# Network Routing

F. P. Kelly

# Network routing

### By F. P. Kelly

*Statistical Laboratory, University of Cambridge, 16 Mill Lane,*
*Cambridge CB2 1SB, U.K.*

How should flows through a network be organized, so that the network responds sensibly to failures and overloads? The question is currently of considerable technological importance in connection with the development of computer and telecommunication networks, while in various other forms it has a long history in the fields of physics and economics. In all of these areas there is interest in how simple, local rules, often involving random actions, can produce coherent and purposeful behaviour at the macroscopic level. This paper describes some examples from these various fields, and indicates how analogies with fundamental concepts such as energy and price can provide powerful insights into the design of routing schemes for communication networks.

## 1. Introduction

Modern computer and telecommunication networks are able to respond to randomly fluctuating demands and failures by rerouting traffic and by reallocating resources. They are able to do this so well that, in many respects, large-scale networks appear as coherent, almost intelligent, organisms. The design and control of such networks require an understanding of a variety of fundamental issues, and this is providing an important stimulus to many areas of mathematics and engineering.

To give an example of current importance, a major practical and theoretical issue concerns the extent to which control can be decentralized. Over a period of time the form of the network or the demands placed on it may change, and routings may need to respond accordingly. It is rarely the case, however, that there should be a central decision-making processor, deciding upon these responses. Such a centralized processor, even if it were itself completely reliable and could cope with the complexity of the computational task involved, would have its lines of com-munication through the network vulnerable to delays and failures. Rather, decision-making should be decentralized and of a simple form: the challenge is to understand how such decentralized decision-making can be organized so that the network as a whole reacts intelligently to outside stimuli.

The behaviour of large-scale systems has been of great interest to mathematicians for over a century, with many examples coming from physics. The behaviour of a gas can be described at the microscopic level in terms of the position and velocity of each molecule. At this level of detail a molecule's velocity appears as a random process, with a stationary distribution as found by Maxwell (and later discussed by Erlang (1925)). Consistent with this detailed microscopic description of the system is macroscopic behaviour best described by quantities such as temperature and pressure. Similarly the behaviour of electrons in an electrical network can be described in terms of random walks, and yet this simple description at the

343

microscopic level leads to rather sophisticated behaviour at the macroscopic level: the pattern of potentials in a network of resistors is just such that it minimizes heat dissipation for a given level of current flow (Thomson & Tait 1879). The local, random behaviour of the electrons causes the network as a whole to solve a rather complex optimization problem.

Of course simple local rules may lead to poor system behaviour if the rules are the wrong ones. Road traffic networks provide a chastening example of this. Braess's paradox describes how, if a new road is added to a congested network, the average speed of traffic may fall rather than rise, and indeed everyone's journey time may lengthen. The paradox may actually have occurred during 'development' in the centre of Stuttgart (Knödel 1969), and counterintuitive consequences of road closures are often reported (*New York Times* 1990). The attempts of individuals to do the best for themselves lead to everyone suffering. It is possible to alter the local rules, by the imposition of appropriate tolls, so that the network behaves more sensibly, and indeed road traffic networks have long provided a key example of the economic principle that externalities need to be appropriately penalized if the invisible hand is to lead to optimal behaviour (Pigou 1920; Walters 1961).

A telephone network provides a fascinating example of a large-scale system where strange effects can occur. For instance, suppose that 'intelligent' exchanges react to blocked routes by rerouting calls along more resource-intensive paths. This in turn may cause later calls to be rerouted, and the cascade effect may lead to a catastrophic change in the network's behaviour. When exchanges strive to be efficient there is a possibility they may overdo it. In some respects the network's behaviour resembles water boiling. Just as a small change in the temperature of a body of water can cause a pronounced macroscopic effect, so a small change in the load on a network can produce an unexpected and massive failure. This discussion indicates the care that must be taken with the development of routing rules for large networks.

There is currently considerable interest in the similarities between complex systems from diverse areas of physics, economics and biology, and it is clear that topics such as noisy optimization and adaptive learning provide mathematical metaphors of value across many fields. The reader is referred to Pines (1987), Anderson *et al.* (1988) and Langton (1989) for the lively and thought-provoking proceedings of a series of workshops, and to Whittle (1986) for a study of statistical equilibrium in systems of interacting components that is both broad and penetrating. Our aim in this paper is, by comparison, much narrower: to describe within a common framework some examples from three areas, and thus to indicate how earlier, often much earlier, work on electrical networks and traffic flow has influenced recent work on routing schemes for communication networks.

In §2 we outline the connections among random walks, electrical networks and variational principles. We describe how an interacting particle system, precisely defined at a microscopic level in terms of simple, randomized, local rules, can also be described at an intermediate level of aggregation in terms of Ohm's Law and Kirchhoff's equations, and at a macroscopic level in terms of energy minimization. The interacting particle system we describe is due to Kingman (1969), and our treatment of the associated optimization problem derives from Whittle (1971). For a beautifully written account of the area, and of its history back to Lords Rayleigh and Kelvin, the reader is referred to Doyle & Snell (1984).

In §3 we consider multiclass flow models, including queueing and road traffic networks. If customers or drivers can choose their routes, then exercise of this choice

may force the system to a competitive equilibrium. If drivers attempt to minimize their own delay the resulting equilibrium flows will minimize a certain objective function defined for the network. However, the objective function is certainly not the average network delay: this is dramatically illustrated by Braess's (1968) paradox, outlined above. We describe a variant of this paradox: if drivers are provided with extra information about random delays ahead, the outcome may well be a new equilibrium in which delays are increased for everyone. Traffic-dependent tolls are sufficient to force the system to an equilibrium which minimizes average network delay: the tolls charge drivers for the delays they cause to others. The study of appropriate tolls has long been a topic of central importance in economics, and it is interesting to note that Pigou (1920, p. 194) used a simple two-node, two-link traffic network to illustrate the possibility that taxation could 'create an "artificial" situation superior to the "natural" one'. (We note in passing that developments in electronics now make feasible the practical application of both route guidance and road pricing (van Vuren & Smart 1990).) Potts & Oliver (1972) survey the characterization of equilibrium flow patterns as extremal values, and Nagurney (1987) provides a recent review of competitive equilibrium problems including more general traffic network models and related models of economic markets.

In the remainder of the paper we outline more recent work on the modelling of telecommunications networks. In §4 we show how the microscopic description of a telephone network, in terms of random arrival streams and rules for accepting and routing calls, leads the network to behave as if it is attempting to minimize a certain potential function. However, just as in road traffic networks, the potential function implicitly minimized may bear little relation to the network performance criteria of interest to system designers. Indeed Wroe *et al.* (1990) have described an example very similar to Braess's paradox, that has arisen in the design of the BT international access network, where the addition of capacity to a network causes the performance to become worse. Various other forms of perverse behaviour can be interpreted in terms of an implicit minimization: for example, in situations where alternative routing causes the potential function to have multiple minima, a slight change in traffic load may cause the network to jump catastrophically between minima. In §5 we discuss how explicit consideration of network performance criteria can lead, through notions of shadow prices and implied costs, to decentralized adaptive routing schemes which are at least attempting to optimize the right function.

Might it be possible to choose the microscopic rules governing the behaviour of a telephone network, the rules for accepting and routing calls, so that the network is implicitly optimizing sensible performance criteria, much as current flow in an electrical network is implicitly minimizing heat dissipation? This is a difficult and wide-ranging question, but at least in some circumstances it can be answered positively. In §6 we describe a scheme which has been developed for the British telephone network by researchers at Cambridge and at British Telecom's laboratories at Martlesham. The scheme, known as Dynamic Alternative Routing (DAR), is now being implemented in the British trunk network (Stacey & Songhurst 1987; Gibbens 1988; Gibbens *et al.* 1988; Key & Whitehead 1988). The scheme will lessen the impact of forecasting errors, make better use of spare capacity and respond robustly to failures and overloads. It will also permit flexible use of network resources enabling, for example, the rapid introduction of new services where demand is often uncertain. DAR uses very simple rules across the network, making constructive use of inherent randomness to search out good routing patterns. In this way the network itself

operates as a distributed computer, executing a highly parallel randomized algorithm to solve a complex optimization problem.

## 2. Random walk and electrical networks

We begin this section by describing a simple flow model due to Kingman (1969) and further discussed by Kelly (1979). The model can be viewed as a naive description of the movement of electrons in a conductor, and we shall phrase our discussion in terms of familiar electrical concepts such as current and potential. We use the model to develop the connections between random walks, electrical networks and variational principles.

Consider, then, the following interacting particle system. There is a set of sites, $J$, and each site $j \in J$ may be empty or may be occupied by a single particle. If site $j$ is occupied and site $k$ is empty then, with probability intensity $\lambda_{jk} (= \lambda_{kj})$, the particle at site $j$ moves to site $k$. If site $j$ is occupied then, with probability intensity $\mu_j$, the particle at site $j$ leaves the system entirely. If site $k$ is empty then, with probability intensity $\nu_k$, a particle arrives at site $k$ from outside the system. These rules define a finite state Markov process, about which we can ask a number of questions. What is the stationary probability $p_j$ that site $j$ is occupied? What is the average net rate of flow of particles from site $j$ to site $k$?

To answer these questions consider the following button model. Append to the set of sites $J$ two further sites, labelled 0 and 1, and let each site contain a single button. The buttons are distinguishable; we can imagine them to be of different colours. The buttons occupying sites $j$ and $k$ interchange positions with probability intensity $\lambda_{jk}$, for $j, k \in J \cup \{0, 1\}$, where

$$\lambda_{0j} = \lambda_{j0} = \mu_j, \quad \lambda_{1j} = \lambda_{j1} = \nu_j, \quad j \in J,$$

and $\lambda_{01} = \lambda_{10} = 0$. We see that any particular button performs a symmetric (and hence reversible) random walk around the sites of the system. Now imagine that buttons entering site 1 are painted black, while buttons entering site 0 are painted white. If $A$ is the set of those sites $j \in J$ which contain a black button then $A$ behaves stochastically just as does the set of occupied sites in the earlier interacting particle system. Thus to find the probability $p_j$ that site $j$ is occupied we need only look backwards through time at the earlier movements of the button which currently occupies site $j$. These movements form a symmetric random walk starting from site $j$ with transition intensities $\lambda_{jk}$ for $j, k \in J \cup \{0, 1\}$; $p_j$ is equal to the probability that this random walk reaches site 1 before site 0. Considering where the first step of the random walk takes the button leads to the equations

$$p_j = \sum_k \frac{\lambda_{jk}}{\sum_i \lambda_{ji}} p_k, \quad j \in J,$$

$$p_0 = 0, \quad p_1 = 1.$$

We can rewrite these equations as

$$\sum_k \lambda_{jk}(p_j - p_k) = 0, \quad j \in J, \tag{2.1}$$

$$p_0 = 0, \quad p_1 = 1. \tag{2.2}$$

Equations (2.1) and (2.2), however, are precisely Kirchhoff's equations for an

electrical network with nodes from the set $J \cup \{0, 1\}$: just interpret $p_j$ as the electrical potential of node $j$, connect nodes $j$ and $k$ by a wire of resistance $\lambda_{jk}^{-1}$, and hold nodes 0 and 1 at potentials 0 and 1 respectively. Equation (2.1) simply states that the total current flowing out of node $j$ is zero.

The average net flow of particles from site $j$ to site $k$ is

$$\lambda_{jk} P\{\text{site } j \text{ occupied and site } k \text{ empty}\} - \lambda_{kj} P\{\text{site } j \text{ empty and site } k \text{ occupied}\}$$
$$= \lambda_{jk}(p_j - p_k).$$

Put more formally, the system is a finite state Markov chain, and hence the net flow of particles from site $j$ to site $k$ over a time interval $[0, T]$ will almost surely converge as $T \to \infty$ to

$$u_{jk} = \lambda_{jk}(p_j - p_k). \tag{2.3}$$

This, however, is precisely the current flowing from node $j$ to node $k$ in the electrical network; equation (2.3) is just Ohm's law. The average flow of particles through the network can be calculated from either the flow out of node 1, or the flow into node 0, and is

$$U = \sum_k \lambda_{1k}(1 - p_k) = \sum_j \lambda_{j0} p_j, \tag{2.4}$$

which is precisely the total current flowing through the electrical network.

The energy dissipated by a potential difference of $p_j - p_k$ across a wire of resistance $\lambda_{jk}^{-1}$ is the product of the current and voltage, and is thus $\lambda_{jk}(p_j - p_k)^2$. The energy dissipation in the above electrical network is then

$$\tfrac{1}{2} \sum_{j, k} \lambda_{jk}(p_j - p_k)^2,$$

where the summation runs over $j, k \in J \cup \{0, 1\}$, and the factor $\tfrac{1}{2}$ is necessary since each edge is counted twice in the sum. Consider the following problem.

$$\text{Minimize} \quad \tfrac{1}{2} \sum_{j, k} \lambda_{jk}(p_j - p_k)^2, \tag{2.5a}$$

$$\text{over} \quad (p_j, j \in J), \tag{2.5b}$$

$$\text{subject to} \quad p_0 = 0, \quad p_1 = 1. \tag{2.5c}$$

The objective function (2.5a) is strictly convex, and a differentiation yields that the unique optimum is given by the solution to Kirchhoff's equations (2.1) and (2.2). This is the *Dirichlet principle*: the potentials taken within the electrical network minimize the total energy dissipation.

The energy dissipated by a current flow $u_{jk}$ through a wire of resistance $\lambda_{jk}^{-1}$ is $u_{jk}^2/\lambda_{jk}$. Consider, then, the following problem, where the objective function is again the energy dissipation of the network, but now expressed in terms of currents rather than potentials.

$$\text{Minimize} \quad \tfrac{1}{2} \sum_{j, k} u_{jk}^2/\lambda_{jk}, \tag{2.6a}$$

$$\text{over} \quad u_{jk}(= -u_{kj}), \quad j, k \in J \cup \{0, 1\}, \tag{2.6b}$$

$$\text{subject to} \quad \sum_k u_{jk} = \begin{cases} 0 & j \in J, \\ -U & j = 0, \\ U & j = 1. \end{cases} \tag{2.6c}$$

The conditions $(2.6c)$ require a flow $U$ of current into the network from node 1, a flow $U$ of current out of the network to node 0, and that flow be balanced at nodes $j \in J$. Again the objective function is strictly convex: hence the optimum is unique and can be obtained by differentiating the *lagrangian form*

$$L(u;p) = \tfrac{1}{2} \sum_{j,k} u_{jk}^2 / \lambda_{jk} - 2 \sum_j p_j \left( \sum_k u_{jk} \right),$$

where $p_j, j \in J \cup \{0,1\}$, are now Lagrange multipliers, and where for later convenience we have introduced an extra factor of 2 before the constraints.

Differentiating $L$ with respect to $u_{jk}$ we obtain

$$\partial L / \partial u_{jk} = 2(u_{jk}/\lambda_{jk} - p_j + p_k),$$

since $u_{kj} = -u_{jk}$. Thus the optimum takes the form

$$u_{jk} = \lambda_{jk}(p_j - p_k), \quad j,k \in J \cup \{0,1\},$$

for a choice of Lagrange multipliers $p_j$ that cause the resulting $u_{jk}$ to satisfy the constraints $(2.6c)$; but the $p_j$s that solve (2.1) and (2.2) *are* such a choice. Thus the Lagrange multipliers should be set equal to the electrical potentials, and the resulting flows $u_{jk}$ are optimal. This is *Thomson's principle*: the flow pattern of current within an electrical network is that which minimizes the energy dissipation over all flow patterns achieving the same total current.

Together the Dirichlet principle and Thomson's principle provide one of the many examples from physics of complementary variational principles (Whittle 1971). From the viewpoint of optimization theory, if the problem (2.6) is recast as one of maximizing flow for a given energy dissipation then it is possible to obtain the problem (2.5) as the formal lagrangian dual.

In this section we have seen an example of a system which can be viewed at various levels of abstraction. At the microscopic level, it can be viewed in terms of particles performing simple random walks. If average rates of particle flow are studied, then these flows are given by Kirchhoff's equations (2.1), (2.2), and Ohm's law (2.3). Finally, at the network level, the flow patterns that emerge solve a constrained optimization problem whose objective function is the network's energy dissipation.

Thomson's principle, in particular, is an extremely useful and suggestive result. It follows directly that if a link is added to an electrical network and the same total current is carried then the energy dissipation is reduced: after all, the old flow pattern remains feasible for the new optimization problem. Might it be possible to design telecommunication networks similarly, so that, for example, additions of capacity are necessarily helpful? Before considering this question further, we look first at queueing networks, and the ways in which implicit optimization can go wrong.

## 3. Queueing networks and traffic flow

In the interacting particle system of the last section the particles behaved similarly, each attempting to perform a symmetric random walk. In this section we study a more complicated system, where particles are of different classes and where class determines the direction of flow through the network.

We begin with a brief description of an open multiclass network of $\cdot/M/1$ queues (for a more leisurely introduction see Kelly (1979)). Let the network have a set $J$ of

queues, and suppose that a customer entering the system is labelled with the route he will follow through the network. More specifically, suppose that customers labelled $r = (r(1), r(2), \ldots, r(M_r))$ arrive at the system in a Poisson stream of rate $\nu_r$ and pass through the sequence of queues

$$r(1), r(2), \ldots, r(M_r),$$

before leaving the system. Thus the queue which a customer labelled $r$ visits at stage $m\,(= 1, 2, \ldots, M_r)$ of his passage through the network is queue $r(m)$. Write $R$ for the set of possible routes, and assume that as $r$ varies over $R$ it indexes independent Poisson arrival streams. Let each queue have a single server operating a first come first served discipline, and suppose that a customer visiting queue $j$ has a service time there which is exponentially distributed with parameter $\phi_j$ and independent of all other service times and of the arrival streams at the network.

Write $j \in r$ if route $r$ passes through queue $j$ at some stage, and, for simplicity of notation, suppose no route passes through the same queue more than once. Let

$$\rho_j = \sum_{r:j\in r} \nu_r,$$

and suppose $\rho_j < \phi_j, j = 1, 2, \ldots, J$. Thus $\rho_j$ measures the throughput of queue $j$.

Let $n_j(t)$ be the number of customers in queue $j$ at time $t$, and let $n(t) = (n_j(t), j \in J)$. Then the (non-Markov) stochastic process $(n(t), t \geq 0)$ has a unique stationary distribution, and under this distribution $\pi(n) = P\{n(t) = n\}$ is given by

$$\pi(n) = \prod_{j \in J} \pi_j(n_j), \tag{3.1}$$

where
$$\pi_j(n_j) = (1 - \rho_j/\phi_j)\,(\rho_j/\phi_j)^{n_j}, \tag{3.2}$$

the geometric distribution familiar as the stationary distribution of an $M/M/1$ queue with arrival rate $\rho_j$. From the distribution (3.2) and Little's formula it follows that the mean sojourn time of a customer in queue $j$ is

$$D_j(\rho_j) = (\phi_j - \rho_j)^{-1}. \tag{3.3}$$

Sometimes we prefer to observe the network from the point of view of customers labelled $r$. When a typical customer labelled $r$ arrives at a queue $j$ on his route, the probability he finds $n_j$ customers already in that queue is given by expression (3.2); the sojourn time of a typical customer labelled $r$ in a queue $j$ on his route is exponentially distributed with parameter $\phi_j - \rho_j$.

A wide range of more general queueing networks share properties of the simple network described above, that the stationary distribution for the numbers in different queues has the product-form (3.1), and that the mean sojourn time of a customer in queue $j$ is a function $D_j(\rho_j)$ of the throughput of queue $j$.

Suppose now that it is possible to vary the parameters $\nu_r$. For example, if the model represents a packet-switched communication network, there may be a variety of possible routes $r$ through the network capable of linking two nodes, and the network may be able to shift traffic to routes with lower mean delays. More formally, let the label $s$ of a customer arriving at the network identify not a single route, but a set of routes, any of which could serve the customer. We can view $s$ as labelling a source-sink node pair. Set $H_{sr} = 1$ if $r \in s$, so that a customer labelled $s$ can be served by the route $r$, and set $H_{sr} = 0$ otherwise. This defines a $0-1$ matrix $H = (H_{sr}, s \in S,$

$r \in R$). For each $r \in R$ let $s(r)$ identify the unique value $s \in S$ such that $H_{sr} = 1$; we thus view $s(r)$ as the source-sink node pair served by route $r$. Let $A_{jr} = 1$ if route $r$ passes through link $j$, and set $A_{jr} = 0$ otherwise. This defines a $0-1$ matrix $A = (A_{jr}, j \in J, r \in R)$.

A *Wardrop equilibrium* is a collection $\nu = (\nu_r, r \in R), \rho = (\rho_j, j \in J)$ of non-negative numbers such that

$$H\nu = b, \tag{3.4a}$$

$$A\nu = \rho, \tag{3.4b}$$

and

$$\nu_r > 0 \Rightarrow \sum_{j \in r} D_j(\rho_j) = \min_{r' \in s(r)} \sum_{j \in r'} D_j(\rho_j), \quad r \in R. \tag{3.4c}$$

Equation $(3.4a)$ states that the traffic over routes $r$ serving the node pair $s$ sums to $b_s$, while equation $(3.4b)$ states that the traffic over routes through link $j$ sums to $\rho_j$. The implication $(3.4c)$ expresses the defining characteristic of a Wardrop equilibrium (Wardrop 1952), that if a route $r$ is actively used, it achieves the minimum delay over all routes serving the node pair $s(r)$.

Does a Wardrop equilibrium exist, and, if so, is it unique? To answer this question, consider the following optimization problem.

$$\text{Minimize} \quad \sum_{j \in J} \int_0^{\rho_j} D_j(z)\, \mathrm{d}z, \tag{3.5a}$$

$$\text{over} \quad \nu, \rho \geqslant 0, \tag{3.5b}$$

$$\text{subject to} \quad H\nu = b, \quad A\nu = \rho. \tag{3.5c}$$

Impose the mild condition that $D_j(z)$ is continuous and strictly increasing. Then the optimum can be found by differentiating the lagrangian form

$$L(\nu, \rho; \lambda, \mu) = \sum_{j \in J} \int_0^{\rho_j} D_j(z)\, \mathrm{d}z + \lambda^{\mathrm{T}}(b - H\nu) - \mu^{\mathrm{T}}(\rho - A\nu),$$

where $\lambda, \mu$ are vectors of Lagrange multipliers. But

$$\frac{\partial L}{\partial \nu_r} = -\lambda_{s(r)} + \sum_{j \in r} \mu_j,$$

and

$$\partial L / \partial \rho_j = D_j(\rho_j) - \mu_j.$$

Hence a maximum of $L$ over $\nu, \rho \geqslant 0$ occurs when

$$\mu_j = D_j(\rho_j)$$

and

$$\lambda_{s(r)} = \sum_{j \in r} \mu_j \quad \text{if} \quad \nu_r > 0$$

$$\leqslant \sum_{j \in r} \mu_j \quad \text{if} \quad \nu_r = 0.$$

The Lagrange multipliers have a simple interpretation: $\mu_j$ is the delay on link $j$, and $\lambda_s$ is the minimum delay over all routes serving the node pair $s$. The minima of the objective function $(3.5a)$ correspond precisely to Wardrop equilibria. Since $D_j(z)$ is strictly increasing the objective function $(3.5a)$ is a strictly convex function of $\rho$.
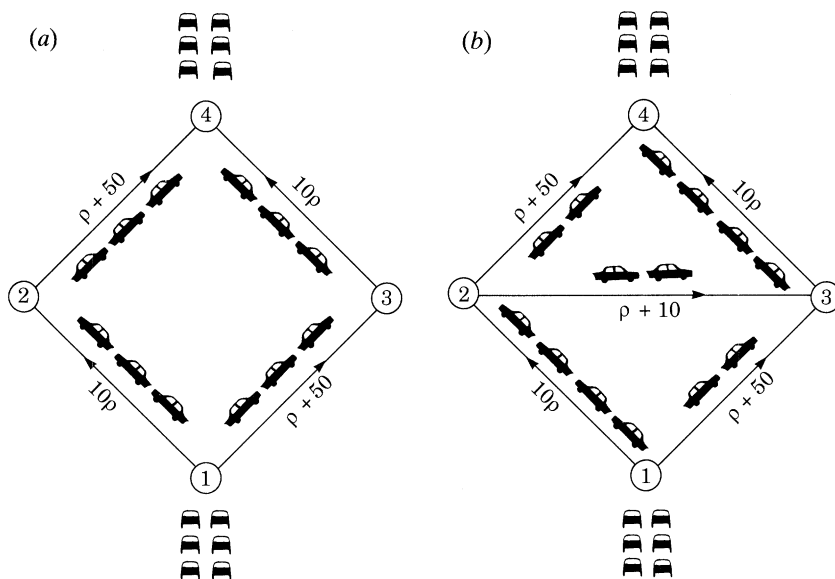
Figure 1. Braess's paradox. The addition of a link causes everyone's journey time to lengthen. (After Braess 1968; Cohen 1988; Cohen & kelly 1990.)

Hence there is a unique optimum for the flow vector $\rho$, although there may be many corresponding values of the non-negative vector $\nu$ satisfying the linear relations (3.5 $c$).

Thus, if traffic in the network distributes itself towards routes with lower mean delays, the equilibrium flows $\rho = (\rho_j, j \in J)$ will solve the optimization problem (3.5). This does *not* mean that average delays in the network will be minimal: and a striking illustration of this fact is provided by Braess's paradox (Braess 1968; Cohen & Kelly 1990). Consider the network illustrated in figure 1 $a$. Cars travel from node 1 to node 4, via either node 2 or node 3. The total flow is 6, and the link delays $D_j(\rho)$ are given next to links in the figure. The Wardrop equilibrium is shown. Now suppose that a new link is added, between nodes 2 and 3, as shown in figure 1 $b$. Traffic is attracted onto the new link, and the new Wardrop equilibrium is shown in figure 1 $b$. Observe that each car incurs a delay of 83 in figure 1 $a$, while each care incurs a delay of 92 in figure 1 $b$. Adding the new link has increased everyone's delay!

The point, of course, is that, while the Wardrop equilibrium minimizes the function (3.5 $a$), this objective function is not closely related to average delays. Consider next the following problem.

$$\text{Minimize} \quad \sum_{j \in J} \rho_j D_j(\rho_j), \tag{3.6 $a$}$$

$$\text{over} \quad \nu, \rho \geqslant 0, \tag{3.6 $b$}$$

$$\text{subject to} \quad H\nu = b, \quad A\nu = \rho. \tag{3.6 $c$}$$

Note that the objective function (3.6 $a$) now measures average network delay. Impose the fairly mild condition that the function (3.6 $a$) is convex and differentiable. The lagrangian for this problem is

$$L(\nu, \rho; \lambda, \mu) = \sum_{j \in J} \rho_j D_j(\rho_j) + \lambda^{\mathrm{T}}(b - H\nu) - \mu^{\mathrm{T}}(\rho - A\nu).$$

Again

$$\frac{\partial L}{\partial \nu_r} = -\lambda_{s(r)} + \sum_{j \in r} \mu_j,$$

but now

$$\partial L / \partial \rho_j = D_j(\rho_j) + \rho_j D_j'(\rho_j) - \mu_j.$$

Hence a maximum of $L$ over $\nu, \rho \geqslant 0$ occurs when

$$\mu_j = D_j(\rho_j) + \rho_j D_j'(\rho_j),$$

and

$$\lambda_{s(r)} = \sum_{j \in r} \mu_j \quad \text{if} \quad \nu_r > 0$$

$$\leqslant \sum_{j \in r} \mu_j \quad \text{if} \quad \nu_r = 0.$$

Again the Lagrange multipliers have a simple interpretation. Suppose that, in addition to the delay $D_j(\rho_j)$, users of link $j$ incur a traffic-dependent toll

$$T_j(\rho_j) = \rho_j D_j'(\rho_j). \tag{3.7}$$

Then $\mu_j$ is the combined cost of using link $j$, and $\lambda_s$ is the minimum cost over all routes serving the node pair $s$. If users select routes in an attempt to minimize the sum of their tolls and their delays, then they will produce a flow pattern which minimizes the average delay in the network $(3.6a)$.

In an electrical network, the addition of an extra wire connecting two nodes cannot increase the energy dissipation associated with a given current flow. We have seen that this follows from Thomson's principle, since the flow pattern before addition of the wire remains feasible afterwards. Similarly, in a traffic network with tolls $(3.7)$ the addition of a new link cannot increase average network delay, since the old flow pattern remains feasible for the new optimization problem $(3.6)$, and so if the new flow pattern is different it must improve the objective function $(3.6a)$.

In our above discussion of equilibrium flow patterns we have assumed that users are aware of the average delays $D_j(\rho_j)$, and perhaps tolls $T_j(\rho_j)$, along different routes. What if, instead, users do not know these quantities precisely, but have to rely on their previous experience along routes? To provide a clear framework for our discussion of this question, consider again the multiclass queueing network described at the beginning of this section. Suppose that a user travelling between a source–sink pair $s$ makes this journey repeatedly, but at times separated by long intervals, and that the user observes queue lengths or sojourn times at queues as he passes through. From these observations a user can estimate average delays $D_j(\rho_j)$ at queues $j \in r$, for each route $r \in s$. If users attempt to minimize their own delays across the network then it is reasonable to expect the stationary behaviour of the network to be described by the multiclass queueing network described earlier, with stationary distribution $(3.1)$–$(3.2)$ where $\rho = (\rho_j, j \in J)$ is the Wardrop equilibrium. This type of equilibrium is often termed *adaptive* or *quasi-static*: on a short timescale arrivals at the network are adequately described by Poisson streams of rates $\nu_r, r \in R$, while the rates themselves adapt over a longer timescale.

Next consider the case where tolls are imposed to encourage traffic towards a system optimal flow pattern. One possibility is that queue $j$ could itself estimate $\rho_j$, the traffic through it, and hence estimate the correct toll $T_j(\rho_j)$. Such an approach has

been carefully investigated in the important paper of Gallager (1977), where an implementation particularly suited to data networks is developed, and adaptive convergence to the system optimal flow pattern is established. Another possibility is that queue $j$ might charge a toll $t_j(n)$ dependent on the number $n$ in queue $j$ just after a user has arrived, or a toll $t_j(D)$ dependent on the actual sojourn $D$ of the user in queue $j$. In either case it will be required that the expected toll be $T_j(\rho_j)$, in order that the system optimal flow pattern be encouraged. Now $n-1$ has the geometric distribution (3.1), and $D$ has an exponential distribution with parameter $(\phi_j - \rho_j)$. Hence we can deduce that the linear tolls

$$t_j(n) = \rho_j n / \phi_j (\phi_j - \rho_j), \quad n = 1, 2 \dots, \tag{3.8a}$$

$$t_j(D) = \rho_j D / (\phi_j - \rho_j), \quad D \geqslant 0, \tag{3.8b}$$

lead to the correct expected tolls. These linear tolls also have the property that they charge a user precisely the externality he causes: for example, the toll (3.8a) can be shown to be the mean additional delay caused to other users of queue $j$ if an additional user is added to queue $j$ bringing its queue size to $n$. Note that linear tolls may be easy to enforce: for example, the toll (3.8b) corresponds to charging each user a flat rate $\rho_j / (\phi_j - \rho_j)$ per unit time he spends in queue $j$. Against this it can be argued that the tolls depend on $\rho_j$ and $\phi_j$; while the queue may know its service rate $\phi_j$, if it is forced to estimate $\rho_j$ then a statistically more efficient feedback signal to users would be the average toll $T_j(\rho_j)$.

Do there exist tolls which do not require queue $j$ to estimate mean levels of traffic? More formally, are there tolls $t_j(n)$ or $t_j(D)$, possibly dependent on $\phi_j$, which have expectation $T_j(\rho_j)$ for all values of $\rho_j$, but which do not depend on $\rho_j$? In fact

$$t_j(n) = n(n-1)/2\phi_j, \quad n = 1, 2 \dots, \tag{3.9a}$$

$$t_j(D) = \tfrac{1}{2}\phi_j D^2 - D, \quad D \geqslant 0, \tag{3.9b}$$

are the unique functions of $n$ or $D$, respectively, which have expectation $T_j(\rho_j)$ for every $\rho_j > 0$. That they have the correct expectations follows by calculation from the geometric and exponential distribution of $n-1$ and $D$, respectively; that they are unique follows since the geometric and exponential distributions are complete (Lehmann 1986). The toll (3.9b) may be negative, although it plus the delay $D$ will be positive. The expected delay on finding $n$ is $n/\phi_j$, and this plus the toll (3.9a) gives the total $n(n+1)/2\phi_j$ proposed by Whittle (1985, 1986). It is interesting to note how rapidly these tolls grow with $n$ or $D$; the quadratic functions (3.9) contrast markedly with the linear functions (3.8).

The multiclass network of $\cdot/M/1$ queues which has led to the forms (3.8) and (3.9) is a very special case, but it is certainly possible to extend the results further. We might, for example, want the delay through a queue to be distributed more like a normal than an exponential random variable. But a series of $m \cdot/M/1$ queues produces a delay $D$ distributed as a gamma random variable with parameters $m$ and $\phi_j - \rho_j$: for such a system

$$T_j(\rho_j) = m\rho_j(\phi_j - \rho_j)^{-2},$$

and

$$t_j(D) = \phi_j D^2 / (m+1) - D, \quad D \geqslant 0,$$

is the unique function of $d$ with expectation $T_j(\rho_j)$ for every traffic level $\rho_j$.

Throughout this section we have been concerned with adaptive routing; many interesting questions arise in connection with *dynamic* routing, where a user arriving at the network may have access to some or full information about the current

network state (see, for example, Laws 1990, 1991). The distinction between adaptive and dynamic routing is not clearcut (see, for example, Kelly 1990), but one point is worth making here. If the system is striving to optimize the wrong function, then providing it with help, in the form of additional information, may damage performance. We illustrate this with an adaptation of Braess's paradox. Suppose the link between nodes 2 and 3 in figure 1$b$ has a delay $\rho+10+X$, where $X$ is a random variable taking the values 0 and 30 with equal probability. The random variable $X$ may indicate the presence or absence of road works, for example. If users do not know the value of $X$ then the expected delay along the link will be $\rho+25$, and the Wardrop equilibrium will be as in figure 1$a$. Each car incurs a delay of 83. If all users know the value of $X$ then the equilibrium flow pattern will be either as in figure 1$a$, if $X=30$, or as in figure 1$b$, if $X=0$. The expected delay, averaged over the two possibilities, will be $\frac{1}{2}(83+92)=87\frac{1}{2}$. Providing all users with knowledge of $X$ increases everyone's delay!

Current developments in electronics make feasible the practical application of both route guidance and road pricing (van Vuren & Smart 1990; Hoffman 1991) and the above discussion further emphasizes the importance of considering these topics together. If a system is encouraged to strive more aggressively towards an implicit objective, it becomes even more important that the objective is not perverse.

## 4. Loss networks

In this section we describe the basic theory of a loss network. The classical example is a telephone network, and we shall phrase our discussion in terms of calls, circuits and routes. However, the readers will observe that our model applies more widely to systems in which before a request (which may be a call, or a task, or a customer) is accepted it is first checked that sufficient resources are available to deal with the request.

Consider then a network with links from a set $J$, and suppose that link $j \in J$ comprises $C_j$ circuits. A subset $r \subset J$ identifies a route. Calls requesting route $r$ arrive as a Poisson stream of rate $\nu_r$, and as $r$ varies it indexes independent Poisson streams. A call requesting route $r$ is blocked and lost if on any link $j \in r$ there is no free circuit. Otherwise the call is connected and simultaneously holds one circuit on each link $j \in r$ for the holding period of the call. The call holding period is randomly distributed with unit mean and independent of earlier arrival and holding times. Write $R$ for the set of possible routes. Set $A_{jr}=1$ if $j \in r$, and $A_{jr}=0$ otherwise. This defines a $0-1$ matrix $A=(A_{jr}, j \in J, r \in R)$.

Let $n_r(t)$ be the number of calls in progress at time $t$ on route $r$, and define the vectors $n(t)=(n_r(t), r \in R)$ and $C=(C_j, j \in J)$. Then the stochastic process $(n(t), t \geqslant 0)$ has a unique stationary distribution, and under this distribution $\pi(n)=P\{n(t)=n\}$ is given by

$$\pi(n) = G(C)^{-1} \prod_{r \in R} \frac{\nu_r^{n_r}}{n_r!}, \quad n \in \mathscr{S}(C), \tag{4.1}$$

where
$$\mathscr{S}(C) = \{n \in \mathbb{Z}_+^R : An \leqslant C\}, \tag{4.2}$$

and $G(C)$ is the normalizing constant (or partition function)

$$G(C) = \left( \sum_{n \in \mathscr{S}(C)} \prod_{r \in R} \frac{\nu_r^{n_r}}{n_r!} \right). \tag{4.3}$$

This result is easy to check in the case where holding times are exponentially distributed: then $(n(t), t \geqslant 0)$ is a Markov process and the distribution (4.1) satisfies the detailed balance conditions

$$\nu_r \pi(n) = (n_r + 1)\pi(n + e_r), \quad n, n + e_r \in \mathcal{S}(C),$$

where $e_r = (I[r' = r], r' \in R)$ is the unit vector describing just one call in progress on route $r$. In this form the result has been known for many years (see Brockmeyer *et al.* 1948): that the form (4.1) is insensitive to the holding time distributions is an example of the modern theory of insensitivity (see, for example, Whittle 1986).

Most quantities of interest can be written in terms of the distribution (4.1) or the partition function (4.3). For example let $L_r$ be the proportion of calls requesting route $r$ that are lost. Since the arrival stream of calls requesting route $r$ is Poisson,

$$1 - L_r = \sum_{n \in \mathcal{S}(C - Ae_r)} \pi(n) = G(C)^{-1} G(C - Ae_r). \tag{4.4}$$

Such simple explicit forms might be thought to provide the complete solution. However, this is far from the case. For all but the smallest networks it is impractical to compute $G$ directly: observe that the number of routes $|R|$ may grow as fast as exponentially with the number of nodes $|J|$, and that in the (otherwise trivial) case when $|R| = |J|$ and $A = I$ the size of the state space $|\mathcal{S}(C)| = \prod_{j \in J} C_j$ grows rapidly with the capacity limitations $C_j, j \in J$. The theory of computational complexity allows these remarks to be stated more formally. Louth (1991) has shown that, even in the restricted case where links have capacity 1 and arrival rates are equal, the task of computing the partition function (4.3) is #$P$-complete. Nevertheless, might a randomized algorithm, for example rejection sampling from the truncated Poisson distribution (4.1) or simulation of the underlying stochastic process, lead to an answer of any required accuracy, without the calculation growing exponentially with either the system size or the required accuracy? Again the answer is in general no: there can be no fully polynomial randomized approximation scheme for the task of computing the partition function (4.3), unless $RP = NP$ (Jerrum & Sinclair 1990; Louth 1991). Simulation is a valuable tool for many particular network structures, but these results indicate its limitations when dealing with arbitrary network topologies. A theme of much recent work has been to find approaches which complement computation and simulation with analytical insights.

Consider the problem of finding the most likely state $n$ under the probability distribution (4.1). This is equivalent to maximizing

$$\sum_r (n_r \ln \nu_r - \ln n_r!),$$

over $n \in \mathcal{S}(C)$, a problem which is complicated by the discrete nature of the state space. To simplify things replace $\ln n!$ by $n \ln n - n$ (recall that by Stirling's formula $\ln n! = n \ln n - n + O(\ln n)$) and replace the integer vector $n$ by a real vector $x$. The resulting problem is the following.

$$\text{Maximize} \quad \sum_r (x_r \ln \nu_r - x_r \ln x_r + x_r), \tag{4.5a}$$

$$\text{over} \quad x \geqslant 0, \tag{4.5b}$$

$$\text{subject to} \quad Ax \leqslant C. \tag{4.5c}$$

Observe that the objective function (4.5a) is differentiable and strictly concave over the cone $x \geqslant 0$ and tends to $-\infty$ as $\|x\| \to \infty$, and the feasible region (4.5c) is a

closed convex set. Hence a maximizing value of $x$ exists and is unique, and can be found by lagrangian methods (Whittle 1971). Consider, then, the lagrangian form

$$L(x, z; y) = \sum_r (x_r \ln \nu_r - x_r \ln x_r + x_r) + \sum_j y_j (C_j - \sum_r A_{jr} x_r - z_j)$$

$$= \sum_r x_r + \sum_r x_r (\ln \nu_r - \ln x_r - \sum_j y_j A_{jr}) + \sum_j y_j C_j - \sum_j y_j z_j,$$

where $z = (z_j, j \in J)$ is the vector of slack variables $z = C - Ax$ and $y = (y_j, j \in J)$ is a vector of Lagrange multipliers. To maximize $L(x, z; y)$ over the cone $x, z \geqslant 0$ we require that $y \geqslant 0, y \cdot z = 0$ and, differentiating with respect to $x_r$,

$$\ln \nu_r - \ln x_r - \sum_j y_j A_{jr} = 0.$$

The maximizing $x_r$ is then

$$\bar{x}_r(y) = \nu_r \exp\left(-\sum_{j \in r} y_j\right), \qquad (4.6)$$

and so

$$\max_{x, z \geqslant 0} L(x, z; y) = \sum_r \bar{x}_r(y) + \sum_j y_j C_j$$

$$= \sum_r \nu_r \exp\left(-\sum_{j \in r} y_j\right) + \sum_j y_j C_j.$$

Hence the lagrangian dual to the primal problem is the following.

$$\text{Minimize} \quad \sum_j \nu_r \exp\left(-\sum_{j \in r} y_j\right) + \sum_j y_j C_j, \qquad (4.7a)$$

$$\text{over} \quad y \geqslant 0. \qquad (4.7b)$$

We may solve the primal problem (4.5) by choosing values for the Lagrange multipliers $y = \bar{y}$ so that $\bar{x}(\bar{y}), \bar{y}$ are primal and dual feasible,

$$\bar{x}(\bar{y}) \geqslant 0, \quad \bar{z} = C - A\bar{x}(\bar{y}) \geqslant 0, \quad \bar{y} \geqslant 0 \qquad (4.8a)$$

and satisfy the complementary slackness conditions

$$\bar{y} \cdot \bar{z}(\bar{y}) = 0. \qquad (4.8b)$$

It is interesting to rewrite these conditions in terms of transformed variables

$$B_j = 1 - \exp(-y_j). \qquad (4.9)$$

Under this transformation the conditions (4.8) on $\bar{y}$ become the following conditions on $B = (B_j, j \in J)$.

$$\sum_{r: j \in r} \nu_r \prod_{i \in r} (1 - B_i) = C_j \quad \text{if} \quad B_j > 0 \qquad (4.10a)$$

$$\leqslant C_j \quad \text{if} \quad B_j = 0, \qquad (4.10b)$$

$$B_1, B_2, \ldots, B_J \in [0, 1). \qquad (4.10c)$$

The convexity properties of the primal problem (4.5) imply that there exist Lagrange multipliers $\bar{y}$ satisfying (4.8), and hence that there exists $B$ satisfying (4.10). Alternatively, observe that the objective function of the dual problem (4.7a) is differentiable and convex over the cone $y \geqslant 0$ and tends to $\infty$ as $\|y\| \to \infty$. Hence an optimum $\bar{y}$ exists; differentiation of the dual objective function with respect to $y_j$

establishes a one-to-one correspondence under the transformation (4.9) between optima of the dual problem and solutions $B$ to conditions (4.10). Finally observe that the mapping $y \mapsto yA$ is one-to-one from the set $y \geqslant 0$ if $A$ has rank $|J|$. The objective function of the dual problem (4.7) is thus strictly convex if $A$ has rank $|J|$. Hence the optimum $\bar{y}$ is unique if $A$ has rank $|J|$.

In summary, we have shown that there exists a unique optimum to the primal problem (4.5), and that it can be expressed in the form

$$x_r = \nu_r \prod_{j \in r} (1 - B_j), \quad r \in R, \tag{4.11}$$

where $B = (B_j, j \in J)$ is any solution to the conditions (4.10) on $B$. There always exists a solution to the conditions on $B$, and it is unique if $A$ has rank $|J|$. There is a one-to-one correspondence between solutions to the conditions (4.10) on $B$ and optima of the dual problem (4.7), given by the transformation (4.9).

Conditions (4.10) have a straightforward interpretation in terms of a continuous, or fluid, flow. Suppose that an offered flow of $\nu_r$ on route $r$ is thinned by a factor $(1 - B_i)$ on each link $i \in r$ so that a flow of

$$\nu_r \prod_{i \in r} (1 - B_i) \tag{4.12}$$

remains. Then conditions (4.10) state that at any link $j$ for which $B_j > 0$ the total capacity of that link, $C_j$, must be completely utilized by the superposition over routes $r$ through link $j$ of the flows (4.12). Conversely no thinning of flow is allowed at a link which is not full.

We have made some approximations in the formulation of the primal problem (4.5), and so the reader may well ask: What, precisely, is the connection between the simple form (4.11) and the distribution (4.1)? This connection we now outline.

The distribution (4.1) is that of $|R|$ independent Poisson random variables conditioned on a collection of linear inequality constraints (4.2). It is natural to look for a limit theorem as the capacities $C_j, j \in J$, and the offered traffics $\nu_r, r \in R$, are increased together (with ratios $C_j/\nu_r$ held fixed). We would expect the distribution (4.1) to approach that of $|R|$ independent normal random variables conditioned on a collection of linear inequality constraints. Limit results of this form are familiar when the linear inequalities are replaced by linear equalities, and arise naturally in the analysis of contingency tables (see, for example, Haberman 1974). The primal problem (4.5) and its solution (4.11) simply establish the form of the centring term in the expected central limit theorem (Kelly 1986; Hunt & Kelly 1989; Hunt 1990).

The central limit theorem implies, as a form of law of large numbers, that

$$L_r \to 1 - \prod_{j \in r} (1 - B_j), \quad r \in R, \tag{4.13}$$

where $B = (B_j, j \in J)$ is a solution to the dual problem (4.7). Here $1 - L_r$ is the exact acceptance probability, given by expression (4.4), and the limit is as capacities $C$ and offered traffics $\nu$ are increased together, with ratios held fixed. Under this limiting régime it is *as if* links block independently, link $j$ blocking with probability $B_j$.

The law of large numbers (4.13) has great appeal as a theoretical limit, but it has disappointing accuracy as a general approximation. In practice $L_r$ is often approximated by the form

$$L_r \approx 1 - \prod_{j \in r} (1 - B_j), \quad r \in R, \tag{4.14}$$

but where $B_j, j \in J$, solve the nonlinear equations

$$B_j = E(\rho_j, C_j), \quad j \in J, \tag{4.15a}$$

$$\rho_j = \sum_{r : j \in r} \nu_r \prod_{i \in r - \{j\}} (1 - B_i), \quad j \in J. \tag{4.15b}$$

Here the function $E(\cdot, \cdot)$ is defined for scalar $\nu$ and $C$ by

$$E(\nu, C) = \frac{\nu^C}{C!} \left[ \sum_{n=0}^{C} \frac{\nu^n}{n!} \right], \tag{4.16}$$

and is thus just Erlang's celebrated formula for the proportion of calls lost at single link of capacity $C$ circuits offered Poisson traffic at rate $\nu$. The idea underlying the approximation is simple to explain. Suppose that a Poisson stream of rate $\nu_r$ is thinned by factor $1 - B_i$ at each link $i \in r - \{j\}$ before being offered to link $j$. If these thinnings could be assumed independent both from link to link and over all routes passing through link $j$ (they clearly are not), then the traffic offered to link $j$ would be Poisson at rate (4.15b), the blocking probability at link $j$ would be (4.15a), and the loss probability on route $r$ would satisfy (4.14) exactly.

A solution $B = (B_j, j \in J)$ to equations (4.15) exists, by the Brouwer fixed point theorem, but is it unique? We shall answer this question, and relate the approximation scheme (4.14)–(4.15) to our earlier limit results, by consideration of an appropriate optimization problem.

Define a function $U(y, C)$ by the implicit relation

$$U(-\ln (1 - E(\nu, C)), C) = \nu(1 - E(\nu, C)). \tag{4.17}$$

Observe that as $\nu$ increases from 0 to $\infty$ the first argument of $U$ increases from 0 to $\infty$ and so this implicit relation defines a function $U : \mathbb{R}_+ \times \mathbb{Z}_+ \to \mathbb{R}_+$. Indeed, for a single link of capacity $C$ circuits the quantity $U(y, C)$ is just the mean number of circuits in use (the *utilization*) when the blocking probability is $1 - \exp(-y)$. Thus $U(y, C)$ is a strictly increasing function of $y$. Consider the following variant of the problem (4.7), which we shall call the revised dual problem.

$$\text{Minimize} \quad \sum_r \nu_r \exp \left( -\sum_{j \in r} y_j \right) + \sum_j \int_0^{y_j} U(z, C_j) \, dz, \tag{4.18a}$$

$$\text{over} \quad y \geqslant 0. \tag{4.18b}$$

Since $U(y, C)$ is a strictly increasing function of $y$, $\int_0^y U(z, C) \, dz$ is a strictly convex function of $y$. Hence the objective function (4.18a) is strictly convex, with a unique minimum. The objective function is also differentiable: hence the stationarity conditions

$$\sum_{r : j \in r} \nu_r \exp \left( -\sum_{i \in r} y_i \right) = U(y_j, C_j), \quad j \in J, \tag{4.19}$$

obtained by differentiating the objective function with respect to $y_j, j \in J$, locate a unique vector $y \geqslant 0$. Now suppose that $(B_j, j \in J)$ is a solution to equations (4.15). Under the equivalence $B = 1 - e^{-y}$ and using the definition (4.16), these equations become precisely equations (4.19). We deduce that equations (4.15) have a unique solution, given in terms of the optimum $y$ of the problem (4.18) by the transformation (4.9).

The conditions (4.10) insist that the carried traffic on a link must equal capacity

before the blocking probability on that link can be positive. Conditions (4.19) are a natural relaxation: as carried traffic approaches capacity, blocking increases, in a manner corresponding to the utilization function $U$. Note that replacing the function $U$ by a function

$$U_f(z, C) = C, \quad z > 0,$$

reduces the revised dual problem (4.18) to the dual problem (4.7). The utilization function $U_f$ would be natural for fluid flow, where if there is any blocking then all capacity is in use. For $C$ large there is not much difference between $U$ and $U_f$; under the limiting régime considered earlier, where $C$ and $\nu$ increase together, the objective function (4.18 $a$) approaches a scaled version of the objective function (4.7 $a$). For non-limiting values of $\nu$ and $C$ the approximation (4.14) is generally much improved when the vector $B$ is defined, via the transformation (4.9), in terms of the optimum to the revised dual problem (4.18) rather than an optimum to the dual problem (4.7). This is intuitively plausible, since, as we have seen, the various equivalent formulations (4.15), (4.18) and (4.19) reflect the fact that as a link's utilization approaches its capacity, the link's blocking increases smoothly. As might be further expected from the informal motivation of the fixed point equations (4.15) in terms of thinned and superimposed Poisson streams, the approximation scheme (4.14)–(4.15) becomes increasingly accurate the more diverse the collection of routes passing through each link (for reviews, see Whitt 1985; Ziedinš 1987; Kelly 1991).

Thus the loss network, defined earlier in terms of Poisson arrival streams and certain rules for accepting calls, can also be viewed as a system implicitly attempting to solve the optimization problems (4.7) or (4.18). It is amusing to note that the term

$$\sum_r \nu_r \exp\left(-\sum_{j \in r} y_j\right) = \sum_r \nu_r \prod_{j \in r} (1 - B_j),$$

appearing in both objective functions corresponds to the average level of traffic carried by the network; we would much prefer the network to be implicitly *maximizing* this term!

The loss network so far considered is a rather simple one, involving just *fixed* routing: if a call fails to be accepted on its fixed route then it is lost. In practice, networks often attempt to improve performance by allowing *alternative* routing, where a call which is blocked on a route may try again on an alternative route. We now consider a simple example of alternative routing in a fully connected network. Suppose that $K$ nodes are linked to form a complete graph. Between any pair of nodes calls arise at rate $\nu$, and there is a link of capacity $C$. If there is a spare circuit on the link joining the end points of a call then the call is accepted and carried by that circuit. Otherwise the call chooses at random a two-link path joining its end points: the call is accepted on that path if both links have a spare circuit, and is lost otherwise.

The generalization of equations (4.15) is based on the same underlying approximation, that links block independently. Let $B$ be the link blocking probability, taken to be the same for each link. The probability that a call overflows from its first choice route is $B$, and the probability it can be accepted at the other link of a two-link alternative route is $1 - B$; the arrival rate of overflowing calls at a link is then $2\nu B(1 - B)$. We should look for a solution to the fixed point equation

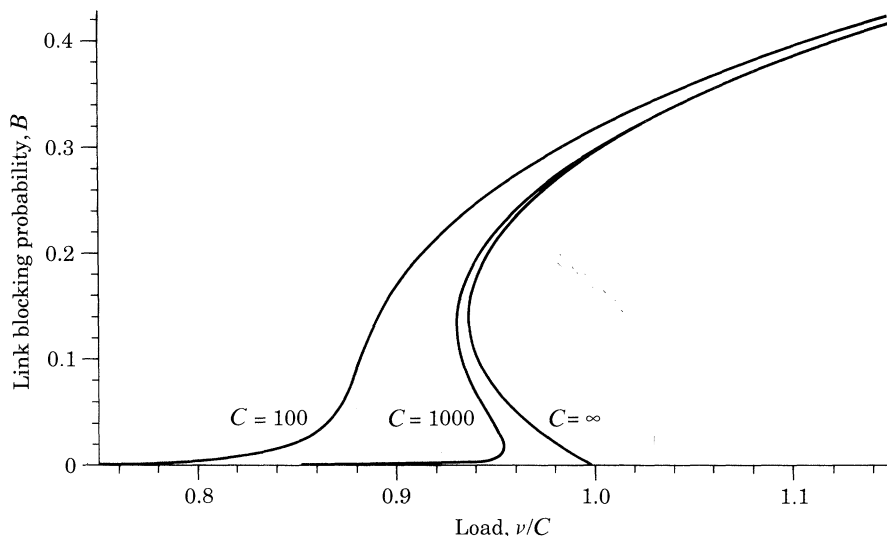$$B = E(\nu + 2\nu B(1 - B), C), \tag{4.20}$$

for $B$.

Figure 2. Instability of blocking probability.

The locus of points satisfying equation (4.20) is illustrated in figure 2. Observe the possibility of multiple solutions for $B$, when $C$ is large enough and for a narrow range of the ratio $\nu/C$. Are the multiple solutions evident in figure 2 a real phenomenon, or are they simply an artefact of the approximation? This question has been tackled by simulation of similar fully connected networks with tens or scores of nodes, and by rigorous analysis of the model as the number of nodes tends to infinity (Akinpelu 1984; Ackerley 1987; Gibbens *et al.* 1990; Crametz & Hunt 1991). The conclusions are clear. For networks with a moderate or large number of nodes the upper and lower solutions for $B$ in figure 2 correspond to distinct, locally stable modes of a stationary distribution. There is also a hysteresis effect: if $\nu$ is varied slowly the mode which obtains may depend not just on the current value of $\nu$ but also upon whether $\nu$ approached this value from above or below. An intuitive explanation is easy to provide. The lower solution corresponds to a mode in which blocking is low, calls are mainly routed directly and relatively few calls are carried on two-link paths. The upper solution corresponds to a mode in which blocking is high and many calls are carried over two-link paths. Such calls use two circuits each, and this additional demand on network resources may cause a substantial number of subsequent calls also to attempt two-link paths. Thus a form of positive feedback may keep the system in the high blocking mode.

It is interesting to reinterpret this discussion in the language of catastrophe theory (Poston & Stewart 1978), as we now briefly indicate. By differentiating the potential function

$$\nu\,\mathrm{e}^{-y} + \nu\,\mathrm{e}^{-2y}(1 - \tfrac{2}{3}\mathrm{e}^{-y}) + \int_0^y U(z, C)\,\mathrm{d}z, \tag{4.21}$$

we find that it is stationary with respect to $y$ when

$$\nu\,\mathrm{e}^{-y} + 2\nu\,\mathrm{e}^{-2y}(1 - \mathrm{e}^{-y}) = U(y, C). \tag{4.22}$$

However, under the equivalence $B = 1 - \mathrm{e}^{-y}$ and using the definition (4.17), the equation (4.22) becomes precisely the equation (4.20). Thus the solutions illustrated

in figure 2 locate the stationary points of the potential function (4.21). If we regard $\nu/C$ as the normal variable and $C$ as the splitting variable, then figure 2 illustrates three cross sections of the cusp catastrophe.

Comparing the potential function (4.21) with the objective function (4.18*a*) of the revised dual problem, we see that alternative routing has led to the introduction of the second term in expression (4.21), and hence to non-convexity and multiple minima. The system is again implicitly minimizing a function, and bistability is to be expected when the function has two local minima.

For the parameter choice $(\nu, C) = (95, 100)$ the network loss probability, $B[1-(1-B)^2]$, takes a value of about 0.12. If alternative routing is not allowed, so that a call blocked on its direct link is lost, then the network loss probability is given by Erlang's formula (4.16) to be 0.05. Thus allowing a blocked call to attempt a two-link alternative route may *increase* the loss probability of the network. We have seen that this is plausible, since if a link accepts an alternatively routed call then it may later have to block a directly routed call, which will then attempt to find two circuits elsewhere in the network. A natural response is to allow a link to reject alternatively routed calls if the number of idle circuits on the link is less than or equal to a certain value, $t$, say. This method of giving priority at a link to certain traffic streams is known as *trunk reservation*, and the parameter $t$ is called the trunk reservation parameter for the link. With trunk reservation in place alternative routing is capable of improving performance, and trunk reservation is widely used in telephone networks (Songhurst 1980).

## 5. Adaptive routing in loss networks

We have seen in the last section that the various potential functions implicitly minimized bear little direct relation to the network performance criteria of interest to system designers. In this section we show that an explicit consideration of such criteria can lead to decentralized adaptive routing schemes which are at least attempting to optimize the right function.

Suppose that each call carried on route $r$ generates an expected revenue $w_r$ (or, equivalently, interpret $w_r$ as the cost of losing a call on route $r$). Then, under the fixed point model (4.15), the rate of return from the network will be

$$W(\nu; C) = \sum_r w_r \lambda_r,$$

where

$$\lambda_r = \nu_r \prod_{j \in r} (1 - B_j),$$

corresponds to the traffic carried on route $r$. We use the notation $W(\nu; C)$ to emphasize the dependence of $W$ on the vectors of offered traffics $(\nu_r, r \in R)$ and capacities $(C_j, j \in J)$. Let

$$\delta_j = \rho_j(E(\rho_j, C_j - 1) - E(\rho_j, C_j)). \tag{5.1}$$

Extend the definition (4.16) to non-integral values of scalar $C$ by linear interpolation, and at integer values of $C_j$ define the derivative of $W(\nu; C)$ with respect to $C_j$ to be the left derivative. Then it is possible to prove (Kelly 1988) that

$$\frac{\mathrm{d}}{\mathrm{d}\nu_r} W(\nu; C) = s_r \prod_{j \in r} (1 - B_j), \tag{5.2}$$

and
$$\frac{\mathrm{d}}{\mathrm{d}C_j} W(\nu;C) = c_j, \tag{5.3}$$

where $s = (s_r, r \in R)$ and $c = (c_j, j \in J)$ are the unique solution to the linear equations

$$s_r = w_r - \sum_{j \in r} c_j, \tag{5.4}$$

$$c_j = \delta_j \sum_{r:j \in r} \lambda_r(s_r + c_j) \Big/ \sum_{r:j \in r} \lambda_r. \tag{5.5}$$

We can interpret $s_r$ as the *surplus value* of a call on route $r$: if such a call is accepted it will earn $w_r$ directly but at an *implied cost* of $c_j$ for each circuit used from link $j$. The implied costs $c$ measure the expected knock-on effects of accepting a call upon later arrivals at the network. From (5.3) it follows that $c_j$ is also a *shadow price*, measuring the sensitivity of the rate of return to the capacity $C_j$ of link $j$. The local character of equations (5.4) and (5.5) is striking. The right-hand side of (5.4) involves costs $c_j$ only for links $j$ on the route $r$, while (5.5) exhibits $c_j$ in terms of an average, weighted over just those routes through link $j$, of $s_r + c_j$.

Expression (5.1) for $\delta_j$ is called Erlang's improvement formula (Brockmeyer *et al.* 1948). Observe that $\delta_j$ is simply the increase in the rate at which calls are blocked if a single link offered Poisson traffic at rate $\rho_j$ has its capacity reduced by one circuit, and that $\delta_j$ increases from zero to one as $\rho_j$ increases from zero to infinity.

Next we describe how equations (5.4) and (5.5) can be used as the basis for decentralized control of routing through the network. As motivation it is helpful to think in terms of the following model of a distributed computation. Suppose there is limited intelligence in the form of arithmetical processing ability available for each link $j$ and for each route $r$. This intelligence may be located centrally or it may be distributed over the nodes of the network; for example the processing for route $r$ might be carried out by the source node for calls on route $r$. Suppose also that there is the possibility of limited communication between the intelligences of link $j$ and route $r$ provided $j \in r$. Consider now equations (5.4) and (5.5). One method for attempting a solution to these equations is repeated substitution. Choose a vector $c$; substitute it in equation (5.4) to obtain a vector $s$; substitute these into equation (5.5) to obtain a revised vector $c$, and repeat. This computation can, however, be distributed over the intelligences of links and routes, since equation (5.4) for $s_r$ involves implied costs $c_j$ only for links $j$ on the route $r$, while equation (5.5) for $c_j$ involves only surplus values $s_r$ for routes $r$ passing through link $j$. Kelly (1988) showed that if

$$\sum_{i \in r-\{j\}} \delta_i < 1, \quad j \in r, \quad r \in R, \tag{5.6}$$

then repeated substitution converges to the unique solution of the equations. When the condition (5.6) is not satisfied it may be necessary to damp the repeated substitution to obtain convergence. This can still be implemented by a distributed computation, but the individual intelligences may require some knowledge of the network beyond that locally available, to damp sufficiently the repeated substitution. In fact it is enough for the intelligences to know just one item of global information, namely $J$, the total number of links in the network.

The quantities $\delta_j$ and $\lambda_r$ appearing in equations (5.4) and (5.5) are not fixed and known. However, they can be estimated by intelligences of links and routes from, for

example, local measurements of carried loads. The estimates can then be used in a distributed computation of the vector $s$. Finally the derivatives (5.2) can be used to implement a decentralized hill-climbing search procedure able to adapt routing patterns in response to changes in the demands on the network.

The inequality (5.6) limits the size of knock-on effects and is a form of light traffic condition. If it is satisfied then the fixed point model (4.15) has a certain stability property: any perturbation of the capacity $C_j$ of link $j$ may cause a change in the entire vector $(\lambda_r, r \in R)$ of carried traffics, but the change in component $\lambda_r$ diminishes rapidly with the extent of the separation between route $r$ and link $j$. Without condition (5.6) perturbations may have influence over arbitrarily great distances. The existence of such effects is *not* an artefact of the approximation: examples have been deduced from the exact stationary distribution (4.1). The discussion of these effects is complicated by *frustration*, the crucial property that lies at the heart of the spin glass problem of statistical mechanics, whereby chains of influence throughout the network may be out of phase and compete with one another. The reader is referred to Kelly (1991) for further discussion and references.

It is possible to define implied costs and surplus values for fixed point models of alternative routing and trunk reservation, and to show that they solve linear relations generalizing (5.4) and (5.5) (see Key 1988; Key & Whitehead 1988; Kelly 1990). The potential for long-range order and instability is more pronounced in networks with alternative routing, since the chains of influence along different paths tend to reinforce one another. The possibility emphasizes the importance of treating a network as a whole: the local benefits of a capacity or routing change may be completely overwhelmed by adverse consequences elsewhere in the network.

## 6. Dynamic Alternative Routing

The routing scheme discussed in the last section attempts to control the network by explicitly computing average traffics, implied costs and surplus values, and deducing derivatives of performance criteria. The scheme resembles the adaptive schemes of §3, and indeed the implied costs precisely parallel the tolls discussed here. Might it be possible to design a network so that it instead resembles the electrical network of §2, with simple rules for call acceptance and routing leading naturally to good behaviour? Early work on this question has led to a scheme, Dynamic Alternative Routing (DAR), now being implemented in British Telecom's British trunk network.

DAR is a simple but effective dynamic routing strategy, which is decentralized and uses only local information. Its definition for a fully connected network is as follows. Suppose there are $K$ nodes in the network, with the link $\{i, j\}$ joining nodes $i$ and $j$ having capacity $C_{ij}$. Each link is assigned a trunk reservation parameter $t_{ij}$, and each source–destination pair $(i, j)$ stores the identity of its current tandem $k(i, j)$ for use in two-link alternative routes. A call between nodes $i$ and $j$ is first offered to the direct link and a call is always routed along that link if there is a free circuit. Otherwise, the call attempts the two-link alternative route via tandem node $k$ with trunk reservation applied to both links. If the call fails to be routed via $k$, this call is lost and, further, the identity of the tandem node is reselected (at random perhaps) from the set $\{1, 2, \ldots, K\} - \{i, j\}$. Note especially that the tandem node is not reselected if the call is successfully routed on either the direct link or the two-link alternative route. Mees (1986) has coined the term *sticky random routing* to emphasize this

property of the scheme. In practice it has been found simpler to reselect a tandem node by cycling around a fixed random permutation; the point is that reselection is not based on any collected data, only the important information that a call has just failed.

Let $p_k(i,j)$ denote the long-run proportion of calls between $i$ and $j$ which are offered to tandem node $k$, and let $q_k(i,j)$ be the long-run proportion of those calls between $i$ and $j$ and offered to tandem node $k$ which are blocked. Then, under uniform reselection,

$$p_a(i,j)\,q_a(i,j) = p_b(i,j)\,q_b(i,j), \quad a, b \neq i, j.$$

Observe that this simple ergodic result is exact for either random reselection or reselection using a fixed permutation. More generally, suppose the DAR mechanism for reselection of the tandem node between $i$ and $j$ chooses node $k$ with long-run frequency $f_k$, where $\Sigma_{k \neq i,j} f_k = 1$. Then each selection of node $k$ is paired with a failed call via node $k$, and so

$$p_a(i,j)\,q_a(i,j) : p_b(i,j)\,q_b(i,j) = f_a : f_b, \quad a, b \neq i, j. \tag{6.1}$$

Observe that if the blocking $q_k(i,j)$ is high on the path through the tandem node $k$, then the proportion of overflow routed via node $k$ will be low. This gives some insight into the means by which the routing scheme implicitly adapts to overloads and failures.

For further insight consider next the effects of mismatches between traffics and capacities. Suppose that traffic is greater than capacity on some links, and less than capacity on others. Is it possible for the excess traffic from overloaded links to be assigned to alternative routes which do not themselves clash with one another; and, if so, is it possible for this to be achieved by a simple algorithm? To proceed with these questions, consider the following random graph problem. Choose $p \in (0, 1)$ and suppose the edges of a fully connected graph on $K$ nodes are independently coloured red with probability $p$ and white with probability $1 - p$. Let a *triangle* be a set of three edges joining each pair from a set of three nodes. Call a triangle *good* if it contains one red and two white edges. Let $P(K)$ be the probability that there exists a set of disjoint good triangles such that each red edge is contained in a triangle. If we interpret the red edges as overloaded links, then $P(K)$ is the probability that there exists a collection of non-overlapping alternative routes. Hajek (1987) has shown that if $p < \frac{1}{3}$ then $P(K) \to 1$ as $K \to \infty$. Hajek's methods are informative about the performance of algorithms as well as the structure of random graphs, and we outline one aspect of his proof. Consider the following very simple greedy algorithm. Suppose disjoint good triangles $T_1, T_2, \ldots, T_k$ have been found already and the algorithm has not yet stopped. If there is no remaining red edge, declare the algorithm successful and stop. Otherwise, call a triangle available (after $k$ steps) if it is a good triangle which is disjoint from $T_1, T_2, \ldots, T_k$. Choose a red edge $e$ at random from the remaining red edges. If no available triangle contains $e$, declare the algorithm unsuccessful and stop. Otherwise, choose $T_{k+1}$ at random from among the available triangles containing $e$. Hajek conjectures that this simple greedy algorithm is successful with probability approaching 1. He proves that when the algorithm stops it has covered at least a proportion $1 - \delta$ of the red edges with probability approaching 1 as $K \to \infty$, for any $\delta > 0$. He establishes the above result by consideration of a modified algorithm. Choose $\epsilon$ with $0 < \epsilon < \frac{1}{3} - p$. Independently consider each white edge and delete it with probability $\epsilon$. Now run the original algorithm (although now only available

triangles with non-deleted edges are counted). Then this modified algorithm is successful with probability approaching 1.

It is informative that a simple greedy algorithm performs so well for the above random graph problem. In many respects DAR resembles a dynamic version of the greedy algorithm. It also has a number of similarities with probabilistic hill-climbing techniques such as simulated annealing, and Gibbens (1988) and Mitra & Seery (1991) have discussed its use as an algorithm to investigate random graph and network routing problems. Consider the multi-commodity flow problem of routing $n(n-1)$ distinct flows between each ordered pair of $n$ nodes. As a linear program this problem has $O(n^3)$ variables, since each of $n(n-1)$ ordered node pairs has one direct route and $(n-2)$ two-link alternatives available, and $O(n^2)$ constraints, associated with the $\frac{1}{2}n(n-1)$ links and distinct flows (cf. Gibbens & Kelly 1990). The problem has a natural mapping onto a computing engine with $O(n^2)$ parallel processors. DAR tackles the nonlinear stochastic version of this problem using a probabilistic hill-climbing technique performed by $O(n^2)$ parallel processors, namely the occupancy levels of the $\frac{1}{2}n(n-1)$ links and the $n(n-1)$ tandem pointers $k(i,j)$.

The formula (6.1) allows the fixed point methods described in §4 to be extended to model DAR; these methods, together with simulation, have been used to investigate the performance of DAR under a wide range of failure and overload conditions. For an account of this work, and extensions of DAR to regular but non-fully connected networks, the reader is referred to Stacey & Songhurst (1987), Gibbens (1988), Gibbens *et al.* (1988), Gibbens & Kelly (1990), Gibbens & Turner (1991).

The routing scheme described in §5 has the advantage that it makes no special assumption about the underlying topology of the network, but it adapts slowly, in response to averaged values. DAR, in comparison, requires some regularity of network structure, but operates on a fast timescale, driven by call arrivals, and uses very simple rules. Because of the different timescales involved it is quite possible that the benefits of both schemes may be achievable: certainly Key (1988) has described an approach to the calculation of implied costs and shadow prices for a fully connected network using DAR, and has shown how the shadow prices can be used interactively to allocate capacity in a network. Current research is developing methods which use sticky random routing on a fast timescale, driven by call arrivals or failure events, with route sets, priorities and trunk reservation selected using implied costs calculated over the longer timescales necessary for averages to be estimated.

# References

Ackerley, R. G. 1987 Hysteresis-type behaviour in networks with extensive overflow. *Br. Telecom Technol. J.* **5**, 42–50.

Akinpelu, J. M. 1984 The overload performance of engineered networks with non-hierarchical routing. *AT&T Tech. J.* **63**, 1261–1281.

Anderson, P. W., Arrow, K. J. & Pines, D. (eds) 1988 *The economy as an evolving complex system.* Redwood City, California: Addison-Wesley.

Braess, D. 1968 Über ein Paradoxon aus der Verkehrsplanung. *Unternehmenforschung* **12**, 258–268.

Brockmeyer, E., Halstrom, H. L. & Jensen, A. 1948 *The life and works of A. K. Erlang.* Copenhagen: Academy of Technical Sciences.

Cohen, J. E. 1988 The counterintuitive in conflict and cooperation. *Am. Scient.* **76**, 576–584.

Cohen, J. E. & Kelly, F. P. 1990 A paradox of congestion in a queueing network. *J. appl. Prob.* **27**, 730–734.

Crametz, J.-P. & Hunt, P. J. 1991 A limit result respecting graph structure for a fully connected network with alternative routing. *Ann. Appl. Prob.* **1**.

Doyle, P. G. & Snell, J. L. 1984 *Random walks and electric networks*. Mathematical Association of America.

Erlang, A. K. 1925 A proof of Maxwell's law, the principal proposition in the kinetic theory of gases. In Brockmeyer *et al.* (1948), pp. 222–226.

Gallager R. G. 1977 A minimum delay routing algorithm using distributed computation. *IEEE Trans. Commun.* **25**, 73–85.

Gibbens, R. J. 1988 Dynamic routing in circuit-switched networks: the dynamic alternative routing strategy. Ph.D. thesis, University of Cambridge.

Gibbens, R. J., Kelly, F. P. & Key, P. B. 1988 Dynamic alternative routing – modelling and behaviour. *Proc. 12th Int. Teletraffic Congress, Turin* (ed. M. Bonatti). Amsterdam: Elsevier.

Gibbens, R. J., Hunt, P. J. & Kelly, F. P. 1990 Bistability in communication networks. In *Disorder in physical systems* (ed. G. R. Grimmett & D. J. A. Welsh), pp. 113–128. Oxford: Oxford University Press.

Gibbens, R. J. & Kelly, F. P. 1990 Dynamic routing in fully connected networks. *IMA J. Math. Cont. Information* **7**, 77–111.

Gibbens, R. J. & Turner, S. R. E. 1991 Sticky random routing in dual-parented networks. *Eighth UK Teletraffic Symposium, Nottingham*. London: IEE.

Haberman, S. J. 1974 *The analysis of frequency data*. University of Chicago Press.

Hajek, B. 1987 Average case analysis of greedy algorithms for Kelly's triangle problem and the independent set problem. In *26th IEEE Conference on Decision and Control*. New York: IEEE.

Hoffman, G. 1991 Up-to-the-minute information as we drive – how it can help road users and traffic management. *Transport Rev.* **11**, 41–61.

Hunt, P. J. 1990 Limit theorems for stochastic loss networks. Ph.D. thesis, University of Cambridge.

Hunt, P. J. & Kelly, F. P. 1989 On critically loaded loss networks. *Adv. appl. Prob.* **21**, 831–841.

Jerrum, M. R. & Sinclair, A. 1990 Polynomial time approximation algorithms for the Ising model. Department of Computer Science, Edinburgh.

Kelly, F. P. 1979 *Reversibility and stochastic networks*. Chichester: Wiley.

Kelly, F. P. 1986 Blocking probabilities in large circuit-switched networks. *Adv. appl. Prob.* **18**, 473–505.

Kelly, F. P. 1988 Routing in circuit-switched networks: optimization, shadow prices and decentralization. *Adv. appl. Prob.* **20**, 112–144.

Kelly, F. P. 1990 Routing and capacity allocation in networks with trunk reservation. *Math. oper. Res.* **15**, 771–793.

Kelly, F. P. 1991 Loss networks. *Ann. appl. Prob.* **1**, 319–378.

Key, P. B. 1988 Implied cost methodology and software tools for a fully connected network with DAR and trunk reservation. *Br. Telecom Technol. J.* **6**, 52–65.

Key, P. B. & Whitehead, M. J. 1988 Cost-effective use of networks employing Dynamic Alternative Routing. In *Proc. 12th Int. Teletraffic Congress, Turin* (ed. M. Bonatti). Amsterdam: Elsevier.

Kingman, J. F. C. 1969 Markov population processes. *J. appl. Prob.* **6**, 1–18.

Knödel, W. 1969 *Graphentheoretische Methoden und ihre Anwendungen*, pp. 56–59. Berlin: Springer-Verlag.

Langton, C. G. (ed.) 1989 *Artificial life*. Redwood City, California: Addison-Wesley.

Laws, C. N. 1990 Dynamic routing in queueing networks. Ph.D. thesis, University of Cambridge.

Laws, C. N. 1992 Resource pooling in queueing networks with dynamic routing. *Adv. appl. Prob.* **24**. (In the press.)

Lehmann, E. L. 1986 *Testing statistical hypotheses*, 2nd edn. New York: Wiley.

Louth, G. M. 1991 Stochastic networks: complexity, dependence and routing. Ph.D. thesis, University of Cambridge.

Mees, A. I. 1986 Simple is the best for dynamic routing of telecommunications. *Nature, Lond.* **323**, 108.

Mitra, D. & Seery, J. B. 1991 Comparative evaluations of randomized and dynamic routing strategies for circuit-switched networks. *IEEE Trans. Commun.* **39**, 102–116.

Nagurney, A. 1987 Competitive equilibrium problems, variational inequalities and regional science. *J. Regional Sci.* **27**, 503–517.

*New York Times* 1990 What if they closed 42nd Street and nobody noticed? 25 December, p. 38.

Pigou, A. C. 1920 *The economics of welfare*. London: Macmillan.

Pines, D. (ed.) 1987 *Emerging synthesis in science*. Redwood City, California: Addison-Wesley.

Poston, T. & Stewart, I. 1978 *Catastrophe theory and its applications*. London: Pitman.

Potts, R. B. & Oliver, R. M. 1972 *Flows in transportation networks*. New York: Academic.

Songhurst, D. J. 1980 Protection against traffic overload in hierarchical networks employing alternative routing. *Telecommunications Network Planning Symposium, Paris*.

Stacey, R. R. & Songhurst, D. J. 1987 Dynamic Alternative Routing in the British Telecom trunk network. *Int. Switching Symposium, Phoenix, Arizona*.

Thomson, W. & Tait, P. G. 1879 *Treatise on natural philosophy*. Cambridge.

van Vuren, T. & Smart, M. B. 1990 Route guidance and road pricing – problems, practicalities and possibilities. *Transport Rev.* **10**, 269–283.

Walters, A. A. 1961 The theory and measurement of private and social cost of highway congestion. *Econometrica* **29**, 676–699.

Wardrop, J. G. 1952 some theoretical aspects of road traffic research. *Proc. Inst. Civil Engng* **1**, 325–378.

Whitt, W. 1985 Blocking when service is required from several facilities simultaneously. *AT&T Tech. J.* **64**, 1807–1856.

Whittle, P. 1971 *Optimization under constraints*. London: Wiley.

Whittle, P. 1985 Scheduling and characterization problems for stochastic networks (with Discussion). *Jl R. statist. Soc.* **47**, 407–428.

Whittle, P. 1986 *Systems in stochastic equilibrium*. Chichester: Wiley.

Wroe, G. A., Cope, G. A. & Whitehead, M. J. 1990 Flexible routing in the BT international access network. *Seventh UK Teletraffic Symposium, Durham*. London: IEE.

Ziediņš, I. B. 1987 Stochastic models of traffic in star and line networks. Ph.D. thesis, University of Cambridge.