

# Load Shedding in Stream Databases: A Control-Based Approach

Yi-Cheng Tu<sup>†</sup> Song Liu<sup>‡</sup> Sunil Prabhakar<sup>†</sup> Bin Yao<sup>‡</sup>

<sup>†</sup>Department of Computer Sciences  
Purdue University  
250 N. University St.  
West Lafayette, Indiana, USA  
{tuyc, sunil}@cs.purdue.edu

<sup>‡</sup>School of Mechanical Engineering  
Purdue University  
140 S. Intramural Drive  
West Lafayette, Indiana, USA  
{liu1, byao}@purdue.edu

## ABSTRACT

Query processing in Data Stream Management Systems (DSMSs) has to meet various Quality-of-Service (QoS) requirements. In many data stream applications, processing *delay* is the most critical quality requirement since the value of query results decreases dramatically over time. The ability to remain within a desired level of delay is significantly hampered under situations of overloading, which are common in data stream systems. When overloaded, DSMSs employ load shedding in order to meet quality requirements and keep pace with the high rate of data arrivals. Data stream applications are extremely dynamic due to bursty data arrivals and time-varying data processing costs. Current approaches ignore system status information in decision-making and consequently are unable to achieve desired control of quality under dynamic load. In this paper, we present a quality management framework that leverages well studied feedback control techniques. We discuss the design and implementation of such a framework in a real DSMS - the Borealis stream manager. We introduce the concept of virtual queue length by which the delays of current incoming data can be effectively controlled. Our data management framework is built on the advantages of system identification and rigorous controller analysis. Experimental results show that our solution achieves significantly fewer QoS (delay) violations with the same or lower level of data loss, as compared to current strategies utilized in DSMSs. It is also robust and bears negligible computational overhead.

## 1. INTRODUCTION

Applications related to processing of data streams have attracted a great deal of attention from the database community. With great social/economical interests, these applications flourish in a number of fields such as environment monitoring, system diagnosis, financial analysis, and mobile services. Unlike traditional data that are mostly static, stream data are produced continuously (e.g. from a sensor network) and are generally too large to be kept in storage

after being processed. Furthermore, most queries against stream data are persistent queries that continuously output results as they are produced. Thus, data stream processing brings great challenges to DBMS design: it imposes a *data-active, query-passive* DBMS model instead of the *data-passive, query-active* model for traditional DBMSs [2]. In recent years, a number of Data Stream Management Systems (DSMSs) have been developed [2, 1, 14, 9].

Query processing in DSMSs has to meet various quality<sup>1</sup> requirements [2]. Similar to those in other real-time applications [22], quality in DSMSs describes the timeliness, reliability, and precision in data processing and service delivery. Important quality parameters in DSMSs include: processing delay, data loss ratio, sampling rate, etc. A salient feature of data stream management is the real-time constraints associated with query processing. In many applications of DSMS, query results are required to be delivered before either a firm (e.g. tracking of stock prices) or soft (e.g. network monitoring for intrusion detection) deadline. Therefore, processing delay is the most critical quality parameter in these applications. On the other hand, users may accept query processing at different levels of accuracy as a result of lost or incomplete data [23, 26]. This provides us with optimization opportunities to trade those quality parameters (e.g., loss ratio, sampling rate) that are less important for shorter delays in case of congestion.

It is difficult to provide delay guarantees in a DSMS due to physical resource limitations and the unpredictable pattern of resource usage by streams. In practice, a DSMS could easily accommodate hundreds or even thousands of streams. Delay requirements may be violated even with careful query optimization and admission control, which are the first line of defense against overloading and generally based on static estimations of each stream's resource consumption. The runtime fluctuations of application resource usage (e.g. bursty arrivals) may cause temporary congestion that interferes with real-time data processing. Under this situation, we need to dynamically adjust application behavior by reducing its non-critical quality parameters. For example, we can increase data loss rate by *load shedding* [26] or reduce the window size for windowed operations [4]. We call such adjustment of application parameters *adaptation*.

<sup>1</sup>In this paper, the words 'QoS' and 'quality' are used interchangeably.

Streaming data are intrinsically dynamic in terms of their bursty arrival patterns and ever-changing tuple processing costs [32, 24]. Thus, it is important to have an adaptation mechanism that promptly adjusts application behavior in response to changes of system and input status [4]. The adaptation architecture should promptly detect the change of quality by continuously monitoring the system and determine whether adaptation should be performed.

While maintaining processing delays under an appropriate level, degradation of other quality should also be controlled. For example, we can always achieve low delays by constantly discarding most of the load. However, query accuracy decreases unnecessarily due to excessive load shedding. It would be desirable to achieve low delays while minimizing data loss. Attempting to solve this problem, current DSMSs employ simple and intuitive strategies to make important adaptation decisions such as the time and magnitude of load shedding. For example, the following load shedding algorithm is used (explicitly) in Aurora [26] and (implicitly) in STREAM [6].

---

```

1  for every  $T$  time units
2      if measured load  $L$  is greater than CPU capacity  $L_0$ 
3          do shedding load with amount  $L - L_0$ 
4      else allow  $L_0 - L$  more load to be admitted

```

---

**Figure 1: Load shedding algorithm in Aurora**

The idea behind this algorithm is: QoS degrades when the load injected into the system is higher than its processing capacity. In dealing with overloading, we only need to make the input load smaller than capacity  $L_0$ . However, in a dynamic environment where the input rate keeps changing, this approach may either make the DSMS unstable (i.e., QoS deviates unboundedly from the desirable value) or overreact by discarding too much load. In Section 4.3.2, we elaborate on this issue.

To remedy the above problems in a systematic way, however, is not trivial. Firstly, we need to understand the nature of the DSMS’s response to changes of inputs. Specifically, a quantitative model that describes how adaptation of stream behavior affects quality (delay) is needed. Secondly, our adaptation algorithm should be robust, meaning that its performance should not be affected by patterns of load fluctuations and cost variations. Another challenge is the design of the monitoring process: it should be light-weight and still able to effectively capture changes of status.

In this paper, we present our approach to address the above challenges. Our solution takes advantage of proven techniques from the field of control theory. Feedback control is extensively utilized in the fields of mechanical, chemical engineering, and aeronautics to deal with systems that bear dynamics that are hard to model [13]. In this work, we view quality-driven load shedding in DSMS as a feedback control problem and solve it with a controller designed from a dynamic DSMS model we develop. Specifically, this paper makes the following contributions:

1. We develop a dynamic model to describe the relationship between average tuple delays and input rate of a DSMS. From this model, we propose the idea of controlling the somewhat unmeasurable delay signal by manipulating the number of outstanding data items;

2. We design a controller to make load shedding decisions via rigorous analysis of the system model. By exploiting results from control theory, our design achieves guaranteed system performance;
3. We implement and evaluate our load shedding framework on a real DSMS. By working on a real system, we achieve better understanding of the DSMS model and obtain more convincing results supporting the validity of our approach; and
4. We identify several problems that are unique in the control of DSMS load shedding and propose system-specific strategies to solve to these problems.

The rest of this paper is organized as follows: we compare our work with related research efforts in Section 2. Section 3 describes the basic DSMS model and problem formulation. Details of our feedback control framework are presented in Section 4. We show experimental results in 5 and conclude the paper in Section 6.

## 2. COMPARISON TO RELATED WORK

Current efforts on DSMSs have addressed system architecture [7, 14], query processing [12, 15], query optimization [28], and stream monitoring [33]. Relatively less attention has been paid to the development of a unified framework to support QoS. An important issue related to QoS control in DSMSs is the development of scheduling policies for query operators. Two relevant efforts present scheduling algorithms that minimize tuple delays [8] and runtime memory consumption [5].

Research on QoS control was first motivated by the real-time requirements of multimedia applications. Most of these efforts emphasize system and network level resource management, which is provided as a service of the operating system [31] or a middleware [21]. The system maps QoS requirements of applications to resource use (system QoS) and QoS control is accomplished by regulating resource allocation to individual applications.

Load shedding has been extensively utilized to deal with overloading in DSMSs [26, 6, 25]. Ref [6] discusses load shedding strategies that minimize the loss of accuracy of aggregation queries. To increase accuracy of arbitrary queries, a *data triage* approach that exploits synopses of the discarded data is proposed in [25]. In the LoadStar system [10, 11], statistical models are utilized to maximize the quality of stream mining results when load shedding has to be performed. Earlier work on QoS-driven load shedding in the context of the Aurora [26] DSMS (now evolving to the Borealis project [3]) is closely related to our study in this paper. In [26], three critical questions about load shedding are raised: *when*, *where*, and *how much* to shed. To answer these questions, Aurora checks system load periodically and triggers shedding when excessive load is detected. A precomputed Load Shedding Roadmap (LSRM) that holds possible shedding plans is used to determine where to shed load. Given the amount of total load to shed, the LSRM finds the best plan to accomplish this such that system utility loss is minimized. The utility is calculated from data loss ratio only.

The Aurora/Borealis work focuses more on the question ‘where to shed load’ (i.e., construction of LSRM) than the questions of ‘when’ and ‘how much’. As shown in Fig. 1, it

uses a heuristic to determine the amount of load shedding and handles processing delays implicitly. The system does not provide information about how the monitoring period  $T$  is set. In this paper, we concentrate on the control of delay QoS under heavy fluctuations/bursts and time-varying processing costs of data inputs, which are common in data stream applications. For this purpose, we need to find a solution that is different from the Aurora load shedder shown in Fig. 1. In other words, our work aims to provide better answers to the questions of *when* and *how much* to shed load under a highly dynamic environment. Our solution can also be used to guide quality adaptation mechanisms other than load shedding. In addition to *statistical shedding* that discards tuples randomly, [26] also explores *semantic shedding* that chooses victim tuples based on a cost/utility analysis of query operators.

Control-theoretic approaches have been used to solve various problems in the areas of networking [18], real-time systems [20], multimedia [19] and event notification [30]. Our work differs significantly from these efforts: first, we address different problems with a different system to be controlled. For example, [20] focuses on deadline misses in a real-time resource scheduler and [19] discusses precision in a video tracking application. While the event notification system in [30] uses DBMS as building blocks, its control target is to balance load across multiple CPUs. Control theory is basically a collection of many mathematical tools for analyzing system features and designing controllers towards guaranteed performance. Therefore, application of control theory to different problems and systems are not straightforward as the choice of appropriate control techniques are essential. Controller design and performance depend heavily on the dynamic model of the system of interest. Derivation of such models is generally non-trivial and involves various techniques for different systems. In a word, the results in above studies give little insights on how our problem can be solved as distinct systems and control targets are involved. Second, we raise several DSMS-specific issues in this paper (Section 4.5) and provide solutions to these problems. From a control theory viewpoint, these issues bring new challenges that are never before addressed in traditional control applications.

We have presented the idea of control-based load shedding in a short paper [27] where we found that even a crude controller outperforms the load shedding method based on static estimations of system status. However, we only validate our idea with a simulator and a simple controller in [27] therefore the real challenges of QoS control in real-world systems (i.e., major contributions of this paper shown in Section 1) are not addressed.

### 3. DSMS MODEL, PROBLEM DESCRIPTION, AND NOTATIONS

In this paper, we study load shedding under a push-based query processing model, which is a generalization of those of the STREAM [4] and Aurora [2] stream managers. In this model, each query plan consists of a number of *operators* connected to form a branched (e.g., I and III in Fig 2A) or unbranched (e.g., II in Fig 2A) *execution path*. Multiple queries form a network of operators so that they can share computations. Multi-stream joins are performed over a *sliding window* whose size is specified by the application either in number of tuples or time. Data from a stream can en-

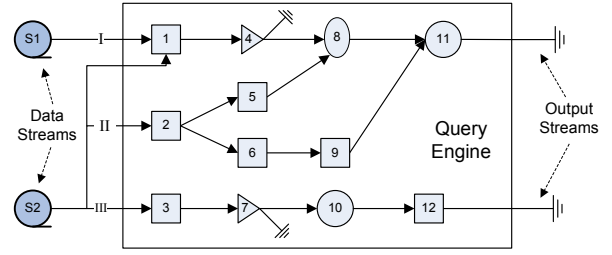


Figure 2: A general system model of DSMS.

ter any number of entry points in the query network. Each operator has its own queue to hold unprocessed inputs and a *scheduler* determines the processing order of operators at runtime.

With respect to a data tuple, *processing delay* is defined as the time elapsed since it arrives at the network buffer of the query engine till it leaves the query network.<sup>2</sup> For example, data from stream source  $S1$  in Fig. 2A departs either after being discarded by the filter operator 4 or entering an output stream after operator 11. For data that could enter multiple execution paths, we can choose the longest path to record its departure time (e.g., 2-6-9-11 or 3-7-10-12 for  $S2$  data in Fig. 2A). Processing delay consists of CPU time spent to execute the operators and time spent in queues.<sup>3</sup> We target a system where data tuples arrive in a dynamic pattern such that future data rates are unpredictable. Furthermore, the expectation of per-tuple CPU cost changes over time. Variations in CPU cost arise from changes in factors such as query network structure (due to addition/deletion of queries), and operator selectivity [26]. In this paper, we assume such variations happen less frequently than the fluctuations of data arrival rates. We believe this is a reasonable assumption as none of the above factors would change abruptly.

To reverse the increase of processing delays due to overloading, the DSMS can perform any of the following adaptations: (i) *load shedding*: discard unprocessed data tuples by placing filters either in the data source or at the entry points to the query network; (ii) *reducing sampling rate*: save costs by changing the frequency of data tuple generation at stream sources; and (iii) *modifying operator features* such as window size of join operators. Although our solution should also work for (ii) and (iii), we focus on load shedding in this paper.

Our quality-driven load shedding framework allows the system administrator to specify a target delay time  $y_d$ . The goal is to maintain the average processing delay of data tuples that arrive within a small time window  $T$  (we will discuss more about the choice of  $T$  later) to be under  $y_d$ . We accomplish this by dynamically dropping load from the system in case of overloading. The problem is how to derive

<sup>2</sup>Here we ignore network delays. This can be justified by the use of networks where transmission delays are either effectively controlled or significantly smaller than our control target.

<sup>3</sup>This implies that CPU power is the bottleneck, which is a reasonable assumption [26]. We understand that limited memory could result in blocking of data processing. However, this should have little effect on our problem because our goal is to control overloading so that the system runs in a zone without such nonlinearities.

the right time and amount of load shedding such that data loss is as low as possible. The selection of shedding locations is not a focal point of this study. However, our framework is designed to work with current strategies that construct shedding plans such as the current load shedder in Borealis. We consider the following metrics in evaluating the adaptation strategy:

- *Delay Violations*, which is the primary goal of the control. Specifically, we record both the *accumulated delay violations* (i.e.,  $\sum y - y_d$  for all data tuples whose processing delay  $y > y_d$ ), and *total delayed tuples*, which is the total number of tuples whose delays are longer than  $y_d$ ;
- *Maximal Overshoot*: the longest delay violation (i.e.,  $y - y_d$ ) recorded. This metric captures transient state performance; and
- *Data Loss Ratio*: the percentage of data tuples discarded. This can be viewed as the *cost* of performing load adaptation.

Symbols used throughout this paper are listed in Table 1.

**Table 1: Notations and symbols.**

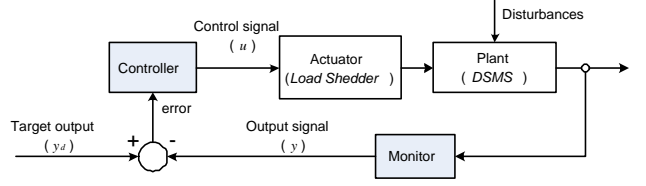
Symbol	Definition	$z$ -domain
$k$	discrete time index	-
$T$	control period	-
$y_d$	target value for delays	-
$H$	CPU power for query processing	-
$y$	processing delay	$Y(z)$
$f_{in}$	data input rate	$F_{in}(z)$
$f_{out}$	data output rate	$F_{out}(z)$
$u$	controller output	$U(z)$
$v$	desired data input rate	-
$c$	per-tuple processing cost	-
$q$	number of outstanding tuples	$Q(z)$
$C(z)$	controller transfer function	$C(z)$
$G(z)$	system (DSMS) transfer function	$G(z)$

## 4. FEEDBACK CONTROL-BASED LOAD SHEDDING FRAMEWORK

In this section we present our quality adaptation framework with the objective of maintaining processing delays.

### 4.1 Overview

The term *control* generally refers to the operations to manipulate particular feature(s) (i.e., output signal) of a *process* by adjusting inputs into the process. The main components of a feedback control system form a *feedback control loop*, as shown in Fig 3. The controlling operations of the feedback control loop are performed as follows: a *monitor* measures the output signal of the *plant*, which is the process to be controlled. The measurements are sent to a *controller*. The controller compares the value of the output signal with a *target* value and maps the *control error*, i.e., difference between the output signal and the target, to a *control signal*. An *actuator* adjusts the behavior of the plant according to the control signal. The goal of the control operations is to



**Figure 3: The feedback control loop.**

overcome the effects of system and environmental uncertainties named *disturbances*. Readers interested in more details on control theory can refer to [13].

The above general model can be translated into a concrete model that serves as a blueprint for our load shedding framework. We still use Fig. 3 to illustrate this. Note that the shaded boxes represent new components that are not found in any existing DSMSs. The plant to be controlled is the query engine of the DSMS and the actuator is the existing load shedding algorithm that adjusts load injected into the plant. In addition, we have a monitor that measures the output signal and a controller to generate the control signal. The unpredictable arrival patterns and processing costs are all treated as disturbances. In this loop, the output signal is the processing delay of tuples, denoted as  $y$  and the control signal (i.e., controller output and system input) is the desirable incoming data rate  $u$ .

We can easily see that the most critical part of the control loop is the controller, which determines the quantity of the input signal (to DSMS) based on the control error. The beauty of control theory is that it provides a series of mathematical tools to design and tune the controller in order to obtain guaranteed performance under disturbances. In the following, we discuss the design of our feedback control, which consists of two phases: system modeling (Section 4.2) and controller design (Section 4.4). We use the open-source Borealis data stream manager [3] as our experimental system. The query engine of Borealis is derived from the Aurora system [2].

### 4.2 System Modeling

An accurate mathematical model of the plant is of great importance to control system design. In this study, the model we are interested in is one that describes the relationship between the delay time  $y$  and the incoming data flow rate  $f_{in}$ . Due to the complexity of the controlled system, we may not be able to derive a model solely based on rigorous analysis. In this case, we can use *system identification* techniques to study system dynamics experimentally.

First of all, the expectation of per-tuple processing cost  $c$  can be precisely estimated in the current Borealis system. Readers can refer to Section 4.2 of [26] for details. For the purpose of system modeling, we treat  $c$  as a constant and relax this assumption in Section 4.5.

The current version of Borealis uses a *round robin* policy to schedule operators and place intermediate results in waiting queues of individual operators. These queues extract input in a first-in-first-out (FIFO) manner therefore we see no priorities assigned to tuples as a result. Let us first consider an ideal situation: all tuples in the network share the same query paths and the system has the same inflow and outflow rates. If there are  $q$  outstanding data

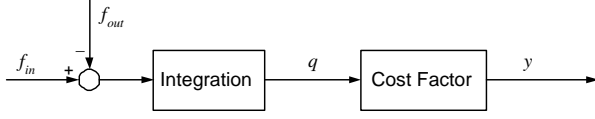


Figure 4: Database model structure.

tuples in the query network when a tuple  $A$  enters, the total processing delay of  $A$  is

$$y = (q + 1)c \quad (1)$$

The reason for this is: when  $A$  sits in the queue of any operator, it will not be processed until all of the  $q$  tuples are cleared from that queue. If the execution path of  $A$  consists of  $n$  operators, a total number of  $nq + n$  operators would have been executed by the system when  $A$  is finished. The cost of the  $n$  operators in a path is  $c$  therefore the total cost becomes  $(q + 1)c$ . Among the total time of  $qc + c$ ,  $qc$  is time spent in waiting queues and  $c$  is the processing time of  $A$  itself. In other words, it is equivalent to processing tuples as a whole (rather than by operators) in the order they arrive. The outstanding tuples can be regarded as entries in a virtual FIFO queue with length  $q$ .

In practice, we cannot use Eq.(1) to model delay time of individual tuples because the real execution paths for different tuples are different. For example, if a tuple is discarded by a *selection* operator in the early part of its possible path, it has a shorter delay as compared to one that passes the selection box (and goes further in the query network). Fortunately, instead of delay time  $y$  of single tuples, we are interested in the average delay time of a series of tuples arriving in a period of time.<sup>4</sup> Let us denote the length of this period, which is called *control period* or *sampling period*, as  $T$  and the average delay of tuples within the  $k$ th period as  $y(k)$ . We propose the following generalization of Eq.(1):

$$y(k) = \frac{c}{H} [q(k - 1) + 1] \quad (2)$$

where  $H$  is a constant named *headroom factor*. i.e., the fraction of processing power used for query processing. We always have  $H < 1$  as resources must be consumed for running the operating system and other maintenance tasks. The intuition behind Eq.(2) is: we can study data tuples with execution paths of the same length in a group. The same reasoning to generate Eq.(1) holds true for each such group: a tuple  $A$  will not leave the network until all other tuples in its group (that entered the network before  $A$ ) are processed. Taking a weighted average of all such groups, each of which can be described by Eq.(1), we get a form that is close to Eq.(2). As the above is an intuitive result, we need to verify it by experiments.

Eq.(2) leads to a system model for Borealis as shown in Fig. 4. The incoming data flow  $f_{in}$  less the data processing rate  $f_{out}$  is accumulated in the virtual queue. Therefore the queue length at the end of period  $k$ ,  $q(k)$ , is equal to the integration of  $f_{in} - f_{out}$  at all times up to the  $k$ -th period.

<sup>4</sup>To guarantee delays for individual tuples, real-time schedulers [17] are generally deployed. Interestingly, in our system, if we can guarantee average delays, those for individual tuples can also be well maintained as the round robin policy is a 'fair' policy.

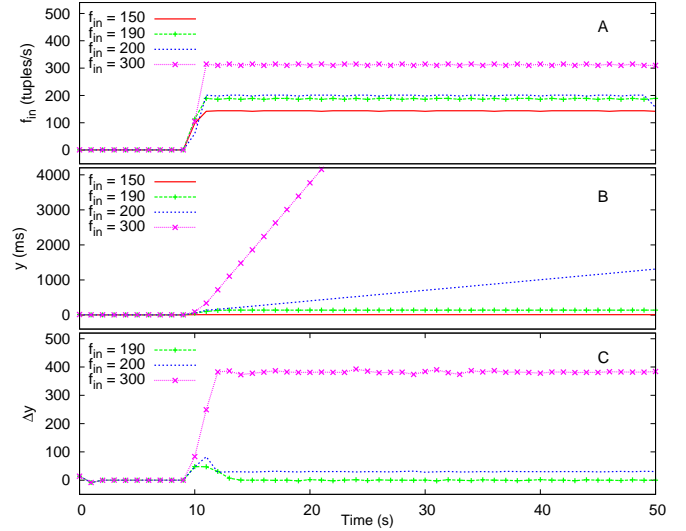


Figure 5: System responses to step inputs.

Eq.(2) becomes:

$$y(k) = \frac{c \cdot T}{H} \sum_{i < k} [f_{in}(i) - f_{out}(i)] + \frac{c}{H}.$$

*Model verification.* The verification of the dynamic model is done experimentally in accordance of system identification techniques. We feed the Borealis system with synthetic data streams having various arrival patterns and record responses in terms of delay time  $y$ . To set the cost factor  $c$  to a constant, we construct a Borealis query network with a number of (14 in this case, details omitted due to space limitations) operators, each of which has a fixed CPU cost. Then we generate stream data whose values follow uniform distributions to fix the selectivity of all filtering operators. By doing these, the average CPU cost of the query network becomes stable. In Fig 5, we report system responses to a stream whose arrival rates follow a step function of time (i.e., rate starts at very low and jumps to a high value at the 10-th second, as shown in Fig. 5A). It is shown in Fig. 5B that, when  $f_{in}$  is less than 190 tuples per second, all data can be processed immediately and a constant processing delay is observed. This implies that the per-tuple CPU cost is approximately  $1000.0/190 = 5.26ms$  as 190 can be viewed as the threshold load that equals the CPU processing capacity (i.e.  $f_{in} = f_{out} = 190$  assuming  $H = 1$ ). On the other hand, when  $f_{in}$  exceeds 190/s, i.e., more data entering the system than the CPU can handle, data accumulates in the virtual queue and delay  $y$  keeps increasing. This is strong evidence of the existence of the integration part in the proposed model. Fig. 5C shows the changing rate of  $y$  (calculated by  $\Delta y = y(k) - y(k - 1)$ ). The fact that  $\Delta y$  converges quickly to a stable value means that there is either no other dynamics or unknown dynamics with insignificant effects in the proposed model.

To further verify the model and determine the model parameter, we compare the real  $y(k)$  values measured and the calculated  $y(k)$  values based on our system model (Eq.(2)). We collect  $q(k)$  values at runtime for the calculation of  $y(k)$ . The results of experiments using the same step inputs as be-



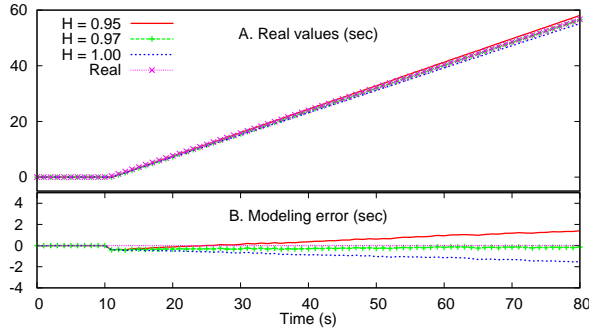


Figure 6: Model verification with step inputs.

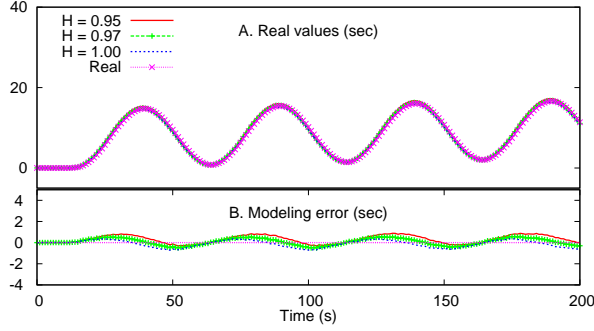


Figure 7: Model verification with sinusoidal inputs.

fore are plotted in Fig. 6. According to Fig. 6A, the values given by our model fits the real data very well for all three choices of  $H$ . However, if we magnify the difference between calculated and real values (Fig. 6B), we can see that, with a  $H = 0.97$ , modeling errors are far less than the other two values of  $H$ .

We also tested the system with sinusoidal inputs and similar results are obtained and are plotted in Fig. 7. In this set of experiments, the incoming data flow rate  $f_{in}$  changes sinusoidally within the range of  $[0, 400]$ . Small, periodical modeling errors can be seen. This means there are probably unknown dynamics that our model fail to capture. This is not surprising due to the complexity of the Borealis system. As we shall see later, feedback controllers, if properly designed, have the power to reduce the effects of modeling errors, especially those that impose small errors such as the one we observe here.

*Model transform.* For the convenience of control analysis, we transform Eq.(2) to a model in the  $z$ -domain:<sup>5</sup>:

$$Y(z) = \frac{c}{H}Q(z) = \frac{c \cdot T}{H(z-1)}[F_{in}(z) - F_{out}(z)] \quad (3)$$

where  $Y(z)$ ,  $Q(z)$ ,  $F_{in}(z)$  and  $F_{out}(z)$  are  $z$ -transforms of signals  $y(k)$ ,  $q(k)$ ,  $f_{in}(k)$  and  $f_{out}(k)$ , respectively. The transfer function of the (Borealis) system in Fig. 4 is:

$$G(z) = \frac{c \cdot T}{H(z-1)}. \quad (4)$$

<sup>5</sup>The  $z$ -transform is a mathematical tool that transforms difference equations to algebraic equations [13], similar to the Laplace transform used for differential equations.

From now on, all control-related analysis will be performed in the  $z$ -domain.

### 4.3 Why feedback control?

Before going into the design of controller, we briefly discuss the basic ideas of feedback control theory and identify some of the problems of non-feedback-control strategies.

#### 4.3.1 Open-loop vs. closed-loop

The unique feature of feedback control is that the output signal is used (as feedback) in generating the control signal. As there exists a complete loop (Fig. 9B) in the system block diagram, feedback control is also called *closed-loop* control. In contrast, strategies such as the one shown in Fig. 1 are *open-loop* control: system output or state information is not used in the controller, therefore it forms an open loop from the reference value to the system output, as shown in Fig. 9A. Here  $r$  is the reference input or desired system output,  $y$  is the actual system output,  $a$  is the system model,  $d_m$ ,  $d_i$  and  $d_o$  represent modeling error, input disturbance and output disturbance. Relating this to our problem, the fluctuations of data arrival rates are modeled as input disturbances and the variable processing costs ( $c$ ) as modeling errors.

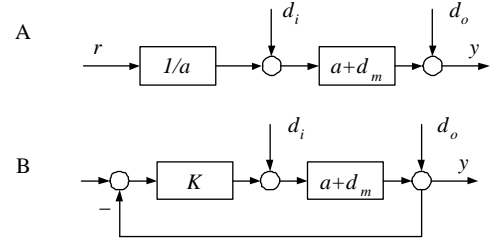


Figure 9: Block diagrams of generic open-loop (A) and closed-loop (B) control systems.

In an ideal case, when there are no model uncertainties (i.e.  $d_m = 0$ ), no input or output disturbances (i.e.  $d_i = d_o = 0$ ), the best open-loop controller would be  $1/a$  given the nominal system model  $a$ . This is because we have  $y = r \frac{1}{a} = r$ , hence the output signal is exactly the reference value. However, in the real world, there are always modeling errors and input/output disturbances, therefore the open-loop system output  $y$  is:

$$y = \left( \frac{r}{a} + d_i \right) (a + d_m) + d_o = r + r \frac{1}{a} d_m + (a + d_m) d_i + d_o \quad (5)$$

From (5), it is obvious that the open-loop system output is subject to modeling error  $d_m$ , input disturbance  $d_i$  and output disturbance  $d_o$ , and there is no way to reduce their effects. On the other hand, in a closed-loop system where the feedback controller  $K$  is also designed based on the nominal system model  $a$ , we have

$$[(r - y)K + d_i](a + d_m) + d_o = y,$$

and the system output  $y$  becomes

$$y = \frac{K(a + d_m)}{1 + K(a + d_m)} r + \frac{(a + d_m)}{1 + K(a + d_m)} d_i + \frac{1}{1 + K(a + d_m)} d_o$$

If the controller  $K$  is chosen large enough, i.e.  $K \gg 1$  and  $K(a + d_m) \gg 1$ , the closed-loop system output  $y$  is

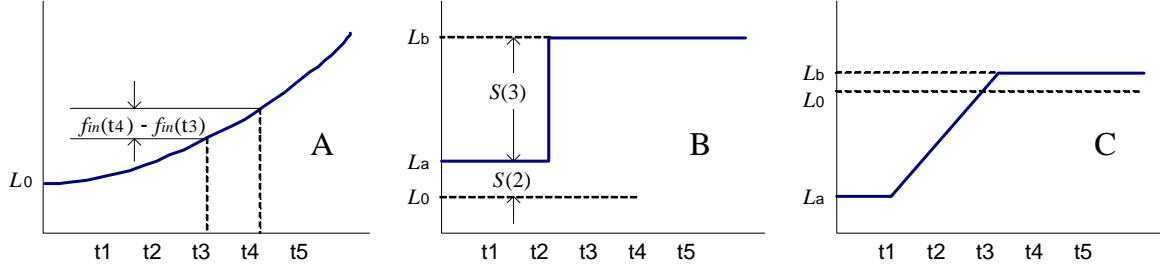


Figure 8: Different cases in which open-loop control has poor performance.

approximately:

$$y \approx r + \frac{1}{K}d_i + \frac{1}{K}d_o \quad (6)$$

It is obvious that the effects of modeling error  $d_m$ , input and output disturbances  $d_i$  and  $d_o$  can be reduced by a factor of  $1/K$ .

The above simple examples show why closed-loop control is better than open-loop control. In summary, the main advantage of closed-loop control over open-loop is the reduction of the effects of modeling error, input and output disturbances therefore it can be exploited in solving our load shedding problem.

#### 4.3.2 Problems of current load shedding solution

As mentioned earlier, the current Aurora method for dynamic load shedding is open-loop in nature: it uses a preset threshold  $L_0$  to adjust the incoming data flow. When the incoming data flow  $L$  is more than  $L_0$ ,  $L - L_0$  amount of data will be discarded. Using our notations, assuming constant processing cost  $c$ ,  $L$  can be replaced by  $f_{in}$ . Thus, the amount of data to be shed in the  $k$ -th sampling period  $S(k)$  is

$$S(k) = f_{in}(k) - L_0 \quad (7)$$

where  $L_0$  is the preset threshold generally set as the processing capacity of the CPU. As  $f_{in}(k)$  is not predictable at the beginning of period  $k$ , we have to use an estimated value such as  $f_{in}(k-1)$ . The algorithm shown in Fig. 1 would result in the following queue length

$$\begin{aligned} q(k) &= q(k-1) - L_0 + [f_{in}(k) - S(k)] \\ &= q(k-1) + f_{in}(k) - f_{in}(k-1), \end{aligned} \quad (8)$$

and average delay time

$$y(k) = q(k) \cdot c = [q(k-1) + f_{in}(k) - f_{in}(k-1)]c. \quad (9)$$

In other words, the queue length at the  $k$ th period is equal to the previous queue length  $q(k-1)$  less the processed data  $L_0$  plus the incoming data with amount  $f_{in}(k) - S(k)$ .

Based on above analysis, we shall see that the open-loop control suffers from poor performance as detailed in the following examples (see Fig. 8 for illustrations).

**Example 1.** *Instability when incoming data rate increases monotonically.* During certain period of time, the incoming data rate may keep increasing as shown in Fig. 8A. This is very typical in dynamic environments. In this case, the shed factor  $S(k)$  is not sufficiently large because it is derived from the incoming data rate  $f_{in}(k-1)$ . According to Eq.(8),

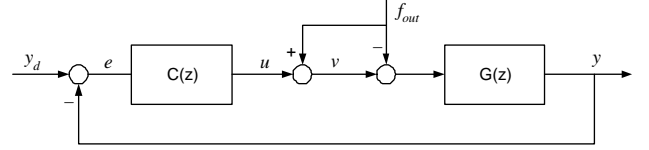


Figure 10: Control system block diagram.

the number of outstanding data tuples will keep growing because  $f_{in}(k) - f_{in}(k-1) > 0$  for all  $k$  in the period. As a result, system output  $y$  also increases unboundedly.

**Example 2.** *Convergence to wrong value in response to step changes.* As illustrated in Fig. 8B, when the incoming data rate undergoes a step change from  $L_a$  to a much larger value  $L_b$ , queue length will increase by  $L_b - L_a$ . If the incoming rate stays on  $L_b$ , no further increase of queue length will occur and system output  $y(k)$  stabilizes. However,  $y(k)$  could converge to a value that is higher than the target value  $y_d$ . And the system is unable to self-correct the deviation due to its open-loop nature (i.e., controller does not know the actual system output).

**Example 3.** *Unnecessary data loss.* When the incoming data rate changes from a stable small value  $L_a$  to a higher value  $L_b$  that is slightly greater than  $L_0$ , the algorithm will discard data with amount of  $L_b - L_0$  (Fig. 8C). However, more data should be allowed to enter the DSMS because the queue is almost empty before the change. In this case, although the delay time  $y(k)$  is smaller or better than the expected target value  $y_d$ , the extra data loss is unnecessary. Again, the reason for this is that the controller does not know the actual system output.

The above three cases do not occur just in the Aurora method. Any open-loop method where the system output does not play any role in the control could face the same or similar problems.

## 4.4 Controller design

In Section 4.3.2 we demonstrated the shortcomings of open-loop solutions. However, even with the system model, it is still not clearly how to use system output to make control decisions and problems such as the one in Example 1 may still exist. In this section, we discuss how to design controllers with guaranteed performance. We start this section by introducing our basic design of controller and continue with Section 4.5 to address some DSMS-specific challenges. The basic control scheme is illustrated in Fig. 10, where  $y_d$  is the preset reference value for delay time,  $e = y_d - y$

is the error signal, and  $u$  represents the controller output (with the same unit as inflow rate  $f_{in}$ ). The meaning of  $u$  is: the *increase* of the number of outstanding tuples (i.e., size of the virtual queue) allowed in the next control period. Therefore, we denote  $v = u + f_{out}$  as the desired data flow rate to the database as  $f_{out}$  tuples will leave the queue.  $C(z)$  is the controller transfer function.

#### 4.4.1 Design based on pole placement

For a dynamic system, continuous or discrete, one can use system poles to determine its dynamic characteristics. System poles are the roots of the denominator polynomial of the transfer function and zeros are the roots of the numerator polynomial. The location of the system poles can tell how fast the system responds to an input and how well the response would be. For example, if a discrete time system has a pole on the real axis between 0 and 1, the system response would not oscillate; if it has a pole outside of the unit circle, the system becomes unstable. The relationship between the location of the system poles and the system response can be found in any control textbook such as [13].

Pole placement design, one of the most important controller design techniques, is to add additional poles and/or zeros into the closed-loop system so that the closed-loop system may have desired performance. If a raw system  $G(z) = \frac{B(z)}{A(z)}$  has poles as the roots of  $A(z) = 0$ , the closed-loop system, after adding a feedback controller  $C(z) = \frac{N(z)}{D(z)}$ , has a closed-loop transfer function  $\frac{C(z)G(z)}{1+C(z)G(z)} = \frac{N(z)B(z)}{D(z)A(z)+N(z)B(z)}$ . Hence the closed-loop system has poles as the roots of  $D(z)A(z) + N(z)B(z) = 0$ . Clearly, the system poles have been moved from  $A(z) = 0$  in the raw system to  $D(z)A(z) + N(z)B(z) = 0$  in the closed-loop system. System performance can be significantly improved by correct selection of  $C(z) = \frac{N(z)}{D(z)}$ .

The closed-loop performance is evaluated by the speed and smoothness, or *convergence rate* and *damping*, of system's response to disturbances. The closer the system poles are to 0, the faster the system response. Although it is theoretically possible to set the closed-loop poles at 0 and make the system respond very fast, it is practically not a good idea due to the large control authority needed for fast response. In our case, it means that if we want the system respond too fast, we may sometimes have to shed a lot of data.

System damping is another important metric to evaluate closed-loop performance. It determines how smooth the system response is. Smaller damping means more severe oscillation, which is not desirable. When damping is less than 0.7, there exist visible oscillations in the system step response; when damping is bigger than 1, there is no oscillation in the system response but the system response becomes slow. Usually we choose the damping between 0.7 and 1.

With the above considerations, we develop the following feedback controller. The detailed design procedure can be found in Appendix A.

$$u(k) = \frac{H}{cT} [b_0 e(k) + b_1 e(k-1)] - au(k-1) \quad (10)$$

where  $a$ ,  $b_0$ , and  $b_1$  are controller parameters that can be easily solved from Eq.(18) and Eq.(19) in Appendix A.

*Handling time-varying characteristics of the plant.* As time goes by, the average processing cost also changes. Let us denote the per-tuple cost at period  $k$  as  $c(k)$ . As our cur-

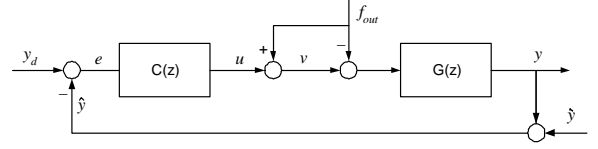


Figure 11: Control system with estimated feedback.

rent controller is designed assuming constant  $c$ , we introduce modeling errors to the closed-loop by allowing  $c$  to change over time. As mentioned in Section 3, we assume the value of  $c(k)$  changes slowly over time, at least compared to moving of data arrival rates. Under this situation, we normally believe the system is still stable with the existing basic controller. Due to its closed-loop nature, the controller should be able to compensate for the effects of such dynamics. Our experimental results (Section 5.1) provides strong evidence favoring this claim. We leave a systematic solution to handle fast-changing  $c$  as future work.

## 4.5 DSMS-specific issues

Whereas the Borealis system model seems to have a simple dynamic structure (Eq.(3)), the control of the system is far from trivial. In addition to the basic controller design, we also have to address the following practical issues.

### 4.5.1 Unavailability of real-time output measurement.

Accurate measurements of system output in real-time is essential in control system design. Unfortunately, this requirement is not met in our system because the output signal is the delay time. The output measurement is not only delayed, but also delayed by an unknown amount (the amount is the output itself!). To be more specific, the output signal of our controller should be the delay of tuples that have just entered the system when we calculate  $u(k)$ . However, at time  $k$ , we can only measure the delay of those that entered the system some time ago. This is a very interesting challenge to control theory as it does not exist in conventional control systems where the controlled signal can always be measured when we need it.

Given the output signal is not measurable when it is needed, can we derive it from the current system status? The answer is 'yes' and it comes right from the system model. We can easily modify the Borealis system to accurately record the number of outstanding data tuples (virtual queue length)  $q(k)$ . This can be done by just counting all the inflow/outflow tuples. We already know that at any time,  $c(k)$  values can be accurately estimated. Therefore, instead of using a measurement of delay  $y$  as the feedback signal, we use an estimation of  $y$  that is derived from Eq.(2):

$$\hat{y}(k) = q(k) \frac{c(k)}{H} + \frac{c(k)}{H}. \quad (11)$$

It is natural that Eq.(11) adds estimation errors to the closed-loop. We denote the estimation error as  $\tilde{y} = y - \hat{y}$ . Fortunately, our controller is still found to be robust by the following argument. When estimated output  $\hat{y}$  is used as feedback signal, the original control system becomes the one shown in Fig. 11. The output of the closed loop system is



hence described by:

$$Y(z) = \frac{C(z)G(z)}{1 + C(z)G(z)}Y_d(z) - \frac{C(z)G(z)}{1 + C(z)G(z)}\tilde{Y}(z) \quad (12)$$

The closed-loop system is still stable as long as  $\tilde{y}$  is bounded, which is always true. The  $Y_d$  term in Eq.(12) shows that the output of the closed-loop system still tracks the target reference signal with designed damping and convergence rate. However, the accuracy is compromised due to the introduction of estimation errors, as represented by the  $\tilde{Y}$  term in Eq.(12).

#### 4.5.2 Load shedder (actuator) design

Given the desired data flow rate  $v(k)$  obtained from the controller, the task of the load adaptor is to cut the incoming data stream (with rate  $f_{in}$ ) such that the actual number of tuples accepted into the system is close to  $v(k)$ . In this paper, we investigate two different ways to accomplish this.

A straightforward way to implement the load shedder is to manipulate the number of data tuples entering the DSMS query network. In other words, we treat the Borealis system as a blackbox by not shedding load within the network. For this purpose, we set a shedding/filtering factor  $\alpha$  ( $0 \leq \alpha \leq 1$ ) to all the data streams. When Borealis receives a tuple, it flips an unfair coin with head probability  $1 - \alpha$ . A tuple is accepted only when the coin shows head. At the end of period  $k$ ,  $\alpha$  should be determined as follows:

$$\alpha = 1 - [v(k)/f_{in}(k+1)]. \quad (13)$$

However,  $f_{in}(k+1)$  is unknown when we calculate  $\alpha$ . We use its value in the current period  $f_{in}(k)$  as an estimation.

Although the above load shedder is simple and works perfectly for the purpose of controlling delays given input  $v$ , it is not used in real-world systems such as Borealis. In Borealis, load can be shed from any queues in the query network. Using the network in Fig. 2 as an example, we can drop tuples in front of any combination of operators from 1 to 12 while the aforementioned load shedder only allows shedding before operators 1, 2, and 3. This difference, however, does not conflict with our system model (therefore controller design). Our model says  $y(k)$  depends on  $q(k)c$ , which is basically the outstanding ‘load’ in the queue. Shedding only intact tuples (outside the network) or partially processed tuples (in the network) makes no difference: the same ‘load’ is being discarded and  $y(k)$  depends on how much load is left in the queue. Given the  $v(k)$  generated by our controller, we know that new load with amount  $L_a = v(k)c(k+1)$  can enter the DSMS during the next period  $k+1$ . However, the outstanding tuples carry a load of  $L_q = q(k)c(k)$  and incoming streams carry a load of  $L_i = f_{in}(k+1)c(k+1)$ , which is approximated by  $f_{in}(k)c(k)$ . Therefore, load with amount of  $L_s = L_q + L_i - L_a$  is to be shed. Pass the  $L_s$  value to the Borealis load shedder, it will find the best plan to bring down the total load by  $L_s$ .

#### 4.5.3 Determination of the control period $T$

The sampling period is an important parameter in digital control systems. An improperly selected sampling period can deteriorate the performance of the closed-loop. In our setup, we consider the following two issues in selecting  $T$ :

1. *Nature of disturbances.* In order to deal with disturbances, our control loop should be able to capture the moving trends of these disturbances. The basic guiding rule for

this is the Nyquist-Shannon sampling theorem [29]. A fundamental principle in the field of information theory, the theorem states that: when sampling a signal, the sampling frequency must be greater than twice the signal frequency in order to reconstruct the original signal perfectly from the sampled version. In our setup, this means the control period should be at most half of the width of the spikes in input rate (as we assume average processing costs changes more slowly). In practice, a sampling frequency that is one order of magnitude larger than the input signal frequency is often used for signal reconstruction. Therefore, a high sampling frequency is preferred to capture the time-varying properties of the system and input data.

2. *Uncertainties in system signals.* In our problem, the output signal  $y(k)$  and processing cost  $c(k)$  are defined as the statistical expectations of a series of tuples. Taking such expectations can eliminate uncertainties brought by the heterogeneity of individual tuples. A larger sampling period (low sampling frequency) is preferred as more smoothing effects can be expected. For example, when tuple processing cost is in the order of milliseconds, setting  $T$  to a fraction of one second level would give us tens to a few hundreds of samples to approximate the real values of  $y(k)$  and  $c(k)$ . For higher sampling frequencies, we get fewer samples to estimate  $y(k)$  and may encounter estimation errors.

We need to make a tradeoff between the above two factors in choosing the right sampling period.

## 5. PERFORMANCE EVALUATION

We implemented a controller and a monitoring module in the Borealis data manager<sup>6</sup> based on our design. As the current release of Borealis does not include the load shedder presented in [26], we also built our own load shedder. The load shedder we built allows shedding from the queue and randomly selects shedding locations. In other words, it is more general than the first load shedder we discuss in Section 4.5.2 but lacks the optimization towards non-delay parameters found in the Borealis load shedder.

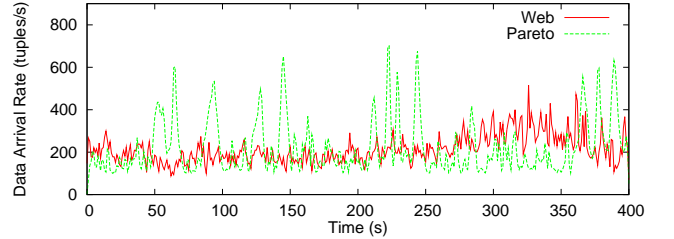


Figure 13: Traces of synthetic and real stream data.

We test our control-based framework with both synthetic and real-world stream data. The real data are traces of requests to a cluster of web servers provided by the Internet Traffic Archive.<sup>7</sup> In this dataset, each record contains a timestamp that shows when the request arrived. The synthetic data are generated in such a way that the number of data tuples per control period follows a long-tailed (Pareto, to be specific) distribution [16]. The skewness of the arrival rates is regulated by a *bias factor*  $\beta$ . The traces of a Pareto stream with  $\beta = 1$  as well as the web access data are plotted

<sup>6</sup><http://nms.lcs.mit.edu/projects/borealis/>

<sup>7</sup>dataset LBL-PKT-4, <http://ita.ee.lbl.gov>

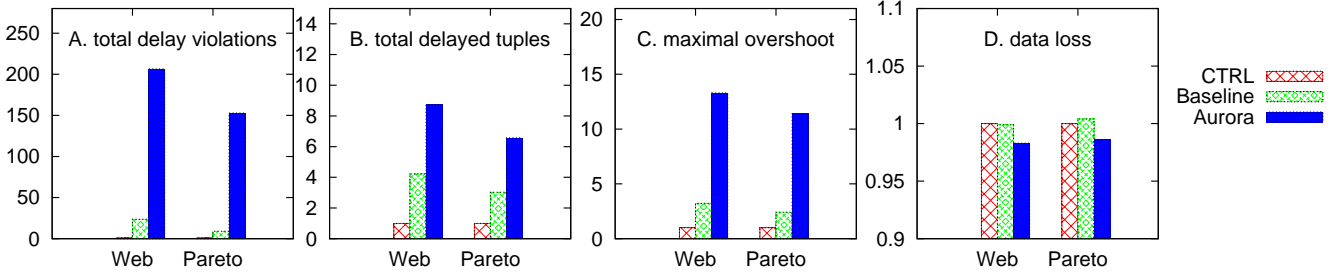


Figure 12: Relative performance of different load shedding strategies.

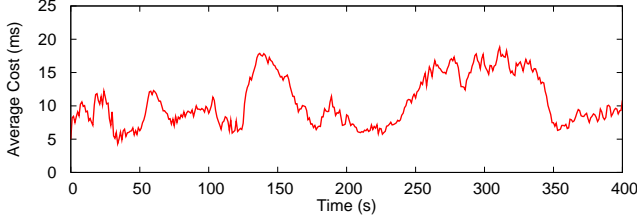


Figure 14: Variable unit processing costs (ms).

in Fig. 13. We can see that the fluctuations in the ‘Pareto’ data are more dramatic than in the ‘Web’ data.

We also use synthetic traces to simulate the variations of per-tuple cost  $c$ . We first generate the cost variations following a Pareto distribution and then modify the trace by adding ‘circumstances’ to it. For example, in the trace plotted in Fig. 14, we have a small peak at the 50th second, a large peak with a sudden jump (starting from the 125th second), and a high terrace with a sudden drop (250th to 350th second). The same network for system identification (Section 4.2) is used for experimental studies.

We compare our control-based framework (referred to as *CTRL* hereafter) with the following two approaches:

- *AURORA*: the algorithm utilized in the current Aurora/Borealis system, as shown in Fig. 1. At the  $k$ -th control period, the measured load  $L$  is  $f_{in}(k-1)$ . To deal with variable per-tuple cost, we define  $L_0 = H/c(k-1)$ . This method represents the current best solution in load shedding in DSMSSs;
- *BASELINE*: a simple feedback control-based method: it takes system status (i.e.,  $q(k), c(k)$  in our case) into account in making decisions. Specifically,  $v(k)$  is obtained from the system model (Eq.(11)): the target value of  $y_d$  would allow  $y_d H/c(k)$  outstanding tuples, therefore  $u(k) = y_d H/c(k) - q(k)$  more tuples can be added to the queue. Consequently, we get  $v(k) = u(k) + f_{out}(k) = -q(k) + \frac{y_d H}{c(k)} + \frac{T H}{c(k)}$ . As  $c(k)$  is unknown, we estimate it with  $c(k-1)$ . This method is used to test the importance of controller design.

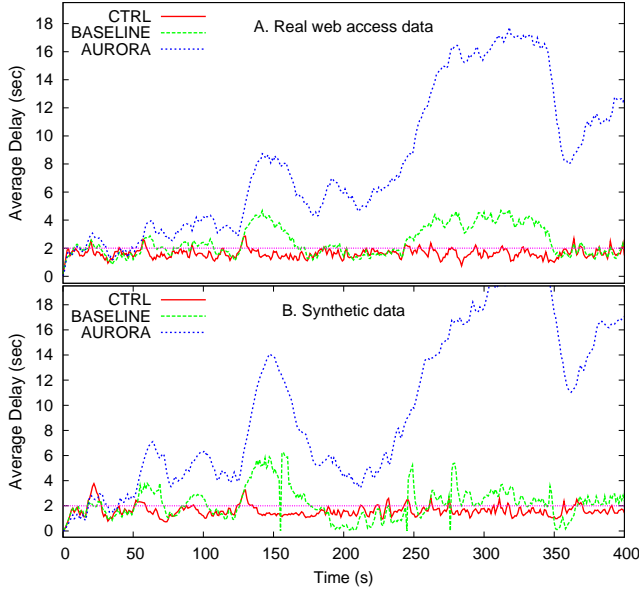
In all the experiments we report in this section, we set target delay value  $y_d$  to 2000 milliseconds unless specified otherwise. We run all tests for 400 seconds. For *CTRL*, the controller parameters value identified by our analysis are:  $b_0 = 0.4$ ,  $b_1 = -0.31$ , and  $a = -0.8$ . Any set of parameters that are solutions to Equations (18) and (19) are supposed

to have the same performance. This is verified by our tests with other set of parameters (details skipped). Following the experiments shown in Fig. 6, we set  $H$  to 0.97. The control period is set to 1000 milliseconds. Going back to Fig. 13, we see that most of the bursts in both traces last longer than a few (4 to 5) seconds therefore a sampling period smaller than two seconds is preferred according to the sampling theorem. The change of costs  $c$  in Fig. 14 has peaks with widths on the order of tens of seconds (with some exceptions) thus one-second period is definitely sufficient. We also test the systems with different choices of  $T$  and  $y_d$ .

## 5.1 Experimental results

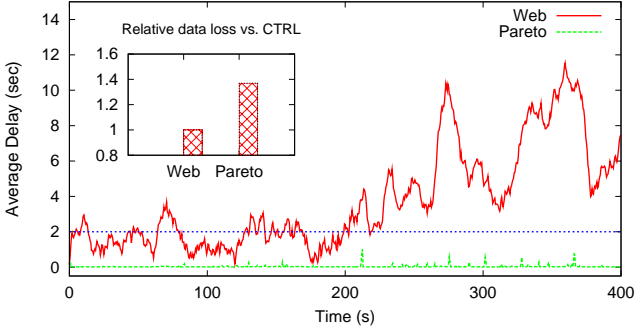
We first compare the long-term performance of *CTRL* with that of the two other algorithms. In Fig.12, we plot the ratios of all four metrics measured (i.e, totals in the 400-second period) from the *AURORA* and *BASELINE* experiments to that of *CTRL*. For example, when injected with the same ‘Web’ data stream, Fig. 12A shows that *AURORA* rendered 205 times more total delay violations than *CTRL* and *BASELINE* had 23 times. Similar results were obtained in total delayed tuples (Fig. 12B) and maximal overshoot (Fig. 12C). Note all data points for *CTRL* are 1.0 in Fig. 12. The data loss ratio for all methods are almost the same with *AURORA* losing slightly fewer tuples (0.986 for ‘Web’ and 0.987 for ‘Pareto’). It is easy to see that, for both real (‘Web’) and synthetic (‘Pareto’) data inputs, *CTRL* is the easy winner in the three delay-related metrics with almost the same amount of data loss. The *BASELINE* method, as a feedback solution, has worse performance than *CTRL* but it also beats *AURORA*.

To better understand the above long-term results, we show the transient performance of all three methods by plotting  $y(k)$  values measured at all control periods in Fig. 15. We can see that, as expected, almost all output in *CTRL* is very close to the target value of two seconds. For *BASELINE* and *AURORA*, we can observe peaks that are large in both height and width. Such peaks are the results of either fluctuations of arrival rate or changes of  $c$  (e.g., those at about 50th second and 125th second, and the high terrace starting from the 230th second). Note the first two peaks of  $c$  also have impact on the *CTRL* system: average delay increases beyond two seconds. However, with the design goal of fast convergence and high damping, the controller in the *CTRL* system can quickly bring the system back to a stable state thus large peaks of  $y$  are avoided. The high terrace has almost no effect on *CTRL*. This is because the value of  $c$  increases gradually before the terrace. Our controller can capture and compensate for this kind of gradual change



**Figure 15: Performance of different load shedding methods.**

while the open-loop system cannot (i.e., Example 2 in Section 4.3.2). From Fig. 15, we can fairly conclude that the design goal of our controller (Section 4.4.1) is achieved.

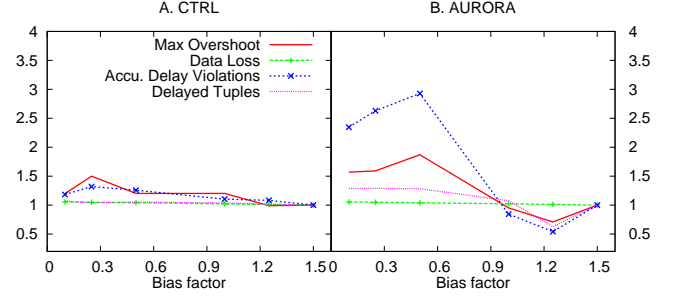


**Figure 16: Performance of 'Aurora' with  $H = 0.96$ .**

*System robustness.* In the above experiments, the *AURORA* method has poor stability: tuple delays increase all the time. A question one might ask is: can we remedy the problem by using a smaller  $L_0$  value (recall the algorithm in Fig. 1) such that more data can be discarded? In our setup, this means the same as changing the  $H$  value (even though  $H = 0.97$  is proved to be correct in Section 4.2) as we define  $L_0$  to be  $H/c$ . Fig. 16 shows the results of the *AURORA* method under both real and synthetic data inputs using a smaller  $H$  value of 0.96. For the 'Web' data inputs, the system is still unstable. Surprisingly, no delay violations can be observed for the 'Pareto' inputs. However, the price for this is huge: it costs 37% more data loss than *CTRL* (small graph in Fig. 16). This result shows the poor robustness of open-loop solutions: it is hard to tune the system as performance depends heavily on the pattern of inputs.

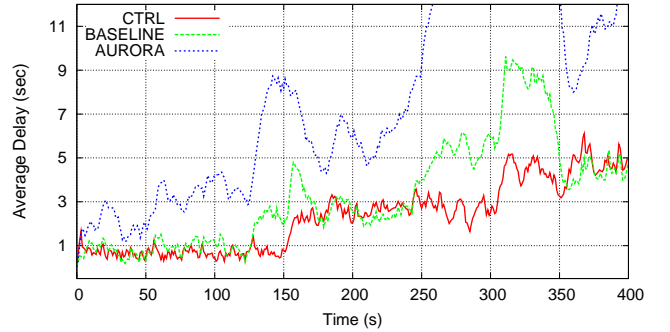
To further study the robustness of the three methods, we test them using data streams with different levels of bursti-

ness. Specifically, we feed the systems with synthetic data streams with bias factors  $\beta$  of 0.1, 0.25, 0.5, 1, 1.25, and 1.5. The smaller the bias factor, the more bursty the input. In Fig. 17, we show the change of all four metrics with respect to the bias factor. All numbers plotted are relative to the corresponding value measured in the case of  $\beta = 1.5$ . As the input stream becomes more bursty, very little difference can be observed in *CTRL* (Fig. 17A) while the changes in *AURORA* (Fig. 17B) are much more dramatic. The performance of *BASELINE* is not significantly affected by the bias factor as well (data not shown).



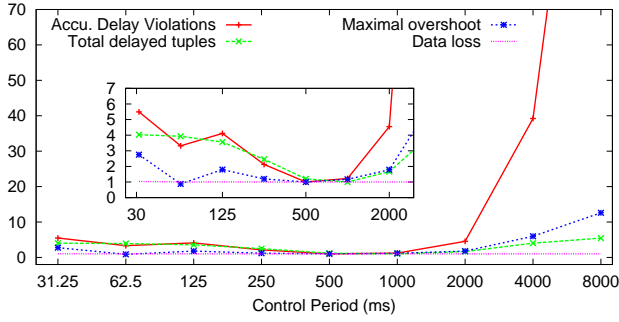
**Figure 17: Effects of input burstiness on performance.**

In Fig. 18 we show how different load shedding methods respond to changes of target value  $y_d$  at runtime. In these experiments, we set  $y_d$  to be 1000 milliseconds initially and change it to 3000 milliseconds at the 150th second and then to 5000 milliseconds at the 300th second. We can see that *CTRL* converges to the new target values very quickly. Furthermore, system stability is not affected by the target value. The *AURORA* method does not respond to the changes of  $y_d$  at all as it is open-loop. When  $y_d$  changes, it takes the *BASELINE* method very long time to converge to the new target value. We use 'Web' data inputs for the experiments in Fig. 18 and using 'Pareto' data gives similar results.



**Figure 18: Responses to change of target value.**

*Effects of control period.* In Section 4.5.3, we discussed the general rules on choosing the right sampling period. We test these rules by running experiments with nine different sampling periods ranging from 31.25 to 8000 milliseconds with the *CTRL* system and 'Web' data stream. In Fig. 19, every data point is the ratio to the lowest corresponding metric in all nine tests and the x-axis has a logarithmic scale. For example, the smallest accumulated delay violations were recorded under  $T = 500ms$  and this value is about



**Figure 19: Performance under different control period.**

40 times as high when  $T = 4000ms$ . Obviously, the magnitude and frequency of delay violations increase desperately when  $T$  is beyond four seconds. As expected, a shorter control period is preferred. This confirms our discussion about the sampling theorem in Section 4.5.3. When  $T$  becomes too small, performance degrades. The best region seems to be  $[250, 1000]$  in these experiments (see small graph in Fig. 19). Similar results are obtained for the ‘Pareto’ data inputs. One thing to point out is: in all experiments with  $T$  smaller than 4000ms, *CTRL* outperforms *BASELINE* and *AURORA* with similar difference shown in Fig. 12.

*Computational overhead.* The operation of our controller only involves several floating point calculations at each control period. In our experimental platform with a Pentium 4 2.4GHz CPU, this time is only about 20 microseconds. This is trivial because the control period is set to be (at least) on the order of hundreds of milliseconds.

## 5.2 Discussion

From the above experiments, we believe our load shedding framework based on a feedback control gives better quantitative guidance to how load shedding should be performed. We see that *CTRL* is the winner in all delay-related metrics. We achieve this by employing a number of techniques. The first lesson we learn from this study is: a thorough understanding of system dynamics is extremely useful in dealing with control-like problems. The idea of controlling delay through virtual queue length provides big advantages over the *AURORA* method. By applying simple rules derived from the model, the *BASELINE* method achieves far better performance than *AURORA*. Decisions based on controller design is another plus for our method. With guaranteed convergence, our controller avoids large and long-time deviations from the desired output while the *BASELINE* method suffers from such deviations.

An important feature of the control-based solution is its robustness. Note that we only use standard inputs to validate the system model and controller tuning is accomplished by mathematical reasoning exclusively. In other words, no training data is needed and performance can be guaranteed for a wide range of inputs. On the contrary, tuning of other methods can be *ad hoc*, as evidenced by the dependence of open-loop solutions on the pattern of data inputs. We have reasons to believe that even the current system model (Fig.4) can be used for DSMSs other than Borealis: we noticed (via experiments) that modifying the query network only changes a parameter ( $c$ ) but not necessarily the struc-

ture of the model. It is highly possible that the model is still applicable to a wide range of scheduling policies that do not consider tuple priorities. Further investigations are needed. In the *CTRL* system, the only thing that’s subject to input/internal uncertainties is the control period  $T$ . However, the proper choice of  $T$  requires very little information about such uncertainties (i.e., signal frequency), which is generally available. For a wide range of  $T$  values, the *CTRL* method still beats the other two algorithms.

## 6. CONCLUSIONS AND FUTURE WORK

This paper argues for the importance of managing data processing delays in data stream management systems. Violations of delay are generally caused by system overloading. Load shedding and other adaptation strategies have been exploited to compensate for degraded delays under overloading. We noticed that the strong dynamics such as bursty arrival pattern of data stream applications require an adaptation strategy with excellent transient-state performance (e.g., fast convergence to steady state), which most of the current works in this area fail to provide. We proposed a load shedding framework that leverages various techniques from the field of control theory. We started by developing a dynamic model of a steam management system. We then construct a feedback control loop to guide load shedding through system analysis and rigorous controller design.

We have implemented our design and performed extensive experiments on a real-world system - the Borealis stream manager. It is evident that our approach achieves better performance in terms of reduced delay violations over current strategies that do not consider system status in decision-making. The control-based strategy is also robust and light-weight. Finally, we believe our explorations can give rise to many opportunities to conduct synergistic research between the database and control engineering communities to extend our knowledge in both fields.

Immediate follow-up work includes more experiments on the Aurora load shedder and more dramatic changes of per-tuple costs (resulting from structure change of the query network, for example). The idea is to use adaptive control techniques to capture the internal variations of the system model and provide better control over the whole system. Our control-based framework can also be extended in a few directions. First of all, there is still room to improve the quality model: we could provide heterogeneous quality guarantees for streams with different priorities; and multiple quality dimensions can be supported at the same time by introducing a multi-in-multi-out control model. Combining stochastic methods such as Kalman Filters with our controller design would yield more powerful adaptation algorithms. Although we did not go into prediction strategies of time series in this paper, we understand that it is a promising direction that is worth serious consideration. We are currently investigating the potential of control theory in a number of other topics in DBMS research such as query re-optimization and dynamic resource allocation in traditional databases.

## 7. REFERENCES

- [1] Niagara Project, <http://www.cs.wisc.edu/niagara/>.
- [2] D. J. Abadi, D. Carney, U. Çetintemel, M. Cherniack, C. Convey, S. Lee, M. Stonebraker, N. Tatbul, and



- S. B. Zdonik. Aurora: a new model and architecture for data stream management. *VLDB Journal*, 12(2):120–139, 2003.
- [3] Daniel J. Abadi, Yanif Ahmad, M. Balazinska, U. Çetintemel, M. Cherniack, J.-H. Hwang, W. Lindner, A. Maskey, A. Rasin, E. Ryvkina, N. Tatbul, and S. Zdonik. The Design of the Borealis Stream Processing Engine. In *Procs. of 2nd CIDR Conference*, January 2005.
- [4] A. Arasu, B. Babcock, S. Babu, M. Datar, J. Rosenstein, K. Ito, I. Nishizawa, and J. Widom. Query Processing, Resource Management, and Approximation in a Data Stream Management System. In *Procs. of 1st CIDR Conf.*, 2003.
- [5] B. Babcock, S. Babu, M. Datar, and R. Motwani. Chain: Operator Scheduling for Memory Minimization in Data Stream Systems. In *Proceedings of ACM SIGMOD '03*, pages 253–264, June 2003.
- [6] B. Babcock, M. Datar, and R. Motwani. Load Shedding for Aggregation Queries over Data Streams. In *Procs. of ICDE Conf.*, 2004.
- [7] D. Carney, U. Çetintemel, M. Cherniack, C. Convey, S. Lee, G. Seidman, M. Stonebraker, N. Tatbul, and S. Zdonik. Monitoring Streams - A New Class of Data Management Applications. In *Procs. of the 28th VLDB Conf.*, pages 84–89, August 2002.
- [8] D. Carney, U. Çetintemel, A. Rasin, S. Zdonik, M. Cherniack, and M. Stonebraker. Operator Scheduling in a Data Stream Manager. In *Procs. of the 29th VLDB Conf.*, pages 838–849, August 2003.
- [9] S. Chandrasekaran, A. Deshpande, M. Franklin, J. Hellerstein, Wei Hong, S. Krishnamurthy, S. Madden, V. Raman, F. Reiss, and M. Shah. TelegraphCQ: Continuous Dataflow Processing for an Uncertain World. In *Proceedings of 1st CIDR Conference*, January 2003.
- [10] Yun Chi, Haixun Wang, and Philip S. Yu. Load Star: Load Shedding in Data Stream Mining. In *Procs. of the 31st VLDB Conf.*, pages 1302–1305, August 2005.
- [11] Yun Chi, Haixun Wang, Philip S. Yu, and Richard R. Muntz. Loadstar: A load shedding scheme for classifying data streams. In *Procs. of SIAM Conf. on Data Mining*, page 3, 2005.
- [12] A. Dobra, M. Garofalakis, J. Gehrke, and R. Rastogi. Processing complex aggregate queries over data streams. In *Procs. of SIGMOD Conf.*, pages 61–72, June 2002.
- [13] G. F. Franklin, J. D. Powell, and A. Emami-Naeini. *Feedback Control of Dynamic Systems*. Prentice Hall, Massachusetts, 2002.
- [14] The STREAM Group. STREAM: The Stanford Stream Data Manager. *IEEE Data Engineering Bulletin*, 26(1):19–26, March 2003.
- [15] M. Hammad, M. Franklin, W. Aref, and A. Elmagarmid. Scheduling for Shared Window Joins Over Data Streams. In *Proceedings of 29th VLDB Conf.*, pages 297–308, August 2003.
- [16] Mor Harchol-Balter, Mark Crovella, and Cristina Murta. On Choosing a Task Assignment Policy for a Distributed Server System. *Journal of Parallel and Distributed Computing*, 59(2):204–228, November 1999.
- [17] Kyoung-Don Kang, Sang H. Son, and John 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.
- [18] S. Keshav. A Control-Theoretic Approach to Flow Control. In *Procs. of ACM SIGCOMM*, September 1991.
- [19] B. Li and K. Nahrstedt. A Control-Based Middleware Framework for Quality of Service Adaptations. *IEEE Journal of Selected Areas in Communications, Special Issue on Service Enabling Platforms*, 17(9):1632–1650, September 1999.
- [20] C. Lu, J. Stankovic, G. Tao, and S. Han. Feedback Control Real-Time Scheduling: Framework, Modeling, and Algorithms. *Journal of Real-Time Systems*, 23(1/2):85–126, September 2002.
- [21] K. Nahrstedt, D. Xu, D. Wichadakul, and B. Li. QoS-Aware Middleware for Ubiquitous and Heterogeneous Environments. *IEEE Communications Magazine*, 2001.
- [22] Klara Nahrstedt and Ralf Steinmetz. Resource Management in Networked Multimedia Systems. *IEEE Computer*, 28(5):52–63, 1995.
- [23] C. Olston, J. Jiang, and J. Widom. Adaptive Filters for Continuous Queries over Distributed Data Streams. In *Proceedings of ACM SIGMOD '03*, pages 563–574, June 2003.
- [24] V. Paxson and S. Floyd. Wide-Area Traffic: The Failure of Poisson Modeling. *IEEE/ACM Transactions on Networking*, 3(3):226–244, 1995.
- [25] Frederick Reiss and Joseph M. Hellerstein. Data Triage: An Adaptive Architecture for Load Shedding in TelegraphCQ. In *Proceedings of ICDE*, pages 155–156, April 2005.
- [26] N. Tatbul, U. Çetintemel, S. Zdonik, M. Cherniack, and M. Stonebraker. Load Shedding in a Data Stream Manager. In *Proceedings of the 29th VLDB Conference*, pages 309–320, August 2003.
- [27] Yi-Cheng Tu, Mohamed Hefeeda, Yuni Xia, and Sunil Prabhakar. Control-based Quality Adaptation in Data Stream Management Systems. In *Proceedings of DEXA*, pages 746–755, August 2005.
- [28] S. Viglas and J. Naughton. Rate-Based Query Optimization for Streaming Information Sources. In *Procs. of SIGMOD Conf.*, pages 37–48, June 2002.
- [29] W.D. Stanley. *Digital Signal Processing*. Reston Publishing Co., 1975.
- [30] W. Xu, J. L. Hellerstein, W. Kramer, and D. Patterson. Control Considerations for Scalable Event Processing. In *Procs. IFIP/IEEE Workshop on Distributed Systems: Operations and Management*, pages 233–244, October 2005.
- [31] D. Yau and S. Lam. Operating System Techniques for Distributed Multimedia. *International Journal of Intelligent Systems*, 13(12):1175–1200, December 1998.
- [32] M. Zhang, T. Madhyastha, N.H. Chan, S. Papadimitriou, and C. Faloutsos. Data Mining Meets Performance Evaluation: Fast Algorithms for Modeling Bursty Traffic. In *Proceedings of the 18th ICDE Conference*, pages 507–516, February 2002.

- [33] Y. Zhu and D. Shasha. StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time. In *Procs. of the 28th VLDB Conf.*, pages 358–369, August 2002.

## Acknowledgments

The authors would like to thank Ms. Nesime Tatbul, Prof. Uğur Çetintemel, and Prof. Stan Zdonik from Brown University for their help on the Aurora/Borealis system.

## APPENDIX

### A. CONTROLLER DESIGN BASED ON POLE PLACEMENT

In this study, we set the desired convergence rate to three sampling periods. This means the system, in response to dynamics, would converge to  $1 - \frac{1}{e} \approx 63\%$  of the desired value in 3 control periods and to 98% in 12 periods. We set the system damping to 1 and set the desired closed-loop poles to be on the real axis, at 0.7. Thus, the desired closed-loop characteristic equation (CLCE) is:

$$(z - 0.7)^2 = z^2 - 1.4z + 0.49 = 0 \quad (14)$$

According to Eq.(4), we have a first-order system. Thus, the controller  $C(z)$  will have one pole and its generic format [13] is

$$C(z) = \frac{H(b_0z + b_1)}{cT(z + a)} \quad (15)$$

where  $b_0, b_1$ , and  $a$  are controller parameters. Therefore, the closed-loop transfer function (CLTF) becomes

$$\frac{C(z)G(z)}{1 + C(z)G(z)} = \frac{b_0z + b_1}{z^2 + (a - 1 + b_0)z + (-a + b_1)} \quad (16)$$

and the actual closed-loop characteristic equation (CLCE) is

$$\text{and } z^2 + (a - 1 + b_0)z + (-a + b_1) = 0. \quad (17)$$

Matching the above CLCE to its desired form shown in Eq.(14), we get the following (Diophantine) equation:

$$z^2 + (a - 1 + b_0)z + (-a + b_1) = z^2 - 1.4z + 0.49 \quad (18)$$

At the steady state, the CLTF should have a static gain that equals one, meaning we want the output  $y$  to be exactly the same as  $y_d$ . This results in the following equality:

$$\left. \frac{b_0z + b_1}{z^2 + (a - 1 + b_0)z + (-a + b_1)} \right|_{z=1} = 1 \quad (19)$$

Solving Equations (18) and (19), one can obtain the controller parameters  $a$ ,  $b_0$ , and  $b_1$ .

In summary, the above design results in a closed loop system having two poles, both of which are on the positive real axis at 0.7.

Now we can generate the control signal  $u$ . Let  $U(z)$  and  $E(z)$  be the  $z$ -transforms of  $u$  and error  $e$ , respectively. According to Fig. 10,  $e$  is the input and  $u$  is the output respect to the controller, we have

$$U(z) = C(z)E(z) = \frac{H(b_0z + b_1)}{cT(z + a)}E(z).$$

Multiplying both sides by  $\frac{(z+a)cT}{z \cdot H}$ , we get

$$U(z) \frac{cT}{H} + \frac{acT}{H} U(z) z^{-1} = b_0 E(z) + b_1 E(z) z^{-1}.$$

By inverse  $z$ -transformation, the above leads to the solution for  $u$  as follows:

$$u(k) = \frac{H}{cT} [b_0 e(k) + b_1 e(k - 1)] - au(k - 1).$$