CHAPTER 11

# Large Deviation and Fluid Approximations in Control of Stochastic Systems

## R. Weber

*University of Cambridge, U.K.*

Large deviation and fluid approximations can provide valuable insight to the behaviour of stochastic systems. This paper summarises some of the key ideas and discusses a number of examples, including a probabilistic analysis of the first-fit decreasing bin packing algorithm, and comments on ideas of Whittle relating to risk-sensitive control of Markov processes. In a final example, a large deviations analysis is used to estimate the frequency of buffer overflow when a number of stationary traffic sources are routed through a switch with a finite buffer. We show that the buffer overflow frequency is asymptotically less that a specified amount if and only if the sum of the effective bandwidths for the different traffic sources is less than the bandwidth of the switch and we suggest a simple approximation for these effective bandwidths.

## 11.1   LARGE DEVIATIONS

The analysis of the large deviation behaviour of a stochastic system can often provide valuable insight to the system's evolution and control. The focus is on the deterministic fluid model. One seeks to estimate the accuracy of the fluid model and to understand the most likely way that the system might depart from the fluid path. Large deviation analysis is less refined than equilibrium or heavy-traffic analysis, but it is attractive in its generality.

Large deviation results have been used in numerous areas, including statistics, physics, and simulation. The rich and beautiful theory is explained in books

by Freidlin and Wentzell (1984), Bucklew (1990) and Dembo and Zeitouni (1993). This section begins with a brief overview, that is intuitive and tutorial. In Sections 11.2–4 we describe some applications under the headings of likely, near-likely and unlikely paths. The final example describes the use of large deviation theory in determining effective bandwidths for bursty traffic sources which share a finitely buffered queue.

**Intuitions**   A large Deviation Principle (LDP) is said to hold when for a family of probability measures, $\mu_V$, indexed by $V$, and defined on the same probability space, we have that for any event $A$

$$\mu_V(A) \sim_L \exp\left(- V \inf_{a\in A} I(a)\right). \tag{1}$$

The parameter $V$, is some measure of system size: for example, the number of particles in the system, the number of observations, or the size of the space in which the system resides. The idea is that the probability of a rare event becomes small as the system size increases. Whittle (1990a) uses the symbol $\sim_L$ to mean 'is asymptotically logarithmically equal to', and this is shorthand for the statement

$$- \inf_{x\in A^\circ} I(x) \leqslant \varliminf_{V\to\infty} \frac{1}{V}\log \mu_V(A^\circ) \leqslant \varlimsup_{V\to\infty} \frac{1}{V}\log \mu_V(\bar{A}) \leqslant - \inf_{x\in\bar{A}} I(x),$$

where the rate function, $I(\cdot)$, is called 'a good rate function' if $\{x:I(x)\leqslant \alpha\}$ is a compact set for all $\alpha < \infty$. In this case we say that $\mu_V$ *obeys an LDP with good rate function* $I(\cdot)$, noticing that this requires separate statements for open and closed sets, here taken as $A^\circ$ and $\bar{A}$, the interior and closure of $A$, respectively.

As an example, suppose $A$ is the event that the number of customers in a stable $M/M/1$ queue exceeds $V$ during a busy period. The queue length can take a variety of different paths on the way to reach $V$; some are more likely than others. Intuitively, equation (1) says that if $A$ occurs then it does so in the most likely way. This can be illuminating, although often the required variational problem is not easy to solve. In summary, the intuitions of large deviation theory are that: *The probability of a rare event decreases to zero exponentially fast as system size increases and this probability can be estimated by the probability of the most likely way the event can occur; given that it has occurred then it almost certainly occurred in the most likely of unlikely ways.*

**The Gärtner–Ellis theorem**   To understand large deviations theory, it is helpful to realise that many of the key results can be derived, at least formally, from the Gärtner–Ellis theorem. Knowing this theorem, one can recall results and search for new ones. The theorem makes a statement about the convergence of a dependent sequence of random vectors, $\{Z_1, Z_2, \ldots\}$ in $\mathscr{R}^d$. It holds under the following assumptions.

1. *The asymptotic logarithmic moment generating function,*

$$\phi(\theta) = \lim_{n \to \infty} n^{-1} \log E[\exp(n\theta^T Z_n)],$$

*exists for all $\theta$, possibly as $\pm \infty$. The set $\{\theta : \phi(\theta) \leqslant k\}$ is closed for every finite $k$.*
2. *The origin is in the interior of the effective domain, defined as $D_\phi = \{\theta : \phi(\theta) < \infty\}$.*
3. *The derivative, $\phi'(\theta)$, exists in the interior of $D_\phi$ and tends to infinity as $\theta$ approaches the boundary of $D_\phi^\circ$.*

**Theorem 1 (Gärtner-Ellis)** *Let $P_n$ be the probability distribution of $Z_n$. Under the assumptions above, $P_n$ satisfies an LDP with good rate function $I(x) = \sup_\theta[\theta^T x - \phi(\theta)]$.*

**Remarks** $I(\cdot)$ is in fact the Legendre transform of $\phi(\cdot)$ and to remember its formula one can recall the derivation of the Chernoff bound for $Z_n \in \mathscr{R}^1$: namely, $P(Z_n > x) \leqslant \inf_{\theta > 0} E \exp(n\theta[Z_n - x]) = \exp(-n \sup_\theta[\theta x - n^{-1} \log E \exp(n\theta Z_n)])$. Notice that since $I(\cdot)$ is convex its Legendre transform is $\phi(\cdot)$ and thus $I(\cdot)$ and $\phi(\cdot)$ are convex duals. From the Gärtner–Ellis theorem we can recall some standard results.

**Example 1 Sums of i.i.d. random variables** When $Z_n$ is the average of i.i.d. random variables, $Z_n = (X_1 + \cdots + X_n)/n$, then $\phi(\theta) = \log E[\exp(\theta X_1)]$ and the large deviation result is known as Cramér's theorem. In the case of univariate Gaussian $N(\mu, \sigma^2)$ random variables, $I(x) = (x - \mu)^2/2\sigma^2$. The observation $I(x) \geqslant 0$, with equality only for $x = \mu$, exemplifies a fact that holds more generally. $P(Z_n > \mu + \varepsilon) \sim_L \exp(-n\varepsilon^2/2\sigma^2)$ agrees with a direct calculation in the Gaussian case that $\log P(Z_n > \mu + \varepsilon) = -n\varepsilon^2/2\sigma^2 - (1/2)\log n + O(1)$. Notice the magnitudes of the terms that are neglected in the large deviation approximation.

**Example 2 Empirical distributions** Suppose $\{X_1, X_2, \ldots\}$ are i.i.d. discrete random variables, taking values $a_1, \ldots, a_d$, with probabilities $\pi_1, \ldots, \pi_d$ respectively. Let $Z_n$ be the vector whose $i$th component is the fraction of the observations $X_1, \ldots, X_n$ that are equal to $a_i$. Then $\hat{\pi}_{(n)} = Z_n$ is the empirical distribution. It follows that $\phi(\theta) = \log \sum_i \pi_i \exp(\theta_i)$, and so after solving the optimization problem posed by the Legendre transform we obtain the following version of Sanov's theorem.

$$P(\hat{\pi}_{(n)} \in A) \sim_L \exp\left(-n \inf_{p \in A} \sum_{i=1}^d p_i \log(p_i/\pi_i)\right). \tag{2}$$

**Example 3 Sample path averages** Suppose $x_1(t), x_2(t), \ldots, x_n(t)$ are realisations of $n$ identically distributed, but possibly dependent, continuous-time Markov processes in $\mathscr{R}^d$. Let $z_n(t) = n^{-1} \sum_{i=1}^n x_i(t)$ be their average. We seek a large deviation result for $P(z_n(t) \in A)$, where $A$ is a set of paths over $[0, T]$.

By seeking to estimate $P(z_n(t + \delta) - z_n(t) \approx \dot{z}(t)\delta \,|\, z_n(t) = z)$ and letting $\delta \to 0$ we are led to consider the derivative characteristic function (d.c.f.)

$$h(z, \theta) = \lim_{\delta \to 0} \lim_{n \to \infty} \frac{1}{n} E \frac{1}{\delta} [\exp(n\theta^T [z_n(t + \delta) - z_n(t)]) - 1 \,|\, z_n(t) = z].$$

The work of Freidlin–Wentzell and Weiss, confirms the guess that the rate function takes the form

$$I(z(\cdot)) = \sup_{\theta(\cdot)} \int_0^T \theta^T \dot{z} - h(z, \theta) dt$$

and

$$P(z_n \in A) \sim_L \exp\left( -n \inf_{z(\cdot) \in A} I(z(\cdot)) \right).$$

The function $h$ exists in the important case that the $n$ Markov processes evolve independently. For a Markov jump model defined on a finite state space $\{1, \ldots, d\}$, we can let the $i$th component of $x_n(t)$ be simply an indicator variable for the event that the process is in state $i$ at time $t$. The $i$th component of $z_n(t)$ is the proportion of the processes $x_1(\cdot), \ldots, x_n(\cdot)$ that are in state $i$, i.e., the empirical distribution of the state. If the processes evolve in a dependent manner the transition rate from state $i$ to $j$ might be allowed to depend on $z$ and we would denote this $q_{ij}(z)$. The d.c.f. is $h(z, \theta) = \sum_{i,j} z_i q_{ij}(z)[\exp(\theta_j - \theta_i) - 1]$. If the processes are independent, $q_{ij}(z) = q_{ij}$.

**Example 4 Fluid approximations**  The fluid approximation for $z(\cdot)$ is the deterministic path $z_\infty(\cdot)$ that starts at $z(0)$ and obeys the differential equation $\dot{z}_\infty = a(z_\infty)$, where $a(z) = \lim_{\delta \to 0} E[z(t + \delta) - z(t) \,|\, z(t) = z]$. From it's definition we have $I(z(\cdot)) \geqslant 0$ and Jensen's inequality applied to the definition of $h$, shows that $\theta a(z_\infty) - h(z_\infty, \theta) \leqslant 0$, and so $I(z(\cdot))$ is minimised to 0 by $z = z_\infty$.

For the Markov jump model discussed in the previous example, Weiss (1983) considers the set of paths $A = \{z(\cdot) : \|z - z_\infty\| \geqslant \varepsilon$ for some $t \in [0, T]\}$ and shows that $\inf_{z \in A} I(x) \geqslant k$ for some $k$, and thus $P(z_n \in A) \sim_L \exp(-nk)$. This characterises the fluid path as the most likely path.

**Varadhan's theorem and tail estimates**  In Sections 11.3 and 11.4 we shall use the following results. The first is known as Varadhan's theorem and is related to the ideal of Laplace that $\sum_i a_i \exp(-n\beta_i)/\sum_i \exp(-n\beta_i) \sim_L a_i$, where $i$ is the index for which $\beta_i$ is smallest. Theorem 3 estimates tail probabilities.

**Theorem 2 (Varadhan's theorem)**  *Let $\{Z_n\}$ be a sequence of random variables defined on a metric space $X$ whose probability distributions satisfy a LDP with good rate function $I(\cdot)$. Suppose $g(\cdot)$ is a real-valued continuous function on $X$, satisfying a moment condition $\lim_{n \to \infty} n^{-1} \log E[\exp(ng(Z_n))] < \infty$. Then $\lim_{n \to \infty} n^{-1} \log E[\exp(ng(Z_n))] = \sup_x [g(x) - I(x)]$.*

**Theorem 3**  *Suppose the same conditions as Theorem 2 hold, and* $\sup_x g(x) > 0$. *Define*

$$\alpha(g, I) = \inf\left\{\alpha > 0 : \sup_x [g(x) - \alpha^{-1}I(x)] \geqslant 0\right\}.$$

*Then for any real $\xi$,*

$$\lim_{\tau \to \infty} \frac{1}{\tau} \log \sup_{n \geqslant 1} n^\xi P(ng(Z_n) > \tau) = \lim_{\tau \to \infty} \frac{1}{\tau} \log \sum_{n=1}^{\infty} n^\xi P(ng(Z_n) > \tau) = -\alpha(g, I).$$

**Proof**  Suppose $\alpha < \alpha_0 = \alpha(g, I)$. This implies $\sup_x [g(x) - \alpha^{-1}I(x)] < -\delta$ for some $\delta > 0$. Then

$$\sum_{n=1}^{\infty} n^\xi P(ng(Z_n) > \tau) \leqslant \sum_{n=1}^{\infty} n^\xi E[\exp\{\alpha(ng(Z_n) - \tau)\}]$$

$$= \exp(-\alpha\tau) \sum_{n=1}^{\infty} n^\xi \exp\left\{n\alpha \sup_x [g(x) - \alpha^{-1}I(x)] + o(n)\right\}.$$

The right-hand equality follows from an application of Theorem 2. The final sum is finite, since $\sum_{n=1}^{\infty} n^\xi \exp(-\delta n) < \infty$, and thus

$$\varlimsup_{\tau \to \infty} \tau^{-1} \log \sum_n n^\xi P(ng(Z_n) > \tau) \leqslant -\alpha.$$

Now we use the large deviation lower bound to write as an asymptotic in $a$,

$$\varliminf_{\tau \to \infty} \frac{1}{\tau} \log \sup_{n \geqslant 1} n^\xi P(ng(Z_n) > \tau) \geqslant \frac{1}{\gamma} \varliminf_{\tau \to \infty} \frac{\gamma}{\tau} \log \left[\frac{\tau}{\gamma}\right]^\xi P([\tau/\gamma]g(Z_{[\tau/\gamma]}) > \tau)$$

$$\geqslant -\frac{1}{\gamma} \inf_{x : g(x) \geqslant \gamma} I(x).$$

The right-hand side is at least $-\alpha$ if there exists a $\gamma$ such that $\alpha > (1/\gamma)\inf_{\{x : g(x) = \gamma\}} I(x)$. This occurs if $\alpha > \alpha(g, I)$.

## 11.2  MOST LIKELY PATHS

**Example 5 An analysis of the FFD bin packing algorithm**  Suppose $n$ items of sizes distributed independently and uniformly over the integers $\{1, 2, \ldots, j\}$ are to be packed into bins of size $k$. A method of packing, called first-fit decreasing (FFD), is to imagine an infinite line of empty bins extending to the right. One takes the items in non-decreasing order of size, scans bins from the left and places each item into the first partially full or empty bin into which it will fit. Suppose $j = 6$ and $k = 13$. If $n$ is a multiple of 144 and there are exactly $n/6$ items of each size, we will have $n/12$ bins packed $[6, 6, 1]$, $n/12$ bins packed $[5, 5, 3]$, $n/18$ bins packed $[4, 4, 4, 1]$, $n/48$ bins packed $[3, 3, 3, 3, 1]$, $n/144$ bins

163

packed $[2, 2, 2, 2, 2, 2, 1]$ and $n/48$ bins packed $[2, 2, 2, 2, 2, 2]$. Note that in the last type of packing there is unused space of 1. We say that the wasted-space in the partially full bins is $n/48$. By assuming $n$ to be divisible by 144 the analysis of FFD is particularly simple. If $n$ had not been divisible by 144 there would have been a few 'transition bins', not conforming to one of the 'repeating' types given above, but the wasted space would have differed from $n/48$ by no more than a constant. Define $w_n(\pi_{(n)})$ as the wasted space within partially full bins that is left after applying FFD when $n$ items' sizes have empirical distribution $\pi_{(n)}$. By an analysis similar to the one illustrated above, we can show $|w_n(\pi_{(n)}) - n\zeta(\pi_{(n)})| < C$. Here $\zeta$ is independent of $n$ and is computed by pretending that $n$ has divisors such that no transition bins are created (even if in fact they would be needed); $C$ is independent of $\pi_{(n)}$ and $n$. By Example 2, the empirical distribution is asymptotically close to uniform. It follows that for $\varepsilon < 1/48$, $P(w(\pi_{(n)}) < \varepsilon n) \sim_L \exp(-nI(\varepsilon))$, where $I(\varepsilon)$ can be found by solving the optimisation problem posed in (2) for $A = \{p : \zeta(p) \leqslant \varepsilon\}$. For other values of $j, k$, cases of linear or sublinear growth in wasted space can be distinguished similarly. However, a finer analysis is needed to distinguish the case in which the expected wasted space is $\Omega(n^{1/2})$. Further discussion and more general cases are considered by Coffman *et al* (1993).

**Example 6 Copies of a MDP operating under a linked control**  Suppose as in Example 3, that the $i$th component of $x$ is 1 or 0 as an underlying Markov decision process, defined on a state space $\{1, \ldots, d\}$, is or is not in state $i$. Suppose that when at time $s$ the underlying process is in state $i$, then action $a$ is taken with probability $u_i^a(s)$, at a cost of $c_i^a$, with transition to state $j$ at rate $q_{ij}^a$. Then the cost, written $c(x, u) = \sum_{i,a} u_i^a c_i^a x_i$, is linear in $x$. Define

$$C(x, u, t) = \int_t^T c(x, u) ds, \qquad (3)$$

where $x, u$ are understood to be evaluated at $t$ on the left-hand side and at $s$ within the integral. If there are $n$ copies of the process evolving independently, but under the same policy and from the same initial state, then clearly $E[C(x_1, u, t)] = E[C(z_n, u, t)]$, where again $z_n(t)$ is the average of $x_1(t), \ldots, x_n(t)$. Since $z_n$ remains close to $z_\infty$, we have $E[C(z_n, u, t)] \sim C(z_\infty, u, t)$. However, this result is trivial since in fact, $z_\infty(t) = Ex_1(t)$ and so the approximation is exact.

Things are more interesting if $n$ copies of the process must be controlled in a manner that introduces dependence. Whittle (1988), has studied such a model, in which in each of $d$ states of a MDP there are two possible actions, called 'active' and 'passive.' As usual a cost is incurred and the state of the process changes in a random fashion depending on the action taken. Consider for the moment a single process, subject to a constraint that the active action must be taken a proportion $\mu$ of the time; the solution to the MDP will be to take this action in a set of states $J_\mu$, (possibly randomising in one state in order to match

the constraint exactly). If $J_\mu$ is monotone non-decreasing in $\mu$, this places a natural priority ordering on the states.

Now suppose that $n$ copies of the process are to be controlled simultaneously, with the constraint that the active action must be taken in exactly $m = n\mu$ processes at each instant. Whittle suggested the reasonable heuristic of taking the active action in precisely those $m$ processes whose states are of highest priority in the order determined above. Observe that under this policy the evolutions of the processes are dependent. However, one hopes that as $n$ increases their evolutions will be nearly independent and that the time-average cost will be the same as that obtained for a single process operating under the constraint that the active action is to be taken a proportion $\mu$ of the time.

A fluid model resolves the question and demonstrates that for some examples the hope may not be realised. Suppose the fluid path, $\dot{z}(t) = a(z, u)$, converges to an asymptotically stable equilibrium point $\bar{z}$. Mitra and Weiss (1988) have shown that for constants $c_1, c_2$, the time-average value of $\| z_n(t) - \bar{z} \|_2$ is less than $c_1 \exp(-nc_2)$. Essentially, this follows from the result quoted at the end of Section 11.4. Hence the time-average cost differs from the cost at $\bar{z}$ by an amount that tends to 0 as $n \to \infty$, and in this case the heuristic policy is asymptotically optimal. However, surprisingly, it can happen that the simultaneously controlled MDPs have a fluid approximation that tends to a stable limit cycle. In this case the time-average cost, per MDP, is asymptotically that of the fluid path integrated around the limit cycle, and this can be more than the cost obtained by the solution of the constrained problem for a single MDP. See Weber and Weiss (1990) for further details.

## 11.3 NEAR-LIKELY PATHS

Whittle (1990a) has given an intriguing elucidation of how a large deviation estimate combined with a risk-sensitive cost function can pose an optimal control problem whose solution is found along what one might call a 'near-likely' path of a stochastic process. The control problem satisfies a stochastic maximum principle and its solution is asymptotically optimal in a deterministic limit. It is also exact in the setting of linear dynamics, exponential quadratic cost and Gaussian noise (LEQG). The following exemplifies these ideas for the Markov jump process, where things are particularly simple.

**Example 7 Risk-sensitive control** Let us take the objective function

$$G(x, u, t) = -\frac{1}{\alpha} \log E \exp\left( -\alpha \int_t^T c(x, u) ds \right). \tag{4}$$

Minimisation of $G(x, u, t)$ for $\alpha > 0$ suggests a degree of risk seeking, in that the coefficient of $-\alpha^2/2$ in a Taylor's series expansion of $G(x, u, t)$ is the variance of the cost in (3). Similarly, minimising $G(x, u, t)$ for $\alpha < 0$ corresponds to risk

aversion. This is the risk-sensitivity criterion that forms the basis of Whittle's (1990b) development of optimal control in the setting LEQG. We shall assume a completely observable state. One of the particularly nice features of (4) is that costs that have already been incurred over an interval $[t, t_1)$ are irrelevant to determining the optimal control over $[t_1, T]$.

As in example 6 we shall consider a Markov jump process, recall that costs are linear and replicate $n$ independent copies, all starting in state $x$, to write

$$G(z_n, u, t) = \frac{1}{n} \sum_{i=1}^{n} G(x_i, u, t) = -\frac{1}{\alpha n} \log E\left[ \exp\left( -\alpha n \int_t^T c(z_n, u) ds \right) \right]. \quad (5)$$

Note that we are assuming that the action taken in each process does not depend on the states of the other processes. At first, one might think that controls, of 'global' character could do better by introducing covariance between the costs incurred by the $n$ processes. However, a dynamic programming argument can be used to show that this is not the case.

An estimate of (5) can be made using Varadhan's theorem; we multiply the cost along each possible path by the large deviation estimate of the probability of that path, and select the product which is greatest. This gives

$$G(z, u, t) = \inf_{z(\cdot)} \sup_{\theta(\cdot)} \left\{ \int_t^T c(z, u) + (1/\alpha)[\theta^T \dot{z} - h(z, u, \theta)] ds \right\} + o(n).$$

In fact, $n$ was arbitrary, so the $o(n)$ term may be ignored. Setting $\lambda = \theta/\alpha$, the problem becomes one of finding a control $u$ to achieve

$$F(z, t) = \inf_{u(\cdot)} \inf_{z(\cdot)} \sup_{\lambda(\cdot)} \left[ \int_t^T c(z, u) + \lambda^T \dot{z} - \alpha^{-1} h(z, u, \lambda \alpha) ds \right]. \quad (6)$$

The above requires stationarity conditions to hold that can be summarised in a maximum principle with Hamiltonian $H(z, u, \lambda) = -c(z, u) + \alpha^{-1} h(z, u, \lambda \alpha)$. The virtue of the maximum principle is that the determination of an optimal policy is assisted by knowing that the optimal control must maximise $H(z, u, \lambda)$, where the adjoint variables, $\lambda = -\partial F/\partial z_i$, satisfy $\dot{\lambda} = -\partial H/\partial z_i$, and that the dynamics for the path are given by $\dot{z} = \partial H/\partial \lambda_i$.

Notice that $\alpha^{-1} h(z, u, \lambda \alpha)$ is the derivative characteristic function for the average of $\alpha^{-1}$ independent processes. What we have done is to rewrite the problem in a new way, in which the risk parameter $\alpha$ enters by way of a change to the dynamics of the path. As $\alpha \to 0$ $\alpha^{-1} h(z, u, \lambda \alpha) \to \lambda^T a(z, u)$ and we recover a maximum principle for an optimal control formulation of the deterministic problem posed in the risk-neutral case. For the Markov jump process,

$$H(z, u, \lambda) = \sum_{i,a} z_i u_i^a \left[ -c_i^a + \alpha^{-1} \sum_j q_{ij}^a (e^{(\lambda_j - \lambda_i)\alpha} - 1) \right],$$

and it turns out that the equation for $\dot{\lambda}_i$ is just the dynamic programming

equation

$$\frac{d}{dt}e^{\lambda_i\alpha} = \sup_u \sum_{i,a} z_i u_i^a \left[ \alpha c_i^a e^{\lambda_i\alpha} - \sum_j q_{ij}^a (e^{\lambda_j\alpha} - e^{\lambda_i\alpha}) \right].$$

Whittle's approach gives the same result, but makes no assumption about linear costs or the form of the Markov process, only that the single process is already close to being an average of $n$ independent copies of some process, and therefore that it deviates only slightly from its fluid limit. Effectively, this is the same as supposing that the process of interest is $z_n$, with d.c.f. $nh(z, \theta/n)$. Taking the initial cost criterion as (5) and applying Varadhan's theorem leads to the risk-sensitive, stochastic maximum principle that is exact to within $o(n)$. What we have seen in this section is that the principle takes an exact form for Markov jump processes, and in fact, whenever $c(x, u)$ is linear in $x$.

Theorem 3 provides a further interpretation. If we take $g(z) = \mu - C(z, u, t)$ and $I(z, u) = \sup_\theta \int \theta^T \dot{z} - h(z, u, \theta)$, then there exists $z(\cdot), u(\cdot)$ such that $-\alpha(g, I) \leqslant -\alpha$, only if $\mu \leqslant F(z, t)$. Thus risk-seeking optimal control aligns with the problem of satisfying $\lim_{\tau \to \infty} \tau^{-1} \log \sum_{n=1}^{\infty} P(n\mu - nC(z_n, u, t) > \tau) \leqslant -\alpha$ for the largest value of $\mu$. Similarly, taking $g(z) = C(z, u, t) - \mu$, risk adverse optimal control aligns with the problem of satisfying $\lim_{\tau \to \infty} \tau^{-1} \log \sum_{n=1}^{\infty} P(nC(z_n, u, t) - n\mu > \tau) \leqslant -\alpha$ for the smallest possible $\mu$. In this case, we are saying that for all $n$, the probability that $C(z_n, u, t)$ exceeds $\mu + \tau/n$, is bounded as $\tau \to \infty$, by $\exp(-\alpha\tau)$. Again, greater $\alpha$ corresponds to greater risk aversion.

## 11.4  UNLIKELY PATHS

In the final section we shall consider the use of large deviations in controlling the probability of a rare event.

**Example 8  Effective bandwidths in ATM**  In a communication network using asynchronous transfer mode (ATM), data is packaged in cells of fixed size and transmitted between switches of the network over high bandwidth links. Because traffic sources are bursty there may be intervals of time over which cells arrive at a switch faster than they can be switched to output links. However, each switch has a buffer and so provided a switch does not attempt to carry too many calls the probability of buffer overflow and resulting cell loss can be kept smaller than some prespecified limit. The question arises as to the numbers and types of calls that can be so carried.

Suppose that a switch handles $M$ classes of traffic, consisting of $N_i$ traffic sources of class $i, i = 1, \ldots, M$. The bandwidth of the switch is the number of cells that it can switch per second and is denoted by $C$. A number of authors have described models in which a quality of service criterion is satisfied if and only if

$$\sum_{i=1}^{M} N_i e_i \leqslant C. \tag{7}$$

Here, $e_i$, is called the *effective bandwidth* of a source in class $i$. Intuitively, the bursty character of a source means that its effective bandwidth should be greater than its average rate. However, because at any moment some sources will deliver cells to the buffer at above their average rate and other sources will do so at below their average rate, there is potential for statistical multiplexing. Thus $e_i$ need not be as great as the peak rate of source $i$.

Of course if effective bandwidths can be associated with bursty ATM sources then problems of admission control and routeing in ATM networks resemble those in circuit-switched networks. Subsequent research can focus on the application of ideas from circuit-switched networks, (such as trunk reservation and dynamic routeing).

Let $W$ have the steady-state distribution of the workload in a queue with an infinite buffer. De Veciana and Walrand (1993) have given an excellent review of effective bandwidth results for models concerned with satisfying a constraint

$$\lim_{B \to /\text{infty}} \frac{1}{B} P(W > B) < -\delta. \tag{8}$$

In a related constraint $P(W > B)$ is replaced in (8) by the steady-state probability that the buffer exceeds level $B$ during a busy period. These constraints have led to similar effective bandwidth formulae in the work of Kelly (1991) (for the $M/GI/1$ and $D/GI/1$ queue), Courcoubetis and Walrand (1991) (for discrete-time Gaussian stationary sources), De Veciana *et al* (1993) and Gibbens and Hunt (1991) (for sources whose rate is modulated by a continuous time Markov process), Kesidis, *et al* (1994) (for general stationary sources) and Courcoubetis and Weber (1995) (in an approximation to effective bandwidths for general stationary sources).

Since the buffer should fill infrequently, the theory of large deviations is useful. We shall adopt the quality of service constraint that there should be only a small proportion of the time that the buffer is full. This form of constraint is more obviously related to a practical quality of service requirement than is (8), though indeed the same effective bandwidths are obtained. The following theorem is based on the derivation of De Veciana and Walrand (1993); it's proof is similar, though we have packaged the key ideas into Theorem 3.

**Theorem 4** *Consider a queue with a finite buffer of size $B$. Assume the number of cells arriving during epoch $n$ is $X_n$, where $\{X_n\}$ is a stationary ergodic process. The service rate is $C$ cells per epoch, with service taking place in such a way that $W_n$, the buffer content at the start of epoch $n$, follows the law $W_n = \min\{(W_{n-1} + X_{n-1} - C)^+, B\}$. Suppose $EX_n < C$ and*

$$\phi(\theta) = \lim_{n \to \infty} \frac{1}{n} \log E \exp\left[ \theta \sum_{i=1}^{n} X_i \right]$$

*exists and satisfies the conditions of the Gärtner–Ellis theorem. Let $L(B)$ be the*

*proportion of epochs during which buffer overflow occurs, defined as*

$$\lim_{n \to \infty} P(W_n + X_n - C > B).$$

*Then*

$$\lim_{B \to \infty} \frac{1}{B} \log L(B) \leqslant -\theta_0 \Leftrightarrow C \geqslant \phi(\theta_0)/\theta_0.$$

**Remark** This implies effective bandwidths. For if the single source imagined in the statement of the theorem actually comprises $N_i$ independent sources of type $i$, each having asymptotic logarithmic generating function $\phi_i(\cdot)$, $i = 1, \ldots, M$, then $\phi(\theta_0)/\theta_0 = \sum_i N_i \phi_i(\theta_0)/\theta_0$, and the effective bandwidth for a type $i$ source is identified as $\phi_i(\theta_0)/\theta_0$.

**Proof** We apply theorem 3, with $Z_n$ taking the stationary distribution of $(X_1 + \cdots + X_n)/n$ and letting $g(x) = x - C$. In this case $I(x) = \sup_\theta [\theta x - \phi(\theta)]$ and $C \geqslant \phi(\theta_0)/\theta_0 \Leftrightarrow \theta_0 C \geqslant \sup_x [\theta_0 x - I(x)] \Leftrightarrow \theta_0 \leqslant \inf_x [I(x)/g(x)]$.

Suppose the buffer is in its stationary distribution at epoch $-1$. If the buffer overflows during epoch $-1$ and was last empty at the start of epoch $-n$, we must have $X_{-n} + \cdots + X_{-1} - nC \geqslant B$. This must hold for some $n$. Hence theorem 3 now applies to bound the right-hand side of the inequality $L(B) \leqslant \sum_{n=1}^{\infty} P(X_{-n} + \cdots + X_{-1} - nC \geqslant B)$.

On the other hand, consider $n$ consecutive periods. By the stationarity of the process, we have $L(B) = E$ (number of epochs of buffer overflow in $n$ consecutive epochs)$/n$. But the numerator is bounded below by the probability that there is buffer overflow in at least one of $n$ epochs and this is at least $P(X_1 + \cdots + X_n - nC \geqslant B)$. So we can apply theorem 3 to $n^{-1} P(X_1 + \cdots + X_n - nC \geqslant B)$, i.e., with $\xi = -1$.

For the performance constraint of $L(B) \leqslant \exp(-\delta)$ we take $\theta_0 = \delta/B$ to be small when $B$ is large. To find an approximation for $\phi(\theta)/\theta$ let us make the following assumption.

**Assumption** *Suppose that* $\phi(\theta) = \mu\theta + \gamma\theta^2/2 + o(\theta^2)$ *where*

$$\gamma = \lim_{n \to \infty} \frac{1}{n} \mathrm{Var}\left( \sum_{i=1}^{N} X_i \right),$$

*and that the stationary process* $\{X_n\}$ *has spectral density*

$$f(\omega) = (1/\pi) \sum_{k=-\infty}^{\infty} \gamma(k) \exp(i\omega k).$$

*Suppose the infinite sum of the autocovariances is absolutely summable and* $f(\cdot)$ *is continuous at $0$. Then* $\gamma = \pi f(0) = \sum_{k=-\infty}^{\infty} \gamma(k)$*, where* $\gamma$ *is called the index of dispersion.*

This condition is plausible if there are no long-term autocorrelations. It is satisfied by Gaussian stationary sources and other processes, such as the Markov modulated fluid.

Assuming the assumption holds for each source type, Courcoubetis and Weber (1993) derived an asymptotic for small $\theta_0$, equivalent to expanding $\phi(\theta_0)/\theta_0$ in powers of $\theta_0$. Doing this, the effective bandwidth of a type $i$ source can be written $\phi_i(\delta/B)/(\delta/B) = \mu_i + \delta\gamma_i/2B + o(\delta/B)$, which suggests use of the approximate effective bandwidth

$$e_i = \mu_i + \frac{\delta\gamma_i}{2B}. \tag{9}$$

Courcoubetis and Walrand (1991) gave this formula for Gaussian sources, in which case $\gamma_i$ is replaced by the variance, and the formula is exact, since the $o(\delta/B)$ term vanishes. The formula is appealing for a number of reasons.

*On-line estimation* The parameters $\mu_i$ and $\gamma_i$ can be estimated on-line. This is attractive since it is unlikely that any theoretical model is rich enough to adequately model all traffic classes.

*Scaling* The bandwidths scale correctly if the division of time into slots is altered. For example, if the definition of a slot doubles, then $c$, $\mu_i$ and $\gamma_i$ all double.

*Filtering* If a source is filtered before it enters the buffer it may become less bursty. Can this help? If we apply a filter with transfer function $T(\cdot)$ and no cells are lost then $T(0) = 1$. The mean of the output is still $\mu$, the index of dispersion is $\pi|T(0)|^2 f(0) = \pi f(0) = \gamma$, and so there is no change in the effective bandwidth due to smoothing. This is not surprising, since in large buffer asymptotics we are already seeing a smoothing in the buffer and the effects of filtering are masked. De Veciana and Walrand (1993) comment that if $T(0) = G < 1$, which happens if the source is thinned, then the bandwidth changes to $G\mu + G^2\delta\gamma/2B$. Thus bandwidths can be reduced by thinning, but not by smoothing.

*Performance* The bandwidths in (9) may be compared with other bandwidth formulae, both by models and by simulation. For example, consider a source whose rate is modulated by a two-state Markov process that alternates between states 1 and 2, with holding times in these states that are exponentially distributed with parameters $\lambda_i$ and $\mu_i$ respectively. The rate of the source is 0 or $a_i$ as the state is 1 or 2 respectively. Then (9) gives

$$e_i = \frac{\lambda_i a_i}{\lambda_i + \mu_i} + \frac{\delta\lambda_i\mu_i a^2}{B(\lambda_i + \mu_i)^3}. \tag{10}$$

Taking $\varepsilon = \delta/B$, $e_i^\dagger = \phi(\varepsilon)/(\varepsilon) = \{-[\lambda_i + \mu_i - a_i\varepsilon] + \sqrt{[\lambda_i + \mu_i - a_i\varepsilon]^2 + 4\lambda_i a_i\varepsilon}\}/2\varepsilon$, which is given by Gibbens and Hunt (1991), and also implicit in the work of De Veciana *et al* (1993). The difference between $e_i$ and $e_i^\dagger$ is small. Consider, for example, a model of a voice call source in which, counting time in seconds

and bits in 1000's, we take $a = 30$ Kbps, $\lambda = 2$, $\mu = 3$, $\delta = 10$. A buffer of 200 ATM cells, each of 54 bytes, of 8 bits, is about $B = 80$. So for $\delta = 1/8$, we find $e = 17.40$ and $e^{\dagger} = 17.47$. For $\lambda = 3$, $\mu = 2$, $e = 23.40$ and $e^{\dagger} = 22.29$. For $\lambda = 2.5$, $\mu = 2.5$, $e = 20.63$ and $e^{\dagger} = 20.00$. Note that $B = 80$ corresponds to buffering about 5 seconds of peak rate from a single source. On this basis a 100 Mbps switch might carry 5700 voice calls of this model class. Courcoubetis *et al* (1994) have made further study use of (9) and (10).

# REFERENCES

Bucklew, J. (1990) *Large Deviation Techniques in Decision, Simulation and Estimation.* John Wiley, New York.

Coffman Jr, E., Johnson, D., Shor, P. and Weber, R. (1993) Bin packing with discrete items sizes, Part II: average-case behavior of FFD and BFD. Technical report, AT&T Bell Laboratories.

Courcoubetis, C., Fouskas, G. and Weber, R. (1994) On the performance of an effective bandwidths formula. Proceedings of the 14th International Teletraffic Conference, 1994.

Courcoubetis, C. and Walrand, J. (1991) Note on the effective bandwidth of ATM traffic at a buffer. unpublished manuscript.

Courcoubetis, C. and Weber, R. (1995) Effective bandwidths for stationary sources. *Probability in the Engineering and Informational Sciences, to* appear.

De Veciana, G., Olivier, C. and Walrand, J. (1993) Large deviations for birth death Markov fluids. *Probability in the Engineering and Informational Sciences* 7, 235–237.

De Veciana, G. and Walrand, J. (1993) Effective bandwidths: call admission, traffic policing and filtering for ATM networks. Technical Report UNB/ERL M93/47, U.C. Berkeley.

Dembo, A. and Zeitouni, O. (1993) *Large Deviations Techniques and Applications.* Jones and Bartlett, Boston.

Freidlin, M. and Wentzell, A. (1984) *Random Perturbations of Dynamical Systems.* New York: Springer-Verlag.

Gibbens, R. and Hunt, P. (1991) Effective bandwidths for the multi-type UAS channel. Statistical Laboratory, University of Cambridge, preprint.

Kelly, F. (1991) Effective bandwidths at multi-class queues. *Queueing Systems* 9, 5–16.

Kesidis, G., Walrand, J. and Chang, C. (1994) Effective bandwidths for multiclass Markov fluids and other ATM sources. *IEEE/ACM Transactions on Networking* 1, 424–428.

Mitra, D. and Weiss, A. (1988) A fluid limit of a closed queueing network with applications to networks. Technical report, AT&T Bell Laboratories.

Weber, R. and Weiss, G. (1990) On an index policy for restless bandits. *J. Appl. Prob.* 27, 637–648.

Weiss, A. (1983) The large deviation of a Markov process which models traffic generation. Technical report, AT&T Bell Laboratories.

Whittle, P. (1988) Restless bandits: activity allocation in a changing world. In J. Gani (Ed.), *A Celebration of Applied Probability,* Applied Probability Special Volume 25A, pp. 287–298. University of Sheffield, UK: Applied Probability Trust.

Whittle, P. (1990a) A risk-sensitive maximum principle. *System and Control Letters* 15, 183–192.

Whittle, P. (1990b). *Risk-Sensitive Optimal Control.* Wiley.

*Queens' College, University of Cambridge, CAMBRIDGE CB3 9ET.*