# Fair Round Robin: A Low Complexity Packet Scheduler with Proportional and Worst-Case Fairness

Xin Yuan, *Member, IEEE* and Zhenhai Duan, *Member, IEEE*

*Abstract*— **Round robin based packet schedulers generally have a low complexity and provide long-term fairness. The main limitation of such schemes is that they do not support short-term fairness. In this paper, we propose a new low complexity round robin scheduler, called Fair Round Robin ($FRR$), that overcomes this limitation. $FRR$ has similar complexity and long-term fairness properties as the stratified round robin scheduler, a recently proposed scheme that arguably provides the best quality-of-service properties among all existing round robin based low complexity packet schedulers. $FRR$ offers better short-term fairness than stratified round robin and other existing round robin schedulers.**

*Index Terms*— **Packet scheduling, proportional fairness, worst-case fairness, round robin scheduler**

## I. Introduction

An ideal packet scheduler should have a *low complexity*, preferably $O(1)$ with respect to the number of flows serviced, while providing *fairness* among the flows. While the definition of the complexity of a packet scheduling algorithm is well understood, the concept of fairness needs further elaboration. Many fairness criteria for packet schedulers have been proposed [10]. In this paper, we will use two well established fairness criteria to evaluate packet schedulers, the *proportional fairness* that was defined by Golestani in [6] and the *worst-case fairness* that was defined by Bennett and Zhang in [2].

Let $S_{i,s}(t_1, t_2)$ be the amount of data of flow $f_i$ sent during time period $[t_1, t_2)$ using scheduler $s$. Let $f_i$ and $f_j$ be any two flows that are backlogged during an arbitrary time period $[t_1, t_2)$. The proportional fairness of scheduler $s$ is measured by the difference between the normalized services received by the two flows, $|\frac{S_{i,s}(t_1,t_2)}{r_i} - \frac{S_{j,s}(t_1,t_2)}{r_j}|$. We will say that *a scheduler has a good proportional fairness property if the difference is bounded by a constant number of packets in each flow*, that is, $|\frac{S_{i,s}(t_1,t_2)}{r_i} - \frac{S_{j,s}(t_1,t_2)}{r_j}| \le c_1 \frac{L_M}{r_i} + c_2 \frac{L_M}{r_j}$, where $c_1$ and $c_2$ are constants and $L_M$ is the maximum packet size. One example scheduler with a good proportional fairness property is the Deficit Round Robin scheduler [17].

A scheduler with a good proportional fairness property guarantees long-term fairness: for any (long) period of time, the services given to any two continuously backlogged flows are roughly proportional to their weights. However, proportional fairness does not imply short-term fairness. Consider for example a scheduler with a good proportional fairness property serving packets from one 1Mbps flow and 1000 1Kbps flows. With a good proportional fairness property, each of the 1000 1Kbps flows can send a

X. Yuan and Z. Duan are with the Department of Computer Science, Florida State University, Tallahassee, FL 32306. E-mail: {xyuan,duan}@cs.fsu.edu.

constant number of packets ahead of the packet that is supposed to be sent by the 1Mbps flow. Hence, during the period when the scheduler sends the few thousand packets from the 1Kbps flows, the 1Mbps flow is under-served: a scheduler with a good proportional fairness property can be short-term unfair to a flow. To better measure the short-term fairness property of a scheduler, worst case fairness is introduced in [2].

A scheduler, $s$, is worst-case fair to flow $f_i$ if and only if the delay of a packet arriving at time $t$ on flow $f_i$ is bounded by $\frac{Q_{i,s}(t)}{r_i} + C_{i,s}$, where $Q_{i,s}(t)$ is the queue size of $f_i$ at $t$, $r_i$ is the guaranteed rate of $f_i$, and $C_{i,s}$ is a constant independent of the queues of other flows. A scheduler is worst-case fair if it is worst-case fair to all flows in the system. If a scheduler, $s$, is worst-case fair, the fairness of the scheduler is measured by the *normalized worst-case fair index* [2]. Let $R$ be the total link bandwidth. The normalized worst-case fair index for the scheduler, $c_s$, is defined as $c_s = \max_i\{\frac{r_i C_{i,s}}{R}\}$. We will say that *a scheduler has a good worst-case fairness property if it has a constant (with respect to the number of flows) normalized worst-case fair index*. One example scheduler with a good worst-case fairness property is the $WF^2Q$ scheduler [1], [2]. A scheduler with **both** a good proportional fairness property and a good worst-case fairness property provides both long-term and short-term fairness to all flows.

Packet schedulers can be broadly classified into two types: timestamp based schemes [1], [2], [5], [6], [12] and round-robin algorithms [7], [8], [13], [17]. Timestamp based schemes have good fairness properties with a relatively high complexity, $O(log\ N)$, where $N$ is the number of flows. Round-robin based algorithms have an $O(1)$ or quasi-$O(1)$ ($O(1)$ under practical assumptions [13]) complexity, but in general do not have good fairness properties. Round robin schemes including Deficit Round Robin (DRR) [17], Smoothed Round Robin (SRR) [7], and STratified Round Robin (STRR) [13] all have good proportional fairness properties. However, none of the existing round-robin schemes is known to have a normalized worst-case fair index that is less than $O(N)$. We will give an example in Section III-A, showing that the normalized worst-case fair index of $STRR$ [13], a recently proposed round-robin based algorithm that arguably provides the best quality-of-service properties among all existing round-robin based schemes, is $\Omega(N)$. It can be shown that the normalized worst-case fair indexes of other round robin schemes, such as Smoothed Round Robin (SRR) [7] and Deficit Round Robin (DRR) [17], are also $\Omega(N)$. Not having a constant worst-case fair index means that the short-term service rate of a flow may significantly deviate from its fair rate, which can cause rate oscillation for a flow [2].

We propose Fair Round-Robin ($FRR$), a round robin based

low complexity packet scheduling scheme that overcomes the limitation of not being able to guarantee short-term fairness. Like $STRR$, $FRR$ employs a two-level scheduling structure and combines the ideas in timestamp based and round-robin schemes. $FRR$ has a similar complexity as $STRR$: both have a low quasi-$O(1)$ complexity. However, unlike $STRR$ and other existing round robin based low complexity packet schedulers that only have a good proportional fairness property, $FRR$ not only has a good proportional fairness property, but also maintains a quasi-$O(1)$ normalized worst-case fair index ($O(1)$ under practical assumptions). To the best of our knowledge, $FRR$ is the only round robin based scheduler with a similar complexity that has such capability.

The rest of the paper is structured as follows. Section II presents related work. Section III gives a motivating example and introduces the background of this work. Section IV describes $FRR$. Section V discusses the QoS properties of $FRR$. Section VI reports the results of the simulation study of $FRR$. Finally, Section VII concludes the paper.

## II. RELATED WORK

We will briefly discuss timestamp based and round-robin packet scheduling schemes since both relate to $FRR$. Some timestamp based schedulers, such as Weighted Fair Queuing ($WFQ$) [12] and Worst-case Fair Weighted Fair Queuing ($WF^2Q$) [1], [2], closely approximate the Generalized Processor Sharing ($GPS$) [5], [12]. These schedulers compute a timestamp for each packet by emulating the progress of a reference $GPS$ server and transmit packets in the increasing order of their timestamps. Other timestamp based approaches, such as Self-Clocked Fair Queuing (SCFQ) [6] and Virtual Clock [22], compute timestamps without referring to a reference $GPS$ server. These methods still need to sort packets according to their timestamps and still have an $O(log\ N)$ per packet processing complexity. The Leap Forward Virtual Clock [18] reduces the sorting complexity by coarsening timestamp values and has an $O(loglog\ N)$ complexity. This scheme requires complex data structures and is not suitable for hardware implementation.

Deficit Round Robin ($DRR$) [17] is one of the round-robin algorithms that enjoy a good proportional fairness property. A number of methods have recently been proposed to improve delay and burstiness properties of $DRR$ [7], [8], [13]. The Smoothed Round Robin (SRR) scheme [7] improves the delay and burstiness properties by spreading the data of a flow to be transmitted in a round over the entire round using a weight spread sequence. Aliquem [8], [9] allows the quantum of a flow to be scaled down, which results in better delay and burstiness properties. The Stratified Round Robin ($STRR$) [13] scheme bundles flows with similar rate requirements, scheduling the bundles through a sorted-priority mechanism, and using a round robin strategy to select flows in each bundle. $STRR$ guarantees that all flows get their fair share of *slots*. It enjoys a single packet delay bound that is independent of the number of flows in the system. However, as will be shown in the next section, the normalized worst-case fairness index for $STRR$ is $\Omega(N)$. $FRR$ is similar to $STRR$ in many aspects: $FRR$ uses exactly the same way to bundle the flows and has the same two-level scheduling structure. $FRR$ differs from $STRR$ in that it uses a different sorted-priority strategy to arbitrate among bundles, and a different round robin scheme to schedule flows within each bundle. The end result is

| $N$ | the number of flows in the system |
|---|---|
| $n$ | the number of classes in the system |
| $R$ | total link bandwidth |
| $r_i$ | guaranteed bandwidth for flow $f_i$ |
| $w_i = \frac{r_i}{R}$ | the weight associated with flow $f_i$ |
| $L_M$ | maximum packet size |
| $S_{i,s}(t_1, t_2)$ | the amount of work received by session $i$ during $[t_1, t_2)$ under the $s$ server |
| $S_{i,s}(t)$ | the amount of work received by session $i$ during $[0, t)$ under the $s$ server |
| $F_{i,s}^k$ | the departure time of the $k$th packet of flow $f_i$ under the $s$ server |
| $F_s^p$ | the departure time of packet $p$ under the $s$ server |
| $Q_{i,s}(t)$ | the queue size of flow $f_i$ at time $t$ under the $s$ server |
| $p_i^k$ | the $k$th packet on flow $f_i$ |

TABLE I

NOTATION USED IN THIS PAPER

that $FRR$ has a similar complexity and a similar proportional fairness property, but a much better worst-case fairness property. Bin Sort Fair Queuing (BSFQ) [4] uses an approximate bin sort mechanism to schedule packets. The worst-case single packet delay of BSFQ is proportional to the number of flows. Hybrid scheduling schemes [14], [15] have also been proposed. While the algorithm components of these schemes are similar to those of $FRR$ and $STRR$, the QoS properties of these schemes are not clear. A recently proposed group round robin scheme [3] uses a two-level scheduling scheme similar to that in [15]. While group round robin is conceptually similar to FRR, it cannot maintain the fairness in cases when not all flows are backlogged continuously.

## III. BACKGROUND

Some notations used in this paper are summarized in Table I. There are $N$ flows $f_1$, $f_2$, ..., $f_N$ sharing a link of bandwidth $R$. Each flow $f_i$ has a minimum guaranteed rate of $r_i$. We will assume that $\sum_{i=1}^N r_i \leq R$. The weight $w_i$ of flow $f_i$ is defined as its guaranteed rate normalized with respect to the total rate of the link, i.e., $w_i = \frac{r_i}{R}$. Thus, we have $\sum_{i=1}^N w_i \leq 1$.

### A. A motivating example

The development of $FRR$ is motivated by $STRR$ [13], a recently proposed round-robin algorithm. In terms of the QoS properties of scheduling results, $STRR$ is arguably the best scheduler among all existing round-robin based low complexity schedulers. We will show that the normalized worst-case fair index of $STRR$ is $\Omega(N)$.

Let $N + 1$ flows, $f_0$, $f_1$, ..., $f_N$, share an output link of bandwidth $2N$. The bandwidth of $f_0$ is $N$ and the bandwidth of each flow $f_i$, $1 \leq i \leq N$, is 1. We will use $R$ to denote the bandwidth of the output link, and $r_i$ to denote the bandwidth of flow $f_i$, $0 \leq i \leq N$. $R = 2N$, $r_0 = N$, and $r_i = 1$, $1 \leq i \leq N$. Let the maximum packet size be $L_M = 1000$ bits. Packets in $f_0$ are of size $L_M = 1000$ bits. Flows $f_1$, $f_2$, ..., $f_N$ are continuously backlogged with packets whose sizes repeat the pattern: $\frac{L_M}{2} = 500$ bits, $L_M = 1000$ bits, $\frac{L_M}{2} = 500$ bits,

500 bits, 1000 bits, 500 bits, and so on. Figure 1 (a) shows the packet arrival pattern assuming $N = 4$. In $STRR$, these flows are grouped into two classes: one class containing only $f_0$ and the other having flows $f_1, ..., f_N$. The bandwidth is allocated in the unit of slots. Let us assume that each of the flows has the minimum weight in its class and the credit assigned to each of the flows in a *slot* is $L_M = 1000$ bits. $STRR$ guarantees that slots are allocated fairly among all flows: $f_0$ is allocated one slot every two slots and each of the flows $f_1$, $f_2$, ..., $f_N$ gets one slot every $2N$ slots. We will use the $DRR$ concept of *round* to describe the scheduling results of $STRR$. In each round, all backlogged flows have a chance to send packets. In the example, each round contains $N$ slots from $f_0$ and one slot from each of the flows $f_i$, $1 \le i \le N$. Due to the differences in packet sizes, the sizes of slots are different. For $f_0$, each slot contains exactly one packet of size $L_M$ and the size of each slot is $L_M$. For $f_i$, $1 \le i \le N$, the size of the first slot is $\frac{L_M}{2} = 500$ bits since the second packet (size $L_M$) cannot be included in this slot. This results in 500-bit credits being passed to the next slot. Hence, the second slot for $f_i$ contains 2 packets (1500-bit data). This pattern is then repeated: the size for an evenly numbered slot for $f_i$, $1 \le i \le N$, is 500 bits and the size of an oddly numbered slot is 1500 bits. The $STRR$ scheduling result is shown in Figure 1 (b), where the rate allocated to flow $f_0$ oscillates between $\frac{4}{3}N$ and $\frac{4}{5}N$ for the alternating rounds. Note that since the size for each round depends on $N$, the duration of a round depends on $N$ and can be large.

Now consider the normalized worst-case fair index of $STRR$, $c_{STRR}$. Let us assume that in the example the last $f_0$ packet in round 1, which is the $2N - 1$-th packet of $f_0$, $p_0^{2N-1}$, arrives at time $a_0^{2N-1}$ right before the starting of round 1 (after the last $f_0$ packet in round 0 departs). At $a_0^{2N-1}$, the queue length in $f_0$ is $Q(a_0^{2N-1}) = N \times L_M$. Following the scheduling results shown in Figure 1 (b), $(N-1) \times L_M$ data in $f_0$ and $(N-1) \times (L_M + \frac{L_M}{2})$ data in flows $f_i$, $1 \le i \le N$, in round 1 are scheduled before $p_0^{2N-1}$. Let $d_0^{2N-1}$ be the departure time of $p_0^{2N-1}$. We have $d_0^{2N-1} - a_0^{2N-1} = \frac{(N-1)(L_M + 1\frac{1}{2}L_M) + L_M}{R}$. $C_{0,STRR} \ge d_0^{2N-1} - a_0^{2N-1} - \frac{Q(a_0^{2N-1})}{r_0} = \frac{(N-1)(L_M + 1\frac{1}{2}L_M) + L_M}{R} - \frac{N \times L_M}{r_0} = \frac{0.25 \times N \times L_M}{r_0} - 0.75 \frac{L_M}{r_0}$. Hence, $c_{STRR} \ge c_{0,STRR} \frac{r_0}{R} = 0.25 \times N \times \frac{L_M}{R} - 0.75 \frac{L_M}{R} = \Omega(N)$.

While we only show the normalized worst-case fair index of STRR in this section, one can easily show that the normalized worst-case fair indexes of other round robin schedulers such as smoothed round robin [7] and deficit round robin [17] are $\Omega(N)$: not having a good bound on worst-case fairness is a common problem with all of these low complexity round robin packet schedulers. $FRR$ overcomes this limitation and grants a much better a short-term fairness property while maintaining a low complexity. The scheduling results of $FRR$ for the example in Figure 1 are shown in Figure 1 (c). As can be seen from the figure, the short-term behavior of $f_0$ is much better than that in Figure 1 (b): counting from the beginning of round 0, for every 2000 bits data sent, exactly 1000 bits are from $f_0$.

### B. Deficit Round Robin (DRR)

Since $FRR$ is built over Deficit Round Robin ($DRR$) [17], we will briefly discuss $DRR$. Like an ordinary round robin scheme, $DRR$ works in rounds. Within each round, each backlogged flow has an opportunity to send packets. Each flow $f_i$ is associated with a quantity $Q_i$ and a variable $DC_i$ (deficit counter). The quantity $Q_i$ is assigned based on the guaranteed rate for $f_i$ and specifies the target amount of data that $f_i$ should send in each round. When flow $f_i$ cannot send exactly $Q_i$ data in a round, $DC_i$ records the quantum that is not used in a round so that the unused quantum can be passed to the next round. To ensure that each backlogged flow can send at least one packet in a round, $Q_i \ge L_M$. Some related properties of $DRR$ are summarized in the following lemmas.

**Lemma 1:** Assuming that flow $f_i$ is continuously backlogged during $[t_1, t_2]$. Let $X$ be the smallest number of continuous $DRR$ rounds that completely enclose $[t_1, t_2]$. The service received by $f_i$ during this period, $S_{i,DRR}(t_1, t_2)$, is bounded by $(X-3)Q_i \le S_{i,DRR}(t_1, t_2) \le (X+1)Q_i$.
*Proof:* See appendix. □

**Lemma 2**: Let $f_1, ..., f_N$ be the $N$ flows in the system with guaranteed rates $r_1, ..., r_N$. $\sum_{i=1}^{N} r_i \le R$. Let $r_{min} = \min_i\{r_i\}$ and $r_{max} = \max_i\{r_i\}$. Let $r_{max} = D * r_{min}$. Assume that $D$ is a constant with respect to $N$ and that $DRR$ is used to schedule the flows with $Q_i = L_M * \frac{r_i}{r_{min}}$. The following statements are true. All constants in this lemma are with respect to $N$.

1) Let packet $p$ arrive at the head of the queue for $f_i$ at time $t$. There exists a constant $c_1$ ($c_1 = O(D^2)$) such that packet $p$ will be serviced before $t + c_1 \times \frac{L_M}{r_i}$.
2) The normalized worst-case fair index of $DRR$ is a constant $c_1$ ($c_1 = O(D^2)$).
3) Let $f_i$ and $f_j$ be continuously backlogged during any given time period $[t_1, t_2]$, there exists two constants $c_1$ and $c_2$ ($c_1 = O(D)$ and $c_2 = O(D)$) such that $|\frac{S_{i,DRR}(t_1,t_2)}{r_i} - \frac{S_{j,DRR}(t_1,t_2)}{r_j}| \le c_1 \frac{L_M}{r_i} + c_2 \frac{L_M}{r_j}$.

*Proof*: See appendix. □

We will call $D = \frac{r_{max}}{r_{min}}$, the *maximum weight difference factor*. Lemma 2 shows that when $D$ is a constant with respect to $N$, $DRR$ is an excellent scheduler with both a good worst case fairness property and a good proportional fairness property. The problem is that when the weights of the flows differ significantly ($D$ is a large number), which is common in practice, the QoS performance bounds, which are functions of $D$, become very large.

$FRR$ extends $DRR$ such that the QoS properties in Lemma 2 hold for any weight distribution, while maintaining a low quasi-$O(1)$ complexity. The basic idea is as follows. $FRR$ chooses a constant $C$ (e.g. $C = 2$) that is independent of $D$ and $N$. $FRR$ groups flows whose weights differ by at most a factor of $C$ into classes and uses a variation of $DRR$ to schedule packets within each class. From Lemma 2, $DRR$ can achieve good QoS properties for flows in each class. Thus, the challenge is to isolate the classes so that flows in different classes, which are flows with significantly different weights, do not affect each other too much. $FRR$ uses a timestamp based scheduler to isolate the classes. As a result, $FRR$ schedules packets in two levels, a timestamp based inter-class scheduling and a $DRR$ based intra-class scheduling.

## IV. FRR: A FAIR ROUND ROBIN SCHEDULER

Like stratified round robin ($STRR$) [13], $FRR$ groups flows into a number of classes with each class containing flows with similar weights. For $k \ge 1$, class $F_k$ is defined as

$$F_k = \{f_i : \frac{1}{C^k} \le w_i < \frac{1}{C^{k-1}}\},$$

(a) Packet arrival pattern

(b) Scheduling results using STRR

(c) Scheduling results using FRR

Fig. 1. A motivating example

where $C$ is a constant independent of $D$ and $N$. Let $r$ be the smallest unit of bandwidth that can be allocated to a flow. The number of classes is $n = \lceil log_C(\frac{R}{r}) \rceil$. In practice, $n$ is usually a small constant. For example, consider an extreme case with $R = 1Tbps$, $r = 1Kbps$ ($D = 10^9$). When $C = 8$, $n = \lceil log_8(10^9) \rceil = 10$. Like [13], we will consider the practical assumption that $n$ is an $O(1)$ constant. However, since $n = \lceil log_C(\frac{R}{r}) \rceil$ in theory, we will derive the bounds on QoS properties and complexity in terms of $n$.

It must be noted that the constant $C$ in $FRR$ is very different from the constant weight difference factor $D$ in Lemma 2. $D$ specifies a limit on the type of flows that can be supported in the system. $C$ is an algorithm parameter that can be selected by the scheduler designer and does not put a limit on the weights of the flows in the system. Consider the case when $R = 1Tbps$ and $r = 1kbps$. $D = \frac{10^{12}}{10^3} = 10^9$. Using $DRR$ to schedule packets may result in extremely poor QoS bounds since $O(D^2)$ can be huge numbers. With $FRR$, one can select a small number $C$ (e.g. $C = 2$) and obtain QoS bounds that are linear functions of $C$ and $n$.

$FRR$ has two scheduling components, intra-class scheduling that determines the order of the packets within each class and the weight of the class, and inter-class scheduling that determines the class, and thus, the packet within the class, to be transmitted over the link. The concept of *weight* will be used in different contexts. A weight is associated with each flow. In intra-class scheduling, the packet stream within a class is partitioned into frames. A frame, which is a logical unit that is similar to a round in $DRR$, contains a set of packets that are scheduled with the same weight in the inter-class scheduling. Notice that a frame in this paper is not the transmission unit in the Medium Access Control (MAC) layer. A weight that represents the aggregate weight for all active flows in a frame is assigned to the frame. The weight of a frame is then used in inter-class scheduling to decide which class is to be served. In the inter-class scheduling, we will also call the weight of current frame in a class the weight of the class. We will use notion $w_i$ to denote the weight of a flow $f_i$ and $W_k$ to denote the weight of a class $F_k$.

*A. Intra-class scheduling*

Assuming that the inter-class scheduling scheme can provide fairness among classes based on their weights, the intra-class scheduler must be able to transfer the fairness at the class level to that at the flow level. To focus on the intra-class scheduling issues, we will assume that $GPS$ is the inter-class scheduling scheme in this sub-section.

The intra-class scheduling scheme in $FRR$, called *Lookahead Deficit Round Robin with Weight Adjustment* ($LDRRWA$), is a variation of $DRR$ with two extensions: a *lookahead* operation and a *weight adjustment* operation. To understand the needs for the two extensions, let us examine the issues when a vanilla $DRR$ scheme is used in intra-class scheduling. In $DRR$, the packet stream within a class is partitioned into *rounds*. Each of active flows is allocated a *quantum* for sending data in a round. To offset the weight differences among the flows, each flow $f_i \in F_k = \{f_i : \frac{1}{C^k} \le w_i < \frac{1}{C^{k-1}}\}$ is assigned a quantum of $Q_i = C^k w_i L_M$. Since $\frac{1}{C^k} \le w_i < \frac{1}{C^{k-1}}$, $L_M \le Q_i < C \times L_M$.

Since inter-class scheduling in $FRR$ is based on based on weights, it is crucial to assign weights to each class such that flows in all classes are treated fairly. In $DRR$, different flows can be active in different rounds. Since the weight assigned to a class must reflect the weights of all active flows, it is natural to assign a different weight to a different round. One simple approach, which is adopted in the group round robin scheme [3], is to assign the sum of weights of all active flows in a round as the weight of the round. This simple approach, however, does not yield a fair scheduler. Consider the case when two flows, $f_1$ and $f_2$, of the same weight $w_1 = w_2 = \frac{1}{C^k}$ are in a class $F_k$. $Q_1 = Q_2 = L_M$. Flow $f_1$ is continuously backlogged and sends $L_M$ data in each round. Flow $f_2$ is active and sends $\frac{Q_2}{2} = \frac{L_M}{2}$ data in each round. The simple approach will assign weight $w_1 + w_2 = 2w_1$ to each round, which results in the guaranteed service rate under $GPS$ for this class to be $2r_1$. Since $\frac{2}{3}$ of the service is used to serve packets in $f_1$, the guaranteed rate for $f_1$ is artificially inflated to $2r_1 \times \frac{2}{3} = \frac{4}{3}r_1$, which is unfair to flows in other classes.

What is the fair weight for a round in class $F_k$? In each round, each active flow $f_i \in F_k$ with a rate $r_i$ is given a quantum of

$Q_i$. The targeted finishing time for $f_i$ is thus $\frac{Q_i}{r_i} = \frac{C^k L_M}{R}$. From Lemma 1, we can see that for a flow $f_i$ that is continuously backlogged in $X$ rounds, the amount of data sent is at most a few packets from $X \times Q_i$. Hence, if the weights for all rounds are assigned such that the service time for each round is $\frac{Q_i}{r_i} = \frac{C^k L_M}{R}$ using the guaranteed service rate, all continuously backlogged flows in the $X$ rounds obtain their fair share of the bandwidth with a small error margin: the fairness at the class level is transferred to the fairness at the flow level. Hence, **the fair weight for a round in class $F_k$ should be one that results in the targeted finishing time of $\frac{Q_i}{r_i} = \frac{C^k L_M}{R}$.**

Let flows $f_1, ..., f_m$ (in a class) be active in a round. Let the data sizes of $f_i$, $1 \le i \le m$, in the round be $s_i$. The size of the round is $round\ size = s_1 + s_2 + ... + s_m$. Let $w'$ be the fair weight for the round and $r'$ be the corresponding guaranteed service rate, $w' = \frac{r'}{R}$. We have $\frac{round\ size}{r'} = \frac{round\ size}{w'R} = \frac{C^k L_M}{R}$. Solving the equation, we obtain
$$w' = \frac{round\ size}{C^k L_M} = \frac{s_1 + s_2 + ... + s_m}{C^k L_M}.$$
$\frac{s_1 + s_2 + ... + s_m}{C^k L_M} = \frac{s_1}{C^k L_M} + \frac{s_2}{C^k L_M} + ... + \frac{s_m}{C^k L_M} = \frac{s_1}{Q_1} w_1 + \frac{s_2}{Q_2} w_2 + ... + \frac{s_m}{Q_m} w_m$: the fair weight for a round can also be interpreted as the sum of the normalized weights of active flows, $\frac{s_i}{Q_i} w_i$: in order to obtain the fair weight for a round, the weight of each active flow $f_i$ must be adjusted from $w_i$ to $\frac{s_i}{Q_i} w_i$ before the adjusted weights are aggregated.

Although the weight adjustment results in the fair weight for a round $w' = \frac{round\ size}{C^k L_M}$, $w'$ may sometimes be more than the sum of the weights of all flows in the class. For example, if a class has only one flow $f_1$ and $Q_1 = L_M$, $f_1$ may send $0.5 L_M$ in one round and $1.5 L_M$ in another round. Using weight adjustment, the weight for the class is $0.5 w_1$ in one round and $1.5 w_1$ in the other round. This temporary raising of weights to $1.5 w_1$ may violate the assumption that the sum of the weights for all classes is less than 1, which is essential for guaranteeing services. $LDRRWA$ uses the lookahead operation to deal with this problem. The lookahead operation moves some currently backlogged packets that are supposed to send in the next round under $DRR$ into the current around. By using the lookahead operation, the size of each round (now called **frame** to be differentiated from the $DRR$ round) is no more than the sum of the quota of all active flows in the round. This guarantees that the fair weight assigned to each frame to be less than the sum of the weights of all active flows in the frame.

*Lookahead Deficit Round Robin with Weight Adjustment ($LDRRWA$)*

In $LDRRWA$, an active flow $f_i \in F_k = \{f_i : \frac{1}{C^k} \le w_i < \frac{1}{C^{k-1}}\}$ is assigned a quantum of
$$Q_i = 2C^k w_i L_M.$$
Since $\frac{1}{C^k} \le w_i < \frac{1}{C^{k-1}}$, $2L_M \le Q_i < 2C \times L_M$. $Q_i$ in $LDRRWA$ is two times the value in $DRR$. The reason is that the deficit counter may be negative in $LDRRWA$, $Q_i = 2C^k w_i L_M$ ensures that a backlogged flow can at least send one packet in a frame. Since $Q_i = 2C^k w_i L_M$, $\frac{Q_i}{r_i} = \frac{2C^k L_M}{R}$ and the fair weight for a frame of size $framsize$ is
$$W_k = \frac{framesize}{2C^k L_M}.$$
Let $f_1, f_2, ..., f_m$ be the flows in class $F_k$. $W_k = \frac{framesize}{2C^k L_M} = \frac{framesize}{\sum_{i=1}^{m} Q_i} \sum_{i=1}^{m} w_i$. $LDRRWA$ employs the lookahead oper-

| variable | explanation |
|---|---|
| $deficitcount_i$ | the deficit count for flow $f_i$ |
| $remaindeficit$ | the sum of quantum not used in the $DRR$ round |
| $lastingflowlist$ | the flows that last to the next frame |
| $framesize$ | the size of the frame |
| $frameweight$ | the weight for the frame |
| $remainsize$ | size of the part of a packet that belongs to current frame |

TABLE II

MAJOR VARIABLES USED IN THE FRAME CALCULATION ALGORITHM

ation to ensure that $framesize \le \sum_{i=1}^{m} Q_i$, which results in $W_k \le \sum_{i=1}^{m} w_i$. A frame in $LDRRWA$ has two parts: the first part includes all packets that are supposed to be sent using $DRR$; the second part includes packets from the lookahead operation. In the lookahead operation, packets from flows that do not use up their quanta and are still backlogged after the current $DRR$ round are sent in the current frame. Notice that in order for a flow to be considered as backlogged after the $DRR$ round, the flow must have at least one currently backlogged packet after the $DRR$ round. Let us denote this packet as the **lookahead packet** for the backlogged flow. Clearly, the size of the lookahead packet is larger than the remaining quota for the flow; and the total remaining credits for the class at the end of the DRR round are no more than the sum of the sizes of all lookahead packets of all backlogged flows. Hence, each flow with a lookahead packet may contribute the lookahead packet in the frame until all remaining credits are consumed. Note that after a flow contributes its packet in the lookahead operation, the deficit counter for this flow has a negative value. The lookahead operation ensures that the aggregate deficit (the sum of the deficits) of all the backlogged flows in every frame is exactly 0 at frame boundaries. In other words, no credit is passed over frame boundaries at the frame level. As a result, the size of each frame is less than or equal to the total credits generated in that frame, which is at most $\sum_{i=1}^{m} Q_i$. Note that the frame boundary may not align with packet boundary: a packet may belong to two frames. Note also that while the aggregate deficit of all backlogged flows is 0 at frame boundaries, each individual flow may have a positive, zero, or negative deficit counter. Allowing a flow to have a negative deficit may potentially cause problems: a flow may steal credits by over sending in a frame (and having a negative deficit at the frame boundary), becoming inactive for a short period of time (so that the negative deficit can be reset), and over sending again. To handle this situation, $LDRRWA$ keeps the negative deficit for one frame when the flow becomes inactive before it resets the negative deficit counter for the flow.

Each frame is decided at the time it starts. Packets arrive during the current frame are sent in later frames. Note that delaying a packet for one frame does not affect the fairness of the scheduler. Next, we will describe the high level logical view of $LDRRWA$. A detailed packet-by-packet implementation of $LDRRWA$ is given in [21].

Figure 2 shows the logic for computing each $LDRRWA$ frame and its weight. Table II summarizes the major variables in the algorithm. Like $DRR$, $deficitcount_i$ is associated with flow $f_i$ to maintain the credits to be passed over to the next $DRR$ round

**Algorithm for computing the next frame for class $F_k$**

(1) $remaindeficit = framesize = 0$
(2) $lastingflowlist = NULL$
(3) **if** $(remainsize > 0)$ **then**
    /*The partial packet belongs to this frame */
(4)     $framesize = framesize + remainsize$
(5) **end if**
    /* forming the $DRR$ round */
(6) **for** each active flow $f_i$ **do**
(7)     $deficitcount_i = deficitcount_i + quantum_i$
(8)     **while** $(deficitcount_i > 0)$ **and** $(f_i$ not empty$)$ **do**
(9)         $pktsize = size(head(f_i))$
(10)         **if** $(pktsize < deficitcount_i)$ **then**
(11)           remove head from $f_i$ and put it in the frame
(12)           $framesize = framesize + pktsize$
(13)           $deficitcount_i = deficitcount_i - pktsize$
(14)         **else** break
(15)         **end if**
(16)     **end while**
(17)     **if** $(f_i$ is empty $)$ **then**
(18)         $deficitcount_i = 0$
(19)     **else**
(20)         $remaindeficit = remaindeficit + deficitcount_i$
(21)         insert $f_i$ to $lastingflowlist$
(22)     **end if**
(23) **end for**
    /* lookahead operation */
(24) $f_i = $ head(lastingflowlist)
(25) **while** $(f_i \neq NULL)$ **and** $(remaindeficit > 0)$ **do**
(26)     $pktsize = size(head(f_i))$
(27)     **if** $(pktsize < remaindeficit)$ **then**
(28)         remove head from $f_i$ and put it in the frame
(29)         $framesize = framesize + pktsize$
(30)         $remaindeficit = remaindeficit - pktsize$
(31)         $deficitcount_i = deficitcount_i - pktsize$
(32)     **else** break
(33)     **end if**
(34)     $f_i = nextflow(f_i)$
(35) **end while**
(36) **if** $(f_i \neq NULL)$ **then**
(37)     $pktsize = size(head(f_i))$
(38)     remove head from $f_i$ and put it in the frame
(39)     $framesize = framesize + remaindeficit$
(40)     $remainsize = pktsize - remaindeficit$
(41)     $deficitcount_i = deficitcount_i - pktsize$
(42) **end if**
    /* computing the weight */
(43) $frameweight = \frac{framesize}{2C^k L_M}$
(44) **if** $(frameweight < \frac{1}{C^k})$ $frameweight = \frac{1}{C^k}$

Fig. 2. The algorithm for computing the next frame for class $F_k$

and decide the amount of data to be sent in one round. After each $DRR$ round, $remaindeficit$ maintains the sum of the quanta not used in the current $DRR$ round, that is, the quanta that cannot be used since the size of the next backlogged packet is larger than the remaining quanta for a flow. In traditional $DRR$, these unused quanta will be passed to the next $DRR$ round. In $LDRRWA$, in addition to passing the unused quanta to the next $DRR$ round, some packets that would be sent in the next $DRR$ round are placed in the current $LDRRWA$ frame so that at frame boundaries $remaindeficit$ is always equal to 0. This is the lookahead operation. The $lastingflowlist$ contains the list of flows that are backlogged at the end of the current $DRR$ round. Flows in $lastingflowlist$ are candidates to supply packets for the lookahead operation. $Frameweight$ is the weight to be used by inter-class scheduling for the current frame. Variable $framesize$ records the size of the current frame. Since $FRR$ needs to enforce that $remaindeficit = 0$ at frame boundaries, frame boundaries may not align with packet boundaries and a packet may belong to two frames. Variable $remainsize$ is the size of the part of the last packet in the frame that belongs to the next frame, and thus, should be counted in the $framesize$ for the next frame.

Let us now examine the algorithm in Figure 2. In the initialization phase, line (1) to line (5), variables are initialized and $remainsize$ is added to $framesize$, which effectively includes the partial packet in the frame to be computed. After the initialization, there are three main components in the algorithm: forming a $DRR$ round, lookahead operation, and weight calculation. In the first component, line (6) to line (23), the algorithm puts all packets in the current $DRR$ round that have not been served into the current frame. In the second component, line (24) to line (42), the algorithm performs the lookahead operation by moving some packets in the next $DRR$ round into the current frame so that $remaindeficit = 0$ at the frame boundary. Notice that each backlogged flow can contribute at most one packet in the lookahead operation. Notice also that while a class as a whole does not pass credits between frames, an individual flow can pass credits from one frame to the next: $deficitcount_i$ may have negative, 0, or positive values at frame boundaries. Finally, lines (43) and (44) compute the weight for the frame.

The complexity of the algorithm in Figure 2 is $O(M)$, where $M$ is the number of packets in a frame. Clearly, this high level algorithm cannot be directly used in a scheduler to determine the next frame and frame weight. Otherwise, it will introduce an $O(M)$ processing complexity, which is larger than $O(N)$. The operations in $LDRRWA$ can be realized in the packet-by-packet operations when packets arrive and depart. By distributing the $O(M)$ operations for determining a frame into $O(M)$ packet arrivals and departures in a frame, $LDDRWA$ only introduces $O(1)$ per packet processing overheads.

A detailed description of the packet-by-packet operations of $LDRRWA$ is given in [21]. While the detailed packet-by-packet operations are rather tedious, the idea is straight-forward. $LDRRWA$ maintains active flows in different queues. To determine a frame and its weight, our scheme determines (1) the total size of the frame, and (2) for each active flow in the frame the size of the data in that flow that belong to the frame. Such information is obtained by maintaining the following information at each packet arrival and departure: the size of the remaining current frame, the size of the next frame, the size of the partial packet in the current frame, the starting time of the current frame,

the deficit counter for each flow, the size of the data for each flow in the current frame, the size of the data for each flow in the next frame, the size of all backlogged data in a flow, and the last time that a flow is serviced. Clearly, bookkeeping for all these variables takes $O(1)$ operations. With such information, the computation of the next frame, as well as the whole $LDRRWA$ can be realized in $O(1)$ operations.

*Properties of LDRRWA*

**Lemma 3**: Assuming that flow $f_i$ is continuously backlogged during $[t_1, t_2)$. Let $X$ be the smallest number of continuous $LDRRWA$ frames that completely enclose $[t_1, t_2)$. The service received by $f_i$ during this period, denoted as $S_{i,LDRRWA}(t_1, t_2)$, is bounded by
$$(X - 3)Q_i \le S_{i,LDRRWA}(t_1, t_2) \le (X + 1)Q_i.$$
*Proof:* See appendix. □

Comparing Lemma 3 and Lemma 1, we can see the similarity between $DRR$ and $LDRRWA$: in both schemes, the amount of data sent from a flow $f_i$ that is continuously backlogged for $X$ frames (rounds) is a few packets from $X \times Q_i$.

**Lemma 4**: In $LDRRWA$, the weight for a frame is less than or equal to the sum of the weights of all flows in the class.
*Proof*: Obvious from the previous discussion. □

At any given time, let $W_k$, $1 \le k \le n$ be the weights for the $n$ classes. Lemma 4 establishes that $\sum_{i=1}^n W_k \le \sum_{i=1}^N w_i \le 1$. Thus, under $GPS$, the bandwidth allocated to class $k$ is given by $\frac{W_k}{\sum_{i=1}^n W_k} R \ge R \times W_k$. We will call $R \times W_k$ the $GPS$ *guaranteed rate*.

**Lemma 5**: Under $GPS$, the time to serve each $LDRRWA$ frame in class $F_k$ is at most $\frac{2C^k L_M}{R}$.
*Proof:* Normally, the frame weight is computed as $W_k = \frac{framesize}{2C^k L_M}$ (line (43) in Figure 2). In cases when $\frac{framesize}{2C^k L_M}$ is less than the smallest weight for a flow in a class, the weight is increased (line (44) in Figure 2) to the smallest weight.

When $W_k = \frac{framesize}{2C^k L_M}$, the $GPS$ guaranteed rate is $R \frac{framesize}{2C^k L_M}$ and the total time to serve the frame is at most $\frac{framesize}{R \frac{framesize}{2C^k L_M}} = \frac{2C^k L_M}{R}$. If $W_k$ is increased, the conclusion still holds. □

**Lemma 6**: Under $GPS$, the time to service $X$ bytes of data in class $F_k$ is at most $\frac{XC^k}{R}$.
*Proof*: The minimum weight assigned to a backlogged class $F_k$ is $\frac{1}{C^k}$. Thus, the $GPS$ guaranteed rate for class $F_k$ is at least $\frac{R}{C^k}$. Thus, the time to serve a queue of size X bytes in class $F_k$ is at most $\frac{X}{\frac{R}{C^k}} = \frac{XC^k}{R}$. □

**Lemma 7**: For a class $F_k$ frame of size no smaller than $2L_M$, the service time for the frame is exactly $2C^k \frac{L_M}{R}$ using the $GPS$ guaranteed rate.
*Proof*: When $framesize \ge 2L_M$, $W_k = \frac{framesize}{2C^k L_M} \ge \frac{1}{C^k}$. Thus, the $GPS$ guaranteed rate for the frame is $R \frac{framesize}{2C^k L_M}$ and the service time for the frame with the guaranteed rate is $\frac{framesize}{R \frac{framesize}{2C^k L_M}} = \frac{2C^k L_M}{R}$. □

**Lemma 8**: Let a class $F_k$ frame contain packets of a continuously backlogged flow $f_i$, the size of frame is no smaller than $2L_M$.
*Proof*: Straight-forward from the fact that no credit is passed from the previous frame and to the next frame and that $Q_i \ge 2L_M$. □

**Lemma 9**: Let $f_i \in F_k$ and $f_j \in F_m$ be continuously backlogged during $[t_1, t_2)$. $k \ge m$. Let $X_k$ and $X_m$ be the smallest numbers

of $F_k$ and $F_m$ frames that completely enclose $[t_1, t_2)$. Assume that classes $F_k$ and $F_m$ are served with the $GPS$ guaranteed rate,
$$(X_k - 1)C^{k-m} \le X_m \le X_k C^{k-m} + 1.$$
*Proof:* Since $f_i \in F_k$ and $f_j \in F_m$ are continuously backlogged during $[t_1, t_2)$, the sizes of all frames during this period are no smaller than $2L_M$ (Lemma 8). From Lemma 7, using the $GPS$ guaranteed rate, the time to service a class $F_k$ frame is exactly $\frac{2C^k L_M}{R}$ and the time for a class $F_m$ frame is exactly $\frac{2C^m L_M}{R}$. Since $X_k$ and $X_m$ are the smallest numbers of $F_k$ and $F_m$ frames that completely enclose $[t_1, t_2)$, we have

$$t2 - t1 \le X_k \frac{2C^k L_M}{R} \le t2 - t1 + \frac{2C^k L_M}{R}$$
$$t2 - t1 \le X_m \frac{2C^m L_M}{R} \le t2 - t1 + \frac{2C^m L_M}{R}$$

Hence, $(X_k - 1)C^{k-m} \le X_m \le X_k C^{k-m} + 1$. □

**Lemma 10**: Let $f_i \in F_k$ and $f_j \in F_m$ be continuously backlogged during $[t_1, t_2)$. $k \ge m$. Let $X_k$ and $X_m$ be the smallest numbers of $F_k$ and $F_m$ frames that completely enclose $[t_1, t_2)$. Assume that the inter-class scheduler is $GPS$,
$$(X_k - 1)C^{k-m} \le X_m \le X_k C^{k-m} + 1.$$
*Proof:* See appendix. □



Fig. 3.  An example

We will use an example to illustrate how $LDRRWA$ interacts with inter-class scheduling to deliver fairness among flows in different classes. Let us assume that $GPS$ is the inter-class scheduling algorithm. Consider scheduling for a link with 4 units of bandwidth with the following settings. $C = 2$ and there are two classes where $F_1 = \{f_i : \frac{1}{2} \le w_i < 1\}$ and $F_2 = \{f_i : \frac{1}{4} \le w_i < \frac{1}{2}\}$. Three flows, $f_1$, $f_2$ and $f_3$, with rates $r_1 = 2$ and $r_2 = r_3 = 1$ are in the system. $w_1 = 1/2$, $w_2 = 1/4$, and $w_3 = 1/4$. Thus, $f_1$ is in $F_1$, and $f_2$ and $f_3$ are in $F_2$. Let $L_M$ be the maximum packet size. The quantum for each of the three flows is $2L_M$. All packets in $f_1$ are of size $L_M$, all packets in $f_2$ are of size $0.99L_M$ and all packets in $f_3$ are of size $0.01L_M$. Flows $f_1$ and $f_2$ are always backlogged. Flow $f_3$ is not always backlogged, its packets arrive in such a way that exactly one packet arrives before a new frame is to be formed. Thus, each $F_2$ frame contains one packet from $f_3$. The example is depicted in Figure 3. As shown in the figure, each $F_1$ frame contains exactly two packets from $f_1$. For $F_2$, the lookahead operation always moves part of the $f_2$ packet in the next $DRR$ round into the current frame, and thus, the frame boundaries are not aligned with packet boundaries.

The weight for $F_1$ is always $1/2$. For $F_2$, the lookahead operation ensures that the size of $f_2$ data in a frame is $2L_M$, and thus, the size of each $F_2$ frame is $2L_M + 0.01L_M = 2.01L_M$. The weight of $F_2$ is computed as $W_2 = \frac{framesize}{2C^k L_M} = \frac{2.01L_M}{8L_M} = \frac{2.01}{8}$. Hence, $F_1$ (and thus $f_1$) is allocated a bandwidth of $4 * \frac{\frac{1}{2}}{\frac{1}{2} + \frac{2.01}{8}} = \frac{16}{6.01} > 2$. $F_2$ is allocated a bandwidth of $4 * \frac{\frac{2.01}{8}}{\frac{1}{2} + \frac{2.01}{8}}$. For each

$F_2$ frame of size $2.01L_M$, $2L_M$ belongs to $f_2$. Thus, the rate allocated to $f_2$ is $4 * \frac{\frac{2.01}{8}}{\frac{1}{2}+\frac{2.01}{8}} * \frac{2L_M}{2.01L_M} = \frac{8}{6.01} > 1$. The rates allocated to $f_1$ and $f_2$ are larger than their guaranteed rates and the worst-case fairness is honored. The ratio of the rates allocated to $f_1$ and $f_2$ is equal to $\frac{\frac{16}{6.01}}{\frac{8}{6.01}} = 2$, which is equal to the ratio of their weights. Thus, the proportional fairness is also honored.

*B. Inter-class scheduling*

$LDRRWA$ assigns different weights for different frames in a class. Moreover, the weight of a frame is decided only after the last packet in the previous frame is sent. Hence, the inter-class scheduling must be able to handle these situations while achieving fair sharing of bandwidth. Although $GPS$ can achieve fair sharing, none of the existing timestamp based schemes can closely approximate $GPS$ under such conditions. We develop a new scheme called Dynamic Weight Worst-case Fair weighted Fair Queuing ($DW^2F^2Q$). $DW^2F^2Q$ has the same scheduling result as $WF^2Q$ [2] when the weights do not change. Theorems presented later show that the difference between the packet departure times under $DW^2F^2Q$ and $GPS$ is at most $(n-1)L_M$, where $n$ is the number of classes. This bound is sufficient for $FRR$ to achieve its QoS performance bounds.

$DW^2F^2Q$ uses the virtual time concept in [12] to track the $GPS$ progress up to the point that it can accurately track, and schedules packets based on their virtual starting/finishing times. Let us denote an event in the system the following: (1) the arrival of a packet to the $GPS$ server, (2) the departure of a packet from the $GPS$ server, and (3) the weight change of a class ($LDRRWA$ may change weight within a packet). Let $t_j$ be the time at which the $j$th event occurs. Let the time of the first arrival of a busy period be denoted as $t_1 = 0$. For each $j = 2, 3, ...$, the set of classes that are busy in the interval $[t_{j-1}, t_j)$ is denoted as $B_{j-1}$. Let us denote $W_{k,j-1}$ the weight for class $F_k$ during the interval $[t_{j-1}, t_j)$, which is a fixed value. Virtual time $V(t)$ is defined to be zero for all times when the system is idle. Assuming that each busy period begins with time 0, $V(t)$ evolves as follows:

$$V(0) = 0$$
$$V(t_{j-1} + \tau) = V(t_{j-1}) + \frac{\tau}{\sum_{k \in B_{j-1}} W_{k,j-1}},$$
$$0 < \tau \leq t_j - t_{j-1}, j = 2, 3, ...$$

As discussed in [12], the rate of change of $V$, $\frac{\partial V(t_j+\tau)}{\partial \tau}$, is $\frac{1}{\sum_{k \in B_j} W_{k,j}}$, and each backlogged class $F_k$ receives service at rate $W_{k,j} \frac{\partial V(t_j+\tau)}{\partial \tau}$. Let us denote the virtual starting time, virtual finishing time, real arrival time, and size of the $i$-th packet $P_k^i$ in $F_k$ as $Vstart(P_k^i)$, $Vfinish(P_k^i)$, $arrive(P_k^i)$, and $size(P_k^i)$, respectively. Let $W_k$ be the weight for the frame that includes $P_k^i$. We have

$$Vstart(P_k^i) = max(Vfinish(P_k^{i-1}), V(arrive(P_k^i)))$$
$$Vfinish(P_k^i) = Vstart(P_k^i) + \frac{size(P_k^i)}{W_k}$$

$DW^2F^2Q$ keeps track of $B_j$ in order to track the progress of the virtual time. When the last packet in a frame is sent later than its virtual finishing time (the packet departed under $GPS$, but not under $DW^2F^2Q$), the weight of the class after the frame virtual finishing time is unknown (until the last packet is sent). Hence, $DW^2F^2Q$ cannot always track the virtual time up to the current time. Before each scheduling decision is made, $DW^2F^2Q$ tracks

the virtual time either to the earliest time when there is a class with an unknown weight or to the current time when the weights of all classes are known up to the current time. We will denote this time (the latest time that $DW^2F^2Q$ can track $GPS$ progress accurately) $T$ and the corresponding virtual time $V(T)$. Hence, either $T$ is the current time, or there exists one class $F_k$ whose current frame virtual finishing time is $V(T)$ and the last packet in that frame has not been sent. $DW^2F^2Q$ only deals with $n$ classes. As a result, it can afford to recompute timestamps for the first packets in all classes in every scheduling step without introducing excessive overheads ($O(n) = O(1)$ under our assumption). To ease composition, we will ignore the issue related to the timing to assign the timestamps to the packets since the timestamps can be recomputed before each scheduling decision is made. $DW^2F^2Q$ has exactly the same criteria as $WF^2Q$ for scheduling packets: (1) only packets whose virtual starting times are earlier than the current virtual time are eligible; and (2) among all eligible packets, the one with the smallest virtual finishing time is selected.

There are potentially two problems in $DW^2F^2Q$. First, determining the virtual finishing time for the last packet in a frame can be a problem when only a part of the packet belongs to the current frame. The weight for the part of the packet in the next frame is unknown until the packet is scheduled. $DW^2F^2Q$ assigns the frame virtual finishing time as the packet virtual finishing time for this type of packets. Although this creates some inaccuracy, the scheduling results are still sufficiently good as will be shown later.

The other problem is that $T$ may not be the current time and the virtual time corresponding to the current time is unknown. In this case, there must exist one class $F_k$ whose current frame virtual finishing time is $V(T)$ and the last packet in that frame, $P$, has not been sent. The virtual finishing time (and the virtual starting time) of $P$ must be less than or equal to $V(T)$. Since $DW^2F^2Q$ has accurate virtual time up to time $T$, the timestamps for all packets with finishing times less than $V(T)$ are available. All unknown virtual finishing times for packets must be larger than $V(T)$. In this case, the packet to be scheduled must have a virtual finishing time less than or equal to $V(T)$. Hence, $DW^2F^2Q$ can simply assign a large timestamp as the virtual finishing time for packets with unknown virtual finishing times (to prevent these packets to be scheduled) and only consider packets whose virtual finishing time is less than or equal to V(T) when $T$ is less than the current time. Not being able to tracking virtual time up to the current time does not prevent $DW^2F^2Q$ from selecting the right packet for transmission. Note that when $T$ equals the current time, $DW^2F^2Q$ schedules packets exactly like $WF^2Q$.

The packet-by-packet operations of $DW^2F^2Q$ are given in [21], where we show that the worst-case per packet scheduling complexity is $O(nlg(n))$. In practical cases, $n = O(1)$ and hence, the per packet complexity of $DW^2F^2Q$ is also $O(1)$. The following theorems shows properties of $DW^2F^2Q$.

**Theorem 1**: $DW^2F^2Q$ is work conserving.
*Proof*: See appendix. □

Since both $GPS$ and $DW^2F^2Q$ are work-conserving disciplines, their busy periods coincide. We will consider packet scheduling within one busy period. Let $F_{i,s}^k$ be the departure time of the $k$th packet in class $i$ under server $s$ in a busy period.
**Lemma 11**: If $F_{i,GPS}^k \leq F_{j,GPS}^m$, $F_{i,DW^2F^2Q}^k < F_{j,DW^2F^2Q}^{m+1}$.
*Proof:* Let $p_i^l$ be the packet at the head of class $i$ at time $t$ when $p_j^{m+1}$ is at the head of class $j$ and is eligible to be transmitted.

Let the timestamp assigned to $p_j^{m+1}$ be $VT$, we have $VT > V(F_{j,GPS}^m)$. This applies even when $p_j^{m+1}$ is the last packet in a frame and is assigned an inaccurate timestamp.

If $l > k$, we have $F_{i,DW^2F^2Q}^k < F_{j,DW^2F^2Q}^{m+1}$ and the lemma is proved. If $l \leq k$, $F_{i,GPS}^l < F_{i,GPS}^{l+1} < ... < F_{i,GPS}^k \leq F_{j,GPS}^m$ and $V(F_{i,GPS}^l) < V(F_{i,GPS}^{l+1}) < ... < V(F_{i,GPS}^k) \leq V(F_{j,GPS}^m) < VT$. For a packet $p_i^X$ that is in the frame boundary, its timestamp is less than or equal to $V(F_{i,GPS}^X)$. Since at time $t$, $p_j^{m+1}$ is eligible for scheduling, $V(t) \geq V(F_{j,GPS}^m)$ and the accurate virtual times for these packets are available, all of these packets have smaller virtual starting and finishing times than $p_j^{m+1}$ and will depart before $p_j^{m+1}$ under $DW^2F^2Q$. Thus, $F_{i,DW^2F^2Q}^k < F_{j,DW^2F^2Q}^{m+1}$. $\square$

Lemma 11 indicates that $DW^2F^2Q$ can at most introduce one packet difference between any two classes in comparison to $GPS$. This leads to the following theory that states that under the assumption that $n$ is a small constant, $DW^2F^2Q$ closely approximates $GPS$. Let $F_s^p$ be the time packet $p$ departs under server $s$.

**Theorem 2**: Let $n$ be the number of classes in the system,
$$F_{DW^2F^2Q}^p - F_{GPS}^p \leq (n-1)\frac{L_M}{R}.$$
*Proof:* Consider any busy period and let the time that it begins be time zero. Let $p_k$ be the $k$th packet of size $s_k$ to depart under $GPS$. We have $F_{GPS}^{p_k} \geq \frac{s_1+s_2+...+s_k}{R}$. Now consider the departure time of $p_k$ under $DW^2F^2Q$. From Lemma 11, each class can have at most one packet whose $GPS$ finishing time is after packet $p_k$ and whose $DW^2F^2Q$ finishing time is before packet $p_k$. Hence, there are at most $n-1$ packets (from the $n-1$ other classes) that depart before packet $p_k$ under $DW^2F^2Q$ and have a $GPS$ finishing time after $F_{GPS}^{p_k}$. Let the $n-1$ packets be $e_1, e_2, ..., e_{n-1}$ with sizes $se_1, se_2, ..., se_{n-1}$. All other packets depart before $p^k$ under $DW^2F^2Q$ must have $GPS$ finishing times earlier than $F_{GPS}^{p_k}$. We have $F_{DW^2F^2Q}^{p_k} \leq \frac{s_1+...+s_k+se_1+...+se_{n-1}}{R}$. Thus, $F_{DW^2F^2Q}^{p_k} - F_{GPS}^{p_k} \leq (n-1)\frac{L_M}{R}$. $\square$

## V. PROPERTIES OF FRR

This session formally analyzes the $QoS$ properties of $FRR$. We will prove that the three statements in Lemma 2 hold for $FRR$ with an arbitrary weight distribution.

**Theorem 3 (single packet delay bound)**: Let packet $p$ arrives at the head of flow $f_i \in F_k$ at time $t$. Using $FRR$, there exists a constant $c_1$ ($c_1 = O(C+n)$), such that $p$ will depart before $t + c_1 * \frac{L_M}{r_i}$.

*Proof*: If class $F_k$ is idle under GPS at time $t$, a new frame that includes packet $p$ will be formed at time $t$. From Lemma 5, under $GPS$, the frame will be served at most at time $t + 2C^k\frac{L_M}{R} \leq t + 2C\frac{L_M}{r_i}$. Hence, from Theorem 2, the frame will be served under $DW^2F^2Q$ before $t + 2C\frac{L_M}{r_i} + (n-1)\frac{L_M}{R} \leq t + (2C+n-1)\frac{L_M}{r_i}$, where $n$ is the number of classes in the system. Thus, there exists $c_1 = 2C + n - 1 = O(C+n)$ such that packet departs before $t + c_1 * \frac{L_M}{r_i}$.

If class $F_k$ is busy under GPS at time $t$, packet $p$ will be included in the next frame. From Lemma 5, $F_{GPS}^p \leq t + 2 * \frac{2L_MC^k}{R} \leq t + \frac{4CL_M}{r_i}$. From Theorem 2, the frame will be served under $DW^2F^2Q$ before $t + 4C\frac{L_M}{r_i} + (n-1)\frac{L_M}{R} \leq t + (4C + n-1)\frac{L_M}{r_i}$. Thus, there exists $c_1 = 4C + n - 1 = O(C+n)$ such that packet $p$ departs before $t + c_1 * \frac{L_M}{r_i}$. $\square$

**Theorem 4 (worst-case fairness)**: $FRR$ has a constant ($O(C+n)$) normalized worst-case fairness index.
*Proof*: Let a packet belonging to flow $f_i \in F_k$ arrive at time $t$, creating a total backlog of $q_i$ bytes in $f_i$'s queue. Let packet $p_1$ be the first packet in the backlog. From the proof of Theorem 3, we have $F_{GPS}^{p_1} \leq t + 4C\frac{L_M}{r_i}$. After the first packet is served under $GPS$, from Lemma 3, at most $\lceil \frac{q_i}{Q_i} \rceil + 3 \leq \frac{q_i}{Q_i} + 4$ frames will be needed to drain the queue. From Lemma 5, under $GPS$, servicing the $\frac{q_i}{Q_i} + 4$ frames will take at most

$$(\frac{q_i}{Q_i} + 4) * 2C^k\frac{L_M}{R} = \frac{q_i}{C^kw_iL_M}\frac{C^kL_M}{R} + 8C^k\frac{L_M}{R} \leq \frac{q_i}{r_i} + 8C\frac{L_M}{r_i}$$

Thus, under $GPS$, the queue will be drained before $t + \frac{q_i}{r_i} + 4C\frac{L_M}{r_i} + 8C\frac{L_M}{r_i}$. From Theorem 2, under $DW^2F^2Q$, the queue will be drained before $t + \frac{q_i}{r_i} + 12C\frac{L_M}{r_i} + (n-1)\frac{L_M}{R}$. Thus, there exists a constant $d = 12C + n - 1 = O(C+n)$ such that the queue will be drained before $t + \frac{q_i}{r_i} + d\frac{L_M}{R}$ and the normalized worst-case fair index for $FRR$ is $\max_i\{\frac{r_i*d\frac{L_M}{r_i}}{R}\} = d\frac{L_M}{R}$. $\square$

The normalized worst-case fair index for $FRR$ is $(12C + n - 1)\frac{L_M}{R}$, which significantly improves that for $STRR$ ($\Omega(N)$). Next we will consider $FRR$'s proportional fairness.

**Lemma 12**: Assuming that $f_i \in F_k$ and $f_j \in F_m$ are continuously backlogged during $[t_1, t_2]$, $k \geq m$. Assume that the inter-class scheduler is $GPS$ and the intra-class scheduler is $LDRRWA$. Let $S_i(t_1, t_2)$ be the services given to flow $f_i$ during $[t_1, t_2]$ and $S_j(t_1, t_2)$ be the services given to flow $f_j$ during $[t_1, t_2]$. There exists two constants $c_1$ and $c_2$ ($c_1 \leq O(C)$ and $c_2 \leq O(C)$) such that

$$|\frac{S_i(t_1,t_2)}{r_i} - \frac{S_j(t_1,t_2)}{r_j}| \leq \frac{c_1 * L_M}{r_i} + \frac{c_2 * L_M}{r_j}.$$

*Proof*: Let $X_k$ and $X_m$ be the smallest numbers of $F_k$ and $F_m$ frames that completely enclose $[t_1, t_2]$. Since $f_i$ and $f_j$ are continuously backlogged during the $[t_1, t_2]$ period, from Lemma 3, the services given to $f_i$ and $f_j$ during this period satisfy:
$$(X_k - 3)Q_i \leq S_i(t_1, t_2) \leq (X_k + 1)Q_i \text{ and}$$
$$(X_m - 3)Q_j \leq S_j(t_1, t_2) \leq (X_m + 1)Q_j.$$
The conclusion follows by manipulating these in-equations and applying Lemma 10, which gives the relation between $X_k$ and $X_m$, $(X_k - 1)C^{k-m} \leq X_m \leq X_kC^{k-m} + 1$.

In the following, we will derive the bound for $\frac{S_i(t_1,t_2)}{r_i} - \frac{S_j(t_1,t_2)}{r_j}$.

$$\frac{S_i(t_1,t_2)}{r_i} - \frac{S_j(t_1,t_2)}{r_j}$$
$$\leq \frac{(X_k+1)Q_i}{r_i} - \frac{(X_m-3)Q_j}{r_j} \leq \frac{Q_iX_k}{r_i} - \frac{Q_jX_m}{r_j} + \frac{Q_i}{r_i} + \frac{3Q_j}{r_j}$$
$$\leq \frac{Q_iX_k}{r_i} - \frac{Q_j(X_k-1)C^{k-m}}{r_j} + \frac{2CL_M}{r_i} + \frac{6CL_M}{r_j}$$
$$= \frac{Q_iX_k}{r_i} - \frac{Q_j(X_k)C^{k-m}}{r_j} + \frac{Q_jC^{k-m}}{r_j} + \frac{2CL_M}{r_i} + \frac{6CL_M}{r_j}$$

We have $\frac{Q_jC^{k-m}}{r_j} = \frac{C^mw_jL_MC^{k-m}}{w_jR} \leq \frac{C*L_M}{r_i}$ and $\frac{Q_iX_k}{r_i} - \frac{Q_j(X_k)C^{k-m}}{r_j} = \frac{C^kw_iL_MX_k}{w_iR} - \frac{C^mw_jL_MX_kC^{k-m}}{w_jR} = 0$. Thus, $\frac{S_i(t_1,t_2)}{r_i} - \frac{S_j(t_1,t_2)}{r_j} \leq \frac{3CL_M}{r_i} + \frac{6CL_M}{r_j}$. The bound for $\frac{S_j(t_1,t_2)}{r_j} - \frac{S_i(t_1,t_2)}{r_i}$ can be derived in a similar fashion. Hence, there exists two constants $c_1$ and $c_2$, $c_1 = O(C)$ and $c_2 = O(C)$, such that $|\frac{S_i(t_1,t_2)}{r_i} - \frac{S_j(t_1,t_2)}{r_j}| \leq \frac{c_1*L_M}{r_i} + \frac{c_2*L_M}{r_j}$. $\square$

Lemma 12 shows that if $GPS$ is used as the inter-class scheduling algorithm, a proportional fairness property is provided. Since $DW^2F^2Q$ closely approximates $GPS$ (Theorem 2), we will show in the next theorem that $FRR$, which uses $DW^2F^2Q$

as the inter-class scheduling algorithm, also supports proportional fairness.

**Theorem 5 (proportional fairness)**: In any time period $[t_1, t_2]$ during which flows $f_i \in F_k$ and $f_j \in F_m$ are continuously backlogged in $FRR$. There exists two constants $c_1 = O(C)$ and $c_2 = O(C)$ such that $|\frac{S_{i,FRR}(t_1,t_2)}{r_i} - \frac{S_{j,FRR}(t_1,t_2)}{r_j}| \leq \frac{c_1 * L_M}{r_i} + \frac{c_2 * L_M}{r_j}$.

*Proof:* See appendix. $\square$

## VI. SIMULATION EXPERIMENTS

We compare $FRR$ with two recently proposed deficit round robin ($DRR$) based schemes, Smoothed Round Robin ($SRR$) [7] and STratified Round Robin ($STRR$) [13]. For reference, we also show the results for two timestamp-based scheduling schemes: $WFQ$ and and $WF^2Q$. All experiments are performed using *ns-2* [11], to which we added $STRR$ and $FRR$ queuing classes. Figure 4 shows the network topology used in the experiments. All links have a bandwidth of $2Mbps$ and a propagation delay of $1ms$. In all experiments, CBR flows have a fixed packet size of 210 bytes, and all other background flows have a fixed packet size uniformly chosen between 128 and 1024 bytes. Except for the experiment summarized in Figure 8 where only CBR and deterministic flows are considered, for all other experiments, we report the results using the confidence interval with a 99% confidence level. The confidence intervals are obtained by running each simulation 50 times with different random seeds and computing from the 50 samples.



Fig. 4. Simulated network.

Figure 5 shows the average end-to-end delays for flows with different rates. Figure 6 shows the worst-case end-to-end delays. In the experiment, there are 10 $CBR$ flows from $S0$ to $R0$ with average rates of $10Kbps$, $20Kbps$, $40Kbps$, $60Kbps$, $80Kbps$, $100Kbps$, $120Kbps$, $160Kbps$, $200Kbps$, and $260Kbps$. The average packet delays of these CBR flows are measured. The background traffic in the system is as follows. There are five exponential on/off flows from $S1$ to $R1$ with rates $60Kbps$, $80Kbps$, $100Kbps$, $120Kbps$, and $160Kbps$. The on-time and the off-time are 0.5 second. There are five Pareto on/off flows from $S2$ to $R2$ with rates $60Kbps$, $80Kbps$, $100Kbps$, $120Kbps$, and $160Kbps$. The on-time and the off-time are 0.5 second. The shape parameter of the Pareto flows is 1.5. Two $7.8Kbps$ FTP flows with infinite traffic are also in the system, one from $S1$ to $R1$ and the other one from $S2$ to $R2$.

In this experiment, all of the three deficit round-robin based schemes, $SRR$, $STRR$, and $FRR$, give reasonable approximation of the timestamp based schemes, $WFQ$ and $WF^2Q$, for all the flows with different rates. In comparison to $SRR$ and $STRR$, $FRR$ achieves average and worst-case end-to-end delays that are closer to the ones with the timestamp based schemes for flows with large rates ($\geq 150Kbps$ in the experiment). In this

experiment, $FRR$ have smaller average end-to-end delays than $SRR$ and $STRR$ for flows whose rates are larger than $40Kbps$, while having larger average packet delays for other flows. In $FRR$, the timestamp based inter-class scheduling mechanism is added on top of $DRR$ so that flows with small rates do not significantly affect flows with large rates. Thus, in a way, $FRR$ gives preference to flows with larger weights in comparison to other $DRR$ bases schemes: the average packet delay with $FRR$ is more proportional to the flow rate than that with $SRR$ and $STRR$.

Figure 7 (a), (b), and (c) shows the short-term throughput achieved by different schemes. Since the results for $SRR$ are very similar to those for $STRR$, we only show the results for $STRR$, similarly, results for $WF^2Q$ are similar to the results for $WFQ$. In this experiment, the short-term throughput in an interval of 100ms is measured. We observe one $300Kbps$ $CBR$ flow and one $600Kbps$ flow from $S0$ to $R0$. In addition, 50 $10Kbps$ $CBR$ flows from $S0$ and $R0$ are introduced. Other background flows are the same as the previous experiment.

Figure 7 (a), (b), and (c) shows the results for the 300 $Kbps$ flow. The results for the 600 $Kbps$ flow have a similar trend. As can be seen from the figure, the short term throughput for $STRR$ (and $SRR$) exhibit heavy fluctuations. On the other hand, $WFQ$ and $FRR$ yield much better short term throughput: within each interval of 100ms, the throughput is always close to the ideal rate. This demonstrates that $FRR$ has a much better short-term fairness property than $SRR$ and $STRR$.

Figure 8 shows the proportional fairness of $FRR$. In this experiment, we observe four deterministic flows from $S0$ to $R0$ with average rates of 100Kbps, 200Kbps, 200Kbps, and 300Kbps. These flows follow an off/on pattern with each off/on period being 10 seconds. Hence, the flows are quiet for 10 seconds and then send in a doubled rate for the next 10 seconds. One 600Kbps CBR flow from $S1$ to $R1$ is introduced in period [10s, 16s] and another 400Kbps CBR flow from $S1$ to $R1$ is introduced in period [12s, 14s]. The CBR flows and the observed flows share the link from $N1$ to $N2$. The bandwidth allocation in the link from $N1$ to $N2$ to each of the flows during period [10s, 19s] is showed in Figure 8. As can be see from the figure, for all periods with different traffics sharing the link, the bandwidth allocation to the four observed flows is exactly proportional to their reserved bandwidths.

The last experiment is designed to study the impacts of $C$, a parameter in $FRR$. The background traffics used in this experiment are the same as those in Figure 5. We observe the worst case end-to-end packet delay for 16 $CBR$ flows from $S0$ to $R0$ with average rates of $10Kbps$, $20Kbps$, $30Kbps$, $40Kbps$, $50Kbps$, $60Kbps$, $70Kbps$, $80Kbps$, $90Kbps$, $100Kbps$, $110Kbps$, $120Kbps$, $130Kbps$, $140Kbps$, $150Kbps$, and $160Kbps$. When $C = 8$, there are two classes in the system, $F_3$ containing flows with rates $10Kbps$, $20Kbps$, and $30Kbps$, and $F_2$ containing the rest of the flows. When $C = 4$, there are three classes in the system, $F_4$ ($10Kbps$ to $30Kbps$), $F_3$ ($40Kbps$ to $120Kbps$), and $F_2$ ($130Kbps$ to $160Kbps$). When $C = 2$, there are 5 classes: $F_8$ ($10Kbps$), $F_7$ ($20Kbps$ and $30Kbps$), $F_6$ ($40Kbps$ to $60Kbps$), $F_5$ ($70Kbps$ to $120Kbps$), and $F_4$ ($130Kbps$ to $160Kbps$).

Figure 9 shows the worst case delay in milli-seconds. We can see that the worst case delay for flows within one class are similar, which is evidenced by the ladder shape curves in the figure. This is expected as the $DRR$ based scheme is used for intra-class

Fig. 5.    Average end-to-end delay.



Fig. 6.    Worst-case end-to-end delay



(a) $STRR$            (b) $FRR$            (c) $WFQ$

Fig. 7.    Short-term throughput



Fig. 8.    Proportional fairness



Fig. 9.    Worst-case end-to-end delay.

Fig. 10.    Delay in terms of numbers of packets

scheduling. The packet delay is directly related to $C$. A smaller $C$ value results in a smaller worst case packet delay.

Figure 10 shows a different view of Figure 9. In this figure, we represent the absolute worst case delay time as the time to send a number of packets (packets are of the same size, $210B$, in this experiment). This allows the delay to be normalized by the flow rate. There are two interesting observations in Figure 10.

First, within each class, $FRR$ biases against flows with larger weights. This is due to the use of a $DRR$ based scheme for intra-class scheduling. Biasing against flows with large weights is a common problem for all $DRR$ based schemes. However, in $FRR$, this problem is limited since the weight difference within a class is bounded. Second, $FRR$ treats different classes fairly. It can be seen that for flows in different classes, the worst case

packet delays fall in ranges with similar lower bounds and upper bounds as shown in the seesaw shape curve (e.g. when $C = 2$).

## VII. CONCLUSION

In this paper, we describe Fair Round Robin ($FRR$), a low quasi-$O(1)$ complexity round robin scheduler that provides proportional and worst-case fairness. In comparison to other $DRR$ based scheduling schemes, $FRR$ has similar complexity and proportional fairness, but better worst-case fairness. The simulation study demonstrates that even in average cases, $FRR$ has better short-term behavior than other $DRR$ based schemes, including smoothed round robin and stratified round robin. The constant factors in the complexity and QoS performance bounds for $FRR$ are still fairly large. Recent improvements on $GPS$ tracking [19], [20], [23] and $DRR$ implementations [8] may be applied to improve $FRR$.

## ACKNOWLEDGEMENT

## REFERENCES

[1] J. Bennett and H. Zhang, "Hierarchical Packet Fair Queueing Algorithms," *ACM/IEEE Trans. on Networking*, 5(5):675-689, Oct. 1997.

[2] J. Bennett and H. Zhang, "$WF^2Q$: Worst Case Fair Weighted Fair Queuing", in *IEEE INFOCOM'96* (1996), pages 120-128.

[3] B. Caprita, J. Nieh, and W. Chan, "Group Round Robin: Improving the Fairness and Complexity of Packet Scheduling." *Proceedings of the 2005 Symposium on Architecture for Networking and Communications Systems* (ANCS'05), pages 29-40, Princeton, New Jersey, 2005.

[4] S.Cheung and C. Pencea, "BSFQ: Bin Sort Fair Queuing," in *IEEE INFOCOM'02* (2002), pages 1640-1649.

[5] A. Demers, S. Keshav, and S. Shenker, "Analysis and Simulation of a Fair Queuing Algorithm," in *ACM SIGCOMM'89* (1989), pages 1-12.

[6] S. Golestani, "A Self-clocked Fair Queueing Scheme for Broadband Applications", in *IEEE INFOCOM'94* (1994), pages 636-646.

[7] C. Guo, "SRR, an O(1) Time Complexity Packet Scheduler for Flows in MultiService Packet Networks," *IEEE/ACM Trans. on Networking*, 12(6):1144-1155, Dec. 2004.

[8] L. Lenzini, E. Mingozzi, and G. Stea, "Aliquem: a Novel DRR Implementation to Achieve Better Latency and Fairness at O(1) Complexity," in *IWQoS'02* (2002), pages 77-86.

[9] L. Lenzini, E. Mingozzi, and G. Stea, "Tradeoffs Between Low Complexity, Low Latency, and Fairness with Deficit Round-Robin Schedulers." *IEEE/ACM Trans. on Networking*, 12(4):681-693, April 2004.

[10] L. Massouli and J. Roberts, "Bandwidth Sharing: Objectives and Algorithms," *IEEE/ACM Trans. on Networking*, 10(3):320-328, June 2002.

[11] "The Network Simulator - ns-2," available at http://www.isi.edu/nsnam/ns.

[12] A. Parekh and R. Gallager, "A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: the Single Node Case," *IEEE/ACM Trans. on Networking*, 1(3):344-357, June 1993.

[13] S. Ramabhadran and J. Pasquale, "The Stratified Round Robin Scheduler: Design, Analysis, and Implementation." *IEEE/ACM Transactions on Networking*, 14(6):1362-1373, Dec. 2006.

[14] J. Rexford and A. Greenberg and F. Bonomi, "Hardware-Efficient Fair Queueing Architectures for High-Speed Networks," *IEEE INFOCOM'96*, (1996), pages 638-646.

[15] J. L. Rexford, F. Bonomi, A. Greenberg, A. Wong, "A Scalable Architecture for Fair Leaky-Bucket Shaping," in *INFOCOM'97* (1997), pages 1056–1064.

[16] D. Stiliadis and A. Varma, "Design and Analysis of Frame-based Fair Queueing: A New Traffic Scheduling Algorithm for Packet-Switched Networks," in *ACM SIGMETRICS'96* (1996), pages 104-115.

[17] M. Shreedhar and G. Varghese, "Efficient Fair Queuing using Deficit Round Robin," in *ACM SIGCOMM'95* (1995), pages 231-242.

[18] S. Suri, G. Varghese, and G Chandranmenon, "Leap Forward Virtual Clock: An $O(loglog\ N)$ Queuing Scheme with Guaranteed Delays and Throughput Fairness," in *IEEE INFOCOM'97* (1997), pages 557-565.

[19] P. Valente, "Exact GPS Simulation with Logarithmic Complexity, and its Application to an Optimally Fair Scheduler," in *ACM SIGCOMM'04* (2004), pages 269-280.

[20] J. Xu and R. J. Lipton, "On Fundamental Tradeoffs between Delay Bounds and Computational Complexity in Packet Scheduling Algorithms," in *ACM SIGCOMM'02* (2002), pages 279-292.

[21] X. Yuan and Z. Duan, "Fair Round Robin: A Low Complexity Packet Scheduler with Proportional and Worst-Case Fairness," *Technical Report TR-080201*, Department of Computer Science, Florida State University, Feb. 2008.

[22] L. Zhang, "Virtual Clock: A New Traffic Control Scheme for Packet Switching Networks", in *ACM SIGCOMM'90* (1990), pages 19-29.

[23] Q. Zhao and J. Xu, "On the Computational Complexity of Maintaining GPS Clock in Packet Scheduling," in *IEEE INFOCOM'04* (2004), pages 2383-2392.

## APPENDIX: PROOFS

**Proof of Lemma 1:** Since $X$ is the smallest number of continuous DRR rounds that completely enclose $[t_1, t_2)$, $f_i$ is served in at least $X-2$ rounds. Thus, $S_{i,DRR}(t_1, t_2) \geq (X-2)*quantum_i - L_M \geq (X-3)quantum_i$ since we assume that $quantum_i \geq L_M$. On the other hand, $f_i$ is served in at most all X rounds, in this case, the total number of data sent should be less than the total quantum generated during the rounds plus the left over from the previous DRR round, which is less than $L_M$. Thus, $S_{i,DRR}(t_1, t_2) \leq X * quantum_i + L_M \leq (X+1)quantum_i$. □

**Proof of Lemma 2**: Since $N * r_{min} \leq r_1 + r_2 + ... + r_N \leq R$, $r_{min} \leq \frac{R}{N}$. $quantum_i = L_M * \frac{r_i}{r_{min}} \leq D * L_M$. Thus, the total size of a round is at most $\sum_{i=1}^{N}\{quantum_i + L_M\} \leq (D+1) * N * L_M$. The time to complete service in a round is at most $\frac{(D+1)N*L_M}{R} \leq (D+1) * \frac{L_M}{r_{min}} \leq (D+1) * \frac{L_M}{r_{max}/D} = D(D+1) * \frac{L_M}{r_{max}} \leq D(D+1) * \frac{L_M}{r_i}$.

Packet $p$ arrives at the head of the queue for $f_i$ time $t$. It takes at most two rounds for the packet to be serviced. There exists a constant $c_1 = 2 * D(D+1) = O(D^2)$ such that packet $p$ will be serviced before $t + c_1 \times \frac{L_M}{r_i}$. This proves the first statement. Next, we will prove the second statement.

Let a packet belonging to flow $f_i$ arrives at time $t$, creating a total backlog of size $q_i$ in $f_i$'s queue. From statement 1., there exists a constant $c_1$ such that the first packet in the queue will be serviced in $t + c_1 \times \frac{L_M}{r_i}$. After the first packet is serviced, there will be at most $\lceil \frac{q_i}{quantum_i} \rceil + 1 \leq \frac{q_i}{quantum_i} + 2$ rounds for the $q_i$ data to be sent. During the $\frac{q_i}{quantum_i} + 2$ rounds, at most $(\frac{q_i}{quantum_i} + 2) * \sum_{j=1}^{N} quantum_j$ quanta are generated, and thus, at most $(\frac{q_i}{quantum_i} + 2) * \sum_{j=1}^{N} quantum_j + N * L_M$ data are sent since each flow can have at most $L_M$ credits left from the previous round. Thus, the total time to complete the $\frac{q_i}{quantum_i} + 2$ rounds is at most

$$\frac{(\frac{q_i}{quantum_i}+2)*\sum_{j=1}^{N} quantum_j + N*L_M}{R_N}$$

$$= (\frac{q_i}{quantum_i} + 2)\frac{\sum_{j=1}^{N} quantum_j}{R} + \frac{N*L_M}{R}$$

$$= (\frac{q_i}{quantum_i} + 2)\frac{\sum_{j=1}^{N} L_M*\frac{r_j}{r_{min}}}{R} + \frac{N*L_M}{R}$$

$$= (\frac{q_i}{quantum_i} + 2)\frac{L_M}{r_{min}}\frac{\sum_{j=1}^{N} r_j}{R} + \frac{N*L_M}{R}$$

$$\leq (\frac{q_i}{quantum_i} + 2)\frac{L_M}{r_{min}} + \frac{N*L_M}{R} \leq (\frac{q_i}{L_M*\frac{r_i}{r_{min}}})\frac{L_M}{r_{min}} + \frac{3L_M}{r_{min}}$$

$$\leq \frac{q_i}{r_i} + \frac{3*L_M}{r_{min}} \leq \frac{q_i}{r_i} + \frac{3*L_M}{r_{max}/D} \leq \frac{q_i}{r_i} + \frac{3D*L_M}{r_i}$$

Thus, there exists a constant $c_2 = c_1 + 3D = O(D^2)$ such that the queue of size $q_i$ will be sent before $t + \frac{q_i}{r_i} + c_2 * \frac{L_M}{r_i}$. This is true for all flows. The normalzied worst case fair index is $c_{DRR} = max_i\{\frac{r_iC_{i,DRR}}{R}\} = \frac{c_2 L_M}{R}$. This proves the second statement.

For any given time period, $[t_1, t_2]$, let $f_i$ and $f_j$ be backlogged during this period that is enclosed by $X$ rounds. From Lemma 1, we have

$$(X-3)quantum_i \leq S_{i,DRR}(t_1, t_2) \leq (X+1)quantum_i$$
$$(X-3)quantum_j \leq S_{j,DRR}(t_1, t_2) \leq (X+1)quantum_j$$

By manipulating these inequations, it can be shown that there exist two constants $c_1 = c_2 = 4D = O(D)$, such that $|\frac{S_{i,DRR}(t_1,t_2)}{r_i} - \frac{S_{j,DRR}(t_1,t_2)}{r_j}| \leq c_1 \frac{L_M}{r_i} + c_2 \frac{L_M}{r_j}$. □

**Proof of Lemma 3:** The notation $S_{i,LDRRWA}(t_1, t_2)$ is abused in this lemma since $LDRRWA$ does not decide the actual timing to service packets. In this lemma, $S_{i,LDRRWA}(t_1, t_2)$ denotes the amount of data for a continuously backlogged flow $f_i$ in $X$ continuous $LDRRWA$ frames (of a particular class) using any inter-class scheduling scheme.

Since $f_i$ is continuously backlogged, it will try to send as many packets as possible in each frame. Since $X$ frames enclose $[t_1, t_2]$, flow $f_i$ will fully utilize at least $X - 2$ frames (all but the first frame and the last frame). In the $X - 2$ frames, $(X-2) \times Q_i$ credits are generated for flow $f_i$. The lookahead operation in the frame prior to the $X - 2$ frames may borrow at most one packet, whose size is less than $L_M$, from $f_i$ in the first of the $X - 2$ frames and flow $f_i$ in the last of the $X - 2$ frames may pass at most $L_M$ credits to the next frame. Note that the lookahead operation borrows at most one packet from each backlogged flow. Thus, $S_{i,LDRRWA}(t_1, t_2) \geq (X - 2) \times Q_i - L_M - L_M$. Since $Q_i \geq 2L_M$, $S_{i,LDRRWA}(t_1, t_2) \geq (X - 3) \times Q_i$.

On the other hand, $f_i$ will be served in at most all the $X$ frames, which produces $X \times Q_i$ credits for $f_i$ during this period of time. Flow $f_i$ in the frame prior to the $X$ frames may have at most $L_M$ left-over credits and the lookahead operation in the last of the $X$ frames may borrow at most $L_M$ credits from $f_i$ in the next frame. Thus,

$$S_{i,LDRRWA}(t_1, t_2) \leq X \times Q_i + L_M + L_M \leq (X+1)Q_i.□$$

**Proof of Lemma 10:** This lemma relaxes the condition in Lemma 10 by not requiring each class to be serviced with its GPS guaranteed rate. Since $f_i \in F_k$ and $f_j \in F_m$ be continuously backlogged during $[t_1, t_2]$, the sizes of all frames during this period are no smaller than $L_M$ (Lemma 9). Let us partition the duration $[t_1, t_2]$ into smaller intervals $[a_1 = t_1, b_1), [a_2 = b_1, b_2),$ ..., $[a_Y = b_{Y-1}, b_Y = t_2)$ such that within each interval $[a_h, b_h)$, $1 \leq h \leq Y$, the weights of all classes are fixed. Let $F_1, ..., F_n$ be the $n$ classes in the system. Let class $F_k$ have weight $w_k^h$ during interval $[a_h, b_h)$, $1 \leq h \leq Y$ (If $F_k$ is not backlogged, $w_k^h = 0$).

The amount of class $F_k$ data sent during $[a_h, b_h)$ is thus,

$$\frac{w_k^h}{\sum_{i=1}^{n} w_i^h} R * (b_h - a_h).$$

Consider a reference scheduling system that contains three classes $RF_k$, $RF_m$, and $RF_o$. Let us use intervals $[aa_1 = t_1, bb_1),$ $[aa_2 = bb_1, bb_2),$ ..., $[aa_Y = bb_{Y-1}, bb_Y)$ to emulate the behavior of classes $F_k$ and $F_m$ during intervals $[a_1 = t_1, b_1),$ $[a_2 = b_1, b_2),$ ..., $[a_Y = b_{Y-1}, b_Y)$ respectively. Let $rw_k^h$ be the weight for class $RF_k$ during interval $[aa_h, bb_h)$, $1 \leq h \leq Y$. Let $rw_m^h$ be the weight for class $RF_m$ during interval $[aa_h, bb_h)$, $1 \leq h \leq Y$. Let $rw_o^h$ be the weight for class $RF_o$ during interval $[aa_h, bb_h)$, $1 \leq h \leq Y$. The weights and the duration of each interval are given as follows:

$$rw_k^h = w_k^h, rw_m^h = w_m^h, rw_o^h = 1 - w_k^h - w_m^h, 1 \leq h \leq Y$$

and

$$bb_h = aa_h + \frac{b_h - a_h}{\sum_{i=1}^{n} w_i^h}, 1 \leq h \leq Y.$$

It can be verified that the amount of classes $RF_k$ and $RF_m$ data sent in an interval $[aa_h, bb_h)$, $1 \leq h \leq Y$, is exactly the same as the amount of classes $F_k$ and $F_m$ data sent in an interval $[a_h, b_h)$, $1 \leq h \leq Y$, respectively. In an interval $[aa_h, bb_h)$, $1 \leq h \leq Y$, let us further assume that Class $RF_k$ has exactly the same sequence of packets as Class $F_k$ has in interval $[a_h, b_h)$ and that Class $RF_m$ has exactly the same sequence of packets as Class $F_m$ has in interval $[a_h, b_h)$. The progress of classes $F_k$ and $F_m$ during $[t_1, t_2]$ is exactly the same as the progress of class $RF_k$ and $RF_m$ during $[aa_1, bb_Y)$

In the reference system, classes $RF_m$ and $RF_k$ are serviced with the $GPS$ guaranteed rate during $[aa_1, bb_Y)$. Let $RX_k$ and $RX_m$ be the smallest numbers of $RF_k$ and $RF_m$ frames that completely enclose $[aa_1, bb_Y)$. From Lemma 10, $(RX_k - 1)C^{k-m} \leq RX_m \leq RX_kC^{k-m} + 1$.

Let $X_k$ and $X_m$ be the smallest number of $F_k$ and $F_m$ frames that completely enclose $[t_1, t_2]$. Since the progress of classes $F_k$ and $F_m$ during $[t_1, t_2]$ is exactly the same as the progress of class $RF_k$ and $RF_m$ during $[aa_1, bb_Y)$, we have $X_k = RX_k$ and $X_m = RX_m$. Thus, $(X_k - 1)C^{k-m} \leq X_m \leq X_kC^{k-m} + 1$. □

**Proof of Theorem 1:** Since $GPS$ is work-conserving, we will prove the theorem by showing that $DW^2F^2Q$ has the same idle and busy periods as $GPS$. Assuming that $DW^2F^2Q$ and $GPS$ have different idle and busy periods. Let $t$ be the first occurrence when $GPS$ and $DW^2F^2Q$ are not in the same state. There are two cases.

Case 1: $GPS$ is idle and $DW^2F^2Q$ is busy, serving packet $p$. Since $t$ is the first occurrence when $GPS$ and $DW^2F^2Q$ are not in the same state, the amount of data served during $[0, t)$ must be the same for the two scheduling schemes. Since $p$ is currently being served under $DW^2F^2Q$, $p$ must be started before $t$ under $GPS$. Since $GPS$ is idle at time $t$, packet $p$ must finish before $t$ under $GPS$. Hence, there must exist a packet $q$ such that $q$ has not been served under $GPS$ during $[0, t)$ and has been served by $DW^2F^2Q$ during $[0, t)$. Since $GPS$ is idle at $t$, packet $q$ should start after $t$ under $GPS$, which indicates that $q$ cannot be served under $DW^2F^2Q$ during $[0, t)$. This is the contradiction.

Case 2: $GPS$ is busy and $DW^2F^2Q$ is idle. Let packets $p_1, p_2, ..., p_i$ be the packets departed under $GPS$ during $[0, t)$ and packets $cp_1, ..., cp_j$ be the packets currently in progress under $GPS$. Since $GPS$ is busy, at least one packet is being serviced