

AD-A232 291

2



RADC-TR-90-101, Vol III (of three)
Final Technical Report
June 1990

**DECENTRALIZED COMPUTING
TECHNOLOGY FOR FAULT-TOLERANT,
SURVIVABLE C3I SYSTEMS
System/Subsystem Specification**

Carnegie-Mellon University

Sponsored by
Strategic Defense Initiative Office



APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED.

The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Strategic Defense Initiative Office or the U.S. Government.

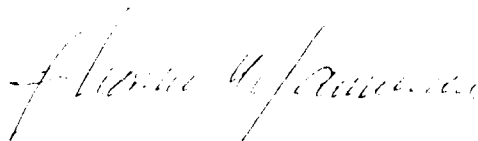
**Rome Air Development Center
Air Force Systems Command
Griffiss Air Force Base, NY 13441-5700**

91 3 01 010

This report has been reviewed by the RADC Public Affairs Office (PA) and is releasable to the National Technical Information Service (NTIS). At NTIS it will be releasable to the general public, including foreign nations.

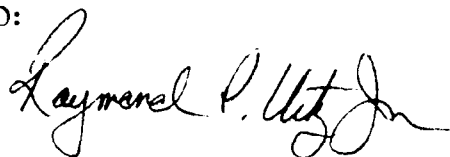
RADC-TR-90-101, Voi III (of three) has been reviewed and is approved for publication.

APPROVED:



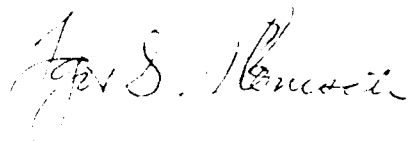
THOMAS F. LAWRENCE
Project Engineer

APPROVED:



Raymond P. Urtz, Jr.
Technical Director
Directorate of Command and Control

FOR THE COMMANDER:



IGOR G. PLONISCH
Directorate of Plans and Programs

If your address has changed or if you wish to be removed from the RADC mailing list, or if the addressee is no longer employed by your organization, please notify RADC (COTD) Griffiss AFB NY 13441-5700. This will assist us in maintaining a current mailing list.

Do not return copies of this report unless contractual obligations or notices on a specific document requires that it be returned.

DECENTRALIZED COMPUTING TECHNOLOGY
FOR FAULT-TOLERANT, SURVIVABLE C³I SYSTEMS
System/Subsystem Specification

J. Duane Northcutt
Edward J. Burke
James G. Hanko
David P. Mavnard
Samuel E. Shipman
Jeffrey E. Trull

E. Douglas Jensen
Raymond K. Clark
Donald C. Lindsay
Franklin D. Reynolds
Jack A. Test

Contractor: Carnegie-Mellon University
Contract Number: F30602-85-C-0274
Effective Date of Contract: 29 Aug 85
Contract Expiration Date: 30 Dec 88
Short Title of Work: Decentralized Computing Tech-
nology for Fault-Tolerant,
Survivable C³I Systems
System/Subsystem Speci-
fication
Period of Work Covered: Aug 85 - Dec 88
Principal Investigator: E. Douglas Jensen
Phone: (508) 393-2989
Project Engineer: Thomas F. Lawrence
Phone: (315) 330-2158

Approved for public release; distribution unlimited.

This research was supported by the Strategic Defense Initia-
tive Office of the Department of Defense and was monitored
by Thomas F. Lawrence (COTD), Griffiss AFB NY 13441-5700,
under contract F30602-85-C-0274.

REPORT DOCUMENTATION PAGE

Form Approved
OPM No. 0704-0188

Public reporting burden for this collection of information is estimated to average 1 hour per response including the time for reviewing the instructions, searching existing data sources, gathering and maintaining the data needed, and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Information and Regulatory Affairs, Office of Management and Budget, Washington, DC 20503

1. AGENCY USE ONLY (Leave Blank)		2. REPORT DATE June 1990	3. REPORT TYPE AND DATES COVERED Final Aug 85 - Dec 88	
4. TITLE AND SUBTITLE DECENTRALIZED COMPUTING TECHNOLOGY FOR FAULT-TOLERANT, SURVIVABLE C ³ I SYSTEMS System/ Subsystem Specification			5. FUNDING NUMBERS C - F 30602-85-C-0274 PE - 63223C PR - 2300 TA - 02 WU - 10	
6. AUTHOR(S) J. Duane Northcutt, E. Douglas Jensen, Edward J. Burke, Raymond K. Clark, James G. Hanko, Donald C. Lindsay, David P. Maynard, (Cont'd)				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Carnegie-Mellon University Pittsburgh PA 15213-3890			8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) Strategic Defense Initiative Office, Office of the Secretary of Defense Wash DC 20301-7100			10. SPONSORING/MONITORING AGENCY REPORT NUMBER RADC-TR-90-101, Vol III (of three)	
11. SUPPLEMENTARY NOTES RADC Project Engineer: Thomas F. Lawrence/COTD/(315) 330-2158				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited.			12b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words) Alpha is an operating system for the mission-critical integration and operation of large, complex, distributed, real-time systems. Such systems are becoming increasingly common in both military (e.g., BM/C ³ , combat platform management) and industrial factory and plant automation (e.g., automobile manufacturing) contexts. They differ substantially from the better-known timesharing systems, numerically-oriented supercomputers, and networks of personal workstations. More surprisingly, they also depart significantly from traditional real-time systems, which are predominately for low-level periodic sampled data monitoring and control.				
14. SUBJECT TERMS Real-Time System Distributed Operating System Distributed Computing			15. NUMBER OF PAGES 196	
			16. PRICE CODE	
17. SECURITY CLASSIFICATION OF REPORT UNCLASSIFIED		18. SECURITY CLASSIFICATION OF THIS PAGE UNCLASSIFIED	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASSIFIED	20. LIMITATION OF ABSTRACT SAR

6 (Cont'd). Franklin D. Reynolds, Samuel E. Shipman, Jack A. Test, Jeffrey E. Trull

Preface

Alpha is an adaptable decentralized operating system for real-time applications, being developed as a part of the Archons project's on-going research into real-time distributed systems. Alpha is the first systems effort of the Archons Project, and an initial version (i.e., Release 1) has been created at Carnegie-Mellon University directly on Sun workstation hardware. It has been demonstrated with a real-time control application written by its first industrial user, General Dynamics. A second version of Alpha (i.e., Release 2) is being produced at Concurrent Computer Corporation. Both versions of Alpha are sponsored by the USAF Rome Air Development Center and are in the public domain for U.S. Government use.

This report consists of four main parts—the first two of which are concerned with Release 1 of Alpha, and the final two parts deal with Alpha Release 2. The first part of this report provides a description of the programmer's interface to the Release 1 kernel of the Alpha operating system. Included in this description is an enumeration of all of the points of entry into the kernel, and the units of kernel functionality provided to the kernel's clients in the form of objects.

The second part of this report describes the design and implementation of the kernel layer of Release 1 of Alpha. The kernel is the lowest layer of functionality in Alpha and it supports the fundamental abstractions of the Alpha programming model. This description provides details on the design and implementation of the major facilities provided by the kernel of Release 1 Alpha.

The third part of this report describes the application-visible interface to Release 2 of the Alpha kernel. This interface includes the operations defined on the kernel objects, kernel supported properties of application objects, object invocations, capabilities, threads, and exceptions. This description does not explicitly deal with details of the design, programming model or behavior of Alpha (all of which are described in other documents). The interface described here is specified in a language- and architecture-independent fashion.

The fourth and final part of this report describes the function and design of the various subsystems within Release 2 of the Alpha kernel. This description includes a discussion of the design and current implementation of the Alpha kernel.



Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By.....	
Distribution/	
Availability Codes	
Dist	Avail and/or Special
A-1	

Contents

Part A: Alpha Release 1 Kernel Interface

Part B: Alpha Release 1 Kernel Internals

Part C: Interim Alpha Release 2 Kernel Interface Specification

Part D: Interim Alpha Release 2 System/Subsystem Specification

Table of Contents

Abstract	A-1
1 Introduction	A-2
2 Kernel Entry Points	A-3
2.1 Traps	A-3
2.1.1 Operation Invocation	A-3
2.1.1.1 Client Invocation.....	A-5
2.1.1.2 Kernel Invocation.....	A-5
2.1.2 Exception Blocks	A-6
2.1.3 Event Logging.....	A-7
2.2 Interrupts.....	A-7
2.2.1 On-Board DUSART.....	A-8
2.2.2 On-Board Timer.....	A-8
2.2.3 Inter-Processor Communications.....	A-8
2.2.4 Spurious Interrupts.....	A-8
3 System Service Objects	A-9
3.1 Manager Objects	A-9
3.1.1 Object Manager.....	A-9
3.1.2 Thread Manager.....	A-11
3.1.3 Semaphore Manager	A-12
3.1.4 Lock Manager	A-13
3.2 Support Objects.....	A-14
3.2.1 Semaphore.....	A-14
3.2.2 Lock	A-15
3.2.3 Time Services.....	A-17
3.3 I/O Objects	A-18
3.3.1 TTY.....	A-18
3.3.2 Network I/O	A-19
3.4 Miscellaneous Object.....	A-20
References	A-22
Appendix I: Kernel Interface Usage Examples	A-23

List of Figures

Figure 1	Invocation Parameter Block Formats	A-4
Figure 2	Object Manager Object.....	A-10
Figure 3	Thread Manager Object	A-11
Figure 4	Semaphore Manager Object	A-12
Figure 5	Lock Manager Object	A-13
Figure 6	Semaphore Object.....	A-14
Figure 7	Lock Object	A-16
Figure 8	Time Services Object.....	A-17
Figure 9	TTY Object.....	A-18
Figure 10	Network I/O Object	A-19
Figure 11	Temporary Services Object	A-20

List of Figures

Table 1	Lock Compatibility Table.....	A-16
----------------	-------------------------------	------

Abstract

This report describes the programmer's interface to the Release 1 kernel of Alpha. Included in this description is an enumeration of all of the points of entry into the kernel (i.e., software generated traps and hardware generated interrupts). Also described in this document are the units of kernel functionality provided to the kernel's clients in the form of objects (known as *system service objects*).

In the design of the Alpha kernel an attempt was made to make the kernel interface be simple and uniform. This document provides an illustration that these goals have been accomplished. There are a minimum number of entry points into the kernel, and most services are provided by a small number of objects (each of which has an equally small number of operations).

Introduction

This document describes the functional interface provided by the Release 1 kernel of Alpha. A major objective in the functional design of the kernel was to provide all system services in a simple and uniform manner. This led to the decision to provide all kernel functions in such a way as to appear to the client as though they are provided by a collection of kernel-defined objects. The functions that are traditionally accessed via traps in more typical operating systems (commonly, one trap for each service), are accessed in Alpha by the single trap that performs the operation invocation function.

A simplified description of the Alpha kernel interface includes only the operation invocation mechanism, but a more complete description includes the definition of all of the services provided by the kernel in the form of objects, and all of the other points of kernel entry.

In this report an enumeration of the entry points into the kernel is given, followed by a list of the services provided by the kernel in the form of objects. The kernel entry mechanisms fall into two categories—traps and interrupts—each of which provide a separate means of entry into the kernel. The kernel services that are provided in the form of objects, known as *system service objects* are accessed through the operation invocation trap.

The bulk of the kernel's interface is involved in providing support for the basic abstractions—i.e., objects, operation invocation, and threads. The operation invocation mechanism is provided to the client as a programming language construct [Shipman 88], and the kernel support for objects and threads is provided by system service objects.

This report represents the actual implemented interface for Release 1 of the Alpha kernel. A description of the intended functional behavior of the Alpha operating system is provided in [Northcutt 88a], while [Northcutt 88b] contains a description of the programming model that the Alpha kernel is intended to support, and [Shipman 88] gives a description of the programming interface that is provided on top of the kernel interface described in this document.

In the following chapters the kernel entry points (traps and interrupts) and system service objects are defined. In an appendix, an example of the manner in which the Alpha object-language preprocessor makes use of the kernel interface mechanisms is presented.

2 Kernel Entry Points

The only means of entering the Alpha kernel is either by software trap or hardware interrupt. Each of these entry points is described in this document—first the system traps are discussed, followed by a discussion of the system interrupts.

The primary entry point into the kernel (and abstractly, the only one) is via the operation invocation facility. Fundamentally, the Alpha kernel interface consists of the object invocation mechanism and a collection of kernel-defined objects. The operation invocation mechanism is made available by the kernel as the only *system call* (implemented by way of a trap instruction), while the kernel-provided objects are kernel routines that present the client with an object interface (and, as part of the kernel's code, are available as soon as the kernel starts up on a node).

In Release 1 of Alpha, all machine exceptions (e.g., divide by zero and unimplemented instructions) are handled by the monitor, and so are not described as part of the kernel interface.

2.1 Traps

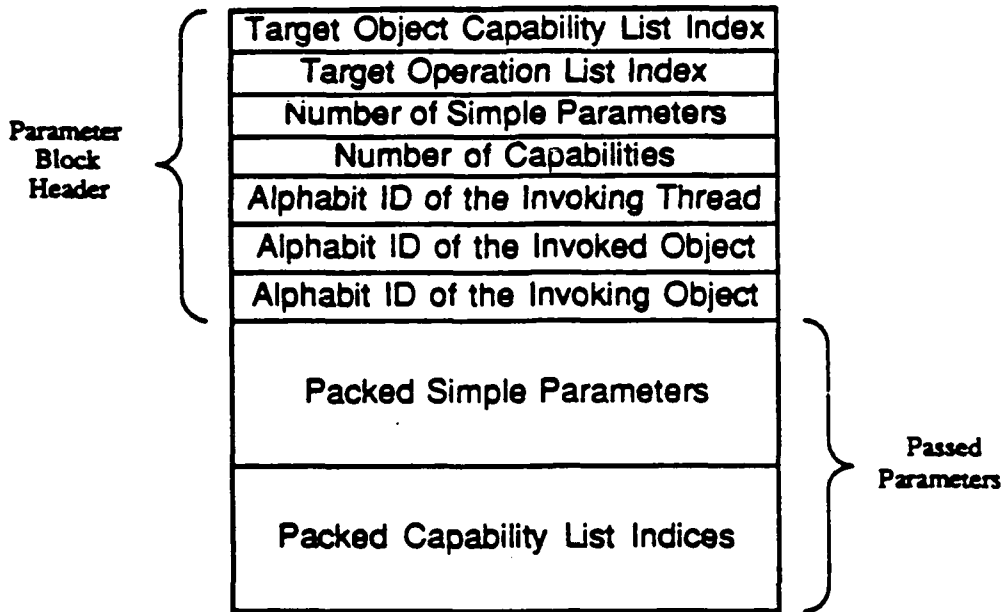
In addition to the primary entry point provided by the operation invocation facility, two other trap entry points are provided in Release 1. These are for exception block management and event logging, and both are intended to be accessed through the operation invocation facility, but were implemented as separate traps for development purposes and will be integrated with the operation invocation facility in a future release.

2.1.1 Operation Invocation

In Alpha, operation invocations are similar to traditional remote procedure calls (RPCs), there is an entry into the kernel and then a return from the kernel to the invoker when the operation is complete. Each thread has a collection of parameters that are brought into a target object by a thread on the invocation of an operation and a set that are returned from the object when the operation completes. Collectively, these parameters are known as the “incoming” parameters of a thread within an object. In addition, a thread that invokes an operation on an object uses a collection of parameters that contain the values sent to, and subsequently returned from, a target object. These parameters are known as “outgoing” parameters. The parameters that pass into an object on an operation invocation are known as *request* parameters, and those that are passed back from the target object when the invocation is complete are known as *reply* parameters. The detailed format of both the parameter blocks used in Release 1 of Alpha is shown in Figure 1, with the request parameters shown in part 1a, and reply parameters in part 1b.

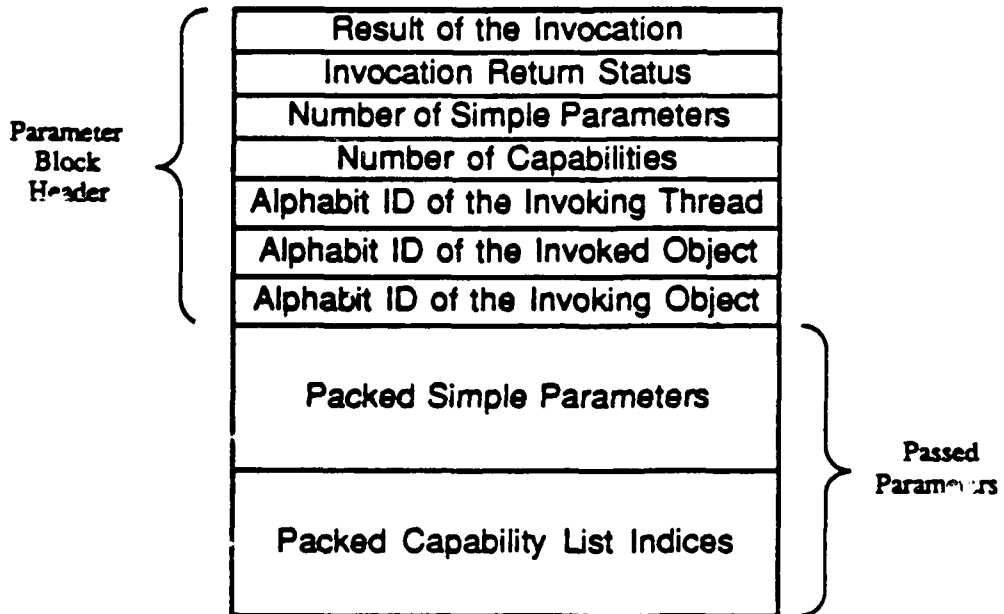
All parameters are packed into the contiguous memory following the parameter block header. The header portion of a parameter block contains all of the control information needed from the client for an invocation, and the parameter portion contains both simple parameters (i.e., arbitrary C data types), and capabilities (i.e., indexes into the invoking object's c-list). The invoker passes capabilities in the form of indexes into the invoking object's c-list. At the beginning of an invocation the kernel validates all passed capability indexes and translates them into object descriptors. Should any of these tests fail, the invocation is not performed and an error indication is returned in the reply parameter block.

Request Parameter Block



a)

Reply Parameter Block



b)

Figure 1: Invocation Parameter Block Formats

When an invoked operation is complete, capability indexes that are to be returned as reply parameters to the invoker are similarly validated and translated by the kernel, and any error indications are likewise returned with the invocation[†].

The logical presentation of the operation invocation facility to the client involves the use of a set traps, and a collection of parameter passing conventions. There are two different system locations from which operation invocations can be performed—i.e., from client objects, and from kernel objects. While operation invocations from each of these points is made to appear the same to the user, for reasons of simplicity and efficiency, a different trap entry point is used for each of these types of invocations. The following discussion describes the behavior of the operation invocation facility in each of these cases.

2.1.1.1 Client Invocation

The most common form of operation invocation emanates from client objects and is known as *client invocation*. This type of invocation involves the movement of parameters between object address spaces. To accomplish this efficiently, the invocation parameters are placed in separate pages so that the memory management hardware can be used to map the parameters among objects (as opposed to copying)[†].

To perform a client invocation, the appropriate values are loaded into the request part of the outgoing parameter block page, and a trap instruction (i.e., trap #10) is executed. On client invocations, the location of the parameter page is known implicitly by the kernel, which accesses the appropriate parameter pages as it carries out invocation.

In a fashion very similar to that of client invocations, the kernel must be entered upon the completion of an operation invoked on a client object. This requires that the target (client) object marshal return parameters and execute a trap instruction (i.e., trap #11) to return to the kernel on its way back to the invoking object. An activity very similar to that performed on the request side of the invocation is carried out as a part of this trap.

When the invoked operation is complete, return values are placed into the reply side of the invoking thread's outgoing parameter block, any passed capabilities are validated, inserted into the invoking object's c-list, and then translated into c-list indices. When this has been done, the operation is allowed to complete and the invoking thread resumes execution within the object from which the invocation was made.

2.1.1.2 Kernel Invocation

The invocation of operation on objects from within a kernel object is performed via a *kernel invocation*. From the client programmer's point of view this form of invocation provides a function identical to that of client invocations, except in that the parameter blocks do not exist in a well-known location in virtual address space, and the invoking object is already within the kernel's address space. Thus, one difference between this form of invocation and client invocation is that pointers to the invoking thread's outgoing, request and reply parameter blocks are pushed onto the kernel stack prior to the execution of the invocation trap (i.e., trap #9).

[†]This is an area where the functionality has been modified in Release 2.

The only other difference is that there is no need to trap back to the kernel when an operation invoked on a kernel object completes. In such a case, the invoking object is already in the kernel's address space, so an RTE from the invoking trap is all that is needed to return from an invocation on a kernel object.

2.1.2 Exception Blocks

The Alpha programming model includes a notion of *exception blocks* which are defined as contiguous regions of code within an object that have special regions of code (known as *exception handlers*) associated with them, that are executed when exceptions occur while threads are executing within the exception blocks. The kernel provides mechanisms to permit the programmer to define the beginning and end of exception blocks and to define the associated exception handlers for these blocks.

While these exception block mechanisms would best be provided by invocations on the thread executing within an exception block, the current implementation provides these functions by an independent trap entry point. The exception block trap (i.e., trap #12) provides the mechanisms that define the start and end of an exception block, as well as the beginning and end of each block's exception handler.

The command types implemented in Release 1 are: *begin/end* a generic exception block, *begin/end* a time constraint block, *begin* the system-defined exception handler, and *end* an exception handler. The exception block *begin* and *end* commands can be executed at any point within an object, the intent is that all *begin* commands have matching *end* commands within the same object. This, however, is not a requirement—if more *begin* than *end* operations are executed within an object, the kernel automatically (effectively) performs the necessary number of *end* operations when an operation completes. In addition, exception blocks (of all varieties) can be nested within an object.

At the point in an object where an exception block command is to be performed, an identifier for the type of operation to be performed is pushed onto the stack (along with any additional parameters needed by the given operation) and the trap instruction is executed. The same trap instruction is used to perform exception block operations from both kernel and client objects.

When a *begin* operation is executed, the kernel saves the state of the executing thread (e.g., stack frame pointer and processor register set) and returns a false condition. When an *end* command is executed, the most recently saved thread state is discarded. Exception blocks manage thread state in a stack-oriented fashion, this is to say that *begin* operations pushes the thread's state and *end* operations pop the saved thread state. Should an exception occur while the thread is executing within an exception block, the kernel returns the thread to each of the previously saved locations in turn, restores the thread's state to what it was when it executed each block's *begin* command, and returns a true condition. With these mechanisms the kernel's clients can define a variety of different exception blocks and customized exception handling code.

The two types of exception handling blocks defined in Release 1 are known as *generic* and *time constraint* exception blocks. Generic exception blocks allow exception handling code to be associated with arbitrary regions of code within an object. Time constraint exception blocks are used to define time constraints on threads, and unlike all of the other commands which only take the command type parameter, time constraint *begin* commands

include thread time constraint parameters. Time constraint parameters are specific to the scheduling algorithm used and are pushed onto the stack following the exception block type code. The exact parameters required for time constraint blocks in Release 1 are defined in [Shipman 88].

2.1.3 Event Logging

An event logging facility was added to Alpha for both system and application program performance monitoring. In Alpha, events are programmer-defined points in a program's execution stream. The kernel logs all (enabled) events in a buffer area that is set aside for this purpose within the kernel's address space.

The kernel's debugging facility permits the interactive examination of event logs, and the event logs can also be uploaded to a development host (via the monitor-provided TFTP facility) and examined with a post-processing tool (described in [Northcutt 88d]).

The kernel's event logging mechanism requires that a parameter indicating the event type (that is used as a index into a event print-string table), and an arbitrary user-defined parameter (that is used to further disambiguate the event). These parameters are pushed onto the stack and then the event trap is executed (i.e., trap #14), following which, the kernel adds a timestamp and logs the event into the kernel's event buffer.

This trap operation can be used from within both kernel and client objects, and it returns promptly following the timestamping and logging of the event.

2.2 Interrupts

In addition to traps, interrupts are used to gain entry to the kernel. The application processor (the processing element within a node on which the kernel proper executes) can receive interrupts from a variety of sources, either from on-board devices or from devices on the node's Multibus.

The general philosophy for the handling of interrupts in Alpha involves their conversion into activities performed by threads, in as prompt a fashion as is possible. This suggests that very few instructions are executed under interrupt—typically all that is performed by an interrupt handler is the initiation of a thread. The reason for this approach is based on the belief that the scheduler should manage all application processing cycles, and traditional interrupt processing subverts this process by asynchronously diverting processing cycles to activities (under the direction of the interrupt handling hardware and not the operating system's scheduler). By converting the activities typically associated with interrupt handlers into a form that the kernel's scheduler can manage, allows all processing cycles according to a globally consistent, time-driven application processor cycle management policy.

The application processor recognizes seven levels of hardware-vectored interrupts, which (in the Alpha testbed) are assigned as follows:

- Level 7: Non-Maskable Interrupt (this interrupt is handled by the monitor)
- Level 6: On-Board Dual USART
- Level 5: On-Board Timer
- Level 4: Inter-processor Communications
- Levels 3-1: Multibus Interrupts

Each of the kernel's interrupt handling routines and their sources are described in the following discussion.

2.2.1 On-Board DUSART

The on-board DUSART is used to provide basic serial communications for each node. The application processors use the A port as the node's console connection, and the B port is used to communicate with the other processors in a node. The monitor uses port A as its (default) console line, and (once loaded and running) the kernel also uses this port for console communications.

The kernel provides a *TTY* object that encapsulates the serial communications unit and provides basic I/O functions to the Alpha system and its applications. The *TTY* object takes characters in as parameters and sends them across the serial line, and also accepts characters from the serial line and unblocks threads that are waiting for incoming characters. In addition, the object provides functionality equivalent to the UNIX raw and cooked I/O modes.

2.2.2 On-Board Timer

The application processor uses its on-board, programmable, counter/timer unit to generate a constantly running, "real-time" system clock for the kernel. Because the timer runs constantly and is of finite length, the counter may "roll-over" during the lifetime of the system. Therefore, the timer hardware is configured to generate interrupts on counter overflows, so that the interrupt handler can increment the software portion of the system clock value when an overflow occurs.

2.2.3 Inter-Processor Communications

The processors within an Alpha testbed node exchange messages via shared memory, and interrupts are used to notify the destination processor of the arrival of a message. Such messages are sent along uni-directional logical links (known as *channels*) among the processors in a node. Inter-processor messages are buffered within channels until they are removed by the destination subsystem.

The kernel executing on the application processor receives interrupts from the various kernel subsystem processors to indicate message arrivals. The application processor's inter-processor communication interrupt handler polls all of the incoming channels, removing and acting on each message it finds, until all channels are empty. The only action that is typically carried out in response to an incoming message is the unblocking of a thread.

2.2.4 Spurious Interrupts

All other interrupts are not used by the Release 1 kernel, and therefore the kernel always runs at minimum priority of three, so as to block out all other interrupts (i.e., Multibus interrupts such as Ethernet controllers, disk controllers, and interrupts directed from the applications processor to the co-processors).

Just to be sure, the kernel installs a spurious interrupt handler for all unused interrupts. This handler does nothing but print a console message and return.

3 System Service Objects

The bulk of the kernel's functionality is provided through services accessed through the operation invocation facility. These services are organized into groups and made to appear as objects, known as "system service objects." These pseudo-objects provide support for the Alpha basic abstractions, as well as a range of supporting facilities.

The following sections describe the interface and function of each of the system service objects that the kernel provides in support of the system's programming model.

3.1 Manager Objects

The object programming paradigm used in Alpha suggests that all system activities be performed by the invocation of operations. For example, all manipulation of objects should be done by invoking operations on the object itself. However, there are a class of operations that do not lend themselves well to this uniform, invocation-based approach—e.g., it is not possible to invoke an operation on an object that does not yet exist, and so this approach cannot be used for object creation. For this reason, the Alpha kernel provides a collection of special *manager* objects to carry out the operations that are not naturally performed by invoking operations on an entity directly.

Once an entity is created, it is possible to perform (almost) all subsequent operations on it directly. Therefore, while it was not strictly necessary to do so, the DELETE operation has been included with the CREATE operations provided by manager objects (for reasons of symmetry).

The manager system service objects for objects, threads, semaphores, and locks are described in the following subsections.

3.1.1 Object Manager

In support of the Alpha object abstraction, the kernel provides a system service object for the dynamic management of object instances. This is known as the *ObjectManager* object, and has operations defined on it for the run-time creation and deletion of object instances, the registration of wellknown capabilities, and the creation of replicated objects.

The replication-oriented operations are a temporary feature in Release 1 of Alpha, that were added in support of some initial application work. These operations do not constitute a proper set of mechanisms, but rather form a complete (if somewhat limited) replication facility. This replication facility includes specific object placement and replica management policies (i.e., a maximum of one replica per node, simple replica regeneration algorithm, and a primary plus backup copy scheme). A different set of operations will be provided in subsequent releases of Alpha; the object placement function belongs above the kernel level, and an additional, higher-level facility is required for the complete management replicated objects.

The interface to the *ObjectManager* object is shown in Figure 2, and its operations are described in the following text.

```

SYSTEM OBJECT  ObjectManager {
    OPERATION  Create(IN CAPA objectType,
                    OUT CAPA instance)
    OPERATION  Delete(IN CAPA instance)
    OPERATION  RegisterWellknown(IN CAPA regularCapa,
                                IN CAPA wellknownCapa)
    OPERATION  ReplicatedCreate(IN CAPA objectType,
                                IN u_long desiredReplicationFactor,
                                OUT CAPA instance,
                                OUT u_long actualReplicationFactor)
    OPERATION  Anchor()
}

```

Figure 2: Object Manager Object

Operation

Create(IN CAPA objectType,
OUT CAPA instance)

This operation is used to create a new instance of an object type from a given object type specification. The 'objectType' parameter is a capability for the object type specification for the type of object that is to be created by this operation, and 'instance' is a capability for the newly created object instance, and is returned when the operation has completed.

The object instance created by this operation is placed local to the node where the operation is performed. The full complement of rights are associated with the capability returned for the new object. This operation will fail if the specified object type is unknown, or if the node currently lacks the resources necessary to support another object.

Operation

Delete(IN CAPA instance)

This operation is used to delete an existing object instance. The 'instance' parameter is a capability for the object to be deleted.

If the passed capability is valid, this operation deletes the specified object and deallocates the resources associated with it. The indicated object is deleted regardless of its physical location in the system.

Operation

RegisterWellknown(IN CAPA regularCapa,
IN CAPA wellknownCapa)

This operation creates an alias for a given capability. The 'regularCapa' parameter is a capability that exists in the local dictionary, while the 'wellknownCapa' parameter is the wellknown capability that is to be bound to the entity indicated by the previous parameter.

Operation

```

ReplicatedCreate(IN CAPA objectType,
                 IN u_long desiredReplicationFactor,
                 OUT CAPA instance,
                 OUT u_long actualReplicationFactor)

```

This operation performs a function similar to that of the CREATE operation defined above, however it creates multiple instances of the new object. The 'objectType' parameter is a capability for the object type specification, and the 'instance' parameter is the returned capability for the newly instantiated (replicated) object. Furthermore, the 'desiredReplicationFactor' parameter is an integer that defines the (maximum) number of replicas that should be created, and the 'actualReplicationFactor' parameter is returned to indicate the actual number of replicas that were created.

This operation creates a replicated object instance, i.e., an object that consists of a number of replicas and a primary copy. The returned capability refers to the primary copy and does not change across the life of the replicated object (even when another replica takes over as the primary copy).

There are several reasons that the number of replicas created may differ from the number requested—e.g., there may not be sufficient resources for the requested number of replicas, and the replication policy dictates that there can be at most one replica per node.

Operation

```

Anchor()

```

This operation causes a thread to block on a semaphore (and otherwise do nothing). This operation is used in extending "lifeline" threads between objects to detect node/communications failures. When an invocation of this operation returns, the node on which this object exists can be considered to have failed.

3.1.2 Thread Manager

In support of threads, the kernel provides an object for the dynamic management of threads. This object is known as the *ThreadManager* object, and it provides functionality for threads that is analogous to what the previously defined *ObjectManager* object provides for objects.

The interface to the *ThreadManager* object is shown in Figure 3, and its operations are described in the following text.

```

SYSTEM OBJECT ThreadManager {
    OPERATION Create( IN CAPA instance,
                    IN u_long initialOperation,
                    IN pamblk initialParameters,
                    OUT CAPA thread)
    OPERATION Delete(IN CAPA thread)
}

```

Figure 3: Thread Manager Object

Operation

```

Create(IN CAPA instance,
      IN u_long initialOperation,
      IN parmblk initialParameters,
      OUT CAPA thread)

```

This operation is used to create threads. The parameters passed to this operation are a capability for the object (i.e., 'instance'), the name of the operation (i.e., 'initialOperation') that the newly created thread should begin execution within, and the initial parameters to be passed into the initial object. This operation returns a capability for the newly created thread (i.e., 'thread'), that contains all of the rights associated with threads. If the specified object does not currently exist on the node at which this operation is invoked, the operation fails.

Operation

```

Delete(IN CAPA thread)

```

This operation is used to delete a specified thread. A capability for a thread is given as the 'thread' parameter. Provided that the given capability is valid, this operation deletes the specified thread and deallocates the resources that are currently associated with it. Furthermore, the thread specified by the given capability is deleted regardless of its current location in the system. In addition, all of the thread's outstanding exception blocks are executed when a thread is deleted.

3.1.3 Semaphore Manager

The kernel provides the *SemaphoreManager* object in order to permit the dynamic creation and deletion of individual instances of *Semaphore* objects. In Release 1 of Alpha, *Semaphore* objects are bound to their creating object—i.e., each newly created *Semaphore* object can only be used by the object that created it, and likewise a *Semaphore* object may only be deleted by its creator.

The interface to the *SemaphoreManager* object is shown in Figure 4, and its operations are described in the following text.

```

SYSTEM OBJECT SemaphoreManager (
  OPERATION AllocSem(IN u_long count,
                    IN boolean pvUndo,
                    IN boolean allocUndo,
                    OUT CAPA Semaphore semaphore)
  OPERATION DeallocSem(IN CAPA Semaphore semaphore)
)

```

Figure 4: Semaphore Manager Object

Operation

```

AllocSem(IN u_long count,
        IN boolean pvUndo,
        IN boolean allocUndo,
        OUT CAPA Semaphore semaphore)

```

This operation instantiates an object of the kernel-defined type *Semaphore*. This operation takes in a parameter to indicate the initial value that the semaphore's counter should take on (i.e., the 'count' parameter). It also takes a pair of boolean parameters that indicate whether the outstanding P and V operations should be undone on exceptions, and whether the allocation should be undone on exceptions (i.e., 'pvUndo' and 'allocUndo', respectively). These parameters are needed in order to permit the system's default exception handling code to restore an object's semaphores to their correct state should an exception occur while a thread is executing within it.

Furthermore, this operation returns a capability for the newly created *Semaphore* object to the invoking object in the 'semaphore' parameter. This capability has the NO_TRANSFER and NO_COPY restrictions applied to it (this is to ensure that semaphores and the objects that use them remain physically co-located with respect to each other).

Operation

DeallocSem(IN CAPA Semaphore semaphore)

This operation deletes the *Semaphore* object whose capability is passed as the 'semaphore' parameter to this operation. Because the NO_TRANSFER and NO_COPY restrictions are associated with the capabilities given to the *Semaphore* object's creator, only the object that created a semaphore has the ability to perform this operation on it.

3.1.4 Lock Manager

The *LockManager* object is provided by the kernel as a means of creating and deleting individual instances of *Lock* objects at run-time. In Release 1 of Alpha, locks are advisory—i.e., the underlying hardware does not enforce the desired synchronization regime. rather programmers must adhere to the proper conventions in order for locks to be effective. Because locks in Release 1 of Alpha are associated with contiguous regions of memory within an object, *Lock* objects, like *Semaphore* objects, remain local to the object that created them. Furthermore, the data regions associated with allocated locks are not permitted to overlap, therefore the *LockManager* returns a failure indication when an attempt is made to allocate a lock that intersects with an existing lock's data region.

The interface to the *LockManager* object is shown in Figure 5, and its operations are described in the following text.

```

SYSTEM OBJECT LockManager {
  OPERATION AllocLock(IN vaddr start,
                      IN u_long size,
                      OUT CAPA Lock lock)
  OPERATION DeallocLock(IN CAPA Lock lock)
}

```

Figure 5: Lock Manager Object

Operation

AllocLock(IN vaddr start,
 IN u_long size,
 OUT CAPA Lock lock)

This operation is used to instantiate and initialize a *Lock* object that is associated with a data region within the object from which the operation is invoked. The data region that the lock is to cover is defined by the starting virtual address given by the 'start' parameter, and a byte count given by the 'size' parameter. A capability for the newly allocated *Lock* object is returned in the 'lock' parameter. As with *Semaphore* objects, the returned capability to the *Lock* object has the `NO_TRANSFER` and `NO_COPY` restrictions applied to it (ensuring that a lock and the object it is bound to remain physically co-located with respect to each other).

Operation

DeallocLock(IN CAPA Lock lock)

This operation is used to delete a lock and return to the system any resources associated with the *Lock* object indicated by the 'lock' parameter. Any threads that are blocked waiting to gain access to the deleted lock are unblocked.

.2 Support Objects

The kernel provides other system service objects to support the Alpha kernel's abstractions. Included among these are objects of the system-defined object types for thread synchronization—i.e., *Semaphore* and *Lock* objects. Instances of these types of objects are generated by their respective manager objects (described in the preceding section). In addition, the kernel supplies a system service object that provides real-time clock services to the client. All of these objects are described in the following subsections.

.2.1 Semaphore

Semaphore objects are instances of a kernel-defined object type, that are created and deleted by the *SemaphoreManager* object. Instances of *Semaphore* objects are used to control the concurrent execution of threads within objects. A *Semaphore* object provides functionality equivalent to the counting semaphore constructs found in many other systems.

Concurrency control is achieved by blocking and unblocking threads with these semaphores. In addition to providing the basic P and V operations, the *Semaphore* object has an operation defined on it that permits all of the threads that are blocked on the semaphore to be released. The latter operation is useful in a number of cases, including exception handling. The interface to the *Semaphore* object is shown in Figure 6, and its operations are described in the following text.

```

SYSTEM OBJECT Semaphore (
  OPERATION PSem()
  OPERATION VSem()
  OPERATION VAllSem()
)

```

Figure 6: Semaphore Object

Operation

PSem()

This operation performs a standard 'P' operation on a counting semaphore. Logically, this operation can be considered to represent an attempt at obtaining a token that corresponds to a resource that the *Semaphore* object manages. If a token is not available, the thread that invokes this operation is blocked until one becomes available. Once a token has been granted and the resource used, the thread then invokes a 'V' operation to, logically, return the synchronization token.

Operation

VSem()

This operation performs the standard 'V' operation on a counting semaphore. Logically, this operation represents the return of a previously granted token to the semaphore's resource pool. If there are any threads blocked waiting on a token, one of the threads that are blocked on the semaphore is granted the token and unblocked. The policy for selecting which of a set of blocked threads to release is definable at system build time. The policy currently in use is compatible with the best-effort scheduling policy used in Alpha.

Operation

VAllSem()

This operation unblocks all threads which may be waiting on the *Semaphore* object, and sets the semaphore's count to zero. If there are no threads queued on the semaphore when this operation is invoked, this operation has no effect.

3.2.2 Lock

Lock objects are instances of a kernel-defined object type that are created and deleted by the *LockManager* object. Instances of *Lock* objects are used to control the concurrent execution of threads within objects. Whereas *Semaphore* objects control the concurrent execution of threads within regions of code, the *Lock* objects control the concurrent access of threads to regions of memory within an object.

Lock objects are associated with specific data regions in an object, and threads executing within an object indicate their intentions to access the data covered by locks through the invocation of LOCK operations on the them. When a thread's manipulation of a locked data item is complete, the UNLOCK operation is invoked to indicate this fact to the kernel.

On LOCK operations, the permission to proceed with the desired operation on the region of data associated with the lock is determined by the system's locking policy (as manifest in the kernel's *lock compatibility table*). Threads are blocked until the lock can be granted in the requested mode. The locking policy currently in use by Release 1 of Alpha has the following modes: *exclusive read*, *concurrent read*, *exclusive write*, *concurrent write*, and *exclusive read/write*. The lock compatibility table for this policy is shown in Table 1. Furthermore, the unblocking of multiple queued, compatible, requests is performed in accordance with the system's time-driven resource management policy.

Requested \ Granted	Concurrent Read	Concurrent Write	Exclusive Read	Exclusive Write	Exclusive Read/Write
Concurrent Read	Compatible	Compatible	Not Compatible	Compatible	Not Compatible
Concurrent Write	Compatible	Compatible	Not Compatible	Not Compatible	Not Compatible
Exclusive Read	Not Compatible	Compatible	Not Compatible	Compatible	Not Compatible
Exclusive Write	Compatible	Not Compatible	Not Compatible	Not Compatible	Not Compatible
Exclusive Read/Write	Not Compatible	Not Compatible	Not Compatible	Not Compatible	Not Compatible



 — Compatible
  — Not Compatible

Table 1: Lock Compatibility Table

When a thread completes execution within an object, the system automatically releases all of the locks that it holds (just as is done with semaphore's). In addition, the default exception handler restores the state of all locks held in *write* modes by a thread that encounters an exception, it also releases all locks acquired by the thread within the exception block.

The interface to the *Lock* object is shown in Figure 7, and its operations are described in the following text.

```

SYSTEM OBJECT Lock {
  OPERATION Lock(IN u_long mode)
  OPERATION ConditionalLock(IN u_long mode,
                           OUT boolean result)
  OPERATION UnLock()
  OPERATION Convert(IN u_long newMode)
  OPERATION Query(OUT u_long currentMode)
}

```

Figure 7: Lock Object

Operation

Lock(IN u_long mode)

This operation is used to acquire a lock—i.e., to announce to the kernel an intention on the part of the invoking thread to access the region of data associated with the lock, in a manner indicated by the 'mode' parameter. If the locking policy that is in place in the system indicates that the requested mode is compatible with the mode that the lock is currently in, then the operation returns and the thread can proceed to manipulate the locked data in the indicated manner. However, if the operation is determined to be incompatible

with the current state of the lock, the invoking thread is blocked until a time when the request becomes compatible.

Operation

```
ConditionalLock(IN u_long mode,
                OUT boolean result)
```

This operation is identical in function to the previously described operation, with the exception that the thread that invokes this operation is never blocked. Rather, a boolean flag is returned in the 'result' parameter, a logical TRUE indicates that the lock has been successfully acquired, and FALSE indicates that the lock could not be acquired in the requested mode.

Operation

```
UnLock()
```

This operation is used to release a previously acquired lock. If the invoking thread does not hold the lock (in any mode) when this operation is invoked, the operation has no effect.

Operation

```
Convert(IN u_long newMode)
```

This operation allows a thread to change the mode in which a lock is held to the mode indicated by the 'newMode' parameter. This operation behaves much the same as the LOCK operation. If the lock is already being held by the thread, this operation does not require that the lock be released first and then reacquired in the new mode, instead it makes these two operations appear as though they were performed atomically. If the lock is not being held at the time this operation is invoked, this acts exactly the same as a LOCK operation.

Operation

```
Query(OUT u_long currentMode)
```

This operation is used by a thread to obtain the state of a lock at the time that the operation is executed. The lock state is returned in the 'currentMode' parameter.

3.2.3 Time Services

This object offers time-related services to clients. Included in these services is a reading of the value of the local node's real-time clock, and an operation that can be used to suspend the execution of a thread until (at least) a specified amount of (real) time has passed. The interface to the *TimeServices* object is shown in Figure 8, and its operations are described in the following text.

```
SYSTEM OBJECT TimeServices {
    OPERATION GetTime(OUT u_long timeval)
    OPERATION Sleep(IN u_long delay)
}
```

Figure 8: Time Services Object

Operation

GetTime(OUT u_long timeval)

This operation reads the system clock, and returns an unsigned integer value which represents the current time (in units of 200 nanoseconds). The time value is returned in the 'timeval' parameter.

Operation

Sleep(IN u_long delay)

This operation suspends the execution of the invoking thread for (at least) a period of time indicated by the 'delay' parameter. In Release 1 of Alpha, the 'delay' parameter is an unsigned integer that represents the amount of time the thread has to be delayed in terms of milliseconds.

3.3 I/O Objects

Release 1 of the Alpha kernel also provides a pair of system service objects that provide the client with the ability to perform I/O with machines external to the testbed's execution environment. See [Northcutt 88c] for a description of the overall testbed environment.

These I/O system service objects provide the basic interfaces necessary for a node to communicate over the application processor's serial interface line, and the node's network connection. The following subsections provide a description of each of these objects and their interfaces.

3.3.1 TTY

The *TTY* system service object provides a simple send/receive interface to a serial line unit on a testbed node's application processor. The object is capable of performing functions equivalent to the "raw" and "cooked" modes of serial I/O found in UNIX. This is to say that the *TTY* object can be used to perform either direct (i.e., unbuffered, single-character) I/O, or buffered I/O with line editing (i.e., complete, carriage-return-terminated strings are returned).

The mode that the *TTY* object operates in is defined by invoking an operation on the object, and it can be changed at run-time. This object only provides the basic I/O functions—functions such as reliable block transfers are intended to be implemented by higher-level objects that make use of this one.

The interface to the *TTY* object is shown in Figure 9, and its operations are described in the following text.

```

SYSTEM OBJECT  TTY {
    OPERATION  Init(IN u_long mode)
    OPERATION  Send(IN char msg)
    OPERATION  Receive(OUT char msg)
}

```

Figure 9: TTY Object

Operation

Init(IN u_long mode)

This operation initializes the object and defines the mode of operation based on the value of the 'mode' parameter. Encoded in this parameter is the desired baud rate, and the mode (i.e., either *raw* or *cooked*).

Operation

Send(IN char msg[])

This operation takes a character string (i.e., the 'msg' parameter), and sends it to the serial port encapsulated by the *TTY* object. In raw mode this operation returns when the characters have all been transmitted, while in cooked mode the operation returns once the characters have been buffered for output.

Operation

Receive(OUT char msg[])

This operation waits for characters to arrive on the object's encapsulated serial port, and then returns them in the 'msg' parameter. In raw mode, a single character is returned, while in cooked mode, the operation does not return until it has received a full input line (i.e., a carriage-return-terminated string). In either case, the invoking thread blocks until characters are available.

3.3.2 Network I/O

This system service object provides an extremely simple message transport interface to the testbed's communications interconnect (i.e., Ethernet). A simple (IP/UDP-based) communications protocol is used to permit short (i.e., less than 256 byte) messages to be exchanged, across the testbed network, with a UNIX process on a gateway machine. This object provides basic, reliable block I/O functions, and higher-level functions can be built on top of it.

The interface to the *NetIO* object is shown in Figure 10, and its operations are described in the following text.

```

SYSTEM OBJECT NetIO {
    OPERATION Init()
    OPERATION Send(IN extmsg msg)
    OPERATION Receive(OUT extmsg msg)
}

```

Figure 10: Network I/O Object**Operation**

Init()

This operation initializes the network I/O object. Following the invocation of this operation a node is capable of sending and receiving messages from the gateway machine. This operation need only be executed once (typically at node start-up).

Operation

Send(IN extmsg msg)

This operation takes a message (in the 'msg' parameter), queues it for transmission, and then returns to the invoker.

Operation

Receive(OUT extmsg msg)

This operation attempts to dequeue a message and place it into the 'msg' parameter. The invoking thread is blocked until a message arrives if there is not one immediately available.

3.4 Miscellaneous Object

As a temporary measure in Release 1, a number of necessary kernel functions have been grouped together into a "catch-all", system service object. This is known as the *TempServices* object, and provides clients with the ability to manage wellknown capabilities, make a thread ready for execution, modify the importance of a thread, and enable the scheduler's run-time statistics gathering facility.

Many of these operations will be subsumed into the thread and object standard operations set in Release 2.

The interface to the *TempServices* object is shown in Figure 11, and its operations are described in the following text.

```

SYSTEM OBJECT TempServices {
  OPERATION AddWellknown(IN CAPA instance,
                        IN CAPA wellknown)
  OPERATION MakeReady(IN CAPA thread)
  OPERATION SetThreadImportance(IN CAPA thread,
                                IN u_long new_importance,
                                OUT u_long old_importance)
  OPERATION SchedStats(IN boolean switch)
}

```

Figure 11: Temporary Services Object

Operation

AddWellknown(IN CAPA instance,
IN CAPA wellknown)

This operation is used to add the wellknown capability given in the 'wellknown' parameter to the c-list of the object indicated by the 'instance' parameter.

Operation

MakeReady(IN CAPA thread)

This operation unblocks a newly created thread (indicated by the capability given by the 'thread' parameter) and makes it ready to run. This operation must be invoked on a newly created thread to start it executing. When a new thread is created it is blocked and cannot begin execution until this operation is invoked.

Operation

SetThreadImportance(IN CAPA thread,
IN u_long newImportance,
OUT u_long oldImportance)

This operation is used to read and modify the importance value of an existing thread. The thread to be modified is given by the capability in the 'thread' parameter. the new importance is an integer value given in the 'newImportance' parameter, and the value of the thread's importance prior to this change is returned in the 'oldImportance' parameter.

Operation

SchedStats(IN boolean switch)

This operation enables and disables the scheduling subsystem's statistics gathering facility. The facility is enabled when the 'switch' parameter is TRUE, and is disabled when it is FALSE.

References

- [Clark 88] Clark, R. K., Kegley, R. B., Keleher, P. J., Maynard, D. P., Northcutt, J. D., Shipman, S. E. and Zimmerman, B. A.
An Example Real-Time Command and Control Application.
Archons Project Technical Report #88032, Department of Computer Science, Carnegie-Mellon University, March, 1988.
- [Northcutt 88a] Northcutt, J. D.
The Alpha Operating System: Requirements and Rationale.
Archons Project Technical Report #88011, Department of Computer Science, Carnegie-Mellon University, January, 1988.
- [Northcutt 88b] Northcutt, J. D. and Clark, R. K.
The Alpha Operating System: Programming Model.
Archons Project Technical Report #88021, Department of Computer Science, Carnegie-Mellon University, February, 1988.
- [Northcutt 88c] Northcutt, J. D.
The Alpha Distributed Computer System Testbed.
Archons Project Technical Report #88033, Department of Computer Science, Carnegie-Mellon University, March, 1988.
- [Northcutt 88d] Northcutt, J. D. and Shipman, S. E.
The Alpha Operating System: Programming Utilities.
Archons Project Technical Report #88061, Department of Computer Science, Carnegie-Mellon University, June, 1988.
- [Shipman 88] Shipman, S. E.
The Alpha Operating System: Programming Language Support.
Archons Project Technical Report #88042, Department of Computer Science, Carnegie-Mellon University, April, 1988.

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

THIS PAGE INTENTIONALLY LEFT BLANK

Table of Contents

Abstract	B-1
1 Introduction	B-2
2 Kernel Structure	B-5
2.1 Alpha Hardware Architecture	B-6
2.2 Kernel Proper	B-7
2.3 Kernel Subsystems	B-8
2.4 IMI Facility	B-9
3 Kernel Proper	B-13
3.1 Support for the Major System Abstractions.....	B-13
3.1.1 Objects	B-13
3.1.2 Threads.....	B-15
3.1.3 Operation Invocations.....	B-19
3.1.4 Access Control.....	B-23
3.1.5 Thread Synchronization	B-24
3.1.5.1 Semaphore Mechanism.....	B-24
3.1.5.2 Lock Mechanism.....	B-25
3.2 Ancillary System Functions	B-26
3.2.1 Critical Resource Preallocation.....	B-26
3.2.1.1 Physical Memory Management Daemon.....	B-26
3.2.1.2 Object and Thread Management Daemons.....	B-27
3.2.2 Exception Handling	B-27
3.2.3 Initialization	B-29
3.3 Key Data Structures	B-31
3.3.1 Control Blocks	B-32
3.3.1.1 Kernel Control Block (KCB).....	B-33
3.3.1.2 Client Object Type Control Block (COTCB).....	B-33
3.3.1.3 Client Object Control Block (COCB).....	B-33
3.3.1.4 Client Thread Control Block (CTCB)	B-33
3.3.1.5 Kernel Object Control Block (KOCB)	B-33
3.3.1.6 Kernel Thread Control Block (KTCB).....	B-34
3.3.1.7 Storage Object Control Block (SOCB).....	B-34
3.3.2 The Dictionary	B-34
3.3.3 Virtual Memory Structures	B-35
3.3.3.1 Memory Management Hardware	B-36
3.3.3.2 Context Management Structures.....	B-37
3.3.3.3 Page Management Structures.....	B-38
3.3.3.4 Auxiliary Memory Management Structures	B-39

4	Scheduling Subsystem	B-41
4.1	Scheduler Requirements	B-41
4.2	Application Time Constraints	B-44
4.2.1	Programming Interface	B-44
4.2.2	Implementation issues	B-47
4.3	Scheduling Facility Implementation	B-49
4.3.1	Internal Interface	B-50
4.3.2	Internal Structure	B-53
4.3.2.1	Scheduling Subsystem Mechanisms	B-55
4.3.2.2	Interrupt Level Routines	B-58
4.3.2.3	Main Loop Level Routines	B-59
4.4	Scheduling Policies	B-60
4.4.1	Policy Definition Modules	B-61
4.4.2	Policy Example	B-61
4.4.2.1	Data Structures	B-61
4.4.2.2	Interrupt Level Routines	B-63
4.4.2.3	Main Loop Level Routines	B-65
4.5	Future Directions	B-71
5	Communication Subsystem	B-73
5.1	Distribution and the Alpha Programming Model	B-73
5.1.1	The Alpha Programming Model	B-73
5.1.2	Distribution in Alpha	B-75
5.1.3	Related Work	B-77
5.2	Communication Requirements	B-78
5.3	Alpha Communication Architecture	B-79
5.3.1	The Communication Virtual Machine	B-80
5.3.2	Communication Protocols	B-82
5.3.3	Related Work	B-83
5.4	Alpha Communication Protocols	B-85
5.4.1	RM Protocol	B-85
5.4.2	RI Protocol	B-86
5.4.3	TMAR Protocol	B-88
5.4.4	PT Protocol	B-91
5.4.5	MAN Protocol	B-91
5.4.6	AUL Protocol	B-92
5.5	Evaluation of Alpha Communication Architecture	B-92
5.5.1	Protocol Definition Model	B-92
5.5.2	Mechanisms	B-93
5.5.3	Flexibility	B-94
5.5.4	Language	B-94
5.5.5	Hardware	B-94
5.6	Future Directions	B-95
5.7	Summary	B-97

6	Storage Management Subsystem	B-98
6.1	Primary Storage Management.....	B-98
6.1.1	Physical Memory	B-100
6.1.2	Virtual Memory	B-101
6.1.2.1	Kernel Region	B-101
6.1.2.2	Client Thread Region.....	B-103
6.1.2.3	Client Object Region	B-104
6.1.2.4	Kernel I/O Region.....	B-105
6.1.3	Address Translation Table	B-105
6.1.4	Virtual Memory Fault Handling	B-107
6.2	Secondary Storage Management.....	B-107
6.2.1	Secondary Storage Facility Interface.....	B-108
6.2.2	Secondary Storage Objects	B-110
6.2.3	Low-Level Storage Management Details	B-112
6.2.3.1	Device Header.....	B-114
6.2.3.2	Space Allocation	B-114
6.2.3.3	Capability Mapping	B-116
7	Acknowledgments	B-119
	References	B-120

List of Figures

Figure 1	System Functional Design Structure.....	B-2
Figure 2	Intra-Node Kernel Structure	B-4
Figure 3	Logical Alpha System Structure	B-5
Figure 4	Alpha Experimental Testbed Node Architecture	B-7
Figure 5	Interprocessor Message Interface (IMI) Model	B-10
Figure 6	Thread Execution States	B-16
Figure 7	Parameter Page Manipulation on Invocation.....	B-18
Figure 8	Operation Invocation Facility	B-20
Figure 9	Types of Operation Invocation	B-22
Figure 10	Application Processor Load Module	B-31
Figure 11	Alpha Dictionary and Control Blocks.....	B-35
Figure 12	Memory Management Hardware	B-37
Figure 13	Virtual Memory Data Structures.....	B-40
Figure 14	Components of a Time-Value Function.....	B-46
Figure 15	Time Constraint Block Example.....	B-47
Figure 16	Application Processor/Scheduling Subsystem Interface	B-51
Figure 17	Scheduling Subsystem Internal Structure	B-54
Figure 18	Thread Control Block	B-62
Figure 19	Time Constraint Block.....	B-62
Figure 20	Example Objects and Threads	B-73
Figure 21	Thread Breaks Due to Node Failure	B-74
Figure 22	Logical View of the Alpha Communication Architecture	B-79
Figure 23	Expanded View of the Communication Architecture.....	B-80
Figure 24	Communication Virtual Machine Token Format.....	B-81
Figure 25	Logical Protocol Structure	B-82
Figure 26	Protocol State Diagram	B-83
Figure 27	Virtual Memory Hierarchy	B-99
Figure 28	Virtual Address Space Layout	B-102
Figure 29	Kernel Virtual Memory Region	B-103
Figure 30	Thread Virtual Memory Region.....	B-104
Figure 31	Object Virtual Memory Region	B-105
Figure 32	Logical View of Secondary Storage in Alpha	B-109

Abstract

This report describes the design and implementation of the kernel layer of Release 1 of the Alpha operating system.

Alpha is an adaptable decentralized operating system for real-time applications, being developed as a part of the Archons project's on-going research into real-time distributed systems. Alpha is unique in two particularly significant ways—first, it is *decentralized*, providing reliable resource management transparently across physically dispersed nodes, so that a distributed application can be implemented and executed as though it were centralized; and second, it provides comprehensive, high technology support for real-time applications, particularly supervisory control systems (e.g., industrial automation applications), which are characterized by predominately aperiodic activities executing under critical time constraints (such as deadlines). Alpha is extremely adaptable so as to be easily optimized for a wide range of problem-specific functionality, performance, and cost.

Alpha has been designed to support real-time, reliable, distributed applications. Since Alpha is, in fact, an example of just such an application, it has been designed in a layered manner. The lowest layer, the *kernel*, supports the fundamental abstractions of the Alpha programming model [Northcutt 88b]. The system layer of the operating system, which is based on the kernel, can then utilize the power of this programming model, as well as the user's application program.

In this report, details are given concerning the design and implementation of the major facilities provided by the Alpha kernel. This includes support for the object, thread and invocation abstractions, as well as the facilities for inter-node communication management, virtual memory management, application processor management, and secondary storage management.

1 Introduction

This report describes the design and implementation of the kernel layer of the Alpha operating system. The kernel of Alpha provides support for the Alpha programming model (as defined in [Northcutt 88b]) by implementing the kernel interface (defined in [Northcutt 88d]). The overall system context within which the kernel layer of Alpha exists and the various functions provided at each layer are illustrated in Figure 1.

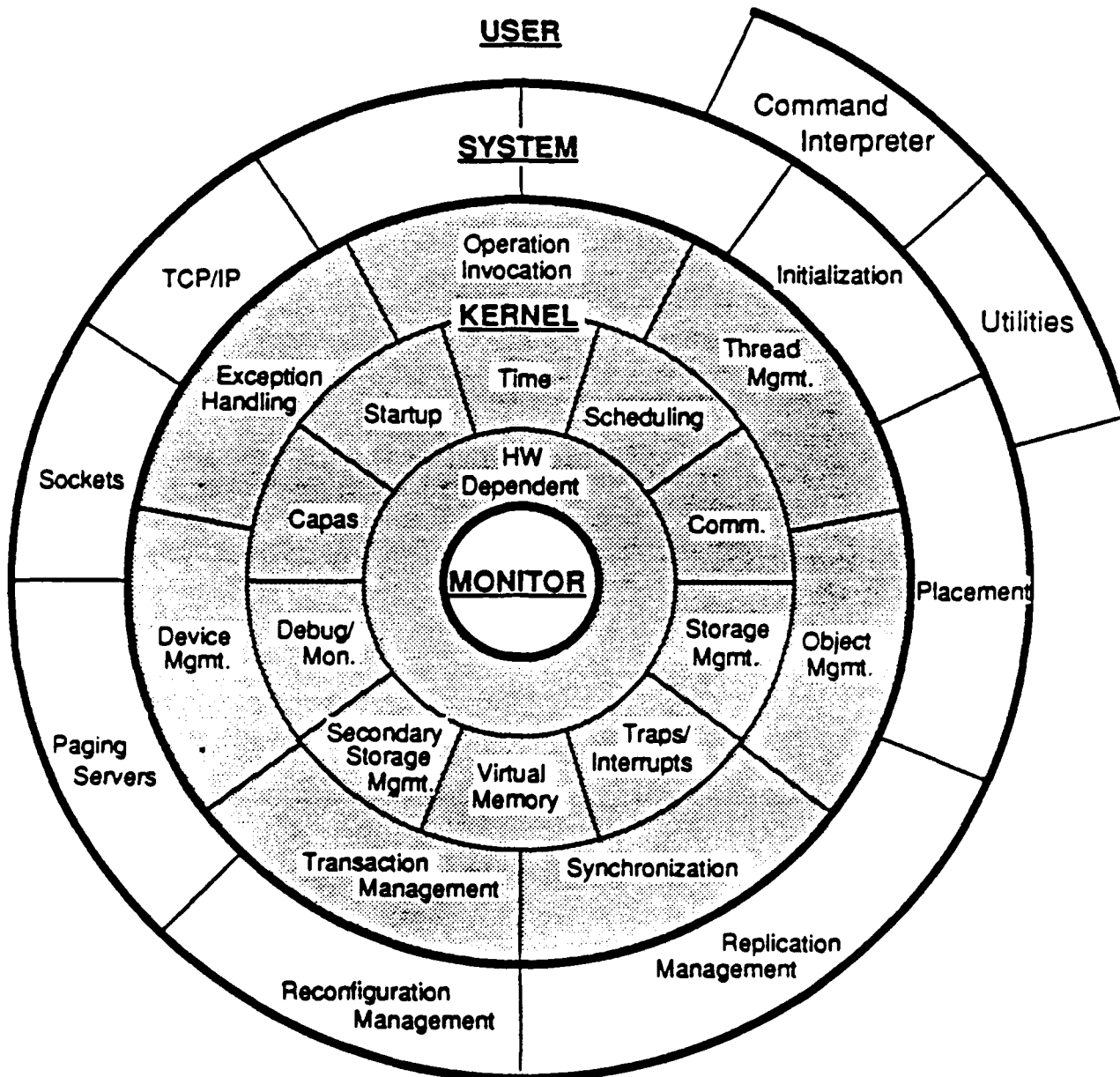


Figure 1: System Functional Design Structure

Alpha has been designed to support real-time, reliable, distributed applications. Since Alpha is, in fact, an example of just such an application, it has been designed in a layered manner. The lowest layer, the *kernel*, supports the fundamental abstractions of the

Alpha programming model [Northcutt 88b]. The system layer of the operating system, which is based on the kernel, can then utilize the power of this programming model, as well as the user's application program.

The Alpha kernel layer is composed primarily of C language routines, with a few assembly language routines. The kernel routines are divided into two groups: routines that provide kernel services to clients by means of an object interface, and basic support routines, that provide services needed for the kernel-level abstractions, but that are not visible to the kernel's clients. Kernel routines are grouped according to their logical function, corresponding to kernel *facilities*. Each kernel facility is composed of one or more *mechanisms* that perform a specific aspect of a facility's functionality. The Alpha kernel facilities include: object management, thread management, operation invocation, concurrency control, replication management, and transaction management. With the exception of the operation invocation facility (whose function spans all of the nodes in the system), the kernel mechanisms perform their functions locally—i.e., on the node at which they exist. Each of these mechanisms provides its particular service to the kernel's clients via an object interface, and all access to kernel-provided functionality is achieved via the operation invocation facility. For the most part, the kernel mechanisms are implemented as (one or more) routines that make use of specific hardware components at a node. In support of the kernel mechanisms, a number of lower-level mechanisms are provided as routines that manage internal data structures, primary memory, inter-node communication, secondary storage, and peripheral devices.

This report describes details of the design and implementation of the major facilities provided by the Alpha kernel. This includes support for the object, thread and invocation abstractions, as well as the facilities for inter-node communication management, virtual memory management, application processor management, and secondary storage management. Because the kernel is replicated at each node, this discussion centers on the functions performed by a replica of the kernel on a single node. Furthermore, the structure of this report parallels that of the kernel itself—i.e., each of the kernel's major structural components are discussed in a separate chapter. Figure 2 provides an illustration of the internal structure of a replica of the kernel.

The following description of the kernel is given in the context of the hardware described in [Northcutt 88c]. Machine dependencies exist in most operating system kernels, and the Alpha kernel is no exception. However, although machine dependencies can be found in the kernel, the techniques used for its construction are relatively machine-independent. Furthermore, those machine dependencies that do exist within the Alpha kernel are clearly marked and isolated from the rest of the code. As a result, Alpha could be modified with relative ease to work with other loosely-coupled architectures, having processors supporting either a linear, or a segmented, virtual address space and demand-paged virtual memory.

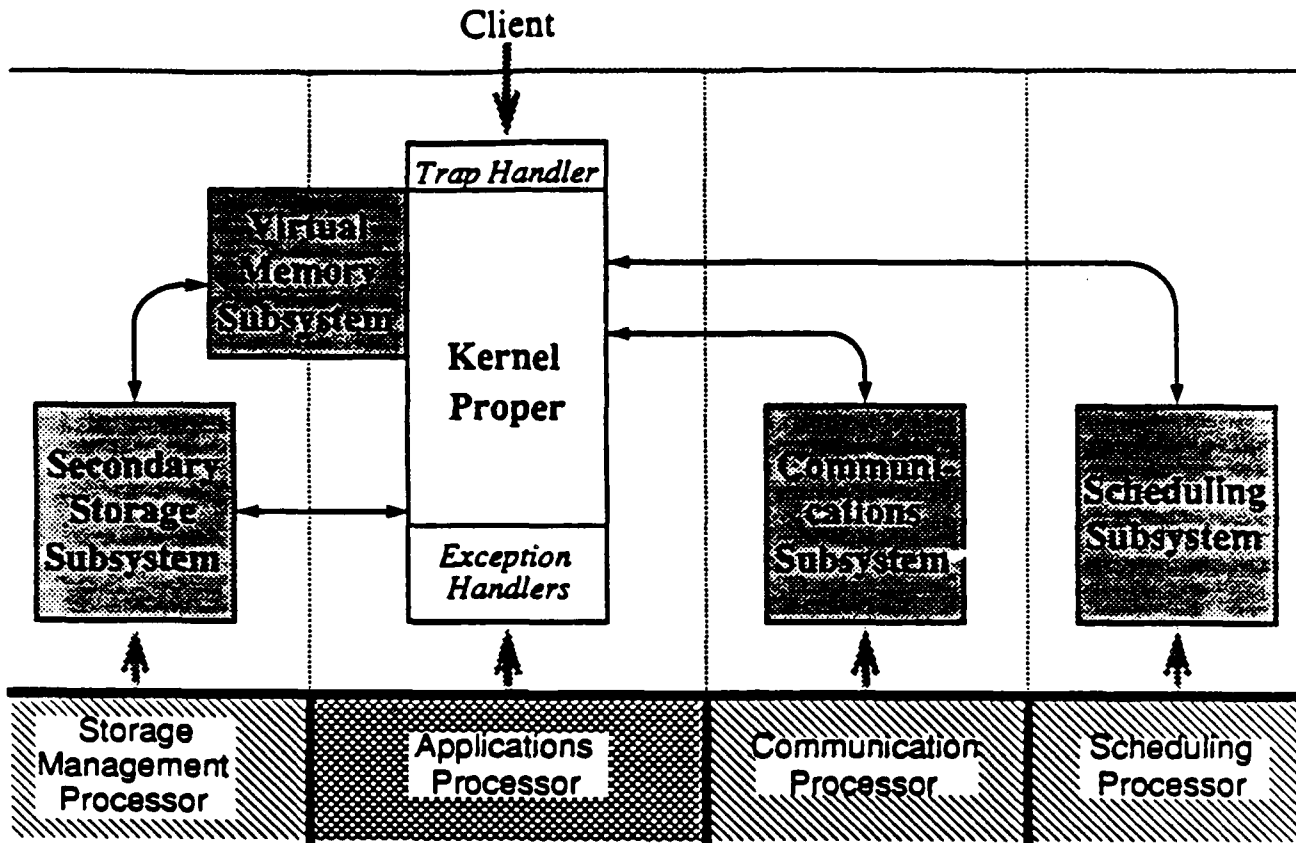


Figure 2: Intra-Node Kernel Structure

2 Kernel Structure

The logical structure of the Alpha Operating System ([Northcutt 87]) is illustrated in Figure 3. Applications run on top of the Alpha Operating System, which internally has the following components:

1. **Kernel** — which supports the basic abstractions of objects, threads, and invocations;
2. **Executive and System** — which offer higher-level operating system functions and embody global resource management policies; both the Executive and System are comprised of objects and threads, and the distinction between them is unimportant for this discussion.

The Alpha Kernel executes on top of a small layer of software, called the Monitor, that provides the lowest level hardware access, allowing the execution of primitive operations and minimal debugging.

The internal structure of the Alpha Kernel, shown in Figure 3b, features two major parts:

- the *Kernel Proper* and
- a set of *Kernel Facilities*.

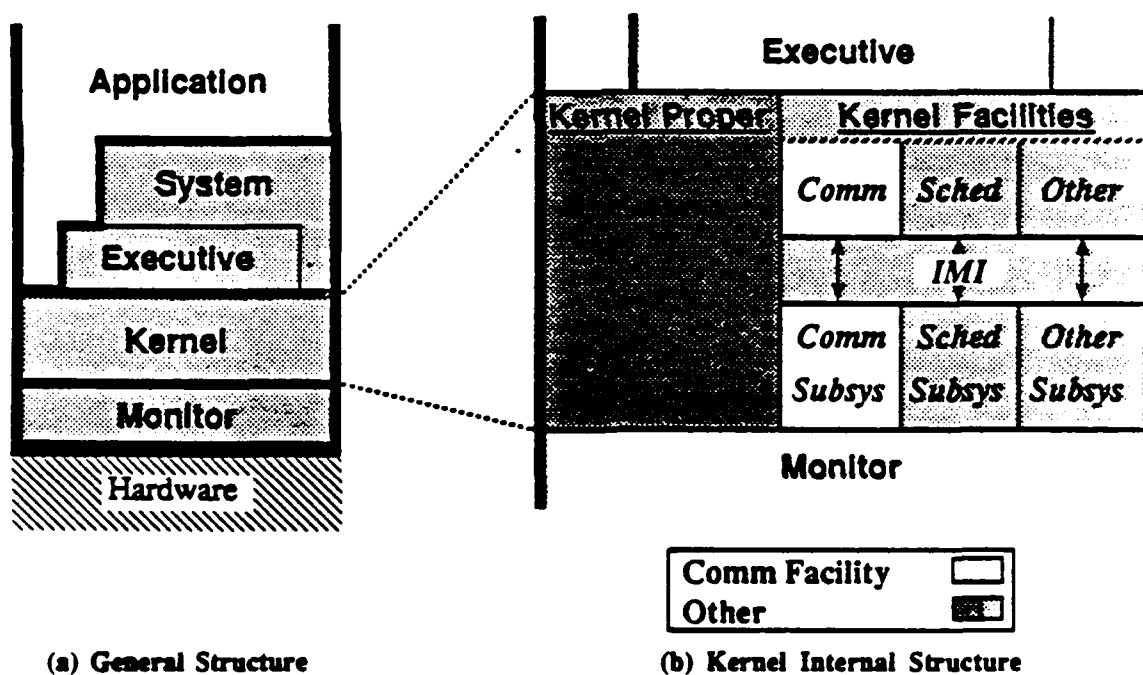


Figure 3: Logical Alpha System Structure

The Kernel Facilities provide important kernel services and are particularly interesting because their services may typically execute concurrently. Consequently, each of the facilities is composed of two parts:

- one of which executes on the Applications Processor along with the Kernel Proper, and
- another which executes on a separate co-processor.

- The part of a kernel facility that executes on a separate co-processor is called a *Kernel Subsystem*. Every kernel facility has a corresponding Kernel Subsystem. For example, in Figure 3b, the Communication Facility, which has been shown in white, while the rest of the Kernel Proper and the other kernel facilities have been shaded, consists of two parts, one of which is labeled the Communication Subsystem..

The Alpha Kernel structure provides actual concurrency for operating system tasks—the Kernel Proper and all of the Kernel Subsystems can operate simultaneously. For instance, the Kernel Proper can perform a local invocation while the Communication Subsystem processes a remote invocation and monitors the current health of all of the threads rooted on the local node. At the same time, the Scheduling Subsystem can revise the execution schedule, and the Storage Management Subsystem can locate and acquire the initial load image for the creation of a new object instance.

The functional partitioning of the kernel into the Kernel Proper and a number of Kernel Facilities that feature independent Kernel Subsystems also facilitates the use of dedicated, special purpose hardware that may be selected and optimized to perform only those functions that are required of a given subsystem, rather than a set of general purpose tasks.

Each Kernel Subsystem communicates with the portion of its facility that resides on the Application Processor with the Kernel Proper by means of a dedicated, two-way, communication channel, known as an *Interprocessor Message Interface (IMI)* channel. (Some subsystems, such as the Communication Subsystem, also share a limited amount of memory with the Kernel Proper.) The IMI channel, which is described in detail below, permits one processor to send a message, containing a command or data, to another processor. In addition, the sender may optionally interrupt the receiver to notify it of the presence of a new message—thus providing an efficient delivery mechanism.

2.1 Alpha Hardware Architecture

Alpha is intended to execute on a set of nodes interconnected by one or more global communication buses. (Currently, the Alpha Experimental Testbed uses an Ethernet as its single global bus.) Since there are a number of operating system activities that can benefit from genuine hardware concurrency—for instance, communication processing and scheduling analysis, in addition to application execution—each testbed node contains multiple processors. As shown in Figure 4, each node has an Application Processor (AP), a Scheduling Co-Processor (SP), a Communication Co-Processor (CP), an Ethernet Controller Board, and a pool of shared memory, all connected by a Multibus backplane. In addition, other co-processors may optionally be added—for instance, to manage primary and secondary storage, a Storage Management Co-Processor (SMP) and one or more Disk Controller Boards may be added. (For this purpose, Alpha currently supports Fujitsu MK2284 168 Megabyte disk drives controlled by a Xylogics 450 disk controller.)

In the current version of the Alpha Experimental Testbed, the AP is a Sun Microsystems Inc. (SMI) version 2.0 CPU board, while all of the co-processors are SMI version 1.5 CPUs. The Ethernet Controller Board, also from SMI, is used to access the Ethernet. It is based on the Intel 82586 Local Area Network (LAN) Co-Processor and has 256

Kilobytes of on-board memory for use by the 82586. In addition, there are 256 Kilobytes of Multibus Shared Memory.

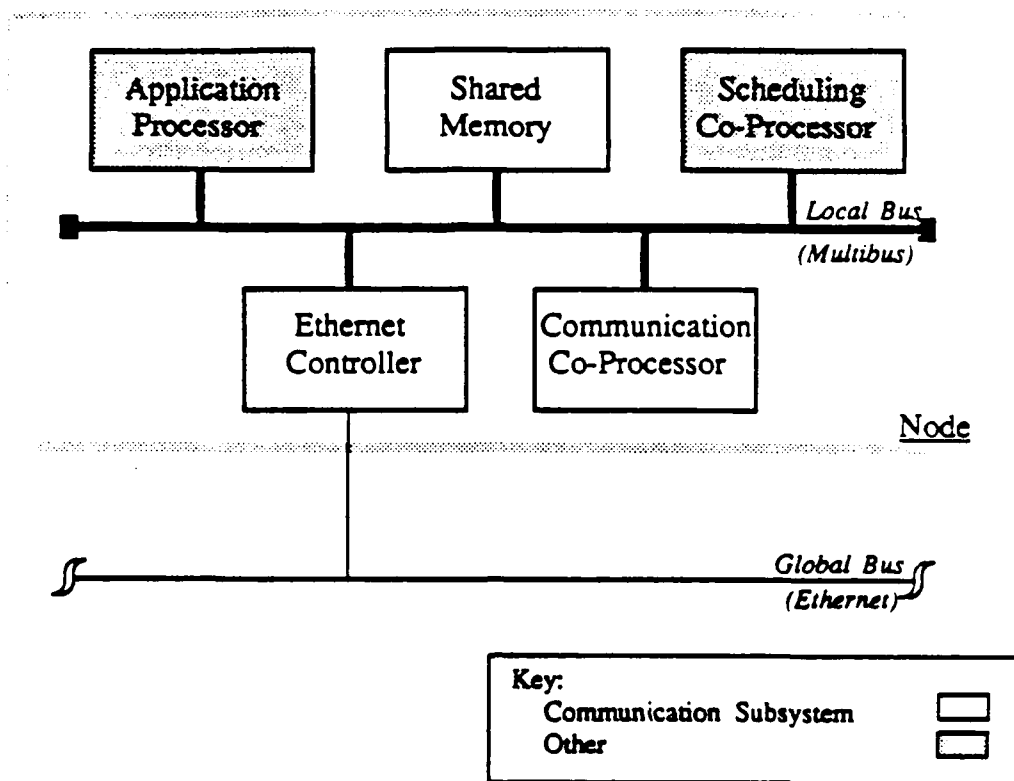


Figure 4: Alpha Experimental Testbed Node Architecture

While the testbed architecture provides a rich base for experimentation, it does impose a set of requirements and constraints on Alpha. Specifically, an Ethernet must interconnect the Alpha nodes, and this Ethernet must be controlled by the SMI Ethernet Controller Board. Furthermore, only off-the-shelf hardware can be considered if any additional hardware is required to provide the desired communication architecture.

As shown in Figure 4, the Communication Subsystem hardware consists of the CP, the Ethernet Controller Board, and the Multibus Shared Memory. All of the other boards are associated with other subsystems or the AP.

2.2 Kernel Proper

Ideally, the Application Processor(s) spends all of its time executing application code—that is, code that is specific to the application at hand, not code that is devoted to maintaining operating system data structures, performing operation invocations, scheduling threads, managing virtual memory mappings, or controlling disks drives or other peripherals. In such an ideal situation, the operating system functions are provided by one or more operating system processors.

Of course, this ideal system cannot be completely realized. At the very least, a relatively small amount of code must reside on the Application Processor to allow it to request operating system services from the operating system processors and to allow it

to react to messages received from those processors. It may also be more efficient to provide some operating system functions directly on the AP if the communication and coordination overhead incurred in using multiple processors was too great.

Alpha aspires to provide an ideal system in the terms just described. However, due to the hardware available when Alpha was begun and the immense technological leap to provide an ideal system, a short-term compromise was reached. Currently, while a number of important operating system functions are run on dedicated processors, many other system functions are provided on the Application Processor. Although this might suggest that the operating system has not been completely separated from the application, it can more appropriately be viewed as a reflection of the fact that the Alpha Experimental Testbed has no Application Processor. To remedy that situation, one or more additional processors, depending on the demands of the application, should be added to each testbed node to act as true, dedicated Application Processors.

For Alpha Release 1, the Kernel Proper executes on the each node's Application Processor, along with a portion of each kernel facility that is provided on that node.

Since a detailed description of the Kernel Proper is offered in the following chapter, a brief example will serve to illustrate the division of function between the software residing on the AP and that residing on a co-processor (corresponding to one of the Kernel Subsystems): the Communication Facility, which must meet the communication and distribution requirements outlined in Chapter 5, is composed of software residing in both the Application Processor and the Communication Co-Processor. The functions performed by the Communication Facility on the Application Processor include local invocation support (that is, support for invocations whose target is on the same node as the invoker), memory mapping manipulations and thread control operations to initiate and complete remote invocations, and exception handling for failed invocations on replicated objects. All of these functions require the use of information maintained by the Kernel Proper and unavailable to the Communication Subsystem directly. Consequently, they are performed by the Communication Facility software on the Application Processor, which does have access to this information.

2.3 Kernel Subsystems

Alpha currently has three subsystems:

- the Scheduling Subsystem,
- the Communication Subsystem, and
- the Storage Management Subsystem.

The standard Alpha node employs a Scheduling Co-Processor and a Communication Co-Processor to support their respective subsystems. Optionally, any of the nodes may also have a Storage Management Co-Processor. In principle, each node could contain a single processor that executes both Alpha and the application.

Future versions of Alpha may support more subsystems than the current version. They may also use more co-processors, possibly with specialized hardware support, to improve and expand on the services provided by the existing subsystems. For example, the services provided by the Storage Management Subsystem may be refined and enhanced so that the Storage Management Subsystem will be replaced by two subsys-

tems—the Primary Storage Management Subsystem, which would manage primary memory, including services such as anticipating the primary memory needs of the local objects and threads to improve performance, and the Secondary Storage Management Subsystem, which would carry out disk management and reliability support functions.

The hardware used for the subsystems was described earlier in Section 2.1 and, in greater depth, in the document, *The Alpha Distributed Computer System Tested* ([Northcutt 88c]).

For the present, in order to contrast the type of function carried out by a co-processor with that carried out by the Application Processor, consider once again the example of the Communication Subsystem: the Communication Subsystem accepts high-level commands from the Application Processor and performs them in parallel with the continued execution of the Application Processor. The Communication Subsystem handles per-packet interrupts and interrupts the Application Processor only to indicate the occurrence of high-level communication events (e.g., an incoming invocation arrival or an outgoing invocation completion). The specific functions provided by the Communication Subsystem include message disassembly and reassembly; participation in the intra- and inter-node thread maintenance protocols, the atomic transaction refresh protocols, the atomic transaction commit and abort protocols, and object replication protocols; and all of the low-level acknowledgment generation and handling.

In the following chapters, both the structure of, and the functions provided by, each of the subsystems will be described in detail.

2.4 IMI Facility

Processors in an Alpha node communicate by means of the Interprocessor Message Interface (IMI) facility. The Application Processor (AP) is connected to each of its node's co-processors by an IMI *channel*. As illustrated by Figure 5, each IMI channel is composed of a pair of unidirectional communication links, called *streams*. For each processor, one stream is designated the *out stream* and transfers IMI messages originated by that processor, while the other stream is designated the *in stream* and receives messages from the other processor on that channel.

Each stream queues IMI messages. A small queue proved to be necessary for Alpha Release 1 in order to avoid deadlock situations that occurred when one channel was used to send multiple messages rapidly between a pair of processors. Despite the presence of the queue, which relaxes the response time requirements that must be observed by a processor to avoid deadlock, it is still desired that the co-processors respond as rapidly as possible to IMI messages, which are typically commands, issued by the AP. Quick co-processor response allows the AP to dedicate itself to application code as much as possible, rather than forcing the AP to spend (relatively) large amounts of its time waiting on responses to IMI messages or waiting to send its next command to a co-processor, because of a full IMI stream queue. Consequently, all of the co-processors have been designed to respond rapidly to AP-originated IMI messages.

Interrupts are employed to provide good responsiveness to IMI message arrival without imposing a high on-going system overhead. The originator of an IMI message may use the *attention line* shown in Figure 5 to cause an interrupt to be sent to the message's receiver once the message has been placed in its out stream.

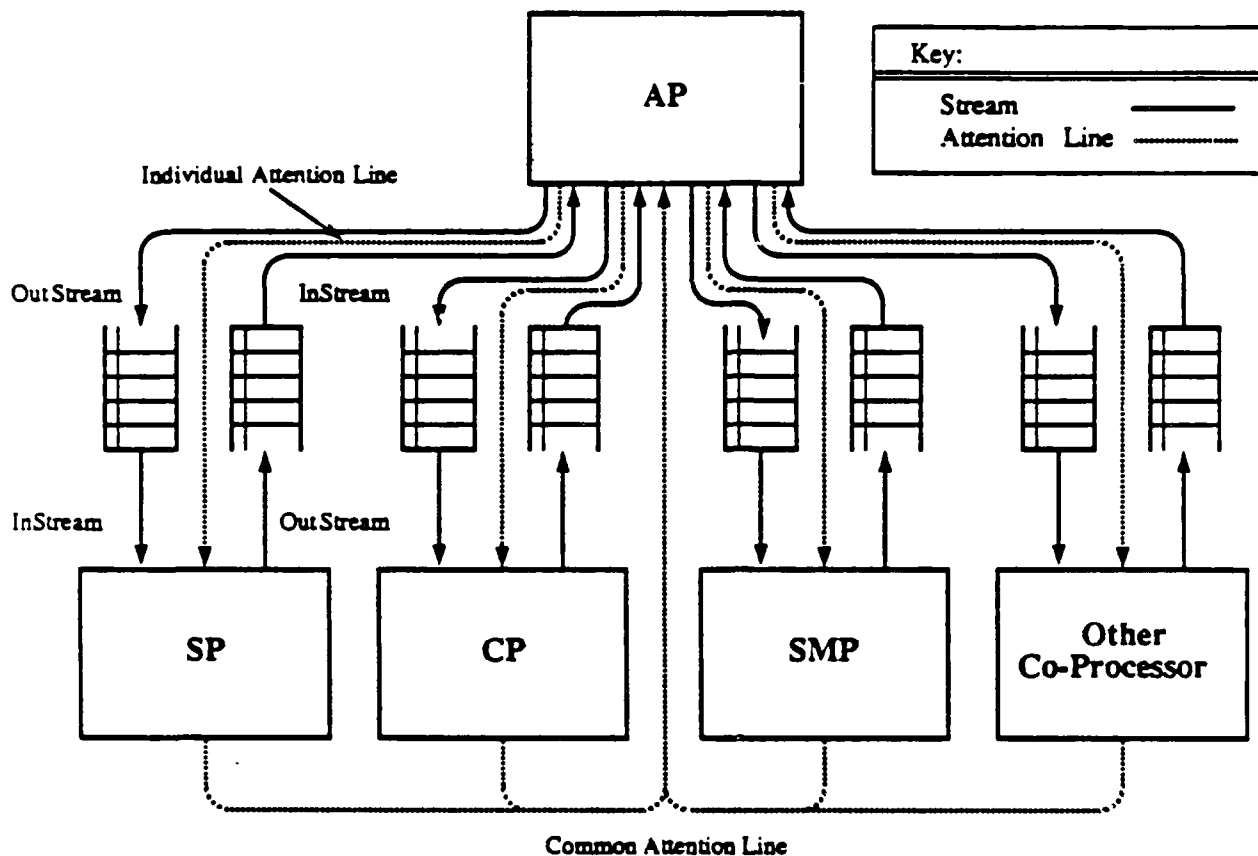


Figure 5: Interprocessor Message Interface (IMI) Model

Multibus interrupt lines are employed to send these interrupts between processors. The AP on each node has an *Interrupt Generator* that can initiate an interrupt on any of four Multibus interrupt lines under program control. An Interrupt Generator is mapped into the Multibus address space for each node so that it may be accessed by any of the processors on that node.

The AP uses a different dedicated Multibus interrupt line to command the attention each of the co-processors, while all of the co-processors use a single Multibus interrupt line to signal the AP. This means that the AP must examine all of the IMI channels when it receives an interrupt in order to find the source of the interrupt, whereas each co-processor assumes that it is receiving an IMI message from the AP. (If there were more interrupt lines available, dedicated lines would have been used to signal the AP as well.) The hardware involved in providing this facility is described in [Northcutt 88c].

The IMI facility is intended to be general-purpose. It supports both synchronous and asynchronous IMI messages, where a synchronous message is one for which the originator of the message "busy waits" until a response (typically, an "answer," that is, a response that includes some data, not just an "acknowledgment") is received, while an asynchronous message is one that is simply sent by the originator without any explicit notion of a response. Following the transmission of an asynchronous message, the IMI channel is once again available for other messages in either direction.

The IMI facility also allows the originator of a message to decide whether or not an interrupt should be sent to the message receiver for each message sent over the channel.

In practice, the full generality of the IMI facility has not been used. To date, interrupts are always used, and synchronous messages are never used. However, this is not because there are no messages that seem to be synchronous in nature. Rather, it has proven to be difficult to use a single IMI channel to handle both synchronous and asynchronous messages simultaneously. In general, the processor on one end of an IMI channel cannot know what the processor on the other end of the channel will next send over the channel. As a result, when a synchronous message is sent over the channel, the originator cannot be sure that the next message received on the channel's in stream will be the response to the synchronous message. Instead, it may easily be an asynchronous message that is unrelated to the synchronous message just sent.

Due to the high volume of data that is typically included in operation invocations, the AP and CP interact by means of shared memory, as well as interprocessor interrupts. In fact, the AP/CP interface is similar to the hardware-implemented *mailbox* schemes used for interprocessor communication in various systems. The interface allows high-bandwidth communication through the use of shared memory, while attaining low overhead and high responsiveness through the use of interrupt-based signaling. This interface also provides for increased communication performance, deriving from the ability of the Application and Communication Processors to effectively move memory through the manipulation of memory mapping tables, rather than performing the costly copying of memory blocks that is normally associated with interprocessor communication activities.

The IMI interconnection architecture has influenced the interface presented by each of the subsystems to the AP in at least two major ways. First, all of the IMI channels connect the AP to some co-processor. There are no IMI channels that connect co-processors directly. As a result, the AP must occasionally act as an intermediary. For example, when a remote operation invocation is performed, the AP queries the Scheduling Co-Processor (SP) to acquire the current scheduling parameters for the invoking thread before sending an IMI message to the Communication Co-Processor (CP) indicating that it should perform the invocation. If the CP and SP were connected by an IMI channel, the scheduling parameters could be acquired by the CP without the involvement of the AP. Since off-loading non-application processing tasks from the AP is one of the goals of Alpha, a richer interprocessor communication architecture would seem to be advantageous.

The fact that the IMI channels act as mailboxes is the second factor that has affected subsystem interfaces. In general, data is sent among processors through the IMI channels. While Alpha has attempted to use shared memory where possible to reduce the amount of data that must be copied among processors, without a true multiprocessor node architecture—where each processor can potentially access all of the memory on the node—this effort cannot succeed completely. For instance, the Kernel Proper on the AP maintains a list of all of the objects and threads that are local to that node. The CP needs to have an identical list to recognize those operation invocations that are broadcast over the Ethernet to which it must respond. However, putting this information in the shared Multibus memory would be impractical—it would take longer to access from the Kernel Proper and it would be severely restricted in size. As a result, in the current implementa-

tion two separate lists must be maintained. This increases the amount and type of information that must be explicitly passed from the AP to the CP, thereby altering the interface that would otherwise exist between the AP and CP. (Similar examples can be offered for the AP/SP interface as well.)

tion two separate lists must be maintained. This increases the amount and type of information that must be explicitly passed from the AP to the CP, thereby altering the interface that would otherwise exist between the AP and CP. (Similar examples can be offered for the AP/SP interface as well.)

3 Kernel Proper

As defined in previous chapter, the Kernel Proper is the portion of the kernel layer of the Alpha operating system that coexists on the Application Processors with both higher-level-system- and user-defined code. The Kernel Proper is at the center of the operating system, surrounded by the kernel subsystems which comprise the remainder of the kernel layer.

The following sections describe the design and implementation of the major functions of the Kernel Proper. First descriptions of the functions provided in support of the system's programming model are given, followed by descriptions of the Kernel Proper's major data structures.

3.1 Support for the Major System Abstractions

This section describes the design and implementation of the portions of the Kernel Proper that provide support for the programming model's major abstractions. The intent of the following discussion is to provide an understanding of the implementation techniques used in Alpha. Throughout this chapter it is assumed that the reader is familiar with the Alpha programming model (as defined in [Northcutt 88b]), and so only design and implementation issues will be discussed.

Specific details relating to the object management, thread management, operation invocation, concurrency control, and access control abstractions are presented here.

3.1.1 Objects

In accordance with the abstract definition of Alpha's objects, objects are simply regions of memory that contain code and data. The data portion of an object consists of several different types of data, with the major division being between statically and dynamically allocated data. An object is created with a unit of static data that is global to all threads that might be executing within the object at any point in time. The size of this data region cannot be increased or decreased, and the static data region consists of both initialized and uninitialized data. The dynamically allocated portion of an object's data is known as the object's *heap* region and can be increased and decreased in size at run-time.

Objects are not permanently associated with any particular context, but rather "float" around in the system's memory until they have operations invoked on them by threads. For this reason, an object has associated with it all of the virtual memory and secondary storage structures necessary for the management of the object. While the data region of an object may be read or written by any thread that has the object mapped into its context, the code portion of an object's memory space is protected in an *execute-only* fashion.

When an invocation is made on an object, the memory that makes up that object is mapped into the invoking thread's context. This is done by manipulating the virtual memory management address maps associated with the invoking thread. The act of mapping an object into a thread's context can be viewed as a "partial context swap," because a portion of the executing thread's context is replaced. However, this activity is considerably simpler (and hence quicker) than a full context swap, because the executing thread does not give up the processor and so the machine state does not have to be saved and restored.

An object can be simultaneously mapped into multiple threads' address spaces on a given node. This permits for concurrent execution of threads within (either the same or different operations of) objects. All threads executing within an object share the object's common (execute-only) code section, as well as the different (read/write) parts of the object's data section. To help manage the synchronization of threads executing within objects, object control blocks contain references to all of the threads that are currently active within an object.

Furthermore, objects provide protected entry points into their code section by way of a software version of hardware "gates" (i.e., capabilities). This is done by way of a simple, hardware-protected indirection table known as an *object entry point table*, which is generated by the object-language pre-processor [Shipman 88]. The object entry point table contains a count of the number of entry points into the object, followed by a list of addresses of the entry points. This is used by the operation invocation facility to validate operation invocation indices and to locate entry points in the objects. Thus, an object control block also contains all of the permanent capabilities that belong to an object at a given point in time.

More detailed descriptions of the memory management and protection aspects of objects are given later in this chapter.

The preceding discussion focussed on the typical type of object found in Alpha—i.e., a *client object*, an object that exists in the machine's user space. There is another type of object in Alpha that is implemented quite differently. These are known as *kernel objects*—they are actually just kernel extensions appearing to the programmer as standard objects. Kernel objects exist within the Kernel Proper's address space, and do not make use of the system's virtual memory features. While the programmer views the kernel object optimization as just another compile-time-specified attribute of an object, the implementation of kernel objects is quite different from that of client objects. Both kinds of object type specifications and definitions are written in the same manner, the only difference is that the KERNEL qualifier is used in the header to indicate that instances of objects of this type are implemented as kernel objects.

Both client and kernel objects are composed of the same major logical components—i.e., code, data, and control information (i.e., virtual memory information, secondary storage information, capability lists, etc.). With client objects the code and data components exist in separate virtual memory extents, whereas with kernel objects, the code and data of all kernel objects are all combined within the code and data portions of the kernel. Kernel objects are an integral part of the kernel, and do not have separate virtual memory structures or their own storage objects associated with them. For this reason, kernel object control blocks do not contain references to virtual memory management information, or the thread context(s) the object is currently mapped into. The kernel object's control block is a quite simple structure that consists mainly of a pointer to the object's entry point block in the kernel's data region, the kernel object's capability list, and its synchronization data structures.

Just as with client objects, the object language pre-processor generates an entry point table for each kernel object. However, entry point tables for client objects are placed at the start of each object's data region, while the entry point tables for kernel objects are placed in various locations in the kernel's data part.

In addition, kernel objects are linked differently than client objects (i.e., kernel objects are linked with the kernel, as opposed to being separate, stand-alone entities), and the object language pre-processor generates different code for the operation invocations made by kernel objects. Since the kernel objects all execute in the kernel context and in supervisor mode, the usual trap instruction used by client objects (executing in user mode) to gain access into the kernel on invocations is not used. Instead, the pre-processor only generates the code to marshal the necessary parameters and then to make a call to the kernel's operation invocation routine.

3.1.2 Threads

There are two varieties of threads in Alpha—*kernel* threads and *client* threads. Kernel threads exist as both kernel optimizations and fundamental units out of which client threads are constructed. In Alpha, a client thread is bound to its own virtual address space, separate from all others. Kernel threads do not have their own virtual address spaces, but rather execute entirely within the kernel; kernel threads are a simplified form of a thread that exists within the kernel's portion of a virtual address space.

A kernel thread is designed to contain all of the basic (non-virtual memory) information necessary for managing a thread. Kernel threads keep the processor's registers when a thread is de-scheduled (i.e., unbound from the processor on which it is executing), along with copies of the thread's user and supervisor stack pointers. Furthermore, kernel threads contain the state information known as the thread's *environment*, that consists of the information necessary for managing threads within nested atomic transactions, and the information needed by the scheduling facility to perform the desired thread scheduling functions. Kernel threads are implemented in a fashion very much like that of standard processes in more traditional systems.

Each client thread in the system is constructed from a kernel thread—i.e., a client thread consists of a kernel thread with additional information. Client threads differ from kernel threads largely in as much as they include virtual memory information to support the virtual address space that is associated with client threads. A client thread contains all the necessary virtual memory information for the thread's address space, and a reference to the currently mapped client object (i.e., the object in which the thread is currently executing). When executing within the kernel, each client thread uses its own kernel stack. This allows each thread to block on its own stack, saving any state necessary, prior to activating another thread. This serves to simplify the kernel's implementation since otherwise, a communal stack would have to be shared among all threads.

Included in the information associated with a thread is an indication of its execution state—i.e., the mode it is currently in with respect to the scheduling facility. Figure 6 provides the state/transition illustration of thread execution states in Alpha. In this figure, the various states that a thread can be in are represented by circles. The legal transitions between thread states are represented by arcs (which are labeled with the names of the scheduling routines that are used to cause the indicated transitions to occur).

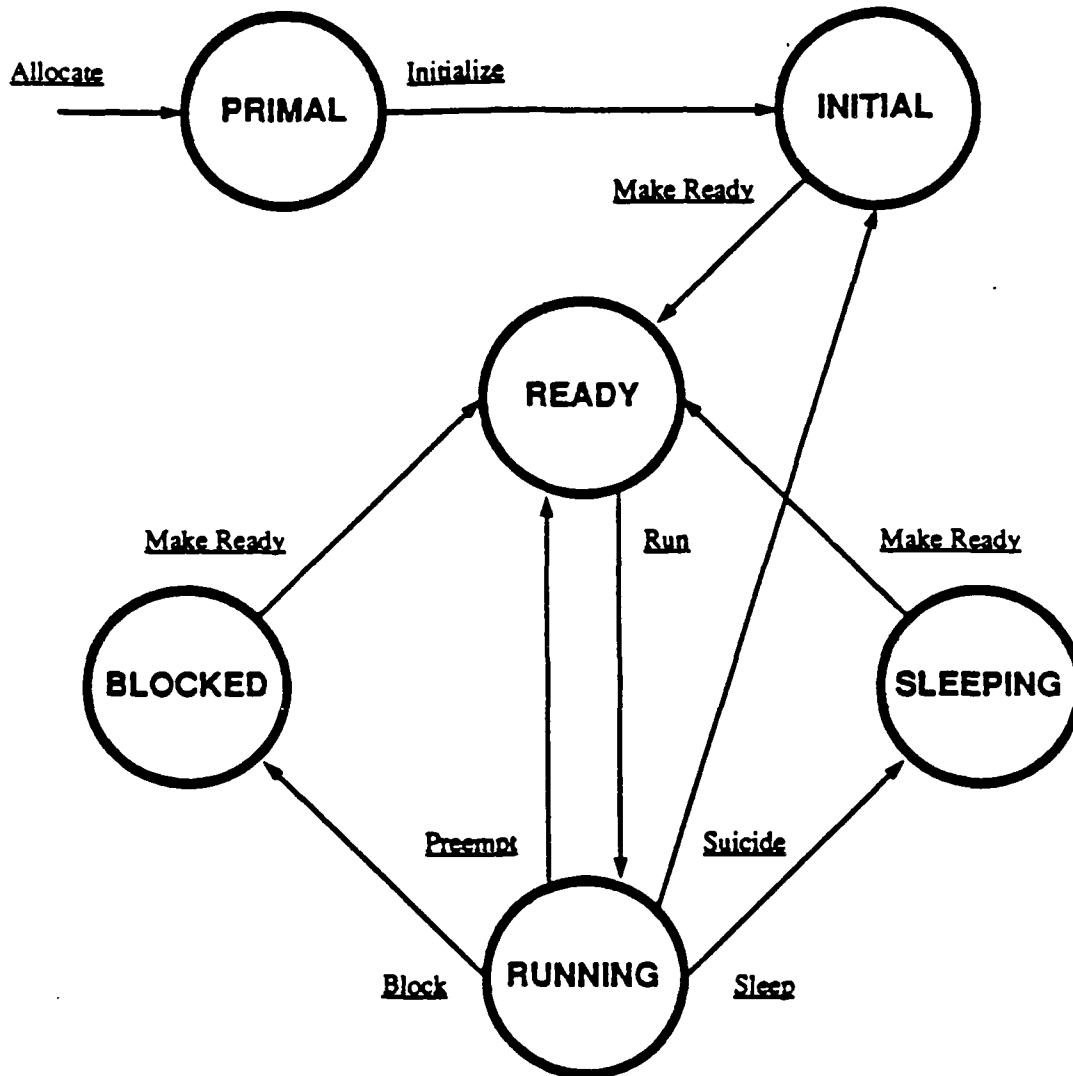


Figure 6. Thread Execution States

Both types of threads may be involved in the invocation of operations on both types of objects (i.e., kernel and client). Each of the four possible combinations have different requirements and characteristics. If a kernel thread invokes an operation on a kernel object, no virtual memory manipulations are required—this type of invocation is the simplest and the most efficient of the four cases. If a client thread makes an invocation on a kernel object, no virtual memory manipulations are required, the client thread just traps from the client thread's context into the kernel—this too is a simple and efficient activity. If a client thread invokes an operation on a client object, the kernel is entered, virtual memory mappings are manipulated to remove the current object from the thread's context, the invoked object is mapped into the thread's context, and the client thread's context is reentered—this is the normal case supported by the kernel, and is moderately complex and moderately costly. If, however, a kernel thread invokes an operation on a client object, a substantial amount of activity has to take place in order to provide the kernel thread with the context necessary to perform the invocation. In order for a kernel

thread to invoke an operation on a client object, it must (temporarily) become a client thread. This activity is straightforward in Alpha, because the design of threads is such that each client thread begins with a kernel thread and assumes the necessary control structures to become a client thread as part of its initialization.

To become a client thread, a kernel thread acquires and initializes a full set of virtual memory management data structures, acquires and initializes a context to execute in, maps into this context the destination object, then enters and begins execution just as if it were a client thread's invocation of a client object. Once the invoked operation on a client object is complete, the client thread's data structures and context are deallocated, and the execution of the kernel thread continues. This type of invocation is the least common, as well as the most complex and costly, of all the types of invocations.

While each client thread is associated with its own separate virtual address space, it does not occupy the full range of virtual addresses defined for its space. The kernel is mapped into a thread's address space, and space is reserved for objects to be mapped into the thread's address space on invocation. The portion of a thread's virtual address space that is occupied by the thread is made up of three major components—i.e., the kernel stack, client stack, and invocation parameter portions. In the current implementation, a thread must be bound to a hardware supported address space (i.e., *context*) in order for the thread to execute. Furthermore, guard pages are placed between each of the major portions of a thread in order to provide a measure of protection from errant pointers and stack over-/underflows.

The thread switches to execute on its kernel stack when it enters the kernel, and the remainder of the time, it executes on its client stack. The thread's client stack is used for the return addresses and local variables associated with the routines that make up an invoked object's operations. The client stack is automatically extended and contracted as the thread's needs changed in the course of its execution.

In Alpha, both capabilities and (simple or structured) data is passed as parameters between objects on operation invocations. In the current implementation, all invocation parameters are placed in the invocation parameter portion of the thread's address space. The invocation parameter passing mechanism takes advantage of the testbed memory management hardware's ability to perform low-cost transfers of physical memory pages among address spaces. The ability to quickly and easily remap physical pages between contexts makes it possible to avoid the high overhead traditionally associated with the copying of parameter blocks among separate address spaces. In addition to mapping invoked objects in and out of the address space of the invoking thread on operation invocation, the kernel must also move parameter blocks between the invoking and invoked objects. In Release 1, a logical stack of invocation parameter pages is maintained within the thread's memory region, with the currently active operation invocation's parameters always on top of the stack. This approach allows the parameter blocks for the currently active invocation to be maintained at a fixed location in a thread's part of an address space, throughout whatever operations the thread may invoke.

The parameters passed into, and removed from, an object as a result of an operation invocation are placed into data structures that exist within virtual memory page boundaries, known as *invocation parameter pages*. Invocation parameter pages are mapped into different locations in a thread's address space in order to effect the passing of param-

eters among objects on operation invocations. Each thread has associated with it a pair of active invocation parameter pages—i.e., an *incoming* page that is shared by the current object and the object that invoked it, and an *outgoing* page that is shared (serially) by the current object and any objects it invokes.

Each invocation parameter page is composed of two main parts—the *request* part that contains the parameters to be passed into the invoked object, and the *reply* part that contains the parameters that are to be passed back to the invoking object when the invoked operation completes. Invocation parameter pages are managed in an *overlapped window* fashion. When an operation is invoked, the incoming parameter page is mapped out of the invoking thread's context, the outgoing parameter page is mapped into the incoming parameter page location for the invoked object, and a new page is allocated and mapped into the invoked object's outgoing parameter page location. When an invocation completes, this process is reversed. The outgoing parameter page is deallocated, the incoming parameter page is mapped back into the outgoing parameter page's location, and the old incoming parameter page is restored. Figure 7 illustrates the behavior of the thread's invocation pages in the course of a (local) operation invocation.

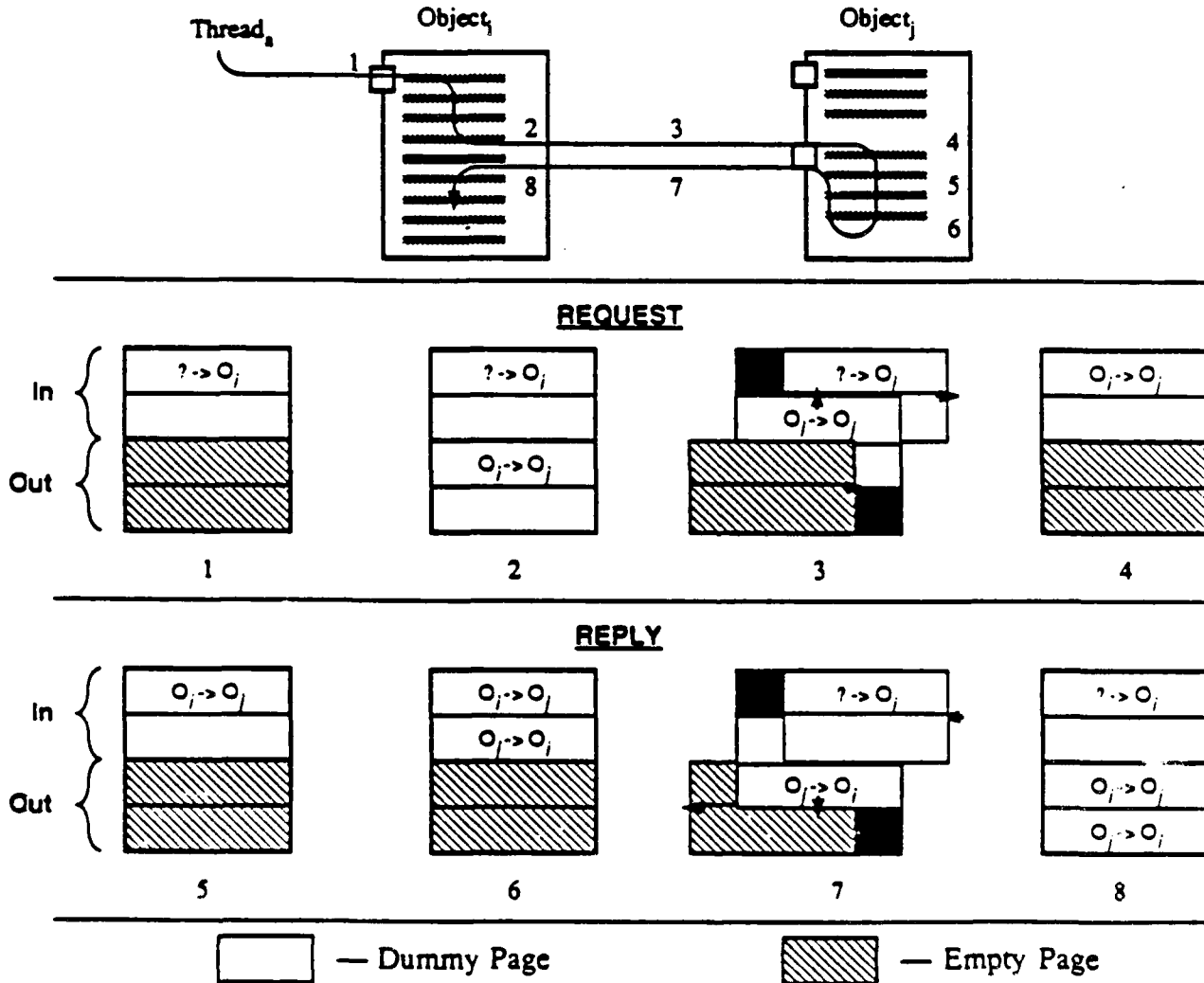


Figure 7: Parameter Page Manipulation on Invocation

3.1.3 Operation Invocations

The operation invocation facility has been implemented to provide a uniform interface to programmers both above and within the kernel. Operation invocations appear the same regardless of the type of object (i.e., client or kernel) making the invocation and the type of object being invoked—i.e., kernel or client object or thread, system service object, replicated object, local or remote.

The object language pre-processor generates the code needed to perform operation invocations. An invocation begins with the marshaling of parameters, which is done by packing all of the passed variables into the invoking object's invocation parameter page, followed by the indices for all of the passed capabilities. The object language pre-processor generates code, at both the source and destination of invocations, to pack and unpack the parameters and capability indices.

Because the current Alpha testbed consists of a collection of homogeneous application processors, there is no need for the transformation of the parameters. However, the object language pre-processor does provide hooks to add such translation routines if the need should arise in the future. If a real object programming language were to be developed for Alpha, the marshaling of parameters could be simplified considerably by allocating the allocating space directly in the parameter pages for the invocation parameters.

When the parameter pages have been loaded with the invocation parameters, the thread traps into the kernel. This causes the application processing element to enter supervisor mode and begin the invocation process. The first thing done by the operation invocation trap handler is the determination of the invocation type, as indicated by information encoded in the target object's Alphasid identifier. There are three different invocation types, which depend on the object that is the destination of the invocation—i.e., *system service*, *local*, or *remote* operation invocations. Figure 8 illustrates the logical structure of a node's operation invocation facility, showing the major components of (and paths through) the facility.

At this point in the operation invocation process, an internal operation invocation interface is used. This interface is provided by the kernel for use both by objects, and the kernel itself, to perform location-independent invocation of operations on a variety of kernel entities (i.e., anything represented by control blocks). This routine takes a pointer to an operation invocation parameter page (in a standard format) and continues with the operation invocation, and is used by all operation invocations from client objects, kernel objects, or functions of the kernel itself. The internal operation invocation routine begins by determining whether the destination of the invocation is a local system service object.

System service operation invocations are used to obtain the services that are provided by the system to the programmer in the form of system-defined objects. For system service invocations, the normal capability translation and validation steps are simplified. This invocation path represents one of a series of "short-circuit"-style optimizations that are provided to compensate for the performance costs of unifying the system's interface by way of the operation invocation facility.

If the destination of an invocation is not a local system service object, a lookup operation is performed on the local Dictionary (see Subsection 3.3.2) to determine if the destination entity is local to this node. If the destination entity's identifier is found in the local

node's Dictionary, the routine that handles invocations on local entities is called. If the destination entity is not found locally, the routine that interacts with the Communication Subsystem (described in Chapter 5) is called to issue a remote invocation. (An exception to this procedure is made in the case of inclusively replicated entities, where the remote invocation is done regardless of whether the destination entity exists locally.)

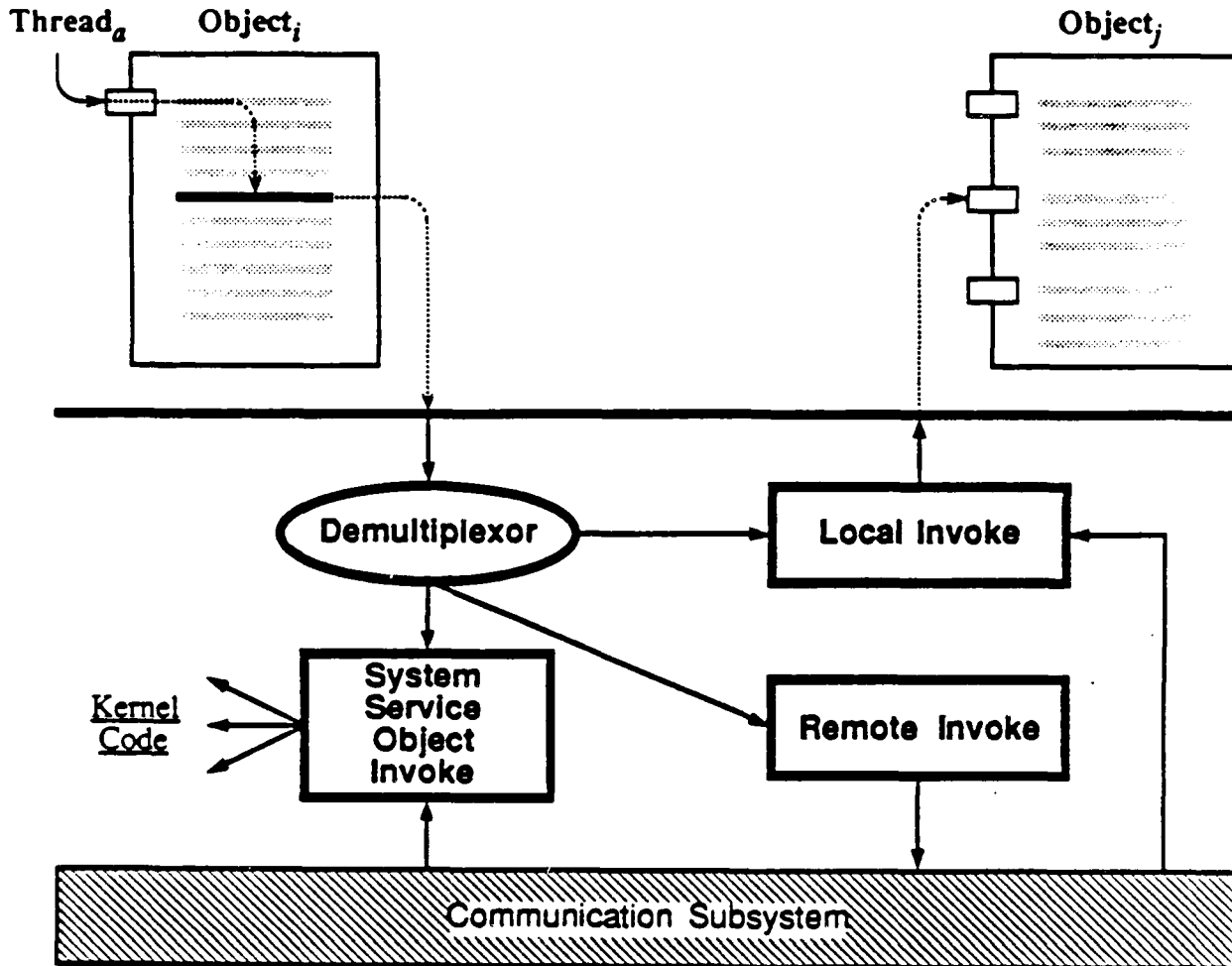


Figure 8: Operation Invocation Facility

In the case of both local and remote operation invocations, the invocation procedure continues with the translation and validation of any capabilities passed as parameters. The use of capabilities is described in detail in Subsection 3.1.4, but briefly, this involves applying the passed capability indices to the invoking object's c-list to obtain the desired *capability descriptors*, and examining them to determine if there are any restrictions applied to the capabilities that would make their desired use illegal. In each invocation there is at least one capability that must be translated and validated—i.e., the capability used to indicate the object that is the destination of the invocation. If a passed capability is invalid (i.e., it is non-existent or its usage is restricted), the invocation is terminated and a failure indication is returned. When all of the capabilities in an invocation are successfully translated, the invocation proceeds by calling the kernel's internal invocation routine.

The code that handles local invocations determines the type of entity on which the operation is to be performed and calls the appropriate routine for each of the potential destination entities. This routine is invoked either as a result of a locally-initiated operation invocation, or as a result of an incoming remote operation invocation (in which case it is initiated by a signal from the Communication Subsystem). Among the routines that handle the local invocations are routines that provide the standard operations on threads and secondary storage objects, in addition to the two main routines that handle invocations on kernel objects and client objects.

The part of the kernel that performs the invocation of operations on kernel objects begins by adding any passed capabilities to the invoked object's c-list, and replacing the capabilities in the invoked object's operation invocation parameter page with their indices relative to the object's c-list. Following this, the entry point of the desired operation is obtained by locating the object's entry point block (as referenced by the object's KOCB), determining if the operation index is within legal bounds, and then indexing into the entry point block to get the address of the operation entry point. Once the invoked operation's entry point has been obtained, the invoked object is entered and execution begun at the start of the desired operation. Note that regardless of the context from which the invocation originated, all operations on kernel objects are performed within the kernel.

The code that is responsible for performing operation invocations on client objects functions similarly to the one that handles kernel object invocations. However, a number of additional activities must be performed in order to initiate an operation on a client object. The kernel routine must first determine whether the operation invocation has a client context associated with it. If not (e.g., if the operation invocation came from a kernel thread in a kernel object), a context must be created before the invocation can be made on the client object; only when the kernel has allocated and initialized a client context, can the invocation precede. When it has ensured that a client context exists, the client object invocation routine locates the virtual memory structures associated with the destination client object, switches to the invoking thread's context, maps the invoking client object out of the invoking thread's context (if there was one), maps the invoked client object into the context, locates and validates the operation's entry point, and begins execution within the invoked object's context. This is done by creating a dummy stack frame, and executing a return from exception instruction to begin execution (in user mode) of the given operation, in the invoked object.

If the destination entity's Alphasid identifier is not found in the local Dictionary, a routine is called that creates a message, puts it in the Communication Subsystem's mailbox, and signals the Communication Co-Processor. This message contains an opcode indicating that this is to be a remote invocation, a pointer to the passed the invocation parameter page, the identifiers of the invoking object and thread, and the invoking thread's execution environment. Once the remote invocation routine passes its request to the Communication Subsystem, the kernel blocks the invoking thread, and continues with the execution of another thread. When the remote invocation completes (successfully or otherwise), the Communication Subsystem places the reply information in shared memory, a response command block is placed in the application processing element's mailbox, and the Communication Subsystem interrupts the application processing element. When the kernel receives such a message from the Communication Subsystem, it locates the invo-

cation reply information, makes the necessary changes to the invoking thread's environment, composes an operation invocation response (to appear as though it came from a local object), and then unblocks the invoking thread. When the invoking thread runs again, the invocation completes in the normal fashion.

Remote invocations create surrogate kernel threads at the destination node, that are obtained from a pool of preallocated resources (as described in Subsection 3.2.1). A surrogate thread inherits the invoking thread's environment and then invokes the operation on the object specified in the incoming parameter page provided by the Communication Subsystem. From that point on, the incoming remote invocation is handled just as if it had originated from a local kernel object and kernel thread. The surrogate thread is added to the remote node's scheduling mix and is scheduled based on its inherited environment information. Once a remote operation invocation completes, the Communication Subsystem is given the results, along with any changes to the thread's environment. This information is passed back to the invoking object, and the surrogate thread is deactivated.

In all of the cases outlined above, once an invoked operation completes, the sequence of steps taken is retraced to return to the invoking object, passing back any results from the invocation. The return parameters are passed back in the reply portion of the parameter block page, and all manipulations of virtual memory are reversed. Figure 9 illustrates the different paths that can be taken through the operation invocation facility.

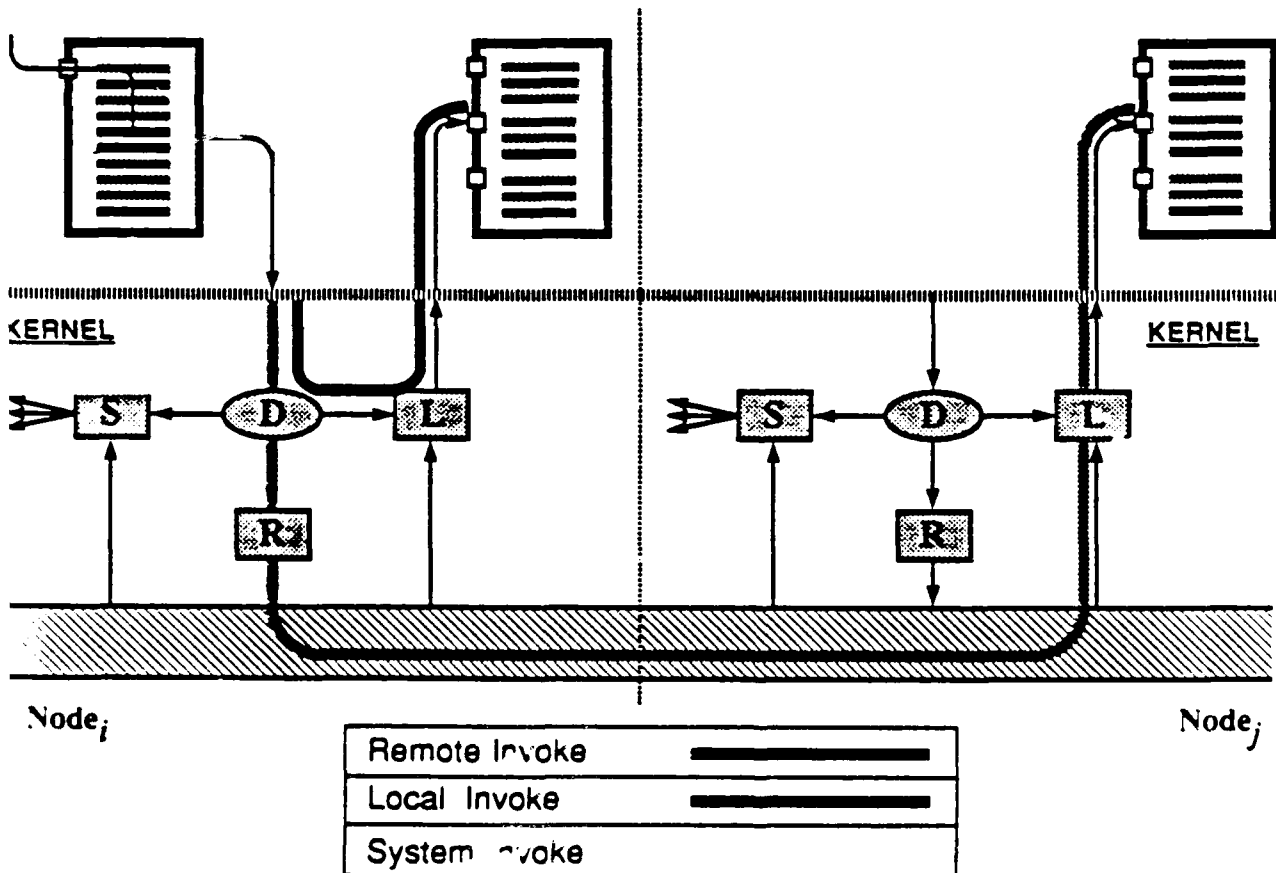


Figure 9: Types of Operation Invocation

3.1.4 Access Control

In some systems (e.g., [Wulf 81] and [Jones 79]), capabilities have a set of *rights* that define the types of operations associated with each capability. In Alpha, the converse is true in that *restrictions* accompany a capability—if a capability is held by an object, that object may perform any operation that is not specifically prohibited by the restrictions associated with the capability.

The manifestations of capabilities local to an object are mapped by the kernel into their internal representations known as *capability descriptors*. A capability descriptor is composed of two components: an object's identifier and the set of restrictions placed on the capability's use. The restrictions associated with a capability are represented by a *restriction vector* containing the per-operation restriction specifications.

An object's capability list (*c-list*) consists of an array of capability descriptors, located in each object's control block, that are managed as a random access list. Since the c-list is part of the object control block, it can be accessed only by the kernel (executing in supervisor mode in the kernel context). Objects do not directly manipulate capabilities, but rather, access their capabilities by way of indices relative to their c-lists. Currently, c-list indices are assigned in sequential, ascending order, although some type of encoding scheme may be employed in the future to reduce the probability of errantly using a valid c-list index.

On operation invocation, capability validation and translation is performed as soon as possible; the destination object capability is checked first, following that, any other capabilities in the parameter list are checked. On return, only the capabilities in the return parameter list must be checked.

To pass a capability as an operation invocation parameter, the kernel takes the c-list index for the destination object and any of the indices of capabilities passed as parameters, and then traps into the kernel. Once in the kernel, the operation invocation facility translates capability indices into their corresponding capability descriptors by looking them up in the client object's c-list. Each capability index is checked for validity during the translation process—i.e., it points to a valid entry in the c-list. Next, the restrictions portion of the capability descriptor is checked to verify that the invoking object has the right to use the capability in the requested manner (i.e., invoke an operation on the destination object, copy the capability, transfer the capability, etc.). Should the validation process fail on any capability, the whole invocation (or invocation return) fails, and a failure indication is passed back to the invoking object. When an invocation fails due to the use of invalid capabilities, the invocation facility does not pass any parameters other than the error indication.

In the course of capability validation, it is determined whether the capability descriptor should be modified. If the capability may only be used once, or if it is being transferred without keeping a copy, the capability is removed after it is validated and translated. If a capability descriptor has had restrictions applied to it to the point of making it unusable, the entire capability is removed from the object's c-list.

When an object is created, a capability is given to the creating object that has only those restrictions which are specified in the object type specification. While capabilities cannot be amplified (i.e., restrictions cannot be removed), restrictions may be added to a

capabilities when they are passed as parameters in operation invocations. This is done by applying the RESTRICT language construct to capabilities in an invocation parameter list. This construct allows the specification of the restrictions that should be added to a particular capability. For each usage of the RESTRICT construct, the object language pre-processor generates a mask that is passed with the parameters to the operation invocation facility and is used in the translation process to modify the restrictions vector of a capability descriptor before passing it to the destination object.

3.1.5 Thread Synchronization

The Alpha kernel provides a set of kernel-defined objects that support the kernel's concurrency control mechanisms—i.e., *semaphore* objects and *lock* objects. Each of these mechanisms are described in turn in the following text.

3.1.5.1 Semaphore Mechanism

The creation of an instance of a *Semaphore* object involves the allocation and initialization of a semaphore data structure. Semaphore data structures consist of a variable containing the semaphore's current state, and the list of threads that are currently blocked waiting on the semaphore. The kernel takes advantage of the fact that capabilities for semaphores cannot be passed to place the semaphore data structures within an object's control block. It similarly takes advantage of the small size of semaphore data structures, and statically allocates a fixed number of semaphore data structures for each object.

When a new instance of a semaphore is created, an unused semaphore data structure is located, initialized with the count value given as a parameter to the create operation, and then a capability is returned to the invoking object. A create operation on the *SemaphoreManager* object returns a failure indication if a semaphore data structure cannot be allocated.

Because *Semaphore* objects are permanently bound to the object that created them, not only can the semaphore data structures be associated with the object's control block, but a special type of capability can be used. To the client, the capabilities for semaphores appear the same as any other ones (i.e., a c-list index), but the capability descriptor does not contain the usual unique identifier and restrictions vector. Rather, it contains an indication of the special nature of the capability and an index into the object's semaphore array for the desired semaphore.

A semaphore data structure contains a count field (that is manipulated by the operations defined on semaphores), a list of the threads that are blocked on the semaphore (waiting for the count to become greater than zero), and a list of the threads that are currently using the semaphore (i.e., threads that have issued P operations on the semaphore and have not yet issued matching V operations).

A P operation decrements the semaphore's count, and should the count become non-positive, the invoking thread is blocked and added to the semaphore's list of blocked threads. A V operation increments the semaphore's count, and should any threads be blocked on the semaphore at this time, a thread from the semaphore's "waiting thread" list is unblocked. The C_P operation first tests if decrementing the semaphore's count would cause it to become nonpositive. If so, the operation returns a parameter indicating this, otherwise the operation is performed as though it were a normal P operation.

Currently, the kernel manages the list of threads blocked on a semaphore in a FIFO fashion. However, semaphores were designed to operate in concert with the Scheduling Subsystem. This was done to allow the selection of appropriate thread to be unblocked to be performed according to the policy used in scheduling the application processor. Work is underway to extend the scheduling subsystem to include the management of blocked threads as part of its functionality.

The kernel maintains, for each semaphore, references to the threads that have P operations outstanding on the semaphore. This is because the kernel must keep track of which threads have outstanding P operations on semaphores, in order to issue matching V operations should a thread encounter an exception (e.g., an aborted atomic transaction, a missed hard-deadline, or as part of the thread repair function) while holding a semaphore. In the course of cleaning up objects following thread exceptions, the kernel checks the list of threads associated with each of the affected objects' semaphores, and for each occurrence of the thread being aborted in a semaphore's list, the kernel may choose to issue a V operation on that semaphore (or not, depending on the mode in which the semaphore is to be used).

3.1.5.2 Lock Mechanism

Lock objects are similar to *Semaphore* objects in their implementation, and instances of *Lock* system service objects are created and accessed in the same manner as *Semaphore* objects. However, each instance of a *Lock* object is associated with a region of data in an object, and each lock has a matching sized area of kernel memory that is used as a write-ahead-log, and is allocated from the kernel's heap region when the lock is instantiated. A collection of data structures similar to the ones used for *Semaphore* objects are associated with the object's control block.

The *Lock* object data structure consists of a field that indicates the mode that the lock is currently in (if it is currently held), a pointer to its log area, and a pointer to the list of threads blocked while waiting to acquire the lock along with an indication of the mode each thread is attempting to acquire the lock in.

When a compatible request is made, the lock is granted the requester, the memory area associated with the lock is copied into its log area, the mode of the lock is changed, and the thread is permitted to proceed. If an incompatible lock request is made, the invoking thread's KTCB is linked into the waiting thread list in the lock's data structure and the lock's request mode is registered with it. When a thread issues an UNLOCK operation on a *Lock* object, the lock mode is restored to its previous state and the blocked thread list is examined.

Like semaphore data structures, those of locks include a list of the threads that currently hold the lock (in any mode) in order to allow the kernel to release all locks held by an aborted thread. The thread to be unblocked when an UNLOCK operation is invoked is chosen from the set of threads in the lock's block list. Currently the kernel selects the first thread with a compatible, outstanding request. However, as with semaphores, the kernel was designed to have the Scheduling Subsystem choose which of the threads waiting in a lock's blocked list should be unblocked.

3.2 Ancillary System Functions

In addition to providing support for the basic abstractions in Alpha, the kernel provides a collections of other services that are needed to round out the functionality of this layer of the system. Included among these functions are features to manage the preallocation of critical system resources, the handling of run-time exceptions, and the initialization of the individual nodes within the system. Each of these functions are defined in the following subsections.

3.2.1 Critical Resource Preallocation

In Alpha, a *daemon* is defined to be a kernel thread and kernel object pair that is provided by the system to perform a system function, independent of any higher-level (system or application) activities. As part of the kernel of Alpha, there are a number of kernel threads and kernel objects that serve to assist in the management of various critical, system resources. The resources managed by daemons include physical memory pages and both client and kernel, thread and object control blocks. Most of the daemons in Alpha serve the same function—i.e., to ensure the supply of a set of preallocated resources in an attempt to speed the expected time for their dynamic allocation and deallocation—and most of the daemons perform their function in a similar fashion.

For each critical resource within the kernel there are a set of routines (generically known as *create* and *delete* routines) that are used for the allocation and deallocation of the desired resources. The *create* routines perform the necessary initialization operations for the desired resource and use routines (known as *get* routines) to acquire the partially initialized, preallocated resources from a pool. Each type of critical resource also has a routine to return these partially initialized resources to their associated pool (these are known as *put* routines).

Should attempts be made to obtain a resource when its associated pool is empty, the requesting object does not wait until the daemon can run and replenish the pool. Instead, whenever pools are exhausted, resources are dynamically allocated via routines that construct them from their basic memory components (known as *alloc* routines). Also, these resources are not completely deallocated (i.e., returned to the kernel's dynamic storage allocator) when they are returned to their pools, but rather the daemons eliminate excess resources when there are too many in a pool. This is done with routines (known as *dealloc* routines), that return the resources to their elemental forms.

3.2.1.1 Physical Memory Management Daemon

The *page daemon* attempts to keep a sma" number of free pages of physical memory in the *free page pool* at all times. This daemon deals with pages of physical memory, and interacts with the Storage Management Subsystem to perform its functions. Whenever the page daemon is activated, it executes the given victim selection algorithm to determine which page(s) should be written out to memory in order to replenish the pool.

The victim selection algorithm used in the current implementation is a simple FIFO-like replacement algorithm. To this end, the kernel maintains a circular list of currently allocated memory pages in order of their allocation. To select a victim, the page daemon begins the searching at the top of the allocated page list (i.e., starting with the page allocated the longest time in the past). The page daemon searches this list in order until it

finds a page that does not have to be written out before being reused (e.g., one that is read-only or has not been modified), or until the search has progressed through N pages in the list. When one of these search termination conditions is met the page pointed to is chosen as a victim, and is then given to the Storage Management Subsystem that is responsible for copying the victim page to backing store. The page daemon continues to select physical pages until enough pages have been selected to bring the free page pool up to its nominal level. The page daemon then blocks until the Storage Management Subsystem indicates that the selected pages have been written to secondary storage, at which point the page daemon adds the pages to the free page pool and blocks again.

The page daemon contains the page replacement policy, and can be easily modified to take the characteristics of the threads that own a page into account when selecting a victim page. In this way, the victim selection algorithm's preference for pages that do not have to be written back to secondary storage could be augmented to also prefer pages of threads with less value to the application. The page replacement algorithm can be based on the same thread environment information (e.g., urgency, importance, or accumulated computation time) used by the Scheduling Facility.

The upper limit associated with the free page pool is not meaningful, as the number of pages in it can never be too large. Thus, the page daemon is unblocked only when the pool's lower limit is crossed. Furthermore, the Storage Management Subsystem is responsible for performing any of the optimizations that might be used to increase paging performance (e.g., caching, managing the physical layout of the disks, or reordering the execution of access requests).

3.2.1.2 Object and Thread Management Daemons

There are four daemons dedicated to managing the pools of preallocated data structures for kernel and client, threads and objects. The resources contained in the pools are known as *skeletons*, and consist of generic control blocks and their associated data structures. Skeletons are initialized with the generic information for each of the types of data structures. When a skeleton is removed from a pool and put into use, the specific information relating to the individual instances of the use of a skeleton is added to it.

These skeletons consist primarily of control blocks, and in the case of client threads or objects, virtual memory information. The major savings in using the preallocated resources derive from not having to dynamically allocate the physical memory space for all of the data structures and from not having to fill in the various fields of the control blocks. In the case of client threads or objects, the saving is somewhat greater in that virtual memory structures need not be obtained and initialized.

3.2.2 Exception Handling

Any system with an application requirements similar to Alpha's demands some form of exception handling facility. These features (i.e., distribution, real-time, and transactions) result in exceptions that are best manifest in the form of asynchronous diversion of the point of program execution control. In Alpha, this means that a mechanism is needed that will allow the programmer to deal with the exceptions which are a fundamental aspect of the system's nature, in a manner that is compatible with the other aspects of the programming model.

Alpha provides an exception handling facility that handles, in a unified fashion, all of the different types of exceptions that may occur. The Alpha exception handling facility is based on the concept of *exception blocks*—i.e., regions of code marked by begin/end operation invocations. The system automatically places an exception block around each operation, time constraint region, and atomic transaction that a programmer generates. In addition, the programmer may explicitly specify exception blocks that consist of arbitrary regions of code.

Each time an exception block is entered, the kernel saves the machine state and bumps the exception block nesting depth counter (both of which are stored in the thread's control block). When an exception block is exited (normally or otherwise), the kernel discards the thread state that was saved at the beginning of the block. The general intent is that when an exception block is entered, the thread's execution state is saved so that should an exception occur, the thread can be vectored back to the start of the block and shunted off to exception handling code. This mechanism behaves much in the same way as the UNIX "longjump" function—the thread is returned to a given point in the code, with its complete data and execution state intact.

The thread's exception block state is managed in a stack-oriented fashion due to the fact that exception blocks are entered and exited in a strictly nested manner. Should an exception occur while a thread is executing within an exception block, the kernel signals the exception, and a depth indication is given with it to indicate how far the thread is to roll back through its exception stack. As a thread rolls back, the previous exception state is discarded, and the new state is used to vector the thread to the desired exception handler code. In addition, a thread's scheduling parameters are popped as necessary while it rolls back through its exception blocks. The thread encountering an exception continues unrolling until its head is positioned to begin executing its normal-flow code at the desired exception level.

Threads include their exception nesting depths as part of their environment—it is passed around with the thread as it moves among nodes, and surrogate threads inherit the exception stack depth (but only maintains its own portion of the exception stack).

In Release 1 of Alpha, exceptions are signalled to threads by jamming an exception return frame onto the thread's kernel stack. This is done so that when the thread attempts to resume its execution, it is not allowed to continue, but is instead forced to execute the kernel's exception handling facility code. In effect, this mechanism inserts the kernel's exception code in-line with the target thread's execution, regardless of where or what state the object is executing in when an exception occurs (e.g., executing within an object, performing a software trap, or handling a hardware interrupt).

In order to jam exception frames onto threads' kernel stacks, it is necessary to know where the top of a stack is when the thread enters the kernel. This is easy to determine in the case of client threads because they execute on their client stack until they trap into the kernel. It is more difficult in the case of kernel threads because they have to execute on their kernel stack at all times. To deal with both of these cases, the system was designed so that a thread's kernel stack frame is marked the first time that it enters the kernel—whether via a trap, interrupt, or procedure call—and it is cleared when the thread leaves the kernel. This is accomplished by using an atomic test and set instruction on

each entry into the kernel in order to determine if it marks an initial transition into the kernel, and if so the thread's kernel stack frame is marked, otherwise nothing is done.

This exception signalling mechanism is idempotent in its effects, and it can therefore be used to repeatedly (e.g., as a result of multiple, "simultaneous" exceptions), and the kernel's exception facility will ensure that the thread will roll back through its exception blocks until it reaches the proper exception level. Furthermore, an exception can be generated at any node, and all affected thread segments will be notified. The exception facilities on separate nodes precede with their operations simultaneously, with the communication subsystem providing the synchronization necessary to ensure that all segment cleanup is complete before the new head of the thread continues its execution.

3.2.3 Initialization

On power-up, all processing elements in a node initialize their hardware resources, perform self-test operations, and then await a load module to begin executing. Typically, each node receives a load module that contains the code and data needed to bring up a replica of the Alpha kernel.

A node's load module is created by the programming language support tools defined in [Shipman 88], and take the form of the example load module shown in Figure 10. The same load module can be used for each node in the system, with per-node differentiation being based on the unique node identifiers and occurring following the initialization of the generic kernel on a node. It is also possible to create node-specific load modules that are tailored to individual nodes (or sets of nodes).

As shown in Figure 10, the canonical Alpha load module consists of a load module header (that contains a first-level loader program and a description of the sub-modules that follow the header), and a set of load sub-modules destined for the node's processing elements. Each of the sub-modules contains the code and data for the various components of the kernel (i.e., the Kernel Proper, the Scheduling Subsystem, the Communication Subsystem, and the Storage Management Subsystem).

As part of the monitor's power-up sequence, the low-core region of the Application Processor's memory is initialized. This region of memory contains the processor's exception vector table, the monitor's data, and an initial stack page for the kernel. The Alpha load module is loaded into the Application Processor's local primary memory above the defined low-core region, and execution begins in the loader program supplied within the load module. The default first-level loader program distributes all of the sub-modules to their intended co-processing element, and relocates the Kernel Proper, and begins execution of the Kernel Proper.

Once the sub-modules are distributed to the co-processors, the initialization of the Kernel Proper begins. The loader relocates the Kernel Proper so that it starts at the second memory segment, and manipulates the memory mapping hardware so that it appears at this location in all contexts. The code and data regions of the Kernel Proper are followed in memory by all of the "wired" (i.e., memory-resident) kernel objects.

The Kernel Proper's initialization begins by modifying the hardware memory translation tables to construct the desired virtual address mappings for the kernel. The kernel's virtual memory initialization function sets aside the physical pages for low-core and the

Kernel Proper's code. These pages are locked in memory (i.e., not available for being paged out) and consist of a contiguous portion of memory at the bottom of the physical memory space that does not participate in virtual memory paging activities. This approach reduces the initial fragmentation of the physical memory space. The remainder of the physical memory is linked together into the *free page list*, from which pages of physical memory can be allocated on demand.

Once this low-level initialization is complete, the Kernel Proper's facilities are initialized and the initialization object is instantiated and the initialization thread begins execution within it. It is up to the initialization thread to instantiate the threads and objects that are to run on the node, based on the information contained in the initialization object.

The co-processing elements in a node follow a power-up initialization sequence similar to that of the Application Processor, and then await the distribution of their sub-load-modules from the Application Processor. The loading of co-processor load modules is accomplished by way of the interprocessor communication facility and the node's shared memory. Each co-processor then begins the execution of the loader program supplied with their load sub-module. In this manner, an entire node is brought up and made ready to begin the execution of application code.

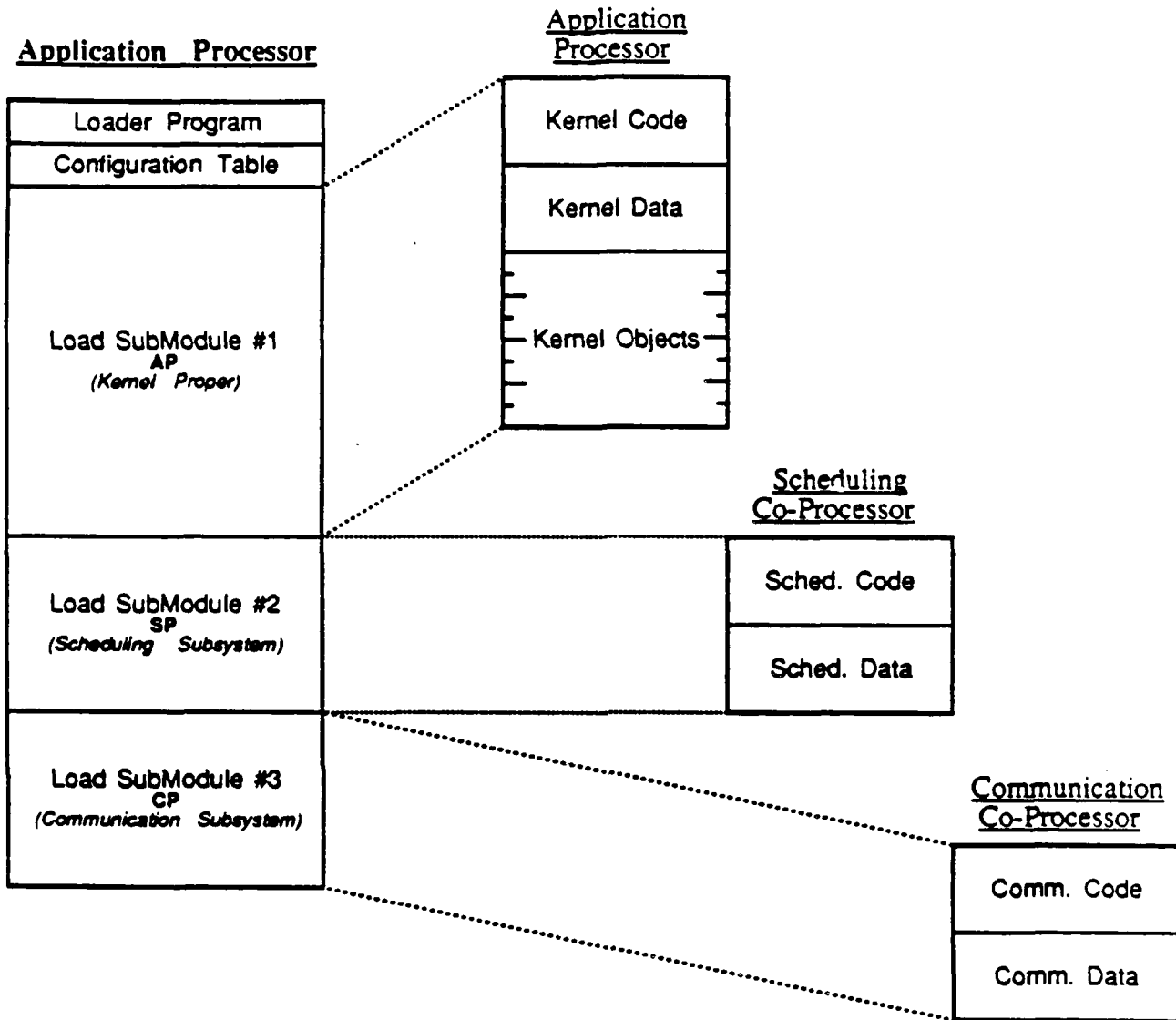


Figure 10: Application Processor Load Module

3.3 Key Data Structures

A large portion of the kernel's data space, as well as a significant portion of the kernel's functionality, is dedicated to the management of a set of internal data structures. These data structures are used in the management of the system's physical resources and in support of the kernel's programming abstractions. These data structures represent the essence of the resources that the kernel manages, and form the nucleus around which the Alpha kernel is implemented.

The fundamental programming entities in Alpha are represented within the kernel by control blocks, all of which are given globally unique identifiers and collected into a directory data structure to allow individual ones to be located by their identifiers. The purpose and contents of each of these data structures are described in the following subsection. In addition to these, a significant portion of the kernel's data structures have to do with virtual memory management. These structures are also described within this section.

3.3.1 Control Blocks

The kernel provides its clients with an object-oriented programming model to enhance the organization and modularity of applications, and the operation invocation facility provides physical location independent access to objects in the distributed system. The benefits that the client derives from these features of the kernel are equally useful for programmers within the kernel. It is frequently the case within the kernel that references are made to entities that exist at remote nodes (e.g., secondary storage images—see Section 6.2). For this reason, the support for distribution provided by the kernel's operation invocation facility is made directly available to the kernel itself by representing all entities as entities for which remote access is required or may be desirable.

To this end, the kernel maintains a variety of control blocks for each such internal, object-like entity in the system. These data structures are known as *control blocks*—i.e., kernel data structures that exist on the same node with the entities they represent. Each kernel-level entity in the system (e.g., a thread or an object) has a corresponding control block maintained by the system. Any entity represented in the kernel with a control block may be accessed uniformly from any node via the kernel's operation invocation facility. This is true whether the entity is stored locally or remotely, and whether the reference originated from within a client object or the kernel itself.

All Alpha control blocks consist of a common header, followed by fields that are specific to the particular entity that is being represented. The control blocks' common header contains: a globally unique identifier used to access the data structure, an indication of the type of entity being represented, a time-stamp that indicates the time at which this structure was last accessed (to be used by the swapping function in selecting entities to be deactivated), and a pointer to another control block (that is used for chaining them together within a common hash table entry).

Each control block has a globally unique, non-reusable identifier in its header known as an *Alphabit identifier*. In the current implementation of Alpha, an Alphabit identifier is constructed using a simple distributed name generation scheme [Lampson 81]. An Alphabit identifier is formed by concatenation of the following four fields:

- an indication of the kind of entity this is (e.g., individual, inclusively replicated, exclusively replicated object, or system service),
- a unique node identifier (provided by the hardware),
- a count that is stored in permanent storage and incremented each time the node crashes (or when explicitly rolled over),
- and a count that is stored within the kernel's primary memory and incremented each time a new identifier is to be generated.

This structure is not visible outside of the kernel, instead they appear to the client as unique, uninterpreted bit strings. The size of Alphabit identifiers was chosen to be large enough to accommodate the number of requests expected during the lifetime of a single execution of the application. An Alphabit identifier size of 48 bits was chosen because it appears to be large enough to perform realistic experiments with, and yet small enough to be used directly as logical addresses on the Ethernet. Also, the non-volatile portion of the Alphabit identifier is provided by a battery backed-up time-of-day clock on the Application Processor board.

The following paragraphs provide a description of each of the different types of control blocks defined in the current release of the Alpha kernel.

3.3.1.1 Kernel Control Block (KCB)

This data structure represents the portions of the kernel that use the virtual memory facility. This control block is used in order to make the kernel appear similar to client objects with respect to its interactions with the virtual memory and secondary storage facilities. Those portions of the kernel required to support virtual memory are permanently kept in primary memory (i.e., the pages are "wired"). The pageable portion of the kernel consists of two memory extents—one for the kernel's heap and another for the kernel's virtual page heap. These regions constitute the great majority of the kernel's memory requirements, making most of the kernel pageable. More detail on the kernel's virtual memory data structures is provided in Subsection 3.3.3.

3.3.1.2 Client Object Type Control Block (COTCB)

This data structure contains the kernel's representation of an object type definition. This structure contains a standard control block header, a print string name for the type description, a flag that indicates if the type is elaborated (i.e., whether an image of it currently exists in primary memory), a file control identifier, and a c-list for all of the well-known capabilities for this type of object.

3.3.1.3 Client Object Control Block (COCB)

This data structure maintains all of the control and status information for a client object. It includes the standard control block header, all of the object's virtual memory information, the identifiers of those threads that are currently active within the object, the capabilities currently in the possession of the object, and the synchronization information for the semaphores and locks associated with the object.

3.3.1.4 Client Thread Control Block (CTCB)

This data structure represents a client thread and includes all of the virtual memory information for the thread, a pointer to the thread's client stack, a pointer to the control block for the object in which this thread is currently active, and a pointer to the control block for the kernel thread associated with this client thread.

Each CTCB has an associated Kernel Thread Control Block (KTCB) which contains the bulk of the information about the thread. All threads begin as kernel threads, and are represented with by KTCBs. When a thread (dynamically) transitions from being a kernel thread to being a client thread, a CTCB is created and attached to the original KTCB. Likewise, if a client thread returns to being a kernel thread, its CTCB is removed and discarded.

3.3.1.5 Kernel Object Control Block (KOCB)

This data structure maintains all of the control and status information for a kernel object. It includes the standard control block header, a pointer to the kernel object's location in the kernel's virtual address space, the current list of the capabilities associated with the kernel object, and the data structures needed for the object's associated semaphores and locks.

3.3.1.6 Kernel Thread Control Block (KTCB)

This data structure provides the internal representation of a kernel thread. In addition to the common header, this control block consists of the thread's identifier, a pointer to its kernel stack, the thread's exception stack, a (mutually exclusive) pair of pointers to either a client thread control block or a kernel object control block, and the thread's *environment*.

In the current implementation, a thread's environment consists of: the current invocation depth of the thread, the thread's current atomic transaction depth, the current state of the thread (e.g., BLOCKED, STOPPED, or RUNNING), the thread's current state with respect to preemption (i.e., whether it is deferred or not), the current number of invocations to roll-back on abort, the thread's global importance, the thread's current time constraint, the shape of the current time constraint's time-value curve, the expected amount of computation time required for the current time constraint block, the amount of time that has elapsed since the start of the current time constraint block, and the amount of processing time that has been accumulated by the computation in its current time constraint.

3.3.1.7 Storage Object Control Block (SOCB)

This data structure maintains all of the control and status information for secondary storage entities (i.e., *instance* and *type* objects). It supports an operation invocation interface between the virtual memory facility and the Secondary Storage Subsystem. These control blocks provide the information necessary to associate an entity that resides in primary memory (e.g., an object or thread), with a secondary storage entity (i.e., a type or instance object).

This control block contains an indication of the kind of entity being represented (e.g., type or instance), the attributes of this entity (e.g., permanent, or atomically-updated), the current state of the entity, and the device-specific addressing information needed to access the entity represented by this control block within secondary storage.

Only those nodes with secondary storage devices will have this type of control block. However, secondary storage objects appear to reside in primary memory, because operations can be invoked on them like any other object.

As a part of the function of the virtual memory facility, pages that exist in primary memory may be moved to their instance objects in secondary storage (and returned on demand). Furthermore, when objects or threads have not been accessed for some period of time, they may be swapped out to their instance objects, leaving behind only a vestigial control block. This is known as *deactivating* an object and is intended mainly for client objects, which can be swapped out to disk and then reconstructed when referenced again. The deactivation function also permits the migration of objects and threads among nodes in the system (by reactivating them on different nodes).

3.3.2 The Dictionary

The *Dictionary* is a global data structure that contains pointers to all of the control blocks in the system at any point in time. The Dictionary is implemented as a partitioned, distributed database, where each node has a partition that contains references to its local control blocks. A local partition of the Dictionary is implemented as a hash table

of pointers to local control blocks, and accessed by Alphabit identifier. As control blocks are instantiated, pointers to them are entered in the local Dictionaries. Once a control block has been found in the Dictionary, a pointer to the desired data structure is returned so that the associated data structures can be manipulated at will.

In support of the form of logical network addressing used in the current release of Alpha, the Communications Subsystem has shared (read-only) access to the Dictionary.

Figure 11 illustrates a local portion of the kernel's Dictionary with instance of various types of control blocks entered into it, and being indexed into by an Alphabit identifier.

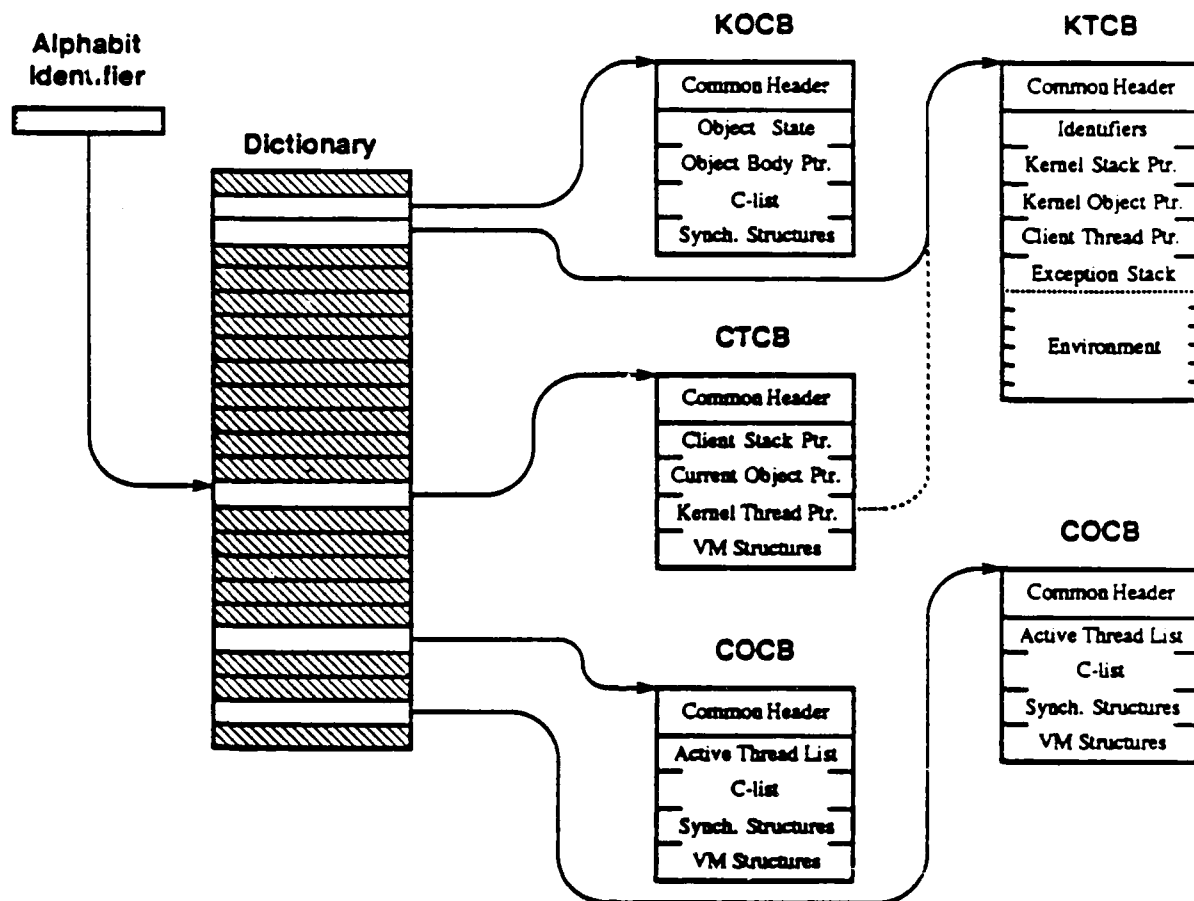


Figure 11: Alpha Dictionary and Control Blocks

3.3.3 Virtual Memory Structures

This Subsection describes the data structures used to support virtual memory in the Alpha kernel. Some of the data structures are provided by the kernel itself, while others are dictated by the underlying hardware. For purposes of modularity and hardware-independence, an effort was made to separate these components in the implementation of Alpha.

This subsection only describes the data structures associated with the system's virtual memory management. The functionality of the virtual memory facility is described in Chapter 7.

Before the virtual memory data structures can be described, it is necessary to be familiar with the memory management hardware found in the Alpha testbed. Following a description of the structures associated with the current testbed's memory management hardware, a description of the additional data structures used to manage entire virtual address spaces (i.e., *contexts*) is given. Finally, the data structures used in the management of individual units of primary memory (i.e., *memory pages*) are described.

4.3.3.1 Memory Management Hardware

Details concerning the Alpha testbed hardware are given in [Northcutt 88c], however the most significant aspects of the Application Processor's memory management structures are summarized here.

The Sun-2 processor board used as the Application Processor in Release 1 of Alpha consists of a MC68010 CPU that issues 24 bit virtual addresses, along with a custom designed Memory Management Unit (MMU) that translates virtual addresses into 24 bit physical addresses and provides a form of access protection.

The structure of the Sun Microsystems version 2.0 MMU is shown in Figure 12. Conceptually, this MMU consists of a two level of indirection look-up table, while at the next level of detail down it consists of:

- a *Context Register* that contains an index into the first of the two translation tables. This index indicates the current address space (i.e., context) that the processor is executing within by referring to a contiguous set of memory segment descriptors.
- a *Segment Map Cache* (SMC) that contains a set of descriptors for a set of memory segments—i.e., contiguous sets of memory pages. Each segment descriptor contains an index into the next-level translation table.
- a *Page Map Cache* (PMC) that contains a set of descriptors for a set of memory pages—i.e., fixed-length, contiguous spans of memory addresses. Each page descriptor contains statistics for the memory page, an indication of the type of page it is, protection information for the page, and an index into the physical memory, pointing to the page of physical memory to which this descriptor refers.

This MMU hardware was designed to provide simple, high-performance virtual address translation, and can be managed in such a fashion as to perform like an address translation cache. Switching between contexts that are loaded in the MMU is an extremely efficient operation (i.e., simply writing to the context register), however only a fixed number of contexts can be loaded into the MMU at a time. Therefore, groups of entries in the MMU tables must be multiplexed among the set of contexts (i.e., threads, in Alpha) that are currently active in the system. Furthermore, all memory protection is concentrated at the page level, and so entries in the SMC do not include any information beyond an index into the PMC. Therefore, a group of entries in the PMC must be reserved for reference by invalid segments (i.e., segments with no pages associated with them).

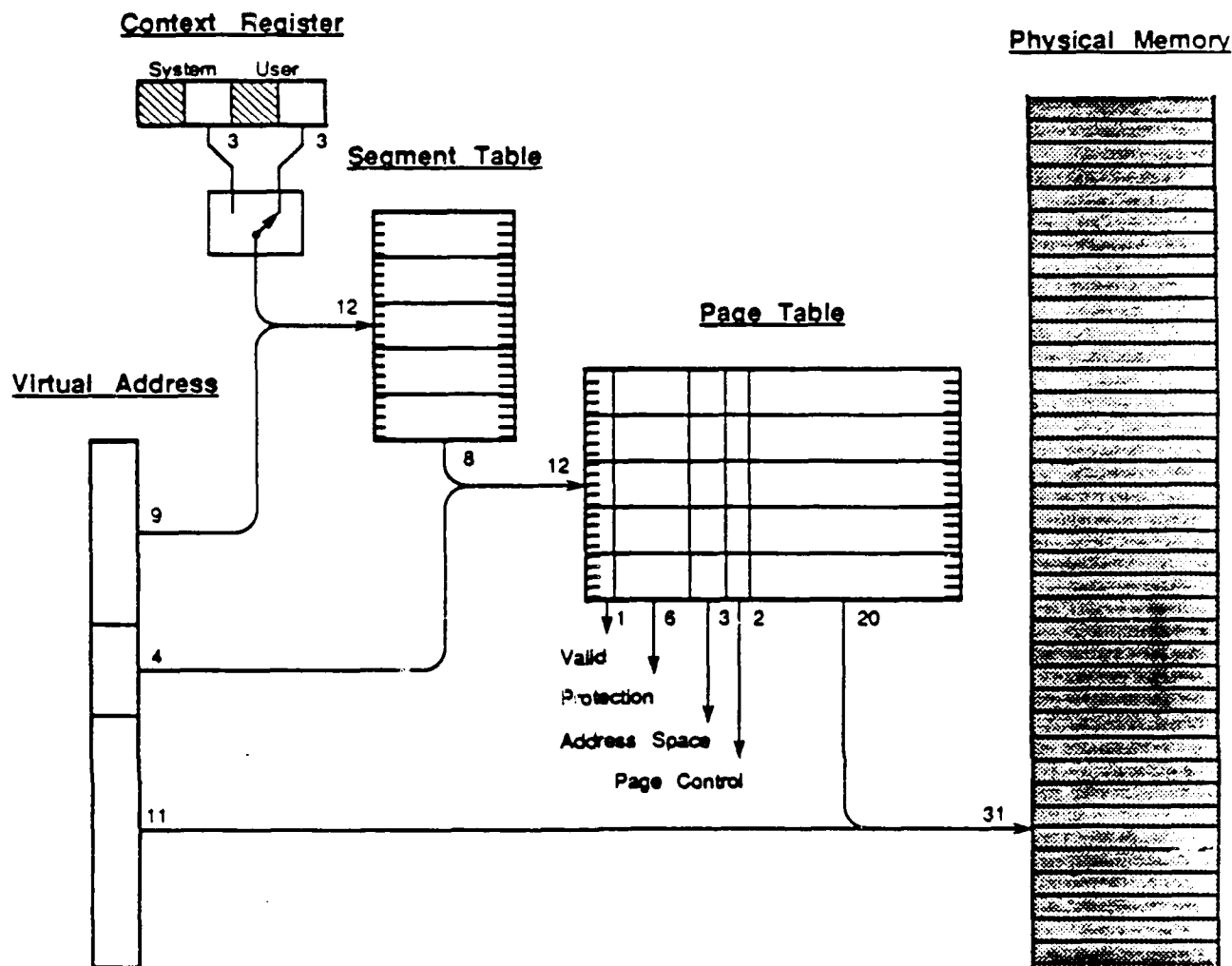


Figure 12: Memory Management Hardware

It should be noted that the Sun-2 MMU provides support for standard, process-oriented programming models (and UNIX in particular), and does not provide any special support for the thread/object model of Alpha. The Alpha programming model might be better supported by an MMU structure similar to that found in the Digital Equipment Corporation's VAX processor line. The VAX MMU manages the caching of translation tables in hardware (via a Translation Lookaside Buffer), and supports a form of partial context swapping (via the SYS, P0, and P1 divisions of contexts) that is compatible with the implementation of operation invocations in Alpha.

3.3.3.2 Context Management Structures

In Alpha, the canonical virtual memory context contains the Kernel Proper, a client thread and a client object. Each of the entities that appear in a context are represented by control blocks within the kernel—i.e., the kernel control block (KCB), a client thread control block (CTCB), and a client object control block (COCB). Each of these control blocks contains a field known as the *segment map image*, that contains the segment

descriptors for the portion of a context that the associated entity occupies. Each segment descriptor contains the current state of the segment (e.g., "invalid" or "loaded in this context"), and the PMC index of the segment's page descriptors (if it is currently loaded). When necessary, the SMC is loaded from the segment map image fields of the appropriate entities. Also, when changes are made to segment descriptors in the SMC, they are also made in the segment map image fields in the affected control blocks.

In addition to the segment map images, these control blocks contain a collection of pointers to the data structures that represent the functionally related collections of memory segments (i.e., *Extents*) that make up an address space. These data structures are known as *Extent Descriptors*, and consist of:

- a pointer to the control block to which the described Extent belongs (this is part of the pointer chain needed to permit the page fault handler to identify the owner of a given page)
- the identifier of the secondary storage control block associated with the Extent,
- the starting and ending virtual addresses of this Extent, relative to its context,
- a pointer to the next Extent belonging to the given control block (i.e., the next contiguous Extent in ascending order of virtual addresses),
- and a pointer to a list of data structures that represent each of the segments that make up the Extent.

Memory Extents and their use in Alpha is discussed further in Chapter 7, while Figure 13 illustrates how they are used in conjunction with other virtual memory data structures.

Sets of linkage pointers assist in the sharing of Extents by objects, and for dealing with the multiple mappings of Extents that occur as the result of threads executing concurrently within a single object. Also, all control blocks referring to a particular Extent are chained together, and a reference count is maintained to determine the number of control blocks that are currently making use of the Extent. The Extent Descriptor also serves to link together the virtual memory data structures of a control block (i.e., the control block itself, its Extent Descriptors, and their segment-level data structures) and the secondary storage images associated with its Extents.

3.3.3.3 Page Management Structures

The data structure that represents individual segments of virtual memory is known as a PMAP, and contains: a pointer to the Extent Descriptor of which this segment is a part, an array that contains the page map descriptors for this segment, another array that contains the state of each page (e.g., RESIDENT, PAGED_OUT, or IN_TRANSIT), and a pointer to the next PMAP in the Extent.

Just as the segment map image fields of the control blocks contain the segment descriptors for the SMC, the page map image fields of the PMAP's contain the page descriptors for the PMC. Similarly, the page map descriptors contained in the PMAP must be loaded into the PMC for any of the pages in that segment to be accessed.

3.3.3.4 Auxiliary Memory Management Structures

Along with the previously described hardware virtual memory structures, the Alpha kernel provides a parallel set of data structures. These kernel-provided *auxiliary* data structures are used to provide additional information for the kernel's virtual memory facility. A decision was made here to trade space for speed, and so additional data structures are provided beyond those logically required, in order to expedite the virtual memory functions of the kernel. The physical structures that have these auxiliary structures associated with them are physical memory pages, the SMC, and the PMC. The *segment map table* is a data structure that parallels the SMC, having an entry for each context's set of segment descriptors. Each segment map table entry contains pointers to the client thread control block and the client object control block currently bound to the context. These pointers provide linkages between a context slot in the SMC and the thread and object that are currently loaded in the slot. In addition, each segment map table entry contains an index to another context slot that is used to link the slots together into *free* and *in-use* context lists. The information maintained in the segment map table entries is used in multiplexing the SMC among all of the client threads that are active at a node.

Similarly, a data structure known as the *page map table* is provided to augment the PMC. Each page map table entry corresponds to a segment's worth of page descriptors in the PMC, and contains a pointer to the PMAP for the segment currently loaded. Each page map table entry also contains an index to permit the free and in-use entries to be linked together by the kernel's virtual memory facility.

A data structure known as the *page state table* parallels the physical memory pages. For each of the application processor's physical memory pages, there exists an entry in the page state table that contains a pointer to the PMAP to which the physical page belongs, and the page's index within the PMAP's page map image. When a physical page is to be manipulated (e.g., for page-out, page-in operations, or on page faults), the page state table provides the information necessary to access the page's virtual memory structures. Additionally, each page state table entry contains an index that is used to link physical pages together for such purposes as maintaining the free page list or the list of pages to be written out to secondary storage.

Figure 13 provides an illustration of each of the virtual memory data structures used in Alpha.

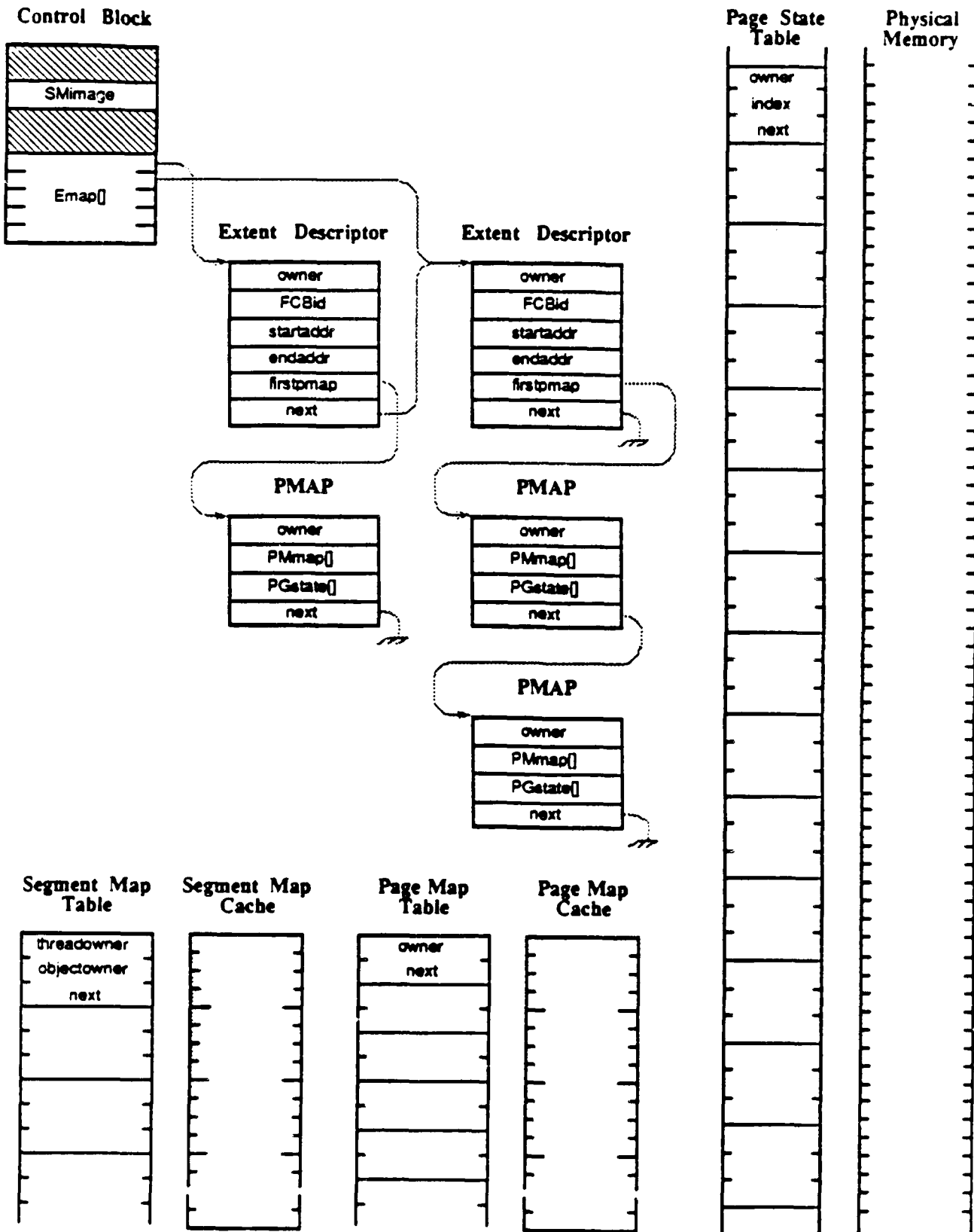


Figure 13. Virtual Memory Data Structures

4 Scheduling Subsystem

The Scheduling Subsystem is one of the principle components of the Alpha kernel, and is related to the other components of Alpha in the manner shown in Figure 2. In the initial release of Alpha, the Scheduling Facility consists of two parts—one that resides in the Kernel Proper and executes on the Application Processor, and another which is known as the Scheduling Subsystem and executes on the Scheduling Co-Processor. These two components execute concurrently and interact via the Interprocessor Message Interface (IMI). The Scheduling Subsystem is different from traditional scheduling facilities in a number of ways—e.g., the Alpha Scheduling Subsystem:

- is cleanly and completely partitioned from the remainder of the kernel by way of a well-defined scheduler interface (enhancing both its potential for parallelism and its modularity),
- is able to operate in parallel with the execution of the application code and Kernel Proper (allowing the execution of computationally demanding scheduling algorithms that would not be practical if the scheduling function were performed in a serial fashion),
- is capable of performing autonomous (i.e., unsolicited) operations on the Kernel Proper (e.g., forcing a context swap—i.e., the suspension of the currently executing thread and the continuation of execution with a different thread),

In the first release of Alpha, the Application Processor is somewhat of a slave, with the Scheduling Subsystem as its master—i.e., the Kernel Proper executes the idle thread and awaits dispatch commands from the Scheduling Subsystem. The Application Processor executes the given thread until another dispatch command arrives from the Scheduling Subsystem, or until the executing thread gives up the processor (i.e., blocks). The Kernel Proper informs the Scheduling Subsystem of all scheduling events which occur on the Application Processor (e.g., when a thread blocks, does a remote invocation, releases a logical synchronization token, or modifies its time constraint attributes). In this configuration, the Kernel Proper is able to dedicate the bulk of the Application Processor's cycles to the execution of application tasks, while the Scheduling Subsystem concurrently carries out the computations necessary to determine how best to manage the Application Processor's cycles.

The following subsections provide the rationale for the Scheduling Subsystem's design, an overview of how the kernel-level programmer makes use of the features provided by the scheduler, a description of the Scheduling Facility's implementation, a discussion concerning the definition of scheduling policies (including an example policy module), and a brief overview of the yet to be completed work on the scheduler.

4.1 Scheduler Requirements

All Application Processor cycles in an Alpha node are intended to be managed by the node's Scheduling Subsystem. The intent here is that the Application Processor execute application code continuously, with the amount of cycles expended on the execution of system code kept to a minimum.

In such a scenario, the scheduler makes all of the decisions concerning the management of Application Processor cycles. In general, resource management decisions must be

made whenever there is contention for a resource. In the case of multiprogramming system, there is almost continuously some form of contention for processor cycles. One form of contention stems from the fact that such a system, by its very nature, multiplexes the Application Processor(s) among multiple software activities. This type of contention persists as long as there is more than one computational activity in the system (which is almost always the case). Another form of contention arises when there are insufficient cycles available to execute the current set of computational activities within their given time constraints. Whereas all systems must deal with the former type of contention, real-time systems must be particularly concerned with the latter type. A scheduler in a real-time system should resolve the first type of contention in such a way as to minimize the second form of contention, and when the latter form of contention occurs the scheduler must be able to resolve this condition in a reasonable way (as defined by the system policy on handling overloads).

Given the great significance of application processing cycle management in real-time systems and the primitive state of understanding that exists in this area, it seemed prudent to design the Alpha Scheduling Subsystem to support experimentation with scheduling policies. This would allow many different scheduling algorithms to be tested and evaluated in the process of gaining an improved understanding of their requirements and characteristics. In addition to this, the Archons project's previous work in the area of scheduling [Jensen 75, Locke 86] has resulted in the generation of a number of requirements for scheduling facilities in operating systems meant to support real-time command and control applications. Included among these are the ability to:

- support a wide range of (possibly unanticipated) scheduling policies, with potentially significant computational demands,
- allow the quick and easy substitution of different scheduling policies into the operating system, transparent to the applications, and
- permit the simple and effective monitoring of the detailed behavior and overall efficacy of a chosen scheduling policy.

In addition to managing application processor cycles, the Scheduling Subsystem should also be involved in the management of other system resources in Alpha. For example, the Storage Management Subsystem might interact with the Scheduling Subsystem in order to decide how best to carry out its function. This is because the Scheduling Subsystem is the repository for all of the system's timeliness information and is optimized for the execution of the policy function. The Scheduling Subsystem contains the time constraint information for each of the computational activities in the system, which includes both the information given by the user and that accumulated by the system at run time. Therefore, the Scheduling Subsystem must be designed to permit other facilities to use its information (concerning timeliness constraints and scheduling policy) to make globally consistent, time-driven resource management decisions.

One function of an operating system is to present an application programmer with a set of abstractions that obscure the unpleasant details of the underlying hardware. In particular, most operating systems permit a programmer to write code with the assumption that it executes alone on a (virtual) hardware base. The more that the realities of a system intrudes on this abstraction (e.g., via interrupts, DMA, etc.), the more difficult the

programmer's job and the harder it is for the system to meet the timeliness demands of its application programs. By centralizing the function of managing application processor cycles, Alpha eliminates the conflicts arising from the disjoint management of cycles by the scheduler and the hardware—cycles are to be given out when the scheduler decides it is appropriate, not when they are requested by the hardware. According to this approach, interrupts that require context swaps to be handled should be routed to the Scheduling Subsystem, allowing it to determine when cycles should be given to the handlers. Taking this concept further, the Scheduling Subsystem should also handle synchronization requests—i.e., when logical synchronization tokens are released (i.e., when V/Unlock operations are performed on semaphore/lock objects) the system should not just select the oldest waiting computation to be allowed to run, but rather it should choose one of the waiting computations based on their current timeliness constraints. This represents one more step towards the unification of all system resource management decisions under a common real-time policy.

The requirements outlined above strongly influenced the design of the Scheduling Subsystem in the Alpha operating system, and a number of features of this subsystem are in direct support of these issues. For example, in order to allow the simple substitution of differing scheduling policies, as well as to facilitate the application of hardware support for scheduling functions, the Scheduling Facility was partitioned from the remainder of the kernel. A clean and well-defined interface has been created between the scheduler and the remainder of Alpha, which has resulted in a number of beneficial side-effects. Included among the benefits of the clear and clean separation of the scheduler from the Kernel Proper are: execution independence—allowing scheduling functions to be performed in parallel with the execution of application activities (in support of the computational demands of real-time scheduling algorithms); a high degree of modularity—permitting the independent modification of the kernel's components (permitting the easy modification of scheduling algorithms); and the codification of the kernel's interface to the scheduler—simplifying the task of monitoring, evaluating, and comparing the behavior of the various algorithms used by the scheduler.

It should be noted that the scheduler is a special operating system facility. Whereas most facilities can be implemented hierarchically (i.e., built out of more primitive, lower-level facilities), the Scheduling Facility is fundamental and presents a cycle that must be broken. The scheduler cannot be implemented as a standard computational unit because it would require that the scheduler be executed in order to run the scheduler—a circularity that must be broken by providing special dispatching services for the scheduling activity.

Because Alpha is a distributed operating system, the local replicas of the scheduler that exist at each node must work together with those at all other nodes in order to ensure the globally consistent management of system-wide processor cycles. This requires that the distributed collection of local/per-node Scheduling Subsystem must be constructed in such a fashion as to share resource management information and a common set of low-level resource management policies.

4.2 Application Time Constraints

The Alpha kernel provide a collection of constructs that permit the programmer to express the timeliness requirements of code regions. These constructs provide the system with the information needed by to perform time-driven management of system resources. The programmer's timeliness needs are expressed to the system in a form known as *time-value functions*, which provide the Scheduling Subsystem with the information necessary to carry out a wide range of scheduling policies. These kernel primitives permit timeliness constraints to be associated with arbitrary regions of code, and in a nested fashion.

The current implementation of Alpha includes time-driven scheduling algorithms that take full advantage of the information provided by time-value functions. Other, more traditional scheduling algorithms make use of a subset of the provided timeliness information.

The distributed nature of the Alpha operating system has a profound impact on the implementation of system timeliness features. This is especially evident with respect to the manner in which time is handled by the system. As is the case with most distributed systems, time must be coordinated among the nodes in Alpha. However, because of Alpha's real-time nature, a special effort must be made to provide nodes with a means of working off a coordinated timebase.

4.2.1 Programming Interface

In Alpha, the programmer indicates to the system his desire that a certain region of code be traversed in a given period of time by enclosing the desired region of code within a pair of begin-/end-time-constraint operations (i.e., via the `BEGIN_TIME_CONSTRAINT` and `END_TIME_CONSTRAINT` standard operations defined on threads).

In order to allow the system to resolve, in a reasonable fashion, situations where the desired time constraints cannot all be met, time constraints in Alpha are defined in terms of time-value functions (which are defined in [Northcutt 88b] and [Shipman 88]). This approach provides the system with application-defined information concerning both the urgency (i.e., the time constraints) and the importance (i.e., the relative value to the system) of regions of code. This permits the system to manage the system resources in such a fashion as to ensure that the maximum number of time constraints can be met given the currently available resources, and when all of the time constraints cannot be met, load is shed in a manner that results in the maximum amount of value being delivered to the system.

Each time a `BEGIN_TIME_CONSTRAINT` operation is invoked, the kernel evaluates the parameters given with the operation and modifies the thread's state to reflect the newly defined constraint. The time constraints associated with a thread are maintained as a part of the thread's state (known as its *attributes*), and consist of an ordered list of the currently active time constraints (i.e., those constraints which have yet to be met or missed). The programmer may use time constraints in a nested fashion, in which case the system manages the thread's time constraint information in a stack-oriented fashion (pushing and popping time constraints as they are initiated and completed in the course of a thread's execution). Threads executing within nested time constraints operate under the most strict of their currently active time constraints—i.e., an operation to begin a time

constraint is effectively ignored should a more demanding time constraint already be in effect. Similarly, when a thread exits a time constraint block (i.e., executes an END_TIME_CONSTRAINT operation or encounters an exception—e.g., fails to meet a hard deadline or aborts a transaction), the kernel determines whether the thread's environment should be altered to reflect a less constraining outer-level time constraint.

An example time value function and its various components is illustrated in Figure 14. The components of a time value function, and their intended use within the system are defined as follows:

- I — **Thread Importance:** the relative value of the computation to system. This value is intended to be used in resolving conflicts during overload conditions (i.e., circumstances where there are not sufficient cycles to meet all of the active thread's time constraints).
- t_c — **Critical Time:** the time at which the execution of the defined block of code should be complete. This value is used as the primary determinant of a code block's time constraint.
- t_e — **Expected Execution Time:** the amount of computation time expected to be required by the code block to complete. This value can be used by the system to predict overload conditions.
- $\sigma(t_e)$ — **Standard Deviation of t_e :** an indication of the degree of variance associated with the given estimate of required execution time. This value is to be used as a measure of confidence in the given execution time estimate.
- S_{pre} — **Shape of Value Function, Pre-Critical Time:** an index into a system-generation-time alterable table of pre-computed value curve shapes. This value is used to select the general shape of the time-value function prior to the critical time.
- S_{post} — **Shape of Value Function, Post-Critical Time:** an index into a system-generation-time alterable table of pre-computed value curve shapes. This value is used to select the general shape of the time-value function following the critical time.
- i_{pre} — **Importance Modifier, Pre-Critical Time:** a multiplier used to scale the portion of the chosen time-value curve prior to the critical time. This value can be used to scale the curve in such a manner as to encode the thread's importance into the time-value curve.
- i_{post} — **Importance Modifier, Post-Critical Time:** a multiplier used to scale the portion of the chosen time-value curve following the critical time. This value can be used to scale the curve in such a manner as to encode the thread's importance into the time-value curve.

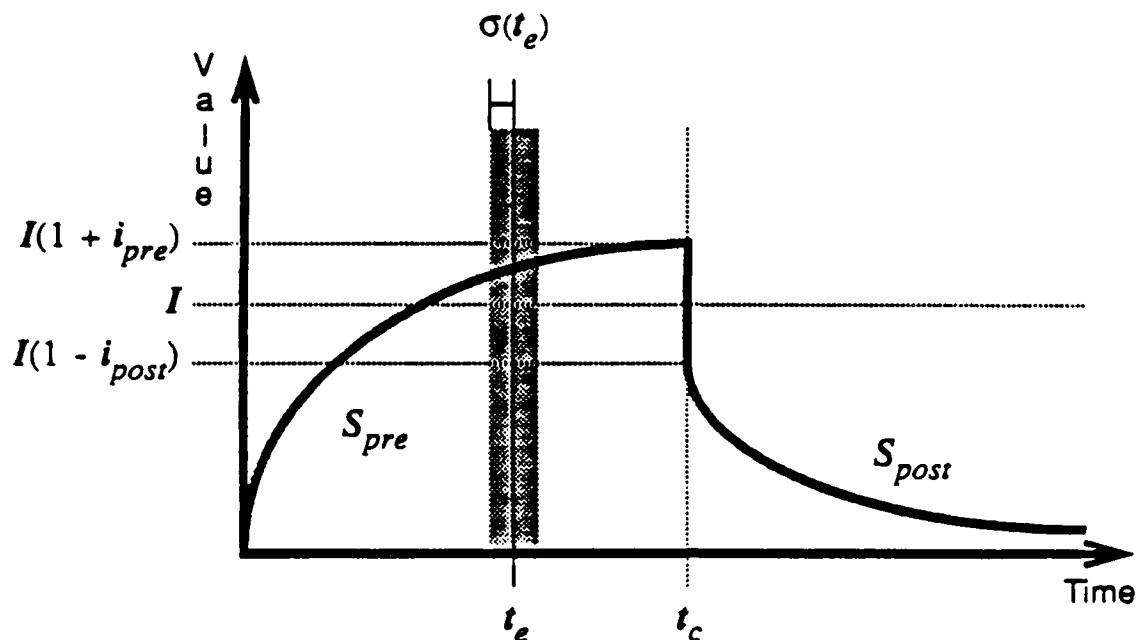


Figure 14: Components of a Time-Value Function

The following time-value function components given as arguments to the `BEGIN_TIME_CONSTRAINT` operation: t_c , t_e , $\sigma(t_e)$, S_{pre} , S_{post} , i_{pre} , i_{post} . The arguments t_c , t_e , and $\sigma(t_e)$ are given in milliseconds, with type `float`. The S_{pre} and S_{post} arguments are indices of system-defined value-function shapes, and are of type `unsigned short`. The shapes of value functions currently defined in Release 1 of Alpha are defined in [Shipman 88]. The arguments i_{pre} and i_{post} are of type `short` and are expressed as percentages of the thread's importance (I), which is determined by invoking an operation on a system service object (as defined in [Northcutt 88d]).

Figure 15 provides an illustration of the use of the kernel's time constraint primitives to implement the Alpha interim object programming language's time constraint block construct. This example illustrates how a time constraint block and its exception handler is implemented—a `BEGIN_TIME_CONSTRAINT` operation is invoked on the thread when the time block is entered, and at the completion of the execution of the code within the time constraint block an `BEGIN_TIME_CONSTRAINT` operation is invoked on the thread. When a time constraint block's time-value function crosses the value origin (or falls to within *epsilon* of zero—where *epsilon* is defined as some small value by the chosen scheduling policy), the thread is said to have missed its *deadline*. When a thread invokes the `BEGIN_TIME_CONSTRAINT` operation, the kernel makes a snapshot of the thread's execution state, and should the thread miss a deadline (or encounter another form of exception) the thread's saved state is restored to the point where the `BEGIN_TIME_CONSTRAINT` operation was invoked. Normally, the `BEGIN_TIME_CONSTRAINT` operation returns a boolean value of `TRUE` when it is invoked, however, when an exception occurs, the operation returns a `FALSE`. In this way, the exception handling code is entered for each of the (potentially nested) exception blocks, with the all of the thread's state available to it for carrying out the desired clean-up functions.

While the thread is executing the code within the body of a time constraint block it will be scheduled according to the time-value function specified by the arguments provided in the `BEGIN_TIME_CONSTRAINT` operation. When a thread exits a time constraint block (either as a result of completing the code within the allocated time, or due to missing a deadline), execution continues in the code following the end of the block, it assumes the time constraint attributes that it had prior to entering the block, and is given application processor cycles appropriately.

```

/* define time constraint block */
TIME CONSTRAINT BLOCK (argument_list) {
    /* block of code to which the time constraint applies */
    --
} ON EXCEPTION {
    /* exception handling code for this time constraint block */
    --
}

```

```

/* time constraints */
if (SELF.BEGIN_TIME_CONSTRAINT(CriticalTime,
                               ExpectedExecutionTime,
                               ExpectedExecutionTimeVariance,
                               PreCriticalTimeShapeIndex,
                               PostCriticalTimeShapeIndex,
                               PreCriticalTimeImportanceModifier,
                               PostCriticalTimeImportanceModifier) == TRUE) {
    /* block of code to which the time constraint applies */
    --
    SELF.END_TIME_CONSTRAINT ();
}
else {
    /* exception handling code for this time constraint block */
    --
    END_EXCEPTION ();
}

```

Figure 15: Time Constraint Block Example

4.2.2 Implementation Issues

When a `BEGIN_TIME_CONSTRAINT` operation is invoked, the critical time given as a parameter is converted into an absolute time within the local node's time-frame. This indication of time is sufficient for the system's needs as long as the thread to which it applies remains on the node where the time constraint was given. However, when threads move among nodes (either as a result of a remote operation invocation, or the migration of a thread/object), the critical time values must be made to remain valid with respect to its local node's real-time clock. In Alpha, maintaining synchronization among clocks in all of the nodes in the system was considered an unreasonable attempt at simulating the behavior of a centralized system. Thus, it was decided that thread time con-

straints must be adjusted to each new local time reference as threads move among nodes. In support of this, the time taken by threads as they move among nodes should be provided to the kernel. Despite the fact that the current communication subnetwork in Alpha (i.e., Ethernet) does not support this functionality, the Communication Subsystem provides reasonable estimates of thread transfer times for use by the kernel. The need for thread transfer time information is example of where a communication subnetwork that is oriented towards real-time applications could provide support for a real-time system. Also, this kind of support is more significant to real-time systems than just the deterministic resolution of contention for communication resources sought after in "real-time" communication subnetworks[†]. The Alpha kernel makes use of measured elapsed times—as opposed to constant time estimates—in many resource management functions (e.g., message transfer times in the Communication Subsystem or device access times in the Storage Management Subsystem). This differs from the approach taken by other real-time systems that attempt to constrain the system's behavior in order to make it conform closely to the designer's static assumptions of execution times. Since real systems cannot be made error-free, and actual delays cannot be kept constant, more adaptive techniques are (in general) better suited for use in practical systems than schemes based on attempts at making *a priori*, absolute guarantees.

Nested time constraints are implemented by the queuing of time constraints on a stack of thread timeliness attributes. In the course of such nesting, less demanding time constraints are ignored by the system, however the time attribute stacks are kept consistent by pushing timeliness attributes when `BEGIN_TIME_CONSTRAINT` operations are invoked, and popping the attributes when `END_TIME_CONSTRAINT` operations are invoked. Furthermore, the kernel ensures that the consistency of thread timeliness attribute stacks is maintained despite thread exceptions (e.g., missed time constraints and thread failures).

Time constraint blocks are standard Alpha exception blocks—they use the same basic mechanisms that save a thread's state at the beginning of a block, restore it if exception occurs before the block is exited (and go to an exception handler), and discarding the thread state when the block is finally exited.

The state associated with a time constraint block includes: the parameters given in `BEGIN_TIME_CONSTRAINT` operation, the run-time statistics accumulated by system, the processor state at point of entry (that can be used for restoring the state of the thread should exceptions occur), and "digested" (i.e., translated/normalized) values passed between nodes on remote operation invocation. Also, the Scheduling Subsystem was implemented in such a manner as to allow the definition of arbitrary time-value curve shapes and assorted probability distributions at system-generation time.

The manner in which both absolute and relative time values are used in Alpha dictates that the need in the Scheduling Subsystem for floating point support. Floating point values are required both for reasons of dynamic range (i.e., very long and very short time constraints) and granularity (i.e., the ability to define times on the desired time boundaries).

[†] Protocol Engine Incorporated's Protocol Engine-based XTP/FDDI communication subnetwork is an example of a network that attempts to meet the needs of real-time systems such as Alpha.

4.3 Scheduling Facility Implementation

The Alpha Scheduling Facility is a major component of the kernel—it not only provides a critical system function (i.e., the distribution of processing cycles to computational tasks), but it is also a facility that interacts strongly with the other kernel-level resource management functions (e.g., remote invocation, synchronization, and storage management). The Scheduling Facility's place in Alpha is illustrated in Figure 2, as is the two main components of the facility. The scheduler consists of a component that exists in the Kernel Proper, and a component that executes independently on the Scheduling Co-Processor and is known as the *Scheduling Subsystem*.

The part of the scheduler within the Kernel Proper provides the system with an interface to the Scheduling Subsystem, where the bulk of the scheduling function is performed. The Kernel Proper's interface to the scheduler is made through a series of routines that control the state of threads.

The routines that comprise the scheduler interface include:

- MakeReady:** Adds a given thread into the scheduler's ready queue, and puts the thread into the READY state.
- Preempt:** Initiates a context swap, putting the currently executing thread into the READY state, and the thread selected from the ready queue by the scheduling algorithm is put into the RUNNING state.
- Run:** Initiates the execution of the thread selected to be placed into the RUNNING state.
- Block:** Removes the currently running thread from the scheduler's ready queue, puts the thread into the BLOCKED state. This routine is typically called as a result of P and LOCK operations.
- Unblock:** Adds a blocked thread to the scheduler's ready queue, forcing the thread from the BLOCKED state back to the READY state. This routine is only called on thread exceptions.
- Sleep:** Performs a Block on the currently executing thread and executes an Unblock on the thread following a specified amount of time.
- Suicide:** Executes a Block on the currently executing thread, deallocates the thread, returning it to the INITIAL state.

These routines provide a well-defined, higher-level interface to the Scheduling Subsystem via the IMI interface between the Application Processor and the Scheduling Co-Processor. All of the Kernel Proper's interactions with the Scheduling Subsystem are handled via this interface. Only simple, serial functions are performed in the Kernel Proper portion of the Scheduling Facility, all other functions that may benefit from the information contained within the scheduling policy modules, or from concurrent execution, are executed within the Scheduling Subsystem. The Kernel Proper part of the Scheduling Facility notifies the Scheduling Subsystem of all scheduling activities (e.g., threads blocking or unblocking and changing timeliness constraints). As a result of a thread relinquishing the processor (i.e., blocking), the Kernel Proper begins the execution of the idle thread and awaits a command from the Scheduling Subsystem to dispatch the next thread to be run. In summary, the Scheduling Subsystem provides concurrency and the bulk of the functionality, while all that the Kernel Proper side does is act as an interface, perform thread dispatch functions when instructed to do so by the Scheduling Subsystem, and executes the idle thread when a thread gives up the CPU (i.e., blocks itself).

The scheduler's clients are not explicitly aware of the Scheduling Subsystem, and it may therefore be modified without having to make changes within the rest of the kernel. The Scheduling Subsystem maintains thread run-time execution statistics instead of the Kernel Proper because of the lack of a fully-shared memory interface between the Application Processor and the Scheduling Co-Processor. Because of this architectural configuration and for reasons of performance, the Scheduling Subsystem duplicates some information maintained by the Kernel Proper.

The interface between the Scheduling Facility's two components is detailed in the following section.

4.3.1 Internal Interface

The interface between the scheduler's parts is significant in that it must be designed in such a fashion as to be transparent to the Kernel-Proper-side of the scheduler, while allowing the implementation of any desired policy within the Scheduling Subsystem. This requires that both the commands and the information exchanged via this interface be sufficient and appropriate to support the needs of a wide range of potential scheduling policies. It is believed that the interface in the current implementation of Alpha meets these needs, and the following provides a definition of the logical interface and details of the actual interface.

The interface to the Scheduling Subsystem was designed so that the Scheduling Facility can perform its function external to the application processor that it is managing. The interface is composed of a collection of commands that are exchanged via a node's inter-processor communication mechanism (i.e., the kernel's IMI service).

In the course of system operation, the Scheduling Subsystem may issue a command to the Kernel Proper, indicating that the currently executing thread should be preempted and a new thread is to be dispatched. The preemption command includes a reference to the Alphabit identifier of the new thread to be executed. The Kernel Proper responds to such a preemption command by suspending the execution of the currently active thread, performing a context swap to allow the execution of the indicated thread, and issuing an acknowledgment to the Scheduling Subsystem indicating that the thread has been suspended. When the Scheduling Subsystem receives a preemption acknowledgment, it adjusts the accumulated computation time for the newly suspended thread, and continues with its scheduling functions while the Application Processor executes the newly dispatched thread.

The scheduler's internal interface also includes commands that allow the Scheduling Subsystem to be notified of changes in a thread's state. These commands are used to indicate when a thread should be added to, or removed from, the scheduler's list of currently active threads (i.e., the *ready queue*) and include information concerning the reason for the indicated manipulations—e.g., a thread has blocked on a semaphore or a lock, a thread has been suspended because of a virtual memory fault, a thread has been frozen or deleted, or a thread has been blocked while it makes an invocation on a remote object. There is also a command that allows the Kernel Proper to inform the Scheduling Subsystem of changes that occur in the attribute information of one of the threads currently in the *ready queue*.

The Scheduling Subsystem's interface has also been designed to permit the Scheduling Facility to execute entirely within the Kernel Proper, should the chosen scheduling algorithm not be able to take advantage of the concurrency provided by the Scheduling Subsystem (such as might be the case with simple round-robin or priority scheduling schemes). Additionally, the Scheduling Subsystem can use the thread run-time information that it maintains to determine which of the waiting threads should be unblocked when a V operation is performed on a semaphore or an UNLOCK operation is performed on a lock.

The actual interface between the Scheduling Facility's components is illustrated in Figure 16.

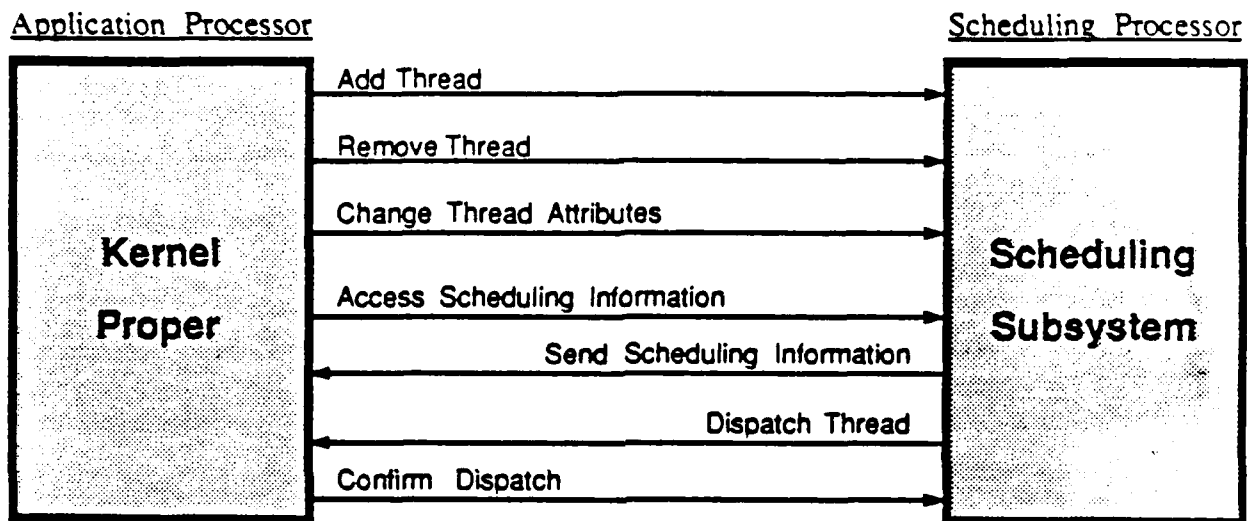


Figure 16: Application Processor/Scheduling Subsystem Interface

The commands sent from the Kernel Proper to the Scheduling Subsystem are:

- **Add Thread:** This command takes as parameters the Alphabit identifier for the thread to be added to the scheduler's ready queue, along with the necessary scheduling information, and sends a message containing these parameters to the Scheduling Co-Processor.

There are three variations of this command—the first being the addition of a simple thread to the ready queue (i.e., due to a local thread unblocking), the second being the addition of a surrogate thread (i.e., a thread created locally to perform a remote invocation from another node), and the third being an add command that does not take effect until after a specified amount of time has elapsed (i.e., a delayed add).

The first time that a thread is added to the ready queue, its scheduling information is passed to the Scheduling Subsystem with the add command. The Scheduling Subsystem retains this information until the thread is destroyed, therefore subsequent add commands need only include the thread's Alphabit identifier.

The add operations performed on surrogate threads take as parameters scheduling information in the Scheduling Subsystem's internal format. This is because

the thread has already been executed (albeit on another node) and has accumulated statistics. Also, the thread's scheduling information has already been given to the Scheduling Subsystem and converted from the format it was provided from the user in, to the format that the subsystem maintains.

- **Remove Thread:** This command takes as a parameter the Alphabit identifier of a thread to be removed from the ready queue, and sends it in a message to the Scheduling Subsystem.

There are two varieties of this command—one that simply removes a given thread from the scheduler's ready queue, and another that removes the thread from the ready queue and deletes the scheduling information associated with it from the Scheduling Subsystem's storage.

- **Change Thread Attributes:** This command alters the time constraint information associated with a thread by the Scheduling Subsystem. There are three types of change commands—one to push new time constraints on a thread's time constraint stack, a second command to pop them off, and yet another command to change the global scheduling attributes of a thread.

The command to push time constraints takes as parameters a thread's identifier and the time constraint information as given by the client, and passes them to the Scheduling Subsystem in messages. While the command to pop time constraints takes as a parameter a count indicating the depth to which the thread's time constraint stack is to be popped to. Both commands take a parameter that indicates the thread which is to be changed.

The command to change a thread's global scheduling attributes takes an identifier for a thread and the new scheduling attributes to be applied to the thread (e.g., the thread's importance), and passes them in a message to the Scheduling Subsystem where they are applied to the given thread's control block.

- **Access Scheduling Information:** There are two different varieties of this command—one to load, and another to store, the Scheduling Subsystem's information related to a thread's current state (i.e., run-time statistics, normalized/transformed scheduling parameters, top of time constraint stack, etc.).

Both types of this command pass as parameters a thread's Alphabit identifier, and a pointer to a shared memory location (outside of the IMI channel's space) where the scheduling information is to be loaded from or stored to.

- **Confirm Dispatch:** This command is a simple message to the Scheduling Subsystem indicating that the requested thread dispatch has been performed. It requires no parameters, and is used to synchronize the Kernel Proper and Scheduling Subsystem. Such synchronization is necessary because the currently running thread may choose to release the processor at the same time that the Scheduling Subsystem issues a preempt command.

This command also allows the Scheduling Subsystem to maintain accurate run-time statistics for the threads. Which is necessary because the time required to perform a preemption may vary—e.g., preemption may be briefly deferred for kernel synchronization purposes, or context swaps may involve differing amounts of effort depending on the type of threads being swapped (i.e., kernel or client).

The commands issued by the Scheduling Subsystem to the Kernel Proper are:

- **Dispatch Thread:** This command sends the Alphabit identifier of the thread that the Scheduling Subsystem's policy determines should be bound to an Application Processor. This command instructs the Kernel Proper to suspend the currently executing thread and run the thread given as a parameter to the command.

A confirmation message is expected from the kernel for this command.

- **Send Scheduling Information:** This command is used to signal that the requested scheduling information has been stored into the buffer indicated in a preceding command to access the scheduling information for a thread.

In effect, this command serves as a synchronization message that permits the Kernel Proper to use the scheduling information that it has requested (e.g., a thread's scheduling information for shipment to another thread on a remote invocation).

As an optimization, the components of the Scheduling Facility exchange local control block pointers, in addition to the globally unique Alphabit identifiers used to indicate the thread which is the target of a particular operation. This enhances the scheduler's performance by simplifying the identification of thread control information on both sides of the interface.

4.3.2 Internal Structure

A primary design objective of the Scheduling Facility in Alpha was the ability to support a wide range of (possibly experimental) real-time scheduling policies. This objective was achieved in Release 1 of Alpha by way of a logical framework for the Scheduling Subsystem within which a wide range of differing scheduling policies can be implemented, through the substitution of policy-specific routines within a well-defined structure. The Scheduling Subsystem framework consists of a collection of scheduling mechanisms and the interface definition for a collection of routines that implement various aspects of scheduling policies. In the current implementation of the Scheduling Subsystem, multiple policies can be implemented by inserting a suite of policy routines into the scheduling mechanism framework, and then selecting the desired policy (by way of link- or run-time settable switches).

The overall design of the Scheduling Subsystem is one in which all activity is performed by policy routines, in response to the receipt of *tokens*. Each scheduling policy is composed of a series of actions that are performed in response to certain (internally or externally generated) stimuli. In this design, all hardware interrupts (which stem from either the scheduling co-processor's on-board devices or messages from the Application Processor) and inter-policy-routine communications are manifest as tokens, each of which is directed to a specific target policy routine. The receipt of a token by a policy routine triggers its execution, the result of which may be the modification of the scheduler's state or the generation of additional tokens.

Figure 17 provides an illustration of the logical structure of the Scheduling Subsystem in Release 1 of Alpha.

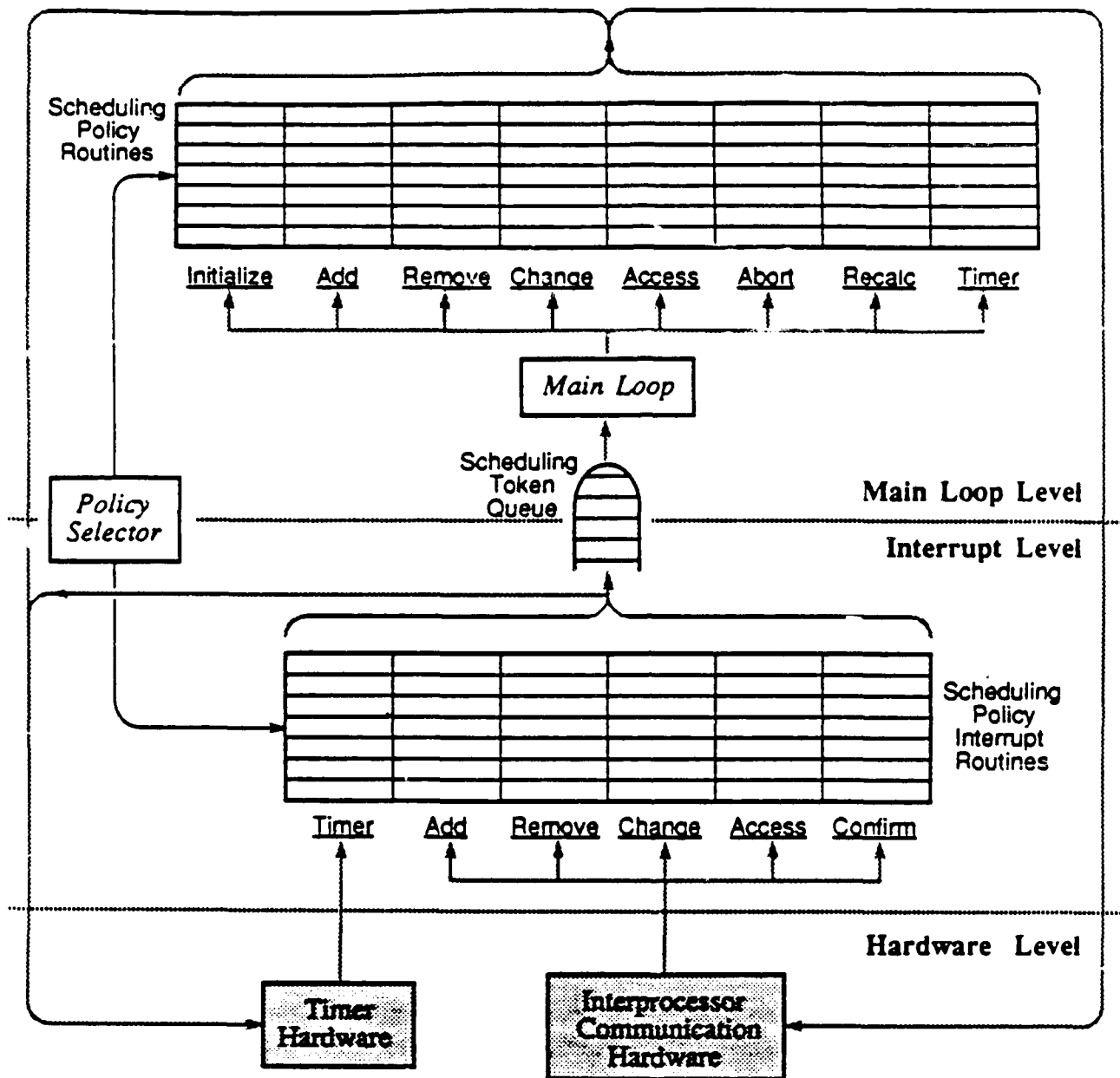


Figure 17: Scheduling Subsystem Internal Structure

The Scheduling Subsystem is composed of two levels—the *interrupt level* and the *main loop level*—each of which is composed of a collection of routines that carry out specific aspects of a total scheduling policy. As its name implies, the interrupt level contains routines that handle the tokens generated (by the underlying framework) as a result of hardware interrupts. The main loop level contains routines that handle tokens that were generated by other routines (at either the main loop or interrupt levels). The substrate upon which the Scheduling Subsystem is constructed provides the basic services of converting interrupts into tokens, managing the queueing and distribution of tokens, and controlling the execution of the routines specified by the targets of tokens.

The Scheduling Subsystem's framework provides: the overall structure to scheduling policies (by defining the token handling routines and their general function); the token handler indirection tables (used to map scheduling tokens to the proper routines for the chosen policy); the interrupt-token generation code (for the IMI interface from the Application Processor and the on-board interval timer); scheduling token manipulation routines (e.g., allocation/deallocation, and scheduling token queue insertion/removal mechanisms); the main loop token interpreter (that removes tokens from the queue and dispatches them to the appropriate token handler routine); a collection of general purpose monitoring and debug mechanisms, hardware management routines (e.g., memory management, timer management, and low-level I/O); the IMI interprocessor communication package (for interacting with the Applications Processor); and a set of stand-alone math libraries (providing support for floating point and probability distribution manipulations).

In order to simplify the synchronization concerns of policy writers, an effort was made in the implementation of the Scheduling Subsystem's mechanisms to ensure that the proper (i.e., FIFO) sequencing of tokens within the main loop queue is maintained in the face of multiple, asynchronous insertions and removals of tokens. Also, default handler routines are placed into indirection table slots where no such scheduling token is expected for the given policy. These default routines signal an exception condition and can optionally be made to generate a system fatal error. In addition to these routines, there can also be a set of custom debug routines defined for each policy. These debug routines can be used by systems developers to print out the (formatted) contents of (both general and policy-specific) data structure, dynamically collected run-time statistics, and other scheduling-policy-specific information.

The following text describes the mechanisms provided by the Scheduling Subsystem's framework and each of the its token handling routines, starting with the interrupt level ones and followed by the main loop level routines.

4.3.2.1 Scheduling Subsystem Mechanisms

A collection of general-purpose mechanisms are provided as part of the Scheduling Subsystem's framework. These mechanisms are provided for use by the routines that make up individual scheduling policies, and are grouped into functionally related sets known as *facilities*. The Scheduling Subsystem's facilities provide mechanisms for the management of memory, time, tokens, and interprocessor communication. The support mechanisms provided by each of these Scheduling Subsystem facilities are defined as follows[†]:

Memory Management Facility

- Dynamic Memory Allocator

INTERFACE: vaddr Falloc()

PARAMETERS: none

RETURNS: Pointer to location in memory of newly allocated block, or a 'NULL' pointer if there is insufficient free memory to honor the request.

[†]The various initialization routines associated with the facilities are omitted here for reasons of brevity.

DESCRIPTION: This routine returns a pointer to a fixed size block of memory. The block size is equal to that of a data structure that is the union of all of the data structures used by scheduling policies. The allocator was designed, above all else, to be fast.

- Dynamic Memory Deallocator

INTERFACE: void Fdealloc(vaddr fptr)

PARAMETERS: The fptr parameter points to the block of memory to be deallocated.

RETURNS: nothing

DESCRIPTION: This routine takes a pointer to a previously allocated block of memory and adds it to the allocator's free list. For performance reasons, no check is made to see if the pointer refers to a valid, currently allocated block of memory.

Time Management Facility

- Create a Timed Event

INTERFACE: short CreateTimedEvent(u_long interval, tokenptr tknptr)

PARAMETERS: The interval parameter indicating an amount of time in terms of scheduling processor clock ticks. The tknptr parameter is a pointer to a scheduling event token.

RETURNS: An identifier for the timed event, or a 'NIL' value if it is unable to allocate a timed event.

DESCRIPTION: This routine takes an indication of the amount of time until a timed event is to occur, and a pointer to the token that is to be issued when the given interval of time has elapsed.

- Cancel a Timed Event

INTERFACE: void CancelTimedEvent(u_short event)

PARAMETERS: The event parameter is an identifier for a timed event.

RETURNS: nothing

DESCRIPTION: This routine deletes a timed event indicated by the identifier for the timed event that was returned when it was created. This routine has no effect if the indicated timed event has already expired when it is invoked.

- Read the System Clock

INTERFACE: u_int GetSystemTime()

PARAMETERS: none

RETURNS: Integer representing current time.

DESCRIPTION: This routine returns a value that indicates the current time in terms of system clock ticks. One tick is equal to 3.2 microseconds.

Token Management Facility

• Create a Scheduling Event Token

INTERFACE: tokenptr CreateToken(tokentype type, tokenargs args)

PARAMETERS: The type parameter is an enumerated type defining the types of tokens that exist. The args parameter is a union of the data structures that comprise the data part of the possible scheduling tokens.

RETURNS: Pointer to the newly created and initialized scheduling event token.

DESCRIPTION: This routine takes an indication of the desired token's type and the data associated with the token, allocates the memory for a token, initializes the token's body with the given data, and then loads the current system time into the token's header.

• Delete a Scheduling Event Token

INTERFACE: void DeleteToken(tokenptr tknptr)

PARAMETERS: The tknptr parameter is a pointer to a scheduling token.

RETURNS: nothing

DESCRIPTION: This routine takes a pointer to a token and deallocates it, returning its memory to the free pool.

• Enqueue Scheduling Event Token

INTERFACE: void EnqueueToken(tokenptr tknptr)

PARAMETERS: The tknptr parameter is a pointer to a scheduling token.

RETURNS: nothing

DESCRIPTION: This routine inserts a token into the tail of the scheduling event token queue, to be processed by the main loop.

Interprocessor Communication Facility

• Send a Message to the Application Processor

INTERFACE: void SendMsg(short opcode, message msg)

PARAMETERS: The opcode parameter is a constant indicating the type of parameter to be transmitted. The msg parameter is a union of the types of messages that can be sent, and it contains the text of the message.

RETURNS: nothing

DESCRIPTION: This routine takes an indication of the type of message to be sent and the body of the message as parameters. It composes a message and transmits it to the Application Processor via the Interprocessor Message Interface.

4.3.2.2 Interrupt Level Routines

Whenever an interrupt occurs at the SP, a token is generated and passed to an interrupt handling routine (as indicated by the policy indirection tables). The scheduler's substrate handles both the generation of these tokens, therefore there is room for alternative interpretations of the function of these routines (e.g., the Timer routine is only invoked as an interrupt generated by the on-board timer).

The functional descriptions for the Scheduling Subsystem's interrupt level scheduling policy routines are defined as follows:

- **Timer:** This routine receives tokens from the on-board counter/timer unit, which is configured in the Scheduling Subsystem to function as an interval timer. The tokens that this routine receives were generated by other routines at an earlier point in time, and inserted into the timed event queue maintained by the interval timer facility.
- **Add:** This routine receives tokens from the IMI facility, which are generated as a result of receiving `Add_Thread` commands from the Kernel Proper. If a delayed add is desired, this routine can generate a timed event (reusing the token that it receives). Otherwise, the token can be inserted into the scheduling event queue for processing at the main loop level.
- **Remove:** This routine receives tokens from the IMI facility, which are generated as a result of receiving `Remove_Thread` commands from the Kernel Proper. Generally, the received token is inserted into the scheduling event queue for processing at the main loop level.
- **Change:** This routine receives tokens from the IMI facility that indicate the desire to change the scheduling information associated with a thread. The typical function of this routine is to enqueue the received token into the scheduling event queue.
- **Access:** This routine receives tokens from the IMI facility that indicate the desire to load or store the currently relevant scheduling parameters associated with a thread. The typical function of this routine is to enqueue the received token into the scheduling event queue.
- **Confirm:** This routine handles the preemption confirmation commands that are returned by the Application Processor in response to commands issued to it earlier. This routine provides the functions required to maintain the synchronization between the components of the Scheduling Facility in the Kernel Proper and the Scheduling Subsystem. Typically, all that this routine does is update the Scheduling Subsystem's state to indicate the identity of the currently executing thread, and update the run time statistics being kept for the thread that was preempted. In order to help detect when the running thread has released the Applications Processor at the same time that a Dispatch command is sent, a distinction is made by the Kernel Proper between preemptions of the idle thread, and all other threads. With this information, synchronization can be maintained among the scheduler parts.

In addition to the example functions defined for each of these routines, the desired statistic gathering activities may also be performed in these routines. Information gathered in this way is typically used by the debug routines added as part of a scheduling policy. It should also be noted that any of these routines can be a *null* routine—i.e., a fatal error is generated should this type of interrupt occur. Furthermore, the scheduling token that initiated the execution of an interrupt routine may either be deallocated (when the interrupt routine handles the event entirely at interrupt level), or reassigned and inserted into the scheduling event queue (to be handled at the main loop level). An scheduling token may alternatively be placed into the interval timer queue, to appear later in the form of a token passed to the specified routine. This action has the effect of delaying the arrival of a token for a specified period of time.

4.3.2.3 Main Loop Level Routines

Tokens are queued up in the scheduling event queue by the interrupt handling routines and the main-loop token policy routines. The function of the main loop is to repeatedly remove tokens from the queue, dispatch them to the appropriate policy routines, and initiate the execution of the target routines.

The functional descriptions for the Scheduling Subsystem's main loop level scheduling policy routines are defined as follows:

- **Initialize:** This routine is intended to provide the initialization functions for the specific policy being used. On system start-up, the Scheduling Subsystem's framework generates a token targeted to this routine. This routine is typically executed once at system start-up time, and is not executed again. The functions that may be performed by this routine include the allocation and initialization of data structures, the creation of initial timed events (e.g., for initiating time-slicing activities), and signalling the Scheduling Subsystem's readiness to the Kernel Proper. It should be noted that the Scheduling Subsystem's mechanisms are initialized at system start-up time, therefore the policy module's routines need only initialize the policy-specific aspects of the subsystem.
- **Add:** The intent of this routine is to add a given thread to the set of active threads maintained by the Scheduling Subsystem. For convenience, this collection of threads is known as the *ready mix*, regardless of the actual manner in which the scheduling policy manages the set of threads. In addition to simply adding a new thread to the ready mix, this routine is frequently involved in the process of ordering the threads in a *ready queue*. When a thread is added to the ready mix, the ready queue may be reordered followed by a check to determine if a currently running thread should be halted and a different thread dispatched. Alternatively, this routine may just add a thread to the ready mix, and then generate a token to perform the ready queue reordering and thread dispatch check.
- **Remove:** This routine is used to remove a thread from the ready mix. This is done when a running thread gives up the processor (e.g., is blocked or destroyed), and may require that the ready queue be reordered and reevaluated to determine the next thread to be dispatched. When the running thread blocks, the Scheduling Subsystem may choose to retain the control blocks allocated for the thread and the scheduling information associated with it. This would be

done in order to expedite the process of unblocking the thread. In the event that the running thread commits suicide, the token received by this routine indicates this, and the routine may choose to deallocate the state associated with the thread within the Scheduling Subsystem.

- **Change:** This routine is to be used to modify the state information kept by the Scheduling Subsystem for the threads it is aware of. The types of modifications that can be performed are the changing of the thread's state (e.g., its importance value), the pushing of new time constraints, and the popping of old time constraints. The changes made to a thread's scheduling state by this operation may require the reordering of the ready queue and the dispatching of new threads.
- **Access:** This routine is intended to be used to load and store the collected internal state of threads. This function is needed to support the movement of thread scheduling state information among nodes in the distributed system. When a thread leaves a node on a remote invocation, the thread's scheduling state is obtained and passed with the thread. At the remote node, a surrogate thread is created and initialized by installing the passed scheduling state into the remote node's Scheduling Subsystem.
- **Abort:** The intended use of this routine is in dealing with an exception condition that relates to the scheduling state of a thread in the ready mix. This routine is typically executed as a result of a token generated by a timed event that signals the fact that a thread's critical time has passed. For example, this routine may handle the case where a thread misses a hard deadline by popping the thread's time constraints, reestablishing the thread's previous time constraints, and initiating a reordering and reevaluation of the ready queue. In addition, this routine may signal the Kernel Proper that a thread has failed to meet a deadline, which may in turn cause the Communications Subsystem to initiate a thread abort sequence among the nodes involved with the thread.
- **Recalc:** This routine is intended to be used to reorder and reevaluate the ready mix in order to determine whether a new thread should be dispatched. This routine is typically executed following a scheduling event that modifies the state of the ready queue, or when a timed event indicates a potentially significant point in time has arrived. This routine is intended to be used when the reordering function cannot be performed in a reasonable amount of time by another routine. With this routine, the primary activities of the other routines can be carried out without the additional performance penalty of a (potentially significant) ready queue reordering activity.
- **Timer:** This routine can be used to perform functions associated with timed events in the main loop. This routine is typically not used, as most timer-related activities can be performed at the interrupt level.

4.4 Scheduling Policies

Scheduling policy modules are composed of a collection of routines and data structures that implement the desired scheduling policy within the framework of the Scheduling Sub-

system's mechanisms. While the previous subsections described the intended function of each of the different types of routines that make up a policy module, there is no a priori restriction on what functions the routines can perform, nor are there any constraints on their size or organization. The scheduling routines can make use of the mechanisms provided by the underlying subsystem, and can be arbitrarily complex in both structure and behavior.

The following text provides a general description of the scheduling policy routines, and an example of an actual policy module.

4.4.1 Policy Definition Modules

A scheduling policy is defined by a pair of arrays of routines (that are inserted into the interrupt- and main-loop-level policy indirection tables), and a set of data structures (that are used within and among the policy routines).

A complete scheduling policy consists of a set of interrupt handling routines that respond to each type of interrupt that might occur in the course of this policy's execution (e.g., a message arrives from the Application Processor, or the interval timer indicates that a given period of time has elapsed). An interrupt handling routine can (at interrupt level) perform the functions required by the policy in response to a particular interrupt, it can simply generate a token and queue it for processing at the main loop, or an interrupt handling routine can perform some part of the desired function at interrupt level and pass the rest off to be handled by the main loop. The general intent is that interrupt handlers perform a minimum amount of necessary function, but rather just queue up a token and allow the main loop's token handling routines to perform the bulk of the work necessary for a given interrupt.

At any one time, the Scheduling Subsystem may contain multiple, complete scheduling policies. This is due to the fact that the routine indirection tables that provide the structure for the creation of scheduling policies, are multidimensional—i.e., there are multiple complete sets of token handling routine slots in each table. A particular policy is selected from the current set by way of a *policy selector* variable which serves as an index into the arrays of indirection tables. The system's entire scheduling policy can be changed at run-time by altering the value of the policy selector.

4.4.2 Policy Example

The following is pseudo-code description of the data structures and routines that make up the policy module for an enhanced dynamic deadline scheduling policy.

4.4.2.1 Data Structures

Thread Control Blocks: per-thread data structures that act as repositories for the information that the scheduler needs to know about a thread. They are created when a thread is first introduced to the scheduler, and are deleted when the thread is killed.

Figure 18 illustrates the structure and contents of a thread control block in Alpha.

AP Thread Control Block Pointer
Thread Alphabit Identifier
Thread Importance
Link to Time Constraint Stack
Time Constraint Depth Count
Blocked Flag
Total Accumulated Computation Time
Last Dispatch Time
Importance Queue Header
Deadline Queue Header

Figure 18: Thread Control Block

Time Constraint Blocks: scheduling parameter control blocks that contain information concerning a thread's time constraints. A time constraint block is created each time a thread enters a time constraint block, and is destroyed when the thread exits the block (either normally or as a result of an exception). Each thread control block contains a list of the time constraint blocks for the thread, which is managed in a stack-oriented fashion.

Time Constraint Stack Header
Time at which Command was Issued
This Constraint's Critical Time
Mean Expected Computation Time
Compute Time Standard Deviation
Previously Available Slack Time
Computation Time Left
Computation Time Received
Beginning Exception Block Level
Timed Event Descriptor

Figure 19: Time Constraint Block

Scheduling Queues: linked list of the control blocks for currently active threads. There are two types of queues in this policy module. The first type of thread queue is known as the *deadline* queue, it contains control blocks for ready threads that currently have time constraints and is kept sorted in shortest-time-constraint-first order. The second type of thread queue is known as the *importance* queue and contains the control blocks of the ready threads that currently do not have time constraints associated with them. Furthermore, the thread control blocks in the importance queue are kept sorted in a greatest-importance-first order.

Scheduling Statistics: a collection of counters and accumulators that keep track of the number of times that an event of interest (e.g., a call to a policy routine) occurs and the accumulated sizes of the thread control block queues. This information is used by the debugging code to print the number of times each type of scheduling event occurs, and to print the average and maximum lengths of the scheduling queues.

4.4.2.2 Interrupt Level Routines

Timer: handles interrupts generated by the interval timer facility. There are two reasons why interval timers are used—delayed Add commands and missed hard deadlines. This is to say that either Add or Abort tokens are queued up for the main loop by this routine.

```
DLTimerIntr() {
    /* add a time-stamp to the token */
    TokenPtr->CurrentTime = GetSystemTime();

    /* queue it up in the event queue */
    EnqueueToken(TokenPtr);

    /* do not dealloc the token yet */
    return(NO_DEALLOC);
}
```

Add: handles interrupts from the Kernel Proper part of the Scheduling Facility that indicate the desire to (either immediately or following a delay) add a thread to the ready mix. If it is a delayed command, then this routine passes the incoming token to the interval timer to be delayed and then issued at a later time. If it is an immediate command, the token is queued up to be handled by the main loop.

```
DLAddIntr() {
    /* check if it is a delayed command */
    if (token is delayed) {
        /* it is, so send it to the interval timer */
        CreateTimedEvent(TokenPtr->DelayInterval, TokenPtr);
    }
    else {
        /* it is not, so queue it up in the event queue */
        EnqueueToken(TokenPtr);
    }

    /* do not dealloc the token yet */
    return(NO_DEALLOC);
}
```

Remove: handles interrupts from the Kernel Proper part of the Scheduling Facility that indicate the desire to remove a thread from the ready mix. This routine just queues up the token to be handled at the main loop level.

```
DLRemoveIntr() {
    /* queue the token up in the event queue */
    EnqueueToken(TokenPtr);

    /* do not dealloc the token yet */
    return(NO_DEALLOC);
}
```

Change: handles interrupts from the Kernel Proper part of the Scheduling Facility that indicate the desire to change a thread's scheduling parameters. This routine just queues up the token to be handled at the main loop level.

```
DLChangeIntr() {
    /* queue the token up in the event queue */
    EnqueueToken(TokenPtr);

    /* do not dealloc the token yet */
    return(NO_DEALLOC);
}
```

Access: handles interrupts from the Kernel Proper part of the Scheduling Facility that indicate a desire to remove a thread from the ready mix. This routine just queues up the token to be handled at the main loop level.

```
DLAcessIntr() {
    /* queue the token up in the event queue */
    EnqueueToken(TokenPtr);

    /* do not dealloc the token yet */
    return(NO_DEALLOC);
}
```

Confirm: handles interrupts from the Kernel Proper which indicate that the requested preemption of the currently running thread has completed. The time is noted when a thread is dispatched, so that when a preempt confirmation command arrives, the amount of computation time received by the thread can be determined. This allows the run-time statistics for the preempted thread to be accurately maintained.

```

DLConfirmIntr() {
    /* check if the preempted thread is not the idle thread */
    If (not preempting idle thread) {
        /* not idling, so update the preempted thread's stats */
        ...
    }

    /* dealloc the token */
    return(DEALLOC);
}

```

4.4.2.3 Main Loop Level Routines

Initialize: initializes the policy module. This includes the creation and initialization of two queues that contain the threads in the ready mix (i.e., the deadline and importance queues), it sets the current thread pointer and preempting thread pointer to NULL, it indicates that a preemption is currently not in progress, and it initializes the timed event facility.

```

DLInit() {
    /* initialize ready queues */
    QInit(IM_Queue);
    QInit(DL_Queue);

    /* initialize global variables */
    currentthreadptr = NULL;
    preemptthreadptr = NULL;
    preemptinprogress = FALSE;

    /* initialize the timed event facility */
    InitTimedEvents();
}

```

Add: adds a thread to the ready mix. This command's token can come directly from a "MakeReady" command in the Kernel Proper, or from the timed event facility as a result of a delayed Add command. This routine can add to the ready mix threads rooted at the local node, or threads rooted at another node (i.e., surrogate threads). If the thread added is locally rooted, it may have already be introduced to the Scheduling Subsystem, or it may be new to the scheduler (i.e., either a newly created thread rooted at the local node, or a surrogate thread). When a thread is first introduced to the Scheduling Subsystem, a thread control block is created and initialized with the thread's initial parameters or those passed with it by the remote node.

If a surrogate thread is under a time constraint, a time constraint control block must also be allocated and initialized with the passed information, and a timed event is created for the thread's currently dominant deadline. If the thread to be added is already known to the scheduler, a pointer to its thread control block is given to the Scheduling Subsys-

tem and the control block is added to the appropriate ready list. Threads are added to the importance or deadline queues (depending on whether they have active time constraints or not) in such a way as to maintain the desired ordering. Finally, a new thread is dispatched if necessary, and the incoming token is deallocated.

```

DLAdd() {
    /* check if this is a new thread */
    if (new thread) {
        /* it is, so allocate and initialize a TCB */
        TCBptr = (tcb)Falloc();
        ...

        /* see if it is a surrogate (with passed scheduling info.) */
        if (surrogate thread) {
            /* it is, so inherit all of its scheduling info */
            ...
        }
    }

    /* prepare the TCB to run */
    TCBptr->blocked = FALSE;

    /* decide if it belongs in the importance queue */
    if (no time constraint) {
        /* no TC, so add it to the importance queue */
        Qinsert(IM_Queue, TokenPtr);
    }
    else {
        /* it has a TC, so add it to the deadline queue */
        Qinsert(DL_Queue, TokenPtr);
    }

    /* check if we need to dispatch a new thread */
    if (new thread is at head of queue) {
        /* the new thread is at the head of the queue */
        DispatchThread();
    }

    /* deallocate the token */
    DeleteToken(TokenPtr);
}

```

Remove: removes a thread from the ready mix. This command can only come from the currently running thread on the Application Processor. There are two cases to be considered—one where the current thread has decided to block, in which case it is just removed from the ready mix and the next thread is dispatched. The other case is where the thread blocks at (approximately) the same time that the Scheduling Subsystem issues a preempt command, in which case the Scheduling Subsystem's state must be kept consistent with that of the Kernel Proper.

In any event, the given thread is removed from the ready mix, the blocked thread's run-time statistics are updated, and the next ready thread is dispatched. Also, there are two variants of this command—one where the indicated thread control block is removed from the ready queue it is in (i.e., blocking the thread by the Kernel Proper's "Block" command), and another where the control block is removed from the ready mix, then deallocated (i.e., the thread is killed by the "Suicide" command in the Kernel Proper). The control blocks for blocked threads are allowed to float in the Scheduling Co-Processor's memory, while references to them are kept in the Application Processor.

```
DLRemove() {
    /* update the current thread's accumulated computation time */
    ...

    /* remove the thread from the ready mix */
    If (under time constraint) {
        /* remove thread from deadline queue */
        QRemove(DL_Queue, TokenPtr);
    }
    else {
        /* remove thread from importance queue */
        QRemove(IM_Queue, TokenPtr);
    }

    /* mark thread as blocked */
    TCBptr->BlockedFlag = TRUE;

    /* check if currently running thread is the one being blocked */
    If (TCBptr is CurrentThreadPtr) {
        /* it is, so dispatch a new thread */
        DispatchThread();
    }

    /* check if the thread control block is to be deallocated */
    If (kill thread) {
        /* it is, so get rid of time constraint stack */
        while (time constraint blocks left) {
            /* deallocate a time constraint block */
            Fdealloc(TCptr);
        }

        /* deallocate the thread control block */
        Fdealloc(TCBptr);
    }

    /* deallocate the token */
    DeleteToken(TokenPtr);
}
```

Change: handles the pushing and popping of a thread's time constraints as well as the changing of the thread's importance. When the thread's scheduling information is changed, its control block is removed from whatever ready queue it might be in, changing the state, and then reinserting it into the appropriate queue (in the proper order). If after the thread's state has been changed and it has been reinserted into a ready queue, a dispatch command may have to be issued to run another thread. When time constraints are pushed, new time constraint control blocks are allocated, initialized with the given time constraint information, and pushed onto the head of the time constraint stack. When a pop time constraint is done, the block at the top of the time constraint stack is removed and deallocated.

```
DLChange() {
    /* remove the thread from the ready mix */
    if (under time constraint) {
        /* remove thread from deadline queue */
        QRemove(DL_Queue, TokenPtr);
    }
    else {
        /* remove thread from importance queue */
        QRemove(IM_Queue, TokenPtr);
    }

    /* decide what kind of change is to be made */
    switch (TokenPtr->SubType) {
        case CHANGE: {
            /* change thread importance */
            ...
        }
        case PUSH_TC: {
            /* push a time constraint */
            TCptr = (tc)Falloc();
            ...
        }
        case POP_TC: {
            /* pop a time constraint */
            ...
            Fdealloc(TCptr);
        }
        ...
    }

    /* reinsert the thread into the appropriate queue and location */
    if (no time constraint) {
        /* no TC, so add it to the importance queue */
        QInsert(IM_Queue, TokenPtr);
    }
    else {
        /* it has a TC, so add it to the deadline queue */
        QInsert(DL_Queue, TokenPtr);
    }

    /* check if we need to dispatch a new thread */
    if (different thread at head of queue) {
        /* the new thread is at the head of the queue */
        DispatchThread();
    }

    /* deallocate the token */
    DeleteToken(TokenPtr);
}
```

Access: allows access to a thread's scheduling information from the Kernel Proper. There are two types of accesses that can be performed as part of this command—i.e., scheduling information reads and writes. For read scheduling information commands, the Scheduling Subsystem places the information for the given thread into the buffer location in common memory indicated in the command from the Kernel Proper. Similarly, writing scheduling information into a thread's control block is done from the shared memory buffer referred to in the command.

```
DLAccessSI() {
    /* check if it is a read command */
    if (dump subcommand) {
        /* it is a read, so copy the information into the buffer */
        ...

        /* send a message to the AP to indicate the dump is done */
        SendMsg(SI_DUMP_DONE, NULL);
    }
    else {
        /* it is a write, so copy the information out of the buffer */
        ...

        /* see if the thread's importance has changed */
        if (importance changed) {
            /* adjust the thread's location in the IM queue */
            QRemove(IM_Queue, TokenPtr);
            QInsert(IM_Queue, TokenPtr);

            /* check if we need to dispatch a new thread */
            if (different thread at head of queue) {
                /* new thread at the head of the queue */
                DispatchThread();
            }
        }

        /* send a message to indicate that the update is done */
        SendMsg(SI_UPDATE_DONE, NULL);
    }

    /* deallocate the token */
    Fdealloc(TokenPtr);
}
```

Abort: this is performed when a thread's hard time constraint expires. It pops time constraint blocks until the desired exception level is reached, propagates computation time statistics forward, and deallocates the popped time constraint blocks. In addition, this routine sends message to Application Processor in order to abort other segments of the thread to the desired level. If no time constraint blocks are left, the thread is moved to the importance queue. In any event, a new thread is dispatched when it is necessary.

```

DLAbort() {
    /* adjust the thread control block's state */
    ...

    /* it is, so get rid of time constraint stack */
    while (not at desired exception level) {
        /* propagate scheduling information */
        ...

        /* deallocate a time constraint block */
        Fdealloc(TCptr);
    }

    /* adjust thread's location in the queues */
    QRemove(DL_Queue, TokenPtr);
    QRemove(IM_Queue, TokenPtr);
    If (no time constraints) {
        QInsert(IM_Queue, TokenPtr);
    }
    else {
        QInsert(DL_Queue, TokenPtr);
    }

    /* check if we need to dispatch a new thread */
    If (different thread at head of queue) {
        /* dispatch new thread at the head of the queue */
        DispatchThread();
    }

    /* send a message to indicate abort */
    SendMsg(ABORT, AbortThread);

    /* deallocate the token */
    Fdealloc(TokenPtr);
}

```

4.5 Future Directions

The Scheduling Facility described in this chapter represents a first attempt at exploring the area of processor scheduling in Alpha. While this work is the first executable manifestation of years of Archons project research in this general area, much of the character of this experiment derives from the underlying testbed's architecture. In particular, the lack of complete, coherent, and symmetrical sharing of primary memory between the Application Processor and the Scheduling Co-Processor resulted in a number of design decisions that are sub-optimal with respect to overall system structure and performance.

Work is currently underway (as a part of Release 2) to rectify the problems introduced by the Release 1 testbed architecture. However, there is still much work that can be

done with the existing system, and an effort has begun to perform a series of experiments with Release 1 of Alpha in order to determine the comparative merits and behaviors of a collection of different scheduling policies. This work involves the creation of a representative real-time application that can impose a realistic load on the system, the implementation of a set of scheduling policies for Alpha, and the generation of a suite of data collection and evaluation utilities.

Given the great computational complexity of the general scheduling problem, the quality of real-time scheduling policies depend heavily on the proper choice of heuristics. Efforts are underway in the Alpha project to refine the existing scheduling algorithms, develop (and experiment with) new heuristics, and develop optimizations for the scheduling computation's inner loops.

In the future, much of the processor scheduling work in Alpha will focus on extending the results obtained so far in order to better deal with issues of *physical distribution*, *dependencies among schedulable entities*, and *multiple schedulable resources*.

An effort must be made to improve the manner in which the notion time is dealt with across the distributed system. This requires support from the Communication Subsystem, so it is being performed in conjunction with the XTP-based real-time network being used in Release 2 of Alpha. In addition, enhancements must be made in order to deal with the effects that distribution has on the information provided by the application to the scheduler. For example, units of computation with time constraints that are presented to a local scheduler may not be able to properly represent their processing needs, and so false overload conditions may be induced. This is because of the fact that it is possible that not all of a computation will be performed on the local node—a thread may obtain services by executing on other nodes. For this reason, it may be necessary to augment the manner in which the Scheduling Subsystem interacts with the Communication Subsystem (e.g., change it to adjust its loading estimates based on a less pessimistic assumption—as opposed to assuming that all computations will be done locally).

Up to this point in time, all Alpha scheduling work has assumed the complete mutual independence of threads. This is a clearly unrealistic assumption, and therefore work (in the form of a Ph.D. thesis) is underway to extend the current work to cope with scheduling with dependency conditions between threads. No existing scheduling algorithm solves the problem of scheduling a number of activities with dynamic dependency relationships in a way that is suitable for real-time systems. This next unit of work is intended to produce an effective scheduling algorithm, a formal model to facilitate the analytic proof of properties of that algorithm, and simulation results that demonstrate the utility of the algorithm for real-time applications.

Another simplifying assumption made in Release 1 of Alpha is that there is a single Application Processor per node. This made it possible to use scheduling policies based on dynamic deadline scheduling techniques, which do not exhibit optimal behavior when there are multiple processors to schedule. Release 2 of Alpha is being designed to extend the scheduler to work with symmetrical multiprocessors. This involves the modification of the scheduler to manage multiple Application Processors, and to allow the Scheduling Subsystem to execute on any processing element within a node (as opposed to being limited to a specially dedicated scheduling processor).

5 Communication Subsystem

The Alpha Communication Subsystem—that part of the Communication Facility responsible for handling distribution—manages the movement of threads among nodes and guarantees the integrity of the threads as they extend across arbitrary nodes. This latter task includes the detection and elimination of orphaned computations.

Alpha's Communication Subsystem has been designed to provide these required functions by means of a group of communication protocols and to facilitate the development and addition of new protocols. It offers a set of fundamental communication mechanisms upon which a variety of communication protocols may be, and have been, built.

This chapter presents Alpha's communication and distribution requirements, explores Alpha's logical and physical communication architectures, and describes the protocols that have been developed to date to meet Alpha's requirements. The most interesting protocols are discussed in some detail. In addition, this work is placed in perspective with other relevant work.

5.1 Distribution and the Alpha Programming Model

The Alpha programming model has been described in detail in [Northcutt 88b] and an understanding of the basic model is assumed here. However, a few points must be made to explain the effects of distribution within the Alpha programming model. These effects impact Alpha's communication requirements, which are described in the following section.

5.1.1 The Alpha Programming Model

Figure 20 shows a small set of objects and a few threads. Notice that the threads pass through objects and that no physical node boundaries are shown. This is because threads move through objects without regard for node boundaries.

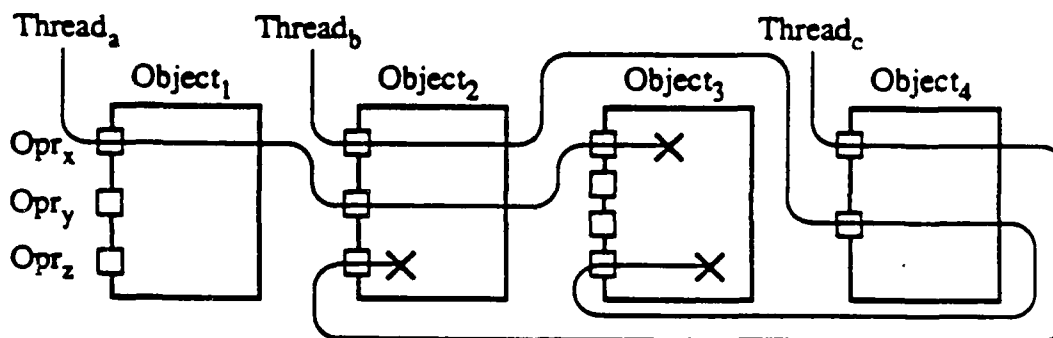


Figure 20: Example Objects and Threads

Every thread has a *root*, and that root is “planted” on the node where the thread was created. Each thread initially executes an invocation on a designated object and may subsequently invoke operations on other objects. Although a thread may span several objects at once (as illustrated by all of the threads in Figure 21), it is only executing instructions in one object at any given time. This active point of control is known as the *head* of the thread, and is indicated by the X on each thread in the figure.

A thread can be decomposed into a number of thread *sections*, where each section is the longest, continuous portion of the thread that is contained entirely on a single node. The entire thread can be produced by concatenating all of its sections. For example, in Figure 21, Thread_a has three sections, one on each node.

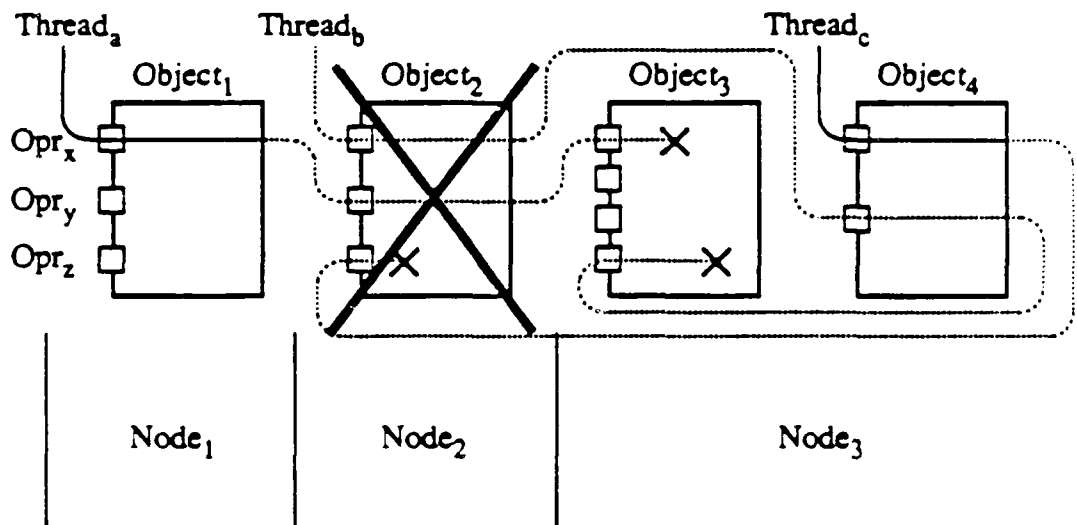


Figure 21: Thread Breaks Due to Node Failure

A thread is *broken* if one or more of its sections are missing due to node or communication failures. Thread breaks often result in orphan computations—computations that can no longer contribute to the overall progress of the thread, and may in fact disrupt its continued progress. To eliminate any potential orphan computations, the thread must be *trimmed* so that the resultant (trimmed) thread has no missing sections between the root and the trimmed head of the thread. The process of trimming the thread is also referred to as *repairing* the thread.

Figure 21 provides examples of thread breaks. The situation shown is identical to that shown in Figure 20. However, in this case, physical node boundaries have been shown; and one node, Node₂, has failed. All three threads have been affected by the node failure and must be repaired. Thread_a, for instance, must be trimmed back to its root section on Node₁, and the orphaned thread section in Object₃ on Node₃ must be eliminated in the process.

Since Alpha is a real-time system, it must execute computations according to application-imposed *time constraints*. A time constraint dictates that a specific group of instructions, possibly spanning several objects, must be executed within a certain interval of time if the program is to behave correctly. Time constraints apply to exception processing—that is, thread trimming due to aborts and failures—as well as to normal processing. Consequently, under normal circumstances, the scheduling parameters for a thread are carried with the thread as it moves through the system making invocations; and under exceptional conditions, scheduling parameters associated with the thread are used to schedule its repair.

At all times a thread is executing under a specific, application-imposed time constraint. In fact, these time constraints may be nested. Whenever a new time constraint is imposed on a thread, the thread executes under the most stringent constraint encountered thus far—it may be the constraint just encountered or it may be the constraint that was already being used for scheduling purposes.

Since Alpha's target applications, and hence their processing requirements, tend to be dynamic and since there are only a finite number of computing resources available, there is a chance that not all of the currently defined time constraints can be satisfied. As a result, some threads may fail to satisfy their time constraints. In that case, these threads are prevented from continuing their normal execution since continuing would mean executing computations at inappropriate times (according to the information provided by the application). Instead, all of the computations that the application desired to complete under the missed time constraint are aborted. Given the nested nature of time constraints, the affected computations are all contiguous, and are all located at or near the thread's head. As pointed out above, the abort processing[†] for these computations is carried out using the most stringent time constraint for the thread that is still satisfiable. Once the designated computations for the thread have aborted, the thread can continue working to satisfy its next time constraint.

Notice the similarity between the abortion of a computation due to an unsatisfied time constraint and the repair of a thread that is broken as a result of a node or communication link failure. In both cases, a continuous portion of the thread, including the thread's head, must be aborted before the repaired thread may continue normal processing. In fact, the coordination that must be done to carry out the processing is identical in the two cases. The only difference is the reason reported to the application for the abort—a missed time constraint in one case, and a broken thread in the other.

5.1.2 Distribution in Alpha

From a programmer's point of view, Alpha has taken a straightforward approach to handling distribution. By employing the object as the fundamental programming unit and a basic system abstraction, and furthermore declaring that objects are indivisible with respect to nodes—that is, an object cannot be divided between nodes—Alpha allows the programmer to focus on the task of writing consistent and correct objects. Operation invocation on other objects is simple—no attention must be paid (by the programmer) to the target object's physical location. Since the entire object is contained on a single node, the author does not have to worry about failure modes where part of the object's state is unavailable due to failures. Each object is written with a well-defined interface consisting primarily of the operations specified for that object's type. Following a node failure or in response to load imbalances, Alpha may move objects among nodes. And since objects do not share memory, they can be moved independently. Node and communication link failures are manifest as the disappearance of object instances and broken

[†]Abort processing in a real-time system is not quite the same as abort processing in other systems, such as transaction processing systems. A real-time computer system typically affects a physical process as part of its normal operation. As a result, abort processing must include not only the restoration of computer memory state, for example, but potentially must also include compensation for actions taken that affected the physical process.

threads. Broken threads are repaired automatically, leaving the programmer to handle the problem of replacing or relocating object instances. This is typically accomplished by creating new instances to replace old ones, or by recovering or migrating old instances.

To simplify the programmer's task and to reduce the bookkeeping required of Alpha, invocations are performed in a location-independent manner—that is, the *invoker* does not need to know the node on which the *invokee* (the target object specified in the invocation) is located. As a result, Alpha must be able to find the *invokee* for any given invocation dynamically.

As was pointed out above, distribution complicates the maintenance of thread integrity (making certain that a thread is entirely intact at all times) and introduces the notion of thread repair. It also demands that Alpha propagate time constraint information along with threads as they move among nodes. When a thread must be repaired due to either a thread break or an abort due to an unsatisfied time constraint (or an aborted transaction), the time constraints for the thread must be circulated in order to schedule the repair correctly. In addition, the affected portion of the thread may span several nodes, so that the abortion of each affected section must be initiated and tracked across the Ethernet.

Not surprisingly, decentralized thread repair is more complicated than centralized thread repair, where centralized thread repair is the thread repair function provided on a single-node system. On a uniprocessor, if the processor fails, all computation stops. On a multiprocessor, if a processor fails, the remaining processors can note that fact and act accordingly, using shared memory to instantaneously and consistently communicate information with one another. In a multi-node system, the loss of a processor or node can be detected by other nodes independently. There is no shared memory to provide a consistent system-wide view for all of the nodes. Rather, each node has its own view of the overall state of the system. Following a node failure, all of the surviving nodes that are affected must detect the failure and coordinate their recovery actions. During the course of handling one failure, another can occur, potentially complicating matters. Although clustered failures may seem unlikely, they may in fact occur often—even during normal operation. For instance, while trimming a thread due to an unsatisfied time constraint, another section of the same thread could be lost due to a node failure, or a second time constraint for the same thread could be missed, necessitating more abort processing (thread trimming).

In addition to timeliness, Alpha is concerned with reliability. In a reliable system, critical data is highly available—that is, despite node and communication link failures, it can be accessed in short order. One way to increase the availability of data in a distributed system is to replicate it. Towards that end, Alpha provides a set of mechanisms that support various replication policies. The major work on the development of replication policies for real-time systems has not yet begun, but a simple initial policy has been put in place in order to begin defining and exercising the necessary replication support mechanisms.

In Alpha, any object can be replicated, and there is no special code inserted in the object to provide replication support. Instead, all of the support is provided by the operation invocation facility and the object naming facility. There are a number of copies of a replicated object, and one of them is designated as the master copy. The others are back-up copies. All invocations on the replicated object are actually performed on the master

copy, and the changes made by each operation are propagated to the backup copies. Should the master copy ever be lost, a simple procedure is employed to promote one of the backup copies to be the new master. Optionally, additional backup copies can also be created. Furthermore, threads are placed in the backup copies by the *ObjectManager* on the master copy's node in order to detect their failure. Should a backup copy ever disappear, the thread that was placed in it will break, and the thread will be trimmed—allowing the thread to execute code in the *ObjectManager* again. This code replaces the backup copy, if possible. Consequently, the *ObjectManager* makes sure that there are always backup copies available in case the master copy should fail.

The Alpha programming model also includes the use of atomic transactions to aid in the construction of reliable applications. The abortion—and, in fact, the completion—of a transaction will require similar coordination to that involved in trimming a thread due to an unsatisfied time constraint. In the case of a transaction, all of the objects participating in the transaction will be involved—even those that no longer host a portion of the thread executing the transaction. As with replication, much of the advanced work related to the development of atomic transactions for real-time systems remains to be done. However, the general demands that transaction support will place on the Communication Facility can already be anticipated.

5.1.3 Related Work

Like Argus ([Liskov 84]), Clouds ([McKendry 84]), Eden ([Almes 85]), CRONUS ([Schantz 85]) and other similar operating systems, Alpha explicitly addresses issues of physical distribution. In addition, all of these systems use the same general abstraction—that of an object—as the basic building block of an application. All but CRONUS feature the use of transactions in order to support the construction of robust applications, as does Alpha. Argus has explored the use of object replication ([Herlihy 84]) as well. Although Alpha has provided mechanisms for transactions and object replication, most of the development of the policies governing their use in a distributed real-time system has yet to be done. CRONUS has emphasized the program development environment and object management issues, areas that have received much less attention in Alpha.

Alpha is distinguished from these systems on two counts. First of all, Alpha explicitly addresses real-time applications. None of the other systems is directly attacking this issue. Timeliness considerations pervade Alpha, and as a result, even functions that are apparently the same in all of these systems—for instance, orphan detection and elimination—differ to a certain extent.

Secondly, all of the other systems—except Clouds—have points of control that are held captive within a single object. In Alpha, these points of control are free to move among objects. This movement among objects can be accomplished in the other systems by means of remote procedure calls ([Nelson 81]) or message passing. The key difference is that Alpha's threads move through the system by themselves—always carrying their own thread-specific scheduling parameters. This scheduling information is not propagated in the other systems, due to the fact that they do not focus on real-time issues. By always propagating the thread's scheduling parameters with the thread, Alpha allows the application writer to create threads that correspond to sequences of activity in the physical world and be assured that the time constraints imposed by the physical world

will be considered at all times as the thread moves throughout the system. In the other systems, intermediate server objects with their captive points of control must be used and will have scheduling parameters that are unrelated to the point of control on whose behalf they are executing.

Clouds comes closer to Alpha in this respect than do the others. It has processes that can move among objects. Its processes are actually equivalent to objects that have a thread embedded in them when they are created. In essence, Clouds has two distinct types of objects—passive and active; while Alpha has only one type, passive, and activity is generated by an entirely different entity—the thread.

Many typical operating systems—Mach ([Accetta 86]) and UNIX ([Ritchie 74]), for example—manage only a single node. Communication with other nodes is carried out by system, or even user, utilities. This contrasts sharply with the approach taken in Alpha where communication, in the form of operation invocation with thread maintenance and repair, lies at the very heart of the system. Each node is designed to behave as a single part of a larger whole, and it is aware that there are other nodes present and that coordinated effort is required to support useful applications.

5.2 Communication Requirements

Alpha has three distinct types of communication requirements. First of all, there are the functional requirements—that is, the set of communication functions that must be performed in order to support the Alpha programming model, as outlined above. Secondly, the experimental nature of Alpha imposes flexibility and adaptability requirements. Finally, the IMI interface and the Alpha hardware architecture constrain communication approaches.

The Alpha programming model dictates the functional requirements for Alpha's communications:

1. implement a location-transparent operation invocation function that propagates real-time scheduling information along with the invoking thread;
2. monitor the condition of active threads — detect thread breaks and repair broken threads as they are detected;
3. support the simple replication scheme described above;
4. provide communication between Alpha threads and objects and UNIX processes — this capability is needed since Alpha, as a research operating system, does not have all of the program development and monitoring facilities, among other things, that UNIX does.

Alpha is a research operating system, and hence its long-term requirements are not precisely known. Since it is a distributed system where communication functions are assumed to be the backbone of the system, Alpha's Communication Facility must be flexible in order to support the changing needs of the system. In addition, Alpha is intended to act as a testbed platform for experimentation with various implementations of a single operating system function or the interaction of designated functions. Once again, the Communication Facility must be flexible in order to meet these requirements.

5.3 Alpha Communication Architecture

The need for a flexible Communication Facility that could be used in a testbed dominated the determination of the Alpha communication architecture.

The functions that must be provided by the Communication Subsystem are quite diverse, and so an approach was taken that employed a number of special-purpose communication protocols, rather than a single general purpose protocol that could support all of the required functions.

Although the protocols were quite different—some were point-to-point, others were broadcast or multicast[†], and yet others used a mixture of these addressing modes; some dealt with only a single invocation, others dealt with many invocations; some were executed periodically, others were not; some used fixed length messages, others used variable length messages—they still demanded some common lower-level support to handle message transmission and receipt, message buffer management and manipulation, time management, and other similar concerns.

To provide this common low-level support, while allowing the development of a number of independent protocols, the resultant architecture (see Figure 22) is layered, providing a set of mechanisms that may be used as primitives by any number of communication protocols. The architecture consists of two key components:

- the Communication Virtual Machine;
- communication protocols.

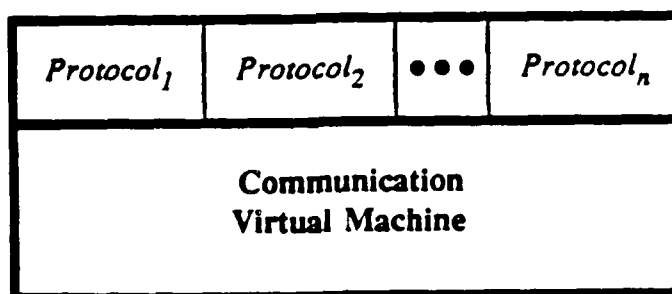


Figure 22: Logical View of the Alpha Communication Architecture

The next two sections describe the Communication Virtual Machine and offer a general discussion of communication protocols within the Alpha communication architecture. Later, the specific protocols that have been implemented to meet Alpha's communication requirements will also be examined.

[†]In a broadcast, a message is sent to all receivers in the system; in a multicast, a message is potentially sent to more than one, but (perhaps) fewer than all, receivers, and those that are not receivers are unaffected by the message. Both broadcast and multicast messages represent instances of one-to-many communication.

5.3.1 The Communication Virtual Machine

The Communication Virtual Machine (CVM) is a simple executive[†] running on the CP, augmented by the communication mechanisms. It allocates the single processor among all of the protocols that desire to use it in a first-come-first-served (FCFS) manner.

As shown in Figure 23, the fundamental data structure within the CVM is the *TokenQueue*. All of the work that must be done is represented in the CVM by tokens. Each token (see Figure 24) has a destination field that indicates the protocol handler for which that token is intended. In addition, each token has fields indicating the source (creator) of the token and its protocol-specific type. As work progresses or arrives, new tokens are added to the *TokenQueue*. As work is dispatched, tokens are removed from the queue.

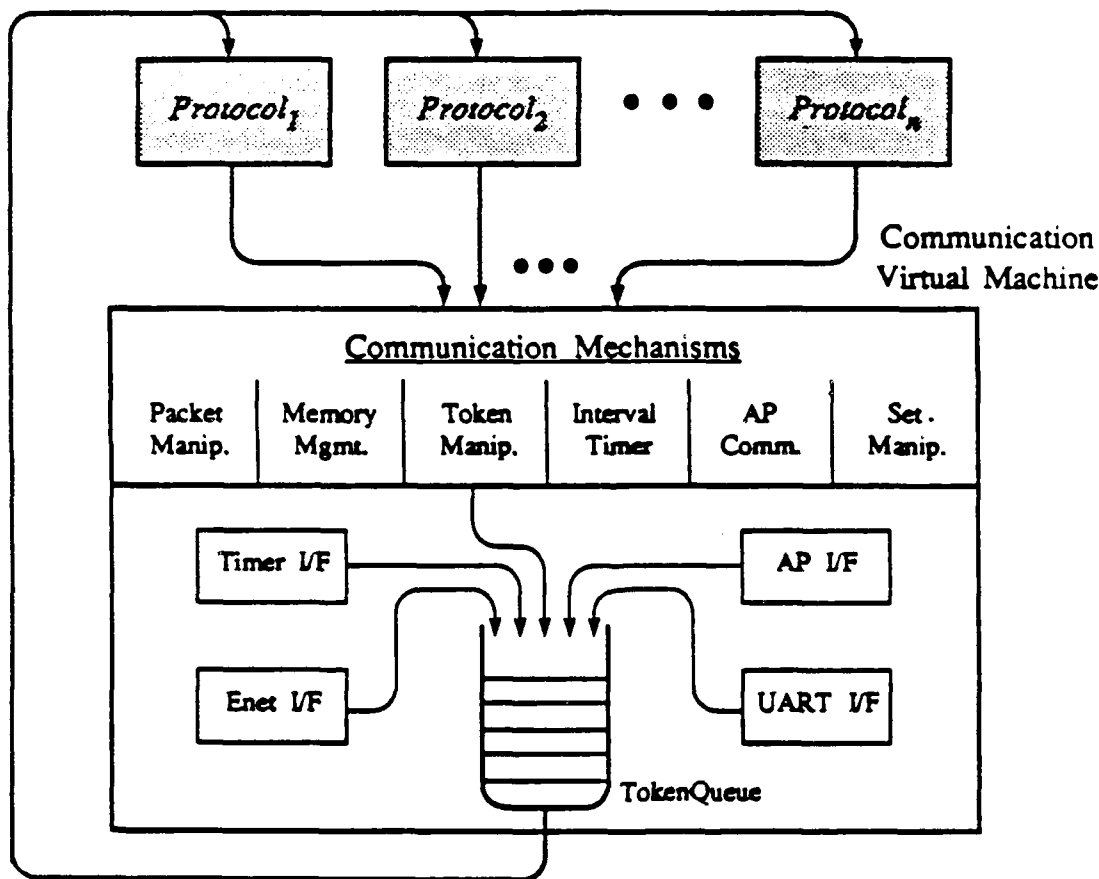


Figure 23: Expanded View of the Communication Architecture

Tokens are added to the *TokenQueue* in response to certain specific *events*. These events may originate from any of a number of physical devices—the Applications Processor, the Ethernet controller, the serial input/output (I/O) controller, or the interval timer—and include:

[†] An executive is a limited operating system. In this case, it is responsible for scheduling the processor and handling interrupts.

- arrival of tokens from the Application Processor Interface;
- arrival or departure of packets from the Ethernet interface;
- arrival or completion of I/O from the serial interface;
- expiration of an interval timer;
- occurrence of specific exception conditions.

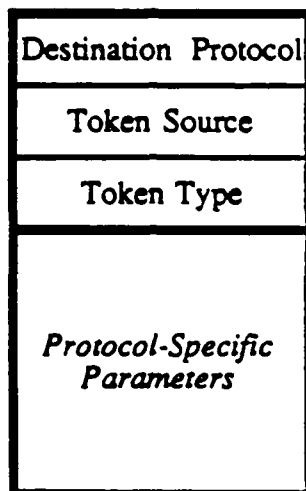


Figure 24: Communication Virtual Machine Token Format

Furthermore, one protocol can generate tokens destined for another protocol, thus allowing the construction of hierarchical protocols. This allows one protocol to use the services provided by others and also facilitates implementations of higher-level protocols that are encapsulated by lower-level protocols.

The CVM employs a simple scheduling algorithm—it simply removes the token from the head of the *TokenQueue* and schedules the protocol handler named in the Destination Protocol field of the token. That token handler is then allowed to execute until it has finished processing the token, at which point it surrenders control of the processor to the CVM executive once again.

Unless they are explicitly disabled, interrupts can occur at any time and will temporarily seize control of the processor. In every case, the interrupt handlers execute for the briefest time possible—they gather whatever state must be recorded in order to perform the bulk of the interrupt handling at a later time and place a token in the *TokenQueue* to initiate that processing.

Protocol writers use a set of communication mechanisms supplied by the CVM as fundamental building blocks for constructing their protocols. These mechanisms include:

- Ethernet packet-handling mechanisms — to create and delete packets; send and receive packets; encapsulate and merge packets;
- interval timer mechanisms — to schedule and cancel timed events;
- memory management mechanisms — to allocate and deallocate shared Multibus memory pages;

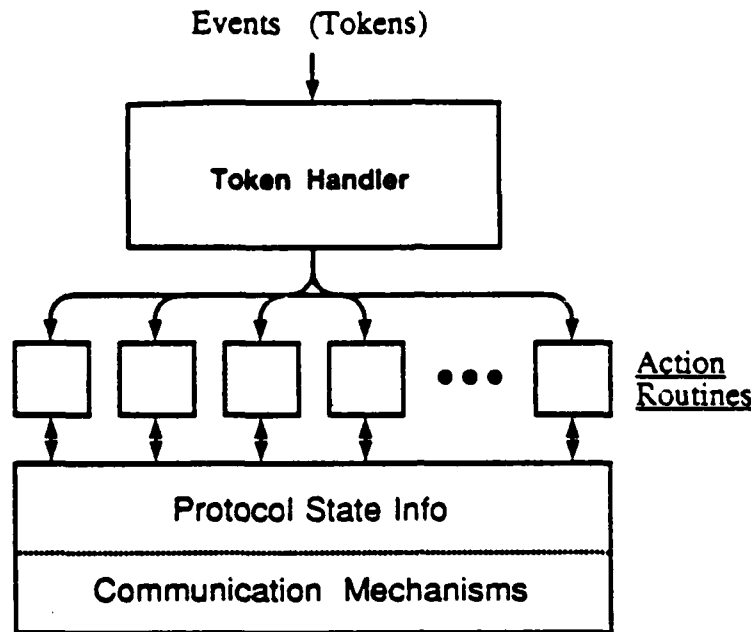


Figure 25: Logical Protocol Structure

- Application Processor interface mechanisms — to send commands to the AP; receive commands from the AP;
- set manipulation mechanisms — to create and delete sets; lookup, add, and remove elements in a set.

5.3.2 Communication Protocols

Protocol definition consists of composing the communication mechanisms to implement the desired protocol. Most of the protocols composed to date consist of a set of relatively short *actions* that are executed at critical times. These actions take the form of routines that are executed in response to the initiating event. To date, these actions have been quite short in terms of the number of steps required to handle the event, which is an encouraging sign. This means that any event for any protocol can be processed fairly quickly, therefore allowing high event-handling throughput to be obtained for the CVM.

Protocols use the CVM to do only what *must* be done at any given time. As a result, they do not spin and wait for another event to occur. Rather, they define an action to be taken in response to that later event and release the CVM for the time being. Of course, many, if not all, protocols will also set a watchdog timer (using the interval timer mechanisms) to guarantee that no work is lost in case the later event never occurs.

Figure 25 shows the logical structure of a protocol. All of the events associated with that protocol are embodied by CVM tokens and are routed to the protocol's *token handler*. Based on the token's source and type, the token handler selects the *action routine* to execute to handle the event. Each action routine has access to the protocol's private state information and to the communication mechanisms provided by the CVM.

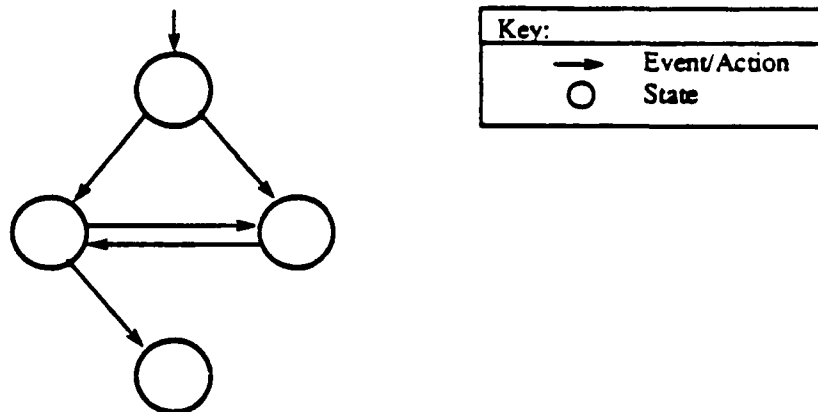


Figure 26: Protocol State Diagram

To specify a protocol, the protocol writer must determine the set of events that may occur during the execution of the protocol, write a token handler that will identify each of them, and specify the steps to be taken by the corresponding action routines for each of them. The flow of control and the interaction of events can often be visualized conveniently in a graphical form, such as the state diagram shown in Figure 26.

The protocols described in this document were all written by hand—that is, they were written in C, without the benefit of a special-purpose, protocol definition language. However, they were written in a very stylized manner in order to explore the feasibility of defining and developing such a language. In the future, a protocol definition language, along with other tools—such as a protocol analyzer—might greatly simplify the design and development of protocols.

5.3.3 Related Work

The Alpha Communication Architecture employs a dedicated processor to provide hardware concurrency in processing and to increase the performance of communication tasks. The use of dedicated hardware for these purposes is by no means new. Similar hardware support dates back at least to the 1960s ([Thornton 64], [Jensen 78]), and continues to the present.

Recently, a number of commercial products have provided hardware communication support with a new twist—the products are communication engines that run specific standard protocols, such as TCP/IP, for specific purposes (like the Micom-Interlan Ethernet NPx00 Network Processor series of protocol processor boards). This is similar to the approach taken in this work, except that, for Alpha, it was deemed important to be flexible enough to support research into a number of special-purpose protocols, not to simply run standard protocols. In fact, some products do support user-supplied protocols ([CMC 85] and the Interphase V-Ethernet 4207 Eagle, among others), and substantial protocol support, much like that desired by Alpha, is being tackled by a commercial venture as well ([Chesson 87]).

There is some debate as to whether it is better to employ general-purpose, standard protocols or to use special-purpose protocols. Typically, the advocates of standard protocols argue that there is little loss in performance introduced by using the standard protocols; there is, almost by definition, wider applicability for applications that use standard protocols; and there is far less time and effort spent on developing redundant functions in a number of different protocols, as might be required for special-purpose protocols ([Sproull 78]). On the other hand, standard protocols do not necessarily provide adequate support for some functions that may be desirable in scalable, distributed systems, such as multicast capabilities. Since Alpha must be free to explore the use of such features, it has chosen to pursue special-purpose protocol development. At the same time, by allowing the construction of hierarchical protocols ([Fleisch 81]), Alpha removes the need for protocols to reimplement functions already available from other protocols.

A considerable number of models for protocol specification have been developed. Some of these take a pragmatic approach to define languages and primitives for communication ([ISO 86], [CCITT 85] and [Nash 83]). Others ([Bochmann 80a], [Gouda 85]) are more interested in defining protocols within a formal model to facilitate verification of correctness and demonstration of properties.

There are many models of protocol specification based on finite state machines (for example, [Bochmann 82] and [Brand 83]), and often protocols are represented graphically in this work. Many of these models differ in details, but the general consensus seems to be that this is a reasonable expressive model for the task at hand. Some finite state models, like the one used for Alpha, do not support non-determinism. That is, they impose a total ordering on events such that, at any time, from any state, only one state transition's precondition is satisfied and, therefore, only one transition is allowed. Other models ([ISO 86]) allow a form of non-determinism where the preconditions of multiple state transitions may be satisfied at once, but still only one transition is chosen (non-deterministically) from among the candidates. In the future, Alpha, inspired by the desire to explore state exchange protocols ([Fletcher 79]) and to develop more general methods to think about protocol specification and behavior, may wish to explore a finite state model that pushes non-determinism one step further—in effect, if several state transition preconditions are satisfied at one instant, then all of them are taken. This is reminiscent of a fork operation and has been used in general finite state machine models before ([Hopcroft 79]). Some of these forked points of control would later cease to exist if no state transition preconditions were satisfied at a later point, while other control points would continue.

The actual development of the Alpha protocols and the CVM mechanisms applied a perspective partially gained through previous work in the area. General design principles of protocols are covered in [Danthine 81]. The notion of using structured building blocks—mechanisms and other protocols—for protocols was influenced by earlier work (such as the Da.akit approach outlined in [Chesson 80]). Finally, the CVM's mechanisms were defined and implemented using optimization techniques, such as those in [Clark 82].

5.4 Alpha Communication Protocols

To meet the functional requirements of the Communication Subsystem, six protocols have been defined. They are:

- **Reliable Message (RM) protocol** — sends packets reliably;
- **Remote Invoke (RI) protocol** — implements remote procedure call (RPC) semantics for invocations whose targets are on another node;
- **Thread Maintenance and Repair (TMAR) protocol** — maintains thread integrity and coordinates repair of broken and aborting threads;
- **Page Transfer (PT) protocol** — transfers memory pages between object instances;
- **Monitor Active Nodes (MAN) protocol** — monitors the “health” of all of the other nodes in the system;
- **Alpha/UNIX Link (AUL) protocol** — provides communication with UNIX systems or applications.

Most of these protocols are quite straightforward and so will be described very briefly. TMAR is the most complex, and the most Alpha-specific, and will be discussed in more depth, as will the RI protocol.

5.4.1 RM Protocol

The RM protocol is quite straightforward. It sends a packet reliably to a single Ethernet destination address[†]. The sender transmits the packet and waits for an explicit acknowledgment from the receiver. The sender sets a watchdog timer in case an acknowledgment is not received within the expected time. If this timer ever goes off, the packet is resent and another watchdog timer is set. This process continues until either an acknowledgment is received or a programmer-specified maximum number of retries has occurred.

The sender portion of the RM protocol always returns a result to its requester—a successful result indicates that an acknowledgment was received in a reasonable amount of time following a transmission of the packet; a failure result indicates that this did not happen. Notice that the protocol could indicate a failure even though the receiver actually got the packet—e.g., due to a heavy load on the remote node, the acknowledgment could arrive late; or a network failure could prevent the successful transmission of an acknowledgment.

The receiver portion of the RM protocol is able to recognize the duplicate packets that may occasionally be sent by keeping track of the sequence numbers for all of the recently received packets from each node in the system. These duplicates are discarded.

The design and implementation of the RM protocol follows the above description nicely. It has four events and four corresponding action routines. The events are:

- the arrival of a request to send a reliable message;
- the arrival of a reliable message;
- the arrival of an acknowledgment for a reliable message;
- the expiration of a watchdog timer.

Although the RM protocol provides a useful communication capability by itself, it is primarily used as a building block for other protocols.

5.4.2 RI Protocol

The RI protocol implements the remote invocation function for Alpha. Operation invocation has essentially *remote procedure call* semantics. In addition, since threads may break at any time, potentially interacting with outstanding invocations, there is a fair amount of interaction between this protocol and the TMAR protocol to be discussed below.

Since each node in the testbed system is a multiprocessor, there are a number of participants involved in carrying out a remote invocation. In particular, the Application Processor, the Communication Co-Processor, and the Scheduling Co-Processor all play a role in remote invocation.

Prior to an invocation, a thread executes code in some object on the invoking node's AP. At some point, it encounters an operation invocation to be performed on an object that is not located on its current node. The Kernel Proper determines that the invocation is directed to an object instance that is not local. The kernel then assumes that it must be intended for a remote object and prepares to perform a remote invocation. To do this, it acquires the scheduling parameters for the invoking thread from the SP, places these parameters with the invocation's request parameters, manipulates memory mapping tables to allow this information to be accessed from the Multibus backplane, and notifies the Communication Subsystem on the CP to perform the remote invocation. At that point, the AP blocks the invoking thread section until it has received the reply for the invocation from the CP. Later, after receiving the reply from the CP, the kernel takes the updated scheduling parameters, which are included with the invocation's reply parameters, and sends them to the SP so that scheduling decisions are made using the most up-to-date information. The invoking thread section is then unblocked and is executed whenever the SP dispatches it.

On the AP of the invoked node, the kernel is notified by the CP that an invocation has been received from another node and is given access to the invocation parameters and the scheduling parameters for the invoking thread, both of which are stored in shared Multibus memory. The AP then passes the scheduling parameters to the SP and schedules a new thread section to be executed using the thread's scheduling parameters. At a time determined by the load of the node and the scheduling parameters under which it is executing, the new thread section will be dispatched. It will execute, possibly performing other invocations, until it completes the invoked operation, at which time it initiates the reply sequence on the AP. This causes the AP to acquire the updated scheduling parameters for the thread from the SP, include those scheduling parameters with the invocation's reply parameters, deallocate the new thread section and notify the CP that the invoked operation has completed.

The CPs on the invoking and invoked nodes coordinate the AP actions just described by means of the RI protocol. The protocol itself has four distinct phases, two on the invoking side and two on the receiving side, that interact with the AP actions just described at the obvious times:

- request leaves the invoking node;
- request arrives at the invokee's node;
- reply leaves the invokee's node;
- reply arrives at the invoker's node.

For each of these steps the RM protocol is used to provide a reliable communication path between the nodes involved.

When an operation invocation is initiated on the invoking node, the invoking CP encapsulates the invocation request parameters and scheduling parameters in a message, addresses the message to the invoked object, and uses the RM protocol to reliably send the message to RI protocol on the node containing the invoked object. A record of the ongoing invocation is created for bookkeeping purposes. After the attempt to reliably send the invocation request, the CP waits to receive a thread synchronization message from the invoked node. This marks the point at which the head of the thread has "officially" moved from the invoking node to the invoked node.

On the invoked node, the RI protocol receives the message requesting the invocation and copies the invocation parameters and scheduling parameters into shared Multibus memory. The AP is notified that the invocation has been received and given the address in the Multibus address space of the parameters. The RI protocol coordinates with the TMAR protocol to register the new section of the invoking thread and then sends the thread synchronization message to the invoking node.[†]

The analogous steps are taken during the reply sequence on the two nodes. Once again, the major complication is making an orderly transition in passing the point of control—that is, the head of the thread—from one node to the other. And in every case, there is always the problem that the thread may have broken and care must be taken to maintain consistent records and to reclaim resources (such as Multibus memory pages) that belong to trimmed threads. This is a theme that occurs again and again in the design of the more complicated protocols: points in the protocols must be identified where asynchronous conditions are resolved. For example, the CVM framework provides an environment in which every action routine is executed atomically, and since events are manifest in the CVM by tokens that are queued in the *TokenQueue*, the queue acts to serialize events. The protocol designer must take advantage of these facts in creating intricate protocols in order to reduce the complexity of the task. The task is already difficult in that much of the work results from the need to coordinate two or more physically dispersed CPs that are truly executing concurrently. Complicated interactions can result.

[†]There is a separate thread synchronization message so that there will always be at least one node responsible for reporting the presence of the thread's head section. In fact, there is a short interval of time during each invocation when both the invoker's node and the invokee's node will report the presence of the thread's head. If, instead, no separate thread synchronization message were sent, then the responsibility for reporting the presence of the thread's head would be implicitly passed from the invoker's node to the invokee's node with the initial invocation packet. As a result, during the time that the packet was in transit and being interpreted by the invokee's node, no node would report the presence of the thread's head. And failure to report the thread's head, under certain circumstances, can result in errantly trimming a section from an unbroken thread. The discussion of the TMAR protocol should make all of this clearer.

To permit a form of logical addressing to be used, the Communication Subsystem maintains a list of the logical names it must recognize. This list is examined each time a logically addressed packet arrives to determine whether the node should accept the packet. Currently, this list is implemented as a hash table whose contents are derived from the Kernel Proper's Dictionary and represent the objects and threads that exist on the node, as well as, other logical entities local to the node (e.g., transactions). When the kernel adds entries to, or removes entries from, the Dictionary, the Communication Subsystem is notified (via an IMI command), in order to keep its logical name table up-to-date.

5.4.3 TMAR Protocol

The TMAR protocol monitors the integrity of threads and repairs threads when a break in the thread is detected. It is also responsible for coordinating the trimming of a thread due to an abort—for instance, due to an unsatisfied time constraint. It is an excellent example of a decentralized protocol executed by a number of identical peers asynchronously on each node in the system.

Every thread has a root, and this root is located on the node containing the object in which the thread was created. Each node is referred to as the *root node* for all of the threads that are rooted on that node. It is the responsibility of the root node to maintain the integrity of all of the threads rooted locally.

To accomplish this, a three phase sequence is executed periodically by each root node. In the first phase, the root node broadcasts a polling message, containing a list of all of the threads rooted at the root node that have made remote invocations, and requesting that any nodes containing sections of these threads report this fact, along with information related to the location of the thread's head, to the root node. In the second phase, each of the nodes containing sections of the listed threads return the requested information to the polling node. This information is examined by the polling node to verify that an unbroken sequence of sections exists between the root and head of each thread in the list. (Some missing section and head reports are tolerated, since communications over the Ethernet are subject to errors. However, a consistent missing section or thread head is certainly a problem.) In the third phase, the polling node broadcasts a refresh message that announces the health of the threads that are rooted locally to all of the other nodes. If a thread is unbroken, the message indicates that the entire thread should be refreshed; if it has broken, only the continuous portion of the thread that begins at the root and ends at the point of the break is refreshed. The sections beyond the break must be eliminated (aborted) before the thread may resume execution at the point of the break. Upon receipt of a refresh message, the other nodes refresh local thread sections as directed.

If a specified time interval is exceeded since the last refresh message for the sections of a thread on a node, then it is assumed that the root node for that thread has failed. In that case, all of the surviving sections are orphans and must be eliminated.

Over time, this three phase poll-respond-refresh sequence keeps all of the unbroken threads intact, without bothering the AP.

In addition, due to the use of broadcast and multicast messages, the amount of information that must be maintained for each thread in the system is small. For instance, the root node never necessarily knows where the head of a given thread is, or in fact, how many sections the thread currently has. It cannot be sure of such information because it

only knows what it received in response to its most recent polling message. Since messages may be lost, it cannot be sure that it received a report about all of the existing sections. Furthermore, even if it did receive all of the responses, the thread may have subsequently executed a new invocation that took it to another node, thereby creating a new thread section.

There is no need for the root node to be aware of such invocations. It is only attempting to ensure thread integrity, and that can be accomplished by detecting thread breaks, which are characterized by missing thread sections, possibly including the head section. Threads are free to move freely without explicitly reporting their movements to the root node as long as they do respond to the root node's polling messages.

Polling messages are sent periodically by each node in the system. Also, the time limit within which a thread must be refreshed is a programmable, integral number of polling cycles, as is the number of responses for any thread that must be missed before it is assumed to be broken. Therefore, the polling frequency establishes the responsiveness of the system—that is, the length of time that passes after a thread break occurs until it is conclusively detected. This allows the protocol to trade communications bandwidth for thread break responsiveness: if necessary, a system can be made more responsive by utilizing adequate bandwidth to use a sufficiently short polling interval to satisfy any responsiveness requirements.

The Alpha programming model dictates that threads must be trimmed according to the proper application-defined scheduling parameters for the section that will be at the head of the trimmed thread. The TMAR protocol handles the circulation of the proper scheduling parameters for trimming a thread.

As a thread moves from one node to another making an invocation, it carries with it only the scheduling parameters under which it is currently executing. No information concerning any other previously established, looser time constraints moves with the thread.[†] Instead, whenever other scheduling parameters are needed to schedule thread section aborts during a thread repair, TMAR locates the proper parameters and circulates them to the nodes that need them.

In the case of a thread break, the proper scheduling parameters are known by the node having the thread section immediately preceding the first break in the thread (that is, the section adjacent to the missing section but nearer to the thread's root). This section will be the head section of the repaired thread. When the root node for the broken thread determines that a break has occurred, it determines the section identifier of the new head section and sends a refresh message that only refreshes the thread up to and including

[†]*This is done to limit the amount of scheduling information that must be propagated with a thread. Since time constraints can be nested arbitrarily deeply, there is no fixed bound that will necessarily hold all of the scheduling information for a thread. Therefore, it cannot be assumed that the scheduling parameters that are needed following an abort or thread break can always be carried with the thread. Consequently, there must be some method to locate the necessary parameters if they are not among those that are propagated with the thread. Since thread breaks will occur at arbitrary locations within a thread, and since there is no required correlation between the tightness of time constraints and the points (with respect to the thread's head) at which they were defined, it makes sense to minimize the amount of scheduling information propagated with the thread, and to use the method that must be provided in any case to locate and circulate the proper scheduling parameters to be used during thread trimming.*

that new head section. All sections beyond that point will be marked for abortion following receipt of the refresh message. If the new head section happens to be located on the thread's root node, then the scheduling parameters under which the abort should be performed are included in the refresh message. Otherwise, when the node holding the new head section receives the refresh message, it notices that it holds that section and sends a message directly to the root node containing the desired scheduling parameters. The root node includes these parameters in all subsequent refresh messages for the thread until the thread repair is complete. Until the proper scheduling parameters are received, the abortion is scheduled on each node using a default time constraint that indicates the abort must be completed in the distant future (literally, by the end of time as far as the system is concerned). In this way, progress can be made on the abort if there is no other work to be done in the meantime; and if there is work to be done, the abort will not interfere with it until the proper scheduling parameters are received.

The situation is somewhat simpler when TMAR is trimming a thread due to an abort, such as an unsatisfied time constraint. In this case, the node on which the unsatisfied time constraint was initially established detects the fact that it has been missed, and immediately notifies the root node of this fact along with the desired scheduling parameters that are to be used to schedule the abort processing. As before, the root node then circulates these parameters to the other nodes containing sections of the thread being trimmed.[†]

During the trimming of a thread, each node trims the sections that are local to it that are beyond the point of the break or abort. After all of these sections have been aborted, the root node for the thread is notified of that fact. The root node continues to perform its normal polling cycle for the thread, refreshing only that portion preceding the break or abort point. It keeps track of the aborted sections as it is notified of them and determines when the entire thread trim operation has been completed—that is, when each section beyond the break or abort point has either been aborted or determined to have failed (presumably because of a node or communication link failure).

Once the entire thread has been trimmed, a notification is sent by the thread's root node to the node containing the trimmed thread's head signalling that the trimmed thread may resume execution.

Since TMAR on a thread's root node continues to perform its normal three phase refresh function even while the thread is being trimmed, the protocol can handle successive failures. That is, while a thread is being trimmed, subsequent thread breaks or aborts will be handled as soon as they are recognized. This makes the protocol somewhat more complicated than it would have to be if subsequent breaks and aborts could be deferred, but it will improve responsiveness and reduce the time that the trimmed thread must wait before proceeding in several cases, most notably when multiple nested time constraints cover a relatively small interval of time.

[†] Actually, each node having a section of the thread operating under the unsatisfied time constraint will independently detect the fact that it has not been met. However, only the node on which the time constraint was established knows the time constraint that was superseded. However, by detecting it themselves, the other nodes have the opportunity to carry out some, or all, of the abort processing (under default scheduling parameters) without having to wait for notification from another node.

Although it is straightforward for the root node to wait for the next polling cycle to take many actions, responsiveness can be enhanced if it acts as soon as possible. Consequently, optimizations have been inserted into the protocol to allow the root node, and in some cases other nodes, to act without waiting for a polling cycle.

The root node's task of determining the state of a given thread is not as simple as it might seem. This is because of the asynchrony among nodes and, hence, TMAR instances. When the root node broadcasts its polling message, the thread is at a single definite place. Each node that is aware of the thread responds to the polling message, but only after processing it, which involves executing a TMAR action routine to determine which threads are of interest. Since all of the nodes have *TokenQueues* of different length and operate asynchronously, their TMAR action routines will execute at different times and, as a result, their views of where the thread's sections are may differ. For example, if the thread is performing an invocation, possibly two different nodes will report having the head of the thread. Added to the possibility that either the polling message or a node's response may be lost due to the nature of the Ethernet, this makes it difficult to be sure about the condition of the thread based on polling information.

5.4.4 PT Protocol

The PT protocol was designed as a general utility, but there was a specific need for it as well—supporting a primitive version of object replication. In the future, it may be employed by the Secondary Storage Subsystem, a debugger or monitoring system, or the invocation facility for large parameters.

The PT protocol transfers arbitrary pages from one object to one or more other objects. The RM protocol is used to send these pages reliably to each destination. Unfortunately, due to Ethernet packet size restrictions, two packets are required to hold a page of data. This means that each transfer requires two individual steps, which complicates coordination somewhat when compared to a single step. (Remember, failures can happen at any time on either end of the transfer.) Of course, the additional coordination and the number of packets required act to decrease the performance of the protocol.

5.4.5 MAN Protocol

The MAN protocol, like the PT protocol, was included in the initial suite of protocols in order to support Alpha's first replication scheme. This protocol allows each node in the system to keep a fairly accurate picture of the status of all of the other nodes in the system at any given time. Once again, the protocol is decentralized, and so each node will have its own view of the other nodes' status.

Each node periodically broadcasts a packet announcing its presence to the other nodes in the system. When such a packet arrives from another node, it indicates that the node should be assumed to be active. If it had previously been assumed to be inactive, then this constitutes a change in state, and the kernel is notified of the change.

If any node has not been heard from for a sufficiently long time, then it is assumed to be inactive. Once again, if this constitutes a change in its assumed state, the kernel is notified of the change.

Finally, the kernel may request a complete report of the current node status information, and the CP will oblige.

5.4.6 AUL Protocol

The AUL protocol provides a link between Alpha and a UNIX environment. Specifically, it allows reliable messages to be sent between the two environments. The reliability is provided by means of positive acknowledgment of messages and the detection and suppression of duplicate messages. Given the UNIX development environment under which Alpha was developed and initially tested, a facility to link them together seems quite advantageous.

In the UNIX environment, AUL presents a standard UNIX interprocess communication (IPC) interface to the Alpha environment. That is, stream sockets, based on DARPA[†] standard TCP/IP ([Postel 81a], [Postel 81b]), provide the connection to Alpha, allowing flexibility and ease of use for the UNIX programmer. Communications from Alpha are also presented to UNIX processes through stream sockets.

In the Alpha system, an object interface is provided for communication to UNIX. A message is sent to UNIX as a result of an invocation of the *Send* operation on the *NetIO* object. Also, there is a single thread that is dedicated to distributing incoming messages from the UNIX environment to the designated recipients.

The two environments access a common Ethernet, although any interconnect that supports the TCP/IP protocols could be used. User Datagram Protocol (UDP, [Postel 80]) packets are actually sent over the Ethernet between the two environments. Under UNIX a network server that sends and receives UDP packets implements the AUL protocol. User processes then communicate with this network server by means of standard UNIX IPC.

5.5 Evaluation of Alpha Communication Architecture

The six protocols just described were necessary and sufficient to meet the requirements of Alpha and its first applications. Therefore, they provide a substantial basis on which to evaluate the value and efficacy of the communication architecture chosen for Alpha.

5.5.1 Protocol Definition Model

First, the model itself—the CVM, the mechanisms, and protocols expressed as state machines—seems appropriate for present purposes. All of the protocols could be broken down into appropriately small steps (action routines) that respond to specific events. The mechanisms have been adequate, but could stand some improvement. They will be discussed in more detail below.

A more general finite state machine model was originally considered for use in Alpha, but since the simpler model seemed capable of providing the required functions more efficiently, the simpler model was chosen. Specifically, the notion of allowing *epsilon-transitions* in a state transition graph was considered and rejected (for the present). These transitions permit specified state transitions to occur “spontaneously”—that is, the transitions are not triggered by explicit events from the *TokenQueue*.

[†]Defense Advanced Research Projects Agency.

Epsilon-transitions would allow a significantly different approach to be used in structuring protocols. For instance, under Alpha's current finite state machine model, if a protocol provides a conversation facility for two communicating peers, then the finite state machine's current state uniquely describes *the* current state of the conversation. If non-determinism—in the form of epsilon-transitions—were permitted, then a given conversation could be in any of a number of states at any given time. Non-determinism would allow each state to assume that a specific set of conditions were satisfied when the state was entered. Actions could then be carried out based on those assumptions, and, if later events invalidate any of the assumed conditions, the protocol state diagram can be defined so that there are no valid state transitions from that state for those events.

Successfully applying this approach might reduce the amount of bookkeeping and overhead involved in writing a protocol. However, this approach would also be costly:

1. more state transitions would have to be examined and action routines executed (since more than a single transition would be possible at any time, compared to exactly one transition in the current model);
2. protocol behavior would be much more difficult to analyze, once again, due to the potentially increased number of states that must be considered at any given time for a single protocol conversation; and
3. since each action routine can take actions that are visible (e.g., sending packets or sending commands to the AP), having more than one action routine execute for any given conversation could be very confusing.

While it may be interesting to explore the use of non-determinism in protocol design and implementation, it does not seem to be urgent, or even necessary, at this time in order to meet Alpha's (non-trivial) requirements.

5.5.2 Mechanisms

The functions provided by the mechanisms were generally sufficient, but some improvements were needed. In particular, although the mechanisms allowed packets to be encapsulated by means of adding headers and trailers[†] to existing packets, there was inadequate provision for stripping this information from received packets. It is clear that there is a hardware advantage that allows the encapsulation to be performed efficiently, and no such advantage exists for stripping the encapsulation. However, in the interest of providing a uniform and complete set of mechanisms that provide well-defined, complementary functions and obey typical rules of abstraction—such as, hiding irrelevant details—the mechanisms should conceal this fact.

Aside from function, the speed of the mechanisms must also be considered. To date, the speed of the mechanisms has been adequate, but it is somewhat slow. The overhead imposed by the CVM model to compose individual mechanisms is minimal, and the mechanisms that manipulate the token queue are efficient. However, the hardware architecture imposes considerable inefficiency in some of the packet handling mechanisms. In particular, the routines that create and receive packets must do excessive copying because the memory used by the Ethernet controller is not directly accessible by the AP

[†]A header is a group of fields preceding the encapsulated data, while a trailer is a group of fields following the encapsulated data.

or CP. In addition, the mechanisms that manipulate the interval timer are surprisingly slow. They provide the semantics that are desired, but the underlying hardware does not match these semantics.

5.5.3 Flexibility

The protocol definition/CVM model provided excellent flexibility. Simple protocols could be defined and integrated into the CPs suite in a matter of hours or days. And simple protocols were easy to write, which is always a useful property because, to date, simple protocols have outnumbered complex protocols. A protocol definition language might reduce the development time for protocols even further.

5.5.4 Language

All of the protocols described were hand-coded. Nonetheless, they were coded in a stylized manner in order to evaluate the feasibility of developing a protocol definition language. The stylized approach did seem useful; it fit the job nicely. And several "idioms" did develop during the course of the definition of these protocols. These can be hidden/subsumed by a programming language.

Overall, a language seems useful, but it does not look like it will improve performance. It will however simplify the specification of the protocol and therefore reduce the time required to write the definition and also reduce the opportunity for minor bookkeeping errors.

5.5.5 Hardware

The hardware architecture was responsible for several problems in developing the mechanisms and the protocols. First of all, Ethernet is inappropriate as a global bus for Alpha, or any real-time system for that matter, because it cannot provide a deterministic, or at least a known, delivery time for any message. If there is no global clock—as is the case with Alpha, a distributed real-time system needs to know how long messages spend in transit. Currently, average delivery times are measured and then used as if they are deterministic. This is an approximation and will occasionally be wrong.

Other problems result from the Intel 82586 LAN co-processor used on the SMI Ethernet Controller Board. The controller is capable of selectively receiving multicast packets over the Ethernet. However, due to the method used to load and change valid multicast addresses for a node, multicast packets for a node may be frequently missed during changes to the multicast address table. Unfortunately, for the protocols implemented so far—in particular, the RI and TMAR protocols—an objectionable number of packets were missed due to this problem (around ten percent for one simple test program)[†]. As a result, the hardware multicast filtering capabilities are not used; instead, the filtering is done in software. This, of course, is costly in terms of performance.

Finally, most, if not all, of the protocols are complicated due to the fact that packets can be lost because of the nature of Ethernet (e.g., two simultaneous packet transmissions will result in a collision, which garbles both packets and which may not be detected by

[†]Intel reports another problem with the 82586 LAN co-processor that also results in missed multicast packets. Due to this problem, an estimated five percent of the multicast packets addressed to a node may be missed.

both transmitters). Typically, the protocols are more complex because states are added to cope with the possibility of losing a packet, and, as states are added, the number of situations (states) that must be considered for each event that occurs increases. In addition, the additional states and checks that are required often reduce the performance of the protocol.

5.6 Future Directions

The future work in Alpha communications can be broken down into four categories. First, the CVM and underlying mechanisms can be improved. The CVM model can be extended to include non-determinism, but as was pointed out above, this would open up a new research area, rather than make a short-term contribution to the Communication Facility.

Second, the Alpha hardware architecture should be improved as much as possible. Most of the changes that could be made in this area have been referred to earlier. Specifically, as many of the following changes as possible should be made:

- replace Ethernet global bus with a bus that can provide deterministic or known delivery times for packets;
- add dual-ported memory that may be accessed by both the Ethernet controller and the AP and CP processors — this will eliminate a significant amount of processing;
- provide a hardware interval timer — the interval timer chip should manage a queue of timed events (that is, events that are requested to occur at specific times) and these events should be ordered according to their event times; it should be possible to specify new events that are correctly inserted into the queue or to remove previously defined events from the queue; one or more parameters should optionally be associated with each event and should be returned by the interval timer when the timer for an event expires; given the events in the queue, the chip should always determine the parameters needed to respond at the right time for the event at the head of the ordered timed event queue; currently, the mechanisms that implement the interval timer functions are time-consuming and could be significantly improved with better hardware support; this is particularly important since every protocol to date has made extensive use of the interval timer mechanisms;
- employ an improved Ethernet controller — better multicast support will reduce the burden on the software and improve performance;
- increase memory bandwidth or buffering on the Ethernet controller — this would reduce or eliminate the underrun and overrun[†] problems that have been observed in the existing implementation; the underruns and overruns do not

[†]The LAN co-processor that is used as an Ethernet controller constructs packets based on information in packet descriptors as they are transmitted. If it cannot read data from memory quickly enough to sustain the transmission, then the packet cannot be completed and the transmission is aborted. This is referred to as an underrun. Similarly, when a packet is received, it is placed directly in memory as specified by a set of packet descriptors. If there is insufficient memory bandwidth available to place the data in memory as it arrives, then some data will be lost. This is referred to as an overrun and results in the loss of all of the information in the partially received packet.

“break” a protocol, but they do result in retransmissions that would otherwise be unnecessary; and they occur often enough to warrant attention.

Third, a protocol definition language can be developed. Enough experience has been gathered to permit a language to be defined. This will simplify the task of specifying a protocol and will also reduce the opportunity for error by removing many of the more mechanical, detailed tasks from the protocol designer (this includes the fact that it will insulate the programmer from details concerning the communication mechanisms).

Fourth, much more work remains to be done on protocols. The existing protocols can be refined, and several different trade-offs can be evaluated. For instance, many protocols can consume more bus bandwidth in order to reduce the response time for some action. Assumptions of the reliability of packet delivery can also be explored in order to trade the right amount of error detection for the required performance, remembering that these trade-offs must be made with the priorities that are important for a real-time systems, not a general-purpose computing system.

In addition, a number of new protocols must be developed. Most prominently, the preliminary support for object replication should be improved and support must be developed to support distributed transactions. Part of the transaction support—to handle orphan detection and elimination—will be very similar to the TMAR protocol, and part of it—to handle coordination of committing and aborting transactions—will probably be traditional.

Fifth, steps should be taken to schedule activities in the Communication Co-Processor in the same manner as threads are scheduled in the AP. The architectural changes described earlier should be made to support this end. If the threads on the AP could efficiently access memory belonging to the CP and the Ethernet controller, then they could carry out many of the action routines in response to communication events themselves. Coordination would need to be provided for access to shared data structures and devices, of course, but otherwise, threads could proceed to carry out their own invocations as best suited the application. In general, the object/thread model should be pushed as far as possible into any of the subsystems in order to take greatest advantage of the sophisticated scheduling done by Alpha. (It is impossible to cheaply imitate this scheduling behavior because it is so compute-intensive. In fact, Alpha dedicates a separate processor to perform this function for the AP. It would be quite possible that the CP, or any other subsystem, would spend most of its time scheduling if an attempt was made to replicate the same algorithm on another subsystem.)

Finally, all of the previous changes define areas for the continued research and development of the Alpha Communication Subsystem. However, at some point in the future, Alpha may require more of a production, as opposed to experimental, facility for communication. This will change a number of requirements and priorities in designing the subsystem. Specifically, flexibility may not be important in a production version of Alpha. If it imposes a sufficiently large cost, it could be eliminated. The strict interfaces that currently exist might also be broken down to some extent in order to improve performance. Some work must be done in the future to anticipate how to transition Alpha's research results into practice in a production system.

5.7 Summary

Alpha is a distributed operating system and must therefore support communication functions among a number of separate processing nodes. The Alpha programming model defines these communication functions. Since Alpha is a research operating system that, among other things, desires to explore several, possibly diverse, approaches to providing these functions, the Alpha communication architecture emphasized flexibility and adaptability. Six specific protocols have been designed and implemented to provide the necessary communication functions. While most of these protocols have been straightforward, one (TMAR) provided fairly unusual functions in a decentralized manner.

The implementation of the first set of Alpha protocols has demonstrated the viability of the communication architecture and has made Alpha a truly distributed system with unique features. Based on the experience gained from this effort, a number of future directions have been identified—most dealing with improvements or modifications to the architecture. A foundation has been laid to support exploration into a number of communication protocols to support the Alpha object/thread paradigm, in general, and the use of atomic transactions and object replication, in particular.

6 Storage Management Subsystem

The final subsystem in Alpha is concerned with the management of storage resources, and can be broken down into two major parts based on the type of storage that is being managed—i.e., *primary storage* and *secondary storage*.

The first section of this chapter describes the services that Alpha provides for the management of the Application Processor's primary memory. This includes the management of physical memory pages, the management of virtual memory spaces, and the handling of virtual to physical address translations. This chapter's second section describes the Alpha secondary storage facility.

6.1 Primary Storage Management

This describes the manner in which the primary memory of an Application Processor is managed by the Alpha kernel. The management strategies applied to physical and virtual memory are described here, as well as the management of the processor's translation tables and the manner in which virtual memory faults are handled. But first, an overview of the testbed's physical memory structure is given, as well as a definition of the requirements that Alpha imposes on its memory subsystem.

As indicated earlier, the Sun-2 processing element supports a collection of physically separate virtual address spaces known as *contexts*. The Sun 2.0 MMU can store address translation information for a maximum of 8 contexts at any given time (corresponding to a virtual address space of 16MB). Each context is subdivided into 512 *segments* (corresponding to 32KB's worth of virtual address space), each of which is further subdivided into 16 *pages* (each of which correspond to 2KB of virtual address space). This MMU design allows quick context swaps among currently loaded contexts to be performed by simply changing the contents of the context register. The lookup tables used to implement the MMU's address translation are fast, but since they are limited in size, they must be managed like a cache by system software. This management is one of the more difficult aspects of the kernel's virtual memory management function.

In Alpha, virtual memory is subdivided into a hierarchy of structures starting, at the highest level, with contexts. Each client thread in Alpha is associated with a separate context, and the Alpha kernel itself coexists in the each context along with a client thread. Client objects are not bound to a specific context but rather float among thread contexts. Also, as a result of the MMU caching activity, client threads can be bound to different hardware contexts in the course of the system's execution. While the kernel itself shares each context with client threads, the kernel region is protected (for *supervisor* mode access only) so as to disallow access from the client objects (that execute in *user* mode).

Each context in Alpha is partitioned into a set of *memory regions*, that are themselves composed of one or more contiguous areas of virtual memory known as *extents*. Each extent is, in turn, composed of one or more of the virtual memory segments defined by the hardware. Each extent in the system has associated with it a secondary storage object, that is used to store an image of the extent for paging or swapping. An extent may exist completely in primary memory, or may be (entirely or in part) in its related secondary storage object. Figure 27 illustrates the decomposition of virtual memory spaces in Alpha.

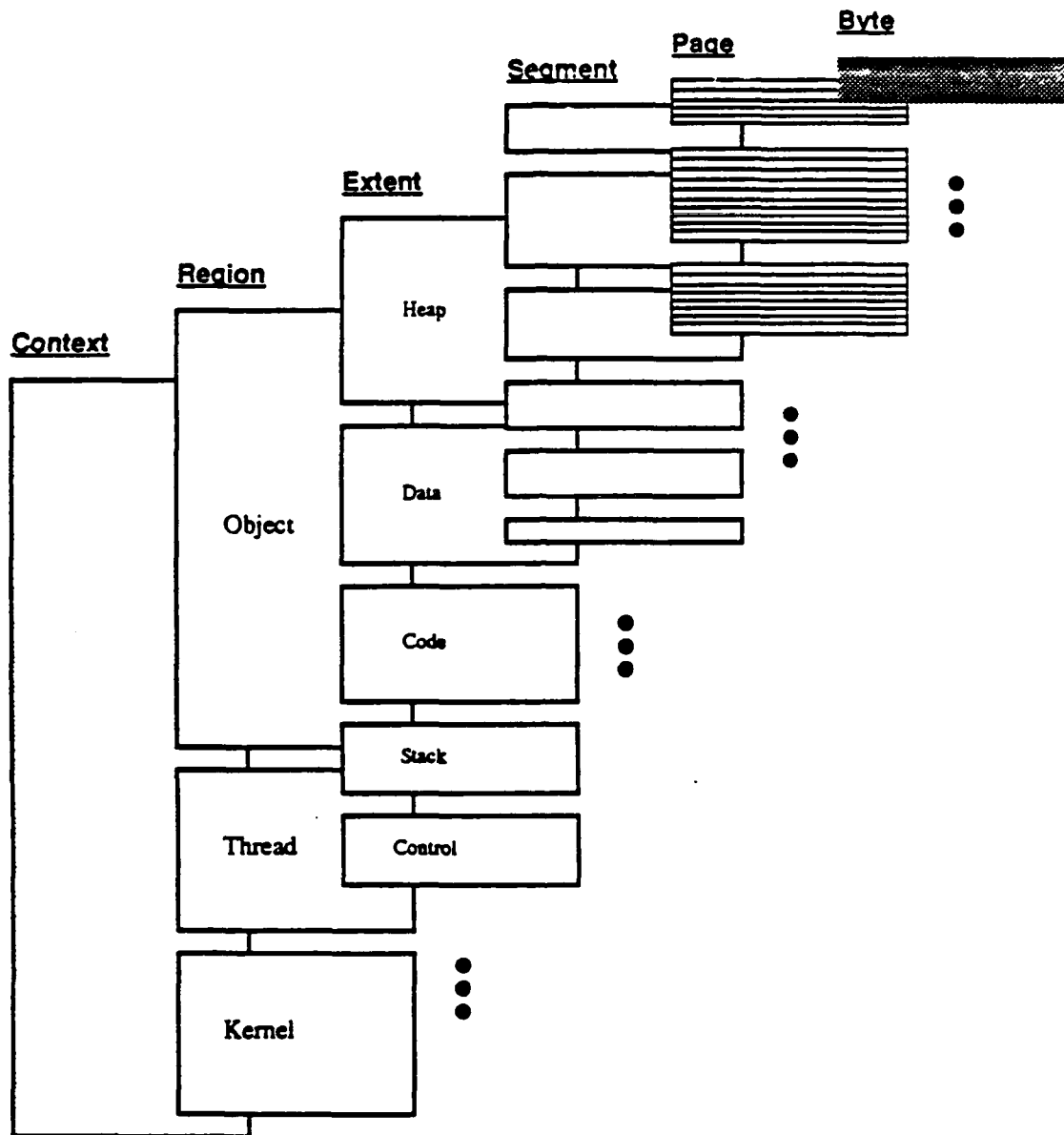


Figure 27: Virtual Memory Hierarchy

The Alpha virtual memory facility is designed to exploit the system's limited physical resources as efficiently as possible. In the case of the Sun Microsystems version 2.0 processing element, the most critically limited resources are physical pages, entries in the SMC, and entries in the PMC. With 1MB of physical memory in the Application Processor, there are only 512 physical pages available. The segment descriptors for only 8 contexts can be contained in the SMC at one time, and only 256 sets of page descriptors can be in the PMC at any one time.

In addition to the effective management of these critical resources, the Alpha virtual memory facility must provide the control information necessary for performing each of the kernel's functions, including the invocation of operations on objects, the paging and swapping of objects and threads, and the dynamic creation or deletion of objects and threads.

The virtual memory structures must be designed to make it possible to share portions of secondary storage objects among extents, and for multiple threads to access the extents that make up an object simultaneously. In order to share them, it must be possible to reclaim critical memory resources preemptively. It must also be possible to share extents, both to maintain the consistency of objects that have multiple threads active in them simultaneously, and to utilize the PMC entries efficiently.

Ideally, there should be only one virtual memory structure on each node for each shared extent, and they should share read-only secondary storage objects wherever possible—e.g., the code extents of instances of the same type of object should make use of a common secondary storage object. Separate secondary storage objects are required for extents that contain instance-specific data, and other extents (e.g., initialized data extents), may use a common secondary storage objects for their initial page-in requests, and their own individual secondary storage objects for paging out all dirty pages and handling any subsequent page-in requests for modified pages.

The MMU is implemented as a two-level lookup table that is managed by Alpha as a cache for the complete collection of address translation information stored in primary memory. All of the threads and objects in Alpha can be bound to and unbound from different contexts as required in the course of their execution. Threads and objects can dynamically come into, and go out of, existence on a given node, and the virtual memory structures must be able to cope with all of this. The virtual memory structures must also cope with the fact that not every node in Alpha will have secondary storage. The virtual memory structures must also permit the replication of an extent's secondary storage objects at several different secondary storage sites for reliability and performance reasons.

6.1.1 Physical Memory

The Sun Microsystems version 2.0 processor board dedicates separate and complete virtual address spaces to local primary memory, memory mapped on-board devices (including the USART's, timers, the monitor PROM's, and the address translation tables that comprise the MMU), Multibus memory, and Multibus I/O. Each page of a context can map into any one of these separate address spaces. The contexts are initialized by the processing element's hardware and firmware, and then the Alpha kernel takes over and completes the initial set-up of the virtual memory facility.

The monitor initializes the processor on power-up, sets the MMU up with an identity mapping of virtual to physical addresses in local memory, invalidating the addresses above the end of the available physical memory, and initializing the necessary memory locations.

The MC68010 microprocessor imposes certain restrictions on the use of the low memory addresses (commonly known as *low-core*). Low-core initialization involves setting up an initial system stack, and installing an initial set of exception vectors (i.e., processor exception vectors, software trap vectors, and interrupt vectors). The monitor also uses portions of low-core to maintain its local variables and a small stack of its own (the *monitor stack*).

Following the power-up initialization performed by the monitor, the kernel is loaded into the system above the defined low-core region (i.e., it is loaded starting at the second segment of all contexts). A page of memory is set aside for the kernel's interrupt stack.

then the code and data regions of the kernel are loaded into memory, followed by all of the “wired” kernel threads and objects. These components are combined into a single load-module that is to be loaded into an application processing element when a node powers up (see Figure 10). The same load module is used for each node in the system; differentiation based on a specific nodes occurs following the initialization of the generic kernel.

The kernel’s initialization code begins by modifying the MMU tables to construct the desired virtual address mappings for the kernel. The kernel’s virtual memory initialization function sets aside the physical pages which contain the low-core information and the load module. These pages are locked in memory (i.e., not available for being paged out) and consist of a contiguous portion of memory at the bottom of the physical memory space that does not participate in virtual memory paging activities. This approach reduces the initial fragmentation of the physical memory space. The remainder of the physical memory is linked together into the *free page list*, from which pages of physical memory are allocated on demand.

The kernel’s pager and swapper daemons (described in Subsection 3.2.1.1) are used in an attempt to keep a small number of pages in the free page list at all times. Individual pages of memory can be temporarily exempted from the paging activity. This is known as “sticking” pages down (as opposed to *wiring* them down, which is permanent), and is accomplished by marking the page as such in its PMAP.

When a physical page is to be written out to secondary storage, it is placed into the page-out list to have its contents written out by the pager daemon, and then be returned to the free list. In a similar fashion, a thread’s invocation parameter pages are linked together and accessed in a stack order. As the thread makes local allocations, it no longer requires its current incoming invocation parameter page. This page is linked into the thread’s list of invocation parameter pages, to be reclaimed when the operation completes and the invocation returns.

6.1.2 Virtual Memory

Each virtual memory context in Alpha has a common layout. Each context is composed of the following regions: the kernel region, the client thread region, the client object region, and the kernel I/O region. Figure 28 illustrates layout of these regions a virtual address space. The following paragraphs provide descriptions of each of these memory regions, their contents and intended use.

6.1.2.1 Kernel Region

A context’s kernel memory region consists of two main parts—the *monitor* and the *Kernel Proper*. The monitor portion contains the processor’s exception vectors, the data needed by the monitor, and a stack used to execute on when executing in extended interrupt handling routines. The monitor’s memory extents are *wired down* (i.e., not available for page replacement), so they are always loaded in the Application Processor’s MMU.

The Kernel Proper portion of the kernel memory region consists of extents for the kernel code, the kernel data, the kernel heap, and the kernel page heap. The pages that make up each of the different kernel extents are protected according to their intended purposes—i.e., supervisor mode execute-only for the kernel code extent and read/write for data, and no access in user mode.

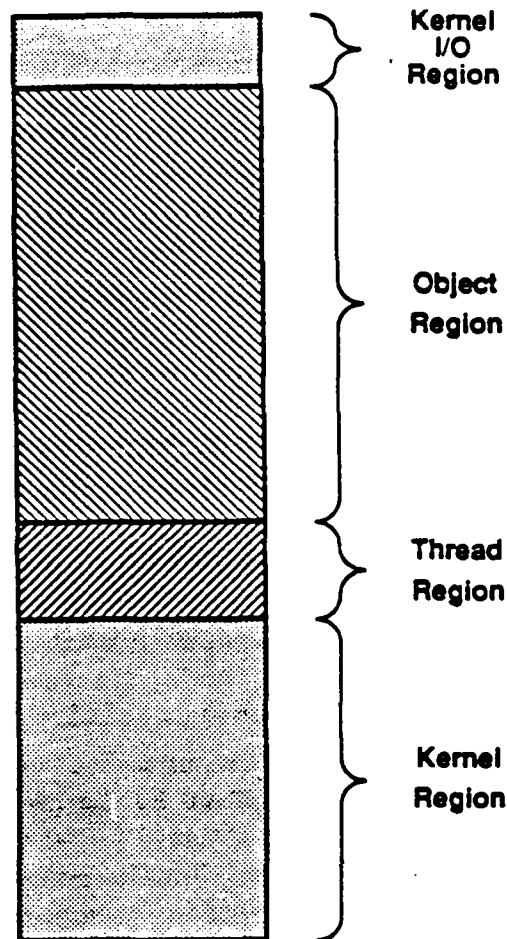


Figure 28: Virtual Address Space Layout

The last two extents of the kernel memory region are known as the *heap extents*, and are where the kernel (dynamically) maintains most of its major data structures, along with the per-object information for the currently existing threads and objects. The memory contained in each of these extents is managed according to a heap discipline. The kernel heap extent provides for the allocation and deallocation of arbitrarily sized units of physical memory within the kernel. The kernel page heap extent is used to allocate and to deallocate storage on the basis of individual pages of memory. The pages of virtual memory obtained from the page heap can be mapped to actual pages of physical memory, or they can be virtual memory place-holders, into which real pages may or may not be mapped. The memory from the kernel heap is used primarily for the creation of control blocks and other kernel data structures. Pages from the page heap are primarily used for the kernel stack and parameter pages used by client threads. The control block that represents the kernel (i.e., the KCB) provides the necessary virtual memory structures to permit the kernel page heap extent to be paged.

The kernel memory region occupies the bottom portion of each context, and this memory is protected to allow access to the kernel region pages only in supervisor mode. All client threads execute in user mode, and the kernel is entered through a processor exception instruction (i.e., a trap, an interrupt, or a bus error). Figure 29 illustrates the virtual memory layout of the kernel region.

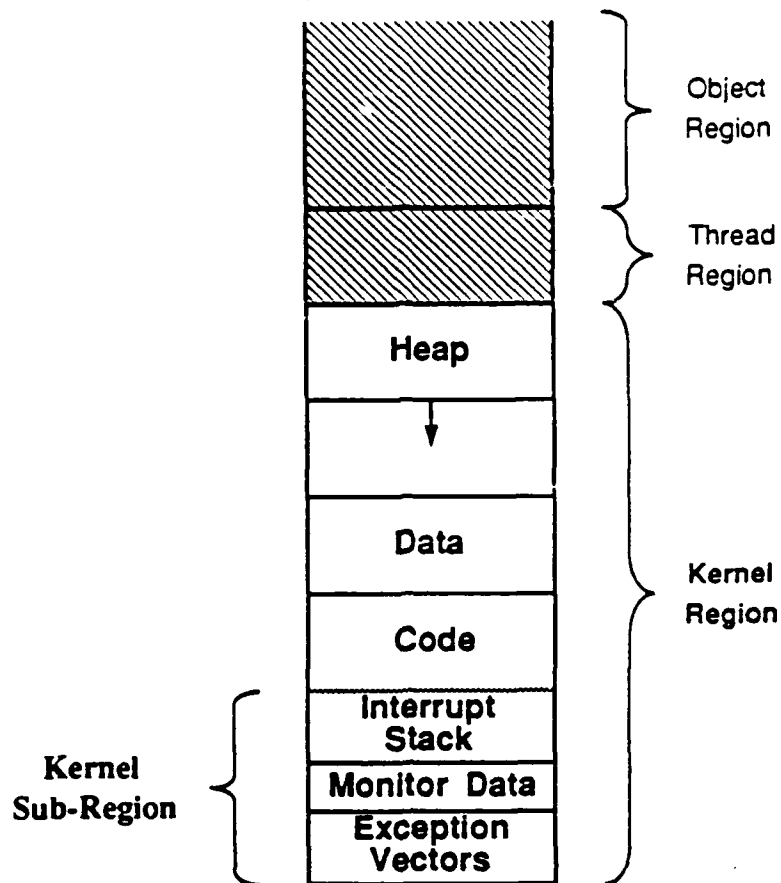


Figure 29: Kernel Virtual Memory Region

6.1.2.2 Client Thread Region

The thread memory region of a context consists of a single extent that provides the thread's user mode stack (i.e., the *client stack*) and the client thread's parameter pages. The thread's client stack is used to provide storage for the automatic variables of threads executing within objects. The client stack is dynamically expandable; it starts out with one page and is expanded each time it overflows, until all pages of the thread's extent is exhausted. Also, to protect the information contained in client stacks across various operation invocations, new client stack pages are allocated on each invocation and the previous ones are protected to prohibit reading and writing. All client threads execute in user mode, therefore the pages in the client thread extent are protected in user read/write mode.

In addition to the client stack, the client thread region contains a number of other pages, including a pair of pages used for the invocation parameter pages, and a *guard page* for the client stack. There is one parameter page for outgoing parameters and one for incoming parameters. When a thread is initialized, there is no highest-level invocation, so the incoming parameter page is provided for each new thread by its creating object. The guard page is a single page, placed between the base of the client stack and the parameter pages, that is used to detect the underflow of the client stack. This is done by protecting the page so that any access to it causes a fault. Figure 30 illustrates the layout of the client thread region of a context.

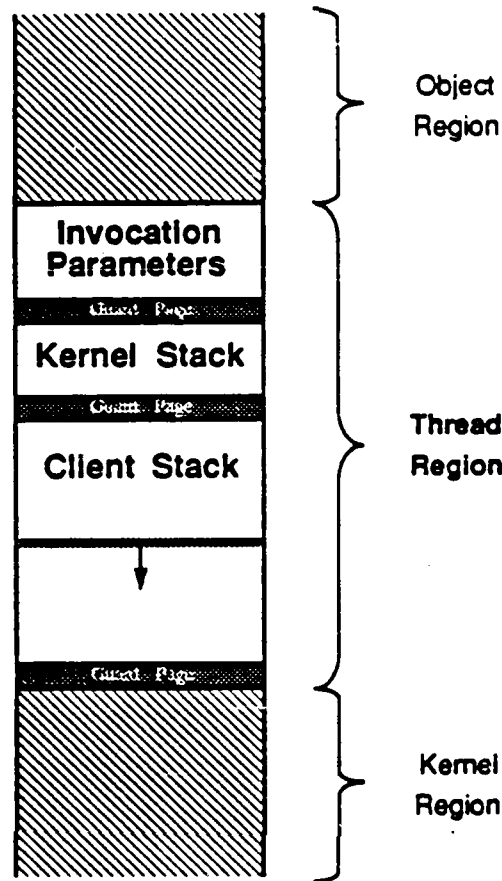


Figure 30: Thread Virtual Memory Region

For a client thread to execute, it must be loaded into a context slot within the MMU. To do this, the kernel must access the thread's control block and obtain the necessary segment descriptors to load into the SMC. The act of acquiring a context for a thread to execute in may involve the unbinding of another thread from its context to make room in the SMC for the thread to be activated. Also, just as pages of physical memory can be wired down, so too can certain contexts be made ineligible for removal from the SMC. In this way it is possible to ensure that a thread's address space can be made ready for execution by only changing the context register.

6.1.2.3 Client Object Region

The client object memory region is the place in a context where a thread's currently active object is located. When an operation is invoked by a thread on a local object, the kernel is entered, the client object currently mapped into the thread's context is mapped out, then the destination object is located and mapped into the invoking thread's context. Then, the appropriate entry point into the object is located and the thread resumes execution at that point within the new object. This process is reversed when an operation completes. While parameter pages are kept in the thread's client stack extent, they are accessible to the client object, and are used to pass parameters on operation invocation.

The client object memory region consists of three extents—the object code extent, object data extent, and the object heap extent. The client object's code extent contains

the object's code and is protected in a user execute-only mode. The client object data extent contains the object's data (both initialized and uninitialized), and is user read/write protected. The object heap extent is where an object's dynamic storage is allocated and this extent is also protected in user read/write mode. Figure 31 provides an illustration of a client thread's region of a context.

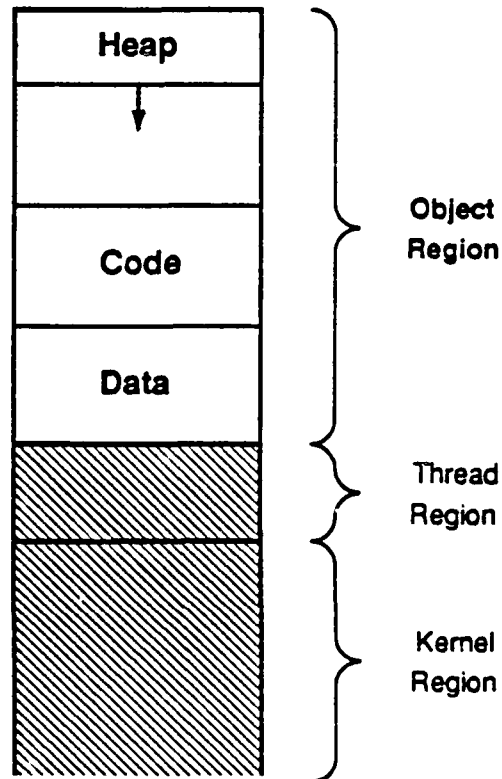


Figure 31: Object Virtual Memory Region

Like client threads, a given client object may at times not be bound to any context. Unlike client threads however, client objects may exist in multiple contexts at any one point in time. Therefore, the client object control blocks must contain references to all of the contexts that the object is currently bound to, as opposed to the reference to the single context required by client thread control blocks.

6.1.2.4 Kernel I/O Region

The final memory region in a context is known as the kernel I/O region. This region is logically a part of the kernel region, but because of hardware constraints, is in a separate location (i.e., exception vectors must be in low-core, and DVMA space is at the top of the virtual address space). This region is managed by the kernel, is not paged, and is protected in supervisor read/write mode. The processor's on-board I/O devices, PROM's, and DVMA accessible memory is mapped into this region.

6.1.3 Address Translation Table

The Application Processor's SMC and PMC are not large enough to contain all of the segment map and page map entries required for each context that must be supported by

the kernel. As a result, the actual virtual memory addressing information is stored in the kernel's virtual memory data structures, and the translation tables are multiplexed among the entities active at a particular node.

In Alpha, the MMU is managed like a cache, and the complete virtual memory information is kept in other data structures within primary memory. The actual virtual memory data structures are always kept up-to-date, and thus it is not necessary to write out MMU descriptors before they are reused. The union of all of the segment map image fields in control blocks, and page map image fields in PMAP's constitute the system's virtual memory data structures. This is a distributed version of the type of memory translation data structure that is usually found in systems where the MMU is implemented as a hardware cache. Because there is no hardware support for the translation cache, there is little point in providing a centralized data structure for the translation information. This distributed approach provides a simple means of managing the necessary mapping information, and simplifies some of the problems associated with the sharing of code.

All modifications of the SMC and PMC are handled by the kernel. The possible reasons for manipulating the MMU include: segment and page faults, the invocation of operations, and binding a thread to a context on activation. When an object or thread is to be loaded into a context, its segment map image is taken from its control block and loaded into the SMC. The necessary page map information is not loaded when a thread is bound to a context, but rather is faulted into the PMC on demand.

In order to set up a client thread's context, the kernel must first enter the new context. This is because hardware restrictions only allow the modification of translation table mappings for the currently active context. Thus a performance penalty is incurred each time the virtual memory mappings of a client context are modified by the kernel.

The MMU caches the address translation information contained in control blocks, and the MMU's tables must be multiplexed among all of the threads active at a given node. To make effective use of the physical resources in the MMU, it must be possible to preemptively remove the segment and page descriptors from the MMU. Similarly, it is possible for the virtual memory facility to restore segment and page descriptors to the MMU on demand (displacing other descriptors where necessary). Whenever translation table entries are removed from the caches in the MMU, they are saved in their primary memory data structures so that they can be restored at a later time.

When a translation table entry is needed in the MMU, and a table is full, one of the existing entries must be selected for replacement. Usually, manipulation of the MMU tables are done on sets of descriptors (i.e., a context's worth of segment descriptors, or a segment's worth of page descriptors). A simple, global (i.e., across the entire cache for a given node) FIFO replacement scheme is used to select the set of descriptors (or *slot*) that is to be replaced from among the set currently eligible for replacement. After a slot has been selected, its information is copied to the necessary control blocks so that the new information can be loaded into its place in the cache.

Some MMU table entries are not eligible for replacement because the information they represent must be mapped at all times. Examples of these types of entries are the segments comprising the monitor, kernel code, and kernel data extents. Other translation table entries must be bound to the MMU for short periods of time, but are otherwise eligi-

ble for replacement. Most notable of these are the descriptors for the segments containing the kernel stack for the currently active thread. Those MMU slots that are never eligible for replacement are removed from the eligible collection at initialization, while the kernel provides a data structure with which those slots that are temporarily ineligible for replacement can be *stuck down* by the kernel.

6.1.4 Virtual Memory Fault Handling

The kernel is responsible for handling all virtual memory faults. When a fault occurs, a bus error trap occurs that switches the processor to supervisor mode and begins execution within the kernel's virtual memory fault handler. Once in the kernel, a determination is made of whether the fault is a segment or a page fault. In the case of a segment fault, the specific type of fault is identified—e.g., an illegal reference, a request to load a segment map, a signal to extend a client stack segment, to fill a page with zeros, or to provide a physical page for a thread's outgoing parameter block page. An illegal memory reference made by a client object terminates the currently active thread, while an illegal reference generated by the kernel causes the application processor to trap to the monitor. In either case, an illegal (page or segment) reference causes a message to be printed on the node's console.

If the bus error indicates a page fault, the nature of the fault is determined—e.g., an illegal reference, a stack extension request, a demand paging request, or a request to provide a parameter page. Illegal page references are handled in a manner similar to how illegal segment references are handled.

Since PMC slots are reclaimed according to a global FIFO replacement policy, a segment fault can result in one PMC slot being emptied in order to accommodate a faulting reference. Also, since a thread's user and supervisor stack pointers must always be valid, the victim selection algorithm of the PMC manager must never select a segment containing the currently executing thread's kernel stack. To insure that this is the case, the notion of a *sticky* PMC slot is introduced—i.e., a slot that is not eligible for replacement. When a thread executes, the kernel always *sticks down* the current thread's stack segments. These segments are then *unstuck* (i.e., made available for replacement) when the thread is descheduled.

In the event of a client stack extension request, the stack is automatically extended (by pages and segments) as long as there is memory space. When there is no longer any (virtual or physical) memory available for the client stack, then the stack extension fault is interpreted as an error condition. Also, when a fault indicates that parameter page is to be supplied by the kernel, a physical page is obtained from the free page list and installed into the appropriate thread's context.

6.2 Secondary Storage Management

In Alpha, secondary memory is managed by the kernel's Secondary Storage Facility (see Figure 2). This facility is internal to the kernel (i.e., does not have an interface that is directly accessible to the clients of the kernel) and provides support for the object abstraction. The major function of the Secondary Storage Facility is to provide and maintain, within secondary storage, images of primary memory extents.

Each memory extent in Alpha has a *storage object* associated with it (that is referenced by an extent descriptor data structure). A storage object acts as a repository for one or more memory extents. Each storage object is physically located in either the transient or permanent area of secondary storage provided by the Secondary Storage Facility, and these images may be updated atomically. These features of secondary storage provide support for the attributes associated with the object programming abstraction.

6.2.1 Secondary Storage Facility Interface

The Secondary Storage Facility provides its services to the rest of the kernel by way of an object-like interface, thereby allowing the operation invocation facility to be used to access the storage objects maintained by the Secondary Storage Facility. The clients of the Secondary Storage Facility therefore need not know the physical location of the storage objects it accesses, they need not be aware of whether a storage object is permanent or atomically updateable, nor must clients know if the storage object is replicated at several secondary storage sites.

Because the kernel uses the operation invocation facility to access storage objects, the Secondary Storage Facility supports the set of operations that may be performed on the facility and its storage objects. The Secondary Storage Facility appears as a kernel object (with a fixed Alphabit identifier, like those of system service objects), and each storage object is given its own unique control block. The Secondary Storage Facility handles the operation invocations directed towards the storage facility as well as those addressed to individual storage objects.

The operations defined on the Secondary Storage Facility are:

- DEFINE_TYPE:** This operation creates a new type storage object, and associates with it the given type identifier. This operation fails if there is insufficient secondary storage to create another storage object, or if there already exists a type storage object with the given name. If successful, this operation creates the control block for this storage object, enters it in the local Dictionary, and returns the Alphabit identifier for the newly created storage object. In Release 1 of Alpha it is not anticipated that new types will be created at run-time. This operation is provided primarily for the adding new types to the secondary storage system in an "off-line" fashion.
- INSTANTIATE:** This operation creates a new instance of a storage object, and associates it with the type storage object that is given as a parameter. This operation returns a failure indication if there is insufficient space to create another storage object, or if there is no type storage object associated with the given identifier. If successful, this operation creates the control block for this storage object, enters it in the local Dictionary, and returns the Alphabit identifier for the newly created storage object.
- DELETE:** This operation is used to remove a specified storage object (of either the type of instance variety). This operation returns a failure indication if the storage object, specified as a parameter is not found. If successful, this operation deallocates the secondary storage associated with the given storage object and deallocates the control block associated with it in the kernel.

In Alpha, replicas of the Storage Management Facility are implemented on those nodes with Secondary Storage Subsystems (i.e., the nodes that have secondary storage devices and Storage Management Co-Processors). The kernel's Secondary Storage Facility behaves as an exclusively replicated client object. Secondary storage operations are performed on a local instance of the Secondary Storage Facility where possible, and from a remotely chosen location otherwise.

Secondary storage devices are connected to nodes in Alpha in much the same way as nodes are connected to the communication subnetwork. Each node's Secondary Storage Subsystem in Alpha consists of a hard disk, a disk controller, and a Storage Management Co-Processor. Taken together, these hardware units can be viewed as a disk with a programmable, caching controller—just as the Communication Co-Processor and the Ethernet interface in the Communication Subsystem can be considered an intelligent network interface. Figure 32 provides an indication of the relationship of the secondary storage devices to the testbed system.

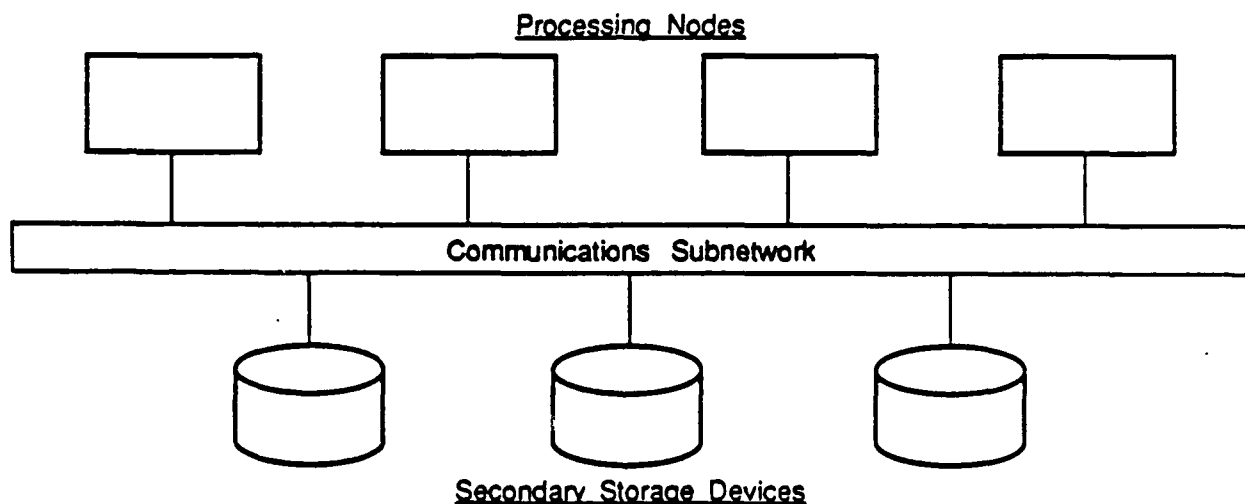


Figure 32: Logical View of Secondary Storage in Alpha

The Secondary Storage Subsystem maintains all of the necessary directory information needed to service the requests made on the facility. The Secondary Storage Subsystem is responsible for managing all of the directory-type information that is necessary to save and retrieve pages to and from storage objects. Also, the Secondary Storage Subsystem is responsible for performing all of the caching of control and data information associated with storage objects. It is responsible for the “short-circuiting” of pages that are queued to be paged-out when a page-in request arrives for them. The Secondary Storage Facility is also responsible for recognizing situations in which pages do not have to be written in order to fulfill page write operations. This may occur when the page to be written is a read only page (e.g., a code page from a type object), or when the page is read/write protected but has not been modified since the last time it was read from its storage object. Additionally, the concurrency provided by the Storage Management Co-Processor makes possible a number of optimizations. For example, the Secondary Storage Subsystem could not only maintain information concerning the physical layout of the information on disks, but could also monitor the position of the heads and the disk's rotational position in order to increase the performance of the subsystem.

6.2.2 Secondary Storage Objects

There are two varieties of storage objects: *type* and *instance* storage objects. Secondary storage objects are registered in the Dictionary of the node at which they exist, just like the objects in primary memory. When the kernel starts up, the only storage objects that exist are those that possess the permanence attribute—type objects and permanent instance objects. Furthermore, when the kernel starts up, updates are completed on all storage objects with the atomic update property that were in the process of being updated when a node failure occurred.

Instance storage objects get their code and initialized data pages from the specified type object, and any number of instance storage objects can share a common type storage object. When the kernel wishes to move portions of an object's extent into or out of primary memory, the Alphabit identifier found in the extent descriptor is used in performing an invocation on the appropriate instance storage object.

The operations defined by the secondary storage facility on storage objects are:

- PAGE_IN:** This operation reads the specified page from a given instance storage object, into the specified physical page, and installs the page in the proper location in virtual memory. This operation fails if the specified storage object cannot be found, or if the specified page is not valid within the given storage object.
- PAGE_OUT:** This operation writes the specified physical page from primary memory into its instance storage object. This operation fails if the specified physical page is not associated with an extent, or if secondary storage is exhausted. If the instance storage object is atomically updatable, the page is not written directly to the storage object, but to a buffer area in the same type of storage (i.e., transient or permanent), and the write is actually completed when an UPDATE operation is performed. Additionally, this operation takes a parameter that indicates whether the page to be written out has locks associated with it. If so, this operation uses the data in the log associated with each lock, in place of that which is actually in the physical page at the time.
- UPDATE:** This operation causes all of the changes made to an object to be performed in such a manner as to not allow the object to be visible in a partially updated state. This operation has no effect if the storage object upon which it is being invoked does not have the atomic update property. This operation deals with all of the pages buffered by the PAGE_OUT operations since the last UPDATE operation was performed. The effects of this operation are achieved by linking in all of the buffered pages, atomically with respect to read operations invoked on the storage object. If the object on which this operation is being performed has the attribute of permanence, the change is registered in a type of *intentions list* and multiple reads and writes are used to ensure that all of the changes are made, even across node failures.

It is less expensive to manage the transient portion of secondary storage than it is to manage the permanent portion. This is because no redundancy is required for the disk directories of transient storage, the disk sectors of transient storage can be heavily cached, and the cached sectors do not have to be written to disk except to make room in the cache. On failure, the transient portion of an object must be reinitialized, but no complicated recovery procedures are required.

The permanent part of secondary storage is managed in a manner similar to that used in traditional file systems. Redundant disk directory information is maintained on the disk to keep from losing objects as a result of a node failure. Care must be taken in performing reads and writes to ensure that once an object has been written, it will retain its state across failures. To this end, a technique similar to that for providing stable storage as described in [Lampson 81] is used. This makes the read and write operations on permanent storage objects more costly than for transient storage objects.

After a failure, the permanent portion of secondary storage must be restored to a consistent state, control blocks must be created in primary memory, initialized, and entered in the local node's Dictionary for all of the permanent objects in the recovering node's secondary storage. All pending atomic update operations on permanent objects must be completed before references to them are placed in the Dictionary.

While the attributes associated with objects are usually defined when they are created, it is possible to change object attributes at run time. These changes must be reflected in the instance storage objects. For example, when a transient object is converted into a permanent object, a permanent storage object must be created, the contents of the transient storage object must be copied into it, the object identifier in the proper extent descriptors must be changed, and the transient object should be deleted.

The operations performed on secondary storage objects involve accessing secondary storage devices, that typically have characteristics that necessitate special treatment (e.g., average access times on the order of 30 milliseconds for disks). To deal with this, special management algorithms are frequently used (e.g., to minimize head motion) which make it necessary to handle secondary storage operations in a sequence different from that in which they were requested. Furthermore, the same thread environment information that is used by the Scheduling Subsystem to resolve contention for processing cycles is used by the Secondary Storage Subsystem to resolve contention for the secondary storage resources. In particular, the order in which disk access requests are serviced is influenced by the characteristics of the threads making the requests.

When an operation is performed on a storage object, the Secondary Storage Facility ensures that the invocation returns as soon as possible. In some cases, the request embodied by the operation invocation has been completed when the invocation returns, while in other cases the return indicates only that the request has been noted and it has been queued. As an optimization, the reply from invocations of operations on storage objects has been decoupled from the indication of the completion of the operation. When an operation on a storage object actually completes, the Secondary Storage Facility executes a V operation on a semaphore that is associated with the thread whose request has been satisfied. For this purpose, each thread control block has associated with it a semaphore dedicated to the invocation of secondary storage operations. Should the thread then wish to wait until the request has completed, it can issue a P operation on its secondary storage semaphore. A C_P operation can be issued by the thread to poll for the completion of the secondary storage operation.

When a thread page faults, the faulting thread enters the kernel via the fault handler and invokes a PAGE_IN operation of its instance object. Upon successful return from the invocation, the thread immediately issues a P operation on its secondary storage semaphore. If the page has not yet been installed, the thread is blocked and will be

unblocked by the kernel when the page has been installed. In this approach, the thread does not automatically relinquish the processor when it makes a system call that could take a long time. Thus, if the requested page is easily available, the Secondary Storage Facility installs it and issues a V operation on the thread's semaphore prior to returning from the invocation, thereby not incurring the overhead of rescheduling the thread. On the other hand, the thread is blocked if the Secondary Storage Facility cannot ensure the prompt installation of the requested page.

The page daemon's thread (see Subsection 3.2.1.1) is typically blocked until the level of the physical page pool crosses its lower limit. When this happens, the page daemon thread is unblocked and goes on to identify a number of victim pages to be written to secondary storage. The page daemon thread invokes a PAGE_OUT operation on the instance object of each victim it chooses. The page daemon need not wait for each request to complete, instead, it issues a C_P operation on all but the last of its page-out requests, issuing a P operation following the last request. In this fashion, the page daemon can queue a large number of page-out requests, have them serviced in arbitrary order, and wait for them all to complete before continuing. Once all of the secondary storage operations have been performed, the page daemon may then become blocked once again until the pool's level crosses its lower threshold.

6.2.3 Low-Level Storage Management Details

While the Secondary Storage Facility was designed to support a wide range of secondary storage technologies, the only one addressed in Release 1 of Alpha is standard rotating magnetic media devices—i.e., Winchester disks. This subsection describes the manner in which information is represented on the physical device(s) in Release 1 of Alpha. Since an Alpha system is expected to be reliable, the issue of representation is particularly important. The chosen format represents a convention, which all software and hardware components must adhere to.

The major fact which gives the disk format conventions importance, is the fact that all the components of a system must occasionally be initialized and go through a start-up phase. This is as true of the data on the secondary storage devices as it is of any other component. However, the data are distinguished by the times at which this occurs. A data collection is initialized at the discretion of a system administrator, and the initialization is performed by a running system. A system is initialized every time that it is started up, and every time that circumstances force a restart, but the data in secondary storage are not initialized at that time. Instead, the data are treated as a special resource whose logical correctness must be assured (or assurable to some degree of confidence).

The considerations that went into the design of the low-level storage formats are (in order of importance): robustness, timeliness, restartability, speed efficiency, and space efficiency. The robustness item reflects the fact that hardware failure is of primary concern of Alpha. It is unacceptable that such failures have a high probability of causing undetectable or unrecoverable damage to data. Since some data are of transitory importance, and since some robustness techniques exact a higher performance penalty than others, it was clear that no single technique offered the most appropriate trade-off for all uses. Therefore, Alpha should allow for several different robustness tradeoffs. In this context, timeliness means that the secondary storage subsystem manages its resources

based on the time constraints of the threads on whose behalf operations are being performed. Furthermore, some attempt has been made to limit the variance of delay associated with the execution of secondary storage operations. In particular, an attempt was made to help avoid "recursive disasters", where a whole chain of events may (with low probability) all have to be serviced before some single request is satisfied. In Alpha, a Secondary Storage Subsystem should not spend an inordinate amount of time performing the various housekeeping activities in order to resolve the consequences of the last shut-down/failure. Instead, these actions are discretionary, or are performed in parallel with other work (potentially as a lower priority activity which can be interrupted/resumed (or aborted/restarted) consistently with the system's higher-level goals).

For Alpha Release 1, it was decided to implement a Secondary Storage Subsystem which was extensible to high sophistication, and to very large disks, but which could be useful to the project at several levels of partial implementation.

The primary robustness technique of the initial implementation is "stable storage" on a non-replicated disk. This technique involves writing two copies of a given datum to the same disk, each copy separately check-summed and time-stamped. Logging techniques are planned for integration with transaction support, but this was not deemed as appropriate for initial low-level requirements. A disk failure model was adopted, which permits a write of one sector to be corrupted, and a single write of N contiguous sectors may leave corrupted data in two successive sectors. The model assumes these effects occur because the write was in progress at the time of the hardware failure. The model assumes that a failure at any other time will leave the stored data uncorrupted. It is also assumed that the failure will leave the subsystem unable to perform further writes. Given this failure model, it was decided that stable storage of a small datum would involve a single write of three contiguous sectors, with null data stored between two identical copies of the datum. All accesses to such a datum must read both copies, and verify both check-sums, and compare the data. If only one datum is acceptable, the other is "healed." If both data are acceptable, but they differ, then the younger or first is assumed to be dominant, and the other is immediately corrected on disk. It was furthermore decided that stable storage of multi-sector data, that is written one sector at a time, would involve two writes. No guard region is used in these cases. Stable storage of multi-sector data, that is written in a single write, has not been needed for low-level purposes.

A basic consequence of the restartability requirement, and of the failure model, is that all low-level data structures on disk must be weakly self-consistent at all times. By this, it is meant that a restart must not have to perform a complete scan of the low-level data, before knowing for certain that given datum can be used in the normal way. (For example, a directory lookup can be performed, without recovery analysis of any components not encountered by a normal lookup.) This is referred to as "weak" self-consistency, because we are willing to allow storage losses. It is an acceptable consequence of a crash that some bounded amount of empty disk space become inaccessible. It is intended that the processing to detect and correct this loss will be discretionary, or low-priority.

Since the Secondary Storage Subsystem is primarily oriented towards objects and paging, rather than towards small text files, it was decided that allocation would be based on extents. That is, device space would not be treated as equal-sized chunks, but would be

treated as varying-length items, potentially as large as an entire disk. This appeared to satisfy several criteria. By allowing compact descriptions of larger entities, access structures could be smaller and therefore traversed in fewer steps, or even completely cached. As will be seen, access times could be bounded, and large data transfers could be performed with considerable directness, a performance consideration, since cached writes have bad failure properties. The structure of objects in memory can (usually) be directly mapped onto disk extents. Since lists of free space are compact, in-memory representations could be kept in address order, facilitating contiguous allocation. A high-efficiency data base subsystem could be allowed low-level control within a dedicated extent.

It was also decided that version codes would be placed into data structures, in such a way that interoperability problems would not be caused by removable disk packs.

6.2.3.1 Device Header

It is important that every writable device be identified as to its history. Therefore, Alpha writes a label at the head of each device (disk) when the device is formatted. This label uses the stable storage convention, and it is intended that it contain a variety of information, such as:

- the version of the conventions holding for the contained data
- the time-stamp of the formatting operation
- a time-stamp bounding the last time that the device was written to
- identification of the disk, and some history, if removable
- information bounding the range(s) of capabilities present on the disk

Information concerning the faults on the device has been stored separately, in the manner which is standard for the device in the Alpha testbed. This was intended to allow the use of a commercially obtained diagnostic program, which deals with "bad blocks" in an acceptable fashion.

Information concerning the geometry of the device is already present in the fault map, and has not yet been represented in the Alpha label.

6.2.3.2 Space Allocation

Since space is dealt with in extents, an empty disk contains exactly one extent. (Actually, it contains exactly two, because some structures are kept in the middle of the disk, to minimize the access time of seeks to these structures.) Therefore, the initial "free list" is only a few bytes long (16 bytes, in our case). When an extent is allocated to some purpose, the initial free list neither grows nor shrinks; instead, it describes less free space, since one of the length counts in the list is decreased.

Matters become more complicated when an extent is freed. The simple solution would be to merge the newly freed extent into one of the existing free extents. However, due to fragmentation, the newly freed extent may not be contiguous to an existing extent, and the free list must become longer. It is in this manner that the list lengthens, and in theory it could lengthen until we are spending 8 bytes to describe every second disk sector: that is, one part in 128 of the disk. In practice, no such explosion is likely, given the expected pattern of demand from dedicated real-time applications, and the time scales involved in that many demands for space. The initial implementation allows the free list to grow to

62 extents, the maximum that will fit into a single disk sector. The data structure is designed to be enlarged, and expansion space has been left on the disk. A background or privileged "scavenger" program is envisaged, which is capable of moving small disk objects around on the disk, thus allowing free extents to be abutted and merged. As long as an implementation allows much more free-list space than the "equilibrium" amount, this approach should have no more problems than any alternative. This has been weighed against the benefits, which include highly contiguous allocation (for fast reads and writes), and fast allocation (since alternative approaches must allocate a larger number of smaller pieces).

Given that there is such a free-list, it was decided to keep it sorted in address order (for easy merging on deallocations), rather than in size order (for best-fit allocation).

Next, it was determined that deallocation could be *lazy*—i.e., the free-list can become disordered. A deallocation merely places the freed extent(s) at the rear of the free-list, and increment a "dirty count" which is internal to the free-list. When the request interpreter has no outstanding work, it will check if the free-list needs to be rationalized, and perform this operation. Rationalizing removes any zero-length entries, and performs the equivalent of a sort. It would also be quite easy to allow rationalization to be interrupted when work comes in.

Ideally, the in-memory free-list would correspond directly to the on-disk free-list. However, it is acceptable to use lazy writes after deallocations, since a crash will at worst lose the deallocated space. The space is not lost forever, of course: it is only lost until a scavenger process computes a correct free-list from the access structures, and writes the correct free-list over the incorrect one. This scavenger does not have to be run immediately after a restart, in accordance with the requirement for restartability.

Lazy writes are also acceptable after free-list rationalization. However, if space was allocated rather than deallocated, lazy writes could cause a problem. After a crash, the free-list on disk might claim the space was free, while the access structures claim that some viable object resides there. This is acceptable if we are willing to run a scavenge procedure on restart, but as has already been stated, we do not want to do this. An alternative is to always write allocation information immediately out to disk, either by writing the free-list, or by using a logging technique. This is speed-limiting, so we prefer to avoid it. The alternative which has been chosen is to write the free-list to disk, but only if an optimization method fails.

The optimization method is simply a pre-allocation scheme. Some number of disk blocks are allocated on start-up, and the extent(s) so allocated are noted in memory. Any sufficiently small allocation requests may be given the pre-allocated space, without causing a disk write, and the pre-allocation state is then marked for lazy renewal. As before, a crash may temporarily lose some amount of disk space. The only time penalty is that start-up/restart must now perform allocations, thus slowing restart by one disk write. When an object is created, the allocation process will involve exactly one disk write, which typically will be performed lazily, that is, at some time after the allocation has been reported as successful.

The free-list data structure is kept in the middle of the disk, to minimize seek times. The middle of the disk can be found because the software knows the disk geometry, either from the disk label, or by other means.

6.2.3.3 Capability Mapping

The "name" of an object on secondary storage, is the Alpha capability which was specified when the object was created. At the level which the kernel deals with, no other style of name is interesting, and it is assumed that many objects would never have any other sort of name. Of course, higher levels of software are free to maintain mapping data, which allows application-specific names (or human-readable names) to be rendered into capabilities. In general, enough tools will exist, and an application could, for example, maintain a tree-structure in the style of UNIX. No specific support for this needs to exist at the kernel level.

The Secondary Storage Subsystem must nevertheless maintain one mapping, which is the mapping from capabilities to objects. This capability mapping information—the access structure—must be kept reliably on the disk. Some systems have used B-trees to do this, however, in Alpha this approach was rejected on the grounds of unpredictability: an update may cause an update, and this may ripple through some large fraction of the B-tree. B-trees do have the advantage of allowing sequential access, but this is not needed: the indexes are capabilities. B-trees are also bulky, in the sense that root nodes must store at least portions of the indexes, whereas the scheme below does not. Also, B-trees are somewhat difficult in an environment where one thread may be updating the tree while others are using it. (The difficulty is another aspect of the ripple problem: deadlock can ensue if a writer does not initially lock almost the whole tree.)

It was decided that we would use a two-level version of a technique called *extendible hashing*. This technique is applicable when the concepts of "next" and "previous" index are not needed. This is because the method uses hashing, which destroys sequential properties, but has the advantage that the index values (here, capabilities) are more uniformly distributed across the access structure. This seemed useful because capabilities tend to have internal structure generated by upcounters, and this means that a capability treated as a bit pattern may have difficult statistical properties. By hashing these bit patterns, a useful statistical property is at hand.

The basic method requires that a *root map* be kept, which contains N *leaf pointers*. N must be a power of two, for reasons which will be apparent shortly. A lookup is done by indexing (with bits from a hashed capability) to select one of the leaf pointers. A leaf pointer is merely enough information to allow us to compute the disk address of any of L *leaves*, where a leaf is a data structure occupying S disk sectors. (In our initial implementation, S is 1.) Each leaf contains a hash-searchable list of capabilities, and attached to each capability is the disk address of an *object descriptor*. (The object descriptor is stored contiguously with (one of) the extent(s) which make up the actual object.)

In the initial implementation, this structure is capable of addressing $128 * 63 = 8064$ object descriptors. This was acceptable as a first cut, since the structure would not saturate before the 64 megabyte disk was full, unless the average object took up less than 7800 bytes. It is not acceptable for larger disks, but there are good techniques for expanding this basic structure.

Several techniques fall under the heading of "efficiency" methods. For example, the leaf pointers could be compressed to B bits, by assuming that all leaves are stored in some bounded region of the disk. (That is, a base pointer having more than B bits is added to the leaf pointer, to convert it to a disk address.) Also, the leaves could have capabilities in sorted order, which allows compression techniques, at some penalty in (in-memory) search time. Collectively, efficiency techniques can expand the structure by a factor of (say) four, which is approximately ten times short of the desired size. In general, these methods seem most appropriate after an acceptable structure has been implemented, and time is found to make small improvements in efficiency.

The crude approach is to use a larger fixed amount of disk space. Direct scaling of each leaf, and of the root map, to four kilobytes each, would be enough. This is actually inefficient, because the structure is likely to be very sparsely populated, leading to inefficient use of disk bandwidth and inefficient caches.

The reasonable thing is to use a compact layout, and expand dynamically to use more disk space as needed. Extendible hashing offers a pleasant algorithm for this, which avoids the potentially high costs of dynamic extension. The article suggests that the root map be initially quite small, and contain identical leaf pointers. (For example, it might be of size 4, with all four pointers pointing to the same leaf.) As objects are created, leaves become more full, and eventually a leaf will overflow. At this time, we determine P , the number of pointers to this leaf. If P is greater than 1, then it will be a power of two. In this case, we split the leaf into two leaves, and adjust half of the P pointers to point to one of the new leaves, and adjust the other half of the pointers so that they point to the other of the new leaves. (The entries in the original leaf will be assigned one way or the other according to where their hash now falls.)

If P was one, then only one slot in the root map yielded that leaf. This case is handled by creating a double sized root map, and copying each pointer in the old map to two adjacent positions in the new map. The P value is now two, and the overflow algorithm given above can be applied. It can now be seen why the root map was required to have a size which is a power of two. The lookup algorithm takes the hash of a capability, and uses several bits of the hash value as an index. When the root map grows (or shrinks), the lookup simply has to start taking one more (or less) bit from the hash value. No other change needs to be made to insure that previously stored capabilities can still be found.

This scheme is logically simple, and should be quite fast, because a leaf split involves at most three disk writes (two leaves and the new root). It also has highly localized disk activity, not only during leaf splits, but at all times, including restart. The root map stays small enough that it can be kept in memory: hence, a lookup performs at most one disk read to obtain the address of an object descriptor. The phrase "at most" is used because caching techniques will be applied to leaves. It should be noted that extendible hashing allows lazy load balancing, so that all leaves can be kept at a "reasonable" fullness. In a system which is operating below overload, it should be possible to assure that all disk splits are performed lazily, as part of the background load balancing. This would mean that object creation could never cause more than one write to the access structure—surely, some sort of minimum, since the object's existence is not committed until some write occurs.

It is intended that leaves will be kept in the middle of the disk, in a statically-allocated region of adequate size. We are not aware of any convincing argument about how to guarantee an advantage from allowing leaves to be scattered. (There are non-extent-based systems, which keep a local map on each cylinder of their disk, but this does not seem appropriate for Alpha.)

It was stated earlier that the object descriptors are not stored in the leaves, but are instead stored at the front of each object (or front portion thereof). Other systems have decided otherwise, and argued that this allowed a lookup to take one less disk read. It has also been argued that in-leaf storage is very efficient for the most common operation on these file systems, namely, a directory scan. However, the Alpha system will not attempt such scans. Also, if the object data are stored in the leaves, then the access structure becomes bulkier.

In fact, an Alpha object descriptor is a fair-sized varying-length structure, since Release 1 objects may consist of up to four segments, each of which may be stored on disk as multiple disk extents. This size, and the algorithmic implications of varying size, persuaded us that object descriptors should be stored separately. This was expected to be useful as the implementation of various higher-level schemes eventually required the augmentation of the object descriptor format.

It was also expected that later implementations would be on machines with enough memory to allow whole-track transfers, and track caching. (This is done for the speed advantage, since a track operation does not necessarily take any longer than a sector operation.) On such a system, the design is expected to show performance improvements that might not be realized by alternate designs.

7 Acknowledgments

Many people contributed to this effort, and many others continue to contribute as the Archons project enters its next phase.

Doug Jensen is the founder of the Archons project. He obtained the support for the project over its 10-year existence, and provided most of the fundamental philosophy and concepts for Alpha. Doug continues to provide support and guidance for this research effort as Release 2 of the Alpha operating system is being developed at Kendall Square Research Corporation.

In his brief visit with the Archons project, Martin McKendry initiated the implementation effort that has become Alpha and provided many of the initial implementation concepts.

Donald Lindsay worked on the secondary storage subsystem and contributed text to Subsection 6.2.3 of this report.

Sam Shipman worked on various subsystems and support tools, and contributed to refining and completing the system design and implementation. David Maynard also worked on various aspects of the system's design and implementation, and acted as the project liaison with General Dynamics during our joint C² demonstration effort. David is now working on a thesis in the area of real-time scheduling for decentralized computers with multiprocessor nodes. Huay-Yong Wang worked on the scheduling subsystem and the design and implementation of various other system functions.

Other project members that contributed to this effort are Jeff Trull, Chuck Kollar, Bruce Taylor, and Dan Reiner. Thanks are due to Tom Lawrence and Dick Metzger, the Archons project's prime sponsors at the Rome Air Development Center.

References

- [Accetta 86] Accetta, M., Barton, R., Bolosky, W., Golub, D., Rashid, R., Tevanian, A. and Young, M.
Mach: A New Kernel Foundation for UNIX Development.
In *Proceedings of Summer Usenix*. July, 1986.
- [Almes 85] Almes, G. T., Black, A. P., Lazowska, E. D. and Noe, J. D.
The Eden System: A Technical Review.
IEEE Transactions on Software Engineering SE-11(1):43-58, January, 1985.
- [Bochmann 78] Bochmann, G. V.
Finite State Description of Communication Protocols.
In *Proceedings of the Computer Network Protocols Symposium*, pages F3-1-F3-11. February, 1978.
- [Bochmann 80a] Bochmann, G. V.
A General Transition Model for Protocols and Communication Services.
IEEE Transactions on Communications COM-28(4):643-650, April, 1980.
- [Bochmann 80b] Bochmann, G. V. and Merlin, P.
On the Construction of Communication Protocols.
In *Proceedings of the Fifth International Conference on Computer Communication*, pages 371-378. October, 1980.
- [Bochmann 82] Bochmann, G. V. et. al.
Experience with Formal Specifications Using an Extended State Transition Model.
IEEE Transactions on Communications COM-30(12):2506-2513, December 1982.
- [Brand 83] Brand, D., and Zafiropulo, P.
On Communicating Finite State Machines.
Journal of the ACM 30(2):323-342, April, 1983.
- [Chesson 80] Chesson, G. L. and Fraser, A. G.
Datakit Network Architecture.
In *COMPCOM 80 — Spring, IEEE Computer Society International Conference*, pages 59-61. February, 1980.
- [Chesson 87] Chesson, G. L.
The Protocol Engine Project.
UNIX Review 5(9):70-77, September, 1987.
- [CCITT 85] CCITT.
Recommendations Z.100 to Z.104: Specification and Description Language.
In *Red Book*. Geneva, Switzerland: ITU, 1985.

- [Clark 82] Clark, D. D.
Modularity and Efficiency in Protocol Implementation.
RFC 817, MIT Laboratory for Computer Science, July, 1982.
- [Clark 88] Clark, R. K., Kegley, R. B., Keleher, P. J., Maynard, D. P., Northcutt, J. D., Shipman, S. E. and Zimmerman, B. A.
An Example Real-Time Command and Control Application.
Archons Project Technical Report #88032, Department of Computer Science, Carnegie-Mellon University, March 1988.
- [CMC 85] Communication Machinery Corporation.
Ethernet Node Processor ENP-30 User's Guide.
Communication Machinery Corporation, April, 1985.
- [Danthine 81] Danthine, A.
Design Principles of Communication Protocols.
In *Data Communications and Computer Networks*, pages 257-273, IFIP, 1981.
- [Fleisch 81] Fleisch, B. D.
An Architecture for Pup Services on a Distributed Operating System.
Operating Systems Review 15(1):26-44, January, 1981.
- [Fletcher 78] Fletcher, J. G. and Watson, R. W.
Mechanisms for a Reliable Timer-Based Protocol.
In *Proceedings of the Computer Network Protocols Symposium*, pages C5-1-C5-17. February, 1978.
- [Fletcher 79] Fletcher, J. G.
Serial Link Protocol Design: A Critique of the X.25 Standard, Level 2.
Technical Report UCRL 83604, Lawrence Livermore Laboratory, August, 1979.
- [Herlihy 84] Herlihy, M. P.
Replication Methods for Abstract Data Types.
PhD thesis, Massachusetts Institute of Technology, 1984.
- [Hopcroft 79] Hopcroft, J. E. and Ullman, J. D.
Introduction to Automata Theory, Languages, and Computation.
Addison-Wesley Publishing Company, 1979.
- [Gouda 85] Gouda, M. G.
Protocol Validation by Fair Progress State Exploration.
Computer Networks and ISDN Systems 9(5):353-361, May, 1985.
- [ISO 86] International Standards Organization.
Estelle — A Formal Description Technique Based on an Extended State Transition Model.
ISO TC 7/SC 21/WG 1 — FDT, Subgroup B, September 1986.
- [Jensen 75] Jensen, E. D.
Time-Value Functions for BMD Radar Scheduling.
Technical Report, Honeywell System and Research Center, June 1975.

- [Jensen 78] Jensen, E. D.
The Honeywell Experimental Distributed Processor — An Overview.
Computer 11(1):28-37, January, 1978.
- [Jones 79] Jones, A., Chansler, R., Durham, I., Schwans, K. and Vegdahl, S.
StarOS, a Multiprocessor Operating System for the Support of Task Forces.
In *Proceedings, Seventh Symposium on Operating System Principles*, pages 117-127. ACM, December 1979.
- [Lampson 81] Lampson, B. W., Paul, M. and Siegert, H. J. (editors).
Lecture Notes in Computer Science. Volume 105: *Distributed Systems—Architecture and Implementation*.
Springer-Verlag, Berlin, 1981.
- [Liskov 84] Liskov, B. H.
Overview of the Argus Language and System.
Programming Methodology Group Memo 40, MIT Laboratory for Computer Science, February, 1984.
- [Locke 86] Locke, C. D.
Best-Effort Decision Making for Real-Time Scheduling.
Ph.D. Thesis, Department of Computer Science, Carnegie-Mellon University, May 1986.
- [McKendry 84] McKendry, M. S.
The Clouds Project: Reliable Operating Systems for Multicomputers.
Project Report, Georgia Institute of Technology, 1984.
- [Nash 83] Nash, S. C.
Automated Implementation of SNA Communication Protocols.
In *Proceedings of IEEE International Conference on Communications*, pages 1316-1322. June, 1983.
- [Nelson 81] Nelson, B. J.
Remote Procedure Call.
PhD thesis, Carnegie-Mellon University, May, 1981.
- [Northcutt 87] Northcutt, J. D.
Mechanisms for Reliable Distributed Real-Time Operating Systems: The Alpha Kernel.
Academic Press, Boston, 1987.
- [Northcutt 88a] Northcutt, J. D.
The Alpha Operating System: Requirements and Rationale.
Archons Project Technical Report #88011, Department of Computer Science, Carnegie-Mellon University, January 1988.
- [Northcutt 88b] Northcutt, J. D. and Clark, R. K.
The Alpha Operating System: Programming Model.
Archons Project Technical Report #88021, Department of Computer Science, Carnegie Mellon University, February, 1988.

- [Northcutt 88c] Northcutt, J. D.
The Alpha Distributed Computer System Testbed.
Archons Project Technical Report #88033, Department of Computer
Science, Carnegie-Mellon University, March, 1988.
- [Northcutt 88d] Northcutt, J. D.
The Alpha Operating System: Kernel Programmer's Interface Manual.
Archons Project Technical Report #88061, Department of Computer
Science, Carnegie-Mellon University, June 1988.
- [Postel 80] Postel, J. B.
User Datagram Protocol.
RFC 768, USC/Information Sciences Institute, August, 1980.
- [Postel 81a] Postel, J. B.
Internet Protocol.
RFC 791, USC/Information Sciences Institute, September, 1981.
- [Postel 81b] Postel, J. B.
Transmission Control Protocol.
RFC 793, USC/Information Sciences Institute, September, 1981.
- [Ritchie 74] Ritchie, D. M. and Thompson, K.
The UNIX Time-Sharing System.
Communications of the ACM 17(7):365-375, July, 1974.
- [Schantz 85] Schantz, R., Schroder, M., Barrow, M., Bono, G., Dean, M., Gurwitz,
R., Lebowitz, K. and Sands, R.
*CRONUS: A Distributed Operating System: Interim Technical Report
No. 5.*
Technical Report 5991, Bolt Beranek and Newman, June, 1985.
- [Shipman 88] Shipman, S. E.
The Alpha Operating System: Programming Language Support.
Archons Project Technical Report #88042, Department of Computer
Science, Carnegie-Mellon University, April 1988.
- [Shrivastava 82] Shrivastava, S. K. and Panzieri, F.
The Design of a Reliable Remote Procedure Call Mechanism.
IEEE Transactions on Computers C-31(7):692-697, July, 1982.
- [Sproull 78] Sproull, R. F. and Cohen, D.
High-Level Protocols.
Proceedings of the IEEE 66(11):1371-1386, November, 1978.
- [Thornton 64] Thornton, J. E.
Parallel Operation in the Control Data 6600.
In *Proceedings of the AFIPS Fall Joint Computer Conference*, vol. 26,
part 2, pages 33-40. 1964.
- [Wulf 81] Wulf, W. A., Levin, R. and Harbison, S. P.
Hydra/C.rump: An Experimental Computer System.
McGraw/Hill, New York, 1981.

[Xerox 81]

Xerox Corporation

Courier: The Remote Procedure Call Protocol.

Technical Report XSIS 038112, Xerox Corporation, December, 1981.

Table of Contents

1	Introduction	C-1
2	Alpha Abstractions	C-2
3	Object Invocation	C-3
4	Capabilities	C-5
5	Exception Handling	C-6
6	Standard Operations.....	C-8
	Update	C-8
	Restore	C-9
	Info	C-9
	Migrate	C-10
	Lock	C-10
	C_Lock.....	C-11
	Unlock.....	C-11
	PreCommit.....	C-12
	Commit	C-12
	Abort	C-13
	AllocateHeap.....	C-13
	DeallocateHeap	C-13
	SaveCapa.....	C-14
	ReleaseCapa.....	C-14
	Delete	C-14
7	Kernel Services.....	C-16
	Schema	C-17
	Create	C-20
	CreateSystem	C-21
	AddExtent	C-22
	AddOperation.....	C-23
	AddRelocation	C-23
	AddSymbol	C-24
	Validate	C-24
	Delete	C-25
	Schema Manager.....	C-26
	Create	C-27
	CreateUnrestricted	C-27
	Semaphore.....	C-29
	P	C-30
	V	C-30
	V_All	C-30

Conditional_P	C-31
Delete	C-31
Semaphore Manager	C-32
Create	C-33
Thread	C-34
Stop	C-35
Delay	C-35
Resume.....	C-35
BeginScope	C-36
EndScope	C-36
Deadline	C-37
Mark.....	C-38
AllocateHeap.....	C-38
DeallocateHeap.....	C-39
Abort	C-39
Thread Manager	C-40
Create	C-41
Transaction Manager	C-42
BeginTransaction	C-43
EndTransaction	C-44
AbortTransaction	C-44
8 Glossary of Terms	C-46
References	C-48

1 Introduction

This document describes the application visible interface to Release 2 of the Alpha kernel. This interface includes the operations defined on the kernel objects, kernel supported properties of application objects, object invocations, capabilities, threads and exceptions. This work was done as part of a subcontract to Carnegie-Mellon University.

The scope of this document is limited to the programmer interface to Alpha. It does not explicitly detail the design, programming model, or behavior of Alpha. For more information see [Northcutt 88a] and [Northcutt 88b].

The interface is specified in a language- and architecture-independent fashion. It is intended that programmers write applications for Alpha in different languages, such as C, Ada, Lisp, Fortran. Each language that is supported on Alpha will have a language specific document describing the interface and run-time model. It is also intended that Alpha run on a variety of architectures. A document describing architecture-specific information (e.g., size of virtual address space, granularity of the clock, and the basic units of addressability) must be produced when Alpha is implemented on a particular architecture.

This document is organized by chapter. Chapter 1 is the introduction to the document. Chapter 2 is a very brief summary of the fundamental Alpha abstractions. Chapter 3 is a discussion of operation invocation. Chapter 4 describes the Alpha notion of a capability. Chapter 5 details exceptions and exception handling. Chapter 6 is a description of standard operations defined by Alpha. Chapter 7 enumerates in detail each of the kernel objects and their operations. Chapter 8 is a glossary of terms.

2 Alpha Abstractions

The abstractions provided by the Alpha kernel are based on a combination of the principles of object-orientation, atomic transactions, replication, and decentralized real-time control. The Alpha kernel interface presents its clients with a set of simple and uniform programming abstractions from which reliable real-time control applications may be constructed. The kernel mechanisms support an object-oriented programming paradigm, where the primary abstractions are *objects*, *operation invocation*, and *threads*.

At the highest level of abstraction, objects in the Alpha kernel are equivalent to abstract data types. Objects are written as individual modules composed of the specific operations that define their interface. Alpha applications can dynamically create and destroy object types and instances of those types, via *Schema objects*. The object abstraction in Alpha extends to all system services, and encapsulates all of the system's physical resources, providing clients with object interfaces to all system-managed services.

All interactions with both client and system objects are via the invocation of operations on objects. The operation invocation mechanism is the fundamental facility on which the remainder of the Alpha kernel is based. The invocation of operations on objects is controlled by the kernel through the use of a capability mechanism. Capabilities provide the only means in Alpha for one object to invoke operations on another object. In this way, access to objects can be controlled by controlling the distribution of their capabilities.

Threads are the run-time manifestations of computations. They are the unit of activity, concurrency, and schedulability in the kernel. The combination of a thread with an object's address space is similar to the commonly used notion of *process*. Threads, however, may move among objects via operation invocations, regardless of the physical nodes on which the individual objects may reside. Threads execute asynchronously with respect to each other, allowing a high degree of concurrency within single objects as well as between objects.

3 Object Invocation

Operation invocation on objects is the mechanism in Alpha that provides the uniform interface between objects and between objects and the kernel. It serves as both system call and communication interface. Operations on other client objects on the same node, on system and kernel supported objects, and on objects on other nodes are all accomplished by invocation.

The general form of the invocation is a kernel call with the capability of the target object, the operation, and a list of input parameters. The invocation return has only a list of the output parameters. The capability of the invoking object and the return entry point are both implicit in the return invocation call.

Copies of the input parameters are passed to the invoked object. The kernel adds two arguments to the front of the list. These are capabilities for the current object and the current thread. These parameters, generated by the kernel from the invocation frame, are the *Object Self* and *Thread Self* capabilities.

The first output parameter is the invocation return value and is set *only* by the kernel just before the invocation returns to the invoking object. The second output parameter is (by convention on kernel supported objects) the operation return value (*status*). The operation return values are specific to particular operations on particular types of objects, but the invocation return values indicate the result of any invocation. If the invocation return value is SUCCESS then as far as the kernel is concerned, the invocation was successful. Otherwise the value indicates some form of error occurred either before the invocation reached the invoked object or after the invoked object performed an invocation return. The following list enumerates the possible invocation return values:

SUCCESS — Invocation was completed successfully.

BAD_TARGET_CAPABILITY — The value used as the target of an operation did not represent a capability.

BAD_INPUT_CAPABILITY — A value passed as a capability input parameter was not a capability.

BAD_OUTPUT_CAPABILITY — A value returned by the target object as a capability output parameter was not a capability.

CAPABILITY_LIMIT — A capability passed as a parameter in the invocation or return exceeded the maximum number of capabilities that can be accommodated on the node.

OBJECT_NOT_FOUND — The capability used as the target of the invocation was no longer bound to an active object.

BAD_OPERATION — The specified operation was not accessible using the target capability.

THREAD_BREAK — The current thread was fragmented at some point due to some event such as object deletion or node failure.

Invocation parameters are always passed by value. It is up to individual implementations of the Alpha kernel how storage for the parameter passing is managed, but invoca-

tions must be able to pass an arbitrary number of arbitrary sized parameters. Invocations can have pointers as parameters only in the restricted case where the pointer was allocated from the *thread heap*. The thread heap is a dynamic storage allocation mechanism that allows memory to be allocated that is associated with the thread. This memory moves with the thread and data stored in it does not have to be copied across invocations.

4 Capabilities

In the Alpha kernel, a capability is a protected naming and access mechanism for objects. When any kernel, system or client object is created, it is assigned a unique identifier which is then bound to a capability. This capability is visible only within the kernel: client objects refer to the capability via a kernel-supplied *capability index*. The capability mechanism, therefore, gives unforgeable access to the object from any object that has a capability index referring to it.

Capabilities may be acquired by objects from *in parameters* of invocations and *out parameters* of invocation returns. When a capability arrives at an object, it is added to the object's kernel-maintained *capability list*; the kernel supplies a capability index by which the object can refer to the capability. The entry made in the capability list is a copy of the capability that was supplied by the source of the invocation (or invocation return). The resulting capability index bears no relationship with the capability index used by the caller.

Capabilities carried into an object as invocation parameters are *transient*, and can only be used by another thread in the object if the capability has explicitly been made *permanent* within the object by executing the *SaveCapa* operation on the object. Saving generates a new capability index for the capability, and this index can then be stored into some instance variable for later use by another thread. The original capability index, representing the transient capability, is still available to the object as long as the thread that brought it is still executing in the object.

When a thread returns from an object, all transient capabilities brought by the thread to the object are removed from the object. That is, the capability indices that referred to the transient capabilities will no longer be valid within the object. Capabilities that are returned as output parameters appear as transient capabilities in the invoking object.

A capability can be *anchored* to a thread. An anchored capability is always transient. It can not be made permanent via the *SaveCapa* operation. Therefore, it can only be used by the thread that brought it as a parameter.

Capabilities can have restrictions placed on them as they are passed as invocation parameters or used as the target of an invocation. The invoked object will receive a copy of the capability with restrictions applied. Restrictions consist of a list of operations that may not be performed on the object that is the target of the capability. Once a restriction is placed on a capability, it cannot be removed; all copies of a restricted capability retain the restrictions.

5 Exception Handling

While a thread is executing within an Alpha application, asynchronous events, or *exceptions*, may occur that prevent processing from proceeding in the normal, sequential manner. For example, a deadline may expire or the thread may encounter a machine fault. Alpha permits an application to specify an *exception scope*, which contains code that is notified of exception conditions that arise when the thread is executing within the scope. This code may take action designed to handle or recover from the exception.

An exception scope is delimited by the invocation of matching pairs of operations: one that begins the exception scope, and one that ends the exception scope. The scope includes all of the thread's activity between the invocation of the operations that define it. Although exception scopes are denoted by operations on objects, they are managed separately for each thread. Therefore, one thread's invocation of a delimiting operation cannot affect any other thread. The specific types of exception scopes are:

General exception: Delimited by *BeginScope* and *EndScope* operations on the Thread object.

Time constraint: Delimited by *Deadline* and *Mark* operations on the Thread object.

Transaction: Delimited by *BeginTransaction* and *EndTransaction* operations on the Transaction Manager object.

Exception scopes may be *nested* within other exception scopes. Exception conditions are reported starting at the most deeply nested exception scope and continuing outward for as long as the exception remains active. The conditions under which an exception ceases to be active are determined by the type of exception.

All types of exception scopes are governed by strict nesting. That is, a nested exception scope of any type must be a subset of the surrounding scope, whatever its type. If the application attempts to execute any scope-delimiting operation that violates the nesting rules, an error condition will be reported and the operation will be refused.

If an exception condition is active within an exception scope (and remains active after any nested exception scopes have reported the exception and completed), it is reported to the exception scope in the following way. The thread stops execution at its current point of control, and resumes execution at the return point of the operation that began the exception scope. This has the effect of appearing to return from that operation a second time (the first was when the exception scope was established). However, this time the operation status value returned will indicate that an exception has occurred. Other than this status value, memory contents are left as modified by the thread before the exception condition arose. The processor state (register contents, etc.) is restored in an implementation-defined way.

If the operation beginning the exception scope has reported an exception condition (by returning a second time), the thread is said to be in the *reporting mode* with respect to the exception scope. The thread remains in reporting mode for the scope until it executes the operation that ends the exception scope. If the thread completes the exception scope while in reporting mode, and the exception condition is still active, its execution will resume at the next enclosing exception scope (in reporting mode with respect to that scope).

A thread that is executing in reporting mode with respect to any exception scope may not execute any operation that attempts to begin another exception scope. This means, for example, that neither transactions nor time constraints may be started within a scope in reporting mode. If the application attempts to execute any scope-delimiting operation that violates this rule, an error condition will be reported and the operation will be refused.

The following are the exception conditions reported by Alpha and the conditions under which they remain active:

Deadline expiration: This exception occurs whenever the value of continuing computation under a time constraint becomes zero. The exception remains active until the thread is no longer within the scope of the time constraint.

Machine check: This exception occurs whenever the processor detects an error, such as a protection violation, illegal instruction, etc. It remains active until a general exception scope is executed, or, if there is no general exception scope, until the thread is deleted.

Thread abort: This exception occurs as the result of an explicit *Abort* operation on the thread or as the result of a thread segment becoming an orphan due to node failure or object deletion. It remains active until the thread or thread segment is deleted.

Transaction exception: This exception occurs as the result of an explicit *AbortTransaction* operation on the Transaction Manager object within the scope of a transaction or as the result of the inaccessibility of any object involved in the transaction (due to node failure or object deletion). It remains active until the thread is no longer within the scope of the transaction.

Whenever a thread is executing in reporting mode within an exception scope, the reporting of other exception conditions, other than machine checks, is deferred until the completion of the exception scope. Machine check exceptions cannot be deferred because they represent an error within the exception scope itself, and, therefore, preclude the continuation of processing. For deadline expiration exceptions, the scheduling parameters controlling the thread's execution immediately revert to those in effect before the expired time constraint, even though the reporting of the exception may be deferred.

6 Standard Operations

The kernel defines a set of operations that is inherited by all system and client objects. These operations define the services the kernel provides for manipulating objects. (The exception to this rule is the *Create* operation which is defined on the Schema object that defines the object type.) Standard Operations may be redefined by the application developer by using the name of a Standard Operation for one of the user defined operations on an object. When a redefined Standard Operation is invoked, the kernel uses the new operation instead of the Standard Operation.

Each Standard Operation has two names—the regular name and the regular name plus the “Kernel_” prefix. An example is *Update* and *Kernel_Update*. The kernel always uses the short name when it invokes a Standard Operation. There is a special relationship between each pair of Standard Operation names. If the regular operation name is invoked and it has not been redefined by the user then the long operation name is invoked.

The purpose of this indirection is to provide a simple mechanism for the user to redefine the semantics of Standard Operations without losing any functionality. As an example, the user can now redefine Info so that additional permissions are checked before the invocation is carried out.

The Standard Operations are:

Update:

target:	Object	
opcode:	Update	
in parameters:	address	Location
	length	Size
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is used to write a portion of the object to secondary storage. The *address* parameter indicates the beginning of the buffer and the *length* parameter indicates its length in bytes.

Operation Status: The value NOT_PERMANENT is returned in the *status* parameter if a new object could not be updated because it did not have the permanence attribute. The value BAD_ADDRESS is returned if any of the addresses specified by the *address*, *length* parameters does not belong to the object. The value NO_SPACE is returned if for some reason there is not enough space on the secondary storage system to perform the update. Otherwise the value SUCCESS is returned.

Restore:

target:	Object	
opcode:	Restore	
in parameters:	address	Location
	length	Size
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is used to read a portion of the object from secondary storage. The *address* parameter indicates the beginning of the buffer and the *length* parameter indicates its length in bytes.

Operation Status: The value NOT_PERMANENT is returned in the *status* parameter if a new object could not be updated because it did not have the permanence attribute. The value BAD_ADDRESS is returned if any of the addresses specified by the *address*, *length* parameters does not belong to the object. Otherwise the value SUCCESS is returned.

Info:

target:	Object	
opcode:	Info	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus
	info	InfoBlock

Description: This operation is used to query the attributes and state of an existing object. The *info* structure contains the following fields:

max_size	— max size of the object
minimum_size	— min size of the object
current_size	— current size of the object
top_of_heap	— this address represents the top of the object heap
active_threads	— number of active threads within the object
thread_segs	— number of thread segments within the object
location	— network address
replication	— replicated ? (true/false), number of elements in set, replication strategy, etc.
atomic	— atomic ? (true/false) transaction strategy
transactions	— number of active transactions affecting the object
permanent	— the object is permanent (true/false)

Note that the information returned may be stale. Other threads executing within the system might change the state of the object after its state is examined but before the state information can be reported.

Operation Status: The value SUCCESS is always returned.

Migrate:

target:	Object
opcode:	Migrate
in parameters:	address NodeLocation
return value:	result InvocationStatus
out parameters:	status OperationStatus

Description: This operation is used by the client to move objects among the nodes in the system. The *location* parameter identifies the node to which the object should be moved. The operation returns only when the requested movement has been completed.

Operation Status: The value BAD_LOCATION is returned in the parameter *status* if the *location* parameter is invalid. A location is invalid if the node referred to by that address cannot be reached by the communications subsystem. Otherwise the value SUCCESS is returned.

Lock:

target:	Object
opcode:	Lock
in parameters:	address Location
	length Size
	mode LockMode
return value:	result InvocationStatus
out parameters:	status OperationStatus

Description: The purpose of the *Lock*, *C_lock*, and *Unlock* operations is to control the concurrent access of threads to individual data regions within an object. These locks are of advisory nature. Locks also figure prominently in Alpha's support of atomic transactions. The logs that track the changes made to an object during a transaction are associated with the locks acquired during the transaction. These features are transparent to the users of locks when no transaction is in progress. The *Lock* operation attempts to add a restriction to a region of the object. The region is specified by the parameters *address* and *length*. The *mode* parameter specifies the restriction or the type of lock requested. There may be existing restrictions on part or all of the region. If this new restriction is consistent with any existing locks then the lock is granted and the operation returns. If the restriction conflicts with an existing lock then the invoking thread is suspended until there are no conflicting locks.

Modes:

- No Readers allowed
- Exclusive Read
- Multiple Read
- no Read restrictions

No Writers allowed
 Exclusive Write
 Multiple Write
 no Write restrictions

The various mode values, one read and one write, can be OR'ed together to form the *mode* parameter.

Operation Status: The kernel will detect when a thread attempts to acquire two locks within an object in an order that would cause deadlock. This could occur if the thread attempted to acquire a Multiple Read lock after acquiring an Exclusive Read lock. If a deadlock is detected the operation fails and the value DEADLOCK is returned in the parameter *status*. The value BAD_LENGTH is returned if the length is invalid. The value BAD_MODE is returned if the *mode* parameter is invalid. Otherwise the value SUCCESS is returned.

C_Lock:

target:	Object	
opcode:	C_Lock	
in parameters:	address	Location
	length	Size
	mode	LockMode
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is similar to the normal *Lock* operation, but never results in the invoking thread being blocked. Instead, a success or failure indication is returned in the locked parameter. The parameter *mode* represents the restriction that the invoking thread wishes to add on the region. The *address* and *length* (in bytes) parameters describe the region to be locked.

Operation Status: The lock compatibility table is consulted, and if no conflicting locks are held the lock is granted and SUCCESS is returned in the *status* parameter. If conflicting locks exist then FAILURE is returned in the *status* parameter. The value BAD_LENGTH is returned in the parameter *status* if the length is invalid. The value BAD_MODE is returned in the parameter *status* if the *mode* parameter is invalid. Otherwise the value SUCCESS is returned.

Unlock:

target:	Object	
opcode:	Unlock	
in parameters:	address	Location
	length	Size
	mode	LockMode
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation attempts to remove a restriction or lock applied to a region of an object. The access restrictions are reduced to the level specified by the *mode* parameter. The region affected is specified the *address* and *length* parameters. An *Unlock* reduces the restrictions owned by the invoking thread. *Unlock* can not be used to raise the level of restrictions by specifying a higher level lock.

Operation Status: The value NO_LOCK is returned in the *status* parameter if the thread does not already possess a lock on the specified range of bytes. The value BAD_LENGTH is returned in the parameter *status* if the length is invalid. The value BAD_MODE is returned in the parameter *status* if the *mode* parameter is invalid. Otherwise the value SUCCESS is returned.

PreCommit:

target:	Object	
opcode:	PreCommit	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is used in the first phase of a two-phase commit function. It performs an UPDATE-like operation that returns once the object's image has written to secondary storage in such a fashion that the commit operation can be completed even if node crashes should occur.

Operation Status: The value NO_TRANSACTION is returned in the parameter *status* if the thread is not involved in a transaction. The value NOT_ATOMIC is returned if the object does not have the ATOMIC attribute. Otherwise the value SUCCESS is returned.

Commit:

target:	Object	
opcode:	Commit	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is used in the second phase of a two-phase commit function. It performs an UPDATE-like function that returns only when the object's image has been completely written to secondary storage, releases any semaphores or locks held by the transaction within the object, and returns a success indication.

Operation Status: The value NO_TRANSACTION is returned in the parameter *status* if the thread is not involved in a transaction. The value NOT_ATOMIC is returned if the object did not possess the ATOMIC attribute. Otherwise the value SUCCESS is returned.

Abort:

target:	Object	
opcode:	Abort	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation is used to abort a transaction, either before or after the pre-commit phase. This operation restores all locked data items to previous state, undoes the pre-commit (i.e., restores the secondary storage to its previous state), releases any semaphores or locks held by the transaction within the object and returns a success indication.

Operation Status: The value NO_TRANSACTION is returned if there was no active transaction. The value NOT_ATOMIC is returned if the object did not possess the ATOMIC attribute. Otherwise the value SUCCESS is returned.

AllocateHeap:

target:	Object	
opcode:	AllocateHeap	
in parameters:	amount	Size
return value:	result	InvocationStatus
out parameters:	status	OperationStatus
	address	Pointer

Description: Allocate memory in the object heap. At least *amount* bytes of memory is allocated and placed into the object heap; the address of the newly allocated area is returned in the *address* parameter. The Alpha kernel may round up the size of this area to a convenient memory boundary, such as an even multiple of the system's page size. The area pointed to by *address* is guaranteed to meet the alignment restrictions of any fundamental data type of the processor.

Operation Status: *Status* is SUCCESS if the heap could be grown by the specified amount, or NO_HEAP if the space could not be made available.

DeallocateHeap:

target:	Object	
opcode:	DeallocateHeap	
in parameters:	address	Pointer
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Deallocate memory from the object heap. The object heap memory located at *address* is removed from the object heap and returned to the system; it is no longer accessible. The *address* parameter must contain a value previously returned by an *AllocateHeap* operation.

Operation Status: *Status* is SUCCESS if the memory is deallocated from the heap, or NO_HEAP if the *address* parameter does not refer to the beginning of an allocated section of the heap.

SaveCapa:

target:	Object
opcode:	SaveCapa
in parameters:	capa Capability
return value:	result InvocationStatus
out parameters:	status OperationStatus
	savedCapa Capability

Description: Makes a permanent copy of a capability in the object's capability list. This allows the capability to be retained by the object for use by other threads.

Operation Status: The value NOT_SAVEABLE is returned if the capability could not be saved by the object, otherwise the value SUCCESS is returned.

ReleaseCapa:

target:	Object
opcode:	ReleaseCapa
in parameters:	capaindex CapabilityIndex
return value:	result InvocationStatus
out parameters:	status OperationStatus

Description: Delete the capability represented by the capability index *capaindex* from the object's capability list. (Note that this parameter is the capability index only; not the capability itself.)

Operation Status: Returns the value BAD_CAPABILITY if *capaindex* was not the index of any valid capability in the object's capability list. Otherwise, returns SUCCESS.

Delete:

target:	Object
opcode:	Delete
in parameters:	threads Boolean
return value:	result InvocationStatus
out parameters:	status OperationStatus

Description: This operation deletes the object. The *threads* parameter indicates whether the object should be deleted even if threads are currently active within it. If *threads* is TRUE then the object is deleted even if there are active threads within it. When the object is deleted the kernel treats any threads within it as if a node failure had been detected.

Operation Status: The value OBJECT_ACTIVE is returned if object has active threads within it and the *threads* parameter was FALSE. Otherwise the value SUCCESS is returned.

7 Kernel Services

The Alpha operating system is composed of three layers: the kernel proper, a set of system objects, and possibly a set of client-level objects. All of the abstractions discussed in this chapter are part of the kernel proper and form the basis upon which all other services in Alpha are built. Alpha kernel services are described in terms of kernel objects and the operations performed on them.

Object Definition

OBJECT:

Schema

Description: In an Alpha system, each object is an instance of a particular object type. An object type is represented in the kernel by a Schema object. Each Schema object is a template containing a description of the structure, operations, and initial state of the data of the object type; it is used by the kernel to create instances of the type of object described by the template. An object instance is created by invoking a *CreateSystem* or *Create* operation on the appropriate Schema object and providing the attributes, such as permanence, that the object instance should be created with.

Each instance of the kernel comes with some predefined Schema objects. For example, the kernel contains Manager objects for both the Thread and Semaphore types objects. These manager objects act as the Schema objects for these object types. There is also a predefined Schema Manager object that can create instances of Schema objects to define other object types.

The user can create new Schema objects to describe new types of client objects. A Schema object is obtained by invoking the *Create* or *CreateUnrestricted* operation on the Schema Manager object. This newly created Schema is blank and cannot be used to instance objects until it is initialized and validated. Any operations and data the programmer wishes to encapsulate within the new type must be explicitly added to the Schema. Once this is done the Schema must be explicitly validated. After a Schema object is validated, it represents a new type; objects of that type can be instanced by invoking the *CreateSystem* or *Create* operation on it. Once a Schema has been validated it cannot be changed.

The information used by the Schema to describe an object is outlined in the following structure.

SCHEMA OBJECT

```
{
  ATTRIBUTES
  {
    Attribute (Validated = TRUE)
    Attribute (Version = 2)
    .
    .
  }
  EXTENTS
  {
    Extent (Type = TEXT, ...)
    Extent (Type = DATA, ...)
    Extent (Type = TEXT, ...)
    .
    .
  }
  OPERATIONS
  {
    Operation (Op_Name = PUSH, ...)
    Operation (Op_Name = POP, ...)
    .
    .
  }
  RELOCATION ENTRIES
  {
    Relocation_entry (Offset = 10, Extent = Text1, Target_Extent = ...)
    .
    .
  }
  SYMBOLS
  {
    Symbol (Name = "stack", Offset = ...)
    .
    .
  }
}
```

The Version attribute of the Schema contains the version number of the Alpha kernel under which this Schema was created. The ability to import and export objects is important for backing up objects and for program development. Therefore, it is assumed that objects and their Schemas may be moved from system to system, Alpha version to Alpha version. The version number field in the Schema structure gives the kernel a way to provide compatibility (when possible) with other Alpha versions, while allowing the Schema structure to evolve as new requirements are identified.

The Validated attribute is a boolean value that indicates whether the *Validate* operation has been successfully invoked on this Schema.

The list of extent structures describes the collection of extents that comprise the object. Each object must have at least one of each type of extent. An extent is a contiguous section of memory whose characteristics are determined by its type. A Size field gives the length in bytes of the extent. An Address field gives the starting Virtual Address (if it is fixed). Extents may sparsely populate an address space but they may not overlap one another within an address space. A Data field gives the initial values to which the memory comprising the extent should be set. Usually this is only relevant to the TEXT and DATA type of extents.

The list of operation structures gives the name of each operation, its entry point or starting address (relative to the beginning of a TEXT extent), the number of IN and OUT Capabilities, and the length of the IN and OUT parameter block. An operation name is a bit string that is OP_NAME_LENGTH bits long.

The list of relocation entries provides information which is used whenever the virtual addresses of the extents are not fixed. In that case, the kernel will assign the extent to virtual addresses at object instance creation time, and use the relocation entries to establish inter-extent references.

The Symbol Table is a list of symbol entries. Each entry has a symbol name and an address (relative to the beginning of an extent) for the symbol. Symbol names are zero terminated byte strings.

Object Operations

Create:

target:	Schema	
opcode:	Create	
in parameters:	max_size	Size
	permanent	Boolean
	atomic	Boolean
	replicate	Boolean
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus
	new_object	Capability

Description: This operation creates a new object instance of the type defined by the target Schema object. The new object will be a client object and will exist in an address space separate from the Alpha kernel and from all other object instances. The *new_object* parameter is the capability of the newly created object.

The *max_size* parameter indicates the maximum size of the object in bytes. If an attempt is made to grow the size of the object past the maximum size the request will fail. If the object is created with a *max_size* of zero then the kernel default *max_size* is used.

If *permanent* is TRUE then the object will have special non-volatile storage allocated to it. The *Update* operation can be used to write the current state of the object to the non-volatile storage. In the case of a node failure, this non-volatile storage is used to restore the object.

If *atomic* is TRUE the object can support atomic transactions. Atomic transactions are a classic technique used in databases to provide consistent, concurrent access to a database. There are special standard operations, including *PreCommit* and *Commit* that are only meaningful for atomic object. Note that an atomic object needs both the permanence and atomic attributes in order for the atomic object to persist despite node failures.

If the parameter *replicated* is FALSE then the remaining parameters are ignored. If it is true, then the object will be replicated.

Operation Status: The value BADVERSION is returned if the target Schema object has a version number that is incompatible with the version of Alpha that is running. The value RESOURCE_LIMIT is returned if the object could not be created because an Alpha resource limit is exceeded (e.g., memory, control structures, etc.). Otherwise the value SUCCESS is returned.

CreateSystem:

target:	Schema	
opcode:	CreateSystem	
in parameters:	max_size	Size
	permanent	Boolean
	atomic	Boolean
	replicate	Boolean
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus
	new_object	Capability

Description: This operation creates a new object instance of the type defined by the target Schema object. The new object is a system object, and shares the address space of the Alpha kernel. The *new_object* parameter is the capability of the newly created object.

The *max_size* parameter indicates the maximum size of the object in bytes. If an attempt is made to grow the size of the object past the maximum size the request will fail. If the object is created with a *max_size* of zero then the kernel default *max_size* is used.

If *permanent* is TRUE then the object will have special non-volatile storage allocated to it. The *Update* operation can be used to write the current state of the object to the non-volatile storage. In the case of a node failure, this non-volatile storage is used to restore the object.

If *atomic* is TRUE the object can support atomic transactions. Atomic transactions are a classic technique used in databases to provide consistent, concurrent access to a database. There are special standard operations, including *PreCommit* and *Commit* that are only meaningful for atomic object. Note that an atomic object needs both the permanence and atomic attributes in order for the atomic object to persist despite node failures.

If the parameter *replicated* is FALSE then the remaining parameters are ignored. If it is true, then the object will be replicated.

Operation Status: The value BADVERSION is returned if the target Schema object has a version number that is incompatible with the version of Alpha that is running. The value BADEXTENTADDRESS is returned if and the Schema contains extents at fixed addresses. The value RESOURCE_LIMIT is returned if the object could not be created because an Alpha resource limit is exceeded (e.g., memory, control structures, etc.). Otherwise the value SUCCESS is returned.

AddExtent:

target:	Schema	
opcode:	AddExtent	
in parameters:	type	ExtentType
	length	Size
	starting_address	Address
	data	Pointer
	d_length	Size
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus
	extent	ExtentID

Description: This operation adds an extent to the Schema object. Objects encapsulate data and the operations that manipulate the data. Objects are composed of extents and entry points. There are three *empty* of extents: TEXT, DATA, and HEAP. Operations are TEXT extents. Data with known initial values are DATA extents. HEAP extents provide for the dynamic allocation and deallocation of memory.

Extents are contiguous sections of memory. The *starting_address* is the starting virtual address for this extent. The *length* plus the *starting_address* specify the section of virtual memory occupied by the extent. Extents may sparsely populate an address space but they may not overlap. If *starting_address* is NULL, then the kernel is free to use any virtual address that it chooses for the extent when an object instance is created. All extents of system objects must be specified with a NULL *starting_address*.

The *data* parameter is a pointer into the thread heap to data that describes the initial state of the extent; the data field extends for *d_length* bytes. If a TEXT extent is specified, then the field contains the code for the operations associated with that extent. If it is a DATA extent, the field holds the initial values for the data associated with that extent. The thread heap memory pointed to by *data* will be deallocated during this operation if the operation returns SUCCESS.

If *d_length* is less than *length*, the unspecified (trailing) portion of the extent will be initialized to zero; if *d_length* is greater than *length*, the extra is ignored.

The return parameter *extent* contains an extent identifier which may be used in subsequent *AddOperation*, *AddRelocation*, and *AddSymbol* operations.

Operation Status: The value BADEXTENTTYPE is returned if the *type* parameter has an illegal value. The value BADEXTENTLENGTH is returned if the *length* parameter has an illegal value. If the length is greater than the maximum length for extents or less than the minimum length for extents then the length is illegal. The value BADEXTENTADDRESS is returned if the *starting_address* parameter has an illegal value or if the new extent overlaps an existing extent. The value ALREADYVALID is returned if the operation cannot be performed because the Schema has previously been validated. Otherwise the value SUCCESS is returned.

Description: This operation adds a relocation entry to the Schema object. It is used when some or all of the extents of the object are specified with a NULL *starting_address*, in which case the virtual address of these extents will be determined at object instance time. During object instantiation, the actual virtual address of the target extent *t_extent* is added to the data at the address *address*, relative to the beginning of extent *extent*. For example, if the data at that address is zero, the result will be the base address of the target extent. If non-zero data is at that address, the result will be an address offset by that amount into the target extent.

The addition performed is ADDRESS_SIZE bytes wide, where ADDRESS_SIZE is the largest natural address size for the processor.

Operation Status: The value BADEXTENTID is returned if *extent* or *t_extent* is an invalid extent identifier. The value BADEXTENTADDRESS is returned if *address* is not a valid offset into the extent *extent*. The value ALREADYVALID is returned if the operation cannot be performed because the Schema has previously been validated. Otherwise the value SUCCESS is returned.

AddSymbol:

target:	Schema	
opcode:	AddSymbol	
in parameters:	name	Symbol
	address	Offset
	extent	ExtentID
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation adds a Symbol to the Symbol table. The *name* field is a zero terminated byte string which is up to SYMBOL_LENGTH bytes in length. This is the name of the symbol. The virtual address, relative to the beginning of extent identified by *extent*, associated with the symbol is *address*.

Operation Status: The symbol table supports MAX_NUM_SYMBOLS. After the table is full, any additional attempts to add symbols will fail and the value SYMBOLTABLEFULL will be returned. The value BADEXTENTID is returned if *extent* is an invalid extent identifier. The value BADEXTENTADDRESS is returned if *address* is not a valid offset into the extent *extent*. The value ALREADYVALID is returned if the operation cannot be performed because the Schema has previously been validated. Otherwise the value SUCCESS is returned.

Validate:

target:	Schema	
opcode:	Validate	
in parameters:	none	
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation validates the current state of the Schema. If successful, additional modifications to the Schema database will be disallowed.

Operation Status: The value BADSCHEMA is returned if the Schema object has an invalid format. The value ALREADYVALID is returned if the Schema object is already VALID. The value NOOPERATIONS is returned if the Schema object has no operations specified. Otherwise the value SUCCESS is returned.

Delete:

target:	Schema	
opcode:	Delete	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation deletes the Schema object.

Operation Status: The value SUCCESS is always returned.

Object Definition

OBJECT:

Schema Manager

Description: The Schema Manager object controls the creation of Schema objects, which are used to dynamically define new object types. A Schema object is obtained by invoking the *Create* or *CreateClient* operation on the Schema Manager object.

Object Operations

Create:

target:	Schema Manager	
opcode:	Create	
in parameters:	permanent	Boolean
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus
	new_schema	Capability

Description: This operation creates a new Schema object of the type defined by the target Schema object. The *new_schema* parameter is the capability of the newly created Schema object.

If *permanent* is TRUE then the Schema object will have special non-volatile storage allocated to it. In the case of a node failure, this non-volatile storage is used to restore the Schema object.

The Schema object created by this operation will have its *CreateSystem* operation restricted. Therefore, only the *Create* operation may be used on the Schema to create object instances, and only client object instances may be created.

Operation Status: The value SUCCESS is always returned.

CreateUnrestricted:

target:	Schema Manager	
opcode:	CreateUnrestricted	
in parameters:	permanent	Boolean
return value:	invoke_status	InvocationStatus
out parameters:	status	OperationStatus
	new_schema	Capability

Description: This operation creates a new Schema object of the type defined by the target Schema object. The *new_schema* parameter is the capability of the newly created Schema object.

If *permanent* is TRUE then the Schema object will have special non-volatile storage allocated to it. In the case of a node failure, this non-volatile storage is used to restore the Schema object.

The Schema object created by this operation has no restrictions on the object creation operations. Therefore, both the *Create* and *Create.System* operations may be used on the Schema to produce client and system object instances, respectively.

Operation Status: The value SUCCESS is always returned.

Object Definition

OBJECT:

Semaphore

Description: The type semaphore represents a kernel administered synchronization mechanism utilizing a standard counting semaphore algorithm. The critical feature of semaphores is that operations on them are performed atomically and are, therefore, mutually exclusive. Semaphores are useful as concurrency control mechanisms between threads within an object or in different objects.

Object Operations

P:

target:	Semaphore	
opcode:	P	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Perform an atomic P operation on the counting semaphore defined as: $\text{count}(\text{Semaphore}) -= 1$; if $\text{count}(\text{Semaphore}) < 0$ then the invoking thread is blocked and a scheduling event is recognized.

Operation Status: The value DELETED is returned in the *status* parameter if the thread was unblocked because the semaphore was deleted. The value V_ALL is returned if the thread was unblocked due to a V_ALL operation. Otherwise, the value SUCCESS is returned.

V:

target:	Semaphore	
opcode:	V	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Perform an atomic V operation on the counting semaphore defined as: $\text{count}(\text{Semaphore}) += 1$; if $\text{count}(\text{Semaphore}) \leq 0$ then select a thread, unblock it, and recognize a scheduling event.

Operation Status: The value SUCCESS is always returned in the out parameter *status*.

V_All:

target:	Semaphore	
opcode:	V_All	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: If there are any threads blocked on the Semaphore, i.e., the $\text{count}(\text{Semaphore}) < 0$, they are unblocked and the count is set to 0. If there are no threads blocked on the Semaphore the *V_All* operation has no effect.

Operation Status: The value SUCCESS is always returned in the out parameter *status*.

Conditional_P:

target:	Semaphore	
opcode:	Conditional_P	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Perform an atomic conditional P operation on the counting semaphore defined as: if $\text{count}(\text{Semaphore}) > 0$ then $\text{count}(\text{Semaphore}) -= 1$, else the invoking thread returns with an error value.

Operation Status: The value OWNED is returned in *status* if the semaphore could not be obtained. The value SUCCESS is returned if the operation succeeded.

Delete:

target:	Semaphore	
opcode:	Delete	
in parameters:	threaded	Boolean
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Delete the semaphore object. The *threaded* parameter indicates whether the Semaphore object can be deleted if there are threads blocked on it. If *threaded* is TRUE then the object can be deleted and the threads will be unblocked. If *threaded* is FALSE then the operation fails.

Operation Status: The value BLOCKED is returned in the *status* parameter if the object to be deleted had blocked threads but the *threaded* parameter was FALSE. The value SUCCESS is returned if the operation succeeded.

Object Definition

OBJECT:

Semaphore Manager

Description: The kernel provides the Semaphore Manager object in order to permit other objects to create and delete individual instances of semaphore objects. The capabilities to Semaphore objects are not passable. This means that operations on Semaphore objects can only be accessed by the objects that created them.

Object Operations

Create:

target:	Semaphore Manager	
opcode:	Create	
in parameters:	count	Size
return value:	result	InvocationStatus
out parameters:	status	OperationStatus
	sema	Capability

Description: Instantiate a kernel object of type semaphore and return a capability to it in *sema*. This operation takes as an input argument the initial *count* value for the semaphore.

Operation Status: The value ERROR is returned in the *status* parameter if a new semaphore object could not be created, otherwise the value SUCCESS is returned.

Object Definition

OBJECT:

Thread

Description: Threads are the embodiment of the Alpha control abstraction. Threads are represented as kernel objects and can be manipulated by invoking operations on them. It is not possible for the client to specify custom versions of these operations.

Description: Resumes the execution of a stopped thread, and has no effect if the thread is already running.

Operation Status: *Status* is SUCCESS if the a new thread is resumed by this invocation and ALREADY_RUNNING if the thread was already running when the invocation occurred.

BeginScope:

target:	Thread	
opcode:	BeginScope	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Establish a general exception scope. The exception scope defined by this operation remains active until a subsequent *EndScope* operation is invoked at the same nesting level in the current operation. If any exception condition occurs while this scope is active, and it remains active after completing any nested exception scope, *BeginScope* will return (a second time) with a value indicating the type of exception. If the thread returns from the operation that performed the *BeginScope* operation without having invoked a corresponding *EndScope* operation, an *EndScope* is automatically supplied by the system.

Operation Status: The value NO_SCOPE is returned in the *status* parameter if an exception scope could not be created. The status NESTING_VIOLATION is returned if the thread is currently executing in the reporting mode of an exception scope. Otherwise, the status SUCCESS is returned. If, while the exception scope established by this operation is active, an exception occurs, the operation returns with the value THREAD_ABORT, DEADLINE_EXPIRE, TRANSACTION_ABORT, or MACHINE_CHECK if the exception was due to a thread abort, the expiration of a deadline, a transaction abort, or a hardware fault, respectively. An implementation of Alpha may define additional values of status to indicate exceptions due to specific machine error conditions.

EndScope:

target:	Thread	
opcode:	EndScope	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: End an exception scope established by a previous *BeginScope* operation. If that *BeginScope* operation has returned the value `THREAD_ABORT`, `DEADLINE_EXPIRE`, or `TRANSACTION_ABORT` (i.e. an exception has been reported and is still active), then the thread will continue execution at the next enclosing exception scope. Otherwise, *EndScope* will return normally.

Operation Status: The value `NO_SCOPE` is returned in the *status* parameter if there is no currently active *BeginScope* operation in the current object. The status `NESTING_VIOLATION` is returned if any nested exception scope has not been closed. Otherwise, the status `SUCCESS` is returned.

Deadline:

target:	Thread	
opcode:	Deadline	
in parameters:	spb	SchedParamBlock
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Establish a time-critical region for the current thread and begin an associated exception scope. This operation takes *spb* as an input argument, which is a scheduling parameter block that defines the time constraint under which the thread will execute until a subsequent *Mark* operation is executed. The interpretation of the information in the scheduling parameter block is determined by the scheduling policy module installed in the Alpha system. Time constraints may be nested, in which case the thread will execute under the most restrictive time constraint specified. If the thread does not execute a *Mark* operation at the same nesting level in the current object before the value of continuing computation (as specified by the scheduling parameter block) becomes zero, a deadline expiration exception will be generated on the thread. This exception remains active until a corresponding *Mark* operation is executed. The exception scope defined by this operation remains active until a *Mark* operation is invoked at the same nesting level in the current object. If any exception occurs while this scope is active, and it remains active after completing any nested exception scope, *Deadline* will return (a second time) with a value indicating the type of the exception. If the thread returns from the operation that performed the *Deadline* operation without having invoked a corresponding *Mark* operation, a *Mark* is automatically supplied by the system.

Operation Status: The value `BAD_SCHEDULING_INFO` is returned in the *status* parameter if a time constraint could not be created. The status `NESTING_VIOLATION` is returned if the thread is currently executing in the reporting mode of an exception scope. Otherwise, the status `SUCCESS` is returned. If, while the exception scope established by this operation is active, an exception occurs, the operation returns with the value `THREAD_ABORT`, `DEADLINE_EXPIRE`, `TRANSACTION_ABORT`, or

MACHINE_CHECK if the exception was due to a thread abort, the expiration of a deadline, a transaction abort, or a hardware fault, respectively. An implementation of Alpha may define additional values of status to indicate exceptions due to specific machine error conditions.

Mark:

target:	Thread	
opcode:	Mark	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: End a time-critical region and its associated exception scope as established by a previous *Deadline* operation. If that *Deadline* operation has returned a value reporting an exception, and the exception is still active after the *Mark* operation, then the thread will continue execution at the next enclosing exception scope. Otherwise, *Mark* will return normally.

Operation Status: The value **NO_DEADLINE** is returned in the *status* parameter if there is no currently active *Deadline* operation in the current object. The value **NESTING_VIOLATION** is returned if any nested exception scope has not been closed. Otherwise, the value **SUCCESS** is returned.

AllocateHeap:

target:	Thread	
opcode:	AllocateHeap	
in parameters:	amount	Size
return value:	result	InvocationStatus
out parameters:	status	OperationStatus
	address	Pointer

Description: Allocate memory in the current thread's thread heap. At least *amount* bytes of memory is allocated and placed into the thread heap; the address of the newly allocated area is returned in the *address* parameter. The Alpha kernel may round up the size of this area to a convenient memory boundary, such as an even multiple of the system's page size. The area pointed to by *address* is guaranteed to meet the alignment restrictions of any fundamental data type of the processor. Memory allocated in the thread heap is visible only to the current thread, and migrates with the thread across object invocations. The addresses of data items allocated in the thread heap can be passed as invocation parameters to other

objects, which will then be able to access the data directly. This can be used by an application to pass large amounts of data between objects in an efficient manner.

Operation Status: *Status* is SUCCESS if the heap could be grown by the specified amount, or NO_HEAP if the space could not be made available.

DeallocateHeap:

target:	Thread	
opcode:	DeallocateHeap	
in parameters:	address	Pointer
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Deallocate memory from the current thread's thread heap. The thread heap memory located at *address* is removed from the thread heap and returned to the system; it is no longer accessible by the thread. The *address* parameter must contain a value previously returned by an *AllocateHeap* operation.

Operation Status: *Status* is SUCCESS if the memory is deallocated from the heap, or NO_HEAP if the *address* parameter does not refer to the beginning of an allocated section of the heap.

Abort:

target:	Thread	
opcode:	Abort	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: Abort the thread. This operation raises a thread abort condition on the *target* thread. If the thread is within the scope of any abort handlers, the thread abort condition will be reported to those handlers. After all abort processing has been completed, the thread is deleted.

Operation Status: *Status* is always SUCCESS.

Object Definition

OBJECT:

Thread Manager

Description: The ThreadManager is the object that controls the creation and deletion of threads.

Object Operations

Create:

target:	ThreadManager	
opcode:	Create	
in parameters:	Target	Capability
	Operation	Operation
	Arg1	Argument
	Arg2	Argument

	ArgN	Argument
return value:	result	InvocationStatus
out parameters:	status	OperationStatus
	newThread	Capability

Description: Create and start a thread executing in an object at the specified operation with the specified arguments and return the capability for it in newThread.

Operation Status: The value SUCCESS is returned in the *status* parameter if the thread was created; otherwise, the value NO_MORE_THREADS is returned (in this case the kernel does not have the resources to create a new thread).

Object Definition

OBJECT:

Transaction Manager

Description: The transaction manager object is provided by the kernel to support the use of atomic transactions within Alpha applications. Operations on this object are used to tell the kernel when a transaction is to begin, to commit, or to abort.

The transaction manager coordinates the actions required to guarantee the atomicity of changes to locked data within each object involved in the transaction. The actual work of updating or rolling back data is performed by standard operations on those objects. The kernel supplied versions of the standard operations may be replaced within an object by application specific implementations, which may effect a different transaction behavior.

Object Operations

BeginTransaction:

target:	TransactionManager	
opcode:	BeginTransaction	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation puts the invoking thread into a new level of atomic transaction and begins an exception scope. If the thread is currently within another transaction level, this invocation begins a nested transaction. Otherwise, this invocation begins a top-level transaction. The transaction and exception scope defined by this operation remains active until a subsequent *EndTransaction* operation is invoked at the same nesting level in the current operation. If any exception occurs while this scope is active, and it remains active after completing any nested exception scope, *BeginTransaction* will return (a second time) with a value indicating the type of exception.

If an *AbortTransaction* operation is invoked within this level of transaction, a transaction abort will be generated on the thread which remains active until an *EndTransaction* operation is successfully executed.

If the thread returns from the operation that performed the *BeginTransaction* operation without having invoked a corresponding *EndTransaction* operation, an *EndTransaction* is automatically supplied by the system.

Operation Status: The value NO_TRANSACTION is returned in the *status* parameter if a new transaction level could not be created. The status NESTING_VIOLATION is returned if the thread is currently executing in the reporting mode of an exception scope. Otherwise the value SUCCESS is returned. If, while the exception scope established by this operation is active, an exception occurs, the operation returns with the value THREAD_ABORT, DEADLINE_EXPIRE, TRANSACTION_ABORT, or MACHINE_CHECK if the exception, abort, or a hardware fault, respectively. An implementation of Alpha may define additional values of status to indicate exceptions due to specific machine error conditions.

EndTransaction:

target:	TransactionManager	
opcode:	EndTransaction	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation causes the current level of atomic transaction to conclude. If the corresponding *BeginTransaction* operation has returned a value reporting an exception, and an exception is still active after the invocation of *EndTransaction* (i.e., the exception was not due to this transaction), then the thread will continue execution at the next enclosing exception scope. If no exception indication has been returned by the corresponding *BeginTransaction*, this operation indicates that the transaction has successfully completed. If the current level is a nested transaction, the commit processing (updating the objects and releasing locks) is deferred until the end of the top-level transaction. For top-level transactions, all modifications to locked data made within this and all nested transactions are committed and become visible to other threads. This entails the following: the *PreCommit* operation is performed on each object traversed by the thread within the scope of the transaction; if all of these report success, the *Commit* operation is invoked on the objects. The system-supplied versions of these operations implement a two phase commit of the modified data to insure that all modifications are made atomically, and the locks and semaphores acquired during the transaction are released. If any of the objects fails or returns an error from the *PreCommit* operation, the *Abort* operation is invoked on all of the objects involved in the transaction.

Operation Status: The value NO_TRANSACTION is returned in the *status* parameter if there is no currently active *BeginTransaction* operation in the current object. The value NESTING_VIOLATION is returned if any nested exception scope has not been closed. Otherwise, the value SUCCESS is returned.

AbortTransaction:

target:	TransactionManager	
opcode:	AbortTransaction	
in parameters:	none	
return value:	result	InvocationStatus
out parameters:	status	OperationStatus

Description: This operation indicates to the kernel that the current level of atomic transaction is to be aborted. The *Abort* operation is performed on each object traversed by the thread within the scope of the transaction. The system-supplied version of the *Abort* operation performs a *V* operation on all Semaphore objects that had a successful *P* operation performed on them within the transaction and releases all locks acquired in the transaction, with all locked data items are restored to their original values. The successful completion of the *AbortTransaction* operation raises an exception condition which remains active until the *EndTransaction* for the current transaction level is invoked. Therefore, in the normal case, this operation will never return to the point of invocation.

Operation Status: The value `NO_TRANSACTION` is returned in the *status* parameter if there has been no matching *BeginTransaction* operation within this object. A successful *AbortTransaction* operation will not return (see above).

Glossary of Terms

- Atomicity** — The property that a sequence of activity appears externally as one uninterrupted action.
- Capability** — Protected names supported by the Alpha kernel. *Client objects* refer to a capability via its *capability index*. References to capabilities are interpreted by the kernel at invocation time.
- Capability Index** — The identifier used by the *client object* to refer to an object represented by a *capability*.
- Capability List** — The per object kernel protected list of *capabilities*.
- Client Object** — An object with its own address space and protection. It communicates with the rest of the system through the *Invoke* and *Invoke Return*. It supports the standard object operations, either explicitly or through kernel-supplied defaults.
- Deadline** — The critical point at which the value of a computation decreases.
- Distributed System** — A collection of communicating nodes whose resources are managed as if for a single entity.
- Exception** — The interruption of normal sequential execution of a *thread*. The types of *abort* in Alpha are *deadline expiration*, *machine fault*, *thread abort*, and *transaction abort*.
- Exception Scope** — A dynamically defined region in which any form of exception can be intercepted.
- Invoke** — The kernel supported mechanism that allows an object to perform an operation on another object.
- Invoke Return** — The kernel supported mechanism that allows an object operation to return results to the object that invoked it.
- Kernel Object** — An object that implements one or more Alpha kernel abstractions. Such objects do not support the standard object operations, and exist within the kernel's address space.
- Node** — A processor or multiprocessor that represents an independent point of failure within a *distributed system* and whose inter-node communication latency significantly exceeds its intra-node communication latency.
- Object** — The fundamental abstraction supported by the Alpha kernel that is the encapsulation of data along with the *operations* for manipulating the data.
- Operation** — The externally identifiable entry into an object that is capable of modifying the data of the object.
- Permanence** — The property of data that ensures its integrity across *node* (or complete system) failure and restart.
- Real-Time** — Activity that is limited by real-world *time constraints*.

- Region** — A contiguous area of memory denoted by an address and length. Locks are specified on regions.
- Serializability** — The property that guarantees that concurrent behavior can be made to appear as if it were performed in a sequential manner.
- System Object** — An object that exists within the kernel's address space. It communicates with the rest of the system through the *Invoke* and *Invoke Return*. It supports the standard object operations, either explicitly or through kernel-supplied defaults.
- Thread** — The fundamental control abstraction supported by the Alpha kernel. *Threads* maintain per invocation *operation* local state and hold the parameters passed in *Invoke* and *Invoke Return*.
- Time Constraint** — A limit on the time to perform some activity.
- Transaction** — An activity that involves a change in the permanent state of an object. The properties of a given *transaction* are some combination of *atomicity*, *serializability*, and *permanence*.

References

- [Northcutt 88a] Northcutt, J. D.
The Alpha Operating System: Requirements and Rationale.
Archons Project Technical Report #88011, Department of Computer
Science, Carnegie-Mellon University, January, 1988.
- [Northcutt 88b] Northcutt, J. D. and Clark, R. K.
The Alpha Operating System: Programming Model.
Archons Project Technical Report #88021, Department of Computer
Science, Carnegie-Mellon University, February, 1988.

Table of Contents

1	Introduction.....	D-1
2	Overview.....	D-2
2.1	Distributed Architecture.....	D-2
2.2	Kernel Layers.....	D-2
2.3	Concurrency and Real-Time.....	D-3
2.3.1	Application Level.....	D-4
2.3.2	Kernel Level.....	D-4
3	Objects.....	D-5
3.1	Object Manager.....	D-5
3.2	Data Structures.....	D-5
3.3	Memory Organization.....	D-7
3.4	Schemas.....	D-8
3.5	Standard Operations.....	D-9
3.5.1	Object Creation.....	D-9
3.5.2	Migration.....	D-10
3.5.3	Information.....	D-11
3.5.4	Permanence.....	D-11
3.5.5	Locking.....	D-12
3.6	Kernel Objects.....	D-13
4	Threads.....	D-14
4.1	Thread Maintenance and Repair.....	D-15
4.2	Exception Tracking and Handling.....	D-15
4.3	Thread Dispatch.....	D-16
4.4	Thread Data Structures.....	D-16
5	Operation Invocation.....	D-18
5.1	Basic Invocation Processing.....	D-18
5.2	Atomic Transaction Invocation.....	D-20
5.3	Replicated Object Invocation.....	D-20
6	Capabilities.....	D-22
6.1	Capability Lists.....	D-22
6.2	Data Structures.....	D-23
7	Semaphores.....	D-26
7.1	Description.....	D-26
7.1.1	Semaphore Manager.....	D-26
7.1.2	Semaphores.....	D-26
7.2	Data Structures.....	D-27

8	Object Replication.....	D-28
8.1	Replicated Object Invocation.....	D-29
8.2	Functional Design.....	D-30
9	Atomic Transactions.....	D-32
9.1	Transaction Management Mechanisms.....	D-32
9.2	General Mechanisms.....	D-34
10	Scheduling Subsystem.....	D-35
10.1	Overview.....	D-35
10.2	Thread Queues.....	D-35
10.2.1	Ready Queue.....	D-36
10.2.2	Wait Queues.....	D-36
10.3	Policy Module.....	D-36
10.3.1	Interface.....	D-37
11	Memory Management.....	D-38
11.1	Design Overview.....	D-39
11.2	Extents.....	D-39
11.2.1	Extent Manager Operations.....	D-40
11.2.2	Extent Operations.....	D-40
11.3	Pager Interface.....	D-41
11.4	Secondary Storage Object Interface.....	D-42
11.5	Page Fault Handling.....	D-42
11.6	Page Reclamation.....	D-43
11.6.1	Working Set.....	D-43
11.6.2	Memory Overcommit.....	D-43
11.6.3	Memory Policy Module.....	D-43
11.7	Machine Independent MMS.....	D-44
11.7.1	Address Maps.....	D-44
11.7.2	Resident Pages.....	D-44
11.7.3	Memory Objects.....	D-44
11.8	Machine Dependent MMS.....	D-44
12	Communications.....	D-46
12.1	Overview.....	D-46
12.2	Kernel Communications.....	D-47
12.2.1	Reliable Message Communication.....	D-47
12.2.2	Remote Invocation and Node Dictionary.....	D-48
12.3	System Communications.....	D-49
12.3.1	Byte Stream.....	D-49
12.3.2	Reliable Datagram.....	D-49
12.3.3	Unreliable Datagram.....	D-49
12.4	Client Communications.....	D-49
12.4.1	Client Communications Interface (Sockets).....	D-50

13 Input/Output.....	D-51
13.1 Generic I/O.....	D-51
13.1.1 The Model.....	D-51
13.1.2 GIO Device Interface.....	D-52
13.2 Alpha I/O Subsystem.....	D-52
13.2.1 GIO Device Driver.....	D-53
13.2.2 Kernel Semaphore Manager.....	D-54
13.2.3 GIO Device Manager.....	D-54
13.2.4 System Events.....	D-55
13.2.5 Special Services.....	D-55
13.2.6 Exception Dispatch.....	D-55
14 Secondary Storage.....	D-56
14.1 Motivation.....	D-56
14.2 Overview.....	D-57
14.3 Secondary Storage Objects.....	D-59
14.4 Partition Objects.....	D-60
14.5 Backing Objects.....	D-61
References.....	D-63

1 Introduction

Alpha is a distributed, object oriented, real-time operating system with mechanisms for reliable computation. This paper describes the function and the design of the various subsystems within the Alpha Kernel Release 2. The design of Release 2 was strongly influenced by the design of Release 1. The earlier version was created at Carnegie Mellon University by J. Duane Northcutt and his colleagues [Northcutt 87]. An introduction to Alpha concepts and its programming model can be found in [Northcutt 88b]. A description of the problems Alpha attempts to solve and the reasoning that led to our decisions can be found in [Northcutt 88a] and [Jensen 88]. This document deals with a realization of the ideas presented in [Jensen 88] and [Reynolds 88].

Flexibility was a major factor in the design of Alpha. As a consequence, a great deal of effort was spent creating simple yet powerful mechanisms that would allow the implementation of a variety of higher level policies. This policy/mechanism separation, similar to [Habermann 76], results in a kernel with a small set of powerful, orthogonal facilities: which is not the same as a trivial operating system or embedded executive.

A number of functions normally associated with contemporary operating systems are not provided by the kernel. These functions are seen as policies that can be implemented via judicious use of the kernel services. One example is a user authentication service. The kernel provides powerful tools for controlling the access to objects. It is left to the application developer to use these tools to build higher level policies to ascertain the identity of a user and whether the user has the right to access a particular object.

The following sections discuss the design of the current implementation of the Alpha kernel. The first section is a very broad overview of the kernel design. The next few sections focus on the design of the services associated with the basic abstractions: objects, threads and invocations. The following sections deal with the ancillary abstractions—capabilities, semaphores, transactions, and the like. The remaining chapters discuss the kernel subsystems such as virtual memory, secondary storage and I/O.

2 Overview

2.1 Distributed Architecture

Alpha is a distributed, object-oriented, real-time operating system. An instance of an Alpha system would be one or more nodes executing a common version of the kernel, typically communicating via a local area network. A node could be any hardware platform with at least one CPU, enough memory and the other resources necessary to execute the kernel. Alpha's distributed nature could, perhaps, be better characterized as decentralized. Centralized design is eschewed wherever possible. Each alpha node executes a fully realized kernel. Each instance of the kernel fully implements the kernel abstractions and provides the facilities for local resource management as well as communication with other nodes. Alpha uses a capability based naming scheme for objects and threads. Capabilities provide a universal, network location independent, name space. Each capability is unique in time and space. Alpha does not use centralized name servers, it uses broadcast protocols and per node lists of local objects. Network location transparent names means that applications executing anywhere on the network see the same name space of objects and that objects can migrate from one node to another without changing their name. Threads are the agents of computation within the Alpha System. All computation in the system, application or kernel, is accomplished by threads. Threads are not bound to a single object. When a thread invokes a operation on an object the thread moves to the invoked object. From the view of an application remote and local invocations have identical semantics. Scheduling, resource allocation and security information is associated with threads and follows them when they move from object to object. All kernel services, even device drivers, are encapsulated within objects. Thus all of the services on a node may be transparently accessed by any node via object invocation.

2.2 Kernel Layers

The Alpha kernel is architecturally organized in layers. Each layer represents a distinct level of abstraction between the hardware and the application. Figure 1 provides an illustration of the logical structure of the Alpha operating system.

The Alpha application layer consists of client objects and system objects. These objects may specified using a variety of programming languages. Client and System objects are described in the chapter on objects. System objects and Client objects share a common interface to the Kernel. The difference between them is that System objects are trusted. This is enforced by the careful distribution of Kernel and System capabilities. System objects provide a simple and elegant method for extending the services provided to Client objects. System objects execute within the Kernel context. This reduces the performance penalty of providing a trusted facility via a System object.

The Services layer supports the abstractions described in the following chapters. These are the application visible features and facilities of the Alpha kernel. The kernel is programmed entirely in the high level systems programming language C++. This layer provides the entirety of the application visible interface. All of the facilities accessible by the application layer are available via object invocation.

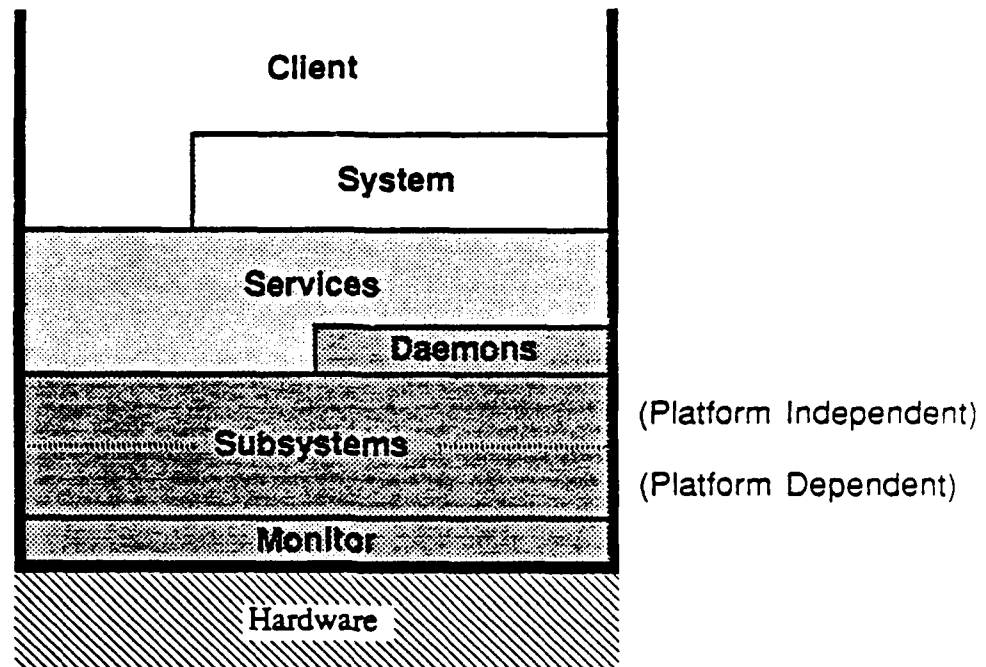


Figure 1: Alpha Logical Structure

The various components of this layer are grouped according to their function. The facilities that support objects are grouped within the object manager, the facilities that support threads are grouped within the thread manager, etc.

The Daemon layer is another user of the Virtual Machine interface but it does not provide an application visible interface. The Daemon layer is the collection of kernel objects and threads that provide services that augment the basic system functions. Some examples are the page-out, garbage collection, and load adjusting daemons.

The Subsystems layer separates the implementation of the Alpha abstractions from the implementation of the services necessary to realize those abstractions on a particular hardware platform. The Subsystems layer deals with those necessary services that are hidden by the Alpha abstractions. These services include Virtual Memory, Secondary Storage, Context Switching, Multiprocessor support, and I/O. Thus the Services layer need not deal with the details of virtual memory, it simply asks for a memory Extent.

This layer is subdivided into platform independent and platform independent layers. This division further isolates hardware dependent behavior and code. This isolation helps to protect the design from being unduly influenced by a particular implementation environment. It also contributes significantly to the portability of Alpha.

2.3 Concurrency and Real-Time

Concurrency or parallel processing is very important to us. The design permits the exploitation of concurrency at the application level and within the kernel itself. The kernel is multi-threaded and fully preemptive. The kernel supports symmetric and asymmetric multiprocessing. These qualities greatly enhance the throughput and the real-time responsiveness of the kernel.

3.1 Application Level

At the application level Alpha supports multiple threads within a single object. This is analogous to light weight processes in RTU [Henize 84]. A rich variety of synchronization services, such as semaphores, locks, transactions, dynamic thread creation and thread local data segments are provided to aid in the development of parallel algorithms. These services are all designed with the needs of real-time applications in mind.

For example, threads waiting on a semaphore are not necessarily unblocked in the order they arrived. Instead, the scheduling subsystem examines the time/value functions associated with each waiting thread and activates the one with the highest value to the system.

3.2 Kernel Level

The kernel is designed to provide the maximum concurrency possible on each hardware platform. The kernel abstractions and all the subsystems are multi-threaded and fully preemptable. At almost any point the current thread of control can be preempted by a more important thread. Multi-threading the kernel provides the structure for supporting multiple processors. It also preserves, within the kernel, Alpha type real-time behavior.

The multi-threaded nature of the kernel allows a single version of the kernel to run on single processor and multi-processor platforms. Maximum concurrency is guaranteed by the judicious use of fine-grain locks. Deadlock in the kernel is avoided by ordering the acquisition of locks. Locks that are used by the kernel and shared with interrupt handlers have the additional characteristic of disabling interrupts when the lock is acquired. Interrupts are enabled when the lock is released.

3 Objects

Alpha objects are passive entities that encapsulate data and the operations that manipulate that data. Objects are named in a network location independent manner via capabilities. Each client object consists of a set of contiguous virtual memory regions called *extents*. Each object's extents exist within a private virtual address domain; the extents of separate objects are never mapped simultaneously into any thread's address space (or *context*). The separation of address spaces for Alpha objects is equivalent to non-intersecting protection domains [Lampson 69]. Each client or system object is an instance of an object type or schema that defines the operations and encapsulated data for the object. In addition, the kernel defines a set of standard operations that are inherited by all client and system objects. Each object instance has a list of capabilities associated with it. The list contains capabilities that have been brought to the object as explicit capability-type invocation parameters. Only the capabilities in this list can be used as the target of an invocation (or passed as a parameter) from that object. Capabilities can be deleted from the list via explicit operations or as the result of a kernel garbage-collection process.

3.1 Object Manager

The Object Manager encapsulates within the kernel the data structures and operations that create and manipulate objects. The Object Manager controls or implements the Object Control Block, the Permanent Object List, the Dictionary, unblocking of locked threads, the recovery daemon, the propagation of updates to replicated objects, the standard operations and other services. The Object Manager does not present an explicit interface to the application developer though it does implement the standard operations that are defined on all objects. The Object Manager works with the Communications Subsystem to control the propagation and utilization of capabilities.

3.2 Data Structures

Within the kernel, objects consist of the various control structures used by the kernel to manage the object and a collection of memory extents. The most important of the object control structures are the Object Control Block and the *C_LIST*. The Object Control Block contains all the status and control information for an object. The *C_LIST* is the kernel maintained list of capabilities owned by an object. The information kept in the Object Control Block includes the memory extents comprising the object, the operation entry points, the threads currently invoking the object, object attributes, etc.

The components of the Object Control Block are:

- CapaPtr — *pointer to its own Capa_ID*
- C_ListPtr — *pointer to the C_LIST for permanent CAPAs*
- EntryPoints — *list of entry-point/operation-name pairs*
- MemoryExtents — *pointer to the list of Address Map entries (memory extent, virtual address) defining the object's address space*
- ThreadFrames — *pointer to the list of thread invocation frames invoking this object*

AttributeFlags

Client / System / Kernel	
Permanent	
Replicated	
Atomic	
Permanent_SSOs	— <i>list of Secondary Storage Objects providing permanence for each of the extents</i>
Locks	— <i>pointer to the list of Lock Control Blocks</i>
LockedThreads	— <i>list of threads waiting for locks</i>
DateCreated	
DateLastModified	
SchemaCapa	— <i>capability of Schema for the object type</i>
MaxSize	— <i>maximum size of the object</i>
MinSize	— <i>minimum size of the object</i>
CurrentSize	— <i>current size of the object</i>
NetworkLocation	— <i>current network location of the object</i>
MigrationList	— <i>List of objects that migrate when this object migrates (usually semaphores)</i>
NextSemaphore	— <i>List of associated semaphores</i>

The C_LIST is a data structure used by the kernel to manage the actual capabilities accessible by the object. When a capability arrives as an invocation parameter, the kernel adds it to the C_LIST and supplies a capability index by which it can be referenced. The C_LIST is maintained in two separate parts. The part of the C_LIST associated with permanent capabilities is kept in the Object Control Block. The part of the C_LIST associated with transient capabilities is kept with the thread's Frame Control Block.

The C_LISTs of all the objects on a node are combined to build the Dictionary. The Dictionary is an ordered list of capabilities that is used by the invocation mechanism to locate objects. For more information about capabilities and the Dictionary see the section on Capabilities.

The Extents of permanent objects are backed up by non-volatile secondary storage. The object control block, including the C_LIST, that is associated with a permanent object is also backed up by non-volatile secondary storage. The Object Manager owns a special Secondary Storage Object (SSO) in the Root Partition that is used to back up the Permanent Object List (POL). This object contains the control blocks and C_LISTs of each permanent object. The entry in the POL SSO is updated each time an UPDATE operation is performed on the matching object.

There are three types of extents within an object: DATA, HEAP, and TEXT. An object always has at least one DATA extent and one TEXT extent. The figure below illustrates the layout of a typical object within the virtual address space of a thread.

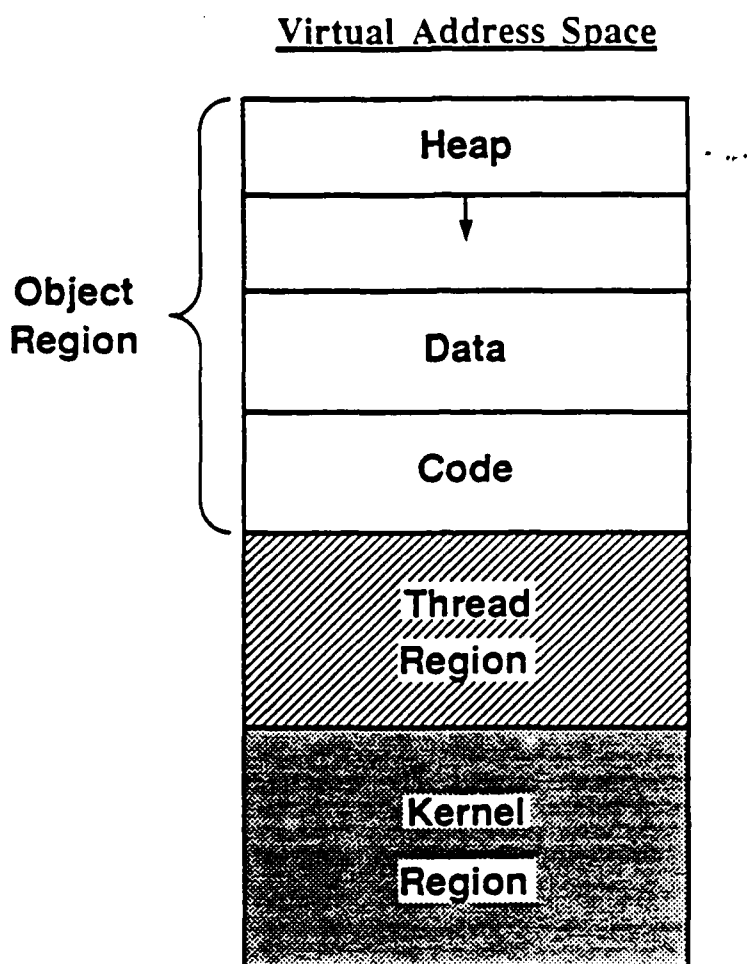


Figure 2: Alpha Virtual Address Space Layout

The data portion of objects consists of three sub-parts. The first two, statically allocated, uninitialized data and statically allocated, initialized data are composed of DATA extents. Dynamically allocated, uninitialized data (i.e., heap storage) is the third sub-part. The data part of an object represents only the global data associated with the object. This global data is accessible by any thread executing within the object. The stacks (e.g., automatic variables of subroutines) are provided on a per-thread basis and are not associated with the object proper.

3.3 Memory Organization

Each client object resides in a separate virtual address domain, allowing the abstraction of disjoint object memory addressing domains to be enforced by a processing node's memory management hardware. Addresses generated by a thread in one object cannot reference memory locations in any other object's memory address domain. Furthermore, portions of an object are similarly protected by the memory management hardware; for example, the code portion is protected from writing.

Note that an address space can be divided into three separate regions, thread, object, and kernel. This is the normal state when an address space is active, i.e., a thread is

executing within an object. The kernel portion of the address space is mapped but protected from unprivileged access. This protection is enforced by the memory management hardware. In order for a thread to access the kernel, a kernel object must be invoked. This is similar to a system call in a traditional operating system.

At any point in time the executable image of a non-replicated object resides entirely on a single node, although at times portions of an object's image may be in secondary storage because of paging or swapping activities. When a thread executing within an object invokes an operation on another object, the target object is mapped into the address space of the invoking thread. When an operation on an object is complete, the invoked object is replaced with the object that made the invocation. The concurrent access of objects by threads makes it possible to have an object's image be shared by more than one context at a time (all at a single node).

System objects reside within the address space of the kernel. They are composed of the same types of memory extents as client objects. Because they are part of the kernel address space, no explicit mapping operation is necessary whenever one is invoked.

Replicated objects may exist on more than one node at a time. With sufficiently relaxed consistency restraints threads may execute within different members of the replicated set at the same point in time.

3.4 Schemas

All objects are instances of type or schema objects. Schema objects are kernel objects that define client or system object types. The user can create new Schema objects to describe new types of client objects. Each Schema object is a template containing a description of the structure, operations, and initial state of the data of the object type; it is used by the kernel to create instances of the type of object described by the template. An object instance is created by invoking a *CreateSystem* or *Create* operation on the appropriate Schema object and providing the attributes, such as permanence, that the object instance should be created with.

The Schema contains a list of extent structures which describes the collection of extents that comprise the object. Each object must have at least one TEXT and DATA extent. Extents may sparsely populate an address space but they may not overlap one another within an address space.

The Schema also contains a list of operation structures which gives the name of each operation, its entry point or starting address (relative to the beginning of a TEXT extent), the number of IN and OUT Capabilities, and the length of the IN and OUT parameter block.

A list of relocation entries provides information which is used whenever the virtual addresses of the extents are not fixed. In that case, the kernel will assign the extent to virtual addresses at object instance creation time, and use the relocation entries to establish inter-extent references.

The Schema Manager is a kernel object that can be used to create new, uninitialized user schema objects. These newly created user schema objects are then initialized with the information they will use to define a type of object. The Schema Manager acts as a meta-schema object.

Schema objects provide a convenient source of initial backing store for the Virtual Memory Subsystem. TEXT extents are treated as read-only objects which allows them to be shared by any number of objects of the same type. Data extents can be shared by a VM system that supports copy-on-write behavior.

The section on Object Creation describes how Schemas are used to build object instances.

3.5 Standard Operations

In addition to the application-specified operations defined on objects (via Schema objects), the system provides a set of standard operations that are inherited by every client and system object. The standard operations provide the means to manipulate objects or to obtain information from the kernel about them.

The code for standard operations is part of the local kernel and is shared by all objects that reside on the node. The implementation of these operations is analogous with that for the kernel-defined objects. That is, while they present an object/operation interface, they may use and manipulate diverse kernel data structures.

The standard operations are:

- | | |
|------------|--------------------|
| 1. UPDATE | 7. UNLOCK |
| 2. RESTORE | 8. PRECOMMIT |
| 3. INFO | 9. COMMIT |
| 4. MIGRATE | 10. ABORT |
| 5. LOCK | 11. ALLOCATEHEAP |
| 6. C_LOCK | 12. DEALLOCATEHEAP |

Each standard operation has two names, the regular name and the reserved name which is the concatenation of the regular name and the suffix `_$Kernel`. The regular and reserved names for the update operation are `update` and `update_$Kernel` respectively. The regular name of a standard operation can be overloaded.

The regular name for a standard operation is overloaded when an object uses the regular name for a locally defined operation. When an object redefines a standard operation the kernel uses the object version of the operation as a target for invocations instead of the kernel version. When an object is created the kernel builds a list of operation names and entry points that is associated with each object control block. When this list is made, during the `CREATE` operation, the kernel protects the reserved operation names from being redefined and it attaches the appropriate entry point the regular operation names.

If the target of the invocation is a client object and the kernel version standard operation is used then the code simply calls the operation. Since the invocation code and the kernel version of all the standard operations live in the kernel address space no context switch is necessary. If a client object version is used then a context switch is necessary.

3.5.1 Object Creation

To instantiate an object, a `CREATE` invocation is made on the appropriate Schema object. There is a section on Schema objects that discusses them in greater detail.

A CREATE invocation takes the following steps when instantiating an object. The Schema object that is serving as the template is checked for correctness. Kernel data structures, such as the Object Control Block, that are necessary to manage the object are allocated. Memory extents are allocated from the virtual memory subsystem for each of the extents specified in the Schema. If the new object is to be Permanent, Secondary Storage objects are allocated for each extent. The operations and other characteristics specified in the Schema are registered in the Object Control Block. The new object is registered with the Dictionary. The Dictionary generates a new Capability and associates it with the new object. This capability is returned as an OUT parameter of the CREATE invocation. Once the object has been created operations can be invoked on it.

A wide spectrum of paging policies, including pre-paged and working set based demand paging, can be supported by Alpha on a per object basis. Paging policy is set by the Memory Management and Secondary Storage Subsystems. The careful selection of the correct virtual memory/paging policy is important because of its effect on system throughput and real-time (e.g., predictable) performance. For addition information on the chapter on Memory Management Subsystem and the chapter on the Secondary Storage Subsystem.

The kernel attempts to make efficient use of primary and secondary storage. This includes the utilization of copy-on-write and copy-on-access techniques where ever reasonable. These techniques help reduce the proliferation of identical copies of data structures. However, these techniques occasionally conflict with the need reliability. When this occurs the decision is always made in favor of reliability.

An example of this conflict is the use of Schema objects to provide the single, sharable source of TEXT pages to the Memory Management Subsystem. This is an appealing idea because the TEXT does not normally change and it must be in the Schema object before the first object of that type can be created. However, reliability problems occur when the objects created by the Schema live on different nodes or have greater reliability attributes than the Schema. In each case it is possible an object created via the Schema to be available at a time when the Schema is not. This means that the Schema would not be available to serve as the source for TEXT pages. In order to avoid compromising the reliability of Alpha applications the source for TEXT pages of permanent and replicated objects is guaranteed to have the same availability as the objects they support.

3.5.2 Migration

When an object is instantiated it is created on the node the parent Schema resides. If it is desirable that the object reside on a different object, it must be migrated. The MIGRATE operation is used to migrate objects. If the object was replicated, each member of the replicated set must be individually migrated to the desired network location.

When a thread attempts to migrate an object the Scheduling subsystem is notified so that it can attempt to ease the effect on other active threads. Once the Scheduler allows the migrate to begin all active threads within the object are paused. Any threads that attempt to invoke operations on the object once the migrate operation has begun are paused before they can actually enter the object. They will be resumed after the migrate operation has completed. If the migration is successful they will migrate with the object, otherwise they will execute in the object at its original location.

If the migration is successful, a message is sent indicating the object's new location and used by all the Dictionaries (for more information about the Dictionary see the chapters on Capabilities and Invocations) that have entries for the object. This message is piggybacked on the success protocol. Next, any paused threads are resumed. Paused threads that were not allowed to enter the object are redirected to the object's new location where they are allowed to retry the invocation.

If the migration fails, the paused threads are all resumed and an indication of the failure is returned to the thread that attempted the MIGRATE.

3.5.3 Information

The INFO operation is essentially a query of the Object Control Block. The following information is returned:

max_size	— <i>max size of the object</i>
minimum_size	— <i>min size of the object</i>
current_size	— <i>current size of the object</i>
top_of_heap	— <i>this address represents the top of the object heap</i>
active_threads	— <i>number of active threads within the object</i>
thread_segs	— <i>number of thread segments within the object</i>
location	— <i>network address</i>
replication	— <i>replicated ? (true/false) number of elements in set replication strategy, etc.</i>
atomic	— <i>atomic ? (true/false) transaction strategy</i>
transactions	— <i>number of active transactions affecting the object</i>
permanent	— <i>the object is permanent (true/false)</i>

Note that the information returned may be stale. Other threads executing within the system might change the state of the object after its state is examined but before the state information can be reported.

3.5.4 Permanence

Alpha supports the notion of object permanence. Permanence is an optional attribute of any system or client object. The goal of object permanence is to provide the facilities for objects to survive the failure of the node they reside on. If the node failure is temporary, then resetting the node should restore all of its permanent objects. If the failure is permanent then it should be possible to recover its permanent objects by migrating them to active nodes.

Alpha does not depend on esoteric non-volatile primary storage to realize object permanence. The UPDATE, RECOVER, LOCK and COMMIT operations work with the Secondary Storage System to provide object permanence.

Among the standard operations defined on objects, the UPDATE and RECOVER operations are the most significant with respect to the object's relationship to its secondary storage image. The UPDATE operation is used to make an object's permanent secondary storage image reflect the current state of a specified range of the object's virtual address space. The RECOVER operation is used to make an object's virtual memory

image reflect the current state of a specified range of the object's permanent secondary storage image.

The Secondary Storage System plays a major role supporting both Atomic Transactions and the Memory Management Subsystem. For more information see the following subsection on Locks, the chapter on Transactions and the chapter on the Secondary Storage Subsystem.

3.5.5 Locking

Locks are a concurrency control mechanism provided to mediate concurrent access to data within an object. Due to the limitations of contemporary hardware platforms, locks are advisory in nature. Enforcement of correct usage of locks is left to the programming environment. A Lock consists of a kernel maintained data structure containing an indicator of the current state of the lock (i.e., its current mode), an indication of which thread (actually thread segment) holds the lock, and a specification of the data region associated with the lock.

The Lock control blocks look like this:

```

range      — the addresses controlled by this lock
mode       — type of lock
thread     — thread segment that owns the lock
next lock  — lock list pointer

```

To perform lock operations, the kernel determines whether the lock request is compatible with the current state of the lock, based on the kernel's lock compatibility table. If a lock request is found to be compatible, the lock is granted and the state of the Lock object is updated. If the lock request is not compatible with the current mode of the Lock object, the thread making the request is blocked and placed into the Lock object's blocked thread list.

Modes:

```

No Readers allowed
Exclusive Read
Multiple Read
no Read Restrictions

```

```

No Writers allowed
Exclusive Write
Multiple Write
no Write Restrictions

```

When a lock is released or a thread waiting for a lock goes away (deadline expired or aborted), the scheduling subsystem is notified. The scheduling subsystem has the responsibility ordering the list of threads waiting on a lock. Once the scheduling system has ordered the list it calls a special operation provided by the Object Manager and passes the list as a parameter. This operation uses the ordered list to create a list of threads that are eligible for unblocking and returns it to the scheduling subsystem. To be eligible the threads must be granted their lock request. The lock request is granted if there are

no conflicting locks and there is no thread of greater importance waiting for a conflicting lock. This prevents important threads from being starved by threads of lesser importance.

In addition to their data structures, Locks in Alpha have special storage areas associated with them. This storage area is provided for the purpose of maintaining a type of write-ahead log in support of the functionality of atomic transactions. When a write-mode lock is acquired by a thread executing within a transaction, the kernel copies the data item associated with the Lock object into its log area. Should a transaction successfully commit, the Lock object's log can be discarded. Should the transaction abort, however, the kernel restores the data item to its previous state through the use of the lock's log. If the object is Permanent the log areas have special secondary storage objects associated with them.

Locks have unique semantics when used within transactions, however, in support of programming-level uniformity, the lock mechanism is designed to appear the same to the programmer, whether or not the threads that make use of them are currently executing atomic transactions. Locks acquired from within transactions are held beyond the point at which clients invoke UNLOCK operations, and released only when the transaction commits. Just as with Semaphore objects, all locks acquired during the execution of a transaction must be released should it abort.

The Alpha kernel makes use of locks as the primary synchronization mechanism for transactions (as opposed to optimistic [Kung 81] or timestamp-based schemes [Bernstein 81]) because other schemes rely on the notions of delays or roll-backs as a *fundamental* part of their normal synchronization activity. These characteristics make such concurrency control schemes undesirable for use in real-time command and control systems, where timeliness is a part of the specification of the correct behavior of the system.

3.6 Kernel Objects

All of the application visible kernel services are presented with an object interface. To the application, Semaphores, Schemas, the Thread Manager, etc. all appear to be objects. However, they are not implemented in the same manner as client or system objects.

4 Threads

Threads are the unit of activity within an Alpha system. Each thread defines an address space, or *context* that is in effect while the thread is executing. The context can have a number of memory *extents*, which are regions of contiguous virtual memory, mapped simultaneously.

When a thread is executing an operation invocation on a client object, it has access to the extents that define the object, the extents that define the invocation frame, and the extent containing the thread heap.

The invocation frame is composed of three extents that define a private execution environment for the thread within the invocation. These extents are:

Client stack: Each invocation frame has a separate, protected stack extent into which the operation code can put thread-private data, such as automatic variables and subroutine return addresses. The stack extent is protected from all other threads and from the same thread when executing within other object invocations. When an invocation is complete, the client stack extent allocated for that invocation is discarded.

In parameters: This extent is used to communicate with the calling operation invocation. It contains the parameters that were passed in for the current operation invocation. It is also used to return parameters to the calling invocation. (Note that the in-parameter extent of an invoked frame is the out-parameter extent of the invoker.)

Out parameters: This extent is used to communicate information to any operations that are invoked from within the current operation invocation. It contains the parameters that are passed out to any nested operation invocations (where it becomes the in-parameter extent). It is also used to receive the return (or reply) parameters from nested invocations. When an invocation is complete, the out-parameter extent allocated for that invocation is discarded.

As a thread invokes an operation, the extents that define the object it was previously executing are removed from the context, and the extents for the new object are mapped in. In addition, the old in-parameter extent is mapped out, the old out-parameter extent is remapped to become the new in-parameter extent, a new (empty) out-parameter extent is allocated, and a new (empty) client stack extent is allocated. On return from the invocation, this process is reversed, and the out-parameter and stack extents of the returning invocation frame are discarded.

The thread heap is an extent that is always mapped into the same virtual address of a thread context. It is used as a high bandwidth communication channel for operation invocations of the same thread. When data is placed within the thread heap, a pointer to it can be passed as an invocation parameter and no data copying is necessary (for all invocations of objects on the same node).

In addition to the client-visible extents, a thread also has a *kernel stack* extent. This extent provides a place to put thread-local data, such as automatic variables and subroutine return addresses, when the thread is executing within the kernel context. This extent is inaccessible when the thread is not executing within the kernel context.

4.1 Thread Maintenance and Repair

When a thread is created, it begins execution by invoking an operation on an object that is specified by the creation parameters. The invocation frame associated with this invocation is called the *root* of the thread. If the thread ever returns from this invocation, the thread is deleted. At any given time, the active invocation frame of the thread is called the *head* of the thread.

Whenever a thread invokes an object on a remote object, the thread head moves to the node containing the object. This is managed within the system by the creation of a surrogate thread information block on each node traversed by the thread.

The notion of *thread segment* is used to manage the movements of a thread among nodes. A thread segment consists of a contiguous sequence of invocations on objects within a single node. (A thread that never invokes a remote object has only one segment.) Note that a thread may have several segments on a single node which are separated by intervening segments on remote nodes. The interrelationship between thread segments is tracked by the Alpha system to deal with node failure or the deletion of objects within a thread segment.

When the kernel on an Alpha node detects the failure of another node, it searches a list of local segments which are attached to segments on the failed node (via invocations to or from an object on the node). When a thread segment has invoked an operation on an object on a failed node, the kernel aborts the invocation and returns a status value indicating the invocation was unsuccessful due to node failure.

Segments of threads that arrived from a failed node are called *orphans*. Orphans represent computations whose results can never be used because they would be reported to a now nonexistent object. Therefore, it is necessary to remove orphans in order to free their resources for other threads to use. Orphans are deleted after the orphan segments are notified via the exception mechanism and are, thereby, provided the opportunity to clean up. Thread segments downstream from orphan segments are also notified.

4.2 Exception Tracking and Handling

During their execution of an object operation, threads can establish exception scopes that contain code that will be notified of exception conditions. Each exception scope must begin and end within the same invocation frame and may be nested within other exception scopes. When an exception occurs, the condition is reported from the innermost nested scope outward until the exception is no longer active (pending).

The kernel tracks exception scopes by keeping a list of exception information records for each invocation frame of the thread. The exception information record contains the type of the exception scope, the nesting level, the machine state to restore when reporting an exception. As exception scopes are nested, new records are pushed on the list; they are popped as each scope ends.

When an exception condition occurs, the top exception information record is used to report the exception condition, then removed. If the condition is still active, the next information record is used to report the condition to that scope. If there are no more records in the current frame, the invocation is aborted and reporting begins at the next higher invocation frame. This continues until the exception is no longer active or the root frame is aborted (thereby deleting the thread).

4.3 Thread Dispatch

The threads on each node are multiplexed onto the node's application processors. Each thread on a node is in one of the following states:

RUNNING: A thread that is currently assigned to an application processor.

READY: A thread that is not assigned to an application processor but that can be run at any time.

BLOCKED: A thread that is waiting for some event (synchronization, I/O, etc.).

STOPPED: A thread that has been halted by an explicit STOP operation.

SLEEPING: A thread that has been halted by an explicit DELAY operation.

The Scheduling subsystem on each node creates and maintains a list of ready threads in decreasing order of value to the system of executing them at the current time. This list is called the *ready queue*. When an application processor becomes idle or the relative value of a thread changes, the ready queue is used to select the proper thread to execute.

Dispatching a thread to an application processor involves the following steps:

- Removing it from the ready queue.
- Mapping its context (including the current invocation frame and invoked object) into the memory management hardware of the application processor.
- Restoring the state of the thread (register contents, etc.) into the application processor.
- Resuming its execution.

4.4 Thread Data Structures

Each thread within an Alpha system has a Thread Control Block (TCB) associated with it that maintains the state of the thread, including the stack of thread segments and invocation frames. When a thread migrates to other nodes to perform invocations on remote objects, a surrogate Thread Control Block is maintained on each node traversed by the thread. The stack of segments for a TCB on each node is maintained as a list of Thread Segment Control Blocks.

The components of the Thread Control Block are:

ThreadID	— <i>the AlphaBitID of the thread</i>
ThreadSegs	— <i>pointer to stack of TSCBs for the thread on this node (most deeply nested first)</i>
CurrentFrame	— <i>pointer to the most deeply nested invocation frame on this node</i>
RootNode	— <i>the node on which the thread is rooted</i>
Context	— <i>list of Address Map entries defining the thread's current context</i>
MMap	— <i>machine dependent memory management information</i>
KernelStack	— <i>extent containing the kernel stack</i>
ThreadHeap	— <i>extent containing the thread heap</i>
Next	— <i>next TCB in ready or wait queue</i>

RunState	— execution state (<i>RUNNING</i> , <i>READY</i> , <i>BLOCKED</i> , <i>STOPPED</i> , or <i>SLEEPING</i>)
Preemptable	— <i>TRUE</i> except during critical section
ExceptionState	— <i>NORMAL</i> or <i>REPORTING</i> mode
RunTime	— cumulative CPU time

The Thread Segment Control Block (TSCB) maintains the relationship between segments of a thread. It also contains the stack of invocation frames for the local invocations constituting the segment as a list of Frame Control Blocks. The components of a TSCB are:

SegID	— system-wide unique ID for this segment
Head	— the invocation frame FCB that is the head of the thread or thread segment; it is the current point of activity of the thread within this thread segment
FromSegID	— the SegID of the next-higher segment in the nesting (NULL if this is the root)
FromNode	— node ID of the next-higher nested segment (NULL for the root)
FromNodeNext	— next TSCB in the per-node 'from' list
FromNodePrev	— previous TSCB in the per-node 'from' list
ToNode	— node ID of the next-lower nested segment
ToNodeNext	— next TSCB in the per-node 'to' list
ToNodePrev	— previous TSCB in the per-node 'to' list
TransNesting	— the atomic transaction nesting level
Thread	— TCB pointer for thread owning this segment
Next	— next TSCB in the per-thread segment list

A Frame Control Block consists of:

Object	— the OCB pointer to the invoked object
ObjectNext	— next FCB in the per-object frame list
ObjectPrev	— previous FCB in the per-object frame list
C_ListPtr	— pointer to the C_LIST for transient CAPAs
Stack	— the client stack extent
InParam	— the in-parameter extent
OutParam	— the out-parameter extent
State	— the machine state (register contents, etc.) of the thread within the frame [TBD]
DeadlineInfo	— list of time constraint records [TBD]
ExceptionInfo	— list of exception and transaction scope records [TBD]
ResourceInfo	— list of resource-use records [TBD]
Thread	— TCB pointer for thread owning this frame
Segment	— TSCB pointer for segment owning this frame
Next	— next FCB in the per-segment frame list

5 Operation Invocation

In the Alpha kernel, the movement of a point of control associated with a computation (i.e., a thread) among objects is accomplished through the invocation of operations on objects. In addition to encapsulating code and data, objects provide the entry points for the invocation of operations. Operations may have one or more parameters associated with them, however, the object model of programming suggests that large amounts of data should not be passed as parameters to invocations.

Because the operation invocation facility is intended to be used for all interactions among objects in the system, the facility is designed to provide, at the lowest possible level in the system, physical-location transparent access to objects. All system-provided services are made available through the operation invocation facility; the kernel routines are designed to provide object-like interfaces to the clients. Various kernel facilities (most notably the virtual memory facility) also make use of the invocation facility as a part of their normal function. This usage of the invocation facility provides the kernel with full visibility of all movements of threads among objects (an important feature in the support of atomic transactions and for system monitoring).

The operation invocation facility has been designed to provide a uniform interface to programmers. Operation invocations appear the same regardless of the type of object (i.e., client or kernel) making the invocation and the type of object being invoked—i.e., kernel or client object, system service object, replicated object, local or remote.

5.1 Basic Invocation Processing

An invocation begins with the marshaling of parameters done by placing the target object's capability index and restriction tuple, the restriction tuple for the thread self capability, a (capability index, restriction tuple) pair for each passed capability, and all of the passed variables into the invoking object's out-parameter extent.

When the out-parameter extent has been loaded with the invocation parameters, the object traps into the kernel. This causes the application processing element to enter supervisor mode and begin the invocation process. The first step of the invocation procedure deals with the translation and validation of any capabilities passed as parameters. This involves applying the passed capability indices to the invoking object's c-list to obtain the desired capability descriptors, applying any additional restrictions found in the parameter list's restriction tuples, and examining the result to determine if there are any restrictions that would make any capability's desired use illegal. In each invocation there is at least one capability that must be translated and validated, and that is the capability used to indicate the object that is the destination of the invocation. If the passed capabilities are in any way improper (e.g., non-existent or usage restricted), the invocation is terminated and a failure indication is returned. When all of the capabilities are successfully translated, the invocation proceeds by calling the internal invocation routine.

At this point in the operation invocation process, an internal operation invocation interface is used. This interface is provided by the kernel for use both by objects, and the kernel itself, to perform location-independent invocation of operations on a variety of kernel entities (i.e., anything represented by capability names). This routine takes a pointer to an operation invocation parameter extent (in a standard format) and continues with the

operation invocation, and is used by all operation invocations from client objects, kernel objects, or functions of the kernel itself. The internal operation invocation routine begins by determining whether the destination of the invocation is a local system service object. If so, the system service object handler is called, otherwise a lookup operation is performed on the local Dictionary to determine if the destination Alphabit is local to this node. If the destination entity's Alphabit Identifier is found in the local node's Dictionary, the routine that handles invocations on local entities is called. If the destination entity is not found locally, the routine that interacts with the communication subsystem is called to issue a remote invocation. (An exception to this procedure is made in the case of inclusively replicated entities, where the remote invocation is done regardless of whether the destination entity exists locally.)

The code that handles local invocations determines the type of entity on which the operation is to be performed and calls the appropriate routine for each of the potential destination types. This routine is invoked either as a result of a locally initiated operation invocation, or as a result of an incoming remote operation invocation (in which case it is initiated by a signal from the communication subsystem). Among the routines that handle the local invocations are routines that provide the standard operations on threads and secondary storage objects, in addition to the two main routines that handle invocations on kernel objects and client objects.

The part of the kernel that performs the invocation of operations on kernel objects begins by adding any passed capabilities to the invoked object's c-list as transient capabilities, and replacing the capabilities in the invoked object's operation invocation parameter extent with indices within the object's c-list. Following this, the entry point of the desired operation is obtained by locating the object's entry point block (as referenced by the object's OCB), determining if the operation name is defined on the object, and then searching the entry point block to get the address of the operation entry point. Once the invoked operation's entry point has been obtained, the invoked object is entered and execution begun at the start of the desired operation. Note that regardless of the context from which the invocation originated, all operations on kernel objects are performed within the kernel.

The code that is responsible for performing operation invocations on client objects functions similarly to the one that handles kernel object invocations. However, a number of additional activities must be performed in order to initiate an operation on a client object. The client object invocation routine locates the virtual memory structures associated with the destination client object maps the invoking client object out of the invoking thread's context (if there was one), maps the invoked client object into the context, locates and validates the operation's entry point, and begins execution in user mode within the thread's context (which now includes the invoked object).

If the destination entity is not found in the local Dictionary, the communication subsystem is entered to perform a remote invocation. When the invoked operation completes, the communications subsystem provides the system with return parameters, along with any changes made to the thread's environment at the remote node. This information is passed back to the invoking object, and the surrogate thread is deactivated.

In all of the cases outlined above, once an invoked operation completes, the sequence of steps taken is retraced to return to the invoking object, passing back any results from the

invocation. The return parameters are passed back in the parameter extent, and all manipulations of virtual memory are reversed.

5.2 Atomic Transaction Invocation

If an invocation is made by a thread within an atomic transaction, the property of thread visit notification is added to the standard invocation service (i.e., timer-based, positive-acknowledged RPC with retries, duplicate suppression, and orphan detection and elimination). Thread visit notification involves having the kernel track invocations made by a thread within an atomic transaction, and should a thread break while the thread is executing within the atomic transaction, all objects visited (i.e., having had operations invoked on them) as a part of that transaction are notified of the transaction's failure. Visit notification takes the form of an unsolicited invocation of the *AbortTransaction* operation on the TransactionManager object instances at the affected nodes. The kernel begins tracking a thread when it invokes a *BeginTransaction* operation on the TransactionManager object, and stops tracking it when a matching *EndTransaction* or *AbortTransaction* operation is invoked on the TransactionManager object. Nesting of atomic transactions is supported by a kernel-maintained count of the depth of transaction nesting. The kernel uses the atomic transaction nesting count to determine when the outermost transaction is exited and tracking can be halted.

5.3 Replicated Object Invocation

The various forms of replication provided in Alpha are supported primarily through features of the invocation facility. The replication support is provided in a manner that is transparent to the client; a client performs operation invocations in the same manner regardless of whether the target object is replicated. The operation invocation facility determines whether or not the object being invoked is replicated, and carries out the appropriate operation invocation protocol for that type of object.

In the case of inclusive replication, the communication mechanism issues a message addressed to the replicated object that serves as the first phase of a two-phase invocation protocol. Each object replica that receives such a first-phase message responds either positively or negatively depending on its current state. The communication subsystem at the invoking node receives the responses and attempts to gather a minimum number of positive responses from the currently existing replicas. This minimum number of responses is specified for each operation when an inclusively replicated object is instantiated. If the initiating node does not obtain the specified minimum number of positive acknowledgments from the replicas, it issues an abort message and retries the first phase of the replicated invoke protocol. If positive acknowledgments are received from all of the replicas, the initiating communication subsystem carries out the second phase of the protocol by issuing an invocation message to all of the instances of the replicated object. In order to serialize the actions of threads within the replicas, all subsequent first phase invocations are negatively acknowledged until there are no more outstanding acknowledged invocations.

When a node's communication subsystem detects an incoming invocation from a replicated object, it waits until the invocations are received from all the necessary replicas before passing the (single) operation invocation signal to the application processor. The

number of replicas involved in any replicated operation invocation is determined by the node initiating the invocation during the first phase of the protocol, and is passed to the replicas in the second phase. All invocations made by a replicated object include this count of the currently involved number of replicas that is used (in a fashion similar to that in [Cooper 84]) to merge the outgoing invocations of a replicated object.

In the case of exclusive replication, a similar type of two phase invocation protocol is used. In the first phase the communication subsystem of the node where the invocation is initiated issues a request message addressed to the replicated object. All of the currently existing replicas respond to this first-phase message and the initiating node selects one of the replicas which responded in a given period of time. The second phase of this invocation protocol involves sending an invocation message to the selected replica. The reply from such an invocation is handled in the normal manner, and outgoing invocations are also handled the same as those from non-replicated objects.

The replica selection algorithm currently in use chooses the first replica to respond to the first phase message. This portion of the operation invocation mechanism was designed to be easily modified or replaced, in order to permit experimentation with differing replica selection algorithms. In order to use more sophisticated selection algorithms, differing types of kernel information (e.g., load statistics or resource utilization estimates) may be required in the replicas' response messages. Thus, the kernel's design permits a node's communication subsystem to have access to the kernel's internal data structures where such information can be obtained.

If more than one replica of an exclusively replicated object exists at a node at some point in time, the replicas are chained together at the same location in the node's Dictionary. Therefore, when an invocation is made on an exclusively replicated object, the first instance of any local replicas that exist on a node is used.

With both of the currently provided forms of replication, the communication subsystems use logical addressing for messages directed to all replicas of a particular object. All logically addressed messages also contain the node's physical address, and so response messages from replicas are sent using the initiating node's physical address, along with the physical address of the replica's node. Subsequent messages from the invoking node (e.g., second phase or retry messages) are sent to the physical addresses of specific nodes to reduce the amount of logically addressed traffic on the communications subnetwork.

6 Capabilities

Capabilities provide system wide logical names or addresses for objects and threads. Access control is provided by the careful distribution of capabilities. The possession of a capability implicitly grants permission to invoke operations on the object. Capabilities are managed by the kernel and cannot be directly accessed by the client. Capabilities in Alpha are unforgeable unlike those used in other systems [Needham 82]. All interactions among objects in Alpha are via invocations and all invocations use system protected capabilities. Thus, capabilities provide a simple, elegant, globally uniform access control mechanism.

6.1 Capability Lists

Each object has a *c_list* that contains the capabilities possessed by the object. The *c_list* is maintained by the system and is not directly accessible by the client. Capabilities in the *c_list* are referred to by their position or their index within the list. This is analogous to the relationship between file descriptors and file description control blocks in UNIX.

Each thread has a *c_list* associated with each invocation frame. This *c_list* contains the thread and object SELF capabilities and any capabilities that were arguments of the invocation. Capabilities that are local to a particular thread are *transient*. The application sees only one list, which is the union of the *c_list* of the object and the *c_list* of the current thread.

Capabilities carried into an object as invocation parameters are *transient*, and can only be used by another thread in the object if the capability has explicitly been made *permanent* within the object by executing the *SaveCapa* operation on the object. Saving generates a new capability index for the capability, and this index can then be stored into some instance variable for later use by another thread. The original capability index, representing the transient capability, is still available to the object as long as the thread that brought it is still executing in the object.

When a thread returns from an object, all transient capabilities brought by the thread to the object are removed from the object. That is, the capability indices that referred to the transient capabilities will no longer be valid within the object. Capabilities that are returned as output parameters appear as transient capabilities in the invoking object.

A capability can be *anchored* to a thread. An anchored capability is always transient. It can not be made permanent via the *SaveCapa* operation. Therefore, it can only be used by the thread that brought it as a parameter.

The protection domain of an object is defined by the list of capabilities possessed by the object, including the current thread capabilities, at any point in time. In systems such as the Cambridge Distributed Computing System [Needham 82] strive to create unforgeable capabilities by utilizing a large, sparsely occupied name space. This approach does not guarantee the security of capabilities. Alpha makes it impossible to forge a capability by removing the means for directly manipulating them. The inability of the client to modify the *c_list* in a covert fashion enhances the utility of the protection domain as a security and reliability feature.

The only way an object can acquire a capability is for it to be passed to the object as an invocation parameter. Capabilities that are passed as parameters must be identified to the kernel so that the kernel can translate them from one object to the next. The kernel uses the index to locate the capability in the `c_list` and make a copy of it in the `c_list` of the target object. The index to the new capability is the value presented as an invocation parameter.

The kernel validates each capability before the invocation completes. The index provided must refer to a currently occupied location in the `c_list`. If the capability is being used to refer to the target object of the invocation then the target operation is checked against the operation restriction list.

Capabilities can have restrictions placed on them as they are passed as invocation parameters or used as the target of an invocation. The invoked object will receive a copy of the capability with restrictions applied. Restrictions consist of a list of operations that may not be performed on the object that is the target of the capability. Once a restriction is placed on a capability, it cannot be removed; all copies of a restricted capability retain the restrictions.

6.2 Data Structures

In a traditional capability based system capabilities are relatively simple bit encoding of a logical name. Alpha uses capabilities to provide network transparent naming and access control. As a consequence the data structures associated with an Alpha capability is more complicated.

The logical structure of a capability is that each object has a `c_list` and each slot in the list is occupied by a capability. Each capability contains the following information:

<capability_name> <anchored_flag> <operation_restriction_list>

The actual kernel data structures are:

```

Object Control Block
    ...
    c_list pointer
    ...
end
C_LIST
    ...
    attributes list pointer
    ...
end
Attributes List
    reference count
    anchored flag
    operation restriction list
    dictionary pointer
end
Dictionary
    ...
    capa_id
        reference count
        capability name (128 bits?)

```

```
                location hint
            end
        ...
    end
```

Since the `c_list` is a dynamic structure that can grow and shrink over time it is not actually allocated within the Object Control Block. Instead, the OCB contains a pointer to the `c_list`.

Each entry in the `c_list` points to the data structures that describe the capability. This is done because of the expense of duplicating the capability. The capability name is 32 bytes and the operation restriction list for a capability can be thousands of bytes long. It is possible for a object to acquire many instances of a capability to a single object. It is also possible for many objects on a node to own identical capabilities. In order to reduce the amount of duplicated data we split the attributes list from the capability name and share the data whenever possible.

The entries in the `c_list` point to Attributes Lists. More than one `c_list` entry can point to a single attribute list. Each attribute list has a count of the total number of `c_list` entries pointing to it. When this count reaches zero the attribute list is deleted. Each list also contains the anchored flag and the operation restriction list. The last field in the attribute list is a pointer to the `capa_id`.

The `capa_id` is the structure that most closely corresponds to the normal idea of a capability. The `capa_id` lives in the Dictionary. The Dictionary is an ordered list of `capa_ids` used by the Communication Subsystem. See section describing the Communication Subsystem for more information about the Dictionary and the Communications Subsystem.

The `capa_id` is the structure that contains the (128 bit?) identifier or logical name. It also contains a reference count and a location hint. The location count is used to determine when it is safe to delete the `capa_id` from the dictionary. There is a location hint field that has different uses depending on whether the capability refers to a local entity or a remote entity. If the entity is local the hint is the address of the control block associated with the entity.

If it is a remote entity the hint is the last known network address of the entity. This hint allows us to reduce the broadcast related burden on the Communications Subsystem. The Communications Subsystem has the responsibility for keeping the hints accurate. Whenever a nameable entity moves from one node to another all the hints in all the dictionaries are updated. If a hint is inaccurate then the target node does a broadcast which, if successful, will result in all the dictionaries being updated. If a capability belongs to a permanent object on a node that fails the hint can become stale. The recovery system will always set hints to NULL. The first use of the capability will cause a broadcast which, if successful, will initialize the hint.

There is a period of time while a object is being migrated that neither the source or the destination is a completely legal address for the object. The communications system treats the source address as the correct address until the migration is complete. Until the migration is complete threads become associated with the version of the object on the source node. Part of the migration involves moving those threads to the destination of the object migration. For more information see the section describing the Communications subsystem.

Note that all capabilities on a node live in its Dictionary, not just the capabilities of objects that live on the node. Since the kernel is not paged and capabilities are kernel data structures and the Dictionary is part of the global shared memory everything has to stay in primary storage anyway. The Dictionary is organized as a balanced tree which keeps the search time to a minimum despite the large size.

7 Semaphores

7.1 Description

The Alpha kernel provides two different concurrency control mechanisms—locks and semaphores. Locks are described in the chapter about objects. Semaphores are intended to solve synchronization problems not directly associated with data locking. Alpha semaphores are kernel objects that implement traditional *counting semaphores*.

7.1.1 Semaphore Manager

Although they are kernel objects, semaphores can be created or destroyed at any time. A special kernel object, the *Semaphore Manager*, is used to create new semaphores. This manager serves as a schema object for semaphores. The only operation defined for the semaphore manager is *CREATE*. This operation returns a capability pointing to the newly created semaphore. Semaphores are associated with object that created them. They can only be accessed by the parent object. They inherit the permanence attribute of the parent and they are migrated whenever the parent is migrated.

7.1.2 Semaphores

The semaphore consists of a count, which is the state of the semaphore, and a list that contains any threads that are blocked. *P* and *V* operations are provided and act as one would expect. Additional operations are *V_ALL*, *C_P*, and *DELETE*.

1. **P** —This is the standard *P* operation on a counting semaphore. It can be considered an attempt to acquire a resource. If the resource is unavailable, the thread blocks.

```

decrement the current_count
if (current_count >= 0) then
    the thread acquires the semaphore resource
else
    the thread blocks

```

2. **C_P** —The *C_P* (Conditional *P*) operation returns a failure without decrementing the *current_count* if it would have blocked. Otherwise it is just like a *P*.

```

if ((current_count - 1) >= 0) then
    decrement the current_count
    the thread acquires the semaphore resource
else
    return a indication of failure

```

3. **V** — This is the standard *V* operation on a counting semaphore. It can be considered an attempt to release a resource. If there are any threads blocked on the semaphore then one of them should be unblocked. The scheduling subsystem is requested to choose the correct thread to unblock. The scheduling subsystem has the responsibility for choosing FIFO or slack time or other criteria for choosing the thread to unblock. A *V* operation never causes the thread to block.

increment the current_count
 if (current_count < 0) then
 a thread is unblocked

Note that V operations are not constrained by the Initial_Count.

4. **V_ALL** — This operation unblocks any threads waiting for this semaphore and resets the current_count to the initial_count.

current_count = initial_count
unblock any waiting threads

5. **DELETE** — This operation unblocks any waiting threads and destroys the semaphore. Any semaphore data structures cleansed.

unblock any waiting threads
 if (permanent == TRUE) then
 free the secondary storage associated with semaphore
 disassociate the semaphore from the parent object
 free semaphore control block

The design of Alpha semaphores was influenced by the desire to provide a quick, predictable synchronization tool. Semaphores could be implemented as client objects but the kernel implementation provides a much faster service. Semaphores are associated with the object that created them. When the parent object is migrated any semaphores associated with it are migrated as well. This allows the invocation cost of using a semaphore to remain constant. The constant cost of invoking a semaphore is useful when designing the time constraints and execution estimates of an application.

In support of atomic transactions, all of an object's semaphores that a thread within a transaction has performed P operations on must have corresponding V operations issued on them should the transaction abort. This requires that the kernel be aware of all of the semaphores that a thread issued P operations on while within a particular atomic transaction. With this information it is possible for the kernel to issue the matching V operations should a transaction abort. This function is simplified in Alpha because all of the semaphores associated with an object are linked to the object's control structures and are therefore easily scanned to determine which ones require V operations to be performed on them.

7.2 Data Structures

Semaphore Control Block:

Initial Count
 Current Count
 Parent Object
 Blocked Threads

The *Initial Count* and the *Current Count* are set to a value specified by the *CREATE* operation. When the semaphore is reset via a *V_ALL* operation the *Initial Count* is used to set the *Current Count*. The *Current Count* is modified by *P*, *V* and *V_ALL* operations as described earlier in this chapter. The Parent Object is a pointer to the parent object. *Blocked Threads* is a list of all the threads blocked on this semaphore.

8 Object Replication

An optional attribute of objects is one that permits some objects to seem reliable, in the sense that they can complete the requested operation in the face of (variously) communication failures, node failures, and system software failures. This attribute is known as availability (i.e., the probability that the object will be available to provide a desired function). The kernel provides mechanisms to support the increased availability of objects through the use of a technique known as replication. When an object is created a specific replication policy is specified, as well as the degree to which the object should be replicated. Replicated objects provide a range of consistency and availability at a range of costs in terms of invocation latency.

In order to meet its goal of availability of services, the Alpha kernel provides support for the replication of objects. The support provided for object replication does not, in itself, constitute a complete data replication scheme. The kernel provides a framework for experimentation in the area of replicated object management. The approach to object replication taken in Alpha is one based on the use of multiple instances of a particular type of object, all of which share a common logical identifier, and consequently appear as replicas of the same object. The general issue of object placement is considered a higher-level issue in Alpha, and therefore the question of how the replicas are to be distributed among physical nodes is to be addressed at a level above the kernel.

As with atomic transactions, the area of object replication is currently being explored as a part of an ongoing thesis project. The mechanisms currently provided by the kernel are meant to be representative of a much more comprehensive set of mechanisms to be developed as a result of this work.

The Alpha kernel is designed to be a framework to support a wide range of replication mechanisms, among which are currently supported two forms of object replication inclusive and exclusive. In the inclusive form of replication the replicas of an object function together as a single object. Therefore, an operation invoked on an inclusively replicated object (in general) results in the operation being performed on all of the existing replicas. The intent of this mechanism is to increase the availability of an object through a redundancy technique similar to an available copies replication scheme [Goodman 83].

To reduce the cost associated with such a replicated invocation, a quorum method may be used [Herlihy 86]. Currently in Alpha, a simplified quorum technique is used, where a quorum is defined to be the minimum number of replicas that must be involved in an operation before an invocation may be completed. By associating a quorum with each operation defined on an object, fewer than the currently existing set of replicas can be involved in a particular operation. This mechanism can be used by the client to create objects with different degrees of availability, response time, and consistency.

A number of issues related to inclusive replication have been deferred in an effort to reduce the scope of this effort. For example, the timestamp or version number mechanisms needed to implement a proper quorum-based replication scheme have not been provided. Nor has the issue of the regeneration of failed replicas been addressed. These issues are being dealt with in a related thesis project.

For exclusive replication, the kernel also provides a number of replicated instances of an object type. In this case, however, an invoked operation need only be performed on any one of the replicas for the operation to be considered complete. This approach provides a form of replication in which any one of a pool of undifferentiated object replicas is chosen, based on a global policy designed to meet certain goals. Examples of replica selection policies are "select the first replica that responds", "select the replica that exists at the least loaded node", "select a replica near the invoking object", etc.

The initial policy used to select a replica chooses the first replica to respond to an invocation. This policy provides the potential for a higher degree of both availability and performance than is achievable with non-replicated objects. However, this is accomplished at the cost of not maintaining the consistency of data across the individual replicas.

8.1 Replicated Object Invocation

The various forms of replication provided in Alpha are supported primarily through features of the invocation facility. The replication support is provided in a manner that is transparent to the client a client performs operation invocations in the same manner regardless of whether the target object is replicated. The operation invocation facility determines whether or not the object being invoked is replicated, and carries out the appropriate operation invocation protocol for that type of object.

In the case of inclusive replication, the communication mechanism issues a message addressed to the replicated object that serves as the first phase of a two-phase invocation protocol. Each object replica that receives such a first-phase message responds either positively or negatively depending on its current state. The communication subsystem at the invoking node receives the responses and attempts to gather a minimum number of positive responses from the currently existing replicas. This minimum number of responses is specified for each operation when an inclusively replicated object is instantiated. If the initiating node does not obtain the specified minimum number of positive acknowledgments from the replicas, it issues an abort message and retries the first phase of the replicated invoke protocol. If positive acknowledgments are received from all of the replicas, the initiating communication subsystem carries out the second phase of the protocol by issuing an invocation message to all of the instances of the replicated object. In order to serialize the actions of threads within the replicas, all subsequent first phase invocations are negatively acknowledged until there are no more outstanding acknowledged invocations.

When a node's communication subsystem detects an incoming invocation from a replicated object, it waits until the invocations are received from all the necessary replicas before passing the (single) operation invocation signal to the application processor. The number of replicas involved in any replicated operation invocation is determined by the node initiating the invocation during the first phase of the protocol, and is passed to the replicas in the second phase. All invocations made by a replicated object include this count of the currently involved number of replicas that is used (in a fashion similar to that in [Cooper 84]) to merge the outgoing invocations of a replicated object.

In the case of exclusive replication, a similar type of two-phase invocation protocol is used. In the first phase the communication subsystem the node where the invocation is initiated issues a request message addressed to the replicated object. All of the current-

ly existing replicas respond to this first-phase message and the initiating node selects one of the replicas which responded in a given period of time. The second phase of this invocation protocol involves sending an invocation message to the selected replica. The reply from such an invocation is handled in the normal manner, and outgoing invocations are also handled the same as those from non-replicated objects.

The replica selection algorithm currently in use chooses the first replica to respond to the first-phase message. This portion of the operation invocation mechanism was designed to be easily modified or replaced, in order to permit experimentation with differing replica selection algorithms. In order to use more sophisticated selection algorithms, differing types of kernel information (e.g., load statistics or resource utilization estimates) may be required in the replicas' response messages. Thus, the kernel's design permits a node's communication subsystem to have access to the kernel's internal data structures where such information can be obtained.

If more than one replica of an exclusively replicated object exists at a node at some point in time, the replicas are chained together at the same location in the node's Dictionary. Therefore, when an invocation is made on an exclusively replicated object, the first instance of any local replicas that exist on a node is used.

With both of the currently provided forms of replication, the communication subsystems use logical addressing for messages directed to all replicas of a particular object. All logically addressed messages also contain the node's physical address, and so response messages from replicas are sent using the initiating node's physical address, along with the physical address of the replica's node. Subsequent messages from the invoking node (e.g., second-phase or retry messages) are sent to the physical addresses of specific nodes to reduce the amount of logically addressed traffic on the communications subnetwork.

8.2 Functional Design

The replication schemes currently supported by the kernel are based on a major underlying assumption that the partitioning of system resources (as it is popularly conceived of in the literature [Herlihy 85]) do not occur in the Alpha environment. This derives from the fact that systems in the real-time command and control application domain are typically provided with sufficient physical redundancy to ensure that the bisection of the distributed computer's interconnection subnetwork could only result from a catastrophic failure of the platform in which it is embedded. In effect, the communication interconnection subnetwork in Alpha is considered equivalent to a backplane in a conventional uniprocessor, and similar failure assumptions are made.

In Alpha, the client creates a replicated object by invoking the `CREATE` operation on a Schema object with the replicated attribute specified. This operation instantiates a specified number of instances of a given object type, each of which share a common object identifier. The multiple instances of a replicated object project to the client a view of a logical object that has different levels of availability, responsiveness, and consistency, than corresponding non-replicated objects.

Each of the individual instances that comprise a replicated object appear similar to all other objects, with the exception of the fact that they all participate (to some degree or other) in operations invoked on the replicated object. An object invoking an operation on

the replicated object does not need to be aware of the fact that the target object is replicated; invocations on replicated object appear identical to the invocation of a non-replicated object. Because the support for replication is provided by the invocation facility, functions involved in the management of replicated objects (e.g., quorums or multiple requests and responses) are performed within each node's communication subsystem and therefore do not incur significant performance penalties on the application.

While the basic mechanisms for replicated object support are provided by the kernel's invocation mechanism, the remainder of the replicated object management policies are performed by objects in the higher layers of the system. The replication mechanisms currently supported in Alpha are representative of a broader collection of replication mechanisms that support a greater range of replication policies.

9 Atomic Transactions

Most existing instances of the use of atomic transactions in operating systems evolved from the database application domain. In many such systems, atomic transactions are supported by migrating portions of a particular database system into the lower levels of the system, resulting in the propagation of a particular database model to the operating system's client. The emphasis on the design of transaction support in Alpha is directed towards more general-purpose atomic transactions. Also, the desire to make use of atomic transactions within the system layer increases the need for higher performance atomic transaction mechanisms.

The transaction approach used in Alpha was designed for flexibility and for real-time performance. To meet the flexibility goal, the transaction mechanisms in Alpha allow the client to determine when transactions are to be used and what characteristics each transaction should have. The real-time goal is addressed through a transaction management algorithm that permits bounds to be placed on the time between machine failure and subsequent transaction abort (since all affected transactions must abort).

An assumption made throughout this development is that the nodes behave in a fail-stop manner i.e., when a failure occurs in any of a node's processing elements, the entire node is considered to have failed and a restart procedure is initiated. This assumption is made in order to avoid having to deal with issues related to Byzantine protocols [Pease 80].

In the Alpha kernel, atomic transactions are supported in part by mechanisms that directly and specifically support atomic transactions, in part by special aspects of mechanisms that are not directly related to atomic transactions, and in part by the normal aspects of kernel mechanisms.

9.1 Transaction Management Mechanisms

The primary kernel mechanism supporting atomic transactions is the TransactionManager object. An atomic transaction is initiated by the invocation of a *BEGIN* operation on the TransactionManager object by the thread wanting to initiate a transaction. An atomic transaction can come to an end in a number of ways should all of the activities complete successfully, the thread that issued the *BEGIN* operation can issue the matching *END* operation; in the event that not all of the activities in the transaction are completed properly, an *ABORT* operation can be issued; alternatively, the kernel could detect that one (or more) of the objects (or nodes) involved in the transaction has failed, and the kernel will then issue an *ABORT* itself.

The TransactionManager object is responsible for performing the transaction coordination function for all of the atomic transactions in Alpha. The activity of transaction coordination involves managing the nesting of atomic transactions by individual threads, tracking all of the objects that a thread invokes operations on while within an atomic transaction, and coordinating the aborting or committing of transactions among all relevant parties. This implies that the TransactionManager object must maintain the information necessary to determine which objects are involved with each transaction. When an atomic transaction ends, the transaction manager must also ensure that all the objects involved in the transaction agree (within some period of time) to commit or to abort.

The transaction manager must also interact with other mechanisms in order to provide the various attributes of atomic transactions. For example, when a node fails, the kernel notifies the TransactionManager, which then becomes responsible for invoking the *ABORT* operation on all objects that were involved in the transaction. The kernel assigns unique identifiers to each transaction; an individual atomic transaction is identified by the identifier of the thread that initiated the transaction along with the current nesting depth of the transaction.

The TransactionManager object is replicated on each node, and the replica on a particular node is responsible for those transactions that involve the objects on that node. All transaction manager replicas interact with each other to commit particular transactions, using a standard two-phase commit protocol. The *BEGIN* operation increments the transaction nesting depth count associated with the invoking thread, and causes the invocation mechanism to begin tracking the invocations made by the thread. The *END* and *ABORT* operations cause the TransactionManager object to discard information about a particular atomic transaction and interact with all of its instances to ensure unanimity in committing or aborting the atomic transaction. The TransactionManager object invokes the transaction-related standard operations on the objects involved in a transaction on a commit or an abort. Furthermore, when an *ABORT* operation is invoked on the TransactionManager, control is returned to the statement in following the *END* operation for the transaction being aborted.

When transactions are aborted at the explicit command of the thread and there are no node failures, the TransactionManager object is able to use the information it has to determine which objects are involved in the transaction and notify all of them of the transaction abort. When a transaction abort is brought about by the failure of nodes, however, a portion of the transaction manager's information is lost. The particular node at which a transaction was initiated (and therefore the node whose TransactionManager object replica serves as the site coordinator for the transaction) may be lost in a failure. Thus, the notification of affected objects must be performed by a decentralized transaction coordination algorithm. The complete abortion of a transaction is necessary to satisfy serializability constraints and to allow failed computations to be "undone." Therefore it is important that the transaction management algorithm guarantee the complete and timely abortion of failed transactions.

To satisfy these requirements, the kernel's transaction abort management is based on the use of a technique known as a dead-man switch, whereby each node autonomously manages (by way of its TransactionManager object replica) the transactions that involve objects existing local to the node. In this approach, an abort time is assigned to each transaction as it is created. This time is never later than the current time by more than the abort interval, thus bounding the maximum time between a failure and the completion of abort. Periodically, the node at which the transaction was initiated executes a two-phase refresh protocol with all of the nodes on which objects involved with the transaction exist, assigning a new abort time to the transaction. If the refresh protocol does not complete successfully (i.e., some node cannot respond), all of the nodes involved in the transaction being refreshed are autonomously aborted at their given abort time. The abort interval may be adjusted by trading a shorter interval for higher communication overhead and processing overhead.

This design of the TransactionManager object has several benefits. The two most significant are that it bounds the abort interval, and that it correctly handles orphans. A proof of this algorithm appears in [McKendry 85], along with a more detailed discussion of this algorithm. This design also permits optimizations that may significantly reduce the required message traffic, although in its current form the algorithm requires message traffic proportional to the square of the number of nodes in the system. The actual amount of message traffic in any implementation can be chosen according to the desired abort interval.

9.2 General Mechanisms

In addition to the kernel mechanisms directly intended for the support of atomic transactions, a number of other kernel mechanisms provide support for threads executing within transactions. Certain kernel mechanisms behave differently (i.e., take on special characteristics) when the thread making use of them is inside an atomic transaction, while the normal characteristics of other mechanisms are used to support atomic transactions. Mechanisms that have special features (or behaviors) to support atomic transactions include the operation invocation mechanism and the thread synchronization mechanisms. The permanence and atomic update attributes of objects, the standard operations associated with client objects, and the thread repair feature of the operation invocation mechanism all contribute to support atomic transactions in Alpha.

10 Scheduling Subsystem

The scheduling subsystem provides the basic mechanisms that: recognize scheduling events, invoke scheduling policy, and perform thread switching. It is important to note that the scheduling subsystem does not implement policy—it does not choose what thread to run or unblock next. The ordering of threads is strictly the responsibility of a scheduling policy module that runs within the kernel. The scheduling subsystem invokes operations on the policy module in order to determine an ordering for thread execution.

The client uses *DEADLINE* and *MARK* operations on the current thread to specify the deadlines and importance information for regions of code. When a node experiences a computational overload, the deadlines of some threads may not be met. To minimize the waste of system resources, the system aborts some threads from the Ready Queue. The same mechanism is used for this as is used for aborting threads in atomic transactions or aborting broken threads.

Deadlines can be nested. The most deeply nested deadline inherits the characteristics of the most stringent time constraint. Information is maintained by the kernel to keep track of which deadline is currently active at any point in time. Deadline constructs are scoped within an operation invocation. If a *DEADLINE* operation is performed within an operation invocation without a *MARK*, the kernel will supply a matching *MARK* on return from the invocation.

10.1 Overview

In Alpha, the binding of threads to the application processor is performed when scheduling events occur (e.g., when it becomes more valuable to the system to run another thread or when the currently executing thread blocks). The scheduling subsystem provides the mechanisms for handling these events in an efficient manner.

The scheduling subsystem itself does not implement any object-level facilities, all entry points are internal to the kernel. Since the scheduling subsystem is entered run on every processor in a node independently, it must be multi-threaded and able to handle a high degree of concurrency with small loss in performance.

It is anticipated that many policy modules will require time-consuming calculations in order to select an optimal thread ordering. To facilitate this, the scheduling subsystem is designed to allow policy modules which use “background” processing for the time-consuming computations. In this mode, one of the application processors is selected to run the thread performing the scheduling work as needed. When no background scheduling is needed, the processor is available to run application threads.

10.2 Thread Queues

Each thread within a node can be in one of the following states: *RUNNING*, *READY*, *BLOCKED*, *STOPPED*, or *SLEEPING*. Threads that are in the *READY* or *RUNNING* states are linked together on the Ready Queue. Threads in the *BLOCKED* state waiting for an event are linked together on a Wait Queue with other threads blocked on the same event (e.g. a semaphore or lock). The Ready Queue and Wait Queues for each node are maintained and managed by the node's scheduling policy module.

0.2.1 Ready Queue

Threads are added to a node's Ready Queue when new threads (roots or surrogates) are created or blocked threads become ready, and threads are removed from the Ready Queue whenever they leave the READY or RUNNING state (i.e., when they are completed, stopped, blocked or deleted). When a thread is added to the Ready Queue, the thread's environment information is used by the scheduling subsystem to provide the necessary inputs to the scheduling policy algorithm. When changes are made to the environment of a thread that is currently in a node's Ready Queue, the modified information is used by the scheduling subsystem. Because the environment information moves with threads among nodes, scheduling decisions are made based on the global attributes of computations, as opposed to purely the local scheduling decisions made in less tightly-integrated systems.

The scheduling subsystem continually examines the information associated with the threads in the Ready Queue, and orders these threads according to the relative value to the system of their completion, as defined by the given scheduling policy. Threads on the Ready Queue are linked together by the scheduling policy module in decreasing current value to the system of their completion. The scheduling subsystem's dispatch mechanism assigns the first n threads on the Ready Queue to the n application processors on the node (the first $(n-1)$ threads if background scheduling is in progress). When it becomes more valuable to execute a thread other than the one currently bound to an application processor, the scheduling subsystem dispatcher preempts the currently executing thread and starts the new thread.

Whenever a thread is executing within the kernel context, it may request that it not be preempted for a short time. This is typically done while the thread holds critical system resources. The dispatcher will defer preemption of such threads until they have indicated that it is once again permissible to preempt them.

10.2.2 Wait Queues

Threads are added to a Wait Queue whenever a thread blocks on an event. There is one Wait Queue associated with each distinct event for which a thread can wait. The scheduling policy module is responsible for maintaining the order of each Wait Queue. This ordering is made by choosing the threads that, if added to the current Ready Queue, would maximize the total value to the system (as defined by the scheduling policy). When the event associated with the Wait Queue occurs, the scheduling subsystem selects one (or more) of the threads on the queue to awaken.

10.3 Policy Module

In Alpha, the time constraints of application activities are expressed as *time-value functions*, with the value to the system of completing each activity given as a function of its completion time (hard deadlines are a simple special case). In addition, activities have relative importance values which are also time-dependent. These time constraints and importances are dynamic and must be continually re-evaluated. Every evaluation is performed for all executing and pending activities collectively so as to maximize the total value to the system across the whole time period represented by the expected durations of all these activities. Time constraints and importance are among the attributes propagat-

ed with computations which cross node boundaries so that resource management can be global.

The design of the scheduling subsystem is such that a wide range of different scheduling policies can be specified with little effort. The policy in effect can be selected at system configuration time. Currently the kernel uses a deadline scheduling algorithm when the computational demands on an application processor can all be met. When the demands for computational resources cannot be met, the threads are scheduled in such a fashion as to maximize the value to the system of the computations that can be performed. The currently used overload handling heuristic is known as a best-effort approach, and is similar to the technique described in [Locke 86].

To support the operations of the scheduling policy module, the scheduling subsystem requires certain information about each thread in the Ready Queue. The information provided by each thread's environment currently includes an indication of the estimated amount of processing time that the thread needs across a specified interval of time, the critical time for the interval, and a time-value function defining the importance of completion. The scheduling subsystem monitors the amount of processing time that each thread has accumulated during the current interval, together with the thread's environment information, and generates a schedule for all of the threads in the Ready Queue.

10.3.1 Interface

The operations that the kernel will call on the scheduling policy module are:

- Initialize:** Provide information on the system resources (number of application processors, etc.).
- Stop:** Remove the current thread from the Ready Queue.
- Ready:** Add a thread to the Ready Queue (remove from Wait Queue).
- Block:** Add the current thread to a Wait Queue (remove from Ready Queue).
- Order:** Order threads in a Wait Queue for selection to run. A module is provided by the kernel which will be given a Wait Queue as ordered by the policy module and which should make ready those threads that meet the selection criteria.
- Release:** Return all threads from a Wait Queue to the Ready Queue.
- Modify:** Modify the current thread's scheduling environment.

The kernel operations that a scheduling policy module can call are:

- Background:** Run the policy module's background scheduler.
- Lock:** Lock a thread queue for modification.
- Unlock:** Unlock a locked thread queue.
- Time:** Get current time.
- Notify:** Request notification when a time interval has elapsed.
- Allocate:** Allocate memory for use by the policy module
- Deallocate:** Free allocated memory.

11 Memory Management

The memory management subsystem in Alpha is responsible for the managing the physical memory and system mapping resources of each Alpha node. The objectives of the memory management system are:

- *Support the memory requirements of client objects and threads.* Objects in Alpha are an encapsulation of data and the code that specifies the method for manipulating the data. The data is further divided into data which exists for the life of the object and whose size and initial values are known when the object is created, and dynamic data which is created and destroyed during the life of the object (the heap). The memory management subsystem must provide memory to hold the initialized data, the dynamic data, and the code. Similarly, threads consist a stack of invocation frames and a thread heap. The memory management subsystem provides memory to hold these structures.
- *Support thread contexts.* Each thread, as the unit of activity in Alpha, operates within an address space, or context, which determines the memory that can be accessed by the operations it performs. The memory management subsystem associates the current invocation frame and thread heap with the thread's context. In addition, the thread's context includes the code, data, and heap of the object in which it is running. Whenever a thread is selected for execution, the memory management system establishes its context within the application processor.
- *Provide protection between thread contexts to enhance fault containment.* Alpha objects represent separate domains which may only interact via invocations. The memory management system must assure that each object is safe from errant memory references from threads in other objects. In the same way, private data of a thread must be protected from access by other threads.
- *Provide mechanisms for efficient invocations.* Invocations are the fundamental means of interaction between objects in Alpha. The memory management subsystem must provide mechanisms that make parameter passing and operation invocation efficient. In particular, it must efficiently support a partial context swap, where only the object portions of a thread's context are changed.
- *Provide mechanisms for efficient memory usage.* Since the code of an object instance cannot be modified, the virtual memory system can create a single area of memory to be used as the code section of all instances of each object type. In addition, many objects may allocate dynamic data in excess of actual use. The virtual memory system may delay actual memory allocation until the data is first used. The use of these techniques can result in improved memory utilization.
- *Provide memory-related services for other kernel subsystems.* Many parts of the Alpha kernel need to allocate and deallocate memory in order to perform their assigned function. In addition, some subsystems (e.g., I/O) need information about memory mappings and the ability to prevent the mapping from changing while they are accessing the memory. The memory management subsystem will provide support for these activities.

- *Provide an extension to primary memory.* Contemporary machine architectures contain a limited size, volatile primary memory in which programs can be run and data can be directly manipulated. It is assumed that there will be Alpha applications large enough to preclude having all parts in primary memory at once. Therefore, the memory management system should, as an application-specific option, support the extension of primary memory into high capacity, low performance, storage.
- *Enhance Alpha portability.* In order to facilitate the use of Alpha Release 2 on wide variety of processor architectures, the memory management system must be designed to be as portable as practical.
- *Separate policies from mechanisms.* The policies which determine the allocation and use of memory should be separated from the mechanisms that are used to implement the policies. In that way, changes in system policy may be made without reimplementing the underlying mechanisms.

11.1 Design Overview

The fundamental entities managed by the Memory Management Subsystem (MMS) are extents and contexts. An extent is an independent parcel of contiguous virtual address space. All memory requirements of client objects and threads are satisfied via extents. Each thread defines a context (or address space) which specifies the set of extents that may be accessed by the thread at any point in time, and the addresses at which they appear. Extents that belong to an object are then "mapped" into a thread's context as it enters the object, and "unmapped" when it leaves.

The MMS is designed to enhance portability by maintaining a separation of machine independent and machine dependent processing. Moving Alpha to a new machine architecture should, therefore, only involve changes to machine dependent part of the MMS. The machine independent portion of MMS deals with extents exclusively, only the machine dependent layer must break up the concept of extent into component page tables, etc. depending on the underlying hardware structure. The architecture of the Alpha MMS is designed to support machines using either a linear address space model or a segmented address space model.

Alpha is designed to run on nodes which consist of tightly-coupled multiprocessors, each of whose processors are running threads concurrently. Therefore, the MMS is designed to support multi-threaded execution and be able to achieve a high degree of parallelism.

11.2 Extents

Extents are independent, contiguous areas of virtual address space with independent protection and secondary storage attributes. Client objects consist of a set of memory extents: the code extent, the data extent, and the heap extent. This group of extents is called the object region. Whenever the client object containing a extent is deleted or garbage collected, the constituent memory extents are also deleted. In a similar manner, a thread region is composed of memory extents for the stack and parameter data associated with each thread, and the thread heap. These extents are deleted whenever the thread is deleted.

Each memory extent is represented as a kernel object (a memory object) which manages a cache of volatile memory chunks that map parts of the extent (the resident set). Each memory extent object contains a paging object capability. The paging object (also referred to as the pager) determines the paging policy of the extent including: pre-paging, updating secondary storage, synchronization, etc.

The kernel uses a default paging object when it creates the memory extents that make up a thread or client object. The default paging object is designed to make use of the capabilities to two secondary storage objects, which are referred to as the initialization object, and the backing store object. The initialization object is an object that supplies the initial contents of pages for the memory extent. This object will not be modified by the MMS. The backing store object is used as a provider of writable paging store for the object. Either secondary storage object may be omitted, depending on the characteristics of the memory extent.

Application specific paging objects (called external pagers) will be supported in future releases of Alpha. These may be created to implement paging policies which are not supported by the default paging object. For example, a special pager could be created to support a different transaction model within an Alpha application. Memory extents obey the same interface as the storage objects used by the default paging object. This can be used, for example, to implement a copy-on-reference sharing of pages by using one memory extent as the initialization object of other extents.

11.2.1 Extent Manager Operations

- The extent manager performs operations relevant to the creation and deletion of extents. These operations are available only within the kernel. The operations defined on the extent-type are:

Create: This operation creates a new extent and associates with it a paging object. Parameters are: a set of initial attributes for the extent, a capability to the paging object and a place-holder for the returned capability to the new extent.

The attributes which can be specified are:

Extent size.

Permissions on the extent.

Whether parts the extent can be paged out when not in use.

Whether the whole extent can be swapped out when not in use.

Delete: This operation deletes an extent. The capability for the extent to be deleted is passed as a parameter. Provided the capability is valid and contains deletion attributes, the extent is deleted.

11.2.2 Extent Operations

The operations defined on a extent instance which are used internally by the kernel in manipulating address spaces are shown below:

Migrate: This operation is used to move the extent among the nodes in the system. The operation takes a parameter that identifies the node to which the object should be moved.

- Modify:** This operation modifies the attributes of the extent. Parameters are: a set of new attributes for the extent.
- Attach:** Map the extent into an address region. Parameters are a region identifier and a requested map address.
- Detach:** Unmap the extent from the region. Parameter is the region identifier.

11.3 Pager Interface

The responsibility for managing the virtual memory of each memory extent is shared by the memory object for the extent, which represents the kernel, and the paging object (pager). The memory object encapsulates the machine dependent and privileged operations required to manage virtual memory. The paging object implements policies on the memory extent, and can control nearly all aspects of memory management for the extent. Unless a policy of pre-loading memory extents is implemented by the pager, the memory object will only request data in response to a page fault.

The Alpha system provides a default pager which it uses for system created memory extents. It is envisioned that additional pagers will be developed as system objects for special needs. The interface to such external pagers is outlined in this section. The operations on memory objects available to a paging object are:

- LoadData:** This operation informs the memory object that the data in the specified range of addresses is to be added to the resident set of the extent. This provides a method to add memory to the extent which has not been requested by the memory object due to a page fault. The memory object will then allocate physical memory and request the data for the pages from the pager. The parameters are: the offset into the object, and the length.
- Restrict:** This operation changes the access permissions and of any resident set pages within the specified region, and whether they may be paged out. The parameters are: the offset into the object, the length of the region, and the access permissions.
- Clean:** This operation forces any resident set pages within the specified region that have been modified to be written out. The parameters are: the offset into the object and the length of the region.
- Invalidate:** This operation forces any resident set pages within the specified region to be invalidated. Any pages that have been modified will be written out. The parameters are: the offset into the object and the length of the region.

The operations on a paging object which are invoked from a memory object are:

- Initialize:** This operation informs the paging object that it will control paging on the specified memory object. The parameters are: the memory object, the secondary storage objects associated with the memory object, the offsets into the secondary storage objects, and the length of the memory object.
- Fill:** This operation requests that the paging object provide the data to fill the specified resident set pages. The parameters are: the offset into the memory object, the length of the region, the place-holder for the data, and requested

access type. If the *Fill* operation on the pager fails, the pages are filled with zeros.

Flush: This operation requests that the paging object write the region to backing store to permit the associated page or pages to be removed from the resident set. The parameters are: the offset into the memory object, the length of the region, and the contents of the region.

Permit: This operation requests that the paging object change the access permissions on the region to permit the specified access or page-out. The parameters are: the secondary storage objects associated with the memory object, the offset into the memory object, the length of the region, and requested access type.

11.4 Secondary Storage Object Interface

The default pager expects a certain set of operations from the secondary storage objects which are used as initialization and backing store objects. Any objects supplying these operations can be used as secondary storage objects with the default pager. External pagers which use this interface will be interoperable with existing secondary storage objects. The operations on a secondary storage object which are invoked from the default pager are:

Read: This operation requests that the paging object provide the data to fill the specified resident set pages. The parameters are: the offset into the memory object, the length of the region, and a place to return the data.

Write: This operation requests that the paging object write the region to backing store to permit the associated page or pages to be removed from the resident set. The parameters are: the offset into the memory object, the length of the region, and the data.

Reserve: This operation requests that the secondary storage object preallocate enough resources to guarantee that the entire memory extent may be accommodated at a later time. The parameter is: the maximum length of the extent.

11.5 Page Fault Handling

When a page fault occurs, the kernel determines which, if any, extent the page fault address maps to, and finds the associated memory object. Processing of the page fault proceeds according to the following rules:

- If the fault address was not within an extent or is an access that is not permitted on the extent, an error exception is generated.
- If the fault was due to a reference to a page that was not in the resident set of the extent, a physical page is allocated and a *Fill* operation is invoked on the pager. If it can, the pager supplies the data for the page; otherwise, the page is filled with zeros. Then the page is added to the resident set, and the faulting operation is resumed.
- If the fault was due to an access that is permitted on the extent but has been temporarily restricted by the pager, a *Permit* operation is invoked on the pager. The faulting operation is then resumed.

11.6 Page Reclamation

The set of memory extents within a system are divided between those that can be swapped out when not in use (swappable), those whose unused pages can be removed (pageable), and those that are locked in memory. Extents that are locked in memory allow for more predictable execution times of the threads using them (because missing page faults cannot happen), but may reduce the system's ability to meet future memory demands. Whenever the amount of free memory falls below a system defined low water mark, the kernel will attempt to reclaim pages by consulting a memory policy module. To support the memory policy module, the kernel uses the concept of working set to determine which pages are in use and which are not. The memory policy module can use the working set information in determining which pages to reclaim.

The "swappable" and "pageable" attributes of extents are specified as a range from zero to MAX (where MAX is a system defined constant). Zero indicates that the extent is not swappable or pageable. Nonzero values indicate that the extent can be paged or swapped. The memory policy object may use the relative nonzero attribute values of extents to determine the order in which their pages should be reclaimed.

11.6.1 Working Set

The working set model of program behavior relies on the locality of reference of programs. That is, most programs repeatedly use the same code and data for extended periods of time, and this can be used to predict the memory requirements of the program. The working set of a memory extent is the set of pages which have been accessed or modified in the previous q seconds of run time of the threads using the extent (q is determined by the machine architecture and speed). Accurate approximations of the working set of a extent can be efficiently calculated via page reference data supplied by the Memory Management Unit (MMU). The calculation of a working set is independent of other activities within the system and, therefore, represents the intrinsic memory requirements of the computation.

11.6.2 Memory Overcommit

When memory is required by a thread, and none is free or can be supplied by removing pages not in any working set, the memory policy module resolves the conflict. It may select a extent to force to secondary storage (swap out), and thereby freeze all threads actively using the extent. It is anticipated that the memory policy object might, for example, use the importance and deadline information for threads using the extent, the pageability attribute of the extent, and the size of the extent to determine which, if any, extent to page out. If the memory policy object does not free any memory, the faulting thread will be suspended until memory is available.

11.6.3 Memory Policy Module

The memory policy module in Alpha Release 2 will follow the following strategy for reclaiming memory: extents that are not being used by any active thread and that may be swapped out have their memory reclaimed. If enough memory cannot be obtained by swapping unused extents, pageable extents are examined for pages that are not currently in the working set, and these are reclaimed. Each page that may be reclaimed as free memory is examined by the kernel to determine whether it had been modified while in

use. If so, the pager for the extent will be invoked with a *Flush* operation to save the modified contents of the page to backing store. If the page is later referenced by a thread, the page will be reloaded by the page fault handler.

The memory policy module uses the swappability and pageability attribute values to order the search through the memory extents for pages to reclaim. Those extents with higher attribute values will have their pages reclaimed before those with lower values. This can be used to control the handling of memory overload conditions.

11.7 Machine Independent MMS

11.7.1 Address Maps

Each thread executing within Alpha is given an address space, or thread context, which is distinct from that of any other thread. As the thread enters and leaves objects, its context is modified to unmap the memory extents associated with the old object and map in the extents of the new object. Each context is specified as a list of address map entries, which relate a range of addresses in the context to a memory extent (as represented by a memory object). The address map entry also specifies a set of permission attributes on the address range. The address map entries are used at page fault time to determine in which object the fault occurred. A context can not have address map entries that overlap address ranges.

11.7.2 Resident Pages

The kernel maintains a page frame data structure that describes the contents of each page of physical memory. Working set information, page modification indications, and access restrictions are maintained within page frame entries. Each page frame can belong to at most one memory object and a memory allocation queue. All page frames belonging to a memory object are linked together in a resident set list which is ordered by offset in the object. In addition, each page is linked on a hash queue which provides efficient lookup of a given offset within an object.

The resident set is the set of pages that have been loaded by software for an object. It is possible that, in some architectures, the pages may be removed by hardware without software intervention. Therefore, the machine dependent MMS will be consulted before assuming that a resident page is actually there.

11.7.3 Memory Objects

Within the kernel, a memory extent object consists of a data structure which contains the attributes of the extent, the number of address map entries pointing to the object, a pointer to its paging object, and a pointer to a list of page frames in the resident set of the object. This data structure is the focus of kernel activities in managing memory allocation and deletion.

11.8 Machine Dependent MMS

Alpha Release 2 is intended to run with limited changes on machine architectures with vastly different models of memory management. These include, but are not limited to, machines using page table mapping and segment mapping. We chose not to impose a memory management architecture that is built around one form to the detriment of oth-

ers. Therefore, the machine dependent interface of Alpha is built at a higher level than that of some other contemporary operating systems. It is envisioned that the machine dependent part could be further broken into an architectural section, which is common for all systems using the same memory management model (e.g., page tables), and another that encapsulates specific information about the hardware.

Because we cannot assume that any page (region of addressible physical memory) can be mapped to any virtual address (due to hardware restrictions), the machine dependent MMS code is responsible for actual physical memory allocation and assignment. During the allocation process, there can be cases where, because of hardware restrictions, another page must be removed to make room. In such cases, the memory policy module will be consulted (with a list of candidates) to make the selection of a page to remove.

The machine dependent MMS is responsible for fielding page faults, determining the fault type, and passing the fault address and type on to the machine independent MMS. It also maintains machine specific data structures which describe the memory allocated to each context. These are used to validate and invalidate context mapping hardware when needed during invocations and thread dispatching. The machine dependent MMS maintains contexts in terms of memory map data structures (*mmaps*). The contents of *mmaps* are not visible outside the machine dependent MMS; they are only used as handles by which the rest of the kernel communicates with the machine dependent MMS. The machine dependent MMS can use *mmaps* to store any information it needs to map contexts in and out of the hardware. The machine dependent MMS can access the data structures of the machine independent part via a pointer to the thread context structure.

12 Communications

12.1 Overview

The communications system provides three levels of service. The first and most important is the support for inter-node kernel communications. This provides for inter-node thread movement, thread maintenance, paging, and node monitoring. The second level provides support for system service communication. It is at this level that remote file system access is provided. The third level is client communications with external systems.

The functionality provided by the communications facility may be quite involved, requiring non-trivial amounts of processing power and the multiple exchanges of packets. For example, the communication facility handles all of the low-level aspects of the operation invocation function. This includes the handling and generation of positive and negative acknowledgments, communication time-outs, and packet retransmissions, in addition to such typical communication functions as message disassembly/reassembly, encapsulation, and page alignment of message data. The communication facility is responsible for executing an inter-node liveness protocol. This protocol requires that each node monitor the transmissions of the other nodes, and send explicit queries to nodes that have not been heard from in some period of time. This is done in order to determine which nodes are working and which are not, and to provide timely support to such invocation functions as thread repair and visit notification. Also, the communication facility includes mechanisms in support of specialized protocols for such functions as logical clock synchronization, two-phase commit protocols, the bidding protocols used to deal with a replicated set of objects, initial program loading, and node reboot functions. In addition to these protocols, the communication facility contains a mechanism that permits a remote node to gain access to the application processor and local memory in a node to perform remote test and diagnostic procedures. This mechanism also supports the remote loading and storing of a node's primary memory for restart, fault location, and debugging.

It has been repeatedly shown that a large proportion of the cost associated with message passing systems is related to the movement of messages among address spaces. To deal with this problem, the communications subsystem and the rest of the kernel interact through a shared virtual memory interface, sharing memory by the manipulation of virtual memory mappings. This provides a significant performance advantage over the copying of blocks of memory, as is typically done in message-passing systems. The communications subsystem moves information to and from the communication subnetwork without copying it among intermediate-level buffers. This requires that the communication subnetwork interface hardware deposit the information it receives directly into the physical memory location where it will ultimately reside, and remove the information it must transmit directly from the location in which it was placed by the applications processor. If this is the case, then all other movement of the communicated data are performed via memory mapping operations.

The location-transparent access to objects provided by the operation invocation facility is based on the use of global object identifiers and is supported by the communication facility through the use of logical destination addressing. For each operation invocation, the kernel performs a lookup operation on the node's Dictionary to determine whether the

destination object is local, and if not, the invocation request is passed to the communication facility. The communication facility addresses all remote communications to destination objects, using their global identifiers. The communication facility's address recognition mechanism uses the local Dictionary to determine which objects are local to the node, and therefore which (logically addressed) packets should be received. This mechanism supports object-level dynamic reconfiguration by simplifying the location of objects i.e., the maintenance of logical to physical address translation tables is not necessary. In addition, the use of logical addressing also supports the use of replicated objects, by having more than one object sharing a logical name.

12.2 Kernel Communications

The Alpha kernel communications system must support six basic services in order for Alpha to be distributed transparently and uniformly across multiple nodes connected by a network or communication bus. The most basic of these is reliable message communication. Next is the remote invocation service. The other three services which depend on remote invocation are thread maintenance and repair, monitor active node status, and page transfer. The last service which is crucial to the operation of remote invocation is the node dictionary. The major differences between Release 1 and Release 2 in the communications area are that it is not required that there be a dedicated communications processor, and that Alpha schedulable threads exist at the lowest possible level in the communications system.

12.2.1 Reliable Message Communication

Reliable message communication is the substrate upon which all the kernel communication protocols are built. Alpha uses the eXpress Transfer Protocol (XTP) protocol [Chesson 88] to provide reliable message transfer. XTP is a lightweight transport protocol with unified internetwork services conforming to OSI Layers 3 and 4. It is designed for implementation in a VLSI architecture called the Protocol Engine (PE).

There are four main goals of an XTP/PE system: high performance, support for real-time, integrability, and full functionality. The performance goal is to have the ability to complete all protocol and related processing for a received packet within the arrival time of the packet. This should have the effect of being able to process a stream of data at the media rate. The initial goal is to support the 100 Mbit/sec sustained transfer rate anticipated by FDDI. The design is intended to be scalable. It should be able to eventually support a Gbit/sec transfer rate.

Support for real-time includes a SORT field and the ability for the list of incoming and outgoing packets to be ordered dynamically. The SORT field is used to indicate the relative importance of each packet. Real-time policies such as those supported by Alpha may use the SORT field to rearrange the order of packets as they are being processed.

The protocol is kept simple and elegant. This permits the implementation in a small number of VLSI packages. The core of XTP is a minimal mechanism, or lightweight transport. It is this simplicity that permits the successful VLSI implementation, thus satisfying the goal of integrability.

Despite the simplicity of the design XTP is a robust, fully functional protocol. It provides:

- traditional stream services
- bulk transfer
- real-time reliable datagram service
- real-time internet gateways
- flow/error/rate control
- message delivery confirmation
- selective retransmission
- message boundary preservation
- multiple addressing plans
- out-of-band signalling
- reliable multicast mechanism
- maintenance packets
- multi-path capability

12.2.2 Remote Invocation and Node Dictionary

At the invocation point, an executing thread enters the kernel. The first parameter of the invocation (after the operation specifier) is the capability for the destination object. At this point the thread executing within the kernel gets the kernel data structure of the object represented by the target capability. The kernel then uses the information in that structure to complete the invocation.

In the case that the capability was for an object local to the current node, the thread then continues executing in the target object operation. However, if the object is on another node, the thread continues execution in the communication system. The remote invoke operation has all the information available to it that the kernel had at the original point of invocation. It has the operation specification, the capability for the destination object, and the invocation parameters. Each of the passed capabilities are registered in the node dictionary if they are local, and are removed from the dictionary if they are remote and are no longer referenced on the local node. A message is constructed out of the destination object capability, the operation, the parameters, and the scheduling information, which is then sent via the reliable message service available to the communication system.

When a network message is received, the scheduling information is retrieved from the message and a local uncommitted thread is acquired. Then the protocol is determined and the thread is started on the "MessageArrived" operation of the appropriate protocol module. In the case of remote invocation, the message arrived operation performs several steps:

- It looks up all the capabilities referenced in the invocation in the dictionary. If the target capability is not on the local node, the invocation is ignored and the thread resources are freed.
- For each look up of a non-target capability that fails, it asks the kernel to create local information structures (Object Control Blocks) for the remote capabilities, and enters them into the dictionary and flags them as remote.

- Using the operation, target capabilities being currently active. If a node has not sent any broadcast traffic for some period of time, then it can then announce explicitly that it is still active. If a thread on some node is "nervous" about the activity on another node, it can explicitly send a request for an activity announcement. Since the reply to this request is broadcast, it should have the side effect of "calming" any other currently "nervous" threads on other nodes.

12.3 System Communications

The protocols that Alpha provides for external communications are here designated the system communications facilities. At a minimum these will consist of TCP/IP, UDP, VMTP and various "raw" network, bus interface, and serial devices. These facilities are provided at system level either for access to kernel services not available to clients, or for efficiency.

12.3.1 Byte Stream

The byte stream service is provided by TCP/IP. This protocol suite is pervasive throughout the scientific and government research communities and thus has high availability and potential for interoperability. Byte stream service is generally assumed to be a reliable virtual end to end circuit. Byte streams rendezvous by means of well known sockets and further communication over other circuits is negotiated via some initial connection protocol at the layer above the byte stream.

12.3.2 Reliable Datagram

The current choice for reliable datagram service (as possibly distinct from the service used by the kernel) is VMTP. Reliable datagram service provide an efficient means to implement process to process asynchronous message communication, remote procedure call, or as in Alpha, operation invocation on remote objects. As a service in Alpha VMTP provides a means to use services of and provide services for other systems that use VMTP as the reliable datagram protocol.

12.3.3 Unreliable Datagram

Unreliable datagram service basically provides process to process (or object to object) addressing for messages. There is no guarantee that the message is delivered. Though it very easy to develop an unreliable datagram protocol, the most widespread protocol within the same community that uses TCP/IP and VMTP is the UDP protocol.

12.4 Client Communications

As much as the user may be able to accomplish within the Alpha environment, at some point the user invariably wants to communicate with some external system. Although many of these needs may be anticipated at system level, in some cases it will still be necessary to give the user explicit communication access to external systems. For purposes of compatibility with other systems, Alpha provides the "socket" system object. This object provides the interface for the various types of communications services such as "byte stream", "reliable datagram", and "unreliable datagram." These services are implemented by TCP/IP, UDP and the other protocols mentioned above.

12.4.1 Client Communications Interface (Sockets)

The socket object as defined here was inspired by that of Berkeley UNIX. The following set of operations are ostensibly available for any socket, but certain operations are meaningful only for certain types of service.

The operations meaningful for all types of service are:

- **CreateSocket** — This operation is invoked on a protocol object to get a socket on that protocol. Subsequent operations are performed on the capability returned by this operation.
- **Shutdown** — Gracefully shut down the protocol with respect to the socket and free the socket.

The operations that are relevant to Byte Stream protocols are:

- **Listen** — Places the socket in a passive state waiting for a remote connect on an address.
- **Connect** — Actively attempts to connect to a remote address.
- **Accept** — Notification to the socket that a connection has been accepted.
- **Read** — Read a specified number of bytes on a socket.
- **Write** — Write a specified number of bytes on a socket.

The operations that can be performed on a datagram (or message) protocol are:

- **Send** — Send a message to an address.
- **Receive** — Return the next message for an address.

13 Input/Output

This chapter discusses the Generic I/O paradigm and the strategy for designing device drivers.

13.1 Generic I/O

The Generic I/O model is simple and powerful. It is adequate to describe the behavior of a wide range of I/O devices ranging from UARTs to framebuffers to disk drives. Providing a model for I/O as a specification of a set of interfaces is consistent with the Alpha philosophy of hierarchical design and data abstraction. This layered approach allows us to isolate hardware specific behaviors behind a consistent, clearly specified interface.

Device drivers that implement the Generic I/O interface will have the useful characteristic of device independence. The usefulness of interposable modules in general and device drivers in particular is reported in the literature [Cox 86], [Ritchie 74], [Cheriton 87a].

13.1.1 The Model

The model for the Generic I/O is an object encapsulating a randomly accessible array of bytes. Buffers of any length can be copied or modified. A failure occurs if a request is made for an address beyond the bounds of the encapsulated byte array. The address to be copied or modified is a device relative address. Address zero is the first addressible byte managed by the device. Devices that are essentially serial, such as a UART, disregard the address and process the IN or OUT buffer as a byte stream. Devices may have a preferred buffer size. This size can be determined via the *info* operation.

Devices also have a STORAGE attribute. This indicates whether you can read back the data you have written to device. When you write to a disk, tape or framebuffer you expect to be able to read the same data at a later time. Network controllers can be read and written but there is not read to believe that what you read is what you wrote.

Information about buffer size, device memory, serial behavior, etc. can be determined via the INFO operation.

The Generic I/O paradigm is simple but powerful enough to model serial (UART, FDDI) and random access (disks, tapes, frame buffers) devices. The object oriented nature of Generic I/O makes implementing other mechanisms for device independent I/O such as UNIX streams [Ritchie 87], V UIO [Cheriton 87b] and sockets [Leffler 87] as a higher level abstraction relatively straightforward.

The GIO model does not explicitly deal with time-outs. There are no time-out parameters in any GIO operation. The functionality of time-outs is provided by normal Alpha time constraints. A GIO device can deal with operation time-outs implicitly or explicitly.

Implicit time out processing is handled local to a particular operation. A time out (deadline) exception block is defined within the operation. When a time out occurs the deadline exception handling routine is invoked. When it completes, control is passed to the outer deadline exception block. The inner deadline block is not part of the GIO interface. It is part of the data hiding that is part and parcel of abstract data types (or objects).

Explicit time out processing is accomplished by providing a operation on the GIO device that can be invoked in the case of time out. Once the outer deadline block is invoked it can in turn invoke a time-out operation on the GIO device object.

13.1.2 GIO Device Interface

GIO device drivers are implemented as standard Alpha objects. Operations are invoked on these objects just like any other object. The buffers used by input and output are located in the thread heap to minimize data copying.

Generic I/O operations:

- *init*
- *attributes*
- *acquire*
- *release*
- *input*
- *output*

The *init* operation has the responsibility for initializing the driver and the device. It also has the responsibility for registering the device object and its interrupt handlers with the kernel. This allows the kernel can notify the driver when certain exceptions, such as power failure, occur. Usually this is the first operation invoked on the object.

The *attributes* operation returns a list of parameters that describe the device. The purpose of the *attributes* operation is provide enough information to permit the creation of device independent I/O abstractions. The behavior of various devices is extremely varied. An application that needs to be device independent can use the *info* operation to find out: the amount of concurrency supported by a device object, maximum and minimum size of I/O operations supports, the type of device and the *acquire* status.

GIO device objects have advisory *acquire* and *release* operations defined on them that act like locks. Implementors of device objects use these operations to control the level of concurrency within the device object. Some devices, like disk drives, allow multiple readers and writers. Others, like most printers, allow only a single writer and no readers. The correctness of the device object is not guaranteed if there are clients who fail to follow the lock protocols.

Reading and writing data to the device is accomplished by the input and output operations. The arguments to these operations include the buffers to be used for the I/O operation.

13.2 Alpha I/O Subsystem

The Alpha I/O Subsystem is a collection of kernel level managers and GIO style device drivers. The managers include the Exception Dispatcher, the GIO Device Manager, the Kernel Semaphore Manager and the System Event Manager. Other managers, such as the VM manager provide special services for the I/O subsystem. GIO devices that are currently supported include disk, tape, UART, and FDDI.

13.2.1 GIO Device Driver

An Alpha device driver is a system object that encapsulates the data structures and text that the kernel uses to control a device. A GIO device driver is a driver that is consistent with the GIO model and interface described in an earlier section. A GIO device driver consists of several important parts: the model conforming interface, the *init* operation, the interrupt handler(s), the main body of the driver, and the System Event operations. Usually there is a physical device that is managed by the driver but it is not necessary. The GIO model supports virtual or pseudo devices as well as "real" devices.

Using objects and threads to build device drivers directly addresses most of our design goals. Since threads are schedule-able entities the real-time characteristics of Alpha can be preserved. Objects are accessed via capabilities which provide network location independent naming. Data is moved to and from device objects via invocations. This design utilizes existing kernel abstractions. No fundamental abstractions needed to be invented.

The kernel does not make any assumptions about which devices exist or how to talk to them. All of that information should be encapsulated within the device object. The kernel does not make any assumptions about the relationship between a particular device object and the number of physical devices it controls. The management of the various control structures, queues, device control blocks, etc. is left to the individual device objects.

The GIO devices are:

- Disk Drive
- Ethernet
- FDDI
- UART
- Framebuffer

DOs are instanced and initialized when the kernel initialized, just like other system objects. This is the time that data structures local to the DO and the physical device are initialized. In addition, the DO needs to register any interrupt handlers with the DMO.

Multiple threads may execute within a device object simultaneously. The kernel provided synchronization primitives, locks and semaphores are available to device objects. This allows the designer of the device object to take advantage of any concurrency the device supports.

A device driver has two different execution states, kernel and interrupt. The kernel state is identical to any other kernel object with a thread executing in it. A driver executing in this state can on any available service. Execution progresses from one operation to another via normal kernel invocations.

Interrupt state occurs as a direct consequence of a interrupt, trap or exception. When an interrupt occurs the system gives control to the Exception Dispatcher. Lower level interrupts and normal thread execution are temporarily suspended. The Exception Dispatcher passes control to the correct interrupt handler. An interrupt handler can be an operation specified by an object when it is registered as a GIO object. Execution in interrupt state is significantly different than any other state. There is no active thread. Most kernel services are unavailable. Services available in interrupt state are CREATE

(thread), C_P (semaphore and kernel semaphore), V (semaphore and kernel semaphore), V_ALL (semaphore and kernel semaphore), C_LOCK (lock an object), and UNLOCK (unlock an object). An object should spend as little time in interrupt as possible. Time spent in an interrupt handler is time unavailable to the scheduler.

13.2.2 Kernel Semaphore Manager

The Alpha kernel is multi-threaded. It supports preemption and multiple processors. In the interest of maximizing concurrency within the kernel we have created a two level locking scheme. The upper level locks are kernel semaphores. They similar to the semaphores available to client objects in that there are P, C_P, V, and V_ALL operations. These semaphores are useful in coordinating multiple threads on a single node regardless of whether the node is multi-processor or not. The lower level semaphores have the same operations but they have the additional effect of blocking and unblocking interrupts. Successful P operations block interrupts and successful V operations unblock interrupts. These lower level semaphores are spin locks when invoked by interrupt handlers.

Threads that attempt P operations spin with interrupts enabled. These are low level semaphores that are used to synchronize with interrupt handlers and the semaphores should not be held very long. The cost of a context switch should greatly outweigh the cost of spinning. If a thread is spinning on a low level kernel semaphore it means that another CPU has the lock. We spin with interrupts enabled so that interrupts local to our CPU can get the lock first. This has the advantage of reducing interrupt latency.

Another unique characteristic of these kernel semaphores is that they are ordered. These semaphores must be acquired in order. Threads carry with them the list of kernel semaphores they have acquired. An attempt to get a semaphore out of order results in a failure.

13.2.3 GIO Device Manager

The Device Manager has the responsibility for providing the services necessary to create and manage device drivers as Alpha objects. When a GIO device is initialized, usually via the *init* operation, it makes itself known to the kernel by registering with the Device Manager. The GIO device sends the Device Manager an object SELF capability, a list of interrupts and interrupt handlers, and a list of notification events. The SELF capability allows the Device Manager the rights to inspect the GIO object via the *info* operation. The Device Manager verifies that the object referred to by the capability can be used as a GIO device object. If the Device Manager encounters any problems the registration operation fails and an error indicating the problem is returned.

Once the device object has been verified the Device Manager attempts to register with the Exception Dispatcher the list of interrupts and interrupt handler addresses. If this is successful then the Device Manager attempts to register the GIO device object and its list of notification events and invocation parameters with the System Event Notifier. If all of this completes successfully then the kernel will treat the object as a GIO device object.

13.2.4 System Events

There is a class of event that is asynchronous and affects all objects on a node. System initialization, power-failure, and node reset are examples of this type of event. System events are detected by the kernel and can affect threads and objects. Normally exceptions are associated with threads. In cases such as power failure, there might be a hard deadline and a great deal of value to the system associated with dealing with the event.

For each type of system event there is a list of object/operation/time value function triples. This is the Notification List. When the kernel detects a system event the list associated with that type of event is scanned. A kernel thread is generated for each entry. The kernel uses the list to provide the target object, operation and time constraint.

Client objects may be migrated from node to node. When an object that has registered for notification of a system event is migrated its notification information follows it. All of its entries in various Notification Lists on the original node are duplicated on the new node. Once the object has been successfully migrated the original Notification entries are removed.

13.2.5 Special Services

Device drivers make special demands on an operating system. Wire pages, Physical addressing, kernel to client address mapping, interrupt masking, spin locks, etc.

13.2.6 Exception Dispatch

All interrupts, traps and exceptions that occur on an Alpha node are processed by the Exception Dispatcher. The Dispatcher has the responsibility for identifying the type and the source of the event. Once this is done the Dispatcher invokes the correct routine to handle the event.

The Dispatcher uses the interrupt dispatch table to match interrupts to their handlers. When a GIO device is registered with the kernel an entry in the dispatch table is made. The table consists of a list of interrupt/address pairs. The interrupt is the name or number of the interrupt. Interrupts are ordered from one to the maximum number of different interrupts supported on a particular hardware platform.

The address is the address of the interrupt handler routine. This routine may be an operation on a GIO device object. All GIO device objects are system objects and therefore exist within the kernel address space. Someday we may support GIO device objects as client objects. Banging the VM hardware to service an interrupt is very expensive so we will still need to be able to map the interrupt handlers into the kernel address space.

The Exception Dispatcher works in conjunction with the scheduling subsystem to preserve context consistency and to reduce unnecessary context switches. Exceptions do not ordinarily involve a context switch. The exception causes the kernel to enter interrupt state and after the exception has been handled the initial context is restored. This is a very quick and simple operation involving the PC, the PSW, the stack pointer, and a few registers. When the Exception Dispatcher is ready to resume the interrupted context it checks to see if it is still at the head of the run list. If the original context has blocked, perhaps it page faulted, or if another thread has become more important to run, then the current context is saved and the DISPATCH operation on the Scheduling Subsystem is called. This operation will run the most important runnable thread.

14 Secondary Storage

This section discusses the design and facilities of the Alpha Secondary Storage Subsystem. The secondary storage subsystem exists to provide support for object permanence, object availability, virtual memory, atomic transactions, and other application specific database needs.

In an typical computer system the primary store is the memory that the CPU acts on directly when executing instructions. Secondary Store is auxiliary, typically non-volatile storage that the CPU references via indirect means. Primary store is usually much faster, efficient and expensive while secondary store is usually slower, larger, cheaper and non-volatile. Disk drives are a well known example of secondary storage.

In the design of Alpha, an attempt was made to provide a conceptually unified model of primary and secondary storage. The traditional interface to secondary storage is via some sort of file system which the application must explicitly manipulate. Files are opened, records are read or written, then the file is closed. The Alpha secondary storage system is integrated into the support for objects. Thus, applications simply deal with objects rather than files and file systems.

14.1 Motivation

The Secondary Storage Subsystem (SSS) provides a powerful, flexible interface. The mechanisms provided by the SSS allow a variety of policies such as demand paging, block caching, etc. to be implemented on top, instead of within the secondary system, without a loss of performance. These policies may be implemented within the kernel, the system layer or by client applications. The system provides all the facilities necessary to support a variety of client defined policies, including file systems.

Clients of the subsystem (VM, Object Manager, Transaction Manager, the kernel, and special database applications) use the mechanisms of the secondary storage subsystem to easily define higher layers which abstract away the details of the secondary storage subsystem. This scheme allows the SSS to support VM behaviors and the object permanence abstractions while preserving the salient features of secondary storage (non-volatility, large address space) without cluttering up the kernel with the artifacts of files and file systems.

- **VM** — We assume that there is not enough physical memory for every object to exist in primary store simultaneously. This requires the VM subsystem to multiplex the primary store. The paging policy of the VM subsystem has the responsibility for preserving the link between the Memory Extents of an object and the backing store for each Extent. Typically an Extent is backed by a Secondary Storage Object. The VM subsystem also manages its own buffer cache.
- **Permanence** — This an optional attribute of objects. The Object Manager provides the object permanence services and maintains the link between an Object and its permanent backing store. This backing store is usually provided by a SSO. To avoid race conditions from arising between Object Manager and the VM subsystem a second backing object is provided for objects with the permanence attribute. The Object manager uses a write through cache, i.e., blocks that are read can be cached for subsequent reads but all writes go immediately

to the backing device. The RESTORE operation reads from this cache, not the VM cache.

- **Transactions (lock logs)** — Support for transactions is another attribute of client and system objects. The Transaction Manager provides support for atomic transactions. A fundamental characteristic of transactions is the ability to abort or undo a transaction. To support this service, the secondary storage subsystem provides a mechanism for associating ordered lists of lock log or shadow blocks with a backing object. The shadow blocks are created when portions of an atomic object are modified using LOCK operations. The lock logs are used to restore the previous state of the object if the transaction is aborted. If the transaction successfully commits, the logs are used to update the permanent backing store. If the node fails during the commit process the logs are used to update the VM state of the object (including the VM backing store), update the permanent backing store for the object, and complete the commit.
- **Kernel** — The Alpha kernel itself makes use of secondary storage to provide permanence for itself and some of its data structures. For instance, the list of permanent objects must, itself, be permanent. The kernel has special dispensation which it takes advantage of for boot strapping. However, it uses the normal interface at all other times.
- **Database clients** — The use of higher levels of abstraction to protect the innocent application from the mean-spiritedness of the physical device has a couple of unfortunate consequences for database people. First, the size of the database is constrained by the size of the virtual address of an object. The current industry standard of thirty two bits of VA space is clearly inadequate for some purposes. Second, there is no easy way for a database designer to optimize the performance by caching key information (the pager by page it out) or by the clever placement of blocks on the device (because the application cannot reference the device, we have cleverly hid it).

To deal with these problems we have promoted the interface to the secondary storage system to the user level. Applications that need access to the low level services provided by the secondary storage subsystem can get to them.

14.2 Overview

The system provides a hierarchy for managing access to the physical devices providing secondary storage. The lowest level is the random access Backing Device. The Backing Device Object encapsulates the hardware specific portion of the SSS. The hardware independent layer consists of Partition Objects and Secondary Storage Objects.

BDs are implemented in terms of GIO devices with memory. Data within the backing device is organized as three lists: the free list of data local to this device, the list of backing objects that use this device, and the list of data associated with each backing object.

Partition Objects (PO) encapsulate the free list and the list of backing objects. A Secondary Storage Object (SSO) encapsulates the list of data that provides backing store and additional control information. POs are similar to traditional file systems and SSOs are similar to individual files. When a request is made of the SSS for a SSO from a type

of PO it may be allocated from any PO of that type. The kernel attempts to find space locally but if that is not possible then the request is made of other nodes with secondary storage. The first node volunteer gets the job.

The object oriented design encourages the precise specification of the interfaces and behaviors of the various components of the subsystem. This allows different flavors of Secondary Storage Objects, Partition Objects and Backing Device objects to coexist. The ability to support multiple versions of the SSS objects allows the system designer to tailor the SSS to suit the application domain. Different applications will have different secondary storage demands and access patterns. It may be desirable for a partition object to allocate blocks to a secondary storage object in contiguous chunks. It may be important for two or more SSO's to be located near one another on the backing device to minimize seek time.

Secondary Storage Objects are used by the Object Manager to help realize object permanence. SSOs themselves have a permanence/transience attribute. Transient SSOs are used by the VM subsystem. Transient SSOs are garbage collected when the node controlling the PO is reset. Permanent SSOs persist across node failures and are used to resurrect permanent objects.

It is less expensive to manage the transient portion of secondary storage than it is to manage the permanent portion. This is because no redundancy is required for the disk directories of transient storage, the disk sectors of transient storage can be heavily cached, and the cached sectors do not have to be written to disk except to make room in the cache. On failure, the transient portion of an object must be reinitialized, but no complicated recovery procedures are required.

Care must be taken in the design of the SSS to provide adequate support for object permanence. Writes must be propagated to the backing device, they may not be cached, to ensure that the written state will be preserved across node failures. Optional device mirroring must be available at the device driver level to protect against media failures. Modification logs are provided in atomic SSOs to support transactions on permanent objects. The portion of the SSS that supports permanence must provide atomic stable storage [Lampson 81].

The objects maintained within the secondary storage facility are registered in the kernel's Dictionary, and may be accessed through the operation invocation facility in the same manner as objects in primary memory. Therefore, the same functionality afforded to client objects by the invocation facility is available to objects in secondary storage. The use of the operation invocation mechanism as the interface between primary and secondary storage supports the construction of a reliable and available secondary storage facility, with location-transparent global access to the objects maintained in secondary storage.

A node can be bootstrapped from the network or the secondary storage subsystem. Location-transparency in conjunction with Partition Groups allows any Alpha node with secondary storage to provide backing store for client objects on other node. Client objects may migrate from node to node without notifying the secondary storage system because all the references between client object and the backing object are via location

transparent capabilities. Thus the sum of all the secondary storage facilities in the system is available to any node.

Each node contains in non-volatile storage (ROM, battery backed up RAM, etc.) the information it needs to each of its boot devices. The boot backing device, if it exists, contains the Root Partition Object for that node. The RPO is used as a repository for information about the kernel and the rest of the SSS. The RPO contains the kernel image, the list of permanent objects and other SSOs used by the kernel.

14.3 Secondary Storage Objects

The Secondary Storage Objects encapsulate the SSO control block, the list of data blocks, the list of shadow block (if any) and the operations and data structures for managing the various lists of blocks. SSOs are roughly analogous to files. Each SSO is associated with a parent Partition Object. They are the backing objects used by the VM system, the Object Manager, and the Transaction Manager. The model assumes that SSO data can be randomly accessed and that I/O operations can use variable length buffers.

The Secondary Storage Object Control Block contains the following information:

permanent;	<i>/* permanence flag */</i>
transactions;	<i>/* transaction flag */</i>
object_type;	<i>/* type of object flag */</i>
dependent_obj;	<i>/* dependent object */</i>
PO;	<i>/* owner Partition Object */</i>
data;	<i>/* list of blocks */</i>
lock_logs;	<i>/* list of blocks associated with locks */</i>

The SSO_control_blk structure contains object attribute information like permanence and transactions, a type field to indicate things like TEXT, DATA, THREAD, etc., a capa for the PO that owns this SSO, a capa to the object it is providing backing store for, and lock bits. The list of blocks is a list of address/length pairs that describe the bytes within the PO that are associated with this SSO. The addresses of these blocks are relative to the parent PO.

The SSO is an indirect object. Each SSO is associated with a parent PO. Each read and write operation is directed to the parent PO which has the responsibility for massaging the request into something a backing device understands. Since the PO provides an intervening layer for each I/O request it can also determine whether the Partition is active before honoring the request.

When used by the VM system to provide backing store for kernel entities and objects the VM system associates these entities with SSOs to act as backing objects. When the SSO is used as backing store by a single entity, the capa for the entity is maintained within the SSO control block. This capa is used in conjunction with the Permanence Bit to give the Garbage Collector Daemon information. If the capa points to a non-existent object, then the SSO can be deleted and its block returned to the Partition's free list.

It is possible that an SSO could act as the backing object for many such entities. For instance, the text is duplicated in each instance of a particular type of client object. Since the text is treated as read_only it would be nice to be able to use the same backing object

rather than duplicating it. In this situation we cannot reasonably track all the objects that might depend on this SSO so we don't even try. The SSO is marked as having multiple users and the garbage collector has to use different rules to determine whether the SSO should be collected.

Secondary Storage Object operations:

- **read(address, size)** — if the range of addresses is not on the allocated list, do not demand allocate buffers and return an error.
- **write(address, size)** — if the range of addresses is not on the allocated list, demand allocate the buffer and write it.
- **set_size(size)** — Set the maximum size, i.e., number of bytes, that is allowed on the backing device for this SSO.
- **get_size(max_size, current_size)** — Return the maximum and actual size of the space used on the backing device by the SSO in bytes.
- **recover** — This operation is invoked to return the SSO to a consistent state following a node failure. It is similar to the UNIX program fsck. This operation is called by the recover operation on the parent Partition object.

14.4 Partition Objects

These objects encapsulate a list of SSOs and a list of free blocks. The list of SSOs are maintained as a flat name space. There is no hierarchy of files or directory structure. However, such a hierarchy could be easily implemented at a higher level. Partition Objects of the same type are said to belong to the same Partition Group.

In addition to encapsulating a free list of blocks, it has algorithms for allocating and deallocating blocks from the list. It creates and destroys the SSOs associated with it. The activate operation enables the PO and the SSOs associated with it. Once activated, the other operations are allowed. The deactivate operation sets a lock that disallows operations on SSOs associated with this PO and any operations other than activate on the PO itself.

POs are roughly analogous to file systems in other operating systems. There is nothing inherent in the interface that prevents the use of a variety of file system architectures from being built on top of POs and SSOs. Further, the various POs within the system could manage their various devices using completely different data structures and algorithms.

The Partition Object Control Block contains the following information:

type	— <i>Partition Group</i>
free blocks	— <i>List of free blocks owned by this PO</i>
transient SSOs	— <i>List of temporary SSOs</i>
permanent SSOs	— <i>List of permanent SSOs</i>
backing device	— <i>Backing device for the PO</i>
next PO	— <i>List of POs (pointer for the kernel)</i>

The SSO logically encapsulates the backing data but the PO does the real work. However, it does not cache the data blocks. The caching of those blocks is left to the clients

of the secondary storage system. Any reads or writes to the data blocks results in I/O to the backing device. The PO does cache the SSO Control Block and the list of data blocks.

The PO defines the following operations:

- **activate** — Allow other operations on this partition.
- **deactivate** — Disallow any operations except activate on this object.
- **create_SSO** — Instance an SSO that belongs to this partition.
- **delete_SSO** — Delete an SSO that belongs to this partition.
- **allocate** — Allocate a block from the free list to the requesting SSO.
- **deallocate** — Deallocate a block from the SSO and return it to the free list within the PO.
- **read(address, size)** — If the range of addresses is not on the allocated list, do not demand allocate buffers and return an error.
- **write(address, size)** — If the range of addresses is not on the allocated list, demand allocate the buffer and write it.
- **get_list_of_SSOs** — Return the list of SSOs that are owned by this PO.
- **free_space** — Return the amount of free bytes in this partition.
- **recover** — This operation is invoked to return the Partition to a consistent state following a node failure. It is similar to the UNIX program fsck. This operation calls the recover operation on each SSO in the Partition.

14.5 Backing Objects

While secondary storage is usually associated with disk and tape drives there are a variety of devices including optical disks, magnetic bubbles, and external RAM that could be useful. Our secondary storage strategy attempts to be as device independent as possible. This desire prompted the invention of an ideal secondary storage device.

Our ideal backing device acts like randomly addressible non-volatile memory with no seek delays. This ideal device can be addressed with the same granularity as primary memory. The size of a I/O operation ranges from a single word to the entire contents of the backing device. One instance of this type of backing device might be a box full of RAM with some sort of UPS. Less than ideal backing devices, such as disk drives, can be supported by a convention of encapsulating them within a driver that provides the ideal behavior via software. Backing devices are Backing Objects (BO).

These objects have the responsibility of presenting an idealized backing store device interface the rest of the secondary storage system. These objects have random access I/O operations. They have the responsibility to translate virtual addresses into device addresses. They have the responsibility for doing any packing and unpacking of I/O requests that have to be done for the particular device. The rest of the secondary storage system does not want to know anything about device specific block sizes. This object either calls a device object which actually controls a device or implements the functions of a device object and controls a physical device itself.

Typical formatted device layout:

```

Device control block
  # of cylinders
  # of platters
  sectors / track
  bytes / sector
  interleave
  bad block map
  partition sizes and pointers
  etc.

Root Partition Control Block
  Free blocks
  SSOs
  Kernel boot image
  List of Permanent Objects
  Other SSOs

Partition Control Block
  Free blocks
  SSOs
  .
  .
  .

```

BDOs define the interface between the Secondary Storage System and the I/O system. The requirements are fairly simple—random access, variable length reads and writes. BDOs do not have to map directly to physical devices. A single BDO might provide access to more than one physical disk. The addresses used by the secondary Storage system could be distributed over a set of disk drives. This technique could be used to provide disk striping. A variation on this theme could be used to provide disk mirroring. Each I/O request could be sent to each disk thereby creating mirror copies.

BDOs could be associated with a section of physical memory. This would allow for the creation of memory disks. Memory disks are most useful in a development environment where the programmer needs to experiment with different secondary storage data structures but does not have easy access to real secondary storage devices. While memory disks might not be very useful in a typical Alpha environment, the system is flexible enough to provide for their easy development. It is this flexibility (rather than the specific facility of memory disks) that the design strives for.

Some operations:

- **read** — get a block of data from a device.
- **write** — put a block of data to a device.
- **format** — write the logical organization of a BDO on to a device, usually a GIO device with storage (i.e., a disk). This operation is similar to the UNIX format program.

References

- [Bernstein 81] Bernstein, P. A. and Goodman N.
Concurrency Control in Distributed Database Systems.
ACM Computing Surveys 13(2):185-221, June 1981.
- [Cheriton 87a] Cheriton, D. R.
The V Distributed System.
Technical Report, Stanford University, March 1987.
- [Cheriton 87b] Cheriton, D. R.
UIO: A Uniform I/O System Interface for Distributed Systems.
ACM Transactions on Computer Systems, 5(1), February 1987.
- [Chesson 88] Chesson, G., Brendan, E., Schryver, V., Cherenon, A. and Whaley, A.
XTP Protocol Definition: Revision 3.1.
Silicon Graphics Technical Report, March 1988.
- [Clark 88] Clark, R. K., Kegley, R. B., Keleher, P. J., Maynard, D. P., Northcutt, J. D., Shipman, S. E. and Zimmerman, B. A.
An Example Real-Time Command and Control Application.
Archons Project Technical Report #88032, Department of Computer Science, Carnegie-Mellon University, March 1988.
- [Cooper 84] Cooper, E. C.
Replicated Procedure Call.
In *Proceedings of Third Annual Symposium on Principles of Distributed Computing*, pages 220-232. ACM, August 1984.
- [Cox 86] Cox, B. J.
Object-Oriented Programming.
Addison-Wesley, Reading, Massachusetts, 1986.
- [Goodman 83] Goodman, N., Skeen, D., Chan, A., Dayal, U., Fox, S. and Ries, D.
A Recovery Algorithm for a Distributed Database System.
In *Proceedings, Second Symposium on Principles of Database Systems*.
ACM, March 1983.
- [Habermann 76] Habermann, A. N., Flon, L. and Coopridner, L.
Modularization and Hierarchy in a Family of Operating Systems.
Communications of the ACM 19(5):266-272, May 1976.
- [Herlihy 85] Herlihy, M. P.
Using Type Information to Enhance the Availability of Partitioned Data.
Technical Report CMU-CS-85-119, Department of Computer Science, Carnegie-Mellon University, April 1985.
- [Herlihy 86] Herlihy, M. P.
A Quorum-Consensus Replication Method for Abstract Data Types.
ACM Transactions on Computer Systems 4(1), February 1986.

- [Henize 84] Henize, J. A.
Understanding Real-Time UNIX.
MASSCOMP Technical Report, January 1986.
- [Jensen 88] Jensen, E. D., Test, J. A., Reynolds F. D., Burke E. J. and Hanco, J. G.
Alpha Release 2 Design Summary.
Archons Project Technical Report #88082, Department of Computer Science, Carnegie-Mellon University, August 1988.
- [Kung 81] Kung, H. T. and Robinson, J. T.
On Optimistic Methods for Concurrency Control.
ACM Transactions on Database Systems 6(2):213-226, June 1981.
- [Lampson 69] Lampson, B. W.
Dynamic Protection Structures.
In *Proceedings of the Fall Joint Computer Conference*, pages 27-38, IFIPS, 1969.
- [Lampson 81] Lampson, B. W., Paul, M. and Siegert, H. J. (editors).
Lecture Notes in Computer Science. Volume 105: *Distributed Systems—Architecture and Implementation.*
Springer-Verlag, Berlin, 1981.
- [Leffler 87] Leffler, S. J. *et. al.*
An Advanced 4.3BSD Interprocess Communication Tutorial.
Unix Programmer's Manual Supplementary Documents 1, June 1987.
- [Locke 86] Locke, C. D.
Best-Effort Decision Making for Real-Time Scheduling.
Ph.D. Thesis, Department of Computer Science, Carnegie-Mellon University, May 1986.
- [McKendry 85] McKendry, M. S. and Herlihy, M. P.
Time-Driven Orphan Elimination.
Technical Report CMU-CS-85-138, Department of Computer Science, Carnegie-Mellon University, 1985.
- [Needham 82] Needham, R. M. and Herbert, A. J.
The Cambridge Distributed Computing System.
International Computer Science Series, 1982.
- [Northcutt 87] Northcutt, J. D.
Mechanisms for Reliable Distributed Real-Time Operating Systems: The Alpha Kernel.
Academic Press, Boston, 1987.
- [Northcutt 88a] Northcutt, J. D.
The Alpha Operating System: Requirements and Rationale.
Archons Project Technical Report #88011, Department of Computer Science, Carnegie-Mellon University, January 1988.

-
- [Northcutt 88b] Northcutt, J. D. and Clark, R. K.
The Alpha Operating System: Programming Model.
Archons Project Technical Report #88021, Department of Computer
Science, Carnegie-Mellon University, February 1988.
- [Northcutt 88c] Northcutt, J. D.
The Alpha Distributed Computer System Testbed.
Archons Project Technical Report #88033, Department of Computer
Science, Carnegie-Mellon University, March 1988.
- [Northcutt 88d] Northcutt, J. D. and Shipman, S. E.
The Alpha Operating System: Programming Utilities.
Archons Project Technical Report #88061, Department of Computer
Science, Carnegie-Mellon University, June 1988.
- [Pease 80] Pease, M., Shostak, R. and Lamport, L.
Reaching Agreement in the Presence of Faults.
Journal of the ACM 27(2):228-234, 1980.
- [Reynolds 88] Reynolds, F. D., Hanks, J. G., Test, J. A., Burke, E. and Jensen, E. D.
Interim Alpha Release 2 Kernel Interface Specification.
Archons Project Technical Report #88122, Department of Computer
Science, Carnegie-Mellon University, December 1988.
- [Ritchie 74] Ritchie, D. M. and Thomson, K.
The UNIX Timesharing System.
Communications of the ACM, 17(7):365-375, July 1974.
- [Ritchie 87] Ritchie, D. M.
A Stream Input-Output System.
UNIX SYSTEM Readings and Applications, Volume II, Prentice-
Hall, 1987.
- [Shipman 88] Shipman, S. E.
The Alpha Operating System: Programming Language Support.
Archons Project Technical Report #88042, Department of Computer
Science, Carnegie-Mellon University, April 1988.