

PATTERNS OF BUFFER OVERFLOW IN A CLASS OF QUEUES WITH LONG MEMORY IN THE INPUT STREAM

DAVID HEATH, SIDNEY RESNICK AND GENNADY SAMORODNITSKY

Cornell University

ABSTRACT. We study the time it takes until a fluid queue with a finite, but large, holding capacity reaches the overflow point. The queue is fed by an *on/off* process, with a heavy tailed *on* distribution, which is known to have long memory. It turns out that the expected time until overflow, as a function of capacity L , increases only polynomially fast, and so overflows happen much more often than in the “classical” light tailed case, where the expected overflow time increases as an exponential function of L . Moreover, we show that in the heavy tailed case, overflows are basically caused by single huge jobs. An implication is that the usual $GI/G/1$ queue with finite but large holding capacity and heavy tailed service times, will overflow about equally often *no matter how much we increase the service rate*. We also study the time until overflow for queues fed by a superposition of k iid. *on/off* processes with a heavy tailed *on* distribution, and show the benefit of pooling the system resources as far as time until overflow is concerned.

1. Introduction.

Traffic on data networks, e.g. Ethernet LANs, has characteristics substantially different from those of traditional voice traffic. An important feature of data traffic lies in its dependence structure; traditional models are based on assumptions of short-range dependence (like Poisson arrivals and exponential call lengths), while recent measurement and analysis of data traffic has produced strong indications of long-range dependence and self-similarity. Several empirical studies present statistical evidence for existence of these non-standard dependence structures. See for example Leland, Taqqu, Willinger and Wilson (1993, 1994); Willinger, Taqqu, Leland and Wilson (1995); Crovella and Bestavros (1995); Cunha, Bestavros and Crovella (1995).

Seeking an explanation for the observed long range dependence and self-similarity, Willinger, Taqqu, Sherman and Wilson, (1995) have modeled traffic between a single source and destination as an *on/off* or *packet train* process. In their model, an idealized source alternates between an *on* state, in which it produces data at a constant rate, and an *off* state in which it produces no data. The durations of the *on* and *off* periods are independent; *on* times are identically distributed, and so are *off* times. The data they present indicates that both *on* and *off* times are reasonably well modeled by heavy tailed distributions with shape parameter governing heaviness represented by the parameter α . In one example, $\alpha = 1.7$ and 1.2 respectively for the *on* and *off* periods. A similar conclusion was drawn by Crovella and Bestavros (1995) who in their study of world wide web use found evidence of heavy tails in such things as file lengths, transfer times, and operator idle periods. Other papers dealing with *on/off* and related models for communication systems are Brichet et al (1996), Kella and Whitt (1992), Choudhury and Whitt (1995).

1991 *Mathematics Subject Classification.* 60K25, 90B15..

Key words and phrases. long range dependence, heavy tails, *on/off* models, $G/G/1$ queue, fluid models, long memory, heavy tailed distribution, regular variation, time to hit a level, buffer overflow, maximum workload, weak convergence.

David Heath, Sidney Resnick and Gennady Samorodnitsky were partially supported by NSA Grant MDA904-95-H-1036 at Cornell University. Resnick and Samorodnitsky also received support from NSF Grant DMS-94-00535 at Cornell University. Samorodnitsky also acknowledges support from the United States – Israel Binational Science Foundation (BSF) Grant 92-00074

Various paradigms for *on/off* models can be kept in mind. One is the storage or fluid queue model where the store is filling at rate 1 during an *on* period and the contents are subject to constant release at rate r when the content level is positive. Another paradigm allows one to imagine work entering the system at rate 1 during *on* periods and a server working at rate r . We use either paradigm as is convenient.

In a previous paper we studied the stationary distributions of the simple *on/off* models. In the present paper we study the behavior of the first time the contents process exceeds level L for large levels L . Since this represents the time until "buffer overflow" in an *on/off* system with limited capacity, it is important in understanding the behavior of traffic networks.

The simplest model consisting of a single *on/off* source feeding a single server queue, is defined as follows. Let $\{X_i, i = 1, 2, \dots\}$ be a sequence of iid nonnegative random variables representing *on* periods, and similarly let $\{Y_i, i = 1, 2, \dots\}$ be iid nonnegative random variables representing *off* periods. The *on* and *off* sequences are independent. Let F_{on} be the common distribution of X_i 's, and F_{off} be the common distribution of Y_i 's. The workload arrives in the system at rate 1 during *on* periods (no workload arrives in the system during *off* periods). The service rate is r . That is, whenever the system is nonempty, workload is leaving the system at rate r . The state of the system at time t (its content at time t , the workload in the system at time t) is denoted by $X(t)$, and can be formally defined as follows. For a $t \geq 0$ let

$$(1.1) \quad Z(t) = \begin{cases} 1, & \text{if } \sum_{i=1}^{n-1} (X_i + Y_i) \leq t < \sum_{i=1}^n (X_i + Y_i) + X_n, \text{ for some } n \geq 1, \\ 0, & \text{otherwise.} \end{cases}$$

So $Z(t)$ is the indicator of the the source "being on" at time t . Defining the service rate at state x by

$$r(x) = \begin{cases} r, & \text{if } x > 0, \\ 0, & \text{if } x = 0, \end{cases}$$

the state process $\{X(t), t \geq 0\}$ is defined by

$$(1.2) \quad dX(t) = Z(t)dt - r(X(t))dt.$$

The analogous *GI/G/1* queue can be thought of as model (1.2) with *on* periods shrunk to zero and workload arriving in the system in lumps of size $\{B_i, i \geq 1\}$. In this context, $\{Y_i, i \geq 1\}$ can be thought of as interarrival times. One can take, for example, $B_i = (1 - r)X_i$, as this is the net increase in the state of the system (1.2) after the i th *on* period. However, the discussion below does not depend on this particular form of the offered work. The service rate is still r , so that the actual service time of the i th customer is $B_i/r, i \geq 1$.

It turns out that the system overflow patterns in the *on/off* model, when the *on* times are heavy tailed, are very similar to those of the *GI/G/1* queue when the amount of work B_i is heavy tailed. The computations describing the structure of system overflows in both cases are similar, and they are easier for the *GI/G/1* queue. We will present the detailed arguments only for the more involved case of the fluid model (1.2).

To emphasize the dramatic effect of heavy tailed distributions on system behavior, we contrast the heavy tailed case with the simplest, classical case in Section 3. Consider the fluid queue with F_{on} being exponential with mean μ_{on} , F_{off} being exponential with mean μ_{off} , such that

$$(1.3) \quad \frac{\mu_{\text{on}}}{\mu_{\text{on}} + \mu_{\text{off}}} < r.$$

Clearly, (1.3) is the necessary and sufficient condition for the system (1.2) to be stable. If

$$(1.4) \quad \tau(L) = \inf\{t \geq 0 : X(t) \geq L\}$$

is the time until the system overflows, then based on martingale and Markov methods, we find in Section 3 that

$$(1.5) \quad E\tau(L) = a \left(\exp \left(L \left(\frac{1}{(1-r)\mu_{\text{on}}} - \frac{1}{r\mu_{\text{off}}} \right) \right) - 1 \right) - bL, \quad L \geq 0,$$

where

$$a = \frac{(\mu_{\text{on}} + \mu_{\text{off}})(r\mu_{\text{off}})^2}{(r\mu_{\text{off}} - (1-r)\mu_{\text{on}})^2},$$

and

$$b = \frac{\mu_{\text{on}} + \mu_{\text{off}}}{r\mu_{\text{off}} - (1-r)\mu_{\text{on}}}.$$

Several conclusions are immediate from (1.5). First of all, in the exponentially distributed *on* and *off* periods case, the expected time till the system overflows increases exponentially fast with the system holding capacity. Secondly, the average *off* time, μ_{off} , critically affects the average time until the system overflows since it affects the multiplicative constant as well as the growth rate. Quite different conclusions are reached in the heavy tailed case.

In Section 2 we show that in the case of *on* times having a heavy tailed distribution, the expected time to exceed L is asymptotically the same as the expected time until a single *on* period would cause the contents to exceed L assuming the contents were empty at the start of the *on* time. This is very different from the exponential result (1.5), for in this case the expected time until a single *on* period causes the system to exceed the capacity L is asymptotic to

$$(\mu_{\text{on}} + \mu_{\text{off}}) \exp \left(\frac{L}{(1-r)\mu_{\text{on}}} \right),$$

which is of a larger order of magnitude than (1.5). In the case of heavy tailed *on* periods the expected time until the system exceeds a level grows much slower than the exponential rate of increase seen in (1.5). Furthermore, the fact that in the heavy tailed case the system overflow is caused by a single long *on* period implies that the mean *off* time μ_{off} affects the expected time until overflow only by its effect on a multiplicative factor but does not otherwise influence the growth rate.

A similar conclusion is valid for the $GI/G/1$ queue with heavy tailed amounts of work $\{B_i, i \geq 1\}$. In this case the offered workload exceeds the system capacity L when a single customer brings amount of work reaching L . In particular, the mean interarrival time affects the time until the overflow only as a multiplicative factor, and it *does not depend on the service rate r* (!). This provides intuition about the "failure modes" of such a system.

Precise arguments showing unusual behavior in the heavy tailed case are presented in the next section where we study the maximum of the fluid queue (1.2) over a single "wet period", and use the findings to obtain functional limit theorems for the maximum process of the queue (1.2) and for the hitting time process of the same queue. Section 3 contrasts in detail the behaviors in cases where *on* and *off* distributions have exponentially bounded tails with the heavy tailed case. Tangentially relevant papers on extremes of queues (which typically emphasize Markovian methods and exponential tails) are Iglehart (1972), Asmussen and Perry (1992), Berger and Whitt (1995); Abate, Choudhury and Whitt (1994).

In Section 4 we study the behavior of models with several *on/off* sources and a single server. We show again that in the case $r < 1$ the asymptotic behavior of the time at which the contents process exceeds L is the same as that of the first time that any of the input processes has an *on* period long enough to achieve level L from an empty initial content level. We then compare the behavior of a system of completely separate *on/off* processes with one in which the inputs are pooled and in which the capacity of the system is the sum of the capacities of the separate systems. Our conclusions quantify the benefits of pooling the system resources.

Other papers on multisource models, usually emphasizing Markovian environments, are Anick, Mitra and Sondhi (1982), Prabhu and Pacheco (1995), Pacheco and Prabhu (1996).

2. Level crossing times in single input models.

In this section we consider the extreme values of the contents process specified in (1.2) and the time for the content to cross a level. The fluid or storage model is generated by an alternating renewal process which feeds a reservoir. We represent the renewal sequence as $\{S_n, n \geq 0\}$ with $S_n = \sum_{i=1}^n (X_i + Y_i)$, $n \geq 1$ and for convenience we suppose $S_0 = 0$. Both F_{on} and F_{off} have finite means μ_{on} and μ_{off} and we set $\mu = \mu_{\text{on}} + \mu_{\text{off}}$. During an *on* period, liquid enters at net rate $1 - r$ and during an *off* period liquid is released at uniform rate r . We assure that neither the input rate nor the output rate overwhelms the other by assuming

$$(2.1) \quad 1 > r > \frac{\mu_{\text{on}}}{\mu}.$$

Define $S_n^{(X)} = \sum_{i=1}^n X_i$ and $S_n^{(Y)} = \sum_{i=1}^n Y_i$ and the stopping time

$$(2.2) \quad \bar{N} = \inf\{n > 0 : (1 - r)S_n^{(X)} - rS_n^{(Y)} \leq 0\}$$

so that

$$[\bar{N} = n] = [(1 - r)S_j^{(X)} - rS_j^{(Y)} > 0, j = 1, \dots, n - 1, (1 - r)S_n^{(X)} - rS_n^{(Y)} \leq 0] \in \mathcal{B}(X_i, Y_i, i = 1, \dots, n).$$

Consider $\{X(S_n), n \geq 0\}$. Comparing $X(S_n)$ with $X(S_{n+1})$ we get

$$(2.3) \quad \begin{aligned} X(S_{n+1}) &= (X(S_n) + (1 - r)X_{n+1} - rY_{n+1})^+ \\ &= (X(S_n) + \xi_{n+1})^+, \end{aligned}$$

where $\{\xi_{n+1} = (1 - r)X_{n+1} - rY_{n+1}\}$ is iid. This equation expresses that the change of contents over a renewal interval is the input during the *on* period and the loss during the *off* period. Of course (2.3) is Lindley's equation (Resnick, 1992, page 270; Asmussen, 1987; Feller, 1971) and since (2.1) implies

$$E\xi_1 = (1 - r)\mu_{\text{on}} - r\mu_{\text{off}} = \mu_{\text{on}} - r\mu < 0,$$

we know from standard theory that the process

$$\{W_n\} := \{X(S_n)\}$$

will be stable and $E\bar{N} < \infty$. As is customary, we call $\{W_n\}$ the *queuing process*.

We suppose that

$$(2.4) \quad 1 - F_{\text{on}}(x) = x^{-\alpha}L(x), \quad \alpha > 1, \quad x \rightarrow \infty,$$

where L is a slowly varying function. Note that the process $\{X(t), t \geq 0\}$ is regenerative (cf. Resnick, 1992; Feller, 1971; Asmussen, 1987). One set of regeneration times is

$$\{C_n\} := \{S_n : X(S_n -) = 0\},$$

which are the times when a dry period ends and input commences to fill the store. In order to understand the behavior of the extremes of $\{X(t)\}$, it is natural to study the extremes over a cycle. For this purpose, it is necessary to understand the tail behavior of the distribution of maximum of the queuing process over one cycle. A result about this is stated next.

Proposition 2.1. *For the stable queuing process $\{W_n\}$ satisfying (2.1) and (2.4), the maximum over a cycle has a distribution tail asymptotic to the tail of the on distribution; that is, as $x \rightarrow \infty$*

$$\begin{aligned}
 P\left[\bigvee_{n=0}^{\bar{N}} W_n > x\right] &\sim P[\xi_1 > x]E(\bar{N}) \\
 &\sim P[(1-r)X_1 > x]E(\bar{N}) \\
 (2.5) \quad &\sim (1-r)^\alpha \bar{F}_{\text{on}}(x)E(\bar{N}).
 \end{aligned}$$

The proof of this critical result is deferred to the end of this section. Note the result depends on F_{on} and r but that F_{off} only affects the answer through the multiplicative factor $E(\bar{N})$.

We now look at the extremes of $\{X(t)\}$ over a cycle and examine the distribution tail of $\vee_{0 \leq s \leq C_1} X(s)$ where

$$C_1 = S_{\bar{N}}.$$

Note that \bar{N} is the first downgoing ladder epoch of the random walk

$$\left\{\sum_{i=1}^n \xi_i, n \geq 0\right\} = \{(1-r)S_n^{(X)} - rS_n^{(Y)}, n \geq 0\}$$

associated to the queuing process $\{W_n\}$ and that it is not the downgoing ladder epoch of $\{S_n\}$ which is determining the time scale.

Corollary 2.2. *Assume the contents process $\{X(t)\}$ satisfies (2.1) and (2.4). The distribution tail of the maximum of the contents process over one cycle is asymptotic to the tail of the on distribution; that is, as $x \rightarrow \infty$*

$$(2.6) \quad P\left[\bigvee_{s=0}^{C_1} X(s) > x\right] \sim (1-r)^\alpha \bar{F}_{\text{on}}(x)E(\bar{N}).$$

Note again that F_{off} only affects the answer through the multiplicative factor $E(\bar{N})$.

Proof. Set $M_1 = \bigvee_{s=0}^{C_1} X(s)$. Because of the sawtooth character of the paths of $X(\cdot)$ we have that

$$M_1 = \bigvee_{j=1}^{\bar{N}} \left((1-r)S_j^{(X)} - rS_{j-1}^{(Y)} \right)$$

and therefore

$$M_1 \geq \bigvee_{j=0}^{\bar{N}} \left((1-r)S_j^{(X)} - rS_j^{(Y)} \right) \stackrel{d}{=} \bigvee_{n=0}^{\bar{N}} W_n.$$

Thus

$$\begin{aligned}
 \liminf_{x \rightarrow \infty} \frac{P[M_1 > x]}{\bar{F}_{\text{on}}(x)} &\geq \liminf_{x \rightarrow \infty} \frac{P[\bigvee_{j=0}^{\bar{N}} W_n > x]}{\bar{F}_{\text{on}}(x)} \geq \\
 &= (1-r)^\alpha E(\bar{N})
 \end{aligned}$$

where the last step uses Proposition 2.1.

To get a reverse inequality, choose K such that

$$E((1-r)X_1 - r(Y_1 \wedge K)) < 0$$

which can always be done since

$$E(Y_1 \wedge K) \uparrow EY_1$$

as $K \uparrow \infty$. Then

$$S_j^{(Y,K)} := \sum_{i=1}^j (Y_i \wedge K) \leq S_j^{(Y)}$$

which obviously gives

$$\bigvee_{j=0}^{\bar{N}} \left((1-r)S_j^{(X)} - rS_{j-1}^{(Y)} \right) \leq \bigvee_{j=1}^{\bar{N}} \left((1-r)S_j^{(X)} - rS_{j-1}^{(Y,K)} \right).$$

Also

$$\bar{N} \leq \bar{N}^{(K)} := \inf\{n > 0 : (1-r)S_n^{(X)} - rS_n^{(Y,K)} \leq 0\}$$

and thus we have

$$\begin{aligned} M_1 &\leq \bigvee_{j=0}^{\bar{N}^{(K)}} \left((1-r)S_j^{(X)} - rS_{j-1}^{(Y,K)} \right) \\ &\leq \bigvee_{j=0}^{\bar{N}^{(K)}} \left((1-r)S_j^{(X)} - rS_j^{(Y,K)} + r(Y_j \wedge K) \right) \\ &\leq \bigvee_{j=0}^{\bar{N}^{(K)}} \left((1-r)S_j^{(X)} - rS_j^{(Y,K)} + rK \right). \end{aligned}$$

We therefore have

$$\limsup_{x \rightarrow \infty} \frac{P[M_1 > x]}{\bar{F}_{\text{on}}(x)} \leq \limsup_{x \rightarrow \infty} \frac{P[\bigvee_{j=0}^{\bar{N}^{(K)}} \left((1-r)S_j^{(X)} - rS_j^{(Y,K)} \right) > x - rK]}{\bar{F}_{\text{on}}(x)}$$

and applying Proposition 2.1 to the random walk $\{(1-r)S_j^{(X)} - rS_j^{(Y,K)}, j \geq 0\}$ we get this equal to

$$=(1-r)^\alpha E(\bar{N}^{(K)}) \rightarrow (1-r)^\alpha E(\bar{N})$$

as $K \rightarrow \infty$. This provides the reverse inequality and completes the proof. \square

We are now in a position to discuss the behavior of the extremes of the contents process and also the behavior of the first passage time over a level. For a non-decreasing function $U : (0, \infty) \mapsto (0, \infty)$ define the (left continuous) inverse

$$U^-(x) = \inf\{s > 0 : U(s) \geq x\}, \quad x > 0.$$

Define the non-decreasing process

$$M(t) = \bigvee_{s=0}^t X(s)$$

and the first passage time ($L > 0$)

$$\begin{aligned}\tau(L) &= \inf\{s > 0 : X(s) \geq L\} \\ &= \inf\{s > 0 : M(s) \geq L\} \\ &= M^{\leftarrow}(L).\end{aligned}$$

Standard inversion techniques from extreme value theory (Resnick, 1987, Section 4.4) allow for the simultaneous treatment of the weak convergence properties of $(M(\cdot), M^{\leftarrow}(\cdot))$ as random elements of $D_r(0, \infty) \times D_l(0, \infty)$ where $D_r(0, \infty)$ is the space of right continuous functions with finite left limits and $D_l(0, \infty)$ is the space of left continuous functions on $(0, \infty)$ with finite right limits. Each space is equipped with the M_1 topology.

Theorem 2.3. *Assume the contents process $\{X(t)\}$ satisfies (2.1) and (2.4). Define the quantile function*

$$b(s) = \left(\frac{1}{1 - F_{\text{on}}} \right)^{\leftarrow}(s).$$

Let $\{Y_\alpha(t), t > 0\}$ be the extremal process (Resnick, 1987, Section 4.3) generated by the extreme value distribution

$$\Phi_\alpha(x) = \exp\{-x^{-\alpha}\}, \quad x > 0$$

so that

$$P[Y_\alpha(t) \leq x] = \Phi_\alpha^t(x).$$

Define

$$S_\alpha(t) = \frac{1-r}{\mu^{1/\alpha}} Y_\alpha(t).$$

Then in $D_r(0, \infty) \times D_l(0, \infty)$ as $u \rightarrow \infty$

$$\left(\frac{M(u \cdot)}{b(u)}, \left(\frac{M(u \cdot)}{b(u)} \right)^{\leftarrow} \right) \Rightarrow (S_\alpha, S_\alpha^{\leftarrow}).$$

In particular we get for the first passage process, as $u \rightarrow \infty$

$$(1 - F_{\text{on}}(u))\tau(u \cdot) \Rightarrow Y_\alpha^{\leftarrow}\left(\frac{\mu^{1/\alpha}}{1-r} \cdot\right)$$

and

$$\lim_{L \rightarrow \infty} P \left[\frac{(1-r)^\alpha}{\mu} (1 - F_{\text{on}}(L))\tau(L) \leq x \right] = P[E(1) \leq x] = 1 - e^{-x}, \quad x > 0,$$

where $E(1)$ is a unit exponential random variable. Furthermore, as $L \rightarrow \infty$

$$(1 - F_{\text{on}}(L))E(\tau(L)) \rightarrow \frac{\mu}{(1-r)^\alpha}.$$

Proof. We let $\{\bar{N}_k, k \geq 1\}$ be the iterates of \bar{N} so that \bar{N}_k is the k th downgoing ladder epoch of the random walk $\{(1-r)S_n^{(X)} - rS_n^{(Y)}, n \geq 0\}$. Then by the strong law of large numbers $\bar{N}_k/k \rightarrow E(\bar{N})$ as $k \rightarrow \infty$. We write

$$M(S_{\bar{N}_k}) = \bigvee_{s=0}^{S_{\bar{N}_k}} X(s) = \bigvee_{i=1}^k \left(\bigvee_{s=S_{\bar{N}_{i-1}}}^{S_{\bar{N}_i}} X(s) \right) := \bigvee_{i=1}^k M_i$$

so that $\{M_i, i \geq 1\}$ is iid. From Corollary 2.2, as $x \rightarrow \infty$

$$P[M_1 > x] \sim (1-r)^\alpha \bar{F}_{\text{on}}(x) E(\bar{N})$$

so that as $u \rightarrow \infty$

$$\begin{aligned} uP[M_1 > b(u)x] &\sim (1-r)^\alpha E(\bar{N}) u \bar{F}_{\text{on}}(b(u)x) \\ &\sim (1-r)^\alpha E(\bar{N}) x^{-\alpha}. \end{aligned}$$

Therefore,

$$(2.7) \quad \bigvee_{i=1}^{[ut]} \frac{M_i}{b(u)} = \frac{M(S_{\bar{N}_{[ut]}})}{b(u)} \Rightarrow (1-r)(E\bar{N})^{1/\alpha} Y_\alpha(t).$$

Observe that as $u \rightarrow \infty$

$$(2.8) \quad \frac{S_{\bar{N}_{[ut]}}}{u} \rightarrow \mu E(\bar{N}) t$$

in $C(0, \infty)$. For the renewal sequence $\{S_{\bar{N}_k}, k \geq 0\}$, let

$$\Theta(t) = \inf\{k : S_{\bar{N}_k} \geq t\}$$

be the associated counting function so that as $u \rightarrow \infty$

$$(2.9) \quad \frac{\Theta(ut)}{u} \rightarrow \frac{t}{ES_{\bar{N}}} = \frac{t}{\mu E(\bar{N})}$$

in $C(0, \infty)$. Note the inequalities

$$\frac{M(S_{\bar{N}_{\Theta(ut)-1}})}{b(u)} \leq \frac{M(ut)}{b(u)} \leq \frac{M(S_{\bar{N}_{\Theta(ut)}})}{b(u)}.$$

Now from Billingsley, 1968, Theorem 4.4 and composition

$$\begin{aligned} \frac{M(S_{\bar{N}_{\Theta(ut)}})}{b(u)} &= \frac{M(S_{\bar{N}_{[u \cdot \Theta(ut)/u]}})}{b(u)} \Rightarrow (1-r)(E\bar{N})^{1/\alpha} Y_\alpha(t/\mu E\bar{N}) \\ &\stackrel{d}{=} (1-r)\mu^{-1/\alpha} Y_\alpha(t) =: S_\alpha(t) \end{aligned}$$

in the J_1 topology and we hope the same result is true in the M_1 topology for the family of processes $M(u \cdot)/b(u)$ as $u \rightarrow \infty$. In order to verify this, we need to show

$$(2.10) \quad \frac{M(S_{\bar{N}_{\Theta(ut)}}) - M(S_{\bar{N}_{\Theta(ut)-1}})}{b(u)} \Rightarrow 0$$

in the M_1 topology. For a fixed t we get for any ϵ and large u that

$$\begin{aligned} P \left[\left| \frac{M(S_{\bar{N}_{\Theta(ut)}}) - M(S_{\bar{N}_{\Theta(ut)-1}})}{b(u)} \right| > \eta \right] &\leq P \left[\left| \frac{M(S_{\bar{N}_{\Theta(ut)}}) - M(S_{\bar{N}_{\Theta(u(t-\epsilon))}})}{b(u)} \right| > \eta \right] \\ &\rightarrow P[|S_\alpha(t) - S_\alpha(t-\epsilon)| > \eta] \end{aligned}$$

which goes to 0 as $\epsilon \rightarrow 0$ by the stochastic continuity of S_α (Resnick, 1987, Proposition 4.7). The multivariate analogue needed to prove M_1 convergence is similar.

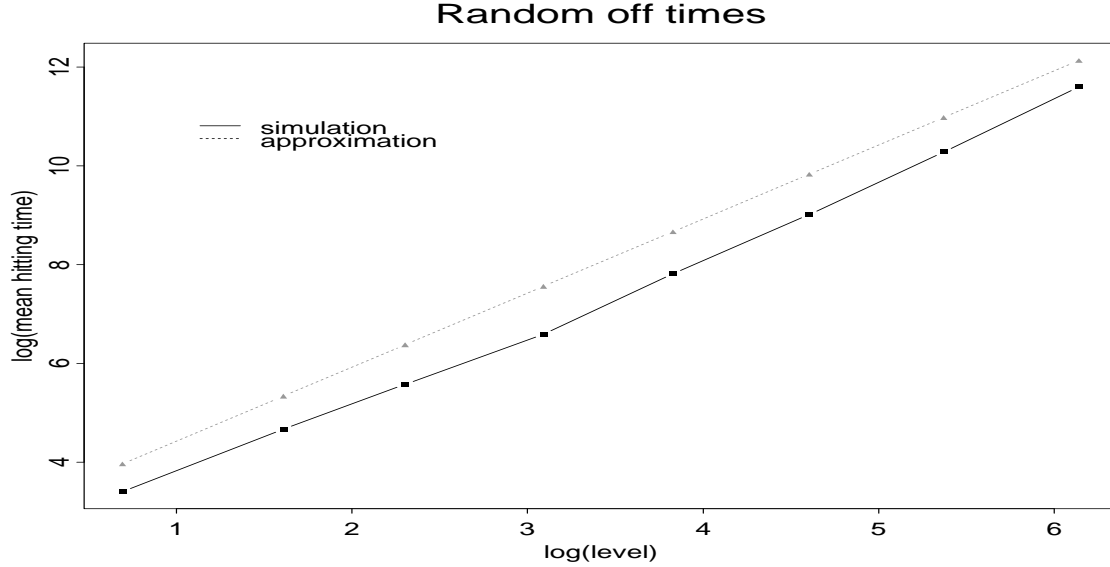


Figure 2.1. Pareto *on/off* periods, $\alpha = 1.5$, $r = .53$.

The weak convergence result for $\tau(\cdot)$ is obtained by taking inverses in the process convergence. Inversion is a continuous operation in the M_1 topology. We note that inverses of extremal processes have exponential marginals (Resnick, 1987) so as $u \rightarrow \infty$

$$\frac{M^\leftarrow(b(u)x)}{u} \Rightarrow S_\alpha^\leftarrow(x)$$

and changing variables $s \mapsto b(u)$ yields

$$\frac{M^\leftarrow(sx)}{\frac{1}{1-F_{\text{on}}(s)}} \Rightarrow S_\alpha^\leftarrow(x).$$

Observe for $y > 0$

$$\begin{aligned} P[S_\alpha^\leftarrow(1) \leq y] &= P[1 \leq S_\alpha(y)] \\ &= P\left[\frac{\mu^{1/\alpha}}{1-r} \leq Y_\alpha(y)\right] \\ &= 1 - \exp\{-y(1-r)^\alpha \mu^{-1}\}. \end{aligned}$$

Finally we consider the result for the expected values. On the one hand, by Fatou's lemma, we get

$$1 \leq \liminf_{L \rightarrow \infty} E \left(\frac{(1-r)^\alpha}{\mu} (1 - F_{\text{on}}(L)) \tau(L) \right).$$

For a reverse inequality, note that

$$\tau(L) \leq S_\nu$$

where

$$\nu := \inf\{n : X_n > L/(1-r)\}$$

so that

$$E\tau(L) \leq E(X_1 + Y_1)E\nu = \mu E\nu.$$

However

$$\begin{aligned} E\nu &= \sum_{n=0}^{\infty} P[\nu > n] = \sum_{n=0}^{\infty} P\left[\bigvee_{i=1}^n X_i \leq L/(1-r)\right] \\ &= \frac{1}{1 - F_{\text{on}}(L/(1-r))} \sim \frac{(1-r)^{-\alpha}}{\bar{F}_{\text{on}}(L)} \end{aligned}$$

and this completes the proof. \square

To illustrate these results we present two modest simulations. For each simulation we supposed F_{on} was Pareto with $\alpha = 1.5$ and $r = .53$. For the first simulation, F_{off} was the same Pareto and for the second simulation F_{off} corresponded to constant *off* times with value 3. We used 500 replications to compute expected hitting times of various levels by simulation and compared these with the approximate mean hitting time given by Theorem 2.3. The levels used for both experiments were 2, 5, 10, 22, 46, 100, 215, 464. The plots use a log scale for both axes. Note that the dotted line appears closer to the solid one when the *off* time is deterministic which may indicate a faster rate of convergence of the the approximation compared to the situation where the approximation has to cope with randomness in the *off* time. However, no systematic investigation has been completed of the rate of convergence.

Deterministic off times

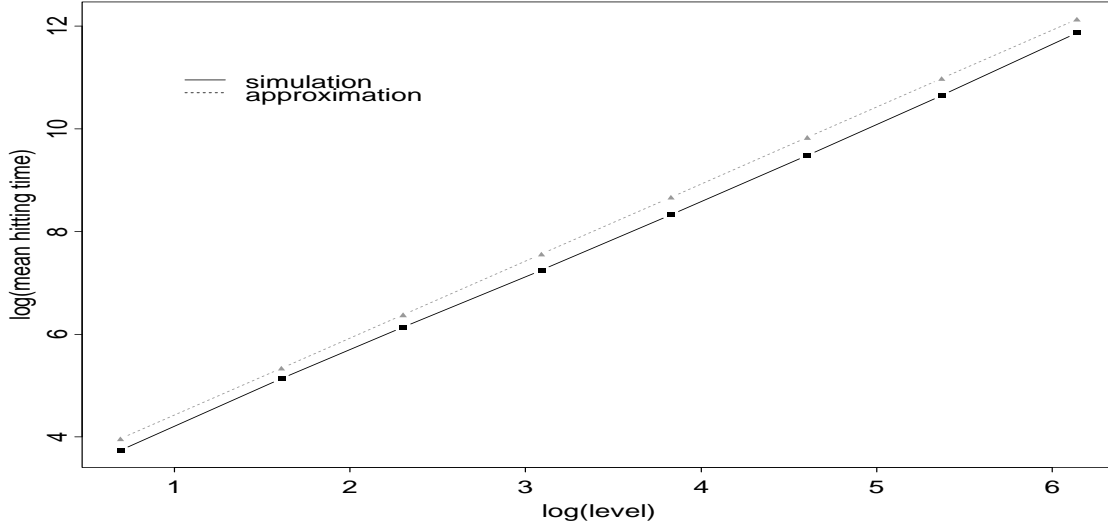


Figure 2.2. Pareto *on* period, deterministic *off* period, $\alpha = 1.5$, $r = .53$.

We also sought experimental evidence to confirm the intuition that in the heavy tailed case the process exceeds a level L because of a very long *on* period. As an additional experiment, we simulated 1000 runs of the process with $\alpha = 1.5$, the *off* distribution concentrated at 3 and $r = .53$. We waited until the process crossed $L = 64$ and then measured the length of the last *on* period X_{on} multiplying by $(1-r)$. We compiled 1000 realizations of

$$\left(\frac{(1-r)X_{\text{on}}}{L}\right) \wedge 3,$$

the truncation by 3 being for the purpose of keeping the data in a comfortable range. The range of the 1000 realizations was $[\cdot 896, 3]$ and 848 observations were at least as large as 1, meaning that in about 85%

of the simulation runs, the process crossed L due to a single large *on* period pushing the process across. A histogram of the data follows showing the preponderance of observations to the right of 1.

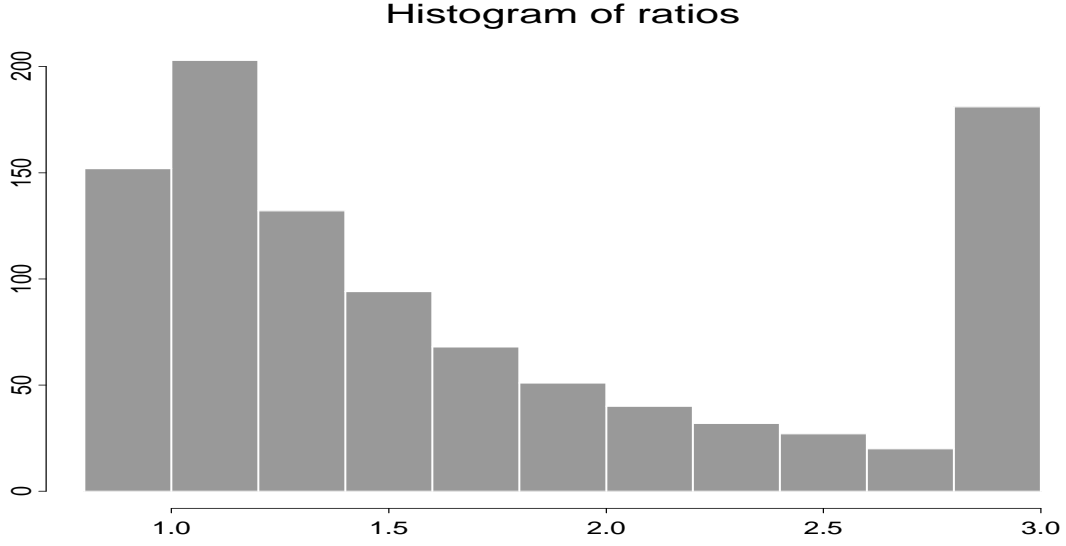


Figure 2.3. $\frac{(1-r)X_{on}}{L} \wedge 3$ for Pareto *on* period, deterministic *off* period, $\alpha = 1.5$, $r = .53$.

It remains to prove Proposition 2.1. We restate it in the following form.

Proposition 2.1'. *Suppose $\{\xi_n, n \geq 1\}$ are iid with $E\xi_1 < 0$ and for $x > 0$*

$$P[\xi_1 > x] =: \bar{F}(x) \sim x^{-\alpha} L(x), \quad \alpha > 1$$

where L is a slowly varying function. Let

$$\{S_n^{(\xi)}, n \geq 0\} = \{0, \sum_{i=1}^n \xi_i, n \geq 1\}$$

be the random walk with negative drift and step distribution F . Suppose

$$\bar{N} = \inf\{n > 0 : S_n^{(\xi)} \leq 0\}$$

is the downgoing ladder epoch of $\{S_n^{(\xi)}\}$. Then as $x \rightarrow \infty$

$$P\left[\bigvee_{n=0}^{\bar{N}} S_n^{(\xi)} > x\right] \sim E(\bar{N})P[\xi_1 > x] = E(\bar{N})\bar{F}(x).$$

Proof. For $x > 0$ let

$$N(x) = \inf\{n > 0 : S_n^{(\xi)} > x\}.$$

Then

$$\begin{aligned}
P\left[\bigvee_{n=0}^{\bar{N}} S_n^{(\xi)} > x\right] &= P\left[\bigvee_{n=0}^{\bar{N}} S_n^{(\xi)} > x, N(x) < \bar{N}\right] \\
&= \sum_{n=1}^{\infty} P[S_i^{(\xi)} \in [0, x], i = 1, \dots, n-1; S_n^{(\xi)} > x] \\
&\geq \sum_{n=1}^{\infty} P[S_i^{(\xi)} \in [0, x], i = 1, \dots, n-1; X_n > x] \\
&= \bar{F}(x) \sum_{n=1}^{\infty} P[S_i^{(\xi)} \in [0, x], i = 1, \dots, n-1] \\
&= \bar{F}(x) E(\bar{N} \wedge N(x)) \\
&\sim \bar{F}(x) E(\bar{N}),
\end{aligned}$$

as $x \rightarrow \infty$. Thus we get

$$(2.11) \quad \liminf_{x \rightarrow \infty} \frac{P[\bigvee_{n=0}^{\bar{N}} S_n^{(\xi)} > x]}{\bar{F}(x)} \geq E(\bar{N}).$$

The reverse inequality is more challenging. Observe that

$$P\left[\bigvee_{n=0}^{\bar{N}} S_n^{(\xi)} > x\right] = P\left[\bigvee_{n=0}^{\bar{N} \wedge N(x)} S_n^{(\xi)} > x\right]$$

and that

$$(2.12) \quad P\left[\bigvee_{n=0}^{\bar{N}} \xi_n > x\right] \sim P\left[\bigvee_{n=0}^{\bar{N} \wedge N(x)} \xi_n > x\right] \sim E(\bar{N}) \bar{F}(x)$$

where the first equivalence follows from Resnick (1986) and second being a minor modification of the first. So we attempt to compare

$$P\left[\bigvee_{n=0}^{\bar{N} \wedge N(x)} S_n^{(\xi)} > x\right] \text{ with } P\left[\bigvee_{n=0}^{\bar{N} \wedge N(x)} \xi_n > x\right].$$

We write as an abbreviation

$$N^* = \bar{N} \wedge N(x).$$

Pick $\lambda > 0$. Then we have

$$\begin{aligned}
(2.13) \quad &P\left[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x\right] - P\left[\bigvee_{n=1}^{N^*} \xi_n > x\right] \\
&= P\left[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x, \bigvee_{n=1}^{N^*} \xi_n \leq x\right] \\
&= P\left[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x, \bigvee_{n=1}^{N^*} \xi_n \leq x - \lambda\right] + P\left[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x, x - \lambda < \bigvee_{n=1}^{N^*} \xi_n \leq x\right] \\
&= I_1(x) + I_2(x).
\end{aligned}$$

From (2.12) we get

$$(2.14) \quad \frac{I_2(x)}{\bar{F}(x)} \leq \frac{P[\bigvee_{i=1}^{N^*} \xi_i > x - \lambda] - P[\bigvee_{i=1}^{N^*} \xi_i > x]}{\bar{F}(x)} \rightarrow 0$$

as $x \rightarrow \infty$. So I_2 is of small order and we turn attention to I_1 .

Pick and fix K such that $[K/2] > \alpha/(\alpha - 1)$ and write

$$\begin{aligned} I_1(x) &\leq P[\bigvee_{i=0}^{N^*} S_i^{(\xi)} > x, \bigvee_{i=1}^{N^*} \xi_i \leq x/K] + P[\bigvee_{i=0}^{N^*} S_i^{(\xi)} > x, x/K < \bigvee_{i=1}^{N^*} \xi_i < x - \lambda] \\ &= I_{11}(x) + I_{12}(x). \end{aligned}$$

For an integer M we further decompose

$$\begin{aligned} I_{12}(x) &= P[\bigvee_{i=0}^{N^*} S_i^{(\xi)} > x, x/K < \bigvee_{i=1}^{N^*} \xi_i \leq x - \lambda, \bigvee_{i=1}^M \xi_i > x/K] \\ &\quad + P[\bigvee_{i=0}^{N^*} S_i^{(\xi)} > x, \bigvee_{i=1}^{N^*} \xi_i \leq x - \lambda, \bigvee_{i=1}^M \xi_i \leq x/K, \bigvee_{i=M+1}^{N^*} \xi_i > x/K, N^* > M] \\ &= I_{121}(x) + I_{122}(x). \end{aligned}$$

For $x > \lambda$

$$\begin{aligned} I_{121}(x) &\leq P[\text{There exists some } n \leq M \wedge N^* \text{ such that } \xi_n > x/K, \bigvee_{m=0}^{N^*} \sum_{\substack{j=1 \\ j \neq n}}^m \xi_j > \lambda] \\ &\leq P[\text{There exists } n \leq M \text{ such that } \xi_n > x/K, \bigvee_{m=0}^{\infty} \sum_{\substack{j=1 \\ j \neq n}}^m \xi_j > \lambda] \\ &\leq M \bar{F}(x/K) P[\bigvee_{j=0}^{\infty} S_j^{(\xi)} > \lambda] \end{aligned}$$

and thus we conclude

$$(2.15) \quad \limsup_{x \rightarrow \infty} \frac{I_{121}(x)}{\bar{F}(x)} \leq M K^\alpha P[\bigvee_{j=0}^{\infty} S_j^{(\xi)} > \lambda].$$

Furthermore for $I_{122}(x)$ we have

$$\begin{aligned} I_{122}(x) &\leq P\{\bigcup_{n \geq M} [\xi_n > x/K, S_i^{(\xi)} \in [0, x], i = 1, \dots, n-1]\} \\ &\leq \sum_{n=M+1}^{\infty} P[S_i^{(\xi)} \in [0, x], i = 1, \dots, n-1; \xi_n > x/K] \\ &= \bar{F}(x/K) E(N^* \wedge N(x) 1_{[N^* \wedge N(x) > M]}) \end{aligned}$$

and thus

$$(2.16) \quad \limsup_{x \rightarrow \infty} \frac{I_{122}(x)}{\bar{F}(x)} \leq K^\alpha E(\bar{N} 1_{[\bar{N} > M]}).$$

Finally we deal with $I_{11}(x)$. Define stopping times

$$\begin{aligned} N_1(x/K) &= \inf\{n \geq 0 : S_n^{(\xi)} \geq x/K\}, \\ N_i(x/K) &= \inf\{n \geq N_{i-1}(x/K) : S_n^{(\xi)} - S_{N_{i-1}(x/K)}^{(\xi)} \geq x/K\}, \quad i \geq 2. \end{aligned}$$

Because all positive steps of the random walk must be small in the piece controlled by $I_{11}(x)$, we have

$$I_{11}(x) \leq P[N_i(x/K) \leq \bar{N}, \quad i = 1, \dots, [K/2]]$$

and by the strong Markov property, this is bounded by

$$\begin{aligned} &\leq P[N_i(x/K) \leq \bar{N}, \quad i = 1, \dots, [K/2] - 1] P\left[\bigvee_{n=0}^{\infty} S_n^{(\xi)} > x/K\right] \\ &\leq \dots \leq (P\left[\bigvee_{n=0}^{\infty} S_n^{(\xi)} > x/K\right])^{[K/2]}. \end{aligned}$$

Since $P[\bigvee_{n=0}^{\infty} S_n^{(\xi)} > z]$ is regularly varying in z with index $-(\alpha - 1)$ (see for example, Bingham, Goldie and Teugels, 1987, page 387) we have

$$\frac{I_{11}(x)}{\bar{F}(x)} \leq \frac{(P[\bigvee_{n=0}^{\infty} S_n^{(\xi)} > x/K])^{[K/2]}}{\bar{F}(x)} \rightarrow 0$$

as $x \rightarrow \infty$.

We now must put the pieces together. We have

$$\begin{aligned} \limsup_{x \rightarrow \infty} \frac{P[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x] - P[\bigvee_{n=0}^{N^*} \xi_n > x]}{\bar{F}(x)} &\leq \limsup_{x \rightarrow \infty} \frac{I_{11}(x) + I_{121}(x) + I_{122}(x) + I_2(x)}{\bar{F}(x)} \\ &\leq K^\alpha M P\left[\bigvee_{n=0}^{\infty} S_n^{(\xi)} > \lambda\right] + K^\alpha E(\bar{N} 1_{[\bar{N} > M]}). \end{aligned}$$

Let $\lambda \rightarrow \infty$ and then $M \rightarrow \infty$ to get

$$\limsup_{x \rightarrow \infty} \frac{P[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x] - P[\bigvee_{n=0}^{N^*} \xi_n > x]}{\bar{F}(x)} = 0.$$

Since $P[\bigvee_{n=0}^{N^*} \xi_n > x] \sim E(\bar{N}) \bar{F}(x)$ we get

$$\limsup_{x \rightarrow \infty} \frac{P[\bigvee_{n=0}^{N^*} S_n^{(\xi)} > x]}{\bar{F}(x)} \leq E(\bar{N})$$

which is the desired reverse inequality. The proof is complete. \square

3. Contrast with exponential tails. If F_{on} has an exponentially bounded tail, known results of Iglehart (1972) can be applied to obtain the analogue of Theorem 2.3. We continue to assume that $\{X_i\}$ and $\{Y_i\}$ are iid, independent of each other, with distributions F_{on} and F_{off} respectively. Set $\xi_i = (1-r)X_i - rY_i$ and for stability continue to suppose $E\xi_1 < 0$. Recall \bar{N} is the first downgoing ladder epoch of the random walk with steps $\{\xi_i\}$. We need to suppose

$$(A1) \quad \text{For some } \gamma > 0: E^{\gamma\xi_1} = 1$$

and for this value of γ we have

$$(A2) \quad E\xi_1 e^{\gamma\xi_1} := \mu_\gamma \in (0, \infty)$$

and

$$(A3) \quad \xi_1 \text{ has a non-lattice distribution.}$$

An analogue of Corollary 2.2 is provided by Iglehart (1972) which suffices to give the tail behavior of the maximum contents level in a cycle. We write $S_n^{(\xi)} = (1-r)S_n^{(X)} - rS_n^{(Y)}$, $n \geq 0$.

Proposition 3.1 (Iglehart, 1972, Lemma 4). *Suppose assumptions A1, A2 and A3 hold. Then for $x > 0$*

$$(3.1) \quad P[M_1 > x] = P\left[\bigvee_{s=0}^{C_1} X(s) > x\right] \sim a(0)e^{-\gamma x},$$

where

$$a(0) = \frac{(1 - E(e^{\gamma S_N^{(\xi)}}))^2}{\gamma \mu_\gamma E(\bar{N})} E(e^{\gamma(1-r)X_1}).$$

We may now follow the line of reasoning of Section 2. The exponential tails given in (3.1) imply

$$(3.2) \quad \gamma \bigvee_{i=1}^{[ut]} M_i - \log a(0)u \Rightarrow Y_0(t)$$

in $D_r(0, \infty)$ where $Y_0(\cdot)$ is the extremal process generated by the Gumbel distribution

$$\Lambda(x) = \exp\{-e^{-x}\}, \quad x \in \mathbb{R}.$$

We then get as $u \rightarrow \infty$

$$\left(\gamma M(ut) - \log a(0)u, M^-\left(\frac{x + \log a(0)u}{\gamma}\right)/u \right) \Rightarrow (Y_0(t/(\mu E(\bar{N}))), \mu E(\bar{N})Y_0^-(x))$$

which leads to

$$a(0) \frac{\tau(x+u)}{e^{\gamma u}} \Rightarrow \mu E(\bar{N})Y_0^-(\gamma x).$$

If $x = 0$ we get as $u \rightarrow \infty$

$$a(0) \frac{\tau(u)}{e^{\gamma u}} \Rightarrow \mu E(\bar{N})Y_0^-(0).$$

Note that $Y_0^-(0)$ is exponentially distributed with mean 1 and we get the final result

$$(3.3) \quad \left(\frac{a(0)}{\mu E(\bar{N})} \frac{\tau(u)}{e^{\gamma u}} \right) = \left(E(e^{\gamma(1-r)X_1}) \frac{(1 - Ee^{-\gamma S_N^{(\xi)}})^2}{\mu(E\bar{N})^2 \gamma \mu_\gamma} \right) \frac{\tau(u)}{e^{\gamma u}} \Rightarrow E(1)$$

as $u \rightarrow \infty$ where $E(1)$ is a unit exponential random variable.

We may also check that $E\tau(u)/e^{\gamma u}$ converges as follows. Observe that

$$(3.4) \quad 0 \leq \tau(s) \leq S_{\bar{N}_{V(s)}}$$

where

$$V(s) = \inf\{n : \bigvee_{i=1}^n M_i \geq s\}$$

since the hitting time of $X(\cdot)$ must occur before the end of the cycle that has a cycle maximum bigger than the level. Since $V(s)$ is geometrically distributed

$$a(0) \frac{V(s)}{e^{\gamma s}} \Rightarrow E(1),$$

where, as usual, $E(1)$ is a unit exponential random variable. Furthermore, as $s \rightarrow \infty$

$$\begin{aligned} ES_{\bar{N}_{V(s)}} &= \mu E\bar{N}_{V(s)} = \mu E(\bar{N})E(V(s)) \\ &= \mu E(\bar{N}) \left(\frac{1}{P[M_1 > s]} \right) \\ &\sim \frac{\mu E(\bar{N})}{a(0)} e^{\gamma s}. \end{aligned}$$

So as $s \rightarrow \infty$

$$\frac{E(S_{\bar{N}_{V(s)}})}{e^{\gamma s}} \rightarrow \frac{\mu E(\bar{N})}{a(0)},$$

and

$$\frac{S_{\bar{N}_{V(s)}}}{e^{\gamma s}} \sim \mu E(\bar{N}) \frac{V(s)}{e^{\gamma s}} \Rightarrow \frac{\mu E(\bar{N})}{a(0)} E(1).$$

These two statements coupled with (3.3) and (3.4) and a variant of Fatou's lemma sometimes called Pratt's lemma (Pratt, 1960) yield the desired result

$$(3.5) \quad E\tau(s) \sim \frac{\mu E(\bar{N})}{a(0)} e^{\gamma s} \quad (s \rightarrow \infty).$$

Note how critically F_{off} enters into formulas (3.3), (3.5) since F_{off} is important in determining the growth rate γ of the hitting time. In the heavy tail case, F_{off} did not play a role in determining the growth rate of $\tau(s)$ since levels were hit basically due just to big upward jumps which were controlled by F_{on} . Recall that in the heavy tailed case, as $s \rightarrow \infty$,

$$E\tau(s) \sim \mu(1-r)^{-\alpha} / \bar{F}_{\text{on}}(s) \sim ES_{\nu} = \mu E\nu$$

where

$$\nu = \inf\{n : (1-r)X_n > s\}.$$

Example: Consider the standard example where both F_{on} and F_{off} are exponential distributions with means μ_{on} and μ_{off} respectively. The negative drift condition is

$$(1-r)\mu_{\text{on}} - r\mu_{\text{off}} < 0$$

and γ must satisfy

$$1 = Ee^{\gamma((1-r)X_1 - rY_1)} = ((1 - \gamma(1-r)\mu_{\text{on}})(1 + \gamma r\mu_{\text{off}}))^{-1}.$$

Solving for γ we get the solutions $\gamma = 0$ and

$$\gamma = \frac{r\mu_{\text{off}} - (1-r)\mu_{\text{on}}}{(1-r)r\mu_{\text{on}}\mu_{\text{off}}} = \frac{-E\xi_1}{(1-r)r\mu_{\text{on}}\mu_{\text{off}}}.$$

The numerator is positive by the drift condition.

We now calculate the coefficient of $\tau(u)/e^{\gamma u}$ in (3.3) in order to compare it to the exact calculation given in (1.5). Since X_1 is assumed exponential with mean μ_{on} we have

$$Ee^{\gamma(1-r)X_1} = \frac{r\mu_{\text{off}}}{(1-r)\mu_{\text{on}}}.$$

To calculate μ_γ we compute $\frac{d}{ds} Ee^{s\xi_1}$ and substitute $s = \gamma$. This calculation is made easier by use of the formulas

$$\begin{aligned} 1 + r\gamma\mu_{\text{off}} &= \frac{r\mu_{\text{off}}}{(1-r)\mu_{\text{on}}} \\ 1 - \gamma(1-r)\mu_{\text{on}} &= \frac{(1-r)\mu_{\text{on}}}{r\mu_{\text{off}}} = \frac{1}{1 + r\gamma\mu_{\text{off}}}. \end{aligned}$$

With these formulas we find

$$\mu_\gamma = -E\xi_1.$$

Next observe that because of exponential tails,

$$S_N^{(\xi)} \stackrel{d}{=} -r\mu_{\text{off}}E(1),$$

where $E(1)$ is a unit exponential random variable and hence

$$\begin{aligned} 1 - Ee^{\gamma S_N^{(\xi)}} &= 1 - \frac{1}{1 + \gamma r\mu_{\text{off}}} \\ &= 1 - \frac{(1-r)\mu_{\text{on}}}{r\mu_{\text{off}}} \\ &= \frac{-E\xi_1}{r\mu_{\text{off}}}. \end{aligned}$$

Knowing the distribution of $S_N^{(\xi)}$ also enables us to compute $E\bar{N}$ since

$$ES_N^{(\xi)} = -r\mu_{\text{off}} = E(\bar{N})E\xi_1$$

and hence

$$E\bar{N} = \frac{r\mu_{\text{off}}}{-E\xi_1}.$$

Putting the ingredients together, (3.2) becomes

$$\frac{(-E\xi_1)^2}{\mu(r\mu_{\text{off}})^2} \frac{\tau(u)}{e^{\gamma u}} \Rightarrow E(1)$$

which agrees with the exact calculation for the expected value in (1.5).

The contrast with the heavy tailed case is very evident. Instead of $E\tau(s)$ being of the same order as ES_ν , the expected time until an *on* period of length at least $s/(1-r)$ occurs, we have in our exponential example

$$\begin{aligned} ES_\nu &= \mu E\nu = \mu/\bar{F}_{\text{on}}(s/(1-r)) \\ &= \mu \exp\{s/(\mu_{\text{on}}(1-r))\}. \end{aligned}$$

However from (3.5)

$$E\tau(s) \sim \frac{\mu(r\mu_{\text{off}})^2}{(-E\xi_1)^2} e^{\gamma s}.$$

Comparing the growth rates γ with $1/(\mu_{\text{on}}(1-r))$ we have

$$0 < \gamma = \frac{1}{(1-r)\mu_{\text{on}}} - \frac{1}{r\mu_{\text{off}}} < \frac{1}{(1-r)\mu_{\text{on}}}$$

so that ES_ν has a faster growth rate which is to be expected since in the exponential case, the process $X(\cdot)$ jumps over a high level as a result of an accumulation of small upward movements and not typically as a result of a single large jump.

To obtain the exact expression for $E\tau(s)$ in this example, proceed as follows. Defining

$$\tilde{X}(t) = (X(t), Z(t)), t \geq 0,$$

we describe the state of the system prior to reaching level s as a Markov process $\{\tilde{X}(t), t \geq 0\}$ with a state space $\mathbb{E} = \{(x, i), 0 \leq x \leq s, i = 0, 1\}$. We can express $\tau(s)$ in terms of the hitting times of the Markov process $\{\tilde{X}(t), t \geq 0\}$ as

$$\tau(s) = T_{(s,1)} := \inf\{t \geq 0 : \tilde{X}(t) = (s, 1)\}.$$

For $\mathbf{x} \in \mathbb{E}$, let $H(\mathbf{x})$ be the expected hitting time $T_{(s,1)}$ starting at \mathbf{x} , and define for an $0 \leq x \leq s$,

$$h_1(x) = H((x, 1)), \quad h_2(x) = H((x, 0)).$$

Then $E\tau(s) = h_1(0)$. Using the natural filtration $\mathcal{F}_t = \sigma(Z(u), 0 \leq u \leq t), t \geq 0$, we observe that for any $t \geq 0$

$$E(T_{(s,1)}|\mathcal{F}_t) = H(\tilde{X}(t \wedge T_{(s,1)})) + \tilde{X}(t \wedge T_{(s,1)}) := M(t).$$

Therefore, $\{M(t), t \geq 0\}$ is a martingale, and its martingale property leads to the following system of ordinary differential equations:

$$(3.6) \quad (1-r)h_1'(x) = -1 + \frac{1}{\mu_{\text{on}}}h_1(x) - \frac{1}{\mu_{\text{on}}}h_2(x),$$

$$(3.7) \quad rh_2'(x) = 1 + \frac{1}{\mu_{\text{off}}}h_1(x) - \frac{1}{\mu_{\text{off}}}h_2(x),$$

with the obvious boundary conditions

$$(3.8) \quad h_2(0) = \mu_{\text{off}} + h_1(0), \quad h_1(s) = 0.$$

The system (3.6–3.8) can be solved in the standard way, and we obtain (1.5).

4. A single server fluid queue fed by several *on/off* processes.

Let $\{X_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$ and $\{Y_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$ be iid copies of the *on* sequence $\{X_i, i = 1, 2, \dots\}$ and the *off* sequence $\{Y_i, i = 1, 2, \dots\}$ correspondingly. We construct k iid *on/off* processes, $\{Z_j(t), t \geq 0\}$, $j = 1, \dots, k$ as in (1.1). Sometimes we will find it convenient to work with stationary versions of $\{Z_j(t), t \geq 0\}$, $j = 1, \dots, k$. Those exist due to the finiteness of μ_{on} and μ_{off} , and can be constructed as follows. Fix a $j = 1, \dots, k$, and let $C_j^{(0)}, X_j^{(0)}, Y_j^{(0)}$ and $Y_0^{(j)}$ be four independent random variables which are independent of $\{X_i^{(j)}, Y_i^{(j)}, i \geq 1\}$ defined as follows: $C_j^{(0)}$ is a Bernoulli random variable with values $\{0, 1\}$ and mass function

$$P[C_j^{(0)} = 1] = \frac{\mu_{\text{on}}}{\mu} = 1 - P[C_j^{(0)} = 0]$$

and ($x > 0$)

$$P[X_j^{(0)} > x] = \int_x^\infty \frac{\bar{F}_{\text{on}}(s)}{\mu_{\text{on}}} ds, \quad P[Y_j^{(0)} > x] = \int_x^\infty \frac{\bar{F}_{\text{off}}(s)}{\mu_{\text{off}}} ds.$$

Finally, $Y_0^{(j)}$ has the F_{off} distribution. Define a delay random variable $D_j^{(0)}$ by

$$D_j^{(0)} = C_j^{(0)}(X_j^{(0)} + Y_0^{(j)}) + (1 - C_j^{(0)})Y_j^{(0)}$$

and a delayed renewal sequence by

$$(4.1) \quad \{S_n^{(j)}, n \geq 0\} := \{D_j^{(0)}, D_j^{(0)} + \sum_{i=1}^n (X_i^{(j)} + Y_i^{(j)}), n \geq 1\}.$$

Then a stationary version of $\{Z_j(t), t \geq 0\}$ is defined by

$$(4.2) \quad Z_j(t) = C_j^{(0)}1_{[0, X_j^{(0)}]}(t) + \sum_{n=0}^\infty 1_{[S_n^{(j)} \leq t < S_n^{(j)} + X_{n+1}^{(j)}]}.$$

See Heath, Resnick and Samorodnitsky (1996) for details. In a similar way we can construct a stationary version $\{Z_j(t), -\infty < t < \infty\}$ defined for all real t . We take, further, the k stationary *on/off* processes to be independent.

In this section we consider a single server queue as above, with service rate r , fed by the k *on/off* processes. That is, the combined inflow rate in the system is given by

$$(4.3) \quad Z^{(k)}(t) = Z_1(t) + \dots + Z_k(t), \quad t \geq 0, \text{ or } -\infty < t < \infty,$$

and, similarly to (1.2), the state $\{X^{(k)}(t), t \geq 0\}$ of the system satisfies

$$(4.4) \quad dX^{(k)}(t) = Z^{(k)}(t) dt - r(X^{(k)}(t)) dt.$$

It is of interest to consider the behavior of a system (4.4) with a general k , first of all as a step towards understanding the queues with more general long memory input streams and, secondly, to understand the effect of pooling the resources in the systems of the type we are considering. The natural rate condition for this system, parallel to (1.3), is

$$(4.5) \quad k \frac{\mu_{\text{on}}}{\mu} < r,$$

saying that the long-term inflow rate to the system (4.4) is less than the potential long-term outflow rate from the system. Of course, we also assume that $r < k$, to make sure that the system is non-degenerate. In the appendix we verify that, under condition (4.5), there is a stationary stochastic process $\{X^{(k)}(t), t \geq 0\}$ satisfying (4.4). This statement, though intuitively obvious, does require an argument due to the lack of renewal structure in the process $\{Z^{(k)}(t), t \geq 0\}$. We assume, as always, that the distribution $F_{\text{on}} * F_{\text{off}}$ is not arithmetic.

When the *on* periods have a heavy tailed distribution, we know from the discussion in Section 2 that, for $k = 1$, the state of a system driven by (4.4) crosses a high level L by increasing to that level almost from 0 within a single *on* period. We expect high level crossing patterns of the system contents to be similar for a general k . Intuitively, the time to reach a high level L should critically depend on

$$(4.6) \quad k_0 = \text{the smallest integer} > r.$$

By the non-triviality assumption, $k_0 \leq k$. If $k_0 = 1$, by analogy to the case $k = 1$ one expects the state of the system to reach a high level L when one of the k *on/off* processes has an *on* period of length $L/(1-r)$. If $k_0 > 1$, the same intuition says that the system will reach a high level L only when k_0 very long *on* periods occur at about the same time, and so it will take much longer until this high level is reached. In this section we prove the above statement for the case $k_0 = 1$, thus generalizing the conclusion reached in Section 2 for $k = 1$. The proof of the Theorem 4.1 is significantly more involved than the argument required for $k = 1$ due, once again, to the lack of renewal structure in the process $\{Z^{(k)}(t), t \geq 0\}$.

A natural way of calculating the time until the system contents reach the level L is by starting from the moment the system is empty, and all k *on/off* processes begin an *on* period. One must realize, however, that for $k > 1$ such a moment in time is far from being “typical”, and, even if we initialize the system in such a way, chances are that such moments will not recur. Therefore, we state our theorem in a more general way, by allowing more general initial conditions. To this end, let H be an arbitrary probability law on R_+^k , whose marginals have finite first moments. Let $(D_1^{(0)} \dots, D_k^{(0)})$ be an H -distributed random vector, independent of the sequences $\{X_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$ and $\{Y_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$. We again define a delayed renewal sequence by (4.1), and, similarly to (4.2), we define $\{Z_j(t), t \geq 0\}$ by

$$(4.7) \quad Z_j(t) = \sum_{n=0}^{\infty} 1_{[S_n^{(j)} \leq t < S_n^{(j)} + X_{n+1}^{(j)}]}.$$

Clearly, this time $\{Z_j(t), t \geq 0\}$ does not have to be stationary. If $\{X^{(k)}(t), t \geq 0\}$ is given now by $X^{(k)}(0) = x_0 \in \{0, 1, 2, \dots\}$ and (4.4), we will denote all probabilities and expectations related to it as P_{H, x_0} and E_{H, x_0} , accordingly. That is, we are allowing the system to start in an arbitrary state x_0 , when all the *on/off* processes are in *off* periods, with H describing the joint distribution of the remainders of the initial *off* periods. Let

$$(4.8) \quad \tau(L) = \inf\{t \geq 0 : X^{(k)}(t) \geq L\}.$$

Theorem 4.1. *Let*

$$\bar{F}_{\text{on}}(x) = x^{-\alpha} L(x), \quad \alpha > 1, \quad x \rightarrow \infty,$$

and assume that for some $p > 1$ we have $EY_1^p < \infty$. If the service rate r satisfies (4.5), and $r < 1$, then for any H and x_0 we have

$$(4.9) \quad \lim_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E_{H, x_0} \tau(L) = \lim_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E_{H, x_0} \tau^*(L) = \frac{1}{k} \mu,$$

where

$$\tau^*(L) = \inf\{t \geq 0 : \text{one of the } k \text{ on/off processes begins at time } t \\ \text{an on period of length at least } \frac{L}{1-r}\}.$$

An argument identical to that of Theorem 2.1 immediately proves the corresponding result for the corresponding $GI/G/1$ queue. Let $\{Y_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$ be iid copies of the sequence of interarrival times $\{Y_i, i = 1, 2, \dots\}$, and let $\{B_i^{(j)}, i = 1, 2, \dots\}$, $j = 1, \dots, k$ be iid copies of the sequence of offered work $\{B_i, i = 1, 2, \dots\}$, so that at time $S_n^{(j)} = Y_1^{(j)} + \dots + Y_n^{(j)}$ an amount of work $B_n^{(j)}$ is brought in the system on the j th channel, $n = 0, 1, \dots$; $j = 1, \dots, k$. Let r be the service rate. Note that the following theorem does not require the assumption $r < 1$.

Theorem 2.2. *Let the distribution F_{on} of B_i satisfy*

$$\bar{F}_{\text{on}}(x) = x^{-\alpha} L(x), \quad \alpha > 1, \quad x \rightarrow \infty,$$

and assume that for some $p > 1$ we have $EY_1^p < \infty$. If

$$(4.11) \quad k \frac{\mu_{\text{on}}}{\mu_{\text{off}}} < r,$$

then

$$(4.12) \quad \lim_{L \rightarrow \infty} \bar{F}_{\text{on}}(L) E\tau(L) = \frac{1}{k} \mu_{\text{off}},$$

where $\tau(L)$ is the first time the amount of work in the system reaches the level L .

In particular, the expected time to reach a high level L in (2.12) *does not depend on the service rate r* . Note that the result of Theorem 2.2 remains true if one initializes in an arbitrary way the state of the system.

Let us look at the implications of our results on the benefits of pooling the system resources. Think of k iid $GI/G/1$ queues, with holding capacity L each, and service rates r each. The queue number j is driven by the sequences $\{Y_i^{(j)}, i = 1, 2, \dots\}$ and $\{B_i^{(j)}, i = 1, 2, \dots\}$ as above, $j = 1, \dots, k$. We take pooling system resources here to mean that we put together the service resources to create a “super-server” with service rate kr , and we feed this “super-server” by a combined stream of the k input processes, as in Theorem 2.2. The holding capacity of the new system is taken to be kL , again, as the result of pooling the resources. Let us look at a generic stream of “customers”, or work (i.e. one of the k original streams of work). One can imagine that when the holding capacity of the system serving these “customers” is reached, the system is blocked for a time to any future arrivals. Under the “ k separate servers” scenario, the expected time until the serving system is blocked is $E\tau(L)$, while when the system resources are pooled, this expected time is $E\tau_{kL}$. By Theorem 2.2 the ratio R of the two expected times is, asymptotically for large L ,

$$R = \lim_{L \rightarrow \infty} \frac{k \bar{F}_{\text{on}}(kL)}{\bar{F}_{\text{on}}(L)} = \frac{1}{k^{\alpha-1}} < 1,$$

which is the expected benefit of pooling the resources. However, this benefit becomes less pronounced if α is close to 1.

We are ready now for the proofs.

Proof of Theorem 2.1. Obviously,

$$(4.13) \quad E_{H, x_0} \tau(L) \leq E_{H, x_0} \tau^*(L) + \frac{L}{1-r}.$$

If

$$(4.14) \quad \tau^{*,j}(L) = \inf \left\{ S_n^{(j)} : n \geq 1, X_{n+1}^{(j)} \geq \frac{L}{1-r} \right\}, \quad j = 1, \dots, k,$$

we have

$$(4.15) \quad \tau^*(L) = \min \left(\tau^{*,1}(L), \dots, \tau^{*,k}(L) \right).$$

Observe that $\tau^{*,1}(L), \dots, \tau^{*,k}(L)$ are, conditionally on $(D_1^{(0)}, \dots, D_k^{(0)})$, independent, and that

$$(4.16) \quad \tau^{*,j}(L) \stackrel{d}{=} D_j^{(0)} + \sum_{i=1}^{G_L^{(j)}} (X_i^{(*,j)} + Y_i^{(j)}),$$

where $G_L^{(j)}$ is a geometric random variables with parameter $\bar{F}_{\text{on}}(\frac{L}{1-r})$, independent of two independent iid sequences, $\{X_i^{(*,j)}, i = 1, 2, \dots\}$ and $\{Y_i^{(j)}, i = 1, 2, \dots\}$, where the latter sequence has, as usual, the F_{off} distribution, and the former has the distribution

$$P(X_i^{(*,j)} \in A) = P\left(X_1^{(j)} \in A | X_1^{(j)} \leq L\right) = \frac{1}{F_{\text{on}}(\frac{L}{1-r})} \int_0^{L/(1-r)} 1(x \in A) F_{\text{on}}(dx).$$

Everything is also independent of the delay random variable $D_0^{(j)}$. In particular, $(X_1^{(j)} | X_1^{(j)} \leq L) \stackrel{\text{st}}{\leq} X_1^{(j)}$, and hence

$$(4.17) \quad \tau^{*,j}(L) \stackrel{\text{st}}{\leq} S_{G_L^{(j)}}^{(j)},$$

where all the random variables appearing in the right hand side of (4.17) are independent. We conclude that

$$\tau^*(L) \stackrel{\text{st}}{\leq} \bigvee_{1 \leq j \leq k} D_0^{(j)} + \sum_{i=1}^{\wedge_{1 \leq j \leq k} G_L^{(j)}} (X_i^{(1)} + Y_i^{(1)}),$$

and so

$$E_{H,x_0} \tau^*(L) \leq E_H \bigvee_{1 \leq j \leq k} D_0^{(j)} + \mu E \bigwedge_{1 \leq j \leq k} G_L^{(j)}.$$

Since $\wedge_{1 \leq j \leq k} G_L^{(j)}$ is, once again, a geometric random variable with parameter $1 - F_{\text{on}}(\frac{L}{1-r})^k$, we obtain immediately that

$$E_{H,x_0} \tau(L) \leq E_{H,x_0} \tau^*(L) + \frac{L}{1-r} \leq E_H \bigvee_{1 \leq j \leq k} D_0^{(j)} + \mu \left(\frac{1}{1 - F_{\text{on}}(\frac{L}{1-r})^k} - 1 \right) + \frac{L}{1-r},$$

which implies that

$$(4.18) \quad \limsup_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E_{H,x_0} \tau(L) \leq \limsup_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E_{H,x_0} \tau^*(L) \leq \frac{1}{k} \mu.$$

For a lower bound, we start with observing that we may assume that $\wedge_{1 \leq j \leq k} D_j^{(0)} = 0$, P_{H, x_0} -almost surely. Indeed, shortening all the initial *off* periods by the same amount can only bring the level crossing time closer. Now, take an $0 < \epsilon < 1$, and observe that

$$\tau(L) \geq \tau^*(L(1-\epsilon))1_{[\tau(L) \geq \tau^*(L(1-\epsilon))]},$$

and so

$$(4.19) \quad E_{H, x_0} \tau(L) \geq E_{H, x_0} \tau^*(L(1-\epsilon)) - E_{H, x_0} \left(\tau^*(L(1-\epsilon))1_{[\tau(L) < \tau^*(L(1-\epsilon))]} \right).$$

The main part of the proof of a lower bound on the expected crossing time is a proof of the fact that

$$(4.20) \quad \lim_{L \rightarrow \infty} \bar{F}_{\text{on}} \left(\frac{L}{1-r} \right) E_{H, x_0} \left(\tau^*(L(1-\epsilon))1_{[\tau(L) < \tau^*(L(1-\epsilon))]} \right) = 0.$$

Suppose that (4.20) has been proved. If we can establish that

$$(4.21) \quad \liminf_{L \rightarrow \infty} \bar{F}_{\text{on}} \left(\frac{L}{1-r} \right) E_{H, x_0} \tau^*(L) \geq \frac{1}{k} \mu,$$

then (4.19), (4.20) and the regular variation of $\bar{F}_{\text{on}}(L)$ will provide the required counterpart to (4.18), and so prove the theorem. To prove (4.21) we may, of course, assume that $H = \delta_{(0, \dots, 0)}$, and $x_0 = 0$. We, therefore, use P and E without any subscripts. We remark at this point this assumption is made only for the purpose of proving of (2.21), and will be removed once the latter has been proved. However, at certain later stages of the proof of the theorem we will find it useful (and possible) to re-impose this assumption once again.

With $G_L^{(j)}$, $j = 1, \dots, k$ as before, let

$$(4.22) \quad \hat{\tau}^*(L) = \min \left(S_{G_L^{(1)}}^{(1)}, \dots, S_{G_L^{(k)}}^{(k)} \right).$$

Let us check first that

$$(4.23) \quad \lim_{L \rightarrow \infty} \bar{F}_{\text{on}} \left(\frac{L}{1-r} \right) \left(E \hat{\tau}^*(L) - E \tau^*(L) \right) = 0.$$

We may assume that the random variables $\{X_i^{(j)}, i \geq 1\}$, $\{X_i^{(*,j)}, i \geq 1\}$, $\{Y_i^{(j)}, i \geq 1\}$ and $G_L^{(j)}$, $j = 1, \dots, k$ are all defined on the same probability space, such that $X_i^{(j)} \geq X_i^{(*,j)}$ for all $i \geq 1$ and $j = 1, \dots, k$, and that the same $\{Y_i^{(j)}, i \geq 1\}$ and $G_L^{(j)}$ are used to define $S_{G_L^{(j)}}^{(j)}$ and $\tau^{*,j}(L)$ for $j = 1, \dots, k$.

We then have:

$$\hat{\tau}^*(L) - \tau^*(L) = \min \left(S_{G_L^{(1)}}^{(1)}, \dots, S_{G_L^{(k)}}^{(k)} \right) - \min \left(\tau^{*,1}(L), \dots, \tau^{*,k}(L) \right) \leq \bigvee_{1 \leq j \leq k} \left(S_{G_L^{(j)}}^{(j)} - \tau^{*,j}(L) \right).$$

Since $S_{G_L^{(j)}}^{(j)} \geq \tau^{*,j}(L)$ for all $j = 1, \dots, k$, we have, therefore,

$$\begin{aligned} E \hat{\tau}^*(L) - E \tau^*(L) &\leq k E \left(S_{G_L^{(1)}}^{(1)} - \tau^{*,1}(L) \right) \\ &= k E G_L^{(1)} (E X_1^{(1)} - E X_1^{(*,1)}) \\ &= k \left(\bar{F}_{\text{on}} \left(\frac{L}{1-r} \right)^{-1} - 1 \right) o(1) \\ &= o \left(\bar{F}_{\text{on}} \left(\frac{L}{1-r} \right)^{-1} \right) \end{aligned}$$

as $L \rightarrow \infty$. This proves (4.23). Therefore, (4.21) will follow once we check that

$$(4.24) \quad \liminf_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E\hat{\tau}^*(L) \geq \frac{1}{k}\mu.$$

Now,

$$E\hat{\tau}^*(L) = E \min\left(S_{G_L^{(1)}}^{(1)}, \dots, S_{G_L^{(k)}}^{(k)}\right) = \int_0^\infty P\left(S_{G_L^{(1)}}^{(1)} > x\right)^k dx.$$

Choose an $N \geq 1$ (large) and $0 < \delta < 1$ (small) and observe that for any $x > 0$,

$$(4.25) \quad \begin{aligned} P\left(S_{G_L^{(1)}}^{(1)} > x\right) &= P\left(\sum_{i=1}^{G_L^{(1)}} (X_i^{(1)} + Y_i^{(1)}) > x\right) \\ &\geq P\left(\cap_{n>N} [G_L^{(1)} > \max(N, \frac{x}{(1-\delta)\mu}), \sum_{i=1}^n (X_i^{(1)} + Y_i^{(1)}) > n(1-\delta)\mu]\right) \\ &= P\left(G_L^{(1)} > \max(N, \frac{x}{(1-\delta)\mu})\right) P\left(\cap_{n>N} [\sum_{i=1}^n (X_i^{(1)} + Y_i^{(1)}) > n(1-\delta)\mu]\right) \\ &=: \theta_N P\left(G_L^{(1)} > \max(N, \frac{x}{(1-\delta)\mu})\right). \end{aligned}$$

Observe that by the strong law of large numbers, $\theta_N \rightarrow 1$ as $N \rightarrow \infty$. We have, therefore,

$$\begin{aligned} E\hat{\tau}^*(L) &\geq \theta_N^k \int_{N(1-\delta)\mu}^\infty P\left(G_L^{(1)} > \frac{x}{(1-\delta)\mu}\right)^k dx \\ &= \theta_N^k \sum_{n=N}^\infty (1-\delta)\mu P\left(G_L^{(1)} > n\right)^k \\ &= \theta_N^k (1-\delta)\mu \frac{F_{\text{on}}\left(\frac{L}{1-r}\right)^{k(N+1)}}{1 - F_{\text{on}}\left(\frac{L}{1-r}\right)^k}. \end{aligned}$$

We conclude that

$$\liminf_{L \rightarrow \infty} \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right) E\hat{\tau}^*(L) \geq \frac{1}{k}\mu\theta_N^k(1-\delta),$$

and (4.24) follows by letting $N \rightarrow \infty$ and $\delta \rightarrow 0$. Therefore, to prove the theorem we only need to prove (4.20).

By the assumptions of the theorem, there is a $p' > 1$ such that both $EX_1^{p'} < \infty$ and $EY_1^{p'} < \infty$. By Hölder's inequality,

$$(4.26) \quad E_{H,x_0}\left(\tau^*(L(1-\epsilon))1(\tau(L) < \tau^*(L(1-\epsilon)))\right) \leq \left(E_{H,x_0}(\tau^*(L(1-\epsilon)))^{p'}\right)^{1/p'} P_{H,x_0}(\tau(L) < \tau^*(L(1-\epsilon)))^{1/q'},$$

where $\frac{1}{p'} + \frac{1}{q'} = 1$. We use the following inequality: for any σ -finite measure spaces $(\Omega_1, \mathcal{F}_1, \mu_1)$ and $(\Omega_2, \mathcal{F}_2, \mu_2)$, a $p' > 1$ and a nonnegative measurable function $f : \Omega_1 \times \Omega_2 \rightarrow R$,

$$(4.27) \quad \left(\int_{\Omega_1} \left(\int_{\Omega_2} f(\omega_1, \omega_2) \mu_2(d\omega_2)\right)^{p'} \mu_1(d\omega_1)\right)^{1/p'} \leq \int_{\Omega_2} \left(\int_{\Omega_1} f(\omega_1, \omega_2)^{p'} \mu_1(d\omega_1)\right)^{1/p'} \mu_2(d\omega_2).$$

See for example Lemma 3.3.1 of Kwapień and Woyczyński (1992). Let $j_0 = \operatorname{argmin}(D_1^{(0)}, \dots, D_k^{(0)})$ (with ties broken in, say, the lexicographical manner), and recall that $D_{j_0}^{(0)} = 0$, P_{H, x_0} -almost surely. We have by (4.27)

$$\begin{aligned}
\left(E_{H, x_0}(\tau^*(L(1-\epsilon)))^{p'}\right)^{1/p'} &\leq \left(E_{H, x_0}(\tau^*(L))^{p'}\right)^{1/p'} \leq \left(E_{H, x_0}(\hat{\tau}_L^*)^{p'}\right)^{1/p'} \\
&\leq \left(E_{H, x_0}(S_{G_L^{(j_0)}}^{(j_0)})^{p'}\right)^{1/p'} \\
&= \left(E\left(\sum_{i=1}^{\infty} (X_i^{(1)} + Y_i^{(1)})1(i \leq G_L^{(1)})\right)^{p'}\right)^{1/p'} \\
&\leq \sum_{i=1}^{\infty} \left(E\left((X_i^{(1)} + Y_i^{(1)})1(i \leq G_L^{(1)})\right)^{p'}\right)^{1/p'} \\
&= \left(E(X_1^{(1)} + Y_1^{(1)})^{p'}\right)^{1/p'} \sum_{i=1}^{\infty} P(G_L^{(1)} \geq i)^{1/p'} \\
&\leq C \bar{F}_{\text{on}}\left(\frac{L}{1-r}\right)^{-1},
\end{aligned}$$

where C is a finite positive constant. It follows from (4.26) that we will prove (4.20) by showing that

$$(4.28) \quad \lim_{L \rightarrow \infty} P_{H, x_0}(\tau(L) < \tau^*(L(1-\epsilon))) = 0.$$

Let us prove first that for every

$$(4.29) \quad N > 2k \frac{\alpha}{\alpha - 1},$$

we have

$$(4.30) \quad \lim_{L \rightarrow \infty} P_{H, x_0}(\tau(L) < \tau^*(L/N)) = 0.$$

To this end, let us “unpool” the system. That is, imagine k separate fluid queuing systems defined by

$$(4.31) \quad dX_j(t) = Z_j(t) dt - \frac{1}{k} r(X_j)(t) dt, \quad t \geq 0$$

where $\{Z_j(t), t \geq 0\}$ is given by (4.7), and $X_j(0) = \frac{1}{k} x_0$, $j = 1, \dots, k$. The k processes $\{X_j(t), t \geq 0\}$, $j = 1, \dots, k$ are, conditionally on the initial delay $(D_1^{(0)}, \dots, D_k^{(0)})$, independent. Let $Y^{(k)}(t) = X_1(t) + \dots + X_k(t)$, $t \geq 0$. The two processes, $\{X^{(k)}(t), t \geq 0\}$ and $\{Y^{(k)}(t), t \geq 0\}$ describe the states of two queuing systems. Obviously, $X^{(k)}(0) = Y^{(k)}(0)$, the two systems have identical inflow streams of work, while the outflow of work from $X^{(k)}(\cdot)$ when the system is not empty is always at rate r , and the rate of outflow of work from $Y^{(k)}(\cdot)$ does not exceed r . Therefore, for every ω ,

$$(4.32) \quad X^{(k)}(t) \leq Y^{(k)}(t), \quad t \geq 0.$$

Note that (4.32) is just an expression of the benefit of pooling the system resources. Define

$$(4.33) \quad \tau^{(Y)}(L) = \inf\{t \geq 0 : Y^{(k)}(t) \geq L\}.$$

Then (4.32) implies that $\tau^{(Y)}(L) \leq \tau(L)$, so that

$$\left\{ \omega : \tau(L) < \tau^*(L/N) \right\} \subseteq \left\{ \omega : \tau^{(Y)}(L) < \tau^*(L/N) \right\},$$

and, therefore,

$$(4.34) \quad P_{H,x_0}(\tau(L) < \tau^*(L/N)) \leq P_{H,x_0}(\tau^{(Y)}(L) < \tau^*(L/N)).$$

Now, let for $j = 1, \dots, k$

$$\tau^{(j)}(L) = \inf \{ t \geq 0 : X_j(t) \geq L \}.$$

Then

$$(4.35) \quad \begin{aligned} P_{H,x_0}(\tau^{(Y)}(L) < \tau^*(L/N)) &\leq P_{H,x_0} \left(\tau^{(j)}(L/k) < \tau^*(L/N) \text{ for some } j = 1, \dots, k \right) \\ &\leq \sum_{i=1}^j P_{H,x_0}(\tau^{(j)}(L/k) < \tau^*(L/N)) \\ &\leq \sum_{i=1}^j P_{H,x_0}(\tau^{(j)}(L/k) < \tau^{*,j}(L/N)). \end{aligned}$$

Therefore, (4.30) will follow from (4.34) and (4.35) once we prove that for every

$$(4.36) \quad M > 2 \frac{\alpha}{\alpha - 1}$$

we have

$$(4.37) \quad \lim_{L \rightarrow \infty} P_{H,x_0}(\tau^{(j)}(L) < \tau^{*,j}(L/M)) = 0,$$

for $j = 1, \dots, k$.

We will prove (4.37) for $j = 1$. Let $T_0 = \inf \{ t > D_1^{(0)} : X_1(t) = 0 \}$. Write

$$P_{H,x_0}(\tau^{(1)}(L) < \tau^{*,1}(L/M)) = P_{H,x_0}(T_0 \leq \tau^{(1)}(L) < \tau^{*,1}(L/M)) + P_{H,x_0}(\tau^{(1)}(L) < \tau^{*,1}(L/M) \wedge T_0).$$

Since for all $L > 2x_0$

$$\begin{aligned} P_{H,x_0}(\tau^{(1)}(L) < \tau^{*,1}(L/M) \wedge T_0) &\leq P_{H,x_0}(\tau^{(1)}(L) < T_0) \\ &\leq P(\tau_{L/2}^{(1)} < T_0) \rightarrow 0 \end{aligned}$$

as $L \rightarrow \infty$ because of the rate condition (4.5) (or recall Corollary 2.2), (4.37) will follow if we prove that

$$(4.38) \quad \lim_{L \rightarrow \infty} P_{H,x_0}(T_0 \leq \tau^{(1)}(L) < \tau^{*,1}(L/M)) = 0.$$

Clearly, at time T_0 the system is in an *off* period. Denote by \tilde{H} the law (under P_{H,x_0}) of the remainder of this *off* period after time T_0 . By the strong Markov property,

$$P_{H,x_0}(T_0 \leq \tau^{(1)}(L) < \tau^{*,1}(L/M)) \leq P_{\tilde{H},0}(\tau^{(1)}(L) < \tau^{*,1}(L/M)).$$

Therefore, (4.38) will follow once we prove (4.37) with $x_0 = 0$, and so (4.37) in its generality will follow if we prove it only for $x_0 = 0$. Assume, therefore, that $x_0 = 0$.

For $K_1 \geq 0$ and $K_2 > 0$ let $\{X_1^{(K_1, K_2)}, t \geq 0\}$ denote the process given by (4.31) (with $j = 1$) when $D_1^{(0)}$ is replaced by $D_1^{(0)} \wedge K_1$, and each $Y_i^{(1)}$ is replaced by $Y_i^{(1)} \wedge K_2$. Let $\tau^{(1)}(L; K_1, K_2)$ and $\tau^{*,1}(L; K_1, K_2)$ be the random times analogous to $\tau^{(1)}(L)$ and $\tau^{*,1}(L)$ correspondingly, defined with respect to the process $\{X_1^{(K_1, K_2)}, t \geq 0\}$. Observe that the event

$$A_{K_1, K_2} = \{\omega : \tau^{(1)}(L; K_1, K_2) < \tau_{L/M}^{*,1}(K_1, K_2)\}$$

increases when K_1 and K_2 decrease. It is enough, therefore, to prove (4.37) in the case when $K_1 = 0$, and K_2 is any finite positive number such that

$$\frac{\mu_{\text{on}}}{\mu_{\text{on}} + E(Y_1^{(1)} \wedge K_2)} < \frac{r}{k}.$$

In other words, we will prove (4.37) in the case when $D_1^{(0)} = 0$, the *off* times are bounded (by K_2), and $x_0 = 0$. We will, therefore, use once again P and E without any subscripts.

We define 3 events, $A_i(L)$, $i = 1, 2, 3$, corresponding to the following 3 possibilities.

- (i) $\tau^{*,1}(L/M) < \tau^{(1)}(L) \wedge T_0$. That is, the process $\{X_1(t), t \geq 0\}$ begins an *on* interval of length at least $\frac{L}{M(1-r)}$ before reaching either level L , or returning to 0.
- (ii) $\tau^{(1)}(L) < \tau^{*,1}(L/M) \wedge T_0$. In other words, the process $\{X_1(t), t \geq 0\}$ reaches level L before starting an *on* interval of length at least $\frac{L}{M(1-r)}$ and before returning to 0.
- (iii) $T_0 < \tau^{(1)}(L) \wedge \tau^{*,1}(L/M)$. In other words, the process $\{X_1(t), t \geq 0\}$ returns to 0 before reaching level L , and without having an *on* interval of the length of at least $\frac{L}{M(1-r)}$.

Clearly,

$$P(\tau^{(1)}(L) < \tau^{*,1}(L/M)) = P(\{\tau^{(1)}(L) < \tau^{*,1}(L/M)\} \cap A_2(L)) + P(\{\tau^{(1)}(L) < \tau^{*,1}(L/M)\} \cap A_3(L)).$$

However, by the strong Markov property,

$$P(\{\tau^{(1)}(L) < \tau^{*,1}(L/M)\} \cap A_3(L)) = P(A_3(L))P(\tau^{(1)}(L) < \tau^{*,1}(L/M)).$$

Therefore,

$$(4.39) \quad P(\tau^{(1)}(L) < \tau^{*,1}(L/M)) = \frac{P(\{\tau^{(1)}(L) < \tau^{*,1}(L/M)\} \cap A_2(L))}{1 - P(A_3(L))}.$$

Let $\xi_i = (1-r)X_i^{(1)} - rY_i^{(i)}$, $i = 1, 2, \dots$. Taking into account that the *off* times are bounded, we can repeat now the argument used to estimate $I_{11}(x)$ in the proof of Proposition 2.1', to conclude that

$$P(\{\tau^{(1)}(L) < \tau^{*,1}(L/M)\} \cap A_2(L)) \leq \left(P\left(\bigvee_{n=0}^{\infty} S_n^{(\xi)} > \frac{L}{M(1-r)} \right) \right)^{[M/2]},$$

and $P\left(\bigvee_{n=0}^{\infty} S_n^{(\xi)} > L\right)$ is regularly varying in L with index $-(\alpha - 1)$. Moreover,

$$1 - P(A_3(L)) \geq P(A_1(L)) \geq P\left(X_1^{(1)} > \frac{L}{M(1-r)}\right) = \bar{F}_{\text{on}}\left(\frac{L}{M(1-r)}\right).$$

We conclude by (4.39) and (4.36) that

$$\limsup_{L \rightarrow \infty} P(\tau^{(1)}(L) < \tau^{*,1}(L/M)) \leq \limsup_{L \rightarrow \infty} \frac{\left(P\left(\bigvee_{n=0}^{\infty} S_n^{(\xi)} > \frac{L}{M(1-r)} \right) \right)^{[M/2]}}{\bar{F}_{\text{on}}\left(\frac{L}{M(1-r)}\right)} = 0.$$

This proves (4.37) and so (4.30) is proven as well. We observe at this point that the above argument that allowed us to assume the initial delays being equal to 0 shows that we have also proved that for every N satisfying (4.29)

$$(4.40) \quad \lim_{L \rightarrow \infty} \sup_H P_{H,0}(\tau(L) < \tau^*(L/N)) = 0.$$

The next step in the proof of (4.28) is to show that two “long” *on* periods are “unlikely to happen simultaneously”. Formally, let

$$(4.41) \quad Q_L = \inf\{t \geq 0 : \text{at time } t \text{ there are two } \textit{on} \text{ periods running, each of length at least } L\}.$$

We claim that there is a function $\gamma_L \rightarrow 0$ as $L \rightarrow \infty$ such that

$$(4.42) \quad \lim_{L \rightarrow \infty} P_{H,x_0}(Q_L \leq \gamma_L^{-1}(\bar{F}_{\text{on}}(L))^{-1}) = 0.$$

Of course, Q_L will only decrease if we assume that all delay times and *off* times are equal to 0, and Q_L is unaffected by x_0 . We will, therefore, once again drop the subscripts from P and E , and assume that all *off* times are equal to 0.

For $j_1, j_2 = 1, \dots, k$, $j_1 \neq j_2$, let

$$Q_L^{(j_1, j_2)} = \inf\{t \geq 0 : \text{at time } t, \text{ the processes } X_{j_1}(\cdot) \text{ and } X_{j_2}(\cdot) \text{ both have } \textit{on} \text{ periods running, each of length at least } L\}.$$

Then

$$Q_L = \bigwedge_{j_1 \neq j_2} Q_L^{(j_1, j_2)},$$

and so for any $q > 0$,

$$(4.43) \quad P(Q_L \leq q) \leq \frac{k(k-1)}{2} P(Q_L^{(1,2)} \leq q).$$

Let

$$(4.44) \quad Q_L^{(1)} = \inf\{t \geq 0 : \text{at time } t, X_1(\cdot) \text{ begins an } \textit{on} \text{ period of length at least } L \text{ during which } X_2(\cdot) \text{ also begins an } \textit{on} \text{ period of length at least } L\}$$

with $Q_L^{(2)}$ defined similarly. Then

$$Q_L^{(1,2)} \geq Q_L^{(1)} \wedge Q_L^{(2)},$$

which means that for any $q > 0$,

$$(4.45) \quad P(Q_L^{(1,2)} \leq q) \leq 2P(Q_L^{(1)} \leq q).$$

Let Z_k , $k = 1, 2, \dots$ be an iid sequence, such that

$$Z_1 \stackrel{d}{=} \sum_{i=1}^{G_L} \hat{X}_i,$$

where G_L is a geometric random variable with parameter $\bar{F}_{\text{on}}(L)$, independent of an iid sequence \hat{X}_i , $i = 1, 2, \dots$, with common law

$$P(\hat{X}_1 \in A) = P(X_1 \in A | X_1 \leq L) = \frac{1}{F_{\text{on}}(L)} \int_0^L 1(x \in A) F_{\text{on}}(dx).$$

Then Z_1 represents the first time $X_1(\cdot)$ starts an *on* period of length at least L . Let H_L be yet another geometric random variable, independent of the sequence Z_k , $k = 1, 2, \dots$, this time with parameter p_L defined as follows. Let W be a random variable with distribution

$$P(W \in A) = P(X_1 \in A | X_1 > L) = \frac{1}{\bar{F}_{\text{on}}(L)} \int_L^\infty 1(x \in A) F_{\text{on}}(dx),$$

and independent of $X_2(\cdot)$. Recall that at time 0 the process $X_2(\cdot)$ starts an *on* interval. Then define

$$(4.46) \quad p_L = P\left(\inf\{S_n^{(2)} : n \geq 1, X_{n+1}^{(2)} \geq L\} \leq W\right).$$

That is, p_L is the probability that $X^{(2)}(\cdot)$ begins an *on* period of length at least L during an *on* period of $X^{(1)}(\cdot)$, whose length is at least L . If $X^{(2)}(\cdot)$ does not start such an *on* period, we then have to wait till the next *on* period of $X^{(1)}(\cdot)$ whose length is at least L . Therefore,

$$(4.47) \quad Q_L^{(1)} \stackrel{\text{st}}{\geq} \sum_{n=1}^{H_L} Z_n.$$

We claim that

$$(4.48) \quad p_L \rightarrow 0 \text{ as } L \rightarrow \infty.$$

Indeed,

$$(4.49) \quad p_L = \frac{1}{\bar{F}_{\text{on}}(L)} \int_L^\infty P(Z_1 \leq t) F_{\text{on}}(dt) = \frac{1}{\bar{F}_{\text{on}}(L)} \int_L^\infty P\left(\sum_{i=1}^{G_L} \hat{X}_i \leq t\right) F_{\text{on}}(dt).$$

Observe that for every $t \geq L$,

$$(4.50) \quad \begin{aligned} P\left(\sum_{i=1}^{G_L} \hat{X}_i \leq t\right) &\leq P\left(G_L \leq \frac{1}{(\bar{F}_{\text{on}}(L))^{1/2}}\right) + P\left(\inf_{n > (\bar{F}_{\text{on}}(L))^{-1/2}} \sum_{i=1}^n \hat{X}_i - n \frac{\mu_{\text{on}}}{2} \leq 0\right) \\ &\quad + P\left(\frac{\mu_{\text{on}}}{2} G_L \leq t\right). \end{aligned}$$

Notice that the first term in the right hand side of (4.50) goes to 0 as $L \rightarrow \infty$, and that by the strong law of large numbers so does the second term in the the right hand side of (4.50). Since neither of these two terms depends on t , (4.48) will follow once we prove that

$$(4.51) \quad \lim_{L \rightarrow \infty} \frac{1}{\bar{F}_{\text{on}}(L)} \int_L^\infty P\left(\frac{\mu_{\text{on}}}{2} G_L \leq t\right) F_{\text{on}}(dt) = 0.$$

However, for all L big enough and all $t \geq L$ we have $[2t/\mu_{\text{on}}] \geq 1$, and so

$$P\left(\frac{\mu_{\text{on}}}{2}G_L \leq t\right) = 1 - \left(F_{\text{on}}(L)\right)^{[2t/\mu_{\text{on}}]} \leq \left[\frac{2t}{\mu_{\text{on}}}\right] \bar{F}_{\text{on}}(L).$$

Therefore,

$$\limsup_{L \rightarrow \infty} \frac{1}{\bar{F}_{\text{on}}(L)} \int_L^\infty P\left(\frac{\mu_{\text{on}}}{2}G_L \leq t\right) dt \leq \frac{2}{\mu_{\text{on}}} \limsup_{L \rightarrow \infty} \int_L^\infty t F_{\text{on}}(dt) = 0$$

since $\mu_{\text{on}} < \infty$. This proves (4.51), and, therefore, we have (4.48).

Define now $\gamma_L = p_L^{1/2}$. Observe that by (4.43) (4.45) and (4.47) we have

$$(4.52) \quad P\left(Q_L \leq \gamma_L^{-1}(\bar{F}_{\text{on}}(L))^{-1}\right) \leq k(k-1)P\left(\sum_{n=1}^{H_L} Z_n \leq p_L^{-1/2}(\bar{F}_{\text{on}}(L))^{-1}\right).$$

However,

$$\sum_{n=1}^{H_L} Z_n \stackrel{d}{=} \sum_{i=1}^{J_L} \hat{X}_i,$$

where J_L is a geometric random variable with parameter $p_L \bar{F}_{\text{on}}(L)$, independent of $\{\hat{X}_i, i = 1, 2, \dots\}$. Since, as $L \rightarrow \infty$,

$$(p_L \bar{F}_{\text{on}}(L))J_L \Rightarrow E(1)$$

where $E(1)$ is a standard exponential random variable and we immediately obtain (4.42) using (4.52) and the law of large numbers.

Let us go back now to the proof of (4.28). Observe, first of all, that for all L big enough,

$$P_{H, x_0}\left(\tau(L) < \tau^*(L(1-\epsilon))\right) \leq P_{H, 0}\left(\tau(L(1-\epsilon/2)) < \tau^*(L(1-\epsilon))\right).$$

Therefore, it is enough to prove (4.28) for $x_0 = 0$, for the general case will follow by making ϵ smaller. We will, therefore, use the notation P_H and E_H , when $x_0 = 0$.

Fix any N satisfying (4.29) and big enough to make the right hand side of (4.53) below positive, and observe that it is enough to prove (4.28) for $\epsilon = 1/N$. Fix, further, a ρ satisfying

$$(4.53) \quad (k-1)\rho \leq \frac{\epsilon}{2} - \epsilon^2.$$

We have

$$(4.54) \quad \begin{aligned} P_H(\tau(L) < \tau^*(L(1-\epsilon))) &\leq P_H\left(\tau(L) < \tau^*(L(1-\epsilon)), \tau(L\epsilon^2) \geq \tau^*(L\epsilon^3), \tau(L) < Q_{\rho L}\right) \\ &\quad + P_H(\tau(L\epsilon^2) < \tau^*(L\epsilon^3)) + P_H(\tau(L) \geq Q_{\rho L}). \end{aligned}$$

Observe that by (2.30)

$$\lim_{L \rightarrow \infty} P_H(\tau(L\epsilon^2) < \tau^*(L\epsilon^3)) = 0.$$

Furthermore, by (4.42) and (4.18)

$$\begin{aligned} P_H(\tau(L) \geq Q_{\rho L}) &\leq P_H\left(Q_{\rho L} \leq (\gamma_{\rho L})^{-1}(\bar{F}_{\text{on}}(\rho L))^{-1}\right) + P_H\left(\tau(L) \geq (\gamma_{\rho L})^{-1}(\bar{F}_{\text{on}}(\rho L))^{-1}\right) \\ &\leq P_H\left(Q_{\rho L} \leq (\gamma_{\rho L})^{-1}(\bar{F}_{\text{on}}(\rho L))^{-1}\right) + \gamma_{\rho L} \bar{F}_{\text{on}}(\rho L) E_H \tau(L) \\ &\rightarrow 0 \end{aligned}$$

as $L \rightarrow \infty$. Therefore, (4.28) will follow once we prove that

$$(4.55) \quad \lim_{L \rightarrow \infty} P_H \left(\tau(L) < \tau^*(L(1-\epsilon)), \tau(L\epsilon^2) \geq \tau^*(L\epsilon^3), \tau(L) < Q_{\rho L} \right) = 0.$$

As a matter of fact, we will prove an even stronger statement. We will prove that

$$(4.56) \quad \lim_{L \rightarrow \infty} \sup_H P_H \left(\tau(L) < \tau^*(L(1-\epsilon)), \tau(L\epsilon^2) \geq \tau^*(L\epsilon^3), \tau(L) < Q_{\rho L} \right) = 0.$$

Let

$$B(L) = \left\{ \tau(L) < \tau^*(L(1-\epsilon)), \tau(L\epsilon^2) \geq \tau^*(L\epsilon^3), \tau(L) < Q_{\rho L} \right\}.$$

We split this event into two events, $B_1(L)$ and $B_2(L)$, accordingly to the following two possibilities.

- (i) After starting the first “wet period” the process $X^{(k)}(\cdot)$ reaches level L before returning to 0.
- (ii) The process $X^{(k)}(\cdot)$ returns to 0 before reaching level L .

Let us look at the event $B_1(L)$ first. Since $B_1(L) \subseteq B(L)$, at time $\tau(L\epsilon^2)$ at most one of the k *on/off* processes has an *on* period whose length is at least ρL . Depending on whether the number of such *on/off* processes is 0 or 1, we split the event $B_1(L)$ into $B_{11}(L)$ and $B_{12}(L)$. Let us look, for example, at $B_{12}(L)$. The treatment of the event $B_{11}(L)$ is similar.

Since $B_{12}(L) \subseteq B(L)$, the single *on* period running at time $\tau(L\epsilon^2)$ of length least ρL , has length not exceeding $\frac{L(1-\epsilon)}{(1-r)}$. Let us now modify the state of the system at time $\tau(L\epsilon^2)$ in the following way. Bring all the work remaining in the presently running *on* periods in one “lump” at time $\tau(L\epsilon^2)$, and attach the remainders of these *on* periods to the subsequent *off* periods. Obviously, this action can only make $\tau(L)$ smaller, and so it can only increase the probability of the event $B_{12}(L)$. Observe that, after this action, the state of the system does not exceed

$$L\epsilon^2 + L(1-\epsilon) + (k-1)\rho L \leq L(1-\epsilon/2)$$

by (4.53). We increase, if necessary, the state of the system to exactly $L(1-\epsilon/2)$ (only increasing in the process the probability of the event $B_{12}(L)$). We conclude that for some H_0 ,

$$(4.57) \quad P_H(B_{12}(L)) \leq P_{H_0, L(1-\epsilon/2)} \left(\sup_{t \geq 0} V^{(k)}(t) \geq L \right),$$

where $\{V^{(k)}(t), t \geq 0\}$ is given by

$$dV^{(k)}(t) = Z^{(k)}(t) dt - r dt,$$

with $\{Z^{(k)}(t), t \geq 0\}$ given by (4.3) and (4.7). However,

$$V^{(k)}(t) = V_1(t) + \dots + V_k(t), \quad t \geq 0,$$

where for $j = 1, \dots, k$ the process $\{V_j(t), t \geq 0\}$ is defined by

$$dV_j(t) = Z_j(t) dt - \frac{r}{k} dt,$$

with $V_j(0) = L(1-\epsilon/2)/k$, and with initial delay governed by the j th marginal law, $H^{(j)}$, of H_0 . We conclude immediately by (4.57) that

$$P_H(B_{12}(L)) \leq k \sup_{H^{(1)}} P_{H^{(1)}, L(1-\epsilon/2)/k} \left(\sup_{t \geq 0} V_1(t) \geq \frac{L}{k} \right) = k P \left(\sup_{t \geq 0} V_1(t) \geq \frac{L\epsilon}{2k} \right),$$

where P without a subscript indicates, as usual, absence of delay and zero initial state. Because of the negative drift, we conclude that

$$(4.58) \quad \lim_{L \rightarrow \infty} \sup_H P_H(B_{12}(L)) = 0.$$

In exactly the same way one can show that

$$(4.59) \quad \lim_{L \rightarrow \infty} \sup_H P_H(B_{11}(L)) = 0.$$

Finally, we consider the event $B_2(L)$ above. Consider the two possibilities that are feasible after the process reaches 0: either after that time the state of the system reaches the level $L\epsilon^2$ before the beginning of the first *on* period of length at least $L\epsilon^3$, or not. Accordingly, by the strong Markov property,

$$P_H(B_2(L)) \leq P_H(\tau^*(L\epsilon^3) < \tau^*(L(1-\epsilon))) \left(\sup_G P_G(\tau(L\epsilon^2) \leq \tau^*(L\epsilon^3)) + \sup_G P_G(B(L)) \right).$$

Taking supremum over H , we obtain

$$\begin{aligned} \sup_H P_H(B(L)) &\leq \sup_H P_H(B_{11}(L)) + \sup_H P_H(B_{12}(L)) \\ &\quad + P(\tau^*(L\epsilon^3) < \tau^*(L(1-\epsilon))) \left(\sup_H P_H(\tau(L\epsilon^2) \leq \tau^*(L\epsilon^3)) + \sup_H P_H(B(L)) \right), \end{aligned}$$

which is the same as

$$(4.60) \quad \sup_H P_H(B(L)) \leq \frac{\sup_H P_H(B_{11}(L)) + \sup_H P_H(B_{12}(L)) + \sup_H P_H(\tau(L\epsilon^2) \leq \tau^*(L\epsilon^3))}{P(\tau^*(L\epsilon^3) = \tau^*(L(1-\epsilon)))}.$$

Now (4.56) follows from (4.60), (4.58), (4.59), (4.40) and the fact that

$$P(\tau^*(L\epsilon^3) = \tau^*(L(1-\epsilon))) = \frac{\bar{F}_{\text{on}}(L(1-\epsilon))}{\bar{F}_{\text{on}}(L\epsilon^3)} \rightarrow \left(\frac{\epsilon^3}{1-\epsilon} \right)^\alpha > 0$$

as $L \rightarrow \infty$ by the regular variation. This completes the proof of the theorem. \square

5. Appendix.

We will construct explicitly a stationary version of $\{X^{(k)}(t), t \geq 0\}$. In fact, we will construct a stationary version of $\{X^{(k)}(t), -\infty < t < \infty\}$ on the whole real line. We proceed as follows. Let $\{Z_j(t), -\infty < t < \infty\}$, $j = 1, \dots, k$ be stationary. Then $\{Z^{(k)}(t), -\infty < t < \infty\}$ defined by (4.3) is itself stationary. We will construct $\{X^{(k)}(t), -\infty < t < \infty\}$ by exhibiting a measurable function $\varphi : D(\mathbb{R}, \{0, 1, \dots, k\}) \rightarrow \mathbb{R}_+$ such that the process

$$(5.1) \quad X^{(k)}(t) = \varphi(Z^{(k)}(\cdot + t)), \quad -\infty < t < \infty$$

satisfies (4.4). Clearly, the process given by (5.1) is stationary.

Let $\mathbf{x} = (x(u), -\infty < u < \infty) \in D(\mathbb{R}, \{0, 1, \dots, k\})$, and let

$$\dots < T_{-2} < T_{-1} < 0 \leq T_1 < T_2 < \dots$$

be the epoch times for $\{x(u), -\infty < u < \infty\}$. That is, at each time T_i , $\{x(u), -\infty < u < \infty\}$ changes its value, and let

$$\mathcal{Z}_{0,1} = \{i : x(T_i - 0) = 0, x(T_i + 0) = 1\}.$$

For $\mathbf{x} = (Z^{(k)}(t), -\infty < t < \infty)$, the points T_i with $i \in \mathcal{Z}_{0,1}$ are the times when one of the *on/off* processes begins an *on* period, following a period when all k processes were off. For every $i \in \mathcal{Z}_{0,1}$ let L_i be the length of the “wet period” commencing at T_i . Formally, let $X^{(k)}(T_i) = 0$, and define $\{X^{(k)}(t), u \geq T_i\}$ by the following analogue of (4.4)

$$(5.2) \quad dX^{(k)}(u) = x(u) du - r(X^{(k)})(u) du.$$

Then

$$L_i = \inf\{u > 0 : X^{(k)}(T_i + u) = 0\}.$$

If we denote

$$G_1 = \left\{ \mathbf{x} : \text{for every } i \in \mathcal{Z}_{0,1}(\mathbf{x}), L_i(\mathbf{x}) < \infty \right\},$$

then the rate condition (4.5) implies that the set $G_1 \in D(\mathbb{R}, \{0, 1, \dots, k\})$ has full measure under the probability measure μ_0 induced by $\{Z^{(k)}(t), -\infty < t < \infty\}$ on $D(\mathbb{R}, \{0, 1, \dots, k\})$. Now, for a $-\infty < t < \infty$ let

$$\mathcal{A}_{0,1}(t) = \{T_i \leq t : i \in \mathcal{Z}_{0,1}\},$$

and define

$$(5.3) \quad T(t) = \inf\{T_i \in \mathcal{A}_{0,1}(t) : T_i + L_i \geq t\}.$$

Denote

$$G_2 = \left\{ \mathbf{x} : \text{for every } -\infty < t < \infty, T_i + L_i \geq t \text{ for only finitely many } i \in \mathcal{A}_{0,1}(t) \right\}.$$

We claim that the set $G_2 \in D(\mathbb{R}, \{0, 1, \dots, k\})$ has full measure under the probability measure μ_0 . That is, we claim that

$$(5.4) \quad P\left(\left\{ \omega : \text{for some } -\infty < t < \infty, T_i + L_i \geq t, \text{ for infinitely many } i \in \mathcal{A}_{0,1}(t)(Z^{(k)}(\cdot)(\omega)) \right\}\right) = 0.$$

Now, concentrating on the rational times only, we see that (5.4) will follow once we prove that for any $-\infty < t < \infty$ we have

$$(5.5) \quad P\left(\left\{ \omega : T_i + L_i > t \text{ for infinitely many } i \in \mathcal{A}_{0,1}(t)(Z^{(k)}(\cdot)(\omega)) \right\}\right) = 0.$$

We check (5.5) for $t = 0$. We have

$$\begin{aligned} & P\left(\left\{ \omega : T_i + L_i > 0 \text{ for infinitely many } i \in \mathcal{A}_{0,1}(0)(Z^{(k)}(\cdot)(\omega)) \right\}\right) \\ & \leq P\left(\limsup_{n \rightarrow -\infty} \left\{ \omega : \sum_{j=1}^k \int_n^0 Z_j(u) du \geq r|n| - r \right\}\right) \end{aligned}$$

and the last probability goes to 0 as $n \rightarrow -\infty$ because of the rate condition (4.5) and the obvious fact that

$$\frac{1}{|n|} \sum_{j=1}^k \int_n^0 Z_j(u) du \rightarrow k \frac{\mu_{\text{on}}}{\mu}$$

as $n \rightarrow -\infty$. This proves (5.5) and so (5.4) holds.

Finally, let

$$G_3 = \left\{ \mathbf{x} : \text{for every } -\infty < t < \infty, \text{ there is a } u \leq t \text{ such that } x(u) = k \right\}.$$

Note that $F_{\text{on}} * F_{\text{off}}$ being non-arithmetic implies that the set $G_3 \in D(\mathbb{R}, \{0, 1, \dots, k\})$ has full measure under the probability measure μ_0 . Finally, we let $G = G_1 \cap G_2 \cap G_3$, and we note that G has also full measure under the probability measure μ_0 .

We will define now the function $\varphi(\mathbf{x})$, $\mathbf{x} \in D(\mathbb{R}, \{0, 1, \dots, k\})$. If $\mathbf{x} \in G^c$, let $\varphi(\mathbf{x}) = 0$ (an arbitrary number). Suppose now that $\mathbf{x} \in G$. Since $\mathbf{x} \in G_3$, we may meaningfully define $s(\mathbf{x}) = \sup\{u \leq 0 : x(u) = k\}$. Since $\mathbf{x} \in G_2$, $T(\mathbf{x}, s(\mathbf{x}))$ is well defined by (5.3). Now define $\varphi(\mathbf{x})$ as $X^{(k)}(0)$ given by solving (5.2) on $(T(\mathbf{x}, s(\mathbf{x})), \infty)$, with $X^{(k)}(T(\mathbf{x}, s(\mathbf{x}))) = 0$. Intuitively, $s(Z^{(k)}(\cdot))$ is the last time before time 0 that all k on/off processes were on, which is identifiable as a time point when the system is in a “wet period”, and $T(Z^{(k)}(\cdot), s(Z^{(k)}(\cdot)))$ is the beginning of that “wet period”, at which time the system must be empty. This constructs the required function $\varphi : D(\mathbb{R}, \{0, 1, \dots, k\}) \rightarrow \mathbb{R}_+$, and thus a stationary version of a process satisfying (3.4).

REFERENCES

- Abate, J., Choudhury, G. and Whitt, W., *Waiting-time tail probabilities in queues with long-tail service-time distributions*, Queueing Systems Theory Appl. **16** (1994), 311–338.
- Anick, D., Mitra, D. and Sondhi, M., *Stochastic theory of a data-handling system with multiple sources*, Bell System Tech. J. **61** (1982), 1871–1894.
- Asmussen, S., *Applied Probability and Queues*, J. Wiley and Sons, New York, 1987.
- Asmussen, S. and Perry, D., *On cycle maxima, first passage problems and extreme value theory for queues*, Stochastic Models **8** (1992), 421–458.
- Avram, Florin and Taqqu, Murad S., *Weak convergence of moving averages with infinite variance*, Dependence in Probability and Statistics, Progr. Probab. Statist., 11, Birkhauser Boston, Boston, MA, 1986, pp. 399–415.
- Avram, Florin and Taqqu, Murad S., *Probability bounds for M -Skorohod oscillations*, Stochastic Processes and their Applications **33** (1989), 63–72.
- Avram, Florin and Taqqu, Murad S., *Weak convergence of sums of moving averages in the α -stable domain of attraction*, Ann. Probability **20** (1992), 483–503.
- Berger, A. and Whitt, W., *Maximum values in queueing processes*, Probab. Engrg. Inform. Sci. **9** (1995), 375–409.
- Billingsley, P., *Convergence of Probability Measures*, Wiley, New York, 1968.
- Bingham, N., Goldie, C. and Teugels, J., *Regular Variation*, Encyclopedia of Mathematics and its Applications, Cambridge University Press, Cambridge, UK, 1987.
- Brichet, F., Roberts, J., Simonian, A. and Veitch, D., *Heavy traffic analysis of a storage model with long range dependent on/off sources*, Preprint (1996).
- Choudhury, G. and Whitt, W., *Long-tail buffer-content distributions in broadband networks*, Preprint, AT&T Bell Laboratories (1995).
- Crovella, M. and Bestavros, A., *Explaining world wide web traffic self-similarity*, Preprint available as TR-95-015 from {crovella,best}@cs.bu.edu (1995).
- Cunha, Bestavros, A. and Crovella, M., *Characteristics of www client-based traces*, Preprint available as BU-CS-95-010 from {crovella,best}@cs.bu.edu (1995).
- Feller, W., *An Introduction to Probability Theory and its Applications*, Vol. II, second edition, Wiley, New York, 1971.
- Heath, D., Resnick, S. and Samorodnitsky, G., *Heavy tails and long range dependence in on/off processes and associated fluid models*, Available as TR1144.ps.Z at <http://www.orie.cornell.edu/trlist/trlist.html> (1996).
- Iglehart, D., *Extreme values in the GI/G/1 queue*, Ann. Math. Statist. **43** (1972), 627–635.
- Kella, O. and Whitt, W., *A storage model with a two-state random environment*, Opns. Res. **40** (1992), 257–262.
- Kwapień, S. and Woyczynski, N.A., *Random Series and Stochastic Integrals: Single and Multiple*, Birkhauser, Boston, 1992.
- Leland, W., Taqqu, M., Willinger, W., Wilson, D., *On the self-similar nature of ethernet traffic*, Proceedings of the ACM/SIGCOMM '93, San Francisco, ACM/SIGCOMM Computer Communications Review **23** (1993), 183–193.
- Leland, W., Taqqu, M., Willinger, W., Wilson, D., *On the self-similar nature of ethernet traffic (extended version)*, IEEE/ACM Transactions on Networking **2** (1994), 1–15.
- Pacheco, A. and Prabhu, N.U., *A Markovian storage model*, Ann. Applied Probability **6** (1996), 76–91.
- Prabhu, N.U. and Pacheco, A., *A storage model for data communication systems*, Queueing Systems **19** (1995), 1–40.
- Pratt, J., *On interchanging limits and integrals*, Ann. Math. Statist. **31** (1960), 74–77.
- Resnick, Sidney, *Point processes, regular variation and weak convergence*, Advances Appl. Probability **18** (1986), 66–138.
- Resnick, Sidney, *Extreme Values, Regular Variation, and Point Processes*, Springer-Verlag, New York, 1987.
- Resnick, S., *Adventures in Stochastic Processes*, Birkhauser, Boston, 1992.
- Willinger, W., Taqqu, M., Leland, W. and Wilson, D., *Self-similarity in high-speed packet traffic: analysis and modeling of ethernet traffic measurements*, Statistical Science **10** (1995), 67–85.

Willinger, W., Taqqu, M., Sherman, R. and Wilson, D., *Self-similarity through high-variability: Statistical analysis of ethernet LAN traffic at the source level*, Preprint (1995).

DAVID HEATH, CORNELL UNIVERSITY, SCHOOL OF OPERATIONS RESEARCH AND INDUSTRIAL ENGINEERING, ETC BUILDING, ITHACA, NY 14853 USA

E-mail address: davidh@orie.cornell.edu

SIDNEY I. RESNICK, CORNELL UNIVERSITY, SCHOOL OF OPERATIONS RESEARCH AND INDUSTRIAL ENGINEERING, ETC BUILDING, ITHACA, NY 14853 USA

E-mail address: sid@orie.cornell.edu

GENNADY SAMORODNITSKY, CORNELL UNIVERSITY, SCHOOL OF OPERATIONS RESEARCH AND INDUSTRIAL ENGINEERING, ETC BUILDING, ITHACA, NY 14853 USA

E-mail address: gennady@orie.cornell.edu