Fair Efficiency, or Low Average Delay without Starvation

Sergey Gorinsky and Christoph Jechlitschek

{gorinsky,chrisj}@arl.wustl.edu

Applied Research Laboratory, Department of Computer Science and Engineering, Washington University One Brookings Drive, Campus Box 1045, Saint Louis, Missouri 63130-4899, United States of America

Abstract-File transfer, web download, and many other applications are primarily interested in minimal delay achievable for their messages. In this paper, we investigate allocating the bottleneck link capacity to transmit messages efficiently but fairly. While SRPT (Shortest Remaining Processing Time) is an optimally efficient algorithm that minimizes average delay of messages, large messages might starve under SRPT in heavy load conditions. PS (Processor Sharing) and ViFi (Virtual Finish Time First) are fair but yield higher average delays than under SRPT. We explore the class of fair algorithms further and prove that no online algorithm in this class is optimally efficient. Then, we derive a fair algorithm SFS (Shortest Fair Sojourn) and report extensive experimental evidence that SFS is consistently more efficient than PS and even ViFi during either temporal overload or steady-state operation, with largest benefits achieved when average load is around the bottleneck link capacity. Furthermore, average delay under the fair SFS remains close to the minimum attained under the unfair SRPT.

I. INTRODUCTION

Distributed applications communicate by sending messages over networks which might use multiple packets to deliver one message. For example, thousands or even millions of packets might carry a long message of web download or file transfer. While such applications are primarily interested in minimal delays for their messages, the network ability to minimize the delays is constrained chiefly by the path capacity from the source to the destination. Specifically, the capacity of a bottleneck link on the path is the main factor determining minimal achievable delays. This paper investigates fair but efficient algorithms for allocating the bottleneck link capacity to minimize message delays.

Shortest Remaining Processing Time (SRPT) schedules messages preemptively in the order of their remaining transmission delays and is optimally efficient [1]. However, the minimal average delay comes at the expense of potential unfairness: in some settings with heavy load, SRPT starves large messages by delaying them without bound [2].

Processor Sharing (PS) is an alternative classic algorithm that instantaneously allocates equal shares of the bottleneck capacity to all pending messages [3]. Consequently, expected delay of a message under PS is proportional to the message size. Also, since PS does not rely on knowledge of message sizes, PS lends itself nicely to implementation in layered network designs. Due to the above reasons, PS has become a traditional ideal in network capacity allocation. Although packet-switching networks do not support instantaneous sharing of a link, a lot of research has been conducted on packet transmission algorithms that approximate the PS ideal. Packetgrained approximations of PS include Weighted Fair Queuing (WFQ) [4], Deficit Round Robin (DRR) [5], and other algorithms for fair queuing at routers as well as fair end-toend congestion control schemes exemplified by Transmission Control Protocol (TCP) [6].

While SRPT is unfair, PS achieves fairness by sacrificing efficiency: average delay of messages under PS is significantly higher. Recent studies reveal remarkable existence of algorithms that have it both ways and combine fairness with SRPT-like efficiency. Virtual Finish Time First (ViFi) is a specific efficient representative of the fair algorithmic class where no message is delayed longer than under PS [7]. ViFi schedules messages preemptively in the order of their finish times under PS and is independently proposed as Fair Sojourn Protocol (FSP) in the context of web servers [8]. Significant reductions in average delay under ViFi versus PS are substantiated both experimentally and analytically [7]–[9].

This paper sheds more light on the class of fair algorithms for network capacity allocation. First, we show that the fair class does not contain an optimally efficient online algorithm. Then, we develop Shortest Fair Sojourn (SFS), a fair algorithm with even lower average delay than under ViFi in most settings. Our extensive simulations over a wide range of network load illustrate efficiency and fairness properties of SFS, ViFi, PS, and SRPT. In particular, we demonstrate that average delay under SFS versus ViFi is consistently lower over the whole range of the experiments.

The rest of the paper is structured as follows. Section II clarifies our model, metrics, and terminology. Section III rules out existence of an optimally efficient algorithm in the fair class. This section also presents SFS and proves its fairness. Section IV reports the experimental comparison of SFS, ViFi, PS, and SRPT. Finally, Section V sums up our findings and discusses future work.

II. MODEL, TERMINOLOGY, AND METRICS

We define a *message* as an atomic data unit meaningful for an application. Messages arrive for network transfer in their entirety. *Delay* of a message is time passed from the message arrival until the whole message reaches its destination. Related studies refer to delay under other names such as transfer time, response time, flow time, or sojourn time. *Transmission delay* of a message represents its communication needs and equals $\frac{S}{C}$, where S is the message size, and C is the capacity of the network bottleneck link shared with all the other messages. Transmission delay is also known as processing time, e.g., as reflected in the name of SRPT. We assume that the communications utilize the entire bottleneck link capacity and experience negligible propagation, node processing, and error recovery delays. Then, besides transmission delay, the only other component of message delay is due to waiting for the bottleneck link to become available.

Arrival times and transmission delays of messages characterize network load. Network service is represented by an algorithm that allocates the bottleneck link capacity to pending messages. In particular, we are interested in *online algorithms* that have no access to information about future messages. Capacity allocation enjoys ideal flexibility that allows both instantaneous link sharing and instantaneous transmission preemption.

Since PS has become synonymous with fairness in network resource allocation, we rely on delays of individual messages under PS as a basis for defining the fair algorithmic class:

Definition 1: Starvation is a scenario where a message finishes later than under PS. An algorithm for capacity allocation is *fair* if and only if no starvation occurs under the algorithm for any network load.

To quantify fairness of an algorithm to a particular message, we introduce a metric of *starvation stretch*:

Definition 2: Starvation stretch $s_x(m)$ of message m under algorithm X is the ratio of message delay $d_x(m)$ under algorithm X to message delay $d_{PS}(m)$ under PS:

$$s_{\rm x}(m) = \frac{d_{\rm x}(m)}{d_{\rm PS}(m)}.$$
(1)

Note that algorithm X is deemed unfair if there exists network load where $s_x(m) > 1$ for at least one message m.

Also note an implicit assumption that network capacity is allocated among messages. We strongly believe that fairness of capacity allocation should be defined with respect to realworld entities, rather than messages or packet flows as in traditional networking. However, since the important "among what" aspect is orthogonal to our main contributions and requires a separate thorough treatment, we do not explore it further in this paper.

To quantify efficiency of network capacity allocation under algorithm X, we measure average delay D_x for all n messages in imposed network load:

$$D_{\rm x} = \frac{\sum\limits_{m=1}^{n} d_{\rm x}(m)}{n}.$$
 (2)

Because SRPT is an optimally efficient algorithm if fairness concerns are put aside, we use average delay under SRPT as a baseline for assessing efficiency of fair algorithms: Definition 3: Average letup L_x under algorithm X is the ratio of average delay D_x under algorithm X to average delay D_{SRPT} under SRPT:

$$L_{\rm X} = \frac{D_{\rm X}}{D_{\rm SRPT}}.$$
(3)

Although a fair algorithm is not always able to match the ideal efficiency of the unfair SRPT, consistent closeness of average letup L_x to 1 is an indicator that fair algorithm X is highly efficient.

III. IMPROVING ON VIFI

While the fair ViFi provides significantly lower average delay than the fair PS [7], one might wonder whether ViFi or any other online algorithm in the fair class minimizes average delay. A simple counterexample proves the inverse theorem:

Theorem 1 (No Optimal Online Algorithm): No online algorithm minimizes average delay without starvation.

Proof: First, let us consider the 12-message network load described in Figure 1. Since SRPT causes no starvation for this load, it is optimal to transmit the messages in an SRPT order: any permutation of 1 through 9 followed by 11, 10, and 12. Figure 2a depicts an optimal schedule with average delay $\frac{618}{12}$. Note that average delay under ViFi, or any other algorithm that transmits message 10 before messages 11 and 12, is higher and equals $\frac{622}{12}$.

Now, let us suppose that message 13 with transmission delay 4 arrives at time 105. By time 105 in the above optimal schedule, messages 1 through 9 and 11 have finished, and message 10 is being transmitted since time 100. Suspending message 10 to transfer message 13 would lead to starvation because the resumed message 10 would complete at time 118 whereas Figure 2b shows that message 10 finishes under PS at time 117. Hence, message 10 has to finish before messages 13 and 12 are transmitted. The resulting schedule provides average delay $\frac{635}{13}$. However, smaller average delay $\frac{634}{13}$ would be achieved if all 13 messages were instead transmitted in a ViFi order: any permutation of 1 through 9 followed by 10, 11 (suspended at time 105 to transfer message 13), 13, 11 (rest of it), and 12. Figure 2c depicts such an optimal ViFi schedule.

Therefore, since optimality of transmitting message 11 (rather than message 10) at time 90 depends on whether message 13 arrives at time 105, there is no optimal online algorithm that minimizes average delay without starvation. ■

Although Theorem 1 precludes existence of a fair online algorithm that minimizes average delay, some algorithms might outperform ViFi with respect to average delay in most

Message	Arrival time	Transmission delay
1 through 9	0	10 each
10	0	14
11	90	10
12	90	20

Fig. 1. Network load in the proof of Theorem 1.



Fig. 2. No online algorithm in the fair class minimizes average delay: messages 11 and 13 are denoted with solid black and stripes respectively.

settings. To design such algorithms, we rely on two insights. First, selecting a message with shorter remaining transmission delay yields lower average delay. Second, while ViFi always transmits messages in the order of their PS finish times, deviation from the PS order to transmit a message with shorter remaining transmission delay does not necessarily lead to starvation.

Shortest Fair Sojourn (SFS) is an algorithm based on the above insights. Let M be a pending message with the shortest remaining transmission delay. If transmitting M first and then scheduling the other pending messages in a ViFi order avoid starvation, SFS transmits M. Otherwise, SFS transmits a message with the smallest PS finish time.

Theorem 2 (Fairness of SFS): SFS is a fair algorithm.

Proof: Messages might starve only because arrival of a new message N increases their SFS finish times past their PS finish times. If N or another message is chosen for transmission despite not having the smallest PS finish time (i.e., the chosen message has the shortest remaining transmission delay), then the updated schedule avoids starvation by SFS definition. Otherwise, SFS schedules N and the other pending messages in a ViFi order. The updated schedule is such that:

- Message N completes by its PS finish time due to fairness of ViFi [7].
- Messages with smaller PS finish times than N's are scheduled before N and hence complete even before their prior PS finish times.
- Messages with the same or larger PS finish times than N's have their PS finish times postponed by the transmission delay of N and hence complete under SFS at or before their extended PS finish times.

In all scenarios, no message starves. SFS is fair.

We also consider SFS+, a computationally more extensive variant of SFS. SFS+ sorts pending messages in the increasing order of their remaining transmission delays and selects the first message F (up to a message with the smallest PS finish time) such that transmitting F and scheduling the other pending messages in a ViFi order avoid starvation. Our experiments show that SFS and SFS+ yield similar average delays.

Moreover, average delay under SFS+ exhibits a surprising tendency of being slightly higher. Hence, the next section compares ViFi, PS, and SRPT with SFS only.

IV. PERFORMANCE EVALUATION

A. Experimental methodology

We simulate transmission of 3,000 messages over a link with capacity C = 10 Tbps under full link utilization and no data loss. Our choices for the link capacity and traffic are dictated by our desire to model high-speed networks of the future. Message sizes and arrival times are drawn from random distributions. For each set of the traffic settings, we repeat the experiment under SRPT, PS, ViFi, and SFS.

To characterize intensity of the traffic, we define a notion of *load* l as

$$l = \frac{m}{C \cdot t} \tag{4}$$

where *m* denotes the average message size, and *t* is the average message interarrival time. Since the number of messages is finite, all message delays remain finite even with l > 100%. This feature of our experimental setup enables us to compare the evaluated algorithms under long-term overload conditions, which is impossible with analytical techniques that target only steady-state algorithmic behaviors.

We experiment with both uniform and Pareto distributions of message sizes. Unless explicitly stated otherwise, we report results for message sizes that are uniformly distributed between 100 GB and 100 TB. For Pareto-sized messages, our experiments confirm the surprising finding by Bansal and Harchol-Balter [10] that even under heavy loads, SRPT starves only few messages and does not inflate starvation stretches much beyond 1. In these settings, the provenly fair ViFi and SFS still provide far superior efficiency in comparison to PS. To illustrate the efficiency gain, we also plot average letups under PS, ViFi, and SFS in scenarios where message sizes are drawn from the Pareto distribution with index 1.5 and minimum message size of 500 GB. In all our experiments, messages arrive according to a Poisson process with such an average rate that yields a desired value of load.



Fig. 3. Fairness of ViFi and SFS versus unfairness of SRPT for uniformly distributed message sizes at 95% load: (a) starvation stretch for every seventh message in a single experiment; (b) average of starvation stretches over 1,000 experiments: the rich get richer under ViFi or SFS, and so do the poor; SFS helps the middle class more; (c) cumulative distribution of individual starvation stretches from 100 experiments.

The code and running instructions for all the reported simulations are available at our web site [11].

B. Fairness: lower delays for all

Figure 3a plots the starvation stretch for every seventh of all 3,000 (ordered by size) messages in a single experiment at 95% load. Under either of SRPT, ViFi or SFS, a small fraction of messages across the whole spectrum of message sizes has starvation stretch 1. These messages finish at exactly the same times as under PS because they conclude a traffic burst by emptying the queue upon their completion (under both PS and SRPT, ViFi or SFS). Small and even midsize messages benefit significantly from SRPT, which delivers them up to 50 times faster than under PS. However, some large messages starve under SRPT. For example, delay for the least lucky message under the unfair SRPT is about 50 times larger than under PS.

For the fair ViFi, Figure 3a shows that 800 smallest messages enjoy similarly low starvation stretches as under SRPT. To explain the similarity, we observe that a small message is likely to possess both the shortest remaining transmission delay and earliest PS finish time among pending messages. For larger messages, the ViFi profile becomes different. Starvation stretches of midsize messages rise significantly closer to 1 than under SRPT. On the other hand, the increase enables ViFi to complete all large messages by their PS finish times.

SFS also schedules small messages similarly to ViFi: respective points in Figure 3a often coincide. The reason for the similarity is the same as for SRPT versus ViFi. Again, SFS and ViFi differ in their treatment of midsize and large message. A dense cluster of points around starvation stretch 1 for large messages under SFS indicates that SFS reduces delays for midsize messages by postponing large messages almost as long as possible without causing starvation. In addition to the across-the-spectrum line at starvation stretch 1, Figure 3a also reveals sparser but still discernible rows of points with starvation stretches $\frac{1}{2}$, $\frac{1}{3}$, and $\frac{1}{4}$. The rows correspond to messages that arrive and finish while 1, 2, or 3 other messages remain pending (under both PS and SRPT, ViFi or SFS).

To expose the discussed trends more clearly, we repeat the experiment 1,000 times and average the 1,000 obtained sets of starvation stretches sorted in the increasing order of message sizes. Figure 3b shows that SRPT substantially decreases delays of small and midsize messages but the largest messages typically starve. Under ViFi, not only small messages (the rich) benefit from abandoning PS but also the largest messages (the poor) have average starvation stretch about 0.7. Hence, ViFi improves upon PS across the board by reducing delays for all classes of messages: rich, middle, and poor! Figure 3b also illustrates strategic differences between SFS and ViFi. By keeping starvation stretches of large messages to enjoy significantly lower delays than under ViFi.

The average starvation stretches reported in Figure 3b blur fates of individual messages. Hence, Figure 3c plots cumulative distributions of all 300,000 individual starvation stretches in 100 instances of our experiment. Comparison of ViFi with SRPT shows that while starvation stretches up to the 85-th percentile are higher under the fair ViFi, the top 5% of starvation stretches under the unfair SRPT exceed 1, i.e., belong to starved messages. Comparison of SFS with ViFi



Fig. 4. Efficiency of SRPT, PS, ViFi, and SFS with uniformly distributed message sizes.

reveals a main divide around 73%. Up to the 73-rd percentile, starvation stretches are lower under SFS. Under either SFS or ViFi, the top 5% of starvation stretches equal 1. Between the 73-rd and 95-th percentiles, ViFi yields smaller starvation stretches. Similarly to Figure 3a, lines in Figure 3c contain horizontal segments at starvation stretches 1, $\frac{1}{2}$, $\frac{1}{3}$, $\frac{1}{4}$, and $\frac{1}{5}$. These flat segments reflect messages that arrive and finish (under both PS and SRPT, ViFi or SFS) while 0, 1, 2, 3, or 4 other messages remain pending.

C. Efficiency: decrease of average delay

To evaluate efficiency of the algorithms, we conduct our experiment for various values of load l. We repeat the experiment 1,000 times for each examined load $l \leq 120\%$, i.e., including all examined instances of underload, but generally less for overloads of l > 120%. Figure 4a illustrates an intuitive expectation that average delays under SRPT, PS, ViFi, and SFS grow as load increases. After load hits and surpasses 100%, the delays remain finite and even decelerate their growth because the number of messages in every experiment is finite. For the extreme of "infinite" load when all 3,000 messages arrive simultaneously, the average delays are analytically expressed by Gorinsky and Rao [7]. In particular, PS yields the following average delay in a single experiment with simultaneous message arrivals:

$$D_{\rm PS}^{\infty} = \frac{\sum_{k=1}^{n} (2(n-k)+1)m_k}{nC}$$
(5)

where m_k is the size of the k-th smallest message, n = 3,000 is the number of messages, and C = 10 Tbps is the link capacity. When the messages arrive simultaneously, SRPT,

SFS, and ViFi produce an identical transmission schedule for the experiment and achieve the same average delay [7]:

$$D_{\text{SRPT}}^{\infty} = D_{\text{SFS}}^{\infty} = D_{\text{ViFi}}^{\infty} = \frac{\sum\limits_{k=1}^{\infty} (n-k+1)m_k}{nC}.$$
 (6)

n

For the considered uniform distribution of message sizes, we derive the expected average delay under PS as:

$$D_{\rm PS}^{\infty} = \frac{(4n+1)m_{\min} + (2n-1)m_{\max}}{6C} \approx 88,118 \text{ seconds}$$

where $m_{\min} = 100$ GB and $m_{\max} = 100$ TB are respectively minimum and maximum message sizes in the distribution. The expected average delay under SRPT, SFS, and ViFi becomes:

$$D_{
m SRPT}^{\infty} = D_{
m SFS}^{\infty} = D_{
m ViFi}^{\infty} = rac{(n+1)(2m_{
m min} + m_{
m max})}{6C} \approx 44,081 \,\, {
m sec}.$$

Figure 4a confirms that experimental average delays converge asymptotically to the above analytical predictions.

Figure 4b plots average letups under PS, ViFi, and SFS. All three letups peak around l = 100%. At this load where the arrival rate matches the link capacity, PS, ViFi, and SFS have respectively 2.8 times, 22%, and 9% larger average delays than under SRPT. Asymptotically, the average letup under PS converges to:

$$L_{\rm PS}^{\infty} = \frac{(4n+1)m_{\rm min} + (2n-1)m_{\rm max}}{(n+1)(2m_{\rm min} + m_{\rm max})} \approx 2$$

while SFS and ViFi converge to the optimal efficiency:

$$L_{\rm SFS}^{\infty} = L_{\rm ViFi}^{\infty} = 1$$

In general, SFS provides SRPT-like efficiency with consistently lower average delays than even under ViFi.



Fig. 5. Efficiency of SRPT, PS, ViFi, and SFS with Pareto-sized messages.

While our results for l > 100% offer interesting insights into behavior of the algorithms in long-term overload conditions, Figure 4c focuses on underload scenarios l < 100% which are the most relevant for steady-state operation. Again, SFS consistently outperforms ViFi. For example, when load equals 80%, average delays under PS, ViFi, and SFS are respectively 2 times, 7%, and only 3% worse then the minimum attained under the unfair SRPT.

Finally, we explore efficiency of PS, ViFi, and SFS for the Pareto distribution of message sizes. As Figure 5 illustrates, Pareto-sized messages reap even greater benefits from abandoning PS in favor of the efficient representatives of the fair class. Average delays under PS, ViFi, and SFS peak around 7.3 times, 25%, and 11% above the minimum provided by the unfair SRPT. Once again, SFS consistently supports the highest efficiency among the examined fair algorithms.

V. CONCLUSION

In this paper, we studied a class of fair algorithms for bottleneck link capacity allocation where no message finishes later than under PS. In addition to PS, the fair class includes ViFi and newly proposed SFS (Shortest Fair Sojourn). While we proved that no online algorithm in the fair class minimizes average delay of messages, our extensive experiments demonstrated that SFS consistently outperforms PS and even ViFi during either temporal overload or steady-state operation, with largest efficiency benefits achieved when average load is around the bottleneck link capacity. Furthermore, average delay under the fair SFS remains close to the minimum attained under the unfair SRPT. Our simulations revealed that SFS and ViFi gain their significant efficiency improvements over PS across the whole spectrum of message sizes, including large messages but primarily due to dramatic delay reductions for small messages. To outperform ViFi, SFS decreases delays for midsize messages by postponing large messages almost as long as possible without causing starvation.

Whereas link scheduling can affect location of bottlenecks in an arbitrary network topology, our focus on one link was clearly an excessive simplification. We intend to tackle the harder general problem in a future study. Extensions of the presented work will also address application diversity. First, our analysis ignored propagation, node processing, error recovery and other delays dominated by bottleneck transmission delays for long messages. For shorter messages, the extra delays contribute more to overall delay and thereby reduce the relative gains from the efficient utilization of the bottleneck link. To handle a more complex model, we will learn from prior research on message-grained transmission over packet-switching networks [12]. Second, some applications are interested in other network performance metrics than minimal delay achievable under current load. For instance, a streaming application might treat each frame of video as a separate message and prefer a delivery service with guaranteed maximum delay between any two subsequent frames. We are designing an integrated allocation framework where one service minimizes message delays, and the other enables applications to reserve and use end-to-end resources as per needed performance.

REFERENCES

- L. E. Schrage, "A Proof of the Optimality of the Shortest Remaining Processing Time Discipline," *Operations Research*, vol. 16, no. 3, pp. 687–690, May-June 1968.
- [2] M. A. Bender, S. Chakrabarti, and S. Muthukrishnan, "Flow and Stretch Metrics for Scheduling Continuous Job Streams," in *Proceedings ACM-SIAM SODA 1998*, January 1998.
- [3] D. Bertsekas and R. Gallager, Data Networks. Prentice-Hall, 1987.
- [4] A. Demers, S. Keshav, and S. Shenker, "Analysis and Simulation of a Fair Queueing Algorithm," in *Proceedings ACM SIGCOMM 1989*, September 1989.
- [5] M. Shreedhar and G. Varghese, "Efficient Fair Queueing Using Deficit Round Robin," in *Proceedings ACM SIGCOMM 1995*, September 1995.
- [6] V. Jacobson, "Congestion Avoidance and Control," in *Proceedings ACM SIGCOMM 1988*, August 1988.
- [7] S. Gorinsky and N. S. V. Rao, "Dedicated Channels as an Optimal Network Support for Effective Transfer of Massive Data," in *Proceedings IEEE INFOCOM High-Speed Networking (HSN 2006)*, April 2006.
- [8] E. J. Friedman and S. G. Henderson, "Fairness and Efficiency in Web Server Protocols," in *Proceedings ACM SIGMETRICS 2003*, June 2003.
- [9] A. Wierman and M. Harchol-Balter, "Bounds on a Fair Policy with Near Optimal Performance," Carnegie Mellon University, Tech. Rep. CMU-CS-03-198, November 2003.
- [10] N. Bansal and M. Harchol-Balter, "Analysis of SRPT Scheduling: Investigating Unfairness," in *Proceedings ACM SIGMETRICS 2001*, June 2001.
- [11] C. Jechlitschek and S. Gorinsky, "Simulation Suite for Comparative Studies of PS, SRPT, ViFi and SFS," May 2007, http://www.arl.wustl.edu/~gorinsky/sfs.
- [12] E. Modiano, "Scheduling Packet Transmissions in a Multi-hop Packet Switched Network Based on Message Length," in *Proceedings ICCCN* 1997, September 1997.