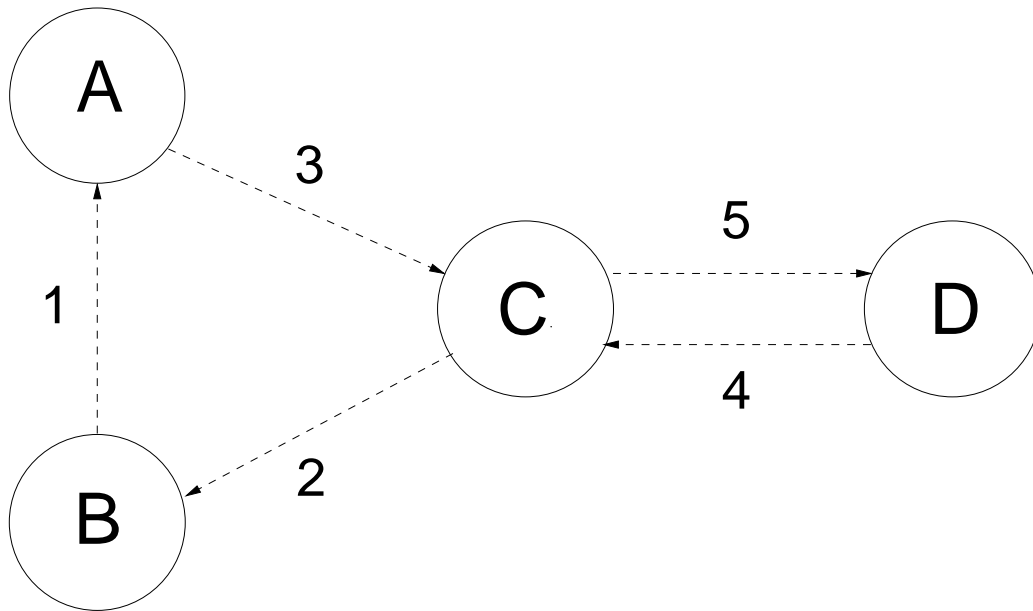# RESOURCE ALLOCATION IN GENERAL QUEUEING NETWORKS

by

Phil Pollett

Department of Mathematics
The University of Queensland

# PACKET SWITCHING NETWORKS



A packet switching network with 4 nodes (labelled A,B,C and D) and 5 links (labelled $1, 2 \ldots, 5$)

# PACKET SWITCHING NETWORKS

$N$ switching nodes (labelled $n = 1, 2, \ldots, N$)

$J$ links (labelled $j = 1, 2, \ldots, J$)

Poisson traffic on route $m \to n$ at rate $\nu_{mn}$

(type-$mn$ traffic)

Common expected message length: $1/\mu$ (bits)

(Message lengths have an arbitrary distribution which does not depend on type)

Transmission rate on link $j$ is $\phi_j$ (bits/sec.)

(There is a first-come first-served (FCFS) discipline at each link)

# ROUTING MECHANISMS

*Fixed routing*

Define $R(m, n)$ to be the collection of (distinct) links used by type-$mn$ traffic:

$$R(m, n) = \{r_{mn}(1), \ldots, r_{mn}(s_{mn})\} \,,$$

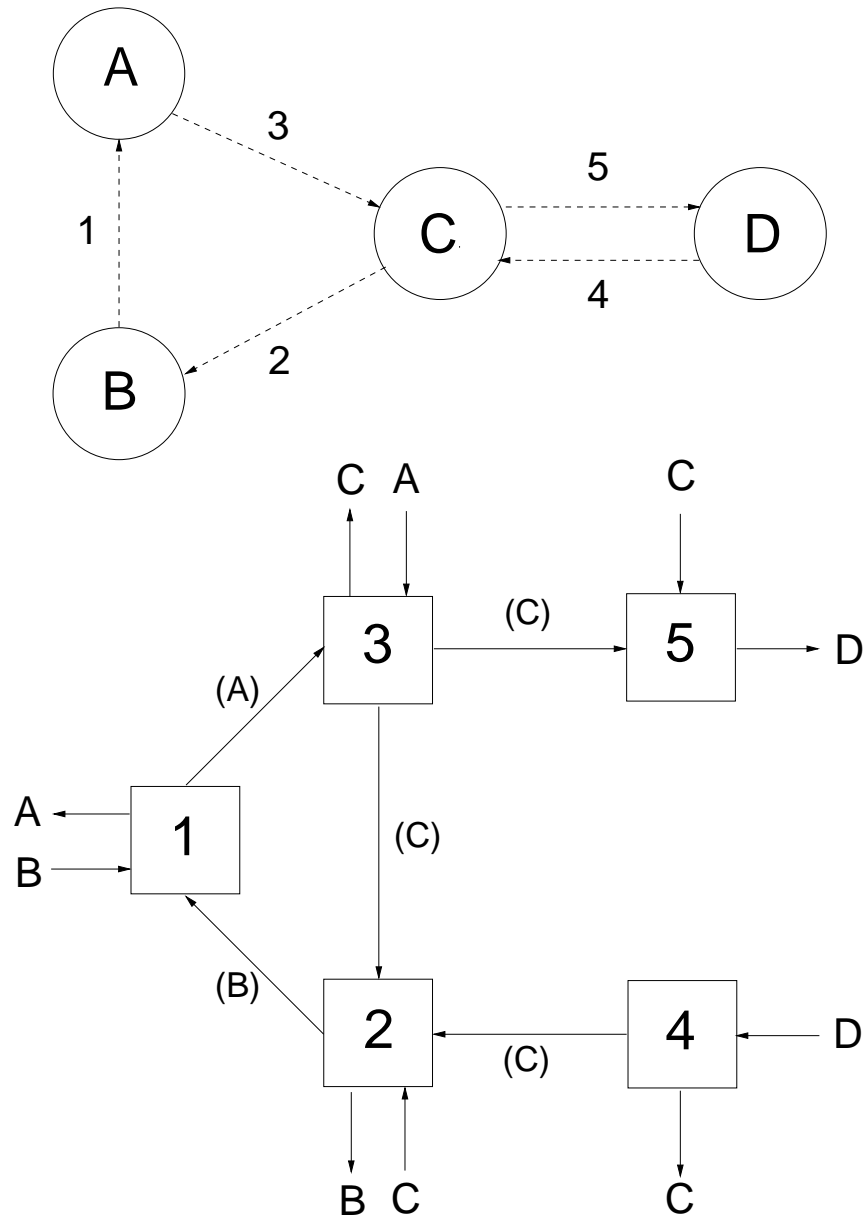where $s_{mn}$ is the number of links on route $m \to n$ and $r_{mn}(s)$ is the link used at stage $s$.

*Random alternative routing*

This can be accommodated within the framework of fixed routing by allowing a finer classification of type $(mni)$:

$$\nu_{mni} = \nu_{mn} q_{mni} \,,$$

where $q_{mni}$ is the probability that alternative route $i$ is chosen; there is fixed set of alternative routes for each OD pair $(m, n)$.

# PACKET SWITCHING NETWORK AND CORRESPONDING QUEUEING NETWORK

# A NETWORK OF QUEUES

Links $\leftrightarrow$ queues

Messages $\leftrightarrow$ customers

$T$ - set of customer types

$\nu_t$ - arrival rate of type-$t$ customers

(Independent Poisson streams)

Route for type-$t$ customers:

$$R(t) = \{r_t(1), \ldots, r_t(s_t)\} \ .$$

# A NETWORK OF QUEUES

If message lengths have an *exponential distribution* the links behave *independently* (indeed, *as if they were isolated*), each with independent streams of Poisson offered traffic (independent among types). For example, if

$$\alpha_j(t,s) = \begin{cases} \nu_t, & \text{if } r_t(s) = j, \\ 0, & \text{otherwise,} \end{cases}$$

so that the arrival rate at link $j$ is given by

$$\alpha_j = \sum_{t \in T} \sum_{s=1}^{s_t} \alpha_j(t,s)$$

and the demand by $a_j = \alpha_j/\mu$ (bits/sec), then, if the system is stable ($a_j < \phi_j$ for each $j$), the expected number of messages at link $j$ is

$$\mathsf{E}(n_j) = \frac{a_j}{\phi_j - a_j}$$

and the expected delay is

$$\mathsf{E}(W_j) = \frac{1}{\alpha_j}\left(\frac{a_j}{\phi_j - a_j}\right) = \frac{1}{\mu\phi_j - \alpha_j}.$$

# THE INDEPENDENCE ASSUMPTION

Kleinrock (1964) proposed the following assumption: that successive messages requesting transmission along *any given link* have lengths which are *independent and identically distributed*, and that message lengths at different links are independent.

Thus, we shall assume that at link $j$ message lengths have a distribution function $F_j(x)$ which has mean $\mu_j^{-1}$ and variance $\sigma_j^2$.

Even under this assumption, the model (now a network of $\cdot/G/1$ queues with a FCFS discipline) is not analytically tractable. To make progress, we shall use the *Residual-life Approximation* (Pollett (1984)).

# THE RESIDUAL LIFE APPROXIMATION

Let $Q_j(x)$ be the distribution function of the *queueing time* at link $j$: the time a message spends in the buffer *before* transmission. The *Residual-Life Approximation* (RLA) provides an accurate approximation for $Q_j(x)$:

$$Q_j(x) \simeq \sum_{n=0}^{\infty} \Pr(n_j = n) G_j^{(n)}(x) , \qquad (1)$$

where

$$G_j(x) = \mu_j \int_0^{\phi_j x} (1 - F_j(y)) \, dy$$

and $G_j^{(n)}(x)$ denotes the $n$-fold convolution of $G_j(x)$. The distribution of $n_j$, the number of messages at link $j$, used in (1) is that of the corresponding *quasireversible network* of symmetric queues obtained by imposing a symmetry condition at each link $j$. In the present context, this amounts to replacing FCFS by a last-come first-served (LCFS) discipline.

# THE RESIDUAL LIFE APPROXIMATION

One immediate consequence of (1) is that the expected queueing time $\bar{Q}_j$ is approximately

$$\frac{1 + \mu_j^2 \phi_j^2}{2\mu_j \phi_j} \, \mathsf{E}(n_j) \,,$$

where $\mathsf{E}(n_j)$ is the expected number of messages at link $j$ in the quasireversible network. Hence, the expected delay at link $j$ is approximated as follows:

$$\mathsf{E}(W_j) \simeq \frac{1}{\mu_j \phi_j} + \frac{1 + \mu_j^2 \phi_j^2}{2\mu_j \phi_j} \, \mathsf{E}(n_j) \,. \qquad (2)$$

In the RLA, it is only $\mathsf{E}(n_j)$ which changes when the service discipline is altered. For the present FCFS discipline $\mathsf{E}(n_j)$ is given by

$$\mathsf{E}(n_j) = \frac{\alpha_j}{\mu_j \phi_j - \alpha_j}$$

# OPTIMAL ALLOCATION OF EFFORT

We shall minimize the average network delay, or equivalently the average number of messages in the network:

$$\bar{m} = \sum_{j=1}^{J} \alpha_j \mathsf{E}(W_j)$$

(using the RLA for $\mathsf{E}(W_j)$).

$F$ - overall network budget

$f_j \phi_j$ (\$-seconds/bit) - cost of operating link $j$

  (The cost of operating link $j$ is proportional to the capacity $\phi_j$)

Thus, we should choose the capacities subject to the cost constraint

$$\sum_{j=1}^{J} f_j \phi_j = F \,.$$

# THE PROBLEM

Let $c_j = \mu_j^2 \sigma_j^2$ be the squared coefficient of variation of $F_j(x)$ and let $a_j = \alpha_j/\mu_j$.

Minimize

$$\bar{m} = \sum_{j=1}^{J} a_j \left\{ \frac{1}{\phi_j} + \frac{a_j(1 + c_j)}{2\phi_j(\phi_j - a_j)} \right\}$$

over $\phi_1, \ldots, \phi_J$ subject to

$$\sum_{j=1}^{J} f_j \phi_j = F.$$

Introduce a lagrange multiplier $\lambda^{-2}$; our problem then becomes one of minimizing

$$L(\phi_1, \ldots, \phi_J; \lambda^{-2}) = \bar{m} + \frac{1}{\lambda^2} \left( \sum_{j=1}^{J} f_j \phi_j - F \right).$$

Setting $\partial L/\partial\phi_j = 0$ yields a quartic polynomial equation in $\phi_j$:

$$2f_j\phi_j^4 - 4a_jf_j\phi_j^3 + 2a_j(a_jf_j - \lambda^2)\phi_j^2$$
$$- 2\epsilon_j a_j^2\lambda^2\phi_j + \epsilon_j a_j^3\lambda^2 = 0, \quad (3)$$

where $\epsilon_j = c_j - 1$.

Find solutions such that $\phi_j > a_j$ (recall that this latter condition is a requirement for stability).

Using the transformation

$$\phi_j f_j/F \to \phi_j, \ a_j f_j/F \to a_j, \ \lambda^2/F \to \lambda^2, \quad (4)$$

the problem reduces to one with unit costs $f_j = F = 1$; equation (3) becomes

$$2\phi_j^4 - 4a_j\phi_j^3 + 2a_j(a_j - \lambda^2)\phi_j^2$$
$$- 2\epsilon_j a_j^2\lambda^2\phi_j + \epsilon_j a_j^3\lambda^2 = 0, \quad (5)$$

and the constraint becomes

$$\sum_{j=1}^{J} \phi_j = 1. \quad (6)$$

# EXPONENTIAL SERVICE TIMES

If transmission times are exponentially distrib-
uted ($\epsilon_j = 0$ for each $j$) it is easy to verify that
(5) has a unique solution on $(a_j, \infty)$ given by

$$\phi_j = a_j + |\lambda|\sqrt{a_j}.$$

Upon application of the constraint (6) we ar-
rive at the optimal capacity assignment

$$\phi_j = a_j + \left(1 - \sum_{k=1}^{J} a_k\right) \frac{\sqrt{a_j}}{\sum_{k=1}^{J}\sqrt{a_k}},$$

for unit costs. In the case of general costs this
becomes

$$\phi_j = a_j + \frac{1}{f_j}\left(F - \sum_{k=1}^{J} f_k a_k\right) \frac{\sqrt{f_j a_j}}{\sum_{k=1}^{J}\sqrt{f_k a_k}},$$

after applying the transformation (4). This is
a result obtained by Kleinrock (1964).

# THE GENERAL CASE

We shall adopt a perturbation approach, assuming that the lagrange multiplier and the optimal allocation take the following forms:

$$\lambda = \lambda_0 + \sum_{k=1}^{J} \lambda_{1k}\epsilon_k + O(\epsilon^2),$$

$$\phi_j = \phi_{0j} + \sum_{k=1}^{J} \phi_{1jk}\epsilon_k + O(\epsilon^2), \tag{7}$$

$$j = 1, \ldots, J,$$

where by $O(\epsilon^2)$ we mean terms of order $\epsilon_i\epsilon_k$. The zero[th] order terms come from Kleinrock's solution:

$$\phi_{0j} = a_j + \lambda_0\sqrt{a_j}, \quad j = 1, \ldots, J,$$

where

$$\lambda_0 = \frac{1 - \sum_{k=1}^{J} a_k}{\sum_{k=1}^{J} \sqrt{a_k}}.$$

# FIRST-ORDER SOLUTION

On substituting (7) into (5) we obtain an expression for $\phi_{1jk}$ in terms of $\lambda_{1k}$, which in turn is calculated using the constraint (6) and by setting $\epsilon_k = \delta_{kj}$ (the Kronecker delta).

To first order, the optimal allocation is

$$\phi_j = a_j + \lambda_0 \sqrt{a_j} - \frac{\sqrt{a_j}}{\sum_{k=1}^{J} \sqrt{a_k}} \sum_{k \neq j} b_k \epsilon_k$$

$$+ \left( 1 - \frac{\sqrt{a_j}}{\sum_{k=1}^{J} \sqrt{a_k}} \right) b_j \epsilon_j \,,$$

where

$$b_k = \frac{1}{4} \lambda_0 a_k^{3/2} \frac{a_k + 2\lambda_0 \sqrt{a_k}}{(a_k + \lambda_0 \sqrt{a_k})^2} \,.$$

# SENSITIVITY

Let $\phi_j$, $j = 1, 2, \ldots, J$, be the new optimal allocation obtained after incrementing $\epsilon_j$ by a small quantity $\delta > 0$. We find that, to first order in $\delta$,

$$\phi'_j - \phi_j = \left(1 - \frac{\sqrt{a_j}}{\sum_{k=1}^{J} \sqrt{a_k}}\right) b_j \delta > 0$$

and, for $i \neq j$,

$$\phi'_i - \phi_i = -\frac{\sqrt{a_i}}{\sum_{k=1}^{J} \sqrt{a_k}}(\phi'_j - \phi_j) < 0\,.$$