

Copyright © 1996, by the author(s).
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

**NETWORKS OF QUEUES WITH LONG-RANGE
DEPENDENT TRAFFIC STREAMS**

by

Venkat Anantharam

Memorandum No. UCB/ERL M96/24

15 April 1996

COVER PAGE

**NETWORKS OF QUEUES WITH LONG-RANGE
DEPENDENT TRAFFIC STREAMS**

by

Venkat Anantharam

Memorandum No. UCB/ERL M96/24

15 April 1996

ELECTRONICS RESEARCH LABORATORY

College of Engineering
University of California, Berkeley
94720

Networks of Queues with Long-Range Dependent Traffic Streams

Venkat Anantharam *
EECS DEPARTMENT
UNIVERSITY OF CALIFORNIA
BERKELEY, CA 94720
`ananth@vyasa.eecs.berkeley.edu`

Abstract

Traditional applications of queuing theory in the analysis and design of communication networks often use the fact that quasi-reversible queues can be interconnected with Bernoulli routing to form network models whose stationary distribution is of product form. To construct such models it is necessary to assume that the external arrival process is a Poisson process (in continuous time) or an i.i.d. sequence of Poisson random variables (in discrete time). Recently, however, several empirical studies have provided strong evidence that the traffic to be carried by the next generation of communication networks exhibits long-range dependence, implying that it cannot be satisfactorily modeled as a Poisson process.

Our purpose in this paper is to exploit the features that make networks of quasi-reversible queues product form so as to construct a self-contained class of queuing network models whose external arrival processes can be long-range dependent. Throughout the paper, we restrict ourselves to discrete time. The long-range dependent arrival process models we consider can be described as follows : Sessions arrive according to a sequence of i.i.d. Poisson random variables. Each session is active for an independent duration that is positive integer valued with a regularly varying distribution having finite mean and infinite second moment. While it is active, a session brings in work at rate 1. The networks we consider are comprised of monotone quasi-reversible S-queues with Bernoulli routing. The queues handle their arrivals at the session level, i.e., once a session enters service at a queue it continues to receive service at rate 1 until all the work it brings in is completed.

Our main observation is that all the internal traffic processes of such networks are long-range dependent. This result provides a family of interesting new examples of long-range dependent processes, in addition to its potential use in studying performance issues for communication networks,

*Research supported in part by NSF grant NCR 94-22513 and by the AT & T Foundation

1 Introduction

In the past few years the telecommunication network arena has witnessed the evolution of networks that can offer a wide variety of services, such as audio, video, data, etc. using a common protocol suite. Such networks are called *integrated services networks*. Since these networks are largely based on transmission over a fibre-optic medium and thus have a bandwidth that is orders of magnitude higher than that of earlier generations of networks, they are also generally called *broad-band networks*. Several recent empirical studies of traffic to be carried by such broad-band networks have established the importance of studying queuing behaviour with traffic models having long-range dependence. For instance, Leland et al. [8] have demonstrated the self-similar nature of Ethernet traffic by a statistical analysis of Ethernet traffic measurements at BellCore; Beran et al. [3] have demonstrated long-range dependence in samples of variable bit rate video traffic generated by a number of different codecs; and Paxson and Floyd [12] have concluded the presence of long-range dependence in TELNET and other wide area network traffic.

Motivated by this empirical evidence a number of studies have recently appeared that investigate single-server queues with long-range dependent arrival models. These include the works of Norros [10] (discussed further by Duffield and O'Connell [5]) who uses a fractional Brownian arrival model; Resnick and Samorodnitsky [13] who use an arrival model derived from a stable integral; Likhanov et al. [9] who use an asymptotically self-similar discrete time arrival model (which also appears in the survey of Cox [4, page 68]); Anantharam [2] who studies extensions of the model of [9]; and Parulekar and Makowski [11] who also study the model of [9].

In the analysis and design of telecommunication networks, the concept of a *quasi-reversible queue*, see Kelly [7], Walrand [14], Walrand [15], and Section 3 of this paper, has been a key in enabling the effective use of queuing models. Networks constructed from quasi-reversible queues with Bernoulli routing and an external arrival process which is a Poisson process (in continuous time) or a sequence of i.i.d. Poisson random variables (in discrete time) offer considerable modeling flexibility and admit *product form* stationary distributions, which makes the computation of stationary performance quantities easy.

Our purpose in this paper is to construct queuing network models with tractable analytical features that allow the modeling of long-range dependent traffic processes. To this end, we bring together the discrete time quasi-reversible S-queue network models of Walrand [15] and the extension of the model of Likhanov et al. [9] studied in [2]. Our networks have external arrival processes which are long-range dependent processes described as follows : Sessions arrive according to a sequence of i.i.d. Poisson random variables. Each session lasts for an independent duration having the distribution of a random variable τ which is positive integer valued with finite mean and infinite variance and having regularly varying tail, i. e. $P(\tau > t) \sim t^{-\alpha} L(t)$ where $1 < \alpha < 2$ and $L(\cdot)$ is a slowly varying function (see appendix). While it is active, a session brings in work at rate 1. Our networks consist of monotone discrete time quasi-reversible S-queues with Bernoulli routing. The queues handle the work *at the session level*, i.e., once a session enters service at a queue it continues to receive service at rate 1 until all the work it brings in is completed. Under the usual stability conditions,

we prove that all the internal traffic processes of such networks are long-range dependent.

The precise definition of the class of external arrival models that we consider is in Section 2. The basic facts about the discrete time quasi-reversible network models constructed from S-queues in [15] are recalled in Section 3, together with the concept of monotonicity for such queues, which was introduced in [1]. The class of network models that we consider is described in Section 4. The main observations we make about such network models are proved in Section 5. Finally, for convenience, elementary definitions and facts about regularly varying functions are collected in an appendix.

We close the introduction by recalling some basic notation and terminology. As usual \mathbf{N} , \mathbf{Z} , and \mathbf{R} denote the set of natural numbers, the set of integers, and the set of real numbers respectively. A discrete time stochastic process will be called *wide sense stationary* if it has constant mean and stationary covariance structure, with finite second moment. For a wide sense stationary stochastic process $(X_l, l \in \mathbf{Z})$, let $r(k) = \text{Cov}(X_l, X_{l+k})$ denote its autocovariance function. The process will be called *long-range dependent* if $\sum_k |r(k)| = \infty$. The process is called *short-range dependent* if it is not long-range dependent, i.e. if $\sum_k |r(k)| < \infty$. For more details about long-range dependence, see Cox [4].

2 Discrete time external arrival process model

In this section we describe a class of discrete time models for long-range dependent arrivals to a network. The duration between time l and $l+1$ is called *slot l* . At time l a random number ξ_l of new sessions arrive. The $(\xi_l, l \in \mathbf{Z})$ are independent and identically distributed (i.i.d.) Poisson random variables with mean λ . Session i arriving at time l is active for a random duration $\tau_l(i)$, i.e. it is active during the slots $l, l+1, \dots, l+\tau_l(i)-1$. Here the $(\tau_l(i), l \in \mathbf{Z}, 1 \leq i \leq \xi_l)$ are i.i.d. random variables taking values in \mathbf{N} . Let p_m denote $P(\tau_l(i) = m)$ and $q_m = \sum_{k=m}^{\infty} p_k$. The session duration random variables are assumed to be regularly varying with finite mean and infinite variance, i.e. there is a slowly varying function $L(\cdot)$ and $1 < \alpha < 2$ such that

$$q_k \sim k^{-\alpha} L(k) . \quad (1)$$

(See the appendix for the basic definitions on regular variation.) Each session active during a slot generates traffic at rate 1 during that slot. Note that, since there are many different slowly varying functions, there is considerable flexibility in the choice of arrival model.

Let $\xi_l(m)$ denote the number of sessions starting at time l that are active for exactly m slots. We have $\xi_l = \sum_{m=1}^{\infty} \xi_l(m)$, with $(\xi_l(m), l \in \mathbf{Z}, m \geq 1)$ being independent Poisson random variables, having mean λp_m for each m .

Let X_l denote the total amount of work brought in by the external arrival process during slot l . Then

$$X_l = \sum_{m=-\infty}^0 \sum_{j=1}^{\infty} \xi_{l+m}(-m+j) .$$

Note that $(X_l, l \in \mathbb{Z})$ is a stationary stochastic process. A straightforward calculation gives

$$r(k) \triangleq \text{Cov}(X_l, X_{l+k}) = \lambda \sum_{j=1}^{\infty} q_{k+j} \sim \frac{k^{-\alpha+1}}{\alpha-1} L(k)$$

Since $\sum_k r(k)$ is divergent when $1 < \alpha < 2$, the process $(X_l, l \in \mathbb{Z})$ is long-range dependent.

3 Discrete time quasi-reversible queuing network model

In this section we describe a class of discrete time queuing network models first introduced by Walrand [15]. We first define a kind of discrete time queue called an *S-queue*. An *S-queue* admits batch arrivals and has batch service. Given an arbitrary arrival sequence $\{\xi_n, n \geq 0\}$ of N valued random variables, the queue length process of an *S-queue* is given by

$$x_{n+1} = x_n + \xi_n - d_{n+1} , \quad (2)$$

where

$$P(d_{n+1} = j \mid x_n, d_m, 0 \leq m \leq n; \xi_k, k \geq 0, x_n + \xi_n = i) = S(i, j)$$

for $0 \leq j \leq i$ and $n \geq 0$. Here x_0 is arbitrary and $d_0 = 0$. Notice from the definition that the operation of an *S-queue* can be visualized as follows : there is a sequence of independent and identically distributed random variables $(d_n(i), i \geq 1)$, $-\infty < n < \infty$, with $P(d_n(i) = j) = S(i, j)$, $0 \leq j \leq i$, which is independent of the arrival process. This sequence can be thought of as a virtual departure process. At time $n+1$, if the the queue size just prior to release of the departures is i , we release $d_{n+1}(i)$ customers.

An *S-queue* is called *quasi-reversible* if, when $\{\xi_n, n \geq 0\}$ is a sequence of i.i.d. Poisson random variables such that the state admits an equilibrium distribution, then the sequence of actual departures $\{d_n, n \geq 0\}$ is also a sequence of i.i.d. Poisson random variables in equilibrium and for all n $\{d_l, l \leq n\}$ and x_n are independent. Walrand [15] gave a necessary and sufficient condition for an *S-queue* to be quasi-reversible. An *S-queue* is quasi-reversible if and only if there is a sequence of numbers $\alpha(j), j \geq 0$, where $\alpha(0) = 1$, $\alpha(j) > 0$ for $j > 0$ and a sequence of normalizing constants $c(i), i \geq 0$ such that $S(i, j)$ has the following form :

$$S(0, 0) = c(0) = 1 , \quad (2.1a)$$

$$S(i, 0) = c(i) , \quad i > 0 , \quad (2.1b)$$

$$S(i, j) = \frac{c(i)}{j!} \alpha(i) \alpha(i-1) \dots \alpha(i-j+1) , \quad 0 < j \leq i , \quad (2.1c)$$

$$\sum_{j=0}^i S(i, j) = 1 . \quad (2.1d)$$

Further, the queue admits an equilibrium distribution π for a Poisson arrival sequence of mean λ if and only if the normalizing constant κ exists such that

$$\pi(i) = \kappa \frac{\lambda^i}{\alpha(0) \dots \alpha(i)}, \quad i \geq 0$$

is a probability distribution.

A discrete time quasi-reversible network of queues can be constructed from J quasi-reversible S -queues (called nodes of the network) fed by an external arrival sequence of i.i.d. Poisson random variables, with Bernoulli routing. In more detail, such a network of queues operates as follows : An external arrival is routed to node p with probability r_{0p} . Such an arrival waits in queue at node p until it receives service; on receiving service it is routed to node q with probability r_{pq} and routed out of the network with probability r_{p0} . If routed to node q it joins the queue there. How many customers are served at each queue at each time is decided according to each queue's individual S -queue parameters, and is independent from queue to queue. The routing decisions are independent from customer to customer, and are also independent of the arrival process and the service decisions at the queues.

Walrand [15] has shown that, as long as the usual rate conditions are met (see below), such a network fed by an i.i.d sequence of Poisson random variables is stable and has a product form stationary distribution, where each queue in the network has a marginal stationary distribution as if it were being fed by an i.i.d. Poisson sequence of the appropriate rate. Here, by the *usual rate conditions* we mean the following : Let γ denote the rate of the external arrival process, and let λ_p , $1 \leq p \leq J$, be the unique solutions of the *flow balance equations* :

$$\gamma r_{0p} + \sum_q \lambda_q r_{qp} = \lambda_p, \quad 1 \leq p \leq J.$$

We require that for each $1 \leq p \leq J$, the S -queue p be stable when fed by an i.i.d. Poisson sequence of random variables of rate λ_p , i.e. that the normalizing constant κ_p can be defined such that

$$\pi_p(i) = \kappa_p \frac{\lambda_p^i}{\alpha_p(0) \dots \alpha_p(i)}, \quad i \geq 0 \quad (3)$$

is a probability distribution, where $\alpha_p(\cdot)$ are the parameters of queue p . By *product form stationary distribution* we mean that the stationary distribution of the number of backlogged customers at the nodes of the network is given by

$$\pi(i_1, \dots, i_J) = \prod_{p=1}^J \pi_p(i_p),$$

with $\pi_p(\cdot)$ as in (3).

To prove our main results we need to impose an additional condition on the parameters of the S -queues used to construct our network. An S -queue is called *monotone* if the sequence $\alpha(i)$, $i \geq 1$ is nondecreasing [1, Section IV]. In Lemma 2 of [1] it is proved that for

a monotone S-queue one has

$$\sum_{j=k}^{i+1} S(i+1, j) \geq \sum_{j=k}^i S(i, j) \text{ for all } 0 \leq k \leq i \text{ and } i \geq 0. \quad (4)$$

In addition, in Lemma 1 of [1] it is proved that for any S-queue (monotone or not) one has

$$\sum_{j=k}^{i+1} S(i+1, i+1-j) \geq \sum_{j=k}^i S(i, i-j) \text{ for all } 0 \leq k \leq i \text{ and } i \geq 0. \quad (5)$$

Together, equations (4) and (5) imply that it is possible to construct the virtual departure variables $((d_n(i), i \geq 1), n \in \mathbf{Z})$ of a monotone S-queue in such a way that one simultaneously has

$$d_n(i+1) \geq d_n(i) \text{ for all } i \geq 0 \text{ and } n \in \mathbf{Z}, \quad (6)$$

and

$$i+1 - d_n(i+1) \geq i - d_n(i) \text{ for all } i \geq 0 \text{ and } n \in \mathbf{Z}, \quad (7)$$

see Lemma 2 of [1]. This means that it is possible to couple the evolution of a network of monotone quasi-reversible S-queues from two different initial conditions, thereby using pathwise techniques to compare the evolution from these different initial conditions. What is meant by this last sentence will be made more clear in Section 5 where such a coupling is used.

4 Network model with long-range dependent traffic

In this section we observe that it is possible to combine the long-range dependent arrival models of the kind considered in Section 2 with the quasi-reversible queuing networks of S-queues described in Section 3 to create a self-contained class of network models that handle long-range dependent arrival processes. In the model we develop, we assume that service decisions are made at the session level without reference to the duration of the session. Further, once a session enters service at a node, we assume that it continues to receive service at rate 1 until all the work that it brings in is completed.

To make this precise, let us first consider what this means for a single quasi-reversible S-queue fed by such an arrival process, cf. Eqn. (2). Immediately after the service decision at time n , let there be x_n backlogged sessions at the queue. Here a backlogged session means one which has arrived at a time prior to time n but has not yet entered service. A session that has already entered service by time n is not considered backlogged even though it might still be currently being served, and indeed the work that it is to bring in might not yet all have arrived. Now at time n , ξ_n new sessions enter the queue as described by the external arrival process model. Then the decision to let d_{n+1} new sessions enter service at time $n+1$ is made according to the conditional probabilities of the S-queue. Finally this leaves behind $x_{n+1} = x_n + \xi_n - d_{n+1}$ sessions still backlogged after the service decision at time $n+1$.

Once a session enters service it continues to receive service at rate 1 until all the work that it brings in is exhausted. To clarify this further, note for instance that a session of duration τ that arrived at time $m < n$ and is still backlogged immediately after the service decision at time n will have brought in a total amount of work of $\tau \wedge (n - m)$ up to time n , which will be sitting in queue throughout slot n . If $\tau > n - m$ this backlogged session will be bringing in further work at rate 1 during slot n , while if $\tau \leq n - m$ no further work will be brought in by this session during slot n . If this session now enters service at time $n + 1$, and if $\tau \geq n + 1 - m$, a total of $n + 1 - m$ units of work will remain in queue until time $m + \tau$ after which the work in queue contributed by this session will decrease at rate one. If, on the other hand $\tau \leq n + 1 - m$, the work in queue contributed by this session will immediately begin to decrease at rate 1 starting from time $n + 1$.

The decision of which sessions to serve is not allowed to refer to the duration of the session. In the case of a single quasi-reversible S -queue fed by a long-range dependent arrival process of the kind considered in Section 2 it is then seen that in equilibrium the departure process is once again a process of the type considered in that section. This is nothing more than the observation that the stationary departure process of a quasi-reversible S -queue is a sequence of i.i.d. Poisson random variables, which is a direct consequence of the definition of quasi-reversibility.

To completely specify the model we also assume that sessions enter service in FCFS order of their arrival at the node, with ties broken by uniform randomization. The particular priority used to serve sessions will be seen to be irrelevant for the general results that we derive, although it would be important in a more detailed analysis of the stationary distribution of the performance quantities of interest.

In our general network model each S -queue functions at the session level as above. The Bernoulli routing between the quasi-reversible queues is also assumed to take place at the session level, i.e. all the work in a session travels along the same route. Heuristically, one can picture a session of duration τ as consisting of a “train” of τ units of work that extends over τ slots, with the work spread out uniformly over this duration. Immediately after the service decisions at all the queues at time n the head of the train of work corresponding to the session will either be at some queue p or will have left the network; in the former case we say that the session is backlogged at node p . It may happen that only a portion of the τ units of work associated with the session may be present in the network. This could happen for two reasons : If the head of the train of work left the network at time $m \leq n$ and $\tau > n - m$, then the initial $n - m$ units of work associated with the session have left the network forever; similarly, if the session first arrived at into the network at time $m \leq n$ and $\tau > n - m$, then the last $\tau - (n - m)$ units of work associated with the session have not even entered the network yet. Of the work associated with the session that is in the network, portions may be sitting at various queues. Rather than introduce burdensome notation, we illustrate the heuristic picture by means of the example of Figure 1.

This is a network of two S -queues, with routing parameters given by

$$r_{01} = 1, \quad r_{02} = 0, \quad r_{10} = r_{11} = 0, \quad r_{12} = 1, \quad r_{20} = r_{21} = \frac{1}{2}, \quad r_{22} = 0.$$

The precise S -queue parameters of the two queues are not relevant to the discussion. Assume

a session of duration 8 slots arrived in the network at time -6 , was served by node 1 at time -5 , served next by node 2 at time -4 , was routed back to node 1, served next by node 1 at time -2 , served next by node 2 at time -1 and routed out of the network. In this case we illustrate in Figure 2 where the work associated with this session is sitting immediately after the service decisions at time 0.

We remark that our network model allows a session to visit several different nodes, with the possibility of feedback. Further, of course, a session may visit several different nodes. Such a session will have the same duration on each visit to each node. It is basically this feature that makes the analysis of the general network model of this kind more interesting than the analysis of a single S-queue that was carried out at the beginning of this section.

To clarify the network model further, and to be able to state and prove the main results, we introduce some notation. Let $(a_n, n \in \mathbf{Z})$ denote the external arrival process of sessions into the network. This is a sequence of i.i.d. Poisson random variables of rate γ . As in Section 2 each session is active for an independent random duration having the distribution of τ , a positive integer valued random variable that is regularly varying with finite mean and infinite variance. Let $(a_n^p, n \in \mathbf{Z})$ denote the process of external arrivals of sessions at node p ; this is a sequence of i.i.d. Poisson random variables of rate γr_{0p} . Let $(\xi_n^p, n \in \mathbf{Z})$ denote the overall arrival process of sessions to node p ; in stationarity this is a stationary sequence of nonnegative integer valued random variables of rate λ_p . Let $(d_n^p, n \in \mathbf{Z})$ denote the process of departures of sessions from node p and $(d_n^{pq}, n \in \mathbf{Z})$ denote the process of departing sessions from node p that are routed to node q , $0 \leq q \leq J$ (with $q = 0$ corresponding to the departing sessions that leave the network). Then we have

$$\xi_n^p = a_n^p + \sum_{q=1}^J d_n^{qp}.$$

Further, with x_n^p denoting the number of backlogged sessions at node p immediately after the service decision at time n , we have

$$x_{n+1}^p = x_n^p + \xi_n^p - d_{n+1}^p. \quad (8)$$

In the following, we will sometimes write d_n^{0p} as an alternative notation for a_n^p .

Some important comments are in order at this point. First of all, note that the arrivals to a node at time n consist of the external arrivals to this node at time n and the departures from all the nodes to this node at time n . Since the arrivals occur after the departures, these sessions will only be available to be served at time $n + 1$. Thus at any time the head of a session can progress at most one node. Secondly, the proof of Theorem 3.1 of [15] shows that in equilibrium for each $n \in \mathbf{Z}$ it holds that

$$x_n^1, \dots, x_n^J, d_n^{10}, \dots, d_n^{1J}, \dots, d_n^{J0}, \dots, d_n^{JJ}, (a_m^1, m \geq n), \dots, (a_m^J, m \geq n)$$

is an independent family of random variables with all of them except perhaps the x_n^p being Poisson random variables of the appropriate means. (Note that this *does not* imply that

($d_n^{pq}, n \in \mathbb{Z}$) is an i.i.d. Poisson sequence.) Finally, the existence and uniqueness of a stationary regime for the network model when the usual rate conditions are met can be argued as follows : Because all service and routing decisions are carried out at the session level, the results of Section 3 ensure that there is a unique stationary regime for the process of the number of backlogged sessions at the nodes. Now, in this stationary regime, let \mathcal{D}_l denote the *set* of sessions that leave the network at time l (or more precisely the set of sessions whose head leaves the network at time l), and let \mathcal{X}_l denote the *set* of backlogged sessions in the network at time l . Further, for any session i let $M(i)$ denote the *number* of nodes the session visits (multiple visits to the same node are counted with their multiplicity). Consider the total amount of work in the network at time n , i.e. the amount of work the network would have to handle if all fresh arrivals from time n onwards were to be deleted. This can be shown to be a.s. finite as follows : It is upper bounded by

$$\sum_{i \in \mathcal{X}_n} M(i)\tau(i) + \sum_{l \leq n} \sum_{i \in \mathcal{D}_l} M(i)(\tau(i) - (n - l))^+ . \quad (9)$$

Since \mathcal{X}_n is a.s. finite, the first term of equation (9) is a.s. finite. Quasi-reversibility ensures that for each $j \geq 1$ ($|\mathcal{D}_l \cap \{i : M(i) = j\}|, l \leq n$) is a sequence of independent Poisson random variables of mean $\gamma\beta_j$, where β_j is the probability that an individual session visits j nodes during its sojourn in the network (this probability depends only on the routing probabilities). Further, since the service decisions and routing do not refer to the session durations, it follows that the number of terms of type $j m \geq 1$ that need to be summed in the second summation on the right hand side of equation (9) is a Poisson random variable of mean $\gamma\beta_j \sum_{m=0}^{\infty} p_{m+j} = \gamma\beta_j q_m$, and that these Poisson random variables are independent. Thus the second summation the right hand side of equation (9) can be rewritten as

$$\sum_{j=1}^{\infty} \sum_{m=1}^{\infty} j m Z_{jm}$$

for appropriately defined independent Poisson random variables ($Z_{jm}, j, m \geq 1$) of means $\gamma\beta_j q_m$ respectively. This sum is a.s. finite, as can be seen by noting that

$$\begin{aligned} P(Z_{jm} = 0 \forall j \geq J \text{ and } \forall m \geq M) &= \left(\prod_{m \geq M} \prod_{j \geq 1} e^{-\gamma\beta_j q_m} \right) \left(\prod_{1 \leq m < M} \prod_{j \geq J} e^{-\gamma\beta_j q_m} \right) \\ &\geq \left(\prod_{m \geq M} e^{-\gamma q_m} \right) \left(\prod_{j \geq J} e^{-\gamma\beta_j} \right) \\ &= e^{-\gamma(\sum_{m \geq M} q_m + \sum_{j \geq J} \beta_j)} \\ &\rightarrow 1 \text{ as } M \rightarrow \infty \text{ and } J \rightarrow \infty . \end{aligned}$$

The main results of this paper, derived in Section 5, refer to this stationary regime.

5 Main Results

For $0 \leq p, q \leq J$, let X_n^{pq} denote the total amount of work routed from node p to node q during slot n . (By convention, X_n^{0q} denotes the total amount of work arriving at node q

during slot n due to external arrivals and X_n^{p0} denotes the total amount of work leaving the network from node p during slot n .) Our main result is the following theorem.

Theorem 5.1 *Suppose all the nodes in the quasi-reversible network of S -queues defined in Section 4 are monotone. Fix any $0 \leq p, q \leq J$. In stationarity the process $(X_n^{pq}, n \in \mathbb{Z})$ is long-range dependent.*

Theorem 5.1 is a consequence of the following two lemmas. The first of these lemmas needs some motivation, so we give it here before formally stating the lemma. Our idea to prove theorem 5.1 is to first decompose the process d_n^{pq} of customers moving from node p to node q as a sum of processes. The individual processes in this decomposition are parametrized by increasing sequences of positive integers starting at zero. Given such a sequence $\mathbf{r} = \{0 = r_0 < r_1 < \dots < r_m\}$ for some $m \geq 0$, the process, say $d_n^{pq, \mathbf{r}}$, that is one of the constituents of this decomposition, consists of those customers who move from node p to node q a total of $m+1$ times of the first is time n , the next time $n + r_1$, and so on, the last being time $n + r_m$. This decomposition of d_n^{pq} also gives a convenient way of decomposing the process X_n^{pq} , enabling the proof of theorem 5.1.

Lemma 5.1 *For $m \geq 0$, let*

$$\mathcal{L}_m = \{0 = r_0 < r_1 < \dots < r_m\},$$

where $r_0, \dots, r_m \in \mathbb{N}$ be the set of all possible increasing sequences of positive integers of length $m+1$, starting at 0. Let $\mathcal{L} = \cup_{m=0}^{\infty} \mathcal{L}_m$. Suppose that for each $\mathbf{r} \in \mathcal{L}$ there is given a wide sense stationary \mathbb{N} -valued sequence $(\zeta_l^{\mathbf{r}}, l \in \mathbb{Z})$ such that these sequences are also jointly wide sense stationary.

Let $\mathcal{Z}_l^{\mathbf{r}}$ denote a set of cardinality $\zeta_l^{\mathbf{r}}$. We think of each $i \in \mathcal{Z}_l^{\mathbf{r}}$ $l \in \mathbb{Z}$, $\mathbf{r} \in \mathcal{L}$ as a session, distinct identical copies of which become active for the first time at $l + r_0, \dots, l + r_m$. By this we mean the following : Let $\tau(i)$ be the duration of such a session i , which has regularly varying distribution as in Section 2; the session duration variables are assumed to be independent and independent of the processes $(\zeta_l^{\mathbf{r}}, l \in \mathbb{Z})$, $\mathbf{r} \in \mathcal{L}$. Then one copy of this session is active during the slots $l + r_0, \dots, l + r_0 + \tau(i) - 1$, another during the slots $l + r_1, \dots, l + r_1 + \tau(i) - 1$, and so on for a total of $m+1$ copies, the last of which is active during the slots $l + r_m, \dots, l + r_m + \tau(i) - 1$.

Let Y_n denote the total number of copies of sessions that become active at time n . Thus

$$Y_n = \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \zeta_l^{\mathbf{r}} \sum_{a=0}^m 1(l + r_a = n). \quad (10)$$

During each slot during which any copy of any session is active, it brings in one unit of work. Let X_n denote the total amount of work brought in during slot n . Thus

$$X_n = \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \sum_{i \in \mathcal{Z}_l^{\mathbf{r}}} \sum_{a=0}^m 1(l + r_a \leq n \leq l + r_a + \tau(i) - 1) \quad (11)$$

Assume that $(Y_n, n \in \mathbf{Z})$ and $(X_n, n \in \mathbf{Z})$ are well defined (a.s. finite) processes. Then, if $(Y_n, n \in \mathbf{Z})$ is short-range dependent, it holds that $(X_n, n \in \mathbf{Z})$ is long-range dependent. \square

Lemma 5.2 Consider a discrete time quasi-reversible network of J monotone S -queues with Bernoulli routing, in stationarity. For each $0 \leq p, q \leq J$, the process $(d_n^{pq}, n \in \mathbf{Z})$ is short-range dependent. \square

We first prove the main theorem using these lemmas, and then prove the lemmas.

Proof of Theorem 5.1 :

Since each $(a_n^p, n \in \mathbf{Z})$ and each $(d_n^{p0}, n \in \mathbf{Z})$, $1 \leq p \leq J$ is an i.i.d. sequence of Poisson random variables, it is immediate that each $(X_n^{0p}, n \in \mathbf{Z})$ and each $(X_n^{p0}, n \in \mathbf{Z})$, $1 \leq p \leq J$ is a process of the kind considered in Section 2, and is therefore long-range dependent. We therefore focus on the case $1 \leq p, q \leq J$.

It suffices to identify the appropriate sequences $(\zeta_l^r, l \in \mathbf{Z})$, $r \in \mathcal{L}$ so as to apply Lemma 5.1. Given fixed $1 \leq p, q \leq J$, for $m \geq 0$ and $r \in \mathcal{L}_m$, define ζ_l^r to be the number of sessions that are routed from node p to node q at precisely the times $l, l+r_1, \dots, l+r_m$. When the network is in stationarity, the sequences $(\zeta_l^r, l \in \mathbf{Z})$, $r \in \mathcal{L}$ are jointly stationary. Further, for each $r \in \mathcal{L}$ and $l \in \mathbf{Z}$ we have $\zeta_l^r \leq d_l^{pq}$, and d_l^{pq} is a Poisson random variable of mean $\lambda_p r_{pq}$. From this it follows that each ζ_l^r has finite second moment. Since the processes $(\zeta_l^r, l \in \mathbf{Z})$, $r \in \mathcal{L}$ are jointly stationary and have finite second moment, they are jointly wide sense stationary, so that Lemma 5.1 can be applied. Note that we also have, following equations (10) and (11), that

$$Y_n = d_n^{pq} \text{ and } X_n = X_n^{pq}.$$

By Lemma 5.2 we see that $(Y_n, n \in \mathbf{Z})$ is short-range dependent. By Lemma 5.1, the claim of the theorem follows. \square

Proof of Lemma 5.1 :

Let τ be a generic \mathbf{N} -valued random variable having the distribution of the session duration random variables, and let $l \in \mathbf{Z}$. We write $C(l, \tau)$ for the indicator of the event that slot 0 is one of the slots $l, \dots, l + \tau - 1$. In other words,

$$\begin{aligned} C(l, \tau) &= 1(l \leq 0 \leq l + \tau - 1) \\ &= 1(l \leq 0) - 1(l + \tau \leq 0) \end{aligned} \tag{12}$$

Note that

$$E[C(l, \tau)] = \begin{cases} 0 & \text{if } l > 0 \\ P(l + \tau \geq 1) = q_{1-l} & \text{if } l \leq 0 \end{cases} \quad (13)$$

For convenience, set $q_k = 0$ if $k \leq 0$. With this notation, we may then write equation (13) as $E[C(l, \tau)] = q_{1-l}$ for all $l \in \mathbb{Z}$.

Let \mathcal{F} denote the sigma-field generated by the variables $\zeta_l^{\mathbf{r}}$, $l \in \mathbb{Z}$, $\mathbf{r} \in \mathcal{L}$. Let $r(k) = \text{Cov}(X_l, X_{l+k})$. Then we may write

$$\begin{aligned} r(k) &= E[X_0 X_k] - E[X_0]E[X_k] \\ &= E[E[X_0 X_k | \mathcal{F}] - E[X_0 | \mathcal{F}]E[X_k | \mathcal{F}]] \\ &\quad + E[E[X_0 | \mathcal{F}]E[X_k | \mathcal{F}] - E[X_0]E[X_k]] \\ &= r^{(1)}(k) + r^{(2)}(k). \end{aligned} \quad (14)$$

In the following we will handle each of the two terms on the right hand side of equation (14) separately.

With the notation introduced in equation (12), we may rewrite equation (11) as

$$X_n = \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \sum_{i \in \mathbb{Z}_l^{\mathbf{r}}} \sum_{a=0}^m C(l - n + r_a, \tau(i)) \quad (15)$$

Using equation (13), we see from equation (15) that

$$E[X_n | \mathcal{F}] = \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \zeta_l^{\mathbf{r}} \sum_{a=0}^m q_{n+1-l-r_a}. \quad (16)$$

We turn to the evaluation of the first term on the right hand side of equation (14). Note that

$$E[X_0 X_k | \mathcal{F}] - E[X_0 | \mathcal{F}]E[X_k | \mathcal{F}] = E[(X_0 - E[X_0 | \mathcal{F}])(X_k - E[X_k | \mathcal{F}]) | \mathcal{F}],$$

is the conditional covariance of X_0 and X_k given \mathcal{F} . We find that

$$\text{Cov}(C(l + r_a, \tau), C(l - k + r_b, \tau)) = q_{1-l-r_a} \wedge q_{k+1-l-r_b} - q_{1-l-r_a} q_{k+1-l-r_b}, \quad (17)$$

which is nonnegative. From equation (15) and the conditional independence of the various $\tau(i)$ given \mathcal{F} , we see that

$$\begin{aligned} E[X_0 X_k | \mathcal{F}] - E[X_0 | \mathcal{F}]E[X_k | \mathcal{F}] &= \\ &= \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \sum_{i \in \mathbb{Z}_l^{\mathbf{r}}} \sum_{a=0}^m \sum_{b=0}^m \text{Cov}(C(l + r_a, \tau(i)), C(l - k + r_b, \tau(i))), \end{aligned}$$

because all the cross-covariances vanish.

Let $\lambda^{\mathbf{r}}$ denote $E[\zeta_l^{\mathbf{r}}]$. From the preceding equation and equations (13) and (17) we have

$$r^{(1)}(k) = \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \lambda^{\mathbf{r}} \sum_{a=0}^m \sum_{b=0}^m (q_{1-l-r_a} \wedge q_{k+1-l-r_b} - q_{1-l-r_a} q_{k+1-l-r_b}) \quad (18)$$

Pick some $\mathbf{r} \in \mathcal{L}$ for which $\lambda^{\mathbf{r}} = \lambda > 0$ (there must be at least one such, otherwise the situation considered in the lemma is vacuous). Noting that every term on the right hand side of equation(18) is nonnegative, we may write

$$r^{(1)}(k) \geq \lambda \sum_{l \in \mathbb{Z}} (q_{1-l} \wedge q_{k+1-l} - q_{1-l} q_{k+1-l}) \quad (19)$$

Noting that $q_{1-l} = 0$ if $l > 0$ and that $q_{k+1-l} \leq q_{1-l}$ if $k \geq 0$, we have

$$\begin{aligned} \sum_{k \geq 0} r^{(1)}(k) &\geq \lambda \sum_{l \leq 0} (1 - q_{1-l}) \sum_{k \geq 0} q_{k+1-l} \\ &= \infty, \end{aligned}$$

where we have used equation (1) for the tail probability asymptotics of the session duration random variables.

We turn next to the evaluation of the second term on the right hand side of equation (14). Note that

$$r^{(2)}(k) = \text{Cov}(E[X_0 | \mathcal{F}], E[X_k | \mathcal{F}]) . \quad (20)$$

Using the relation

$$q_j = \sum_{t \in \mathbb{Z}} q_t 1(t = j)$$

and the definition of Y_n in equation (10) we may rewrite equation (16) to get

$$\begin{aligned} E[X_n | \mathcal{F}] &= \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \zeta_l^{\mathbf{r}} \sum_{a=0}^m q_{n+1-l-r_a} \\ &= \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \zeta_l^{\mathbf{r}} \sum_{a=0}^m \sum_{t \in \mathbb{Z}} q_t 1(t = n+1-l-r_a) \\ &= \sum_{t \in \mathbb{Z}} q_t \sum_{l \in \mathbb{Z}} \sum_{m \geq 0} \sum_{\mathbf{r} \in \mathcal{L}_m} \zeta_l^{\mathbf{r}} \sum_{a=0}^m 1(l+r_a = n+1-t) \\ &= \sum_{t \in \mathbb{Z}} q_t Y_{n+1-t} . \end{aligned} \quad (21)$$

Now let

$$\Gamma_j = \text{Cov}(Y_0, Y_j) . \quad (22)$$

From equations (20), (21), and (22) we have

$$\begin{aligned}
 r^{(2)}(k) &= \text{Cov}(E[X_0 | \mathcal{F}], E[X_k | \mathcal{F}]) \\
 &= \text{Cov}\left(\sum_{t \in \mathbb{Z}} q_t Y_{1-t}, \sum_{s \in \mathbb{Z}} q_s Y_{k+1-s}\right) \\
 &= \sum_{t \in \mathbb{Z}} \sum_{s \in \mathbb{Z}} q_t q_s \Gamma_{k+t-s} \\
 &= \sum_{j \in \mathbb{Z}} \Gamma_j \sum_{t \in \mathbb{Z}} q_t q_{k+t-j} .
 \end{aligned}$$

Thus

$$\begin{aligned}
 \sum_{k \in \mathbb{Z}} |r^{(2)}(k)| &= \sum_{k \in \mathbb{Z}} \left| \sum_{j \in \mathbb{Z}} \Gamma_j \sum_{t \in \mathbb{Z}} q_t q_{k+t-j} \right| \\
 &= \sum_{k \in \mathbb{Z}} \sum_{j \in \mathbb{Z}} |\Gamma_j| \sum_{t \in \mathbb{Z}} q_t q_{k+t-j} \\
 &= \sum_{j \in \mathbb{Z}} |\Gamma_j| \sum_{t \in \mathbb{Z}} q_t \sum_{k \in \mathbb{Z}} q_{k+t-j} \\
 &\leq (E[\tau])^2 \sum_{j \in \mathbb{Z}} |\Gamma_j| \\
 &< \infty ,
 \end{aligned}$$

where we have used the assumption that $\sum_{j \in \mathbb{Z}} |\Gamma_j| < \infty$.

This completes the proof of the lemma.

□

Proof of Lemma 5.2 :

First note that for any $1 \leq p \leq J$, each of the processes $(d_n^{0p}, n \in \mathbb{Z})$ and $(d_n^{p0}, n \in \mathbb{Z})$ is an i.i.d. sequence of Poisson random variables, and is therefore short-range dependent. Thus we focus on the case $1 \leq p, q \leq J$.

It is in the proof of this lemma that we exploit the assumption that each of the quasi-reversible S-queues used to construct our network is monotone. Fix $1 \leq p, q \leq J$. Our goal is to demonstrate that

$$\sum_{k \geq 0} |\text{Cov}(d_0^{pq}, d_k^{pq})| < \infty . \tag{23}$$

Note that

$$\text{Cov}(d_0^{pq}, d_k^{pq}) = E[d_0^{pq} d_k^{pq}] - E[d_0^{pq}] E[d_k^{pq}]$$

$$= \sum_{a \in \mathbb{N}} aP(d_0^{pq} = a) (E[d_k^{pq} | d_0^{pq} = a] - E[d_k^{pq}]) ,$$

so that

$$\sum_{k \geq 0} |\text{Cov}(d_0^{pq}, d_k^{pq})| \leq \sum_{a \in \mathbb{N}} aP(d_0^{pq} = a) \sum_{k \geq 0} |E[d_k^{pq} | d_0^{pq} = a] - E[d_k^{pq}]| .$$

We now proceed to estimate $\sum_{k \geq 0} |E[d_k^{pq} | d_0^{pq} = a] - E[d_k^{pq}]|$. By the independence of

$$x_0^1, \dots, x_0^J, d_0^{10}, \dots, d_0^{1J}, \dots, d_0^{J0}, \dots, d_0^{JJ}, (a_m^1, m \geq 0), \dots, (a_m^J, m \geq 0)$$

in stationarity, conditioning on $d_0^{pq} = a$ is equivalent to setting up an initial condition at the queues just prior to the service decisions at time 1 that differs from the equilibrium situation *only* in that node q was fed a customers from node p instead of the Poisson mean $\lambda_p r_{pq}$ customers that would have been fed in stationarity. But now we may couple the two initial conditions (the stationary situation and the situation where we have conditioned on $d_0^{pq} = a$) by the following simple expedient : If there are more customers than a fed back from node p to node q in the stationary situation, we color the extra customers in the stationary situation red, while if there are less customers than a fed back from node p to node q in the stationary situation we color the extra customers in the conditioned situation yellow. We then run the network (which consists entirely of either colorless customers and red customers or of colorless customers and yellow customers), with the convention that colorless customers have priority over colored customers.

The virtual departure variables $d_n^p(i)$, $i \geq 0$, $n \in \mathbb{Z}$, $1 \leq p \leq J$ obey the conditions (6) and (7). Consider first the case where there are less customers than a fed from node p to node q in the stationary situation. Then the evolution of colorless customers will be precisely as in the stationary situation, and the evolution of the sum of colored and colorless customers will be exactly as in the situation where we condition on $d_0^{pq} = a$. To understand why this follows from the properties (6) and (7), consider a node r during slot m , and suppose there are k colorless customers and l colored customers in queue just prior to time $m+1$, and the decision needs to be made as to how many customers to release at time $m+1$, in both the stationary and the conditioned situation. One needs to refer to the virtual departure variables $d_{m+1}^r(k)$ and $d_{m+1}^r(k+l)$ in order to make these decisions. Property (6) of the virtual departure variables ensures that $d_{m+1}^r(k+l) \geq d_{m+1}^r(k)$, and property (7) ensures that $d_{m+1}^r(k+l) - d_{m+1}^r(k) \leq l$, so that, giving priority to colorless customers over colored customers, we can release $d_{m+1}^r(k)$ colorless customers and $d_{m+1}^r(k+l) - d_{m+1}^r(k)$ colored customers, and simultaneously meet the requirements of both the stationary situation and the conditioned situation.

Similarly, in the case where there are less customers than a fed from node p to node q in the stationary situation, the evolution of the sum of colored and colorless customers will be precisely as in the stationary situation, while the evolution of colorless customers will be precisely as in the situation where we condition on $d_0^{pq} = a$.

In either case, one sees that $\sum_{k \geq 0} |E[d_k^{pq} | d_0^{pq} = a] - E[d_k^{pq}]|$ equals the mean number of times colored customers move from node p to node q at times $k \geq 0$. This is

upper bounded by

$$\kappa E[|Z - a|] \quad (24)$$

for a constant κ related to the routing variables and where Z is a Poisson random variable of mean $\lambda_p r_{pq}$. To see this, note that the number of movements made from node p to node q by the individual customers that start in node p are independent random variables, because routing decisions are independent from customer to customer. The mean number of such movements can be taken to be the constant κ , which is finite.

The expression in (24) can be further upper bounded by

$$\kappa \lambda_p r_{pq} + \kappa a$$

so that we get

$$\sum_{k \geq 0} |\text{Cov}(d_0^{pq}, d_k^{pq})| \leq \sum_{a \in \mathbb{N}} a(\kappa \lambda_p r_{pq} + \kappa a) P(d_0^{pq} = a) .$$

But d_0^{pq} is itself a Poisson random variable of mean $\lambda_p r_{pq}$. This gives the desired equation (23).

□

References

- [1] ANANTHARAM, V., The Input-Output Map of a Monotone Discrete time Quasi-reversible Node. *IEEE Transactions on Information Theory*, Vol. 39, No. 2, pp. 543-552, 1993. Correction published in Vol. 39, No. 4, pg. 1466, 1993.
- [2] ANANTHARAM, V., On the Sojourn Time of Sessions at an ATM Buffer with Long-Range Dependent Input Traffic. *Proceedings of the 34th IEEE Conference on Decision and Control*, New Orleans, December 1995.
- [3] BERAN, J., SHERMAN, R., TAQQU, M. S., AND WILLINGER, W., Long-Range Dependence in Variable-Bit-Rate Video Traffic. *IEEE Transactions on Communications*, Vol. 43, No. 2/3/4, pp. 1566-1579, 1995.
- [4] COX, D. R., Long-Range Dependence : A Review. In *Statistics : An Appraisal*, edited by H. A. David and H. T. David, Iowa State University Press, pp. 55-74, 1984.
- [5] DUFFIELD, N. G., AND O'CONNELL, N., Large Deviations and Overflow Probabilities for the General Single-server Queue, with Applications. To appear in *Proceedings of the Cambridge Philosophical Society*, 1995.
- [6] FELLER, W., *An Introduction to Probability Theory and its Application*. John Wiley and Sons, New York, 1971.
- [7] KELLY, F., *Reversibility and Stochastic Networks*. John Wiley and Sons, New York, 1979.

- [8] LELAND, W. E., TAQQU, M. S., WILLINGER, W., AND WILSON, D. V., On the Self-Similar Nature of Ethernet Traffic (Extended Version). *IEEE/ACM Transactions on Networking*, Vol. 2, No. 1, pp. 1 -15, 1994.
- [9] LIKHANOV, N., TSYBAKOV, B., AND GEORGANAS, N. D., Analysis of an ATM Buffer with Self-Similar ("Fractal") Input Traffic. *Proceedings of the 14 th Annual IEEE Infocom*, pp. 985 -992, 1995.
- [10] NORROS, I., A storage model with self-similar input. *Queueing Systems : Theory and Applications*, Vol. 16, pp. 387 -396, 1994.
- [11] PARULEKAR, M. AND MAKOWSKI, A. M., Buffer Overflow Probabilities for a Multiplexer with Self-similar Traffic. *Proceedings of the 34th IEEE Conference on Decision and Control*, New Orleans, December 1995.
- [12] PAXSON, V. AND FLOYD, S., Wide Area Traffic : The Failure of Poisson Modeling. *IEEE/ACM Transactions on Networking*, Vol. 3, No. 3, pp. 226 -244, 1995.
- [13] RESNICK, S. AND SAMORODNITSKY, G., The Effect of Long Range Dependence in a Simple Queuing Model. *Preprint*, Cornell University, 14 pp., 1994.
- [14] WALRAND, J., *An Introduction to Queuing Networks*. Prentice Hall, Englewood Cliffs, N.J., 1988.
- [15] WALRAND, J., A Discrete-Time Queuing Network. *Journal of Applied Probability*, Vol. 20, pp. 903-909, 1983.

Appendix

Regular Variation

For convenience we reproduce the basic definitions on nonnegative random variables with regularly varying tails. For more information, consult Feller [6, pp. 275-284].

Definition 1 A positive (not necessarily monotone) function $L(\cdot)$ defined on an interval $[a, \infty)$ is said to be slowly varying if for every $x > 0$ it holds that

$$\lim_{t \rightarrow \infty} \frac{L(tx)}{L(t)} = 1 .$$

Note that there are many slowly varying functions. For instance, any power of any iterated logarithm is slowly varying; any positive function that approaches a strictly positive limit at ∞ is slowly varying.

Definition 2 A positive function $G(\cdot)$ defined on an interval $[a, \infty)$ is said to be regularly varying with exponent $-\alpha$ if it is of the form

$$G(x) = x^{-\alpha} L(x)$$

where $L(\cdot)$ is a slowly varying function.

In this paper we are only interested in the case $1 < \alpha < 2$.

Definition 3 *A nonnegative random variable X having cumulative distribution function $F(x) = P(X \leq x)$ is said to have a regularly varying tail if $1 - F(x)$ is regularly varying.*

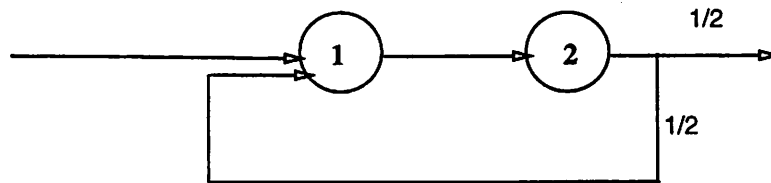


Figure 1: An example network.

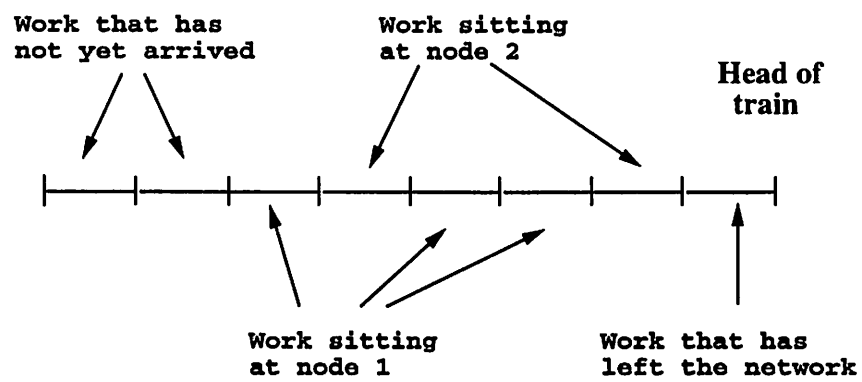


Figure 2: An example illustrating where the work associated with a session might be located.