots, yielding fifteen possible system states. We refer to each state using the notation $f : s$, where f is the number of fast robots whose last target was on the high-value circuit. Correspondingly, s is the number of slow robots whose last target was on the high-value circuit. The columns with subscript VC in Table 1.6 show the mean and standard deviation of the time the system spent in each of the states while running the TAVC algorithm. The values are percentages of the total stable period. The columns with subscript R describe the same values for a set of 15 trials using randomly tasks allocations.

$f : s$	$\hat{\mu}_{VC}$	s_{VC}	$\hat{\mu}_R$	s_R	T	$\hat{\mu}_{VC} - \hat{\mu}_R$	$\frac{\hat{\mu}_{VC} - \hat{\mu}_R}{\hat{\mu}_R}$
0:0	0.2	0.2	1.4	0.5	1.5	-1.3	-0.88
0:1	1.5	1.2	3.4	1.7	4.7	-3.2	-0.67
0:2	2.0	2.6	5.0	1.4	4.7	-2.7	-0.58
0:3	0.5	0.7	2.0	1.1	1.5	-1.1	-0.70
1:0	2.9	1.6	3.6	1.4	4.7	-1.8	-0.38
1:1	12.1	4.9	13.5	3.3	14.1	-2.0	-0.14
1:2	14.4	6.1	14.4	2.4	14.1	0.3	0.02
1:3	4.1	2.7	4.0	1.4	4.7	-0.6	-0.13
2:0	7.0	5.0	5.1	1.2	4.7	2.3	0.48
2:1	19.4	7.4	15.7	2.3	14.1	5.30	0.37
2:2	18.5	5.0	14.7	3.0	14.1	4.4	0.3
2:3	5.7	3.6	4.4	2.4	4.7	1.0	0.20
3:0	2.7	4.1	1.7	0.6	1.5	1.2	0.78
3:1	5.2	4.6	5.2	1.6	4.7	0.5	0.10
3:2	3.3	2.6	4.2	1.2	4.7	-1.4	-0.29
3:3	0.7	0.7	1.5	0.8	1.5	-0.9	-0.56

Table 1.6: State/time distribution for six heterogeneous robots

The column labeled T lists the probability of choosing a sample of size f from a population of $g = 3$ fast robots as well as choosing a sample of size s from a population of $h = 3$ slow robots according to Equation 1.7. The time distribution produced by the six random controllers is closely aligned with this theoretical estimate, though the differences are statistically significant.

$$T = \frac{100g!h!}{f!(g-f)!s!(h-s)!2^g2^h} \quad (1.7)$$

The difference between the state-time distribution produced by the TAVC algorithm and the distribution produced by random task allocation is presented in the column labeled $\hat{\mu}_{VC} - \hat{\mu}_R$. This difference is presented as a percentage of the mean of the random distribution, $\hat{\mu}_R$, for each state in the last column, labeled $\frac{\hat{\mu}_{VC} - \hat{\mu}_R}{\hat{\mu}_R}$.

The time distributions for the random allocation and the TAVC algorithm, i.e., columns 2 and 4 of Table 1.6, are presented graphically in Figure 1.12.

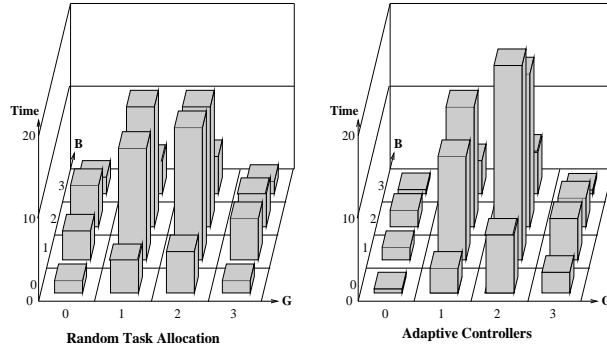


Figure 1.12: State/time distributions for six homogeneous robots

The differences between the distributions produced by the TAVC algorithm and the random task allocation, i.e., the last two columns, are also presented graphically in Figure 1.13.

Over 15 experiments, the difference in time spent in state 0 : 3 is statistically significant. The second histogram in Figure 1.13 shows the differences relative to an underlying random distribution. The optimal state, 0 : 3 stands out as the state with the highest relative increase. This validates the TAVC model by confirming that the group behavior has converged to promote the state predicted by the TAVC model. The performance data shows that the TAVC algorithm produces a total target value throughput of $0.81 * 10^{-2}$ per second. This is significantly higher than the performance of a group of robots where tasks are allocated randomly, $0.78 * 10^{-2}$ per second. Together, the time distribution data and the performance data show that the TAVC algorithm improves the group's performance by promoting the dedicated service structure predicted by the TAVC model. The transition functions also indicate this. The fast robots, on average, have higher estimated utilities for servicing the high-value circuit

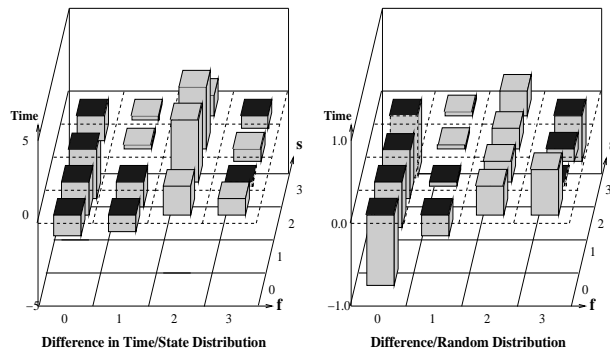


Figure 1.13: Difference in the state/time distributions produced by TAVC and random task allocation

while the slow robots have higher estimated utilities for servicing the low-value circuit. Finally, as predicted, the mean difference between the estimated utilities are greater for the fast robots than for the slow robots.

1.6.4 Conclusions

Our results demonstrate that when using a Boltzmann-softmax action selection function, there is an observable migration through persistence effect. This validates our hypothesis that a performance-related persistence affects the group dynamics in a way that promotes task allocations where high-value tasks are allocated to high-performance robots and low-value tasks are allocated to low-performance robots.

Our results demonstrated a significant improvement in group performance over a control group using random task allocation. This shows that the TAVC algorithm is able to produce task allocations that are sensitive to the underlying group dynamics, including differences in individual robot performance levels.

1.7 Summary and conclusions

We have discussed three main challenges in producing high-quality solutions to MRTA problems: modeling group dynamics, problem complexity, and communication overhead. We then presented a learning algorithm that used individual learning and feedback that reflected the effects of the prevailing group dynamics. We presented experimental evidence that the algorithm produces effects that improve group performance. We then showed how one such effect, specialization, can be used to construct a MRTA algorithm based on a model of MRTA as TAVC. The TAVC algorithm is an important contribution to existing MRTA algorithms in that it is demonstrably sensitive to the effects of group dynamics, including the effects of different individual performance levels among the robots in a group. We presented experimental evidence that validated the

TAVC model of MRTA and showed that the TAVC algorithm can out-perform pre-programmed solutions. Finally we showed how using a Boltzmann-softmax action selection function in the TAVC algorithm can relate a robot's performance level to its persistence level. We also presented experimental evidence that a performance-related persistence level affects the group dynamics and produces task allocations where the robots with the highest performance levels are allocated the high-value tasks. Therefore, we hope that the main contributions of this work, a new model of group dynamics in MRTA and associated algorithm, will help address the three main difficulties outlined at the start.

Acknowledgments

This work was supported in part by a Department of Energy (DOE) Robotics and Intelligent Machines (RIM) grant DE-FG03-01ER45905 and in part by an Office of Naval Research (ONR) Defense University Research Instrumentation Program (DURIP) grant 00014-00-1-0638.

The Player/Stage robot device server and simulator was developed originally by Brian Gerkey, Kasper Støy, Andrew Howard and Richard Vaughan at the Center for Robotics and Embedded Systems (CRES) at the University of Southern California. The laser bar-code technology was developed by Andrew Howard. These technologies are freely available under the GNU General Public License.

Bibliography

- [1] ARKIN, Ronald C., *Behaviors-based robotics*, MIT Press Cambridge, Massachusetts (1998).
- [2] BALCH, Tucker R., "The impact of diversity on performance in multi-robot foraging", *The proceedings of the Third International Conference on Autonomous Agents (Agents'99)* (Seattle, WA, May 1 - 5, 1999) (O. ETZIONI, J. P. MÜLLER, AND J. M. BRADSHAW eds.), ACM Press, 92–99.
- [3] BALCH, Tucker R., Zia KHAN, and Manuela M. VELOSO, "Automatically tracking and analyzing the behavior of live insect colonies", *Proceedings of the Fifth International Conference on Autonomous Agents (Agents'01)* (Montreal, Canada, May 31 - June 4, 2001) (J. P. MÜLLER, E. ANDRE, S. SEN, AND C. FRASSON eds.), ACM Press, 521–528.
- [4] BOTELHO, Silvia, and Rachid ALAMI, "M+ : a scheme for multi-robot cooperation through negotiated task allocation and achievement", *Proceedings of the 1999 IEEE International Conference on Robotics and Automation (ICRA'99)* (Detroit, Michigan, May 10-15, 1999), IEEE Press, 1234–1239.
- [5] BOWLING, Michael, and Manuela M. VELOSO, "Rational and convergent learning in stochastic games", *Proceedings of the 17th International Joint*

- Conference on Artificial Intelligence (IJCAI'01)* (B. NEBEL ed.), Morgan Kaufmann, 1021–1026.
- [6] BRUCKER, Peter, *Scheduling Algorithms* Second ed., Springer Verlag Berlin, Germany (1998).
- [7] CHASE, Ivan D., Marc WEISSBURG, and Theodore H. DEWITT, “The vacancy chain process: a new mechanism of resource distribution in animals with application to hermit crabs”, *Animal Behavior* **36** (1988), 1265–1274.
- [8] CORMEN, Thomas H., Charles E. LEISERSON, Ronald L. RIVEST, and Clifford STEIN, *Introduction to algorithms* second ed., MIT Press Cambridge, Massachusetts (2001).
- [9] DAHL, Torbjørn S., Maja J. MATARIĆ, and Gaurav S. SUKHATME, “Adaptive spatio-temporal organization in groups of robots”, *Proceedings of the 2002 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'02)* (Lausanne, Switzerland, September 30 - October 4, 2002), IEEE Press, 1044–1049.
- [10] DAHL, Torbjørn S., Maja J. MATARIĆ, and Gaurav S. SUKHATME, “Scheduling with group dynamics: A multi-robot task allocation algorithm based on vacancy chains”, *Tech. Rep. no. CRES-002-07*, Center for Robotics and Embedded Systems, University of Southern California, Los Angeles, CA (2002).
- [11] DAHL, Torbjørn S., Maja J. MATARIĆ, and Gaurav S. SUKHATME, “Multi-robot task-allocation through vacancy chains”, *Proceedings of the 2003 IEEE International Conference on Robotics and Automation (ICRA'03)* (Taipei, Taiwan, September 9 - 14, 2003), IEEE Press, 2293–2298.
- [12] GERKEY, Brian P., and Maja J. MATARIĆ, “Sold!: Auction methods for multi-robot coordination”, *IEEE Transactions on Robotics and Automation* **18**, 5 (October 2002), 758–768.
- [13] GERKEY, Brian P., and Maja J. MATARIĆ, “A formal framework for the study of task allocation in multi-robot systems”, *Tech. Rep. no. CRES-03-013*, University of Southern California, Center for Robotics and Embedded Systems, (2003).
- [14] GERKEY, Brian P., and Maja J. MATARIĆ, “A framework for studying multi-robot task allocation”, *Multi-Robot Systems: From Swarms to Intelligent Automata, Volume II, Proceedings of the 2nd International Naval Research Laboratory Workshop on Multi-Robot Systems* (Washington, DC, March 17 - 19, 2003) (A. C. SCHULTZ ET AL. eds.).
- [15] GERKEY, Brian P., Richard T. VAUGHAN, Kasper STØY, Andrew HOWARD, Gaurav S. SUKHATME, and Maja J. MATARIĆ, “Most valuable player: A robot device server for distributed control”, *Proceedings*

- of the 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'01) (Wailea, Hawaii, October 29 - November 3, 2001), 1226–1231.
- [16] GOLDBERG, Dani, and Maja J MATARIĆ, “Learning multiple models for reward maximization”, *Proceedings of the 17th International Conference on Machine Learning (ICML'00)* (Stanford, California, June 29 - July 2, 2000) (P. LANGLEY ed.), Morgan Kaufmann, 319–326.
 - [17] GOLDBERG, Dani, and Maja J MATARIĆ, “Design and evaluation of robust behavior-based controllers for distributed multi-robot collection tasks”, *Robot Teams: From Diversity to Polymorphism*, (T. R. BALCH AND L. E. PARKER eds.). A K Peters Ltd (2001), pp. 315–244.
 - [18] HAN, Kwun, and Manuela M. VELOSO, “Automated robot behavior recognition applied to robotic soccer”, *Proceedings of the 9th International Symposium on Robotics Research (ISRR'99)* (Snowbird, Utah, October 9 - 12, 1999), 199–204.
 - [19] LERMAN, Kristina, Asram. GALSTYAN, Alcherio MARTINOLI, and Auke J. IJSPEERT, “A macroscopic analytical model of collaboration in distributed robotic systems”, *Artificial Life* **7**, 4 (2001), 375–393.
 - [20] LITTMAN, Michael L., “Markov games as a framework for multi-agent reinforcement learning”, *Proceedings of the 11th International Conference on Machine Learning (ICML'94)* (New Brunswick, New Jersey, July 10 - 13, 1994) (W. W. COHEN AND H. HIRSH eds.), Morgan Kaufmann, 157–193.
 - [21] MATARIĆ, Maja J., *Interaction and Intelligent Behavior*, PhD thesis Massachusetts Institute of Technology Cambridge, Massachusetts (1994).
 - [22] MATARIĆ, Maja J, “Behavior-based control: Examples from navigation, learning, and group behavior”, *Journal of Experimental and Theoretical Artificial Intelligence, special issue on Software Architectures for Physical Agents* **9**, 2-3 (1997), 323–336.
 - [23] MATARIĆ, Maja J, “Reinforcement learning in the multi-robot domain”, *Autonomous Robots* **4**, 1 (1997), 73–83.
 - [24] MITCHELL, Tom M., *Machine Learning*, McGraw-Hill New York (1997).
 - [25] MYERSON, Roger B., *Game Theory: Analysis of Conflict*, Harvard University Press Cambridge, Massachusetts (1991).
 - [26] PARKER, Lynne E., “L-ALLIANCE: Task-Oriented Multi-Robot Learning in Behaviour-Based Systems”, *Advanced Robotics, Special Issue on Selected Papers from IROS'96* **11**, 4 (1997), 305–322.
 - [27] SETH, Anil K., “Modelling group foraging: Individual suboptimality, interference, and a kind of matching”, *Adaptive Behavior* **9**, 2 (2001), 67–91.

- [28] SUTTON, Richard S., and Andrew G. BARTO, *Reinforcement learning: an introduction*, MIT Press Cambridge, Massachusetts (1998).
- [29] WERGER, Barry B., and Maja J MATARIĆ, “Broadcast of local eligibility for multi-target observation”, *Distributed Autonomous Robotic Systems 4, Proceedings of the 5th International Symposium on Distributed, Autonomous Robotic Systems (DARS’00)* (Knoxville, Tennessee, October 4-6, 2000) (L. E. PARKER, G. BEKEY, AND J. BARHEN eds.), Springer, 347–356.
- [30] YAN, Helen, and Maja J MATARIĆ, “General spatial features for analysis of multi-robot and human activities from raw position data”, *Proceedings of the 2002 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS’02)* (Lausanne, Switzerland, September 30 - October 4, 2002), IEEE Press, 2770–2775.