# THE THEORY OF NETWORKS OF SINGLE SERVER QUEUES AND THE TANDEM QUEUE MODEL

PIERRE LE GALL
*France Telecom, CNET*
*4 Parc de la Bérengère*
*F-92210 Saint-Cloud, FRANCE*

We consider the stochastic behavior of networks of single server queues when successive service times of a given customer are highly correlated. The study is conducted in two particular cases: 1) networks in heavy traffic, and 2) networks in which all successive service times have the same value (for a given customer), in order *to avoid the possibility of breaking up the busy periods*. We then show how the *local queueing delay* (for an **arbitrary** customer) can be derived through an *equivalent tandem queue* on the condition that one other local queueing delay is added: the **jitter delay** due to the independence of partial traffic streams.

We consider a practical application of the results by investigating the influence of long packets on the queueing delay of short packets in modern packet switched telecommunication networks. We compare these results with the results given by traffic simulation methods to conclude that there is good agreement between results of calculation and of traffic simulation.

**Key words:** Queueing Networks, Tandem Queues, Local Queueing Delay, Jitter Delay.

**AMS subject classifications:** 60J, 60K, 94C.

## 1. Introduction

We consider a network of **single server** queues where (at each stage) the customers only enter the downstream queue when they have fully completed service. The service discipline at all queues is "*first come-first served*". It has become rather usual to characterize the stochastic behavior of such a network by means of the so-called "*product form*" solution, whether the network be open or closed.

There is, however, a fundamental difference in the notion of traffic source. In this paper, we evaluate the case where successive service times (at the various network stages) of a given customer are highly correlated. The above mentioned theories no longer apply, although we will show how the local queueing delay (for an **arbitrary** customer) can be derived from an **equivalent tandem queue**, thereby eliminating (and integrating) the impact of intermediate arrivals and "premature departures".

However, this equivalence is only possible in the following two important cases: 1) *heavily loaded networks avoiding to break up the busy periods*, and 2) where the *successive* service times (for a *given* customer) are *identical*. It will then be necessary to add (for an **arbitrary** traffic stream) another supplementary local queueing delay: a **jitter delay** due to the mutual independence of the partial traffic streams constituting the local traffic stream.

This supplementary delay appears when we consider the new *local* arrival order of customers, which is different to the arrival order at the network input. This type of problem arises since we define traffic sources at the network input in opposition to the usual concept of local traffic source. The frequently adopted principle of considering a departure from a previous stage as a traffic source offered to the following stage leads to ignoring the jitter effect.

We conduct the study in three parts. After introducing assumptions, notation, and some preliminary theory in Section 2, the study (to evaluate the overall queueing delay) neglects the case of "premature departures" in Section 3. In other words, we study the case of a "*concentration tree*". The distribution of the **local** queueing delay for an **arbitrary** customer (including the jitter delay) appears as that of the difference between two *overall queueing delays* in equivalent tandem queues. In Section 4, "*premature departures*" are taken into account, along with the jitter effect. Finally, in Section 5, we apply our results to the case of two distinct populations of packets (of different lengths) handled in a modern packet switched telecommunication network. We make a number of observations by traffic simulation to check our theoretical results, which appear to be in good agreement with traffic simulations. In Sections 3 and 4, the processes of *non-simultaneous* arrivals are governed by some *general probability distribution* in case of a *stationary regime*. Section 5 is restricted to the case of Poisson arrivals.

## 2. Preliminary Theory, Assumptions and Notation

### 2.1 The Equivalent Tandem Queue

We assume that there are no "premature departures" and no local exogenous arrivals. To understand the basic phenomenon, we consider the simple case of two traffic streams each carried by a single server at the first stage before merging at a single second stage server. (See Figure 1.)

In Le Gall [6], Section II.2 we have already considered this case assuming successive service times are identical for a given customer. Consider the overall process (carried by the two first stage servers) and two arrivals (at the network input) $X_{n_1}$ and $X_{n_2}$, which successively find the two first stage servers idle. We showed that the number, size, and length of the busy periods of the second stage server during the interval $(X_{n_1}, X_{n_2})$ would be the same (subtracting the possibility of some jitter of the busy periods) if the two first stage servers were replaced by one server. We deduced that *the queueing delay distribution* (subtracting any jitter delay) *at the second stage would be unchanged for an* **arbitrary** *customer* (independently of the considered partial traffic stream). The result was extended by successive extrapolations to the case of a concentration tree. Therefore, to evaluate the local queueing delay at the final stage of the concentration tree, we may replace the tree by an "*equivalent tan-*

*dem queue*" carrying the same traffic streams with the same service times, provided we add a *jitter delay* generated by the mutual independence of the various branches of the concentration tree. This jitter delay takes into account the impact of the variations in the local arrival order, since the equivalent tandem queue corresponds to a local arrival order identical to the overall arrival order at the network input. In the case of branches with unequal lengths, we will define later the *parameter* $m_0$ (see equation (17)) to represent the number of stages of the equivalent tandem queue.

In general, this equivalence property is not true when successive service times (of the same customer) are different and varying, or are mutually independent. However, in the case of heavily loaded networks, the property does remain true *if it may avoid breaking up the busy periods*, since there is just one busy period of the second stage server during the interval $(\mathbf{X_{n_1}}, \mathbf{X_{n_2}})$. The size and the length of this busy period is therefore unchanged due to the extended delays at the second stage, which tends to cause busy periods to accumulate. Finally, in the second case of heavily loaded networks, it appears the local queueing delay at the final stage can be also evaluated by introducing the above concept of an "*equivalent tandem queue*". We formulate the following hypothesis:

**Hypothesis 1 (Busy Periods not Broken Up):** *We make the following hypothesis concerning the only two possibly different cases which we consider in this paper.*

a)    ***For arbitrary traffic intensity***, *we assume all successive service times are* ***identical*** *for a* ***given*** *customer, or*

b)    ***For heavy load*** (*at successive stages*), *successive service times are arbitrary, provided it may avoid breaking up the busy periods* (*see condition* (3) *below*).

*In particular, the second condition is satisfied when successive service times are not decreasing, introducing a* ***strong correlation*** *between these service times.*
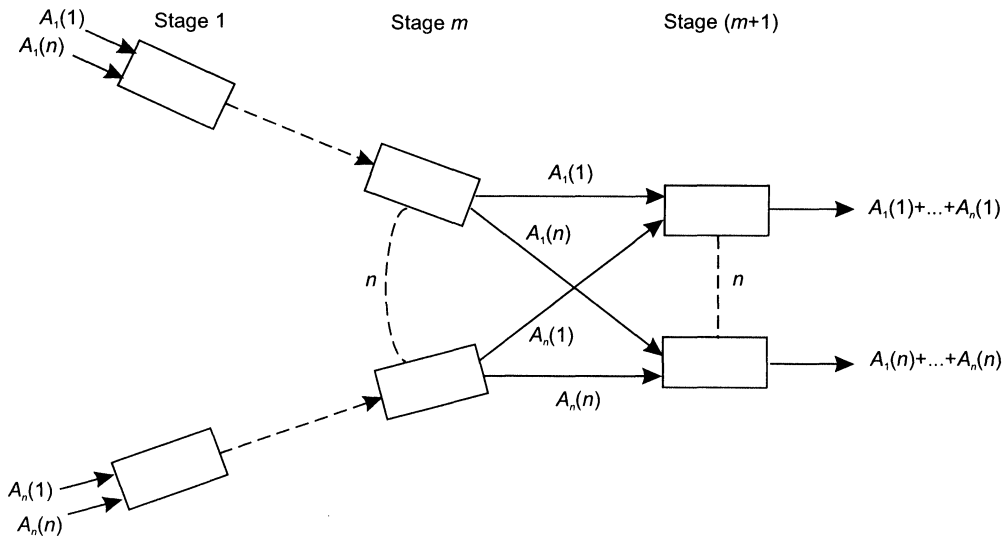


**Figure 1.** The full network

## 2.2  First Example of a Symmetrical Network

In Figure 1, we consider a first example of symmetrical network: $\mathbf{n}$ identical single server $(m+1)$-stage tandem queues, carrying statistically identical traffic streams which are interconnected at the last queueing stage $(m+1)$. In other words, the last stage consists of $\mathbf{n}$ queues, each of which is *dedicated* to a particular type of work.

Customers of a given tandem queue are forwarded to the last stage server dedicated to the type of work required. Naturally, this server can be accessed by all $\mathbf{n}$ tandem queues for the type of work in question, as illustrated in Figure 1.

From a theoretical point of view, this example is very interesting to evaluate the local queueing delay at the final stage. When $\mathbf{n}$ is large enough so that each final server receives only a small fraction of the traffic from each tandem queue, the *local* arrival process at the considered last stage dedicated server becomes practically Poisson. It is thus tempting to consider that this server constitutes a queueing system of the M/G/1 type. In fact, the upstream part of the network generates variations in the local arrival order. There exists a strong influence in the local queueing process giving rise to two kinds of phenomenon: a **G/G/1** (and not a M/G/1) **queue**, and **a jitter effect** (globally for busy periods).

## 2.3  Terminology

In the following, we avoid using the term *"waiting time"*, which has acquired a number of different meanings in Operations Research. At each local server, we identify the following three *local* times:
  1)  the **local** *queueing delay* due uniquely to the queueing phenomenon;
  2)  the *service time* $\mathbf{T_n}$ corresponding to the working time of the server $[\mathbf{T} = E(\mathbf{T_n})]$; and
  3)  the **local** *sojourn time* $\mathbf{S}$ corresponding to the sum of the first two times, since we assume the customer leaves the server at the end of its service to join the queue at the next server.

In opposition to the term *"local delay"*, we also use the term *"overall delay"* meaning the sum of delays over a number of successive stages. We also use the usual concepts of:
  1)  *"arrival rate"* $\boldsymbol{\lambda}$,
  2)  *"load"* (i.e., traffic intensity) $\rho = \lambda \mathrm{T}$,
and consider the *length* $\mathbf{L}$ of the local queue at an arbitrary arrival instant at the considered server. To verify the validity of the method of evaluating local queueing delays by traffic simulation, we will make use of Little's formula for mean values in the stationary regime:

$$\bar{\mathrm{L}} = \lambda \bar{\mathrm{S}} . \tag{1}$$

This illustrates the need to clearly define terminology since the formula is usually written:

$$\bar{\mathrm{L}} = \lambda \bar{\mathrm{W}},$$

where $\bar{\mathrm{W}}$ denotes the mean *"waiting time"*. This results in confusion between the mean queueing delay and the mean sojourn time $\bar{\mathrm{S}}$.

## 2.4 Notation and Assumptions

Recall that the queueing discipline (in each queue) is "*first come-first served*". We assume the system is in the *stationary regime*.

For each of the **n** *identical and independent tandem queues*, at stage $k$ ( $= 1, \ldots, m$) and for the $j^{th}$ customer at the considered queue, we set:

1)    *local queue delay*: $w_j^k$;

2)    *service time*: $\mathbf{T}_j^k$;

3)    *sojourn time*: $s_j^k = w_j^k + \mathbf{T}_j^k$;

4)    *arrival epoch at stage* k: $\mathbf{X}_j^k$;

5)    *interarrival interval* between customers $(j-1)$ and $j$: $\mathbf{Y}_{j-1}^k = \mathbf{X}_j^k - \mathbf{X}_{j-1}^k$;

6)    $\mathrm{E} \exp\left(-z\mathbf{Y}_{j-1}^1\right) = \varphi_0^1(z) = \int_0^\infty e^{-zt} dF_0^1(t)$;

7)    $\mathrm{E} \exp(-z\mathrm{T}_j) = \varphi_1(z) = \int_0^\infty e^{-zt} dF_1(t)$;   $\mathrm{R}(z) \leq 0$.

$F_0^1(t)$ and $F_1(t)$ are therefore the distribution functions of the interarrival time (at the tandem queue input) and the service time $\mathbf{T}_j$, respectively. Recall that $\boldsymbol{\lambda}$ is the *arrival rate* at each (identical) tandem queue. This process of non-simultaneous arrivals (at the entry to the network) is governed by some *general* probability distribution in case of a *stationary regime*.

We assume *traffic is distributed uniformly over the* **n** *single server queues of the final* $(m+1)^{th}$ *stage*, so *the arrival rate at each server is still* $\lambda$. The *overall delays* from stage **i** to stage **k** ( $= 2, \ldots, m+1$) for the $j^{th}$ customer may be written:

$$W_j(i;k) = \sum_{h=i}^k w_j^h \text{ and } S_j(i;k) = \sum_{h=i}^k s_j^h.$$

Lastly, *for a **given** server of the final* $(m+1)^{th}$ *stage*, we consider the following interarrival intervals between arrivals $(j'-1)$ and $j'$:

1)    interarrival interval at the input to all **n** tandem queues for the customers offered to the considered final server: $\mathbf{Y}_{j'-1}$;

2)    interarrival interval at the considered final server: $\mathbf{Y}'_{j'-1}$; and

3)    $\mathrm{E} \exp(-z\mathbf{Y}_{j'-1}) = \varphi_0(z) = \int_0^\infty e^{-zt} dF_0(t)$.

Note the difference between the arrival process $\mathbf{Y}_{j-1}^1$ for a given tandem queue and the process $\mathbf{Y}_{i'-1}$ relating to all tandem queues, but only for customers destined for the considered final stage server. The same distinction applies between $\varphi_0^1(z)$ and $\varphi_0(z)$. When process $\mathbf{Y}_{j'-1}$ may be assumed to be a *renewal process*, process $\mathbf{Y}'_{j'-1}$ is more general, due to the jitter effect (see Section 3.1), as well as the upstream interference delay (see Sections 4.2 and 4.3). But we recall that, *usually,* $\mathbf{Y}_{j'-1}$ *is governed by some general probability distribution*.

## 2.5  Symmetry and Equivalence

Due to the *hypothesis* 1 leading to the possibility of an "*equivalent tandem queue*", and if we neglect the jitter effect provisionally, the interarrival interval at the considered final server may be written:

$$\mathbf{Y}'_{j'-1} = \mathrm{T}^{m+1}_{j'} + e^{m}_{j'}, \tag{2}$$

where $e^{m}_{j'}$, is the *occasional idle period* (at stage **m** of the equivalent tandem queue) between arrivals $(j'-1)$ and $j'$. Since we assume that we avoid breaking up the busy periods, we have to satisfy the condition

$$\boxed{\mathbf{T}^{k-1}_{j'} \leq s^{k}_{j'-1}}, \quad k = 2, 3, \ldots, m+1, \tag{3}$$

from Theorem (A.1) in Le Gall [6], leading to the equalities

$$e^{m}_{j'} = e^{1}_{j'}. \tag{4}$$

This means that the customer initiating a busy period at stage 1 (of the equivalent tandem queue) initiates an analogous busy period at each following stage. Consequently, we may use the equation of the G/G/1 server at stage $(m+1)$:

$$w^{m+1}_{j'} = \mathrm{Max}[0, s^{m+1}_{j'-1} - (\mathrm{T}^{m}_{j'} + e^{1}_{j'})]. \tag{5}$$

During a busy period of the final server, we have $e^{1}_{j'} = 0$. Equation (5) becomes:

$$w^{m+1}_{j'} = \mathrm{Max}[0, s^{m+1}_{j'-1} - \mathrm{T}^{m}_{j'}]. \tag{6}$$

The influence of the arrival process disappears during the busy period. This explains the possibility of introducing the symmetrical network of Section 2.2, even if the traffic streams are not identical. Moreover, condition (3) may be satisfied even in the case of mutually independent, successive service times *for heavily loaded networks*, when the extended delays at the final stage tend to cause busy periods to amalgamate. Finally, *this property of equivalence may be satisfied in a general way, for the same bit rate at each stage.*

### 2.6  Preliminary Results for Identical Successive Service Times

We recall that to analyze the queue at a given final stage server in the stationary regime, we need to define an equivalent tandem queue to which process $\mathbf{Y}_{j'-1}$ is offered. We must also consider the queue at the first stage of this tandem queue defined by the couple $[\mathbf{Y}_{j'-1}, \mathbf{T}_{j}]$ for which the characteristic function of the delay is $\Phi_1(z)$, and the probability of no delay is $\mathbf{Q}_1$, the service time $\mathbf{T}_{j}$ characteristic function being $\varphi_1(z)$. In the case of *identical successive service times*, we give the main results already presented in Le Gall [8].

It will be necessary to consider the queue $[\mathbf{Y}_{j'-1}, \mathbf{T}_{j}; T_j < t]$ corresponding to the service time characteristic function:

$$\varphi(z; t) = \int_0^t e^{-zu} dF_1(u), \tag{7}$$

when the service time $\mathbf{T}_{j}$ is less than $t$. The *probability of no delay* in this case will be written $\mathbf{Q}(t)$, and we have:

$$Q(\infty) = \mathbf{Q}_1. \tag{8}$$

After this first stage, let $S_m(t)$ denote the distribution function of the *overall sojourn time over the m other stages* of this equivalent tandem queue. From Formula 4.21 in Le Gall [8], we have, in the stationary regime and for a renewal arrival process:

$$S_m(t) = F_1(t/m)R_m(t),$$ (9)

with

$$R_m(t) = \frac{1}{2\pi i} \int_{+0} \varphi_0(-u)\left[\frac{\Phi_1(u)}{Q(t/m)}\right] \text{Exp}\left(-u\int_t^\infty \frac{1 - F_1(v/m)}{Q(v/m)}dv\right)\frac{du}{u},$$

where the first integral is a (Cauchy) contour integral with a contour consisting of a straight line just to the right of the imaginary axis in the complex plane (**u**), running from bottom to top and being closed at infinity on the right-hand side.

If $Y_{j'-1}$ is a **Poisson** process, Formula 5.2 in Le Gall [8] gives {see also Boxma [1] for $m = 1$}:

$$R_m(t) = \frac{1-\rho}{Q(t/m)} \text{Exp}\left(-\lambda\int_t^\infty \frac{1 - F_1(v/m)}{Q(v/m)}dv\right),$$ (10)

with

$$1 - \rho = 1 - \lambda\bar{T} = Q(\infty); \quad Q(t) \cong 1 - \lambda\int_0^t u dF_1(u).$$

If **t** is large enough, and if $T_j$ has a concentrated finite support, expression B.10 in Le Gall [6] for $m = 1$ allows us to approximate this expression by:

$$R_1(t) \cong \frac{1-\rho}{1 - \rho(t/T)}.$$ (11)

If we now subtract the overall service time from the overall sojourn time, we deduce from (9) the distribution function of the *overall queueing delay* relative to the *m* other stages of the equivalent tandem queue:

$$U(t/m) = \int_0^\infty R_m(t+\theta)dF_1(\theta/m).$$ (12)

In particular, approximation (11) leads to the following approximation expression for $m = 1$:

$$U_1(t) = \int_0^\infty \frac{1-\rho}{1 - \rho\frac{t+\theta}{T}}dF_1(\theta).$$ (13)

### 2.7 Extension for Arbitrary Successive Service Times

In the case of arbitrary successive service times (for a *given* customer) satisfying condition (3) of hypothesis 1, Theorem A.1 in Le Gall [6] leads to the following recurrence relation for the $(m+1)$-stage equivalent tandem queue:

$$(T_{j'}^1 + w_{j'}^2) + \ldots + (T_{j'}^m + w_{j'}^{m+1})$$

$$= \text{Max}[T^1_{j'} + \ldots + T^m_{j'}, S_{j'-1}(2; m+1) - e^1_{j'}], \qquad (14)$$

which may be rewritten:

$$(T^2_{j'} + w^2_{j'}) + \ldots + (T^{m+1}_{j'} + w^{m+1}_{j'}) + (T^1_{j'} - T^{m+1}_{j'})$$

$$= \text{Max}[T^1_{j'} + \ldots + T^m_{j'}, S_{j'-1}(2; m+1) - e^1_{j'}],$$

Let us introduce a second hypothesis for this general case:

**Hypothesis 2 (Arbitrary Successive Service Times):** *If $\epsilon$ is an arbitrary small positive number, we suppose that the following condition*

$$\text{Lim}_{m \to \infty} \text{Prob}\left( \left| \frac{T^1_{j'} - T^{m+1}_{j'}}{T^2_{j'} + \ldots + T^{m+1}_{j'}} \right| < \epsilon \right) \to 1, \qquad (15)$$

*concerning the service times of the arrival process* $\mathbf{Y}_{j'-1}$, *is satisfied for any* $j' > 0$.

In this case, the recurrence relation in equation (14) is equivalent (in probability for $m$ large enough) to the following relation:

$$\mathbf{S}_{j'}(2; m+1) = \text{Max}[\mathbf{T'}_{j'}(m), \mathbf{S}_{j'-1}(2; m+1) - e^1_{j'}], \qquad (16)$$

with

$$\mathbf{T'}_{j'}(m) = T^1_{j'} + \ldots + T^m_{j'}.$$

From a property presented in Le Gall [7], Formula 2, this tandem queue may be assimilated "in distribution", and for $m$ large enough, to a tandem queue with $(m_0 + 1)$ successive **identical** single servers, where $m_0$ is defined by the relation:

$$\text{Var}(\mathbf{m_0} T^{m+1}_{j'}) = \text{Var} \mathbf{T'}_{j'}(\mathbf{m}), \qquad (17)$$

if the local service time $T^{m+1}_{j'}$ and $\mathbf{T'}_{j'}(\mathbf{m})$ are not constant. In the case of mutual independence between successive service times (for the same customer), we have:

$$\mathbf{m_0} = \text{Var} \mathbf{T'}_{j'}(\mathbf{m}) / \text{Var} T^{m+1}_{j'}. \qquad (18)$$

In the case of mutual dependence for service times highly varying, we have usually:

$$(ET^{m+1}_{j'})^2 < < E(T^{m+1}_{j'})^2 \text{ and } [E\mathbf{T'}_{j'}(\mathbf{m})]^2 < < E(\mathbf{T'}_{j'}(\mathbf{m}))^2.$$

In that case, we have practically: $\text{Var } T \cong ET^2$, and consequently:

$$\mathbf{m_0^2} = E(\mathbf{T'}_{j'}(\mathbf{m}))^2 / E(T^{m+1}_{j'})^2. \qquad (19)$$

Finally, the *hypotheses* 1 and 2 allow us to reduce the study, for $m$ large enough, to the case of $(m_0 + 1)$ successive identical single servers, *without any intermediate arrivals*, that is to evaluate the local queueing delay for an *arbitrary* customer at the final stage. Note an important property depending on these hypotheses: the busy

periods of successive stages are corresponding since we avoid breaking them up. During these periods, $e_{j'}^1 = 0$, and (16) gives:

$$S_{j'}(2; m+1) = S_{j'-1}(2; m+1) = S_{j_0'}(2; m+1), \tag{20}$$

where $j_0'$ is the number of the customer initiating the busy period. Consequently, the overall upstream sojourn time is stable during the busy period. If we except the phenomenon of the jitter delay for the busy periods, the arrival process at the entry to the network (in the symmetrical case) is translated and reproduced at the final stage, explaining a local queueing phenomenon observation similar to a G/G/1 (and not M/G/1) queue with some jitter.

## 3. The Overall Queueing Delay

Since we abandon the usual assumption of a local traffic source in favor of the real traffic source at the network input, it is necessary to define the *local* queueing delay at stage $(m+1)$ as the difference between two *overall* queueing delays: 1) the queueing delay $W_{j'}(1; m+1)$ from stage 1 to stage $(m+1)$, to be evaluated in this section; and 2) another queueing delay from stage 1 to stage $m$, to be studied in the following section. We note that these delays should include the influence of upstream delays due to the mutual interference between partial traffic streams offered to the considered final stage server, only , corresponding to the arrival process $Y_{j'-1}$. The impact of "premature departures" (not offered to the final stage) appears in the upstream busy period lengths *which cannot be split*.

For instance, in Le Gall [4], Sections 6 and 7 for identical traffic streams, the interference (upon the mean downstream delay) of an upstream M/D/1 server is $T(\rho/n)/(1-\rho)$ and not $T(\rho/n)/[1-(\rho/n)]$. It is a phenomenon quite different of the lost call model. Consequently, the upstream load of the equivalent tandem queue should be the same as with "premature departures" in the original network. This is the case for the symmetrical network in Figure 1, which may be considered as a reference network for the equivalent tandem queue even if the original network is not symmetrical [due to the fact we consider an *arbitrary* customer at stage $(m+1)$, only]. This has already been seen in Section 2.5.

### 3.1 The Local Jitter Delay

As we already mentioned in Section 2.1, the equivalent tandem queue corresponds to a local arrival order [at stage $(m+1)$] identical to the order at the network input: process $Y_{j'-1}$. To take account of some variations in the local order (process $Y'_{j'-1}$) not attributable to "premature departures", we noted the need to evaluate some jitter effect globally for the local busy period. This jitter effect has already been considered in Keilson [3], p. 150, §8.6. This effect is only significant in the case of heavy traffic when the server busy period is long, due to our hypotheses. This effect only depends on the service time $T_{j_0'}^{m+1}$ initiating the busy period, and not on the following service times (see equation (6)). As this quantity is considered only at the successive starting epochs of these busy periods, it can be considered to be constant during the busy periods and equal to the mean value $T = E(T_{j'}^{m+1})$.

This additional **jitter delay** $J$ comes from the mutual independence of the partial traffic streams offered to the considered final stage server, and is generated by the load $\rho''$ of the considered server, excluding the load related to the incoming tandem queue corresponding to the customer which initiates the busy period. In the symmetrical case of Figure 1, we have:

$$\rho'' = \left(1 - \frac{1}{n}\right)\rho,$$

with

$$\rho = \lambda T.$$

Furthermore, following Fisher and Stanford [2], the service times generated by the load $\rho''$ create an apparent service time $T''_j$ for the considered customer equal to the busy period generated by $\rho''$. As for $T''^{m+1}_{j0}$, we have to consider only the mean value $T''$, since $T''_j$ may be considered to be constant during the busy periods. Finally, the jitter delay $J$ can be considered as the local queueing delay of a tandem queue offered load $\rho''$ with successive service times practically equal to:

$$\boldsymbol{T''} = \frac{T}{1 - \rho''},$$

with

$$\rho'' = \left(1 - \frac{1}{n}\right)\rho. \tag{21}$$

This is due to the Poisson approximation for the local arrival process, when $\boldsymbol{n}$ is large enough (traffic simulation gives $\boldsymbol{n} > 5$ approximately). The fact that service times may be considered as practically constant and equal to $\boldsymbol{T''}$ means that the value of $\boldsymbol{m}$ has no impact. Consequently, we have: $S_m(t) \equiv S_1(t)$.

For large $\boldsymbol{n}$, the mutual independence of the $\boldsymbol{n}$ original tandem queues tends to produce a Poisson local arrival process. We can therefore consider that the previous tandem queue is offered Poisson traffic at stage 1. Equation (11) then becomes:

$$R_1(t) = \frac{1 - \rho''}{1 - \rho'' \frac{t}{T''}}. \tag{22}$$

By applying (13), it should be noted that $F_1$ corresponds to the real service time $T$. This expression allows us to derive the following expression for the distribution function of the **local jitter delay** $J$ in the case of large $\boldsymbol{n}$ ( $> 5$):

$$J(t) = \frac{1 - \rho''}{1 - \rho'' \frac{t + T}{T''}}. \tag{23}$$

We have already used this simple expression in Le Gall [9].

### 3.2 The Overall Queuing Delay Without Premature Departures

Following our comments in Section 2.1, we derive the following theorem for *general input* in case of a *stationary regime*:

Theorem 3 (Overall Queuing Delay in the Concentration Tree): *For large $\boldsymbol{n}$ and using hypotheses 1 and 2 (if necessary), the **overall queueing delay** [from stage 1 to*

*stage* $(m_0 + 1)$], *denoted* $W_{j'}(1; m_0 + 1)$ *for the* $j^{th}$ *arbitrary customer with the arrival process* $Y_{j'-1}$ *offered to the concentration tree, is the sum of two **independent** delays:*

1) **the overall queueing delay** $V_{j'}(1; m_0 + 1)$ *relative to an* $(m_0 + 1)$-*stage equivalent tandem queue (with identical successive service times) offered arrival process* $Y_{j'-1}$, *the* $m_0$ *supplementary stages being defined by equation* (17) *in order to integrate the intermediate arrivals; the distribution function of the additional queueing delay due to these* $m_0$ *stages is given by equation* (12), *with equation* (9) *if* $Y_{j'-1}$ *is a renewal process, and with equation* (10) *if* $Y_{j'-1}$ *is Poisson; and*

2) *a **local jitter delay** $J$, the distribution function of which is given by equation* (23). *In other words:*

$$W_{j'}(1; m_0 + 1) = V_{j'}(1; m_0 + 1) + J. \tag{24}$$

**Notes:**

1) The network is no longer necessarily symmetrical, provided equation (21) of $\rho''$ is changed.

2) Without any "premature departure" and with no local exogenous arrivals, the *local* queueing delay at the final stage, in the stationary regime, and by deleting the suffix $j'$, is

$$w = W(1; m_0 + 1) - V(2; m_0 + 1). \tag{25}$$

This quantity may be observed directly.

3) The observation (by simulation) and the calculation of the *overall queueing delay* $W(1; m_0 + 1)$ may be extended to the case with "premature departures", just before the final stage, by offering them to this stage with zero service times. The access time at the considered final stage server is not changed, and for its evaluation, we may *avoid distinguishing the normal traffic* (offered at the final stage) *and "premature departures"*. But, in this way, we evaluate $W(1; m_0 + 1)$ **as perceived from the entry** to the network and not from the considered local stage.

## 4. The Local Queueing Delay with Premature Departures

Now, we take "premature departures" into account in the *stationary regime*. In equation (25) we cannot isolate the term $V(2; m_0 + 1)$, which now corresponds to a part of the traffic handled in the considered incoming tandem queue. When we consider (at the final stage) an *arbitrary* customer, he comes from a *given* incoming $m_i$-stage tandem queue $i$ $(i = 1, ..., n)$, corresponding to an overall queueing delay $V'(1; m_i)$ from stage 1 to stage $m_i$. If the traffic streams of this tandem queue $i$ are identical, we only have a probability $(1/n)$ that this customer is offered to the considered final stage server in Figure 1. When the traffic streams are not identical, we set for the *incoming tandem queue i*:

1) $a_{m_i}$: *total load* at this stage $m_i$;

2) $a'_{m_i}$: *part of the total load* $a_{m_i}$ *corresponding to the customers offered to*

*the considered final stage server*; and

3)     $a''_{m_i}$: *part of the total load $a_{m_i}$ corresponding to "premature departures"*;
$$a_{m_i} = a'_{m_i} + a''_{m_i}.$$

Based on the following argument, the preceding probability $(1/n)$ has to be replaced by $(a'_{m_i}/a_{m_i})$. Among the mean number of customers arriving (at stage $m_i$) during a mean service time ($= a_{m_i}$), the mean number of customers offered to the considered final stage server is $\boldsymbol{a'_{m_i}}$. The proportion $(a'_{m_i}/a_{m_i})$ gives the proportion of customers concerned inside an upstream busy period. This argument is valid for a *general input*. This rule has been presented and checked by traffic simulation successively in:

1)     Le Gall [4], Formula 16 and Section 10 for Poisson arrivals;
2)     Le Gall [5], Formula 5.14 and Section 7 for renewal arrivals; and
3)     Le Gall [9], Section 2.3 for identical successive service times with Poisson arrivals.

See also the simple case of constant and identical service times at each stage, as considered in the first paragraph of Section 3 above. If we except the influence of the jitter, the use of a probability equal to 1, instead of $(a'_{m_i}/a_{m_i})$, would lead to a local queueing delay equal to zero! This nonsense proves the existence of the reduction factor $(a'_{m_i}/a_{m_i})$.

Let $\boldsymbol{h_i}$ be the random variable $= 1$ with probability $(a'_{m_i}/a_{m_i})$, and $= 0$ with probability $[1 - (a'_{m_i}/a_{m_i})]$. For a non-symmetrical network, we deduce in the stationary regime:

$$w = W(1; m_0 + 1) - \boldsymbol{h_i} V'(1; m_i)$$

$$= \boldsymbol{h_i}[W(1; m_0 + 1) - V'(1; m_i)] + (\boldsymbol{1 - h_i}) W(1; m_0 + 1), \qquad (26)$$

$$\mathrm{E} \exp(zw) = \left(\frac{a'_{m_i}}{a_{m_i}}\right) \mathrm{E} \exp\{z[W(1; m_0 + 1) - V'(1; m_i)]\}$$

$$+ \left(\frac{a''_{m_i}}{a_{m_i}}\right) \mathrm{E} \exp(zW(1; m_0 + 1)).$$

As a consequence, equation (24) has to be revised in case of "premature departures". The equivalent tandem queue is not changed, where:

$$W(1; m_0 + 1) = V(1; 2) + [J + V(2; m_0 + 1)],$$

$$V'(1; m_i) = V'(1; 2) + V'(2; m_i), \qquad (27)$$

are independent of premature departures. It is a theoretical tandem queue giving a certain overall queueing delay which has to be considered as a whole, without distinguishing the components. But the *"additional queueing delay"* $[\boldsymbol{V(2; m_0 + 1) + J}]$, related to the $\boldsymbol{m_0}$ supplementary stages, *is perceived from the considered local stage as affected by the probabilities* $(a'_{m_i}/a_{m_i})$ for each branch $\boldsymbol{i}$. It is not perceived as an

observer at the entry to the network (see Note 3 of Theorem 3).

To evaluate this overall queueing delay for an *arbitrary* customer we do not have to distinguish these branches. Also, from equation (6), we already noted that the equivalent tandem queue may correspond to a symmetrical network, as illustrated in Figure 1. We set:

$$a = \sum_i a_{m_i}, \quad a' = \sum_i a'_{m_i}, \quad a'' = \sum_i a''_{m_i}. \tag{28}$$

The parameters of this hypothetical and symmetrical network, defining the equivalent tandem queue, become:

$$m' = \text{integer part of } m_0 \text{ [as defined by equation (17)]},$$
$$n = \text{integer part of } (a/a'). \tag{29}$$

Note that the loads in the equivalent tandem queue (premature departures excluded), and in the $n$ branches of the symmetrical network of Figure 1 (premature departures included), are identical. In the stationary regime, the delay $V(2; m_0 + 1)$ in expression (24) becomes $V(2; m' + 1)$, but

$$W_0 = V(1; 2) = V'(1; 2) \tag{30}$$

is not changed, as perceived locally (and globally for the arrival process $Y_{j'-1}$) without any impact of the upstream stage. Finally, relations (26) and preceding considerations lead to the following theorem:

**Theorem 4 (Overall Queueing Delay with Premature Departures):** *In the conditions of Theorem 3, but with premature departures in the stationary regime, the characteristic function of the **overall queueing delay** in the concentration tree for an **arbitrary** customer, **as perceived at the considered local stage** with the impact of premature departures at the upstream stage, is:*

$$\mathrm{E}\exp[zW(1; m' + 1]$$
$$= \left(\frac{a''}{a}\right)\mathrm{E}\exp(zW_0) + \left(\frac{a'}{a}\right)\mathrm{E}\exp[z(W_0 + J + V(2; m' + 1)], \tag{31}$$

*where $a, a', a'', m',$ and $W_0$ are defined by equations (28), (29) and (30), successively, with $W_0$ corresponding to a G/G/1 queue of an isolated server offered arrival process $Y_{j'-1}$.*

**Notes:**

1) In other words, *when the considered customer is delayed by a premature departure* [probability: $(a''/a)$], the overall upstream queueing delay (and the jitter) are not perceived from the local stage. It is considered *as a "fresh" arrival* not depending on the network. Finally, our concept of traffic source at the entry to the network leads again to the classical concept of local traffic source!

2) In the example of Note 3 (after Theorem 3), in which *"premature departures"* become *"normal"* customers (but with a zero service time at the final stage), the overall access time did not change for a number of *"normal"* customers $n$ ( $= a/a'$) times higher. In Le Gall [9], substitution (9), we recognized the need of a reduction factor leading to expression (26) above, but we did not apply it in a

term of expression (10) of this reference, evaluating the mean overall queueing delay.

Now, we do not distinguish the partial traffic streams, and we recall our hypothesis of *general input* in case of a **stationary regime**. Considerations from equations (26) and (31) lead to the following theorem, taking into account the same load in the equivalent tandem queue and in the branches of the hypothetical and symmetrical concentration tree:

**Theorem 5 (Local Queueing Delay for an Arbitrary Customer):** *In the conditions of Theorem 3, but with premature departures at the upstream stage, the characteristic function of the* **local queueing delay** *for an* **arbitrary** *customer is:*

$$\boxed{\operatorname{E}\exp\left(zw\right) = \left(\frac{a'}{a}\right)\operatorname{E}\exp\left(zD\right) + \left(\frac{a''}{a}\right)\operatorname{E}\exp\left(zW_0\right)},$$

with

$$D = J + [V(2;m'+1) - V(2;m')]. \tag{32}$$

$W_0$ *and* $V(2;m'+1)$ *are defined by Theorem 4, with* $V(2;m')$ *corresponding to the additional queueing delay of the* $m'$-*stage equivalent tandem queue* [*instead of* $(m'+1)$ *stage*]. $D$ *is the* **local** *queueing delay* (*jitter included*) *at the last stage* $(m'+1)$ *of the* **equivalent tandem queue**; $a, a'$ *and* $a''$ *are defined by equations* (28); *and* $m'$ *is defined by equations* (29).

**Notes:**

1) Another proof of Theorems 4 and 5 may be given by the example of Note 3 (after Theorem 3). In that case, we have to add the delay $(1 - h_i)[J + V(2;m'+1)]$ to $W(1;m'+1)$ [see expression (31)]; and to $V(2;m')$, simultaneously, without changing $w$ (Theorem 5) to be in the condition of this particular example.

2) Equation (32) is the *key formula* of this paper. The *local* queueing delay for an *arbitrary* customer may be deduced by interpolating the two extreme cases: 1) $D$ (equivalent tandem queue) when there are no "premature departures"; and 2) $W_0$ (isolated G/G/1 server) when there are plenty of "premature departures". In this latter case, we may find again the product form theory for highly meshed networks.

## 5. Case of Single Link Packet Switched Networks

We apply the preceding considerations to the case of single link packet switched networks with *Poisson arrivals*, as usually considered for modern telecommunication networks.

### 5.1 The Traffic Model

We consider the symmetrical network with $n$ branches, as illustrated in Figure 1. We define a traffic stream by deleting any suffix $i$, provisionally. The *arrival rate* (in each branch) is $\lambda$, and the bit rate is the same on each link, (i.e., the successive packet lengths ( = service times) are identical: $T_n^1 = T_n^2 = \ldots = T_n^m = T_n$). The distribution function is: $F(t) = \operatorname{Prob}(T_n < t)$.

The global traffic stream is the mixture of two partial traffic streams of category

$j$ $(j = 1, 2)$, corresponding to *packets of constant length* $T_j$ $(T_1 < T_2)$, with distribution function $F_j(t)$, $\lambda_j$ being the arrival rate. We let:

$$\rho_j = \lambda_j T_j, \lambda = \lambda_1 + \lambda_2, \rho = \rho_1 + \rho_2, \bar{T} = (\rho/\lambda);$$

$$dF_j(t) = (\lambda_j/\lambda) \text{ if } t = T_j, = 0 \text{ if } t \neq T_j. \tag{33}$$

From equation (10) we deduce (at the final stage) for $T_1 < t \le T_2$:

$$Q(t) = 1 - \rho_1, R_1(t) = \frac{1-\rho}{1-\rho_1} \text{Exp}[-c(T_2 - t)],$$

with

$$c = \frac{\lambda_2}{1 - \rho_1}. \tag{34}$$

From M/G/1 queue, we have:

$$\overline{W}_0 = \frac{1}{2} \frac{\rho}{1-\rho} \frac{E(T)^2}{E(T)},$$

with

$$E(T)^2 = \rho_1\left(\frac{T_1}{\lambda}\right) + \rho_2\left(\frac{T_2}{\lambda}\right);$$

$$\text{Var } W_0 = \frac{1}{3} \frac{\rho}{1-\rho} \frac{E(T)^3}{E(T)} + (\overline{W}_0)^2,$$

with

$$E(T)^3 = \rho_1\left(\frac{T_1^2}{\lambda}\right) + \rho_2\left(\frac{T_2^2}{\lambda}\right). \tag{35}$$

From equations (21) and (23), we deduce the approximate expression for the moments of the *jitter delay* distribution $(\alpha = 1, 2)$:

$$E(J^\alpha) = \alpha \int_{\bar{T}}^{T''} (t - \bar{T})^{\alpha - 1}\left[1 - \frac{1-\rho''}{1-\rho'' \frac{t}{T''}}\right] dt$$

with

$$\rho'' = \left(1 - \frac{1}{n}\right)\rho, T'' = \frac{\bar{T}}{1 - \rho''}. \tag{36}$$

From equation (12) we deduce the moments of the *additional queueing delay* of the $(m + 1)$-stage $(m' = m)$ equivalent tandem queue for an *arbitrary* customer $(\alpha = 1, 2)$:

$$E[V(2; m+1)]^\alpha = \frac{\lambda_1}{\lambda}[m(T_2 - T_1)]^\alpha\left[1 - \alpha\frac{1-\rho}{1-\rho_1} \frac{(K-1)^{\alpha - 1} + (-1)^\alpha \exp(-K)}{K^\alpha}\right]$$

with

$$K = mc(T_2 - T_1). \tag{37}$$

Note, for long packets, that the additional queueing delay is equal to zero. Let $w_L$ and $w_A$ be the local queueing delay (at the final stage) of *long* and *arbitrary* packets, respectively. Expressions (32), (35), (36), and (37) give for the *mean* values:

$$W_L = \left(1 - \tfrac{1}{n}\right)\overline{W_0} + (\tfrac{1}{n})\overline{J},$$

$$W_A = W_L + \left(\tfrac{1}{n}\right)\left[\overline{V(2;m+1)} - \overline{V(2;m)}\right]. \tag{38}$$

## 5.2 Traffic Simulations

Some simulations were carried out (see Romoeuf [10]) for the complete and symmetrical 2 and 4-stage network of Figure 1 ($m = 1, 3$), with $n = 10$, $T_1 = 1$, and $T_2 = 10$. Usually, it may be sufficient to simulate the truncated network shown in Figure 2. The number of packets ran per final stage server was approximately $10^7$, giving an excellent accuracy.
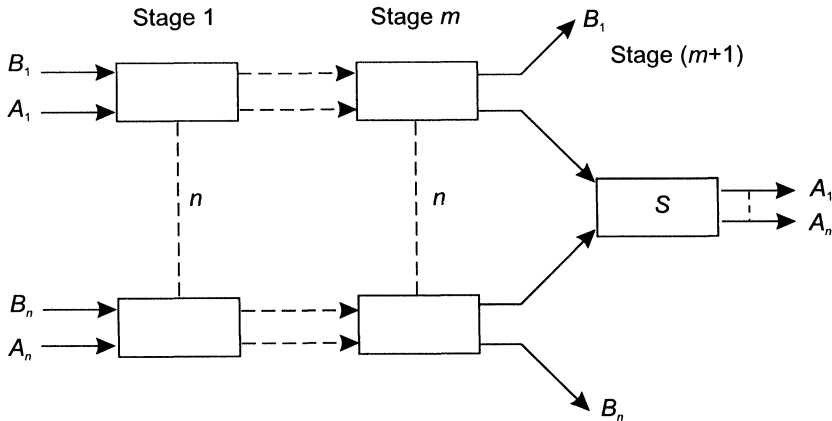


**Figure 2.** The truncated network

**Table 1** gives (**for $n = 10$**) comparative results for $W_L$ and $W_A$ and the difference $\Delta = W$ (simulated value) $- W$ (calculated value). Here, $\rho_1$ and $\rho_2$ are the loads (traffic intensities) in each server due to short and long packets, respectively. $\Delta$ appears equal to zero in the case $\rho_2 = 0.3\rho_1$ (low "long" load), and equal to 0.7 ($\rho = 0.9$) and 0.08 ($\rho = 0.6$) in the case $\rho_2 = \rho_1$ (high "long" load). It means some discrepancy for the mean jitter delay only, in case of $\rho_2 = \rho_1$: $\overline{J} = 11.4$ instead of 4.4 for $\rho = 0.9$, and $\overline{J} = 1.3$ instead of 0.5 for $\rho = 0.6$. In that case, the slight number of short packets per busy period (at the final stage) is not appropriate to the assumption of stable short service times during the busy period (see Section 3.1) generated by a long service time. But the discrepancy is not significant for $W_A$: 24.2 instead of 23.5 ($\rho = 0.9$), and 4.5 instead of 4.4 ($\rho = 0.6$). For $\rho_2 = 0.3\rho_1$ (low "long" load) there is no discrepancy for the evaluation of the jitter delay, and in every case, the mean additional queueing delay (due to the tandem queue model) is correctly evaluated.

| $\rho = 0.9$ | | | | | | | |
|---|---|---|---|---|---|---|---|
| | $W_0$ | $W_L$ | | | $W_A$ | | |
| $m = 1$ | | $C$ | $S$ | $\Delta$ | $C$ | $S$ | $\Delta$ |
| $\rho_2 = \rho_1$ | 24.7 | 22.7 | 23.4 | 0.7 | 23.4 | 24.1 | 0.7 |
| $\rho_2 = 0.3\rho_1$ | 14.0 | 12.9 | 12.9 | 0 | 13.5 | 13.5 | 0 |
| $m = 3$ | | | | | | | |
| $\rho_2 = \rho_1$ | 24.7 | 22.7 | 23.4 | 0.7 | 23.5 | 24.2 | 0.7 |
| $\rho_2 = 0.3\rho_1$ | 14.0 | 12.9 | 12.9 | 0 | 13.7 | 13.7 | 0 |

**Table 1a**

| $\rho = 0.6$ | | | | | | | |
|---|---|---|---|---|---|---|---|
| | $W_0$ | $W_L$ | | | $W_A$ | | |
| $m = 1$ | | $C$ | $S$ | $\Delta$ | $C$ | $S$ | $\Delta$ |
| $\rho_2 = \rho_1$ | 4.12 | 3.76 | 3.84 | 0.08 | 4.19 | 4.24 | 0.05 |
| $\rho_2 = 0.3\rho_1$ | 2.33 | 2.13 | 2.13 | 0 | 2.42 | 2.42 | 0 |
| $m = 3$ | | | | | | | |
| $\rho_2 = \rho_1$ | 4.12 | 3.76 | 3.84 | 0.08 | 4.40 | 4.47 | 0.07 |
| $\rho_2 = 0.3\rho_1$ | 2.33 | 2.13 | 2.13 | 0 | 2.64 | 2.60 | -0.04 |

**Table 1b**

**10** identical tandem queues of **m** successive queues converging
on one final stage server (among 10) [see Figure 1 or 2]

1) **In each upstream tandem queue:**
   traffic mix: 2 populations of packet lengths ( = service times).

| | short packets | long packets | Total |
|---|---|---|---|
| packet length | $T_1 = 1$ | $T_2 = 10$ | |
| arrival rate | $\lambda_1$ | $\lambda_2$ | $\lambda = \lambda_1 + \lambda_2$ |
| load* | $\rho_1 = \lambda_1 T_1$ | $\rho_2 = \lambda_2 T_2$ | $\rho = \rho_1 + \rho_2$ |

*( = traffic intensity)

2) **In the final stage queue**
   *a)* **Mean local queueing delay:**
   M/G/1 server for total load $\rho$:    $W_0$;
                long packet:    $W_L$;
           arbitrary packet:    $W_A$.

   *b)* **Columns:**
   1)   $C$: calculated value $W$ (formulae (38));
   2)   S: simulated value $W'$;
   3)   $\Delta = W' - W$.

   **(Symmetrical case)**   $n = 10$

For $m = 3$ and $\rho = 0.6$, **Table 2** gives comparative results for $W_A$ *when* $n$ *varies.* The number $n$ of tandem queues converging vary from $n = 1$ to $n = 10$. We may see again a good agreement between the results given by calculation and by traffic simulation.

| $n$ | $\rho_2 = \rho_1$ $W_A$ | | $\rho_2 = 0.25\rho_1$ $W_A$ | |
|:---:|:---:|:---:|:---:|:---:|
| | $C$ | $S$ | $C$ | $S$ |
| 1 | 6.4 | 6.4 | 4.8 | 4.8 |
| 2 | 5.3 | 5.6 | 3.4 | 3.3 |
| 3 | 4.9 | 5.2 | 3.0 | 2.8 |
| 5 | 4.6 | 4.7 | 2.7 | 2.6 |
| 7 | 4.5 | 4.6 | 2.5 | 2.4 |
| 10 | 4.4 | 4.5 | 2.4 | 2.2 |
| $W_0$ | 4.12 | | 2.10 | |

**Table 2:** $n$ identical tandem queues of $(m = )$ 3 successive queues
converging on one final stage server (among $n$).
[See Figure 1 or 2, and notation and data in Table 1]

1) **total load** ( = traffic intensity) in each server:  $\rho = \rho_1 + \rho_2 = \mathbf{0.6}$.
2) **mean local queueing delay** (at the **final** queue):
         for an **arbitrary** packet:     $W_A$;
         case of an M/G/1 server:   $W_0$.
3) **columns:**    $C$: calculated value (formulae (38));
          S:  simulated value.

         (**Symmetrical case**)  **Influence of** $n$.


# 6.  Conclusion

We have shown that the existence of an (approximately evaluated) *jitter effect* allows us to connect the theory of queueing networks to the tandem queue model in that case of single servers *avoiding the break up of busy periods*, which may be of some general application for heavily loaded networks. It appears there is a strong influence of *"premature departures"* at the upstream stage (in meshed networks) leading to a local queue, *approximately identical to that of an isolated G/G/1 server in case of many "premature departures"*, even if we put aside the concept of local traffic source! In the other cases, the local queue may be deduced *by interpolating* the preceding case with the case of no "premature departures" ( = *equivalent tandem queue*).


# References

[1]    Boxma, O.J., On a tandem queueing model with identical service times at both counters, *Adv. Appl. Probab.* **11** (1979), 616-659.
[2]    Fisher, W. and Stanford, D., Approximations for the per-class waiting time and interdeparture time in the $\sum_i \text{GI}_i/\text{GI}_i/1$ queue, *Perf. Eval.* **14**:2 (1992), 65-78.

[3]  Keilson, J., *Markov Chain Models-Rarity and Exponentiality*, Springer-Verlag, New York 1979.

[4]  Le Gall, P., Packetized traffic engineering for new services, *Proc. ITC Seminar*, Adelaide (Australia, Sept. 1989), *Computer Networks and ISDN Systems J.* **20** (1990), 425-433.

[5]  Le Gall, P., Networks of single server queues with first in-first out service discipline, *Proc. ITC-13* (Copenhagen, June 1991), *Queueing, Performance and Control in ATM* **15** (1991), 161-172.

[6]  Le Gall, P., Traffic modeling in packet switched networks for single links, *Annales des Telecom.* **49**:3-4 (1994), 111-126.

[7]  Le Gall, P., Bursty traffic in packet switched networks, *Proc. ITC-14* (Antibes, France, June 1994), *The Fund. Role of Teletraffic in the Evolution of Telecom. Networks* **1a** (1994), 535-549.

[8]  Le Gall, P., The overall sojourn time in tandem queues with identical successive service times and renewal input, *Stoch. Proc. and their Applications* **52** (1994), 165-178.

[9]  Le Gall, P., Traffic simulation and traffic modeling in the control plane, *Proc. 10$^{th}$ ITC-Specialist Seminar* (Lund, Sweden, Sept. 1996).

[10] Romoeuf, L., Traffic simulations in the control plane, Note *CNET/DT/LAA/ EIA/EVP/95*-08LR, 22 May 1995.