

I/O-Aware Deadline Miss Ratio Management in Real-Time Embedded Databases

Woochul Kang, Sang H. Son, and John A. Stankovic
Department of Computer Science
University of Virginia
{wk5f,son,stankovic}@cs.virginia.edu

Mehdi Amirijoo
Department of Computer and Information Science
Linköping University
meham@ida.iuu.se

Abstract

Recently, cheap and large capacity non-volatile memory such as flash memory is rapidly replacing disks not only in embedded systems, but also in high performance servers. Unlike disks, the access time of flash memory is not affected by mechanical parts, thus the access time is highly predictable. However, in real-time embedded databases deadline misses may occur if data objects in flash memory are not properly managed. Buffer cache can be used to mitigate the problem. However, since the workload of a real-time database cannot be precisely predicted, it may not be feasible to provide enough buffer space to satisfy all timing constraints. Several deadline miss ratio management schemes have been proposed, but they do not consider I/O activities. In this paper, we present an I/O-aware deadline miss ratio management scheme in real-time embedded databases whose secondary storage is flash memory. We propose an adaptive I/O deadline assignment scheme, where I/O deadlines are derived from up-to-date system status. We also present a deadline miss ratio management architecture where a control theory-based feedback control loop prevents resource overload both in I/O and CPU. A simulation study shows that our approach can effectively cope with both I/O and CPU overload to achieve the desired deadline miss ratio.

1 Introduction

A number of emerging applications require real-time data services to handle large amounts of data in timely fashion. Examples include online stock trading, traffic control, target tracking, and network management. Unlike traditional applications, they require transactions to be completed within their deadlines. Several approaches using real-time database systems (RTDBS) have been proposed to handle these time-constrained transactions [4][17][11]. Most of them are based on main-memory databases due to the inherent unpredictability of disk I/O. In main-memory databases, it is assumed that the database is small enough to fit into main memory, thus eliminating most of I/O operations. However, this assumption does not hold when the database size is bigger than main memory. Another problem with main-memory databases is that they are vulnerable to system failures. Recent advancement in semiconductor technology enables us to overcome the unpredictability of disks with cheap and large capacity non-volatile memories. In particular, flash memory is rapidly replacing disks not only in real-

time embedded systems, but also in high-performance servers [2]. Unlike disks, the access time to flash memory is not affected by mechanical parts, thus, flash memory is several orders of magnitude faster and highly predictable. These desirable characteristics make flash memory more suitable for RTDBS. However, flash memory still has high overhead in access time compared to volatile memory such as SRAM. Table 1 shows the overhead of flash memory in comparison to SRAM.

Type	Read	Write	Erase
SRAM	10ns	10ns	N/A
NAND Flash	10 μ s	200 μ s	2ms

Table 1. Characteristics of memory devices [14].

Such gaps in access time is critical for RTDBS because high I/O overheads may incur deadline misses if data objects are not properly managed. System designers may provide buffer memory to mitigate the high overhead of flash memory accesses. However, in many cases, the workloads are unknown at design time and they can change dynamically. Therefore, it may not be feasible to provide enough buffer space to satisfy all timing constraints. In particular, resource-constrained embedded systems are extremely costly and space sensitive, and cannot afford to have large buffers.

In this paper, we propose an I/O-aware deadline miss ratio management scheme for a real-time embedded database system (RTEDBS) whose secondary storage is flash memory. The contributions of this paper are three-fold:

1. an adaptive I/O deadline assignment scheme,
2. a model of deadline miss ratio in terms of I/O and CPU workloads, and
3. a feedback control architecture to satisfy a given deadline miss ratio.

To the best of our knowledge, this is the first paper on deadline miss ratio management of RTEDBS with flash storage using feedback control to consider both I/O and CPU workloads. Previous approaches on deadline miss ratio management in RTDBSs assumed main-memory databases [11][5]. As a result they only considered CPU workloads, which is insufficient when managing the Quality of Service (QoS) of modern RTEDBS consisting of volatile memory as well as non-volatile flash memory. Several I/O-aware approaches [3][7] have been proposed. However, they assumed disks as a secondary storage, and their primary research focus was deadline-driven disk scheduling to mitigate the unpredictability of disk operations.

A key issue in deadline miss ratio management with different kinds of resources, e.g. I/O and CPU, is to find out which resource is the bottleneck that causes deadline misses. In RTEDBS, a deadline miss can happen either because of an overload in I/O or CPU, or both. By properly setting deadlines for each resource request and observing its deadline misses, we can tell which resource is the bottleneck. However, deriving a deadline for each resource is not straightforward because the deadlines are set for transactions, not for individual resource requests. In this paper, we assume a 2-phase transaction model, where each transaction consists of an I/O phase and a computation phase. The I/O deadlines and subsequent CPU deadlines for respective I/O phases and CPU phases are derived from transaction deadlines with up-to-date system overload status. When either I/O or CPU is overloaded, the related deadline is adjusted accordingly to give more time to the overloaded resource to complete operations.

Having I/O deadlines and CPU deadlines enables us to measure and model the system in terms of I/O and CPU workloads and their respective deadline miss ratios. A straightforward approach would be to build separate models for I/O and CPU. However, our experiments show that CPU and I/O deadline miss ratios are coupled and affect each other, thus, necessitating multiple-input/ multiple-output (MIMO) modeling of the system. In this paper, the RTEDBS is modeled as a MIMO system to capture this coupling of control inputs and system outputs.

Using feedback controllers has shown to be effective for real-time systems with unpredictable workload [5][11][13]. Therefore, a feedback control architecture is proposed in this paper to guarantee the desired deadline miss ratio. At each sampling instant, the feedback control loop measures I/O and CPU miss ratios and computes control signals, e.g., the required I/O and CPU workload adjustment. In particular, our approach controls both I/O and CPU workloads at the same time because of close interactions between them.

A set of experiments were performed to evaluate the performance of the proposed scheme. Our evaluation results show that our approach gives robust and controlled behavior under a variety of workloads and access patterns compared to the baseline algorithm.

The rest of the paper is organized as follows. Section 2 describes our system and data model for RTEDBS. In Section 3, our deadline miss ratio management architecture is described. In section 4, the performance evaluation results are presented. In section 5, we present the related work. Finally, Section 6 concludes the paper and discusses future work.

2 System and Data Model for RTEDBS

2.1 System Model

A real-time embedded system having a CPU, main memory, and flash memory is considered. Figure 1 shows a typical H/W configuration of a real-time embedded system. The flash memory is used for persistent data storage.

The buffer pool is a cache between the flash memory and the CPU. It is shared between transactions to reduce the data storage access time. The maximum size of the buffer pool is

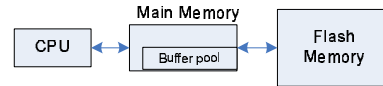


Figure 1. A H/W configuration.

configured at deployment time and does not change over its lifetime. However, the proportion of allocated buffer pool size between different classes is not fixed since a fixed allocation can incur inefficient use of the buffer resource.

The I/O load between main memory and flash memory occurs only when an explicit I/O request is issued for a data object not present in the buffer pool. Implicit I/O requests such as page faults are not considered because virtual memory is typically not supported in embedded systems.

2.2 Data and Transaction Model

In our data model, data objects can be classified into two classes, temporal and non-temporal data. Temporal data objects are updated periodically by update transactions. A temporal data object d_i is considered valid and fresh as long as the current time is not later than the timestamp of d_i plus the validity interval of d_i . Update transactions are write-only. In contrast to update transactions, user transactions may read both temporal and non-temporal data objects and modify non-temporal data objects.

Because embedded platforms are assumed for our RTEDBS, transactions are *canned transactions*, whose characteristics including data requirement and worst-case computation time is known at the design time. However, workload and data access patterns of the whole RTEDBS can be unpredictable and change dynamically because the invocation frequency of each transaction is unknown. Since data requirements are known for each transaction, data requests of each transaction can be gathered before its computation to improve the response time. To this end, we model each transaction as a two-phase operation, an I/O phase and a computation phase. In the I/O phase, data objects for the transaction are brought to the buffer pool from the flash memory. If all data objects are already present in the buffer pool, the I/O phase is skipped. In a single transaction, the computation phase can begin only after all its required data objects are present in the buffer pool. However, the I/O and the computation phase of different transactions can overlap. For example, while transaction i is under the I/O operation, transaction j can perform its computation. A transaction can commit after the computation phase by updating a copy of the data object in the memory buffer. The time required to update the data object in memory buffer is ignored in this transaction model because it is relatively small compared to flash memory accesses. The buffer manager will *eventually* write the updated buffers back to flash memory when the buffer manager is running out of buffers.

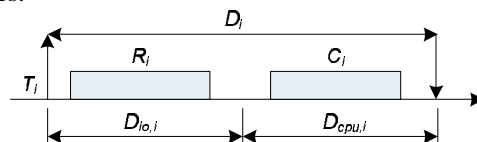


Figure 2. The timing of a typical transaction.

Figure 2 shows the timing of a typical transaction i , where T_i is a release time, D_i is a relative deadline, R_i is a I/O time, C_i is a computation time, $D_{io,i}$ is a relative I/O deadline, and $D_{cpu,i}$ is a relative CPU deadline. While the deadline of each transaction i is set by the application in consideration of the worst case I/O time and computation time, $D_{io,i}$ and $D_{cpu,i}$ are dynamically derived from the relative deadline D_i and current deadline miss ratio. When the I/O deadline is missed, the transaction is aborted and its computation phase is not initiated. This dynamic I/O deadline assignment scheme is explained in the section 3.2.

Transactions are classified into classes by their importance. Two service classes, a guaranteed service class and a best effort service class, are assumed in this study for simplicity. Resources including the buffer space is shared between service classes. Instead of providing each service class fixed amount of resources, the allocation of resources to each service class is adaptive in our RTEDBS. For example, when the I/O of a guaranteed service class is overloaded, the buffer space of the best-effort service class can be utilized to relieve the overload. Our adaptive buffer space allocation scheme is presented in Section 3.4.

3 Approach

A system overload due to a transient surge of the workload is the main source of deadline misses in soft real-time systems. Overload can occur in either I/O or CPU, or both. When the number of deadline misses increases because of the scarcity of a specific resource, the feedback control loop for the resource reduces the workload by adjusting the system parameters. In this section, we present a QoS management architecture which controls both CPU and I/O workloads with a feedback control.

3.1 Performance Metrics

In our approach, the main performance metric is the deadline miss ratio of real-time transactions. A deadline miss can happen due to either an I/O deadline miss or a CPU deadline miss.

CPU deadline miss: A transaction misses its deadline and all data objects needed by the transaction are present in the buffer pool by the I/O deadline.

I/O deadline miss: Some data objects needed by the transaction are not present in the buffer pool at the I/O deadline. In this case, the computation phase is not initiated.

When a deadline miss happens, it belongs to either a CPU or an I/O deadline miss, but not both. Accordingly, CPU and I/O deadline miss ratios are obtained as follows,

$$m_{cpu} = \frac{\text{\# of CPU deadline misses in admitted transactions}}{\text{\# of admitted transactions}} \quad (1)$$

$$m_{io} = \frac{\text{\# of I/O deadline misses in admitted transactions}}{\text{\# of admitted transactions}}, \quad (2)$$

$$\text{where, total miss ratio} = m_{cpu} + m_{io}. \quad (3)$$

By observing the deadline miss ratios for CPU and I/O, we can tell which resource is the current bottleneck in the RTEDBS.

The desired levels of miss ratio, m_{ref} , for the guaranteed service class is expressed in the QoS specification. The desired level of deadline miss ratio for CPU, $m_{ref,cpu}$, and I/O, $m_{ref,io}$, are set separately such that their sum is equal to m_{ref} . The reference deadline miss ratios for I/O and CPU are weighted according to system overload status as follows,

$$m_{ref,cpu}(t) = \rho_{cpu}(t) \times m_{ref} \quad (4)$$

$$m_{ref,io}(t) = \rho_{io}(t) \times m_{ref} \quad (5)$$

where,

$$\rho_{cpu}(t) = \frac{\text{Total number of CPU deadline misses}}{\text{Total number of deadline misses}} \quad (6)$$

$$\rho_{io}(t) = \frac{\text{Total number of I/O deadline misses}}{\text{Total number of deadline misses}}. \quad (7)$$

Note that $\rho_{cpu}(t)$ and $\rho_{io}(t)$ are defined only when the RTEDBS has deadline misses, and the sum of $\rho_{cpu}(t)$ and $\rho_{io}(t)$ equals to one.

3.2 Adaptive I/O Deadline Assignment

Because a transaction deadline is set for the transaction, not for I/O or computation phases, deriving an I/O deadline is not straightforward. In this paper, instead of setting I/O deadlines statically for each transaction, we define I/O deadlines recursively as a time-varying function of the deadline miss ratio of the past sampling period as follows,

$$\text{I/O deadline for transaction } i = T_i + D_i - C_i \times \left(1 + m_{cpu}(t) \times \frac{D_i - C_i}{C_i}\right), \quad (8)$$

where $m_{cpu}(t)$ is the CPU deadline miss ratio during the past sampling period. Because of (3), the definition of I/O deadline is recursive. The initial value for $m_{cpu}(0)$ is set to zero. The I/O deadline reflects up-to-date resource status and how much slack time is required for the computation phase to finish within the transaction deadline; In our two-phase transaction model, the deadline for the computation phase is equal to the transaction deadline.

An I/O deadline can be as long as $T_i + D_i - C_i$ when the CPU is not overloaded, thus, m_{cpu} is close to zero; this implies that the computation phase requires less slack time to meet the transaction deadline. In contrast, an I/O deadline can be as short as T_i when the CPU resource is the bottleneck; this implies that the computation phase requires more slack time. Overall, I/O deadlines are set inverse proportionally to the CPU deadline miss ratio of the latest sampling period.

Setting I/O deadlines can serve two purposes; time-cognizant I/O scheduling and I/O workload control. In this paper, we use I/O deadlines only for I/O workload control purpose. A *First-come/first-service* scheduling policy is used for I/O scheduling. We reserve the study on the impact of the time-cognizant I/O scheduling as our future work.

3.3 Deadline Miss Ratio Management Architecture

Figure 3 shows the deadline miss ratio management architecture for RTEDBS. Transactions issued by sensors (update transactions) and users (user transactions) are placed on the ready queue. The transaction queue can have several service classes. The figure shows the two classes, a guaranteed service queue and a best effort service queue. The transactions in the best effort service queue are dispatched only if the ready queue for the guaranteed service class is empty. The dispatched transactions are managed by the transaction handler which consists of buffer manager (BM), freshness manager (FM), concurrency control (CC), and scheduler (S). Transactions are monitored by the monitor and the statistics of monitored transactions including deadline miss ratios, (m_{io} and m_{cpu}), and utilizations (u_{io} and u_{cpu}), are reported to the QoS controllers on every sampling period.

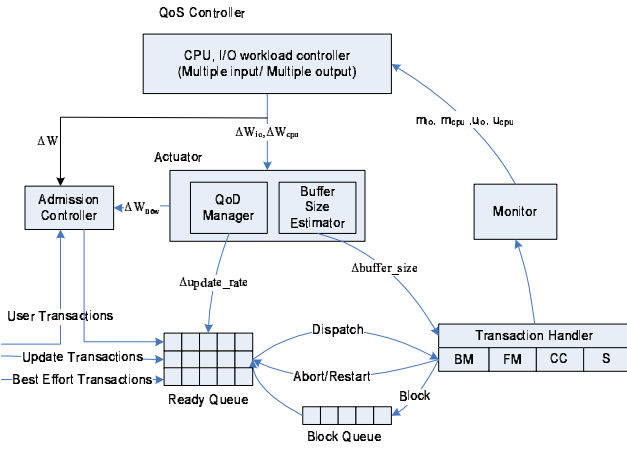


Figure 3. Feedback control architecture

Controllers calculate the workload adaptations required to meet the reference miss ratios, $m_{ref,io}$ and $m_{ref,cpu}$, by comparing them to current miss ratios. Once the workload changes are determined, they are enforced by two actuators independently. For CPU workload adaptation, the update rates of temporal data are adjusted. For I/O workload adaptation, the buffer size of the guaranteed service is adjusted. If the workload needs to be adjusted more than the actuators can handle, the admission controller adapts the workload by allowing or denying user transactions.

3.4 I/O and CPU Workload Adjustment

In I/O-aware deadline miss ratio management, the workload should be separately defined for I/O and CPU. The CPU workload, W_{cpu} , can be measured as the amount of requested computation to do at a given time. However, we need to be more specific about I/O workload because the I/O workload of a transaction depends on the buffer cache; in the presence of buffer cache, I/O requests from a transaction incur I/O operations only if data objects are not found in the buffer. In our study, I/O workload, W_{io} for the past sampling period is estimated as follows,

$$W_{io} = \frac{\# \text{ of buffer misses} \times \text{buffer page size}}{\text{bandwidth} \times \text{sampling period}}. \quad (9)$$

Note that we consider explicit I/O requests and subsequent buffer misses as the only source of I/O workload. Because I/O workload is determined only by the number of buffer misses, we can adjust the I/O workload by adapting the buffer hit ratio, bhr . When a random access pattern is assumed, the bhr is a linear function of the buffer size and can be expressed as follows,

$$bhr(b) = c_1 \times b + c_2, \quad (10)$$

where b is the buffer size, and c_1 and c_2 are constants. The constants c_1 and c_2 are obtained by profiling. When the target I/O workload is $W_{target,io}$, (9) can be rewritten with bhr as follows,

$$W_{target,io} = \frac{\# \text{ of buffer access} \times (1 - bhr(b)) \times \text{buffer page size}}{\text{bandwidth} \times \text{sampling period}}. \quad (11)$$

Substituting bhr to (10) gives the desired buffer size,

$$b_{desired} = \frac{1}{c_1} \times \left(1 + c_2 - \frac{W_{target,io} \times \text{bandwidth} \times \text{sampling period}}{\# \text{ of buffer access} \times \text{buffer page size}} \right). \quad (12)$$

Therefore, the difference between the current buffer size and $b_{desired}$ is the amount of the buffer size to be adjusted to achieve the target I/O workload.

It may be argued that adaptive buffer adjustment is unnecessary because providing larger buffer space will solve I/O overload problem. However, memory is a scarce resource in most embedded systems and we cannot afford to have large buffer memory. Therefore, we need to increase the total utilization of the buffer memory by adaptively allocating buffer space to each service class. For instance, by reducing the buffer space of a service class when its I/O is under-utilized, the RTEDBS can provide more buffer space to the other service classes that is in need of more buffer space.

In contrast to the I/O workload, the CPU workload can be adjusted by changing the precision of transactions [5] or the freshness of temporal data [11]. In this paper, we change the freshness of temporal data by changing update intervals of temporal data. For details, readers are referred to [11].

If further workload adaptation is not possible by changing update intervals and the buffer size, admission control is applied. By allowing or denying more user transactions, workloads are adjusted. However, we should be careful in applying admission control because it changes both I/O and CPU workload together. Aggressive admission control can make system unstable. For instance, if I/O allows 50% additional user transaction while CPU needs only 10% additional user transactions, then allowing 50% additional user transaction can cause overshoot in CPU miss ratio. To prevent excessive overshoot, we take conservative approach by taking the average of two. Another interesting issue in applying admission control occurs when each resource wants to adjust workload in different directions; e.g., CPU wants to increase its workload by allowing 20% additional user transactions while I/O wants to decrease its workload by 20%. In this case, the admission rate is determined by the resource that has higher deadline miss ratio.

3.5 Control Loop Design

In this section, we model RTEDBS in terms of I/O and CPU deadline miss ratios and build a controller to manage deadline misses.

3.5.1 System Modeling

The first step in the design of a feedback control loop is the modeling of the controlled system [9]; the RTEDBS in our study.

Unlike previous work [5][11], which have single-input, single-output (SISO), the RTEDBS in this paper has multiple inputs (W_{cpu} and W_{io}) and multiple outputs (m_{cpu} and m_{io}). We may choose to use two separate SISO models for each pair of control input and system output; (W_{cpu} , m_{cpu}) and (W_{io} , m_{io}), respectively. However, if an input of the system is highly affected by another input, then a Multiple Inputs/Multiple Outputs (MIMO) model should be considered [8]. Having two SISO models does not capture the interaction between different control inputs and system outputs; for example, the interaction between I/O workload and CPU deadline miss ratio cannot be modeled with two SISO models.

To understand the interaction between control inputs and system outputs in the RTEDBS, we performed a series of experiments by applying a discrete sine wave input while the other input was fixed. We did it for both CPU workload and I/O workload. The workloads were adjusted by controlling the buffer size and update rates of temporal data for I/O and CPU respectively. Figure 4(a) shows the result of applying sine wave CPU workload with 40% amplitude while the I/O workload was fixed at 120%. While CPU workload changes between 80% and 120%, I/O workload stays at around 110% without a significant deviation. This result shows that I/O workload is not affected by CPU workload change. On the contrary, an I/O workload change has high impact on CPU workload as shown in Figure 4(b). In Figure 4(b), a discrete sine wave input was applied to I/O while the CPU workload was fixed at 130%. Even though we fixed the CPU workload, it was affected by changes to the I/O workload, and consequently the CPU deadline miss ratio was also affected by it. Interestingly enough, the results show that CPU workload is inversely proportional to I/O workload. This is because transactions are aborted when their I/O deadlines are missed. The higher the I/O deadline miss ratio is, the more transactions are aborted by I/O deadline misses. Therefore, CPU workload decreases, thus, decreasing the CPU deadline miss ratio. This experiment shows that a MIMO model is more appropriate for RTEDBS than SISO models.

Another issue in modeling a computing system is its non-linearity and time-variant characteristics. Complex systems such as RTEDBS can show a non-linear response to inputs. For example, the CPU deadline miss ratio develops quite differently when the CPU is saturated from when it is not saturated. However, the system can be approximated quite closely with linear time invariant models such as the ARX model by choosing an operating region where the system's response is approximately linear [9]. In case of RTEDBS, the operating region of the controller is set so that the deadline miss ratio is bigger than 0

and less than 50%; the operating region of inputs, CPU and I/O workloads, can be obtained by inverting the MIMO model. The form of linear time invariant model for the RTEDBS is shown in (13), with parameters \mathbf{A} and \mathbf{B} .

$$\begin{pmatrix} m_{io}(k+1) \\ m_{cpu}(k+1) \end{pmatrix} = \mathbf{A} \cdot \begin{pmatrix} m_{io}(k) \\ m_{cpu}(k) \end{pmatrix} + \mathbf{B} \cdot \begin{pmatrix} W_{io}(k) \\ W_{cpu}(k) \end{pmatrix} \quad (13)$$

Because the RTEDBS is modeled as a MIMO system, \mathbf{A} and \mathbf{B} are 2x2 matrices. A RTEDBS simulator which will be introduced in Section 4 was used for *system identification* [12]. In the system identification, relatively prime sine wave workloads for CPU and I/O were applied simultaneously to get the parameters. In our study, the RTEDBS model has $\mathbf{A} = \begin{pmatrix} 0.5403 & 0.1740 \\ 0.2405 & 0.4500 \end{pmatrix}$, and $\mathbf{B} = \begin{pmatrix} 0.2400 & -0.1208 \\ -0.0100 & 0.1774 \end{pmatrix}$ as its parameters.

In terms of system order, note that we model the RTEDBS as a first-order system; the current outputs are determined by their inputs and outputs of the last sample. As we will show later, the accuracy of the model is satisfactory and, hence, the chosen model order is sufficient for our purposes.

The model can be validated by comparing the experimental result to what the model predicts. Figure 5 plots the experimental response of the RTEDBS and the prediction of the model. We can see that the model gives highly accurate prediction.

3.5.2 Controller Design

For RTEDBS, we choose to use a proportional integral (PI) control function given by,

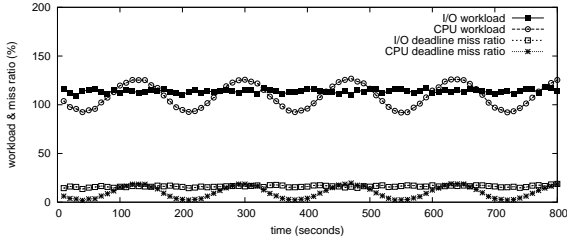
$$U(k) = K_p \cdot E(k) + K_I \cdot \sum_{j=1}^{k-1} E(j). \quad (14)$$

At each sampling instant k , the controller computes the control input $U(k) = [W_{IO}(k) \ W_{CPU}(k)]^T$ by monitoring the control error $E(k) = [m_{ref,IO}(k) - m_{IO}(k) \ m_{ref,CPU}(k) - m_{CPU}(k)]^T$. K_p and K_I are controller gains.

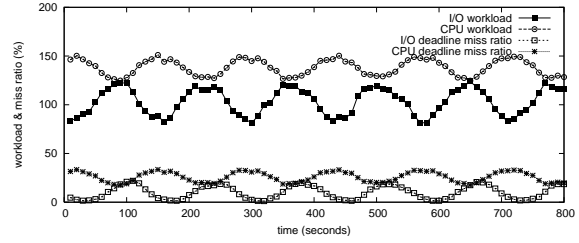
One important design consideration in computing systems such as RTEDBS which have a stochastic nature is to control the trade-off between short settling times and overreacting to random fluctuations. If a controller is too aggressive, then the controller over-reacts to this random fluctuation. To this end, we choose to use the linear quadratic regulator (LQR) technique to find control gains, which is accepted as a more general technique for MIMO systems [9]. In LQR, control gains are set to minimize the quadratic cost function,

$$J = \sum_{k=0}^{\infty} [E(k) \ V(k)] \cdot Q \cdot \begin{pmatrix} E(k) \\ V(k) \end{pmatrix} + U(k)^T \cdot R \cdot U(k). \quad (15)$$

The cost function includes the control errors $E(k)$, accumulated errors $V(k)$, and weighting matrices Q and R . LQR allows us to better negotiate the trade-offs between speed of response and over-reaction to random fluctuation by selecting appropriate Q and R matrices. Q quantifies the cost of control errors and R quantifies the cost of control effort. Since controlling I/O workload by changing the buffer size incurs higher cost than controlling CPU workload by changing update interval, we impose higher weight on I/O. We choose $R = \text{diag}(1/3,$



(a) sine wave CPU workload



(b) sine wave I/O workload

Figure 4. SISO inputs to RTEDBS.

1/8) where 1/3 is the cost of I/O control and 1/8 is the cost of CPU control. Then, we choose $Q = \text{diag}(0.1, 0.1, 0.001, 0.001)$ to weight the control errors more heavily than the integrated control errors. For more details on LQR technique, readers are referred to [9].

Finally, in terms of sampling interval, we sample every 10 seconds. In RTEDBSs, the buffer management affects the choice of sampling interval in particular because the buffer hit ratio changes slowly after adjusting the buffer size. If the sampling interval is too short, controlling buffer size may not make effect until the next sampling period, thus, wasting the control effort.

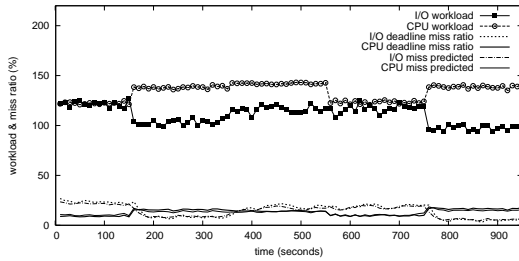


Figure 5. Model validation.

4 Experiment

The main objective of the experiment is to test the effectiveness of controlling I/O and CPU workloads together instead of considering only one in the presence of unpredictable workloads. A scheme which is not I/O-aware is compared to our scheme. For experiments, we developed a simulator that models the proposed RTEDBS. Various workloads were applied to the simulator to test its performance.

4.1 Simulation Settings

The simulated workload consists of sensor data update transactions and user transactions. User transaction workloads are synthesized to test wide range of workloads and patterns. Update transaction workloads follow similar settings as in [11]. I/O components including NAND flash memory are modeled to follow the performance characteristics typically found in commercial products [2].

4.1.1 Data and Update Transactions

3,000 temporal data objects and another 3,000 non-temporal data objects reside in the database. Among them, the up-

Parameter	Value
# of temporal data objects	3000
Update interval (P_i)	$Uniform(100ms, 50sec)$
EET_i	$Uniform(2ms, 4ms)$
Actual exec. time	$Normal(EET_i, \sqrt{EET_i})$
Relative deadline	$2 \times P_i$
# data object access/update	1
Update CPU load	$\approx 50\%$
Update I/O load	$\ll 1\%$

Table 2. Update transaction settings.

Parameter	Value
# of non-temporal data objects	3000
$EECT_i$	$Uniform(3ms, 5ms)$
Actual exec. time	$Normal(EECT_i, \sqrt{EECT_i})$
EET_i	$NUM_{data} \times ReadAccessTime/page$
Relative deadline	$(EECT_i + EET_i) \times \text{slack factor}$
Slack factor	$Uniform(5, 10)$
$NUM_{data}(I/O \text{ intensive})$	$Normal(150, 30)$
$NUM_{data}(\text{balanced})$	$Normal(100, 20)$
$NUM_{data}(\text{CPU intensive})$	$Normal(50, 10)$

Table 3. User transaction settings.

Parameter	Value
Read access time / page	$30\mu s$
Write access time	$300\mu s$
Erase time	N/A
Page size	512 bytes
Bank interleaving	N/A
Flash memory max. bandwidth	16MB/sec
BUS bandwidth	∞
Buffer size	30 - 3000 data objects
Buffer replacement algorithm	LRU

Table 4. Flash memory and buffer settings.

date stream updates only temporal data. The update period, p_i , follows a uniform distribution $Uniform(100ms, 50sec)$. The expected execution time (EET) of an update transaction is uniformly distributed in the range (3ms, 6ms). The actual execution time is given by normal distribution $Normal(EET_i, \sqrt{EET_i})$. The relative deadline of an update is set to $2 \times p_i$. An update transaction incurs I/O operations if its target data object is not found in the buffer. However, I/O deadlines for update transactions are not set because update transactions are not I/O-intensive. The default settings shown in Table 2 generate about 50% CPU load and less than 1% I/O load when the buffer can hold 50% of total temporal data objects.

4.1.2 User Transaction

A user transaction accesses both temporal and non-temporal data and updates only non-temporal data. The arrival rate of

user transactions to the database follows the Poisson distribution. In terms of the number of data accesses per transaction (NUM_{data}), three settings are tested; I/O intensive settings where NUM_{data} is given by $Normal(150, 30)$, balanced settings whose NUM_{data} follows $Normal(100, 20)$, and CPU intensive settings where NUM_{data} is given by $Normal(50, 10)$. The execution of user transactions consists of an I/O phase and a computing phase. The expected execution time ($EECT_i$) of the computation phase is given by the $Uniform(3ms, 5ms)$. The actual execution time follows Normal ($EECT_i, \sqrt{EECT_i}$). The expected execution time of the I/O phase ($EEIT_i$) is given by the multiplication of NUM_{data} and the read access time to flash memory per page. The actual execution time of the I/O phase is determined by the number of data objects found in the buffer. The deadline of a user transaction is $(EECT_i + EEIT_i) \times \text{slack factor}$. The slack factor is uniformly distributed ranging from 5 to 10. In I/O intensive settings, user transactions result in more I/O workload than CPU workload; when 100 user transactions arrive per second in addition to default update transactions, about 55% and 30% additional I/O and CPU loads are incurred respectively. In CPU intensive settings, 18% I/O load and 30% CPU load are incurred. In balanced settings, 35% I/O load and 30% CPU load are incurred. In experiments, the arrival rates are adjusted to increase or decrease I/O and CPU load.

4.1.3 Flash Memory and Buffer

When a transaction accesses data objects, the buffer pool in the main memory is first searched, and if not found, the data objects in the persistent storage are brought to the buffer in main memory. *Least Recently Used (LRU)* buffer replacement scheme is used for buffer management. The maximum size of buffer pool is set to hold 3000 data objects (1/2 of total data objects). The buffer pool is shared between the guaranteed service and the best-effort service. The ratio of buffer pool sizes between two services is adjusted dynamically. A NAND flash memory is assumed for persistent data storage. Read operations occur in the unit of a page. The size of a page is set to 512 bytes. Read time per page is set to $30\mu s$. Write time is set to $300\mu s$. Erase time is not considered because we assume the size of flash memory is big enough to ignore the effect of erases and garbage collection [6]. The flash memory is assumed to have only one bank and read/write requests are serialized to the bank. Because we assume no interleaving between banks, the maximum bandwidth of the flash memory accesses is determined only by the access latencies. The obtained maximum bandwidth is about 16MB/sec. The bandwidth of the interconnection bus between the main memory and the flash memory is assumed to be much greater than the flash memory bandwidth, thus avoiding interferences from other bus operations; this assumption is reasonable when we consider the two most common bus technologies, USB [1] and PCI [16], that have 40 MB/sec and 133.3 MB/sec maximum bandwidth respectively.

4.2 Baseline

To our best knowledge, the issues of simultaneous control of I/O and CPU workloads for deadline miss ratio management

have hardly been studied in real-time databases. Therefore, we compare our scheme (**I/O-CPU**) with the following baseline scheme which was introduced in [11].

CPU-ONLY: This scheme is not I/O-aware; the I/O deadline is not set for transactions, and I/O and CPU deadlines are not distinguished. When deadline miss ratio deviates from the desired miss ratio, only CPU workload is adjusted by changing update rates of temporal data and applying admission control; I/O workload is not dynamically controlled. This scheme is originally designed for main-memory real-time database that has no or negligible I/O workload. For comparison to our approach, the RT-EDBS was modeled by first-order SISO model; The CPU workload is the control input and the deadline miss ratio is the system output. A PI controller is used.

4.3 Results

Each simulation is run at least 5 times and their average is taken. 90% confidence intervals are drawn for each data point. For deadline miss ratios, confidence intervals are not shown because they are no more than 0.5%. For experiments, the reference miss ratio is set to 3%. Two metrics are used to compare our approach to the baseline; average miss ratio and throughput. The average miss ratio tells if the miss ratio requirement is satisfied, and the throughput tells if underutilization is occurred to achieve the miss ratio requirement. The throughput is defined as the percentage of timely transactions over total number of submitted transactions.

4.3.1 Experiment 1: Varying Loads

Computational systems usually shows different behavior for different workloads, especially when overloaded. In this experiment, workloads are varied by applying increasing number of user transactions. Overload can result from either I/O or CPU, or both. We apply three different set of workloads by applying three different set of user transactions.

Balanced I/O and CPU: NUM_{data} , the number of accesses data per transaction, follows $Normal(100, 20)$. The user transactions incurs almost same amount of I/O and CPU workload. In this setting, the I/O workload varies from 50% to 190% by applying more user transactions. The CPU workload varies accordingly.

CPU intensive: NUM_{data} follows $Normal(50, 10)$. Each user transaction incurs about 1.5 times more CPU load than I/O load. In this setting, the CPU workload varies from 50% to 190%. The I/O workload varies accordingly.

I/O intensive: NUM_{data} follows $Normal(150, 30)$. Each user transaction incurs about 100% more I/O load than CPU load. In this setting, the I/O workload varies from 50% to 190%. The CPU workload varies accordingly.

Note that the *workload* in the three settings indicates the amount of workload applied to the simulated RTEDBS when all transactions are admitted and no workload control is applied.

Furthermore, because the measured I/O workload changes with the size of buffer space and its subsequent buffer hit ratio, the same set of transactions can incur different I/O workloads in different RTEDB settings. In our experiments, sets of transactions are prepared to incur $x\%$ workload in a RTEDBS setting, whose buffer size is 1000 data objects. The same sets of transactions are applied to each experiment to generate $x\%$ workload. The actual measured workload in each experiment may be different due to different settings, e.g, different buffer sizes, and applying admission control and controllers

The result is shown in Figure 6-8. The results show that both I/O-CPU and CPU-ONLY effectively achieve the desired miss ratio, that is 3%, in all three different workloads. However, throughputs of two schemes to achieve the desired miss ratio are quite different. In all three workloads, I/O-CPU shows much higher throughputs. For instance, when 190% I/O workload (or CPU workload in CPU-intensive) is applied, I/O-CPU achieves 17%-28% higher throughput than CPU-ONLY. The throughput gap between I/O-CPU and CPU-ONLY increases as the workload increases. I/O-CPU shows especially better throughput than CPU-ONLY when the workload is I/O intensive as in Figure 8. This is because I/O-CPU can effectively use more buffer space as shown in Figure 8-(b). For CPU-ONLY, the buffer size is fixed at 1000 data objects. As the I/O workload increases, it incurs more buffer cache misses. For I/O-CPU, the I/O workload is effectively reduced by utilizing unused buffer space of best-effort service class, thus, increasing buffer hit ratio.

Interestingly enough, I/O-CPU consumes less buffer than CPU-ONLY both in balanced workload and CPU-intensive workload even though it achieves higher throughput than CPU-ONLY as in Figure 6 and 7. For instance, in Figure 7, I/O-CPU achieves 20% more throughput when 190% CPU workload is applied even though it consumes almost zero buffers. At first, this seems counter-intuitive because providing more buffer cache should reduce I/O load, thus reducing deadline misses due to I/O overload. This result can be attributed to the close interaction between I/O workload and CPU workload. When I/O is not the bottleneck causing transaction deadline misses, having larger buffer cache does not improve the overall miss ratio. In fact, a larger buffer cache can deteriorate the overall miss ratio. As shown in Figure 4-(b), I/O workload does not only control the I/O miss ratio but it also controls CPU miss ratio; increasing I/O workload decreases CPU workload by aborting transactions, and in contrary, decreasing I/O workload increases CPU workload, thus incurring more CPU deadline misses if CPU is already overloaded. In I/O-CPU, the transactions which are less likely to meet the deadline are selectively aborted by having small buffer cache. In other words, when CPU is overloaded while I/O is not, CPU workload can be effectively controlled by adjusting buffer cache size. Furthermore, I/O workload control via buffer cache adjustment is fine-grained, thus it prevents overshooting. In contrast, CPU-ONLY achieves desired miss ratio only through admission control when adjusting update intervals of temporal data is not available. With admission control, the amount of workload adjustment is coarse-grained. Therefore, overshoot and subsequent underutilization can happen. Actually, CPU-ONLY suf-

fers from underutilization in Figure 6-(a); The miss ratio is far lower than the desired value.

These results demonstrate that controlling I/O workload via buffer management not only fosters effective use of buffer cache resource but also enables fine-grained CPU workload control. Overall, our approach guarantees the desired QoS with much higher throughput; this implies resources are more effectively used in our approach.

4.3.2 Experiment 2: Varying Data Access Patterns

I/O workload is highly affected by data access patterns. By default, we assumed a uniform access pattern. However, the data access patterns can be different from a uniform access pattern. Moreover, the data access patterns can change at runtime. Therefore, the deadline miss ratio management scheme should be robust enough to cope with different data access patterns. In this section, the effect of data contention is tested using $x-y$ access scheme as described in [11]. In the $x-y$ access scheme, $x\%$ of data accesses are directed to $y\%$ of the data in the database. For instance, with 90-10 access pattern, 90% of data accesses are directed to 10% of data in the database, thus, incurring data contention on 10% of entire data. 50-50 access pattern is essentially equal to uniform access pattern.

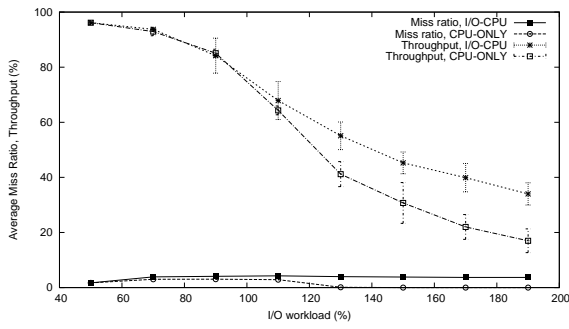
We test the robustness of our approach by applying three different $x-y$ access patterns; 90-10, 70-30, and 50-50 data access patterns. Because of space limitation, we only show the result when balanced I/O and CPU workload is applied in Figure 9. As shown in Figure 9-(a), our deadline miss ratio scheme achieves the desired miss ratio in all three different access patterns. Furthermore, the throughput of the three access patterns are not significantly different; 90-10 achieves no more than 5% higher throughput than 50-50. However, the buffer size to achieve the same performance with different access patterns are quite different as shown in Figure 9-(b). As the degree of data contention increases, the smaller size of buffer is enough to achieve the same degree of miss ratio and throughput. This suggests that miss ratio management should have flexible buffer management scheme to dynamically adjust to different data access patterns. Our results demonstrate that the proposed miss ratio management scheme is robust enough to cope with different data access patterns.

4.3.3 Experiment 3: Transient Performance

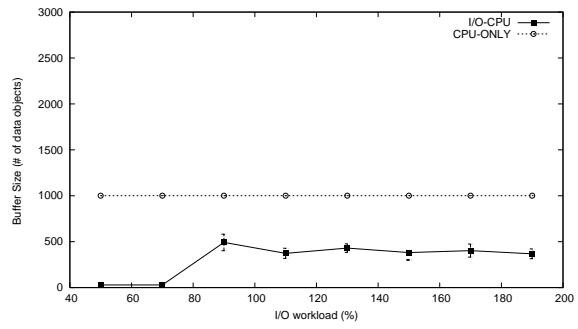
Average performance is not enough to show the performance of dynamic systems like RTEDBSs. Transient performance including settling times and overshoots should be small enough to satisfy the requirements of applications. Three different workloads are applied as in experiment-1 to observe the transient behavior of our scheme. Due to the space limitation, we only show the result when the balanced workload was applied in Figure 10. The results are similar when other workloads are applied.

Originally the system is set to have 30% I/O and 70% CPU loads. At 200 seconds, user transactions surge to increase the I/O workload by 190%. CPU load increases accordingly.

we can see that the CPU workload increases instantly to 330% after I/O workload is reduced by the controller at 230

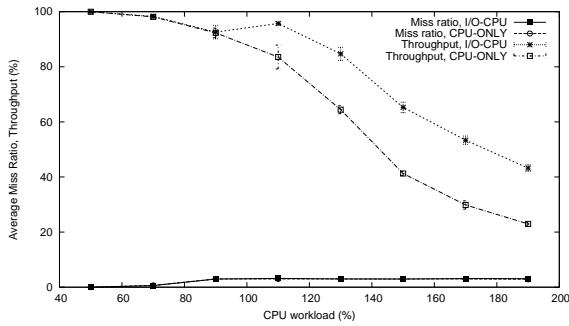


(a) Miss ratio and Throughput

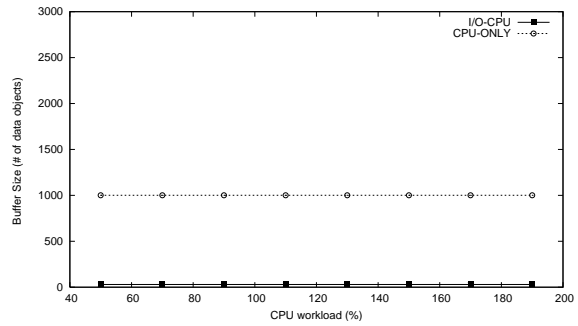


(b) Buffer size

Figure 6. Average performance when varying balanced I/O and CPU workload.

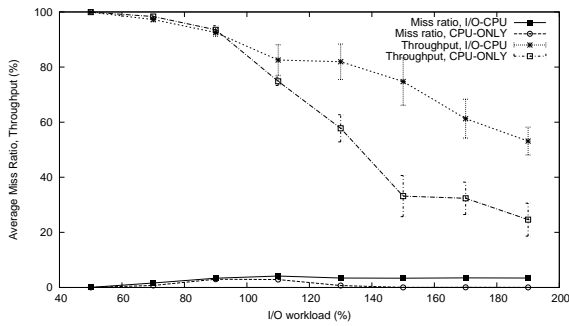


(a) Miss ratio and Throughput

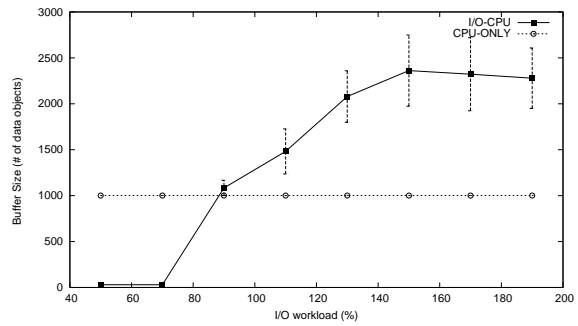


(b) Buffer size

Figure 7. Average performance when varying CPU intensive workload.

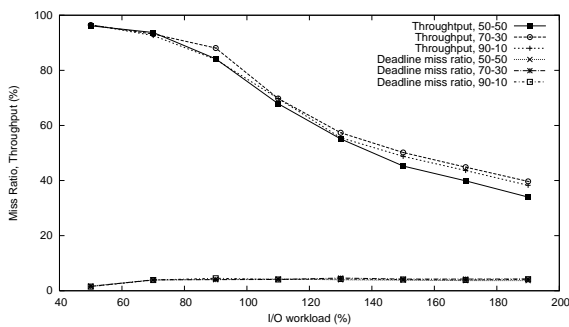


(a) Miss ratio and Throughput

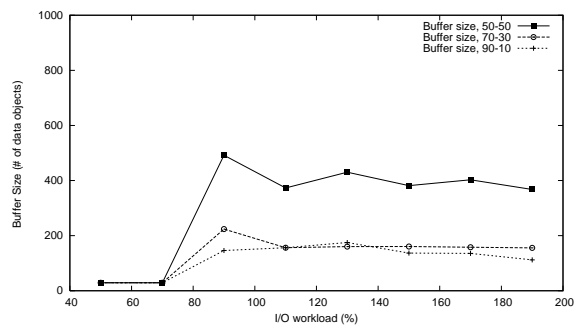


(b) Buffer size

Figure 8. Average performance when varying I/O intensive workload.



(a) Miss ratio and Throughput



(b) Buffer size

Figure 9. $x - y$ data access patterns with varying workload.

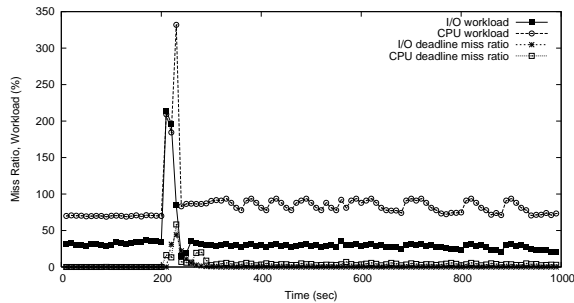


Figure 10. Sudden surge of I/O and CPU workload.

seconds. This is because reducing I/O workload by increasing the buffer cache size allows more transactions to be issued for the CPU phase without being aborted at the I/O phase, thus, causing a sudden increase in the CPU workload. The miss ratio settles down within 70 seconds. 70 seconds settling time may not be satisfactory for some real-time systems whose workloads are highly bursty. However, it satisfies the requirements of wide range of real-time applications.

5 Related Work

In the past two decades, the research in RTDBS has received lots of attention [15][10]. Most of them assume main memory databases [4] due to the inherent unpredictability of disk I/O. However, several I/O-aware approaches have been proposed, especially in terms of deadline-driven disk scheduling [3][7]. Unlike these approaches, our approach assumes flash memory as a non-volatile storage due to predictable performance.

Feedback control has been applied to QoS management in real-time systems due to its robustness against unpredictable operating environments. Lu et al. [13] proposed a feedback control real-time scheduling framework where they presented algorithms for managing miss ratio and utilization. Kang et al. [11] and Amirijoo et al. [5] both proposed feedback control-based QoS management architectures for main memory RTDBS to support the desired QoS. All these works consider only CPU resource as the source of deadline misses and I/O is not considered. Consequently, their approaches are based on SISO models. Unlike these approaches, we consider both I/O and CPU resources for deadline miss ratio management. Furthermore, due to the close interaction between I/O and CPU load, we use a MIMO technique.

6 Conclusions and Future Work

Despite the abundance of flash memory as a non-volatile secondary storage in modern real-time embedded systems, the problem of managing data in flash memory for real-time applications has not been well addressed. To address this problem, in this paper we presented an I/O-aware deadline miss ratio management scheme in RTEDBS whose secondary storage is a flash memory.

We showed that I/O and CPU workloads are closely related and a MIMO technique is required to capture the interaction between them. Furthermore, a MIMO feedback control loop was

designed to control I/O and CPU workload simultaneously. Our approach gives robust and controlled behavior in terms of guaranteeing the desired miss ratio and achieving high throughput in diverse workloads, access patterns, and even in the presence of transient overloads. The proposed algorithm outperforms the baseline algorithm where only CPU overload is considered. As one of the first work on I/O-aware deadline miss ratio management, the significance of our work will increase as flash memory increasingly replaces disks in real-time embedded systems.

We will extend this work in several ways. One direction is to investigate the impact of applying different I/O scheduling algorithm such as EDF; in this paper, we used FIFO policy for its simplicity. Secondly, we plan to implement our scheme in real embedded platforms.

References

- [1] <http://www.everythingusb.com/usb2/faq.htm>.
- [2] <http://www.samsung.com/products/semiconductor/>.
- [3] R. Abbott and H. Garcia-Molina. Scheduling real-time transactions with disk resident data. In *VLDB '89: Proceedings of the 15th international conference on Very large data bases*, pages 385–395, San Francisco, CA, USA, 1989. Morgan Kaufmann Publishers Inc.
- [4] B. Adelberg. *STRIP: A Soft Real-Time Main Memory Database for Open Systems*. PhD thesis, Stanford University, 1997.
- [5] M. Amirijoo, J. Hansson, and S. H. Son. Specification and management of QoS in real-time databases supporting imprecise computations. *IEEE Transactions on Computers*, 55(3):304–319, March 2006.
- [6] L.-P. Chang, T.-W. Kuo, and S.-W. Lo. Real-time garbage collection for flash-memory storage systems of real-time embedded systems. *Trans. on Embedded Computing Sys.*, 3(4):837–863, 2004.
- [7] S. Chen, J. A. Stankovic, J. F. Kurose, and D. Towsley. Performance evaluation of two new disk scheduling algorithms for real-time systems. *The Journal of Real-Time Systems*, 3(3):307–336, Sept. 1991.
- [8] Y. Diao, N. Gandhi, and J. Hellerstein. Using MIMO feedback control to enforce policies for interrelated metrics with application to the Apache web server. In *Network Operations and Management Symposium*, pages 219–234, 2002.
- [9] J. L. Hellerstein, Y. Diao, S. Parekh, and D. M. Tilbury. *Feedback Control of Computing Systems*. Wiley IEEE press, 2004.
- [10] T.-W. K. E. Kam-yiu Lam. *Real-Time Database Systems: Architecture and Techniques*. Kluwer Academic Publishers, 2001.
- [11] K.-D. Kang, S. H. Son, and J. A. Stankovic. Managing deadline miss ratio and sensor data freshness in real-time databases. *IEEE Transactions on Knowledge and Data Engineering*, 16(10):1200–1216, October 2004.
- [12] L. Ljung. *Systems Identification: Theory for the User 2nd edition*. Prentice Hall PTR, 1999.
- [13] C. Lu, J. A. Stankovic, S. H. Son, and G. Tao. Feedback control real-time scheduling: Framework, modeling, and algorithms. *Real-Time Syst.*, 23(1-2):85–126, 2002.
- [14] C. Park, J. Seo, D. Seo, S. Kim, and B. Kim. Cost-efficient memory architecture design of nand flash memory embedded systems. In *21st International Conference on Computer Design*, 2003.
- [15] K. Ramamritham. Real-time databases. *Distrib. Parallel Databases*, 1(2):199–226, 1993.
- [16] T. Shanley and D. Anderson. *PCI System Architecture(4th Edition)*. Addison-Wesley Professional, 1999.
- [17] J. A. Stankovic, S. H. Son, and J. Liebeherr. BeeHive: Global multimedia database support for dependable, real-time applications. *Lecture Notes in Computer Science*, 1553:51–72, 1998.