



Long Term Dependences and Heavy Tails in Traffics and Queues Generated by Memoriless on/off Sources in Series

Philippe Jacquet

► **To cite this version:**

Philippe Jacquet. Long Term Dependences and Heavy Tails in Traffics and Queues Generated by Memoriless on/off Sources in Series. [Research Report] RR-3516, INRIA. 1998. <inria-00073168>

HAL Id: inria-00073168

<https://hal.inria.fr/inria-00073168>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*Long term dependences and heavy tails in
traffics and queues generated by memoriless
on/off sources in series*

Philippe Jacquet

No 3516

———— THÈME 1 ————



*Rapport
de recherche*

Long term dependences and heavy tails in traffics and queues generated by memoriless on/off sources in series

Philippe Jacquet

Thème 1 — Réseaux et systèmes
Projet Hipercom

Rapport de recherche n° 3516 — — 20 pages

Abstract: This report presents the analytical study of a surprising case where a denumerable set of independent and memoriless on/off sources in series creates a traffic with long term dependences. We use Mellin transform to characterize the asymptotic behavior of those dependences. We use similar tools to analyse the queue size of a server subject to these sources in series. This architecture roughly models a shared network like internet. We show that the queue size distribution has a heavy tail. These results agree with the recent experimental datas collected on Web activity.

Key-words: on/off sources, long term dependences, heavy tails, Mellin transform, networks, Web

(Résumé : tsvp)

Longues dépendances et queues lourde dans les trafics et les files d'attente issus d'une série de sources "on/off" sans mémoire

Résumé : Ce rapport présente l'étude analytique d'un cas d'école surprenant où un ensemble dénombrable de sources on/off en série, indépendantes et sans mémoire, créent un trafic à longues dépendances. On utilise la transformée de Mellin pour mettre en évidence le comportement asymptotique de ces longues dépendances. On utilise les mêmes outils pour étudier que la file d'attente d'un serveur soumis à ces sources en série. Ce modèle approche le modèle d'un réseau à ressources partagées comme internet. On montre que la taille de la file d'attente présente une queue lourde de distribution. Ces résultats sont en accord avec les données expérimentales récoltées récemment sur l'activité sur le Web.

Mots-clé : Sources on/off, longues dépendances, queues lourdes, transformée de Mellin, réseaux, Web

1 Introduction

Long term dependences and self-similarities have recently been detected on Web traffic. This phenomenon is of theoretical importance since it seems to contradict the Law of Large Numbers and its Poisson derivations. Indeed on large time scale the aggregation of independent sources should look like a Poisson flow. Measured Web traffic contradicts this statement [6].

1.1 The Law of Large Numbers and the Poisson law

Before entering in the details of this study, we make a very brief incursion in History. In the very beginning of the XIXth century, Napoleon ordered the mathematician Poisson to find an explanation about the statistical occurrence of high rank officers in the *Grande Armée* killed by accident by falling from their horses. Poisson did not find any explanation, and in fact *proved* that there were no explanation at all: the high rank officers have the same odds to fall down their horses as any other horseman in the army, and that the fluctuations detected in the statistical occurrence of such events were just consequence of the geometric distribution induced by the Law of Large Numbers. Napoleon, happy enough that his high rank officers were proved to be not more stupid than the high rank officers of any other army, went ahead for new military adventures.

The Poisson distribution survived Napoleon and found a new confirmation in the atomic physics in the early 1900s. The spontaneous emission of particles in the radioactivity process follows Poisson law. In this case the atoms play the role of horses and the neutrons, the role of high rank officers.

After such nice achievements, they were no reason why Poisson law should not also apply to telecommunication traffics. In fact it did as long communication traffics were telephone-voice oriented. The problems eventually arise when the telecommunication traffic turn out to be data oriented as it prevails on Web traffic.

1.2 Long term dependences and Internet packet loss

What is a long term dependence? There are many definitions, which are reviewed in details in other chapters. Let us stay on a very basic flavour. We call $I(t)$ the intensity of traffic at time t . We denote $C(x)$ the covariance of $I(t)$ and $I(t+x)$ when t varies:

$$C(x) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T I(t+x)I(t)dt - \left(\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T I(t)dt \right)^2 \quad (1)$$

When the traffic has long term dependences, then $C(x) \asymp Bx^{-\beta}$ with $\beta < 1$. When the traffic fits the Law of Large Numbers, then $\beta \geq 1$. Quantity β is frequently referred as the Hultz parameter. When the traffic is pure Poisson, *i.e.* completely memoriless, then $C(x) = 0$ (or abusively, $\beta = \infty$).

Long term dependences are interesting not only because they contradict Poisson law, they are also interesting because they lead to very practical but annoying effects. One of

these effects is that it dramatically increases packet loss in Internet. The Internet important actor is the router. In some simple model, the router can be seen as a buffer served by a network server. When the buffer overflows, then some packets are lost. The lost packets must be resent following TCP-IP, therefore adding extra delays and traffics.

If we simply model the router by a $M/M/1$ queue with input rate λ and service rate 1, then the probability p_n that the queue length be greater than n is exactly λ^n . In a first order approximation quantity p_n can be identified with the packet loss rate in a buffer of size n . Therefore to block packet loss below some acceptable level ε it suffices to make buffer capacity greater than $\frac{\log \varepsilon}{\log \lambda}$, *i.e.* a logarithmic function of $1/\varepsilon$. In general telecom designers plan on $\varepsilon = 10^{-6}$ and $\lambda < 0.8$.

The Web traffic experiences long term dependences [6], and such profiles lead to loss and retry rates $p_n \asymp B' x^{-\beta}$ for some β . In other words, the queue size distribution has a heavy tail. Under this condition it is clear that buffer capacity would need to be raised to $(\frac{B'}{\varepsilon})^{1/\beta}$ to keep packet loss under the acceptable level, which no longer leads to a logarithmic function of $1/\varepsilon$, but to a *polynomial* function of $1/\varepsilon$. Indeed this minimal size would be of several magnitude orders higher than the capacity obtained with the Poisson model. In fact actual router capacities reveal to be dangerously underestimated with regards to this new traffic conditions.

1.3 The origin of long term dependences

There is a question still unanswered: why the Web data traffic does not fit Poisson law? We can retain three main basic sources of plausible explanations:

1. The sources of traffic are not independent;
2. The communication protocol creates the long term dependences;
3. The sources have long term dependences profile.

Each of these explanations is plausible. Indeed if a large number of sources are correlated, then it is clear that arbitrary long dependences will be very easy to create. For example it is known that the telephone network is overloaded just after a soccer contest. Secondly the protocol TCP-IP is known to create great disturbance in the retry process with its binary exponential load split process. Thirdly, it suffices that only one source has long term dependences to create long term dependences in the aggregated traffic.

These explanations however have their own limitations. In the context of the first explanation, there is no reason to have correlated sources in the Web traffic similar to those of the *end of soccer contest* syndrom. In the second explanation if the original data traffic were Poisson with $\lambda < 1$, then there would be not enough packet loss to let TCP-IP making a decisive impact. In the third explanation, if there is only one source with long term dependences, then quantity $C(x)$ will basically only reflect the long term dependence of this very source, thus a very marginal effect. Conversely if all sources have identically

long term dependences, the aggregation of several i.i.d sources will converges to a Poisson traffic according to the Law of Large Numbers.

The author of the present chapter risks a fourth explanation:

4. Independent sources with divergent profiles create the long term dependence.

The author will analyse this option in the simplified case where each source has a memoriless profile. The *paradoxal* result is that a potentially infinite set of independent sources with memoriless profiles can create long term dependences. The fact that the source have memoriless profile is not a fundamental assumption in this model but it considerably simplifies the model to make it tractable.

2 Statements of the model and results

One of our primary aim is to illustrate the use of some tools borrowed from Complex Analysis and more specifically the use of the Mellin transform in order to estimate the performance of queueing models involving aggregations of independent “on/off” input sources. This model finds application in networking and in particular has interesting properties with regards to the analysis of long term dependences which arise in Web traffic.

In particular we use the Mellin transform. The Mellin transform is a derivation of the Laplace transform but instead of tracking asymptotic exponential bounds it tracks polynomial bounds. Therefore it is a very powerful tool for long term dependence analysis.

We would like to coin here the term *Analytical Information Theory* to describe problems of information theory that are solved by analytical methods borrowed from complex analysis [7, 10]. Look at [11] for a survey and at [1] for a detailed description of the tools and proofs used in the present paper.

2.1 The “on/off” sources as traffic models

Let us consider a large multiple-access network like the Internet. It can also be an Ethernet link or an ATM switch. By “large” we mean a large population of users and by “multiple-access” we mean that a single resource of communication has to be shared. We assume that users are independent from each other, which is a commonly accepted assumption. We also assume that every user behaves like an “on/off” input source.

“On/off” sources are very useful for realistic networking modelisation. User activity is characterized by *wake* periods interleaved with *sleep* periods. During a wake period, the user has intense activity, *i.e.* creates messages with constant peak rate, λ , (*e.g.* for files retrieval, network browsing, *etc*). When the user goes for a coffee break, or works on local device, or simply when he goes home, he is in sleep period. During a sleep period, the user has very low activity, to simplify we will assume that it has no activity at all.

To simplify we also model the multiple access communication resource as a queueing with a single server, and messages plays the role of customers stored in this queueing waiting for service.

The challenge here is to consider the coexistence of several on/off sources with different sleep/wake parameters and to let the number of on/off sources increasing. Of course we impose the restriction that the sum of the average loads over all on/off sources does not exceed the capacity of the network server. In the case where the on/off sources were i.i.d., then each on/off source would independently offer a fraction of the total load. Thus the aggregated input flow would straightforwardly converge to a continuous Poisson flow. More generally when sources have very similar profiles, the asymptotic process would show renewal properties without long term dependences.

Each on/off source is used to model one user on the Net. The coexistence of several on/off sources with different profiles is a realistic way to model the variety of users on the web. Indeed, on one end we have the *surfer* users who browse the Net, continually downloading and discarding short pieces of information. The surfer will experience short wake and sleep periods resulting to high bandwidth demand in average. On the other end the *wise* user downloads much larger pieces of information (long wake periods) and spend more time in think state (longer sleep period), resulting to a smaller average bandwidth demand. It must be noted that in general the surfer user and the wise user will have both similar peak rates in wake periods, because of similar protocol and CPU.

2.2 The results

In the sequel we consider that we have an infinite denumerable set of on/off sources. We assume that the infinite sequence of individual average loads sums below the network-server capacity. We are not interested in a sequence with a geometric decay. We consider the case where the sequence of individual average loads does not decrease too fast. More precisely we consider a sequence of individual loads which polynomially decreases like in a Parreto sequence. Under this condition, we will prove *via* analytic methods the two following main results:

1. The overall arrival process has long term dependences.
2. The queue distribution has a tail at least polynomial.

The first result may provide a (very superficial) surprise, since the aggregation of independent sources with renewal properties should not produce long term correlations with such strong amplitude. The second result confirms statistics about buffer occupancy in network under real traffic conditions (Bellcore statistical data [6]).

It is not discussed how the Parreto profile would be the best descriptor of real profile distribution of users. Why users are different but not as different as are the coefficients in a geometric sequence? Maybe because people are all different, but not *that* different. However this kind of profile distribution seems to receive at least partial confirmation from experiment.

Long terms dependences in a network has a non negligible impact on network performance. As we mentionned in introduction, when traffic models are under pure continuous Poisson, the queue distributions in buffer has exponential tails. When this tail is known,

engineer can specify the minimum buffer size which fits a maximum overflow probability requirement. But if the tail turns out to be polynomial, then the designer would need to rise the previous minimum buffer size estimate to some several magnitude orders to fit the same overflow probability requirement.

2.3 The parameters

In the following we assume that the wake and sleep period durations in each single on/off source are exponentially distributed with respective parameter Poisson-wake and Poisson-sleep. This makes the transition between state “on” and state “off” a memoriless process. This is a simplifying assumption in order to keep the model tractable. We don’t assert that real life users *must* have memoriless behaviour.

We then consider a set of on/off sources such that the sequence of the ratio of Poisson-wake and Poisson-sleep decays like polynomials of degree β . We will show that under some conditions that the arrival process obtained by the aggregation of the on/off sources is β -long term dependent (section 2) and the queue length has a polynomial tail of degree at least β (section 3).

3 Long term dependences with aggregated on/off sources

3.1 Definition of the on/off sources

As we said above, an on/off process is a Poisson arrival process with a time varying rate. The Poisson rate oscillates between 0 and λ . The transitions times between these two values are exponentially distributed with respective values ν_1 and ν_0 . Quantity ν_1 is the transition rate from state “on” to state “off” and quantity ν_0 is the transition rate from state “off” to state “on”. In other words, ν_1 is the Poisson-wake parameters and ν_0 is the Poisson-sleep parameter. See [2, 1] for a detailed description.

3.2 Time covariance in a single on/off source

Let \mathcal{I}_θ be a time interval of size θ and for x non-negative let $\mathcal{I}_\theta + x$ the same interval translated of $\theta + x$. We define $A(\mathcal{I}_\theta)$ as the number of arrivals that occurs during interval \mathcal{I}_θ . We define the quantity $\text{Cov}(\theta, x)$ as the covariance of $(A(\mathcal{I}_\theta), A(\mathcal{I}_\theta + x))$, *i.e.* $\text{Cov}(\theta, x) = \text{E}(A(\mathcal{I}_\theta) \times A(\mathcal{I}_\theta + x)) - \text{E}(A(\mathcal{I}_\theta))\text{E}(A(\mathcal{I}_\theta + x))$. Our aim is to give a closed formula for $\text{Cov}(\theta, x)$.

To this end we introduce the following matrix \mathbf{T} that we call the *transition* matrix:

$$\mathbf{T} = \begin{bmatrix} -\nu_1 & \nu_0 \\ \nu_1 & -\nu_0 \end{bmatrix} \quad (2)$$

The eigenvalues of \mathbf{T} are 0 and $-(\nu_1 + \nu_0)$.

Let us suppose that at time 0, we start the process with a state randomly distributed between “on”s and “off”s according to a given probability vector \mathbf{u} . Via an easy Markov model we compute the probabilistic repartition between states “down” and states “off” of the arrival process at a further time x . The probability vector is exactly $\exp(x\mathbf{T})\mathbf{u}$. Notice that the eigenvalues of $\exp(x\mathbf{T})$ are the exponential of the eigenvalues of \mathbf{T} time x , *i.e.* 1 and $\exp(-(\nu_1 + \nu_0)x)$.

Going further, we establish that the expected proportion of “up”s periods and “down” periods intercepted by the interval $[0, \theta]$ will be $1/\theta \int_0^\theta \exp(x\mathbf{T})\mathbf{u}dx$. Notice that the eigenvalues of $\int_0^\theta \exp(x\mathbf{T})\mathbf{u}dx$ are respectively θ and $(1 - \exp(-(\nu_1 + \nu_0)\theta))/(\nu_1 + \nu_0)$. If we define ℓ as the covector $[\lambda, 0]$, then the average arrival on interval $[0, \theta]$ will be the scalar product $\langle \ell, \int_0^\theta \exp(x\mathbf{T})\mathbf{u}dx \rangle$.

Theorem 1 *We have the following explicit formula:*

$$\begin{aligned} \text{Cov}(\theta, x) &= \frac{[\lambda(1 - \exp(-(\nu_1 + \nu_0)\theta))]^2 \times}{\times \exp(-(\nu_1 + \nu_0)x)\nu_1\nu_0(\nu_1 + \nu_0)^{-4}} . \end{aligned} \quad (3)$$

Proof: If we fix θ , then the function $\text{Cov}(\theta, x)$ is a linear function of $\exp(x\mathbf{T})$ and therefore a linear combination of the two eigenvalues 1 and $\exp(-(\nu_1 + \nu_0)x)$. Identification of the coefficients easily comes. ■

By simple limit argument when $\theta \rightarrow 0$, we establish the following corollary.

Corollary 1 *Quantity $C(x)$ has expression:*

$$C(x) = \lambda^2 \frac{\nu_1\nu_0}{(\nu_1 + \nu_0)^2} \exp(-(\nu_1 + \nu_0)x) . \quad (4)$$

Figure 1 displays the function $C(x)$ obtained by simulation of a single on/off source with $\lambda = 2$, $\nu_1 = 0.09$ and $\nu_0 = 0.21$ (therefore average load is 0.6). The dramatic exponential decay is confirmed.

3.3 Aggregations of on/off arrival sources

In the following we consider a general system where the arrival process is a superposition of several independent on/off sources. We consider a denumerable set of on/off sources indexed from 1 to ∞ . See figure 2. Let \mathcal{I} be an interval, we call $A_j(\mathcal{I})$ the number of arrivals coming from the j th source during interval \mathcal{I} . Let $X(\mathcal{I})$ denotes the total arrival numbers queued during interval \mathcal{I} we obviously have $X(\mathcal{I}) = \sum_{j=1}^{\infty} A_j(\mathcal{I})$.

In the sequel we will show that such an aggregation can create an arrival process with long term dependences. An arrival process is β long term dependent if for any fixed interval \mathcal{I} , quantity $\text{Cov}(X(\mathcal{I}+x)X(\mathcal{I}))$ is of order $x^{-\beta}$, when x tends to infinity. We say the process is at most β -long term dependent when $\text{Cov}(X(\mathcal{I}+x)X(\mathcal{I}))$ is in magnitude order at most

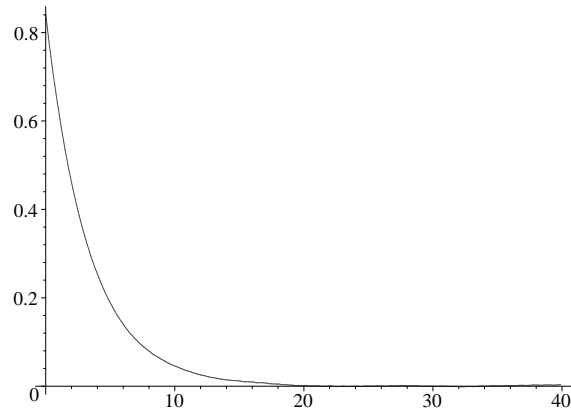


Figure 1: Simulated covariance function $C(x)$ for a single on/off source, $\lambda = 2$, mean load 0.6, $\nu_1 = 0.09$, $\nu_0 = 0.21$

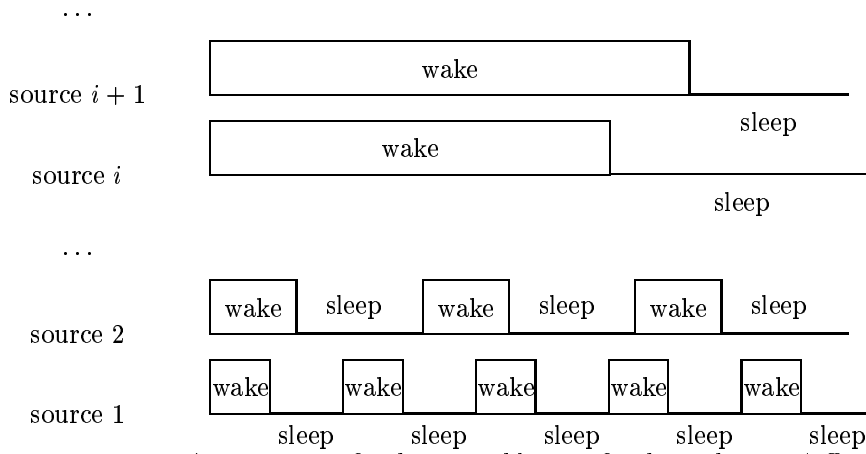


Figure 2: Aggregation of a denumerable set of independent on/off sources

$x^{-\beta}$. To simplify our computations we will just consider the asymptotics over function $C(x)$, *i.e.* when the size of interval \mathcal{I} tends to zero.

To simplify further the computations throughout this section we assume that:

- for every i , the i th on/off source has peak rate λ_i ;

- the λ_i are all identical and equal to a given $\lambda > 0$;
- for every on/off source we have $\nu_0 = \varepsilon_i^2$ and $\nu_1 = \varepsilon_i - \varepsilon_i^2$, for some sequence $\varepsilon_i > 0$.

Consequently, the average arrival per unit time is $\lambda \times \sum_{j=1}^{\infty} \varepsilon_j$.

We fix the Hultz parameter $\beta < 1$. In the sequel we suppose that the series $\sum \varepsilon_i$ is convergent. We call $\eta(s)$ the Dirichlet series $\sum_{j=1}^{\infty} \varepsilon_j^s$ and we assume that $\eta(s)$ is absolutely convergent for all complex numbers s with real part strictly less than $1 - \beta$. For example $\varepsilon_j = j^{1/(\beta-1)}$. Notice that in this case the Dirichlet series $\eta(s) = \zeta(s/(1 - \beta))$, *i.e.* can be identified with the Riemann *zeta* function. Figure 3 displays the ν_1 (off/on rates) and ν_0 (on/off rates) parameters of 60 such on/off sources with $\beta = 0.5$.

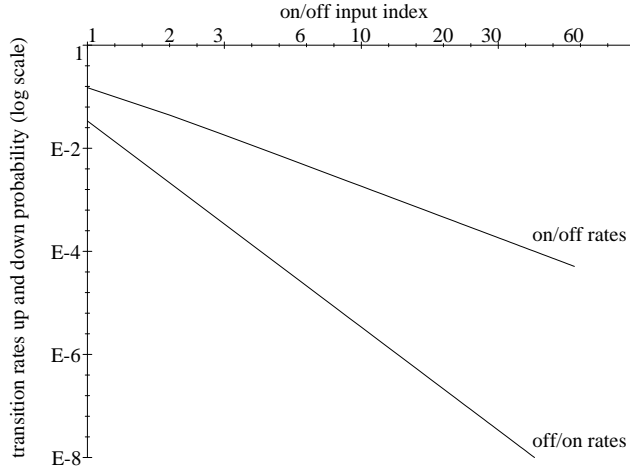


Figure 3: Transition rates of 60 on/off sources $\beta = 0.5$. Transition rates from on state to off state: ν_i (on/off rates); transition rate from off state to on state: δ_i (off/on rates)

Theorem 2 *If the serie $\eta(1 - \beta) = \sum \varepsilon_j^{1-\beta}$ converges, then the aggregated process is at most β long term dependent. Furthermore, if the Dirichlet series $\eta(s)$ has a simple pole on $s = 1 - \beta$ on a vertical strip, then the aggregated process is β long term dependent and the coefficients are exactly determined.*

Proof: Quantity $C_j(x)$ denotes the intensity covariance function of source number j . Due to the independence between on/off sources we immediately have $C(x) = \sum_{j=1}^{\infty} C_j(x)$. Referring to Theorem 1 we obtain:

$$C(x) = \lambda^2 \sum_{j=1}^{\infty} (1 - \varepsilon_j) \varepsilon_j \exp(-\varepsilon_j x) \quad (5)$$

We identify in the above a harmonic sum [5]. The asymptotics analysis of such a sum is classically obtained via Mellin transform. In particular it is important to determine the

definition domain and the singularity set of the Mellin transform $C^*(s)$, that is, defined for appropriate complex numbers s by:

$$C^*(s) = \int_0^\infty C(x)x^{s-1} dx . \tag{6}$$

From 5 we thus obtain (cf. [5]):

$$\begin{aligned} C^*(s) &= \lambda^2 \Gamma(s) \sum_{j=1}^\infty (1 - \varepsilon_j) \varepsilon_j^{1-s} \\ &= \lambda^2 \Gamma(s) (\eta(1-s) - \eta(2-s)) \end{aligned}$$

where $\Gamma(s)$ is the classical notation for the Mellin transform of function e^{-x} .

The expression of $C^*(s)$ converges for all s such that $\Gamma(s)$ and the series ε_j^{1-s} converges, i.e. for $0 < \Re(s) \leq \beta$.

The inverse Mellin transform, (cf. [5]):

$$C(x) = \frac{1}{2i\pi} \int_{c-i\infty}^{c+i\infty} C^*(s)x^{-s} ds , \tag{7}$$

holds for all c in the definition domain. Therefore $C(x) = O(x^{1-c})$ for all $c < 1 - \beta$, which proves the first point of the theorem.

Furthermore let us suppose that $\eta(s)$ has a single pole on $s = 1 - \beta$ with residue μ and can be continued on an extended vertical strip including $s = 1 - \beta$ where it can be bounded except on any compact neighborhood of $s = 1 - \beta$. Under this assumption we can move the line of integration in (7) to pass the singularity $s = \beta$ and to stand ε further on the right, as illustrated on figure 4. The residue theorem gives:

$$C(x) = \mu \lambda^2 \Gamma(\beta) x^{-\beta} + \int_{\beta+\varepsilon-i\infty}^{\beta+\varepsilon+i\infty} C^*(s)x^{-s} ds \tag{8}$$

The last integration on the left hand side is $O(x^{-\beta-\varepsilon})$ via obvious majoration under the sign “ \int ”. ■

Figure 5 displays function $C(x)$ obtained by simulation of the 60 on/off sources described in figure 3, $\beta = 0.5$, $\lambda = 2$ and the average load is 0.6. Since it is not obvious to depict that $C(x) = O(x^{-\beta})$ when $x \rightarrow \infty$, figure 6 displays function $C(x)x^\beta$ ($\beta = 0.5$) obtained from figure 5.

3.4 Playing with long term dependence asymptotics

By tuning system parameters one can give rise to *freak* source systems where $C(x) \asymp B(x)x^{-\beta}$ with $B(x)$ oscillating between two values: $\liminf B(x) \neq \limsup B(x)$.

For example, let us analyse the case $\lambda_j = 2^{j/2}\lambda$ and $\varepsilon_j = 2^{-j/(1-\beta)}$, with $\beta < 1$. It is probably very difficult to imagine practical on/off sources with such a drastic increase in

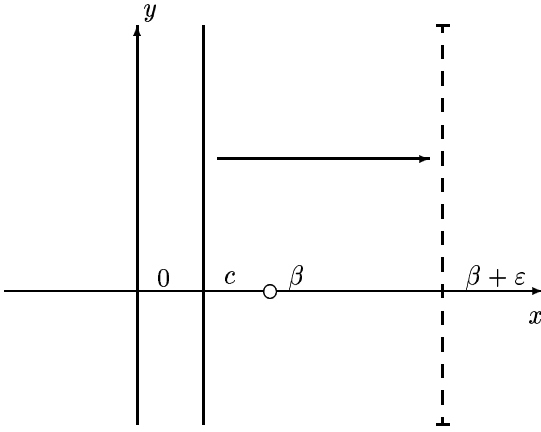
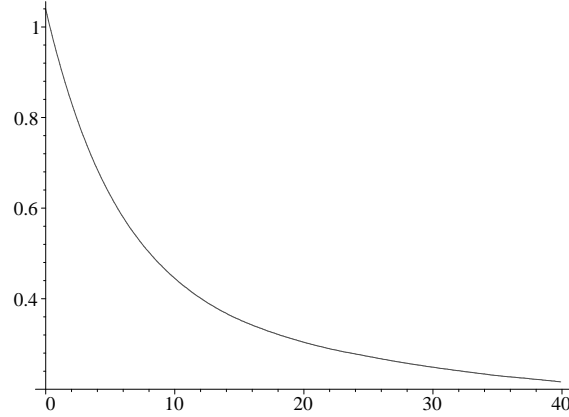


Figure 4: Move of the integration line in reverse Mellin transform

Figure 5: Simulated covariance function $C(x)$ for the aggregated 60 on/off sources, $\beta = 0.5$, $\lambda = 2$, mean load 0.6.

peak rates (notice that we still have the average load $\sum_j \lambda_j \varepsilon_j < \infty$), but we refer to it as a mathematical study on a limiting case.

By the way we obtain

$$\begin{aligned} C^*(s) &= \lambda^2 \Gamma(s) \left(\sum_{j=0}^{j=\infty} 2^j \varepsilon_j^{1-s} - 2^j \varepsilon_j^{2-s} \right) \\ &= \lambda^2 \Gamma(s) \left(\frac{1}{1 - 2^{\frac{s-\beta}{1-\beta}}} - \frac{1}{1 - 2^{\frac{s-\beta-1}{1-\beta}}} \right). \end{aligned}$$

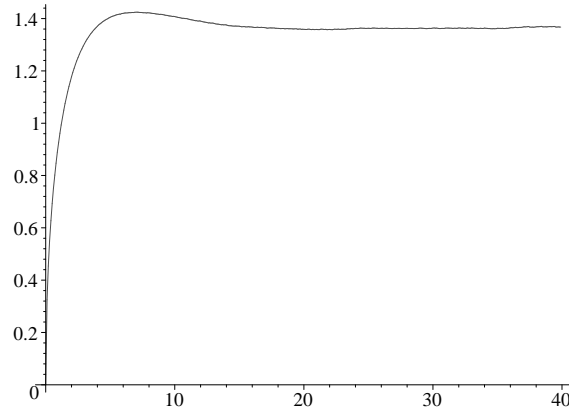


Figure 6: Simulated covariance function $C(x)x^\beta$ for the aggregated 60 on/off sources, $\beta = 0.5$ $\lambda = 2$, mean load 0.6.

We notice that $C^*(s)$ has a sequence of simple poles $s_k = \beta + \frac{2i(1-\beta)k\pi}{\log 2}$, for k integer, which creates a singularity set regularly spaced on the vertical axis $\Re(s) = \beta$ (see figure 7). Applying residu theorem on the Mellin inverse transform we obtain:

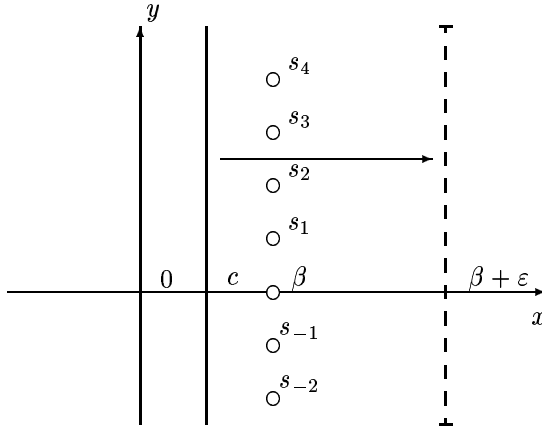


Figure 7: Reverse Mellin transform with multiple poles

$$C(x) = \lambda^2 \frac{1-\beta}{\log 2} \left(\sum_k \Gamma(s_k) \exp(-2ik\pi(\frac{1-\beta}{\log 2}) \log x) \right) x^{-\beta} + O(x^{-\beta-\epsilon}). \quad (9)$$

In other word we have proved that $C(x) \asymp B(\log x)x^{-\beta}$ where $B(\cdot)$ is a periodic function, of period $\log 2/(1-\beta)$, whose Fourier coefficients are proportional to $\Gamma(s_k)$. Since function

$B(\cdot)$ is not constant (indeed Fourier coefficients are all non zero), we have $\liminf B(x) \neq \limsup B(x)$.

Figure 8 displays function $C(x)x^\beta$ computed for $\beta = 0.5$. Quantity $C(x)x^\beta$ indefinitely bounces between its lower bound at 1.275 and its upper bound at 1.281.

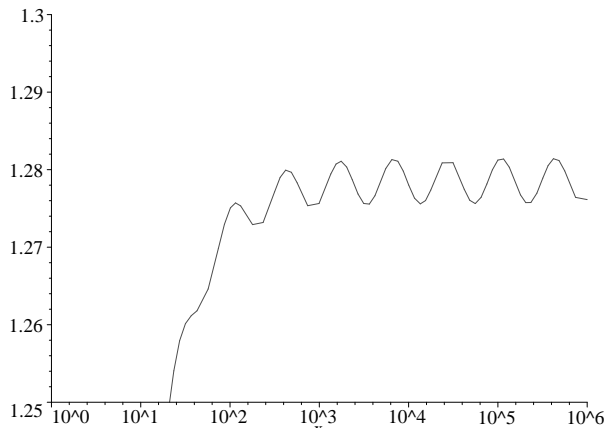


Figure 8: Computed covariance function $\check{C}(x)x^\beta$ for the fluctuating case with $\beta = 0.5$, logarithmic scale for x .

4 Queueing with on/off input sources

In the following, the network is modelled via a queue with a single server which receives input from a set of on/off sources. The service time is exponential with mean value of 1 time unit. Our aim is to find the steady state distribution of the queue length and in particular to find asymptotics of the probabilities p_n that the queue length exceeds n customers, when $n \rightarrow \infty$.

4.1 Queueing with a single on/off source

We refer to [2] to state the following theorem:

Theorem 3 *The queue length generating function with the exponential server and an on/off input source, with $\nu_1 = \varepsilon^2$ and $\nu_0 = \varepsilon - \varepsilon^2$ satisfies:*

$$q(z) = \frac{(1 - \lambda\varepsilon z_1)z - (1 - \lambda\varepsilon)z_1}{z - z_1} \quad (10)$$

with

$$z_1 = 1/2 \frac{1 + \lambda + \varepsilon + \sqrt{1 - 2\lambda + 2\varepsilon + \lambda^2 + 2\lambda\varepsilon + \varepsilon^2 - 4\lambda\varepsilon^2}}{\lambda + \lambda\varepsilon^2} . \quad (11)$$

Proof: This is a straightforward adaptation of [2]. ■

Corollary 2 *When $n \rightarrow \infty$, quantities p_n exponentially decrease with $p_n = \lambda \varepsilon z_1^{1-n}$.*

4.2 Queueing under an infinite number of on/off sources

We adopt the same hypotheses as in the section about traffic covariance, and we add a new condition:

- The quantity λ is strictly greater than 1.

Under this very condition it can be seen that each on/off stream in state “up” ask for an instantaneous workload larger than the server capacity. Nevertheless the global stability of the system is preserved because we keep $\lambda \sum_{j=1}^{\infty} \varepsilon_j < 1$. This condition is a kind of artefact but it simplifies computations.

In the following we say that a non-negative random variable X has a *polynomial tail of degree $-\beta$* if the probability of the event $X > x$ is of order $x^{-\beta}$ when $x \rightarrow \infty$. We say that the random variable has a polynomial tail of degree at least $-\beta$ if it exists α such that the probability that $X > x$ tends to be greater than $\alpha x^{-\beta}$ for some $\alpha > 0$ when $x \rightarrow \infty$.

Theorem 4 *If the Dirichlet series $\sum_i \varepsilon_i^s$ has a simple pole on $s = 1 - \beta$ with residue μ , then the queue length has a polynomial tail of degree at least $-\beta$ and coefficients for a lower bound can be exactly determined.*

Proof: To simplify, we assume the same hypotheses as in the proof of Theorem 2: the Dirichlet series can be continued to a vertical strip $1 - \beta < \Re(s) < 1 - \beta + \varepsilon$.

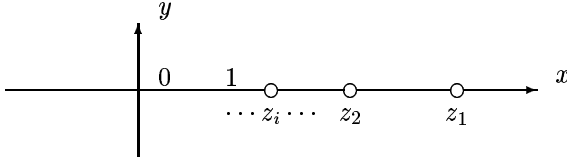
Instead of directly looking on the queue length in the single server mode we consider a modified system called “multi-server” system. In the multi-server system we assume that every on/off stream has a separate queue and a distinct server. It is clear that the sum of the queue length in the multi-server system is always smaller than the queue length of the single server system. Therefore if the multi-server system shows a polynomial tail of degree $-\beta$ then the single server system will show a polynomial tail of degree at least $-\beta$.

Our aim is to show that the multi-server model has a polynomial tail of degree $-\beta$. Let $g(z)$ the probability generating function of the sum of the queue lengths of the multi-server mode. We have the obvious identity: $g(z) = \prod_{i=1}^n q_i(z)$, where $q_i(z)$ is the probability generating function of the queue size of the server attached to the on/off stream number i .

Generating functions $q_i(z)$ are computed according to Theorem 3. Let z_i be the equivalent of z_1 in theorem 3 for on/off source i . We have formally:

$$g(z) = \prod_{i=1}^{\infty} \frac{(1 - \lambda \varepsilon_i z_i)z - (1 - \lambda \varepsilon)z_i}{z - z_i} \quad (12)$$

Quantity $g(z)$ is defined as long as z does not belong to the set of singularities z_i . Under this restriction, the expression $g(z)$ converges because the series in $z_i - 1$ converges. Indeed,

Figure 9: Locations of the poles z_i in the complex plan

$z_n = 1 + \varepsilon_n/(\lambda - 1) + O(\varepsilon_n^2)$ when $n \rightarrow \infty$, and the series in ε_i converges (see figure 9). In passing $g(z) = a + O(1/z)$ with $a = \prod_{i=1}^{\infty} (1 - \lambda \varepsilon_i z_i)$, when $|z| \rightarrow \infty$.

In the following we will establish that $g(1-x) = 1 + \alpha x^\beta + O(x^{\beta+\varepsilon})$ when x converges to 0 with the condition that $|\arg(z)| < \pi - \theta$ for some α and $\theta < \pi/2$. We define $\ell(z) = \log(g(1-x)) - \log(a)$. We have

$$\ell(x) = \sum_{i=1}^{\infty} \left[\log\left(1 + \frac{1 - \lambda \varepsilon_i z_i}{z_i - 1} x\right) - \log\left(1 + \frac{x}{z_i - 1}\right) \right] - \log(a). \quad (13)$$

Since $\ell(x) = -\log(a) + o(1)$ when $x \rightarrow 0$ and that $\ell(x) = O(1/x)$ when $x \rightarrow \infty$, the Mellin transform $\ell^*(s)$ of $\ell(x)$ is defined on the strip $0 < \Re(s) < 1$ and has expression (cf. [5]):

$$\sum_{i=1}^{\infty} \left[\left(\frac{1 - \lambda \varepsilon_i z_i}{z_i - 1}\right)^{-s} - \left(\frac{1}{z_i - 1}\right)^{-s} \right] \frac{-\pi}{s \sin \pi s}. \quad (14)$$

We identify in $(s \sin \pi s)^{-1}(-\pi)$ the Mellin transform of function $\log(1+x)$.

To get asymptotics of $\ell(x)$ when $x \rightarrow 0$ we identify the singularities of $\ell^*(s)$ located on the left of the definition domain, i.e for negative $\Re(s)$ in order to use the residu theorem in the reverse Mellin transform.

The first encountered pole is at $s = 0$ since $(s \sin \pi s)^{-1} \pi$ has double root at $s = 0$. This pole is finally of degree 1 since the factor in front of it has also root at $s = 0$, its residue is $\sum_{i=1}^{\infty} \log(1 - \lambda \varepsilon_i z_i) = \log(a)$. There is a less obvious pole on $s = -\beta$. Indeed:

$$\sum_{i=1}^{\infty} \left[\left(\frac{1 - \lambda \varepsilon_i z_i}{z_i - 1}\right)^{-s} - \left(\frac{1}{z_i - 1}\right)^{-s} \right] = \sum_{i=1}^{\infty} s \varepsilon_i^{1+s} (1 - \lambda)^{-s} (1 + O(\varepsilon_i)) \quad (15)$$

since $z_i = 1 + \varepsilon_i/(\lambda - 1) + O(\varepsilon_i^2)$. Therefore the pole at $s = -\beta$ has residu $\mu \beta (1 - \lambda)^\beta$.

Applying the residu theorem in the reverse Mellin transform by moving the line of integration from c to $-\beta - \varepsilon$:

$$\begin{aligned} \ell(x) &= \frac{1}{2i\pi} \int_{c-i\infty}^{c+i\infty} \ell^*(s) x^{-s} ds \\ &= -\log(a) - \mu \beta (1 - \lambda)^\beta \pi (\beta \sin \pi \beta)^{-1} x^\beta + \\ &\quad + \frac{1}{2i\pi} \int_{-\beta-\varepsilon-i\infty}^{-\beta-\varepsilon+i\infty} \ell^*(s) x^{-s} ds \end{aligned}$$

A trivial majorization under the above integral, translates immediately into the correct asymptotics of $g(z)$ around $z = 1$: $g(1 - x) = 1 - \alpha x^\beta + O(x^{\beta+\varepsilon})$ with $\alpha = a\mu\pi\beta(1 - \lambda)^\beta(\beta \sin \pi\beta)^{-1}$. Of course, this asymptotic is valid under the condition that x varies in an open set such that $|\arg(x)| < \pi - \theta$ for a θ arbitrarily chosen $\theta < \pi/2$.

The latter conditions suffice to translate asymptotics of $g(z)$ into asymptotics over its coefficients. Let b_n be the coefficients of $g(z)$: $\sum_{n=0}^{\infty} b_n z^n = g(z)$. Since $g(1 - x) = 1 - \alpha x^\beta + O(x^{\beta+\varepsilon})$. We immediately deduce from the singularity analysis of [4] that

$$b_n = \frac{\alpha}{\Gamma(-\beta)} n^{-1-\beta} + O(n^{-1-\beta-\varepsilon}) . \tag{16}$$

The cumulative coefficients $p_n = b_n + b_{n+1} + \dots$ observe similar asymptotics:

$$p_n = \frac{\alpha}{\Gamma(1-\beta)} n^{-\beta} + O(n^{-\beta-\varepsilon}) . \tag{17}$$

The above expression proves the theorem. ■

Figure 10, displays the distribution of queue size. The results are obtained by sampling 20,000 windows of 10 time unit size. If we select a logarithmic scale for horizontal axis, then the curve $\log(p_n)$ versus $\log n$ is asymptotically linear as a kind of illustration of a polynomial tail distribution.

Remark: We can extend this result in a similar manner as for the covariance analysis, *i.e.* making the singularity set a little more complicated. In principle, it is possible to obtain tail distributions with oscillating terms. We can also treat more general cases. The artefact condition $\lambda > 1$ can also be eliminated. But the most promising workitem left in the present study is the estimation of a suitable *upper* bound for the queue length tail distribution, since the present work is limited to lower bounds.

5 Simulation Results

We analyse the traffic performance of the aggregation of an infinite sequence of independent on/off sources. In fact we use an aggregation of 60 on/off sources. The parameters of the on/off input with index i are: $\lambda_i = 2.0$, $\varepsilon_i = \alpha i^{1/(\beta-1)}$; we choose $\beta = 0.5$ and β is tuned in order to keep the average input at 0.6. Figure 3 displays the sequence of on/off parameters on a logarithmic scale.

Figures 11, 12 and 13 respectively shows the histograms of arrivals and buffer size for time scale varying from 1 time unit, 100 time units and 10,000 time units. We notice interesting self-similarities between congestion periods.

6 Conclusion

We have presented a simple model of arrivals based on the aggregation of independent memoriless on/off sources. We show that appropriate tuning of parameters produces long

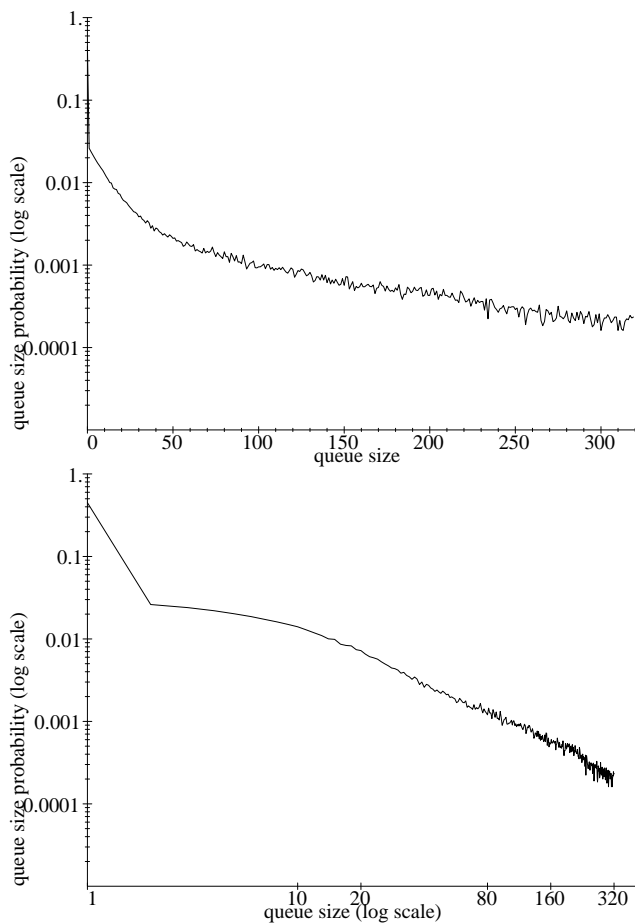


Figure 10: Simulated queue size distribution (bottom: logarithmic scale), on/off aggregation

term dependences and polynomial tail distributions of buffer occupancy. Such model may give rise to a new source of long term dependence in Web traffics. Future work can progress in two different directions. The first direction is to find less restrictive conditions on parameter tuning. The second direction is to find an exact estimate of the asymptotic tail distribution instead of a straightforward lower bound.

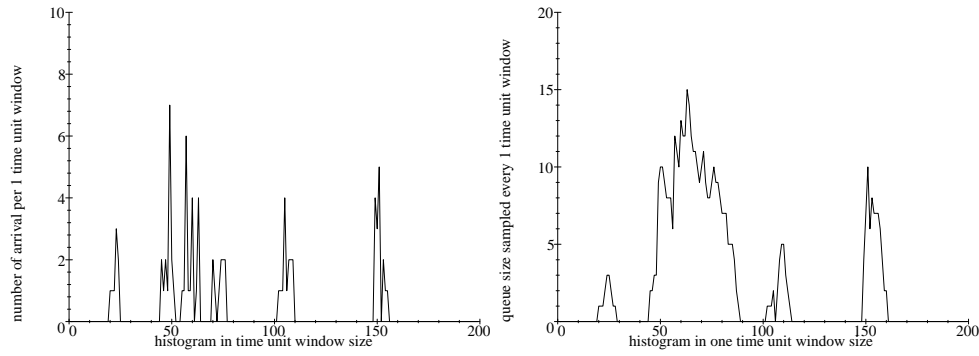


Figure 11: On/off aggregation arrival, histogram of arrivals (left), and queue sizes (right), sampled in one time unit window size

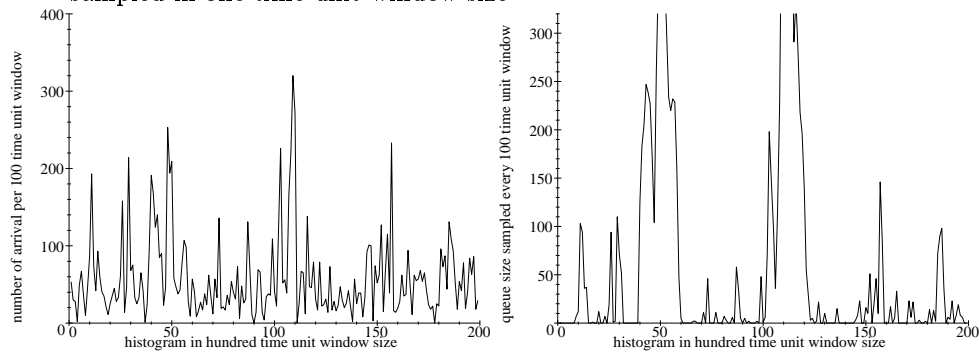


Figure 12: On/off aggregation arrival, histogram of arrivals (left), and queue sizes (right), sampled in one hundred time unit window size

References

- [1] P. Jacquet, Analytic information theory in service of queueing with aggregated exponential on/off arrivals, in *proc. 35th annual Allerton conf. on Communication, Control and Computing*, pp. 242-251, 1997.
- [2] M. Neuts, *Matrix-geometric solutions in stochastic models : an algorithmic approach*, the John Hopkins University press, 1981.
- [3] N. Likhanov, B. Tsybakov, N. Georganas, Analysis of an ATM buffer with self-similar (“fractal”) input, submitted, 1994.
- [4] P. Flajolet, A. Odlyzko, Singularity analysis of generating functions, in *SIAM J. Disc. Math.*, Vol 3, No 2, pp. 216-240, 1990.

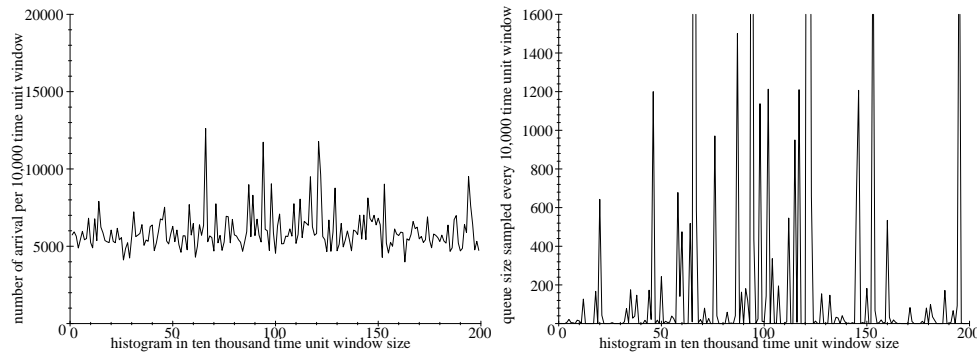


Figure 13: On/off aggregation arrival, histogram of arrivals (left), and queue sizes (right), sampled in ten thousand time unit window size

- [5] P. Flajolet, X. Gourdon, P. Dumas, Mellin transform and asymptotics: Harmonic sums, in *Theoretical Computer Science*, Vol 144, No 1-2, pp. 3-58, 1995.
- [6] W. Willinger, M. Taqqu, W. Leland, D. Wilson, Self-similarity in high-speed packet traffic: analysis and modeling of Ethernet traffic measurement, in *Statistical Science*, Vol 10, No 1, pp. 67-85, 1995.
- [7] P. Jacquet, W. Szpankowski, Toward analytical information theory: entropy computations, submitted, 1997.
- [8] P. Glynn, W. Whitt, Logarithmic asymptotics for steady-state tail probabilities in a single-server queue, *J. Appl. Prob.*, to appear.
- [9] N. Duffield, N. O'Connell. Large deviations and overflow probabilities for the general single server queue, with applications. DIAS Technical Report No DIAS-STP-93-30, 1993.
- [10] W. Szpankowski, On Certain Recurrences Arising in Universal Coding, preprint.
- [11] W. Szpankowski, Techniques of the Average Case Analysis of Algorithms, in *Handbook on Algorithms and Theory of Computation* (Ed. M. Atallah), CRC 1997.



Unit e de recherche INRIA Lorraine, Technop le de Nancy-Brabois, Campus scientifique,
615 rue du Jardin Botanique, BP 101, 54600 VILLERS L ES NANCY
Unit e de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex
Unit e de recherche INRIA Rh ne-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN
Unit e de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex
Unit e de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

 diteur
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399