

Exploring the Reliability and Robustness of HEAF(2) for Quantifying the Intensity of Long-Range Dependent Network Traffic

Karim Mohammed Rezaul and Vic Grout

karim@cair-uk.org

v.grout@newi.ac.uk

Centre for Applied Internet Research (CAIR), University of Wales, NEWI Plas Coch Campus, Wrexham, UK

Abstract

The intensity of Long-Range Dependence (LRD) for communications network traffic can be measured using the Hurst parameter. LRD characteristics in computer networks, however, present a fundamentally different set of problems in research towards the future of network design. There are various estimators of the Hurst parameter, which differ in the reliability of their results. Getting robust and reliable estimators can help to improve traffic characterization, performance modelling, planning and engineering of real networks. Earlier research [1] introduced an estimator called the Hurst Exponent from the Autocorrelation Function (HEAF) and it was shown why lag 2 in HEAF (i.e. HEAF (2)) is considered when estimating LRD of network traffic. This paper considers the robustness of HEAF(2) when estimating the Hurst parameter of data traffic (e.g. packet sequences) with outliers and also the reliability of HEAF(2).

Key words:

ACF, HEAF(2), LRD, Self-similarity.

Introduction

The Long-Range Dependence (LRD) property of traffic fluctuations has important implications for the performance, design and dimensioning of a network [2]. A simple, direct parameter characterizing the degree of long-range dependence is the *Hurst parameter*. The Hurst exponent (or Hurst parameter, H), which more than a half-century ago was proposed for the analysis of long-term storage capacity of reservoirs [3], is used today to measure the intensity of LRD in network traffic. A number of methods have been proposed to estimate the Hurst parameter. Some of the most popular include the aggregated variance time (V/T) [4], Rescaled-range (R/S) [2, 3], and the Higuchi [5] and wavelet-based methods [6, 7], although there are numerous others. In all these methods, H is calculated by taking the slope from a log-log plot. Over time, the wavelet-based Hurst parameter has acquired popularity in estimating LRD traffic. However the study [8] explores the advantages and limitations of wavelet estimators and found that a traffic trace with a number of deterministic shifts in the mean rate results in a steep wavelet spectrum which leads to overestimating the Hurst parameter. The intensity of long-range dependence is measured for file size or document size [9], packet counts (number of packets per unit time) [10, 11, 12], inter-arrival time [13, 14], frame size [15], connection size [16], packet length [17], number of bytes per unit time [2], bit or byte rate [18] and so on.

This paper continues the work on the new estimator introduced earlier named the *Hurst Exponent* by

Autocorrelation Function (HEAF) [1]. HEAF estimates H by a process which is simple, quick and reliable. In order to investigate the robustness of HEAF(2), two different types of simulation studies were performed. The first one uses fractional Gaussian noise (fGn) sequences generated by the Dietrich-Newsam algorithm [19, 20], which generates exact self-similar sequences. The second one uses a *fractional autoregressive moving average (FARIMA)* process [21, 22]. Stationarity is assumed for these kinds of classical models (FARIMA and fGN) because it is convenient from a theoretical point of view, particularly to check the validity of any hypothesis. The sequences generated by these models show a bell-shaped (i.e. Gaussian) curve either exactly or with small variation. However, our concern in this research is to determine whether, in the case that the underlying process is not FARIMA or fGN, HEAF(2) can still capture the long-range dependency of the traffic. We investigate what role HEAF(2) can play to yield an estimate with a good degree of accuracy if the traffic is non-stationary. For instance, if the data traffic possesses outliers, we consider how to estimate H by eliminating these outliers to give satisfactory and reliable information.

The paper is organised as follows. Section 2 describes the definitions of self-similarity, long-range dependence and the autocorrelation function. Section 3 introduces the HEAF estimator. Section 4 discusses robust versions of the autocorrelation function. Finally the results are presented in section 5.

2. Self-Similarity, Long-Range Dependence and the Autocorrelation Function

In general, two or more objects having the same characteristics are called self-similar. A phenomenon that is self-similar looks the same or behaves the same when viewed at different degrees of magnification or different scales on a dimension and is bursty over all time scales. Self-similarity is the property of a series of data points to retain a pattern or appearance, regardless of the level of granularity used, and is the result of long-range dependence in the data series. If a self-similar process is bursty on a wide range of timescales, it may exhibit long-range-dependence. In general, lagged autocorrelations are used in time series analysis for empirical stationary tests. Self-similarity manifests itself as long-range dependence (i.e., long memory) in the time series of arrivals. The evidence of very slow, linear decay in the sample lag *autocorrelation function (ACF)* indicates nonstationary behaviour [23]. The research in [24] shows that Internet traffic is nonstationary.

Long-range-dependence means that all the values, at any time, are correlated in a positive and non-negligible way with values at all future instants. A continuous time process, $Y = \{Y(t), t \geq 0\}$, is self-similar if it satisfies the following condition [25]:

$$Y(t) \stackrel{d}{=} a^{-H} Y(at), \quad \forall a > 0, \text{ and } 0 < H < 1 \quad (2.1)$$

where H is the index of self-similarity, called the *Hurst parameter*. and the equality is in the sense of finite-dimensional distributions.

The stationary process X is said to be a long-range dependent process if its autocorrelation function is non-summable [26], meaning that $\sum_{k=-\infty}^{\infty} r_k = \infty$ (2.2)

The details of how ACF decays with k are of interest because the behaviour of the tail of the ACF completely determines its summability. According to [2], X is said to exhibit long-range dependence if

$$r_k \sim L(t)k^{-(2-2H)}, \text{ as } k \rightarrow \infty \quad (2.3)$$

where $\frac{1}{2} < H < 1$ and $L(\cdot)$ slowly varies at infinity, i.e.,

$$\lim_{t \rightarrow \infty} \frac{L(xt)}{L(t)} = 1, \text{ for all } x > 0 \quad (2.4)$$

Equation (2.3) implies that LRD is characterized by an autocorrelation function that decays hyperbolically rather than exponentially fast.

3. HEAF: A ‘Hurst Exponent by Autocorrelation Function’ Estimator

A new estimator has been introduced [1] by extending the approach of Kettani and Gubner [27]. As in [27], for given observed data X_i (i.e. X_1, \dots, X_n), the sample autocorrelation function can be calculated by the following method:

$$\text{Let } \hat{m}_n = \frac{1}{n} \sum_{i=1}^n X_i \quad (3.1)$$

$$\text{and } \hat{g}_n(k) = \frac{1}{n} \sum_{i=1}^{n-k} (X_i - \hat{m}_n)(X_{i+k} - \hat{m}_n), \quad (3.2)$$

where $k=0, 1, 2, \dots, n$,

$$\text{with } \hat{s}_n^2 = \hat{g}_n(0). \quad (3.3)$$

Then the sample autocorrelations of lag k are given by

$$\hat{r}_k = \frac{\hat{g}_n(k)}{\hat{s}_n^2} \quad (3.4)$$

(Equations (3.1), (3.2), (3.3) and (3.4) denote the sample mean, the sample covariance, the sample variance and the sample autocorrelation, respectively). A second-order stationary process is said to be exactly second-order self-similar, with Hurst exponent $1/2 < H < 1$, if

$$r_k = 0.5 [(k+1)^{2H} - 2k^{2H} + (k-1)^{2H}] \quad (3.5)$$

From equation (3.5), Kettani and Gubner suggest a moment estimator of H . They consider the case $k=1$ and replace r_1 by its sample estimate \hat{r}_1 , as defined in equation (3.4). This gives an estimate for H of the form

$$\hat{H} = \frac{1}{2} + \frac{1}{2 \log_e 2} \log_e (1 + \hat{r}_1) \quad (3.6)$$

Clearly, this estimate is straightforward to evaluate, requiring no iterative calculations. For more details of the properties of this estimator, see Kettani and Gubner [27].

An alternative estimator of H is proposed based upon equation (3.5), by considering the cases where $k > 1$. Note that the sample equivalent of equation (3.5) can be expressed as

$$f(H) = \hat{r}_k - 0.5\{(k+1)^{2H} - 2k^{2H} + (k-1)^{2H}\} = 0. \quad (3.7)$$

Thus, for a given observed \hat{r}_k , $k > 1$, a suitable numerical procedure can be used to solve this equation, and find an estimate of H . This is denoted as a HEAF(k) estimate of H .

To solve equation (3.7) for H , the well-known Newton-Raphson (N-R) method is used. This requires the derivative of $f(H)$. Here note that $k \geq 1$,

$$f'(H) = -0.5 \left\{ \begin{array}{l} (2 \log(k+1))(k+1)^{2H} \\ - (4 \log(k))k^{2H} + \\ (2 \log(k-1))(k-1)^{2H} \end{array} \right\} \quad (3.8)$$

Hence, the algorithm to estimate HEAF(k), for any lag k , consists of the following steps:

1. Compute the sample autocorrelations for lag k of a given data set by equation (3.4). (X_i can be denoted as the number of bits, bytes, packets or bit rates observed during the i th interval. If X_i is a Gaussian process, it is known as fractional Gaussian noise).
2. Make an initial approximation for H , e.g. $H_1 = 0.6$, then calculate H_2, H_3, H_4, \dots , successively using $H_{r+1} = H_r - f(H_r) / f'(H_r)$, until convergence, to find the estimate \hat{H} for the given lag k . An initial consideration is the case where $k = 2$ in equation (3.2); i.e. HEAF(2) is considered first.

One of the major advantages of the HEAF estimator is speed, as the N-R-method converges very quickly to a root. There is no general convergence criterion for N-R. Its convergence depends on the nature of the function and on the accuracy of the initial approximation. Fortunately, the form of the function (i.e., equation (3.7)) appears to converge quickly (within at most four iterations) for any initial approximation in the range of interest, namely H in $(0.2, 1)$. If an iteration value, H_r is such that $f'(H_r) \cong 0$, then one can face ‘‘division by zero’’ or a near-zero number. This will give a large magnitude for the next value, H_{r+1} which in turn stops the iteration. This problem can be resolved by increasing the tolerance parameter in the N-R program. All HEAF(k), for $k = 2, \dots, 11$, have been considered and no difficulty in finding the root in $(0.5, 1)$ has been encountered.

4. A Robust Autocorrelation Function

The forecasting of network traffic and Quality of Service (QoS) can be affected by additive outliers. The sample ACF used in HEAF (2) is somewhat

controversial. In this research, we test the performance of HEAF (2) by using three robust ACFs: Trimmed ACF (TACF) [28], variance-ratio of differences and sums, known as the D/S variance estimator [29, 30] and the weighted sample autocorrelation function (shortened to WACF) [31]. Polasek [32] showed how to eliminate these additive outliers by different robust ACFs. According to his findings, the sample ACF (i.e. moment based) was surprisingly ranked second after TACF for eliminating additive outliers. Due to space limitations we only present the results from Trimmed ACF.

The Trimmed ACF can be calculated by the following procedure:

Let $z_{(1)} \leq z_{(2)} \leq \dots \leq z_{(n)}$ be the ordered observations of the given time series z_1, z_2, \dots, z_n . Chan and Wei [28] introduced the \mathbf{a} -trimmed sample autocorrelation function (TACF) defined by

$$\hat{\mathbf{r}}_T(k) = \frac{\hat{\mathbf{g}}_T(k)}{\hat{\mathbf{g}}_T(0)}$$

where

$$\hat{\mathbf{g}}_T(k) = \frac{\sum_{t=k+1}^n (z_{t-k} - \bar{z}^{(a)}) (z_t - \bar{z}^{(a)}) L_{t-k}^{(a)} L_t^{(a)}}{\sum_{t=k+1}^n L_{t-k}^{(a)} L_t^{(a)}}$$

$$\bar{z}^{(a)} = \frac{\sum_{t=1}^n z_t L_t^{(a)}}{\sum_{t=1}^n L_t^{(a)}} \quad \text{and}$$

$$L_t^{(a)} = \begin{cases} 0, & \text{if } z_t \leq z_{(g)} \text{ or } z_t \geq z_{(n-g+1)} \\ 1, & \text{otherwise} \end{cases}$$

where $g = [\mathbf{a}n]$ is the integer part of $\mathbf{a}n$ and $0 \leq \mathbf{a} \leq 0.05$. Chan and Wei showed that TACF is, in general, very successful in removing the adverse effect of outliers in the estimation of the ACF. The parameter, called automatic alpha can be estimated by the trimmean filter (TMF) [33, 34].

The procedure for estimating alpha by TMF is as follows:

1. Sort the data in ascending order.
2. Calculate the parameter Q according to the equation below

$$Q = \frac{[U(20\%) - L(20\%)]}{[U(50\%) - L(50\%)]}$$

- where $U(x\%)$ is the average of the upper $x\%$ of the ordered sample and $L(x\%)$ is the average of the lower $x\%$ of the ordered sample.

- Q is a measure of the departure of the distribution contained in the sample from a normal distribution.

3. Trim off each tail of the ordered distribution according to the value of the trimmean parameter alpha.

$$\mathbf{a}(Q) = \begin{cases} 0.04 & Q \leq 1.75 \\ 0.04 + 0.01 * \frac{(Q-1.75)}{0.25} & 1.75 < Q < 2.0 \\ 0.05 & Q \geq 2.0 \end{cases}$$

Note that the alpha parameter given in [33, 34] is modified here for estimating a good degree of accuracy when considering network traffic data. The TMF assumes the distribution to be symmetric, but not necessarily Gaussian. For a pure Gaussian distribution of data, 4 percent of the data is trimmed from each tail of the original sorted distribution. For a given segment of time, a maximum of 5 percent of the data is trimmed off each tail.

5. Results and Discussion

This section describes the robustness and reliability of the HEAF(2) estimator.

5.1 Robustness of the HEAF(2) estimator

In [1, 35, 36], the results show that HEAF(2) is an estimator of H with relatively good bias and mean square errors (MSE) when estimating fractional Gaussian noise or FARIMA processes. Because of its simplicity and reliability, it is believed that HEAF (2) can be used for real time network traffic control. Of course, a real process will be unlikely to be exactly an fGn process or even FARIMA process. Indeed, a real process may suffer from 'noise', discrepancy values or other outliers. This section considers the robustness of the proposed estimator, HEAF(2), against departures from ideal assumptions.

In order to test the robustness of HEAF (2), we generate some noisy sequences by mixing the data sequences generated by FARIMA (0, d, 0) and fGn processes for a particular Hurst parameter (H). Obviously H will be changed when using noisy data, meaning that the process (FARIMA or fGN) no longer exists as it holds additive outliers. Figure 1 shows a pictorial view of noisy data to be analysed in order to explore the robustness of the HEAF(2).

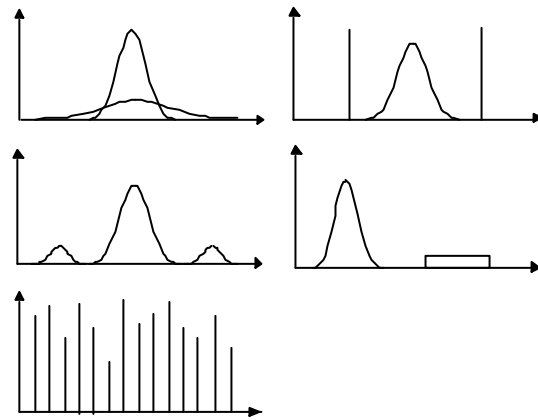


Fig. 1. Pictorial view of noisy samples (i.e. data with additive outliers)

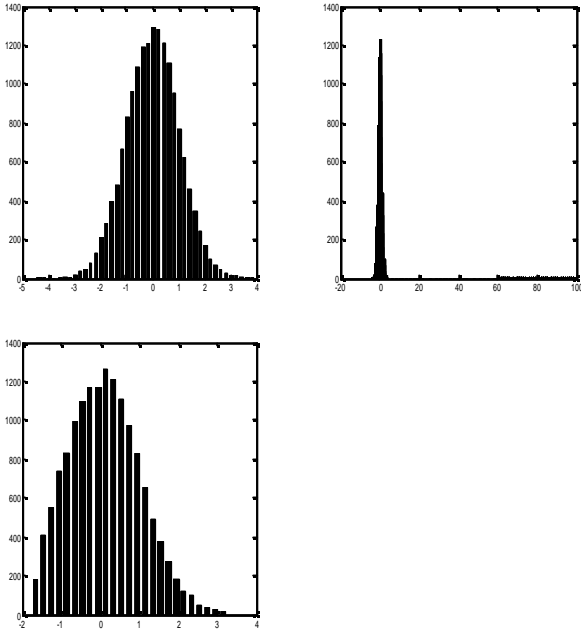


Fig. 2. Top left figure - Data generated by the FARIMA (0,d, 0) process for $H=0.6$ (H measured by HEAF (2) = 0.585), $N = 16384$. Top right figure – generated noisy sample (measured $H = 0.8993$). Bottom figure - $H = 0.576$ (after elimination of the outliers) where $\alpha = 0.048$

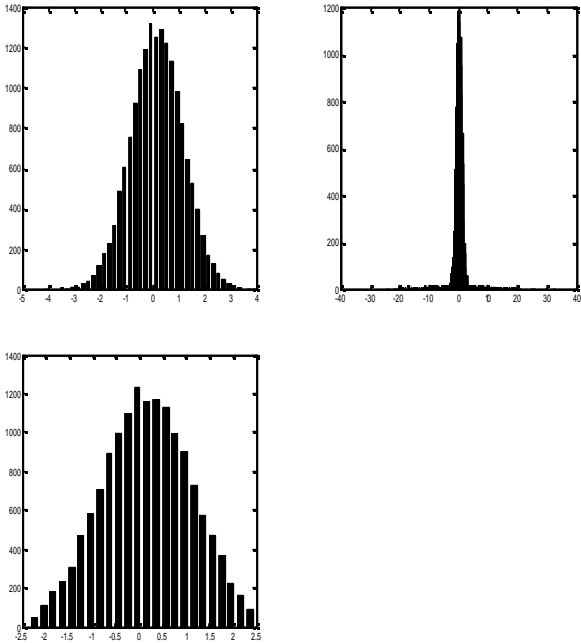


Fig. 3. Top left figure - Data generated by fGN process for $H=0.6$ (H measured by HEAF (2) = 0.795), $N = 16384$. Top right figure – generated Noisy sample (measured $H = 0.556$). Bottom figure - $H = 0.77$ (after elimination of the outliers) where $\alpha = 0.048$

In Figure 2, the top left figure shows a histogram of a FARIMA (0, d, 0) process with $H = 0.6$. The estimated H by HEAF(2) is 0.585. The top right figure gives a histogram for noisy samples generated by mixing with a FARIMA (0,d, 0) process for $H = 0.6$, having the sample length, $N = 16384$. The histogram at the bottom of Figure 2 is plotted after elimination of the outliers shown in the top right figure. The estimated H for noisy samples (top right figure) and samples after elimination are 0.8993 and

0.576 respectively. The outliers from noisy samples are eliminated by automatic alpha and then the alpha value is used in TACF which in turn applied in HEAF(2).

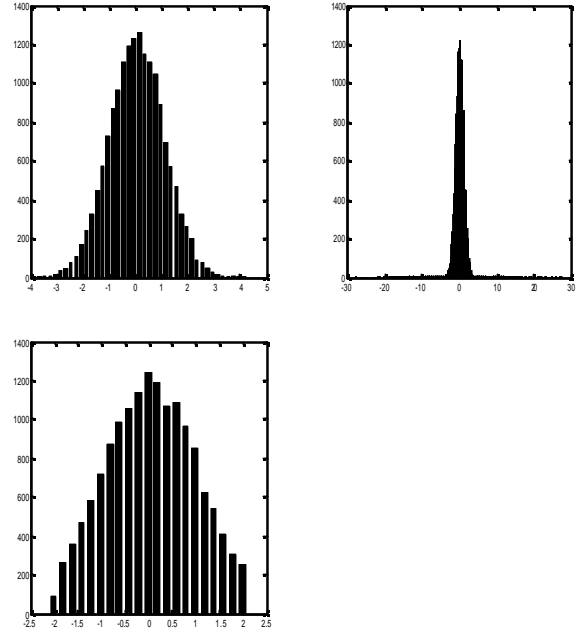


Fig. 4. Top left figure - Data generated by FARIMA (0,d, 0) process for $H=0.7$ (H measured by HEAF (2) = 0.683), $N = 16384$. Top right figure – generated Noisy sample (measured $H = 0.573$). Bottom figure - $H = 0.65$ (after elimination of the outliers) where $\alpha = 0.041$

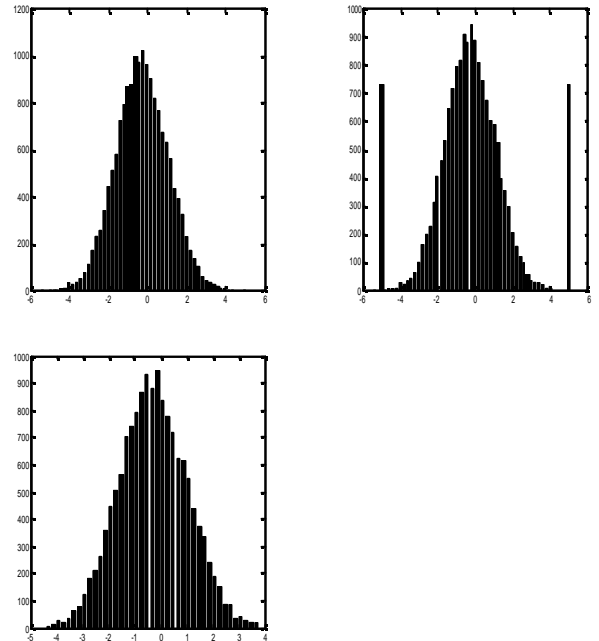


Fig. 5. Top left figure - Data generated by FARIMA (0,d, 0) process for $H=0.9$ (H measured by HEAF (2) = 0.858), $N = 16384$. Top right figure – generated Noisy sample (measured $H = 0.701$). Bottom figure - $H = 0.857$ (after elimination of the outliers) where $\alpha = 0.045$

In Figure 6, uniform random numbers are chosen to generate FARIMA (0,d,0) sequences for various Hurst parameters. Due to the uniform random function used in the process, FARIMA (0,d,0) generates only positive sequences, which can imitate real Internet packet sequences. It is clear from the results presented in this

research that additive outliers can be removed by applying robust ACF and that, after elimination of these outliers from different case studies, it is evident that HEAF(2) yields a reliable value for H.

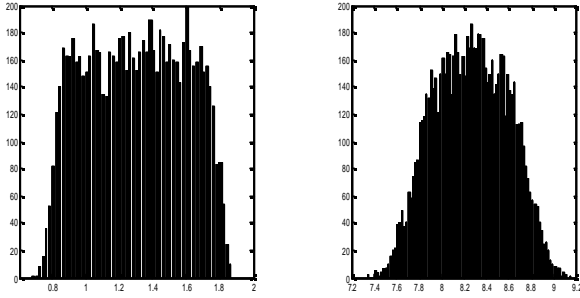


Fig. 6. Left figure - FARIMA (0,d,0) for H = 0.6 (measured H by HEAF(2) = 0.579), N = 8192. Right figure - FARIMA (0, d, 0) for H = 0.8 (measured H by HEAF(2) = 0.774), N = 8192.

5.2 Reliability of HEAF(2) estimator

In this section, we discuss how to determine a reliable estimator based on simulation experiments. In most cases, researchers find bias, MSE and the confidence interval (CI) of the estimator to explore its reliability and robustness. Sometimes it is hard to make a decision by looking at the CI of the estimator. For instance, for H = 0.7, 100 different realisations of self-similar sequences have been generated, each with sample length N = 10000. Now, for an estimator, the CI is found to be (0.583, 0.605). Looking at such a CI for an estimator, one can easily conclude that the estimator performs better for that particular Hurst parameter. However, the real scenario can be observed when looking at the 20 lowest (say) and the 20 highest (say) values of the Hurst parameter from the realisations out of those 100 realisations. Earlier research [35] investigates the properties of some existing estimators. Here we show a comparison of the reliability of the estimators such as rescaled-range analysis (R/S), variance-time analysis (V/T), the wavelet-based estimator and the Higuchi method in conjunction with the HEAF(2) estimator.

Table I: Lowest H-values for 15 different realizations. Sample length, N = 16384. FARIMA (0,d,0) process. H = 0.6.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.561	0.5728	0.379	0.301	0.369
2	0.5677	0.5767	0.4	0.338	0.438
3	0.5714	0.577	0.444	0.341	0.467
4	0.5734	0.578	0.45	0.36	0.468
5	0.5758	0.5782	0.452	0.366	0.473
6	0.5761	0.57845	0.461	0.398	0.48
7	0.5767	0.5787	0.464	0.403	0.482
8	0.577	0.57885	0.465	0.405	0.485
9	0.5774	0.57975	0.467	0.412	0.497
10	0.5775	0.5802	0.481	0.417	0.499
11	0.5795	0.58105	0.484	0.427	0.499
12	0.5807	0.5812	0.501	0.455	0.503
13	0.5807	0.5819	0.512	0.458	0.505
14	0.5819	0.5821	0.518	0.463	0.509
15	0.5822	0.5828	0.526	0.466	0.514

With the intention of understanding the long-range dependence, we started our simulation study with FARIMA(0,d,0), which is the simplest and most

fundamental of the fractionally differenced FARIMA processes. We generated 100 different realisations, each with sample length, N = 16384. Note that each realisation implies a set of data that contains the sample length, N = 16384. Here, we wish to observe how H-values vary for different realisations (i.e. data set)s generated by the process for a particular H (e.g. H = 0.6). This is the reason for choosing the 15 lowest and 15 highest H values to find the stable range of estimators for a particular H.

Table II: Highest H-values for 15 different realizations. Sample length, N = 16384. FARIMA (0,d,0) process. H = 0.6.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.616	0.6006	0.769	1.043	0.777
2	0.6124	0.59975	0.757	1.011	0.728
3	0.6116	0.59915	0.753	0.986	0.721
4	0.6107	0.5986	0.741	0.963	0.715
5	0.6095	0.5976	0.74	0.944	0.693
6	0.6093	0.5975	0.738	0.905	0.686
7	0.6079	0.597	0.725	0.901	0.678
8	0.6079	0.59695	0.722	0.89	0.667
9	0.6066	0.59665	0.708	0.888	0.666
10	0.6059	0.59585	0.701	0.887	0.652
11	0.6053	0.5958	0.694	0.878	0.65
12	0.605	0.5957	0.685	0.877	0.648
13	0.6048	0.59535	0.682	0.852	0.645
14	0.6035	0.59495	0.682	0.852	0.642
15	0.6032	0.59475	0.68	0.851	0.641

Table III: Lowest H-values for 15 different realizations. Sample length, N = 16384. FARIMA (0,d,0) process. H = 0.7.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.6668	0.66385	0.497	0.422	0.507
2	0.6669	0.66575	0.502	0.48	0.526
3	0.6679	0.6659	0.543	0.496	0.542
4	0.6689	0.66665	0.553	0.497	0.542
5	0.669	0.667	0.554	0.517	0.547
6	0.6702	0.66785	0.556	0.542	0.563
7	0.6702	0.66825	0.567	0.545	0.565
8	0.6705	0.66855	0.571	0.563	0.566
9	0.6714	0.66855	0.572	0.564	0.573
10	0.6714	0.66955	0.576	0.565	0.574
11	0.6723	0.6696	0.582	0.594	0.584
12	0.6731	0.67025	0.585	0.607	0.596
13	0.6733	0.6709	0.587	0.608	0.599
14	0.6738	0.6715	0.588	0.633	0.606
15	0.6743	0.67155	0.588	0.649	0.606

Table IV: Highest H-values for 15 different realizations. Sample length, N = 16384. FARIMA (0,d,0) process. H = 0.7.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.7002	0.69075	0.863	1.099	0.904
2	0.6989	0.6893	0.86	1.047	0.885
3	0.6975	0.6892	0.834	1.028	0.871
4	0.6969	0.68785	0.832	1.022	0.812
5	0.6945	0.6865	0.814	1.014	0.81
6	0.6938	0.6859	0.809	0.993	0.804
7	0.6933	0.68535	0.808	0.96	0.789
8	0.6932	0.68505	0.801	0.959	0.782
9	0.693	0.68495	0.791	0.958	0.777
10	0.692	0.6845	0.79	0.947	0.774
11	0.6919	0.6844	0.779	0.936	0.773
12	0.6915	0.6843	0.779	0.936	0.769
13	0.6914	0.684	0.773	0.934	0.767
14	0.6912	0.6838	0.771	0.929	0.767
15	0.6911	0.6837	0.771	0.926	0.765

Table I, Table III, Table V and Table VII, illustrate the H-values from 15 different realisations (out of 100

realisations) for $H = 0.6$, $H = 0.7$, $H = 0.8$ and $H = 0.9$ respectively. It is shown clearly in these tables that the HEAF(2) and wavelet-based methods are more stable than other estimators and that the values of H estimated are in an acceptable range for the corresponding simulated H . Also, when considering the highest H -values, HEAF(2) outperforms the other estimators for corresponding $H = 0.6$, $H = 0.7$, $H = 0.8$ and $H = 0.9$, as shown in Tables II, IV, VI and VIII. Note that R/S analysis, V/T analysis and the Higuchi method are not reliable estimators as they sometimes underestimate or overestimate the degree of long-range dependence.

Table V: Lowest H -values for 15 different realizations. Sample length $N = 16384$. FARIMA (0,d,0) process. $H = 0.8$.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.7545	0.7492	0.579	0.362	0.591
2	0.7551	0.75445	0.589	0.38	0.604
3	0.7576	0.7551	0.609	0.394	0.613
4	0.7578	0.75525	0.611	0.432	0.636
5	0.7589	0.7564	0.617	0.536	0.644
6	0.7605	0.75665	0.617	0.563	0.645
7	0.761	0.75685	0.626	0.565	0.645
8	0.7616	0.75755	0.628	0.602	0.65
9	0.7619	0.7577	0.628	0.608	0.662
10	0.762	0.75815	0.634	0.61	0.666
11	0.7624	0.75825	0.642	0.639	0.677
12	0.7634	0.75845	0.653	0.643	0.679
13	0.7634	0.75915	0.653	0.647	0.683
14	0.7638	0.7594	0.655	0.65	0.686
15	0.7639	0.76095	0.662	0.665	0.687

Table VI: Highest H -values for 15 different realizations. Sample length $N = 16384$. FARIMA (0,d,0) process. $H = 0.8$.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.7935	0.7816	1.059	1.138	0.93
2	0.7928	0.7801	0.968	1.119	0.912
3	0.7925	0.7791	0.96	1.109	0.899
4	0.7882	0.77565	0.922	1.098	0.897
5	0.7872	0.77565	0.921	1.054	0.89
6	0.7867	0.77535	0.905	1.036	0.883
7	0.786	0.7753	0.901	1.035	0.882
8	0.7859	0.7751	0.895	1.034	0.878
9	0.785	0.7751	0.887	1.01	0.874
10	0.7842	0.7749	0.881	1.01	0.872
11	0.7839	0.77425	0.878	1.009	0.87
12	0.7837	0.77395	0.878	0.997	0.867
13	0.7836	0.77385	0.878	0.995	0.861
14	0.7835	0.77375	0.871	0.994	0.859
15	0.7834	0.77355	0.871	0.991	0.851

Table VII: Lowest H -values for 15 different realizations. Sample length $N = 16384$. FARIMA (0,d,0) process. $H = 0.9$.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.8401	0.8451	0.596	0.396	0.654
2	0.8413	0.8459	0.603	0.535	0.655
3	0.8426	0.84615	0.622	0.583	0.681
4	0.8441	0.84615	0.647	0.593	0.721
5	0.845	0.84695	0.654	0.611	0.723
6	0.8452	0.8473	0.658	0.624	0.73
7	0.8467	0.8473	0.673	0.625	0.745
8	0.8468	0.8486	0.685	0.645	0.745
9	0.847	0.84905	0.692	0.646	0.751
10	0.8471	0.8491	0.706	0.67	0.751
11	0.8471	0.84945	0.707	0.68	0.768
12	0.8479	0.84955	0.71	0.695	0.779
13	0.8489	0.8497	0.71	0.697	0.781
14	0.8499	0.84975	0.712	0.71	0.782
15	0.8501	0.85	0.719	0.726	0.782

Table VIII: Highest H -values for 15 different realizations. Sample length $N = 16384$. FARIMA (0,d,0) process. $H = 0.9$.

Number of realisations	HEAF(2)	Wavelet	R/S	V/T	Higuchi
1	0.8922	0.8666	1.104	1.215	0.992
2	0.8893	0.86645	1.025	1.149	0.978
3	0.8829	0.86635	1.018	1.144	0.976
4	0.8778	0.865	1.014	1.116	0.968
5	0.8772	0.86495	1.012	1.116	0.963
6	0.8755	0.8648	1.012	1.11	0.958
7	0.8752	0.86455	1.009	1.109	0.955
8	0.8745	0.86455	1.008	1.104	0.948
9	0.8744	0.8642	0.995	1.101	0.947
10	0.8732	0.8636	0.991	1.098	0.938
11	0.8721	0.8635	0.985	1.094	0.936
12	0.872	0.86315	0.979	1.093	0.935
13	0.8713	0.8628	0.976	1.092	0.933
14	0.8712	0.8628	0.97	1.091	0.932
15	0.871	0.8621	0.965	1.09	0.929

6. Conclusion

It is possible to derive wrong conclusions and wrong models when measuring the intensity of the LRD with unreliable estimators. In this research we have shown that the plausible H for given data can be overestimated or underestimated due to additive outliers occurring in the data. These outliers can be removed by applying robust ACF in HEAF(2) and in this case HEAF(2) yields consistent and reliable results. Also, based on the comparison of simulation experiments for both fGn and FARIMA (0, d, 0) processes, it is evident that HEAF(2) is a stable method that quantifies the reliable degree of long-range dependence. Through its simplicity, robustness and reliability, we believe that HEAF(2) can be used to estimate the intensity of LRD in real time network traffic.

Acknowledgments

The authors would like to acknowledge Professor Robert Gilchrist, Department of CCTM, London Metropolitan University, UK, for valuable discussions about the robustness of the estimators.

References

- [1] Karim M. Rezaul, Algirdas Pakštas, Robert Gilchrist, Thomas M. Chen, HEAF: A Novel Estimator for Long-Range Dependent Self-similar Network Traffic, Y. Koucheryavy, J. Harju, and V.B. Iversen (Eds.): *Next Generation Teletraffic and Wired/Wireless Advanced Networking (NEW2AN)*, May 29 - June 2, 2006, LNCS 4003, pp. 34 – 45.
- [2] Leland Will E. Taqqu M. S., Willinger W. and Wilson D. V., On the Self-similar nature of Ethernet Traffic (Extended version), *IEEE/ACM Transactions on Networking*, February 1994, Vol. 2, No. 1, pp. 1-15.
- [3] Hurst H. E., Long-term storage capacity of reservoirs, *Transactions of the American Society of Civil Engineers*, 1951, vol.116, pp 770-808.
- [4] Ton Dieker, *Simulation of Fractional Brownian Motion*, Masters Thesis, Department of Mathematical Sciences, University of Twente, The Netherlands, 2004.
- [5] T. Higuchi, Approach to an irregular time series on the basis of the fractal theory, *Physica D*, 1988, vol.31, pp. 277-283.
- [6] P. Abry, P. Flandrin, M. S. Taqqu and D. Veitch, Wavelets for the Analysis, Estimation, and Synthesis of Scaling Data, K. Park and W. Willinger (editors), *Self-Similar Network Traffic and Performance Evaluation*. John Wiley & Sons, New York, 2000,

- pp. 39-88.
- [7] P. Abry and D. Veitch, Wavelet Analysis of Long-Range Dependent traffic. *IEEE Transactions on Information Theory*, 1998, vol.44, No 1, pp.2-15.
- [8] Stilian Stoev, Murad Taqqu, Cheolwoo Park and J.S. Marron, *Strengths and Limitations of the Wavelet Spectrum Method in the Analysis of Internet Traffic*, Technical Report #2004-8, March 26, 2004.
- [9] Kihong Park, Gitae Kim, Mark Crovella, On the relationship between file sizes, transport protocols, and self-similar network traffic, *Fourth International Conference on Network Protocols (ICNP'96)*. 1996, pp. 171-180.
- [10] Vern Paxson and Sally Floyd, Wide-Area Traffic: The Failure of Poisson Modeling, *IEEE/ACM Transactions on Networking*, June 1995, Vol. 3 No. 3, pp. 226-244.
- [11] *UNC Network Data Analysis Study Group*, University of North Carolina, http://www-dirt.cs.unc.edu/net_lrd/ (visited on May 23, 2005).
- [12] W. Willinger, M. S. Taqqu, W. E. Leland, and D. V. Wilson, Self-similarity in high speed packet traffic: Analysis and modeling of ethernet traffic measurements, *Statistical Science*, 1995, Vol. 10, pp. 67-85.
- [13] Dario Rossi, Luca Muscariello, Marco Mellia, On the Properties of TCP Flow Arrival Process, *IEEE International Conference on Communications (ICC 2004)*, Paris - France, June 20-24, 2004.
- [14] Ashok Erramilli, Onuttom Narayan, and Walter Willinger, Experimental Queueing Analysis with Long-range Dependent Packet Traffic, *IEEE/ACM Transactions on Networking*, April 1996, Vol. 4, No. 2, pp. 209-223.
- [15] O. Rose, *Estimation of the Hurst Parameter of Long-Range Dependent Time Series*, Report No. 137, February 1996, Institute of Computer Science, University of Wurzburg.
- [16] E. Willekens and J. Teugels, Asymptotic expansions for waiting time probabilities in an M/G/1 queue with longtailed service time, *Queueing Systems 10*, 1992, pp. 295-312.
- [17] An Ge, Franco Callegati, and Lakshman S. Tamil, On Optical Burst Switching and Self-Similar Traffic, *IEEE Communications Letters*, March 2000, Vol. 4, No. 3, pp. 98-100.
- [18] R.J. Gibbens, Traffic characterisation and effective bandwidths for broadband network traces, *Stochastic Networks, Theory and Applications*, 1996, pp. 169-179, Oxford Science Pub.
- [19] Dietrich C. R. and Newsam G. N., A fast and exact method for multidimensional Gaussian Stochastic simulations, *Water Resources Research*, 1993, vol. 29, No. 8, pp.2861-2869.
- [20] C.R. Dietrich and G.N. Newsam, Fast and exact simulation of stationary Gaussian processes through circulant embedding of the covariance matrix, *SIAM Journal on Scientific Computing*, 1997, vol.18, pp.1088-1107.
- [21] Kokozka, P. S. and Taqqu, M. S., Fractional ARIMA with stable innovations, *Stochastic Processes and their Applications* 1995, vol.60, pp.19-47.
- [22] Stilian Stoev and Murad S. Taqqu, Simulation methods for linear fractional stable motion and FARIMA using the Fast Fourier Transform, *Fractals*, 2004, vol.12, No 1, pp.95-121.
- [23] Brocklebank J. and D. Dickey, SAS System for Forecasting Time Series. *SAS Institute Inc.* Cary NC. 1986.
- [24] Jin Cao, William S. Cleveland, Dong Lin, and Don X. Sun, On the Nonstationarity of Internet Traffic, *Proc. ACM SIGMETRICS '01*, 102-112, 2001.
- [25] Walter Willinger, Vern Paxson, and Murad Taqqu, Self-similarity and Heavy Tails: Structural Modeling of Network Traffic, Adler, R., Feldman, R., and Taqqu, M.S., (editors), *In A Practical Guide to Heavy Tails: Statistical Techniques and Applications*, Birkhauser, 1998.
- [26] Cox D., Long-Range Dependence: a Review. H. A. David and H. T. David (eds.), *In Statistics: An Appraisal*, Iowa State Statistical Library, The Iowa State University Press, 1984, pp.55-74.
- [27] H. Kettani and J. A. Gubner, A Novel Approach to the Estimation of the Hurst Parameter in Self-Similar Traffic, *Proceedings of the 27th Annual IEEE Conference on Local Computer Networks (LCN 2002)*, Tampa, Florida, November, 2002, pp.160-165.
- [28] Chan, W. S. and Wei, W. W. S., A comparison of some estimators of time series autocorrelations, *Computational Statistics & Data Analysis*, vol. 14, 1992, pp. 149-163.
- [29] Gnanadesikan, R. and Kettenring, J.R. (1972) Robust estimates, residuals, and outlier detection with multiresponse data, *Biometrics* 28, 81-124.
- [30] Polasek W., *Robust and Resistant measures for the autocorrelation function*, mimeo, Institute of Statistics and Informatics, University of Vienna, 1982.
- [31] Wang W. and Wei W.S.W., *ASA proceedings of Business and Economic Statistics Section*, pp. 175-180, 1993.
- [32] Polasek, W. - Mertl, R. *Robust and jackknife estimators for the autocorrelation function*. Österreichische Zeitschrift für Statistik und Informatik 20, 1990, pp. 351-364.
- [33] Hogg Robert V., Adaptive Robust procedures: A practical Review and some suggestions for future applications and theory, *Journal of the American Statistical Association*, vol. 69, December 1974, pp.909-927
- [34] Vic Barnett and Toby Lewis, *Outliers in Statistical data*, 1994, 3rd Edition, John Wiley & Sons Ltd.
- [35] Karim M. Rezaul, Algirdas Pakštas, Robert Gilchrist, "Investigation of the Properties of the HEAF Estimator Using Simulation Experiments and MPEG-encoded Video Traces", *10th IEEE International Conference on Intelligent Engineering Systems (INES 2006)*, London, UK, June 26-28, 2006, pp. 276-281.
- [36] Karim M. Rezaul, Robert Gilchrist and Algirdas Pakštas, "Long-range Dependent Self-similar Network Traffic: A Simulation Study to Compare Some New Estimators", *7th Annual PostGraduate Symposium on the Convergence of Telecommunications, Networking and Broadcasting (PGnet 2006)*, Liverpool, UK, June 26-27, 2006, pp. 323-329.



Karim Mohammed Rezaul received BSc. degree in the field of Naval Architecture and Marine Engineering from Bangladesh University of Engineering and Technology (BUET), Dhaka, Bangladesh in 1998. In 2001, he was awarded an MSc degree in Marine Technology from Norwegian University of Science and Technology, Trondheim, Norway. He is a Member of Royal Institution of

Naval Architects (RINA), UK and Institution of Engineers Bangladesh (IEB), Bangladesh.

In February 2002, Mr. Karim was appointed as visiting lecturer in the department of computing, communications Technology and Mathematics at London Metropolitan University and continued until June 2005. Now he is a visiting lecturer in the department of computing at Central College London. His research interests include Network Traffic Engineering, Statistical analyses of data, Heavy tail distribution of network traffic, Long-range dependent network traffic, Modelling network traffic by Fractal and wavelet method, Traffic control mechanism and Stochastic process and probability distribution.



Vic Grout was awarded the BSc(Hons) degree in Mathematics and Computing from the University of Exeter (UK) in 1984 and the PhD degree in Communication Engineering from Plymouth Polytechnic (UK) in 1988.

He has worked in senior positions in both academia and industry for twenty years and has published and presented over 100 research papers. He is currently a Reader in Computer Science at the University of Wales NEWI, Wrexham in the UK, where he leads the Centre for Applied Internet Research (CAIR). His research interests and those of his research students span several areas of computational mathematics, particularly the application of heuristic principles to large-scale problems in network design and management.

Dr. Grout is a Chartered Engineer, Chartered Scientist, Chartered Mathematician and Chartered IT Professional, a Member of the IMA, IET(IEE), ACM, IEEE, IEEE Computer and Communications Societies and a Fellow of the British Computer Society (BCS). He chairs the biennial international conference series on Internet Technologies and Applications (ITA 05 and ITA 07).