

Combinatorics, Probability and Computing (1999) 8, 473–482. Printed in the United Kingdom
© 1999 Cambridge University Press

Studying Balanced Allocations with Differential Equations[†]

MICHAEL MITZENMACHER[§]

Digital Systems Research Center, 130 Lytton Avenue, Palo Alto, CA 94301, USA
(e-mail: michaelm@pa.dec.com)

Received 4 November 1997; revised 25 June 1998

Using differential equations, we examine the GREEDY algorithm studied by Azar, Broder, Karlin and Upfal for distributed load balancing [1]. This approach yields accurate estimates of the actual load distribution, provides insight into the exponential improvement GREEDY offers over simple random selection, and allows one to prove tight concentration theorems about the loads in a straightforward manner.

1. Introduction

Suppose that n balls are placed into n bins, with each ball being placed into a bin chosen independently and uniformly at random. Let the *load* of a bin be the number of balls in that bin after all balls have been thrown. It is well known that, with high probability, the maximum load upon completion will be approximately $\frac{\log n}{\log \log n}$ [8]. (We use *with high probability* to mean with probability at least $1 - O(1/n)$, where n is the number of balls. Also, \log will always mean the natural logarithm, unless otherwise noted.)

Azar, Broder, Karlin and Upfal considered how much more evenly distributed the load would be if each ball had two (or more) choices [1]. Suppose that the balls are placed sequentially, and each ball is placed into the less full of *two* bins chosen independently and uniformly at random with replacement (breaking ties arbitrarily). In this case, they showed that the maximum load drops to $\frac{\log \log n}{\log 2} + O(1)$ with high probability. If each ball instead

[†] A preliminary version of this work appeared in the *Proceedings of the 37th Annual IEEE Symposium on the Foundations of Computer Science*, October 1996.

[§] Much of this work was done at U.C. Berkeley, supported by a fellowship from the Office of Naval Research and by NSF grant CCR-9505448.

has d choices, then the maximum load will be $\frac{\log \log n}{\log d} + O(1)$ with high probability. Having two choices hence yields a qualitatively different type of behaviour from the single choice case, leading to an exponential improvement in the maximum load; having more than two choices further improves the maximum load by only a constant factor. This result has important implications for distributed load balancing, hashing, and PRAM simulation [1].

Following Azar, Broder, Karlin and Upfal, we refer to the algorithm in which each ball has d random choices as $\text{GREEDY}(d)$. In this paper, we develop an alternative method of studying the performance of $\text{GREEDY}(d)$ using differential equations. The differential equations describe the limiting performance of $\text{GREEDY}(d)$ as the number of balls and bins tends to infinity. As we will demonstrate, the description of the limiting performance proves highly accurate, even when n is relatively small. In particular, we note that from our results one can determine the fraction of bins of any fixed load at the end of the process for the limiting case as n goes to infinity. These limiting quantities provide accurate estimates for the fraction of bins of each fixed load for sufficiently large systems. This analysis therefore adds significantly more detail to the previous analysis of [1]. Moreover, besides giving more insight into the actual performance of GREEDY , our methods provide a great deal of intuition behind the behavioural difference between one and two choices.

Our motivation in studying this problem is twofold. First, we wish to demonstrate and highlight this methodology, and encourage its use for studying random processes. While this methodology is by no means new, its uses have been surprisingly limited. The technical results justifying the relationship between families of Markov processes and differential equations date back at least to Kurtz [13, 14]. Karp and Sipser provided an early use of this technology to analyse an algorithm for finding maximum matchings in sparse random graphs [11]. Other past applications in the analysis of algorithms include [9, 12], and more recently many more have been found (see, for example, [2, 3, 15, 17, 21], to name a few).

Our second motivation is to demonstrate the power of using two choices. This idea dates back at least as far as the work of Eager, Lazowska and Zahorjan [7], who examined a dynamic load balancing model based on viewing processors as single server queues. In the static setting, this idea was also studied by Hajek [9], who used the same approach we undertake to determine the fraction of empty bins. The aforementioned exponential improvement in behaviour was noted and proved first in a paper by Karp, Luby and Meyer auf der Heide [10]. The work by Azar, Broder, Karlin and Upfal examined a simpler model that clarified the argument and provided many new results. Related work by the author [16, 17], as well as by others [20], examines the power of two choices in dynamic settings. Continued work in the area includes recent work by Stemmann [19] and Czumaj and Stemmann [6].

In the rest of the paper, we explain the derivation of the differential equations that describe the GREEDY strategy of [1] and compare the results from the differential equations with simulations. We also demonstrate how the equations give more insight into the behaviour of GREEDY and how the equations relate to the work in [1].

2. The differential equations

In this section, we demonstrate how to establish a family of differential equations that can be used to model the behaviour of the GREEDY strategy of [1]. We begin the process with m balls and n bins. We shall require for our analysis that $m = cn$ for some constant c . Balls arrive sequentially, and, upon arrival, each ball chooses d bins independently and uniformly at random (with replacement); the ball is then placed in the least loaded of these bins (with ties broken arbitrarily).

We first ask how many bins remain empty after the protocol GREEDY(d) terminates. This question has a natural interpretation in the task-processor model: how many of our processors are not utilized? The question can also be seen as a matching problem on random bipartite graphs: given a bipartite graph with n vertices on each side such that each vertex on the left has d edges to vertices chosen independently and uniformly at random on the right, what is the expected size of the greedy matching obtained by sequentially matching vertices on the left to a random unmatched neighbour? Our attack, again, is to consider this system as $n \rightarrow \infty$. This question has been previously solved by Hajek using entirely similar techniques [9]. We shall briefly repeat his argument with some additional insights. Once we show how to answer the question of the number of empty bins, we shall extend it to the more general load balancing problem.

2.1. The empty bins problem

We set up the problem of the number of empty bins by developing a Markov chain with a simple state that describes the balls and bins process. We first establish a concept of time. Let $Y(T)$ be the number of non-empty bins after T balls have been thrown. Then $\{Y(i)\}, i = 0 \dots m$, is clearly a Markov chain. Moreover,

$$\mathbb{E}[Y(T+1) - Y(T)] = 1 - \left(\frac{Y(T)}{n}\right)^d, \quad (2.1)$$

since the probability that a ball finds all non-empty bins among its d choices is $(Y(T)/n)^d$.

The notation becomes somewhat more convenient if we scale by a factor of n . We let t be the time at which exactly nt balls have been thrown, and we let $y(t)$ be the fraction of non-empty bins. Then equation (2.1) becomes

$$\frac{\mathbb{E}[y(t + 1/n) - y(t)]}{1/n} = 1 - (y(t))^d. \quad (2.2)$$

We claim the random process described by equation (2.2) is well approximated by the trajectory of the differential equation

$$\frac{dy}{dt} = 1 - y^d, \quad (2.3)$$

where this equation has been obtained from equation (2.2) by replacing the right-hand side with the appropriate limiting value as n tends to infinity, dy/dt . That the random

process given by the Markov chain closely follows the trajectory given by the differential equation follows easily from known techniques, such as, for example, Kurtz's theorem, or the similar work on random graphs by Wormald [21]. (As previously mentioned, the balls and bins process has a natural interpretation in terms of random bipartite graphs.)

To clarify this connection, here we state a form of Kurtz's theorem, as given in [18, Theorem 5.3]. We provide the necessary notation, and then relate the notation back to the underlying balls and bins process. Suppose we are given a finite set of vectors $\{\vec{e}_1, \dots, \vec{e}_k\}$ in \mathbf{R}^d . We consider an initial process $\vec{x}(t)$ with generator

$$Lf(\vec{x}) = \sum_{i=1}^k \lambda_i(\vec{x})(f(\vec{x} + \vec{e}_i) - f(\vec{x}))$$

and a scaled process $\vec{z}_n(t)$ with generator

$$L_n f(\vec{x}) = \sum_{i=1}^k n \lambda_i(\vec{x}) \left(f\left(\vec{x} + \frac{\vec{e}_i}{n}\right) - f(\vec{x}) \right).$$

The limiting operator L_∞ satisfies

$$L_\infty f(\vec{x}) = \sum_{i=1}^k n \lambda_i(\vec{x}) \langle \nabla f(\vec{x}), \vec{e}_i \rangle,$$

and corresponds to the deterministic solution \vec{z}_∞ of the equation

$$\frac{d}{dt} \vec{z}_\infty(t) = \sum_{i=1}^k \lambda_i(\vec{z}_\infty(t)) \vec{e}_i. \quad (2.4)$$

We relate the above notation to the problem of keeping track of the fraction of non-empty bins. For convenience, think of the times balls enter the system as being determined by a Poisson arrival process at a rate of one per unit time. (This does not affect the result; it merely simplifies the discussion. See, for example, [13].) In the scaled process with n balls and n bins, n balls arrive per unit time. The process is one-dimensional, and hence $\vec{e}_1 = (1)$. The rate function λ_1 satisfies $\lambda_1(\vec{x}) = 1 - x_1^d$: this is just the probability that an incoming ball lands in an empty bin. The limiting process is then given by the differential equation (2.4), which in this case is equivalent to equation (2.3); note that for convenience we have dropped the vector notation and simply used the variable y .

Kurtz's theorem states that the scaled processes approach the limiting process, with error bounds similar to Chernoff-like bounds.

Theorem 2.1 (Kurtz's theorem [18], Theorem 5.3). *Let $\lambda_i(\vec{x}) : \mathbf{R}^d \rightarrow \mathbf{R}^+$ be uniformly bounded and Lipschitz continuous, and let \vec{z}_∞ be the unique solution of (2.4) with $\vec{z}_\infty(0) = \vec{x}(0)$. For each finite T there exist a positive constant C_1 and a function C_2 with*

$$\lim_{\epsilon \downarrow 0} \frac{C_2(\epsilon)}{\epsilon^2} \in (0, \infty) \quad \text{and} \quad \lim_{\epsilon \uparrow \infty} \frac{C_2(\epsilon)}{\epsilon} = \infty$$

such that, for all $n \geq 1$ and $\epsilon > 0$,

$$\Pr \left(\sup_{0 \leq t \leq T} |\tilde{z}_n(t) - \tilde{z}_\infty(t)| \geq \epsilon \right) \leq C_1 e^{-nC_2(\epsilon)}.$$

Moreover, C_1 and C_2 can be chosen independently of \tilde{x} .

The connection between the balls and bins process and the differential equation (2.3) yields the following theorem.

Theorem 2.2. *Suppose cn balls are thrown into n bins according to the GREEDY(d) protocol for some constant c . Let Y_{cn} be the number of non-empty bins when the process terminates. Then $\lim_{n \rightarrow \infty} \mathbb{E}[\frac{Y_{cn}}{n}] = y_c$, where $y_c < 1$ satisfies*

$$c = \sum_{i=0}^{\infty} \frac{y_c^{id+1}}{(id+1)}.$$

Proof. The preconditions for Kurtz's theorem (the above or [14, Chapter 8]) are easily checked for the one-dimensional system described by (2.3), so by Kurtz's theorem we have that this differential equation is the correct limiting process.[†] Instead of solving (2.3) for y in terms of t , we solve for t in terms of y : $\frac{dt}{dy} = \frac{1}{1-y^d} = \sum_{i=0}^{\infty} y^{id}$. We integrate up to some time t , yielding

$$t = \sum_{i=0}^{\infty} \frac{y(t)^{id+1}}{(id+1)}. \quad (2.5)$$

From equation (2.5), given d we can solve for $y(t)$ for any value of t using, for example, binary search. One can also attempt to find an equation for y in terms of d and t ; standard integral tables give such equations when $d = 2, 3$ and 4 , for example. When $t = c$, all of the balls have been thrown, and the process terminates. Plugging $t = c$ into equation (2.5) yields the theorem, with $y_c = y(c)$. \square

We may actually use Kurtz's theorem to obtain a concentration result.

Theorem 2.3. *In the notation of Theorem 2.2, $|\frac{Y_{cn}}{n} - y_c|$ is, with high probability,*

$$O \left(\sqrt{\frac{\log n}{n}} \right),$$

where the constant depends on c . \square

[†] Again, it appears that there might be a problem here since we consider events occurring at discrete time-steps, instead of according to random times from a Poisson process. One can always adopt the convention that each discrete time-step corresponds to an amount of time given by an exponentially distributed random variable. In the limiting case, this distinction disappears.

One can also obtain entirely similar bounds for Y_{cn} using more straightforward martingale arguments. In the following, we assume familiarity with basic martingale theory: see, for example, [4, Chapter 7] for more information. We use the following form of the martingale tail inequality due to Azuma [5]:

Lemma 2.1 (Azuma [5]). *Let X_0, X_1, \dots, X_m be a martingale sequence such that, for each k ,*

$$|X_k - X_{k-1}| \leq 1.$$

Then, for any $\alpha > 0$,

$$\Pr(|X_m - X_0| > \alpha\sqrt{m}) < 2e^{-\alpha^2/2}.$$

Theorem 2.4. *In the notation of Theorem 2.2, $\Pr(|Y_{cn} - \mathbb{E}[Y_{cn}]| > \alpha\sqrt{cn}) < 2e^{-\alpha^2/2}$ for any $\alpha > 0$.*

Proof. For $0 \leq j \leq cn$, let \mathcal{F}_j be the σ -field of events corresponding to the possible states after j balls have been placed, and $Z_j = \mathbb{E}[Y_{cn} | \mathcal{F}_j]$ be the associated conditional expectation of Y_{cn} . Then the random variables $\{Z_j\}_{j=0}^{cn}$ form a Doob martingale, and it is clear that $|Z_j - Z_{j-1}| \leq 1$. The theorem now follows from Lemma 2.1. \square

Theorem 2.4 implies that Y_{cn} is within $O(\sqrt{n \log n})$ of its expected value with high probability; however, the martingale approach does not immediately lead us to the value to which Y_{cn}/n converges. This is a standard limitation of the martingale approach: in contrast, the differential equations approach allows us to find the mean as well as prove concentration around the mean.

2.2. Bins with fixed load

We can extend the previous analysis to find the fraction of bins with load at least (or exactly) k for any constant k as $n \rightarrow \infty$. We first establish the appropriate Markov chain. Let $s_i(t)$ be the fraction of bins with load at least i at time t , where again at time t exactly nt balls have been thrown. Then the corresponding differential equations regarding the growth of the s_i (for $i \geq 1$) are easily determined:

$$\begin{cases} \frac{ds_i}{dt} &= (s_{i-1}^d - s_i^d) \text{ for } i \geq 1; \\ s_0 &= 1. \end{cases} \quad (2.6)$$

The differential equations have the following simple interpretation: for there to be an increase in the number of bins with at least i balls, the d choices of a ball about to be placed must all be bins with load at least $i-1$, but not all bins with load at least i .

In the context of the notation established for Kurtz's theorem, we let \vec{x} be a k -dimensional vector, with \vec{e}_i being a standard unit vector in the i th dimension. Then

Table 1 Predicted behaviour for GREEDY(d) and average results from 100 simulations with 1 million balls.

	$d = 2$ prediction	1 million simulation	$d = 3$ prediction	1 million simulation
s_1	0.7616	0.7616	0.8231	0.8230
s_2	0.2295	0.2295	0.1765	0.1765
s_3	0.0089	0.0089	0.00051	0.00051
s_4	0.000006	0.000007	$< 10^{-11}$	0
s_5	$< 10^{-11}$	0	$< 10^{-11}$	0

$\lambda_1(\tilde{x}) = 1 - x_1^d$ and $\lambda_i(\tilde{x}) = x_{i-1}^d - x_i^d$ for $i > 1$. The corresponding limiting solution given by (2.4) is then equivalent to the system (2.6), where again we have removed the vector notation for convenience.

In contrast to Section 2.1, where we could derive a formula for the fraction of empty bins, we are not aware of how to determine explicit formulae for $s_i(t)$ in general. These systems of differential equations can be solved numerically using standard methods, however; for up to any fixed k and t , we can accurately determine $s_k(t)$. By applying Kurtz's theorem to the k -dimensional process (as in Theorem 2.3) or martingale arguments (as in Theorem 2.4), one can show that these results will be accurate with high probability.

We also demonstrate that our technique accurately predicts the behaviour of the GREEDY(d) algorithm by comparing with simulation results. The first and third columns of Table 1 shows the values of $s_i(1)$ for $d = 2$ and $d = 3$ as calculated from the differential equations. We use these values as predictions for the process where we throw n balls into n bins. From the predictions with $d = 2$, one would not expect to see bins with load five until billions of balls have been thrown. Similarly, choosing $d = 3$ one expects a maximum load of three until billions of balls have been thrown. These results match our simulations of several hundred runs with up to thirty-two million balls, the largest simulation we have attempted. We also present the averages from one hundred simulations of one million balls for $d = 2$ and $d = 3$, which demonstrate the accuracy of the technique in predicting the behaviour of the system. Further simulations reveal that, in general, the solution given by the limiting system of differential equations becomes more accurate as n grows, and the deviation from this solution is small, as one would expect. This accuracy is a marked advantage of this approach: previous techniques have not provided ways of concretely predicting actual performance.

2.3. Relationship to $O(\log \log n)$ bounds

We can also use the differential equations to provide an alternative derivation of a key idea from the proof of the upper bounds on the maximum load of GREEDY(d). The approach of looking at the underlying differential equations provides insight into how the s_k decrease, which is essential to determining the $O(\log \log n)$ bounds.

We begin by focusing on the case where the number of balls n equals the number of bins n , and consider the limiting description given by the differential equations as $n \rightarrow \infty$.

Lemma 2.2. *For the family of differential equations given by (2.6), $s_i(1) \leq [s_{i-1}(1)]^d$.*

Proof. We wish to know the values of $s_i(1)$. Because the s_i are all non-decreasing over time and nonnegative, from (2.6),

$$\frac{ds_i}{dt} = s_{i-1}^d - s_i^d \leq [s_{i-1}(1)]^d$$

for all $t \leq 1$ and hence by integrating

$$s_i(1) \leq [s_{i-1}(1)]^d. \tag{2.7}$$

□

One can calculate $s_1(1)$ directly from Theorem 2.2, and it follows from a simple induction on (2.7) that

$$s_i(1) \leq [s_1(1)]^{d^{i-1}}.$$

In other words, the $s_i(1)$, which represent the limiting fraction of bins with load at least i after all balls have been thrown, decrease *doubly exponentially*, as the i is in the second level of the exponent. Using Kurtz's theorem one obtains a high probability result for the case of a finite number of balls n .

Theorem 2.5. *Let w_i^n be the fraction of bins with load at least i when n balls are thrown into n bins using GREEDY(d). Then, for any $\epsilon > 0$ and fixed i , $w_i^n \leq (1 + \epsilon)[s_1(1)]^{d^{i-1}}$ with probability at least $1 - e^{-c(\log i)nc^2}$ for some suitable constant c .*

Proof. The proof is a direct application of Kurtz's theorem, using the appropriate error bounds. □

This doubly exponential decrease in the $s_i(1)$ (or, equivalently, of the w_i^n) is a key step of the proof of Azar, Broder, Karlin and Upfal in [1], where it is proved via an inductive use of Chernoff bounds. Theorem 2.5 shows that this induction can be replaced by applying Kurtz's theorem to the differential equations, at least up to any fixed constant value of i . This approach has some advantages. Most importantly, Lemma 2.2 and the corresponding inductive bound for $s_i(1)$ seem quite natural and make transparent how the s_i decrease. Additionally, using this approach can improve the additive $O(1)$ term in Theorem 4 of [1].

Intuitively, this doubly exponential decrease suggests that, if we look at bins with load

at least $i^* = \lceil \frac{\log \log n}{\log d} + \gamma \rceil$, where $\gamma = (\log 2 - \log \log(1/s_1(1)))/\log d$, then

$$s_{i^*+1} \leq [s_1(1)]^{d^{\log \log n / \log d + \gamma}} \leq \frac{1}{n^2},$$

and hence with high probability there will be no bins with load at least $i^* + 1$.

Note, however, that as stated in Theorem 2.5 the high probability bounds only hold up to any fixed i , and not for values of i up to $\Omega(\log \log n)$. This weakness in the result stems from directly applying the differential equations approach and Kurtz's theorem, which requires the underlying state space to be finite-dimensional, and hence loads up to only some fixed constant can be considered. By reverting back to the explicit martingale argument that underlies Kurtz's theorem, we can circumvent this restriction somewhat (up until the point where s_i is $1/n^{1/2+\delta}$ for any $\delta > 0$), but at some point when s_i is sufficiently small it seems we have to explicitly handle this case directly, as is done in Theorem 4 of [1]. We omit details of these more extensive arguments, since their form would be almost entirely a restatement of Theorem 4 of [1], replacing their use of Chernoff bounds with an equivalent martingale argument. A general framework that would allow us to apply the differential equations in this instance in a more straightforward manner would clearly be appealing, since the differential equations do make clear the behaviour of the s_i .

In the case of m balls and n bins, a similar argument to Lemma 2.2 and Theorem 2.5 holds when $m = cn$ for some constant c . As in Section 2.1, when $m = cn$, the infinite process runs until time c ; if m is not a linear function of n , the time until the process terminates is dependent on n , and Kurtz's theorem cannot be applied. For Lemma 2.2, the appropriate result becomes

$$s_i(c) \leq c[s_{i-1}(c)]^d.$$

By noting that $s_{c^2+c}(c) \leq 1/(c+1)$ (the fraction of bins with load at least x cannot be more than c/x), one may again inductively show as before that the tails of the loads are doubly exponentially decreasing. Improvements can be made in the constants by using the differential equations to find better starting points than $s_{c^2+c}(c) \leq 1/(c+1)$ for the induction.

3. Conclusion

There are significant advantages to using differential equations to study randomized load balancing problems. The insight one gains about the problem and the numerical accuracy one obtains are quite convincing. Moreover, when the corresponding state spaces are finite-dimensional, applying Kurtz's theorem can yield simple proofs of the limiting behaviour. A general framework for dealing with spaces that are not necessarily finite-dimensional would greatly simplify using this approach in developing bounds such as the $O(\log \log n)$ bounds of [1]. We expect this approach will find a great deal of further use in the analysis of load balancing schemes, as well as other algorithmic areas.

References

- [1] Azar, Y., Broder, A., Karlin, A. and Upfal, E. (1994) Balanced allocations. In *Proceedings of the 26th ACM Symposium on the Theory of Computing*, pp. 593–602.
- [2] Achlioptas, D. and Molloy, M. (1997) The analysis of a list-coloring algorithm on a random graph. In *Proceedings of the 38th IEEE Symposium on Foundations of Computer Science*, pp. 204–212.
- [3] Alanyali, M. and Hajek, B. (1997) Analysis of simple algorithms for dynamic load balancing. *Mathematics of Operations Research* **22** 840–871.
- [4] Alon, N., Spencer, J. and Erdős, P. (1992) *The Probabilistic Method*, Wiley.
- [5] Azuma, K. (1967) Weighted sums of certain dependent random variables. *Tôhoku Mathematical Journal* **19** 357–367.
- [6] Czumaj, A. and Stemmann, V. (1997) Randomized allocation processes. In *Proceedings of the 38th IEEE Symposium on Foundations of Computer Science*, pp. 194–203.
- [7] Eager, D. L., Lazowska, E. D. and Zahorjan, J. (1986) Adaptive load sharing in homogeneous distributed systems. *IEEE Trans. Software Engineering* **12** 662–675.
- [8] Gonnet, G. (1981) Expected length of the longest probe sequence in hash code searching. *J. Assoc. Comput. Mach.* **28** 289–304.
- [9] Hajek, B. (1988) Asymptotic analysis of an assignment problem arising in a distributed communications protocol. In *Proceedings of the 27th IEEE Conference on Decision and Control*, pp. 1455–1459.
- [10] Karp, R. M., Luby, M. and Meyer auf der Heide, F. (1992) Efficient PRAM simulation on a distributed memory machine. In *Proceedings of the 24th ACM Symposium on the Theory of Computing*, pp. 318–326.
- [11] Karp, R. M. and Sipser, M. (1981) Maximum matchings in sparse random graphs. In *Proceedings of the 22nd IEEE Symposium on Foundations of Computer Science*, pp. 364–375.
- [12] Karp, R. M., Vazirani, U. V. and Vazirani, V. V. (1990) An optimal algorithm for on-line bipartite matching. In *Proceedings of the 22nd ACM Symposium on the Theory of Computing*, pp. 352–358.
- [13] Kurtz, T. G. (1970) Solutions of ordinary differential equations as limits of pure jump Markov processes. *J. Appl. Probab.* **7** 49–58.
- [14] Kurtz, T. G. (1981) *Approximation of Population Processes*, SIAM.
- [15] Luby, M., Mitzenmacher, M., Shokrollahi, M. A., Spielman, D. and Stemmann, V. (1997) Practical loss-resilient codes. In *Proceedings of the 29th ACM Symposium on the Theory of Computing*, pp. 150–159.
- [16] Mitzenmacher, M. (1996) Load balancing and density dependent jump Markov processes. In *Proc. of the 37th IEEE Symp. on Foundations of Computer Science*, pp. 213–222.
- [17] Mitzenmacher, M. (1996) The power of two choices in randomized load balancing. PhD thesis, University of California, Berkeley.
- [18] Shwartz, A., Weiss, A. (1995) *Large Deviations for Performance Analysis*, Chapman and Hall.
- [19] Stemmann, V. (1996) Parallel balanced allocations. In *Proceedings of the Eighth Annual ACM Symposium on Parallel Algorithms and Architectures*, pp. 261–269.
- [20] Vvedenskaya, N. D., Dobrushin, R. L. and Karpelevich, F. I. (1996) Queueing system with selection of the shortest of two queues: an asymptotic approach. *Problems of Information Transmission* **32** 15–27.
- [21] Wormald, N. C. (1995) Differential equations for random processes and random graphs. *Ann. Appl. Probab.* **5** 1217–1235.