

Computing Call Admission Capacities in Linear Networks

E. G. Coffman, Jr.[†], Anja Feldmann[‡], Nabil Kahale[‡], and Bjorn Poonen^{#1}

[†]Bell Labs, Lucent Technologies, Murray Hill, NJ 07974

[‡]AT&T Labs-Research, Florham Park, NJ 07932

[#]Dept. of Mathematics, University of California, Berkeley, CA 94720

April 8, 1999

ABSTRACT

We study call admission rates in a linear communication network with each call identified by an arrival time, duration, bandwidth requirement, and origin-destination pair. Network links all have the same bandwidth capacity, and a call can be admitted only if there is sufficient bandwidth available on every link along the call's path. Calls not admitted are held in a queue, in contrast to the protocol of loss networks. We determine maximum admission rates (capacities) under greedy call allocation rules such as First Fit and Best Fit for several baseline models, and prove that the natural necessary condition for stability is sufficient. We establish the close connections between our new problems and the classical problems of bin packing and interval packing. In view of these connections, it is surprising to find that Best Fit allocation policies are inferior to First Fit policies in the models studied.

1 Introduction

An effective, easily implemented call admission policy is a key element in the construction of many communication networks. However, the design and analysis of such policies typically poses difficult problems, particularly within multimedia networking technologies such as ATM (Asynchronous Transfer Mode). To study these problems in a linear network with given link bandwidths, we adopt a model in which a call request is identified by its source and destination, its arrival time, its bandwidth requirement, and its duration. Calls can be scheduled across a given link at the same time if and only if their cumulative bandwidth requirement does not exceed the link bandwidth. The term 'bandwidth' need not be taken literally; the model can be applied to other interpretations such as an 'effective bandwidth' concept (see e.g., [EM93]).

Most modern research on call admission policies has concentrated on loss systems where the decision to accept or reject a call is made once and for all at the time of

¹The fourth author is partially supported by an NSF Mathematical Sciences Postdoctoral Research Fellowship.

the call's arrival. Considerable effort has been put into the competitive analysis of policies in this setting (see e.g., [BFL96] for a recent paper with many references). For the more relevant stochastic analysis, see [Kel91] for an extensive survey of the research on loss networks, including one-dimensional networks corresponding to those introduced here. In contrast, our model assumes that calls can be delayed (placed into a queue) and admitted later under more favorable traffic conditions. In the recent work of Feldmann et al [FMS⁺95, Fel95], such systems are studied at length. To justify the added mechanism of call admission with delay, Feldmann et al argue that, by eliminating the call retrial/resubmission traffic, especially that generated by 'patient' computers, admission policies exert greater control over a smaller offered load. (A more thorough discussion can be found in the papers just cited.)

Combinatorial analysis [FMS⁺95, Fel95] and simulations [Fel95] indicate that call admission with delay has many advantages. However, very little is known about the performance of such systems within stochastic models, which is our point of departure here. We impose probability laws on call characteristics and then compute call admission capacities for certain one dimensional (linear) networks; in particular, we determine in each case the maximum call arrival rate such that the number of delayed calls at any time remains finite in expected value. Our results contribute to mathematical foundations by establishing the tractability of baseline models of communication systems. In addition, together with simulations, our results give practical insights into the behavior, sometimes unexpected, of greedy call admission policies.

The general call-admission problem is defined on a linear network, where there are $n+1$ nodes and n links, each with bandwidth capacity 1. Call requests arrive in a Poisson stream at total rate λ_0 , with the m -th call having bandwidth requirement b_m , $0 \leq b_m \leq 1$, and duration $d_m > 0$. Bandwidth requirements and durations are independent, and both $\{b_m\}$ and $\{d_m\}$ are i.i.d. sequences. Time is scaled so that $E[d_m] = 1$. A call's source and destination are given by a pair (i, j) , where i , the source, is a uniform random sample from the set $\{1, \dots, n+1\}$ of nodes, and j , the destination, is a uniform random sample from $\{1, \dots, n+1\} - \{i\}$. Note that, in an equivalent set-up, we could have calls arrive at the $n+1$ nodes in independent Poisson streams, each at rate $\lambda_0/(n+1)$; as before, a call would have a destination chosen independently at random from $\{1, \dots, n+1\} - \{i\}$, where i is the call's source. Queues of delayed calls might be allowed to form at individual nodes, but our interest would still be restricted to the total length of all queues.

We assume that the network has a duplex structure in that two calls can use any link at the same time so long as their sources are on opposite sides of the link. Thus, the left-to-right traffic (calls with sources to the left of their destinations) does not impede the right-to-left traffic, and both types of calls have total arrival rate $\lambda = \lambda_0/2$.

In the next three sections, we examine special cases of the general model, as follows. We begin in Section 2 by analyzing a greedy algorithm for the case where $b_m = d_m = 1$, $m = 1, 2, \dots$, and thus where bandwidth partitioning (packing) problems

are avoided. The problem reduces to scheduling calls under the constraint that none of those scheduled for the same time unit overlap; we call this *interval scheduling*, since calls are completely defined by discrete intervals $\{i, i + 1, \dots, j\}$. Lagarias, Odlyzko, and Zagier [LOZ85] aptly refer to this set-up as a ‘disjointly shared network’; they study the loss version and extend the analysis to the case where any number up to some fixed $K > 1$ of calls can overlap.

Section 3 replaces the unit durations with independent exponential durations (with parameter 1 by our convention $E[d_m] = 1$), but to compensate for the greater difficulty, it restricts the network to $n = 3$ nodes (2 links). Section 4 introduces bandwidth packing problems by taking the b_m as independent random samples from the uniform distribution on $\{1/k, \dots, (k-1)/k\}$. The concession to greater difficulty in this case is the restriction to a single link ($n = 1$) and unit duration calls. Stability theorems along with discussions of applications, simulation results, and conjectures for more general models also appear in Sections 2-4. Proofs of the stability theorems appear in Sections 5-7.

2 Greedy Interval Scheduling; Constant Durations

We assume here that $b_m = d_m = 1$ for all calls. Call admissions and hence call departures occur only at integer times and follow a simple left-to-right greedy policy. At the end of each time unit, the greedy policy searches the left-to-right calls in order of increasing source index admitting each call encountered, if any, that does not overlap calls already admitted. (Two calls overlap if and only if they have a link in common and have sources on the same side of the link.) The policy then finishes with a similar scan of the right-to-left calls in order of decreasing source index.

To determine conditions for finite expected queue lengths, we focus on just the left-to-right traffic; the independent right-to-left traffic must satisfy the same conditions by symmetry. Calls arrive at rate λ and each is equally likely to be any one of the $\binom{n+1}{2}$ source-destination pairs (i, j) , $1 \leq i < j \leq n + 1$. Consider the traffic on the middle link if n is odd or on either of the two middle links if n is even. A routine count shows that a fraction $(\lfloor n/2 \rfloor + 1)/(2\lfloor n/2 \rfloor + 1)$ of the calls needs this most heavily used link. Thus,

$$\lambda < \frac{2\lfloor n/2 \rfloor + 1}{\lfloor n/2 \rfloor + 1} \tag{1}$$

is an obvious necessary condition for stability of the left-to-right subsystem and hence the entire system. But is (1) sufficient? By mapping this problem into one of Kahale and Leighton’s [KL95] seemingly quite different packet routing problems, we prove the following affirmative answer in Section 5.

Theorem 1. *Let $b_m = d_m = 1$, $m \geq 1$. Then under the greedy policy the expected queue length at time t is uniformly bounded for all t if and only if (1) holds.*

The structure of the interval scheduling problem is a discretization of interval *parking* problems [JSW90, CMP94] in which some subset of a given collection \mathcal{C} of n subintervals $[a_i, b_i] \subseteq [0, 1]$, $1 \leq i \leq n$, is to be parked in $[0, 1]$. (An item will be parked in $[0, 1]$ when it specifies a subinterval $[a, b]$ in which it must be placed; it will be *packed* in $[0, 1]$ if it specifies only a length and can be placed into any available subinterval of the same length.) Algorithms have been studied for parking a maximum disjoint subset of \mathcal{C} into $[0, 1]$ and for parking a disjoint subset that minimizes wasted space, i.e., the Lebesgue measure of the points in $[0, 1]$ not covered by any parked subinterval. The given subintervals are independent and determined by two independent uniform random draws from $[0, 1]$. The objectives of the algorithm analyses are large- n estimates of the expected number parked and the expected wasted space. Here, our more general stochastic setting, in which calls (intervals) arrive and depart at random, leads to stability issues that have no counterpart in the earlier models.

3 Greedy Interval Scheduling; Exponential Durations.

Let the constant (unit) durations of Section 2 be replaced by i.i.d. exponentially distributed durations with mean 1. The following greedy policy, which we call First Fit, is adopted for scheduling calls. When a call arrives, it is admitted immediately if it does not overlap a call already in progress; otherwise, it joins the end of a queue of waiting calls. When a call departs, the policy scans the queue in arrival order admitting calls whenever one is encountered that does not overlap any call already admitted.

The new problem is open for general n , but for $n = 2$ we prove in Section 3 that (1), which becomes $\lambda < 3/2$, again yields the desired behavior.

Theorem 2. *Let $b_m = 1$, $m \geq 1$, let $n = 2$, and let the d_m be independent, exponentially distributed (mean-1) durations. Then under the First Fit scheduling rule, the expected queue length at time t is uniformly bounded for all t if and only if $\lambda < 3/2$.*

Our proof technique is based on drift analysis; it reduces chiefly to finding a suitable potential (test or Lyapunov) function with negative drift (see e.g., [Haj82, MT93], which covers the theory needed here and has many references to applications). Results of extensive simulations suggest that (1) is also sufficient for every $n \geq 3$. The continuous limit, $n \rightarrow \infty$, normalized on $[0, 1]$ was simulated directly, i.e., we simulated the limit $n \rightarrow \infty$ of the linear network put on the real line at points $0, 1/n, \dots, (n-1)/n, 1$. Random calls were the intervals bounded by two independent uniform random draws from $(0,1)$. The results are illustrated in column 1 of Table 1 and suggest strongly that First Fit was stable for any $\lambda < 2$.

Another algorithm of obvious interest for call admission is Best Fit. Best Fit always admits the largest waiting call (i.e., a call needing the largest number of links) that can fit in an available sequence of links; at each decision point Best Fit iterates this rule until no further calls can be admitted. For general n , Best Fit needs a tie breaking rule to decide between intervals of the same length. Tie breaking rules that minimize congestion favor intervals closer to the middle of $[1, n + 1]$, i.e., intervals that overlap (interfere with) a greater number of other intervals.

For the continuous relaxation under Best Fit, which is as above for First Fit, the tie breaking rule is unimportant. In a surprising comparison with the First Fit rule, simulations of the continuous relaxation gave convincing evidence that $\lambda < 2$ is **not** sufficient for stability under Best Fit; the results are illustrated in the second column of Table 1. On the other hand, simulations also suggested that (1) was indeed sufficient under Best Fit when $n = 2$ or 3 , irrespective of the tie breaking rules in effect; we intend to prove these properties in a future paper. The question is then: what is the smallest n for which (1) is no longer sufficient for stability under Best Fit?

According to simulations, the answer is: (1) is not sufficient for stability for any $n \geq 4$ (under any tie breaking rule). Table 1 illustrates the experimental results that support the claim for $n = 4$ (see the middle two columns). We remark that in the Best Fit simulations for $n = 4$, every sample of the queue had at most a few dozen waiting calls specifying intervals other than the middle length-2 interval (2,4); it was only the number of waiting (2,4) calls that grew without bound.

4 Bandwidth Packing

Consider next the problem with unit call durations, but with bandwidths b_m drawn independently and uniformly at random from $\{1/k, \dots, (k-1)/k\}$ for some given $k > 2$. Since calls sharing a link must have a cumulative bandwidth not exceeding the link capacity 1, we have a stochastic bin-packing problem combined with the stochastic parking problem of Sections 2 and 3. The packing problem remains nontrivial even if we remove the parking problem by taking $n = 1$. Under this assumption, our goal is the call admission capacity of the link when calls are packed (admitted) according to the following Best Fit rule.

Packing decisions are made only at integer times. Best Fit begins by packing a waiting call with the largest bandwidth. Best Fit then iteratively packs waiting calls with the largest bandwidth no larger than the available link bandwidth left over by calls already packed. This step is repeated until no waiting calls remain, or all such calls are too big to fit in the remaining bandwidth. The average bandwidth is $1/2$ so $\lambda < 2$ is an obvious necessary condition for finite expected queue lengths when $n = 1$. Section 7 proves, again using drift analysis, the corresponding sufficient condition.

Events	interval				bandwidth	
	First Fit $n = \infty$	Best Fit $n = \infty$	First Fit $n = 4$	Best Fit $n = 4$	First Fit $k = \infty$	Best Fit $k = \infty$
1000000	1674	6863	211	5603	2500	3131
2000000	1645	12842	227	11719	3378	9360
3000000	941	18906	249	17411	3347	14078
4000000	1015	25799	102	21776	2992	16396
5000000	1205	31526	166	26392	3091	20321
6000000	744	38440	405	32610	2766	22389
7000000	793	44111	88	38540	3554	25784
8000000	1317	49889	347	43221	2900	28338
9000000	1665	55953	179	49267	3178	31729
10000000	1097	62774	226	54385	2808	35222

Table 1: Queue lengths after every million events for two interval packing algorithms (columns 1–4) and two bandwidth packing algorithms (columns 5 and 6). Call durations were mean-1 exponentials for all cases. The cases $n = \infty$ refer to the continuous relaxation where intervals are random subintervals of $[0, 1]$; the arrival rate was 1.95 for each such case, just less than 2 as required by (1). For the $n = 4$ cases, the arrival rate was 1.65, just less than the $5/3$ required by (1). Section 4 discusses the last two columns; the case $k = \infty$ for the bandwidth parameter means that bandwidth requirements are i.i.d. uniform random draws from the continuous interval $[0, 1]$. An arrival rate of 1.95 was also chosen for these simulations.

As a partial check of our results the simulations in the last four columns were done independently by two different authors, each using a different approach and coding in a different language with a different random number generator; the results were indistinguishable statistically.

Theorem 3. *Assume a single link ($n = 1$) and unit call durations. Let the bandwidths b_m be i.i.d. samples from the uniform distribution on $\{1/k, \dots, (k-1)/k\}$. The expected queue length under Best Fit is uniformly bounded for all t if and only if $\lambda < 2$.*

Interestingly, this $n = 1$ case gives a new result in an equivalent model of slotted communication systems [CHJR93]. In the latter interpretation, calls are messages and call bandwidths (fractions of a unit bandwidth capacity) are message durations (fractions of a time unit); the link bandwidths in successive unit intervals become the unit-duration slots in which subsets of messages are assigned (or packed) and transmitted. With arriving messages modeled by a discretized Markov process, the analysis in [CHJR93] focuses on the Next Fit algorithm: When a message arrives and finds no messages waiting, it is assigned to (will be sent in) the next time slot. If a message arrives and finds other waiting messages, it is assigned to the latest time slot already allocated at least one message, if it fits in the remaining unallocated time of that slot; otherwise, the message is assigned to the next, as yet unused, time slot (and hence eventually transmitted one time unit later). Our analysis adds to the earlier work by proving a capacity (stability) result for the much more efficient Best Fit packing algorithm, which assigns messages to slots where they fit best; the message-rate capacity under Next Fit is only $3/2$, whereas it is 2 under Best Fit.

Research on the packing problem of this section originated in the work of Kipnis and Robert [KR90], who studied a strict FIFO version of the problem: a call can not be admitted before an earlier arriving call even if there is sufficient bandwidth for it. Strong results are obtained in [KR90], including formulas for throughput and invariant measures; general call-duration distributions are considered and call durations are allowed to depend on bandwidth requirements. (Their results were expressed in the terminology of computer storage problems, which motivated their work; bandwidth was storage, calls were jobs,)

Note that, in a Markov chain underlying the FIFO process, the queue can be specified by giving only the queue length and the size of the call at the head of the queue. However, the sizes of *all* waiting calls, and their arrival order, must be specified in a Markovian state of the First Fit version of the model. This version was recently studied in [CS98]; our stability condition for Best Fit was also proved to be the stability condition for First Fit. The approach in [CS98] is based on the fluid-limit techniques originated by Rybko and Stolyar [RS92] and developed over the past few years by Dai [Dai95] and others (see [DM95, Stol95] for key references).

The case of exponential call durations yields intriguing and important open problems. Simulations illustrated in the last two columns of Table 1 suggest that $\lambda < 2$ remains a necessary and sufficient condition for stability under First Fit, although congestion is increased by the greater variability of call durations. Interestingly, we again find that Best Fit is inferior to First Fit; the arrival rate $\lambda = 1.95$ puts Best Fit “over the edge”

into the instability region. We remark that the size distribution of waiting calls is very different under First Fit and Best Fit. Under First Fit, the probability mass function is increasing, with the large majority of waiting calls being larger than 0.5. Under Best Fit, the sizes of waiting calls typically concentrated in the middle range, but the degree of concentration varied widely, and the distribution seemed not to tend to a limit at all; Best Fit would empty the queue of calls of some size while calls of other sizes accumulated, but later it would work on calls of some other size while calls of the first size accumulated, so that the size distribution would be constantly changing over a long period of time.

The analytical difficulties encountered under exponential call durations are brought out in a very recent article of Dantzer, Haddani, and Robert [DHR99]. They use a fluid-limit approach in an analysis of stability where bandwidth requirements are assumed to be generated by a two-point distribution. Besides stability conditions, a characterization of transient behavior in the unstable case is given.

5 Proof of Theorem 1.

Recall that we have unit durations and unit bandwidth requirements for each call. We map the interval scheduling problem on a linear network of $n + 1$ nodes into the following *packet routing* problem on the same network. Packets arrive at the $n + 1$ nodes, joining FIFO queues, according to independent Poisson processes, each at rate $\lambda_0/(n + 1)$. Each packet arriving at node i specifies a destination uniformly at random from $\{1, \dots, n + 1\} - \{i\}$, independently of all other arrivals. Packets are admitted to the network at integer times only. For each i , the packet, if any, at the head of queue i and destined for a node $j > i$ (respectively $j < i$) is admitted as soon as the link $(i, i + 1)$ (respectively $(i - 1, i)$) is not needed by a packet already on the network and going in the same direction. (Recall that in all our duplex models, left-to-right traffic is independent of right-to-left traffic.) Once admitted, a packet moves from its origin i to its destination j , taking exactly $|j - i|$ time units, one time unit for each link traversed. The motion is nonpreemptive (also called “hot potato”) in contrast to other models where packets can be removed from the network and put back on later.

We analyze only the left-to-right traffic in the packet routing problem, just as in the interval scheduling problem; the total arrival rate of packets that move right is $\lambda = \lambda_0/2$. For convenience, assume that arrivals in the interval scheduling problem are sorted into queues $1, \dots, n$ according to the left end of the interval, i.e., an interval with leftmost node i is added to queue i . Thus, in each time unit the left-to-right greedy algorithm scans the queues in increasing order of index. A trajectory of the interval scheduling process is defined by a sequence of triples (a_i, s_i, d_i) , $i = 1, 2, \dots$, where a_i is the time of arrival of the i -th interval and s_i, d_i are its source and destination, with

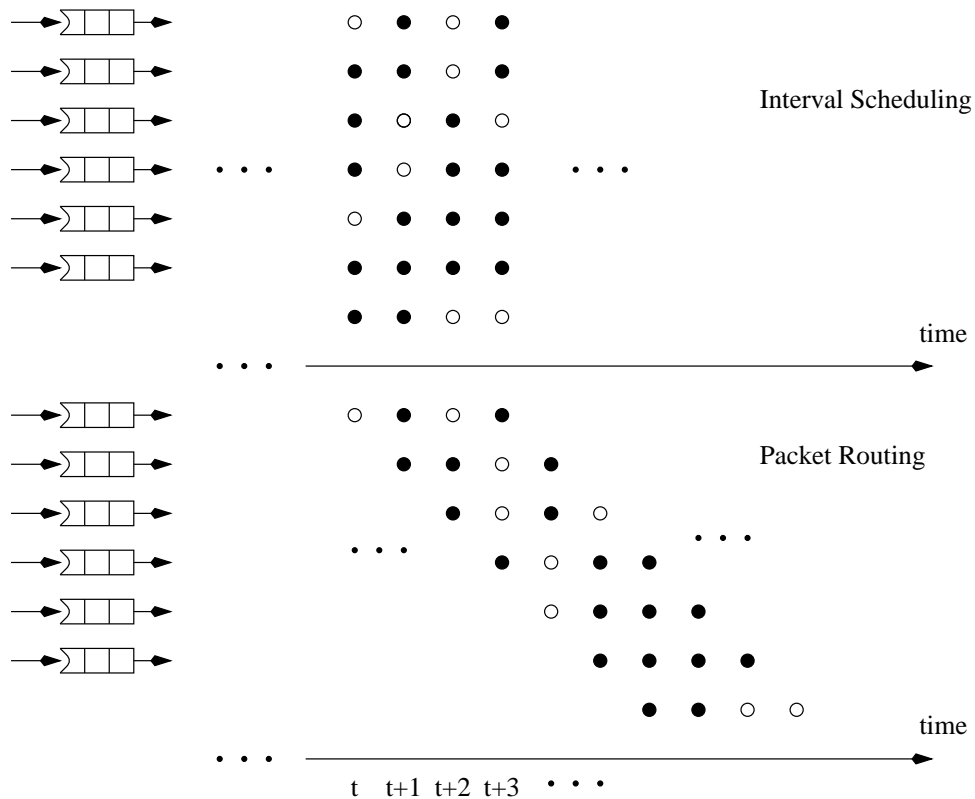


Figure 1: Reduction of Interval Scheduling to Packet Routing, $n = 6$.

$1 \leq s_i < d_i \leq n + 1$. We map this trajectory into a trajectory of the packet routing problem by replacing (a_i, s_i, d_i) with $(a_i + s_i - 1, s_i, d_i)$, $i = 1, 2, \dots$, where the latter triple now denotes the arrival of a packet at time $a_i + s_i - 1$ to be sent from node s_i to node d_i . (No packets arrive in $[0, s_i - 1]$ at queue i in the corresponding packet routing problem.) Figure 1 illustrates the mapping of one trajectory into another; for clarity the linear network has been rotated 90° to the vertical (left-to-right has become top-down), and time is as usual the horizontal dimension.

In the upper half of Figure 1, an isolated vertical sequence of bullets denotes an interval scheduled at the time (column) where it appears. In the lower half, an isolated sequence of bullets going down and rightward along a diagonal traces the motion in time and space of a packet admitted at the beginning of the sequence and delivered at the end. The top-down greedy scheduling of intervals at time t in the upper half of the figure becomes the nonpreemptive motion of packets along the *same* pattern of paths on the diagonal starting at time t . From the mapping $(a_i, s_i, d_i) \rightarrow (a_i + s_i - 1, s_i, d_i)$, an easy induction establishes this observation and the fact that the length of queue

i at time t in the interval scheduling process is equal to the length of queue i at time $t+i-1$ in the modified packet routing process, where the $n+1$ Poisson processes are not initially "turned on" until time $i-1$ for queue i . It follows that expected queue-lengths are bounded in the interval scheduling problem if and only if they are bounded in the modified packet routing problem.

But the time shifts (of at most n) of the Poisson arrival processes for the packet routing problem have only a transient effect. After time n , we have $n+1$ Poisson processes each at rate $\lambda_0/(n+1)$; limiting queue-length distributions, if they exist, will be the same as in the version of the packet routing problem defined in [KL95], where arrival processes all begin at time 0. Moreover, it follows from results of Kahale and Leighton [KL95] that $\lambda < \frac{2\lfloor n/2 \rfloor + 1}{\lfloor n/2 \rfloor + 1}$ is necessary and sufficient for bounded expected queue lengths in the packet routing problem, so Theorem 1 is proved. ■

6 Proof of Theorem 2.

Recall that unit bandwidths, exponential durations, and 3 nodes ($n = 2$ links) are assumed. The 3 possible call requests become the intervals $(1, 2)$, $(2, 3)$, and $(1, 3)$. We say that these calls are type L (left), R (right), and B (big) calls, respectively.

We consider the greedy algorithm that always immediately schedules available links in the 2-link network, if possible, by taking the call closest to the head of the queue that fits without overlapping calls currently scheduled. The problem is to show that when $\lambda < 3/2$, the queue has bounded expected length.

Roughly speaking, our approach will be to determine a potential function ϕ mapping states into the reals such that ϕ stays within a constant factor of the queue length and has negative drift outside some finite set. After proving that the magnitude of a jump in ϕ has an exponential tail probability, we will invoke a result of Hajek [Haj82] to complete the proof that expected queue lengths remain finite in expected value.

The following lemma will help us construct the desired potential function.

Lemma 1. *For any $\delta > 0$ there exists a positive integer r and a continuous function $g(x)$ such that*

1. $g(x) = |x|/2$ for $|x| \geq r$
2. $g(x)$ is increasing for $x \geq 0$ and decreasing for $x \leq 0$
3. $0 \leq g(x+1) - 2g(x) + g(x-1) \leq \delta$ for all x .

Proof. Clearly there exists an even C^∞ bump function $h(x)$ satisfying

1. $0 \leq h(x) \leq \delta$ for all x
2. $h(x) = 0$ for $|x| \geq r$, for some positive integer r
3. $\int_0^\infty h(x) = 1/2$.

Let $j(x) = \int_0^x h(t)dt$ and let $g(x) = \int j(x)dx$ with the constant of integration chosen to make $g(x) = |x|/2$ for $|x| \geq r$. Then by repeated application of the Fundamental Theorem of Calculus,

$$g(x+1) - 2g(x) + g(x-1) = \int_{y=0}^1 \int_{z=-1}^0 h(x+y+z)dy dz,$$

which is between 0 and δ , as desired. ■

Let $\alpha = 3/2 - \lambda$, which we assume to be positive. Choose a function g and a positive integer r as in Lemma 1 for $\delta = \alpha/10$. Choose $\epsilon < \frac{\alpha}{10(r+1)}$, and for convenience adjust ϵ so that ϵ^{-1} is an integer.

For the purpose of the bookkeeping, imagine that each call in progress remains in position in the queue until it has been completed, at which time it is removed (and those calls behind it move forward one notch). A state $X(t)$ consists of the list of calls in the queue, together with the specification of which calls are currently in progress. We now associate various quantities with a state. Let n_B , n_L , and n_R denote the number of calls of each type in the queue (including those currently in progress), and let $\ell = n_B + n_L + n_R$ denote the queue length. For $1 \leq i \leq n_L + n_R$, a_i denotes the position of the i -th call of type L or R in the queue, where position 1 denotes the head of the queue. For $i > n_L + n_R$, we define $a_i = \ell$. Let s equal the length of the run of type B calls at the head of the queue, if a type B call is currently in progress, and let $s = 0$ otherwise. The potential of a state is defined by the formula

$$\phi = n_B + \frac{n_L + n_R}{2} + g(n_L - n_R) - \min(\epsilon^{-1}, s)\epsilon - \sum_{i=1}^r \min(2\epsilon^{-1}, a_i)\epsilon.$$

A glance at the function $g(x)$ shows that the part $\frac{n_L + n_R}{2} + g(n_L - n_R)$ can be thought of as a smoothing of the function $\max(n_L, n_R) = \frac{n_L + n_R}{2} + \frac{|n_L - n_R|}{2}$. The ϵ -terms subtracted at the end are intended to give bonus reductions in potential to certain favorable states, such as those with a long run of type B calls that can be scheduled one after the other with no waste in the network, or those which are close to having such a run of type B calls in the sense that there are few calls of type L or R close to the head of the queue.

Proposition 1. *There is a constant τ such that*

$$\ell/2 - \tau \leq \phi \leq \ell + \tau$$

for every state.

Proof. The continuous function $g(x) - |x|/2$ has compact support, so it is bounded in absolute value by some constant τ_0 . The ϵ -terms are bounded in absolute value by 3. Thus ϕ is within $\tau := \tau_0 + 3$ of

$$\phi' := n_B + \max(n_L, n_R).$$

To complete the proof, simply note that

$$\ell/2 \leq n_B + \frac{n_L + n_R}{2} \leq \phi' \leq n_B + n_L + n_R = \ell.$$

■

The instantaneous drift of the potential ϕ , i.e., the expected rate of change of ϕ at a given time, is denoted by $E(d\phi/dt)$.

Proposition 2. *Let \mathcal{S} be the (finite) set of states for which $\ell \leq 2\epsilon^{-1}$. Whenever the current state is not in \mathcal{S} , $E(d\phi/dt) \leq -3\alpha/10$.*

Proof. First let us calculate an upper bound on the increase in potential that arises from arrivals at the tail of the queue. Note that s and the a_i can only increase from arrivals, so we may disregard their contributions to the change in potential. If a type B call arrives, n_B increases by 1, so the potential increases by at most 1. If a type L call arrives, the potential increases by at most $1/2 + g(x+1) - g(x)$, where $x = n_L - n_R$. If a type R call arrives, the potential increases by at most $1/2 + g(x-1) - g(x)$. The expected number of arrivals of each type is $\lambda/3$, so the total expected rate of increase in potential due to arrivals is at most

$$\begin{aligned} & \frac{\lambda}{3} \left[1 + \left(\frac{1}{2} + g(x+1) - g(x) \right) + \left(\frac{1}{2} + g(x-1) - g(x) \right) \right] \\ &= \frac{\lambda}{3} [2 + (g(x+1) - 2g(x) + g(x-1))] \\ &\leq \frac{\lambda}{3} \left(2 + \frac{\alpha}{10} \right) \\ &= \frac{(3/2 - \alpha)}{3} \left(2 + \frac{\alpha}{10} \right) \\ &\leq 1 - \frac{2}{3}\alpha + \frac{\alpha}{20} \\ &\leq 1 - 3\alpha/5. \end{aligned}$$

To calculate the rate of change in potential due to call completions, we subdivide into cases according to the type of calls currently in progress.

Case 1. A type B call is currently in progress.

If the type B call finishes, n_B drops by 1, s drops by at most 1, and the a_i each drop by at most 1, so the potential drops by at least $1 - (1 + r)\epsilon$, which is at least $1 - \alpha/10$.

Case 2. Calls of type L and R are currently in progress.

If the type L call finishes, then n_L drops by 1, s remains 0, and the a_i drop by at most 1 (they may increase, but this will only help us), so the potential drops by at least $1/2 + g(x) - g(x - 1) - r\epsilon$. (Here again $x = n_L - n_R$.) Similarly, if the type R call finishes, the potential drops by at least $1/2 + g(x) - g(x + 1) - r\epsilon$. Thus, the expected rate of decrease of potential is at least

$$\begin{aligned} & [1/2 + g(x) - g(x - 1) - r\epsilon] + [1/2 + g(x) - g(x + 1) - r\epsilon] \\ &= 1 - [g(x + 1) - 2g(x) + g(x - 1)] - 2r\epsilon \\ &\geq 1 - \alpha/10 - 2\alpha/10 \\ &= 1 - 3\alpha/10. \end{aligned}$$

Case 3. Only a call C of type L is currently in progress.

This case is possible only if $n_R = 0$. If C finishes, then n_L drops by 1, s possibly jumps from 0 to something positive, and each a_i drops by at most 1 (they may jump in the positive direction as well). Let p denote the position of C in the queue. A call ahead of C in the queue cannot be of type L since otherwise it would have been completed before C , and also cannot be of type R , since $n_R = 0$. Thus the first $p - 1$ calls in the queue must all be of type B .

Case 3a. $n_L > r$.

Then if C finishes, the potential drops by at least

$$1/2 + g(n_L) - g(n_L - 1) - r\epsilon = 1/2 + |n_L|/2 - |n_L - 1|/2 - r\epsilon \geq 1 - \alpha/10.$$

Case 3b. $n_L \leq r$ and $p > \epsilon^{-1}$.

Then if C finishes, s jumps from 0 to at least $p - 1 \geq \epsilon^{-1}$ and $\min(\epsilon^{-1}, s)\epsilon$ jumps from 0 to 1, so the potential drops by at least

$$1/2 + [g(n_L) - g(n_L - 1)] + 1 - r\epsilon \geq 3/2 + 0 - r\epsilon \geq 3/2 - \alpha/10 > 1.$$

Case 3c. $n_L \leq r$ and $p \leq \epsilon^{-1}$ and $\ell > 2\epsilon^{-1}$.

Let q be the largest integer with $1 \leq q \leq r$ such that $a_q \leq 2\epsilon^{-1}$. For convenience let $b_i = \min(2\epsilon^{-1}, a_i)$. (These are the coefficients of some of the ϵ -terms.) If C finishes, the new values of b_i , which we denote b'_i are as follows: $b'_i = b_{i+1} - 1$ for $i \leq q - 1$, and

$b'_i = 2\epsilon^{-1} = b_{i+1}$ for $i \geq q$. Then

$$\begin{aligned} \left(\sum_{i=1}^r b'_i \right) - \left(\sum_{i=1}^r b_i \right) &= \left[\sum_{i=1}^{q-1} (b_{i+1} - 1) + \sum_{i=q}^r 2\epsilon^{-1} \right] - \left[\sum_{i=1}^q b_i + \sum_{i=q+1}^r 2\epsilon^{-1} \right] \\ &= 2\epsilon^{-1} - b_1 - (q-1) \cdot 1 \\ &= 2\epsilon^{-1} - p - (q-1) \\ &\geq \epsilon^{-1} - r, \end{aligned}$$

so the potential drops by at least

$$1/2 + [g(n_L) - g(n_L - 1)] + (\epsilon^{-1} - r)\epsilon \geq 1/2 + 0 + 1 - r\epsilon \geq 3/2 - \alpha/10 > 1.$$

Case 3d. $n_L \leq r$ and $p \leq \epsilon^{-1}$ and $\ell \leq 2\epsilon^{-1}$.

Then the state is in \mathcal{S} .

To summarize, we have found that in Case 3, the expected rate of decrease in potential due to call completions is at least $1 - \alpha/10$ in all subcases where the state is not in \mathcal{S} .

Case 4. Only a type R call is currently in progress.

This is exactly like Case 3.

Case 5. No calls are in progress.

This is possible only if the queue is empty. Then the state is in \mathcal{S} .

To summarize, the expected rate of decrease in potential due to call completions is at least $1 - 3\alpha/10$ if the current state is not in \mathcal{S} . Combining this with the expected rate of increase created by arrivals, we have, whenever the current state is not in \mathcal{S} ,

$$E(d\phi/dt) \leq (1 - 3\alpha/5) - (1 - 3\alpha/10) = -3\alpha/10.$$

■

We need also to control the jumps in ϕ , but this is easy: the change in potential $\Delta\phi$ when an event (arrival or completion) occurs is clearly uniformly bounded.

The proof of Theorem 2 will conclude by applying a well-known stability-type result for processes with a discrete time parameter. To prepare, we now construct a Markov chain which for our purposes is equivalent to the Markov process $\{X(t)\} = \{X(t), t \geq 0\}$ of Propositions 1 and 2. In particular, we *uniformize* $\{X(t)\}$ to obtain a Markov chain $\{X_i\} = \{X_i\}_{i \geq 0}$ on the same state space (see e.g. [Kei79]). The two processes evolve under the same Poisson arrivals of types L , R , and B calls at total rate λ , and they begin in the same initial state. The state transitions $X_i \rightarrow X_{i+1}$ are the same as in $\{X(t)\}$ except that now transitions from a state to itself are allowed. The transitions occur as

follows: if X_i is empty, then with probability $\lambda/(\lambda+2)$ state X_{i+1} results from an arrival, and with probability $2/(\lambda+2)$ there is no change and X_{i+1} is also empty. If X_i has two calls in progress, then with probabilities $\lambda/(\lambda+2)$, $1/(\lambda+2)$, and $1/(\lambda+2)$ the event creating X_{i+1} is respectively an arrival, the departure of the L call, or the departure of the R call. Finally, if only one call is in progress, then with probabilities $\lambda/(\lambda+2)$, $1/(\lambda+2)$, and $1/(\lambda+2)$, state X_{i+1} results respectively from an arrival, completion of the call, or no change at all.

It is easily verified that $X(t) \stackrel{d}{=} X_{N(t)}$, where $N(t)$ is a Poisson (counting) process with rate $\lambda+2$ (see e.g. [Kei79, p. 20]). Thus, Theorem 2 will be proved if we can show that the expected queue length in $\{X_i\}$ stays bounded. For this, we obtain a key property of $\{X_i\}$ from Proposition 2 and the construction of $\{X_i\}$: a bound on the expected change $E[\Delta\phi]$ in the potential function in one step of $\{X_i\}$ is given by the bound in Proposition 2 on the expected rate of change times the expected time $1/(\lambda+2)$ between events in $N(t)$, i.e., in states where the queue length is greater than $2\epsilon^{-1}$,

$$\begin{aligned} E[\Delta\phi] &\leq -\left(\frac{1}{\lambda+2}\right)\frac{3}{10}\alpha \\ &= -\left(\frac{1}{7/2-\alpha}\right)\frac{3}{10}\alpha \\ &\leq -\alpha/12. \end{aligned} \tag{2}$$

This result together with Theorem 13.0.1 in [MT93] shows that the chain X_i is ergodic with a unique stationary distribution. The remainder of the proof shows that this distribution has a finite mean, the approach being the drift analysis of Hajek [Haj82].

Lemma 2. *Let $\{\psi_i\}$ be a sequence of real-valued random variables with $E[\psi_0] < \infty$ and suppose there exist constants $\xi, \gamma > 0$ such that we have the drift condition*

$$E[\psi_{i+1} - \psi_i | \psi_i = z] \leq -\gamma, \text{ for all } z > \xi,$$

and there exist constants $\beta, \eta > 0$ such that we have the jump condition²

$$E \left[e^{\beta|\psi_{i+1}-\psi_i|} | \psi_i = z \right] \leq \eta \text{ for all } z.$$

²It is natural to ask whether the absolute value signs around $\psi_{i+1} - \psi_i$ in the jump condition could be removed, since *a priori* one might expect that it is only large jumps in the *positive* direction that could make the $E[\psi_i]$ unbounded. But in fact the absolute value *is* necessary, even for Markov chains having the nonnegative integers as state space, as the following example shows.

Let $\psi_0 = 0$, and let the transition probabilities be as follows. If $\psi_i = 0$, then ψ_{i+1} is 0 or 1, each with probability $1/2$. If $\psi_i = z \geq 1$, then ψ_{i+1} equals 0 or $z+1$, with probabilities $2/(z+2)$ and $z/(z+2)$, respectively. We have

$$E[\psi_{i+1} - \psi_i | \psi_i = z] = -z/(z+2) < -1/4$$

for $z \geq 1$, and $E[e^{\beta|\psi_{i+1}-\psi_i|}] \leq e$ for all $z \geq 0$. One can then check that the ψ_i converge in law to ψ_∞

Then

$$\limsup_{i \rightarrow \infty} E[\psi_i] < \infty.$$

Proof. Our drift and jump conditions correspond to conditions C1 and C2 of [Haj82], which imply his conditions D1 and D2. Then (2.6) of [Haj82] implies

$$E[e^{\alpha\psi_i} | \psi_0 = z] \leq Ae^{\alpha z} + B \quad (3)$$

for some positive constants A , B , and α depending only on ξ , γ , β , and η . Hence

$$E[\alpha\psi_i | \psi_0 \leq 0] \leq E[e^{\alpha\psi_i} | \psi_0 \leq 0] \leq A + B,$$

$$E[\psi_i | \psi_0 \leq 0] \leq \frac{A + B}{\alpha}. \quad (4)$$

Rewriting (3) as $E[e^{\alpha(\psi_i - \psi_0)} | \psi_0 = z] \leq A + Be^{-\alpha z}$ shows

$$E[\alpha(\psi_i - \psi_0) | \psi_0 > 0] \leq E[e^{\alpha(\psi_i - \psi_0)} | \psi_0 > 0] \leq A + B,$$

$$E[\psi_i | \psi_0 > 0] \leq \frac{A + B}{\alpha} + E[\psi_0 | \psi_0 > 0], \quad (5)$$

but $E[\psi_0 | \psi_0 > 0]$ is finite by assumption. Combining (4) and (5) completes the proof. \blacksquare

Clearly, by (2) and Proposition 1, the drift condition is satisfied by ϕ with $\xi = 2\epsilon^{-1} + \tau$ and $\gamma = \alpha/12$. The jump condition follows trivially from the fact that $\Delta\phi$ is uniformly bounded, so

$$\limsup_{i \rightarrow \infty} E[X_i] < \infty$$

follows from Proposition 1 and Lemma 2. Finally, the uniform boundedness stated in Theorem 2 then follows from the ergodicity of $\{X_i\}$.

where

$$\text{Prob}(\psi_\infty = z) = \begin{cases} \frac{1}{2}, & \text{if } z = 0 \\ \frac{1}{2z(z+1)}, & \text{if } z > 0, \end{cases}$$

and hence $E[\psi_i] \rightarrow \infty$.

7 Proof of Theorem 3.

Recall that communication is now across a single link ($n = 1$) of bandwidth 1, calls have unit durations, and bandwidth requests are chosen independently and uniformly at random from $\{1/k, \dots, (k-1)/k\}$. Time is slotted with slots having length 1, and calls being scheduled cannot straddle different time slots. Calls are scheduled according to best fit: the waiting call requesting the largest bandwidth is scheduled first, then the largest of those that can be accommodated by the remaining bandwidth, etc. Calls can be scheduled simultaneously if and only if the sum of their requested bandwidths is less than or equal to 1; i.e., we do not think of the call bandwidths as having a *location* within the available bandwidth, as we would in a parking problem formulation. The problem is to show that if $\lambda < 2$, then the queue has bounded expected length. Drift analysis will again be our approach.

Before getting into the body of the proof, we need a few elementary probability estimates for the Poisson process and the $M/D/1$ queue (the single-server queue with Poisson arrivals and constant service times). The estimates will all apply to events that occur over intervals with durations proportional to T . We say that such an event occurs with *high probability* if for some constant $\beta > 0$ it occurs with probability $1 - O(e^{-\beta T})$ as $T \rightarrow \infty$. Correspondingly, a *low probability* event has probability $O(e^{-\beta T})$ as $T \rightarrow \infty$ for some $\beta > 0$.

First, estimates of the Poisson distribution prove the following standard result.

Fact 1. *Fix $\lambda, \epsilon > 0$. The number of Poisson arrivals at rate λ that occur in a time interval $[t, t + T]$ is between $(\lambda - \epsilon)T$ and $(\lambda + \epsilon)T$ with high probability.*

The remaining facts that we need all concern the $M/D/1$ queue. Later, we will take $\lambda = 2 - \alpha$ as the arrival rate in the bandwidth packing problem. With this interpretation of α , it is convenient to take the arrival rate $\lambda/2 = 1 - \alpha/2$ for the $M/D/1$ queue in the facts below. The constant service times are assumed to have unit durations. As will be seen, Fact 1 is at the heart of each of the facts below. The facts have two parts. We prove the first part of each fact; similar arguments prove the second parts which we leave to the reader.

Fact 2. *Fix $\delta, \epsilon > 0$ such that $\delta < \alpha/2$.*

(i) Suppose that the queue length at time t is $M \leq \delta T$. Then with high probability the queue length at time $t + T$ is less than ϵT .

(ii) Suppose that $\epsilon < \delta$ and the queue length is $M \geq \delta T$ at time t . Then with high probability the queue length at time $t + T$ is less than $M - \epsilon T$.

Proof of (i) By Fact 1 and the inequality $\delta < \alpha/2$, the queue will empty with high probability in the interval $[t, t + \frac{\alpha}{2}T]$. Thus, (i) will follow if we can show that the net

increase in the queue length over an interval $[t', t' + \gamma T] \subseteq [t, t + T]$ is less than ϵT with high probability for any fixed γ , $0 < \gamma < 1$.

A net increase of ϵT during $[t', t' + \gamma T]$ implies that, for some γ' , $0 < \gamma' \leq \gamma$, the queue length was strictly positive throughout a subinterval of length $\gamma' T$ during which the number of arrivals exceeded the number of departures by at least ϵT . The arrival rate is less than 1 so by Fact 1, this event has low probability. But there are at most T^2 different intervals of length $\gamma' T$ in $[t', t' + \gamma T]$. Thus, for some $\beta > 0$ the net increase in queue length during $[t', t' + \gamma T]$ is less than ϵT with probability $1 - O(T^2 e^{-\beta T})$ and hence with high probability, since $T^2 e^{-\beta T} = O(e^{-\beta' T})$ for any $\beta' < \beta$. ■

Fact 3. Fix $\delta, \epsilon > 0$.

(i) Suppose that $\epsilon\alpha/2 > \delta$ and the queue length at time t is $M \leq \delta T$. Then with high probability the queue empties at some time before $t + \epsilon T$.

(ii) Suppose that $\delta < \alpha/2$ and the queue length at time t is $M \leq \delta T$. Then with high probability the queue will be empty at least once in $[t + T - \epsilon T, t + T]$.

Proof of (i) The number of arrivals during $[t, t + \epsilon T]$ is approximately $(1 - \alpha/2)(\epsilon T)$, in the sense of Fact 1. The number of departures during $[t, t + \epsilon T]$ is exactly ϵT , if the queue does not become empty during this time interval. In this case, the net decrease in queue length would be approximately $(\alpha/2)\epsilon T$, which is impossible, because the original queue length was less than this. ■

Recall that $\alpha = 2 - \lambda$, which we assume is small but positive. In the course of the proof, we will introduce various other constants $\kappa_1, \kappa_2, T, \ell_0$, and assume certain inequalities between them. At the end, we will show that all of these hypotheses can be satisfied.

Let n_i denote the number of calls in the queue of bandwidth exactly i/k . A state, which we again denote by X_i , is given by a vector $(n_1, n_2, \dots, n_{k-1})$. Let $\ell = \sum_{i=1}^{k-1} n_i$ be the queue length. We say that a call is *large* if its requested bandwidth is greater than $1/2$, and we let $M = \sum_{i=(k+1)/2}^{k-1} n_i$ be the number of large calls. (We assume k is odd for simplicity, although the same method of proof works for even k .) Let $k_* = (k - 1)/2$ and for $1 \leq i \leq k_*$ define

$$m_i = \frac{n_{k_*} + n_{k_*-1} + \dots + n_i}{k_* - i + 1}.$$

Finally, let

$$\begin{aligned} \phi_1 &= \max(0, \kappa_1(M - \kappa_2)) \\ \phi_2 &= \max(m_1, m_2, \dots, m_{k_*}) \end{aligned}$$

and

$$\phi = \phi_1 + \phi_2, \tag{6}$$

where κ_1 and κ_2 are positive constants to be specified later.

Proposition 3. *There exist positive constants σ_1 , σ_2 , and τ such that*

$$\sigma_1 \ell - \tau \leq \phi \leq \sigma_2 \ell$$

for all states.

Proof: We have $\phi_1 \leq \kappa_1 M$, and $\phi_2 \leq n_1 + n_2 + \dots + n_{k_*}$, so $\phi \leq \max(\kappa_1, 1)\ell$. On the other hand,

$$\phi \geq \kappa_1(M - \kappa_2) + \frac{n_1 + n_2 + \dots + n_{k_*}}{k_*} \geq \min(\kappa_1, 1/k_*)\ell - \kappa_1 \kappa_2$$

is a lower bound of the desired form.

Proposition 4. *There exist choices for the parameters $\kappa_1, \kappa_2, T, L_0, \gamma > 0$ such that the following negative drift condition holds: for any state in which $\ell > \ell_0$, the expected change in ϕ from time t to time $t + T$ is less than or equal to $-\gamma$.*

Proof: We consider λ (and hence α) as fixed. Define

$$\kappa_1 = \frac{400}{\alpha^2}$$

and

$$\kappa_2 = \frac{\alpha^2}{200}T, \tag{7}$$

where T is tending to infinity.

The increase of ϕ during $[t, t + T]$ is in any case bounded by the constant $\max(\kappa_1, 1)$ times the number of arrivals during that period. The contribution to the expected number of arrivals from a situation occurring with low probability will tend to zero as $T \rightarrow \infty$, since the arrivals are Poisson. Therefore, in our analysis of the expected change of ϕ , any situations occurring with low probability can be disregarded.

If we look only at the large calls in the queue, they arrive with rate $\lambda/2$, and are admitted one at a time as long as the queue contains any, as in the $M/D/1$ queue. Since $\lambda/2 = 1 - \alpha/2 < 1$, the system is stable if we forget about the small calls.

Case 1. $M > 2\kappa_2$ in the state at time t .

Let M' denote the new M at time $t + T$. By Fact 2(ii) with $\delta = \alpha^2/100$ and $\epsilon = \alpha^2/200$, $M' \leq M - \kappa_2$ with high probability, so the expected decrease in ϕ_1 is at least (approximately) $\kappa_1\kappa_2 \geq 2T$. On the other hand, the expected number of arrivals of small calls between time t and $t + T$ is less than T , so the expected increase in ϕ_2 is at most T . Thus the net expected change in ϕ is at most $-T$, which is below $-\gamma$, even for $\gamma = 1$, say, since T is going to be large.

Case 2. $M \leq 2\kappa_2$ in the state at time t .

By Fact 2(i) with $\delta = \alpha^2/100$ and $\epsilon = \alpha^2/200$, we will be in a state with $M' \leq \kappa_2$ at time $t + T$, with high probability. Thus the expected change in ϕ_1 will not make a significant positive contribution to the expected change in ϕ .

By Fact 3(i) with $\delta = \alpha^2/100$ and $\epsilon = \alpha/20$, the time between t and the first emptying of the queue of large calls happens in time less than $(\alpha/20)T$ with high probability. By Fact 3(ii) with $\delta = \alpha^2/100$ and $\epsilon = \alpha/20$, the last time before $t + T$ when the queue is emptied of large calls occurs after time $t + T - (\alpha/20)T$, with high probability. Therefore, the time between the first emptying of the queue of large calls and the last emptying in the time interval $[t, t + T]$ is at least $(1 - \alpha/10)T$ with high probability.

Thus it remains to be shown that ϕ_2 has a negative drift given that the time between the first and last emptyings of the queue of large calls is at least $(1 - \alpha/10)T$. Between these emptyings, the distribution of the large calls processed is the same as the distribution of the large calls arriving. In particular, for each i , $1 \leq i \leq k_*$, the number of slots of bandwidth i/k available alongside these large calls while they are in progress is approximately $\frac{2-\alpha}{2k_*}(1 - \alpha/10)T$, and more precisely, is at least $\frac{2-1.01\alpha}{2k_*}(1 - \alpha/10)T \geq \frac{1-0.6\alpha}{k_*}T$ with high probability, by Fact 1. In addition, there are approximately $\frac{\alpha}{2}(1 - \alpha/10)T$ time intervals when no large call is in progress, and more precisely, at least $\frac{0.99\alpha}{2}(1 - \alpha/10)T \geq 0.49\alpha T$ with high probability, again because of Fact 1.

Let j_1 be such that $\phi_2 = m_{j_1}$. Let $m_{j_2}, m_{j_3}, \dots, m_{j_r}$ be the other m_i which are within $10T$ of ϕ_2 . Then, by definition of the m_i ,

$$\begin{aligned} n_{j_i} &= (k_* - j_i + 1)m_{j_i} - (k_* - j_i)m_{j_i+1} \\ &\geq (k_* - j_i + 1)m_{j_i} - (k_* - j_i)(m_{j_i} + 10T) \\ &\geq m_{j_i} - k_*(10T) \\ &\geq \phi_2 - (10k_* + 10)T \\ &> kT, \end{aligned}$$

the last inequality holding if ℓ_0 is sufficiently large as a function of α and T . (If the last inequality did not hold, then $\phi_2 \leq (k + 10k_* + 10)T$, so

$$\phi \leq \kappa_1 M + (k + 10k_* + 10)T \leq 2\kappa_1\kappa_2 + (k + 10k_* + 10)T$$

which by Proposition 3 would give an upper bound on ℓ .) Note that the string of inequalities giving $n_{j_i} > kT$ is valid when $j_i = k_*$ even though m_{j_i+1} is undefined when $j_i = k_*$, since the term involving it is then zero. Since $n_{j_i} > kT$, there is no chance of the calls in the queue of bandwidth j_i/k being depleted during the time interval $[t, t + T]$, because at most k calls can be in progress during each unit of time.

The expected number of arrivals requesting bandwidth i/k is $\left(\frac{2-\alpha}{2k_*}\right)T$, so with high probability the number of such arrivals is less than $\left(\frac{2-\alpha/2}{2k_*}\right)T$, by Fact 1. Therefore the increase in m_{j_i} due to arrivals is (with high probability) at most $\left(\frac{2-\alpha/2}{2k_*}\right)T$. On the other hand, each slot (including the bandwidth 1 slots) of bandwidth at least j_i/k will be filled in best-fit by one small call of bandwidth at least j_i/k (since the calls of bandwidth j_i/k are not depleted). The number of such slots is at least

$$0.49\alpha T + \sum_{h=j_i}^{k_*} \frac{1 - 0.6\alpha}{k_*} T \geq \frac{k_* - j_i + 1}{k_*} (0.49\alpha T + (1 - 0.6\alpha)T)$$

so the decrease in m_{j_i} due to completions is at least

$$\frac{(1 - 0.11\alpha)T}{k_*}.$$

Thus the net change in each m_{j_i} is a decrease by at least

$$\frac{(1 - 0.11\alpha)T}{k_*} - \left(\frac{2 - \alpha/2}{2k_*}\right)T = \frac{0.14\alpha}{k_*}T.$$

The other m_i , those that were not within $10T$ of the maximum, will not increase by more than $2T$ (because of our high probability assumption on the distribution of arrivals), so ϕ_2 will decrease by at least $\frac{0.14\alpha}{k_*}T$, as desired.

Let us summarize the choices of parameters that should be made: first set $\kappa_1 = 400/\alpha^2$ and $\kappa_2 = \frac{\alpha^2}{200}T$. Next choose T (and hence also κ_2) large enough (as a function of α) so that the events we claim occur with high probability actually do occur with sufficiently high probability so that in the cases in which the events fail to take place, the contribution to the expected drift in potential is negligible. Finally choose ℓ_0 large (as a function of α and T) and choose $\gamma > 0$ small. ■

Recall that the increase of ϕ in $[t, t + T]$ is at most a constant times the number of arrivals during $[t, t + T]$. Since the arrivals are Poisson, it follows that increases in ϕ over $[t, t + T]$ have an exponential tail. Since kT bounds the number of completions in $[t, t + T]$, the magnitude of a decrease in ϕ is bounded. Together with Proposition 5, these observations mean that Lemma 2 applies to ϕ at a sequence of times that are

multiples of T , and so the expected size of ϕ is uniformly bounded at times which are multiples of T . By Proposition 3, the same is true for the expected queue length. Also, the expected increase in the queue length *within* time intervals of length T is bounded, so the expected queue length is uniformly bounded at all times, and Theorem 3 is proved. ■

References

- [BFL96] Yair Bartal, Amos Fiat, and Stefano Leonardi. Lower bounds for on-line graph problems with application to on-line circuit and optical routing. *Proceedings 28th Ann. ACM Symp. Th. Comput.*, 1996, pages 530–540
- [CHJR93] E. G. Coffman, Jr., S. Halfin, A. Jean-Marie, and P. Robert. Stochastic analysis of a slotted, FIFO communication channel. *IEEE Trans. Inf. Th.*, **39**(1993), pages 1555–1566.
- [CMP94] E.G. Coffman, Jr., C.L. Mallows, and B. Poonen. Parking arcs on the circle with applications to one-dimensional communication networks. *Ann. Appl. Prob.*, **4**(1994), pages 1098–1111.
- [CS98] E. G. Coffman, Jr. and A. L. Stolyar. Bandwidth packing. (submitted for publication) 1999.
- [Dai95] J. G. Dai. On the positive Harris recurrence for open multiclass queueing networks: A unified approach via fluid limit models. *Ann. Appl. Prob.*, **5**(1995), pp. 49–77.
- [DHR99] Jean-Francois Dantzer, Mostafa Haddani, and Philippe Robert. On the stability of a bandwidth packing algorithm. Technical Report, INRIA, Domaine de Voluceau, B.P. 105, 78153 Le Chesnay Cedex, France.
- [DM95] J. G. Dai and S. P. Meyn. Stability and convergence of moments for multiclass queueing networks via fluid limit models. *IEEE Trans. Auto. Cont.*, **40**(1995), pages 1899–1904.
- [EM93] Anwar I. Elwalid and Debasis Mitra. Effective bandwidth of general markovian traffic sources and admission control of high speed networks. *IEEE/ACM Transactions on Networking*, 1(3):329–343, 1993.
- [Fel95] Anja Feldmann. *On-line Call Admission for High-Speed Networks*. Ph.D. thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, October 1995.

- [FMS⁺95] Anja Feldmann, Bruce M. Maggs, Jiří Sgall, Daniel D. Sleator, and Andrew Tomkins. Competitive analysis of call admission algorithms that allow delay. Technical Report CMU-CS-95-102, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, January 1995.
- [Haj82] Bruce Hajek. Hitting-time and occupation-time bounds implied by drift analysis with applications. *Adv. Appl. Prob.*, **14**(1982), pages 502–525.
- [JSW90] J. Justicz, E.R. Scheinerman, and P. Winkler. Random intervals. *Amer. Math. Monthly*, **97**(1990), pages 881–889.
- [KL95] N. Kahale and F. T. Leighton. Greedy dynamic routing on arrays. *Proc. 6-th Ann. ACM–SIAM Symp. on Disc. Alg.*. SIAM Press, 1995, pages 558–566.
- [Kei79] Julian Keilson. *Markov Chain Models—Rarity and Exponentiality*. Springer-Verlag, New York, 1979.
- [Kel91] F. P. Kelly. Loss networks. *Ann. Appl. Prob.*, **1**(1991), pages 319–378.
- [KR90] C. Kipnis and Ph. Robert. A dynamic storage process. *Stoch. Proc. Applic.*, **34**(1990), pages 155–169.
- [LOZ85] J. C. Lagarias, A. M. Odlyzko, and D. B. Zagier. Realizable traffic patterns and capacity of disjointly shared networks. *Comput. Networks*, **10**(1985), pages 275–285.
- [MT93] S.P. Meyn and R.L. Tweedie. *Markov Chains and Stochastic Stability*. Springer-Verlag, New York, 1993.
- [RS92] A. N. Rybko and A. L. Stolyar. Ergodicity of stochastic processes describing the operation of open queueing networks. *Problems of Information Transmission*, **28**(1992), pp. 199–220.
- [Stol95] A. L. Stolyar. On the stability of multiclass queueing networks: A relaxed sufficient condition via limiting fluid processes. *Markov Processes and Related Fields*, **1**(1995), pp. 491–512.