



The Design and Performance of a Real-Time CORBA Scheduling Service

CHRISTOPHER D. GILL*

Department of Computer Science, Washington University, St. Louis, MO 63130, USA

cdgill@cs.wustl.edu

DAVID L. LEVINE

Department of Computer Science, Washington University, St. Louis, MO 63130, USA

levine@cs.wustl.edu

DOUGLAS C. SCHMIDT

Electrical and Computer Engineering Department, University of California, Irvine

schmidt@uci.edu

Abstract. There is increasing demand to extend CORBA middleware to support applications with stringent quality of service (QoS) requirements. However, conventional CORBA middleware does not define standard features to dynamically schedule operations for applications that possess deterministic real-time requirements. This paper presents three contributions to the study of real-time CORBA operation scheduling strategies.

First, we document our evolution from static to dynamic scheduling for applications with deterministic real-time requirements. Second, we describe the flexible scheduling service framework in our real-time CORBA implementation, TAO, which supports core scheduling strategies efficiently. Third, we present results from empirical benchmarks that quantify the behavior of these scheduling strategies and assess the overhead of dynamic scheduling in TAO. Our empirical results using TAO show that dynamic scheduling of CORBA operations can be deterministic and can achieve acceptable latency for operations, even with moderate levels of queuing.

Keywords: middleware and APIs, quality of service issues, mission critical/safety critical systems, dynamic scheduling algorithms and analysis, distributed systems

1. Introduction

1.1. Motivation

Supporting the quality of service (QoS) demands of next-generation real-time applications requires object-oriented (OO) middleware that is flexible, efficient, predictable, and convenient to program. Applications with deterministic real-time requirements, such as process control and avionics mission computing systems (Levine, Gill, and Schmidt, 1998), impose severe constraints on the design and implementation of real-time OO middleware. For example, avionics mission computing applications typically manage sensors and operator displays, and control on-board equipment. Middleware for such applications must support deterministic real-time QoS requirements.

* This work was supported in part by Boeing, DARPA contract 9701516, Lucent, Motorola, NSF grant NCR-9628218, Siemens, and Sprint.

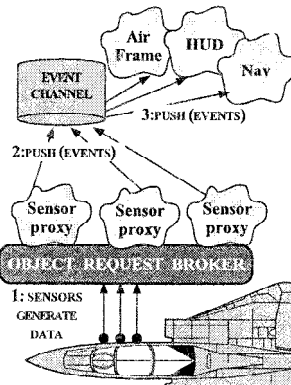


Figure 1. Example avionics mission computing application.

1.2. Design and Implementation Challenges

Figure 1 illustrates the architecture of a representative real-time application—an OO avionics mission computing platform (Harrison, Levine, and Schmidt, 1997)—developed and deployed using OO middleware components and services based on CORBA (Object Management Group, 1998). CORBA Object Request Brokers (ORB), allow clients to invoke operations on target object implementations without concern for where the object resides, what language the object implementations are written in, the OS/hardware platform, or the types of communication protocols, networks, and buses used to interconnect distributed objects (Vinoski, 1997). However, achieving these benefits for deterministic real-time applications requires the resolution of the following design and implementation challenges:

Scheduling assurance prior to run-time: In some real-time applications, the consequences of missing a deadline at run-time can be catastrophic. For example, failure to process an input from the pilot by a specified deadline can be disastrous in an avionics application, especially in mission critical situations. Therefore, it is essential to validate *prior to run-time* that all critical processing deadlines will be met.

Severe resource limitations: Many real-time applications must minimize processing due to strict resource constraints, such as cost, weight, and power consumption restrictions. A consequence of using off-line scheduling analysis is that worst-case processing requirements drive the schedule. Therefore, resource allocation and scheduling must always accommodate the worst case, even if non-worst case scenarios are common. For example, an application that relies on real-time image processing (Pyarali, Harrison, and Schmidt, 1996) may have to assume that such processing will take some maximum amount of time, when often it may take much less.

Distributed processing: Clients running on one processor must be able to invoke operations on servants on other processors. Likewise, the allocation of operations to processors should be flexible. For instance, it should be transparent to the application design and implementation whether an operation resides on the same processor as the client that invokes it.

Testability: Real-time software is complex, critical, and long-lived. Therefore, maintenance is often problematic and expensive (Newport, 1994). A large percentage of software maintenance involves testing. Current scheduling approaches are validated by extensive testing, which is tedious and non-comprehensive. Therefore, analytical assurance is essential to help reduce validation costs by focusing the requisite testing on the most strategic system components.

Adaptability across product families: Some real-time systems are custom-built for specific product families. Development and testing costs can be reduced if large, common components can be factored out. In addition, validation and certification of components can be shared across product families, amortizing development time and effort.

1.3. Applying CORBA to Deterministic Real-Time Applications

Our experience using CORBA on telecommunication (Schmidt, 1996) and medical imaging projects (Pyarali et al., 1997) illustrates that it is well-suited for conventional request/response applications with “best-effort” QoS requirements. Moreover, CORBA addresses issues of distributed processing and adaptation across product families by promoting the separation of interfaces from implementations and supporting component reuse (Vinoski, 1997).

However, conventional CORBA ORBs are not yet suited for demanding real-time applications because they do not provide features or optimizations to schedule operations that require deterministic real-time QoS (Object Management Group, 1998). To meet these requirements, we have developed a real-time CORBA Object Request Broker (ORB) called TAO (Schmidt, Levine, and Mungee, 1998). TAO is an open source, implementation of standard CORBA whose ORB and services support efficient and predictable real-time, distributed object computing. Our prior work on TAO has explored many dimensions of high-performance and real-time ORB design and performance, including event processing (Harrison, Levine, and Schmidt, 1997), request demultiplexing (Gokhale et al., 1999), I/O subsystem integration (Kuhns et al., 1999), concurrency and connection architectures (Schmidt et al., 1999), and IDL complier stub/skeleton optimizations (Gokhale and Schmidt, 1998). This paper extends our previous work on a real-time CORBA static scheduling service (Schmidt, Levine, and Mungee, 1998) by incorporating a *strategized scheduling service framework* into TAO. This framework allows the configuration and empirical evaluation of multiple static, dynamic, and hybrid static/dynamic scheduling strategies, such as Rate Monotonic Scheduling (RMS) (Liu and Layland, 1973), Earliest Deadline First (EDF) (Liu and Layland, 1973), Minimum Laxity First (MLF) (Stewart and Khosla, 1992), and Maximum Urgency First (MUF) (Stewart and Khosla, 1992).

To maintain scheduling guarantees and to simplify testing for demanding real-world real-time applications, we have extended our prior work on TAO incrementally. In particular, our approach focuses on applications with the following characteristics: (1) bounded executions—operations stay within the limits of their specified execution times; (2) bounded rates—dispatch requests arrive within the limits of their specified execution times; (2) bounded rates—dispatch requests arrive within the specified period; (3) known operations—all operations are known to the scheduler before run-time, or are reflected entirely within the execution times of other specified operations. These types of applications are historically configured and scheduled *statically*, which enables TAO to minimize run-time overhead that would otherwise stem from mechanisms that enforce operation execution time limits (Harrison, Levine, and Schmidt, 1997) or perform dynamic admission control.

Within these constraints, the work on TAO's strategized scheduling service framework described in this paper allows applications to specify custom static and/or dynamic scheduling strategies. This framework increases adaptability across application families and operating systems, while preserving the rigorous scheduling guarantees and testability offered by our previous work on statically scheduled CORBA operations.

1.4. Paper Organizations

The remainder of this paper is organized as follows: Section 2 reviews the drawbacks of off-line, static scheduling and introduces the dynamic and hybrid static/dynamic scheduling strategies our research is evaluating. Section 3 discusses the design and implementation of TAO's scheduling service framework, which supports a range of static, dynamic, or hybrid static/dynamic real-time scheduling strategies. Section 4 presents results from benchmarks that evaluate the dynamic scheduling strategies empirically to compare the run-time dispatching overhead of static and dynamic scheduling strategies. Section 5 discusses related work and Section 6 presents concluding remarks.

2. Overview of Scheduling Strategies

This section describes the limitations of purely static scheduling and outlines the potential benefits of applying dynamic scheduling. In addition, we evaluate the limitations of purely dynamic scheduling strategies. This evaluation motivates the hybrid static/dynamic strategy used by TAO to schedule real-time CORBA operations, which is described in Section 3.

2.1. Synopsis of Scheduling Terminology

Precise terminology is necessary to describe and evaluate static, dynamic, and hybrid scheduling strategies. Figure 2 shows the relationships between the key terms defined below.

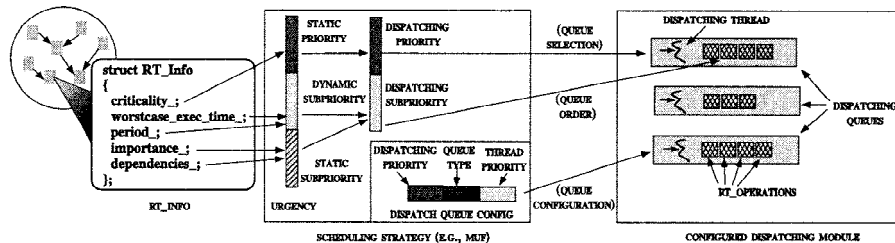


Figure 2. Relationships between operation, scheduling, and dispatching terminology.

RT-Operation and RT-Info: In TAO, an *RT-Operation* is a scheduled CORBA operation (Schmidt, Levine, and Mungee, 1998). In this paper, we use *operation* interchangeably with *RT-Operation*. An *RT-Info* struct is associated with each operation and contains its QoS parameters. Section 3.4.1 describes these concepts in more detail. The *RT-Info* structure contains the following operation characteristics shown in Figure 2 and described below:

- **Criticality:** Criticality is an application-supplied value that indicates the significance of a CORBA operation’s completion prior to its deadline. Higher criticality should be assigned to operations that incur greater cost to an application if they fail to complete execution before their deadlines. Some scheduling strategies, such as MUF, give greater priority to more critical operations than to less critical ones.

- **Worst-case execution time:** This is the longest time required to execute a single dispatch of an operation. Worst case execution times may be determined through techniques like simulation, instruction counting, or benchmarking on the target platform.

- **Period:** Period is the interval between dispatches of an operation.

- **Importance:** Importance is a lesser indication of a CORBA operation’s significance. Like its criticality, an operation’s importance value is supplied by an application. Importance is used as a “tie-breaker” to assign a unique static subpriority for each operation.

- **Dependencies:** An operation may *depend* on data produced by another operation. An operation that depends on the data from another operation may execute only after the other operation has completed.

Scheduling strategy: A scheduling strategy transforms the information from an operation’s *RT-Info* by (1) assigning an *urgency* to the operation based on its static priority, dynamic subpriority, and static subpriority values, (2) mapping urgency into dispatching priority and dispatching subpriority values for the operation, and (3) providing dispatching queue configuration information so that each operation can be dispatched according to its assigned dispatching priority and dispatching subpriority. These concepts are discussed in more detail in Sections 3.4.3 and 3.4.4. The key elements of this transformation are shown in Figure 2 and defined as follows:

- **Urgency:** Urgency (Stewart and Khosla, 1992) is an ordered tuple consisting of (1) static priority, (2) dynamic subpriority, and (3) static subpriority. Static priority is the highest ranking priority component in the urgency tuple, followed by dynamic subpriority and then static subpriority, respectively.

- **Static priority:** Static priority assignment establishes a fixed number of priority partitions into which all operations must fall. The number of static priority partitions is established off-line. An operation's static priority value is often determined off-line. However, the value assigned a particular dispatch of the operation could vary at run-time, depending on which scheduling strategy is employed.

- **Dynamic subpriority:** Dynamic subpriority is a value generated and used at run-time to order operations *within* a static priority level, according to the run-time and static characteristics of each operation. For example, a subpriority based on the operation with the "closest deadline" must be computed dynamically.

- **Static subpriority:** Static subpriority values are determined prior to run-time. Static subpriority acts as a tie-breaker when both static priority and dynamic subpriority are equal.

- **Dispatching priority:** An operation's dispatching priority corresponds to the real-time priority of the thread in which it will be dispatched. Operations with higher dispatching priorities are executed in threads with higher real-time priorities.

- **Dispatching subpriority:** Dispatching subpriority is used to order operations within a dispatching priority level. Operations with higher dispatching subpriority are executed ahead of operations with the same dispatching priority, but with lower dispatching subpriority.

- **Queue configuration:** A separate queue must be configured for each distinct dispatching priority. The scheduling strategy assigns each queue a dispatching type, e.g., static, deadline, or laxity; a dispatching priority; and a thread priority.

Together, urgency and dispatching (sub)priority assignment specify requirements that certain operations will meet their deadlines. To support end-to-end QoS requirements, operations with higher dispatching priorities *should not* be delayed by operations with lower dispatching priorities. Two research challenges must be resolved to achieve this goal: (1) strategies must be identified to correctly specify end-to-end QoS requirements for different operations and (2) dispatching modules must enforce these end-to-end QoS specifications. The following two definitions are useful in addressing these challenges:

- **Critical set:** The critical set consists of all operations whose completion prior to deadline is crucial to the integrity of the system. If all operations in the critical set can be assured of meeting their deadlines, a schedule that preserves the system's integrity can be constructed.

- **Minimum critical priority:** The minimum critical priority is the lowest dispatching priority level to which operations in the critical set are assigned. Depending on the scheduling strategy, the critical set may span multiple *dispatching priority* levels. To ensure that the critical set is schedulable, all operations at the minimum critical priority level must be schedulable.

TAO's scheduling strategies rely primarily on priority- and subpriority-based dispatching, which can be enforced efficiently either by mechanisms available in the OS kernel (e.g., preemptive thread priorities) or can be implemented efficiently in middleware (e.g., dynamic subpriorities). Other scheduling strategies, such as Time-based Scheduling (Wang, Wang, and Lin, 1999) and FIFO-r (Tyan and Hou, 1999), use additional characteristics to order the dispatches of operations. These characteristics include:

- **Resource share:** Resource share is a measure of an operation's appropriate share of a resource (e.g., CPU time), and is used to ensure fairness among operations that are not otherwise prioritized. For example, a share-based scheduling strategy might maintain information about each operation's past execution time. This information could be used to dispatch operations so that within every priority level each operation consumes CPU time proportional to its fair share.
- **Timing constraints:** Timing constraints capture explicit requirements for operation dispatch and completion times. For example, a timing constraint might specify that an operation must be dispatched within T time units after another operation completes.
- **Dispatching module:** A dispatching module is responsible for (1) constructing the appropriate type of queue for each dispatching priority and (2) assigning each dispatching thread's priority to the value provided by the scheduling strategy. A TAO ORB endsystem can be configured with dispatching modules at several layers, including the I/O subsystem (Kuhns et al., 1999), ORB Core (Schmidt et al., 1999), and/or the Event Service (Harrison, Levine, and Schmidt, 1999). TAO's dispatching modules are discussed in Section 3.4.6.

2.2. *Limitations of Static Scheduling*

Many hard real-time systems, such as those for avionics mission computing and manufacturing process controllers, have traditionally been scheduled statically using RMS (Klein et al., 1993). Static scheduling provides schedulability assurance prior to run-time and can be implemented with low run-time overhead (Schmidt, Levine, and Mungie, 1998). However, static scheduling has the following disadvantages:

Utilization phasing penalty for non-harmonic periods: In statically scheduled systems, achievable utilization can be reduced if the periods of all operations are *not* related harmonically. Operations are related harmonically if their periods are integral multiples of one another. When periods are not harmonic, the phasing of the operations produces unscheduled gaps of time. This reduces the maximum schedulable percentage of the CPU, i.e., the *schedulable bound*, to below unity. The *utilization phasing penalty* is the difference between the value of the schedulable bound and 100%. Liu and Layland established a least upper utilization bound of $n(2^{1/n} - 1)$ (1973) for a set of operations to be statically schedulable, where n is the number of distinct non-harmonic operation periods in the system. Recent research has shown that this bound may be overly pessimistic in some cases (Han and Tyan, 1997). However, the fact remains that with static priority assignment *some* unschedulable gaps may be created by non-harmonic periods.

Inflexible handling of invocation-to-invocation variation in resource requirements:

Because priorities cannot be changed easily at run-time, allocations must be based on worst-case assumptions.¹ Thus, if an operation usually requires 5 msec of CPU time, but under certain conditions requires 8 msec, static scheduling analysis must assume that 8 msec will be required for every invocation. Again, utilization is effectively penalized because the resource will be idle for 3 msec in the usual case.

Impact of situational factors on resource requirements: Recent advances in static priority analysis (Mok and Chen, 1997; Han, 1998) have shown that the schedulable bound for statically prioritized operations can be improved dramatically in some cases. These techniques rely, however, on advance knowledge of (1) arrival patterns of operation dispatch requests and (2) sequences of operation execution times. In many distributed real-time applications, such as those for avionics mission computing and image processing, variation in load on the system is largely due to *situational* factors. Thus, such detailed information may not be available accurately prior to run-time.

In general, static scheduling limits the ability of real-time systems to adapt to changing conditions and changing configurations. In addition, static scheduling provides resource access guarantees at the cost of lower resource utilization. To overcome the limitations of static scheduling, therefore, we have investigated the use of dynamic strategies to schedule CORBA operations for applications with real-time QoS requirements.

2.3. *Overcoming Static Scheduling Limitations with Dynamic Scheduling*

Dynamic scheduling offers a way to address the drawbacks described in Section 2.2. In particular, dynamic scheduling strategies offer optimal utilization capabilities (Liu and Layland, 1973) and handle invocation-to-invocation variations in execution times efficiently. If these drawbacks can be alleviated without incurring excessive overhead or non-determinism, dynamic scheduling can be beneficial for real-time applications with deterministic QoS requirements.

Demanding real-time applications, such as avionics mission computing, cannot tolerate unnecessary overhead and non-determinism at run-time. Therefore, we restrict our attention in this paper to scheduling approaches that do not perform schedulability analysis at run-time. In particular, we do not consider strategies that require run-time admission control for dynamic scheduling. Rather, we only consider scheduling strategies where it is possible to select the set of operations critical to the application *statically*. Among such strategies, we are most interested in those whose *dynamic* run-time behavior allows maximal resource utilization.

Unfortunately, many dynamic scheduling strategies do not offer the *a priori* guarantees of static scheduling. For instance, purely dynamically scheduled systems, i.e., those without any form of admission control for dynamically generated operations, can behave non-deterministically under heavy loads. Thus, operations that are critical to an application may miss their deadlines because they were (1) delayed by non-critical operations or (2) delayed by an excessive number of critical operations.

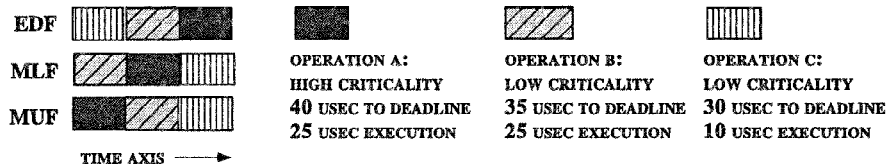


Figure 3. Dynamic scheduling strategies.

The remainder of this section reviews several strategies for dynamic and hybrid static/dynamic scheduling, using the terminology defined in Section 2.1. These scheduling strategies include purely dynamic techniques, such as EDF and MLF, as well as the hybrid MUF strategy.

2.3.1. Purely Dynamic Scheduling Strategies

This section briefly reviews two well known purely dynamic scheduling strategies, EDF (Liu and Layland, 1973; Klein et al., 1993) and MLF (Stewart and Khosla, 1992). These strategies are illustrated in Figure 3 and discussed below. In addition, Figure 3 depicts the hybrid static/dynamic MUF (Stewart and Khosla, 1992) scheduling strategy discussed in Section 2.3.2.

Earliest Deadline First (EDF): EDF (Liu and Layland, 1973; Klein et al., 1993) is a dynamic scheduling strategy that orders dispatches of operations based on time-to-deadline, as shown in Figure 3.² Operation executions with closer deadlines are dispatched before those with more distant deadlines. The EDF scheduling strategy is invoked whenever a dispatch of an operation is requested. Depending on the mapping of priority components into thread priorities, the new dispatch may or may not preempt the currently executing operation, as discussed in Section 3.4.4.

A key limitation of EDF is that an operation with the earliest deadline is dispatched, whether or not there is sufficient time remaining to complete its execution prior to the deadline. Therefore, the fact that an operation cannot meet its deadline will not be detected until *after* the deadline has passed. Moreover, that operation will continue to consume CPU time that could otherwise be allocated to other operations that could still meet their deadlines.

Minimum Laxity First (MLF): MLF (Stewart and Khosla, 1992) refines the EDF strategy by taking into account operation execution time. It dispatches an operation whose *laxity* is least, as shown in Figure 3. Laxity is defined as the time-to-deadline minus the remaining execution time.

Using MLF, it is possible to detect that an operation will not meet its deadline *prior* to the deadline itself. If this occurs, a scheduler can reevaluate the operation before allocating the CPU, as discussed for the MUF scheduling strategy in Section 2.3.2.

Evaluation of EDF and MLF:

- **Advantages:** From a scheduling perspective, the main advantage of EDF and MLF is that they overcome the utilization limitations of RMS. In particular, the utilization phasing penalty described in Section 2.2 cannot occur with EDF and MLF, because they assign priorities based on run-time characteristics. In addition, EDF and MLF handle harmonic and nonharmonic periods comparably. Moreover, they respond flexibly to invocation-to-invocation variations in resource requirements, allowing CPU time unused by one operation to be reallocated to other operations. Thus, they can produce schedules that are optimal in terms of CPU utilization (Liu and Layland, 1973). Finally, both EDF and MLF can dispatch operations within a single static priority level (Liu and Layland, 1973; Stewart and Khosla, 1992) which is useful for non-preemptive single-threaded environments.

- **Disadvantages:** From a performance perspective, a disadvantage of purely dynamic scheduling approaches like MLF and EDF is that their scheduling strategies require higher overhead to evaluate at run-time. In addition, these purely dynamic scheduling strategies offer no control over *which* operations will miss their deadlines if the schedulable bound is exceeded. As operations are added to the schedule to achieve higher utilization, the margin of safety for *all* operations decreases. As the system becomes overloaded, therefore, the risk of missing a deadline increases for every operation.

2.3.2. Maximum Urgency First (MUF)

The MUF (Stewart and Khosla, 1992) scheduling strategy supports the deterministic rigor of the static RMS scheduling approach *and* the flexibility of dynamic scheduling approaches like EDF and MLF. MUF is the default scheduler for the Chimera real-time operating system (Stewart, Schmitz, and Khosla, 1992). TAO supports a variant of MUF in its strategized CORBA scheduling service framework, which is presented in Section 3. MUF can assign both static *and* dynamic priority components. In contrast, RMS assigns all priority components statically based on fixed rates and EDF/MLF assign all priority components dynamically based on deadlines/laxities. The hybrid priority assignment in MUF overcomes the drawbacks of the individual scheduling strategies by combining techniques from each, as described below:

Criticality: In MUF, operations with higher *criticality* are assigned to higher static priority levels. Assigning static priorities according to criticality prevents operations critical to the application from being preempted by non-critical operations. Ordering operations by application-defined criticality reflects a subtle and fundamental shift in the notion of priority assignment. In particular, RMS, EDF, and MLF exhibit a rigid mapping from *empirical* operation characteristics to a *single* priority value.

Moreover, EDF and MLF offer little or no control over which operations will miss their deadlines under overload conditions. In contrast, MUF affords applications the ability to distinguish operations arbitrarily, giving them *explicit* control over *which* operations will miss their deadlines under conditions of overload. Therefore, it can protect a critical *subset* of the entire set of operations. This fundamental shift in the notion of priority assignment leads to the generalization of scheduling techniques discussed in Section 3.

Dynamic Subpriority: An operation's dynamic subpriority is evaluated whenever it is enqueued in or dequeued from a dynamically ordered dispatching queue. At the instant of evaluation, dynamic subpriority in MUF is a function of the laxity of an operation. By assigning dynamic subpriorities according to laxity, MUF offers higher utilization of the CPU than the statically scheduled strategies. MUF also allows deadline failures to be detected *before* they actually occur, except when an operation that would otherwise meet its deadline is preempted by a higher criticality operation. Moreover, MUF can apply various types of error handling policies when deadlines are missed (Stewart and Khosla, 1992). For example, if an operation has negative laxity prior to being dispatched, it can be diverted from the dispatching queue. This allows operations that can still meet their deadlines to be dispatched instead.

Static Subpriority: In MUF, *static subpriority* is a static, application-specific, optional value. It is used to order the dispatches of operations that have the same criticality and the same dynamic subpriority. Thus, static subpriority has lower precedence than either criticality or dynamic subpriority. Assigning a unique static subpriority allows a total dispatch ordering of operations at run-time. For a given arrival pattern of operation requests, the total ordering ensures that the dispatch order will always be identical. This assurance improves system predictability, reliability, and testability. The variant of MUF used in TAO's strategized scheduling service enforces a total dispatch ordering by providing an importance field in the TAO RT-Info CORBA operation QoS descriptor (Schmidt, Levine, and Mungee, 1998), which is described in Section 2.1. TAO's scheduling service uses importance, as well as a topological ordering of operations, to assign a unique static subpriority for each operation within a given criticality level.

Incidentally, the original definition of MUF in Stewart and Khosla (1992) uses the terms *dynamic priority* and *user priority*, whereas we use the term *dynamic subpriority* and *static subpriority* for TAO's scheduling service, respectively. We selected different terminology to indicate the subordination to static priority. These terms are interchangeable when referring to the MUF strategy, however.

3. The Design of TAO's Strategized Scheduling Service

TAO's scheduling service provides real-time CORBA applications the flexibility to specify and use different scheduling strategies, according to their specific QoS requirements and available OS features. This flexibility allows CORBA applications to extend the set of available scheduling strategies *without* impacting strategies used by other applications. Moreover, it shields application developers from unnecessary details of their scheduling strategies.

This section outlines the design goals and architecture of TAO's strategized scheduling service framework. After briefly describing TAO in Section 3.1, Section 3.2 discusses the design goals of TAO's strategized scheduling service. Section 3.3 offers an overview of its architecture and operation. Finally, Section 3.4 discusses the resulting architecture in detail.

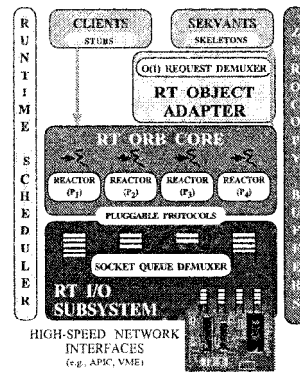


Figure 4. Components in the TAO real-time ORB endsystem.

3.1. Overview of TAO

TAO is a high-performance, real-time ORB endsystem targeted for applications with deterministic QoS requirements, as well as “best-effort” requirements. The TAO ORB endsystem contains the network interface, OS, communication protocol, and CORBA-compliant middleware components and features shown in Figure 4. TAO supports the standard OMG CORBA reference model (Object Management Group, 1998), with the following enhancements designed to overcome the shortcomings of conventional ORBs (Schmidt et al., 1999) for high-performance and real-time applications:

Optimized IDL Stubs and Skeletons: IDL stubs and skeletons perform marshaling and demarshaling of application operation parameters, respectively. TAO’s IDL compiler generates stubs-skeletons that can selectively use highly optimized compiled and/or interpretive marshaling/demarshaling (Gokhale, Schmidt, Levine, and Mungee, 1998). This flexibility allows application developers to selectively trade off time and space, which is crucial for high-performance, real-time, and/or embedded distributed systems.

Real-time Object Adapter: An Object Adapter associates servants with the ORB and demultiplexes incoming requests to servants. TAO’s real-time Object Adapter (Gokhale et al., 1999) uses perfect hashing (Schmidt, 1990) and active demultiplexing (Gokhale et al., 1999) optimizations to dispatch servant operations in constant $O(1)$ time, regardless of the number of active connections, servants, and operations defined in IDL interfaces.

Run-time Scheduler: A real-time scheduler (Object Management Group, 1998) maps application QoS requirements, such as bounding end-to-end latency and meeting periodic scheduling deadlines, to ORB endsystem/network resources, such as CPU, memory, network connections, and storage devices. TAO’s run-time scheduler supports static (Schmidt, Levine, and Mungee, 1998) real-time scheduling, as well as the dynamic and hybrid static/dynamic real-time scheduling strategies described in this paper.

Real-time ORB Core: An ORB Core delivers client requests to the Object Adapter and returns responses (if any) to clients. TAO's real-time ORB Core (Schmidt et al., 1999) uses a multi-threaded, preemptive, priority-based connection and concurrency architecture (Gokhale and Schmidt, 1998) to provide an efficient and predictable CORBA protocol engine. TAO's ORB core allows customized protocols to be plugged into the ORB without affecting the standard CORBA application programming model.

Real-time I/O subsystem: TAO's real-time I/O (RIO) subsystem (Kuhns et al., 1999) extends support for CORBA into the OS. RIO assigns priorities to real-time I/O threads so that the schedulability of application components and ORB endsystem resources can be enforced. When integrated with advanced hardware, such as the high-speed network interfaces described below, RIO can (1) perform early demultiplexing of I/O events onto prioritized kernel threads to avoid thread-based priority inversion and (2) maintain distinct priority streams to avoid packet-based priority inversion. TAO also runs efficiently and relatively predictably on conventional I/O subsystems that lack advanced QoS features.

High-speed network interface: At the core of TAO's I/O subsystem is a "daisy-chained" network interface consisting of one or more ATM Port Interconnect Controller (APIC) chips (Dittia, Parulkar, and Cox, 1997). The APIC is designed to sustain an aggregate bi-directional data rate of 2.4 Gbps using zero-copy buffering optimization to avoid data copying across endsystem layers. In addition, TAO runs on conventional real-time interconnects, such as VME backplanes and multiprocessor shared memory environments.

TAO is developed using lower-level middleware called ACE (Schmidt and Suda, 1994), which implements core concurrency and distribution patterns (Schmidt et al., 2000) for communication software. ACE provides reusable C++ wrapper facades and framework components that support the QoS requirements of high-performance, real-time applications and higher-level middleware like TAO ACE and TAO run on a wide range of OS platforms, including Win32, most versions of UNIX, and real-time operating systems like Sun/Chorus ClassiX, LynxOS, and VxWorks.

3.2. *Design Goals of TAO's Scheduling Service*

To alleviate the limitations with existing scheduling strategies described in Section 2, our real-time scheduling research focuses on developing a CORBA-based framework that enables applications to (1) maximize total utilization, (2) preserve scheduling guarantees for critical operations, and (3) adapt flexibly to different application and platform characteristics. These goals are illustrated in Figure 5 and summarized below:

Goal 1—Higher utilization: The upper pair of timelines in Figure 5 demonstrates our first research goal: *higher utilization*. This timeline shows a case where a critical operation execution did not, in fact, use its worst-case execution time. With dynamic scheduling, an additional noncritical operation could be dispatched, thereby achieving higher resource utilization.

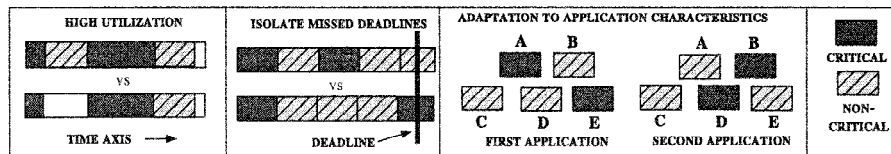


Figure 5. Design goals of TAO's dynamic scheduling service.

Goal 2—Preserving scheduling guarantees: The lower pair of timelines in Figure 5 demonstrates our second research goal: *preserving scheduling guarantees for critical operations*. In the lower timeline, priority is based only on traditional scheduling parameters, such as rate and laxity. In the upper timeline, criticality is also included. Both timelines depict schedule overrun. When criticality is considered, only non-critical operations miss their deadlines.

Goal 3—Adaptive scheduling: The sets of operation blocks at the bottom of Figure 5 demonstrate our third research goal: *providing applications with the flexibility to adapt to varying application requirements and platform features*. In this example, the first and second applications use the same five operations. However, the first application considers operations A and E critical, whereas the second application considers operations B and D critical. By allowing applications to select which operations are critical, it is possible to provide scheduling behavior that is appropriate to each application's individual requirements.

These three goals motivate the design of TAO's strategized scheduling service framework, described in Section 3.3. For the real-time systems (Harrison, Levine, and Schmidt, 1997; Schmidt, Levine, and Mungee, 1998; Kuhns et al., 1999; Schmidt et al., 1999; Lachenmaier, 1998) to which TAO has been applied, it has been possible to identify a core set of operations whose execution before deadlines is *critical* to the integrity of the system. Therefore, the TAO's scheduling service is designed to ensure that critical CORBA operations will meet their deadlines, even when the total utilization exceeds the schedulable bound.

If it is possible to ensure missed deadlines will be isolated to non-critical operations, then adding non-critical operations to the schedule to increase total CPU utilization will not increase the risk of missing critical deadlines. The risk will only increase for those operations whose execution prior to deadline is *not* critical to the integrity of the system. In this way, the risk to the whole system is minimized when it is loaded for higher utilization.

3.3. TAO's Strategized Scheduling Service Framework

TAO's scheduling service framework is designed to support a variety of scheduling strategies, including RMS, EDF, MLF, and MUF. In addition, this framework provides a common environment to systematically compare both existing and new scheduling strategies. This flexibility is achieved in TAO via the *Strategy* pattern (Gamma et al., 1995), which allows

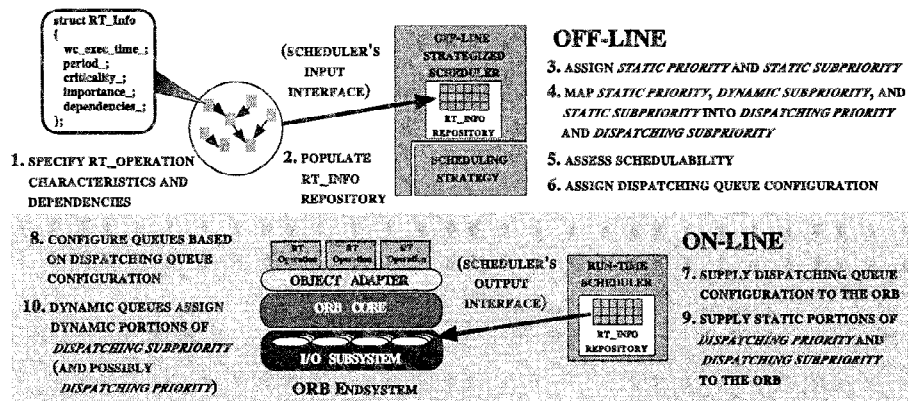


Figure 6. Processing steps in TAO's dynamic scheduling service architecture.

parts of the sequence of steps in an algorithm to be replaced, thus providing interchangeable variations within a consistent algorithmic framework. TAO's scheduling service framework uses the Strategy pattern to encapsulate a family of scheduling algorithms within a fixed CORBA IDL interface, thereby enabling different strategies to be configured *independently* from applications that use them.

The architecture and behavior of TAO's strategized scheduling service is illustrated in Figure 6. This architecture evolved from our earlier work on a CORBA scheduling service (Schmidt, Levine, and Mungee, 1998) that supported purely static RMS for avionics mission computing applications (Harrison, Levine, and Schmidt, 1997; Levine, Gill, and Schmidt, 1998; Lachenmaier, 1998). Based on this work, as well as our experience prototyping dynamic scheduling strategies, we have identified the following set of common steps that are necessary to configure and process requests for a broad range of scheduling strategies:

Step 1: A CORBA application specifies QoS information and passes it to TAO's scheduling service, which is implemented as a CORBA object, i.e., it implements an IDL interface.

Step 2: At configuration time, which can occur either off-line or on-line, the application passes this QoS information into TAO's scheduling service via its *input interface*, which is described in Section 3.4.1. TAO's scheduling service stores the QoS information in its repository of `RT_Info` descriptors. TAO's scheduling service constructs operation dependency graphs based on `RT_Info`s registered with it by the application. The scheduling service then identifies threads of execution by examining the terminal nodes of these dependency graphs. TAO's scheduling service can then infer information induced by the dependency graph, such as the effective periods of execution of dependent operations.

Step 3: Next, TAO's scheduling service assigns static priorities and subpriorities to operations. These values are assigned according to the specific strategy used to configure the

scheduling service. For example, when the TAO scheduling service is configured with the MUF strategy, static priority is assigned according to operation criticality.

Step 4: Based on the specific strategy used to configure it, TAO's scheduling service divides the dispatching priority and dispatching subpriority components into statically and dynamically assigned portions. The static priority and static subpriority values are used to assign the static portions of the dispatching priority and dispatching subpriority of the operations. These dispatching priorities and subpriorities reside in TAO's `RT_Info` repository. Performing this step at configuration time helps minimize run-time overhead and non-determinism.

Step 5: In this step, TAO's scheduling service assesses schedulability. A set of operations is considered *schedulable* if all critical operations will meet their deadlines. Schedulability is assessed according to whether all operations within and above the minimum critical static priority level will be able to meet their deadlines, based on the worst case simultaneous arrival of all operations, i.e., the *critical instant* (Liu and Layland, 1973). Operations are augmented with a dynamic subpriority based on the critical instant, and their resulting dispatching priority and dispatching subpriority are used to assess worst case feasibility of the critical operations. This *static* analysis can provide a worst-case schedulability assessment for static, dynamic, and hybrid strategies alike.

Step 6: Based on the assigned dispatching priorities, and in accordance with the specific strategy used to configure the off-line scheduling service, the number and types of dispatching queues needed to dispatch the generated schedule are assigned. For example, when the scheduling service is configured with the MLF strategy, there is a single queue, which uses laxity-based prioritization. As in Step 3, this static information is cached in the `RT_Info` repository until it is needed at run-time.

Step 7: When TAO's ORB endsystems and applications are initialized at run-time, the configuration information in the `RT_Info` repository is used by the scheduling service's run-time scheduler component, which is collocated within an ORB endsystem. The ORB uses this run-time scheduler to retrieve (1) the thread priority at which each queue dispatches operations and (2) the type of dispatching prioritization used by each queue. The scheduling service's run-time component provides this information to the ORB via its *output interface*, as described in Section 3.4.2. By encapsulating the thread priority and dispatching type information behind its output interface, TAO's strategized scheduling service decouples the *policies* for dispatching behavior from the *mechanisms* used to enforce those policies.

Step 8: In this step, the ORB configures its *dispatching modules*, i.e., in its I/O subsystem, ORB Core, and/or Event Service, as described in Section 3.4.6. Dispatching modules use the information from the scheduling service's output interface to create the correct number and types of queues and associate them with the correct thread priorities that service the queues.

Step 9: When an operation request arrives from a client at run-time, the appropriate dispatching module must identify the dispatching queue to which the request belongs and initialize the request's dispatching subpriority. To accomplish this, the dispatching module queries TAO's scheduling service's output interface, as described in Section 3.4.2. The run-time scheduler component of TAO's scheduling service (1) retrieves the static portions of the dispatching priority and dispatching subpriority from the `RT_InfO` repository and (2) supplies them to the dispatching module. By caching static information that was computed at configuration time, TAO's strategized scheduling service minimizes run-time overhead and non-determinism for each operation invocation.

Step 10: If the dispatching queue where the operation request is placed was configured as a *dynamic queue* in step 8, the dynamic portions of the request's dispatching subpriority (and possibly its dispatching priority) are assigned. The queue first does this when it enqueues the request and then updates these dynamic portions only as necessary when other operations are enqueued or dequeued. By efficiently managing updates to dynamic information, TAO's dynamic queues minimize the amount of overhead they introduce.

Steps 3–6 represent the strategized portion of the scheduling framework, which varies with each distinct scheduling strategy. Steps 1–2 and 7–10 represent the fixed portion of the framework, which remains the same for all scheduling strategies. Steps 1–6 typically occur off-line during the schedule configuration process, while steps 7–10 typically occur on-line. Support for these static and dynamic steps underscore the hybrid nature of TAO's scheduling architecture.

The remainder of this section describes TAO's strategized scheduling service framework in detail. Section 3.4 motivates why TAO allows applications to vary their scheduling strategy and shows how TAO's framework design achieves this flexibility.

3.4. *Enhancing Flexibility in TAO's Strategized Scheduling Framework*

The QoS requirements of applications and the hardware/software features of platforms and networks in which they are hosted often vary significantly. For instance, a scheduling strategy that is ideal for telecommunication (Schmidt, 1996) or medical imaging (Pyarali, Harrison, and Schmidt, 1996) applications may be poorly suited for avionics mission computing (Harrison, Levine, and Schmidt, 1997). Therefore, TAO's scheduling service framework allows applications to vary their scheduling strategies selectively. The flexibility of TAO's strategized scheduling service architecture is motivated by the following two design goals: (1) shield application developers from unnecessary implementation details of alternative scheduling strategies—this improves the system's reliability and maintainability; (2) decouple the strategy for priority assignment from the dispatching model so the two can be varied independently—this increases the system's flexibility to adapt to varying application requirements and platform features.

TAO supports this flexibility by decoupling the *fixed* portion of its scheduling framework from the *variable* portion, as follows:

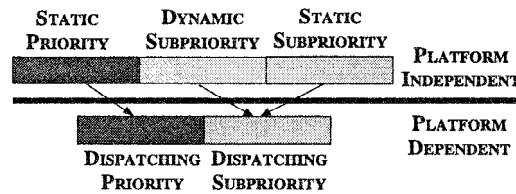


Figure 7. Distinct mapping levels within each strategy.

Fixed CORBA IDL interfaces: The fixed portion of TAO's strategized scheduling service framework is defined by the following two CORBA IDL interfaces:

- **Input Interface:** As discussed in Section 3.4.1, the input interface consists of the three operations shown in Figure 8(A). Applications use these operations to manipulate QoS characteristics expressed with TAO's `RT_InfO` descriptors (Schmidt, Levine, and Mungee, 1998) (Steps 1 and 2 of Figure 6).

- **Output Interface:** As discussed in Section 3.4.2, the output interface consists of the two operations shown in Figure 8(B). One operation returns the dispatching module configuration priority components assigned to an operation (step 7 of Figure 6). The other returns the dispatching priority and dispatching subpriority components assigned to an operation (step 9 of Figure 6). Section 3.4.6 describes how TAO's dispatching modules use information from TAO's scheduling service's output interface to configure and manage dispatching queues, as well as dispatch operations according to the generated schedule.

Variable mappings: By decoupling (1) the strategy for urgency to dispatching priority and dispatching subpriority, TAO allows the scheduling strategy and the underlying dispatching model to vary independently. In addition, this decoupling enables the systematic comparison of scheduling strategies over a range of dispatching models, from fully preemptive-by-urgency, through preemptive-by-priority-band, to entirely non-preemptive, which are discussed in Section 3.4.5. As illustrated in figure 7, the variable portion of TAO's scheduling service framework is implemented by the following two distinct levels of mapping from input interface to output interface:

- **Input mapping:** The input mapping assigns *platform-independent* urgency components, i.e., static priority, dynamic subpriority, and static subpriority, based on (1) the operation characteristics specified to the input interface and (2) the selected scheduling strategy, e.g., MUF, MLF, EDF, or RMS. Section 3.4.3 describes how each of the strategies implemented in TAO maps from the input interface to urgency values.

- **Output mapping:** The output mapping assigns dispatching priority and dispatching subpriority according to based on the urgency components assigned in the first level, and the underlying dispatching model. Section 3.4.4 describes how the output mapping translates the assigned urgency values into the appropriate dispatching priority and dispatching subpriority values for the output interface. Section 3.4.5 describes alternatives to the output mapping used in TAO and discusses key design issues related to these alternatives.

The remainder of this section describes how TAO's strategized scheduling service framework supports these fixed CORBA IDL interfaces and variable mappings. For each aspect of TAO's scheduling framework, we outline the key design challenges and explain how our solution addresses these challenges.

3.4.1. TAO's Scheduling Service Input Interface

Real-time applications must specify their QoS information to their selected scheduling strategy. The scheduling strategy then uses this application-specific QoS information to ensure that the QoS received by the application conforms to this information. The key design issues for QoS specification, and how TAO's strategized scheduling framework addresses them, are as follows:

Decoupling QoS specification and strategy details: Although the application must specify its QoS information to the instantiated scheduling strategy, it is essential that it not tightly couple the application to any specific scheduling strategy. This flexibly allows the scheduling strategy to be varied independently of the application-specific QoS information. Thus, changing the scheduling strategy need not require changes to the application.

TAO's scheduling strategy framework is designed to minimize unnecessary constraints on the values that application developers specify to the input interface described in Section 3.4.1. For instance, one (non-recommended) way to implement the RMS, EDF, and MLF strategies in TAO's scheduling service framework would be to implement them as variants of the MUF strategy. This can be done by manipulating the values of the operation characteristics (Stewart and Khosla, 1992). However, this approach would tightly couple applications to the MUF scheduling strategy and the strategy being emulated. There is a significant drawback to tightly coupling the behavior of a scheduling service to the characteristics of application operations. In particular, if the value of one operation characteristic used by an application changes, developers must remember to manually modify other operation characteristics specified to the scheduling service in order to preserve the same mapping. In general, TAO's scheduling service framework shields application developers from such unnecessary details.

create(): This operation takes a string with the operation name as an input parameter. It creates a new `RT_Info` descriptor for that operation name and returns a handle for that descriptor to the caller. If an `RT_Info` descriptor for that operation name already exists, `create` raises the `DUPLICATE_NAME` exception.

add_dependency(): This operation takes two `RT_Info` descriptor handles as input parameters. It places a dependency on the second handle's operation in the first handle's `RT_Info` descriptor. If either the handles refers to an invalid `RT_Info` descriptor, `add_dependency` raises the `UNKNOWN_TASK` exception.

```

interface Scheduler
{
    // ...

    // Create a new RT_Info descriptor for entry_point
    handle_t create ( in string entry_point )
        raises ( DUPLICATE_NAME );

    // Add dependency to handle's RT_Info descriptor
    void add_dependency ( in handle_t handle,
                        in handle_t dependency )
        raises ( UNKNOWN_TASK );

    // Set values of operation characteristics
    // in handle's RT_Info descriptor
    void set ( in handle_t handle,
              in Criticality criticality,
              in Time_worstcase_exec_time,
              in Period_period,
              in Importance_importance )
        raises ( UNKNOWN_TASK );

    // ...
}

interface Scheduler
{
    // ...

    // Get configuration information for the queue that will dispatch all
    // RT_Operations that are assigned dispatching priority d_priority
    void dispatch_configuration ( in Dispatching_Priority d_priority,
                                out OS_Priority os_priority,
                                out Dispatching_Type d_type )
        raises ( UNKNOWN_DISPATCH_PRIORITY,
              NOT_SCHEDULED );

    // Get static dispatching subpriority and dispatching
    // priority assigned to the handle's RT_Operation
    void priority ( in handle_t handle,
                  out Dispatching_Subpriority d_subpriority,
                  out Dispatching_Priority d_priority )
        raises ( UNKNOWN_TASK,
              NOT_SCHEDULED );

    // ...
}

```

Figure 8. TAO scheduling service IDL interfaces: (A) input (B) output.

set(): This operation takes an `RT_Info` descriptor handle and values for several operation characteristics as input parameters. The `set` operation assigns the passed input values to the corresponding operation characteristics in the `RT_Info` descriptor. If the passed handle refers to an invalid `RT_Info` descriptor, `set` raises the `UNKNOWN_TASK` exception.

3.4.2. TAO's Scheduling Service Output Interface

An ORB must obtain QoS enforcement information generated by the scheduling strategy. This information is then used by an ORB to enforce the QoS specified by the scheduling strategy. The key design issues for QoS enforcement, and how TAO's strategized scheduling framework addresses them, are as follows:

Decoupling QoS enforcement and strategy detail: While an ORB must obtain QoS enforcement information generated by the instantiated scheduling strategy, it must not tightly couple the ORB to any specific scheduling strategy. This allows the scheduling strategy to be varied independently from the ORB, so that changing the scheduling strategy does not require any changes to the ORB.

Defining a fixed output interface: TAO's scheduling service framework decouples QoS enforcement from any specific scheduling strategy by providing a *fixed* CORBA IDL *output interface*, behind which the scheduling strategy details are hidden. As illustrated in steps 7 through 10 of Figure 6, the ORB uses TAO's scheduling service output interface to obtain QoS enforcement information and configure its dispatching modules accordingly. The output interface for TAO's scheduling service consists of the corba IDL interface operations shown in Figure 8(B) and described below:

dispatch_configuration(): This operation provides configuration information for queues in the dispatching modules used by the ORB endsystem (step 7 of Figure 6). It takes a dispatching priority value as an input parameter. It returns the OS thread priority and dispatching type corresponding to that dispatching priority level. The run-time component of TAO's scheduling service retrieves these values from the `RT_Info` repository, where they were stored by TAO's off-line scheduling component (step 6 of Figure 6). The `dispatch_configuration` operation will raise the `UNKNOWN_DISPATCH_PRIORITY` exception if it is passed a dispatching priority parameter that is not in the schedule. Likewise, if a schedule has not been generated, the `dispatch_configuration` operation raises the `NOT_SCHEDULED` exception.

priority(): This operation provides dispatching subpriority information for an operation request (step 9 of Figure 6). It takes an `RT_Info` descriptor handle as an input parameter and returns the assigned dispatching subpriority and dispatching priority as output parameters. The run-time component of TAO's scheduling service retrieves the dispatching priority and dispatching subpriority values stored in the `RT_Info` repository by its off-line component (step 4 of Figure 6). If the passed handle does not refer to a valid `RT_Info` descriptor, `priority` raises the `UNKNOWN_TASK` exception. If a schedule has not been generated, `priority` raises the `NOT_SCHEDULED` exception.

3.4.3. Input Mappings Implemented in TAO's Scheduling Service

Each scheduling strategy must provide a platform-independent assignment of priority values to each operation. This allows an ORB to leverage commonality among scheduling strategies, while preserving appropriate variations among them. The key design issues for platform-independent priority assignment, and how TAO's strategized scheduling framework resolves them, are as follows:

Decoupling priority from OS-specific mechanisms: It is important that the priorities and subpriorities assigned by the scheduling strategies remain independent of specific priority enforcement capabilities of the underlying OS platform, at least at some level. This allows the same scheduling strategy to be implemented, with only minor modifications, on platforms with diverse priority enforcement capabilities.

Defining a platform-independent input mapping: TAO decouples priority from the specific priority enforcement capabilities of the underlying OS platform by providing a separate, platform-independent, level of priority assignment. TAO's strategized scheduling service leverages the commonality among these mappings to make its implementation more uniform. The variations between these mappings provide hooks for adaptation to the requirements of specific applications. Furthermore, TAO's strategized scheduling service framework simplifies development and experimentation with *new* scheduling strategies within TAO's standards-compliant real-time CORBA middleware framework.

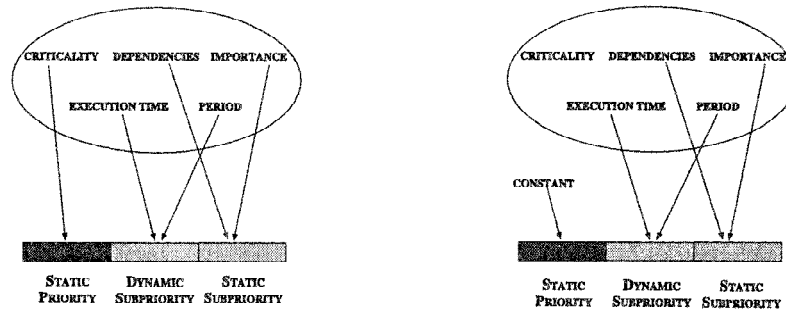


Figure 9. Input mappings: (A) MUF (B) MLF.

Input mappings for MUF, MLF, EDF, and RMS have been implemented in TAO's stratigized scheduling service, as described below. In each mapping, static subpriority is assigned first using importance and second using a topological ordering based on dependencies. The canonical definitions of MLF, EDF, and RMS do not include a minimal static ordering. Adding it to TAO's strategy implementations for these strategies has no adverse effect, however. This is because MLF, EDF, and RMS require that *all* operations are guaranteed to meet their deadlines for the schedule to be feasible, under *any* ordering of operations with otherwise identical priorities. Moreover, static ordering has the benefit of ensuring determinism for each possible assignment of urgency values.

Defining a platform-independent mapping for MUF: TAO provides a platform-independent mapping from operation characteristics onto urgency for MUF as shown in Figure 9(A). Static priority is assigned according to criticality in this mapping. Dynamic subpriority is assigned in the MUF input mapping according to *laxity*. Laxity is a function of the operation's period, execution time, arrival time, and the time of evaluation.

Defining a platform-independent mapping for MLF: TAO provides a platform-independent mapping from operation characteristics onto urgency for MLF as shown in Figure 9(B). The mapping for MLF assigns a constant (zero) value to the static priority of each operation. This results in a single static priority. The minimum critical priority is this lone static priority. The MLF strategy assigns the dynamic subpriority of each operation according to its laxity.

Defining a platform-independent mapping for EDF: TAO provides a platform-independent mapping from operation characteristics onto urgency EDF as shown in Figure 10(A). Like the MLF mapping, the EDF mapping also assigns a zero value to the static priority of each operation. Moreover, the EDF strategy assigns the dynamic subpriority of each operation according to its *time-to-deadline*, which is a function of its period, its arrival time, and the time of evaluation.

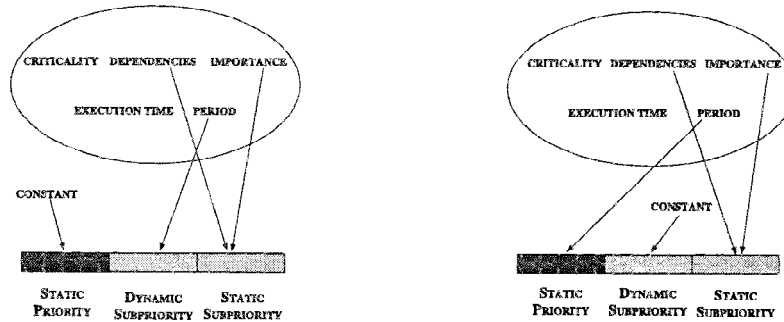


Figure 10. Input mappings: (A) EDF (B) RMS.

Defining a platform-independent mapping for RMS: TAO provides a platform-independent mapping from operation characteristics onto urgency for RMS as shown in Figure 10(B). The RMS mapping assigns the static priority of each operation according to its *period*, with higher static priority for each shorter period. The period for aperiodic execution must be assumed to be the worst case. In RMS, all operations are critical, so the minimum critical priority is the minimum static priority in the system. The RMS strategy assigns a constant (zero) value to the dynamic subpriority of each operation.

3.4.4. Output Mapping Implemented in TAO's Scheduling Service

Each scheduling strategy assigns platform-independent urgency values to operations, which must then be used to dispatch operations using each endsystem's OS-specific dispatching model. The key design issues for platform-specific dispatching, and how TAO's strategized scheduling framework resolves them, are as follows:

Enforcing priority through platform-specific dispatching: The input mappings described in Section 3.4.3 specify priorities and subpriorities for operations. However, there is no mechanism to enforce these priorities, independent of the platform-specific dispatching models. Therefore, each scheduling strategy must provide a mapping from platform-independent urgency values into platform-dependent dispatching priorities and subpriorities.

Defining platform-specific values for TAO's dispatching modules: In each of TAO's scheduling strategies, an output mapping transforms the platform-independent priority and subpriority values into dispatching priority and subpriority requirements that can be enforced by the specific dispatching models in real systems. Figure 11 illustrates the output mapping used by the scheduling strategies implemented in TAO. Each part of the mapping is described below.

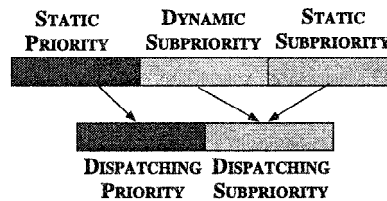


Figure 11. Output mapping implemented in TAO.

- **Dispatching priority:** In this mapping, static priority maps directly to dispatching priority. This mapping corresponds to the priority band dispatching model described in Section 3.4.2. Each unique static priority assigned by the input mapping results in a distinct thread priority in TAO's ORB request dispatching modules, which are described in Section 3.4.6. Thus, an operation with higher static priority will always preempt one with lower static priority.

- **Dispatching subpriority:** Dynamic subpriority and static subpriority map to dispatching subpriority. TAO's strategized scheduling service performs this mapping efficiently at run-time by transforming both dynamic and static subpriorities into a binary representation. TAO's preemption subpriority mapping scheme preserves the ordering of operation dispatches according to their assigned *urgency* values. Operations with the same static priority are ordered first by dynamic subpriority and second by static subpriority.

3.4.5. Alternative Output Mappings

The scheduling strategies implemented in TAO strike a balance between preemption granularity and run-time overhead. This design is appropriate for the hard real-time avionics applications we have developed. However, it is important to consider the consequences of the specific output mappings. Key design issues, and how TAO's strategized scheduling framework addresses them, are as follows:

Adapting to alternative OS dispatching models: TAO's strategized scheduling architecture is designed to adapt to the needs of a range of applications, not just hard real-time avionics systems. Different types of applications and platforms may require different resolutions of key design forces. Depending on (1) whether the OS supports thread preemption, (2) the number of distinct thread priorities supported, and (3) the preemption granularity desired by the application, several dispatching models can be supported by the output interface of TAO's scheduling service. Below, we examine three canonical variations supported by TAO, which are illustrated in Figure 12 and described below.

- **Preemptive-by-urgency:** The output mapping currently implemented in TAO only supports preemption *between* static priority levels. Thus, a newly arrived operation will not be dispatched until the operation executing currently at its same preemption priority

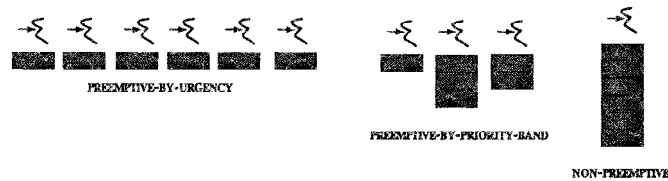


Figure 12. Dispatching models supported by TAO.

level has run to completion, even if the new operation has greater urgency. By assigning dispatching priority according to urgency, all scheduling strategies can be made fully preemptive, modulo OS dispatch latency overhead (Khanna et al., 1992). This model incurs greater complexity in the dispatching module implementation, including locking during enqueue and dequeue operations, which in turn increases run-time overhead.

- Preemptive-by-priority-band:** This model divides the range of all possible urgencies into fixed priority bands. It is similar to the non-preemptive dispatching model used by message queues in the UNIX System V STREAMS I/O subsystem (Rago, 1993; Kuhns et al., 1999). This dispatching model maintains a slightly weaker invariant than the preemptive-by-urgency model: at any given instant, an operation from the highest fixed-priority band that has operations able to execute is executing. The strategies implemented TAO's strategized scheduling service use a form of this model, as described in Section 3.4.4.

- Non-preemptive:** This model uses a single priority queue and is non-preemptive. Unlike the previous model, this models can be used on platforms that lack preemptive multi-threading.

3.4.6. Integrating TAO's Scheduling Service with its Dispatching Modules

Applications can benefit from strategized scheduling at a variety of points along an end-to-end request-response path. TAO's dispatching modules can be integrated at a number of different points within TAO's ORB endsystem architecture. This section (1) motivates the key design issues, (2) shows how the dispatching modules fit within TAO's ORB endsystem architecture, (3) describes the internal queueing mechanism of TAO's dispatching modules, and (4) discusses the issue of run-time control over dispatching priority within these dispatching modules.

Enforcing end-to-end QoS requirements: As noted in Section 2.1, one of our key research challenges is to implement dispatching modules that can enforce end-to-end QoS requirements. By designing dispatching modules that can enforce dispatching priority and dispatching subpriority of operations at arbitrary points in the TAO ORB endsystem architecture, we have increased TAO's ability to adapt to varying QoS enforcement capabilities of the endsystem OS platforms along an end-to-end request-response path.

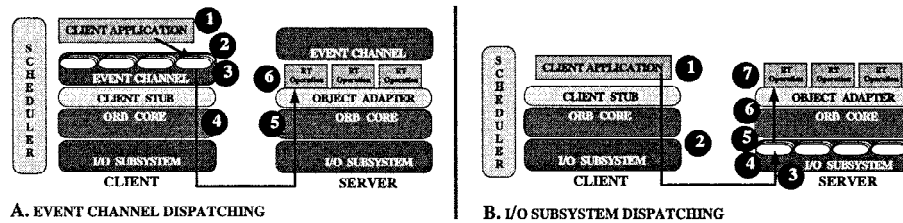


Figure 13. Alternative placement of dispatching modules.

Supporting alternative dispatching module configurations: The output interface of TAO's scheduling service is designed to work with dispatching modules in one or more layers in the TAO ORB endsystem architecture. For example, TAO's real-time extensions to the CORBA Event Service (Harrison, Levine, and Schmidt, 1997) use the scheduler output interface, as does its I/O subsystem (Kuhns et al., 1999). Figure 13(A) illustrates a configuration where a dispatching module resides in TAO's real-time Event Service (Harrison, Levine, and Schmidt, 1997). In Figure 13(A), the client application (1) pushes an event to TAO's Event Service according to dispatching priority and then dispatching subpriority. Each dispatched event results in (2), (3), and (4) a flow of control down through the ORB layers on the client and (5) back up through the ORB layers on the server, where (6) the operation is dispatched.

Figure 13(B) illustrates an alternative configuration where a dispatching module resides in TAO's I/O subsystem (Kuhns et al., 1999). The client application (1) makes a direct operation calls to the ORB, which (2) passes requests down through the ORB layers on the client and (3) back up the I/O subsystem layer on the server. The I/O subsystem's dispatching module (4) enqueues operation requests and (5) dispatches them according to their dispatching priority and dispatching subpriority, respectively. Each dispatched operation request results in (6) a flow of control up through the higher ORB layers on the server, where (7) the operation is dispatched.

Figure 13 illustrates two alternatives for configuring a dispatching module within a TAO ORB endsystem. However, TAO supports other configurations, as well. For example, a TAO ORB endsystem can be configured with dispatching modules in *both* the I/O subsystem and the Event Service. This configuration helps avoid priority inversions on server ORB endsystems via *early demultiplexing* (Kuhns et al., 1999) of events to prioritized threads in the I/O subsystem. Likewise, it helps avoid priority inversions on client ORB endsystems via prioritized dispatching of events in a collated client-side Event Service.

Internal architecture of TAO's dispatching modules: Figure 14 illustrates the general queueing mechanism used by the dispatching modules in TAO's ORB endsystem. In addition, this figure shows how the output information provided by TAO's scheduling service is used to configure and operate a dispatching module. During system initialization, each dispatching module obtains the thread priority and dispatching type for each of its queues from the scheduling service's output interface, as described in Section 3.4.2. Next,

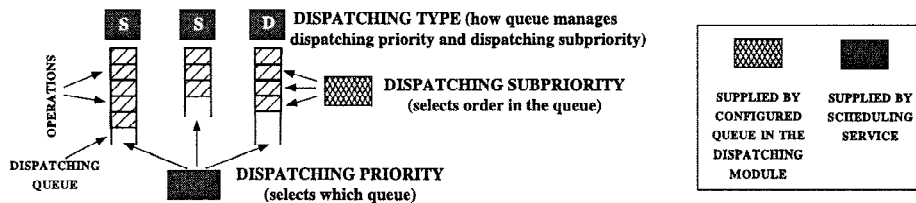


Figure 14. Example queuing mechanism in a TAO dispatching module.

each queue is assigned a unique dispatching priority number, a unique thread priority, and an enumerated dispatching type. Finally, each dispatching module has an ordered queue of pending dispatches per dispatching priority.

To preserve QoS guarantees, operations are inserted into the appropriate dispatching queue according to their assigned dispatching priority. Operations within a dispatching queue are ordered by their assigned dispatching subpriority. To minimize priority inversions, operations are dispatched from the queue with the highest thread priority, preempting any operation executing in a lower priority thread (Harrison, Levine, and Schmidt, 1997). To minimize preemption overhead, there is no preemption within a given priority queue.

The following three values are defined for the dispatching type: (1) `STATIC_DISPATCHING`—this type specifies a queue that only considers the static portion of an operation’s dispatching subpriority; (2) `DEADLINE_DISPATCHING`—this type specifies a queue that considers the dynamic and static portions of an operation’s dispatching subpriority, and updates the dynamic portion according to the time remaining until the operation’s deadline; (3) `LAXITY_DISPATCHING`—this type specifies a queue that considers the dynamic and static portions of an operation’s dispatching subpriority, and updates the dynamic portion according to the operation’s laxity. The deadline- and laxity-based queues update operation dispatching subpriorities whenever an operation is enqueued or dequeued.

4. The Performance of TAO’s Strategized Scheduling Service

Real-time scheduling behavior can only be achieved if the scheduling and dispatching mechanisms perform efficiently and predictably. Therefore, empirical benchmarks are needed to validate the framework. To ensure that TAO’s strategized scheduling service framework is efficient and predictable, we measured the dispatching overhead in TAO’s strategized scheduling service. We measured latency, which is the amount of time an operation is delayed, using time stamps. The run-time overheads for the static and dynamic scheduling strategies can be compared based on this measured latency.

We conducted two experiments. The first determined the run-time cost of dynamic dispatching for end-to-end performance. The second assessed the potential increase in dispatching overhead as varying loads were placed on the dispatching queues described in Section 3.4.6. These tests demonstrated that TAO’s dispatching modules can enforce dynamic end-to-end QoS requirements within acceptable levels of overhead.

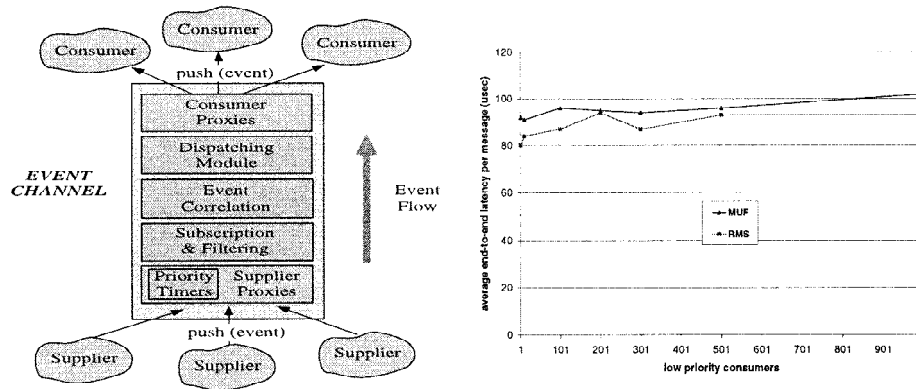


Figure 15. (A) Event service architecture (B) End-to-end run-time overhead.

The remainder of this section (1) describes an experiment to measure the minimum achievable end-to-end overhead for both static and dynamic scheduling strategies using TAO's Event Service over the TAO ORB, (2) describes an experiment to measure the overhead for static and dynamic dispatching queues as the load on these queues increases, and (3) draws conclusions about dynamic scheduling from the results of these experiments.

4.1. Dynamic Scheduling Overhead in TAO's Real-Time Event Service

The first experiment quantified the dynamic scheduling overhead in TAO's Event Service (Harrison et al., 1998), shown in Figure 15(A). This experiment consisted of a single high-priority supplier/consumer pair, and a varied number of low-priority event supplier/consumer pairs, ranging from 1 to 1,000 pairs. By varying the number of low-priority suppliers and consumers, this experiment measured (1) the effect of increasing low-priority load on high-priority performance, and (2) the minimum relative overhead associated with dynamic operation dispatching.

We measured the latency in event delivery between the high-priority supplier and consumer. This latency included (1) the time required for the TAO run-time scheduler to satisfy the Event Service dispatch module scheduling request plus (2) the time the request spent enqueued in the dispatch module. The test was run for two different scheduling strategies on a Sun Ultra 30 uni-processor 300 MHz Ultra SPARC CPU with 256 MB of memory, running SunOS5.5.1 and using the real-time (RT) scheduling class (Kuhns et al., 1999).

TAO's stratigized scheduling service was configured with an off-line RMS strategy and an $O(1)$ table lookup at run-time. The dynamic strategy used MUF and therefore required an additional run-time laxity calculation. The high-priority supplier and consumer were paced so that each high-priority operation was dequeued before the next was enqueued. This design remove any queueing effect from the high-priority queue, so its minimum relative overhead could be measured accurately.

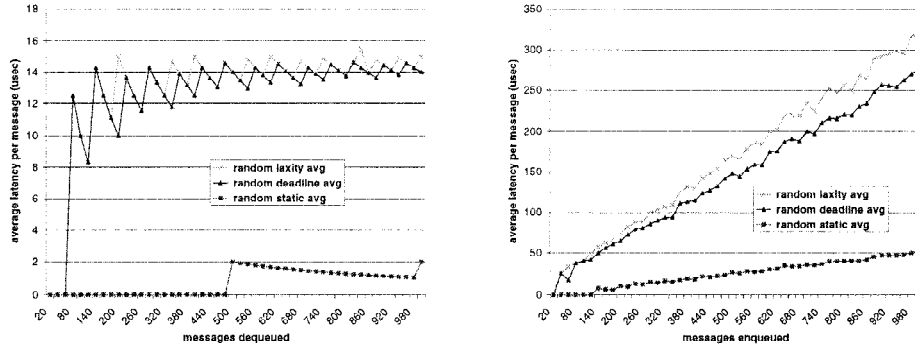


Figure 16. Static and dynamic dispatch overhead: (A) dequeue (B) enqueue.

The results of this experiment are shown in Figure 15(B). This figure illustrates that there was no significant change in high-priority performance with increasing low-priority load. Likewise, there appears to be only a small (< 10 percent) overhead end-to-end for dynamic dispatching with no queuing effect. In addition, the absolute overhead was between 80 and 100 μ secs.

4.2. Dynamic Scheduling Overhead in TAO's Dispatching Modules

The experiment described in Section 4.1 established the minimum relative end-to-end overhead for dynamic scheduling in TAO. Our second experiment gauged the potential impact of an increasing number of enqueued messages on this overhead. To measure this queuing effect accurately, we eliminated as many sources of constant overhead as possible. For instance, the queues were tested in isolation from TAO's Event Service and only the overhead of the enqueue and dequeue operations was measured.

The test was run on Windows NT 4.0 (SP3), in the real-time scheduling class on a dual-CPU Intel 333 MHz Micron Powerdigm with 256 MB of memory. The test used time stamps to measure the latency added by enqueue and dequeue operations for an increasing number of messages in the queue. A separate iteration of the test was run for each of an increasing number of enqueued messages. Messages were enqueued in random order. The same order was used for all queues in a given test iteration.

The test was run with three different kinds of dispatching queues. We tested static-, deadline-, and laxity-based queues. The static queue, which was used by the RMS scheduling strategy, used a $O(1)$ table lookup at run-time. The deadline-based queue, which was used by the EDF scheduling strategy, required an additional deadline calculation at run-time. The laxity-based queue, which was used by the MUF and MLF scheduling strategies, required an additional laxity calculation at run-time. The dequeue overhead for the laxity-based queue was highest, followed by the deadline-based queue, and then the static queue. As shown in Figure 16(A), there was an initial increase in overhead for dequeue operations in the laxity and deadline-based queues as the number of enqueued messages increases.

However, the overhead per-dequeue operation rapidly saturated at $\sim 14 \mu\text{secs}$ per operation for these queues. Thus, as the number of enqueued operations increased, the overhead for dequeue operations for the laxity- and deadline-based queues remained within a constant factor of seven times the overhead of the static queue.

The overhead for randomly ordered enqueue operations was highest for the laxity-based queue, followed by the overhead for deadline-based queue, and last for the static queue. As shown in Figure 16(B), the overhead per-enqueue operation increased linearly with the number of enqueued operations for all three kinds of queues. The overhead for enqueue operations for the laxity- and deadline-based queues remained within a constant factor of roughly six of the static queue overhead as the number of enqueued operations increased.

4.3. Analysis of Empirical Results

The tests described in Section 4.1 and Section 4.2 were run independently and in different experimental settings. Taken together, their results confirm empirically that dynamic scheduling strategies can be used effectively in real-time systems. Further, these results identify potential targets for optimization in cases where application requirements, such as heavy queue loading, degrade performance.

The remainder of this section (1) considers the implications of these results for systems with either moderate or heavy queueing, and (2) discusses alternative dispatching implementations and the conditions under which each may be preferable.

4.3.1. Moderately-loaded systems

Figure 15(B) shows that the minimal end-to-end latency for the laxity-based MUF scheduling strategy was only slightly higher than for the static RMS scheduling strategy. For systems where the maximum number of messages that can be enqueued at one time remains very small, the additional end-to-end overhead for dynamically scheduled dispatching should be relatively low.

If the number of messages that can be enqueued at one time increases, however, the effects of dynamic queue management become more prevalent, assuming a randomized enqueueing order. This dynamic queue management overhead is distributed between the enqueue and dequeue operations, so the measured overhead for both must be considered.

As shown in Figure 16(A), the overhead of dequeue operations does not appear significant for systems with fewer than 50 messages enqueued at one time. As the number of enqueued messages reached 100 messages, however, the overhead per-dequeue operation jumped to $\sim 12 \mu\text{secs}$ in the benchmark configuration described in Section 4.2. Even with a large number of enqueued messages, this overhead remained around $14 \mu\text{secs}$ per dequeue operation, roughly a factor of six times the overhead per-dequeue operation in the static queue. Thus, the overhead from dequeue operations in the laxity- and deadline-based queues remains reasonable, even as the number of enqueued operations increases significantly.

As shown in Figure 16(B), the overhead for laxity- and deadline-based enqueue operations does not appear to be significant if fewer than 20 messages are enqueued at one time. As the

number of enqueued messages reached 50 in the experiment described in Section 4.2, the overhead per-enqueue operation for the dynamic queues jumped to $\sim 20 \mu\text{secs}$. Although the laxity- and deadline-based enqueue overhead remained within a constant factor of six times the static enqueue overhead when more than 50 messages were enqueued, the significance of this constant factor increased with the number of enqueued messages.

4.3.2. *Heavily-Loaded Systems*

Depending on the characteristics of the specific application, the overhead for laxity- or deadline-based dispatching may reach unacceptable levels as the number of enqueued messages increases. Figure 16(B) shows that as the number of enqueued messages reached 1,000, the average overhead *per enqueue operation* exceeded $300 \mu\text{secs}$ for messages enqueued in randomized order. Thus, the total CPU time needed to enqueue these 1,000 messages was above 0.3 seconds.

For systems with such a large queueing effect, the overhead from dequeue operations will be minimal compared to the overhead for enqueue operations in the dispatching queues. Section 4.3.4 discusses two alternative dispatching priority queue implementations and describes when each are optimal for different numbers of enqueued messages and different application characteristics.

4.3.3. *Alternative Dispatching Mechanisms*

The dispatching queues described in Section 3.4.6 are implemented as linked lists. This minimizes the dequeue overhead for the static-, deadline-, and laxity-based dispatching queues, even as the number of enqueued messages becomes large. For the statically dispatched queues, the dispatching overhead remains reasonable as well, even as the number of enqueued messages approaches 1,000. However, for the laxity- and deadline-based queues, the enqueue overhead grows significantly as the number of enqueued messages increases.

One alternative to a linked list message queue implementation is to use a *heap*. A heap is a partially-ordered, almost-complete binary tree that ensures the average- and worst-case time complexity for enqueueing or dequeueing is $O(\lg n)$. The trade-off is that in the linked list priority queue implementation, enqueue operations are $O(n)$ and dequeue operations are $O(1)$. Conversely, in the heap-based priority queue implementation, both enqueue and dequeue operations are $O(\log n)$.

4.4. *Conclusions from Empirical Experiments*

The following conclusions can be drawn from the empirical results of our experiments with TAO's stratagized scheduling service:

Minimal end-to-end overhead: The minimal end-to-end overhead for dynamic scheduling strategies is comparable to that for static scheduling strategies, with only a small increase

due to dynamic priority computations. This indicates that dynamic end-to-end QoS requirements can be enforced within acceptable levels of overhead, assuming other sources of system overhead are minimized.

Range of acceptable performance: The range of acceptable performance is sustained for dynamic scheduling strategies, up to a load of ~ 150 messages enqueued at one time. TAO's strategized scheduling service and dispatching modules can adapt flexibly to alternative queuing implementations, so that for heavier loads, heap-based queues may be preferable.

Our empirical results provide strong evidence for the efficacy of our dynamic scheduling approach. The overhead of enforcing dynamic end-to-end QoS requirements remains within acceptable limits for systems with light to moderate queue loading. Further, the empirical results suggest alternative queuing implementations to give optimal performance under increasing loads. Thus, dynamic scheduling using TAO's strategized scheduling service framework can be achieved both efficiently and predictably.

5. Related Work

Real-time middleware is an emerging field of study. An increasing number of research efforts are focusing on integrating QoS and real-time scheduling into middleware like CORBA. This section compares our work on TAO with related QoS middleware integration research.

CORBA-related QoS research:

- **Mitre Real-time CORBA:** Krupp et al., (1994) at MITRE Corporation were among the first to elucidate the requirements of real-time CORBA systems. A system consisting of a commercial off-the-shelf real-time OS, a CORBA-compliant ORB, and a real-time object-oriented database management system is under development ("Statement of Work," 1997). Similar to TAO's original static scheduling service (Schmidt, Levine, and Mungee, 1999), their initial static scheduling approach used RMS, though a strategy for dynamic deadline monotonic scheduling support has been designed (Cooper et al., 1997).

- **URI TDMI:** Wolfe et al. are developing a real-time CORBA system at the US Navy Research and Development Laboratories (NRaD) and the University of Rhode Island (URI) (1997). The system supports expression and enforcement of dynamic end-to-end timing constraints through timed distributed operation invocations (TDMIs) (Fay-Wolfe et al., 1995). A TDMI corresponds to TAO's `RT.Operation` (Schmidt, Levine, and Mungee, 1998). Likewise, an `TR.Environment` structure contains QoS parameters similar to those in TAO's `RT.Info`. One difference between TAO and the URI approaches is that TDMIs express required timing constraints, e.g., deadlines relative to the current time, whereas `RT.Operations` publish their resources, e.g., CPU time, requirements.

- **BBN QuO:** The *Quality Objects* (QuO) distributed object middleware is developed at BBN technologies (Zinky, Bakken, and Schantz, 1997). QuO is based on CORBA

and provides the following support for agile applications running in wide-area networks: (1) provides *run-time performance tuning and configuration* through the specification of operating regions, behavior alternatives, and reconfiguration strategies that allows the QuO run-time to adaptively trigger reconfiguration as system conditions change (represented by transitions between operating regions), (2) gives *feedback* across software and distribution boundaries based on a control loop in which client applications and server objects request levels of service and are notified of changes in service, and (3) supports *code mobility* that enables QuO to migrate object functionality into local address spaces in order to tune performance and to further support highly optimized adaptive reconfiguration.

The QuO model employs several *QoS definition languages* (QDLs) that describe the QoS characteristics of various objects, such as expected usage patterns, structural details of objects, and resource availability. QuO's QDLs are based on the separation of concerns advocated by Aspect-Oriented Programming (AoP) (Kiczales, 1997). The QuO middleware adds significant value to adaptive real-time ORBs such as TAO. We are currently collaborating with the BBN QuO team to integrate the TAO and QuO middleware as part of the DARPA Quorum integration project.

- **UCSB Realize:** The Realize project at UCSB (Kalogeraki, Melliar-Smith, and Moser, 1997) supports soft real-time resource management of CORBA distributed systems. Realize aims to reduce the difficulty of developing real-time systems and to permit distributed real-time programs to be programmed, tested, and debugged as easily as single sequential programs. Realize integrates distributed real-time scheduling with fault-tolerance, fault-tolerance with totally-ordered multicasting, and totally-ordered multicasting with distributed real-time scheduling, within the context of OO programming and existing standard operating systems. The Realize resource management model can be hosted on top of TAO (Kalogeraki, Melliar-Smith, and Moser, 1997).

- **UIUC Epiq:** The Epiq project (Feng, Syyid, and Liu, 1997) defines an open real-time CORBA scheme that provides QoS guarantees and run-time scheduling flexibility. Epiq explicitly extends TAO's off-line scheduling model to provide on-line scheduling. In addition, Epiq allows clients to be added and removed dynamically via an admission test at run-time. The Epiq project is work-in-progress and empirical results are not yet available.

- **UCITMO:** The Time-triggered Message-triggered Objects (TMO) project (Kim, 1997) at the University of California, Irvine, supports the integrated design of distributed OO systems and real-time simulators of their operating environments. The TMO model provides structured timing semantics for distributed real-time object-oriented applications by extending conventional invocation semantics for object methods (i.e., CORBA operations) to include (1) invocation of time-triggered operations based on system times and (2) invocation and time bounded execution of conventional message-triggered operations.

TAO differs from TMO in that it provides a complete CORBA ORB, as well as CORBA ORB services and real-time extensions. Timer-based invocation capabilities are provided through TAO's Real-Time Event Service (Harrison et al., 1998). Where the TMO model creates new ORB services to provide its time-based invocation capabilities (Kim and Shokri, 1999), TAO provides a subset of these capabilities by extending the standard CORBA COS Event Service. We believe TMO and TAO are complementary technologies since (1) TMO

extends and generalizes TAO's existing time-based invocation capabilities and (2) TAO provides a configurable and dependable connection infrastructure needed by the TMO CNCM service. We are currently collaborating with the UCI TMO team to integrate the TAO and TMO middleware as part of the DARPA Quorum integration project.

Non-CORBA-related QoS research:

- **ARMADA:** The ARMADA project (Mehra, Indiresan, and Shin, 1997; Abdelzaher et al., 1997) defines a set of communication and middleware services that support fault-tolerant and end-to-end guarantees for real-time distributed applications. ARMADA provides real-time communication services based on the X-kernel and the Open Group's MK micro-kernel. This infrastructure provides a foundation for constructing higher-level real-time middleware services.

TAO differs from ARMADA in that most of the real-time infrastructure features in TAO are integrated into its ORB Core (Schmidt et al., 1999) and I/O subsystem (Kuhns et al., 1999), rather than in a micro-kernel. In addition, TAO implements the OMG CORBA standard, while also providing the hooks necessary to integrate with an underlying real-time I/O subsystem and OS. Thus, the real-time services provided by ARMADA's communication system can be utilized by TAO to support standards-based applications running over a vertically and horizontally integrated real-time systems.

- **CMU Publisher/Subscriber:** Rajkumar et al., (1995) at CMU developed a real-time Publisher/Subscriber model that is similar to the TAO's Real-time Even Service (Harrison, Levine, and Schmidt, 1997) e.g., it uses real-time threads to prevent priority inversion within its communication framework. The CMU model does not utilize any QoS specifications from publishers (event suppliers) or subscribers (event consumers), however. Therefore, scheduling is based on the assignment of request priorities, which is not addressed by the CMU model.

In contrast, TAO's Scheduling Service and real-time Event Service utilize QoS parameters from suppliers and consumers to assure resource access via priorities. One interesting aspect of the CMU Publisher/Subscriber model is the separation of priorities for subscription and data transfer. By handling these activities with different threads, with possibly different priorities, the impact of on-line scheduling on real-time processing can be minimized.

- **UCI RED-Linux Scheduling Framework:** Wang et al., (1999) at the University of California, Irvine, have proposed a general scheduling framework to unify three distinct kinds of scheduling approaches: *priority-based*, *time-based*, and *share-based*. Wang et al., decompose scheduling behavior into policy (*allocator*) and mechanism (*dispatching*) components, which are similar to the TAO scheduling service framework. They have implemented the dispatching portion of this framework in their real-time extensions to the Linux kernel, called RED-Linux.

While the RED-Linux approach to scheduling relies on special-purpose extensions to the OS kernel, TAO's scheduling service relies only on commonly available OS features, such as preemptive thread priorities. Therefore, TAO's dispatching mechanisms can leverage standards-based CORBA middleware and it can perform effectively on a wide range of commonly available real-time and general-purpose OS platforms.

- **OSU Share-based Scheduling:** Tyan et al., (1999) at Ohio State University, have developed a general framework for share-based scheduling. They demonstrate their framework's ability to implement a number of well-known fair queueing algorithms, as well as its ability to implement new kinds of share-based scheduling algorithms.

TAO's strategized scheduling service differs in that it uses priority based scheduling approaches, in order to address applications with hard real-time requirements. In our future research, we are investigating share-based scheduling and its interaction with priority-based scheduling for various classes of real-time applications.

6. Concluding Remarks

Many hard real-time systems, such as avionics mission computing and manufacturing process control systems, have traditionally been scheduled statically using variants of rate monotonic scheduling (RMS). Static scheduling provides assurance of schedulability prior to run-time and can be implemented with low run-time overhead. However, static scheduling handles nonperiodic processing inefficiently and treats invocation-to-invocation variations in resource requirements inflexibly. As a consequence, scheduled resources are underutilized and the resulting systems are hard to adapt to meet worst-case processing requirements.

Dynamic scheduling alleviates many limitations of static scheduling. However, purely dynamic scheduling strategies offer little or no control over which operations will miss their deadlines in an overloaded schedule. In addition, dynamic scheduling has a higher run-time cost because certain computations must be performed on-line, so it is necessary to measure this additional overhead and assess its significance.

To quantify the tradeoffs between static and dynamic scheduling algorithms, we have developed a *strategized scheduling service framework* and integrated this framework with TAO (Schmidt, Levine, and Mungee, 1998), which is our high-performance, real-time ORB. TAO offers applications the flexibility to specify and use different scheduling strategies, according to their specific needs. A C++ implementation of TAO's strategized scheduling service framework is available with the TAO ORB at the following URL: www.cs.wustl.edu/~schmidt/TAO.html.

This paper describes how we then used TAO's Event Service and run-time Scheduling Service to empirically measure end-to-end latency with and without queueing. Our results indicate that hybrid static/dynamic scheduling strategies can be used in real-time CORBA applications to (1) offer higher resource utilization than purely static scheduling strategies with acceptable run-time cost, (2) preserve the scheduling guarantees for critical operations even under an overloaded schedule, and (3) provide applications the flexibility to adapt to varying application requirements and platform features. Our empirical measurements provide a foundation upon which we will develop practical guidelines for configuring and using appropriate scheduling strategies for real-time CORBA applications. We are currently exploring the following areas in our future research on dynamic scheduling of real-time CORBA operations: (1) varying operation characteristics—additional empirical measurements are needed to assess the impact of varying the values of different operation characteristics on the performance of the scheduling strategies; (2) distributed scheduling

behavior—further empirical measurements are needed to determine the impact of factors such as network latency on the end-to-end performance of dynamically scheduled distributed systems; (3) available platform features—we are exploring the impact of various platform-specific features, such as preemptive multi-threading, on run-time scheduling behavior; (4) application requirements—a detailed examination of the impact of application specific requirements, such as policies for handling missed deadlines, will help guide the development of additional protocols for dynamically scheduled systems.

Acknowledgments

This work was funded in part by Boeing. We gratefully acknowledge the support and direction of the Boeing Principal Investigator, Bryan Doerr. In addition, we would like to thank Fred Kuhns for his guidance on the details of TAO's real-time I/O subsystem, and Priya Narasimhan for her extensive comments on this paper.

Notes

1. Priorities can be changed via *mode changes* (Schmidt, Levine, and Mungee, 1998), but that is too coarse to capture invocation-to-invocation variations in the resource requirements of complex applications.
2. A *dispatch* is a particular execution of an *operation*.

References

- Abdelzaher, T., Dawson, S., Feng, W.-C., Jahanian, F., Johnson, S., Mehra, A., Mitton, T., Shaikh, A., Shin, K., Wang, Z., and Zou, H. 1997. ARMADA middleware suite. *Proc. of the Workshop on Middleware for Real-Time Systems and Services, IEEE*. San Francisco, CA.
- Cooper, G., DiPippo, L.C., Esibov, L., Ginis, R., Johnston, R., Kortman, P., Krupp, P., Mauer, J., Squadrito, M., Thuraisingham, B., Wohlever, S., and Wolfe, V.F. 1997. Real-time CORBA development at MITRE, NRaD, Tri-Pacific and URI. *Proc. of the Workshop on Middleware for Real-Time Systems and Services, IEEE*. San Francisco, CA.
- Dittia, Z.D., Parulkar, G.M., and Cox, J.S.R. 1997. The APIC approach to high performance network interface design: Protected DMA and other techniques. *Proc. of INFOCOM '97, IEEE*. Kobe, Japan.
- Fay-Wolfe, V., Black, J.K., Thuraisingham, B., and Krupp, P. 1995. Real-time method invocations in distributed environments. University of Rhode Island, Department of Computer Science and Statistics. Tech. Rep. 95-244.
- Feng, W., Syyid, U., and Liu, J.-S. 1997. Providing for an open, real-time CORBA. *Proc. of the Workshop on Middleware for Real-Time Systems and Services, IEEE*. San Francisco, CA.
- Gamma, E., Helm, R., Johnson, R., and Vlissides, J. 1995. *Design Patterns: Elements of Reusable Object-Oriented Software*. Reading, MA: Addison-Wesley.
- Gokhale, A., and Schmidt, D.C. 1998. Optimizing a CORBA IIOP Protocol Engine for Minimal Footprint Multimedia Systems. *Journal on Selected Areas in Communications*. Special issue on *Service Enabling Platforms for Networked Multimedia Systems*.
- Han, C.-C.J. 1998. A better polynomial-time schedulability test for real-time multiframe tasks. *IEEE Real-Time Systems Symposium, IEEE*. Madrid, Spain.
- Han, C.-C.J., and Tyan, H.Y. 1997. A better polynomial-time schedulability test for real-time fixed-priority scheduling algorithms. *IEEE Real-Time Systems Symposium, IEEE*. San Francisco, CA.
- Harrison, T.H., Levine, D.L., and Schmidt, D.C. 1997. The design and performance of a real-time CORBA event service. *Proc. of OOPSLA '97, ACM*. Atlanta, GA.

- Harrison, T.H., O’Ryan, C., Levine, D.L., and Schmidt, D.C. 1998. The design and performance of a real-time CORBA event service. Submitted to *The Journal on Selected Areas in Communications*. Special issue on *Service Enabling Platforms for Networked Multimedia Systems*.
- Kalogeraki, V., Melliar-Smith, P., and Moser, L. 1997. Soft real-time resource management in CORBA distributed systems. *Proc. of the Workshop on Middleware for Real-Time Systems and Services, IEEE*. San Francisco, CA.
- Khanna, S., et al. 1992. Realtime scheduling in SunOS 5.0. *Proc. of the USENIX Winter Conference*. USENIX Association, pp. 375–390.
- Kiczales, G. 1997. Aspect-oriented programming. *Proc. of the 11th European Conference on Object-Oriented Programming*.
- Kim, K.H.K. 1997. Object structures for real-time systems and simulators. *IEEE Computer*, pp. 62–70.
- Kim, K., and Shokri, E. 1999. Two CORBA services enabling TMO network programming. *Fourth International Workshop on Object-Oriented, Real-Time Dependable Systems, IEEE*.
- Klein, M.H., Ralya, T., Pollak, B., Obenza, R., and Harbour, M.G. 1993. *A Practitioner’s Handbook for Real-Time Analysis: Guide to Rate Monotonic Analysis for Real-Time Systems*. Norwell, MA: Kluwer Academic Publishers.
- Kuhns, F., Schmidt, D.C., Levine, D.L., and Bector, R. 1999. The design and performance of RIO—A real-time I/O subsystem for ORB endsystems. *Proc. of the 5th IEEE Real-Time Technology and Applications Symposium, IEEE*. Vancouver, British Columbia, Canada.
- Lachenmaier, R. 1998. Open systems architecture puts six bombs on target. www.cs.wustl.edu/~schmidt/TAO-boeing.html.
- Levine, D.L., Gill, C.D., and Schmidt, D.C. 1998. Dynamic scheduling strategies for avionics mission computing. *Proc. of the 17th IEEE/AIAA Digital Avionics Systems Conference (DASC)*.
- Liu, C., and Layland, J. 1973. Scheduling algorithms for multiprogramming in a hard-real-time environment. *JACM* 20: 46–61.
- Mehra, A., Indiresan, A., and Shin, K.G. 1997. Structuring communication software for quality-of-service guarantees. *IEEE Transactions on Software Engineering* 23: 616–634.
- Mok, A.K., and Chen, D. 1997. A multiframe model for real-time tasks. *IEEE Transactions of Software Engineering* 23: 635–645.
- Newport, J.R. 1994 *Avionics Systems Design*. Boca Raton, FL: CRC Press.
- Object Management Group. 1998. *The Common Object Request Broker: Architecture and Specification*. 2.2 ed.
- Object Management Group. 1998. *Realtime CORBA 1.0 Joint Submission*. OMG Document orbos/98-12-05 ed.
- Pyarali, I., Harrison, T.H., and Schmidt, D.C. 1996. Design and performance of an object-oriented framework for high-performance electronic medical imaging. *USENIX Computing Systems*, vol. 9.
- Pyarali, I., O’Ryan, C., Schmidt, D.C., Kachroo, V., Arulanthu, A., Wang, N., and Gokhale, A. 1999. Applying Optimization Patterns to Design Real-time ORBs. *Proc. of the 5th Conference on Object-Oriented Technologies and Systems, USENIX*. San Diego, CA.
- Rago, S. 1993. *UNIX System V Network Programming* Reading MA: Addison-Wesley.
- Rajkumar, R., Gagliardi, M., and Sha, L. 1995. The real-time publisher/subscriber inter-process communication model for distributed real-time systems: Design and implementation. *First IEEE Real-Time Technology and Applications Symposium*.
- Schmidt, D.C. 1996. A family of design patterns for application-level gateways. *The Theory and Practice of Object Systems* (Special issue on *Patterns and Pattern Languages*) 2(1).
- Schmidt, D.C. 1990. GPERF: A perfect hash function generator. *Proc. of the 2nd C++ Conference, USENIX*. San Francisco, CA, pp. 87–102.
- Schmidt, D.C., Levine, D.L., and Mungee, S. 1998. The design and performance of real-time object request brokers. *Computer Communications* 21: 294–324.
- Schmidt, D.C., Mungee, S., Flores-Gaitain, S., and Gokhale, A. 1999. Software architectures for reducing priority inversion and non-determinism in real-time object request brokers. *Journal of Real-time Systems*. To appear.
- Schmidt, D.C., Stal, M., Rohnert, H., Buschmann, F. 2000. *Pattern-oriented Software Architecture: Vol. 2, Patterns for Concurrent and Networked Objects*. U.K.: Wiley.
- Schmidt, D.C., and Suda, T. 1994. An object-oriented framework for dynamically configuring extensible distributed communication systems. *IEE/BCS Distributed Systems Engineering Journal* (Special issue on *Configurable Distributed Systems*) 2: 280–193.
- Stewart, D.B., and Khosla, P.K. 1992. Real-time scheduling of sensor-based control systems. *Real-Time Programming* (W. Halang and K. Ramamritham, eds.) Tarrytown, NY: Pergamon Press.

- Stewart, D.B., Schmitz, D.E., and Khosla, P.K. 1992. Implementing real-time robotic systems using CHIMERA II. *Proc. of 1990 IEEE International Conference on Robotics and Automation* Cincinnati, OH.
- Statement of work for the extend sentry program, CPFF Project, ECSP replacement phase II. 1997. Submitted to OMG in response to RFI ORBOS/96-09-02.
- Thuraisingham, B., Krupp, P., Schafer, A., and Wolfe, V. 1994. On real-time extensions to the common object request broker architecture. *Proc. of the Object Oriented Programming, Systems, Languages, and Applications (OOPSLA) Workshop on Experiences with CORBA*, ACM.
- Tyan, H.-Y., and Hou, J.C. 1999. A rate-based message scheduling paradigm. *Fourth International Workshop on Object-Oriented, Real-Time Dependable Systems*, IEEE.
- Vinoski, S. 1997. CORBA: Integrating diverse applications within distributed heterogeneous environments. *IEEE Communications Magazine* 14.
- Wang, Y.-C., and Lin, K.-J. 1999. Implementing a general real-time scheduling framework in the RED-linux real-time kernel. *IEEE Real-Time Systems Symposium*, IEEE.
- Wolfe, V.F., DiPippo, L.C., Ginis, R., Squadrito, M., Wohlever, S., Zyk, I., and Johnston, R. 1997. Real-time CORBA. *Proc. of the Third IEEE Real-Time Technology and Applications Symposium*. Montréal, Canada.
- Zinky, J.A., Bakken, D.E., and Schantz, R. 1997. Architectural support for quality of service for CORBA objects. *Theory and Practice of Object Systems* 3(1).