# INRIA

# *Investigating self-similarity and heavy tailed distributions on a large scale experimental facility*

Patrick Loiseau — Paulo Gonçalves — Pascale Primet Vicat-Blanc — Pierre Borgnat — Patrice Abry — Guillaume Dewaele

**N° ????**

March 2008

Thème NUM

*R*apport
*de recherche*

# Investigating self-similarity and heavy tailed distributions on a large scale experimental facility

Patrick Loiseau*, Paulo Gonçalves* , Pascale Primet Vicat-Blanc* , Pierre Borgnat† , Patrice Abry† , Guillaume Dewaele‡

**Abstract:** After seminal work by Taqqu et al. relating self-similarity to heavy tail distributions, a number of research articles verified that aggregated Internet traffic time series show self-similarity and that Internet attributes, like WEB file sizes and flow lengths, were heavy tailed. However, the validation of the theoretical prediction relating self-similarity and heavy tails remains unsatisfactorily addressed, being investigated either using numerical or network simulations, or from uncontrolled web traffic data. Notably, this prediction has never been conclusively verified on real networks using controlled and stationary scenarii, prescribing specific heavy-tail distributions, and estimating confidence intervals. In the present work, we use the potential and facilities offered by the large-scale, deeply reconfigurable and fully controllable experimental Grid5000 instrument, to investigate the prediction observability on real networks. To this end we organize a large number of controlled traffic circulation sessions on a nation-wide real network involving two hundred independent hosts. We use a FPGA-based measurement system, to collect the corresponding traffic at packet level. We then estimate both the self-similarity exponent of the aggregated time series and the heavy-tail index of flow size distributions, independently. Comparison of these two estimated parameters, enables us to discuss the practical applicability conditions of the theoretical prediction.

**Key-words:** Computer networks, Grid5000, metrology, self-similarity, heavy-tail distributions

\* INRIA RESO, ENS Lyon, Université de Lyon.
† SiSyPhe, CNRS, ENS Lyon, Université de Lyon.
‡ SiSyPhe, ENS Lyon, Université de Lyon.

# Vérification du lien entre auto-similarité et distributions à queues lourdes sur un dispositif grande échelle

**Résumé :** À la suite du travail théorique de Taqqu et de ses collaborateurs, reliant l'auto-similarité aux distributions à queues lourdes, quantité d'articles de recherche ont vérifié que les séries temporelles de trafic internet présentent en effet un caractère auto-similaire, et qu'en effet aussi, certaines variables d'internet, telle que par exemple les tailles de flux, étaient à queue lourde. Cependant, la validation de cette prédiction théorique liant auto-similarité et distributions à queues lourdes, reste peu satisfaisante dans la mesure où elle n'a été expérimentalement vérifiée que sur des simulateurs numériques de réseaux, ou sur des données de trafic réel dont on ne maîtrise aucun des paramètres. En particulier, cette relation n'a jamais été formellement validée sur des réseaux réels en situation contrôlée de scénarios stationnaires, dans lesquels des distributions à queues lourdes spécifiques sont prescrites, et des intervalles de confiances estimés. Dans ce travail, nous exploitons le potentiel et les capacités offertes par Grid5000, une plate-forme à grande échelle, profondément reconfigurable et totalement controlée, pour confronter cette prédiction théorique au contexte d'un véritable réseau. Pour ce faire, nous avons procédé à un grand nombre d'expériences *in situ*, où nous avons généré entre deux cents nœuds indépendants, différents profils d'un trafic entièrement contrôlé. Pour collecter les données correspondantes, nous utilisons un système à base de FPGA capable de traiter des flux de 1Gb/s avec une granularité à l'échelle du paquet. À partir de ces données, nous estimons indépendamment l'exposant d'auto-similarité du débit aggrégé et l'indice de queue lourde des distributions de taille de flux. La mise en correspondance de ces deux estimations nous permet alors de définir en pratique, les contours d'application du théorème.

**Mots-clés :** Réseaux informatiques, Grid5000, métrologie, auto-similarité, distributions à queues lourdes

# 1   Motivations

Comprehension and prediction of the network traffic is a constant and central preoccupation for internet service providers. Challenging questions, such as the optimization of network resource utilization that respect the application constrains, the detection (and ideally the anticipation) of anomalies and congestion, contribute to guarantee a better quality of service (QoS) to users. From a statistical viewpoint, this is a challenging and arduous problem that encompasses several components: network design, control mechanims, transport protocols and the nature of traffic itself. In the last decade, great attention has been devoted to the statistical study of time series and random variables, which collected at the core of networks, are valuable fingerprints of the system state and of its evolution. With this in mind, the pioneering work by [**?**] and [**?**] evidenced that the Poisson hypothesis, a relevant and broadly used model for phone networks, failed at describing computer network traffic. Instead, self-similarity was shown a much more appropriate paradigm, and since then, many authors have reported its existence in a wide variety of traffics [**?**, **?**, **?**, **?**]. Following up this prominent discovery, the theoretical work by Taqqu and collaborators constituted another major breakthrough in computer network traffic modeling, identifying a plausible origin of self-similarity in traffic time series [**?**, **?**, **?**]. It is stated that the heavy-tail nature of some probability distributions, mainly that of flow size distributions, suffice to generate traffic exhibiting long range dependence, a particular manifestation of self-similarity [**?**]. To support their claim, they established a close form relation connecting the heavy tail thickness (as measured by a tail index) and the self-similarity exponent.

Notwithstanding its mathematical soundness, pragmatic validity of this model has been corroborated with real world traffic data only partially, so far. First pitfall lies in the definition of long range dependence itself, which, as we will see, is a scale invariance property that holds only asymptotically for long observation durations. Its consistent measurement requires that experimental conditions maintain constant, and that no external activity perturbs the traffic characteristics. In those conditions, finding a scale range that limits itself to stationary data, and that is sufficiently wide to endorse reliable self-similarity measurements, is an intricate task.

Secondly, even though real traffic traces had led to check concordance between tail index and self-similarity exponent, only was it perceived for a given network configuration that necessarily corresponded to a single particular value of the parameters set. An extensive test, to verify that self-similarity exponent obeys the same rule when the tail index is forced to range over some interval of interest, was never performed on a large scale real network plate-form.

Finally, the exact role of the exchange protocol, viewed as a subsidiary factor from this particular model, is still controversial [**?**, **?**, **?**]. Due to the lack of flexible, versatile, while realistic experimental environments, part of this metrology questioning has been addressed by researchers of the network community, using simulators, emulators or production platforms. However, these tools have limitations on their own, which turn difficult the studies, and yield only incomplete results.

In the present work, we use the potential and the facilities offered by the very large-scale, deeply reconfigurable and fully controlable experimental Grid5000 instrument to empirically investigate the scope of applicability of Theorem pro-

posed by Taqqu et al. [**?**, **?**, **?**]. Under controlled experimental conditions, we first prescribe the flow size distribution to different tail indices and compare the measured traffic self-similar exponents with their corresponding theoretical predictions. Then, we elucidate the role of the protocol and of the rate control mechanism on traffic scaling properties. In the course, we resort to efficient estimators of the heavy-tail index and of the self-similarity exponent derived from recent advances in wavelet based statistics and time series analysis.

The sequel is organized as follows. Section **??** summarizes related works. Section **??** elaborates on theoretical foundations of the present work, including a concise definition of parameters of interest. In section **??** we develop the specifities of our experimental testbed, and we describe our experimental designs. Section **??** presents and comments the results. Conclusions and perspectives are itemized in section **??**.

## 2   Related Work

Without giving full bibliography on the subject (many can be found in [**?**, **?**, **?**]), there have been extensive reports on self-similarity in network traffic. As most of them are based on measurements and on analysis of real-world traces from the Internet, they only permit experimental validation of a single point on the curve, corresponding to one particular configuration. As its was presented before, the question here is more on the relation between these two properties, which is rooted in the seminal work by [**?**, **?**] about the M/G/N queueing models with heavy-tail distributions of ON periods. Nonetheless, first experimental works by Crovella and co-authors [**?**, **?**], hinted that this theoretical relation holds for internet traffic, and later on, also for more general types of traffic [**?**, **?**]. However, due to the impossibility of controlling important parameters when monitoring the Internet, only compatibility of the formula could be tested against real data , but there is no statistically grounded evidences that self-similarity measured in network traffic is the work of this sole equality. On the other hand, study of self-similarity at large scales is very sensitive to inevitable non-stationnarities (day and week periodicities for instance) and to fortuitous anomalies existing on the Internet (see for instance [**?**]). It seems that the question has, since, never received a full experimental validation. In order to obtain such a validation, an important feature is to be able to make the heavy tail index vary, and there is only few attempts to validate the relation under these conditions. One is conducted in [**?**], that uses a network simulator, and where some departure from the theoretical prediction is reported (Fig. 3 in this article). This deviation is probably caused by the limited length of the simulation and also by the bias introduced by the used scaling estimator (R/S and Variance Time) on short traces. Actually, the main restriction of simulators lies in their scalability limitation, and in the difficulty of their validation. Indeed, the network is an abstraction, protocols are not production code, and the number of traffic sources or bitrates you can simulate depends on the computing power of the machine. Large-scale experimental facilities are alternatives that may overcome both Internet and simulators limitations as they permit to control network parameters and traffic generation, including statistics and stationarity issues.
Emulab [**?**] is a network experimental facility where network protocols and services are run in a fully controlled and centralized environment. The emulation

software runs on a cluster where nodes can be configured to emulate network links. In an Emulab experiment, the user specifies an arbitrary network topology, having a controllable, predictable, and reproducible environment. He has full root access on PC nodes, and he can run the operating system of his choice. However, the core network's equipments and links are emulated. The RON testbed [**?**] consists of about 40 machines scattered around the Internet. These nodes are used for measurement studies and evaluation of distributed systems. RON does not offer any reconfiguration capability at the network or at the nodes' level. The PlanetLab testbed [**?**] consists of about 800 PCs on 400 sites (every site runs 2 PCs) connected to the Internet (no specific or dedicated link). PlanetLab allows researchers to run experiments under real-world conditions, and at a very large scale. Research groups are able to request a PlanetLab slice (virtual machine) in which they can run their own experiment .

Grid5000, the experimental facility we use in the present work, proposes a different approach where the geographically distributed resources (large clusters connected by ultra high end optical networks) are running actual pieces of software in a real wide area environment. Grid5000 proposes a complimentary approach to PlanetLab, both in terms of resources and of experimental environment. Grid5000 allows reproducing experimental conditions, including network trafic and CPU usage. This feature warrants that evaluations and comparisons are conducted according to a strict and scientific method.

## 3 Theory

Taqqu's Theorem relates two statistical properties that are ubiquitously observed in computer networks: On the one hand, self-similarity that is defined at the level of aggregated time-series of the traffic, and on the other hand, heavy-tailness that involves grouping of packets (such as TCP connections). Simplistically, network traffic is described as a superposition of flows (without notions of users, or sessions,...) that permits us to adopt the following simple two-level model: *(i)* Packets are emitted and grouped in flows whose length (or number of packets) follows a heavy tailed distributed random variable [**?**, **?**, **?**]; *(ii)* the sum over those flows approximates network traffic on a link or a router. This crude description is coherent with current (yet more elaborate) statistical model for Internet traffic [**?**, **?**].

After a succinct definition of these two statistical properties, we present the corresponding parameter estimation procedures that we use in our simulations, and chosen amongst those reckoned to present excellent estimation performance.

### 3.1 Self-similarity and long range dependence

#### 3.1.1 Definition

Taqqu's Theorem implies that Internet time series are relevantly modeled by fractional Brownian motion (fBm), the most prominent member of a class of stochastic processes, referred to as *self-similar processes with stationary increments* ($H$-sssi, in short). Process $X$ is said to be $H$-sssi if and only if its satisfies [**?**]:

$$X(t) - X(0) \overset{fdd}{=} X(u+t) - X(u), \forall t, u \in \mathbb{R}, \tag{1}$$

$$X(t) \overset{fdd}{=} a^H X\left(\frac{t}{a}\right), \forall t, a > 0, \ 0 < H < 1, \tag{2}$$

where $\overset{fdd}{=}$ means equality for all finite dimensional distributions. Eq. (**??**) indicates that the increments of $X$ form stationary processes (while $X$ itself is not stationary). Essentially, self-similarity, Eq. (**??**), means that no characteristic scale of time can be identified as playing a specific role in the analysis or description of $X$. Corollarily, Eq. (**??**) implies that $\mathbb{E}X(t)^2 = \mathbb{E}X(1)^2 t^{2H}$, underlining both the scale free and the non stationary natures of the process.

It turns out that the covariance function of the increment process, $Y(t) = X(t+1) - X(t)$, of a $H$-sssi process $X$ satisfies, for $|\tau| \to +\infty$:

$$\mathbb{E}Y(t)Y(t+\tau) \sim \mathbb{E}X(1)^2 H(2H-1)|\tau|^{2H-2}, \tag{3}$$

When $1/2 < H < 1$, hence $0 < 2 - 2H < 1$, such a power law decay of the covariance function of a stationary process is referred to as long range dependence [**?**, **?**].

Long range dependence and self-similarity designate two different notions, albeit often confused. The latter is associated to non stationary time series, such as fBms, while the former is related to stationary time series, such as fBm's increments. In the present work, given that Taqqu's Theorem predicts that the cumulated sums of aggregated Internet time series are self-similar, we adopt here the same angle and discuss the results in terms of self-similarity of the integrated traces.

### 3.1.2 Self similarity parameter estimation

In [**?**], it was shown that wavelet transforms provide a relevant procedure for the estimation of the self-similarity parameter. This procedure revealed particularly efficient at analyzing Internet time series in [**?**, **?**] and has then been massively used in this context.

Let $d_X(j,k) = \langle \psi_{j,k}, X \rangle$ denote the (Discrete) Wavelet Transform coefficients, where the collection $\{\psi_{j,k}(t) = 2^{-j/2}\psi_0(2^{-j}t - k), k \in \mathbb{Z}, j \in \mathbb{Z}\}$ forms a basis of $L^2(\mathbb{R})$ [**?**]. The reference template $\psi_0$ is termed mother-wavelet and is characterized by its number of vanishing moments $N_\psi > 1$, an integer such that $\int t^k \psi_0(t)\, dt \equiv 0, \forall k = 0, ..., N_\psi - 1$. Then, decomposing a $H$-sssi process, the variance of the wavelet coefficients verifies [**?**]:

$$\mathbb{E}|d_X(j,k)|^2 = \mathbb{E}|d_X(0,0)|^2 2^{j(2H+1)}, \tag{4}$$

and, provided $N > H + 1/2$, the sequence $\{d_X(j,k), k = \ldots, -1, 0, 1, \ldots\}$ form a stationary and weakly correlated time series [**?**]. These two central properties warrant to use the empirical mean $S(j) = n_j^{-1}\sum_k |d_X(j,k)|^2$, ($n_j$ being the number of available coefficients at scale $2^j$) to estimate the ensemble average $\mathbb{E}|d_X(j,k)|^2$. Eq. (**??**) indicates that self-similarity transposes to a linear behavior of $\log_2 S(j)$ vs. $\log_2 2^j = j$ plots, often referred to as Logscale Diagrams (LD) in the literature [**?**, **?**]. A (weighted) linear regression of the LD within a proper range of octaves $j_1, j_2$ is used to estimate $H$.

In [**?**, **?**, **?**], the estimators performance are both theoretically and practically quantified, and are proved to compare satisfactorily against the best parametric techniques. Moreover, this estimator is endowed with a practical robustness that comes from its extra degree of freedom $N_\psi$. Its main use issue lies in the correct choice of the regression range $j_1 \leq j \leq j_2$. This will be discussed in Section **??**, in the light of actual measurements.

## 3.2 Heavy Tail

### 3.2.1 Definition

A (positive) random variable $\mathbf{w}$ is said to be heavy tailed, with tail exponent $\alpha > 0$ (and noted $\alpha$-HT) when the tail of its cumulative distribution function, $F_{\mathbf{w}}$, is characterized by an algebraic decrease [**?**]:

$$P(\mathbf{w} > w) = 1 - F_{\mathbf{w}}(w) \sim cw^{-\alpha} \quad \text{for } w \to \infty. \tag{5}$$

A $\alpha$-HT random variable $\mathbf{w}$ has finite moments up to order $\alpha$. For instance, when $1 < \alpha < 2$, $\mathbf{w}$ has finite mean but infinite variance. A paradigm for $\alpha$-HT positive random variable is given by the Pareto distribution:

$$F_{\mathbf{w}}(w) = 1 - \left(\frac{k}{w+k}\right)^\alpha, \tag{6}$$

with $k > 0$ and $\alpha > 1$. Its mean reads: $\mathbb{E}\mathbf{w} = k/(\alpha - 1)$.

### 3.2.2 Tail exponent estimation

Estimation of the tail exponent $\alpha$ for $\alpha$-HT random variables is an intricate issue that received considerable theoretical attention in the statistics literature: measuring the tail exponent of a HT-distribution amounts to evaluate from observations, how fast does the probability of rare events decrease in Eq. (**??**). Once random variables are know to be drawn from an a priori distribution, such as the Pareto form (**??**) for example, parametric estimators exist and yield accurate estimates of the tail index $\alpha$ (see e.g. [**?**]). However, if the actual distribution of observations does not match the a priori expected $\alpha$-HT model, parametric estimators eloquently fail at measuring the tail decay.

For this reason, the non-parametric empirical estimator of $\alpha$ proposed in [**?**] will be preferred. The principle of this estimator is simple and relies on the Fourier mapping between the cumulative distribution function $F_{\mathbf{w}}(w)$ and the characteristic function $\chi_{\mathbf{w}}(s)$ of a random variable:

$$\chi_{\mathbf{w}}(s) = \int e^{-isw} \, dF_{\mathbf{w}}(w). \tag{7}$$

By a duality argument, the tail exponent $\alpha$ that bounds the order of finite moments of $F_{\mathbf{w}}$,

$$\alpha = \sup_r \{r > 0 : \int |w|^r \, dF_{\mathbf{w}}(w) < \infty\}, \tag{8}$$

transposes to the local Lipschitz regularity of the characteristic function $\chi_{\mathbf{w}}$ at the origin, according to:

$$\alpha = \sup_r \{r > 0 : 1 - \Re\chi_{\mathbf{w}}(s) = \mathcal{O}(s^r) \text{ as } s \to 0^+\}, \tag{9}$$

where $\Re$ stands for the real part. It is easy to recognize in this power law behavior of $\Re\chi_{\mathbf{w}}(s)$, a scale invariance property of the same type of that of relation (**??**), which is conveniently identifiable with wavelet analyses. Hence, computing the discrete wavelet decomposition of $\Re\chi_{\mathbf{w}}$, and retaining only the wavelet coefficients that lie at the origin $k = 0$, yields the following multiresolution quantity:

$$d_{\chi_{\mathbf{w}}}(j,0) = \mathbb{E}\Psi_0\left(2^j\mathbf{w}\right) \le C2^{j\alpha} \text{ for } j \to -\infty, \tag{10}$$

where $\Psi_0(\cdot)$ denotes the Fourier transform of analyzing wavelet $\psi_0(\cdot)$. Now, let $\{w_0, \cdots, w_{n-1}\}$ be a set of i.i.d. $\alpha$-HT random variables, and replace the ensemble average in Eq. (**??**) by its empirical estimator, the estimate $\widehat{\alpha}$ simply results from a linear regression of the form

$$
\begin{aligned}
\log \widehat{d}_{\chi_{\mathbf{w}}}^{(n)}(j,0) &= \log n^{-1} \sum_{i=o}^{n-1} \Psi(2^j w_i) \\
&\approx \widehat{\alpha} j + \log C, \text{ as } j \to -\infty.
\end{aligned}
\tag{11}
$$

The estimator was proven to converge for all heavy tail distributions, and also it has a reduced variance of estimation in $\mathcal{O}(n^{-1})$, where $n$ is the sample size. We refer the interested reader to [**?**] where robustness and effective use of this estimator are thoroughly studied. Yet, let us mention the existence of a theoretical scale range where the linear model, Eq. (**??**), holds, and which shows very helpful for practitioners to adequately adjust their linear fitting over a correct scale range.

## 3.3   Taqqu's Theorem

A central result for interpreting statistical modeling of network traffic is a celebrated Theorem due to M. Taqqu and collaborators [**?**, **?**, **?**], in which heavy-tailness of flow sessions has been put forward as a possible explanation for the occurrence of self-similarity of Internet traffic.

The original result considers a M/G/N queueing model served by $N$ independent sources, whose activities $Z_i(t)$, $i \in \{1, ..., N\}$, are described as a binary ON/OFF processes. The durations of the ON periods (corresponding to a packet train emission by a source) consists of i.i.d. positive random variables $\tau_{\mathrm{ON}}$, distributed according to a heavy-tail law $P_{ON}$, with exponent $\alpha$. Intertwined with the ON periods, the OFF periods (a source does not emit traffic), have i.i.d. random durations $\tau_{\mathrm{OFF}}$ drawn from another possibly heavy-tailed distribution $P_{\mathrm{OFF}}$ with tail index $\beta$. Thus, the $Z_i(t)$ consist of a 0/1 reward-renewal processes with i.i.d. activation periods.

Now, let $Y_N(t) = \sum_{i=1}^{N} Z_i(t)$ denote the aggregated traffic time series and define the cumulative process $X_N(Tt)$:

$$X_N(tT) = \int_0^{Tt} Y_N(u)\mathrm{d}u = \int_0^{Tt} \left(\sum_{i=1}^{N} Z_i(u)\right) \mathrm{d}u. \tag{12}$$

Taqqu's Theorem (cf. [**?**]) states that when taking the limits $N \to \infty$ (infinitely many users) and $T \to \infty$ (infinitely long observation duration), in this order, then $X_N(tT)$ behaves as:

$$X_N(tT) \sim \frac{\mathbb{E}\tau_{\mathrm{ON}}}{\mathbb{E}\tau_{\mathrm{ON}} + \mathbb{E}\tau_{\mathrm{OFF}}} NTt + C\sqrt{N}T^H B_H(t). \tag{13}$$

In this relation, $C$ is a constant and $B_H$ denotes a fractional Brownian motion with Hurst parameter:

$$H = \frac{3 - \alpha^*}{2}, \text{ where } \alpha^* = \min(\alpha, \beta, 2). \tag{14}$$

The order of the limits is compelling to obtain this asymptotic behavior; this has been discussed theoretically elsewhere and is beyond the issues we address here. The main conclusion of Taqqu's Theorem is that, in the limit of (infinitely) long observations, fractional Brownian motions superimposed to a deterministic linear trend, are relevant asymptotic models to describe the cumulated sum of aggregated traffic time series. Moreover, Eq. (**??**) shows that only heavy tailed distributions with infinite variance (i.e., $1 < min(\alpha, \beta) < 2$) can generate self-similarity associated to long range dependence (i.e. $H > 1/2$). Conversely, when both activity and inactivity periods have finite variance durations, $\alpha^* = 2$ and consequently $H = 1/2$, which means no long range dependency.

## 4   Experimental setup

To study the practical validity of Taqqu's result, we use the potential and facilities offered by the very large-scale, deeply reconfigurable and fully controllable experimental Grid5000 instrument, so as to overcome the limitations previously exposed of emulations, simulations or measurements in production networks. After a general overview of Grid5000, the metrology platform is described first. Design of a large set of experiments, aimed at studying the actual dependence between the network traffic self-similarity and the heavy-tailness of flow size distributions, is finally detailed.

### 4.1   Grid5000 instrument overview

Grid5000, is a 5000 CPUs nation-wide Grid infrastructure for research in Grid computing [**?**], providing a scientific tool for computer scientists similar to the large-scale instruments used by physicists, astronomers, and biologists. It is a research tool featured with deep reconfiguration, control and monitoring capabilities designed for studying large scale distributed systems and for complementing theoretical models and simulators. Up to 17 french laboratories involved and 9 sites are hosting one or more cluster of about 500 cores each. A dedicated private optical networking infrastructure, provided by RENATER, the French NREN is interconnecting the Grid5000 sites. Two international interconnection are also available: one at 10 Gb/s interconnecting Grid5000 with DAS3 in Netherlands and one at 1 Gb/s with Naregi in Japan. In the Grid5000 platform, the network backbone is composed of private 10 Gb/s Ethernet links connected to a DWDM core with dedicated 10 Gb/s lambdas with bottlenecks at 1 Gb/s in Lyon and Bordeaux (see Figure **??**).

Grid5000 offers to every user full control of the requested experimental resources. Its uses dedicated network links between sites, allows users reserving the same set of resources across successive experiments, allows users running their experiments in dedicated nodes (obtained by reservation) and lets users install and run their proper experimental condition injectors and measurements software. Grid5000 exposes two tools to implement these features : OAR is