

Year 1 Progress Monitoring Report

Geraint Palmer

Supervisors: Professor Paul Harper & Dr. Vincent Knight

June 16, 2015

1 Project Overview & Plan

This PhD is supervised by Professor Paul Harper and Dr Vincent Knight, and is match-funded between the University and Aneurin Bevan University Health Board (ABUHB). The research will investigate the workforce needs required to care for the frail and the elderly within the ABUHB. A systems-wide view will be considered to account for the complexity and connectedness of the healthcare system. Stochastic aspects will need to be captured, such as the variability in patient pathways and resource consumption.

The whole system will be modelled as a queueing network, with nodes representing departments in order to capture general patient flows. The model parameters will be obtained by exploring and analysing data acquired from the SAIL data bank. An agent-based simulation model will be created from this model, populated with the obtained data, where further exploration can be undertaken.

Frail and elderly patients are currently a particular priority for many healthcare planners and managers. An aging population means an increasing amount of elderly people are accessing healthcare service. Frail patients tend to have prolonged hospital length of stay partly due to insufficient care at home or in the community. Patients who remain in hospitals when they could otherwise have returned home block beds, increasing congestion and pressure on other areas of the healthcare system. This becomes a strain on the workforce, both in hospitals and community care teams. An understanding of the situation could alleviate this pressure.

The next subsection will discuss how patient flows through the healthcare system will be analysed.

1.1 Patient Flows

Queueing theory, and in particular queueing networks are techniques used to model the stochastic flow of customers around a system of service nodes. Note that throughout this report the terms node, service centre, service node and service station will be used interchangeably.

These concepts have been successfully adapted to healthcare systems, both in detailed individual level systems, e.g. [1], [10] and in examples with coarser granularity, for example [17].

An example of the sort of queueing network and patient flows that will be modelled is shown in Figure 1.

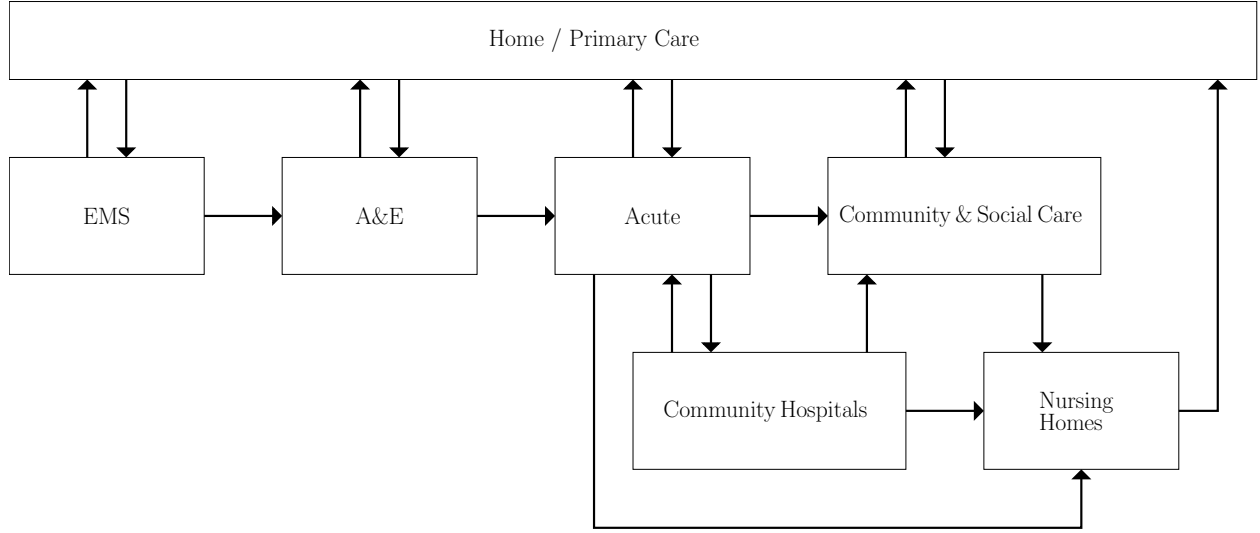


Figure 1: Patient flows around a healthcare system.

From exploring data and discussing with professionals at the ABUHB, the general architecture of the health system can be constructed, and patient flows built up. This will then be modelled as an open queueing network. Using results from the theory of Jackson networks [15], performance measures and steady-state probabilities can be obtained. This model will then be extended to include blocking between nodes so that the model represents reality, in which patients may be unable to progress to their next destination due to lack of capacity. Again performance measures and steady-state probabilities will be obtained.

A simulation framework has been written in Python (<http://www.python.org/>), analytic results will be compared to results obtained through simulation. A user friendly Django application of this simulation will be built which will feature simple parameter inputs and graphical outputs.

An extension of analysing current and potential patient flows is to derive optimal patient flows, which will be discussed in the next subsection.

1.2 Optimal Routing of Patients

One interesting avenue to explore is the optimal routing of patients through the system. That is, when there is a choice of where to send a patient, what is the optimal action to take? This question will be investigated using reinforcement learning algorithms.

Optimal analytical approaches will be investigated as well as agent based models. One particular implementation that will be used is **reinforcement learning**.

Reinforcement learning is a machine learning technique in which an agent learns through exploration of their environment, by completing actions and subsequently receiving rewards or punishments. Overviews

are given in [26] and [27]. The q -learning algorithm will be used to find the optimal action to take given a specific state.

In order to implement this for the queueing network simulation model:

- nodes or service centres will become the agents
- the state will be the number of patients at that node
- the reward will be some function of expected wait

Two potentially interesting areas to investigate are the effect of social vs selfish rewards, and how much each agent knows about the system. Social vs selfish rewards corresponds to node receiving rewards according to the system's state vs receiving rewards from just its own state. Agents could be observing the whole system's state or only observing its own state.

An extension of the above problem could be finding the optimal routing of patients when the agents are blind to their state, that is when they are not aware of the system state. Some results from game theory, and techniques such as genetic algorithms and linear programming will be explored.

The purpose of investigating patient flows in this project is to plan and supply the workforce appropriately, which will be discussed in the next subsection.

1.3 Workforce Planning

The performance measures of the queueing network models will be used to compare the workforce needs and skill mix required to care for frail and elderly patients across the ABUHB. With the user-friendly application, future trends and potential scenarios may be explored, with their effect on workforce investigated.

Using methods such as acuity ratios and nurses-per-bed, workforce needs can be estimated over time and across seasons. Various workforce planning techniques are reviewed in [14]. The aim of workforce planning is to find the best courses of actions to take in order to reduce the workforce gap that will extend over time. It is also important to gain an understanding of how changes to patient flows will affect the required skills and workforce for frail and elderly patients.

The next section will give an overview of my progress so far.

2 Progress & Learning

The Gantt chart in Figure 2 shows what has been achieved since October 2014.

I have participated, or will be participating in the following activities:

- Oct 2014: Gave a talk on my MSc project at a SWORDS event, Cardiff.
- Feb 2015: Gave a talk at [Python Namibia 2015](#) at the University of Namibia, Windhoek.
- Mar 2015: NATCOR Stochastic Modelling course, Lancaster.

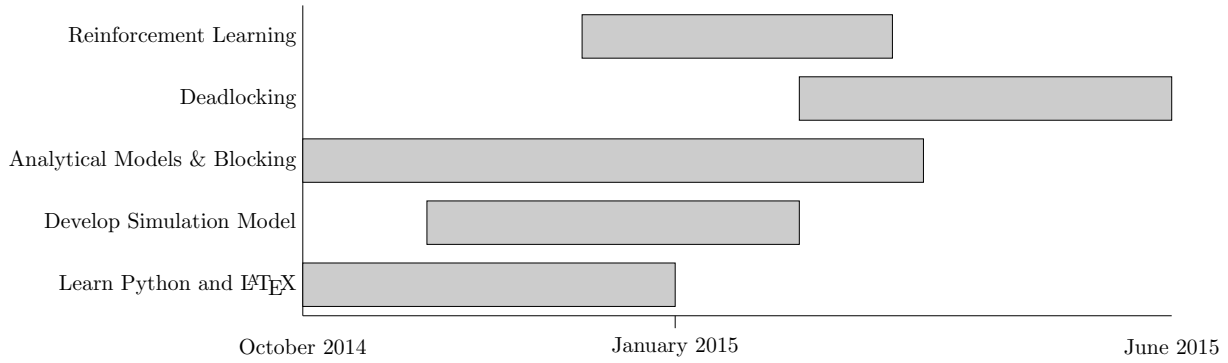


Figure 2: Gantt chart of research project to-date.

- Apr 2015: Data Mining course, Cardiff.
- May 2015: Gave a lightning talk on Linear Programming at the ‘[Developing Mathematical Models in Healthcare](#)’ seminar, Cardiff.
- May 2015: Gave a talk at the Wales Mathematics Colloquium, Gregynog.
- June 2015: Part of the organising committee of [DjangoCon Europe 2015](#), also gave a lightning talk at the event, Cardiff.
- June 2015: Will be part of the [1000 Lives National Learning Event](#) with the ABCi, Cardiff.
- July 2015: Will be giving a talk at [EURO 2015](#), the European Conference on Operational Research, Glasgow.
- Sept 2015: Am co-stream-organisers on the Stochastic Modelling stream of the [Young OR 19](#) conference, Birmingham.

In addition to the above, I have been involved in teaching, assessing and marking for a number of modules on the MSc programme including an object-orientated programming hackathon, the second year module MA0261 and the first year module MA1003. I have participated in two hackathon weekends in order to develop my programming skills, and am a regular tutor for maths and stats support.

The rest of this report is structured as follows: Section 3 will review the relevant literature, Section 4 will discuss methods of analytically investigating queueing networks, and Section 5 will discuss deadlocking properties of queueing networks.

3 Literature Review

This literature review is structured as follows: The first section gives an overview of how operational research is used in healthcare. The second, third and fourth sections give overviews on various aspects of queueing theory, including restricted queueing networks. The fifth section discusses how queueing networks have been used in healthcare. The sixth section introduces the concept of reinforcement learning.

3.1 Literature on Operational Research in Healthcare

Operational research techniques are well suited to be applied to the kinds of problems that arise in the health care system. This view is supported in [6]. This article describes the problems experienced in health care settings as mostly dealing with uncertainty or variability in the demand of services and resources, which may be analysed using standard operational research methods including queueing theory and simulation. Other health care issues lending themselves to OR listed include scheduling and capacity planning problems. The nature of healthcare systems is described in [12], with key features listed below:

- Complexity: for example different flow rules for different patients at different times
- Uncertainty: for example in demand, which can vary with month of the year, day of the week and even hour of the day
- Variability: for example the attributes and length of stay distributions patient classes can be vastly different
- Integration: and the all parts of the system are highly integrated with other parts of the hospital and other hospitals in the region

The paper builds a framework for building OR tools and models, which should be flexible and versatile to be applied to more than one specific area of the system, easy to use and able to answer "what if?" scenarios, integrated with all other parts of the system, and be validated with data.

Whole hospital and whole systems approaches are sometimes necessary when modelling and optimising healthcare systems in order to capture and investigate patient flows, up- and downstream resources, and identifying bottlenecks in the structures. A review of the operational research literature that involves multi-departmental analysis is presented in [30]. Here the authors conclude that there is an emphasis in the literature on the interaction between the wards, the emergency department and the surgery departments.

A whole system simulation is built in [31], bringing together eye clinics and social care to model elderly patients with age-related macular degeneration. This model is novel as it combines discrete event simulation, systems dynamics and agent-based modelling in one model in order to model patient flow through the eye clinic, patients' sight decay, and their social care respectively.

One OR technique that is used in healthcare is queueing networks, which will be discussed in the next subsection.

3.2 Queueing Networks

Queueing theory is an important branch of operational research, and the area of queueing networks have shown great promise in modelling computer and telecommunication systems and manufacturing procedures. An introduction to some of these results is given in [25]. The theory of a single queue has been widely investigated and developed, with Kendall's notation becoming the standard method of describing a queue. Kendall's notation describes a queue as $A/B/C/K/m/Z$: A is the arrival process, B the service process, C

the number of servers, K the maximum capacity of the queue, m the size of the population served, and Z the service discipline. Service disciplines may be varied, however the main ones include:

- First In First Out (FIFO),
- Last Come First Served (LCFS),
- First In Random Out (FIRO),
- Process Sharing (PS) in which each customer in the queue is served for a proportion of the time until their service time is complete.

Note that if the last three parameters are omitted, then it is assumed that $K = \infty$, $m = \infty$, and $Z = \text{FIFO}$.

Queueing networks comprise of a number of service centres with queues before them arranged into a network, with routing probabilities r_{ij} such that a customer finishing service at node i joins a queue at node j with probability r_{ij} . In [22] a history of the development of queueing networks is given, and an overview of the different types of networks studied.

The simplest, open networks are those that have at least one node to which customers arrive from the exterior, and at least one node from which customers leave to the exterior. Closed networks are self-contained networks where no customers arrive or depart to or from the exterior, and as a result always contain a constant total of N customers in the system.

A fundamental result for the study of queues in series is shown in both [7] and [23], although derived in different ways. The result, known as Burke's Theorem, states that the departure rate of an $M/M/c$ service station is equal to its arrival rate. This result is fundamental to the study of queues arranged in series or networks. The latter paper notes that this result holds for more general service disciplines too, including LCFS, FIRO and systems where the number of servers varies with the number of customers present.

The study of a system of queues in series, or routed into a network, is given in [15]. This paper investigates the simplest of cases, where all arrivals and service times are Markovian, the service discipline is FIFO, there is room for an infinite length queue between service centres, and the system is open. It is shown that in this case, each service centre, or node, behaves as an independent $M/M/c$ queue, which can be analysed as such. This approach is known as the decomposition method. The following product form solution is proved:

$$P(k_1, k_2, \dots, k_M) = \prod_{i=1}^M P_i(k_i) \quad (1)$$

where $P_i(k_i)$ is the probability of node i containing k_i customers, and $P(k_1, k_2, \dots, k_M)$ represents the probability of the system being in the state (k_1, k_2, \dots, k_M) , having k_j customers at node j for all j .

The basic model is extended in [16] to include different classes of customer, each with their own routing probabilities around the network and own service times at each node. The analysis concentrates on obtaining steady-state probabilities for the system state defined as $C = (c_1, c_2, \dots, c_J)$. Here for each node i , c_i is a specific ordering of classes of customer in the queue. A few extensions to this model are also presented, including changing the service discipline to PS and LCFS, prohibiting service from beginning until there is a

certain amount of customers in the queue, and making service effort rely on the current system state rather than the number of people currently in that queue.

This section will be expanded upon later. The next subsection will discuss more general feature in queueing theory that are useful in modelling.

3.3 General Features of Queueing Theory

More complicated customer behaviours have been studied in queueing theory, such as balking and reneging. Balking is defined as a customer arriving but deciding not to join the queue, therefore getting lost to the system. Reneging is when a customer entering a queue but then leaving before entering service. Queues with reneging are usually denoted with an extension to Kendall's notation, by $A/B/C/K/m/Z + D$, where D denotes the reneging method. In [2] and [3] the authors derive steady-state probabilities and mean values for an $M/M/1$ queue where customers exhibit balking and reneging. These papers give two separate customer behaviours for balking, first where customers balk with probability n/N when there are n customers in the system, $N - 1$ being an upper bound for queue size; and another where customers balk with probability $(1 - \beta/n)$ when there are n customers in the system, with $0 \leq \beta \leq 1$ being a measure of willingness to join the queue. In these papers a customer reneges after waiting time t with probability $\alpha e^{-\alpha t}$.

In the next subsection, queues with blocking properties will be discussed.

3.4 Queueing Networks with Blocking

Restricted queueing networks, or queueing networks with no or limited intermediate queues between service stations are more complicated to analyse. In these systems, if a customer finishes service at one node but is unable to join a queue at his destination node due to lack of queueing space, that customer remains in the current node, restricting his servers from starting the next customer's service. This is known as blocking. Since blocking introduces interdependencies between nodes, the product form solution of unrestricted networks is not appropriate. One of the first papers to consider these sorts of systems was [13]. Results are derived by writing out and solving the systems' difference equations. The same method is used in [5]. The thesis investigates two and three node systems, as well as systems with one service centre with infinite queue routing into a number of these two and three node systems.

A two node system with no intermediate queue and blocking is studied in [4]. In this paper the moment generating functions of waiting time and number of customers in the system are derived, from which further performance measures can be obtained.

Two features of blocking are described in [17]: (i) patients completing service at a blocked station remain there until there is sufficient queueing space at the next station, and (ii) these patients block other patients from entering that station. If a station only has characteristic (i) then it is referred to as 'classic congestion', and if a station has both characteristics it is referred to as 'blocking'. Three blocking situations are studied in [20], defined by rules on the system reaching 'full blocking' and their 'unblocking rule'. If the node that is subject to blocking has r parallel servers, then once r^* ($1 \leq r \leq r^*$) servers become blocked all remaining

unblocked servers stop service and the node becomes fully blocked. Given that there is a room for M customers to wait between nodes, once there are only k^* ($0 \leq k^* \leq M + r^* - 2$) customers waiting between the nodes all services may start again and the node becomes unblocked.

An approximation method for solving queueing networks with downstream blocking only is presented in [28]. The algorithm finds the mean values of a queueing network with feed forward flows and single server nodes. Iteratively working from the node furthest downstream and working backwards, if that station does exhibit blocking it finds the effective service time, that is the weighted sum of service time and the mean time blocked and waiting to transition to the next node, and computes the effective service time for the next upstream node with a recursive formula. This method is adapted to multi-server service stations with no intermediate queues in [17], and a similar iterative method was used in [18].

Restricted queueing networks can give rise to the phenomenon of deadlock []. In the simplest of cases, deadlock occurs when a customer finishes service at node i and is blocked from transitioning to node j ; however the individuals in node j are all blocked, directly or indirectly, by the blocked individual in node i . This can cause problems for both analytical and simulation models as in reality deadlock can be avoided or resolved quickly, and most of the literature on blocking conveniently assumes the networks are deadlock-free.

For closed networks of K customers with only one class of customer, [19] proves the following condition to ensures no deadlock: for each minimum cycle C , $K < \sum_{j \in C} B_j$, the total number of customers cannot exceed the total queueing capacity of each minimum subcycle of the network. The paper also presents algorithms for finding the minimum queueing space required to ensure deadlock never occurs, for closed cactus networks, where no two cycles have more than one node in common. This result is extended to multiple classes of customer in [21], with more restrictions such as single servers and each class having the same service time distribution. Here a integer linear program is formulated to find the minimum queueing space assignment that prevents deadlock. The literature does not discuss deadlock properties in open restricted queueing networks.

The next subsection will review how some modellers have applied the theory of queueing networks to a healthcare system.

3.5 Queueing Networks Applied to Healthcare

This section will give an overview of particular application of queueing network analysis to healthcare.

A geriatric ward is modelled as an $M/PH/c$ queue in [11], a multi server queue with phase-type distribution where no queue is allowed to build up, i.e. patients are sent elsewhere and lost to the system if all c beds are occupied. This model is used to find the optimal number of beds in order to reduce the cost of turning away patients and maintaining empty beds.

Queueing networks have been used to model healthcare systems. In [1] a health centre is modelled as an open queueing network with eight nodes representing five care-givers and three areas prior to seeing a doctor or nurse, in order to assess how to improve waiting times. The model proved successful and disproved a widely held belief that the front desk was the system bottleneck, concluding that including time spent prior to seeing a care-giver in appointment times would reduce waits. In [10] the orthopaedics department of

the Middelheim hospital in Antwerpen was modelled as an open queueing network with five service centres and eighteen patient classes. Both preemptive and nonpreemptive interruptions were incorporated into the service times of one node. Three analytical methods of finding flow times were compared: two formulas applied to the decomposition method, Kingman and Whitt, and then by using the Brownian model. These were compared to a simulation model of the system, where it was concluded that the decomposition method far outperformed the Brownian model, and the Kingman formula exerted slightly better results than the Whitt formula.

In [1] the authors list a number of reasons why queueing network models do not match real-life healthcare situations. Queueing network analysis generally focus on steady-state statistics, though in situations where service starts and needs each day the system is reset every day and rarely, if ever reach steady-state. Queueing analysis also assumed that arrival and service times are independent of one another, yet if an appointment system is used this is not the case. However, the model used yielded realistic results and some useful conclusions and recommendations. In [17] a queueing network with upstream blocking is used to model the flow of mental health patients through a care system in Philadelphia, where service stations are housing facilities: extended acute hospitals, residential facilities and supported housing. Its aim was to find an alternative approach to the needs assessment approach of capacity planning, however it was found that the queueing network model with static arrival rates could not effectively model the real-world as customer behaviour exhibited renegeing based on waiting times.

The next subsection will introduce reinforcement learning.

3.6 Reinforcement Learning

The concept of machines learning from experience stems from the first papers on computing, notably [29]. In this early paper the idea is discussed that actions chosen at random will be repeated more often if that action resulted in a reward, and less often if that action resulted in punishment. This forms the basis of reinforcement learning.

Comprehensive overviews are given in [26] and [27].

The next section will give overviews and examples on solving the steady-state probabilities and performance measures of queueing networks.

4 Analytical Models of Open Queueing Networks

4.1 Queueing Networks

Queueing networks are a set of service centers, with or without a queueing space before them. They are connected by possible routes such that when a customer finishes service at one service centre he or she joins the queue of another service centre according to a probability distribution. Throughout this report queueing networks will be represented as in Figure 3 unless stated otherwise.

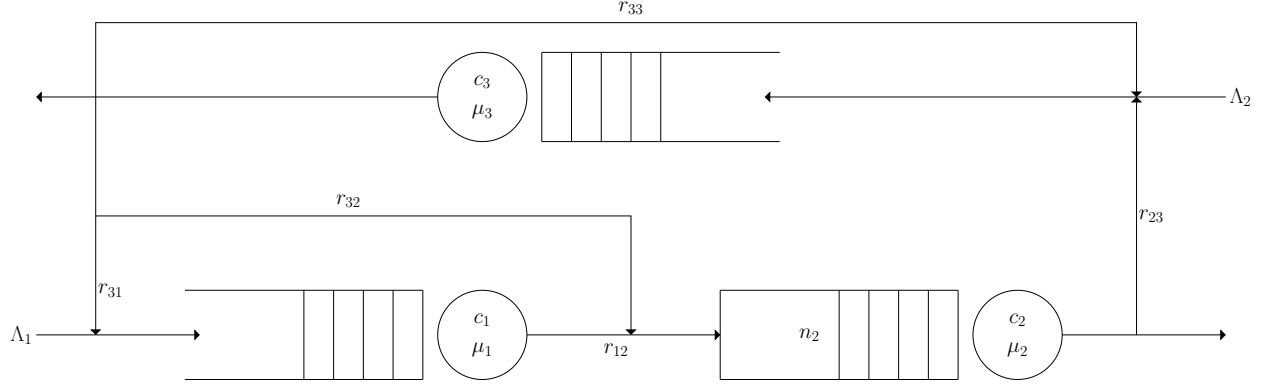


Figure 3: A queueing network

Here the service centres are shown as circles and all possible routings are shown as arrows. The external arrival rate to node i is denoted by Λ_i , there are c_i servers serving at a rate of μ_i at node i . Routing probabilities are described by the routing matrix R , with elements r_{ij} denoting the probability of transitioning from node i to node j once service is complete. The probability of leaving the system after service at node i is given by $1 - \sum_j r_{ij}$. The routing matrix for the network in Figure 3 is given by

$$R = \begin{pmatrix} 0 & r_{12} & 0 \\ 0 & 0 & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix}$$

Nodes 1 and 3 have open queues, denoted with open rectangles, and have infinite queueing capacity. Node 2 is restricted, and its queue can only accommodate n_2 waiting customers. If a customer finishes service and wants to transition to node 2, but node 2 has a full queue, then that customer remains at its current node until space becomes available at node 2. While remaining at its node, that customer blocks waiting customers from entering service until they have left that node.

As well as being classified as either restricted (containing at least one node with limited capacity) and unrestricted, queueing networks are classified as either open or closed. In an open queueing network customers can arrive externally to the network, and can leave the network once service is finished. In closed queueing networks, there is a fixed amount of customers in the network, no customers can enter the network from the outside, and no customers can leave the network.

This report will only look at open queueing networks.

The next subsection will explain one method of solving the steady-state probabilities of open unrestricted queueing networks.

4.2 Open Unrestricted Queueing Networks

In order to solve open unrestricted queueing networks, we make use of the following theorem due to [7] and [23]. Known as Burke's Theorem, it states that the departure rate of customers from a service station is equal to the arrival rate.

Therefore, each node's effective arrival rate can be calculated by solving the following traffic equations

$$\lambda_i = \Lambda_i + \sum_j r_{ji} \quad \forall i \quad (2)$$

We can now make use of Jackson's Theorem [15] in order to use the decomposition method to solve the steady-state probabilities of a queueing network. It states that, given a network with M service stations, each service station can be solved independently of the others. That is, if $P(k_1, k_2, \dots, k_M)$ is the probability of there being k_1 customers at node 1 and k_2 customers at node 2, etc., and $P_i(k_i)$ is the probability of there being k_i customers at node i , then the following is true

$$P(k_1, k_2, \dots, k_M) = \prod_{i=1}^M P_i(k_i) \quad (3)$$

From (3) finding the steady-state probabilities of the networks requires decomposing the network into individual nodes. Performance measures such as expected number of customers in the system and expected time spent in the system can straightforwardly be found in the same manner:

$$L = \sum_i L_i \quad (4)$$

$$W = \sum_i \frac{L_i}{\lambda_i} \quad (5)$$

The next subsection will give an example of this.

4.2.1 Example of Solving Using Jackson's Theorem

Consider the queueing network shown in Figure 4.

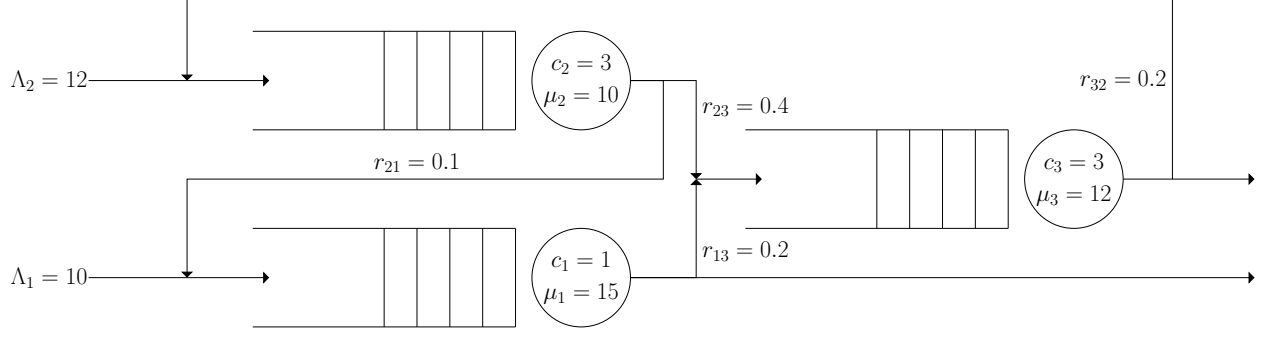


Figure 4: An open queueing network.

The traffic equations yield the following set of simultaneous equations:

$$\begin{aligned}\lambda_1 &= 10 + 0.1\lambda_2 \\ \lambda_2 &= 12 + 0.2\lambda_3 \\ \lambda_3 &= 0.2\lambda_1 + 0.4\lambda_2\end{aligned}$$

Solving the above equations gives following effective arrival rates to each node, $\lambda_1 \approx 11.354$, $\lambda_2 \approx 13.537$, $\lambda_3 \approx 7.686$. Then, using the following equations from standard queueing theory, each node can be analysed independently.

$$P(0) = \left[\sum_{j=0}^{c-1} \frac{(\lambda/\mu)^j}{j!} + \frac{(\lambda/\mu)^c}{c! (1 - (\lambda/\mu c))} \right]^{-1} \quad (6)$$

$$P(j) = \begin{cases} \frac{(\lambda/\mu)^j}{j!} P(0) & \text{for } j < c \\ \frac{(\lambda/\mu)^j}{c! c^{j-c}} P(0) & \text{for } j \geq c \end{cases} \quad (7)$$

Now, if we wish to find $P(3, 1, 1)$, the probability that node 1 has 3 customers, node 2 has 1 customer and node 3 has 1 customer, by Jackson's Theorem we can look at each node separately. The above equations yield $P_1(3) \approx 0.10541$, $P_2(1) \approx 0.22172$ and $P_3(1) \approx 0.29827$. And so:

$$\begin{aligned}P(k_1 = 3, k_2 = 1, k_3 = 1) &= P(k_1 = 3)P(k_2 = 1)P(k_3 = 1) \\ &\approx 0.10541 \times 0.22172 \times 0.29827 \\ &\approx 0.00697\end{aligned}$$

Similarly, we can find the expected number of customers in the system. Standard queueing theory give the

following formula:

$$L = \left[\frac{(\lambda/\mu)^c \mu}{(c-1)!(c\mu - \lambda)^2} \right] P(0) + \frac{\lambda}{\mu} \quad (8)$$

For each node we get $L_1 = 0.96453$, $L_2 = 1.36482$ and $L_3 = 0.64097$. From which we get:

$$\begin{aligned} L &= L_1 + L_2 + L_3 \\ &\approx 0.96453 + 1.36482 + 0.64097 \\ &\approx 2.97034 \end{aligned}$$

The next subsection will discuss methods of deriving steady-state probabilities of open restricted queueing networks.

4.3 Restricted Queueing Networks

Restricted queueing networks are those where at least one node has limited capacity and blocking rules apply. These networks are not as simple to solve as open unrestricted queueing networks, as Burke's Theorem no longer applies. These become even more difficult to solve when there exist feedback loops in the networks. We will concentrate on restricted networks with feedforward flows only.

4.3.1 Using Differential Difference

This method was used in [5] to investigate the behaviour of queues with blocking. In this method equations are derived for every state of the queueing network, and so networks with very many states become very difficult to solve. This method is entirely equivalent to the method described in the next subsection, Markov chains. In this case, we are modelling the queue as a Markov chain and solving from first principals, rather than manipulating transition matrices in the next subsection.

In this example, a network with only two single-server service centres, and no queueing space, is solved. This queueing network is shown in Figure 5.

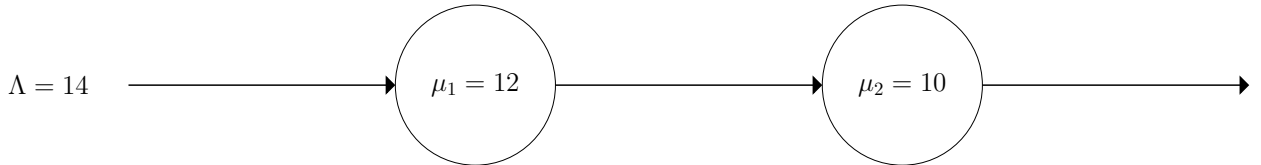


Figure 5: A two-node queueing network with no queueing capacities.

Define $P_{i,j}(t)$ as there being i customers in the first node, j customers in the second node at time t . Define

$P_{1b,1}(t)$ as there being 1 blocked customer in the first node, 1 customer in the second node at time t . Then we get the following difference equations:

$$\begin{aligned}
P_{0,0}(t + \delta t) &= (1 - \Lambda\delta t)P_{0,0}(t) + (1 - \Lambda\delta t)\mu_2\delta tP_{0,1}(t) \\
P_{1,0}(t + \delta t) &= \Lambda\delta tP_{0,0}(t) + (1 - \mu_1\delta t)P_{1,0}(t) + (1 - \mu_1\delta t)\mu_2\delta tP_{1,1}(t) \\
P_{0,1}(t + \delta t) &= (1 - \Lambda\delta t)(1 - \mu_2\delta t)P_{0,1}(t) + \mu_1\delta tP_{1,0}(t) + \mu_2\delta tP_{1b,1}(t) \\
P_{1,1}(t + \delta t) &= (1 - \mu_1\delta t)(1 - \mu_2\delta t)P_{1,1}(t) + \Lambda\delta t(1 - \mu_2\delta t)P_{0,1}(t) \\
P_{1b,1}(t + \delta t) &= (1 - \mu_2\delta t)P_{1b,1}(t) + \mu_1\delta t(1 - \mu_2\delta t)P_{1,1}(t)
\end{aligned}$$

Differentiating leaves:

$$\begin{aligned}
\frac{dP_{0,0}(t)}{dt} &= -\Lambda P_{0,0}(t) + \mu_2 P_{0,1}(t) \\
\frac{dP_{1,0}(t)}{dt} &= \Lambda P_{0,0}(t) - \mu_1 P_{1,0}(t) + \mu_2 P_{1,1}(t) \\
\frac{dP_{0,1}(t)}{dt} &= -\mu_2 P_{0,1}(t) - \Lambda P_{0,1}(t) + \mu_1 P_{1,0}(t) + \mu_2 P_{1b,1}(t) \\
\frac{dP_{1,1}(t)}{dt} &= -\mu_2 P_{1,1}(t) - \mu_1 P_{1,1}(t) + \Lambda P_{0,1}(t) \\
\frac{dP_{1b,1}(t)}{dt} &= -\mu_2 P_{1b,1}(t) + \mu_1 P_{1,1}(t)
\end{aligned}$$

We are interested in steady-state probabilities. Steady states occur when the differentials are set to 0:

$$\begin{aligned}
\Lambda P_{0,0} &= \mu_2 P_{0,1} \\
\mu_1 P_{1,0} &= \Lambda P_{0,0} + \mu_2 P_{1,1} \\
(\Lambda + \mu_2) P_{0,1} &= \mu_1 P_{1,0} + \mu_2 P_{1b,1} \\
(\mu_1 + \mu_2) P_{1,1} &= \Lambda P_{0,1} \\
\mu_2 P_{1b,1} &= \mu_1 P_{1,1}
\end{aligned}$$

We also have that all probabilities must sum to 1:

$$\sum_i \sum_j P_{i,j} = 1$$

This set of equations can be solved to give the following expressions for the desired probabilities:

$$\begin{aligned}
P_{0,0} &= \frac{\mu_1 \mu_2^2 (\mu_1 + \mu_2)}{\mu_2^3 (\Lambda + \mu_1) + \mu_2^2 (\Lambda + \mu_1)^2 + \mu_2 (\Lambda + \mu_1) (\Lambda \mu_1) + \Lambda^2 \mu_1^2} \\
P_{0,1} &= \frac{\Lambda \mu_1 \mu_2 (\mu_1 + \mu_2)}{\mu_2^3 (\Lambda + \mu_1) + \mu_2^2 (\Lambda + \mu_1)^2 + \mu_2 (\Lambda + \mu_1) (\Lambda \mu_1) + \Lambda^2 \mu_1^2} \\
P_{1,1} &= \frac{\Lambda^2 \mu_1 \mu_2}{\mu_2^3 (\Lambda + \mu_1) + \mu_2^2 (\Lambda + \mu_1)^2 + \mu_2 (\Lambda + \mu_1) (\Lambda \mu_1) + \Lambda^2 \mu_1^2} \\
P_{1b,1} &= \frac{\Lambda^2 \mu_1^2}{\mu_2^3 (\Lambda + \mu_1) + \mu_2^2 (\Lambda + \mu_1)^2 + \mu_2 (\Lambda + \mu_1) (\Lambda \mu_1) + \Lambda^2 \mu_1^2} \\
P_{1,0} &= \frac{\Lambda (\Lambda + \mu_1 + \mu_2) \mu_2^2}{\mu_2^3 (\Lambda + \mu_1) + \mu_2^2 (\Lambda + \mu_1)^2 + \mu_2 (\Lambda + \mu_1) (\Lambda \mu_1) + \Lambda^2 \mu_1^2}
\end{aligned}$$

Substituting in the problem's values given in Figure 5, $\Lambda = 14$, $\mu_1 = 12$ and $\mu_2 = 10$, the probabilities become:

$$P_{0,0} \approx 0.15951$$

$$P_{0,1} \approx 0.22332$$

$$P_{1,1} \approx 0.14211$$

$$P_{1b,1} \approx 0.17053$$

$$P_{1,0} \approx 0.30452$$

4.3.2 Using Markov Chain Transition Matrices

Queueing networks with larger state-spaces can be solved numerically. The computations of the transition matrices and matrix inversions are all done using Python. Note again that this is the same model used in the previous subsection, however here we are making use of matrix algebra rather than solving from first principals

The queueing network solved using the difference equations method can be generalised to include finite queues before each server. This network is shown in Figure 6.

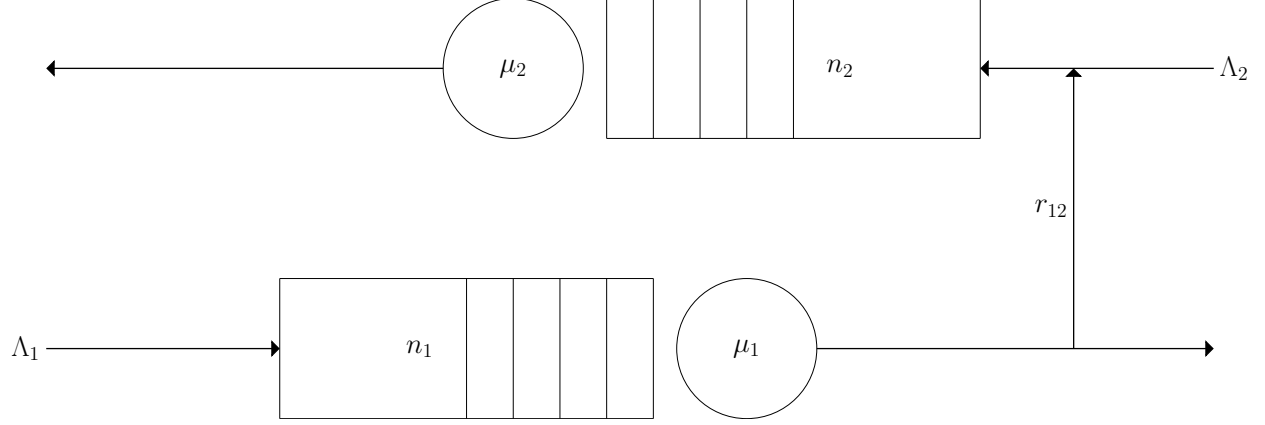


Figure 6: A two-node queueing network with finite queues.

Substituting in $n_1 = 0$, $n_2 = 0$, $r_{12} = 1$ and $\Lambda_2 = 0$ the problem is equivalent to that solved using difference equations, and identical numerical solutions are found.

The state space for this network is defined as follows:

$$S = \{(i, j) \in \mathbb{N}^{(n_1+1 \times n_2+2)} | 0 \leq i + j \leq n_1 + n_2 + 2\}$$

where i denotes the number of customers waiting or in service at node 1, and j denotes the number of customers waiting or in service at node 2 plus the number of customers at node 1 blocked from entering node 2. The transition rates $q_{(i_1, j_1), (i_2, j_2)}$, the rate at which the system transitions from state (i_1, j_2) to state (i_2, j_2) , are defined as follows:

$$q_{(i_1, j_1), (i_2, j_2)} = \begin{cases} \left. \begin{array}{l} \Lambda_1 \text{ if } i_1 \leq n_1 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (1, 0) \\ \left. \begin{array}{l} \Lambda_2 \text{ if } j_1 \leq n_2 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (0, 1) \\ \left. \begin{array}{l} (1 - r_{12})\mu_1 \text{ if } j_1 \leq n_1 + 1 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (-1, 0) \\ \left. \begin{array}{l} \mu_2 \text{ if } \delta = (0, -1) \\ r_{12}\mu_1 \text{ if } j_1 \leq n_2 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (-1, 1) \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

where $\delta = (i_2 - i_1, j_2 - j_1)$.

The transition matrix is constructed numerically using Python, discretized, and then solved by finding the eigenvectors corresponding to the eigenvalue 1.

This method becomes more difficult when some infinite queues are allowed, as the state-space becomes

infinite. Truncating the state-space is one way of overcoming this.

The representation of the states will have to be altered to accommodate more nodes, more routing possibilities and multi-server queues.

An approximation method for restricted networks with downstream flows is discussed in the next section.

4.3.3 An Approximate Approach

An approximation method was developed by [28], and developed and implemented in [17] for queueing networks with blocking with no intermediate queues. This approximation makes use of the effective service time, illustrated in Figure 7. The effective service time at node i , denoted $1/\tilde{\mu}_i$, is the time spent from entering service at node i to departing that node. It includes the time spent blocked at that node before being able to move on to another node.

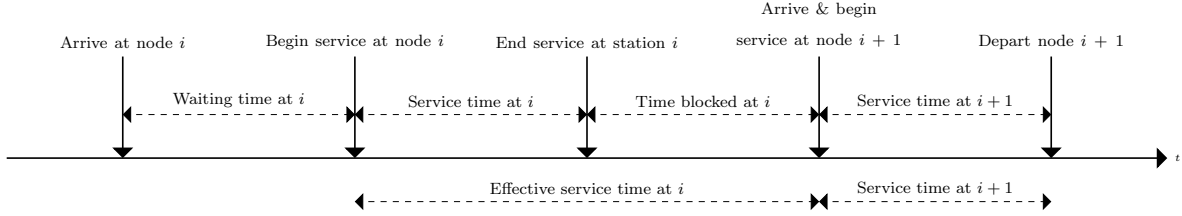


Figure 7: Effective service time.

All service stations directly upstream from a service station that has limited capacity will have an effective service rate. The effective service rate is calculated as follows:

$$\frac{1}{\tilde{\mu}_i} = \left(1 - \sum_j r_{ij}\right) \left(\frac{1}{\mu_i}\right) + \sum_j r_{ij} \left(\frac{1}{\mu_i} + W_j\right) \quad (10)$$

The number of customers at each node, can then be solved recursively, beginning at the most downstream nodes and working upstream, using the following equations:

$$L_i = \left[\sum_{n=0}^{c_i-1} \frac{\omega_i^n}{n!} + \frac{\omega_i^{c_i}}{(1-\rho_i)c_i!} \right]^{-1} \frac{\rho_i \omega_i^{c_i}}{(1-\rho_i)^2 c_i!} \quad (11)$$

$$W_i = \frac{L_i}{\lambda_i} \quad (12)$$

where $\omega_i = \frac{\lambda_i}{\mu_i}$, $\rho_i = \frac{\omega_i}{c_i}$. The service rate μ_i should be replaced with the effective service rate $\tilde{\mu}_i$ where appropriate, and effective arrival rates are computed by solving the traffic equations as before.

As an example, we will approximate the expected number of customers, and expected wait, at each node in the queueing network shown in Figure 8.

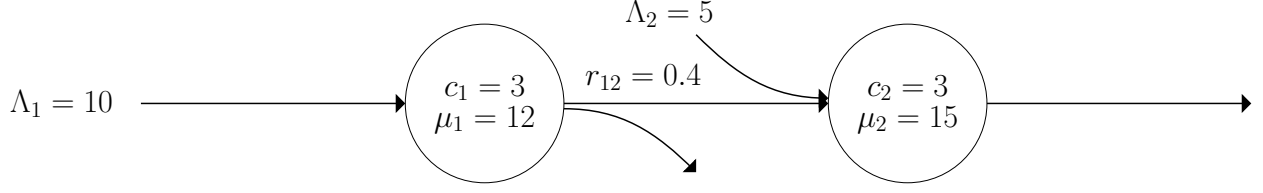


Figure 8: A two node queueing network with no intermediate queues.

First, we solve the traffic equations in order to obtain the effective arrival rates.

$$\lambda_1 = \Lambda_1$$

$$\lambda_2 = \Lambda_2 + r_{12}\lambda_1$$

yielding $\lambda_1 = 10$ and $\lambda_2 = 9$.

We begin at node 2, as it is furthest downstream, it has no immediate downstream nodes. Therefore its effective service rate is the same as its service rate. By equation (11) we obtain $L_2 \approx 0.0061643$, and by Little's Law we get $W_2 \approx 0.0006849$.

We now have all we require to find the effective service rate of node 1.

$$\begin{aligned} \frac{1}{\tilde{\mu}_1} &= \left(1 - \sum_j r_{1j}\right) \left(\frac{1}{\mu_1}\right) + \sum_j r_{1j} \left(\frac{1}{\mu_j} + W_j\right) \\ &\approx \left(0.6 \times \frac{1}{12}\right) + \left(0.4 \times \left(\frac{1}{15} + 0.0006849\right)\right) \\ &\approx 0.0836073 \end{aligned}$$

giving $\tilde{\mu}_1 \approx 11.9606772$. Whence we get $L_1 \approx 0.0224831$ and $W_1 \approx 0.0022483$.

The next subsection will discuss an alternative to analytical methods, simulation.

4.3.4 Simulation

The methods described in the report can only be applied to restricted networks with downstream flows, that is, no networks with cycles. This is due to the breakdown of Markovian arrivals and deadlock possibilities. In order to investigate networks with these features computer simulation will have to be used.

A simulation framework has been developed in Python that handles restricted and unrestricted open queueing networks of any architecture with multiple classes of patients and a variety of service distributions. The framework uses the three-phase-simulation-approach described in [24]. This is shown in Figure 9.

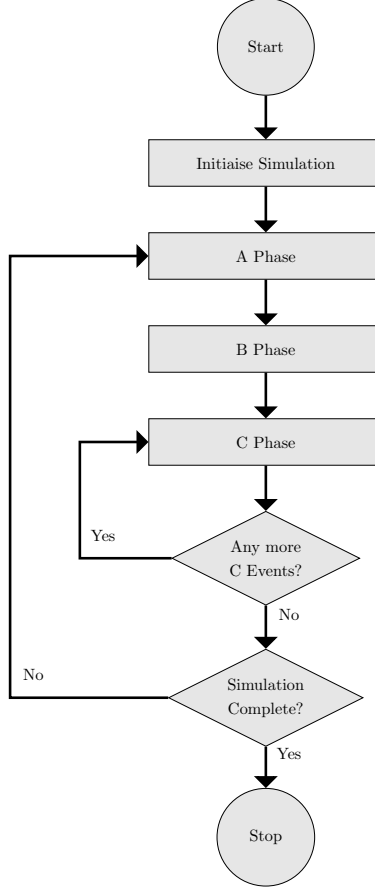


Figure 9: The Three Phase Simulation Approach [24]

In the A-phase the simulation clock is moved to the time of the next event. In the B-phase all scheduled events for that time are carried out. In the C-phase we check whether some of the events just carried out has caused any more conditional events to need to occur, and these are carried out. At the decision node ‘Any more C?’, we check whether the C-events that just occurred has caused any more C events to need to occur.

In the simulation framework built, B-events correspond to either an external arrival, or a customer finishing service. C-events correspond to a customer starting service, or a customer being released from a node.

Table 1 below shows some of the possible conditional events that arise from other events.

<i>Event</i>	<i>Conditional Event</i>	<i>Reason</i>
External arrival	Nothing	The customer arrives and needs to queue before starting service.
External arrival	Start service	The customer arrives and can begin service immediately as there is no queue.
Finish service	Release/Move customer	A customer finishes service, there is sufficient queueing space at the next node, and so moves.
Finish service	Nothing	A customer finishes service, but there is insufficient queueing space at the next node, and so the customer becomes blocked.
Release/Move customer	Start service	A customer moves to the next node and can begin service immediately as there is no queue.
Release/Move customer	Start service	A customer moves from the current node, allowing a customer at the front of the queue to start service.
Release/Move customer	Release/Move customer	A customer moves from the current node, allowing enough queueing space from a blocked customer to now enter that node.
Release/Move customer	Nothing	A customer moves to the next node where there is a queue, and so service does not start immediately. There is no one in the current node's queue to start service.
Release/Move customer	Nothing	A customer exits the system, and there is no queue at the current node and so no-one starts service.

Table 1: Table listing possible conditional event triggers.

Figure 10 illustrates the general structure of the objects and classes of the simulation code.

There is one **Simulation** object. This contains all the global variables used, information about the queueing network itself, and the methods that run the simulation and write data to file. It also contains one **ArrivalNode**, one **ExitNode**, and as many **Nodes** the queueing network requires.

The **ArrivalNode** object spawns new **Individuals**, that is customers arriving externally. They never remain here, immediately being released to the relevant node when spawned, or lost to the system if there is not enough queueing capacity at that node.

The **ExitNode** is a dummy node, that simply acts as a bucket to send and store **Individuals** that leave the system.

The **Nodes** are where **Individuals** queue for service, receive service, and remain when blocked. Nodes contain all the methods required to give **Individuals** `arrival_date`, `service_start_date`, `service_end_date`, `service_time` and `exit_date`. This object also contains all methods to move **Individuals** from **Node** to

Node, block **Individuals**, and write the **DataRecords** of **Individuals**.

The **Individual** objects are the entities that are traversing the queueing network, waiting, being served, getting blocked, transitioning from node to node, and being served. They contain a store of every **DataRecord** that **Individual** has had written.

Each **DataRecord** contains all the information about one **Individual's** visit to one **Node**.

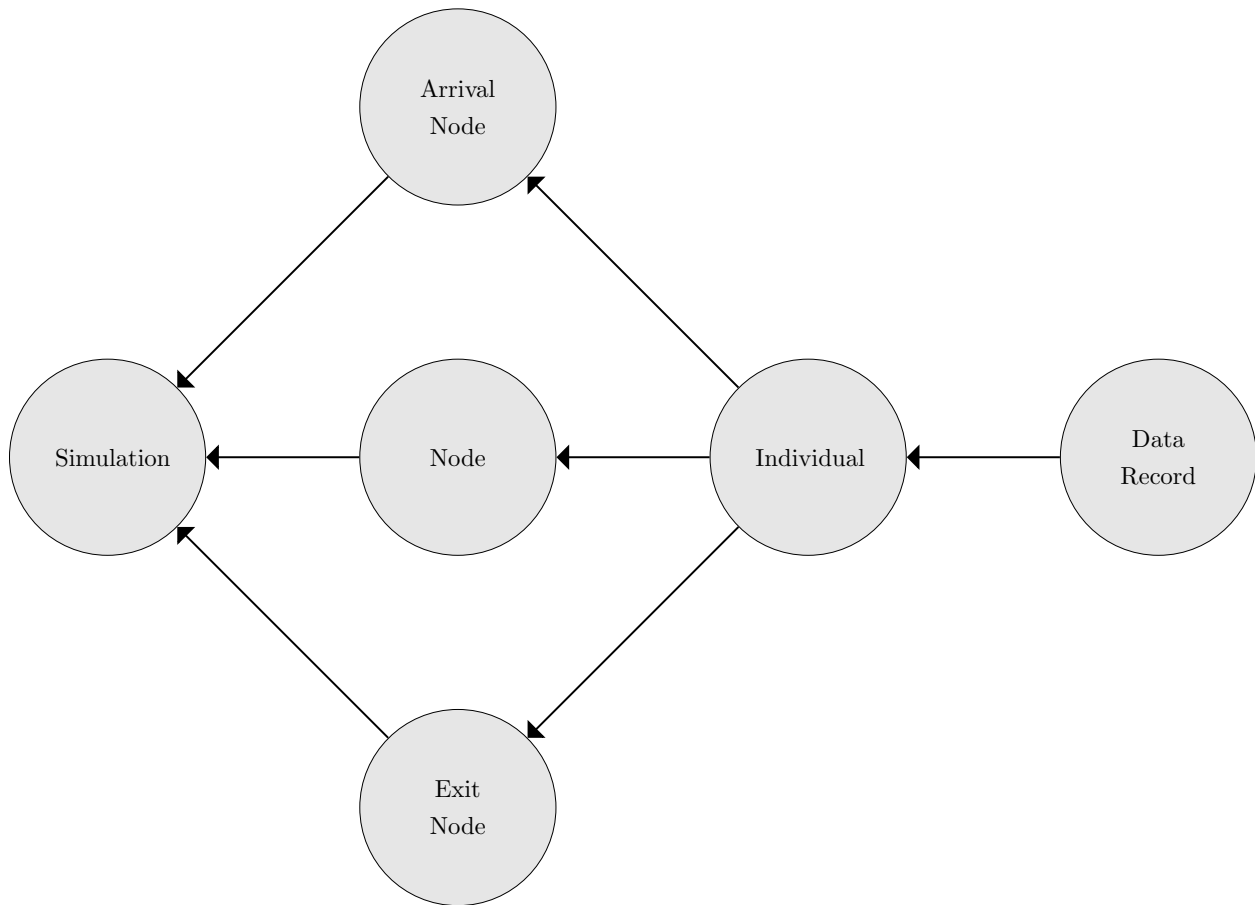


Figure 10: Structure of objects and classes of the simulation code.

The following pseudocode describes the main loop of the simulation:

```
Find next active node
while current time < max time do
    Next active node has event
    for all nodes do
        | update next event date
    end

    next active node  $\leftarrow$  node with closest next event date
    current time  $\leftarrow$  next event date of next active node
end
```

Here most the simulation process is carried out in the command ‘have event’. If the next active node is an **ArrivalNode** then the following pseudocode describes the method:

```
Create a new Individual()
Determine this individual’s next node
if enough capacity at next node then
    | Join queue at next node
else
    | The individual is lost to the system
end
```

ExitNode will never have an event, as its next event date is always set to infinity. For all other **Nodes**, the ‘have event’ command is described in the pseudocode below:

```
next individual  $\leftarrow$  individual that is finishing service
Determine next individual’s next node
if enough capacity at next node then
    | release next individual
else
    | block next individual
end
```

where the command release acts as follows:

```

Move individual to next node
if queue length at next node = 0 then
    | start service
end

if any individuals need to join this queue then
    | release blocked individual
end

if individual waiting to start service at this node then
    | start service
end

```

Notice that the release command is recursive, and so the process of unblocking individuals can continue until no more individuals can be unblocked.

5 Deadlocking Properties of Open Restricted Queueing Networks

This section will define and discuss the properties and detection of deadlock in queueing networks. Throughout the section, when discussing queueing networks, it is assumed that the queueing network is open and connected. Open queueing networks are those networks that have at least one node to which customers arrive from the exterior, and at least one node from which customers leave to the exterior.

5.1 Deadlock

Definition 1. *When a simulation is in a situation where at least one service station, despite having arrivals, ceases to finish any more services due to recursive upstream blocking the system is said to be in deadlock.*

Deadlock can be experienced in any open queueing network that experiences blocking, with at least once cycle containing all service stations with restricted queueing capacity.

Deadlock occurs when a customer finishes service at node i and is blocked from transitioning to node j ; however the individuals in node j are all blocked, directly or indirectly, by the blocked individual in node i . That is, deadlock occurs if every individual blocking individual X , directly or indirectly, are also blocked.

In Figure 11 a simple two node queueing network is shown in a deadlocked state. Customer occupying server A_1 has finished service at node A , but remains there as there is not enough queueing space at node B to accept them. We say the customer at server A_1 is blocked by the customer at server B_1 , as he is waiting for that customer to be released. Similarly, the customer occupying server B_1 has finished service at node

B , but remains there as there is not enough queueing space at node A , and so the customer at server B_1 is blocked by the customer at server A_1 .

When there are multiple servers, individuals become blocked by all customers in service at the destination service station. Figure 12a shows two nodes in deadlock, the customer occupying server A_1 is blocked by customers at both B_1 and B_2 , who are both blocked by the customer at A_1 . However in 12b, customer at A_1 is blocked by customers at both B_1 and B_2 , but the customer at B_2 isn't blocked, and so there is no deadlock.

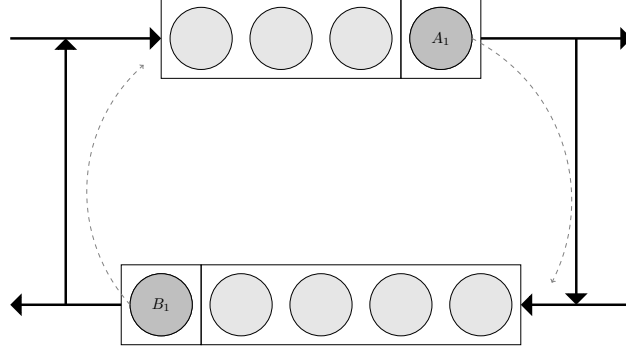


Figure 11: Two nodes in deadlock.

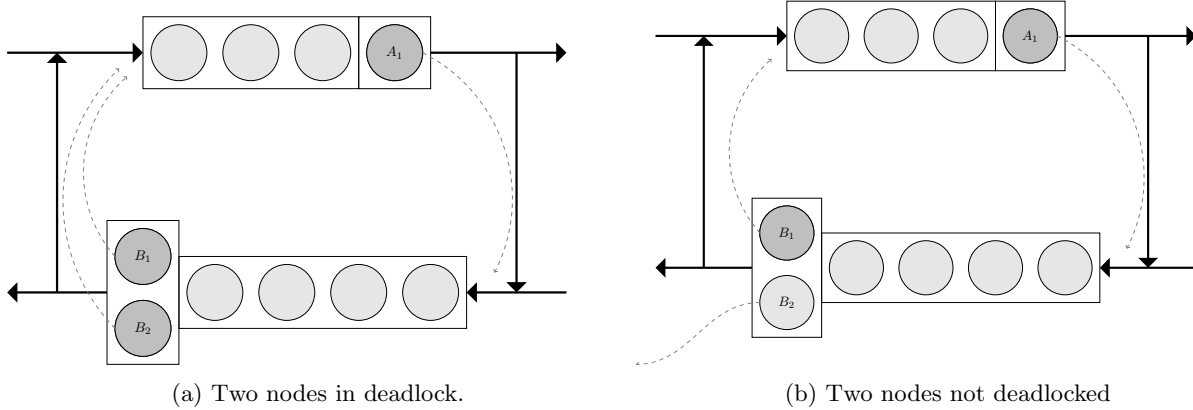


Figure 12: Two nodes: a) in deadlock and b) not in deadlock.

Note that the whole queueing network need not be deadlocked, only a part of it. If one section of the network is in deadlock, then the system is deadlocked, even though customers may still be able to have services and transitions in other areas of the network. An example is shown in Figure 13a. This idea is expanded on in the next subsection.

5.2 Types of Deadlock

The previous subsection introduced the idea that parts of a queueing network can be in a deadlocked state, although other parts will continue to flow. The different configurations of which nodes experience deadlock can be thought of as different types of deadlock. Each different type of deadlocked state can be denoted (i_1, i_2, i_3, \dots) where $i_k = -1$ if node k node is participating in that deadlocked state. The amount of different types of deadlock that a queueing network can experience is equal to the number of directed cycles in the queueing network's routing matrix.

For connected queueing networks, these deadlocks can be classified into transient deadlocked states and the absorbing deadlocked state.

Definition 2. *A transient deadlock state is when there are still some changes of state whilst a subgraph of the queueing network is itself in deadlock.*

Definition 3. *The absorbing deadlock state is when all subgraphs of the queueing network are in deadlock.*

Figure 13 shows a three nodes network in a transient deadlocked state, and an absorbing deadlocked state. In Figure 13a the occupants of servers B_1 and B_2 are blocked from entering node A ; and the occupant of server A_1 is blocked from entering node B , and so these two nodes are in deadlock. However, node C can continue with regular services, until the occupants of every server of C attempt to join a deadlocked node. At which point, the whole system is deadlocked, and so has reached absorbing deadlock, show in Figure 13b.

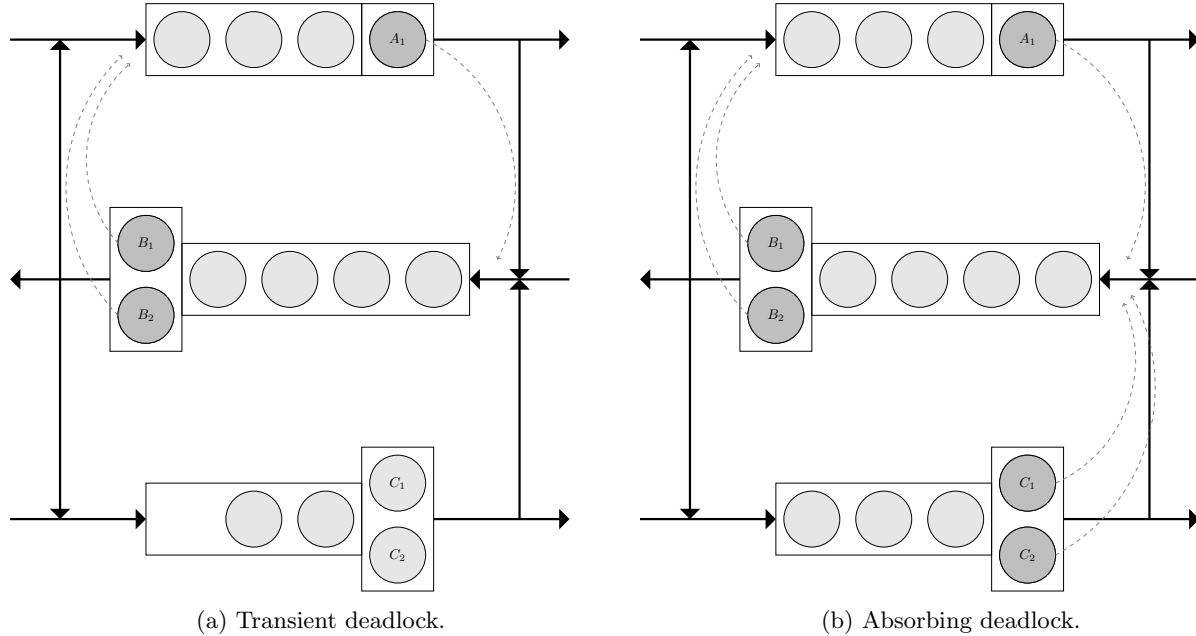


Figure 13: Types of deadlock: a) transient and b) absorbing.

For a queueing network Q with N service stations, the absorbing deadlocked state corresponds to $(i_1 = -1, \dots, i_N = -1)$, the state where all service stations experience deadlock. It should be clear that if the queueing network is connected, then there is a non-zero probability that once one part of the network is in deadlock, the whole system will fall into a deadlocked state, simply by the individuals in the non-deadlocked nodes attempting to transition into a deadlocked node. That is, once Q falls into one of the transient deadlocked states, it will eventually transition, either directly or through other transient deadlocked states, into the absorbing deadlocked state.

If the routing matrix of Q is complete, that is there is a possible route from every service station to every other service station, then there are $\sum_{i=1}^N \binom{N}{i}$ possible deadlock types.

5.3 Literature Review

Most of the literature on blocking either does not consider networks with feedback loops, ignores the possibility of deadlock, or conveniently assumes the networks are deadlock-free. For closed networks of K customers with only one class of customer, [19] proves the following condition to ensure no deadlock: for each minimum cycle C , $K < \sum_{j \in C} B_j$, the total number of customers cannot exceed the total queueing capacity of each minimum subcycle of the network. The paper also presents algorithms for finding the minimum queueing space required to ensure deadlock never occurs, for closed cactus networks, where no two cycles have more than one node in common. This result is extended to multiple classes of customer in [21], with more restrictions such as single servers and each class having the same service time distribution. Here a integer linear program is formulated to find the minimum queueing space assignment that prevents deadlock. The literature does not discuss deadlock properties in open restricted queueing networks.

General deadlock situations that are not specific to queueing networks are discussed in [9]. Conditions for this type of deadlock, also referred to as deadly embraces, to potentially occur are given:

- Mutual exclusion: Tasks have exclusive control over resources.
- Wait for: Tasks do not release resources while waiting for other resources.
- No preemption: 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.

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 [8] the vertices and edges of the state graph are given labels in relation to a reference node. Using these labels *simple bounded circuits* are defined whose existence within the state graph is sufficient to detect deadlock.

5.4 Definitions

$ V(D) $	The order of the directed graph D is its number of vertices.
Weakly connected component	A weakly connected component of a digraph containing X is the set of all nodes that can be reached from X if we ignore the direction of the edges.
Direct successor	If a directed graph contains an edge from X_i to X_j , then we say that X_j is a direct successor of X_i .
Ancestors	If a directed graph contains a path from X_i to X_j , then we say that X_i is an ancestor of X_j .
Descendants	If a directed graph contains a path from X_i to X_j , then we say that X_j is a descendant of X_i .
$\deg^{\text{out}}(X)$	The out-degree of X is the number of outgoing edges emanating from that vertex.
Subgraph	A subgraph H of a graph G is a graph whose vertices are a subset of the vertex set of G , and whose edges are a subset of the edge set of G .
Sink vertex	A sink vertex is a vertex in a directed graph that has out-degree of zero.
Knot	In a directed graph, a knot is a set of vertices with out-edges such that while traversing the directed edges of that directed graph, once a vertex in the knot is reached, you cannot reach any vertex that is not in the knot.

5.5 State Digraph

Presented is a method of detecting when deadlock occurs in an open queueing network Q with N nodes, using a dynamic directed graph, the state graph.

Let the number of servers in node i be denoted by c_i . Define $D(t) = (V(t), E(t))$ as the state graph of Q at time t .

The vertices at time t , $V(t)$ correspond to servers in the queueing system. Thus, $|V(D(t))| = \sum_{i=1}^N c_i$ for all $t \geq 0$.

The edges at time t , $E(t)$ correspond to a blocking relationship. There is a directed edge at time t from vertex $X_a \in V(t)$ to vertex $X_b \in V(t)$ if and only if an individual occupying the server corresponding to vertex X_a is being blocked by an individual occupying the server corresponding to vertex X_b .

The state graph $D(t)$ can be partitioned into N service-station subgraphs, $D(t) = \bigcup_{i=1}^N d_i(t)$, where the vertices of $d_i(t)$ represent the servers of node i . The vertex set of each subgraph is static over time, however their edge sets may change.

The state graph is dynamically built up as follows. When an individual finishes service at node i , and this individual's next destination is node j , but there is not enough queueing capacity for j to accept that

individual, then that individual remains at node i and becomes blocked. At this point c_j directed edges between this individual's server and the vertices of $d_j(t)$ are created in $D(t)$.

When an individual is released and another customer who wasn't blocked occupies their server, that server's out-edges are removed. When an individual is released and another customer who was previously blocked occupies their server, that server's out-edges are removed along with the in-edge from the server who that previously blocked customer occupied. When an individual is released and there isn't another customer to occupy that server, then all edges incident to that server are removed.

This general process of building up the state graph as the queueing network is simulated will now be shown. Customers are labelled (i, j, k) where i denotes the server that customer is occupying, j denotes that individual's i.d. number, and k denotes the service station that customer is waiting to enter. As an example, a customer labelled $(A_2, 10, C)$ would have an i.d. number of 2, is occupying server A_2 and is currently waiting to join node C . If a customer isn't occupying a server the notation \emptyset is used. Similarly for customers occupying a server and still in service, their next destination is yet undecided, so \emptyset is used.

The simulation starts with full queues, and every server occupied by a customer in service. This is shown in Figure 14.

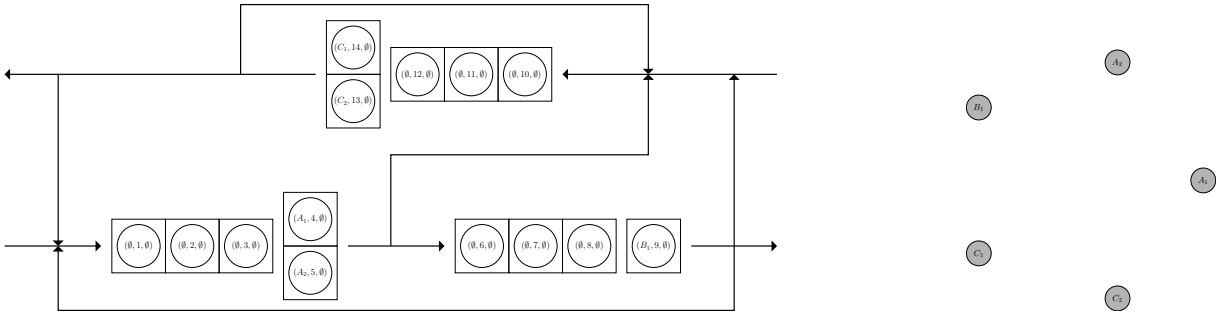


Figure 14

Customer 13 finishes service, and is blocked from entering node A . Figure 15.

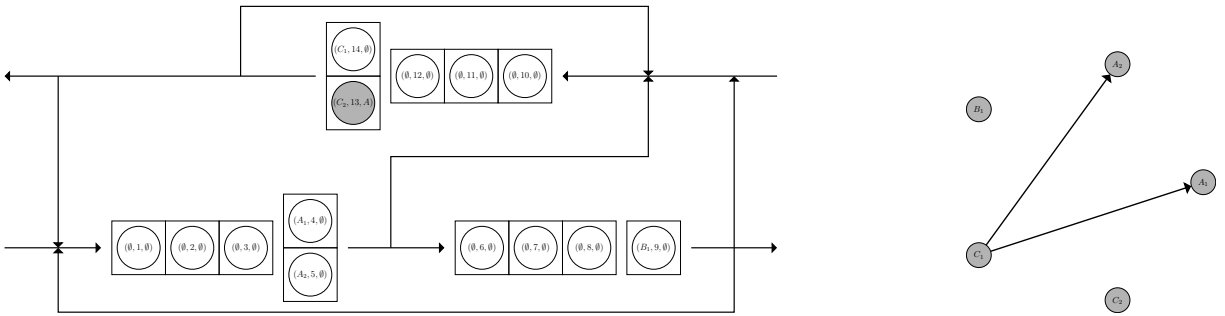


Figure 15

Then customer 4 finishes service and is blocked from entering node B . Figure 16.

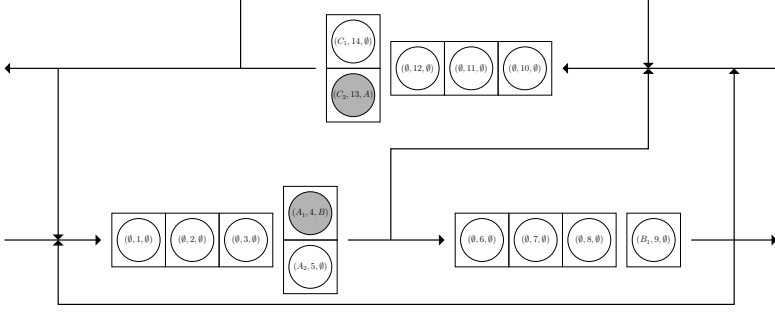


Figure 16

Then customer 9 finishes service and is blocked from entering node A . Figure 17.

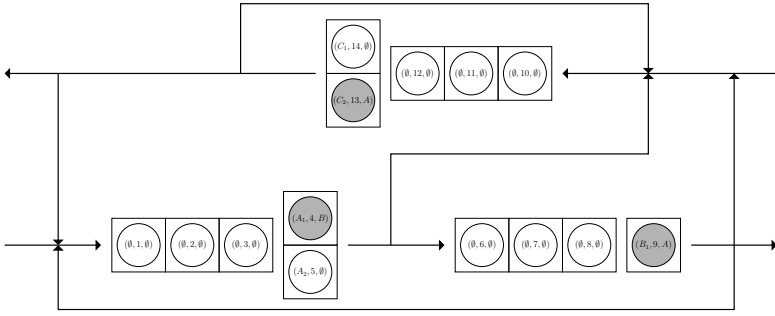


Figure 17

Finally, in Figure 18 customer 5 finishes service and wants to reenter the queue for node A but is blocked. A deadlock situation arises as customer 5 is waiting for customer 4 to move, who is waiting for customer 9 to move, who is waiting for either customer 4 or 5 to move.

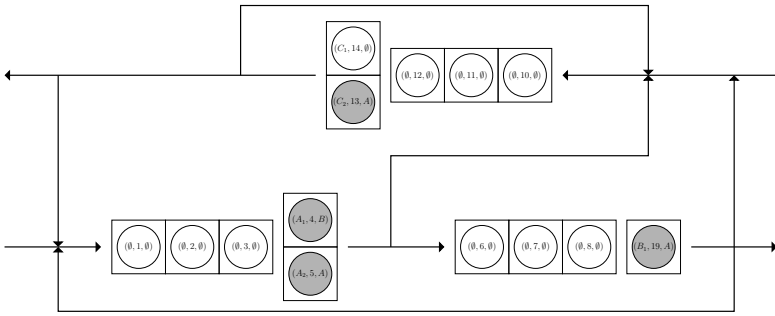
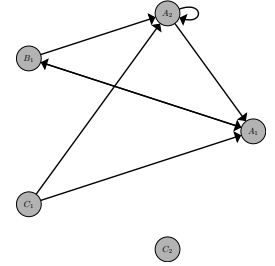
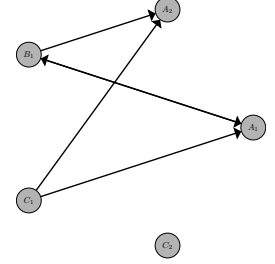
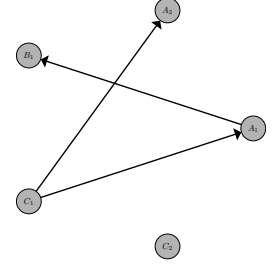


Figure 18

The rules on how edges are removed from the state graph will now be shown. For illustrative purposes the queueing network here is a different queueing network than discussed above.

Here the simulation begins with four customers occupying servers; those at node A blocked to node B ,



the customer at node C blocked to node A , and the customer at node B still in service. This is shown in Figure 19.

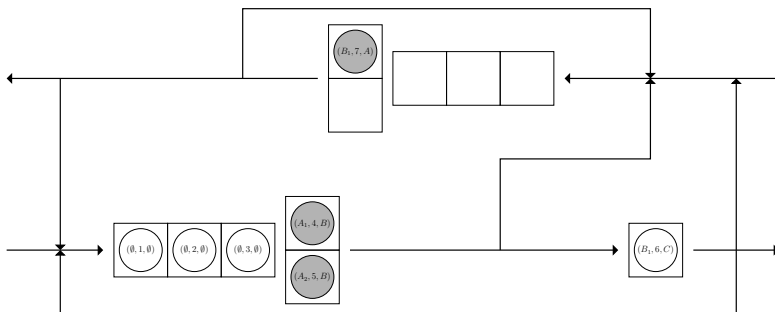
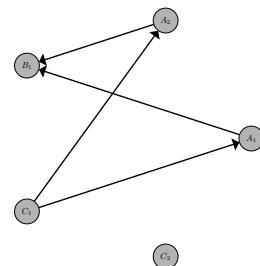


Figure 19



Customer 6 finishes service and immediately joins service at node C . Figure 20.

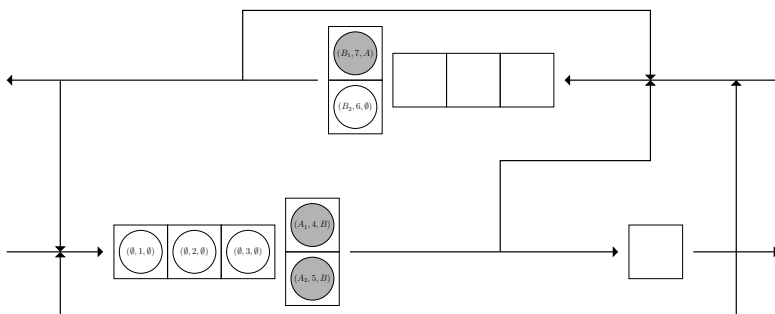
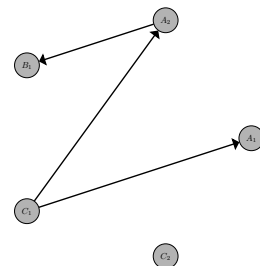


Figure 20



Now there is room for customer 4 to move into service at node B . Figure 21. Notice that the edge $A_2 \rightarrow B_1$ remains in the state graph, as customer 5 is still blocked by that server.

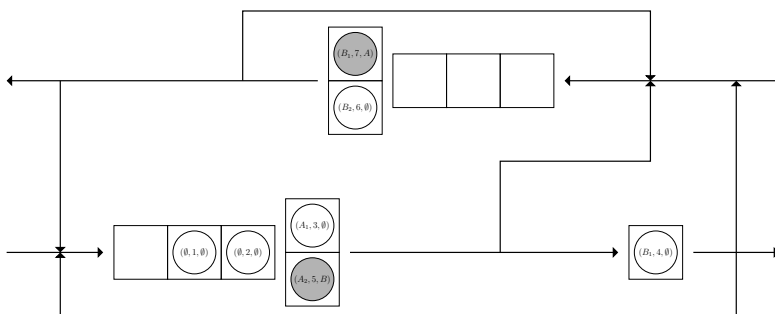
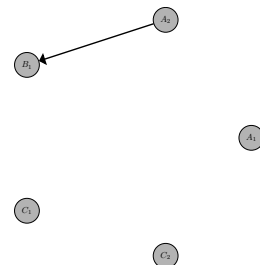


Figure 21



The customers queueing at node A move along the queue, with customer 3 beginning service. This leaves enough room for customer 7 to join the back of the queue at A . Figure 22.

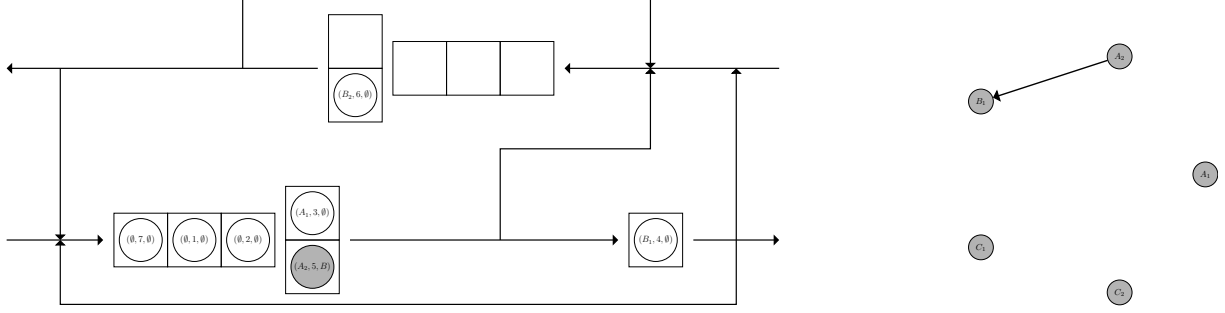


Figure 22

Observations

Consider one weakly connected component $G(t)$ of $D(t)$. Consider the node $X_a \in G(t)$. If X_a is unoccupied, then X_a has no incident edges. Consider the case when X_a is occupied by individual a , whose next destination is node j . Then X_a 's direct successors are the servers occupied by individuals who are blocked or in service at node j . We can interpret all X_a 's descendents as the servers whose occupants are directly or indirectly blocking a , and we can interpret all X_a 's ancestors as those servers whose individuals who are being blocked directly or indirectly by a .

Note that the only possibilities for $\deg^{\text{out}}(X_a)$ are being 0 or c_j . If $\deg^{\text{out}}(X_a) = c_j$ then a is blocked by all its direct successors. The only other situation is that a is not blocked, and $X_a \in G(t)$ because a is in service at X_a and blocking other individuals, in which case $\deg^{\text{out}}(X_a) = 0$.

It is clear that if all of X_a 's descendents are occupied by blocked individuals, then the system is deadlocked at time t . We also know that by definition all of X_a 's ancestors are occupied by blocked individuals.

Also note that if a service-station subgraph $d_i(t)$ contains edges, then there is an individual in $X_a \in d_i(t)$ that is being blocked by himself. This does not necessarily mean there is deadlock.

5.6 Results on the State Digraph

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

Proof. Consider one weakly connected component $G(t)$ of $D(t)$ at time t . All vertices of $G(t)$ are either descendents of another vertex and so are occupied by an individual who is blocking someone; or are ancestors of another vertex, and so are occupied by someone who is blocked.

Assume that $G(t)$ contains a vertex X such that $\deg^{\text{out}}(X) = 0$, and there is a path from every other non-sink vertex to X . This implies that X 's occupant is not blocked and is a descendent of another vertex. Therefore Q is not deadlocked as there does not exist a vertex whose descendents are all blocked.

Now assume that we have deadlock. For a vertex X who is deadlocked, all descendents of X are occupied by individuals who are blocked, and so must have out-degrees greater than 0. And so there is no path from

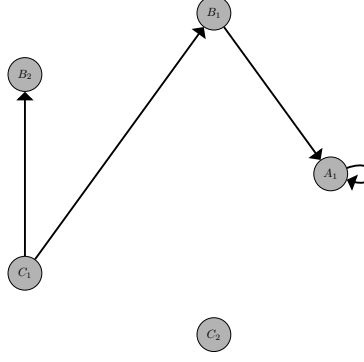


Figure 23: A Counter-Example State Digraph.

X to a vertex with out-degree of 0. □

Lemma 1. *For a queueing network with two nodes or less, a deadlocked state arises if and only if there exists a weakly connected component without a sink node.*

Proof. Consider a one node queueing network Q_1 .

If there is deadlock, then all servers are occupied by blocked individuals, and so all servers have an out-edge.

Consider a two node queueing network Q_2 .

If both nodes are involved in the deadlock, so there is a customer in node 1 blocked from entering node 2, and a 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, let's say 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 the server in node 2 has an out-edge, and so there is still no sink.

For the case of a queueing network with more than two nodes, the following counter-example proves the claim:

Begin with all servers occupied by customers in service. The customer at server B_1 is blocked from entering node A . Then the customer at server C_1 is blocked from entering node B . Then the customer at server A_1 is blocked from entering node A . The resulting state digraph in Figure 23 has a weakly connected component with a sink. □

Lemma 2. *An absorbing deadlocked state arises at time t if $D(t)$ doesn't contain a sink vertex.*

Proof. A vertex with out-degree greater than zero represents an occupied server whose occupant has finished service and is blocked. If all vertices have out-degree greater than zero, then all servers are occupied by blocked individuals. A release at vertex X_a can only be triggered by one of X_a 's descendants finishing service. As all servers are occupied by blocked individuals, no server can finish service, and so no server can release their occupant, implying an absorbing deadlocked state. \square

5.7 Finding Knots

Here is an initial draft of an algorithm for dynamically finding a knot in the state digraph of the queueing network.

Theorem 2. *The existence of a knot in the state digraph is equivalent to the existence of a vertex with an out-edge but without a path to a sink.*

Proof. Consider a vertex X that has an out-edge, but does not have a path to a sink. All of X 's descendants also do not have paths to sinks.

As the number of vertices in $D(t)$ are finite ($|V(D(t))| = \sum_{i=1}^N c_i$ for all $t \geq 0$), then all the possible paths leading out from X must either be part of a cycle or lead to a cycle. And so all paths from X must terminate in a knot. \square

Keeping track of the list of sinks in D will be useful. Let's call this list `sink_list`. Initially all vertices of $D(0)$ are sinks, and belong to `sink_list`. At the point when a vertex gains an out-edge, that vertex is removed from `sink_list`. At the point when a vertex has an out-edge removed, if $\deg^{\text{out}} = 0$ then that vertex is added to `sink_list`. Therefore, at any one point we have knowledge of all the sinks in $D(t)$.

Now, the only action that can lead to a vertex transitioning from having a path to a sink and not having a path to a sink is adding an edge. Removing an edge of $D(t)$ will not cause deadlock, as removing an edge implies that a customer has moved, and a customer moving cannot cause deadlock. Therefore, we only need to check the vertex's paths to sinks when an edge is added.

5.8 Markov Chain Model

It is interesting to build an analytical model of the system's behaviour to deadlock. As a Markov chain model, the deadlocking state is an absorbing state, and so any queueing network that can experience deadlock is guaranteed to experience deadlock.

We can however find the expected time until deadlock is reached. It is shown in [25] that for a discrete transition matrix of the form $P = \begin{pmatrix} T & U \\ 0 & V \end{pmatrix}$ then the expected number of time steps until absorption starting from state i is the i th element of the vector

$$(I - T)^{-1}e \tag{13}$$

where e is a vector of 1s.

5.8.1 One Node Network

Consider the one node network with feedback loop shown in Figure 24. There is room for n customers to queue at any one time, customers arrive at a rate of Λ and served at a rate μ . Once a customer has finished service he rejoins the queue with probability r_{11} , and so exits the system with probability $1 - r_{11}$.

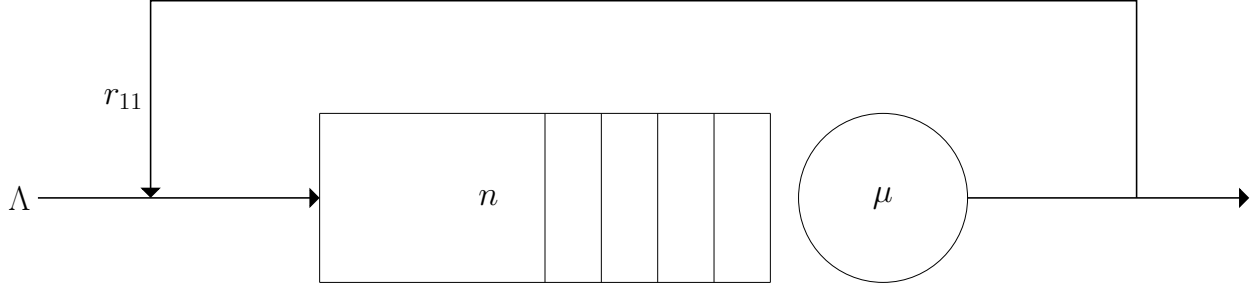


Figure 24: A one node queueing network.

State space:

$$S = \{i \in \mathbb{N} | 0 \leq i \leq n + 1\} \cup \{-1\}$$

where i denotes number of individuals in service or waiting.

If we define $\delta = i_2 - i_1$ for all $i_k \in S | i_k \geq 0$, then the transitions are given by:

$$q_{i_1, i_2} = \begin{cases} 0 & \text{if } i = n + 1 \\ \Lambda & \text{otherwise} \end{cases} \quad \text{if } \delta = 1$$

$$\begin{cases} (1 - r_{11})\mu & \text{if } \delta = -1 \\ 0 & \text{otherwise} \end{cases} \quad \text{if } \delta = -1$$

$$0 \quad \text{otherwise}$$
(14)

$$q_{i, -1} = \begin{cases} r_{11}\mu & \text{if } i = n + 1 \\ 0 & \text{otherwise} \end{cases}$$
(15)

and

$$q_{-1, i} = 0$$
(16)

The Markov chain is shown in Figure 25.

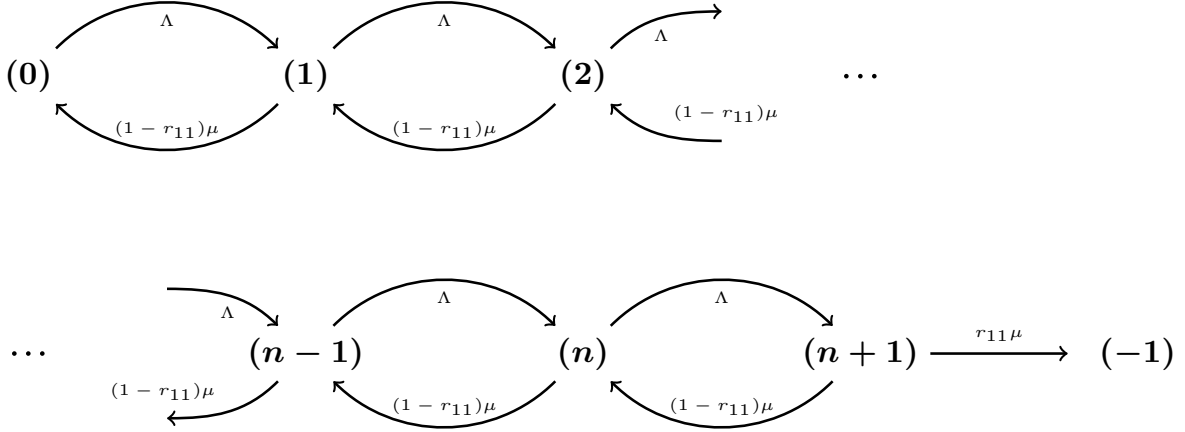


Figure 25: Markov chain of the one node system.

5.8.2 Two Node Network without Self Loops

Consider the queueing network shown in Figure 26. This shows two $M/M/1$ queues, with n_i queueing capacity at each service station and service rates μ_i . Λ_i is the external arrival rates to each service station. All routing possibilities except self loops are possible, where the routing probability from node i to node j is denoted by r_{ij} .

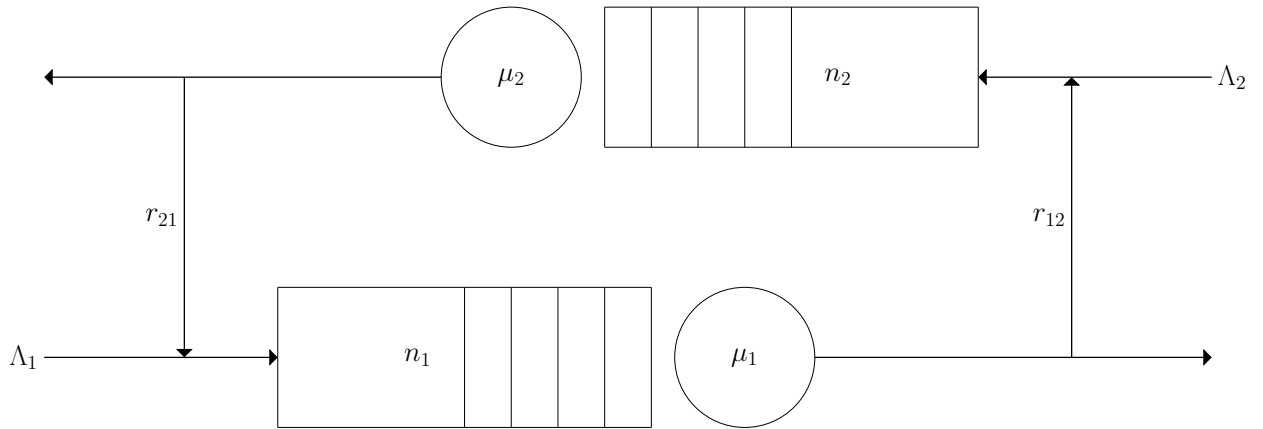


Figure 26: A two node queueing network.

- State space:

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

where i denotes number of individuals:

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

where j denotes number of individuals:

- In service or waiting at the second node.
 - Occupying a server but having finished service at the first node waiting to join the first
- and the state $(-1, -1)$ denotes the deadlocked state.

If we define $\delta = (i_2, j_2) - (i_1, j_1)$ for all $(i_k, j_k), s \in S | i_k, j_k \geq 0$, then the transitions are given by:

$$q_{(i_1, j_1), (i_2, j_2)} = \begin{cases} \left. \begin{array}{l} \Lambda_1 \text{ if } i_1 \leq n_1 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (1, 0) \\ \left. \begin{array}{l} \Lambda_2 \text{ if } j_1 \leq n_2 \\ 0 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (0, 1) \\ \left. \begin{array}{l} 0 \text{ if } j_1 = n_1 + 2 \\ (1 - r_{12})\mu_1 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (-1, 0) \\ \left. \begin{array}{l} 0 \text{ if } i_1 = n_1 + 2 \\ (1 - r_{21})\mu_2 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (0, -1) \\ \left. \begin{array}{l} 0 \text{ if } j_1 = n_2 + 2 \\ r_{12}\mu_1 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (-1, 1) \\ \left. \begin{array}{l} 0 \text{ if } i_1 = n_1 + 2 \\ r_{21}\mu_2 \text{ otherwise} \end{array} \right\} & \text{if } \delta = (1, -1) \\ 0 & \text{otherwise} \end{cases} \quad (17)$$

$$q_{(i_1, j_1), (-1, -1)} = \begin{cases} r_{21}\mu_2 & \text{if } (i, j) = (n_1, n_2 + 2) \\ r_{12}\mu_1 & \text{if } (i, j) = (n_1 + 2, n_2) \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

and

$$q_{-1, s} = 0 \quad (19)$$

For $n_1 = 1$ and $n_2 = 2$, the resulting Markov chain is shown in Figure 27.

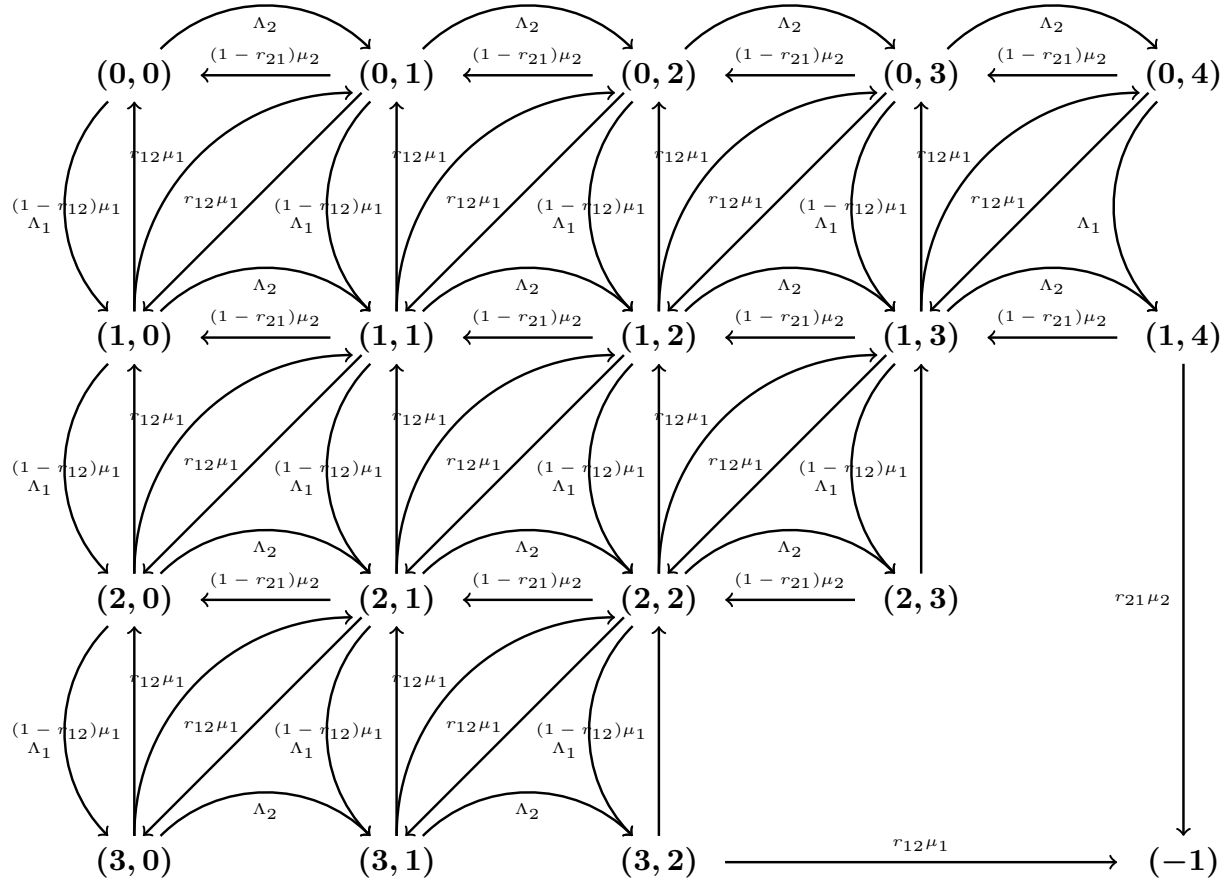


Figure 27: Markov chain of the two node system with $n_1 = 1$ and $n_2 = 2$.

Figure 28 shows the effect of varying the parameters of the above Markov model. Only parameters for one node are shown, as the other node's parameters will have the same affect. We can see that increasing the arrival rate Λ_1 and the transition probability r_{12} results in reaching deadlock faster. This is intuitive as increasing these parameters results in the first node's queue filling up quicker. Increasing the service rate μ_1 and the queueing capacity n_1 results in reaching deadlock slower. Again these are intuitive, as increasing the service rate results in moving customers out of the system quicker, and increasing the queueing capacity allows more customers in the system before becoming deadlock.

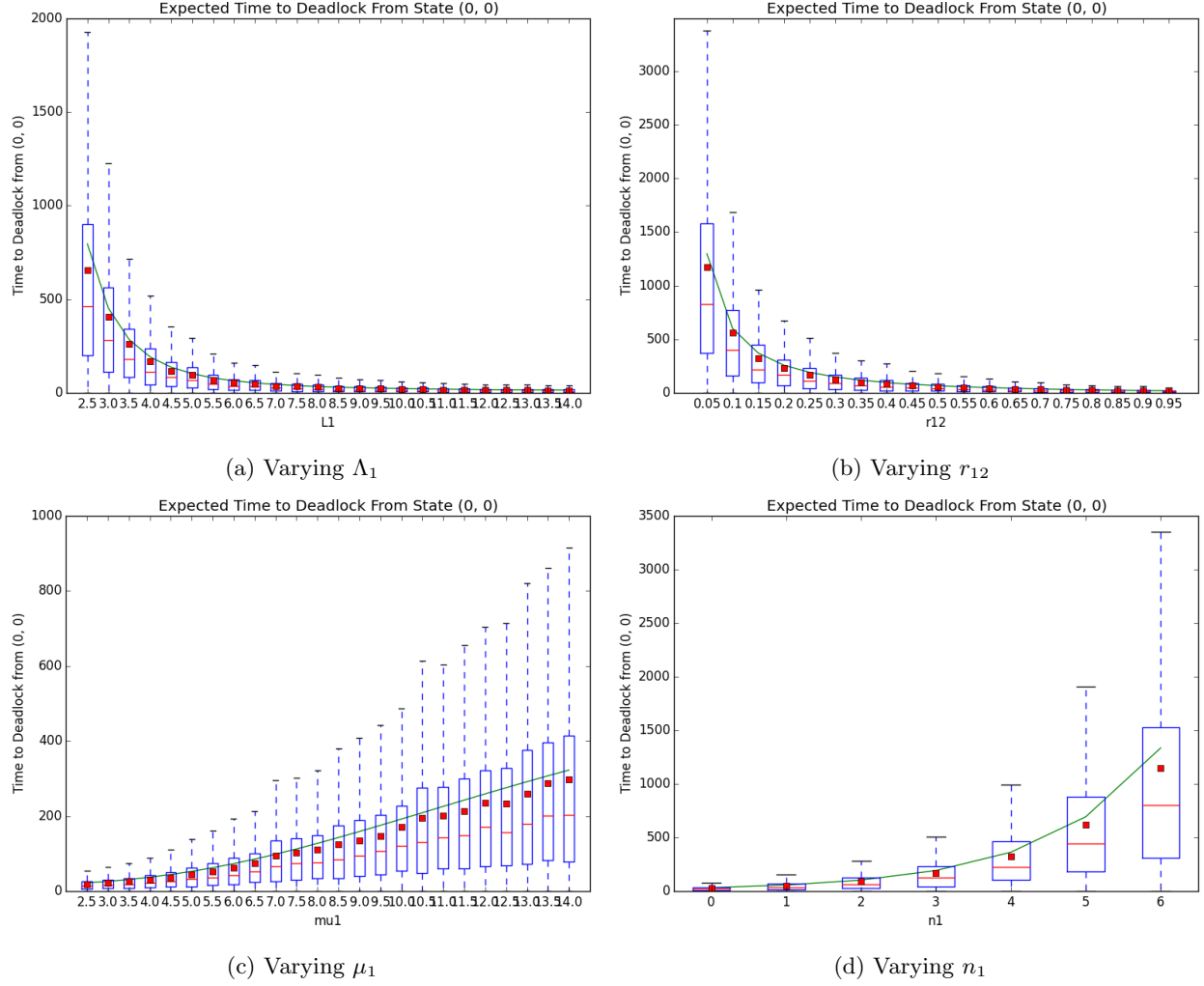


Figure 28: Analytical & Simulation Results of Times to Deadlock (1000 iterations)

5.9 Resolving Deadlock

Once the system falls into a deadlocked state for the first time, that is the first transient deadlocked state since the last resolution, the simulated needs to automatically resolve the deadlock and allow services to resume again. This is not necessarily as simple as moving a blocked customer to his next node, as we need to conserve the numbers of customers at each service station. Closer inspection of the state digraph is required in order to find a way to resolve deadlock whilst conserving this property.

At deadlock, the service stations can be classified into the following three mutually exclusive categories:

- Nodes that are not deadlocked: These are nodes that do not contain any blocked individuals.
- Causation nodes, nodes causing deadlock: These are nodes where every server is occupied by a blocked individual, and there is at least one blocked individual waiting to enter that node who is directly or

indirectly being blocked by an individual in this node.

- Affected nodes, nodes affected by deadlock: These are nodes containing at least one blocked individual who is directly or indirectly being blocked by an individual at a node that is causing deadlock, but is not classified as causing deadlock itself.

At the first instance of deadlock, there will only be one knot in $D(t)$. Let us denote the knot as K . The vertices of K correspond to servers. As there is no sink node, all vertices of K have an out-edge, and so all vertices in K contain a blocked individual. Therefore, there are no vertices in K belonging to nodes that are not deadlocked. All vertices of K correspond to servers of causation nodes, and a causation node has all its servers belonging to K . An affected node has servers belonging to the same weakly connected component as K , but does not have servers in K .

When choosing which customer to move in order to resolve deadlock, we must be careful to conserve the number of customers at each service station. Causation nodes have full queues, and a customer may only be moved into a full queue if this causes another customer to simultaneously move from this node. Another complication arises due to the blocking mechanism used, in which those customers who have been blocked longer must be moved first. This property may have to be broken in order to ensure the conservation property is not. Assume that we have weighted the edges of the digraph with the time that they were created.

The following algorithm is proposed in order to resolve deadlock:

```

Find the knot  $K$  in  $D(t)$ 
Find the cycle  $C \in K$  whose average edge weight is minimum
Start at  $V_0$ 
for  $V_i$  in  $C$  do
    | Move the individual who is waiting to get to  $V_{i+1}$ 
end
Redraw  $D(t)$ 

```

References

- [1] S. Albin, J. Barrett, D. Ito, and J.E. Mueller. A queueing network analysis of a health center. *Queueing systems*, 7(1):51–61, 1990.
- [2] C. Ancker Jr and A. Gafarian. Some queueing problems with balking and reneging i. *Operations research*, 11(1):88–100, 1963.
- [3] C. Ancker Jr and A. Gafarian. Some queueing problems with balking and reneging ii. *Operations research*, 11(6):928–937, 1963.
- [4] B. Avi-Itzhak and M. Yadin. A sequence of two servers with no intermediate queue. *Management science*, 11(5):553–564, 1965.
- [5] J. Baber. *Queues in series with blocking*. PhD thesis, Cardiff University, 2008.

- [6] H. Buhaug. Long waiting lists in hospitals: operational research needs to be used more often and may provide answers. *BMJ: British medical journal*, 324(7332):252, 2002.
- [7] P. Burke. The output of a queueing system. *Operations research*, 4(6):699–704, 1956.
- [8] H. Cho, T. Kumaran, and R. Wysk. Graph-theoretic deadlock detection and resolution for flexible manufacturing systems. *IEEE transactions on robotics and automation*, 11(3):413–421, 1995.
- [9] E. Coffman and M. Elphick. System deadlocks. *Computing surveys*, 3(2):67–78, 1971.
- [10] S. Creemers and M. Lambrecht. Modelling a healthcare system as a queueing network: the case of a belgian hospital. *Available at SSRN 1093618*, 2007.
- [11] F. Gorunescu, S.I. McClean, and P.H. Millard. A queueing model for bed-occupancy management and planning of hospitals. *The journal of the operational research society*, 53(1):19–24, 2002.
- [12] P. Harper. A framework for operational modelling of hospital resources. *Health care management science*, 5(3):165–173, 2002.
- [13] G. Hunt. Sequential arrays of waiting lines. *Operations research*, 4(6):674–683, 1956.
- [14] K. Hurst. *Selecting and applying methods for estimating the size and mix of nursing teams: A Systematic Review of the Literature Commissioned by the Department of Health*. Nuffield Institute for Health, 2003.
- [15] J. Jackson. Networks of waiting lines. *Operations research*, 5(4):518–521, 1957.
- [16] F. Kelly. Networks of queues with customers of different types. *Journal of applied probability*, 12(3):542–554, 1975.
- [17] N. Koizumi, E. Kuno, and T.E. Smith. Modeling patient flows using a queueing network with blocking. *Health care management science*, 8(1):49–60, 2005.
- [18] R. Korporaal, A. Ridder, P. Klopogge, and R. Dekker. An analytic model for capacity planning of prisons in the netherlands. *The journal of the operational research society*, 51(11):1228–1237, 2000.
- [19] S. Kundu and I. Akyildiz. Deadlock buffer allocation in closed queueing networks. *Queueing systems*, 4(1):47–56, 1989.
- [20] G. Latouche and M. Neuts. Efficient algorithmic solutions to exponential tandem queues with blocking. *SIAM journal on algebraic discrete methods*, 1(1):93–106, 1980.
- [21] J. Liebeherr and I. Akyildiz. Deadlock properties of queueing networks with finite capacities and multiple routing chains. *Queueing systems*, 20(3-4):409–431, 1995.
- [22] K. Rege. Multi-class queueing models for performance analysis of computer systems. *Sadhana*, 15(4-5):355–363, 1990.
- [23] E. Reich. Waiting times when queues are in tandem. *The annals of mathematical statistics*, 28(3):768–773, 1957.
- [24] S Robinson. *Simulation: the practice of model development and use*. John Wiley & Sons, Ltd, 2007.
- [25] W. Stewart. *Probability, markov chains, queues, and simulation*. Princeton university press, 2009.

- [26] R. Sutton and A. Barto. *Reinforcement learning: an introduction*. MIT press, 1998.
- [27] C. Szepesvri. *Algorithms for reinforcement learning*. Morgan & Claypool Publishers, 2010.
- [28] Y. Takahashi, H. Miyahara, and T. Hasegawa. An approximation method for open restricted queueing networks. *Operations research*, 28(3):594–602, 1980.
- [29] A. Turing. Computing machinery and intelligence. *Mind*, 59(236):243–260, 1950.
- [30] P.T. Vanberkel, R.J. Boucherie, E.W. Hans, J.L. Hurink, and N. Litvak. A survey of health care models that encompass multiple departments. 2009.
- [31] J. Viana, S. Rossiter, A. Channon, S. Brailsford, and A Lotery. A multi-paradigm, whole system view of health and social care for age-related macular degeneration. In *Proceedings of the Winter Simulation Conference*, 2012.