

# Scalable Real-Time Bandwidth Fairness in Switches

Robert MacDavid, Xiaoqi Chen<sup>1</sup>, and Jennifer Rexford<sup>1</sup>

**Abstract**—Network operators want to enforce fair bandwidth sharing between users without solely relying on congestion control running on end-user devices. However, in edge networks (e.g., 5G), the number of user devices sharing a bottleneck link far exceeds the number of queues supported by today’s switch hardware; even accurately tracking per-user sending rates may become too resource-intensive. Meanwhile, traditional software-based queuing on CPUs struggles to meet the high throughput and low latency demanded by 5G users. We propose **Approximate Hierarchical Allocation of Bandwidth (AHAB)**, a per-user bandwidth limit enforcer that runs fully in the data plane of commodity switches. AHAB tracks each user’s approximate traffic rate and compares it against a bandwidth limit, which is iteratively updated via a real-time feedback loop to achieve max-min fairness across users. Using a novel sketch data structure, AHAB avoids storing per-user state, and therefore scales to thousands of slices and millions of users. Furthermore, AHAB supports network slicing, where each slice has a guaranteed share of the bandwidth that can be scavenged by other slices when under-utilized. Evaluation shows AHAB can achieve fair bandwidth allocation within 3.1ms, 13x faster than prior data-plane hierarchical schedulers.

**Index Terms**—Network slicing, fair queuing, packet scheduling, admission control, P4, programmable data plane.

## I. INTRODUCTION

**F**AIR bandwidth allocation between users is an important goal for network operators, since a minority of users demanding too much bandwidth should not negatively affect other users’ quality of service. Yet, leaving bandwidth allocation entirely to congestion control running on end hosts may lead to unfair allocation between different congestion control algorithms. Fair bandwidth allocation is, therefore, a necessary function of the core network. As modern networks scale to higher speed and more users, implementing per-user fair bandwidth allocation becomes increasingly more challenging.

*Network slicing* is a network feature that allows an operator to divide its network resources into many virtualized networks. Slicing enables operators to rapidly create new service offerings for different markets, while achieving performance isolation and quality-of-service guarantees between different slices. To support slicing, the network needs to implement both *intra-slice* fairness where different users within the same slice gets a fair share of the slice’s bandwidth, as well as

*inter-slice* fairness where each slice gets its share of bandwidth proportional to its specified weight. Meanwhile, the idle capacity from underutilized slices must also be fairly distributed to other over-subscribed slices.

One real-life example of a sliced network is the mobile access network. As IoT and 5G becomes prevalent, we face scalability challenges in implementing fairness. A base station may serve 100-1000 user devices, which belong to different classes of services (IoT, smartphones, mobile broadband, first responders, etc.) and have different usage patterns. Each slice (class of service) gets its guaranteed share of bandwidth; when a slice has few active users, its unused bandwidth can be distributed to users in other slices. Meanwhile, we want different users within the same slice to fairly share the limited physical-layer bandwidth: every user in the same slice should be allocated the same maximum bandwidth limit, which should be increased or decreased in real time based on both the number of active users in the slice and the total bandwidth allocated to the slice.

The slice-based fairness paradigm also exists in other scenarios. A data-center network operator may slice its network capacity into multiple classes of service (free tier, spot instances, enterprise customers, etc.) and allocate bandwidth fairly between different tenants within the same slice. Likewise, a network-layer DDoS mitigation mechanism might slice the network to serve different websites, and fairly allocate the bandwidth between all (potentially malicious) clients visiting a particular website. We illustrate an example two-layer hierarchy in Figure 1, where many users (mobile devices, virtual machines, or clients) are grouped into slices, and different slices divide the total bandwidth equally. As user bandwidth demand changes constantly, the fair allocation also changes.

In all of these example use cases, the number of users within each network slice (from thousands to millions) far exceeds the number of hardware queues available on today’s networking hardware, which commonly supports 8-32 queues per port. In today’s mobile network, client rate-limiting and scheduling are sometimes implemented as a virtual network function running on server CPUs [9]. Such a setup supports versatile scheduling policies, yet it requires many CPU cores to serve high-speed traffic and often adds latency and jitter to the traffic. Meanwhile, maintaining ultra-low latency for latency-sensitive applications is one of the most important features in 5G and next-generation 6G networks, which already achieves sub-10ms end-to-end latency [16], [19].

The emergence of high-speed programmable network devices had enabled implementing Active Queue Management (AQM) algorithms directly in the switch data plane [14], [18], [21], [22]. Although recent works [10], [13] have offloaded

Manuscript received 11 July 2023; accepted 14 September 2023; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor Y. Zhang. This work was supported by the Defense Advanced Research Projects Agency (DARPA) under Grant HR001120C0107. An earlier version of this paper appeared at [DOI: 10.1109/INFOCOM53939.2023.10228997]. (Corresponding author: Jennifer Rexford.)

The authors are with the Department of Computer Science, Princeton University, Princeton, NJ 08544 USA (e-mail: macdavid@cs.princeton.edu; xiaoqi@cs.princeton.edu; jrex@cs.princeton.edu).

Digital Object Identifier 10.1109/TNET.2023.3317172

1558-2566 © 2023 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.  
See <https://www.ieee.org/publications/rights/index.html> for more information.

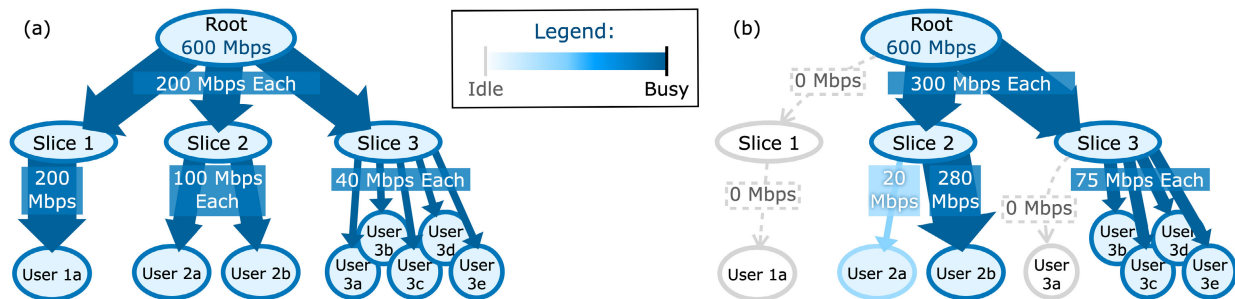


Fig. 1. In (a), all links are fully utilized. In (b), Users 1a, 2a, and 3a do not need their fair share and the surplus is redistributed.

many mobile core network functionalities onto programmable switches, traffic scheduling is a notable exception. To the best of our knowledge, no existing work has attempted to offload scalable slice-based fair bandwidth allocation to high-speed programmable switches. Cebinae [20] enforces long-term fair bandwidth allocation but takes seconds to converge. HCSFQ [22] supports slice-based fair bandwidth allocation but requires per-user memory to monitor each user’s sending rate; this not only adds control-plane overhead for adding and removing users, but also leads to scalability challenges given the limited amount of memory in the data plane.

There are two main challenges for running fair bandwidth allocation directly within the data plane of high-speed programmable switches. Firstly, the available memory is insufficient for maintaining per-user state. We therefore need to use approximate data structures, whose memory footprint scales sub-linearly with the number of users (as discussed in § IV). Secondly, we can only perform a limited set of arithmetic operations in the data plane. We use lookup tables to implement approximated multiplication and division, which is then used for calculating linear interpolation. This enables us to implement real-time, closed-loop iterative update for the per-user bandwidth limit (as discussed in § V). Finally, without using separate queues for each user, we enforce per-user bandwidth limits via probabilistic packet dropping, achieving *approximate* fair bandwidth allocation.

In this paper, we present Approximate Hierarchical Allocation of Bandwidth (AHAB), a hierarchical per-user bandwidth limit enforcer directly implemented in the data plane of programmable switch hardware. AHAB dynamically adjusts the per-user bandwidth limit for each slice in real time, calculated using max-min fairness with the bandwidth demand of all users across all slices. The novelty of AHAB can be summarized as follows:

- **Scalability:** By using the novel CMS-LPF approximate data structure, AHAB avoids maintaining per-user state in data-plane memory, thereby supporting millions of simultaneous users.
- **Fast Convergence:** When user traffic changes, AHAB’s interpolation-based iterative bandwidth limit update converges to fair bandwidth allocation within 3.1ms, 13x faster than prior work [22].
- **Precise Enforcement:** We use probabilistic dropping to precisely enforce bandwidth limits. This allows users to steadily send at the fair rate observing the bandwidth

limit, without requiring hardware queues to pace packets as needed by prior work.

- **One-stop Bandwidth Allocation:** AHAB supports an arbitrary number of hierarchy levels. Therefore, a single instance of AHAB in the core network can rate-limit traffic correctly to adhere to all downstream bandwidth bottlenecks. This is highly useful when downstream devices do not support sophisticated scheduling policies (e.g., legacy routers or thin Wi-Fi access points), or when the network operator is unable to arbitrarily adjust device configurations, possibly because the core and downstream networks are managed between different administrative entities (e.g., MVNOs and wireless carriers).

The rest of this paper is structured as follows. § II defines the hierarchical fair bandwidth allocation problem. § III presents an overview of AHAB’s division of labor between control and data plane. § IV discusses how AHAB overcomes the scalability challenge by avoiding per-user memory using a customized approximate data structure, while § V describes how AHAB approximately calculates an interpolation-based bandwidth limit update given the arithmetic constraints in the data plane. Evaluation in § VI demonstrate that AHAB converges to a fair bandwidth allocation quickly within 5 ms, achieving both fairness and throughput stability. We discuss related work in § VII and conclude in § VIII.

## II. HIERARCHICAL FAIR BANDWIDTH ALLOCATION

AHAB needs to allocate a network slice’s available bandwidth fairly between all users in different slices based on max-min fairness. In this section, we present the same problem definition as in earlier works [3], [7], [15], [22]. For simplicity of discussion, for now we assume all users and slices in our system have equal weight, although it is trivial to add weights and allocate bandwidth proportionally.

### A. Motivating Example

In Figure 1, we illustrate a two-level scheduling hierarchy. The root has a total downlink capacity of 600Mbps and serves three slices. There are 1, 2, and 5 users in each slice, respectively.

Figure 1(a) depicts the “busy” scenario where all users are downloading as fast as possible; the total bandwidth of 600Mbps is equally divided into 200Mbps per slice. The sole user in the first slice gets 200Mbps bandwidth, the users in

TABLE I  
SYMBOL TABLE

| Symbol                     | Definition  |
|----------------------------|---|
| $C_n$                      | Bandwidth capacity allocated to slice $n$         |
| $D_n$                      | Total bandwidth demand of slice $n$               |
| $T_n$                      | Per-user bandwidth limit of slice $n$             |
| $R_m$                      | Bandwidth demand of user $m \in \text{child}(n)$  |
| $T_{low}, T_{mid}, T_{hi}$ | Candidate bandwidth limit for next epoch          |
| $f(T_*)$                   | Hypothetical total sending if the limit was $T_*$ |

the second slice both get 100Mbps, while the five users in the last slice each get 40Mbps.

When user bandwidth demand changes, the fair allocation also changes. In Figure 1(b) we show a scenario where user 1a and 3a leave the network (e.g., powered off) and user 2a has a low bandwidth demand (e.g., audio streaming); all other users are still downloading as fast as possible. In this case, slice 1 does not use any bandwidth, so the total bandwidth is equally divided to the other two slices (300Mbps each). User 2b can use all the remaining bandwidth in slice 2, which is 280Mbps. Users 3b-3e each can use 75Mbps of fair share.

Hierarchical fair bandwidth allocation achieve fairness in both levels: busy users in the same slice get the same bandwidth, after considering the underutilized users in the same slice; different busy slices also get the same aggregated bandwidth, after re-allocating the unused bandwidth from underutilized slices.

### B. Max-Min Fairness Allocation

Now we formally define the fair hierarchical bandwidth allocation based on max-min fairness and work-conserving scheduling.

Let us first denote a slice  $n$ 's total bandwidth capacity as  $C_n$ , and its set of users as  $\text{child}(n)$  (its "children" in the scheduling hierarchy). Each user  $m \in \text{child}(n)$  has a bandwidth demand  $R_m$ , which is the maximum bandwidth it would like to consume if no bandwidth limit is enforced. We can therefore calculate slice  $n$ 's total bandwidth demand as the sum of all user's demand:

$$D_n = \sum_{m \in \text{child}(n)} R_m, \quad (1)$$

and we call slice  $n$  underutilized if  $D_n \leq C_n$ . In this case, a bandwidth limit is unnecessary as the slice's capacity is greater than the total demand of all users.

However, when a slice is busy, the demand exceeds capacity and we need to impose a per-user bandwidth limit, such that the *actual* total bandwidth usage of this slice equals to its capacity. Based on max-min fairness, we can define the per-user bandwidth limit as

$$T_n = \arg \max_T \sum_{m \in \text{child}(n)} \min(T, R_m) \leq C_n. \quad (2)$$

Enforcing this bandwidth limit would have no effect for users requiring less bandwidth ( $R_m < T_n$ ) and only affects the users using more bandwidth. We can prove that such a limit  $T_n$  always exist for a busy node  $n$ , and the resulting enforcement achieves the unique max-min fairness allocation

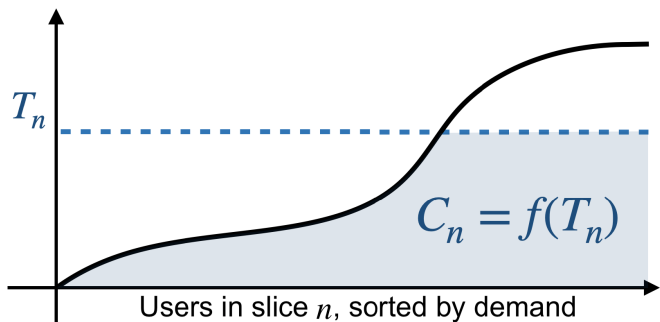


Fig. 2. Limit  $T_n$  is enforced against those users in slice  $n$  whose demand exceeds than  $T_n$ , so that the total bandwidth consumed is exactly  $C_n$ .

between users. For completeness, we define  $T_n = \infty$  for underutilized slices.

In Figure 2, we illustrate the relationship between the bandwidth limit  $T_n$  and the total bandwidth used by users in slice  $n$ . The demands  $\{R_m \mid m \in \text{child}(n)\}$ , when sorted from least to greatest, form a demand curve. After enforcing per-user limit  $T_n$ , each user uses  $\min(T_n, R_m)$ , and the total bandwidth used by all user is  $\sum_{m \in \text{child}(n)} \min(R_m, T_n)$  which is represented the size of the shaded area, under the intersection of demand curve  $R_m$  and the horizontal line of limit  $T_n$ . This area increases as we increase  $T_n$ , and the ideal limit  $T_n$  means the shaded area has size exactly equal to the available bandwidth capacity  $C_n$  of this slice. Thus, our goal for finding the right bandwidth limit  $T_n$  under max-min fairness is equivalent to finding the right "horizontal cut" across the demand curve.

Note that the fair allocation is not fixed: As the user's bandwidth demand constantly changes, it is necessary to recalculate the fair allocation and update  $C_n$  and  $T_n$ . Calculating the hierarchical fair bandwidth allocation requires a two-step process: aggregating the demands upwards, then allocating the capacity downwards. First, we calculate the sum of all slice's demand,  $D = \sum_n D_n$ , which represents the total bandwidth demand at the root of the entire scheduling hierarchy. We then allocate the total capacity  $C$  at the root of the scheduling hierarchy for different slices. When the total demand exceeds capacity ( $D > C$ ), we can allocate per-slice capacity  $C_n$  for each slice, using exactly the same max-min fairness principle as in Equation 2. Subsequently, using  $C_n$ , we can obtain the per-user rate limit  $T_n$  for each slice. When these limits are enforced for every user, we implement max-min fair bandwidth allocation across the entire scheduling hierarchy, and the total bandwidth used by all slices will equal to the root's total capacity.

The discussion here focused on a two-level scheduling hierarchy; the same definition applies to more levels.

### III. AHAB SYSTEM OVERVIEW

Figure 3 illustrates the basic design of AHAB. At a high level, we split the bandwidth allocation process into a fast-reacting data plane component and a more sophisticated control-plane component for hierarchical updates.

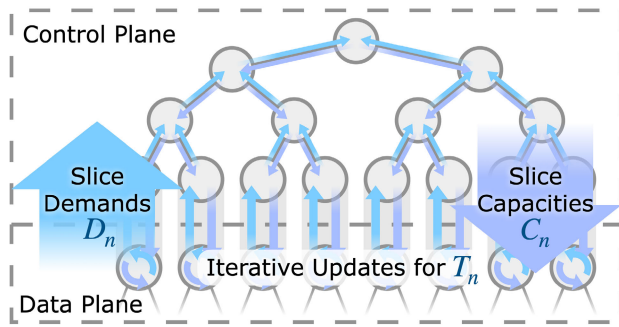


Fig. 3. The control plane maintains inter-slice fairness by periodically reading each slice bandwidth demands  $D_n$  and writing fair capacities  $C_n$ ; the data plane keeps intra-slice fairness by iteratively updating bandwidth limits  $T_n$ .

#### A. Data Plane: Intra-Slice Fairness

To quickly react to changes in individual user's traffic, AHAB calculates iterative updates for the per-user bandwidth limit  $T_n$  fully within the data plane, using approximated linear regression. This allows the intra-slice bandwidth allocation to converge within milliseconds after a user starts or stops sending, much faster than updating using the switch control plane.

Since it is impossible to perfectly predict the constantly changing per-user traffic demand, AHAB splits the traffic into very small time epochs (on the order of milliseconds) and uses the demand distribution in the past epoch as a prediction for the next epoch. At the end of each epoch, we use this demand distribution (as illustrated in Figure 2) to iteratively refine the per-user bandwidth limit  $T_n$ , such that the total bandwidth used by all users in the slice will be approximately equal to the capacity  $C_n$ .

#### B. Control Plane: Inter-Slice Fairness

The control plane is responsible for implementing fair allocation for all the other layers of the scheduling hierarchy beyond the slice level. It periodically reads the per-slice total bandwidth demand  $D_n$  (by reading registers populated by the data plane), and calculate the aggregate demand at higher level of the hierarchy. Then, it calculates the bandwidth allocation according to max-min fairness (§II-B) starting from the root, until it obtains the per-slice fair allocation capacity  $C_n$ . These values are then written to registers that will be read by the data plane.

The periodic update happens every 10-20ms, mostly due to communication bottleneck (reading demands and writing capacities). Thanks to statistical multiplexing, the aggregated bandwidth demand of different slices changes on a longer timescale. Therefore, the slightly slower update of capacities has little impact on maintaining intra-slice fairness. We also note that when all slices are busy, their fair allocation is constant; the allocated capacities only changes when some slices are underutilized.

### IV. SCALING BEYOND MEMORY LIMITS

In this section, we discuss how AHAB overcomes the scalability challenge imposed by hardware memory size limits.

We first discuss how the bandwidth limit  $T_n$  is enforced on each user using their estimated sending rates. Subsequently, we show how AHAB avoids allocating per-user memory, using a novel approximate data structure that combines Count-Min Sketch with Low-Pass Filters to estimate per-user sending rates. Finally, we discuss how we share one approximate data structure across all slices using weight-based normalization.

#### A. Enforcing Bandwidth Limits

For the entire scheduling hierarchy to achieve bandwidth fairness, we must properly enforce bandwidth limit  $T_n$  on all users. Naively, we can allocate one queue per user and assign the bandwidth limit as the queue's drain rate. However, the number of users (thousands to millions) far exceeds the number of queues available in hardware switches (8)-32 queues per port). Instead, we can enforce bandwidth limits using active queue management, or more specifically probabilistic dropping as discussed in [17] and [12], as long as we know the user's sending rate. This approach does not require a traffic scheduler, and can be performed even if the switch has only a single queue.

For a user  $m$  in slice  $n$  with bandwidth limit  $T_n$  and sending rate  $R_m$ , we can enforce the bandwidth limit  $T_n$  by dropping its packets with probability

$$\Pr[\text{drop}] = 1 - \min\left(1, \frac{T_n}{R_m}\right) \quad (3)$$

as described in [17] and [12]. If a user uses less than the limit  $T_n$ , no packet will be dropped; otherwise, after probabilistic dropping the user's remaining packets will use bandwidth equal to  $T_n$ .

We also observe that probabilistic dropping is unfriendly for TCP flows, as TCP achieves low goodput if we perform probabilistic dropping whenever its instantaneous sending rate exceeds  $T_n$ . This is partly because TCP congestion control slows down severely upon two consecutive packet drops, by reducing its congestion window back to slow-start; meanwhile, randomized packet dropping will quite often produce consecutive packet drops.

Therefore, we specifically optimize the rate-limiting for TCP flows by using *periodic* instead of *probabilistic* dropping, and drop one packet approximately every  $T_n \cdot \text{const}$  bytes, corresponding to the desired congestion window size for achieving goodput equal to  $T_n$ . We also adapt various TCP shaping techniques discussed in Nimble [18], such as adding Early Congestion Notification (ECN) flags. This way, much fewer TCP packets are dropped compared to non-TCP flows; well-behaving TCP flows enjoy relatively steady goodput while AHAB can enforce the bandwidth limit  $T_n$  effectively. We leave the open question of how to optimally enforce bandwidth limits on TCP flows with various congestion control algorithms as a future work.

#### B. Avoiding Per-User Memory

Knowing a user's sending rate  $R_m$  is vital for correctly enforcing the bandwidth limit. As discussed in [17] and [22], asking the sender of all traffic to attach their traffic rate to

each packet is an easy yet unrealistic solution, as the sender might belong to a different administrative entity and may not honestly report the rate. Therefore, AHAB needs to measure each user's sending rate directly.

Recent works [18], [22] in queue scheduling within high-speed programmable switches rely on using the onboard memory to maintain per-user sending rate statistics, by allocating one traffic counter per user. However, programmable switches only have a limited amount of onboard memory in the data plane, limiting its scalability. At any given time, a core network switch may be servicing millions of users across thousands of base stations, making it infeasible to store any per-user state in memory, not to mention the hassle of keeping the memory allocation up-to-date when users constantly join or leave the network.

Instead, we build a customized memory-efficient approximate data structure to track per-user sending rate, by combining two techniques: Low-Pass Filters (LPF) and Count-Min Sketch (CMS) [6]:

- The LPF is a self-decaying counter, available as an advanced feature of the Tofino switch hardware. If we add value  $x$  at time  $t$  to a LPF with previous value  $v_0$  and last update time  $t_0$ , its new value becomes

$$v = x + v_0 e^{-(t-t_0)/\tau}, \quad (4)$$

where  $\tau$  is its decay time constant. As discussed in [17], if we aggregate the packet sizes of a single user's traffic in a LPF, the LPF will report an exponentially-decayed moving sum of recent packet sizes, which is proportional to a good estimate of the user's instantaneous sending rate.

- The CMS [6] is an approximate data structure that answers frequency queries, using  $r$  rows of hash-indexed arrays each having  $c$  counters. Given an "insertion" with a size and a user ID, we find one location per row by applying  $r$  different random hash functions over its ID, and increment the counters at those location by the size; when querying the total size of a particular ID, we find the same  $r$  locations and report the minimum of the  $r$  counters.

When used to estimate size of flows in traffic, CMS is good at reporting heavy flows, as it never underestimates flow sizes. However, a vanilla CMS can only track the total number of bytes sent by a user since the CMS is initialized, not the user's instantaneous sending rate. Although it is possible to run multiple instances of CMS in a round-robin fashion to query moving-window flow rates [5], such an arrangement adds complexity and requires 2x-4x more memory.

AHAB combines CMS with LPF by replacing individual counters in the CMS structure with LPF counters. When inserting a packet with its size and user ID, we apply the  $r$  hash functions over the user ID to locate one LPF counter per row, and add the packet size to these  $r$  LPF counters. When querying the instantaneous sending rate of the same user ID, we read the LPF counters in the same  $r$  locations, and use the minimum across their reported rate as the estimated sending rate of this user. This allows us to estimate per-user sending rate without the need to allocate per-user memory.

We note that the combined CMS-LPF estimator retains the following additive-error guarantee:

*Theorem 1:* Let  $R_i$  be user  $i$ 's sending rate reported by an ideal LPF counter, and  $\sum_m R_m$  be the total sending rate across all users, again reported by ideal per-user LPF counters. When querying a CMS-LPF estimator of size  $r \times c$ , the estimated sending rate  $\tilde{R}_i$  satisfies  $\tilde{R}_i \geq R_i$  and

$$\Pr \left[ \tilde{R}_i \leq R_i + \epsilon \sum_m R_m \right] \leq \delta, \quad (5)$$

with  $\epsilon = e/r$  and  $\delta = e^{-c}$ .

*Proof.* Note that CMS is a linear sum in the ID dimension, where each counter is the sum of a random subset of user IDs. Meanwhile, LPF is a linear sum of packet sizes in the time dimension, where the reported rate is the inner product of past packet sizes and a time-decaying constant function. The two linear operators are interchangeable. Thus, we can naturally derive the error bound using the additive error property of an unmodified  $r$ -row,  $c$ -column CMS.  $\square$

The CMS-LPF also inherits the no-underestimation guarantee from CMS. Therefore, in our context of enforcing bandwidth limit:

- 1) A user's traffic rate is never underestimated, ensuring that a user exceeding bandwidth limit will always be rate-limited.
- 2) With a small bounded probability, a user's traffic rate may be overestimated and appear higher than the limit, resulting in rate-limiting. We can limit this probability by adjusting the memory size of CMS-LPF.

### C. Sharing One Rate Estimator Across Slices

Naively, AHAB would allocate one CMS-LPF estimator for each slice. However, due to the natural skewness of traffic, not all slices will have lots of "heavy" users sending at high rates. Some slices may be underutilized and have no heavy user at all, and the memory dedicated for their estimators is wasted.

Instead, we share a single CMS-LPF estimator across all slices. We can then exploit statistical multiplexing, as the heavy users and busy slices are now effectively using the unused memory sacrificed by the underutilized slices with no heavy user.

However, we note that CMS provides an additive error guarantee, meaning that the error of each user's estimated rate is of similar magnitude regardless of the true sending rate of the user. This is not a problem for intra-slice comparison, as we only care about enforcing bandwidth limits for heavy users and can safely ignore the underutilized users. Yet, different slices may have vastly different bandwidth allocations. If two slices of capacity 100Mbps and 10Gbps naively share the same CMS-LPF structure, the 10Gbps slice will dominate; "small" users of 200Mbps in the heavy slice will overwhelm the CMS while the "heavy" users of 30Mbps in the small slice become a rounding error.

To ensure the estimation error is scaled proportionally with the bandwidth of different slices, we perform *pre-update normalization*: before packet sizes are fed into the CMS-LPF, we scale the packet size inversely proportional to the

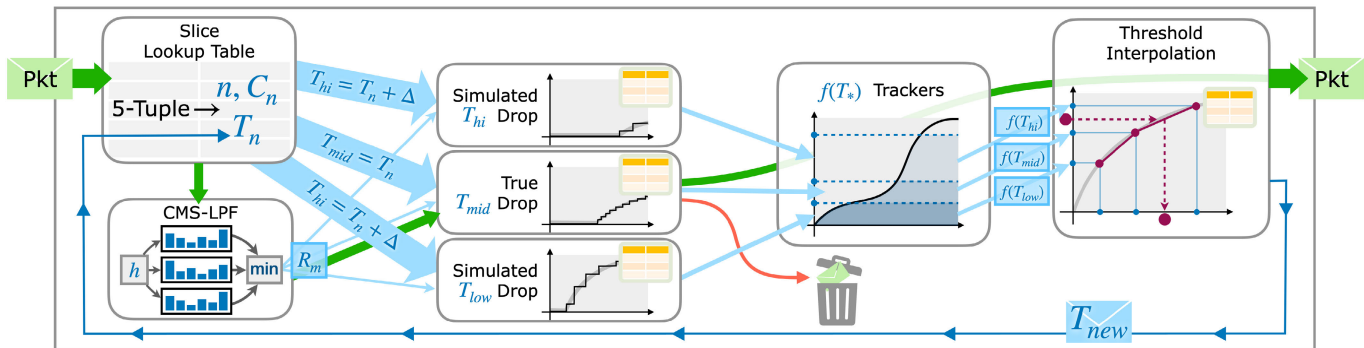


Fig. 4. When a packet arrives, AHAB first maps it to a slice  $n$  and estimates its user's sending rate  $R_m$  using the CMS-LPF estimator (§IV-B), then uses probabilistic dropping to enforce slice  $n$ 's bandwidth limit  $T_n$  (§V-A). AHAB also maintains two bandwidth limit candidates  $T_{low}$ ,  $T_{hi}$  and tracks the hypothetical total bandwidth usage  $f(\cdot)$  (§V-B), which is used to derive a more accurate bandwidth limit  $T_{new}$  via approximated linear interpolation (§V-C).

weight of their parent slice. This guarantees that we can track heavy users in each slice effectively, with the estimation error proportional to its particular slice-level sending rate, regardless of how slow or fast other slices are sending. Subsequently, when we enforce per-user bandwidth limits, the estimated sending rates we read from CMS-LPF are also compared to scaled versions of bandwidth limits.

## V. APPROXIMATE ARITHMETIC IN THE DATA PLANE

To achieve line-rate packet processing and low forwarding latency, high-speed programmable switches like Intel Tofino support a limited set of arithmetic operations, and we can only perform a constant number of computational steps per packet. Thus, it is infeasible to exactly track the bandwidth demands, precisely calculate the fair allocation, or accurately compute per-user drop probabilities.

Figure 4 shows an overview of the AHAB data plane, where we use *approximate* arithmetic heavily to implement probabilistic dropping and interpolation-based iterative update to the bandwidth limit. We also note that the approximate arithmetic techniques presented here are widely applicable to other applications running in programmable switches, beyond fair bandwidth allocation.

### A. Approximated Probabilistic Dropping

For a given user, we obtain its estimated sending rate  $R_m$  from the CMS-LPF estimator and compare it against the per-user bandwidth limit  $T_n$  of its slice. For non-TCP traffic, we need to enforce the bandwidth limit by dropping the packet with probability  $1 - \frac{T_n}{R_m}$ .

Since we cannot calculate exact division in the data plane, we perform approximate arithmetic using a TCAM lookup table, similar to the approximate multiplication technique used in Nimble [18].

The first step of approximate division is to truncate the numerator  $T_n$  and denominator  $R_m$  in binary form. We focus on the highest  $b$ -bits of the denominator starting from the leading 1. For example, let  $b = 5$ , if  $R_m$  is  $0b0011010011010$ , the leading  $b$ -bit number is  $j=11010$ . This means we can approximate  $R_m \approx 0b11010 \times 2^6$ . We also look at the same bits in the numerator: if  $T_n$  is  $0b0001111110100$ ,

we extract the  $b$ -bit number  $i=01111$  and approximate it as  $T_n \approx 0b01111 \times 2^6$ . Note that, since the denominator is always larger than the numerator ( $R_m > T_n$ ), we will not miss any 1 in the higher bits in the nominator.

We now observe that the fraction  $\frac{T_n}{R_m}$  can be approximated using these  $b$ -bit numbers, as

$$\frac{0b0001111110100}{0b0011010011010} \approx \frac{0b01111}{0b11010} \approx 0.576.$$

More formally, we use

$$\frac{I}{J} \approx \frac{i \cdot 2^N}{j \cdot 2^N} \quad (6)$$

to approximate division, with  $i$  and  $j$  both being  $b$ -bit numbers. We build an approximate division lookup table that maps  $(i, j)$  to an approximation of  $i/j$ . After picking the appropriate  $b$ -bit substrings  $i, j$  of the input numerator and denominator, the final step of approximate division is simply look up  $(i, j)$  in the table.

Since  $i$  and  $j$  are truncated from  $I$  and  $J$ , they both have a one-sided bias compared with the original. To reduce the error bias of lookup table entries, the entries are computed as  $\frac{i+0.5}{j+0.5}$  instead of  $\frac{i}{j}$ , which changes the  $x\%$  worst-case one-sided error to  $x/2\%$  two-sided error.

Here we make two observations. First, the most-significant 1-bit in  $J$  always end up in  $j$ , so we only need to generate lookup table entries for half of the possible  $j$ s starting with 1. Second, notice that the numerator is always smaller than the denominator in our use case, lookup table entries do not need to be installed for any division results greater than 1, further reducing the number of entries. For example, for  $b = 5$  bits, we only need to generate entries for  $j=0b10000$  to  $j=0b11111$ , and for each  $j$  we generate  $j + 1$  rules for  $0 \leq i \leq j$ . In total, the number of rules in the lookup table is

$$\sum_{j=2^{(b-1)}}^{(2^b)-1} (j + 1) = \sum_{j=16}^{31} j + 1 = 392.$$

This table achieves an approximation error of 8.2% on average. With  $b = 6$ , the lookup table grows to 1552 entries and we can reduce the error to 4.6%. One such lookup table only costs approximately 9KB of data-plane memory.

The final step of implementing probabilistic dropping is to sample a number uniformly at random between  $[0, 1]$  using the random number generator, and comparing it with the approximated division. We drop a packet if

$$U(0, 1) > \frac{T_n}{R_m} \approx \frac{i + 0.5}{j + 0.5}. \quad (7)$$

Separately, for TCP traffic, we need to perform periodic rather than probabilistic dropping. We maintain a hash-indexed array of “count-down” counters and map each TCP flow to a counter using its 5-tuple. When a TCP packet arrives and the estimated rate exceeds the bandwidth limit  $T_n$ , we check its corresponding counter: if the counter is already zero, we drop this packet and reset the countdown counter to  $T_n \cdot \text{const}$ ; otherwise, we spare the packet from being dropped and subtract its size from the counter. We also use another set of counters to add ECN marking periodically. This way, TCP flows can quickly converge to the desired rate limit and enjoy steady goodput.

### B. Tracking the Bandwidth Demand

For a given slice, it is infeasible for switches to track its entire bandwidth demand curve (shown in Figure 2) representing all user’s sending rates, which we can neither store nor sort. However, we can track the actual bandwidth used by all users  $f(T_n)$ , which is a function of the currently-enforced bandwidth limit  $T_n$  and represented by the shaded area under the demand curve intersected with  $T_n$ . To get  $f(T_n)$ , we simply need to use a LPF to track the size of all packets that are not dropped.

Still, comparing  $f(T_n)$  with  $C_n$  only tells us whether we are over- or under-utilizing the capacity  $C_n$ , i.e., whether we should increase or decrease  $T_n$ . This does not say much about what is the ideal limit or how much should we change  $T_n$ .

Adjusting  $T_n$  using a fixed step size or fixed proportion is problematic: if the steps are too large, we cannot precisely converge to the exact allocation. Yet, small step sizes mean we need to wait for many iterations before converging, when the fair rate changes dramatically.

To better analyze how to update  $T_n$ , we further specify two *candidate* bandwidth limits, a lower candidate

$$T_{low} = T_n - \Delta$$

and a higher candidate

$$T_{hi} = T_n + \Delta$$

where  $\Delta$  is the maximum step-size we want to change  $T_n$ . In practice, we choose  $T_{low} \approx 0.5 \cdot T_n$  and  $T_{hi} \approx 1.5 \cdot T_n$ . From now on, we also refer to  $T_{mid} = 1.0 \cdot T_n$  as the middle candidate.

We now track two more hypothetical total transmitted bandwidth  $f(T_{low})$  and  $f(T_{hi})$ , by generating two hypothetical probabilistic dropping decisions in addition to the real dropping decision. Using the same lookup table technique discussed in §V-A, we approximately calculate  $\frac{T_{hi}}{R_m}$  and  $\frac{T_{lo}}{R_m}$  and track the packets that are hypothetically not dropped under  $T_{low}$  or  $T_{hi}$  respectively. As illustrated in Figure 5,  $f(T_{low})$ ,

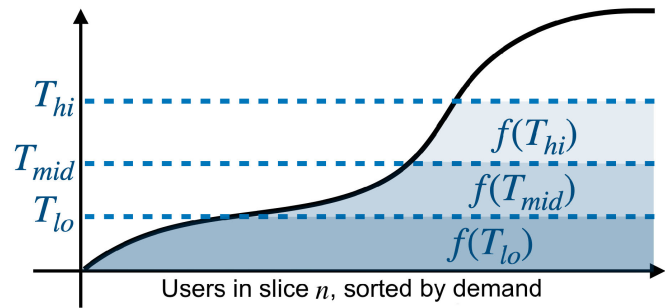


Fig. 5. Relationship between bandwidth limit candidates  $T_{low}$ ,  $T_{mid}$ ,  $T_{hi}$  and total bandwidth consumption  $f(T_{low})$ ,  $f(T_{mid})$ ,  $f(T_{hi})$ .

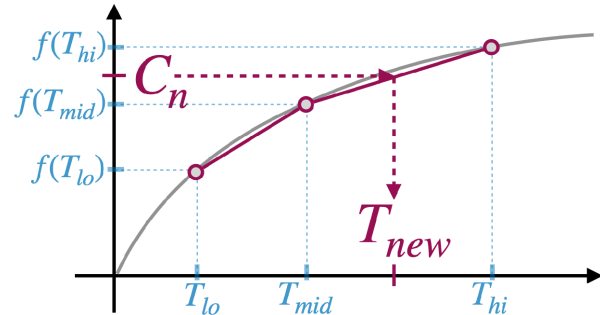


Fig. 6. Finding new bandwidth limit using interpolation. When we plot  $f(T_*)$ , x-axis in this figure refers to the bandwidth limit  $T_*$  (y-axis in Fig. 5), while y-axis in this figure refers to the area under line  $T_*$  in Fig. 5.

$f(T_{mid})$ , and  $f(T_{hi})$  are the shaded area under the demand curve intersected with different horizontal lines.

Figure 6 plots the monotonically-increasing function  $f$ . The optimal bandwidth limit  $\tilde{T}$  satisfies  $f(\tilde{T}) = C_n$ , thus we need to calculate a new limit  $T_{new}$  that is as close to  $\tilde{T}$  as possible.

### C. Update Bandwidth Limit via Approximate Interpolation

Instead of using a fixed step size, we can adjust the bandwidth limit much more accurately using linear interpolation, given the three candidate points on the  $f$  curve.

Let us first assume the ideal bandwidth limit lies between the lower and higher candidate points, i.e.,  $f(T_{low}) < C_n < f(T_{hi})$ . Without loss of generality, assume we need to adjust to a higher limit, i.e.,  $f(T_{mid}) < C_n < f(T_{hi})$ . As illustrated in Figure 6, we can calculate the new bandwidth limit using linear interpolation:

$$T_{new} = T_{mid} + \frac{C_n - f(T_{mid})}{f(T_{hi}) - f(T_{mid})} \times (T_{hi} - T_{mid}). \quad (8)$$

In Figure 7 we illustrate how to calculate the approximated linear interpolation using an example. To calculate

$$\frac{C_n - f(T_{mid})}{f(T_{hi}) - f(T_{mid})} = \frac{1012}{1690}, \quad (9)$$

we first use the same approximate division lookup table technique discussed earlier in §V-A, except the division results are now stored as a (mantissa, exponent) pair. We truncate the numerator and denominator to get most significant non-zero bits  $i = 011111/j = 11010$ , and retrieve the approximate

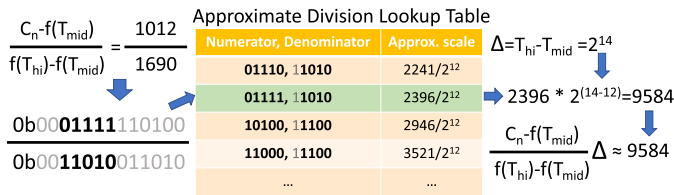


Fig. 7. We use a lookup table to implement approximate linear interpolation. We first match on the highest binary bits of numerator and denominator to get a scaled division result, then multiply  $\Delta$  via bit shifting for the final result.

division result

$$\frac{i + 0.5}{j + 0.5} = \frac{2396}{2^{12}}. \quad (10)$$

After the approximate division, we need to multiply the result by  $\Delta$ . To make this calculation easier, we choose  $\Delta$  to be a power of 2, reducing the multiplication into a bit-shift. In practice, we set  $\Delta = 2^{\lfloor \log_2(\frac{1}{2}T_{mid}) \rfloor}$ , meaning  $T_{low} = T_{mid} - \Delta \approx 0.5 T_{mid}$  and  $T_{hi} = T_{mid} + \Delta \approx 1.5 T_{mid}$ .

The divide-then-multiply calculation can be applied as a single bit shift. In the example in Figure 7, we have  $\Delta = 2^{14}$  and need to calculate  $\frac{2396}{2^{12}} \cdot 2^{14}$ , which can be simplified into a left shift:

$$2396 \ll (14 - 12) = 9584. \quad (11)$$

Finally, we finish the last addition operation in the approximate linear interpolation, and obtain the new bandwidth limit

$$T_{new} = T_{mid} + 9584. \quad (12)$$

Similarly, when adjusting towards a lower limit, we use

$$T_{new} = T_{mid} - \frac{f(T_{mid}) - C_n}{f(T_{mid}) - f(T_{low})} \times (T_{mid} - T_{low}). \quad (13)$$

Notice that we use subtraction from  $T_{mid}$  instead of adding up from  $T_{low}$  to interpolate. This is because the approximate division has a constant relative error proportional to the result. By subtracting the result from  $T_{mid}$ , we can make more accurate fine-grained adjustments near  $T_{mid}$  to better converge towards the optimal bandwidth limit. Instead, if we use

$$T_{new} = T_{low} + \frac{C_n - f(T_{low})}{f(T_{mid}) - f(T_{low})} \times (T_{mid} - T_{low}), \quad (14)$$

the approximated interpolation is more accurate near  $T_{low}$  and has a larger error near  $T_{mid}$ .

When  $C_n$  falls out of the range  $[f(T_{low}), f(T_{hi})]$ , our estimate candidates are too far off from the ideal bandwidth limit, and we clip the update by choosing  $f(T_{low})$  or  $f(T_{hi})$  directly. Clipping prevents overshooting caused by using linear interpolation outside of the two candidate points.

We further note that although CMS-LPF will introduce over-estimation errors across the board for all estimated rates, our closed-loop bandwidth limit update process will naturally adapt to this error. When all rates are slightly over-estimated while the bandwidth limit  $T_n$  is not yet over-estimated, users will suffer from an unnecessarily high drop probability, leading to less than  $C_n$  total traffic; AHAB will then automatically raise  $T_n$  to account for such global over-estimation.

#### D. Iterative Update Using Worker Packets

To achieve fast convergence towards intra-slice fairness, AHAB updates the bandwidth limit  $T_n$  fully within data plane. At the end of every epoch, AHAB calculates a new bandwidth limit  $T_{new}$  for each slice using approximate interpolation, and use it as the new bandwidth limit for the next epoch.

However, as shown in Figure 4,  $T_n$  is stored in a register memory lookup table near the beginning of the switch's packet-processing pipeline while the new limit  $T_{new}$  is only available in later pipeline stages; the pipeline's memory access constraint does not allow us to write  $T_{new}$  back to the same register memory directly. Therefore, at the end of every epoch, we generate one worker packet per slice by packet cloning, and use packet recirculation to let the worker packet go through the pipeline twice. The worker packet reads all the candidate bandwidth limits, performs the approximated linear interpolation calculation to derive the new bandwidth limit  $T_{new}$ , and carries it to the beginning of the packet-processing pipeline to be written into register memory.

Although the update is slightly delayed due to packet recirculation (about  $0.65\mu s$ ), only a very small fraction of packets near the beginning of the epoch are affected, therefore the actual difference in enforcement due to the delayed update is negligible. As we show in §VI, this closed-loop update process rapidly converges to the fair bandwidth allocation.

#### E. Supporting Weighted Allocation

A network operator sometimes needs to allocate bandwidth in proportion to a pre-assigned weight, for example when implementing differentiated services. AHAB supports weighted fair allocation at both the slice level and the user level.

To support weighted fair bandwidth allocation between users in the same slice, we scale each packet's length using the user's weight: if a user  $m$  has weight  $w_m$ , a packet with size  $S$  is scaled into  $\frac{S}{w_m}$  before being used to calculate the user's scaled sending rate  $R_m$  in the CMS-LPF estimator. This way, we can directly compare different user sending rates  $R_m$  against the same per-user bandwidth limit  $T_n$ .

Meanwhile, the control plane is more flexible and trivially supports allocating bandwidth to different slices based on their weight. We simply divide each slice's demand and capacity by its weight before computing the max-min fairness allocation.

We also note that the weights assigned to slices / users can be easily updated at run time. To adjust the weight for a subset of users, we adjust the rules installed in the slice lookup table in the data plane; to adjust the weight of a slice, we modify it directly from the control plane.

## VI. EVALUATION

Using a prototype implementation running in a hardware testbed, we show that AHAB can quickly achieve fair and stable bandwidth allocation between users. Compared to the prior state-of-the-art, HCSFQ [22], AHAB not only converges to the target fair bandwidth allocation faster (in 3.1ms), but also achieves comparable or better fairness and throughput stability. Subsequently, we use real-world traffic traces in a



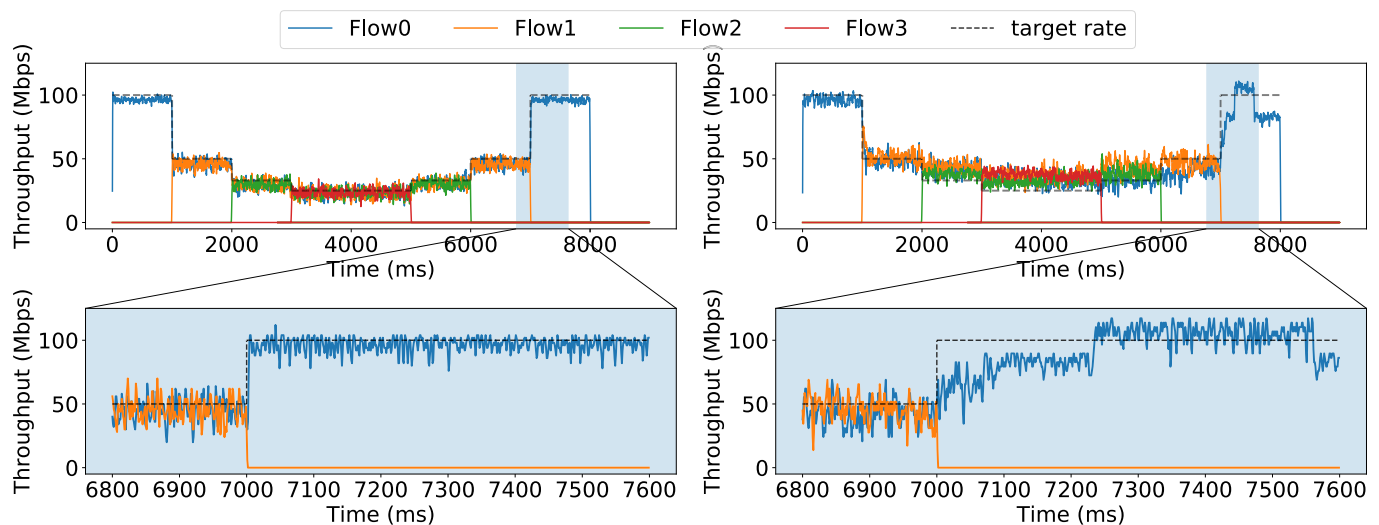


Fig. 8. AHAB (left) converges to fair bandwidth allocation within 3.1ms on average, while HCSFQ [22] (right) needs 42.3ms.

simulation-based experiment to show that AHAB scales well to 5.9-23.9 million users with a reasonable memory footprint, and CMS-LPF has a minimal impact on scheduling fairness.

#### A. Testbed Experiment Setup

We evaluate AHAB’s real-world scheduling fairness using a hardware testbed with two sender and receiver servers connected via an Intel Tofino Wedge32-X programmable switch, which runs a prototype implementation of AHAB written in 2,000 lines of P4 [1]. Both servers have a 20-core CPU and a Mellanox ConnectX-5  $2 \times 100$ Gbps NIC, and run Ubuntu 20.04. The sender sends TCP flows using `iperf3` with Linux’s default congestion control (cubic), and send UDP flows using either `iperf3` or a customized Go script that performs millisecond-level throughput measurement. We set the iterative update epoch time to 1ms and configure the LPF rate estimator’s time constant to  $\tau=4$ ms. Unless otherwise noted, we use a CMS-LPF estimator with size  $3 \times 2048$ .

In all experiments, we treat each flow as a unique user, using its 5-tuple (source and destination IP/port pairs) as the user ID. We note that real-world traffic may be grouped more coarsely; for example, one user in a mobile network may include all flows destined for the same device (same destination IP).

#### B. Fast Convergence

We now compare AHAB to the state-of-the-art of hierarchical fair queuing based on programmable switch: HCSFQ [22]. HCSFQ iteratively converges to the fair rate via Additive Increase Multiplicative Decrease, limiting its convergence speed when the number of users decreases and the fair rate increases.

To demonstrate the difference in convergence time, we program both AHAB and HCSFQ to enforce fairness between four UDP flows in a single slice, with a fixed 100Mbps capacity. All four flows have the same 100Mbps constant sending rate, but have different starting and ending time: they run between  $T = 0-8$ s,  $T = 1-7$ s,  $T = 2-6$ s, and  $T = 3-5$ s, respectively.

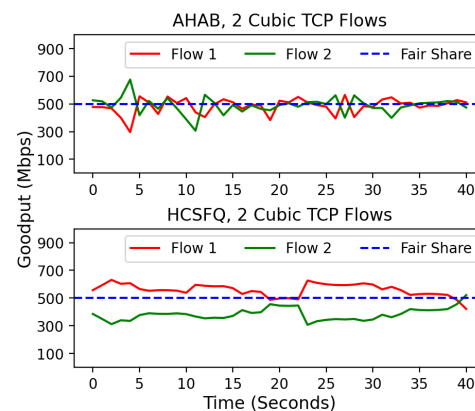


Fig. 9. Goodput of two competing TCP flows over time.

Figure 8 shows the actual bandwidth used by the four flows over time, after bandwidth limit enforcement done by AHAB (left) or HCSFQ (right). At a longer timescale (top), the two schedulers behaved similarly. However, if we zoom in to a smaller timescale and plot the millisecond-level per-flow throughput (bottom) immediately after  $T=7$ s (where flow 1 stopped), we can see AHAB converges much faster than HCSFQ to allow flow 0 to use the full 100Mbps bandwidth. We measured the time for flow throughput to converge to within 10% of ideal fair bandwidth limit. AHAB’s interpolation-based update only needs around three iterations to converge, taking only 3.1ms on average (at most 5ms). In comparison, HCSFQ needs on average 42.3ms (at most 234ms).

#### C. Fairness and Goodput Stability

We first demonstrate that AHAB can effectively enforce fair bandwidth allocation for TCP flows, by simultaneously running 2, 4, or 8 flows sharing a slice with fixed 1Gbps capacity.

Figure 9 shows the goodput over time when we run two TCP flows; although both systems achieve fairness between

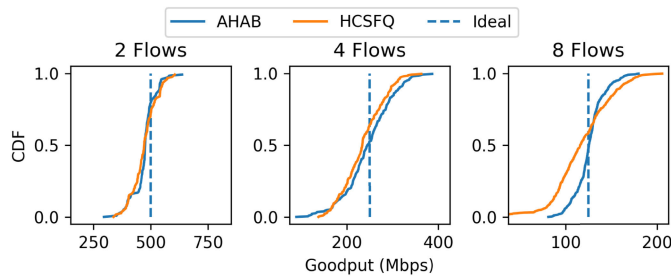


Fig. 10. Cumulative distribution of competing TCP flows's goodput when using AHAB versus HCSFQ.

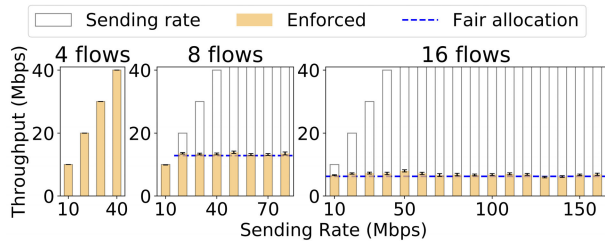


Fig. 11. Given flows with various sending rates, AHAB's approximate probabilistic dropping achieved fair bandwidth allocation within 6% error.

two flows in terms of average goodput over time, AHAB achieve better goodput stability over time. For analysis and comparison, in Figure 10 we plot the cumulative distribution function (CDF) of the TCP goodput of all flows, reported by `iperf3` in 1-second intervals across 60 seconds. In the ideal case, all flows exhibit the same goodput across time, leading to a steeper CDF. The stability achieved by AHAB is comparable to that of HCSFQ: on average, the goodputs of flows enforced by AHAB are within 12.1% of ideal fair share, while those enforced by HCSFQ exhibits 15.5%.

Meanwhile, we also show approximate probabilistic dropping can effectively achieve fair bandwidth allocation for non-TCP traffic that does not react to bandwidth limiting (no congestion control), even with very different sending rates. We let multiple UDP flows share the same slice with fixed 100Mbps capacity, and configured their sending rate to be 10Mbps, 20Mbps, 30Mbps, and so on. In Figure 11, we show the throughput achieved by these flows, when 4, 8, and 16 flows are sent simultaneously. In the latter two cases, the slice is over-utilized and approximate probabilistic dropping kicks in. Although AHAB needs to apply vastly different dropping probabilities for the wide range of sending rates, the resulting allocation is quite fair. On average, the mean throughput achieved is within 4% and 6% of the fair bandwidth allocation target, for 8 and 16 flows respectively. This corresponds to the error of the approximated division using the lookup table.

Figure 12 demonstrates AHAB's support of weighted fairness. We start three groups of flows with weight 1x, 2x, and 4x respectively, with four flows per group, all sharing one slice with 1Gbps capacity. Flows with the same weight achieve the same throughput, proportional to their allocated weight, and their attained throughput averages within 15% and 1% of the weighted fair allocation for TCP and UDP, respectively.

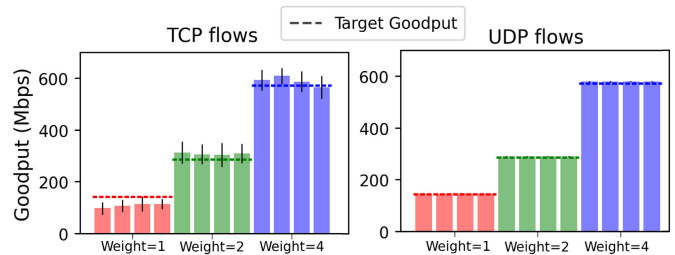


Fig. 12. Goodput of weighted TCP and UDP flows sharing one 1Gbps slice.

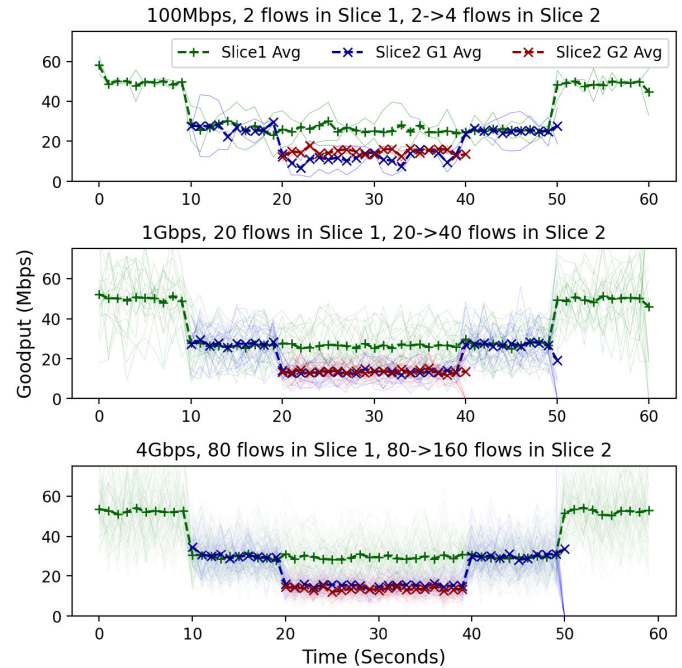


Fig. 13. Three experiments showing two slices sharing a common bottleneck. Slice 1 uses half the bandwidth even when Slice 2 has twice as many flows.

#### D. Inter-Slice Fairness

In Figure 13 we demonstrate that AHAB rapidly adjusts to the changing bandwidth demands of different slices. We set up an experiment where Slice 1 always has  $x$  users (TCP flows), and Slice 2 is initially idle with no user. At  $T = 10$ s  $x$  users in Slice 2 that starts sending, lasting until  $T = 50$ s. At  $T = 20$ s another  $x$  users join Slice 2 and start sending until  $T = 40$ s. In the ideal case, all bandwidth is fully allocated to Slice 1 between  $T = 0-10$ s as well as  $T = 50-60$ s, fairly shared between  $x$  users; the total bandwidth is split in half between Slice 1 and Slice 2 during  $T = 10-50$ s. When more users are added to Slice 2 during  $T = 20-40$ s, users in Slice 2 each get a lower share while users in Slice 1 are not affected.

Figure 13 shows three scenarios: the total bandwidth shared by the two slices are 100, 1000, and 4000 Mbps, respectively, and we also have  $x = 2, 20, \text{ and } 80$  users proportionally. We plot and compare the average goodput attained by the users in each slice, which is also a good indicator of fairness between slices. The bandwidth allocation between slices always quickly converged to fairness (within a few RTTs). When users send UDP traffic instead, AHAB instantly achieves near-perfect fair allocation for all three cases; the result is omitted here.

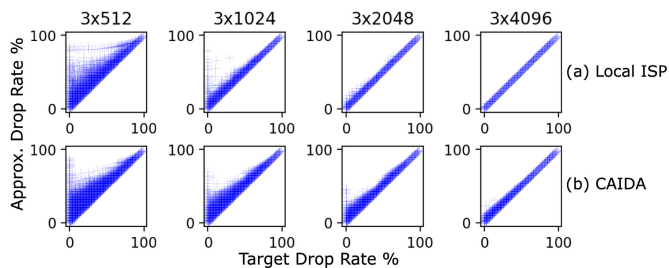


Fig. 14.  $3 \times 4096$  CMS-LPF estimator is sufficient for serving millions of users, with negligible error in drop probability.

### E. Scalability

To evaluate AHAB’s performance at scale, we run trace-based simulation experiments to understand how much memory is needed to support a large number of users.

We collected a 15-minute anonymized traffic trace from the core network of a local Internet Service Provider and played the trace through a Python-based simulator. We treat each of the 5,980,000 unique source-destination IP pairs as an user, and let all users share a single slice with capacity  $C_n$  set to 0.84Gbps, equal to the average throughput of the trace. Due to natural fluctuations in traffic rate, the instantaneous bandwidth demand often exceeds  $C_n$ . The simulator calculates the fair per-user bandwidth limit  $T_n$  for each epoch, and then calculates the “target” probabilistic drop rate  $1 - \min(\frac{T_n}{R_m})$  using the ground-truth per-user sending rate  $R_m$ .

Meanwhile, we also simulate the per-user estimated sending rate  $\widehat{R}_m$  reported by CMS-LPF estimators of different sizes, and use  $\widehat{R}_m$  to calculate the “approximated” drop probability  $1 - \min(\frac{T_n}{\widehat{R}_m})$ . Shrinking the size of CMS-LPF estimator reduces the accuracy of rate estimation, which in turns leads to more error in the drop probability.

As shown by Figure 14(a), using a CMS-LPF estimator with size  $3 \times 512$  led to a fraction of packets with drop rate higher than the target; although most errors lie in over-utilized users, some users with a target drop rate of 0% (under-utilized users) also experience significant packet drops. Meanwhile, CMS-LPF size  $3 \times 4096$  is sufficient to reduce errors to negligible level. We conclude that a CMS-LPF estimator with size  $3 \times 4096$  is sufficient for AHAB to accurately produce per-user rate estimate for the 5,980,000 unique users in our trace.

We also run the same simulation using five minutes of CAIDA Anonymized Internet Trace 2018 [4]. The trace has an average throughout of 3.5Gbps and has 7,300,000 unique flow 5-tuples. We obtain similar results, as shown in Figure 14(b).

Now we analyze the switch hardware resources used by AHAB, and specifically focus on the memory used by the CMS-LPF estimator. As shown in Table II, a small  $2048 \times 3$  CMS-LPF estimator only costs a small fraction (1.56%) of all stateful memory available on the switch hardware, yet it already supports accurately enforcing bandwidth limit for 3 million users. We can fit a much larger sketch than what we used in the prototype: allocating a  $16384 \times 3$  CMS-LPF estimator costs 8.85% of the available memory. Assuming similar traffic skewness as in our ISP trace, a single programmable switch can support 23.9 million devices across all slices, which is sufficient for many application scenarios.

TABLE II  
MEMORY UTILIZATION AND SUPPORTED NUMBER OF USERS W.R.T. DIFFERENT SIZES OF THE CMS-LPF ESTIMATOR

| CMS-LPF Dimension | Hardware Memory Utilization | Supported # of Users |
|-------------------|-----------------------------|----------------------|
| 2048x3 (24KB)     | 1.56%                       | 2,990,000            |
| 4096x3 (48KB)     | 2.60%                       | 5,980,000            |
| 16384x3 (768KB)   | 8.85%                       | 23,920,000 (est.)    |

TABLE III  
UTILIZATION OF OTHER SWITCH HARDWARE RESOURCES

| Resource    | Instr. Words | Hash Units | TCAM |
|-------------|--------------|------------|------|
| Utilization | 28.6%        | 37.5%      | 5.2% |

As for the number of slices, our prototype program supports up to 16,000 slices. The Tofino v1 switch supports 3.2Tbps aggregated throughput, which can be shared among 2,000 downstream base stations. It is possible to expand further by adding more entries to the slice lookup table and allocating more per-slice bandwidth demand trackers, as they only occupy a small fraction of the total data-plane memory usage (with the majority being the CMS-LPF rate estimator). The primary limiting factor on the number of slices is control-plane speed, as supporting  $N$  slices requires the control plane to read  $N$  demands and write  $N$  capacities per update.

We also report other resource utilization of AHAB data-plane program in Table III. These resource usages are constant regardless of the number of slices and users.

## VII. RELATED WORK

### A. Fair Queuing Using Estimated Rate

Core Stateless Fair Queuing [17] is a network architecture where edge nodes estimate the rate of incoming flows and attach the rate to packets, while core nodes in the network choose a fair per-flow rate and enforce it using probabilistic dropping. It requires maintaining per-flow state to estimate sending rates. Approximate Fairness through Differential Dropping [12] uses a shadow buffer that holds recent packets to approximately derive per-flow rates, and similarly performs probabilistic dropping. It is not straightforward to implement a large shadow buffer given the computational constraints present in today’s high-speed programmable switches. Also, both works require one dedicated hardware queue per “slice” (group of flows).

### B. Rank-Based Scheduling

In Push-In, First-Out (PIFO) queues, each packet is pushed in with a certain rank, and packets with the highest rank are transmitted first. Admit-In, First-Out [21] and SP-PIFO [2] both approximate the behavior of a PIFO queue on commodity programmable switches. AIFO uses a sample of recently admitted packets to estimate the rank distribution of packets in the queue, which is used to decide a threshold and reject low-ranked packets from being admitted. Meanwhile, SP-PIFO uses an array of strict-priority queues and dynamically adjust the mapping from ranks to queues using estimated quantile distribution of ranks. These works both assume an oracle which assigns ranks to packets.

### C. Fair Queuing in the Data Plane

Approximate Fair Queuing [14] implements scalable per-flow fair queuing by splitting traffic into calendar epochs. This design requires rapidly rotating the priority between multiple queues to serve different future epochs, and does not support a multi-layer scheduling hierarchy. Gearbox [8] proposed a new hardware design specifically supporting multi-level calendar queuing. Meanwhile, Hierarchical Core-Stateless Fair Queuing [22] extends CSFQ [17] and uses Addictive Increase, Multiplicative Decrease (AIMD) to iteratively find a fair per-user (per-tenant) sending rate limit using queue congestion status feedback. Although it can support multiple layers of scheduling hierarchy, its dependency on per-user memory for estimating per-user sending rates hurts scalability. The AIMD process also takes a relatively long time to adapt when the fair rate increases. Cebinae [20] uses leaky-bucket filters to estimate per-flow rate and enforce fairness by “taxing” the heavy flows, however it takes several seconds to converge to fair allocation. Nimble [18] implements precise TCP flow rate limiting by simulating queue draining in the data plane. However, it only supports fixed rates set by the control plane and requires per-flow memory. Instead, our work automatically adjusts and enforces fair per-user bandwidth limit within milliseconds timescale for millions of users.

### VIII. CONCLUSION

We present AHAB, a data-plane hierarchical fair bandwidth limit enforcer. Using a novel approximate data structure, AHAB scales to millions of users across thousands of network slices. AHAB exploits approximate arithmetic to implement interpolation-based bandwidth limit update fully within the data plane, leading to fast convergence. Evaluation shows that AHAB converges to a fair allocation within 3.1ms, 13x faster than prior work, without sacrificing fairness or stability. The authors have provided public access to their code at <https://github.com/Princeton-Cabernet/AHAB>.

### ACKNOWLEDGMENT

The authors sincerely Zhuolong Yu for his assistance in evaluation experiments. They also thank Henry Birge-Lee and Yufei Zheng for their helpful comments and feedback for the project.

### REFERENCES

- [1] (2021). *P4—Programming Protocol-Independent Packet Processors*. [Online]. Available: <https://p4.org>
- [2] A. G. Alcoz, A. Dietmüller, and L. Vanbever, “SP-PIFO: Approximating push-in first-out behaviors using strict-priority queues,” in *Proc. 17th USENIX Symp. Networked Syst. Design Implement.*, Feb. 2020, pp. 59–76.
- [3] J. C. R. Bennett and H. Zhang, “Hierarchical packet fair queuing algorithms,” in *Proc. Appl., Technol., Archit., Protocols Comput. Commun.*, New York, NY, USA, Aug. 1996, pp. 143–156.
- [4] CAIDA. (Jul. 19, 2018). *The CAIDA UCSD Anonymized Internet Traces 2018*. [Online]. Available: [https://www.caida.org/data/passive/passive\\_dataset.xml](https://www.caida.org/data/passive/passive_dataset.xml)
- [5] X. Chen et al., “Fine-grained queue measurement in the data plane,” in *Proc. 15th Int. Conf. Emerg. Netw. Exp. Technol.*, New York, NY, USA, Dec. 2019, pp. 15–29.
- [6] G. Cormode and S. Muthukrishnan, “An improved data stream summary: The count-min sketch and its applications,” *J. Algorithms*, vol. 55, no. 1, pp. 58–75, Apr. 2005.
- [7] S. Floyd and V. Jacobson, “Link-sharing and resource management models for packet networks,” *IEEE/ACM Trans. Netw.*, vol. 3, no. 4, pp. 365–386, Aug. 1995.
- [8] P. Gao, A. Dalleggio, Y. Xu, and H. J. Chao, “GearBox: A hierarchical packet scheduler for approximate weighted fair queuing,” in *Proc. 19th USENIX Symp. Networked Syst. Design Implement.*, Apr. 2022, pp. 551–565.
- [9] S. Hasan et al., “Building flexible, low-cost wireless access networks with Magma,” in *Proc. 20th USENIX Symp. Networked Syst. Design Implement. (NSDI)*, Apr. 2023, pp. 1667–1681.
- [10] R. MacDavid et al., “A P4-based 5G user plane function,” in *Proc. ACM SIGCOMM Symp. SDN Res. (SOSR)*, Oct. 2021, pp. 162–168.
- [11] R. MacDavid, X. Chen, and J. Rexford, “Scalable real-time bandwidth fairness in switches,” in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, May 2023, pp. 1–10.
- [12] R. Pan, L. Breslau, B. Prabhakar, and S. Shenker, “Approximate fairness through differential dropping,” *ACM SIGCOMM Comput. Commun. Rev.*, vol. 33, no. 2, pp. 23–39, Apr. 2003.
- [13] R. Shah, V. Kumar, M. Vutukuru, and P. Kulkarni, “TurboEPC: Leveraging dataplane programmability to accelerate the mobile packet core,” in *Proc. Symp. SDN Res.*, Mar. 2020, pp. 83–95.
- [14] N. K. Sharma, M. Liu, K. Atreya, and A. Krishnamurthy, “Approximating fair queuing on reconfigurable switches,” in *Proc. 15th USENIX Symp. Networked Syst. Design Implement.* Renton, WA, USA: USENIX Assoc., Apr. 2018, pp. 1–16.
- [15] A. Sivaraman et al., “Programmable packet scheduling at line rate,” in *Proc. ACM SIGCOMM Conf.*, New York, NY, USA, Aug. 2016, pp. 44–57.
- [16] G. Soós, D. Ficzer, P. Varga, and Z. Szalay, “Practical 5G KPI measurement results on a non-standalone architecture,” in *Proc. IEEE/IFIP Netw. Oper. Manage. Symp.*, Apr. 2020, pp. 1–5.
- [17] I. Stoica, S. Shenker, and H. Zhang, “Core-stateless fair queuing: Achieving approximately fair bandwidth allocations in high speed networks,” in *Proc. ACM SIGCOMM*, Oct. 1998, pp. 118–130.
- [18] V. S. Thapeta, K. Shinde, M. Malekpourshahraki, D. Grassi, B. Vamanan, and B. E. Stephens, “Nimble: Scalable TCP-friendly programmable in-network rate-limiting,” in *Proc. ACM SIGCOMM Symp. SDN Res. (SOSR)*, Oct. 2021, pp. 27–40.
- [19] D. Xu et al., “Understanding operational 5G: A first measurement study on its coverage, performance and energy consumption,” in *Proc. Annu. Conf. ACM Special Interest Group Data Commun. Appl., Technol., Archit., Protocols Comput. Commun.*, New York, NY, USA, Jul. 2020, pp. 479–494.
- [20] L. Yu, J. Sonchack, and V. Liu, “Cebinae: Scalable in-network fairness augmentation,” in *Proc. ACM SIGCOMM Conf. (SIGCOMM)*, Aug. 2022, pp. 219–232.
- [21] Z. Yu et al., “Programmable packet scheduling with a single queue,” in *Proc. ACM SIGCOMM Conf.*, Aug. 2021, pp. 179–193.
- [22] Z. Yu, J. Wu, V. Braverman, I. Stoica, and X. Jin, “Twenty years after: Hierarchical core-stateless fair queuing,” in *Proc. 18th USENIX Symp. Networked Syst. Design Implement.*, Apr. 2021, pp. 29–45.