

Using Markovian Models to Replicate Real ATM Traffics

Åke Arvidsson* and Christer Lind^{†‡}

University of Karlskrona/Ronneby* and Telia Research AB[†]

*Dept. of Telecommun. and Maths., Univ. of Karlskrona/Ronneby,
S-371 79 Karlskrona, Sweden. Email: akear@itm.hk-r.se

Tel: +46 455 78053. Fax: +46 455 78057.

[†]Telia Research AB, Commun. Sys., Box 85,

S-201 20 Malmö, Sweden. Email: Christer.Lind@malmo.trab.se

Tel: +46 40 105137. Fax: +46 40 307029.

Abstract

Among the more commonly employed models for performance analysis of ATM networks, *e.g.* to dimension buffers in switches, we find Markov modulated Poisson processes (MMPPs) and Markov modulated Bernoulli processes (MMBP). These models are often used with the only motivation that they are capable of producing bursty traffic. Although this is true in a general sense, little is known about whether that capability extends to the particular case of real traffics.

We report on an investigation where these models are tried in the latter sense. More precisely, we review and try a number of methods proposed for fitting MMPPs (MMBPs) to observed traffic data. The data consists of sixty traces which are extracted from the Bellcore Ethernet measurements according to length (short, medium, and long) and local average load (light and heavy). We then compare the performance of the buffer of a single server system when subject to the real traffic and the fitted model respectively.

It is found that the two cases differ significantly in terms of buffer occupancy, and that these differences are caused by deficiencies in the different fitting methods and possibly also by limitations in the models themselves. Nevertheless, some fitting methods are identified which, with further development, might work as models of burstiness within limited time spans on the order of two seconds. We also briefly comment the relationship between our results and recent works on fractal traffic characteristics.

Keywords

Bursty traffic model, ATM cell level traffic model, accuracy, Markov modulated Poisson Process, Markov modulated Bernoulli Process, MMPP, MMBP.

[‡]The major part of this work was carried out while Christer Lind was with the Department of Communication Systems, Lund Institute of Technology, Sweden.

1 MARKOVIAN MODELS FOR REAL ATM TRAFFICS

1.1 Markovian Models

Models of bursty traffic are frequently used in the context of performance analysis of ATM networks, *e.g.* to dimension buffers in switches. Among the more commonly employed models we find the Markov modulated Poisson processes (MMPPs) and Markov modulated Bernoulli processes (MMBPs). The two processes are doubly stochastic point processes where the rate of a Poisson (Bernoulli) process is governed by an underlying Markov chain in continuous (discrete) time. Arrivals and state transitions of the modulating chain are statistically independent. The processes are fully characterised by the number of states in the modulating chain s , the transition rates (probabilities) $q_{u,v}$, and the arrival rates (probabilities) r_u , $u, v \in \{1, \dots, s\}$.

To restrict the number of parameters, the number of states s is often set equal to two, in which case the model is referred to as a Switched Poisson (Bernoulli) Process, or an SPP (SBP). In the special case of the SPP (SBP) having an arrival rate of zero in one of its states, the process is called an Interrupted Poisson (Bernoulli) Process.

The main reasons why these models are frequently employed are probably their ability to match various burstiness characteristics, and their mathematical tractability. However, little is known about their actual relevance when it comes to producing a traffic that is not only generally bursty, but that in some meaning is equivalent to real traffic.

The current work is a preliminary attempt to investigate this aspect of simple Markovian models, typically SPPs and SBPs. The emphasis of the work is on their suitability for performance analysis, in particular with respect to buffer dimensioning. The general idea is to produce cell arrivals to an infinite buffer which is emptied by a single server, and study the number of cells present in the buffer at each arrival instant. A model that in our sense is equivalent to real traffic, would result in a buffer occupancy that is statistically identical to that of a real traffic.

To our knowledge, very few papers have been published where models are verified against real traffics in terms of buffer occupancy. Instead it appears that most researchers who verify models tend to do this against other models (!). One notable exception from this is the paper on video modelling published by Frater *et al.* (1994) and Rose (1994), where the queuing behaviour of a real traffic is compared to that of a model.

1.2 Replicating Real Traffics

Users wishing to establish a connection over an ATM network are required to declare a number of parameters characterising the traffic they wish to submit. These parameters include peak rate, sustainable rate, burst size, and possibly others. Typical factors affecting the choice of parameter settings include the nature of the application, characteristics of the user premises equipment, the access medium, and the tariff structures.

Testing models under this scenario, the model should represent the traffic actually submitted, *i.e.*, after possible shaping by the policing device. For a given trace, the user could declare virtually any set of parameters, and deliver the traffic in a number of conforming and non-conforming ways. To avoid restrictive presumptions regarding these parameters and delivery, we assume that the parameters are set such that the traffic can be passed transparently to the network, and therefore simply match the models directly to the traces.

The procedure of matching a model to a traffic trace is referred to as a fitting method. The fitting methods considered in this work can be classified in three categories: Sequence fitting, direct metrics fitting, and indirect metrics fitting. It is pointed out that the choice of modelling in discrete or continuous time is more a matter of mathematical convenience than of replication accuracy, Arvidsson *et al.* (1991).

Sequence Fitting

The idea of sequence fitting is based on the presumption that the trace is in fact produced by a specific model the parameters of which are unknown. Fitting a model to a trace therefore means to find the set of parameters of this particular model that have the highest likelihood of producing that sequence. The typical procedure is to start from an initial guess of the parameter set and successively improve it with respect to the likelihood of obtaining the trace until no further improvement can be obtained.

We have used two methods of this class, one due to Meier-Hellstern (1987) (KMH) and another one due to Rydén (1992) (TR). Both are developed for MMPPs with any number $s > 1$ of states, but are here applied to the case $s = 2$.

Direct Metrics Fitting

Direct metrics fitting does not presume that a certain model is actually valid, but simply aims at making the model in question reproduce certain “important” and mathematically tractable properties of the trace. Typical such properties fitted to are moments and correlations of inter arrival times and of the number of arrivals within intervals of length t .

We have considered four such methods, Rossiter (1987) (MR), Heffes *et al.* (1986) (HL), Gusella (1991) (RG), and Park *et al.* (1994) (DP). The three former are developed for and applied to MMPPs with $s = 2$ states and the latter to MMBPs with $s = 2$ states.

Indirect Metrics Fitting

Indirect metrics fitting means that the observed process is first transformed into another process which is then dealt with as for direct metrics fitting. The transformation procedure we have considered is the identification of “active periods” and “passive periods”, an idea first proposed by Jain *et al.* (1986). Active periods refer to uninterrupted sequences of one or more short inter arrival times, and passive ones to uninterrupted sequences of one or more long inter arrival times. Properties of interest in the transformed process include moments and correlations of the lengths of the two periods and of the activity within each of them.

We have used four such approaches, Solé *et al.* (1990) (SDG/1) and (SDG/4), Bonomi *et al.* (1994) (BMMP), and Lee *et al.* (1992) (LL). Both SDG/ x -methods refer to MMBPs with $s = 2$ states and allow for activities between zero and one during both periods. BMMP and LL refer to MMBPs and MMPPs with $s = 3$ and $s = 2$ states respectively, and both prescribe strictly no activity during passive periods and strictly full activity during active ones.

1.3 Preliminaries of the Investigation

It is well known that traffic characteristics depend heavily both on the source (*e.g.* video or data) and on the content (*e.g.* drama, sports, file transfer and www-retrievals). It is also clear that not even for a given source and content, there is such a thing as a “typical behaviour”. A general investigation of traffic replicating properties would therefore require tremendous amounts of recorded traffic traces. We have restricted ourselves to one class of traffic which could be labelled “LAN interconnect”. The motivation

for our particular choice is twofold: LAN interconnect is expected to be one of the first traffics to be sent over ATM, and LAN traffics measurements were readily available to us through the Bellcore (1989) measurements.

Numerous papers, *e.g.* Leland *et al.* (1994), Paxon *et al.* (1994), Pruthi (1995), and others, have reported on the self similar properties of these traffic traces, the presence of variations on all time scales, and the heavy tailed buffer occupancy distributions resulting from them. These findings raise fundamental questions regarding the relevance of Markovian models, in particular for those with small numbers of states s , the variabilities of which span a strictly limited time scale, *cf.* Andersen (1995).

The present work is, however, restricted to model variations within certain time scales. This is motivated by engineering aspects of buffer dimensioning, where loss constraints for slowly varying traffics may call for very large buffers, quite possibly large enough to violate delay constraints and even beyond reasonable physical limitations. (This becomes obvious when looking at buffer sizes and performance for systems that store excess traffic generated during working hours and transmit it during the nights.) Generally speaking, we can thus identify two kinds of variations: Fast variations which can be smoothed by a buffer, and slow variations which cannot. We are only interested in the former.

For slow variations we can see at least three possible ways: the first one is to multiplex a very large number of independent sources in which case even slow variations can be statistically multiplexed; the second one is to provide enough transmission capacity to handle the peaks and simply put up with the resulting poor utilisation in the valleys; and the third one is to trace the slow variations and dynamically adjust the allocated transmission capacity in accordance with the variations. Our work is based on the last approach, and we presume the presence of a control mechanism that dynamically adjusts the capacity of the server to the long term average of the traffic load. We do not develop such a mechanism here, but only mention that it could be driven by user initiated requests for more or less network resources following the opening or closing of new applications (ftp, telnet, netscape *etc.*), with signals from system initiated monitoring of traffics and/or buffers as an alternative or supplement.

In the language of ATM traffic control variations are often said to take place in the cell scale (typically on the order of μ s), burst scale (α ms), activity scale (α s), session scale (α min) *etc.*, *e.g.* Bagnoli *et al.* (1994), Hui (1988), Key (1995), Ramamurthy *et al.* (1994), and others. Clearly, buffers are intended only for the cell-, burst- and possibly activity scale, hence modelling of these scales is sufficient from a buffer dimensioning point of view.

2 EXPERIMENTS WITH REAL TRAFFIC AND MODELS

2.1 Background

We defined an experimental test bed based on the following scenario: A user wishes to convey LAN data over an ATM network. The LAN is a 10 Mbps Ethernet, and the user is connected transparently to the ATM network via a 34 Mbps link. Before the LAN packets are delivered over this link to the network, the Ethernet overhead is stripped of, and the remaining data packed into cells. Each 53 octet cell can take 44 octets of Ethernet data, since 4 octets "pay load" are used for AAL3/4 overhead, and the last 5 octets constitute the ATM header.

We implemented a simulator with a server of capacity C and an infinite buffer. Cell arrivals follow sample traces from the Bellcore (1989) material converted to ATM as above, or are drawn from a mathematical model. The traces used were chosen according to length and load: Three time scales were chosen,

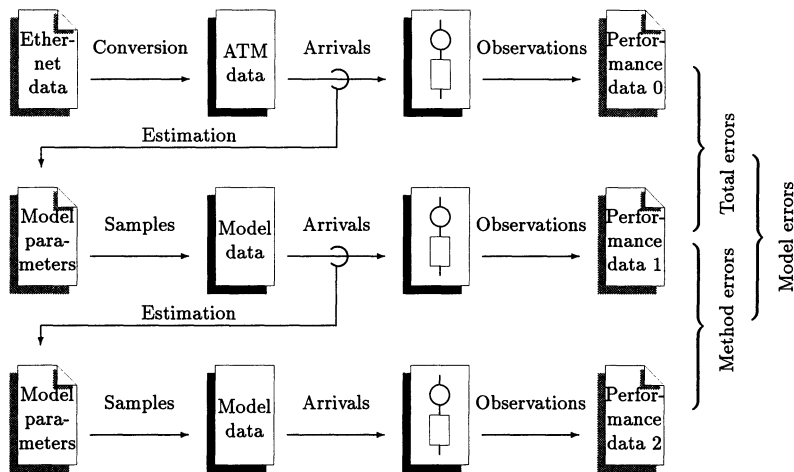


Figure 1 Experiments carried out.

viz. 0.20, 2.00 and 20.0 seconds, and two long term loads, *viz.* 42% and 85% of the actual peak value observed in intervals of those lengths in the entire material made available to us. Finally, for each time scale and load condition were 10 distinct traces selected, resulting in a total of $3 \times 2 \times 10 = 60$ traces. It is observed that the chosen time scales should well cover the normal scope of buffer modelling, *i.e.*, cell scale and burst scale variations.

For each trace, the transmission capacity C was set according to the formula for equivalent bandwidth given by Vakil (1993) $C = a(1 - \log a/p)$, where a is the long term average rate (in our case over the entire trace) and p is the peak rate (in our case 34 Mbps). This setting is high enough to ensure that the system is not overloaded, while at the same time it is low enough to let queues build up during the peaks.

2.2 Experiments

In an initial series of runs, each real trace was used as arrival generator in our simulator, figure 1. The trace was run repeatedly in order to emulate a local “steady state”. At each arrival instant we noted the number of cells present in the buffer, which was taken as the sole performance metric for the single server system, denoted in the figure as “Performance data 0” (PD0).

Next, each of the traces were fed into each of the parameter fitting procedures mentioned above. This resulted in one set of model parameters for each trace and each fitting method. The models thus obtained were then used as traffic sources in our simulator, and the performance of the buffer was monitored as before. The observations are shown as “Performance data 1” (PD1) in figure 1.

Noting that the models are fit directly to the traces, one would ideally expect that the models give the same buffer performance as the real traces, *i.e.* PD0 to be equal to PD1. However, it must be remembered

Table 1 Number of infeasible fits.

Model used	0.20 sec.		2.00 sec.		20.0 sec.	
	42% load	85% load	42% load	85% load	42% load	85% load
KMH	—	—	—	—	—	—
TR	—	1	—	—	—	—
MR	6	10	—	5	—	3
HL	4	9	—	4	—	—
RG	3	10	—	8	1	6
DP	5	8	—	1	—	—
BBMP	—	—	—	—	—	—
LL	—	—	—	—	—	—
SDG/1	—	—	—	—	—	—
SDG/4	—	—	—	—	—	—

that the models themselves cannot take all the blame of any differences detected, but some may be due to deficiencies in the parameter estimation *etc.* We may thus say that any difference obtained between a PD0 and PD1 consists of two components: One which is due to the model, and an one which is due to the fitting method and our implementation thereof. We call the former component “model error”, the latter part “method error” and refer to the observed sum as “total errors”.

In order to estimate the two components separately, a new set of experiments was conducted: The above runs for each model and fitting method were monitored and fed to the same fitting procedure as the one used for the model under study, *i.e.*, we fitted each of model to themselves. The resulting set of models were then taken as arrival generators in our simulator, and the buffer performance again monitored as before. The results are indicated as “Performance data 2” (PD2) in figure 1. The fact that the models fitted to are valid by definition in this series of runs means that there are no model errors, but any differences between PD1 and PD2 relate to method errors only. Loosely speaking, we may then obtain the model error by subtracting the method error from the total error.

3 RESULTS

3.1 Validity

For a set of model parameters to be *feasible*, we require that arrival rates are ≥ 0 and transition rates > 0 for MMPPs, and that arrival probabilities are ≥ 0 and ≤ 1 and transition probabilities > 0 and ≤ 1 for MMPPs. The requirements follow from physical interpretations with the added condition that the modulating chain must not be absorbing. Not all fitting methods came up with feasible parameters for all samples. The number of failures are shown in table 1 for each model respectively.

The table shows that these anomalies occur almost solely for direct metrics fitting. The only exception from this rule is one sequence fit, where the modulating chain turned out to be absorbing. It is also

Table 2 Number of abnormal fits.

<i>Model used</i>	<i>0.20 sec.</i>		<i>2.00 sec.</i>		<i>20.0 sec.</i>	
	<i>42% load</i>	<i>85% load</i>	<i>42% load</i>	<i>85% load</i>	<i>42% load</i>	<i>85% load</i>
KMH	1	—	—	—	6	—
TR	—	—	2	4	7	—
MR	—	—	—	—	—	—
HL	—	—	—	—	—	—
RG	—	—	—	—	1	—
DP	—	—	—	—	2	—
BBMP	—	—	—	—	5	2
LL	1	7	—	—	—	—
SDG/1	—	—	—	—	9	—
SDG/4	—	—	5	—	7	1

seen that infeasible parameters are more often obtained at high loads and when fitting to short intervals. Notably, MR and RG failed for all ten traces for the most extreme case in this respect.

Infeasible parameters are explained as follows: Direct fitting methods employ four equations in four metrics from which the four MMPP (MMBP) parameters are found. The output of an MMPP (MMBP) has certain limits regarding the relations between various metrics, and infeasible parameters from a certain trace therefore indicate that the model is incapable of exactly reproducing the metrics of that trace. In this case, one could alternatively find the nearest feasible solution as some kind of best fit. However, our work does not aim at developing or improving fitting methods, but is restricted to testing existing proposals.

Furthermore, for a fitting to be *meaningful*, the resulting traffic model must produce an average queue length that is in the vicinity of the one obtained for the real traffic. We have rather arbitrarily stated that non-meaningful results are those that differ by a factor of 10 or more from the target values. The occurrence of such cases is shown in table 2.

It is seen that abnormal fits almost only occur for sequence fitting and indirect fitting. A closer look at the numbers behind the table reveals mismatches resulting in permanent overloads of the simulated system for the entries referring to sequence fitting. This means that the considered methods, which are iterative, sometimes converge towards a solution that is not correct in terms of average arrival rate. Again, it is beyond the scope of this work to solve the problems behind this phenomenon. For the indirect fitting methods, the abnormal values are less severe, but simply point at weaknesses in the methods as such.

3.2 Accuracy

Metrics

We now remove the infeasible and abnormal fits from our data set and investigate the accuracy of the models with respect to the remaining runs. More precisely, we consider how well the various models and fitting methods can mimic real traffics with respect to dimensioning buffers over the selected time scales.

Let the occupancy of the buffer at an arrival instant be denoted by a stochastic variable Q and define two primary metrics of system performance, *viz.* $E\{Q\}$, the mean buffer occupancy over the entire distribution, and $E\{Q'\}$, the mean occupancy over the tail of the distribution,

$$E\{Q\} = \sum_{k=0}^{\infty} kp(k); \quad E\{Q'\} = \sum_{k=k'}^{\infty} kp'(k)$$

where $p(k)$ refer to the probability of an arriving customer finding k customers already in the queue, and $p'(k)$ is $p(k)$ renormalised over the tail. The tail is defined as all states $k \geq k'$, where k' is the smallest k' such that $\sum_{\kappa=k'}^{\infty} p(\kappa) \leq 10^{-2}$. Note that the latter metric does not refer to a single point, which would have made it very sensitive, but to the *rescaled average* of the last percent of the distribution and thus captures the tail in a wider sense. This number was chosen as a compromise between tail probabilities relevant to buffer dimensioning, typically on the order of 10^{-9} , and simulation feasibility and accuracy.

Adding to the notation, we let Q_i be the performance metric observed from the i th data set in figure 1, *i.e.*, Q_0 refers to the real trace, Q_1 to the fit to the real trace, and Q_2 refers to the fit to the fit. Finally, we define two metrics of the *total error* mentioned in figure 1 as

$$\epsilon_{\text{tot}}(Q) = 1 - E\{Q_1\}/E\{Q_0\}; \quad \epsilon_{\text{tot}}(Q') = 1 - E\{Q'_1\}/E\{Q'_0\}$$

two metrics of the *method error* in the same figure as

$$\epsilon_{\text{met}}(Q) = 1 - E\{Q_2\}/E\{Q_1\}; \quad \epsilon_{\text{met}}(Q') = 1 - E\{Q'_2\}/E\{Q'_1\}$$

and two metrics of the *model error* in the same figure as

$$\epsilon_{\text{mod}}(Q) = E\{Q_2\}/E\{Q_1\} - E\{Q_1\}/E\{Q_0\}; \quad \epsilon_{\text{mod}}(Q') = E\{Q'_2\}/E\{Q'_1\} - E\{Q'_1\}/E\{Q'_0\}$$

Total Errors

Tables 3 and 4 show $\epsilon_{\text{tot}}(Q)$ and $\epsilon_{\text{tot}}(Q')$ respectively for each combination of time scale and load. The numbers shown refer to the average over all valid traces. Rather than providing standard deviations as a supplement, each entry in the table was subject to a *t*-test, *i.e.*, we tested whether the observed average, given the variations between the various traces, could in fact be an observation of a distribution with zero average. The results are depicted in tables 5 and 6 respectively: The number of stars indicate the confidence by which the hypothesis is rejected: three stars mean 99.9% certainty, two stars 99% certainty and one star 95% certainty. No stars thus indicate that the hypothesis cannot be rejected with an error probability below 5%, but *not* that the hypothesis is correct.

It is seen that large errors are frequent, and generally more so for the tail than for the whole distribution. We also note that while many models tend to over estimate the mean of the queue length, they still underestimate the tail, an observation in accordance with observations from heavy tailed traffic.

Table 3 might give the impression that some of the methods based on direct metrics fitting perform reasonably well for short intervals with small, non-significant errors. However, it must be remembered that these values are based on very few actual observations because of the large number of infeasible fits, *cf.* table 1. A similar observation holds for the results of TR in the cases of longer traces.

The same is true for the mean of the tail of the distribution: The only positions with small average errors which pass a test for zero, are those that contain few entries, in particular the case with a time span

Table 3 Total errors observed for the mean of the whole of the distribution.

<i>Model used</i>	<i>0.20 sec.</i>				<i>2.00 sec.</i>				<i>20.0 sec.</i>			
	<i>42% load</i>		<i>85% load</i>		<i>42% load</i>		<i>85% load</i>		<i>42% load</i>		<i>85% load</i>	
KMH	36	**	-136	***	79	***	53	***	88	***	77	***
TR	48	***	-117	**	-6		55	***	2		10	
MR	14		—	—	42	***	34	**	58	***	52	***
HL	23	*	—	—	43	***	41	**	56	***	56	***
RG	11		—	—	43	***	36	*	53	***	53	**
DP	18		-144		75	***	42	***	82	***	67	***
BBMP	52	***	-34	**	82	***	68	***	86	***	78	***
LL	-409	***	-246		-226	**	-232	*	-207	***	-178	***
SDG/1	7		-168	***	65	***	39	***	—	—	68	***
SDG/4	46	***	-70	***	76	***	53	***	83	***	74	***

Table 4 Total errors observed for the mean of the tail of the distribution.

<i>Model used</i>	<i>0.20 sec.</i>				<i>2.00 sec.</i>				<i>20.0 sec.</i>			
	<i>42% load</i>		<i>85% load</i>		<i>42% load</i>		<i>85% load</i>		<i>42% load</i>		<i>85% load</i>	
KMH	-138	**	-700	***	39	**	14		87	***	74	***
TR	-102	**	-648	***	-54		22		3		41	
MR	-204		—	—	-31		3		53	***	48	***
HL	-163	**	—	—	-42		-6		37	**	54	***
RG	-187	**	—	—	-30		7		48	***	64	**
DP	-177	*	-685		32	**	-1		76	***	66	***
BBMP	-54		-336	***	51	***	50	***	80	***	76	***
LL	-1305	***	-1014	*	-640	**	-377	*	-291	***	-205	***
SDG/1	-223	***	-842	***	7		-10		—	—	67	***
SDG/4	-86	*	-532	***	40	**	22	*	81	**	74	***

of 2 seconds and with a long term average load of 85%. An overall conclusion is that the generally large errors make it hard to find a “best model”, and selecting a “worst model” appears equally meaningless.

Method Errors

We will now attempt to get a better idea of the origin of the errors: *i.e.*, if these should be attributed to the models themselves, or if it is just as likely that it is the fitting procedure and our implementation thereof that are to be blamed. Tables 5 and 6 show $\epsilon_{\text{met}}(Q)$ and $\epsilon_{\text{met}}(Q')$ in the same way as above. It is immediately seen that the method errors are by no means small or insignificant for any of the models

Table 5 Method errors for the mean of whole of the distribution.

Model used	0.20 sec.				2.00 sec.				20.0 sec.			
	42% load		85% load		42% load		85% load		42% load		85% load	
KMH	-16	***	-19	***	-14	***	-16	***	-12	***	-15	***
TR	-1	**	-3	***	-3		-2	***	-4		-2	**
MR	17		—	—	18		4		28	***	14	*
HL	32	**	—	—	15	**	38	**	7	*	32	***
RG	31	**	—	—	28	*	2		25	**	13	
DP	0		1		1	***	1	*	1	**	1	**
BBMP	-2	*	-1	*	-1		-1	**	-2		-1	***
LL	-9		-46		-16		-40	**	-10		24	**
SDG/1	-26	***	-19	***	-13	*	-23	***	—	—	-28	***
SDG/4	23	***	15	***	31	**	26	***	18	**	33	***

Table 6 Method errors for the mean of tail of the distribution.

Model used	0.20 sec.				2.00 sec.				20.0 sec.			
	42% load		85% load		42% load		85% load		42% load		85% load	
KMH	-15	**	-14	***	-15	*	-11	***	-16	**	-8	***
TR	1		0		0		-1		-2		-5	
MR	14		—	—	15		3		25	***	16	**
HL	29	*	—	—	11	*	35	**	7	*	29	***
RG	30	*	—	—	29	**	4		25	**	10	
DP	22	***	19		24	***	14	***	9	*	3	
BBMP	0		-4		-1		-2		-7		-2	*
LL	-28		-18		-18		-45	**	-7		27	**
SDG/1	-42	***	-21	**	-9		-40	***	—	—	-37	***
SDG/4	28	***	20	***	34	**	24	***	23	**	29	***

but TR and BBMP. If only the mean is considered, the method errors of DP could also be referred to as small.

It is clear that the fitting methods are not particularly stable when it comes to estimating parameters of a model which they essentially have created themselves. This does not mean to say that the formulae provided in the various papers where the methods are put forward are incorrect. What it does say, however, is that the traffic characteristics we are concerned with results in parameters which are hard to estimate. That is, when the trace is fitted to a model for the first time, we get parameters which are in range of hard-to-estimate MMPP-parameters (MMBPs-parameters). The existence of such cases

Table 7 Model errors for the mean of whole of the distribution.

Model used	0.20 sec.			2.00 sec.			20.0 sec.					
	42% load	85% load		42% load	85% load		42% load	85% load				
KMH	52	**	-117	***	93	***	69	***	100	***	92	***
TR	49	***	-115	**	-3		57	***	7		13	
MR	-2		—	—	25		30		31	**	38	***
HL	-9		—	—	28	*	3		49	***	23	**
RG	-21		—	—	15		33		27	**	40	
DP	18		-144		74	***	42	***	81	***	66	***
BBMP	54	***	-34	**	82	***	69	***	89	***	80	***
LL	-401	**	-200		-210	**	-191		-196	**	-202	***
SDG/1	33	*	-150	***	78	***	62	***	—	—	96	***
SDG/4	23		-84	***	45	**	27	**	65	***	41	***

Table 8 Model errors for the mean of tail of the distribution.

Model used	0.20 sec.			2.00 sec.			20.0 sec.					
	42% load	85% load		42% load	85% load		42% load	85% load				
KMH	-123	*	-686	***	54	**	25	**	103	***	82	***
TR	-103	**	-648	***	-54		23		5		46	
MR	-219		—	—	-45		0		28	**	31	**
HL	-192	***	—	—	-53		-41		30	*	25	*
RG	-217	**	—	—	-60	*	3		22	*	55	*
DP	-198	*	-704		8		-15		67	***	63	***
BBMP	-54		-331	***	51	***	52	***	87	***	78	***
LL	-1277	***	-996		-622	**	-332	*	-285	***	-232	***
SDG/1	-181	**	-822	***	15		31	**	—	—	104	***
SDG/4	-114	**	-552	***	6		-2		58	**	44	***

is mentioned already by many of the authors behind the fitting methods, see *e.g.* Meier-Hellstern (1987), Rossiter (1987), and Rydén (1992).

Model Errors

We will finally try to get an idea of the applicability of the models themselves, without respect to the particular fitting method used. Tables 7 and 8 show $\epsilon_{\text{mod}}(Q)$ and $\epsilon_{\text{mod}}(Q')$ in the same way as above.

It is noted that the model errors are of the same order as the total errors and larger than the method errors. Comparing to the former accuracy measures, the differences between the various fitting methods

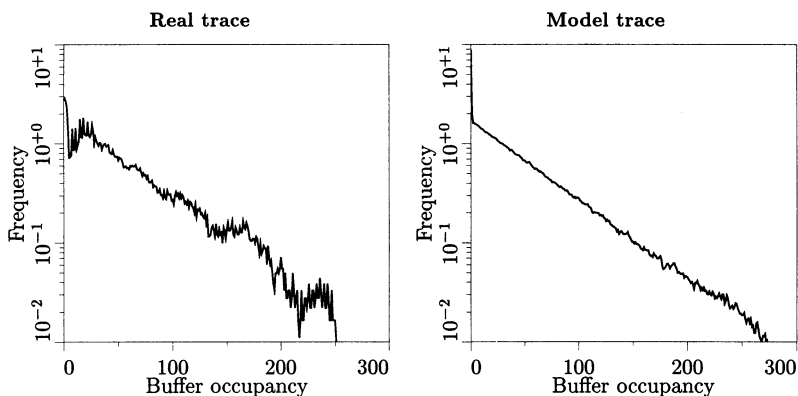


Figure 2 Comparison of buffer occupancies resulting from real and artificial traces for a trace of 2.00 seconds with 85% load. The artificial trace is produced by an SPP fitted by means of the HL-method. The upper plots refer to the whole trace and the lower ones to the first 2% of the trace.

remain, hence our attempt to separate the fitting method from the model is not entirely successful. The tables clearly show that no model succeeds in accurately predicting both the mean of the whole distribution and of its tail. As before, low values are almost exclusively noted in conjunction with a large number of failed fits. This makes it hard to point at any particularly successful or promising model.

Some Detailed Results

To get a deeper understanding of the results, we have arbitrarily selected a case for which reasonably good agreement was obtained in the study above, *viz.* direct fitting for 2 second intervals.

Two plots in figure 2 show the buffer occupancy distributions for the real and artificial traffics respectively. The two curves clearly appear quite similar at a first glance. On the other hand, at a closer look, the two differ around zero and in their tails: The real data has a lower value at the origin and exhibits a knee at the tail, while the model data has a higher value in the origin and the tail is straight. These findings are in agreement with what has been suggested by Pruthi (1995) and others: Markov-type models result in buffer occupancy distributions with exponential tails, while many real traffics result in power-law tails.

Looking at the similarity of the curves, these differences might be regarded as minor details, but it must be remembered that the models are to be used for determining loss probabilities on the order of 10^{-9} . Using our performance metrics, the similar shapes are reflected by the means, $E\{Q_0\} = 53.46$ for the real data and $E\{Q_1\} = 49.39$ for the model, while the different tails results in $E\{Q'_0\} = 251.8$ and $E\{Q'_1\} = 462.2$ respectively.

Figure 3 shows plots of the number of arrivals within intervals, $N(t, t + \Delta t)$ *vs.* t for the real and

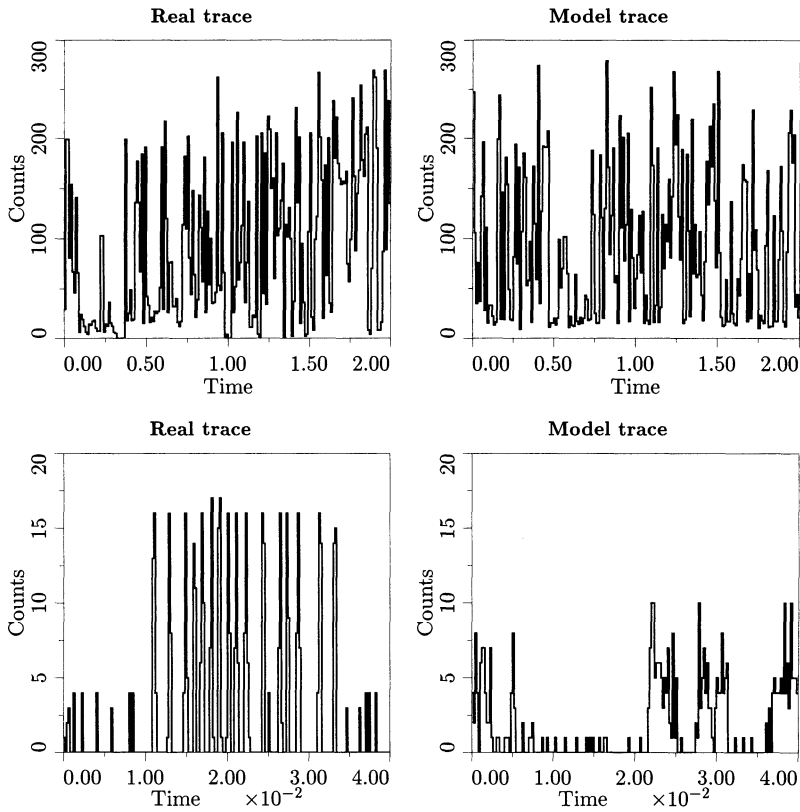


Figure 3 Comparison of real and artificial traces for a trace of 2.00 seconds with 85% load. The artificial trace is produced by an SPP fitted by means of the HL-method. The upper plots refer to the whole trace and the lower ones to the first 2% of the trace.

artificial traces. The upper plots refer to the whole trace and the lower ones show the first 2% of the trace in more detail.

It is noted that there are no apparent, fundamental differences in the large time scale between the real trace and the artificial one. However, it is also seen from the two lower plots that this statement does not seem to hold in the higher frequencies. This again confirms the results by Pruthi (1995) and others, that

the lower frequencies rule the average queue length, and hence the average delay, while higher frequencies are critical to the tail of the queue and therefore to the loss probability in case of an infinite buffer.

Noting a reasonable agreement for the overall mean, but a less good for the tail, it is tempting to conclude the models covered by our investigation might be more useful for calculating delays than losses. Tables 3–8 do, however, not support such a conclusion in general.

4 CONCLUSIONS AND FURTHER WORK

We have tried a number of methods proposed for fitting an MMPP (MMBP) to observed traffic data. Sixty data sets were extracted from the Bellcore Ethernet measurements according to length and local average load, so that short, medium and long periods of both light and heavy loads were tried. We then compared the performance of the buffer of a single server system when subject to the real traffic and when subject to traffic from the fitted model.

Several cases of infeasible parameters were recorded. A simple solution to this problem might be to restate the various methods as constrained optimisation problems, where a best fit under the condition of feasible parameters is determined.

It was found that the two cases differ significantly in terms of buffer occupancy, and that these differences are caused by deficiencies in the different fitting methods and possibly also by limitations in the models themselves.

Restricting ourselves to shorter time spans of up to two seconds, it was noted that direct metric fitting methods produced the smallest errors and resulted in traces that appeared identical to the real ones on a large time scale. It is therefore concluded that the most promising candidates for a “good model and fitting method” are found in the group, though further work is needed to clarify the importance of and methods for fitting a wider range of frequencies before “safe” models can be devised.

Moreover, if time scales of two seconds can be models, there is no reason why shorter time spans could not be mastered too if the problem of fitting to a small data set can be solved. On the other hand, it also seems clear from the tables that the chances of finding small MMPPs (MMBP) that remain valid for time scales of 20 seconds and above are fairly slim.

This work is different from what is normally published on modelling and fitting bursty traffics in that we use *real* traffic. This fact means that we have had to develop new practices and faced difficulties in ending up with neat conclusions regarding a perfect model and fitting method.

We believe, however, that there is enough real data available to stop validating models against models, but actually use real data instead. This work constitutes a first step in this direction, and we hope to have inspired others than ourselves to continue this important work. We can identify a large number of issues that need to be looked into, for instance

- Finding traffic characteristics which are relevant from the point of view of buffer dimensioning and for which simple and robust estimation techniques can be devised. Some theoretical proposals are given in *e.g.* Andrade *et al.* (1991) and Grünenfelder *et al.* (1994).
- Finding fitting methods for these characteristics which always come up with the best physically feasible fit. A first approach is to somewhat modify the methods tried here.
- A repeat of our investigation but with much more than ten samples per time scale and utilisation level.
- Repeating our investigation as above for other traffic sources than the Bellcore Ethernet.

We also note that if modelling short time scale variations shall be useful, a number of issues must be resolved, for example how to handle the long term variations in practice. Possible candidates includes reallocations of transmission capacity according to network predictions (e.g. by monitoring cell flow or buffer contents) or users' requests (e.g. when opening or closing particular applications or application modes).

5 ACKNOWLEDGEMENT

Thanks to Dr. Parag Pruthi at the Royal Institute of Technology (Sweden) for providing the extracts from the Bellcore measurements, to Dr. Tobias Rydén at the Lund Institute of Technology (Sweden) for providing the code for the TR method, and to Dr. Claes Jogr eus at the University of Karlskrona/Ronneby (Sweden) for discussions on statistical testing.

REFERENCES

- Andrade, J., Burakowski, W., and Villen-Altamirano, M. (1991) Characterization of Cell Traffic Generated by an ATM Source, in *Teletraffic and Datatraffic in a Period of Change*, Elsevier.
- Arvidsson,  ., Berry, L. and Harris, R. (1991) Performance Comparison of Bursty Traffic Models, in *Proc. Austr. Broadb. Switching and Services Symp.*, Sydney.
- Andersen, A., Jensen, A. and Friis Nielsen, B. (1995) Modelling and Performance Study of Packet-Traffic with Self-Similar Characteristics over Several Time Scales with Markovian Arrival Processes, in *Proc. 12th Nordic Teletraffic Sem.*, Helsinki.
- Bagnoli, G., Listanti, M. and Winkler, R. (1994) Cell Level and Frame Level Performance of Traffic Control Schemes for No Resource Reservation Data Communications in ATM Networks, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Bellcore, Measurements made on August, 29 1989 at 11.25 a.m. at *Bellcore Research and Engineering Centre*, Morristown.
- Bonomi, F., Meyer, J., Montagna, S. and Paglino, R. (1994) Minimal On/Off Source Models for ATM Traffic, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Frater, M., Tan, P. and Arnold, J. (1994) Variable Bit Rate Video Traffic on the Broadband ISDN: Modelling and Verification, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Gr unfelder, R. and Robert, S. (1994) Which Arrival Law Parameters are Decisive for Queueing System Performance, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Gusella, R. (1991) Characterizing the Variability of Arrival Processes with Indexes of Dispersion. *IEEE J. Sel. Areas in Commun.*, **9**, 203–211.
- Heffes, H. and Lucantoni, D. (1986) A Markov Modulated Characterization of Packetized Voice and Data Traffic and Related Statistical Multiplexer Performance. *IEEE J. Sel. Areas in Commun.*, **4**, 856–8.
- Hui, J. (1988) Resource Allocation for Broadband Networks. *IEEE J. Sel. Areas in Commun.*, **6**, 1598–608.
- Jain, R. and Routhier, S. (1986) Packet Trains — Measurements and a New Model for Computer Network Traffic. *IEEE J. Sel. Areas in Commun.*, **4**, 986–95.

- Key, P. (1995) Modelling, Measurement, and Connection Admission Control in ATM Networks, in *Proc. 9th ITC Specialists Seminar on Teletraffic Modelling and Measurement in Broadband and Mobile Communications*, Leidschendam.
- Lee, J. and Lee, B. (1992) Performance Analysis of ATM Cell Multiplexer with MMPP Input. *IEICE Trans. on Commun.*, **E75-B**, 709–14.
- Leland, W., Taqqu, M., Willinger, W. and Wilson, D. (1994) On the Self-Similar Nature of Ethernet Traffic (Extended Version). *IEEE/ACM Trans. on Networking.*, **2**, 1–15.
- Meier-Hellstern, K. (1987) A Fitting Algorithm for Markov-Modulated Poisson Processes Having Two Arrival Rates. *Eur. J. Op. Res.*, **29**, 370–7.
- Park, D. and Perros, H. (1994) m-MMBP Characterization of the Departure Process of an m-MMBP/Geo/1/K Queue, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Paxson, V. and Floyd, S. (1994) Wide-Area Traffic The Failure of Poisson Modeling, in *Proc. ACM Sigcomm 94*, London.
- Pruthi, P. (1995) An Application of Chaotic Maps to Packet Traffic Modeling. *Ph.D. dissertation*, Dept. of Teleinformatics, Royal Inst. of Tech., Stockholm.
- Ramamurthy, G. and Dighe, R. (1994) Analysis of Multilevel Hierarchical congestion Controls in B-ISDN, in *The Fundamental Role of Teletraffic in the Evolution of Telecommun. Netw.*, Elsevier.
- Rose, O. and Frater, M. (1994) A Comparison of Models for VBR Video Traffic Sources in B-ISDN, in *Proc. IFIP TC6/WG6.2 Broadband Commun. '94*, Paris.
- Rossiter, M. (1987) A Switched Poisson Model for Data Traffic. *Austr. Telecommun. Res.*, **21**, 53–7.
- Rydén, T. (1992) Parameter Estimation for Markov Modulated Poisson Processes. *Technical report (TFMS-LUTNFD23083)*, Department of Mathematical Statistics, Lund Institute of Technology, Lund.
- Solé, J., Domingo, J. and Garcia, J. (1990) Modelling the Bursty Characteristics of ATM Cell Streams, in *Proc. IEE Int. Conf. on Integr. Broadb. Services and Netw.*, London.
- Vakil, F. (1993) A Capacity Allocation Rule for ATM Networks, in *Proc. IEEE Globecom '93*, Houston.