# On achieving fairness in the joint allocation of buffer and bandwidth resources: Principles and algorithms

Yunkai Zhou [a,1], Harish Sethu [b,*]

[a] *MSN Division, Microsoft Corporation, 1 Microsoft Way, Redmond, WA 98052-8300, United States*
[b] *Department of Electrical and Computer Engineering, Drexel University, 3141 Chestnut Street, Philadelphia, PA 19104-2875, United States*

## Abstract

Fairness in network traffic management can improve the isolation between traffic streams, offer a more predictable performance, eliminate certain kinds of transient bottlenecks and may serve as a critical component of a strategy to achieve certain guaranteed services such as delay bounds and minimum bandwidths. While fairness in bandwidth allocation over a shared link has been studied extensively, the desired eventual goal is overall fairness in the use of all the resources in the network. This paper is concerned with achieving fairness in the joint allocation of buffer and bandwidth resources. Although a large variety of buffer management strategies have been proposed in the research literature, a provably fair and practical algorithm based on a rigorously defined theoretical framework does not exist. In this paper, we describe such a framework and a new, provably fair, and practical strategy for the joint allocation of buffer and bandwidth resources using the max–min notion of fairness. Through simulation experiments using real gateway traffic and video traffic traces, we demonstrate the improved fairness of our strategy in comparison to several popular buffer management algorithms. Joint management of buffer and bandwidth resources involves both an entry policy into the buffer and an exit policy through the output link. Our study reveals that, even though algorithms such as WFQ and DRR that can serve as fair exit policies have received significantly more attention, a fair entry policy is more critical than a fair exit policy to the overall fairness goal when buffer resources are constrained.
© 2005 Elsevier B.V. All rights reserved.

*Keywords:* Fairness; Fair scheduling; Resource allocation; Buffer management; Max–min; RED

## 1. Introduction

### 1.1. Background and motivation

Fairness is an intuitively desirable property in the allocation of resources in a network shared among multiple flows of traffic from different users. Even

---

* Corresponding author. Tel.: +1 215 895 5876; fax: +1 215 895 1695.
*E-mail addresses:* yunkaiz@microsoft.com (Y. Zhou), sethu @ece.drexel.edu (H. Sethu).
[1] Tel.: +1 425 722 7049; fax: +1 425 936 7329.

when the network is over-provisioned, as is the case in parts of the Internet core today, strict fairness in traffic management can improve the isolation between traffic streams, offer a more predictable performance and also improve performance by eliminating some transient bottlenecks. Fair allocation of network resources is especially critical in wireless networks and access networks where the demand for resources is frequently greater than the availability. Fair scheduling policies can also be used to guarantee certain quality-of-service (QoS) requirements such as delay bounds and minimum bandwidths. These policies are likely to play a critical role in future packet-switched networks in supporting applications such as video conferencing and Internet TV stations through controlling the interactions among various traffic streams with different requirements.

Several formal notions of fairness have been proposed to address the question of what is fair in the allocation of a single shared resource among multiple requesting entities. These include, among others, max–min fairness [1–3], proportional fairness [4], utility max–min fairness [5] and minimum potential delay [6]. During the last several years, a variety of algorithms that seek to realize these formal notions have been proposed and implemented to achieve fair allocation of bandwidth on a shared output link [1,2,5–10]. However, bandwidth on a link is only one among several kinds of resources shared by multiple flows in a typical network. As flows of traffic traverse through a network, they share with other flows a variety of resources such as links, buffers and router CPUs in their path. The allocation policies with respect to each of these resources can have a significant impact on the overall performance and QoS achieved by flows. Even though fair scheduling of bandwidth over a link has received the most attention, overall fairness in the use of all the resources in the network is ultimately the desired goal. Several researchers, for example, have already recognized the importance of joint allocation of buffer and bandwidth resources [11–14]. Buffer allocation policies in switches and routers are directly related to congestion avoidance and flow control policies with a direct impact on end-user applications. Fair allocation of buffer resources in routers and switches takes on additional significance with the increasing prevalence of multimedia applications that use UDP instead of TCP and choose to avoid end-to-end congestion avoidance policies.

This paper is concerned with achieving fairness in the joint allocation of buffer and bandwidth resources in a network. A management policy for a shared buffer consists of two components. The *entry scheduler* determines which data from which flows are permitted into the buffer and which are not. The entry scheduler is also responsible for pushout, i.e., the discarding of data from the shared buffer in order to accommodate new arriving traffic. The *exit scheduler* dequeues traffic from the shared buffer and transmits them onto the output link. It is the combination of both the entry and the exit schedulers that determines the overall fairness in the allocation of the buffer and bandwidth resources.

Over the last couple of decades, researchers have proposed and analyzed a variety of entry policies [15–22]. For example, Random Early Detection (RED) [17] is widely used in current Internet routers. A review of various RED-based buffer management algorithms may be found in [23]. Another popular class of entry schedulers are based on modifications to Fair Buffer Allocation (FBA) [19], a review of which may be found in [24]. Most of these have attempted to maximize performance or achieve congestion avoidance although several of them have also tried to be fair by one measure or another. A precise and formal notion of fairness in buffer allocation, however, has not yet been developed. Thus, there is currently no theoretical framework around which one can design practical and fair buffer allocation algorithms, and there also are no formal means of evaluating the various buffer allocation policies already proposed. This paper seeks to provide such a framework to define fairness in the joint allocation of buffer and bandwidth resources, and to facilitate the design of provably fair buffer management strategies.

### 1.2. Contributions

The primary contribution of this paper is a new, provably fair and practical algorithm for the joint allocation of buffer and bandwidth resources. This algorithm is based on a framework that provides a simple but powerful generalization of any of several notions of fairness previously defined for the allocation of a single shared resource or a set of resources viewed as a single entity. Our contribution begins with the definition of an ideally fair strategy, *Fluid-flow Fair Buffering* (FFB), for the joint allocation of buffer and bandwidth resources using the max–min notion of fairness. FFB is an ideally fair