# Modeling local stationary behavior of Internet traffic

D. Moltchanov

*Abstract*—**Non-stationary behavior of aggregated IP traffic patterns was demonstrated in a number of studies. However, none of those did either consider practical aspects of this phenomenon or propose suitable model to capture it. Searching for model for IP traffic aggregates we introduce the concept of local stationarity and demonstrate that it allows to model traffic patterns measured in high-speed operational networks. The proposed model is on-line in nature and suitable for real-time estimation of the traffic state in terms of piecewise covariance stationary stochastic process. As a basic tool of the model we use change-point statistical test allowing us to dynamically and automatically determine whether statistical characteristics of the traffic pattern changes and, if so, estimate new parameters of the traffic pattern. We provide numerical examples and discuss applications of the proposed model that include but not limited to dynamic resource reservation, routing with guaranteed bandwidth, etc.**

## I. Introduction

In the last decade the Internet has experienced tremendous growth in the number of users that resulted in dramatic increase of the traffic volume carried by the network. The nature of Internet traffic has been a hot research topic since then. One of the most important findings was that traffic aggregates in IP networks may exhibit high degree of variability [1]. High variability manifests itself requiring significantly more network resources in order to serve the traffic with given performance metrics compared to classic Poisson traffic. It means that there must be enough resources to serve very large bursts in traffic pattern. However, there are also long time spans during which local average of the traffic aggregate may stay well below the mean of the whole process.

High-variability of traffic observations may stem from a number of phenomena including self-similar, long-range dependent and/or non-stationary behavior. Despite of significant efforts in the area there is still no common agreement which of these properties are fundamental for a traffic to exhibit high variability. Since the discovery of long-range dependence and self-similarity in Internet traffic observations almost all long-term variations in the mean value of traffic observations were attributed to their effect. However, high variability can also be caused by non-stationary behavior of observations when one of more statistical characteristics gradually or abruptly change in time.

Non-stationarity of traffic aggregates was only recently started to be considered as a possible reason for high variability of aggregated traffic patterns. However, it was known for a while that aggregated traffic is characterized by deterministic trends similar to those found in telephone traffic. For example,

it was demonstrated that there are clear daily variations in the traffic patterns [2]. It was also noticed in [2] that there is a clear indication of busy hour in link usage patterns. When statistical characteristics change in time, non-stationarity can be important factor contributing to high variability of traffic aggregates [3].

Although non-stationary behavior of aggregated IP traffic patterns was demonstrated in many studies none of those discussed potential implications of this phenomena neither proposed practical models to capture this behavior. Contrarily to those studies, in this paper we consider the problem of modeling traffic aggregates in modern IP networks. We do not assume stationarity for traffic aggregates as it rarely hold in practice. We justify our assumptions observing statistical characteristics of real traffic traces. Attempting to model the traffic we introduce the concept of local stationarity and demonstrate how to exploit it to describe traffic aggregates on-the-fly. We also discuss applications scenarios and limitations of the proposed model.

The rest of the paper is organized as follows. In Section II we discuss results and highlight recent development in the area of broadband traffic analysis and modeling. Next, in Section III, we consider statistical characteristics of real traffic traces highlighting their non-stationary nature. Model for covariance stationary behavior of traffic observations is presented in Section IV. In Section V this model is extended to piecewise covariance stationary behavior. Numerical examples are demonstrated in Section VI. Usage of the model in traffic engineering applications is discussed in Section VII. Conclusions are drawn in the last section.

## II. Related Work

### A. Long-range Dependence and Self-similarity

In the middle of $90th$ non-stationarity was not considered as a possible cause of high-variability of traffic observations. Many authors either implicitly or explicitly assumed stationarity for Internet traffic observations. Self-similarity, long-range dependence and heavy-tail distributions have all been claimed to be fundamental properties of aggregated traffic [4]. All these three notions are related to each other. A process is said to be asymptotically second-order self-similar when the structure of its ACF is preserved under time-aggregation. A process is said to be long-range dependent when its ACF is not summable meaning that dependence continues to infinity. Note that asymptotically second-order self-similar processes are always long-range dependent implying non-summable ACF [4]. In this case, descriptor of the degree of self-similarity, Hurst parameter $(0.5 < H < 1)$, is also meaningful for long-range dependent processes. Heavy-tailed distributions (e.g. Pareto) or nearly heavy-tailed distributions (e.g. Weibull) are

claimed to be the main source for self-similar behavior of traffic aggregates [5].

It is important to note that observations of long-range dependent processes are characterized by the so-called 'wave' behavior where observations tend to be higher or lower than the mean of the process for long durations of time. We note that such behavior may also be inherent for observations of non-stationary processes when the mean changes in time. In practice we always deal with limited number of observations. As a result, it is difficult to distinguish between observations exhibiting second-order self-similar and non-stationary behaviors [3].

There are a number of estimators used to detect whether given observations are taken from self-similar process or not. Variance considered as the function of the level of traffic aggregation is often used to determine degree of self-similarity. Well-known rescaled-adjusted (R/S) and variance-time statistics are based on this relation. Periodogram-based and and Whittle's estimators are based on spectrum estimates. In practice, applying these estimators to traffic observations one gets a range of Hurst parameters that may significantly vary. Wavelet-based Abry-Veicth estimator was recently claimed as being reliable even in presence of non-stationarity [6]. In [7] authors compared performance of Abry-Veicth, Periodogram-based, Whittle and R/S estimators for artificially generated self-similar traffic. They demonstrated that none of these estimators provide reliable estimate of Hurst parameter. Moreover, it was shown that Abry-Veicth estimator sometimes behaved worse compared to others resulting in absurd results.

We note that Hurst parameter is the most important parameter for modeling self-similar traffic. Moreover, it is self-similarity that was often claimed to produce a significant impact on performance of the traffic service process in the network. As a result, small deviations in the Husrt parameter may lead to different performance predicted using self-similar models provided to traffic aggregates. We believe that modeling of traffic observations exhibiting high degree of variability is still an open problem.

### B. Non-stationarity

Recently, many authors stated to question whether self-similarity is the only reason for high variability of traffic aggregates. In [2], [8] it was demonstrated that traffic pattern has almost deterministic daily variations resulting in clear non-stationary behavior on a day timescale. Authors in [9] demonstrated that multiplexed traffic on a high-speed link may have non-stationary behavior and discussed possible causes of non-stationarity of traffic observations. They argue that this could be due a number of reasons including time-varying number of aggregated sources, routing changes, specific aggregation of constant number of stationary sources. They suggested that at multi-seconds timescales aggregated network traffic is characterized by piecewise-linear non-stationary behavior with possible long-range dependent features. Authors in [10] studied interarrival times of packets at high-speed links. They claim that packet and connection arrival processes may exhibit non-stationary behavior arguing that time-varying nature of the

number of active sources is responsible for variations in traffic patterns. It was also reported in [11] that approximately $10\%$ of commercial Internet routes have lifetimes of few hours of even less. This could result in abrupt changes in the number of sources bottlenecked at the link.

Modeling traffic aggregates in IP networks many authors did not perform any tests for stationarity assuming that their traces are realizations of weakly stationary stochastic process [1], [12], [13], [5]. Other authors consider non-stationary behavior of traffic observations as the factor that may hamper a clear conclusion about properties of traffic observations [3], [9], [10]. The common approach was to choose a sufficiently small blocks of observations such that observations in separate blocks are expected to be at least weakly stationary. For example, when testing for applicability of the Gaussian model to traffic modeling authors in [14] neglected a part of their trace claiming that it may introduce 'undesirable' non-stationary behavior. Authors in [10] assumed that $5$ minute blocks of their traffic observations is sufficient to ensure intra-block stationarity. Analyzing the loss process of IP packets Yajnik *et al.* [15] carried out a simple test for stationarity. They assumed that loss observations are realization of the weakly stationary process when they are within $E[Y] \pm 0.05$, where $E[Y]$ is the mean calculated using sliding window of $2000$ observations. Unfortunately, this procedure results in high number of false indications of non-stationarity and does not provide any statistical guarantees.

We note that testing for stationary behavior having limited number of observations is extremely complicated task and no general procedures are available. To our knowledge [11] was the first paper where authors implicitly tried to represent non-stationary behavior of the traffic using piecewise non-stationary process. To discriminate between covariance-stationary segments authors proposed to use change-point detection algorithm. The main problem is that their test explicitly assumes that the time-series are uncorrelated at all lags. This assumption may not hold in practice. Authors in [9] also used change-point estimators to determine boundaries of stationary segments. The another pitfall of approaches taken in [11], [9] is off-line nature of change detection algorithms. This limits usability of these tests to off-line analysis of data.

### C. Marginal Distribution of Aggregated Traffic

Another questions which is strictly related to the traffic nature is the marginal distribution. Inspired by the central limit theorem, a number of authors argues for Gaussian behavior of traffic aggregates. For example, in [16] authors highlighted that aggregation of a large number of independent sources leads to Gaussian. In [17] it was suggested that aggregation of even fairly small number of sources may lead to Gaussian behavior. Authors in [14] tested Gaussian approximation for traffic aggregates. They demonstrated that for this assumption to hold a certain traffic aggregate should have sufficient vertical (number of sources) or horizontal (timescale) aggregation.

On the contrary, there are a number of publications stating that packet counts tend to have non-gaussian marginal distribution (see [13], [5] among others). First of all, it was noted in

[14] that Gaussian distribution can only serve as approximation of empirical traffic data as it always has probability mass on negative axis. Secondly, starting from Leland *et al.* [18] heavy-tailed marginal distribution was conventionally considered to be inherent for traffic observations. However, another reason for marginal distribution to be heavy-tailed is non-stationary nature of traffic aggregates. Indeed, even small increase in the mean value may lead to long tail of the distribution.

For network operators and end users it does not matter what mathematical model the traffic actually follows – they have to deal with negative effects of packet losses. The aim of this paper is not to prove that Internet traffic on high-speed links is stationary or not but to propose a tool to deal with time-varying effects of aggregated arrival processes. Note that these effects can be due to a very different phenomena. For example, it could be due to self-similar properties of arriving traffic contributing more traffic than average during some duration of time or due to change in the mean value of observations. The overall task is to estimate the *current state* of the traffic in terms of the model. This model can be further used by network equipment to take appropriate actions.

## III. STATISTICAL CHARACTERISTICS

In this paper, we use traffic traces from NLANR passive measurement and analysis (PMA) project [19]. Since our main focus is on real-time traffic aggregates we consider UDP traffic only. To ensure that the required level of horizontal aggregation is achieved the timescale of interest was set to 1 second. To demonstrate and analyze statistical characteristics we choose day 2 and day 3 Auckland VIII traces. Our choice is motivated by usage of these traces in related studies. We check that observations and conclusions stated in this section hold for other traces from PMA archive.

### A. 24-*hours Traces*

Let $\{Y(k), k = 0, 1, \ldots, N\}$ denote empirical observations of the number of bytes seen at the link in successive intervals of one second length. Time-series of day 2 and day 3 Auckland VIII traces are shown in Fig. 1. According to the classic approach adopted in analysis of traffic patterns we consider histogram of relative frequencies of bytes in time intervals of constant length and corresponding ACF.

Normalized ACFs (NACF) and histograms are shown in Fig. 2. Note that both histograms have a bimodal structure. This structure comes from long durations of time when local mean stays above or below the mean of the whole trace and may serve as an indicator of possible non-stationary behavior. For sufficiently large lags NACFs of both traces are well above zero implying that observations may have long memory. However, at some instants of time NACFs become zero. Then, for Day 3 trace NACF takes on negative values. This may also serve as an indicator of non-stationary behavior where mean of the trace varies in time.

To get visual impression how statistical characteristics of traffic aggregates vary in time we consider cumulative sum statistics (CUSUM, [20]), $C_Y(k) = \sum_{i=0}^{k}(Y(i) - \mu_Y)$, where $\mu_Y$ is the global mean of observations. Usually, CUSUM statistics is plotted on chart. If during a period of time most of the values are greater that the mean of the whole trace the CUSUM statistics is increasing. Therefore, a segment of the CUSUM statistics with positive slope indicates a period where the values tend to be above the mean. Similarly, a segment with negative slope corresponds to a period of time where the values tend to be below the mean of the trace. If the slope of CUSUM statistics is constant in time the average of the process stays relatively the same. A sudden change in direction of the CUSUM statistics indicates a sudden change in the mean.

CUSUM charts for considered traces are shown in Fig. 3. As one observes, there are a number of drastic changes in considered traces and there are durations of time when average of the process remains relatively constant. This illustrates that these traffic aggregates follow non-stationary behavior. We also note that there are many smaller changes in CUSUM statistics that may impose non-stationary behavior on much shorter timescales.

### B. 1-*hour Traces*

Consider now statistical characteristics of two 1-hour sub-traces chosen from day 3 trace and shown in Fig. 4. Note that these subtraces were chosen completely arbitrarily. However, results presented later are inherent for any 1-hour trace from Auckland VIII set of traces. NACFs and histograms of these subtraces are shown in Fig. 5. Both histograms have heavy-tails. Histogram of subtrace 2 also has three modes. Observing corresponding time-series, one may conclude that the reason for this structure is non-homogeneity of samples rather than self-similar or long-range dependent nature of observations. NACFs are also characterized by complex behavior. Both traces have NACF values greater than zero for large lags. Then, NACFs drop below zero and stay negative for sufficient time.

CUSUM charts for each hour of day 3 trace are shown in Fig. 6. Note that mean values of traces often change even during hour-long time spans. On the other hand, there are long intervals during which the mean value remains constant. For example, CUSUM chart for 22:00-23:00 trace illustrated in Fig. 4(b) clearly demonstrates that there are five intervals with different mean values. These are $0 - 920$, $920 - 1180$, $1180 - 2230$, $2230 - 2890$ and $2890 - 3600$. Obviously, these changes cannot be attributed to long-range dependent or self-similar behavior.

### C. Piecewise Stationary Behavior

Observing Fig. 1, Fig. 4 and Fig. 6 one may suggest that the traffic may experience so-called piecewise stationary behavior. At this point we can only use visual observations of traffic patterns and CUSUM chart to predict points at which statistical characteristics change. For example, for both subtraces in Fig. 4 stationary behavior is expected for first 600 observations. Consider statistical characteristics of these segments.

Histograms of chosen parts of subtraces and their approximations by normal distributions are shown in Fig. 7(a) and Fig. 7(b). These approximations suggest that the distribution of aggregated traffic patterns may converge to normal. NACFs
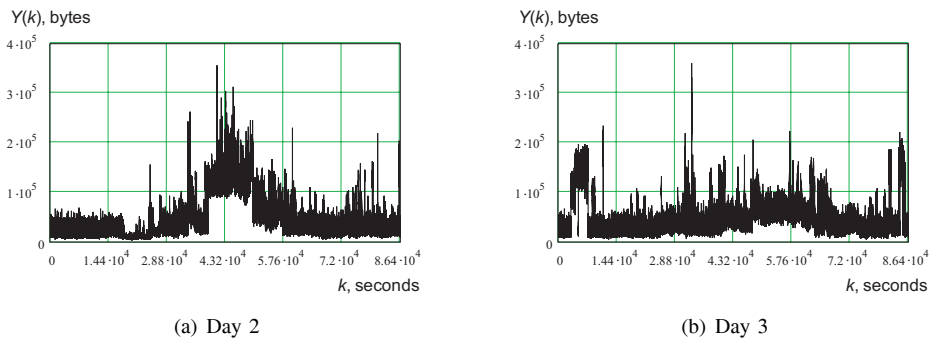
(a) Day 2



(b) Day 3

Fig. 1.    Time-series of 24-hours Auckland VIII traces.



(a) Histogram: Day 2



(b) NACF: Day 2



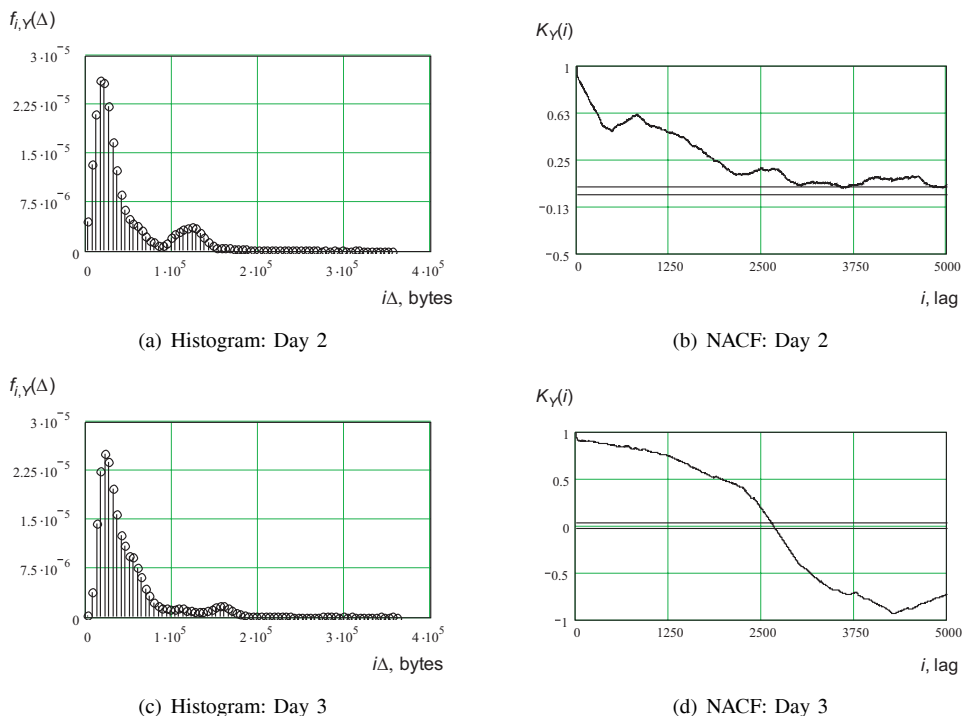(c) Histogram: Day 3



(d) NACF: Day 3

Fig. 2.    Histograms and NACFs of 24-hours Auckland VIII traces.
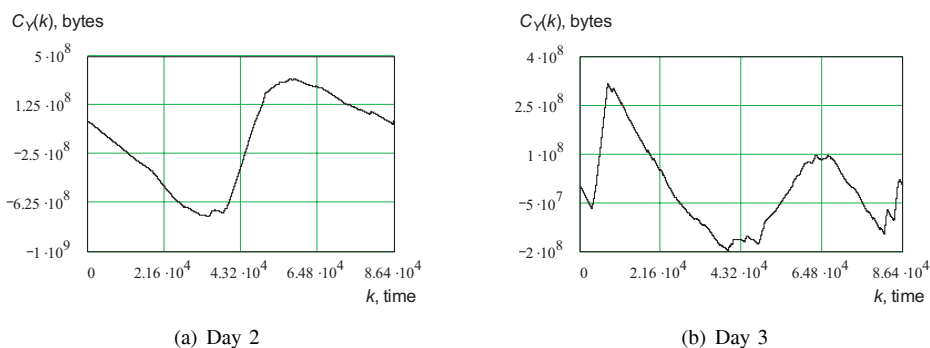


(a) Day 2



(b) Day 3

Fig. 3.    CUSUM statistics for 24-hours Auckland VIII traces.

of chosen parts are shown in 7(c) and Fig. 7(d) using solid lines with circles. One may observe that the memory of the process is short and limited to first few lags. We approximate this behavior using a single geometrical term in the form $y(i) = [K_Y(1)]^i$, $i = 0, 1, \ldots$, where $K_Y(1)$ is lag-1 NACF

value of observations. In Fig. 7(c) and 7(d) approximating functions are shown by solid thick lines.

Note that carried out statistical tests do not allow us to be statistically strict with our conclusion regarding covariance stationarity of considered subtraces. Unfortunately, there are
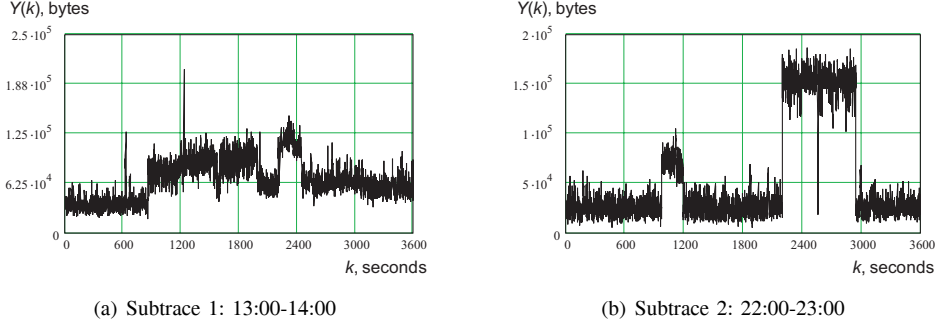
(a) Subtrace 1: 13:00-14:00



(b) Subtrace 2: 22:00-23:00

Fig. 4.   Time-series of 1-hour subtraces chosen from Day 3.



(a) Histogram: subtrace 1



(b) NACF: subtrace 1



(c) Histogram: subtrace 2
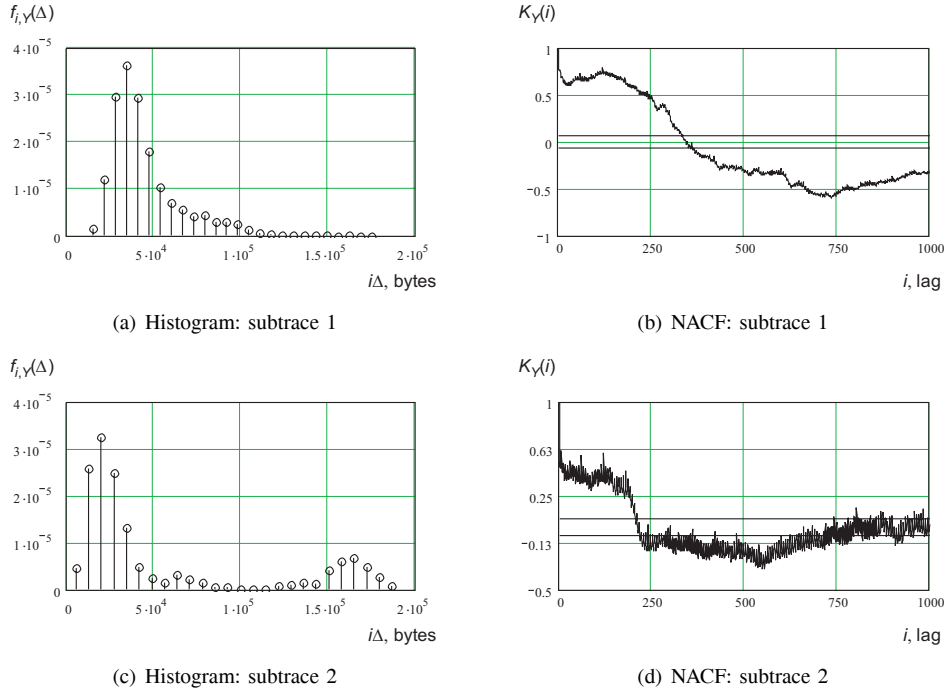


(d) NACF: subtrace 2

Fig. 5.   Histogram and NACF of 1-hour subtraces.

no effective methods to statistically test whether a limited set of observations is stationary or not. However, carried out tests statistically confirm that the marginal distribution is normal. NACF tends to zero as time progresses confirming that underlying processes are ergodic. These conditions are necessary for a underlying process to be covariance stationary.

## IV. MODEL FOR STATIONARY TRAFFIC

To model stationary parts of aggregated traffic we use autoregressive process of order one, AR(1). It has normal distribution and geometrically decaying ACF. These properties allow us to assume that AR(1) process may produce fair approximation of statistical data presented in Fig. 7.

A process is said to be autoregressive of order one when it is given by

$$X(n) = \phi_0 + \phi_1 X(n-1) + \epsilon(n), \qquad n = 0, 1, \ldots, \quad (1)$$

where $\phi_0$ and $\phi_1$ are some constants, $\{\epsilon(n), n = 0, 1, \ldots\}$ are iid random variables having the same Normal distribution with zero mean and variance $\sigma^2[\epsilon]$.

For (1) to be weakly stationary it is sufficient to have $\phi_1 \neq 1$. In this case

$$E[X(n)] = \mu_X, \qquad Cov(X_0(n), X_i(n+i)) = \gamma_X(i), \quad (2)$$

where $\mu_X$, $\gamma_X(i)$, $i = 0, 1, \ldots$, are some constants.

Mean, variance and covariance of AR(1) are related to $\phi_0$, $\phi_1$ and $\sigma^2[\epsilon]$ as

$$\mu_X = \frac{\phi_0}{1 - \phi_1}, \ \ \sigma^2[X] = \frac{\sigma^2[\epsilon]}{1 - \phi_1^2}, \ \ \gamma_X(i) = \phi_1^i \gamma_X(0). \quad (3)$$

Parameters of AR(1) models are related to statistical data as $\phi_1 = K_Y(1)$, $\phi_0 = \mu_Y(1 - \phi_1)$ and $\sigma^2[\epsilon] = \sigma^2[Y](1 - \phi_1^2)$, where $K_Y(1)$, $\mu_Y$ and $\sigma^2[Y]$ are estimates of lag-1 autocorrelation coefficient, mean and variance, respectively.

## V. DETECTING CHANGES IN TRAFFIC PATTERNS

### A. The Methodology

To differentiate between fluctuations of the traffic around the mean and changes in the mean value of the traffic process

(a) 00:00 - 06:00



(b) 06:00 - 12:00



(c) 12:00 - 18:00
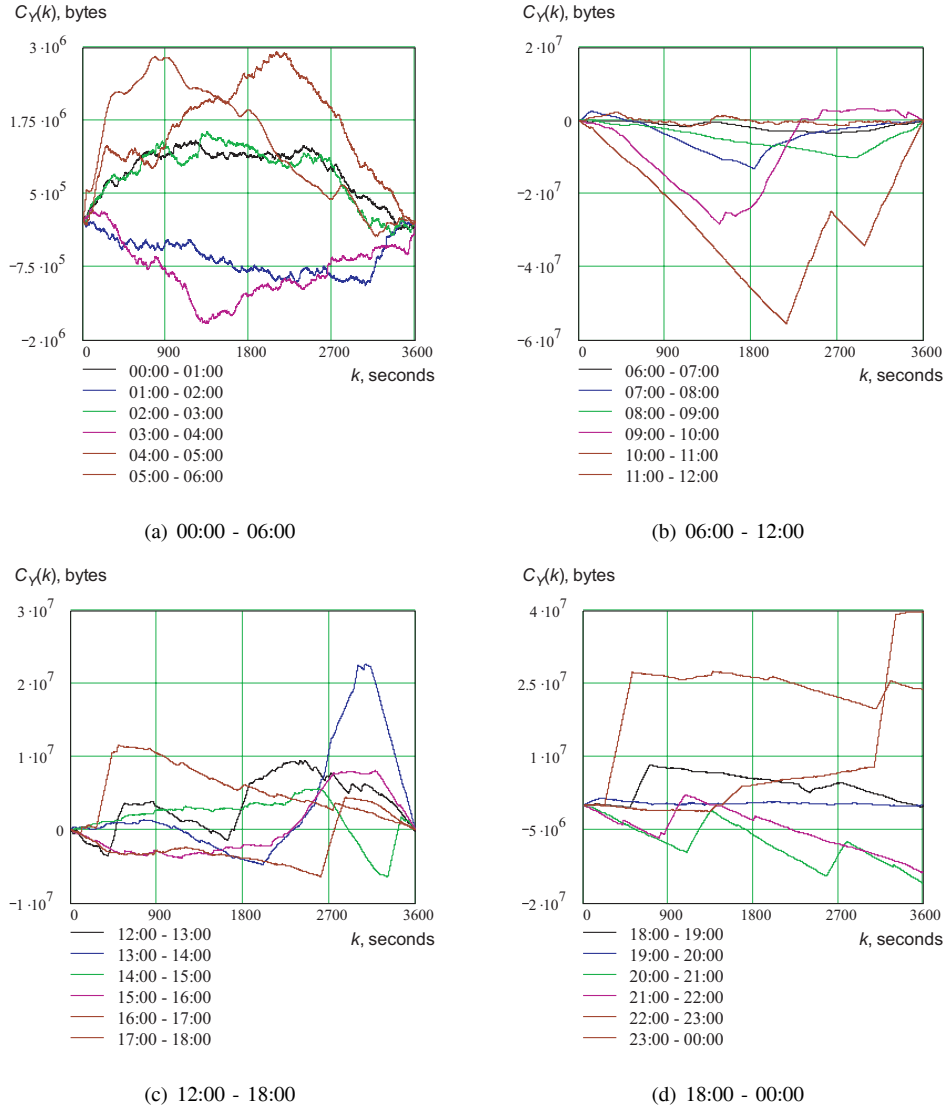


(d) 18:00 - 00:00

Fig. 6.    CUSUM statistics for 1-hour Auckland VIII traces.

we propose to use change-point statistical tests. There are a number of change detection algorithms developed to date. The common approach to deal with this task is to use control charts including Shewhart, CUSUM, and exponentially weighted moving average (EWMA) charts [21]. These charts originally came from statistical process control (SPC) where they are successfully used to monitor quality of production.

The idea of control charts is to classify causes of deviation from the target value into two groups. These are common and special causes of deviation. Deviation due to common causes is the effect of numerous causes affecting the process. They are inherent part of a process. Special causes are not the part of the process. Control charts signals the point at which special causes occur using two control limits. If observations are between them, process is 'in-control'. If some observations fall outside, the process is classified as 'out-of-control'. To detect changes in aggregated traffic patterns we assume that common causes are those resulting in inherent stochastic nature of the aggregated traffic. Special causes are those causing changes in

the mean value of 'in-control' process.

*B. Change in the Mean Value*

Assume that $k$ observations $\{Y(n), n = 0, 1, \ldots, k-1\}$ of a certain stochastic process have the same distribution $F_0$. The change-point statistical test refers to testing the null hypothesis, $H_0$, that a currently observed observation $k$ has distribution $F_0$ against alternative hypothesis, $H_1$, that its distribution is $F_1$. The latter case is when a change occurs in the distribution

$$H_1 : F_{Y,i} = \begin{cases} F_0 & i = 0, 1, \ldots, k-1, \\ F_1 & i = k. \end{cases} \qquad (4)$$

where $F_{Y,i}$ is the distribution of observation $i$.

The situation (4) is illustrated in Fig. 8, where changes in mean and variance led to respective changes in distribution. In Fig. 8(a) first 40 and last 20 observations were generated using normal distribution $N(\mu, \sigma) = N(200, 30)$, where $\mu$ and $\sigma$ are mean and standard deviation, respectively. Observations

(a) Histogram: subtrace 1



(b) Histogram: subtrace 2
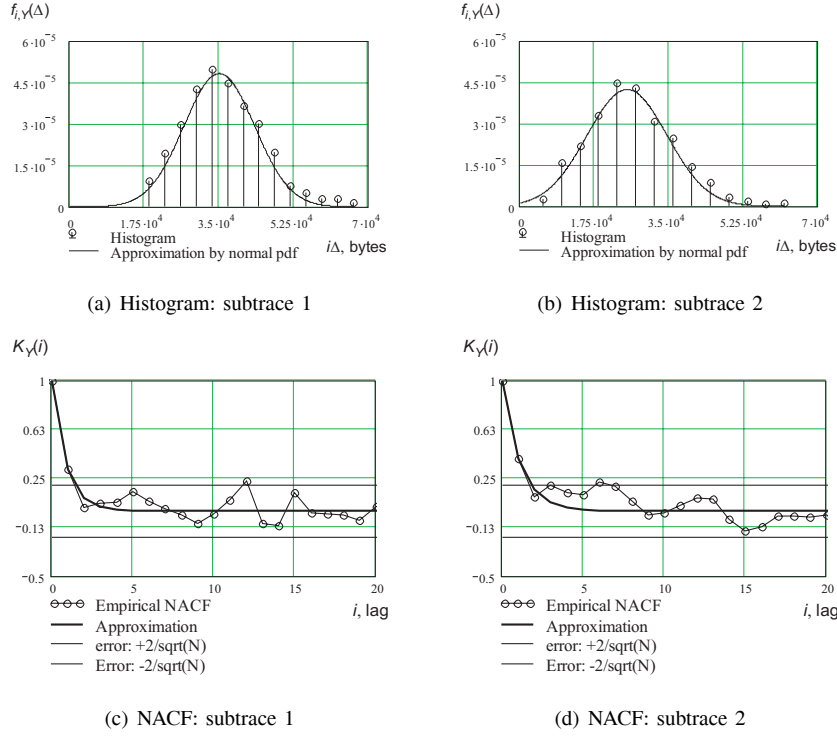


(c) NACF: subtrace 1



(d) NACF: subtrace 2

Fig. 7. Histograms and NACFs of first 600 observation of subtraces.

in the middle were generated using $N(250, 30)$. Fig. 8(b) demonstrates change in the variance that may also occur in traffic observations. Here first 40 and last 20 observations came from $N(200, 30)$ while those in the middle were drawn from $N(200, 60)$. Note that it may not be easy to visually detect changes in statistical characteristics. The change-point statistical test should be able to automatically detect these shifts.

It is often assumed that $F_0$ and $F_1$ are known except for some parameters of $F_1$. Control charts are used to detect changes in these unknown parameters. In our case the form of the distribution is known in advance, the whole task is to detect a change in the mean resulting in the following test

$$H_0 : E[Y] = \mu,$$
$$H_1 : E[Y] = \begin{cases} \mu_0 & i = 0, 1, \dots, k-1, \\ \mu_1 & i = k. \end{cases} \qquad (5)$$

*C. Change-point Statistical Tests*

The major shortcoming of change-point statistical tests is that they often assume that observations are independent. Traffic observations are not necessary independent but can be correlated. Autocorrelation makes classic control charts less sensitive to changes in the mean. For detecting changes in the mean value of autocorrelated processes two approaches have been proposed [22]. The first approach is to modify classic control charts. Control limits of these charts are widened to take into account correlational properties of empirical observations. The idea of the second approach is to fit observations to the time-series model and subsequently test residuals. If the

model fits data well, the residuals are uncorrelated and classic control chars can be used.

Performance of change-point tests for autocorrelated data has been compared in [22]. It was shown that residuals-based approach performs well when the autocorrelation is negative. When the autocorrelation is positive, modified control charts on initial observations perform better. This is because changes in the mean value are differently transferred to residuals for positive and negative autocorrelations [23]. Additionally, accuracy of residuals-based approach strictly depends on accuracy of fitting of the model to statistical data. Taking into account positive autocorrelation found in traffic tarces, modified Shewhart, CUSUM and EWMA charts provide an attractive option for on-line change detection.

Note that traffic traces may also contain outliers. They are of local significance and usually does not affect further service process of traffic. Control charts for detecting shifts in traffic observations should therefore contain mechanisms to deal with these outliers. Modified Shewhart chart performs on original observations and, therefore, does not incorporate such mechanism. As a result, every time an outlier occurs a change is signalled worsening performance of the algorithm. Modified CUSUM and EWMA charts operate on smoothed statistics protecting against outliers. However, on-line implementation of CUSUM chart requires knowledge of the global mean of the monitored process. In our task this information is not available. Thus, for detecting changes in traffic observations EWMA chart is chosen.
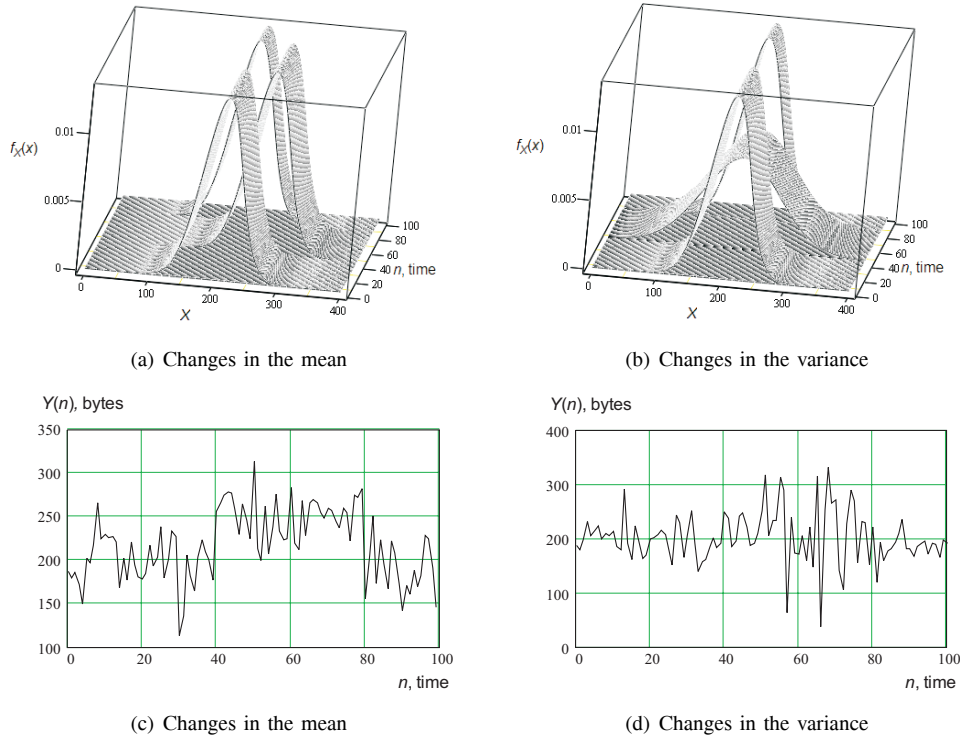
(a) Changes in the mean



(b) Changes in the variance



(c) Changes in the mean



(d) Changes in the variance

Fig. 8.   Changes in mean and variance of normal distribution.

### D. EWMA Control Charts

Let $\{Y(n), n = 0, 1, \dots\}$ be a sequence of observations. The value of EWMA statistic at the time $n$, denoted by $L_Y(n)$, is given by

$$L_Y(n) = \gamma Y(n) + (1 - \gamma)L_Y(n - 1), \qquad (6)$$

where parameter $\gamma \in (0, 1)$ is constant.

In (6) $L_Y(n)$ extends its memory not only to the previous value but weights values of previous observations according to constant coefficient $\gamma$. In (6) this previous information is completely included in $L_Y(n - 1)$. To show it, let us rewrite $\{L_Y(n), n = 0, 1, \dots\}$ statistics recursively, starting from $L_Y(0) = Y(0)$

$$
\begin{aligned}
L_Y(0) &= Y(0), \\
L_Y(1) &= \gamma Y(1) + (1 - \gamma)Y(0), \\
L_Y(2) &= \gamma Y(2) + \gamma(1 - \gamma)Y(1) + (1 - \gamma)^2 Y(0), \\
&\dots .
\end{aligned}
\qquad (7)
$$

Since for any constant $n$ the following holds

$$\gamma \sum_{i=0}^{n-1} (1 - \gamma)^i + (1 - \gamma)^n = 1, \qquad (8)$$

one can see that (7) converges to

$$L_Y(n) = \gamma \sum_{i=0}^{n-1} (1 - \gamma)^i Y(n - i) + (1 - \gamma)^n Y(0). \qquad (9)$$

The EWMA charts takes central part among other control charts. Although, according to (6), the most recent value always receives more weight in computation of $L_Y(n)$, the choice of $\gamma$ determines the effect of previous observations of the process on the current value of EWMA statistics. Indeed, when $\gamma \to 1$ all weight is placed on the current observation, $L_Y(n) \to Y(n)$, and EWMA statistics degenerate to initial observations. As a result, EWMA control chart behaves like Shewhart one. Contrarily, when $\gamma \to 0$ the current observation gets only a little weight, but most weight is assigned to previous observations. In this case, EWMA control chart behaves similar to CUSUM one. Summarizing, EWMA charts give more flexibility at the expense of additional complexity in determining one more parameter $\gamma$.

Assume now that given observations $\{Y(n), n = 0, 1, \dots, N\}$ are taken from strictly stationary process whose all sections are independently and identically distributed random variables with the same distribution, mean $E[Y]$ and variance $\sigma^2[Y]$. In fact, (6) defines a new stochastic process as a function of initial observations and this process has different statistical characteristics compared to those of $\{Y(n), n = 0, 1, \}$. Given that $L_Y(0) = E[Y]$, the mean of the process is

$$E[L_Y(n)] = E[Y](1 - (1 - \lambda)^n) + (1 - \lambda)^n E[Y], \quad (10)$$

that converges to constant $E[L_Y] = E[Y] = \mu$ as $n \to \infty$.

The variance of $\{L_Y(n), n = 0, 1, \dots\}$ is given by

$$\sigma^2[L_Y] = \sigma^2[Y]\left(\frac{\gamma}{2 - \gamma}\right)(1 - (1 - \gamma)^{2n}). \qquad (11)$$

Using (11) the control limits for EWMA charts ($E[L_Y] \pm$

$C(n)$) are computed as follows

$$E[L_Y] \pm k\sigma[Y]\sqrt{\left(\frac{\gamma}{2-\gamma}\right)(1-(1-\gamma)^{2n})}, \quad (12)$$

where $k$ is a design parameter whose values are tabulated in the literature [22], $C(n)$ denotes the deviation of control limits from the EWMA value at time $n$.

According to (12) an out-of-control behavior is signaled when $L_Y(n)$ at some point in time is less than $(E[L_Y] - C(n))$ or greater than $(E[L_Y] + C(n))$. Note that in (12) upper and lower control limits are time-varying in nature. However, when $n \to \infty$ it is easy to see that

$$\lim_{n\to\infty}\left(\frac{\gamma}{2-\gamma}\right)(1-(1-\gamma)^{2n} = \frac{\gamma}{2-\gamma}, \quad (13)$$

and constant limits $E[L_Y] \pm k\sigma[Y]\sqrt{\frac{\gamma}{2-\gamma}}$ can be used instead.

Assume now that given observations $\{Y(n), n = 0, 1, \ldots, N\}$ are taken from covariance stationary process with mean $E[Y]$ and variance $\sigma^2[Y]$ and can be well represented by AR(1) process. If $L_Y(0) = E[Y]$, observing (10) it is easy to see that $E[L_Y] = E[Y] = \mu$ when $n \to \infty$. The approximation of variance of $\{L_Y(n), n = 0, 1, \ldots\}$ for $n \to \infty$ is given by [22]

$$\sigma^2[L_Y] = \sigma^2[Y]\left(\frac{\gamma}{2-\gamma}\right)\left(\frac{1+\phi_1(1-\gamma)}{1-\phi_1(1-\gamma)}\right), \quad (14)$$

where $\phi_1$ is the parameter of AR(1) process.

The control limits are given by

$$E[L_Y] \pm k\sigma[Y]\sqrt{\left(\frac{\gamma}{2-\gamma}\right)\left(\frac{1+\phi_1(1-\gamma)}{1-\phi_1(1-\gamma)}\right)}. \quad (15)$$

*E. Details of the Algorithm*

The proposed algorithm is intended to operate as follows. When change detection procedure starts a warm-up time is spent gathering current traffic statistics. During this period the process is uncontrolled. At the end of warm-up period control limits are estimated. These control limits as well as the current value of EWMA statistics is used to control the process in what follows. While the current value of EWMA statistics is in between control limits the process is classified to be covariance stationary. Once the change is detected a new process is considered to be in-control and the control chart has to be re-parameterized according to statistics of this process. To do so a new warm-up period is started.

To parameterize the EWMA control chart a number of parameters have to be provided. Firstly, parameter $\gamma$ determining the decline of weights of past observations should be set. The values of $k$ and $\gamma$ determine the wideness of control belts for a given process with a certain $\sigma^2[Y]$ and $\phi_1$. These two parameters affect behavior of the so-called average run length (ARL) value that is usually used to determine efficiency of a certain change detection procedure. ARL is defined as the average number of observation of the in-control process up to the first out-of-control signal. The ARL is the function of both $k$ and $\gamma$. Different parameters of $k$ and $\gamma$ for a given ARL, $\sigma^2[Y]$ and $\phi_1$ are provided in [22], [24]. In this paper

to derive values of ARL we use computer simulations. Finally, $E[Y]$ and $\sigma^2[Y]$ are not usually known in practice and must be estimated from empirical data. Therefore, estimates of $E[Y]$ and $\sigma^2[Y]$ should be used in (15).

The algorithm can be summarized as follows
- during warm-up period estimate $E[Y]$, $\sigma^2[Y]$, $K_Y(1)$;
- choose $\gamma$ such that ARL matches a chosen value;
- estimate control limits of the chart;
- update EWMA statistics once new observation is ready;
- alarm once EWMA statistics is outside control limits;
- start new warm-up period.

To estimate statistics of a new 'in-control' process there must always be a certain warm-up period. There is a trade-off between accuracy of estimates and duration of the warm-up period. When the warm-up period is too small statistical estimates can be severely biased. Choosing warm-up period too large results in long periods of uncertainty during which performance of the traffic is uncontrolled. We recommend to set the length of warm-up periods to 50-200 observations. This choice provides reasonable resistance to outliers. Since the number of observations used to compute statistics of 'in-control' processes is relatively small, unbiased point estimators should be used. Finally, to keep storage requirements as low as possible recursive estimators should be used.

Note that changes in arrival statistics may also occur during warm-up periods. In this case, statistical estimates can be biased resulting in very large control limits. To deal with this problem the following procedure is suggested. Once $m_L$ observations during warm-up period are obtained, the 'local' EWMA chart is parameterized. This chart must have such ARL that inaccuracies in statistical estimates are taken into account. At the same time, control limits must allow to detect abrupt changes in arrival statistics. If change is detected, warm-up period is restarted. Our experiments with traffic aggregates from [19] demonstrated that values of $m_L$ and in-control ARL set to 5-10 and 2000-4000, respectively, are reasonable choice. It is also important to note that only few changes during warm-up periods were observed in our experiments. However, when they do occur, performance of the change-point detection algorithm decreases substantially.

## VI. NUMERICAL EXAMPLES

*A. Artificial Traces*

To demonstrate applicability of the proposed change-point detection algorithm we firstly use artificially generated traces. These traces consist of segments each of which follows AR(1) model with possibly different means. Lengths of segments were chosen randomly between 50 and 1000 observations. Each time new segment was generated, parameter $\phi_0$ of AR(1) models was chosen randomly out of the following vector $\vec{\phi_0} = \{100, 120, \ldots, 220\}$. Parameters $\phi_1$ and $\sigma^2[\epsilon]$ were kept constant and set to 0.3 and 30, respectively. Since $\phi_0$ is the only parameter that affects mean of the AR(1) model (1), mean of segments varied while variance and lag-1 autocorrelation coefficient were kept constant. We note that these parameters allow negative observations with non-negligible probability. Segments was concatenated to obtain the whole trace. The

whole trace is then realization of non-stationary process whose mean varies in time.

First 10000 observations from two generated traces are presented in Fig. 9(a) and Fig. 9(c). Results of EWMA change-point test for these traces are demonstrated in Fig. 9(b) and Fig. 9(d), respectively. The warm-up period was set to 100 observations. Parameters of EWMA test were chosen such that the target ARL value is 500. As one may notice, changes in non-stationary traces have been successfully detected. We note that in this particular case EWMA test detected all changes that occurred in considered traces. However, there can still be false alarms. To increase accuracy of the test the target ARL values can be increased. However, it will worsen reactive properties of the chart.

### B. Traffic Measurements

Consider how the proposed change-point detection algorithm performs for real traffic patterns. Time-series and corresponding results of EWMA change-point test for a number of traces from [19] are demonstrated in Fig. 10. Warm-up period was set to 100 observations. Parameters of EWMA test were chosen such that the ARL value is 500. As one may observe significant changes in the mean value of observations occur in all traces. EWMA test successfully detects those changes. For example, changes in IPLS3 are evident and all of them were detected.

## VII. APPLICATION SCENARIOS

The proposed modeling framework can be used whenever there is a need to describe traffic patterns in real-time. This situation occurs in QoS-aware networks such as multi-protocol label switching (MPLS) or differentiated services (DiffServ) when the user does not know exactly the traffic volume it may pose to the network but still requires network operator to provide performance guarantees. In this scenario network operator should be able to estimate on-the-fly the amount of resources required to serve the *current* traffic pattern with given performance metrics. Some of foreseen applications within this scope are discussed below.

### A. Dynamic Resource Reservation in DiffServ

Traffic patterns in DiffServ networks are often parameterized statically using the token bucket mechanism [25]. However, usually it is impossible to know in advance how much resources will be needed for a given traffic aggregate. In this case resource allocation based on token bucket parameters may lead to ineffective usage of network resources where during long durations of time traffic fluctuates under the advertised level. When servicing systems in a network are allowed to change resources assigned to a traffic aggregate on-the-fly, this shortcoming can be effectively avoided.

To deal with abovementioned problem we propose the dynamic resource reservation algorithm. It is based on measurements of local behavior of traffic aggregates and successive estimation of the amount of resources required to serve the traffic with given performance parameters. According to our

proposal, a customer notifies the network operator about the maximum load it may impose on the network. Network operator estimates the amount of resources required for advertised maximum load, chooses a path that providing these resources but does not reserve resources along this path.

Resource reservation is made dynamically as characteristics of traffic evolves in time. In this case service provider will pay exactly for the amount of traffic transferred in a network while the network operator is not required to statically devote a fraction of resources to the service that may not fully use them. Free resources can be temporarily assigned to other traffic patterns.

The proposed dynamic resource reservation algorithm should be implemented in all DiffServ routers. Basic elements of ingress and interior network nodes in DiffServ domain are shown in Fig. 11, where control activities are denoted by dashed lines, packet flows are shown by solid lines. We assume that ingress routers conform to DiffServ specifications including traffic conditioning and buffering [25]. The only important difference compared to DiffServ functionality is the resource controller. The purpose of this controller is to manage the resource allocation for a given behavior aggregate. Depending on the type of the node, controller takes the input from the conditioner or behavior aggregate classifier and monitors arriving traffic for possible changes in its statistical characteristics. Another responsibility of the controller is to estimate the amount of resources required to serve incoming traffic with given performance metrics. These actions concern in-profile traffic only. It is important to note that the controller does not change statistical characteristics of the traffic pattern allowing it to proceed unaltered to the output port.

The ingress node is intended to operate as follows. Traffic metering unit is statically parameterized such that the maximum load a customer intends to pay for is allowed to enter the network. Those packets that conform to this specification proceed to the resource controller. Resource controller monitors conforming traffic for possible changes in its statistical characteristics and estimates the amount of resources required to serve it with given performance metrics. When a change in traffic statistics occurs, resource controller estimates a new resource allocation and advertise it to the resource allocation mechanism at the output port and to the DiffServ node. Note that functionality related to dynamic resource reservation is similar for both ingress and interior DiffServ nodes.

As a descriptor of the amount of resources required by traffic aggregates token bucket is used. The actual amount of resources in terms of the buffer space and outgoing link rate share are inferred from token bucket parameters at each node. According to the token bucket mechanism, the amount of traffic allowed in the time interval $[0, t)$ is upper bounded by $A(t) = rt + b$, $t \geq 0$, where token rate $r$ is related to the outgoing link rate share, bucket size $b$ is related to the buffer space. To parameterize a token bucket we have to find a pair $(r, b)$ satisfying performance criteria. Usually, there are infinite number of pairs $(r, b)$ satisfying the required loss performance. However, there is always upper bound on $b$ that also satisfies delay requirements. There are a number of approaches to estimate token bucket parameters using statistics

(a) Artificial trace



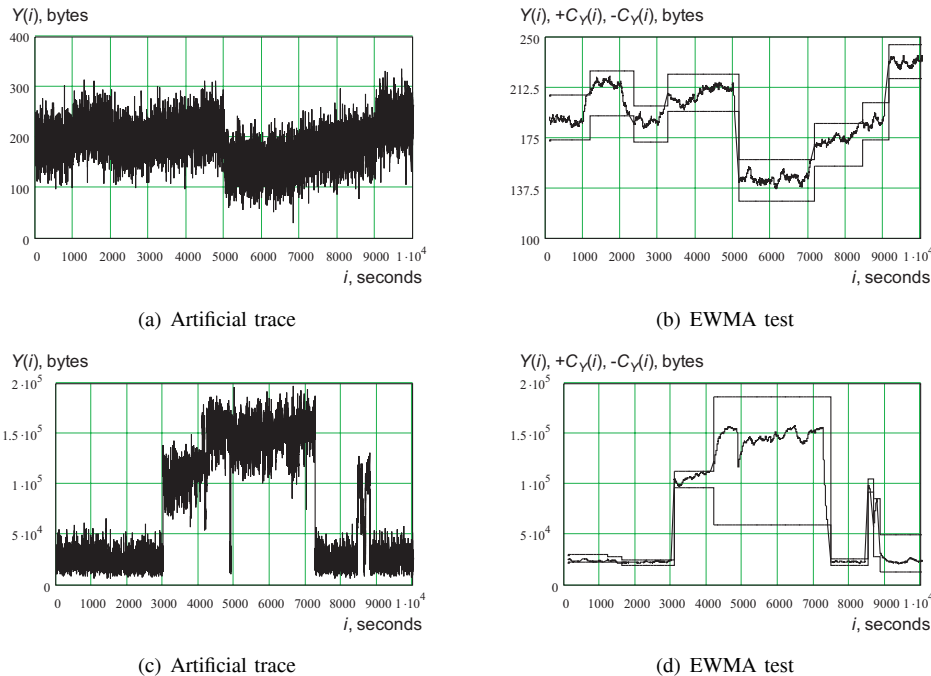(b) EWMA test



(c) Artificial trace



(d) EWMA test

Fig. 9.   Artificial traces and corresponding EWMA change-point tests.

of arrival process. They range from approximations providing token bucket parameters for worst case traffic scenarios to exact ones involving solution of equivalent queueing systems. Since we are interested in simple approach providing a feasible option for on-line implementation to estimate parameters of the token bucket we use overflow theory.

Note that the bucket size, $b$, should be set such that the delay is equal to or less than the maximum allowable delay in a network element. Thus, the task reduces to finding a suitable link rate share that should be assigned to a traffic aggregate. Let $\gamma(t)$ be the cumulative arrival process and denote $F_t(x) = Pr\{\gamma(t) \leq x\}$. When $F_t(x)$ is normal the following approximation can be used [16]

$$r = E[X] + \alpha\sigma[X], \qquad \alpha = \sqrt{-2\ln\epsilon - \ln 2\pi}, \qquad (16)$$

where $\epsilon$ is the target buffer overflow probability.

### B. Routing with Bandwidth Guarantees in MPLS

Nowadays many vendors already implemented both Diff-Serv and MPLS capabilities in their networking equipment. DiffServ provides tight performance bounds to traffic aggregates but requires all packets to follow the same path in a DiffServ domain. As one may see this requirement contradicts the hop-by-hop forwarding scheme of IP protocol and DiffServ must use an external forwarding scheme to implement QoS aware networks. MPLS complements DiffServ in this category providing a forwarding mechanism that allows to dynamically and automatically choose and reserve resources on the best available path for a given class of traffic [26]. For this purposes resource reservation protocol with traffic engineering support (RSVP-TE) or constrained routing label distribution protocols (CR-LDP) is used in MPLS. MPLS and DiffServ complement

each other providing all required features for QoS-aware networks [27].

To effectively route traffic in MPLS/DiffServ network it is important for network operators to know users' traffic demands in advance. The straightforward way is to derive these demands from service level agreements (SLA). When SLA is established network operator should provide bandwidth guarantees for the amount of traffic it agrees to serve. However, it is rarely feasible to know exact amount of traffic forcing users to guess. In this case, the network operator has rights to provide no guarantees to excessive traffic.

If the LSP arrival pattern as well as their traffic demands are known in advance off-line routing algorithms can be used to determine the best possible paths. When one or both characteristics are not known we have to resort to on-line routing. However, even in the case of on-line routing partial knowledge of requests arrival pattern and their traffic demands is extremely beneficial. It has been recently demonstrated that given the same performance guarantees in terms of the available bandwidth better resource utilization of the network can be obtained using on-line routing when daily traffic pattern is known in advance [28].

When information about statistical characteristics of LSPs is not available in advance we propose to use our algorithm to detect changes in arrival patterns. The proposed algorithm should be implemented at ingress routers only. We assume that ingress routers implement all the features specified by DiffServ and MPLS including traffic conditioning and buffering. The only important difference compared to DiffServ/MPLS is the resource controller implemented at ingress nodes. The main purpose of this controller is to control the resource allocation for LSP. Controller takes the input from the traffic conditioner and monitors the traffic for possible changes in its statistical
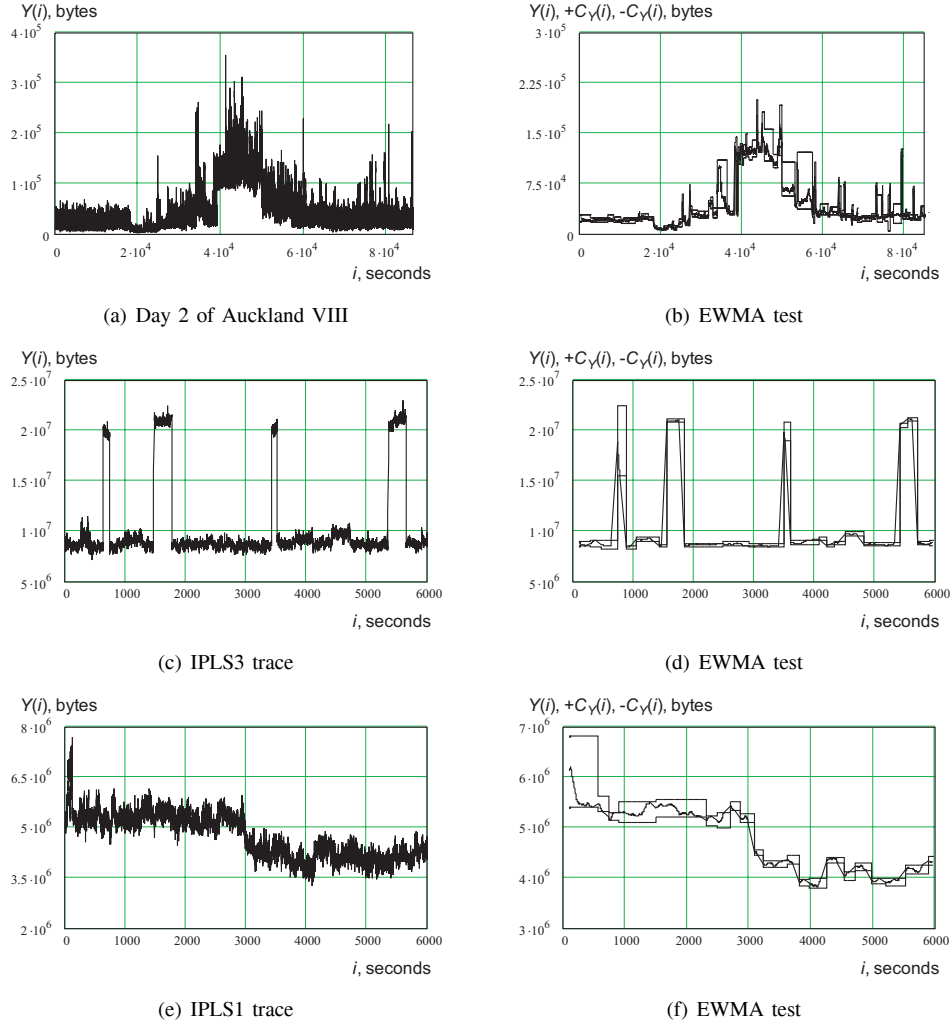
(a) Day 2 of Auckland VIII

(b) EWMA test

(c) IPLS3 trace

(d) EWMA test

(e) IPLS1 trace

(f) EWMA test

Fig. 10.   Time-series and results of EWMA change-point test.


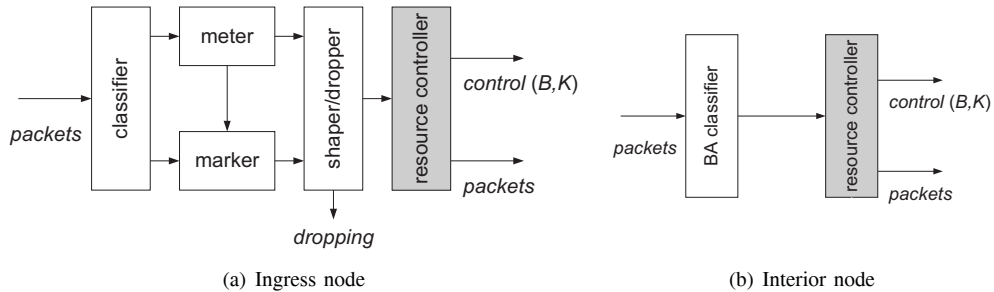
(a) Ingress node

(b) Interior node

Fig. 11.   Basic elements of ingress and interior network nodes in DiffServ.

characteristics. When change in traffic characteristics occur, resource controller estimates new resource allocation. RSVP-TE is then used to update resource allocation for a given traffic aggregate at all nodes along the path of associated label switched path (LSP). This update procedure is performed using PATH and RESV RSVP messages that maintain soft reservation states in interior nodes.

We would like to highlight that our algorithm requires ability of operator's network to dynamically re-route traffic. Indeed, those resources that are not currently used by the considered serve can be redistributed between other services

(if they require them) or assigned to the best effort serve. When the traffic of the considered service starts to build up again, the amount of traffic that used excessive available resources must be re-routed to another, possibly, 'worse' path that has enough available resources. This feature can be implemented using traffic engineering features of MPLS [26] or MPLS/DiffServ [29].

## VIII. CONCLUSIONS

We proposed the model for aggregated IP traffic. Firstly, we demonstrated that IP traffic observation can be repre-
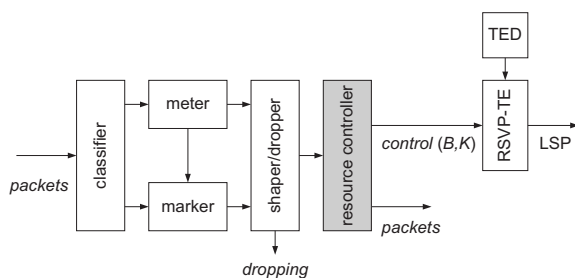
Fig. 12. Traffic conditioning functions.

sented using locally stationary processes. Using change-point statistical test as a basic tool of the algorithm we proposed to divide traffic observations into covariance stationary parts. This procedure should be implemented in real-time and the resulted model is inherently on-line in nature. We highlighted that the proposed approach can be of high practical importance in many areas of QoS-aware networks. This includes establishment of SLAs between network operators, dynamic bandwidth allocation in DiffServ networks, routing with bandwidth guarantees in MPLS or MPLS/DiffServ networks, etc.

We note that the ultimate goal of this paper was not to prove that Internet traffic on high-speed links is stationary or not but to propose a tool to deal with time-varying effects found in arrival processes. We did not claim that the traffic characteristics are in excellent agreement with local stationary assumption neither we claimed that those studies revealing self-similar and long-range dependent nature of traffic patterns failed to correctly interpret statistical properties. Instead, we stress that traffic patterns may exhibit high variability due to a number of phenomena including long-range dependent, self-similar and/or non-stationary properties. The overall task was reduced to estimating the *current state* of the traffic in terms of the covariance-stationary model. The current model remains valid for $i$, $i = 0, 1, \ldots$ future observations and can be used by networking equipment to take appropriate actions related to optimal traffic engineering. We also note that the model for covariance-stationary segments can be extended to include self-similar and long-range dependent behavior.

## REFERENCES

[1] M. Leland, M. Taqqu, W. Willinger, and D. Wilson. On the self-similar nature of ethernet traffic. *IEEE Trans. Netw.*, 2(1):1–15, Feb. 1994.
[2] K Thompson, G. Miller, and R. Wilder. Wide-area internet traffic patterns and characteristics. *IEEE Network*, 11:10–23, Nov./Dec. 1997.
[3] N. Duffield, J. Lewis, N. O'Connell, R. Russell, and F. Toomey. Statistical issues raised by the bellcore data. In *Proc. of 11th IEE UK Teletraffic Symposium*, March 1994.
[4] K. Park and W. Willinger. *Self-similar network traffic and performance evaluation*. John Wiley & Sons, 2000.
[5] M. Crovella and A. Bestavros. Self-similarity in world wide web traffic: Evidence and possible causes. *IEEE Trans. Netw.*, 5(6):835–846, Dec. 1997.
[6] M. Roughan and D. Veitch. Measuring long-range dependence under changing traffic conditions. In *Proc. IEEE INFOCOM*, pages 1513–1521, March 1999.
[7] T. Bohnert and E. Monteiro. A comment on simulating LRD traffic with pareto ON/OFF sources. In *Proc. ACM CoNEXT*, pages 228–229, Oct. 2005.
[8] V. Paxson and S. Floyd. Wide-area traffic: The failure of poisson modeling. In *ACM Sigcomm*, pages 257–268, 1994.

[9] T. Karagiannis, M. Molle, M. Faloutsos, and A. Broido. A nonstationary Poisson view of Internet traffic. In *Proc. IEEE INFOCOM*, pages 1558–1569, Hong Kong, March 2004.
[10] J. Cao, W. Cleveland, D. Lin, and D. Sun. On the nonstationarity of Internet traffic. In *Proc. ACM SIGMETRICS*, pages 102–112, 2001.
[11] Y. Zhang, N. Duffield, V. Paxson, and S. Shenker. On the constancy of Internet path properties. In *Proc. ACM SIGCOMM Internet Measurement Workshop*, pages 197–211, Nov. 2001.
[12] J. Beran, R. Sherman, M. Taqqu, and W. Willinger. Long-range dependence in variable bit rate video traffic. *IEEE Trans. Comm.*, 5:1566–1579, 1995.
[13] W. Willinger, M. Taqqu, R. Sherman, and D. Wilson. Self-similarity through high-variability: statistical analysis of ethernet LAN traffic at the source level. *IEEE/ACM Trans. Netw.*, 5(1):71–86, 1997.
[14] J. Kilpi and I. Norros. Testing the Gaussian approximation of aggregate traffic. In *Proc. 2nd Internet Measurement Workshop*, Marseille, France, Nov. 2002, Available at: http://www.imconf.net/imw-2002/proceedings.html.
[15] M. Yajnik, S. Moon, J. Kurose, and D. Towsley. Measurment and modeling of the temporal dependence in the packet loss. In *IEEE Infocom*, March 1999.
[16] R. Guerin, H. Ahamadi, and M. Naghshieh. Equivalent capacity and its application to bandwidth allocation in high-speed networks. *IEEE JSAC*, 9(7):968–981, Sep. 1991.
[17] D. Eun and N. Shroff. A measurement-analytic approach for QoS estimation in a network based on the dominant time scale. *IEEE Trans. Netw.*, 11(2):222–235, Apr. 2003.
[18] W. Leland, M. Taqqu, Willinger W., and D. Wilson. On the self-similar nature of Ethernet traffic (extended version). *IEEE/ACM Trans. on Netw.*, 2:1–15, February 1994.
[19] Passive Measurement and Analysis (PMA) project. Available at: http://pma.nlanr.net/, Accessed on: 26.11.2006, National Laboratory for Applied Network Research (NLANR).
[20] E. Page. A test for a change in a parameter occurring at an unknown point. *Biometrika*, 42:523–527, 1955.
[21] D. Montgomery. *Introduction to statistical quality control*. John Wiley & Sons, New York, 3rd edition, 1996.
[22] J. Wieringa. Control charts for monitoring the mean of AR(1) data. Available at: http://www.ub.rug.nl/eldoc/som/a/98a09/98a09.pdf, University of Groningen, Department of Econometrics, Faculty of Economic Sciences, Accessed on 06.07.2005.
[23] D. Apley and J. Shi. The GLRT for statistical process control of autocorrelated processes. *IIE Transactions*, 31:1123–1134, 1999.
[24] J. Wieringa. Statistical process control for serially correlated data, PhD Thesis. Available at: http://dissertations.ub.rug.nl/files/faculties/eco/1999/j.e.wieringa/, University of Groningen, Department of Econometrics, Faculty of Economic Sciences, Accessed on 18.10.2005.
[25] S. Blake, D. Black, E. Davies, Z. Wang, and Weiss W. An architecture for Differentiated Services. RFC 2475, IETF, 1998.
[26] D. Awduche, J. Malcolm, J. Agogbua, M. O'Dell, and J. McManus. Requirements for traffic engineering over MPLS. RFC 2702, IETF, 1999.
[27] F. Le Faucheur, L. Wu, B. Davie, S. Davari, P. Vaananen, R. Krishnan, P. Cheval, and J. Heinanen. Multi-protocol label switching (MPLS) support of differentiated services. RFC 3270, IETF, 2002.
[28] F. Ricciato and U. Monaco. Routing demands with time-varying bandwidth profiles on a MPLS network. *Comp. Netw.*, 47:47–61, 2005.
[29] F. Le Faucheur and W. Lai. Requirements for support of differentiated services-aware MPLS traffic engineering. RFC 3564, IETF, 2003.

**Dmitri Moltchanov** Dmitri Moltchanov is a Senior Research Scientist in the Institute of Communications Engineering at the Tampere University of Technology, Finland. He received the M.Sc. and Candidate of Science degrees from State University of Telecommunications, St.Petersburg, Russia in 2000 and 2002, respectively, and the PhD degree from Tampere University of Technology, Tampere, Finland in 2006. His research interests include performance evaluation and optimization issues of wireless and wired IP networks, ad hoc and sensor networks and P2P networks. Dmitri Moltchanov serves as a TPC member in a number of international conferences. He authored more than 40 publications.