

Agent-Based Routing in Queueing Systems

THÈSE N° 4598 (2010)

PRÉSENTÉE LE 26 FÉVRIER 2010

À LA FACULTÉ SCIENCES ET TECHNIQUES DE L'INGÉNIEUR

LABORATOIRE DE PRODUCTION MICROTECHNIQUE 1

PROGRAMME DOCTORAL EN SYSTÈMES DE PRODUCTION ET ROBOTIQUE

ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

POUR L'OBTENTION DU GRADE DE DOCTEUR ÈS SCIENCES

PAR

Olivier GALLAY

acceptée sur proposition du jury:

Prof. A. Billard, présidente du jury
Prof. M.-O. Hongler, directeur de thèse
Prof. D. Armbruster, rapporteur
Prof. C. Pfister, rapporteur
Prof. A. Van Ackere, rapporteur



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Suisse
2010

AGENT-BASED ROUTING IN QUEUEING SYSTEMS

LAUSANNE, EPFL
2010

IN PART SUPPORTED BY THE "FONDS NATIONAL SUISSE POUR LA
RECHERCHE SCIENTIFIQUE"

To Isabelle and to my parents

Preface

Before diving into the main core of this thesis, it is useful to understand in which context and under which circumstances this work has been fulfilled. The present book is the outcome of a research activity led in a micro-engineering laboratory, under the supervision of a theoretical physicist and worked out by an engineer in communication systems with a specialization in probabilities. It is therefore not surprising that, in the end, this thesis intrinsically exhibits a strong interdisciplinary flavour. I hope that this work reflects some virtues of the people I have been in close contact with during the last four years, namely the curiosity of theoretical physicists, the pragmatism of engineers and the formalism of mathematicians.

This thesis deals with the decentralization of routing mechanisms in queueing networks, with a view on potential applications in service, manufacturing and supply systems. In this context, two antagonistic approaches could have been followed, either to develop stylized, analytically solvable models or to construct idiosyncratic, most often heavy, simulation frameworks. Due to the complete lack of available theory concerning the main topic considered in this work, we have chosen to adopt the first type of methodology. Nevertheless, whatever the assumptions we have taken within our models, they always have been clearly enunciated such that our modelling choices remain transparent and fully open to criticism. Due to their intrinsic computational tractability, our collection of idealized models constitutes a perfect starting point for understanding and, ultimately, possibly improving the complex processes under investigation. It is likely that practitioners will find our models in a sense oversimplified, but they will however be able to extract useful information from the analytical descriptions provided in this thesis about various collective phenomena emerging in such context.

The introduction of history-based routing mechanisms in queueing networks, as considered throughout this thesis, is a new topic in itself. For this reason, lots of attention has been paid to emphasize the originality of our work as well

as to confine it clearly by drawing some relationships with several potentially related domains. We have furthermore, all along this work, adopted an approach that is as didactic as possible. We hope that this choice of presentation will help non-specialists to go through the arguments, but will nevertheless let specialists appreciate new advances in their field. In this regard, the stylized, solvable models proposed in this thesis allow for synthetic, mostly analytical resolutions and they are thus well-tailored to be possibly taught to students in master or post-graduate courses.

Lausanne, November 2009

Olivier Gallay

Summary

Waiting time in any network is often a costly and hence a bad experience. Therefore, to avoid jamming regions becomes essential in the optimization of traffic flows. In this regard, the conception and the control of complex networks supporting flows of units are key issues in various strategic engineering and service areas ranging from manufacturing systems, supply chains, retail stores, transportation and communication networks to only quote a few. Flow dynamics depend jointly on the routing rules defining the ways the units are dispatched at the network vertexes and on the behaviour of the servers which process these items. Due to stochastic customer demands, to fluctuations in the raw material of supply chains, to failures arising in production devices, to uncertainties in operator availability and to ubiquitous financial volatility steadily affecting optimization objectives, the flow dynamics are always affected by random fluctuations. The need to model, to study and to quantify the characteristics of such complex stochastic dynamics has strongly stimulated the development of the so-called queueing network (QN) theory. Under fairly general hypothesis (including the possibility to describe the underlying dynamics by general Markov processes), powerful methods are available to calculate the time-invariant probability densities ultimately describing the system state and therefore to obtain useful quantitative information on the corresponding stationary regimes. However, the Markovian character imposed to the dynamics obviously limits not only the behaviour of the servers but also lays down strong restrictions on the allowable routing rules followed by the circulating items.

While the classical QN theory assumes that the circulating items composing the flows are mainly passive entities, it is however mandatory in numerous applications that each circulating item possesses its own identity as well as the ability to take individual decisions. Indeed, the inherent complexity characterizing nowadays production and/or service networks strongly favours decentralized and self-organizing mechanisms to regulate the circulating flows of humans, matter and/or information. This clearly shows the importance

to study the flow dynamics of QNs visited by autonomous travelling agents which select their routing according to individual historical data collected during their past progression in the network. When such kind of history-based (HB) routing mechanisms is taken into account, the resulting flow dynamics are intrinsically non-Markovian and hence one of the fundamental assumption of classical QN theory is explicitly violated. In particular, the existence of stationary regimes can not be anymore guaranteed. As it will be unveiled in this thesis, the joint action of HB features and of feedback loop topologies in the QNs opens wide the door for the emergence of entirely new dynamical features. The decentralized autonomous dispatching rules considered in the present work require the capability for each circulating entity to monitor, to memorize and to process data. As a consequence, we actually consider agents that autonomously adapt to environmental changes and that interact to produce intelligent global behaviours. As it is typical for self-organized systems, the numerous stigmergic interactions between the agents and their environment open possibilities for the emergence of various collective structures. In this context, we show in the present thesis how QNs in which agents are allowed to modify autonomously their routing strategies according to measures of *(i)* their personal waiting times and/or *(ii)* real-time observations of the queue contents might give rise to the creation of dissipative collective spatio-temporal patterns. By restricting our analysis to simple network topologies and despite the intrinsically non-Markovian character of the dynamics, we are able to explicitly provide analytical results characterizing the macroscopic structures that might appear in such stylized multi-agent QNs. Since the different emerging collective patterns that will be described in this work are solely induced by the agents' local actions and interactions, our particular QNs reveal themselves to be specific instances of complex adaptive systems (CAS).

For the purpose of modelling recurrent services, we first consider the dynamics of a single server feedback queueing system where the agents' autonomous decision to take the feedback loop depends on their individual experimented waiting time. This is an idealization of the dynamics induced by a collection of customers who remain loyal to a server provided their service time stays below a critical threshold. The emerging self-organized behaviour takes the form of a cyclo-stationary regime which exhibits a periodical purging of the queue content. This ensures a bounded queue length as well as a maximum server utilization. Further, we increase the complexity of the network and consider a topology with two parallel feedback queues in which the agents, besides their ability to monitor waiting times, also possess a vision capability (hence, the agents "intelligence" is enhanced). Concerning the agents' initial choice between the two service providers, we introduce several different routing policies. These autonomous rules, which are wholly based on specific agent capabilities, differ in their level of complexity and with respect to the available amount of information about the current system state. In this context,

we are able to analytically characterize nonlinear collective behaviours such as synchronization of oscillations and noise-induced stabilization.

Next, we consider the market partition dynamics between two recurrent service providers in a closed topology, where the customers' satisfaction depends as before in a nonlinear fashion on the elapsed waiting time to receive service. We describe the periodic cannibalization effects that might emerge in this system. Further, we introduce spatial considerations within the dynamics of a system composed of two servers competing for a stationary incoming flow of customers. More particularly, we study explicitly the market sharing dynamics between these two service providers when customers are not only sensitive to costs related to waiting time but also to transportation costs. In this context, we fully characterize the phase transition that occurs between regimes where waiting time considerations are predominant and regimes where transportation aspects dominate.

While the multi-agent type of dynamics considered in this work is in essence inspired by human behaviour, the present modelling framework is by far not only restricted to social systems. Indeed, to attach RFID or "smart tags" on any type of circulating units (raw materials, finished goods, carrying pallets,...) allows for the implementation of such type of decentralized dynamics in logistics, supply or manufacturing networks dealing with highly customized products. To illustrate this assertion, we propose a new dynamic fully decentralized load sharing policy (LSP) which optimally dispatches the incoming workload according to the current availability of a set of operators. The underlying decentralized decisions rely on a "smart tasks" paradigm in which each incoming task is endowed with specific autonomous routing decision mechanisms. By optimal LSP, we mean here that our "smart tasks" based algorithm permanently requires the engagement of a minimum number of operators while still respecting due dates. The numerous stigmergic interactions between these autonomous tasks solely cause the optimal LSP to emerge.

Keywords: Queueing Systems, Multi-Agent Dynamics, Autonomous History-Based Routing, Nonlinear Delayed Feedback, Spatio-Temporal Patterns, Self-Organization, Complex Adaptive Systems, Production and Service Networks.

Résumé

Dans tout système industriel ou de service, attendre représente souvent pour l'utilisateur une désagréable, voire une coûteuse expérience. Par conséquent, il devient alors essentiel, lorsque l'on cherche à optimiser les flux de trafic dans de tels réseaux, de chercher à éviter la formation de régions congestionnées. Pour ces raisons, la conception et le contrôle de réseaux complexes, impliquant des flux d'éléments déterminés, donne naissance de nos jours à des problèmes incontournables dans des domaines stratégiques variés tels que la gestion de chaînes logistiques, le commerce de détail ou encore les réseaux de communication et de transport. La dynamique de tels flux dépend non seulement du comportement des serveurs traitant les différents éléments en circulation, mais également des règles de routage définissant le parcours de ces unités au sein du réseau. L'aspect souvent stochastique de la demande, la variabilité affectant la disponibilité de la matière première au sein des chaînes logistiques ainsi que les incessantes fluctuations économiques impliquent que les dynamiques de flux à considérer sont clairement de nature aléatoire. Dès lors, le besoin toujours plus grand de pouvoir modéliser, étudier et quantifier le comportement de ces flux stochastiques a hautement stimulé au fil des ans le développement d'une théorie spécifique dédiée à l'étude des réseaux de files d'attente. Cette théorie offre, sous des hypothèses relativement générales (qui incluent notamment que la dynamique puisse être décrite à l'aide de processus de Markov généraux), des méthodes puissantes dans le but de calculer les densités de probabilité, invariantes dans le temps, décrivant de manière ultime l'état du système. Ces méthodes classiques permettent par conséquent d'obtenir des informations importantes concernant les régimes stationnaires correspondants. Toutefois, il est important de remarquer que la nature markovienne que l'on se doit d'imposer à la dynamique induit des limites non seulement sur la marche des serveurs, mais cela impose également des restrictions fortes sur les règles de routage possibles pour les éléments en circulation.

La théorie des réseaux de files d'attente classique suppose que les flux d'éléments en circulation sont composés d'entités essentiellement passives. Il

est cependant nécessaire, pour de nombreuses applications, que les unités en circulation possèdent une identité propre et qu'elles aient ainsi la capacité de pouvoir prendre des décisions de manière individuelle. En effet, de nos jours, l'inhérente complexité qui caractérise généralement les réseaux de production et de service parle en faveur de l'implémentation de mécanismes décentralisés et auto-organisés pour la régulation de flux humains, de matière ou d'information. Par conséquent, on comprend aisément l'importance que revêt l'étude dédiée à la dynamique de flux dans des réseaux de files d'attente visités par des agents autonomes capables de sélectionner leur chemin en fonction de données historiques, collectées individuellement lors de leur évolution passée au sein du système. Lorsque l'on considère ce type de mécanismes de routage, basés sur l'historique, il s'ensuit inévitablement que la dynamique de flux résultante est intrinsèquement non-markovienne et, par conséquent, que l'une des hypothèses fondamentales de la théorie des réseaux de files d'attente classique est explicitement violée. En particulier, il n'est plus possible de garantir l'existence de régimes stationnaires. A cet égard, il sera montré dans cette thèse que la prise en compte de données historiques dans les décisions de routage, associée à la présence de non-linéarités topologiques, permet l'émergence de dynamiques totalement nouvelles, ceci même dans des réseaux de files d'attente simples. Les règles de routage décentralisées considérées dans ce travail impliquent que chaque entité en circulation soit capable d'acquiescer, de mémoriser ainsi que de traiter certaines données. Nous serons donc en présence d'agents aptes à s'adapter de manière autonome aux changements de l'environnement dans lequel ils évoluent. De plus, ces agents interagissent de sorte à produire des comportements collectifs intelligents. En effet, il est générique, pour tout type de système auto-organisé, que les nombreuses interactions stochastiques existant entre les agents et leur environnement donnent naissance à des structures macroscopiques variées et potentiellement complexes. Dans ce contexte, nous exhiberons dans cette thèse comment des réseaux de files d'attente dans lesquels les agents en circulation sont capables de modifier de manière autonome leur stratégie de routage en fonction (i) de mesures de leurs temps d'attente personnels et/ou (ii) d'observations en temps réel de longueurs de file d'attente peut provoquer l'émergence, au niveau collectif, de structures spatio-temporelles dissipatives. En restreignant notre analyse à des topologies de réseau simples, nous fournirons explicitement des résultats analytiques pour caractériser les phénomènes macroscopiques pouvant apparaître dans de tels systèmes d'attente multi-agents, ceci malgré le caractère intrinsèquement non-markovien des dynamiques considérées. Ces différentes structures collectives étant uniquement dues aux actions ainsi qu'aux interactions des agents en circulation, nos modèles de réseaux de files d'attente idéalisés se révèlent appartenir de façon manifeste à la classe des systèmes complexes adaptatifs (*Complex Adaptive Systems*).

Dans le but de modéliser un service récurrent, nous considérons premièrement

le comportement d'un unique serveur avec possibilité de retour, où les agents autonomes basent leur décision d'emprunter ou non la boucle de retour sur le temps d'attente qu'ils ont personnellement expérimenté pour recevoir ledit service. Ce modèle représente une idéalisation de la dynamique induite par une population de clients restant fidèles à un fournisseur pour autant que leur temps de service expérimenté soit inférieur à un seuil critique. On observe pour ce système l'émergence d'un régime cyclo-stationnaire et plus particulièrement d'une purge périodique du contenu de la queue. Cette structure auto-organisée assure une longueur de file d'attente bornée ainsi qu'une utilisation maximale du serveur. Nous augmentons par la suite la complexité de la topologie étudiée et nous nous intéressons à un réseau ouvert impliquant deux serveurs avec boucle de retour placés en parallèle. Dans ce nouveau modèle, "l'intelligence" des agents est accrue et ils possèdent à présent, en plus de leur capacité à mesurer des temps d'attente, un système de vision leur permettant d'observer le contenu des files d'attente qu'ils rencontrent lors de leur parcours. Nous introduisons différentes règles de décision concernant le choix initial des agents entre les deux serveurs disponibles. Ces mécanismes décentralisés, qui font uniquement appel aux capacités spécifiques des agents, diffèrent dans leur complexité ainsi que dans la quantité d'information à disposition sur l'état actuel du système. Nous donnons une caractérisation analytique des structures collectives non-linéaires qui émergent dans ce contexte, à savoir des phénomènes de synchronisation d'oscillateurs ainsi que de stabilisation par le bruit.

Nous étudions ensuite la dynamique régissant la répartition de marché entre deux fournisseurs de service sis dans une topologie fermée. Comme précédemment, la satisfaction des clients, et par conséquent leur fidélité à un fournisseur de service donné, dépend à nouveau du temps d'attente expérimenté pour recevoir le service concerné, cette dépendance étant supposée non-linéaire. Nous décrivons explicitement les effets de cannibalisation qui émergent de façon périodique dans un tel système. Plus loin, nous introduisons des considérations spatiales dans l'étude de la dynamique d'un système ouvert composé de deux serveurs en compétition pour un unique flux de clients entrants. Plus particulièrement, nous étudions analytiquement l'évolution temporelle de la répartition de marché entre ces deux fournisseurs de service lorsque les clients sont non seulement sensibles aux temps d'attente, mais également aux coûts de transport. Dans ce contexte, nous caractérisons entièrement la transition de phase qui se produit entre des régimes où les considérations de temps sont prédominantes et des régimes pour lesquels les aspects spatiaux sont, au contraire, prépondérants.

Le type de dynamique multi-agents considéré tout au long de ce travail est clairement inspiré du comportement humain, mais il n'est de loin pas limité aux réseaux sociaux. En effet, il est possible d'attacher des puces RFID (*i.e.* des étiquettes "intelligentes") sur n'importe quelle sorte d'unités en circu-

lation (pièces détachées, produits finis, palettes de transport,...) et de permettre ainsi l'implémentation de mécanismes décentralisés dans des réseaux logistiques, d'approvisionnement ou encore de production. Une telle gestion décentralisée permet notamment de traiter de manière efficace les processus de production de produits possédant un nombre important de variantes, car étant hautement personnalisables. Nous proposons un modèle qui illustre parfaitement l'efficacité que peut revêtir l'implémentation de telles méthodes multi-agents dans des réseaux de production. Plus précisément, nous introduisons un nouveau mécanisme optimal de répartition de charge, dynamique et décentralisé, qui permet de distribuer une charge de travail entrante en fonction de la disponibilité actuelle d'un ensemble d'opérateurs. Par optimalité, nous entendons ici que notre algorithme assure de façon permanente que le nombre d'opérateurs engagés est minimal, tout en garantissant le respect de dates d'échéance données. Dans notre modèle, chaque tâche entrant dans le système est dotée de mécanismes lui permettant de prendre des décisions autonomes spécifiques. Les nombreuses interactions stigmergiques entre ces tâches autonomes provoquent à elles seules l'émergence d'une répartition optimale de la charge entrante entre les différents opérateurs, ceci de manière permanente et avec haute reactivité.

Mots clés: réseaux de files d'attente, dynamique multi-agents, routage autonome basé sur l'historique, réactions non-linéaires et différées, structures spatio-temporelles, auto-organisation, systèmes complexes adaptifs, réseaux de production et de service.

Zusammenfassung

Die Wartezeit in einem Netzwerk ist oftmals kostenintensiv und daher eine schlechte Erfahrung. Deshalb ist es für die Optimierung von Verkehrsströmen grundlegend, blockierte Regionen zu vermeiden. In dieser Hinsicht unterstützen die Konzeption und Kontrolle von komplexen Netzwerken im Kern die Ströme von Einheiten in verschiedenen strategischen Ingenieur- und Dienstleistungsgebieten, die von Fertigungssystemen, zu Liefer- und Versorgungsketten, Warengeschäften, Transport- und Kommunikationsnetzwerken reichen, um nur einige Beispiele zu nennen. Die Strömungsdynamik hängt gemeinsam von den Regeln der Routenwahl, die den Weg der Einheiten im Netzwerk abwickeln und von der Serverdynamik ab, die das Voranschreiten der Einheiten antreibt. Infolge stochastischer Schwankungen der Verbrauchsnachfrage, Fluktuationen der Ausgangsmaterialien in der Lieferkette, anfallenden Ausfällen in den Produktionselementen, Unsicherheiten der Verfügbarkeit der Anlagen und allgegenwärtigen finanziellen Preisschwankungen, die stetig auf die zu optimierenden Ziele einwirken, ist die Strömungsdynamik immer von zufälligen Fluktuationen beeinflusst. Die Erfordernis, die Charakteristiken einer solchen komplexen stochastischen Dynamik zu modellieren, zu studieren und zu quantifizieren, hat stark die Entwicklung der Queueing Network Theory vorangetrieben. Nach einer ziemlich generellen Hypothese (die die Möglichkeit beinhaltet, die zugrundeliegende Dynamik durch generelle Markov-Prozesse zu beschreiben), sind leistungsfähige Methoden verfügbar, um Zeit-invariante Wahrscheinlichkeitsdichten abzuleiten, welche letztlich den Systemzustand beschreiben und daher nützliche Informationen der korrespondierenden stationären Regime liefern. Allerdings begrenzt eine Anwendung der Markov-Eigenschaft auf die Dynamik offensichtlich nicht nur das Verhalten der Zusteller, sondern legt auch der erlaubten Auswahl von Routen starke Einschränkungen auf, nach denen die Güter umlaufen können.

Während die klassische Queueing Network Theory annimmt, dass sich im Umlauf befindliche Güter aus Strömen zusammensetzen und hauptsächlich passive Größen sind, ist es jedoch zwangsläufig in zahlreichen Anwendungen

so, dass jede im Umlauf befindliche Einheit sowohl seine eigene Identität als auch die Fähigkeit besitzt, individuelle Entscheidungen zu treffen. Tatsächlich bevorzugt die heutzutage inhärente und/oder Produktions- und Service-Netzwerke charakterisierende Komplexität eindeutig dezentralisierte und selbstorganisatorische Mechanismen, um die im Umlauf befindlichen Einheiten von Mensch, Materie und/oder Energie zu regulieren. Dies begründet deutlich die Bedeutung des Studiums der Strömungsdynamik von Queueing Networks, die von autonom reisenden Agenten bevölkert werden, die ihre Route entsprechend ihrer individuellen Vorgeschichte wählen, welche sie durch ihren vergangenen Verlauf im Netzwerk gesammelt haben. Wenn solche Art der von der Vorgeschichte abhängigen Routenmechanismen in Betracht gezogen werden, ist die sich ergebende Strömungsdynamik an sich Nicht-Markovisch, wodurch eine der fundamentalen Annahmen der klassischen Queueing Network Theory ausdrücklich verletzt ist. Besonders die Existenz von stationären Regimen kann nicht mehr länger garantiert werden. Wie in der Thesis enthüllt werden wird, öffnen die gemeinsamen Handlungen der geschichtsabhängigen Eigenschaften und die Strukturen der Rückkopplungsschleifen in den Queueing Networks das Tor für das Auftreten von vollkommen neuen dynamischen Eigenschaften. Die dezentralisierten autonomen Regeln der Verteilung, die in der vorliegenden Arbeit betrachtet werden, erfordern die Fähigkeiten des umlaufenden Agenten, zu beobachten, zu erinnern und Daten zu verarbeiten. Als direkte Konsequenz beobachten wir tatsächlich Agenten, die sich autonom auf Veränderungen der Umwelt anpassen und die interagieren, um intelligentes Verhalten auf globaler Ebene zu produzieren. Wie es für selbstorganisierende Systeme typisch ist, eröffnen die zahlreichen stigmergischen Interaktionen zwischen den Agenten und ihrer Umgebung Möglichkeiten für das Aufkommen verschiedener kollektiver Muster. In diesem Zusammenhang zeigen wir in dieser Arbeit wie Queueing Networks, in denen den Agenten erlaubt ist, selbstständig ihre Routenstrategie entsprechend der Einschätzung der *(i)* persönlichen Wartezeit und/oder *(ii)* der Echtzeit der Beobachtungen der Warteschlangendauer zu modifizieren, Aufschluss über das Aufkommen von dissipativen, kollektiven und räumlich-zeitlichen Mustern geben können. Durch die Einschränkung unserer Analyse auf einfache Netzwerk-Topologien und trotz des intrinsischen Nicht-Markovischen Charakters der Dynamik, sind wir ausdrücklich in der Lage, analytische Ergebnisse bereitzustellen, die die makroskopische Struktur in solch einem stilisierten Multi-Agenten Queueing Network auftreten lässt. Zumal die unterschiedlichen auftretenden kollektiven Muster, die in dieser Arbeit beschrieben werden, ausschließlich durch die Aktionen der Agenten und Interaktionen bedingt sind, und unser Queueing Networks selbst spezifische Belegstellen von komplexen adaptiven Systemen (*Complex Adaptive Systems*) aufweist.

Um rekurrente Dienstleistungen zu modellieren, betrachten wir zuerst die Dynamik eines einzelnen Server Feedback Queueing Systems, bei dem die autonomen Entscheidungen der Agenten, eine Rückkopplungsschleife zu be-

nutzen, von der individuell erlittenen Wartezeit abhängt. Dies stellt eine Idealisierung der Dynamik dar, die nur die loyal zu einem Server verbleibenden Verbraucher erfasst, so dass ihre Servicezeit unterhalb einer kritischen Schwelle bleibt. Das Auftreten einer selbstorganisierenden Dynamik nimmt die Gestalt eines wechselstationären Regimes an, welches eine periodische Durchspülung der in der Warteschlange befindlichen Einheiten bewirkt. Dies sichert eine begrenzte Dauer der Warteschlange sowie eine maximale Nützlichkeit des Servers. Desweiteren steigern wir die Komplexität des Netzwerks und betrachten eine Topologie mit zwei parallelen Queueing Feedback Loops, in denen die Agenten, neben ihrer Möglichkeit, ihre Wartezeit im Auge zu behalten, ebenfalls eine Fähigkeit besitzen, in die Zukunft zu blicken, so dass von hier an die Intelligenz der Agenten erweitert ist. In Bezug auf die ursprüngliche Wahl der Agenten zwischen zwei Service-Anbietern, führen wir einige unterschiedliche Routenregeln ein. Diese autonomen Regeln, welche komplett auf den spezifischen Eigenschaften der Agenten basieren, unterscheiden sich in ihrem Level der Komplexität und bezüglich der verfügbaren Menge an Informationen über den laufenden Systemzustand. In diesem Zusammenhang sind wir in der Lage, analytisch das nicht-lineare, kollektive Verhalten als eine Synchronisation von Oszillationen und geräuschinduzierte Stabilisierungseffekte zu charakterisieren.

Als nächstes wird die Dynamik des Marktes zwischen zwei periodisch auftretenden Service-Anbietern in einer geschlossenen Struktur, in der die Zufriedenheit der Verbraucher wie schon zuvor auf nicht-lineare Art und Weise von der bereits verstrichenen Wartezeit für den Erhalt des Service abhängt, aufgeteilt. Wir beschreiben die periodischen Kannibalisierungseffekte, die in diesem System entstehen können. Weiterhin führen wir räumliche Betrachtungen innerhalb der Dynamik eines Systems ein, das aus zwei Servern besteht, die für einen stationären eingehenden Fluss von Verbraucher wetteifern. Noch genauer untersuchen wir ausdrücklich die gemeinsame Benutzung der Marktdynamik zwischen den beiden Dienstleistungsanbietern, wenn die Verbraucher nicht nur empfindlich auf Kosten, die mit der Wartezeit verknüpft sind, sondern auch auf Transportkosten reagieren. In diesem Zusammenhang beschreiben wir komplett den Phasenübergang, der zwischen zwei Regimen auftritt, bei denen zum einen Betrachtungen zur Wartezeit vorherrschend sind und Regime zum anderen, bei denen die Transportaspekte dominieren.

Obwohl die in dieser Arbeit betrachtete Dynamik von Multi-Agenten seinem Wesen nach von menschlichem Verhalten inspiriert ist, ist dieser Modellansatz bei weitem nicht auf soziale Systeme beschränkt. Um tatsächlich RFID oder "smart tags" in jeder beliebigen Art an umlaufenden Stückzahlen (Rohstoffen, Endprodukten, Stapelpaletten,...) einzubringen, sind für solche Art von Anwendungen einer dezentralisierten Dynamik in Logistik, Liefer- oder Netzwerke umgehend höchst maßgeschneiderte Produkte erlaubt. Zur Veranschaulichung beabsichtigen wir eine neue Dynamik einer vol-

lkommen dezentralisierten Load Sharing Policy (LSP) einzuführen, welche optimal die anfallende Arbeitslast von einem Satz an Benutzern entsprechend der laufenden Verfügbarkeit abfertigt. Die grundlegenden dezentralisierten Entscheidungen beruhen auf einem “smart task”-Paradigma, in welchem jede einfallende Aufgabe mit spezifischen unabhängigen Mechanismen zur Routenauswahl ausgestattet ist. Mit einer optimalen LSP meinen wir hier, dass unsere auf smart tasks basierenden Algorithmen permanent die Beschäftigung einer minimalen Anzahl von Betreibern erfordern, während dabei stets das Fälligkeitsdatum berücksichtigt wird. Die zahlreichen stigmergischen Interaktionen zwischen diesen autonomen Aufgaben sind einzig und allein die Ursache für das Eintreten einer optimalen LSP.

Schlüsselwörter: Warteschlangensysteme, Dynamik von Multi-Agenten, autonome von der Vorgeschichte abhängige Ablaufplanung, nichtlinear verzögerte Rückkopplung, räumliche und zeitliche Muster, Selbstorganisation, Komplexe Adaptive Systeme, Produktions- und Dienstleistungsnetzwerke.

Acknowledgments - Remerciements

Obviously, a thesis is never the matter of a single person and I would like to thank here the numerous people that have, directly or indirectly, contributed to the success of this unforgettable experience.

First of all, there are no words to express the deepest gratitude I have to my thesis advisor, Max-Olivier Hongler. After such four rewarding years, I cannot but realize how lucky I have been, as an apprentice researcher, to have Max as a guide. I am very grateful to him to have let me evolve, in complete freedom, in such a highly challenging and fulfilling atmosphere. Throughout the last four years, I have deeply appreciated his communicative curiosity and the numerous inventive ideas he has constantly shared with me. Max's everyday availability, his valuable help, his reassuring guidance and his tremendous scientific knowledge were essential in the elaboration of the present thesis. Beside that, Max has always been sharing more than only science with me, making this period of life a very enriching human experience. I have had the chance to discover a very endearing and generous personality and I am happy to have you now, Max, as a lasting friend.

I would like to thank Jacques Jacot to have hosted me in his laboratory during the time of my PhD thesis. I would also like to warmly express to him my gratitude for his helpful expertise concerning production research and industrial considerations, domains that were previously unknown to me due to my mostly theoretical background. I am also grateful to him to have let me supervise a master project in collaboration with an industrial partner, I have been learning a lot in this experience. I warmly thank Ann van Ackere, Dieter Armbruster and Charles Pfister for having accepted to be members of my thesis committee and for having read so carefully the manuscript. It led to a fruitful feedback and I am very grateful for their valuable comments and suggestions which helped to improve the quality of this manuscript. I wish to thank Aude Billard for having held the position of president of my thesis committee.

I still want to thank here all the colleagues I have been in touch with throughout the last four years for interesting and fruitful conversations. Beside work-mates at EPFL, my gratitude also goes to colleagues in Funchal, Bielefeld, Bremen, Grenoble and Paris. I would particularly like to thank Roger Filliger for his close collaboration with our research group, his always constructive comments and for numerous enlightening discussions. And also thank you, Roger, for having shared with me, during this discussion four years ago, your great experience with Max and for having thus even strengthened my decision to undertake a thesis with him.

To close the first part of these acknowledgements, I would like to thank Panagiotis Tzieropoulos for helpful conversations on transportation issues that ultimately led to the writing of Chapter 10. My gratitude also goes to François Genoud for fruitful discussions, in Knokke-Le-Zoute, concerning the appendix of Chapter 7. I am grateful to Olivier Lévêque for having indicated to me the open PhD position in Max's group. Finally, I warmly thank Richard Colmorn for his salutary help to produce the german abstract.

Je tiens à remercier sincèrement Julio Rodriguez, avec qui j'ai eu beaucoup de plaisir à partager le même bureau ces deux dernières années. Merci pour ton amitié, ta disponibilité et pour les nombreuses discussions, scientifiques ou autres, que nous avons eues tout au long de notre cohabitation. Je te passe à présent le flambeau et te souhaite tout le succès possible pour la suite de ta thèse. Merci également à mes autres autres collègues de bureau successifs, Christophe Vignat, Ludovic Charvier et Niklaus Johner (un merci particulier à toi pour m'avoir guidé dans les méandres administratifs de fin de thèse), pour les bons moments passés en votre compagnie, les rires et les discussions variées et enrichissantes. Je remercie vivement mes collègues du Laboratoire de Production Microtechnique pour l'excellente ambiance qui règne en son sein et pour l'organisation de nombreuses activités extra-scolaires chatoyantes. Tout travail scientifique n'est rendu possible qu'avec le soutien d'une équipe logistique et administrative de qualité. Ma gratitude va donc à Karine Genoud, Diane Morier-Genoud et Morgane Katz pour leur aide, leur disponibilité et leur précieux travail de tous les jours.

Je remercie tout spécialement mes parents, ma maman Anne et mon papa Jean-Pierre, à qui, est-il nécessaire de l'écrire, je dois tout et sans qui cette thèse n'aurait pas vu le jour. Je leur suis reconnaissant de tout coeur de m'avoir toujours soutenu tout au long de mon parcours académique et de m'avoir laissé à chaque étape de ma vie libre choix de la voie que je souhaitais suivre. Un grand merci également ma soeur, ma tante, mon oncle et ma grand-mère pour leur soutien affectueux. J'ai une pensée pour ma chère grande-tante Anne-Marie, qui m'a toujours exprimé son sincère soutien dans mon entreprise académique et qui doit être fière de moi aujourd'hui. J'aimerais également

remercier la famille Genoud, pour son soutien sans faille et son chaleureux accueil. Je remercie Alain, Alexandre, Claes, Jérémie, Steve, Tomacz et Valentine pour les belles années que nous avons passées ensemble à l'EPFL, qui ont probablement conforté mon choix de vouloir y rester quatre années supplémentaires. Je ne peux mentionner ici tous les amis auxquels j'aimerais dire merci pour les bons moments passés au cours des quatre dernières années et qui ont contribué au fait que cette période ne me laissera que de merveilleux souvenirs. Je remercie Isabelle pour son soutien et son amour de tous les jours. Tu as profondément contribué à la réussite de cette entreprise et tu as fait que ce travail a été réalisé dans le bonheur le plus absolu. Merci également de m'avoir aéré l'esprit en me proposant de partager tes activités architecturales, telles que le bricolage de nuit, dans un atelier-maquette se trouvant à deux pas de mon bureau qui plus est. Je tiens enfin à remercier Michael Jackson, Marvin Gaye et Roger Federer, pour leur contribution, non-quantifiable et purement émotionnelle, au présent manuscrit.

Contents

Part I Introduction and Literature Review

1	Introduction	3
1.1	On the Relationship Between Waiting Time and Customer Satisfaction - An Introductory Example	3
1.2	Fundamental Motivations - Multi-Agent Dynamics in Queueing Networks	4
1.3	Essential Differences with Classical Approaches	7
1.4	Reducing and Handling the Complexity - On the Importance of Stylized Models	10
1.5	Original Contributions Exposed in this Thesis	12
1.5.1	Fundamental Point of View	12
1.5.2	Conceptual Point of View	12
1.6	Organization	13
2	A Brief Review of Related Literature	15
2.1	Directly Related Literature	15
2.1.1	Queueing Networks with History-Based Features	15
2.1.2	Analytically Tractable Manufacturing Systems	17
2.2	A Brief Historical Review on Complex Systems	19
2.3	Related Literature on Multi-Agent Systems	21
2.4	Related Literature on Dynamical Systems Involving Time Delays	23

Part II Autonomous Agents in a Single-Server Feedback Queueing System - Modelling a Recurrent Service

3	A First Model - Single-Server Queueing System with Feedback Loop	29
3.1	Introduction	29

3.2	Model	31
3.3	Classical Analysis - Assuming Markovian Evolution	32
3.4	Micro-Scale Agent-Based Analysis - Siphon Dynamics	35
3.5	Multi-Agent Induced Limit-Cycle - Stigmergic Relaxation Oscillator	39
3.6	On Relationships with Well-Known Stable Limit-Cycles	41
3.6.1	The Unique Nonlinear Oscillator - Amplitude-Dependent Frequency	41
3.6.2	Supply Chains - Reorder Point Policy	43
3.7	Centralized Versus Decentralized Control to Achieve Queue Length Stability - Distinct Dynamics	44
3.8	Relaxing Some Assumptions	45
3.8.1	Delayed Feedback	45
3.8.2	Heterogeneous Agents	50
3.9	Implementation Within Logistics Networks - The Smart Parts Paradigm	55
3.10	A Solvable Occurrence of a Complex Adaptive System	56
3.11	Concluding Remarks	59
3.12	Contributions of Chapter 3	60
4	Extended Model for Recurrent Services - Introducing Weariness Aspects	61
4.1	Introduction	61
4.2	Model - Considering Weariness	63
4.3	Experimental Observations	65
4.4	Analytical Discussion	67
4.5	Concluding Remarks	72
4.6	Contributions of Chapter 4	74

Part III Multiple-Stage Feedback Queueing Systems - Networks with Competing Servers

5	Parallel Servers with Feedback Loop - Stabilization by Noise	77
5.1	Introduction	77
5.2	Model	79
5.3	Fixed Entrance Dispatching Rule	80
5.3.1	Deterministic Polling	80
5.3.2	Random Dispatching Rule	80
5.4	Entrance Dispatching Based on a Partial Observation of the Queue Contents - Noise Induced Stabilization	81
5.4.1	Experimental Observations	83
5.4.2	Analytical Approach	83

5.5	Flow Dispatching Based on Fully Observable Queues - Synchronization of Oscillations	86
5.6	Concluding Remarks	86
5.7	Contributions of Chapter 5	87
6	Closed Network of Two Servers with Feedback Loop - Competing Server Dynamics	89
6.1	Introduction	89
6.2	Model	90
6.3	Market Partition Dynamics	92
6.4	Concluding Remarks	95
6.5	Contributions of Chapter 6	95

Part IV Introducing Spatial Aspects - Market Sharing Spatio-Temporal Dynamics

7	Spatial Market Sharing Dynamics Between Two Service Providers	99
7.1	Introduction	99
7.2	Model	101
7.3	Symmetric Configurations	105
7.3.1	Explicit Illustration - Symmetric Case	109
7.4	Asymmetric Configurations	113
7.4.1	Heterogeneous Servers	114
7.4.2	Asymmetric Positions and Different Prices	116
7.5	Concluding Remarks	117
7.6	Contributions of Chapter 7	118
8	Spatial Market Sharing Dynamics in Presence of Customers' History-Based Decision Policy	119
8.1	Introduction	119
8.2	Model	120
8.3	Exploration via Simulation Experiments - Prospective Results	122
8.4	Concluding Remarks	126
8.5	Contributions of Chapter 8	126

Part V Towards Possible Applications

9	Agent-Based Optimal Real-Time Load Sharing - Application to Manufacturing Systems	131
9.1	Introduction	131
9.2	Basic Modelling Framework	134
9.3	Multi-Agent Type Algorithm	136

9.4 Emergence of Optimal Load Sharing Dynamics 136
9.5 Queue Content Oscillatory Behaviour - Siphon Dynamics 140
9.6 Concluding Remarks 143
9.7 Contributions of Chapter 9 145

10 Other Possible Applications 147
10.1 Transportation Networks 147
10.2 Smart Parts Driven Supply-Chains 151
10.3 Contributions of Chapter 10 156

Part VI Conclusion and Perspectives

11 Conclusion and Perspectives 159

Part VII Appendix

12 Appendix Chapter 2 - Typical Delayed Dynamics 165
12.1 Introduction 165
12.2 Crowded Day at the Pantheon 166

13 Appendix Chapter 7 171

References 175

Index 183

Curriculum Vitae 187

Introduction and Literature Review

Introduction

1.1 On the Relationship Between Waiting Time and Customer Satisfaction - An Introductory Example

Nowadays, the situations where one is subject to wait in queues are plentiful, ranging from transportation issues to leisure activities. Furthermore, time is in our globalized society a rare commodity that people try to preserve their best. In this regard, customers asking for any service turn out to be more and more demanding concerning the delay to receive it. Numerous studies have indeed clearly identified the great influence that waiting time plays on customer satisfaction. In the case of recurrent services, the actual satisfaction felt by the customers reveals itself to be a key factor as it determines whether they will stay loyal to a service provider in the future. Such relationship between the customers' satisfaction and their loyalty will manifestly affect the dynamics governing the queuing processes in front of the service providers. Let us now illustrate this assertion with the help of a simple yet realistic example which describes a typical situation that could arise in many various service systems.

We consider here a particular ski track, like ones we find at any ski resort. As it is often the case, many skiers use simultaneously the same track and hence, after each run, a skier possibly has to wait in order to take the ski lift that will convey him (her) again to the top of the slope. Among other possible satisfaction criteria (weariness due to successive uses of the same track, quality of the snow, etc.), the waiting time suffered at the ski lift majorly influences the skier's decision to stay on the same ski track or to go farther and try another one. Basically, each skier possesses a personal patience threshold below which he (she) will be satisfied with the service at the ski lift. This threshold clearly governs the skiers' decision to leave a given track or not. Hence, skiers' routing decisions within ski track networks sensibly depend on their individual waiting time history within the networks, giving thus an explicit non-Markovian character to the dynamics arising in such systems. The addition of these numerous local history-based routing decisions creates global collective patterns



Fig. 1.1. Samivel, *Les Chenilles Processionnelles*, 1970, aquarelle. *Copyright: Musée d'Ethnographie de Genève.*

that might take, for example, the form of temporal oscillations of the queue length at ski lifts, a phenomenon that is very likely to be observed at any ski resort. These emerging collective structures being fully induced by the skiers' individual and autonomous actions, it becomes thus essential to consider such class of microscopic local decisions in order to get an accurate description of the global system behaviour.

1.2 Fundamental Motivations - Multi-Agent Dynamics in Queueing Networks

The optimal control of the flow dynamics of matter, information and money feeding complex network structures is a classical topic in operational research. This general problem arises naturally in several strategic areas such as production/supply chains, passengers/cargo transport and computerized communication systems. The flow dynamics depend jointly on the routing rules defining the ways the flows are dispatched at the network vertexes and on the dynamics of the servers which process the various items in circulation. Due to stochastic customer demands, to fluctuations in the raw material in supply chains, to failures arising in the production devices, to uncertainties

in operator availability and to ubiquitous financial volatility steadily affecting optimization objectives, the flow dynamics are always affected by random fluctuations. The need to model, to study and to quantify the characteristics of such complex stochastic dynamics has strongly stimulated the development of the queueing network (QN) theory. Nowadays, the QN theory offers a wealth of reliable mathematical tools for calculating most relevant performance measures of such dynamics. The basic hypothesis behind any QN modelling is the possibility to describe the underlying dynamics by general Markov processes. When this is realized, very general results (pioneered by the widely known Jackson factorization theorem) are available to characterize time-invariant flow regimes and hence to compute stationary performance measures (a more detailed description of these classical concepts is given in Section 1.3). Imposing such a Markovian character obviously limits not only the dynamics of the servers but also lays down strong restrictions on the allowable routing rules followed by the circulating items. In the present work, we will study networks for which the Markovian character of the dynamics has to be abandoned. More precisely, the non-Markovian features will originate from the routing decisions which will be based on the items' individual experience collected during their journey through the QNs. In other words, we will consider history-based (HB) routing laws along this thesis. The presence of memory mechanisms in routing decision rules explicitly precludes the Markovian character of the underlying dynamics and this opens wide the door for the emergence of entirely new dynamical features. More particularly, the very existence of stationary regimes can no longer be taken for granted when HB rules are implemented. As we will see, the joint presence of feedback loop topologies in the QNs (*i.e.* possibilities of flow re-injections) and HB routing decisions is responsible for the emergence of collective spatio-temporal patterns in the flow dynamics. As a typical example of such global structures, the implementation of autonomous HB routing rules in QNs actually implies the existence of delay effects in the dynamics and as a result oscillatory behaviours are likely to be observed. As it will be pointed out in this thesis, note that depending on the specific objective functions to be fulfilled within the QNs, the possible emerging collective patterns might be either seen as positive or negative effects and should hence be respectively favoured or suppressed.

In general, numerous HB routing rules could be considered. In this thesis we will mainly focus on situations where the time spent in specific sections of the QNs will determine the criterion used to select a specific route on which to engage at a bifurcating node of the network. Implicitly, such waiting time criteria require a real-time capability for each circulating agent¹ to monitor, to memorize and then to process data to ultimately form an individual routing decision. As a direct consequence, one realizes that we actually deal with

¹ From now on, we will speak of agents, customers and “smart” parts interchangeably.

QNs roamed by autonomous, decision making (*i.e.* “intelligent”) agents, with *stigmergic*² mutual interactions. Furthermore, since the different emerging collective patterns that will be described in this work are solely induced by the agents’ autonomous routing decisions, our particular QNs might be seen as particular instances of complex adaptive systems.

While the multi-agent type of dynamics considered here is in essence inspired by human behaviour, our modelling framework is by far not only restricted to social (service) networks. Indeed, attaching RFID or smart tags on any kind of circulating units allows for the implementation of such type of decentralized dynamics in logistics, supply or manufacturing networks. To emphasize the actual importance of developing such type of decentralization techniques, it is enlightening to know that the German Research Foundation (DFG) is currently massively funding a research project³ which goes exactly into that direction. As the dynamic and structural complexity of logistics networks makes it very difficult to provide all information necessary for fully central planning and control, this research project, via a highly holistic and cross-disciplinary approach, studies in detail the theoretical and practical dimensions concerning the implementation of autonomous processes within logistics systems.

In regard to the above considerations, it follows that the present work is strongly interdisciplinary as it stands at the frontier between various active research fields, namely QN theory, agent-based modelling, complexity science and management of logistics networks. Note in addition that the approach chosen in this thesis is to avoid idiosyncrasy and consequently to construct stylized models that remain analytically tractable thanks to a limited number of considered parameters. The development of such idealized models is important since they potentially allow for the calibration and validation of more realistic and idiosyncratic frameworks that can be characterized solely by simulation techniques.

² Stigmergic qualifies indirect communication in a self-organizing system where individual parts (here the circulating agents) interact with one another by modifying their local environment (here the environment is basically the queue length). The concept of stigmergy was originally introduced by Pierre-Paul Grassé in the context of zoology to describe the behaviour of social insects (more particularly termites), [49]. Stigmergy occurs when the behaviour of a subject is determined or influenced by the consequences of the other subjects’ previous actions.

³ Collaborative Research Centre 637 “Autonomous Logistics Processes - A Paradigm Shift and its Limitations” (SFB 637), active from 2004 and at least until 2011 at the University of Bremen, Germany.

1.3 Essential Differences with Classical Approaches

Stochastic networks, the other appellation of QNs, refer to queueing systems in which customers move between different stations (the network nodes) where they receive services. One encounters such networks in many various areas such as transportation, telecommunication, distributed computing, manufacturing or supply chains. This high potential for applicability has engendered a prolific research activity in the last fifty years, which has resulted in the construction of a dedicated theory that provides powerful tools useful for the description of general QNs. As exposed in [120], classical QN theory mainly consists in approximating the dynamics of queueing systems by time-homogeneous Markov chains. The ultimate goal of this modelling technique is to compute the time-invariant distribution of these Markov chains by solving a linear system of equations (the so-called *routing balance equations*). Provided the chain is irreducible⁴ and ergodic, a unique stationary distribution exists and this measure is also the limiting distribution of the chain (*i.e.* the system is asymptotically stationary). This time-invariant distribution, which fully characterizes the stationary state of the associated QN, might then be used to compute valuable information about the system behaviour, such as the average throughput of the servers or the customers' mean sojourn time at the network nodes.

It is possible to describe a QN with a Markov chain when the three following properties are satisfied:

- (1) The external arrivals feeding the network can be modelled by general renewal processes.
- (2) Similarly, the successive service times at each node of the network also form a general renewal process.
- (3) Customers move between the nodes of the network according to independent Markovian routing.

Under the above hypothesis, a stochastic network can be represented by a Markov process whose dynamics $\{\zeta_t, t \in \mathbb{R}\}$ at time t start afresh from ζ_t . In other words, the state of the queueing system at time t contains all the relevant information about past evolution that is mandatory to predict the future behaviour. While the memoryless properties (1) and (2) of the renewal processes describing the services and arrivals will be satisfied in the present work, the assumption (3) on the Markovian character of the routing will in our context be explicitly violated. Indeed, due to their HB characteristics, the routing mechanisms considered in this work will obviously not be independent of the past evolution of the network. Consequently, it will make the description of our queueing systems by Markov chains impossible, as well as it will preclude classical computation of stationary measures with the help of linear

⁴ The probability of transition between any two possible states of the chain is strictly positive.

routing balance equations.

To emphasize the novelty of our approach, we describe in the following the two most famous classes of QNs, namely Jackson and Whittle networks. Contrary to the models that will be introduced in this work, the whole three above Markov properties are valid for both of these classes and hence general results are available for the description of stationary regimes. These two wide classes of networks, which are consequently both Markov processes, have notably become prominent because their explicit stationary distributions exhibit comfortable closed-form (*i.e.* product-form) expressions.

(1) *Jackson networks.*

One of the simplest, yet general class of QNs is universally known as Jackson networks. Indeed, this first significant development in QNs has first been introduced in the sixties by Jackson, [65, 66]. A Jackson network consists in a finite number of nodes, each one representing a service facility. The customers are assumed homogeneous (at each node, the service time distribution and the routing mechanism is the same for all customers) and the queue discipline is FIFO. In the case of an open network, customers arrive from outside following a Poisson process. While the service times at each node are assumed to be exponentially distributed, the service rate can be both node-dependent and state-dependent. More particularly, in a Jackson network, the service rate and routing at each node might solely depend on the current congestion state of that node (*i.e.* the dependency only involves the number of customers waiting at the moment at that service facility). The management procedures are hence only based on local (incomplete) information about the system state. It is widely established that Jackson networks are ergodic and that it is possible to solve easily the routing balance equations (also called traffic equations) in order to find the unique invariant measure, which possesses a simple product form, [29, 48, 109].

Between 1960 and 1975, Jackson networks were the only ones used for modelling QNs on account of their simplicity of use. The emergence of computer systems restarted the research into simple solutions for new types of networks and gave the impulsion for the development of a class of QNs that is extensively known nowadays under the name of Whittle networks, [123].

(2) *Whittle networks.*

A significant extension of Jackson networks has been developed over the years and constitutes the class of Whittle networks. In the Whittle modelling framework, the service rate and the routing at a particular node might not only depend on the current congestion state of that node but on the global congestion state of the network (*i.e.* the dependency involves here the number of customers being at the moment at all the nodes). Concerning the routing mechanisms, while they still have to be assumed Markovian, they can now be based on an enhanced (complete) knowledge

of the current system state. Moreover, Whittle processes are perfectly suitable to model networks with multiple types of units. In that situation, each circulating unit carries a class label which can be permanent or temporary and subject to change as the unit moves through the network. Following this setting, the routing mechanisms and the service rates might now depend on the customer type. In essence, the only difference is that one has to construct now a Markov process that keeps track not only of the global number of units at the nodes of the network but that characterizes the number of units of each type at all the nodes. Obviously, the number of different customer types has to be finite to allow for such a construction. While the queue policy was necessarily FIFO for Jackson networks, the introduction of different classes of customers allows for the implementation of policies that give priority to specific customers. Despite their enhanced complexity, Whittle networks still allow for the straightforward computation of product-form stationary distributions, [29, 109]. The two following sub-classes of Whittle networks are pointed out as they are especially prominent.

- *Kelly networks*. These particular Whittle processes, introduced by Kelly in [73], are multiclass networks in which the customers are divided into classes with respect to their specific route through the network. Accordingly, the type of a customer is allowed to influence its routing decisions as well as its exponential service time distribution at each queue.
- *Baskett, Chandy, Muntz and Palacios (BCMP) networks*. In these specific networks, the service times depend on both the number of customers at the node and the number of customers of the same type as that being served, [16]. Moreover, the service times are generalized in BCMP networks and might follow in that context Cox distributions. Note that these particular distributions still possess the Markov property, but in a multidimensional space, [48].

Formally, following the classical modelling frameworks described above, we would be along this thesis in presence of multiclass QNs with an infinite number of customer classes. Indeed, after each service completion, the circulating agents should be categorized with respect to their individual experimented waiting times, as these HB measures determine their future routing through the network. As the distribution of the waiting times is assumed continuous in our modelling framework (arrival and service times follow renewal processes), an infinite number of classes would necessarily be needed to classify the agents and for this reason we would violate one of the pillar hypothesis of classical multiclass QNs.

1.4 Reducing and Handling the Complexity - On the Importance of Stylized Models

In the last decade there has been a large interest in the development of agent-based models (ABMs) aimed at reproducing and understanding the emergence of collective patterns from the lower level of complex systems to a higher level. In other words, multi-agent frameworks are built to allow for a microscopic description of the many local actions and interactions governing the global behaviour of various dynamical systems (a more detailed discussion on ABMs as well as a succinct state of the art in that field is given in Section 2.3). In this perspective, two opposite approaches might be considered.

The mainstream way (that composes the great majority of the devoted literature) to treat multi-agent systems tends towards idiosyncrasy and, accordingly, towards the development of extremely detailed modelling frameworks. In that regard, the number of considered parameters becomes rapidly large, in order to represent agents' heterogeneity and local possibly complex behaviours. Such a modelling approach obviously leads to models with high intrinsic complexity and an analytical resolution becomes in that context very unlikely to obtain. The ever increasing computational power available nowadays is hence mandatory to handle such complex multi-agent systems and specific powerful simulation frameworks are developed in this sense. The main advantage of such pure numerical description techniques is that one can imagine, model and simulate almost any level of complexity (within the limits of the available computing capacity). However, the major resulting drawback is a lack of potential for generalization.

On the other hand, the point of view adopted in this thesis is to handle the complexity by reducing drastically the number of considered parameters and to propose consequently stylized (*i.e.* idealized) models that remain analytically tractable. Such kind of models nevertheless permits to obtain a detailed understanding of the origin and nature of the emerging collective patterns. Moreover, the simplicity of the models implies that changes and variants can be addressed in a simple way, in order to eventually incorporate additional features. The same choice of modelling approach has been recently adopted in economics, [3]. In this contribution, a minimal ABM for financial markets is proposed to understand the nature of several global self-organized patterns that appear in this context. While real markets are obviously extremely complex systems, this stylized model follows the alternative of incorporating only the essential ingredients to reproduce the most important deviations that might be observed in price evolution. As a consequence, the understanding of these collective phenomena becomes in this framework more limpid. In spite of the limited number of considered parameters, this study on financial markets however represents a solid basis on which one might eventually add more realistic features thereafter. Note that stylized models naturally fulfil

the logical principle of the *Occam's razor*⁵ (also known as the *law of parsimony*), which states that the explanation of any phenomenon should make as few assumptions as possible and should consider the smallest number of parameters, eliminating those that make no difference in the observed evolution.

It is interesting to wonder about the pertinence of the two antagonistic modelling approaches described above and the potential synergy that exists between them. When idiosyncratic simulation frameworks, that tend to be a very detailed representation of reality, are without contest powerful and close to applications (and thus potentially easily transferable to the private sector), their complexity might make their analysis becomes very tricky. Some questions arise: how to validate the simulated results and how to determine if they are representative? Furthermore, how to determine whether the considered model is the “right” one and whether it is correctly calibrated? Due to the large range of considered parameters, validation methods are rare (they are in most cases nonexistent), as they are very complicated to establish. In the matter of consistency, it remains however absolutely essential to be able to determine the sensitivity of the models, their robustness as well as their level of generality. As an example, this modelling problem has recently arisen in economics, more precisely in several successive attempts to model the French labour market. The stylized model proposed in [26], which analyzes the introduction of a new job contract into the labour market, is missing some features that appear at the microscopic level (*i.e.* oscillations of the unemployment rate, which describe job precariousness) in the idiosyncratic multi-agent simulation framework developed in [83]. However, the calibration and validation of this simulation modelling framework would have been impossible without the former mean-field idealized model introduced in [26]. This example perfectly shows the mutual need that exists between the two antagonistic modelling approaches and, although stylized models are in a sense less complete and obviously truncated representation of real-life applications, there is clearly a synergy between them and the more realistic microscopic models obtained by pure simulation techniques. Indeed, the *workability* of the first ones (*i.e.* their handled complexity) might potentially allow for the essential validation and calibration of the second ones.

In [90], B. McKelvey makes a virulent attack on the idiosyncratic trend that widely predominates nowadays in organization science and he carefully argues for the importance of the development in parallel of idealized models. History has shown that such idealizations, which assume uniform rather than complex heterogeneous microstates and hence allow for tractable studies, have often proven to be useful in the life cycle of a science. Indeed, McKelvey defends in his diatribe that, although stylized models are not exact representations of

⁵ “*Entities should not be multiplied unnecessarily*”, statement apocryphally attributed to William of Ockham, 14th century.

real phenomena, they have in many situations provided a necessary basis to give way thereafter to more realistic modelling frameworks.

1.5 Original Contributions Exposed in this Thesis

In this section, we briefly expose the main contributions contained in this work, first from a fundamental angle then from a conceptual point of view. Note that a detailed list of all the new results and contributions enclosed in each chapter is provided in the sections concluding every chapter.

1.5.1 Fundamental Point of View

The study of spatio-temporal flow patterns due to autonomous agents traveling with HB (hence non-Markovian) routing rules remains, despite its truly strong interdisciplinary integration, an almost unexplored topic in QN theory. On the other hand, with a view to potential applications (besides service systems), the highly flexible modern production and supply networks, able to satisfy extreme ranges of customized products, rely more and more on decentralized mechanisms able to self-organize the material flows visiting intricate network topologies. In this context, we introduce in this thesis several solvable models that exhibit some of the collective phenomena that could emerge within the dynamics of such systems. We describe mostly analytically these emerging self-organized patterns. As a result of this analysis, we actually transfer and introduce in the field of service and production systems several phenomena that are classical in basic sciences, namely sustained stable oscillations, synchronization of oscillators, stabilization by noise phenomena and noise-induced phase transitions.

1.5.2 Conceptual Point of View

When dealing with decentralized, agent-based management in logistics systems, apart simulation experiments devoted to very specific case studies, only rather prospective and mostly conceptual contributions are available in the existing literature. To the best of our knowledge, there barely exists any available mathematical modelling studies where the potentiality of the “smart parts” (*i.e.* “intelligent” circulating units) concept for production, service and supply chains networks is analyzed. In this regard, the simple yet paradigmatic, idealized and didactic models proposed in this thesis are highly welcome. They prove that solvable instances of complex (logistics) adaptive systems might exist and be thoroughly studied. Furthermore, we provide a tractable example that clearly illustrates how emerging self-organized patterns in multi-agent manufacturing systems can be used to optimize the dynamics. Note finally that the selection of our models is strongly based on their potential industrial

relevance and therefore they should retain the attention of a multidisciplinary audience composed of the community of complex systems scientists as well as production managers.

1.6 Organization

This thesis is divided into seven parts, namely:

- I: Introduction to the fundamental motivations of this thesis and review of the related literature.
- II: Single-server feedback queueing system roamed by autonomous agents with HB routing decisions.
- III: Network topologies involving two competing servers with feedback loops.
- IV: Introducing spatial aspects into the queueing dynamics of two competing servers.
- V: Towards possible applications.
- VI: General conclusion and perspectives.
- VII: Appendices.

These different parts are, in turn, subdivided into chapters that have been conceived to be as self-contained as possible with dedicated introductory parts and concluding remarks.

Part I starts with Chapter 1 (the one you are actually reading now) which presents the essential motivations that led to the writing of this thesis, namely the fundamental reasons to study the circulation of autonomous agents in QNs. Chapter 2 is dedicated to a brief review of the literature related to the present work.

In Part II, we consider a single-server feedback queueing system as a stylized way to model a recurrent service. In Chapter 3, we introduce the HB routing rule, based on the circulating agents' individual elapsed waiting times, that will govern their autonomous decision whether to take the feedback loop or not. We describe the emerging self-organized collective pattern (*i.e.* the stable temporal oscillations of the queue content) that results from the numerous agents' local actions and interactions. In Chapter 4, we propose an extension to the model previously studied in Chapter 3, where we consider weariness aspects.

Different network topologies involving two competing servers with feedback loop are discussed in Part III, which is composed of Chapters 5 and 6. While Chapter 5 is dedicated to an open network with two parallel servers, Chapter 6 is devoted to a closed market topology.

In Part IV, which is composed of Chapters 7 and 8, spatial aspects are considered in the dynamics of queueing systems. More particularly, we study in Chapter 7 the market partition between two distinct providers that deliver services to customers whose behaviour is sensitive to waiting time and transportation costs. In Chapter 8, we extend somehow the configuration previously studied in Chapter 7 and we study how such a spatial market is dynamically shared in presence of recurrent customers that base their satisfaction on experienced waiting times.

Several possible applications of the concepts developed in this thesis are discussed in Part V. In Chapter 9, focusing on manufacturing systems, we propose a new fully decentralized dynamic load sharing policy, that optimally dispatches the global incoming workload according to the current availability of a set of operators. Chapter 10 is devoted to a discussion on how the different aspects considered in the present work could be extended to the domains of transportation and supply chains.

A general conclusion as well as perspectives arising from the present thesis are given in Part VI, and more particularly in Chapter 11.

Part VII is devoted to appendices. In Chapter 12, we present an explicit illustration, inspired of a real-life situation, of how the existence of time delays in the dynamics of queueing systems might produce oscillations phenomena. Chapter 13 provides some technical details related to Chapter 7.

This book ends with an extensive bibliography, a subject index and the author's brief curriculum vitae and publications.

A Brief Review of Related Literature

This chapter is dedicated to a brief review of the literature that is related to the multidisciplinary aspects developed in this thesis. In the following, we mention the most prominent contributions that share some similarities and relationships with the present work. First we expose the rather scarce available contributions that are directly in the vein of the modelling framework proposed in this thesis, namely queueing networks (QNs) with history-based (HB) routing mechanisms. Then we provide a succinct state of the art in the following areas that are related to the interdisciplinary topic of this work:

- Complex Systems
- Multi-Agent Systems
- Dynamical Systems Involving Time Delays

2.1 Directly Related Literature

2.1.1 Queueing Networks with History-Based Features

Networks in which HB routing decisions are present can be easily identified in various important contexts such as transportation (pedestrian, car, train and air traffic issues), production and supply chains, leisure and hospitality (theme parks, ski resorts, hotels and food industry management) and health care industry. Despite this wealth of applicability, the literature devoted to formal models studying the flow dynamics involving HB routing rules remains so far rather scarce. Obviously, this is partly due to the technical difficulties inherent in the non-Markovian and nonlinear features of the underlying dynamics. Nevertheless, several recent illustrations where directly related (yet mostly experimental) situations are handled can be pointed out.

(1) Leisure and Hospitality.

In [17], the influence of waiting time on the satisfaction and loyalty of customers using recurrent services is exhibited and explicitly studied. More

particularly, following a survey conducted in the medical care industry (which is in many aspects in close connection with the hospitality business sector), contribution [17] aims to investigate how customers use their waiting time satisfaction in order to determine whether to remain loyal (*i.e.* whether they will come back for future services) or alternatively to change their service provider. As discussed in [79], the waiting time also influences significantly the satisfaction and return decision of customers visiting fast-food facilities. By analogy with the modelling framework considered in this thesis, general models of recurrent services where customer satisfaction and loyalty are driven solely by the perceived waiting times while queueing are presented in [56, 114, 115]. In [115], commuters having the choice between alternative roads base their routing decisions on their neighbours' most recent waiting experience and on their own complete waiting history (*i.e.* on their exponentially weighted average experimented waiting times for each road). This contribution explicitly exhibits the striking feature that a self-organizing system based on local information and locally rational agents might outperform (*i.e.* the average travel time of commuters is reduced) the Nash equilibrium (that assumes full information). A general recurrent service where customers establish their loyalty to the service provider in function of a long-term satisfaction measure (*i.e.* on their average waiting time) is proposed in [56]. This contribution shows that, even in a purely deterministic framework, a queueing system with long-term feedback can exhibit queue length stabilization, periodic evolution or chaotic behaviour. In [114], a short-term feedback is added to the model considered in [56]. More precisely and in addition to the long-term feedback previously introduced, customers base their short-term service quality satisfaction on their most recent experimented waiting times and they adapt their visit frequency to that particular measure. While periodic behaviour of the queue content might still be observed in this case, it is interesting to notice that no chaotic dynamics will emerge here and that, in one sense, the additional short-term feedback stabilizes the system. The models introduced in [114, 115] find natural applications in sports clubs, supermarkets and internet access management. The theme park industry is another sector to which the models presented in this thesis are closely related, [69, 71, 72]. Roughly speaking, customers will decide to line up again for a new run at the same attraction after an exciting roller-coaster ride providing their waiting time remains acceptable for them. Ski traffic management offers another world-wide illustration where customer satisfaction and hence future (HB) routing decisions are directly related to suffered waiting time, [104]. Indeed, as exposed in the introductory example of Section 1.1, the waiting time spent at a ski lift clearly affects the customers' future decision whether to leave a ski track or not.

(2) *Supply Chains and Production Management.*

The need, in supply chain management, for coordination strategies leading to adaptive, flexible and collective behaviours motivated a recent contribution proposed by A. Surana et al. [111], in which the authors show how a coherent global behaviour can be generated by using only elementary components with local interactions. This paper explains how basic concepts and operational tools of Complex Adaptive Systems¹ (CAS) fit naturally and efficiently for characterizing the supply chains dynamics. In this context, contributions [76, 111] expose the dynamics of simple QNs for which feedback loops and delays coexist and yield temporal oscillations of the queue contents. Originally introduced in the framework of telephone switching systems [42], these particular QNs have been further applied in the context of supply chains in [76, 111]. Note, however, that, contrary to the various waiting time criteria to be studied in this thesis and which mostly characterize service and production systems, the HB features in [42, 76, 111] are due to HB changes of the agents' priority status in re-entrant queueing systems.

In close connection to actual production issues, we also mention here the recent contributions devoted to *Real-Time Queueing Systems* (RTQS), [12, 38, 82]. Unlike standard queueing theory, RTQS focus on the ability of a queue discipline to meet production task timing requirements, for instance the distribution of lateness. In the RTQS described in [82], each incoming task in a queueing system is endowed with a *due date* before which it has to be completed. To reduce the potential lateness, a *dynamic scheduling policy* in which waiting tasks are placed by decreasing *leadtime* (*i.e.* the more urgent being the closest to the server) is implemented. This dynamical scheduling rule can be viewed as a special illustration of HB routing. Indeed, this *Earliest-Deadline-First priority rule* implies that each incoming task triggers a rearrangement (*i.e.* a real-time routing) dependent on the waiting history of all tasks present in the queue.

2.1.2 Analytically Tractable Manufacturing Systems

In Chapter 9, we propose a fully decentralized (agent-based) load sharing policy, based on HB data collected individually by the circulating items, that optimally dispatches the incoming workload according to the current availability of a (variable) set of operators. Beside a manifest relevance for applications, most particularly in production and manufacturing systems, our model is analytically tractable, a rather uncommon feature when dealing with multi-agent dynamics and complex adaptive logistics systems. To the best of our knowledge, the only other instance of an agent-based (entirely decentralized) manufacturing system allowing for an analytical description is the so-called

¹ An introductory discussion on Complex Adaptive Systems as well as a description of their principal characteristics can be found in Section 3.10.

bucket brigades model, which consists in a technique of organizing the different workers along an assembly line (more generally along a linear production line) so that the line balances itself. More specifically, following Bartholdi and Eisenstein, the situation under consideration is the following: each worker carries a product towards completion; when the last worker finishes his product, he walks back upstream to take over the work of his predecessor, who walks back and takes over the work of his predecessor and so on, until, after relinquishing his product, the first worker walks back to the start to begin a new product. In the mid-nineties, Bartholdi and Eisenstein introduced bucket brigades modelling in the context of assembly lines and were the first to give an analytical study of the dynamics that emerge in such systems, [13]. More precisely, it is shown in [13] that, if the workers are sequenced from slowest to fastest, it leads to the spontaneous generation of a stable partition of work (*i.e.* every worker repeatedly executes the same respective portion of work content on each produced item). This production arrangement moreover maximizes the overall throughput of the system, among all ways of organizing the workers and the respective portions of work. The emerging self-organized optimal balance of the work assignment is shown to be the unique fixed point to which the system will converge to and this independently of the initial positions of the workers. Note that similarly to the models developed in this thesis, and more generally typical of multi-agent systems, bucket brigades resist to the death or to the birth of a worker as the system again optimally organizes itself after such events. While the original contribution [13] has been developed for a deterministic framework, stochastic environments are considered in [14] and it appears that bucket brigades remain effective even in the presence of variability in the work content. Bucket brigades can be extended to in-tree assembly networks and it is shown in [15] that self-balancing will also emerge in these more complex systems. While it is assumed in original bucket brigades that it is possible to order the workers with respect to their speed and that this ordering remains valid over time, Armbruster and Gel propose in [6] a relevant extension where they study the dynamics arising when the workers' speed might vary drastically from one portion of the line to the other. In this case, a static ordering of the workers with respect to their speed might not be anymore reachable (*i.e.* the velocity of one worker does not uniformly dominate that of another along the whole line) and the self-organized emerging features might consequently disappear. For such environments, it is shown analytically in [6] that bucket brigades remain efficient in terms of self-balancing behaviour and throughput performance when one enables passing between the workers: when a worker is caught up by its predecessor (meaning that its predecessor has become faster), these two workers invert their position on the assembly line (in traditional bucket brigades, the order of the workers is preserved all the time). More precisely, the implementation of this passing rule may solely drive the system from an unstable regime to one where the bucket brigade self-balances. The use of the same passing rule is also appropriate when workers' speed on specific portions of the line increase due to

learning and it is shown in [7] that enabling such passing policy leads to very robust assembly lines, the dynamics of which also exhibit self-organized optimal production arrangement and overall throughput. To conclude, note that the direct interactions that take place between the agents in bucket brigades differ in essence from the agents' stigmergic interactions that we will consider along this thesis.

2.2 A Brief Historical Review on Complex Systems

This section provides the reader with a short historical review on the development of *complexity science*² over the last century and on the recent transfer of these concepts to logistics and supply networks. In this regard, we follow the main lines drawn both in the clear and concise historical survey on complex systems given in [127] and in the more detailed and complete reviews provided in [40, 86]. Complexity science is far from being a single theory: seeking to answer to some fundamental questions about living, adaptable, changeable systems, it encompasses more than one theoretical framework and is thus highly interdisciplinary.

The notion of complex systems was born at the beginning of the 20th century with the early works of H. Poincaré on the trajectory of planets. Poincaré showed that it was mathematically impossible to obtain an exact solution to the equations that describe a somehow simple system containing three planets with intrinsic nonlinear interactions. Thus, he revealed to the world that a completely causal system can exhibit an indeterminate (chaotic) behaviour. In other words, Poincaré unveiled that even a system that appears to be simple can explode into complex and unpredictable behaviour. This new conceptual difficulty generated an important research activity in the last fifty years that, following several distinct but nevertheless complementary axes, gave birth to what is called nowadays complexity science. In the following, we give a non-exhaustive list of some of the famous ancestors of the present day research on complex systems.

(1) *Game Theory.*

J. Von Neumann originally developed game theory in the 1920s as a set of tools to analyze economical behaviour by mathematical techniques, [119]. More precisely, Von Neumann proposed a formalization that economics is the outcome of the interaction of many competing agents (*i.e.* players). These agents behave according to universal rules in order to react to what the others do. This modelling framework allows thus for the description of various forms of competitive behaviours that might be found in many social complex systems, not only in the field of economics but also in numerous

² This term refers to the specific field of science dedicated to the analysis of complex systems.

applications ranging from political to military. Further developments in game theory were proposed over the years, the most famous one being the contribution of J. Nash, who discriminated between cooperative and non-cooperative games and found an equilibrium point in the case of non-cooperation (the so-called *Nash equilibrium*). Game theory has played over the years an important role in social sciences and in biology, as it can explain the emergence of various global structures within a group from simple interaction rules.

(2) *Neural Networks.*

Artificial neural network (NN) research started in the early 1950s in order to gain a better understanding of biological NNs as well as to provide a metaphor for various cognitive processes. Indeed, the idea of spontaneous order in the brain due to decentralized networks of simple neurons emerged at that time and rapidly becomes widely accepted. An artificial NN involves a network of simple processing elements (*i.e.* the artificial neurons) whose global behaviour might exhibit complex patterns. The emerging collective structure is determined by the neurons' parameters and by the strength of the links that exist between them. Usually, NNs operate in two phases. The first one consists in an unsupervised (decentralized) learning with the ultimate goal of matching a desired output. Once the learning phase is complete, the NN is able to classify incoming information and to work hence as a pattern recognition system. NNs are highly robust and adaptive since their global output is fully induced by the neurons' simple actions and interactions and consequently the system overall performance would be only smoothly affected by the removal of a neuron.

(3) *Theory of Dissipative Structures.*

Instigated by Prigogine and his "Brussel school" in the mid 1950s, the theory of dissipative structures (TDP) was a cornerstone in the later development of the theory dedicated to self-organizing systems. Grown out from the thermodynamics of open systems, the TDP intends to describe the formation of temporal, spatial and/or spatio-temporal structures in physical complex systems operating far from thermodynamic equilibrium. Originally developed in the context of physico-chemical systems, the concepts peculiar to dissipative structures have then been transferred from the late 1960s to biological and social sciences with the objective of establishing general principles about the conditions under which particles locally interact to spontaneously produce newly ordered complex collective patterns.

These distinct and highly interdisciplinary research axes (we could also mention here cybernetics and synergetics) provided the impetus for the development from the 1980s of a subclass of complex systems named *complex adaptive systems* (CAS). These particular self-organizing systems are considered complex since they are diverse and made up of multiple interconnected elements

and adaptive because these elements have the capacity to learn from experience (a more detailed discussion on CAS and connections to the present work is given in Section 3.10). Originated from biology and first pertained to living entities, CAS have progressively been transferred to logistics and supply networks [30, 111, 121] thanks to the availability of new forms of communication and information technologies (such as RFID or smart tags). Indeed, one might also observe for these systems various emergent phenomena even if they are not composed of living entities *stricto sensu* and these collective patterns solely follow from the units' autonomous nonlinear (local and causal) actions and interactions. The relevance and legitimacy of CAS in logistics, as well as their practical implications, have been recently strongly emphasized in [63, 127]. In the same vein, the actual impact of decentralized management and of the resulting self-organized features on process and product quality is addressed in [88].

2.3 Related Literature on Multi-Agent Systems

The study of the social behaviour of various species of insects has been one of the major inspiring influence for the former development of agent-based models (ABMs). First developed as a relatively simple concept in the late 1940s, agent-based modelling has waited until the 1990s and the ever growing availability of computational power (which provides powerful simulation frameworks) to play the important role it has nowadays in many research areas. We recall here the following classical definition of (autonomous) agents:

An agent is an autonomous decision-making entity which possesses the ability to observe and act upon its environment, without any external intervention. The agent individual actions result from these collected observations, but possibly also from personal historical data. An agent is rational, as its behaviour is directed towards achieving specific goals. Depending on the situation, agent complexity (concerning behaviour rules and goals to achieve) might vary from very low to very high.

The ultimate goal of ABMs is to consider in detail the local actions and interactions of autonomous individuals and to analyze the emerging (from the agent micro-scale to the collective behaviour macro-scale) global effects which might be observed in the system as a whole. This goal reflects and formalizes the popular assertion according to which *the whole is greater than the sum of its parts*, which illustrates that very simple local agent rules might result in far more complex and interesting global behaviours.

The very prolific research on multi-agent systems is at the edges of complex systems, sociology, game theory and computer science. ABMs have been applied and developed in a very wide variety of social, business and technological

domains, as well as in particle physics, chemistry and cell biology. Particular applications include word of mouth propagation, spread of epidemics, analysis of traffic congestion, portfolio management, supply chain optimization and distributed computing. ABMs might help to explain the emergence of various global phenomena like power law distributions driving the length of traffic jams, stock market crashes and the apparition of bullwhip effect in supply chains. In the following, we give some examples of works that are somehow related to the topics developed in this thesis. For a more complete overview of the numerous applications and implications of multi-agent systems, the reader is encouraged to consult contributions [18] (business applications of ABMs to human systems), [80, 98] (applications in manufacturing systems and supply chains), [4] (agent-based computational economics) and [126] (general overview of possible applications).

Internet has become more and more omnipresent in our everyday life and in that sense e-commerce promises to become very popular in the near future. In that context, ABMs have been developed as it is possible to implement agents that replace us to select a service provider and that find the products that best fit to our needs and expectations, [98]. Multi-agent systems have been also extensively used in the framework of supply chains, [80]. Indeed, ABMs allow for the consideration of the high heterogeneity and autonomy existing at the different echelons of the chains. As an example, multi-agent techniques have been used to study and reduce the effects of demand amplification (*i.e.* the bullwhip effect) within supply chains networks, [97]. Another natural application of agent-based modelling is the management of customer flows in theme parks. Indeed, the numerous customers interact via the environment, since their waiting time at an attraction manifestly depends on the other people's choices. In [69, 71, 72], a simulation framework is proposed for this situation and the global patterns that might emerge are analyzed. In contribution [74], the authors treat distributed computing systems as a self-organizing collection of autonomous agents cooperating as peers. In that case, the joint effect of nonlinearities and delays affects the underlying dynamics by creating macro-scale oscillatory behaviours.

Due to the often extreme complexity resulting from the numerous agent local actions and interactions, agent-based modelling often requires highly powerful computational procedures. At the same time and for the same reason, multi-agent systems become very quickly analytically intractable and hence a very large majority of the contributions available in the dedicated literature avoid completely mathematical methods and solely use simulation techniques to explore the emerging dynamics. In this regard, the models proposed in this thesis belong to the very few existing ABMs that allow for an entire analytical description. To the best of our knowledge, the only other example of a multi-agent production system that allows for an analytical study is the so-called bucket brigades described previously in Section 2.1.2.

2.4 Related Literature on Dynamical Systems Involving Time Delays

When a dynamical system involves any form of feedback control, it often requires a finite time to sense information and then take the appropriate reaction to it. Accordingly, the great majority of these systems, which now occupy a place of central importance in many areas of science, will imply time delays in their evolution. These kind of feedback control processes are often used to maintain the system in a stable state. Resulting dynamics are characterized by delay differential equations (DDEs), in which the evolution of the system at time t depends on its state at time $t - \tau$, $\tau > 0$. Note that, from a mathematical point of view, it is a very difficult task to solve DDEs and consequently many recent contributions use the rapid advances in computational power to provide strong numerical techniques in this context. One might be tempted to ignore small time lag effects and use ordinary differential equations (ODEs) model as a substitute for DDEs, but as it is underlined in [75], “*small delays can have large effects*” and such a simplification might hence be risky. Oscillatory behaviours are frequently associated with dynamical systems involving time delays. In this regard, the models introduced in this thesis manifestly involve such intrinsic time lags (see Section 3.4) and it is hence not surprising that their dynamics might exhibit oscillatory features.

Without the intention of being exhaustive, we point out here various application domains where we can find dynamical systems with time delays (which can hence be described by DDEs). While there exists a collection of books that treat mathematical and practical implications of DDEs (see in particular [39, 51, 75]), the reader is encouraged to consult the recent contribution of T. Erneux [41]; this book provides a very rich and nicely illustrated overview of the different possible application areas for DDEs. Note that, depending on the particular situation, the oscillatory instabilities resulting from the presence of time delays might be either seen as positive or negative effects and should therefore be respectively favoured or suppressed.

(1) Ecology.

The population densities of many different species can fluctuate nearly periodically over time, even if there is no predatory interaction of other species, [41]. Ecologists have come to the conclusion that the period of these oscillations cannot be explained simply by seasonal variations. Consequently, the incorporation of a delay, characterizing the gestation times as well as the environmental conditions such as supply of food, in the *logistic equation*³ has been unveiled to be the cause of such population density

³ This famous population model was originally considered by P. Verhulst in [118]. It describes the sigmoidal growth of population densities. The classical model assumes that organisms’ birth and death rates respond instantaneously to changes in population size and thus do not exhibit periodic dynamics.

oscillatory dynamics. Initiated in [64], the study of these particular DDEs explains perfectly the sustained oscillations that might appear in the evolution of a single species population. More specifically, one might show that the particular value of the delay pilots a *Hopf bifurcation* between a regime characterized by a stable steady state (small delays) and a regime that admits a stable harmonic solution (larger delays). Delayed logistic equations have been widely considered in the literature, with applications ranging from blowfly population evolution [50] to microbial growth [124], and their stability issues are still discussed nowadays [110].

(2) *Biology.*

Various complex processes appearing in the domains of physiology, immunology, epidemiology and neural networks are described by specific classes of DDEs. Blood pressure is well established to potentially exhibit sustained stable oscillations that are due to the delayed action of the sympathetic nervous system on the vasculature, [106]. Another important DDE, originally introduced by M. Mackey in [85], describes an autoimmune disease that causes periodic crashes in circulating red blood cells. This particular equation produces limit-cycle oscillations that are entirely due to a time-delayed negative feedback. In another context, human pupil size does not change immediately in response to a change in illumination (the delay is about 300 msec). This delayed mechanism (whose aim is to regulate the retinal flux by changing the pupil size) and the resulting sustained oscillations in the pupil area are described in [94]. Finally, in genetics, number of genes change their expression pattern dynamically by displaying stable oscillations. As shown in [67], this oscillatory behaviour can be explained by the relatively large time delay that characterizes the transport between the different cellular compartments.

(3) *Economics.*

Various economic activities (prices, inflation,...) exhibit recurring fluctuations over time. This leads to the idea that business cycles are actually self-sustained oscillations due to the joint effect of nonlinearities and delayed reactions of economic factors, [41]. In 1935, M. Kalecki gave the first detailed analytical study of such periodic business cycles appearing in economics, [68]. The key feature of this work is the introduction of a time lag between the moment an investment decision is taken and the moment investment goods are installed. Since the pioneer works of M. Kalecki, the consideration of systems of DDEs as well as the consecutive emergence of limit-cycle solutions have progressively become classical features in economics. In this regard, many recent contributions (see [20, 122] among others) study business dynamics affected by both lagged effects and nonlinearities.

(4) *Traffic.*

The never-ceasing activity in the literature devoted to transportation sys-

tems is proportional to the great impact that these systems have in our every day lives. It is classical in traffic flow theory to include a delay due to the driver reaction times. As shown by the car-following model considered in [36], a delay of 1 second, which is typical for most drivers, leads to damped oscillations in the velocity of a vehicle encountering a slower circulating entity. Accordingly, this kind of oscillatory instability is also encountered in the spacing between two vehicles after the leading one reduces or increases its speed.

(5) *Mechanical Engineering.*

In the context of fabrication and freight-transfer, it is likely that cranes lift several hundred of tons and it is thus important that the payload is moved rapidly and smoothly, avoiding too large oscillatory moves that could make the operator lose any control of it. In this regard, a time-delayed feedback control mechanism has been developed, in particular in [57], in order to reduce the oscillations possibly emerging from the operator manipulations. As it is shown in [41], the implementation of such a control mechanism is not without any danger. Indeed, while the oscillations are damped for small perturbations around the equilibrium position, the time-delayed feedback control mechanism leads to the creation of sustained oscillations for larger initial perturbations.

(6) *Chemical Processes.*

Experiments show that oscillatory phenomena occur in the production process of 1,3-propanediol (1,3-PD), an important industrial chemical. During this process, which consists of the fermentation of glycerol by microbial, there exists a *lag phase* (which is due to a large change of environmental conditions and metabolic response of cells) between the moment glycerol is inserted and the time the chemical reaction effectively starts. In [84], the production process of 1,3-PD is modelled with a system of DDEs and the authors thus conclude that time delay is probably one of the major reasons for which the microbial population, and consequently the amount of produced 1,3-PD, oscillates. From a purely industrial point of view, these oscillations are often ignored since they do not affect fundamentally the chemical reaction process.

(7) *Telecom.*

As cited previously in this chapter, contribution [42] highlights and studies (with the help of a flow model) the possible emergence of queue content sustained oscillations in the context of telephone switching systems. Moreover, for certain parameter ranges, this particular class of queueing systems might also exhibit a chaotic behaviour. In this context, the authors see the possible resulting oscillatory behaviours (which is a consequence of the joint effect of nonlinearities in the network and intrinsic delay phenomena) as a limiting factor, since it can significantly lower the real-time capacity of the switching system, and they propose a new service discipline

to reduce them. Along the same lines, but in the framework of multi-agent distributed computing systems, contribution [74] exposes how delayed status information can lead to inefficient resource use oscillations and how the utilization of misinformation might be used to damp this emerging oscillatory behaviour.

Note that the two major differences between the models developed in this thesis and the former contributions mentioned above in this brief literature review are:

- (i) While the time delays will possess in the present work the particularity to be self-induced by the circulating agents themselves (*i.e.* a time gap elapses between the moment an agent begins to observe the current state of the system in order to make a routing decision and the moment it will apply this routing choice and hence modify the future system evolution), they are mainly produced by external sources in the papers cited in this section.
- (ii) As exposed in Chapter 3, the intrinsic time delays will actually be variable in this work (*i.e.* the time lag between the agents' observation of the system state and their resulting feedback action varies periodically over time, in function of the queue length). Hence, the inherent equations governing the dynamics of our systems are actually DDEs with variable (random) and state-dependent delays. The theory regarding such type of equations is still far to be complete and lots of attention has been recently carried on this topic in the dedicated literature.

In Appendix 12, we give an illustration of how the presence of time delays in the dynamics of queueing systems might, also in this context, leads to oscillations phenomena. This illustrative example, observed in a real-life context and hence directly derived from a typical service management situation, helps to get an intuitive understanding of the intrinsic mechanisms leading to such a periodic behaviour.

**Autonomous Agents in a Single-Server
Feedback Queueing System - Modelling a
Recurrent Service**

A First Model - Single-Server Queuing System with Feedback Loop

Summary. *With the intention of modelling recurrent services, we consider the dynamics of a single-server feedback queueing system where the circulating items are autonomous agents able to take individual routing actions. The decision of an agent to use the feedback loop or not is based on its personal waiting time in the system. The joint action of the nonlinearity (i.e. the feedback loop) and of the non-Markovian character of the dynamics (i.e. the agents' autonomous history-based routing mechanism) causes the emergence of a collective temporal structure which takes the form of a periodic purging of the queue content. We show that this agent-induced cyclostationary behaviour fundamentally differs from the stationary state that would be obtained following a classical Markovian analysis. For regimes where the law of large numbers holds, the emerging self-organized dynamics becomes quasi-deterministic and can be analytically discussed. In addition to this analytical tractability, it is remarkable that the class of models considered in this chapter reveals itself to belong to the realm of complex adaptive systems, being thus one of the very rare solvable occurrences of such systems existing in the available literature.*

3.1 Introduction

Whatever the type of service, waiting in queues before being attended reduces the utility perceived by the customers. As time is valued for both the server and the customers, the complex relationships between the waiting times and the consumers' satisfaction is a central topic in marketing (see [17] and the references therein). Accordingly, when for a given task, competing facilities are available, a server's ability to reduce the actual or perceived waiting time of incoming customers increases its attractiveness and may drastically modify the market sharing proportions (situations with two competing servers will be presented later in this work). For example, lowering the service duration, to enhance customers' satisfaction, does require investments which have to be counterbalanced by an extra inflow of new customers.

When the service is used recurrently by customers, the utility they have per-

ceived during their last visit plays a major role in their decision to remain loyal to a service provider. Often in recurrent service systems, the cost of retaining an existing customer is comparatively less than the cost of acquiring a new one and hence the customers' loyalty is a central issue in optimizing gains. Our modelling approach has been stimulated by several recent (yet mostly experimental) studies of systems where recurrent service requirements occur. Among them, actual situations ranging from medical care [17] to leisure and hospitality facilities such as fast-food restaurants [79], ski resorts [104] and theme parks [69, 71, 72] are some well-known illustrations.

In this chapter, we consider the simplest possible topology of a queueing system with feedback loop, namely a single server with possibility of reinjection. Indeed, this configuration allows for the modelling of a general recurrent service. The items circulating in this network will be endowed with an elementary form of intelligence that enables them to individually decide either to stay in the system for another service or to leave. More particularly, the items will base their routing decisions on their elapsed sojourn times to receive service. The presence of such type of history-based (HB) routing policies (*i.e.* mechanisms in which memory enters) within our particular queueing network (QN) confers to the dynamics a manifest non-Markovian character, precluding thus the use of classical QN theory to describe the dynamics. Beside that, the individual routing decisions taken autonomously by the circulating items (humans or "smart parts") give to this queueing mechanism the basic feature characterizing multi-agent systems. Indeed, we are manifestly here in a situation where a service facility is visited by a population of autonomous decision-making agents individually assessing their situation and making decision on the basis of a HB rule.

The chapter is organized as follows. In Section 3.2, we describe the single-server feedback queueing model considered throughout this chapter. In particular, we introduce the HB routing rule, based on individual elapsed waiting time, that will govern the behaviour of the autonomous agents circulating in the network. In Section 3.3, we first follow classical QN theory techniques and describe the corresponding stationary solution. We show in Section 3.4 that such a classical approach does not capture an essential feature of the emerging dynamics. Indeed, a more careful micro-scale analysis of the self-organized flow dynamics that results from the agents' stigmergic interactions show that the queue content exhibits quasi-deterministic oscillations. We give an analytical description of this oscillatory behaviour. As motivated in Section 3.5, the oscillations observed for our feedback queueing system reproduce the signature of so-called relaxation oscillators and it follows that our model actually exhibits an agent-induced limit-cycle. In Section 3.6, we draw some relationships with other well-known limit-cycles, namely the Mathews-Lakshmanan oscillator and the periodic dynamics resulting from the implementation of reorder point policies in supply chains. In Section 3.7, we emphasize explicitly that

centralized and decentralized controls might give rise to distinct dynamics in QNs. In Section 3.8, some assumptions of our basic modelling framework are relaxed. Despite these relaxations, the main features of the emerging collective dynamics are preserved. More particularly, a delay between the time a customer leaves service and the time he/she asks to be served once more is introduced in Section 3.8.1. Then, the customers' patience parameter, which controls their HB routing behaviour, is heterogenized in Section 3.8.2. As emphasized in Section 3.9, the type of dynamics exposed in this chapter is not only restricted to social (human) systems but might also be observed in any "smart parts" logistics network thanks to the nowadays wide availability of RFID tags. In Section 3.10, we justify that the class of models considered in this chapter can be viewed as particular cases of complex adaptive systems. Finally, Section 3.11 is devoted to conclusions and perspectives.

3.2 Model

The simplest possible network composed of a single queue with the presence of a feedback routing node is sketched in Fig. 3.1. An incoming flow of customers,

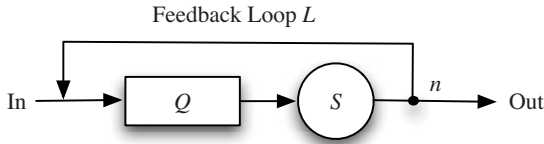


Fig. 3.1. A single-stage queueing system with feedback loop.

described by a renewal process with mean inter-arrival time $\frac{1}{\lambda}$ and probability distribution $A(x)$ with density $dA(x)$, is served by a processing unit whose service times are i.i.d. random variables with mean $\frac{1}{\mu}$, probability distribution $B(x)$ and density $dB(x)$. Accordingly, the parameters λ and μ are respectively the incoming and service rates of the renewal processes. We assume that the distributions $A(x)$ and $B(x)$ have finite moments. Here, we suppose that the traffic intensity $\rho = \frac{\lambda}{\mu} < 1 \Leftrightarrow \lambda < \mu$, which ensures the stability of the queueing system in absence of feedback loop. Assume also that the waiting room capacity is unlimited and that the service discipline is first-in-first-out (FIFO). After being served, the routing of each customer at the QN decision node n will be either

- (i) to leave the system definitively, or
- (ii) to follow the feedback loop and line up again to be served once more.

Several well-known contributions [35, 101, 112] consider the classical situation arising when the decision between the choices (i) and (ii) is taken randomly. When this is the case, by imposing a stationary flow balance (*i.e.* incoming equals outgoing flow), the system is driven into a self-consistent stationary regime (which can be described by time-independent distributions). As we will see in Section 3.4, such stationary flows strongly differ from the queue dynamics observed when “intelligent” agents, able to base their routing on historical data (here the time spent while queueing and being served), circulate in the network. Specifically, assume now that each agent is able to record the total waiting time W he/she spent to receive service (*i.e.* W is the sum of the queueing and processing times; in queueing theory, W is commonly known as the sojourn time). Assume further that W controls the routing decision between the alternatives (i) and (ii), namely the following history-based (HB) routing rule \mathcal{R} is implemented:

$$\mathcal{R} = \begin{cases} \text{follow alternative (i)} & \text{if } W > P, \\ \text{follow alternative (ii)} & \text{if } W \leq P, \end{cases} \quad (3.1)$$

where P is called the patience parameter of the agents in circulation. When alternative (ii) is chosen, we will speak of loyal customers, as the agents are pleased with the server and then return to it for another service. Note that the flow of customers taking the feedback loop is added to the flow of fresh customers entering the system (with incoming rate λ). We assume that, when joining the queue, a loyal customer behaves as a fresh customer (*i.e.* only the last recorded sojourn time will be determining for its routing at node n). Note that the routing is now clearly HB - it is determined by the sojourn time that each customer spent in order to be served. The underlying idea behind this type of model is to study how, by modifying the quality of service (here the sojourn time), it is possible to enhance the circulation of loyal customers (*i.e.* those taking the feedback loop). In the following, we will for a start focus on homogeneous agents for which P is a common value. This assumption will thereafter be relaxed in Section 3.8.

3.3 Classical Analysis - Assuming Markovian Evolution

We first follow in this section classical QN techniques and we proceed to the analysis of the system described in Section 3.2 by assuming a Markovian evolution as well as the existence of a stationary regime. But as we will see further in Section 3.4, while such analysis could represent a satisfactory mean-field macroscopic approximation¹, it however does not capture an essential feature of the emerging dynamics. Indeed, we will see that the considered multi-agent

¹ The mean-field approach consists in replacing the agents’ numerous interactions by an effective external field. Like this, a many-body problem is reduced into a one-body problem.

behaviour leads not to the convergence to a pure stationary regime but, due to the actual non-Markovian character of the dynamics and the presence of an intrinsic delay mechanism, yields the emergence of a periodic evolution of the queue content. As described later in Section 3.7, it is furthermore interesting to note that the stationary regime that we will compute in the present section would actually appear if the management (*i.e.* the routing of the items) was fully centralized to the server - this clearly differs from the periodic dynamics that emerges when the decision-making is decentralized to the agents (see Section 3.4).

Regarding QN terminology, the dynamics induced by our particular routing rule \mathcal{R} implies that we are here in presence of a feedback queueing system with state-dependent reentering flow. Note that, while state-dependent QNs are abundantly studied in the literature, little attention has been devoted so far to the present type of state-dependent feedback queueing systems. Following [60], we define λ_P as the flow rate of agents taking the feedback loop and λ_{out} as the rate at which the agents leave the system (remember that λ denotes the arrival rate of new customers). We write the probability density function of the waiting time W as:

$$\text{Prob}\{x \leq W \leq x + dx\} = \pi_W(x)dx.$$

Define now $\gamma(P) \in [0, 1]$ as the proportion of agents who remain in the system at decision node n (*i.e.* those taking the feedback loop):

$$\gamma(P) = \text{Prob}\{0 \leq W \leq P\} = \int_0^P \pi_W(x)dx.$$

Obviously, $\gamma(P)$ is an increasing function of P (*i.e.* the proportion of remaining agents raises with P). Note that, in the limit case where $P \rightarrow 0$, the behaviour of the feedback queueing system will tend to the one of an ordinary open queue.

To simplify the presentation and to allow for a full analytical description, we restrict in the following to exponentially distributed inter-arrival and service times², *i.e.*

$$A(x) = 1 - e^{-\lambda x} \quad \text{and} \quad B(x) = 1 - e^{-\mu x}.$$

In view of this assumption, we recall that the following properties directly hold, [93]:

- (1) Arrivals and services occur according to Poisson processes.
- (2) The process resulting from the multiplexing or the parting of Poisson processes is itself a Poisson process.

² A further discussion for general distributions can be found in [60].

(3) In the stationary regime, the departure process of a $M/M/1$ queue with a waiting room of infinite capacity is a Poisson process with the same rate as the input Poisson process.

Using this last property and assuming the existence of a stationary regime, we may write

$$\lambda = \lambda_{\text{out}}. \quad (3.2)$$

Following properties (2) and (3) above, the flow that enters decision node n is a Poisson process with intensity $\lambda + \lambda_P$ and we can write:

$$\lambda_P = \gamma(P) (\lambda + \lambda_P). \quad (3.3)$$

Combining Eqs. (3.2) and (3.3), we get:

$$\lambda = \lambda_{\text{out}} = (1 - \gamma(P)) (\lambda + \lambda_P).$$

For an $M/M/1$ queue with arrival rate $(\lambda + \lambda_P)$ and service rate μ , it is well-known that the probability density function $\pi_W(x)$, characterizing the waiting time W in the stationary regime, is given by:

$$\pi_W(x) = (\mu - (\lambda + \lambda_P)) e^{-(\mu - (\lambda + \lambda_P))x}. \quad (3.4)$$

In view of Eq. (3.4), the branching ratio $\gamma(P)$ hence satisfies:

$$\begin{aligned} \gamma(P) &= \int_0^P (\mu - (\lambda + \lambda_P)) e^{-(\mu - (\lambda + \lambda_P))x} dx \\ &= 1 - e^{-(\mu - (\lambda + \lambda_P))P}. \end{aligned} \quad (3.5)$$

Using Eqs. (3.3) and (3.5), we directly get the following transcendent equation:

$$\lambda_P = \left(1 - e^{-(\mu - (\lambda + \lambda_P))P}\right) (\lambda + \lambda_P). \quad (3.6)$$

Hence, given the set of control parameters (λ, μ, P) , the solution of Eq. (3.6) determines the effective feedback rate λ_P in the stationary regime. An illustration of how the parameter P governs the value of λ_P is given in Fig. 3.2. Without giving rigorous proofs, we observe the following facts:

- (1) The feedback flow rate λ_P is increasing with P (*i.e.* the more patient are the agents, the more they are likely to take the feedback loop).
- (2) $\lambda_P < \mu - \lambda \Leftrightarrow \rho_{\text{tot}} = \frac{\lambda + \lambda_P}{\mu} < 1$, where ρ_{tot} is the total traffic rate (*i.e.* the routing rule \mathcal{R} drives the system into a stable state).

Introducing the input traffic rate $\rho = \frac{\lambda}{\mu}$, we can rewrite Eq. (3.6) as:

$$\rho = \rho_{\text{tot}} e^{-\mu P (1 - \rho_{\text{tot}})} =: F(\rho_{\text{tot}}). \quad (3.7)$$

It is immediate to verify that $F(\rho_{\text{tot}})$ is monotonically increasing, with $F(0) = 0$ and $F(1) = 1$. Hence, a single effective traffic value $\rho^* \in [0, 1]$ solves the

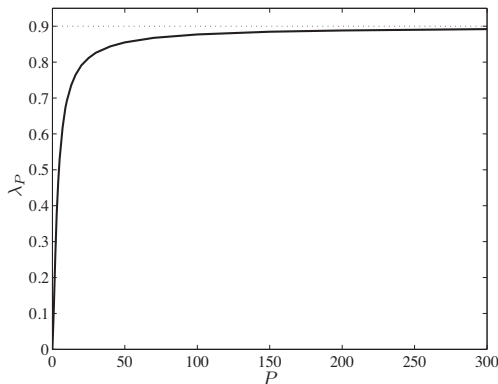


Fig. 3.2. Flow rate λ_P in the feedback loop in function of the value of the parameter P , when $\lambda = 0.1$ and $\mu = 1$.

transcendent Eq. (3.7). According to this solution, the queue content will ultimately fluctuate around a time-independent average

$$\bar{Q}_{\text{stat}} = \frac{\rho^*}{1 - \rho^*}. \quad (3.8)$$

Fig. 3.3 illustrates how the input traffic rate ρ influences the effective traffic rate ρ^* that actually feeds the server. Note that the solution ρ^* , which assumes the existence of a stationary regime, does not actually take into account the agent character of the dynamics and is thus the solely result of a *rate (i.e. flow) analysis*. In the next section, we reconsider the same feedback queueing system from a microscopic level point of view (*i.e.* we explicitly take into account the agents' local numerous routing actions) and we show that the type of stationary solutions computed above is missing key features of the emerging dynamics.

3.4 Micro-Scale Agent-Based Analysis - Siphon Dynamics

Let us now show that the Markovian analysis described in the preceding section is missing an essential aspect of the system behaviour. Indeed, as discussed in [43], when the circulating items apply the HB routing rule \mathcal{R} stated in Eq. (3.1) and when P is large enough, quasi-deterministic cyclo-stationary regimes emerge, *i.e.* stable temporal oscillations of the queue length $Q(t)$ are observed and this independently of the detailed nature of the probability laws $A(x)$ and $B(x)$. Actually, despite the presence of fluctuations, this robust and

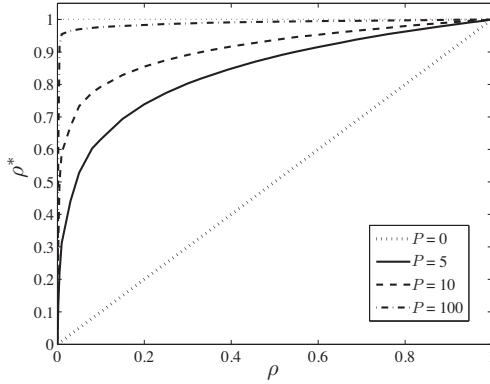


Fig. 3.3. Effective traffic ρ^* in the feedback queueing system in function of the input traffic ρ , when $\mu = 1$.

quasi-deterministic behaviour is directly reminiscent from the law of large numbers (LLN). The importance of the relative fluctuations around the associated average sojourn time $\langle W \rangle$ (which is the sum of individual processing times) decreases for large queue content $Q(t)$ (a quantitative characterization is given in [43]). Accordingly, for large P , the dynamics can be discussed via a deterministic approach (involving a constant service time $\frac{1}{\mu}$), [43, 60]. Hence, for a given queue length Q_c and a given corresponding patience parameter $P = \frac{Q_c}{\mu}$, an incoming tagged customer ζ lining behind Q_c other customers, will, when reaching node n , choose the alternative (i) (*i.e.* to leave the system). Indeed, for such a deterministic regime, the measured total sojourn time $W = \frac{Q_c}{\mu} + \frac{1}{\mu} > P$. However, before ζ makes its way through the queue and reaches the node n , the queue content $Q(t)$ still increases at the (deterministic) rate λ (as nobody leaves the system during this time interval), implying a delay mechanism in the draining of the queue content. As soon as ζ reaches n , and thus leaves the system, a second dynamical phase is triggered. During this second phase, the customers arriving immediately after ζ do also experiment a waiting time exceeding P and will hence also leave the system. As $\lambda < \mu$, the queue population $Q(t)$ decreases during this second dynamical phase and the depletion lasts until a satisfied customer (and hence his/her immediate successors) reach the node n . When this happens, the first dynamical phase starts again and $Q(t)$ fills up at rate λ . The alternation between these two dynamical phases produces a cyclo-stationary behaviour whose very existence is entirely due to the elementary “intelligence” attached to the circulating agents. Elementary “intelligence” refers here to the agents’ capability to monitor and to memorize the time while queueing and to take an individual routing decision

accordingly. It is enlightening to visualize the queue dynamics by using the hydrodynamic analogy sketched in Fig. 3.4. Indeed, one can convince oneself

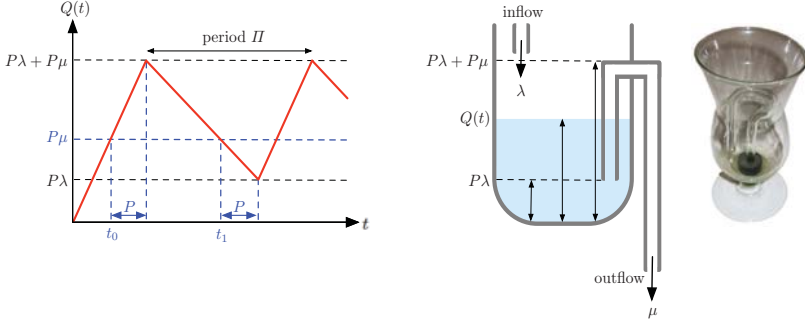


Fig. 3.4. Hydrodynamic analogy. *Left:* The agent entering at t_0 is the first one of a whole cluster U of unsatisfied customers and triggers the alternation of $Q(t)$ from the increasing to the decreasing state at $t_0 + P$. The last agent belonging to the cluster of unsatisfied customers U is the one entering just before t_1 and triggers the switch of $Q(t)$ from the decreasing to the increasing state at $t_1 + P$. This simple delay dynamics repeats and creates stable oscillations. *Right:* The “Tantalus glass” siphon model. The queue length corresponds to the water level $Q(t)$. The inflow and outflow rates are λ respectively μ . The siphon leaves a water residue of height $P\lambda$ due to the constant inflow during P . The effective siphon length is $P\mu$.

that the time-dependent queue content level is fully analogous to the liquid oscillations arising in a self-siphoning³ “Tantalus glass” (sketched in Fig. 3.4.*Right*). This simple hydrodynamic system, also considered and described in [99, 100], reveals itself to belong to the well-known class of relaxation oscillators (see Section 3.5). In addition, for large P , the purely deterministic context ensuing from the LLN enables an elementary derivation of both the amplitude Δ and the period Π of the queue population $Q(t)$. Following [43] for further analytical details, we obtain (see also Fig. 3.4):

$$\Delta = P\mu, \quad (3.9)$$

$$\Pi = P \left[2 + \frac{\lambda}{\mu - \lambda} + \frac{\mu - \lambda}{\lambda} \right] \quad (3.10)$$

and provided $P \gg \max\left(\frac{1}{\lambda}, \frac{1}{\mu}\right)$ both Eqs. (3.9) and (3.10) are in perfect agreement with simulation experiments (see Fig. 3.5), as discussed in [43, 60],

³ Note that siphon effects have also been considered in the context of Petri nets (see [32] among others). However, contrary to the effect occurring in the framework of Petri nets, which leads to infinite delays (“*any empty siphon remains empty*”), our particular queue dynamics deals with cyclo-stationary siphon effects.

and this for any possible choice of the probability distributions $A(x)$ and $B(x)$. As shown in Fig. 3.5, the influence of the LLN explicitly grows as P increases and causes the curves to become smoother and (quasi-)deterministically pe-

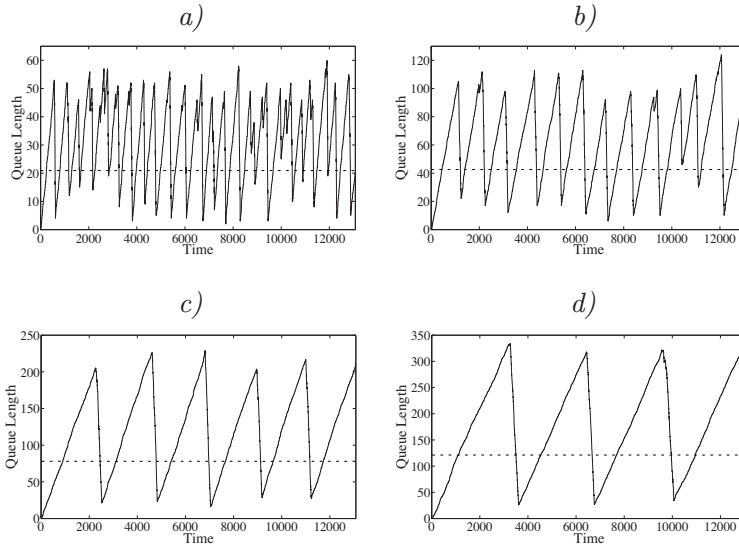


Fig. 3.5. Queue length oscillations obtained by simulation for exponentially distributed inter-arrival and service times with $\lambda = 0.1$, $\mu = 1$ and *a*) $P = 50$, *b*) $P = 100$, *c*) $P = 200$, *d*) $P = 300$. These simulated behaviours (and more precisely the influence of the law of large numbers, which is increasing with P) remain qualitatively the same for any other possible inter-arrival and service times probability distributions. In each of the four graphs, the dotted line represents the stationary queue content given by Eq. (3.8), that would be predicted by the purely Markovian analysis described in Section 3.3.

riodic with growing P . This can be observed explicitly when looking at the spectrum component (*i.e.* the Fourier transform) of the queue dynamics for increasing values of P (see Fig. 3.6).

It can furthermore be emphasized that, as illustrated in Figs. 3.4 and 3.5, the emerging periodic structure of the queue dynamics directly implies that the queue content $Q(t)$ is above bounded (by $P\lambda + P\mu$) and also never vanishes ($Q(t) > 0, \forall t$). The queue length self-organized behaviour thus ensures:

- the stabilization of the queueing system (which results from the periodic purging of the queue),

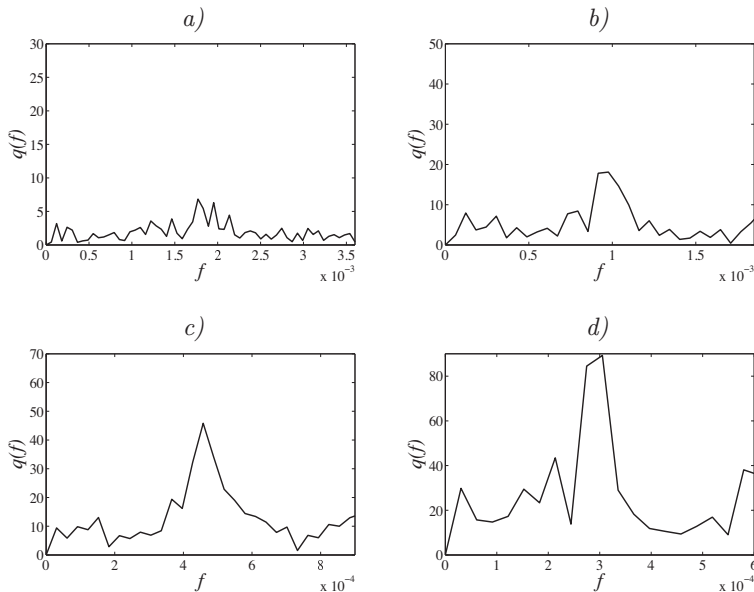


Fig. 3.6. Spectrum of the queue content dynamics obtained by simulation for exponentially distributed inter-arrival and service times with $\lambda = 0.1$, $\mu = 1$ and a) $P = 50$, b) $P = 100$, c) $P = 200$, d) $P = 300$.

- a maximum busy period (*i.e.* equal to 1) for the server and hence resource use maximization.

Such self-organization phenomena are solely induced by the autonomous agents' individual actions, a feature which is typical for complex adaptive systems (see Section 3.10). Note finally that, as pointed out at the beginning of this section and as illustrated in Fig. 3.5, the emerging cyclo-stationary regime observed here clearly differs from the pure stationary evolution that follows from the solely Markovian analysis described in Section 3.3.

3.5 Multi-Agent Induced Limit-Cycle - Stigmergic Relaxation Oscillator

As it is underlined in [102], the hydrodynamic self-siphoning device illustrated in Fig. 3.4 belongs to a particular class of self-sustained systems made popular by the early works of B. van der Pol [116] and generally referred as the class of *relaxation oscillators*. These systems are in essence characterized by the presence of two time scales in their dynamics. More particularly, each cycle

is composed of intervals of slow and fast motion. These different velocities define the two distinct phases of the oscillations: accumulation and “firing” (see Fig. 3.7). Consequently, the style of the oscillations in a relaxation oscillator is a sequence of pulses and thus clearly differs from a simple sine wave.

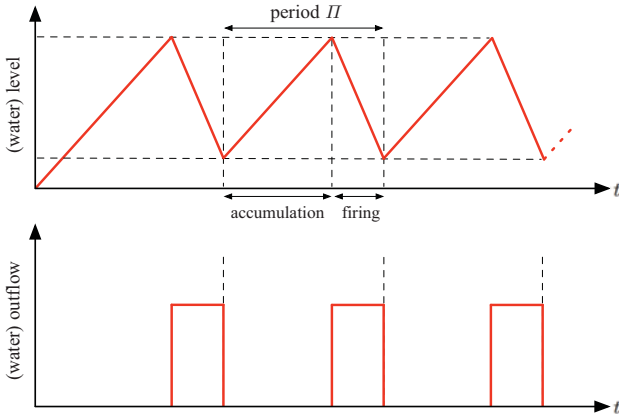


Fig. 3.7. General time evolution of a relaxation oscillator.

The complete analogy existing between the queueing and siphon dynamics exposed in the preceding section (see Fig. 3.4) explicitly suggests that our particular feedback queueing system belongs also to the class of relaxation oscillators. Furthermore, the multi-agent aspect of the dynamics enables us to speak here of *stigmatic relaxation oscillator* exhibiting an *agent-induced limit-cycle*.

The *accumulate-and-fire* dynamics of relaxation oscillators is of course not restricted to the hydrodynamic device presented previously in this work. Indeed, as unveiled in [102], this particular dynamics captures the main features of several various real systems. A first widely famous example is the electronic generator considered by B. van der Pol while he was working, in the 1920s, as a Head Physicist at Philips Physical Laboratory in Eindhoven [116]. This oscillator, which operates in a very similar way, is composed of a battery, a capacitor, a resistance and a neon tube. First, the capacitor is being charged (the capacitance and resistance determining the characteristic time of this phase); the growth of the voltage corresponding to the increase of the water level in the siphon model. When the voltage reaches a critical threshold value, electric conduction is initiated in the neon tube, the capacitor consequently quickly discharges through the lamp and its voltage drops until the tube becomes

nonconductive again. This process repeats itself endlessly. Like the outflow of water in the hydrodynamic device, the flow of electric current through the lamp is a sequence of short pulses (and therefore, the lamp exhibits periodic flash lights). Such generators, originally used in the framework of radio communication, are still used nowadays to build oscilloscopes, cathodic TV sets or computer displays.

Relaxation oscillations are encountered in many biological systems, like spontaneous firing neurons or heartbeat [102], for which we also observe limit-cycles composed of an increasing phase (that can be linear or not) followed by a resetting at a threshold. The mechanism that takes place in firing neurons (a constant current is injected into the cells and produce spikes periodically) is very similar to the one of the van der Pol generator (see [102] for more details). Among other examples, in the case of sensor networks, the firing rate determines the intensity of the perceived stimulus. Concerning the heartbeat, the whole cardiovascular system can be seen as a single oscillator, [117]. While this system might be affected by other physiological rhythms, including respiration, these external rhythms do not fundamentally produce the heartbeat but rather disturb it (indeed, it can be experimented that an isolated heart preserves its capability of periodic contraction *in vitro*). Despite these external perturbations, the heartbeat is able to produce its rhythm by itself and can hence be considered as a self-sustained oscillator, [102]. The heartbeat being the macroscopic result of the interactions between the numerous entities forming the cardiovascular system, it might moreover be seen as a multi-agent limit-cycle. Furthermore, due to its ability to adapt to fluctuating environments, the cardiovascular system can be naturally classified as a complex adaptive system (see Section 3.10).

3.6 On Relationships with Well-Known Stable Limit-Cycles

In this section, we first show that the agent-induced oscillations described in this chapter share the similarity with the Mathews-Lakshmanan oscillator to exhibit an amplitude-dependent frequency. Secondly, we draw a relationship between the periodic behaviour emerging in our context and the self-sustained oscillations of the inventory level that appear, in the domain of supply chains, when reorder point policies are implemented.

3.6.1 The Unique Nonlinear Oscillator - Amplitude-Dependent Frequency

The Hamiltonian of a classical (undamped) simple harmonic oscillator is given by

$$H(x, y) = \frac{1}{2}x^2 + \frac{1}{2}y^2,$$

and, accordingly, the associated dynamics satisfies the following Hamilton equations:

$$\begin{cases} \dot{x} = y \\ \dot{y} = -x. \end{cases}$$

The simple harmonic motion is hence described by the following equation:

$$x(t) = \Delta \sin(t + \phi_0),$$

which consists in sinusoidal oscillations about the equilibrium point (the amplitude Δ and phase ϕ_0 are determined by the initial conditions). Note that these limit-cycle oscillations have a constant amplitude and a constant frequency. Furthermore, the frequency does not depend on the amplitude, which is typical for linear systems. This last property implies that the simple harmonic oscillator would make a perfectly robust clock, since possible random perturbations in the amplitude would not affect its frequency (and period).

Now, we consider the following Hamiltonian:

$$H(x, y) = \log(\cosh(y)) + \frac{1}{2}\log(1+x^2),$$

which is associated to the *unique nonlinear oscillator* introduced by Mathews and Lakshmanan in [89]. The corresponding Hamilton equations

$$\begin{cases} \dot{x} = \tanh(y) \\ \dot{y} = -\frac{x}{1+x^2} \end{cases}$$

clearly exhibit here a highly nonlinear pattern. Direct calculations show that this system admits the periodic solution

$$x(t) = \Delta \sin(\omega(\Delta)t + \phi_0), \quad \omega(\Delta) = \frac{1}{\sqrt{1+\Delta^2}}.$$

Despite the strong nonlinear character of the Mathews-Lakshmanan oscillator, it is remarkable that its bounded periodic motion reveals itself to be simple harmonic. Moreover, the oscillations exhibit an amplitude-dependent frequency which is the direct signature of the underlying nonlinearity of the dynamics. Note that the accuracy of the clock associated with this oscillator might suffer from this property since its frequency would vary with randomly fluctuating amplitude.

The agent-induced limit-cycle oscillations studied in Section 3.4 are in essence similar to those of the Mathews-Lakshmanan oscillator since they also possess an amplitude-dependent frequency (which is due in this case to the agents' nonlinear feedback rule). For both oscillators, the frequency is monotonically

decreasing with growing amplitude. For our multi-agent queueing system, we find, using Eqs. (3.9) and (3.10), that the amplitude Δ and the period Π (respectively the frequency f) of the queue content oscillations are connected with the following equations:

$$\Pi = \frac{\Delta\mu}{(\mu - \lambda)\lambda} = \frac{\Delta}{(\mu - \lambda)\rho} \iff f = \frac{1}{\Pi} = \frac{(\mu - \lambda)\rho}{\Delta}. \quad (3.11)$$

Following Eq. (3.11), the dependency between the amplitude and the period is linear and consequently the ratio \mathcal{K} between these two quantities

$$\mathcal{K} = \frac{\Delta}{\Pi} = (\mu - \lambda)\rho$$

is independent of P .

3.6.2 Supply Chains - Reorder Point Policy

In any supply chain, stocks of items must be held to constantly meet future random demand. The management of these inventories is not trivial since there always exists a time lag (called lead time in this context) between the date of placing an order for material and the date on which the items are actually received. Among the numerous existing management procedures, the *reorder point policy* represents a classical and convenient way to control the evolution of inventories, [33]. The reorder point (ROP) refers to the level of inventory at which a fresh order of material must be placed to refurbish the stock. Taking explicitly into account the existence of the lead time, the ROP is chosen such that new items will arrive before the firm runs out of stock. It is usually computed with the following formula:

$$\begin{aligned} \text{Reorder Point} &= \text{Mean Consumption per Unit of Lead Time} \times \text{Lead Time} \\ &+ \text{Safety Stock,} \end{aligned}$$

where the safety stock is a minimum inventory level that acts as a protection against shortages due to fluctuating demand. The determination of the safety stock involves a trade-off between the risk of stock-out (which results in lost sales due to customer dissatisfaction) and the increased costs associated with carrying additional inventory. Classical ROP policies assume a fixed order quantity that is often referred as the Economic Order Quantity (EOQ).

As illustrated in Fig. 3.8, the implementation of a ROP policy leads not surprisingly to a stable oscillatory evolution of the inventory level, the amplitude of which is equal to the EOQ. A stationary distribution for the demand is of course an essential condition for the stability of the oscillations. The limit-cycle observed in this context is somehow very similar to the agent-induced limit-cycle described in the preceding sections. Just like the effect induced

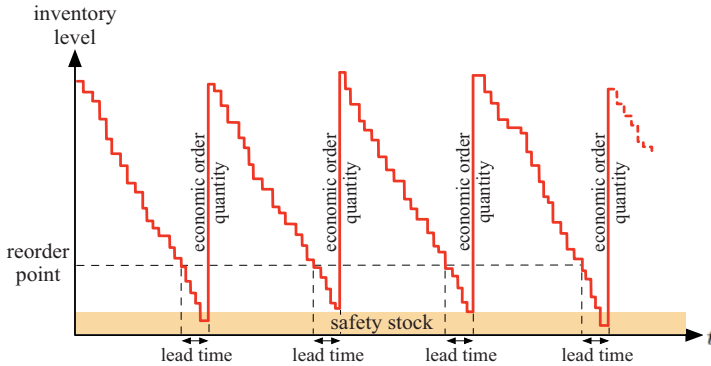


Fig. 3.8. Inventory level dynamics when a reorder point policy is implemented.

by the patience parameter P in our feedback queueing system, here the oscillations are getting smoother with increasing EOQ. Note that contrary to the fully decentralized agent-based mechanism exposed earlier in this chapter, here the stable limit-cycle is fully induced by a central controller that places orders when the reorder point is reached.

3.7 Centralized Versus Decentralized Control to Achieve Queue Length Stability - Distinct Dynamics

To achieve system stability (*i.e.* to ensure that the queue length does not ultimately explode) within the QN with feedback loop considered throughout this chapter, one could actually rely either on centralized or on decentralized mechanisms. It is interesting to observe that, in the present context, these two antagonistic management approaches would lead to distinct system dynamics.

For one thing, we remark that a management fully centralized to the server would in fact drive the system into the stationary regime described previously (by assuming a Markovian evolution) in Section 3.3. Indeed, let us consider the situation where the server solely determines the items' routing choices whether to take or not the feedback loop. More particularly, the server observes, $\forall t$, the current queue content $Q(t)$ and determines accordingly the ratio of items

$$p(t) = \mathbb{P}(W \leq P \mid Q(t))$$

that are likely to be unsatisfied with service at that time (*i.e.* when their waiting time W exceeds their patience parameter P). As already emphasized in this chapter, such a purging of unsatisfied entities leads to the stabilization of the system. Under the above type of centralized management, an item

entering decision node n (see Fig. 3.1) at time t would be commanded with probability $p(t)$ to take the feedback loop and with probability $1 - p(t)$ to leave the system. Without calling here for rigorous proofs, it is possible to understand easily and intuitively how the system will evolve in this situation. Starting from an empty queue, the proportion $p(t)$ of items taking the feedback loop will decrease over time from 1 to ultimately reach the ratio $\gamma(P)$ (see Section 3.3) that will be effective in the stationary regime. At the same time, the queue length $Q(t)$ will progressively increase until it finally evolves around the time-independent average \bar{Q}_{stat} given by Eq. (3.8) (see Section 3.3).

As it has already been extensively studied in Section 3.4, when an agent-based fully decentralized mechanism is implemented within our feedback queueing system (*i.e.* all routing decisions at node n are taken autonomously by the circulating entities), we observe the emergence of stable oscillations of the queue content, which ensures intrinsically a bounded queue length and hence the stabilization of the system.

To summarize, it is very enlightening to compare the behaviours arising on one hand when the system control is solely left to the server and when it is on the other hand ensured by the circulating entities themselves. Indeed, in the present context, in order to achieve the same ultimate objective (*i.e.* to stabilize the system thanks to the departure of unsatisfied agents), the dynamics would be driven into a purely stationary regime for centralized management while a cyclo-stationary state would emerge in case of full decentralization.

3.8 Relaxing Some Assumptions

In this section, we reconsider the agent-based feedback queueing system introduced previously in this chapter and we abandon some of the hypothesis assumed in the model respectively described and analyzed in Sections 3.2 and 3.4. We first introduce, in Section 3.8.1, a time delay between the moment an agent decides to ask for another service (*i.e.* to take the feedback loop) and the time it will effectively ask for this additional service. Then, in Section 3.8.2, we individualize the agents' patience parameter and hence consider heterogeneous agents with distinct individual reactions to suffered waiting time.

3.8.1 Delayed Feedback

While there clearly exist practical cases of recurrent services where coming back customers line instantaneously for another service (*e.g.* theme parks, ski lifts, etc.), there also exist numerous situations where a time delay occurs before a customer comes back to be served once again. This happens for example when people visit restaurants (or fast-foods, cafeterias), retail shops or

parking facilities. One might also think to commuter drivers using a traffic road to do their daily journey. As illustrated in Fig. 3.9, we consider here a

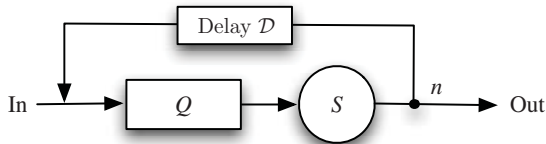


Fig. 3.9. A single stage queueing system with delayed feedback mechanism.

single server queueing system with a delayed feedback. Now, when a satisfied agent decides, at decision node n , to take the feedback loop, it will ask for another service after a time delay $\mathcal{D} \geq 0$. Beside that, all the assumptions taken in Sections 3.2 and 3.4 remain valid here.

As illustrated in Fig. 3.10, simulation experiments show that two different

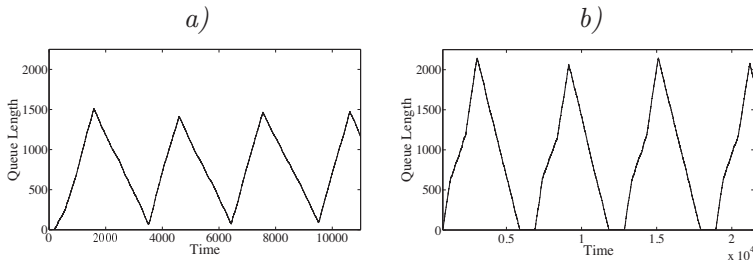


Fig. 3.10. Queue length dynamics when inter-arrival times are uniformly distributed in $[0.1, 1.5]$ with mean $\frac{1}{\lambda} = 0.8$, service times are uniformly distributed in $[0.1, 0.9]$ with mean $\frac{1}{\mu} = 0.5$, $P = 300$ and $a) \mathcal{D} = 400$, $b) \mathcal{D} = 1500$.

generic queue dynamics might emerge depending on the value of \mathcal{D} . The first situation is likely to arise in presence of small delay before the feedback is effective while the second one will emerge for larger values of \mathcal{D} .

(1) “Small” delay $\mathcal{D} \leq P\lambda / (\mu - \lambda)$.

When the delay in the feedback loop remains relatively small, the queue dynamics is qualitatively similar to the ones observed when there is no delay. The presence of the delay \mathcal{D} acts upon the length of the delay mechanism (*i.e.* the inherent siphon process) that creates the switches between the increasing and decreasing phases of the queue oscillations.

More precisely, from the moment $Q(t)$ reaches level $P\mu$, the time delay before the switch to the next oscillatory phase is here not anymore equal to P (corresponding to the situation where $\mathcal{D} = 0$), but to $P + \mathcal{D}$ (see Fig. 3.11). To get more understanding of that phenomenon, let us consider the

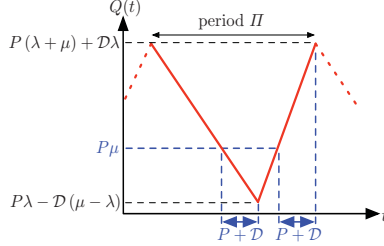


Fig. 3.11. Oscillatory queue length dynamics when the agents' feedback is delayed by $\mathcal{D} \leq P\lambda/(\mu - \lambda)$. The emerging queue dynamics is qualitatively similar to the one observed when $\mathcal{D} = 0$ (see Fig. 3.4), but the siphon delay mechanism is here of length $P + \mathcal{D}$.

agent ζ joining the queue just after $Q(t) = P\mu$, during the decreasing phase. Consequently, ζ will suffer a sojourn time equal to P and, from the moment ζ leaves service, a time delay \mathcal{D} will still elapse before ζ comes back in the queue, triggering thus the switch from the decreasing to the increasing state of the oscillations. An explanatory sketch of the emerging queue dynamics observed when $\mathcal{D} \leq P\lambda/(\mu - \lambda)$ is illustrated in Fig. 3.11. For large P (*i.e.* when the LLN influences the system behaviour by manifestly smoothing the queue length dynamics), the corresponding amplitude Δ and period Π of the quasi-deterministic oscillations are now respectively given by:

$$\Delta = (P + \mathcal{D})\mu,$$

$$\Pi = (P + \mathcal{D}) \left[2 + \frac{\lambda}{\mu - \lambda} + \frac{\mu - \lambda}{\lambda} \right],$$

which corresponds to Eqs. (3.9) and (3.10), with an inherent siphon delay mechanism of length $P + \mathcal{D}$ (instead of P). Note that the condition $\mathcal{D} \leq P\lambda/(\mu - \lambda)$ implies that $Q(t) > 0, \forall t$.

(2) “Large” delay $\mathcal{D} > P\lambda/(\mu - \lambda)$.

When the delay grows and satisfy $\mathcal{D} > P\lambda/(\mu - \lambda)$, the queue is likely to decrease until being completely empty. More precisely, the queue will remain empty during a strictly positive interval of time $T_E = \mathcal{D} - P\lambda/(\mu - \lambda)$. This period during which $Q(t) = 0$ will have a direct influence on the magnitude and length of the increasing phase observed for the queue oscillatory dynamics (see Fig. 3.12). To understand that, let us

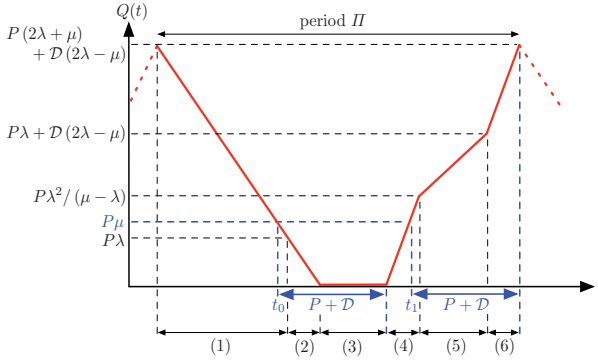


Fig. 3.12. Oscillatory queue length dynamics when the agents' feedback is delayed by $\mathcal{D} > P\lambda/(\mu - \lambda)$. *Phase (1):* During this phase, all the agents leave the system after service and hence the queue decreases at rate $\lambda - \mu < 0$. Indeed, all these agents face a queue $Q(t) > P\mu$ and thus experiment a sojourn time $W > P$. *Phase (2):* $Q(t)$ continues to decrease at rate $\lambda - \mu$ until $Q(t) = 0$ but the agents do not leave the system definitely as they will take the feedback loop after a time delay \mathcal{D} . Indeed, the agent ζ_1 who joins the queue at t_0 triggers, at time $t_0 + P$, a batch of satisfied agents who will ask for another service after a delay \mathcal{D} will have elapsed. In this phase, $\xi_{out,t} = \mu$ ($Q(t) > 0$). *Phase (3):* During this phase, the queue remains empty as ζ_1 and the following batch of satisfied agents are still in the time delay \mathcal{D} before they ask for another service. In this phase, $\xi_{out,t} = \lambda$ ($Q(t) = 0$). *Phase (4):* When customer ζ_1 joins again the queue, at time $t_0 + P + \mathcal{D}$, it triggers the switch to the increasing phase of the dynamics. Since $\xi_{out,t} = \mu$ in Phase (2), $Q(t)$ increases here at rate $\xi_{in,t} = \lambda$. The agent ζ_2 who joins the queue at time t_1 will trigger, at time $t_1 + P + \mathcal{D}$, the switch to the decreasing phase of the dynamics (a batch of unsatisfied agent will leave definitely the system). *Phase (5)* During this phase, $Q(t)$ increases at rate $\xi_{in,t} = 2\lambda - \mu$, since $\xi_{out,t} = \lambda$ during Phase (3). *Phase (6):* $Q(t)$ increases again at rate $\xi_{in,t} = \lambda$ (during Phase (4), $\xi_{out,t} = \mu$). Accordingly, we find that the amplitude and period of the oscillations are respectively given by $\Delta = P(2\lambda + \mu) + \mathcal{D}(2\lambda - \mu)$ and $\Pi = P\frac{\mu(\mu + \lambda)}{\lambda(\mu - \lambda)} + \mathcal{D}\frac{\mu}{\mu - \lambda}$. This figure corresponds to the situation where t_1 (triggering of the inherent siphon delay mechanism) happens during Phase (4). This situation arises when $\lambda^* \leq \lambda < \mu$, where λ^* is the unique solution of $f(\lambda) = \lambda^2/(\mu - \lambda) = \mu$ when $\lambda < \mu$ ($f(\lambda)$ is a strictly increasing function of λ). Note that $\lambda^* > \mu/2$, implying thus that $2\lambda - \mu > 0$ in the present situation. The other typical situations where t_1 takes place during Phases (5) and (6) are sketched in Figs. 3.13 and 3.14.

consider the rate $\xi_{out,t}$ at which the agents leave the server at time t . When $Q(t) = 0$, $\xi_{out,t} = \lambda$, as the incoming rate λ is smaller than the service rate μ (i.e. $\rho = \lambda/\mu < 1$). On the other hand, when $Q(t) > 0$, the agents leave the server at rate $\xi_{out,t} = \mu$. As we will see, this rate $\xi_{out,t}$ has an influence on how the queue will be fed after the time delay \mathcal{D} . Indeed, a satisfied agent that leaves the server at time t will join the queue at time

$t + \mathcal{D}$. The rate $\xi_{in,t}$ at which the queue evolves at time t is given by:

$$\xi_{in,t} = \lambda - \mu + \xi_{out,t-\mathcal{D}}.$$

Depending on the value of $\xi_{out,t-\mathcal{D}} \in \{\lambda; \mu\}$, $\xi_{in,t}$ can take the two following values:

$$\xi_{in,t} = \begin{cases} \lambda & \text{when } \xi_{out,t-\mathcal{D}} = \mu, \\ 2\lambda - \mu & \text{when } \xi_{out,t-\mathcal{D}} = \lambda. \end{cases}$$

In other words, while the queue always increases at rate λ when $\mathcal{D} \leq P\lambda/(\mu - \lambda)$, the increasing phase of $Q(t)$ is in this case composed of different parts where the queue increases at rate λ or at rate $2\lambda - \mu$. More precisely, every time the queue remains empty for a period of time T_E , it implies that, after a delay \mathcal{D} , the queue will increase at rate $2\lambda - \mu$ during a period of length T_E . The typical queue dynamics arising when $\mathcal{D} > P\lambda/(\mu - \lambda)$ is explained in more details in Fig. 3.12. Even if the self-induced mechanism leading to the oscillatory behaviour of the queue content, described in Fig. 3.12, remains qualitatively the same for any value of the triplet of parameters $(P, \mathcal{D}, \rho = \frac{\lambda}{\mu})$, one might observe slightly different queue dynamics depending on the value of these parameters. More particularly, the shape of the dynamics depends on when the siphon purging mechanism is triggered during the increasing phase (*i.e.* when $Q(t) = P\mu$). Without being exhaustive, two other typical situations are exposed in Figs. 3.13 and 3.14. Note that in any situation, it remains possible, for large P , to characterize analytically the amplitude and period of the emerging quasi-deterministic oscillations of the queue content.

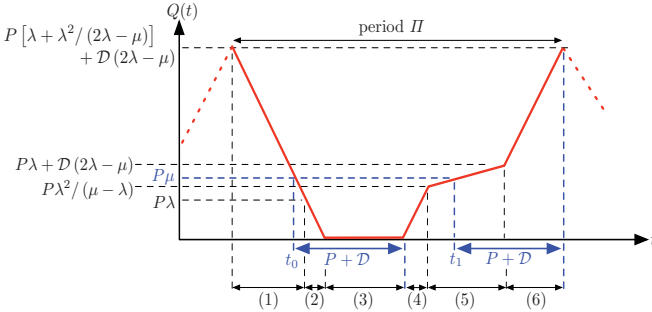


Fig. 3.13. Oscillatory queue length dynamics when the agents' feedback is delayed by $\mathcal{D} > P\lambda/(\mu - \lambda)$. The illustrated situation arises when $\frac{\mu(P+\mathcal{D})}{P+2\mathcal{D}} \leq \lambda < \lambda^*$. Note that $2\lambda - \mu > 0$ in the present situation. In this case, the amplitude and period of the queue length oscillations are respectively given by $\Delta = P \left(\lambda + \frac{\lambda^2}{2\lambda - \mu} \right) + \mathcal{D}(2\lambda - \mu)$ and $\Pi = P \left(\frac{\lambda\mu}{(2\lambda - \mu)(\mu - \lambda)} \right) + \mathcal{D} \frac{\mu}{\mu - \lambda}$.

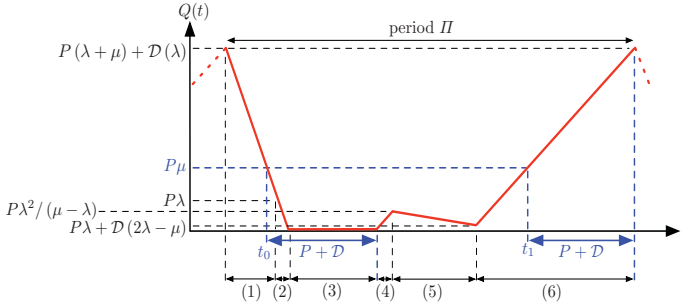


Fig. 3.14. Oscillatory queue length dynamics when the agents’ feedback is delayed by $\mathcal{D} > P\lambda/(\mu - \lambda)$. The illustrated situation arises when $\frac{\mu}{P+2\mathcal{D}} \leq \lambda < \frac{\mu(P+\mathcal{D})}{P+2\mathcal{D}}$. Note that, in the present situation, $2\lambda - \mu$ might be either positive or negative. In this case, the amplitude and period of the queue length oscillations are respectively given by $\Delta = P(\lambda + \mu) + \mathcal{D}\lambda$ and $\Pi = P(2 + \frac{\mu-\lambda}{\lambda}) + \mathcal{D}(3 + \frac{\mu-2\lambda}{\lambda})$.

3.8.2 Heterogeneous Agents

Until that point, we have considered homogeneous agents, sharing a common patience parameter P . As exposed previously in this chapter, such a stylized model is very useful in the sense that it allows for an in depth understanding of the emerging phenomena and it moreover enables fully analytical considerations. In this section, we relax this assumption and study how the flow dynamics is affected in the more realistic case where this patience parameter is individualized. More precisely, we consider again the feedback queueing system introduced in Section 3.2 but now, each heterogeneous agent $\zeta \in \mathbb{N}$ possesses an individual patience parameter $P_\zeta > 0$, which is an i.d.d. continuous random variable taking values in $[P_{\min}; P_{\max}]$ with cumulative distribution function $F(x) = \mathbb{P}(P_\zeta < x)$, mean $\bar{P} = \mathbb{E}(P_\zeta)$ and variance σ_P^2 .

As we will see, the main difference to the former case (see Section 3.4) is that now the amplitude Δ of the queue content oscillations depends not only on μ but also on λ . This complicates somehow the dynamics and makes analytical quantitative results more tedious to derive, but it does not qualitatively affect the essence of the emerging phenomenon (*i.e.* cyclo-stationary behaviour of the queue content).

It can be observed by simulations (see Fig. 3.15) that the more the customers’ patience parameter fluctuates, the smaller are the oscillations (*i.e.* Δ is monotonically decreasing with σ_P^2 from a mean maximum value $\Delta_{\max} = \bar{P}\mu$ to a mean minimum value Δ_{\min}).

We first derive very basic upper and lower bounds that confine the queue

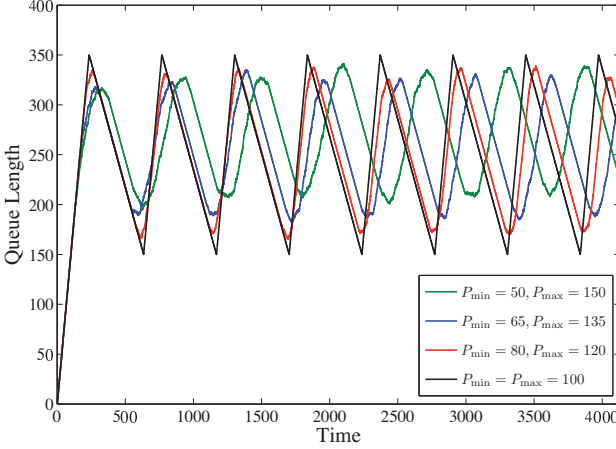


Fig. 3.15. Queue length dynamics for heterogeneous agents with P_ζ uniformly distributed in $[P_{\min}; P_{\max}]$, when inter-arrival times are uniformly distributed in $[\frac{1}{3}, 1]$, service times are uniformly distributed in $[0.1, 0.9]$ ($\rho = \frac{2}{3}$).

content oscillations obtained for heterogeneous patience parameter. Trivially, the oscillations are bounded from above by $P_{\max}(\lambda + \mu)$, which corresponds to the peak of the oscillatory behaviour when all the agents share a common patience parameter $P_\zeta = P_{\max}, \forall \zeta$ (see Fig. 3.16). Likewise, the oscillations are bounded from below by $P_{\min}\lambda$ (corresponding to the situation where $P_\zeta = P_{\min}, \forall \zeta$). These trivial bounds allow us to conclude that, when the agents' patience parameter is randomized, the amplitude Δ of the queue content oscillations is bounded such that:

$$\Delta < \lambda(P_{\max} - P_{\min}) + \mu P_{\max}. \quad (3.12)$$

As we will now see, these trivial bounds can be further noticeably improved. In this respect, we focus on the queue content purging mechanism, which is, as shown in the explanatory graph given in Fig. 3.17, more complex in the present situation than in the case of homogeneous agents. Let t_0 be the time at which the queue content reaches level $Q(t_0) = P_{\min}\mu$, during the increasing phase. A customer ζ_0 entering the queue just after t_0 will possibly leave the system after service, as it will suffer, with high probability, a sojourn time $W > P_{\min}$. More precisely, ζ_0 will not take the feedback loop provided its individual patience parameter P_{ζ_0} is equal to P_{\min} , which happens with probability $F(P_{\min})$. Consequently, from time $t_0 + P_{\min}$, a part of the customers will be unsatisfied with service. We define

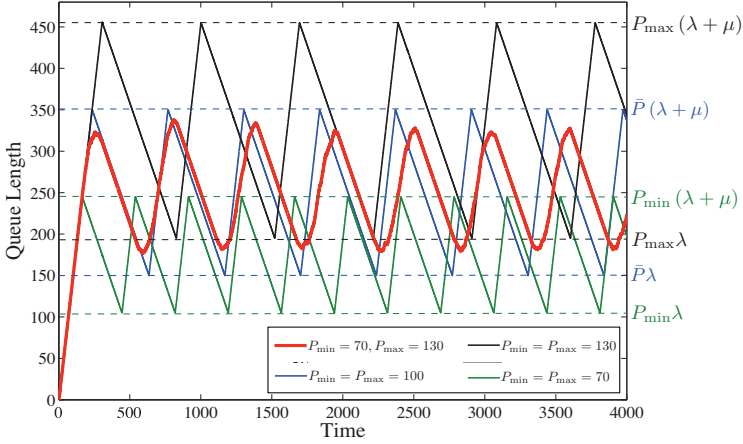


Fig. 3.16. Queue length dynamics for heterogeneous agents with P_C uniformly distributed in $[P_{\min}; P_{\max}]$, when inter-arrival times are uniformly distributed in $[\frac{1}{3}, 1]$, service times are uniformly distributed in $[0.1, 0.9]$ ($\rho = \frac{2}{3}$).

$$\begin{aligned} \mu_{\text{eff}}\left(t + \frac{Q(t)}{\mu}\right) &= \mu \mathbb{P}\left(P_C \leq \frac{Q(t)}{\mu}\right) \\ &= \mu F\left(\frac{Q(t)}{\mu}\right) \end{aligned} \quad (3.13)$$

as the effective rate at which the agents leave the system after service (remember that there is a delay between the time t a customer enters the queue and the time $t + \frac{Q(t)}{\mu}$ it leaves service). Note that $\mu_{\text{eff}}(t) \in [0, \mu]$ is equal to 0 for $t < t_0 + P_{\min}$ ($\frac{Q(t)}{\mu} < P_{\min} \Rightarrow F\left(\frac{Q(t)}{\mu}\right) = 0, \forall t < t_0$). Furthermore, $\mu_{\text{eff}}(t)$ monotonically increases from 0 to μ during the time interval $[t_0 + P_{\min}, t_2 + P_{\max}]$, where t_2 is the time at which the queue reaches level $Q(t_2) = P_{\max}\mu$. Indeed, the larger is the queue content $Q(t)$, the less likely are the customers to take the feedback loop. At time $t_2 + P_{\max}$, any customer leaving service will have suffered a sojourn time $W > P_{\max}$ and will hence quit the system (*i.e.* $\mu_{\text{eff}}(t) = \mu, t \geq t_2 + P_{\max}$).

When, in the case of homogeneous agents, $Q(t)$ was increasing at a constant rate λ and instantaneously switched to decrease at rate $\lambda - \mu$, the transition between the increasing and decreasing phases is here smoother. Indeed, the queue length evolves at rate $\lambda - \mu_{\text{eff}}(t)$, which progressively transits from λ to $\lambda - \mu$ during the time interval $[t_0 + P_{\min}, t_2 + P_{\max}]$.

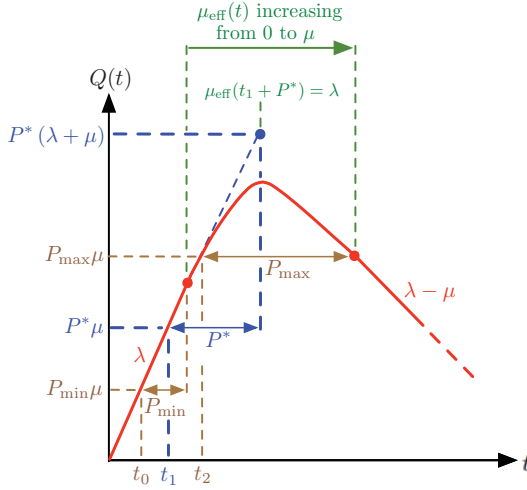


Fig. 3.17. Explanatory graph of the queue length dynamics observed for heterogeneous agents with individual patience parameter randomly distributed in $[P_{\min}; P_{\max}]$.

Let $t_1 \in [t_0, t_2]$ be the time such that

$$\lambda - \mu_{\text{eff}}(t_1 + P^*) = 0, \quad (3.14)$$

with

$$P^* = \frac{Q(t_1)}{\mu}. \quad (3.15)$$

At time t_1 , the queue length is large enough such that the proportion of agents joining the queue and that will leave the system after service (*i.e.* those with $P_\zeta < P^*$) balances the incoming rate λ of agents. Consequently, the switch between the increasing and decreasing phases takes place at time $t_1 + P^*$. Using Eq. (3.13), Eq. (3.14) becomes:

$$F(P^*) = \frac{\lambda}{\mu} = \rho$$

and hence P^* is given by

$$P^* = F^{-1}(\rho), \quad (3.16)$$

with the following definition:

$$F^{-1}(y) = \inf_{x \in \mathbb{R}} \{F(x) \geq y\}.$$

Note that $P^* \in [P_{\min}; P_{\max}]$. In the time interval $[t_1, t_1 + P^*]$, as some customers might already leave system after service, $Q(t)$ increases at rate

$$\lambda - \mu_{\text{eff}}(t) \leq \lambda. \quad (3.17)$$

Consequently, the peak of the oscillation (*i.e.* the switch between the increasing and decreasing phases) satisfies

$$\begin{aligned} Q(t_1 + P^*) &= Q(t_1) + \int_{t_1}^{t_1 + P^*} (\lambda - \mu_{\text{eff}}(t)) dt \\ &\leq P^* \mu + P^* \lambda \\ &= P^* (\lambda + \mu), \end{aligned} \quad (3.18)$$

where we have used Eqs. (3.15) and (3.17) to get the second inequality.

Likewise, now considering the second part of the oscillations, $Q(t)$ also switches smoothly from the decreasing to the increasing phases and the growing rate progressively goes from $\lambda - \mu$ (all customers leaving system after service) to λ (all customers taking the feedback loop), as more and more agents are satisfied with service with decreasing queue length. Let t_3 be the time such that

$$\lambda - \mu_{\text{eff}}(t_3 + P^*) = 0,$$

during the decreasing phase. Repeating the same lines as the above reasoning, it follows that the queue length is bounded from below such that:

$$Q(t_3 + P^*) \geq P^* \lambda. \quad (3.19)$$

From Eqs. (3.18) and (3.19), it thus results that the amplitude Δ of the queue content oscillations satisfies the following condition:

$$\Delta \leq P^* \mu. \quad (3.20)$$

Note that this bound not only depends on μ but has also, through the definition of P^* given by Eq. (3.16) (remember that $\rho = \frac{\lambda}{\mu}$), an implicit dependency on the external incoming traffic λ . Indeed, as stipulated at the beginning of this chapter, the shape of the emerging self-organized oscillations depends on both λ and μ in the case of heterogeneous agents and hence the refined bound given by Eq. (3.20) directly considers this specific feature. While the trivial, easily derived, bounds (for the minimum, maximum and amplitude of the oscillations) provided previously in this chapter and summarized in Eq. (3.12) were solely the result of a worst-case analysis, the more refined bounds given by Eqs. (3.18), (3.19) and (3.20) take into account the whole probability distribution that governs the value of the agents' patience parameter (and consequently not only the average but the higher moments of the distribution are considered). The improvement in terms of accuracy resulting from the implementation of these more refined bounds with respect to the trivial ones is illustrated in Fig. 3.18. To close, note that the performance of the refined bounds also depends on the particular shape of the distribution

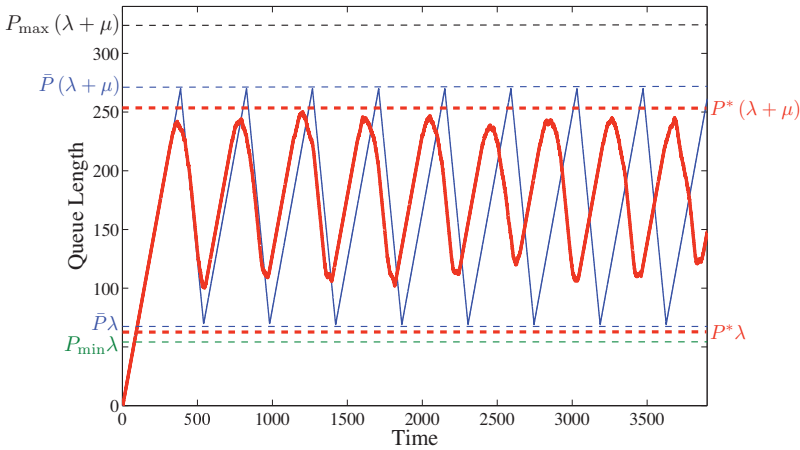


Fig. 3.18. Queue length dynamics for heterogeneous agents with P_c uniformly distributed in $[80; 120]$, when inter-arrival times are uniformly distributed in $[0.5, 2.4]$, service times are uniformly distributed in $[0.1, 0.9]$ ($\rho = 0.35$).

of the patience parameter (skewness and kurtosis). Restricting to symmetric distributions, it appears that the upper bound $P^*(\lambda + \mu)$ is more accurate than the lower bound $P^*\lambda$ when $\rho < \frac{1}{2}$ (see Fig. 3.18) and furthermore the performance of the upper bound increases monotonically as ρ tends towards 0. When $\rho > \frac{1}{2}$, the situation is reversed: the lower bound is the better and its accuracy is enhanced as ρ goes to 1.

3.9 Implementation Within Logistics Networks - The Smart Parts Paradigm

Mainly inspired by human behaviour, the particular routing mechanisms proposed in this chapter, based on historical data, might possibly also be encountered in logistics systems. Indeed, new forms of communication and information technologies enable the development of autonomous cooperating processes within logistics networks. More concretely, the ever growing availability of new technologies such as Radio-Frequency-Identification-Devices (RFID) (*i.e.* “smart tags”) allows for the wide implementation of local “intelligence” to physical items circulating in various logistics systems like distribution networks, production lines or supply chains. These autonomous decision-making units are generally known as *smart parts*, [127]. The main characteristic of these smart parts is their capability to control themselves by taking non-human decisions and by individually solving problems such as planning and

production choices. Depending on the application, smart parts could be raw materials, components or finished products, but they could also be the transit or transportation equipments themselves, such as pallets, conveyors or trucks, [127]. Following this paradigm, the type of dynamics presented in this chapter might hence naturally emerge in so-called smart parts logistics networks. An illustration of how the type of emerging collective dynamics studied in this chapter can actually be used to implement fully decentralized management in multi-server production lines with smart parts, and how it can accordingly achieve optimal dynamics, is given in Chapter 9.

3.10 A Solvable Occurrence of a Complex Adaptive System

Coming from biology, the concept of Complex Adaptive Systems (CAS) originates in the field of complexity science and is often referred to living systems. It has been formalized, according to John H. Holland [58, 59, 95], with the help of the following definition:

A Complex Adaptive System (CAS) is a dynamic network made of many agents (which may represent cells, species, individuals, households, firms, nations) evolving in parallel. Each agent is constantly acting and reacting, autonomously, to what the other agents are doing. Consequently, the agents find themselves in an environment produced by the successive interactions made with the other agents in the system. The control of a CAS tends to be highly dispersed and decentralized, without any singular entity deliberately managing and controlling it. If there is to emerge any coherent behaviour in the system, it has to arise from competition and cooperation among the agents themselves. The overall behaviour of the system is the result of a huge number of decisions made at every moment by many individual agents.

In other words, the highly interdisciplinary study on CAS focuses its modelling activities on how microstate actions self-organize into emergent collective patterns. Provided all the above properties are satisfied within a logistics system, one speak then of a complex adaptive logistics system (CALs). As stipulated in Section 3.9, the smart parts concept opens wide the door for the existence of CALs in many different domains, including production, distribution and supply chains. Depending on the specific implementation, note that the feedback queueing system considered in this chapter might either be called a CAS in the case of humans asking for services or a CALs when smart parts are circulating within a logistics network.

Our aim here is to show that the class of models considered in this chapter might be identified to belong to the realm of CAS. More precisely, do the

presently considered service (respectively logistics) systems possess all the required properties to be formally classified as CAS (respectively CALS)? In this prospect, we will base our study on [127], a recent contribution where it is carefully testified that supply networks can be validly treated as CAS. More specifically, in [127], the authors define, primarily following [59, 70], a list of the essential properties that characterize CAS. In the following, we look over that list in our particular framework and we consequently show that our class of models fulfils all the generally accepted features to be categorized as CAS. Properties (1) to (4) characterize the elements evolving in the complex system while properties (5) to (7) qualify the system global behaviour.

- (1) *Autonomy.* In our context, each agent possesses the capability to take its own routing decisions within the network. Hence, we are in presence of agents that act autonomously, without any influence or help of a central controller. All routing decisions within the network are self-initiated and decentralized to the circulating items.
- (2) *Interaction.* The interactions between the different agents constituting our feedback queueing system are purely stigmergic. Indeed, there is no direct communication between the circulating items, but they interact via the state of the system. More particularly, when an agent ζ decides to take the feedback loop and hence to stay in the system, its action will directly influence the other agents' future behaviour. Indeed, these agents will suffer a larger waiting time due to the presence of the agent ζ in the system, which might accordingly influence their future routing decisions within the network. All agents share thus a stable and unique degree of interaction, while they automatically transmit information to the others about the system state when leaving or staying in the system. Note furthermore that the different agents do not share the same knowledge about the system state. In fact, each agent ζ only knows the congestion state of the queueing system at a specific time t_ζ (that corresponds to the time when ζ entered the system).
- (3) *Heterogeneous Agents.* In the stylized feedback queueing system considered in this chapter, although the same structural rule \mathcal{R} governs the agents' individual routing decision mechanisms, the agents' behaviour might differ with respect to their individual patience parameter P_ζ , as considered in Section 3.8.2. As we have noticed in Section 3.8.2, the essence of the emerging global phenomena (*i.e.* the periodic structure of the queue content evolution) is preserved despite the introduction of such heterogeneity in the agents' patience parameter. In that situation, we hence deal with heterogeneous agents, each of them acting differently (*i.e.* deciding to take the feedback loop or not) with report to the current, individually observed, state of the system. Indeed, when all the agents follow qualitatively the same individual goal (*i.e.* to stay loyal to the service provider when service is satisfactory), the achievement of this goal might quantitatively differ (with respect to P_ζ) from agent to agent. Note that, while a large degree

of heterogeneity (*i.e.* highly fluctuating individual patience parameter P_ζ) does not affect the essence of the emerging self-organized structure, it could however diminish the analytical tractability of our models.

- (4) *Ability to Learn.* The items circulating in our particular queueing system possess the capability to learn from the system state, which consists more precisely in assimilating the permanently updated result of the successive feedbacks (*i.e.* to leave or to stay in the system) left by the preceding agents. This very basic learning and adapting ability, which ultimately enables the implementation of the agents' autonomous routing mechanisms based on personal waiting memory, will be enhanced in the following chapters. Note that our agents fulfil the essence of what is any learning process; they make an experiment (*i.e.* they experiment the sojourn time in the system), they assimilate that measure and then they modify their future behaviour (*i.e.* they take their next autonomous routing decision accordingly).
- (5) *Self-Organization.* As we have seen in Section 3.4, the collective structure arising from the agents' numerous individual actions takes here the form of a periodic purging of the queueing system (see Fig. 3.5). From a managerial point of view, the coevolution of the autonomous agents within the network leads to the self-organized stabilization of the queueing system and to the maximization of the resource utilization (*i.e.* the server busy period is maximized). The emergence of this global pattern is entirely due to the autonomous circulating agents themselves, since no central entity exerts any authority on their behaviour. Note that the system is adaptive as it reacts to an excessive queue length by triggering a purging of the queue.
- (6) *Melting Zone.* Introduced by Stuart A. Kaufmann [70], the concept of melting zone refers, [127], to a region between the *edge of order* and the *edge of chaos*. These edges define a *region of emergent complexity*, [91, 92], where the door is open for self-organized structures and emergent system global behaviour. In our model, this phase transition results from the value of the agents' patience parameter P (or the average value when P_ζ is individualized). For sufficiently large values of this parameter (see Fig. 3.5), our queueing system evolves and stays within the melting zone and hence maintains itself in state of *self-organized criticality*, [11]. The capability of our system to stay, for sufficiently large P , in the melting zone is essential for its survival, as it ensures intrinsically a bounded queue length and hence the stability of the queueing process. As it is stipulated in [127], while it has been possible to identify the melting zone in various living systems, the existence of the melting zone is generally only assumed for self-organized logistics systems and there does not exist at this time any other examples of such a parallel in the logistics literature.
- (7) *Coevolution.* In our queueing system, the agents are competing for a limited resource (*i.e.* the service). Each agent leaves a feedback on its personal service satisfaction when it decides to leave or to stay in the system, thus building up a stigmergic interaction mechanism between the agents. In

other words, the agents sequentially respond to the others' actions, but they also successively influence the environment by leaving a feedback on their individual experience with service. According to [30], coevolutionary processes within logistics networks result from the joint presence of non-linearity and path dependence. In our model, the nonlinearity takes the form of the feedback loop topology and the path dependence is inherent to the non-Markovian routing decision rule \mathcal{R} adopted by the agents. Note that in the situation where the environment would be affected by external changes (for example in the case where the incoming or service rate would be modified), the system dynamics would directly adapt to the new configuration (*i.e.* the period and amplitude of the queue content stable oscillations would be modified accordingly).

According to the above list of required properties, we can conclude that the bridge between the feedback queueing system considered in this chapter and CAS is manifest. Consequently, we are hence in presence of a specific instance of a CAS (or CALS when smart parts are implemented within logistics systems) in the context of QNs, which reveals furthermore itself to possess the striking feature to remain analytically tractable.

3.11 Concluding Remarks

We have seen in this chapter that the presence of autonomous decision making agents in queueing systems fundamentally modifies the intrinsic flow dynamics and hence the resulting queueing processes. More particularly, a management fully decentralized to the circulating items can, among others, leads to oscillatory behaviour of the queue content. Depending on the application, this cyclo-stationary pattern might obviously have either positive or negative implications. As the stylized models presented in this chapter allow to understand in detail the origins of the emerging phenomena, they could potentially help, in practical situations, managers who want either to benefit from these phenomena or on the contrary who try to reduce their emergence. As an illustration, think of a manager aiming to reduce the creation of temporal oscillations in a service network. Following numerous studies, the environment manifestly influences the waiting perception and hence to provide for example comfortable seats, free drinks or discounts will reduce the waiting time perceived by the customers. Based on this assertion, a strategy according to which discounts would be periodically given when the waiting time has been particularly large would help to reduce the amplitude of the queue content oscillations. There also exist situations where one could benefit from the emerging periodic structure, an example in the domain of production lines is provided in Chapter 9.

3.12 Contributions of Chapter 3

- We present a new type of dynamics arising in queueing networks when one of the most important hypothesis of classical queueing theory, the Markovian character, is broken. We construct a stylized model for recurrent services where the circulating items are autonomous agents able to decide their routing based on their past history within the network. While its non-Markovian aspect precludes the existence of a stationary state *stricto sensu*, our model reveals itself to possess a robust self-organized structure which takes here the form of a cyclo-stationary state. The emerging dynamics leads to oscillatory behaviour of the queue length, which ensures bounded queue length and maximum resource utilization. For regimes where the law of large number holds, we discuss analytically the emerging periodic structure of the queue content.
- We emphasize that even if the type of dynamics exposed in this chapter is mainly inspired by human behaviour, it might also be encountered in logistics systems. Indeed, the increasing availability of RFID tags allows for the implementation of local “intelligence” on the items circulating in any logistics network. Often referred as the “smart parts” paradigm, such kind of decentralized management opens great potentialities in the domains of production, distribution and supply chains.
- The accumulate-and-fire dynamics that emerges in our particular feedback queueing model is the classical signature of relaxation oscillators and hence we point out that we are actually in presence of a multi-agent limit-cycle. Furthermore, we show that the models considered throughout this chapter can be classified as complex adaptive systems, counting thus among the very rare solvable occurrences of such systems existing in the available literature.

Extended Model for Recurrent Services - Introducing Weariness Aspects

Summary. *For recurrent service providers (fast-food, entertainment, medical care, etc.), retaining loyal customers is obviously a key issue. The customers' loyalty essentially depends on their service satisfaction which is defined via an ad-hoc utility function. As seen previously in this work, the utility function strongly depends, among several other criteria, on the past perceived waiting time. Moreover, it is likely that the patience that customers consent to allow in waiting often decreases as a function of the successive uses of the service (i.e. lassitude). In this chapter, we introduce this weariness aspect into the feedback queueing system considered in Chapter 3. Consequently, we propose an idealized queueing model with feedback loop in which the customers' loyalty is determined solely by the individual experience gained during the successive visits to a service (i.e. the waiting time and the number of services yet received). Again, for regimes where the law of large numbers holds, a deterministic approach enables to analytically discuss the resulting multi-agent dynamics governing the customers' flows. One is able, in particular, to fully calculate the characteristics of the emerging complex patterns (i.e. structured temporal oscillations of the queue content) which are observed to be strongly structurally stable.*

4.1 Introduction

In the preceding chapter, only the last recorded waiting time was used by the agents in their routing decision criterion. Relaxing this assumption by, for instance, allowing the customers to use the information collected during their successive visits to a server is also relevant in numbers of real-life applications. Thinking to traffic issues for example, a driver often considers not only his/her last trip but a more elaborate (and often irrational) function of his/her successive trips to decide his/her upcoming daily route (see Section 10.1 for an illustration). A synthetic and quantitative formulation of such a heuristic observation requires a set of mathematical tools available from a "fusion" between game and queueing theoretical approaches, a general point of view adopted by R. Hassin and M. Haviv in [55]. One of the fundamental lessons taught by game theory is that the distinction between games played

only once with those played repeatedly is mandatory, as it leads to drastically different optimal strategies, [9]. Similarly, for competition between queueing nodes, the fact that customers pay a single or several successive visits to the servers (and the fact that history matters in the competition process) does strongly influence the resulting traffic flows. In [55], customers always pay a single visit to the servers and the resulting stationary equilibria (*i.e.* the Nash equilibria) are thoroughly discussed. In this chapter, we will focus on the new dynamical features emerging from repeated visits to a server, when these successive visits influence the customers' routing decisions.

As an illustration, let us consider a theme park in which, among other attractions, a roller-coaster is offered. This roller-coaster entertains people at a limited flow rate which, due to the high demand, is responsible for the formation of a queueing process. We assume here that the park entrance fare offers to visitors the possibility to attend any attraction repeatedly and without limitation. Due to the exciting sensations generated by the roller-coaster, customers agree to line-up and are fairly patient when attending the coaster for the first time. Repeated runs however does weary their patience. The trade-off between the excitement delivered by a roller-coaster trip and the waiting burden incurred before boarding, can be quantified by a (usually individual) utility function. When the utility is negative, the customers are deterred and leave the roller-coaster for another spot. This roller-coaster example belongs to the highly profitable leisure and hospitality (L&H) sector, which includes the entertainment and recreation, the tourism and accommodation as well as the food services. Far from being exceptional, note that the previous customers' behaviour is in fact quite common in the L&H sector. Already exposed in the introductory example of Section 1.1, ski traffic management offers another world-wide illustration. Generically, people will change slope either when they have suffered a large waiting-time at the ski lift and/or when they are bored of having done the same ski run several times.

Despite the apparent simplicity of the above roller-coaster illustration, the resulting queueing dynamics is highly complex. It depends indeed simultaneously on a routing feedback loop (*i.e.* an intrinsic nonlinearity) due to the customers lining for a new trip and on an enhanced (with regard to Chapter 3) history-based (HB) routing decision mechanism. The decision to come back (*i.e.* to remain loyal) or to quit would now be taken according to a patience threshold, itself depending on the number of previous runs already achieved. At first sight, little hope is left regarding the possibility to characterize analytically the traffic flows resulting from this complex dynamics. Keeping the central features, namely (*i*) the nonlinearity and (*ii*) the enhanced HB decision-making policy, we are nevertheless able, for a somehow simplified class of models, to describe analytically the resulting dynamics. Basically, two simplifications of the original situation exposed above are introduced. On one hand, we restrict to agents having a common utility function. On the

other hand, we separate the role played by the waiting time and the number of repeated visits to the server by introducing two distinct utility thresholds. When exceeded, these thresholds trigger the loss of the agents' loyalty. This simplified multi-agent dynamics generates the emergence of generically stable time-dependent periodic queue contents. The self-organized oscillatory behaviour of the queue length observed in the present case reveals itself to be more complex in comparison to the one studied in Chapter 3 and possesses inherent features that are specific to the weariness aspect introduced here.

The chapter is organized as follows. The single-stage feedback queueing model that will be considered throughout this chapter is described in Section 4.2. More particularly, we introduce the new HB routing rule, enhanced in complexity in comparison to the one considered in Chapter 3, that incorporates now lassitude aspects. Section 4.3 is devoted to experimental results. Illustrations of the self-organized collective patterns (characterized by self-sustained oscillations) that emerge when the new routing rule is implemented are provided. In Section 4.4, following a deterministic approach, we describe analytically this oscillatory behaviour. We explain in particular the apparition of an extra peak during the increasing phase of the oscillations. The chapter ends with Section 4.5, in which concluding remarks and further perspectives are given.

4.2 Model - Considering Weariness

In this chapter, we start with the same framework as in Chapter 3 and we model again the customers' behaviour faced with a recurrent service with a single queueing network (QN) composed, as illustrated in Fig. 3.1, of a single server and a feedback queue. An incoming flow of customers, described by a renewal process with mean inter-arrival time $\frac{1}{\lambda}$ and probability distribution $A(x)$ with density $dA(x)$, is served by a processing unit which service times are i.i.d. random variables with mean $\frac{1}{\mu}$, probability distribution $B(x)$ and density $dB(x)$. Accordingly, the incoming and service rates of the renewal processes are respectively denoted by the parameters λ and μ . The distributions $A(x)$ and $B(x)$ are assumed to have finite moments. We suppose again that the traffic intensity $\rho = \frac{\lambda}{\mu} < 1 \Leftrightarrow \lambda < \mu$, which ensures the stability of the queueing system when there is no feedback loop. We assume finally that the waiting room capacity is unlimited and that the service discipline is first-in-first-out (FIFO). After being served at the decision node n , each customer has to choose among two possibilities, namely:

- (i) either to quit the system definitively, or
- (ii) to follow the feedback loop and line up again to be served once more.

The present simplified model of recurrent services assumes that the customers' loyalty is based on their individual experience with the service provider. More

particularly, the decision of a customer, at n , either to come back for another service or to leave the system depends here on

(1) the last sojourn time W it has spent in the system in order to be served and on

(2) the number of services N_{it} it has already received (we say that the customer is at its i^{th} iteration).

We suppose furthermore that the influence played by these two measures is uncoupled. The customers, who share a common utility function (*i.e.* we consider homogeneous customers), consider two separate thresholds. They possess first, as in Chapter 3, a common patience parameter P to which they will compare their last experimented sojourn time W . Secondly, they will check that they have received less services than a common weariness parameter N_{\max} . It leads to the introduction of two independent rules \mathcal{R}_1 and \mathcal{R}_2 , that the customers will apply when they decide their routing at decision node n . The first one is controlled by the sojourn time and is given by:

$$\mathcal{R}_1 = \begin{cases} \text{follow alternative (i)} & \text{if } W > P, \\ \text{follow alternative (ii)} & \text{if } W \leq P. \end{cases}$$

The second rule is driven by the number of already received services, it is defined as:

$$\mathcal{R}_2 = \begin{cases} \text{follow alternative (i)} & \text{if } N_{it} \geq N_{\max}, \\ \text{follow alternative (ii)} & \text{if } N_{it} < N_{\max}. \end{cases}$$

Combining these two independent rules, the customers will hence choose their routing at n following:

$$\begin{aligned} \mathcal{R}_{\text{tot}} &= \mathcal{R}_1 \cap \mathcal{R}_2 \\ &= \begin{cases} \text{follow alternative (ii)} & \text{if } W \leq P \text{ and} \\ & N_{it} < N_{\max}, \\ \text{follow alternative (i)} & \text{otherwise.} \end{cases} \end{aligned} \quad (4.1)$$

As already stipulated in Chapter 3, we speak of loyal customers when alternative (i) is chosen, as they are pleased with the server and then return to it for another service. The rule \mathcal{R}_{tot} states that, providing its sojourn time remains below its patience parameter and providing it has already received less than a limiting number of services, a customer comes back for another service.

Remember that the dynamics involving \mathcal{R}_1 alone is discussed in Chapter 3. In this case and when P is large enough, quasi-deterministic cyclo-stationary regimes emerge, *i.e.* stable temporal oscillations of the queue level $Q(t)$ are observed (see Fig. 3.5) and this independently of the detailed nature of the probability laws $A(x)$ and $B(x)$.

Note that our simplified model assumes that P_i , the customers' common patience parameter when they receive their i^{th} service, has the following form:

$$P_i = \begin{cases} P & \text{if } i \leq N_{\max}, \\ 0 & \text{if } i > N_{\max}. \end{cases}$$

The next step to get closer to real-life customers' behaviour would be to consider more general forms for P_i . Typically, this patience parameter could be monotonically decreasing with i , denoting that the customers' loyalty suffers from a progressive weariness over time. Likewise, N_{\max} is here common to all customers. A natural generalization would be to consider that each customer possesses its own behaviour when faced with lassitude.

4.3 Experimental Observations

Figs. 4.1 and 4.2 show the typical dynamics of the queue length when customers follow the HB rule \mathcal{R}_{tot} given by Eq. (4.1) to choose their routing at decision node n .

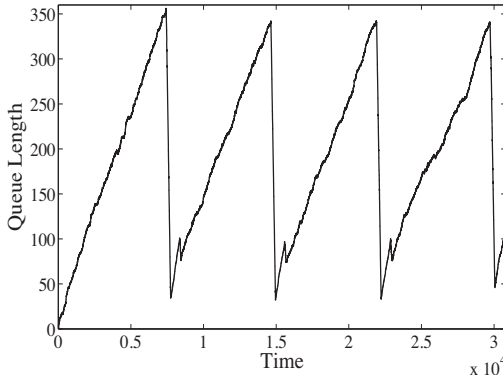


Fig. 4.1. Queue length dynamics when inter-arrival times are uniformly distributed in $[1, 17]$ with mean $\frac{1}{\lambda} = 9$ and coefficient of variation $CV = 0.51$, service times are uniformly distributed in $[0.1, 1.7]$ with mean $\frac{1}{\mu} = 0.9$ and $CV = 0.51$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 12$.

Independently of the inter-arrival and service times distributions and when P is large enough (see [43] for a more detailed discussion on the parameter P), we also observe here the emergence of quasi-deterministic cyclo-stationary regimes, *i.e.* stable temporal oscillations of the queue content. Indeed, in spite

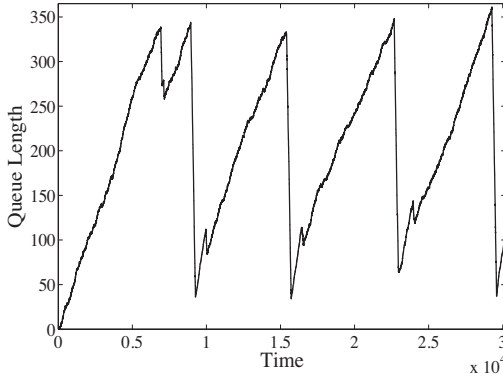


Fig. 4.2. Queue length dynamics when inter-arrival times are Erlang(3) with mean $\frac{1}{\lambda} = 9$ and CV = 0.57, service times are Erlang(3) with mean $\frac{1}{\mu} = 0.9$ and CV = 0.57 ($\rho = 0.1$), $P = 300$ and $N_{\max} = 12$.

of the presence of strong fluctuations, this robust and quasi-deterministic behaviour is a direct consequence of the law of large numbers (LLN), whose smoothing effect is manifestly observable in Figs. 4.1 and 4.2. According to this, for sufficiently large P , the dynamics can be approximately discussed via a deterministic approach. The dynamics resulting from deterministic inter-arrival and service times is illustrated in Fig. 4.3. It is remarkable that, compared with the dynamics exposed in Chapter 3 (involving rule \mathcal{R}_1 alone), the oscillations exhibit here an extra peak during their increasing phase. This peak is entirely due to rule \mathcal{R}_2 .

A restricted range of the values of the control parameters (*i.e.* λ , μ , P and N_{\max}) produces the fully complex dynamics visible in Figs. 4.1, 4.2 and 4.3. Indeed, when N_{\max} is large, the customers remain in the system for a long time before getting wearied. Hence, the queue length increases (new customer arrivals) and eventually reaches a level implying sojourn times larger than P . All the customers hence leave the system due to routing rule \mathcal{R}_1 (*i.e.* waiting-time) and none following \mathcal{R}_2 (*i.e.* maximum number of received services). The resulting dynamics is illustrated in Fig. 4.4 and 4.5 (random and deterministic cases). As we will see later, the regime where no extra peak appears emerges when:

$$\rho(1 + \rho)^{N_{\max} - 1} \geq 1. \quad (4.2)$$

Whenever the condition given by Eq. (4.2) holds, the model is similar to the one where only rule \mathcal{R}_1 is implemented (see Chapter 3).

We have focused here on homogeneous agents behaviour. However, note that

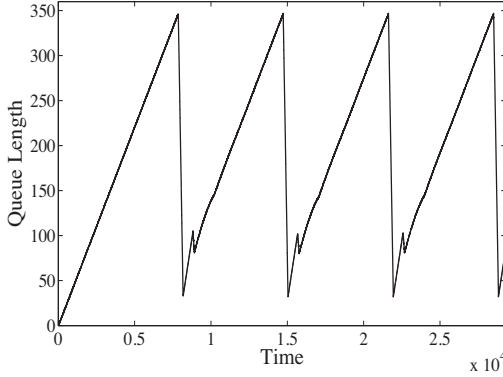


Fig. 4.3. Queue length dynamics for deterministic inter-arrival times $\frac{1}{\lambda} = 9$, deterministic service times $\frac{1}{\mu} = 0.9$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 12$.

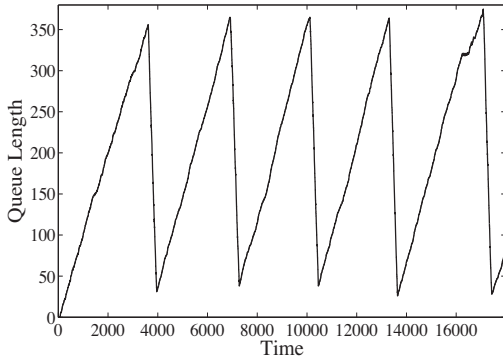


Fig. 4.4. Queue length dynamics when inter-arrival times are Erlang(3) with $\frac{1}{\lambda} = 9$, service times are Erlang(3) with $\frac{1}{\mu} = 0.9$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 25$.

the emergence of a macroscopic stable collective structure would also arise for agents with individualized patience parameter (see Section 3.8.2).

4.4 Analytical Discussion

As illustrated in Section 4.3, experimental results show that the dynamics, insofar as P is large enough, exhibit stable temporal oscillations of the queue content even in presence of strong fluctuations. Hence from now on, we focus

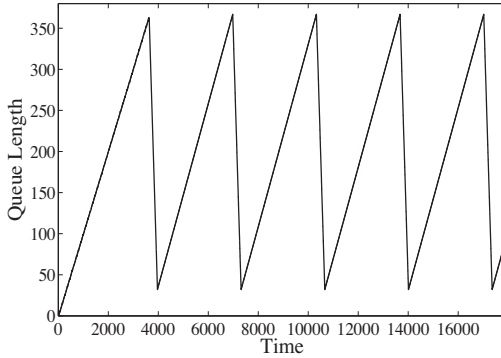


Fig. 4.5. Queue length dynamics for deterministic inter-arrival times $\frac{1}{\lambda} = 9$, deterministic service times $\frac{1}{\mu} = 0.9$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 25$.

on deterministic dynamics. We decompose the oscillatory dynamics in five distinct phases. Fig. 4.6 gives a sketch of these five phases on a generic oscillation of the queue length.

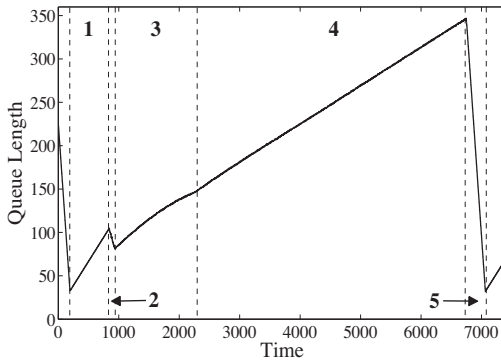


Fig. 4.6. A generic oscillation of the queue content dynamics for deterministic inter-arrival times $\frac{1}{\lambda} = 9$, deterministic service times $\frac{1}{\mu} = 0.9$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 12$.

- **First Phase: Pure Feeding**

As we will see in the fifth phase below, the queue is initially populated with an offset of $P\lambda$ fresh customers, who haven't received any service yet. During the first phase, which begins without loss of generality at time $t = 0$, the queue length remains small enough so that the customers wait less than their patience parameter P (and hence rule \mathcal{R}_1 is satisfied). Furthermore, since there is initially only fresh customers in the queue, no customer leaves the system due to the maximum number N_{\max} of iterations. The queue length hence increases at rate λ (arrival rate of new customers in the system). This first phase ends at time T_1 when the original $P\lambda$ customers have received N_{\max} services and thus will start to leave the system. To compute T_1 , we focus on $Q(k)$, the number of customers in the queue when all the $P\lambda$ original customers have completed k iterations, $k = 0, \dots, N_{\max}$. We first have that:

$$Q(1) = P\lambda \frac{1}{\mu} \lambda + P\lambda = P\lambda(1 + \rho).$$

Indeed, initially, the necessary time to serve the $P\lambda$ original customers is equal to $T(0) = P\lambda \frac{1}{\mu}$. During that time, $T(0)\lambda$ fresh customers join the queue, which is still populated by the $P\lambda$ original customers, who have now received one service. Following this iterative reasoning, we find that:

$$Q(k) = P\lambda(1 + \rho)^k, \quad k = 0, \dots, N_{\max}.$$

Accordingly, we find that:

$$T_1 = \frac{1}{\mu} \sum_{k=0}^{N_{\max}-1} Q(k) = P \left[(1 + \rho)^{N_{\max}} - 1 \right].$$

- **Second Phase: Offset Purging**

During the second phase, beginning at T_1 , a proportion of the customers receiving service leaves the system because they have been provided the maximum number of services (*i.e.* wearied customers). At the end of the first phase, $Q(N_{\max})$ customers populate the queue. Among them, $P\lambda$ have already done N_{\max} iterations (the original offset of customers). Accordingly, they will leave the system after the next service. We define thus the effective offset purging rate

$$\mu_{\text{eff}}(N_{\max}) = \mu \frac{P\lambda}{Q(N_{\max})} = \frac{\mu}{(1 + \rho)^{N_{\max}}}$$

as the rate at which the original customers leave the system due to the maximum number of iterations. In the second phase, the slope of the queue content is $s = \lambda - \mu_{\text{eff}}(N_{\max})$. Note that, in the case where $s \geq 0$, the extra

peak disappears. This phase ends when the $Q(N_{\max})$ initial customers have been served. The duration of the second phase is hence given by:

$$T_2 = \frac{Q(N_{\max})}{\mu}.$$

- **Third Phase: First Residual Offset Purging**

Again, in the third phase, wearied customers leave the system. In this phase, the presence of the first $P\lambda$ original customers still influences the dynamics. At the beginning of Phase 3, there are

$$Q(N_{\max} + 1) = \overbrace{Q(N_{\max}) \frac{1}{\mu} \lambda}^{\text{freshly arrived customers}} + \underbrace{Q(N_{\max}) - P\lambda}_{\text{leaving customers}} + \underbrace{P\lambda}_{\text{loyal customers}}$$

populating the queue. After the passage of these $Q(N_{\max} + 1)$ initial customers, there is

$$Q(N_{\max} + 2) = \overbrace{Q(N_{\max} + 1) \frac{1}{\mu} \lambda}^{\text{freshly arrived customers}} + \underbrace{Q(N_{\max} + 1) - P\lambda \frac{1}{\mu}}_{\text{leaving customers}} + \underbrace{P\lambda}_{\text{loyal customers}}$$

customers in the queue. Iteratively, we find that:

$$Q(N_{\max} + k) = P\lambda \left[(1 + \rho)^{N_{\max} + k} - (1 + \rho)^{k-2} (1 + k\rho) \right], \quad 1 \leq k \leq N_{\max} + 2,$$

where $Q(N_{\max} + k)$ is the queue content after all the $Q(N_{\max} + k - 1)$ customers previously in the queue have received service, $1 \leq k \leq N_{\max} + 2$. The resulting effective purging rate (here piecewise linear) reads as:

$$\mu_{\text{eff}}(N_{\max} + k) = \mu \frac{\rho(1 + \rho)}{(1 + \rho)^{N_{\max} + 2} - (1 + k\rho)}, \quad 1 \leq k \leq N_{\max} + 1,$$

with $\lambda - \mu_{\text{eff}}(N_{\max} + k)$ being the rate at which the queue raises between levels $Q(N_{\max} + k)$ and $Q(N_{\max} + k + 1)$. The length of each iteration, during which the queue content increases from $Q(N_{\max} + k)$ to $Q(N_{\max} + k + 1)$, is given by:

$$T(N_{\max} + k) = \frac{Q(N_{\max} + k)}{\mu}, \quad 1 \leq k \leq N_{\max} + 1.$$

Hence, the total duration of the third phase is equals to:

$$T_3 = \sum_{k=1}^{N_{\max}+1} T(N_{\max} + k)$$

$$= P \{ (1 + \rho)^{2N_{\max}+2} - [1 + \rho(N_{\max} + 2)] (1 + \rho)^{N_{\max}} \}.$$

• **Fourth Phase: Next Residual Offset Purging, Constant Growth**

The influence played by the $P\lambda$ original customers is further reduced and becomes, in a first order approximation, negligible. This leads to the apparition of a quasi-constant queue content increase rate $\lambda - \mu_{CG}$. Note that, if required, higher order approximations would be analytically conceivable. Accordingly, the rate of customers having reached the maximum number of iterations and thus leaving the system is quasi-constant over time during Phase 4. The purging rate $\mu_{\text{eff}}(2N_{\max} + 2)$ (the exact value of the rate at the beginning of this phase) yields a good approximation for μ_{CG} . It is given by:

$$\mu_{\text{eff}}(2N_{\max} + 2) = \frac{\mu\rho \left[(1 + \rho)^{N_{\max}+1} - 1 \right]}{(1 + \rho)^{2N_{\max}+2} - (1 + \rho)^{N_{\max}} [1 + (N_{\max} + 2)\rho]}.$$

Table 4.1 gives the accuracy of this approximation for several values of the control parameters.

External Parameters	$\mu_{\text{eff}}(2N_{\max} + 2)$	μ_{CG}	Error
$\frac{1}{\lambda} = 9, \frac{1}{\mu} = 0.9,$ $P = 300, N_{\max} = 12$	0.0632	0.0656	4%
$\frac{1}{\lambda} = 9, \frac{1}{\mu} = 0.9,$ $P = 300, N_{\max} = 17$	0.0312	0.0297	5%
$\frac{1}{\lambda} = 9, \frac{1}{\mu} = 0.8,$ $P = 300, N_{\max} = 12$	0.0788	0.0836	6%
$\frac{1}{\lambda} = 9, \frac{1}{\mu} = 1.1,$ $P = 300, N_{\max} = 12$	0.0419	0.0425	1.5%

Table 4.1. Comparison between the approximation and the experimental value of the quasi-constant purging rate of the fourth phase.

The queue length approximately raises at rate $\lambda - \mu_{\text{eff}}(N_{\max} + 2)$ until it reaches a P -dependent siphoning threshold. This phase ends when the queue length reaches the level (see Section 3.4):

$$Q(2N_{\max} + 3) = \overbrace{P\mu}^{\text{critical level}} + \underbrace{P(\lambda - \mu_{\text{eff}}(2N_{\max} + 2))}_{\text{delay mechanism}}.$$

Its duration is hence given by:

$$T_4 = \frac{Q(2N_{\max} + 3) - Q(2N_{\max} + 2)}{\lambda - \mu_{\text{eff}}(2N_{\max} + 2)}.$$

- **Fifth Phase: Siphon Mechanism**

At time

$$\tau = T_1 + T_2 + T_3 + T_4 - P(\lambda - \mu_{\text{eff}}(2N_{\max} + 2)),$$

the queue length is large enough ($= P\mu$) to get sojourn times W larger than P . As a consequence, at time $\tau + P$, a siphon purging with rate $\mu - \lambda$ is triggered. More precisely, when the queue content exceeds a critical level, it autonomously releases its emptying. This behaviour is fully analogous to the hydrodynamic self-siphoning device discussed in Section 3.4. When the siphon purging happens before there are wearied customers leaving the system, which happens whenever the condition given by Eq. (4.2) holds, regimes where only Phases 1 and 5 are visible emerge (see Figs. 4.4 and 4.5). Note that, in Phase 5, all customers leave the system due to rule \mathcal{R}_1 . At the end of the fifth phase, there remain $P\lambda$ customers in the queue (see Section 3.4). All these $P\lambda$ customers are new incomers, who have never been served yet. The duration of this last phase is given by:

$$\begin{aligned} T_5 &= \frac{Q(2N_{\max} + 3) - P\lambda}{\lambda - \mu} \\ &= \frac{P[\mu - \mu_{\text{eff}}(2N_{\max} + 2)]}{\lambda - \mu}. \end{aligned}$$

Summarizing the situation and grouping the above information, it is hence possible to compute the period

$$II = T_1 + T_2 + T_3 + T_4 + T_5$$

and the amplitude

$$\Delta = Q(2N_{\max} + 3) - P\lambda = P[\mu - \mu_{\text{eff}}(2N_{\max} + 2)]$$

of the stable temporal oscillations that are observed for the considered dynamics.

Table 4.2 gives a comparison, for the numerical values used in Fig. 4.6, between the analytical results given by the formulae derived in this section and simulation results.

4.5 Concluding Remarks

The present model is, by many aspects, oversimplified. In particular, to assume that all agents share a common patience threshold is obviously a pale

Analytical Results	Matching With Simulation
$T_1 = 641.46$	below 1% error
$Q(N_{\max}) = 104.61$	below 1% error
$\lambda - \mu_{\text{eff}}(N_{\max}) = -0.2429$	below 1% error
$T_2 = 94.15$	below 1% error
$Q(N_{\max} + 1) = 81.73$	below 1% error
$\lambda - \mu_{\text{eff}}(N_{\max} + 1) = 0.0658$	below 1% error
$\lambda - \mu_{\text{eff}}(2N_{\max} + 1) = 0.0295$	below 1% error
$Q(2N_{\max} + 2) = 146.18$	below 1% error
$T_3 = 1315.81$	below 1% error
$\lambda - \mu_{\text{eff}}(2N_{\max} + 2) = 0.0479$	below 5% error
$Q(2N_{\max} + 3) = 347.70$	below 1% error
$T_4 = 4207.09$	below 3% error
$T_5 = 314.37$	below 1% error
$\Pi = 6572.88$	below 3% error
$\Delta = 314.37$	below 1% error

Table 4.2. Comparison between analytical and simulation results for deterministic inter-arrival times $\frac{1}{\lambda} = 9$, deterministic service times $\frac{1}{\mu} = 0.9$ ($\rho = 0.1$), $P = 300$ and $N_{\max} = 12$.

reflect of reality. Everybody has its own perception of waiting time, which will directly affect the associated utility function. To further approach real situations and therefore to confer a more direct practical impact to our present modelling framework, it would now be required to actually characterize the underlying utility functions, a task which would obviously strongly depend on the particular situations to be investigated. Nevertheless, our model has so far the merit to allow for an analytical approach to a noticeably complex dynamics. It shows, once more, the structural emergence of macroscopic temporal patterns resulting from elementary, though nonlinear, individual interactions between autonomous agents. Similarly to ants which act according to the concentration of pheromones, here our agents decide in view of their waiting times (which actually depend on the other agents' previous actions within the system) and personal history with the server, which confers to our dynamics its stigmergic self-organizing character.

To close, let us emphasize that the very strong structural stability (*i.e.* the high insensitivity to external noise sources) of the oscillations reported in Figs. 4.1, 4.2 and 4.3 definitely increases the modelling power offered by this class of multi-agent nonlinear dynamics. Models that enjoy strong structural stability evolution are the cornerstones of a *synergetic* approach which, with a limited number of salient relevant features, are able to encompass under a common modelling framework a wide range of transdisciplinary situations.

4.6 Contributions of Chapter 4

- We explicitly consider the weariness that affects customers using a recurrent service. In that regard, we provide a stylized model of a queueing system with feedback where the decision to remain loyal to the service provider depends on both the suffered waiting time and the number of already received services. Like the models studied in Chapter 3, the framework proposed here enables an analytical description of the dynamics. The emerging collective pattern takes again here the form of self-sustained oscillations of the queue content, but they possess now an additional structure (*i.e.* an extra peak during the increasing phase).

Multiple-Stage Feedback Queueing Systems -
Networks with Competing Servers

Parrallel Servers with Feedback Loop - Stabilization by Noise

Summary. *Similarly to Chapters 3 and 4, we consider here queueing networks (QNs) with feedback loops roamed by “intelligent” agents, able to select their routing through the network on the basis of their measured waiting times at the QN nodes. Remember that it is an idealized modelling framework to discuss the dynamics of customers who base their loyalty to a service supplier on their individual waiting time satisfaction (i.e. they remain loyal provided they have waited less than a critical threshold). In this chapter, we consider a more complex network topology composed of two parallel service branches with feedback loop of the type previously introduced in Chapter 3. Concerning the customers’ initial choice between the two service providers, we introduce several different routing policies. These autonomous rules, which are wholly based on specific agent capabilities, differ in their level of complexity and with respect to the available amount of information about the current system state. Depending on the implemented routing policy at entrance, we show that the present two parallel servers multi-agent system possesses traffic flows that might exhibit various types of collective patterns. Even for this somehow simple network topology, the emergent cooperative behaviours manifest themselves via stable macroscopic temporal oscillations, synchronization of the queue contents or stabilization by noise phenomena. As already emphasized in the preceding chapters, for a wide range of control parameters, the underlying presence of the law of large numbers enables us to use deterministic evolution laws to analytically characterize such cooperative evolutions of our multi-agent system. In particular, we study in detail the case where the servers are sporadically subject to failures altering their ordinary behaviour.*

5.1 Introduction

In many real-life circumstances, customers face the situation of having to choose between different firms providing the same service. When this choice is taken by fresh customers, who haven’t any a priori ideas about the distinct service providers, the expected waiting time before receiving service would definitely represent an important decision criterion. In this chapter, we will consider such type of decision processes with customers having to choose among

two vendors providing a recurrent service. After this initial choice, the customers' loyalty to the chosen service provider depends as in the preceding chapters on the experimented waiting times. Inherent to such type of markets with different firms offering the same service activity, aspects such as competition or market partition between the service providers are obviously of great importance.

In this chapter as well as in Chapter 6, extending the modelling framework originally introduced in Chapter 3, we analyze more complex networks involving two servers with feedback loop topologies. More precisely, we will consider here an open network composed of two parallel servers of that particular type. While the agents' initial choice between the two servers will be based on a priori measures (*i.e.* expected value of the waiting times, based on real-time observations of the current queue contents), their later loyalty to this server will be based on a posteriori observations of the system state (*i.e.* individually perceived waiting times). Focusing on the study of the traffic load partition between the two service providers, we show that the behaviour of the present multi-agent system leads to the emergence of various types of collective temporal patterns. Note that in Chapter 7, we will treat the situations where the choice between the two service providers is taken not only with respect to expected waiting time considerations but also in function of spatial aspects.

This chapter is organized as follows. In Section 5.2, we introduce the network formed by two parallel feedback queues and a bifurcation point that will be considered throughout this chapter. Any incoming agent has to choose, on entry, between one of the two servers but once the choice made it can neither renege nor jockey between the queues. This initial routing decision can be either deterministic, random, and/or guided by a partial or a full observation of the real-time content of the queues. The new ability of the agents to observe queue contents introduced in this chapter, added to their capability to monitor waiting times (as previously described in this work), offers the possibility to generate new cooperative time evolutions. We first describe in Section 5.3 the system global behaviour when the choice between the two servers requires no special agent features and follows a fixed dispatching rule (deterministic or random). In Section 5.4, we focus on situations where the agents can observe the queue content of a single server and the decision to join the observed queue is correspondingly based on its content (the agent enters the queue if the content is below a critical population threshold). In this case, we show how the presence of random fluctuations in the service times can stabilize a flow dynamics which is otherwise unstable for purely deterministic service times. In particular, we provide an analytical study for the case where the servers are randomly and sporadically subject to failures altering their ordinary behaviour. In Section 5.5, we allow the agents to observe both queue contents. In this situation, when a "shortest-queue-first" scheduling rule is adopted at

the bifurcation node, a full synchronization of the queue contents oscillations is observed. Finally, Section 5.6 is devoted to concluding remarks.

5.2 Model

We increase in this chapter the complexity of the queueing network (QN) studied in the preceding chapters and pay attention to the configuration \mathcal{D} , formed by a dipole of feedback queues in an open network topology, as sketched in Fig. 5.1. Two feedback queueing systems of the type introduced in Chapter

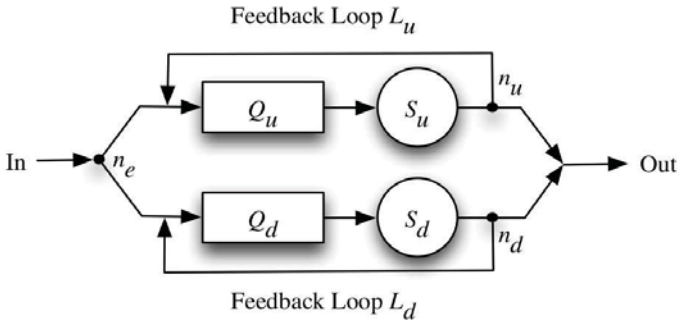


Fig. 5.1. A bifurcation of queueing systems with feedback loop.

3 are placed in parallel. The total incoming external customers feeding this system is a renewal process with rate Λ . At a first decision node (DN) n_e (where e stands for *entry*), the agents face two routing possibilities: to either join server S_u , or to join S_d . In front of S_u and S_d , the agents wait in queues whose respective contents will be denoted by $Q_u(t)$ and $Q_d(t)$ (from Fig. 5.1, the indexes u and d stand for *up* and *down* respectively). We will write by μ_u and μ_d the respective service rates of S_u and S_d . The capacities of both queues $Q_u(t)$ and $Q_d(t)$ are assumed to be unlimited and the service policy is first-in-first-out (FIFO). The presence of the feedback loops introduces two corresponding DNs n_u and n_d . As in Chapter 3, at n_u and n_d the decision to enter into the feedback loop depends on the sojourn time W individually measured by each customer. More precisely, the history-based (HB) routing rule \mathcal{R} given by Eq. (3.1) drives the agents' behaviour at these decision nodes. Concerning the initial choice between the two servers S_u and S_d on entry, we will separately consider in the following three typical scenarios depending on the agents' ability to gather information:

- (1) Fixed entrance dispatching rule requiring no special agent capabilities (Section 5.3).

- (2) Agent-based entrance dispatching rule based on a partial (individual) observation of the queue contents (Section 5.4).
- (3) Entrance dispatching rule based on a complete observation of the queue contents (Section 5.5).

In other words, the agents' capability, when choosing one of the two servers, to observe the state of the system will be progressively enhanced in the three above routing rules.

5.3 Fixed Entrance Dispatching Rule

Let us start with blind agents only being able to record the total waiting time spent to receive service (*i.e.* queueing + processing times) but unable to observe the queue contents $Q_u(t)$ and $Q_d(t)$. Hence, in this case, the routing decision at node n_e will not depend on any complex agent capability (“intelligence”) and an incoming customer will thus simply select between the servers S_u and S_d by using either a deterministic or a random rule, completely independent from the queue lengths $Q_u(t)$ and $Q_d(t)$. In the following, we describe the dynamics arising for two such typical dispatching rules, namely *deterministic polling* and *Bernoulli sampling*.

5.3.1 Deterministic Polling

In this case, the time horizon is divided into deterministic intervals T_u and T_d during which S_u and S_d respectively are alternatively fed with the total incoming traffic Λ . The conditions $\rho_u = \frac{T_u}{T_u+T_d} \cdot \frac{\Lambda}{\mu_u} < 1$ and $\rho_d = \frac{T_d}{T_u+T_d} \cdot \frac{\Lambda}{\mu_d} < 1$ ensure the stability of the system. In view of Chapter 3, it is not surprising that stable oscillations of the queue contents will, here again, be observed. However, instead of being smooth, the alternative feeding of the servers creates indentations in the time evolutions of $Q_u(t)$ and $Q_d(t)$. The frequency of the alternations, given by T_u and T_d , determines the indentation structure. Qualitatively, increasing the frequency of the alternations decreases the roughness of the curve. For large P , the amplitudes and frequencies of the two uncoupled oscillations can be determined using Eqs. (3.9) and (3.10) with the parameters $(\lambda_u = \frac{T_u\Lambda}{T_u+T_d}, \mu_u)$ on one hand and $(\lambda_d = \frac{T_d\Lambda}{T_u+T_d}, \mu_d)$ on the other hand. This is in perfect agreement with the simulation experiments given in Fig. 5.2.

5.3.2 Random Dispatching Rule

Here, we typically consider a Bernoulli sampling of the incoming flow, where the Bernoulli random variable is determined by a parameter r , ($0 \leq r \leq 1$). Each agent, upon arrival, simply draws a Bernoulli random variable and

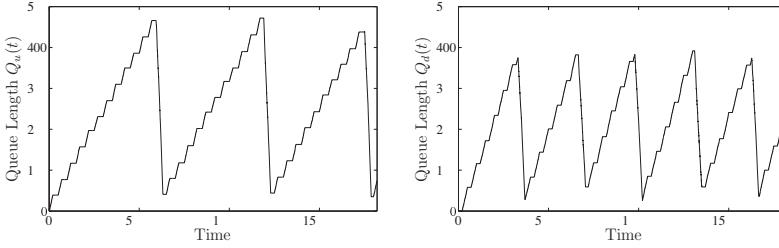


Fig. 5.2. Deterministic polling entrance rule: queue length indented oscillations obtained by simulation for exponentially distributed inter-arrival and service times with $\lambda = 0.2$, $\mu_u = 1.25$, $\mu_d = 1$, $T_u = 200$, $T_d = 300$ and $P = 350$.

chooses between the two servers according to this particular realization. From a global perspective, a partial traffic with rate $r\lambda$ enters into server S_u while a traffic with rate $(1-r)\lambda$ enters into S_d . The Bernoulli sampling implies that both systems S_u and S_d evolve independently and individually follow the dynamics exposed in Chapter 3. For large P , the two uncoupled, quasi-deterministic, cyclo-stationary behaviours have amplitudes and frequencies that are given by Eqs. (3.9) and (3.10) with parameters $(r\lambda, \mu_u)$ and $((1-r)\lambda, \mu_d)$ respectively.

5.4 Entrance Dispatching Based on a Partial Observation of the Queue Contents - Noise Induced Stabilization

Besides chronometers to record W , each customer is now endowed with a “visual system”¹ enabling him/her to observe, in real-time, the instantaneous queue content $Q_u(t)$. Assume however that $Q_d(t)$ always remains hidden to the incoming agents, although they know the average service rate μ_d . At time t , an incoming agent at node n_e first observes the queue content $Q_u(t)$ and, based on this observation, decides either to enter S_u or to join S_d . Once entered into a queue, neither reneging nor jockeying (*i.e.* jumping between S_u and S_d) is allowed. Note that except for the presence of feedback loops, this network configuration is fully similar to the two gas stations network studied in [54]. In this contribution, two gas stations are located one after the other on a main road. A driver who needs to refuel is only able to observe the queue length $Q_u(t)$ at the first station (which would be here S_u). Then, he/she compares $Q_u(t)$ to the conditional expected queue content at the second station (here S_d) and decides either to enter into the first station or to postpone

¹ Like the individual computation of the waiting times, such vision capability could be concretely realized with an ability of the agents to read the (permanently updated) value of dedicated registers.

his/her refuelling and enter into the second one.

Returning to our present model, we assume from now on that an incoming agent decides:

- (i) either to enter S_u whenever $Q_u(t)$ strictly stays below a threshold value N^* (*i.e.* when $Q_u(t) < N^*$)
- (ii) or to enter S_d otherwise.

At the DNs n_u and n_d , the routing rule \mathcal{R} given by Eq. (3.1) depends, as in Chapter 3, on a patience parameter P which is again assumed to be common to all agents. Assume that the patience P and the threshold control parameter N^* are adjusted as:

$$P \geq \frac{N^* + 1 + \delta}{\mu_u}, \quad \delta \in \mathbb{N}^+, \quad (5.1)$$

where δ denotes a tolerance level above the expected sojourn time. We can interpret (further details will be given later) the routing decision at n_u to be a formal illustration of the H. Maister's first principle of the psychology of waiting lines [87], namely: "*Satisfaction equals perception minus expectation*". Indeed, at the DN n_e , the level N^* defines, via P as given by Eq. (5.1), an expected admissible sojourn time. Later, when reaching n_u , each agent compares its measured sojourn time (playing the role of the *perceived* sojourn time) with P (playing the role of the *expected* sojourn time) and then takes its routing decision accordingly.

We consider first the deterministic dynamics where S_u operates with a fixed service time $\frac{1}{\mu_u}$. When, at a given time t , $Q_u(t) = N^*$ agents are waiting in front of S_u , they will remain loyal to S_u forever (*i.e.* these agents will loop forever and ever). Indeed, their measured sojourn time W never exceeds P and, the dynamics being deterministic, no perturbation will alter this dynamically "frozen" situation. As a result, once $Q_u(t) = N^*$, the server S_u remains definitively unavailable for any external incomer and the global incoming traffic with rate Λ is entirely dispatched to S_d . Whenever $\frac{\Lambda}{\mu_d} > 1$, the queueing system will thus be unstable (*i.e.* $\lim_{t \rightarrow \infty} Q_d(t) = \infty$).

Assume now that random fluctuations affect the service times of S_u . While Eq. (5.1) is still satisfied on average, service time noise triggers, at node n_u , a random flow of unsatisfied customers, who will definitively leave the system. Hence, with the presence of noise (in the service time), the availability (for external traffic) of S_u effectively increases - remember that this availability is null in absence of noise. Consequently, a part of the global incoming traffic is now processed by S_u . For a selected range of control parameters, we may simultaneously have:

$$\rho_u = \frac{\alpha A}{\mu_u} < 1 \quad \text{and} \quad \rho_d = \frac{(1 - \alpha)A}{\mu_d} < 1, \quad 0 \leq \alpha \leq 1, \quad (5.2)$$

where αA and $(1 - \alpha)A$ stand for the stationary average rates of the partial traffic flows feeding S_u and S_d ($\alpha = 0$ corresponds to the purely deterministic case considered before). Whenever Eq. (5.2) holds, both queueing branches are dynamically stable. The previous qualitative reasoning suggests that there exists a critical variance $\sigma_{u,c}^2$ of the service times of S_u (and hence a critical value α_c) such that:

- (a) for $\sigma_u^2 \geq \sigma_{u,c}^2$, the queueing system is stable.
- (b) for $\sigma_u^2 < \sigma_{u,c}^2$, the queueing system is unstable.

Hence, we can speak here of a *noise-induced stabilization of the dynamics*, which is studied below in more details both experimentally and analytically.

5.4.1 Experimental Observations

The above dynamical behaviour can be explicitly observed in simulation experiments where the incoming flow of customers is an exponential process with parameter Λ and the S_u service times are drawn from a probability density $dB_u(x)$ being:

- (1) uniform with support $\left[\frac{1}{\mu_u} - \xi, \frac{1}{\mu_u} + \xi\right]$ with $\xi \geq 0$ (thus $\sigma_u^2 = \frac{\xi^2}{3}$). The following numerical values were used: $\Lambda = 1.11$, $\frac{1}{\mu_u} = \frac{1}{\mu_d} = 1$, $N^* = 28$ and $P = 30$ (*i.e.* $\delta = 1$ in Eq. (5.1)). We observe that for $\xi \geq 0.118 \Rightarrow \sigma_{u,c}^2 \geq 0.0046$, the queueing system remains stable, while it becomes unstable (*i.e.* $\lim_{t \rightarrow \infty} Q_d(t) = \infty$) for smaller values of ξ .
- (2) a Normal law $\mathcal{N}\left(\frac{1}{\mu_u}, \sigma_u^2\right)$. For the same numerical values as above, we observe that for $\sigma_u^2 \geq \sigma_{u,c}^2 = 0.0046$, the queueing system remains stable, while it becomes unstable for $\sigma_u^2 < \sigma_{u,c}^2$.

5.4.2 Analytical Approach

To analytically discuss the stability issue, let us consider the situation where the service times of S_u are independent Bernoulli random variables with values $\left\{\frac{1}{\mu_u}, \frac{1}{\mu^+}\right\}$ and corresponding probabilities $(1 - q)$ and q respectively, $0 \leq q \ll 1$. We assume that $\mu^+ < \mu_u$ and interpret $\frac{1}{\mu^+}$ (with $\frac{1}{\mu^+} > \frac{1}{\mu_u}$) as the effective service time occurring when a failure alters the ordinary behaviour of the server S_u . Remember that the agents follow the FIFO rule and are homogeneous in their patience parameter P , chosen here to fulfil:

$$P < \frac{N^*}{\mu_u} + \frac{1}{\mu^+} \quad \text{and} \quad P > \frac{N^* + 1}{\mu_u}, \quad (5.3)$$

where the second expression is actually Eq. (5.1) with $\delta = 1$ (compared to Eq. (5.1), note here the use of a strict rather than a weak inequality). When,

at a given time t , $Q_u(t) = N^* - 1$, an incoming tagged customer ζ at DN n_e will decide to enter S_u . Later on when ζ reaches n_u , he/she will, according to Eq. (5.3) and the routing rule \mathcal{R} given by Eq. (3.1), choose:

- (i) either to follow the feedback loop, whenever no failure occurred during the service of the N^* customers who were directly in front of him (including the customer who was served when ζ joined $Q_u(t)$) and during his/her own service
- (ii) or to leave the system, whenever one or more failures occurred during the service of the N^* customers who were directly in front of him/her and during his/her own service.

Hence, in absence of failures and when $Q_u(t) = N^*$, the agents will remain in the feedback loop forever and, at DNs n_e and n_u , neither an externally new incomer nor a leaving customer will be observed. However, as soon as failures occur in S_u , Eq. (5.3) implies that one or more customers will definitively leave the system after the decision at n_u . Hence, this implies that the global incoming traffic will now be shared between S_u and S_d . Assume that:

$$\mu_d < \Lambda \quad \iff \quad \rho_d = \frac{\Lambda}{\mu_d} > 1. \quad (5.4)$$

Thus, S_d cannot sustain alone the full traffic load without being in an unstable regime ($\rho_d > 1 \Rightarrow \lim_{t \rightarrow \infty} Q_d(t) = \infty$). Remember that $\alpha\Lambda$ and $(1 - \alpha)\Lambda$ denote the rates of the average partial traffics processed respectively by S_u and S_d . There exists a critical incoming flow, defined by $(1 - \alpha_c)\Lambda$, above which the queue $Q_d(t)$ becomes unstable. For the associated traffic intensities, this implies that:

$$\rho_u = \frac{\alpha\Lambda}{\mu_u} < 1 \quad \text{and} \quad \rho_d = \frac{(1 - \alpha)\Lambda}{\mu_d} < 1, \quad \forall \alpha > \alpha_c, \quad (5.5)$$

$$\rho_{d,c} = \frac{(1 - \alpha_c)\Lambda}{\mu_d} = 1, \quad (5.6)$$

where $\rho_{d,c}$ is the critical traffic load driving the queue $Q_d(t)$ to its marginal stability regime.

To proceed further with analytical considerations, let us now focus on rare events regimes (RER), for which more than a single failure during $N^* + 1$ consecutive ordinary services is a highly improbable event. As N^* is the threshold value governing the decision at node n_e and P fulfils Eq. (5.3), the RER is expected when $N^* + 1 \ll \frac{1}{q}$. Under the RER, each failure triggers the purging of the queue $Q_u(t)$. Indeed, due to the FIFO scheduling rule, when a failure occurs at time t , all the N^* agents at the moment in $Q_u(t)$ will experiment a sojourn time larger than P when arriving at n_u (*i.e.* these are the loyal customers currently travelling in the feedback loop and that will feed server S_u for $t' > t$). As it has been discussed in Section 3.4, this produces a *siphon*

purging, here of size N^* . In the RER, the succession of these siphon events will be approximately uncorrelated. Hence, in the stationary regime, we can simply estimate the outgoing flow rate λ_u at DN n_u as being given by:

$$\lambda_u = \text{Prob}\{a \text{ single failure occurs}\} N^* \mu_u = q N^* \mu_u. \quad (5.7)$$

When Eq. (5.7) holds, the partial traffic on S_d is given by:

$$\rho_d = \frac{\lambda_d}{\mu_d} = \frac{\Lambda - \lambda_u}{\mu_d} = \frac{\Lambda - q N^* \mu_u}{\mu_d}. \quad (5.8)$$

The marginal stability of queue $Q_d(t)$ is attained at the critical traffic $\rho_d = \rho_{d,c} = 1$, which implies:

$$q \geq q_c := \frac{\Lambda - \mu_d}{N^* \mu_u}. \quad (5.9)$$

In terms of α_c , we can write:

$$\alpha_c = 1 - \frac{\mu_d}{\Lambda}. \quad (5.10)$$

Finally, we can also express the stability condition given by Eq. (5.9) in terms of the critical variance $\sigma_{u,c}^2$ of the underlying Bernoulli random variable. We obtain:

$$\sigma_u^2 \geq \sigma_{u,c}^2 = q_c(1 - q_c) \left(\frac{1}{\mu^+} - \frac{1}{\mu_u} \right)^2. \quad (5.11)$$

The numerical experiments reported in Tab. 5.1 are in perfect agreement with Eqs. (5.9) to (5.11).

Global incoming traffic Λ	Simulated stability condition on q	Simulated stability condition on σ_u^2
1.05	0.0017	0.00075
1.1	0.0034	0.0015

Table 5.1. Stability conditions obtained when using a discrete events simulator with the following parameters: $N^* = 28$, $\frac{1}{\mu_d} = \frac{1}{\mu_u} = 1$, $\frac{1}{\mu^+} = 3$ and $P = 30$. No discrepancy between simulated and theoretical results have been observed up to the shown precision.

While the concept of stabilization by noise is currently discussed in the context of stochastic differential equations [8, 53, 107], the present class of models exemplifies clearly that such a random stabilization can also be encountered in multi-agent systems where a nonlinearity (in our case, the feedback loop) is present.

5.5 Flow Dispatching Based on Fully Observable Queues - Synchronization of Oscillations

Here, we assume that both queues $Q_u(t)$ and $Q_d(t)$ can be observed simultaneously by the incoming agents. Thus, compared with Section 5.4, the information gathering process has been further increased. Based on such real-time observation of both queue contents, several dispatching policies at the DN n_e can be constructed. Among the simplest and most natural rules, let us here focus on the policy sending a new externally incoming customer to the shortest observed queue. This *Shortest-Queue-First* (SQF) rule yields the natural emergence, for large common patience parameter P , of *synchronized stable temporal oscillations* of the queue contents $Q_u(t)$ and $Q_d(t)$. This happens for any initial conditions of the queue populations. As before, when P is large and common to all agents, a purely deterministic approach is perfectly suitable to describe the emerging oscillatory behaviours. We assume that $\frac{\lambda}{\mu_u + \mu_d} < 1$ to ensure the stability of the system. Let us consider, without loss of generality, that $\frac{1}{\mu_u} \geq \frac{1}{\mu_d}$. The two following cases may arise:

(1) *Non-generic case* : two identical servers (*i.e.* $\frac{1}{\mu_u} = \frac{1}{\mu_d}$).

The total incoming traffic is evenly divided between the two servers, both receiving a partial traffic with rate $\frac{\lambda}{2}$. The amplitude and period of the common synchronized stable temporal oscillations of the queue contents $Q_u(t)$ and $Q_d(t)$ are given by Eqs. (3.9) and (3.10) with parameters $\frac{\lambda}{2}$ and μ .

(2) *Generic case* : two servers with service rate ratio $\frac{1}{\mu_u} > \frac{1}{\mu_d}$.

Even though the servers do not work at the same speed, the queue contents $Q_u(t)$ and $Q_d(t)$ are driven to be equal at any time, provided $\frac{\lambda}{\mu_d} > 1$ (*i.e.* upon the condition that S_d is not able to handle alone the total incoming flow). The greater speed of S_d implies that the customers joining this server will remain satisfied for a longer queue length than with S_u . As a consequence of the SQF rule, there will be more unsatisfied customers with server S_u and this server will thus process a greater part of the global incoming traffic than S_d (*i.e.* S_u will absorb more fresh customers, but these customers will stay less time in the system than those joining S_d). As shown in Fig. 5.3, two distinct dynamics may emerge depending on the arrival and service rates.

5.6 Concluding Remarks

To characterize and to quantify how an incoming flow of customers is effectively shared between providers offering the same service is a topic of great importance in many different businesses. In this chapter, we have focused on the influence that expected waiting times before being served could have

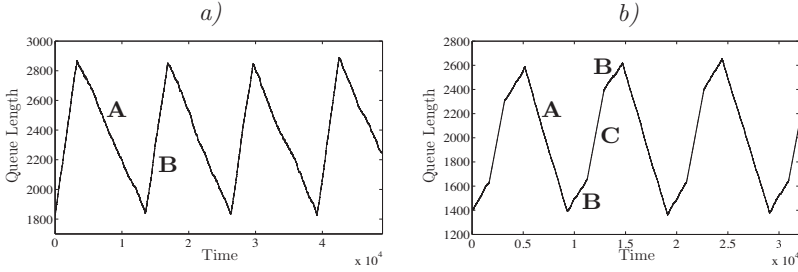


Fig. 5.3. The SQF policy implies that we have with large probability $0 \leq |Q_u(t) - Q_d(t)| \leq 1, \forall t$. Here, we only show the state of $Q_u(t)$ in the above figures. Fig.5.3.a): Queue content $Q_u(t)$ when $P = 2500$ and the inter-arrival and service processes are exponential with parameters $\Lambda = 1.25$, $\frac{1}{\mu_u} = 1.6$ and $\frac{1}{\mu_d} = 1.2$. The amplitude and period of the common synchronized stable temporal oscillations are given by $\Delta = \frac{P\mu_d}{2}$ and $\Pi = P \left(\frac{\mu_d}{\mu_d + \mu_u - \Lambda} + \frac{\mu_d}{\Lambda - \mu_u} \right)$ respectively. The two different slopes are given by **A** = $\frac{\Lambda - \mu_u - \mu_d}{2}$ and **B** = $\frac{\Lambda - \mu_u}{2}$. Fig.5.3.b): Queue content $Q_u(t)$ when $P = 2500$ and the inter-arrival and service processes are exponential with parameters $\Lambda = 0.9$, $\frac{1}{\mu_u} = 1.6$ and $\frac{1}{\mu_d} = 1.2$. The dynamics differs from the case a) by the presence of a time interval with slope **C** = $\frac{\Lambda}{2}$. During this interval, customers in S_u and S_d are all satisfied. On the other hand, during the time intervals with slope **A** and **B**, the customers in S_u are unsatisfied (the customers in S_d being unsatisfied only during the interval with slope **A**). For instance, in the configuration a), all the customers joining S_u are unsatisfied, because $Q_u(t)$ always remains above the critical threshold. The complexity of the dynamics in the case b) requires more involved computations, which precludes to give simple and compact expressions for the amplitude and the period of the synchronized oscillations. However, due to the quasi-deterministic nature of the dynamics (when P is sufficiently large), an analytical characterization is still feasible.

on the customers' decision processes. In particular, we have exhibited possible collective patterns that might emerge due to such kind of autonomous customers' routing rules, namely stable quasi-deterministic oscillations of the queue contents, stabilization by noise phenomena and synchronization of oscillations. While the fully analytical study of Section 5.4.2 is restricted to regimes where undesired events (such as failures) are rare (*i.e.* sporadic temporary unavailability of the servers), simulation experiments however clearly indicate that the described self-organized phenomena remain observable for very general situations.

5.7 Contributions of Chapter 5

- Extending the modelling framework introduced in the preceding chapters, we consider here an open network composed of two parallel servers with

feedback loop like those previously described in Chapter 3. We introduce various (with respect to their specific level of complexity as well as the available amount of individually collected information about the current system state) agent-based policies for the customers' initial choice between the two servers. Added to the agents' capability to autonomously choose their routing with regard to waiting time considerations (hence in function of history-based data), we show that, depending on the routing rule implemented at entrance, the emerging dynamics might exhibit collective phenomena like stable quasi-deterministic oscillations of the queue lengths, noise-induced stabilization and synchronization of oscillations.

Closed Network of Two Servers with Feedback Loop - Competing Server Dynamics

Summary. *We study the market partition dynamics between two recurrent service providers in a closed market topology. As generally assumed all along this thesis, the customers' satisfaction (and hence their loyalty to a service provider) depends in a nonlinear fashion on the elapsed waiting time to receive service. In the present context, when unsatisfied with a service provider, a customer leaves it and join the other vendor. Without giving tedious computational details, we provide in this chapter a qualitative overview of the variety of temporal collective phenomena that might emerge in the dynamics of the considered multi-agent two servers queueing system.*

6.1 Introduction

Managers take strategic and investment decisions in order to anticipate the future evolution of the market. In that respect, they can rely either on statistical analysis or on dynamical modelling. The first approach is suitable, comfortable and effective in order to perform long-term strategic planning. On the other hand, to focus on the dynamical evolution of the market allows for a deeper understanding of the inherent business processes, which consequently allows for the implementation of strategic short-term policies that lead to quicker and more reactive adaptation to possible market fluctuations. In this chapter, and more generally in this thesis, we follow the second modelling approach and hence we try to describe the temporal evolution of our particular classes of multi-agent queueing systems.

When a market is composed of several service providers, they will compete to become more attractive than the others and consequently attract a larger number of customers. When the gain of a new customer at a server coincides with the loss of a consumer for an other service provider, we speak of *cannibalization* effect. This phenomenon is likely to happen in closed market configurations. Moreover, when a server attracts and cannibalizes customers

such that it holds a market share that is larger than the expected equilibrium¹ (*i.e.* there exists an asymmetry between the number of customers at each vendor with respect to their respective service speeds), we speak in that case of *super-cannibalization*. The closed network with two recurrent service facilities considered in this chapter, whose dynamics will exhibit such cannibalization and super-cannibalization effects, idealizes the quality of service competition which arises between two distinct vendors within the same market.

Let us take the example of two internet service providers, namely S_1 and S_2 , competing for a fixed pool of potential customers. If, at a given moment, S_2 proposes a faster access to the internet and/or a larger bandwidth, it will start to attract customers from the slower service provider S_1 , triggering thus a cannibalization phenomenon. However, due to the presence of an informational delay in the dynamics, it is likely that too many customers will actually leave S_1 for S_2 ; this will yield the creation of a super-cannibalization effect. The above equilibrium market share held by S_2 will imply that its processing rate will noticeably decrease and S_2 will ultimately become less attractive than S_1 . As a consequence, this will trigger the apparition of a flow of customers that take the inverse way from S_2 to S_1 , creating thus an antagonistic cannibalization effect. Following this simple reasoning, the dynamics of such type of systems might hence exhibit successive and periodic super-cannibalization phenomena.

The chapter is organized as follows. In Section 6.2, we introduce the closed network with two recurrent service facilities that will be considered throughout this chapter. In Section 6.3, we give a qualitative panorama of the range of collective patterns that might appear due to the agents' actions and interactions within our particular queueing system. More precisely, we observe cannibalization as well as super-cannibalization effects in the dynamics, which ultimately lead either to oscillations between the respective number of agents at each server or to the stabilization of the number of customers at each node. Several perspectives and conclusions are given in Section 6.4.

6.2 Model

Consider the closed network sketched in Fig. 6.1, formed by two feedback queueing systems as those discussed in Chapter 3. The servers S_i ($i = 1, 2$) composing the network have an average service time $\frac{1}{\mu_i}$, where μ_i ($i = 1, 2$) stand for the respective service rates. Without loss of generality, we assume that $\mu_1 \leq \mu_2$. The total number $N \in \mathbb{N}^+$ of agents circulating in the network is fixed and we allow the capacities C_i ($i = 1, 2$) of both queues to be large

¹ In equilibrium, each server should hold a market share that is proportional to its processing rate.

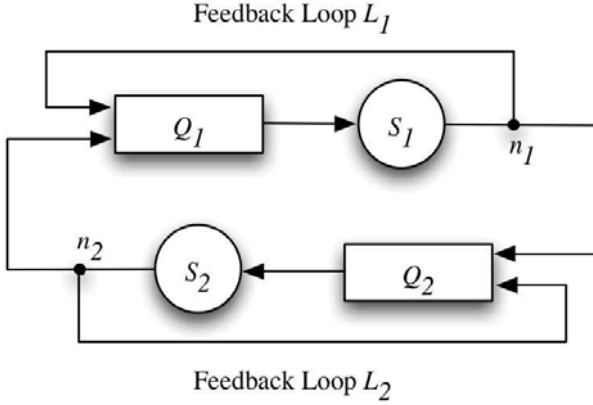


Fig. 6.1. Closed network of two servers with feedback loop.

enough to accommodate the entire population (*i.e.* one assumes that $C_i > N$ for $i = 1, 2$). Directly inspired from Chapter 3, each circulating customer is again equipped with a clock and monitors its total waiting time spent to receive service (*i.e.* its sojourn time) - the clock is reset to $t = 0$ each time a customer enters into a queue, the clock time value W_i is obtained when reaching a node n_i ($i = 1, 2$). The measured value W_i of each customer is then compared with a fixed and common to all customers patience parameter P . Thanks to the time monitoring, the history-based (HB) routing rule \mathcal{R} , originally introduced in Eq. (3.1), can be implemented here and becomes in the present context:

$$\mathcal{R} = \begin{cases} \text{Go to the feedback loop} & \text{if } W_i \leq P, \\ \quad \text{and hence go to server } S_i & \\ \text{Avoid the feedback loop} & \text{if } W_i > P, \\ \quad \text{and hence go to server } S_{\lceil i+1 \rceil} & \end{cases} \quad (6.1)$$

with the notation:

$$S_{\lceil i+1 \rceil} = \begin{cases} S_1 & \text{if } i = 2, \\ S_2 & \text{if } i = 1. \end{cases}$$

Writing $N_1(t)$ and $N_2(t) = N - N_1(t)$ for the number of customers (including the one being served) waiting in $Q_1(t)$ and $Q_2(t)$ respectively and using a fluid queueing picture to describe the population dynamics, we can write:

$$N_1(t) = N_1(0) +$$

$$\int_0^t [\mu_2 \mathbb{I}(W_2(s) \geq P) \mathbb{I}(N_1(s) < N) - \mu_1 \mathbb{I}(W_1(s) \geq P) \mathbb{I}(N_1(s) > 0)] ds, \quad (6.2)$$

where $W_i(t)$ is the sojourn time of customers leaving queue Q_i , $i = 1, 2$, at time t . The function $\mathbb{I}\{E\}$ is the indicator of the event E (i.e. $\mathbb{I}\{E\} \equiv 1$ when the event $\{E\}$ is realized and 0 otherwise). We assume, from now on, that the common patience parameter P is chosen such that $P_{\min} \leq P \leq P_{\max}$, with P_{\min} being large enough to safely allow, as in Chapter 3, a deterministic analysis (thanks to the influence of the law of large numbers). Furthermore, P_{\max} fulfils:

$$\mu_1 P_{\max} < N, \quad (6.3)$$

which, in view of the policy \mathcal{R} stated in Eq. (6.1), rules out the trivial situations occurring when all the customers are systematically satisfied and therefore stay loyal to their initial server.

6.3 Market Partition Dynamics

The total number of customers being fixed, the dynamical state of the system is entirely determined by the single variable $N_1(t)$. While, for the deterministic evolution, a complete analytical characterization of the emerging dynamics is possible [77], we present here only the most relevant qualitative features exhibited by this dynamical system. More particularly, four separated regimes, summarized in Fig. 6.2, can be characterized in function of the initial queue content $N_1(0)$:

- **N -regime:** $N_1(0) > \mu_1 P$ and $N_1(0) \geq (N - \mu_2 P)$

Starting at time $t = 0$ with $N_1(0) > \mu_1 P$ customers lining in Q_1 , we conclude that after time $t = P$, customers reaching the node n_1 will be unsatisfied and therefore leave to populate Q_2 . To study the role played by the second condition: $N_1(0) \geq (N - \mu_2 P) \Leftrightarrow (N - N_1(0)) = N_2(0) \leq \mu_2 P$, two sub-cases have to be examined:

(1) *Generic case:* $\mu_1 < \mu_2$

(a) When $\mu_2 P \geq N$, the server S_2 alone is able to accommodate, with satisfaction, all customers. Consider the situation where $N_1(0)$ customers initially populate Q_1 . As $N_1(0) > \mu_1 P$, the first $\mu_1 P$ customers reaching the node n_1 will be satisfied and therefore return to line again in Q_1 . The remaining $N_1(0) - \mu_1 P$ customers, being unsatisfied, go to line in Q_2 . At their second visit to S_1 , the $\mu_1 P$ customers initially satisfied will, when reaching n_1 for the second time, be unsatisfied. Indeed, they did effectively wait $\frac{N_1(0)}{\mu_1} > P$ in Q_1 to receive their second service. This mechanism implies that ultimately all customers leave the first node and populate Q_2 and stay there forever. This behaviour can be used to mimic how an efficient service can ultimately monopolize an entire market sector.

(b) In the case where $\mu_2 P < N$, the server S_2 alone is not able to accommodate, with satisfaction, all customers. Hence, temporal oscillations of the queue contents will be sustained (*i.e.* periodic antagonistic super-cannibalization phenomena).

(2) *Non-generic case: $\mu_1 = \mu_2$*

In this case, none of the servers is able to accommodate, with satisfaction, the entire population. Hence, only oscillating regimes are generated.

- **\mathcal{W} -regime: $N_1(0) > \mu_1 P$ and $N_1(0) < (N - \mu_2 P)$**

Starting with $N_1(0) > \mu_1 P$ customers in Q_1 implies that at time $t = P$, customers leave Q_1 to enter into Q_2 . Again, to study the role played by the second condition: $N_1(0) < (N - \mu_2 P) \Leftrightarrow (N - N_1(0)) = N_2(0) > \mu_2 P$, two sub-cases have to be examined:

(1) *Generic case: $\mu_1 < \mu_2$*

The second condition stated above implies similarly that, at time $t = P$, customers leave Q_2 to enter into Q_1 . Hence, unsatisfied customers are systematically generated.

(2) *Non-generic case: $\mu_1 = \mu_2$*

Here, one can show that the queue contents remain constant, as customers travel from Q_1 to Q_2 in a similar way, exactly as they would do in a closed tandem fluid queue without feedback.

- **\mathcal{S} -regime: $N_1(0) \leq \mu_1 P$ and $N_1(0) < (N - \mu_2 P)$**

Starting with $N_1(0) \leq \mu_1 P$ implies that customers initially in Q_1 are satisfied and stay in that queue. The second condition: $N_1(0) < (N - \mu_2 P) \Leftrightarrow (N - N_1(0)) = N_2(0) > \mu_2 P$ implies that, for $t = P$, customers in Q_2 are unsatisfied and therefore leave to populate Q_1 . The discussion of this case is very similar to the \mathcal{N} -regime. Here, however, due to the fact that $N \geq N_2(0) > \mu_2 P$, the server S_2 will never be able to attract the entire market and therefore, only oscillating behaviours of the queue contents are observable.

- **\mathcal{E} -regime: $N_1(0) \leq \mu_1 P$ and $N_1(0) \geq (N - \mu_2 P)$**

Starting with $N_1(0) \leq \mu_1 P$ implies that customers initially in Q_1 are satisfied and stay in Q_1 . The second condition: $N_1(0) \geq (N - \mu_2 P) \Leftrightarrow (N - N_1(0)) = N_2(0) \leq \mu_2 P$ implies similarly that customers are satisfied and therefore stay in Q_2 . In this regime, no exchange of customers is

observed and the system effectively behaves as if it were formed by two independent servers.

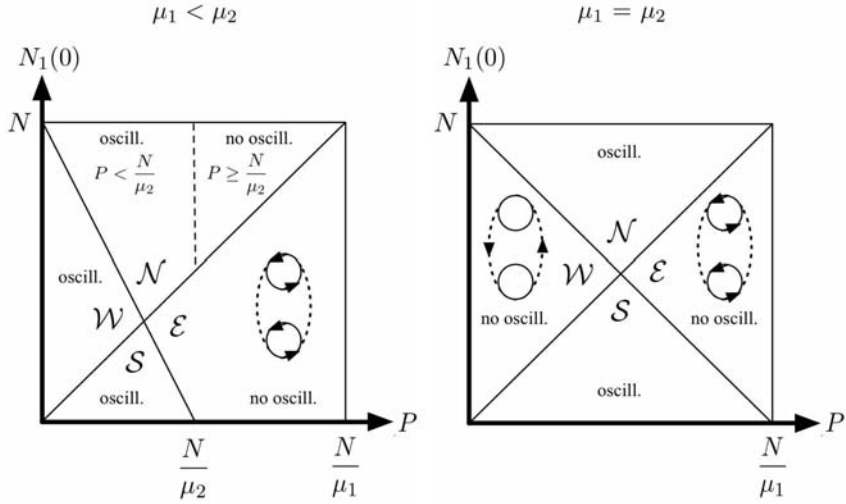


Fig. 6.2. Summary of the different regimes obtained in function of the initial condition and the patience parameter.

More refined analytical results and the rigorous proofs of the above results rely on taking into account the ability for the system (*i.e.* both servers) to be able to accommodate, with satisfaction, the entire population of customers. This can be discussed by introducing the parameter $\Theta = N - (\mu_1 + \mu_2)P$. For $\Theta > 0$, a systematic flow of unsatisfied customers (and hence a recurrent cannibalization effect) will always be generated in the system. On the other hand, when $\Theta < 0$, the presence of such systematic flow of unsatisfied customers depends on the initial conditions (this explicit ergodicity violation is due to the intrinsic non-Markovian property of the underlying dynamics) and thus distinct types of regimes might emerge. This competing servers dynamics illustrates explicitly how rich spatio-temporal structures can be generated by HB routing in queueing networks.

In the present context, the various collective emerging patterns take the form either of stable oscillations between the respective number of customers populating each server or of a stabilization of the number of consumers at each node. These phenomena are the consequence of inherent antagonistic cannibalization and super-cannibalization effects between the two servers. Note that these periodic phenomena would be removed in presence of less reactive

agents that base their satisfaction with a service provider not only on their last visit but on a larger history window. In particular, when the customers do not react directly (*i.e.* they do not leave a service provider) after a single bad experience but wait until having received a given number of unsatisfactory experiments before changing of vendor, then no periodic flow of unsatisfied customers will be generated and the number of agents in each node will be stabilized.

6.4 Concluding Remarks

The class of models considered in this chapter fits particularly well to e-commerce applications. Indeed, fickle customers, with immediate reactions to short-term satisfaction measures, are likely to be observed in this framework. The possibilities that consumers have to make a rapid comparison between different available service providers is enhanced within such market channels. In particular, customers often have a direct access to information about the current expected delay to receive service and they can thus easily establish whether or not it is profitable to change of vendor. Among the potential directions to extend the idealized modelling framework proposed in this chapter, one could introduce a cost induced by a change of service provider. Depending on the relative importance of the costs (*i.e.* vendor switch and waiting time costs), it would be interesting to study the resulting modified market partition dynamics.

6.5 Contributions of Chapter 6

- We study the market partition dynamics of a closed network composed of two recurrent services (*i.e.* two servers with feedback loop). We qualitatively describe the various collective patterns that might emerge in such a queueing system, namely cannibalization phenomena that lead to periodic flows of unsatisfied customers changing of service provider (*i.e.* oscillations of the respective number of customers at each server) or stabilization of the number of consumers at each node.

**Introducing Spatial Aspects - Market Sharing
Spatio-Temporal Dynamics**

Spatial Market Sharing Dynamics Between Two Service Providers

Summary. *In this chapter, we study the market partition between two distinct firms that deliver services to customers whose behaviour is sensitive to waiting time. In our model, the incoming customers select a firm on the basis of its posted price, the expected waiting time and its brand. More specifically, we quantify by a cost any departure from the ideal brand expected by each incoming customer. Considering that the two underlying queueing processes operate under high traffic regimes, we analyze the market sharing dynamics by using a diffusion process. As a function of control parameters, such as the waiting and brand departure costs or the incoming traffic intensity, we are able to analytically characterize a transition between an Hotelling-like regime (dominated by brand considerations) and a deadline type regime (dominated by waiting time considerations). The market sharing dynamics is described by the time evolution of a boundary point, the dynamics of which belongs to the class of noise-induced phase transitions, a type of phenomena so far widely discussed in physics, chemistry and biology. Explicit illustrations for both symmetric (i.e. identical servers) and asymmetric cases are worked out.*

7.1 Introduction

In his original contribution [62], H. Hotelling considered the case where two vendors supply an identical product that is perceived homogeneous by incoming customers. However, the vendors being separated in geographical space, transportation costs to be added to the mill prices charged by the vendors are generated. In presence of two vendors, there exists an inner market *boundary point*, for which the mill price plus the transportation costs from both suppliers break even. This seminal modelling framework has stimulated a wealth of contributions with the goal to relax some of the oversimplifying hypothesis of the original model. In particular, the introduction of *elastic demands* (i.e. customers are not prepared to pay “prohibitive prices” for the product) has been discussed in [105]. Note that the original Hotelling’s model is basically deterministic - it indeed does not incorporate random perturbations which actually may corrupt the prices and then affect the customers’ decision process.

Among the numerous potential noise sources, one of the simplest and most natural way to incorporate randomness is to consider the situations where the customers' decision to select one of the vendors depends on the expected time delay before service. This simple and realistic situation has been recently proposed by G. Cachon and P. Harker in [24] and [25]. As these authors clearly emphasized in [24], the resulting inherent analytical complexity implies that rather seldom are the models dealing with firms that simultaneously compete with both prices and processing rates. In this regard, the aim of the present work is to investigate a class of simple models for which explicit analytical considerations can be worked out. While in [24] the firms are assumed to adjust their processing rates to guarantee a fixed expected time cost, the class of models developed in this chapter takes into account the fluctuations of the waiting times and therefore keeps full track of the randomness induced by the underlying queueing processes. Note that the adjusting processing rates assumption proposed in [24] allows a discussion based only on averages. Contrary to [24], where no variance effects enter into the description of the model (*i.e.* this is effectively a “pseudo-stochastic” model), our approach explicitly emphasizes the role played by the variance of the fluctuations - also called in the sequel the “volatility” of the underlying noise sources. As discussed in [55], the introduction of waiting costs in the dynamics of queueing systems leads to the concept of *externalities* (*i.e.* the costs induced on later incomers by a customer who is just joining the queue). In the class of models to be discussed here, these externalities trigger the random dynamics controlling the boundary point which defines the market partition. While, for classical Hotelling-like models available in the literature, the interest is often paid directly on the competition aspects existing between the service providers (see for instance [25, 34, 105, 108]), in the present study we exclusively focus on the market sharing dynamics.

Service models where both distance and quality of service enter into consideration find, among others, a practical framework in the secondary health care market. More precisely, let us consider patients who wait for non-urgent operations, that can be mid-term planned. As said in [96], where an application of the standard Hotelling model to the secondary health care market is proposed, patients may accept meeting monetary and non-monetary costs inherent to distance, if they expect a positive return in terms of enhanced quality of service. Furthermore, while the quality of service perceived by the patients combines several different aspects, it clearly includes the time to wait for the operation to take place. Another typical situation will be met when car drivers entering into a city centre are offered alternative choices between several parking lots (here we focus on two lots). It is nowadays common to post in real-time, at the entrance of the city, the number of available parking spaces of each parking lot. The actual time required to complete a parking action, which here plays the role of the waiting time, is clearly monotonously decreasing with the number of available spaces of the parking lot. Hence, the

selection of the best parking lot does not only depend on its location, but also on its current content.

The chapter is organized as follows. In Section 7.2, we introduce the linear market with two service providers that will be considered throughout this chapter. In Section 7.3, attention is restricted to the simplest case where symmetric configurations are discussed. We show that, for heavy traffic regimes of the underlying queueing processes, the boundary point partitioning the market interval is governed by a scalar stochastic differential equation with multiplicative noise source. For this dynamics, it is possible to explicitly calculate the associated stationary probability measure. The multiplicative character of the noise source offers the possibility to observe a so-called *noise-induced phase transition*, which manifests itself by a change of the modal character of the stationary probability measure - namely a transition from uni- to a bimodal probability density. In the present context, the transition between these two regimes relates to a transition between a regime where the Hotelling's spatial (*i.e.* the brand) character dominates in the decision taken by the incoming customers and a regime where the time delays dominate. We explicitly work out a simple, though fully representative, illustration belonging to our class of models. For this particular case, we are able to fully calculate the relaxation rate (*i.e.* the transient regimes) characterizing the approach towards the final statistical equilibrium. The relaxation process is strongly dependent on the relative importance of the externalities arising in the associated queueing processes. A short account devoted to simulation experiments explicitly comforts our analytical findings. The dynamics arising for asymmetric configurations is discussed in Section 7.4. Following the technique used for the symmetric case, we compute the stationary probability density function of the boundary point when the two servers work at different service rates (*i.e.* dynamic asymmetry). We also consider the cases where the two service providers charge non-equal prices and the configurations where the two servers do not have symmetric positions with respect to the centre of the market. We show that while these static asymmetries strongly influence the transient regime, they however do not affect the emerging stationary regime. Finally, Section 7.5 is devoted to conclusions.

7.2 Model

As in [24], our starting point will be a two servers Hotelling's model where two service providers S_1 and S_2 are located in a (linear) market confined on a segment $\Omega := [-\Delta, +\Delta] \subset \mathbb{R}$, $\Delta > 0$. The positions of the service providers are denoted respectively by x_1 and x_2 and satisfy $x_1 < x_2$. Let $L = x_2 - x_1$ denotes the distance between S_1 and S_2 . The servers S_1 and S_2 charge respectively prices p_1 and p_2 . Departing now from the original Hotelling's model, we add queueing processes in front of S_1 and S_2 and following [55], we will attach

waiting costs to any customer lining in the queues before being served. Taking into account waiting costs thus confers a dynamical character to the original Hotelling's model. Specifically, our dynamical model exhibits the following features and obeys to the following rules:

- (1) *Arrivals dynamics.* Incoming customers follow a Poisson rule with rate Λ , hence the average time between two arrivals will be $\frac{1}{\Lambda}$.
- (2) *Spatial distribution of the arrivals.* Incoming customers arrive at a random location $x \in \Omega$ drawn from a uniform probability density $U(\Omega)$ with support on Ω .
- (3) *Services dynamics.* Both servers S_i , $i = 1, 2$, have generally distributed service times with rate μ_i , hence the average service time will be $\frac{1}{\mu_i}$, $i = 1, 2$.
- (4) *Traffic intensity.* The traffic into the system is limited to $\rho := \frac{\Lambda}{\mu_1 + \mu_2} < 1$. This ensures that the system is globally stable (*i.e.* the global incoming rate is less than the global service rate).
- (5) *Queueing processes.* When an incoming customer finds servers S_1 or S_2 busy, he/she will wait for service and line-up in a queue. The capacity of the queue is assumed to be unlimited and the service discipline is first-in-first-out (FIFO). In view of points (1) and (2), we hence consider M/G/1 queues.
- (6) *Customer information gathering.* Upon his/her arrival at $x \in \Omega$, each incoming customer knows:
 - (i) his/her relative distance $|x - x_1|$ and $|x - x_2|$ to the servers S_1 and S_2 .
 - (ii) the contents $N_1(t)$ and $N_2(t)$ of both queues ($t \in \mathbb{R}^+$ being the arrival time). In other words, both queue contents are observable to any incoming customer.
- (7) *Cost structures.* There are two types of costs incurred by any customer, namely:
 - (i) the waiting time cost (WTC), characterized by a cost parameter c_w with physical unit $\left[\frac{\text{dollar}}{\text{time unit}}\right]$.
 - (ii) the brand departure cost¹ (BDC), quantified by a cost parameter c_t with physical unit $\left[\frac{\text{dollar}}{\text{brand distance unit}}\right]$.

¹ The modelling framework proposed in this chapter would also perfectly fit in the case of a duopoly with customers sensitive to transportation costs (in place of the brand considerations that are considered here). In that situation, the brand departure cost would simply be replaced by an analog transportation cost. Indeed, the present modelling framework has directly been inspired by the Hotelling's duopoly model, where two firms are located within a linear market and where relative geographical distances to these vendors influence the customers' decisions. Extending this situation to brand aspects, the idea here is, when choosing a service provider, to quantify with a distance any departure from the ideal brand expected by a customer. Using this analogy, the Hotelling's spatial duopoly model can thus be used to characterize a market with two differently branded firms.

(8) *Decision policy.* Upon arrival, an incoming customer is aware of:

- the queue contents $N_1(t)$ and $N_2(t)$,
- his/her relative position to S_1 and S_2 ,
- the values of the costs c_w and c_t ,
- the service rates μ_1 and μ_2 ,
- the posted prices p_1 and p_2 .

Based on this information set, the incoming customer autonomously decides which server S_1 or S_2 he/she will join.

(9) *Demand structure.* Following the original Hotelling's case, we assume an inelastic demand, *i.e.* a customer will purchase the service at any price, even if the proposed price is arbitrarily large.

A sketch of our modelling framework can be found in Fig. 7.1. Extending

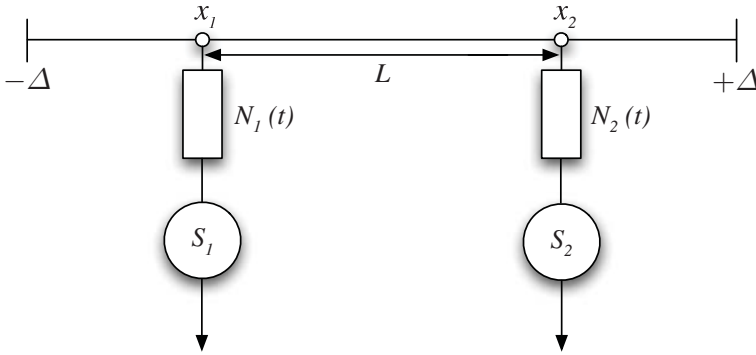


Fig. 7.1. Bounded market with two vendors and time sensitive customers.

original Hotelling's original configuration, the consideration of waiting times confers to the present class of multi-agent queueing systems an explicit dynamical character. While in the agent-based models described in the preceding chapters the items in circulation were basing their autonomous routing choices mostly on perceived waiting times (and hence on history-based, a posteriori, data), here the customers take into account the expected waiting times (and thus a priori observations) to construct their individual decisions.

When served by S_i , an incoming customer has a utility function $U_i(x)$, $i = 1, 2$, where x is the customer's initial position which enters into the decision policy. In words, the functions $U_i(x)$ quantify the gain realized by a customer choosing server S_i when he/she enters into the system at location x . Specifically, for linear waiting and brand departure (transportation) costs, the utility functions read as:

$$U_i(x) = a - p_i - c_t|x - x_i| - c_w\mathbb{E}(W_i|N_i(t)), \quad i = 1, 2, \quad (7.1)$$

with a being a systematic reward due to the service and $\mathbb{E}(W_i|N_i(t))$ standing for the conditional expected waiting time at S_i when $N_i(t)$ already waiting customers are observed. As $\frac{1}{\mu_i}$ is the average service time at server S_i , this last conditional expectation is readily given by:

$$\mathbb{E}(W_i|N_i(t)) = \frac{N_i(t)}{\mu_i}.$$

We obviously assume that any customer maximizes his/her utility function when choosing one of the two servers. This suggests to introduce a time-dependent *market partition boundary point* $Y_t \in [-\Delta, +\Delta]$ implicitly defined by:

$$U_1(Y_t) = U_2(Y_t). \quad (7.2)$$

Hence, our strictly increasing (linear) BDC which we assume from now on imply that Y_t dynamically separates the two monopolies held by S_1 and S_2 . A sketch of the situation is given in Fig. 7.2. As Y_t is a function of the two stochastic processes $N_1(t)$ and $N_2(t)$, it will be itself a stochastic process.

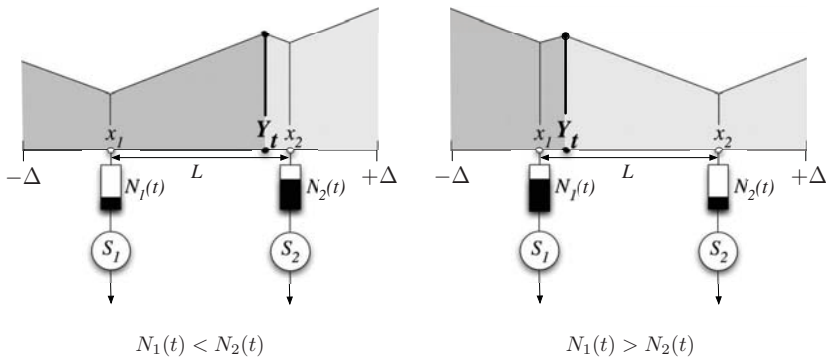


Fig. 7.2. Cost structure as a function of the customers' location. The total costs for a customer located at position x are the sum of the selling price p_i , the waiting time cost $c_w\mathbb{E}(W_i|N_i(t))$ (identical service rates are assumed in this figure) and the brand departure cost $c_t|x - x_i|$. Any customer chooses the service provider minimizing his/her total costs (*i.e.* it corresponds to maximize his/her utility function). As a consequence, all the customers standing on the left of Y_t will choose S_1 , those on the right go to S_2 . The difference between the two figures is the current queue contents. These contents determine the position of the market partition boundary point Y_t , which separates the respective market shares held by S_1 and S_2 .

Let $\lambda_i(t, Y_t)$ denote the partial incoming rate of customers feeding S_i at time t and hence:

$$\lambda_1(t, Y_t) + \lambda_2(t, Y_t) = \Lambda, \quad \forall t \in \mathbb{R}^+. \quad (7.3)$$

In view of the assumption (2) (*i.e.* spatially uniform arrivals on $\Omega = [-\Delta, +\Delta]$), the partial traffic flows feeding respectively S_1 and S_2 result from the Bernoulli “thinning” of the incoming Poisson flow with global rate Λ . The branching probability is given by $P_{\text{branch}} = \frac{\Delta - Y_t}{2\Delta}$ and it is established (see [28] for instance) that such a thinning produces two independent Poisson processes with partial rates:

$$\lambda_1(t, Y_t) = \frac{\Delta + Y_t}{2\Delta} \Lambda \quad \text{and} \quad \lambda_2(t, Y_t) = \frac{\Delta - Y_t}{2\Delta} \Lambda. \quad (7.4)$$

For the utility functions given by Eq. (7.1), the time-dependent market partition boundary point will obey, $\forall t \in \mathbb{R}^+$:

$$Y_t = \begin{cases} (a) \frac{c_w}{2c_t} \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right) + \frac{x_1 + x_2}{2} + \frac{p_2 - p_1}{2c_t} & \text{if } c_t L \geq |\Psi|, \\ (b) + \Delta & \text{if } c_t L < \Psi, \\ (c) - \Delta & \text{if } c_t L < -\Psi, \end{cases} \quad (7.5)$$

where:

$$\Psi := \Psi(N_1(t), N_2(t), \mu_1, \mu_2, p_1, p_2) = p_2 - p_1 + c_w \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right).$$

In case (a), $Y_t \in [x_1, x_2] \subset [-\Delta, +\Delta]$. Indeed in this case, the BDC from one server to the other (*i.e.* $c_t L$) is greater than the global difference between the prices and the WTCs of the two servers (*i.e.* $|\Psi|$). Hence, a customer located near the server having the longest queue will choose this server anyway. In cases (b) and (c), any customer in the whole interval $[-\Delta, +\Delta]$ joins the server having the shortest queue. Indeed, the gain in WTC (due to the difference between the queue contents) and in price exceeds the BDC incurred by the distance from one server to the other. A representation of the dynamics induced by Eq. (7.5) for a particular choice of the control parameters is found in Fig. 7.3.

We now separately discuss fully symmetric configurations (Section 7.3) and various types of asymmetric situations (Section 7.4).

7.3 Symmetric Configurations

The positions of the service providers are assumed to satisfy $-\Delta \leq x_1 < 0$ and $0 < x_2 \leq +\Delta$ and they are located symmetrically with respect to the center of the market, *i.e.* $x_1 = -x_2$. Furthermore, the servers S_1 and S_2 offer homogeneous services $\mu_1 = \mu_2 = \mu$ and both charge an equal price $p_1 = p_2 = p$.

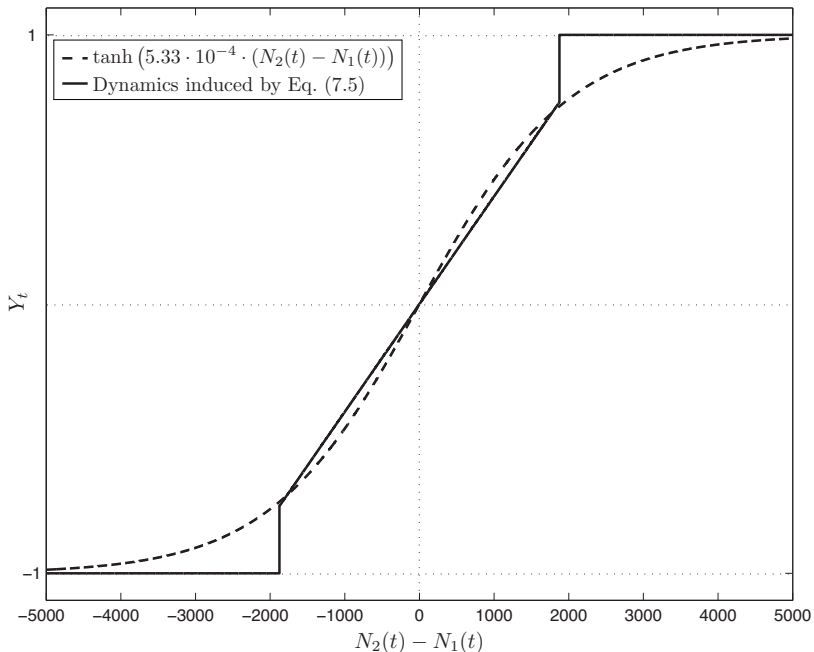


Fig. 7.3. Typical representation of the boundary point dynamics Y_t for $\Delta = 1$. The solid line shows the dynamics induced by Eq. (7.5) when $c_t = 10$, $c_w = 8 \times 10^{-3}$, $\mu_1 = \mu_2 = 1$, $p_1 = p_2 = 1$ and $x_2 = -x_1 = \frac{3}{4}$. The dashed line shows the approximate dynamics given by Eq. (7.21) when $\gamma = 5.33 \times 10^{-4}$.

Let $A_i(t)$, $D_i(t)$ and $N_i(t)$ respectively denote the numbers of arrivals, departures and the population in S_i at time t . From now on, we restrict ourselves to *heavy traffic* regimes characterized by $\rho = \frac{\lambda}{2\mu} = 1 - \epsilon$, with ϵ small. Writing

$$N_i(t) = A_i(t) - D_i(t),$$

in heavy traffic the server S_i has very long busy periods and hence the process $N_i(t)$ does almost never vanish, $i = 1, 2$. This implies that the departure and arrival processes are almost independent. In heavy traffic regimes, it is well established (see in particular [93]) that both queue contents at time t are well approximated by diffusion processes of the form:

$$N_i(t) = \int_0^t [\lambda_i(s, Y_s) - \mu] ds + \int_0^t V_i(s, Y_s) dB_{i,s} \quad i = 1, 2, \quad (7.6)$$

where $B_{1,t}$ and $B_{2,t}$ are independent standard Brownian motions and the terms $V_i(t, Y_t)$ denote the state-dependent “volatilities” (*i.e.* the standard

deviations) given by:

$$V_i(t, Y_t)^2 = \lambda_i(t, Y_t)^3 \sigma_{a,i}^2 + \mu^3 \sigma_{s,i}^2 \quad i = 1, 2, \quad (7.7)$$

with $\sigma_{a,i}^2$ (resp. $\sigma_{s,i}^2$) being the variance of the inter-arrival times (resp. the variance of the service times) for server S_i . Using Eqs. (7.3) to (7.7) and the fact that $B_{1,t}$ and $B_{2,t}$ are independent, we therefore can write:

$$N_2(t) - N_1(t) = -\frac{\Lambda}{\Delta} \int_0^t Y_s ds + \int_0^t V(s, Y_s) dB_s, \quad (7.8)$$

with B_t being a standard Brownian motion and $V^2(t, Y_t) = V_1(t, Y_t)^2 + V_2(t, Y_t)^2 = \Lambda + \mu^3 (\sigma_{s,1}^2 + \sigma_{s,2}^2)$ - remember that for Poisson processes, we have $\sigma_{a,i}^2 = \lambda_i(t, Y_t)^{-2}$.

In this symmetric configuration, Eq. (7.5) reduces to, $\forall t \in \mathbb{R}^+$:

$$Y_t = \begin{cases} \frac{c_w}{2\mu c_t} (N_2(t) - N_1(t)) & \text{if } c_t L \geq |\tilde{\Psi}|, \\ +\Delta & \text{if } c_t L < \tilde{\Psi}, \\ -\Delta & \text{if } c_t L < -\tilde{\Psi}, \end{cases} \quad (7.9)$$

where, for this symmetric case, we define:

$$\tilde{\Psi} = \Psi(N_1(t), N_2(t), \mu, \mu, p, p) = \frac{c_w}{\mu} (N_2(t) - N_1(t)).$$

To proceed further with analytical calculations, we approximate the dynamics implied by Eq. (7.9) by introducing an odd (due to the symmetry of the problem), effective monotonously increasing one-to-one, $C^2(\mathbb{R})$ function²

$$f(\cdot) : \mathbb{R} \rightarrow [-1, +1] \quad (7.10)$$

fulfilling

$$Y_t = \Delta f(\gamma(N_2(t) - N_1(t))), \quad (7.11)$$

with:

$$\gamma := \frac{c_w}{\mu L c_t}. \quad (7.12)$$

The control parameter γ is dimensionless and quantifies the respective importance of the different costs. Note that in Eq. (7.12), the time unit is measured on average service time.

² In the following, we will be interested in describing a noise-induced phase transition that occurs between regimes where the market sharing dynamics is respectively characterized by a unimodal and by a bimodal stationary probability density function. For the approximation considered in Eq. (7.11), a formal characterization of the class of functions that effectively leads to such noise-induced phase transition is given in Appendix 13.

As f is invertible, Eq. (7.11) can be written as:

$$f^{-1} \left(\frac{Y_t}{\Delta} \right) = \gamma (N_2(t) - N_1(t)). \quad (7.13)$$

Using Eq. (7.8), Eq. (7.13) becomes:

$$f^{-1} \left(\frac{Y_t}{\Delta} \right) = -\frac{\gamma\Lambda}{\Delta} \int_0^t Y_s ds + \gamma \int_0^t V(s, Y_s) dB_s. \quad (7.14)$$

Differentiating, we obtain:

$$(f^{-1})' \left(\frac{Y_t}{\Delta} \right) dY_t = -\gamma\Lambda Y_t dt + \Delta\gamma V(t, Y_t) dB_t, \quad (7.15)$$

which can be written as:

$$dY_t = -\frac{\gamma\Lambda Y_t}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} dt + \frac{\Delta\gamma V(t, Y_t)}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} dB_t. \quad (7.16)$$

In our settings (remember that we deal with M/G/1 queues), $V(t, Y_t) = V = \sqrt{\Lambda + \mu^3 (\sigma_{s,1}^2 + \sigma_{s,2}^2)}$ does not depend on Y_t nor on t . We can thus write Eq. (7.16) as:

$$dY_t = -\frac{\gamma\Lambda Y_t}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} dt + \frac{\Delta\gamma V}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} dB_t. \quad (7.17)$$

The stochastic differential equation (SDE) given by Eq. (7.17) describes the effective dynamics of the boundary position Y_t . The White Gaussian noise dB_t being merely the limit of finitely correlated processes, we assign to the underlying stochastic integral relative to Eq. (7.17) the Sratonovitch's interpretation, [61]. Hence, the transition probability density $P(y, t \mid y_0, t_0)$ describing the solution of the SDE (7.17) reads as:

$$\frac{\partial}{\partial t} P(y, t \mid y_0, t_0) = \mathcal{F}P(y, t \mid y_0, t_0), \quad (7.18)$$

with Fokker-Planck operator taking here the form, [61]:

$$\mathcal{F}(\cdot) := \frac{\partial}{\partial y} \left[\frac{\gamma\Lambda y}{(f^{-1})' \left(\frac{y}{\Delta} \right)} (\cdot) \right] + \frac{1}{2} \frac{\partial}{\partial y} \left[g(y) \frac{\partial}{\partial y} g(y) (\cdot) \right], \quad g(y) = \frac{\Delta\gamma V}{(f^{-1})' \left(\frac{y}{\Delta} \right)}.$$

The stationary probability density function $P_s(y)$ solving Eq. (7.18), with vanishing left hand side, reads as:

$$P_s(y) = \mathcal{N} (f^{-1})' \left(\frac{y}{\Delta} \right) \exp \left\{ -\frac{2\Lambda}{\gamma \Delta^2 V^2} \int^y u (f^{-1})' \left(\frac{u}{\Delta} \right) du \right\}, \quad (7.19)$$

for $y \in [-\Delta, +\Delta]$, with $\mathcal{N} < \infty$ a normalization constant.

Symmetry (*i.e.* our assumptions that $x_1 = -x_2$, $\mu_1 = \mu_2$ and $p_1 = p_2$) implies that $P_s(y) = P_s(-y)$. In particular, studying the curvature $\mathcal{R}(0)$ of $P_s(y)$ at $y = 0$ directly provides information regarding the modularity of $P_s(y)$. From Eq. (7.19), we directly obtain:

$$\text{sign} \{ \mathcal{R}(0) \} = \text{sign} \left\{ -\gamma V^2 f'''(0) - 2\Lambda (f^{-1})'(0) (f'(0))^3 \right\}. \quad (7.20)$$

For given functions f , we observe that the sign of the curvature $\mathcal{R}(0)$ directly depends on the values of the (control) external parameters (here c_w , c_t , L , Λ and μ) solely. A curvature sign change exhibits a transition of regime triggered by the presence of fluctuations. This is referred as a *noise-induced phase transition*, [61], and an explicit illustration is now worked out.

7.3.1 Explicit Illustration - Symmetric Case

Belonging to the previous class of models, the particular choice

$$Y_t = \Delta \tanh(\gamma(N_2(t) - N_1(t))) \quad (7.21)$$

leads to very simple algebra. A particular representation of Eq. (7.21), in comparison with the dynamics induced by Eq. (7.5), is found in Fig. 7.3.

For this particular case, the SDE (7.17), describing the effective boundary point dynamics, becomes:

$$dY_t = -\gamma \Lambda Y_t \left(1 - \left(\frac{Y_t}{\Delta} \right)^2 \right) dt + \Delta \gamma V \left(1 - \left(\frac{Y_t}{\Delta} \right)^2 \right) dB_t. \quad (7.22)$$

In view of Eq. (7.19), the corresponding stationary probability density function simply becomes:

$$P_s(y) = \mathcal{N} \left(1 - \left(\frac{y}{\Delta} \right)^2 \right)^{\frac{\Lambda}{\gamma V^2} - 1} \quad \text{for } y \in [-\Delta, +\Delta], \quad (7.23)$$

where \mathcal{N} is the normalization constant given here by:

$$\mathcal{N}^{-1} = \Delta \int_0^1 t^{-\frac{1}{2}} (1-t)^{\frac{\Lambda}{\gamma V^2} - 1} dt = \Delta \Upsilon \left(\frac{1}{2}, \frac{\Lambda}{\gamma V^2} \right),$$

where $\Upsilon(x, y) := \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$ and $\Gamma(x)$ stands for the Gamma function. An illustration of the probability density function given by Eq. (7.23) for different

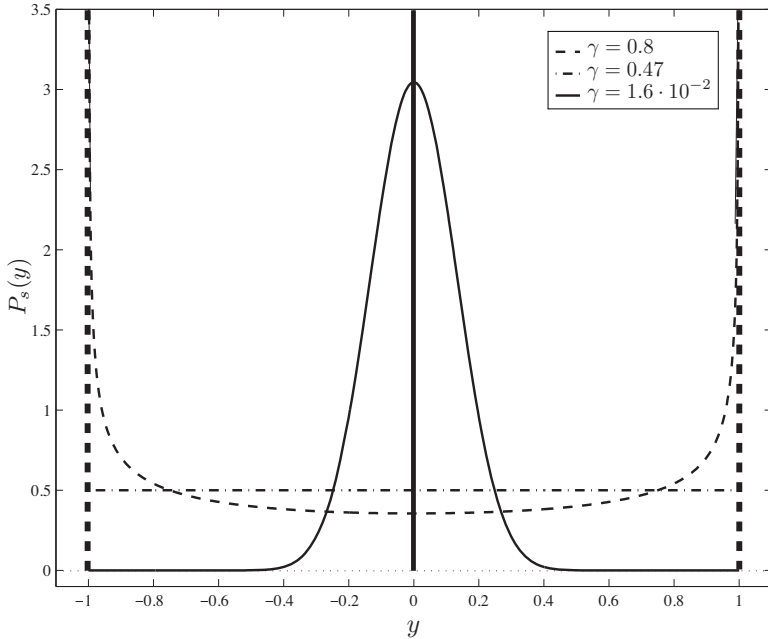


Fig. 7.4. Stationary probability density function of the time-dependent boundary position Y_t when $f(x) = \tanh(x)$, $\Delta = 1$, $\Lambda = 1.8$, $\mu = 1$ ($\rho = 0.9$) and the service time processes are Poisson. This density is drawn for three different values of $\gamma = [0.8; 0.47; 1.6 \times 10^{-2}]$. Furthermore, when $\gamma \rightarrow \infty$ (it corresponds to purely deadline type regimes, which appear when the customers solely focus on waiting time issues), the density is sharply peaked at $y = -\Delta = -1$ and $y = +\Delta = +1$. In the other limit, $\gamma \rightarrow 0$ (corresponding to purely Hotelling-like regimes, which emerge when the customers' choices are wholly driven by brand considerations), the density is restricted to a single peak at $y = 0$. This graph clearly exhibits the noise-induced phase transition arising in our dynamical model.

values of γ and $\Delta = 1$ is found in Fig. 7.4. Regarding Eq. (7.20), the sign of the curvature $\mathcal{R}(0)$ of $P_s(y)$ at $y = 0$ is here given by:

$$\mathcal{R}(0) \begin{cases} > 0 & \text{when } \frac{\Lambda}{\gamma V^2} < 1, \\ = 0 & \text{when } \frac{\Lambda}{\gamma V^2} = 1, \\ < 0 & \text{when } \frac{\Lambda}{\gamma V^2} > 1. \end{cases} \quad (7.24)$$

The information given by Eq. (7.24) (which is in perfect agreement with what we would expect with regard to the form of $P_s(y)$ given by Eq. (7.23)) perfectly describes the modularity of $P_s(y)$ and the underlying noise-induced phase transition.

Transient Behaviour

For the choice given in Eq. (7.21), we can also study the rate of approach to the equilibrium state. Indeed, by introducing the change of variables:

$$t \mapsto \tau = \gamma^2 V^2 t, \quad X_t \mapsto Y_t = \Delta \tanh(X_t), \quad (7.25)$$

the dynamics given by Eq. (7.22) reduces to:

$$\begin{aligned} dX_\tau &= - \left(\frac{\Lambda}{\gamma V^2} + 1 \right) \tanh(X_\tau) d\tau + dW_\tau \\ &:= -2K \tanh(X_\tau) d\tau + dW_\tau \end{aligned} \quad (7.26)$$

and the time-dependent solution $P(y, t \mid y_0, 0)$ of the associated Fokker-Planck operator is known for long (see for instance [125]). As an illustration, let us mention that for the situations where the dimensionless parameter $K := \frac{\Lambda}{2\gamma V^2} + \frac{1}{2} \in \mathbb{N}$, the explicit form simplifies somewhat and is given by, [125]:

$$\begin{aligned} P(y, t \mid y_0, 0) &= \frac{1}{(1+z^2)^{K+\frac{1}{2}}} \left[(1+z_0^2)^{\frac{K}{2}} (1+z^2)^{\frac{K}{2}} \frac{1}{2\sqrt{\pi t}} e^{-K^2 t} e^{-g(z, z_0, t)^2} \right] + \\ &\frac{1}{\pi(1+z^2)^{K+\frac{1}{2}}} \sum_{n=0}^{K-1} \frac{(K-n)}{n! \Gamma(2K+1-n)} e^{-n(2K-n)t} \theta_n(z_0) \theta_n(z) f_n(z, z_0, t), \end{aligned} \quad (7.27)$$

with the definitions:

$$\begin{aligned} \sinh(z) &:= y, \quad f_n(x, x_0, t) := \frac{1}{\sqrt{\pi}} \int_{g(x, x_0, t) - (K-n)\sqrt{t}}^{g(x, x_0, t) + (K-n)\sqrt{t}} e^{-u^2} du, \\ g(x, x_0, t) &:= \frac{\operatorname{arcsinh}(x) - \operatorname{arcsinh}(x_0)}{2\sqrt{t}} \end{aligned}$$

and the polynomials:

$$\theta_n(x) := (-1)^n 2^{K-n} \Gamma(K-n + \frac{1}{2}) (1+x^2)^{K+\frac{1}{2}} \frac{d^n}{dx^n} (1+x^2)^{n-K-\frac{1}{2}}.$$

In particular, the long time scale t_{relax} governing the approach to the stationary state given by Eq. (7.23) is determined by the spectral gap between 0 and the first non vanishing eigenvalue of the Fokker-Planck equation (7.18) (remember that the vanishing eigenvalue corresponds to the stationary probability measure given by Eq. (7.19)). It follows that:

$$1/t_{relax} = \begin{cases} (2K-1)\gamma^2 V^2 = \Lambda\gamma & \text{if } K \geq 1, \\ K^2 \gamma^2 V^2 = \frac{\gamma\Lambda + \gamma^2 V^2}{2} & \text{if } K < 1. \end{cases} \quad (7.28)$$

From Eq. (7.28), we can draw the following remarks:

(1) *Spectral characteristics of the Fokker-Planck equation.*

In view of Eq. (7.28), there are two relaxation regimes governed by the spectral properties of the associated Fokker-Planck equation (7.18). As discussed in [125], for $K \geq 1$ the spectrum exhibits both discrete and continuum parts whereas for $K < 1$ only the continuum part survives.

(2) *Regime transitions.*

Note that the transition from unimodal to bimodal densities given in Eq. (7.23) by $\left(\frac{A}{\gamma V^2} - 1\right) = 0$ coincides with the transition in the relaxation regimes given by Eq. (7.28)

(3) *Rate of approach to the equilibrium.*

When discrete eigenvalues exist, the asymptotic relaxation time towards the single mode stationary probability density function (given by Eq. (7.23)) is faster compared to the relaxation rate associated with the purely continuum spectrum which drives the system to the bimodal stationary density. This can be intuitively understood in limiting regimes. Indeed, for the pure Hotelling's case, a situation arising when $c_t \rightarrow \infty$, the boundary position probability density is delta-peaked in the middle of the market interval (remember that we have focused in this section on symmetric configurations) and the relaxation time to reach this equilibrium is vanishingly small - this corresponds to the deterministic scheduling rule which commands to "*join the closest server*". For dominating Hotelling's type regimes (*i.e.* when brand aspects are predominant), the externalities (*i.e.* the waiting costs affecting incomers arriving behind a customer entering into service) have little influence on the equilibrium probability density which describes the boundary point - this produces a fast relaxation towards the statistical equilibrium, which will be close to the limiting delta-peaked density. In the contrary, when the deadline type regime dominates (*i.e.* when waiting time considerations are predominant), a new incomer strongly modifies the dynamical state of the system and hence strongly perturbs the underlying probability measure, thus implying long relaxation times to the statistical equilibrium. Note that for $K = 0$ in Eq. (7.28), a situation realized when $c_w \rightarrow \infty$, the relaxation time diverges to infinity, meaning that no statistical equilibrium exists - this corresponds to the purely deterministic scheduling policy which commands to always "*join the server exhibiting the shortest queue*".

Simulation Experiments

We have simulated the dynamics of the market position boundary position Y_t in the particular case where Y_t fulfils Eq. (7.21). Each customer, upon arrival, determines on which side of the boundary point Y_t (dynamically given by Eq. (7.21), with regard to the current content of the queues) is his/her

(uniformly distributed) position and he/she joins the queue hence chosen. We have computed an estimation of the stationary probability density function of the boundary position Y_t after 10^5 customers have passed through the system. The simulation experiments performed for different values of the control parameters (here γ , Λ and μ) confirm the presence of the noise-induced phase transition. Fig. 7.5 illustrates how the simulation results fit perfectly well with the theoretical findings given by Eq. (7.23).

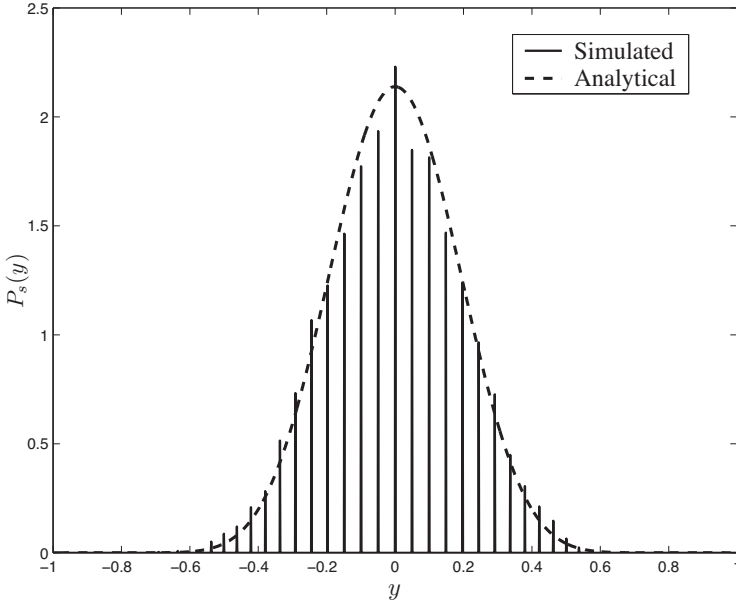


Fig. 7.5. Simulated and theoretical stationary probability density function of the time-dependent boundary position $Y_t = \Delta \cdot \tanh(\gamma(N_2(t) - N_1(t)))$ when $\Delta = 1$, $\Lambda = 1.9$, $\mu = 1$ ($\rho = 0.95$), $\gamma = 5 \cdot 10^{-2}$ and the service time processes are Poisson.

7.4 Asymmetric Configurations

Different sources of asymmetry are possible:

(1) *Dynamic Asymmetry.*

This situation is encountered when heterogeneous servers are operating, it is treated in Section 7.4.1.

(2) *Static Asymmetry.*

This arises in presence of non-symmetric server locations and/or non-equal prices. Configurations where the servers have asymmetric positions with respect to the centre of the market and situations where the posted prices are different lead to dynamics that are analogous to the symmetric case and they are discussed in Section 7.4.2.

7.4.1 Heterogeneous Servers

We treat in this section the situations where the two service providers work at different service rates $\mu_1 \neq \mu_2$ but post the same prices $p_1 = p_2 = p$ and have symmetric locations $x_2 = -x_1$. For utility functions satisfying Eq. (7.1), the dynamic boundary point Y_t here obeys, $\forall t \in \mathbb{R}^+$:

$$Y_t = \begin{cases} \frac{c_w}{2c_t} \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right) & \text{if } c_t L \geq c_w \left| \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right) \right|, \\ +\Delta & \text{if } c_t L < c_w \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right), \\ -\Delta & \text{if } c_t L < c_w \left(\frac{N_1(t)}{\mu_1} - \frac{N_2(t)}{\mu_2} \right). \end{cases} \quad (7.29)$$

In view of Eq. (7.29), the dynamics is here driven by the difference between the normalized numbers of customers $\frac{N_i(t)}{\mu_i}$, $i = 1, 2$, waiting in the different queues. While we have used the approximation given by Eq. (7.11) in the symmetric case, we approximate here the dynamics implied by Eq. (7.29) with:

$$Y_t = \Delta f \left(\tilde{\gamma} \left(\frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} \right) \right), \quad (7.30)$$

where f is a function satisfying the same hypothesis as in the symmetric case and:

$$\tilde{\gamma} := \frac{c_w}{L c_t}. \quad (7.31)$$

Note that $\tilde{\gamma} = \mu \gamma$ (with γ being given by Eq. (7.12)). Following the same methodology used to derive Eqs. (7.6) to (7.8), it ensues that:

$$\begin{aligned} \frac{N_2(t)}{\mu_2} - \frac{N_1(t)}{\mu_1} &= \frac{\Lambda}{2} \int_0^t \left[\left(\frac{1}{\mu_2} - \frac{1}{\mu_1} \right) - \frac{1}{\Delta} \left(\frac{1}{\mu_2} + \frac{1}{\mu_1} \right) Y_s \right] ds \\ &\quad + \int_0^t \tilde{V}(s, Y_s) dB_s, \end{aligned} \quad (7.32)$$

where, for Poisson arrival and service processes:

$$\begin{aligned} \tilde{V}(t, Y_t)^2 &= \frac{V_1(t, Y_t)^2}{\mu_1^2} + \frac{V_2(t, Y_t)^2}{\mu_2^2} \\ &= \frac{\Lambda Y_t}{2\Delta} \left(\frac{1}{\mu_1^2} - \frac{1}{\mu_2^2} \right) + \frac{\Lambda}{2} \left(\frac{1}{\mu_1^2} + \frac{1}{\mu_2^2} \right) + \left(\frac{1}{\mu_1} + \frac{1}{\mu_2} \right). \end{aligned}$$

Starting from Eq. (7.30) and following the lines used to derive Eq. (7.11) to Eq. (7.17), we obtain:

$$dY_t = \left[\frac{\frac{\tilde{\gamma}\Lambda\Delta}{2} \left(\frac{1}{\mu_2} - \frac{1}{\mu_1} \right) - \frac{\tilde{\gamma}\Lambda Y_t}{2} \left(\frac{1}{\mu_1} + \frac{1}{\mu_2} \right)}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} \right] dt + \frac{\Delta \tilde{\gamma} \tilde{V}}{(f^{-1})' \left(\frac{Y_t}{\Delta} \right)} dB_t. \quad (7.33)$$

Setting $\mu_1 = \mu_2 = \mu$ in Eq. (7.33), we directly recover the dynamics valid in the symmetric case, given by Eq. (7.17). The stationary probability density function ensuing from the dynamics stated in Eq. (7.33) is given by:

$$P_s(\mu_1, \mu_2; y) = \mathcal{N} \frac{(f^{-1})' \left(\frac{y}{\Delta} \right)}{\tilde{V}(t, y)} \times \exp \left\{ \frac{\Lambda}{\tilde{\gamma}\Delta^2} \int^y \left[\frac{\Delta \left(\frac{1}{\mu_1} - \frac{1}{\mu_2} \right) - u \left(\frac{1}{\mu_1} + \frac{1}{\mu_2} \right)}{\tilde{V}(t, u)^2} \right] (f^{-1})' \left(\frac{u}{\Delta} \right) du \right\}, \quad (7.34)$$

where \mathcal{N} is a normalization constant. Note that the structure of the dynamics obviously implies that $P_s(\mu_1, \mu_2; y) = P_s(\mu_2, \mu_1; -y)$. Remark that we consistently recover Eq. (7.19) when we fix $\mu_1 = \mu_2 = \mu$ in Eq. (7.34).

Explicit Illustration - Asymmetric Case

For the particular choice $f(x) = \tanh(x)$, we find:

$$P_s(\mu_1, \mu_2; y) = \mathcal{N} \left(1 + \frac{y}{\Delta} \right)^{-1 - \frac{\beta - \alpha}{\xi - \delta}} \left(1 - \frac{y}{\Delta} \right)^{-1 - \frac{\beta + \alpha}{\xi + \delta}} \times (\delta\Delta + \xi y)^{\frac{2(\beta\xi - \alpha\delta)}{\xi^2 - \delta^2} - \frac{1}{2}}, \quad (7.35)$$

where:

$$\alpha = -\frac{\tilde{\gamma}\Lambda\Delta}{2} \left(\frac{1}{\mu_1} + \frac{1}{\mu_2} \right), \quad \beta = \frac{\tilde{\gamma}\Lambda\Delta}{2} \left(\frac{1}{\mu_1} - \frac{1}{\mu_2} \right), \quad \xi = \frac{\tilde{\gamma}^2\Lambda\Delta}{2} \left(\frac{1}{\mu_1^2} - \frac{1}{\mu_2^2} \right)$$

and $\delta = \frac{\tilde{\gamma}^2\Lambda\Delta}{2} \left(\frac{1}{\mu_1^2} + \frac{1}{\mu_2^2} \right) + \tilde{\gamma}^2\Delta \left(\frac{1}{\mu_1} + \frac{1}{\mu_2} \right).$

A sketch of the stationary distributions arising for heterogeneous services is given in Fig. 7.6. We observe that P_s is shifted (and biased) to the opposite side of the most effective server. This server obviously attracts more customers than the slowest one. This illustrates the fact that the most effective server does enhance its market share.

Finally, observe that for $\mu_1 = \mu_2 = \mu$ in Eq. (7.35), we consistently recover Eq. (7.23).

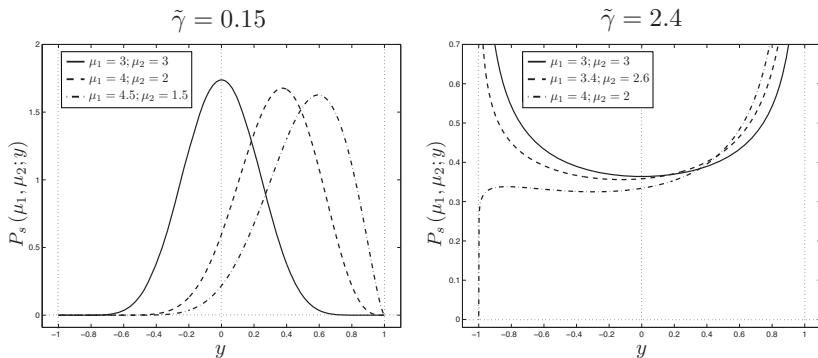


Fig. 7.6. Stationary probability density function of the time-dependent boundary position Y_t for heterogeneous service rates $\mu_1 \neq \mu_2$ and $f(x) = \tanh(x)$. Here, $\Delta = 1$ and $\Lambda = 5.7$. *Left:* $\tilde{\gamma} = 0.15$ (Hotelling-like regime). *Right:* $\tilde{\gamma} = 2.4$ (deadline type regime), note that the mixed boundary behaviour of the dash-dot curve can only arise for asymmetric configurations.

7.4.2 Asymmetric Positions and Different Prices

When the two service providers are not symmetrically located with respect to the centre of the market (*i.e.* $x_1 \neq -x_2$), but have equal service rates and prices, the utility functions felt by the customers are modified such that the boundary position Y_t obeys:

$$Y_t = \begin{cases} \frac{c_w}{2\mu c_t} (N_2(t) - N_1(t)) + \frac{x_1 + x_2}{2} & \text{if } c_t L \geq |\tilde{\Psi}|, \\ +\Delta & \text{if } c_t L < \tilde{\Psi}, \\ -\Delta & \text{if } c_t L < -\tilde{\Psi}. \end{cases} \quad (7.36)$$

Similarly, when the service providers differ only in their posted price (*i.e.* $p_1 \neq p_2$), the time-dependent boundary position obeys:

$$Y_t = \begin{cases} \frac{c_w}{2\mu c_t} (N_2(t) - N_1(t)) + \frac{p_2 - p_1}{2c_t} & \text{if } c_t L \geq |\hat{\Psi}|, \\ +\Delta & \text{if } c_t L < \hat{\Psi}, \\ -\Delta & \text{if } c_t L < -\hat{\Psi}, \end{cases} \quad (7.37)$$

where:

$$\hat{\Psi} = \Psi(N_1(t), N_2(t), \mu, \mu, p_1, p_2) = p_2 - p_1 + \frac{c_w}{\mu} (N_2(t) - N_1(t)).$$

For both cases (or a combination of them), the addition of a static asymmetry contribution is required and the dynamics can be now approximated by:

$$Y_t = \Delta f(\gamma(N_2(t) - N_1(t) + \eta)), \quad (7.38)$$

where $\eta \in \mathbb{R}$ is a new parameter that quantifies the asymmetry. Using the same methodology as for the symmetric case, we get:

$$f^{-1}\left(\frac{Y_t}{\Delta}\right) = -\frac{\gamma\Lambda}{\Delta} \int_0^t Y_s ds + \gamma \int_0^t V(s, Y_s) dB_s + \gamma\eta. \quad (7.39)$$

When taking the time-derivative of Eq. (7.39), the asymmetric contribution (*i.e.* $\gamma\eta$) disappears and we get back to the symmetric dynamics given by Eq. (7.15). Hence for static asymmetry, the stationary probability density coincides with the symmetric case, given by Eq. (7.19). This should in fact not come as a surprise. Indeed, the static asymmetry manifests itself only during the transient regime. Starting with empty queues for both servers, the boundary point is initially located closer to the less attractive (in terms of price and/or position) server, implying thus a larger feeding rate to the most attractive one. The (static) lack of attractivity of one server will gradually be counterbalanced by a (dynamic) larger number of customers visiting the most attractive server. Asymptotically, the stationary regime for Y_t behaves as in the symmetric case. Note however that, contrary to the symmetric case, the stationary queue contents will however not be equal anymore.

7.5 Concluding Remarks

When alternative choices between services are offered to customers, several criteria enter into their final selection decision. There are namely static criteria such as posted prices and server locations as considered in the original Hotelling's model but there also exist dynamic aspects typified by the waiting times before service. While we have mainly considered perceived waiting times as a decision criterion in the preceding chapters (*i.e.* a posteriori, history-based observations), here the customers autonomously choose their routing with regard to expected waiting times (*i.e.* a priori measures). As already emphasized in this thesis, it is intuitively clear that the negative aspects of the expected, actual and/or perceived waiting times strongly affect the final customers' satisfaction and hence their individual decisions. Furthermore, as in generic situations the waiting time is an intrinsically random quantity, it is thus naturally described in the context of queueing theory. Focusing on a simple duopoly configuration, we have studied here the (stochastic) dynamics of the frontier which defines the partition between the market shares held by the two service providers. For heavy traffic regimes of the underlying queueing processes, the market partition boundary point can be described by a (random) diffusion dynamics (*i.e.* a differential equation driven by a White Gaussian noise) with a multiplicative noise source (*i.e.* a state-dependent diffusion term). It is remarkable that the stationary probability measure characterizing the market partition boundary position dynamics exhibits a noise-induced phase transition triggered by the values of the external control parameters

(brand departure cost, waiting time cost, service rate and spatial separation between the servers). Note that multiplicative noise processes are not confined to physics, chemistry and biology domains where they first have been applied - they also naturally occur in operational research, in economics and more generally in social sciences. One of the most popular illustration is clearly found in financial mathematics - the Black-Scholes model, which is based on the geometric Brownian motion (hence a multiplicative noise process). Note however that contrary to the market sharing dynamics considered here, no noise-induced phase transition occurs in this financial context.

7.6 Contributions of Chapter 7

- We introduce an explicit dynamical character into the famous Hotelling's duopoly model. More precisely, the dynamical nature of our model reflects that customers do not choose their service provider based only on spatial (brand) aspects, but also in function of the expected waiting times at each server.
- Considering the queueing processes that are created at each server, we formulate and solve the stochastic differential equation that drives the dynamics of the market partition between the two service providers. Accordingly, the stationary probability density function reveals itself to be either uni- or bimodal depending on the values of a tuple of external control parameters. In other words, we explicitly show that a noise-induced phase transition occurs between regimes where brand consideration dominates and regimes where the time delays are predominant.

Spatial Market Sharing Dynamics in Presence of Customers' History-Based Decision Policy

Summary. *In this chapter, we propose some prospective results concerning an extension for the spatial queueing system previously introduced and studied in Chapter 7. More precisely, we are here concerned about recurrent customers that autonomously modify their position within the market (which denotes their relative distance to the two considered service providers) after each received service. The customers take into account their last experimented waiting time when they compute their updated position in the market interval. Remember that the customers' location within the market ultimately determines their choice between the two servers. Despite the inherent complexity that is due to the customers' history-based numerous local actions (indeed the agents actually consider here their complete history within the system into their decision processes, as their current location within the market interval is the ultimate result of all their experimented waiting times until now), we observe emerging collective structures that reveal themselves to be highly robust. We provide in this chapter a preliminary study of this particular queueing system and we give a first exploration, via simulation experiments, of the self-organizing dynamics appearing in this system. More particularly, we unveil the agent-induced collective spatio-temporal patterns that might emerge in the present context, namely the temporal oscillations driving the market sharing between the two service providers and the periodic wave governing the customers' spatial dispersion within the market interval.*

8.1 Introduction

Following the example of Chapter 7, we will again focus in this chapter on queueing networks (QNs) for which the spatial dimension enters explicitly into the dynamic modelling. As defined in [23], the main new feature of spatial queueing systems, which is not covered by classical QN theory, is the direct importance of the agents' location within the area that is covered by the different servers composing the network. As a perfect illustration, note that a very natural extension of classical queueing models towards queues with a structured space in which users are served finds for example a relevant application in mobile communication systems. The model previously

introduced and studied in Chapter 7 (remember that it has been inspired by the famous Hotelling duopoly model, [62]), which involves a configuration with two distinct servers sharing a spatially distributed market, is already a basic spatial queueing system. In the present chapter, we extend somehow this model and we study the same market configuration but in presence of history-based (HB) routing mechanisms. In particular, we consider recurrent customers that sequentially modify their location within the market interval (and hence their upcoming choice between the two service providers) after each service in function of their suffered waiting time. As we will see, to consider this type of agents' HB routing rules within our two servers spatial queueing system leads to the generation of self-organized spatio-temporal collective patterns. More particularly, the emerging global structures take here the form of stable oscillations of the market partition respectively held by the two servers (*i.e.* periodic cannibalization effects) and a periodic wave governing the customers' spatial dispersion within the market interval.

The chapter is organized as follows. In Section 8.2, we describe the extended model considered throughout this chapter. In particular, we introduce the HB decision mechanism that governs the customers' recurrent choice between the two servers. In Section 8.3, we describe the dynamics arising for this particular queueing system and we unveil the emerging spatio-temporal collective structures that are solely due to the customers' numerous local actions and interactions. The chapter ends with Section 8.4, in which some concluding remarks and further perspectives are given.

8.2 Model

In this chapter, we modify and extend in a way the configuration introduced in Chapter 7. Our starting point is again a two servers Hotelling's model where the two service providers S_1 and S_2 are located in a (linear) market confined on a segment $\Omega := [-\Delta, +\Delta] \subset \mathbb{R}$, $\Delta > 0$. The positions of the service providers are denoted respectively by x_1 and x_2 and satisfy $x_1 < x_2$. The servers S_1 and S_2 charge respectively prices p_1 and p_2 . Following the example of Chapter 7, queueing processes are considered in front of S_1 and S_2 and their service times are generally distributed with rate μ_1 and μ_2 respectively. To simplify the presentation, we restrict in the following on fully symmetric configurations (*i.e.* $x_1 = -x_2$, $\mu_1 = \mu_2 = \mu$ and $p_1 = p_2 = p$) and we suppose, without loss of generality, that $\Delta = 1$. Now, departing from the configuration studied in Chapter 7, we consider here recurrent customers that proceed to successive visits to the servers. The number of these customers remains constant and is fixed to N . We suppose furthermore that the initial position $y_{\zeta,1}$ of each customer $\zeta \in \{1, 2, \dots, N\}$ is uniformly distributed over the whole market interval. While in Chapter 7 the expected (*i.e.* a priori) waiting times were driving the customers' selection between the two servers, here this choice

will take into account the experimented waiting times (*i.e.* the customers will hence consider individual, a posteriori, HB data).

Starting from an initial position $y_{\zeta,1} \in \Omega$ that denotes his/her personal a priori brand departure to the two service providers, a customer ζ updates his/her location within the market interval from $y_{\zeta,n}$ to $y_{\zeta,n+1}$ (and hence his/her preference profile with report to both servers) after his/her n^{th} received service in the following way:

$$y_{\zeta,n+1} = y_{\zeta,n} - (1 - |y_{\zeta,n}|) \text{sign}(y_{\zeta,n}) \text{tgh}(\gamma\mu W_{\zeta,n}) \in \Omega, \quad n \in \mathbb{N}^*,$$

where $W_{\zeta,n}$ is the waiting time experimented by ζ to receive his/her n^{th} service and γ is a dimensionless parameter that quantifies the customers' loss of satisfaction due to waiting. Note that the market sharing dynamics considered here will differ from the one described in Chapter 7 in the following manner:

- In Chapter 7, the customers' locations were fixed (*i.e.* they were uniformly distributed over the whole market interval Ω) and the market partition boundary point Y_t was dynamically moving with respect to the evolution of the queue contents.
- Here the market partition boundary point is fixed (it is equal to 0 in the fully symmetric configuration considered here) but the recurrent customers are constantly travelling within the market interval in function of their experimented waiting times.

In the present context, a customer ζ chooses for its n^{th} service, $n \in \mathbb{N}^*$, between the two servers with respect to the following rule:

$$\mathcal{R}_s = \begin{cases} \text{go to } S_1 & \text{if } y_{\zeta,n} < 0, \\ \text{go to } S_2 & \text{if } y_{\zeta,n} \geq 0. \end{cases}$$

According to these autonomous decisions, the evolution of the customers within the market interval ultimately determines the market sharing dynamics between the two service providers. A sketch of the present modelling framework can be found in Fig. 8.1.

In one sense, the present configuration can be interpreted as a situation where the customers remain loyal to a service provider insofar as their experimented waiting times with this server until now remain somehow more satisfactory than the ones experimented with the other provider. More precisely, the customers react to any suffered waiting time with a service provider by modifying accordingly their location within the market interval to get closer to the other one (*i.e.* they update their preference profile in favour of the other service provider).

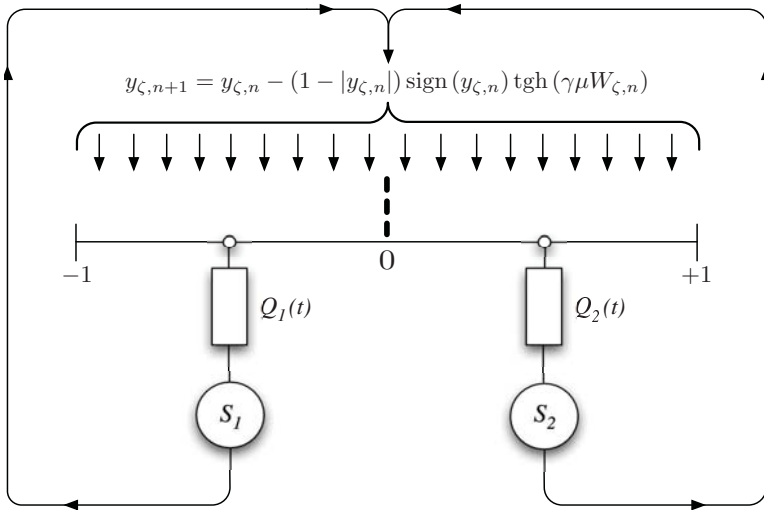


Fig. 8.1. Linear market with recurrent customers who update, based on waiting time satisfaction, their position within the market interval after each completed service.

8.3 Exploration via Simulation Experiments - Prospective Results

While it was possible to compute the stationary probability distribution of the boundary point defining the market partition (between the two service providers) in the original model studied in Chapter 7, the introduction of HB routing mechanisms implies that no stationary state exists here for the market sharing dynamics. Indeed, as it is illustrated in Fig. 8.2, the evolution of the market partition always remains time-dependent in the present context. More exactly, the underlying presence of delay effects in the customers' autonomous HB mechanisms¹ leads to the emergence of a stable cyclo-stationary behaviour of the queue contents². As the number of customers N in the system is fixed, the two queue lengths have a perfectly antagonistic behaviour and the state of the system is thus entirely given by the difference between the two queue contents. Simulation experiments indicate the following facts:

¹ Indeed, there exists a time lag between the moment a customer joins a queue and the moment he/she will modify, after having received the service, his/her location within the market interval and hence his/her upcoming choice between the service providers.

² In this configuration, the queue contents wholly determine the respective market shares held by the two service providers.

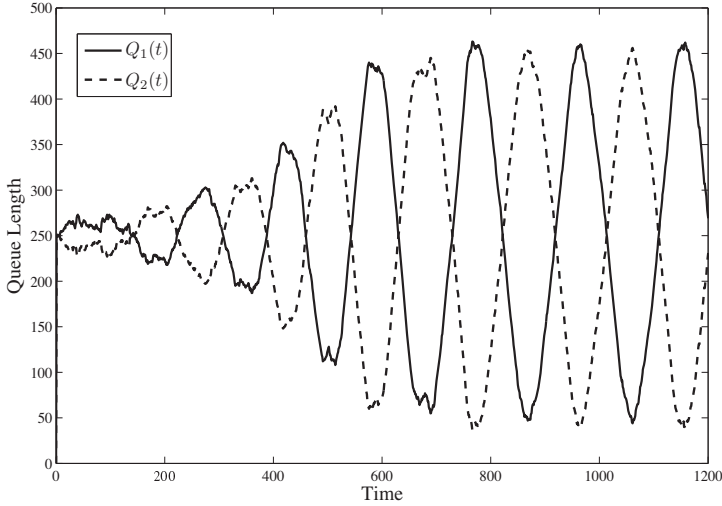


Fig. 8.2. Queue content dynamics when the service times of both S_1 and S_2 are uniformly distributed in $[0.01, 0.19]$, $N = 500$ and $\gamma = 0.002$.

- The amplitude Δ of the oscillations does not depend on the service rate μ (which is here supposed common to both servers) and shows to satisfy the following form:

$$\Delta = \beta N, \quad \beta \in [0, 1], \quad (\text{simulation results indicate a value } \beta = 0.75).$$

- The relaxation time t_{relax} towards the stable cyclo-stationary regime shows to have no dependency with the number N of customers in the system. On the other hand, t_{relax} is obviously a function of both μ and γ . In more details:

- (1) t_{relax} clearly decreases when μ increases,
- (2) t_{relax} grows when γ goes down (*i.e.* when the customers become less reactive to waiting times).

The local dynamics of a typical customer in our particular queueing system is illustrated in Fig. 8.3. Each customer ζ , $\zeta \in \{1, 2, \dots, N\}$, autonomously modifies his/her location within the market interval following a stochastic process that depends (with a variable time delay) on the random evolution of the queue contents (indeed, the experimented waiting times are directly related to the queue contents). It is interesting to observe that such highly irregular customer moves within the market interval give rise to a stable periodic wave

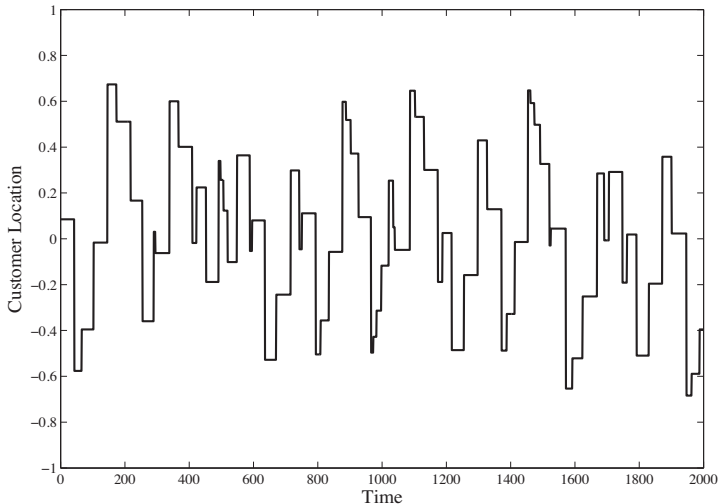


Fig. 8.3. Position dynamics of a typical customer $\zeta \in \{1, 2, \dots, N\}$ within the market interval when the service times of both S_1 and S_2 are uniformly distributed in $[0.01, 0.19]$, $N = 500$ and $\gamma = 0.002$.

which governs the dynamics of the customers' dispersion on the market segment, as illustrated in Fig. 8.4.

The underlying stability of the emerging spatio-temporal patterns (*i.e.* the queue length oscillations and the periodic wave governing the customers' repartition within the market interval) is a direct consequence of the law of large numbers (LLN). Indeed, when the population N of customers in the system is large, it is obvious, besides a manifest averaging effect, that the queue contents are likely to be important. Consequently, for this type of regimes, the waiting times become an almost deterministic function of the queue length, which consequently removes some of the fluctuations that affects the dynamics. In that sense, the regimes where the LLN holds leave some hope for the future derivation of an analytical description of the self-organized structures described in this chapter. Note finally that these emerging spatio-temporal patterns are here again entirely due to the agents' HB local actions and stigmergic interactions. Indeed, the customers, thanks to their individual decisions taken over time, autonomously stabilize both queues by triggering periodic purgings of their content. More precisely, when a queue content becomes large, it is likely that customers will leave the corresponding service provider because they will be unsatisfied with service (*i.e.* they will massively

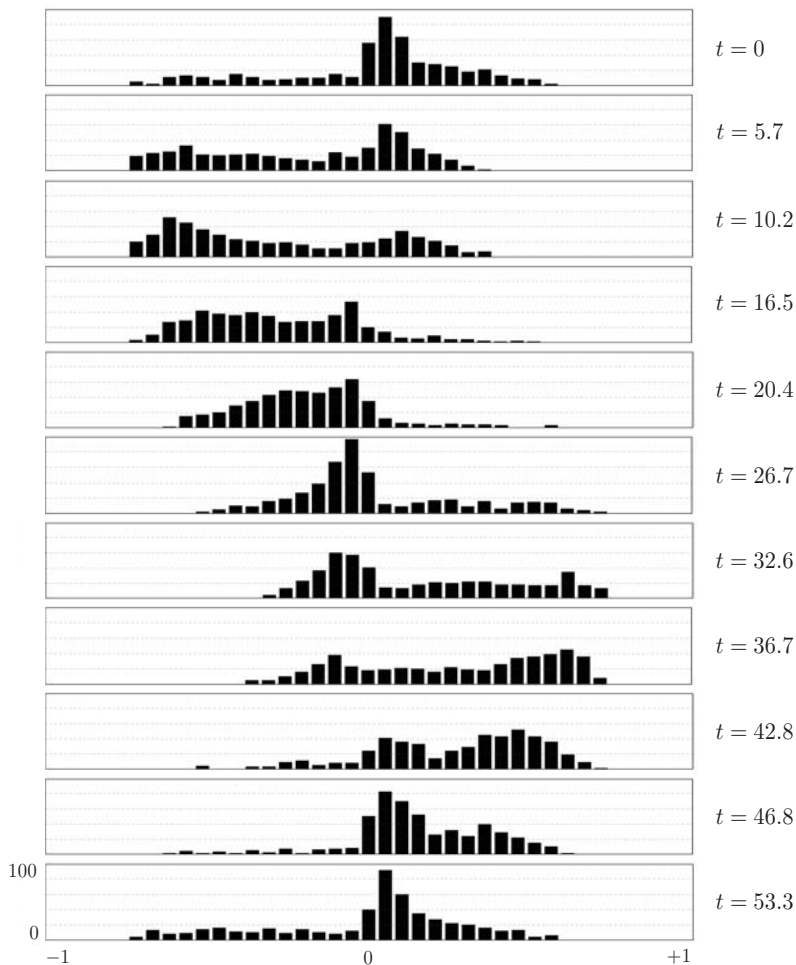


Fig. 8.4. Dynamics of the customers' spatial dispersion within the market interval. The service times of both S_1 and S_2 are uniformly distributed in $[0.01, 0.19]$, $N = 500$ and $\gamma = 0.002$.

modify their location within the market interval towards the other server). It creates thus a flow of unsatisfied customers that join the other service provider. Due to the time delay that exists in the customers' local dynamics (*i.e.* a customer reacts to an excessive queue length only when leaving the server), this oscillatory behaviour lasts forever and will never converge to a stable stationary equilibrium. As a consequence, the emerging global dynamics hence

exhibit periodic cannibalization effects (see Chapter 6 for more details) that are successively provoked by the two servers.

8.4 Concluding Remarks

The dynamics presented in this chapter are obviously the result of a stylized modelling of reality. However, the emerging spatio-temporal structures turn out to be very rich even for such an idealized class of models. Moreover, it is likely that analytical techniques could be further derived in order to describe precisely the collective patterns observed here only with the help of simulation techniques. In order to get closer to realistic situations, the configuration considered in the present chapter could be extended as follows. The customers' local behaviour could be individualized such that their personal location within the market interval would now be autonomously computed in the following manner:

$$y_{\zeta,n+1} = y_{\zeta,n} - (1 - |y_{\zeta,n}|) \operatorname{sign}(y_{\zeta,n}) \operatorname{tgh}(\gamma_{\zeta} \mu(W_{\zeta,n} - P_{\zeta})) \in \Omega, \quad n \in \mathbb{N}^*,$$

where P_{ζ} (resp. γ_{ζ}) is a patience threshold (resp. a reactivity parameter) that would be randomized and specific to each customer ζ . Under this new configuration, a customer ζ would be satisfied when his/her waiting time $W_{\zeta,n}$ would be smaller than his/her patience parameter P_{ζ} and ζ would accordingly modify his/her location within the market interval in order to get closer to the experimented service provider. On the other hand, when the waiting time $W_{\zeta,n}$ would become larger than P_{ζ} , ζ would modify his/her position in order to depart from the experimented server. For this new configuration, we still observe, via simulation experiments, queue content oscillatory behaviours for small average values of the patience parameter $\mathbb{E}(P_{\zeta})$ but the amplitude of the oscillations decreases when this average value increases. This can be easily understood as the larger is the average value of the patience parameter, the greater is the chance that the customers remain satisfied with a service provider. For large values of $\mathbb{E}(P_{\zeta})$, we might even observe overall stabilization of the system to a stationary state (*i.e.* all the customers are attached to a service provider and remain loyal to it forever).

To conclude and on a different note, it would be interesting to let the servers have the possibility to move within the market interval, thus allowing for competition processes between the two servers. In particular, lasting cannibalization effects (*i.e.* a server attracts permanently the whole market) might be expected in this case.

8.5 Contributions of Chapter 8

- We modify and extend the spatial queueing system introduced in Chapter 7 by considering now recurrent customers that update their preference profile

with respect to the two servers after each received service in function of individually experimented waiting times. We unveil the possible collective spatio-temporal structures that might emerge in this context due to this autonomous agents' HB local behaviour.

Towards Possible Applications

Agent-Based Optimal Real-Time Load Sharing - Application to Manufacturing Systems

Summary. *In this chapter, we propose a new real-time load sharing policy (LSP), which optimally dispatches the incoming workload according to the current availability of a set of operators. Optimality means here that the global service permanently requires the engagement of a minimum number of operators while still respecting due dates. To cope with inherent randomness due to operator failures as well as non-stationary fluctuating incoming workload, any optimal LSP rule will necessarily rely on real-time updating mechanisms. Accordingly, a permanent monitoring of the traffic workload, of the queue contents and of other relevant dynamic state variables is often realized by a central workload dispatcher. In this contribution, we abandon such a “classical” approach and we propose a fully decentralized algorithm which fulfils the optimal load sharing process. The underlying decentralized decisions rely on a “smart tasks” paradigm in which each incoming task is endowed with an autonomous routing decision mechanism. Incoming jobs hence possess, in this work, the status of autonomous agents endowed with “local intelligence”. Stigmergic interactions between these agents cause the optimal LSP to emerge. We emphasize that beside a manifest strict relevance for applications, our class of models is analytically tractable, a rather uncommon feature when dealing with multi-agent dynamics and complex adaptive logistics systems.*

9.1 Introduction

The reduction of manpower or other resource costs is an everlasting managerial challenge in any production and service network. Such contraction of the operating costs obviously relies on an optimized workload sharing between the available operators. Processing the full incoming load by using the minimum number of available operators, while still respecting given due dates, is clearly the basic optimization objective. The operator random failures as well as the non-stationary fluctuating incoming workload force the optimal load sharing policy (LSP) to be based on a permanent monitoring of the system state (*i.e.* queue contents, instantaneous traffic, etc). While this information updating process, on which our adaptive optimal LSP will be based, is often fulfilled by

a central dispatcher, our present contribution shows how fully decentralized mechanisms, of multi-agent type, are also perfectly suitable to achieve the same objective.

A large body of the available related literature pays attention to the customer side. Following that, the problem consists in minimizing the customers' average waiting time and, therefore, one tries to balance the incoming work such that the maximum total load on each server is minimized. Referred as adaptive load balancing, this classical problem has first been addressed using centralized management (see [19] for instance) and more recently by using decentralized mechanisms (see [22] among others). In this work, we adopt the complementary point of view of the service provider and try to minimize the number of engaged operators while nevertheless respecting due dates. Note that although our model does not, *stricto sensu*, optimize customer satisfaction, it allows however to impose an upper-bound to the maximum waiting time in the system by an ad-hoc tuning of the control parameters. While numerous aspects of load balancing and load sharing have been abundantly discussed over the last three decades, it is surprising that relatively little attention has been devoted to information gathering costs. Along these lines, let us mention contribution [21], where with the aim to minimize the average waiting time, the authors take explicitly into account the monitoring costs. The ultimate goal in [21] is to find a trade-off between the benefit and the costs of information gathering needed for any adaptive load sharing mechanism. To that purpose, an autonomous load sharing mechanism is derived, which adapts optimally the number of monitored servers to the current workload. This study is hence somehow related to the present work, where our aim is to optimally determine the number of servers to engage in order to face the current load.

The present study shares several similarities with well-known congestion control problems arising in the Internet, [1, 10, 44, 52, 103], where one tries to regulate the data flows to avoid congestion at servers (*i.e.* the so-called gateways in the Internet, which correspond to the operators in the present case). In both cases, the ultimate goal is to simultaneously ensure queue stability and maximization of resource utilization (busy period here and throughput in the Internet framework). To that purpose, the usual technique is to implement feedback information flows to warn about server congestion. While the presence of randomness definitely favours flexible and decentralized management in both contexts, there exist however manifest differences between the two modelling frameworks. Indeed, the agent character is in the present paper carried by the circulating tasks themselves while it is managed for the most part by the servers in the Internet. While in congestion control problems, fairness between the different users (in terms of throughput and/or delay) is essential, this feature is not required in the present case. Furthermore, it is common in congestion control mechanisms to use randomization to discard packets when

a buffer gets congested and hence ensure such kind of fairness between users. On the contrary, noise is used in our framework to dispatch the tasks between the servers, hence forming a noise-induced stabilization mechanism. Several congestion control rules [44, 52, 103] generate stable oscillations of the queue content and these will also be observed here. In both cases, the oscillations are due to the presence of information delay in the underlying controller, a phenomena thoroughly investigated for self-interested agents competing for common resources [74].

To implement decentralized control mechanisms, as those introduced in this chapter, within complex manufacturing, supply or production systems becomes nowadays an essential objective that is widely accepted and it gives thus rise to a prolific dedicated research activity. In this context and in close connection to the present work, Armbruster et al. introduce in [5] an autonomous control approach to manage dynamic and fluctuating production networks. As it is proposed in this chapter, the flow traffic management is in this contribution fully decentralized to the circulating parts. More precisely and quite similarly to the present situation, the parts autonomously compute their routing decisions through the network based on backward propagated information about the throughput times for the different possible routes, reproducing thus the way social insects communicate using pheromones¹. The approach proposed in [5] and the present one share several essential similarities but they nevertheless differ in the ultimate goal they are intended to achieve. Indeed, while in [5] the goal is to reduce the parts' throughput time within the network, the objective here is to minimize the number of engaged operators (while however ensuring simultaneously bounded throughput times).

The emerging dynamics to be discussed here exhibit optimal load sharing that is entirely due to decentralized history-based routing mechanisms implemented by the circulating tasks. As already emphasized in Section 3.10, the complete lack of central management, the agents' autonomy, their ability to learn and the stigmergic agent type of interactions between them imply that our class of models belongs to the field of Complex Adaptive Logistics Systems (CALs). Not restricted to manufacturing systems, remember that instances of CALs might also appear in the framework of supply chains. In that regard, the need in supply chain management for coordination strategies leading to adaptive, flexible and collective emerging behaviours is exhibited in [111] and it is showed furthermore that coherent global system dynamics can be generated by using only elementary components with local interactions. This contribution explains thus how basic concepts and operational tools of Complex Adaptive Systems (CAS) fit naturally and efficiently to characterize

¹ A pheromone is a chemical signal left by an individual on its way that will trigger a natural response in other members of the same species. As a typical example, ants mark their paths with the help of pheromones.

as well the supply chains dynamics. While the relevance and legitimacy of CAS have since long been emphasized in basic sciences (*i.e.* physics, chemistry and biology), it is remarkable that CAS also strongly enter into the engineering world, for example in logistics [63, 127], in traffic issues [78] and in production and service systems [45, 46, 47] as shown in this thesis. Note in addition that, besides its direct relevance to load sharing problems, the analytical tractability of the present class of models contributes to enrich the, so far, short list of analytically solvable CALS. In closing, we emphasize that among several possibilities to implement our algorithm in practice, Radio-Frequency-Identification-Devices (RFID) attached to the incoming tasks provide a natural solution (see Section 3.9). Indeed, the RFID technology available nowadays directly allows for the implementation of local intelligence to circulating items in production systems, as it is testified in [31, 81, 113], which explore how this technology leads to an effective and efficient management of various business processes.

The chapter is organized as follows. In Section 9.2, we describe our basic modelling framework, namely a general multi-server production centre fed by an incoming flow of “smart tasks”. In Section 9.3, we introduce our multi-agent type dynamic load sharing algorithm. In Section 9.4, we study the emergence of self-organized stable load sharing, by using analytical considerations as well as simulation results. In Section 9.5, we describe in more details the oscillatory behaviour that appears in the queue content dynamics. To finish, Section 9.6 is devoted to several conclusions and perspectives.

9.2 Basic Modelling Framework

We consider a production centre fed by an incoming flow of tasks modelled by a non-stationary, random renewal process with rate $\lambda(t)$. The production center is therefore a generic queueing system with N parallel servers. Each incoming task ζ_j ($j \in \mathbb{N}^*$) requires a specific amount $U_j > 0$ of processing time. The U_j 's are characterized by i.i.d. random variables with general probability distribution, the mean of which is fixed to 1.

The objective is to realize an optimal load sharing policy defined by:

(\mathcal{O}): *Optimal Load Sharing Policy (LSP)*

- (i) **“Process the global incoming workload by permanently engaging the minimal number of available servers”** or equivalently, using queueing theory terminology, **“maximize the busy period of the engaged servers”**.
- (ii) **“Keep the average waiting time below a given level”**.

To achieve the objective \mathcal{O} , one can rely either on a centralized solution (*i.e.* a central dispatcher) or on ad-hoc decentralized control mechanisms. Our aim here is to realize \mathcal{O} by using a multi-agent fully decentralized framework. Such decentralization permanently ensures strong reactivity and high flexibility to cope with random and non-stationary environments. In the sequel, we assume non-preemptive LSP (*i.e.* a task cannot be transferred from one server to another after its execution has started).

Server Parameters.

Each server M_α , $\alpha \in \{1, 2, \dots, N\}$, is characterized by:

- (1) its processing rate μ_α ,
- (2) its queue capacity parameter $\mathcal{C}_\alpha > 0$, that plays the role of a congestion threshold,
- (3) a two states (“open” or “closed”) warning semaphore S_α , whose aim is to relay information about possible congestion states at M_α and whose sensitivity is tuned by the queue capacity parameter \mathcal{C}_α ,

“Smart Task” Agent Character of Incoming Jobs.

Each incoming task ζ_j , $j \in \mathbb{N}^*$, has the capability to:

- (1) identify the server M_α , $\alpha \in \{1, 2, \dots, N\}$, that did process the task,
- (2) record its sojourn time $\tau_{j,\alpha}$ spent into the system while served by M_α (*i.e.* $\tau_{j,\alpha} = W_{j,\alpha} + V_{j,\alpha}$, $W_{j,\alpha}$ being the waiting time in the queue and $V_{j,\alpha}$ the processing time),
- (3) compute a set of individual dispatching probabilities $p_\alpha(t)$, $\alpha \in \{1, 2, \dots, N-1\}$, that characterize for each task an autonomous routing strategy within the network,
- (4) read the state of the semaphores S_α attached to each server.

As we have assumed the incoming task size to be i.i.d. with mean 1, the service times of M_α inherit the randomness and are hence also i.i.d. random variables with identical distribution and mean $\frac{1}{\mu_\alpha}$.

9.3 Multi-Agent Type Algorithm

A “smart task” behaves as follows.

(1) *On entry.*

An incoming task ζ_j ($j \in \mathbb{N}^*$) at time t first reads the state of semaphore S_1 (which gives information on the current congestion state of server M_1). If S_1 is open, the task enters M_1 . If S_1 is closed, the task routing is: enter M_1 with probability $p_1(t)$ and, with probability $(1 - p_1(t))$, read the state of S_2 to tentatively join M_2 . If S_2 is open then M_2 processes the task. If S_2 is closed, the task enters with probability $p_2(t)$ into M_2 and with probability $(1 - p_2(t))$ reads the state of S_3 to tentatively join M_3 . The rule is then applied iteratively.

(2) *On exit.*

If the sojourn time $\tau_{j,\alpha}$ of the outgoing task ζ_j , processed by server M_α , exceeds \mathcal{C}_α (*i.e.* $\tau_{j,\alpha} > \mathcal{C}_\alpha$) then ζ_j sets S_α to the state “closed”. In words, when a task processed by M_α has spent \mathcal{C}_α in the system, congestion is detected and this task triggers immediately the closing of S_α even if the processing is not yet completed. Conversely, on exit, provided $\tau_{j,\alpha} \leq \mathcal{C}_\alpha$, ζ_j sets S_α to the state “open”, which indicates that M_α is not congested anymore.

Note that, even when a semaphore S_α is closed, a $p_\alpha(t)$ -based partial incoming traffic continues to be processed by a congested server M_α . The tasks joining such an overloaded server ultimately enable the reopening of the associated semaphore as soon as the workload becomes undercritical (*i.e.* when $\tau_{j,\alpha} \leq \mathcal{C}_\alpha$). The ability for each travelling task to monitor information left by predecessors and to process this information to autonomously decide its routing strategy confers to this dynamics a manifest adaptive multi-agent character (see Fig. 9.1 for a summarizing sketch of the present modelling framework).

9.4 Emergence of Optimal Load Sharing Dynamics

From now on, we assume that the number N of potentially available servers is sufficient to always handle the offered incoming workload, *i.e.*

$$\frac{\lambda(t)}{\sum_{\alpha=1}^N \mu_\alpha} < 1.$$

By the construction of the multi-agent dynamics, given in Section 9.3, our LSP automatically ensures permanently (*i*) the engagement of the minimal number of servers and (*ii*) queue stability. As exposed later in Section 9.5, we observe (see Fig. 9.2) that for large enough \mathcal{C}_α 's, $\alpha \in \{1, \dots, N - 1\}$, the queue contents exhibit stable temporal oscillations whose maximum values are given by:

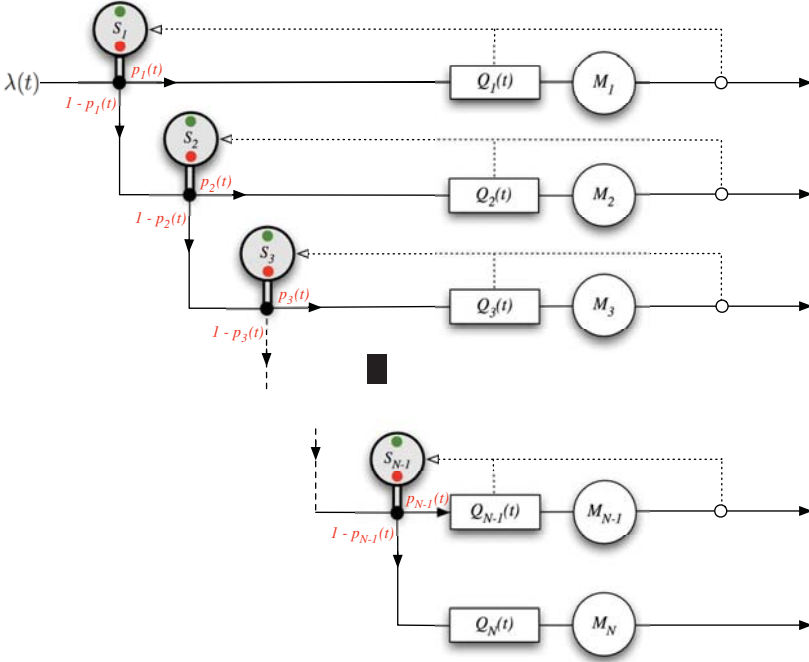


Fig. 9.1. N -parallel servers queuing system with decentralized load sharing mechanism.

$$\begin{cases} Q_{\max,1}(t) = \mathcal{C}_1 \lambda(t) - 1; \\ Q_{\max,\alpha}(t) = \mathcal{C}_\alpha \lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) - 1, \alpha \in \{2, \dots, N - 1\}. \end{cases} \quad (9.1)$$

Tuned by the control parameters \mathcal{C}_α , $\alpha \in \{1, \dots, N - 1\}$, these maximum possible queue contents can be used to calibrate the waiting room sizes and hence, due to Little’s law, to limit the task waiting times. In particular, tasks subject to deadlines can hence be handled within the present framework. As shown later in Section 9.5, it is here worth to emphasize that the queue content of the engaged servers never vanishes, thus ensuring maximum busy period and hence optimal load sharing.

Let us now discuss in more details the role played by the dispatching probabilities $p_\alpha(t)$ introduced in Section 9.2. Remember that a congested server continues to be fed by a reduced incoming flow with rate:

$$\begin{cases} p_1(t)\lambda(t) & \text{for server } M_1, \text{ and} \\ p_\alpha(t)\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) & \text{for server } M_\alpha, \alpha \in \{2, \dots, N - 1\}. \end{cases}$$

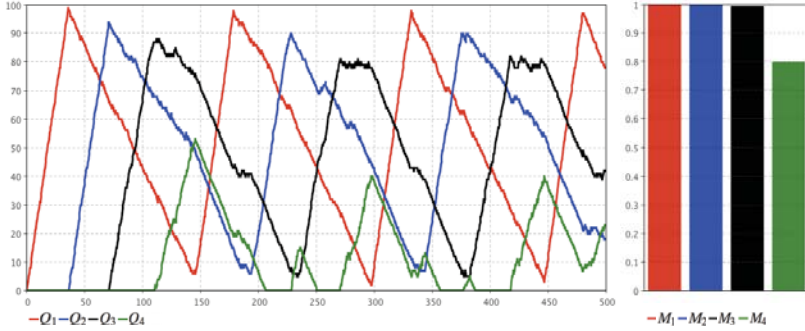


Fig. 9.2. *Left:* Temporal evolution of the queue contents for interarrival times uniformly distributed in $[0.16; 0.36]$ ($\lambda(t) = 3.8$), $K = 10$, service times uniformly distributed in $[0.5; 1.5]$ ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = 1$), $C_1 = C_2 = C_3 = C_4 = 26$ and $\epsilon_\alpha = 0.22 - 0.7(\alpha - 1)$, $\alpha \in \{1, 2, 3\}$. *Right:* Corresponding server utilization.

Consequently, after a congestion occurs, the queue in front of the congested server will effectively decrease iff the following condition is satisfied:

$$\begin{cases} p_1(t)\lambda(t) - \mu_1 < 0 & \text{for } M_1, \text{ and} \\ p_\alpha(t)\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) - \mu_\alpha < 0 & \text{for } M_\alpha, \alpha \in \{2, \dots, N-1\}. \end{cases} \quad (9.2)$$

To fulfil condition (9.2) the dispatching probabilities $p_\alpha(t)$, $\alpha \in \{1, \dots, N-1\}$, have to be chosen as:

$$p_1(t) = \frac{\mu_1}{\lambda(t)} - \epsilon_1$$

and

$$p_\alpha(t) = \frac{\mu_\alpha}{\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t))} - \epsilon_\alpha, \quad \alpha \in \{2, \dots, N-1\},$$

with $\frac{\mu_1}{\lambda(t)} > \epsilon_1 > 0$ and $\frac{\mu_\alpha}{\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t))} > \epsilon_\alpha > 0$, $\alpha \in \{2, \dots, N-1\}$. This choice strictly ensures the stability of the queue contents. The smaller is the value of the ϵ_α 's, the closer to optimality is the load sharing (*i.e.* the busy period of the engaged server M_α converges to 1 when $\epsilon_\alpha \rightarrow 0$). A too drastic reduction of the partial traffic (ϵ_α large $\Rightarrow p_\alpha(t)$ small) yields poor reactivity of the system. Indeed, the resulting long delay before the reopening of the semaphore is likely to empty the queue, thus leading to a decrease of the busy period.

As noted in [21], the incoming flow rate $\lambda(t)$ can itself be estimated, in real-time, by elementary agent-to-agent interactive mechanisms. This ultimately enables each task ζ_j ($j \in \mathbb{N}^*$) entering the system at time t_j , to estimate autonomously the ad-hoc dispatching probabilities $\bar{p}_{\alpha,j} = p_\alpha(t_j)$, $\alpha \in \{1, \dots, N-1\}$, which characterize its routing strategy through the network. With this, our load sharing algorithm becomes fully decentralized, all

routing decisions being taken by the circulating items themselves. One basic possibility to implement such a decentralized traffic estimation reads as follows.

Multi-Agent Type Traffic Load Estimator.

Each task ζ_i ($i \in \mathbb{N}^*$) stores, upon arrival, its entry time t_i in the system on a register permanently accessible to the other tasks. The traffic estimator $\bar{\lambda}_j(t)$, computed by the incoming task ζ_j ($j \in \mathbb{N}^*$), relies on an observation window of size K . ζ_j reads the entry-time t_{j-K} of the K^{th} preceding task and estimates the instantaneous traffic by:

$$\bar{\lambda}_j(t) = \frac{t_j - t_{j-K}}{K}.$$

As clearly illustrated in Fig. 9.2, this very rough estimation of the global incoming traffic is sufficient to complete the overall objective \mathcal{O} and this despite the underlying randomness. There is obviously an optimal trade-off to select an appropriate value for the observation window K . Valuing mostly the reactivity, we prefer small values of K for the observation window. As illustrated in Fig. 9.3, small values of K lead indeed to highly reactive response (*i.e.* the

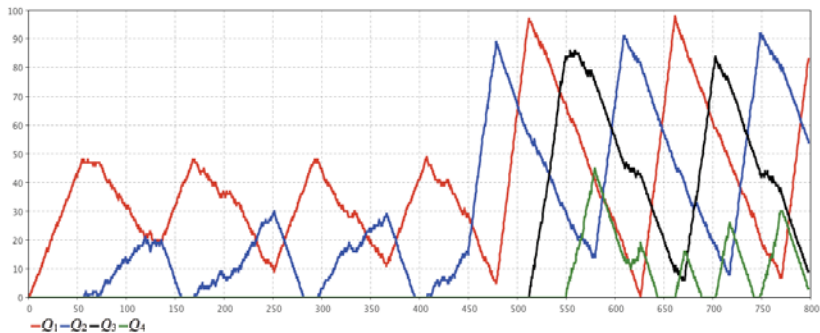


Fig. 9.3. Temporal evolution of the queue contents for service times uniformly distributed in $[0.5; 1.5]$ ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = 1$), $C_1 = C_2 = C_3 = C_4 = 26$ and $\epsilon_\alpha = 0.22 - 0.7(\alpha - 1)$, $\alpha \in \{1, 2, 3\}$. Interarrival times are uniformly distributed in $[0.32; 0.72]$ ($\lambda(t) = 1.9$) for $0 \leq t < 450$ and uniformly distributed in $[0.16; 0.36]$ ($\lambda(t) = 3.8$) for $t \geq 450$, $K = 10$.

length of the transient adaptive phase is almost negligible). This is in particular perfectly suitable for non-stationary incoming traffic loads. Quantitatively, the length of the adaptive phase thus only depends on the effective delay between the time a congestion effectively occurs and the time it is detected. This delay, for server M_α , is equal to C_α (see Section 9.5 for more details).

Note that depending on the specific management issues, larger observation windows K could also be selected whenever smooth reactions are required for system reliability or to avoid large set-up costs. Observe that most Internet congestion control mechanisms rely on relatively large values of K for the observation window to smoothly react to bursty traffic.

9.5 Queue Content Oscillatory Behaviour - Siphon Dynamics

We discuss here in more details the emergence of the temporal oscillations observed for the different queue contents, as illustrated in Fig. 9.2. As in Section 3.4, we focus for the discussion of this temporal oscillatory behaviour on a deterministic approach. This approach, due to the influence of the law of large numbers (LLN), is also relevant in presence of fluctuations when the C_α 's, $\alpha \in \{1, \dots, N-1\}$, are sufficiently large (*i.e.* quasi-deterministic stable cyclo-stationary queue oscillations emerge independently of the inter-arrival and service time distributions). Note qualitatively that the relative importance of the fluctuations around the task average sojourn time (which is the sum of the preceding tasks individual processing times) decreases for large queue content $Q_\alpha(t)$ (a quantitative characterization is given in [43]). Fig. 9.4 (in comparison with Fig. 9.2) explicitly exhibits that the larger are the C_α 's, the smoother are the oscillations.

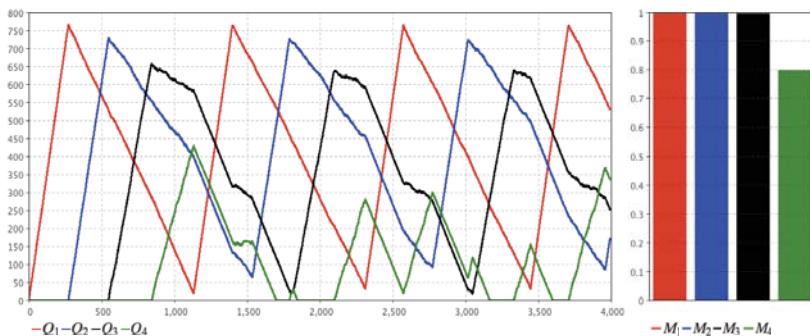


Fig. 9.4. *Left:* Temporal evolution of the queue contents for interarrival times uniformly distributed in $[0.16; 0.36]$ ($\lambda(t) = 3.8$), $K = 10$, service times uniformly distributed in $[0.5; 1.5]$ ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = 1$), $C_1 = C_2 = C_3 = C_4 = 200$ and $\epsilon_\alpha = 0.22 - 0.7(\alpha - 1)$, $\alpha \in \{1, 2, 3\}$. The smoothing effect due to the underlying LLN is manifestly observable by comparing Figs. 9.2 and 9.4. *Right:* Corresponding server utilization.

Along the lines exposed in Section 3.4, we start by characterizing the oscillations of the first queue content $Q_1(t)$. During an initial phase, $Q_1(t)$ increases at rate $\lambda(t) - \mu_1$, the whole traffic $\lambda(t)$ is indeed dispatched to M_1 , which is not yet overloaded. M_1 is considered as congested when $Q_1(t)$ reaches the level $\mathcal{C}_1\mu_1 - 1$. Indeed, at this time, a newly incoming task ζ_j will spend on average \mathcal{C}_1 in the system (*i.e.* its mean waiting time will be equal to $(\mathcal{C}_1\mu_1 - 1) \frac{1}{\mu_1} = \mathcal{C}_1 - \frac{1}{\mu_1}$ and its processing time will be equal on average to $\frac{1}{\mu_1}$). The queue $Q_1(t)$ reaches its auto-siphoning threshold when the congestion is first detected, which happens when ζ_j has waited \mathcal{C}_1 in the system. This starts a second operating phase during which $Q_1(t)$ decreases at rate $p_1(t)\lambda(t) - \mu_1$. This second phase lasts until a task detects that M_1 is not congested anymore; this happens after a time delay \mathcal{C}_1 initiated when $Q_1(t)$ reached again the level $\mathcal{C}_1\mu_1 - 1$. The alternation between these two operating phases creates queue content stable oscillations whose amplitude Δ_1 and period Π_1 are respectively given by :

$$\Delta_1(t) = (1 - p_1(t)) \mathcal{C}_1 \lambda(t)$$

and

$$\Pi_1(t) = \mathcal{C}_1 \left[2 + \frac{\mu_1 - p_1(t)\lambda(t)}{\lambda(t) - \mu_1} + \frac{\lambda(t) - \mu_1}{\mu_1 - p_1(t)\lambda(t)} \right].$$

The maximum and minimum values of these oscillations are given respectively by Eq. (9.1) and

$$Q_{\min,1}(t) = p_1(t)\mathcal{C}_1\lambda(t) - 1.$$

To understand the underlying delay mechanism, it is enlightening to visualize the queue dynamics by using the hydrodynamic analogy sketched in Fig. 9.5. We emphasize that contrary to the flow dynamics discussed in Section 3.4 where there is a feedback loop fed by physical items, here the feedback is purely informational (*i.e.* the items leave in any case the system after service but they deliver an indicative feedback to the following tasks).

The oscillation frequency of $Q_\alpha(t)$, $\alpha \in \{2, \dots, N\}$, are identical to $Q_1(t)$, hence

$$\Pi_\alpha = \Pi_1, \quad \alpha \in \{2, \dots, N\}.$$

This is illustrated in Figs. 9.6 and 9.7 (respectively corresponding to Figs. 9.2 and 9.4), where we exhibit the Fourier components (*i.e.* the spectrum) of the queue dynamics obtained by simulation for small, respectively large, values of \mathcal{C}_α , $\alpha \in \{1, \dots, N - 1\}$. As expected, for large values of \mathcal{C}_α , the spectrum exhibits a sharp mode (*i.e.* the signal-to-noise ratio is enhanced).

For servers M_α , $\alpha \in \{2, \dots, N - 1\}$, the oscillations exhibit an additional structure. Namely for $Q_\alpha(t)$, $\alpha \in \{2, \dots, N - 1\}$, there exists an alternation between three distinct operating phases:

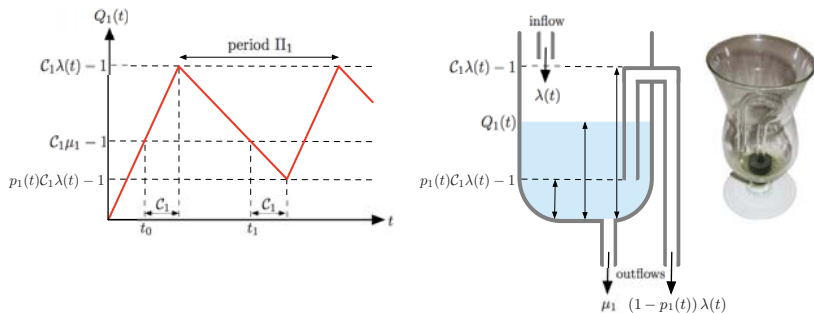


Fig. 9.5. Hydrodynamic analogy. *Left:* The task entering at t_0 is the first one of a whole cluster \mathcal{G} of tasks that will detect congestion. This task triggers the alternation of $Q_1(t)$ from the increasing to the decreasing state at $t_0 + C_1$. The last task belonging to cluster \mathcal{G} is the one entering the system just before t_1 and triggers the switch of $Q_1(t)$ from the decreasing to the increasing state at $t_1 + C_1$. This simple delay dynamics repeats and creates stable oscillations of the queue content. *Right:* The “Tantalus glass” siphon model. The water level corresponds to the queue length $Q_1(t)$. The continuous inflow and outflow rates are respectively given by $\lambda(t)$ and μ_1 . The periodic alternate siphoning outflow is $(1 - p_1(t)) \lambda(t)$. The siphon leaves a water residue of height $p_1(t)C_1\lambda(t) - 1$, due to the continuous inflow during C_1 . The effective siphon length is $(1 - p_1(t)) C_1\lambda(t)$.

(1) When all the servers $M_\beta, \beta \in \{1, \dots, \alpha - 1\}$, are overloaded and hence their semaphores are closed, M_α receives a traffic with rate:

$$\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)).$$

Indeed, from the full incoming workload $\lambda(t)$, one has to subtract the $p_\alpha(t)$ -based partial traffics that feed the congested servers. As a consequence, during this phase, $Q_\alpha(t)$ increases at rate $\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) - \mu_\alpha$.

(2) Once server M_α becomes congested, $Q_\alpha(t)$ starts to empty. Provided all servers $M_\beta, \beta \in \{1, \dots, \alpha - 1\}$, remain congested, the incoming traffic continues to be dispatched to M_α . As M_α is congested, it only receives a $p_\alpha(t)$ -based part of this traffic and $Q_\alpha(t)$ hence decreases at rate $p_\alpha(t)\lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) - \mu_\alpha$.

(3) Whenever one among the servers $M_\beta, \beta \in \{1, \dots, \alpha - 1\}$, is no longer congested (therefore its semaphore has been reopened), this server attracts the full incoming workload and hence M_α is not fed anymore. Thus, $Q_\alpha(t)$ decreases at rate $-\mu_\alpha$.

The alternation between these three phases is completely determined by the queue dynamics of the yet engaged servers. Basically, the time at which a queue starts to empty triggers the feeding of the next server to be engaged.

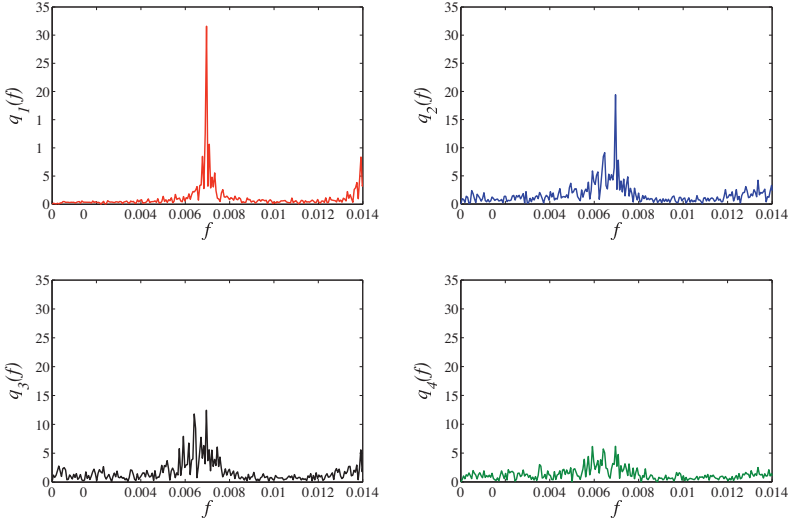


Fig. 9.6. Spectrum of the queue content dynamics for interarrival times uniformly distributed in $[0.16; 0.36]$ ($\lambda(t) = 3.8$), $K = 10$, service times uniformly distributed in $[0.5; 1.5]$ ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = 1$), $C_1 = C_2 = C_3 = C_4 = 26$ and $\epsilon_\alpha = 0.22 - 0.7(\alpha - 1)$, $\alpha \in \{1, 2, 3\}$.

The oscillatory behaviour of $Q_\alpha(t)$, $\alpha \in \{2, \dots, N - 1\}$, is characterized by maximum and minimum values given respectively by Eq. (9.1) and:

$$Q_{\min,\alpha}(t) = C_\alpha p_\alpha(t) \lambda(t) \prod_{i=1}^{\alpha-1} (1 - p_i(t)) - 1, \quad \alpha \in \{2, \dots, N - 1\}.$$

Consequently, the oscillation amplitude for these queue contents is given by:

$$\Delta_\alpha(t) = C_\alpha \lambda(t) \prod_{i=1}^{\alpha} (1 - p_i(t)), \quad \alpha \in \{2, \dots, N - 1\}.$$

9.6 Concluding Remarks

In a competitive environment, to attract new and to keep loyal customers is the basic concern of any service provider or manufacturer, which definitely requires a high service customization to match all specific demands. Service or product customization affects both the quantity and the nature of the incoming demands. Focusing on quantitative aspects, one should clearly expect that high service customization necessarily leads to non-stationary and highly

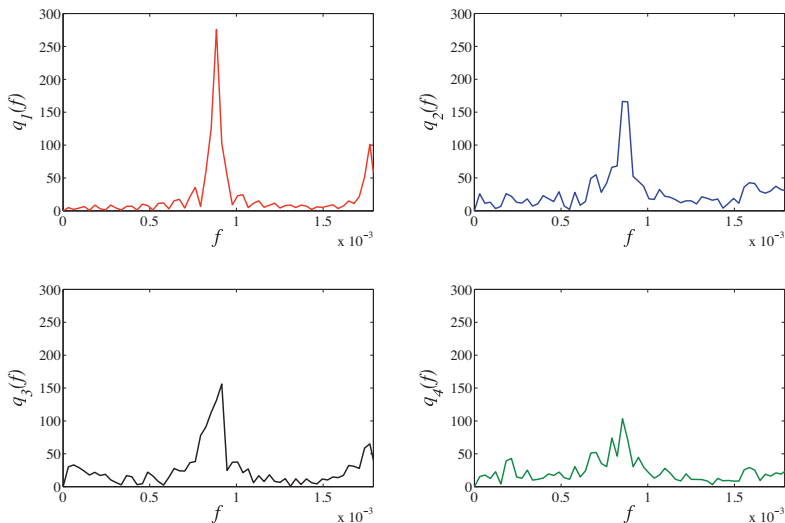


Fig. 9.7. Spectrum of the queue content dynamics for interarrival times uniformly distributed in $[0.16; 0.36]$ ($\lambda(t) = 3.8$), $K = 10$, service times uniformly distributed in $[0.5; 1.5]$ ($\mu_1 = \mu_2 = \mu_3 = \mu_4 = 1$), $\mathcal{C}_1 = \mathcal{C}_2 = \mathcal{C}_3 = \mathcal{C}_4 = 200$ and $\epsilon_\alpha = 0.22 - 0.7(\alpha - 1)$, $\alpha \in \{1, 2, 3\}$.

fluctuating rate of the service demand. Hence, constructing efficient service policies able, in such random and time-dependent environments, to entirely fulfil customer satisfaction while maintaining the operating costs at the lowest possible level is definitely a complex challenge. The ubiquitous non-stationary and fluctuating nature of the underlying demand imposes flexibility and reactivity to be key characteristics of any efficient algorithm required by the system management. Among other classical problems, we focus in this chapter on the construction of an optimal service policy which enables the sharing of the global incoming workload between a set of available servers. In general, such a load sharing policy will be achieved by a central operator which dispatches, by an on-going real-time information gathering of a set of relevant system state variables, the incoming jobs among the available servers. Leaving aside such a centralized point of view and hence following a vigorous recent trend emerging in production and service systems, we explicitly show how the same task can also be perfectly realized using a fully decentralized algorithm. Our basic idea relies on a multi-agent perspective implying the service management to be performed by the jobs themselves (*i.e.* “smart tasks” paradigm). We are able to explicitly show how autonomous smart tasks can ensure, in real-time, that the global load is processed by the minimum number of engaged operators while permanently avoiding the system to get congested. The intrinsic

simplicity of our algorithm, the analytical tractability of our models which offers an intimate understanding of the operating dynamics and, finally, the possibility, using RFID, to confer an autonomous agent character to incoming jobs, clearly suggest how the implementation could actually be concretely realized. The basic modelling framework proposed in this chapter offers several possibilities for refinements to cover different realistic situations, namely one could consider the following extensions:

- (1) To address “fairness” issues by adding an additional decentralized mechanism which would ultimately ensure that all incoming tasks wait in average the same time before being served.
- (2) To allow different task types as well as distinct server capabilities with multiple working levels (each working level of a server would correspond to the required service time to proceed with a specific type of task) and to introduce a simple matchmaking mechanism between task type and service capability.
- (3) To provide direct communication skills between agents (*i.e.* stigmergic interactions and direct interactions would hence coexist).
- (4) To introduce learning capabilities inside the system (repeated services would decrease the required service times, etc...).

9.7 Contributions of Chapter 9

- We propose a fully decentralized dynamic load sharing policy that optimally dispatches the incoming workload according to the current availability of a set of operators. The proposed algorithm permanently ensures the engagement of a minimum number of operators while still respecting due dates.
- Constructed on a “smart parts” paradigm, our agent-based model nevertheless allows for analytical analysis, a rather uncommon feature when dealing with fully decentralized manufacturing or logistics processes. In that regard, our particular model is one of the very rare solvable instances of complex adaptive logistics systems available so far in the literature.

Other Possible Applications

Summary. *The modelling framework developed in this thesis could definitely be extended to cover a very wide range of logistics problems in industrial and service applications. In this chapter, we give several ideas on how some of the original concepts approached in the present work could be transferred to two research fields of central importance, namely transportation and supply chains.*

10.1 Transportation Networks

The agent characteristic of the models presented in this work, as well as their specific history-based (HB) features, naturally indicate that similar emerging phenomena might be observed in various transportation networks (cars, trains, buses, metros,...). Indeed, the users in such systems clearly use, when they have the choice between different routes, their personal experience within the network to take their future routing decisions. This inherent path dependence implying delay effects in the dynamics, is it hence likely to observe any oscillatory behaviour in transportation networks?

Talking to practitioners, it comes out that permanent oscillatory behaviours are generally not observed in traffic systems. Because of that, it is usually assumed in the transportation area that a stationary state (which can be related, in terms of game theory, to a Nash equilibrium) exists. In real situations, it is possible to observe this stable behaviour thanks to field measurements, where one measures for example every 15 minutes the current congestion state at a specific point of the network. While there exist cyclic phenomena, that can be attributed to the seasons or to the different days of the week, as well as an increase of the mobility due to the economical progress (*i.e.* the economical sharings actually increase in our globalized society), practitioners are not aware of permanent oscillatory dynamics that could result from the users' local routing decisions. As we will see below, this has to be connected to the fact that users in transportation systems, due to their frequent and recurrent

(daily) use of the road network, are highly sensitive and reactive to waiting times. Consequently, their individual routing mechanisms turn out to be very complex, heterogeneous and highly path dependent (*i.e.* taking into account a large personal history within the network).

While permanent oscillations are not likely to emerge in transportation networks, it is nevertheless possible to observe oscillatory behaviours during the transient regime following a modification of the network (roadworks, opening or closing of a road, capacity increase). However, these oscillations turn out to be unstable and diminish over time. This deadening of oscillatory behaviours is probably due to the strong repetitiveness that governs the users' decisions in transportation systems. Indeed, we are nowadays more and more subject to frequent and long journeys the relative costs of which are increasingly important. For this reason, the optimization of transportation aspects has become a very sensitive topic to which people pay a great attention. The users' routing decisions are hence carefully thought through and often take into account their global experience within the transportation system. The time window used to determine the HB routing choices is doubtlessly larger in reality than what is assumed in the stylized models presented in this work. These more complex and refined decision mechanisms imply a reduction of the oscillations over time and lead to the emergence of self-regulation of traffic phenomena. Indeed, the oscillatory features of the dynamics are in one sense taken into account by the experienced users after some time. However, it is enlightening to note that the possibly emerging oscillatory transient phases are mainly due to the users' HB individual routing decisions and, in this sense, the models introduced in this thesis help to get a phenomenological understanding of these particular dynamics.

Besides the convergence to global stable states, transportation practitioners also observe, at the user level, various cyclic features in the individual routing decisions. These individual oscillations (*i.e.* at the microscopic level) are however not anymore observable at the macroscopic level. They are indeed counterbalanced between the different users, which are often in large number and highly heterogeneous in transportation networks. Globally, while there obviously remain in real situations fluctuations of small amplitude around the stationary state, these fluctuating components do in general not possess a cyclic structure. It is moreover enlightening to know that, contrary to the stylized models presented in this work, the intelligence is nowadays not only decentralized to the circulating items in real transportation systems. For example, the traffic lights autonomously adapt their control following real-time measurements of the traffic and will increase the *green* time slots in case of congestion. Note finally that some of the oscillatory phenomena that can nevertheless be observed in transportation networks are due to the influence that available traffic information (radio, GPS, internet,...) has on the pool of users. This information on the network status is often delayed and might hence give

rise to oscillatory behaviours. Novel approaches based on the selective use of misinformation given to users might help to reduce this kind of oscillations, [74].

In the following, we present an illustrative model of a transportation network, the dynamics of which exhibit the same queue length oscillations softened over time as those that might be observed in real traffic systems. Compared to the HB routing rules previously developed in this thesis, the present routing decision mechanism assumes that the users take into account a larger history window to determine their future route through the network.

We consider the network sketched in Fig. 10.1, which consists of two par-

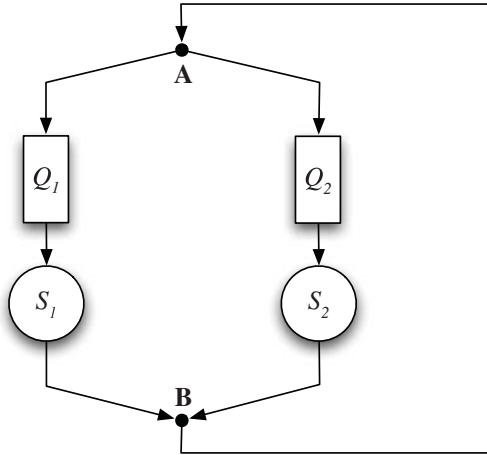


Fig. 10.1. Queueing representation of a transportation network with two alternative routes and recurrent users.

allel queues in a closed topology. Following classic concepts available in the existing literature (see [2, 27] among others), this particular queueing model represents a transportation network composed of two different routes linking two different locations **A** and **B**. Each route $i \in \{1; 2\}$ is characterized by its own processing rate μ_i . As our aim is to model the behaviour of users doing recurrently (typically daily) the same journey, we assume here a fixed pool of N commuters evolving in a closed network configuration. In addition to the capability of autonomously computing suffered waiting times, each user $\zeta \in \mathbb{N}$ is now also able to calculate its personal mean waiting time \bar{W}_ζ^1 (resp. \bar{W}_ζ^2) with respect to its complete history with server S_1 (resp. S_2). Initially set to 0, \bar{W}_ζ^1 (resp. \bar{W}_ζ^2) is updated by ζ after each run in route 1 (resp. route

2). Every day, before each run from **A** to **B**, each user ζ chooses its route according to the following rule:

$$\tilde{\mathcal{R}} = \begin{cases} \text{if } \bar{W}_\zeta^1 < \bar{W}_\zeta^2 & \Rightarrow \begin{cases} \text{choose route 1 with probability } 1 - \epsilon \\ \text{choose route 2 with probability } \epsilon \end{cases} \\ \text{if } \bar{W}_\zeta^1 \geq \bar{W}_\zeta^2 & \Rightarrow \begin{cases} \text{choose route 1 with probability } \epsilon \\ \text{choose route 2 with probability } 1 - \epsilon \end{cases} \end{cases}$$

with $\epsilon \in [0, 1]$. In words, ζ chooses with probability $1 - \epsilon$ the route for which it has, on average, experimented until now the smaller waiting time. In order to express human irrationality, it is nevertheless possible that ζ decides however, with probability ϵ , to take the other route. We assume that αN users (resp. $(1 - \alpha)N$), $\alpha \in [0, 1]$, initially choose route 1 (resp. route 2). In the case of a network modification (closed route, introduction of a new route,...), this a priori knowledge of the system status (*i.e.* the value of α) will strongly depend on the quantity and quality of the information given to the users. Note that in comparison to the model considered in [115], where commuters having to choose between alternative roads take also into account in their routing decisions some of their neighbours' most recent waiting experience, the choice is here based only on the individual waiting history.

When the users follow rule $\tilde{\mathcal{R}}$ to individually choose their route through the network, simulation experiments (see Fig. 10.2) show that the emerging dynamics exhibit oscillatory behaviour of the queue contents (*i.e.* the congestion state of the two routes possesses an oscillating nature) that are reduced over time. This reduction leads to the convergence to a stationary state characterized by

$$\lim_{t \rightarrow \infty} Q_1(t) = \frac{\mu_1}{\mu_1 + \mu_2} \quad \text{and} \quad \lim_{t \rightarrow \infty} Q_2(t) = \frac{\mu_2}{\mu_1 + \mu_2}.$$

The softened oscillatory behaviour and the resulting convergence to this stationary state is a direct consequence of the larger HB structure defining the users' autonomous routing decisions. Indeed, this structure implies a manifest averaging over time that creates the resulting oscillations smoothing effect. The illustration given in Fig. 10.3 testifies this assertion. More precisely, Fig. 10.3 shows that, when the waiting time averages (*i.e.* \bar{W}_ζ^1 and \bar{W}_ζ^2 , $\zeta \in \mathbb{N}$) are computed by the users using only their last three experimented waiting times, the oscillations are not killed and the system remains indefinitely in a cyclo-stationary regime (*i.e.* the oscillations are sustained).

When the number of samples to compute the waiting time averages is limited, simulation experiments show that the value of ϵ governs the phase transition between regimes with convergence to stationary or cyclo-stationary states. Indeed, ϵ directly influences the length of the transient phase possibly leading to a purely stationary regime. More precisely, the relaxation time diminishes when ϵ is decreasing. When ϵ is small enough, the corresponding re-

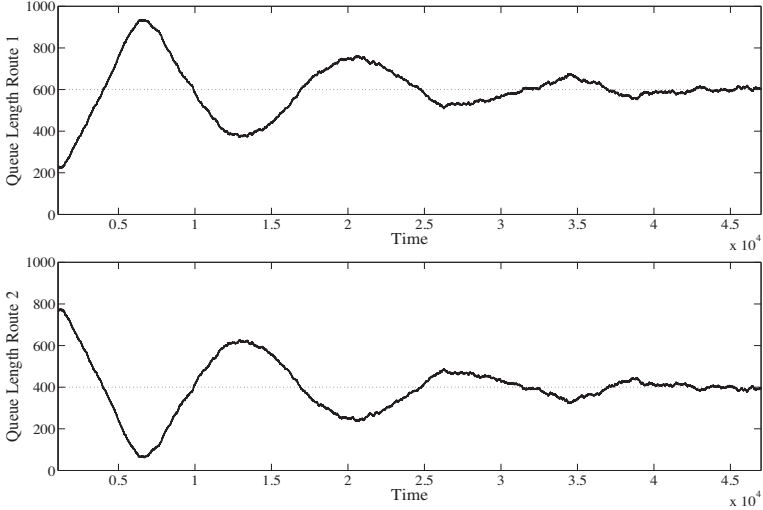


Fig. 10.2. Emerging dynamics in a transportation network with two parallel routes when the $N = 1000$ users apply routing rule $\tilde{\mathcal{R}}$. The service times of routes 1 and 2 are respectively uniformly distributed in $[1, \frac{7}{3}]$ ($\mu_1 = 0.6$) and $[2, 3]$ ($\mu_2 = 0.4$), $\alpha = 0.25$ and $\epsilon = 0.25$.

laxation time t_{relax} is small with report to the time window T_{HB} considered by the users to compute their waiting time averages. Hence, the oscillations are smoothed and the system is driven into a purely stationary state. On the other hand, for larger values of ϵ , t_{relax} gets larger than T_{HB} and the oscillations are consequently not averaged since they are in one sense not taken into account by the users within their routing decision mechanisms. In this case, a cyclo-stationary state emerges and remain indefinitely. These two different possible regimes are illustrated respectively in Figs. 10.4 (small ϵ) and 10.3 (large ϵ).

10.2 Smart Parts Driven Supply-Chains

Whatever the nature of the industrial activity, it invariably relies on the generation and the maintenance of flows of items and/or goods circulating in production, service and supply chains networks. This shows that, thermodynamically speaking, industrial activities can be viewed as off-equilibrium thermodynamic processes in which matter, information, energy and money feed complex entangled and topologically time-dependent networks. The underlying flow dynamics are generally subject to random fluctuations generated by dynamic environments (failures, prices and demand fluctuations, geo-political

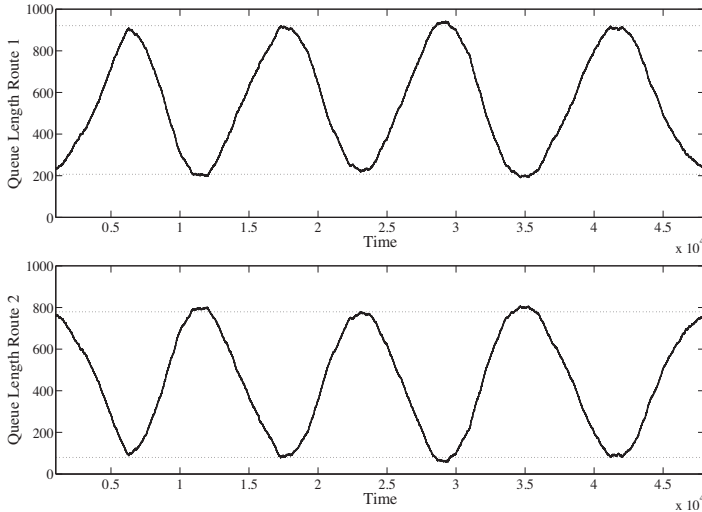


Fig. 10.3. Emerging dynamics in a transportation network with two parallel routes when the $N = 1000$ users apply routing rule $\tilde{\mathcal{R}}$, but with an **average computed only with the last 3 experimented waiting times**. The service times of routes 1 and 2 are respectively uniformly distributed in $[1, \frac{2}{3}]$ ($\mu_1 = 0.6$) and $[2, 3]$ ($\mu_2 = 0.4$), $\alpha = 0.25$ and $\epsilon = \mathbf{0.25}$.

issues,...). A significant part of the engineering work consists in controlling these complex flows to simultaneously maximize income and minimize environmental impacts, thus ensuring overall sustainability. By its very nature, flow control dynamics always raises a wealth of management problems which reflect the latest progresses in mathematical modelling, information processing and operational research. Due to their complexity and high specificity, one is tempted, at least at first sight, to assert that actually relevant flow control problems offer, despite their off-equilibrium thermodynamics nature, little perspectives for a synthetic approach necessary to formulate basic research actions. Yet such an assertion is denied by several recent contributions, including this thesis, that exhibit explicit and fruitful crossovers between basic science paradigms and management of production, service and supply chains networks. Specifically and in direct connection with the present work, focus should in the future be paid on the possible ways to implement efficient self-generated (fully decentralized) flow control algorithms that use emergent collective dynamics resulting from autonomous interacting agents. Indeed, it becomes now commonly acknowledged that only decentralized decision mechanisms will be able to permanently regulate with celerity and agility flows of raw materials, manufactured goods, food, services and energy required by non-stationary and random demand patterns generated by the economical global-

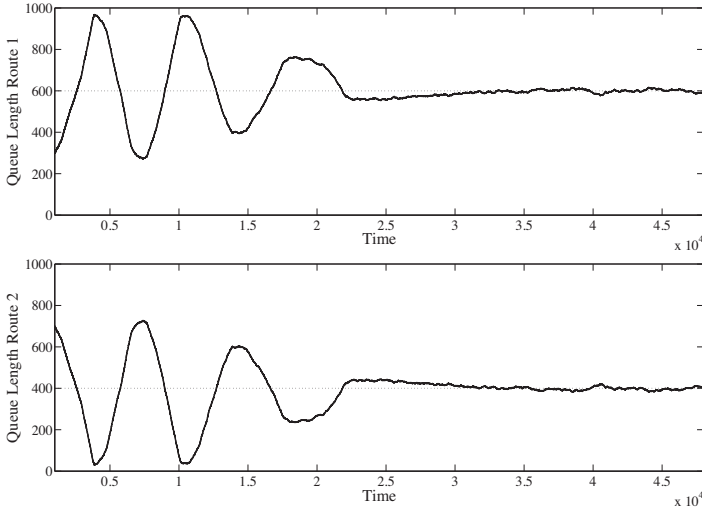


Fig. 10.4. Emerging dynamics in a transportation network with two parallel routes when the $N = 1000$ users apply routing rule $\tilde{\mathcal{R}}$, but with an **average computed only with the last 3 suffered waiting times**. The service times of routes 1 and 2 are respectively uniformly distributed in $[1, \frac{7}{3}]$ ($\mu_1 = 0.6$) and $[2, 3]$ ($\mu_2 = 0.4$), $\alpha = 0.25$ and $\epsilon = 0.05$.

ization. As previously unveiled in this thesis, the signature of living systems and societies, namely their unique capability to generate real-time adaptive behaviours with high resilience against random perturbations should definitely inspire new approaches to the problem of optimally regulating industrial flows.

In all supply chains there exist various forms of local feedback control, including the type of time delay satisfaction measure extensively studied in this thesis. Among the several noise sources affecting the dynamics of supply chains, non-steady flows typically arise in this context from coherent superpositions of tiny local fluctuations which, due to the numerous echelons composing the chains, give ultimately rise to very significant variations. Research actions should be directed to the implementation of the “smart parts” concept in simple supply chains to generate robust supply strategies with (1) self-reducing “bullwhip” effect and (2) efficient supplier selection for supply chains with branching topologies. Indeed, one should use and explore, via mathematical modelling and statistical physics tools, such a general self-organizing perspective with the ultimate goal of using these self-organizing mechanisms to construct resilient algorithms able to optimize the flow dynamics of supply chains networks. The numerous potentialities offered by “smart parts” regulated supply chains are reviewed in the purely conceptual, though

engrossing and enlightening contribution [127]. As a motivating illustration of this paradigm, consider the clear trade-off existing in our global economy between a steadily increasing demand for resources as well as for transportation capabilities and the antagonistic essential focus on the minimization of energy and environmental impacts. Complexity and randomness in the transportation processes within supply networks strongly reduce the efficiency of centralized shipping management. In this context, “smart marine containers” equipped with real-time positioning capabilities (via GPS technology) offer extremely appealing potentialities to built fully decentralized logistics solutions with short-time reactivity and low energy consumption characteristics. Being also equipped with RFID chips indicating their final destination, the containers effectively become “smart parts” able to autonomously plan, in real-time, their best routing from their current position to the final customer location.

To define future research directions, one should start here from the basic aspects of supply chains dynamics theory and use the Daganzo’s dedicated formalism, [33]. This formulation relies on systems of coupled, nonlinear, finite-difference equations to describe the time evolution of the order flows and the inventory contents. Thanks to this mathematical modelling, the dynamical origins of the so-called “bullwhip” effect, ubiquitous in supply chains dynamics, are clearly identified. Let us here recall that the “bullwhip” phenomenon expresses the strong tendency in supply chains that small fluctuations at the customers’ flow of demands are amplified along the echelons of the chain and ultimately give rise to huge demand variations in the upper echelons. Using Daganzo’s formalism, the emergence of the “bullwhip” effect can be explained via a dynamical systems stability theory. More particularly, various non-equivalent supply chains stability concepts (*i.e.* stability on the small and stability on the large) are introduced in [33] to exhibit how the “bullwhip” phenomenon actually occurs. Linear stability around steady order flows and harmonic analysis reveals that the “bullwhip” amplification effects are due to the effective presence of delays in the dynamics. This allows to conclude that “bullwhip” phenomena could be reduced by decentralized commitment policies in which the order sizes to be passed in future should be partly anticipated. The length of the commitment time interval and its connection with “bullwhip” reduction would then have to be studied, even for nonlinear supply policies, via kinematic waves analysis.

Contrary to common practice where intelligence is spread over the various suppliers of the chain, one should investigate adaptive supply algorithms which, due to “smart parts” control mechanisms, would lead to flows with attenuated “bullwhip” behaviours and optimal supplier selection.

(1) *Attenuation of the “bullwhip” effect.*

To practically proceed, attention should first be paid to the so-called autonomous supply chain models characterized by interactions restricted to the nearest successive echelons. In discrete time ($n \in \mathbb{N}$), the supply chains dynamics take the form of a coupled map lattice:

$$\begin{cases} Q_{j,n} &= H_j(K_{j,n}, Q_{j-1,n-1}, Q_{j-1,n-2}, \dots, Q_{j-1,n-B}) - K_{j,n}, \quad n \in \mathbb{N}, \\ K_{j,n+1} &= H_j(K_{j,n}, Q_{j-1,n-1}, Q_{j-1,n-2}, \dots, Q_{j-1,n-B}) - Q_{j-1,n}, \quad n \in \mathbb{N}, \end{cases}$$

where $Q_{j,n}$ and $K_{j,n}$ are respectively the order size and the inventory content of the echelon number $j = 0, 1, 2, \dots, N$ (note that 0 is the customer level in this formalism), at the discrete time $n \in \mathbb{N}$. The functional $H_j(K_{j,n}, Q_{j-1,n-1}, Q_{j-1,n-2}, \dots, Q_{j-1,n-B})$ characterizes the dynamics for each echelon. More precisely, this is the kernel of the order policy and it is usually given a priori. Future research actions should propose to allow the set of kernels H_j to be adaptive in time, this thanks to “smart parts” characteristics of the items in circulation. Such an idea should first be investigated by considering the elementary (uniform) (s, S) -order-point-policy with a kernel $H_j^{(s,S)}$ characterized by:

$$H_j^{(s,S)} \underbrace{=}_{\text{uniformity}} H^{(s,S)} = \begin{cases} 0 & \text{if } K_{j,n} > s, \\ S & \text{if } K_{j,n} \leq s, \quad \forall j = 1, 2, \dots, N. \end{cases}$$

In other words the discontinuous functional $H_j^{(s,S)}$ is actually, as long as the inventory contents fulfil $K_{j,n} > s$, a *lot-for-lot* strategy for which no “bullwhip” effect develops. At that point, it would still remain to optimally select the threshold parameters (s, S) that characterize the implemented supply policy. To this aim, we would propose to implement a “smart orders” based mechanism able, for non-stationary demand patterns, to permanently select the required optimal (s, S) thresholds (optimality meaning here to minimize the apparition of potential “bullwhip” phenomena).

(2) *Self-routing supply paths.*

It is common in real situations to encounter supply networks with complex topologies. At the branching nodes of such supply chains, orders have to be efficiently dispatched to the available supply paths of the network. To cope in real-time with ubiquitous non-stationarity affecting the supply nodes, future research actions should investigate how “smart parts” concepts could be implemented to yield self-routing and individualized optimal supply paths - for example paths might be either quick and expensive or cheap and slow. To this aim, decentralized feedback mechanisms as those developed along this thesis or pheromone approaches should be considered.

10.3 Contributions of Chapter 10

- We present a simple yet realistic model that establishes a bridge between the models introduced in this work and transportation networks. This new stylized transportation model helps to understand some of the emerging phenomena that can be observed in real traffic systems. The consideration of more refined users' history-based mechanisms, which are common in transportation systems, leaves the door open for the observation of oscillatory phenomena, but the resulting oscillations reveal themselves to be softened over time.
- We give some prospective ideas for an investigation on how the concepts developed in this thesis could be transferred to the domain of supply chains.

Conclusion and Perspectives

Conclusion and Perspectives

The main stream in queueing network (QN) theory is to consider the circulation of classes of items sharing all the same set of permanent attributes. In this context, a single item is fully representative of all the members belonging to its class. The present thesis definitely differs from such a classical point of view and considers the flow dynamics in networks where each circulating item is an autonomous agent able to adapt its routing according to historical data monitored during its past journey through the network. Accordingly, a single circulating item is not a copy of the others and the resulting dynamics are not covered by the ordinary tools of QN theory. As studied in the present work, global dynamics of interacting autonomous agents within QNs explicitly belong to the vast realm of complex systems, for which collective properties emerge from the individual “intelligence” endowed to each agent. Indeed, for such systems, emergent macroscopic properties are more than the sum of their microscopic components. As stipulated in the present work, the underlying history-based (HB) routing decision mechanisms considered along this thesis violate the basic hypothesis of classical queueing models. Accordingly, emerging spatio-temporal flow patterns due to non-Markovian routing decisions individually taken by circulating agents remain, despite their truly strong interdisciplinary integration, an almost unexplored topic in the available literature on QNs. In this regard, QN theory involving the circulation of autonomous agents can definitely be viewed as a new topic in itself. While the generic character and the synergetic modelling potential offered by multi-agent systems dynamics have already been abundantly explored in basic sciences (physics, chemistry, biology) and in social sciences (economics, finance, psychology, car traffic), this thesis opens some doors towards the investigation of such concepts in production, supply chains and service QNs. With a view to applications, the highly flexible modern production and supply networks, able to satisfy extreme ranges of customized products, rely more and more on decentralized mechanisms able to self-organize the material flows visiting intricate networks topologies. In this context, the increasing availability of RFID technology offers the possibility for a wide implementation of such

local intelligence to circulating items in various logistics networks. Such “intelligent” devices should be able, in real-time, to select autonomously, according to ad-hoc historical individually measured data and real-time observations, the best possible routing alternative through the network. No matter the application domain, we will probably face in a near future mixed situations in which part of the dynamics will still obey to a classical centralized control and part of it will be left to self-organization. This therefore raises a basic question: how much to control and how much to let spontaneously emerge? With the aim to contribute to give an answer to this basic issue, the present work studies several possible self-organized structures emerging in stylized yet paradigmatic queueing systems. The ultimate goal of determining an efficient compromise between pure interventionism (due to centralized controls) and self-organization (due the swarm intelligence of the agents) is not confined to issues in production and supply chains networks but answers a fundamental question in the basic management of any logistics system.

The catalogue of techniques exposed along this thesis, which has ultimately yielded the elaboration of several classes of solvable complex adaptive systems, clearly proves that analytical considerations can be handled even in presence of multi-agent systems with explicit non-Markovian dynamics. To the best of our knowledge, there barely exist no available similar mathematical modelling studies where the potentiality of “smart parts” concept for production, service and supply chains networks is analyzed. In this regard, the simple and didactic models provided in this thesis are more than highly welcome, in particular for teaching activities. While possessing a wide potential for applications and hence strong industrial relevance, the inherent tractability of our stylized models allows at the same time for a deep understanding of the emerging self-organized phenomena. Accordingly, the collection of solvable models presented in this thesis should retain the attention of a multidisciplinary audience ranging from scientists of the community of complex systems to production managers. At the very beginning of this thesis, the ultimate (but naive) goal was to develop an approach that was at the same time theoretical (providing new analytical results) and pragmatic (paying attention to real-life problems). One cannot but notice that such a very ambitious perspective is only partially achieved. Practitioners would probably say that our stylized models would have to be refined to be implemented in real situations and theoreticians would probably emphasize that no universal rules about non-Markovian processes in QNs are actually available in this thesis. This is obviously true and further future research has definitely to be carried out towards these two antagonistic directions. This thesis however, thanks to its inherent compromise between theoretical and practical considerations, opens wide the doors, gives some prospective insights and unveils some potentialities towards both of these future research directions.

The intrinsic complexity relative to the presence of agent-based non-Markovian

processes in QNs leaves at first sight little hope for the establishment of general analytical considerations and it seems obviously difficult to yield the elaboration of a dedicated universal theoretical framework. While this work does not provide in fine such type of general rules for QNs roamed by agents with HB decision capabilities (and by extension for the treatment of non-Markovian processes within QNs), vain attempts to establish embryos of a universal theory have been conducted first *(i)* by using the theory of delay differential equations then *(ii)* by considering Polya urn models (random processes with reinforcement). While these two modelling techniques obviously possess several aspects that would be in favour of an extension to the present framework, our guess is that attention should be rather paid to the establishment of an interdisciplinary dedicated theory that might include also techniques from the dynamical systems area. From a global point of view, besides the specific perspectives already drawn in the concluding remarks that close the different chapters composing this thesis (which notably propose various refinements for our models), the two following general research directions should now be considered, namely *(1)* to refine the proposed stylized models and to extend them to fit to real practical situations and *(2)* to eventually develop a general theoretical framework for the analysis of QNs with non-Markovian autonomous routing mechanisms. In order to maximize the chances of success, these two research directions should probably be mostly uncoupled in a short-term horizon.

- (1)* On one hand, one should follow a pragmatic and problem centered approach the aim of which would be to come up with practical solutions to specific logistics situations. To cope with real cases, new parameters would be progressively added to the stylized models presented here and more complex network topologies would be considered, while still trying to preserve the tractability of the models.
- (2)* General theoretical concepts would obviously help to gain deeper and more generic understanding of the processes driving the dynamics of non-Markovian QNs. Such detailed understanding of these processes and the corresponding theoretical knowledge could possibly be transferred, in a later stage of the research process, to the industrial community to potentially implement new innovative practical solutions.

Part VII

Appendix

Appendix Chapter 2 - Typical Delayed Dynamics

Summary. *In this appendix, we give an explicit illustration of how the existence of time delays in the dynamics of queueing systems might produce oscillations phenomena. This illustrative example, observed in a real context and hence directly derived from a typical service management situation, helps to get an intuitive understanding of the inherent mechanisms leading to such a periodic behaviour.*

12.1 Introduction

As briefly exposed in Chapter 2, the presence of time delays (or more generally non-local in time effects) into dynamical systems is far from being innocent and opens wide the door for the birth of various complex phenomena, [51]. Delayed dynamics might in particular exhibit regimes characterized by oscillatory behaviours - this can be illustrated even in the simple linear relaxation

$$\frac{d}{dt}X(t) = -X(t), \quad X(0) = X_0,$$

which admits the solution $X(t) = X_0 e^{-t}$. Indeed, this simple time evolution rule is able to produce oscillations, when a delay $\tau = \frac{\pi}{2}$ is introduced (the ad-hoc initial function on $t \in [-\frac{\pi}{2}, 0]$ has obviously to be implemented):

$$\frac{d}{dt}X(t) = -X(t - \tau); \quad X(0) = X_0$$

admits the solution $X(t) = X_0 \cos(t)$ when $\tau = \frac{\pi}{2}$.

After having observed time oscillations for this simple linear delayed equation, one clearly suspects that queueing networks (QN's) with intrinsic time delay features might also give rise to oscillatory dynamics. In the next section, we corroborate this assertion with a practical illustration arising in the field of service management.

12.2 Crowded Day at the Pantheon

The story takes place on January 19th 2007. On this day, during the “Homage aux Justes de la Nation”, the entrance to the Pantheon in Paris was exceptionally free. The costless admission attracting an unusually large incom-



Fig. 12.1. Pantheon, Paris. *Copyright: Jean-Christophe Benoist.*

ing flow of visitors, the access to the building had to be specially regulated on this day. The regulation process that was implemented in this particular situation is described in the next paragraph.

After having possibly lined up for some time outside, the visitors are allowed to enter into the monument. Once inside, they first stay on the ground floor of the dome to visit an exhibition paying tribute to the “Justes de la Nation”. Then, the majority of the visitors proceed with their visit, whose second step is to go down to the famous crypt, where are buried numerous *great men* (Curie, Monge, Rousseau and Voltaire to mention only a few). Due to the relative shortness of the crypt, the flow of visitors has to be carefully regulated at its entrance (which is actually located inside the dome, opposite to the main entrance). The aim of this regulation is to limit the number of visitors being simultaneously in the crypt. As a consequence of this regulating process, a queue of waiting visitors might appear in front of the access to the

crypt. The size of this queue has to be controlled since it is located inside the dome and thus encroaches upon the exhibition area. To take care of that, the management strategy of the person in charge, January 19th 2007, was:

“When the queue at the entrance of the crypt reaches a given threshold (which is represented by a particular pillar in the building), tell the colleague in charge at the entrance of the building (instantaneously, via walkie-talkies) to stop letting visitors enter into the dome. Visitors will be allowed to enter into the building only when the queue will be again below the given threshold.”

In practice, it was possible to observe that this strategy leads to an oscillatory behaviour of the content of the queue located in front of the crypt. In order to understand the source of these oscillations, we now present a queueing model corresponding to this situation.

Following classical queueing concepts (see [93] among others), we consider the particular queueing network \mathcal{N} sketched in Fig. 12.2, which is composed

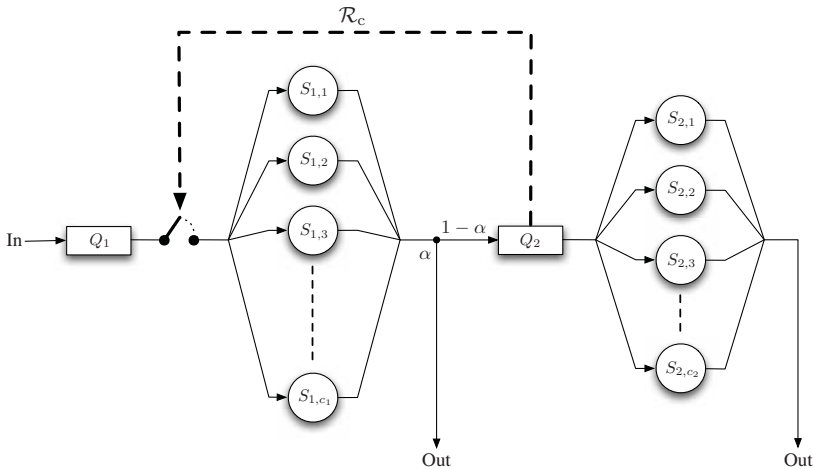


Fig. 12.2. Queueing network \mathcal{N} , a mathematical representation of the queueing processes characterizing the flow of visitors in the Pantheon in Paris.

of two successive service facilities, S_1 and S_2 , set in line. The first service facility characterizes the exhibition located on the ground floor and the second one the visit of the crypt. While a part α of the visitors leave the system after S_1 , the remaining ones follow their way to S_2 (this choice is represented by a Bernoulli random variable). $Q_1(t)$ and $Q_2(t)$ respectively denote the queue contents in front of the first and second service facilities. S_1 (resp. S_2)

possesses c_1 (resp. c_2) parallel service channels, which corresponds to the maximum number of visitors allowed in the exhibition area (resp. in the crypt). In the present application, c_1 reveals itself to be greater than c_2 . The visitors arrive to the first service facility (*i.e.* to the main entrance of the building) following a general renewal process with mean inter-arrival time $\frac{1}{\lambda}$. The service times of S_1 (resp. S_2) are characterized by generally distributed i.i.d. random variables with mean $\frac{1}{\mu_1}$ (resp. $\frac{1}{\mu_2}$). These service times denote the length of the successive visits in the exhibition area and in the crypt. Finally, while the capacity of the first queue (*i.e.* the one outside the building) is assumed to be infinite, the capacity of the second one (*i.e.* the one located in front of the crypt) is limited to K . To resume, the first service facility is represented by a $G/G/c_1/\infty$ queueing model and the second one with a $\cdot/G/c_2/K$ model.

The system management obeys to the following control rule:

$$\mathcal{R}_c = \begin{cases} \text{Admission to the first service} & \text{when } Q_2(t) < N_c, \\ \quad \text{facility } S_1 \text{ is open} & \\ \text{Admission to the first service} & \text{when } Q_2(t) \geq N_c, \\ \quad \text{facility } S_1 \text{ is closed} & \end{cases}$$

where $N_c < K$ is a fixed threshold used to handle congestion in front of the crypt. The transmission of this management information between the two different queues is instantaneous.

As illustrated in Fig. 12.3, when we simulate the dynamics of the queue-

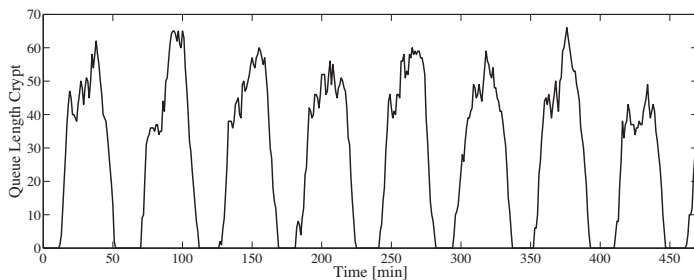


Fig. 12.3. Queue content dynamics in front of the second service facility (*i.e.* the crypt) when $c_1 = 200$, $c_2 = 100$, $\alpha = 0.3$, $N_c = 30$ and the service times of S_1 and S_2 are uniformly distributed in $[10, 30]$ (*i.e.* the visitors spend on average 20 minutes in the exhibition area, then in the crypt).

ing system \mathcal{N} , we clearly observe queue content oscillations in front of the second service facility (*i.e.* the crypt). The period of these oscillations, a little bit less than one hour (see Fig. 12.4), corresponds to what we could observe

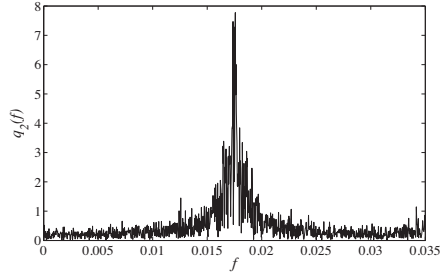


Fig. 12.4. Spectrum corresponding to the queue content dynamics illustrated in Fig. 12.3 (*i.e.* the queue located in front of the crypt).

on site, January 19th 2007. The oscillatory behaviour, which is entirely due to the presence of a delay effect in the dynamics, can be easily understood as follows. Following rule \mathcal{R}_c , when the manager observes that $Q_2(t) = N_c$, he/she warns of a congestion in front of the crypt and the admissions into the building are immediately stopped. However, $Q_2(t)$ will continue to be fed for some time with the visitors currently in the first service facility (*i.e.* those still visiting the exhibition located on the ground floor). There exists hence a delay between the time at which the managerial decision to close the admissions is taken and the time at which this decision really takes effect (*i.e.* when no visitors join $Q_2(t)$ anymore). During this time delay, $Q_2(t)$ continues to increase and will obviously exceed the fixed critical level N_c . Consequently, from a managerial point of view, it is important that N_c is chosen smaller than the effective waiting room capacity K , that should in any case not be surpassed.

An extension of this particular queueing system with intrinsic time delay effects and an application to after sales processes in the watchmaker industry is provided in [37].

Appendix Chapter 7

In this appendix, we describe the class of functions that could be considered to approximate the market sharing dynamics considered in Chapter 7 (see Section 7.3) when one wants to observe a noise-induced phase transition between uni- and bimodal stationary probability density functions (PDFs). More precisely, we are interested in the functions for which there exists a particular value of γ (call it $\bar{\gamma}$) such that the stationary PDF $P_s(y)$ given by Eq. (7.19) satisfies:

$$P_s(y) = \frac{1}{2\Delta}, \quad \forall y \in [-\Delta, +\Delta]. \quad (13.1)$$

The condition given in Eq. (13.1) is necessary in order to be able to observe a phase transition between Hotelling-like (unimodal stationary PDF) and deadline type regimes (bimodal stationary PDF). The unique class of functions fulfilling this condition is given in the following lemma.

Lemma 13.1. *The functions $f : \mathbb{R} \rightarrow [-1, +1]$ of the form:*

$$(f^{-1})' \left(\frac{y}{\Delta} \right) = \frac{1}{\alpha \left(\frac{y}{\Delta} \right)^2 + \beta}, \quad \alpha < 0, \quad \beta \geq 0, \quad -\frac{\beta}{\alpha} \leq 1, \quad (13.2)$$

are the only ones for which there exists a particular value $\bar{\gamma} = -\frac{A}{\alpha V^2} > 0$ such that:

$$\begin{aligned} P_s(y) &= \mathcal{N} \left((f^{-1})' \left(\frac{y}{\Delta} \right) \exp \left\{ -\frac{2A}{\bar{\gamma} \Delta^2 V^2} \int^y u (f^{-1})' \left(\frac{u}{\Delta} \right) du \right\} \right) \\ &= \frac{1}{2\Delta}, \quad \forall y \in [-\Delta, +\Delta]. \end{aligned} \quad (13.3)$$

For such functions f , $P_s(y)$ exhibits a single mode when $\gamma < \bar{\gamma}$ and two modes when $\gamma > \bar{\gamma}$.

Proof. First, note that Eq. (13.3) is satisfied iff $P'_s(y) = 0, \forall y \in [-\Delta, +\Delta]$. Fixing $\bar{\gamma}$ in Eq. (7.19) and differentiating, we obtain the following condition on f :

$$(f^{-1})'' \left(\frac{y}{\Delta} \right) - \frac{2A}{\bar{\gamma}\Delta V^2} y \left((f^{-1})' \left(\frac{y}{\Delta} \right) \right)^2 = 0, \quad \forall y \in [-\Delta, +\Delta]. \quad (13.4)$$

The unique set of solutions of this ordinary differential equation is given by:

$$(f^{-1})' \left(\frac{y}{\Delta} \right) = \frac{1}{\alpha \left(\frac{y}{\Delta} \right)^2 + \beta}, \quad \beta \geq 0, \quad \forall y \in [-\Delta, +\Delta], \quad (13.5)$$

where $\alpha = -\frac{A}{\bar{\gamma}V^2} < 0$. Furthermore, in order that $f \in [-1, +1]$, we suppose that $-\frac{\beta}{\alpha} \leq 1$. Inserting Eq. (13.5) into Eq. (7.19) and using the fact that $\alpha = -\frac{A}{\bar{\gamma}V^2}$, we find that the stationary probability density function $P_s(y)$ is equal, $\forall \gamma > 0$, to

$$\begin{aligned} P_s(y) &= \mathcal{N} \frac{1}{\alpha \left(\frac{y}{\Delta} \right)^2 + \beta} \exp \left\{ -\frac{2A}{\gamma\Delta^2 V^2} \int^y u \frac{1}{\alpha \left(\frac{u}{\Delta} \right)^2 + \beta} du \right\} \\ &= \mathcal{N} \left(-\frac{A}{\bar{\gamma}V^2} \left(\frac{y}{\Delta} \right)^2 + \beta \right)^{\frac{\bar{\gamma}}{\gamma}-1}, \quad \forall y \in [-\Delta, +\Delta] \end{aligned} \quad (13.6)$$

for functions f satisfying Eq. (13.5). In regard to Eq. (13.6), $P_s(y)$ is hence either a unimodal or a bimodal distribution for $\gamma \neq \bar{\gamma}$. \square

Note that for the function \tanh ($\alpha = -1$ and $\beta = 1$), $\bar{\gamma} = \frac{A}{V^2}$ and Eq. (13.6) thus becomes:

$$P_s(y) = \mathcal{N} \left(1 - \left(\frac{y}{\Delta} \right)^2 \right)^{\frac{A}{\gamma V^2}-1}, \quad \forall y \in [-\Delta, +\Delta],$$

which is in perfect agreement with Eq. (7.23).

Integrating Eq. (13.2) gives the following equivalent condition on f :

$$f^{-1}(y) = -\frac{1}{\sqrt{-\alpha\beta}} \text{Arctanh} \left(\sqrt{\frac{-\alpha}{\beta}} y \right) + C, \quad C \in \mathbb{R}. \quad (13.7)$$

Then, determining the inverse of the functions satisfying Eq. (13.7), we find that the following monotonic increasing functions:

$$f(x) = \sqrt{\frac{\beta}{-\alpha}} \left(\frac{K e^{\sqrt{-\alpha\beta}x} - e^{-\sqrt{-\alpha\beta}x}}{K e^{\sqrt{-\alpha\beta}x} + e^{-\sqrt{-\alpha\beta}x}} \right), \quad K > 0, \quad (13.8)$$

satisfy Eq. (13.2) and hence, by Lemma 13.1, compose the class of functions we are interested in. Observe that the function \tanh well belongs to this class

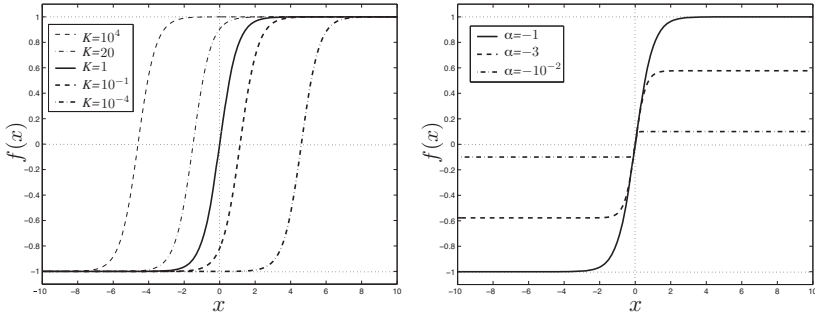


Fig. 13.1. Several functions belonging to the class given by Eq. (13.8). *Left:* $\alpha = -1$, $\beta = 1$. *Right:* $\beta = 1$, $K = 1$.

($\alpha = -1$, $\beta = 1$ and $K = 1$). Several instances of functions belonging to the class characterized by Eq. (13.8) are drawn in Fig. 13.1.

In addition to the large scope of dynamics covered by the possible choice of γ in Eq. (7.11), the diversity of the functions resulting from the arbitrary choice of the 3-tuple (α, β, K) in Eq. (13.8) actually allows us to cover the dynamics induced by a still wider range of customers' decision policy. Although we have considered odd functions that are perfectly suitable for the purely symmetric configurations of Section 7.3 as well as for the asymmetric situations considered in Section 7.4, note that it would also be possible to model static non-symmetric configurations by playing on the value of K (see Fig. 13.1.*Left*). However, as discussed in Section 7.4.2, this kind of static asymmetry drives the system to the same stationary behaviour as in the symmetric case (indeed, $(f^{-1})'$ does not depend on K) in Eq. (13.8). Furthermore, the class of functions given by Eq. (13.8) would also fit to model situations where the boundary point is confined to a smaller interval $\mathcal{I} \subset \Omega$ (see Fig. 13.1.*Right*). This would correspond to configurations where a part of the customers is automatically bound to a service provider, no matter the current state of the queue at this server. For example, one could consider that “competition” occurs only in the interval $\mathcal{I} = [x_1, x_2]$ between the two service providers and that the customers in $[-\Delta, x_1]$ (resp. $[x_2, +\Delta]$) always choose S_1 (resp. S_2).

References

1. H. Abdel-Jaber, M. Woodward, F. Thabtah, and A. Abu-Ali. Performance Evaluation for DRED Discrete-Time Queueing Network Analytical Model. *Journal of Network and Computer Applications*, 31:750–770, 2008.
2. H. Afimeimounga, W. Solomon, and I. Ziedins. The Downs-Thomson Paradox: Existence, Uniqueness and Stability of User Equilibria. *Queueing Systems*, 49:321–334, 2005.
3. V. Alfi, M. Cristelli, L. Pietronero, and A. Zaccaria. Minimal Agent Based Model for Financial Markets I: Origin and Self-Organization of Stylized Facts. *The European Physical Journal B*, 67:385–397, 2009.
4. H. M. Amann and L. Tesfatsion. *Handbook of Computational Economics (Volume 2)*. Elsevier, 2006.
5. D. Armbruster, C. de Beer, M. Freitag, T. Jagalski, and C. Ringhofer. Autonomous Control of Production Networks Using a Pheromone Approach. *Physica A*, 363:104–114, 2006.
6. D. Armbruster and E. S. Gel. Bucket Brigades Revisited: Are They Always Effective? *European Journal of Operational Research*, 172:213–229, 2006.
7. D. Armbruster, E. S. Gel, and J. Murakami. Bucket Brigades with Worker Learning. *European Journal of Operational Research*, 176:264–274, 2007.
8. L. Arnold, H. Crauel, and V. Wihstutz. Stabilization of Linear Systems by Noise. *SIAM Journal on Control and Optimization*, 21(3):451–461, 1983.
9. R. J. Aumann. Game theory. In: *The New Palgrave, A Dictionary of Economics (Volume 2)*, pp. 460–482, J. Eatwell, M. Milgate and P. Newman (Eds.), Macmillan, London and Basingstoke, 1987.
10. J. Aweya, M. Ouellette, and D. Y. Montuno. A Control Theoretic Approach to Active Queue Management. *Computer Networks*, 36:203–235, 2001.
11. P. Bak. *How Nature Works: The Science of Self-Organized Criticality*. Copernicus, New York, NY, 1996.
12. R. O. Baldwin, N. J. Davis IV, J. E. Kobza, and S. F. Midkiff. Real-Time Queueing Theory: A Tutorial Presentation with an Admission Control Application. *Queueing Systems*, 35:1–21, 2000.
13. J. J. Bartholdi III and D. D. Eisenstein. A Production Line That Balances Itself. *Operations Research*, 44(1):21–34, 1996.
14. J. J. Bartholdi III, D. D. Eisenstein, and R. D. Foley. Performance of Bucket Brigades when Work is Stochastic. *Operations Research*, 49(5):710–719, 2001.

15. J. J. Bartholdi III, D. D. Eisenstein, and Y. F. Lim. Bucket Brigades on In-Tree Assembly Networks. *European Journal of Operational Research*, 168:870–879, 2006.
16. F. Baskett, K. M. Chandy, R. R. Muntz, and F. G. Palacios. Open, Closed and Mixed Networks of Queues with Different Classes of Customers. *Journal of the ACM*, 22(2):248–260, 1975.
17. F. Bielen and N. Demoulin. Waiting Time Influence on the Satisfaction-Loyalty Relationship in Services. *Managing Service Quality*, 17(2):174–193, 2007.
18. E. Bonabeau. Agent-Based Modeling Methods and Techniques for Simulating Human Systems. *Proceedings of The National Academy of Sciences USA*, 99 (supplement 3):7280–7287, 2002.
19. F. Bonomi and A. Kumar. Adaptive Optimal Load Balancing in a Nonhomogeneous Multiserver System with a Central Job Scheduler. *IEEE Transactions on Computers*, 39(10):1232–1250, 1990.
20. R. Boucekkine, D. de la Croix, and O. Licandro. Modelling Vintage Structures with DDEs: Principles and Applications. *Mathematical Population Studies*, 11(3):151–179, 2004.
21. D. Breitgand, R. Cohen, A. Nahir, and D. Raz. Cost Aware Adaptive Load Sharing. *Lecture Notes in Computer Science*, 4725:208–224, 2007.
22. D. Breitgand, R. Cohen, A. Nahir, and D. Raz. On Fully Distributed Adaptive Load Balancing. *Lecture Notes in Computer Science*, 4785:74–85, 2007.
23. L. Breuer. *From Markov Jump Processes to Spatial Queues*. Kluwer Academic Publishers, Dordrecht, 2003.
24. G. P. Cachon and P. T. Harker. Service Competition, Outsourcing and Co-Production in a Queueing Game. *Preliminary Report*, <http://knowledge.wharton.upenn.edu/papers/889.pdf>, 1999.
25. G. P. Cachon and P. T. Harker. Competition and Outsourcing with Scale Economies. *Management Science*, 48(10):1314–1333, 2002.
26. P. Cahuc and S. Carcillo. The Shortcomings of a Partial Release of Employment Protection Laws: The Case of the 2005 French Reform. *IMF Working Paper*, 06/301 (January), 2007.
27. B. Calvert. The Downs-Thomson Effect in a Markov Process. *Probability in the Engineering and Informational Sciences*, 11:327–340, 1997.
28. E. Çinlar. *Introduction to Stochastic Processes*. Prentice-Hall International Editions, Englewood Cliffs, 1975.
29. H. Chen and D. D. Yao. *Fundamentals of Queueing Networks*. Springer-Verlag, 2001.
30. T. Y. Choi, K. J. Dooley, and M. Rungtusanatham. Supply Networks and Complex Adaptive Systems: Control Versus Emergence. *Journal of Operations Management*, 19(3):351–366, 2001.
31. H. K. H. Chow, K. L. Choy, and W. B. Lee. A Dynamic Logistics Process Knowledge-Based System - An RFID Multi-Agent Approach. *Knowledge-Based Systems*, 20:357–372, 2007.
32. F. Chu and X. L. Xie. Deadlock Analysis of Petri Nets Using Siphons and Mathematical Programming. *IEEE Transactions on Robotics and Automation*, 13(6):793–804, 1997.
33. C. F. Daganzo. *A Theory of Supply Chains*. Springer, Lecture Notes in Economics and Mathematical Systems, 2003.
34. C. D’Aspremont, J. J. Gabszewicz, and J. F. Thisse. On Hotelling’s “Stability in Competition”. *Econometrica*, 47(5):1145–1150, 1979.

35. G. R. D'Avignon and R. L. Disney. Single-Server Queues with State-Dependent Feedback. *INFOR Journal*, 14(1):71–85, 1976.
36. L. C. Davis. Modifications of the Optimal Velocity Traffic Model to Include Delay due to Driver Reaction Time. *Physica A*, 319:557–567, 2003.
37. F. Dias. Gestion des flux dans le service après-vente en milieu horloger. *Semester Project*, EPFL, 2008.
38. B. Doytchinov, J. Lehoczy, and S. Shreve. Real-Time Queues in Heavy Traffic with Earliest-Deadline-First Queue Discipline. *The Annals of Applied Probability*, 11(2):332–378, 2001.
39. R. D. Driver. *Ordinary and Delay Differential Equations*. Springer, Applied Mathematics Series 20, 1977.
40. P. Erdi. *Complexity Explained*. Springer, 2008.
41. T. Erneux. *Applied Delay Differential Equations*. Springer, 2009.
42. A. Erramilli and L. J. Forys. Oscillations and Chaos in a Flow Model of Switching System. *IEEE Journal on Selected Area in Communications*, 9(2):171–178, 1991.
43. R. Filliger and M.-O. Hongler. Syphon Dynamics - A Soluble Model of Multi-Agents Cooperative Behavior. *Europhysics Letters*, 70(3):285–291, 2005.
44. S. Floyd and V. Jacobson. Random Early Detection Gateways for Congestion Avoidance. *IEEE/ACM Transactions on Networking*, 1(4):397–413, 1993.
45. O. Gallay and M.-O. Hongler. Cooperative Dynamics of Loyal Customers in Queueing Networks. *Journal of Systems Science and Systems Engineering*, 17(2):241–254, 2008.
46. O. Gallay and M.-O. Hongler. Weariness and Loyalty Loss in Recurrent Service Models. In: *Proceedings of MOSIM'08*, volume 2, pp. 1011–1018, Paris, France, 2008.
47. O. Gallay and M.-O. Hongler. Circulation of Autonomous Agents in Production and Service Networks. *International Journal of Production Economics*, 120:378–388, 2009.
48. E. Gelenbe and G. Pujolle. *Introduction to Queueing Networks*. John Wiley and Sons, 1987.
49. P.-P. Grassé. *Termitologia*. Masson, Paris, 1984.
50. W. S. C. Gurney, S. P. Blythe, and R. M. Nisbet. Nicholson's Blowflies Revisited. *Nature*, 287:17–21, 1980.
51. I. Györi and G. Lads. *Oscillation Theory of Delay Differential Equations*. Clarendon Press, 1991.
52. E. Hashem. Analysis of Random Drop for Gateway Congestion Control. *Report TR-465*, MIT Laboratory of Computer Science, Boston, 1989.
53. R. Z. Has'minskii. *Stochastic Stability of Differential Equations*. Sijthoff & Noordhoff, Alphen aan den Rijn, 1980.
54. R. Hassin. On the Advantage of Being the First Server. *Management Science*, 42(4):618–623, 1996.
55. R. Hassin and M Haviv. *To Queue or not to Queue*. Kluwer Academic Publishers, 2003.
56. C. Haxholdt, E. R. Larsen, and A. van Ackere. Mode Locking and Chaos in a Deterministic Queueing Model with Feedback. *Management Science*, 49(6):816–830, 2003.
57. R. J. Henry, Z. N. Masoud, A. H. Nayfeh, and D. T. Mook. Cargo Pendulation Reduction on Ship-Mounted Cranes via Boom-Lu Angle Actuation. *Journal of Vibration Control*, 7:1253–1264, 2001.

58. J. H. Holland. *Hidden Order*. Perseus Books, Cambridge, MA, 1995.
59. J. H. Holland. Complex Adaptive Systems and Spontaneous Emergence. In: *Complexity and Industrial Clusters*, pp. 25-34, A. Q. Curzio and M. Fortis (Eds.), Physica, Heidelberg, 2002.
60. M.-O. Hongler, N. Cheikhrouhou, and R. Glardon. An Elementary Model for Customer Fidelity. In: *Proceedings of MOSIM'04*, volume 2, pp. 899-906, Lavoisier Editions, Nantes, France, 2004.
61. W. Horsthemke and R. Lefever. *Noise-Induced Transitions (Second Edition)*. Springer-Verlag, 2006.
62. H. Hotelling. Stability in Competition. *The Economic Journal*, 39(153):41-57, 1929.
63. M. Hülsmann, J. Grapp, and Y. Li. Strategic Adaptivity in Global Supply Chains - Competitive Advantage by Autonomous Cooperation. *International Journal of Production Economics*, 114:14-26, 2008.
64. G. E. Hutchinson. Circular Causal Systems in Ecology. *Annals of the New York Academy of Sciences*, 50:221-246, 1948.
65. J. R. Jackson. Networks of Waiting Lines. *Operations Research*, 5(4):518-521, 1957.
66. J. R. Jackson. Jobshop-Like Queueing Systems. *Management Science*, 10(1):131-142, 1963.
67. M. H. Jensen, K. Sneppen, and G. Tian. Sustained Oscillations and Time Delays in Gene Expression of Protein Hes1. *FEBS Letters*, 541:176-177, 2003.
68. M. Kalecki. A Macrodynamic Theory of Business Cycle. *Econometrica*, 3:327-344, 1935.
69. T. Kataoka, H. Kawamura, K. Kurumatani, and A. Ohuchi. Distributed Visitors Coordination System in Theme Park Problem. *MMAS 2004*, LNAI 3446:335-348, 2005.
70. S. A. Kaufmann. *The Origins of Order: Self-Organization and Selection in Evolution*. Oxford University Press, New York, NY, 1993.
71. H. Kawamura, T. Kataoka, K. Kurumatani, and A. Ohuchi. Investigation of Global Performance Affected by Congestion Avoiding Behavior in Theme Park Problem. *IEEEJ Transactions EIS*, 124(10):1922-1929, 2004.
72. H. Kawamura, K. Kurumatani, and A. Ohuchi. Modeling of Theme Park Problem with Multiagent for Mass User Support. *Lecture Notes in Computer Science*, 3012:48-69, 2004.
73. F. P. Kelly. Networks of Queues with Customers of Different Types. *Journal of Applied Probability*, 12(3):542-554, 1975.
74. M. Klein, R. Metzler, and Y. Bar-Yam. Handling Emergent Resource Use Oscillations. *IEEE Transactions on Systems, Man, and Cybernetics Part A: Systems and Humans*, 35(3):327-336, 2005.
75. Y. Kuang. *Delay Differential Equations with Applications in Population Dynamics*. Academic Press, Mathematics in Science and Engineering 191, 1991.
76. S. R. T. Kumara, P. Ranjan, A. Surana, and V. Narayanan. Decision Making in Logistics: A Chaos Theory Based Analysis. *Annals of the CIRP*, 52(1):381-384, 2003.
77. G. Labouchère. Comportement d'auto-organisation de la dynamique des clients en marché clos induit par des décisions basées sur leur historique. *Diploma Work*, EPFL, 2005.

78. S. Lämmer and D. Helbing. Self-Control of Traffic Lights and Vehicle Flows in Urban Road Networks. *Journal of Statistical Mechanics: Theory and Experiment*, 4, art. no. P04019, 2008.
79. A. K. Y. Law, Y. V. Hui, and X. Zhao. Modeling Repurchase Frequency and Customer Satisfaction for Fast Food Outlets. *International Journal of Quality & Reliability Management*, 21(5):545–563, 2004.
80. J. H. Lee and C.-O. Kim. Multi-Agent Systems Applications in Manufacturing Systems and Supply Chain Management: a Review Paper. *International Journal of Production Research*, 46(1):233–265, 2008.
81. L. S. Lee, K. D. Fiedler, and J. S. Smith. Radio Frequency Identification (RFID) Implementation in the Service Sector: A Customer-Facing Diffusion Model. *International Journal of Production Economics*, 112:587–600, 2008.
82. J. P. Lehoczky. Using Real-Time Queueing Theory to Control Lateness in Real-Time Systems. *Performance Evaluation Review*, 25(1):158–168, 1997.
83. Z. Lewkowicz and J.-D. Kant. A Multiagent Simulation of a Stylized French Labor Market: Emergences at the Micro Level. *Advances in Complex Systems*, 11(2):217–230, 2008.
84. H. Lian, E. Feng, X. Li, J. Ye, and Z. Xiu. Oscillatory Behavior in Microbial Continuous Culture with Discrete Time Delay. *Nonlinear Analysis: Real World Applications*, 10:2749–2757, 2009.
85. M. C. Mackey. Mathematical Models of Hematopoietic Cell Replication and Control. In: *The Art of Mathematical Modelling: Case Studies in Ecology, Physiology and Biofluids*, pp. 149-178, H. G. Othmer, F. R. Adler, M. A. Lewis and J. C. Dallon (Eds.), Prentice-Hall, Englewood Cliffs, 1997.
86. S. Maguire, B. McKelvey, L. Mirabeau, and N. Oztas. Complexity Science and Organization Studies. In: *Handbook of Organization Studies (Second Edition)*, pp. 165-214, S. Clegg, C. Hardy, T. Lawrence and W. Nord (Eds.), Sage, Thousand Oaks, 1996.
87. D. H. Maister. The Psychology of Waiting Lines. *The Service Encounter*, Chap. 8, J. A. Czepiel, M. R. Solomon and D. F. Surprenant (Eds.), D.C. Heath, Lexington, Mass, 1985.
88. P. Massotte and R. Bataille. Future Production Systems: Influence of Self-Organization on Approaches to Quality Engineering. *International of Production Economics*, 64:359–377, 2000.
89. P. M. Mathews and M. Lakshmanan. On a Unique Nonlinear Oscillator. *Quarterly of Applied Mathematics*, 32:215–218, 1974.
90. B. McKelvey. Quasi-Natural Organization Science. *Organization Science*, 8(4):352–380, 1997.
91. B. McKelvey. Self-Organization, Complexity Catastrophe, and Microstate Models at the Edge of Chaos. In: *Variations in Organization Science - In Honor of Donald T. Campbellin*, pp. 279-307, J. A. C. Baum and B. McKelvey (Eds.), Sage, Thousand Oaks, CA, 1999.
92. B. McKelvey. Emergent Strategy via Complexity Leadership: Using Complexity Science and Adaptive Tension to Build Distributed Intelligence. In: *Complexity and Leadership Volume I - Conceptual Foundations*, pp. 225-268, M. Uhl-Bien and R. Marion (Eds.), Information Age Publishing, Charlotte, NC, 2007.
93. J. Medhi. *Stochastic Models in Queueing Theory (Second Edition)*. Academic Press, 2003.

94. J. G. Milton and A. Longtin. Evaluation of Pupil Constriction and Dilation from Cycling Measurements. *Vision Research*, 30(4):515–525, 1990.
95. M. Mitchell Waldrop. *Complexity: The Emerging Science at the Edge of Order and Chaos*. Simon and Schuster, New York, NY, 1992.
96. M. Montefiori. Spatial Competition for Quality in the Market for Hospital Care. *European Journal of Health Economics*, 6:131–135, 2005.
97. T. Moyaux, B. Chaib-draa, and S. D’Amours. Multi-Agent Coordination Based on Tokens: Reduction of the Bullwhip Effect in a Forest Supply Chain. In: *Proceedings of AAMAS 2003*, Melbourne, Australia, 2003.
98. T. Moyaux, B. Chaib-draa, and S. D’Amours. Supply Chain Management and Multiagent Systems: An Overview. In: *MultiAgent-Based Supply Chain Management*, Chapter 1, B. Chaib-draa and J. P. Müller (Eds.), Springer, 2006.
99. J. I. Neimark. *Mathematical Models in Natural Science and Engineering*. Springer, Berlin, 2003.
100. Y. G. Panovko and I. I. Gubanova. *Stability and Oscillation of Elastic Systems (in Russian)*. Nauka, Moscow, 1964.
101. E. A. Peköz and N. Joglekar. Poisson Traffic Flow in a General Feedback Queue. *Journal of Applied Probability*, 39(3):630–636, 2002.
102. A. Pikovsky, M. Rosenblum, and J. Kurths. *Synchronization: A Universal Concept in Nonlinear Sciences*. Pages 41–44, Cambridge University Press, 2001.
103. S. Pingali, D. Tipper, and J. Hammond. The Performance of Adaptive Window Flow Controls, in a Dynamic Load Environment. In: *Proceedings of IEEE INFOCOM ’90*, pp. 55–62, 1990.
104. M. Pullman and G. M. Thompson. Evaluating Capacity - and Demand - Management Decisions at a Ski Resort. *Cornell Hotel & Restaurant Administration Quarterly*, 43(6):25–36, 2002.
105. T. Puu. Hotelling’s “Ice Cream Dealers” with Elastic Demand. *The Annals of Regional Science*, 36(1):1–17, 2002.
106. J. V. Ringwood and S. V. Malpas. Slow Oscillations in Blood Pressure via a Nonlinear Feedback Model. *American Journal of Regulatory, Integrative and Comparative Physiology*, 280(4):R1105–R1115, 2001.
107. P. S. Ruszczynski and L. B. Kish. Noise Enhanced Efficiency of Ordered Traffic. *Physics Letters A*, 276:187–190, 2000.
108. H. Sanner. Instability in Competition: Hotelling Re-considered. *Diskussionsbeitrag Nr. 79*, Universität Postdam, September 2005.
109. R. Serfozo. *Introduction to Stochastic Networks*. Springer, 1999.
110. C. Sun, M. Han, and Y. Lin. Analysis of Stability and Hopf Bifurcation for a Delayed Logistic Equation. *Chaos, Solitons and Fractals*, 31(3):672–682, 2007.
111. A. Surana, S. R. T. Kumara, M. Greaves, and U. N. Raghavan. Supply-Chain Networks: A Complex Adaptive Systems Perspective. *International Journal of Production Research*, 43(20):4235–4265, 2005.
112. L. Takács. A Single-Server Queue with Feedback. *The Bell System Technical Journal*, 42:505–519, 1963.
113. S. F. Tzeng, W.-H. Chen, and F.-Y. Pai. Evaluating the Business Value of RFID: Evidence From Five Case Studies. *International Journal of Production Economics*, 112:601–613, 2008.
114. A. van Ackere, C. Haxholdt, and E. R. Larsen. Long-Term and Short-Term Customer Reaction: A Two-Stage Queueing Approach. *System Dynamics Review*, 22(4):349–369, 2006.

115. A. van Ackere and E. R. Larsen. Self-Organising Behaviour in the Presence of Negative Externalities: A Conceptual Model of Commuter Choice. *European Journal of Operational Research*, 157:501–513, 2004.
116. B. van der Pol. On Relaxation Oscillation. *Philosophical Magazine*, 2:978–992, 1926.
117. B. van der Pol and J. van der Mark. The Heartbeat Considered as a Relaxation Oscillation, and an Electrical Model of the Heart. *Philosophical Magazine*, 6:763–775, 1928.
118. P. F. Verhulst. Notice sur la Loi que la Population Suit dans son Accroissement. *Correspondance Mathématique et Physique*, 10:113–125, 1838.
119. J. von Neumann and O. Morgenstern. *Theory of Games and Economic Behavior*. Princeton University Press, 1944.
120. J. Walrand. *An Introduction to Queueing Networks*. Prentice-Hall International Editions, Englewood Cliffs, 1988.
121. H. Warnecke. *The Fractal Company: A Revolution in Corporate Culture*. Springer, 1993.
122. S. Whitby, D. Parker, and A. Tobias. Non-Linear Dynamics of Duopolistic Competition: A R and D Model and Simulation. *Journal of Business Research*, 51:179–191, 2001.
123. P. Whittle. *Systems in Stochastic Equilibrium*. John Wiley and Sons, 1986.
124. K. H. Wolf and J. Venus. Description of the Delayed Microbial Growth by an Extended Logistic Equation. *Acta Biotechnologica*, 12(5):405–410, 1992.
125. E. Wong. The Construction of a Class of Stationary Markoff Processes. *Sixteenth Symposium in Applied Mathematics - Stochastic Processes in Mathematical Physics and Engineering*, pp. 264–276, R. Bellman (Ed.), Providence, RI: American Mathematical Society, 1964.
126. M. Woolridge. *An Introduction to Multiagent Systems (Second Edition)*. Chapter 10, John Wiley and Sons, 2009.
127. C. Wycisk, B. McKelvey, and M. Hülsmann. “Smart Parts” Supply Networks as Complex Adaptive Systems: Analysis and Implications. *International Journal of Physical Distribution & Logistics Management*, 38(2):108–125, 2008.

Index

- accumulate-and-fire dynamics, 40
- agent
 - agent-based model, 10, 21
 - agent-induced limit-cycle, 40
 - autonomous agent, 5, 21, 30, 39, 55, 58, 123, 136
 - heterogeneous agents, 50, 57
 - homogeneous agents, 64
 - multi-agent dynamics, 17, 63
 - multi-agent system, 21, 30, 78, 103, 135
- assembly line, 18
- autonomy, 57
- BCMP network, 9
- Bernoulli sampling, 80, 105
- brand departure cost, 102
- Brownian motion
 - standard Brownian motion, 106
- bucket brigade, 18
- bullwhip effect, 22, 153
- busy period
 - maximum busy period, 39
- cannibalization
 - cannibalization effect, 89, 126
 - super-cannibalization effect, 90, 93
- centralized management, 44
- coevolution, 58
- collective
 - collective behaviour, 21, 78
- competition
 - competing servers, 94
 - quality of service competition, 90
- complex systems, 10, 19
 - complex adaptive logistics systems, 56, 133
 - complex adaptive systems, 17, 56
 - complexity science, 19
- congestion control, 132
- decentralization
 - decentralized management, 44, 135
- decision node, 31
- delay
 - delay differential equation, 23, 26
 - delay mechanism, 37, 46, 71, 141, 154
 - delayed feedback, 46
 - random delay, 26
 - state-dependent delay, 26
 - time delay, 17, 23, 43, 165
- deterministic polling, 80
- differential equation
 - delay differential equation, 23, 26
 - ordinary differential equation, 23, 172
 - stochastic differential equation, 108
- diffusion process, 106
- dipole, 79
- dispatching policy
 - autonomous dispatching policy, 81, 135
 - fixed dispatching policy, 80
 - random dispatching policy, 80
 - shortest-queue-first dispatching policy, 86
- dissipative structures
 - theory of dissipative structures, 20
- dynamical system, 23

- economic order quantity, 43
- electronic generator, 40
- externality, 100
- feedback
 - delayed feedback, 46
 - feedback control, 23, 153
 - feedback loop, 17, 31, 46, 51, 62, 79
 - informational feedback, 141
 - state-dependent feedback, 33
- FIFO (first-in-first-out), 31, 63, 79, 102
- flow
 - flow reinjection, 30
 - state-dependent reentering flow, 33
- Fokker-Planck
 - Fokker-Planck equation, 108
 - Fokker-Planck operator, 108, 111
- game theory, 19, 61
- Hamiltonian, 41
- heavy traffic
 - heavy traffic regime, 106
- history-based
 - history-based routing decisions, 3, 5, 15, 30, 62, 120, 133, 148
 - history-based routing rule, 32, 91
- Hopf bifurcation, 24
- Hotelling's model, 99, 101, 120
- hydrodynamic analogy, 37, 141
- interaction, 57
 - stigmergic interactions, 6, 57, 58
- Jackson network, 8
- Kelly network, 9
- law of large numbers, 36, 66, 92, 124, 140
- learning
 - ability to learn, 58
- leisure and hospitality sector, 15
- limit-cycle, 43
 - agent-induced limit-cycle, 40
 - limit-cycle oscillations, 24
 - stable limit-cycle, 41
- Little's law, 137
- load sharing
 - load sharing policy, 131, 134
- logistics
 - complex adaptive logistics systems, 56, 133
 - logistics network, 55
- loyalty
 - customer loyalty, 3, 15, 16
 - loyal customer, 32, 121
- manufacturing system, 17
- marginal stability regime, 84
- market partition boundary point, 104, 121
- market sharing dynamics, 92, 107, 121
- Markov
 - Markov chain, 7
 - Markovian dynamics, 32
 - Markovian routing policy, 7
 - non-Markovian dynamics, 15, 30
 - non-Markovian routing policy, 5
- Mathews-Lakshmanan oscillator, 41
- melting zone, 58
- multiplicative noise process, 117
- network topology
 - closed network topology, 90, 149
 - open network topology, 79
- neural network, 20
- nonlinear dynamics, 15
- Occam's razor, 11
- oscillations
 - cyclo-stationary regime, 36, 45, 65, 122
 - market sharing oscillations, 120
 - oscillatory dynamics, 23, 149
 - queue content oscillations, 17, 36, 49, 50, 93, 136, 140, 149, 150, 168
 - self-sustained oscillations, 24, 41
 - softened oscillatory dynamics, 150
 - stable temporal oscillations, 35, 72
 - synchronized oscillations, 86
 - uncoupled oscillations, 80
- path dependency, 148
- patience parameter, 32, 82
 - homogeneous patience parameter, 32, 83
 - individualized patience parameter, 50
- periodic wave, 123

- phase transition, 58, 150
 - noise-induced phase transition, 109
- pheromone, 133
- Poisson process, 105
- production
 - production centre, 134
 - production line, 18
 - production management, 17
- purging
 - purging rate, 69
 - queue content periodic purging, 38, 124
- quasi-deterministic regime, 35
- queue content
 - expected queue content, 81
 - observed queue content, 81
- queueing systems
 - queueing network, 7, 15
 - queueing network theory, 5
 - real-time queueing system, 17
 - single-server queueing system, 31
 - spatial queueing system, 99, 119
 - stable queueing system, 83
 - two-servers queueing system, 79, 90
 - unstable queueing system, 82
- rare events regime, 84
- recurrent service, 16, 30, 63, 78, 120
- relaxation oscillator, 37, 39
 - stigmergic relaxation oscillator, 40
- relaxation time, 111, 112, 123, 150
- reorder point policy, 43
- resource use maximization, 39
- RFID, 55, 134
- routing balance equations, 7
- satisfaction
 - customer satisfaction, 3, 16
 - long-term satisfaction, 16
 - service satisfaction, 16
 - short-term satisfaction, 16
 - waiting time satisfaction, 16
- self-organization, 58
 - self-organized criticality, 58
 - self-organized structure, 38, 87, 124
- semaphore, 135
- server
 - heterogeneous servers, 86, 92, 114
 - homogeneous servers, 86, 93, 105
- signal-to-noise ratio, 141
- simulation framework, 10
- siphon mechanism, 37, 46, 49, 72, 85, 140
- smart parts, 55, 135, 153
- stabilization, 38
 - noise-induced stabilization, 83
 - system stabilization, 45
- stationarity
 - stationary probability density function, 108
 - stationary regime, 34, 44, 150
- stigmergy
 - stigmergic interactions, 6, 57, 58
 - stigmergic relaxation oscillator, 40
- stylized model, 10, 11
- supply chains, 17, 43, 153
- Tantalus glass, 37, 142
- traffic load
 - traffic load estimator, 139
- transient regime, 111, 148, 150
- transportation
 - transportation cost, 102
 - transportation system, 25, 147
- utility function, 64, 103
- volatility
 - state-dependent volatility, 106
- waiting time
 - expected waiting time, 78, 82, 103
 - experimented waiting time, 32, 78, 82, 121, 150
 - perceived waiting time, 16, 78, 82
 - waiting time cost, 102
- weariness, 63
- Whittle network, 8

Curriculum Vitae

OLIVIER GALLAY

Born: 7th July 1981 in Lausanne (VD)

Nationality: Swiss

Contact: olivier.gallay_1@a3.epfl.ch

Education

- 2000: Swiss Federal Maturity in Mathematics and Sciences and Baccalaureate ès Science, Pully (VD).
- 2003: Bachelor of Science in Communication Systems from the Ecole Polytechnique Fédérale de Lausanne (EPFL).
- 2005: Master of Science in Communication Systems from the Ecole Polytechnique Fédérale de Lausanne (EPFL).

Recent Professional Experience

- 2005: Internship at the IBM Zürich Research Laboratory, Rüschlikon (ZH). Project: *Robustness and Regulatory Risk in Supply Chains*.
- 2005-2009: Part-time tutorial work at the Laboratory of Microengineering for Manufacturing (LPM-EPFL). Supervised and guided students' master and semester projects.
- 2006-2009: Reviewing activities for the following journals:
- European Journal of Operational Research
 - International Journal of Production Economics
 - IEEE Transactions on Industrial Informatics

Publications

Journal Articles

- O. Gallay and M.-O. Hongler. Cooperative Dynamics of Loyal Customers in Queueing Networks. *Journal of Systems Science and Systems Engineering*, 17(2):241-254, 2008.
- O. Gallay and M.-O. Hongler. Market Sharing Dynamics Between Two Service Providers. *European Journal of Operational Research*, 190:241-254, 2008.
- O. Gallay and M.-O. Hongler. Circulation of Autonomous Agents in Production and Service Networks. *International Journal of Production Economics*, 120:378-388, 2009.
- O. Gallay and M.-O. Hongler. Multi-Agent Adaptive Mechanism Leading to Optimal Real-Time Load Sharing. To appear in the *Journal of Computing and Information Technology*, 2010.
- O. Gallay and M.-O. Hongler. Stylized Models for Recurrent Services. To appear in the *Supply Chain Forum: an International Journal in Supply Chain Management*, 2010.
- M.-O. Hongler, O. Gallay, M. Hülsmann, P. Cordes and R. Colmorn. Centralized Versus Decentralized Control - A Solvable Stylized Model in Transportation. To appear in *Physica A*, 2010.

Articles in Refereed International Conference Proceedings

- O. Gallay and M.-O. Hongler. Cooperative Dynamics of Loyal Customers in Queueing Networks. In: *Proceedings of IC SSSM'06, International Conference on Service Systems and Service Management*, Troyes, France, October 2006.
- O. Gallay and M.-O. Hongler. Market Partition in a Dynamic Linear City Game Model. In: *Proceedings of EURO XXII, 22nd European Conference on Operational Research*, Prague, Czech Republic, July 2007.
- O. Gallay and M.-O. Hongler. Weariness and Loyalty Loss in Recurrent Service Models. In: *Proceedings of MOSIM'08, 7th International Conference on Modelization and Simulation*, Paris, France, April 2008.
- O. Gallay and M.-O. Hongler. Multi-Agent Adaptive Mechanism Leading to Optimal Real-Time Load Sharing. In: *Proceedings of MATHMOD 2009, 6th Vienna International Conference on Mathematical Modelling*, Vienna, Austria, February 2009.

Selected Talks

- O. Gally. Compétition entre Deux Serveurs avec Boucle de Retour. *Rencontre Franco-Suisse sur le Bruit et les Non-Linéarités*, Grenoble, France, November 2005.
- O. Gally. Cooperative Dynamics of Loyal Customers in Queueing Networks. *Madeira Math Encounters XXXI (Stochastics, Networks, Infinite Particle Systems)*, CCM, Funchal, Madeira, October 2006.
- O. Gally. Adding Waiting Time Penalties to the Hotelling Model. *Madeira Math Encounters XXXI (Stochastics, Networks, Infinite Particle Systems)*, CCM, Funchal, Madeira, October 2006.
- O. Gally. Stylized Models of Queueing Networks Roamed by Autonomous Agents. *International Workshop on Evolution and Structure of Complex Systems and Networks: Basic Techniques and Applications*, ZiF, Bielefeld, Germany, February 2008.
- O. Gally. Multi-Agent Adaptive Mechanism Leading to Optimal Real-Time Load Sharing. *Rencontre Franco-Suisse sur le Bruit et les Non-Linéarités*, EPFL, Lausanne, November 2008.
- O. Gally. Collective Flow Patterns Generation in Queueing Networks Roamed by Autonomous Agents. *Madeira Math Encounters XXXVII (Particle Systems, Feynman Integrals, Stochastic Analysis)*, CCM, Funchal, Madeira, July 2009.
- M.-O. Hongler and O. Gally. Collective Flow Patterns Generation in Queueing Networks Roamed by Autonomous Agents. *Collaborative Research Centre 637 (SFB 637): Autonomous Cooperating Logistic Processes - A Paradigm Shift and its Limitations*, Jacobs University, Bremen, Germany, September 2009.
- O. Gally. Market Sharing Dynamics Between Two Service Providers. *Collaborative Research Centre 637 (SFB 637): Autonomous Cooperating Logistic Processes - A Paradigm Shift and its Limitations*, Jacobs University, Bremen, Germany, September 2009.