



Stochastics and Statistics

Modelling deadlock in open restricted queueing networks



Geraint I. Palmer, Paul R. Harper*, Vincent A. Knight

School of Mathematics, Cardiff University, Senghennydd Road, Cardiff CF24 4AG, United Kingdom

ARTICLE INFO

Article history:

Received 15 June 2016

Accepted 16 October 2017

Available online 7 November 2017

Keywords:

Queueing

Queueing networks

Deadlock

Markov models

ABSTRACT

Open restricted queueing networks give rise to the phenomenon of deadlock, whereby some customers may be unable to ever leave a server due to mutual blocking. This paper explores deadlock in queueing networks with limited queueing capacity, presents a method of detecting deadlock in discrete event simulations, and builds Markov chain models of these deadlocking networks. The three networks for which Markov models are given include single and multi-server networks for one and two node systems. The expected times to deadlock of these models are compared to results obtained using a simulation of the stochastic process, together with the developed deadlock detection method. This paper aims to be of value to simulation modellers of queues.

© 2017 The Author(s). Published by Elsevier B.V.

This is an open access article under the CC BY license. (<http://creativecommons.org/licenses/by/4.0/>)

1. Introduction

The study and modelling of queueing networks with blocking is an important tool in many aspects of operational research, both analytically and through simulation. These models have applications in many varied settings such as healthcare, supply chains, manufacturing and communications systems. However, these types of models have their limitations, due to their potential to become permanently blocked in deadlock, or a deadly embrace of resources. These deadlocks can be real and observed in reality, in which case accurate modelling of deadlock is needed; or they can be a symptom of a model unable to capture certain behaviours. This may occur in models where deadlock situations are easily adjusted in reality. In this case, such as by swapping two customers, a good understanding of deadlock is needed in order to model the adjusted reality.

Queueing networks are described as open if customers can enter and leave the system from the exterior. Restricted networks are those where at least one service centre has limited queueing space or capacity before it. Deadlock is caused by blocking. This paper considers Type I blocking: after service a customer will be blocked from joining a queue at another node if that node's queueing capacity is full. While blocked, that customer remains with its server until space becomes available at its destination. During this time that server is unavailable to begin another customer's service.

For the purposes of this paper, deadlock is defined as follows.

Definition 1. When there is a subset of blocked customers who are blocked directly or indirectly by customers in that subset only, then the system is said to be in deadlock.

This implies that a system is in deadlock when at least one service station permanently ceases to begin or finish any more services, due to circular blocking. Fig. 1 shows an open two node restricted queueing network in deadlock. The customer at the top server is blocked from entering the bottom node as there is a full queue, and similarly the customer at the bottom server is blocked from entering the top node as there is a full queue. It is clear that by following the rules of blocking defined above, no more natural movement can happen. This system is in deadlock as there is a subset of blocked customers (the customer with server A_1 and the customer with server B_1) who are only being blocked by each other.

This paper is concerned with open restricted queueing networks that experience Type I blocking. Exponential service times and Poisson arrivals are assumed. First in first out, or FIFO service discipline is also assumed. Throughout the paper service centres will be referred to as nodes, and for the i th node of a queueing network the following notation is used:

- Λ_i denotes the external arrival rate.
- μ_i denotes the service rate.
- c_i denotes the number of parallel servers.
- n_i denotes the queueing capacity.
- r_{ij} denotes the routing probability from node i to node j upon completion of service at node i .

The main contribution of this work to the literature is a formal, rigorous, and analytical study of deadlock in queueing processes, which has never been done before. Methodologies for the

* Corresponding author.

E-mail address: harper@cardiff.ac.uk (P.R. Harper).

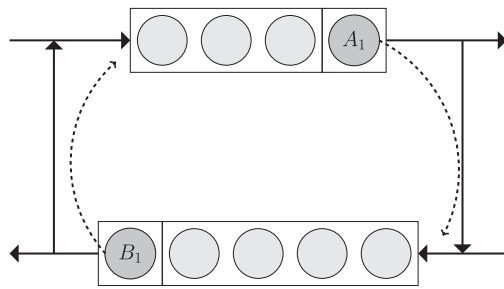


Fig. 1. Example of an open two node restricted queueing network in deadlock.

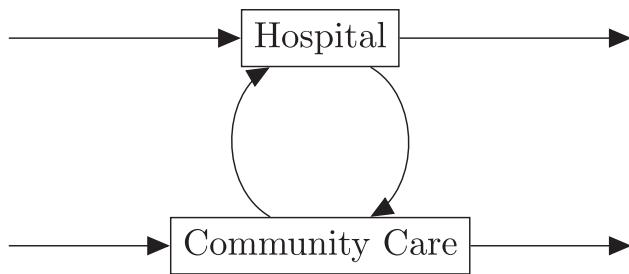


Fig. 2. Diagram of patient flows at an interface between secondary care services at a hospital and community care services.

detection of deadlock as well as Markovian models of deadlock are presented. In particular, this paper looks at detecting when deadlock occurs, and the time until a deadlock occurs from an empty system. First, a method for detecting deadlock in simulations of queueing networks is presented. Then, Markov models of simple deadlocking queueing networks are built. This not only contributes a novel theoretic advancement of the study of queueing networks, but also has the potential for impact in real world queueing networks, as discussed in Section 2.

The remainder of this paper is structured as follows: Section 2 gives a motivating example to put the work in context. Section 3 gives an overview of existing literature on the subject. Section 4 presents a method of detecting deadlock in simulations of queueing networks. Section 5 presents Markov models of three deadlocking queueing networks, finds their expected time to deadlock, and compares these with results obtained through simulation models.

2. A motivating example

Here we present a motivating example of a healthcare system. In this example deadlock may be easily resolved in reality, however analytical stochastic models and simulations may be restricted by deadlock. Therefore an understanding of this phenomenon, and an ability to overcome this effect in discrete event simulations, is essential for modelling this system.

Consider the interface between secondary care services at a hospital and community care services. Patients can be admitted to hospital via a variety of routes (e.g. through emergency services, or outpatients), and via referral from community care services. Patients can begin receiving community care packages after referral from GP, or via referral from the hospital. Considering only the hospital and community care services as nodes, this system is shown in Fig. 2.

If there are no free hospital beds, then patients being referred from community care services will be sustained by community care workers until beds become available. If there are no community care packages available, then patients requiring packages but unfit to return home after a hospital stay will remain in hospital, blocking beds until a community care package becomes avail-

able. Type I blocking occurs here, as patients and staff do not know the future capacity of their next destination prior to service. This type of bed blocking is well known (Manzano-Santaella, 2010). This causes problems for patients as they are being cared for in an inappropriate setting for their condition, and also for the health care providers as secondary care may be more expensive than primary care, and resolution of this causes administrative stress.

In this model there is a non-zero probability of everyone at the hospital blocking beds waiting for community care packages, and everyone at community care being sustained waiting for beds at the hospital. Thus the model will exhibit deadlock. In reality, there is communication between these services and patients can swap places. This ensures no deadlock.

Restricted feedback loops that exhibit mutual blocking such as this one have been observed in real healthcare systems, as described in a case study in Osorio and Bierlaire (2009). However the authors here state that this type of blocking “may be irrelevant in practice given that the swapping of patients can be identified and carried out easily”. In Koizumi, Kuno, and Smith (2005) a health and community care system is described as having restricted feedback loops. However due to ease of modelling, and to avoid the restrictions caused by deadlock, these feedback loops are omitted from the model. This emphasises the discrepancies that occur between common modelling techniques and reality in systems that may reach deadlock.

An understanding of how deadlock behaves in these models will aid the modelling process. A deadlock detection method for the simulation model will be invaluable in modelling realistic deadlock resolution methods, thus ensuring correct models can be built of systems like this with circular blocking.

3. Literature review

Restricted queueing networks that exhibit blocking are well discussed in the literature, both exact (Avi-Itzhak & Yadin, 1965; Baber, 2008; Gordon & Newell, 1967; Hunt, 1956; Koizumi et al., 2005; Latouche & Neuts, 1980; Perros, Nilsson, & Liu, 1988) and approximate methods (Allon, Deo, & Lin, 2013; Dallery & Frein, 1993; Korporaal, Ridder, Klopprogge, & Dekker, 2000; Onvural, 1990; Osorio & Bierlaire, 2009; Perros et al., 1988; Takahashi, Miyahara, & Hasegawa, 1980). Discussions on restricted queueing networks with feedback loops, that may exhibit deadlock, are sparse however. In fact the problem of deadlock in queueing networks has either been ignored, not studied, or assumed resolved in much of the literature (Onvural, 1990; Osorio & Bierlaire, 2009; Perros et al., 1988).

Central to the study of deadlock in queueing networks is the concept of blocking. In Onvural and Perros (1986) three types of blocking are described. Type I blocking occurs when a customer is blocked after completing service, and remains with the server until capacity at their destination node becomes available. Type II blocking occurs when a customer declares their destination before beginning service, and is only granted service if there is available capacity at their destination node. In Type III blocking, instead of getting blocked, a customer is required to repeat their service if there is no capacity at their destination. This type of blocking comes in two forms, fixed destination where the customer's destination does not change at each repetition of service, and random destination, where the customer's destination is re-sampled from a probability distribution after each repetition.

There has been a body of research around deadlock which does not consider the underlying stochastic structure of the system (Coffman & Elphick, 1971; Reveliotis, 2015a; 2015b). This type of deadlock, also referred to as deadly embraces (Coffman & Elphick, 1971), can potentially occur under the following conditions:

- Mutual exclusion: Tasks have exclusive control over resources.
- Wait for: Tasks do not release resources while waiting for other resources.
- No pre-emption: Resources cannot be removed until they have been used to completion.
- Circular wait: A circular chain of tasks exists, where each task requests a resource from another task in the chain.

In open restricted queueing networks the mutual exclusion condition is satisfied as customers cannot share servers; the wait for condition is satisfied due to the rules of Type I blocking; the no pre-emption condition is satisfied in networks that have no or non-pre-emptive priority (this paper only considers networks with no priority); and the circular wait condition is satisfied if the queueing network contains a cycle where all nodes have limited queueing capacity.

Allowing a system to reach deadlock can be problematic in cases where automated systems cannot continue operations, or where simulations cannot accurately model reality. In general there are three strategies for dealing with the problem of deadlock (Elmagarmid, 1986; Kawadkar, Prasad, & Dwivedi, 2014; Venkatesh & Smith, 2005; Vis, 2006):

- Avoidance, in which decisions are made as time unfolds to avoid reaching deadlock.
- Prevention, in which the system is designed such that it cannot possibly deadlock.
- Detection and recovery.

Note that (Holt, 1972) lists the three strategies as prevention, detection and crashing, which is equivalent to having no deadlock strategy. Allowing the system to crash now and again may be more economical in some systems where deadlocks do not occur often enough to justify the investment and effort of implementing an avoidance/resolution strategy.

Prevention and avoidance strategies have been used extensively in an area known as Discrete Event Systems (Reveliotis, 2015a; 2015b). A number of techniques and methods have been used to implement deadlock avoidance (Belik, 1990; Dijkstra, 1982; Ezpeleta, Tricas, Garca-Valls, & Colom, 2002; Kawadkar et al., 2014; Marchetti & Munier-Kordon, 2009; Vis, 2006; Viswanadham, Nara-hari, & Johnson, 1990). These techniques generally determine when resources cannot be allocated as that allocation would lead to deadlock. In Florian, Mahut, and Tremblay (2008) a priority based deadlock avoidance algorithm is implemented in a traffic simulation model. The purpose of the avoidance scheme here is not to reflect deadlock avoidance in reality, but to avoid deadlocks that will occur in the simulation due to missing information or incomplete models.

The literature has discussed deadlock prevention in closed queueing networks under Type I blocking (Kundu & Akyildiz, 1989; Liebeherr & Akyildiz, 1995; Onvural, 1990; Schmidt & Jackman, 2000). These have involved determining the minimum queueing space assignment that prevents deadlock for a given population size, or turning customers away if certain nodes are full. For simulation modelling however, prevention and avoidance techniques may not be appropriate as they can potentially inhibit realism in the simulation by taking actions that do not occur in the system being modelled (Venkatesh, Smith, Deuermeyer, & Curry, 1998).

A popular method of detecting general deadlock is the use of wait-for graphs, state-graphs and other variants (Cheng, 1990; Cho, Kumaran, & Wysk, 1995; Coffman & Elphick, 1971; Deuermeyer, Curry, Duchowski, & Venkatesh, 1997; Elmagarmid, 1986; Holt, 1972; Venkatesh & Smith, 2003; 2005; Venkatesh et al., 1998). These wait-for graphs, keep track of all circular wait relations between tasks. In Coffman and Elphick (1971) dynamic state-graphs are defined with resources as vertices and requests as edges. For

scenarios where there is only one type of each resource, deadlock arises if and only if the state-graph contains a cycle. In Cho et al. (1995) 'simple bounded circuits' are defined by giving the vertices and edges of the state graph labels in relation to a reference node. The existence of these circuits within the state graph indicates if the system is in deadlock. A strategy of this type is developed in this paper to detect deadlock in general queueing systems.

Bipartite entity-resource graphs are used in Holt (1972), Deuermeyer et al. (1997) and Venkatesh and Smith (2003) to detect deadlock in systems with both consumable and reusable resources. Two different types of deadlock are detected, transient deadlock and permanent deadlock. A deadlock resolution procedure is proposed that attempts to break cycles in the entity-resource graph. This work is furthered in Venkatesh et al. (1998) where deadlock is detected and resolved for situations where entities may request more than one resource.

Deadlock detection and recovery in closed queueing networks through swapping customers is assumed in Perros et al. (1988), with zero transition time assumed between deadlocked states and the corresponding resolved state. Time to resolve deadlock may not be negligible in reality. Deadlock detection and recovery is listed as one of the two possible solutions for handling deadlock in queueing networks in Akyildiz (1989), although there is no further discussion.

Note that a number of deadlock types are defined in Venkatesh et al. (1998). The terminology of that paper differs greatly to that used here. Notably the concept described by their term 'Transient deadlock' is not considered a deadlock situation at all in this paper according to Definition 1.

4. Deadlock detection

In order to detect when deadlock has occurred in a queueing network simulation, a state digraph is used, a form of wait-for graph. In previous literature on wait-for graphs these are bespoke graphs that represent system states, where edges denote some form of waiting or blockage relationships. Here we present a generic state digraph that is defined for all FIFO queueing networks that exhibit Type I blocking:

Definition 2. The state digraph $D(t)$ of a queueing network is defined by that network's state at any time t . Vertices of the state digraph correspond to servers of the network. A directed edge denotes a blockage relationship in the following manner: if a customer at the k th server of node i is blocked from entering node j , then there are directed edges from the vertex corresponding to node i 's k th server to every vertex corresponding to the servers of node j .

To illustrate this concept Fig. 3 shows examples of queueing networks in and out of deadlock, and the corresponding state digraph in each case.

These graph theoretic terms will be used throughout this paper (Bang-Jensen & Gutin, 2008; Gibbons, 1985; Wilson, 1970):

- Two vertices v_1 and v_2 are said to be *weakly connected* if there is a directed path from v_1 to v_2 or a directed path from v_2 to v_1 .
- Two vertices v_1 and v_2 are said to be *strongly connected* if there is a directed path from v_1 to v_2 and a directed path from v_2 to v_1 . Note that strongly connected vertices are also weakly connected.
- A *weakly connected component* of a digraph is a subgraph induced by a maximal subset of weakly connected vertices.
- A *strongly connected component* of a digraph is a subgraph induced by a maximal subset of strongly connected vertices.
- The *out-degree* of v_1 is the number of out-edges incident to v_1 .

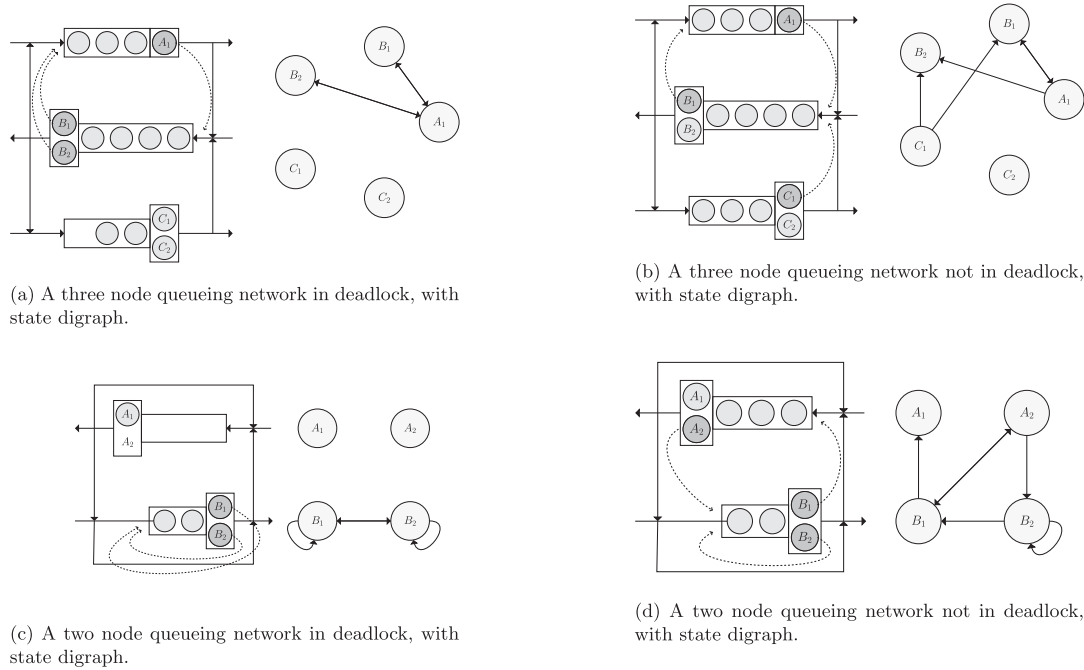


Fig. 3. Examples of state digraphs with their corresponding queueing networks.

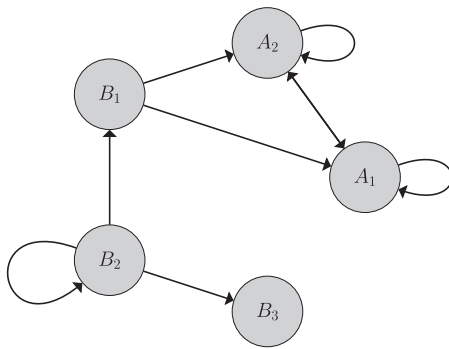


Fig. 4. State digraph of the counter-example.

- A *sink* is a vertex whose out-degree is zero.
- A *knot*, or terminal strong component, is a strongly connected component containing no vertices with a path to any vertices outside that component.

Consider a vertex v in $D(t)$. Some observations:

- If the server corresponding to v is unoccupied, then v has no incident edges.
- It can be interpreted that all vertices with a path to v correspond to servers whose individuals are being blocked directly or indirectly by the customer at the server corresponding to v .
- Similarly it can be interpreted that all vertices that v has a path to correspond to servers whose occupants are directly or indirectly blocking the customer at the server corresponding to v .
- It is clear that if all vertices that v has a path to correspond to servers occupied by blocked individuals, then the system is in deadlock at time t .

The following results are used detect deadlock for open restricted queueing networks.

Theorem 1. A deadlocked state arises at time t if and only if $D(t)$ contains a knot.

Proof. Consider a queueing network with set of servers S . Consider the state digraph $D(t) = (V, E(t))$. Note that there is a 1–1 pairing between the elements of S and the elements of V , by Definition 2 of $D(t)$.

- Assume the system is in deadlock at time t . By Definition 1 (of deadlock) there exists $\mathcal{S} \subseteq S$ a subset of servers with blocked customers blocked only by customers at servers in \mathcal{S} . Consider $\mathcal{V} \subseteq V$ corresponding to \mathcal{S} .

For each $s \in \mathcal{S}$ the corresponding $v \in \mathcal{V}$ has at least one out-edge in $E(t)$ because the customer at s is blocked (by Definition 2).

By Definition 1 every $v \in \mathcal{V}$ has at least one path in $E(t)$ to a vertex in \mathcal{V} , and has no path in $E(t)$ to any vertex outside of \mathcal{V} . By definition of a knot, there exists $G = (\mathcal{V}, \mathcal{E})$, where $\mathcal{E} \subseteq E(t)$, such that G is either a knot or a collection of knots.

- Assume that $D(t)$ contains a knot $G = (\mathcal{V}, \mathcal{E})$. Consider $\mathcal{S} \subseteq S$ corresponding to \mathcal{V} .

As G is a knot, every $v \in \mathcal{V}$ has an out-edge, thus every customer at $s \in \mathcal{S}$ is blocked (by Definition 2).

As G is a knot, there is no path in $E(t)$ from any $v \in \mathcal{V}$ to any vertex outside of \mathcal{V} . Therefore every customer at $s \in \mathcal{S}$ is blocked directly or indirectly by customers at servers in \mathcal{S} .

By Definition 1 this implies the system is in deadlock at time t . □

The knot condition can be simplified for specific cases. Theorem 2 may offer computational advantages in cases where it is easier to identify weakly connected components to knots.

Theorem 2. For queueing networks:

1. with one node
2. with two nodes, each with two or fewer parallel servers
3. with a finite amount of nodes, each with a single-server

a deadlocked state arises at time t if and only if $D(t)$ contains a weakly connected component without a sink.

Proof. To prove the result in one direction, each case is considered separately:

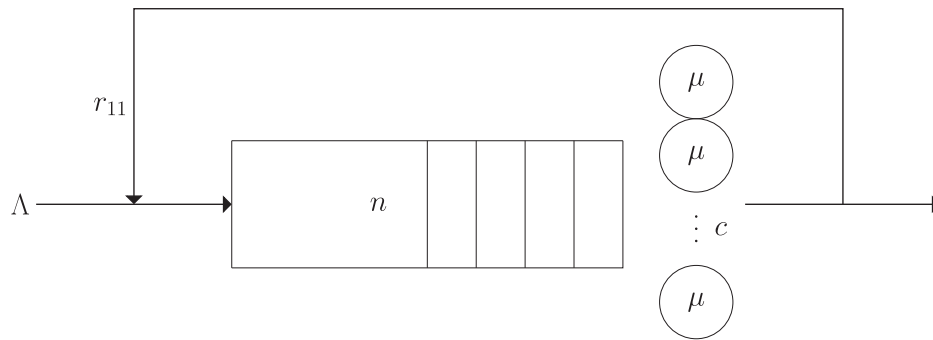


Fig. 5. An open one node multi-server restricted queueing network.

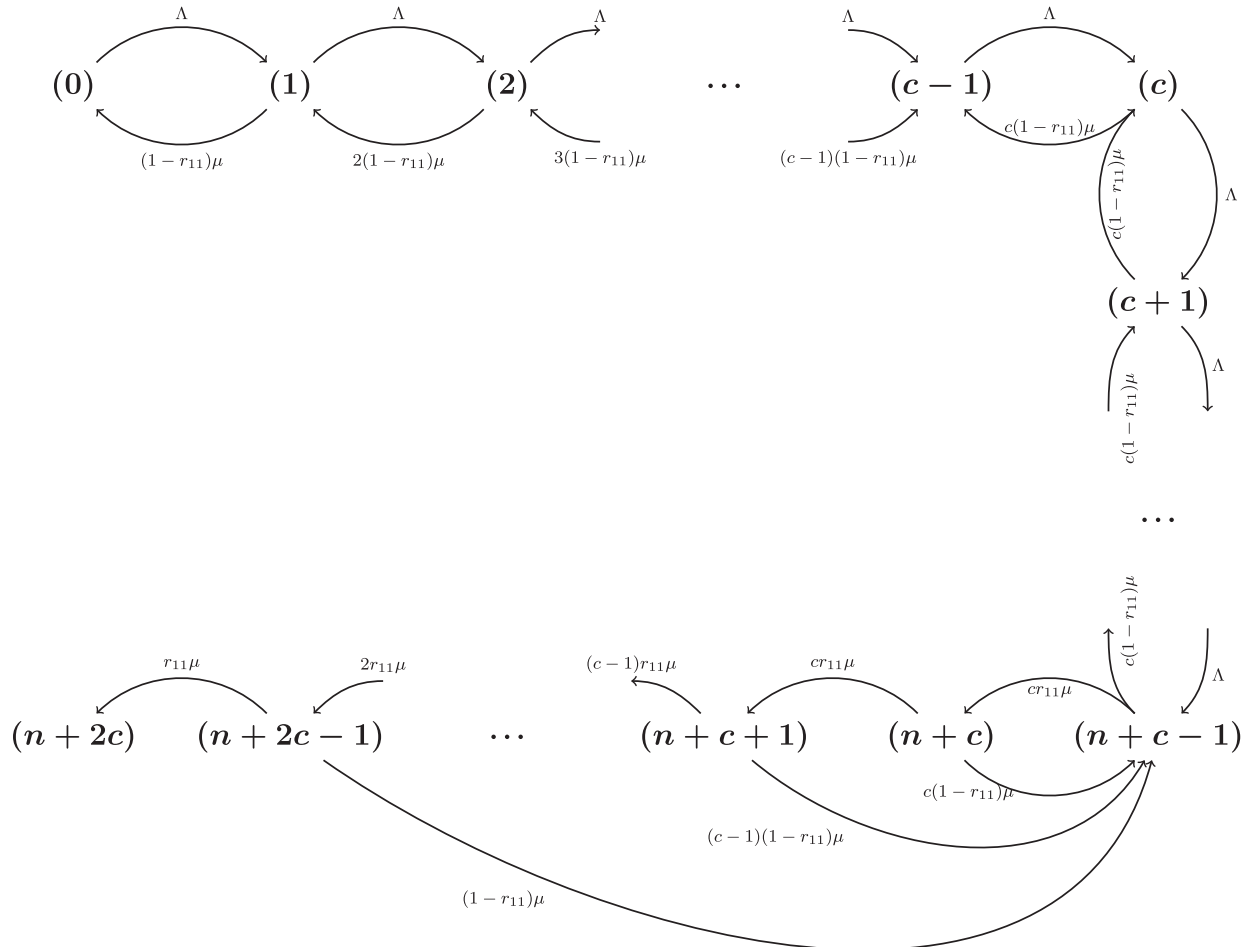


Fig. 6. Diagrammatic representation of the Markov chain for a multi-server one node system.

1. Consider a one node queueing network.
If there is deadlock, then all servers are occupied by blocked individuals, and so all vertices corresponding to those servers have an out-edge. Thus there are no sinks.
2. Consider a two node queueing network, each node with 2 or fewer parallel servers.
If both nodes are involved in the deadlock, so there is at least one customer in node 1 blocked from entering node 2, and at least one customer from node 2 blocked from entering node 1, then all servers in node 1 and node 2 in $D(t)$ will have out edges as they are occupied by a blocked individual. The servers of node 1 and 2 consist of the entirety of $D(t)$, and so there is no sink nodes.

Now consider the case when only one node is involved in the deadlock. Without loss of generality, consider that node 1 is in deadlock with itself, then the servers of node 1 have out-edges. For the servers of node 2 to be part of that weakly connected component, there either needs to be an edge from a server in node 1 to a server in node 2, or an edge from a server in node 2 to a server in node 1. An edge from a server in node 1 to a server in node 2 implies that a customer from node 1 is blocked from entering node 2, and so node 1 is not in deadlock with itself. An edge from a server in node 2 to a server in node 1 implies that a customer in node 2 is blocked from entering node 1. In this case one server in node 2 has an out-edge. Now either the other server of node 2 is empty or still in service, and so is not part of that weakly connected component, or the other

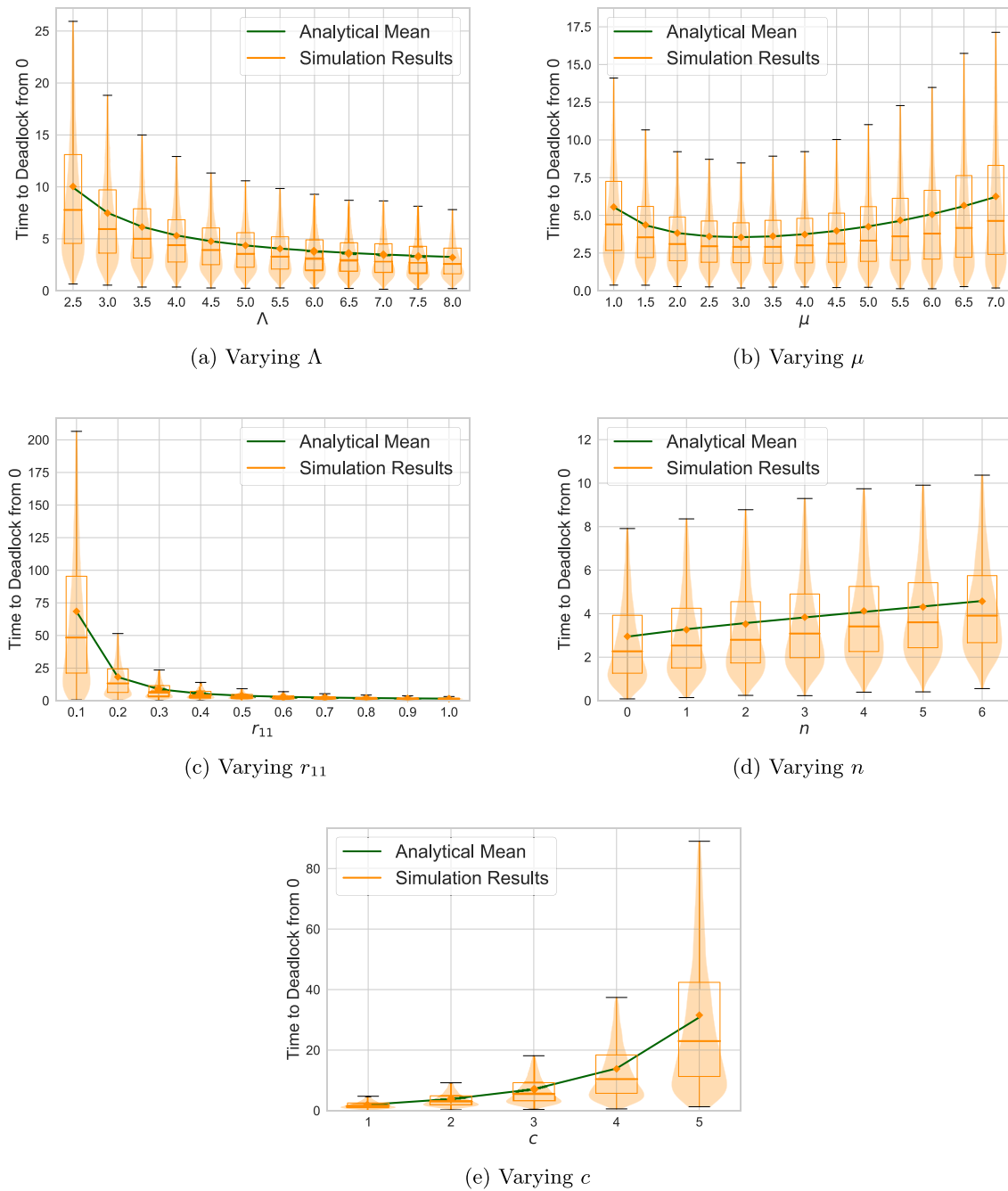


Fig. 7. Time to deadlock in the multi-server one node system, analytical and simulation results (10,000 repetitions).

server's customer is blocked and so has an out edge. Thus there are no sinks.

3. Consider a queueing network with N nodes, each with a single-server.

If $1 \leq n \leq N$ nodes are involved in the deadlock, then each server in those n nodes has a blocked customer, and so the corresponding vertex in the state digraph has an out-edge.

Of the nodes not involved in that deadlock, the vertices corresponding to their servers can only be in the same weakly connected component if:

- They contain a blocked individual that is blocked to the nodes involved in the deadlock.
- Individuals in the deadlocked nodes are blocked to those nodes not involved in the deadlock.

In the first case, the vertices corresponding to the servers in those nodes will all have an out-edge. In the second case it is implied that the customers at servers in deadlocked nodes are blocked to both a node in deadlock and a node not in deadlock, which is not possible (customers can only be blocked to one location at a time).

Thus there are no sinks.

Proving the result in the other direction is equivalent to proving that a weakly connected component without a sink contains a knot:

- Consider a weakly connected component, G , of $D(t)$.
- Assume G contains no knots. By definition of a knot, this implies:
 - G contains a sink; or

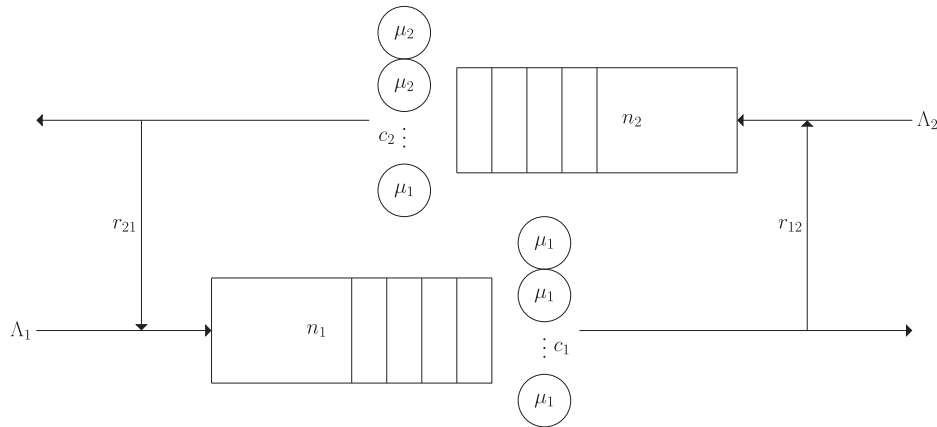


Fig. 8. An open two node multi-server restricted queueing network.

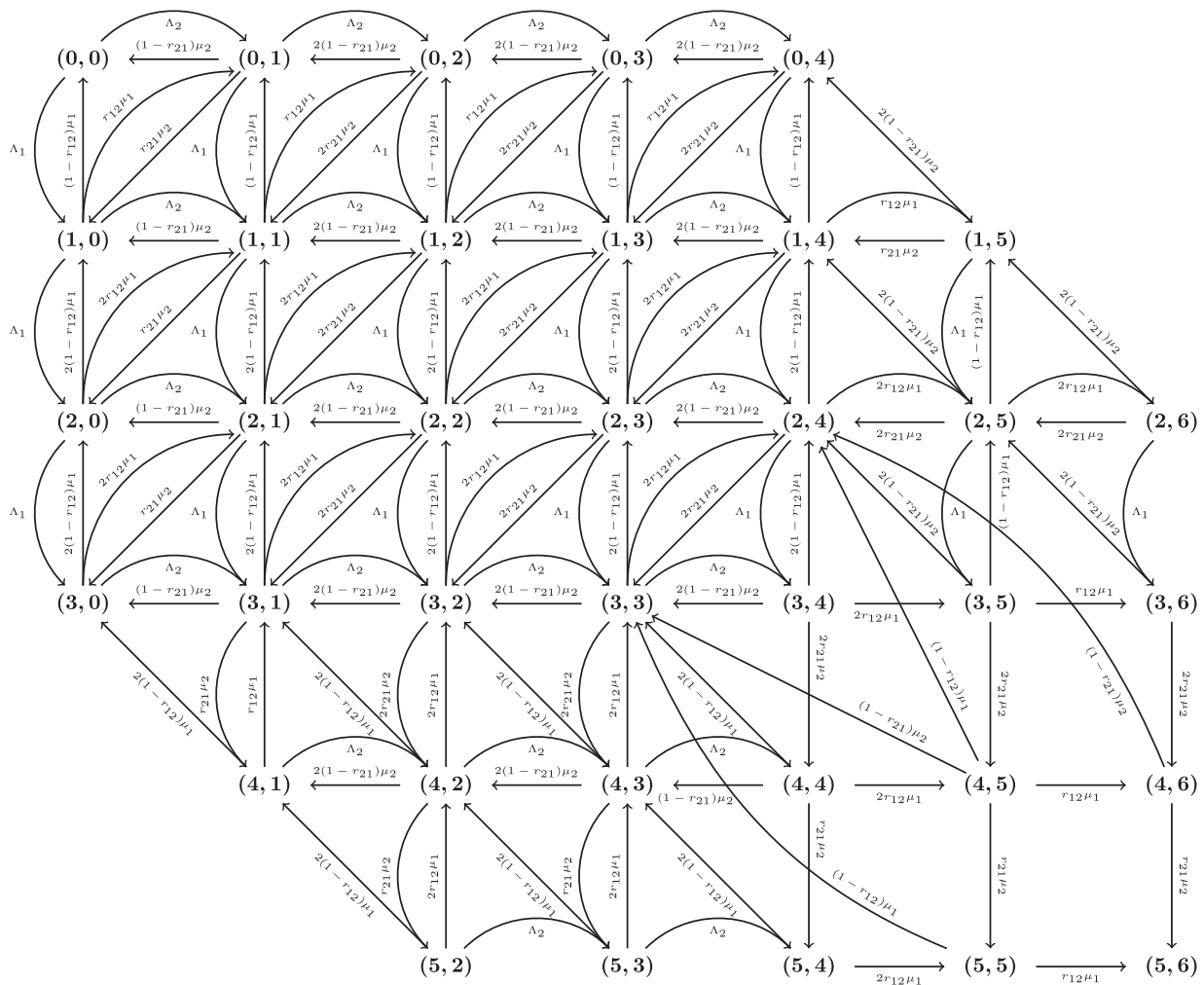


Fig. 9. Diagrammatic representation of the Markov chain for a multi-server two node system without self-loops with $n_1 = 1, n_2 = c_1 = c_2 = 2$. The deadlock state is (5, 6).

- G contains a vertex with a path to another vertex outside of G (contradicting the fact that G is a weakly connected component).

Thus the existence of a weakly connected component without a sink in $D(t)$ implies that there is a knot in $D(t)$, and the result follows by applying [Theorem 1](#). \square

In the general case using the result of [Theorem 2](#) is not sufficient to detect deadlock. In order to illustrate this, consider the

following counter-example of a two node queueing network, where node A has two servers, node B has three servers. Beginning with all servers occupied by customers in service and full queues. The customer at server A_1 becomes blocked to node A . The customer at server B_1 becomes blocked to node A . The customer at server B_2 becomes blocked to node B . The customer at server A_2 becomes blocked to node A . Node A is now in deadlock. The resulting state

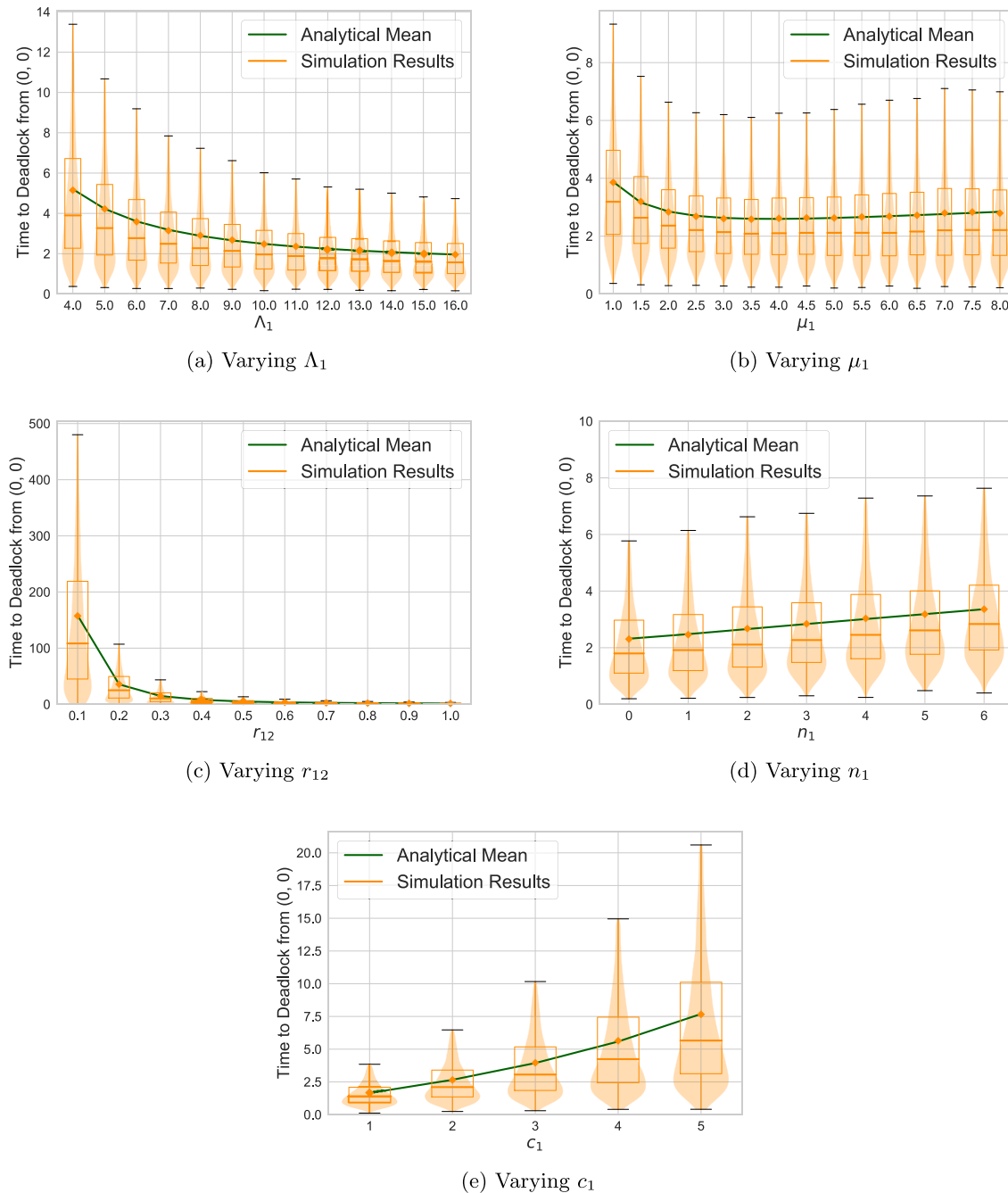


Fig. 10. Time to deadlock in the multi-server two node system without self loops, analytical and simulation results (10,000 repetitions).

digraph, shown in Fig. 4, has a weakly connected component with a sink.

For the purposes of this paper, a simulation model is used to verify that the results of this section and the analytical model in Section 5 are in agreement. Specifically the time taken to reach deadlock from an empty system is investigated, and the simulation model gives information on the distribution of the times to deadlock. The model is built using Ciw (The Ciw library developers, 2017). This is an object oriented framework in Python (The Python Software Foundation, 2015), with care taken to ensure reproducibility of the results (Hong, Crick, Gent, & Kotthoff, 2015).

The digraph $D(t)$ is implemented as an attribute of the queueing network and is updated at the appropriate events. Note that a brute force algorithm is used to check whether any strongly

connected component of $D(t)$ is a knot in order to implement Theorem 1. More efficient algorithms could be used for other specific use cases.

In the next section Markov models of three queueing networks are built, and their deadlock properties discussed.

5. Markovian models of deadlocking queueing networks

The following three networks describe all possible configurations of deadlocking queueing networks with two or fewer nodes:

1. Open one node, multi-server restricted queueing network with feedback loop (Section 5.1).
2. Open two node, multi-server restricted queueing network with routes between nodes (Section 5.2).

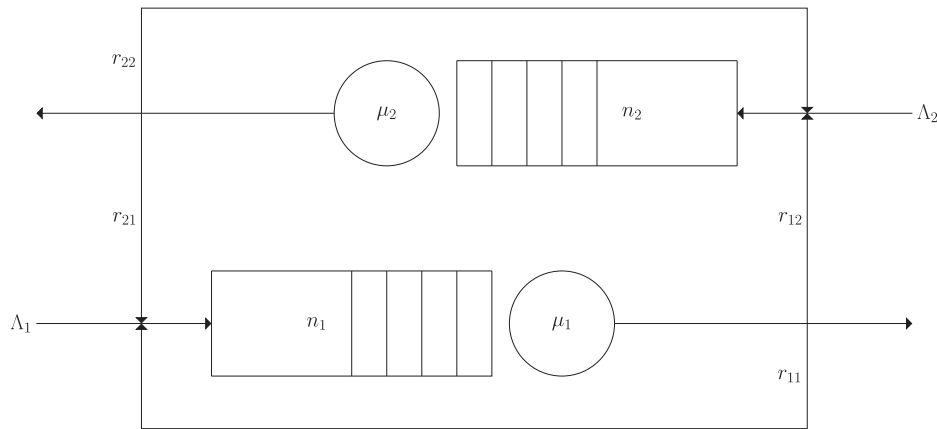


Fig. 11. An open two node single-server restricted queueing network.

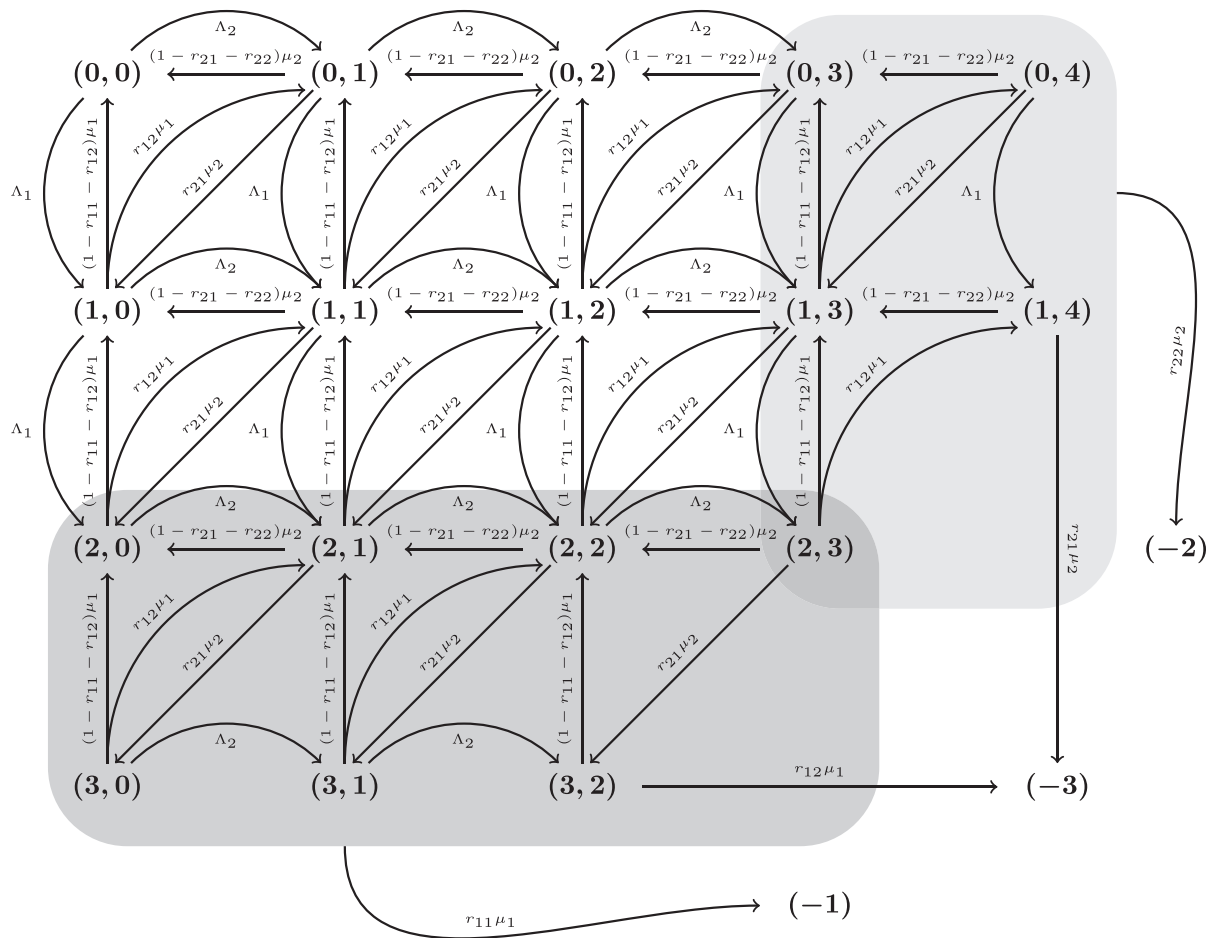


Fig. 12. Diagrammatic representation of the Markov chain for single server two node system with $n_1 = 1$ and $n_2 = 2$.

3. Open two node, multi-server restricted queueing network with routes between nodes and self-loops.

In this section Markov models are built for networks 1 and 2, and their expected time to deadlock found. The state space for network 3 is too large to model in a similar way to the others, and so isn't considered in this paper. A single server version is modelled however (Section 5.3), and the multi-server system is briefly discussed in Section 6.

In general a continuous Markov chain model of a deadlocking queueing network is defined by a set of states S and the transition rates between these states q_{s_1, s_2} . Each state $s \in S$ uniquely

defines a configuration of customers around the queueing network. Deadlock states are also present, either denoted by that specific configuration of customers, or by negative numbers, for example (-1) . Deadlocked states cannot transition to any other state, and so are absorbing states of the Markov chain. Therefore any queueing network that can experience deadlock is guaranteed to experience deadlock, as absorbing Markov chains are guaranteed to enter one of its absorbing states.

The expected time until deadlock is reached is equivalent to the expected time to absorption of the Markov chain, which can be found using classic results (Stewart, 2009). The continuous Markov

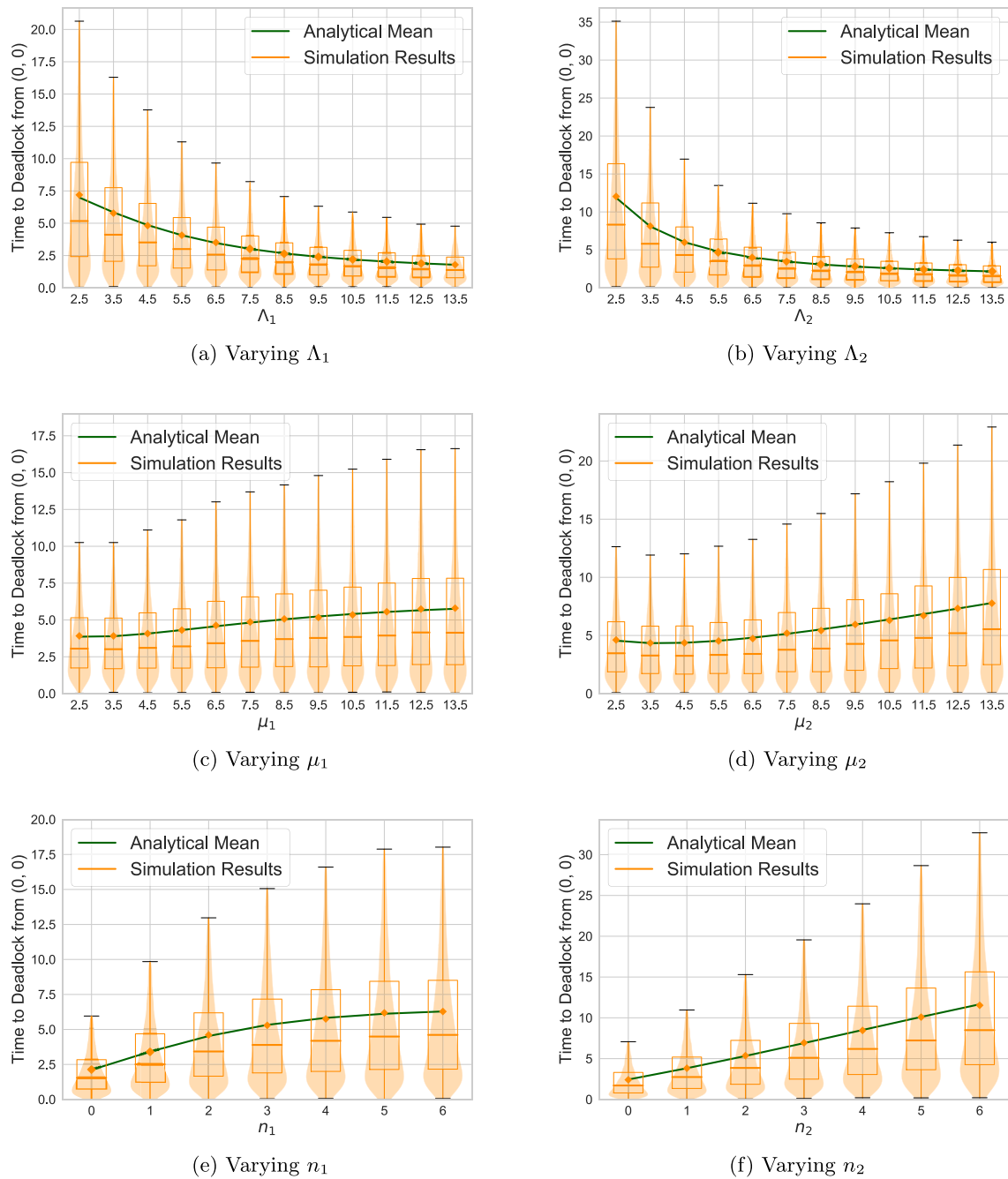


Fig. 13. Time to deadlock in the single-server two node system, analytical and simulation results (10,000 repetitions).

chain is converted to a discrete time absorbing Markov chain with canonical form:

$$P = \begin{pmatrix} T & U \\ 0 & I \end{pmatrix}$$

where I is the identity matrix. Now the expected number of time steps until absorption starting from the i th state is the i th element of the vector

$$(I - T)^{-1}e \tag{5.1}$$

where e is a vector of 1s.

Therefore by discretising the continuous Markov chain and ensuring the correct order of states, the expected number of time steps to absorption, which corresponds to deadlock, can be found.

This can be converted back to continuous time by multiplying by the time step used in the discretisation process.

When there is more than one deadlock state, there is more than one absorbing state in the Markov chain. Here the expected time to absorption is the expected time to a deadlock state, whichever one that may be.

5.1. One node multi-server

Consider the open one node multi-server restricted queueing network shown in Fig. 5. This shows an $M/M/c/n$ queue where customers arrive at a rate of Λ and served at a rate μ . Once a customer has finished service they rejoin the queue with probability r_{11} , and so exit the system with probability $1 - r_{11}$.

The state space is given by:

$$S = \{i \in \mathbb{N} \mid 0 \leq i \leq n + 2c\}$$

where i denotes the number of individuals in the system plus the number of individuals who are blocked. For example, $i = n + c + 2$ denotes a full system, $n + c$ individuals in the node, and 2 of those individuals are also blocked. The state $i = n + 2c$ denotes the deadlock state, that is every customer with a server is blocked.

Define $\delta = i_2 - i_1$ for all $i_k \in S$. The transitions are given by Eqs. (5.2) and (5.3).

$$q_{i_1, i_2} = \begin{cases} \Lambda & \text{if } \delta = 1 \\ (1 - r_{11})\mu \min(i, c) & \text{if } \delta = -1 \\ 0 & \text{otherwise} \end{cases} \quad \text{if } i_1 < n + c \quad (5.2)$$

$$q_{i_1, i_2} = \begin{cases} (c - b)r_{11}\mu & \text{if } \delta = 1 \\ (1 - r_{11})(b - k)\mu & \text{if } \delta = -b - 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{if } i_1 = n + c + b \quad \forall \quad 0 \leq b \leq c \quad (5.3)$$

where b denotes the number of blocked customers. The Markov chain is shown in Fig. 6.

Fig. 7 shows the effect of varying the parameters of the above Markov model. Base parameters of $\Lambda = 6$, $n = 3$, $\mu = 2$, $r_{11} = 0.5$ and $c = 2$ were used. It can be seen that increasing the arrival rate Λ and the routing probability r_{11} results in reaching deadlock faster. This is intuitive as increasing these parameters results in the queue filling up quicker. Increasing the queueing capacity n results in reaching deadlock slower. Again this is intuitive, as increasing the queueing capacity allows more customers in the system before reaching deadlock.

Increasing the amount of servers has a similar effect to increasing the queueing capacity, as there are now more transient states to go through before reaching the deadlock state. Varying the amount of servers has a greater effect on the time to deadlock however, as any state in which customers are blocked, $i \in [n + c + 1, n + 2c]$, can jump back to state $i = n + c - 1$ simply with a service where the customer does not rejoin the queue. Increasing the amount of servers also increases the rate at which customer leave the system, but not the rate at which customers enter the system. This means that the rate of increase of the number of customers in the system increases, however the rate of decrease of the number of customers in the system does not change, thus it would take longer to reach a full system, a requirement of deadlock.

The behaviour as the service rate μ varies is not monotonic, as the service rate contributes towards both moving customers from the system and allowing customers to rejoin the queue, causing blockages and deadlock. If the function $\omega(\mu)$ describes the expected time to deadlock of this system as the service rate μ varies, and all other parameters are fixed, then it is observed that $\omega(\mu)$ has one critical point and is a local minimum for $\mu \in (0, \infty)$.

The observed bowl shaped curve can be explained by considering the effect of varying the service rate. At $\lim_{\mu \rightarrow 0} \omega(\mu)$ there is infinite service time, and so there is infinite time until deadlock. At $\lim_{\mu \rightarrow \infty} \omega(\mu)$ there is zero service time, the queue can never form, and so there is infinite time to deadlock. At low service rates, below a certain threshold $\hat{\mu}$, the arrival rate is relatively large compared to the service rate, a saturated system can be assumed. At this point services in which a customer exits the system do not have much of an effect on the system state, another immediate arrival can be assumed. However services in which a customer wishes to rejoin the queue results in a blockage as the system is saturated. Therefore, increasing the service rate here increases the chance of a blockage, and so the chance of deadlock. Above $\hat{\mu}$ the service rate is large enough that a saturated system cannot be assumed, and so services in which the customer exits the system do

have an affect on the number of customers in the system. Thus increasing the service rate increases the rate at which customers are removed from the system, and so there is less chance of reaching deadlock.

5.2. Two node multi-server without self-loops

Consider the open two node multi-server restricted queueing network shown in Fig. 8. This shows two $M/M/c_i/n_i$ queues, with service rates μ_i and external arrival rates Λ_i . All routing probabilities r_{ij} may be positive apart from self-loops r_{ii} , for each node i . Note that this system is equivalent to the one described in Section 2.

The state space is given by:

$$S = \{(i, j) \in \mathbb{N}^2 \mid i \leq n_1 + c_1 + j, j \leq n_2 + c_2 + i\}$$

where i denotes the number of individuals at node 1 plus the number of individuals blocked waiting to enter node 1, and j denotes the number of individuals at node 2 plus the number of individuals blocked waiting to enter node 2. For example, $(i, j) = (n_1 + c_1 + 2, n_2 + c_2 + 1)$ denotes a full system, $n_1 + c_1$ individuals at node 1, two of whom are blocked waiting to enter node 2; $n_2 + c_2$ individuals at node 2, one of whom is blocked waiting to enter node 1. The state $(i, j) = (n_1 + c_1 + c_2, n_2 + c_2 + c_1)$ denotes the deadlock state. The Markov chain is shown in Fig. 9.

Define $\delta = (i_2, j_2) - (i_1, j_1)$, $b_1 = \max(0, i_1 - n_1 - c_1)$, $b_2 = \max(0, i_2 - n_2 - c_2)$, $s_1 = \min(i_1, c_1) - b_2$ and $s_2 = \min(i_2, c_2) - b_1$ for all $(i_k, j_k) \in S$. Then the transitions $q_{(i_1, j_1), (i_2, j_2)}$ are given in Table 1.

The values b_1 and b_2 correspond to the number of customers blocked to node 1 and node 2, respectively. The values s_1 and s_2 correspond to the amount of customers currently in service at node 1 and node 2, respectively.

Fig. 10 shows the effect of varying the parameters of the above Markov model. Base parameters of $\Lambda_1 = 9$, $\Lambda_2 = 7.5$, $n_1 = 2$, $n_2 = 1$, $\mu_1 = 5.5$, $\mu_2 = 6.5$, $r_{12} = 0.7$, $r_{21} = 0.6$, $c_1 = 2$ and $c_2 = 2$ were used. Only plots for the parameters corresponding to node 1 are shown, node 2 shows similar behaviour. Similar behaviour is observed to that seen in Fig. 7.

5.3. Two node single-server with self-loops

Consider the open two node single-server restricted queueing network shown in Fig. 11. This shows two $M/M/1/n_i$ queues with service rates μ_i and external arrival rates Λ_i . All routes are possible, where the routing probability from node i to node j is denoted by r_{ij} .

The state space is given by:

$$S = \{(i, j) \in \mathbb{N}^2 \mid 0 \leq i + j \leq n_1 + n_2 + 2\} \cup \{(-1), (-2), (-3)\}$$

where i denotes the number of individuals:

- In service or waiting at the first node.
- Occupying a server but having finished service at the second node and waiting to join the first.

where j denotes the number of individuals:

- In service or waiting at the second node.
- Occupying a server but having finished service at the first node and waiting to join the second.

and the state (-3) denotes the deadlock state caused by both nodes; (-1) denotes the deadlock state caused by the first node only; and (-2) denotes the deadlock state caused by the second node only.

Table 1
Table of transitions $q_{(i_1, j_1), (i_2, j_2)}$ for a multi-server two node network.

	$j_1 < n_2 + c_2$	$j_1 = n_2 + c_2$	$j_1 > n_2 + c_2$
$i_1 < n_1 + c_1$	Λ_1 if $\delta = (1, 0)$ Λ_2 if $\delta = (0, 1)$ $r_{12}s_1\mu_1$ if $\delta = (-1, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, -1)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	Λ_1 if $\delta = (1, 0)$ $r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, -1)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	Λ_1 if $\delta = (1, 0)$ $r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (0, -1)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (-1, -1)$
$i_1 = n_1 + c_1$	Λ_2 if $\delta = (0, 1)$ $r_{12}s_1\mu_1$ if $\delta = (-1, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	$r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	$r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, 0)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (-1, -1)$
$i_1 > n_1 + c_1$	Λ_2 if $\delta = (0, 1)$ $r_{12}s_1\mu_1$ if $\delta = (-1, 0)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, -1)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	$r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-1, -1)$ $(1 - r_{21})s_2\mu_2$ if $\delta = (0, -1)$	$r_{12}s_1\mu_1$ if $\delta = (0, 1)$ $r_{21}s_2\mu_2$ if $\delta = (1, 0)$ $(1 - r_{12})s_1\mu_1$ if $\delta = (-\min(b_1 + 1, b_2 + 1), -\min(b_1, b_2 + 1))$ $(1 - r_{21})s_2\mu_2$ if $\delta = (-\min(b_1 + 1, b_2 + 1), -\min(b_1 + 1, b_2 + 1))$

Define $\delta = (i_2, j_2) - (i_1, j_1)$ for all $(i_k, j_k) \in S$. The transitions are given by Eqs. (5.4)–(5.7).

explained by the existence of more than one deadlock state for this system. Deadlock state (-1) involves node 1 only, deadlock state

$$q_{(i_1, j_1), (i_2, j_2)} = \left\{ \begin{array}{ll} \Lambda_1 & \text{if } i_1 < n_1 + 1 \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (1, 0)$$

$$\left\{ \begin{array}{ll} \Lambda_2 & \text{if } j_1 < n_2 + 1 \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (0, 1)$$

$$\left\{ \begin{array}{ll} (1 - r_{11} - r_{12})\mu_1 & \text{if } j_1 < n_2 + 2 \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (-1, 0)$$

$$\left\{ \begin{array}{ll} (1 - r_{21} - r_{22})\mu_2 & \text{if } i_1 < n_1 + 2 \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (0, -1)$$

$$\left\{ \begin{array}{ll} r_{12}\mu_1 & \text{if } j_1 < n_2 + 2 \text{ and } (i_1, j_1) \neq (n_1 + 2, n_2) \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (-1, 1)$$

$$\left\{ \begin{array}{ll} r_{21}\mu_2 & \text{if } i_1 < n_1 + 2 \text{ and } (i_1, j_1) \neq (n_1, n_2 + 2) \\ 0 & \text{otherwise} \end{array} \right\} \text{ if } \delta = (1, -1)$$

$$0 \text{ otherwise} \tag{5.4}$$

$$q_{(i_1, j_1), (-1)} = \left\{ \begin{array}{ll} r_{11}\mu_1 & \text{if } i > n_1 \text{ and } j < n_2 + 2 \\ 0 & \text{otherwise} \end{array} \right\} \tag{5.5}$$

$$q_{(i_1, j_1), (-2)} = \left\{ \begin{array}{ll} r_{22}\mu_2 & \text{if } j > n_2 \text{ and } i < n_1 + 2 \\ 0 & \text{otherwise} \end{array} \right\} \tag{5.6}$$

$$q_{(i_1, j_1), (-3)} = \left\{ \begin{array}{ll} r_{21}\mu_2 & \text{if } (i, j) = (n_1, n_2 + 2) \\ r_{12}\mu_1 & \text{if } (i, j) = (n_1 + 2, n_2) \\ 0 & \text{otherwise} \end{array} \right\} \tag{5.7}$$

$$q_{-1, s} = 0 \tag{5.8}$$

$$q_{-2, s} = 0 \tag{5.9}$$

$$q_{-3, s} = 0 \tag{5.10}$$

Note that there are now three different deadlock states, thus two more ways to reach deadlock, Eqs. (5.5) and (5.6).

For $n_1 = 1$ and $n_2 = 2$, the resulting Markov chain is shown in Fig. 12.

Fig. 13 shows the effect on the time to deadlock of varying the parameters of the above Markov model. Base parameters of $\Lambda_1 = 4$, $\Lambda_2 = 5$, $n_1 = 3$, $n_2 = 2$, $\mu_1 = 10$, $\mu_2 = 8$, $r_{11} = 0.1$, $r_{12} = 0.25$, $r_{21} = 0.15$ and $r_{22} = 0.1$ are used.

In general, similar behaviour is observed to that seen in Figs. 7 and 10. A notable difference however is that the increase or decrease in the time to deadlock flattens as the parameter in question increases or decreases. This is observable in Fig. 13(e). This is

(-2) involves node 2 only, while deadlock state (-3) involves both nodes. Therefore if a parameter at node 1 is increased or decreased such that the time to a deadlock involving node 1, states (-1) and (-3) , approaches infinity, then the overall time to deadlock of the system will become unchanging, as varying that parameter will not effect the time to deadlock state (-2) .

6. Conclusions

This paper has explored deadlock in open restricted queueing networks. It has been shown that analysing a queueing network's corresponding state digraph is sufficient to detect when deadlock occurs in queueing networks. In general the presence of a knot in the state digraph will highlight that deadlock has occurred in the network, however for special cases the presence of a weakly connected component with no sink is sufficient. Incorporating this into a simulation model, time to deadlock can be observed.

Markov models of three deadlocking queueing networks have been built. Using linear algebraic techniques the expected time to deadlock from each state was found, and its behaviour as system parameters are varied was explored. These analytical results were compared with results obtained from the simulation model.

Further research is needed to build a Markov model of the open two node, multi-server restricted queueing network with routes between nodes and feedback loops, that is network 3 from Section 5. In networks 1 and 2 customers only have one potential destination, and so customers may only get blocked from moving to one destination. In network 3 with single servers, although

customers have two destinations, a blockage to the same node immediately results in deadlock. In all these cases, the unblocking mechanism is simple, as there is only ever one option of which node a customer joins when unblocked. However in network 3 with multiple servers, there are two destination nodes to which a customer may join when unblocked. Therefore, any representations of any states with blocked customers also need to hold information about these customers' destination nodes.

In addition to this, the order in which customers become blocked is important. In networks 1 and 2 when space become available at a node there is only one other node from which a blocked customer can become unblocked, however in network 3 a node that has space available must accept the customer that has been blocked longest to that node. Therefore all states with blocked customers are also required to record the order in which the customers become blocked. Combining the two requirements above, it is clear that as the number of servers increases, the size of the state space for this queueing network quickly grows combinatorially. Therefore it is not possible to consider this state space in the same way as for networks 1 and 2.

For the Markov models built in this paper Poisson arrivals and exponential service rates are assumed, and only blocking of Type I is considered. A future research direction could be to model other service and arrival distributions using phase-type distributions, and incorporating these into the Markov models of deadlocking queueing networks. Blocking of Type II and III should also be considered, both in the analytical models and whether the deadlock detection method presented here still holds. Systems under Type III blocking with random destination will not reach deadlock, as there is a non-zero probability of a blocked customer leaving the system. This type of blocking may be considered a deadlock prevention mechanism.

Acknowledgements

The authors would like to thank and express gratitude to all anonymous referees whose comments and suggestions have greatly improved the work. The following software libraries have been used in this work: the matplotlib library (Hunter, 2007) and the NumPy library (Walt, Colbert, & Varoquaux, 2011).

References

- Akyildiz, I. (1989). Product form approximations for queueing networks with multiple servers and blocking. *IEEE Transactions on Computers*, 38(1), 99–114.
- Allon, G., Deo, S., & Lin, W. (2013). The impact of size and occupancy of hospital on the extent of ambulance diversion: Theory and evidence. *Operations Research*, 61(3), 544–562.
- Avi-Itzhak, B., & Yadin, M. (1965). A sequence of two servers with no intermediate queue. *Management Science*, 11(5), 553–564.
- Baber, J. (2008). *Queues in series with blocking*. Cardiff University (Ph.D. thesis).
- Bang-Jensen, J., & Gutin, G. Z. (2008). *Digraphs: Theory, algorithms and applications*. Springer Science & Business Media.
- Belik, F. (1990). An efficient deadlock avoidance technique. *IEEE Transactions on Computers*, 39(7), 882–888.
- Cheng, J. (1990). Task-wait-for graphs and their application to handling tasking deadlocks. In *Proceedings of the conference on tri-ada'90* (pp. 376–390). ACM.
- Cho, H., Kumaran, T., & Wysk, R. (1995). Graph-theoretic deadlock detection and resolution for flexible manufacturing systems. *IEEE Transactions on Robotics and Automation*, 11(3), 413–421.
- Coffman, E., & Elphick, M. (1971). System deadlocks. *Computing Surveys*, 3(2), 67–78.
- Dallery, Y., & Frein, Y. (1993). On decomposition methods for tandem queueing networks with blocking. *Operations Research*, 41(2), 386–399.
- Deuermeier, B., Curry, G., Duchowski, A., & Venkatesh, S. (1997). An automatic approach to deadlock detection and resolution in discrete simulation systems. *INFORMS Journal on Computing*, 9(2), 195–205.
- Dijkstra, E. (1982). The mathematics behind the bankers algorithm. In *Selected writings on computing: A personal perspective* (pp. 308–312). Springer.
- Elmagarmid, A. (1986). A survey of distributed deadlock detection algorithms. *ACM Sigmod Record*, 15(3), 37–45.
- Ezpeleta, J., Tricas, F., Garc eja-Vall s, F., & Colom, J. (2002). A banker's solution for deadlock avoidance in FMS with flexible routing and multiresource states. *IEEE Transactions on Robotics and Automation*, 18(4), 621–625.
- Florian, M., Mahut, M., & Tremblay, N. (2008). Application of a simulation-based dynamic traffic assignment model. *European Journal of Operational Research*, 189(3), 1381–1392.
- Gibbons, A. (1985). *Algorithmic graph theory*. Cambridge University Press.
- Gordon, W., & Newell, G. (1967). Cyclic queueing systems with restricted length queues. *Operations Research*, 15(2), 266–277.
- Holt, R. (1972). Some deadlock properties of computer systems. *ACM Computing Surveys*, 4(3), 179–196.
- Hong, N., Crick, T., Gent, I., & Kotthoff, L. (2015). Top tips to make your research irreproducible. preprint arXiv:1504.00062.
- Hunt, G. (1956). Sequential arrays of waiting lines. *Operations Research*, 4(6), 674–683.
- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90–95.
- Kawadkar, P., Prasad, S., & Dwivedi, A. (2014). Deadlock avoidance based on bankers algorithm for waiting state processes. *International Journal of Innovative Science and Modern Engineering*, 2(12), 1–2.
- Koizumi, N., Kuno, E., & Smith, T. (2005). Modeling patient flows using a queueing network with blocking. *Health Care Management Science*, 8(1), 49–60.
- Korporaal, R., Ridder, A., Klopprogge, P., & Dekker, R. (2000). An analytic model for capacity planning of prisons in the Netherlands. *The Journal of the Operational Research Society*, 51(11), 1228–1237.
- Kundu, S., & Akyildiz, I. (1989). Deadlock buffer allocation in closed queueing networks. *Queueing Systems*, 4(1), 47–56.
- Latouche, G., & Neuts, M. (1980). Efficient algorithmic solutions to exponential tandem queues with blocking. *SIAM Journal on Algebraic Discrete Methods*, 1(1), 93–106.
- Liebeherr, J., & Akyildiz, I. (1995). Deadlock properties of queueing networks with finite capacities and multiple routing chains. *Queueing Systems*, 20(3–4), 409–431.
- Manzano-Santaella, A. (2010). From bed-blocking to delayed discharges: precursors and interpretations of a contested concept. *Health Services Management Research*, 23(3), 121–127.
- Marchetti, O., & Munier-Kordon, A. (2009). A sufficient condition for the liveness of weighted event graphs. *European Journal of Operational Research*, 197(2), 532–540.
- Onvural, R. (1990). Survey of closed queueing networks with blocking. *ACM Computing Surveys*, 22(2), 83–121.
- Onvural, R., & Perros, H. (1986). On equivalencies of blocking mechanisms in queueing networks with blocking. *Operations Research Letters*, 5(6), 293–297.
- Osorio, C., & Bierlaire, M. (2009). An analytic finite capacity queueing network model capturing the propagation of congestion and blocking. *European Journal of Operational Research*, 196(3), 996–1007.
- Perros, H., Nilsson, A., & Liu, Y. (1988). Approximate analysis of product-form type queueing networks with blocking and deadlock. *Performance Evaluation*, 8(1), 19–39.
- Reveliotis, S. (2015a). Coordinating autonomy. *IEEE Robotics & Automation Magazine*, 22(2), 77–94.
- Reveliotis, S. (2015b). Real-time management of complex resource allocation systems: Necessity, achievements and further challenges. *IFAC-PapersOnLine*, 48(7), 50–57.
- Schmidt, L. C., & Jackman, J. (2000). Modeling recirculating conveyors with blocking. *European Journal of Operational Research*, 124(2), 422–436.
- Stewart, W. (2009). *Probability, markov chains, queues, and simulation*. Princeton University Press.
- Takahashi, Y., Miyahara, H., & Hasegawa, T. (1980). An approximation method for open restricted queueing networks. *Operations Research*, 28(3), 594–602.
- The Cw library developers (2017). Cw: 1.0.0. <https://doi.org/10.5281/zenodo.461865>.
- The Python Software Foundation (2015). Python 3.5.1. www.python.org.
- Venkatesh, S., & Smith, J. (2003). A graph-theoretic, linear-time scheme to detect and resolve deadlocks in flexible manufacturing cells. *Journal of Manufacturing Systems*, 22(3), 220–238.
- Venkatesh, S., & Smith, J. (2005). An evaluation of deadlock-handling strategies in semiconductor cluster tools. *IEEE Transactions on Semiconductor Manufacturing*, 18(1), 197–201.
- Venkatesh, S., Smith, J., Deuermeier, B., & Curry, G. (1998). Deadlock detection and resolution for discrete-event simulation: multiple-unit seizures. *IIE Transactions*, 30(3), 201–206.
- Vis, I. F. A. (2006). Survey of research in the design and control of automated guided vehicle systems. *European Journal of Operational Research*, 170(3), 677–709.
- Viswanadham, N., Narahari, Y., & Johnson, T. (1990). Deadlock prevention and deadlock avoidance in flexible manufacturing systems using petri net models. *IEEE Transactions on Robotics and Automation*, 6(6), 713–723.
- Walt, S. v. d., Colbert, S. C., & Varoquaux, G. (2011). The numpy array: A structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2), 22–30.
- Wilson, R. (1970). *An introduction to graph theory*. Pearson Education India.