

On Fundamental Tradeoffs between Delay Bounds and Computational Complexity in Packet Scheduling Algorithms

Jun Xu and Richard J. Lipton
College of Computing
Georgia Institute of Technology
{jx,rjl}@cc.gatech.edu

ABSTRACT

In this work, we clarify, extend and solve an open problem concerning the computational complexity for packet scheduling algorithms to achieve tight end-to-end delay bounds. We first focus on the difference between the time a packet finishes service in a scheduling algorithm and its virtual finish time under a GPS (General Processor Sharing) scheduler, called *GPS-relative delay*. We prove that, under a slightly restrictive but reasonable computational model, the lower bound computational complexity of any scheduling algorithm that guarantees $O(1)$ GPS-relative delay bound is $\Omega(\log_2 n)$ (widely believed as a “folklore theorem” but never proved). We also discover that, surprisingly, the complexity lower bound remains the same even if the delay bound is relaxed to $O(n^a)$ for $0 < a < 1$. This implies that the delay-complexity tradeoff curve is “flat” in the “interval” $[O(1), O(n)]$. We later extend both complexity results (for $O(1)$ or $O(n^a)$ delay) to a much stronger computational model. Finally, we show that the same complexity lower bounds are conditionally applicable to guaranteeing tight end-to-end delay bounds. This is done by untangling the relationship between the GPS-relative delay bound and the end-to-end delay bound.

Categories and Subject Descriptors

F.2 [Theory of Computation]: Analysis of Algorithms and Problem Complexity; C.2 [Computer System Organization]: Computer Communication Networks

General Terms

Algorithms, Performance, Theory

*The work of both authors was supported in part by NSF grant no. ITR/SY 0113933. The second author was also supported in part by NSF grant no. CCR 0002299 and by Telcordia Research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGCOMM'02, August 19-23, 2002, Pittsburgh, Pennsylvania, USA.
Copyright 2002 ACM 1-58113-570-X/02/0008 ...\$5.00.

Keywords

Computational complexity, packet scheduling, Quality of Service, delay bound, decision tree

1. INTRODUCTION

Packet scheduling is an important mechanism in providing QoS guarantees in data networks [4, 7, 19]. The fairest algorithm for packet scheduling is General Processor Sharing (GPS) [4, 12]. However, GPS is not a realistic algorithm since in a packet network, service is performed packet-by-packet, rather than “bit by bit” as in GPS. Nevertheless, GPS serves as a reference scheduler that real-world packet-by-packet scheduling algorithms (e.g., WFQ [4]) can be compared with in terms of end-to-end delay bounds and fair bandwidth allocation.

In a link of rate r served by a GPS scheduler, each session $i = 1, 2, \dots, n$ is assigned a weight value ϕ_i . Each backlogged session j at every moment t is served simultaneously at rate $r_j = r\phi_j / (\sum_{j \in B(t)} \phi_j)$, where $B(t)$ is the set of sessions that are backlogged at time t . One important property of GPS, proved in [12], is that it can guarantee tight end-to-end delay bound to traffic that is leaky-bucket [17] constrained.

It is interesting to look at the *GPS-relative delay* of a packet served by a scheduling algorithm *ALG* as compared to *GPS*. For each packet p , it is defined as $\max(0, F_p^{ALG} - F_p^{GPS})$, where F_p^{ALG} and F_p^{GPS} are the times when the packet p finishes service in the *ALG* scheduler and in the *GPS* scheduler, respectively. It has been shown in [12] and [2] respectively that *WFQ* and *WF²Q* schedulers both have a worst-case GPS-relative delay bound of $\frac{L_{max}}{r}$, where L_{max} is the maximum packet size in the network and r is the total link bandwidth. That is, for each packet p ,

$$F_p^{WFQ} - F_p^{GPS} \leq \frac{L_{max}}{r} \quad (1)$$

$$F_p^{WF^2Q} - F_p^{GPS} \leq \frac{L_{max}}{r} \quad (2)$$

We simply say that the delay bound is $O(1)$ since L_{max} and r can be viewed as constants independent of the number of sessions n . *WFQ* and *WF²Q* achieve this $O(1)$ delay bound by (a) keeping perfect track of the GPS clock and (b) picking among all (in *WFQ*) or all eligible (in *WF²Q*) head-of-session packets, the one with smallest GPS virtual finish time to serve next. The per-packet worst-case computational complexity of the second part ((b) part) in both

WFQ and WF^2Q is $O(\log_2 n)$. In other words, the computational cost to “pay” for the $O(1)$ GPS-relative delay bound in both WFQ and WF^2Q is $O(\log_2 n)$ ¹.

On the other hand, round-robin algorithms such as DRR (Deficit Round Robin) [13] and WRR (Weighted Round Robin) [10] have a low implementation cost of $O(1)$. However, they in general cannot provide the tight GPS-relative delay bound of $\frac{L_{max}}{r}$. In fact, the best possible delay bound they can provide is $O(n)$. This is illustrated in Fig. 1. We assume that these n sessions share the same link and have the same weight. Without loss of generality, we also assume that these sessions are served in the round-robin order $1, 2, \dots, n$. At time 0, packets of length M have arrived at sessions $1, 2, \dots, n-1$, and a packet of length $m < M$ has arrived at session n . Suppose M is no larger than the *service quantum* size used in round-robin algorithms so that all these packets are in the same service frame. Then clearly the short packet in session n will be served behind $n-1$ long packets. So the GPS-relative delay of the short packet can be calculated as $\frac{(n-1)(M-m)}{r}$, which is $O(n)$.

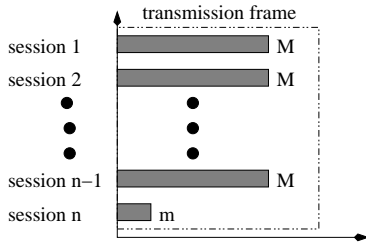


Figure 1: How round robin algorithms incur $O(n)$ GPS-relative delay

We have just shown that algorithms with $O(\log_2 n)$ complexity (GPS time tracking overhead excluded) such as WFQ and WF^2Q can provide $O(1)$ GPS-relative delay bound, while $O(1)$ round-robin algorithms such as DRR and WRR can only guarantee a delay bound of $O(n)$. An open problem proposed in Sigcomm’01 by Guo (author of [8]) is whether this represents indeed the fundamental tradeoff between computational complexity of the scheduling algorithm and the GPS-relative delay bound they can achieve. More specifically, Guo asks whether it is possible to design an $o(\log_2 n)$ (ideally $O(1)$) algorithm to guarantee $O(1)$ GPS-relative delay bound. Our work clarifies and extends this question, and answers it in a comprehensive way.

The first major result of this paper is to show that $\Omega(\log_2 n)$ is indeed the complexity lower bound to guarantee $O(1)$ GPS-relative delay², excluding the cost of tracking GPS time. This bound is established under the decision tree computation model that allows direct comparisons between its inputs (in our context between GPS virtual finish times of the packets). This model seems slightly restrictive but is

¹Here the cost of the GPS clock tracking ((a) part) is not included.

²Leap Forward Virtual Clock (LFVC) scheduler has a low implementation complexity of $O(\log(\log n))$ using timestamp discretization, but may incur $O(n)$ GPS-relative delay in the worst case. This is because, with small but positive probability, the “discretization error” may add up rather than cancel out.

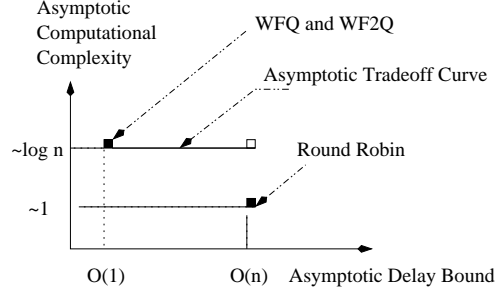


Figure 2: The asymptotic tradeoff curve between delay bound and computational complexity

reasonable for our context, since such comparisons are indeed sufficient for assuring $O(1)$ GPS-relative delay bound in WFQ and WF^2Q [12, 2]. This result granted for the moment, we now have two points on the complexity-delay tradeoff curve, as shown in Fig. 2. One is $O(n)$ delay at the complexity of $\Omega(1)$ and the other is the $O(1)$ delay at the complexity of $\Omega(\log_2 n)$. One interesting question to ask is how do other parts of the “tradeoff curve” look. More specifically, to guarantee a delay bound that is asymptotically between $O(1)$ and $O(n)$, say $O(\sqrt{n})$, can the complexity of packet scheduling be asymptotically lower than $\Omega(\log_2 n)$, say $\Omega(\sqrt{\log_2 n})$? The result we discover and prove is surprising: for any fixed $0 < a < 1$, the asymptotic complexity for achieving $O(n^a)$ delay is always $\Omega(\log_2 n)$. As shown in Fig. 2, this basically says that the asymptotic tradeoff curve is “flat” and has a “jump” at $O(n)$.

The second major result of this paper is to strengthen the aforementioned lower bounds by extending them to a much stronger computational model: decision tree that allows *linear comparisons*. However, under this computational model, we are able to prove the same complexity lower bounds of $\Omega(\log_2 n)$ only when the scheduling algorithm guarantees $O(1)$ or $O(n^a)$ ($0 < a < 1$) *disadvantage delay* bound. *Disadvantage delay* is a slightly stronger type of delay than the GPS-relative delay, since for each packet, its *disadvantage delay* is no smaller than its GPS-relative delay. Nevertheless, the second result is provably stronger (by Theorem 5) than our first result (for both $O(1)$ and $O(n^a)$ cases).

Our third and final result is to show that the same complexity lower bounds can to a certain extent be extended to guaranteeing tight end-to-end delay bounds. This is done by studying the relationship between the GPS-relative delay bound and the end-to-end delay bound. In particular we show that, providing tight GPS-relative delay bound of $\frac{L_{max}}{r}$ in Latency Rate (LR) schedulers (introduced in [15]) is conditionally equivalent to providing the tight *latency bound* of $\frac{L_{max}}{r} + \frac{L_{max,i}}{r_i}$, where $L_{max,i}$ is the maximum size of a packet in session i and r_i is the guaranteed rate of session i . This allows us to prove that the per packet computational complexity for guaranteeing a tight latency bound of $O(n^a \frac{L_{max}}{r} + \frac{L_{max,i}}{r_i})$ for all $L_{max,i} > 0$ is $\Omega(\log_2 n)$, under certain reasonable conditions.

Though it is widely believed as a “folklore theorem” that scheduling algorithms which can provide tight end-to-end delay bounds require $\Omega(\log_2 n)$ complexity (typically used for maintaining a priority queue), it has never been care-

fully formulated and proved. To the best of our knowledge, our work is the first major and successful step in establishing such complexity lower bounds. Our initial goal was to show that the $\Omega(\log_2 n)$ delay bounds hold under the decision tree model that allows linear comparisons. Though we are not able to prove this result in full generality, our rigorous formulation of the problem and techniques introduced in proving slightly weaker results serve as the basis for further exploration of this problem.

The rest of the paper is organized as follows. In Section 2, we introduce the computational models and assumptions we will use in proving our results. The aforementioned three major results are established in Section 3, 4, and 5 respectively. Section 6 concludes the paper.

2. ASSUMPTIONS AND COMPUTATIONAL MODELS

In general, complexity lower bounds of a computing problem are derived based on problem-specific assumptions and conditions, and a computational model that specifies what operations are allowed in solving the problem and how they are “charged” in terms of complexity. In Section 2.1, we describe a network load and resource allocation condition called CBFS (continuously backlogged fair sharing) under which all later lower bounds will be derived. In Section 2.2, we introduce two computational models that will be used in Section 3 and 4, respectively. Finally in Section 2.3, we discuss why decision tree computational models are chosen for studying complexity lower bounds.

2.1 CBFS condition

All lower bounds in this paper will be derived under a network load and resource sharing condition called *continuously backlogged fair sharing* (CBFS). Let n be the number of sessions and r be the total bandwidth of the link. In CBFS,

- (Fair Sharing) Each session has equal weight, that is, for any $1 \leq i < j \leq n$, $\phi_i = \phi_j$.
- (Continuously Backlogged) Each session has a packet arrival at time 0. Also, for any $t > 0$ and $1 \leq i \leq n$, $A_i(t) \geq \frac{r}{n}t$. Here $A_i(t)$ is the amount of session i traffic that has arrived during the interval $[0, t]$.

We call the second part of the condition “continuously backlogged” because if these sessions are served by a GPS scheduler, they will be continuously backlogged from time 0. This is proved in the next proposition.

PROPOSITION 1. *For any packet arrival instance that conforms to the CBFS condition, each and every session will be continuously backlogged when served by a GPS scheduler.*

PROOF. The proof is by contradiction. Suppose some sessions may become unbacklogged at certain points of time. We can view packet scheduling as an event-driven system in which the events are the arrivals and departures of the packets. Since all sessions are backlogged at time 0, the following is the only possible way that session i may become unbacklogged at time t : a packet departs from session i at time t , and its next packet does not arrive until time $\tau > t$ ($\tau = \infty$ if there is no such arrival). Let t^* be the time that

the *earliest* such packet departure event happens. Suppose this happens to session i^* , and session i^* does not become backlogged until $\tau^* > t^*$. By the definition of t^* , all sessions are continuously backlogged between $[0, t^*]$. So, under the GPS scheduler, the amount of service each session receives during this period is the same, which is $\frac{r}{n}t^*$. Let $\tau^* > t' > t^*$ (to avoid the case $\tau^* = \infty$). Then the amount of service session i^* receives during the interval $[0, t']$ is $\frac{r}{n}t^* < \frac{r}{n}t'$, which violates the second part of the CBFS condition. \square

Since our lower bounds are on the computational complexity in the worst case, the general lower bounds can only be **higher than or equal to** the bounds derived under the CBFS condition (i.e., we don’t “gain” from this condition). The significance of this condition is profound:

- First, computing the GPS virtual finish time of a packet p becomes an $O(1)$ operation (see remark after Proposition 2). So CBFS condition allows us to “naturally exclude” the cost of tracking GPS clock.
- Second, we will show that under the CBFS condition, many existing scheduling algorithms such as Virtual Clock (VC) [20], Frame-based Fair Queuing (FFQ) [14] and WF^2Q+ [3] are equivalent to either WFQ or WF^2Q (Proposition 3). So whenever we need to relate our results to these scheduling algorithms, we only need to study WFQ and WF^2Q .
- Third, the complexity lower bounds that are proved under this condition are still tight enough. In other words, we are not “losing” too much ground on complexity lower bounds when restricted by this condition.

In our later proofs, we assume that the size of a packet can take any real number between L_{min} and L_{max} , which denote the minimum and the maximum packet sizes respectively. This is, in general, not true for packet networks. However, it can be shown that if we remove part one (fair sharing) of the CBFS condition and instead allow weighted sharing (with part two adjusted accordingly), we do not need to insist on such freedom in packet size. In fact, our proofs will work even for ATM networks where fixed packet size is used. Since this proof is less interesting, we omit it here to save space.

For simplicity of discussion in our later proofs, we further assume that $L_{min} = 0$, i.e., the minimum packet size can be as small as zero. It can be shown that all later proofs remain true when $L_{min} > 0$, with adjustments on the constants used in the proofs. The proof of this claim is also omitted due to space limitations.

In the following, we prove two important results (Proposition 2 and 3) concerning the equivalence between scheduling algorithms under the CBFS condition.

DEFINITION 1. *We say that two scheduling algorithms are equivalent under a condition \mathcal{C} if given any arrival instance conforming to condition \mathcal{C} , these two algorithms will generate the same packet service schedule.*

NOTATION 1. *For the k ’th packet in session i , let $L_{i,k}$, $T_{i,k}$, and $F_{i,k}$ denote its length, arrival time, and GPS virtual finish time, respectively. Let $V(t)$ denote the GPS virtual time as a function of real time t .*

PROPOSITION 2. Under the CBFS condition,

- (a) $F_{i,k} = F_{i,k-1} + \frac{L_{i,k}}{r_i}$, $1 \leq i \leq n$ and $k > 0$. Here we let $F_{i,0} = 0$ by definition.
- (b) $V(t) \equiv t$

PROOF. (a) In GPS,

$$F_{i,k} = \max(F_{i,k-1}, T_{i,k}) + \frac{L_{i,k}}{r_i} \quad (3)$$

It is clear that $T_{i,k} \leq F_{i,k-1}$, since otherwise, during the time period $[F_{i,k-1}, T_{i,k}]$ the session i is idle under GPS, violating the *continuously backlogged* (Proposition 1) property of CBFS. So the formula (3) becomes $F_{i,k} = F_{i,k-1} + \frac{L_{i,k}}{r_i}$.

(b) Recall that $V(t)$ is defined as follows:

$$V(0) = 0 \quad (4)$$

$$V(t + \tau) = V(t) + r\tau / \left(\sum_{i \in B(t)} r_i \right) \quad (5)$$

where $B(t)$ is the set of sessions that are active during the interval $[t, t + \tau]$. Here $t + \tau$ can be anytime before the occurrence of the first *event* (the arrival or departure of a packet) after t . Since all sessions are backlogged all the time under GPS (Proposition 1), $B(t)$ is exactly the set of all sessions. Therefore, $\sum_{i \in B(t)} r_i = r$ and consequently (5) becomes $V(t + \tau) = V(t) + \tau$. This, combined with (4), implies that $V(t) \equiv t$. \square

Remark: It is clear from (a) that the calculation of GPS virtual finish time is an $O(1)$ operation (under the CBFS condition) per packet, as the program can store the result of $F_{i,k-1}$ from the prior computation.

PROPOSITION 3. Under the CBFS condition, Virtual Clock (VC) [20] and Frame-based Fair Queuing (FFQ) [14] are equivalent to WFQ, and WF^2Q+ [3] is equivalent to WF^2Q .

PROOF. Note that a scheduling algorithm is determined by the following two components: (a) the calculation of the *estimated* virtual finish time of a packet and (b) the policy in selecting the next packet for service. Our first step is to show that, under the CBFS condition, $\tilde{F}_{i,k}$, the estimated virtual finish time, agrees with $F_{i,k}$, the actual one, in VC, FFQ, and WF^2Q+ (i.e., the equivalence of (a) part). The equivalence of (b) part is shown in the second step.

To show $F_{i,k} \equiv \tilde{F}_{i,k}$, it suffices to show that (I) $\tilde{F}_{i,k} = \tilde{F}_{i,k-1} + \frac{L_{i,k}}{r_i}$. This is because $F_{i,k} = F_{i,k-1} + \frac{L_{i,k}}{r_i}$ due to Proposition 2a. Here $\tilde{F}_{i,0} = 0$ by definition. Let $\tilde{V}(t)$ be the estimation of the virtual time as a function of real time t . It also suffices to show (II) $\tilde{V}(t) \equiv t$, since $V(t) \equiv t$ under CBFS (Proposition 2b), and $V(t) \equiv \tilde{V}(t)$ implies (I). In the following, we show in all three algorithms, either the assertion (I) or (II) holds:

- VC: $\tilde{V}(t) \equiv t$ (aforementioned assertion (II)) by definition [20].
- FFQ: In FFQ, the approximation of GPS virtual time is based on a concept called *frame*, which represents the maximum amount of service any flow may receive

during a frame period. A frame period ends when all traffic belonging to the frame has been serviced and the next frame period immediately starts. This maximum can be reached if and only if the flow is continuously backlogged during the frame period. All packets within a frame period will be served based on the increasing order of their *estimated* GPS virtual finish time, and packets belonging to the future frames will be served only after all packets within the current frame finish service. In the following, the virtual time estimation function in FFQ is denoted as P (instead of $\tilde{V}(t)$), following the notations used in [14]. There are two program statements in FFQ that will change the value of P :

(a) $P \leftarrow P + \text{length}(j)/f$ (line 1 in Fig. 3 of [14])

(b) $P \leftarrow \max(\text{Frame}, P)$ (line 9 in Fig. 3 of [14])

It can be shown that under the CBFS condition, $\text{Frame} \leq P$ as always, using similar arguments as used in the WF^2Q+ proof below. The rigorous proof of this requires a detailed description of the algorithm and involved invariant-based induction steps, which is omitted here. So statement (b) above never changes the value of P . Therefore, all the changes to P come from (a), which is equivalent to the aforementioned assertion (I).

- WF^2Q+ : We would like to show the aforementioned assertion (II) $\tilde{V}(t) \equiv t$. We prove this by induction on the packet arrival and departure events, since the virtual time estimation is triggered only by and for these events. Note that $\tilde{V}(t) \equiv t$ at the time the 0'th event happens (i.e., $t=0$). Suppose at (real) time t^* when the j 'th event ($j \geq 0$) happens, $\tilde{V}(t^*) = t^*$. We need to show that $\tilde{V}(t^* + \tau) = t^* + \tau$, in which $t^* + \tau$ is the time when the $(j+1)$ 'th event happens. In WF^2Q+ , $\tilde{V}(t^* + \tau) = \max(\tilde{V}(t^*) + \tau, \min_{i \in B(t^*)} S_i(t^*))$, where $B(t^*)$ is the set of backlogged sessions during the interval $[t^*, t^* + \tau]$ and $S_i(t^*)$ is the virtual start time of (backlogged) session i 's head-of-line (HOL) packet at time t^* . It suffices to show (a) $\tilde{V}(t^*) + \tau \geq \min_{i \in B(t^*)} S_i(t^*)$, since this implies $\tilde{V}(t^* + \tau) = \tilde{V}(t^*) + \tau$, which combined with the induction hypothesis, implies $\tilde{V}(t^* + \tau) = t^* + \tau$. We prove (a) by contradiction. Suppose (b) $t^* + \tau < \min_{i \in B(t^*)} S_i(t^*)$ holds. Due to the CBFS condition, each session i is continuously receiving service up to the virtual time $S_i(t^*)$. Since $S_i(t^*) > t^* + \tau$ for any i according to (b), this implies that the amount of service rendered during the period $[0, t^* + \tau]$ is $\sum_{i=1}^n S_i(t^*) > \sum_{i=1}^n (t^* + \tau) r_i = (t^* + \tau)r$. However, $(t^* + \tau)r$ is the maximum possible amount of service rendered within $t^* + \tau$ seconds.

Finally, note that the policy of selecting next packet for service is the same in FFQ and VC as in WFQ: choosing the one with the smallest estimated virtual finish time. Also, such policy is the same in WF^2Q+ as in WF^2Q , since both select the packet that has the lowest timestamp among those that should have started service in GPS. \square

2.2 Decision tree models

We adopt a standard and commonly-used computational model in proving lower bounds: the *decision tree*. A decision tree program in general takes as input a list of real

variables $\{x_i\}_{1 \leq i \leq n}$. Each internal and external (leaf) node of the tree is labeled with a predicate of these inputs. The algorithm starts execution at the root node. In general, when control is centered at any internal node, the predicate labeling that node is evaluated, and the program control is passed to its left or right child when the value is “yes” or “no” respectively. Before the control is switched over, the program is allowed to execute *unlimited number* of sequential operations such as data movements and arithmetic operations. In particular, the program is allowed to store all results (i.e., no constraint on storage space) from prior computations. When program control reaches a leaf node the predicate there is evaluated and its result is considered as the output of the program. The complexity of such an algorithm is defined as the depth of the tree, which is simply the number of predicates that needs to be evaluated in the *worst case*. Fig. 3 shows a simple decision tree with six nodes. Each P_i ($1 \leq i \leq 6$) is a predicate of the inputs.

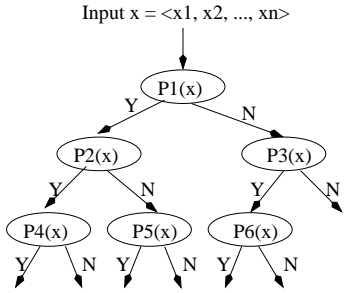


Figure 3: Decision tree computational model

The decision tree was originally proposed for *decision problems*, in which the output is binary: simply “yes” or “no”. The model can be extended to handling more general problems the output of which is not necessarily binary. For example, in the context of this work, the output will be the sequence in which packets get scheduled.

Allowing different types of predicates to be used in the decision tree results in models of different computational powers. On one extreme, if the decision tree program allows the magic predicate $P(x_1, x_2, \dots, x_n)$ that exactly solves the problem, then the complexity of the problem is precisely 1. On the other extreme, if the decision tree program only allows constant predicates, then nontrivial (nonconstant) decision problems are simply not solvable under this model, no matter how much computation is performed. In this work, we consider predicates that are reasonable in the sense that existing scheduling algorithms are able to provide $O(1)$ GPS-delay bounds using only such predicates.

The first computational model we consider is the decision tree that allows *linear tests* [5]. In this model, each predicate allowed by the decision tree is in the form of “ $h(x_1, x_2, \dots, x_n) \geq 0$ ”, where h is a linear function (defined below) of the inputs $\{x_i\}_{1 \leq i \leq n}$.

DEFINITION 2 (LINEAR FUNCTION). A linear function f of the variables $\{x_i\}_{1 \leq i \leq n}$ is defined as $f(x_1, x_2, \dots, x_n) = \sum_{i=1}^n a_i x_i + a_0$, where $\{a_i\}_{0 \leq i \leq n}$ are real constants.

This model will be used in our proofs in Section 4. In the context of this work, the inputs will be the lengths and the

arrival times of the packets. Note that the linear comparison model is quite generous: functions like f in the above definition may take up to $O(n)$ steps to compute, but the model charge only “1” for it. However, one may still argue that linear function can be restrictive for packet scheduling since it does not offer an efficient way to calculate GPS virtual finish times from the inputs. Note that GPS virtual finish time is *in general* not a linear function (actually piece-wise linear) of the inputs. Fortunately, as we proved in Proposition 2a, under the CBFS condition, the GPS virtual finish time of any packet is indeed a linear function of these inputs! So under the CBFS condition, this model is not restrictive.

Under the CBFS condition, this model is reasonable also in the sense that many existing scheduling algorithms, including WFQ , VC , FFQ , WF^2Q , and WF^2Q+ , use only the operations allowed in the model. Due to Proposition 3, under the CBFS condition, we only need to consider WFQ and WF^2Q . Note that in both WFQ and WF^2Q , (1) GPS time estimation is an $O(1)$ operation and does not require branching statements under the CBFS condition (see remark after Proposition 2), and (2) comparisons between virtual finish times (shown to be the linear functions of the inputs) are all that is needed in making scheduling decisions. Careful readers would point out that WF^2Q also involves comparisons with virtual start times. However, note that under the CBFS condition, the virtual start time of a packet is exactly the virtual finish time of the previous packet in the same session!

The second computational model we introduce is the decision tree that allows comparisons only between its inputs. It has been used in proving the $\Omega(n \log_2 n)$ lower bound for comparison-based sorting algorithms [1]. It is strictly weaker than the previous model since the set of predicates that are allowed in this model is a proper subset of what is allowed in the previous model. However, for the particular class of instances that are used in establishing our lower bounds, the second computational model is also reasonable. We will show that under the CBFS condition, allowing comparisons among inputs is equivalent to allowing comparisons between GPS virtual finish times of the packets in that instance class. Since both WFQ and WF^2Q are able to provide $O(1)$ GPS-relative delay bounds using such comparisons only, this model is not restrictive either. In summary, both computational models are practical and nonrestrictive, in the sense that they are actually being used by existing scheduling algorithms.

2.3 Remarks on the decision tree model

A decision tree program allowing certain branching predicates is computationally stronger than a computer program that allows the same types of branching predicates and is memory-constrained. This is because (1) the decision tree can be totally different when the size of input changes, and (2) the computational complexity counted in the decision tree model is only the *depth* of the tree, not the *size* of the tree. Neither is true about a regular computer program. So the lower bound derived under the decision tree model can be no larger than the lower bound achievable by a computer program! For example, Knapsack³, a well-known

³Among a set $T = \{x_1, x_2, \dots, x_n\}$ of n real numbers, decide whether there exists $S \subseteq T$ such that $\sum_{x \in S} x = 1$.

NP-complete problem, has an $O(n^5 \log^2 n)$ algorithm⁴ in the decision tree model that allows linear comparisons [9]. Despite the fact that a decision tree algorithm can be computationally stronger than a computer program, when allowing the same branching predicates, many lower bound proofs are based on decision tree. This is because (1) they provide powerful tools for proving lower bounds, and (2) so far there is no model that exactly captures the computational power of a computer program and at the same time provides such powerful tools.

3. COMPLEXITY-DELAY TRADEOFFS WHEN ALLOWING COMPARISONS BETWEEN INPUTS

In this section, we prove that if only comparisons between inputs are allowed, the complexity to assure $O(1)$ or $O(n^a)$ ($0 < a < 1$) GPS-relative delay bound is $\Omega(\log_2 n)$. In Section 3.1, we introduce two general lemmas used in later proofs. Section 3.2 and 3.3 proves the $\Omega(\log_2 n)$ complexity lower bounds for the case of $O(1)$ and $O(n^a)$ respectively.

3.1 Preliminaries

A reduction argument similar to those used in NP completeness proofs is used in proving the complexity lower bounds throughout this paper. The basic idea is to convert a problem P_1 , the complexity lower bound of which is known (say B), to the problem P_2 , for which we would like to determine the complexity lower bound. If it can be shown that the conversion cost is no more than C , then the complexity lower bound of problem P_2 is at least $B - C$. In this paper, sorting problems with known complexity lower bound of $n \log_2 n - o(n \log_2 n)$ (for instance of size n), are reduced to scheduling problems with $O(n^a)$ ($0 \leq a < 1$) delay or disadvantage guarantees. We show that the reduction cost is no more than $an \log_2 n$. Therefore, the scheduling complexity has to be at least $(1 - a)n \log_2 n - o(n \log_2 n)$ since otherwise the resulting sorting algorithm beats its proven complexity lower bound, which is impossible. Here we write $n \log_2 n - o(n \log_2 n)$ instead of $\Omega(n \log_2 n)$ in order to emphasize that the constant factor of the main term is precisely 1.

In the following, we state without proof the well-known complexity lower bound for comparison-based sorting [1]. Its proof can be found in several algorithm textbooks, including [1]. It is clear from the proof that this lower bound holds even if all the real numbers are between two numbers m and M ($0 \leq m < M$).

LEMMA 1 (SORTING LOWER BOUND [1]). *To sort a set of n numbers $\{x_i\}_{1 \leq i \leq n}$ using only comparisons among them, requires $n \log_2 n - o(n \log_2 n)$ steps in the worst case.*

Reduction to the sorting problem is sufficient for proving the lower bounds (when allowing direct comparisons among inputs) for scheduling throughout this section. However, to prove stronger results (when allowing linear tests) in Section 4, we need to reduce them to a stronger version of the sorting lower bound (Lemma 3). Since the reduction steps for proving stronger lower bounds in Section 4 can be “reused”

⁴This, however, does not imply $P = NP$, since a decision tree algorithm can be more powerful than a computer program.

for proving the weaker results in this section, for the overall succinctness of the proofs, reductions in this section will also be based on Lemma 3 (stronger version) instead of Lemma 1 (weaker version).

DEFINITION 3. *The set membership problem is to determine whether the input $\{x_i\}_{1 \leq i \leq n}$, viewed as a point (x_1, x_2, \dots, x_n) in the Euclidean space R^n , belongs to a set $L \subseteq R^n$.*

In the following, we state a general lemma concerning complexity of set membership problems (defined above) under the decision tree model that allows linear tests. This lemma, due to Dobkin and Lipton [5], has been used extensively in lower bound proofs (e.g., [6]). In complexity theory, lower bound for solving a set membership problem is closely related to the geometric properties of the set. The following lemma essentially states that if the set consists of N disconnected open sets, determining its membership requires at least $\log_2 N$ complexity.

LEMMA 2. *Any linear search tree that solves the membership problem for a disjoint union of a family $\{A_i\}_{i \in I}$ of open subsets of R^n requires at least $\log_2 |I|$ queries in the worst case [5].*

PROOF (ADAPTED FROM [5]). Consider the decision tree algorithm for deciding membership in a set $L \subseteq R^n$. At any leaf node, the algorithm must answer “yes” or “no” to the questions of whether the inputs x_1, x_2, \dots, x_n are coordinates of a point in L . Let the set of points “accepted” at leaf p be denoted by T_p (i.e., T_p is the set of points for which all tests in the tree have identical outcomes and lead to leaf node p , for which the algorithm answers “yes”). The leaf nodes of the tree partition R^n into disjoint convex regions because all comparisons are between linear functions of the coordinates of the input point, so in particular each of the accepting sets T_p is convex.

We prove the lemma by contradiction. Suppose that the level of the tree is less than $\log_2 |I|$. Then the number of leaf nodes must be strictly less than I . Now since L consisting of $|I|$ disjoint regions, some accepting node T_p must accept points in two regions due to the pigeon-hole principle, say L_α and L_β . Choose any points $P_1 \in T_p \cap L_\alpha$ and $P_2 \in T_p \cap L_\beta$. Note that the linear comparisons (viewed as hyperplanes) dissect R^n into convex polytopes. By the convexity of T_p , every point on the line $P_1 P_2$ is in T_p . So for every such point the algorithm answers “yes”. However, L_α and L_β are disjoint open sets, so the line $P_1 P_2$ contains points not in L . This contradicts the correctness of the membership algorithm. \square

Now we are ready to introduce the aforementioned stronger result, concerning sorting complexity lower bound when allowing linear tests. Let $0 \leq m < M$ be two real numbers. The following Lemma (Lemma 3) essentially states that, when linear tests are allowed, the same sorting complexity lower bound ($n \log_2 n - o(n \log_2 n)$) still holds when these n numbers are evenly distributed in the following n neighborhoods: $\{(m + \frac{i(M-m)}{n+1} - \epsilon, m + \frac{i(M-m)}{n+1} + \epsilon)\}_{1 \leq i \leq n}$ (i.e., there exists a permutation π of n elements such that $m + \frac{i(M-m)}{n+1} - \epsilon < x_{\pi(i)} < m + \frac{i(M-m)}{n+1} + \epsilon$, $i = 1, 2, \dots, n$). To see this, we show that this sorting problem is at least asymptotically “as hard as” the membership problem for

the following set L : $L = \{(y_1, y_2, \dots, y_n) \in R^n : \text{there exists a permutation } \pi \text{ of } 1, \dots, n \text{ such that } m + \frac{i(M-m)}{n+1} - \delta < y_{\pi(i)} < m + \frac{i(M-m)}{n+1} + \delta, i = 1, 2, \dots, n\}$. Here $0 < \delta < \frac{M-m}{3(n+1)}$ is a “small” real constant. “Sorting” is at least asymptotically “as hard”, since if there is an algorithm for sorting with computational complexity B , then there is a $B + O(n)$ algorithm for the membership problem (just sort the numbers using B time and check using $O(n)$ time if they are in the corresponding neighborhoods).

LEMMA 3. *Under the decision tree model that allows linear tests, given the inputs $\{x_i\}_{1 \leq i \leq n}$, determining whether $(x_1, x_2, \dots, x_n) \in L$ requires at least $n \log_2 n - o(n \log_2 n)$ linear tests.*

Note that this result is stronger than Lemma 1 since here the computational model (allowing linear tests) is stronger and there are restrictions on the values that these n numbers can take.

PROOF. Let Π be the set of permutations on the set $\{1, 2, \dots, n\}$. Then by the definition of L , $L = \bigcup_{\pi \in \Pi} L_\pi$. Here $L_\pi = \{(y_1, y_2, \dots, y_n) : m + \frac{i(M-m)}{n+1} - \delta < y_{\pi(i)} < m + \frac{i(M-m)}{n+1} + \delta, i = 1, 2, \dots, n\}$. Each L_π is obviously an open set. Also L_{π_1} and L_{π_2} are disjoint if $\pi_1 \neq \pi_2$. To see this, note that if $\pi_1(i) \neq \pi_2(i)$ for some i , then each point in L_{π_1} and each point in L_{π_2} must have a minimum distance of δ between their i 'th coordinates.

The number of such regions $\{L_\pi\}_{\pi \in \Pi}$ is $n!$ because $|\Pi| = n!$. So by Lemma 2, the number of comparisons must be at least $\log_2(n!)$, which by Stirling's formula ($n! \sim \sqrt{2\pi n} (\frac{n}{e})^n$), is equal to $n \log_2 n - o(n \log_2 n)$. \square

Remark: We emphasize that the floor (and equivalently the ceiling) function is not allowed in the decision tree. Otherwise, an $O(n)$ algorithm obviously exists for deciding L -membership based on bucket sorting. Note that the floor function is not a linear function (piecewise linear instead). The linearity of the test is very important in the proof of Lemma 2 since it relies on the fact that the linear tests dissect the space R^n into *convex* regions (polytopes). These regions are no longer convex when the floor function is allowed. For this reason, the floor function⁵ is disallowed in almost all lower bound proofs. Nevertheless, despite the fact that the floor function will “spoil” our lower bound proofs (and many other proofs), no existing scheduling algorithm (certainly allowed to use “floor”) is known to have a *worst case* computational complexity of $o(n \log_2 n)$ and guarantee $O(1)$ or $O(n^a)$ ($0 < a < 1$) *worst-case* GPS-relative delay. Studying the computation power of “floor” on this scheduling problem can be a topic for future research.

3.2 $\Omega(\log_2 n)$ complexity for $O(1)$ delay

In this section, we prove that $\Omega(\log_2 n)$ complexity is required to guarantee $O(1)$ GPS-relative delay, when only comparisons between inputs (equivalently GPS virtual finish times) are allowed. A naive argument for this would be that it takes $\Omega(\log_2 n)$ per packet to schedule the packets according to the sorted order of their GPS virtual finish times. However, this argument is not a proof since it can be shown that to be sorted is not a necessary condition (although it

is sufficient [12]) to assure $O(1)$ GPS-relative delay. For example, if a GPS-relative delay bound of 10 maximum size packets needs to be assured, then given a service schedule sorted according to their GPS virtual finish times, any 9 packets can be relocated (intra-session packet service order should however be preserved) without violating this delay bound.

Before stating the lower bounds and their proofs, we would like to explain the intuition behind them. The core idea is to reduce the problem of scheduling with $O(1)$ delay bound to the problem of sorting. Given any sorting instance, we reduce it to a scheduling instance in $O(n)$ time and run the scheduler with $O(1)$ delay bounds on it. Then we can show that the resulting output can be sorted in $O(n)$ time. Since the sorting complexity is $n \log_2 n - o(n \log_2 n)$, the scheduling complexity has to be at least $n \log_2 n - o(n \log_2 n)$. Otherwise, we have an algorithm that asymptotically beats the complexity lower bound, which is impossible. The proof is split into two parts. In the first part (Theorem 1), we explain the reduction algorithm and establish the complexity equations. In the second part (Theorem 2), we show that this reduction program is correct in the sense that the resulting program (output of the reduction process) indeed performs sorting correctly. This is proved using standard invariant-based techniques for establishing program correctness, and an assertion that a scheduling program should satisfy (Lemma 4), when comparisons are only allowed among inputs.

In proving the following theorem, we assume that there is a $O(1)$ -Delay-Scheduler procedure which guarantees that the GPS-relative delay of any packet will not exceed $K \frac{L_{max}}{r}$ (i.e., $O(1)$). Here $K \geq 1$ is a constant integer independent of the number of sessions n and the total link bandwidth r . We also assume that the CBFS condition is satisfied.

THEOREM 1 (COMPLEXITY). *The computational complexity lower bound of the procedure $O(1)$ -Delay-Scheduler is $\Omega(\log_2 n)$ per packet.*

PROOF. Our proof uses the reduction method in computational complexity, similar to those used in NP completeness proofs. We construct a procedure for solving L -membership (defined in the previous section) as follows. Recall that $L = \{(y_1, y_2, \dots, y_n) : \text{there exists a permutation } \pi \text{ of } \{1, 2, \dots, n\} \text{ such that } m + \frac{i(M-m)}{n+1} - \delta < y_{\pi(i)} < m + \frac{i(M-m)}{n+1} + \delta, i = 1, 2, \dots, n\}$, where $0 < \delta < \frac{M-m}{3(n+1)}$. Here we let $m = L_{min} = 0$ ($L_{min} = 0$ as mentioned before) and $M = L_{max}$, where L_{max} and L_{min} are the maximum and minimum packet sizes respectively. We proved in Lemma 3 that the number of linear tests that are needed in determining L -membership is $n \log_2 n - o(n \log_2 n)$. Now, given the inputs $\{x_i\}_{1 \leq i \leq n}$ to the L -membership problem, we convert it to an instance of packet arrivals. We then feed the packet arrival instance to the procedure $O(1)$ -Delay-Scheduler. Finally, we process the output from the procedure to solve the L -membership problem. Since the total number of comparisons for solving L -membership are $n \log_2 n - o(n \log_2 n)$ in the worst case, a simple counting argument allows us to show that $O(1)$ -Delay-Scheduler must use $\Omega(n \log_2 n)$ comparisons in the worst case. This reduction is illustrated in Fig. 4.

The procedure in Fig. 4 is divided into three parts. In the first part (line 5 through 20), the program first checks if

⁵Its computational power is discussed in [18] in detail.

```

1. Procedure L-Membership I
2. input:  $x_1, x_2, \dots, x_n$ 
3. output: ‘‘yes’’ if  $(x_1, x_2, \dots, x_n) \in L$  and ‘‘no’’ otherwise
4. begin
5.   /* Part I: Create a packet arrival instance and feed it to scheduler */
6.   if  $0 < x_i < L_{max}$  for  $1 \leq i \leq n$  then proceed
7.   else answer ‘‘no’’ endif
8.   for  $i=1$  to  $n$  begin
9.     create (first) packet arrival  $A_{i,1}$  to session  $i$  of length  $x_i$  at time 0
10.    create (second) packet arrival  $A_{i,2}$  to session  $i$  of length  $L_{max} - x_i$  at time 0
11.  end /* for */
12.  call Procedure  $O(1)$ -Delay-Scheduler with
13.    input: arrival instance  $A = \{A_{i,j}\}_{1 \leq i \leq n, 1 \leq j \leq 2}$ 
14.    output: sorted schedule  $S = \{S_i\}_{1 \leq i \leq 2n}$  with  $O(1)$  delay guarantee
15.   $j:=1$ 
16.  for  $i=1$  to  $2n$  begin
17.    if  $S_i$  is the first packet of a session then  $T[j] = S_i$  endif
18.     $j:=j+1$ 
19.    /*  $T$  will only have  $n$  elements: first packets of the  $n$  sessions */
20.  end /* for */

21.  /* Part II: ‘‘sort’’ the output schedule from the scheduler */
22.  for  $i:= 2$  to  $K+2$  begin
23.    perform binary insertion of  $T_i$  into the list  $T_1, T_2, \dots, T_{i-1}$  according to their lengths
24.    /* sort the first  $K+2$  packets using binary insertion according to their lengths */
25.  end /* for */
26.  for  $i:= K+3$  to  $n$ 
27.    perform binary insertion of  $T_i$  into the list  $T_{i-K-2}, T_{i-K-1}, \dots, T_{i-1}$  according to lengths
28.    /* binary insertion into a ‘‘window’’ of size  $K+2$  */
29.  end {for}

30.  /* Part III: check if the ‘‘sorted’’ list, viewed as a point in  $R^n$ , is in  $L$  */
31.  if  $\frac{iL_{max}}{n+1} - \delta < \text{length}(T_i) < \frac{iL_{max}}{n+1} + \delta$  for  $i = 1, 2, \dots, n$  then answer ‘‘yes’’
32.  else answer ‘‘no’’ endif
33. end /* procedure */

```

Figure 4: Algorithm I for L-Membership Test.

all the inputs are in the legitimate range $(0, L_{max})$. It then generates two packets for each session i that arrive at time 0. The first and second packets of session i are of length x_i and $L_{max} - x_i$, respectively. Clearly, between time 0 and $\frac{nL_{max}}{r}$, the CBFS condition holds. The arrival instance is fed as input to the procedure $O(1)$ -delay-scheduler that guarantees a delay bound of $K \frac{L_{max}}{r}$. The output is the schedule of these $2n$ packet by the scheduling procedure. Then the second packet of each session is removed from the schedule (line 16 through 20). In the second part (lines 21 through 29), these packets are sorted according to their lengths, if $(x_1, x_2, \dots, x_n) \in L$ and the procedure $O(1)$ -Delay-Scheduler indeed guarantees $O(1)$ GPS-relative delay. In the third part (line 30 through 32), the processed (sorted) sequence is checked to see if it is indeed in L .

Recall that the procedure $O(1)$ -Delay-Scheduler is allowed to perform comparisons between its inputs, which are arrival times (0) and lengths of the packets. In addition, the constant L_{max} is allowed to be compared with any input⁶. Note that this is equivalent to allowing comparisons between GPS

virtual finish times of the packets, which are in the form of either $\frac{nx_i}{r}$ (first packet of session i), $i = 1, 2, \dots, n$, or $\frac{nL_{max}}{r}$ (second packets of all sessions). Both are linear functions of the inputs which can be used in L -membership without compromising its $n \log_2 n - o(n \log_2 n)$ complexity lower bound (by Lemma 3). Now it is straightforward to verify that excluding the procedure $O(1)$ -Delay-Scheduler, a total of $O(n)$ linear comparisons/tests are performed throughout the L -membership procedure. They include (a) comparisons in line 17 between the GPS virtual finish time of T_i and $\frac{nL_{max}}{r}$, (b) comparisons between GPS virtual finish times of packets from line 21 through 29, and (c) comparisons in line 31 to check if the (sorted) input is in L . So the number of comparisons used in the procedure $O(1)$ -Delay-Scheduler must be $\Omega(n \log_2 n)$. Otherwise, L -membership uses only $o(n \log_2 n)$ comparisons, which contradicts Lemma 3. Therefore, the **amortized complexity per packet** is $\Omega(\log_2 n)$.

We have yet to prove the correctness of the L -membership procedure, i.e., it solves the L -membership correctly for any inputs. This is shown next in Theorem 2. \square

⁶We can artificially create a dummy session which has a packet arrival of length L_{max} at time 0.

Remark: An interesting question is whether $\Omega(\log_2 n)$ per packet cost is paid only at the very beginning (e.g., first

2n packets) and the cost per packet will be amortized to $o(\log_2 n)$ or even $O(1)$ over the long run. The answer is that this cost cannot be amortized. Its proof is omitted here due to space limitations and will be included in the later version of the paper.

THEOREM 2 (CORRECTNESS). *The procedure in Fig. 4 will return yes if and only if $(x_1, x_2, \dots, x_n) \in L$.*

PROOF. The “only if” part is straightforward since line 30 through 32 (validity check) will definitely answer “no” if $(x_1, x_2, \dots, x_n) \notin L$. We only need to prove the “if” part.

Note that after the execution of line 20, $\{\text{length}(T_i)\}_{1 \leq i \leq n}$ is a permutation of the inputs $\{x_i\}_{1 \leq i \leq n}$. Right after the execution of line 25, the lengths of T_1, T_2, \dots, T_{K+2} are in increasing order. We prove by induction that the lengths of all packets are sorted in increasing order after the execution of the loop from line 26 to 29. We refer to the iterations in the loop as $I_{K+3}, I_{K+4}, \dots, I_n$, indexed by the value of i in each iteration. We prove that the first i numbers are sorted after iteration i , $i = K+3, \dots, n$. This is obviously true for $i = K+3$. Suppose it is true for $i = q \geq K+3$. We prove that it is also true for $i = q+1$.

We claim that, right after the execution of line 20, in the schedule $\{T_i\}_{1 \leq i \leq n}$, for $K+3 \leq i \leq n$, there can be no more than $K+2$ elements among T_1, T_2, \dots, T_{i-1} that are longer than T_i . This is proved below in Lemma 4. Then since the lengths of T_1, T_2, \dots, T_q are sorted in increasing order after iteration q by the induction hypothesis, we know that $\text{length}(T_{q-K-2}) \leq \text{length}(T_{q+1})$. Otherwise, there are at least $K+3$ packets ($T_{q-K-2}, T_{q-K-1}, \dots, T_q$) that are longer than T_{q+1} . So for correct binary insertion, the program only needs to search between the index $q-k-2$ and q , as the program does in line 27. So the length of the first $q+1$ packets remain sorted after the insertion: the $i = q+1$ case proved. Finally, note that line 31 correctly checks for L -membership if the numbers $\{\text{length}(T_i)\}_{1 \leq i \leq n}$ are sorted in increasing order. \square

The following lemma states that no packet will be scheduled behind more than $K+2$ packets that have larger GPS virtual finish time. The intuition behind its proof is the following. Suppose a packet P is scheduled behind $K+2$ packets that have larger timestamps. We convert the current packet arrival instance into another instance in which (a) all timestamps that are no larger than P 's (including P itself) are changed to small positive numbers that are close to 0 and all timestamps that are larger are changed to large numbers that are close to L_{max} , and (b) the order of any pair of timestamps remain unchanged. The condition (b) guarantees that the resulting schedule will be the same if only direct comparisons among inputs are allowed. However, P is scheduled behind that $K+1$ packets under the new service schedule, which can be shown to violate the $O(1)$ GPS-relative delay guarantee for P .

LEMMA 4. *Suppose that $(x_1, x_2, \dots, x_n) \in L$. Then for any i , $1 \leq i \leq n$, there can be no more than $K+2$ packets among T_1, T_2, \dots, T_{i-1} that are longer than T_i , in the scheduler output right after the execution of line 20 in Fig. 4.*

PROOF. Note that $\text{length}(T_k) \neq \text{length}(T_l)$ when $k \neq l$, since $(x_1, x_2, \dots, x_n) \in L$. So there exists a unique permutation π of $\{1, 2, \dots, n\}$, such that $\text{length}(T_{\pi(1)}) < \text{length}(T_{\pi(2)}) < \dots < \text{length}(T_{\pi(n)})$. We prove the lemma by contradiction. For any $i > K+3$, suppose there are more than $K+2$ packets that are scheduled before T_i and are longer than T_i . Suppose $\pi(j) = i$, i.e., T_i is the j 'th smallest packet among $\{T_k\}_{1 \leq k \leq n}$. We argue that $i \leq j + K + 2$. In other words, T_i should not be displaced backward by more than $K+2$ positions. To see this, we generate two arbitrary sets of real numbers $\{\alpha_k\}_{1 \leq k \leq n}$ and $\{\beta_k\}_{1 \leq k \leq n}$, where $0 < \alpha_1 < \alpha_2 < \dots < \alpha_n < \delta$ and $0 < \beta_n < \beta_{n-1} < \dots < \beta_1 < \delta$. Here $\delta < \frac{L_{max}}{3(n+1)}$ as before. We consider what happens if we modify the inputs $\{x_k\}_{1 \leq k \leq n}$ to the L -membership in the following way: x_k is changed to $\alpha_{\pi(k)}$ if $x_k \leq x_i$ and is changed to $L_{max} - \beta_{\pi(k)}$ if $x_k > x_i$. It is not hard to verify that the relative order of any two numbers x_l and x_m is the same after the change. Note that the procedure $O(1)$ -Delay-Scheduler is only allowed to compare among the inputs, which are $\{x_i\}_{1 \leq i \leq n}$, 0, and L_{max} . Clearly, with the modified inputs, the decision tree of the procedure $O(1)$ -Delay-Scheduler will follow the same path from the root to the leaf as with the original inputs, since all predicates along the path are evaluated to the same value! Consequently, the output schedule of the packets remain the same with the modified inputs. In the new schedule with the modified inputs, since there are $K+2$ packets that are scheduled before T_i and are longer than $L_{max} - \delta$, the actual finish time of T_i is larger than $(K+2) \frac{L_{max} - \delta}{r} > (K+1) \frac{L_{max}}{r}$. However, its GPS virtual finish time is no larger than $\frac{n\delta}{r} < \frac{L_{max}}{r}$. So the GPS-relative delay of the packet T_i must be larger than $(K+1) \frac{L_{max}}{r} - \frac{L_{max}}{r} = K \frac{L_{max}}{r}$. This violates the assumed property of $O(1)$ -Delay-Scheduler. \square

Remark: The ideas contained in the proof bear some similarity to that of Knuth's 0-1 law [11], which states the following. If a sorting network can correctly sort inputs consisting of any arbitrary combinations of 0's and 1's, it must be able to correctly sort all inputs. In our proof, $\{\alpha_i\}_{1 \leq i \leq n}$ and $\{L_{max} - \beta_i\}_{1 \leq i \leq n}$, to a certain extent, can be viewed as such 0's and 1's.

3.3 $\Omega(\log_2 n)$ complexity for $O(n^a)$ delay

In this section, we prove that the tradeoff curve is flat as shown in Fig. 2: $\Omega(\log_2 n)$ complexity is required even when $O(n^a)$ delay ($0 < a < 1$) can be tolerated. Its reduction proof is mostly similar to that of Theorem 1. The main difference is that the constant factor before the asymptotic term ($n \log_2 n$) becomes critical in this case.

THEOREM 3. *Suppose we have a procedure $O(n^a)$ -Delay-Scheduler that guarantees a GPS-relative delay of no more than $K n^a \frac{L_{max}}{r}$. Here $K \geq 1$ is an integer constant and $0 < a < 1$ is a real constant. Then the complexity lower bound of $O(n^a)$ -Delay-Scheduler is $\Omega(\log_2 n)$ if it is allowed to compare only between any two inputs.*

PROOF (SKETCH). The proof of this theorem is very similar to that of Theorem 1 and 2. We construct a procedure L -membership-II, which makes “oracle calls” to $O(n^a)$ -Delay-Scheduler, shown in Fig. 5. Since it is mostly the same as the program shown in Fig. 4, we display only the lines that are different.

```

1. Procedure L-Membership II
   ..... same as in Fig. 4 .....
2.   call Procedure  $O(n^a)$ -Delay-Scheduler with
   ..... same as in Fig. 4 .....
21.  /* Part II: ‘‘sort’’ the output schedule from the scheduler */
22.  for  $i := 2$  to  $Kn^a + 2$  begin
23.    perform binary insertion of  $T_i$  into the list  $T_1, T_2, \dots, T_{i-1}$  according to their lengths
24.    /* sort the first  $Kn^a + 2$  numbers using binary insertion */
25.  end /* for */
26.  for  $i := Kn^a + 3$  to  $n$  begin
27.    perform binary insertion of  $T_i$  into the list  $T_{i-Kn^a-2}, \dots, T_{i-1}$  according to their lengths
28.    /* binary insertion into a ‘‘window’’ of size  $Kn^a + 2$  */
29.  end {for}
   ..... same as in Fig. 4 .....

```

Figure 5: Algorithm II for L-Membership Test.

Analysis of the complexity is similar to the proof of Theorem 1. The number of comparisons that are used in line 21 through line 29 is no more than $n \log_2(Kn^a + 2) = an \log_2 n + o(n \log_2 n)$. Note that the number of operations performed from line 26 through 29 is actually n^{2a} if the data movements are also counted. However, as we have explained earlier in Section 2.2, we ‘‘charge’’ only for the comparisons. So the number of comparisons used in $O(n^a)$ -Delay-Scheduler must be at least $(1 - a)n \log_2 n - o(n \log_2 n)$ since otherwise L-membership II uses less than $n \log_2 n - o(n \log_2 n)$ comparisons in the worst case. This would contradict Lemma 3.

Proof of correctness for the procedure L-membership II is also similar to that of Theorem 2. We only need to show the following lemma. We omit its proof here since it is similar to that of lemma 4. \square

LEMMA 5. *Suppose that $(x_1, x_2, \dots, x_n) \in L$. Then for any i , $1 \leq i \leq n$, there can be no more than $Kn^a + 2$ packets among T_1, T_2, \dots, T_{i-1} that are longer than T_i , in the scheduler output (right after line 20) in Fig. 5.*

4. COMPLEXITY-DELAY TRADEOFFS WHEN ALLOWING LINEAR TESTS

In the previous section, we have established the lower bound of $\Omega(\log_2 n)$ for guaranteeing $O(n^a)$ GPS-relative delay for $0 \leq a < 1$. However, the computational model is slightly restrictive: we only allow the comparisons among the inputs (equivalently the GPS virtual finish times). In this section, we extend the complexity lower bounds to a much stronger computational model, namely, the decision tree that allows comparisons between linear combinations of the inputs. However, to be able to prove the same complexity bounds in the new model, we require that the same $(O(n^a))$ for $0 \leq a < 1$ delay bounds are achieved for a different and stronger type of delay called *disadvantage delay*. Despite this restriction, the overall result is provably stronger (by Theorem 5) than results (Theorems 1 and 3) in the last section. Whether the same complexity lower bound holds when linear comparisons are allowed and $O(1)$ or $O(n^a)$ GPS-relative delay bound needs to be guaranteed is left as an open problem for further exploration.

With respect to a service schedule of packets $T = T_1, T_2, \dots, T_n$, we define *disadvantage* of a packet T_i (denoted as

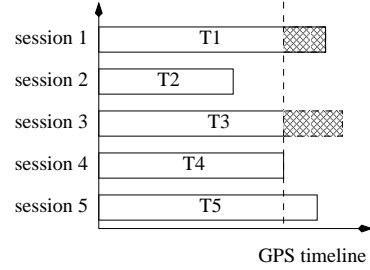


Figure 6: Disadvantage of packet T_4 (the shaded area)

$disadv(T_i)$) as the amount of traffic that has actually been served in the schedule T , which should have been served after the virtual finish time of T_i in GPS. The *disadvantage delay* is defined as *disadvantage* divided by the link rate r .

In Fig. 6, the shaded area adds up to the disadvantage of the packet T_4 when the service schedule is in the order T_1, T_2, \dots, T_5 . Recall that F_p^{GPS} denotes the virtual finish time of the packet p served by GPS scheduler. Formally, the disadvantage of the packet T_i ($i = 1, 2, \dots, n$) is

$$disadv(T_i) = \sum_{j=0}^{i-1} \max(0, F_{T_j}^{GPS} - F_{T_i}^{GPS}) \quad (6)$$

So $disadv(T_i)$ can be viewed as the total amount of ‘‘undue advantage’’ in terms of service other packets have gained over the packet T_i . The following lemma that the *disadvantage delay* of a packet is always no more than its GPS-relative delay.

LEMMA 6. *Under the CBFS condition, for any packet service schedule $T = T_1, T_2, \dots, T_n$, the disadvantage delay of T_i ($1 \leq i \leq n$) is always no larger than its GPS-relative delay.*

PROOF. Let $B(T_i)$ and $A(T_i)$ be the set of bits that should have been served before and after $F_{T_i}^{GPS}$ in a GPS scheduler, respectively. Then, under the CBFS condition, the GPS-relative delay of T_i can be written as $\max(0, \frac{1}{r} \sum_{j=0}^{i-1} \dots)$

$\text{length}(A(T_i) \cap T_j) - \frac{1}{r} \sum_{j=i+1}^n \text{length}(B(T_i) \cap T_j)$. The disadvantage delay of a packet, on the other hand, is $\frac{1}{r} \sum_{j=0}^{i-1} \text{length}(A(S_i) \cap T_j)$. Obviously, the latter is no larger than the former. \square

Remark: The above lemma implies that, for the same amount, guaranteeing disadvantage delay bound is stronger (harder) than guaranteeing GPS-relative delay bound. However, it is only slightly stronger: the disadvantage delay bound of WFQ is zero and that of WF^2Q is also zero if all packets arrive at the same time (so *eligibility test* [2] is no longer an issue).

Now we are ready to state and prove the main theorem of this section: $\Omega(\log_2 n)$ per packet is needed to guarantee a disadvantage delay bound of no more than $O(n^a)$ ($0 \leq a < 1$). The proof is again based on reduction technique. An L-membership (i.e., stronger version of sorting) instance is reduced to a scheduling instance and fed to a scheduler that guarantees $O(n^a)$ ($0 \leq a < 1$) disadvantage bound as input. It can be shown that L-membership test can be performed on the output within $\frac{a+1}{2} n \log_2 n$ time. So the scheduling has to “pay” the difference between the L-membership lower bound $n \log_2 n - o(n \log_2 n)$ and $\frac{a+1}{2} n \log_2 n$, which is $\frac{1-a}{2} n \log_2 n$ (i.e., $\Omega(n \log_2 n)$), or $\Omega(\log_2 n)$ per packet.

In proving the following theorem, we assume that there is a $O(n^a)$ -Disadvantage-Scheduler ($0 \leq a < 1$) that guarantees a disadvantage delay bound of $K n^a \frac{L_{max}}{r}$ (i.e., $O(n^a)$), where $K \geq 1$ is an integer constant.

THEOREM 4. *The number of linear tests used in the procedure $O(n^a)$ -Disadvantage-Scheduler ($0 \leq a < 1$) have a lower bound of $\Omega(n \log_2 n)$ in the worst case.*

PROOF. The framework of the proof is the same as those of theorem 1, 2, and 3. A procedure for L-membership test is shown in Fig. 7. Since it is very similar to the program shown in Fig. 4, we show only the lines that are different. The comparisons used in the procedure include (a) comparisons used in $O(n^a)$ -Disadvantage-Scheduler, (b) no more than $n \log_2(\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil) = \frac{a+1}{2} n \log_2 n + o(n \log_2 n)$ comparisons used in line 21 through 29, and (c) $O(n)$ comparisons used line 16 through 20. Since (a) + (b) + (c) = $n \log_2 n - o(n \log_2 n)$, we know that the (a) part must be at least $(1 - \frac{a+1}{2}) n \log_2 n = \Omega(n \log_2 n)$.

It remains to prove the correctness of L-membership III. Again its proof is quite similar to that of Theorem 2. We claim that the $\{T_i\}_{1 \leq i \leq n}$ are sorted after the execution of line 29. Similar to the proof of theorem 2, it suffices to show that the following lemma holds. \square

LEMMA 7. *Suppose $(x_1, x_2, \dots, x_n) \in L$. Then right before the execution of line 22 in Fig. 7, for any packet T_i , there can be no more than $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$ packets, among T_1, T_2, \dots, T_{i-1} , that are longer than T_i .*

PROOF. We prove by contradiction. Let $\Gamma = \{T_j : 1 \leq j \leq i-1, \text{length}(T_j) > \text{length}(T_i)\}$ and let $N = |\Gamma|$. Since $(x_1, x_2, \dots, x_n) \in L$, we know that each interval $(\frac{jL_{max}}{n+1} - \delta, \frac{jL_{max}}{n+1} + \delta)$, $1 \leq j \leq n$, must contain the length of one and exactly one packet among $\{T_j\}_{1 \leq j \leq n}$. So in the sorted order of their lengths, packets in Γ must be longer than T_i for at least $\frac{L_{max}}{n+1} - 2\delta$, $\frac{2L_{max}}{n+1} - 2\delta$, \dots , and $\frac{NL_{max}}{n+1} - 2\delta$,

respectively. Suppose $N > \lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$. Then $\text{disadv}(T_i) > \frac{1}{r} \sum_{j=1}^N (\frac{jL_{max}}{n+1} - 2\delta) = \frac{1}{r} (\frac{\frac{1}{2}N(N+1)L_{max}}{n+1} - 2N\delta) > K n^a \frac{L_{max}}{r}$. This contradicts the guarantee provided by the procedure $O(n^a)$ -Disadvantage-Scheduler. Therefore, N must be no more than $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$. \square

Compared to Theorem 1 and 3, Theorem 4 allows for a much stronger computational model. However, it has to enforce a slightly stronger type of delay (disadvantage delay) than GPS-relative delay to maintain the same lower bounds. Nevertheless, the overall result of Theorem 4 is provably stronger than that of Theorem 1 and 3, shown next.

THEOREM 5. *If a scheduler assures $O(n^a)$ GPS-relative delay bound using only comparisons between inputs (equivalently GPS virtual finish times), it also necessarily assures $O(n^a)$ disadvantage delay bound.*

PROOF. Proof of Lemma 4 can be adapted to show that among $\{T_j\}_{1 \leq j \leq i-1}$ there can be no more than $Kn^a + 2$ packets that are longer than T_i . So the disadvantage delay of T_i is no more than $(Kn^a + 2) \frac{L_{max}}{r}$, which is $O(n^a)$. \square

We conclude this section by the following conjecture that was the initial goal of this work.

CONJECTURE 1. *The complexity lower bound for a scheduling algorithm to achieve a delay bound of $O(1)$, under the decision tree model that allows linear tests, is $\Omega(\log_2 n)$ per packet. A stronger result would be to generalize it further to the case of $O(n^a)$ ($0 < a < 1$) delay bound.*

5. LINKING GPS-RELATIVE DELAY WITH END-TO-END DELAY

In the previous two sections, we obtain complexity lower bounds for achieving $O(n^a)$ ($0 \leq a < 1$) GPS-relative or disadvantage delay bounds. However, it is more interesting to derive complexity lower bounds for scheduling algorithms that provide end-to-end delay bounds. In this section, we show that our lower bound complexity results can indeed be put into the context of providing tight end-to-end delay bounds. This is done by studying the relationship between the GPS-relative delay and the end-to-end delay.

In [15], Stiliadis and Varma defined a general class of latency rate (LR) schedulers (called *servers* in [15]) capable of describing the worst-case behavior of numerous scheduling algorithms. From the viewpoint of a session i , any LR scheduler is characterized by two parameters: *latency bound* Θ_i and *minimum guaranteed rate* r_i . We further assume that the j 'th busy period of session i starts at time τ . Let $W_{i,j}(\tau, t)$ denote the total service provided to packets in session i that arrive after time τ and until time t by the scheduler. A scheduler S belongs to the class LR if for all times t after time τ and until the packets that arrived during this period are serviced,

$$W_{i,j}(\tau, t) \geq \max(0, r_i(t - \tau - \Theta_i)) \quad (7)$$

It has been shown that, for a large class of LR schedulers (including WFQ [12], FFQ [14], VC [20], WF^2Q [2], WF^2Q+ [3]), the latency bound of session i , denoted as Θ_i , is

$$\Theta_i = \frac{L_{max}}{r} + \frac{L_{max,i}}{r_i} \quad (8)$$

```

1. Procedure L-Membership III
   ..... same as in Fig. 4 .....
12.   call Procedure  $O(n^a)$ -Disadvantage-Scheduler with
   ..... same as in Fig. 4 .....
21.   /* Part II: ‘‘sort’’ the output schedule from the scheduler */
22.   for i:= 2 to  $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$  begin
23.     perform binary insertion of  $T_i$  into the list  $T_1, T_2, \dots, T_{i-1}$  according to their lengths
24.     /* sort the first  $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$  numbers using binary insertion */
25.   end /* for */
26.   for i:=  $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil + 1$  to  $n$  begin
27.     binary insertion of  $T_i$  into the list  $T_{i-\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil}, \dots, T_{i-1}$  according to lengths
28.     /* binary insertion into a ‘‘window’’ of size  $\lceil \sqrt{2(n+1)(Kn^a+1)} \rceil$  */
29.   end {for}
   ..... same as in Fig. 4 .....

```

Figure 7: Algorithm III for L-Membership Test.

Here $L_{max,i}$ is the maximum size of a packet in session i and r_i is the service rate guaranteed to session i . For (8) to hold, it should be true that the link is not over-subscribed, i.e., $\sum_{i=1}^n r_i \leq r$. Note in (8) that the first term in RHS is the GPS-relative delay bound in both WFQ and WF^2Q .

One important property of the latency bound Θ_i , shown in [15] is that it can be viewed as the worst-case delay seen by a session i packet arriving into an empty session i queue. It has been shown in [15] that the latency bound is further connected to the end-to-end delay bound of session i , denoted as D_i^N , by the following inequality:

$$D_i^N \leq \frac{\sigma_i}{r_i} + \sum_{j=1}^N \Theta_i^j \quad (9)$$

Here N is the number of nodes (routers) that traffic in session i traverses and Θ_i^j is the latency bound of session i in j 'th scheduler. Also, traffic in session i is leaky-bucket constrained and σ_i is the size of the leaky bucket. This result is strong and important since different routers on the path may use different LR schedulers, but (9) still holds in this heterogeneous setting.

We show, in the following theorem, that under a special CBFS condition called $CBFS+$, providing tight GPS-relative delay bound is equivalent to providing tight *latency bound* in any LR scheduler. In $CBFS+$, the j 'th packet ($j \geq 2$) in session i arrives just at the time the $(j-1)$ 'th packet finishes service under the GPS scheduler. In other words, each packet in session i arrives just in time to satisfy the CBFS condition. The following theorem is one major step in connecting our complexity results to the complexity of providing tight end-to-end delay bounds.

THEOREM 6. *Under the $CBFS+$ condition, an LR scheduler is able to guarantee a GPS-relative delay bound of B for all packets, if and only if, for all $L_{max,i} > 0$, it guarantees a latency bound of $B + \frac{L_{max,i}}{r_i}$ for session i .*

PROOF. (if part): Given any packet p , let p' be the very first packet of the session i busy period where p is in. Let W be the total amount of session i traffic that arrive between p and p' (p and p' included). Note that if p arrives to see an empty queue i , then p and p' are the same packet. Let

$\epsilon > 0$ be a constant. Suppose p' arrives at time τ . Then according to (7), p must finish service at the time $t = \tau + B + \frac{L_{max,i} + W}{r_i} + \epsilon$ since by time t the server must have accomplished $r_i(t - \tau - B - \frac{L_{max,i}}{r_i} + \epsilon) = W + r_i\epsilon > W$ amount of service. This includes the service of p since W is the total session i traffic between p and p' . Now under the CBFS condition, the GPS virtual finish time of the packet p will be no smaller than $\tau + \frac{W}{r_i}$. So the GPS-relative delay of the packet p will be no larger than $t - \tau - \frac{W}{r_i} = B + \frac{L_{max,i}}{r_i} + \epsilon$. By making $L_{max,i}$ and ϵ arbitrarily close to 0^+ , we get the GPS-relative delay bound of B . Note that we only need the CBFS condition (instead of $CBFS+$) in this part.

(only if part): Suppose that a packet p of size l arrives at time τ to see an empty session i queue. Then under the $CBFS+$ condition, τ is exactly the time that its previous packet finishes service in GPS. So the virtual finish time of p is exactly $\tau + \frac{l}{r_i}$, since its GPS virtual start time is τ . Since its GPS relative delay is no more than B , it must finish service by the time $\tau + \frac{l}{r_i} + B$. So the latency of p is at most $\frac{l}{r_i} + B \leq B + \frac{L_{max,i}}{r_i}$. \square

Remark: We have just shown that, under the CBFS condition, an LR scheduler that provides a tight latency bound also provides a tight GPS-relative delay bound (the ‘‘if’’ part). This in general (without CBFS) is not true: LR schedulers such as FFQ [14], VC [20], and WF^2Q+ [3] all provide a tight latency bound of $\frac{L_{max}}{r} + \frac{L_{max,i}}{r_i}$, but it can be shown that none of them can provide $O(1)$ or even $O(n^a)$ ($0 < a < 1$) GPS-relative delay in the worst case. However, under the CBFS condition, we show in Proposition 3 that FFQ and VC become equivalent to WFQ , and WF^2Q+ becomes equivalent to WF^2Q , and we have shown earlier that WFQ and WF^2Q guarantee $O(1)$ GPS-relative delay bounds.

The following corollary is the **major result of this section**.

COROLLARY 1. *Under the CBFS condition, when only comparisons between inputs are allowed, the per packet computational complexity for an LR scheduler to guarantee a latency bound of $O(n^a \frac{L_{max}}{r} + \frac{L_{max,i}}{r_i})$ ($0 \leq a < 1$) for all $L_{max,i} > 0$ is $\Omega(\log_2 n)$.*

PROOF. Combine Theorem 1, 2, and 3, and the “if part” of Theorem 6. \square

The implications of this corollary are profound. Note that for all schedulers on the path to be *LR* servers with tight delay bounds is a sufficient rather than necessary condition for achieving tight overall end-to-end delay bounds. Therefore, Corollary 1 does not establish in full generality the complexity lower bounds for achieving tight end-to-end delay bounds. However, there is substantial evidence [15] that this is a “fairly tight” sufficient condition, as most existing scheduling algorithms that can collectively achieve tight end-to-end delay bounds are *LR* servers. Corollary 1 essentially states that such lower bounds hold if the end-to-end delay bounds are established through “good” *LR* servers⁷.

Finally, we identify one open problem that we feel very likely to be solvable and its solution can be a very exciting result, stated as Conjecture 2 below. Note that Conjecture 2 is strictly weaker than Conjecture 1, as it can be shown that the latter implies the former.

CONJECTURE 2. *The complexity lower bound for an LR scheduler (introduced in [15]) to achieve a tight latency bound of $O(1) \frac{L_{max}}{r} + \frac{L_{max,i}}{r_i}$ is $\Omega(\log_2 n)$ per packet, under the decision tree model that allows linear tests.*

Remark: *FFQ*, *VC*, and *WF²Q+* all achieve this latency bound at the complexity of $O(\log_2 n)$ per packet, without the restriction of the CBFS condition. If this conjecture is true, it implies that these algorithms are asymptotically optimal for this purpose, which is an exciting result! Note that Corollary 1 proves this complexity lower bound under the weaker model that allows only comparisons among inputs.

6. CONCLUSIONS

In this work, we clarify, extend and solve an open problem concerning the computational complexity for packet scheduling algorithms to achieve tight delay bounds. To the best of our knowledge, this is the first major step in establishing the complexity lower bounds for packet scheduling algorithms. Our three major results can be summarized as follows:

1. We prove that $\Omega(\log_2 n)$ is indeed the per packet complexity lower bound to guarantee $O(1)$ GPS-relative delay (excluding the cost of tracking GPS time), if a scheduling algorithm is only allowed to compare among inputs (equivalently among GPS virtual finish times) in its decision tree. Moreover, we prove that the complexity lower bound remains the same even if the GPS-relative delay bound is relaxed to $O(n^a)$ for $0 < a < 1$, thus establishing the complete “tradeoff curve”.
2. We are able to extend our complexity results to a much stronger computational model: a decision tree that allows linear tests. However, this comes at the cost of having to enforce a slightly stronger type of delay (disadvantage delay) in the same asymptotic amount ($O(n^a)$, $0 \leq a < 1$). Nevertheless, we show that the overall results remain stronger.

⁷One possible way not to use such “good” *LR* servers to establish tight end-to-end delay bounds may be to use dynamic packet state (DPS) (introduced first in SCORE [16]) to convey scheduling state information from one scheduler to another.

3. We show that in a Latency Rate (LR) [15] scheduler, providing a tight GPS-relative delay bound of $\frac{n^a L_{max}}{r}$ ($0 \leq a < 1$) is equivalent to providing a tight latency bound of $\frac{n^a L_{max}}{r} + \frac{L_{max,i}}{r_i}$, under the CBFS+ condition. This allows us to conclude that the per packet complexity to guarantee $\frac{n^a L_{max}}{r} + \frac{L_{max,i}}{r_i}$ latency bound using only comparisons among inputs is $\Omega(\log_2 n)$. This, to a certain extent, connects the complexity lower bounds for guaranteeing tight GPS-relative delay bounds in a single scheduler to the complexity that is needed to guarantee end-to-end delay bounds.

7. ACKNOWLEDGMENTS

We thank Dr. Chuanxiong Guo for helpful discussions with the first author. We thank Prof. George Varghese for encouraging the first author to work on this important problem. We thank Dr. Scott Shenker, the shepherd of this paper, Prof. Ellen Zegura, Prof. Yechezkel Zalcstein, Mr. Shashidhar Merugu and anonymous referees for their insightful comments and suggestions that help improve the quality and accessibility of the paper.

8. REFERENCES

- [1] A. Aho, J. Hopcroft, and J. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1973.
- [2] J. Bennett and H. Zhang. *WF²Q*: worst-case fair weighted fair queuing. In *IEEE INFOCOM'96*, Mar. 1996.
- [3] J. Bennett and H. Zhang. Hierarchical packet fair queuing algorithms. *IEEE/ACM Transactions on Networking*, 5(5):675–689, 1997.
- [4] A. Demers, S. Keshav, and S. Shenker. Analysis and simulation of a fair queuing algorithm. *Internetworking: Research and Experience*, pages 3–26, 1990. Also in Proceedings of ACM SIGCOMM'89.
- [5] D. Dobkin and R. Lipton. A lower bound of $\frac{1}{2}n^2$ on linear search programs for the knapsack problem. *J. Comput. Syst. Sci.*, 16:413–417, 1978.
- [6] M. Fredman and B. Weide. On the complexity of computing of measure of $\bigcup[a_i, b_i]$. *Communications of the ACM*, 21(7), July 1978.
- [7] A. Greenberg and N. Madras. How fair is fair queuing? *Journal of the ACM*, 39(3):568–598, 1992. Also in Proc. of Performance 1990.
- [8] C. Guo. SSR: an $o(1)$ time complexity packet scheduler for flows in multi-service packet networks. In *Proc. of Sigcomm'01*, Sept. 2001.
- [9] F. Heide. A polynomial linear search algorithm for the n -dimensional knapsack problem. *J. of the ACM*, 31:668–676, 1984.
- [10] M. Katevenis, S. Sidiropoulos, and C. Courcoubetis. Weighed round-robin cell multiplexing in general-purpose ATM switch chip. *IEEE Journal on Selected Areas in Communications*, 9:1265–1279, Oct. 1991.
- [11] D. Knuth. *The Art of Computer Programming, Sorting and Searching (Volume 3)*. Addison-Wesley, 1998.

- [12] A. Parekh and R. Gallager. A generalized processor sharing approach to flow control in integrated services networks: the single node case. *IEEE/ACM Transaction on Networking*, 1(3):344–357, June 1993.
- [13] M. Shreedhar and G. Varghese. Efficient fair queuing using deficit round robin. In *Proc. of ACM SIGCOMM'95*, pages 231–242, Aug. 1995.
- [14] D. Stiliadis and A. Varma. Design and analysis of frame-based fair queuing: A new traffic scheduling algorithm for packet switched networks. In *Proc. of ACM Sigmetrics'96*, pages 104–115, May 1996.
- [15] D. Stiliadis and A. Varma. Latency-rate servers: a general model for analysis of traffic scheduling algorithms. In *Proc. of Infocom'96*, Mar. 1996.
- [16] I. Stoica and H. Zhang. Providing guaranteed services without per flow management. In *Proc. of ACM SIGCOMM*, Sept. 1999.
- [17] J. Turner. New directions in communications (or which way to the information age?). *IEEE Communications Magazine*, 24:8–15, Oct. 1986.
- [18] G. Yuval. Finding nearest neighbors. *Information Processing Letters*, 5(3):63–65, Aug. 1976.
- [19] H. Zhang. Service disciplines for guaranteed performance service in packet switching networks. *Proceedings of the IEEE*, 83(10), Oct. 1995.
- [20] L. Zhang. Virtualclock: a new traffic control algorithm for packet switching networks. *ACM Transactions on Computer Systems*, 9:101–124, May 1991.