
Department of Applied Mathematics
Faculty of EEMCS



University of Twente
The Netherlands

P.O. Box 217
7500 AE Enschede
The Netherlands

Phone: +31-53-4893400
Fax: +31-53-4893114

Email: memo@math.utwente.nl
www.math.utwente.nl/publications

Memorandum No. 1708

**Throughputs in processor sharing models
for integrated stream and elastic traffic**

R. LITJENS¹, J.L. VAN DEN BERG AND
R.J. BOUCHERIE

January, 2004

ISSN 0169-2690

¹Knowledge Innovation Center, TNO Telecom, The Netherlands

Throughputs in processor sharing models for integrated stream and elastic traffic

Remco Litjens*, Hans van den Berg*,[‡] and Richard J. Boucherie[†]

* Knowledge Innovation Center, TNO Telecom, The Netherlands

[†] Stochastic Operations Research, University of Twente, The Netherlands

[‡] DACS, University of Twente, The Netherlands

Abstract

We present an analytical study of throughput measures in Processor Sharing queueing systems with randomly varying service rates, modelling a communication link in an integrated services network carrying prioritised stream traffic and elastic traffic. A number of distinct throughput measures for the elastic traffic are defined and analysed. In particular, the differences between the various throughput measures and the impact of the elastic call size distribution are investigated. It is concluded that the call-average throughput, which is most relevant from the user point of view but typically hard to analyse, is very well approximated by the newly proposed so-called expected instantaneous throughput, which can easily be obtained from the system's steady state distribution.

Keywords: throughput, processor sharing, random environment.

AMS Subject classifications: primary 90B18, 90B22; secondary: 60K25.

1 INTRODUCTION

Processor sharing (PS) queueing models are widely applicable to situations where a common resource is shared by a varying number of concurrent users. In particular, PS models have been fruitfully applied in the field of the performance evaluation of computer systems and telecommunication networks. For instance, the PS service discipline appropriately models the design principle of fair resource sharing by TCP controlled elastic data flows or packet level scheduling schemes in e.g. IP, GPRS, UMTS and Wireless Local Area Networks (WLANS) [1, 2, 21, 22, 24].

The 'classical' PS model consists of a single server fairly sharing its fixed capacity among the varying number of present calls. A relevant extension is the PS queue with randomly varying service capacity, which models e.g. the impact of prioritised (speech) traffic on (low priority) video or data flows sharing a common network link. Important performance measures for PS queues are sojourn times and throughputs. In the queueing literature, the analyses of PS models are generally focussed towards the (conditional) expected *sojourn times* and their distribution, and many analytical results are available. In contrast, although the relevance is apparent from practical applications, *throughput* analyses are rare and only a few results are known. We therefore concentrate on the analysis and comparison of a variety of relevant throughput measures in PS models with fixed or randomly varying capacity.

LITERATURE In the literature the analysis of PS systems is primarily targeted towards *sojourn times* of jobs with a given service requirement (job size). Well-known results for the $M/G/1/PS$ queue are the linearity and insensitivity properties, i.e. the expected sojourn time of a tagged job is proportional to its service

requirement and independent of the service requirement distribution of the other jobs (see e.g. [20]). The sojourn time distribution for the $M/G/1/PS$ queue has been derived by Yashkov [39] and Ott [30]. Cohen [8] considers a generalisation of the $M/G/1/PS$ queue, viz. the so-called Generalised Processor Sharing (GPS) model, in which the service rate of the jobs is an arbitrary function of the number of jobs in the system. Note that e.g. the multiple server $M/G/c/PS$ queue and the classical Erlang loss model are special cases of the GPS model, which also possesses the linearity and insensitivity properties mentioned above. The reader is referred to [40] and [41] for overviews of the available results on ‘classical’ PS systems, see also the more recent paper by Zwart and Boxma [42] focusing on sojourn time asymptotics for the $M/G/1/PS$ queue with heavy tailed service requirement distributions (e.g. Pareto).

In the present paper PS systems with *randomly varying* service rates (e.g. due to the presence of higher priority jobs consuming part of the total service capacity) play a particularly important role. Randomly varying service rates severely complicate the analysis, and the nice properties of the expected sojourn time do not hold anymore. Nunez-Queija [27] analyses an $M/M/1/PS$ model with an on/off server, and derives closed-form expressions for several sojourn time statistics. In [29] Nunez-Queija et al. consider a GPS model with two priority classes, where each of the high priority jobs takes a fixed amount of the server capacity and the low priority jobs utilise the (fluctuating) remaining service capacity in a PS fashion. For this model, expressions for the (conditional) expected sojourn times of the low priority customers are derived. A generalisation and more extensive treatment of this work can be found in [26] and [28]. [23] presents and analytically supports the remarkable phenomenon that in the PS model with randomly varying capacity, the expected sojourn times are smaller if the job sizes are more variable, which is a relevant insight in light of the commonly acknowledged property that e.g. WWW pages are heavy tailed.

Throughput analyses of PS systems are rare in the literature. The only references known to the authors are by Kherani and Kumar [16, 17], who assess the PS service discipline as a model to evaluate the performance of TCP-controlled elastic (data) traffic in the Internet (cf. [25, 29, 33]). For the $M/G/1/PS$ model, [16, 17] compare different throughput measures by means of analysis and simulations. In a number of other papers with a larger scope, e.g. dimensioning of large IP networks, a seemingly arbitrary throughput measure is selected as a basis for the performance analysis, without substantiating the validity of such a measure. These papers are further discussed below; see also Section 2.3, where various throughput measures are mathematically defined.

CONTRIBUTION The principal objective of this paper is to derive and compare, both analytically and numerically, a variety of throughput performance measures in Processor Sharing models with fixed and varying service capacity. Aside from a substantial original contribution in the definition, analysis and comparison of throughput measures, the few known results have been included in order to also establish a survey character of the paper.

While from the customer’s perspective, the *call-average throughput* is the most relevant throughput measure, in PS systems the call-average throughput may be hard to determine analytically [16, 17], which is an important reason to assess the closeness of a number of other throughput measures. In several papers [11, 16, 17] the *time-average throughput*, defined as the expected throughput the ‘server’ provides to an elastic call at an arbitrary (non-idle) time instant, is applied to approximate the call-average throughput. Many other papers use the *ratio* of the expected transfer volume and the expected sojourn time as an approximation [1, 2, 3, 4, 9, 31]. In this paper we introduce the *expected instantaneous throughput*, i.e. the throughput an admitted call experiences immediately upon admission to the system, as a new throughput measure, which can be analysed relatively easily. The experiments demonstrate that the expected instantaneous throughput is the *only* one among these throughput measures which excellently approximates the call-average throughput for each of the investigated PS models and over the entire range of elastic traffic loads.

OUTLINE Section 2 specifies the traffic models, the investigated PS systems and the various throughput measures. An analytical evaluation is presented in Section 3. Section 4 presents and discusses the results of an extensive set of numerical experiments carried out to compare the different throughput measures for the different PS models. Although the analysis is of a generic nature and the results are certainly broadly applicable, the applied terminology in the numerical experiments is associated with the example context of a single cell in an integrated services GSM/GPRS network. The concluding remarks in Section 5 end the main body of this paper.

2 MODELS AND MEASURES

We consider a single server equipped with a channel pool of C traffic channels which can be assigned to stream calls, characterised by a fixed channel assignment (e.g. speech telephony), and elastic calls of two distinct types: (i) elastic calls whose sojourn time is unaffected by the (dynamically) assigned service rate (e.g. video telephony); and (ii) elastic calls whose sojourn time is affected by the assigned service rate (e.g. data transfer). In the remainder of the paper the three service types will be referred to by means of the given typical example services. The defining characteristics of the different services are given below, followed by the specification of the call handling procedures in four distinct performance models. An overview of the considered performance measures ends the section.

2.1 TRAFFIC MODELS

Three distinct service types are considered in the investigated processor sharing models:

SPEECH SERVICE Speech calls arrive according to a Poisson process with arrival intensity λ_{speech} and have a generally distributed duration with mean $1/\mu_{\text{speech}}$. A speech call requires a fixed assignment of one traffic channel. The speech traffic load is given by $\rho_{\text{speech}} \equiv \lambda_{\text{speech}}/\mu_{\text{speech}}$ and is expressed in Erlangs.

VIDEO SERVICE Video calls arrive according to a Poisson process with arrival intensity λ_{video} , have a generally distributed duration with mean $1/\mu_{\text{video}}$, and are elastic (*scalable*) in the ideal sense that the assigned number of traffic channels and thus the video quality can instantaneously and with perfect granularity adapt to the varying network load. The number of traffic channels that can be assigned to a video call is constrained by a maximum denoted $\beta_{\text{video}}^{\max}$. On the other hand, acceptable video quality is guaranteed by means of a minimum channel assignment of $\beta_{\text{video}}^{\min} \in [0, \beta_{\text{video}}^{\max}]$ traffic channels, corresponding to a bit rate of $r_{\text{video}}\beta_{\text{video}}^{\min}$ kbits/s, with r_{video} the effective video bit rate per traffic channel. Although effectively the video traffic load is influenced by $\beta_{\text{video}}^{\min}$, the definition of choice is $\rho_{\text{video}} \equiv \lambda_{\text{video}}/\mu_{\text{video}}$, as it allows the absence of an absolute QOS guarantee ($\beta_{\text{video}}^{\min} = 0$).

DATA SERVICE Data calls arrive according to a Poisson process with arrival intensity λ_{data} . A data call is assumed to be the transfer of a file with a generally distributed size, which is expressed in its nominal transfer time assuming a single dedicated traffic channel. The mean call size and effective data bit rate per traffic channel are denoted by $1/\mu_{\text{data}}$ and r_{data} (in kbits/s) respectively, corresponding with an actual mean transfer volume of $r_{\text{data}}/\mu_{\text{data}}$ kbits. Data calls are elastic in the sense that they are delay tolerant and can therefore tolerate a dynamic channel assignment, which affects the experienced throughput and thus the data call's sojourn time. As for the video calls, a maximum assignment denoted $\beta_{\text{data}}^{\max}$ is enforced to incorporate the terminals' technical limitations, while a possible QOS requirement is modelled by means of a minimum channel assignment $\beta_{\text{data}}^{\min}$. The data traffic load is given by $\rho_{\text{data}} \equiv \lambda_{\text{data}}/\mu_{\text{data}}$, while the normalised data traffic load is denoted $\rho_{\text{data}}^* \equiv \rho_{\text{data}}/C$.

Observe from the service specifications above that the key difference between video and data calls is the impact of the channel assignment on the calls' presence in the system. For video calls, the channel assignment influences the perceived audio and image quality experienced on the video terminal, while it does not affect the autonomously sampled video call duration. In case of data calls, the channel assignment affects the rate at which the file is transferred and thus the data call's sojourn time, which aside from the data throughput is a key performance measure in itself.

2.2 SYSTEM MODELS

Four distinct performance models are investigated, concentrating on one of the specific elastic services to be handled according to a PS service discipline with a fixed or varying capacity. In the latter case, the considered elastic service shares the server's aggregate capacity with a speech service, which utilises the server's capacity with preemptive priority, and thus implicitly leaves a time-varying residual capacity for the elastic calls. The different models are denoted V, D, SV and SD, and specified below. Let $S(t)$, $V(t)$ and $D(t)$ denote the process following the number of speech, video and data calls present at time $t \geq 0$, with states denoted s , v and d , respectively.

V MODEL In the V model video calls share the available C channels (fixed) in a PS fashion, i.e. given a presence of v video calls, each video call enjoys an instantaneous channel assignment of $\beta_{\text{video}}(v) \equiv \min\{C/v, \beta_{\text{video}}^{\max}\}$, obeying the constraint enforced by $\beta_{\text{video}}^{\max}$. The effective bit rate per channel is r_{video} kbits/s. In case of a positive minimum QOS requirement $\beta_{\text{video}}^{\min} > 0$, call admission control enforces a maximum presence of $v_{\max} \equiv \left\lfloor C/\beta_{\text{video}}^{\min} \right\rfloor$ video calls.

SV MODEL In the SV model the C traffic channels are dynamically shared by speech and video calls. Aside from the channels that are assigned to present video calls in order to meet their QOS requirement, the server's capacity is available with preemptive priority for speech calls. In other words, an arriving speech call is admitted if and only if $s + 1 \leq s_{\max}(v) \equiv \left\lfloor C - v\beta_{\text{video}}^{\min} \right\rfloor$, given a presence of s speech and v video calls. Analogously, if $\beta_{\text{video}}^{\min} > 0$, the condition for the admission of a video call is given by $v + 1 \leq v_{\max}(s) \equiv \left\lfloor (C - s)/\beta_{\text{video}}^{\min} \right\rfloor$. At any given time, the capacity that is not assigned to speech calls, is fairly shared by the present video calls in a PS fashion, i.e. each video call is assigned an instantaneous channel assignment of $\beta_{\text{video}}(s, v) \equiv \min\{(C - s)/v, \beta_{\text{video}}^{\max}\}$, which is guaranteed to exceed the minimum QOS requirement due to effects of the call admission control. Observe that the SV model is an example of a multi-rate model (see e.g. [14, 32]) incorporating speech and video calls with respective capacity requirements of 1 and $\beta_{\text{video}}^{\min}$ traffic channels.

D MODEL The D model is equivalent to the $M/G/1/d_{\max}/GPS$ queueing model with state-dependent aggregate service rates (due to $\beta_{\text{data}}^{\max}$) treated in [8], i.e. given a presence of d data calls, each such data call is assigned an instantaneous channel assignment of $\beta_{\text{data}}(d) \equiv \min\{C/d, \beta_{\text{data}}^{\max}\}$, with a per channel bit rate of r_{data} kbits/s. The call admission control threshold $d_{\max} \equiv \left\lfloor C/\beta_{\text{data}}^{\min} \right\rfloor$ is enforced if $\beta_{\text{data}}^{\min} > 0$.

SD MODEL In the SD model the C traffic channels are dynamically shared by speech and data calls. In line with the above specification of the SV model, the call admission control conditions for the admission of a speech or data call are given by $s + 1 \leq s_{\max}(d) \equiv \left\lfloor C - d\beta_{\text{data}}^{\min} \right\rfloor$ and $d + 1 \leq \left\lfloor (C - s)/\beta_{\text{data}}^{\min} \right\rfloor$ (only if $\beta_{\text{data}}^{\min} > 0$), respectively, given a presence of s speech and d data calls. At any given time, the capacity that is not assigned to speech calls, is fairly shared by the present data calls, i.e. each data call is assigned an instantaneous channel assignment of $\beta_{\text{data}}(s, d) \equiv \min\{(C - s)/d, \beta_{\text{data}}^{\max}\} \geq \beta_{\text{data}}^{\min}$.

We note that an integrated SVD model considering speech, video and data services jointly can also be analysed, see e.g., [22] where mean sojourn times are discussed. For our present purposes, however, the more focused models specified above are selected to convey the principal results most transparently.

2.3 THROUGHPUT MEASURES

In this subsection the definitions of the different performance measures are given. The definitions apply to both video and data calls. Denote with a_k (d_k) the arrival (departure) time of the k^{th} admitted elastic call, with $\tau_k \equiv d_k - a_k$ the call's sojourn time and with x_k the associated information volume (in kbits) transferred during its sojourn. Recall that for the video service the durations τ_k are autonomously sampled and the transfer volumes x_k are determined by the system dynamics, while for the data service the reverse holds. Let τ and x be the corresponding random variables with expected values $\mathbf{E}\{\tau\}$ and $\mathbf{E}\{x\}$. The *call-average* throughput is defined as

$$\mathbf{R}^c \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \frac{x_k}{\tau_k} = \mathbf{E}\left\{\frac{x}{\tau}\right\}. \quad (1)$$

With $N(t)$ the number of elastic calls present in the system and $C(t)$ the aggregate number of channels assigned to the elastic service at time $t \geq 0$, the *time-average* throughput is defined as

$$\mathbf{R}^t \equiv \lim_{t \rightarrow \infty} \frac{\frac{1}{t} \int_0^t \frac{rC(u)}{N(u)} \mathbf{1}\{N(u) \geq 1\} du}{\frac{1}{t} \int_0^t \mathbf{1}\{N(u) \geq 1\} du}, \quad (2)$$

where r denotes the effective information bit rate per traffic channel. Note that $N(t)$ is given by $V(t)$ in the (s)V model or $D(t)$ in the (s)D model, while $C(t)/N(t)$ is given by the channel assignment functions $\beta(\cdot)$. The time-average throughput is used to approximate the call-average throughput in e.g. [11, 16, 17]. We introduce the expected *instantaneous* throughput as

$$\mathbf{R}^i \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \frac{rC(a_k)}{N(a_k^+)}, \quad (3)$$

where $N(a_k^+)$ denotes the number of present elastic calls immediately after the k^{th} elastic call arrival and thus includes the new call. The *ratio* of the expected transfer volume and the expected sojourn time is defined as

$$\mathbf{R}^r \equiv \lim_{n \rightarrow \infty} \frac{\frac{1}{n} \sum_{k=1}^n x_k}{\frac{1}{n} \sum_{k=1}^n \tau_k} = \frac{\mathbf{E}\{x\}}{\mathbf{E}\{\tau\}}, \quad (4)$$

which is applied in e.g. [1, 2, 3, 4, 9, 31]. Note that \mathbf{R}^r can also be written as

$$\mathbf{R}^r = \frac{\lambda(1 - \mathbf{P}) \mathbf{E}\{x\}}{\lambda(1 - \mathbf{P}) \mathbf{E}\{\tau\}} = \lim_{t \rightarrow \infty} \frac{\frac{1}{t} \int_0^t rC(u) du}{\frac{1}{t} \int_0^t N(u) du},$$

where λ denotes the elastic call arrival rate and \mathbf{P} the elastic call blocking probability (see also below). This alternate expression for \mathbf{R}^r is given by the ratio of the long-term average aggregate system throughput and the long-term average number of elastic calls in the system. Its equivalence to expression (4) is due to the fact that in equilibrium the aggregate admitted bit rate must be equal to the aggregate processed bit rate (numerator)

and Little's law (denominator). As a final measure, the (unitless) call-average *stretch* (or the *normalised sojourn time*) is given by

$$\mathbf{S} \equiv \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \frac{\tau_k}{\left(\frac{x_k}{rC}\right)} = rC \mathbf{E} \left\{ \frac{\tau}{x} \right\}, \quad (5)$$

which is relevant for the data service only and is used as a performance measure in e.g. [13, 33]. For the special case of unrestricted channel assignments, i.e. $\beta_{\text{data}}^{\min} = \beta_{\text{video}}^{\min} = 0$ and $\beta^{\max} \geq C$, let $\tilde{\mathbf{R}}^c$, $\tilde{\mathbf{R}}^t$, $\tilde{\mathbf{R}}^i$, $\tilde{\mathbf{R}}^r$ and $\tilde{\mathbf{S}}$ denote the associated performance measures corresponding to the more general measures specified above.

In addition to these throughput measures, the included Call Admission Control schemes imply the occurrence of call blocking and thus the need to determine the speech, video and data call blocking probability (denoted \mathbf{P}), defined as the probability that an arriving call of a given type is denied admission to the system. Clearly, video or data calls experience blocking only if $\beta_{\text{video}}^{\min} > 0$ or $\beta_{\text{data}}^{\min} > 0$, respectively.

3 PERFORMANCE ANALYSIS

In this section we derive analytical expressions for the relevant performance measures in the four models specified above.

3.1 SV MODEL

Consider the SV model with generally distributed speech and video call durations. The evolution of the system in the SV model can then be described by the continuous-time stochastic process $(S(t), V(t))_{t \geq 0}$, with states denoted (s, v) . The process' state space is given by $\mathbb{S} \equiv \{(s, v) \in \mathbb{N}_0 \times \mathbb{N}_0 : s + v\beta_{\text{video}}^{\min} \leq C\}$. The unique equilibrium probability vector $\boldsymbol{\pi}$ of the stochastic process, given by

$$\pi(s, v) = \left(\sum_{(s,v) \in \mathbb{S}} \frac{\rho_{\text{speech}}^s}{s!} \frac{\rho_{\text{video}}^v}{v!} \right)^{-1} \frac{\rho_{\text{speech}}^s}{s!} \frac{\rho_{\text{video}}^v}{v!}, \quad (s, v) \in \mathbb{S},$$

is *insensitive* to the specific form of the speech and video call distributions, depending on their means only (see e.g. [14, 15, 32]). For the special case of unrestricted channel assignments to the video service, the state space is equal to $\tilde{\mathbb{S}} \equiv \{(s, v) \in \mathbb{N}_0 \times \mathbb{N}_0 : s \leq C\}$, and the equilibrium distribution is given by

$$\tilde{\pi}(s, v) = \exp(-\rho_{\text{video}}) \left(\sum_{s=0}^C \frac{\rho_{\text{speech}}^s}{s!} \right)^{-1} \frac{\rho_{\text{speech}}^s}{s!} \frac{\rho_{\text{video}}^v}{v!}, \quad (s, v) \in \tilde{\mathbb{S}}.$$

Using the PASTA property [38], the call blocking probabilities are readily derived from the equilibrium distribution:

$$\mathbf{P}_{\text{speech}} = \sum_{v=0}^{v_{\max}(0)} \pi(s_{\max}(v), v) \quad \text{and} \quad \mathbf{P}_{\text{video}} = \sum_{s=0}^C \pi(s, v_{\max}(s)).$$

In the case of unrestricted channel assignments to the video service, the speech call blocking probability is simply given by the Erlang loss probability, since speech traffic does not 'see' video traffic in the absence of video QoS guarantees, while the video call blocking probability equals zero.

3.1.1 CALL-AVERAGE THROUGHPUT

We start with a conditional analysis of the *call-average* throughput of a video call of a given duration τ which is admitted to the system in state (s, v) . In this conditional analysis we first confine ourselves to the case of *exponentially* distributed speech and video call durations.

For each state $(s, v) \in \mathbb{S}_{\text{video}}^+ \equiv \{(s, v) \in \mathbb{S} : v > 0\}$, denote with $\hat{x}_{s,v}(\tau)$ the conditional expected transfer volume of an admitted video call of duration τ , arriving at a given system state (s, v) , where v includes the new video call. The derivation involves a modified version of the Markov chain that is readily specified to describe the evolution of the SV model's stochastic process under the exponentiality assumption (see also [22]). Characterised by the presence of one permanent video call, the modified Markov chain consequently has the reduced state space $\mathbb{S}_{\text{video}}^+$. The video call departure rates in the associated infinitesimal generator $\mathcal{Q}_{\text{video}}^*$ reflect the presence of the permanent video call, i.e. $\mathcal{Q}_{\text{video}}^*((s, v); (s, v-1)) = (v-1)\mu_{\text{video}}$. The equilibrium distribution vector $\boldsymbol{\pi}_{\text{video}}^* \equiv (\pi_{\text{video}}^*(s, v), (s, v) \in \mathbb{S}_{\text{video}}^+)$, lexicographically ordered in (s, v) , of the modified Markov chain is, invoking reversibility and truncation of a reversible process [15], readily obtained as

$$\pi_{\text{video}}^*(s, v) = \frac{\pi(s, v-1)}{\sum_{(s', v') \in \mathbb{S}_{\text{video}}^+} \pi(s', v'-1)}, \quad (s, v) \in \mathbb{S}_{\text{video}}^+, \quad (6)$$

i.e. the equilibrium probabilities $\pi_{\text{video}}^*(s, v)$ corresponding to the modified Markov chain with one permanent video call are *equal* to the conditional probabilities that a newly admitted video call brings the system in state (s, v) in the original Markov chain. The equilibrium distribution $\boldsymbol{\pi}_{\text{video}}^*$ can readily be seen to be insensitive to the specific form of the speech and video call distributions [14, 15, 32]. Let $\mathcal{B}_{\text{video}} \equiv \text{diag}(\beta_{\text{video}}(s, v), (s, v) \in \mathbb{S}_{\text{video}}^+)$ denote the diagonal matrix of video channel assignments, lexicographically ordered in (s, v) .

As a special case of Theorem 1 in [22], for exponentially distributed video call durations the conditional expected video transfer volume vector $\hat{\mathbf{x}}(\tau) \equiv (\hat{x}_{s,v}(\tau), (s, v) \in \mathbb{S}_{\text{video}}^+)$, lexicographically ordered in (s, v) , is then given by

$$\hat{\mathbf{x}}(\tau) = \tau r_{\text{video}} (\boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}) \mathbf{1} + [\mathcal{I} - \exp\{\tau \mathcal{Q}_{\text{video}}^*\}] \boldsymbol{\gamma}_{\text{video}},$$

where $\boldsymbol{\gamma}_{\text{video}} \equiv (\gamma_{\text{video}}(s, v), (s, v) \in \mathbb{S}_{\text{video}}^+)$ is the unique solution to

$$\begin{aligned} \mathcal{Q}_{\text{video}}^* \boldsymbol{\gamma}_{\text{video}} &= r_{\text{video}} \{(\boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}) \mathbf{1} - \mathcal{B}_{\text{video}} \mathbf{1}\}, \\ \boldsymbol{\pi}_{\text{video}}^* \boldsymbol{\gamma}_{\text{video}} &= \mathbf{0}. \end{aligned} \quad (7)$$

The conditional expected (call-average) video throughput $\mathbf{R}_{\text{video}}^c(s, v, \tau)$ of a video call admitted to the system in state (s, v) and with a given holding time τ is then given by

$$\mathbf{R}_{\text{video}}^c(s, v, \tau) = \frac{\hat{x}_{s,v}(\tau)}{\tau} \quad (8)$$

(recall (1)), while deconditioning on the system state upon admission yields the conditional expected (call-average) video throughput of an admitted video call with duration τ , given by

$$\begin{aligned} \mathbf{R}_{\text{video}}^c(\tau) &= \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \left(\frac{\pi(s, v-1)}{\sum_{(s', v') \in \mathbb{S}_{\text{video}}^+} \pi(s', v'-1)} \right) \mathbf{R}_{\text{video}}^c(s, v, \tau) \\ &= \boldsymbol{\pi}_{\text{video}}^* \left\{ r_{\text{video}} (\boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}) \mathbf{1} + \frac{1}{\tau} [\mathcal{I} - \exp\{\tau \mathcal{Q}_{\text{video}}^*\}] \boldsymbol{\gamma}_{\text{video}} \right\} \\ &= r_{\text{video}} (\boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}) + \frac{1}{\tau} \boldsymbol{\pi}_{\text{video}}^* \left(\boldsymbol{\gamma}_{\text{video}} - \sum_{k=0}^{\infty} \frac{(\tau \mathcal{Q}_{\text{video}}^*)^k}{k!} \boldsymbol{\gamma}_{\text{video}} \right) \\ &= r_{\text{video}} \boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1} \\ &= r_{\text{video}} \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \left(\frac{\pi(s, v-1)}{\sum_{(s', v') \in \mathbb{S}_{\text{video}}^+} \pi(s', v'-1)} \right) \beta_{\text{video}}(s, v), \end{aligned}$$

using (7) and $\boldsymbol{\pi}_{\text{video}}^* \mathbf{Q}_{\text{video}}^* = \mathbf{0}$. Observe that $r_{\text{video}} \boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}$ is equal to the time-average video throughput in the SV model with one permanent video call (see also below). Comparing the first and last expression in the above derivation might confuse the reader into thinking that $\mathbf{R}_{\text{video}}^c(s, v, \tau)$ is simply equal to $r_{\text{video}} \beta_{\text{video}}(s, v)$, which is however readily seen to be not the case. Observe that $\mathbf{R}_{\text{video}}^c(\tau)$ does not depend on τ , so that the call-average video throughput is given by

$$\mathbf{R}_{\text{video}}^c = \int_{\tau=0}^{\infty} \mathbf{R}_{\text{video}}^c(\tau) \mu_{\text{video}} \exp\{-\tau \mu_{\text{video}}\} d\tau = \mathbf{R}_{\text{video}}^c(\tau) = r_{\text{video}} \boldsymbol{\pi}_{\text{video}}^* \mathcal{B}_{\text{video}} \mathbf{1}. \quad (9)$$

Whereas the above derivations utilised the exponentiality of the speech and video call durations, Theorem 1 claims that the obtained expressions for both $\mathbf{R}_{\text{video}}^c$ and $\mathbf{R}_{\text{video}}^c(\tau)$ (not $\mathbf{R}_{\text{video}}^c(s, v, \tau)$) are *insensitive* to the distributions of the speech and video call durations, apart from their means.

Theorem 1 *The call-average video throughput $\mathbf{R}_{\text{video}}^c$ and the conditional call-average video throughput $\mathbf{R}_{\text{video}}^c(\tau)$ are insensitive to the speech and video call duration distributions apart from their means.*

The proof of this theorem is presented in Appendix 1.

3.1.2 TIME-AVERAGE THROUGHPUT

Using the theory of regenerative processes (e.g. [37, 38]), the *time-average* video throughput is given by

$$\begin{aligned} \mathbf{R}_{\text{video}}^t &= \lim_{t \rightarrow \infty} \frac{\frac{1}{t} \int_0^t r_{\text{video}} \beta_{\text{video}}(S(u), V(u)) \mathbf{1}\{V(u) \geq 1\} du}{\frac{1}{t} \int_0^t \mathbf{1}\{V(u) \geq 1\} du} \\ &= r_{\text{video}} \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \left(\frac{\pi(s, v)}{\sum_{(s', v') \in \mathbb{S}_{\text{video}}^+} \pi(s', v')} \right) \beta_{\text{video}}(s, v), \end{aligned} \quad (10)$$

(cf. (2)), where $\pi(s, v) / \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v)$ is the equilibrium probability that the system is in state (s, v) , conditioned on the presence of at least one video call. The involved Césaro limits are derived using the renewal reward theorem [37, 38]. For the special case without channel assignment restrictions this yields

$$\tilde{\mathbf{R}}_{\text{video}}^t = \frac{r_{\text{video}}}{(\exp(\rho_{\text{video}}) - 1)} \left(\sum_{v=1}^{\infty} \frac{\rho_{\text{video}}^v}{v v!} \right) (C - \rho_{\text{speech}} (1 - \mathbf{P}_{\text{speech}})),$$

where $\mathbf{P}_{\text{speech}}$ is the Erlang loss probability. Note that the derivation of (10) does not require information on the specific form of the equilibrium distribution π . As this equilibrium distribution is insensitive to the call duration distribution (except for its mean), this property is inherited by the time-average video throughput.

3.1.3 EXPECTED INSTANTANEOUS THROUGHPUT

The expected *instantaneous* video throughput as defined in (3) is obtained as

$$\begin{aligned} \mathbf{R}_{\text{video}}^i &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n r_{\text{video}} \beta_{\text{video}}(S(a_k), V(a_k^+)) \\ &= r_{\text{video}} \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \left(\frac{\pi(s, v-1)}{\sum_{(s', v') \in \mathbb{S}_{\text{video}}^+} \pi(s', v'-1)} \right) \beta_{\text{video}}(s, v), \\ &= r_{\text{video}} \sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi_{\text{video}}^*(s, v) \beta_{\text{video}}(s, v), \end{aligned} \quad (11)$$

once again applying the theory of regenerative processes. As for the time-average throughput, the expected instantaneous video throughput measure inherits its insensitivity with respect to the specific form of the video call duration distribution from the insensitivity of π_{video}^* . Observe that the expected instantaneous video throughput is equal to the call-average video throughput, and hence so is the special case with unrestricted channel assignments.

3.1.4 RATIO THROUGHPUT MEASURE

The *ratio* of the expected video call transfer volume and the expected video call duration is given by

$$\mathbf{R}_{\text{video}}^r = \frac{\mathbf{E}\{\tau \mathbf{R}_{\text{video}}^c(\tau)\}}{\mu_{\text{video}}^{-1}} = \mathbf{R}_{\text{video}}^c$$

(cf. (4)), where the numerator is indeed equal to the expected transfer volume of a video call, using the fact that $\mathbf{R}_{\text{video}}^c(\tau) = \mathbf{R}_{\text{video}}^c$ does not depend on τ . It is readily seen that also for the special case of unrestricted channel assignments, the ratio throughput measure is equal to the corresponding call-average video throughput.

3.1.5 COMPARISON OF THROUGHPUT MEASURES

The *call-average* video throughput, the expected *instantaneous* video throughput and the *ratio* of the expected video call transfer volume and the expected video call duration appear to be identical, i.e.

$$\mathbf{R}_{\text{video}}^c = \mathbf{R}_{\text{video}}^i = \mathbf{R}_{\text{video}}^r,$$

and hence what remains is to compare these measures with the *time-average* throughput. Based on the explicit expressions (9) and (10), we will show in Theorem 2 for the case of $\beta_{\text{video}}^{\min} \in \{0, 1, \dots, C\}$ that the time-average throughput exceeds the call-average throughput: $\mathbf{R}_{\text{video}}^c \leq \mathbf{R}_{\text{video}}^t$. As an interesting corollary, we obtain that the time-average video throughput is monotonous in the offered video traffic load, i.e. $\frac{\partial \mathbf{R}_{\text{video}}^t}{\partial \rho_{\text{video}}} \leq 0$, which is noted to be non-trivial. While for $\rho_{\text{speech}} = 0$ (v model) this monotonicity can readily be concluded via stochastic monotonicity, for $\rho_{\text{speech}} > 0$ speech calls may take the place of video calls thus destroying stochastic monotonicity.

Theorem 2 *In the SV model with $\beta_{\text{video}}^{\min} \in \{0, 1, \dots, C\}$, the call-average video throughput is less than or equal to the time-average video throughput: $\mathbf{R}_{\text{video}}^c \leq \mathbf{R}_{\text{video}}^t$.*

The proof of this theorem and the proof of the following Corollary are given in Appendix 2 and Appendix 3 respectively.

Corollary 3 *The time-average video throughput is non-increasing in the video traffic load for $\beta_{\text{video}}^{\min} \in \{0, 1, \dots, C\}$, i.e.*

$$\frac{\partial \mathbf{R}_{\text{video}}^t}{\partial \rho_{\text{video}}} \leq 0$$

3.2 V MODEL

Since all relevant video throughput measures have been derived in closed-form for the SV model, including those for the case of unrestricted channel assignments, an explicit consideration of the v model would be superfluous, as it is merely a special case of the SV model with $\rho_{\text{speech}} = 0$. Also the ordering of the different throughput measures is as under the SV model.

3.3 SD MODEL

Consider the SD model with *exponentially* distributed speech call durations and data call sizes. The evolution of the system in the SD model can then be described by an irreducible two-dimensional continuous-time Markov chain $(S(t), D(t))_{t \geq 0}$, with states denoted (s, d) . The state space of the Markov chain is given by $\mathbb{S} \equiv \{(s, d) \in \mathbb{N}_0 \times \mathbb{N}_0 : s + d\beta_{\text{data}}^{\min} \leq C\}$, while its infinitesimal generator \mathcal{Q} is readily specified in terms of the speech and data call arrival and departure rates (see e.g. [21]). The irreducibility of the finite state space Markov chain $(S(t), D(t))_{t \geq 0}$ ensures the existence of a unique probability vector $\boldsymbol{\pi}$ that satisfies the system of global balance equations $\boldsymbol{\pi}\mathcal{Q} = \mathbf{0}$, with $\mathbf{0}$ the vector with all entries zero. The equilibrium distribution is *not* insensitive to the specific form of the speech call duration and data call size distributions. For the Markovian case, the equilibrium distribution can be determined numerically, e.g. by a successive overrelaxation procedure [37].

Using PASTA, the speech and data call blocking probabilities are given by

$$\mathbf{P}_{\text{speech}} = \sum_{d=0}^{d_{\max}^{(0)}} \pi(s_{\max}(d), d) \quad \text{and} \quad \mathbf{P}_{\text{data}} = \sum_{s=0}^C \pi(s, d_{\max}(s)).$$

In the special case of unrestricted channel assignments to the data service, the speech call blocking probability becomes equal to the Erlang loss probability, as speech traffic does not ‘see’ data traffic in the absence of data QoS guarantees, while the data call blocking probability becomes zero.

3.3.1 CALL-AVERAGE THROUGHPUT

Compared to other data throughput measures, obtaining explicit expressions for the *call-average* data throughput $\mathbf{R}_{\text{data}}^c$ is more involved. We first concentrate on the distribution of the data call sojourn times, conditional on the data call size. For each state $(s, d) \in \mathbb{S}_{\text{data}}^+ \equiv \{(s, d) \in \mathbb{S} : d > 0\}$ define $\tau_{s,d}(x)$ as the random time it takes to transfer a file of size x , arriving at a given system state (s, d) , where d includes the new data call. Define the Laplace-Stieltjes transform of the distribution of $\tau_{s,d}(x)$ by

$$T_{s,d}(\zeta, x) \equiv \mathbf{E} \{ \exp \{ -\zeta \tau_{s,d}(x) \} \}, \quad \text{Re}(\zeta) \geq 0, \quad (s, d) \in \mathbb{S}_{\text{data}}^+$$

and let $\mathbf{T}(\zeta, x) = (T_{s,d}(\zeta, x), (s, d) \in \mathbb{S}_{\text{data}}^+)$ be lexicographically ordered in $(s, d) \in \mathbb{S}_{\text{data}}^+$.

In an analogous manner as used to determine the conditional expected transfer volumes of video calls in the SV model, the derivation of an explicit expression for $\mathbf{T}(\zeta, x)$ involves a modified version of the original Markov chain, governed by infinitesimal generator $\mathcal{Q}_{\text{data}}^*$, characterised by the presence of one permanent data call, and with state space $\mathbb{S}_{\text{data}}^+$. The data call departure rates in the modified chain reflect the presence of the permanent data call, and are equal to $\mathcal{Q}_{\text{data}}^*((s, d); (s, d-1)) = \beta_{\text{data}}(s, d)(d-1)\mu_{\text{data}}$. Denote with $\boldsymbol{\pi}_{\text{data}}^*$ the unique equilibrium distribution of the modified Markov chain and let $\mathcal{B}_{\text{data}} \equiv \text{diag}(\beta_{\text{data}}(s, d), (s, d) \in \mathbb{S}_{\text{data}}^+)$ be the diagonal matrix of data channel assignments, lexicographically ordered in (s, d) . Partition $\mathbb{S}_{\text{data}}^+$ into $\mathbb{S}_{\text{data},0}^+ \equiv \{(s, d) \in \mathbb{S}_{\text{data}}^+ : \beta_{\text{data}}(s, d) = 0\}$ and its complement $\mathbb{S}_{\text{data},+}^+ \equiv \mathbb{S}_{\text{data}}^+ \setminus \mathbb{S}_{\text{data},0}^+$, and reorder the rows and columns in $\mathcal{Q}_{\text{data}}^*$, $\mathcal{B}_{\text{data}}$, $\boldsymbol{\pi}_{\text{data}}^*$ and $\mathbf{T}(\zeta, x)$ in accordance with the introduced state space partitioning, in order to allow the partitioning

$$\mathcal{Q}_{\text{data}}^* = \begin{bmatrix} \mathcal{Q}_{++}^* & \mathcal{Q}_{+0}^* \\ \mathcal{Q}_{0+}^* & \mathcal{Q}_{00}^* \end{bmatrix}, \quad \mathcal{B}_{\text{data}} = \begin{bmatrix} \mathcal{B}_+ & \mathcal{O} \\ \mathcal{O} & \mathcal{O} \end{bmatrix},$$

and

$$\boldsymbol{\pi}_{\text{data}}^* = (\boldsymbol{\pi}_{\text{data},0}^*, \boldsymbol{\pi}_{\text{data},+}^*), \quad \mathbf{T}(\zeta, x) = (\mathbf{T}_0(\zeta, x), \mathbf{T}_+(\zeta, x)),$$

where we omit the ‘data’ subscript in the submatrices of $\mathcal{Q}_{\text{data}}^*$ and $\mathcal{B}_{\text{data}}$ for enhanced readability. We note that in case $\beta_{\text{data}}^{\min} > 0$, this implies that $\mathbb{S}_{\text{data},0}^+ = \emptyset$, leading to a slightly simplified analysis (see [26, Section 4.2]).

As shown in [26, Section 4.4], for $x \geq 0$ and $\text{Re}(\zeta) \geq 0$, a closed-form expression for $\mathbf{T}(\zeta, x)$ is given by

$$\mathbf{T}_0(\zeta, x) = -(\mathcal{Q}_{00}^* - \zeta\mathcal{I})^{-1} \mathcal{Q}_{0+}^* \mathbf{T}_+(\zeta, x),$$

and

$$\mathbf{T}_+(\zeta, x) = \exp \left\{ x \mathcal{B}_+^{-1} \left(\mathcal{Q}_{++}^* - \mathcal{Q}_{+0}^* (\mathcal{Q}_{00}^* - \zeta\mathcal{I})^{-1} \mathcal{Q}_{0+}^* - \zeta\mathcal{I} \right) \right\} \mathbf{1}.$$

The conditional expected throughput $\mathbf{R}_{\text{data}}^c(s, d, x)$ of a data call admitted to the system in state (s, d) and with a given size x is given by

$$\begin{aligned} \mathbf{R}_{\text{data}}^c(s, d, x) &= r_{\text{data}} \mathbf{E} \left\{ \frac{x}{\tau_{s,d}(x)} \right\} \\ &= r_{\text{data}} \int_{\tau=0}^{\infty} \frac{x}{\tau} d\Phi_{s,d,x}(\tau) \\ &= r_{\text{data}} x \int_{\tau=0}^{\infty} \left(\int_{\zeta=0}^{\infty} \exp\{-\zeta\tau\} d\zeta \right) d\Phi_{s,d,x}(\tau) \\ &= r_{\text{data}} x \int_{\zeta=0}^{\infty} \left(\int_{\tau=0}^{\infty} \exp\{-\zeta\tau\} d\Phi_{s,d,x}(\tau) \right) d\zeta \\ &= r_{\text{data}} x \int_{\zeta=0}^{\infty} T_{s,d}(\zeta, x) d\zeta, \end{aligned}$$

where $\Phi_{s,d,x}(\tau)$ denotes the cumulative distribution function of $\tau_{s,d}(x)$ given data call size x and system state (s, d) upon the considered data call's admission. Deconditioning on the system state (s, d) upon admission yields

$$\mathbf{R}_{\text{data}}^c(x) = \sum_{(s,d) \in \mathbb{S}_{\text{data}}^+} \left(\frac{\pi(s, d-1)}{\sum_{(s',d') \in \mathbb{S}_{\text{data}}^+} \pi(s', d'-1)} \right) \mathbf{R}_{\text{data}}^c(s, d, x),$$

while subsequently deconditioning on the exponentially distributed data call size x gives the call-average data throughput:

$$\mathbf{R}_{\text{data}}^c = \mu_{\text{data}} \sum_{(s,d) \in \mathbb{S}_{\text{data}}^+} \left(\frac{\pi(s, d-1)}{\sum_{(s',d') \in \mathbb{S}_{\text{data}}^+} \pi(s', d'-1)} \right) \int_{x=0}^{\infty} \exp(-\mu_{\text{data}}x) \mathbf{R}_{\text{data}}^c(s, d, x) dx.$$

3.3.2 TIME-AVERAGE THROUGHPUT

By analogy with the derivation of the time-average video throughput in (10), the *time-average* data throughput is obtained as

$$\mathbf{R}_{\text{data}}^t = r_{\text{data}} \sum_{(s,d) \in \mathbb{S}_{\text{data}}^+} \left(\frac{\pi(s, d)}{\sum_{(s',d') \in \mathbb{S}_{\text{data}}^+} \pi(s', d')} \right) \beta_{\text{data}}(s, d).$$

Since the equilibrium distribution can only be numerically obtained, the above expression does not simplify for the special case of unrestricted channel assignments.

3.3.3 EXPECTED INSTANTANEOUS THROUGHPUT

Similar to the derivation of the corresponding measure (11) for the SV model, the expected *instantaneous* data throughput is given by

$$\mathbf{R}_{\text{data}}^i = r_{\text{data}} \sum_{(s,d) \in \mathbb{S}_{\text{data}}^+} \left(\frac{\pi(s, d-1)}{\sum_{(s',d') \in \mathbb{S}_{\text{data}}^+} \pi(s', d'-1)} \right) \beta_{\text{data}}(s, d).$$

3.3.4 RATIO THROUGHPUT MEASURE

The *ratio* of the expected data call size and the expected data call sojourn time is equal to

$$\mathbf{R}_{\text{data}}^r = \left(\frac{r_{\text{data}}}{\mu_{\text{data}}} \right) / \left(\frac{\sum_{(s,d) \in \mathbb{S}} d\pi(s, d)}{\lambda_{\text{data}}(1 - \mathbf{P}_{\text{data}})} \right) = r_{\text{data}} \frac{\rho_{\text{data}}(1 - \mathbf{P}_{\text{data}})}{\sum_{(s,d) \in \mathbb{S}} d\pi(s, d)},$$

where Little's formula (see e.g. [37]) is applied to derive the expected data call sojourn time.

3.3.5 CALL-AVERAGE STRETCH

Using

$$\mathbf{E} \left\{ \frac{\tau_{s,d}(x)}{x} \right\} = -\frac{1}{x} \frac{\partial}{\partial \zeta} T_{s,d}(\zeta, x) \Big|_{\zeta=0},$$

with $T_{s,d}(\zeta, x)$ as defined above, the expected (call-average) data stretch is given by

$$\begin{aligned} \mathbf{S}_{\text{data}} &= C \mathbf{E} \left\{ \frac{\tau_{s,d}(x)}{x} \right\} = \\ &= -C \mu_{\text{data}} \sum_{(s,d) \in \mathbb{S}_{\text{data}}^+} \left(\frac{\pi(s, d-1)}{\sum_{(s',d') \in \mathbb{S}_{\text{data}}^+} \pi(s', d'-1)} \right) \times \\ &\quad \times \left\{ \int_{x=0}^{\infty} \frac{1}{x} \exp(-\mu_{\text{data}} x) \left(\frac{\partial}{\partial \zeta} T_{s,d}(\zeta, x) \Big|_{\zeta=0} \right) dx \right\}, \end{aligned}$$

conform the definition given by (5), noting that in the above analysis the data call size x is expressed in units of r_{data} kbits (see also Section 2.1).

3.3.6 COMPARISON OF MEASURES

The expressions for the various throughput measures derived above for the SD model do not allow an analytical comparison. A numerical comparison is presented in Section 4.

3.4 D MODEL

The D model is a special case of the SD model with $\rho_{\text{speech}} = 0$. Moreover, the D model is equivalent to the $M/G/1/d_{\text{max}}/GPS$ queueing model with state-dependent aggregate service rates given by $dr_{\text{data}}\beta_{\text{data}}(d) = dr_{\text{data}} \min\{C/d, \beta_{\text{data}}^{\text{max}}\}$, see [8]. For this model, the equilibrium distribution is known to be *insensitive* to the specific form of the data call size distribution, and is given by

$$\pi(d) = \frac{(\rho_{\text{data}}^*)^d \phi(d)}{\sum_{d'=0}^{d_{\text{max}}} (\rho_{\text{data}}^*)^{d'} \phi(d')} \quad \text{with} \quad \phi(d) \equiv \left(\prod_{d'=1}^d \frac{d' \beta_{\text{data}}(d')}{C} \right)^{-1},$$

$d = 0, \dots, d_{\max}$, where $\rho_{\text{data}}^* \equiv \rho_{\text{data}}/C$ denotes the normalised data traffic load and $\phi(0) \equiv 1$ by convention. For the special case of unrestricted channel assignments, the D model reduces to the standard $M/G/1/PS$ queueing model, which has a geometric equilibrium distribution:

$$\tilde{\pi}(d) = (1 - \rho_{\text{data}}^*) (\rho_{\text{data}}^*)^d, \quad d \geq 0,$$

requiring $\rho_{\text{data}}^* < 1$ for stability.

Using PASTA, the data call blocking probability is equal to

$$\mathbf{P}_{\text{data}} = \pi(d_{\max}),$$

while it is equal to zero in the case of unrestricted channel assignments.

3.4.1 CALL-AVERAGE THROUGHPUT

In this section we assume exponentially distributed data call sizes. We first derive a closed-form expression for $\mathbf{T}(\zeta, x) \equiv (T_d(\zeta, x), d = 1, \dots, d_{\max})$ with $T_d(\zeta, x)$ the Laplace-Stieltjes transform of the distribution of $\tau_d(x)$, i.e. the random sojourn time of a data call of size x admitted to the system in the presence of $d - 1$ other data calls. Recall that x is expressed in the nominal sojourn time (in seconds). By analogy with the similar analysis presented for the SD model, $\mathcal{B}_{\text{data}}$ is the diagonal matrix of channel assignments and $\mathcal{Q}_{\text{data}}^*$ is the infinitesimal generator corresponding the D model's modified Markov chain with one permanent data call. In this data-only model, $\beta_{\text{data}}(d) > 0$ for all $d \geq 1$, so that no partitioning of $\mathbf{T}(\zeta, x)$ is required. As a specific instance of the result presented in [26, Section 4.2], for $x \geq 0$ and $\text{Re}(\zeta) \geq 0$, $\mathbf{T}(\zeta, x)$ is given by the closed-form expression

$$\mathbf{T}(\zeta, x) = \exp \{ x \mathcal{B}_{\text{data}}^{-1} (\mathcal{Q}_{\text{data}}^* - \zeta \mathcal{I}) \} \mathbf{1}.$$

By analogy with the analysis for the SD model, expressions for the conditional expected throughput measures $\mathbf{R}_{\text{data}}^c(d, x)$ and $\mathbf{R}_{\text{data}}^c(x)$ are readily derived. We limit ourselves here to stating the (unconditional) call-average data throughput:

$$\mathbf{R}_{\text{data}}^c = \mu_{\text{data}} \sum_{d=1}^{d_{\max}} \left(\frac{\pi(d-1)}{\sum_{d'=1}^{d_{\max}} \pi(d'-1)} \right) \int_{x=0}^{\infty} \exp(-\mu_{\text{data}} x) \left(r_{\text{data}} x \int_{\zeta=0}^{\infty} T_d(\zeta, x) d\zeta \right) dx.$$

For the case of unrestricted channel assignments, $\tilde{\mathbf{R}}_{\text{data}}^c(x)$ can be obtained using the following closed-form expression for the deconditioned Laplace-Stieltjes transform $\tilde{T}(\zeta, x)$ as derived in [6]:

$$\begin{aligned} \tilde{T}(\zeta, x) &\equiv \mathbf{E} \{ \exp \{ -\zeta \tau(x) \} \} = \sum_{d=1}^{\infty} \left(\frac{\pi(d-1)}{\sum_{d'=1}^{\infty} \pi(d'-1)} \right) \tilde{T}_d(\zeta, x) \\ &= \frac{(1 - \rho_{\text{data}}^*) (1 - \rho_{\text{data}}^* r^2) \exp \{ -(\lambda_{\text{data}} (1 - r) + \zeta) x \}}{(1 - \rho_{\text{data}}^* r)^2 - \rho_{\text{data}}^* (1 - r)^2 \exp \{ -\mu x (1 - \rho_{\text{data}}^* r^2) / r \}}, \end{aligned}$$

with $\text{Re}(\zeta) \geq 0$ and r given by

$$r = \frac{(\lambda_{\text{data}} + \mu_{\text{data}} + \zeta) - \sqrt{(\lambda_{\text{data}} + \mu_{\text{data}} + \zeta)^2 - 4\lambda_{\text{data}}\mu_{\text{data}}}}{2\lambda_{\text{data}}},$$

so that the conditional expected (call-average) data throughput is given by

$$\begin{aligned}
\tilde{\mathbf{R}}_{\text{data}}^c(x) &= \sum_{d=1}^{\infty} \pi(d-1) \left(r_{\text{data}} x \int_{\zeta=0}^{\infty} \tilde{T}_d(\zeta, x) d\zeta \right) \\
&= r_{\text{data}} x \int_{\zeta=0}^{\infty} \left(\sum_{d=1}^{\infty} \pi(d-1) \tilde{T}_d(\zeta, x) \right) d\zeta = r_{\text{data}} x \int_{\zeta=0}^{\infty} \tilde{T}(\zeta, x) d\zeta \\
&= r_{\text{data}} x \int_{\zeta=0}^{\infty} \frac{(1-\rho)(1-\rho r^2) \exp\{-(\lambda(1-r)+\zeta)x\}}{(1-\rho r)^2 - \rho(1-r)^2 \exp\{-\mu x(1-\rho r^2)/r\}} d\zeta.
\end{aligned}$$

3.4.2 TIME-AVERAGE THROUGHPUT

The *time-average* data throughput is given by

$$\mathbf{R}_{\text{data}}^t = r_{\text{data}} \sum_{d=1}^{d_{\max}} \left(\frac{\pi(d)}{\sum_{d'=1}^{d_{\max}} \pi(d')} \right) \beta_{\text{data}}(d),$$

while in the case of unrestricted channel assignments, the time-average data throughput is given by

$$\begin{aligned}
\tilde{\mathbf{R}}_{\text{data}}^t &= r_{\text{data}} \sum_{d=1}^{\infty} \left(\frac{(1-\rho_{\text{data}}^*)(\rho_{\text{data}}^*)^d}{\sum_{d'=1}^{\infty} (1-\rho_{\text{data}}^*)(\rho_{\text{data}}^*)^{d'}} \right) \frac{C}{d} \\
&= r_{\text{data}} C \left(\frac{1-\rho_{\text{data}}^*}{\rho_{\text{data}}^*} \right) \sum_{d=1}^{\infty} \left(\frac{(\rho_{\text{data}}^*)^d}{d} \right) \\
&= r_{\text{data}} C \left(\frac{1-\rho_{\text{data}}^*}{\rho_{\text{data}}^*} \right) \ln \left(\frac{1}{1-\rho_{\text{data}}^*} \right),
\end{aligned}$$

requiring $\rho_{\text{data}}^* < 1$ for stability. Note that due to the insensitivity of the equilibrium distribution, these expressions for the time-average throughput are also insensitive to the specific form of the data call size distribution.

3.4.3 EXPECTED INSTANTANEOUS THROUGHPUT

The expected *instantaneous* data throughput is given by

$$\mathbf{R}_{\text{data}}^i = r_{\text{data}} \sum_{d=1}^{d_{\max}} \left(\frac{\pi(d-1)}{\sum_{d'=1}^{d_{\max}} \pi(d'-1)} \right) \beta_{\text{data}}(d). \quad (12)$$

In the special case of unrestricted channel assignments, the expected instantaneous data throughput is equal to the time-average data throughput:

$$\tilde{\mathbf{R}}_{\text{data}}^i = r_{\text{data}} C \left(\frac{1-\rho_{\text{data}}^*}{\rho_{\text{data}}^*} \right) \ln \left(\frac{1}{1-\rho_{\text{data}}^*} \right),$$

requiring $\rho_{\text{data}}^* < 1$ for stability. Once again, the above expressions for the expected instantaneous throughputs inherit the insensitivity property of the equilibrium distribution.

3.4.4 RATIO THROUGHPUT MEASURE

The *ratio* of the expected data call size and the expected data call sojourn time is equal to

$$\mathbf{R}_{\text{data}}^r = r_{\text{data}} \frac{\rho_{\text{data}}(1 - \mathbf{P}_{\text{data}})}{\sum_{d=0}^{d_{\text{max}}} d\pi(d)},$$

while in the case of unrestricted channel assignments we have

$$\tilde{\mathbf{R}}_{\text{data}}^r = r_{\text{data}} C (1 - \rho_{\text{data}}^*),$$

requiring $\rho_{\text{data}}^* \leq 1$. Both expressions are insensitive to the data call size distribution aside from its mean.

3.4.5 CALL-AVERAGE STRETCH

The call-average *stretch* is given by

$$\begin{aligned} \mathbf{S}_{\text{data}} &= \mathbf{E}\{\mathbf{S}_{\text{data}}(x)\} = C \mathbf{E}\left\{\frac{\mathbf{T}_{\text{data}}(x)}{x}\right\} \\ &= C \mathbf{E}\left\{\frac{1}{x} \left(x \frac{\sum_{d=0}^{d_{\text{max}}} d\pi(d)}{\rho_{\text{data}}(1 - \mathbf{P}_{\text{data}})}\right)\right\} = \frac{\sum_{d=0}^{d_{\text{max}}} d\pi(d)}{\rho_{\text{data}}^*(1 - \mathbf{P}_{\text{data}})}, \end{aligned}$$

using the known linearity in x of the conditional expected sojourn time $\mathbf{T}_{\text{data}}(x)$ of a data call of size x [8, 37]. The call-average stretch for the case of unrestricted channel assignments is readily derived to be equal to

$$\tilde{\mathbf{S}}_{\text{data}} = \frac{1}{1 - \rho_{\text{data}}^*},$$

requiring $\rho_{\text{data}}^* < 1$ for stability. Note that the effect of the channel rate r_{data} is captured only in the definition of the data traffic load ρ_{data}^* .

3.4.6 COMPARISON OF MEASURES

We now present a number of results on relations between the different throughput measures derived above. Our first result relates the call average throughput and the ratio throughput measure.

Theorem 4 *For the D model,*

$$\mathbf{R}_{\text{data}}^c \geq \mathbf{R}_{\text{data}}^r. \quad (13)$$

Proof. The result is a straightforward extension of the equivalent result given in [16] for the case of unrestricted channel assignments. Applying Jensen's inequality (see e.g. [34]) with convex mapping $\psi(x) \equiv 1/x$:

$$\begin{aligned} \mathbf{R}_{\text{data}}^c &= r_{\text{data}} \mathbf{E}\left\{\psi\left(\frac{\mathbf{T}_{\text{data}}(x)}{x}\right)\right\} \\ &\geq r_{\text{data}} \psi\left(\mathbf{E}\left\{\frac{\mathbf{T}_{\text{data}}(x)}{x}\right\}\right) = r_{\text{data}} \left(\mathbf{E}\left\{\frac{1}{x} \left(x \frac{\sum_{d=0}^{d_{\text{max}}} d\pi(d)}{\rho_{\text{data}}(1 - \mathbf{P}_{\text{data}})}\right)\right\}\right)^{-1} \\ &= r_{\text{data}} \frac{\rho_{\text{data}}(1 - \mathbf{P}_{\text{data}})}{\sum_{d=0}^{d_{\text{max}}} d\pi(d)} = \mathbf{R}_{\text{data}}^r. \end{aligned}$$

■

We further adopt the following result for the case of unrestricted channel assignments and deterministic data call sizes.

Theorem 5 (Kherani and Kumar [16]) *In case of deterministic data call sizes, the following inequality holds:*

$$\tilde{\mathbf{R}}_{data}^t > \tilde{\mathbf{R}}_{data}^c. \quad (14)$$

Lastly, the explicitly derived expressions above revealed that, *only* for the case of unrestricted channel assignments, the time-average throughput is equal to the expected instantaneous throughput:

$$\tilde{\mathbf{R}}_{data}^t = \tilde{\mathbf{R}}_{data}^i,$$

while in general it holds that

$$\mathbf{R}_{data}^r \mathbf{S}_{data} = \tilde{\mathbf{R}}_{data}^r \tilde{\mathbf{S}}_{data} = r_{data} C.$$

4 NUMERICAL EXPERIMENTS

In this section we present the results from a set of numerical experiments, carried out in order to provide further insight in the throughput performance of elastic (video or data) calls in a system with a fixed or varying service capacity. A numerical comparison of throughput measures requires numerical values for the relevant parameters. To this end, we have selected the wireless environment of a GSM/GPRS cell with 22 channels. Section 4.1 below provides a more detailed motivation of these parameters. Next, Section 4.2 contains a numerical evaluation of the conditional call-average throughput in the V and D models as a function of the (exponentially distributed) elastic call size, the number of competing elastic calls found upon admission and the CAC threshold. Subsequently, in Section 4.3, an extensive numerical comparison is presented of the various (unconditional) throughput measures in the V, D, SV and SD models, considering different elastic call size distributions where relevant. As the results will demonstrate, the expected instantaneous throughput is the only throughput measure that closely approximates the call-average throughput for all considered scenarios.

4.1 PARAMETER SETTING

The system and traffic parameter settings are summarised in Table 1. The choice of the number of traffic channels C in the integrated services SV/SD models is based on a GSM/GPRS cell with 22 traffic channels (corresponding to 3 GSM frequencies *minus* 2 control channels). The capacity selected for the single service V/D models is equal to the average number of idle traffic channels in the SV/SD models, i.e. $22 - \rho_{speech} (1 - \mathbf{P}_{speech})$, where ρ_{speech} is chosen such that the corresponding speech call blocking probability is 1%. The speech call durations are exponentially distributed. An average call duration of 50 seconds is assumed for both the speech and video service. The average data file transfer is set at 320 kbits, which normalises to the given expected duration of μ_{data}^{-1} seconds. The video (data) bit rate per traffic channel is set to 13.4 (9.05) kbits/s, based on an assumed GPRS coding scheme CS-2 (CS-1). The video and data traffic loads are varied between 0 and the applicable value of C . Potential practical upper bounds on the channel assignment are disregarded. In the conditional throughput analyses for the V/D models, the minimum QoS requirements are varied within the range $[0, C]$, so that corresponding CAC thresholds between 1 and ∞ are considered, while no such restrictions are imposed for the unconditional throughput analyses.

4.2 CONDITIONAL THROUGHPUT RESULTS (V AND D MODELS)

We now present the results of the numerical conditional throughput analyses that have been carried out for the single service V and D models, respectively.

Table 1: Summary of the parameter settings assumed for the numerical experiments, based on the chosen context of a single cell in a GSM/GPRS network.

	SV model	V model	SD model	D model
C	22	8.486	22	8.486
μ_{speech}^{-1}	50 seconds	-	50 seconds	-
ρ_{speech}	13.651 Erlang	-	13.651 Erlang	-
μ_{video}^{-1}	50 seconds	50 seconds	-	-
ρ_{video}	$\in (0, C)$	$\in (0, C)$	-	-
r_{video}	13.4 kbits/s	13.4 kbits/s	-	-
$\beta_{\text{video}}^{\min}$	0 channels	$\in [0, C]$ channels	-	-
μ_{data}^{-1}	-	-	35.359 seconds	35.359 seconds
ρ_{data}	-	-	$\in (0, C)$	$\in (0, C)$
r_{data}	-	-	9.05 kbits/s	9.05 kbits/s
$\beta_{\text{data}}^{\min}$	-	-	0 channels	$\in [0, C]$ channels
$\beta_{\text{GPRS}}^{\max}$	C	C	C	C

V MODEL For the V model $\mathbf{R}^i = \mathbf{R}^c = \mathbf{R}^s$, so that the relevant throughput measures are \mathbf{R}^t and \mathbf{R}^c . Consider the case $\rho_{\text{video}} = \frac{1}{2}C = 11$, and a CAC threshold $v_{\max} = 10$, which is in our parameter setting achieved by taking $\beta_{\text{video}}^{\min} \in (0.7715, 0.8486]$. The left graph of Figure 1 shows the conditional call-average video throughputs (in kbits/s) for the case of exponentially distributed video call durations. A logarithmic scale is used for the video call duration τ (expressed in seconds). The results in the left chart assume a CAC threshold of $v_{\max} = 10$, which is achieved by setting $\beta_{\text{video}}^{\min} \in (0.7715, 0.8486]$, and leads to a video call blocking probability of $\mathbf{P}_{\text{video}} = 0.0075$. The depicted curve for $\mathbf{R}_{\text{video}}^c(v, \tau)$ is obtained using a special case of the result shown by (8), i.e. without speech traffic. As $\tau \downarrow 0$, the call-average throughput conditional on the system state v upon admission approaches $r_{\text{video}}\beta_{\text{video}}(v) = 113.7023/v$. As τ increases the impact of the system state upon admission vanishes and for each v the call-average throughput converges towards the time-average video throughput in a system with one permanent video call. Observe that for low (high) v , convergence is from above (below), in accordance with intuition. Further note that the qualitative results of these graphs are not affected by the actual parameter settings.

The right chart of Figure 1 depicts $\mathbf{R}_{\text{video}}^c(\tau)$ as function of the CAC threshold $v_{\max} \in \{1, 2, \dots, \infty\}$ established by varying $\beta_{\text{video}}^{\min} \in [0, C]$. As indicated in Section 3, $\mathbf{R}_{\text{video}}^c(\tau)$ is independent of the video call duration τ , with values exponentially decreasing from $r_{\text{video}}C = 113.7023$ for $v_{\max} = 1$ to $r_{\text{video}}C(1 - \exp(-\rho_{\text{video}}))/\rho_{\text{video}} = 26.4149$ for $v_{\max} = \infty$.

D MODEL For the D model all throughput measures considered in Section 3 are different. First, by analogy with the results above, consider the conditional call-average data throughput $\mathbf{R}_{\text{data}}^c(d, x)$. For the case of exponentially distributed call sizes, $\rho_{\text{data}} = \frac{1}{2}C = 11$ ($\rho_{\text{data}}^* = 0.5$), and a CAC threshold $d_{\max} = 10$ which is in our parameter setting achieved by setting $\beta_{\text{data}}^{\min} \in (0.7715, 0.8486]$ the left chart of Figure 2 depicts $\mathbf{R}_{\text{data}}^c(d, x)$. The profile of this chart is very similar to the left chart of Figure 1: for $x \downarrow 0$, $\mathbf{R}_{\text{data}}^c(x, d)$ is the instantaneous data throughput $r_{\text{data}}\beta_{\text{data}}(d) = 76.7915/d$, while $\lim_{x \rightarrow \infty} \mathbf{R}_{\text{data}}^c(d, x)$ is independent of d and given by the

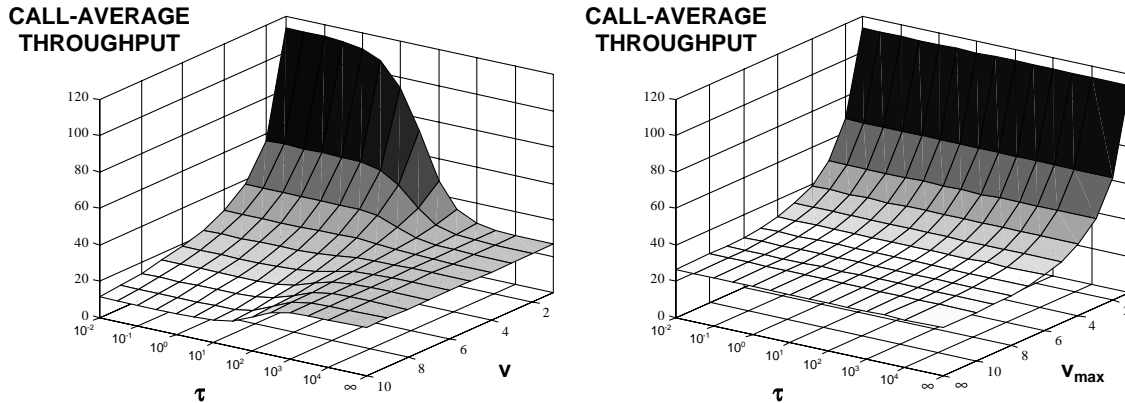


Figure 1: Conditional expected throughput performance in the v model. The left (right) chart shows the call-average throughput of a tagged video call as a function of its duration τ and the number of video calls v found upon admission (the CAC threshold v_{\max}).

time-average data throughput in data-only system with one permanent call, readily derived to be

$$r_{\text{data}} C \frac{(1 - \rho_{\text{data}}^*) (1 - (\rho_{\text{data}}^*)^{d_{\max}})}{(1 - (\rho_{\text{data}}^*)^{d_{\max}+1}) - (d_{\max} + 1) (\rho_{\text{data}}^*)^{d_{\max}} (1 - \rho_{\text{data}}^*)} = 38.5843. \quad (15)$$

In contrast with the v model, in the D model the time-average throughput in the adjusted Markov chain with one permanent data call is *not* equal to the call-average throughput in the original Markov chain.

The right chart shows $\mathbf{R}_{\text{data}}^c(x)$ for various CAC thresholds $d_{\max} \in \{1, 2, \dots, \infty\}$. In the trivial case of $d_{\max} = 1$, the call-average data throughput is equal to the aggregate service rate $r_{\text{data}} C = 76.7915$, independent of the data call size x . As d_{\max} increases, not only does $\mathbf{R}_{\text{data}}^c(x)$ decrease due to an increased carried data traffic load and hence a greater competition for resources, it is also no longer independent of x . For a given CAC threshold of d_{\max} , $\mathbf{R}_{\text{data}}^c(x)$ decreases from the corresponding expected instantaneous data throughput $\mathbf{R}_{\text{data}}^i$ (cf. expression (12)) to the expected time-average data throughput in the associated modified Markov chain with one permanent data call (cf. expression (15)). Unlike in the v model, in the D model small calls experience a higher average throughput than large calls. It is stressed, however, that the expected sojourn time is proportional in the data call size, so that the expected stretch is insensitive to the data call size. The potential confusion is due to the fact that the reciprocal of the expectation of a random variable is typically unequal to the expectation of the reciprocal of that random variable. Observe that the expected instantaneous throughput is an upper bound for the call-average throughput.

4.3 UNCONDITIONAL THROUGHPUT RESULTS (v , sv , D AND SD MODELS)

The remainder of this numerical section concentrates on the unconditional throughput as a function of the elastic traffic load, with a principal focus on the proximity of the various throughput measures in the different PS models.

(S)V MODEL Consider the sv and v models. Figure 3 depicts the various (unconditional) throughput performance measures as a function of the normalised elastic traffic load. In all considered cases channel

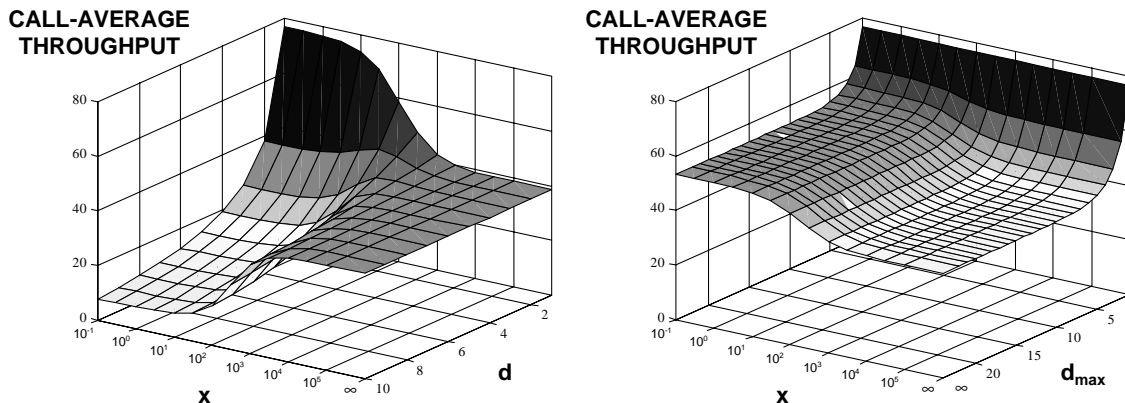


Figure 2: Conditional expected throughput performance in the D model. The left (right) chart shows the call-average throughput of a tagged data call as a function of its size x and the number of data calls d found upon admission (the CAC threshold d_{\max}).

assignment restrictions have been imposed on the elastic services. The left chart covers both the SV and the V models, for which all throughput measures are identical for any given normalised video traffic load $\rho_{\text{video}}^* \equiv \rho_{\text{video}}/C$, with C appropriately chosen in each model (see Table 1). The chart reveals both the demonstrated equality of $\mathbf{R}_{\text{video}}^c$, $\mathbf{R}_{\text{video}}^i$ and $\mathbf{R}_{\text{video}}^r$, and the proven ordering of $\mathbf{R}_{\text{video}}^t \geq \mathbf{R}_{\text{video}}^c$. It can be observed from the numerical results that $\mathbf{R}_{\text{video}}^t$ may exceed $\mathbf{R}_{\text{video}}^c$ by more than 36%.

D MODEL The right chart of Figure 3 concentrates on the D model. Since (only) the call-average throughput measure $\mathbf{R}_{\text{data}}^c$ is sensitive to the data call size distribution and no explicit expression could be derived, three distinct curves have been obtained via dynamic simulations for deterministic (zero variance), exponential and Pareto (with shape parameter $\alpha = 1.35$: infinite variance) data call size distributions. Sufficient numerical accuracy is ensured in the simulation experiment, indicated by a relative precision of the constructed 95% confidence intervals that is no worse than 5%. Observe that the call-average throughput is higher for more variable data call sizes, as also observed in [16], although the discrepancies are extremely small. This is probably due to the fact that a more variable data call size distribution features a relatively large number of small data calls, which appear to experience higher throughputs than large data calls (cf. the right chart of Figure 2).

As shown in Section 3, the insensitive time-average and expected instantaneous throughput measures are identical, and turn out to be a very good, only slightly overestimating (cf. (14)), approximation for the call-average throughput. Finally, $\mathbf{R}_{\text{data}}^r$ significantly underestimates the call-average throughput (cf. (13)), for high data traffic loads even by a factor exceeding 2.

SD MODEL For the SD model all the throughput measures are more or less sensitive to the data call size distribution, so that for reasons of clarity the numerical results are presented in the two separate charts of Figure 4 (for each marker in the legend, the left (right) throughput measure is depicted in the left (right) chart). In all cases, observe again that a more variable data call size distribution appears to lead to higher expected throughputs, which is in agreement with the sojourn time results of [23]. In this data model with varying service capacity, both the time-average throughput ($\mathbf{R}_{\text{data}}^t$) and the ratio of the expected data call size and the expected sojourn time ($\mathbf{R}_{\text{data}}^r$) are significantly lower than the call-average throughput ($\mathbf{R}_{\text{data}}^c$), in particular for lower data traffic loads. In contrast, the expected instantaneous throughput ($\mathbf{R}_{\text{data}}^i$) remains to

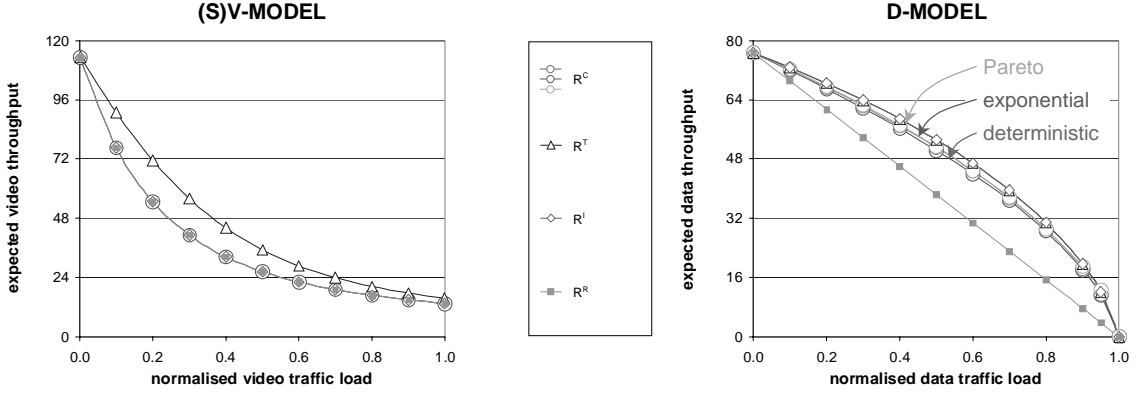


Figure 3: Comparison of different throughput measures in the SV, V and D models. The (insensitive) throughput measures in the left chart are identical for the SV and V models, given an appropriately normalised video traffic load. The right chart depicts for the D model the insensitive $\mathbf{R}_{\text{data}}^t$, $\mathbf{R}_{\text{data}}^i$ and $\mathbf{R}_{\text{data}}^r$ measures, along with the sensitive $\mathbf{R}_{\text{data}}^c$ measure for three distinct data call size distributions.

be a very good and fairly insensitive approximation for $\mathbf{R}_{\text{data}}^c$, across the entire range of data traffic loads. The *slight* overestimation of the call-average throughput seems to be not significant enough to lead to perilously loose Call Admission Control schemes or planning guidelines.

Comparing the throughput results for the D and SD models, observe that the call-average data throughput appears to be fairly insensitive to the variability of the available capacity, as also observed in [9] (recall that for the SV and V models, the call-average video throughputs were identical). Only for heavy data traffic loads, the call-average data throughput is non-negligibly higher for the fixed capacity D model.

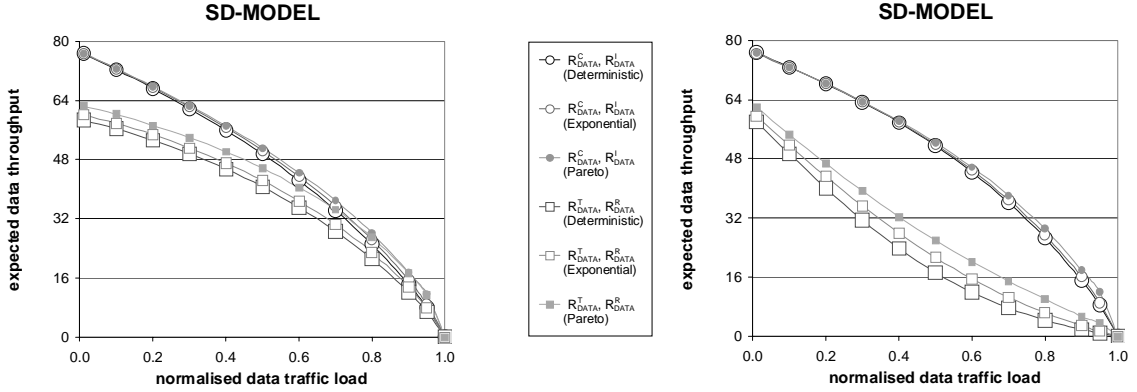


Figure 4: Comparison of different throughput measures in the SD model. All throughput measures are sensitive to the data call size distributions. The performance induced by three distinct distributions is shown.

In order to get a better grasp on the large discrepancy between e.g. the time- and call-average data throughputs in the SD model, the left chart of Figure 5 shows the time-average data throughput versus the normalised data traffic load for various degrees of acceleration of the speech call arrival and departure process. Keeping ρ_{speech} fixed at 13.651 Erlang, we multiply both λ_{speech} and μ_{speech} by the acceleration factor $\vartheta \in \{1, 10, 100, \infty\}$. The case of $\vartheta = 1$ refers to the original model and the associated curve is identical to the one for $\mathbf{R}_{\text{data}}^t$ in Figure

4 (left chart). At the other extreme, in the case of $\vartheta \rightarrow \infty$ the speech calls arrive and depart so quickly, that from the perspective of the data traffic, the available capacity is deterministic at $C - \rho_{\text{speech}}(1 - \mathbf{P}_{\text{speech}})$, and hence the accelerated model corresponds with the D model. As a consequence, the associated curve is identical to the one for $\mathbf{R}_{\text{data}}^t$ in Figure 3 (right chart). Observe that as the capacity fluctuation process is accelerated, i.e. when ϑ is increased from 1 to ∞ , the time-average throughput curves gradually approach the one corresponding to the extreme case of the D model, and the time-average throughput thus approximates the call-average throughput more and more closely.

The right chart of Figure 5 shows the expected stretch of a data call for both the SD and D models. As noted in Section 3, the expected stretch in the D model is insensitive to the data call size distribution. For the SD model, such insensitivity does not hold, as is demonstrated by the three expected stretch curves for deterministic, exponential and Pareto (with shape parameter $\alpha = 1.35$) data call size distributions. In correspondence with the throughput performance, the expected stretch appears to be smaller (better) for more highly variable data call sizes. A noteworthy observation from the numerical experiments that is not included in the figure, is that the expected stretch turns out to be infinitely large for the considered subexponential Weibull data call size distributions, i.e. with coefficient of variation greater than 1, for any data traffic load. In contrast, for highly variable Pareto distributions such as the one included in the figure, the expected stretch was nicely finite within the stable regime of data traffic loads. The probable reason for this phenomenon is that a subexponential Weibull distribution features many very small data calls, which may suffer from excessively large relative sojourn times in the case of a varying service capacity that is even equal to zero at times. Pareto distributions are inherently truncated at the lower end, however, so that extremely small data calls do simply not occur. In any case, the expected stretch thus appears to be less useful as a measure for throughput performance.

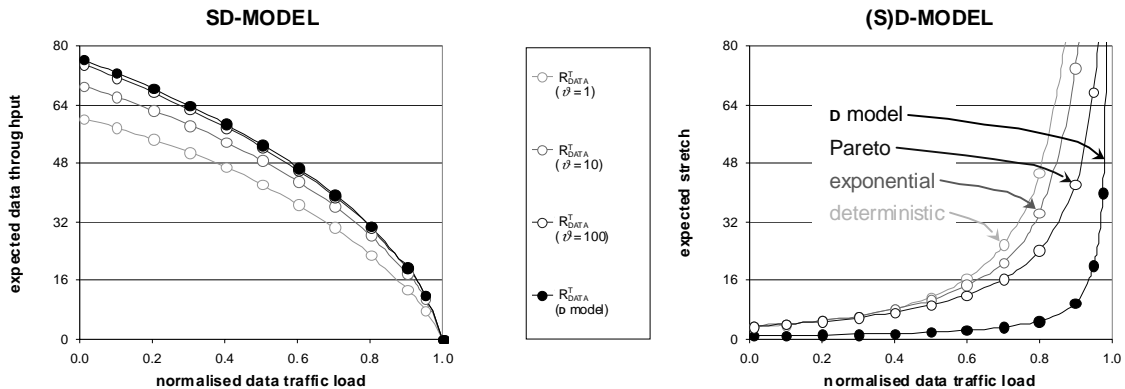


Figure 5: The impact of acceleration of the speech call arrival and departure process on $\mathbf{R}_{\text{data}}^t$ in the SD model (left chart). The expected stretch performance for different data call size distributions (SD model) as well as the insensitive values for the D model.

5 CONCLUDING REMARKS

In this paper we have specified, derived and compared, both analytically and numerically, a set of throughput measures in telecommunication systems serving elastic video or data calls according to a Processor Sharing service discipline. The available capacity was either fixed or randomly varying, corresponding with e.g. an integrated services network link, where the elastic calls utilise the capacity left idle by prioritised speech traffic.

Among the considered throughput measures, the call-average throughput is the most appropriate indicator of the experienced Quality of service. However, for models involving elastic calls of the data type, it is hard to determine this measure analytically. Among the alternative throughput measures, the newly proposed and readily analytically derived expected instantaneous throughput is the only measure which excellently approximates (or is even equal to) the call-average throughput in all considered system models and across the entire range of considered elastic traffic loads. In particular for the practically most relevant model integrating speech and data traffic, other typically applied throughput measures such as the time-average throughput or the ratio of the expected call size and the expected sojourn time, significantly underestimate the call-average throughput. An intuitive reasoning for the generally (near-)perfect fit of the expected instantaneous throughput is that apparently, the throughput an elastic call experiences immediately upon arrival is an excellent predictor of what the call is likely to experience throughout its lifetime. Moreover, among the considered throughput measures, the expected instantaneous throughput is the *only* approximate measure that is truly *call-centric*.

The analytical evaluation further revealed that the expected call-average throughput of elastic video calls in the considered PS models is *insensitive* to both the variability of the available capacity and the call duration distribution, while the numerical experiments indicated that this insensitivity property also holds for the data service to a considerable degree. As seen in [23], the latter insensitivity does not hold if the data performance is measured by the (conditional) expected sojourn time.

ACKNOWLEDGMENTS

This research has been carried out partly within the project BEYOND 3G supported by the Dutch Ministry of Economic Affairs under the program Technologische Samenwerking ICT-Doorbraakprojecten, project number TSIT1025 . The research is partly supported by the Technology Foundation STW, Applied Science Division of NWO and the technology programme of the DutchMinistry of Economic Affairs, The Netherlands.

References

- [1] S. Ben Fredj, T. Bonald, A. Proutiere, G. Régnié and J.W. Roberts, “Statistical bandwidth sharing: a study of congestion at flow level”, *Proceedings of SIGCOMM '01*, San Diego, USA, 2001.
- [2] N. Benameur, S. Ben Fredj, F. Delcoigne, S. Oueslati-Boulahia and J.W. Roberts, “Integrated admission control for streaming and elastic traffic”, *Proceedings of the 2nd International workshop on Quality of future Internet services*, Coimbra, Portugal, 2001.
- [3] T. Bonald and L. Massoulié, “Impact of fairness on Internet performance”, *Proceedings of SIGMETRICS '01*, Cambridge, USA, 2001.
- [4] T. Bonald and J.W. Roberts, “Performance of bandwidth sharing mechanisms for service differentiation in the Internet”, *Proceedings of the ITC specialist seminar on IP traffic measurement, modelling and management*, Monterey, USA, pp. 22.1-22.10, 2000.
- [5] E. Chlebus, “Empirical validation of call holding time distribution in cellular communication systems”, *Proceedings of ITC 15*, Washington DC, USA, pp. 1179-1188, 1997.
- [6] E.G. Coffman, R.R. Muntz and H. Trotter, “Waiting time distributions for processor-sharing systems”, *Journal of the Association for Computing Machinery*, vol. 17, pp. 123-130, 1970.

- [7] J.W. Cohen, "Some results on regular variation in queueing and fluctuation theory", *Journal of applied probability*, vol. 10, pp. 343-353, 1973.
- [8] J.W. Cohen, "The multiple phase service network with generalized processor sharing", *Acta informatica*, vol.12, pp. 245-284, 1979.
- [9] F. Delcoigne, A. Proutière and G. Régnié, "Modelling integration of streaming and data traffic", *Proceedings of the ITC specialist seminar on Internet traffic engineering and traffic management*, Würzburg, Germany, 2002.
- [10] N.M. van Dijk, "*Queueing networks and product forms: a system's approach*", John Wiley & Sons, New York, USA, 1993.
- [11] C. Douligeris, "Multiobjective flow control in telecommunication networks", *Proceedings of INFOCOM '92*, Florence, Italy, 1992.
- [12] P.R. Jelenković and P. Momčilović, "Resource sharing with subexponential distributions", *Proceedings of INFOCOM '02*, New York, USA, pp. 179-190, 2002.
- [13] N. Joshi, S.R. Kadaba, S. Patel and G.S. Sundaram, "Downlink scheduling in CDMA data networks", *Proceedings of MOBICOM '00*, Boston, USA, 2000.
- [14] J. S. Kaufman, "Blocking in a shared resource environment", *IEEE Transactions on communications*, vol. 29, no. 10, pp. 1474-1481, 1981.
- [15] F.P. Kelly, "*Reversibility and stochastic networks*", John Wiley & Sons, New York, USA, 1979.
- [16] A.A. Kherani and A. Kumar, "Performance analysis of TCP with nonpersistent sessions", *Proceedings of the Workshop on Modelling of flow and congestion control*, Paris, France, 2000.
- [17] A.A. Kherani and A. Kumar, "Stochastic models for throughput analysis of randomly arriving elastic flows in the Internet", *Proceedings of INFOCOM '02*, New York, USA, 2002.
- [18] L. Kleinrock, "Analysis of a time-shared processor", *Naval research logistics quarterly*, vol. 11, pp. 59-73, 1964.
- [19] L. Kleinrock, "Time-shared systems: a theoretical treatment", *Journal of the Association for Computing Machinery*, vol. 14, no. 2, pp. 242-261, 1967.
- [20] L. Kleinrock, "*Queueing systems, volume II*", John Wiley & Sons, New York, USA, 1976.
- [21] R. Litjens and R.J. Boucherie, "Radio resource sharing in a GSM/GPRS network", *Proceedings of the ITC Specialist seminar on Mobile systems and mobility*, Lillehammer, Norway, pp. 261-274, 2000.
- [22] R. Litjens and R.J. Boucherie, "Quality-of-service differentiation in an integrated services GSM/GPRS network", *to appear in: Advances in Performance Evaluation*, 2004.
- [23] R. Litjens and R.J. Boucherie, "Elastic calls in an integrated services network: the greater the call size variability the better the Quality-Of-Service", *Performance evaluation*, vol. 52, no. 4, pp. 193-220, 2003.
- [24] R. Litjens, F. Roijers, J.L. van den Berg, R.J. Boucherie and M. Fleuren, "Performance analysis of WLANs: an integrated packet/flow level approach", *Proceedings of ITC18*, Berlin, Germany, September 2003.

- [25] L. Massoulié and J.W. Roberts, “Arguments in favour of admission control for TCP flows”, *Proceedings of ITC 16*, Edinburgh, Scotland, 1999.
- [26] R. Núñez Queija, “Processor-sharing models for integrated-services networks”, Ph.D. thesis, Technische Universiteit Eindhoven, 1999.
- [27] R. Núñez Queija, “Sojourn times in a processor sharing queue with service interruptions”, *Queueing systems*, vol. 34, pp. 351-386, 2000.
- [28] R. Núñez Queija, “Sojourn times in non-homogeneous QBD processes with processor sharing”, *Stochastic models*, vol. 17, pp. 61-92, 2001.
- [29] R. Núñez Queija, J.L. van den Berg, and M.R.H. Mandjes, “Performance evaluation of strategies for integration of elastic and stream traffic”, *Proceedings of ITC 16*, Edinburgh, Scotland, pp. 1039-1050, 1999.
- [30] T. Ott, “The sojourn time distribution in the $M/G/1$ queue with processor sharing”, *Journal of applied probability*, vol. 21, pp. 360-378, 1984.
- [31] S. Rácz, B.P. Gerö and G. Fodor, “Flow level performance analysis of a multi-service system supporting elastic and adaptive services”, *Performance evaluation*, vol. 49, no. 1/4, pp. 451-469, 2002.
- [32] J. W. Roberts, “A service system with heterogeneous user requirements - application to multiservice telecommunications systems”, *Performance of data communications systems and their applications*, G. Pujolle (editor), North Holland, New York, USA, 1981.
- [33] J.W. Roberts and L. Massoulié, “Bandwidth sharing and admission control for elastic traffic”, *Proceedings of the ITC specialist seminar on Teletraffic issues related to multimedia and nomadic communications*, Yokohama, Japan, 1998.
- [34] S.M. Ross, “*A first course in probability*”, Macmillan Publishing Company, New York, USA, 1984.
- [35] M. Sakata, S. Noguchi and J. Oizumi, “Analysis of a processor-based queueing model for time-sharing systems”, *Proceedings of the 2nd Hawaii international conference on System sciences*, USA, pp. 625-628, 1969.
- [36] M. Sakata, S. Noguchi and J. Oizumi, “An analysis of the $M/G/1$ queue under round robin scheduling”, *Operations research*, vol. 19, pp. 371-385, 1971.
- [37] H.C. Tijms, “*Stochastic modelling and analysis: a computational approach*”, John Wiley & Sons, Chichester, England, 1986.
- [38] R.W. Wolff, “*Stochastic modeling and the theory of queues*”, Prentice-Hall, Englewood Cliffs, USA, 1989.
- [39] S.F. Yashkov, “A derivation of response time distribution for a $M/G/1$ processor sharing queue”, *Problems in control and information theory*, vol. 12, pp. 133-148, 1983.
- [40] S.F. Yashkov, “Processor sharing queues: some progress in analysis”, *Queueing systems*, vol. 2, pp. 1-17, 1987.
- [41] S.F. Yashkov, “Mathematical problems in the theory of processor sharing queueing systems”, *Journal of Soviet mathematics*, vol. 58, pp. 101-147, 1992.

[42] A.P. Zwart and O.J. Boxma, "Sojourn time asymptotics in the $M/G/1$ processor sharing queue", *Queueing systems*, vol. 35, pp. 141-166, 2000.

Appendix 1: Proof of Theorem 1

Proof. The stationary joint distribution $\pi(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}})$ of the number of speech (S) and video calls (V) present in the system and the associated residual call durations $\boldsymbol{\Theta}_{\text{speech}} \equiv (\Theta_{\text{speech}}(1), \dots, \Theta_{\text{speech}}(S))$ and $\boldsymbol{\Theta}_{\text{video}} \equiv (\Theta_{\text{video}}(1), \dots, \Theta_{\text{video}}(V))$ is given by (see e.g. [10])

$$\begin{aligned} \pi(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}}) &= \Pr\{S = s, V = v, \boldsymbol{\Theta}_{\text{speech}} \in [\boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{speech}} + d\boldsymbol{\vartheta}_{\text{speech}}], \\ &\quad \boldsymbol{\Theta}_{\text{video}} \in [\boldsymbol{\vartheta}_{\text{video}}, \boldsymbol{\vartheta}_{\text{video}} + d\boldsymbol{\vartheta}_{\text{video}}]\} \\ &= G(\rho_{\text{speech}}, \rho_{\text{video}}, C) \left\{ \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s! v!} \prod_{s'=1}^s \left(\frac{\bar{\Phi}_{\text{speech}}(\vartheta_{\text{speech}}(s'))}{\mu_{\text{speech}}^{-1}} d\vartheta_{\text{speech}}(s') \right) \times \right. \\ &\quad \left. \prod_{v'=1}^v \left(\frac{\bar{\Phi}_{\text{video}}(\vartheta_{\text{video}}(v'))}{\mu_{\text{video}}^{-1}} d\vartheta_{\text{video}}(v') \right) \right\}, \end{aligned}$$

for $(s, v) \in \mathbb{S} = \mathbb{S}(C) \equiv \{(s, v) \in \mathbb{N}_0 \times \mathbb{N}_0 : s + v\beta_{\text{video}}^{\min} \leq C\}$, $\boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}} \geq \mathbf{0}$, where the vectors $d\boldsymbol{\vartheta}_{\text{speech}}$ and $d\boldsymbol{\vartheta}_{\text{video}}$ consist of infinitesimally small elements,

$$G(\rho_{\text{speech}}, \rho_{\text{video}}, C) \equiv \left(\sum_{(s,v) \in \mathbb{S}(C)} \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s! v!} \right)^{-1},$$

and where $\bar{\Phi}_{\text{speech}}$ and $\bar{\Phi}_{\text{video}}$ denote the complementary cumulative distributions of the speech and video call durations, respectively.

Using PASTA, the joint distribution $\pi_{\text{video}}^{\bullet}(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}})$ of $(S, V, \boldsymbol{\Theta}_{\text{speech}}, \boldsymbol{\Theta}_{\text{video}})$ upon *admission* of a tagged video call is readily given by

$$\begin{aligned} \pi_{\text{video}}^{\bullet}(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}}) &= \Pr\{S = s, V = v, \boldsymbol{\Theta}_{\text{speech}} \in [\boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{speech}} + d\boldsymbol{\vartheta}_{\text{speech}}], \\ &\quad \boldsymbol{\Theta}_{\text{video}} \in [\boldsymbol{\vartheta}_{\text{video}}, \boldsymbol{\vartheta}_{\text{video}} + d\boldsymbol{\vartheta}_{\text{video}}] \mid s + v\beta_{\text{video}}^{\min} \leq C - \beta_{\text{video}}^{\min}\} \\ &= G(\rho_{\text{speech}}, \rho_{\text{video}}, C - \beta_{\text{video}}^{\min}) \times \\ &\quad \left\{ \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s! v!} \prod_{s'=1}^s \frac{\bar{\Phi}_{\text{speech}}(\vartheta_{\text{speech}}(s'))}{\mu_{\text{speech}}^{-1}} \prod_{v'=1}^v \frac{\bar{\Phi}_{\text{video}}(\vartheta_{\text{video}}(v'))}{\mu_{\text{video}}^{-1}} \right\}, \end{aligned}$$

for $(s, v) \in \mathbb{S}(C - \beta_{\text{video}}^{\min})$, where v excludes the newly admitted tagged video call.

Observe that $\pi_{\text{video}}^{\bullet}(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}})$ is equal to the stationary joint distribution of the number of speech and video calls and their residual call durations in a corresponding system with capacity $C - \beta_{\text{video}}^{\min}$ instead of C , or equivalently, in the original system but with one *permanent* video call (where v excludes this call). Hence the system state remains stochastically identical throughout the duration of the tagged video call. The associated (partially deconditioned) system state distribution $\pi_{\text{video}}^{\bullet}(s, v)$ is given by

$$\begin{aligned} \pi_{\text{video}}^{\bullet}(s, v) &= \int_{\vartheta_{\text{speech}}(1)=0}^{\infty} \int_{\vartheta_{\text{speech}}(s)=0}^{\infty} \cdots \int_{\vartheta_{\text{video}}(1)=0}^{\infty} \int_{\vartheta_{\text{video}}(v)=0}^{\infty} \pi_{\text{video}}^{\bullet}(s, v, \boldsymbol{\vartheta}_{\text{speech}}, \boldsymbol{\vartheta}_{\text{video}}) \\ &= G(\rho_{\text{speech}}, \rho_{\text{video}}, C - \beta_{\text{video}}^{\min}) \left\{ \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s! v!} \right\}, \end{aligned} \tag{16}$$

for $(s, v) \in \mathbb{S}(C - \beta_{\text{video}}^{\min})$. Since the throughput of the tagged video call is completely determined by the distribution of the number of speech and *other* video calls present during its lifetime, as given in (16), it is now immediately clear that the conditional call-average throughput $\mathbf{R}_{\text{video}}^c(\tau)$ of the tagged video call is *independent* of its duration τ , i.e. $\mathbf{R}_{\text{video}}^c(\tau) = \mathbf{R}_{\text{video}}^c$, for all $\tau \geq 0$. In particular, it is equal to the expected instantaneous video throughput experienced upon admission, which inherits its insensitivity from the insensitivity of π_{video}^* (see also Section 3.1.3 below). ■

Appendix 2: Proof of Theorem 2

Proof. For the extreme cases of infinitesimally small or infinitely large video traffic loads, it is readily argued that the call- and time-average video throughput measures are identical. Under an extremely *light* video traffic load ($\rho_{\text{video}} \downarrow 0$), a (rarely) occurring system state $(s, v) \in \mathbf{S}_{\text{video}}^+$ must have $v = 1$, almost surely, for both the original stochastic process, and the modified process with one permanent video call. As a consequence, the time-average video throughputs of both processes are identical, and hence so are the call- and time-average video throughputs of the original process. We thus have that

$$\lim_{\rho_{\text{video}} \downarrow 0} \mathbf{R}_{\text{video}}^t = \lim_{\rho_{\text{video}} \downarrow 0} \mathbf{R}_{\text{video}}^c$$

as can readily be verified from (9) and (10).

Alternatively, an infinitely *heavy* video traffic load ($\rho_{\text{video}} \rightarrow \infty$, assuming $\beta_{\text{video}}^{\min} > 0$ for stability) leads to a (complete or near) crowding out of speech calls, and implies the everlasting presence of $v_{\max}(0) = \lceil C_{\text{total}}/\beta_{\text{video}}^{\min} \rceil \geq 1$ video calls, and hence again the performance of the original and the modified process are the same. In particular, all video throughput measures are identical, so that

$$\lim_{\rho_{\text{video}} \rightarrow \infty} \mathbf{R}_{\text{video}}^t = \lim_{\rho_{\text{video}} \rightarrow \infty} \mathbf{R}_{\text{video}}^c$$

Now assume that $0 \leq \rho_{\text{video}} < \infty$. Then, from (9) and (10),

$$\begin{aligned} & \mathbf{R}_{\text{video}}^c \leq \mathbf{R}_{\text{video}}^t \\ \Leftrightarrow & r_{\text{video}} \sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \left(\frac{\pi(s, v-1)}{\sum_{(s',v') \in \mathbf{S}_{\text{video}}^+} \pi(s', v'-1)} \right) \beta_{\text{video}}(s, v) \\ & \leq r_{\text{video}} \sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \left(\frac{\pi(s, v)}{\sum_{(s',v') \in \mathbf{S}_{\text{video}}^+} \pi(s', v')} \right) \beta_{\text{video}}(s, v) \\ \Leftrightarrow & \left(\sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^{v-1}}{s!(v-1)!} \beta_{\text{video}}(s, v) \right) \left(\sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s!v!} \right) + \\ & - \left(\sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^v}{s!v!} \beta_{\text{video}}(s, v) \right) \left(\sum_{(s,v) \in \mathbf{S}_{\text{video}}^+} \frac{\rho_{\text{speech}}^s \rho_{\text{video}}^{v-1}}{s!(v-1)!} \right) \leq 0 \\ \Leftrightarrow & \left(\sum_{v=0}^{v_{\max}-1} \rho_{\text{video}}^v \sum_{s=0}^{C-\beta^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!v!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{w=0}^{v_{\max}-1} \rho_{\text{video}}^w \sum_{s=0}^{C-\beta^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!(w+1)!} \right) + \\ & - \left(\sum_{v=0}^{v_{\max}-1} \rho_{\text{video}}^v \sum_{s=0}^{C-\beta^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!(v+1)!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{w=0}^{v_{\max}-1} \rho_{\text{video}}^w \sum_{s=0}^{C-\beta^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!w!} \right) \leq 0 \end{aligned}$$

Recognising that the LHS is a polynomial in ρ of degree $2(v_{\max} - 1) = 2\left(\left\lfloor C/\beta_{\text{video}}^{\min} \right\rfloor - 1\right)$, the above condition can be written in the following form:

$$\sum_{k=0}^{2(v_{\max}-1)} \rho_{\text{video}}^k \sum_{v+w=k} \zeta_{v,w} \leq 0, \quad (17)$$

where the coefficients $\zeta_{v,w}$, $v, w, = 0, \dots, v_{\max} - 1$, are given by

$$\begin{aligned} \zeta_{v,w} &= \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!v!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!(w+1)!} \right) + \\ &\quad - \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!(v+1)!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!w!} \right) \\ &= \frac{1}{v!w!} \left(\frac{1}{w+1} - \frac{1}{v+1} \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!} \right). \end{aligned}$$

Note that $\zeta_{v,v} = 0$, $v = 0, \dots, v_{\max} - 1$, so that the coefficients for ρ^0 and $\rho^{2(v_{\max}-1)}$ vanish.

Observe that since $\rho_{\text{video}} \geq 0$, a sufficient condition for (17) is that all coefficients $\sum_{v+w=k} \zeta_{v,w} \leq 0$, $k = 1, \dots, 2v_{\max} - 1$. To this end, we will show that $\zeta_{v,w} + \zeta_{w,v} \leq 0$, where we take $v < w$ without loss of generality, i.e.

$$\begin{aligned} \zeta_{v,w} + \zeta_{w,v} \leq 0 &\iff \\ &\left(\frac{1}{w+1} - \frac{1}{v+1} \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!} \beta_{\text{video}}(s, v+1) \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!} \right) + \\ &+ \left(\frac{1}{v+1} - \frac{1}{w+1} \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!} \beta_{\text{video}}(s, w+1) \right) \left(\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!} \right) \leq 0 \end{aligned}$$

or, equivalently,

$$\frac{\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!} \beta_{\text{video}}(s, v+1)}{\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(v+1)} \frac{\rho_{\text{speech}}^s}{s!}} \geq \frac{\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!} \beta_{\text{video}}(s, w+1)}{\sum_{s=0}^{C-\beta_{\text{video}}^{\min}(w+1)} \frac{\rho_{\text{speech}}^s}{s!}},$$

i.e.,

$$\mathbf{E} \left\{ \beta_{\text{video}} \left(S_{C-\beta_{\text{video}}^{\min}(v+1)}, v+1 \right) \right\} \geq \mathbf{E} \left\{ \beta_{\text{video}} \left(S_{C-\beta_{\text{video}}^{\min}(w+1)}, w+1 \right) \right\},$$

where S_x is a random variable distributed as the queue length in a standard Erlang loss model with capacity x and traffic load ρ_{speech} . Observe that effectively we have reduced the inequality $\mathbf{R}_{\text{video}}^c \leq \mathbf{R}_{\text{video}}^t$ for the SV model to a set of inequalities for a speech-only model, i.e. for the standard Erlang loss model.

To complete the proof, we will show that $\beta_{\text{video}} \left(S_{C-\beta_{\text{video}}^{\min}(v+1)}, v+1 \right)$ is almost surely non-increasing in v , for $v = 0, \dots, v_{\max} - 1$. Substituting $y = C - \beta_{\text{video}}^{\min}(v+1)$ we have that

$$\beta_{\text{video}} \left(S_{C-\beta_{\text{video}}^{\min}(v+1)}, v+1 \right) = \beta_{\text{video}} \left(S_y, \frac{C-y}{\beta_{\text{video}}^{\min}} \right),$$

which we will demonstrate to be almost surely non-decreasing in y , by comparing the above expression for $y, y + \beta_{\text{video}}^{\min} \in [0, C - \beta_{\text{video}}^{\min}]$, where the lower (upper) bound corresponds with $v = v_{\max} - 1$ ($v = 0$). First

observe that the sample paths of the Erlang loss model with capacity y and $y + \beta_{\text{video}}^{\min}$ can readily be compared. Clearly, for an identical input of interarrival times and call lengths it must be that the sample path of the system with capacity $y + \beta_{\text{video}}^{\min}$ is never below that of the system with capacity y . In fact, starting with an empty system, the sample paths coincide until a call is blocked in the system with capacity y . Then, during the period that the system with capacity y is full, it may be that one or more additional calls are admitted to the system with capacity $y + \beta_{\text{video}}^{\min}$. Note that at most $\beta_{\text{video}}^{\min}$ additional calls can be accepted. The sojourn times of the additional calls are independent of the sojourn times of the other calls in the system with capacity $y + \beta_{\text{video}}^{\min}$, which are also present in the system with capacity y . Hence, with probability 1,

$$S_y \leq S_{y+\beta_{\text{video}}^{\min}} \leq S_y + \beta_{\text{video}}^{\min} \text{ and } S_y \leq y.$$

Combining these results with the fact that $y + \beta_{\text{video}}^{\min} \leq C$ and, in general, for $a, b \in \mathbb{R}$ it holds that if $a \geq b > \epsilon$ then $\left(\frac{a-\epsilon}{b-\epsilon}\right) \geq \frac{a}{b}$, implies that

$$\frac{C - S_{y+\beta_{\text{video}}^{\min}}}{C - (y + \beta_{\text{video}}^{\min})} \geq \frac{C - (S_y + \beta_{\text{video}}^{\min})}{C - (y + \beta_{\text{video}}^{\min})} \geq \frac{C - S_y}{C - y},$$

with probability 1. Recall that

$$\beta_{\text{video}} \left(S_y, \frac{C - y}{\beta_{\text{video}}^{\min}} \right) = \min \left\{ \beta_{\text{video}}^{\max}, \beta_{\text{video}}^{\min} \frac{C - S_y}{C - y} \right\},$$

so that

$$\beta_{\text{video}} \left(S_{y+\beta_{\text{video}}^{\min}}, \frac{C - (y + \beta_{\text{video}}^{\min})}{\beta_{\text{video}}^{\min}} \right) \geq \beta_{\text{video}} \left(S_y, \frac{C - y}{\beta_{\text{video}}^{\min}} \right),$$

with probability 1, which completes the proof. ■

Appendix 3: Proof of Corollary 3

Proof. The proof follows from manipulating the inequality proven in Theorem 2, using expressions (9) and (10), and relating it to the derivative of the time-average video throughput expression (10) with respect to

ρ_{video} :

$$\begin{aligned}
& \mathbf{R}_{\text{video}}^c \leq \mathbf{R}_{\text{video}}^t \\
\iff & \left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v-1) \beta_{\text{video}}(s, v) \right) \left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v) \right) + \\
& - \left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v) \beta_{\text{video}}(s, v) \right) \left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v-1) \right) \leq 0 \\
\iff & \frac{\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v-1) \beta_{\text{video}}(s, v)}{\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v)} + \\
& \frac{\left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v) \beta_{\text{video}}(s, v) \right) \left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v-1) \right)}{\left(\sum_{(s,v) \in \mathbb{S}_{\text{video}}^+} \pi(s, v) \right)^2} \leq 0 \\
\iff & \frac{\partial \mathbf{R}_{\text{video}}^t}{\partial \rho_{\text{video}}} \leq 0.
\end{aligned}$$

■