

Emergent Robot Differentiation for Distributed Multi-Robot Task Allocation

Torbjørn S. Dahl¹, Maja J Matarić², and Gaurav S. Sukhatme²

¹ Norwegian Defence Research Establishment (FFI), Kjeller, Norway
Torbjorn-Semb.Dahl@ffi.no

² Center for Robotics and Embedded Systems
University of Southern California, Los Angeles, California
mataric|gaurav@usc.edu

Summary. We present a distributed mechanism for automatically allocating tasks to robots in a manner sensitive to each robot's performance level without hand-coding these levels in advance. This mechanism is an important part of improving multi-robot task allocation (MRTA) in systems where communication is restricted or where the complexity of the group dynamics makes it necessary to make allocation decisions locally. The general mechanism is demonstrated as an improvement on our previously published task allocation through vacancy chains (TAVC) algorithm for distributed MRTA. The TAVC algorithm uses individual reinforcement learning of task utilities and relies on the specializing abilities of the members of the group to produce dedicated optimal allocations. Through experiments with realistic simulator we evaluate the improved algorithm by comparing it to random allocation. We conclude that using softmax action selection functions on task utility values makes algorithms responsive to different performance levels in a group of heterogeneous robots.

1 Introduction and Motivation

Multi-robot task allocation (MRTA) algorithms for heterogeneous groups of robots have to be able to differentiate between robots based on their performance in order to optimize allocation. Existing MRTA algorithms [?, ?] generally do this based on hand-coded information about the task utilities relative to each robot. Using hand-coded task utilities, however, these algorithms are typically not sensitive to the effects of group dynamics, such as interference and synergy. These effects typically have to be estimated at runtime as they are difficult to model due to their volatility and complexity.

We have previously [?] presented a model of distributed MRTA as task distribution through *vacancy chains*, a distribution process common in animal and human societies [?]. We have also presented an algorithm based on task allocation through vacancy chains (TAVC) that validated our TAVC model and

demonstrated improvements in performance over existing MRTA algorithms in dynamic domains. The TAVC algorithm is a partial solution to the problem of modeling the effects of group dynamics. The original TAVC algorithm, however, only optimized allocation for homogeneous groups of robots and was not designed to differentiate between robots based on individual performance.

In this paper we show how softmax action selection [?], on an inter-robot level, has the effect of reliably allocating high-performance robots to high-value tasks. This provides us with a general mechanism for producing performance sensitive task allocations in a distributed manner. We demonstrate this mechanism by using it to extend our original TAVC algorithm. We also present experimental evidence that the resulting allocation leads to an improved group performance.

2 The Prioritized Transportation Problem

Cooperative transportation is a multi-robot MRTA problem where group dynamics can have a critical impact on performance. In the basic transportation problem, a group of robots traverse a given environment in order to transport items between the sources and the sinks. We call the time taken to traverse the environment once from a sink via a source and back to a sink the *traversal time*. To perform optimally on this task the robots must minimize the total traversal time. The basic transportation problem is one of the sub-problems of *foraging* [?]. If the locations of sources and sinks are given, the foraging problem is reduced to a problem of transportation. The prioritized transportation problem extends the basic transportation problem by dividing the sinks into sets of different priority.

When sources and sinks are spatially distributed into distinct pairs or *circuits*, the optimal allocation will have to dedicate each robot to one of these circuits in order to avoid the increased processing implied by crossing between circuits. We call transportation tasks in such environments *spatially classifiable* [?]. To optimize its performance on these transportation problems, a group of robots must strike the correct balance between the different values of tasks within a class and the different traversal times as defined by the current levels of interference and synergy on each circuit.

3 Task Allocation through Vacancy Chains

The inspiration for our TAVC algorithm is the vacancy chain process through which resources are distributed to consumers. The typical example of resource distribution through a vacancy chain is a bureaucracy where the retirement of a senior employee creates a vacancy that is filled by a less senior employee. This promotion, in turn, creates a second vacancy to be filled, and so on, resulting 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.

3.1 The TAVC Model

Here we generalize our previous TAVC model to include groups of *heterogeneous* robots. According to our model, any number of robots can be assigned to tasks from a given *class*. Two tasks are in the same class if doing one affects the efficiency with which it is possible to do the other. We restrict ourselves to problems where these classes are disjoint.

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 filled. A particular set of robots, J , servicing the same task, i , will have a mean task processing frequency, $c_{i,J}$, dependent on the degree to which the robots are able to work concurrently without interfering with each other. The difference in mean task processing frequency together with the task value, w_i , define the contribution of the last robot added or the last service-slot filled. We call this contribution, which can be negative, the *slot-value*, $s_{i,J,j}$. The formal definition is given in Equation 1. When assigning an additional robot to a task leads to a decrease in the task processing frequency, the slot-value correspondingly becomes negative.

$$s_{i,J,j} = w_i(c_{i,J \cup \{j\}} - c_{i,J}) \quad (1)$$

In a scenario where the service-slots are allocated optimally, i.e., where the total value of the filled service slots are maximized, a failure in a robot servicing a high-value task will result in an empty high-value service-slot that must be re-allocated to reestablish optimality. Expressed in the vacancy chain framework, a vacant, high-value service-slot is a resource to be distributed between the robots.

The TAVC model formalizes the system level performance-related effects of group dynamics in terms of individual time measurements. Hence, it can be used as a model for optimizing MRTA problems that satisfy the *greedy property* [?] such as homogeneous MRTA. To optimize problems that do not satisfy this property, such as heterogeneous MRTA, it is necessary to introduce elaborations and additional allocation mechanisms such as the emergent performance-based robot differentiation mechanism presented here.

3.2 The TAVC Algorithm

In the TAVC algorithm, each robot keeps a local estimate of task utilities and chooses its next task according to an action selection function. The TAVC algorithm uses Q-learning for task utility estimation. As Q-learning is time insensitive, it is necessary to make the temporal aspect of performance explicit in the reward function in order to use Q-learning to improve performance over

time. We used a reward function based on the last task processing time, t , and task value, w_i as presented in Equation 2

$$r = \frac{w_i}{t} \quad (2)$$

This reward function promotes the tasks with the highest value because these will on average provide a higher reward. However, if a robot consistently occupies a service-slot that is sub-optimal due to too much interference, the increased average traversal time will reduce the average reward for that slot below the average reward of the optimal service-slots. This change in average reward will attract the robot back to an optimal slot.

3.3 Softmax Action Selection as a Performance Differentiator

When robots keep a set of task related utilities, different functions can be used to select the next task to undertake based on these utilities [?]. With a Greedy- ϵ function, all the tasks, apart from the one with the highest utility, have equal probability, ϵ , of being explored. With a softmax function however, the probability of trying a suboptimal task is correlated with the relative estimated utility of that task. Our contribution in this paper is to demonstrate that the use of a Boltzmann softmax function on a task-selection level, has reliable effects on an inter-robot level, where it functions as a mechanism for allocating high-value tasks to high-performance robots. A robot that on average can service tasks in time \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 fast robot with a lower average service time \bar{p} will have a correspondingly higher difference between the estimated utility of high- and low-value tasks. Using a softmax action selection function, this greater difference in estimated utility theoretically translates into a greater probability of servicing high-value tasks. Such persistence will lead to the fast robots servicing high-value tasks and may, depending on the group dynamics and the task values, lead to the slow robots servicing low-value tasks.

4 Controller Architecture

All the robots in the experiments presented here used the same adaptive, behavior-based controller [?]. However, our TAVC algorithm and the performance sensitive robot differentiation mechanism we present here are independent of the underlying control architecture, being defined purely in terms of task utilities and the action selection function.

The robots were divided into two groups in order to make them heterogeneous. The first group was made to operate at a default speed of 300 mm/sec. The second group had a default speed of 200 mm/sec. The speeds were chosen

to be equidistant from the speed used in our homogeneous robot experiments [?] in order to preserve the distribution where three robots served the high-value circuit and three robots serving the low-value circuit. This general difference in task processing speed encompasses more specific differences such as differences in robot morphology and task competence. The results here therefore generalize to all heterogeneous multi-robot systems where the differences between the participating robots can be expressed in terms of differences in task processing speed.

The input- or state-space reflected which circuit the robot used for its last traversal. The action space corresponded to the available tasks. For two tasks and six robots this resulted in six two-by-two Q-tables or 24 Q-values. The robots used temporal difference Q-learning [?] to associate the different input states with one of the high level approach behaviors. 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 either a Greedy- ϵ or a Boltzmann softmax function. For the Greedy- ϵ function, ϵ was set empirically to 0.1 . For the Boltzmann softmax function, the temperature parameter τ was set empirically to 0.005 . With these values the experimentation rate was significant without being overwhelming. Because we wanted the system to remain adaptive to changes in the environment we did not decrease ϵ or τ over time, as is common.

5 Experimental Design

Having previously performed experiments to demonstrate how the TAVC algorithm with a Greedy- ϵ action selection function allocates tasks in a group of homogeneous robots [?], the work focused on testing our hypothesis that using a softmax action selection function would make a MRTA algorithm sensitive to differences in robot skill levels.

The experiments were done in simulation on the *Player/Stage* software platform. From experience, controllers written for the Stage simulator work with little or no modification on real Pioneers. The experiment used simulated Pioneer 2DX robots with SICK laser range-finders and PTZ color cameras and took place in a twelve by eight meter environment with two sets of sources and sinks. Figure 1 shows a graphical rendering of the simulated environment, with the sources and sinks labeled.

The sources and sinks were simulated bar-codes made from high-reflection material and recognizable by the laser range finder. 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 previous experiments with homogeneous robots [?], in order to produce a task allocation where three robots serviced one circuit and three robots serviced the other circuit, it was necessary that it was less attractive to the robots to be one of four robots servicing the high-value circuit than to be one

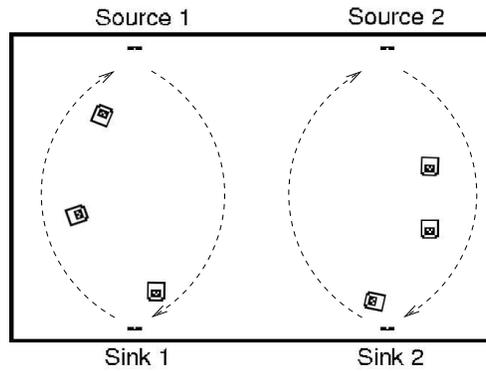


Fig. 1. The Simulated Environment with Circuits Indicated

of three servicing the low-value circuit. We empirically estimated the relevant average traversal times. To satisfy the given constraints we chose the circuit values as used in the reward-function given in Equation 2 to be $w_1 = 2200$ and $w_2 = 2000$. Assuming that the allocation with three robots on each circuit persists in the case of heterogeneous robots, the optimal allocation in terms of the TAVC model is to have the three fast robots on the high-value circuit and the three slow robots on the low-value circuit.

We had two aims for the heterogeneous robot experiment. First, to test whether the modified TAVC algorithm would produce the allocation predicted to be optimal by the TAVC model, where the three fast robots serviced the high-value circuit and the three slow robots serviced the low-value circuit. Second, to demonstrate that this allocation improved the performance of the system to a level significantly above the performance level of a group where tasks were allocated randomly.

6 Results

We defined a convergence period of 15 hours based on the stability of the system performance. The current allocations was identified by looking at 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 labeled $\hat{\mu}_a$ and s_a in Table 1 show the mean and standard deviation of the time the system spent in each of the states while running the modified TAVC algorithm. The values are percentages of the total stable period. The columns labeled $\hat{\mu}_r$ and s_r give the mean and standard deviation of the time the system spent in each of the states for a set of 15 trials using a group of robots that randomly chose

between tasks. The column labeled T lists the combinatorial probability of

$f : s$	$\hat{\mu}_a$	s_a	$\hat{\mu}_r$	s_r	T	$\hat{\mu}_a - \hat{\mu}_r$	$\frac{\hat{\mu}_a - \hat{\mu}_r}{\hat{\mu}_r}$
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. State-Time Dist. for Heterogeneous Robots

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. This probability is given in Equation 3. It is worth noticing that the time distribution produced by random allocation is closely aligned with the theoretical estimate, though the differences are statistically significant.

$$T = \frac{100g!h!}{f!(g-f)!s!(h-s)!2^g2^h} \quad (3)$$

The difference between the state-time distribution produced by the modified TAVC algorithm and the distribution produced by random allocation is presented in the column labeled $\hat{\mu}_a - \hat{\mu}_r$. This difference is presented as a percentage of the mean times from the random distribution, $\hat{\mu}_r$, for each state in the last column, labeled $\frac{\hat{\mu}_a - \hat{\mu}_r}{\hat{\mu}_r}$. The difference between the distributions produced by the adaptive controllers and the random controllers, i.e., the last two columns, are also presented graphically by the two histograms in Figure 2.

Over these 15 experiments, the increase in time spent in state 0 : 3 is statistically significant. In the second histogram in Figure 2 the optimal state, 0 : 3, stands out as the state with the highest relative increase in time. This confirms that the group's set of Q-tables have converged to promote the state defined as optimal according to the TAVC model. The performance data also show that the performance of a group of robots controlled by the TAVC algorithm, 0.081 of target value per 10 seconds, is significantly higher than the

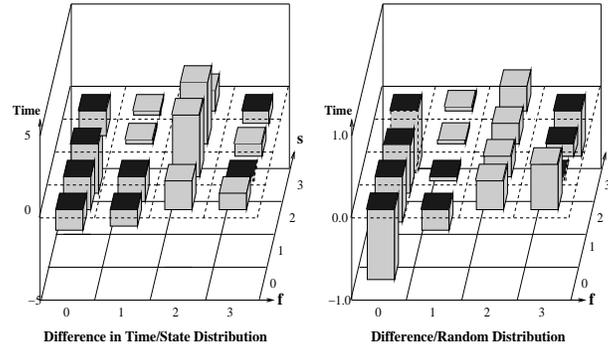


Fig. 2. Difference in Distributions

performance of random allocation, 0.078 of target value per 10 seconds. Together, the state-time distribution data and the performance data show that the modified TAVC algorithm improve the group’s performance by adopting a dedicated service structure that conforms to the predictions of the TAVC model. The learned Q-tables show that, as predicted, the fast robots, on average, have higher estimated utilities for high-value tasks and higher average differences between the estimated utilities of high- and low-value tasks.

7 Related Work

In addition to existing MRTA algorithms that rely on hand-coded information about performance levels, task utility and group dynamics [?, ?], there is a large body of work on modeling group dynamics.

Goldberg and Matarić [?] developed Augmented Markov Models, or transition probability matrices with additional temporal information, to learn statistical models of interaction in a space of abstract behaviors. Yan and Matarić [?] have used multi-level modeling of group behaviors from spatial data to describe both human and robot activity. Lerman *et al.* [?] have studied three different types of models of cooperation and interference in groups of robots: sensor-based simulations, microscopic numerical models, and macroscopic numerical models. Currently there are no models of the effects of group dynamics with the speed, generality, and predictive accuracy necessary to specify the effects of interaction on task processing times in task allocation problems when the aim is to constructing optimal schedules on the fly. Simulation-based models are in general too slow while macroscopic mathematical models make too many simplifying assumptions in order to be of predictive use.

In the absence of models, learning has been used to increase the applicability of both centralized and distributed MRTA and related scheduling algorithms. In the L-ALLIANCE work by Parker [?], each robot explicitly estimates its own performance and the performance of other robots on se-

lected tasks and uses these values to reallocate tasks by taking them over or acquiescing. The L-ALLIANCE algorithm uses local utility estimates to make local allocation decision, but needs hand-coded estimation procedures that reduce the general applicability of this algorithm. Brauer and Weiss [?] use a distributed RL mechanism for Multi-Machine Scheduling (MMS) where each machine estimates locally the optimal receiver of the material it has processed. This approach, like ours, uses local action selection and utility estimates. The MMS problem however does not contain the complex group dynamics of the transportation problem. Balch [?] studied performance-related reward functions for robot using Q-learning. Our results build on Balch's work and further explore the use of local performance based reward functions for optimizing group performance.

8 Conclusions

Our experiments show that when using a Boltzmann softmax function, the modified TAVC algorithm was sensitive to the different operating speeds of the robots and promoted the optimal state as defined by the TAVC model. As a percentage of the total system time, the increase in time the system spent in the optimal state is too small to create significant improvements in the performance of cooperative multi-robot systems. Also, a comparison with random allocation does not demonstrate the performance of the modified TAVC algorithm relative to existing MRTA algorithms. For such purposes a basic greedy allocation algorithm would have provided a better basis for comparison. Finally, we do not know how the TAVC algorithm would scale up to problems that are too complex for each robot to evaluate every circuit, but it seems feasible that satisfactory global solutions could emerge through migration based on local interactions. In spite of these limitations, our results are important as they show that the theoretical probability of using a softmax action selection function to make distributed MRTA sensitive to different robot abilities does translate into observable effects in realistic multi-robot simulations.

Acknowledgments

This work is 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.

References

1. T. R. Balch. Reward and Diversity in Multirobot Foraging. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence (IJCAI'99)*

- Workshop: Learning About, From and With other Agents*, Stockholm, Sweden, July 31 - August 6 1999.
2. Tucker 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 (Agents'99)*, pages 92–99, Seattle, Washington, May 1 - 5 1999. ACM Press.
 3. Wilfried Brauer and Gerhard Weiß. Multi-machine scheduling - a multi-agent learning approach. In *Proceedings of the 3rd International Conference on Multi-Agent Systems (ICMAS'98)*, pages 42–48, Paris, France, July 4 -7 1998. IEEE Press.
 4. Ivan D. Chase, 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*, 36:1265–1274, 1988.
 5. Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to algorithms*. MIT Press, Cambridge, Massachusetts, second edition, 2001.
 6. Torbjørn S. Dahl, Maja J Matarić, and Gaurav S. Sukhatme. Scheduling with group dynamics: A multi-robot task allocation algorithm based on vacancy chains. Technical Report CRES-002-07, Center for Robotics and Embedded Systems, University of Southern California, Los angeles, CA, 2002.
 7. Torbjørn S. Dahl, Maja J. Matarić, and Gaurav S. Sukhatme. Multi-robot task-allocation through vacancy chains. In *Proceedings of the 2003 IEEE International Conference on Robotics and Automation (ICRA'03)*, pages 2293–2298, Taipei, Taiwan, September 9 - 14 2003. IEEE Press.
 8. Brian P. Gerkey and Maja J Matarić. Sold!: Auction methods for multi-robot coordination. *IEEE Transactions on Robotics and Automation*, 18(5):758–768, October 2002.
 9. Dani Goldberg and Maja J Matarić. Learning multiple models for reward maximization. In Pat Langley, editor, *Proceedings of the 17th International Conference on Machine Learning (ICML'00)*, pages 319–326, Stanford, California, June 29 - July 2 2000. Morgan Kaufmann.
 10. Kristina Lerman, Asram. Galstyan, Alcherio Martinoli, and Auke J. Ijspeert. A macroscopic analytical model of collaboration in distributed robotic systems. *Artificial Life*, 7(4):375–393, 2001.
 11. Maja J. Matarić. 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):323–336, 1997.
 12. Lynne E. Parker. L-ALLIANCE: Task-Oriented Multi-Robot Learning in Behaviour-Based Systems. *Advanced Robotics, Special Issue on Selected Papers from IROS'96*, 11(4):305–322, 1997.
 13. Richard S. Sutton and Andrew G. Barto. *Reinforcement learning: an introduction*. MIT Press, Cambridge, Massachusetts, 1998.
 14. Helen Yan and Maja J Matarić. General spatial features for analysis of multi-robot and human activities from raw position data. In *Proceedings of the 2002 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'02)*, pages 2770–2775, Lausanne, Switzerland, September 30 - October 4 2002. IEEE Press.