

Load Sharing in Distributed Real-Time Systems with Broadcast of State Changes

Kang G. Shin and Yi-Chieh Chang¹

TR-88-006

October 31, 1988

¹Kang G. Shin is with the International Computer Science Institute, 1947 Center St., Suite 600, Berkeley, CA 94704-1105, on leave from Real-Time Computing Laboratory, Dept. of Electrical Engineering and Computer Science, The University of Michigan, Ann Arbor, MI 48109-2122, which Yi-Chieh Chang is currently with.

The work reported in this paper was supported in part by the Office of Naval Research under contract N00014-85-K-0122 and the NASA under grant NAG-1-887. Any opinions, findings, and recommendations expressed in this publication are those of the authors and do not necessarily reflect the view of the funding agencies.

LOAD SHARING IN DISTRIBUTED REAL-TIME SYSTEMS WITH BROADCAST OF STATE CHANGES

Kang G. Shin and Yi-Chieh Chang ¹

October 31, 1988

ABSTRACT

If task arrivals are not uniformly distributed over the nodes in a distributed real-time system, some nodes may become overloaded while others are lightly-loaded or even idle. Consequently, some tasks cannot be completed before their deadlines, even if the overall system has the capacity to meet all deadlines. Load sharing (LS) is one way to alleviate this difficulty.

In this paper, we propose a decentralized, dynamic LS method for a distributed real-time system. Under this LS method, whenever the state of a node changes from lightly-loaded to overloaded and vice versa, the node broadcasts this change to a set of nodes, called a *buddy set*, in the system. An overloaded node can select, without probing other nodes, the first available node from its *preferred list*, an ordered set of nodes in its buddy set. Preferred lists are so constructed that the probability of more than one overloaded node "dumping" their loads on a single lightly-loaded node may be made very small. Performance of the proposed LS policy is evaluated with both analytic modeling and simulation. Analytic models are used to derive the distribution of queue length at each node, the probability of meeting task deadlines, and analyze the effects of buddy set size, the frequency of state change, and the average system sojourn time of each task. On the other hand, simulation is used to verify analytic results. The proposed LS method is shown to meet task deadlines with a very high probability.

Index Terms: Distributed real-time systems, load sharing, deadlines, missing probability, buddy set, preferred list, state-change broadcast.

¹Kang G. Shin is with the International Computer Science Institute, 1947 Center Street, Suite 600, Berkeley, CA 94704-1105 on leave from Real-Time Computing Laboratory, Department of Electrical Engineering and Computer Science, The University of Michigan, Ann Arbor, MI 48109-2122, which Yi-Chieh Chang is currently with.

The work reported in this paper was supported in part by the Office of Naval Research under contract N00014-85-K-0122 and the NASA under grant NAG-1-887. Any opinions, findings, and recommendations expressed in this publication are those of the authors and do not necessarily reflect the view of the funding agencies.

1 Introduction

Failure to complete a real-time task before its deadline could cause a disaster [1,2]. Due to their potential for high performance and reliability with the multiplicity of processors, distributed systems are natural candidates to implement real-time applications. However, if task arrivals are unevenly distributed over the nodes in a distributed real-time system, some nodes may become overloaded, and thus, unable to complete all their tasks in time, while other nodes are lightly-loaded. In such a case, even if the total processing power of the system is sufficient to complete all incoming tasks in time, some tasks arriving at overloaded nodes may not complete in time. One way to alleviate this problem is load sharing (LS); some of those tasks arriving at overloaded nodes are transferred to lightly-loaded nodes for execution.

LS in general-purpose distributed systems has been studied extensively by numerous researchers [3,4,5,6,7,8]. Decisions on how to share loads among the nodes are either *static* or *dynamic*. A static decision is independent of the current system state, whereas a dynamic decision depends on the system state at the time of decision. Static LS can also be viewed as the nondeterministic allocation of tasks in a system [8,9,10,11], where an overloaded node i will transfer its task to node j with probability P_{ij} , which is independent of the current system state. Although static LS is simple and easy to analyze with queueing models, its potential benefit is limited since it does not adapt itself to the time-varying system state [4]. For example, even when node i is overloaded, it still has to receive tasks from other nodes with the same probability as if it were lightly-loaded. On the other hand, when dynamic LS is used, an overloaded node can transfer its task(s) to other node(s) using the information on the current system state [4,6,7,12,13,14]. Since any dynamic policy requires each node to know states of the other nodes, it is inherently more complex than any static policy. The advantage of a dynamic policy is that it adapts itself to the time-varying system state, and thus, can ease the difficulty associated with static LS.

LS algorithms can be source-initiated or server-initiated, depending on which node initiates the transfer of task(s). The node at which external tasks arrive is the source (sender) node, and the node that processes these tasks is the server (receiver) [3]. In the source-initiated approach, an overloaded source node initiates the transfer of a newly arriving external task based on some strategies, while in the server-initiated approach, a

lightly-loaded or idle server will probe each of the potential source nodes to share its load with. LS algorithms are further divided into several levels according to the amount of information required for them. After analyzing and comparing the performance of these algorithms, Wang and Morris [3] concluded that an algorithm that collects more information will generally produce better results and that the server-initiated approach will usually outperform the source-initiated approach if the task transfer cost is not significant.

As was discussed in a recent paper by Eager *et. al* [4], LS is composed of a *transfer policy* and a *location policy*. The transfer policy determines when a node should transfer its task(s), i.e., when a node becomes overloaded. The location policy determines where a source node should send its task(s) to. Their objective was to minimize the average system response time by moving tasks from overloaded nodes to lightly-loaded ones. A simple threshold was used in the transfer policy; that is, whenever the queue length of a node exceeds this threshold, it will attempt to transfer its incoming task to another node. Three different location policies were simulated and compared: random, threshold, and shortest. Their simulation results indicated that the random policy improves system performance significantly, as compared to the system without LS. The threshold policy can further improve system performance, as compared to the random policy. The improvement of the shortest policy is about the same as that of the threshold policy, although it requires more system state information than that of the threshold policy.

LS in distributed real-time systems is addressed far less in literature than that in general-purpose distributed systems. Kurose *et. al* proposed a quasi-dynamic LS algorithm in a soft real-time system [5]. A job is considered to be *lost* if its completion time exceeds its deadline. Their primary objective was to reduce the probability of losing a job with LS. A node will transfer some of its jobs to another node if the unfinished workload exceeds its time constraints. The server was selected on a probabilistic basis which is independent of the current system state.

The location policy in most early work can be viewed as sender-initiated, since an overloaded node (sender) selects another node (a candidate receiver) and check whether or not this node can share its load. If it can, the sender will transfer some of its tasks to that node; otherwise, the sender will probe another node. This process will repeat until a receiver is found or a prespecified limit is reached. There are two major drawbacks associated with

this approach. First, the sender needs to probe other nodes before transferring any of its tasks. This will introduce an additional delay in completing the tasks to be transferred. Second, if only a few nodes in the system are lightly-loaded, the sender may not be able to locate a receiver by probing only a limited number of nodes. In such a case, overloaded nodes must execute all their tasks locally, missing the deadlines of some of these tasks.

In a real-time system, the probability of missing a job must be kept as low as possible because the loss of a job may lead to a disastrous circumstance. Note that almost all LS methods known to date are concerned only with the *average* system performance, rather than the performance of each individual task. To alleviate this weakness, we propose a new LS method in which each node needs to maintain state information of only a small set of nodes, called a *buddy set*. Whenever a node becomes overloaded (lightly-loaded) due to the arrival and/or transfer (completion) of tasks, it will broadcast its change of state to all the other nodes in its buddy set. Every node that receives this information will update its state information by eliminating the overloaded node from, or adding the lightly-loaded node to, its list of available receivers. An overloaded node can select the first node in its list of available receivers and transfer a task to that node. Notice that our LS method is completely different from conventional receiver-initiated LS methods that are characterized in [3]. A lightly-loaded server in the conventional receiver-initiated approach probes other nodes to share their work with. By contrast, our method transfers tasks from overloaded nodes to other lightly-loaded nodes using the state changes broadcast within their respective buddy sets.

The rest of this paper is organized as follows. Section 2 discusses the problem of implementing the proposed LS method. Collection of state information and construction of preferred lists are detailed in Section 3. Section 4 presents one exact model, one approximate model, and an approximate solution to the exact model for the proposed LS method. In Section 5, the performance of the proposed LS method is evaluated with the models derived in Section 4 and is also simulated to verify the analytic results. Finally, the paper concludes with Section 6.

2 Problem Statement

In the proposed LS method, each node must maintain and update the state information of other nodes. An overloaded node can transfer a task to another node based on state information without any probing delay. To implement this method, one must resolve the following issues:

- Efficient means of collecting and updating state information; collection of state information must not hamper normal communications, such as task transfer or I/O operation.
- Selection of a server node in case there are more than one lightly-loaded node.

These issues are addressed below in some detail.

2.1 State Information

To collect state information, one must decide from which nodes the information should be collected and how often this information should be updated. One straightforward method is for each node to collect and update state information from all other nodes in the system at a fixed time interval. However, it is very difficult to determine an appropriate collection and update interval which ensures the accuracy of state information while keeping below an acceptable level the I/O overhead caused by the collection of state information. Although a short interval (i.e., frequent state update) ensures the accuracy of state information, this will introduce an $O(n^2)$ I/O overhead each time for collecting state information, where n is the number of nodes in the system. This may, in turn, severely delay normal I/O operation and task transfer, thus degrading (rather than improving) system performance. On the other hand, the I/O overhead decreases as the frequency of state collection and update decreases. But, this may cause the state information recorded in a node to be obsolete. For example, if the state of a node has changed from lightly-loaded to overloaded before the next update, other nodes may transfer their tasks to this already overloaded node based on the obsolete state information.

Ideally, each node should keep state information as accurate and as up-to-date as possible while keeping the associated I/O overhead as low as possible. To achieve this goal, we

propose a *state-change broadcast* algorithm to collect and update the state information at each node. Under this algorithm, each node maintains only the state information of a small set (i.e., a buddy set) of nodes, e.g., neighbors of the node. A node is said to be *lightly-loaded* if its queue length or cumulative execution time (CET) is below a threshold, and *fully-loaded* if the queue length or CET is larger than another threshold. Each node will broadcast the change of state to all the other nodes in its buddy set only if it switched from lightly-loaded to fully-loaded, and vice versa. Since a node will receive new information only when the state of a node in its buddy set has changed, each node will have the exact state information of all the other nodes in its buddy set, as long as there is no significant delay in broadcasting state changes. (See Section 5.1.7 for more on the broadcasting delay.) Moreover, the I/O traffic for broadcasting state changes can be controlled by adjusting the two thresholds. More on this will be discussed later.

2.2 Preferred List

In the sender-initiated location policy, an overloaded node will transfer a task to the lightly-loaded node found first during the probing [3,4,5]. In our proposed location policy, an overloaded node may find more than one lightly-loaded node. One must therefore establish a rule for selecting a receiver to avoid the situation where two overloaded nodes simultaneously transfer tasks to the same lightly-loaded node. A *preferred list* is proposed to counter this situation. Since each node maintains the state information of the nodes in its buddy set, one can order these nodes in each node's preferred list. The first node in this list is the *most preferred* and the second the *second most preferred*, and so on. Note that order of preference changes with time, e.g., if the most preferred node becomes overloaded, then the second most preferred node, if not overloaded, becomes the most preferred. An overloaded node will transfer its task(s) to its most preferred node in the list. Based on the system topology, the "static" order of nodes in each node's preferred list is so permuted that a node is the most preferred of one and only one other node in the corresponding buddy set. (This order does not change with time, although some of nodes will drop out of the list of available receivers when they get overloaded, and regain their spots when they become lightly-loaded again.) Since each overloaded node is most likely to select the first node in its preferred list, the problem of more than one overloaded node "dumping" their loads on

one node is unlikely to occur. Nevertheless, dumping could occur since, for example, the third most preferred node, say N_x , of an overloaded node N_o can become its most preferred while N_x is also the most preferred of another overloaded node. But, the probability of this to happen is small, since it will occur only after overloading all the nodes ahead of N_x in the N_o 's preferred list.

3 State Information and Preferred List

It is assumed that all nodes in the system can communicate with one another via an arbitrary interconnection network. Also, each node is assumed to be stable, i.e., each node's load density is less than one. Each node has a network processor to handle the usual communications and task transfers between nodes without burdening the node processor. Each node has two sources of task arrivals, external tasks and transferred-in tasks, and one server (single node processor). The tasks arriving at each node may be executed locally or remotely at any other nodes in the system.

3.1 Collection of State Information

To determine the state of a node, three thresholds, TH_l , TH_f , and TH_h , are defined. These thresholds can be queue length or cumulative execution time (CET), depending on task characteristics. For example, if every task has the same (identically distributed) execution time, one can use the (average) queue length to measure the workload of each node. However, if task execution times are neither identical nor identically distributed, one must use the CET to measure the workload of each node. By comparing each node's current workload with these thresholds, the state of the node is determined to be in one of three states: light (L), full (F), and heavy (H). Queue length (QL) will be used to measure a node's workload throughout this paper. (Use of CET is much more involved, as pointed out in Section 6, and will be treated in a forthcoming paper.) A node is in L state if $QL \leq TH_l$, and in F state if $TH_f \leq QL \leq TH_h$, and in H state if $QL > TH_h$. An L-state node can accept one or more tasks from other nodes and complete them before their deadlines. A node in F state cannot accept tasks from any other nodes but can complete all of its own tasks in time. A node in H state cannot complete all of its own tasks in time, and thus,

must transfer, if possible, some of its task(s) to other node(s). Since a node in L state can share other nodes' loads, it is said to be in *share mode*. A node in F state will neither accept tasks from other nodes nor transfer tasks to others, and is said to be in *independent mode*. A node in H state must transfer some of its tasks, and is said to be in *transfer mode*. Note that H is usually a transient state because, if arrival of a new task at a node moves the node to H state, the node will transfer this task to another lightly-loaded node and then move back to F state. However, if an H-state node cannot find an L-state node from its buddy set, it will be forced to remain in H state and may miss some of its task deadlines.

According to our state-change broadcast algorithm, each node will broadcast the change of state to all the other nodes in its buddy set only when it moves from L to H and/or H to L. Upon receiving a state-change broadcast, every node in the corresponding buddy set will update its state information accordingly.

Two different thresholds, TH_l and TH_f , are used in the proposed LS method for the following reason. If only one threshold were used, a node would be in L (F) state when its queue length is less (greater) than this threshold. In such a case an L-state node may switch to F state after receiving a task from another node, and an F-state node may switch back to L state after completing a task. Since an L-state node will receive tasks from other nodes, it is likely to move to F state. On the other hand, an F-state node only accepts its own external tasks and is likely to switch back to L state, since every node is assumed to be stable. Thus, every node in the system will frequently switch between L and H states, thereby increasing the I/O traffic to broadcast state changes. Change of state occurs infrequently (frequently) when the difference between TH_l and TH_f is large (small).

The third threshold, TH_h , is used to avoid unnecessary task transfers. If we combine F and H states into one threshold, then acceptance of one transferred task may make a node fully- as well as heavily- loaded. In this case, the fully- and heavily- loaded node must transfer its own newly arriving task to another node. Had it not accepted the transferred task, the node would not have to transfer the newly arriving task, and thus, one of the two task transfers would not have been needed. By introducing another threshold $TH_h > TH_f$, each node will broadcast the change of state when it moves to F state, preventing other nodes from transferring tasks to that node. Since there is a non-zero range between TH_f and TH_h and every node is assumed to be stable, a node is unlikely to become heavily-

loaded with its own arriving tasks. Thus, this range can be used to control unnecessary task transfers.

The above three thresholds greatly influence system performance, such as the average task execution time, the probability of missing task deadlines, and the overhead of broadcasting state changes. Thus, these thresholds must be determined to meet the system performance requirement. For example, in a real-time system, TH_h is a critical point below which a node processor can complete all queued tasks before their deadlines with a probability higher than a specified value. The range between TH_l and TH_f must be chosen to keep the I/O overhead induced by state-change broadcasts below a specified value. These thresholds are also sensitive to system load and have to be adjusted as system load varies. (More on this will be discussed in Section 5.)

3.2 List of Preferred Nodes

As mentioned in Section 2, the purpose of constructing a preferred list at each node is to avoid the probing delay and the problem of more than one overloaded node dumping tasks on that node. The cost of task transfer is an increasing function of the physical distance between the sender and receiver nodes. To reduce this cost, the receiver node should be located as closely to the source node as possible. The preferred list of each node is thus structured based on the number of hops between the source and receiver nodes. The first entry of a node's preferred list consists of those nodes one hop away from the node, and the second entry consists of those nodes two hops away from the node, and so on. When there are more than one node in each entry, these nodes must be ordered to minimize the dumping problem.

To demonstrate the ordering of nodes in each buddy set, consider a regular² system with n nodes, N_1, N_2, \dots, N_n , where the degree of N_i is $k \forall i$. Link j of N_i is assigned a direction d_j , $0 \leq j \leq k-1$. The N_i 's "static"³ preferred list is then constructed as follows. The set of N_i 's immediate neighbors, denoted by P_1^i , is placed in the first entry of the N_i 's preferred list. The N_i 's second entry, denoted by P_2^i , consists of the nodes in the first entry

²A system is said to be *regular* if all node degrees are identical.

³This list is determined by the system topology and remains unchanged. However, the availability of each node in this list changes with time.

of every node in P_1^i , excluding the duplicated nodes. Generally, P_l^i is the set of nodes which are listed in the first entry of every node in P_{l-1}^i , excluding the duplicated nodes.

Among the nodes in P_1^i , the node in direction d_0 is chosen to be the N_i 's most preferred node in this entry of N_i , denoted by N_1^{i1} , and the node in direction d_1 is the N_i 's second most preferred node in this entry list, denoted by N_2^{i1} , and so on. The nodes in P_2^i are ordered as follows. The nodes in the N_1^{i1} 's first entry are checked according to their order in the entry. If a node in the N_1^{i1} 's first entry did not appear at any previous entry of N_i , it will be copied into the second entry of N_i in the same order as in the N_1^{i1} 's first entry. After all nodes in the first entry of N_1^{i1} are checked and copied, the nodes in the first entry of node N_2^{i1} will be checked and copied by the same procedure. This procedure will repeat until P_2^i is completed. The ordering of nodes in $P_l^i \forall l > 2$ can be determined similarly.

As an example, consider how the preferred list of each node in a 4-cube system (Fig. 1) is constructed. The identity (ID) of each node is expressed by a 4-bit number, $b_3b_2b_1b_0$. The direction d_i of node N_k is the link that connects N_k to a node whose ID differs from the N_k 's ID in bit position i , where $0 \leq i \leq 3$. One can now apply the algorithm described above to construct the preferred list for each node in the 4-cube system as shown in Fig. 2.

Once each node's preferred list is constructed, a heavily-loaded node N_i can select a lightly-loaded node as follows. Check node N_1^{i1} first; if it is lightly-loaded, N_i will transfer a task to N_1^{i1} , otherwise N_2^{i1} is checked, and so on. (This checking can easily be implemented with a pointer which is made to point to the first available node in the list.) If all nodes in P_1^i are heavily-loaded, N_i will sequentially check the nodes in P_2^i . If, albeit rare, a heavily-loaded node cannot find any lightly-loaded node from its preferred list, all of its tasks will be forced to execute locally.

The preferred list constructed above has the following advantages. First, since each node is the most preferred node of one and only one node, the probability of a lightly-loaded node being selected by more than one heavily-loaded node is very small. Second, the cost of task transfer is minimal, since a receiver node is selected, with a high probability, from the physical proximity of the source node. Moreover, the time overhead for selecting a lightly-loaded node is negligibly small, because the time-consuming probing procedure used in most known methods [4,5,7] is not needed.

Since the size of the preferred list (buddy set) will affect the probability of a task missing

its deadline, it must be chosen to ensure that this probability is lower than the specified limit. However, a buddy set must not be too large because the larger the size of buddy set, the higher I/O overhead will result. Thus, there is a trade-off between the capability of meeting task deadlines and the I/O overhead associated with state-change broadcasts. More on this will be discussed in Section 5.

4 Models for the Proposed LS Method

A modified embedded Markov chain is proposed to model the performance of the proposed LS method. We begin with an exact model from which an approximate solution and an approximate model, called the *upper bound model*, will be derived. The real solution will be shown to be (i) always upper bounded by the solution derived from the upper bound model, and (ii) very close to the approximate solution. Note that the embedded Markov chain is commonly used to analyze arbitrarily arriving tasks. Let each task take one unit of time to complete⁴, k_t be the number of task arrivals during the interval $[t, t + 1)$, and α_{k_t} be the probability of k_t arrivals in $[t, t + 1)$. For example, when the interarrival time of external tasks is exponentially distributed with rate λ , α_{k_t} can be calculated by [15]:

$$\alpha_{k_t} = \frac{\lambda^{k_t}}{k_t!} e^{-\lambda}. \quad (4.1)$$

Let x_t and x_{t+1} denote the queue lengths at time t and $t + 1$, respectively. Then,

$$x_{t+1} = \begin{cases} k_t & \text{if } x_t = 0 \text{ and } k_t \leq TH_h \\ x_t + k_t - 1 & \text{if } x_t > 0 \text{ and } x_t + k_t \leq TH_h + 1 \\ TH_h & \text{if } x_t + k_t > TH_h. \end{cases} \quad (4.2)$$

The above relation represents the case of ideal load sharing, since overloaded nodes (with more than TH_h tasks) can always find other nodes to transfer their “surplus”⁵ tasks to.

Using Eq. (4.2), one can derive the probability distribution of queue length. Two modifications must be made to include the effects of transferring and receiving tasks among the

⁴Since (average) queue length is used to measure workloads, without loss of generality, one can assume (average) task execution time to be one unit of time.

⁵Those tasks arriving after first TH_h unfinished tasks.

nodes in a buddy set. The first modification is to adjust the task arrival rate to include transferred-in tasks when a node is in L state. As shown in Fig. 3, the total arrival rate becomes $\lambda^* = \lambda^e + \lambda^t$, where λ^e and λ^t are the arrival rates of external and transferred-in tasks, respectively. A node's state transition probability depends on the total task arrival rate (λ^*) when the node is in L state, and thus, α 's must be recalculated accordingly. Let α^* 's represent the transition probability corresponding to λ^* , whereas α 's represent that corresponding to λ^e only. The second modification is made to the maximum queue length. Since a node will transfer tasks to other nodes when $QL > TH_h$, the queue length of a node with ideal LS is bounded by TH_h .

To illustrate these modifications, consider the threshold pattern "1 2 3" (i.e., $TH_l = 1, TH_f = 2, TH_h = 3$) as an example. Let q_i represent the probability of $QL = i$. Then,

$$\begin{aligned}
q_0 &= \alpha_0^* q_0 + \alpha_0^* q_1 \\
q_1 &= \alpha_1^* q_0 + \alpha_1^* q_1 + \alpha_0 q_2 \\
q_2 &= \alpha_2^* q_0 + \alpha_2^* q_1 + \alpha_1 q_2 + \alpha_0 q_3 \\
q_3 &= (1 - \alpha_0^* - \alpha_1^* - \alpha_2^*) q_0 + (1 - \alpha_0^* - \alpha_1^* - \alpha_2^*) q_1 \\
&\quad + (1 - \alpha_0 - \alpha_1) q_2 + (1 - \alpha_0) q_3 \\
q_k &= 0 \text{ for all } k > 3.
\end{aligned} \tag{4.3}$$

Note that the assumption that a task takes one unit of time to complete is used in the above equation.

As mentioned earlier, Eq. (4.3) represents ideal load sharing, i.e., an overloaded node can always locate L-state nodes to which tasks can be transferred. In reality however, an overloaded node may not be able to find any lightly-loaded node from its buddy set to share its load with. An embedded Markov chain model is developed below to handle this realistic case. In our LS method, the tasks in a node will be transferred to other nodes if its queue length exceeds $TH_h + 1$ (TH_h) upon (before) completion of a task. A node can receive tasks from other nodes only when $QL < TH_f$. To transfer surplus tasks, the sharing capacity of each buddy set must be greater than or equal to the total number of surplus tasks in that buddy set. If this condition does not hold, the queue length of an overloaded node could grow larger than TH_h . To calculate the probability of a node's queue length growing larger than TH_h , the following parameters are introduced. Let ε_i (θ_i) be the probability of

exactly (at least) i nodes being available to share the surplus tasks within a buddy set. So, $\theta_i = 1 - \sum_{k=0}^{i-1} \varepsilon_k$, for $1 \leq i \leq n$, and $\theta_i = \varepsilon_i = 0$, for $i > n$, where n is the size of buddy set. Assuming $x_t > 0$ for the previous example with threshold pattern "1 2 3", the number of surplus tasks in a node is $k_{ov} \equiv x_t + k_t - TH_h$. When $k_{ov} = 1$, $x_{t+1} = TH_h$ and the node will not transfer any task. When $k_{ov} = 2$, $x_{t+1} = TH_h$ if there is at least one node available for load sharing in its buddy set, and $x_{t+1} = TH_h + 1$ if none of the nodes in the buddy set is available for load sharing. Similarly, when $k_{ov} = \ell > 2$, $x_{t+1} = TH_h$ if there are at least $\ell - 1$ ($\leq n$) nodes available in its buddy set, or $x_{t+1} = TH_h + 1$ if there are exactly $\ell - 2$ nodes available, or, in general, $x_{t+1} = TH_h + j$ when there are exactly $\ell - (j + 1)$ nodes available. Then, the state transition relation can be rewritten as:

$$x_{t+1} = \begin{cases} k_t & \text{if } x_t = 0 \text{ and } k_t \leq TH_h \\ TH_h \text{ with prob. } \theta_{k_t-TH_h}, \text{ or } (TH_h + 1) \text{ with prob.} \\ \quad \varepsilon_{k_t-TH_h-1}, \dots, \text{ or } k_t \text{ with prob. } \varepsilon_0 & \text{if } x_t = 0 \text{ and } k_t > TH_h \\ x_t + k_t - 1 & \text{if } x_t > 0 \text{ and } x_t + k_t - 1 \leq TH_h \\ TH_h \text{ with prob. } \theta_{x_t+k_t-TH_h-1}, \text{ or } (TH_h + 1) \text{ with prob.} \\ \quad \varepsilon_{x_t+k_t-TH_h-2}, \dots, \text{ or } (x_t + k_t - 1) \text{ with prob. } \varepsilon_0 & \text{if } x_t > 0 \text{ and } x_t + k_t - 1 > TH_h. \end{cases} \quad (4.4)$$

Using the above relation, q_k 's can be derived. For example, when the threshold pattern "1 2 3" is chosen, one can derive:

$$\begin{aligned} q_3 &= \left(\alpha_3^* + \sum_{i=1}^{\infty} \theta_i \alpha_{i+3}^* \right) (q_0 + q_1) + \sum_{i=2}^4 \left(\alpha_{4-i} + \sum_{j=1}^{\infty} \theta_j \alpha_{j+4-i} \right) q_i \\ &\quad + \sum_{i=5}^{\infty} \left(\sum_{j=i-4}^{\infty} \theta_j \alpha_{j-i+4} \right) q_i \\ q_k &= \sum_{i=0}^n \varepsilon_i \alpha_{i+4}^* (q_0 + q_1) + \sum_{i=2}^{k+1} \left(\sum_{j=0}^{\infty} \varepsilon_j \alpha_{j+k-i+1} \right) q_i + \sum_{i=k+2}^{\infty} \left(\sum_{j=0}^{\infty} \varepsilon_{j+i-k-1} \alpha_j \right) q_i \\ &\quad \text{for } k = 4, \dots, \infty. \end{aligned} \quad (4.5)$$

Note that the q_k 's for $k < TH_h = 3$ are the same as shown in Eq. (4.3). q_k 's for other threshold patterns can be derived similarly.

Although the above equations can be used to calculate the distribution of queue length, ε_k 's and θ_k 's are in practice too complex to derive. For example, ε_k is the probability of having k nodes available for load sharing in an n -node buddy set, the calculation of which

requires to consider $n!/(n-k)!k!$ different possibilities. The total number of possibilities that need to be considered for the calculation of ε_k for $k = 1, \dots, n$ is 2^n . Our analysis shows this number to be over 1,000 patterns when each buddy set contains 10 to 15 nodes. Furthermore, each of these patterns needs to be analyzed separately, since the probability of a node being in L state depends on the state of other nodes in the buddy set. Thus, it is extremely tedious to derive these parameters. To alleviate this difficulty, we develop an approximate model, called the *upper bound model*, and an approximate solution to Eq. (4.4). The former is used to derive ε_k 's and a rough idea on the goodness of our LS method, while the latter to obtain (close) approximate q_k 's from Eq. (4.4) using the parameters derived from the upper bound model. The real solution will be shown to be (i) upper bounded by the solution to the upper bound model and (ii) very close to the approximate solution.

4.1 Upper Bound Model and Solution

4.1.1 Upper Bound Model

This model is derived under the assumption that every node can always transfer only one surplus task to another node. (The rest of surplus tasks are forced to queue at that node.) Since on the average a half of the computation capacity in each buddy set has to be available for load sharing,⁶ the real probability of a node being unable to transfer a surplus task is always less than that derived from this model. Moreover, the state transition relation in this model is much simpler than the exact model. If $k_{ov} = 1$ at a node, then $x_{t+1} = TH_h$ and the node will not transfer any task. If $k_{ov} = 2$, then $x_{t+1} = TH_h$ with probability $\theta_1 = 1 - \varepsilon_0$; otherwise, $x_{t+1} = TH_h + 1$ with probability ε_0 . However, when $k_{ov} = \ell$, $x_{t+1} = TH_h + \ell - 2$ with probability θ_1 , or $x_{t+1} = TH_h + \ell - 1$ with probability

⁶Otherwise, load sharing is usually infeasible and thus should not be considered.

ε_0 . Summarizing the above leads to:

$$x_{t+1} = \begin{cases} k_t & \text{if } x_t = 0 \text{ and } k_t \leq TH_h \\ (k_t - 1) \text{ with prob. } \theta_1, \text{ or } k_t \text{ with prob. } \varepsilon_0 & \text{if } x_t = 0 \text{ and } k_t > TH_h \\ x_t + k_t - 1 & \text{if } x_t > 0 \text{ and } x_t + k_t - 1 \leq TH_h \\ (x_t + k_t - 2) \text{ with prob. } \theta_1, \\ \text{or } (x_t + k_t - 1) \text{ with prob. } \varepsilon_0 & \text{if } x_t > 0 \text{ and } x_t + k_t - 1 > TH_h. \end{cases} \quad (4.6)$$

The distribution of queue length can now be derived from this equation as follows:

$$\begin{aligned} q_0 &= \alpha_0^* q_0 + \alpha_0^* q_1 \\ q_1 &= \alpha_1^* q_0 + \alpha_1^* q_1 + \alpha_0 q_2 \\ q_2 &= \alpha_2^* q_0 + \alpha_2^* q_1 + \alpha_1 q_2 + \alpha_0 q_3 \\ q_3 &= [\alpha_3^* + (1 - \varepsilon) \alpha_4^*] q_0 + [\alpha_3^* + (1 - \varepsilon) \alpha_4^*] q_1 + [\alpha_2 + (1 - \varepsilon) \alpha_3] q_2 \\ &\quad + [\alpha_1 + (1 - \varepsilon) \alpha_2] q_3 + [\alpha_0 + (1 - \varepsilon) \alpha_1] q_4 + (1 - \varepsilon) \alpha_0 q_5 \\ q_4 &= [\varepsilon \alpha_4^* + (1 - \varepsilon) \alpha_5^*] (q_0 + q_1) + [\varepsilon \alpha_3 + (1 - \varepsilon) \alpha_4] q_2 + [\varepsilon \alpha_2 + (1 - \varepsilon) \alpha_3] q_3 \\ &\quad + [\varepsilon \alpha_1 + (1 - \varepsilon) \alpha_2] q_4 + [\varepsilon \alpha_0 + (1 - \varepsilon) \alpha_1] q_5 + (1 - \varepsilon) \alpha_0 q_6 \\ q_k &= [\varepsilon \alpha_k^* + (1 - \varepsilon) \alpha_{k+1}^*] (q_0 + q_1) + \sum_{j=2}^{k+1} [\varepsilon \alpha_{k-j+1} + (1 - \varepsilon) \alpha_{k-j+2}] q_j \\ &\quad + (1 - \varepsilon) \alpha_0 q_{k+2} \text{ for all } k \geq 4, \end{aligned} \quad (4.7)$$

where $\varepsilon = \varepsilon_0$. The above equations can be rewritten in vector form: $Q = A Q$, where $Q = [q_0, \dots, q_n]^T$, A is an $n \times n$ coefficient matrix, and n is the size of buddy set. Using the above equations and $\sum_{i=0}^n q_i = 1$, one can solve for Q , which is called the *upper-bound solution*.

Solutions to the upper bound model can be shown to bound the real solution as follows. The only difference between the exact model Eq. (4.4) and the upper bound model Eq. (4.6) is that transitions to queue lengths $TH_h, TH_h + 1, \dots, (x_t + k_t - 2)$ in Eq. (4.4) are combined into a single transition to $QL = x_t + k_t - 2$ in Eq. (4.6). Since $x_t + k_t - 2 \geq TH_h$, the transition to a queue length greater than TH_h is exaggerated in Eq. (4.6). Thus, the solution to the exact model will be bounded by that to the upper bound model when $k > TH_h$. Note that the upper bound model is identical to the exact model when $k \leq TH_h$.

4.1.2 Solving Upper Bound Model

The upper bound model is analyzed first to get a rough idea on the goodness of our LS method. λ^* 's and ε must be known before solving the upper bound model for q_k 's. On the other hand, these parameters depend on q_k 's, and thus, the model cannot be solved for q_k 's without knowing λ^* 's and ε . A two-step approximation approach is taken to handle the difficulty associated with this recursion problem. In the first step, ε is set to zero, and the model is then solved for λ^t and q_k 's. The resulting q_k 's are still an upper bound for the true solution. The second step is to compute ε by using the q_k 's obtained in the first step.

By setting ε 's to zero, q_k 's for $k \geq 3$ in the upper bound model become:

$$\begin{aligned} q_3 &= (\alpha_3^* + \alpha_4^*) (q_0 + q_1) + (\alpha_2 + \alpha_3) q_2 + (\alpha_1 + \alpha_2) q_3 \\ &\quad + (\alpha_0 + \alpha_1) q_4 + \alpha_0 q_5 \\ q_4 &= \alpha_5^* (q_0 + q_1) + \alpha_4 q_2 + \alpha_3 q_3 + \alpha_2 q_4 + \alpha_1 q_5 + \alpha_0 q_6 \\ q_k &= \alpha_{k+1}^* (q_0 + q_1) + \sum_{j=2}^{k+1} \alpha_{k-j+2} q_j + \alpha_0 q_{k+2} \quad \text{for all } k > 4. \end{aligned} \quad (4.8)$$

The above equation can be solved by using an iterative method. Initially, λ^t is set to zero. One can compute q_k 's and then λ^t from:

$$\lambda^S \equiv \left[\sum_{k=4}^{\infty} (k-3) \alpha_k^* \right] (q_0 + q_1) + \sum_{i=2}^4 \left[\sum_{k=5-i}^{\infty} (k-4+i) \alpha_k \right] q_i + \sum_{i=5}^{\infty} \left[\sum_{k=0}^{\infty} k \alpha_k \right] q_i. \quad (4.9)$$

Note that λ^S is the rate of task transfer out of a node. If all nodes' external task arrival rates are identical, then $\lambda^t = \lambda^S$. Otherwise, λ^t must be calculated by Eq. (4.11). After calculating λ^t , λ^* is obtained by adding λ^e to λ^t , and then q_k 's are recalculated with the new λ^* , which will, in turn, change λ^t . This procedure will repeat until q_k 's and λ^t converge to fixed values. (The convergence will be proved later.)

Lemma 1: $\frac{d q_k}{d \lambda^*}$ satisfies the following properties:

1. $\frac{d q_0}{d \lambda^*} < 0$
2. $\sum_{k=0}^{\infty} \frac{d q_k}{d \lambda^*} = 0$

$$3. \left| \frac{d q_k}{d \lambda^*} \right| < 1 \quad \forall k.$$

Proof: Since the probability of a system being idle (q_0) will decrease as the task arrival rate increases, the first property holds. The second property holds because the sum of all q_k 's is equal to one, and thus, the sum of the variations of all q_k 's must be equal to zero. The last property can be proved by contradiction. Suppose $\left| \frac{d q_k}{d \lambda^*} \right| \geq 1$. Then q_k may become negative or greater than one, if the variation of λ^* exceeds one, a possible event when an L-state node is surrounded by more than one H-state node. Since q_k is the probability of $QL = k$, it can be neither negative nor greater than one. Contradiction. \square

Lemma 2: $0 \leq \frac{d \lambda^t}{d \lambda^*} < 1$.

Proof:

$$\begin{aligned} \frac{d \lambda^t}{d \lambda^*} &= \left[\sum_{k=4}^{\infty} (k-3) \frac{d \alpha_k^*}{d \lambda^*} \right] (q_0 + q_1) + \left[\sum_{k=4}^{\infty} (k-3) \alpha_k^* \right] \left(\frac{d q_0}{d \lambda^*} + \frac{d q_1}{d \lambda^*} \right) \\ &\quad + \sum_{i=2}^4 \left[\sum_{k=5-i}^{\infty} (k-4+i) \alpha_k \right] \frac{d q_i}{d \lambda^*} + \sum_{i=5}^{\infty} \left[\sum_{k=0}^{\infty} k \alpha_k \right] \frac{d q_k}{d \lambda^*} \\ &= (1 - \alpha_0^* - \alpha_1^* - \alpha_2^*) (q_0 + q_1) + \left[\sum_{k=4}^{\infty} (k-3) \alpha_k^* \right] \left(\frac{d q_0}{d \lambda^*} + \frac{d q_1}{d \lambda^*} \right) \\ &\quad + \sum_{i=2}^4 \left[\sum_{k=5-i}^{\infty} (k-4+i) \alpha_k \right] \frac{d q_i}{d \lambda^*} + \sum_{i=5}^{\infty} \left[\sum_{k=0}^{\infty} k \alpha_k \right]. \end{aligned}$$

According to the definitions of α_k and α_k^* , each summation in the above equation is equal to or less than the average task arrival rate which is less than one in a stable system. (Recall that each node's service rate is assumed to be unity.) By the second and third properties of Lemma 1, the sum of the last three terms will be less than one. Furthermore, the first term will be much less than one, because the first three α^* 's usually dominate the determination of transition probabilities. Thus, the lemma follows. \square

Theorem 1: q_k 's and λ^t derived from the above iterative method converge to fixed values in a finite number of steps.

Proof: Let $(i)d \lambda^*$ and $(i)d \lambda^t$ be the variations of λ^* and λ^t at the i^{th} iteration, respectively. These parameters at the $i+1^{th}$ iteration are related to those at the i^{th} iteration

by:

$$\begin{aligned} {}^{(i+1)}d\lambda^t &= \frac{d\lambda^t}{d\lambda^*} {}^{(i)}d\lambda^* \\ {}^{(i+1)}d\lambda^* &= {}^{(i+1)}\lambda^t - {}^{(i)}\lambda^t = {}^{(i+1)}d\lambda^t. \end{aligned}$$

Since $|\frac{d\lambda^t}{d\lambda^*}| < 1$ by Lemma 4, $d\lambda^t$ at the $(i+1)^{th}$ iteration will be smaller than $d\lambda^*$ at the i^{th} iteration. Since the variation of λ^* at the $(i+1)^{th}$ iteration is equal to that of λ^t at the $(i+1)^{th}$ iteration, we get ${}^{(i+1)}d\lambda^t < {}^{(i)}d\lambda^t$ and ${}^{(i+1)}d\lambda^* < {}^{(i)}d\lambda^*$. Thus, the variation of λ^t will decrease to zero after a finite number of iterations, and so is λ^* . Substituting the convergent λ^t and λ^* into Eq. (4.8), unique q_k 's can be determined, i.e., q_k 's converge, too.

□

In fact, λ^t and λ^* are shown to converge after only two to three iterations in our analysis. This indicates that the derivatives of λ^t and q_k with respect to λ^* are much smaller than 1.

4.1.3 Derivation of ε

The main difficulty in deriving ε_k 's lies in the fact that the queue lengths in a buddy set depend on one another. Thus, the dependent LS environment is converted to an independent environment by using the Bayes theorem. To facilitate the description of our approach for an $(n+1)$ -node buddy set, it is necessary to introduce the following variables.

- N_j^i : the j^{th} preferred node of N_i .
- x_i : the N_i 's queue length.
- x_{ij} : the queue length of N_j^i .
- x_{jk}^i : the queue length of the k^{th} preferred node of N_j^i .
- λ_i^S : the rate of task transfer out of N_i .
- λ_{ij}^S : the rate of task transfer out of N_j^i .
- $\lambda_i^{S_b}$: the rate of task transfer out of N_i given that N_i is not in sharing mode.
- $\lambda_{ij}^{S_b}$: the rate of task transfer out of N_j^i given that N_j^i is not in sharing mode.
- λ_i^t : the rate of task transfer into N_i .

It is easy to see that $\lambda^{S_b} > \lambda^S$, since tasks are not actually transferred out of a node unless the node is in H state and λ^S is the average transfer-out rate over the entire time period of interest.

Let N_0 be the node under consideration, then

$$\varepsilon = P(x_{01} \geq TH_f, \dots, x_{0n} \geq TH_f) \quad (4.10)$$

$$\begin{aligned} \lambda_0^t &= \lambda_{01}^S P(x_0 \leq TH_l) + \lambda_{02}^S P(x_0 \leq TH_l, x_{21}^0 \geq TH_f) \\ &\quad + \lambda_{03}^S P(x_0 \leq TH_l, x_{31}^0 \geq TH_f, x_{32}^0 \geq TH_f) + \dots \\ &\quad + \lambda_{0n}^S P(x_0 \leq TH_l, x_{n1}^0 \geq TH_f, x_{n2}^0 \geq TH_f, \dots, x_{nn}^0 \geq TH_f). \end{aligned} \quad (4.11)$$

Note that λ_0^t derived from Eq. (4.11) is equivalent to that derived from Eq. (4.9) if all nodes have the same external task arrival rate. Using Eq. (4.9) in such a case, for $j = 1, \dots, n$

$$\lambda_{0j}^{S_b} = \sum_{i=2}^4 \left[\sum_{k=5-i}^{\infty} (k-4+i) \alpha_k \right] \frac{q_i}{P^{nsh}} + \sum_{i=5}^{\infty} \left[\sum_{k=0}^{\infty} k \alpha_k \right] \frac{q_i}{P^{nsh}}, \quad (4.12)$$

where $P^{nsh} = 1 - q_0 - q_1$. Using the Bayes formula, the probability of both N_1^0 and N_2^0 not being in sharing mode can be calculated by

$$P(x_{01} \geq TH_f, x_{02} \geq TH_f) = P(x_{01} \geq TH_f) P(x_{02} \geq TH_f | x_{01} \geq TH_f). \quad (4.13)$$

Since the dependence between queue lengths is included in λ^* , its effect can be reflected by adjusting the rate of task transfer into N_2^0 given that N_1^0 is not in sharing mode. So, the conditional probability $P(x_{02} \geq TH_f | x_{01} \geq TH_f)$ can be equated to $P(x_{02} \geq TH_f)$, while λ^* of N_2^0 must be adjusted to reflect the effect of N_1^0 's unavailability. As shown in Eq. (4.12), such an adjustment will increase the rate of task transfer out of N_1^0 given that it is not in sharing mode, which will, in turn, increase N_2^0 's task transfer-in rate. Moreover, N_0 will select N_2^0 as the most preferred node given that N_1^0 is not in sharing mode, and thus, the task transfer-in rate of N_2^0 should be recalculated. For convenience, let N_2 represent the node N_2^0 under consideration, then

$$\begin{aligned} \lambda_2^t &= \lambda_{21}^S P(x_2 \leq TH_l) + \lambda_{22}^S P(x_2 \leq TH_l) + \lambda_{23}^S P(x_2 \leq TH_l, x_{31}^2 \geq TH_f, x_{32}^2 \geq TH_f) \\ &\quad + \dots + \lambda_{2n}^S P(x_2 \leq TH_l, x_{n1}^2 \geq TH_f, x_{n2}^2 \geq TH_f, \dots, x_{nn}^2 \geq TH_f). \end{aligned} \quad (4.14)$$

The first two terms of Eq. (4.14) represent the transferred-in tasks from the N_2 's most preferred node and N_0 ($= N_2^2$). Since N_1^0 is unavailable, N_2 becomes the most preferred

node of both N_1^2 and N_0 . Clearly, N_2 's λ^* will be larger than those of N_0 and N_1^0 . Hence, it is likely to switch to no-sharing mode when N_1^0 is in no-sharing mode. Similarly, the probability of all N_1^0 , N_2^0 , and N_3^0 not being in sharing mode can be calculated as:

$$\begin{aligned}
P(x_{01} \geq TH_f, x_{02} \geq TH_f, x_{03} \geq TH_f) &= P(x_{01} \geq TH_f, x_{02} \geq TH_f) \times \\
&\quad P(x_{03} \geq TH_f \mid x_{01} \geq TH_f, x_{02} \geq TH_f) \\
&= P(x_{03} \geq TH_f) P(x_{01} \geq TH_f, x_{02} \geq TH_f) \\
&= P(x_{03} \geq TH_f) P(x_{01} \geq TH_f) P(x_{02} \geq TH_f).
\end{aligned}$$

λ^* of N_2^0 and N_3^0 must be recalculated as described above. The correctness of Eq. (4.11) can be verified as follows. When all nodes in the system have the same distribution of queue length and the same λ^S , Eq. (4.11) can be simplified as:

$$\begin{aligned}
\lambda_0^t &= \lambda_{01}^S P(x_0 \leq TH_l) + \lambda_{02}^S P(x_0 \leq TH_l) P(x_{21}^0 \geq TH_f) + \dots \\
&\quad + \lambda_{0n}^S P(x_0 \leq TH_l) P(x_{n1}^0 \geq TH_f) P(x_{n2}^0 \geq TH_f) \dots P(x_{nn}^0 \geq TH_f) \\
&= \lambda^S P(x_0 \leq TH_l) \left[1 + P(x_{21}^0 \geq TH_f) + P(x_{31}^0 \geq TH_f)^2 + \dots + P(x_{n1}^0 \geq TH_f)^n \right] \\
&= \lambda^S P(x_0 \leq TH_l) \frac{1 - P(x_{21}^0 \geq TH_f)^{n+1}}{1 - P(x_{21}^0 \geq TH_f)} \simeq \lambda^S P(x_0 \leq TH_l) \frac{1}{1 - P(x_{21}^0 \geq TH_f)} \\
&= \lambda^S P(x_0 \leq TH_l) \frac{1}{P(x_{21} \leq TH_l)} \simeq \lambda^S.
\end{aligned}$$

Consider a 4-cube system as an example, in which, without loss of generality, N_0 can be viewed as the center node for the derivation of ϵ . From Fig. 2, the N_0 's preferred list is $N_1 N_2 N_4 N_8 N_6 N_{10} N_{12} N_3 N_5 N_9 N_{14} N_{13} N_{11} N_7$. The first 4 nodes are N_0 's immediate neighbors. Since the nodes near the end of the list are unlikely to be selected for load sharing, the adjusted task transfer-in rate of these four nodes can be approximated by adding $\lambda_0^S P(x_0 \leq TH_l)$ to the λ^t of these nodes. Since increasing task transfer-in rate will change queue length, the q_k 's of these nodes need to be recalculated for

$$\begin{aligned}
P(x_1 \geq TH_f, x_2 \geq TH_f, x_4 \geq TH_f, x_8 \geq TH_f) &= P(x_1 \geq TH_f) P(x_2 \geq TH_f) \times \\
&\quad P(x_4 \geq TH_f) P(x_8 \geq TH_f) \\
&\equiv \epsilon^{(4)}. \tag{4.15}
\end{aligned}$$

Note that the states of these four nodes are different from that of N_0 , because their task transfer-in rates are higher than that of N_0 . Thus, these nodes are more likely to be in

H-state than N_0 . Similarly, one can calculate the adjusted task transfer-in rates for N_5 — N_{10} . As shown in Fig. 2, each of these nodes has two of the previous four nodes in its entry-1 list. Furthermore, as the number of H-state nodes increases, tasks will be transferred to a less preferred node of N_0 . The adjusted task transfer-in rates of these nodes are:

$$\begin{aligned}
\lambda_6^t &= \lambda^t + \lambda_2^{S_b} P(x_6 \leq TH_l) P(x_7 \geq TH_f) + \lambda_4^{S_b} P(x_6 \leq TH_l) P(x_7 \geq TH_f) \\
&\quad + \lambda_0^S P(x_6 \leq TH_l) \\
\lambda_{10}^t &= \lambda^t + \lambda_8^{S_b} P(x_{10} \leq TH_l) P(x_{11} \geq TH_f) + \lambda_0^S P(x_{10} \leq TH_l) \\
&\quad + \lambda_2^{S_b} P(x_{10} \leq TH_l) P(x_{11} \geq TH_f) P(x_{14} \geq TH_f) \\
&\quad + \lambda_6^{S_b} P(x_{10} \leq TH_l) P(x_{11} \geq TH_f) P(x_{12} \geq TH_f) P(x_{14} \geq TH_f) \\
\lambda_{12}^t &= \lambda^t + \lambda_8^{S_b} P(x_{12} \leq TH_l) P(x_{13} \geq TH_f) P(x_{14} \geq TH_f) + \lambda_0^S P(x_{12} \leq TH_l) \\
&\quad + (\lambda_4^{S_b} + \lambda_6^{S_b} + \lambda_{10}^{S_b}) P(x_{12} \leq TH_l) P(x_{13} \geq TH_f) P(x_{14} \geq TH_f) \\
\lambda_3^t &= \lambda^t + \lambda_1^{S_b} P(x_3 \leq TH_l) + \lambda_2^{S_b} P(x_3 \leq TH_l) + \lambda_0^S P(x_3 \leq TH_l) \\
&\quad + (\lambda_5^{S_b} + \lambda_9^{S_b}) P(x_3 \leq TH_l) P(x_7 \geq TH_f) P(x_{11} \geq TH_f) \\
&\quad + (\lambda_6^{S_b} + \lambda_{10}^{S_b}) P(x_3 \leq TH_l) P(x_7 \geq TH_f) P(x_{11} \geq TH_f) P(x_{15} \geq TH_f) \\
\lambda_5^t &= \lambda^t + \lambda_4^{S_b} P(x_5 \leq TH_l) + \lambda_1^{S_b} P(x_5 \leq TH_l) P(x_7 \geq TH_f) + \lambda_0^S P(x_5 \leq TH_l) \\
&\quad + \lambda_3^{S_b} P(x_5 \leq TH_l) P(x_7 \geq TH_f) P(x_{13} \geq TH_f) \\
&\quad + (\lambda_6^{S_b} + \lambda_{12}^{S_b}) P(x_5 \leq TH_l) P(x_7 \geq TH_f) P(x_9 \geq TH_f) P(x_{13} \geq TH_f) P(x_{15} \geq TH_f) \\
\lambda_9^t &= \lambda^t + \lambda_8^{S_b} P(x_9 \leq TH_l) + \lambda_1^{S_b} P(x_9 \leq TH_l) P(x_{11} \geq TH_f) P(x_{13} \geq TH_f) + \lambda_0^S P(x_9 \leq TH_l) \\
&\quad + (\lambda_3^{S_b} + \lambda_5^{S_b} + \lambda_{10}^{S_b} + \lambda_{12}^{S_b}) P(x_9 \leq TH_l) P(x_{11} \geq TH_f) P(x_{13} \geq TH_f) P(x_{15} \geq TH_f).
\end{aligned}$$

Once the task transfer-in rates of entry-2 nodes are adjusted, the probability of having all entry-1 and entry-2 nodes in no-sharing mode can be calculated as:

$$\begin{aligned}
P(x_1 \geq TH_f, x_2 \geq TH_f, \dots, x_9 \geq TH_f) &= P(x_1 \geq TH_f) \times \\
&\quad P(x_2 \geq TH_f) \cdots P(x_9 \geq TH_f) \quad (4.16) \\
&\equiv \epsilon^{(10)}.
\end{aligned}$$

Similarly, one can calculate the probability of all other nodes in a 4-cube system being unavailable as $\epsilon^{(15)} \equiv \epsilon$.

4.2 Approximate Solution

Although q_k 's and ε ($= \varepsilon_0$) can be derived from the upper bound model, it is still very tedious to calculate $\varepsilon_k \forall k > 0$, because there are too many possibilities to be considered and each of them is difficult to analyze due to the mutual dependence of each node's load sharing in a buddy set. Moreover, the solution to the upper bound model fails to include the effects of buddy set size and threshold patterns on the capability of meeting tasks' deadlines, while the simulation results in Section 5 did show significant differences when these parameters were changed. So, it is necessary to derive a solution which is simple but closer to the real solution to Eq. (4.4) than the upper bound solution. Since there are $n!/(n-k)!k!$ possibilities in calculating ε_k in an n -node buddy set, these possibilities can be approximated by only one possibility in which a node is in no-sharing mode with the largest probability. This possibility occurs when all other nodes in the buddy set are in no-sharing mode.

Consider the N_0 's preferred list in Fig. 2 again. The probabilities of N_2 to N_9 being in no-sharing mode are different from one another due to the adjustment of task transfer-in rates given that more preferred nodes are in no-sharing mode. As the number of no-sharing nodes increases, the adjusted task transfer-in rate of the next preferred node increases. Eventually, N_7 , the least preferred node of N_0 , will receive the largest number of transfer-in tasks, thus moving it in no-sharing mode with the highest probability within N_0 's buddy set. For convenience, let P^{sh} and P^{nsh} denote the probabilities of N_7 being in sharing and no-sharing mode, respectively. Then, ε_k 's can be approximated by:

$$\varepsilon_k = \frac{n!}{(n-k)!k!} (P^{nsh})^{n-k} (P^{sh})^k. \quad (4.17)$$

Substituting ε_k 's derived from Eq. (4.17) into Eq. (4.5) and applying the iterative method discussed in the previous subsection, we can easily obtain an approximate solution. The calculated results are listed in Tables 1 and 3 in comparison with the results derived from the upper bound model and simulations (to be discussed in the next section).

Note that the ε_k 's derived from Eq. (4.8) are essentially the same as those derived from Eq. (4.7) since both have the same queue state equations for $QL < TH_h$ which are dominant in the probability calculation.

5 Performance Analysis

The performance of the proposed LS policy is evaluated with the upper bound model, the approximate solution, and simulation. The first two are used to derive the distribution of queue length at each node, the probability of meeting task deadlines, and analyze the effects of buddy set size, the frequency of state changes, and the average system sojourn time of each task. On the other hand, simulation is used to verify analytic results.

5.1 Analytic Results

The proposed queueing model can be applied to any arrival process, but the transition probability α_k must be given prior to the calculation of q_k 's with Eqs. (4.4) and (4.8). To demonstrate the main idea of our LS method, we present some numerical results for the case when both arrivals of external and transferred-in tasks follow exponential distributions. (Note, however, that our LS method and models are not restricted to exponential distributions.)

5.1.1 Distribution of Queue Length

The distributions of queue length for two different external task arrival rates in a 16-node system are calculated from the upper bound model and the approximate solution, and compared with simulation results as well as with the case of no load sharing (Table 1). The q_k 's calculated from the upper bound model and the approximate solution are very close to each other when $k \leq TH_h$. This was expected because the two differ only when $k > TH_h$. This fact also ensures the accuracy in calculating ε , since it was computed with the q_k 's derived from the upper bound model and then used to derive approximate q_k 's from Eq. (4.4). Moreover, the distribution of queue length obtained via simulation is shown to be very close to the approximate solution for all k and is bounded by that obtained from the upper bound model when $k > TH_h$. Since the approximate solution is always very close to the real solution, we will use it in all the following discussions unless stated otherwise.

5.1.2 Probability of Meeting Deadlines

A task is said to be *missed* if its system sojourn time⁷ exceeds a given deadline. According to our queueing model, the completion time of a newly arriving task is equal to the current queue length plus one unit of time. Since the probability of $QL > TH_h$ is quite small, one can choose TH_h to be one less than the given deadline such that the probability of missing deadlines, or simply called the *missing probability*, becomes the probability of encountering $QL > TH_h$ at the time of a task arrival. Clearly, the missing probability depends on the given deadline and system load. However, by selecting proper threshold pattern and buddy set size, it is possible to minimize the missing probability. Figs. 4 and 5 show plots of missing probabilities vs. task deadlines for different threshold patterns.

Generally, the missing probability increases as system load gets heavier (Fig. 6) and/or the deadline gets shorter. By choosing an appropriate threshold pattern, e.g., "1 2 3" in Figs. 4–6, the missing probability can be reduced to a small value even when system load changes (except when the system is heavily-loaded, e.g., $\lambda^e \geq 0.9$). The analytic results also show that the choice of a threshold pattern is sensitive to system load. For example, threshold pattern "0 1 1" results in a small missing probability when the system is lightly-loaded, while resulting in a much higher missing probability as system load increases. Threshold pattern "1 2 3" is found to yield a reasonably small missing probability for a wide range of load density ($0 < \lambda^e < 0.8$). Fig. 4 shows an interesting result of the upper bound model: missing probabilities for different threshold patterns are quite close to each other, and thus, difficult to distinguish which pattern is better over the others. This is opposite to what has been shown by the approximate solution in Figs. 5 and 6. That is, the upper bound model exaggerates the probability of switching to a queue length greater than TH_h , and thus, the effect of threshold pattern becomes unimportant. Since threshold pattern "1 2 3" exhibits the best performance among the three patterns considered, the performance with this pattern is further compared with simulation results. As shown in Figs. 7 and 8, the missing probability obtained from the simulation is always upper bounded by those obtained from the upper bound model and is very close to the approximate solution.

⁷The system sojourn time of a task is composed of its execution time, queueing time, and task transfer time.

5.1.3 Average System Sojourn Time vs. Missing Probability

The average system sojourn time can be obtained by dividing the sum of all tasks' system sojourn times by the total number of tasks processed. Mathematically, the average system sojourn time is equal to the expected task execution time, $\sum_{k=0}^{\infty} (k+1)q_k$. The average system sojourn time is calculated for several different threshold patterns and buddy set sizes as presented in Table 2. One interesting result found in this calculation is that the lower TH_l and TH_h , the smaller the average system sojourn time results, and that buddy set size shows only minor effects on the average system sojourn time. This is in sharp contrast with the results reported in [4], where the average system sojourn time under the shortest queue policy was shown to be only slightly smaller than that under the threshold policy. In our algorithm, the shortest queue (threshold) policy is equivalent to selecting $TH_l = 0$ and $TH_f = 1$ ($TH_l > 0$). As shown in Table 2, the threshold pattern with $TH_l = 0$ and $TH_f = 1$ always results in a substantially smaller average system sojourn time than the pattern with $TH_l > 0$. This is the advantage resulting from our state-change broadcast since the I/O overhead for collecting state information in the case of $TH_l = 0$ is essentially the same as the case of $TH_l > 0$. However, the I/O overhead associated with the shortest queue policy is higher than that of the threshold policy due to its required probing of other nodes [4], offsetting the potential gain to be made by transferring tasks to a node with the shortest queue. Consequently, our LS algorithm outperforms sender-initiated LS algorithms even when the average system time is used to measure their performance.

Another important result is that a threshold pattern that results in a lower average system sojourn time does not always yield a lower missing probability. For example, consider the buddy sets of size 10 in Tables 2 and 3. Pattern "0 1 2" results in a smaller average task system sojourn time than "1 2 2", but a larger missing probability than "1 2 2" when the deadline is greater than 2. Moreover, some thresholds may result in almost the same average task system sojourn time but yield quite different missing probabilities, e.g., "0 2 3" and "1 2 3" when the deadline is greater than 3 in Table 3. Hence, those approaches based on minimizing the average task system sojourn time alone may not be applicable to the analysis of real-time systems.

5.1.4 System Utilization

The system utilization is defined as the ratio of external task arrival rate (λ^e) to the system service rate, which is unity in our LS model. (Thus, the system utilization is simply λ^e .) Since the missing probability depends on system workload (Fig. 8), we can solve Eq. (4.8) to derive λ^e as a function of q_k 's and then the maximum system utilization can be obtained by equating q_{TH_k+1} to the specified missing probability. Some of calculated results are plotted in Fig. 9. This is in sharp contrast to the common notion that real-time systems have to be designed to sacrifice utilization for a lower missing probability.

5.1.5 Buddy Set Size and Preferred List

The effect of changing buddy set size on the missing probability can best be explained by the approximate solution as shown in Fig. 10. Buddy set size affects the missing probability significantly when it grows from 4 to 10 and $\lambda^e > 0.7$, but its incremental effect becomes insignificant beyond 10. Actually, there is little notable decrease in the missing probability when buddy set size grows beyond 15. Surprisingly, the missing probability for a 4-node buddy set is about three orders of magnitude less than those without LS when the system is lightly loaded ($\lambda^e = 0.5$), and is about the same as those for buddy sets of size larger than 10 when the system is heavily loaded ($\lambda^e > 0.8$). So, buddy set size can be chosen to range from 10 to 15, regardless of the system size. The most interesting result is found to be that the missing probability in a large system (of 64 nodes in Table 4) is much smaller than that of a small system (of 16 nodes in Table 3). For example, consider threshold "1 2 3" with a 10-node buddy set at $\lambda^e = 0.8$. The missing probability of a 64-node system is about 3, 4, and 20 times smaller than that of the 16 nodes system when the deadline is 4, 5, and 6, respectively. This significant improvement was found for all other threshold patterns, and thus, it is concluded that the larger the system size, the better the performance of the proposed LS method will result. Note that the I/O overhead for broadcasting state changes remains unchanged and independent of system size because buddy set size is fixed (to 10–15). Furthermore, the incremental decrease in missing probability becomes insignificant when buddy set size is over 15 (Table 4).

Use of buddy sets and preferred lists in our LS algorithm plays a major role in lowering

the missing probability in a large system. As discussed in Section 3, the buddy set of a node consists of the nodes in its physical proximity, and each node in the buddy set is selected according to the order of its preference. Moreover, preferred lists are constructed in such a way that each node is the i^{th} ($i = 1 \dots n$) preferred node of only one other node and the preferred lists of the nodes in the same buddy set are completely different from each other. As a result, the surplus tasks within each buddy set will be evenly shared by all lightly-loaded nodes in the system, rather than overloading a few lightly-loaded nodes within the buddy set. As the system size increases, the percentage of common nodes in the preferred lists of a buddy set gets smaller, and thus, the surplus tasks are more evenly distributed in the system, resulting in a better performance.

5.1.6 Frequency of State Change

In our LS method, each node needs to broadcast change of state to all the other nodes in its buddy set. Since a state change occurs when a node switches from L-state to H-state and vice versa, the probability of a state change becomes:

$$P(x_{k_i+1} \leq TH_l | x_{k_i} \geq TH_f) + P(x_{k_i+1} \geq TH_f | x_{k_i} \leq TH_l).$$

The computation results of Table 5 showed that the frequency of state change can be reduced to 10 – 15% of the total number of arrived tasks by setting $TH_f - TH_l = 2$. Note that this frequency becomes about 100% of the number of external task arrivals when the threshold pattern is set to $TH_l = 0$ and $TH_f = 1$. The resulting high frequency of state change should rule out this type of threshold patterns.

The I/O overhead for collecting state information in our LS method is determined by the frequency of state change and buddy set size, while it was determined by the number of task transfers and probing in [4]. Since the frequency of state change can be controlled by adjusting the difference between TH_l and TH_f , this frequency with threshold “1 3 3” and $\lambda^e > 0.7$ is found to be about the same as the percentage of external arrivals that are transferred out. Moreover, transferring one task may require to probe 5 to 6 other nodes [4] and each probe generates two I/O messages (one for request and one for response) in sender-initiated methods, whereas each state-change broadcast in our LS method generates n messages, where n is the buddy set size. The I/O overhead for broadcasting state changes

for threshold "1 3 3" and a 10-node buddy set is about the same as that in a sender-initiated approach. However, the time for selecting a destination node in our method is much shorter than that in any sender-initiated approach, because, in our approach, transfer of a task can be made upon its arrival without probing any other nodes. Besides, each task transfer in our LS method will take less time than other LS methods, because use of a preferred list will usually locate a receiver in the sender's physical proximity.

5.1.7 Task Transfer and Broadcasting Delays

When a node selects, and transfers a task to, an L-state node, the L-state node may receive an external task and move to H-state before receiving the transferred task. In this case, the transferred task will actually arrive at a node in H-state. Thus, the transition probability with non-zero task arrivals (α_k for $k > 0$) to H-state is larger than that used in Eq. (4.5), where transferred tasks are assumed to be accepted only when a receiving node is in L-state. One can estimate this probability and adjust the corresponding α_k 's. The broadcasting delay has the same effect as task transfer delay.

Although these delays may affect the queue distribution, the missing probability can be made insensitive to them by properly choosing a threshold pattern. For example, a task which arrives when $QL = TH_f$ will not be transferred again if $TH_h > TH_f$, e.g., threshold "1 2 3", but it will be retransferred if $TH_h = TH_f$, e.g., threshold "1 2 2". Since retransferring tasks induces higher I/O overheads without improving the capability of meeting the deadline, the threshold patterns that are sensitive to these delays may result in a higher missing probability than those that are not. The effect of these delays on the missing probability is investigated further in our simulation.

5.2 Simulation Results

For our simulation the average load density is changed from 0.5 to 0.9, and buddy sets of size 4, 10, and 15 are considered. Ten threshold patterns are chosen out of all possible combinations for the simulation of a 4-cube system and the results are given in all tables except for Table 2. A few selected thresholds for a 6-cube system are also simulated and listed in Table 4. The time for transferring a task between two nodes within a buddy set is

assumed to be 10% of the task execution time and the time for broadcasting a state change to one of the nodes in a buddy set is assumed to be 1% of the task execution time.

In most cases, simulation results are consistent with, and close to, the approximate solution. However, the analytically derived q_k 's for $k > TH_h$ are always less than those obtained from simulation when the system is lightly-loaded ($\lambda^e \leq 0.5$) or heavily-loaded ($\lambda^e \geq 0.9$), especially in the threshold patterns with $TH_l > 0$. This discrepancy may have been caused by the delays of transferring a task and broadcasting state changes. The effect of setting TH_h to be larger than TH_f is also observed in the simulation. The percentage of task retransfers is higher when $TH_f = TH_h$, but lower when $TH_f < TH_h$. Since retransferring tasks will not improve performance but increase I/O overheads, the threshold pattern with $TH_f < TH_h$ is a better choice than the pattern with $TH_f = TH_h$. This observation also explains why the missing probability associated with threshold "1 2 3" is smaller than that of "1 2 2" when $\lambda^e > 0.7$, deadline > 3 , and buddy set size > 10 .

To study the effect of changing task transfer costs, we ran simulations with task transfer costs 5, 10, 20, and 30% of the task execution time. As shown in Table 6, the missing probability of threshold "1 2 3" remains almost unchanged. Summing up all previous results, threshold "1 2 3" appears to be suitable for a wide range of system load. Note that, although the missing probability of threshold "1 2 2" is usually close to that of threshold "1 2 3", the task transfer rate associated with "1 2 2" is much higher than that with "1 2 3". Thus, considering cost-performance efficiency, threshold "1 2 3" is a better choice than "1 2 2".

5.3 Remarks

There are several advantages of using the upper bound model and the approximate solution, as compared to simulations. First, the results derived from the upper bound model can be used to guarantee the specified system reliability, because the actual missing probability is always less than that derived from the upper bound model. Second, system utilization can be analyzed by using the analytic models. Third, our analytic models provide, at almost no cost, many pieces of useful information with high accuracy. For example, any meaningful simulation of our LS method requires hundreds of CPU hours (in a computer as powerful as VAX-11/780) to get an accuracy of 10^{-6} in the calculation of q_k 's for a system of moderate

size. Moreover, simulation may be able to provide information only for a particular system workload; it is too costly to generate q_k 's with simulation as a function of system workload.

6 Conclusions

We have proposed and analyzed a new LS method based on the broadcast of state changes. By selecting an appropriate threshold pattern and buddy set, one can reduce the missing probability to a small number, and thus, the proposed LS method has high potential use for various real-time applications. The I/O overhead for broadcasting state changes can be controlled to an acceptable level by selecting an appropriate threshold pattern, thereby making the proposed LS method be cost-effective.

There are several issues worth further investigation. First, it is necessary, but difficult, to derive an exact analytic formula for the probability of a task missing its deadline. Second, if each task has a different execution time, queue length is not sufficient to determine the workload at each node. In such a case, one must consider the actual task execution times and use the cumulative execution time to determine the load of each node. Furthermore, if a node thinks itself to be lightly-loaded and broadcasts its availability to other nodes, it may receive a task whose computation is too involved for the node to complete in time. Thus, the state of a node must contain a sufficient amount of information to ensure that the lightly-loaded node can process all transferred tasks in time. Third, if the task execution time is a random variable, a continuous-time Markov model must be used to simulate and analyze system performance.

Optimization of the tradeoffs existing in the proposed LS method is an interesting design problem of its own. For example, there is a tradeoff between the buddy set size and load sharing capability. The I/O overhead associated with state-change broadcasts can be reduced by shrinking the buddy set size, but this will limit the load sharing capability. All of these issues are matters of our future inquiry.

References

- [1] K. G. Shin, C. M. Krishna, and Y.-H. Lee, "A unified method for evaluating real-time computer controllers and its application", *IEEE Trans. on Auto. Contr.*, vol. AC-30,

no. 4, pp. 357-366, April 1985.

- [2] D. W. Leinbaugh, "Guaranteed response times in a hard-real-time environment", *IEEE Trans. Software Engr.*, vol. SE-6, no. 1, pp. 85-93, January 1980.
- [3] Y.-T. Wang and R. J. T. Morris, "Load sharing in distributed systems", *IEEE Trans. Comput.*, vol. C-34, no. 3, pp. 204-217, March 1985.
- [4] D. L. Eager, E. D. Lazowska, and J. Zahorjan, "Adaptive load sharing in homogeneous distributed systems", *IEEE Trans. Software Engr.*, vol. SE-12, no. 5, pp. 662-675, May 1986.
- [5] J. F. Kurose, S. Singh, and R. Chipalkatti, "A study of quasi-dynamic load sharing in soft real-time distributed computer systems", *IEEE Real-Time Systems Symposium*, pp. 201-208, 1986.
- [6] P. S. Yu, S. Balsamo, and Y.-H. Lee, "Dynamic transaction routing in distributed database systems", *IEEE Trans. Software Engr.*, vol. SE-14, no. 9, pp. 1307-1318, September 1988.
- [7] P. Krueger and R. Finkel, "An adaptive load balancing algorithm for a multicomputer", *computer science technical report # 539, University of Wisconsin-Madison*, 1987.
- [8] A. Kratzer and D. Hammerstorm, "A study of load levelling", *IEEE Real-Time Systems Symposium*, pp. 647-652, 1980.
- [9] A. N. Tantawi and D. Towsley, "Optimal static load balancing in distributed computer systems", *Journal of ACM*, pp. 445-465, April 1985.
- [10] L. M. Ni and K. Hwang, "Optimal load balancing in a multiple processor system with many job systems", *IEEE Trans. Software Engr.*, vol. SE-11, no. 5, pp. 491-496, May 1985.
- [11] J. F. Kurose and S. Singh, "A distributed algorithm for optimum static load balancing in distributed computer systems", *IEEE Real-Time Systems Symposium*, pp. 458-467, 1985.
- [12] Y.-C. Chow and W. H. Kohler, "Models for dynamic load balancing in a heterogeneous multiple processor system", *IEEE Trans. Comput.*, vol. C-28, no. 5, , May 1979.
- [13] H. S. Stone, "Multiprocessor scheduling with the aid of network flow algorithms", *IEEE Trans. Software Engr.*, vol. SE-3, no. 1, pp. 85-93, January 1977.
- [14] H. S. Stone, "Critical load factors in two-processor distributed systems", *IEEE Trans. Software Engr.*, vol. SE-4, no. 3, pp. 254-258, May 1978.
- [15] L. Kleinrock, *Queueing Systems Vol. I: Theory.*, New York: Wiley, 1975.

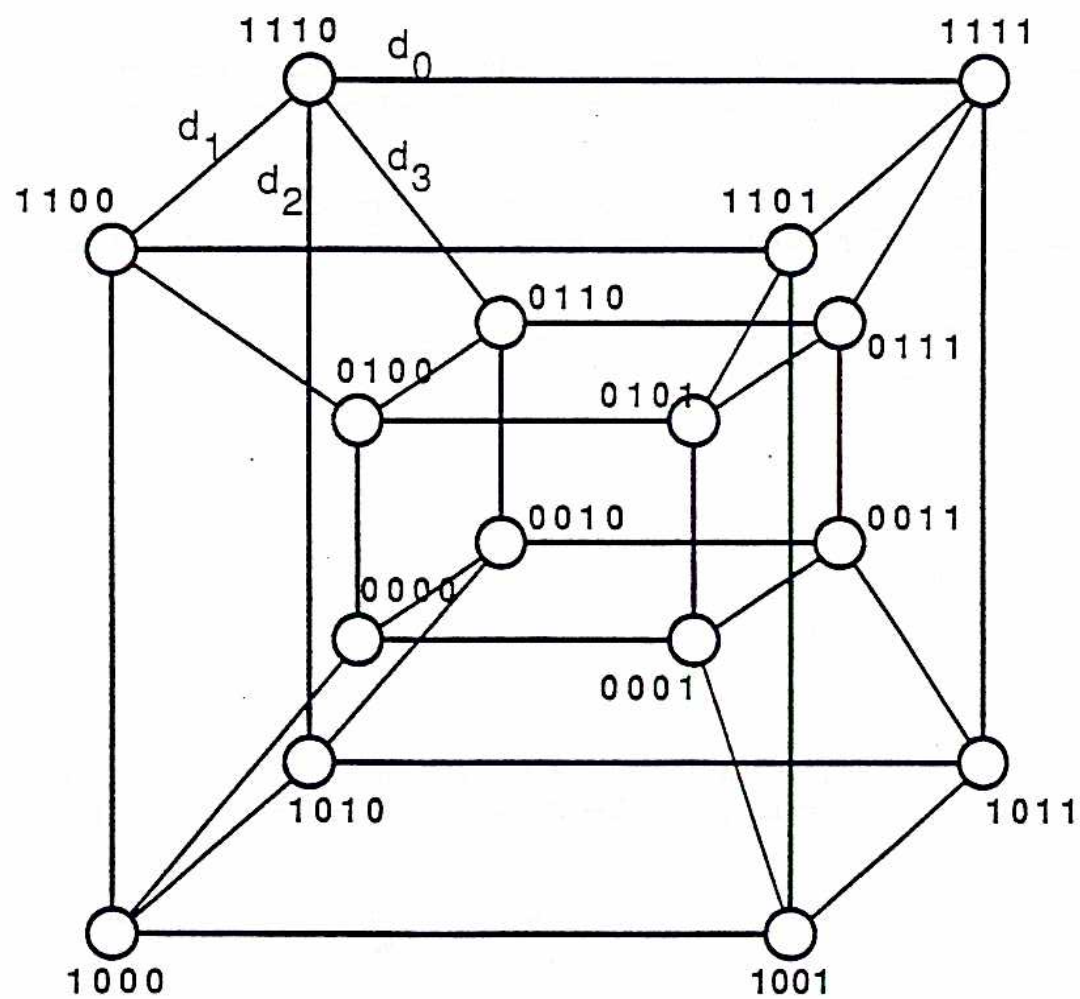


Figure 1. A 4-cube system.

| Order of preference | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---------------------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| node 0 | 1 | 2 | 4 | 8 | 6 | 10 | 12 | 3 | 5 | 9 | 14 | 13 | 11 | 7 | 15 |
| node 1 | 0 | 3 | 5 | 9 | 7 | 11 | 13 | 2 | 4 | 8 | 15 | 12 | 10 | 6 | 14 |
| node 2 | 3 | 0 | 6 | 10 | 4 | 8 | 14 | 1 | 7 | 11 | 12 | 15 | 9 | 5 | 13 |
| node 3 | 2 | 1 | 7 | 11 | 5 | 9 | 15 | 0 | 6 | 10 | 13 | 14 | 8 | 4 | 12 |
| node 4 | 5 | 6 | 0 | 12 | 2 | 14 | 8 | 7 | 1 | 13 | 10 | 9 | 15 | 3 | 11 |
| node 5 | 4 | 7 | 1 | 13 | 3 | 15 | 9 | 6 | 0 | 12 | 11 | 8 | 14 | 2 | 10 |
| node 6 | 7 | 4 | 2 | 14 | 0 | 12 | 10 | 5 | 3 | 15 | 8 | 11 | 13 | 1 | 9 |
| node 7 | 6 | 5 | 3 | 15 | 1 | 13 | 11 | 4 | 2 | 14 | 9 | 10 | 12 | 0 | 8 |
| node 8 | 9 | 10 | 12 | 0 | 14 | 2 | 4 | 11 | 13 | 1 | 6 | 5 | 3 | 15 | 7 |
| node 9 | 8 | 11 | 13 | 1 | 15 | 3 | 5 | 10 | 12 | 0 | 7 | 4 | 2 | 14 | 6 |
| node 10 | 11 | 8 | 14 | 2 | 12 | 0 | 6 | 9 | 15 | 3 | 4 | 7 | 1 | 13 | 5 |
| node 11 | 10 | 9 | 15 | 3 | 13 | 1 | 7 | 8 | 14 | 2 | 5 | 6 | 0 | 12 | 4 |
| node 12 | 13 | 14 | 8 | 4 | 10 | 6 | 0 | 15 | 9 | 5 | 2 | 1 | 7 | 11 | 3 |
| node 13 | 12 | 15 | 9 | 5 | 11 | 7 | 1 | 14 | 8 | 4 | 3 | 0 | 6 | 10 | 2 |
| node 14 | 15 | 12 | 10 | 6 | 8 | 4 | 2 | 13 | 11 | 7 | 0 | 3 | 5 | 9 | 1 |
| node 15 | 14 | 13 | 11 | 7 | 9 | 5 | 3 | 12 | 10 | 6 | 1 | 2 | 4 | 8 | 0 |

Figure 2. Preferred lists of a 4-cube system.

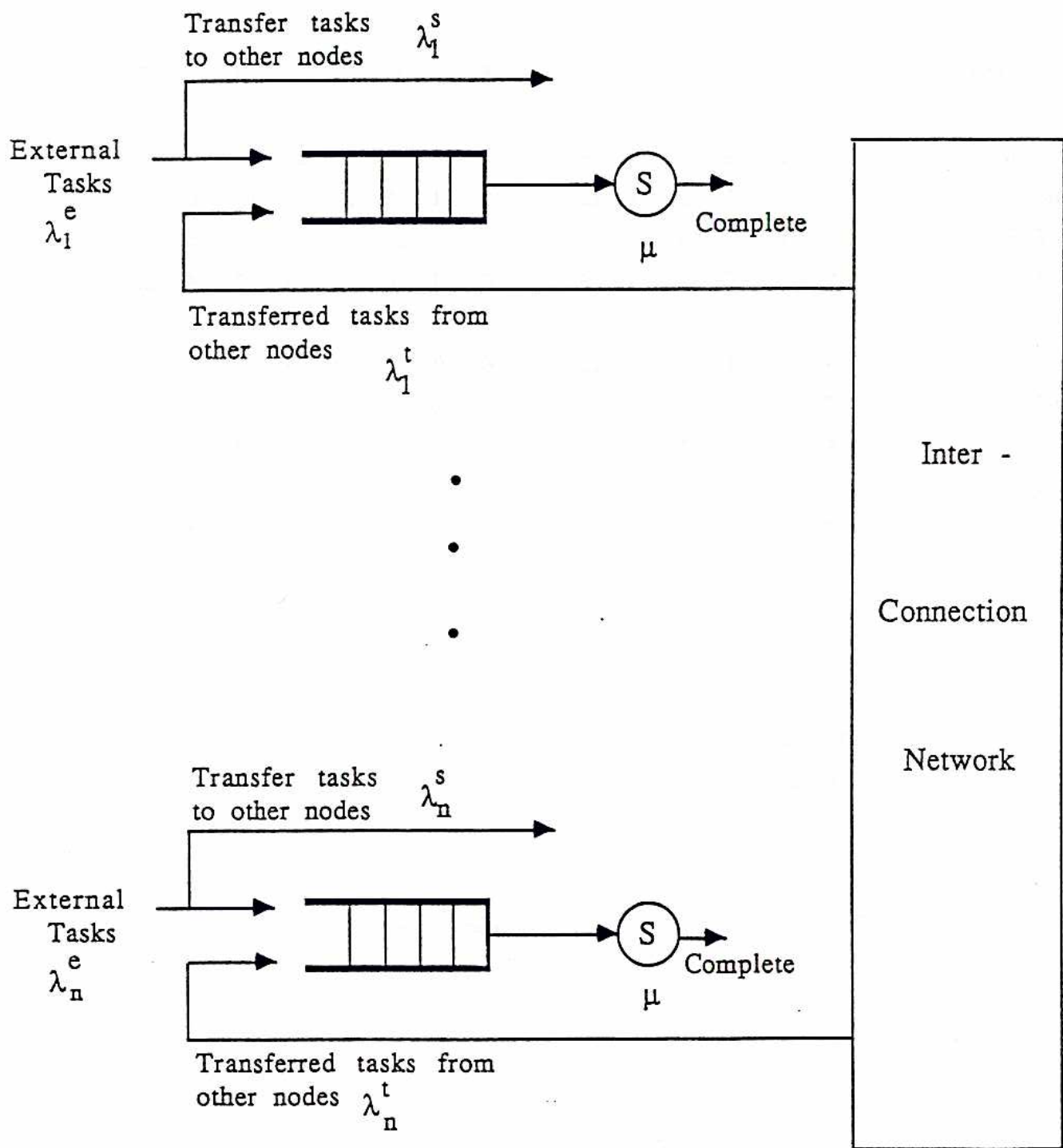


Figure 3. System model.

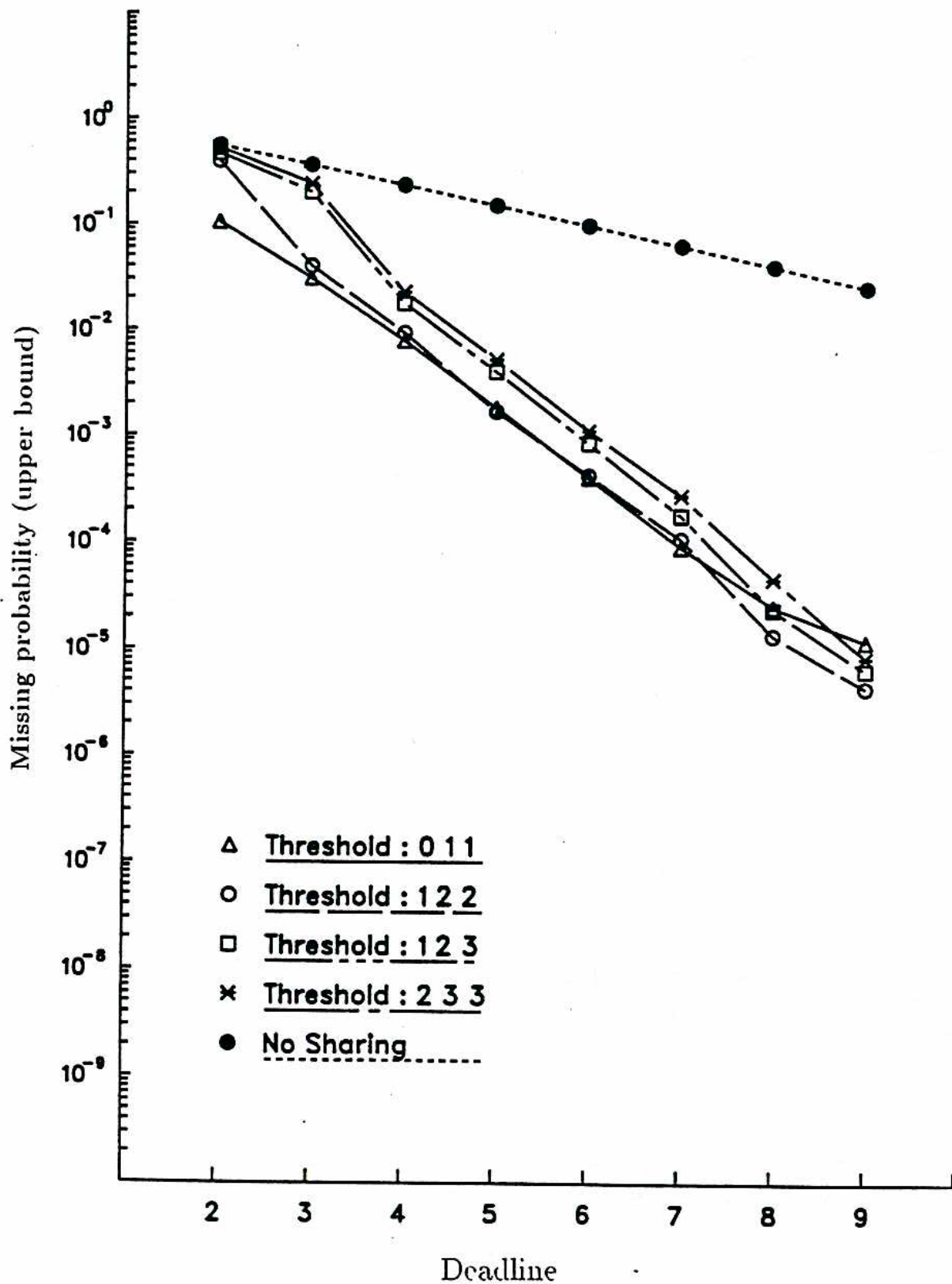


Figure 4. Upper bound missing probabilities vs. deadlines for different thresholds when $\lambda^e = 0.8$.

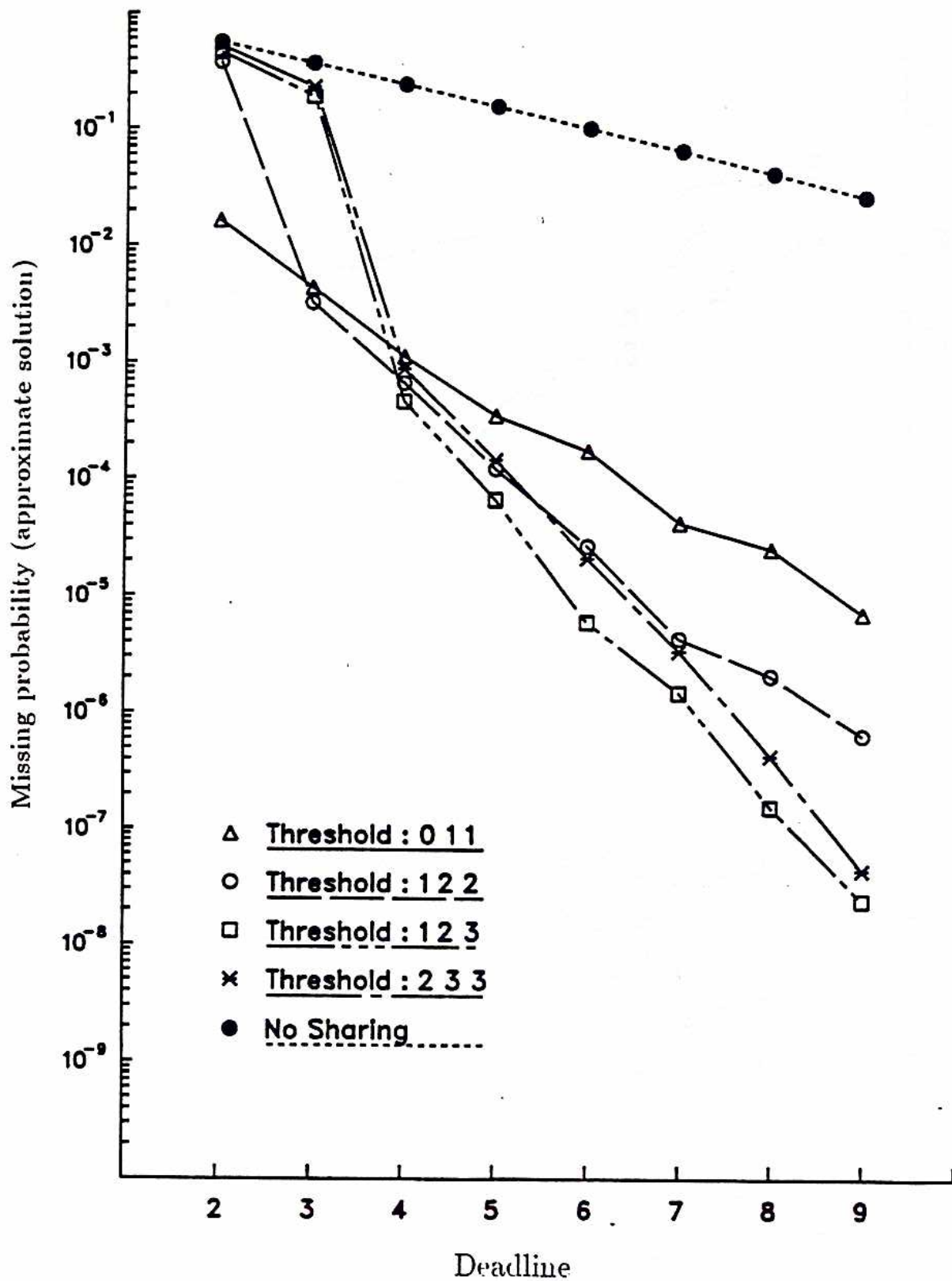


Figure 5. Approximate missing probabilities vs. deadlines for different thresholds when $\lambda^e = 0.8$.

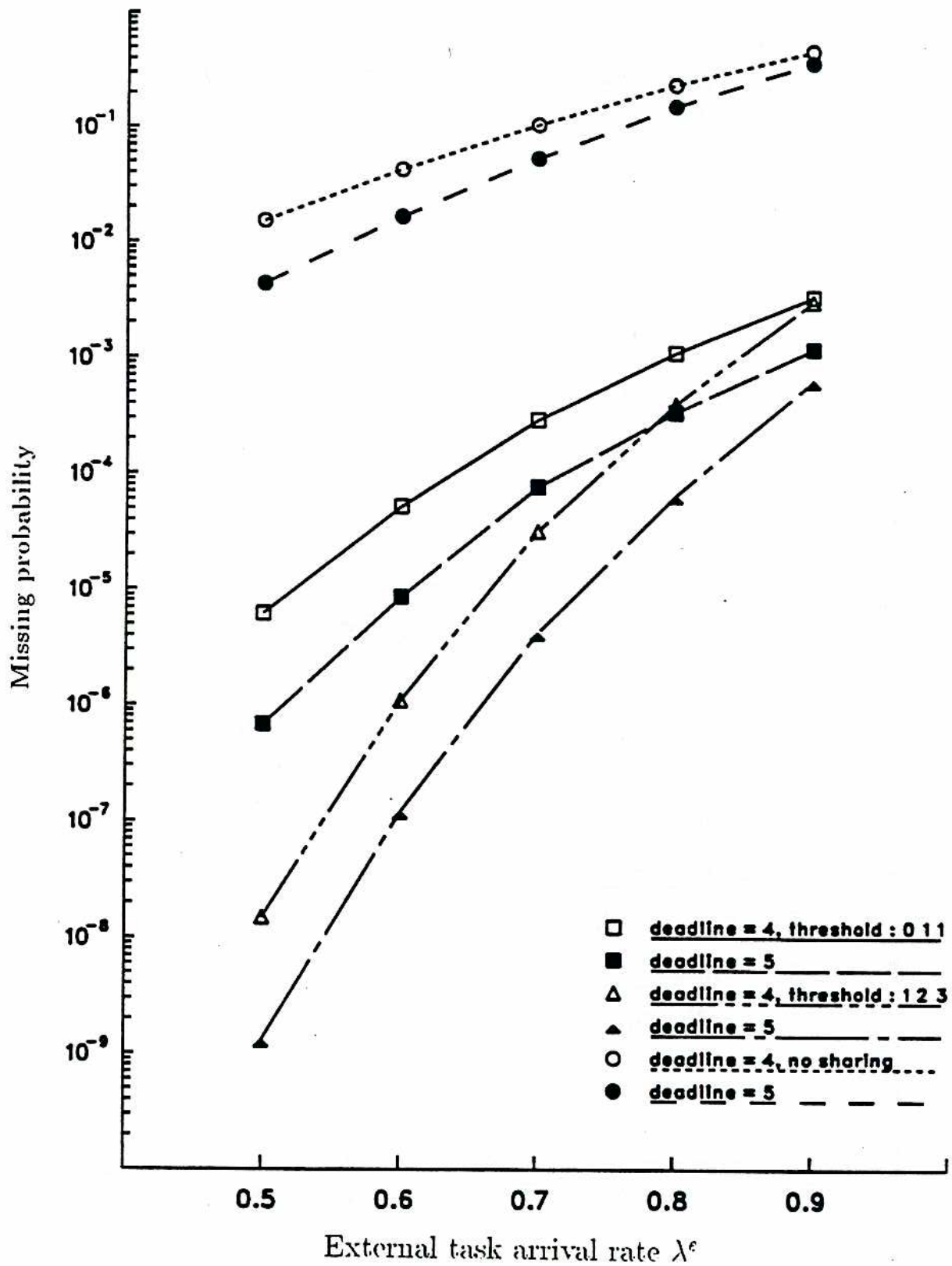


Figure 6. Approximate missing probabilities vs. load density for different thresholds and deadlines.

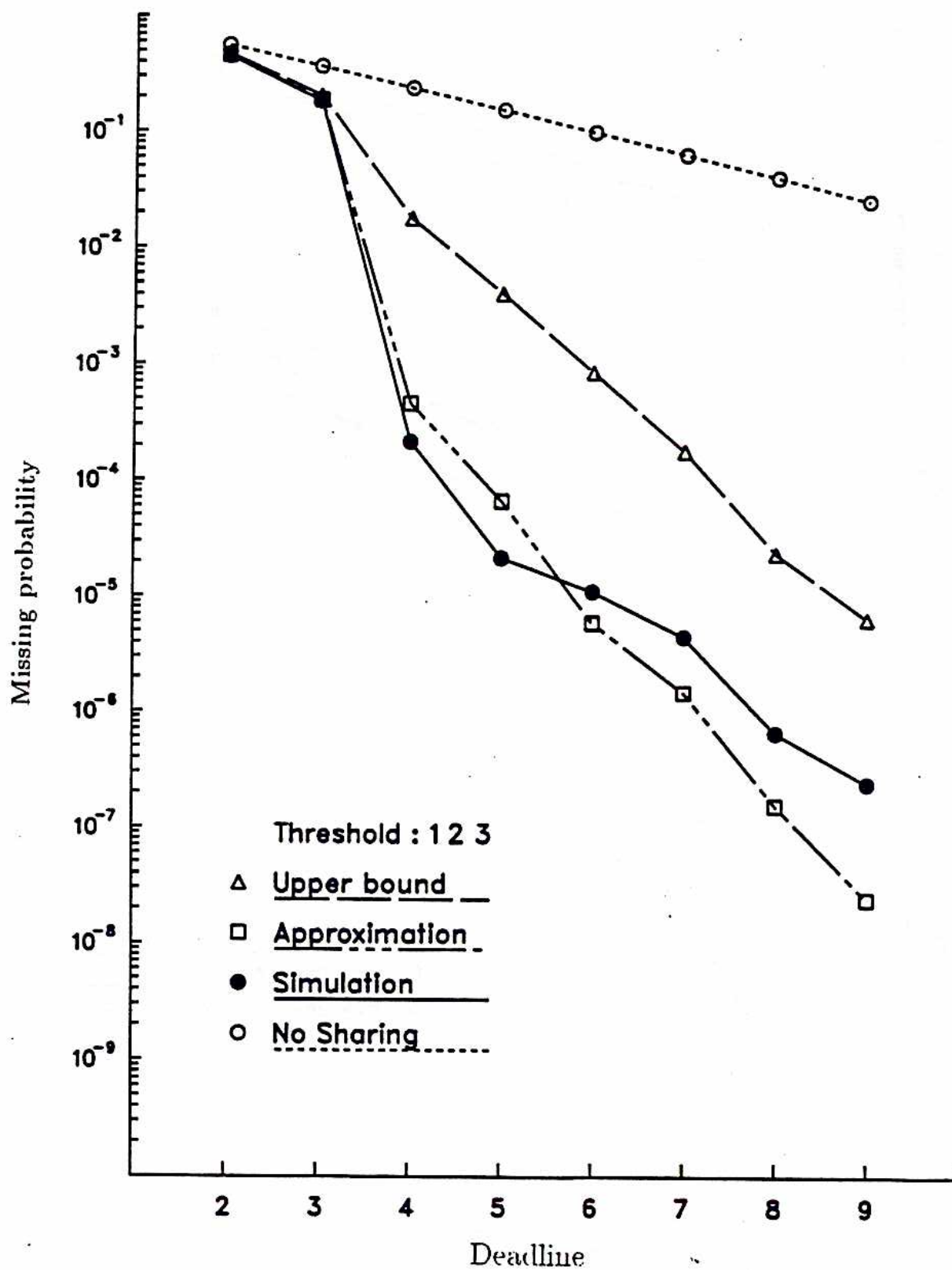


Figure 7. Simulated, approximate and upper bound missing probabilities vs. deadlines in when $\lambda^e = 0.8$.

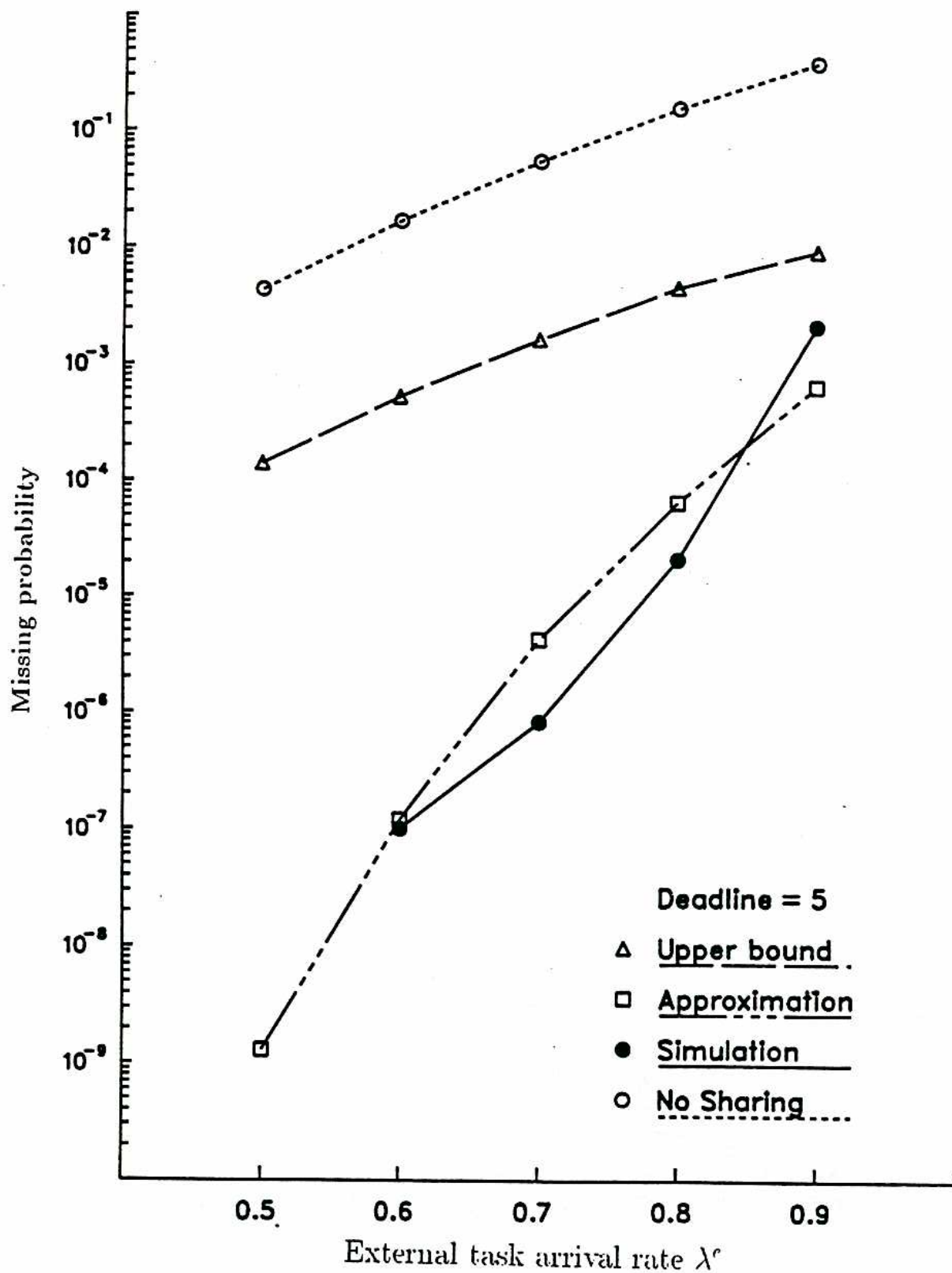


Figure 8. Approximate, upper bound, and simulated missing probabilities vs. load density with threshold pattern "1 2 3".

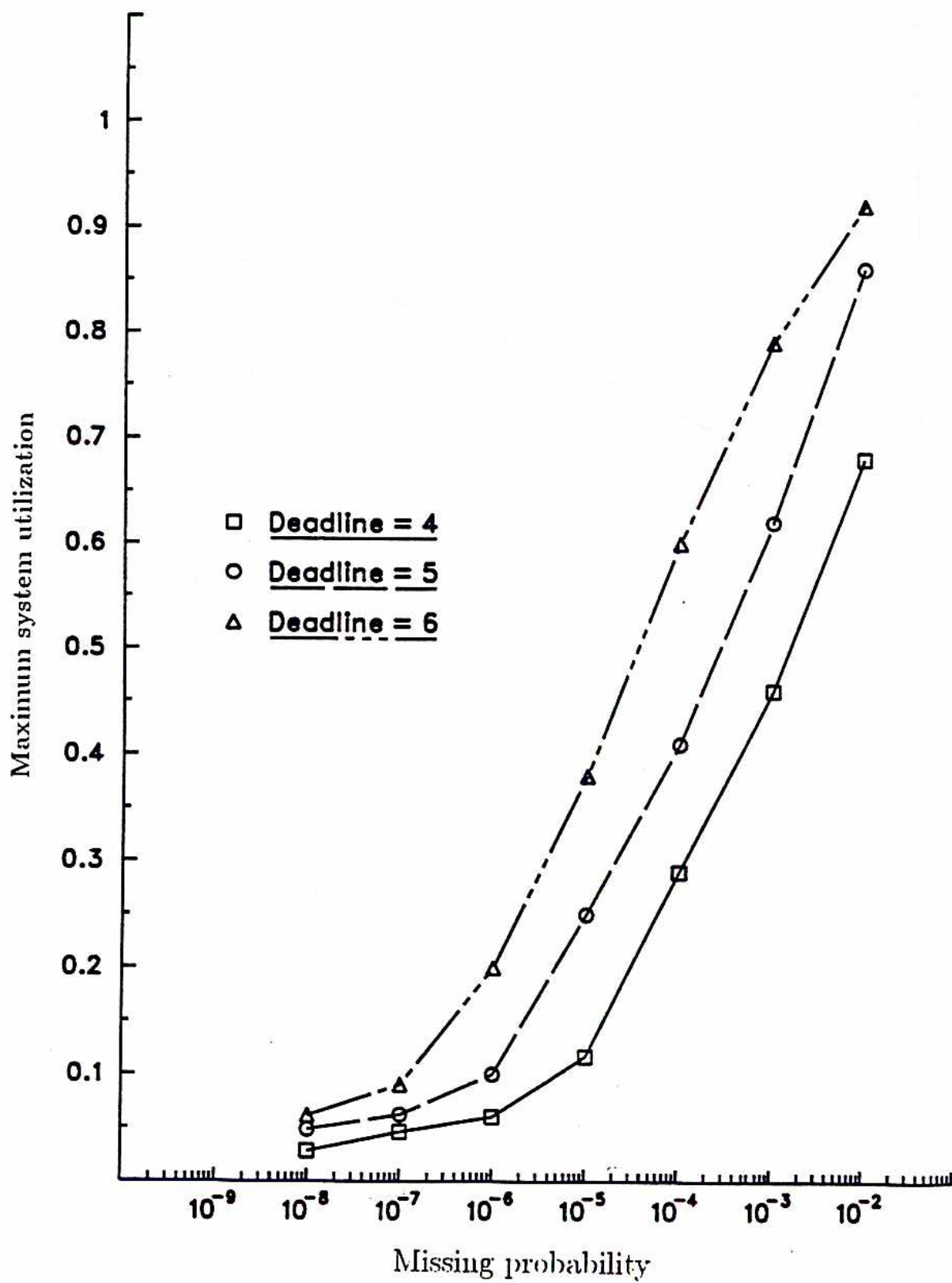


Figure 9. Maximum system utilization vs. missing probability for different deadlines.

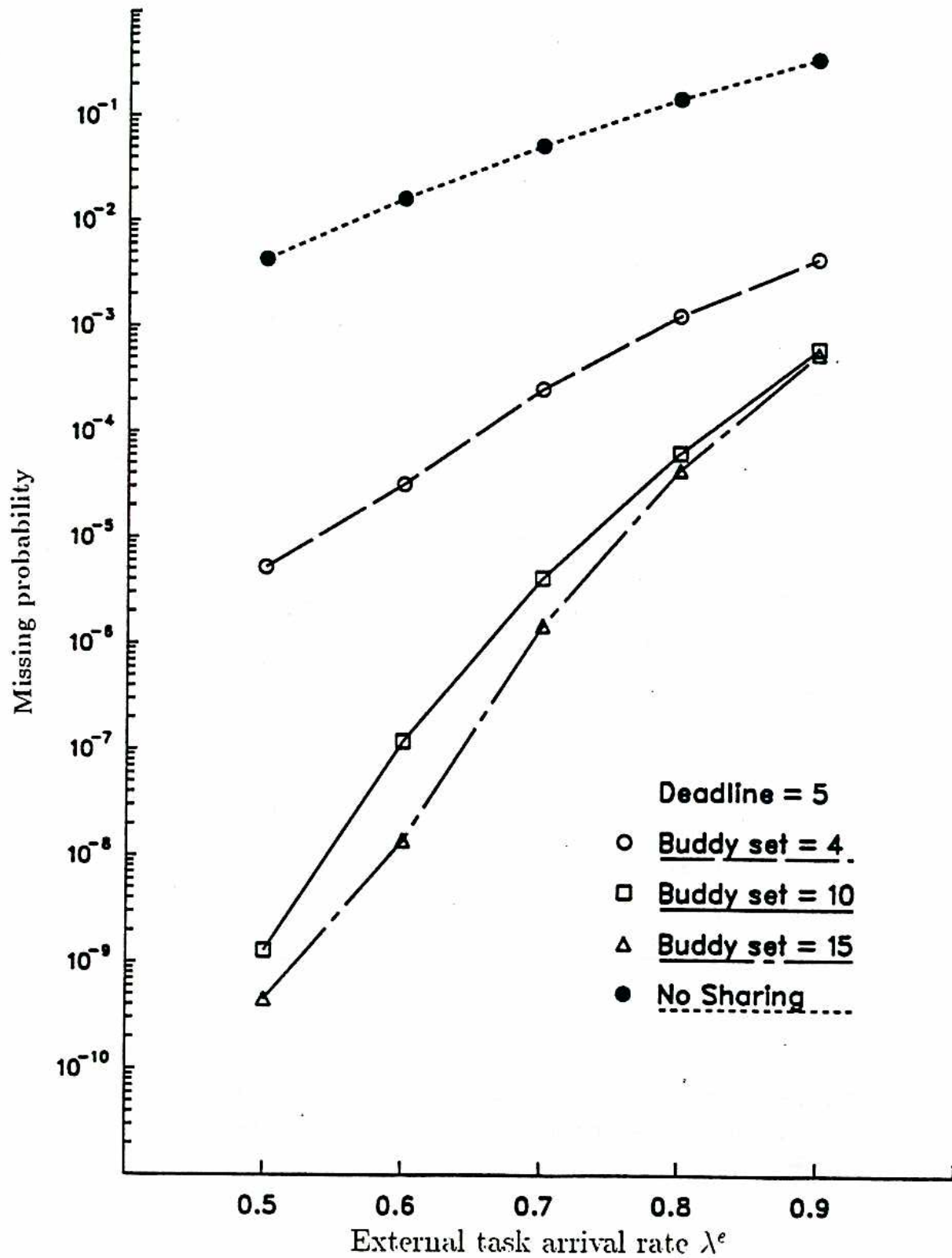


Figure 10. Approximate missing probabilities vs. load density with threshold "1 2 3" and different buddy set sizes.

| $(\lambda^e = 0.5)$ Queue Length \ Model | Simulation | Approximation | Upper bound | no load sharing |
|---|-----------------------|------------------------|-----------------------|-----------------------|
| 0 | 0.5037 | 0.4987 | 0.4992 | 0.5000 |
| 1 | 0.3316 | 0.3317 | 0.3321 | 0.3244 |
| 2 | 0.1254 | 0.1289 | 0.1278 | 0.1226 |
| 3 | 0.0387 | 0.0406 | 0.0398 | 0.0377 |
| 4 | 3.13×10^{-7} | 1.43×10^{-8} | 9.66×10^{-4} | 0.0109 |
| 5 | — | 1.21×10^{-9} | 1.24×10^{-4} | 0.0031 |
| 6 | — | 9.09×10^{-11} | 1.3×10^{-5} | 8.84×10^{-4} |
| 7 | — | 5.96×10^{-12} | 1.4×10^{-6} | 2.52×10^{-4} |
| 8 | — | 3.35×10^{-13} | 1.33×10^{-7} | 7.16×10^{-5} |
| 9 | — | 1.46×10^{-14} | 1.28×10^{-8} | 2.04×10^{-5} |

| $(\lambda^e = 0.8)$ Queue Length \ Model | Simulation | Approximation | Upper bound | no load sharing |
|---|-----------------------|-----------------------|-----------------------|-----------------|
| 0 | 0.2213 | 0.2264 | 0.2185 | 0.2004 |
| 1 | 0.3317 | 0.3194 | 0.3136 | 0.2456 |
| 2 | 0.2656 | 0.2675 | 0.2651 | 0.1898 |
| 3 | 0.1810 | 0.1862 | 0.1853 | 0.1278 |
| 4 | 1.69×10^{-4} | 3.34×10^{-4} | 0.0135 | 0.0834 |
| 5 | 2.01×10^{-5} | 5.94×10^{-5} | 0.0032 | 0.0542 |
| 6 | 8.98×10^{-6} | 8.99×10^{-6} | 0.0007 | 0.0353 |
| 7 | 3.21×10^{-6} | 1.18×10^{-6} | 0.0001 | 0.0229 |
| 8 | 9.76×10^{-7} | 1.36×10^{-7} | 1.36×10^{-6} | 0.0149 |
| 9 | 6.83×10^{-7} | 1.41×10^{-8} | 1.38×10^{-7} | 0.0097 |

Table 1. Comparison of distributions of queue length derived from the upper bound model, the approximate solution, and simulation with threshold pattern "1 2 3" and a buddy set of size 10.

| $(\lambda^e = 0.5)$ Threshold \ Buddy set | 4 | 10 | 15 |
|--|-------|-------|-------|
| 0 1 1 | 1.414 | 1.394 | 1.401 |
| 0 1 2 | 1.609 | 1.602 | 1.605 |
| 0 2 2 | 1.608 | 1.608 | 1.622 |
| 0 2 3 | 1.696 | 1.695 | 1.700 |
| 1 2 2 | 1.618 | 1.618 | 1.625 |
| 1 2 3 | 1.698 | 1.698 | 1.701 |
| 1 3 3 | 1.700 | 1.700 | 1.703 |
| 2 3 3 | 1.701 | 1.701 | 1.703 |
| No sharing | 1.750 | 1.750 | 1.750 |

| $(\lambda^e = 0.8)$ Threshold \ Buddy set | 4 | 10 | 15 |
|--|-------|-------|-------|
| 0 1 1 | 1.918 | 1.692 | 1.651 |
| 0 1 2 | 2.236 | 2.065 | 1.033 |
| 0 2 2 | 2.120 | 2.075 | 2.071 |
| 0 2 3 | 2.426 | 2.382 | 2.380 |
| 1 2 2 | 2.123 | 2.109 | 2.110 |
| 1 2 3 | 2.423 | 2.407 | 2.405 |
| 1 3 3 | 2.427 | 2.422 | 2.421 |
| 2 3 3 | 2.475 | 2.472 | 2.473 |
| No sharing | 3.370 | 3.370 | 3.370 |

Table 2. Average task system time when $\lambda^e = 0.5$ and $\lambda^e = 0.8$.

| $(\lambda^e = 0.5)$ | Threshold | 0 1 1 | 0 1 2 | 1 2 2 | 1 2 3 | 2 3 3 | No Sharing |
|-----------------------------|-----------|-----------------------|-----------------------|------------------------|------------------------|------------------------|------------|
| | Deadline | | | | | | |
| simulation | 2 | 0.0007 | 0.1296 | 0.1339 | 0.1642 | 0.1661 | — |
| | 3 | 4.81×10^{-5} | 8.82×10^{-5} | 1.73×10^{-5} | 0.0388 | 0.0392 | — |
| | 4 | 3.10×10^{-6} | 5.32×10^{-6} | 1.45×10^{-7} | 3.13×10^{-7} | 6.59×10^{-7} | — |
| | 5 | 2.85×10^{-7} | $< 10^{-7}$ | $< 10^{-7}$ | $< 10^{-7}$ | $< 10^{-7}$ | — |
| | 6 | $< 10^{-7}$ | $< 10^{-7}$ | $< 10^{-7}$ | $< 10^{-7}$ | $< 10^{-7}$ | — |
| analytic (approximation) | 2 | 0.0004 | 0.1370 | 0.1446 | 0.1681 | 0.1691 | 0.1756 |
| | 3 | 5.10×10^{-5} | 8.41×10^{-5} | 7.46×10^{-7} | 0.0398 | 0.0406 | 0.0530 |
| | 4 | 6.25×10^{-6} | 9.73×10^{-6} | 7.11×10^{-8} | 1.46×10^{-8} | 1.87×10^{-9} | 0.0152 |
| | 5 | 6.81×10^{-7} | 9.51×10^{-7} | 6.01×10^{-9} | 1.31×10^{-9} | 1.44×10^{-10} | 0.0043 |
| | 6 | 7.61×10^{-8} | 8.8×10^{-8} | 9.11×10^{-11} | 1.31×10^{-11} | 1.01×10^{-11} | 0.0012 |

| $(\lambda^e = 0.8)$ | Threshold | 0 1 1 | 0 1 2 | 0 2 3 | 1 2 2 | 1 2 3 | 2 3 3 | No Sharing |
|-----------------------------|-----------|--------|--------|-----------------------|-----------------------|-----------------------|-----------------------|------------|
| | Deadline | | | | | | | |
| simulation | 2 | 0.0745 | 0.3258 | 0.4414 | 0.3461 | 0.4468 | 0.4858 | — |
| | 3 | 0.0176 | 0.0305 | 0.1805 | 0.0019 | 0.1812 | 0.2007 | — |
| | 4 | 0.0045 | 0.0067 | 0.012 | 0.0003 | 0.0002 | 0.0002 | — |
| | 5 | 0.0014 | 0.0015 | 0.0003 | 8.37×10^{-5} | 2.11×10^{-5} | 2.28×10^{-5} | — |
| | 6 | 0.0006 | 0.0004 | 0.0001 | 3.61×10^{-5} | 1.11×10^{-5} | 6.18×10^{-6} | — |
| analytic (approximation) | 2 | 0.0161 | 0.3315 | 0.4342 | 0.3723 | 0.4541 | 0.5151 | 0.5539 |
| | 3 | 0.0043 | 0.0105 | 0.1772 | 0.0032 | 0.1866 | 0.2256 | 0.3641 |
| | 4 | 0.0011 | 0.0024 | 0.0013 | 0.0007 | 0.0004 | 0.0009 | 0.2363 |
| | 5 | 0.0004 | 0.0005 | 0.0003 | 0.0001 | 6.54×10^{-5} | 0.0002 | 0.1528 |
| | 6 | 0.0002 | 0.0002 | 3.95×10^{-5} | 2.64×10^{-5} | 1.2×10^{-5} | 2.58×10^{-5} | 0.0986 |

Table 3. Missing probabilities for various threshold patterns and task deadlines when $\lambda^e = 0.5$ and $\lambda^e = 0.8$, and buddy set size = 10.

| $(\lambda^e = 0.8)$ | Threshold | | | | | |
|---------------------|-----------|----------------------|----------------------|----------------------|----------------------|----------------------|
| Buddy set | Deadline | 0 1 2 | 1 2 2 | 1 2 3 | 1 3 3 | 2 3 3 |
| 4 | 2 | 0.3702 | 0.3491 | 0.4506 | 0.4555 | 0.4873 |
| | 3 | 0.0950 | 0.0071 | 0.1856 | 0.1885 | 0.2018 |
| | 4 | 0.0320 | 0.0018 | 0.0047 | 0.0041 | 0.0011 |
| | 5 | 0.0107 | 0.0004 | 0.0008 | 0.0006 | 0.0002 |
| | 6 | 0.0034 | 9.1×10^{-5} | 0.0002 | 0.0002 | 5.8×10^{-5} |
| 10 | 2 | 0.3181 | 0.3448 | 0.4463 | 0.4526 | 0.4856 |
| | 3 | 0.0229 | 0.0012 | 0.1806 | 0.1860 | 0.2008 |
| | 4 | 0.0043 | 8.9×10^{-5} | 6.4×10^{-5} | 0.0004 | 0.0002 |
| | 5 | 0.0008 | 8.1×10^{-6} | 5.2×10^{-6} | 1.6×10^{-5} | 6.6×10^{-6} |
| | 6 | 0.0001 | 6.4×10^{-7} | 4.7×10^{-7} | 2.9×10^{-6} | 5.8×10^{-7} |
| 15 | 2 | 0.3073 | 0.3448 | 0.4463 | 0.4532 | 0.4856 |
| | 3 | 0.0089 | 0.0012 | 0.1806 | 0.1863 | 0.2006 |
| | 4 | 0.0013 | 5.7×10^{-5} | 3.9×10^{-5} | 0.0004 | 0.0002 |
| | 5 | 0.0002 | 5.1×10^{-6} | 3.2×10^{-6} | 1.1×10^{-5} | 7.6×10^{-6} |
| | 6 | 2.0×10^{-5} | 2.4×10^{-7} | 1.6×10^{-7} | 8.4×10^{-7} | 2.6×10^{-7} |
| 21 | 2 | 0.3026 | 0.3445 | 0.4463 | 0.4526 | 0.4856 |
| | 3 | 0.0033 | 0.0012 | 0.1806 | 0.1860 | 0.2008 |
| | 4 | 0.0004 | 6.0×10^{-5} | 3.7×10^{-5} | 0.0004 | 0.0002 |
| | 5 | 9.2×10^{-5} | 4.1×10^{-6} | 2.2×10^{-6} | 1.1×10^{-5} | 4.6×10^{-6} |
| | 6 | 1.1×10^{-5} | 1.7×10^{-7} | $< 10^{-7}$ | 7.4×10^{-7} | 2.9×10^{-7} |

Table 4. Missing probabilities in a 6-cube for different thresholds and buddy set sizes.

| $(\lambda^e = 0.5)$ Threshold | Task Transfers | | Frequency of State Change | |
|----------------------------------|----------------|-------------|---------------------------|-------------|
| | Simulation | Approximate | Simulation | Approximate |
| 0 1 1 | 0.1071 | 0.1300 | 1.0792 | 1.1128 |
| 0 1 2 | 0.0304 | 0.0330 | 1.0203 | 1.0266 |
| 0 2 2 | 0.0484 | 0.0501 | 0.1327 | 0.1471 |
| 1 2 2 | 0.0332 | 0.0392 | 0.1574 | 0.1754 |
| 1 2 3 | 0.0092 | 0.0098 | 0.1511 | 0.1552 |
| 1 3 3 | 0.0097 | 0.0101 | 0.0407 | 0.0370 |

| $(\lambda^e = 0.8)$ Threshold | Task Transfers | | Frequency of State Change | |
|----------------------------------|----------------|-------------|---------------------------|-------------|
| | Simulation | Approximate | Simulation | Approximate |
| 0 1 1 | 0.2241 | 0.2230 | 0.6277 | 0.8118 |
| 0 1 2 | 0.1145 | 0.1600 | 0.5283 | 0.5984 |
| 0 2 2 | 0.2274 | 0.2015 | 0.1613 | 0.2531 |
| 1 2 2 | 0.1677 | 0.1891 | 0.2576 | 0.3316 |
| 1 2 3 | 0.0877 | 0.0811 | 0.2182 | 0.2404 |
| 1 3 3 | 0.1064 | 0.0934 | 0.1050 | 0.1204 |

Table 5. Number of task transfers vs. frequency of state change for different thresholds.

| $(\lambda^e = 0.8)$ | Threshold | 0 1 2 | 1 2 2 | 1 2 3 | 2 3 3 |
|---------------------|-----------|-----------------------|-----------------------|-----------------------|-----------------------|
| Transfer Cost | Deadline | | | | |
| 5% | 2 | 0.3227 | 0.3460 | 0.4473 | 0.4860 |
| | 3 | 0.0307 | 0.0013 | 0.1812 | 0.2007 |
| | 4 | 0.0067 | 1.86×10^{-4} | 0.0047 | 1.04×10^{-4} |
| | 5 | 0.0015 | 3.51×10^{-5} | 4.26×10^{-5} | 6.65×10^{-6} |
| | 6 | 3.51×10^{-4} | 7.43×10^{-6} | 9.11×10^{-6} | 5.98×10^{-7} |
| 10% | 2 | 0.3258 | 0.3461 | 0.4468 | 0.4858 |
| | 3 | 0.0305 | 0.0019 | 0.1812 | 0.2007 |
| | 4 | 0.0067 | 2.8×10^{-4} | 2.13×10^{-4} | 2.07×10^{-4} |
| | 5 | 0.0015 | 8.37×10^{-5} | 2.11×10^{-5} | 2.28×10^{-5} |
| | 6 | 0.0015 | 8.37×10^{-5} | 2.11×10^{-5} | 2.28×10^{-5} |
| 20% | 2 | 0.3300 | 0.3530 | 0.4504 | 0.4847 |
| | 3 | 0.0295 | 0.0046 | 0.1918 | 0.2017 |
| | 4 | 0.0065 | 5.55×10^{-4} | 0.0012 | 7.25×10^{-4} |
| | 5 | 0.0015 | 8.69×10^{-5} | 1.2×10^{-4} | 5.31×10^{-5} |
| | 6 | 3.50×10^{-4} | 2.0×10^{-5} | 1.36×10^{-5} | 6.03×10^{-6} |
| 30% | 2 | 0.3488 | 0.3817 | 0.4681 | 0.4914 |
| | 3 | 0.0320 | 0.0168 | 0.1956 | 0.2177 |
| | 4 | 0.0072 | 0.0023 | 0.0016 | 0.0050 |
| | 5 | 0.0016 | 3.2×10^{-4} | 1.93×10^{-5} | 5.29×10^{-4} |
| | 6 | 3.39×10^{-4} | 4.34×10^{-5} | 2.93×10^{-5} | 4.95×10^{-5} |

Table 6. Missing probabilities vs. task transfer costs for different deadlines and thresholds.

