MODELLING CIRCUIT-SWITCHED MULTI-STAGE INTERCONNECTION NETWORKS

P HARRISON

Rapporteur: S Caughey



MODELLING CIRCUIT-SWITCHED MULTI-STAGE INTERCONNECTION NETWORKS[†]

Peter G Harrison and Naresh M Patel Department of Computing Imperial College London

Abstract

A major component of a parallel machine is its interconnection network (IN), which provides concurrent communication between the processing elements. It is common to use a multi-stage interconnection network (MIN) which is constructed using crossbar switches and introduces contention not only for destination addresses but also for internal links. Both types of contention are increased when non-local communication across a MIN becomes concentrated on a certain destination address, the *hot-spot*. This paper considers analytical models of asynchronous, circuit-switched INs in which partial paths are held during path building, beginning with a single crossbar and extending recursively to MINs. Since a path must be held between source and destination processors before data can be transmitted, switching networks are passive resources and queueing networks which include them do not therefore have product-form solutions. Using decomposition techniques, the flow equivalent server (FES) that represents a bank of devices transmitting through a switching network is determined, under mild approximating assumptions. In the case of a full crossbar the FES can be solved directly and the result can be applied recursively to model the MIN. Two cases are considered: one in which there is uniform routing and the other where there is a hot-spot at one of the output pins. Validation with respect to simulation for MINs with up to six stages (64-way switching) indicates a high degree of accuracy in the models.

⁺ The majority of this paper is to appear in the Journal of the ACM as "The representation of multi-stage interconnection networks in queueing models of parallel systems". It may only be reproduced subject to the ACM copyright.

1 Introduction

In this paper we consider the case of a circuit-switched multi-stage interconnection network. Here, the network is modelled as a *passive* resource in a queueing network; in order to operate a data-transmission server must hold a path comprising internal links and crossbars at each stage through to the destination addressed by the task at the front of its queue. For packet switching, such paths need not be held throughout a message's transmission, and buffered internal crossbars can be modelled as conventional servers since their transmission times are non-negligible; we saw this in the previous paper. Similar problems to those of circuit-switching also arise if the buffers can become exhausted and some form of blocking ensues.

Our approach to the problem relies on the recursive structure of the MINs we consider which is in contrast to analogous research reported by Kelly on telephone networks, [KE86]. In either case, we could in principle consider all routes through the network individually and solve the associated Markov process directly to obtain the Erlang loss formula - assuming, that is, that partial paths are not held but released on a collision, resulting in lost transmissions. Because of the large number of routes, such a direct solution is computationally intractable. Kelly's analysis relies on the network behaving as if the steady state probabilities that each link is blocked were independent, which property is shown to hold asymptotically as link capacities and network traffic approach infinity jointly so that the traffic offered per link remains constant. In fact the approximation is very accurate for telephone networks because of the large capacity links and large number of nodes. Our analysis differs from this in that links have small capacities (viz. one) and also that partial paths are held during path building. Moreover, partly as a result of the former difference, there is a strong dependence between individual links which do not behave remotely as if they were independent, and it is our recursive analysis which captures this dependence in a simple way. In this sense, our analysis addresses the opposite end of the spectrum of switching networks and their control protocols.

In the next section, we define recursively the network topology we will be using and introduce some new terminology. Our general approach to modelling a queueing subnetwork that represents a collection of DMA servers together with a circuit-switched MIN is considered in Section 3. First we apply decomposition techniques to obtain the service rate of a *flow equivalent server* (FES) by determining the throughput of the short-circuited subnetwork as a function of its population, according to the method of

[CHW75]. This throughput is obtained by defining a simple Markov process, under appropriate assumptions, in which some function, μ_n , for throughput is assumed to be available if the number of active *input* pins, n, is known - i.e. the number of DMA servers currently wishing to transmit. In Sections 4 and 5 we derive expressions for μ_n in the cases that the interconnection network is respectively a *full crossbar switch* and a *delta network*, which is itself constructed from and analysed in terms of crossbars. The degree of the degradation introduced by these networks on the DMA servers' throughput is illustrated by numerical predictions in Section 7. In these analyses, the networks are *uniform* which means that all output destination addresses (pins) are selected with equal probability by any transaction, but in Section 6 we consider the case where the destination address space may have a *hot-spot*, [PN85], i.e. one pin which is selected with a higher probability than the rest, which are still all selected with the same (reduced) probability. Our conclusions and suggestions for future work are laid out in Section 8.

2 The partial shuffle topology and terminology

For our recursive analysis we adopt the *partial shuffle* variant of the *cube-topology*, [PA81], which may be defined recursively as follows for networks with 2-way crossbars (the definition can be generalised easily to networks constructed from b-way crossbars with b>2):

- (i) A one-stage network, Π_1 , is the single 2-way crossbar
- (ii) An s-stage network, Π_s , (s>0) has 2^s inputs and outputs, i.e. 2^{s-1} switches in each stage, and is constructed by connecting a *head stage* of 2^{s-1} switches to the right of 2 *tail networks* of (s-1) stages according to the partial shuffle topology shown in Figure 2. The *i*th switch in stage s takes its top input and bottom input from the *i*th pin of the upper subnetwork and the *i*th pin of the lower subnetwork, respectively. This property of the topology is particularly useful in obtaining the recurrence formulae in Sections 5 and 6.

In order to associate various features of the delta network with more familiar terms, used for example in describing graphs and trees, we use certain words synonymously: the full crossbar switch is sometimes referred to as a *node* and connections between switches in adjacent stages are called *links*. A sequence of connected nodes and links between the network input pins and output pins is called a *path*, and any contiguous portion of a path which begins at an input pin is called a *partial path*. The *decode tree* of a network *output* pin comprises all paths to that pin, although in this paper we only need to consider the

nodes in those paths. Likewise, the *decode tree* of a network *input* pin is composed of all paths originating from that pin. The pins in any stage of the network are numbered consecutively, starting at zero for the top pin, as shown in Figure 2.





In addition to this numbering, each output pin belongs to a *class* such that class 0 contains the pin numbered 0 and class k contains all pins numbered 2^{k-1} to 2^{k-1} inclusively (k>0), as shown in Figure 3. These classes will be used to distinguish different degrees of hot-spot contention in Section 6.



Figure 3 Pin Classes in a 8-way Network

In the analysis of non-uniform routing within delta networks, the top pin is taken to be the hot-spot, i.e. the top pin address is chosen by an arriving customer with some probability ρ and all other addresses are chosen with the same probability $\frac{1-\rho}{w-1}$ in the case of w-way connection. The top pin will be referred to as the *hot pin* and the other pins as *cool pins*.

3 Interconnection Networks in Queueing Models

The interconnection network is a passive resource in that paths in it need to be held by active resources (DMA processors) for them to provide their service. This is an example of service blocking and is typically found in circuit-switched networks. There are several types of blocking that can occur in queueing networks and approximate solutions for such networks can be found in [AK87]. However, such methods become impractical when there are large numbers of servers and in the following sections we describe how aggregation techniques can be applied to a queueing network with an IN.

3.1 Flow Equivalent Server for the Interconnection Network

In a locally balanced queueing network, it is a standard result that any subnetwork can be replaced by a FES with service rate equal to the throughput of the short-circuited subnetwork when its population is equal to the FES queue length, [CHW75]. We apply the same method to subnetworks which comprise a bank of DMA servers together with some interconnection network, the remainder of the enclosing network satisfying the conditions for local balance. Under appropriate assumptions, the behaviour of the shortcircuited subnetwork can be represented by a stationary Markov process, and for each valid population we find its throughput and hence fes-rate by solving the balance equations of this process. The interconnection network is represented by a function which determines the effective service rate of the bank of DMA devices when they compete for network paths, and the appropriate functions for full crossbar and delta networks are derived in Sections 4.1 and 5.1 respectively; in the form of a recurrence relation that exploits the recursive structure of the network in the latter case. In Section 3.3, by considering the limiting case in which there is no degradation, i.e. every active input pin of a network is always connected to a destination, the well-known expression for the throughput of a multiple server is obtained, and later we also derive a simple formula in the case that the interconnection network is a crossbar (Section 4.2).

3.2 Underlying Markov Process

The flow equivalent server for b parallel DMA channels connected through an interconnection network is defined by the short-circuited network shown in Figure 4 for populations N>1. When N=1 the customer experiences no contention for links and so the throughput is the same as the DMA service rate. The shaded box represents the switching network which limits the number of active outputs, m, to some value between 1 and the number of active inputs, n, according to its internal connection topology.

The switching network may be a full crossbar or a delta network with circuit switching and partial paths held. In the steady state, when there are n active inputs to the network we assume that they are uniformly distributed over the network input pins. We also make the simplifying assumption that path set up and release times are negligible compared to service durations, i.e. the path across the network or partial path to the first blockage will be established instantly. In practice, this is a reasonable assumption in networks where message lengths (and so DMA transfer times) are large when compared to crossbar switching times. For example, in the ALICE machine the crossbar can switch in 85ns whereas the mean transfer time is $14.4\mu s$.



Figure 4 The short-circuited subnetwork defining the fes rate

The state space, Ω_1 , for this closed system is defined by $\Omega_1 = \{(\underline{n}, \underline{m}) | \sum n_i = N, #(\underline{n}) \ge \underline{m} \ge 1\}$, where the population is N>1, $\underline{n} = (n_1, ..., n_b)$ is a vector of queue lengths at the source servers $(n_i \ge 0, 1 \le i \le b)$, m is a corresponding number of active output pins, and $#(\underline{n})$ is the number of non-zero components of \underline{n} ; we will use n to denote $#(\underline{n})$ where there is no confusion. Now let the random variable $N_i(t)$ denote the number of tasks at server i $(1 \le i \le b)$, and M(t) the number of active outputs at time $t\ge 0$. $X(t) = \{(\underline{N}(t), M(t)) | t\ge 0\}$ has finite state space Ω_1 and is a non-null recurrent Markov process, with generators defined by appropriate state transition probabilities, under the assumption that all servers have negative exponentially distributed service times; this assumption can be

relaxed for various queueing disciplines in standard ways, e.g. [BCMP75].

We could now try to obtain the balance equations for the state space probabilities of X(t)in the steady state by considering every possible state transition which may occur on a service completion. This results in a large number of equations and the approach is impractical even in the simple case of a full crossbar. Moreover, in a similar analysis of a delta network, the state space Ω_1 must be extended to represent every path and partial path established, vastly increasing computational complexity.

We therefore consider a simpler process, Y(t), in which the switching network is represented by an expression v_n for the expected number of active outputs when there are n active inputs in the steady state and we assume (approximately) that this relationship holds at all times. Thus, in the case of a direct connection (no switching network) every 'input' will always be connected to an 'output' and we have $v_n=n$ which does always hold exactly. We will also use $\mu_n = v_n\mu$ to denote the *effective service rate function* of a bank of servers, each with rate μ , connected through such a network when n servers have non-empty queues. Thus it is no longer necessary to include the number of active outputs in the state of the process Y(t), since we can approximate this by v_n when there are n active inputs. When the population of the short-circuited subnetwork is N, the state space is $\Omega_2 = \{\underline{n} \mid \Sigma n_i = N\}$ and we denote the stationary probability of $\underline{n} \in \Omega_2$ by $\pi(\underline{n})$.

3.3 Underlying Birth and Death Process

We now aggregate all states <u>n</u> with the same number of non-zero components and consider the integer-valued, stationary Markov process Z(t) which describes the number of active inputs to the switching network at time t and has state space $\Omega_3 = \{n \mid 1 \le n \le \min(b,N)\}$ with stationary probability p_n for $n \in \Omega_3$, i.e. p_n is the steady state probability that n input pins are active. The aggregate state n therefore represents the subset of states $S_n = \{\underline{n} \mid \#(\underline{n})=n, \sum n_i=N\}$ of the process Y(t), where the sets S_n form a partition of Ω_2 .

For a given n, each state in S_n has the same steady state probability with respect to the process Y(t), i.e. each arrangement of customers in any n non-empty queues is equally likely. This follows because the generators of the Markov process Y(t) depend only on n and not directly on the component values n_i. The balance equations for Z(t) then follow by expressing $\pi(\underline{n})$ in terms of p_n for each $\underline{n} \in \Omega_2$, [HA87,PA89]. Of course a stronger result holds in a simple Markovian queueing network (with no switching network) when

Modelling circuit-switched multi-stage interconnection networks

the visitation rate to service rate ratio is the same for all servers, namely that *all* states are equi-probable (cf. the multiple server result discussed below).

We use the following elementary combinatorial results in several places in this paper:-

- (i) The number of ways of arranging n balls in m boxes is $^{n+m-1}C_{m-1}$, (1) where the combinatorial function, C, is defined by $^{n}C_{m} = \frac{n!}{m!(n-m)!}$
- (ii) The number of ways of arranging n balls in m *non-empty* boxes, that is so that there is at least one ball in each of the m boxes, is the same as the number of ways of arranging n-m balls in m boxes (m having already been accounted for), i.e. $^{n-1}C_{m-1}$ (2)
- (iii) $\frac{n}{m}^{n-1}C_{m-1} = {}^{n}C_{m}$ for integers $n \ge m > 0$.

Thus we deduce from (ii) that for $n \ge m$, $\sum_{k=1}^{m} {}^{m}C_{k} {}^{n-1}C_{k-1} = {}^{m+n-1}C_{m-1}$ since if n balls

are arranged in m boxes, they must occupy k non-empty boxes for some k, $1 \le k \le m$, and there are ${}^{m}C_{k}$ ways of selecting k boxes from m. This result follows formally from Lemma 4.1, which with its corollary we will find useful a number of times in this paper.

We can model the original network (Section 3.2) by Z(t) provided we assume that:

- the effective service rate function is known
- all arrangements of customers on n non-empty queues are equally likely, which has been shown to hold exactly when the IN is a crossbar as discussed above.

By analysing process Z(t) in the steady state we can obtain the following theorem:

THEOREM 3.1

Consider a b-way IN with known effective service rate function μ_n in a cyclic network comprising a parallel bank of DMA servers with population N. Under the assumption that all arrangements of customers on n non-empty queues are equally likely, the throughput is given by:-

(3)

$$T(N) = \sum_{n=1}^{b} \mu_n p_n \tag{4}$$

where
$$p_n = \frac{p_1 \mu_1 \prod_{j=1}^{n-1} \{(b-j)(N-j)\}}{\mu_n (n-1)!^2}$$
 (5)
and $p_1^{-1} = \sum_{n=1}^{b} \frac{\mu_1 \prod_{j=1}^{n-1} \{(b-j)(N-j)\}}{\mu_n (n-1)!^2}$ (6)

PROOF

State transitions of the process Z(t) can only occur on completion of service by a task which then recycles to join one of the b queues. Thus there are only two possible types of state transitions: $n \rightarrow n+1$ and $n+1 \rightarrow n$ ($1 \le n < b$) since transitions $\underline{n} \rightarrow \underline{n}'$ of the process Y(t) with $\#(\underline{n})=\#(\underline{n}')$ do not cause a change of state in Z(t). For the first type of transition to be possible, the completing task must depart from a queue in which there is at least one other task (length > 1), and must then join one of the (b-n) empty queues. Let $q_{i,j}$ be the probability that the state changes from i to j on a service completion. Since states (of Y(t)) with the same number of non-empty queues are equi-probable, the transition $n \rightarrow n+1$ occurs with probability $q_{n,n+1}$ given by:-

 $= \frac{(b-n)}{b} Pr(\text{queue length at a given active input} > 1 \text{ when n are active})$

$$= \frac{(b-n)}{b} \frac{\# \text{ arrangements of } N-1 \text{ tasks in } n \text{ non-empty queues}}{\# \text{ arrangements of } N \text{ tasks in } n \text{ non-empty queues}}$$

which by (2) becomes:

$$= \frac{(b-n)}{b} \frac{N-2C_{n-1}}{N-1C_{n-1}} = \frac{(b-n)}{b} \frac{(N-n)}{(N-1)}$$
(7)

Conversely, for the second type of transition, the completing task must leave from a queue in which there are no other tasks and must then join one of the n non-empty queues. The transition $n+1 \rightarrow n$ therefore has probability $q_{n+1,n}$ given by:-

$$q_{n+1,n} = \frac{n}{b} \frac{\# \text{ arrangements of } N-1 \text{ tasks in } n \text{ non-empty queues}}{\# \text{ arrangements of } N \text{ tasks in } n+1 \text{ non-empty queues}}$$

 $q_{n,n+1}$

11

$$\frac{n}{b} \frac{N^{-2}C_{n-1}}{N^{-1}C_{n}} = \frac{n^{2}}{b(N-1)}$$
(8)

Now from (7) and (8), given that the effective service rate function of the switching network is μ_n , we have the following balance equations for the process Z(t) with 1≤n
the:

$$\mu_n p_n (b-n)(N-n) = \mu_{n+1} p_{n+1} n^2$$
(9)

Solving this for p_n gives (5), with normalising constant given by (6). Equation (4) follows from this distribution and the effective service rate function.

When the effective service rate function μ_n is a simple expression, it is not necessary to compute the distribution p_n directly to obtain the throughput T(N) by Theorem 3.1. In such cases an expression for throughput may be obtained by the following theorem:

THEOREM 3.2

=

For a b-way IN with effective service rate function μ_n , in a closed system of a parallel bank of DMA servers with population N, the throughput is given by:-

$$T(N) = \frac{bN}{b+N-1} H(1)$$
(10)

where

$$H(z) = \sum_{n=1}^{\infty} \frac{t_n z^n}{n}$$
(11)

and t_n is given by the following recurrence formula and boundary condition, respectively:-

$$(b-n)(N-n)t_n = n^2 t_{n+1}$$
 (1≤n

$$\sum_{n=1}^{b} t_n / \mu_n = 1$$
(13)

and

PROOF

Substituting $t_n = \mu_n p_n$ in the balance equations (9) and the normalising condition for p_n we obtain the recurrence formula (12) and the boundary condition, (13). We can solve

this by defining the generating function

$$G(z) = \sum_{n=1}^{\infty} t_n z^n \qquad (t_n = 0 \text{ for } n > b)$$

G(z) has first derivative G'(z) = $\sum_{n=1}^{\infty} nt_n z^{n-1}$ and we define H(z) = $\int_{0}^{z} \frac{G(u)}{u} du$

We now have T(N) = G(1) and we can rewrite the recurrence formula (12) as

$$Nb(t_n/n) - (N+b)t_n + nt_n + t_{n+1} - (n+1)t_{n+1} = 0$$

Multiplying by zⁿ and summing yields

$$NbH(z) - (N+b)G(z) + zG'(z) + z^{-1}(G(z)-t_1z) - G'(z) + t_1 = 0$$

which simplifies to $(1-z)G'(z) + (N+b-z^{-1})G(z) - NbH(z) = 0$

Thus, setting z=1, we obtain (10), the required expression for the throughput T(N).

In Theorem 3.2, H(1) is derived from the boundary condition (13) which depends on the specific form of μ_n given by the characteristics of the switching network. In the simplest exact case where we have direct connection (no switching network), $\mu_n=n\mu$ and the

boundary condition (13) becomes $\sum_{n=1}^{b} \frac{t_n}{n\mu} = 1$, i.e. $H(1) = \mu$, giving the well-known 'multiple server' formula for $T(N) = \frac{bN}{b+N-1}\mu$ derived by a number of authors, e.g. [HF86]. There is a similarity between the expression for T(N) when there is a direct connection and the expression for μ_n when there is a full crossbar (given in the next section). In fact, the two expressions can be unified by the substitution, N=n. As we will see shortly, this is to be expected because both expressions can be derived by purely probabilistic arguments which use the fact that each arrangement of tasks over a fixed number of non-empty queues occurs with equal probability in the steady state.

3.4 The Exact Result for Full Crossbars

Rather than assume that the output pin service rate is always the mean value μ_n and solving the simple process Z(t), we may be a little more precise and consider the process X'(t), defined by X'(t) = {#(N(t)),M(t)}, instead of X(t). The state space of X'(t) is then

 $S' = \{(n,m) \mid 1 \le m \le n \le b\}$ and for full crossbars, we may solve for the throughput of the short-circuited subnetwork considered in the previous sections by solving for the equilibrium state space probabilities of S' directly for any population N.

We define the steady-state probability distribution of the stationary Markov process X'(t), P(n,m) for $(n,m) \in S'$, by

$$P(n,m) = \lim_{t \to \infty} Pr(\#(\underline{N}(t))=n, M(t)=m)$$
$$= \sum_{\substack{\#(\underline{n})=n \\ (\underline{n},m) \in \Omega_1}} P(\underline{n},m)$$

where $P(\underline{n},m)$ is the corresponding equilibrium probability for the state (\underline{n},m) of the process X(t).

Now, the number of arrangements of N tasks over n given non-empty queues is $^{N-1}C_{n-1}$ (by (2)) and we know that given $(n,m) \in S'$, each state $(\underline{n},m) \in \Omega_1$ with $\#(\underline{n}) = n$ has the same steady-state probability,

$$\frac{P(\underline{n},m)}{N-1}C_{n-1}C_{n}$$

for a b-way crossbar since all transition probabilities depend only on n and m (recall the argument in the previous section).

In this way we can obtain the following equations for the stationary probabilities, P(n,m), of the Markov process X'(t), for $(n,m) \in S'$:

 $P(n,m)mb(N-1) = (b-n+1)(n-1)(N-n+1) \{P(n-1,m+1)\theta_1(n-1,m+1) + P(n-1,m)\theta_2(n-1,m) + P(n-1,m-1)\theta_3(n-1,m-1) + P(n-1,m-2)\theta_4(n-1,m-2)\} + n^2(n+1)\{P(n+1,m+1)\phi_1(n+1,m+1) + P(n+1,m)\phi_2(n+1,m)\}$

+ $n\{n(N-n) + (n-1)(b-n+1)\}\{P(n,m+1)\phi_1(n,m+1) + P(n,m)\phi_2(n,m) + P(n,m-1)\phi_3(n,m-1)\}$

for appropriately defined parameters $\underline{\theta}$, $\underline{\phi}$, $\underline{\phi}$ given in [HA87] and [PA89]. Notice that in general, after a service completion, the state (n,m) can transit to a state (n',m') with m' = m-1 (recycling task joins a non-empty queue and all queues are blocked or transmitting already through the m-1 servers), m, m+1 or m+2 (another task becomes unblocked by

the service completion, the next task in the completing task's queue enters service at a *different* free server, and the recycling customer enters an empty queue and is not blocked). The state space S' is only of order b^2 and so it is quite feasible to compute performance measures for full crossbars by solving the above equations directly. Moreover, as we shall see in later sections, once computed, these results may be used as base cases in the recursion which yields the corresponding performance measures for delta networks.

In the case of a 2-way crossbar switch, b=2, with N>1, the values of these parameters are $\theta_2(1,1)=\theta_3(1,1)=1/2$, $\phi_1(2,2)=1$, $\phi_2(2,1)=1/2$, $\phi_1(2,2)=1/2$, $\phi_2(1,1)=1$, $\phi_2(2,1)=1/4$, $\phi_2(2,2)=1/2$, $\phi_3(2,1)=1/4$, and the (dependent) equations become:

(N-1)P(1,1)	= 2P(2,2) + P(2,1)
(2N-1)P(2,1)	= (N-1)P(1,1) + 2(2N-3)P(2,2)
2(2N-1)P(2,2)	= (N-1)P(1,1) + (2N-3)P(2,1)

These equations have normalised solution:

$$P(1,1) = \frac{4}{(3N+1)}, \qquad P(2,1) = \frac{2(N-1)}{(3N+1)}, \qquad P(2,2) = \frac{(N-1)}{(3N+1)}$$
(14)

giving a throughput $T(N) = 1.P(1,1) + 1.P(2,1) + 2.P(2,2) = \frac{4N}{(3N+1)}$ (15) Thus throughput is always less than 4/3, and approaches this value as $N \rightarrow \infty$. This is as expected for a full 2-way crossbar with both of its inputs active - see Section 4.1 for example. For finite N, there is always a non-zero probability that an input is inactive, giving lower throughput.

4 FES for DMA Servers with Full Crossbar

4.1 Effective Service Rate Function for Crossbars

For an a×b crossbar, where outputs are randomly selected and inputs are equally utilised, given n active inputs, the probability that m outputs are active, $p_b(m|n)$, is easily determined by ball-in-box arguments since by symmetry every valid arrangement of the n and m (active) input and output pins over b pins is equally likely in the steady state. From this follows the mean number of active outputs, conditional on n active inputs, and hence conditional throughput, μ_n .

Modelling circuit-switched multi-stage interconnection networks

We will require the following lemma which has appeared in [KN68].

LEMMA 4.1

For
$$n \ge 1$$
, $0 \le h \le X-n+1$, $\sum_{k=1}^{n} {}^{n-1}C_{k-1} {}^{X}C_{k+h-1} = {}^{n+X-1}C_{X-h}$

PROOF

For n=1, the left hand side, $h = {}^{X}C_{h} = {}^{X}C_{X-h} = rhs$, the right hand side if $0 \le h \le X$. Now assume inductively that the result is true for $1 \le n \le N$. Then,

$$\sum_{k=1}^{N+1} {}^{N}C_{k-1} {}^{X}C_{k+h-1} = \sum_{k=2}^{N} ({}^{N-1}C_{k-2} + {}^{N-1}C_{k-1}) {}^{X}C_{k+h-1} + {}^{X}C_{N+h} + {}^{X}C_{h}$$
$$= \sum_{k=1}^{N} ({}^{N-1}C_{k-1} {}^{X}C_{k+h} + {}^{N-1}C_{k-1} {}^{X}C_{k+h-1})$$

(we have replaced the dummy summation variable, k, by k+1 in the first sum, and extended each summation domain to include the two loose terms.) Therefore, using the inductive hypothesis twice, for $0 \le h+1 \le X-N+1$ and $0 \le h \le X-N+1$, i.e. for $0 \le h \le X-N$,

$$\ln s = N + X - 1C_{X-h-1} + N + X - 1C_{X-h} = N + X - C_{X-h} = rhs \square$$

COROLLARY 4.2

For $n \ge 1$, $X \ge 1$, $0 \le h \le X-n$, $\sum_{k=1}^{n} {n-1}C_{k-1} X C_{k+h} (k+h) = X X^{k+n-2}C_{X-h-1}$

PROOF

Since ${}^{X}C_{k+h}(k+h) = X \cdot {}^{X-1}C_{k+h-1}$, lhs = X · $\sum_{k=1}^{n} {}^{n-1}C_{k-1} \cdot {}^{X-1}C_{k+h-1}$ = rhs by the lemma.

Now let the random variables M, N denote the numbers of active outputs and inputs respectively. Then,

$$p_b(m \mid n) = Pr(M=m \mid N=n) = \frac{{}^{b}C_m{}^{n-1}C_{m-1}}{{}^{b+n-1}C_{b-1}}$$

Modelling circuit-switched multi-stage interconnection networks

Then by Lemma 4.1, $\sum_{m=1}^{\min(n,b)} p_b(m|n) = 1$ (rearranging as above when b<n), and by

Corollary 4.2 with h=0, the expected number of active outputs when there are n active inputs is given by:-

$$\sum p_{b}(m|n).m = b.(b+n-2C_{b-1} / b+n-1C_{b-1}) = \frac{bn}{(b+n-1)}$$

which, in the terminology of the previous section, gives a degraded service rate function

$$\mu_n = \frac{bn}{(b+n-1)}\mu \tag{16}$$

4.2 Throughput Function for a Crossbar

Now that we have a simple expression for μ_n in the case for a full crossbar, we can use Theorem 3.2 to obtain a simple expression for the throughput T(N).

COROLLARY 4.3

For closed system with population N, consisting of b parallel DMA servers transmitting across a $b \times a$ crossbar network the throughput is given by:-

$$T(N) = \frac{abN\mu}{(a+b-1)N + (a-1)(b-1)}$$

PROOF

For a b×a crossbar network, we obtain the following conditional throughput function from (15), $\mu_n = \frac{an}{a+n-1}\mu$ from which the boundary condition (13) becomes

$$G(1) + (a-1)H(1) = a\mu$$

Substituting T(N) = G(1) and using (10) from Theorem 3.2 we obtain $\left\{\frac{bN}{b+N-1} + (a-1)\right\}H(1) = a\mu$

from which the result follows by substituting H(1) in (10).

In the case of a 2×2 crossbar for example, substituting b=2 and (16) in (5) and simplifying gives $p_1 = 4/(3N+1)$ and $p_2 = 3(N-1)/(3N+1)$ which are the marginal probabilities of (14), the solution obtained by solving the balance equations of the process X'(t) directly (see Section 3.4). Also using these probabilities we derived in Section 3.4, $T(N) = \frac{4N}{3N+1}\mu$ which agrees with the above corollary. Hence our simpler analysis of INs using process Z(t) and effective service rate function μ_n is exact with respect to the more precise analysis in Section 3.4 for a 2×2 crossbar.

In the case of a multi-stage delta network, the expression for μ_n is given by a recurrence formula as derived in the following section, and so a compact expression for throughput would appear impossible to derive. We have obtained our numerical predictions (see Section 7) of throughput T(N) in this case from the degraded rates μ_n by directly computing the set of probabilities p_n (using (5) and (6)) and hence T(N) from (4) in Theorem 3.1.

5 Conditional Throughput Function for a Uniform Delta Network

We analyse delta networks by considering the behaviour of the uppermost rightmost crossbars in each of the subnetworks involved in the recursive definition. This analysis is based upon applications of Little's result, [LI61], and so initially requires no assumptions about the distributions of task service times, only that the crossbar is in stochastic equilibrium. In order to obtain a recursive solution, however, we will make certain approximating assumptions to determine the blocking probability of an arriving task and the probability that it "sees" the other input pin of the crossbar already held by another task. We will see that these assumptions are very mild, especially in the case of a saturated delta network.

First we define some terms, noting that a crossbar connected to servers with appropriate service time distributions can always be described by a stationary stochastic process (i.e. one with a steady state or equilibrium state space distribution).

5.1 Definitions for a 2 by 2 Crossbar

Let the state of a 2×2 crossbar be the binary 4-tuple (I₀, I₁, U₀, U₁) describing the states (active or inactive) of its upper (I₀) and lower (I₁) input pins and output pins U₀, U₁ similarly. In the steady state, we define the following:

- A pin's *mean holding time* (MHT) is the expected elapsed time between the pin becoming active and the departure of the task holding it after its service completion (i.e. the pin next becoming inactive, even if instantaneously)
- An active output pin's *mean residual holding time* (MRHT) is the expected elapsed time between an arrival at the crossbar and the pin next becoming inactive (possibly instantaneously)

• The equilibrium state probabilities are denoted by

 $\pi(i_0 i_1 u_0 u_1) \qquad (i_0, i_1, u_0, u_1 \in \{0, 1\})$

We make the following abbreviations for the marginal probabilities over an (n+1)-dimensional state space (here n=3):

$$\pi(x_0 \dots x_n) = \sum_{\substack{y_i \in \{0,1\} \ (x_i = *) \\ y_i = x_i \text{ otherwise}}} \pi(y_0 \dots y_n)$$

 $\pi(x_0 \dots x_m) = \pi(x_0 \dots x_m * \dots *)$

for 0≤m≤n

(where there are n-m *'s on the right hand side)

- The *blocking probability*, b_i, for input pin i (i=0,1) is the steady state probability that, at the instant a task arrives on that pin, the output pin it requires is already held by a task currently holding the other input pin
- The crossbar is said to possess the *arriving observer property* (AOP) if the steady state probability that an input pin is active at the instant a task arrives on the other input pin is equal to the equilibrium (marginal) probability that the former input pin is active.

5.2 A 2 by 2 Crossbar in Equilibrium

LEMMA 5.1

For a 2×2 crossbar with output pins that

- (a) are selected with the same probability by incoming traffic
- (b) have the same MHT, m, and MRHT, d,

the probability that either output pin is active is

$$p = \frac{m}{2} \left\{ \frac{\pi_0}{m + b_0 d} + \frac{\pi_1}{m + b_1 d} \right\}$$

where π_0 , π_1 are the utilisations of the upper and lower inputs respectively.

PROOF

Let the arrival rates (i.e. reciprocals of the mean inter-arrival times) on the upper and

lower input pins be λ_0 and λ_1 respectively. Three applications of Little's result then yields:

- (i) For an output pin, $p = \left[\frac{\lambda_0 + \lambda_1}{2}\right]m$, since the probability that the pin is active is the same as its mean queue length;
- (ii) For the upper input pin, $\pi(1) = \lambda_0(m+b_0d)$ similarly;
- (iii) For the lower input pin, $\pi(* 1) = \lambda_1(m+b_1d)$ similarly.

Since $\pi(1) = \pi_0$ and $\pi(* 1) = \pi_1$, the result follows by eliminating λ_0 and λ_1 .

COROLLARY 5.2

(a) Assuming the crossbar has the AOP and that m = d

$$p = \frac{\pi_0}{2+\pi_1} + \frac{\pi_1}{2+\pi_0}$$

(b) For a crossbar with $\pi = \pi_0 = \pi_1$, e.g. one in which the input processes are the same for each input pin, which has the AOP and m = d,

$$p = \frac{2\pi}{2+\pi}$$

(c) If further the crossbar is *saturated*, i.e $\pi = 1$,

$$p = \frac{2}{3}$$

PROOF

Since output pins are selected with equal probability, $b_i = \frac{1}{2}\pi_{1-i}$ (i=0,1) by the AOP.

LEMMA 5.3

The steady state probabilities for a single crossbar with output pins selected with equal probability may be written in terms of marginal probabilities as follows:

(a) $\pi(0\ 1\ 0\ 1) = \pi(0\ 1\ 1\ 0) = \frac{1}{2}\pi(0\ 1)$

(b)
$$\pi(1\ 0\ 0\ 1) = \pi(1\ 0\ 1\ 0) = \frac{1}{2}\pi(1\ 0)$$

(c)
$$\pi(1\ 1\ 0\ 1) = \pi(1\ 1\ 1\ 0) = \frac{d'}{2m'+d'}\pi(1\ 1)$$

(d)
$$\pi(1\ 1\ 1\ 1) = \frac{2m'-d'}{2m'+d'}\pi(1\ 1)$$

where m', d' are respectively the MHT and MRHT of each output pin, conditional on both the crossbar's input pins being active.

PROOF

Cases (a) and (b) follow by symmetry. For the other cases, we apply Little's result in the steady state at times when both inputs are active, i.e. the total queue length is 2. Let the throughput of each output pin be τ (the same for each by hypothesis), and the steady state probability that it is active, conditional on both inputs being active, be p'.

Then Little's result applied to an output pin gives $p' = \tau m'$ and to the whole switch $2 = 2\tau \{m' + \frac{1}{2} d'\}$

Thus $p' = \frac{2m'}{2m'+d'}$ and since $1-p' = \frac{\pi(1\ 1\ 0\ 1)}{\pi(1\ 1)}$ the result (c) follows. Result (d) follows from $p' = \frac{\pi(1\ 1\ 1\ 0) + \pi(1\ 1\ 1\ 1)}{\pi(1\ 1)}$

COROLLARY 5.4

If the MHT = MRHT,

 $\pi(1\ 1\ 0\ 1) = \pi(1\ 1\ 1\ 0) = \pi(1\ 1\ 1\ 1) = \frac{1}{3}\pi(1\ 1)$

and the throughput of the crossbar is $\frac{4}{3}$ when both its inputs are active.

5.3 Exact Analysis of the Delta-2 Network

The steady state probability that a given output pin in any stage of a delta network is active satisfies a recurrence relation derived by considering a single crossbar, using Lemma 5.1, and the network's recursive structure.

Let m_s , d_s be the MHT and MRHT respectively for output pins at stage s in a uniform delta network with J stages ($1 \le s \le J$). Thus $m_J = d_J = \mu^{-1}$ if the servers connected to the last stage are exponential with parameter μ . Let b_s be the blocking probability for an arrival on any input pin at stage s. Then we have the following:

PROPOSITION 5.5

The equilibrium probability π_s that an output pin (of any crossbar) in stage s of a J-stage delta network is active satisfies

 $\pi_{s} = \frac{\pi_{s-1}}{1 + b_{s}\alpha_{s}} \qquad (1 \le s \le J)$

where $\alpha_s = \frac{d_s}{m_s}$ and π_0 is the equilibrium probability that an input pin in stage 1 is active.

PROOF

Any crossbar in stage s $(1 \le s \le J)$ satisfies the conditions of Lemma 5.1 since its output pins are stochastically identical and clearly a steady state exists since the Markov process representing the state of *every* pin in the network is ergodic.

5.4 Simplifying Assumptions

In the sequel, apart from the next section, we make two simplifying assumptions that, for all pins in the network:

- (i) MHT = MRHT i.e. that $m_s = d_s$
- (ii) the AOP holds so that $b_s = \frac{1}{2}\pi_{s-1}$

Under these intuitively reasonable assumptions, the recurrence in Proposition 5.5 becomes solvable in closed form for the saturated, uniform case. In general, these assumptions are both approximations in that they do not hold at every stage of the network. This is easily seen in assumption (i) which requires *all* pin holding time distributions to be exponential, although this property does hold in the final stage of a network connected to exponential servers.

If we assume the servers connected to the final stage of the network are exponential, then all holding time distributions are mixtures of Erlang distributions with the same parameter. Thus if arrivals at any crossbar's input pin are *random* with respect to the holding times of the other input pin, $\frac{1}{2} \le \frac{d_s}{m_s} \le 1$ since the mean residual service time for an Erlang-k distribution with parameter μ is $\frac{k+1}{2\mu}$. In fact arrivals are not random and the MRHT is much closer to the bound given by our approximation. This is because the service completion that caused the arrival may have also just unblocked another task and so the new arrival may be blocked at a pin which has just become active with significantly non-zero probability. In this case we have $d_s=m_s$, but we may also have $d_s>m_s$ if an arrival at one pin tends to occur more frequently during long intervals of activity of the other pin. Since the pin holding time distribution seen by an arrival is not known, it is difficult to determine whether $d_s < m_s$ or $d_s > m_s$. However, the assumption that $d_s = m_s$ certainly appears very reasonable.

The more critical assumption is that the AOP holds. Intuitively it should not hold exactly, but it should provide a good approximation. In the next section we prove these two conjectures for a two-stage network with exponential servers. This suggests that the AOP cannot hold everywhere in an arbitrary sized delta network which includes two-stage subnetworks, but this would appear difficult to establish for arbitrary holding time distributions, and not an important issue in view of the other approximating assumption.

In Section 5.5 we compare the approximate method with the exact solution for a twostage saturated, uniform network and show that the AOP does not hold by deriving explicitly the steady state probability that an input pin is active at the instant a task arrives on the other input pin of a crossbar in the final stage. The approximating assumptions do turn out to be very mild and yield extremely accurate results when compared with both the exact solution in the two stage case and with simulations of larger networks.

5.5 Analysis of Arriving Observer

In this section we consider a two-stage delta-2 network connected to four exponential servers. In the steady state, the probability that a crossbar's input pin is active at the instant a task arrives on the other input pin can be determined by first finding the proportion of time in the long term that the system spends in states, s' say, that *can* provide an arrival to the upper (say) input pin, via a transition into an *input state*, s say. By the Key Renewal Theorem, see [CI75] for example, the unnormalised equilibrium probability that the system is in state s immediately after the arrival instant is the expected number of transitions into state s in unit time, which can be determined from the equilibrium probabilities of the states s' and the transition rates from s' to s. This is the approach taken by [SM81] which considers the distribution of similar input (and output) states in product-form queueing networks.

5.5.1 Notation

Let $S = \{\underline{b} \mid b_j=0,1; 0 \le j \le 3\}$ be the set of binary quadruples describing the (marginal) states of the output pins (numbered j from the top) in the *first* stage of the network. An active pin is denoted by 1 and an inactive one by 0.

Let $A = \{s \in S \mid s_0=1\}$ be the set of states that can exist immediately after an arrival at the

top input pin of the second stage.

For $s \in A$, let $B(s) = \{s \in S \mid \exists a \text{ one-step transition } s' \rightarrow s\}$

Let T(s',s) denote the server at which a service completion can cause the transition $s' \rightarrow s$

Since we are considering a renewal process, the expected number of transitions into state $s \in A$ in unit time is

$$\Phi_{s} = \sum_{s' \in B(s)} \pi(s') p_{s's}$$

where we assume without loss of generality that the servers have rate 1 and

- $\pi(x)$ is the equilibrium probability for state $x \in S$, i.e. the proportion of time in the long term that the system spends in state x
- $p_{s's}$ is the probability that a transition from state s' enters state s, conditional on that transition being the result of a service completion at server T(s',s)

5.5.2 State seen by an Arriving Observer

The steady state probability that an arrival at the top input pin of the second stage finds the other input pin of the same (top) crossbar active is therefore $\frac{\Phi^*}{\Phi}$ where

$$\Phi^* = \sum_{\substack{s \in A \\ s_2 = 1}} \Phi_s \quad \text{and} \quad \Phi = \sum_{s \in A} \Phi_s$$

By calculating Φ_s individually for every $s \in A$, we can indeed determine $\frac{\Phi^*}{\Phi}$ but the

method is extremely laborious and we can do better by considering departures rather than arrivals and exploiting the symmetry of the model as follows.

 Φ^* is the long term arrival rate on input pin 0 of the upper crossbar in the second stage when the other input pin of that crossbar is active. By symmetry, this is a half of the *total* long term rate of arrivals to the crossbar that result in a state with both input pins active. In the steady state, this is the same as half the departure rate (i.e. throughput) from the crossbar's output pins in states with both its input pins active. Thus by Corollary 5.4, we have

$$\Phi^* = \frac{1}{2} \frac{4}{3} \pi (1 * 1 *)$$

Similarly, $\Phi^0 \triangleq \Phi - \Phi^*$ is the total long term rate of arrivals to the same crossbar that result in only the top input pin being active. In the steady state this is equal to the departure rate from the state with the upper input pin active and the lower one inactive. Thus we have

 $\Phi^0 = \pi(1 * 0 *) = \pi(1 * 0 1)$ for a saturated network

The ratio $\frac{\Phi^*}{\Phi}$ can therefore be computed using the following equilibrium marginal probabilities (which use the symmetry of this network):

$$p \triangleq \pi(1 \ 0 \ 1 \ 0) = \pi(0 \ 1 \ 0 \ 1)$$

$$p' \triangleq \pi(1 \ 0 \ 0 \ 1) = \pi(0 \ 1 \ 1 \ 0)$$

$$q \triangleq \pi(1 \ 1 \ 1 \ 0) = \pi(1 \ 1 \ 0 \ 1) = \pi(1 \ 0 \ 1 \ 1) = \pi(0 \ 1 \ 1 \ 1)$$

$$r \triangleq \pi(1 \ 1 \ 1 \ 1) = 1 - 2p - 2p' - 4q$$

This gives

$$\frac{\Phi^*}{\Phi} = \frac{\frac{2}{3}\pi(1*1*)}{\frac{2}{3}\pi(1*1*) + \pi(1*01)}$$
$$= \frac{2(p+2q+r)}{2p+3p'+7q+2r}$$

Now, if the AOP held, we would have $\frac{\Phi^*}{\Phi} = \pi(**1*) = p + p' + 3q + r$ which we cannot establish since there are no more independent identities. In fact, exact analysis of the Markov process (see [PA89]) yields $\frac{\Phi^*}{\Phi} = \frac{1454}{2179} = 0.667279$ to 6 decimal places. If we assume that MHT = MRHT for the first stage as well as the second, as in the approximate analysis for delta networks, this ratio is $\frac{2}{3}$. We might therefore expect that our approximation will be very good, and this is borne out in section 7. In fact for the two stage case, the exact throughput is $\frac{17432}{8719} = 1.999312$ to 6 d.p. compared with the approximate result of 2.

5.6 Recursive Analysis of the Delta Network

In this section we derive the throughput of a delta network with uniform routing, conditional on the number of active inputs to the network. For an s-stage delta-2 network (s≥1) we define the set $V_s = \{(n_0,...,n_{2^{s-1}}) \mid n_i \in \{0,1\}, 0 \le i < 2^s\}$ and the following random variables:

- (i) <u>N</u>∈ V_s which represents the state of the input pins to the network (N_i=0 means that input pin i is inactive and N_i=1 means that it is active, i.e. the DMA server connected to it is wishing to transmit).
- (ii) $\underline{Z} \in V_s$ which represents the state of each network output pin.

We also define the function # for which $\#(\underline{V})$ is the number of non-zero components in the vector \underline{V} and the analogous function $\#_1$ for which $\#_1(\underline{V})$ is the number of non-zero components in the *first half* of the vector \underline{V} .

We will require the following probability distribution: $Q_s(i|n) = Pr(\#_1(\underline{N})=i | \#(\underline{N})=n)$ which is the probability that i active inputs are in the upper part of the s-stage network (i.e. the first 2^{s-1} pins) given that there are n active inputs altogether. By symmetry, when there are n active inputs to the network, they are uniformly distributed over the 2^s input pins. Since all arrangements of the n active inputs are equally likely:-

$$Q_{s}(i \mid n) = \frac{kC_{i}kC_{n-i}}{2kC_{n}}$$
 where $k = 2^{s-1}$ (17)

We also define the following probability distribution:

$$T_{s}(n) = Pr(Z_{0}=1 | \#(\underline{N})=n)$$

which represents the probability that the topmost output pin of the s-stage network is active given that there are n active inputs to the network. This distribution can be defined by using the recursive structure of the delta network with the partial shuffle topology (Section 2).



Figure 5 Recursive structure of s-stage delta network

Consider the rightmost, topmost switch in an s-stage network. This takes its inputs from two separate (s-1)-stage subnetworks (Figure 5). When we condition on the number of active inputs to each subnetwork (their sum being #(N)=n), the output pin utilisation of either network can be found. Hence, we can determine the utilisation of both the inputs to the rightmost, topmost crossbar by using the result obtained for the next smaller networks in the recursion. The distribution of the state of the top output pin then follows by considering a single crossbar and the results of the previous sections. This gives the recurrence formula in the following theorem:

THEOREM 5.6

Under the assumptions that every crossbar in the network has the AOP and outputs with MRHT = MHT, the output pin utilisation of a s-stage delta-2 network conditional on the number of active inputs being n is given by the following:-

For s>1,
$$T_s(n) = \sum_{i=\max(0,n-2^{s-1})}^{\min(n,2^{s-1})} Q_s(i|n) U(T_{s-1}(i), T_{s-1}(n-i))$$
 (18)

$$T_1(0) = 0, \quad T_1(1) = U(0,1), \quad T_1(2) = U(1,1)$$
 (19)

where

$$U(\pi_0,\pi_1) = \frac{\pi_0}{2+\pi_1} + \frac{\pi_1}{2+\pi_0}$$

Modelling circuit-switched multi-stage interconnection networks

PROOF

For an s-stage delta network (s>1) we define the following random variables (see Figure 5):-

 $Z \in \{0,1\}$ the state of output pin 0 of an s-stage network

 $I_0 \in \{0,1,2,..2^{s-1}\}$ the number of active inputs to the upper (s-1)-stage subnetwork $I_1 \in \{0,1,2,..2^{s-1}\}$ the number of active inputs to the lower (s-1)-stage subnetwork

$$Pr(Z=1, I_0+I_1=n) = \sum_{i=0}^{n} Pr(Z=1, I_0=i, I_1=n-i)$$

=
$$\sum_{i=0}^{n} Pr(I_0=i, I_1=n-i) Pr(Z=1 | I_0=i, I_1=n-i)$$

$$Pr(Z=1 | I_0+I_1=n) = \sum_{i=0}^{n} Pr(I_0=i, I_1=n-i | I_0+I_1=n) Pr(Z=1 | I_0=i, I_1=n-i)$$

The recurrence (18) follows from the following substitutions:-

$$\begin{split} T_{s}(n) &= \Pr(Z=1 \mid I_{0}+I_{1}=n) \\ Q_{s}(i \mid n) &= \Pr(I_{0}=i, \ I_{1}=n-i \mid I_{0}+I_{1}=n) \\ U(T_{s-1}(i), \ T_{s-1}(n-i)) &= \Pr(Z=1 \mid I_{0}=i, \ I_{1}=n-i) \end{split}$$

where $U(\pi_0,\pi_1)$ (the crossbar output pin utilisation when the crossbar inputs have utilisations π_0 and π_1) is p in Corollary 5.2.

COROLLARY 5.7

The expected number of active output pins (m) conditional on n active inputs to a J-stage network is given by:-

$$E(m|n) = 2^{J} T_{J}(n)$$
⁽²⁰⁾

From the theorem and its corollary, the required conditional throughput function $\mu_n = \mu E(m|n)$ can be used to solve the birth and death process defined in Section 3.3.

5.7 Saturated Delta Network

The delta network becomes saturated when the queue lengths at each of its inputs grows very large. This occurs as the number of customers in the closed system of DMA servers and delta network approaches infinity.

COROLLARY 5.8

For a saturated delta-2 network, p_s, the probability that the top pin in stage s is active is

Modelling circuit-switched multi-stage interconnection networks

given by the following:-

$$p_{s} = \frac{2}{s+2} \tag{21}$$

and the expected number of active outputs in a J-stage network is given by:-

$$E(m \mid 2^{J}) = \frac{2^{J+1}}{I+2}$$
(22)

PROOF

Since all input pins to the network are active the recurrence of Section 5.1 simplifies considerably in the absence of splitting probabilities $Q_s(iln)$ since $Q_s(2^{s-1} | 2^s)=1$ (from (17))

Consider an s-stage delta network with all 2^{s} inputs active. The (s-1)-stage subnetworks are also saturated (see Figure 5) and so the recurrence (18) in Theorem 5.6 simplifies to give the following:-

$$p_{s} = \frac{2p_{s-1}}{2+p_{s-1}}$$
 where $p_{s} = T_{s}(2^{s})$

This recurrence has solution (21), and (22) follows from Corollary 5.7.

For a 2-stage, saturated delta-2 network with μ =1, we find $p_2 = 1/2$ using (21). Thus, the throughput of the saturated network is 2. This compares with a figure of 4.4/(4+4-1) = $2^2/7$ given by equation (16) for a full crossbar network.

Although the recurrence (18) has only been simplified under these rather restricted conditions, these are of course the conditions of greatest interest, viz high loading. In the following section, we derive similar recurrence formulae for the case when the output pins are non-uniformly selected, which allows us to investigate the effect of hot-spots on performance.

6 Conditional Throughput Function for Non-uniform Delta Network

So far we have assumed that all outputs from the MIN are equally utilised, but the use of memory elements, say, on the output side may not be uniform for all elements. For example, if a shared lock or a frequently-accessed part of a data structure were residing in a particular memory element then requests to that memory address would be more frequent than to other destination addresses. We can characterise this by assigning a routing probability to each memory for all requests. Furthermore, each processing element (and so input pin) may have a different memory access pattern, i.e. it may have a favoured set of memories which it frequently accesses. However, in this section, we

Modelling circuit-switched multi-stage interconnection networks

only consider the simplest case in which there is a single hot-spot, viz. the top output pin, and all other memories are equally utilised. With more computational effort, more general hot-spot contention can be modelled by a direct extension of the approach presented here.

The basis of the analysis of a non-uniform delta network is the same as that of a uniform one, viz. using Theorem 3.1 to determine the throughput function for the FES given some effective service rate function μ_n . Hence we derive a recursive formula to determine μ_n for the case where there is a hot-spot at the top output pin. This entails the use of reachability properties for any task at an arbitrary switch in the network attempting to reach a particular output pin. We first show that paths to output pins of the same class (defined in Section 2) have the same steady state probability of being blocked to an arrival, i.e. have the same degree of contention. We can then obtain recurrence formulae for the throughput of each output pin class in much the same way as we did for the uniform case in Theorem 5.6. This result is parameterised in terms of unknown pin mean holding time ratios, and these are finally determined iteratively using the fact that in the network, an output pin's throughput must be proportional to its selection probability (which is the same for all but the top pin).

6.1 Pin Class Reachability and Blocking Probability

Before analysing hot-spot contention in delta networks, we first prove some reachability properties for tasks in a network in terms of pin classes, i.e. to what classes of network output pins can a task connect from a given class of pin within the network. We then consider the probability that a path to an output pin of a given class is blocked.

LEMMA 6.1

For any s-stage delta-2 network with a partial shuffle topology, only class k+1 ($1 \le k < s$) output pins are reachable from class k (and only from class k) output pins in the (s-1)-stage subnetworks, and only the class 0 and class 1 pins are reachable from the class 0 (and only from class 0) pins in the subnetworks.

PROOF

The class 0 pins of the (s-1)-stage subnetworks connect to the top switch in stage s, so only class 0 and class 1 pins of the s-stage network are reachable from the class 0 (and only from class 0) pins of the subnetworks.

Now consider the 2^{k-1} class k pins labelled i in the range $2^{k-1} \le i \le 2^{k}-1$ ($1 \le k < s$) in the (s-1)-stage subnetworks. Since switch i in stage s takes its inputs from output pin i in stage s-1 and has outputs labelled 2i and 2i+1, only pin j in the range $2(2^{k-1}) \le j \le 2(2^{k}-1)-1$, i.e. $2^k \le j \le 2^{k+1}-1$ is reachable from (and only from) subnetwork pin of class k. In other words, only class k+1 pins in stage s are reachable from class k (and only from class k) pins in the subnetworks, where $1 \le k < s$. \Box

PROPOSITION 6.2

For any s-stage subnetwork $(0 < s \le J)$ in the recursive definition of a J-stage delta-2 network with a partial shuffle topology, only network output pins up to and including class J-s are reachable from the class 0 (and only from class 0) output pins of the subnetworks. Likewise, only network output pins of class J-s+k ($0 < k \le s$) are reachable from class k (and only from class k) output pins of the subnetworks.

PROOF

The proof is by induction on t=J-s for s=J,J-1,..,1. For the base case (t=0), the proposition is trivially true.

For the case t>0, let us assume the proposition holds for J-s=t-1. By Lemma 6.1, only stage s output pins of class 0 and 1 can be reached from class 0 (and only from class 0) pins in the (s-1)-stage subnetworks, so only network output pins up to and including class t are reachable from class 0 (and only from class 0) pins of the subnetworks (by the inductive hypothesis).

Also by Lemma 6.1, only stage s output pins of class j+1 can be reached from class j (and only from class j) pins ($1 \le j < s$) in the (s-1)-stage subnetworks, so only network output pins of class t+j are reachable from class j (and only from class j) pins of the subnetworks (by the inductive hypothesis).

Hence, the proposition holds for J-s=t and this completes the proof. \Box

For each inter-stage link there is an associated steady state probability that an arriving task wishing to use that link will be blocked: this is called the *link blocking probability*. Each path contains a set of links across the network and has a similar blocking probability. Two paths have the same blocking probability if all links between corresponding stages have the same blocking probability. Now suppose that when a task arrives at a particular switching node it selects either output pin with equal probability because the routing probabilities to all possible destination addresses from this node are the same; then the link blocking probability for both output links will be identical. By applying this observation to appropriate nodes in the network, we can show that paths to certain output

Modelling circuit-switched multi-stage interconnection networks

pins have the same blocking probability despite the presence of a hot-spot at the top destination address. This is formalised by the following proposition.

PROPOSITION 6.3

For an s-stage delta-2 network with a partial shuffle topology in which all output pins but the topmost one are uniformly utilised, the path blocking probability is the same for all paths to pins of the same class (i.e. paths to 2^k class k+1 pins ($0 \le k \le s$) have the same blocking probability).

PROOF

The base case of our inductive proof (s=1) holds trivially because there is only one pin of class 0 and one of class 1. For the multi-stage case (s>1), let us assume that the result holds for an (s-1)-stage delta network. Now since both upper and lower subnetworks are identical by symmetry and customers are uniformly distributed over the network input pins, the path blocking probabilities for partial paths to a pin of class k ($1 \le k \le s-1$) in both subnetworks are the same (by the inductive hypothesis). By Lemma 6.1, we know that pins of class k>0 in the subnetworks connect only to pins of class k+1 in stage s. Since the hot-spot is not reachable from these stage s pins, their link blocking probabilities are the same (by class k +1 in stage s. Since the same. Hence path blocking probabilities are the same for all paths to pins of class k+1 ($1 \le k \le s-1$ or $2 \le k+1 \le s$).

Finally, we note that when k=0, the proposition is satisfied vacuously which completes the proof. \Box

When there is a single hot-spot at the top output pin of a delta network, paths to output pins of the same class experience the same degree of contention. If different classes of network output pins are reachable from a given switch in the network then paths from the upper switch output may have different blocking probabilities to paths from the lower switch output. Consequently, when the upper and lower switch outputs become active, the mean time taken for them to become inactive, their mean holding time, which we will also call their *release time*, will be different. Hence modelling the presence of a hot-spot differs from the uniform routing case in two important aspects: the different routing probabilities and the different release times for outputs of a given switch.

6.2 Hot-spot Traffic Model

The hot-spot traffic model is characterised by the hot-spot routing probability ρ which is the probability that a new arrival selects the top network output pin as its destination.

Each of the other network outputs are chosen with probability $(1-\rho)/(2^{J}-1)$ in a J-stage network.

In our analysis of hot-spots, we will only need to consider the nodes in the decode tree of the top network input pin because all input pins are identical. Firstly, consider an arrival to the topmost switching node in stage s of a J-stage network. We can determine the probability that the task selects the upper switch output pin by summing the routing probabilities to the network pins which are reachable from the upper pin. It has been shown in Proposition 6.2 that, from the upper output, network pin number 0 is reachable as are all pins up to and including class t, and from the lower output, only pins of class t+1 are reachable, where t=J-s. Since there are 2^{k-1} pins of class k, the switch routing probability $\omega(s)$ for the top switch in stage s is given by :-

$$\omega(s) = \frac{\rho + \sum_{k=1}^{t} 2^{k-1} q}{\rho + \sum_{k=1}^{t+1} 2^{k-1} q} = \frac{\rho + (2^{t}-1) q}{\rho + (2^{t}-1) q}$$
(23)

where $q = (1-\rho)/(2^{J}-1)$.

The expression for $\omega(s)$ is applicable to the top row of switches across the network. For the other switches in the decode tree of the top network input, the switch output pins are of identical class and so, by Proposition 6.2, only network output pins of identical class are reachable. Hence, the routing probabilities to the switch outputs are the same and so the value of ρ is a half.

6.3 Non-uniform 2 by 2 Crossbar

Since all network inputs are treated identically, we only need to consider a single representative pin (say the topmost input). We base our recursive analysis, to determine the conditional throughput function, on the decode tree of the top input pin. Hence, we will require knowledge of the properties of the crossbars in this decode tree only. In particular, the topmost switches in each stage of this decode tree will be non-uniform whereas other switches will be uniform because their output pins have the same class; the hot-spot is not reachable from these switches. As with the uniform routing case, we wish to determine the probability that a particular switch output is active given the utilisation of the inputs to the switch. However, in this case, not only do we have to consider the non-uniform routing but also the fact that when the upper and lower switch

outputs become active, their release times are different. In fact, in our analysis of the crossbar, we only need to know the *ratio* of the release times for the upper and lower outputs.

6.3.1 Definitions

The following definitions for the non-uniform crossbar are analogous or extensions to those for the uniform crossbar considered in Section 5.1.

- The equilibrium probabilities π_0 , π_1 are abbreviations of $\pi(1)$, $\pi(* 1)$ respectively (as in Section 5.1).
- The release time ratio r is the ratio of the MHT of the lower output and the MHT of the upper output.
- The blocking probability b_{ij} for input pin i and output pin j (i,j∈ {0,1}) is the steady state probability that, at the instant a task arrives on input pin i, the output pin it requires (pin j) is already held by a task currently holding the other input pin.

6.3.2 A Non-uniform Crossbar in Equilibrium

LEMMA 6.4

For a non-uniform 2×2 crossbar with output pins 0 and 1 that have

- selection probability ρ and $(1-\rho)$
- MHT m and rm, where r is the release time ratio
- MRHT do and d1

the probability that the upper output is active is given by:

$$p = \rho m \left[\frac{\pi_0}{\rho(m+b_{00}d_0) + (1-\rho)(rm+b_{01}d_1)} + \frac{\pi_1}{\rho(m+b_{10}d_0) + (1-\rho)(rm+b_{11}d_1)} \right]$$

and the probability that the lower output is active is given by:

$$q = \frac{(1-\rho)r p}{\rho}$$

where π_i is the utilisation of input pin i (i=0,1)

PROOF

Let the arrival rate on input pin i be λ_i (i=0,1). Four applications of Little's result then yields

(i) For input pin 0: $\pi_0 = \lambda_0 \left[\rho(m + b_{00}d_0) + (1-\rho)(m + b_{01}d_1) \right]$

(ii) For input pin 1: $\pi_1 = \lambda_1 \left[\rho(m+b_{10}d_0) + (1-\rho)(rm+b_{11}d_1) \right]$

(iii) For output pin 0: $p = \rho(\lambda_0 + \lambda_1)m$

(iv) For output pin 1: $q = (1-\rho)(\lambda_0 + \lambda_1)rm$

The result follows by eliminating λ_0 and λ_1 .

COROLLARY 6.5

Assuming that the non-uniform crossbar has the AOP and that the MRHT = MHT for both output pins,

$$p = \rho(\rho + (1-\rho)r) \left[\frac{\pi_0}{G(\pi_1)} + \frac{\pi_1}{G(\pi_0)} \right]$$
 and $q = \frac{(1-\rho)r p}{\rho}$

where $G(\pi_i) = (1+\pi_i)(\rho^2 + (1-\rho)^2 r^2) + 2\rho(1-\rho)r$

PROOF

Conditional on only the lower input being active, let $P(\uparrow)$ be the probability that the task is holding the upper output and let $P(\downarrow)$ be the probability that it is holding the lower output. Applying Little's result to the crossbar at times when only its lower input is active is equivalent to using Lemma 6.4 with $\pi_0=0$ and $\pi_1=1$ (so that $b_{10}=b_{11}=0$). This gives:

$$p' = \frac{\rho}{\rho + (1-\rho)r}$$
 and $q' = \frac{(1-\rho)r}{\rho + (1-\rho)r}$

where p' and q' are the utilisations of the upper and lower output pins when only the lower input pin is active. Also we have:

$$p' = \frac{\pi(0 \ 1 \ 1 \ 0)}{\pi(0 \ 1)} \quad \text{and} \quad q' = \frac{\pi(0 \ 1 \ 0 \ 1)}{\pi(0 \ 1)}$$
$$o \quad P(\uparrow) = p' = \frac{\rho}{\rho + (1 - \rho)r} \quad \text{and} \quad P(\downarrow) = q' = \frac{(1 - \rho)r}{\rho + (1 - \rho)r}$$

and so

Given that the AOP holds and that one input pin is active (by symmetry), a new arrival on the other input sees the active input connected to the upper output with probability $P(\uparrow)$ and to the lower output with probability $P(\downarrow)$. Thus we have the following blocking probabilities:

$b_{00} = \pi_1 P(\uparrow)$	and	$b_{01} = \pi_1 P(\downarrow),$
$b_{10} = \pi_0 P(\uparrow)$	and	$b_{11} = \pi_0 P(\downarrow)$

The result follows by substituting these blocking probabilities and $d_0=m$, $d_1=rm$ in Lemma 6.4.

Let $U_k(\pi_0,\pi_1,\rho,r)$ be the utilisation of output pin k of a crossbar with input utilisations π_0 and π_1 ; upper pin routing probability ρ ; and release time ratio r. Thus $U_0(\pi_0,\pi_1,\rho,r)$

=p and $U_1(\pi_0,\pi_1,\rho,r) = q$ in Corollary 6.5.

6.4 Recursive Analysis of a Delta Network with a Hot-spot

In this sub-section, we derive the throughput of a delta network with non-uniform routing, conditional on the number of active inputs to the network. Although the approach presented here is applicable to full non-uniform routing in which all output pins are selected with different probabilities, for simplicity, we only consider the problem of a single hot-spot which is selected by a newly arriving customer with probability ρ , and all other output pins are uniformly utilised.

Suppose the top pin is the hot-spot and so $\rho > 1/b$ for a b-way MIN. Then all nodes in the decode tree of the top output pin will be busier than in the uniformly utilised case, so we refer to it as the hot decode tree. The degree of overlap with another output pin's decode tree will determine the extent of the extra contention caused by the hot-spot. For example, customers destined for pin number one, adjacent to the hot pin, suffer the greatest contention because pin one's decode tree nodes are identical to the hot-pin's decode tree nodes. However, the customers destined for the bottom output pin suffer less hot-spot contention because its decode tree nodes only overlap the hot decode tree at the leaves i.e. the input pins of the first stage. Consequently, in equilibrium, a greater proportion of the customers will build up for pin number one than for the bottom pin. It is easy to see, by an argument analogous to the one used to prove Proposition 6.3, that certain output pins have the same build up of customers and that these pins can be grouped together to form the classes defined earlier in Section 2.

Given the <u>N</u> and <u>Z</u> defined in Section 5.1, we now define the following probability distributions:-

 $T_{s}^{(0)}(n) = Pr(Z_{0}=1 | \#(\underline{N})=n)$

and for k>0,

 $T_{s}^{(k)}(n) = Pr(Z_{2k-1}=1 | \#(\underline{N})=n)$

 $T_s^{(k)}(n)$ is the probability that an output pin of class k is active in an s-stage network with n active inputs and typically we use pin number 2^{k-1} as the representative pin for all class k pins.

In this recursive analysis, we consider the decode tree of the top input pin of the network. For the base case (s=1) of the recursion, the interconnection is a 2-way crossbar with

switch routing probability for the top switch output given by $\omega(1)$. For the s-stage case (s>1), the upper and lower (s-1)-stage subnetworks are separate so that we can determine utilisations of the inputs to the rightmost crossbars by using the result obtained for the smaller network in the recursion. This gives the recurrence formulae in the following theorem.

THEOREM 6.6

Suppose there are n_J active inputs to a J-stage delta network with a partial shuffle topology. Assuming that every crossbar in the network has the AOP and outputs with MRHT = MHT, the utilisation of a class k output pin in stage s is given by:- for $1 \le J$, $k \in \{0,1\}$, $1 \le n \le n_J$,

$$T_{s}^{(k)}(n) = \sum_{i=\max(0,n-2^{s-1})}^{\min(n,2^{s-1})} Q_{s}(i|n) U_{k}(T_{s-1}^{(0)}(i), T_{s-1}^{(0)}(n-i), \omega(s), r(s,n_{J}))$$
(24)

for $1 \le J$, $1 \le k \le s$, $1 \le n \le n_J$,

$$T_{s}^{(k)}(n) = \sum_{i=\max(0,n-2^{s-1})}^{\min(n,2^{s-1})} U_{0}(T_{s-1}^{(k-1)}(i), T_{s-1}^{(k-1)}(n-i), 1/2, 1)$$
(25)

and for
$$k \in \{0,1\}$$
, we have:-

$$T_1^{(k)}(0) = 0$$

$$T_1^{(k)}(1) = U_k(0, 1, \omega(1), r(1,n_J))$$

$$T_1^{(k)}(2) = U_k(1, 1, \omega(1), r(1,n_J))$$
(26)

where $\omega(s)$ is the switch routing probability and $r(s,n_J)$ is the release time ratio for the top switch in stage s.

PROOF

The proof is an analogous to that of Theorem 5.6. In addition it uses Lemma 6.1 for (24) and (25), i.e. a pin of class k in stage s is reachable from a pin of class k-1 only.

COROLLARY 6.7

The expected number of active output pins (m) conditional on nJ active inputs to a J-stage network is given by:-

$$E(mln_{J}) = T_{J}^{(0)}(n_{J}) + (2^{J}-1) T_{J}^{(1)}(n_{J})$$
(27)

From the theorem and its corollary with $n_J = 1, 2, ... 2^J$, the required conditional throughput function $\mu_n = \mu E(m|n)$ can be used to solve the birth and death process

defined in section 3.3. This gives T(N), the throughput of the closed system for a given population N. However, the release time ratios $r(s,n_J)$ are unknown in Theorem 6.6. It now remains to obtain these ratios and this is the subject of the next section.

6.5 Fixed-point Equation for Release Time Ratios

First we define the following notation:

- Suppose there are nJ active inputs to a J-stage delta network, so that the release time ratios r(s,nJ) can be abbreviated to r_s for the topmost crossbar in stage s (s=1,2,...,J), where $r_s>0$.
- Let <u>t</u> describe the utilisations of the network output pins, i.e. t_k is the utilisation of a class k pin (k=0,1,...,J). From Theorem 6.6, we have $t_k = T_J^{(k)}(n_J)$.
- Let ρ ($0 \le \rho < 1$) be the hot-spot routing probability that describes the network traffic and ρ_k be the routing probability to a class k output pin, so that

$$\rho_0 = \rho$$
 and $\rho_k = \frac{1-\rho}{2^J-1}$ (k = 1,2,...,J)

- Let $\omega(s)$, which we abbreviate to ω_s , be the corresponding top switch routing probability function, where $0 \le \omega_s < 1$ (s=1,2,...,J).
- Let $\underline{\mathbf{r}} = (\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_J), \underline{\mathbf{t}} = (\mathbf{t}_0, \mathbf{t}_1, \dots, \mathbf{t}_J), \underline{\boldsymbol{\rho}} = (\rho_0, \rho_1, \dots, \rho_J)$ and $\underline{\boldsymbol{\omega}} = (\omega_1, \omega_2, \dots, \omega_J)$

We now derive a fixed-point equation for the pair $(\underline{\omega},\underline{r})$ using the functions A, B and C defined below. Each function takes a pair of vectors and returns a pair of vectors in which the second component (\underline{r}) is the same in both pairs. However, in function A, the \underline{r} influences the value of the first vector in the result.

- Firstly, from Theorem 6.6, the utilisation vector is a function of the switch routing probabilities and the release time ratios. Suppose F is this function and define $(\underline{t},\underline{r}) = A(\underline{\omega},\underline{r}) = (F(\underline{\omega},\underline{r}),\underline{r})$.
- Now <u>t</u> induces a set of network routing probabilities since the proportion of the total throughput flowing out of a pin in equilibrium must be equal to its selection probability,

i.e.

$$\rho_k = \frac{t_k}{t_0 + \sum_{j=1}^J 2^{j-1} t_j} = [G(\underline{t})]_k, \text{ say.}$$

Now let $(\underline{\rho},\underline{\mathbf{r}}) = \mathbf{B}(\underline{\mathbf{t}},\underline{\mathbf{r}}) = (\mathbf{G}(\underline{\mathbf{t}}),\underline{\mathbf{r}}).$

The switch routing probability is a function of the network routing probability, and given by a more general form of (23):

$$\omega_{s} = \frac{\rho_{0} + \sum_{k=1}^{J-s} 2^{k-1} \rho_{k}}{\rho_{0} + \sum_{k=1}^{J-s+1} 2^{k-1} \rho_{k}} = [H(\underline{\rho})]_{s}, \text{ say.}$$

Now let $(\underline{\omega},\underline{\mathbf{r}}) = \mathbf{C}(\underline{\rho},\underline{\mathbf{r}}) = (\mathbf{H}(\underline{\rho}),\underline{\mathbf{r}}).$

Thus we have the following fixed-point equation: $(\underline{\omega},\underline{r}) = C(B(A(\underline{\omega},\underline{r})))$. Although for all $(\underline{\omega},\underline{r})$, $C(B(A(\underline{\omega},\underline{r}))) = (\underline{\omega}',\underline{r})$ for some $\underline{\omega}'$, the choice of \underline{r} determines $\underline{\omega}'$ and hence the fixed-point pair.

We require fixed-points ($\underline{\omega},\underline{r}$) of the above equation. Typically, $\underline{\omega}$ is given by the traffic pattern and a guess is required for \underline{r} . From the analysis, we know that $r_J=1$ because complete paths suffer no further contention and DMA transfer rates are the same regardless of the output pin. The proof of uniqueness of \underline{r} with respect to the fixed-point equation is not given here but we conjecture that \underline{r} is unique because in the 2-stage network different release time ratios in the first stage crossbar give output pin utilisations that induce different switch routing probabilities. This means that both cannot satisfy the fixed-point equation.

6.6 Iterative Method

Let $r_s^{(i)}$ be the approximation for r_s after the ith iteration. As an initial guess, we choose the following: $r_s^{(0)}=1$ for s=1,2,...,J. This guess is exact when $n_J=1$ and also for uniform traffic. Furthermore, in all approximations to <u>r</u> we know that $r_J=1$.

The induced switch routing probabilities before the ith iteration are given by the following:

$$\underline{\omega}^{(i)} = \text{first}(C(B(A(\underline{\omega}, \mathbf{r}^{(i)})))) \quad \text{where } \text{first}((p,q)) = p$$

Modelling circuit-switched multi-stage interconnection networks

The difference between the required and the observed switch routing probabilities can be quantified as follows:

$$d_{s}(i) = \frac{\omega_{s}(i) - \omega_{s}}{\omega_{s}}$$

Now if $d_s^{(i)}$ is positive, then the utilisation of the upper pin is higher than expected and so the release time ratio for the top switch in stage s needs to be increased. Thus we use the following to update the value of r:

$$r_s^{(i+1)} = r_s^{(i)} [1 + Dd_s^{(i)}]$$

where D (typically between 1 and 4) is a damping factor used to reduce the number of iterations. The stopping condition is $|d_s^{(i)}| < \varepsilon$, s=1,2,...,J for some arbitrarily small ε .

7 Validation and Numerical Results

7.1 Validation of Analytical Model

The model is validated against a simulation model of the closed system of parallel DMA servers (Section 3) connected to a circuit-switched delta network in which partial paths are held. Once a complete path is established, it is held for a random time interval which has a negative exponential time interval with unit mean.

Four types of traffic models are considered:

- Saturated network with uniform traffic
- · Saturated network with a hot-spot
- Network with small population and uniform traffic
- Network with small population and a hot-spot

The hot-spot model considered is one in which the hot pin has routing probability twice that of a cool pin.

The throughput results of the validation are tabulated in Tables 7.1, 7.2, 7.3 and 7.4. As can be seen from the first two tables, the analytical model is very accurate (less than 1% error) when the network is saturated. The relative errors are noticeably higher (up to 2.9%) in the case of the non-saturated network because of the additional assumptions required in modelling the closed system by the birth and death process (Section 3.3), i.e.

that all arrangements of customers on a given number of active inputs are equally likely.

No. of Stages	Analytical model	Simulation model	95% confidence interval	Relative Error (%)
2	2.000	1.992	(1.952, 2.032)	0.40
3	3.200	3.185	(3.143, 3.228)	0.47
4	5.333	5.375	(5.313, 5.437)	-0.78
5	9.143	9.163	(9.101, 9.225)	-0.22
6	16.00	15.97	(15.85,16.08)	0.19

Table 7.1 Saturated, uniform delta network

No. of Stages	Hot-spot probability	Analytical model	Simulation model	95% confidence interval	Relative Error (%)
2	0.400000	1.896	1.892	(1.866, 1.917)	0.21
3	0.222222	3.055	3.057	(3.017, 3.097)	-0.07
4	0.117647	5.174	5.193	(5.115, 5.271)	-0.37
5	0.060606	8.996	8.989	(8.898, 9.079)	0.08
6	0.030769	15.88	15.84	(15.71,15.97)	0.25

2

Table 7.2 Saturated, non-uniform delta network

No. of Stages	Network Population	Analytical model	Simulation model	95% confidence interval	Relative Error (%)
2	4	1.612	1.644	(1.603, 1.685)	-1.9
3	8	2.548	2.543	(2.498, 2.567)	0.20
4	16	4.283	4.227	(4.172, 4.283)	1.3
5	32	7.460	7.248	(7.198, 7.299)	2.9
6	64	13.28	12.98	(12.89, 13.08)	2.3

Table 7.3 Non-saturated, uniform delta network

No. of	Network	Hot-spot	Analytical	Simulation	95% confidence	Relative
Stages	Population	probability	model	model	interval	Error (%)
2	4	0.400000	1.564	1.579	(1.559, 1.598)	-0.95
3	8	0.222222	2.479	2.485	(2.440, 2.531)	-0.24
4	16	0.117647	4.206	4.174	(4.104, 4.244)	0.77
5	32	0.060606	7.385	7.216	(7.139, 7.293)	2.3
6	64	0.030769	13.21	12.88	(12.77,12.99)	2.6
	Table 7	A Non ast			dalle matemanle	

 Table 7.4
 Non-saturated, non-uniform delta network

The confidence intervals were estimated from the simulation results by using the batch means approach. Each sample run modelled 5000 units of time and 5 samples were taken, giving a total simulation time of 25,000 units. For large networks, the whole simulation run took over an hour to complete on a SUN 3 workstation. Notice that in the saturated, uniform network with 2 stages, the simulator's 95% confidence limits do include the exact result of 1.9993 (to 4 d.p.), but that the analytical model provides the more accurate estimate.

7.2 Numerical Results

The graph of throughput against the number of customers, N, for a closed system compares three types of 16-way interconnections.

- The full 16-way crossbar with uniform routing (Section 4)
- The 4-stage delta-2 network with uniform routing (Section 5)
- The 4-stage delta-2 network with hot-spot routing probability 0.2 (Section 6)

As shown in Figure 6, the full crossbar gives the highest throughput and approaches its asymptote closely only when N is large (greater than 100). The uniform delta-2 network, becomes saturated much faster (when N is about 100) because tasks face path conflicts on top of memory conflicts. In the presence of a hot-spot this network gives even lower throughput and becomes saturated even faster as the hot pin's decode tree saturates causing more path conflicts. The uniform interconnections only show similar throughput when the population is very small (less than 5). This indicates that even at low loads the blocking caused by path conflicts significantly reduces the throughput of a delta network.



Figure 6 Comparison of performance of 16-way networks

Figure 7 shows the effect of hot-spot contention in a 16-way, 4-stage delta-2 network connected to a bank of DMA servers in a closed system with various populations, N. The graph of throughput against the routing probability displays the classic hot-spot effect for populations N=8, 16 and the saturated case in which all input pins are always active. This phenomenon is caused by the additional internal contention placed by the hot-spot on other decode trees. For the saturated case, as ρ increases the throughput reaches its peak when all the output pins are uniformly utilised (ρ =1/16), drops sharply as more of the traffic is routed to the hot-spot, and then follows the curve 1/ ρ very closely, ultimately giving the throughput of a serial link, when ρ =1. In fact we can easily show that the 1/ ρ curve provides an upper bound for *all* protocols as follows.





Suppose the routing probability to network output pin i is ρ_i and the probability that pin i is active (i.e. pin i's throughput when $\mu=1$) is t_i . Now since ρ_i of the total throughput T(N) is routed to output pin i, we have $t_i = \rho_i T(N)$ which implies $T(N) \le 1/\rho_i$ for all i, i.e. $T(N) \le 1/\max(\rho_i)$. Hence if ρ is the hot-spot routing probability $T(N) \le 1/\rho$. When ρ is large, the throughput curves tend towards the curve $1/\rho$, particularly for large N, because the hot output pin is almost always active. These curves relate to a circuit-switched network in which partial paths are held. When paths are released and the

server retries immediately, throughput increases and so would lie somewhere between the curve $1/\rho$ and the corresponding curve shown in Figure 7; $1/\rho$ is still an upper bound by the same reasoning. In this idealised case, when there is uniform routing ($\rho=1/16$), the throughput is 16μ because all output pins are busy, so the throughput would drop considerably faster if ρ increased slightly, thereby giving a more dramatic hot-spot effect. Similarly, packet-switched networks would give an even more drastic reduction in throughput for a slight increase in traffic to a particular output.

For all of these protocols, in larger networks and with large N, the hot-spot effect is more pronounced because the throughput is greater in the uniformly utilised case and the throughput curve closely follows the $1/\rho$ curve as ρ increases past about 0.2. However, it may be unlikely that very high hot-spot routing probabilities are encountered in practice, especially when there are a large number of output pins.

8 Conclusion

1

The work described in this paper extends the flow equivalent server aggregation technique to incorporate passive resources which represent switching networks of various types. In this way the modelling of large-scale parallel computer architectures is greatly simplified, and we have presented results predicting the throughput of circuit-switched networks in which partial paths are held by blocked transactions during path building. We can therefore see the degradation in performance suffered by a MIN compared with the equivalent full crossbar offering the same connectivity (if it were possible to fabricate such a device). We also derived for the first time by analytical methods, we believe, results showing the effect of "hot-spots" in asynchronous circuit-switched networks, where one destination address is more frequently selected than the others and causes performance to suffer, ultimately giving the throughput achieved by a single serial link. Previously, this effect has only been shown by simulation, e.g. [PN85].

The circuit switching communication protocol considered is an important one, being the simplest to fabricate and the one to be used in a new generation of optical switches, [JO86] for example. However, there are a number of other protocols that should be considered. The simplest of these is infinitely buffered packet switching which has already been studied as we noted in the Introduction, but many hybrids are also possible; perhaps the simplest being packet switching with limited buffering and hence blocking. In addition a model should be developed in which switching times are not neglected.

Modelling circuit-switched multi-stage interconnection networks

Certainly, for contemporary electronic crossbars and even quite small message lengths, switching time is negligible compared with data transmission time. But for the optical devices referred to the converse holds for the present - unless whole files of data are normally transmitted as single messages. In such a model, the path building process cannot be assumed instantaneous.

The work presented in this paper includes the modelling of hot-spots parameterised by the probability that an arriving task selects the destination address of the hot-spot. One relatively simple extension to the model presented here, is to allow full non-uniform routing so that the effect of more than one hot-spot and the effect of proximity of hotspots can be studied quantitatively. A further extension would be to model systems in which certain output addresses are favoured, depending on the input pin on which a task arrives. Here, of course, the arrival process would no longer be the same at all inputs.

Finally, our analysis should also be adapted to represent circuit switching in which transactions which are blocked during path building *release* their partial paths, effectively being "lost" and having to retry after some (randomly distributed) delay. As indicated in Section 7, for this protocol, the hot-spot effect on throughput appears to be more dramatic than when partial paths are held, especially for fast retry rates. Thus it is important to model this retry protocol and a pilot study which uses fixed-point methods may be found in [CO87]. This analysis would then be close to the classical modelling of telephone networks. We would also expect packet-switched MINs to show a more pronounced hot-spot effect because when there is no hot-spot they achieve higher throughput. This would then allow comparison of the effectiveness of protocols in a given MIN with respect to the extent of hot-spot contention.

References

- [AK87] I.F. Akyildiz "General Closed Queueing Networks with Blocking" in Proc. 12th Annual International Symposium on Computer Performance Modelling, December 1987, Brussels, Belgium.
- [BA83] L.N. Bhuyan, D.P. Agrawal, "Design and Performance of Generalised Interconnection Networks", *IEEE Transactions on Computers* C-32, No. 12, (1983), pp. 1081-1090.
- [BCMP75] F. Baskett, K. M. Chandy, R.R. Muntz, F.G. Palacios, "Open, Closed and Mixed Networks of Queues with Different Classes of Customers", *Journal of the ACM* 22, 2 (1975)

Modelling circuit-switched multi-stage interconnection networks

- [CHW75] K. M. Chandy, U. Herzog, L. Woo, "Parametric Analysis of Queueing Networks", *IBM J. Res. Develop.*, January, 1975.
- [CI75] E. Cinlar, Introduction to Stochastic Processes, Engelwood Cliffs, N.J.: Prentice-Hall (1975)
- [CO87] G.A. Cope, *The Modelling of Circuit-switched Multi-stage* Interconnection Networks, MSc Thesis, Imperial College, 1987.
- [CO77] P. J. Courtois, *Decomposability*, Academic Press, 1977.
- [DR81] J. Darlington, M. J. Reeve, "ALICE : A Multiprocessor Reduction Machine for the Parallel Evaluation of Applicative Languages", *ACM/MIT Conference on Functional Programming Languages and Computer Architecture*, 1981.
- [GH88] U. Garg, Y.P. Huang, "Decomposing Banyan Networks for Performance Analysis", *IEEE Transactions on Computers* C-37, No. 3, 1988, pp.371-376.
- [GL73] L.R. Goke, G.L. Lipovski, "Banyan Networks for Partitioning Multiprocessor Systems", *Proc. of the First Annual Symposium on Computer Architecture*, 1973, pp. 21-28.
- [HA87] P.G. Harrison, "Queueing Models of the Large-scale Parallel Computer Architectures", ORSA/TIMS Workshop on Queueing Networks and their Applications, New Jersey, January 1987.
- [HF86] P.G. Harrison, A.J. Field "Performance Modelling of Parallel Computer Architectures", in *Proc. Performance '86 and ACM Sigmetrics 1986*, May 1986. pp. 18-27.
- [HR86] P. G. Harrison, M. J. Reeve, "The Parallel Graph Reduction Machine, ALICE", in *Proc. Workshop on Graph Reduction*, Santa Fe, September, 1986, also published in LNCS series, Springer-Verlag.
- [JD81] J.R. Jump, D.M. Dias, "Analysis and Simulation of Buffered Delta Networks", *IEEE Trans on Computers*, Vol. C-29, No.9, September 1980, pp. 791-801.
- [JO86] Digital Optics : Application to Computing and Communications. DTI JOERS (Joint Opto-Electronic Research Scheme) proposal, Issue 2, March 1986.
- [KE86] F.P. Kelly, "Blocking Probabilities in Large Circuit-switched Networks", *Advances in Applied Probability* 18, pp. 473-505.
- [KN68] D. Knuth, Fundamental Algorithms, Vol. 1, Addison Wesley, 1968.
- [KP86] M. Kumar, G.F. Pfister, "The Onset of Hot spot Contention", *Proc.1986* International Conference on Parallel Processing, August 1986, pp. 28-34
- [KU79] P.J. Kuehn, "Approximate Analysis of General Queueing Networks by Decomposition", *IEEE Trans. Comm.*, COM-27,1, 1979.

Modelling circuit-switched multi-stage interconnection networks

- [LI61] J.D.C Little, (1961), "A Proof of the Queueing Formula $L = \lambda W$ ", Operations Research 9 (3), pp. 383-387.
- [LI87] Y.S. Liu, 'The Delta Network Performance and the "Hot Spot" Traffic', Ultracomputer Project Technical Report (1987), Courant Institute, New York University.
- [MC87] D. Mitra, R. Cieslak, "Randomized Parallel Communications on an Extension of the Omega Network", *J.ACM* Vol 34 No.4, October 1987, pp. 802-824.
- [PA81] J.H. Patel, "Performance of Processor-Memory Interconnections for Multiprocessors", *IEEE Transactions on Computers*, October 1981, pp. 771-780.
- [PA89] N.M. Patel, *Models of Circuit-switched Interconnection Networks*, PhD Thesis, Imperial College, August 1989.
- [PH88] A. Pombortsis, C. Halatsis, "Performance of Crossbar Interconnection Networks in Presence of "Hot Spots"", *Electronics Letters* 24, No.3 (1988), pp. 182-184.
- [PN85] G.F. Pfister, V.A. Norton, "Hot spot Contention and Combining in Multistage Interconnection Networks", *IEEE Transactions on Computers*, Vol. C-34, No. 10, October 1985, pp. 943-948.
- [SI85] H.J. Siegel, Interconnection Networks for Large-Scale Parallel Processing: Theory and Case Studies, Lexington Books, 1985.
- [SM81] K.C. Sevcik, I. Mitrani, "The Distribution of Queueing Network States at Input and Output Instants", J. ACM, Vol 28, No. 2, April 1981, pp. 358-371.

and a standard standard and a standard a standard and standard and a standard and a standard and a standard and International Additional Additional Annuality of Antifathian adult (1999) (2019) and a feel of a standard and a Additional Ad

DISCUSSION

Rapporteur : Steve Caughey

During his presentation Dr. Harrison explained that one of the assumptions he had made concerning a two by two crossbar was that the mean residual holding time equals the mean holding time. Professor B. Littlewood asked if the holding time function was exponential and Dr. Harrison replied that yes, that needed to be true for this assumption to hold.

After the lecture Dr. A. Burns asked what conclusions could be drawn from this work. Dr. Harrison said that he had produced a realistic model for circuitswitched delta networks. He said it was possible to apply the decomposition method he had described to a large system and to use the results to find the service rate of a simple aggregate server. Also, his work shows that the hot-spot effect is less dramatic under the protocol used on the ALICE machine because the maximum throughput is low compared to other protocols. This is because partial paths are held within that protocol. Improved throughput could be provided if partial paths were released and further work is needed on the ALICE machine to investigate this.

Professor Littlewood asked why the blocking probabilities are treated so deterministically. He suggested that the probabilities might change with time. Dr. Harrison said that his model isn't transient, it assumes steady state and so is an approximation. He confirmed that Professor Littlewood was referring to situations such as `bursty' traffic and said that this was a complex problem on which further work was required.

Dr. I. Mitrani asked how long it took to reach steady state in the simulations. Dr. Harrison replied "quickly" but had no details available. He asked Dr. Mitrani what he thought the behaviour would be. Dr. Mitrani said that in a single queue model steady state is reached exponentially quickly but he would expect the time taken to grow with the number of pins on the switch. Dr. Harrison confirmed that they had observed this happening.

Professor C. Girault asked if it was true that the hypotheses which had been presented were more convenient for applications such as telephonic switching rather than distributed Operating Systems where processes might have to synchronise therefore creating favoured pairs of input and output pins. Dr. Harrison confirmed this to be true.

Professor Girault asked if the results which Dr. Harrison had displayed at the end of his presentation were for a heavily loaded network. Dr. Harrison replied that, no, the results referred to a network which was not saturated. He stressed that the results for a saturated network were more accurate as, for a nonsaturated network, the distribution of jobs over the set of input pins needs to be taken into account.

A second seco

Production of the second s a second se