



---

**Poseidon House  
Castle Park  
Cambridge CB3 0RD  
United Kingdom**

TELEPHONE:  
INTERNATIONAL:  
FAX:  
E-MAIL:

**Cambridge (01223) 515010  
+44 1223 515010  
+44 1223 359779  
apm@ansa.co.uk**

---

## **ANSA Phase III**

# **An Overview of Real-Time ANSAware 1.0**

**Guangxing Li**

### **Abstract**

Distributed computing and real-time computing are well established areas of research, but their integration is yet to be studied because they seldom use compatible techniques. This paper provides an overview of an ANSA based system environment (named ANSAware/RT) for distributed real-time applications. The focus of this article is the engineering mechanisms necessary to extend ANSAware to support real-time computing. ANSAware/RT incorporates real-time tasks and communication channels as its basic programming components. It synthesises aspects of resource requirements, allocation and scheduling into an object-based programming paradigm.

This is an external paper for the Journal of Distributed Systems Engineering.

---

APM.1285.00.01

**Draft**  
External Paper

5th December 1994

---

**Distribution:**  
**Supersedes:**  
**Superseded by:**



---

# Contents

---

1	1	<b>An Overview of Real-Time ANSAware 1.0</b>
1	1.1	Introduction
1	1.2	ANSA
2	1.2.1	ANSA object model
2	1.2.2	ANSA engineering model and ANSAware
3	1.3	Distributing real-time objects
4	1.4	ANSAware/RT design
4	1.5	ANSAware/RT tasking model
4	1.5.1	Real-time objects
5	1.5.2	Real-time object invocation
5	1.5.3	Scheduling
6	1.6	ANSAware/RT communication model
7	1.6.1	A parallel protocol stack
7	1.6.2	A timed RPC protocol
9	1.7	ANSAware/RT implementation
10	1.8	ANSAware tasking
10	1.9	ANSAware/RT tasking
11	1.9.1	Global data protection
11	1.9.2	Thread private state
11	1.9.3	Stacked threads
11	1.9.4	Thread
12	1.9.5	Entry
12	1.9.6	Synchronous I/O
13	1.9.7	Communication tasks and system tasks
13	1.10	ANSAware communication system
14	1.10.1	Interface reference
14	1.10.2	Message passing protocol
14	1.10.3	Execution protocol
14	1.10.4	Channel and session
14	1.10.5	Bindings
15	1.11	ANSAware/RT communication system
15	1.11.1	QoS and explicit binding
15	1.11.2	State-full message passing protocol
16	1.11.3	State-full execution protocol
16	1.11.4	Timed remote execution protocol
16	1.11.5	In-band QoS
16	1.11.6	Session overridden
17	1.12	Performance evaluation
19	1.13	Summary
19	1.14	Related work
19	1.15	Acknowledgement



---

# 1 An Overview of Real-Time ANSAware 1.0

---

## 1.1 Introduction

---

Distributed real-time computing is spreading. The need of a supporting system environment is inevitable. Even though distributed computing environment and real-time computing environment are established areas of research, their integration is yet to be studied, because they seldom use compatible techniques.

The need for an integrated system environment for distributed real-time applications is driven by two technology trends. Firstly, advances in digital communication networks and in personal workstations are beginning to allow the simultaneous processing of real-time data, voice, and video. This requires distribution and real-time control functionality to be intrinsic elements of the system, rather than bolted on after thoughts. Secondly, the size of real-time systems is increasing: one-million-line real-time software systems in telecomms, manufacturing, transportation and other areas are common today [Gopinath and Bihari 93]. Such systems are very large and distributed by nature. There is an increasing need to adopt an open architectural approach so that real-time systems engineering can be augmented with other techniques to address scale, reuses, evolution, distribution and other issues.

Real-time processing places unique requirements on distributed systems including predictability, programmer control, timeliness, mission orientation and performance [Li and Bacon 94]. Such features are yet to be provided by the current distributed computing environments.

This paper provides an overview of the design, implementation and performance evaluation of the real-time ANSAware (ANSAware/RT) 1.0. ANSAware/RT is based on the ANSA architecture [Herbert 94] and its example implementation ANSAware (AW) 4.1 [APM 92]. The focus of this article is the engineering mechanisms; other issues such as architecture and application programming interfaces can be found in [Li and Otway 94, Li 94a].

This paper is organised as follows. Section 2 gives a short introduction about ANSA. Section 3 discusses the important aspects of real-time objects. Section 4-6 presents the engineering designs required to extend ANSA for real-time systems. Section 7-11 outlines the implementation techniques used for ANSAware/RT in comparison with ANSAware. Section 12 gives the performance evaluation of the system by using of the Distributed Hartstone Benchmark [Mercer et al. 90]. Section 13 gives a summary and finally section 14 discusses related work.

## 1.2 ANSA

---

ANSA is an Architecture for Open Distributed Processing (ODP), which provides new ways of thinking about the design and construction of object oriented client/server distributed systems. ANSA uses five complementary

models (enterprise, information, computational, engineering and technology) to describe architectural components, of which the computational model and engineering model are most relevant to this work. The overall ANSA framework has been captured in the joint ISO/IEC and ITU-T Basic Reference Model of ODP (RM-ODP) [ISO 94].

### 1.2.1 ANSA object model

The core of the ANSA computational model (ACM) is the use of *objects* as units of distribution for management and replacement. An object has one or more *interfaces* that are the points of provision and use of services. Interfaces are first class entities in their own right and references to them may be freely passed around the system.

An interface contains a set of named *operations* (i.e. procedures or methods) which defines its type. Interfaces have the usual remote procedure call style of interaction: operations are invoked with a set of arguments and a response is returned. Arguments and results to invocations consist of references to other interfaces. The effect of an interaction is that the client and server share access to the argument and result interfaces. This model makes each interface an abstract data type. ANSA also has a stream interface, which is beyond the scope of this paper.

### 1.2.2 ANSA engineering model and ANSAware

ANSAware is an implementation of the ANSA computational model and an example of the ANSA engineering model (AEM). ANSAware is a suite of software for building ODP systems, providing a basic platform and software development support in the form of program generators and system management applications. ANSAware provides a uniform view of a multi-vendor world, allowing system builders to link together distributed components into network wide applications.

The AEM provides a framework for the specification of mechanisms to support distribution of application programs that conform to ACM. The main components of the AEM are:

- *transparency mechanisms* which automates aspects of distribution such as object migration and object replication.
- *a nucleus* which provides minimal and sufficient support for the implementation of distribution. It encapsulates all of the heterogeneity of processor and memory architecture.
- *capsules*: which are collections of application objects, transparency mechanisms and nucleus objects forming a virtual node of a network.
- *threads*: sequences of instructions modelling a computational model activity within a capsule. A thread represents a unit of potentially concurrent activity that can be evaluated in parallel with other threads, subject to synchronization constraints.
- *tasks*: virtual processors which provide threads with the resources (e.g. stacks) they require to progress. Tasks<sup>1</sup> provide the resources for real concurrency. An ANSA task is conceptually equivalent to an operating system *thread*.

- *interface references*: identifiers which contains sufficient information to allow their holder (client) to establish communication with the interface denoted by the reference (server).
- *channels*: resources required to enable end to end communication. The initiating side (client) end-point of a channel is called a *plug*. The receiving side (server) end-point is called a *socket*.

### 1.3 Distributing real-time objects

The essence of a real-time object model is to provide the basic abstractions so that stringent timing constraints of real-time activities are respected (guaranteed ideally). The main difficulty is that the actual timing characteristics of software are determined not only by the ANSAware/RT processor speed, but also by the sharing policy for scarce resources. In most high level languages, this dependency is considered as non-essential detail that is to be hidden from the programmer. As a result the performance of software implemented in these languages becomes sensitive to resource allocation strategies (in a dynamic system, this means performance depends on system load), and outside the control of individual programmers. More complex resources such as the communication subsystem of distributed systems further accentuate the problem with the introduction of (sometimes distributed) resource allocation algorithms which are usually inaccessible to the application programmer.

Object interdependence can be classified into two categories:

- *static* interdependence --- the structural relationships between objects,
- *dynamic* interdependence --- the interactions (execution views) between objects.

Many useful results are known about the static relationships between distributed objects. Related concepts, such as abstract data typing, type checking and subtyping, are accepted and used widely. On the other hand, little consensus has been achieved on the execution view of objects. Many approaches to object execution have been proposed, some of which are the *active object* model, the *passive object* model, and the *actor object* model.

For real-time applications, this execution aspect is of vital importance --- it has fundamental impact on the *predictability* of computational activities. Real-time object execution models are required to address not only how the computational activities are carried out, but also how shared resources are used (i.e. the manner in which contention for system resources is resolved taking into account timing constraints of real-time activities). The latter issue is often neglected and considered irrelevant engineering detail in non-real-time computing. Distributed real-time systems must provide support for the specialized requirements of real-time communication, tasking, scheduling, and control. These requirements must be explicitly addressed in an object execution model, if the object-oriented approach is expected to be applicable to a real-time world.

---

1. ANSA threads are cheap resources (each requires less than one hundred bytes of memory); whereas ANSA tasks are expensive resources (each requires several kilobytes of memory). In a distributed application there may be many threads (e.g. 100's or 1000's); it is important only to allocate a task to execute a thread when there is a processor available to run it.

## 1.4 ANSAware/RT design

The collective effect of ACM and AEM defines the ANSA Object Execution Model. This model is designed for object distribution, but not for real-time applications. It lacks real-time predictability in the following sense:

- multiplexing both tasks and communication channels whenever possible
- both thread/task scheduling and communication scheduling are implicit
- no abstraction is provided to express urgency and resource requirement for application programmers.

A real-time object model can be obtained by extending the execution model with explicit resource allocation, real-time scheduling and real-time communication support. The ANSAware/RT real-time object model has two parts: a tasking model and a communication model

## 1.5 ANSAware/RT tasking model

### 1.5.1 Real-time objects

A real-time object is composed of data, one or more tasks of execution, and a set of interfaces. A new abstraction, *scheduling entry*, or shortly *entry*, is introduced as the basic mechanism for real-time scheduling.

An entry is a thread queue with a record of control data. An entry may be created dynamically, and interfaces of an object may be bound to it. When an interface is bound to an entry, each operation request on the interface will be transferred (by the infrastructure) to a thread enqueued on the entry. Any thread representing a computational activity is also spawned on an entry. The entry is an engineering concept which is confined within a capsule.

Figure 1.1: Real-time object illustration

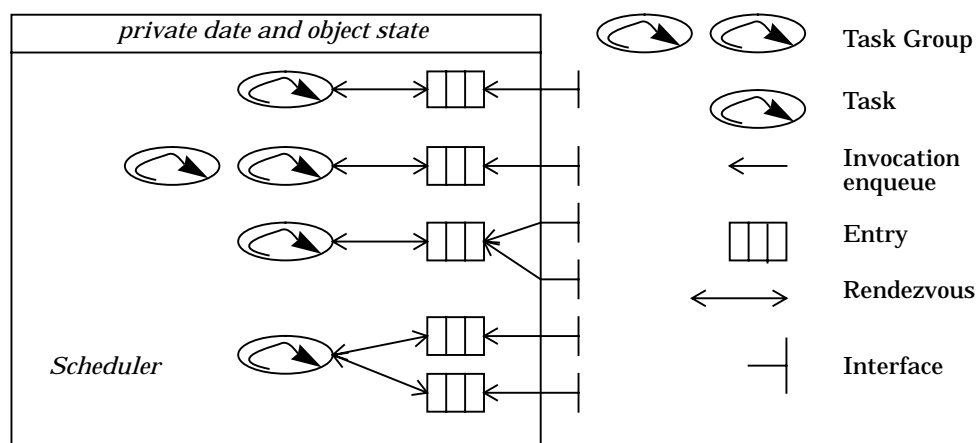


Figure 1.1 gives a graphical illustration of a real-time object.

Flexible tasking is based on the entry abstraction. System tasks may be allocated for each individual entry. The tasks are dedicated to execute the threads on the entry they support. When executing a thread, a task is also allowed to *rendezvous* with other entries dynamically. A *rendezvous* of a task with an entry means that the task waits to accept and execute one thread on



the entry. Different control parameters may be selected for each entry to choose a thread queuing policy, a task/entry rendezvous policy, and to enforce concurrency controls.

In such an object model with data, interface, entry and tasks encapsulated within a capsule, there is a choice of how many entries are allocated, which interface is attached to which entry, how many tasks are allocated to an entry, whether a task can rendezvous with a specific entry, and what kinds of resource scheduling policies are used.

The choice to allocate a separate entry for some interfaces reflects the need to separate these interfaces from others for the purpose of resource management. The number of tasks allocated to an entry not only enforces the real concurrency allowed for the execution of threads on the entry, but also affects the real-time scheduling properties, for example, *preemptivity* and *priority inversion* [Li 93]. The flexibility for allowing a task to rendezvous with an entry enables an application to have complete control over its virtual processor(s) based on its knowledge of the system state. These resource management activities can all be done dynamically, increasing the flexibility and usefulness in an open dynamic environment.

### 1.5.2 Real-time object invocation

ANSAware/RT allows the association of an optional *priority* (criticality) and/or *deadline* with each invocation. As the invocation crosses from one object to another, this priority and/or deadline is passed and becomes a property of the executing thread on the server site.

### 1.5.3 Scheduling

The main goal of the real-time tasking design is to allow the maximum control of scheduling at the application level. Care has been taken to achieve the balance between flexible and deterministic scheduling. A policy/mechanism separation approach has been taken to address the diversity of real-time programming. Scheduling is defined in layers as:

- thread scheduling --- the rendezvous scheduler on each entry.
- task scheduling --- the nucleus scheduler on tasks.

Task scheduling and thread scheduling are two separate, but related scheduling domains. Task scheduling is defined in the nucleus or the underlying operating system kernel. Preemption is used together with task scheduling parameters to order (either partially or completely) the otherwise non-deterministic behaviour of the task execution. Thread scheduling is defined per entry. Task scheduling manages multiplexing of task executions over processor(s). Thread scheduling manages the multiplexing of requests (thread) over tasks. The primary function performed by multiplexing is the sharing of processor resources, which is similar to the multiplexing in communications systems and protocols for sharing communication resources. The use of separate entries to process requests on separate interfaces offers a number of (potential) advantages

- allows the use of a specific scheduling policy (thread scheduling policy) suitable for each interface or each interface class,
- allows the possibility of using interface specific tasks to serve requests, and thus allows for more efficient resource utilisation,

- separate entries may be processed in parallel, thus increasing performance,
- allows the possibility of end-to-end scheduling and guarantees,
- preserves the modularity and separation of service interfaces.

There are two issues in thread scheduling management. One is how a thread is enqueued in an entry (with the assumption that the first thread in the queue is executed first). Such a policy may be a system defined one, like invocation priority based, or an application provided one. Some typical thread enqueue policies are (1) first come first service (FCFS), (2) priority based (PB), (3) deadline based (DB), (4) priority and deadline based (PDB).

Another issue is how a serving task rendezvous with a thread in an entry, i.e. how the thread scheduling parameters (priority and/or deadline) are used/ inherited by the task. This is defined by a task/thread rendezvous policy. Such a policy affects how the serving task competes for processor resources with other tasks. Some typical task/thread rendezvous policies are (1) null --- the priority/deadline of a thread has no effect on the serving task, (2) priority inheritance, (3) transitive priority inheritance, (4) priority ceiling and (5) deadline inheritance.

Detailed examination of some typical real-time scheduling schemes, such as priority based scheduling and deadline based scheduling, can be found in [Li 93].

## 1.6 ANSAware/RT communication model

---

Real-time applications present more complicated functional requirements to the underlying communication systems. ANSAware/RT provides abstractions to express the individual Quality of Service (QoS) constraints for a communication channel and the selection of in-band QoS parameters of a communication channel.

The following abstraction is provided in ANSAware/RT

- the allocation of a separate communication channel for each client/server binding
- the association of a communication QoS to a channel
- the association of an in-band QoS to an invocation.

The in-band QoS object supports the specification of a priority, a timeout, a deadline, a deadline type, and any other QoS parameters a communication channel may support, this depends on the individual network system.

Priority and deadline are used to convey the urgency of a real-time activity across a network. The combination of a timeout, a deadline and a deadline type can be used to bound the expected execution time of an invocation, which is further discussed in section 1.6.2.

From engineering point of view, the following mechanisms are required:

- a parallel protocol stack
- a timed RPC protocol

### 1.6.1 A parallel protocol stack

A parallel communication protocol stack allows the preallocation of communication resources (a separate channel, for example) and the removal of layered multiplexing. The main gain is that it allows the application to explore the communication QoS that the low-level operating environment can provide. For example, in an ATM environment, a channel (or a virtual circuit) is allowed to associate with various transportation QoS, such as jitter, delay, priority etc. Even in an ordinary operating environment, this design allows the choose of communication protocols, such as TCP, UDP, IPC etc.

### 1.6.2 A timed RPC protocol

Arbitrary delays associated with synchronous invocation cannot be tolerated due to the time-dependent nature of real-time applications. A dependable protocol is desirable to provide a timeliness service for real-time RPC, or timed RPC (TRPC).

Invocations in ANSAware/RT can attach deadline constraints to their communication requests. Such TRPC calls raise the following three issues:

- the management of time in a networked environment. Intuitively, a deadline is an upper bound, which is placed on the time duration for the invocation to occur. Therefore both the server and client must have the same sense of time --- the deadline. It is thus necessary to assume a common sense of time is provided by the infrastructure between a client and a server
- the interpretation of deadlines
- a communication protocol to implement reasonable meanings of deadlines.

There are two goals one might try to accomplish with the deadline of a TRPC:

- to establish a bound on the time at which the delay in awaiting a TRPC call expires
- to establish a bound on the time at which a TRPC call is either scheduled to execute and finish or is unschedulable and cancelled.

The design of a TRPC is complicated by the fact that the end-to-end delay of messages can be arbitrary or even infinite (messages can get lost). It can be shown that the two goals are not mutually compatible. In its simplest form, in which each request takes zero service time, the TRPC problem is equivalent to the *timed synchronous communication* problem [Lee and Davidson 90]. In the case of message loss, timed synchronous communication is the well known *Two Generals* problem, in which the two generals are trying to agree upon a *common time* of attack before a deadline but can only communicate via unreliable messengers. Such a protocol does not exist.

The design of a TRPC is further complicated by the fact that making the decision of whether a request is schedulable at the server side is often *unattainable* --- a guarantee scheduler often makes many of the impossible assumptions such as that the invocation service time is known, operations are independent etc.

Because using one deadline value to accomplish the two goals in a TRPC may result in incompatible situations, two arguments --- a *timeout* and a *deadline* -- are used instead. Each is aimed at one goal only. The timeout is used to specify the first goal --- how long the client is willing to wait for its result. It affects a client side of the TRPC protocol only. The deadline is used for the

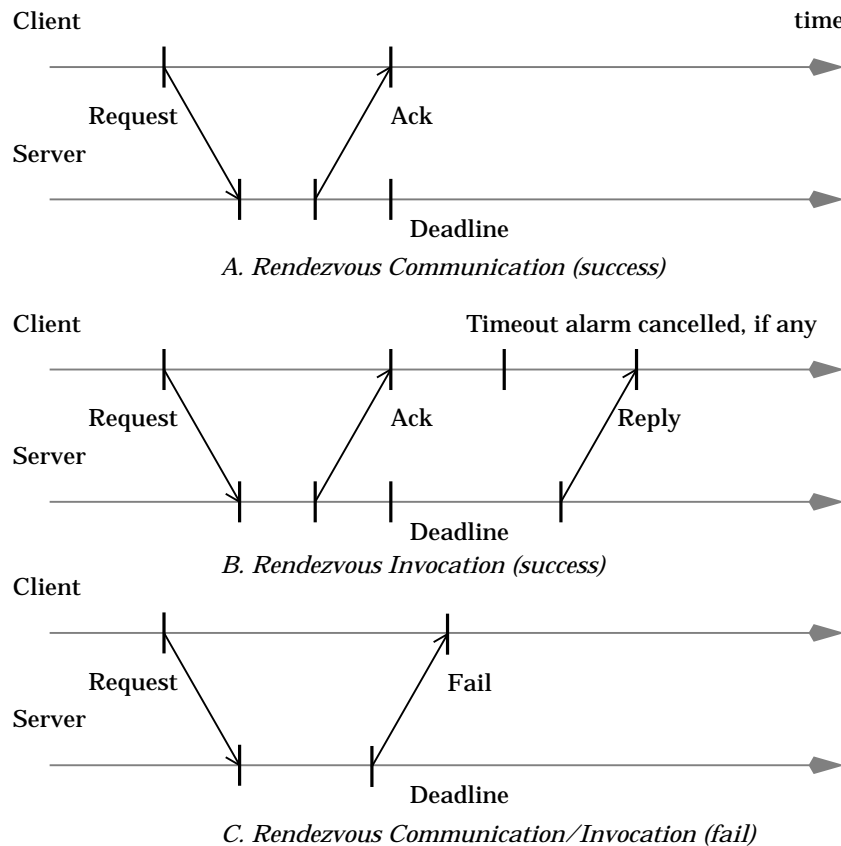
second goal --- within which the request should be executed on the server. It affects the server side of the TRPC protocol only.

It should be pointed out that using the two separate arguments does not solve the TRPC consistency problem. Rather, the two arguments give the problem a more realistic definition, allowing different relaxations be explored.

The first relaxation is using a *timeout* to enforce the client's *absolute* deadline. The client decides that the request is unsuccessful if it does not get a reply/acknowledgement from the server by the timeout. There is a possibility for *inconsistent decisions* --- the client believes the request is failed, while the server knows the request is successful. Deadlines may or may not be used in this situation. The timeout expiration presents the client an exception situation of *don't know*. It is up to the client to take further rescue actions.

The second relaxation is using a deadline to specify a client's objective time value by which the request should be finished. Whether this deadline can be guaranteed or not is purely a matter of server scheduling and message passing delays. In this relaxation, the client waits until a reply/acknowledgement is received from the server. Therefore, the client deadline is not *absolute*. This relaxation allows a client and its server to reach a consistent decision.

**Figure 1.2: Rendezvous Communication/Invocation Interaction**



The second relaxation can be further extended by relaxing the meaning of a deadline. Instead of bounding the finishing time of a request, a deadline can be used to bound the start time of a request in the server --- to bound the start time by which the request is rendezvoused with a server task. If the rendezvous is issued before the deadline, then the request is successful and a

success acknowledgement is sent back to the client, otherwise the request is cancelled and a fail acknowledgement is returned. At the client side, there are two possible actions to be taken when it receives a success acknowledgement. One is that the client thinks the request is finished, and control is returned so that it can continue. This is defined by the *RendezvousCommunication* deadline type. Another is that the client cancels its timeout, if any, and waits until a reply is returned later by the server. This is defined by the *RendezvousInvocation* deadline type. The two resulting interaction patterns are illustrated in Figure 1.2.

The deadline type parameter is introduced to choose how the deadline can be used in the server side of a call. It can be used to control the latest start or latest finish time of an invocation.

In summary, an invocation may be associated with an optional timeout, an optional deadline, and an optional deadline type. The collective choice of the three parameters determines the behaviour the TRPC protocol. The result of such a TRPC call can be a *timeout* --- possibly an inconsistent state, a *success* or a *failure*.

---

## 1.7 ANSAware/RT implementation

---

ANSAware/RT was first implemented as RIDE [Li 93] in the Cambridge University Distributed Systems Environment. It used a real-time microkernel named WANDA, and ANSAware 3.0. The environment was composed of interconnected 680x0, VAX, ARM and MIPS machines over an ATM and Ethernet network.

ANSAware/RT 1.0 described in this paper is an evolution of RIDE to run over a standard real-time environment. ANSAware/RT 1.0 also uses binding and QoS concepts taking from the RM-ODP to provide a more general vehicle for resource management and requirement specification. Binding is the process (operation) by which an activity in one object establishes the ability to invoke operations at an interface to some other object. A binding establishes and controls the communication sessions involving multiple objects so that their interactions are possible.

ANSAware/RT 1.0 has been implemented over the DEC/Alpha OSF1 system, the HP/RT target system and LynxOS system. ANSAware/RT 1.0 achieved the following design goals

- compatible with a new version of ANSAware 4.1
- running over a de-facto industry standard: real-time POSIX threads (pthread)
- full pthread real-time scheduling and threading capabilities
- selective communication multiplexing by QoS specification and explicit binding operations
- application controlled resource allocation
- multiple RPC protocols: use different protocols for real-time and non-real-time transportation
- interoperation between different real-time platforms while retaining their real-time features.

In the following sections, AW tasking is explained first and then ANSAware/RT tasking is described in detail. It then proceeds to the discussion of AW communication system and ANSAware/RT communication system.

---

## 1.8 ANSAware tasking

---

AW threads represent points of execution and provide the notion of logical concurrency. AW tasks represents the resources required (stacks) to execute an AW thread and provide the actual concurrency.

Logically, AW starts with several threads and one or more tasks. There is a receiver thread for receiving messages on the communication endpoints, a time thread to execute time-related activities, and an application program thread to execute the user program code.

AW tasks are user level entities implemented through a coroutine package. Additional tasks may be created to provide extra physical concurrency. Tasks are shared by all threads. All threads waiting to execute are queued on one FCFS queue (named entry in ANSAware/RT). The AW nucleus scheduler assigns free tasks to execute queued threads. The scheduler is non-preemptive and is only entered when the current thread/task blocks or terminates. If a thread/task is resumed, the scheduler will return control to it.

AW takes several advantages of the coroutine nature of its tasking package:

- use global, continuous and extensible memory area to store the shared data structures holding the capsule state. AW increases memory for shared data structures in a dynamic manner but requires that the existing memory and the newly allocated memory be contiguous. This requirement has been achieved by copying the existing data to a new location where contiguous memory is available. In this way, memory space is allocated on demand, resulting downsize of AW package.
- use global variables to carry context information. The variables that form the context of an ANSA task are global variables which are shared by all ANSA tasks. Thus, context information is passed to all the procedures through global variables. This allows fast inter task context switch and fast procedure execution (i.e. there is no need to pass context information through procedure parameters).
- shared data are is accessible without a synchronization mechanism. AW task scheduler is non-preemptive, a task, while execution, will not relinquish control until it blocks or terminates. Therefore, it guarantees exclusive access of the shared data in a single processor environment, and there is no need for access protection of shared data.

---

## 1.9 ANSAware/RT tasking

---

In ANSAware/RT, each task is mapped into a pthread, and task scheduling is done by the underlying operating system. All pthread attributes also apply to tasks, allowing the exploitation of preemptive real-time scheduling, multiple scheduling policies, kernel supported synchronization objects, task private data, task exception handling, task synchronous I/O etc. pthread features. The original AW task schedule was made redundant.

### 1.9.1 Global data protection

Because of the real concurrency<sup>1</sup> and preemptive nature of the pthread system, synchronisation is needed to ensure the safe access to shared data. A pessimistic synchronisation approach is taken: all data structures are protected by a single lock. To perform any ANSAware/RT operation, a task must first acquire the lock, then operate on the shared data, and finally release the lock when finished.

### 1.9.2 Thread private state

Each thread has a few private state variables, such as `exception_code`, `exception_state`, `memory_list` etc. These thread private state variables are stored in the global data area in AW. Thread private state are used frequently in both application program and AW operations.

In ANSAware/RT, global data area needs to be protected by a synchronization mechanism. Therefore, the AW thread private data may introduce a significant performance overhead if no change to AW is made.

The solution adopted is to use pthread per-thread state to store AW thread private state (rather than using the global area). Such state information are then accessible by using of the `pthread_getspecific` procedure without a synchronization operation.

The thread private state are actually part of a task private data area. When a task is created, it allocates a private state area as pthread per-thread data, and part of this area is used as thread private state when the task is executing a thread.

### 1.9.3 Stacked threads

AW threads are non-stacked: a task will not execute another thread before it finishes the current one.

ANSAware/RT introduces the dynamic rendezvous mechanism: a thread may rendezvous with another thread while execution. The required extensions are to allow a task to execute another thread when it is executing a thread.

This stacked thread mechanism is implemented by pushing the thread private state area into the task's stack before executing the new thread (so that the new thread still can use the same task private data as its thread private state), and restore the thread private state from stack when the new thread finishes.

### 1.9.4 Thread

Threads are created in two cases: (1) an application may create new threads for additional concurrency; (2) a communication task may create one additional thread for each RPC request from a client. In AW, a new thread is queued on the capsule-wide FCFS thread queue, waiting to be executed by a free ANSA system task. In ANSAware/RT, a new thread is queued on an entry instead of the capsule FCFS queue. In case (1), the application gives an additional entry argument when a new thread is created. In case (2), the binding between an interface and an entry determines on which entry the new thread should be queued.

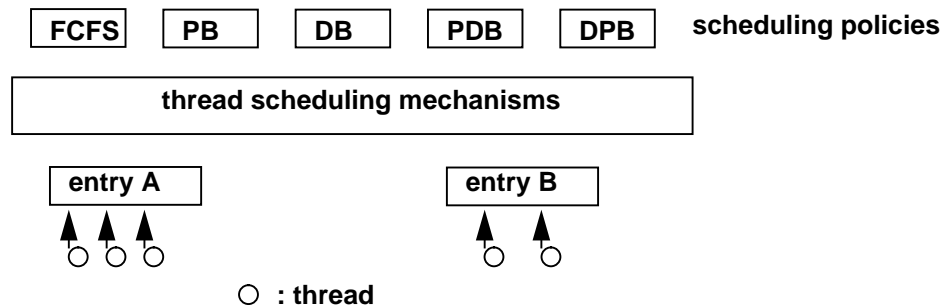
---

1. pthreads can be executed in parallel, for example, in a multiprocessor environment.

### 1.9.5 Entry

Each entry is associated with a thread queue and a thread queuing policy.

Figure 1.3: Thread scheduling: policies and mechanisms



Policy/mechanism separation is used for efficient coding. A common set of thread queuing/dequeuing mechanisms is provided, and on top of the mechanisms a set of scheduling policy objects are placed. Figure 1.2 illustrates such a design.

Each entry is also associated with a rendezvous policy. Each such policy provides two functions: `rendezvous_inheritance` and `rendezvous_deinheritance`. The `rendezvous_inheritance` function is executed before a task executes a thread so that the task can take the thread scheduling parameters into consideration. For example, it allows the task to inherit the thread's priority. The `rendezvous_deinheritance` function is executed after a task finishes the execution of a thread to eliminate any scheduling effect on the task caused by the `rendezvous_inheritance` function.

### 1.9.6 Synchronous I/O

AW assumes a totally asynchronous I/O model because

- it allows the tight combination of communication scheduler and task scheduler for efficient AW activity scheduling
- it prevents a capsule from blocking because of an otherwise synchronous I/O operation.

The asynchronous I/O approach separates out the indication that data is available from the actual reading of the data.

The asynchronous I/O model in AW is supported by

- a UNIX `pin(3)` programming interface. An application can register an interrupt handler to be invoked when input occurs on a `pin` and that handler is then able to spawn a thread to read any input data
- a non-blocking keyboard input library
- a library for supporting X11 applications.

With `pthread` implementation of the tasking system, the asynchronous I/O model is no longer necessary because

- task scheduling is done by OS, there is no tight integration of tasking scheduling and communication scheduling



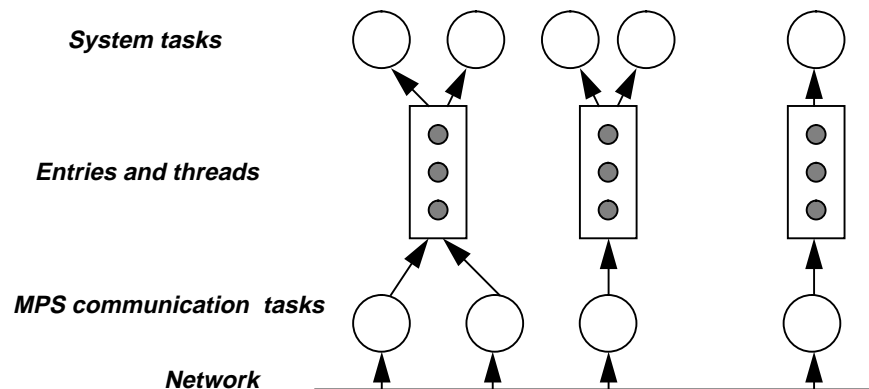
- a capsule will not block when a thread is doing a synchronous I/O

In other words, ANSAware/RT does not need to assume the asynchronous I/O model, and a complete synchronous I/O model is more natural and easy to programming. Therefore, ANSAware/RT removes the pin interface, the non-blocking keyboard input library and the X11 library, and assumes the application programmer will access the equivalent synchronous I/O operations supported by the pthread interface.

### 1.9.7 Communication tasks and system tasks

Dedicated communication tasks are spawned to process incoming messages and the corresponding protocol by using of synchronous I/O operations. For each message passing service (MPS, seeing next section) endpoint (a socket), a task is spawned for handle messages from it. The communication task generates a thread corresponding to each invocation request. The thread is queued on an entry to which the called interface is bound. In this vein, the communication task is actually both a thread generator and a thread scheduler. The threads are executed by system tasks of an entry which are allocated by an application or the capsule. The scenario is shown in Figure 1.3.

Figure 1.4: Communication tasks and system tasks



When a thread makes a synchronous invocation to a server, it blocks at a condition variable which is defined on a task's private data area. When a reply is back and processed by a communication task, the condition variable is signalled and therefore the calling thread is woken up.

### 1.10 ANSAware communication system

AW communication system implements four protocol layers:

- Message passing protocol (MPS): an interface to the transport protocols provided by the underlying operating system
- Execution protocol (EX): implement the invocation of ANSA operations. AW 4.1 supports the Remote Execution Protocol (REX) for point to point invocations and the Group Execution Protocol (GEX) for group invocations

- **Channels/sessions:** used to store the end-to-end state required for a remote invocation and to synchronise the execution of the tasking and the communication systems
- **Stubs:** marshal host language level variables into (and out of) linear communications buffers.

AW communication design is for efficient resource utilization by multiplexing the channels provided by each of these between those of the next layer.

Real-time communication is not a concern of current AW.

#### 1.10.1 Interface reference

An interface reference (*ifref*) contains sufficient information to allow the holder (client) to establish communication with the interface denoted (server). An interface reference has a set of address records, each of which in turn consists a channel id and the network address of the underlying MPS.

#### 1.10.2 Message passing protocol

The MPS interface is stateless and defines an unreliable datagram service. There is no mechanism for QoS based selection/setup of a MPS module. High-level protocols multiplex MPS endpoint whenever possible (in a capsule wide basis).

#### 1.10.3 Execution protocol

The EX interface is also stateless and there is also no mechanism for QoS based selection/setup of a EX module.

REX is designed for asynchronous communication optimized for either low-latency or high throughput. REX provides a rate-based transportation service. The execution reliability semantics of REX calls are exactly-once in the absence of total communication failure.

#### 1.10.4 Channel and session

Channels (i.e. sockets and plugs) are used to store static communication information; sessions duplicate channel state and store additional dynamic information for each invocation.

There is an one-to-one correspondence between channels and ifrefs. Clients send/receive invocation requests/replies through plugs. Servers receive/transmit invocation requests/replies over sockets. The channel id provides an extra layer of demultiplexing on top of the MPS address, so that the capsule can locate the right dispatcher for an specific interface.

There is no interaction between the channel and MPS modules when channels are created and destroyed; the two are independent of each other.

#### 1.10.5 Bindings

A service provider fabricates an ifref by an interface creation operation. A client holding the ifref must then bind to the service in order to communicate with it.

The ifref is created by an implicit binding operation at the server side. The binding operation allocates a socket and concatenates its id with the default

communication addresses (address hint) supported by all MPS in the server to form the interface reference.

Client side binding (the creation of a plug) is performed the first invocation of a service; the first invocation is detected by the absence of the ifref from the ifref to plug cache. Removing an ifref from the cache will force a rebinding.

The cache provides a mapping from an ifref to a plug. The bind operation which allocates a new plug also adds the plug to the cache.

At no stage is there any interaction between the binding process and the EX and MPS modules; it is assumed that all communications between channels is multiplexed over a single MPS address with in a capsule.

---

## 1.11 ANSAware/RT communication system

---

### 1.11.1 QoS and explicit binding

QoS objects are introduced to express different communication requirement and are used by explicit binding operations to create different communication channels.

When creating a service instance, a QoS object is allowed to be associated with the interface creation operation. The operation uses an explicit binding operation to fabricate the result ifref. The explicit binding operation calls the corresponding explicit binding operations in each protocol layer (EX and MPS). The binding operation at each protocol layer interprets the relevant QoS attributes, setups the protocol related binding states, and returns a result that may be used to build the ifref as the result.

In comparison with the implicit binding, the ifref created by an explicit binding contains only these information deduced from the QoS, rather than the default information that provides the maximum communicability and also the maximum multiplexing.

Client side explicit binding is performed by an explicit binding operation, which is also associated with a QoS object. The binding operation, like the server side explicit binding operation, calls the explicit binding operations at each protocol layer with the ifref and the QoS object as arguments. The explicit binding operation at each protocol layer executes a complimentary operation to the relevant QoS and the ifref to create the client site binding. The binding operation creates a plug and adds it in the plug cache, so that later calls on the interface are guaranteed an established channel.

### 1.11.2 State-full message passing protocol

The MPS interface is extended to be state-full: it supports a connection-oriented communication paradigm. The connection is encapsulated as binding informations of a channel. Each channel is associated with a binding data structure which represents end-to-end state establishment with some known channel-specify properties (deduced from a QoS object).

The MPS interface is extended with three operations: server explicit binding operation, client explicit binding operation and binding release operation. The original message *send* and *receive* operations are also extended to make use of the binding informations.

A default binding is established at MPS initialisation time to be used as the default communication channel for the implicitly bound interfaces.

### 1.11.3 State-full execution protocol

EX is modified to use the state-full MPS, and itself is redesigned to state-full as well. This allows the addition of extra binding information to the binding created by MPS to include EX dependent data and state. For example, the binding contains extra information about the header size of an EX protocol which can be used by MPS to fetch the correct EX-dependent packet headers.

The EX interface is extended with three operations: server explicit binding operation, client explicit binding operation and binding release operation. The original message *send* and *receive* operations are also extended to make use of the binding informations.

### 1.11.4 Timed remote execution protocol

The generic design of the binding and state-full protocols allows the insertion of new EX protocols. The TRPC is implemented as one example.

Timed Remote Execution Protocol (TRES) is a cut-down version of REX as follows:

- no fragmentation, this has significantly reduced the size of the protocol
- at-most-once semantics, no timeout controlled retry
- no security check, no passing and checking of nonce
- small header size, the result of the above three design choose

TRES also extends REX in the following aspects to implement the semantics of TRPC:

- the header of packages is expanded to include the information about the priority, deadline and deadline type
- extended session functions to process timeout at client side and deadline expires at server side
- extra message types for handling deadline exception and confirmation.

TRES supports only explicit binding operations, i.e. interfaces created by implicit binding operations will not be able to use TRES.

### 1.11.5 In-band QoS

In-band QoS is allowed to be associated to each invocation to select the dynamic QoS parameters of a channel. Currently, if a channel uses TRES as its EX, the in-band QoS can select a priority, a timeout, a deadline and a deadline type.

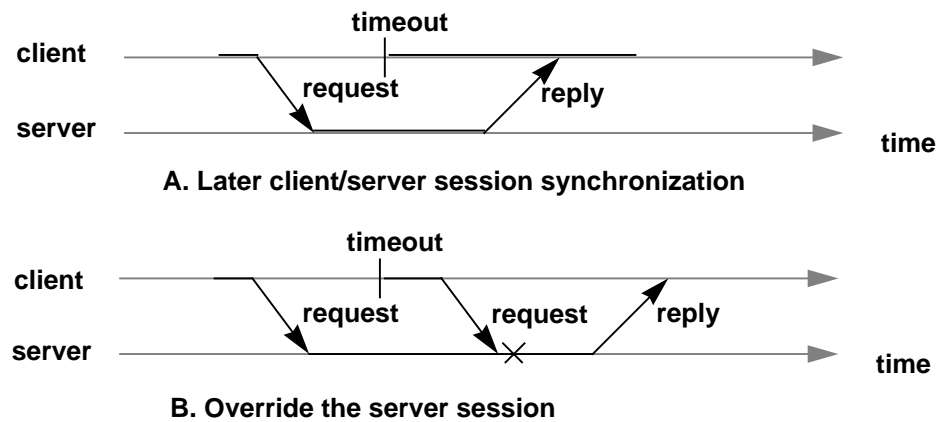
In-band QoS are associated with sessions by the interpreter at the client side, and by the TRES at the server side.

### 1.11.6 Session overridden

Timeout at client sides are a mixed blessing: the desired semantics of a timeout is when it expires the client should resume control (so that the client can take some immediate recovery actions). However, the operation is still carried on at the server side and extra packet exchange is required to synchronise the client and server sessions. If the packet exchange takes place

at the timeout expiry time, the extra overhead of synchronisation may lead to uncertain timeout semantics. Therefore, an alternative approach is pursued.

**Figure 1.5: Session timeout recovery**



The ANSAware/RT approach for session timeout recovery is illustrated in Figure 4.1. With this approach, the client continues immediately after the timeout, and the client session is set to idle. No synchronisation packet exchange is initiated by the client. It allows the existence of inconsistency between a client session and its server side session. Should the server returns an obsolete result later, synchronisation of the client and server sessions are taken then. The approach also allows the server side session to be aware that its client side may timeout, and the client side session may be used for another invocation. A possible effect (caused by the ANSAware approach to session management) is that a later invocation from the same client side session may override a server side session representing an obsolete invocation.

### 1.12 Performance evaluation

There are several standard synthetic benchmarks for real-time computing systems, including Hartstone Benchmark (HB) [Weiderman 89], Distributed Hartstone Benchmark (DHB) [Mercer et al. 90] and Hartstone Distributed Benchmark (HDB) [Kamenoff and Weiderman 91]. The HB is a set of timing requirements for testing a system's ability to handle hard real-time applications. It is specified as a set of tasks with well-defined workload and timing constraints. It is a benchmark for single processor machines. The DHB and HDB are both extensions of HB for distributed real-time systems. They are designed to give figures of merit for the complex end-to-end scheduling and timing behaviour of the system. In comparison, the HDB gives a broader definition and merit of real-time distributed systems' behaviour, while the DHB has a concrete definition of the series of tests.

DHB was chosen to measure and evaluate ANSAware/RT performance. The intention of DHB is to measure the real-time performance of the processor scheduling, the communication network scheduling and the coordination between these scheduling domains. It is argued that since more sophisticated scheduling algorithms may require more overhead for low-level operations, a system which offers better schedulability for its applications and thus better overall performance may not have the best times for low-level operations. A system which is leaner and faster in terms of low-level operations may not be

capable of scheduling a task set to meet all of its deadlines. DHB is thus designed to factor all of these attributes into the overall evaluation of a system.

DHB defines five sets of experiments based on a typical client/server interaction model. The task workload is expressed in Kilo-Whetstone (KWS) or milliseconds. A KWS is one execution of a mathematical library, which factors out the effect of a typical arithmetic computing. Each experiment starts with five periodic client tasks and one or two server tasks. Each task is required to execute a specific amount of workload within its period. The experiment continues with added workload until any task's deadline cannot be met. The five sets of experiments are:

- DSHcl, a Distributed, Synchronized, and Harmonic task set which tests the communication latency of the system. The server computation time is increased in milliseconds to squeeze the time left for message passing.
- DSHpq, a Distributed, Synchronized, and Harmonic task set. The server computation time is increased in KWS to measure how well the system is in priority queuing of communication packets
- DSNpp, a Distributed, Synchronized, and Non-harmonic task. The number of low-priority client tasks are increased to test the degree of preemptability of the protocol engines.
- DSHcb, a Distributed, Synchronized, and Harmonic task set. The number of high-priority client tasks is increased to test the bandwidth for real-time communications.
- DSHmc series, a task set for measuring media contention, which does not apply to the Ethernet.

To achieve comparable results, DHB was executed on ANSAware/RT by using of a 10 Mbps Ethernet and two DEC Alpha 3000/300 workstations.

To make a comparison, the relevant performance of the ARTS distributed real-time operating system and RIDE system are also given in Table 1.1. The ARTS performance is copied from [Mercer et al. 90] which was measured by using SUN3/140s and a private 10 Mbps Ethernet.

**Table 1.1: RIDE, ANSAware/RT vs ARTS performance**

Series	ARTS SUN 3/140	RIDE DEC Firefly	ANSAware/RT DEC Alpha 3000/300
DSHcl	35 ms	26 ms	41 ms
DSHpq	18 KWS	16 KWS	2010 KWS
DSHpp	(13) 20 tasks	18 tasks	105 tasks
DSHcb	14 tasks	15 tasks	23 tasks

Comparison of the performance of ANSAware/RT, RIDE and ARTS is, however, not as simple as it looks. The ARTS system uses kernel supported objects, object invocations, and preemptive protocol processing; while ANSAware/RT/RIDE uses a relatively heavyweight user level RPC mechanism. In RPC systems, the marshalling and un-marshalling of arguments, the overhead of an RPC protocol, the multiplexing of a required operation within an interface, and the demultiplexing of replies for clients are time consuming. Taking these into account, it is reasonable that RIDE is 9 ms

less efficient in the DSHcl series test (which tests communication latency). On the other hand, RIDE performs as well as ARTS in the DSHpq, DSNpp and DSHcb series tests. That is, RIDE can achieve about the same performance as ARTS in the priority queuing of communication packets, in the preemptability of the protocol engine, and in the provision of communication bandwidth.

The much better performance of ANSAware/RT reflects the combination effects of the superiority of a commercial real-time operating system, a much powerful processor and a carefully tuned mechanisms based on the practical experience of RIDE.

---

### 1.13 Summary

---

The paper reviews some of the main features of the real-time ANSAware (ANSAware/RT) version 1.0. ANSAware/RT provides a framework to facilitate the enforcement of stringent timing constraints found in distributed real-time applications. The ANSAware/RT design incorporates tasks and communication channels (the two most important resources in real-time distributed computing) as its basic programming components. It synthesises aspects of resource requirements, resource allocation and resource scheduling into an object-based programming paradigm. Predictability, user control and mission criticality are the main characteristics of the model.

The performance of the ANSAware/RT implementation is compared to that of some typical systems by using of the Distributed Hartstone Benchmarks, and has shown that the design is viable.

---

### 1.14 Related work

---

The general discussion of an open system architecture for real-time processing can be found in [Li and Otway 94]. The description of the ANSAware/RT programming interfaces can be found in [Li 94a]. Comparison of the two ANSA based real-time environments ANSAware/RT and RIDE is described in [Li 94].

Current research at CNET [Hazard et al. 93] and Lancaster University [Coulson et al. 92] are all converging on a common architecture for distributed multimedia and real-time processing relevant to the ANSAware/RT system.

---

### 1.15 Acknowledgement

---

The author of this article would like to acknowledge the contribution of his colleagues in the ANSA core team, particularly Andrew Herbert, Dave Otway and John Warne for their valuable comments and suggestions.





---

## References

---

[APM 92]

APM Ltd., ANSAware Version 4.1 Manual, Architecture Projects Management Ltd., Cambridge U.K., May 1992.

[Coulson et al. 92]

G Coulson et al. Extensions to ANSA for Multimedia Computing, Computer Networks and ISDN Systems, 25, pp 305 - 323, 1992.

[Gopinath and Bihari 93]

P Gopinath and T Bihari, Concepts and Examples of Object-Oriented Real-Time Systems, In Readings in Real-Time systems, Y H Lee and C M Krishna ed., pp 123-136, IEEE CS Press, June 1993.

[Hazard et al. 93]

L Hazard et al. Towards the Integration of Real Time and QoS Handling in ANSA Architecture, ANSA Phase 3 Project Report CNET/RC.ARCADÉ.01, June 1993.

[Herbert 94]

A Herbert, An ANSA Overview, IEEE Network, pp 18-23, January 1994.

[ISO 94]

ISO/IEC 10746-3, ITU-TS Recommendation X.903: Basic Reference Model of Open Distributed Processing: Prescriptive Model, (2nd CD draft) April 1994.

[Kamenoff and Weiderman 91]

N I Kamenoff and N H Weiderman, Hartstone Distributed Benchmark: Requirements and Definitions, Proc. of Twelfth IEEE Real-Time Systems Symposium, 1991.

[Lee and Davidson 90]

I Lee and S B Davidson, A Performance Analysis of Timed Synchronous Communication Primitives, IEEE Transactions on Computers, Vol. 39, No. 9, pp 1117 -1131, September 1990.

[Li and Bacon 94]

G Li and J Bacon, Supporting Distributed Real-Time Objects, in IEEE Proceedings of the Second Workshop on Parallel and Distributed Real-Time Systems, pp 138 - 143, Cancun, Mexico, April, 1994.

[Li and Otway 94]

G Li and D Otway, An Open Architecture for Real-Time Processing, to appear in ICL Technical Journal, November, 1994.

[Li 94]

G Li, Distributing Real-Time Objects: the ANSA Approach, to appear in Proceedings of IEEE CS 1st Workshop on Object-Oriented Real-Time Dependable Systems, Dana Point, California, October, 1994.

[Li 94a]

G Li, Real-Time ANSAware Version 1.0: Programming and System Overview, APM document 1207, Architecture Projects Management Ltd., Cambridge U.K., May 1994.

[Li 93]

G Li, Supporting Distributed Realtime Computing, PhD thesis, University of Cambridge Computer Laboratory, Technical Report 322, August 1993.

[Mercer et al. 90]

C W Mercer and Y Ishikawa and H Tokuda, Distributed Hartstone: A Distributed Real-Time Benchmark Suite, International Conference on Distributed Computing Systems, 1990.

[Weiderman 89]

N Weiderman, Hartstone: Synthetic Benchmark Requirements for Hard Real-Time Applications, Software Engineering Institute, Carnegie Mellon University, Technical Report CMU/SEI-89-TR-23, June, 1989.