

An Infrastructure for Large-Scale Human-Robot Interaction

Ashley D. Tews, Gaurav S. Sukhatme, Maja J. Matarić

Abstract— While there is much research in human-robot interaction on the personal, or one-to-one level, there is little on the large scale, or many-to-many. We propose an interaction infrastructure for many human operators to utilize a large number of robots. The infrastructure operates in two phases allowing the tasks and services offered by many robots to be exposed to many users, as well as the personal interaction between a user and a robot. We assume that each robot has some degree of autonomy and that the humans assume a supervisory role. The approach is analyzed in detail on a single robot system with many users to establish its low level operation and then on a large-scale experiment involving several robot systems offering many tasks and services to many users. In the latter experiment, the infrastructure managed the interaction between hundreds of users and several robot and embedded systems.

Index Terms—Humans, Interaction, Robots, Scalability, Systems.

I. INTRODUCTION

The interaction between a robot and a human can occur at many levels. There is much research on one-to-one (single robot and single human) interaction but the approaches typically do not scale to large groups. This is mainly due to the communication paths required for interaction between interacting entities being dedicated to the type of interaction being carried out. Also, for tightly coupled interaction, the communication bandwidth required is typically high.

In this paper, we investigate *many-to-many interaction: what is required to facilitate it*. We intend our method to scale to the order of hundreds of robots and users. It is anticipated that multitudes of robots will soon be used to achieve goals beyond their current small-scale limits. These goals may be building large structures on Earth or more compellingly in orbit around Earth, or on other planets in preparation for human habitation. They may also be used in a human-coordinated de-mining effort, where a few human supervisors can coordinate hundreds of robots. Finally, they may carry out military operations to reduce the danger to humans. It is assumed that all of the above applications would use robots that have some degree of autonomy, such that they could be tasked by human supervisors rather than intimately controlled. The latter approach would require too much effort

for the interaction to be productive. An interaction mechanism for these scenarios must also allow for heterogeneous robots, since many different types of robots with different capabilities would likely be needed.

Two issues to be considered for scaling interaction to larger groups and for more general tasks are: 1) how to expose and connect a large number of entities for potential interaction to each other, and 2) how to allow for the possible types of personal interaction. For hundreds of robots and hundreds of people, the intricate and tightly-coupled approach typical of one-to-one interaction is infeasible due to the logistics involved in interconnecting them and their typically large, dedicated communication bandwidth requirements. A reduced level of coupling and communication is required for the mechanism to be scalable. Additionally, when personal interaction between a human and robot *is* undertaken, the same mechanism should allow for the maximum amount of information exchange between them. Therefore, for an interaction approach to scale, it needs to provide as much functionality as possible for individual (one-to-one) human-robot interaction, as well as providing access and exposure of the robot's services to other humans (many-to-many interaction).

For large-scale interactions, it is conceivable that a single person can supervise or control multiple robots or robot groups concurrently in a one-to-many interaction. The interaction mechanism also is required to accommodate this.

We propose a two-phase approach aligned with the goals of allowing many-to-many, one-to-many and one-to-one interaction. This is accomplished by maximizing access between multitudes of robots and humans by minimizing the implementation and communication overhead required, and maximizing the interaction bandwidth between small numbers of robots and humans. The robots in such an approach are assumed to have some level of autonomy and, as such, allow humans to take on a supervisory role in the interaction. This allows for a lower bandwidth of communication for one-to-one and one-to-many type interaction [1].

This paper is organized as follows. Section 2 provides an insight into research relating to the fundamentals of the presented infrastructure presented. Section 3 discusses the criteria for designing a large-scale interaction mechanism as the foundation for the infrastructure. Section 4 introduces the

infrastructure and discusses its two phases of operation. Section 5 provides details of an experiment that analyzes the infrastructure operation in terms of how well it manages interaction requests and the effect and implications of buffering requests. A second experiment investigates the scalability of the infrastructure with several systems and many users. Section 6 elucidates the philosophy behind the infrastructure design aligned with the criteria outlined in Section 3. Section 7 describes the next stage of development for the infrastructure and concludes the paper.

II. RELATED WORK

The discussion of scalability in human-robot interaction can be broadened to include other systems since in a large-scale application, not all of the entities will necessarily be humans or robots. A system can be considered to consist of either physical (such as robots or embedded devices) or virtual entities (such as the services provided by computer programs). As such, human-system interaction covers several large areas of research. In [2], a taxonomy is presented categorizing human and intelligent system research. *An intelligent system* is a system that allows user interaction (e.g., tele-robotics[1]). The categories in the taxonomy include the application of the research, research approach, system autonomy, interaction distance, and interaction media, with each category divided into subcategories. Numerous references are grouped under these categories, illustrating the vastness of the field. Another key axis that can be added to the classifications is the scalability of the approaches in terms of users and systems.

TABLE 1
USER/SYSTEM SCALABILITY CLASSIFICATION

	Users	Systems
1	One	One
2	One	Multiple
3	Multiple	One
4	Multiple	Multiple

In this section of the paper, we classify human-system architectures in terms of their scalability, *assuming a system consists of a set of nodes that provide a set of services to users*. These classifications are shown in Table 1 and are discussed in the following sections, separating them into single (non-shaded table entries) and multi-user architectures (shaded table entries).

Single-user System Architectures

There are many examples relating to the entries in the table above. The issues concerning scalability are highlighted when examining these examples. An important issue that affects the scalability of an approach is how *tightly coupled* the users are with the system. This concerns both access to system resources and levels of control over system services. Single-user systems can allow the user full access to system resources

and control of the services. Scaling these to more users can be difficult due to these attributes (such as communication bandwidth) having to be shared. Tele-operation of a single system node is an example where this is highlighted [1][3].

Multi-robot research also provides examples where multiple robots can be controlled individually or as a group by a single user [4][5][6][7]. In each of these systems, the nodes have a certain degree of autonomy and the user has either high-level group control, or low-level individual control. This indicates that the higher the number of users an interaction architecture has to deal with, the looser the coupling between users and systems.

The concept of a single user for multiple systems indicates a hierarchical command structure where there is a top-level entity that controls the systems below its level.

Multi-user System Architectures

Multi-user and multi-service systems are forced to address issues of scalability, system resource allocation, and control. A key concern is how to control the information and communication pathways for multiple users. Common approaches involve using multi-agent architectures [8][9][10][11] and the Internet.

Multi-agent architectures offer services to users with suitable access to the systems they serve. Examples include interactive workspaces [12][13] which consist of dedicated environments with embedded devices providing services for users. The users generally must be in the environment to access the services.

Multi-user/multi-robot architectures are demonstrated by [14][15]. These systems use the Internet as a communication medium for commanding robots from remote locations.

Many of the examples presented in this section are scalable only in their local operating environment or application domain. The infrastructure presented in Section IV is intended to operate beyond these limits, as well as contend with the diversity of coupling and resource-sharing required by individual systems.

III. THE REQUIREMENTS FOR HIGH SCALABILITY

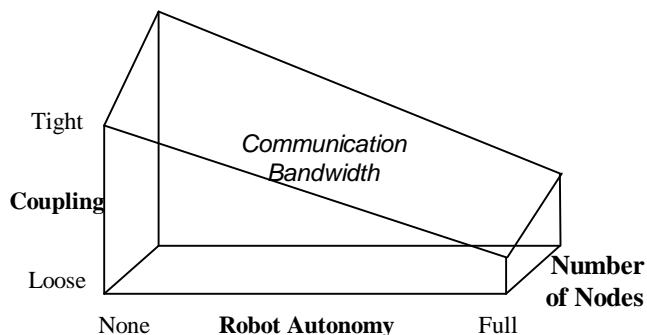


Figure 1. The relationship between the communication bandwidth and the number of interacting entities, robot autonomy and coupling.

In its simplest form, human-robot interaction consists of

the human, the robot, and the interaction mechanism. The robot and human interact through a communication medium, whether it has an explicit form such as spoken language or message passing, or implicit such as a robot observing a human's deliberate actions. To extend the numbers of entities capable of interacting, this communication bandwidth required between them must be kept low. There are several issues that affect the amount of communication required. Figure 1 shows an abstract view of those considered the most critical for scalability.

The level of autonomy affects the communication required between a human supervisor and a robot. As robots become more autonomous, they require less intervention or guidance by a supervisor, and therefore less communication [1]. Consider the difference in user control required between a mobile robot that is capable of GPS-based navigation with path planning and obstacle avoidance, and the same robot that is totally controlled by a remote user. Navigating to a destination requires more user input for the remotely-controlled robot and therefore, more communication bandwidth. The autonomous robot would require only the destination location and would traverse its own path.

Interaction can exist on many levels of coupling between a robot and a human. For example, a robot observing a human with the purpose of imitation [16] needs to place much effort into sensing and interpreting what has been sensed. The human, in turn, is aware that the robot is imitating, and acts accordingly. Consequently, during the interactive period, both the robot and human are tightly coupled and much of their communication bandwidth is dedicated to this task. The communication in this case is implicit, which requires the user to perform actions understandable and distinguishable to the robot's sensing and interpreting systems. A more loosely coupled example is if a robot requires a simple signal to start its task, such as a user pressing a key on a keyboard. The robot's dedication to the human and the communication bandwidth in this case are much less.

The number of nodes also affects the communication bandwidth required for human-robot interaction. If it is assumed that each of the preceding examples requires a certain amount of communication bandwidth between the user and the robot for one-to-one interaction, then as their numbers increase, so does the bandwidth required. However, homogeneous communication between a human and several robots may allow for on-to-many synchronous interaction. For example, if several robots are observing a human for the purposes of mimicking in a teacher-student type of relationship, then they are all using the same information implicitly communicated from the human.

Apart from the communication bandwidth issues, the capabilities of the robots also affect the design of a large-scale interaction mechanism. In an application such as building a structure, the robots may not have the same design or capabilities. Therefore, the interaction between a human and

two different robots may require different methods. The interaction mechanism should allow and manage this diversity.

Each of the factors discussed above influence the design of a large-scale interaction mechanism in the following manner:

- it should allow many humans to interact with many robots (many-to-many),
- it should allow for a single person to interact with many robots (one-to-many),
- it should allow for personal interaction (one-to-one),
- it should allow for the different levels of coupling,
- it should maximize the communication bandwidth between an interacting user and robot, and
- it should manage the heterogeneity of capabilities of different robots or robot groups.

To carry this out for potentially hundreds of users and robots also means the many-to-many aspect of the mechanism must be simple enough not to burden the communication paths or processing required to allow interaction. It also must be robust in order to handle contingencies such as robot system failure. The design of our interaction infrastructure is influenced by these requirements and is discussed next.

IV. THE INFRASTRUCTURE

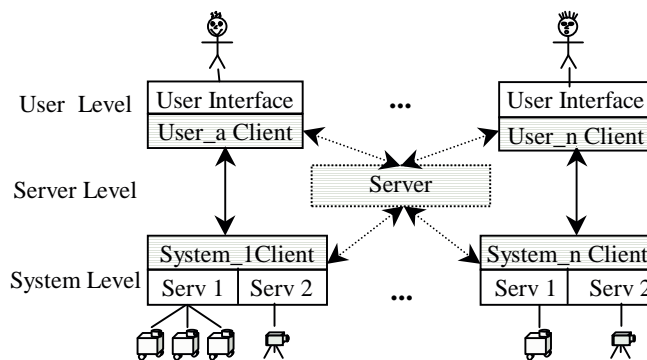


Figure 2. The infrastructure for connecting multiple users to multiple system services.

A. Overview

The infrastructure is based on a server/client architecture as shown in Figure 2. Each system is typically a robot, robot group, or embedded system (e.g., a security camera network). Each entity (system or user) has a client that interfaces with the server and manages the low-level signal and message passing. It is assumed that the server's contact information is known to all systems and users wishing to use the infrastructure. The server is the centralized hub of operations that manages the allocation of systems' services and users requests. The data flow throughout the infrastructure is by UDP communication.

The server contains a list of all capabilities - termed 'services' that each connected system provides. The users select whichever services they require and are connected to the corresponding systems. The appropriate level (coupling)

of interaction is then established between the user and the selected service, and is system- and service-dependent. The communication between the user and the service is decoupled from the rest of the infrastructure to maximize bandwidth for interaction and not inhibit the server from carrying out its higher level operations. The bandwidth may be reduced if the user selects multiple services. The operation of the infrastructure is discussed next.

B. First Phase Operation

The primary goal of the first phase is to provide a scalable method that enables many users to access the services provided by many systems. Since it is intended to scale to hundreds of users and systems, it must be simple and fast. To achieve this, a minimalist approach is taken. Each system client contains a list of service descriptions its system is capable of. These descriptions have little required format but are assumed to be informative enough for the users to understand their semantics. They also contain information classifying the service as a command or data type and whether to set up a task queue. A *data* service provides information about its progress whereas a *command* service requires a user to input information for its operation as well as providing progress information. Task queuing can be applied at the service or system level. System level queuing is when a system is capable of several services, but can only carry out one at a time. Request queuing increases the functionality of the infrastructure and allows users access to services they could otherwise be unaware of if the system or service was removed from the server's global service list after it was selected. For example, adding a task queue to a room-cleaning robot's services allows it to proceed to the next task as soon as it has finished cleaning the current room instead of waiting to be tasked *after* it has finished.

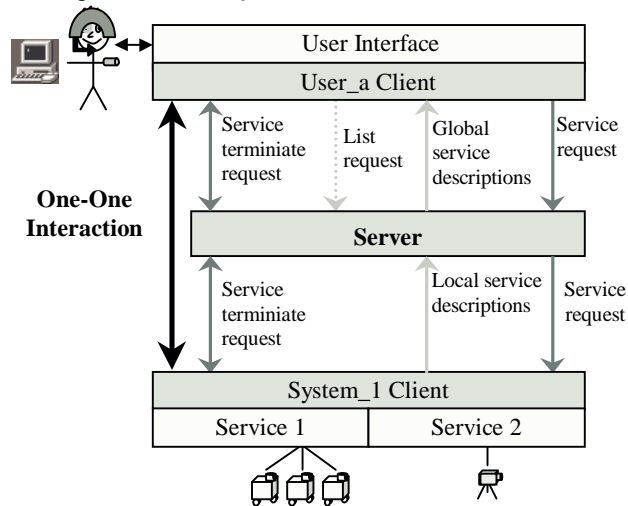


Figure 3. Signal and message paths are shown in Figure 3. Through each system's client, the service descriptions are sent

to the server. The system client then waits for requests.

The server collects the service descriptions along with their source's contact information. It generates appropriate service or system task request queues as required. Each system is added to the system database if it is new, or checked to see if it already exists. If it exists, a check is carried out to compare the services to those already known. If they are new, they are added to the system's service list. Once a system is in the database, it is routinely contacted to ensure it is still on line. If the server receives no response after several attempts, it removes the system from the database. This way, any delinquent systems can be removed to avoid cluttering the database or providing unfulfillable services to users. The only time a service is removed from a system listed on the database is when the system client sends a service termination signal. This effectively removes the service until it is resent by the system client and processed by the server as a new service.

Users connect to the server via their clients and are sent the global service list. They then have the options of service selection, refreshing the list, linking services, or exiting as shown in Figure 4.

The service and system information stored in the server is dynamic since users and systems go on and off line at various times. Consequently, the global service list is a 'snapshot' in time of the information on the server. Therefore, one of the options is to request the latest list.

Selecting the service linking option requires the user to select a semantically compatible data and command service. At this stage in development, it is assumed that the users know which services can be linked with others. For the experiments in the following sections, service linking is not utilized however, it is described in more detail in Appendix A.

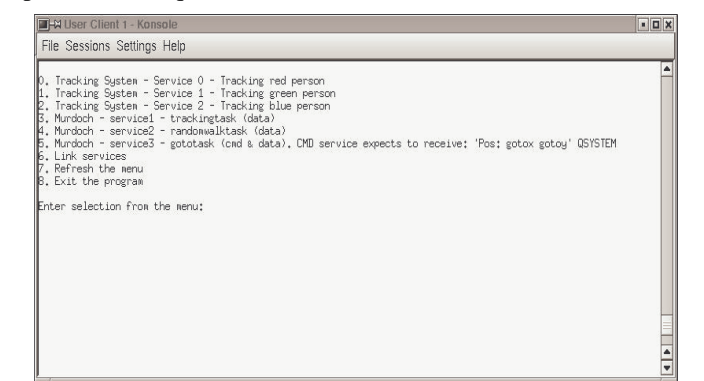


Figure 4. Example user menu.

If a service has been selected, the second phase of the infrastructure operation commences, whereby the server couples the user to the service for more direct interaction. For each service a user selects, a window is created to manage the interaction and display the service information. The user may select many services at any time, which enables them to interact on a one-to-many scale.

For each service a user selects, a window is created to manage the interaction and display the service information. The user may select many services at any time, which enables them to interact on a one-to-many scale.

C. Second Phase Operation

The server receives a user's service request in the form of the service description and the contact information for the system providing the service. These factors are considered to be unique for any service to avoid errors in service selection. An alternative is to index the system and service information on the server, but since the systems can be dynamically added and removed, their index will not be unique over time and their slot may be given to a new system. So any user request from an outdated list could result in them being connected to an un-requested service.

If the service or its system has a task queue, the request is added to the end of the queue. If the queue is full, the user is informed. If there are no tasks in the queue or there is no task queue for this service, the server sends the user's contact information and service description to the appropriate system. In the case of service linking, this occurs for both the command and data service in the same manner. The applicable system client receives the request and commences the service, communicating directly with the user. This alleviates the communication load on the server and allows for a greater bandwidth of communication between the user and the system. If the service requires a user's input (such as a command service), it is annotated in the communication. The system sends task-relevant information to the user, which is displayed in the designated service window. Eventually, this one-to-one interaction will be managed on system-dependent, user-modifiable [17] or semantically-selected user interfaces [11], but these concepts are non-trivial and are currently beyond the scope of this research. The user and system are now linked until either terminates the service. The user can terminate the service by either closing down the service's window, or logging out, in which case all connected services are notified. If a system terminates a service, it disconnects that service from the user and notifies them.

By separating the operation of the infrastructure into two phases, the server is dedicated to managing the many systems and users connected to it during the first phase. The second phase allows for more dedicated communication between the user and the selected services. The following section provides an in-depth analysis of how the infrastructure works by micro and macro-inspection of response times for requests and scalability issues.

V. EXPERIMENTAL RESULTS

To validate the design of the infrastructure, we performed two experiments, described in this section. The first presents a micro-view of the infrastructure response times for a robot system with and without user queuing. The second presents a macro-view of the scalability aspects of the infrastructure using several systems and hundreds of users over hours of operation.

All experiments are carried out in the Stage¹ simulator using Player² devices [18]. Stage has been demonstrated on many projects to provide a valid representation of the real operation of the Pioneer robots. Player provides a suite of software clients that represent the handlers for the devices on the robots such as sonar, SICK laser range-finders and color cameras. The combination of Player and Stage allows programs developed and tested in Stage to transfer to the real robots.

Stage is used in these experiments to simulate robots operating in a realistically scaled environment of the second floor of our building. This provides a realistic office space environment with corridors and rooms for the robots to carry out their tasks, and allows the systems' parameter settings to be used for testing in the real world at a later stage.

A. Experiment 1 – Micro View of the Infrastructure

TABLE 2
PARAMETERS FOR EXPERIMENT 1

Number of users generated	25
User generation frequency	15 seconds
Experiment running time	500 seconds
Task running time	20 seconds
Minimum possible requests	12 or 13
Maximum possible requests	25
Queue size	0, 3, 5, or 7

This experiment shows typical user and system response times in a many-to-one interaction scenario. It also demonstrates the benefits of using task queuing. A single system is used with a single *random-wander* service representing a robot that simply randomly wanders around the environment and provides information about its current location. The task continues for 20 seconds before it terminates and becomes available for the next request. The system, users, and server are run on separate machines.

The parameters for the experiment are shown in Table 2. The users are automated and generated upon system start-up and every 15 seconds thereafter, until 25 have been generated. Once a user is on line, they immediately select the service if it is available. If it is not, they do not make another attempt. The number of tasks fully that can be completed range from 12 to 25, depending on how much task buffering is applied. The minimum possible number of requests that can be fulfilled is determined by the user generation rate versus the task running time.

¹ Stage was developed jointly at the USC Robotics Research Labs and HRL Labs and is freely available under the GNU General Public License from {<http://playerstage.sourceforge.net>}.

² Player was developed jointly at the USC Robotics Research Labs and HRL Labs and is freely available under the GNU General Public License from {<http://playerstage.sourceforge.net>}.

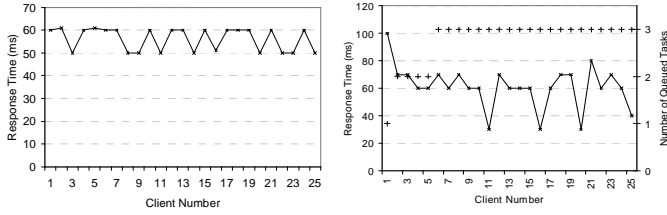


Figure 5. Server response times for no queue (left) and a 3 buffered queue (right)

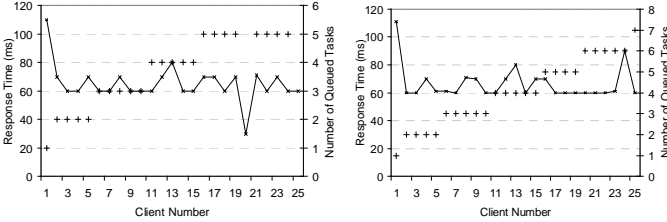


Figure 6. Server response times for a 5 (left) and 7 (right) buffered queue

Figure 5 and Figure 6 show the server’s response times (solid lines) for processing a request and the number of tasks in the task queue (crosses) if used. A request is considered processed by either sending it to the system client or placing it in a task queue.

For all buffered approaches, the response time is generally between 60 and 80 ms, regardless of the buffer size. The high exceptions for the first user client are due to the time taken to add the user to the buffer and then remove them again since their request is executed immediately. The low points (typically 30 ms) are due to the user’s request being denied since the queue is full, as indicated by the crosses denoting the number of tasks in the queue at the time the user’s request was placed. The approach with no buffering is indifferent to whether the user’s requests are fulfilled or denied and the response times are between 50 and 60 ms. The request statistics are shown in Table 3.

TABLE 3
USERS’ REQUESTS FULFILLED

Buffer Size	Requests Fulfilled
0	13 (52%)
3	21 (84%)
5	23 (92%)
7	25 (100%)

These results demonstrate that a buffer size of seven allows all users’ requests to be fulfilled within the constraints of the experiment parameters. It also demonstrates the load on the server for the extra processing involved with adding task queues to a system is negligible.

The user’s response times reflect the time from when they initiate a request to when a response is received from the system. For the 0-buffered system, this consistently averages 75 ms. The response times and task queue sizes for the buffered approaches are more variable. They are shown in Figure 7 as solid and broken lines, respectively. Unsuccessful attempts are recorded as a response time of 0 to identify the corresponding users.

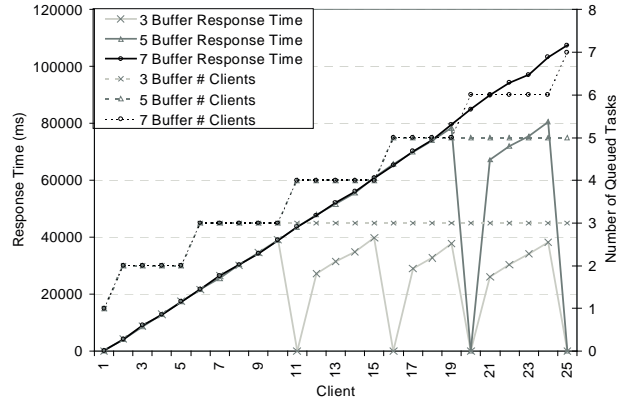


Figure 7. User client service request response time (solid lines corresponding to the left y axis) with task buffering (broken lines corresponding to the right y axis)

As can be seen in Figure 7, the response times varied depending on the size of the task queue. In all approaches, the first 10 users had their task requests fulfilled. After that, the times increase as more users get placed on the buffer and each has to wait their turn. Once each buffer is full, the subsequently generated user has their request denied. Since the task length is 20 seconds and the first user on the queue is currently connected to the task, each approach has a maximum theoretical response time for the users of $(\text{buffer_size} - 1) \times \text{task_length}$. This gives 40, 80 and 120 seconds for the three, five and seven-buffer approach, respectively. These upper limits are apparent on the graph when the buffers are full.

B. Experiment 2 – Macro View of the Infrastructure

While the previous experiment provides an in-depth view of the infrastructure’s timing and service buffering, the experiment in this section investigates scalability. This is demonstrated with five systems offering many services to many users.

Each system provides varying numbers and types of *data* service. The service types are classified as discrete or continuous. Discrete services provide the status of a sensor while continuous services provide streaming updates of information about their task status. There are three continuous service systems and two discrete systems. The continuous service systems consist of Pioneer robots that carry out tasks jointly or individually. Each system is briefly described next.

1) Murdoch – continuous services

There are three robots under control of a task-allocation system known as Murdoch [19][20]. Murdoch uses an auction-based method for assigning incoming tasks to any robot that fulfils the task requirements; auction communication is performed using anonymous publish-subscribe messaging. For this experiment, Murdoch is used as a central contact point to allocate tasks to robots, where each robot is capable of carrying out three tasks: random walk, target tracking, and delivery. Random walking requires a robot to wander randomly and provides location information

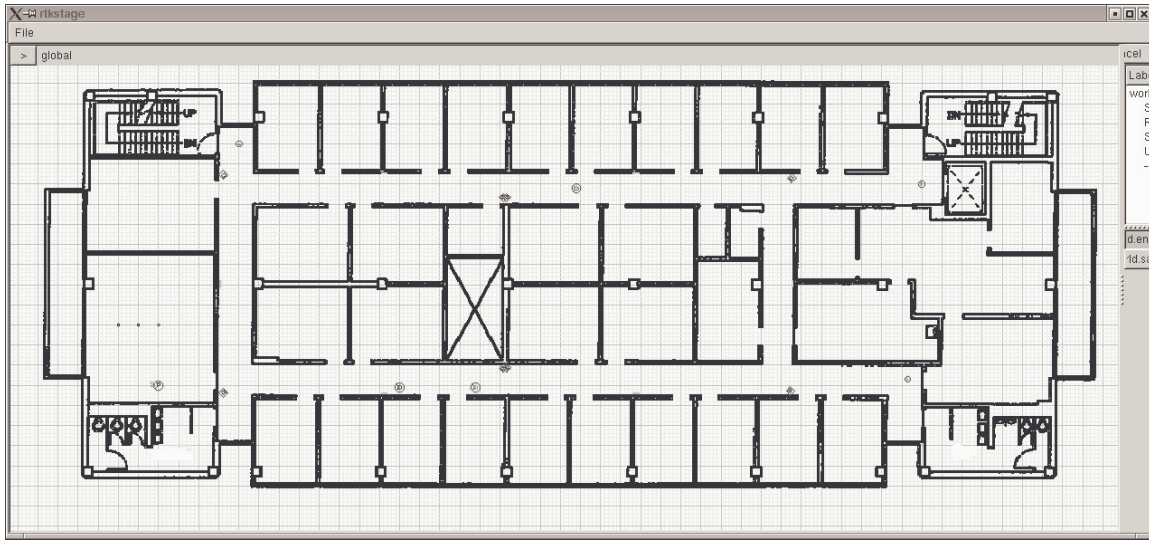


Figure 8. The simulated environment showing robots in their initial positions

to the user. For target tracking, the robot acquires and follows a coloured mobile object and provides its current global location. The delivery service requires the user to input the destination coordinates for the robot to navigate to. Once a robot is tasked, it does not carry out another task until the current one is terminated. Therefore, upon service selection, the relationship between humans and robots is one-to-one.

2) Tracking – continuous services

The second robot-based system uses a group of robots that provide intuitive tracking of several targets [21]. This system positions robots dynamically to provide the best coverage for tracking several targets concurrently. The displacement of each robot is based on several factors, including target and robot density in pre-defined designated areas of the environment. The system has several configurations depending on the application requirements including using the robots as stationary sensors. To prevent congestion in the corridors of the simulated office floor caused by too many mobile entities, this is the configuration chosen for the experiment. The robots are positioned to provide nearly complete coverage of the corridors in the simulated environment. Eight robots are required for this operation. Each is equipped with a color camera and reports the global positions of simulated humans with red, green, or blue markers on them. This system provides three services concurrently to a maximum of ten users.

3) Room Cleaning/Trashing – continuous service

The third system involves a robot simulating a room cleaning or trashing operation. Cleaning involves moving all coloured objects to a certain area in a room while trashing involves dispersing these objects around the room. The cleaning task simulates a system that is on call for a user to select when a room needs to be cleaned. This can be applied to a vacuuming robot or one that can pick up trash. The robot is in a designated room and provides one-to-one interaction relationship.

4) Room Lights and Room Occupied - discrete services

Many of the rooms in the simulated building floor have two sensors that determine whether the room is occupied and the status of its lights. The lights are automated to turn on and off at random times and the services provide the light status of each room. The room-occupied system provides services that ascertain whether a room is currently occupied by a simulated human. Both of these systems provide discrete services for most rooms on the floor and add many services to the global list.

a) Setup

TABLE 4
EXPERIMENT ENTITIES AND SERVICES

System	No. of Nodes	Services per node	Total System Services	Total Users per System
Murdoch	3	3	9	3 (+6)
Tracking	8	3	3	10
Cleanup	1	2	2	1
Light Condition	33	1	33	N/a
Room Occupied	33	1	33	N/a
Total	78		80	

The environment for this experiment is the second floor of the USC Computer Science building, simulated in Stage along with the robots and devices as shown in Figure 8. Apart from the robots providing services, there are also three simulated humans (implemented as robots) that wander randomly through rooms and corridors. Each human has a differently colored marker for the tracking systems to detect. The numbers of entities and services is shown in Table 4.

There are 80 services in total to be used for this experiment. Note that the tracking system has eight nodes capable of three services each, but only provides three services to the experiment. This is due to each robot operating in parallel with its peers, so they form a global, rather than a

local sensor network.

The numbers of users/system indicate how many users can be using the system at the same time. Murdoch has one for each robot and a task queue, so even though a service may be unavailable at the time a user requests it, it can be added to the queue. The discrete services do not have a limit to the number of users they can manage, since the requests are handled with a single response. All systems except Cleanup remain on line when their services are selected. Once a Cleanup service (room cleaning or trashing) is selected, it removes itself from the server's database. The other systems inform users if they cannot process their requests due to operating at their maximum user limit.

The users consist of real and simulated people. The simulated users are generated randomly and automatically on several different computers and make random service selections up to an arbitrarily chosen maximum of two. They select any of the options on the menu (except service linking) and manage the resulting interaction as a real user would. At any time, there can be a maximum of 32 users on line.

Since many of the services are discrete, they will have a greater probability of being selected when all systems are on line. Therefore, they are started after a significant period of time during which, only the robot systems are connected to the infrastructure.

The systems, users, and server are distributed across several computers connected to the local area network operating at 10Mbs. The experiment is conducted over four continuous hours of operation.

b) Results

TABLE 5
STATISTICS FOR EXPERIMENT 2

Total users	985
Maximum users on line	32
Total requested services	916
Total fulfilled services	730
Total requested service lists	2022

The overall statistics of the experiment are recorded in Table 5. From this table, it can be seen that the user rate is high (over 200/hr) for the period of time the experiment ran. Not all users requested services and of those that did, most were fulfilled (approximately 80%). Requests were rejected for several reasons depending on the system:

- Murdoch* - the request queue lengths are full,
- Tracking* - the maximum number of user requests are being managed,
- Cleanup* - a request is made from an outdated user service menu when the clean/trash service is requested.

Other rejections occur if a user selected a service and logs out before notification is received from the system that the request is being processed.

Whenever a user connects to the infrastructure, they are

given the most recent service list; hence the high number of service list refresh requests.

System Usage

TABLE 6
SYSTEM USAGE STATISTICS

System	No. of Services	Running Time (%)	No. Requests	Request Ratio (%)
Murdoch	9	100	151	16.4
Tracking	3	100	142	15.5
Cleaner	2	77.9	24	2.6
Room Occupied	33	65.8	347	37.9
Light Status	33	53.3	253	27.6

The system usage rate is shown in Table 6. The continuous service systems were operating for the longest period of time but had fewer requests due to the number of discrete services providing most of the services, and most continuous services capable of running for arbitrary lengths of time. Nevertheless, a significant number of requests were made to these systems.

During the experiment, only the Murdoch systems went off line. This occurred three times due to Murdoch timing out a task that has been running for a pre-set time (approximately 17 minutes). When the systems went off line, the server detected it and they were removed from the global service list until they were restarted.

User Clients

Figure 9 provides an example of the typical response times for the services to respond to users' requests. Generally, the times are low, except for Murdoch, due to its task queuing. There were no time limits imposed for service usage, so a user with a task in the queue could wait for an arbitrary amount of time, as manifested by the long wait times. The other spuriously high wait times are due to the time taken by the host machine to process the task while carrying out internal housekeeping.

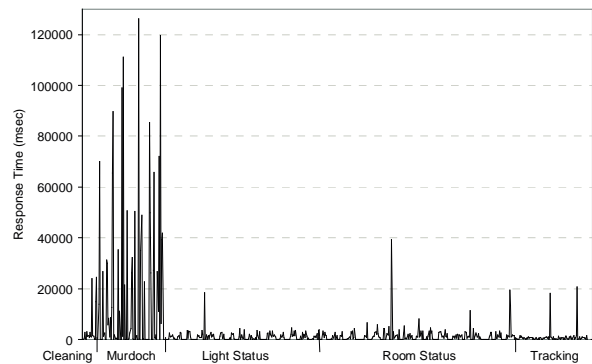


Figure 9. Example response times for users requesting services

Server Statistics

Apart from the demonstrated number of systems, services, and users the infrastructure managed, another important

aspect of scalability is the amount of information flow for operations. The largest transaction for the server is when a user requests the global service list. Figure 10 shows the throughput of information during these requests. The requests are shown in gray (bytes received) with the service list size (bytes sent) in black. The steps in the service list size indicate that a new system and its services have been added to the server's database. The small fluctuations indicate that some of the services are unavailable. This occurred for the Cleanup service and when the Murdoch robots went off line.

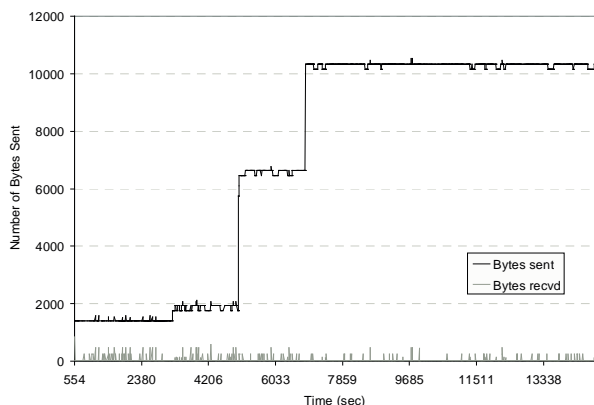


Figure 10. The number of bytes sent and received by the server in response to a service list request

The other requests sent to the server are service requests and service termination messages. A service termination message can be sent from either the system or a user connected to the service. When sent from the user, it means the user has finished with the service. When sent from the system, it means the service is to be removed from the global service list.

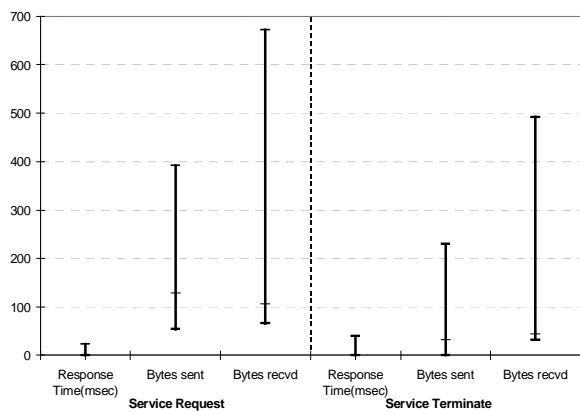


Figure 11. The statistics for service request and terminate shown with maximum, average and minimum marks

Figure 11 shows the response times and data statistics for managing these requests. 'Bytes received' is the size of the message received by the server for the request and 'Bytes sent' the resultant message sent from the server if required. The average processing time for both types of message is less than a millisecond and the longest time is less than 40 ms.

As can be seen from the figure, the average amount of information managed in order to process (send and receive) a

service request and service terminate messages is typically 100 bytes or less. The theoretical maximum number of these requests that can be managed is limited by the UDP buffer size, which is set at 32kb. This would allow approximately 40 simultaneous requests of average size to be received.

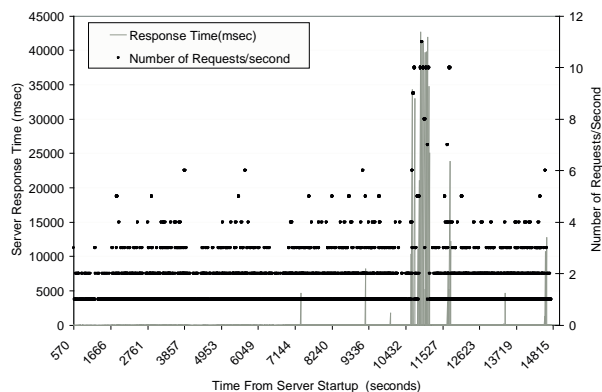


Figure 12. Number of responses/second versus response time for the server

Figure 12 illustrates how the rate of requests affects the server's response times for all messages it received. When the request rate is consistently above nine requests/second for a period of time, the server's response time rises uncharacteristically. However, the server was still able to process all requests since this burst rate did not persist. Further testing of this limit will determine both the maximum burst rate of requests the server can handle and the prolonged consistent rate. Upon analyzing the requests during the high burst time, most were found to be list refreshes, which require the most bandwidth.

VI. DISCUSSION

There are two important issues to be considered when developing a large-scale human-robot interaction mechanism. The first is how to allow possibly hundreds of humans to interact with hundreds of robots: the many-to-many problem. The second issue is how to allow personal interaction between a single human and single robot system with the same mechanism that addresses the first issue. There are many other issues to address for one-to-one interaction, including the modality of the interface and the coupling required for the interaction. These issues are important and constitute distinct areas of research. They are obviously still relevant to the many-to-many interaction discussed in this paper, but are considered to be lower level issues to be integrated and addressed once a large-scale interaction mechanism has been developed.

Our approach to addressing the primary issues for large-scale interaction consists of two phases aligned with these issues. The first phase depicts the many-to-many interaction, whereby every user is indirectly in contact with every system through a centralized server. This is beneficial for large-scale robotics applications such as building structures in human-adverse environments, and multi-robot military missions such

as terrain-mapping and scouting. In these domains, there can be many robots involved that have various abilities and are capable of carrying out tasks with some degree of autonomy. Humans have a supervisory role and can monitor and interact as required. When more direct contact with a robot is required, the second phase of our approach is activated. The user selects the relevant service and is put in direct contact with the system. At this point, the interaction is system-dependent and assumes a one-to-one role.

To achieve large-scale interaction, it has to be assumed that the participating robot systems can have unknown abilities and require different levels of interaction. Therefore, the infrastructure must be general enough to allow the heterogeneity of the systems and robots using it. In the experiments described in this paper, there are several robot and embedded device systems used to demonstrate this feature. The services are either continuous, in the sense that they require constant interaction, or discrete in providing the instantaneous status of a situation. For most of the demonstrated services, the interaction is mainly in the form of human supervision and monitoring, with the exception of the delivery task which requires the user to provide a delivery destination. These are simple examples of one-to-one and one-to-many interactions, but they allow the important aspects of the infrastructure's scalability to be demonstrated. Further development of the infrastructure will allow for more complicated personal interaction to be undertaken.

In Section 3, the requirements for scalability are outlined. They state that an interaction mechanism should be able to: manage heterogeneous robot capabilities, allow for different levels of user-robot coupling, maximize bandwidth between interacting entities, and allow many users to interact with many robots. To achieve this, the approach should be simple, to reduce response times, and robust, to handle contingencies. Our infrastructure has achieved many of these design goals as explained above and supported by the initial experiments presented in this document. There are many enhancements that can be added to the current version of the infrastructure. The major ones are described in the next section.

VII. FUTURE WORK

The major issues to be addressed with the current infrastructure are to add redundancies to the server, develop an information semantics checking mechanism for service linking, and perform further testing of its scalability and bandwidth limits.

The server is currently a centralized mechanism. If it fails, no new requests will be processed. Any existing one-to-one interactions will continue, however. To address this possible shortcoming, an investigation will be undertaken to distribute the server. This can be carried out by having multiple versions of the server running in parallel and, if one fails, another will become the primary server and continue operation. Another approach in the theme of enhancing

scalability is to adaptively distribute the server's operation over several computers, so each handles a part of the server's load. This way, any failure should only affect only part of the operation.

There is currently little format imposed on the service descriptions and information sent from the systems. As such, the descriptions need to contain enough information for the user to understand them. The information sent from a selected service may also need to have a semantic protocol if it is to be linked. This is particularly relevant when the linked services rely on local and possibly streaming information such as the locations of mobile objects. To enforce this, a protocol will be developed for the service descriptions to state their location and data type if applicable. The server can parse the descriptions to extract this information and store it with the service to be checked if a service link is requested. Generally, the service linking option should be kept as unconstrained as possible to allow more creative combinations of services.

Since the main focus of this paper is scalability, the limits of the infrastructure need to be established. The second experiment described in this paper provides an initial set of tested parameters but does not test the limits. The components of the infrastructure are the user and system clients, and the server. Many of these components' limits, such as response times and bandwidth, are dependent on the host computer and network architecture. From the user's perspective, any number of users can be on line, since they only affect the infrastructure's scalability when they issue requests. From the system perspective, scalability *is* affected since the systems' services and contact information is stored on the server, sometimes with task buffering. Therefore, many of the scalability limits can be analyzed on the server, such as the maximum number of systems, services, and task buffers it can handle, and also user information such as request burst rate and request response times.

The infrastructure we developed has demonstrated itself to be a viable approach to large-scale human-robot interaction. The experiments in this paper are carried out in simulation, which is currently the most logical method of developing and testing an infrastructure for large-scale interaction. However, an application intended for the real world needs to be tested in the real world. This will be part of the next testing phase of the infrastructure using the same robot systems as the experiments in this paper. This will further validate our claim that this infrastructure is a realistic approach to large-scale human robot interaction.

BIBLIOGRAPHY

- [1] R. Murphy, and E. Rogers, "Cooperative Assistance for Remote Robot Supervision", *Presence, Special Issue*, vol. 5, no. 2, pp. 224-240, 1996.
- [2] A. Agah, "Human Interactions with Intelligent Systems: Research Taxonomy", *Computers and Electrical Engineering*, 2001, no. 27, pp. 71-107.
- [3] K. Kawamura, R. Peters, C. Johnson, P. Nilas, P. and S. Thongchai, "Supervisory Control of Mobile Robots Using Sensory Egosphere", *IEEE*

- Symposium on Computational Intelligence in Robotics and Automation*, Jul. 29-Aug. 01, Banff, Canada, 2001.
- [4] K. Ali, and R. Arkin, "Integration of Reactive and Telerobotic Control in Multi-agent Robotics Systems", *From Animals to Animats 3: SAB 1994*, pp. 473-478.
 - [5] R. Arkin, T. Collins, and Y. Endo, "Tactical Mobile Robot Mission Specification and Execution", *Mobile Robots XIV*, Boston, MA, September, 1999, pp. 150-163.
 - [6] H. Jones, and M. Snyder, M. "Supervisory Control of Multiple Robots Based on a Real-time Strategy Game Interaction Paradigm", *Proceedings of the 2001 IEEE International Conference on Systems, Man, and Cybernetics*, Tucson, AZ, October 2001, pp. 383-388.
 - [7] P. Rybski, I. Burt, T. Dahlin, M. Gini, D. Hougen, D. Krantz, F. Nageotte, N. Papanikolopoulos, and S. Stoeter, "System Architecture for Versatile Autonomous and Teleoperated Control of Multiple Miniature Robots", *Proceedings of the 2001 IEEE International Conference on Robotics and Automation*, Seoul, Korea, 2001.
 - [8] T. Finin, R. Fritzson, D. McKay, and R. McEntire, "KQML as an Agent Communication Language", *Proceedings of the Third International Conference on Information and Knowledge Management*, ACM Press, 1994, pp. 456-463.
 - [9] S. Goldsmith, S. Spires, and L. Phillips, "Object Frameworks for Agent System Development", Sandia National Laboratories, Tech. Rep. SAND98-0742A, 1998.
 - [10] L. Phillips, S. Goldsmith, and S. Spires, "CHI: A General Agent Communication Framework", Sandia National Laboratories, Tech. Rep. SAND98-2825C, 1998.
 - [11] S. Ponnekanti, B. Lee, A. Fox, P. Hanrahan, and T. Winograd, "ICrafter: A Service Framework for Ubiquitous Computing Environments", *Proceedings of the Ubiquitous Computing Conference. Lecture Notes in Computer Science*, vol. 2201, 2001, pp. 56-75.
 - [12] A. Fox, B. Johanson, P. Hanrahan, and T. Winograd, "Integrating Information Appliances into an Interactive Workspace", *IEEE Computer Graphics and Applications*, May/June, 2000, pp. 54-65.
 - [13] J. Lee, and H. Hashimoto, "Intelligent Space - Its concept and contents", *Proceedings of the 2000 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Takamatsu, Japan, 2000, vol. 2. pp. 1358-1363.
 - [14] I. Elhajj, J. Tan, N. Xi, W. Fung, Y. Liu, T. Kaga, Y. Hasegawa, and T. Fukuda, "Multi-site Internet-based Cooperative Control of Robotic Operations", *Proceedings of the 2000 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Takamatsu, Japan, 2000, pp.826-831.
 - [15] H. Surmann, and M. Theißinger, "ROBODIS: A Dispatching System for Multiple Autonomous Service Robots", *Proceedings of the FSR'99 Robotics Applications for the Next Millenium*, Pittsburgh, PA, 1999, pp. 168-173.
 - [16] M. Nicolescu, and M. Mataric, "Learning and Interacting in Human-Robot Domains", *Special Issue of IEEE Transactions on Systems, Man, and Cybernetics, Part A: Systems and Humans*, vol. 31, no. 5, pp. 419-430, September, 2001.
 - [17] B. Dorhonceanu, B. Sletterink, and I. Marsic, "A Novel User Interface for Group Collaboration", *Proceedings of the 33rd Hawaii International Conference on System Sciences*, Maui, Jan. 4-7, 2000.
 - [18] B. Gerkey, R. Vaughan, K. Støy, A. Howard, G. Sukhatme, and M. Mataric, "Most Valuable Player: A Robot Device Server for Distributed Control", *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2001)*, Wailea, Hawaii, Oct. 29-Nov. 3, 2001, pp. 1226-1231.
 - [19] B. Gerkey, and M. Mataric, "Principled Communication for Dynamic Multi-Robot Task Allocation", *Experimental Robotics VII, LCNIS 271*, Springer-Verlag, 2001, pp. 353-362.
 - [20] B. Gerkey, and M. Mataric, "Sold!: Auction methods for multi-robot coordination", *Accepted for publication in IEEE Transactions on Robotics and Automation Special Issue on Multi-robot Systems*, Tamio Arai, Enrico Pagello and Lynne Parker Eds., 2002.
 - [21] B. Jung, and G. Sukhatme, "A Region-based Approach for Cooperative Multi-Target Tracking in a Structured Environment," accepted, *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2002.

APPENDIX A - THE SERVICE LINKING FEATURE

Service linking allows one service to provide information necessary for another service to operate. These types of

service are termed *data* and *command*. A data service provides the information for a command service. An example of a command service is the service provided by a delivery robot that needs a destination for delivering a note or package. A data service can provide the location of a person wandering around a building occupied by the command-service robot. Linking these two services allows the delivery robot to deliver the note to the person (if it can keep up) even though the person's location is constantly changing. This assumes that the semantic information required between the two systems is understood, which may be difficult to determine. Any data-providing service to be linked to it needs to supply coordinates in the same local coordinate reference frame and of the same type. This is a relatively low-level example and the linked data could easily be more symbolic, such as supplying the identification number of an RFID tag that a person is wearing, to facilitate the service robot's search for that person. The data can also be heterogeneous if the link is to be used as a sequential trigger, such as once service A completes, start service B. To allow more freedom with service linking, the semantics should only be an issue when the services are running concurrently and not sequentially. Services capable of this will be identified in their service description.