УДК 004.62:[004.732:004.738.5.057.4]

# PRINCIPAL COMPONENT ANALYSIS OF NETWORK TRAFFIC MEASUREMENTS: THE «CATERPILLAR»-SSA APPROACH

*I. Antoniou* [a,b], *V. V. Ivanov* [a,c], *Valery V. Ivanov* [c], *P. V. Zrelov* [c]

[a] International Solvay Institutes for Physics and Chemistry, Brussels, Belgium
[b] Department of Mathematics, Aristoteles University of Thessaloniki, Thessaloniki, Greece
[c] Joint Institute for Nuclear Research, Dubna

We applied the Principal Component Analysis (PCA), namely, the «Caterpillar»-SSA approach [1, 2], to the network traffic measurements. This approach proved to be very efficient for understanding the main features of terms forming the network traffic. The statistical analysis of leading components has demonstrated that even a few first components form the main part of information traffic. The residual components play the role of small irregular variations which do not fit in the basic part of the network traffic and can be interpreted as a stochastic noise. Based on the feature characteristics of residual components, we developed a statistical method for the selection and elimination of residuals from the whole set of principal components.

Мы применили метод главных компонентов, а именно подход «Caterpillar»-SSA, для анализа измерений информационного трафика. Этот подход оказался очень эффективным для понимания основных особенностей членов, ответственных за формирование сетевого трафика. Статистический анализ главных компонентов показал, что уже несколько первых компонентов формируют основную часть информационного трафика. Остаточные компоненты, которые не следуют основному закону сетевого трафика, могут быть интерпретированы как случайный шум. Используя характерные особенности остаточных компонентов, мы разработали статистический метод, с помощью которого можно отбирать такие компоненты с целью их исключения из всего набора главных компонентов.

## INTRODUCTION

We applied [3] a nonlinear time series analysis [4] to the traffic measurements obtained at the input of the intermediate-size Local Area Network (LAN). We have demonstrated that nonlinear techniques can be successfully used for a deeper understanding of main features of the traffic data. At the same time, we found that due to a very complicated character of traffic series the traditional algorithms of nonlinear analysis do not give reliable estimations of the analyzed time series. For instance, the Grassberger–Procaccia algorithm gives a very high dimension for original traffic measurements. However, after filtering out a high-frequency component, which can be considered as a noise, we obtained a more realistic result for the embedding dimension of the underlying process. This result has been confirmed independently by the Principal Component Analysis (PCA) method [3] in the framework of the «Caterpillar»–Singular Spectrum Analysis (SSA) approach [1, 2].

The Principal Component Analysis is a well-known technique in multivariate data analysis [5–9]. The PCA method consists in applying a linear transformation of the original data

space into a *feature space*, where the data set may be represented by a reduced number of «effective» features and yet retains most of the intrinsic information content of the data. The «Caterpillar»-SSA approach is a novel scheme, which is very efficient for the analysis of time series corresponding to any arbitrary signal [1, 2].

In our study we use the traffic measurements obtained at the input of Dubna University [10] LAN, which includes approximately 200–250 interconnected computers. We describe in Sec. 1 the data acquisition system of this LAN, realized on the basis of a standard PC. In Sec. 2 we present the basic concept of the «Caterpillar»-SSA scheme. In Sec. 3 we apply this technique to the traffic measurements and analyze leading components responsible for the main part of the network traffic. In Sec. 4 we study the residual components and propose a statistical method for their selection and elimination from the whole set of principal components.

## 1. DATA ACQUISITION SYSTEM

The measurements of network traffic have been realized at the external side of the input lock of LAN. The performance of the data acquisition system is based on realization of an open mode driver [11] (see Fig. 1).
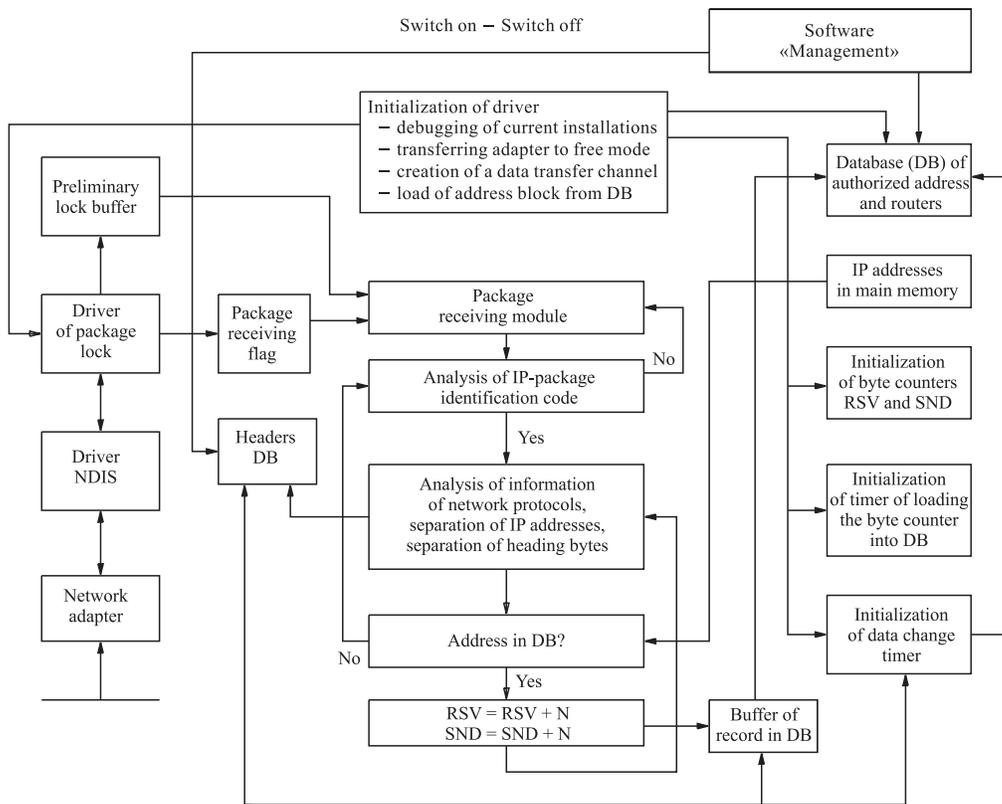


Fig. 1. Scheme of a data acquisition system

In standard conditions the network adapter of a computer is in a mode of detecting a carrying signal (main harmonic 4–6 MHz). After appearance in the cable bits of the package preamble, the network adapter comes to a mode of 1 bit and 1 byte synchronization with the transmitter and starts receiving first bytes of the package heading. As soon as one succeeds in extracting the MAC-address of the shot receiver from the first bytes taken by the adapter, the network adapter compares it to its own. If the result of the comparison is negative, the network adapter ceases to record the shot's bytes into its internal buffer and cleans its contents and then waits until the next package appears.

In order to provide conditions for reception and analysis of all the packages transmitted over the network, it is necessary to move the adapter devices to a free mode when all possible shots are recorded in the buffer. This operation is executed through the instructions of the NDIS driver.

The free-mode driver records the accepted packages in the preliminary capture buffer and displays the flag of receiving the package. Then the receiving package module is activated and analysis of the margin of the package's type is carried out to extract TCP/IP packages from the whole stream.

After identification it is possible to separate and delete the data block as well as to record the headers to the SQL-server database. The recording is performed together with the time data with a frequency up to 10 kHz. Although the recording is performed with buffering, the mode of saving the packages' headers requires enormous server's resources, as in this case there is a permanent procedure of recording with small portions to the hard disk. That is why this mode is switched on if required at the management system's instruction.
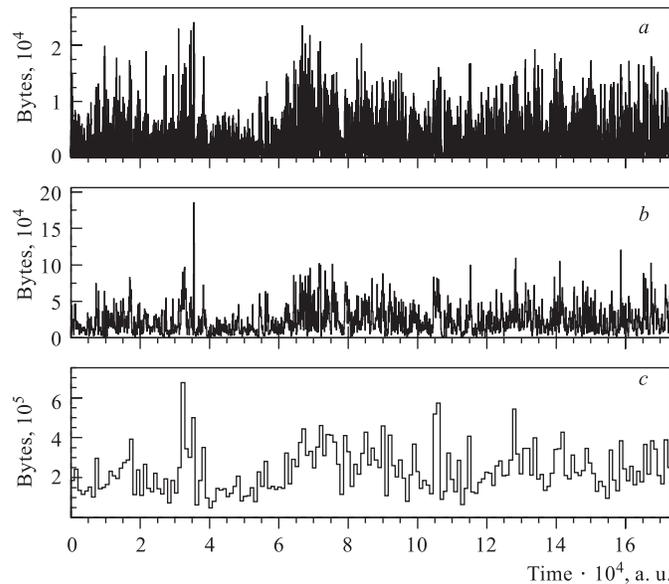


Fig. 2. Traffic measurements aggregated with different bin sizes: 0.1 s (*a*), 1 s (*b*) and 10 s (*c*)

The system also provides control over the external traffic of the local area network on the basis of checking records in the router table. Initial information on the legal IP addresses is

saved in the database of the LAN computers from which data on legal addresses are loaded into the main memory array. The users which do not participate in forming the external traffic are not taken into account when calculating the number of transferred and received bytes. In order to decrease the number of sessions of recording the information on the external traffic in the database, a timer of load out of the buffer and a timer of changing a current date have been introduced into the system.

The recorded traffic data correspond approximately to 20 h (1 600 000 records with a frequency up to 10 kHz, which corresponds to 1-ms bin size) of measurements. The part of this series corresponding approximately to one hour of measurements and aggregated with different bin sizes is presented in Fig. 2. Two protocols are used in the «Dubna» LAN. The NetBEUI protocol is applied only for internal exchanges, and the TCP/IP for external communications. The contribution of the NetBEUI traffic has been estimated around 1–6 packages per second during daily working hours. This is negligibly small compared to the TCP/IP traffic. In this connection, we may neglect the influence of non-IP traffic on the TCP/IP traffic.

## 2. BASIC CONCEPT OF THE «CATERPILLAR»-SSA TECHNIQUE

The «Caterpillar»-SSA approach [1,2] is applied to the analysis of time series corresponding to any arbitrary signal $f(t)$, with $t > 0$ determined in equidistant points. The basic «Caterpillar»-SSA scheme includes four main steps:

1) transformation of one-dimensional series into multidimensional form,
2) singular value decomposition of the multidimensional series,
3) principal component analysis and selection of feature components,
4) reconstruction of one-dimensional series using the selected components.

The transformation of one-dimensional series

$$x_i = f(t_i) = f[(i-1)\Delta t], \quad i = 1, 2, \ldots, K \tag{1}$$

into a multidimensional series is realized by representing (1) in matrix form:

$$X = (x_{ij})_{i,j=1}^{k,L} = \begin{pmatrix} x_1 & x_2 & x_3 & \ldots & x_L \\ x_2 & x_3 & x_4 & \ldots & x_{L+1} \\ x_3 & x_4 & x_5 & \ldots & x_{L+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_k & x_{k+1} & x_{k+2} & \ldots & x_K \end{pmatrix}, \tag{2}$$

where $L < K$ is called the caterpillar length and $k = K - L + 1$.

Then the eigenvalues $\lambda_i$ $(i = 1, 2, \ldots, L)$ and eigenvectors $\mathbf{V}_i$ $(i = 1, 2, \ldots, L)$ of the covariance matrix $C = (1/k)XX^T$ are determined. The matrix of eigenvectors $V$ is used for transition to principal components:

$$Y = V^T X = (Y_1, Y_2, \ldots, Y_L), \tag{3}$$

where $Y_i$ $(i = 1, 2, \ldots, L)$ are rows of $k$ elements.

The equality

$$\sum_{i=1}^{L} \frac{\lambda_i}{L} = \sum_{i=1}^{L} \alpha_i = 1$$

permits one to estimate the contribution $\alpha_i$ (in decreasing order) of the $i$th principal component into the analyzed series.

This contribution can be interpreted as fraction of information related to a single component, and it helps, together with analytical and visual analysis of eigenvectors and principal components, to select feature components for reconstruction of one-dimensional series. Usually the selection of specific components depends on a goal which we pursue and the information content of particular components (see, for example, [12–14, 1, 2]).

## 3. PCA OF TRAFFIC MEASUREMENTS: ANALYSIS OF LEADING COMPONENTS

The «Caterpillar»-SSA approach foresees a preliminary centering and normalization of time series to be analyzed. Depending on the character of the series this procedure can be applied separately, altogether or even not applied at all. Unfortunately, the authors of [1, 2] do not propose exact recommendations for such a procedure. Based on our preliminary analysis, the traffic series has been centered but not normalized.

The caterpillar length (or window) $C_L$ has been chosen based on the analysis of the autocorrelation function for traffic measurements [3]. In this study we used different values of $C_L$, starting from the minimal value $C_L = 12$ up to $C_L = 20$.

Figure 3 shows part of the daily traffic measurements aggregated with the bin size 1 s, which has been used in this study.
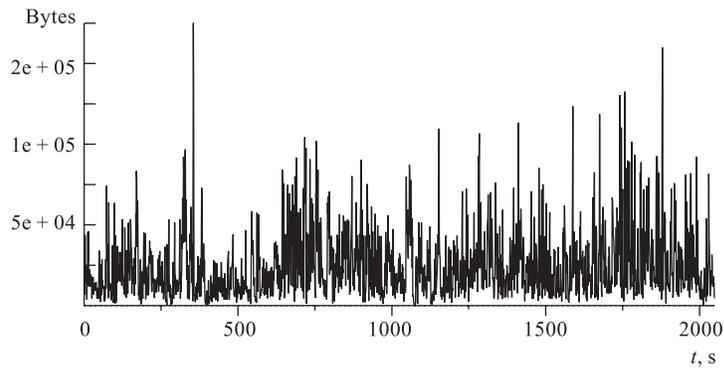


Fig. 3. Traffic measurements aggregated with the bin size 1 s

One of the main results of the application of the «Caterpillar»-SSA technique to the analyzed series is presented in Fig. 4, where the contribution of the eigenvalues on a percentage basis for $C_L = 12$ and 20 is shown. This information permits one to estimate the number of principal components, which effectively contribute to the analyzed series.
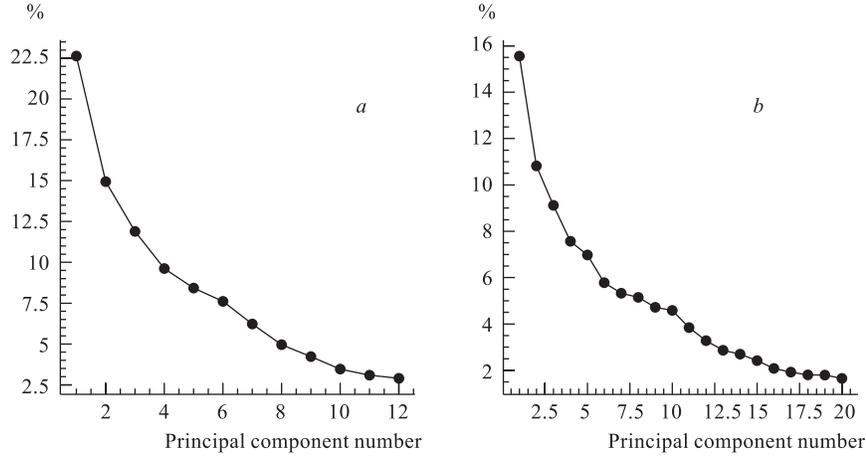
Fig. 4. Contributions of eigenvalues on a percentage basis for original traffic measurements: $C_L = 12$ (*a*) and $C_L = 20$ (*b*)

*Table* 1. **Results of fitting the packet size distributions, corresponding to** $N$ **leading components, by the log-normal function (2)**

| $N$, leading comp. | $\sigma$ | $\mu$ | $\nu$ | $\chi^2$ |
|---|---|---|---|---|
| 1 | $0.273 \pm 0.009$ | $10.44 \pm 0.01$ | 47 | 87.49 |
| 2 | $0.304 \pm 0.005$ | $10.40 \pm 0.01$ | 44 | 66.82 |
| 3 | $0.349 \pm 0.007$ | $10.38 \pm 0.01$ | 47 | 53.10 |
| 4 | $0.377 \pm 0.008$ | $10.37 \pm 0.01$ | 47 | 63.52 |
| 5 | $0.420 \pm 0.011$ | $10.35 \pm 0.01$ | 47 | 68.50 |
| 6 | $0.432 \pm 0.012$ | $10.34 \pm 0.01$ | 46 | 59.12 |
| 7 | $0.426 \pm 0.008$ | $10.35 \pm 0.01$ | 47 | 49.03 |
| 8 | $0.444 \pm 0.007$ | $10.34 \pm 0.01$ | 47 | 34.39 |
| 9 | $0.463 \pm 0.008$ | $10.33 \pm 0.01$ | 43 | 38.94 |
| 10 | $0.482 \pm 0.009$ | $10.32 \pm 0.01$ | 47 | 37.76 |
| 11 | $0.489 \pm 0.008$ | $10.31 \pm 0.01$ | 47 | 55.64 |
| 12 | $0.500 \pm 0.009$ | $10.32 \pm 0.01$ | 47 | 59.00 |
| 13 | $0.506 \pm 0.008$ | $10.32 \pm 0.01$ | 43 | 51.97 |
| 15 | $0.518 \pm 0.009$ | $10.31 \pm 0.01$ | 46 | 55.16 |
| 17 | $0.516 \pm 0.008$ | $10.30 \pm 0.01$ | 47 | 78.59 |
| 19 | $0.513 \pm 0.008$ | $10.30 \pm 0.01$ | 44 | 101.6 |

Taking into account [16], it is reasonable to assume that the packet size distributions, corresponding to leading components, may be described by the log-normal distribution. In order to check whether these distributions follow the log-normal form, we fitted them by the log-normal function [17]:

$$f(x) = \frac{A}{\sqrt{2\pi}\sigma} \frac{1}{x} \exp\left[ -\frac{1}{2\sigma^2} (\ln x - \mu)^2 \right], \tag{2}$$

where $\sigma$ and $\mu$ are parameters and $A$ is a normalizing factor. The fitting procedure has been realized with the help of the MINUIT package [18] in the framework of the well-known Physical Analysis Workstation (PAW, see details in [19]).

We present in Table 1 the results of fitting the packet size distributions, corresponding to different number $N$ of leading components (the results presented here are for $C_L = 20$), by function (2). Here $\nu$ is the number of degrees of freedom for the $\chi^2$ test.



Fig. 5. The dependence of $\chi^2/\nu$ versus the number of leading components $N$

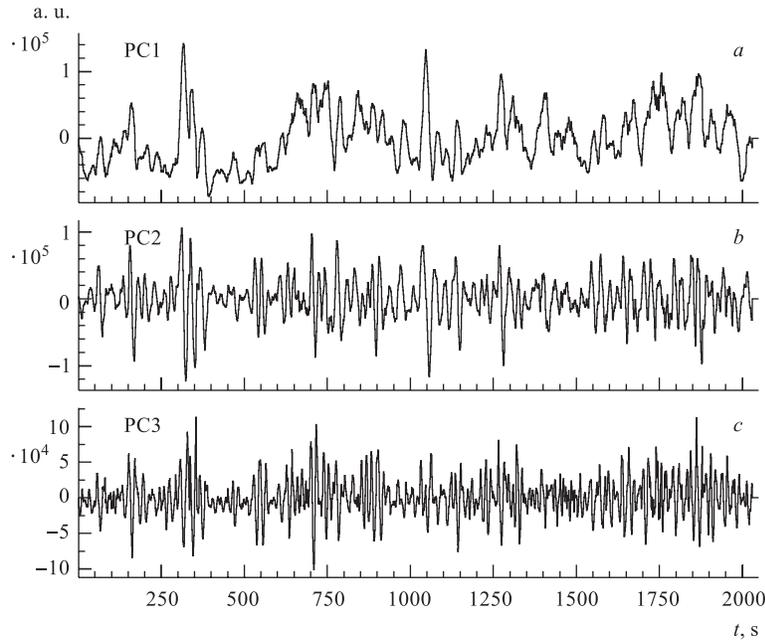Fig. 6. Fitting the distribution corresponding to $N = 3$ leading components by function (2)



Fig. 7. Time series corresponding to three leading components (after subtraction of the caterpillar average value)

Figure 5 shows the dependence of $\chi^2/\nu$ on $N$ (for $C_L = 20$). Two lines parallel to the abscissa axis show the significance levels (or the probability that the observed chi squared will exceed the value $\chi^2$ by chance even for a correct model: see, for instance, [15, 17]) $\alpha = 10\%$ (the top line, $\chi^2/\nu = 1.247$) and $\alpha = 89.5\%$ (the bottom line, $\chi^2/\nu = 0.732$) corresponding to the $\chi^2$ test for $\nu = 47$.

This dependence demonstrates that the testing distribution does not pass the null hypothesis (2), when only the first leading component is taken into account. Then, with the increase of $N$, the value of $\chi^2$ is rapidly decreasing, and for $N = 3$ one can see a quite good level of correspondence ($\alpha = 22\%$) of the distribution to the null hypothesis (Fig. 6).
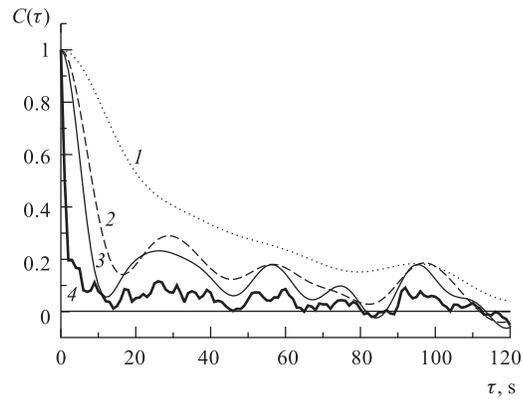


Fig. 8. Autocorrelation function $C(\tau)$ of reconstructed series corresponding to a different number of leading components: *1* — one leading component; *2* — two leading components; *3* — three leading components; *4* — original data
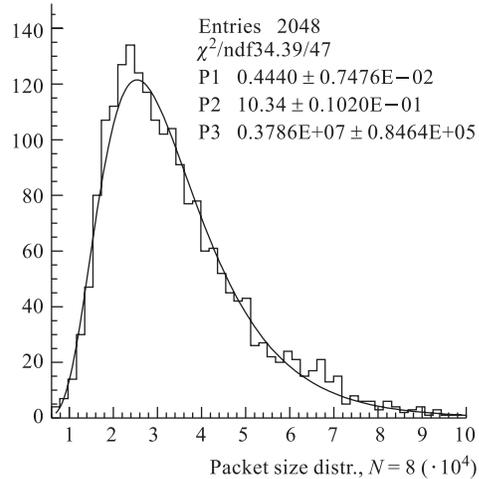
This result is of great interest because only three first components (of 20) form the fundamental part of the information traffic. Figure 7 shows the series reconstructed on the basis of the first, second and third leading component, correspondingly, after the subtraction of the caterpillar average value. Figure 8 presents the dependence of the autocorrelation function

$$C(\tau) = \frac{\sum_{i=1}^{K}(x_{i+\tau} - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^{K}(x_i - \bar{x})^2}, \quad \bar{x} = \frac{1}{K}\sum_{i=1}^{K}x_i.$$

(3)

One can see from these figures that the autocorrelation function corresponding to the sum of three leading components is close to the autocorrelation function for the original data. Their summary contribution to the general dispersion is around 40% (see Fig. 4 for $C_L = 20$).

This result has been confirmed for the shorter caterpillar length, $C_L = 12$. In this case only two leading components, their lump contribution approximately coinciding with the contribution



Fig. 9. Fitting the distribution corresponding to eight leading components by function (2)

of the three leading components for $C_L = 20$ (see Fig. 4), reproduce the log-normal form of the traffic.

Further increase of $N$ leads to unexpected increase of $\chi^2$ (for $N = 4$ and 5) together with the decrease of the significance level below 10%. Then the value of $\chi^2/\nu$ rapidly

decreases and reaches its record minimal value 0.732 for $N = 8$. The corresponding statistical distribution is presented in Fig. 9. It demonstrates both a very good level of correspondence of the reconstructed distribution to the null hypothesis ($\alpha = 89.5\%$) and a reliable accuracy of approximation for all regions of the analyzed distribution. The summary contribution of eight leading components into the general dispersion is around 66%.

Figure 10 shows the reconstructed series using the «Caterpillar»-SSA method (for $C_L = 20$) on the basis of eight leading components. One can clearly see that it reproduces characteristic features of the original series presented in Fig. 3.
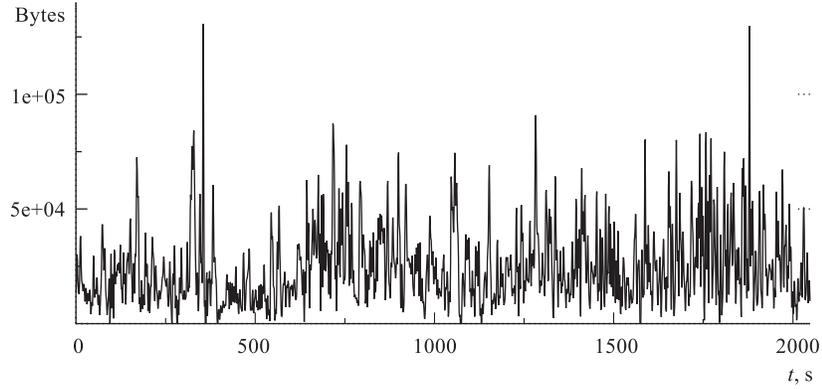


Fig. 10. Traffic series measurements reconstructed by the caterpillar method (for $C_L = 20$) on the basis of eight leading components

## 4. PCA OF TRAFFIC MEASUREMENTS: ANALYSIS OF RESIDUAL COMPONENTS

In the region of large $N$ there is a growth of $\chi^2$, especially noticeable at $N \geqslant 15$ (see Fig. 5). Such a tendency may be caused by the influence of the residual components related to small irregular variations, which do not fit in the basic model of the network traffic (2) and can be interpreted as stochastic noise.

Figure 11 shows a series reconstructed on the basis of the smallest residual component, namely, component 20. One can clearly see that this series has a significantly different character compared to the original traffic measurements. It looks like a nonstationary process symmetric against zero mean value.

Figure 12 shows the statistical distribution corresponding to the series presented in Fig. 11. It quite well follows the Gaussian distribution which is confirmed by the $\chi^2$ test (see Fig. 12). The autocorrelation function of the corresponding series shows that it behaves like noise.

However, when increasing the number of residual components, their summary distribution quickly starts losing the symmetric form together with growth of correlations between the series terms.

In order to estimate the amount of residual components that can be eliminated from the original time series without influence on its fundamental part, we divide all principal components into two parts:
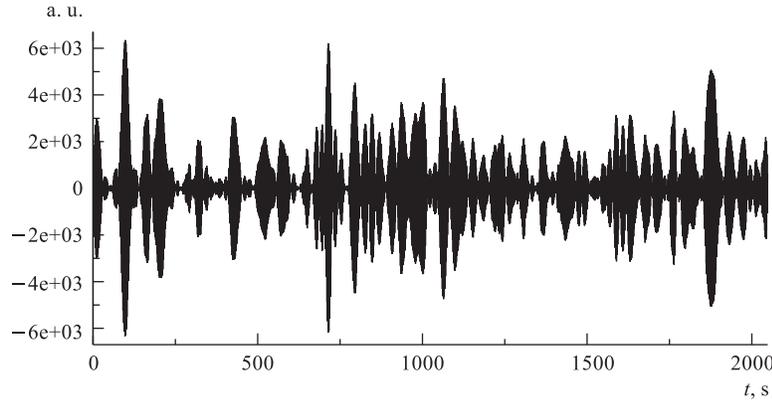
Fig. 11. Traffic series reconstructed by the caterpillar method ($C_L = 20$) on the basis of the smallest component

● first part corresponding to the leading components and responsible for the log-normal form of the packet size distribution,

● second part related to residual components, which is described by a symmetric statistical distribution and behaves like a stochastic noise.

As criterion for selection of the second part we used the «moment» of the symmetry violation for the series corresponding to the residual components. A well-known sign test has been used for testing the symmetry against zero of residual distributions. The sign test has the following form:

$$\mu = \sum_{i=1}^{n} \Theta(X_i), \qquad (4)$$

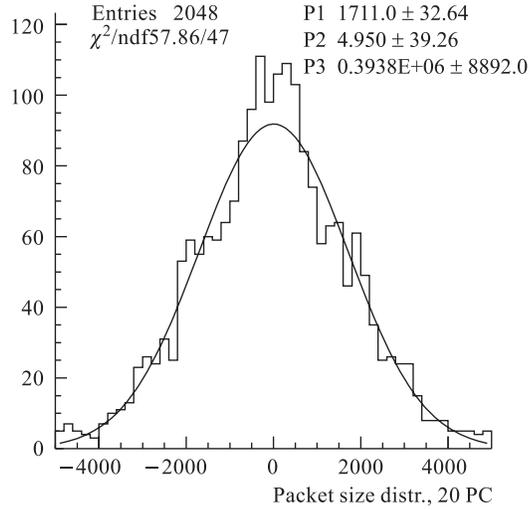where $X_1, \ldots, X_n$ are observables; $n$ is the sample size, and $\Theta$ is the Heaviside function:



Fig. 12. Statistical distribution of the time series presented in Fig. 11; the fitting curve corresponds to the Gaussian distribution

$$\Theta(x) = \begin{cases} 1, & x > 0, \\ 0, & x \leqslant 0. \end{cases}$$

When the null hypothesis is valid, the $\mu$ distribution is approximated (in case of large $n$) by

$$P\{\mu \leqslant m \mid n, p\} \approx \Phi\left( \frac{m - np + 0.5}{\sqrt{np(1-p)}} \right),$$

where $\Phi$ is the distribution function of the normal distribution; $p = 0.5$ and $n = 2048$ (in our case).

Figure 13 shows the dependence of the $\mu$ value versus the number of residual components (for caterpillar lengths 12 and 20). It is clearly seen that the $\mu$ value exceeds the reliable confidential level, when the number of residual components is greater than 6 for $C_L = 12$ and 11 for $C_L = 20$.
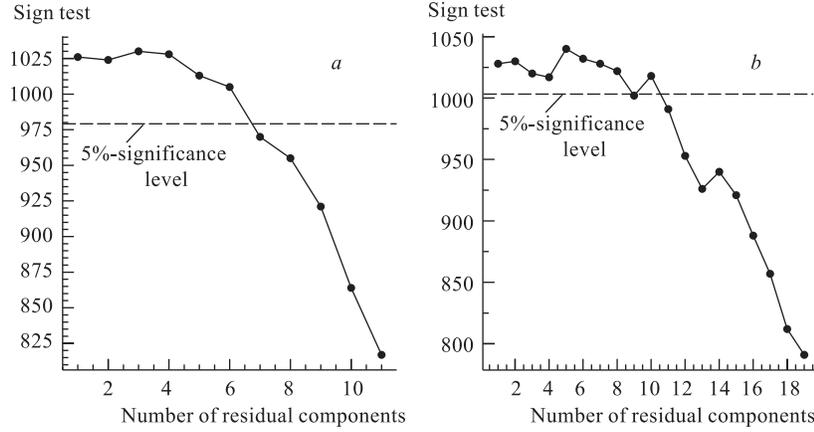


Fig. 13. The values of the sign test $\mu$ versus the number of residual components for the caterpillar length $C_L = 12$ (*a*) and $C_L = 20$ (*b*)

In order to confirm the results obtained by the sign test, we applied a more powerful criterion based on the $\omega_n^2$ statistics [20]. This criterion tests the symmetry against $x = 0$ of the distribution function $F(x)$ for observables $X_1, \ldots, X_n$, i.e., the null hypothesis $H_0$: $F(x) = 1 - F(x)$. The corresponding $\omega_n^2$ statistics has the form:

$$\omega_n^2 = n \int\limits_{-\infty}^{\infty} [F_n(x) + F_n(-x) - 1]^2 \, dF_n(x), \tag{5}$$

where $F_n(x)$ is the empirical distribution function. It is more convenient to calculate the values of the statistics (5), using the following formula:

$$\omega_n^2 = \sum_{j=1}^{n} \left[ F_n(-X_{(j)}) - \frac{n - j + 1}{n} \right]^2,$$

where $X_{(1)} \leqslant \ldots \leqslant X_{(n)}$ is the variational series constructed on the basis of observables.

Figure 14 shows the dependences of $\omega_n^2$ values on the number of residual components for two cases of the caterpillar length: $C_L = 12$ and 20. These dependences have distinct characteristic features at $k = 4$ for $C_L = 12$, and $k = 7$ for $C_L = 20$ (one can see that the number of such components approximately equals to one third of the caterpillar length), after which, when $k$ is increasing, there is a quick rise of $\omega_n^2$. This rise means that the residual
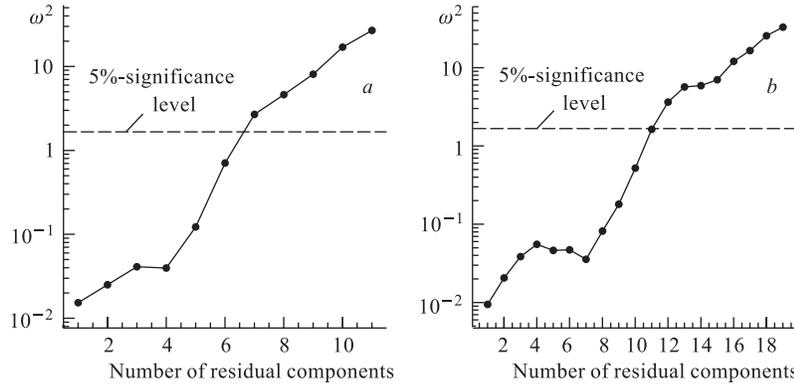
Fig. 14. The dependences of $\omega_n^2$ values on the number of residual components for two cases of caterpillar length: $C_L = 12$ (*a*) and $C_L = 20$ (*b*)

series loses its symmetric character, because in the second part the components responsible for the log-normality are involved.

One can see from Fig. 14 that the number of residual components $k = 6$ for $C_L = 12$ and $k = 11$ for $C_L = 20$ correspond to the 5%-significance level for the $\omega^2$ criterion. This coincides with the result obtained for the sign test (Fig. 13). These estimates of the number of components, which do not noticeably influence the fundamental part of traffic, qualitatively coincide with the result obtained in Sec. 3 applying the $\chi^2$ test (Fig. 5).

## CONCLUSION

We applied the «Caterpillar»-SSA approach [1, 2], which is an extension of the principal component analysis, to network traffic measurements, in order to understand the main features of the principal components of the network traffic. Our analysis of the leading components has shown that a few first components alone form the fundamental part of the information traffic and that the correspondence to the log-normal distribution remains valid up to intermediate values of $N$. In the region of large $N$ there was found a noticeable growth of $\chi^2$, which can be explained by the influence of residual components related to small irregular variations. Based on feature characteristics of residual components, we developed a statistical method that permits one to estimate the number of components which do not play a noticeable role in the fundamental part of traffic and can be eliminated from the whole set of components.

Thus, the statistical analysis of traffic measurements based on the joint application of $\chi^2$ and $\omega^2$ tests gives the possibility of splitting the whole set of components into two classes. The first class includes the leading components responsible for the main contribution to the traffic, and the second class involves residual contributions that can be interpreted as noise. A more detailed analysis of the boundary region between these two groups may provide additional information on traffic components and thus simplify the understanding of traffic dynamics.

## REFERENCES

1. Principal Components of Time Series: Caterpillar Method / Eds. D. L. Danilov, A. A. Zhigljavsky. St. Petersburg, 1997 (in Russian).

2. *Golyandina N., Nekrutkin V., Zhigljavsky A.* Analysis of Time Series Structure: SSA and Related Techniques. Chapman & Hall/CRC, 2001.

3. *Akritas P. et al.* Nonlinear Analysis of Network Traffic // Chaos, Solitons & Fractals. 2002. V. 14(4). P. 595–606.

4. *Abarbanel H. D. I.* Analysis of Observed Chaotic Data. Springer-Verlag New York, Inc., 1996.

5. *Preizendorfer R. W.* Principal Component Analysis in Meteorology and Oceanography. Elsevier, 1988.

6. *Jolliffe I. T.* Principal Component Analysis. Springer-Verlag, 1986.

7. *Jackson J. E.* A User's Guide to Principal Component Analysis. N. Y., 1992. P. 26–62.

8. *Karhunen K.* Uber lineare methoden in der Wahrscheinlichkeitsrechnung // Ann. Acad. Scient. Fennicae. Series A1: Mathematica–Physica. V. 37. P. 3–79 (Transl.: RAND corp. Santa Monica, CA, 1960. Rep. T-131).

9. *Loéve M.* Probability Theory. 3rd ed. N. Y.: Van Nostrand, 1963.

10. The State University «Dubna»: http://www.uni-dubna.ru

11. *Vasiliev P. V. et al.* System for Acquisition, Analysis and Control of Network Traffic for the JINR Local Network Segment: the «Dubna» University Example. JINR Commun. D11-2001-266. Dubna, 2001.

12. *Broomhead D. S., King G. P.* Extracting Qualitative Dynamics from Experimental Data // Physica D. 1986. V. 20. P. 217–236.

13. *Broomhead D. S., King G. P.* Time-Series Analysis // Proc. Roy. Soc. London. 1989. V. 423. P. 103–110.

14. *Vautard R., Yiou P., Ghil M.* Singular Spectrum Analysis: A Toolkit for Short, Noisy Chaotic Signals // Physica D. 1992. V. 58. P. 95–126.

15. *Press W.H. et al.* Numerical Recipes in C: The Art of Scientific Computing. 2nd ed. Cambridge University Press, 1988; 1992.

16. *Antoniou I. et al.* On the Log-Normal Distribution of Network Traffic // Physica D. 2002. V. 167. P. 72–85.

17. *Eadie W. T. et al.* Statistical Methods in Experimental Physics. Amsterdam; London: North-Holland Pub. Comp., 1971.

18. *James F.* MINUIT — Function Minimization and Error Analysis. Reference manual, version 94.1. CERN Program Library D506. 1998.

19. *Brun R. et al.* PAW — Physics Analysis Workstation. CERN Program Library Q121. 1989.

20. *Martinov G. V.* Omega-Squared Criteria. M.: Nauka, 1978 (in Russian).