

ControlWare: A Middleware Architecture for Feedback Control of Software Performance *

Ronghua Zhang, Chenyang Lu, Tarek F. Abdelzaher, John A. Stankovic
Department of Computer Science
University of Virginia
Charlottesville, VA 22903
e-mail: rz5b, chenyang, zaher, stankovic@cs.virginia.edu

Abstract

Attainment of software performance assurances in open, largely unpredictable environments has recently become an important focus for real-time research. Unlike closed embedded systems, many contemporary distributed real-time applications operate in environments where offered load and available resources suffer considerable random fluctuations, thereby complicating the performance assurance problem. Feedback control theory has recently been identified as a promising analytic foundation for controlling performance of such unpredictable, poorly modeled software systems, the same way other engineering disciplines have used this theory for physical process control.

In this paper, we describe the design and implementation of ControlWare, a middleware QoS-control architecture based on control theory, motivated by the needs of performance-assured Internet services. It offers a new type of guarantees we call convergence guarantees that lie between hard and probabilistic guarantees. The efficacy of the architecture in achieving its QoS goals under realistic load conditions is demonstrated in the context of web server and proxy QoS management.

1. Introduction

To achieve predictable behavior in distributed, poorly modeled, uncertain environments of today's open performance-assured applications, traditional approaches, such as resource reservation [23] and *a priori* knowledge of worst case execution conditions [27], are no longer applicable. Several recent research efforts have suggested the use of control theory [16, 30, 21, 18]. This theory offers a new type of guarantee that lies between hard and average (e.g., probabilistic), which we call *convergence guarantees* [19]. A basic convergence guarantee states that, upon any pertur-

bation, the performance variable of choice will converge to its desired value within a specified bounded time and that its deviation from that value is always bounded. The success of control theory is, in large part, due to its robustness in the face of modeling errors and external disturbances, which reduces the need for accurate system models — a much welcome property when accurate models are difficult to construct. Recent results have shown that control theory can also be successfully applied to the control of software performance [19].

The authors have applied control theory successfully in several example case studies involving computing applications. These case studies include performance isolation in web servers [5], web server delay control [18], proxy cache relative hit ratio control [21], network-layer active queue management for delay and loss differentiation [10], and microprocessor thermal management [29]. In this paper, we leverage the underlying insights to develop a middleware layer for QoS control that provides control-theoretic performance guarantees under uncertainty. The middleware is specifically targeted for Internet services. It allows the user to express QoS specifications off-line, maps these specifications into appropriate feedback control loop sets, tunes loop controllers analytically to guarantee convergence to specifications, and connects loops to the right performance sensors and actuators in the application such that the desired QoS is achieved. One main novelty of the middleware lies in isolating the application programmer from control-theoretic concerns while utilizing this theory to achieve the desired QoS guarantees.

The rest of the paper is organized as follows. Section 2 introduces the control theoretical approach in more detail and explains how QoS specifications are mapped into control loops. Section 3 presents the middleware architecture. Section 4 describes the resource management component. An evaluation of our QoS control functionality using the implemented middleware prototype in a web application scenario is presented in Section 5. Section 6 presents related work. The paper concludes with Section 7.

*The work reported in this paper was supported in part by the National Science Foundation under grants CCR-0093144 and CCR-0098269, and DARPA grants number F33615-01-C-1905 and N00014-01-1-0576.

Report Documentation Page				Form Approved OMB No. 0704-0188	
Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.					
1. REPORT DATE 2002		2. REPORT TYPE		3. DATES COVERED 00-00-2002 to 00-00-2002	
4. TITLE AND SUBTITLE ControlWare: A Middleware Architecture for Feedback Control of Software Performance				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of Virginia, Department of Computer Science, 151 Engineer's Way, Charlottesville, VA, 22094-4740				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES 10	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

2. Middleware-Based QoS Control

The application of feedback control theory to open distributed QoS-sensitive applications, such as mail servers, web servers and proxy caches, encounters three main challenges. First, from a control-theory perspective, there should exist a general way to convert QoS specifications of a computing system such as an Internet server, into feedback loops with known set points and feedback control parameters. This is achieved via multiple software tools and libraries that we describe in this section. Second, from a systems perspective, a convenient interface needs to be found between the service software and the middleware control loops that manage its performance. In our system, this interface is implemented by an entity called a *SoftBus*, which is described in Section 3. Third, appropriate software performance sensors and actuators must be designed. We elaborate on this challenge in Section 4. Figure 1 shows an overall picture of the middleware components.

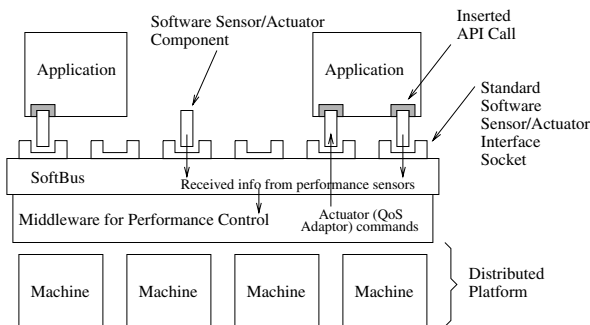


Figure 1. Overall Architecture

2.1. Service Development with ControlWare

To illustrate the main features of ControlWare, we first overview the development methodology of ControlWare-based performance-assured software. An application designer using our middleware for QoS provisioning would typically undergo the process shown in Figure 2:

- **QoS specification:** The required QoS guarantees should be specified for the system. ControlWare provides a simple Contract Description Language (CDL), which is used to describe the desired QoS guarantee. A partial syntax of CDL is presented in Appendix A.
- **QoS to control-loop mapping:** A tool called the *QoS mapper* interprets the CDL description offline and maps the required QoS guarantees to a set of feedback control loops and their set points. The QoS mapper specifies the feedback control loops using a *topology description language* and stores it in a configuration file.

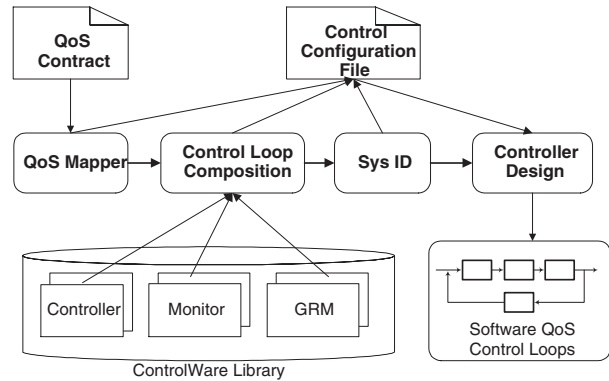


Figure 2. Development Methodology

- **Control loop composition:** The *loop composer* configures QoS monitors (also called sensors), actuators, and controllers in the manner described by the topology description language. These components can come from the library of ControlWare, and can also be supplied by users.
- **System identification:** To design a controller that can achieve the desired QoS guarantees, the mathematic model of the system must be known beforehand. ControlWare provides a system identification service that automatically derives difference equation models based on system performance traces [7]. Our prior publications [13, 18, 1] have already shown the feasibility and success of system identification of software systems.
- **Controller configuration and tuning:** Based on the model derived by system identification, ControlWare's controller design service can automatically tune the controllers to guarantee stability and desired transient response to load variations [19]. The resultant controller parameters are written into a configuration file.

The above process is somewhat similar to the way control engineers configure a distributed physical process control system. What is new here is that the controlled system is a software service, and the control goal is to provide convergence guarantees on QoS. With ControlWare, software engineers can easily add performance assurances to their systems without the need for a control-engineer's background. Our middleware automates that part of the feedback loop configuration process.

2.2. QoS Mapping

The cornerstone of a control theoretic paradigm for QoS guarantees in software systems lies in our ability to convert common resource management and software performance

assurance problems into feedback control problems. Our middleware contains a library of templates written in our topology description language, each formulating a particular type of QoS guarantees as a feedback control problem. The library is extendible in that a control engineer can transform a new guarantee type into a macro that describes the corresponding loop interconnection topology and store that macro in the middleware's library. Currently the library includes template for absolute convergence guarantees, performance isolation, statistical multiplexing, prioritization, relative differentiated service guarantees, and optimization guarantees. As an example, we describe the implementation of the basic (absolute) convergence guarantee, and its use in formulating relative guarantees, prioritization, and optimization as feedback control problems.

2.3. The Absolute Convergence Guarantee

Since it is impossible to achieve absolute guarantees in a system where load and resources are not known *a priori*, we define the absolute guarantee problem as one of *convergence* to a specified performance. The statement of the problem is to ensure that a performance metric, R , (i) converges within a specified exponentially decaying envelope to a fixed value, $R_{desired}$, and that (ii) the maximum deviation $R_{desired} - R$ be bounded at all times, as shown in Figure 3.

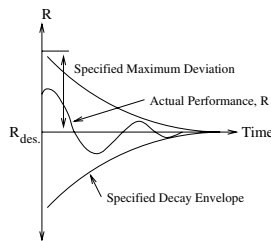


Figure 3. The Absolute Guarantee Specification

The problem requires that R be both *measurable* and *controllable*. *Measurability* requires that in the steady state, the measured value of R asymptotically approach its true value. *Controllability* refers to the requirement that the application must have some adaption mechanism, $A(R)$ that affects the value of R . For example, if R is CPU utilization, $A(R)$ can be an admission control mechanism.

The absolute convergence guarantee is translated into the control loop shown in Figure 4. The loop samples the measured performance, compares it to the desired value $R_{desired}$, and uses the difference to induce changes in resource allocation via the actuator $A(R)$.

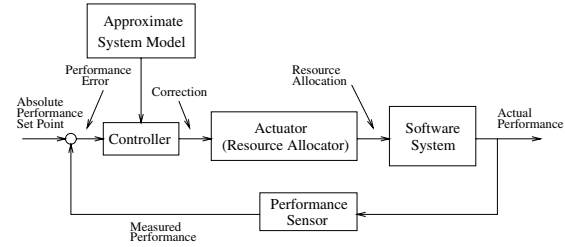


Figure 4. Absolute Guarantees

The absolute convergence guarantee loop is the elementary building block of our middleware from which all other assurances follow, as described below.

2.4. The Relative Guarantee Problem

The relative guarantee is to keep the ratio fixed between the performance (such as delay, throughput, etc) of two traffic classes. In general, let n be the number of classes in the system, and H_i be the measured performance of class i . The differentiation policy requires that the performance of different classes be related by the expression: $H_1 : H_2 : \dots : H_n = C_1 : C_2 : \dots : C_n$, where C_i is a proportionality constant or weight of class i . This kind of guarantee can be translated into n control loops, one for each class. In i_{th} control loop, the sensor measures the *relative performance*, $R_i = H_i / (H_1 + H_2 + \dots + H_n)$. This value is compared with the set point $\frac{C_i}{\sum C_j}$ to get the performance error e_i . The resource allocation of the class is then altered by $f(e_i)$. It is trivial to prove that $\sum_{1 \leq i \leq n} f(e_i) = 0$, for any linear function f . Hence, the feedback loops can operate independently with one loop per class, while the total amount of allocated resource remains constant if the controller is a linear function of error.

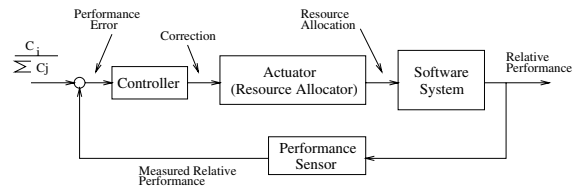


Figure 5. Relative Differentiated Service

2.5. The Prioritization Problem

The prioritization problem is defined as one where all service clients are partitioned into n classes, such that for every class i , it is desired that clients of that class do not suffer any contention over some shared resource r from any clients of classes $i + 1, \dots, n$.

We implement these semantics by a composition of instances of the elementary block described in Section 2.3. First, we make the entire server capacity available to the highest priority class using the basic convergence guarantee loop of Figure 4 with a set point equal to total server capacity. If that set point is not reached (because there is not enough demand), the unused capacity of each class is measured and treated as the set point for the resource allocation to the lower priority class. This requires a sensor array such that $S(R_i)$ measures the fraction of resource r consumed by clients of class i , as well as an actuator array where $A(R_i)$ controls the resource allocation of class i . The arrays are implemented by a set of per class performance counters and admission control limits.

The feedback loop architecture for prioritization is described for a two class server in Figure 6. One control loop is needed per class. Application performance converges to that of a strictly prioritized system. The approach may be used to implement logical priorities in middleware when the controlled server itself does not support priorities by design, such as the Apache [12] web server. An example and evaluation of this use is presented in [3].

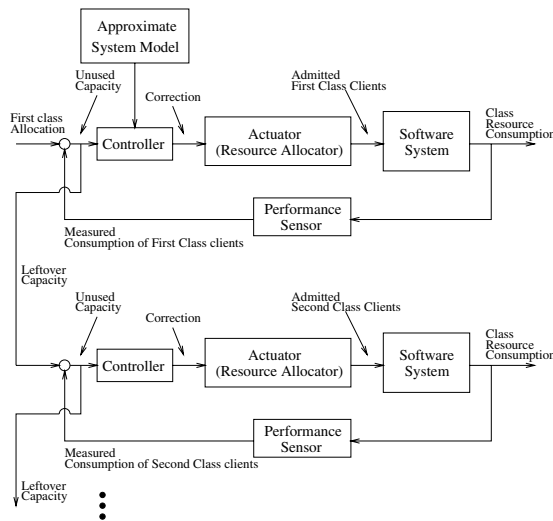


Figure 6. Prioritization

2.6. The Utility Optimization Problem

Another type of performance problems addressable using a control-theoretic framework is that of utility optimization. Following a microeconomic model [22], consider a computing service which produces an amount of work w . Let the benefit per unit of work be k . Hence, the total utility U produced by the service is $U = kw$. Let the resource consumption of the service be some nonlinear function, $g(w)$, which represents a measure of cost. It is desired to achieve

the maximum net profit, i.e., maximize $kw - g(w)$. Assuming a concave cost function, $g(w)$, the profit is maximized when the marginal utility is equal to the marginal cost, or when $\frac{dg(w)}{dw} = k$. The equation can be solved for w which then becomes the control set point, R . In a computing example, w , may be the desired server utilization, the desired workload size, or other metrics depending in the problem formulation. The approach is illustrated in Figure 7.

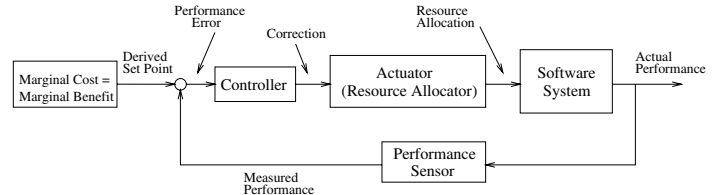


Figure 7. Utility Maximization

As shown above, ControlWare can express the most common guarantee types required in performance-assurance software by casting them appropriately as feedback control problems. Once the control loops are designed from the QoS specification, the middleware uses textbook techniques [28] to estimate system models and determine appropriate feedback controller parameters for guaranteed convergence of the control loops to the specified performance.

3. SoftBus: The Distributed Interface

To promote interoperability, the control engineering community standardized open layered interface architectures such as the Fieldbus [9] greatly simplifying the interconnection of sensors, actuators, and controllers in a digital control system. Similarly, ControlWare implements a SoftBus whose main purpose is to provide a common interface for efficient information exchange between software performance sensors, actuators and controllers across machines and address spaces. The sensors, actuators and controllers need not know each other's locations and need not worry about distributed communication. Underneath the common API, different information exchange mechanisms are developed for different situations.

3.1. Interface Modules

We support two types of software sensors and actuators: *passive* and *active*. A *passive* sensor or actuator is just a function call that returns sample data or accepts a command when called by the controller. An *active* sensor or actuator, in contrast, is a process or thread which may be running in its own address space. It is usually awakened periodically by the operating system scheduler to perform sensing

or actuation. For example, an idle CPU-time sensor may be implemented as an active sensor process which runs at the lowest priority and computes the percentage of time it has been executing to infer processor utilization.

Correspondingly, two interface modules are provided to facilitate the communication between sensors/actuators and SoftBus. Each component is attached to the SoftBus via an appropriate interface module. Internally, communication with local passive components is implemented as a direct function call, while communication with local active ones is through shared memory. Section 3.4 describes the transparent distributed communication between components.

3.2. Registrar

Configurability is achieved through the registration and deregistration of control loop components. Registration API is exported by an entity, called the registrar. Internally, the registrar maintains a cache. For each local component, it records in the cache the component's type (sensor/actuator or controller), a callback function pointer if it is passive, or a shared memory address if it is active. For remote components, it will record their location. Originally, only local components have entries in the registrar's cache. When some component's information is needed but can not be found in the cache, the registrar contacts an external directory server and caches the received information. When caching information on remote components, the registrar also creates a daemon to wait for invalidation messages from the directory server. When it receives a message notifying it of the deregistration of some components from the directory server, the registrar will purge the corresponding entries from the cache accordingly.

3.3. Directory Server

The directory server maintains the location and properties of all control loop components. To maintain cache consistency, the directory server keeps track of all machines that cache its information and notifies them when data has changed. When all the components are on one machine, the directory server is no longer needed. In this case, SoftBus optimizes itself automatically by shutting down the unnecessary daemons, and inhibiting communication between the registrars and the directory server. In the present implementation, the number and identities of the machines which run SoftBus is stored in a static configuration file. It is reasonably straightforward to extend this architecture to allow new machines to subscribe to the SoftBus dynamically using a group membership service such as [25, 6, 2].

3.4. Data Agent

The data agent abstracts away remote communication between sensors, actuators, and controllers. When an interface module of some component has data to send to another, the data agent first queries the registrar for information about the target component. If it is a remote one, the local agent forwards the request to the data agent on the destination machine. If the destination is local, data is passed to that component's interface module via shared memory.

Figure 8 depicts the above components and their interactions. The architecture allows easy and flexible configuration of control loops in which the controlled system is a software process.

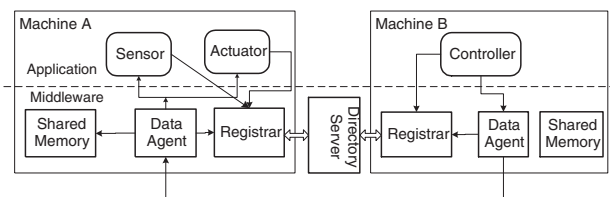


Figure 8. The Software Bus

4. Resource Management

An important challenge in applying a control-theoretic paradigm to software QoS control lies in finding appropriate sensors and actuators for software services. Sensors typically amount to a modest instrumentation of application code. For example, a sensor measuring the request rate on a particular site can be implemented as a simple counter that is reset periodically. A sensor measuring delay can be implemented as a moving average of the difference between two timestamps. Often the measured metric is already available as a variable maintained by the controlled software service (e.g., some queue length) or the operating system (such as CPU utilization). All one needs to do to implement a sensor and pass the value to the middleware.

The main application interface challenge, therefore, lies at the application/actuator boundary. To meet this challenge, our middleware includes a generic resource manager that serves as a multipurpose actuator. The manager exports a uniform API to the application and has a back-end interface to the machine's native resource allocation mechanisms. In this section, we present the design of our multipurpose actuator and the interface between it and the application. Note that, custom-made actuators that are not based on our generic resource manager can still be interfaced to SoftBus as described in Section 3, since the latter is oblivious to the type of actuator used.

Our generic resource manager (GRM) is designed for use with Internet servers such as web servers, DNS servers, mail

servers, and proxy cache servers. It understands the notion of *traffic classes*, and exports the abstraction of *resource quota* to represent the amount of logical resources allocated to a particular class. The action of the manager lies in controlling resource quota allocations.

The structure of the generic resource manager is shown in Figure 9. In the figure, the Classifier and Resource Allocator are provided by the application. The Resource Allocator does resource allocation. The Queue Manager maintains one queue for each class, governed by a certain queuing policy. The Quota Manager maintains a resource quota for each class.

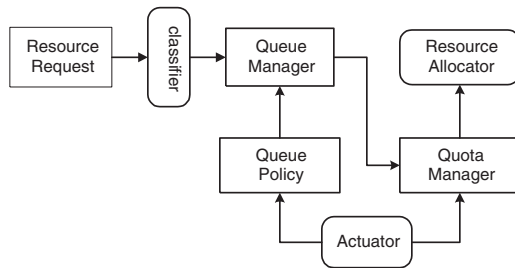


Figure 9. Structure of GRM

4.1. Customizing the GRM

To make this manager general and flexible, we try to expose as many tunable ‘knobs’ as possible so that the application can control the behavior of the manager as needed. These knobs are exposed to the outside world as *policies*. Some of the policies are:

1. *Space Policy*: This policy controls the total space used by the managed resource queues and the space allocation among the queues. The total space can be unlimited (limited only by available memory) or limited. The application can also set a limit on some(all) queues and let the remaining queues share the remaining space.
2. *Overflow Policy*: This policy takes effect only when some queues are sharing limited space and the space gets used up. Two options are supported: *reject* and *replace*. When the policy is reject, current request will be simply rejected. If the policy is replace, the last request of the lowest priority queue that shares the limited space will be evicted from the queue (application will be notified via a callback function) and current request will occupy its space.
3. *Enqueue Policy*: Apart from the queue for each class, the queue manager also maintains an ordered list of the

requests in all the queues. This policy influences the order of the requests in the list. System default policy is FIFO.

4. *Dequeue Policy*: This policy influences the dequeuing of the request. Currently three choices are available: FIFO, priority and proportional. FIFO means dequeue the request according to its position in the ordered list. Priority means always processing the high-priority queue before processing the low-priority queue. If proportional policy is chosen, the application can specify the dequeue ratio among the classes. For example, by setting the ratio to be 2 : 1, the queue for the class 0 will be dequeued twice as fast as the queue for class 1.

4.2. Interaction Between GRM and Application

Figure 10 summarizes the interaction between GRM and the application. When some resource is requested by the application, the request is first classified by the Classifier. After that, the request is passed to GRM by calling *insertRequest*. GRM controls resource allocation by checking the request against two constraints: (i) the queue length, and (ii) the quota constraint. If the queue for the given class is empty and the class has quota, the request is satisfied immediately via the function call *allocProc* to the resource allocator, and the quota is updated accordingly. If the request can’t be satisfied immediately, it will be buffered in the its queue. When some resource becomes available, the application calls *resourceAvailable* to notify GRM, which will try to satisfy as many pending requests as possible.

It’s important to mention that quota is a purely logical concept. Unlike the traditional resource reservation system, in our middleware the mapping of quota to physical resource consumption need not be known.

In effect, the GRM is a logical queuing, admission control, and resource allocation policy interface with a backend that is capable of executing a primitive service function such as assigning a request to a service process. The GRM generalizes the expression of various resource allocation policies in a common framework and makes it possible to control logical quota allocations in a trial-and-error fashion until performance constraints are met. The trial and error is guaranteed to converge because of the way controllers are designed which is the advantage of using a control-theoretic approach. Most importantly, the physical mapping of quota to actual resource consumption need not be known for correct operation, which separates this approach from resource reservation systems.

5. Evaluation

To test ControlWare, we instrument the Apache [12] server and Squid [26] server to interface to the middleware. We

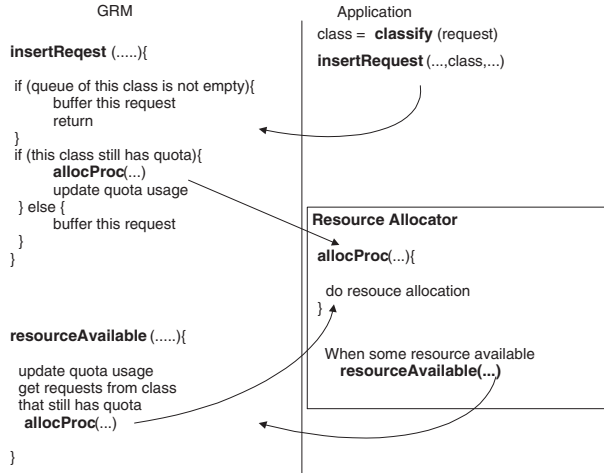


Figure 10. Resource allocation procedure

specify relative service differentiation as the QoS objective. On the proxy cache, we require differentiation in terms of hit ratio achieved to different content classes. On the server, we require differentiation in terms of service delay. While the semantics of the performance variable being differentiated are not interpreted by the middleware, our choice of variable is implicitly expressed in the choice of sensors. Hence, we instrument Squid to measure hit ratio and instrument Apache to measure service delay. These sensors are interfaced to SoftBus. The two servers chosen are quite different in terms of the resource types managed to achieve performance differentiation. In Squid, we manage cache size allocated to each content class. In Apache we manage the number of processes allocated to serve requests of each class. By applying our middleware to these servers, we demonstrate its versatility and ability to satisfy diverse performance guarantees. We describe our experiments in more detail next.

5.1. Providing Hit Ratio Differentiation in Squid

Figure 11 depicts the structure of the instrumented Squid. Cache space is shared by several classes and each class has a quota of the space. Generally, the space used by some class will directly affect its hit ratio. The sensor, actuator and controller provided by ControlWare constitute the control loops. ControlWare creates one for each class. Each sensor $S(i)$ returns the relative hit ratio of class i , i.e., $\frac{HR_i}{\sum_{k=0}^n HR_k}$. Each actuator changes the space allocated to its class by a value proportional to the error. The sensor and actuator are of the passive type. Periodically, ControlWare invokes the controller, which reads data from the sensor via SoftBus, calculates the resource change to be applied, and writes the result to the actuator via SoftBus.

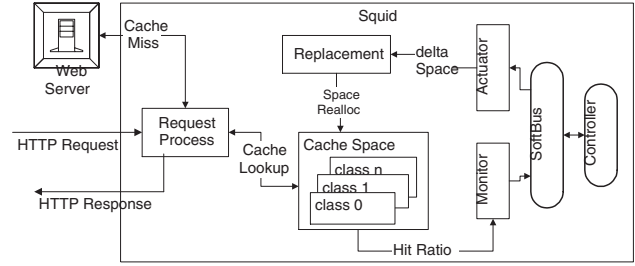


Figure 11. Structure of the modified Squid

The above proxy cache prototype was tested on a 100Mbps Ethernet LAN of nine Linux PCs. Each machine has a 450MHz AMD K6-2 processor and a 256MB RAM. Three client machines run Surge[8] to generate workload. Surge is a web workload generation tool known for its realistic reproduction of real web traffic patterns such as manifestation of a heavy-tailed request arrival and file-size distributions, a Zipf requested file popularity distribution, and proper temporal locality of accesses. Each client machine simulates 100 users. Three machines were used to run Apache. Each client machine generates requests for the content located at one of the Apache machines. In our experiment, there are 3 content classes. We specified that the hit ratio of the three classes satisfy the condition $H_0 : H_1 : H_2 = 3 : 2 : 1$. Squid is configured to use 8M bytes as its cache. Figure 12 shows the observed hit ratio differentiation during the experiment.

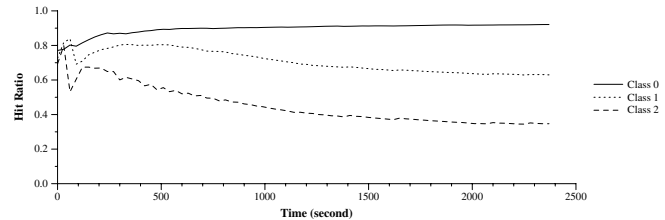


Figure 12. Hit Ratio of three classes

As we can see from Figure 12, the instrumented squid server successfully provides the specified hit ratio differentiation, illustrating the success of the middleware in providing performance guarantees.

5.2. Providing Delay Differentiation in Apache

We interfaced the Apache web server to SoftBus as depicted in Figure 13. We implemented a request classifier, and a delay sensor. The generic resource manager described in Section 4 was used as the actuator. The GRM was interfaced to a resource allocator which passed accepted requests (socket descriptors) to background Apache processes when instructed by the GRM.

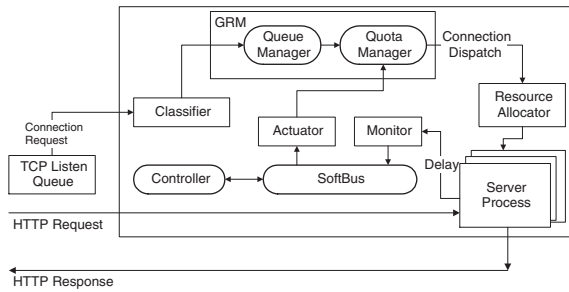


Figure 13. Structure of the modified Apache

Four client machines were used to run Surge and generate realistic server workload. As before, each client simulates 100 users. We divided the client machines into two classes with two machines per class. In the first half of the experiment, only one machine from class 0 generates requests. The second one is turned on after 870 seconds. We specified that the connection delay D_i of class 0 and class 1 should satisfy $D_0 : D_1 = 1 : 3$ at all times. Figure 14 shows the results of this experiment.

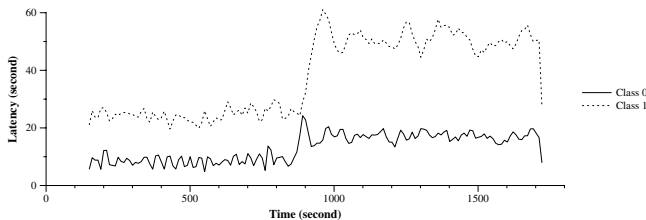


Figure 14. Relative Delay between two classes

From the figure, we can see that before 870 seconds, the delay of class 1 is about 3 times of the delay of class 0 as specified. When the second machine of class 0 is turned on, the delay of class 0 is increased suddenly. The controller reacts by allocating more processes to class 0. At about 1000 seconds, the delay ratio converge to around 3 again.

The evaluation clearly demonstrates that the middleware is capable of providing the specified performance guarantees with only a modest instrumentation cost and no control-theoretical experience required from the software developer. The middleware is versatile in that it is not tailored for a specific software service or a specific performance metric. Not only does the middleware allow new services to be efficiently augmented with QoS provisioning, but also it makes it easy to retrofit delivery of QoS assurances into services that were not designed with this purpose in mind. This paper provides a proof of concept of the utility of ControlWare as an embodiment of a general control-theoretic paradigm for QoS guarantees in software systems. Relative guarantees were used as an example in our evaluation. The authors

will report on a detailed evaluation of other types of guarantees and other Internet services in a subsequent publication.

5.3. Performance Evaluation

One concern of applying ControlWare is the overhead it introduces. We argue that the overhead is actually very small. First, when the application is not distributed, ControlWare can optimize itself by shutting down unnecessary daemons, and stop communication between the registrars and directory servers. Even in the distributed case, the directory server only needs to be contacted when the location of some component is unknown. After that, this information is cached locally. Since the control loop structure is very stable, the overhead of maintaining the cache consistency is almost neglectable. So the overhead is just the round trip time over the network for fetching data from remote components.

To quantify the overhead, we design a small program based on ControlWare, and test it on the testbed. The control loop spans two machines. Sensor and actuator are located at one machine, and controller resides at the other. The directory server runs on a third machine. All the components are reactive. Each invocation of the feedback control costs 4.8ms. Since in a typical application based on ControlWare, the invocation of control loop will not be so frequent, usually at the magnitude of second. Hence, the overhead of ControlWare will be relatively even smaller.

6. Related Work

The control theoretical approach has been successfully applied to in several computer system projects. At the network layer, Hollot et. al. [16] applied control theory to analyze the RED congestion control algorithm on IP routers. Recently feedback control is applied to active queue management to provide loss and delay differentiation [10].

In the area of CPU scheduling, Steere et. al. [30] developed a feedback based CPU scheduler that synchronizes the progress of consumers and supplier. In [20], feedback control real-time scheduling algorithms were developed to provide deadline miss ratio guarantees to real-time applications with unpredictable workloads.

Recently Internet server software has become a focus area of feedback control because the unpredictabilities of the workload. Examples of such QoS control includes delay and bandwidth control in Web servers [4, 18], hit ratio differentiation in web caches [21] and queue management in e-mail servers [24].

The above projects demonstrated that feedback control provides a powerful theoretical foundation to provide robust QoS guarantees in a wide range of software systems. However, their feedback control loops are implemented as

individual cases from scratch. Significant efforts are needed to develop and tune the feedback control loop in every case. No middleware has been developed in the above projects to provide general support for composition and tuning of software feedback control loops.

The SoftBus architecture in ControlWare is related to distributed middleware such as CORBA [15] and DCOM [11]. Similar to the location-transparent method invocation in CORBA and DCOM, SoftBus supports plug-and-play control loop components (i.e., monitors, controllers, and actuators). However, unlike CORBA and DCOM, SoftBus provides direct support for feedback control by supporting active and passive interaction mechanisms among monitors, controllers, and actuators, as well as general resource managers for server systems. Furthermore, ControlWare provides system identification and controller tuning services that are not supported by other common middleware services.

The SWiFT project [14] at OGI and the Agilos project [17] at UIUC share some similar goals with ControlWare. SWiFT is a toolkit for constructing feedback control loops from libraries. It also supports the visualization and simulation of software control loops. The Agilos project constructs a middleware to support QoS control and adaptation. ControlWare is different from SWiFT and Agilos in its unique SoftBus architecture that enables flexible plug-and-play control components in a location independent fashion (e.g., components of a same control loop can be from different address spaces and even remote nodes). In comparison, Agilos has a fixed two level control structure and a fixed set of monitors (e.g., CPU and bandwidth monitors). More importantly, ControlWare is the first middleware that provides end-to-end support of the whole development process of creating QoS control in software systems. This process includes defining QoS contracts, mapping contracts to feedback control loops, system identification, controller tuning, and implementation. For example, neither SWiFT or Agilos supports the mapping from QoS contracts to feedback control loops or system identification. They also do not provide the generic resource manager as in ControlWare.

7. Conclusions

In this paper, we described a new middleware architecture for QoS guarantees in distributed environments such as the Internet. The architecture implements a new paradigm for QoS control, which is especially suitable for systems operating in highly uncertain environments or when accurate system load and resource models are not available. Our preliminary evaluation of the ability of this middleware to provide advertised guarantees in the context of selected different applications illustrates the promise of this approach. While prior efforts have been made to apply control theory

to QoS control, ours is the first comprehensive middleware service that incorporates these principles under a clear well defined set of APIs which substantially reduces the development effort of performance assured applications and Internet services.

Future work of the authors will focus on understanding the limitations of a control-theoretic approach and deriving new guarantee semantics. A possible disadvantage of using feedback only as a means to correct performance is the need for a performance error to occur first before a feedback controller can respond. In the future, we shall focus on mechanisms that combine prediction with feedback to improve convergence to specifications in a highly dynamic unpredictable system. We shall also extend the middleware to allow fully dynamic online re-configuration during normal system operation, and investigate other types of performance guarantees that might be achievable in a feedback control context. For example, it may be interesting to cast adaptive guarantees on service availability, security, and fault-tolerance as feedback control problems.

Acknowledgments

The authors would like to thank Sang Son, Gang Tao, and Ying Lu for providing useful comments.

References

- [1] T. Abdelzaher. An automated profiling subsystem for qos-aware services. In *IEEE Real-Time Technology and Applications Symposium*, Washington, D.C., June 2000.
- [2] T. Abdelzaher, A. Shaikh, F. Jahanian, and K. Shin. RT-CAST: Lightweight multicast for real-time process groups. In *IEEE Real-Time Technology and Applications Symposium*, Boston, Massachusetts, June 1996.
- [3] T. F. Abdelzaher and N. Bhatti. Web server QoS management by adaptive content delivery. In *International Workshop on Quality of Service*, London, UK, June 1999.
- [4] T. F. Abdelzaher and N. T. Bhatti. Web content adaptation to improve server overload behavior. *WWW8 / Computer Networks*, 31(11-16):1563–1577, 1999.
- [5] T. F. Abdelzaher, K. G. Shin, and N. Bhatti. Performance guarantees for web server end-systems: A control-theoretical approach. *IEEE Transactions on Parallel and Distributed Systems*, 2001. Accepted.
- [6] Y. Amir, L. Moser, P. Melliar-Smith, D. Agarwal, and P. Ciarfella. The Totem single-ring ordering and membership protocol. *ACM Transactions on Computer Systems*, 13(4):311–342, November 1995.
- [7] K. J. Astrom and B. Wittenmark. *Adaptive Control*, chapter 2. Addison Wesley, 2nd edition, 1995.
- [8] P. Barford and M. Crovella. Generating representative web workloads for network and server performance evaluation. In *Measurement and Modeling of Computer Systems*, pages 151–160, 1998.
- [9] A. Chatha. Fieldbus: the foundation for field control systems. *Control Engineering*, 41(6):47–50, May 1994.

- [10] N. Christin, J. Liebeherr, and T. Abdelzaher. A quantitative assured forwarding service. Technical Report CS Technical Report 2001-21, University of Virginia, 2001.
- [11] M. Corporation. Distributed component object model protocol - dcom/1.0, 1998.
- [12] A. S. Foundation. <http://www.apache.org>.
- [13] N. Gandhi, S. Parekh, J. Hellerstein, and D. Tilbury. Feedback control of a lotus notes server: Modeling and control design. In *American Control Conference*, 2001.
- [14] A. Goel, D. Steere, C. Pu, and J. Walpole. Swift: A feedback control and dynamic reconfiguration toolkit, 1999.
- [15] O. GROUP. The common object request broker: Architecture and specification, 1995.
- [16] C. Holot, V. Misra, D. Towsley, and W. Gong. A control theoretic analysis of red, 2000.
- [17] B. Li and K. Nahrstedt. A control-based middleware framework for quality of service adaptations, 1999.
- [18] C. Lu, T. Abdelzaher, J. Stankovic, and S. Son. A feedback control approach for guaranteeing relative delays in web servers. In *IEEE Real-Time Technology and Applications Symposium*, Taipei, Taiwan, June 2001.
- [19] C. Lu, J. A. Stankovic, T. F. Abdelzaher, G. Tao, S. H. Son, and M. Marley. Performance specifications and metrics for adaptive real-time systems. In *IEEE Real-Time Systems Symposium*, Orlando, FL, December 2000.
- [20] C. Lu, J. A. Stankovic, G. Tao, and S. H. Son. Feedback control real-time scheduling: Framework, modeling, and algorithms. *Real-Time Systems Journal*, Special Issue on Control-Theoretical Approaches to Real-Time Computing, March-April, 2002.
- [21] Y. Lu, A. Saxena, and T. F. Abdelzaher. Differentiated caching services; a control-theoretical approach. In *International Conference on Distributed Computing System*, Phoenix, Arizona, April 2001.
- [22] W. A. McEachern. *Economics*. South-Western College Publishing, 5th edition, 1999.
- [23] C. Mercer, S. Savage, and H. Tokuda. Processor capacity reserves: Operating system support for multimedia applications. In *Proceedings of the IEEE International Conference on Multimedia Computing and Systems*, May 1994.
- [24] S. Parekh, N. Gandhi, J. L. Hellerstein, D. Tilbury, T. S. Jayram, and J. Bigus. Using control theory to achieve service level objectives in performance management. In *IFIP/IEEE International Symposium on Integrated Network Management*, 2001.
- [25] L. Rodrigues, P. Veríssimo, and J. Rufino. A low-level processor group membership protocol for LANs. In *Proc. Int. Conf. on Distributed Computer Systems*, pages 541–550, 1993.
- [26] S. W. P. Server. <http://www.squid-cache.org>.
- [27] L. Sha, R. Rajkumar, and S. S. Sathaye. Generalized rate monotonic scheduling theory: A framework for developing real-time systems. *Proceedings of the IEEE*, 82(1):68–82, January 1994.
- [28] F. G. Shinskey. *Process control systems: application, design, and tuning*. McGraw-Hill, New York, 4th edition edition, 1996.
- [29] K. Skadron, T. Abdelzaher, and M. Stan. Control-theoretic techniques and thermal rc modeling for accurate and localized dynamic thermal management. In *International Symposium on High Performance Computer Architecture*, Cambridge, MA, February 2002.
- [30] D. C. Steere, A. Goel, J. Gruenberg, D. McNamee, C. Pu, and J. Walpole. A feedback-driven proportion allocator for real-rate scheduling. In *Operating Systems Design and Implementation*, pages 145–158, 1999.

Appendix A: Contract Description Language

The Contract Description Language (CDL) is used to describe desired convergence guarantees. Its main syntax is as follows.

```

GUARANTEE NAME {
    GUARANTEE_TYPE = type;
    TOTAL_CAPACITY = capacity;
    CLASS_0 = QoS_0;
    CLASS_1 = QoS_1;
    . . . . .
    CLASS_num = QoS_num;
}

```

GUARANTEE_TYPE: The GUARANTEE_TYPES currently supported by ControlWare include ABSOLUTE, RELATIVE, and STATISTICAL_MULTIPLEXING. Different guarantee types need different interpretations of the QoS values and are mapped to different feedback control loops as described in Section 2.2. Although utility optimization is not listed as a guarantee type, it is equivalent to absolute guarantees because it is mapped to single feedback control loop per class.

TOTAL_CAPACITY: total capacity is only useful when GUARANTEE_TYPE = STATISTICAL_MULTIPLEXING. The set point of the best effort server is the total capacity minus the capacity allocated to all guaranteed service classes.

CLASS.i: Each service class represents a category of requests with a guarantee depending on the application requirements. For example, a service class can be all the HTTP requests from premium clients. The assignment CLASS.i = QoS.i specifies the guaranteed QoS for class i. Note that the guaranteed QoS have different meanings for different guarantee types. For ABSOLUTE and STATISTICAL_MULTIPLEXING guarantees, QoS.i represents the absolute value for desired QoS, while RELATIVE guarantees are only interested only the relative value (ratio) between the QoS.i's of different classes.