

A Framework for Structured Distributed Object Computing

K. Mani Chandy, Joseph Kiniry, Adam Rifkin, and Daniel Zimmerman*

`infospheres@cs.caltech.edu`

Computer Science 256-80

California Institute of Technology

Pasadena, California 91125

<http://www.infospheres.caltech.edu/>

Abstract

This paper presents a four-faceted framework for distributed applications that use worldwide networks connecting large numbers of people, software tools, monitoring instruments, and control devices. We describe a class of applications, identify requirements for a framework that supports these applications, and propose a design fulfilling those requirements. We discuss some initial experiences using the framework, and compare our design with other approaches.

*The Caltech Infospheres Project is sponsored by the Air Force Office of Scientific Research under grant AFOSR F49620-94-1-0244, by the CISE directorate of the National Science Foundation under Problem Solving Environments grant CCR-9527130, by the Center for Research in Parallel Computing under grant NSF CCR-9120008, by the Advanced Research Projects Agency, and by Parasoft, Inc. and Novell, Inc. This paper is partially based upon work supported under a National Science Foundation Graduate Fellowship.

1 Personal Command and Control Applications

The global information infrastructure will soon connect large numbers of processes that manage devices and human interfaces. Interprocess communication will allow processes to respond to events on such devices as medical monitoring equipment, scientific instruments, home appliances, and security systems, and on such software as scheduling programs, document management systems, Web browsers, and complex computation engines.

The contribution of this paper is a simple, generic framework for developing distributed systems for personal applications. By employing our framework, developers can quickly build interactive command and control processes that run over the Internet. Our framework is composed of four facets: (i) *processes* are persistent communicating objects; (ii) *personal networks* provide wiring diagrams and behaviors for these connected processes; (iii) *sessions* are transactions performed by the processes participating in a personal network; and (iv) *infospheres* are custom collections of processes for use in personal networks.

Infospheres and Personal Networks. *Warfighter's infosphere* is a term coined by the military to represent the electronic interface between a military unit and its environment. This human-computer interface is provided by the military *C4I* (command, control, communications, computers, and intelligence) infrastructure.

Our goal is to provide a framework for transitioning the concepts of infospheres and C4I to individuals and small organizations to create analogous lightweight command and control systems. *Personal networks* are roadmaps for such systems, specifying processes arranged in a topology, with a specified cooperative behavior. For example, a person in Nevada may have a *emergency notification personal network* that incorporates processes for medical monitoring devices in her parents' home in Florida, security and utility systems in her home and car in Reno, a global position sensing device on her teenage son's car in Montréal, a Nikkei Market stock ticker tape, and software programs that monitor urgent pages and e-mails.

Task Forces for Organizations. Personal networks can also be used by institutions and businesses to create task forces to handle short-term situations. The structure of personal networks comprises the organizational, informational, and workflow structures of the corresponding task force. *Workflow* describes the manner in which jobs are processed in stages by different processes [37].

One example of a task force is a panel that reviews proposals submitted to the National Science Foundation (NSF). Panel members come from a variety of institutions, and the panel has an organizational structure with a general chair, subcommittees, primary reviewers, and secondary reviewers. The panel's informational structure includes the hierarchy of proposals and reviews, and the panel's workflow is the flow of proposal and review copies. The panel has its own organizational, informational, and workflow structures that coexist with those of NSF. In this sense, NSF's organizational and informational structures adapt in a dynamic, but systematic, way to include new people and resources as needed.

Collaborative Computational Science. Consider a collaboration in which researchers at different sites work together on a computational experiment that requires the use of scarce devices, such as the Keck Telescope or an extremely large multiprocessor computer, and software at different locations to be composed together. Now consider a scientist, Samantha, who joins the team after a specific experiment and would like to repeat the experiment to reconstruct the sequence of events that generated its final results. To do so, she must have access to an archive of all the tools used to conduct the experiment, including all the input and output data and all the annotations added by the participants.

Ideally, Samantha should enter the archived virtual laboratory and find it exactly as it was when the experiment was conducted. This gives her the option of either conducting a modification of the experiment herself or witnessing the original experiment as it unfolds before her. As Samantha explores the research team history, she can use the Web to view shared documents and other archived distributed experiments.

Additionally, if the experiment uses scarce resources, Samantha should be able to run the experiment again in one of two different ways. Either she can “virtually” re-run the computation, reusing the data from

the archived experiment without having to schedule time on the scarce resource. Or, Samantha might wish to run the experiment again for real. In that case she would have to schedule time on the telescope and the supercomputer.

An important aspect of this archived virtual laboratory is that the components of the laboratory are geographically distributed: input data sets are generated at one site, a meshing computation is conducted at a different site, and the output data is post-processed prior to visualization at a third site. What is being archived is a distributed system consisting of components from all three sites.

Another important point to make is the the professional scientist, Samantha in this case, does not need to run the archived computation using the same resources and software as the previous researchers did. She is able to engineer her own new software components and devices, either from scratch or through inheritance, to plug in to her computation. The participant can play as active or passive a role in component creation, as dictated by need and expertise. In particular, the incorporation of basic standard communication devices, such as pagers, fax machines, and printers, is key to reifying generic archived experiments in a manner in which is locally contextually correct.

After Samantha explores one archived experiment, she may study the annotations of the experimenters and then follow links to related experiments. She can follow links to later experiments by the same collection of experimenters, explore attempts by other groups to replicate the experiments with different tools, or jump to a document discussing the public policy issues raised by the experiments. She can even reuse data or components from a previous experiment in a new experiment of her own design.

2 Requirements Analysis

A framework to support personal networks (and their components) should satisfy four main requirements: scalability, simplicity, security, and adaptability.

Scalability. Personal networks should scale to include devices, tools, and people connected to the Internet. The critical scaling issue is not the number of processes connected in a personal network, but rather the size of the pool from which the processes in personal networks are drawn. The only limit to the number of processes connected in a personal network is the number of activities that can be managed effectively. However, issues of scaling in naming, connections, and services depend on the size of the global set of processes and resources.

Personal networks should be tolerant of wide ranges of quality of service because the processes in a personal network can exist on a single system or span several continents. The framework should both support large numbers of concurrent personal networks and provide a core set of services for creating and using personal networks.

Simplicity. The usage and programming model for personal networks should be simple enough to be usable by anyone. The simplicity of dialing a telephone led to the widespread use of telephones despite the complexity of the worldwide telecommunications network. If personal networks are to become effective tools, their use should be similarly intuitive. So, the model's API should be easy for programmers to learn quickly, and the accompanying visual tools should allow non-programmers to use palettes of existing constructs to customize their personal networks.

Security. A research instrument shared by several people may have one interface for setting control parameters and a different interface, accessible by a small set of authorized personnel, for accessing the data recorded by the instrument. Also, instruction messages sent to the "modify-parameter" interface may be of a different type than instructions to the "read-data" interface. Therefore, the framework should allow processes to have multiple typed interfaces and provide the ability to set security restrictions on at least a per-interface basis.

Adaptability. It should be possible to create and modify personal networks rapidly and flexibly, because task forces often need to be set up quickly and in an ad hoc manner. Network topologies should be emergent

rather than static, so processes should be able to create and delete connections during a session. Additionally, personal network processes should be able to communicate with applications and devices that were unknown or nonexistent prior to the creation of the personal network. So, the framework should be extensible enough to support interoperability with other distributed technologies.

3 Design of an Extensible Framework

Our framework employs three structuring mechanisms: personal networks, to facilitate long-term collaborations between people or groups; sessions, to provide a mechanism for carrying out the short-term tasks necessary within these personal networks; and infospheres, to allow customization of processes and personal networks.

To illustrate these structuring mechanisms, consider a consortium of institutions carrying out research on a common problem. It has a personal network composed of processes that belong to the infospheres of the consortium members. This personal network is a structured way to manage the collection of resources, processes, and communication channels used in distributed tasks such as simulating financial scenarios, determining meeting times, and querying distributed databases. Each session of this personal network handles the acquisition, use, and release of resources, processes, and channels for the life of a specific task.

Infospheres are discussed in our framework user's guide [17]. This paper focuses on the conceptual models for processes, personal networks, and sessions.

3.1 Conceptual Model: Processes

Processes are the persistent communicating objects that manage devices and interfaces. In our framework, we call these processes *djinnns*.

Process States. A given process can be in one of three states. An *active* process is a process that has at least one executing thread; it can change its state and perform any tasks it has pending, including communications. A *waiting* process has no executing threads; its state remains unchanged while it is waiting, and it remains in the waiting state until one of a specified set of input ports becomes nonempty, at which point it becomes active and resumes execution. Active and waiting processes are collectively referred to as a *ready* process.

Ready processes occupy process slots and can make use of other resources provided by the operating system. By contrast, processes in the third state, *frozen*, do not occupy process slots. In fact, frozen processes do not use any operating system resources except for the persistent storage, such as a file or a database, that is used to maintain process state information.

Freezing, Summoning, and Thawing Processes. Associated with each process is a *freeze* method, that saves the state of the process to a persistent store, and a *thaw* method, that restores the process state from the store. Typical processes remain in the frozen state nearly all the time, and therefore require minimal resources. In our framework, only a waiting process can be frozen, and it can only be frozen at process-specified points. When its freeze method is invoked, a process yields all the system resources it holds.

A ready process can *summon* a frozen process. The act of summoning instantiates the frozen process, causes its thaw method to be invoked, and initiates a transition to the ready state. If a process is ready when it is summoned, it remains ready. In either case, a summoned process remains ready until it receives at least one message from its summoner or a specified timeout interval elapses.

Mobile Processes. Frozen processes can move from one machine to another, but ready processes cannot. This restriction allows ready processes to communicate using our framework’s underlying fast transport layer, that requires unchanging addresses for communication resources. All processes have a permanent “home address” from which summons can be forwarded. Once a process becomes ready at a given location,

it remains at that location until the process is next frozen. The persistent state of a process is always stored at the home address of that process.

3.2 Conceptual Model: Personal Networks

Conceptually, a personal network is a wiring diagram, analogous to a home entertainment system, with directed wires connecting device outputs to the inputs of other devices. We chose this model for its simplicity [8]. A personal network consists of an arrangement of processes and a set of directed, typed, secure communication channels connecting process output ports to the input ports of other processes; its topology can be represented by a labeled directed graph. Note that, unlike home entertainment system components, processes can freely create input ports, create output ports, and change wire connections.

Communication Structures. Processes communicate with each other by passing messages. Associated with each process is a set of *inboxes* and a set of *outboxes*. Inboxes and outboxes are collectively called *mailboxes*. Every mailbox has a type and an access control list, both of which are used to enforce personal network structure and security. These mailboxes correspond to the device inputs and outputs used in the wiring diagram conceptual model.

Process interconnections are asymmetric; a process can connect any of its outboxes to any set of inboxes for which it has references. A connection is a first-in-first-out, directed, secure, error-free broadcast channel from the outbox to each connected inbox. Our framework contains support for message prioritization, available through standard multithreading techniques.

Message Delivery. Our framework communication layer works by removing the message at the head of a nonempty outbox and appending a copy to each connected inbox. If the communication layer cannot deliver a message, an exception is raised in the sender containing the message, destination inbox, and specific error condition. Our system uses a sliding window protocol [29] to manage the messages in transit.

Every message at the head of an outbox will eventually be handled by the communication layer. The conceptual model uses asynchronous messages rather than remote procedure calls, to be tolerant of the range of message delays experienced along different links of the Internet. As a result, we can think about message delivery from an outbox to inboxes as a simple synchronous operation even though the actual implementation is asynchronous and complex.

Dynamic Structures. A process can create, delete, and change mailboxes. The operation of creating a mailbox returns a global reference to that mailbox. This reference can then be passed, in messages, to other processes. Since a process can change its connections and mailboxes, the topology of a personal network can evolve over time as required to perform new tasks.

As long as a process remains ready, references to its mailboxes are valid; when a process is frozen, all references to its mailboxes become invalid. Since all references to the mailboxes of frozen processes are invalid, frozen processes can move and then be thawed, at which point the references to their mailboxes need be refreshed via a summons. Because no valid references to their mailboxes exist, frozen processes cannot participate in sessions.

3.3 Conceptual Model: Sessions

Operationally, a session is a task carried out by (the processes in) a personal network [9]. It is *initiated* by a process in the personal network, and is *completed* when the task has been accomplished. A later session may use the same processes to carry out another task. Thus, a personal network consists of a group of processes in a specified topology, interacting in sessions to perform tasks.

The Session Constraint. We adopt the convention that sessions must satisfy the two part *session constraint*:

1. As long as any process within the session holds a reference to a mailbox belonging to another process

within the session, that reference must remain valid.

2. A mailbox's access control list cannot be constricted as long as any other process in the session holds a reference to that mailbox.

The session constraint ensures that, during a session, information flows correctly between processes.

A session is usually started by the process initially charged with accomplishing a task. This *initiator process* creates a session by summoning the processes that will initially participate. It then obtains references to their mailboxes, passes these references to the other processes, and makes the appropriate connections of its outboxes to the inboxes of the participating processes. We discuss session implementation and reasoning issues in Section 4.

There are many ways of satisfying the session constraint. One simple way is to ensure that once a process participates in a session it remains ready until the session terminates, and that once a process sends its mailbox references to other processes it leaves these mailboxes unchanged for the duration of the session. Another approach is to have the initiating process detect the completion of the task through a diffusing computation, after which it can inform the other session members that the session can be disbanded.

An Example Session. An example of a session is the task of determining an acceptable meeting time and place for a quorum of committee members. Each committee member has an infosphere containing a calendar process that manages his or her appointments. A personal network describes the topology of these calendar processes. A session initiator sets up the network connections of this personal network. The processes negotiate to find an acceptable meeting time or to determine that no suitable time exists. The task completes, the session ends, and the processes freeze. Note that the framework does not *require* that processes freeze when the session terminates.

During a session, the processes must receive the quality of service they need to accomplish their task. Therefore, communication is routed directly from process to process, rather than through object request

brokers or intermediate processes as in client-server systems. Once a session is constructed, our framework's only communication role is to choose the appropriate protocols and channels. A session can negotiate with the underlying communication layer to use the most appropriate process-to-process mechanism. The current framework supports only UDP, but we plan in future releases to support a range of protocols such as TCP and communication layers such as Globus [11].

4 Structuring Mechanisms

Personal networks and sessions can be used not only as structuring mechanisms, but also for reasoning about the services provided to distributed systems.

4.1 Reasoning About Sessions

Consider a consortium of institutions working together on a research project. From time to time, people and resources of the consortium carry out a collaborative task by initiating a session, setting up connections using the personal network, performing the necessary machinations for the task, disbanding the connections, and terminating the session. Furthermore, several sessions initiated by the same consortium may be executing at the same time. For instance, a session to determine a meeting time for the executive committee and a session that reads measurements from devices in order to carry out a distributed computation could execute simultaneously. Moreover, the same process may participate concurrently in sessions initiated by different consortia or task forces. For example, a calendar manager may participate concurrently in sessions determining meeting times for a scout troop and a conference program committee. Our framework allows processes to participate in concurrent sessions [9].

A resource may be requested by a session in either *exclusive mode* or *nonexclusive mode*. For example, a visualization engine may need to be in exclusive mode for a task: while the task is executing, no other task can access it. However, a process managing a calendar can be useful in nonexclusive mode: several sessions

can not only read the calendar concurrently, but also modify different parts of the calendar concurrently.

Because we cannot predict a priori the applications and sessions that will run concurrently, we restrict access to modify the states of the processes participating in a given session, to reason about that session's behavior. Such restrictions are currently provided in thread libraries by mutexes and monitors; our challenge is to provide similar constructs with our framework for use in distributed systems in a generic, extensible, and scalable manner.

4.2 Services for Sessions

New capabilities are added to our framework either by subclassing existing processes or by extending the framework. A *service* is a framework extension that is applicable to an assortment of distributed algorithms. Examples include mechanisms for locking, deadlock avoidance, termination detection, and resource reservation.

Locking Mechanisms. Even if a process participates concurrently in several sessions, there are points in a computation when one session needs exclusive access to certain objects. For example, at some point, the session determining the meeting time for a program committee needs to obtain exclusive access to the relevant portions of the calendars of all the committee members. Therefore, one service our framework should provide is the acquisition of locks on distributed objects accessed during a session. A great deal of work exists relating to locking in distributed databases and distributed transaction systems [14, 24]. Presently, our framework provides only an exclusive lock on an object, but the framework can be extended to include other types of locks, such as read and write locks.

Deadlock Avoidance. If sessions lock objects in an incremental fashion, deadlock can occur. For instance, if one session locks object A and then object B , and another session locks B and then A , the sessions may deadlock because each session holds one object while waiting for the other. Therefore, our framework deals

only with the case where a session requests locks on a set of objects only when it holds no locks; a session must release all locks that it holds before requesting locks on a different set of processes. An alternative solution would be to allow incremental locking in some total ordering, but we are not exploring this solution because it does not scale to distributed systems drawn from a worldwide pool of objects.

Termination Detection and Resource Reservation. Other services that can be extended into our framework include session termination detection and resource reservation. Termination detection can be used by an initiating process of a session to, for instance, determine when the states of the processes involved in the session need to be “rolled back” in the event of a failure. Resource reservation is a generic service through which the resources required by a session can be reserved for some time in the future. For instance, one might reserve the visualization engine at location X and the monitoring instrument at location Y for the earliest time after 5:00 PM today.

5 Experience With Our Framework

Two examples illustrate the ease with which programmers have used our framework to develop distributed systems. Using our model and middleware packages, a programmer was able to specify, design, reason about, and implement a distributed calendar application in under a week. Since our infrastructure handled the communication layer, the programmer could concentrate his skills on the high-level design and implementation.

Also, using our framework, given a specification of the processes and communication protocols, the students in an undergraduate class at Caltech were able to write processes that participated in a five-card draw poker tournament session. The students were given a week to design, reason about, and implement their poker-playing processes; we spent approximately the same amount of time specifying those processes and their interactions.

In both these cases, patterns helped the programmers develop their code quickly. Patterns encapsulate software solutions to common problems [13], and our framework has incorporated some applications of concurrency patterns in Java [21]. Initial experience with our framework has suggested several other patterns, both for collaborations between processes and for state-transition systems.

5.1 A General Resource Reservation Framework for Scientific Computing

This work makes three contributions for distributed resource allocation in scientific applications. First, we developed an *abstract model* in which different resources are represented as tokens of different colors; in this model, processes acquire resources by acquiring tokens of appropriate colors. Second, we developed *distributed scheduling algorithms* that allow multiple resource managers to determine custom policies to control allocation of the tokens representing their particular resources. These algorithms allow multiple resource managers, each with its own resource management policy, to collaborate in providing resources for the whole system. Third, we developed an actual implementation of a distributed resource scheduling algorithm framework using our abstract model. This implementation using Infospheres showcases the benefits of distributing the task of resource allocation to multiple resource managers.

Often in scientific computing, a user needs access to several distributed heterogeneous resources. For instance, consider a scientist conducting a distributed experiment [7] that requires a supercomputer in one location, a visualization unit in another location, and a special high quality printer in still another location. All three resources are essential to the experiment; so, the scientist needs to *synchronously* lock and use all three *distributed resources* for the same time period to complete the computing task. The distributed heterogeneous resources together form a *networked virtual supercomputer* or *metacomputer* [5]. Furthermore, the scientist wants resources to be scheduled automatically as a service of the appropriate software, whether with or without the inclusion of specific supplemental information such as the times the user is available to perform the experiment.

Traditional Metacomputing: Central Scheduling. Traditional metacomputing resource allocation [23, 1] uses a central authority for scheduling, usually for efficiency. As a simple example, the IBM SP2 employs a scheduling algorithm [22] that reduces the wait time of jobs requiring only a few nodes, if these can be scheduled without delaying more computationally intensive jobs.

Alternative Approach: Distributed Scheduling. By contrast, consider the computation needs of users requiring resources managed by different groups in different geographic regions. Scheduling in this case is more complicated because it is impractical for individual sites to “know” global information that would help them to do more efficient scheduling [11].

The owner of a set of resources may have resource management policies that are different from those of owners of other resource sets. Our challenge is two-fold: (i) to establish methods of cooperation among owners so that the collection of owners offers system-wide resources to users, and (ii) to make the algorithms scalable so that new providers of resources can enter the common resource pool quickly and semi-autonomously.

Resources as Tokens. The infrastructure for reserving resources in a distributed system is required in many applications. Our research deals with designs and implementations of distributed resource management schemes that coordinate different resource management schemes for different sets of resources. Though this paper addresses resources used in metacomputing, our research deals with resources in many distributed applications.

A convenient abstraction for such applications represents each indivisible resource by an indivisible *token* of some color [6]; different types of resources have different colors. For instance, a node of an IBM SP2 can be represented as a token of the IBM SP color. Likewise, a room in a hotel can be represented by a token of the hotel color.

Our model deals with time explicitly. So, a reservation can be made for 64 nodes of an IBM SP2 for 10 contiguous hours, or a hotel for seven nights. We also deal with locality in the sense that requests can be

made for collections of tokens that are “near” each other.

At this point, we do not deal with arbitrary relationships and constraints between token colors. For instance, we do not automatically handle the relationship that a 200 MHz Pentium Pro is 1.8 times as fast as a 100 MHz Pentium Pro. Each kind of machine is represented by a token of a different color, and we do not deal automatically with relationships between tokens of different colors. Methods for registering relationships between colors and for automatically exploiting these relationships is a subject for future research.

The centralized IBM SP2 scheduling algorithm relies on knowledge of how many nodes each process needs to “promote” less computationally-intensive tasks as necessary. On the other hand, if each node in a supercomputer was to be scheduled independently in a distributed way, efficient scheduling becomes much more difficult. We believe that as metacomputing applications use distributed heterogeneous systems, they will need algorithms for efficient distributed resource scheduling.

Contributions of this work. This work represents a general framework for heterogeneous resource reservation. Within this framework, we developed the simple object-oriented Java implementation using Infospheres. Specific contributions of this work [31] are:

- *We provide an abstraction for distributed resource management problems that fits many, but not all, applications.* Our implementation is suitable for all applications that fit our abstraction.
- *Our implementation is distributed, and it coordinates multiple resource managers, each with its own policy.* In essence, our research deals with program composition: we show how to compose multiple resource managers to obtain a distributed resource manager for a distributed system.
- *We take into account user preferences efficiently.* Our implementation automatically reserves resources without significant user intervention and without resorting to excessive message-passing through a novel technique of passing a short Java program to the resource manager that controls access to each resource. The mobile program performs its resource scheduling more quickly on the server than on the client

when network latency is high.

5.2 Networked Data Objects

More and more electronic documents are being created by groups of people; however, very few programs are designed to accommodate groups of people editing a document. Sometimes people sit around the computer screen and share the keyboard to work on a document together, but this has a locality problem: it is difficult for everyone to see the screen, and it can be difficult for people to add ideas to the document as soon as they have them. Other times, each person modifies his own copy of the document, and then these documents are merged into a single document at a later date. Unfortunately, attempts to merge these documents can introduce interference problems among different edits of the same material. A fresh approach would be nice: we would like to allow many people edit the same file at the same time, with changes made by one person on his or her computer being immediately visible to everyone else. With this goal in mind, we set out to use the Infospheres packages to develop a shared distributed editor of shared documents [35].

Our original design incorporated a strict client/server approach, which had several advantages: clients never had to see the entire file they were editing, and small client/server messages could describe users' modifications to the documents. A server allowed an easy handling of race conditions as well, since the concurrent changes of the same sections of text can be resolved by the order of their arrival time at the server.

However, when we redesigned the editor as a true peer-to-peer application, the document data could belong to the entire network of editors running on distributed machines, instead of belonging to one of the individual editors. This provides fault-tolerant characteristics: if any editor were to suddenly drop from the network, the rest would be able to recover easily.

An Object with Data Distributed over a Network. From this application, we developed the idea of a generic *Network Data Object (NDO)*. A network data object, in essence, is an object (with component

data and associated methods) that can be shared among multiple processes communicating over a network. The NDO exists only as long as those processes exist, and the data in an NDO can be manipulated or read by any process that is sharing the NDO. Changes made by one process participating in the NDO session are immediately visible to all other processes, and processes can join and leave the session as their respective users desire. The document's data disappears when all of the processes disappear, although any process can save a local copy of the document before exiting.

A text editor is an interesting distributed application to write, mainly because it has a variety of requirements. The cursors move quickly, and require rapid-response synchronization across all the communicating processes. The lines of text typed by distributed users are generally independent from each other, and each user process is likely to only modify one or two lines at a time. The file as a whole does not change much from ascertainable modification to ascertainable modification. In addition, a user scrolling through the file without making changes expects scrolling to flow naturally, as if all the data resided locally (despite the fact that modifications to the data may exist at different locations in the network and these changes are only communicated when needed).

Distributing the Data Object. A Network Data Object is a collection of data that is shared across processes connected by a network. From the standpoint of the processes involved, this is analogous to having access to these variables as stored in shared (local) memory. In our implementation, all of the shared data is stored in types derived from a base class, `SharedVar`. No `SharedVar` exists standalone; every `SharedVar` in an application is contained within an NDO.

Using an NDO instead of more traditional networking methods, such as message passing, simplifies distributed application writing considerably. Instead of viewing the program as communicating data between multiple computers, programming with an NDO is similar to writing a multithreaded, non-networked application. Any Java primitive type or class can be stored in a `SharedVar`, and the NDO will determine what messages need to be sent to keep all of the communicating processes up-to-date. In addition, because NDOs

are lightweight and there is no need for interface compilers, skeletons, or stubs, there is less complexity for an application developer to manage.

Automating the Data Object's Network Facilities. The NDO model makes the networking decisions straightforward: a programmer no longer needs to worry about deadlock because the NDO ensures that the application has the latest copy of the data using locks. The price for this convenience is that the application cannot lock one variable and then lock another without unlocking the first, meaning that the application cannot lock one variable on the basis of the value of another.

For example, when a user presses a key, that character should be inserted after the cursor on the current line. But, the application cannot know the cursor's location for sure, until it is locked. Once the cursor is locked, the application now knows which line to lock, but nothing can be locked without first releasing the lock on the cursor. However, once the application releases the lock on the cursor, the cursor may move to another line. The application cannot assume that the line the cursor was originally on is still the correct one: perhaps another user has inserted a line and all the line numbers have been shifted up by one. If the application naively inserts a character into the line where the cursor used to be, it would seem to the user as if the character was inserted into the line above the cursor!

Our solution to this problem is to look at the current cursor position to estimate which line to lock. Then, the application jointly locks both the cursor and the line corresponding to this estimate. If the cursor is not on the line that was locked, this process can be repeated with the cursor's new position as the estimate of the current line. Fortunately, the NDO handles data with efficient mechanisms for locking and re-locking. Of course, success is not guaranteed against a demonic scenario with this method: the unlikely possibility that every time the application releases the cursor to try again, another process could insert a new line, moving the cursor. Overall, the NDO provides a good solution for typical data objects shared across a network.

5.3 DALI: The Distributed Artificial Life Infrastructure

DALI [30] is a Java framework built on top of the Infospheres Infrastructure designed for the creation of complex, asynchronous artificial-life simulations that run in a distributed environment. Since DALI is written in Java, DALI simulations can run on most available operating systems. Also, since DALI uses the Infosphere Infrastructure, participating systems can be located anywhere on the Internet and simulations can scale to a very large number of active objects.

There are two fundamental aspects of DALI that differentiate it from other discrete simulation systems. First, DALI's object models reflect the hierarchical relationships that exist in many complex systems; thus creation, organization, management, and communication in complex simulations is anthropomorphic, natural, and efficient. Second, all aspects of DALI are asynchronous. No object synchronization is enforced or required, unless necessary for the simulation, in which case the degree of synchronization is up to the developer. Several aspects of the Infospheres Infrastructure can assist with such synchronization, including the service mechanism and synchronous message passing. Thus, distributed simulations are efficient and scale well.

While there are several software frameworks available for building artificial life simulations [2, 25], few are flexible enough to build arbitrary new scenarios. In many cases, simulations are synchronous, run on a single computer, and the simulation entities must be either entirely passive, completely reactive, or are required to operate in a fixed time quantum. Needless to say, for distributed simulations, the software must be ported to numerous architectures and there is a potential for host of installation and maintenance woes.

The DALI framework consists of seven base classes. Four classes consist of the DALI hierarchical simulation model: *Agent*, *Gang*, *Region*, and *Simulation*. Note *Simulation* inherits from *Region*, *Region* inherits from *Gang*, and *Gang* inherits from *Agent*. Behavior and characteristics, whether they are an *Agent's* color, a *Gang's* aggressiveness, a *Region's* rainfall, or a *Simulation's* database location are all encoded in classes derived from the *Genome* base class. If we wish to model non-active components (like water, trees, grass,

etc.), we inherit from the *Resource* class. Finally, watching and recording a simulation are handled by classes derived from the *Observer* base class. Together, these seven base classes provide a flexible, sophisticated simulation structure. An example will help clarify how these classes are used in a real simulation.

An Example. Imagine modeling two chimpanzee colonies. One needs to simulate the animals and their environment. Let's consider the environment first, since it is simpler than the animals.

Assume that the colonies reside in a ten kilometer square area of jungle. We will represent each square kilometer of jungle with a *JungleRegion* class, which inherits from *Region*. Thus, we will have one hundred instances of *JungleRegion* running in our simulation, and these one hundred distributed objects can be spread across any number of systems.

Within each region of jungle, there exist a number of natural resources that are part of the simulation. Trees, water, and food of various kinds are all part of the simulation. These are all entities with which the active agents of the simulation can interact (they can be climbed, drunk, and eaten, respectively). Each entity is represented by an instance of the appropriate subclass of *Resource*. These resources are real, physical elements within the simulation (*i.e.* they have a location, size, weight, value, *etc.*), thus their instances initially reside on the same CPU as the *JungleRegion* in which they are found. In each *JungleRegion* we might have several thousand unique instances of such resources.

Generally, an individual agent's characteristics are represented by a subclass of *Genome*. Variations in the genome determine how one agent differs from another. Thus, general behavior trends can be controlled and varied through genome modification, and such trends can be inherited from generation to generation by genome transfer via reproduction. Each of the resources above (*Tree, Water, etc.*) has a genome in which it encodes that resource's characteristics (how tall a tree can grow, how clean the water is, *etc.*).

Each chimp in the simulation is represented by an instance of the *Chimp* class which inherits directly from *Agent*. (If one wished to be more precise, *Mammal* and *Primate* classes could be used, but we will not be that detailed here.) Each chimp has a genome as well, encoding the characteristic's particular to that

chimp.

What about *Gang*? There are many times in a simulation when one wishes to group agents together, literally or figuratively. For instance, let's assume that chimps have sets of characteristics depending on the chimp's sex. To simulate this, we will create *MaleChimp* and *FemaleChimp* classes which inherit from *Gang*. Thus "gang", in this context, means "collection of like things". All males chimps in the simulation will be a member of (referred to) by a singular (possibly replicated) instance of *MaleChimp*; likewise for females. Other similar *Gang*-derived classes include *ChimpColony*, *InfantChimp*, *AdolescentChimp*, *AdultChimp*, *PassiveChimp*, and *AggressiveChimp*.

The gang-derived classes provide built in capabilities for group-wise communication and simulation control. For an example of group-wise communication, suppose the leader of a colony wishes to communicate something to all the members of his clan (say, where some food is located). Within the system, sending a message to the agent "Colony1" (an instance of *ChimpColony*) would automatically propagate that message to all members of "Colony1". Of course, there are spatial and physical constraints on such communication, so a message would instead be sent to the "Colony1" instance associated with the *JungleRegion* in which the lead chimp is currently residing.

Group-wise simulation control is continually utilized through the use of the *Observer* class's derivatives, something of a distributed system *View* (a part of the Model-View-Controller model summarized in [19]). For example, suppose we wish to record the weight of all infant chimps through the course of a simulation. An appropriate subclass of *Observer* is written to either poll the top-most *InfantChimp* instance for each chimp within its domain's weight. Or, a publish/subscribe model can be used in which each *InfantChimp* class is told to send its current weight to a specific instance of *Observer* every specific time interval. Either of these communications is accomplished through a single send on the part of the simulation or observer responsible for initiating the interaction which is then propagated in a multicast-ish fashion through the simulation.

Use of the Infosphere Infrastructure. The Infosphere Infrastructure was extremely helpful in building the DALI system in several ways. First and foremost, we wished to be able to use thousands (if not tens or hundreds of thousands) of entities in a simulation. Most single system models, and all distributed system models, require objects to be active in memory for their entire lifecycle. In this case, this is effectively the lifecycle of the simulation, which can range from hours to months. For example, assume each object only takes up 512 bytes of memory (a conservative estimate, to be sure). If we were using a single system, a moderately sized simulation (100,000 objects) would require almost fifty megabytes of memory just for data; a large simulation might require several gigabytes of main memory.

Even if we distribute this load, only the tiniest fraction of the entities in the system are active at any point in time. (If a tree falls in the forest and there is no one there to hear it...). Assume there are fifty chimps in each colony. Only those fifty objects, the objects with which they are directly interacting (say, another couple of hundred), and the meta-objects (*Gangs*) associated with them (another few dozen) need be in system memory. Why use up all of those resources (memory, CPU time, disk for swapping virtual memory) if it is simply unnecessary? One might argue that this type of resource management could be handled by a single system with a great deal of virtual memory, but then the simulation has little or no control over the policies of its active objects.

Distributed simulation systems have advantages and disadvantages. The primary benefit of such a system is that of increased resource availability. The simulation can potentially use more processors, memory and disk storage than it could when it run on a single machine. Some of the secondary benefits include the ability for multiple collaborators to participate in a simulation, the possibility of taking advantage of scarce resources (like multiprocessor machines or special purpose hardware), and the capability of researchers, instructors, and students to watch and participate in a simulation as it progresses.

The primary sacrifice made in using a distributed system is that of communication performance across nodes. Communication over a network is very expensive when compared to local communication. Therefore

if a simulation is structured without taking into account network issues, it can perform poorly. Normally, if one maps the simulation structure to the network (as we did in the above example), this problem is alleviated.

Contributions of DALI. DALI provides an innovative architecture not just for artificial life experimentation, but for asynchronous distributed simulation in general. Many scientific models fit into the patterns supported by DALI. Its capabilities for group-wise communication and management naturally lends itself to complex simulations that involve collections of entities, and the Infosphere Infrastructure’s capability in dealing with extremely large numbers of distributed objects means that simulations can be extremely large. Finally, the distributed nature of the system lends itself to supporting the scheduling and use of rare network resources like multiprocessor machines for compute-intensive tasks and high power graphics workstations for visualization, as described in section 5.1 above.

More information on DALI can be found in [30].

6 Patterns in the Infosphere Infrastructure

Collaboration Patterns. Several patterns of collaboration network topologies have emerged from our exploration of personal networks. A personal network consisting of a “master” process maintaining all modifications to an object shared by the other objects of the personal network fits the *Personal Network Star* pattern. For example, a concurrent document editing system with a single process responsible for maintaining changes during a personal network would match this pattern. This pattern roughly corresponds to a system with a single server with a set of clients, though more sophisticated systems (such as a hierarchy with multiple servers and multiple clients) could also be developed. The Personal Network Star pattern was employed in both the calendar and poker applications mentioned above.

A personal network in which each of the processes collaborate without a master, with all modifications announced to the entire group, fits the *Personal Network Full Connection* pattern. For example, a concurrent

document editing system in which every process sends every modification to every other process, and every process is responsible for updating the local view of the shared object, would match this pattern. This pattern roughly corresponds to a peer-to-peer distributed system, though more sophisticated systems (such as different priorities for different peers) could also be developed.

A personal network in which messages are propagated in a ring during collaboration fits the *Personal Network Ring* pattern. For example, a document editing system in which the session-initiator process has a document and makes changes to it, then sends the modified document to the next process for it to make changes, and so on until the document is returned to the session-initiator process, would match this pattern. This pattern roughly corresponds to a workflow distributed system, though more elaborate workflow templates could also be developed. The Personal Network Full Connection and Personal Network Ring patterns were used in the poker applications mentioned above.

We are investigating other middleware patterns as well, such as hierarchical broadcast using publishing and subscribing processes, and dataflow using waiting and notification processes.

State-Transition System Patterns. In addition to collaboration patterns among the processes in a personal network, our experiences with user interfaces for describing network topologies has given rise to a pair of state-transition system patterns. Using these patterns, developers can design and reason about the changes of state in the processes participating in a session.

One pattern is the *Transition on Modes* pattern, in which the processes change their states based on a combination of their respective modes and the messages they receive on their inboxes. For example, in a distributed accounting system, a money receipt message would cause different ledgers to be modified, based on whether the controlling process was in “accounts receivable” or “accounts payable” mode.

Another pattern is the *Transition on Functions* pattern, in which the processes change their states based on a function of the information contained within the messages they receive on their inboxes. For example, in a distributed accounting system, an income transfer may require different actions based on how much

money is being transferred, for tax shelter purposes.

7 Framework Implementation

Version 1.0 of our tools and models, released in January 1998, is classified in the “white box framework” level of the taxonomy given by the framework pattern language [32]. With the addition of more applications, services, visual builders, and language tools, we are developing a “black box framework.” To guarantee widespread, unrestricted use, our initial revisions of the framework (release from August, 1996 through August, 1997) were developed using Sun’s Java Developer’s Kit (JDK) 1.0.2. The final 1.0 version uses Java 1.1.

The next generation of the work, version 2.0, scheduled for release in spring 1998, is a complete re-design and reimplementaion of the infrastructure, and takes advantage of the following newly standardized packages:

- *Remote Method Invocation* (RMI) for a proxy-based distributed object model.
- *Object Serialization* facilities for packing and unpacking objects and messages (both for communication and for persistent storage).
- *Java Database Connectivity* support for persistent storage of, and queries on, process, state, and interface data.
- *Interface Definition Language* (IDL) packages for interoperability with CORBA distributed objects.
- *Security* packages for communication encryption and process authentication.
- *Reflection* packages for innovative structuring of emergent personal networks and process behavior.

An overview of this innovative new work is available on the web [18] and has been published as a workshop position paper.

8 Related Work

Frameworks are reusable designs for software system processes, described by the combination of a set of objects and the way those objects can be used [32]. Our framework consists of some middleware APIs, a model for using them, and services and patterns that are helpful not only in inheriting from objects, but extending them as well. These features allow the reuse of both design and code, reducing the effort required to develop an application. In this sense, our framework is comparable to other metacomputing, component, and concurrency frameworks.

Metacomputing Frameworks. Our framework efforts are similar to recent metacomputing endeavors in that we use the Internet as a resource for concurrent computations. *Globus* provides the infrastructure to create networked virtual supercomputers for running applications [11]. Similarly, *NPAC at Syracuse* seeks to perform High Performance Computing and Communications (HPCC) activities using a Web-enabled concurrent virtual machine [12]. *Javelin* is a Java-based architecture for writing parallel programs, implemented over Internet hosts, clients, and brokers [4]. *Legion* is a C++-based architecture and object model for providing the illusion of a single virtual machine to users for wide-area parallel processing [15]. Although our framework could be used for metacomputing applications, we provide neither seamless parallelism, nor facilities for developing extremely high-performance applications. Rather, we provide mechanisms for programmers to develop distributed system components and personal networks quickly, and we plan to provide mechanisms for non-programmers to customize their components and their personal networks easily.

Component Frameworks. Many other framework systems also have the goal of creating distributed system components. The *ADAPTIVE Communication Environment* (ACE) provides an integrated framework of reusable C++ wrappers and components that perform common communications software tasks [33]; this framework is amenable to a design pattern group useful to many object-oriented communication systems [34]. *Hector* is a Python-based distributed object framework that provides a communications transparency layer

enabling negotiation of communication protocol qualities, comprehensive support services for application objects, and a four-tiered architecture for interaction [3]. *Aglets* provide a Java-based framework for secure Internet agents that are mobile, moving state along with the program components themselves [20]. We differ from these efforts because our emphasis is on reasoning about global compositional distributed systems.

Concurrency Frameworks. We have considered several previous approaches to concurrent communicating processes in developing our framework. The *Communicating Sequential Processes* (CSP) model assumes each process is active for the entire duration of the computation [16]. Like Fortran M [10], we implement this model, adding such implementation artifacts as dealing with process setup and removal, and permitting prioritized waits to resolve resource contention. Unlike Fortran M, sessions provide a hybrid technique for running communicating distributed processes that are frozen when they are not performing any work, yet have persistent state that can be revived whenever a new session is initiated.

This persistence model is similar to mechanisms provided as recent ORB services [36]. However, the *CORBA* process model, implemented using the Basic Object Adaptor (BOA) of a given Object Request Broker (ORB), maintains that only the broker stay active for the entire duration of the computation [26]. Like *Client-Server*, *Remote Procedure Call*, and *Remote Method Invocation* systems, CORBA only spawns remote processes to perform isolated remote tasks. In our framework, the model supports interaction not just through a broker or server, but also directly between the ports of distributed processes in a peer-to-peer fashion.

9 Summary

In this paper, we have presented a framework for developing personal networks and their component processes, for using those processes in sessions to perform distributed tasks, and for reasoning about those processes and sessions. Our approach is novel in its simplicity, scalability, and flexibility; new system pro-

cesses can be developed by inheriting from framework processes or by extending framework services.

In further research, we plan on using the framework to develop more substantial personal networks, including several task forces: research consortia that use instruments, computation engines, and visualization devices at different sites; oversight committees for conferences or journal publications; and working groups whose members hail from different organizations. In addition, we plan to investigate an array of services for use with the framework, including tools for active process mobility, distributed collaboration, termination detection, resource management, and session coordination.

References

- [1] D. Abramson, R. Sasic, J. Giddy, and B. Hall, ‘Nimrod: A Tool for Performing Parameterised Simulations Using Distributed Workstations’, *Proceedings of the Fourth IEEE Symposium on High Performance Distributed Computing*, 1995.
- [2] A. Adami and C. T. Brown, ‘Evolutionary Learning in the 2D Artificial Life System Avida’, *Artificial Life IV: Proceedings of the Fourth International Workshop on the Synthesis and Simulation of Living Systems*, 1994.
- [3] D. Arnold, A. Bond, M. Chilvers, and R. Taylor, ‘Hector: Distributed Objects in Python’, *Proceedings of the Fourth International Python Conference*, Livermore, California, June 1996.
- [4] B. Christiansen, P. Capello, M. F. Ionescu, M. O. Neary, K. E. Schausser, and D. Wu, ‘Javelin: Internet-Based Parallel Computing Using Java’, *ACM Workshop on Java for Science and Engineering Computation*, 1997.
- [5] C. Catlett and L. Smarr, ‘Metacomputing’, *Communications of the ACM*, Volume 35, Pages 44–52, 1992.

- [6] K. M. Chandy and J. Misra, *Parallel Program Design: A Foundation*, Addison-Wesley, 1988.
- [7] K. M. Chandy, J. Kiniry, A. Rifkin, and D. Zimmerman, ‘Webs of Archived Distributed Computations for Asynchronous Collaboration’, *Journal of Supercomputing*, Volume 11(2):101–118, 1997.
- [8] K. M. Chandy, A. Rifkin, P. A. G. Sivilotti, J. Mandelson, M. Richardson, W. Tanaka, and L. Weisman, ‘A World-Wide Distributed System Using Java and the Internet’, *Proceedings of the Fifth IEEE International Symposium on High Performance Distributed Computing*, Syracuse, New York, August 1996.
- [9] K. M. Chandy and A. Rifkin, ‘Systematic Composition of Objects in Distributed Internet Applications: Processes and Sessions’, *Proceedings of the Thirtieth Hawaii International Conference on System Sciences*, Maui, Hawaii, January 1997.
- [10] I. T. Foster and K. M. Chandy”, ‘Fortran M: A Language for Modular Parallel Programming’. *Journal of Parallel and Distributed Computing*, Volume 26, Number 1, Pages 24–35, April 1995.
- [11] I. Foster and C. Kesselman, ‘Globus: A Metacomputing Infrastructure Toolkit’, *International Journal of Supercomputer Applications*, 1997. To appear.
- [12] G. Fox and W. Furmanski, ‘Towards Web/Java based High Performance Distributed Computing – An Evolving Virtual Machine’, *Proceedings of the Fifth IEEE International Symposium on High Performance Distributed Computing*, Syracuse, New York, August 1996.
- [13] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, *Design Patterns: Elements of Reusable Object-Oriented Software*, Addison-Wesley, 1995.
- [14] J. Gray and A. Reuter, *Transaction Processing: Concepts and Techniques*, Morgan Kaufmann, 1993.
- [15] A. S. Grimshaw, W. A. Wulf, and the Legion team, ‘The Legion Vision of a Worldwide Virtual Computer’, *Communications of the ACM*, Volume 40, Number 1, Pages 39–45, January 1997.

- [16] C. A. R. Hoare, ‘Communicating Sequential Processes’, *Communications of the ACM*, Volume 21, Number 8, Pages 666–677, August 1978.
- [17] The Infospheres Research Group, ‘The Infospheres Infrastructure User’s Guide’, Technical Report, California Institute of Technology, 1997.
- [18] The Infospheres Research Group, ‘Executive Summary’, 1998
<http://www.infospheres.caltech.edu/infospheres.html> Also published as a position paper at the OMG-DARPA-MCC Workshop on Compositional Software Architectures, January, 1998.
- [19] G. E. Krasner S. T. Pope, ‘A cookbook for using model-view-controller user interface paradigm in smalltalk-80’, *Journal of Object-oriented Programming*, August/September 1988. Vol. 1(3), pp. 26-49.
- [20] D. B. Lange and M. Oshima, *Programming Mobile Agents in Java — With the Java Aglet API*, IBM Research, 1997.
- [21] D. Lea, *Concurrent Programming in Java: Design Principles and Patterns*, Addison-Wesley, 1996.
- [22] D. A. Lifka, M. W. Henderson, and K. Rayl, *Users Guide to the Argonne SP Scheduling System*, Argonne National Laboratory Mathematics and Computer Science Division Technical Memorandum, Number ANL/MCS-TM-201, May 1995.
- [23] M. Litzkow, M. Livney, and M. Mutka, ‘Condor – A Hunter of Idle Workstations’, *Proceedings of the Eighth International Conference on Distributed Computing Systems*, Pages 104–111, 1988.
- [24] N. A. Lynch, M. Merritt, W. E. Weihl, and A. Fekete, *Atomic Transactions*, Morgan Kaufmann, 1994.
- [25] N. Minar, R. Burkhart, C. Langton, and M. Askenzai, *The Swarm Simulation System: A Toolkit for Building Multi-Agent Simulations*, Santa Fe Institute, 1996.
- [26] Object Management Group, *The Common Object Request Broker: Architecture and Specification (CORBA)*, revision 2.0, 1995.

- [27] Open Group, *The Distributed Computing Environment's Cell Directory Service*, The Open Group, 1997.
- [28] Open Group, *The Distributed Computing Environment's Global Directory Service*, The Open Group, 1997.
- [29] L. L. Peterson and B. S. Davie, *Computer Networks: A Systems Approach*, Morgan Kaufmann, 1996.
- [30] D. Pinkston, A. Nicholson, and J. Kiniry, *DALI: A Distributed Artificial Life Infrastructure*, Technical Report (to be submitted to *Artificial Life* from MIT Press), California Institute of Technology, 1997.
- [31] R. Ramamoorthi, A. Rifkin, B. Dimitrov, and K. M. Chandy, 'A General Resource Reservation Framework for Scientific Computing', Proceedings of the the first International Scientific Computing in Object-Oriented Parallel Environments (*ISCOPE*) Conference, 1997.
- [32] D. Roberts and R. Johnson, 'Evolving Frameworks: A Pattern Language for Developing Object-Oriented Frameworks', *Proceedings of Pattern Languages of Programs*, Allerton Park, Illinois, September 1996.
- [33] D. C. Schmidt, 'ACE: an Object-Oriented Framework for Developing Distributed Applications', *Proceedings of the Sixth USENIX C++ Technical Conference*, Cambridge, Massachusetts, April 1994.
- [34] D. C. Schmidt, 'A Family of Design Patterns for Application Level Gateways', *Theory and Practice of Object Systems*, Wiley and Sons, Volume 2, Number 1, 1996.
- [35] L. Thomas, S. Suchter, and A. Rifkin, 'Developing Peer-to-Peer Applications on the Internet: the Distributed Editor, SimulEdit', *Dr. Dobbs Journal* 281, Pages 76–81, January 1998.
- [36] R. Sessions, *Object Persistence Beyond Object-Oriented Databases*, Prentice Hall, 1996.
- [37] Workflow Management Coalition, International Organization for the Development and Promotion of Workflow Standards, *Workflow Glossary*, Workflow Management Coalition, Belgium, 1995.