# Modeling the impact of interference on wireless ad hoc network performance

Tom Coenen

# Modeling the impact of interference on wireless ad hoc network performance

Tom Coenen

**Graduation committee:**

**Chairman:**
Prof. dr. P.M.G. Apers          University of Twente
**Promotors:**
Prof. dr. R.J. Boucherie          University of Twente
Prof. dr. J.L. van den Berg          University of Twente
**Co-promotor:**
dr. ir. M. de Graaf          Thales Netherlands B.V., University of Twente

**Members:**
Prof. dr. C. Blondia          University of Antwerp
Prof. dr. ir. S.M. Heemstra - de Groot          Eindhoven University of Technology
Prof. dr. ir. H.J. Broersma          University of Twente
dr. ir. G.J. Heijenk          University of Twente
dr. ir. J. Goseling          University of Twente

# MODELING THE IMPACT OF INTERFERENCE ON WIRELESS AD HOC NETWORK PERFORMANCE

PROEFSCHRIFT

ter verkrijging van
de graad van doctor aan de Universiteit Twente,
op gezag van de rector magnificus,
prof. dr. T.T.M. Palstra,
volgens het besluit van het College voor Promoties,
in het openbaar te verdedigen
op vrijdag 9 juni 2017 om 14:45

door

Tom Johannes Maria Coenen

geboren op 24 maart 1980
te Venray, Nederland

Dit proefschrift is goedgekeurd door:

Prof. dr. R.J. Boucherie (promotor)
Prof. dr. J.L van den Berg (promotor)
dr. ir. M. de Graaf (co-promotor)

*Voor mijn ouders*
*voor jullie onvoorwaardelijke steun*

# Voorwoord

Een PhD traject begint officieel wanneer je aangesteld wordt aan de universiteit, maar eigenlijk is het een stap die volgt op de vele stappen die je in je leven ervoor al hebt gezet. Het gaat misschien te ver om bij het basisonderwijs, laat staan de kleuterklas te beginnen, maar ik denk dat voor mij het voortgezet onderwijs toch wel het begin is geweest van deze carrièrestap. Hier kwam de interesse in de wetenschap voor mij tot leven met in het bijzonder mijn liefde voor wiskunde, al zal de eigenlijk niet bestaande "uitmuntend" voor rekenen op mijn rapport op de basisschool hier misschien ook wel iets mee te maken hebben gehad. De keuze voor Toegepaste Wiskunde aan de Universiteit Twente kwam pas na het bezoeken van ongelofelijk veel verschillende studies, want mijn interesse was breed. Die brede interesse werd tactisch ingezet om aan het eind van mijn studie een afstudeerproject te accepteren dat meerdere gebieden van de wiskunde omvatte. Na dit project was de stap richting het promotietraject dat resulteerde in dit proefschrift snel gezet, het leek er bijna naadloos in door te vloeien. Maar hoe soepel de weg tot dan toe gelopen was bleek geen voorbode voor hoe mijn promotietraject zou verlopen. We zijn inmiddels zo'n 13 jaar verder dan toen het begon. In de tussentijd zijn er naast het werken om tot dit proefschrift te komen ook werkzaamheden geweest als docent op de UT, als leraar in het voortgezet onderwijs en als vakdidacticus aan de lerarenopleiding, waarbij de laatste twee dat ook in de toekomst nog zullen zijn. In deze lange tijd zijn er vele mensen geweest die een rol hebben gespeeld in het tot stand komen van dit proefschrift en deze mensen wil ik hiervoor hartelijk danken, een aantal zal ik specifiek noemen.

Om te beginnen wil ik mijn promotoren Richard J. Boucherie en Hans van den Berg en mijn co-promotor Maurits de Graaf bedanken voor hun begeleiding en ondersteuning gedurende mijn promotietraject. Richard, na mijn afstudeerproject bij jou vroeg je of ik door wou gaan in een promotietraject, waar je hoge verwachtingen bij had. Ik denk dat het anders is gelopen dan je toen had verwacht en ik ben je dankbaar dat je ondanks dit andere verloop mij het vertrouwen hebt gegeven dat dit toch tot een goed einde kon komen. Je directe aanpak en de door je commentaar compleet rood gekleurde drafts van mijn papers hebben me veel geleerd over de academische wereld. Ik hoop in mijn toekomstige baan ook nog van je expertise en visie gebruik te mogen maken. Hans, het was een voorrecht om jou als rustige en constante factor als promotor te hebben. Je duidelijke aanwijzingen en gerichte commentaar zijn altijd van grote waarde geweest. Maurits, als co-promotor heb je mij in een periode dat de vaart eruit begon te raken op de goede weg weten te zetten. Ik bewonder hoe je het voor elkaar krijgt op de ene dag op de UT zoveel werk gedaan te krijgen, inclusief het begeleiden van PhD studenten. Samen met jou heb ik een

ander vlak van de wiskunde kunnen toevoegen aan dit proefschrift, waarvoor mijn dank. De leden van mijn promotiecommissie, prof. dr. Chris Blondia, prof. dr. ir. Sonia Heemstra - de Groot, prof. dr. ir. Hajo Broersma, dr. ir. Geert Heijenk en dr. ir. Jasper Goseling dank ik voor de bereidheid mijn proefschrift te beoordelen.

Met iedereen van de vakgroep SOR is het altijd fijn samenwerken geweest in een vriendelijke en open sfeer. Jasper, het einde van mijn promotietraject was nooit in zicht gekomen en ook niet bereikt zonder jouw ondersteuning. Vooral in mijn sabbatical de laatste drie maanden stond je altijd klaar om commentaar te geven op mijn vele drafts, mee te denken als ik mezelf weer eens in de war had gebracht en positiviteit uit te stralen dat dit eindpunt bereikt zou worden. Jan-Kees en Werner, jullie deur, die altijd open stond om te praten over waar ik maar tegenaan liep in mijn onderzoek, heb ik vaak dankbaar gebruik van gemaakt. Alle PhD's, dat zijn er zo veel dat ik niet ga proberen alle namen te noemen, het was fijn om met jullie samen dit traject te doorlopen.

Iedereen van de vakgroep OMPL, waar ik een aantal jaar als docent werkzaam heb mogen zijn. Wat een ontzettend fijne sfeer heerste er altijd bij jullie, met de vele koffiemomenten samen en de bijeenkomsten buiten de UT. Professor de Smit en professor Zijm, bedankt voor het geven van de mogelijkheid om deel uit te maken van deze groep. Erwin, Matthieu, Ahmad, Marco en Martijn, de samenwerking met jullie maakte mijn tijd daar een waardevolle en plezierige periode uit mijn leven. Ook de vele PhD's uit deze groep, waar ik weer niet ga proberen alle namen te noemen, bedankt voor de gezellige tijd!

Collega's van Reggesteyn, de school waar ik als wiskundeleraar aan de slag ben gegaan, ik voelde me meteen welkom bij jullie. Iedereen was vriendelijk en behulpzaam om mij als startend leraar wegwijs te maken. Inmiddels ben ik de status van startend leraar wel voorbij, maar die vriendelijkheid en behulpzaamheid zijn nooit veranderd. Wendy, bedankt dat je mee wilde helpen om mij een sabbatical te laten nemen om dit proefschrift te kunnen voltooien. Wim, Erik, Hans, Christiaan, Gerrit-Jan, Esther, Liset, Ina en Erik, jullie zijn een erg leuke vakgroep om mee samen te werken. Ook jullie bedankt voor het opvangen van mijn uren zodat ik mijn sabbatical op kon nemen. Marcia, Marco en Aniek, jullie zijn echt leuk volk!

Collega's van ELAN, de lerarenopleiding en de nieuwste stap in mijn carrière, door jullie is het duidelijk dat ik een baan heb waar ik me helemaal in mijn element voel. Nellie, bedankt voor het mij introduceren in deze wereld middels de CoL en het uitspreken van het vertrouwen dat ik samen met Mark uiteindelijk jouw positie zou over kunnen nemen. Met jouw kennis en enorme netwerk zal dat geen eenvoudige taak zijn, maar je schijnbaar oneindige enthousiasme en inzet werken aanstekelijk. Hopelijk mag ik nog lang met je samenwerken. Mark, het samen beginnen als vakdidacticus en de soepele samenwerking vanaf de start is fantastisch. Ik denk dat we samen een mooie tijd tegemoet gaan met deze nieuwe uitdaging. Gerard, het is geruststellend jou als ervaren vakdidacticus bij ons te hebben. Adri, Susan en Jan, bedankt voor het bieden van deze mogelijkheid en het vertrouwen dat jullie uitstraalden bij het aanvaarden van mijn sollicitatie ondanks dat mijn proefschrift nog niet voltooid was.

Zonder vriendschappen komt niemand ver en dat geldt zeker ook voor mij.

De ontspannende momenten naast het werk zijn onontbeerlijk geweest. Het sporten bij ENTAC is daar een goed voorbeeld van. Ik haal nog steeds veel energie uit de inspannende ontspanning van het trainen en de competitie en ben trots deel uit te mogen maken van het bestuur en het eerste team van deze mooie club. Thijs, vooral van jou als trainer en teamgenoot heb ik ongelofelijk veel kunnen leren en ik ben blij dat ik jou en Rianne tot mijn vriendengroep mag rekenen zodat we ook naast het tafeltennis veel mooie momenten hebben kunnen beleven. Jeroen en Marjan en Gert en Annemarie, bedankt voor de gezellige tijden samen waarin vele verhalen over ouderschap van grote waarde zijn geweest. Jaap en Marianne, René en Hanneke en Chris en Esther, het is fantastisch dat onze vriendschap die begon in het voortgezet onderwijs nog steeds bestaat. De sporadische bijeenkomsten met al onze geweldige kids zijn altijd weer iets om naar uit te kijken. Marie Jose, de avondjes squash zijn al een lange tijd terug gestopt maar gelukkig vinden we zo nu en dan de tijd om bij te kletsen en zetten we nog regelmatig een degelijk resultaat neer in de pubquiz als vast duo in een steeds wisselende samenstelling van het team. Peter and Connie, you are simply amazing. I don't think I can imagine people that are more friendly and give so much hospitality. You were there for me during a difficult time, opening your home to me even with your baby on the point of arriving. You also opened your home in Canada to me as a base for my explorations of this beautiful country. I will always treasure these memories and keep hoping I can return the favour sometime. Matthias und Elli, ich freue mich unglaublich das aus meine Zeit in Münster so ein schöne Freundschaft gewachsen ist und das wir uns mit unseren Familien für sehr tolle Wochenenden regelmäßig treffen. Floris en Suzan, een goede buur is beter dan een verre vriend. In dit geval hebben we goede vrienden als buren en dat is onbetaalbaar. Bedankt voor de vele spontane en gezellige avonden, spelletjes, filmpjes, etentjes, uitmuntende bbq's, uitjes en een mooi festivalweekend.

Pap en mam, ik zeg het vaak genoeg tegen anderen: "Ik heb de beste ouders van de wereld". Hopelijk weten jullie hoe belangrijk jullie voor me zijn. De basis voor dit proefschrift werd al jong gelegd door jullie aanmoediging om aan de toekomst te denken en mijn best te doen op school. Altijd kan ik op jullie ondersteuning rekenen, op welk vlak dan ook. Ik draag mijn proefschrift met liefde aan jullie op. Nicole en Ludo, ik ben blij met onze band die we de laatste jaren verder hebben zien groeien en de heerlijke momenten die we samen beleven met onze kids. Herbert en Annie, je vriendin kies je, je schoonouders krijg je er bij. Gelukkig heb ik het daarbij uitstekend getroffen.

Marlies, ik vind je lief. Die woorden zeggen je denk ik alles en dat is maar goed ook want verder kan ik niet in woorden uitdrukken hoe heerlijk het is je in mijn leven te hebben. Samen met ons prachtige mannetje Sven maak je het leven geweldig. Ik hou van jullie!

Tom

# Table of Contents

# Introduction

Wireless communication has been developing rapidly in the past decades and the world relies on it in a still increasing fashion. Pioneering work was done in the late 1800's and the 20th century brought many new developments and devices. This development culminated into the appearance of the mobile phone. Yet the development did not stop there. Next to communication between people, many different devices are now communicating, ever collecting, analysing and interpreting data and taking appropriate actions as a result. One can think of the use of GPS in navigation systems and the real-time traffic information that is included in it, a home cinema setup with TV, speakers and hub, all controlled with a mobile phone, or the monitoring of agricultural areas and starting irrigation when the fields get too dry. The world is relying heavily on the reliability and stability of such communication.

To ensure successful wireless communication between devices several aspects need to be taken into account: The devices must be equipped with appropriate hardware to transmit and receive the signals and must have enough battery power to complete communication. As wireless signals fade over distance, the distance between the devices must be limited to ensure that they can reach each other. With multiple devices trying to communicate at the same time, the transmitted signals may collide and disrupt the reception of these signals. This phenomenon is known as interference. Data to be transmitted needs to be stored at the devices, which must have sufficient capacity to do so. Also, the time needed to complete a transmission has to be limited and the network capacity has to be sufficient to transmit all the data.

Next to infrastructure-based networks, networks without an infrastructure are becoming more common. These networks are enoted as ad hoc networks. Ad hoc networks are characterized by a group of (mobile) users who communicate with each other without the use of dedicated network nodes and without any centralized control, i.e. these networks are self-configuring.

This thesis focuses on the impact of interference on the performance of wireless ad hoc networks. Mathematical models are presented that analyse the impact on the capacity of the network, the delay packets experience and the throughput the network can achieve. Various views are adopted to take interference into account, such as an interference graph showing which devices can and cannot transmit at the same time. Also the lifetime of the network is considered, as the battery capacity of the devices in ad hoc networks is often limited. The models presented in this thesis contribute to the understanding of the impact of interference and provide insights that are of interest when

Figure 1.1: A Wireless Local Area Network

designing or deploying ad hoc networks.

As a complete description of the world of wireless communication is impossible to give, this chapter first presents a short overview of the characteristics of wireless ad hoc networks addressed in this thesis and provides an introduction in the terminology that is used. The second section presents some basic graph and queuing theory of interest. The third section discusses the addressed research questions and the fourth section the contributions of this thesis. The fifth and final section presents an outline of the thesis.

## 1.1   Wireless ad hoc networks

Under the term network we understand a collection of devices that want to exchange data. The different devices in the network are called the nodes of the network and two nodes are connected by a so called link when direct communication between these nodes is possible. As the word wireless literally says, no wires (or cables) are involved in a wireless network, but communication takes place over radio waves. This provides a number of advantages over the wired network, such as the ability to move around with the device without losing the connection to the other devices and lower costs.

Many different types of wireless networks exist, the most commonly known being the Wireless Local Area Network (WLAN). Such a network links a couple of devices over a short distance, often to an access point that connects the devices to the Internet. A cellular network is a mobile network with so called base stations, each serving a certain area ('cell') around it. These cells together provide coverage over larger geographical areas. Devices such as mobile phones are therefore able to communicate even if the user is moving through cells during transmission. The first generation, 1G, made it possible to carry analogue voice over channels. With the introduction of 2G, data services such as SMS became possible. Its successor, 3G, offered faster rates, making video calls possible. 4G

Figure 1.2: Example of an ad hoc network, connected to the Internet

and the upcoming 5G improve the data transfer rates even further.

Wireless ad hoc networks are characterized by their decentralized nature. Where most wireless networks have a configured infrastructure with centralized control, ad hoc networks are self-configuring and dynamic. An example of an ad hoc network, connected to the Internet, is shown in Figure 1.2. Due to the high mobility of the users, which are the nodes of the network, the topology of the network constantly changes. This calls for dynamic routing, which is capable of taking these frequent changes into account. Communication between the users takes place over multiple hops, as other users forward messages to deliver them to the right recipient.

Ad hoc networks are easy and quick to deploy. No specific tasks are assigned to the nodes of the network and no routing is prescribed, making ad hoc networks very suitable for situations where infrastructure no longer exists such as when natural disasters have destroyed the infrastructure or in war situations. As the nodes in the network can be very simple, the costs of such a network can be low. The decentralized nature of the network increases the mobility of the network, nodes can move around without destroying the infrastructure. Ad hoc networks are also robust, as the failure of a single node generally does not influence the overall connectivity of the network.

Within the group of ad hoc networks there are again different types. A mobile ad hoc network (MANET) consists of continuously self-configuring mobile devices connected without wires. A vehicular ad hoc network (VANET) is an ad hoc network between vehicles, which for example can be used in traffic to warn cars for upcoming congestions or accidents. Wireless sensor networks consist of sensors deployed in an area they need to monitor. Data that is collected is then forwarded through the other sensors to some collection point, for example to monitor a forest [KNB$^+$06]. On a smaller scale, all devices close to a user which can communicate wirelessly are considered a Personal Area Network (PAN) [CGJ$^+$06].

Wireless (ad hoc) networks face a number of challenges that don't play

Figure 1.3: The hidden node problem

in wired networks [Pet06],[Wil06]. This thesis will particularly focus on the impact of interference. As radio signals use a certain frequency, signals sent over the same channel can collide, meaning that two communications arriving at a receiving user at the same time disrupts the reception of these signals. The information that has been transmitted is not received correctly by the node, which may then be retransmitted or is lost. Even a a single flow of packets through a network can cause self-interference, as multiple nodes may be involved simultaneously in the transmission of packets.

Even though an advantage of ad hoc networks is that they are more flexible, this also creates a disadvantage. When the nodes in the network are very mobile, the topology of the network constantly changes, making it hard to set up a stable communication session between nodes. Dynamic routing protocols have been developed to tackle this problem. With the absence of an infrastructure, information can be sent to nodes that do not need it, making the use of the network less efficient. Especially with the challenge of interference that wireless networks face, this impact can be large. Another challenge is the lifetime of the network. As most devices are equipped with a battery, their lifetimes are limited, especially in sensor networks where there is only space for a small battery.

The communication over multiple nodes also poses problems. The hidden node problem occurs when a node is visible from one node, but not from other nodes of the network. Figure 1.3 shows an example of the hidden node problem. Node A transmits to node C, but node B cannot detect this transmission. Node B might also start transmitting to node C or another node, causing a collision at node C. The hidden node problem is a specific example of interference.

To diminish the impact of interference, protocols are active during communication between devices. The Medium Access Control (MAC) protocols determine which of the users of a network are allowed to use the medium. An example of a MAC protocol is CSMA/CA, Carrier Sense Multiple Access with Collision Avoidance. Using this protocol, a node that wants to transmit first senses if the network is free, i.e. no other transmissions are taking place. If this is the case, it starts transmitting. If there is a transmission going on, the node waits until the transmission is completed. Before transmitting, the node first sends a request-to-send (RTS) message. This is then received by all nearby nodes, so that they know they cannot transmit until this node is done. The receiving node sends a clear-to-send (CTS) message back so that the node knows it can start transmitting and hidden nodes also know this transmission will take place, even though they did not receive the original RTS. This way collisions can be prevented (the hidden node problem is avoided), at the cost of overhead. Also,

this approach presents the exposed node problem, as nodes might receive a CTS message and refrain from transmitting, even though there is need for them to do so. A different protocol is Time Division Multiple Access (TDMA) where time is divided into slots and these slots are assigned to different users. Many protocols and mechanisms are discussed in [Toh02]. The most commonly used specifications and settings for the MAC layer stem from the IEEE [IEE], where the IEEE 802.11 protocols are the best known for use in WLANs and ad hoc networks.

## 1.2 Methodologies for performance modeling and analysis of wireless ad hoc network

This section presents some basics of interest for the remainder of this thesis. The subsections describe the two fields that are used to model (wireless ad hoc) networks: graph theory and queuing theory. In the chapters that follow, the definitions presented here are in general not repeated but assumed to be known to the reader.

### 1.2.1 Graph theory

Graphs are used as an abstract representation of many different types of networks, including communication networks, transport networks, biological networks and social networks. Such an abstract representation of networks is very useful in order to identify and analyse all kind of structural properties, like connectivity and shortest paths. In the case of an ad hoc network, the users/devices of the network are the nodes in the graph and the communication links are depicted as edges of the graph. When communication is only possible in a certain direction, these edges are depicted as arrows, known as arcs or directed edges of the network. Connectivity in the network can be considered using the graph representation, where nodes are connected if they are within each other's transmission range. A path is a collection of edges that lead from one node to another.

Other characteristics can also be modelled as a graph, such as interference. An interference graph again uses the nodes to depict the users/devices, but now connects a node to another node when a transmission of the node causes interference for the other node. Additional information can be included in graphs, like assigning a value that depicts the capacity to the edges of the network. Or nodes can be given a value stating the number of radios it has available for transmission over different channels. For an end-to-end transmission over multiple hops, each edge in the path between the communicating nodes needs to have enough capacity and each node an available radio set to the appropriate channel to allocate the communication.

In this thesis we use graph theory in particular to study the maximum throughput that can be achieved between two nodes, a source and destination node, in the network by considering a graph where each edge has a certain capacity. The max-flow min-cut theorem of Ford and Fulkerson [FF56] provides this maximal throughput. It makes use of an imaginary line, a cut through the edges of the network, dividing the network into two parts, each part containing

either the source or destination node. As the capacity that can be achieved between the source and destination node is limited by the sum of the capacities of the edges that are cut, an upper bound on the throughput is acquired. By finding the cut that gives the lowest value (min-cut), you find the highest throughput (max-flow) that can be achieved by the network.

When multiple users want to communicate, this problem extends to the multi commodity flow problem (MCFP). The MCFP states a number of sources and destinations with their demands and poses the question if these demands can be accommodated by the edges with their given capacities. To solve this problem in an integer setting is extremely hard (NP-complete), but using linear programming it is possible to solve the problem for fractional flows. More constraints can be added to include other limitations, such as interference, that occur in wireless networks. Chapter 5 presents an approach to include interference constraints into the MCFP.

### 1.2.2   Queuing theory

Communication between wireless devices in a network takes place by packets being sent from one user to another. By modeling each user in an ad hoc network as a queue for these packets and the network as a server or servers that process these packets, we can identify and analyse many properties of the network. The order in which the packets are served and the time it takes to serve/transmit a packet are input parameters of the system. The state of a network is described by a vector with the number of packets in each of the queues and the state space of the system consists of all possible vectors. When the queues have a limited capacity to store packets, this state space is bounded, otherwise it is infinite. The system changes from one state to another due to arrivals and departures of packets after service. The queue lengths, the time it takes for a packet to reach its destination, the waiting time of packets before service, the busy time of the server and the throughput, which is the total amount of data the network can process per time unit, are performance metrics that can be calculated and all fall into the domain of queuing theory. We now briefly introduce and discuss some specific queuing models which are used in this thesis.

This thesis will consider Markov chains, where the transition from one state to another only depends on the current state, not on previous states. The transition from one state to another state in the discrete time Markov chain is given by the transition probability, or in the continuous time Markov chain by transition rates. The stationary (or steady state) distribution can be seen as the long run probability distribution of finding the system in a certain state.

The M/M/1 queue is the most basic queuing model where packets arrive according to a Poisson process and the service time is exponential. For this queue the performance measures noted earlier are well known. Jackson networks (cf. [Jac57],[Kel79]) are well known for their product-form stationary distribution, meaning that the stationary distribution of the system is the product of the stationary distribution of each of the nodes. These networks play an important role in Chapters 7 and 4. In a Processor Sharing queue (cf. [NnQ00],[KMR71],[FMI80]), a server does not serve one packet at a time, but its

capacity is distributed over multiple packets. A different amount of capacity can be allocated to different packets. Processor sharing plays an important role in Chapter 6. In a polling system (cf. [Lev90]) a server does not stay at a queue, but travels from queue to queue to process packets. As due to interference in an ad hoc network not all users can transmit simultaneously, this corresponds to users taking turns as is the case in a polling model. In a system with server vacations (cf. [FC85],[Kra89]), a server does not continually serve packets but may stop for an amount of time. From the perspective of a user in an ad hoc network, this corresponds to the user being allowed to transmit a certain amount of time, whereas due to interference vacations are imposed on the user, during which the user has to wait. Models that incorporate these properties are considered in Chapter 3.

## 1.3 Research questions and contribution

In a wireless ad hoc network, devices can transmit messages with a higher power to reach devices at a longer distance in one transmission or they can transmit with lower power and let other nodes forward their messages, which increases the number of transmissions that are needed. The lifetime of a network, the time until the first node depletes its battery, is modelled in Chapter 2. Using mean value analysis, we provide models for the lifetime distribution of a network where either nodes transmit at a power that ensures that all nodes receive the transmission or at a power such that only the nearest node receives the transmission. In the latter case, this node forwards the message to the next nearest node until the transmission is broadcasted over the complete network. In addition, networks where a number of nodes are denoted as master nodes, are analysed. In these networks nodes transmit to their designated master node, which forwards the message to the other master nodes. These master nodes then complete the final step by forwarding the message to all nodes in their designated section of the network. The models provide insight in the trade-off between power usage per transmission and the number of transmissions needed to distribute messages over the network. We show that the network size has an impact on the optimal choice, as for very small networks direct transmission provides a longer lifetime of the network than full routing.

Regardless of using direct communication between devices or letting other devices forward messages, the time it takes to complete communication, the end-to-end delay, has to remain limited. The network needs to be able to distinguish between different types of communications. This is why priorities can be set in a network for different users or different applications, for example by the use of parameters in protocols or by reserving channels for a certain type of communication. Chapter 3 addresses the aspect of delay and the impact of traffic prioritization. Considering the nodes as queues and the network as a server that visits these queues, the network is modelled by a polling system. The probability that a queue is visited differs due to the priority the traffic of the queue is given, which is considered to be either high or low. Using an iterative algorithm the average number of customers in each queue is calculated. This result is then used to determine the waiting time that packets of each

type of queue experiences. This provides valuable insight in the impact of QoS differentiation in networks and the level of prioritization that is needed to ensure a timely delivery of packets.

The impact that interference has on a wireless ad hoc network can be seen as a limiting factor on the rate at which nodes can transmit their data. When multiple nodes are active, the service rate of each active node decreases. Chapter 4 researches which arrival rates a two node ad hoc network can handle and how the different service rates affect the performance of the network. Building on known results, we provide insightful expressions for the stability range of the two node network of coupled queues. By providing conditions for which the network has a product-form distribution, we construct networks that are similar to the coupled queue network. Using a Markov reward approach, this enables us to provide bounds on the performance of the network. In addition, we show that allocating all capacity to one of the nodes provides better performance measures over sharing of the network capacity between the nodes.

As transmissions on the same frequency can cause interference and collisions, several approaches to prevent users from transmitting on the same frequency at the same time can be used to diminish the impact of interference. Dividing time into small frames or slots and assigning these slots to different users is one of them. The portion that is assigned to a user then defines the capacity allocated to this user. In Chapter 5 we research the maximum capacity that a network can achieve from a graph theoretic viewpoint, both for networks with one frequency channel and for networks where the nodes have multiple radios so that different channels can be used. By extending the multi commodity flow problem to include the impact of interference, we provide a theorem that gives sufficient and necessary conditions for a network to have enough capacity to satisfy a given demand of traffic to be transmitted from a number of sources to designated destinations. The use of the theorem provides insight in the location of bottlenecks in the network due to interference, enabling smarter channel allocation and network design.

A single flow of packets in an ad hoc network can also experience interference when it travels over multiple hops. The different nodes involved have to compete with each other to obtain the channel, meaning that part of the time nodes are waiting their turn. This influences the throughput the network can achieve. Chapter 6 researches the impact of the CSMA/CA protocol on the throughput of an ad hoc network from a queuing theoretic point of view. Taking into account the impact of the protocol on a packet level, the capacity allocated to a flow is determined. Considering the network on a flow level, we show that processor sharing models provide a good approximation of the throughput.

Routing has a large impact on the performance of the network. When too much traffic is routed through a single node, it may not be able to cope. Such a node is then labelled a 'bottleneck'. Even a single flow of packets with a large amount of data to be transmitted over multiple hops can cause a bottleneck to appear. The issue of bottlenecks in an ad hoc network is addressed in Chapter 7. Starting from a discrete time model that incorporates the contention between the active nodes of the network, a continuous time approximation is constructed with state dependent service rates. Considering long term average behaviour,

we determine state independent rates and show that in this case the network has a product-form distribution. This enables us to analyse the average queue length at each node, showing accurately where the bottlenecks of the network are located and at what offered load they appear. The model also correctly predicts the surprising result that increasing the offered load can change the location of the bottleneck. Predicting where bottlenecks occur plays a vital role in the deployment of ad hoc networks.

Overall, this thesis shows the high complexity of wireless ad hoc network analysis, even for small networks. Due to interference, which is shown to play a role in many different ways, the performance of wireless ad hoc networks is hard to analyse. Starting in Chapter 2 with the tradeoff between the number of hops used versus the power used per transmission, we show that the number of transmissions that a network can accomodate depends on the network design. The time it would take to actually perform all these transmissions depends on the impact of interference. As we show in Chapter 3, different types of traffic need to be considered as the total time need for a complete transmission may have to be limited. Even with only two types of traffic, the analysis is quite involved. Focussing in more detail on interference, Chapters 4,5,6 and 7 present different approaches to take the impact of interference into account. Where Chapter 4 provides a way to approximate many relevant performance measures, Chapter 5 uses graph theory to obtain bounds on the throughput of the network. Chapter 6 suggests that letting go of the intricate details on packet level of the effect that interference causes may be needed to make sure results can be obtained. Finally Chapter 7 uses numerous approximation steps to pinpoint the location where interference has the biggest impact. All in all, the wide variety of approaches presented in this thesis provide a good basis for further research, showing the difficulties that can be expected, providing interesting and relevant insights and obtaining results on important performance measures of wireless ad hoc networks.

## 1.4   Outline of the thesis

This chapter is closed by an outline of the remainder of this thesis, summarizing the results presented per chapter.

Chapter 2 analyses the lifetime of a network, which is defined as the time it takes until the battery of the first node is depleted. Two situations are considered: Direct transmissions between the source and destination or full routing where neighbouring nodes relay the traffic for each communication. For these settings the distribution of the network lifetime is determined. The trade-off between the number of transmissions and the distance bridged by each transmission is analysed. The nodes of the network are considered to be on a one dimensional grid or are uniformly distributed. We show that for nodes on a grid it is beneficial to use full routing. For uniformly distributed nodes, the number of nodes in the network determines which approach is better. For small networks, direct transmission outperforms the full routing approach. In this case, the longer distance that needs to be bridged weighs up against the increased number of transmissions that are needed. An intermediate approach,

choosing master nodes that forward data to other master nodes is simulated. Models for the expected lifetime are provided that give approximations which are close to the simulated results. The content of this chapter is based on the following paper:

- T.J.M.Coenen, J.C.W. van Ommeren and M. de Graaf. Routing versus energy optimization in a linear network, Workshop proceedings of the 23th International Conference on Architecture of Computer Systems, ARCS 2010, pp. 253-258, 2010.

Chapter 3 models the delay in a wireless ad hoc network using a polling model to take into account QoS differentiation in ad hoc networks. Traffic can have either high or low priority, determining the probability that a node is serving a packet. The delay experienced by packets of each class is analysed by considering each queue separately as being served by a server that takes holidays. The length of these holidays depends on the state of the system, making it hard to analyse them. An iteration algorithm, which is proven to monotonically converge, is presented to compute the waiting time distribution of a queue that uses the steady state for all other queues. Iterating over all queues provides de delay for packets at all queues, which gives accurate results for low to moderately loaded networks. The content of this chapter is based on the following paper:

- T.J.M.Coenen, J.L. van den Berg and R.J. Boucherie. Analysis of a polling system modeling QoS differentiation in WLANs, ValueTools'08 - Proceedings of the 3rd International Conference on Performance Evaluation Methodologies and Tools, 2008.

Chapter 4 combines results on product-form networks with a Markov reward approach to find bounds on any performance measure that is linear in each of the components of the state space. A two node network is considered where traffic can be forwarded from the first to the second node. When both nodes are active, the interference causes a lower service rate than when only one node is active. The stability range of the system is analysed, showing that increasing the rate at the boundaries of the system expands the stability range. Conditions for a geometric product-from solution are given which are used for comparison with the network under consideration. The Markov reward approach provides bounds for several performance measures, where we show that comparison with different product-form networks obtains different bounds. The content of this chapter is based on the following paper:

- T.J.M.Coenen, R.J. Boucherie and J. Goseling. Bounds on a two node network, submitted, 2016.

Chapter 5 analyses whether a network with a given traffic demand, capacities on each link and ranges of interference between the nodes can accommodate all the traffic demand. In the first part only one channel is available, so interference plays a large role in determining the throughput of the network. The network is modelled using a multi commodity flow problem and a theorem is stated that gives sufficient and necessary conditions for the problem to be solvable. For a

single source and destination pair the maximal throughput is computed using the max-flow min-cut theorem. The second part extends the results of the first part by including the option of using different channels. The theorem is extended to include these channels, giving a basis for an algorithm for channel allocation in wireless networks. The content of this chapter is based on the following papers:

- T.J.M.Coenen, M. de Graaf and R.J. Boucherie. An upper bound on multi-hop wireless network performance, Proceedings of the International Teletraffic Congress, ITC-20, 2007.

- T.J.M.Coenen, M. de Graaf and R.J. Boucherie. An upper bound on multi-hop multi-channel wireless network performance, Proceedings of Mobility'08, 2008.

Chapter 6 considers the throughput of ad hoc networks, taking into account the parameters involved in the CSMA/CA protocol with RTS-CTS in a wireless network. First, considering the packet level details, the aggregate system throughput is determined. Next, taking the flow level dynamics into account, the throughput is divided over all flows, taking into account the impact of multiple hops used in flows. This leads to two Processor Sharing models: Batch arrival processor sharing (BPS) and Discriminatory processor sharing (DPS). Simulation shows that the models provide an accurate estimation of the throughput for small networks. The content of this chapter is based on the following papers:

- T.J.M.Coenen, J.L. van den Berg and R.J. Boucherie. A flow level model for wireless multihop ad hoc network throughput, Proceedings of the 3rd International Working Conference on Performance modelling and Evaluation of Heterogeneous Networks HET-NETs '05, pp. 1-10, 2005.

- T.J.M.Coenen, J.L. van den Berg and R.J. Boucherie. Flow transfer times in wireless multihop ad hoc networks, Performance Modelling and Analysis of Heterogeneous Networks, pp. 113-132, 2009.

Chapter 7 considers the impact of node contention on the throughput in an ad hoc network. During each time slot the nodes of the network contend for the channel, depending on the protocol in use. Starting with a discrete time Markov chain we model the behaviour in the slotted time. To facilitate further analysis, we use long term average behaviour to model the discrete time Markov chain as a continuous time Markov chain, taking into account that certain nodes may be bottleneck nodes. The transition rates in this chain are state dependent, making it hard to analyse the network, so that further approximation is needed to obtain results on the throughput of the network. We approximate the continuous time Markov chain by a product-form network. This enables us to find the bottlenecks for a wireless network of any size and topology and to approximate its throughput. As the main result, an algorithm is provided that incorporates all these steps and gives very accurate results for the maximal throughput of the network. For a multihop tandem network a limiting result is obtained for the rate allocated to the first couple of nodes when all nodes continually want to transmit packets. The content of this chapter is based on the following paper:

- T.J.M.Coenen, J.L. van den Berg, R.J. Boucherie, M. de Graaf and A.M. Al Hanbali. Bottlenecks and stability in networks with contending nodes, International journal of electronics and communications (AEÜ), vol. 67, pp. 88-97, 2013.

# Routing versus energy optimization in a linear network

In wireless networks, devices (or nodes) often have a limited battery supply to use for the sending and reception of transmissions. By allowing nodes to relay messages for other nodes, the distance that needs to be bridged can be reduced, thus limiting the energy needed for a transmission. However, the number of transmissions a node needs to perform increases, costing more energy. Defining the lifetime of the network as the time until the first node depletes its battery, we investigate the impact of routing choices on the lifetime. In particular we focus on a linear network with two extreme cases where nodes send messages directly to all other nodes, or use 'full routing' where transmissions are only sent to neighbouring nodes. We distinguish between networks with nodes on a grid or uniformly distributed and with full or random battery supply. Using simulation we validate our analytical results on the lifetime distribution and discuss intermediate options for relaying of transmissions. We show that the size of the network is of influence on the optimal approach, as for very small networks it is optimal to use direct transmission over full routing.

## 2.1 Introduction

Mobile wireless networks are often battery powered which makes it important to maximize the network lifetime: batteries are (relatively) heavy, large, and sometimes difficult to replace. Here, the network lifetime is defined as the time until the first node depletes its battery. The broadcast network lifetime problem asks for settings of transmit powers and (node-dependent) sets of relay nodes, that maximize the network lifetime, under the assumption that all nodes originate broadcast traffic.

Literature in this area considers the lifetime maximization in mobile ad-hoc networks (MANETs). Often, the complexity is reduced by assuming transmissions originate from a single source (Kang and Poovendran [KP05], Pow and Goh [LG05] and Park and Sahni [PS07]). The related problem of minimizing the total energy consumption for broadcast traffic has also been widely studied, because it provides a crude upper bound to the lifetime of the network. Liang [Lia02] and Cagalj et al. [CHE02] have proven independently that minimizing the total transmitted power is NP-hard.

The contribution of this chapter is an (approximate) mean value analysis of two specific cases of this problem, for nodes located on a straight line. The

analyzed algorithms are the following: (1) direct transmissions (in which each nodes simply broadcasts its messages to all the other nodes, and no relaying takes place) and (2) full routing, where each message coming from a node is relayed by the neighbor(s) of that node. For these algorithms, we provide a framework for calculation of the probability distribution and expectation of the network lifetime. Through simulation we also consider the intermediate option of a fixed number of nodes that relay traffic, called master nodes, in designated sections of network. These master nodes receive transmissions and relay it to all nodes within their section and to neighbouring master nodes, thus distributing the transmission over the complete network.

This chapter answers a question that arose when considering the impact of routing on the network lifetime. With direct transmissions each node has few transmissions over a large distance. With full routing nodes perform a lot of transmissions over short distances. A priori, it is not clear which of the two approaches is the best for the network lifetime. This analysis provides insight in the network lifetime that can be gained by introducing (a form) of routing or master node selection which is directly relevant for radio networks. A more general interest lies in applications to Wireless Personal Area Networks (WPANs), and sensor networks. Here one could envisage a distinction between very simple devices (clients), and more powerful devices (eligible routers). From a theoretical viewpoint this analysis provides a stepping stone for further generalizations, mainly to the two dimensional case. Our results show that the network size influences the optimal choice for routing regarding the network lifetime.

## 2.2   General model and notation

In this chapter we investigate the effect routing has on the lifetime of the network. In [GO09] an analysis of networks with a single master node was presented under different master selection algorithms, including random selection, most centered, highest battery and optimal. For the random selection algorithm, we extend the work presented in [GO09] for different scenarios in a linear network. We distinguish the following scenarios:

1. *Direct transmission (DT)*: Each node transmits its message to all other nodes

2. *Full routing (FR)*: Each node transmits all messages only to its neighbouring nodes.

Next to analysing these scenarios analytically, we also investigate a scenario with master nodes (MN) by simulation. In this scenario a limited number of nodes are selected as master nodes to relay the transmissions over the network. The different scenarios are depicted in Figure 2.1, showing possible transmissions between nodes. In the case of direct transmission, the complete distance is bridged by a direct transmission, whereas in the case of full routing, multiple transmissions are made using direct neighbours to relay the transmission. In case master nodes are chosen (denoted by an M under the node), a node first transmits to the master node of it region, which relays the transmission to its

Figure 2.1: Possible transmissions for the direct transmission, full routing and master node scenarios

neighbour master nodes, which again will relay to its own neighbour master node and all nodes within its own region.

The lifetime of the network depends on the number of transmissions a node has to make, the distance it has to bridge and its battery supply. We distinguish between networks where all nodes have an equal (full) battery supply and where the battery has a random supply. To analyze the different scenarios, we use the following notation:

Consider a network with nodes $V$ which are distributed uniformly on the line [0,1] and let $|V| = n$. For a set $M \subseteq V$ of potential master nodes, a power assignment is a function $p : V \to \mathbb{R}$. Following the notation of [L+05], to each ordered pair $(u, v)$ of transceivers we assign a transmit power threshold, denoted by $c(u, v)$, with the following meaning: a signal transmitted by transceiver $u$ can be received by $v$ only when the transmit power is at least $c(u, v)$. We assume that $c(u, v) = \|u - v\|^2$ for all pairs $\{u, v\} \in V$. In the case of full routing, transmissions are only towards neighbouring nodes, whereas for direct transmission the transmission goes as far the furthest node. Each vertex is equipped with battery supply $b_v$, which is reduced by an amount $\lambda p(v)$ for each message transmission by $v$ with transmit power $p(v)$. Similarly, $b_v$ is reduced by amount $\mu r(v)$ for each message reception by $v$.

In our simplified analysis, we assume $\mu = 0$ (receive power is negligible), $\lambda = 1$ (by scaling), $E$ corresponds to a complete graph and each node transmits one message. In this case the only variables are the node locations and the initial battery levels: $G = (V, b)$.

For a node $v \in V$, let $p(v)$ denote the power assignment $p(v) : V \to \mathbb{R}$ defined as:

$$p(v) = \begin{cases} c(u, v) \text{ with } u = \arg\max_{w \in V}(|w - v|) \text{ for DT} \\ c(u, v) \text{ with } u = \arg\max_{w \in N(v)}(|w - v|) \text{ for FR}, \end{cases} \qquad (2.1)$$

where $N(v)$ denotes the neigbouring node(s) of node $v$.

Let $T_1, T_2, T_3, \ldots$ denote the time periods under consideration, where we assume that in each period, each node transmits once and all time periods have equal length. During a transmission all other nodes are silent until completion, so interference is not taken into account. We call a series of transmissions were each node transmits once a *round* and measure the lifetime of the network in rounds. As the order of transmission may not be known, the message lifetime, the number of messages sent until the first node depletes its battery, cannot be calculated exactly. The notion of rounds allows us to disregard the order in which the transmissions take place. Based on the stated assumptions, we obtain that after a round $r$ the battery supply is as follows:

$$b_v^{(r+1)} = \begin{cases} b_v^{(r)} - p(v) & \text{for all } v \in V \text{ for DT} \\ b_v^{(r)} - np(v) & \text{for all } v \in V \text{ for FR.} \end{cases} \tag{2.2}$$

Note that in case of full routing, we do not take into account the direction a transmission has come from. As of a received transmissions it may not always be known what the origin was, a node can not determine which nodes still need to receive it. Therefore, in our model, the node will always transmit to both neighbouring nodes. The network lifetime $L$, expressed in the number of rounds, can now be found as:

$$L = \begin{cases} \min_{v \in V}\left(\frac{b_v}{p(v)}\right) \text{ for DT} \\ \min_{v \in V}\left(\frac{b_v}{np(v)}\right) \text{ for FR.} \end{cases} \tag{2.3}$$

Summarizing, one can see that in full routing, each node only needs to transmit as far as its furthest direct neighbour, but the number of times this transmission needs to take place each round is equal to the number of nodes. Opposed to this, each node transmits only once in the case of direct transmission, but over a longer distance, using more energy per transmission. In the following we analyze and compare these two scenarios to determine their impact on the network lifetime.

## 2.3   Nodes on a grid

As an example of what we like to achieve, we first present an analysis of a network where all nodes are located on a grid. We consider both scenarios, direct transmission and full routing, both with nodes having a full battery capacity or a random one. The analysis of nodes situated on a grid provides insight in the impact of routing on the network lifetime when nodes can be tactically placed. Later we discuss the situation where nodes are uniformly distributed over the area to be covered. Obviously, the network lifetime is infinite when al nodes are positioned at the same location and thus no upper bound exists. Assuming that the complete network should be covered, the nodes are positioned on a grid, with equal distances between the nodes. Assuming that the first node is positioned at location 0 and the last one at 1, the remaining $n-2$ nodes are positioned with a distance of $\frac{1}{n-1}$ between them.

### 2.3.1 Direct transmission

When each node uses a direct transmission to all other nodes, the longest distance determines the lifetime of the network when each node has a full battery. The outer nodes have the longest distance to bridge and will deplete their battery supply in one round, which gives a lower bound on the network lifetime. When the battery supply at each node is randomly, i.i.d. distributed, it is not necessarily one of the outer nodes that depletes its battery supply first. The lifetime $L$ of the network, when the battery supply is uniformly distributed on $[c, 1]$, is then given by

$$
\begin{aligned}
P(L \geq t) &= P(\min_{v \in V} \frac{b_v}{D_v^2} \geq t) \\
&= \frac{1}{(1-c)} (\prod_{i=1}^{\lfloor \frac{n}{2} \rfloor} \min(1 - \frac{(n-i)^2 t}{(n-1)^2}, 1) \times \\
&\quad \prod_{i=\lfloor \frac{n}{2} \rfloor +1}^{n} \min(1 - \frac{(i-1)^2 t}{(n-1)^2}, 1)).
\end{aligned}
\tag{2.4}
$$

where the first (second) product denotes the first (second) set of nodes that has the longest distance to the last (first) node, which is a distance of $\frac{n-i}{n-1}$ (distance of $\frac{i-1}{n-1}$). This formula follows from the insight that node $i$ has a lifetime $L_i$ larger than $t$ that is given by

$$
P(L_i > t) = \frac{1}{1-c} \min(1 - \frac{(n-i)^2 t}{(n-1)^2}, 1)
\tag{2.5}
$$

as the probability that the battery has a certain capacity is given by $\frac{1}{i-c}$ and this capacity is used depending on the distance which is given by $\frac{n-i}{n-1}$. As the lifetime depends on the first node to deplete its battery, the lifetime of the network exceeds $t$ when all nodes have a lifetime that exceeds $t$.

### 2.3.2 Full routing

In case of full routing, each node has to bridge a distance of $\frac{1}{n-1}$, leading to a network lifetime of $\frac{(n-1)^2}{n}$ when all nodes have a full battery supply. With random battery supply, the node with the lowest battery supply determines the lifetime, which for $\frac{(n-1)^2}{n} c \leq t \leq \frac{(n-1)^2}{n}$ has the following distribution:

$$
\begin{aligned}
P(L \leq t) &= P(\min_{v \in V} b_v \leq \frac{nt}{(n-1)^2}) \tag{2.6} \\
&= 1 - (\frac{(n-1)^2 - nt}{(n-1)^2(1-c)})^n
\end{aligned}
$$

as the minimum battery capacity of $n$ nodes is distributed as

$$P(\min_{v \in V} b_v \le b) = 1 - (\frac{1-b}{1-c})^n,$$ (2.7)

which leads to an expected lifetime of

$$EL = \frac{(n-1)^2(nc+1)}{n(n+1)}.$$ (2.8)

## 2.4   Uniformly distributed nodes

As the position of the nodes often can not be chosen, we will analyze the network where all nodes are uniformly distributed over the region $[0, 1]$.

### 2.4.1   Direct transmission

When all nodes have a full battery, the lifetime of the network depends only on the distance the nodes have to bridge. When there are no master nodes, each node transmits its own message to all other nodes. The distance that needs to be bridged for this depends on the position of the nodes. Obviously, the outer nodes have the largest distance $D$ to bridge and will hence have the lowest lifetime. The probability density function of the distance $D$ between the outer nodes is given by

$$f_D(d) = n(n-1)d^{n-2}(1-d)$$ (2.9)

with an expected distance of $ED = \frac{n-1}{n+1}$. The lifetime distribution is given by

$$
\begin{aligned}
P(L \le t) &= P(\frac{1}{D^2} \le t) \\
&= 1 + (n-1)(\frac{1}{t})^{\frac{n}{2}} - (\frac{1}{t})^{\frac{n-1}{2}}
\end{aligned}
$$ (2.10)

and an expected lifetime $EL$ of

$$EL = \frac{(n-1)}{(n-2)(n-3)}$$ (2.11)

for $n \ge 4$ and infinity for smaller networks. The expected lifetime in number of rounds hence is decreasing in $n$, but the number of messages sent per round is increasing.

When the nodes in the network do not have the same battery supply, the nodes that have the longest distance to bridge will not necessarily be the ones to deplete their battery first. Even though the battery supply at each node is random and independent, the correlation between the distances between the nodes makes the analysis of this scenario much harder. Let $D_i$ denote the distance to the farthest node for node $i$ and $b_i$ it's battery power, then the

lifetime of the node is given by

$$L_i = \frac{b_i}{D_i^2} \tag{2.12}$$

and the lifetime of the network is given by

$$L = \min_i(L_i). \tag{2.13}$$

All the $B_i$ are independent, but the $D_i$ are not, complicating the analysis of $L$. We therefore analyze a worst case scenario, providing a lower bound on the network lifetime, by assuming that the node with the longest distance to bridge also has the lowest battery supply of all nodes in the network. As was known from the analysis with nodes having an equal capacity, the maximal distance $D$ between any two nodes is distributed as (2.9) and the minimum battery capacity as (2.7). As a bound on the network lifetime we hence obtain (for $\frac{b}{t} < 1$)

$$
\begin{aligned}
P(L \le t) \quad &\ge \quad P(\frac{\min(B_1, ..., B_n)}{D^2} \le t) \\
&= \quad \int_c^1 P(D \ge \sqrt{\frac{b}{t}}) \frac{n}{1-b}(\frac{1-b}{1-c})^n db \\
&= \quad \int_{\sqrt{\frac{b}{t}}}^1 \int_c^1 \frac{n^2(n-1)}{1-b} l^{n-2}(1-l)(\frac{1-b}{1-c})^n db dl
\end{aligned}
\tag{2.14}
$$

### 2.4.2   Full routing

**Theorem 2.1.** *The distribution of the lifetime of a network with full routing is given by*

$$
P(L \le t) = \begin{cases} \sum_{i=1}^{n-1}(-1)^{i-1}\binom{n-1}{i}(1-i\sqrt{\frac{1}{nt}})^n & for\ \frac{(n-1)^2}{n} \le t \le \infty, \\ \sum_{i=1}^{m-1}(-1)^{i-1}\binom{n-1}{i}(1-i\sqrt{\frac{1}{nt}})^n & for\ \frac{(m-1)^2}{n} \le t \le \frac{m^2}{n} \\ & and\ 1 < m < n. \end{cases}
\tag{2.15}
$$

*Proof.* When all nodes use full routing, the lifetime of the network is determined by the largest distance $D$ that needs to be bridged between two nodes. The probability that a gap of size $d$ exists between nodes can be found as follows. First, let $d \ge \frac{1}{2}$, so that only one such gap can be present. In this case we have that no nodes can be in an interval $d$. The probability that all nodes are not in this interval is given by $(1-d)^n$. This interval has to be somewhere between the nodes, for which there are $n-1$ choices (between $1^{st}$ and $2^{nd}$ until between $n-1^{st}$ and $n^{th}$). This gives for $d \ge \frac{1}{2}$ the probability of

$$P(D \ge d) = (n-1)(1-d)^n \tag{2.16}$$

Now let $\frac{1}{3} \le d \le \frac{1}{2}$. In this case there may be one or two gaps of size $d$. Using the reasoning above for there being at least one such gap gives expression (2.16), but we have to subtract all the situations where there are two such gaps as these

are counted double (once for each gap). When two gaps exist, all nodes are in an area $(1 - 2d)$ with probability $(1 - 2d)^n$. The two gaps need to be placed between the nodes, but not both between the same nodes (as otherwise $d \geq \frac{1}{2}$), which can be done in $\binom{n-1}{2}$ ways, leading for $\frac{1}{3} \leq d \leq \frac{1}{2}$ to the probability

$$P(D \geq d) = (n - 1)(1 - d)^n - \binom{n - 1}{2}(1 - 2d)^n, \tag{2.17}$$

assuming that $n \geq 3$, otherwise there couldn't be two gaps. In general this reasoning leads to

$$P(D \geq d) = \begin{cases} \sum_{i=1}^{n-1}(-1)^{i-1}\binom{n-1}{i}(1 - id)^n & \text{for } 0 \leq d \leq \frac{1}{n-1}, \\ \sum_{i=1}^{m-1}(-1)^{i-1}\binom{n-1}{i}(1 - id)^n & \text{for } \frac{1}{m} \leq d \leq \frac{1}{m-1} \text{ and } 1 < m < n. \end{cases} \tag{2.18}$$

Using this result, we get for the distribution of the lifetime $L$ of the network that

$$P(L \leq t) = P(D \geq \sqrt{\frac{1}{nt}}) \tag{2.19}$$

which leads to the lifetime distribution as stated in (2.15) in the theorem. $\quad\square$

The expected lifetime of the network follows from the distribution and is given by

$$EL = \frac{1}{n}\sum_{m=2}^{n-1}\int_{(m-1)^2}^{m^2}\sum_{i=1}^{m-1}\frac{(-1)^{i-1}\binom{n-1}{i}\sqrt{\frac{1}{t}}}{n}i(1 - \sqrt{\frac{1t^n}{)}}2dt \tag{2.20}$$

$$+\frac{1}{n}\int_{(n-1)^2}^{\infty}\sum_{i=1}^{n-1}\frac{(-1)^{i-1}\binom{n-1}{i}\sqrt{\frac{1}{t}}ni(1 - i\sqrt{\frac{1}{t}})^n}{2}dt.$$

When the battery supply is random, this again has a big impact on the analysis of the expected network lifetime and its dependence on the number of nodes. With a network consisting of more nodes, more messages will be sent per round and the probability of a node having a very low battery increases, which deteriorates the lifetime of the network. However, the distance that needs to be bridged may decrease, thus improving the lifetime of the network. We again analyze a worst case scenario. The battery supply of the node with the lowest supply is distributed as given in (2.7) and the lower bound on the lifetime of the network is given by

$$P(\frac{\min(b_1, .., b_n)}{\max(D_1, .., D_{n-1})^2} \leq t), \tag{2.21}$$

where the $D_i$ denote the distance between the $i^{th}$ and $i+1^{st}$ node in the network. The lifetime distribution is thus given by (for $\frac{b}{nt} < 1$)

$$P(L \leq t) \tag{2.22}$$

Figure 2.2: Expected lifetime for a network with nodes on a grid and uniformly distributed battery supply

$$
= \int_c^1 P(\max(D_1,..,D_{n-1}) \geq \sqrt{\frac{b}{nt}}) \frac{n}{1-b} \left(\frac{1-b}{1-c}\right)^n db \qquad (2.23)
$$

$$
= \begin{cases} \int_c^1 \sum_{i=1}^{n-1} (-1)^{i-1} \binom{n-1}{i} (1 - i\sqrt{\frac{b}{nt}})^n \frac{n}{1-b} (\frac{1-b}{1-c})^n db \\ \text{for } \frac{(n-1)^2}{n} \leq t \leq \infty, \\ \int_c^1 \sum_{i=1}^{m-1} (-1)^{i-1} \binom{n-1}{i} (1 - i\sqrt{\frac{b}{nt}})^n \frac{n}{1-b} (\frac{1-b}{1-c})^n db \\ \text{for } \frac{(m-1)^2}{n} \leq t \leq \frac{m^2}{n} \text{ and } 1 < m < n. \end{cases}
$$

## 2.5 Validation and discussion

For nodes situated on a grid, with uniformly distributed battery levels between $[0,1]$, the expected lifetime (in rounds) is as depicted in Figure 2.2 for both the scenario of direct transmission and full routing.

The figure shows that for the direct transmission scenario the lifetime of the network in general decreases as the number of nodes grows. This obviously is the case as adding nodes to the two nodes at the edge of the network can only decrease the lifetime of the network, as the outer nodes still need to transmit over the same distance. The increase in lifetime when going from 3 to 4 nodes is due to the change in the grid. It is better to have two nodes bridging a gap of $\frac{2}{3}$, than one node bridging a gap of $\frac{1}{2}$.

For the full routing, the addition of nodes is beneficial. In this scenario, adding a node decreases the distance that needs to be bridged, yet increases the number of transmissions. Apparently, the increase in number of transmissions is of lesser effect compared to the gain by decreasing the distance. The result for a network with two nodes takes into account that for full routing a node

Figure 2.3: Comparison of the model with simulation for the expected lifetime of the network for the scenarios of direct transmission and full routing

always resends a received transmission, thus the lower lifetime in the full routing scenario compared to the direct routing.

Plotting the results for uniformly distributed networks for the scenarios with direct transmission and full routing and comparing to simulation gives Figure 2.3. For readability of the upcoming plots, we from now on show an approximation of the message lifetime, that is $nEL$, with $EL$ the expected round lifetime as discussed.

The model and simulation are very close together, thus validating our results. For small networks ($n < 7$), it is better to use direct transmission than to use full routing, whereas for larger networks the opposite holds. When the network is very small, addition of a node will increase the number of transmissions per node, but the maximal distance that needs to be bridged needs not to be decreased significantly. The probability of all nodes being close together in a small network is high, leading to an infinite expected lifetime of networks smaller than four nodes. As the network gets bigger, the maximal distance to be bridged will go to 1 for the direct transmission, shown by the almost linear growth of the graph for larger $n$. For the full routing scenario, the increase is steeper as decrease in distance that needs to be bridged has a quadratic impact and the impact of the increase of messages to be sent is cancelled by considering the message lifetime. This reasoning already shows that for very large networks, full routing will always outperform any other scenario.

The lower bounds calculated for the scenarios with random battery supply (with $c = 0$) are depicted in 2.4. The lower bound calculated is not a good approximation for the expected lifetime, but a lot closer to the simulated result than for example the upper bound where all nodes have a full battery capacity. For the scenario with direct transmission approximating the expected lifetime

Figure 2.4: Lower bounds and simulation of the expected network lifetime for the scenarios of direct transmission and full routing with uniformly distributed battery supply

with the lower bound is more suitable than for the scenario with full routing. Interesting is the observation that full routing now outperforms direct transmission for any network size.

Next to the analyzed scenarios, one could also argue that an intermediate approach may be more suitable, chosing a set number of so-called master nodes that will relay the transmission for a certain region. In [GO09], the authors analyze networks with one master node, that receives all transmissions and then broadcasts them to all other nodes. The optimal choice of the master node is discussed, as well as randomly chosing a master node, chosing the node with the highest battery supply and the most centered node. When more master nodes are used, it makes sense to divide the network into sections, where the master nodes broadcasts received transmissions to all nodes in it's section and relays transmissions to neighbouring nodes as in the full routing scenario. Simulating networks with a fixed number of masters gives results as depicted in Figure 2.5 and Figure 2.6 for nodes with random and full battery supply.

As can be seen from the figures, the lifetime when using a fixed number of master nodes hardly depends on the size of the network. This is due to the fact that the master nodes have the largest distance to bridge and the most transmissions to send. Adding a (non-master) node hence has hardly any impact on the number of transmissions the master node can do. Only for a small network using as little master nodes as possible is optimal. For larger networks it holds that more master nodes results in a longer network message lifetime. For comparison, the results for direct transmission and full routing are included in the figures. Note, however, that the results for these settings assume a completely uniform distribution of the nodes over the interval $[0, 1]$, whereas

Figure 2.5: Comparison of the expected lifetime of a linear network with full battery supply for direct transmission, master node selection and full routing



Figure 2.6: Comparison of the expected lifetime of a linear network with uniformly distributed battery supply for direct transmission, master node selection and full routing

when using multiple masters, the assumption is taken that each section contains at least one node to be selected as master node. This explains for example why chosing 8 masters in a 8 node network gives a different result than using full routing, as the expected maximal distance to be bridged by a node is smaller

using the assumption that each section contains a node. For a network with 4 nodes, this assumption causes the expected maximal distance to be bridged to be higher than in a uniform distribution, causing the full routing scenario to give a better expected lifetime for this setting.

## 2.6 Conclusion

For two scenario's, direct transmission and full routing, the distribution of the network lifetime has been determined. When nodes can be placed on a grid, it is beneficial to use full routing as the increase in messages to be sent per round is compensated for by the decrease in the distance that needs to be bridged by each node. When nodes are uniformly distributed over the interval $[0, 1]$ however, the network size is of influence. In full routing, each node only needs to transmit as far as its furthest direct neighbour. However, the number of times this transmission needs to take place each round is equal to the number of nodes, opposed to only once in the case of direct transmission. For very small networks, it therefore is more energy efficient to use direct transmission opposed to full routing, as in this case the longer distances that may need to be bridged weigh up against the reduction of messages that need to be sent each round in comparison to full routing.

By simulation also scenarios with a limited number of relaying (master) nodes have been analyzed. A similar trend is then visible, that for smaller networks, it is optimal to use as little relaying nodes as possible, whereas for larger networks, the optimal choice is to use as many relaying nodes as possible. Possible future extensions are to analyze two dimensional networks, making the results presented here a good stepping stone for further generalizations.

# Analysis of a polling system modeling QoS differentiation in WLANs

This chapter models WLANs with QoS differentiation capability using a polling system with a random polling scheme, a 1-limited service discipline and deterministic service requirement. The system contains high and low priority queues that are distinguished via the probability of being served next. We propose a new iteration algorithm to approximate the waiting time of customers in the high and low priority queues. As shown by simulation results, our approximation is accurate for light to moderately loaded networks.

## 3.1 Introduction

Wireless Local Area Networks (WLANs) have become widely available for internet access and there is currently a growing demand for the support of other applications, in particular speech and video. Specific mechanisms then need to be deployed in order to provide appropriate QoS to the various applications. A typical approach to provide such QoS differentiation is for example by giving a larger share of the available capacity to preferred users, or giving priority to preferred classes. Introduction of such mechanisms requires insight into their performance. This chapter investigates the influence of prioritization of the packet delay handling at the Medium Access Control (MAC) layer in WLANs.

In IEEE 802.11 WLAN prioritization appears in the support of different QoS classes. These QoS classes are implemented via different settings of MAC layer parameters, like their access time, the maximum and minimum value for their back-off counter or the number of consecutive packets that may be transmitted, see [IEE05] for an overview of IEEE 802.11e that incorporates these mechanisms. QoS provisioning for IEEE 802.11 systems has been investigated mainly via discrete event simulations. Analytical models yielding robust insight into system behaviour are scarce. To a large extent, such models are based on the pioneering work of Bianchi [Bia00], in which a basic 802.11 system with persistent sources, i.e. sources that always have packets ready to be transmitted, is modeled and analysed using a Markov chain approach and validated via simulation showing excellent agreement with actual system behaviour. Extensions to include physical layer details are given in e.g. [HVS01],[W$^+$02]. The extension to non-persistent sources is provided in [CBvB05a],[L$^+$03], where a flow level model is introduced that is analysed using a Processor Sharing queueing model. Comparison with discrete event simulation shows that indeed the MAC layer can be adequately

modeled via the Processor Sharing mechanism. Extensions to multiple traffic classes with different QoS requirements, as e.g. in 802.11e, are among others presented in [Xia05],[XM06], [ZC03].

Although the flow level modeling of [CBvB05a],[L$^+$03],[Xia05], [XM06],[ZC03] captures the resource sharing behaviour of the MAC layer of 802.11 protocols, the essential behaviour at the packet level is not captured. At that level a flow consists of a series of packets that are transmitted one by one, where transmissions of different flows are intertwined. Especially for real time applications, such as speech/telephony, the packet level is of high importance. In [EOs05], a packet level analysis for non-persistent sources is presented, extending the Markov model of Bianchi to include the probability of the node going into an empty backoff state. We take a further step to analyze the packet level by modeling the MAC layer as a polling model where the server works off packets at different queues. The essential characteristics of the QoS aware MAC protocol are incorporated via the frequency at which the server visits the different nodes. In particular, we give the server a high probability of visiting a node with high priority packets.

In our polling model, we consider two types of queues, viz. high and low priority queues, each type with a different probability of the server moving to it. Upon departure from a queue, either after service of a packet or at the arrival of a packet to an empty system, the server randomly selects a queue according to these probabilities, which mimics the behaviour of the MAC layer in 802.11 systems. Note that we do not claim to accurately model the behaviour of the IEEE 802.11e protocol, but analyse a mathematically interesting model that provides insight into the effect of prioritization such as used in the IEEE 802.11 MAC layer. In our model, we will take the probability of moving to a high priority (HP) queue to be $\alpha$ times as high as moving to a low priority (LP) queue. The service time of a packet is considered to be deterministic as the packet sizes in the system are equal for all queues and the channel speed is assumed to be constant at all times. As a queue is only allowed to transmit one packet when obtaining the channel, the service discipline is 1-limited. This chapter analyzes the steady state waiting time for this 1-limited polling system with random polling.

For the 1-limited polling model, general results are available in literature. In [FC85] Fuhrmann and Cooper derive the well known decomposition result for queues with server vacations, which is very useful for analyzing polling models. For symmetric queues, so with identical arrival and service rates at the queue, and a cyclic polling order, [Fuh85] extends this result to give analytical results on the average waiting time of packets in the queues. In [BW89], Boxma gives a pseudoconservation law for the mean waiting time in a polling system with Markovian polling, that includes random polling. This law provides an exact expression for a weighted sum of the mean waiting times at all queues, which need not be symmetrical. However, results for individual queues cannot be derived from this law when the network is not symmetric.

The main contribution of this chapter is an analysis of the steady state marginal distribution of the waiting time of packets for different types of queues in a 1-limited asymmetric polling model. We consider the different queues in

the system individually and model a particular queue as a queue with server vacations, where these vacations depend on the state of the other queues. To obtain the steady state waiting time distribution, we propose an iteration algorithm. The algorithm computes the marginal steady state distribution of the number of packets at a tagged queue, assuming a steady state at all other queues. Iterating this approach over the queues, for various settings, we obtain the steady state waiting time distribution for packets at the different queues.

The remainder of this chapter is organised as follows. Section 2 describes the queueing networks under consideration and the analytical approach for determining the distribution of the waiting time of customers per queue. Numerical results of the proposed algorithm are compared with simulation in Section 3, and Section 4 concludes the chapter.

## 3.2 Model description and Analysis

Consider a polling model consisting of queues $Q_1, ..., Q_n$ with finite buffer $B$ and a single server $S$ visiting the queues. Customers arrive at a queue $Q_i$ according to a Poisson process with rate $\lambda_i$. The service process at the queues is deterministic with service time $\tau$ per customer and there is no switchover time between the queues. The routing policy for the server is random, meaning there is a probability $p_i$ that the server moves to queue $Q_i$ upon departure from queue $Q_j$, $j = 1, ..., n$. For a high priority queue, this probability is $\alpha$ times as high as for a low priority queue, that is $p_{HP} = \alpha p_{LP}$. The service policy is assumed to be 1-limited, meaning at most one customer is served at each visit of the server, and customers are served FCFS at each queue. When the server reaches an empty queue, it will immediately proceed to the next. When all queues are empty, the server waits at the last queue to instantly move to the first queue that receives a customer. To ensure stability of the system we assume that $\rho = \sum_{i=1}^{n} \lambda_i \tau < 1$.

In the following, we derive expressions for the average waiting time of a packet for both types of queue. We start by considering one high priority queue surrounded by $n$ low priority queues. The server will move to the HP queue with probability $\frac{\alpha}{n+\alpha}$ and to a certain LP queue with probability $\frac{1}{n+\alpha}$. We present an algorithm to approximate the waiting time of a packet for both types of queue. This algorithm considers queues separately as served by a server with vacations. The length of the vacations depends on the number of customers at the other queues. Starting with an arbitrary distribution of the number of customers at the other queues, the steady state of the number of customers in the considered queue is determined, using the vacation time distribution. This process is iterated over the different types of queues repeatedly, until convergence occurs. For specific cases, being that either the HP or LP queues are saturated, meaning they always have packets ready to be transmitted, exact results are presented. Exact results are also given for the case where all queues have equal priority.

### 3.2.1   General case

To determine the average waiting time of a packet in the queue, we consider the queues separately, as if they are in isolation. From the point of view of a queue, the server is either present and serving a packet, or away while serving an other queue. We thus can consider a queue as an M/D/1/B queue with vacations (c.f. [Fuh84],[Kra89],[Lee89]), where the absence of the server while serving other queues are the vacations. The length of these vacations, which depends on the number of customers at the other queues, influences the waiting time of the packets in the queue. For illustratory reasons, we first give the analysis for the scenario where there are two queues, one high priority and one low priority queue, which as we show in the subsequent subsection can be extended to any number of queues.

**Two queues**

In the two queue scenario, each queue can be considered separately as a queue with a server that goes on vacation. The duration of a vacation now depends on the state of the other queue. We approximate the distribution of the length of the vacation $V_x$, given the number of customers $N_y$ at the other queue (HP or LP) using the following recursion:

$$P(V_x = k\tau | N_y = i) = \tag{3.1}$$

$$q_y \sum_{j=i-1}^{B} P(V_x = (k-1)\tau | N_y = j) P(A_y = j - i + 1), \forall k \geq 1$$

$$P(V_x = 0 | N_y = i) = \begin{cases} 1, & i = 0 \\ (1 - q_y), & i = 1, ..., B \end{cases}$$

where $V_x$ is the length of the vacation seen by the queue $x$, $N_y$, $q_y$ and $A_y$ are the number of customers at queue $y$, the probability of the server polling queue $y$ and the number of arriving customers at queue $y$ during a service time, respectively. Note that the length of a service period is known to be $\tau$ due to the 1-limited service discipline, hence we will denote this as a service time. The variable $x$ can be the $HP$ or $LP$ queue and $y$ is the other type of queue. The vacation length distribution is then determined using

$$P(V_x = k\tau) = \sum_{i=0}^{B} P(V_x = k\tau | N_y = i) P(N_y = i), \ k \geq 0 \tag{3.2}$$

As the steady state distribution $P(N_y = i)$, is not known, we start with an arbitrary distribution, for example an always empty queue. Using this distribution, the vacation distribution for the other queue is obtained.

We derive the steady state distribution of the number of customers in the queue using the vacation time distribution, so that by using Little's law we acquire the expected waiting time of a packet. The queue under consideration can be seen as an M/D/1/B queue with vacations (c.f. [Fuh84],[Kra89],[Lee89]).

To analyze the steady state of this queue, we first focus on the state of the system at embedded points, which are after the departure of a customer or the end of a vacation. The probability $p_n$ that an embedded point is the completion of a service and the departing customer leaves $n$ customers behind, and the probability $q_n$ that an embedded point is a vacation termination with $n$ customers in the system are related in the following manner

$$p_n = \sum_{k=1}^{n+1} g_{n-k+1} q_k, \ n = 0, 1, .., B-2 \tag{3.3}$$

$$p_{B-1} = \sum_{k=1}^{B} g_{B-k}^C q_k \tag{3.4}$$

$$q_n = \sum_{k=0}^{n} h_{n-k} p_k + h_n q_0, \ n = 0, 1, .., B-1 \tag{3.5}$$

$$q_B = \sum_{k=0}^{B-1} h_{B-k}^C p_k + h_B^C q_0 \tag{3.6}$$

$$\sum_{n=0}^{B-1} p_n + \sum_{n=0}^{B} q_n = 1 \tag{3.7}$$

where $g_j$ and $h_j$ denote the probability of $j$ customers arriving during a service and vacation time, respectively, $g_j^C$ and $h_j^C$ denote the probability of $j$ or more customers arriving. As these probabilities are known, this set of equations can be solved, giving the steady state distribution at the end of an interval (either a service or vacation). To determine the continuous time steady state distribution, we note that the number of times a departing customer leaves a certain number of customers behind equals the number of times an arriving customer finds this number of customers in the system. We have to take into account, however, that an arriving customer can find $B$ customers in the system in which case the customer is discarded and leaves. Let $P_B$ denote the probability that an arriving customer finds the system full. To evaluate this expression, observe that

$$P_B = \frac{\rho - \rho'}{\rho} \tag{3.8}$$

where $\rho = \lambda\tau$, $\lambda = \sum_i \lambda_i$ is the offered load and $\rho'$ is the carried load,

$$\rho' = \frac{(1-b)\tau}{bEV + (1-b)\tau} \tag{3.9}$$

where $EV$ denotes the expected vacation time and $b$ denotes the probability that an embedded point is a vacation termination point,

$$b = \sum_{n=0}^{B} q_n. \tag{3.10}$$

Let $\sigma$ denote the multiplicative inverse of the average interval between consecutive embedded points, that is

$$\sigma^{-1} = bEV + (1 - b)\tau \qquad (3.11)$$

then

$$P_B = 1 - \frac{(1 - b)\sigma}{\lambda}. \qquad (3.12)$$

The queue length distribution at arrival epochs, $\pi_n$, $n = 0, ..., B$ is

$$
\begin{aligned}
\pi_n \;&=\; P(\text{Arrival sees } n \text{ packets—Arrival is accepted})(1 - P_B) &(3.13)\\
&\quad + P(\text{Arrival sees } n \text{ packets—Arrival not accepted})P_B &(3.14)\\
&=\; p_n(1 - P_B) + P_B 1(n = B) &(3.15)
\end{aligned}
$$

where

$$1(n = B) = \begin{cases} 1 \text{ if } n = B \\ 0 \text{ otherwise} \end{cases} \qquad (3.16)$$

Combining these results, we obtain

$$
\begin{aligned}
\pi_n \;&=\; \frac{(1 - b)\sigma}{\lambda} p_n, \; n = 0, 1, ..., B - 1 &(3.17)\\
\pi_B \;&=\; 1 - \frac{(1 - b)\sigma}{\lambda}.
\end{aligned}
$$

From PASTA we obtain that the continuous time steady state queue length distribution is given by $\pi_n$, $n = 0, ..., B$. Note that (3.17) requires the average vacation time $EV$, and (3.2) the distribution of the other queue to determine the vacation time distribution. We may iterate (3.2) and (3.17) to obtain an approximation of the steady state queue length distribution.

The algorithm approximates in each iteration the number of customers found at the other queue to determine the vacation time for the tagged queue. When this vacation time is underestimated, the server switches back early to the queue and starts servicing a packet at the considered queue (when available), thus leaving the server busy. When however the vacation time is overestimated, the approach leaves the server at the other queue for too long a period, where this queue might actually have become empty, thus leaving the server idle while it could process jobs in the tagged (non-empty) queue. The presented approach hence underestimates the capacity of the server, but equally for both queues. The average queue length of all customers in the total system, which for larger values of $B$ approximately can be seen as an M/D/1 queue as it is work conserving, is known and given by

$$EN_{total} = \frac{\rho(2 - \rho)}{2(1 - \rho)} \qquad (3.18)$$

where $\rho = (\lambda_{LP} + \lambda_{HP})\tau$, the load of the total system. The results obtained by the iteration give a higher average queue length due to underestimation of the server capacity. The queue length of each type of customer should hence be

---

**Algorithm 3.1** Algorithm to calculate the average number of customers at each queue

---

Calculation of $EN_x(It)$

1. Initialize
   $It := 1$, $x := 1$, $y := 2$
   $P(N_y = i) = \gamma_i$, $EN_i(0) = 0$ for $i = 0, ..., B$
   where $EN_i(j)$ denotes the average queue length of queue $i$ in iteration $j$

2. Determine the vacation time distribution at queue $x$ from (3.2), and $EV := EV_x$

3. Determine the queue length distribution $P(N_x = n) = \pi_n$, $n = 0, ..., B$, from (3.17) and determine the average queue length $EN_x(It)$

4. Set $y := x$, $x := 3 - y$ and repeat steps 2 and 3 for this setting

5. If $\frac{EN_x(It) - EN_x(It-1)}{EN_x(It)} < 0.01$ for both $x = 1, 2$, then STOP
   Else $y := x$, $x := 3 - y$, $It := It + 1$ Go to Step 2

---

scaled down, so that the average queue length of all customers in the system is correct. This leads to an improved estimation of the average queue length of a customer per type of queue. Using Little's law, we obtain the average waiting time for each type of queue.

The algorithm can start with an arbitrarily chosen steady state distribution for the queue length of the HP queue. From this, a new steady state is computed for the same queue. Starting from each initial distribution for the HP queue, Algorithm 3.1 converges to the steady state distribution. Theorem 3.1 below states that this convergence is monotone starting from either an empty or full HP queue using stochastic ordering. Let $X$ and $Y$ be random variables with distribution $F_X(.)$ and $F_Y(.)$, respectively. We say that $X \leq_{st} Y$ iff $F_X(x) \geq F_Y(x)$ for all $x \geq 0$ (c.f. [Ros96], p.410).

**Theorem 3.1.** *For each initial distribution, Algorithm 3.1 converges monotonically.*

*Proof.* Let $X_i^{HP}$ and $X_i^{LP}$ denote random variables for the queue length distributions of the HP and LP queue after the $i^{th}$ iteration and let $Y_i^{LP}$ and $Y_i^{HP}$ denote the random variables for the corresponding vacation length distributions. From (3.2) it follows that if $X_0^{HP} \leq_{st} X_1^{HP}$ also $Y_0^{LP} \leq_{st} Y_1^{LP}$ as a higher queue length for the HP queue leads to a longer vacation length for the LP queue. From (3.17) it follows that if $Y_0^{LP} \leq_{st} Y_1^{LP}$ also $X_0^{LP} \leq_{st} X_1^{LP}$ as a longer vacation for the server of the LP queue leads to a higher number of packets in the LP queue. Following the same reasoning for the LP node, we have that $X_0^{LP} \leq_{st} X_1^{LP}$ leads to $Y_0^{HP} \leq_{st} Y_1^{HP}$ and $Y_0^{HP} \leq_{st} Y_1^{HP}$ leads to $X_1^{HP} \leq_{st} X_2^{HP}$. It thus follows that $X_0^{HP} \leq_{st} X_i^{HP}$ for any $i \geq 1$ as long as $X_0^{HP} \leq_{st} X_1^{HP}$. Similarly we have that $X_0^{HP} \geq_{st} X_i^{HP}$ for any $i \geq 1$ as long as $X_0^{HP} \geq_{st} X_1^{HP}$.

From Theorem 3.1, an obvious approach is to start with

$$P(X_0^{HP} = n) = \begin{cases} 1 & \text{for } n = 0 \\ 0 & \text{for } n > 0 \end{cases} \qquad (3.19)$$

since it then holds that $X_0^{HP} \leq_{st} X$ for any $X$ with a non-negative distribution. Let $X_i^*$ denote the random variable following an equilibrium distribution, that is $X_i^* = X_{i+1}^*$. We then have that as $X_0^{HP} \leq_{st} X_i^*$, also $X_i^{HP} \leq_{st} X_i^*$, so the iteration process cannot jump past an equilibrium. In every iteration, the distribution may change, moving closer towards the equilibrium distribution. Similarly, we can start with the distribution

$$P(X_0^{HP} = n) = \begin{cases} 1 & \text{for } n = B \\ 0 & \text{for } n < B \end{cases} \qquad (3.20)$$

where $B$ is the maximum number of customers in the queue. It then follows that $X_0^{HP} \geq_{st} X$ for any $X$, so that after every step we have that $X_i^{HP} \geq_{st} X_i^*$ as $X_0^{HP} \geq_{st} X_i^*$. In this case every iteration takes a step closer to the equilibrium from above. Using Algorithm 3.1 starting from both (3.19) and (3.20), we find our approximation.

**Multiple queues**

The approach for two queues can easily be extended to multiple queues of any priority class. The vacation length of a considered queue then depends on the state of all the other queues, and can be computed by analogy to (3.2). The vacation length distribution in this case is given by

$$P(V_x = k\tau | N_y = i_y, y \neq x) = \qquad (3.21)$$

$$\sum_{y \neq x} \quad q_y \sum_{\substack{a_z = i_z, z \neq x, y \\ a_y = i_y - 1}}^{B} P(V_x = (k-1)\tau | N_z = a_z, z \neq x) \cdot$$

$$\prod_{z \neq x, y} P(A_z = a_z - i_z) \cdot P(A_y = a_y - i_y + 1)$$

$$P(V_x = 0 | N_y = i_y, y \neq x) = \frac{q_x}{q_x + \sum_{y, i_y > 0} q_y}$$

Here $q_x$ denotes the probability of the server jumping to queue $x$. The vacation length distribution is found using

$$P(V_x = k\tau) = \qquad (3.22)$$

$$\sum_{i_y = 0, y \neq x}^{B} P(V_x = k\tau | N_y = i_y, y \neq x) P(N_y = i_y, y \neq x)$$

where again the steady state queue length distribution of the other queues is needed. Starting again with a random distribution for all but one queue we

find the vacation time for this tagged queue and hence the corresponding steady state queue length distribution of this queue. This distribution can now be used for all queues of the same class and the other class can be analyzed using the steps of the algorithm. Note that the proof of convergence remains the same, as the analysis is done for each type of queue. In the case of multiple queues with balanced load, that is with identical arrival rates at the queues, the random variables $X_i^{HP}$, $X_i^{LP}$, $Y_i^{HP}$ and $Y_i^{LP}$ can be used for all queues of the same type as they are indentical. When arrival rates at the queues are different, the same reasoning can be used for all separate variables $X_i^{HP_j}$, $X_i^{LP_j}$, $Y_i^{HP_j}$ and $Y_i^{LP_j}$, where the subscript $j$ denotes a specific queue of the type HP or LP.

### 3.2.2   Special cases

For a high priority queue, it may be needed that a certain average waiting time can be guaranteed. To obtain the maximal average waiting time in a network with one HP queue and $n$ LP queues, we give results for the situation with saturated LP queues. To analyze the impact of prioritizing the high priority queue on the low priority queues, we compare the average waiting time at the LP queues without an HP queue in the system, with the case where the HP queue is saturated. For these special cases, exact results are available, which are given in this section.

#### Saturated LP queues

Consider one high priority queue with Poisson $(\lambda_{HP})$ packet arrivals and $n$ saturated low priority queues, i.e. $\lambda_{LP} \to \infty$. Let the probability $q$ of visiting the high priority queue be

$$q = \frac{\alpha}{n + \alpha} \tag{3.23}$$

where $\alpha$ denotes the factor of importance given to the high priority queue, meaning the probability of visiting the HP queue compared to the LP queue is $\alpha$ times as high. For the HP queue, the vacation length distribution is then given by the geometric distribution

$$P(V = k\tau) = (1 - q)^k q \tag{3.24}$$

as any time the server does not jump to the HP queue, it will service exactly one packet at an LP queue. As the average time between arrivals of the server is $\frac{\tau}{q}$ and the server only serves one customer at each visit, the HP queue is stable when $q > \lambda\tau$. With the exact distribution of the vacation length known, we can use the pgf of the number of customers in the queue as given by (3.17) to determine the average number of customers in the HP queue. The average waiting time then easily follows from Little's law.

#### Empty and saturated HP queue

We now consider the case where the low priority queues are no longer saturated, but each have an arrival process of rate $\lambda_{LP}$ and a deterministic service time of

value $\tau$. Let queue $n+1$ be the HP queue, the conservation law (c.f. [Gro90]) then states that

$$\sum_{i=1}^{n+1} \rho_i EW_i^q = \rho \frac{\sum_{i=1}^{n+1} \lambda_i \beta_i^{(2)}}{2(1-\rho)} \tag{3.25}$$

where $EW_i^q$ denotes the average waiting time in the queue (not including service) and $\rho_i = \lambda_{LP}\tau$ for $i = 1...n$ and $\rho_{n+1} = \lambda_{HP}\tau$, so that and $\rho = n\rho_{LP} + \rho_{HP}$. As the service time distribution is deterministic for any queue, we have that $\beta_i^{(2)} = \tau^2$ and the total waiting time of a customer is $EW_i = EW_i^q + \tau$.

Consider the case where there are only $n$ LP queues, so the arrival rate at the HP queue is set equal to zero. The stability condition is that $\rho = n\lambda_{LP}\tau < 1$ and it immediately follows that

$$\sum_{i=1}^{n} \rho_{LP} EW_{LP}^q = \rho \frac{\sum_{i=1}^{n} \lambda_{LP}\tau^2}{2(1-\rho)} \tag{3.26}$$

$$EW_{LP}^q = \frac{\rho\tau}{2(1-\rho)} \tag{3.27}$$

$$EW_{LP} = \frac{(2-\rho)\tau}{2(1-\rho)} \tag{3.28}$$

Now consider the case where the HP queue is saturated. We have $n$ identical LP queues, and from the perspective of the LP queues the server incurs a switchover time when it visits the HP queue. The stability condition for this system is that $\frac{n\lambda\sigma}{(1-\rho)} < 1$, where $\sigma$ denotes the mean switchover time, as this is the number of arriving customers during the average cycle time of a queue. Let $p_i$ denote the probability of jumping to queue $i$ and $s_i$ the average time it takes to switch to queue $i$. We have a pseudo-conservation law stating that (c.f. [BW89])

$$\sum_{i=1}^{n} \rho_i [1 - \frac{\lambda_i}{p_i} \frac{\sigma}{1-\rho}] EW_i = \tag{3.29}$$

$$\rho \frac{\sum_{i=1}^{n} \lambda_i \beta_i^{(2)}}{2(1-\rho)} + \frac{\sigma}{1-\rho} \sum_{i=1}^{n} \frac{\rho_i}{p_i} - \sum_{i=1}^{n} \rho_i s_i + \frac{\rho}{2\sigma} \sum_{i=1}^{n} p_i s_i^{(2)}, \tag{3.30}$$

where for our model we have that $\lambda_i = \lambda_{LP} = \lambda$, $\rho_i = \lambda\tau$, $\rho = n\lambda\tau$, $\beta_i^{(2)} = \tau^2$, $p_i = \frac{1}{n}$, $s_i = \frac{q\tau}{1-q}$, $s_i^{(2)} = \frac{q(q+1)\tau}{(1-q)^2}$ and $\sigma = s_i$ as all switchover times are equal. Here $q$ denotes the probability of the server polling the HP queue. As the LP queues are statistically identical, the expression simplifies to

$$EW_{LP} = \frac{\frac{n\lambda\tau^2}{2(1-n\lambda\tau)} + \frac{q\tau n}{(1-q)(1-n\lambda\tau)} + \frac{q+1-2q\tau}{2(1-q)}}{[1 - \frac{n\lambda q\tau}{(1-q)(1-n\lambda\tau)}]} \tag{3.31}$$

and applying Little's law the average total number of customers in the queue is obtained. Note that this approach can easily be extended to a case with multiple high priority queues, as only the probability of the server being on vacation

| | Rates | | Simulation | | Algorithm | | Error (%) | |
|---|---|---|---|---|---|---|---|---|
| $\alpha$ | $\lambda_{LP}$ | $\lambda_{HP}$ | LP | HP | LP | HP | LP | HP |
| 2 | 0.1 | 0.1 | 0.1129 | 0.1120 | 0.1125 | 0.1125 | 0.3455 | 0.4437 |
| 3 | 0.1 | 0.1 | 0.1132 | 0.1118 | 0.1166 | 0.1084 | 3.0128 | 2.9663 |
| 4 | 0.1 | 0.1 | 0.1133 | 0.1117 | 0.1174 | 0.1076 | 3.6018 | 3.6286 |
| 2 | 0.2 | 0.2 | 0.2723 | 0.2609 | 0.2829 | 0.2504 | 3.9085 | 4.0431 |
| 3 | 0.2 | 0.2 | 0.2753 | 0.2580 | 0.2911 | 0.2422 | 5.7375 | 6.1312 |
| 4 | 0.2 | 0.2 | 0.2771 | 0.2569 | 0.2961 | 0.2373 | 6.8574 | 7.6399 |
| 2 | 0.3 | 0.3 | 0.5623 | 0.4888 | 0.5918 | 0.4582 | 5.2475 | 6.2384 |
| 3 | 0.3 | 0.3 | 0.5792 | 0.4714 | 0.6253 | 0.4247 | 7.9612 | 9.9101 |
| 4 | 0.3 | 0.3 | 0.5881 | 0.4624 | 0.6454 | 0.4046 | 9.7323 | 12.5021 |

Table 3.1: Average waiting time in a two node network with balanced load

changes, so only the values of $s_i$ and $s_i^{(2)}$ need to be adjusted.

## 3.3 Validation

In the following we validate our approximation approach by comparison with simulation results. For a wide variety of settings, varying the load of the system and the grade of prioritization, the average waiting times of packets at the individual queues are determined. Note that the approach presented calculates the distribution of the waiting time, but only the averages are used in the following for comparison with simulation. Results for the scenario with one high priority and one or two low priority queues are considered, together with the special cases.

### 3.3.1 General case

**Two queues**

Table 3.1 shows the average waiting time of packets in a queue computed by the algorithm compared with simulation results for different loads of the system in the case of two queues, one HP and one LP queue. The table shows the impact of varying $\alpha$, the relative importance of the HP queue compared to a LP queue. The load at the queues is balanced, i.e. each queue has the same arrival rate of packets. The probability of moving to the HP queue is $q = \frac{\alpha}{n+\alpha}$, which is $\alpha$ times as high as for the LP queue and the buffer size is set to 15 for all cases. The impact of the differentiation appears to be higher when the load of the system increases. For a low load, the queues are often empty, thus making it possible for the server to attend to packets directly upon arrival. As the load increases, the queues will be fuller and the waiting time depends more on the frequency at which the server visits the queues. We observe that the accuracy of the algorithm deteriorates as the load of the system increases. For a highly loaded system, the queues will at times be fully loaded, causing arriving packets to be lost. This effect is not taken into account when using the pseudoconservation law to scale the obtained results. Simulation however

| Rates | | Simulation | | Algorithm | | Error (%) | |
|---|---|---|---|---|---|---|---|
| $\lambda_{LP}$ | $\lambda_{HP}$ | LP | HP | LP | HP | LP | HP |
| 0.2 | 0.5 | 0.4724 | 1.0469 | 0.5068 | 1.0098 | 7.2785 | 3.5384 |
| 0.5 | 0.2 | 1.1693 | 0.3475 | 1.2053 | 0.3113 | 3.0780 | 10.4126 |
| 0.1 | 0.01 | 0.1061 | 0.0107 | 0.1061 | 0.0106 | 0.0416 | 0.5174 |
| 0.01 | 0.1 | 0.0106 | 0.1062 | 0.0111 | 0.1057 | 3.9743 | 0.4905 |
| 0.4 | 0.1 | 0.6120 | 0.1384 | 0.6176 | 0.1324 | 0.9038 | 4.3436 |
| 0.1 | 0.4 | 0.1554 | 0.5934 | 0.1712 | 0.5788 | 10.1452 | 2.4608 |
| 0.1 | 0.3 | 0.1373 | 0.3966 | 0.1475 | 0.3858 | 7.3913 | 2.7269 |
| 0.3 | 0.1 | 0.4081 | 0.1286 | 0.4093 | 0.1240 | 0.3084 | 3.5340 |

Table 3.2: Average waiting time in a two node network with unbalanced load

shows that the impact of this approximation is limited, as the average number of packets in the system remains close to a system with infinite queues.

In a similar fashion Table 3.2 shows results for unbalanced arrival rates, with the probability $q = \frac{2}{3}$ ($\alpha = 2$) of visiting the HP queue kept constant. For more unbalanced situations, the results deteriorate, especially for higher loads. For the node with the lower arrival rate, the error made by the algorithm is bigger, as the average queue length is smaller. Comparing the impact of increasing the load of the LP queue on the HP and vice versa shows that the increase in load of the HP queue has a bigger impact on the average waiting time at the LP queue than increasing the load of the LP queue has on the HP queue. As an increase of the load will cause the queue to be non-empty for a larger fraction of the time, the impact it has on the other queue by causing the server to go on a vacation becomes larger. As a HP has a higher probability of being visited, increasing the load of this queue has a bigger impact than increasing the load at the LP queue.

**Three queues**

In Table 3.3 we consider the scenario with three queues, one HP queue and two LP queues. The table shows the average waiting time of packets computed by the algorithm compared with simulation results for the situation with balanced load. As for the situation with two nodes, we observe that for higher loads, the impact of the prioritization increases. Again, the results deteriorate as the load of the system increases. Comparison with the results of Table 3.1 furhter shows that the impact of prioritization is higher when more nodes are active in the network. The decrease in the average waiting time of customers for the HP queue is stronger relative to the decrease for the two node situation. With more queues present, the relative increase in probability of being visited is higher when the value of $\alpha$ is increased. For example, increasing the value of $\alpha$ from 2

|  | Rates | | Simulation | | Algorithm | | Error | |
|---|---|---|---|---|---|---|---|---|
| $\alpha$ | $\lambda_{LP}$ | $\lambda_{HP}$ | LP | HP | LP | HP | LP | HP |
| 2 | 0.1 | 0.1 | 0.1217 | 0.1205 | 0.1245 | 0.1151 | 2.386 | 4.465 |
| 3 | 0.1 | 0.1 | 0.1228 | 0.1188 | 0.1261 | 0.1122 | 2.639 | 5.556 |
| 4 | 0.1 | 0.1 | 0.1229 | 0.1185 | 0.1269 | 0.1104 | 3.323 | 6.868 |
| 2 | 0.2 | 0.2 | 0.3654 | 0.3202 | 0.3766 | 0.2928 | 3.054 | 7.300 |
| 3 | 0.2 | 0.2 | 0.3711 | 0.3072 | 0.3882 | 0.2736 | 4.602 | 10.949 |
| 4 | 0.2 | 0.2 | 0.3749 | 0.2986 | 0.3946 | 0.2608 | 5.250 | 12.678 |
| 2 | 0.3 | 0.3 | 1.9029 | 0.9133 | 2.0479 | 0.8543 | 7.621 | 6.461 |
| 3 | 0.3 | 0.3 | 2.0098 | 0.7526 | 2.1639 | 0.6221 | 7.669 | 17.337 |
| 4 | 0.3 | 0.3 | 2.0653 | 0.6844 | 2.2162 | 0.5177 | 7.304 | 24.356 |

Table 3.3: Average waiting time in a three node network with balanced load

to 3 for both situations gives the following relative increase (r.i.):

$$
\begin{array}{cccc}
 & \alpha = 2 & \alpha = 3 & \text{r.i.} \\
\text{2 nodes} & q = \dfrac{2}{3} & q = \dfrac{3}{4} & 12.5\% \\
\text{3 nodes} & q = \dfrac{1}{2} & q = \dfrac{3}{5} & 20\%
\end{array}
\tag{3.32}
$$

For all settings, no more than 15 iterations were needed by the algorithm with the accuracy set in such a way that the last step gave an improvement less that 1%. Longer runs with higher accuracy did not improve the results significantly. To run the iterations, the values of $P(V_x = k\tau | N_y = i)$ for the two node case and $P(V_x = k\tau | N_y = i_y, y \neq x)$ for the three node case had to be computed once using the iterations given in (3.1) and (3.21), which is time consuming for large values of the buffer sizes. For highly filled buffers however, the geometric distribution can be used, as the probability of the vacation having a duration of $k\tau$ is then very close to the probability of first visiting $k$ other queues before visiting the considered queue, as the other queues will not become empty during the process. The time needed for the iteration itself is very limited, as (3.2) (or (3.22)) only encompasses the addition over all possible values of queue lenghts and (3.17) is a small enough system of equations to be solved within seconds.

### 3.3.2 Special cases

**Saturated low priority queues**

For a user with important traffic, the QoS differentiation is of high importance. To get an idea of the impact of the settings for the differentiation, a worst case scenario can be analysed to see the minimal prioritization that is needed to obtain a certain average waiting time for the high priority packets. The worst case scenario is when all other (low priority) queues always have traffic to transmit. Figure 3.1 (left) shows the average waiting time of a packet in the HP queue, for different values of $n$, the number of saturated low priority queues in the system. The arrival rate at the HP queue is set to $\lambda_{HP} = 0.01$. The three

Figure 3.1: Average waiting time for the scenario with saturated LP and saturated HP queues

lines represent the results of the model for $\alpha = 2...4$, the grade of prioritization. It clearly follows from the figure that where for a sparse network (low number of LP queues) the differentiation has a limited effect and that for a dense network (high number of LP queues) giving more priority has a much bigger impact.

**Empty and saturated high priority queue**

The differentiation between users is primarily done to provide better performance for more important traffic. However, it also has to be taken into account what the impact is on performance of the less important traffic. If the prioritization of the high priority queue is too high, the low priority queues might be starved. To analyse the impact on the low priority queues, we compare the situation without the HP queue (or an empty HP queue) with the situation that the HP queue always has packets to transmit. In the latter case, we vary the grade of prioritization. Figure 3.1 (right) shows the average waiting time of a packet in an LP queue, for different values of $n$, for different settings of the HP queue. The arrival rate $\lambda_{LP}$ is set to 0.01 for each of the $n$ LP queues. In this case the HP queue is either absent (or empty) in which case the complete network behaves as a standard M/D/1 queue where each separate queue has the same average behaviour or the HP is saturated, with different values for $\alpha$, the grade of prioritization. For higher values of $\alpha$ the server will more often be processing HP packets, leaving less capacity for the LP queues. This shows from the figure as the waiting time reaches high values already for lower values of $n$. When the network is sparse, we see there is already a substantial impact of the differentiation on the waiting time of the low priority packets.

## 3.4   Conclusion

In this chapter we analyzed the impact of QoS differentiation on the delay of packets for different classes of queues using a 1-limited polling model with a random scheduling policy and deterministic service times, capturing the random nature of the MAC layer protocol. The model gives insight in the effect of the parameter settings on the QoS in a WLAN for the individual classes of queues. We developed an approximation approach for the packet delay in a network with high and low priority queues. Comparison with simulation results shows that for low to moderately loaded systems, the approach works well. Our results provide a tool for network designers to determine the level priority that is needed to ensure a certain expected waiting time of a customer.

# Bounds for linear performance measures in a two node network

We consider a queueing network with two nodes in which the servers of these nodes are coupled in the sense that the service rate of a server is determined by the presence or absence of packets in the other queue. We present stability conditions and conditions for which the network has a geometric product-form stationary distribution. Using a Markov reward approach we establish error bounds on various steady-state performance measures of this network. The basic idea is to compare the performance with the performance of a perturbed process that has a geometric product-form stationary distribution. Additionally, the impact of the allocation of the service capacity to the nodes is analysed, showing that it is optimal to allocate all of the capacity to one node.

## 4.1 Introduction

We consider a queueing network with two nodes in which the servers of these nodes are coupled in the sense that the service rate of a server is determined by the presence or absence of packets in the other queue. More precisely, let $\mu_1^*$ and $\mu_2^*$ denote the service rate at queue 1 and 2, respectively, if the other queue is empty. If both queues are non-empty these service rates reduce to $\mu_1 < \mu_1^*$ and $\mu_2 < \mu_2^*$. The queueing network with two coupled nodes was first studied in [FIM99] by considering two parallel M/M/1 queues with coupled service rates. In this paper we generalize this model by allowing for forwarding between these queues, if a packet completes service at one queue it joins the other queue with some probability.

Applications of this model arise, for instance, in a wireless network with two nodes whose transmissions cause mutual interference, reducing the service rate if both nodes are active. More generally, in a communication network where several nodes need to share resources, service rates will be reduced if several nodes are active. The reduction is caused either directly by interference, or by the overhead introduced by protocols to mitigate this interference. Note that such protocols leave considerable freedom in how to allocate rates $\mu_1$ and $\mu_2$ to the servers. Therefore, part of our interest will be in how to allocate $\mu_1$ and $\mu_2$ to the servers in order to minimize certain steady-state performance measures.

We will model the network as a random walk in the quarter-plane, a model which has been extensively studied in, for instance, [FIM99, CB83]. The approach taken in [FIM99, CB83, FI79] is to find an expression for the generating function

of the stationary distribution of the process by formulating it as the solution of a boundary value problem. The form of the results that can be obtained using this approach do not provide the insight that is required to analyze, for instance, optimal allocation of service rates. This is illustrated by [RO03] in which a tandem queue with coupled processors (a special case of our model in which there is forwarding with probability one in one direction only) is analyzed using this analytical approach, but no additional insights about the system behavior are reported. Therefore, we follow another path in this paper.

We use the Markov reward approach to establishing error bounds and comparison results as developed by van Dijk [vD11] to develop bounds on performance as well as to obtain insight into optimal rate allocation strategies. Van Dijk and Puterman introduced the Markov reward approach in [vDP88, vD88]. The method was further refined by van Dijk and applied to, for instance, Erlang loss models in [BvD09]. An overview of this method is presented in [vD11]. The Markov reward approach compares two queueing networks modeled as Markov chains, one for which the stationary distribution is not known, and another one, which is a modification of the first one, for which the stationary distribution is known. In [GBvO13] an approach is presented to establish error bounds and comparison results as the solution of a linear program, making use of the Markov reward approach of van Dijk.

The contributions of this paper are as follows. First, we provide necessary and sufficient conditions on the arrival rates for which there exists a rate allocation scheme that results in a stable network. These results are based on existing stability results for random walks in the quarter-plane [FIM99, Miy11]. We present results that provide significantly more insight for the special case of a coupled queue with forwarding.

Second, conditions are given for which the network has a geometric product-form stationary distribution. This extends results of Bayer and Boucherie [BB02] who consider very specific boundary behavior. A product-form characterization for random walks in the quarter-plane is given by Latouche and Miyazawa in [LM14]. We demonstrate that a coupled queue with forwarding has a geometric product-form distribution if and only if $\mu_1 = \mu_1^*$ and $\mu_2 = \mu_2^*$. This means that the Jackson network where the rates at the boundary are equal to the rates in the interior is the only setting that ensures this product-form. Also it means that the result of Fayolle et al that $\mu_1 + \mu_2 = \mu_1^* + \mu_2^*$ is sufficient and necessary if there is no forwarding does not extend to the case with forwarding.

Third, we bound the performance of networks that do not have a geometric product-form stationary distribution by establishing an error bound with a slightly perturbed network that does have a geometric product-form stationary distribution. This provides bounds on several performance measures, where examples are given for the probability that the system is empty and for the average number of customers in each queue. Finally, the impact of the allocation of the service capacity to the nodes is analysed, showing that it is optimal to allocate all capacity to one node.

Our work is related to [BJL08] where parallel systems with coupled service rates are studied (a special case of our model in which there is no forwarding) in which the service rate is a function of the number of packets in the other queue.

Sufficient and necessary conditions for the stability are derived and it is shown that these conditions are sharp when the service rate at each queue is decreasing in the number of customers in other queues, and has uniform limits as the queue lengths tend to infinity. Also, in [YH91] a packet radio network with two nodes is considered and it is analyzed which control structure leads to a geometric product-form stationary distribution and, therefore, the expected throughput and the expected packet delay. Yeh [Yeh02] incorporates both network layer and physical layer issues in his model to find a delay-optimal rate allocation in an M-user additive Gaussian noise channel. He extends the work of Telatar and Galager [TG95] by allowing packets to queue at the transmitters and a Poisson arrival process of packets instead of all packets being present at the start. It is shown that in the symmetric case, the Longer-Queue-Higher-Rate (LQHR) allocation is optimal.

The remainder of this paper is organized as follows. In Section 2 we introduce the model and Section 3 discusses the stability of the network. Next, Section 4 gives examples of useful reference networks that have a product-form stationary distribution and in Section 5 we use a Markov reward approach to obtain bounds on performance measures of the network. Section 6 analyses the impact the allocation of the capacity over the two nodes has on the performance of the network. Section 7 finally concludes the paper.

## 4.2   Model and problem statement

Consider a two node wireless queueing network with packets arriving at node $i = 1, 2$ according to a Poisson process with rate $\lambda_i$. Both nodes have an infinite size buffer. If both nodes have packets in their buffer packets are served according to a first come first served (FCFS) discipline from node $i = 1, 2$ at rate $\mu_i$. If only node $i$ has a non-empty buffer it serves at rate $\mu_i^* > \mu_i$. All service times are exponentially distributed. After service completion at node $i$ a packet is forwarded to the other node with probability $p_i$ and it leaves the system with probability $1 - p_i$.

We model the network as a continuous time Markov chain on state space $S = \{0, 1, \dots\}^2$, where $n = (n_1, n_2) \in S$ represents the state in which there are $n_i$ packets at node $i$. Let $e_i$ denote the unit vector with a 1 on position $i$ and for later use let $d_1 = (-1, 1)$ and $d_2 = (1, -1)$. The transition rates for this model are as follows:

$$
\begin{aligned}
n \to n + e_i \quad &\text{with rate} \quad \lambda_i \text{ for } i \in \{1, 2\}, \\
n \to n - e_i \quad &\text{with rate} \quad (1 - p_i)\mu_i \text{ for } i \in \{1, 2\}, n_1 > 0 \text{ and } n_2 > 0, \\
n \to n - e_i + e_j \quad &\text{with rate} \quad p_i\mu_i \text{ for } i \neq j, n_1 > 0 \text{ and } n_2 > 0, \qquad (4.1) \\
n \to n - e_i \quad &\text{with rate} \quad (1 - p_i)\mu_i^* \text{ for } i \neq j, n_i > 0 \text{ and } n_j = 0, \\
n \to n - e_i + e_j \quad &\text{with rate} \quad p_i\mu_i^* \text{ for } i \neq j, n_i > 0 \text{ and } n_j = 0.
\end{aligned}
$$

The transition rates are depicted in Figure 4.1. We refer to this model as the coupled queue with forwarding $R$. For future reference, we also introduce a generalized version of the network as depicted in Figure 4.2 with different rates

Figure 4.1: Transition rates for the coupled queue with forwarding $(R)$



Figure 4.2: Transition rates for the generalized
coupled queue with forwarding $(G)$

at the boundaries of the network. We will refer to this network as $G$. We assume that all the networks that we consider are aperiodic and irreducible.

In the following we denote the transition rate from state $n$ to $n + ke_1 + le_2$ as $q_{k,l}(n)$, $k, l \in \{-1, 0, 1\}$ and let $\pi(n)$ denote the steady state probability of the network being in state $n$. For notational convenience, let $\gamma_1$ and $\gamma_2$ denote the overall arrival rate at nodes 1 and 2, respectively, i.e. $\gamma_1$ and $\gamma_2$ satisfy

$$\gamma_1 = \lambda_1 + \gamma_2 p_2, \quad \text{and} \quad \gamma_2 = \lambda_2 + \gamma_1 p_1. \tag{4.2}$$

It follows that $\gamma_1 = (\lambda_1 + \lambda_2 p_2)/(1 - p_1 p_2)$ and $\gamma_2 = (\lambda_2 + \lambda_1 p_1)/(1 - p_1 p_2)$.

Let $\mu_c$ denote the total capacity of the system

$$\mu_1 + \mu_2 = \mu_c. \tag{4.3}$$

We analyse the stability of coupled queue network $R$ and as performance measures

we are interested in the probability $\pi(0,0)$ that the system is empty and the average number of customers in each queue. To obtain results we will establish error bounds between R and an appropriately chosen G network. In particular, we find conditions for $G$ to have a product-form stationary distribution. We also analyse the impact of the allocation of the capacity $\mu_c$ to the two nodes of $R$ on the performance measures, i.e. the probability of the system being empty and the average queue length at each node.

## 4.3 Stability

In this section we analyze stability of $R$, the coupled queue with routing. More precisely, we consider necessary and sufficient conditions under which this process is positive recurrent. Note that we do not consider stability of the generalized coupled queue process. We will see in Section 4.4 that stability of the generalized processes that we will consider will follow from our other results.

Necessary and sufficient stability conditions for $R$ follow from the general results for random walks in the quarter-plane [FIM99, Miy11]. Before presenting the result we define $(M_x, M_y)$, $(M'_x, M'_y)$ and $(M''_x, M''_y)$ as the drift in the interior of the state space, at the horizontal axis and at the vertical axis, respectively. More precisely, let

$$M_x = \lambda_1 + p_2\mu_2 - \mu_1, \qquad M'_x = \lambda_1 - \mu_1^*, \qquad M''_x = \lambda_1 + p_2\mu_2^*,$$
$$M_y = \lambda_2 + p_1\mu_1 - \mu_2, \qquad M'_y = \lambda_2 + p_1\mu_1^*, \qquad M''_y = \lambda_2 - \mu_2^*.$$

The result below appears in [Miy11].

**Theorem 4.1 ([Miy11], Lemma 6.4).** *R is positive recurrent if and only if one of the following three conditions holds:*

1. $M_x < 0$, $M_y < 0$, $M_xM'_y - M_yM'_x < 0$ and $M_yM''_x - M_xM''_y < 0$,

2. $M_x \geq 0$, $M_y < 0$ and $M_xM'_y - M_yM'_x < 0$ and $M''_y < 0$ for $M''_x = 0$,

3. $M_y \geq 0$, $M_x < 0$ and $M_yM''_x - M_xM''_y < 0$ and $M'_x < 0$ for $M'_y = 0$.

Theorem 4.1 extends the results in [FIM99] to include $M''_x = 0$ and $M'_y = 0$ in the second and third case, respectively.

Our first result deals with stability of $R$ for a given service rate allocation $\mu_1$ and $\mu_2$. It is expressed in terms of $\gamma_1$ and $\gamma_2$, which we recall are defined in Section 4.2 as the overall arrival rate at node 1 and 2, respectively.

**Theorem 4.2.** *R is positive recurrent if and only if*

$$\gamma_1 \leq \mu_1 \quad and \quad \frac{\gamma_1}{\mu_1}\left(1 - \frac{\mu_2}{\mu_2^*}\right) + \frac{\gamma_2}{\mu_2^*} < 1 \tag{4.4}$$

*or*

$$\gamma_1 > \mu_1 \quad and \quad \frac{\gamma_1}{\mu_1^*} + \frac{\gamma_2}{\mu_2}\left(1 - \frac{\mu_1}{\mu_1^*}\right) < 1. \tag{4.5}$$

*Proof.* First we express the stability conditions of Theorem 4.1 in terms of the following functions in $\gamma_1$ and $\gamma_2$:

$$f_1(\gamma_1, \gamma_2) = \gamma_1 - \mu_1 - p_2(\gamma_2 - \mu_2), \qquad (4.6)$$

$$f_2(\gamma_1, \gamma_2) = \gamma_2 - \mu_2 - p_1(\gamma_1 - \mu_1), \qquad (4.7)$$

$$f_3(\gamma_1, \gamma_2) = \gamma_2(\mu_1^* - \mu_1) + \mu_2(\gamma_1 - \mu_1^*), \qquad (4.8)$$

$$f_4(\gamma_1, \gamma_2) = \gamma_1(\mu_2^* - \mu_2) + \mu_1(\gamma_2 - \mu_2^*). \qquad (4.9)$$

Theorem 4.1 gives that $R$ is positive recurrent if and only if one of the following three conditions holds:

C1 $f_1(\gamma_1, \gamma_2) < 0$, $f_2(\gamma_1, \gamma_2) < 0$, $f_3(\gamma_1, \gamma_2) < 0$ and $f_4(\gamma_1, \gamma_2) < 0$,

C2 $f_1(\gamma_1, \gamma_2) \geq 0$, $f_2(\gamma_1, \gamma_2) < 0$ and $f_3(\gamma_1, \gamma_2) < 0$,

C3 $f_1(\gamma_1, \gamma_2) < 0$, $f_2(\gamma_1, \gamma_2) \geq 0$ and $f_4(\gamma_1, \gamma_2) < 0$,

i.e. positive recurrence is defined in terms of the half spaces induced by $f_1, \ldots, f_4$. In the above conditions we have excluded the cases $M_x'' = 0$ and $M_y' = 0$, because $M_x'' > 0$ and $M_y' > 0$ due to $\lambda_1 > 0$ and $\lambda_2 > 0$, respectively. The result of this theorem states that the above three conditions reduce to

$$f_3(\gamma_1, \gamma_2) < 0 \text{ if } \gamma_1 > \mu_1 \text{ and } f_4(\gamma_1, \gamma_2) < 0 \text{ if } \gamma_1 \leq \mu_1 \qquad (4.10)$$

and this remains to be shown.

For sufficiency of (4.10) observe that

$$f_1(\mu_1, \mu_2) = f_2(\mu_1, \mu_2) = f_3(\mu_1, \mu_2) = f_4(\mu_1, \mu_2) = 0. \qquad (4.11)$$

$$f_1(0, \mu_2 - p_1\mu_1) = f_2\left(0, \mu_2 - \frac{\mu_1}{p_2}\right) = f_3\left(0, \frac{\mu_1^*\mu_2}{\mu_1^* - \mu_1}\right) = f_4(0, \mu_2^*) = 0. \qquad (4.12)$$

Furthermore

$$\mu_2^* \geq \mu_2 - p_1\mu_1 \geq \mu_2 - \mu_1/p_2. \qquad (4.13)$$

It follows that for $\gamma_1 \leq \mu_1$ and $f_4(\gamma_1, \gamma_2) < 0$ we always satisfy one of the Conditions C1–C3, see also Figure 4.3. For $\gamma_1 > \mu_1$ and $f_3(\gamma_1, \gamma_2) < 0$ we satisfy Condition C2.

To show neccessity of (4.10) let $\gamma_1 \leq \mu_1$ and $f_4(\gamma_1, \gamma_2) \geq 0$. Since $f_4(\gamma_1, \gamma_2) \geq 0$, Conditions C1 and C3 are not satisfied. Also $f_2(\gamma_1, \gamma_2) \geq 0$ by (4.11)–(4.13) and, therefore, Condition C2 is not satisfied. For $\gamma_1 > \mu_1$ and $f_3(\gamma_1, \gamma_2) \geq 0$ Conditions C1 and C2 are not satisfied. Also, for $\gamma_1 > \mu_1$, $f_2(\gamma_1, \gamma_2) \geq 0 \implies f_4(\gamma_1, \gamma_2) \geq 0$ by (4.11)–(4.13) and Condition C2 cannot be satisfied.  □

Figure 4.3 shows the arrival rates for which the system is stable for a set value of the service rates $\mu_1$ and $\mu_2$. The numbered sections I,II and III represent where Conditions C1, C2 and C3 are met respectively.

Considering all possible values of $\mu_1$ and $\mu_2$ provides the stability range $SR$, i.e. the set of all rates for which the system can be stable provided the optimal service rate allocation is chosen. The stability range is visualised in Figure 4.4.

Figure 4.3: Set of stable rates for set service rates $\mu_1$ and $\mu_2$



Figure 4.4: Set of stable arrival rates

**Corollary 4.3.** *The stability range of the network is given by $\gamma_1 < \mu_1^*$ and $\gamma_2 < \mu_2^*$ and*

$$\frac{\gamma_1}{\max(\mu_1^*, \mu_c)} + \frac{\gamma_2}{\min(\mu_2^*, \mu_c)}\left(1 - \frac{\max(\mu_c - \mu_2^*, 0)}{\mu_1^*}\right) < 1 \ or \quad (4.14)$$

$$\frac{\gamma_1}{\min(\mu_1^*, \mu_c)}\left(1 - \frac{\max(\mu_c - \mu_1^*, 0)}{\mu_2^*}\right) + \frac{\gamma_2}{\max(\mu_2^*, \mu_c)} < 1 \ or \quad (4.15)$$

$$\frac{\gamma_1}{\mu_c} + \frac{\gamma_2}{\mu_c} < 1. \quad (4.16)$$

Corollary 4.3 follows from Theorem 4.2 in a straightforward fashion when considering that due to symmetry we can substitute $\gamma_1 \leq \mu_1$ by $\gamma_2 > \mu_2$ and $\gamma_1 > \mu_1$ by $\gamma_2 \leq \mu_2$ in (4.4) and (4.5). It obviously follows that $\gamma_i < \mu_i^*$ must hold for the second part of both inequalities to hold. Conditions (4.14) and (4.15) follow when considering the extreme allocations for $\mu_1$ and $\mu_2$. The extreme allocations $(\mu_1, \mu_2)$ are $(0, \mu_c)$ and $(\mu_c, 0)$, but when $\mu_i^* < \mu_c$ it needs to be taken into account that $\mu_i < \mu_i^*$. Assuming that one of the boundary rates is lower than the system capacity, say $\mu_2^* < \mu_c$ and $\mu_1^* \geq \mu_c$ then the extreme allocations are $(\mu_c - \mu_2^*, \mu_2^*)$ and $(\mu_c, 0)$. When both boundary rates are lower than the system capacity the extreme allocations are $(\mu_c - \mu_2^*, \mu_2)$ and $(\mu_1^*, \mu_c - \mu_1^*)$. This situation calls for the addition of condition (4.16), as (4.14) and (4.15) do not consider the possible allocations with $\mu_1 + \mu_2 = \mu_c$ for $\mu_c - \mu_2^* < \mu_1 < \mu_1^*$ which show that the rates $\gamma_1 + \gamma_2 < \mu_c$ are feasible as long as $\gamma_1 < \mu_1^*$ and $\gamma_2 < \mu_2^*$.

## 4.4   Product-form characterization

For the situation that $\mu_i^* = \mu_i$, $i = 1, 2$, network $R$ is a Jackson network, which has a geometric product-form stationary distribution

$$\pi(n_1, n_2) = \prod_{i=1}^{2} r_i^{n_i}(1 - r_i), \tag{4.17}$$

for $r_1 = \frac{\gamma_1}{\mu_1}$ and $r_2 = \frac{\gamma_2}{\mu_2}$. It is known from [FI79] that if $p_1 = p_2 = 0$ network $R$ has a product-form stationary distribution if and only if $\mu_1^* + \mu_2^* = \mu_1 + \mu_2$. A natural question is whether such a condition can be generalized to the case that $p_1 \neq 0$ or $p_2 \neq 0$.

In this section we generalize this question and investigate necessary and sufficient conditions for the generalized coupled queue with forwarding $G$, as shown in Figure 4.2, to have a geometric product-form stationary distribution. Recall that at the boundaries of the state space the transition rates are denoted by $a_{k,l}$ for the horizontal axis, by $b_{k,l}$ for the vertical axis and by $c_{k,l}$ at the origin. Also, in the $G$ network we have $a_{1,1} = b_{1,1} = c_{1,1} = 0$. In [LM14] necessary and sufficient conditions are provided for a more general network in which $a_{1,1} \neq 0$, $b_{1,1} \neq 0$ or $c_{1,1} \neq 0$. In our first results below we derive such conditions from first principles, because it enables us to obtain these conditions in the form that is most suitable for follow-up results in this section.

**Theorem 4.4.** *$G$ has a geometric product-form stationary distribution if and only if unique $r_1, r_2 \in (0, 1)$ exist such that*

$$r_2 b_{1,-1} - a_{1,0} + c_{1,0} = \mu_2 p_2 r_2, \tag{4.18}$$

$$r_1 a_{-1,1} - b_{0,1} + c_{0,1} = \mu_1 p_1 r_1, \tag{4.19}$$

$$r_1 a_{-1,0} + r_2 b_{0,-1} - c_{1,0} - c_{0,1} = 0, \tag{4.20}$$

$$r_2 b_{1,-1} + b_{1,0} = \lambda_1 + \mu_2 p_2 r_2, \tag{4.21}$$

$$r_1 a_{-1,1} + a_{0,1} = \lambda_2 + \mu_1 p_1 r_1, \tag{4.22}$$

$$b_{0,1} - r_2 b_{0,-1} - r_2 b_{1,-1} = -\mu_1 p_1 r_1 - \frac{r_2}{r_2 - 1}(\lambda_1 + \mu_2 p_2 r_2 - \mu_1 r_1), \tag{4.23}$$

$$\lambda_1 + \lambda_2 + \mu_1 + \mu_2 = \mu_1(1 - p_1)r_1 + \mu_2(1 - p_2)r_2$$
$$+ \mu_1 p_1 \frac{r_1}{r_2} + \mu_2 p_2 \frac{r_2}{r_1} + \lambda_1 \frac{1}{r_1} + \lambda_2 \frac{1}{r_2}. \tag{4.24}$$

*Proof.* The balance equations must hold for all states $(n_1, n_2)$:

$$\sum_{i,j \in \{-1,0,1\}} \pi(n_1, n_2)q_{ij}(n_1, n_2) = \sum_{i,j \in \{-1,0,1\}} \pi(n_1 - i, n_2 - j)q_{ij}(n_1, n_2). \tag{4.25}$$

We adopt a matrix notation with all new rates in a column vector denoted by $x$:

$$x = [a_{10}, a_{01}, a_{-11}, a_{-10}, b_{10}, b_{01}, b_{0-1}, b_{1-1}, c_{10}, c_{01}]^T, \tag{4.26}$$

and let the rates in the interior be denoted by $q_{i,j}$, omitting the dependency on

the state $n$ as these rates are equal for all states in the interior. The balance equations, apart from the one for the interior, are given by

$$Ax = y, \tag{4.27}$$

where $A$ is given by

$$
\begin{pmatrix}
0 & 0 & 0 & 0 & 1 & 1-\frac{1}{r_2} & 1-r_2 & 1 & 0 & 0 \\
1-\frac{1}{r_1} & 1 & 1 & 1-r_1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & -\frac{r_1}{r_2} & 0 & 1 & 1 & 1-r_2 & 1 & 0 & -\frac{1}{r_2} \\
1 & 1 & 1 & 1-r_1 & 0 & 0 & 0 & -\frac{r_2}{r_1} & -\frac{1}{r_1} & 0 \\
0 & 0 & 0 & 0 & -\frac{1}{r_1} & 0 & 0 & -\frac{r_2}{r_1} & 0 & 0 \\
0 & -\frac{1}{r_2} & -\frac{r_1}{r_2} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & -\frac{1}{r_2} & -\frac{r_1}{r_2} & 0 & -\frac{1}{r_1} & 0 & 0 & -\frac{r_2}{r_1} & 0 & 0 \\
0 & 0 & 0 & -r_1 & 0 & 0 & -r_2 & 0 & 1 & 1
\end{pmatrix}
$$

and $y = [y_1, ..., y_8]^T$ is given by

$$y_1 = r_1 q_{-10} + \frac{r_1}{r_2} q_{-11},$$

$$y_2 = r_2 q_{0-1} + \frac{r_2}{r_1} q_{1-1},$$

$$y_3 = r_1 q_{-10},$$

$$y_4 = r_2 q_{0-1},$$

$$y_5 = (r_1 - 1)q_{-10} + (r_2 - 1)q_{0-1} + (\frac{1}{r_2} - 1)q_{01} + (\frac{r_1}{r_2} - 1)q_{-11} - q_{10} - q_{1-1},$$

$$y_6 = (r_1 - 1)q_{-10} + (r_2 - 1)q_{0-1} + (\frac{r_2}{r_1} - 1)q_{1-1} + (\frac{1}{r_1} - 1)q_{10} - q_{01} - q_{-11},$$

$$y_7 = (r_1 - 1)q_{-10} + (r_2 - 1)q_{0-1} - q_{1-1} - q_{10} - q_{01} - q_{-11},$$

$$y_8 = 0.$$

Using Gauss-Jordan elimination on the matrix $[A|y]$ for $G$ we obtain (omitting empty rows and combining a few rows for convenient notation)

$$
\left(
\begin{array}{cccccccccc|c}
1 & 0 & 0 & 0 & 0 & 0 & 0 & -r_2 & -1 & 0 & -\mu_2 p_2 r_2 \\
0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & -1 & \lambda_2 \\
0 & 0 & 1 & 0 & 0 & -\frac{1}{r_1} & 0 & 0 & 0 & \frac{1}{r_1} & \mu_1 p_1 \\
0 & 0 & 0 & 1 & 0 & 0 & \frac{r_2}{r_1} & 0 & -\frac{1}{r_1} & -\frac{1}{r_1} & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & r_2 & 0 & 0 & \lambda_1 + \mu_2 p_2 r_2 \\
0 & 0 & 0 & 0 & 0 & -1 & r_2 & r_2 & 0 & 0 & \mu_1 p_1 r_1 + \frac{r_2}{r_2-1}(\lambda_1 + \mu_2 p_2 r_2 - \mu_1 r_1)
\end{array}
\right)
$$

Equations (4.18)-(4.23) now follow from linear combinations of these equations.$\square$

From Theorem 4.4 we can straightforwardly derive a known result from [FI79] for our coupled queue network $R$. We prove only that the conditions are neccessary. It is shown in [FI79] that these conditions are also sufficient.

**Theorem 4.5.** *For the coupled queue network $R$ without forwarding of packets, i.e. $p_1 = p_2 = 0$, a necessary condition to have a product-form stationary distribution is*

$$\mu_1^* + \mu_2^* = \mu_1 + \mu_2. \tag{4.28}$$

*Proof.* Recall that for $R$ we have by definition that $a_{1,0} = b_{1,0} = c_{1,0} = \lambda_1$ and $a_{0,1} = b_{0,1} = c_{0,1} = \lambda_2$. When no forwarding of packets occurs, i.e. $p_1 = p_2 = 0$, we have that $a_{-1,1} = b_{1,-1} = 0$, $a_{-1,0} = \mu_1^*$ and $b_{0,-1} = \mu_2^*$. Equations (4.18),(4.19),(4.21) and (4.22) automatically hold. For the network to be of geometric product-form, the service rates at the boundary have to be chosen such that

$$r_1\mu_1^* + r_2\mu_2^* = \lambda_1 + \lambda_2, \tag{4.29}$$

as is given by Equation (4.20). We combine this equation with the balance equations at both boundaries

$$\lambda_1 + \lambda_2 + \mu_1^* = \frac{\lambda_1}{r_1} + \mu_1^* r_1 + \mu_2 r_2, \tag{4.30}$$

$$\lambda_1 + \lambda_2 + \mu_2^* = \frac{\lambda_2}{r_2} + \mu_2^* r_2 + \mu_1 r_1, \tag{4.31}$$

which gives

$$\mu_1^* + \mu_2^* = \frac{\lambda_1}{r_1} + \frac{\lambda_2}{r_2} + \mu_1 r_1 + \mu_2 r_2 - \lambda_1 - \lambda_2. \tag{4.32}$$

Finally, considering the balance equation (4.24) of the interior, the right hand side of (4.32) can be simplified to $\mu_1 + \mu_2$ completing the proof. $\qquad\square$

Our next result deals with the case that there is forwarding between the nodes, i.e. we have $p_1 \neq 0$ or $p_2 \neq 0$. Our results states that a $R$ network has a product-form stationary distribution if and only if $\mu_1^* = \mu_1$ and $\mu_2^* = \mu_2$. This means that the Jackson network is the only $R$ network that has a product-form stationary distribution. The result shows that generalizations to, for instance $\mu_1^* + \mu_2^* = \mu_1 + \mu_2$, are not possible as soon as $p_1 \neq 0$ or $p_2 \neq 0$.

**Theorem 4.6.** *The coupled queue $R$ with forwarding of packets, i.e. $p_1 \neq 0$ or $p_2 \neq 0$, has a product-form stationary distribution if and only if*

$$\mu_1^* = \mu_1 \text{ and } \mu_2^* = \mu_2. \tag{4.33}$$

*Proof.* Without loss of generality, assume that $p_1 \neq 0$. Equation (4.22) shows that with arrival rate $a_{0,1} = \lambda_2$ the rate $a_{-1,1} = \mu_1^* p_1$ must equal $\mu_1 p_1$. This shows that the rate at the boundary must equal the rate at the interior, i.e. $\mu_1^* = \mu_1$. We now use a similar approach as for the case without forwarding of packets, starting from Equation (4.20):

$$r_1\mu_1 + r_2\mu_2^* = \lambda_1 + \lambda_2. \tag{4.34}$$

Combining (4.34) with the balance equation at the vertical boundary

$$\lambda_1 + \lambda_2 + \mu_2^* = \frac{\lambda_2}{r_2} + \mu_2^* r_2 + \mu_1 r_1 (1 - p_1) + \frac{r_1}{r_2} \mu_1 p_1 \tag{4.35}$$

gives

$$\mu_2^* = \frac{\lambda_2}{r_2} - \mu_1 p_1 r_1 + \frac{r_1}{r_2} \mu_1 p_1. \tag{4.36}$$

Using balance equation (4.24) of the interior this equals

$$\mu_2^* = \lambda_1 + \lambda_2 + \mu_1 + \mu_2 - \frac{\lambda_1}{r_1} - \mu_1 r_1 - \mu_2 r_2. \tag{4.37}$$

Finally the balance equation at the horizontal boundary

$$\lambda_1 + \lambda_2 + \mu_1 = \frac{\lambda_1}{r_1} + \mu_1 r_1 + \mu_2 r_2 \tag{4.38}$$

shows that the right hand side of (4.37) equals $\mu_2$, completing the proof. $\quad\square$

In the next section we will analyze $R$ with $\mu_i^* > \mu_i$, $p_1 > 0$ and $p_2 = 0$, which by Theorem 4.6 does not have a product-form stationary distribution. Our approach will be to bound the performance in terms of a perturbed network that is obtained by changing some of the rates at the boundary of the state space. These changes will be made such that the resulting perturbed network has a product-form stationary distribution with known parameters. The details of the performance bounding method will be given in the next section. Here we present the two perturbed networks, denoted by $G_1$ and $G_2$ that will be used. The transition rates of the perturbed networks will be denoted with a bar.

$G_1$ The first perturbed network is an $R$ network in which we take $\bar{\mu}_1^* = \mu_1$ and $\bar{\mu}_2^* = \mu_2$. Recall from above that this is an instance of a Jackson network with $r_1 = \gamma_1/\mu_1$ and $r_2 = \gamma_2/\mu_2$.

$G_2$ The second perturbed network generalizes to a $G$ network. We don't perturb the rates of $R$ at the origin and the vertical boundary, i.e. $\bar{b}_{1,0} = \bar{c}_{1,0} = \lambda_1$, $\bar{b}_{0,1} = \bar{c}_{0,1} = \lambda_2$, $\bar{b}_{1,-1} = 0$ (recall that $p_2 = 0$) and $\bar{b}_{0,-1} = \mu_2^*$. On the horizontal axis (verifying that we satisfy Equations (4.18)–(4.22)) we take

$$\bar{a}_{1,0} = \lambda_1, \quad \bar{a}_{0,1} = \lambda_2, \quad \bar{a}_{-1,1} = \mu_1 p_1 \quad \text{and} \quad \bar{a}_{-1,0} = \frac{\lambda_1 + \lambda_2 - \mu_2 r_2}{r_1}, \tag{4.39}$$

where $r_1$ and $r_2$ are the solution in $(0,1)$ of

$$r_1 = \frac{\lambda_1}{\mu_1 + (r_2 - 1)(\mu_2^* - \mu_2)}, \quad \text{and} \quad r_1 = \frac{\mu_2^* r_2 - \lambda_2}{(\mu_2^* - \mu_2) r_2 + \mu_1 p_1}. \tag{4.40}$$

The network $G_1$ always exists, i.e. we can always construct rates $\bar{\mu}_1^* = \mu_1$ and $\bar{\mu}_2^* = \mu_2$. and obtain a positive recurrent $R$ network with $r_1 = \gamma_1/\mu_1$ and $r_2 = \gamma_2/\mu_2$. We have no guarantee that the $G_2$ network can always be

constructed. We will show in the next section that there are examples of $R$ networks for which the $G_2$ network does not exist. Sufficient conditions for existence of $G_2$ in an explicit form do not seem to follow straightforwardly from (4.39) and (4.40).

The network $G_2$ is obtained by only perturbing the horizontal axis. One could also consider perturbing only the vertical axis. However, it follows from (4.22) that this is only possible when $\mu_1^* = \mu_1$ and this will, therefore, not be considered in the remainder.

## 4.5  Markov reward approach and bounds

In the previous section we provided necessary and sufficient conditions under which an $R$ network has a product-form stationary distribution. Our approach to analyzing networks that do not have a product-form stationary distribution is to establish upper and lower bounds on the various performance measures like the probability that the system is empty or the expected number of packets in a node. The bounds will be established by comparing the network of interest with a perturbed network with a product-form stationary distribution. We follow the Markov reward approach to establishing error bounds as developed by van Dijk. An overview of this method is presented in [vD11]. More precisely, we use the method presented in [GBvO13] to establish error bounds as the solution of a linear program.

The Markov reward approach compares two queueing networks modeled as Markov chains, one for which the stationary distribution is not known, and another one, which is a modification of the first one, for which the stationary distribution is known. The key element of the approach is to analyze steady state performance measures by means of a cumulative reward structure. In particular, we express our performance measure of interest as $\mathcal{F} = \sum_n \pi(n)F(n)$, where $\pi(n)$ is the (unknown) stationary distribution and $F(n)$ is a reward function. We adapt $F(n)$ to the desired performance measure: If $F(n) = n_1$, then $\mathcal{F}$ corresponds to the expected number of packets in the first node; if $F(n) = \mathbf{1}\{n = 0\}$, then $\mathcal{F}$ corresponds to the probability that the system is empty. We present the Markov reward approach in terms of a general function $F(n)$ of the form

$$F(n) = \begin{cases} f_{1,0} + f_{1,1}n_1, & \text{if } n \in C_1, \\ f_{2,0} + f_{2,2}n_2, & \text{if } n \in C_2, \\ f_{3,0}, & \text{if } n \in C_3, \\ f_{4,0} + f_{4,1}n_1 + f_{4,2}n_2, & \text{if } n \in C_4, \end{cases} \tag{4.41}$$

where $f_{m,n}$ are the constants that define the function and $C_1, \ldots, C_4$ denote the horizontal axis, the vertical axis, the origin and the interior of the state space, respectively. Our functions $F(n)$ are linear in each of the components of the state space. For notational convenience, let $N_k$ denote the set of possible transitions from a state in $C_k, k = 1, .., 4$. Moreover, let $k(n)$ denote the index of the component of the state space that state $n$ is in. To illustrate: $k((3,0)) = $ and $k((4,6)) = 4$. Finally, let $q_{i,j}(n)$ and $\bar{q}_{i,j}(n)$ denote the transition rate from

$n$ to $n + (i, j)$ in the original and in the perturbed network, respectively.

Since we are interested in steady state behavior only we can uniformize our network to obtain a discrete-time Markov chain for which the stationary distribution is the same as for our continuous-time network. The Markov reward approach is most conveniently presented in terms of this discrete-time equivalent. More precisely, we assume that both the original and the perturbed network can be uniformized with the same uniformization constant $w$, i.e. we let

$$w \geq \max \left\{ \max_{n} \sum_{i,j \in \{-1,0,1\}} q_{i,j}(n), \quad \max_{n} \sum_{i,j \in \{-1,0,1\}} \bar{q}_{i,j}(n) \right\}. \tag{4.42}$$

Let $p_{k(n),u}$ and $\bar{p}_{k(n),u}$ denote the probability of making a transition from a state $n$ in $C_k(n)$ to $n + u$ in the original and in the perturbed network, respectively. The uniformization of both networks gives

$$p_{k(n),u} = \frac{q_{i,j}(n)}{w}, \quad \text{and} \quad \bar{p}_{k(n),u} = \frac{\bar{q}_{i,j}(n)}{w}, \tag{4.43}$$

for $u \neq 0$ and

$$p_{k(n),0} = 1 - \sum_{u \in N_k} p_{k(n),u}, \quad \text{and} \quad \bar{p}_{k(n),0} = 1 - \sum_{u \in N_k} \bar{p}_{k(n),u}. \tag{4.44}$$

Define the difference in transition probabilities to be given by

$$\delta_{k(n),u} = \bar{p}_{k(n),u} - p_{k(n),u}. \tag{4.45}$$

The starting point for the Markov reward approach is to consider $F^t(n)$, the expected cumulative reward at time $t$ when the network starts in state $n$ at time 0, for which

$$F^t(n) = \begin{cases} F(n) + \sum_{u \in N_{k(n)}} p_{k(n),u} F^{t-1}(n + u), & \text{if } t > 0, \\ 0, & \text{if } t = 0. \end{cases} \tag{4.46}$$

Note that $\mathcal{F} = \lim_{t \to \infty} \sum_n F^t(n)/\pi(n)$. The key idea is that bounds on terms of the form $F^t(n+u) - F^t(n)$ can be used to bound $\mathcal{F}$ in terms of $\bar{\pi}(n)$. Therefore, we introduce

$$D_u^t(n) = F^t(n + u) - F^t(n) \tag{4.47}$$

and refer to $D_u^t(n)$ as bias terms.

**Theorem 4.7 ([vD11], Result 9.3.5).** *Let* $\bar{F} : S \to [0, \infty)$ *and* $G : S \to [0, \infty)$ *satisfy*

$$\left| \bar{F}(n) - F(n) + \sum_{u \in N_{k(n)}} \delta_{k(n),u} D_u^t(n) \right| \leq G(n), \tag{4.48}$$

| $k$ | $j$ | $u$ | $c_{1,k,j,u}$ |
|---|---|---|---|
| 1 | 1 | $N_1$ | $p_{1,u}$ |
| 2 | 1 | $\{d_1, e_1, d_2\}$ | $p_{4,u}$ |
| 2 | 1 | $e_2$ | $p_{4,e_2} - p_{2,d_1} + c_{1,2,1,d_1}$ |
| 2 | 1 | 0 | $p_{4,0} - p_{2,e_1} + c_{1,2,1,e_1}$ |
| 2 | 1 | $-e_2$ | $p_{4,-e_2} - p_{2,d_2} + c_{1,2,1,d_2}$ |
| 2 | 2 | 0 | $p_{4,-d_2} - p_{2,e_2} + c_{1,2,1,e_2}$ |
| 2 | 2 | $-e_2$ | $p_{4,-e_1} - p_{2,0} + c_{1,2,2,0} + c_{1,2,1,0}$ |
| 3 | 1 | $\{e_1, d_1\}$ | $p_{1,u}$ |
| 3 | 1 | $e_2$ | $p_{1,e_2} - p_{3,d_1} + c_{1,3,1,d_1}$ |
| 3 | 1 | 0 | $p_{1,0} - p_{3,e_1} + c_{1,3,1,e_1}$ |
| 3 | 2 | 0 | $p_{1,-d_2} - p_{3,e_2} + c_{1,3,1,e_2}$ |
| 4 | 1 | $N_4$ | $p_{4,u}$ |

| $k$ | $j$ | $u$ | $c_{2,k,j,u}$ |
|---|---|---|---|
| 1 | 2 | $\{d_1, e_2, -d_2\}$ | $p_{4,u}$ |
| 1 | 2 | $e_1$ | $p_{4,e_1} - p_{1,d_1} + c_{2,1,2,d_1}$ |
| 1 | 2 | 0 | $p_{4,0} - p_{1,e_2} + c_{2,1,2,e_2}$ |
| 1 | 2 | $-e_1$ | $p_{4,-e_1} - p_{1,-d_2} + c_{2,1,2,-d_2}$ |
| 1 | 1 | 0 | $p_{4,d_2} - p_{1,e_1} + c_{2,1,2,e_1}$ |
| 1 | 1 | $-e_1$ | $p_{4,-e_2} - p_{1,0} + c_{2,1,1,0} + c_{2,1,2,0}$ |
| 2 | 2 | $N_2$ | $p_{2,u}$ |
| 3 | 2 | $\{d_1, e_2\}$ | $p_{2,u}$ |
| 3 | 2 | $e_1$ | $p_{2,e_1} - p_{3,d_1} + c_{2,3,2,d_1}$ |
| 3 | 2 | 0 | $p_{2,0} - p_{3,e_2} + c_{2,3,2,e_2}$ |
| 3 | 1 | 0 | $p_{2,d_2} - p_{3,e_1} + c_{2,3,2,e_1}$ |
| 4 | 2 | $N_4$ | $p_{4,u}$ |

Table 4.1: Values for constants $c_{i,k,j,u}$.

*for all $n \in S$ and $t \geq 0$ then*

$$\sum_{n \in S} [\bar{F}(n) - G(n)]\bar{\pi}(n) \leq \mathcal{F} \leq \sum_{n \in S} [\bar{F}(n) + G(n)]\bar{\pi}(n). \qquad (4.49)$$

The difficulty in applying Theorem 4.7 is that it is meticulous to find an appropriate function $G(n)$. In [GBvO13] a linear programming approach is presented that provides a general optimization problem that establishes an error bound between any two random walks and can bound the difference which has $\delta_{k,u}$ and $\bar{\pi}$ as input parameters and a lower (or upper) bound as output. This linear programming approach will be used in this section to establish performance bounds for our $R$ network. We will not provide all details of the approach, but

we will present its main ideas. The key result in [GBvO13] is the following.

**Lemma 4.8 ([GBvO13], Lemma 3).** *Let constants $c_{i,k,j,u}$, $i,j = 1,2$, $k = 1,\ldots,4$, $u \in \{-1,0,1\}^2$ take the values given in Table 4.1. If $A_i : S \to [0,\infty)$ and $B_i : S \to [0,\infty)$, $i = 1,2$ satisfy*

$$F(n + e_i) - F(n) + \sum_{j=1,2} \sum_{u \in N_k} \max\{-c_{i,k,j,u} A_j(n + u), c_{i,k,j,u} B_j(n + u)\} \le B_i(n),$$

$$F(n) - F(n + e_i) + \sum_{j=1,2} \sum_{u \in N_k} \max\{-c_{i,k,j,u} B_j(n + u), c_{i,k,j,u} A_j(n + u)\} \le A_i(n),$$

*where $k = k(n)$, for all $n \in S$ and $i = 1,2$, then*

$$-A_i(n) \le D_{e_i}^t(n) \le B_i(n), \tag{4.50}$$

*for all $t \ge 0$, $n \in S$ and $i = 1,2$.*

The intuition behind this result is that it can be shown that for the values of $c_{i,k,j,u}$ given in Table 4.1 we have

$$D_i^{t+1}(n) = F(n + e_i) - F(n) + \sum_{j=1,2} \sum_{u \in N_{k(n)}} c_{i,k(n),j,u} D_j^t(n + u). \tag{4.51}$$

The result of Lemma 4.8 then follows by induction. Lemma 4.8 establishes bounds on the bias terms $D_{e_1}^t(n)$ and $D_{e_2}^t(n)$, i.e. on the bias terms in the unit directions. It is shown in [GBvO13] that these bounds can be extended to provide bounds in the other (diagonal) directions and that linear constraints can be formulated that capture the requirements on $G(n)$ itself. Finally, it is shown in [GBvO13] that it is possible to reduce the resulting linear program to a program with a finite number of variables and constraints.

We first compare the $R$ network with the product-form network $G_1$ where $\mu_i^* = \mu_i$. Using the results of Section 4.4 we also compare with network $G_2$, where we set the rates of the vertical axis and the origin equal to $R$, having calculated the remaining rates at the vertical boundary by solving the equations of Theorem 4.4. The rates are presented in Table 4.2 (where the rates equal to $R$ are omitted from the table for readability).

As a first performance measure we consider $\pi(0,0)$, the probability that the system is empty, i.e. $f_{3,0} = 1$ in (4.41) while all other values are 0. The bounds and simulated values are presented in Figure 4.5. The comparison of $R$ and $G_1$ provides an upper bound that is close to the simulated values, yet the lower bound does not increase for higher values of $\mu_i^*$. The comparison of $R$ and $G_2$ does not provide an improvement of the bounds.

As a second performance measure, we provide bounds on $EN_1$ and $EN_2$, the expected number of customers at queue 1 and 2 respecitvely, i.e. $f_{i,1} = 1$ or $f_{i,2} = 1$ for all $i$, as presented in Figure 4.6. For $EN_1$, bounds obtained by comparing $R$ and $G_2$ again do not give an improvement. However, for $EN_2$ we see an improvement on the upper bound for the lower values of $\mu_i^*$. For future research this shows that it would be interesting to find more product-form

| $\mu_2^*$ | $a_{1,0}$ | $a_{0,1}$ | $a_{-1,0}$ | $a_{-1,1}$ | $r_1$ | $r_2$ |
|------|------|------|------|------|------|------|
| 0.4 | 0.1 | 0.1 | 0.28 | 0.12 | 0.25 | 0.325 |
| 0.45 | 0.1 | 0.1 | 0.23 | 0.12 | 0.2738 | 0.3045 |
| 0.5 | 0.1 | 0.1 | 0.18 | 0.12 | 0.3039 | 0.2906 |
| 0.55 | 0.1 | 0.1 | 0.13 | 0.12 | 0.3420 | 0.2828 |
| 0.6 | 0.1 | 0.1 | 0.08 | 0.12 | 0.3902 | 0.2813 |
| 0.65 | 0.1 | 0.1 | 0.03 | 0.12 | 0.4510 | 0.2869 |

Table 4.2: Rates for which network $G_2$ has a product-form stationary distribution for various values of $\mu_2^*$ with $\lambda_1 = \lambda_2 = 0.1$, $\mu_1 = \mu_2 = 0.4$ and $p_1 = 0.3$



Figure 4.5: Probability $\pi(0,0)$ for $R$ for various values of $\mu_i^*$ with $\lambda_i = 0.1, \mu_i = 0.4$ and $p_1 = 0.3$

networks to compare with, possibly improving the bounds.

The Markov reward approach starts from the performance measure of the perturbed network and provides the bounds by adding (subtracting) the error bounds of equation (4.48) to obtain the upper (lower) bound, as can be seen in (4.49). This causes that for an increasing performance measure as the probability of the system being empty that the lower bound does not increase. For the decreasing value of the expected number of customers in the queue the upper bound doesn't decrease, but even increases. This explains why for the probability of the state being empty the upper bound is close to the simulation, whereas the lower bounds are closer for the expected of number of customers in a queue.

## 4.6   Optimal rate allocation

The previous section provides bounds on the perfomance measures, but we are also interested in the impact of allocating capacity to each server in $R$. Recall

Figure 4.6: Expected number of customers in each queue
for various values of $\mu_i^*$ with $\lambda_i = 0.1, \mu_i = 0.4$ and $p_1 = 0.3$

that the system capacity is given by $\mu_c = \mu_1 + \mu_2$ and the service rates at the boundary are higher than in the interior, i.e. $\mu_i^* \geq \mu_i$. Van Dijk [vD11] provides a means to compare networks which differ only in transition probabilities.

**Theorem 4.9.** *(cf. [vD11]) Let $\hat{F} : S \to [0, \infty)$ and $F : S \to [0, \infty)$ satisfy*

$$\hat{F}(n) - F(n) + \sum_{u \in N_{k(n)}} \delta_{k(n),u} D_u^t(n) \leq 0, \tag{4.52}$$

*for all $n \in S$ and $t \geq 0$ then*

$$\sum_n \hat{F}(n)\hat{\pi}(n) \leq \sum_n F(n)\pi(n). \tag{4.53}$$

In the following we compare two $R$ networks with different service rates in the interior and show that in this case only one bias term is relevant in Theorem 4.9. We provide recursive expressions for this bias term and use proof by induction to show the sign of the bias term is negative under certain conditions, proving it is optimal to allocate all system capacity to one of the nodes.

**Theorem 4.10.** *Assume that $\mu_1^* \leq \mu_c$ and $\mu_2^* \geq \mu_c$, then the steady state probability $\pi(0,0)$ for the system $R$ without forwarding, i.e. $p_1 = p_2 = 0$, is minimized by allocating*

$$\mu_1 = \mu_c \text{ and } \mu_2 = 0. \tag{4.54}$$

*Proof.* We compare an $R$ network with service rates $(\mu_1, \mu_2)$ in the interior with and $R$ network wiith rates $(\mu_1 - \epsilon, \mu_2 + \epsilon) = (\hat{\mu_1}, \hat{\mu_2})$, where both nodes do not forward their packets to the other one, i.e. $p_1 = p_2 = 0$. We uniformize both

networks using the uniformization constant as given in (4.42), which in this case is equal to $w = \lambda_1 + \lambda_2 + \mu_1^* + \mu_2^*$. As is shown in [vD11], this constant can be omitted from the equations and for ease of notation we will denote the transition probabilities by the transition rates. The difference $\delta_{k,u}$ between the transition probabilities of the systems that are being compared is given by

$$\delta_{k,u} = \begin{cases} -\epsilon, & \text{if } n \in C_4, u = -e_1, \\ \epsilon, & \text{if } n \in C_4, u = -e_2, \\ 0 & \text{otherwise.} \end{cases} \tag{4.55}$$

and the same reward functions are used for both networks, so that $\hat{F}(n) = F(n)$. Recalling that $d_2 = (1, -1)$, the left hand side of Equation (4.52) of Theorem 4.9 is given by

$$\begin{aligned} -\epsilon D^t_{-e_1} + \epsilon D^t_{-e_2} &= \\ -\epsilon(F^t(n - e_1) - F^t(n)) + \epsilon(F^t(n - e_2) - F^t(n)) &= \\ \epsilon(F^t(n - e_2) - F^t(n - e_1)) &= \\ \epsilon(F^t(n - e_1 + d_2) - F^t(n - e_1)) &= \\ \epsilon D^t_{d_2}(n - e_1). \end{aligned} \tag{4.56}$$

To show that $D^t_{d_2}(n - e_1) < 0$ so that Theorem 4.9 can be applied, we use the recurrence relations for $D^t_u(n)$. As an example we give this recurrence relation for one value of $n$ with $u = d_2$. First note the definition of the bias term

$$D^t_{d_2}(n) = F^t(n + d_2) - F^t(n).$$

Starting from state $(0, 1)$ and considering that the step $d_2$ would lead to state $(1, 0)$, the definition and the possible transitions as given in (4.46) from each of these states gives

$$\begin{aligned} D^{t+1}_{d_2}(0, 1) &= F^{t+1}(1, 0) - F^{t+1}(0, 1) \\ &= F(1, 0) - F(0, 1) + \lambda_1 F^t(2, 0) + \lambda_2 F^t(1, 1) \\ &\quad + \mu_1^* F^t(0, 0) + \mu_2^* F(1, 0) - \lambda_1 F^t(1, 1) - \lambda_2 F^t(0, 2) \\ &\quad - \mu_1^* F^t(0, 1) - \mu_2^* F^t(0, 0) \\ &= F(1, 0) - F(0, 1) + \lambda_1(F^t(2, 0) - F^t(1, 1)) \\ &\quad + \lambda_2(F^t(1, 1) - F^t(0, 2)) + \mu_1^*(F^t(0, 0) - F^t(0, 1)) \quad (4.57) \\ &\quad + \mu_2^*(F^t(1, 0) - F^t(0, 0)) \\ &= F(1, 0) - F(0, 1) + \lambda_1 D_{d_2}(1, 1) + \lambda_2 D_{d_2}(0, 2) \\ &\quad + \mu_1^* D_{-e_2}(0, 1) + \mu_2^* D_{d_2}(0, 1) - \mu_2^* D_{-e_2}(0, 1) \\ &= F(1, 0) - F(0, 1) + \lambda_1 D_{d_2}(1, 1) + \lambda_2 D_{d_2}(0, 2) \\ &\quad + (\mu_1^* - \mu_2^*) D_{-e_2}(0, 1) + \mu_2^* D_{d_2}(0, 1). \end{aligned}$$

Similarly, all other expressions can be derived.

We now prove that the assumptions $\mu_1^* \leq \mu_c$ and $\mu_2^* \geq \mu_c$ are sufficient for

the bias term $D_{d_2}^t(n - e_1)$ to be non-positive. Let $\nu = \mu_1^* + \mu_2^* - \mu_1 - \mu_2$ then it can be verified that the recursion of the bias terms can be expressed as

$$
D_{d_2}^{t+1}(k) = \begin{cases}
F(1,0) - F(0,1) + \lambda_1 D_{d_2}^t(1,1) + \lambda_2 D_{d_2}^t(0,2) \\
+ \mu_2^* D_{d_2}^t(0,1) + (\mu_1^* - \mu_2^*) D_{-e_2}^t(0,1), & \text{if } k = (0,1) \\
F(n+1,0) - F(n,1) + \lambda_1 D_{d_2}^t(n+1,1) + \lambda_2 D_{d_2}^t(n,2) \\
+ \mu_1 D_{d_2}^t(n-1,1) + \mu_2^* D_{d_2}^t(n,1) \\
+ (\mu_1^* - \mu_1 - \mu_2) D_{-e_2}^t(n,1), & \text{if } k = (n,1) \\
F(1,m-1) - F(0,m) + \lambda_1 D_{d_2}^t(1,m) + \lambda_2 D_{d_2}^t(0,m+1) \\
+ \mu_2 D_{d_2}^t(0,m-1) + \nu D_{d_2}^t(0,m) \\
+ (\mu_1 + \mu_2 - \mu_2^*) D_{-e_2}^t(0,m), & \text{if } k = (0,m) \\
F(n+1,m-1) - F(n,m) + \lambda_1 D_{d_2}^t(n+1,m) \\
+ \lambda_2 D_{d_2}^t(n,m+1) + \mu_1 D_{d_2}^t(n-1,m) \\
+ \mu_2 D_{d_2}^t(n,m-1) + \nu D_{d_2}^t(n,m), & \text{if } k = (n,m)
\end{cases}
\tag{4.58}
$$

and

$$
D_{-e_2}^{t+1}(k) = \begin{cases}
F(0,0) - F(0,1) + \lambda_1 D_{-e_2}^t(1,1) \\
+ \lambda_2 D_{-e_2}^t(0,2) + \mu_1^* D_{-e_2}^t(0,1) & \text{if } k = (0,1) \\
F(n,0) - F(n,1) + \lambda_1 D_{-e_2}^t(n+1,1) \\
+ \lambda_2 D_{-e_2}^t(n,2) + \mu_1^* D_{-e_2}^t(n-1,1) \\
+ (\mu_1 - \mu_1^*) D_{d_2}^t(n-1,1) + \nu D_{-e_2}^t(n,1) & \text{if } k = (n,1) \\
F(0,m-1) - F(0,m) + \lambda_1 D_{-e_2}^t(1,m) \\
+ \lambda_2 D_{-e_2}^t(0,m+1) + \mu_1^* D_{-e_2}^t(0,m) + \mu_2^* D_{-e_2}^t(0,m-1) & \text{if } k = (0,m) \\
F(n,m-1) - F(n,m) + \lambda_1 D_{-e_2}^t(n+1,m) \\
+ \lambda_2 D_{-e_2}^t(n,m+1) + \mu_1 D_{-e_2}^t(n-1,m) \\
+ \mu_2 D_{-e_2}^t(n,m-1) + \nu D_{-e_2}^t(n,m), & \text{if } k = (n,m)
\end{cases}
\tag{4.59}
$$

The remainder of the proof is by use of induction. First note that $F(0,0) = 1$ and $F(n,m) = 0$ for all other values when considering the probability that the system is empty and that $\mu_1^* - \mu_2^* \le 0$, which follows immediately from the assumptions. As a base for the recursion we have that $D_{-e_2}^0(k) = 0 \ge 0$ and $D_{d_2}^0(k) = 0 \le 0$. For the inductive step, note that as $\lambda_i, \mu_i, \mu_i^*$ and $\nu$ are positive, it follows from (4.59) that as $\mu_1 - \mu_1^*$ is negative, $D_{-e_2}^{t+1}(k)$ is positive when $D_{-e_2}^t(k)$ is positive, as long as $D_{d_2}^t(k)$ is negative. From (4.58) it follows that $D_{d_2}^{t+1}(k)$ is negative when $D_{d_2}^t(k)$ is negative, $D_{-e_2}^t(k)$ is positive and the assumptions hold, as this ensures that the expressions $\mu_1^* - \mu_2^*$, $\mu_1^* - \mu_1 - \mu_2$ and $\mu_1 + \mu_2 - \mu_2^*$ become negative. Together with the base this completes the proof.

Concluding, we have that $D_{d_2}^t(n - e_1)$ is non positive and $\hat{F}(n) = F(n)$, so that

$$
\hat{F}(n) - F(n) + \epsilon D_{d_2}^t(n - e_1) \le 0. \tag{4.60}
$$

It follows from Theorem 4.9 that the network performs better without adjusting

Figure 4.7: Performance measures for $\mu_c = 1, \mu_1^* = \mu_2^* = 3, \lambda_1 = 0.5$ and $\lambda_2 = 0.4$

the service rates. Starting with the premise that all capacity is allocated to node 1, this analysis holds, proving that it is suboptimal to move capacity to node 2. Hence an optimal situation is acquired, completing the proof.                □

The line of proof followed here does not apply for other values of $\mu_i^*$ compared to $\mu_c$ as the signs of the expressions $\mu_1^* - \mu_2^*$, $\mu_1^* - \mu_1 - \mu_2$ and $\mu_1 + \mu_2 - \mu_2^*$ become positive. This prohibits the use of induction on the signs of $D_{d_2}^{t+1}(k)$ and $D_{-e_2}^{t+1}(k)$. For the values of $\mu_i^*$ not included in Theorem 4.10, we postulate that it is still optimal to allocate all the capacity to one node, even when considering other performance measures. The node to which all capacity should be allocated depends on the arrival rate at each of the nodes. In the following we provide simulation results to support our postulation.

As Figure 4.7 shows, allocating all capacity to one of the nodes provides better results than allocating the capacity evenly as on average less packets remain at both nodes. Due to the lower arrival rate at node 2, it is better to allocate all capacity to this node, as this increases the probability of reaching a state at the boundary, so that the much higher service rate can be used.

When the arrival rates for each node are more unbalanced, as shown in Figure 4.8, the best choice is clearly to allocate all capacity to the node with the lowest arrival rate, in this case node 2. This now even holds for the performance measures of node 1 itself. If 90% of the capacity would be allocated to node 1, the system would even become unstable. Note that even though $\lambda_1 + \lambda_2 > \mu_c$, the system is stable when the capacity is allocated in a correct manner, in accordance to our findings in Section 4.3.

## 4.7   Conclusion

This paper analysed a two node network where the two nodes interfere, causing a lower service rate when both nodes are active. Starting with the stability of the system, it was shown that increasing the rate at the edge of the state space expands the stability region. Conditions were given for which the system

Figure 4.8: Performance measures for $\mu_c = 1, \mu_1^* = \mu_2^* = 3$, $\lambda_1 = 0.95$ and $\lambda_2 = 0.1$

has a product-form stationary distribution, providing networks for comparison with the network under consideration using a Markov reward approach. For any performance measure that is linear in each of the components of the state space an approach is given to provide bounds. Examples have been provided for the average number of customers in each queue and the probability of the system being empty, showing that some of the provided bounds are close to simulated results. This gives a promising start to further investigate the strength of the Markov reward approach to obtain results for a network, for which it is well known that analytical results are hard to obtain. We also analysed the impact of the allocation of the capacity to each of the nodes when they both are busy, showing that allocating all the capacity to one of the nodes provides better results than evenly distributing the available capacity. The node with the lowest arrival should receive all capacity, as this way the probability of reaching a state at the boundary, with a higher service rate, is highest. As our examples point out, this may even be beneficial to the node that does not receive any of the capacity.

# Upper bounds on multi-hop multi-channel wireless network performance

Given a placement of wireless nodes in space and a traffic demand between pairs of nodes, can these traffic demands be supported by the resulting network? A key issue in answering this question is interference between nodes. This chapter presents a generic model for sustainable network load in a multi-hop wireless network under interference constraints, and recasts this model into a multicommodity flow problem with interference constraints. Using Farkas' Lemma, we obtain a necessary and sufficient condition for feasibility of this multicommodity flow problem, leading to a tight upper bound on network throughput. The results are illustrated by examples such as a serial network and a network taken from literature.

## 5.1   Introduction

Interference is an important aspect of wireless networks that seriously affects the capacity of the network. This is especially so in a wireless multi-hop network, where a transmission on one link interferes with transmissions on links in the vicinity. On a multi-hop path self-interference may result in substantial degradation of end-to-end network performance. In this respect, due to interaction among hops, multi-hop wireless networks differ considerably from wired networks thus calling for new modelling and analysis techniques that take into account interference constraints.

In the absence of interference, but including capacity constraints on the transmission rate of nodes, feasibility of a set of traffic demands between pairs of nodes can be determined by considering the flow allocation in the network as a multicommodity flow problem. The network is modelled as a graph, with the vertices representing the nodes where traffic is originated, terminated or forwarded. There is an edge between vertices if the corresponding nodes are within each others transmission range. Each edge has a capacity, representing the throughput that is possible over that edge. The multicommodity flow problem then addresses the question whether there exists a set of paths and real numbers (fractions) so that: (1) for each traffic demand, there is a set of paths from the traffic source to the destination; (2) fractions of the traffic demand can be allocated to each path so that for each source-destination pair the total traffic

demand is realized and (3) the capacity constraints are taken into account.

This chapter considers a generic wireless network configuration specified via parameters such as nodes, transmission and interference ranges, as well as a traffic matrix indicating the demands between source nodes and sink nodes. We make no assumptions about the homogeneity of nodes with regard to transmission range or interference range, nor the capacity of the links. This is in contrast to previous work [GK00] that has focussed on asymptotic bounds under assumptions such as node homogeneity and random communication patterns. In [J⁺03] a conflict graph is used to address the problem of finding a feasible flow allocation to realize demands between pairs of nodes. While the conflict graph provides a more comprehensive modelling of the scheduling problem, it is also more complicated to deal with. A detailed discussion on the relation of our work with [J⁺03] is presented in Section 5.4, showing that we obtain a tight upper bound for an example provided in [J⁺03]. In [J⁺04] the question of a routing algorithm to find paths satisfying the traffic demands in a distributed setting is addressed. For an LP relaxation of the interference problem, [KN03] presents necessary conditions for link flow feasibility. This yields an upper bound similar to that of [J⁺03]. In addition, [KN03] introduces an edge colouring problem in which each colour at an edge represents a time slot for transmission. This problem is solved using a FPTAS, yielding a lower bound for the link flow allocation. In [KN05] this work is extended to multi-radio and multi-channel networks. Our work does not solve any LP's, but provides a good characterization for the feasibility of the fractional multiflow problem with interference constraints which provides a fast way of finding upper bounds for a slightly different interference constraint setting.

In this chapter we introduce a new approach to model interference in a carrier sensing multi-hop wireless network. To this end, we transform the sustainable load problem into a multicommodity flow problem that we extend with interference constraints. The main theorem of this chapter states a condition for the feasibility of the multicommodity flow problem with interference constraints, given the demands between source and destination nodes. Using this theorem, we compute the maximal throughput between a single source and a single destination. We consider the following elements to be the key contribution of our work:

- The use of polyhedral combinatorics (Farkas' Lemma) to obtain a structural expression for feasibility of the multicommodity flow problem with interference constraints, in terms of a 'generalized cut condition' analogous to the 'max-flow min-cut' theorem of Ford and Fulkerson [FF56].

- The generality of our framework which incorporates the following realistic effects: the transmission range is not necessarily equal to the interference range, the network may consist of mixed wired and wireless connections and wireless links have different capacities depending on distance, obstacles or transmission power.

The remainder of this chapter is organized as follows. In Section 2 we introduce interference constraints to model the wireless parts of the network

while taking the impact of the capacity of the links into account. Our main theorem is stated in Section 3, followed by examples and applications in Section 4. Section 5 extends the network to include multiple channels. In Section 6 we introduce interference constraints to model the wireless parts of the network while taking the impact of the capacity of the links and the now included available radios and channels into account. Our main theorem for this setting is stated in Section 7, followed by examples in Section 8. Section 9 concludes the chapter and indicates how the results can be used for a channel allocating algorithm.

## 5.2 Ad hoc interference model

Ad-hoc networks use transmissions over a wireless channel to communicate between users. However, if multiple transmissions take place at the same time over the same wireless channel, transmissions may collide and the data will be lost. This interference limits the throughput of ad-hoc networks. Adopting the model of [HBO06], we will model the interference constraints. To do so, we define the transmission range and the interference range of a node. When nodes in a wireless network want to communicate, they need to be close enough to receive each others signals. The *transmission range* of a node is the maximum distance from that node to where its received signal strength is sufficient for maintaining communications. Even though a signal may be too weak to be received correctly outside the transmission range, the signal can still cause interference preventing nodes from receiving other signals correctly. The *interference range* is the maximum distance from a node to where it prohibits other nodes to maintain communications. Note that in general the transmission and interference range are not equal. In the following we adopt a graph representation (see e.g. [GMW04],[J$^+$03]) in which these ranges will be represented by arcs.

Let $V$ denote the set of nodes, $A$ the set of arcs and let $\delta^+(u)$ denote the arcs leaving node $u$. In a carrier sensing network, nodes within each others interference range will avoid transmitting at the same time. We will model this as follows: Let $R(v)$ denote the set of nodes within the interference range of node $v$ (which includes $v$ itself), that is: if one of the nodes in $R(v)\backslash\{v\}$ is transmitting, then $v$ cannot receive other transmissions, nor can $v$ transmit. Let $\rho(u)$ denote the fraction of time that a node $u$ is transmitting, then

$$\sum_{u \in R(v)} \rho(u) \leq f(v) \qquad (5.1)$$

where $f(v)$ denotes the interference capacity of the node. The interference capacity denotes the amount of interference a node can handle and still transmit data itself. For a wired network one could set $f(v) = \infty$, whereas $f(v) = 1$ ensures that no two nodes within each others interference range transmit at the same time.

Consider the set $I(v)$ of all arcs leaving $v$, entering $v$ and leaving the nodes that are in $v$'s interference range:

$$I(v) = \{a | a \in \delta^+(u), u \in R(v)\}. \qquad (5.2)$$

Figure 5.1: $I(v)$ and $J(a)$

It follows that for all the arcs that are within $I(v)$ the interference capacity $f(v)$ may not be exceeded. In particular, if $f(v) = 1$, simultaneous transmissions cannot take place over arcs $a_1, a_2 \in I(v)$ . For later use, we also introduce here a dual notion of $I(v)$, viz. $J(a)$ , the set of vertices that experience interference by a transmission over arc $a$:

$$J(a) = \{v \in V | a \in I(v)\}. \tag{5.3}$$

Consider node $v \in V$. The interference arc set $I(v)$ is denoted using bold arcs in Figure 5.1(a) and the set of nodes $J(a)$ affected by arc $a$ by the grey nodes in Figure 5.1(b).

To each arc $a$ a capacity $b(a) > 0$ is assigned. In actual networks, due to e.g. unequal distances among nodes or external disturbances such as noise, the link capacities may be different. Consider a set of source and destination pairs $(r_1, s_1), ..., (r_k, s_k)$. When the net amount of flow between a source node $r_i$ and a destination node $s_i$ is $d_i$, we say that the value of the $(r_i, s_i)$ flow is $d_i$. For an allocation, let $x_i(a)$ denote the amount of traffic for source destination pair $(r_i, s_i)$ over link $a$. We will call this the flow of commodity $i$ over arc $a$. To take the link capacities into account, we use the following interference constraints:

$$\sum_{i=1}^{k} \sum_{a \in I(v)} \frac{x_i(a)}{b(a)} \leq f(v) \quad \forall v \in V, \tag{5.4}$$

Note that the interference constraints indicate whether a node can receive correctly and that we assume that collisions are fatal. However, as we impose condition (5.4) on all nodes, for a transmission over link $(u, v)$ to be successful, we have that (5.4) must hold for both $u$ and $v$, so both sender and receiver must be free of interference. This closely resembles the behaviour of IEEE 802.11 under RTS-CTS, where both sender and receiver must be free of interference, see e.g. [L+04]. For a successful communication, the sender must be able to hear the link layer acknowledgement transmitted by the receiver.

## 5.3 Multicommodity flow problem with interference constraints

The multicommodity flow problem (MCFP) describes the problem of finding an allocation of flows over links such that all flows are transferred from their source to their destination, without exceeding the capacity of the links. The *multicommodity flow problem with interference constraints* is as follows.

Given a graph $G(V, A)$, with link capacities $b : A \to \mathbb{R}^+$, interference capacities $f : V \to \mathbb{R}^+$ and source and destination pairs $(r_1, s_1), ..., (r_k, s_k)$ with demands $d_1, ..., d_k \in \mathbb{R}^+$, find for each $i = 1, ..., k$ an $(r_i, s_i)$ flow $x_i \in \mathbb{R}_+^{|A|}$ of value $d_i$, where $x_i(a)$ is the amount of traffic of commodity $i$ sent via arc $a$, and so that for each arc $a \in A$ and vertex $v \in V$ the capacity and interference constraints are met. Let $\delta^+(U) = \{a = (u, v) \in A | u \in U, v \notin U\}$ and $\delta^-(U) = \{a = (u, v) \in A | u \notin U, v \in U\}$ so that $\delta^+(v)$ and $\delta^-(v)$ denote the arcs leaving and entering node $v$ respectively. Our multicommodity flow problem with interference constraints has the following feasibility constraints:

$$\sum_{i=1}^{k} x_i(a) \leq b(a), \quad \forall a \in A \tag{5.5}$$

$$\sum_{a \in \delta^+(v)} x_i(a) = \sum_{a \in \delta^-(v)} x_i(a), \quad \forall v \in V, v \neq r_i, s_i \tag{5.6}$$

$$\sum_{a \in \delta^+(r_i)} x_i(a) - \sum_{a \in \delta^-(r_i)} x_i(a) = d_i, \quad \forall i \tag{5.7}$$

$$\sum_{a \in \delta^+(s_i)} x_i(a) - \sum_{a \in \delta^-(s_i)} x_i(a) = -d_i, \quad \forall i \tag{5.8}$$

$$\sum_{i=1}^{k} \sum_{a \in I(v)} \frac{x_i(a)}{b(a)} \leq f(v), \quad \forall v \in V \tag{5.9}$$

Equation (5.5) shows the capacity constraints on the arcs, equation (5.6) assures flow conservation, i.e. for each node the flow in must equal the flow out, and (5.7) and (5.8) define that the demands leaving the source and entering the destination. Note that (5.8) is redundant as it follows from (5.6) and (5.7), but is included here for completeness. Equation (5.9) is our interference constraint. Equations (5.5)-(5.8) define the multicommodity flow problem in its standard form, that is included in our formulation by setting the interference capacities of all nodes to infinity, that is $f(v) = \infty$ for all $v \in V$.

We now formulate a generalized cut condition for the multicommodity flow problem with interference constraints, so including (5.9). To this end, define length functions $l : A \to \mathbb{R}^+$ on all arcs and interference functions $w : V \to \mathbb{R}^+$ on all nodes.

For a given length function $l : A \to \mathbb{R}^+$ and interference function $w : V \to \mathbb{R}^+$, let $\mathbf{dist}_{l,w}(r_i, s_i)$ denote the distance function that incorporates both the length

and interference, where the distance of a path is built up of the distance $\mathbf{q}_{l,w}(a)$ of the arcs in the path as follows

$$\mathbf{q}_{l,w}(a) = l(a) + \sum_{t \in J(a)} \frac{w(t)}{b(a)} \tag{5.10}$$

$$\mathbf{q}_{l,w}(P) = \sum_{a \in P} \mathbf{q}_{l,w}(a) \tag{5.11}$$

$$\mathbf{dist}_{l,w}(r,s) = \min_{P \in P_{r,s}} \mathbf{q}_{l,w}(P). \tag{5.12}$$

with $J(a)$ as in (5.3) and $P_{r,s}$ the set of all paths from $r$ to $s$.

**Theorem 5.1.** *The multicommodity flow problem with interference constraints is feasible, if and only if for all length functions $l : A \to \mathbb{R}^+$ and node interference functions $w : V \to \mathbb{R}^+$ it holds that*

$$\sum_{i=1}^{k} d_i \mathbf{dist}_{l,w}(r_i, s_i) \leq \sum_{a \in A} l(a)b(a) + \sum_{v \in V} w(v)f(v). \tag{5.13}$$

*Proof.* Given a directed graph $G(V, A)$, arc capacities $b : A \to \mathbb{Q}^+$, node interference constraints $f : V \to \mathbb{Q}^+$, disjoint pairs $(r_1, s_1)...(r_k, s_k)$ and demands $d_1...d_k \in \mathbb{Q}^+$, the multicommodity flow problem with interference constraints can be written as an LP problem as follows:
Find $x = (x_1, ..., x_k)$ where $x_i : A \to \mathbb{Q}^+$ denotes the values of flow $r_i \to s_i$ assigned to an arc $a$ s.t.

$$Ax \leq b \tag{5.14}$$
$$Cx = \hat{d} \tag{5.15}$$
$$Ex \leq f \tag{5.16}$$

for $A = m \times mk$, $C = nk \times mk$, $E = n \times mk$ where $m = |A|$ and $n = |V|$, with $b = m \times 1$ denoting the capacity constraints, $\hat{d} = (\hat{d}_1, ..., \hat{d}_k) = nk \times 1$ denoting the flow constraints and $f = n \times 1$ denoting the interference constraints with:

$$A = [I_m, I_m, ..., I_m] \tag{5.17}$$

$$C = \begin{bmatrix} M & 0 & 0 \\ 0 & ... & 0 \\ 0 & 0 & M \end{bmatrix} \tag{5.18}$$

where $M$ is an $n \times m$ matrix defined by

$$M_{v,a} = \begin{cases} 1 & \text{when } a \text{ leaves } v \\ -1 & \text{when } a \text{ enters } v \\ 0 & \text{otherwise} \end{cases} \tag{5.19}$$

$$\hat{d}_i(v) \;\;=\;\; \begin{cases} d_i & \text{when } v \text{ is } r_i \\ -d_i & \text{when } v \text{ is } s_i \\ 0 & \text{otherwise} \end{cases} \tag{5.20}$$

$$E \;\;=\;\; [S, S, ..., S] \tag{5.21}$$

with $S$ an $n \times m$ matrix defined by

$$S_{v,a} = \begin{cases} \frac{1}{b(a)} & \text{when } a \in I(v) \\ 0 & \text{otherwise} \end{cases}. \tag{5.22}$$

According to Farkas' Lemma there exists a solution $x \geq 0$ satisfying (5.5)-(5.9) if and only if for all vectors $y, w \geq 0$ and $z$, with $y \in \mathbb{R}^m$, $w = (w_1, ..., w_n) \in \mathbb{R}^n$ and $z = (z_1, ..., z_k) \in n\mathbb{R}^{kn}$ where $z_i : V \to \mathbb{R}^n$ :

$$yA + zC + wE \geq 0 \Rightarrow yb + z\widehat{d} + wf \geq 0 \tag{5.23}$$

From the definitions of $A, C, E$ we find

$$yA \;\;=\;\; [y, y, ..., y] \tag{5.24}$$
$$zC \;\;=\;\; [z_1 M, z_2 M, ..., z_k M] \tag{5.25}$$
$$wE \;\;=\;\; [wS, wS, ..., wS], \tag{5.26}$$

where $z_i M$ is given by

$$z_i M = \sum_{v \in V} z_i(v) M_{v,a} = [z_i(u_1) - z_i(v_1), ..., z_i(u_m) - z_i(v_m)] \tag{5.27}$$

with $(u_j, v_j)$ denoting the starting and ending node of an arc $a_j$, and where the element of $wS$ corresponding to arc $a$ is given by

$$(wS)_a = \sum_{v \in V} w(v) S_{v,a} = \sum_{\substack{v \in V \\ a \in I(v)}} \frac{w(v)}{b(a)}. \tag{5.28}$$

As a consequence, (5.23) reads for all $y \geq 0, w \geq 0$ and $z_i \in \mathbb{R}^n$

$$z_i(v) - z_i(u) \leq y(a) + \sum_{t \in J(a)} \frac{w(t)}{b(a)}, \quad \forall i = 1...k, \forall a = (u, v) \in A \Rightarrow \tag{5.29}$$

$$\sum_{i=1}^{k} d_i(z_i(s_i) - z_i(r_i)) \leq \sum_{a \in A} y(a) b(a) + \sum_{v \in V} w(v) f(v).$$

We now show that there exists a feasible solution $x \geq 0$ if and only if for all length functions $l : a \to \mathbb{Q}^+$ and node interference functions $w : V \to \mathbb{R}^+$ it

holds that

$$\sum_{i=1}^{k} d_i \mathbf{dist}_{l,w}(r_i, s_i) \leq \sum_{a \in A} l(a)b(a) + \sum_{v \in V} w(v)f(v). \tag{5.30}$$

Suppose there is a feasible solution, then (5.30) holds, now choose $l(a) = y(a)$ as the length function, and (for all $i$) $z_i(s) - z_i(r)$ as the distance between $r$ and $s$, yielding

$$\sum_{i=1}^{k} d_i \mathbf{dist}_{l,w}(r_i, s_i) \quad = \quad \sum_{i=1}^{k} d_i(z_i(s_i) - z_i(r_i)) \tag{5.31}$$

$$\leq \quad \sum_{a \in A} l(a)b(a) + \sum_{v \in V} w(v)f(v) \tag{5.32}$$

Next suppose that (5.30) holds, we will now show that also (5.30) holds. Let the minimizing path use the arcs $(a_1, ..., a_p)$, then

$$\sum_{a \in A} l(a)b(a) + \sum_{v \in V} w(v)f(v) \quad \geq \quad \sum_{i=1}^{k} d_i \mathbf{dist}_{l,w}(r_i, s_i) \tag{5.33}$$

$$= \quad \sum_{i=1}^{k} d_i \sum_{j=1}^{p} \mathbf{q}_{l,w}(a_j)$$

$$= \quad \sum_{i=1}^{k} d_i \left( \sum_{j=1}^{p} l(a_j) + \sum_{t \in J(a_j)} \frac{w(t)}{b(a_j)} \right)$$

Let $z_i : V \to \mathbb{R}$ be so that $z_i(v) - z_i(u) \leq y(a) + \sum_{t \in J(a)} \frac{w(t)}{b(a)}$, which in combination with $l(a) = y(a)$ gives that

$$\sum_{i=1}^{k} d_i \left( \sum_{j=1}^{p} l(a_j) + \sum_{t \in J(a_j)} \frac{w(t)}{b(a_j)} \right) \geq \sum_{i=1}^{k} d_i \sum_{j=1}^{p} z_i(v_j) - z_i(u_j) \tag{5.34}$$

where $a_j = (u_j, v_j)$ and $u_1 = r_i$, $v_j = u_{j+1}$ and $v_p = s_i$ so that the right hand side of the expression simplifies to

$$\sum_{i=1}^{k} d_i \sum_{j=1}^{p} z_i(v_j) - z_i(u_j) = \sum_{i=1}^{k} d_i(z_i(s_i) - z_i(r_i)) \tag{5.35}$$

which taken together with (5.33) gives

$$\sum_{a \in A} l(a)b(a) + \sum_{v \in V} w(v)f(v) \geq \sum_{i=1}^{k} d_i(z_i(s_i) - z_i(r_i)) \tag{5.36}$$

completing the proof.                                                              □

The interference function $w(v)$ is a dual variable that can be interpreted as the price paid for the interference capacity of a node $v$, which gives a weighted cut of the nodes (as the length function $l(a)$ can be interpreted as the price paid for the capacity of an arc $a$, which also gives a weighted cut of the arcs). The distance of an arc is the price paid for the capacity of that arc, together with the price paid for the interference capacity used by that arc. Note that our theorem is an extension of the cut condition for the multicommodity flow problem (without interference). This can be seen as follows. Setting all interference capacities to a very large value the condition of Theorem 5.1 reduces to

$$\sum_{i=1}^{k} d_i \mathbf{dist}_l(r_i, s_i) \leq \sum_{a \in A} l(a) b(a) \tag{5.37}$$

as the inequality only makes sense for $w = 0$. The cut condition by setting a length of 1 to all arcs in a cut and zero otherwise for the multicommodity flow problem states that

$$\sum_{a \in \delta^+(U)} b(a) \geq \sum_{\substack{r_i \in U \\ s_i \notin U}} d_i . \tag{5.38}$$

This cut condition is necessary for the existence of a solution, but not sufficient. The max-flow min-cut theorem states that for the specific case that $k = 1$, for every network, there exists a flow (max-flow) for which the amount is equal to the total capacity of the smallest cut in the network (min-cut), see Ford and Fulkerson [FF56].

A direct consequence of Theorem 5.1 is that when there is only one commodity, a bound on the throughput $d$ of the network can be determined by

$$d = \frac{\sum_{a \in A} l(a) b(a) + \sum_{v \in V} w(v) f(v)}{\mathbf{dist}_{l,w}(r, s)}. \tag{5.39}$$

When there are multiple commodities with demands $d_1, ..., d_k$, Theorem 5.1 can determine the maximal value $0 \leq \lambda \leq 1$ such that for all commodities a throughput of $\lambda d_i$ can be achieved using

$$\lambda = \frac{\sum_{a \in A} l(a) b(a) + \sum_{v \in V} w(v) f(v)}{\sum_{i=1}^{k} d_i \mathbf{dist}_{l,w}(r_i, s_i)}. \tag{5.40}$$

## 5.4   Examples

We can consider the network as depicted in Figure 5.2 as nodes in a network, connected by links using 802.11b with a maximum transmission rate of 11 Mbit/s, but with link 3 having a bad connection, due to distance or a disturbance, only reaching the minimal transmission rate of 5.5 Mbit/s.

We want to transmit data from node 1 to node 5. If we solve for the best solution without considering interference, it is clear that we can send at a speed of 5.5 Mbit/s, as this is limited by the slowest link. The interference constraints

Figure 5.2: Series of nodes with capacity constraints

for the network with identical slowest link capacities would imply that all links can be used one third of the time, leading to an overall throughput of 1.83 Mbit/s when considering the constraints separately.

For the example of Figure 5.2 we have the interference constraints

$$\frac{x(1)}{11} + \frac{x(2)}{11} + \frac{x(3)}{5.5} \leq 1 \tag{5.41}$$

$$\frac{x(2)}{11} + \frac{x(3)}{5.5} + \frac{x(4)}{11} \leq 1. \tag{5.42}$$

From a direct solution of (5.5)-(5.9) it follows that $x(a) = 2.75$ is a feasible solution, higher than the earlier claimed 1.83 Mbit/s. Using Theorem 5.1, we find that

$$\begin{aligned} d_1(l(1) + l(2) + l(3) + l(4) + \\ \frac{2w(1) + 4w(2) + 4w(3) + 3w(4) + w(5)}{11}) \leq \\ 11l(1) + 11l(2) + 5.5l(3) + 11l(4) + \\ w(1) + w(2) + w(3) + w(4) + w(5) \end{aligned} \tag{5.43}$$

which gives by taking the cut $w(3) = 1$ and all other values (including $l(i)$) equal to zero

$$\frac{4}{11}d_1 \leq 1. \tag{5.44}$$

So $x(a) = 2.75$ is also the optimal solution. The value now found for $x(a)$ can be interpreted as the fraction of time a link is in use multiplied by the transmission rate of the link. This shows that links one, two and four get $\frac{1}{4}^{th}$ of the time, whereas link three gets $\frac{1}{2}$ of the time. So when considering four slots, link one, two and four each get one slot (where link one and four use the same slot!) and link three the other two. This way we have an accurate representation of the network incorporating both the capacity and interference constraints, together with the flow conservation laws.

We now consider the more sophisticated network used in [J$^+$03] as depicted in Figure 5.3, where all arcs have capacity 1. The upper bound on the throughput from node 0 to node 8 for this network obtained in [J$^+$03] is 0.667, opposed to the optimal 0.5, even though their algorithm has discovered all possible cliques in the conflict graph. Using Theorem 1 and taking the cut $w_1 = w_3 = 1$ or $w_4 = 1$ gives the lowest upper bound that can be achieved, resulting in $d_1 \leq 0.5$, which is tight. The reason we obtain a different result than in [J$^+$03] is that we use constraints for all nodes, so that considering for example node 0 we have

Figure 5.3: A 3x3 grid

that $x(3) + x(9) \leq 1$, as signals transmitted from node 1 and 3 reach node 0. In the approach of Jain et al., arcs 3 and 9 are not connected in the conflict graph, as a simultaneous transmission over both arcs is possible.

There is an interesting relation between the approach presented in this chapter and the results of [J$^+$03]. Jain et al. use a *conflict graph* to determine lower and upper bounds for the throughput of the network. The conflict graph $C$ has vertices corresponding to the arcs in the transmission graph, where there is an edge between two arcs if and only if the arcs are not allowed to transmit simultaneously. In our approach, $C$ has as vertex set $A$, where there is an edge between $a_i$ and $a_j$ (for some $1 \leq i, j \leq |A|$) if and only if $\exists v \in V$ s.t. $a_i, a_j \in I(v)$. Note that $I(v)$ defines a clique in $C$ for each $v$ in $V$. In fact, our interference model adopted here resembles the *protocol model* of [J$^+$03], but it is 'stricter' in the sense that for the same network, we have more interference constraints (edges in the conflict graph) than [J$^+$03].

In [J$^+$03] it is shown that a vector $x_i : A \rightarrow \mathbb{R}^+$ (corresponding to a flow $i$), can be scheduled without interference conflicts if and only if $x_i$ lies in the stable set polytope of $C$. (The stable set polytope is the convex hull of the incidence vectors of the stable sets in the graph). It is well-known that the stable set polytope is contained in the *fractional* stable set polytope. (The fractional stable set polytope is defined by all constraints indicating that the total flow in a maximal clique in the conflict graph is at most 1.)
In this chapter, instead of first defining the conflict graph $C$ and then discovering its cliques, we directly formulate inequalities corresponding to the cliques $I(v)$ for all $v \in V$. The polytope defined by these inequalities will therefore contain the fractional stable set polytope. As a result, the upper bound obtained here cannot always be achieved using a flow allocation without interference conflicts.

## 5.5   Extension to multi-channel

We now extend the problem taking into consideration the availability of multiple radios, so that different channels can be used for simultaneous transmissions that do not interfere with each other. As such, this part of the chapter presents a generic model for sustainable network load in a multi-hop multi-channel setting by again recasting the model in a multicommodity flow problem with interference constraints and stating a theorem which gives a necessary and sufficient condition for the feasibility of this multicommodity flow problem. From this theorem an upper bound is derived for the throughput that can be achieved by the network, which is illustrated by examples. We indicate how the results can be used as a basis for a channel allocating algorithm.

One way to overcome the loss in capacity due to interference is the use of different channels. This use of different channels requires the nodes of the network to be equipped with multiple radios. To optimize the performance of the network, an allocation of the channels to these radios needs to be found that maximizes the throughput, the main performance measure of the network under consideration. These aspects call for a model that takes into consideration the interference and the channel allocation in a multi-hop multi-channel wireless network.

In this part we provide a theorem that can be used to show bounds similar to those in [GK00] and the extension to multi-channel networks as presented in [KV05]. For completeness and autonomicity, we repeat some of the definitions. We consider a generic wireless network configuration and traffic load specified via parameters such as nodes, radios, channels, transmission and interference ranges, as well as the traffic matrix indicating the demands between source and sink nodes. We make no assumptions about the homogeneity of nodes with regard to transmission range, interference range or number of radios, nor the capacity of the links. By introducing interference capacity as a parameter, we model the impact of interference. The choice of this parameter can in some degree model different mechanisms such as TDMA and CSMA. In the current setting, we assume transmissions can follow each other instantly, thus lying close to TDMA. Less efficient mechanisms or mechanisms with collision avoidance such as CSMA can be modeled to some extent as well, but this lies outside the scope of this chapter.

We use the same definitions of the transmission range and interference range as for the case with one channel, which now holds for each separate channel. Also, the graph representation with these ranges represented by arcs is adopted. We let $V$ denote the set of nodes, $A$ the set of arcs, $G$ the set of available channels, and $K$ the set of commodities, that is the pairs of nodes that want to communicate.

In a carrier sensing network, nodes within each others interference range will avoid transmitting at the same time. We will model this similar to the one channel case, but include the different channels: Let $I^g(v)$ denote the set of arcs which use will prevent node $v$ from transmitting on the same channel $g$, that is: if one of the arcs in $I^g(v)$ is in use for a transmission on channel $g$, then $v$ cannot receive any other transmissions, nor can $v$ transmit on the same channel.

Let $\rho^g(a)$ denote the fraction of time that an arc $a$ is transmitting over channel $g$, then

$$\sum_{a \in I^g(v)} \rho^g(a) \leq f^g(v) \quad \forall g \in G \tag{5.45}$$

where $f^g(v)$ denotes the interference capacity of the node for a channel $g$. The interference capacity denotes the amount of interference a node can handle and still transmit data itself. For a wired network one could set $f^g(v) = \infty$. The constraint for $f^g(v) = 1$ can be derived from the fact that no two arcs in $I^g(v)$ can be simultaneously used on the same channel for transmission. For later use, we also introduce here a dual notion of $I^g(v)$, viz. $J^g(a)$, the set of vertices that experience interference by a transmission over arc $a$ using channel $g$:

$$J^g(a) = \{v \in V | a \in I^g(v)\}. \tag{5.46}$$

This situation is equal to the situation depicted in Figure 5.1.

In literature, different interference models are used, e.g. the protocol model in [GK00] where a transmission from node $u$ to node $v$ is successful if and only if $|w - v| \geq (1 + \Delta)|u - v|$ for any other node $w$ transmitting on the same channel at the same time, where the quantity $\Delta$ models a guard zone preventing other nodes from transmitting at the same time over the same channel and $|u - v|$ is the distance between nodes $u$ and $v$. So all arcs starting within a distance of $(1 + \Delta)|u - v|$ around node $v$ are included in $I^g(v)$, whereas in $J^g(a)$ all nodes within distance $(1 + \Delta)|u - v|$ of $u$ are included. In [GK00] a node can use an arbitrary transmission power, thus influencing the interference range, we however assume that the set $I^g(v)$ only depends on $v$ and not on the node it is receiving from, which is the case when all nodes transmit at the same power. As then there is no dependence upon the distance between the nodes, this setting is easily modelled.

To each arc $a$ a capacity $b^g(a) > 0$ is assigned for all channels $g \in G$, which can only be used if both endnodes have radios set to this channel. In actual networks, due to e.g. unequal distances among nodes or external disturbances such as noise, the link capacities may be different.

Consider a set of source and destination pairs $(r_1, s_1), ..., (r_K, s_K)$. When the net amount of flow between a source node $r_k$ and a destination node $s_k$ is $d_k$, we say that the value of the $(r_k, s_k)$ flow is $d_k$. For an allocation, let $x_k^g(a)$ denote the amount of traffic for source-destination pair $(r_k, s_k)$ over link $a$, using channel $g$. We will call this the flow of commodity $k$ over arc $a$ using channel $g$. To take both the link capacities and interference capacities into account, we use the following interference constraints:

$$\sum_{k \in K} \sum_{a \in I(v)} \frac{x_k^g(a)}{b^g(a)} \leq f^g(v) \quad \forall g, \forall v \in V, \tag{5.47}$$

Note that the interference constraints indicate whether a node can receive correctly and that we assume that collisions are fatal. However, as we impose condition (5.47) on all nodes, for a transmission over arc $a = (u, v)$ to be

successful, we have that (5.47) must hold for both $u$ and $v$, so both sender and receiver must be free of interference. This closely resembles the behaviour of IEEE 802.11 under RTS-CTS, where both sender and receiver must be free of interference. For a successful communication, the sender must be able to hear the link layer acknowledgement transmitted by the receiver.

## 5.6    Multi-channel multicommodity flow problem with interference constraints

In the first part of the chapter, we assumed that only one channel is available so that all transmissions interfere with each other (within the interference range). Here we assume there are multiple channels available over which transmissions can take place and that these channels are such that no interference occurs between transmissions on different channels. All nodes have a number of radios that can each be set to a different channel and can all be active at the same time, either sending or receiving a transmission. The constraints for the feasibility of a flow through a network are now extended taking into account the use of multiple channels. To this end, we introduce the variables $1^g(v) \in \{0, 1\}$ to denote whether a node has one of its radios set to channel $g$. Likewise, we define $1^g(a) \in \{0, 1\}$ to denote whether both endpoints of an arc have radios set to channel $g$. The flow variables are given by $x_k^g(a) \geq 0$, the amount of flow for commodity $k$ sent over link $a$ using channel $g$. Each link $a$ has a capacity $b^g(a)$ when using channel $g$. Adopting the notation of [KN05], a node $v$ has $\kappa(v)$ radios available and has an interference capacity $f^g(v)$ for channel $g$. The problem of finding an allocation $x_k^g(a)$ satisfying demand $d_k$ then has the following constraints for feasibility (using $\delta(v)$ for all arcs connected to node $v$, and $\delta^+(v)$ and $\delta^-(v)$ for the out- and ingoing arcs respectively)

$$\sum_{g \in G} 1^g(v) \leq \kappa(v) \qquad \forall v \in V \tag{5.48}$$

$$\sum_{a \in \delta(v)} 1^g(a) \leq \delta(v) 1^g(v) \qquad \forall v \in V, \forall g \in G \tag{5.49}$$

$$\sum_{k \in K} x_k^g(a) \leq b^g(a) 1^g(a) \qquad \forall a \in A, \forall g \in G \tag{5.50}$$

$$\sum_{g \in G} \left( \sum_{a \in \delta^+(v)} x_k^g(a) - \sum_{a \in \delta^-(v)} x_k^g(a) \right) = 0 \qquad \forall v \notin \{r_k, s_k\}, \forall k \in K \tag{5.51}$$

$$\sum_{g \in G} \left( \sum_{a \in \delta^+(r_k)} x_k^g(a) - \sum_{a \in \delta^-(r_k)} x_k^g(a) \right) = d_k \qquad \forall k \in K \tag{5.52}$$

$$\sum_{g \in G} \left( \sum_{a \in \delta^+(s_k)} x_k^g(a) - \sum_{a \in \delta^-(s_k)} x_k^g(a) \right) = -d_k \qquad \forall k \in K \tag{5.53}$$

$$\sum_{k \in K} \sum_{a \in I^g(v)} \frac{x_k^g(a)}{b^g(a)} \leq f^g(v) \qquad \forall g \in G, \forall v \in V \tag{5.54}$$

Constraints (5.48) and (5.49) respectively assure that a node does not exceed the number of radios it has available and that an arc can only use a certain channel if both endnodes have a radio set to this channel. Constraint (5.50) is the capacity constraint, which also incorporates the constraint that if a channel is not used for a certain arc, there will be no flow on this channel over the arc. Constraints (5.51)-(5.53) are the flow conservation constraints and make sure the demanded flow leaves the source nodes and enters the sink nodes for all commodities. Note that flows entering an intermediate node over a certain channel may leave this node over a different channel, hence the summation over all channels. Finally, (5.54) gives the interference constraints, where all channels are considered independently as they do not interfere with each other. Important to note is that we assume that all different radios at the same node can transmit at the same time (as long as there are enough links available for these channels), so multiple transmissions from the same node do not interfere if the used channels are different. This approach can be used as it is our intention to derive an upper bound on the network throughput.

The problem stated is a mixed integer programming problem and is in general hard to solve. In this problem, we need to assign channels to the radios of all the nodes and find a flow value for all the arcs. When the channel allocation is given, the problem reduces considerably in complexity, as it becomes a linear programming problem. In the remainder of this chapter, we assume that the allocation of channels to radios for all nodes is such that constraints (5.48) and (5.49) hold, so that the remaining constraints are (5.50)-(5.54) where the only unknown variables are the $x_k^g(a)$'s.

To formulate a generalized cut condition as in [FF56], we define length functions $y = [y^1, ..., y^G]$ with $y^g : A \to \mathbb{R}^+$ on all arcs for each separate channel and interference functions $w = [w^1, ..., w^G]$ with $w^g : V \to \mathbb{R}^+$ on all vertices. Let $\mathbf{dist}_{y,w}(r, s)$ be the distance between two nodes $r$ and $s$, taking into account both the length and interference functions, defined as

$$q_{y,w}(a) = \min_{g \in G}(y^g(a) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)}) \quad (5.55)$$

$$q_{y,w}(P) = \sum_{a \in P} q_{y,w}(a) \quad (5.56)$$

$$\mathbf{dist}_{y,w}(r, s) = \min_{P \in P_{r,s}} q_{y,w}(P) \quad (5.57)$$

The distance of a path is built up of the distances of it's arcs, defining the total distance between two nodes $r$ and $s$ by the path with the smallest distance of all paths $P_{r,s}$ between $r$ and $s$. We now arrive at our main theorem:

**Theorem 5.2.** *The MCFP problem with interference constraints and given channel allocation has a solution iff for all length functions $y^g : A \to \mathbb{R}^+$ and node interference functions $w^g : V \to \mathbb{R}^+$ it holds that*

$$\sum_{k \in K} d_k \mathbf{dist}_{y,w}(r_k, s_k) \leq \quad (5.58)$$

$$\sum_{g \in G} (\sum_{a \in A} y^g(a) b^g(a) 1^g(a) + \sum_{v \in V} w^g(v) f^g(v)).$$

*Proof.* We define the vector $x$ as

$$x = [x_1^1(a_1), ..., x_1^1(a_m), x_1^2(a_1), ..., x_1^G(a_1), ..., x_K^G(a_m)] \tag{5.59}$$

where $m = |A|$. According to Farkas' Lemma (c.f. [Sch03]) there exists a solution $x \geq 0$ satisfying (5.50)-(5.54) iff for all vectors $y \in \mathbb{R}^{Gm}, y \geq 0$, $w \in \mathbb{R}^{Gn}, w \geq 0$ and $z \in \mathbb{R}^{Kn}$, it holds that:

$$y^g(a) + z_k(p) - z_k(q) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)} \geq 0 \tag{5.60}$$

$$\forall g \in G, \ a \in A, \ \forall k \in K \Rightarrow$$

$$\sum_{a \in A} \sum_{g \in G} y^g(a) b^g(a) 1^g(a) + \sum_{k \in K} d_k(z_k(r_k) - z_k(s_k))$$

$$+ \sum_{g \in G} \sum_{v \in V} w^g(v) f^g(v) \geq 0.$$

where we use that $a = (p, q)$. Assume that (5.58) holds, we will now show that there exists a solution to the MCFP. Suppose that a path of minimal distance for commodity $k$ uses the arcs $(a_k^1, ..., a_k^P)$, then

$$\sum_{g \in G} (\sum_{a \in A} y^g(a) b^g(a) 1^g(a) + \sum_{v \in V} w^g(v) f^g(v)) \tag{5.61}$$

$$\geq \sum_{k \in K} d_k \sum_{j=1}^{P} q_{y,w}(a_k^j)$$

$$= \sum_{k \in K} d_k \sum_{j=1}^{P} \min_{g \in G} (y^g(a_k^j) + \sum_{v \in J^g(a_k^j)} \frac{w^g(v)}{b^g(a_k^j)})$$

Let the first part of the condition in Farkas' Lemma hold, we have to show that also the right hand side of the implication in (5.60) holds. As the first part holds (by assumption) we have that $z_k(h) - z_k(g) \leq y^g(a) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)}$ for all $g \in G$, $a \in A$, and $k \in K$. This implies that for any arc $a = (p, q)$

$$z_k(q) - z_k(p) \leq \min_{g \in G} (y^g(a) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)}) \tag{5.62}$$

so then

$$\sum_{k \in K} d_k \sum_{j=1}^{P} \min_{g \in G} (y^g(a_k^j) + \sum_{v \in J^g(a_k^j)} \frac{w^g(v)}{b^g(a_k^j)}) \tag{5.63}$$

$$\geq \sum_{k \in K} \sum_{j=1}^{P} d_k(z_k(q_k^j) - z_k(p_k^j))$$

$$= \sum_{k \in K} d_k(z_k(s_k) - z_k(r_k))$$

taking together (5.61) and (5.63) gives

$$\sum_{g \in G}(\sum_{a \in A} y^g(a)b^g(a)1^g(a) + \sum_{v \in V} w^g(v)f^g(v)) \qquad (5.64)$$

$$\geq \sum_{k \in K} d_k(z_k(s_k) - z_k(r_k))$$

showing the right hand side of Farkas' Lemma holds whenever the left hand side holds, proving the existence of a solution.

We have shown that if (5.58) holds, a solution exists, we now prove the reverse. Assuming that a solution exists and so (5.60) holds, we show that (5.58) holds. Hence, we have that if

$$y^g(a) + z_k(p) - z_k(q) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)} \geq 0, \qquad (5.65)$$

$$\forall g \in G, \ a \in A, \ \forall k \in K,$$

then also

$$\sum_{k \in K} d_k(z_k(s_k) - z_k(r_k)) \qquad (5.66)$$

$$\leq \sum_{g \in G}(\sum_{a \in A} y^g(a)b^g(a)1^g(a) + \sum_{v \in V} w^g(v)f^g(v))$$

Now we define $z_k(v) := \mathbf{dist}_{y,w}(r_k, v)$, so that the right hand side of (5.60) reduces to

$$\sum_{k \in K} d_k(\mathbf{dist}_{y,w}(r_k, s_k) - \mathbf{dist}_{y,w}(r_k, r_k)) \qquad (5.67)$$

$$= \sum_{k \in K} d_k \mathbf{dist}_{y,w}(r_k, s_k)$$

$$\leq \sum_{g \in G}(\sum_{a \in A} y^g(a)b^g(a)1^g(a) + \sum_{v \in V} w^g(v)f^g(v))$$

which shows that (5.58) holds, if the left hand side of (5.60) is satisfied. To show that this is the case, consider the distance between two nodes $(p, q)$, which are connected by the arc $a$. We then have that

$$\mathbf{dist}_{y,w}(r_k, p) + q_{y,w}(a) \geq \mathbf{dist}_{y,w}(r_k, q) \qquad (5.68)$$

as the distance is defined by the minimizing path which for node $r_k$ does not

necessarily lead through node $q$. It now follows that

$$
\begin{aligned}
q_{y,w}(a) &\geq \mathbf{dist}_{y,w}(r_k, q) - \mathbf{dist}_{y,w}(r_k, p) & (5.69) \\
&= z_k(q) - z_k(p)
\end{aligned}
$$

Using the definition of $q_{y,w}(a)$ this leads to

$$
\min_{g \in G}(y^g(a) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)}) \geq z_k(q) - z_k(p). \tag{5.70}
$$

As the inequality holds for the minimum over $g$, it will hold for any $g \in G$, so we have shown that

$$
y^g(a) + \sum_{v \in J^g(a)} \frac{w^g(v)}{b^g(a)} \geq z_k(q) - z_k(p), \tag{5.71}
$$

$$
\forall g \quad \in \quad G,\ a \in A,\ \forall k \in K,
$$

being the left hand side of (5.60), holds, thus completing the proof. $\qquad\square$

Any choice for the length function $y$ and interference function $w$, which together we call a cut, defines a bound on the possible throughput. In absence of interference constraints, for a single commodity this theorem states that there exists a flow (max-flow) for which the amount is equal to the total capacity of the smallest cut in the network (min-cut), c.f. [FF56],[FF62].

## 5.7    Examples

Consider again the simple situation where there is one path from source node 1 to destination node 5 as depicted in Figure 5.2 and each node, irrespective which channel is in use, can hear transmissions from its direct neighbours so that $I^g(v) = \{a | a \in \delta^+(u), u \in N(v)\}$, where $N(v)$ is the set of all neighbours of $v$. We assume that each node has two radios, there are three available channels and $b^g(a) = 1$ for all arcs and channels. At first, we set all radios to use channel one as if its a single channel problem. An achievable throughput for this setting is $d_1 = \frac{1}{3}$, using each arc for a third of the time, that is $x_1^1(a) = \frac{1}{3}$ for all arcs. The constraint given by the theorem (c.f. [CdGB07]) becomes (ignoring the other channels)

$$
\begin{aligned}
d_1(y_1 + y_2 + y_3 + y_4 + 2w_1 + 3w_2 + 3w_3 + 2w_4 + w_5) & \quad (5.72) \\
\leq \quad y_1 + y_2 + y_3 + y_4 + w_1 + w_2 + w_3 + w_4 + w_5
\end{aligned}
$$

which has as optimal cut $w_2 = 1$ or $w_3 = 1$ and all other values set to 0, leading to an upper bound for the throughput $d_1$ of $\frac{1}{3}$. This proves immediately that this is also the optimal throughput that can be achieved by this network. To try to increase the throughput, we choose to set link 2 to channel two, meaning that radios at node 2 and 3 are set to channel 2, as it interferes with both node 2 and 3 (but we could have chosen link 3 as well). This seems an obvious choice as the

cut showed that this is where the 'bottleneck' is situated. For the new situation, a possible allocation is $x_1^1(a_1) = x_1^2(a_2) = x_1^1(a_3) = x_1^1(a_4) = \frac{1}{2}$ obtaining a throughput of $d_1 = \frac{1}{2}$. The constraint given by the theorem becomes

$$d_1(\min \left\{ \begin{array}{l} y_1^1 + w_1^1 + w_2^1 \\ y_1^2 + w_1^2 + w_2^2 \end{array} \right. + \min \left\{ \begin{array}{l} y_2^1 + w_1^1 + w_2^1 + w_3^1 \\ y_2^2 + w_1^2 + w_2^2 + w_3^2 \end{array} \right. \tag{5.73}$$

$$+ \min \left\{ \begin{array}{l} y_3^1 + w_2^1 + w_3^1 + w_4^1 \\ y_3^2 + w_2^2 + w_3^2 + w_4^2 \end{array} \right. + \min \left\{ \begin{array}{l} y_4^1 + w_3^1 + w_4^1 + w_5^1 \\ y_4^2 + w_3^2 + w_4^2 + w_5^2 \end{array} \right. )$$

$$\leq y_1^1 + y_2^2 + y_3^2 + y_3^1 + y_4^1 + w_1^1 + w_2^1$$

$$+ w_3^1 + w_4^1 + w_5^1 + w_1^2 + w_2^2 + w_3^2 + w_4^2 + w_5^2$$

which has as an optimal cut $w_2^1 = w_3^1 = y_1^2 = y_4^2 = 1$, $y_3^2 = 2$ and all other values 0, giving an upper bound of $d_1 \leq \frac{1}{2}$, again proving this to be the optimal throughput. Now we set link 3 to channel 3 (as this appears to be the bottleneck point) by changing the radios at node 3 and 4 from channel 1 to 3. An obvious allocation now is

$$x_1^1(a_1) = x_1^2(a_2) = x_1^3(a_3) = x_1^1(a_4) = 1 \tag{5.74}$$

giving a throughput of $d_1 = 1$. The constraint given by the theorem (now including all possible channels) is

$$d_1(\min \left\{ \begin{array}{l} y_1^1 + w_1^1 + w_2^1 \\ y_1^2 + w_1^2 + w_2^2 \\ y_1^3 + w_1^3 + w_2^3 \end{array} \right. + \min \left\{ \begin{array}{l} y_2^1 + w_1^1 + w_2^1 + w_3^1 \\ y_2^2 + w_1^2 + w_2^2 + w_3^2 \\ y_2^3 + w_1^3 + w_2^3 + w_3^3 \end{array} \right. \tag{5.75}$$

$$+ \min \left\{ \begin{array}{l} y_3^1 + w_2^1 + w_3^1 + w_4^1 \\ y_3^2 + w_2^2 + w_3^2 + w_4^2 \\ y_3^3 + w_2^3 + w_3^3 + w_4^3 \end{array} \right. + \min \left\{ \begin{array}{l} y_4^1 + w_3^1 + w_4^1 + w_5^1 \\ y_4^2 + w_3^2 + w_4^2 + w_5^2 \\ y_4^3 + w_3^3 + w_4^3 + w_5^3 \end{array} \right. )$$

$$\leq y_1^1 + y_2^2 + y_3^3 + y_4^1 + w_1^1 + w_2^1 + w_3^1 + w_4^1 + w_5^1$$

$$+ w_1^2 + w_2^2 + w_3^2 + w_4^2 + w_5^2 + w_1^3 + w_2^3 + w_3^3 + w_4^3 + w_5^3$$

with optimal cut $y_a^g = 1$ for all arcs and channels and $w_k^g = 0$ for all nodes and channels, giving an upper bound and thus, as this is also an achievable throughput, optimal throughput of $d_1 = 1$. It turns out that for this setting with three channels and two radios, no higher throughput can be obtained, so the assignment of channel 1 to link 1 and 4, channel 2 to link 2 and channel 3 to link 3, is (one of) the optimal channel allocations. Note however that we do assume that the network transmits over link 1 and 4 over the same channel at the exact same time. Whether a protocol can assure this to be the case in a real network remains open, but the found value anyway holds as an upper bound for the achievable throughput.

Another example, based on the approach of Gupta and Kumar [GK00] for which they obtain their well known upper bound on the capacity of wireless networks, is a network consisting of $n$ nodes that all transmit at a rate $\lambda$ to a neighbouring node. In our model, when considering only one channel, we

hence set the number of commodities $K$ equal to $n$, the number of nodes, and all capacities $b(a)$ to one. The demand for each commodity $d_k$ equals $\lambda$ and the source-destination pairs are neighbouring nodes. First consider the one dimensional case with one channel where nodes are distributed over a line. For any arc $a$ and a fixed transmission range $r$, the set $J(a)$ can be seen as all the nodes within a distance $\Phi = (1 + \Delta)r$ of the startnode $v^+$ of the arc, which is an interval around the node. Now we can use the following approach to find an upper bound on the throughput: Divide the line of length $L$ into equal parts of length $\Phi$. In each part, select one node $v$ (when a node is available) for which the value $w(v)$ is set to $\frac{1}{n}$, whereas we keep the value of all $y(a)$ at 0. As any other node within the same part has a distance less than $\Phi$ to a chosen node, it is always within the interference range of the chosen node and corresponding arc. For each commodity, which only uses one arc $a$ as all transmissions are to a neighbour, we thus have chosen at least one node $v$ in $J(a)$ to be set to $\frac{1}{n}$, showing that the left hand side of (5.58) is larger than $\lambda$. As the number of nodes for which the value of $w(v)$ has been set to $\frac{1}{n}$ is smaller than or equal to $\frac{L}{\Phi}$, the theorem states that

$$\lambda \leq \frac{L}{\Phi n} \tag{5.76}$$

giving an upper bound on the possible throughput per node $\lambda$. However, in the work of Gupta and Kumar, it is not just the rate $\lambda$ that is optimized, but the distance covered is taken into account as well. As the average distance $\bar{L}$ between nodes is given by $\frac{L}{n}$, and each node transmits at rate $\lambda$, we have for the total throughput of the network $T = \lambda n \bar{L}$ that

$$T \leq \frac{L^2}{\Phi n}. \tag{5.77}$$

In the two dimensional case, again the interference set $J(a)$ can be seen as all the nodes within a distance $\Phi$ of the startnode $v^+$ of the arc, being a circle around the node. We assume that the total area $A$ under consideration is a square with side $\sqrt{A}$. We divide this area into small squares with side $\frac{\Phi}{\sqrt{2}}$, and when available, choose a node in each of these squares for which we set the value of $w(v)$ to $\frac{1}{n}$. Again it holds that any other node also contained in a square is covered by the chosen node, as the maximum distance between two nodes in a square is $\Phi$, so within interference range. The number of chosen nodes is limited by $\frac{2A}{\Phi^2}$, and taking the average distance over which a transmission takes place into account we get from the theorem that

$$T \leq \frac{2A}{\Phi^2} \bar{L}. \tag{5.78}$$

Comparing to the setting used in Gupta and Kumar, consider the situation where all nodes are ordered in a grid, with equal distances between all nodes. In this case, the distance travelled by a transmission is equal to $\bar{L} = \frac{\sqrt{A}}{\sqrt{n}}$ and the radius of the circles around the sending nodes is given by $\Phi = (1 + \Delta)\bar{L}$, as in

[GK00] this depends on the transmission distance. We hence obtain

$$T \leq \frac{2A}{((1+\Delta)\bar{L})^2}\bar{L} = \frac{2\sqrt{A}\sqrt{n}}{(1+\Delta)^2} \tag{5.79}$$

showing, apart from a constant value, the same results as obtained in Gupta and Kumar. For multiple channels, it is shown that when each node has as many radios as there are channels, the same bound holds. This can easily be seen as each channel can be considered separately with a capacity that is $\frac{1}{G}^{th}$ of the total capacity and then adding the throughput for all $G$ channels together. Results in [KV05] show similar bounds for multichannel networks, stating that when there are less radios than channels available, there is a degradation in the capacity of the network, depending on the ratio of channels versus radios. The bound presented using our model for the single channel case can be obtained for each channel $g$ separately, choosing in each interval/square one node to set its value $w^g(v)$ to $\frac{1}{n}$. As the nodes are now divided into groups by the channels they have their radios set to, it more often will occur that an interval/square does not contain a node set to the channel under consideration. This will lower the bound, but so far no exact results have been obtained as were found in [KV05]. Practical application of channel allocation settings are presented in i.e.[L$^+$01],[Ste07] using a simulation program, showing similar results as obtained in the analytical papers of [GK00] and [KV05].

## 5.8 Conclusion

In this chapter we consider multi-hop multi-channel wireless networks, for which we have stated a theorem giving a necessary and sufficient condition for the existence of a solution for the multicommodity flow problem with interference constraints, given a required throughput between nodes of the network. The use of the theorem is illustrated by examples and similar bounds are obtained as presented in the well known paper by Gupta and Kumar. The applicability of this work for wireless multi-hop multi-channel networks can for example be seen when designing a network that must be able to support a certain throughput for each user. Also, the theorem provides insight in the bottleneck in the network, the location where the interference suffered is at its limit or the capacity of the link is fully used. This information gives a basis for developing an algorithm that allocates the channels to the available radios at the nodes. Starting with a network where all radios are set to the same channel, the bottleneck can be determined using the theorem, and the channel allocation can be adjusted at this location. Due to the general nature of the theorem, any type of network with heterogeneous users and protocols can be modelled.

# A flow level model for wireless multihop ad hoc network throughput

A flow level model for multihop wireless ad hoc networks is presented in this chapter. Considering different scenarios, a multihop WLAN and a path with a TCP-like flow control protocol, we investigate how capacity is allocated between the users of a network. This leads us to two different Processor Sharing models, BPS and DPS, which are discussed and compared. Simulation is used to validate the proposed models. The flow level view leads to new insights into the impact of interference on the capacity of ad hoc networks, where we show that the different queueing models provide good approximations for the troughput that a network can achieve.

## 6.1 Introduction

The multihop property and the interference this creates in wireless ad hoc networks presents new challenges to make this type of network effective. A commonly used MAC protocol to deal with these challenges is the IEEE 802.11 MAC protocol, which uses Carrier Sense Multiple Access with Collision Avoidance (CSMA/CA).

The specific characteristics of multihop ad hoc networks calls for new models to analyze the performance. In this chapter important performance measures, like the system throughput and the transfer time of flows, are investigated. The allocation of the capacity over the different network nodes plays an important role in this.

We extend the successful approach for analyzing flow transfer times in (single hop) WLANs as presented in [L+03] to multihop networks. Two different network scenarios are considered. In the first scenario flows may have different path lengths (in terms of number of hops), but follow disjunct routes. The other scenario deals with a path, which we denote as a serial network, in which multiple flows may travel through a particular node. Considering the capacity allocation in both scenarios, we propose two processor-sharing (PS) models for describing the behaviour of the network at flow level. The first model, called Batch Processor Sharing (BPS), deals with a queueing system where batches of jobs arrive. All jobs in the BPS model are served at the same time and are given an equal share of the capacity of the server. In the second model, called Discriminatory Processor Sharing (DPS), jobs arrive one by one and all jobs in the queue are served at the same time, but some jobs get a bigger share of

the capacity of the server than others. The batch sizes in the BPS model and the capacity shares in the DPS model reflect the different path lengths of the flows in the ad hoc network scenarios. The modeling results are compared with results obtained by simulation.

The rest of this chapter is constructed as follows. In the remaining part of this section, we will give a review of related literature on the subject. Next, in Section 2 the IEEE 802.11 protocol is described. Section 3 presents the two main ad hoc network scenarios under consideration, and investigates the distribution of the capacity over the users in the network. The resulting processor-sharing models for analyzing flow transfer times are shown in Section 4. These models are validated by simulation in Section 5. Finally, Section 6 summarizes and concludes the chapter.

### 6.1.1    Literature review

Many papers have been devoted to the capacity and throughput of wireless (multihop) networks. Most of them use results from simulation to describe the characteristics of ad hoc networks, whereas analytical studies are scarce. The impact of MAC layer interference on the capacity of ad hoc networks as addressed in this chapter has been studied in several settings. For instance [L+01] uses simulation to show that capacity can be very low in ad hoc networks. Scaling appears only to be possible if the distance between the source and destination remains small as the network grows. An analytical approach for determining the capacity is presented in [GK00], where it is shown how the throughput depends on the number of nodes in the network (when this number becomes large). A paper that focuses more on the multihop property of ad hoc networks is [GV02], which gives asymptotic results for a wireless network under a relay traffic pattern, whereas [JS03] considers mesh networks, slightly different from ad hoc networks. Both focus on the throughput of a chain of users processing flows in one direction over multiple hops. A bottleneck is found which determines the throughput that the network can achieve.

In the work of Litjens et al. [L+03], an integrated packet/flow level approach is used to analyze flow transfer times in a single hop WLAN scenario. Considering, the system throughput at the packet level, and taking the system dynamics at flow level into account, leads to a processor-sharing (PS) type of queueing model for the flow level. This PS model captures the equal allocation of transmission capacity among the active flows. Using known results for this PS model an approximation for the mean flow transfer time is proposed. Simulation results show that the approximation is very accurate.

Modeling bandwidth sharing in fixed communication networks by PS systems has been done by amongst others Bonald and Nunez-Queija. In the papers of Bonald [BP02],[BP03], the main notion is that modeling the network with processor-sharing can lead to balanced fairness, which means that each user in the network receives an equal share of the available network resources. This type of PS network is analyzed by considering the bottleneck node and distributing the capacity there first. All nodes servicing the same flow adjust their capacity allocation accordingly, avoiding congestion. The capacity allocated to each flow

is determined analytically. In his dissertation, Nunez-Queija discusses many different PS models for integrated services networks [NnQ00].

Batch arrival processor-sharing models have been investigated extensively in the past. It was Kleinrock who started with this approach. In his paper with Muntz and Rodemich [KMR71] a start was made in giving a complete analytical approach to determine the throughput of a PS network. A discriminatory processor-sharing model has also been used for modeling networks. Kleinrock [Kle67] already started with this in 1967 which created much interest in this type of network. In 1980, Fayolle, Mitrani and Iasnogorodski [FMI80] built on the work of Kleinrock. New results have been obtained in [CBvB05a],[CBvB05b] where the queue length distribution and sojourn times for PS models are determined. The results presented hold for general service requirements and are used to analyze WLANs with Quality of Service support [C+05].

## 6.2   IEEE 802.11 MAC Layer Protocol

As signals transmitted by a user will not only be heard by the receiver, but also by all other nodes in the vicinity of the sender, this interference limits the capacity of the network. When multiple signals reach a node at the same time, a collision occurs and the signals cannot be received correctly and packets are lost. To reduce the number of transmissions that fail and the impact this has on the throughput of the network, IEEE 802.11 uses Carrier Sense Multiple Access with Collision Avoidance (CSMA/CA). IEEE 802.11 can function in infrastructure or ad hoc mode, depending if an access point is being used. This has no implications on the MAC layer.

When a node wants to transmit, it will first listen to find out if other nodes are already transmitting: carrier sensing. If other nodes are transmitting, the node will not transmit. When the network becomes available, the node waits for a certain time (DIFS) and if the network is then still free, a timer is started to avoid collisions. This timer is paused as soon as the node senses a transmission from another node. When the network becomes free again and stays free for a DIFS, the timer continues. When the timer ends, transmission starts. This approach does not make it completely sure that collisions will not occur. Therefore, instead of sending the packets of the data immediately, a node first transmits a request-to-send (RTS). The receiver replies to this RTS by sending a clear-to-send message (CTS). The time between these transmissions (SIFS) is smaller than DIFS. After receiving the CTS, transmission of the data starts. This approach is used so that in case of a collision, only the RTS is lost, and not a much bigger packet containing data. This way the impact of a collision on the throughput of the network is reduced. When a collision occurs, a timer starts again but is set to a time taken from a window that is twice as big as before. When the transmission was succesful, the procedure repeats as long as there still are packets that need to be transmitted. The operation of CSMA/CA is shown in Figure 6.1.

Under CSMA/CA, all nodes that want to transmit compete for network resources. In the multihop wireless ad hoc network we are considering, packets from a multihop flow are present at multiple nodes. Such a flow is competing for
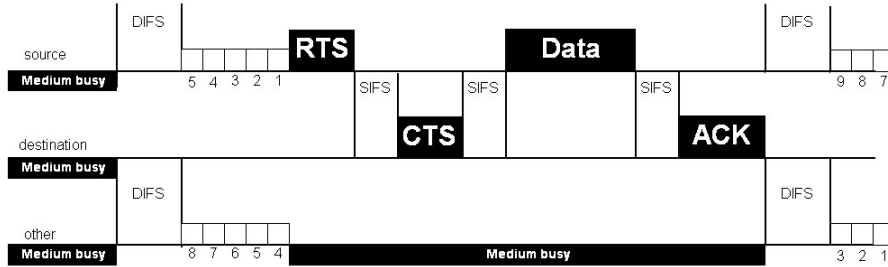
Figure 6.1: CSMA/CA with RTS-CTS

the network resources through multiple nodes at the same time. Hence even if there is only one multihop flow in the network, there is interference between the different nodes that are involved in the transmission. All flows in the network will have to share the MAC layer capacity, and the capacity allocation over these flows will influence the throughput of the network.

## 6.3   Scenarios

In this section, first, the analysis of the single hop WLAN considered in [L+03] is shortly reviewed, after which two multihop scenarios are described. For the analysis of these scenarios we use a similar approach as used for the WLAN in [L+03], extending the approach to inlude the multihop aspects involved in these scenarios.

### 6.3.1   Single hop WLAN scenario

In [L+03] a single hop WLAN is considered, in which new flow transmissions are initiated according to a Poisson process. Flow sizes are random variables with general distributions. The network operates under the IEEE 802.11 MAC protocol as described in Section 6.2. First an analysis is made on the packet level of the aggregate system throughput that can be reached in a WLAN with a fixed number of persistent flows. Using Markov-chain analysis, the probability that a node is transmitting is computed, as well as the probability that a transmission fails. From this, the aggregate system throughput is derived, including the influence of the headers and control packets. Simulation validates the result that the average system throughput is about 87% of the total capacity on the MAC layer, and slightly dependent on the number of present flows. Next, using the results from the first step, the transfer time is analyzed, taking the flow level dynamics into account. The assumption is made that the service rate per flow is found by giving each flow an equal share of the aggregate data throughput computed for the persistent number of flows in the network. This leads to a processor-sharing (PS) model with state dependent service rates which is analytically tractable. Simulation shows that the results attained through this approach are very accurate.
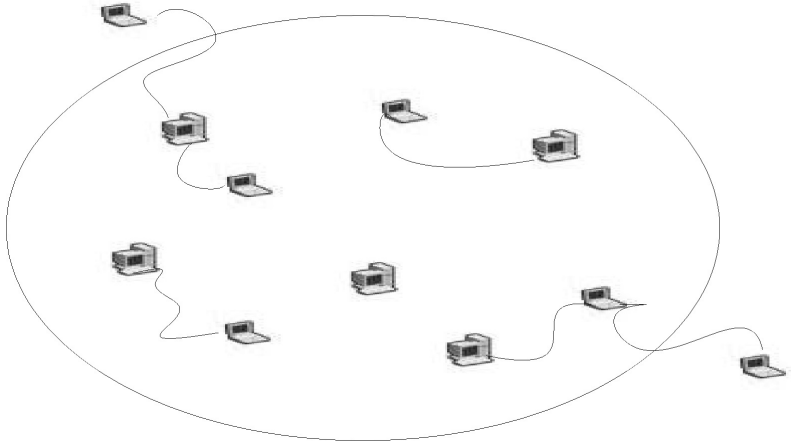
Figure 6.2: A multihop ad hoc scenario

### 6.3.2 Multihop Ad Hoc scenario

The model presented in [L+03] only considers single hop flows. This model is expanded by also allowing multihop flows. Where at first a flow was completed if the packets were sent from a node to the access point or vice versa, now there is the possibility that all packets are forwarded to another node before completing the transfer.

A transmission of any of the nodes in the cell can be heard by all other nodes in the cell. This means that no two transmissions can take place at the same time, since the data will be lost due to a collision. The situation in Figure 6.2 is a WLAN cell in ad hoc mode with connections of only one or two hops. A second hop is then used to connect to a user outside the cell. Whenever a node is relaying a flow, this node will only compete for the network resources if there are packets available that need to be sent. If at one point there are no packets available, the node will remain idle until new packets arrive that need to be forwarded.

Following the approach presented in [L+03], we first determine the aggregate system throughput of the network. The number of persistent users is taken to be the number of nodes active in sending the flows. This means that a flow over two hops, which needs two nodes, is counted as two users. However, there might be moments that the relaying nodes are not competing for the network since there are no packets available. This differs from the single hop WLAN situation described above where each node always has packets that need to be transmitted. Through simulation the aggregate system throughput of the multihop network is determined.

To take the flow level dynamics into account, the capacity allocation needs to be known. As in the WLAN model, we assume that every node receives an equal share of the aggregate system throughput. In general this is the case since all nodes behave according to the IEEE 802.11 protocol. For the flow level dynamics, we assume that flows are big enough to have packets available at
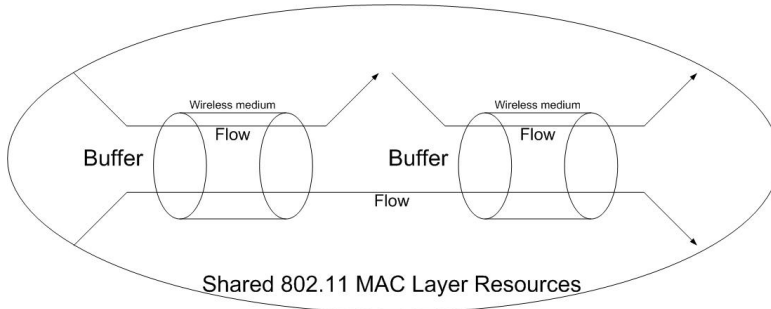
Figure 6.3: Serial model

all nodes they are going through, for most of the time. Then these nodes are continuously competing for the network and all nodes get an equal share of the capacity. A flow over two hops will get a share of the capacity for both the nodes it uses. Hence the capacity allocated to a flow over two hops will be twice the amount of the capacity allocated to a flow over one hop, but note that each packet has to be sent twice. Just as for the WLAN model, this approach leads to a processor-sharing type of model which will be discussed in Section 6.4.

### 6.3.3   Multihop serial network scenario

A different scenario is where nodes can serve more than one flow at the same time. Assuming that a flow will at most need two hops to reach its destination, such a network can be modeled as a network consisting of three nodes, with two connecting links. This network is represented by the model shown in Figure 6.3.

The flows through the wireless medium (the links), represented by the arrows through the tubes, will compete for the channel. There are three types of flows. Flows of type 0 will go over both links (as depicted by the lower arrow going through both tubes); flows of type 1 (2) will only use link 1 (2) (the upper left (right) arrow going through the left (right) tube). All flows consist of packets which first arrive in a buffer before being sent over the wireless medium. Because of interference, the links have to share the MAC layer resources.

Assume that the arrival process of the flows is according to a Poisson distribution and the flows consist of many packets, so that flows over multiple hops will again usually have packets in the queue of every node it passes. The packets of all flows join the same queue at a node. With all flows arriving simultaneously, the packets will be in the queue in a mixed order and are serviced according to a FCFS discipline. Hence all flows are serviced, packet by packet, in a more or less random order, as shown in Figure 6.4.

This way of servicing is approximated by a processor-sharing (PS) service discipline for the flows, where each link is assumed to have the same capacity. Just as in the previous scenarios, the aggregate throughput needs to be determined, which is done through simulation.

An important aspect regarding the throughput of the different types of flows in the network is how the capacity is allocated among them. This depends on the

Packet level view

| 1 | 0 | 1 | 2 | 2 | 0 |
|---|---|---|---|---|---|

Flow level view

| 0 | 0 | 0 | 0 |
|---|---|---|---|

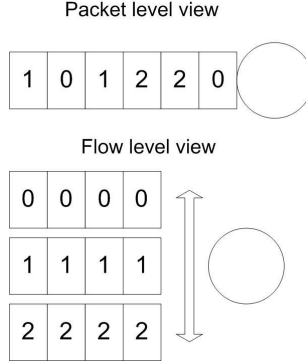| 1 | 1 | 1 | 1 |
|---|---|---|---|

| 2 | 2 | 2 | 2 |
|---|---|---|---|

Figure 6.4: Processor sharing on flow level

flow control protocol used in the network. A commonly used transport protocol in (wired) networks is TCP. This protocol tries to avoid congestion by fairly sharing the available resources among active flows. A similar protocol can be used in wireless networks. Suppose that in the reference model there are $n_0$, $n_1$ and $n_2$ flows of type 0,1 and 2 respectively. At first, the available capacity will be shared in a fair way over both nodes, each node receives half of the total capacity. Both nodes will then use this capacity to process the flows that are in their queues. This is assumed to be done so that each flow gets an equal share, which is called egalitarian processor-sharing. Now the situation can occur that flows of type 0 get a different amount of the capacity over the first link than over the second link, as will be discussed later. If the capacity at the first link is higher, the queue at the second server will build up, since it cannot serve the flow as fast as it arrives. If the capacity at the first link is lower, then the queue at node two will often not contain any packets of the flow of type 0, since these are processed faster than the rate at which they arrive. These are unwanted situations, and so the flow control protocol will notify the sources to transmit at different rates for the flows that are causing the congestion. We assume a TCP-like flow control protocol to be active, which alters the use of capacity in case of queues building up or being empty most of the time. The flow control protocol resembles TCP, but assumes a perfect version with instantaneous rate adaptation, so it will be e.g. independent of the round trip times (RTT) of the flows.

In our model there can be a loss of packets due to a build up in the queue of node 2. This happens if the flows of type 0 get more capacity in node 1 than in node 2 for a long time. This is for example the case when there are flows of type 2, while there are no flows of type 1 in the system and both links get half of the total capacity. If the flows of type 0 keep being processed at this rate, packets will be lost. However, the flow control protocol will make sure that the rate at which packets are transmitted over link 1 is lowered. The rate to lower it to is the rate at which the flow is processed at node 2. The capacity used by node 1 hence drops due to the lowering of the transmission rate. This capacity can then be used by node 2.

**Theorem 6.1.** *The capacity a flow receives at a link is equal for any type of flow at any link, namely $\frac{1}{2n_0+n_i}$ when there are only flows of type 0 and one other type (i).*

*Proof.* Consider the network where only flows of type 0 and 1 are present. Let there be $n_0$ ($n_1$) flows of type 0 (1). If both links get half of the total capacity, which we set to be one, then the capacity allocated to a flow at node 1 will be equal to $\frac{1}{2}\frac{1}{n_0+n_1}$. At node 2, the flows of type 0 will receive a capacity of $\frac{1}{2}\frac{1}{n_0}$. If this situation persists, the queue at node 2 will often be empty, since the flows of type 0 are processed faster than the rate at which they arrive. The flow control protocol therefore lowers the rate at node 2 and sets the capacity of a flow to $\frac{1}{2}\frac{1}{n_0+n_1}$. Node 1 can now use the residual capacity and a flow will get a capacity of $\frac{1}{n_0+n_1}(1-\frac{1}{2}\frac{n_0}{n_0+n_1})$. This rate however is higher than the rate at node 2 and so the buffer will fill. The rate at node 2 is adjusted by the flow control protocol to $\frac{1}{n_0+n_1}(1-\frac{1}{2}\frac{n_0}{n_0+n_1})$, which leaves a capacity for a flow at node 1 of $\frac{1}{n_0+n_1}(1-\frac{n_0}{n_0+n_1}(1-\frac{1}{2}\frac{n_0}{n_0+n_1}))$. This process continues (where each step is instantaneous by assumption) and it can easily be seen that this converges to a capacity allocation of $\frac{1}{2n_0+n_1}$ for each flow on any link. In total, a flow of type 0 will receive $\frac{2}{2n_0+n_1}$ of the capacity, whereas a flow of type 1 will get $\frac{1}{2n_0+n_1}$. The proof for the situation with only type 0 and type 2 flows follows in the same manner. $\square$

If however there are flows of both type 1 and type 2 in the system, this situation will not occur. Supposing that there are more flows of type 2 than of type 1, a flow of type 0 will get a rate of $\frac{1}{n_0+n_2}$ at node 2, but receives a higher rate of $\frac{1}{n_0+n_1}$ at node 1. The protocol will hence lower the rate for all type 0 flows to $\frac{1}{n_0+n_2}$ at node 1. The capacity that now becomes available is not given to node 2, but the flows of type 1 will claim this extra capacity since the node has the right to use half of the total capacity. This is a different and interesting situation for further research, but we will not consider it any further in the analysis in this chapter.

We see that in the first situation the flow control protocol leads to the situation in which all flows get the same share of capacity per link. Hence of the total capacity, a flow over two links will get twice as much capacity as a flow over one link. This is equivalent to what we found for the two-hop WLAN model. Hence we need to analyze the same situation, which as noted before leads to a processor-sharing model, which will be discussed in the next section.

The analogy between the scenario as shown in Figure 6.2 and the scenario of Figure 6.3 can be seen as follows. A flow over two hops coincides with a flow of type 0 and flows over one hop coincide with a flow of type 1 or 2 (but not both), which is arbitrary. In the serial network, just as in the two-hop WLAN scenario when a node relays, a flow taking the second hop will only compete for the network resources if it has packets available to be sent. If at one point no such packets are available in the buffer, no capacity is allocated to this flow at that node.

## 6.4 Flow level models

This section presents models that approximate the flow level dynamics of the multihop scenarios presented in the previous section. For the multihop ad hoc scenario, the equal share given to each of the nodes determines the capacity allocation. For the serial scenario, the transmission control protocol determines the capacity allocation at the MAC layer. This section presents two analytical models that capture the flow level dynamics of both scenarios. The two models take the capacity allocation into account by either varying the amount of jobs in an egalitarian processor-sharing queue or the priority and size of jobs in a discriminatory processor-sharing queue.

### 6.4.1 Batch Arrival Processor Sharing model

We can consider the network as a server with one queue, where all flows enter the queue, independent of the link(s) they have to be transmitted over. Since all flows are processed at the same time, we can consider the flows to be processed according to a processor-sharing discipline. As a flow over two hops gets the double amount of capacity, we can consider a flow over two hops as asking for capacity twice. Hence we can see the arrival of a flow over two hops as the arrival of two flows at the same time. A flow over two hops will be at both servers at the same time, so it will ask for capacity as if it were two different flows, which is captured in this abstract view. We thus arrive at a batch arrival processor-sharing model (BPS) with egalitarian processor-sharing, since the capacity for a flow is equal at every hop. A flow over a single hop is then equivalent to a single arrival, whereas an arriving flow over two hops is equivalent to two jobs arriving as a batch. It is important to note here is that we do need to consider that all jobs in a batch should not only have the same arrival time but also the same departure time, i.e. jobs in a single batch have the same service demand, since they represent only one flow.

Consider the $M^X/G/1$ PS queue where $\lambda$ is the batch arrival rate, $a$ is the average batch size, $b$ is the average number of jobs that arrive in addition to the tagged job and $\overline{F}(x)$ is the complementary distribution function of the job size. The conditional response time of a job with service requirement $x$, $T(x)$, has to satisfy the system of differential equations ([KMR71]):

$$T'(x) = \lambda a \int_0^\infty T'(y)\overline{F}(x+y)dy + \lambda a \int_0^x T'(y)\overline{F}(x-y)dy + b\overline{F}(x) + 1 \quad (6.1)$$

The load in the system is given by:

$$\rho = \lambda a E[X]. \quad (6.2)$$

When flows have an exponential service requirement, solutions can be found ([AAB03]). For the $M^X/M/1$ PS queue, this leads to:

$$T(x) = \frac{x}{1-\rho} + \frac{b(2-\rho)E[X]}{2(1-\rho)^2}(1 - e^{\frac{-(1-\rho)}{E[X]}x}) \quad (6.3)$$

and bounds are given by:

$$\frac{x}{1-\rho} \leq T(x) \leq \min(\frac{b+1}{1-\rho}x, \frac{x}{1-\rho} + \frac{b(2-\rho)E[X]}{2(1-\rho)^2}) \tag{6.4}$$

where the bounds coincide when $x^* = \frac{(2-\rho)E[X]}{2(1-\rho)}$.

In these models, the departure moments of jobs inside a batch will not be the same. Only when the service times are deterministic, is this model applicable. Therefore, we also propose a more appropriate model.

### 6.4.2 Discriminatory Processor Sharing model

A flow over two hops receives more capacity than a flow over one hop, hence we can instead consider the processor-sharing not to be egalitarian. The jobs are then processed at the same time, but not all jobs get an equal share. As a flow over two hops takes twice the amount of capacity, it can be seen as being serviced twice as fast as a single hop flow. A flow over two hops however has an expected service requirement that is twice the expected service requirement of a single hop flow. We thus arrive at a discriminatory processor-sharing model (DPS). In this type of model, all jobs get processed at the same time, but not all jobs get the same amount of service. Customers are given a certain weight which shows how much more service they receive in comparison to other users. In our case, a job that represents a two-hop flow will get a weight twice as high as a job representing a single-hop flow.

Consider the $M/G/1$ DPS queue where $\lambda_j$ denotes the arrival intensity of class $j$ jobs, $g_j$ denotes the 'weight' of class $j$ customers and $F_j(x)$ the distribution of the required service with mean $1/\mu_j$ and a total of $M$ classes. The conditional response time of a job of class $k$, given its size $t$, $W_k(t)$, satisfies the system of differential equations ([Kle67]):

$$W_k'(t) = 1 + \sum_{j=1}^{M} \int_0^\infty \frac{\lambda_j g_j W_j'(u)}{g_k}(1 - F_j(u + \frac{g_j t}{g_k}))du + \tag{6.5}$$

$$\int_0^t W_k'(u) \sum_{j=1}^{M} \frac{\lambda_j g_j (1 - F_j(g_j(t-u)/g_k))}{g_k} du, \ k = 1...M.$$

For the $M/M/1$ DPS queue, we know that (for the derivation see [FMI80]):

$$W_k(t) = \frac{t}{1-\rho} + \sum_{j=1}^{m} \frac{g_k c_j \alpha_j + d_j}{\alpha_j^2}(1 - e^{-\alpha_j t/g_k}), \tag{6.6}$$

where the $\alpha_j'$s are the distinct roots of $1 - \Psi^*(s) = 1 - \sum_{j=1}^{M} \frac{\lambda_j g_j}{\mu_j g_j + s} = 0$ and $c_j$ and $d_j$ are given by:

$$c_j = (s + \alpha_j)a^*(s)|_{s=-\alpha_j} = \frac{\prod_{k=1}^{m}(g_k \mu_k - \alpha_j)}{-\alpha_j \prod_{k \neq j}^{m}(\alpha_k - \alpha_j)} \tag{6.7}$$

$$
\begin{aligned}
d_j &= (s^2 + \alpha_j^2)\theta(s)|_{s^2=-\alpha_j^2} \\
&= \frac{[\sum_{k=1}^{M} \lambda_k g_k^2/(\mu_k^2 g_k^2 - \alpha_j^2)][\prod_{k=1}^{m}(g_k^2 \mu_k^2 - \alpha_j^2)]}{\prod_{k \neq j}^{m}(\alpha_k^2 - \alpha_j^2)}.
\end{aligned} \tag{6.8}
$$

Since each flow is represented by only one job in the system, the departure of a job is equivalent to the departure of a flow from the network. Therefore, this approach gives a better approximation of the situation that we want to model. When considering deterministic service requirements, we see that both models give equivalent results.

## 6.5 Numerical results

To verify that the proposed model is an accurate approximation of the network under consideration, a simulation model has been constructed to obtain data on the sojourn time of a flow in the network. The simulation model mimics the transmissions as they occur in the scenario depicted in Figure 6.2. The simulation model uses the following standard settings for the parameters:

| parameter | value | parameter | value | parameter | value |
|-----------|-------|-----------|-------|-----------|-------|
| PHY | 192 bit | Payload size | 12 kbit | SIFS | 10 $\mu s$ |
| MAC | 272 bit | $r_{net}$ | 1 Mbit/s | DIFS | SIFS + $2\tau$ |
| RTS | PHY+160 bit | $n_{max}$ | 100 | cw$_{min / max}$ | 31/1023 |
| CTS | PHY+112 bit | $\delta$ | 1 $\mu s$ | r$^*$ | 5 |
| ACK | PHY+112 bit | $\tau$ | 20 $\mu s$ | r$_{max}$ | 6 |

Here $r_{net}$ is the rate at which the network can transmit data, $n_{max}$ is the maximum number of users, PHY, MAC, RTS, CTS and ACK give the sizes of the headers and interframe spaces, $\delta$ is the propagation delay, $\tau$ is the slot duration, cw are the values for the contention windows, $r^*$ is the maximum number of times the contention window may be doubled and $r_{max}$ is the maximum number of retransmissions for one packet. The payload size is set at 12 kbit, the maximum amount of data that can be sent within a packet. The probability that a flow is over two hops is given as input to the simulator. All flows are either single or double hop flows. Files arrive at the system according to a Poisson process, where users try to transmit the packets according to the IEEE 802.11 protocol discussed earlier. When the first packet of a double hop flow has been sent over the first hop, a free user is assigned as the relaying node and this user will also start transmitting the packets that it receives from the first user. When the second user has no packets in its queue to relay, it will go into waiting, meaning that he will not compete for the channel. The DCF function of IEEE 802.11 is incorporated in the simulation, where collisions are considered to be fatal, meaning that all packets involved in the collision are lost and retransmitted after backing of.

The results of the simulation and the $M^X/M/1$ BPS and $M/M/1$ DPS model are compared, we hence are considering Poisson arrivals and exponentially distributed file sizes. The model considers the network in a situation that the full capacity of the network can be used for the files, which is not the case for
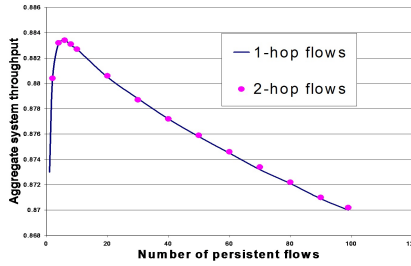
Figure 6.5: Aggregate system throughput

the simulation program. Here headers are added to all the packets and the number of active users influence the aggregate system throughput as described in [L+03]. Following their approach, the aggregate throughput for persistent flows is determined. For the determination of this aggregate throughput the amount of double hop flows can be of influence. The interference that the flow causes for itself deteriorates the throughput of the network. However, the second node in a double hop flow may not always have packets to transmit, at which point it will not cause interference. Simulation is used to compare the results of single and double hop flows. In Figure 6.5, the aggregate system throughput is computed for single hop flows, and compared with the aggregate system throughput of double hop flows, where the number of persistent flows is taken to be twice the amount of double hop flows. Figure 6.5 clearly shows that the aggregate throughput is hardly influenced by double or single hop flows. This shows that we can use the results for single hops, but counting the double hop flows as if there are two users in the system.

Under the RTS/CTS mode, the aggregate throughput is roughly constant as was also found for the WLAN situation in [L+03]. As can be read from the figure only about 88% of the capacity can really be used. This is taken into account in the calculation of the average transfer time in the models. First we compare the average transfer time of a job in the system with the results from the BPS model. A drawback of this model is that it is not possible to make a distinction between the single and double hop flows in this model. Results are shown in Figure 6.6 for the situation where 30% and 70% of all the flows are double hop flows and the file sizes are exponentially distributed with a mean of 150 kbytes.

Figure 6.6 shows that for a lower amount of double hop flows, the approximation is better. When 70% of all the flows are double hop flows, the difference becomes bigger. Next, the comparison is made between simulation and the DPS model, where we can distinguish between the different type of flows, which is shown in Figure 6.7.

The approximations are accurate, independent of the amount of double hop flows in the system. It can be seen that for a higher load the approximation is slightly worse. Interesting is to see that the model overestimates the transfer time for single hop flows, whereas is underestimates the transfer time of double hop flows. This is due to the fact that the assumption that a double hop flow
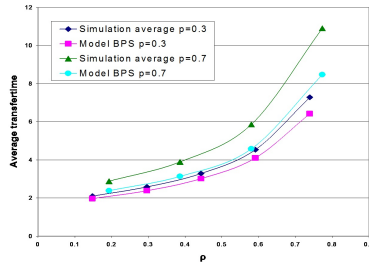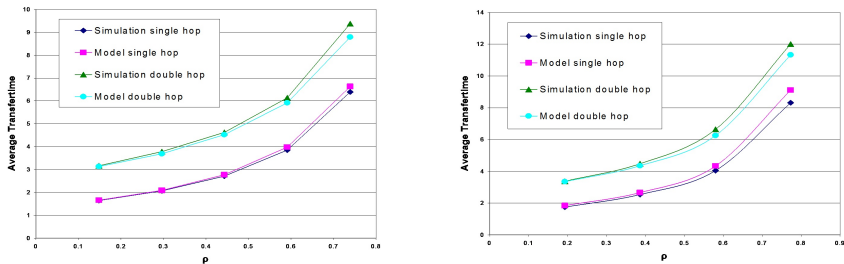
Figure 6.6: BPS vs simulation



Figure 6.7: DPS vs simulation for $p = 0.3$ and $p = 0.7$

receives twice the capacity of a single hop flow is not exact, since packets are not always available at all the different nodes on a multihop path. Therefore the capacity allocated to double hop flows is slightly less than the double of single hop flows, resulting in the differences seen in the figures.

## 6.6 Conclusion

Using a flow level approach for modeling an ad hoc network, many difficulties might be avoided that would occur when using a packet level view. Results from the past have shown that simulation is often the only possible approach to get results and for analytical approaches to be possible many assumptions have to be made. However, the flow level view has a promising future, even for analytical approaches.

This chapter discusses the network from a flow level point of view, with two types of models for approximating the throughput of such a network, namely the BPS and DPS models. The $M^X/G/1$ PS and $M/G/1$ DPS queues can be used, since these models take the allocation of the capacity over the different flows in the network into account. Much work on these types of models has already been done, and some analytical results have been presented. When considering the network as a BPS queue, the problem arises that all flows inside a batch should have the same service requirement, which is not the case for the $M^X/G/1$ PS queue, making it impossible to distinguish between different classes of flows. Therefore, it is more accurate to use the $M/G/1$ DPS queue for modeling the ad hoc network.

The flow level view as presented in the chapter opens new opportunities for modeling ad hoc networks, taking into account interference at the MAC layer, especially self-interference within a flow. Through simulation the model has been validated, showing that the results obtained using the $M/M/1$ DPS queue gives a good approximation of the transfer time in an ad hoc network using IEEE 802.11 in RTS/CTS mode, independent of whether many or few flows are double hop flows.

# Bottlenecks and stability in networks with contending nodes

Motivated by interference, this chapter considers a new class of queueing network models, where nodes have to contend with each other to serve their customers. In each time slot, a node with a non-empty queue either serves a customer or is blocked by a node in its vicinity. The focus of our study is on analyzing the throughput and identifying bottleneck nodes in such networks. Our modeling and analysis approach consists of two steps. First, considering the slotted model on a longer timescale, the behaviour is described by a continuous time Markov chain with state dependent service rates. In the second step, the state dependent service rates are replaced by their long run averages resulting in an approximate product form network. This enables us to determine the bottleneck nodes and the stability condition of the system. Numerical results show that our approximation approach provides very accurate results with respect to the maximum throughput a network can support. It also reveals a surprising effect regarding the location of bottlenecks in the network when the offered load is increased.

## 7.1  Introduction

Inspired by wireless ad hoc networks where interference prohibits neighbouring nodes to simultaneously transmit packets, this chapter considers a class of open queueing network models in which servers contend for service slots. In each time slot nodes that have packets available for transmission try to obtain the channel to transmit their packets. As nodes within each others interference range cannot transmit at the same time, an allocation mechanism, i.e. a medium access control protocol, is used to decide which nodes get the opportunity to transmit, i.e. to serve a packet. Once a server in a node is allowed to transmit a packet, it blocks the servers in a specified set of other nodes corresponding to an interference neighbourhood. Upon service completion, a packet either moves to a next node for further service, or leaves the network. The network is called stable when for each node the average service rate exceeds the average arrival rate of packets. When multiple or large flows pass through a node, the service rate of the node may not suffice, making this node a bottleneck. This chapter investigates the stability range, the arrival rates of flows at which nodes become bottlenecks, and the throughput of the network.

The behaviour of the system under consideration can be described by a state dependent discrete time Markov chain as we assume that in each slot

Figure 7.1: Approximation steps

the contention between nodes takes place independent of previous outcomes. Inspired by results obtained for loss-networks, we make a two step approximation to analyze this network, see Figure 7.1. As a first step, we consider the long term average behaviour, which neglects the effect of the slotted time and leads to a continuous time Markov chain. However, with the transition rates of this chain still being state dependent, analysis remains cumbersome and a further approximation is needed. Using a long term average service rate, we introduce a product form network approximation which enables us to find the bottlenecks in networks of arbitrary size and topology and determine the maximal throughput. Interestingly, it turns out that when the load of the network is increased, a bottleneck node can become stable again as a different node becomes the bottleneck. This surprising behaviour is predicted correctly by our product form model.

The remainder of the chapter is organized as follows. First, Section 2 gives a literature overview, after which Section 3 introduces the discrete model and contention process. Section 4 describes the first approximation step resulting in the continuous time model with state dependent service rates, followed by the second approximation step in Section 5. Section 6 gives the results for the stability analysis and Section 7 presents results from simulation to illustrate the accuracy of the model presented in this chapter. Finally, Section 8 concludes the chapter.

## 7.2    Literature and contribution

The stability of networks, as considered in this chapter, has received considerable interest in literature. Inspired by wireless networks, [RE88] analyzes a discrete time slotted ALOHA system. Bounds on the stability region are found using the concept of dominance. A different approach is presented in [JvdMvdW07] where the rate stability and output rates are calculated for shared resource networks. Stability conditions for separate nodes are derived for general allocation functions under mild assumptions. The model discussed in this chapter however does not fall under the set of allocation functions, as the overall capacity of the network is not constant due to interference. For a network of parallel servers with coupled service rates, necessary and sufficient conditions for stability are derived in [BJL08]. Stability and performance of networks where the service rate depends on the network state is also analyzed in [VLK01], where transmissions over links with a fixed capacity are considered. Opposed to the work presented in these papers, the rate allocated to a server does not depend on the number of packets present in the queue, but on the number of nodes competing.

Similar assumptions regarding the contention between nodes are made in [D$^+$08], where alive nodes block other nodes as discussed in this chapter. The throughput in a multihop tandem network is considered both under saturation, where each node generates its own traffic and under a single flow over all nodes. The authors conjecture that a random access scheme severely degrades the throughput of the network.

Analytic results for a multihop network with two contending queues are presented in [RO03]. Using the theory of Riemann-Hilbert boundary value problems, the generating function of the stationary distribution is obtained. In [LR06] some performance measures of this system are analyzed, focussing on the computational issues that occur. Even for such a small network as considered in these papers, a complex analysis is needed to obtain analytical results. The approach we present is applicable for general size networks, however we do not obtain results on the stationary distribution, but on stability and throughput.

The optimal throughput a network can support, often referred to as the capacity of the network, is discussed in [GK00], which however does not focus on multi-hop networks. This aspect is addressed in [GV02], where for a single multi-hop flow a new capacity limit is derived. These results are limiting results for large networks. More detailed models are discussed in [NL02] for a tandem and lattice network with saturated nodes. They calculate the optimal offered load, preventing packet loss in a network with hidden nodes. This work is extended for multiple crossing flows in [F$^+$06]. Instead of focussing on the specific parameters of the MAC protocol, as presented in these papers, we take a higher level view, providing valuable insights on bottleneck locations for general networks.

Next to limiting the capacity of a network, contention between nodes has an impact on the fairness of protocols, as in the equality in rate allocated to nodes or the throughput of flows. In [DDT08] the authors describe the border effects in a CSMA/CA network and its impact on fairness. The stability and throughput for a weighted fair queueing model with saturated nodes is discussed in [EKEA07] showing that the throughput, while taking into account the topology, routing and random access in the MAC layer, does not depend on the load in the intermediate nodes as long as the network is stable.

Different aspects of importance for the stability and throughput of networks have also received much attention. Focussing on the impact of routing, [KEAA08] investigates the stability and throughput of static wireless networks with slotted time. The authors show that routing has a large impact on the stability properties and that as long as the intermediate queues in a network are stable, the throughput does not depend on the traffic generated at these intermediate nodes. In [HvM04] the focus is on the calculation of the interference to noise ratio and show the influence of the network size and the data rate on this ratio and link this to the throughput of the network.

The contribution of this chapter is that we provide a comprehensible model that accurately predicts the bottlenecks and maximal throughput of a network, which also is applicable for networks with unstable nodes. The results provide insight in the impact that contention between nodes has on the performance of the network, without the need of a complex analysis.
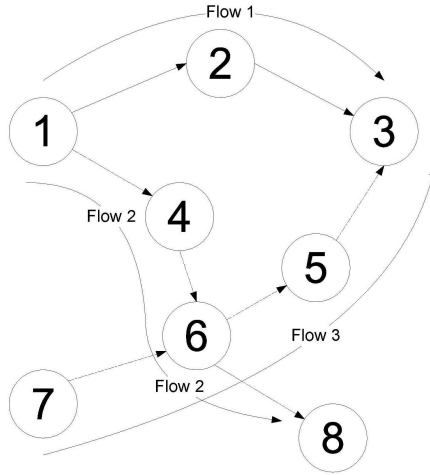
Figure 7.2: Network with three traffic flows

## 7.3   Discrete time model

### 7.3.1   General model

Consider a network consisting of $n$ queues with infinite buffers. Due to contention between nodes not all nodes can transmit their packets at the same time. We define the *contention set* $I(i)$, $i = 1, \ldots, n$, of a node $i$ as the set of nodes blocked from transmission when node $i$ is transmitting. A typical example is the set of nodes within a certain interference range. However, for the model there need not be a relation between the network structure and the contention set. The way contention between nodes takes place will be elaborated upon below. A set of $J$ traffic flows $f(t_j)$, $j = 1, \ldots, J$, travel over multihop paths, denoted by the ordered sets $t_j$, from node $t_j(1)$ to $t_j(m_j)$, where we assume that no loops are made within a path, i.e. paths are simple and packets automatically follow their path. Traffic consists of equally sized packets that are transmitted one packet per time slot. An example of such a network is depicted in Figure 7.2. In this example a network of 8 nodes is depicted. In the figure there are three flows: $f(t_1)$ from node 1 through node 2 to node 3 (i.e. we have that $t_1 = \{1, 2, 3\}$), $f(t_2)$ from node 1 through nodes 4 and 6 to node 8 and $f(t_3)$ from node 7 through nodes 6 and 5 to node 3.

Packets arrive at the origin nodes $t_j(1)$ (i.e. nodes 1 and 7 in Figure 7.2) according to a Poisson process with rate $\lambda_j$ for flow $j$ and are served first come first served. A node is called *stable* when its average service rate exceeds the average arrival rate of packets at the node, and a network is called stable when all its nodes are. A node that is unstable is called a *bottleneck* node. The average number of packets of a flow that reach the destination node per time unit is the throughput of this flow, which is limited by the service rate of the bottleneck nodes of the network. The main interest in this chapter is the throughput of the

flows and the identification of bottleneck nodes.

A node is called *alive* when it has packets to transmit and thus participates in the contention. In each time slot, all alive nodes contend to be allowed to transmit a packet. The probability of a node being allowed to transmit depends on the set of nodes contending. We focus on this aspect in the following section. In each time slot a node is either not contending, blocked or allowed to transmit. In each time slot this process is repeated, where we assume the selection of nodes being allowed to transmit to be independent between time slots.

Let $p_i$ denote the probability that node $i$ is alive and let $\pi$ be a liveliness vector, such that $\pi_i = 1$ if node $i$ is alive and $\pi_i = 0$ otherwise. The set of all $2^n$ possible liveliness vectors is denoted by $\Pi$. The probability that a liveliness vector $\pi$ occurs is denoted by $q_\pi$. The probability that node $i$ transmits under liveliness vector $\pi$ is denoted by $r_{i,\pi}$. The network can be represented by a discrete time Markov chain with the queue lengths at each node as the state of the system. Actually, as packets are forwarded to a next node depending on the flow they belong to, also the type of the packets, stating the flow they belong to, in the queue needs to be included in the state description. However, in our steady state description these types will not play a role and are therefore omitted from the state description. The transition probabilities depend on the state of the system via the liveliness vector only, i.e. the number or type of packets in a queue does not affect the probability of a node transmitting in a slot, unless it is empty.

### 7.3.2 Contention

Multiple nodes can only be transmitting simultaneously in the same time slot when they are outside of each others contention set. If multiple nodes within each others contention set are alive, the contention protocol decides which nodes may transmit. The probability that a node is allowed to transmit a packet in the following slot can be determined when the contention sets, the protocol in use and the competing nodes are known. We assume an ideal contention protocol, where no collisions will occur and hence no packets will be lost.

As we are not interested in the details of the contention protocol but only the corresponding probabilities for nodes to transmit, we will use a simple protocol giving each node an initial equal probability of winning a contention. For other protocols, transmission probabilities can also be calculated. Using

$$|\pi| = \sum_{i=1}^{n} \pi_i, \tag{7.1}$$

i.e. $|\pi|$ equals the number of alive nodes and (taking $e_m$ as the unit vector of length $n$, with all zeros except a 1 on location $m$),

$$\tilde{\pi}(k) = \pi - \sum_{m \in I(k):\pi_m=1} e_m \tag{7.2}$$

as the liveliness vector remaining after a node $k$ blocks all nodes in its contention

set (as it won the contention), the probability $r_{i,\pi}$ that a node $i$ may transmit a packet under liveliness vector $\pi$ can be calculated using the following recursion:

$$r_{i,\pi} = \begin{cases} 0 & \text{for } \pi_i = 0 \\ (1 + \sum_{k \neq i : \pi_k = 1} r_{i,\tilde{\pi}(k)})/|\pi| & \text{for } \pi_i = 1 \end{cases} \tag{7.3}$$

with $r_{i,\mathbf{0}} = 0$, where $\mathbf{0}$ denotes a liveliness vector with no alive nodes. This can be seen as follows: With equal probability of $\frac{1}{|\pi|}$ any non-empty node (so node $i$ itself or any other alive node $k$) wins the direct contention. Assuming node $k$ wins the contention, it blocks all nodes in its contention region, reducing the liveliness state to $\tilde{\pi}(k)$, after which all remaining nodes contend again. Any node that did not win the contention, but was not blocked hence can compete again and might win the new contention, with probability $\frac{1}{|\tilde{\pi}(k)|}$. This process continues until all non-empty nodes either are allowed to transmit or are blocked.

As an example, consider the network as depicted in Figure 7.2 with contention sets chosen such that nearby nodes contend: $I(1) = \{2, 4\}$, $I(2) = \{1\}$, $I(4) = \{1, 5, 6\}$, $I(5) = \{4, 6\}$, $I(6) = \{4, 5, 7\}$, $I(7) = \{6\}$. Assume that all 6 nodes have packets to transmit (as nodes 3 and 8 do not transmit packets they are never alive), so that $\pi = (1, 1, 0, 1, 1, 1, 1, 0)$. The probability $r_{4,\pi}$ that node 4 will be allowed to transmit by directly winning the contention is $\frac{1}{|\pi|} = \frac{1}{6}$. If for example node 1 wins the contention, node 4 is blocked, as it is in its contention set. As $\tilde{\pi}(1) = (0, 0, 0, 0, 1, 1, 1, 0)$, we get $r_{4,\tilde{\pi}(1)} = 0$ as $\tilde{\pi}_4(1) = 0$. The same holds if node 5 or 6 wins the contention, as $r_{4,\tilde{\pi}(5)} = 0$ and $r_{4,\tilde{\pi}(6)} = 0$. If node 2 or 7 wins the contention, node 4 still could be allowed to transmit. The probability that node 4 wins contention after node 2 has won the contention is given by $r_{4,\tilde{\pi}(2)}$, where $\tilde{\pi}(2) = (0, 0, 0, 1, 1, 1, 1, 0)$. This probability can be calculated by calculating $r_{4,\pi}$, but with $\pi = (0, 0, 0, 1, 1, 1, 1, 0)$, showing the recursion. As after each step, but with the new value of $\pi$, the number of zeroes in the liveliness vector increases, the recursion will stop when $\pi = (0, 0, 0, 0, 0, 0, 0, 0)$. For this example, the probability the nodes are allowed to transmit are given by $[\frac{19}{48}, \frac{29}{48}, 0, \frac{14}{48}, \frac{20}{48}, \frac{19}{72}, \frac{53}{72}, 0]$. We further analyze this network in Section 7.7.2.

Note that the overall probability of being allowed to transmit is not equal for all nodes. A similar analysis to obtain $r_{i,\pi}$ can be done for any network, with any contention sets and protocol. More extensive calculations will be needed for larger networks with different topologies, but the principle will not change. In the remainder of this chapter, we will assume that the contention regions and protocol are known, such that all conditional rates $r_{i,\pi}$ of the nodes can be calculated.

## 7.4    Approximation step 1: Continuous time

We are interested in the long term average behaviour of the network, especially the throughput and stability issues. Considering the system on a higher level and a larger time scale, the discrete character due to the time slots fades and the model can be seen as a continuous time Markov process. The state of the system consists of the number and type of packets at each queue, but as the state of the

system only influences the transition rates through the liveliness of the network, we do not focus on the queue lengths. The flow a packet being served belongs to determines the direction in which it will be forwarded. We incorporate this into the model as described below. In the following we will denote parameters used for the continuous time approximation by adding a hat to the equivalent parameter in the original discrete time model.

When a queue has packets available and the liveliness is given by $\pi$, the probability of a packet being sent is given by $r_{i,\pi}$. On average, the number of packets sent per slot under state $\pi$ hence is $r_{i,\pi}$. For the continuous time Markov chain, we approximate the service rate under state $\pi$ of the node using the exponential distribution with rate $\hat{r}_{i,\pi} = r_{i,\pi}$. The probability of a node being alive or not depends on the arrival rate of packets and the service rate at the node. We first focus on the arrival rate of packets.

Whenever the nodes on a multihop path $t_j$ preceding a node $t_j(i)$ are stable, the arrival rate from this flow will be $\lambda_j$, the external arrival rate of the flow. The total arrival rate of traffic $a_i$ at node $i$ is given by

$$a_i = \sum_{j:i\in t_j} \lambda_j(i) \tag{7.4}$$

where $\lambda_j(i)$ is the arrival rate at node $i$ for flow $f(t_j)$. When the network is stable, this simplifies to $a_i = \sum_{j:i\in t_j} \lambda_j$. When there are unstable nodes in the network, the arrival rate of packets at each queue can be determined as follows. Due to the multihop feed forward structure of the network we have that the arrival rate $\lambda_j(i)$ is determined by its preceding nodes. If one or more of the preceding nodes are unstable, the average arrival rate for the nodes after the bottleneck on this path will depend on the service rate of the unstable nodes. The probability $p_{t_j(i-1)t_j(i)}$ that a served packet at node $t_j(i-1)$ continues to node $t_j(i)$, the packet is of flow $f(t_j)$, is given by

$$p_{t_j(i-1)t_j(i)} = \frac{\lambda_j(t_j(i-1))}{a_{t_j(i-1)}}. \tag{7.5}$$

The arrival rate $\lambda_j(t_j(i))$ from flow $f(t_j)$ at node $t_j(i)$ is given by

$$\lambda_j(t_j(i)) = \min(\lambda_j(t_j(i-1)), p_{t_j(i-1)t_j(i)}\hat{r}_{t_j(i-1)}), \tag{7.6}$$

where $\lambda_j(t_j(1)) = \lambda_j$, the external arrival rate of packets at the first node in path $t_j$. This can be seen as follows: either the preceding node can serve all its incoming traffic, or its service rate is too low. In the latter case, the fraction of the service rate of node $t_j(i-1)$ that is used for flow $f(t_j)$, equal to $p_{t_j(i-1)t_j(i)}$, determines the arrival rate at the next node for this flow. Here $\hat{r}_{t_j(i-1)}$ denotes the average state independent service rate of node $t_j(i-1)$, which will be determined in the next section. Assuming this rate is known, equations (7.5) and (7.6) give a system of equations that can easily be solved, giving the arrival rate per flow at each node. We use these arrival rates in the analysis of the liveliness of the system, which influences the service rate of the nodes.

## 7.5   Approximation step 2: Product form network

The Markov chain with state dependent service rates is not amenable for analysis. For a network with only two queues in tandem, this equals the model presented in [RO03] under deterministic service times. Even for such a small network, a complex analysis is needed to obtain analytical results. Therefore, for an arbitrary network, we approximate the continuous time approximation by obtaining an appropriate *state independent* service rate for each node to analyze the behaviour of the network.

   The state independent service rate $\hat{r}_i$ is obtained by considering the long term average percentage of time the system is in a state with liveliness vector $\pi$. The probability of node $i$ being alive is given by

$$\hat{p}_i = \min(\frac{a_i}{\hat{r}_i}, 1). \tag{7.7}$$

For the final approximation step, let $\hat{q}_\pi$ denote the steady state probability that the liveliness vector is $\pi$ (to be calculated later) and assume the state independent average service rate of a node $i$ in the network to be given by,

$$\hat{r}_i = \sum_{\pi \in \Pi} \frac{\hat{r}_{i,\pi} \hat{q}_\pi}{\hat{p}_i} \tag{7.8}$$

We obtain equation (7.8) by considering a large time scale and weighing the service rate over the possible liveliness of the system, i.e. by unconditioning on the liveliness, but conditioning on the node being alive. The state independent service rate $\hat{r}_i$ can be seen as the average rate at which a node services packets, given that it is alive.

**Theorem 7.1.** *The steady state probability $\hat{q}_\pi$ that the system is in a state with liveliness vector $\pi$ is given by*

$$\hat{q}_\pi = \prod_{i=1}^{n} (1 - \hat{p}_i)^{(1-\pi_i)} \hat{p}_i^{\pi_i}. \tag{7.9}$$

*Proof.* Summarizing the above, we have the following assumptions for the state independent continuous time approximation:

1. The external arrival process of traffic at queues is a Poisson process.

2. There is infinite waiting space at all the queues.

3. The service time at the queues has an exponential distribution and is independent of the state of the system and arrival process.

4. After completion of service at queue $i$ a packet instantaneously moves to the next queue $k$ with probability $p_{ik}$, $k = 1, \ldots, n$, for additional service or with probability $p_{i0}$ the packet completes service and leaves the system, where we have that $\sum_{k=0}^{n} p_{ik} = 1$. The routing probabilities are independent of the history of the system.

A network for which the assumptions 1) to 4) hold is a product form network (c.f. [Kel79]). Hence, the probability of a certain state of the system occuring is the product of the probabilities of nodes containing a certain number of packets. As the state of the system directly implies a certain liveliness, also the liveliness vector can be found as the product of the liveliness of separate nodes, showing (7.9) holds.  □

We will now use equations (7.4)-(7.9) as an approximation for the discrete time model. This rather coarse approximation will provide quite accurate results, as we are interested in the influence of the load on *average* behaviour of the network.

## 7.6  Stability

The average service rate $\hat{r}_i$ at which each node operates determines the load under which the network is stable. As presented earlier, the average service rate of a node $i$ is given by (7.8) and the probability that a node is alive by (7.7). Writing out the expression for $\hat{q}_\pi$ and inserting (7.7) into (7.8), we obtain $n$ equations, with $2n$ unknowns, which are the $a_i$ and the $\hat{r}_i$. Assuming that all nodes are stable, that is when all $a_i < \hat{r}_i$, the arrival rate at each node is known. The values for $\hat{r}_i$ can hence be calculated for a stable system. However, it is still to be determined for which values of $\lambda_j$ (and thus $a_i$) the network is stable.

As presented in Section 7.4, the arrival rate for a certain flow $j$ at node $t_j(i)$ is given by (7.6) and the total arrival rate by (7.4). Using the $n$ equations (7.4) for $a_i$, it is possible to solve the system of $2n$ unknown variables, which entails solving polynomials of degrees that increase exponentially with the network size. Solutions can be obtained numerically, however, using for instance the Algorithm 7.1 to obtain the values of $\hat{r}_i$.

To analyze the convergence of Algorithm 7.1, we consider the separate steps and the recursion. The initial value of $\hat{r}_i = 1$ corresponds to a network without contention, immediately giving an indication whether the network is stable or not. To calculate all $\lambda_j(k)$'s in step 2), the equations (7.6), (7.4) and (7.5) need to be combined, giving $Jm$ equations with equally many unknown variables which can be solved. From these values, obviously steps 3) through 6) can be calculated, leading to the recursion.

Let $g(r)$ denote the function that calculates the new value of $r$ using the steps described. The function $g(.) : \mathbb{R}^n \to \mathbb{R}^n$ is a continuous function on the convex compact subset $[0,1]^n$. Following Brouwers fixed point theorem (c.f. [Ist81]), we consider the equation $g(r) = r$, which has a solution, which we need to show to be the unique fixed point. To achieve this, we use the Contraction Mapping Theorem (CMT, c.f. [Ist81]), saying that the equation $g(r) = r$ has a unique solution if and only if

- The function $g(.)$ maps $[0,1]^n$ to $[0,1]^n$

- There is a constant $G < 1$ such that $||g(x) - g(y)|| \leq G||x - y||$ for all $x, y \in [0,1]^n$

---

**Algorithm 7.1** Algorithm to calculate the service rate of all nodes

Calculation of $\hat{r}_i, i = 1, \ldots, n$

1. Set all values $\hat{r}_i$ to 1, $i = 1, \ldots, n$

2. Calculate $\lambda_j(k)$, $j = 1, \ldots, J$
   and $k = t_j(1), \ldots, t_j(m_k)$

3. Calculate $a_i = \sum_{j : i \in t_j} \lambda_j(i)$, $i = 1, \ldots, n$

4. Calculate $\hat{p}_i = min(\frac{a_i}{\hat{r}_i}, 1)$, $i = 1, \ldots, n$

5. Calculate $\hat{q}_\pi = \prod_{i=1}^{n} (1 - \hat{p}_i)^{(1-\pi_i)} \hat{p}_i^{\pi_i}$, $\pi \in \Pi$

6. Calculate new $\hat{r}_i = \sum_{\pi \in \Pi} \frac{\hat{r}_{i,\pi} \hat{q}_\pi}{\hat{p}_i}$, $i = 1, \ldots, n$

7. Calculate the difference $\epsilon_i = \hat{r}_i(new) - \hat{r}_i(old)$, $i = 1, \ldots, n$

8. Repeat step 2 till 7 until convergence occurs, that is $|\epsilon| \leq \delta$ for an appropriate value of $\delta$.

---

First, the algorithm needs to be shown to map any starting value for $r$ to another value of $r$ that is within the possible range of $[0, 1]^n$. For this to be the case, we need that

$$0 \leq \sum_\pi \hat{r}_{i,\pi} \hat{q}_\pi \leq \hat{p}_i.$$

The first inequality is obvious, for the second one we note that $\hat{r}_{i,\pi} = 0$ for all $\pi$ such that $\pi_i = 0$ and that $\hat{r}_{i,\pi} \leq 1$. This gives that

$$
\begin{aligned}
\sum_\pi \hat{r}_{i,\pi} \hat{q}_\pi \quad &\leq \quad \sum_{\pi : \pi_i = 1} \hat{q}_\pi \\
&= \quad \sum_{\pi : \pi_i = 1} \prod_{j=1}^{n} (1 - \hat{p}_j)^{1-\pi_j} \hat{p}_j^{\pi_j} \\
&= \quad \hat{p}_i \sum_{\pi : \pi_i = 1} \prod_{j \neq i} (1 - \hat{p}_j)^{1-\pi_j} \hat{p}_j^{\pi_j} = \hat{p}_i,
\end{aligned}
$$

where the last equality holds as we sum over all possible liveliness states for the network without node $i$, proving the first part of the contraction mapping theorem.

The second part is more involved. We provide a complete proof for a two node network and indicate why the second condition is conjectured to hold for larger networks. When following the steps of the algorithm for a two node tandem network, we have that

$$a(1) \quad = \quad \lambda(1) = \lambda \text{ and } a(2) = \lambda(2) = \min(\lambda, \hat{r}_1)$$

$$\hat{p}_1 = \min(\frac{\lambda}{\hat{r}_1}, 1) \text{ and } \hat{p}_2 = \min(\frac{\min(\lambda, \hat{r}_1)}{\hat{r}_2}, 1)$$

$$\hat{r}_1 = 1 - \frac{1}{2}\hat{p}_2 \text{ and } \hat{r}_2 = 1 - \frac{1}{2}\hat{p}_1.$$

Note that as $0 \le \hat{p}_i \le 1$ we have that $\hat{r}_i \in [\frac{1}{2}, 1]$. First assuming we are dealing with a stable network, the arrival rate at both nodes equals $\lambda$. By substituting $p_i$, we obtain the functional vector

$$g(\hat{r}) = (1 - \frac{\lambda}{2\hat{r}_2}, 1 - \frac{\lambda}{2\hat{r}_1}).$$

This gives, for $x = (x_1, \ldots, x_n)$,

$$||g(x) - g(y)||^2 = (\frac{\lambda}{2x_2 y_2})^2 (x_2 - y_2)^2$$
$$+ (\frac{\lambda}{2x_1 y_1})^2 (x_1 - y_1)^2,$$

and for this to be smaller than $||x - y||^2$ we need to have that $(\frac{\lambda}{2x_i y_i})^2 < 1$. As we assumed a stable network, we have that $\lambda < x_i$, so that

$$\frac{\lambda}{2x_2 y_2} < \frac{1}{2y_i} \le 1$$

since $y_i \in [\frac{1}{2}, 1]$ and so indeed the second condition holds proving that for a stable system the algorithm converges. If the system would be unstable, we have that

$$g(\hat{r}) = \left(1 - \frac{\min\left(\frac{\min(\lambda, \hat{r}_1)}{\hat{r}_2}, 1\right)}{2}, 1 - \frac{\min\left(\frac{\lambda}{\hat{r}_1}, 1\right)}{2}\right),$$

where the following situations can occur: $\lambda \ge \hat{r}_1$ or $\hat{r}_2 \le \lambda < \hat{r}_1$. In the first case we have that

$$g(\hat{r}) = (1 - \frac{1}{2}\min(\frac{\hat{r}_1}{\hat{r}_2}, 1), \frac{1}{2})$$

which within two steps of the algorithm leads to $g(\hat{r}) = (\frac{1}{2}, \frac{1}{2})$ and thus converges to this unique solution. In the second case we have that

$$g(\hat{r}) = (\frac{1}{2}, 1 - \frac{\lambda}{2\hat{r}_1}),$$

$$||g(x) - g(y)||^2 = (\frac{\lambda}{2x_1 y_1})^2 (x_1 - y_1)^2$$

and $(\frac{\lambda}{2x_1 y_1})^2 < 1$ as shown earlier, completing the proof that the algorithm converges for this two node network.

Considering a three node network, we obtain the following function (ommiting

the hat in the notation):

$$g(r) = (1 - \frac{\min(\frac{\min(\lambda,r_1)}{r_2}, 1)}{2}$$
$$+ \frac{\min(\frac{\min(\lambda,r_1)}{r_2}, 1) \min(\frac{\min(\min(\lambda,r_1),r_2)}{r_3}, 1)}{6},$$
$$1 - \frac{\min(\frac{\lambda}{r_1}, 1)}{2} - \frac{\min(\frac{\min(\min(\lambda,r_1),r_2)}{r_3}, 1)}{2}$$
$$+ \frac{\min(\frac{\lambda}{r_1}, 1) \min(\frac{\min(\min(\lambda,r_1),r_2)}{r_3}, 1)}{3},$$
$$1 - \frac{\min(\frac{\min(\lambda,r_1)}{r_2}, 1)}{2} + \frac{\min(\frac{\lambda}{r_1}, 1) \min(\frac{\min(\lambda,r_1)}{r_2}, 1)}{6}).$$

As we have that $g(p) = (1 - \frac{p_2}{2} + \frac{p_2 p_3}{6}, 1 - \frac{p_1}{2} - \frac{p_3}{2} + \frac{p_1 p_3}{3}, 1 - \frac{p_2}{2} + \frac{p_1 p_2}{6})$, starting in $([\frac{1}{2}, 1], [\frac{1}{3}, 1], [\frac{1}{2}, 1])$, $g(.)$ will also project on this range. For the CMT to hold, we first consider the stable system again, so that $\lambda < r_i$. In this case we have that

$$g(r) = (1 - \frac{1}{2}\frac{\lambda}{r_2} + \frac{1}{6}\frac{\lambda^2}{r_2 r_3}, 1 - \frac{1}{2}\frac{\lambda}{r_1} - \frac{1}{2}\frac{\lambda}{r_3} + \frac{1}{3}\frac{\lambda^2}{r_1 r_3},$$
$$1 - \frac{1}{2}\frac{\lambda}{r_2} + \frac{1}{6}\frac{\lambda^2}{r_1 r_2}).$$

Checking whether $||g(x) - g(y)|| < ||x - y||$ proves to be cumbersome, even for such a small network. Therefore we numerically analyzed the function $h(x, y) = ||g(x) - g(y)||(||x - y||)^{-1}$ which proved to be smaller than one for all values of $x$ and $y$. As in the two node network, it is easy to show that for an instable network, either there is an obvious direct convergence to the rates $(\frac{2}{3}, \frac{1}{3}, \frac{2}{3})$ or convergence is proven by using parts of the approach for the stable case. We postulate that for any network a similar analysis will show that the algorithm constitutes a contraction, and thus converges.

We have numerically established that Algorithm 7.1 converges to a unique solution $\hat{r}_i$ for any values of $\lambda_j$, $j = 1, \ldots, J$. Using Algorithm 7.1, the service rate of all nodes can be calculated for any set of flows through the network. The corresponding arrival rates at the destination nodes of the flows give the throughput of the network. Whenever the network is stable, the total throughput will equal $\sum_j \lambda_j$. For a general network, the calculation of the throughput, independent of the topology of the network, involves solving $n$ equations in $n$ unknowns. Using Algorithm 7.1, the arrival rate(s) can be chosen arbitrarily. To determine the stability range of the network, we separately consider each flow in the network. Fixing the arrival rates of all but one flow (such that the system with these flows is stable), there exists a value $\lambda_{opt}$ for the remaining flow such that $a_k = \hat{r}_k$ for at least one $k \in 1, \ldots, n$, which provides the maximal throughput $\lambda_{opt}$ of this flow. Node $k$ is then the bottleneck of the network. In this manner the stability range of the network can be calculated (examples are shown in the following section).

## 7.7 Examples and validation

### 7.7.1 Multihop tandem network

In the following we analyze a multihop tandem network. When considering a general network, the analysis of the stability region involves considering flows separately. First, we show how for a specific contention protocol the transmission probabilities $r_{i,\pi}$ can be calculated in this network, which corresponds to a single multihop transmission in a network. Next, we use simulation to validate results obtained by our algorithm for different sizes of the network. Some surprising results are obtained, which are correctly predicted by our model.

Consider a tandem network of size $n$. The average service rate at which a node transmits depends on the position in the tandem network. As indirectly all nodes in the network influence each other, the total length of the network has an impact. This impact when all nodes are alive is shown, using a contention protocol selecting a node to transmit with equal probability among all alive nodes.

Consider the tandem network such that nodes cannot transmit and receive at the same time. A node that is allowed to transmit hence blocks its direct neighbour(s). When all $n$ nodes are alive, each node has a probability $\frac{1}{n}$ of obtaining the channel directly and blocking its neighbour(s). The remaining nodes continue contending for the channel until they are either blocked or allowed to transmit. The rate $r_{i,\mathbf{1}}(n)$ for a node at position $i$ in a fully alive tandem network of length $n$ can be calculated using

$$
\begin{aligned}
r_{i,\mathbf{1}}(n) \quad = \quad & \frac{1}{n}[\sum_{k=1}^{i-2} r_{i-k-1,\mathbf{1}}(n-k-1) \\
& +1 + \sum_{k=i+2}^{n} r_{i,\mathbf{1}}(k-2)].
\end{aligned} \tag{7.10}
$$

The right hand side of (7.10) follows from the node winning the contention: If the first node in the network wins the contention, it blocks the second node and the remaining $n-2$ nodes compete, with node $i$ now at position $i-2$. Otherwise, in a similar manner, a node before (but not a neighbouring) node $i$ wins the contention, node $i$ wins the contention itself, either of node $i$'s neighbours wins the contention or a node $k$ behind node $i$ wins the contention. Each of these events occurs with a probability of $\frac{1}{n}$, together giving the recursive formula.

Note that a multihop tandem network (in this setting) with nodes that are not alive can be decomposed into many smaller multihop networks. For a fully alive tandem network where nodes cannot transmit and receive at the same time, Table 7.1 shows the rates for different lengths of the network.

**Theorem 7.2.** *For the multihop tandem network with all alive nodes, the rate allocated to the nodes converges when the network size increases, where in particular*

$$
\lim_{n\to\infty} r_{1,\mathbf{1}}(n) = 1 - \frac{1}{e} \quad and \quad \lim_{n\to\infty} r_{2,\mathbf{1}}(n) = \frac{1}{e}. \tag{7.11}
$$

| Size Node | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | - | - | - | - | - | - | - | - | - | - | - |
| 2 | 0.5 | 0.5 | - | - | - | - | - | - | - | - | - | - |
| 3 | 0.6666 | 0.3333 | 0.666 | - | - | - | - | - | - | - | - | - |
| 4 | 0.625 | 0.375 | 0.375 | 0.625 | - | - | - | - | - | - | - | - |
| 5 | 0.6333 | 0.3667 | 0.4667 | 0.3667 | 0.6333 | - | - | - | - | - | - | - |
| 6 | 0.6319 | 0.3681 | 0.4444 | 0.4444 | 0.3681 | 0.6319 | - | - | - | - | - | - |
| 7 | 0.6321 | 0.3679 | 0.4488 | 0.4262 | 0.4488 | 0.3679 | 0.6321 | - | - | - | - | - |
| 8 | 0.6321 | 0.3679 | 0.4481 | 0.4297 | 0.4297 | 0.4481 | 0.3679 | 0.6321 | - | - | - | - |
| 9 | 0.6321 | 0.3679 | 0.4482 | 0.4291 | 0.4334 | 0.4291 | 0.4482 | 0.3679 | 0.6321 | - | - | - |
| 10 | 0.6321 | 0.3679 | 0.4482 | 0.4292 | 0.4328 | 0.4328 | 0.4292 | 0.4482 | 0.3679 | 0.6321 | - | - |
| 11 | 0.6321 | 0.3679 | 0.4482 | 0.4292 | 0.4329 | 0.4322 | 0.4329 | 0.4292 | 0.4482 | 0.3679 | 0.6321 | - |
| 12 | 0.6321 | 0.3679 | 0.4482 | 0.4292 | 0.4329 | 0.4323 | 0.4323 | 0.4392 | 0.4292 | 0.4482 | 0.3679 | 0.6321 |

Table 7.1: Transmission probability for a fully alive tandem network

*Proof.* The formula for the rate $r_{i,\mathbf{1}}(n)$ of a node on position $i$ in an $n$ node network that is fully alive is given by

$$nr_{i,\mathbf{1}}(n) = \sum_{k=1}^{i-2} r_{i-k-1,\mathbf{1}}(n-k-1) + 1 + \sum_{k=i}^{n-2} r_{i,\mathbf{1}}(k) \qquad (7.12)$$

as described in the paper. Due to symmetry of the network we also have that

$$r_{i,\mathbf{1}}(n) = r_{n-i+1,\mathbf{1}}(n) \qquad i = 1, \dots, n.$$

The rate of a node can never exceed one, but will be one if the node is the only alive node within its interference region, i.e. its neighbours are not alive. The minimal rate of a node is $\frac{1}{n}$ as with this probability it wins the contention over all other nodes.

In the following we omit the $\mathbf{1}$ denoting the fully alive network. To find an expression for $r_i(n)$, note that

$$nr_i(n) - (n-1)r_i(n-1) = r_i(n+2)$$
$$+ \sum_{k=1}^{i-2} [r_{i-k-1}(n-k-1) - r_{i-k-1}(n-k-2)]$$

and letting $c_i(n) = r_i(n) - r_i(n-1)$ this gives

$$c_i(n) = \frac{1}{n}[\sum_{k=1}^{i-2} c_k(n+k-i) - c_i(n-1)].$$

As $-1 \le c_i(n) \le 1$ for any value of $i$ and $n$, we have that

$$c_i(n) \le \frac{1}{n}[(i-2) - c_i(n-1)] \le \frac{1}{n}(i-1)$$

so that for each $i$ we have that $\lim_{n \to \infty} c_i(n) = 0$, proving that $r_i(n)$ converges for $n \to \infty$.

For $i = 1$ this leads to

$$c_1(n) = -\frac{1}{n}c_1(n-1)$$

which gives

$$c_1(n) = \frac{(-1)^{n-1}}{n!} \ , \ r_1(n) = \sum_{i=1}^{n} \frac{(-1)^{i-1}}{i!}.$$

Similarly, we have that

$$c_2(n) = \frac{(-1)^n}{n!} \ , \ r_2(n) = \sum_{i=1}^{n} \frac{(-1)^i}{i!}.$$

Taking the limit shows that

$$\lim_{n\to\infty} r_1(n) = 1 - \frac{1}{e} \ , \ \lim_{n\to\infty} r_2(n) = \frac{1}{e}$$

Unfortunately, for larger values of $i$, no nice expressions are found for $c_i(n)$ or $r_i(n)$, but the limiting values can be calculated using the same approach. The results are presented in Table 7.1. $\qquad\square$

Other limits are observed in Table 7.1, showing that the border effects fade for the middle nodes as the length of the network increases, in accordance with [DDT08]. This border effect already starts to fade for networks of size 12.

We note that the calculation of the rates $r_{i,\pi}$ for the linear setting has the pleasant property that the rate of a certain node $i$ under liveliness $\pi$ is only dependent on the number of nodes that are alive and directly connected to each other. When considering different contention sets and protocol or network layout, this property however may no longer be present.

To validate the results presented in this chapter, a simulation model has been constructed that mimics the behaviour of the discrete time network under consideration. The arrival and processing of the packets is modeled, with a simulation for each parameter setting lasting one million simulated time slots after a warm up period of 100.000 slots. The results are compared with the stability ranges and the throughput of the network calculated with the state independent continuous time approximation, using the provided algorithm. For some settings, we provide the exact derivation of the results.

Consider the multihop tandem network for $n = 3$. The average service rates at which the nodes operate are given by (using (7.8) and (7.9))

$$
\begin{aligned}
\hat{r}_1 &= (1 - \hat{p}_2) + \frac{1}{2}\hat{p}_2(1 - \hat{p}_3) + \frac{2}{3}\hat{p}_2\hat{p}_3 & (7.13) \\
\hat{r}_2 &= (1 - \hat{p}_1)(1 - \hat{p}_3) + \frac{1}{2}\hat{p}_1(1 - \hat{p}_3) \\
&\quad + \frac{1}{2}(1 - \hat{p}_1)\hat{p}_3 + \frac{1}{3}\hat{p}_1\hat{p}_3 \\
\hat{r}_3 &= (1 - \hat{p}_2) + \frac{1}{2}(1 - \hat{p}_1)\hat{p}_2 + \frac{2}{3}\hat{p}_1\hat{p}_2.
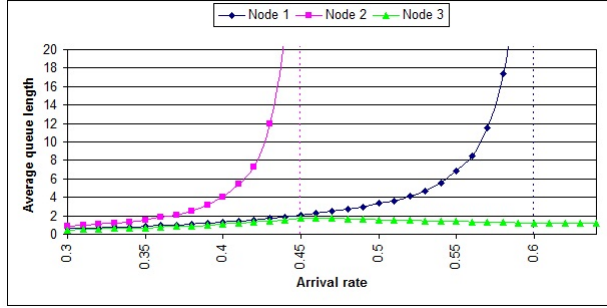\end{aligned}
$$

Figure 7.3: Simulated average queue length in a 3 hop network with the instability rates calculated by the model

Obviously, the second node will be the bottleneck of the network as $\hat{r}_2$ is smaller than $\hat{r}_1$ and $\hat{r}_3$, as it is the only node contending with two neighbours. When node 2 is unstable, we have that $\hat{p}_2 = 1$. To determine at what arrival rate $\lambda$ this will occur, we use that $\lambda = a_i = \hat{r}_2$, so that

$$\hat{p}_1 = \frac{\hat{r}_2}{\hat{r}_1} \text{ and } \hat{p}_3 = \frac{\hat{r}_2}{\hat{r}_3}. \tag{7.14}$$

Combining the equations (7.13) and (7.14) with $\hat{p}_2 = 1$ we find that $\hat{p}_1 = \hat{p}_3 = \frac{9-\sqrt{57}}{2}$, resulting in the critical arrival rate of $\lambda = \hat{r}_2 = 8 - \sqrt{57}$. From this value of $\lambda$ on the second node will be unstable. If we increase the arrival rate even more, the first node will also become unstable. The third node however will always remain stable, as its service rate will always be higher than the service rate at the second node, which determines the arrival rate at the third node. To find from which value of $\lambda$ on the first node will also be unstable, we substitute $\hat{p}_1 = \hat{p}_2 = 1$ in (7.13) which leads to $\hat{p}_3 = 0.6$, and the rate at which node 1 becomes unstable equals $\lambda = \hat{r}_1 = 0.6$. Also note that the rate of the second node has now fallen to a value of $\hat{r}_2 = 0.4$, so that the throughput of the network has decreased.

For the three node tandem network, Figure 7.3 shows the average queue length at the three nodes for increasing load of the system and Figure 7.4 shows the throughput of the system. The calculated values of arrival rates for which queues become unstable are depicted as dotted vertical lines in the figures.

As can be seen in Figs. 7.3 and 7.4, the arrival rates at which the first and second node become unstable coincide with the calculated values. Additional simulations for the arrival rates near the ones causing instability of nodes were performed to confirm the results, but are not shown in the figures to maintain readability. The throughput, which reaches a maximum of $8 - \sqrt{57} \approx 0.4501$ when the second node becomes unstable, decreases after this value. This decrease in throughput is caused by the decrease in service rate at the second node, as the first node becomes more highly loaded. This causes the first queue to be alive a larger fraction of the time, blocking the second node. The throughput settles at 0.4 after the first node has become unstable at an arrival rate of 0.6,
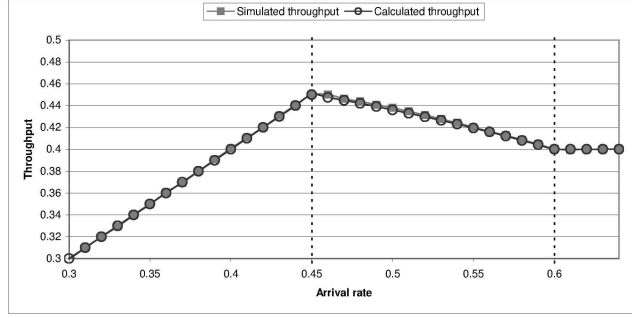
Figure 7.4: Simulated and calculated throughput of a 3 hop network with the instability rates calculated by the model

which is in agreement with the values calculated.

Next, considering a larger network with 5 hops, one might expect that it is the second node that becomes the bottleneck. Using the presented model and setting $\hat{p}_2 = 1$ however shows that no real valued solution exists, meaning that node 2 cannot be the node to become unstable first. It actually is the third node that becomes the bottleneck first at an arrival rate of 0.4323, which is the maximum throughput of the network. Increasing the arrival rate to 0.4448 causes the second node to become unstable as well. Increasing the arrival rate further, the third node becomes stable again. The presented model also determines the arrival rate at which this occurs by making a small adjustment to the equations. As the third queue will become stable as soon as its average service rate is lower than the second queue's rate, we now set $\hat{r}_2 = \hat{r}_3$. As both queues are still unstable we have that $\hat{p}_2 = \hat{p}_3 = 1$ and that $\hat{p}_4 = \frac{\hat{r}_2}{\hat{r}_4}$ and $\hat{p}_5 = \frac{\hat{r}_2}{\hat{r}_5}$. Using the standard equations for the $\hat{r}_i$'s and setting $\hat{p}_1 = \frac{\lambda}{\hat{r}_1}$, we solve the system to obtain $\lambda = 0.4803$ and $\hat{r}_2 = \hat{r}_3 = 0.4306$. Finally increasing the arrival rate to 0.6108 causes the first node to become unstable, resulting in a throughput of 0.3892. Simulation of the network under consideration provided the results as presented in Figures 7.5 and 7.6 where the vertical lines show the calculated values for which nodes become (un)stable.

That node 3 is the first node to become unstable can be called surprising. When all queues are alive, the average service rate of queue 2 is lower than that of queue 3. However, when queue 1 and/or queue 5 are empty, the third queue has the lowest rate (see Table 7.1 for a 3 to 5 node network). As can be seen in Figure 7.5, the average queue length at nodes 1 and 5 are low for the load when queue 2 and 3 are already reaching instability. This indicates that they frequently will not be alive, which is in the disadvantage of the third node, making it the bottleneck node. However, as the arrival rate increases, nodes 1 and 5 will be alive more often, which is beneficial for node 3, resulting in the queue becoming stable again. Surprising as this behaviour may be, it is predicted correctly by the model.
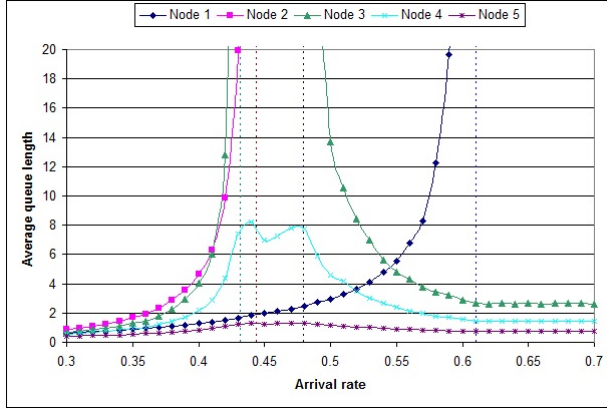
Figure 7.5: Simulated average queue length in a 5 hop network with the instability rates calculated by the model
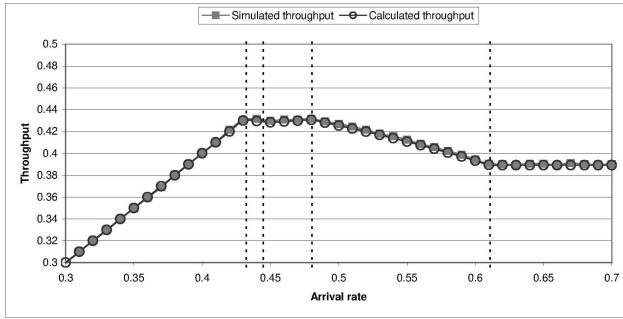


Figure 7.6: Simulated and calculated throughput of a 5 hop network with the instability rates calculated by the model

### 7.7.2   General eight node network

Consider the network as depicted in Figure 7.2. Note that any set of interfering nodes can be used, mimicking the behaviour of any acces control protocol, i.e. to mimic an RTS/CTS protocol all nodes within transmission range of the sending and receiving node can be used as the contention set. To avoid trivial results we set the interference ranges for this example to be (only showing the nodes that need to transmit) $I(1) = \{2, 4\}$, $I(2) = \{1\}$, $I(4) = \{1, 5, 6\}$, $I(5) = \{4, 6\}$, $I(6) = \{4, 5, 7\}$, $I(7) = \{6\}$. First flow $f(t_1)$ is set up, with rate $\lambda_1 = 0.1$. Obviously the network can handle this flow. Second, flow $f(t_3)$ is set up, with rate $\lambda_3 = 0.1$ as well. Again, the network remains stable (note that even though both flows have node 3 as endpoint, this does not cause problems as we assume perfect reception of all transmissions). Now flow $f(t_2)$ is initiated and the open question is which rate can be achieved for this flow. The arrival rates of traffic
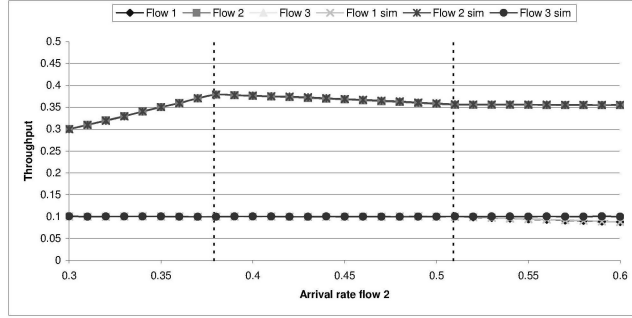
Figure 7.7: Simulated and calculated throughput of a 8 node network with the instability rates calculated by the model

at the nodes, as long as the network is stable, is given by

| Node | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Arrival rate | $\lambda_2 + 0.1$ | 0.1 | 0.2 | $\lambda_2$ |
| Node | 5 | 6 | 7 | 8 |
| Arrival rate | 0.1 | $\lambda_2 + 0.1$ | 0.1 | $\lambda_2$ |

and the probabilities $q_\pi$ of all possible liveliness vectors can easily be calculated. Using these values in equations (7.7) and (7.8) gives a set of 8 equations with 9 unknowns (all the $\hat{r}_i$ and $\lambda_2$), which can be solved when it is known which node becomes the bottleneck. Using $\lambda_2 = \hat{r}_i$ for $i = 1, 4, 6$ and solving shows that it is node 4 that becomes the bottleneck at an arrival rate of $\lambda_2 = 0.3789$. Increasing the arrival rate $\lambda_2$ further causes node 1 to become unstable as well, influencing the throughput of the first flow. Using the model, this is calculated to happen at an arrival rate of $\lambda_2 = 0.5092$. Figure 7.7 shows the throughput of the separate flows for an increasing arrival rate of the second flow. Both the values calculated by the model and the simulation results are shown.

Numerical evaluation shows that the model gives very accurate predictions of the throughput, where the error at each calculated point stays below 1%. The load at which node 1 and 4 become unstable can be recognised as the points where the slope of the graph changes, where the simulation again shows that this is at the point predicted by the model.

## 7.8 Conclusion

As interference limits the capacity of wireless ad hoc networks, networks with contending nodes are analyzed in this chapter. Each time slot, nodes compete to transmit a packet from their queue, where a winning node blocks other nodes in its neighbourhood. The time slot system is approximated in two steps. First, by considering the long run average behaviour of the discrete time system, a continuous time model is obtained. As the second step, appropriate state independent service rates for the nodes in the network are determined. Combining

relations between the arrival and service rates of the nodes, bottleneck nodes are identified which determine the throughput of a multihop wireless network. Using the two rather coarse approximation steps, we propose a product form network approximation. Taking advantage of the properties of product form networks, equations for the liveliness vector (whether nodes have packets in their queues or not) and the average service rates of the nodes are derived and solved using a simple algorithm. Surprisingly, the continuous approximation for the long term average behaviour turns out to give accurate results concerning the stability and throughput of the network. Other performance measures, as the queue length and waiting time, have not been considered.

Our approach provides very accurate results for the lowest arrival rate of a flow at which one of the nodes becomes unstable, thus giving the maximal throughput for this flow. Also, increasing the arrival rate further, instability of the rest of the nodes is analyzed. Our model correctly predicts surprising behaviour in a multihop tandem network, where a queue at first turning out to be the bottleneck, returned to stability again after increasing the arrival rate. The approach presented is applicable for general networks, with various contention settings and protocols. Using simulations of the discrete time system, results were compared with the continuous time model, showing that the model provides very accurate results.

# References

[AAB03]    K. Avrachenkov, U. Ayesta, and P. Brown. Batch arrival M/G/1 processor sharing with application to multilevel processor sharing scheduling. *INRIA Technical Report 5043*, 2003.

[BB02]     N. Bayer and R. J. Boucherie. On the structure of the space of geometric product-form models. *Probability in Engineering and Informational Science*, 16:241–270, 2002.

[Bia00]    G. Bianchi. Performance analysis of the IEEE 802.11 distributed coordination functions. *IEEE Journal on Selected Areas in Communications*, 18:535–547, 2000.

[BJL08]    S. C. Borst, M. Jonckheere, and L. Leskelä. Stability of parallel queueing systems with coupled service rates. *Discrete Event Dynamic Systems*, 18:447–472, 2008.

[BP02]     T. Bonald and A. Proutière. Insensitivity in processor-sharing networks. *Proceedings of Performance '02*, 2002.

[BP03]     T. Bonald and A. Proutière. Insensitive bandwidth sharing in data networks. *Queueing Systems*, 44:69–100, 2003.

[BvD09]    R. J. Boucherie and N. M. van Dijk. Monotonicity and error bounds for networks of Erlang loss queues. *Queueing systems*, 62:159–193, 2009.

[BW89]     O. J. Boxma and J. A. Weststrate. Waiting times in polling systems with markovian server routing. *Informatik-Fachberichte*, 218:89–104, 1989.

[C⁺05]     S. K. Cheung et al. An analytical packet/flow-level modelling approach for wireless LANs with Quality-of-Service support. *Lecture Notes in Computer Science, Proceedings of the 19th International Teletraffic Congress*, 2005.

[CB83]     J. W. Cohen and O. J. Boxma. *Boundary value problems in queueing system analysis*. North Holland Publishing Company, Amsterdam, 1983.

[CBvB05a]  S. K. Cheung, J. L. Berg v.d., and R. J. Boucherie. Decomposing the queue length distribution of processor-sharing models into queue lengths of permanent customer queues. *Proceedings of IFIP's Performance 2005*, 2005.

[CBvB05b]  S. K. Cheung, J. L. Berg v.d., and R. J. Boucherie. Insensitive bounds for the moments of the sojourn time distribution in the M/G/1 processor-sharing queue. *Research Memorandum 1766 of University of Twente*, 2005.

[CdGB07]   T. J. M. Coenen, M. de Graaf, and R. J. Boucherie. An upper bound on multi-hop wireless network performance. *Lecture Notes on Computer Science*, 4516:335–347, 2007.

[CGJ+06]    T. J. M. Coenen, P. T. H. Goering, A. Jehangir, J. L. van den Berg, R. J. Boucherie, S. M. Heemstra de Groot, G. J. Heijenk, S. S. Dhillon, W. Lu, A. Lo, P. F. A. van Mieghem, and I. G. M. M. Niemegeers. Architectural and QoS aspects of personal networks. *Proceedings of First International Workshop on Personalized Networks, PerNets 2006*, pages 1–3, 2006.

[CHE02]     M. Cagalj, J. Hubaux, and C. Enz. Minimum-energy broadcast in all-wireless networks, NP-completeness and distribution issues. *Proceedings of the Annual International Conference on Mobile Computing and Networking (MOBICOM)*, pages 172–182, 2002.

[D+08]      D. Denteneer et al. IEEE 802.11s and the philosophers' problem. *Statistica Neerlandica*, 62:283–298, 2008.

[DDT08]     M. Durvy, O. Dousse, and P. Thiran. Border effects, fairness, and phase transition in large wireless networks. *Proceedings of INFOCOM 2008*, pages 601–609, 2008.

[EKEA07]    R. El-Khoury and R. El-Azouzi. Stability-throughput analysis in a multi-hop ad hoc networks with weighted fair queueing. *45th Annual Allerton Conference*, pages 1066–1073, 2007.

[EOs05]     P. E. Engelstad and O. N. Ø sterbø. Non-saturation and saturation analysis of IEEE 802.11e EDCA with starvation prediction. *Proceedings of ACM MSWiM '05*, 2005.

[F+06]      T. Fujiwara et al. Throughput analysis on string multi-hop networks with multiple flow. *NCSP '06*, pages 417–420, 2006.

[FC85]      S. W. Fuhrmann and R. B. Cooper. Stochastic decompositions in the M/G/1 queue with generalized vacations. *Operations Research*, 33:1117–1129, 1985.

[FF56]      L. R. Ford and D. R. Fulkerson. Maximal flow through a network. *Canadian Journal of Mathematics*, 8:399–404, 1956.

[FF62]      L. R. Ford and D. R. Fulkerson. *Flows in Networks*. Princeton University Press, Princeton, NJ, 1962.

[FI79]      G. Fayolle and R. Iasnogorodski. Two coupled processors: The reduction to a riemann-hilbert problem. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 47:325–351, 1979.

[FIM99]     G. Fayolle, R. Iasnogorodski, and V. Malyshev. *Random walks in the quarter plane: algebraic methods, boundary value problems and applications*. Springer, 1999.

[FMI80]     G. Fayolle, I. Mitrani, and R. Iasnogorodski. Sharing a processor among many job classes. *Journal of the Association for Computing Machinery*, 27:519–532, 1980.

[Fuh84]     S. W. Fuhrmann. A note on the M/G/1 queue with server vacations. *Operations Research*, 32:1368–1373, 1984.

[Fuh85]     S. W. Fuhrmann. Symmetric queues served in cyclic order. *Operations Research Letters*, 4:139–144, 1985.

[GBvO13]    J. Goseling, R. J. Boucherie, and J. C. W. van Ommeren. Linear programming error bounds for random walks in the quarter-plane. *Under review*, 2013.

[GK00]      K. Gupta and P. R. Kumar. The capacity of wireless networks. *IEEE Transactions on Information Theory*, 46:388–404, 2000.

[GMW04]     R. Gupta, J. Musacchio, and J. Walrand. Sufficient rate constraints for QoS flows in ad-hoc networks. *UCB/ERL Technical Memorandum M04/42*, 2004.

[GO09]      M. de Graaf and J. C. W. van Ommeren. Increasing network lifetime by battery-aware master selection in radio networks. *Technical Report of the University of Twente*, 2009.

[Gro90]     W. P. Groenendijk. *Conservation laws in polling systems*. PhD thesis, 1990.

[GV02]      M. Gastpar and M. Vetterli. On the capacity of wireless networks: The relay case. *Proceedings of IEEE INFOCOM '02*, 2002.

[HBO06]     R. de Haan, R. J. Boucherie, and J. C. W. van Ommeren. Modelling multi-path signal interference in ad-hoc networks. *Technical Memorandum Dep. of Applied Mathematics, University of Twente*, 2006.

[HvM04]     R. Hekmat and P. van Mieghem. Interference in wireless multi-hop ad-hoc networks and its effect on network capacity. *Wireless Networks*, 10:389–399, 2004.

[HVS01]     Z. Hadzi-Velkov and B. Spasenovski. Capture effect in IEEE 802.11 wireless LANs. *Proceedings of IEEE ICWLHN '01*, 2001.

[IEE]       IEEE. IEEE website. https://www.ieee.org/index.html.

[IEE05]     IEEE standard for wireless LAN medium access control (MAC) and physical layer (PHY) specifications, medium access control (MAC) quality of service enhancements, 2005.

[Ist81]     V. I. Istratescu. *Fixed point theory: An introduction*. D. Reidel Publishing, Dordrecht, Holland, 1981.

[J$^+$03]   K. Jain et al. Impact of interference on multihop wireless network performance. *Proceedings of ACM Mobicom 2003*, pages 66–80, 2003.

[J$^+$04]   Z. Jia et al. Bandwidth guaranteed routing for ad-hoc networks with interference constraints. *Proceedings of ISCC*, 2004.

[Jac57]     J. Jackson. Network of waiting lines. *Operations Research*, 5:518–521, 1957.

[JS03]      J. Jun and M. L. Sichitiu. The nominal capacity of wireless mesh networks. *IEEE Wireless Communications*, pages 8–14, 2003.

[JvdMvdW07] M. Jonckheere, R. D. van der Mei, and W. van der Weij. Rate stability and output rates in queueing networks with shared resources. *Technical Report of CWI*, 2007.

[KEAA08]    A. Kherani, R. El Azouzi, and E. Altman. Stability-throughput tradeoff and routing in multihop wireless ad-hoc networks. *Computer Networks*, 52:1365–1389, 2008.

[Kel79]     F. P. Kelly. *Reversibility and Stochastic Networks*. Wiley, Chichester, 1979.

[Kle67]     L. Kleinrock. Time-shared systems: A theoretical treatment. *Journal of the Association for Computing Machinery*, 14:242–261, 1967.

[KMR71]   L. Kleinrock, M. M. Muntz, and E. Rodemich. The processor sharing queueing model for time shared systems with bulk arrivals. *Networks: an international journal*, 1:1–13, 1971.

[KN03]    M. Kodialam and T. Nandagopal. Characterizing achievable rates in multihop wireless networks: The joint routing and scheduling problem. *Proceedings of ACM Mobicom 2003*, pages 42–54, 2003.

[KN05]    M. Kodialam and T. Nandagopal. Characterizing the capacity region in multi-radio multi-channel wireless mesh networks. *Proceedings of ACM Mobicom 2005*, 2005.

[KNB+06]  P. Korteweg, M. Nuyens, R. H. Bisseling, T. J. M. Coenen, H. van den Esker, B. Frenk, R. de Haan, B. Heydenreich, R. van den Hofstad, J. C. H. W. in 't Panhuis, L. Spanjers, and M. van Wieren. Math saves the forest. *Mathematics in Industry; Scientific Proceedings of the 55th European Study Group Mathematics with Industry*, pages 117–140, 2006.

[KP05]    I. Kang and R. Poovendran. Maximizing network lifetime of broadcasting over wireless stationary ad hoc networks. *Mobile Networks and Applications*, 10:879–889, 2005.

[Kra89]   M. Kramer. Stationary distributions in a queueing system with vacation times and limited service. *Queueing Systems*, 4:57–68, 1989.

[KV05]    P. N. Kyasanur and N. H. Vaidya. Multi-channel wireless networks: Capacity and protocols. *Technical Report of University of Illinois at Urbana-Champaign*, 2005.

[L+01]    J. Li et al. Capacity of ad hoc wireless networks. *Proceedings of ACM Mobicom 2001*, pages 61–69, 2001.

[L+03]    R. Litjens et al. Performance analysis of wireless LANs: An integrated packet/flow level approach. *Proceedings of the 18th International Teletraffic Congress 18*, 2003.

[L+04]    R. Litjens et al. Analysis of flow transfer times in IEEE 802.11 wireless LANs. *Annales de télécommunications*, 59:1407–1432, 2004.

[L+05]    E. Lloyd et al. Algorithmic aspects of topology control problems for ad-hoc networks. *Mobile Networks and applications*, 10:19–34, 2005.

[Lee89]   T. T. Lee. M/G/1/N queue with vacation time and limited service discipline. *Performance Evaluation*, 9:180–190, 1989.

[Lev90]   H. Levy. Polling systems: Applications, modeling and optimization. *IEEE Transactions on Communications*, 38, 1990.

[LG05]    C. P. Low and L. W. Goh. On the construction of energy-efficient maximum residual battery capacity broadcast trees in static ad-hoc wireless networks. *Computer Communications*, 29:93–103, 2005.

[Lia02]   W. Liang. Constructing minimum-energy broadcast trees in wireless ad hoc networks. *Proceedings of the International Symposium on Mobile Ad Hoc Networking and Com-puting (MobiHoc)*, pages 112–122, 2002.

[LM14]    G. Latouche and M. Miyazawa. Product-form characterization for a two-dimensional reflecting random walk. *Queueing Systems*, 77:373–391, 2014.

[LR06]     J. S. H. van Leeuwaarden and J. A. C. Resing. A tandem queue with coupled processors: computational issues. *Queueing Systems*, 51:29–52, 2006.

[Miy11]    M. Miyazawa. Light tail asymptotics in multidimensional reflecting processes for queueing networks. *TOP 2011*, 19:233–299, 2011.

[NL02]     P. C. Ng and S. C. Liew. Offered load control in IEEE802.11 multi-hop ad-hoc networks. *ICC 2002*, 2:1074–1079, 2002.

[NnQ00]    R. Núñez Queija. *Processor-sharing models for integrated-services networks*. PhD thesis, 2000.

[Pet06]    A. Petcher. QoS in wireless data networks. www-docs/cse574-06/ftp/wirelessqos/index.html, 2006.

[PS07]     J. Park and S. Sahni. Maximum lifetime broadcasting in wireless networks. *3rd ACS/IEEE International Conference on Computer Systems and Applications*, pages 1–8, 2007.

[RE88]     R. R. Rao and A. Ephredimes. On the stability of interacting queues in a multiple access system. *IEEE Transactions of Informatics Theory*, 43:918–930, 1988.

[RO03]     J. A. C. Resing and L. Örmeci. A tandem queueing model with coupled processors. *Operations Research Letters*, 31:383–389, 2003.

[Ros96]    S. M. Ross. *Stochastic Processes*. Wiley, New York, 1996.

[Sch03]    A. Schrijver. *Combinatorial Optimization:Polyhedra and efficiency*. Springer, Berlin, 2003.

[Ste07]    J. Stemerdink. A simulation approach to a robust multi-radio multi-channel wireless ad hoc network. *Proceedings of the IST Mobile Summit 2007*, 2007.

[TG95]     I. E. Telatar and R. Gallager. Combining queueing theory with information theory and multiaccess. *IEEE Journal on Selected Areas in Communications*, 13:963–969, 1995.

[Toh02]    C. K. Toh. *Ad hoc mobile wireless networks: protocols and systems*. Prentice Hall, 2002.

[vD88]     N. M. van Dijk. Simple bounds for queueing systems with breakdowns. *Performance Evaluation*, 8:117–128, 1988.

[vD11]     N. M. van Dijk. Error bounds and comparison results: The markov reward approach for queueing networks. In R. J. Boucherie and N. M. van Dijk, editors, *Queueing Networks: A Fundamental Approach*, volume 154. Springer, 2011.

[vDP88]    N. M. van Dijk and M.L. Puterman. Perturbation theory for Markov reward processes with applications to queueing systems. *Advances in Applied Probability*, 20:79–98, 1988.

[VLK01]    G. de Veciana, T. J. Lee, and T. Konstantopoulos. Stability and performance analysis of networks supporting elastic services. *IEEE/ACM Transactions on Networking*, 9:2–14, 2001.

[W+02]     H. Wu et al. Performance of reliable transport protocol over IEEE 802.11 wireless LAN: Analysis and enhancement. *Proceedings of IEEE INFOCOM '02*, 2002.

[Wil06]     A. Willig. Wireless sensor networks: concept, challenges and approaches. *Elektrotechnik und Informationstechnik*, 123:224–231, 2006.

[Xia05]     Y. Xiao. Performance analysis of priority schemes for IEEE 802.11 and IEEE 802.11e wireless LANs. *IEEE Transactions on Wireless Communications*, 4, 2005.

[XM06]      L. Xiong and G. Mao. Saturated throughput analysis of IEEE 802.11e using two-dimensional markov chain model. *Proceedings QShine '06*, 2006.

[Yeh02]     E. Yeh. Delay-optimal rate allocation in multiaccess communications: a cross-layer view. *2002 IEEE workshop on Multimedia Signal Processing*, pages 404–407, 2002.

[YH91]      H-C. Yu and R. L. Hamilton. A buffered two node packet radio network with product form solution. *IEEE transactions on communications*, 39:62–75, 1991.

[ZC03]      H. Zhu and I. Chlamtac. An analytic model for IEEE 802.11e EDCF differential services. *Proceedings of IEEE ICCCN '03*, 2003.

# Summary

It's impossible to picture the world without the vast amount of wireless networks used by society nowadays. People are using their smartphones not only to communicate with eachother, but also to control many devices in their home environment, from the heating and lighting in their homes to the volume and program on their TVs. It's not just people with a smartphone or a laptop that are making use of wireless connections. Sensor networks are used for the monitoring of geographical areas, for example for security reasons or environmental data collection. Vehicular networks communicate to warn for upcoming congestions on the road due to accidents and military forces communicate over secure networks that need to be deployed from scratch in warzones. In the latter case, but often also in sensor networks, no centralized control is possible or available. The network that is deployed needs to configure itself for the situation at hand, which is why these networks are denoted as ad hoc networks.

The wireless nature of ad hoc networks poses multiple challenges that need to be faced for the network to operate properly. One of the main challenges is the problem of interference, which is the effect that when multiple signals are received at the same time, they collide and prohibit correct reception of the signals. This thesis focusses on the impact that interference has on ad hoc networks, in particular on the capacity, lifetime and throughput of the network and the congestion and delay in the network.

Chapter 2 analyses the lifetime of a network, which is defined as the time it takes until the battery of the first node is depleted. Two situations are considered: Direct transmissions between the source and destination or full routing where neighbouring nodes relay the traffic for each communication. For these settings the distribution of the network lifetime is determined. The trade-off between the number of transmissions and the distance bridged by each transmission is analysed. The nodes of the network are considered to be on a one dimensional grid or are uniformly distributed. We show that for nodes on a grid it is beneficial to use full routing. For uniformly distributed nodes, the number of nodes in the network determines which approach is better. For small networks, direct transmission outperforms the full routing approach. In this case, the longer distance that needs to be bridged weighs up against the increased number of transmissions that are needed. An intermediate approach, choosing master nodes that forward data to other master nodes is simulated. Models for the expected lifetime are provided that give approximations which are close to the simulated results.

Chapter 3 models the delay in a wireless ad hoc network using a polling model to take into account QoS differentiation in ad hoc networks. Traffic can have either high or low priority, determining the probability that a node is serving a

packet. The delay experienced by packets of each class is analysed by considering
each queue separately as being served by a server that takes holidays. The length
of these holidays depends on the state of the system, making it hard to analyze
them. An iteration algorithm, which is proven to monotonically converge, is
presented to compute the waiting time distribution of a queue that uses the
steady state for all other queues. Iterating over all queues provides de delay for
packets at all queues, which gives accurate results for low to moderately loaded
networks.

Chapter 4 combines results on product-form networks with a Markov reward
approach to find bounds on any performance measure that is linear in each of
the components of the state space. A two node network is considered where
traffic can be forwarded from the first to the second node. When both nodes are
active, the interference causes a lower service rate than when only one node is
active. The stability range of the system is analysed, showing that increasing
the rate at the boundaries of the system expands the stability range. Conditions
for a geometric product-from solution are given which are used for comparison
with the network under consideration. The Markov reward approach provides
bounds for several performance measures, where we show that the choice of the
product-form network to compare with obtains different bounds.

Chapter 5 analyses whether a network with a given traffic demand, capacities
on each links and ranges of interference between the nodes can accommodate all
the traffic demand. In the first part only one channel is available, so interference
plays a large role in determining the throughput of the network. The network is
modeled using a multi commodity flow problem and a theorem is stated that
gives sufficient and necessary conditions for the problem to be solvable. For a
single source and destination pair the maximal throughput is computed. The
second part expands the results of the first part by including the option of using
different channels. The theorem is expanded to include these channels, giving a
basis for an algorithm for channel allocation in wireless networks.

Chapter 6 considers the throughput of ad hoc networks, taking into account
the parameters involved in the CSMA/CA protocol with RTS-CTS in a wireless
network. Setting different priorities to flows over a different number of hops takes
into account these parameters. First, considering the packet level details, the
aggregate system throughput is determined. Next, taking the flow level dynamics
into account, the throughput is divided over all flows, taking into account the
impact of multiple hops used in flows. This leads to two Processor Sharing
models: Batch arrival processor sharing (BPS) and Discriminatory processor
sharing (DPS). Simulation shows that the models provide an accurate estimation
of the throughput for small networks.

Chapter 7 considers the impact on the throughput of the contention that
happens between nodes in an ad hoc network. During each time slot the nodes of
the network contend for the channel, depending on the protocol in use. Starting
with a discrete time Markov chain we model the behaviour in the slotted time.
Using long term average behaviour, we then model this discrete time Markov
chain as a continuous time Markov chain, taking into account that certain nodes
may be bottleneck nodes. The transition rates in this chain are state dependent
so that further approximation is needed to obtain results on the throughput of

the network. Now considering a long term average service rate, we approximate the continuous time Markov chain by a product-form network. This enables us to find the bottlenecks for a wireless network of any size and topology and to approximate its throughput. As the main result, an algorithm is provided that incorporates all these steps and gives very accurate results for the maximal throughput of the network. For a multihop tandem network a limiting result is obtained for the rate allocated to the first nodes of a fully alive network. Additionally, a surprising effect is observed where the location of the bottleneck changes as the load of the network increases, which is correctly predicted by the algorithm.

# Samenvatting

Het is onmogelijk om de wereld voor te stellen zonder het enorma aantal draadloze netwerken dat hedendaags gebruikt wordt. Mensen gebruiken hun smartphone niet alleen om met elkaar te communiceren, maar ook om meerdere apparaten te besturen, van de verwarming en belichting in hun huis tot het volume en het kanaal op de TV. Het zijn niet alleen mensen met hun smartphone of laptop die gebruik maken van draadloze verbindingen. Sensor netwerken worden gebruikt voor het monitoren van geografische gebieden, bijvoorbeeld voor beveiliging of het verzamelen van ecologische data. In voertuignetwerken wordt gecommuniceerd om waarschuwingen af te geven voor aankomende files door ongelukken en militaire troepen communiceren over beveiligde netwerken die ter plekke van de grond af opgebouwd moeten worden. In het laatste geval, maar ook in sensornetwerken, is geen gecentralizeerde controle mogelijk of beschikbaar. Het netwerk dat wordt opgezet moet zichzelf configureren voor de situatie die voorhanden is, waardoor deze netwerken als ad hoc netwerken bestempeld worden.

De draadloze aard van ad hoc netwerken brengt een aantal uitdagingen met zich mee die aangepakt moeten worden om het netwerk goed te laten opereren. Een van de grootste uitdagingen is het probleem van interferentie, het effect dat wanneer meerdere signalen tegelijkertijd ontvangen worden, er een botsing ontstaat en een correcte ontvangst van de signalen verstoord wordt. Dit proefschrift focust op de invloed die interferentie heeft op ad hoc netwerken, in het bijzonder op de capaciteit en de levensduur van het netwerk en de doorstroom, ophoping en vertraging in het netwerk.

Hoofdstuk 2 analyseert de levensduur van een ad hoc netwerk, wat gedefinieert wordt als de tijd die nodig is totdat de eerste batterij in het netwerk leeg is. Twee verschillende situaties worden bekeken: directe verzending tussen de bron en bestemming en volledige routering waarbij aangrenzende knopen data doorsturen voor elke communicatie. Voor deze configuraties wordt de verdeling van de netwerklevensduur bepaald. De wisselwerking tussen het aantal transmissies en de afstand die overbrugd wordt per transmissie wordt geanalyseerd. De knopen van het netwerk liggen op een één-dimensionale grid of zijn uniform verdeeld. We laten zien dat voor knopen op een grid het gunstiger is om volledige routering te gebruiken. Voor uniform verdeelde knopen is het aantal knopen in het netwerk van invloed op de optimale keuze. Voor kleine netwerken werkt directe transmissie beter dan volledige routering. In dit geval weegt de grotere afstand die overbrugd moet worden op tegen de toename van het aantal transmissies dat nodig is. Een tussenliggende aanpak, waarbij hoofdknopen aangewezen worden die data doorsturen via de andere hoofdknopen wordt gesimuleerd. Modellen voor de verwachte levensduur worden gegeven en de benaderingen die ze opleveren liggen dicht bij de gesimuleerde waarden.

Hoofdstuk 3 modelleert de vertraging in een draadloos ad hoc netwerk door gebruik te maken van een polling model. Hierbij wordt het onderscheid in QoS tussen de verschillende knopen in het netwerk meegenomen. Verkeer in het netwerk kan een hoge of een lage prioriteit krijgen, wat bepaalt wat de kans is dat een knoop een pakket mag verzenden. De vertraging die ervaren wordt door pakketten van elke klasse wordt geanalyseerd door elke wachtrij apart te zien alsof die geholpen wordt door een server die vakanties neemt. De lente van de vakanties hangt af van de staat waarin het systeem zich bevindt, wat het moeilijk maakt om dit model te analyseren. Door een iteratief algoritme, waarvan bewezen wordt dat deze monotoon convergeert, wordt gepresenteerd om de verdeling van de wachttijd in een wachtrij te berekenen. Hierbij wordt gebruik gemaakt van de evenwichtsverdeling van de andere wachtrijen. Door te itereren over alle wachtrijen wordt de vertraging voor de pakketten van iedere wachtrij bepaald, wat een nauwkeurig resultaat oplevert voor wachtrijen met een lichte tot matige belasting.

Hoofdstuk 4 combineert resultaten over productvorm netwerken met een Markov beloning aanpak om grenzen te vinden voor de prestatie van een ad hoc netwerk. Een netwerk bestaande uit twee knopen wordt bekeken waarbij paketten doorgestuurd kunnen worden van de eerste naar de tweede knoop. Wanneer beide knopen tegelijk actief zijn, zorgt de interferentie voor een lager tempo waarin de pakketten verzonden kunnen worden dan wanneer er slechts één actief is. Het domein waarbinnen het netwerk stabiel is wordt bepaald, wat laat zien dat bij een hoger tempo van de knopen als ze alleen actief zijn het domein vergroot. Voorwaarden waarbinnen het netwerk een geometrische productvorm evenwichtsverdeling heeft worden gegeven en worden gebruikt om te vergelijken met het netwerk dat onderzocht wordt. De Markov beloning aanpak geeft grenzen aan voor verschillende prestatiematen, waarbij we laten zien dat de keuze van het productvorm netwerk waarmee vergeleken wordt verschillende grenzen oplevert.

Hoofdstuk 5 analyseert of een netwerk met een gegeven vraag, capaciteit van elke link en het bereik van interferentie tussen de knooppunten aan de complete vraag kan voldoen. In het eerste deel wordt een netwerk met één kanaal bekeken, zodat interferentie een grote rol speelt op de doorstroom van het netwerk. Het netwerk wordt gemodelleerd door een multi commodity flow problem en een stelling wordt geformuleerd die voldoende en noodzakelijke voorwaarden geeft waaronder het probleem oplosbaar is. Voor een enkele bron en bestemming wordt de maximale doorvoer berekend. Het tweede deel van het hoofdstuk breidt het resultaat uit van het eerste deel door meerdere kanalen waarover data kan worden verzonden mee te nemen. De stelling wordt uitgebreid om deze kanalen mee te nemen, wat een basis verschaft voor een algoritme om kanalen toe te kennen in een draadloos ad hoc netwerk.

Hoofdstuk 6 bekijkt de doorstroom in ad hoc netwerken, rekening houdend met de parameters van het CSMA/CA protocol met RTS-CTS in een ad hoc netwerk. Door verschillende prioriteiten in te stellen voor transmissies over meerdere hops worden de parameters gemodelleerd. In eerste instantie wordt het netwerk op packet niveau bekeken en wordt de gemiddelde doorstroom van het netwerk bepaald. Hierna wordt het netwerk op flow niveau bekeken en wordt de

doorstroom verdeeld over alle flows, waarbij de impact van de meerdere hops meegenomen wordt. De leidt tot twee Processor Sharing modellen: Batch arrival processor sharing (BPS) en Discriminatory processor sharing (DPS). Simulatie laat zien dat de modellen een nauwkeurige schatting van de doorstroom opleveren voor kleine netwerken.

Hoofdstuk 7 bekijkt de invloed op de doorstroom van de strijd die tussen de knopen van het ad hoc netwerk plaatsvindt om het recht te krijgen een signaal te verzenden. Tijdens elk tijdslot strijden de knopen om het kanaal te mogen gebruiken, afhankelijk van het protocol dat gebruikt wordt. Beginnend met een Markov keten in discrete tijd modelleren we het gedrag in de tijd die opgedeeld wordt in tijdsloten. Door gebruik te maken van lange termijn gemiddeld gedrag zetten we het discrete model om in een Markov keten met continue tijd. Hierbij wordt er rekening mee gehouden dat sommige knopen als knelpunt op kunnen treden. De overgangssnelheden in deze keten zijn afhankelijk van de staat van het systeem zodat verdere benaderingen nodig zijn om resultaten over de doorstroom te verkrijgen. Door de lange termijn gemiddelde behandeltijd te gebruiken benaderen we de Markov keten in continue tijd als een productvorm netwerk. Dit brengt ons in staat om de knelpunten van een ad hoc netwerk van willekeurige grootte en indeling te vinden de doorstroom te benaderen. As hoofdresultaat wordt een algoritme gegeven dat al de benaderingsstappen bevat en zeer nauwkeurige resultaten geeft over de maximale doorstroom dat een netwerk aankan. Voor een mutlihop tandem netwerk is een limiet resultaat gevonden voor de capaciteit van het netwerk dat aan de eerste knopen toegekend wordt in een volledig actief netwerk. Verder wordt een verrassend effect, waar de locatie van een knelpunt verandert door de toename van de drukte op het netwerk, correct voorspeld door het algoritme.