

Monotonicity properties of user equilibrium policies for parallel batch systems

Yizheng Chen, Mark Holmes and Ilze Ziedins

July 8, 2011

Abstract

We study a simple network with two parallel batch service queues, where service at a queue commences when the batch is full and each queue is served by infinitely many servers. A stream of general arrivals observe the current state of the system on arrival and choose which queue to join to minimize their own expected transit time. We show that for each set of parameter values there exists a unique user equilibrium policy and that it possesses various monotonicity properties. User equilibrium policies for probabilistic routing are also discussed and compared with the state-dependent setting.

1 Introduction

It is well known that networks where individuals can choose their own route through the system may experience far worse performance than that seen under system optimal routing (for instance, Bell and Stidham [6], Cohen and Kelly [15] and Whitt [33]). This is an important issue for both communication and transportation networks – see Patriksson [24] for an overview of traffic assignment problems and Roughgarden and Tardos [25] for a discussion of selfish routing in communication networks.

The model we consider in this paper is motivated primarily by the problem of selfish routing choices in transportation applications and is a simple parallel queueing network with two routes from a source to a common destination (see Figure 1). The network consists of two batch service queues with batch sizes N_i , $i = 1, 2$ and each queue has an infinite number of servers, i.e. both routes are $M/M^{(N_i)}/\infty$ queues with $N_i \geq 2$. There are three independent Poisson arrival streams: two intrinsic arrival streams (at rate σ_1 and σ_2 to queue 1 and queue 2 respectively), and a stream of general arrivals at rate λ who may choose to join either of the two queues. In this system users may prefer to join a longer queue with more customers awaiting service if the free capacity is lower than in the other queue, since the expected waiting time for service may then be shorter. The service time at each queue is exponentially distributed with rate μ_i , $i = 1, 2$. Service commences at a queue when its batch is full. Neither jockeying nor reneging is allowed, so that once a customer has joined a queue, they must stay in their batch of their queue until their batch service is completed. All inter-arrival times and service times are independent of each other. Figure 1 shows the structure of the system. Applications of batch service queues

include examples of public transportation, such as airport shuttle buses that leave when full (Afimeimounga et al. [1], [2] and Jia [20]). A modified model where the batch service may commence before the batch is full is more generally applicable to bus and shared taxi services.

In this paper we concentrate mostly on state-dependent routing where routing decisions made by general arrivals depend on both the parameter values of the system (arrival and service rates and batch sizes) and the instantaneous state of the system when a general arrival occurs. However, we do also consider probabilistic routing for general arrivals, which corresponds to the case where an individual's knowledge of the system is limited to the parameter values, and information about the instantaneous state of the system is not available when routing decisions are made. This may be a more realistic model for systems where individuals accumulate knowledge of the system over a period of time, and routing patterns change gradually in response to that. We give examples showing that state-dependent routing may lead to shorter delays – that is, increased knowledge of the instantaneous state of the system may lead to improved performance overall – although unlike for system optimal policies, it is not always the case that increasing the state information on which decisions are based leads to improved performance.

Routing problems such as these fall within a class of optimization and control problems known as dynamic games (see e.g. [4, 3] and the references therein). We are interested here in minimizing delay, but other forms of performance measure may also be considered (for instance, blocking probabilities in loss networks). In a game with a fixed finite number of players, a strategy is a Nash equilibrium [23] if no player can benefit by changing her strategy, provided all other players keep their strategies unchanged. Wardrop [31] studied traffic assignment problems in the 1950s and defined equilibria in the context of road networks with infinitely many road users such that each individual has an infinitesimal effect upon the system. Under a Wardrop equilibrium, the journey times on all routes in use are equal, and less than those which would be experienced by a single vehicle on any unused route [31]. We are interested in types of decision policies, known as user equilibrium or user optimal policies, under which no individual can improve his/her perceived travel time by unilaterally adopting a different decision policy.

In those models where the delay experienced by an arrival depends only on the current state of the system, user equilibrium policies are relatively easy to find (see e.g. Winston [34] and Spicer and Ziedins [26]). In the model considered here, however, a customer's delay may depend not only on the current number of customers in the system, but also on routing decisions of subsequent arrivals. The queueing network in which Braess's paradox is commonly studied is an example of a system where this is the case (Braess [9], translated in Braess et al. [10], Cohen and Kelly [15] and Calvert et al. [13]). More recently, Altman and Shimkin [5] considered a system where customers choose between a processor-sharing queue and an infinite server queue. They used a coupling argument to show that a unique user equilibrium policy exists and possesses a certain natural property. Ben-Shahar et al. [7] extended their results to multiple customer types, and Brooms [11] studied a related system. Hassin and Haviv [19] studied user decisions in a network with an infinite server queue. For their models a one-dimensional state representation of the system suffices to calculate expected delays, since delays at the infinite server

queue do not depend on the length of the queue. Afimeimounga et al. [1, 2] studied a system that is perhaps most closely related to the one we consider here, consisting of a batch-service queue and an $M|M|1$ queue, under both probabilistic and state-dependent routing, which did require a two-dimensional state representation. For the state dependent case they showed, as we will do here, the existence and uniqueness of the user equilibrium. In work in progress [14] the authors study a parallel system with two processor sharing queues.

Since performance under selfish routing may be considerably worse than at the system optimum, it is perhaps not surprising that changing parameters with the intention of improving performance (for instance, by increasing service rates, or adding routes to a network) may instead have a negative effect on system performance. This was first observed with Braess's paradox ([9, 10]) with the addition of an extra route to the system, and later authors have shown similar paradoxical effects can occur when service rates are increased (see, for instance, Calvert [12], Downs [16], Thomson [28], and Afimemimounga et al. [1, 2] for a discussion of the Downs-Thomson paradox). Numerical examples show that for our network, increasing service rates may lead to greater delays under the user equilibrium. As in [2], in our examples this non-monotonicity effect is considerably reduced under state-dependent routing as opposed to probabilistic routing – increasing information given to the users may improve performance in this respect.

Finally, we mention briefly that there is a long history of studying optimal policies for parallel queues, whether they be the user equilibrium policies that we study here, or system optimal policies (see e.g. Winston [34], Weber [32], Koole et al. [21] and Whitt [33] and the references therein). Walrand [30] and Gelenbe et al. [18] give excellent general introductions to queueing networks and Boxma et al. [8] provide a more recent overview of solution methods for performance analysis of parallel and distributed systems. The results in this paper rely heavily on use of the coupling method. Excellent introductions to these can be found in Lindvall [22], Thorisson [29], and El-Taha and Stidham [17] .

The structure of the paper is as follows. In Section 2 we introduce our basic notation, formally define the (expected) delay and the user equilibrium, and state our main results. In Section 3, using an explicit construction and coupling arguments we prove our main results on existence and monotonicity properties of state-dependent user equilibrium policies for the parallel batch service systems under consideration. In Section 4 we discuss probabilistic routing and give numerical results comparing user equilibria in the stationary regime under probabilistic routing with the state-dependent setting. The results suggest that arrivals who incorporate full knowledge of the current state of the system in their individual decision making can lead to a reduction in non-monotonicity effects and reduce the overall delay in the system. We conclude with possible extensions to the model and further discussion.

2 Notation and definitions

The network of interest consists of two $M/M^{(N_i)}/\infty$ queues in parallel, with batch sizes N_1 and N_2 , service rates μ_1 and μ_2 and dedicated arrival rates σ_1 and σ_2

for queues 1 and 2 respectively. In addition general arrivals occur at rate λ , and each such arrival makes an instantaneous decision of which queue to join based on knowledge of the state of the system immediately prior to their arrival, the parameter values and the queueing mechanism. All inter-arrival and service times are independent of each other and exponentially distributed; and the decisions made by general arrivals are also independent and time-homogeneous, given the state of the system.

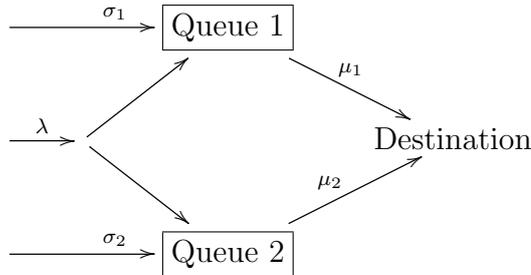


Figure 1: The network diagram for two $M/M^{(N_i)}/\infty$ queues

To define the process associated with this system, we fix N_1, N_2 and $\Gamma = \{\lambda, \sigma_1, \mu_1, \sigma_2, \mu_2, \tilde{p}\}$, and denote the state of the system at a given time t by $Z(t) = (Z_1(t), Z_2(t))$ where $Z_i(t)$ is the number of customers awaiting service in queue i at time t , $i = 1, 2$. Suppose that a customer/user arrives at time t when the system is in state $\mathbf{n} = (n_1, n_2)$ (we write this as $Z(t^-) = (n_1, n_2)$). If that user joins queue 1, then $Z(t) = ((n_1 + 1) \bmod N_1, n_2)$, while if the user joins queue 2, then $Z(t) = (n_1, (n_2 + 1) \bmod N_2)$. Note that if $Z_i(t^-) = N_i - 1$ and the arrival joins queue i , then since service commences immediately once the batch is complete, the number of customers waiting for service in queue i drops to zero, i.e. $Z_i(t) = 0$. Hence, we define the state space of the network to be $\mathbf{S} = \{\mathbf{n} = (n_1, n_2) : 0 \leq n_i \leq N_i - 1, n_i \in \mathbb{Z}, \text{ for } i = 1, 2\}$. For notational convenience, we set $\mathbf{e}_1 = (1, 0)$, and $\mathbf{e}_2 = (0, 1)$. For $\mathbf{n} \in \mathbf{S}$ we write

$$\mathbf{n} + \tilde{\mathbf{e}}_i = \begin{cases} ((n_1 + 1) \bmod N_1, n_2), & \text{if } i = 1 \\ ((n_1, (n_2 + 1) \bmod N_2), & \text{if } i = 2. \end{cases}$$

In every state $\mathbf{n} \in \mathbf{S}$, an arriving general customer must decide which queue to join. Hence, we make the following definitions:

Definition 2.1. A decision policy $D = \{D_0, D_1, D_2\}$ for \mathbf{S} is a partition of \mathbf{S} such that if a general customer arrives to find the system in state $\mathbf{n} = (n_1, n_2) \in D_i$, $i = 1, 2$, then he/she joins queue i ; if $\mathbf{n} \in D_0$, then the general customer joins queue 1 with probability \tilde{p} , and queue 2 with probability $1 - \tilde{p}$. The set of all decision policies is denoted by $\mathcal{D} = \mathcal{D}(\mathbf{S})$. For fixed Γ , the process operating under $D \in \mathcal{D}$ is denoted by $Z_D = \{Z_D(t)\}_{t \geq 0}$.

Let Q_D be the Q -matrix of the process Z_D operating under decision policy $D \in \mathcal{D}$ with parameters Γ . The transition rates from state $\mathbf{n} = (n_1, n_2)$ to $\mathbf{n}' =$

$(n'_1, n'_2) \neq \mathbf{n}$ are

$$Q_D(\mathbf{n}, \mathbf{n}') = \begin{cases} \sigma_1 + \lambda \cdot I_{\mathbf{n} \in D_1} + \tilde{p}\lambda \cdot I_{\mathbf{n} \in D_0}, & \text{if } \mathbf{n}' = \mathbf{n} + \tilde{\mathbf{e}}_1 \\ \sigma_2 + \lambda \cdot I_{\mathbf{n} \in D_2} + (1 - \tilde{p})\lambda \cdot I_{\mathbf{n} \in D_0}, & \text{if } \mathbf{n}' = \mathbf{n} + \tilde{\mathbf{e}}_2 \\ 0, & \text{otherwise,} \end{cases}$$

where I_A is an indicator function taking the value 1 if A is true, and 0 otherwise.

Let $H_i = \{\mathbf{n} = (n_1, n_2) : n_i = 0\} \subset \mathbf{S}$, $i = 1, 2$. Then the first hitting time or reaching time to H_i from state \mathbf{n} is defined as:

Definition 2.2. Given a process Z_D such that $Z_D(t) = \mathbf{n}$ let

$$T_{i;D}^Z(\mathbf{n}; t) = \inf\{s \geq 0 : Z_D(t+s) \in H_i\}, \quad i = 1, 2. \quad (1)$$

be the first hitting time to H_i given that $Z_D(t) = \mathbf{n}$.

Note that $T_{i;D}(\mathbf{n}; t)$ represents the waiting time (from time t) until service commences at queue i , if the system is in state $\mathbf{n} \in \mathbf{S}$ at time t . The distribution of $T_{i;D}^Z(\mathbf{n}; t)$ does not depend on t and we will often drop this from the notation. We also write $T_{i;D}(\mathbf{n}) = T_{i;D}^Z(\mathbf{n})$ when the process Z_D is clear.

Definition 2.3. Denote by $m_{i;D}(\mathbf{n}) = \mathbb{E}[T_{i;D}(\mathbf{n})]$, $i = 1, 2$, the expected waiting time via queue i until service commences when the system is in state $\mathbf{n} \in \mathbf{S}$, i.e. $Z_D(t) = \mathbf{n}$.

Definition 2.4. Let $z_{i;D}(\mathbf{n})$ be the expected transit time through the system (including the service time) for a general customer arriving at time t if $Z_D(t-) = \mathbf{n} \in \mathbf{S}$ and $Z_D(t) = \mathbf{n} + \tilde{\mathbf{e}}_i$, that is, for a general arrival who joins queue i when the system is in state $\mathbf{n} \in \mathbf{S}$.

If a general arrival joins queue i when $n_i = N_i - 1$, then that arrival does not wait for service since s/he has completed the current batch and her/his arrival triggers an immediate service commencement for that batch at queue i . Thus the expected transit time for the last arrival in a batch is just the expected service time for queue i . Now suppose that a general arrival joins queue i when the process Z_D is in state \mathbf{n} with $n_i < N_i - 1$. Then the process jumps to state $\mathbf{n} + \mathbf{e}_i$, and thus $z_{i;D}(\mathbf{n})$ is the expected hitting time to H_i from state $\mathbf{n} + \mathbf{e}_i$ plus the expected service time for queue i . Hence, for $i = 1, 2$,

$$z_{i;D}(\mathbf{n}) = \frac{1}{\mu_i}, \quad \text{if } n_i = N_i - 1; \quad (2)$$

$$z_{i;D}(\mathbf{n}) = m_{i;D}(\mathbf{n} + \mathbf{e}_i) + \frac{1}{\mu_i}, \quad \forall n_i < N_i - 1. \quad (3)$$

Let $\Lambda = \sigma_1 + \sigma_2 + \lambda$. If $n_i < N_i - 1$, in order to find $z_{i;D}(\mathbf{n})$, we condition on the first jump of the process after a general arrival has joined queue i . The time to the next event is exponentially distributed with mean $\frac{1}{\Lambda}$. Given $D \in \mathcal{D}$, the next change of state can be due to an intrinsic arrival to queue 1 or queue 2, or a general arrival, in which case the resulting state depends on D . Therefore if $0 \leq n_1 < N_1 - 1$,

$$\begin{aligned}
z_{1;D}(\mathbf{n}) &= \frac{1}{\Lambda} + \frac{\sigma_1 + \lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_1 \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_1 \in D_0\}}}{\Lambda} z_1(\mathbf{n} + \mathbf{e}_1) \\
&\quad + \frac{\sigma_2 + \lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_1 \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_1 \in D_0\}}}{\Lambda} z_1(\mathbf{n} + \tilde{\mathbf{e}}_2).
\end{aligned} \tag{4}$$

Similarly if $0 \leq n_2 < N_2 - 1$,

$$\begin{aligned}
z_{2;D}(\mathbf{n}) &= \frac{1}{\Lambda} + \frac{\sigma_1 + \lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_2 \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_2 \in D_0\}}}{\Lambda} z_2(\mathbf{n} + \tilde{\mathbf{e}}_1) \\
&\quad + \frac{\sigma_2 + \lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_2 \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{n}+\mathbf{e}_2 \in D_0\}}}{\Lambda} z_2(\mathbf{n} + \mathbf{e}_2).
\end{aligned} \tag{5}$$

It can be shown that this finite set of linear equations has a unique solution for all $\lambda, \sigma_1, \sigma_2 > 0$. We now define user equilibrium policies D^* formally as follows.

Definition 2.5. A policy $D^* \in \mathcal{D}$ is a user equilibrium policy for (N_1, N_2, Γ) if $\forall \mathbf{n} \in \mathbf{S}$,

$$\mathbf{n} \in \begin{cases} D_1^*, & \iff z_{1;D^*}(\mathbf{n}) < z_{2;D^*}(\mathbf{n}) \\ D_0^*, & \iff z_{1;D^*}(\mathbf{n}) = z_{2;D^*}(\mathbf{n}) \\ D_2^*, & \iff z_{1;D^*}(\mathbf{n}) > z_{2;D^*}(\mathbf{n}). \end{cases}$$

We now give the formal definition of monotonicity as follows.

Definition 2.6. A policy $D = \{D_0, D_1, D_2\}$ is monotone, if it satisfies the following two conditions for each $i, j \in \{1, 2\}$, $i \neq j$:

1. $\mathbf{n} \in D_i \cup D_0 \Rightarrow \mathbf{n} + \mathbf{e}_i \in D_i$, $\forall \mathbf{n} \in \mathbf{S}$ such that $n_i < N_i - 1$;
2. $\mathbf{n} \in D_i \cup D_0 \Rightarrow \mathbf{n} - \mathbf{e}_j \in D_i$, $\forall \mathbf{n} \in \mathbf{S}$ such that $0 < n_j$.

The following two theorems are the main results of this paper.

Theorem 2.7. For each set of (non-zero) parameter values (N_1, N_2, Γ) there exists a unique user equilibrium policy, D^* .

Theorem 2.8. Any user equilibrium D^* is monotone.

Further monotonicity results (in terms of varying Γ) appear in Section 3.3.

3 Uniqueness and monotonicity

In this section, we prove a number of important lemmas and use these to prove the main results of this paper. We use the notation $X \stackrel{st}{\leq} Y$ if $\mathbb{P}(X \leq x) \geq \mathbb{P}(Y \leq x)$ for every $x \in \mathbb{R}$. Throughout this section N_1, N_2 , and Γ are fixed.

3.1 Preliminary lemmas

Definition 3.1. Let $i, j \in \{1, 2\}$, $i \neq j$. A policy $D = \{D_0, D_1, D_2\}$ is m -level monotone for queue i , $m \in \{1, 2, \dots, N_i\}$, if it satisfies the following two conditions:

1. $\mathbf{n} \in D_i \cup D_0 \Rightarrow \mathbf{n} + \mathbf{e}_i \in D_i$, $\forall N_i - m \leq n_i < N_i - 1$;
2. $\mathbf{n} \in D_i \cup D_0 \Rightarrow \mathbf{n} - \mathbf{e}_j \in D_i$, $\forall 0 < n_j \leq N_j - 1$.

Definition 3.2. We say $D \in \mathcal{D}$ is monotone for queue i , if D is N_i -level monotone for queue i , $i = 1, 2$.

Definition 3.3. A policy $D \in \mathcal{D}$ is (ξ, τ) -level monotone, if it is ξ -level monotone for queue 1 and τ -level monotone for queue 2.

Clearly a policy $D \in \mathcal{D}$ is monotone (Definition 2.6) if and only if it is (N_1, N_2) -level monotone.

Lemma 3.4. Let $\mathbf{u} = (u_1, u_2), \mathbf{v} = (v_1, v_2) \in \mathbf{S}$ satisfy one of the following conditions.

- (I) $u_2 - v_2 \geq v_1 - u_1 \geq 0$;
- (II) $u_2 - v_2 \geq N_1 - (u_1 - v_1) > 0$.

Suppose that $D \in \mathcal{D}$ is $\tau = N_2 - v_2$ level monotone for queue 2. Then $T_{2,D}(\mathbf{u}) \stackrel{st}{\leq} T_{2,D}(\mathbf{v})$. Furthermore, if $u_2 > v_2$, then $m_{2,D}(\mathbf{u}) < m_{2,D}(\mathbf{v})$.

Proof. The proof of this lemma, as for several of the results below, relies on a coupling argument. We construct a joint process $\{(Y, W)(t)\}_{t \geq 0}$ with state space $\mathbf{S} \times \mathbf{S}$ such that both Y and W follow the law of Z_D . Let

$$(Y, W)(0) = ((Y_1, Y_2), (W_1, W_2))(0) = (\mathbf{u}, \mathbf{v}).$$

The transitions of this joint process are such that Y and W make the same transition whenever possible, thus automatically preserving whichever of the conditions (I) or (II) holds. In those cases where that is not possible, the transition is chosen so that the new state still satisfies either (I) or (II). The joint process has the following transitions for $(\mathbf{a}, \mathbf{b}) \in \mathbf{S} \times \mathbf{S}$:

- (i) At rate $\sigma_1 + \lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_0 \cup D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

- (ii) At rate $\sigma_2 + \lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_0 \cup D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

- (iii) At rate $\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_2\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

(iv) At rate $\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_1\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

Since we are interested in the hitting time to the set H_2 for $X \in \{Y, W\}$, let $T_2^X(\mathbf{a}, \mathbf{b}; t) = \inf\{s \geq 0 : X_2(s+t) \in H_2\}$ be the marginal hitting time to H_2 at time t if the joint system is in state (\mathbf{a}, \mathbf{b}) at time t . We will often write $T_2^X(\mathbf{a}, \mathbf{b})$ for $T_2^X(\mathbf{a}, \mathbf{b}; t)$ since the distribution of $T_2^X(\mathbf{a}, \mathbf{b}; t)$ does not depend on t . We will show that for any pair of states (\mathbf{a}, \mathbf{b}) satisfying either of conditions (I) or (II), $T_2^Y(\mathbf{a}, \mathbf{b}; 0) \leq T_2^W(\mathbf{a}, \mathbf{b}; 0)$ almost surely. We do this by showing that with probability 1, $Y_2(t) \geq W_2(t)$, $\forall 0 \leq t < T_2^Y(\mathbf{a}, \mathbf{b}; 0)$.

Notice that condition (I) includes the case $\mathbf{a} = \mathbf{b}$. If the joint process jumps into state (\mathbf{a}, \mathbf{a}) , then $Y = W$ with probability 1, and for any $\mathbf{a} \in \mathbf{S}$, $T_2^Y(\mathbf{a}) = T_2^W(\mathbf{a})$. Henceforth, we assume $\mathbf{a} \neq \mathbf{b}$.

Denote by $(\mathbf{a}', \mathbf{b}')$ the resulting state of the joint system immediately after a transition from state (\mathbf{a}, \mathbf{b}) . We discuss each of the transitions (i) to (iv) in turn and show that when (\mathbf{a}, \mathbf{b}) satisfies either of conditions (I) and (II), so does $(\mathbf{a}', \mathbf{b}')$.

Case (i): if $\max(a_1, b_1) < N_1 - 1$, then $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_1, \mathbf{b} + \mathbf{e}_1)$ and whichever condition (\mathbf{a}, \mathbf{b}) satisfied is preserved (i.e. is also satisfied by $(\mathbf{a}', \mathbf{b}')$). If $\max(a_1, b_1) = N_1 - 1$, there are three cases: if $a_1 = b_1 = N_1 - 1$, then (\mathbf{a}, \mathbf{b}) satisfies condition (I) and $(\mathbf{a}', \mathbf{b}') = ((0, a_2), (0, b_2))$, which clearly preserves condition (I); if $a_1 < b_1 = N_1 - 1$, then $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_1, (0, b_2))$, so $a'_2 - b'_2 = a_2 - b_2 \geq N_1 - (1 + a_1) = N_1 - (a'_1 - b'_1) > 0$ and condition (II) is satisfied; if $b_1 < a_1 = N_1 - 1$, then $(\mathbf{a}', \mathbf{b}') = ((0, a_2), \mathbf{b} + \mathbf{e}_1)$, so $a'_2 - b'_2 = a_2 - b_2 \geq N_1 - (a_1 - b_1) = N_1 - (N_1 - 1 - b_1) = b'_1 - a'_1$, and thus condition (I) is satisfied.

Case (ii): if $a_2 = N_2 - 1$, $\mathbf{a}' = (a_1, 0) \in H_2$ and $T_2^Y(\mathbf{a}, \mathbf{b}) \leq T_2^W(\mathbf{a}, \mathbf{b})$. If $a_2 < N_2 - 1$, then $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_2, \mathbf{b} + \mathbf{e}_2)$ and it is trivial that whichever condition held for (\mathbf{a}, \mathbf{b}) is preserved.

Case (iii): since D is $\tau = N_2 - b_2$ level monotone for queue 2, by Definition 3.1, if (\mathbf{a}, \mathbf{b}) satisfies condition (I), then $I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_2\}} = I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_2\}} = I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_0\}} = 0$. If (\mathbf{a}, \mathbf{b}) satisfies condition (II), then there are two cases: if $a_1 < N_1 - 1$, then $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_1, \mathbf{b} + \mathbf{e}_2)$ and $N_1 - (a_1 - b_1) > N_1 - (N_1 - 1 - b_1) = b_1 + 1 \geq 1$, so that $a'_2 - b'_2 = a_2 - (b_2 + 1) \geq N_1 - (a_1 - b_1) - 1 = N_1 - (a'_1 - b'_1) > 0$, and condition (II) is preserved; if $a_1 = N_1 - 1$, then $(\mathbf{a}', \mathbf{b}') = ((0, a_2), \mathbf{b} + \mathbf{e}_2)$ and $a'_2 - b'_2 = a_2 - (b_2 + 1) \geq N_1 - (N_1 - 1 - b_1) - 1 = b_1 = b'_1 - a'_1 \geq 0$, so that condition (I) is satisfied.

Case (iv): if $a_2 = N_2 - 1$, then $\mathbf{a}' = (0, a_2) \in H_2$ and $T_2^Y(\mathbf{a}, \mathbf{b}) \leq T_2^W(\mathbf{a}, \mathbf{b})$. If $b_2 < a_2 < N_2 - 1$, suppose (\mathbf{a}, \mathbf{b}) satisfies condition (I), if $b_1 < N_1 - 1$, $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_2, \mathbf{b} + \mathbf{e}_1)$, then $a'_2 - b'_2 \geq b_1 - a_1 + 1 = b'_1 - a'_1 \geq 0$, thus condition (I) is preserved; if $b_1 = N_1 - 1$, $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_2, (0, b_2))$, then $a'_2 - b'_2 = (a_2 + 1) - b_2 \geq (N_1 - 1) - a_1 + 1 = N_1 - (a_1 - 0) = N_1 - (a'_1 - b'_1) > 0$, so condition (II) is satisfied. Suppose that (\mathbf{a}, \mathbf{b}) satisfies condition (II), $(\mathbf{a}', \mathbf{b}') = (\mathbf{a} + \mathbf{e}_2, \mathbf{b} + \mathbf{e}_1)$, then $a'_2 - b'_2 = (a_2 + 1) - b_2 \geq N_1 - (a_1 - b_1) + 1 = N_1 - (a'_1 - b'_1) > 0$, so condition (II) is preserved.

In summary, if (\mathbf{a}, \mathbf{b}) satisfies condition (I) or (II), so does the resulting state $(\mathbf{a}', \mathbf{b}')$ and we observe that any state (\mathbf{a}, \mathbf{b}) satisfying condition (I) or (II) also

automatically satisfies $a_2 \geq b_2$. Thus we have shown that if $(Y, W)(0) = (\mathbf{u}, \mathbf{v})$ satisfies either condition (I) or (II) then for each t , $0 \leq t < T_2^Y(\mathbf{u}, \mathbf{v}; 0)$, the joint process $(Y, W)(t)$ satisfies either condition (I) or (II), and hence $Y_2(t) \geq W_2(t)$ for $0 \leq t < T_2^Y(\mathbf{u}, \mathbf{v}; 0)$ with probability 1, which implies that $T_2^Y(\mathbf{u}, \mathbf{v}; 0) \leq T_2^W(\mathbf{u}, \mathbf{v}; 0)$ almost surely, i.e. $T_2(\mathbf{u}) \stackrel{st}{\leq} T_2(\mathbf{v})$.

Let $T_Y = T_2^Y(\mathbf{u}, \mathbf{v}; 0)$ and $T_W = T_2^W(\mathbf{u}, \mathbf{v}; 0)$ and consider the case where $u_2 > v_2$. We observe that there is a positive probability that the next $N_2 - u_2$ customers are all intrinsic arrivals to queue 2. Since Y has more customers in queue 2, with positive probability it will complete the queue 2 batch before W does, thus $\mathbb{P}(T_Y < T_W) > 0$. Hence $m_2(\mathbf{u}) = \mathbb{E}[T_Y] < \mathbb{E}[T_W] = m_2(\mathbf{v})$. ■

Lemma 3.5. *Let $\mathbf{a} = (a_1, a_2)$, $\mathbf{a}' = (a_1, a'_2) \in \mathbf{S}$ be such that $a_2 > a'_2$. Suppose that $D \in \mathcal{D}$ is τ -level monotone for queue 2, where $\tau = N_2 - (a'_2 + 1)$, then $z_{2,D}(\mathbf{a}) < z_{2,D}(\mathbf{a}')$.*

Proof. When $a_2 = N_2 - 1$, by (2) and (3), we have $z_2(\mathbf{a}') = m_2(\mathbf{a}' + \mathbf{e}_2) + \frac{1}{\mu_2} > \frac{1}{\mu_2} = z_2(\mathbf{a})$. When $a'_2 < a_2 < N_2 - 1$, condition (I) holds for $(\mathbf{a}, \mathbf{a}')$ and so by Lemma 3.4 we have $z_2(\mathbf{a}) = m_2(\mathbf{a} + \mathbf{e}_2) + 1/\mu_2 < m_2(\mathbf{a}' + \mathbf{e}_2) + 1/\mu_2 = z_2(\mathbf{a}')$. ■

Remark: Both Lemma 3.4 and 3.5 refer to queue 2, however corresponding results hold for queue 1 by relabeling the queues.

If in Lemma 3.4, we strengthen the conditions on the starting states, then the assumption of monotonicity is no longer required as per the following lemma.

Lemma 3.6. *Let $D \in \mathcal{D}$. Then $z_{1,D}(\mathbf{u}) \geq z_{1,D}(\mathbf{v})$ and $z_{2,D}(\mathbf{u}) \leq z_{2,D}(\mathbf{v})$ if either of the following hold:*

- (A) $u_2 - v_2 = v_1 - u_1 \geq 0$;
- (B) $u_2 - v_2 = N_1 - (u_1 - v_1) > 0$.

Furthermore, if $u_2 - v_2 = v_1 - u_1 > 0$, then the inequalities are strict.

Proof. We omit the details of this proof as it follows from a similar coupling argument as in Lemma 3.4. ■

Lemma 3.7. *Let $k, l, i \in \{1, 2\}$, $k \neq l$. For any user equilibrium policy $D^* \in \mathcal{D}^*$, each of the following conditions implies $\mathbf{n} \in D_i^*$.*

1. $\mathbf{n} + \mathbf{e}_k, \mathbf{n} + \mathbf{e}_l \in D_i^*$;
2. $\mathbf{n} + \mathbf{e}_k \in D_0^*, \mathbf{n} + \mathbf{e}_l \in D_i^*$;

Proof. Without loss of generality, assume $i = 1$. From Definition 2.5, the user equilibrium decision at a given state $\mathbf{n} \in \mathbf{S}$ depends on the expected transit times via each of the two queues.

Condition 1: By (4) and (5):

$$\begin{aligned} \Lambda \cdot z_{1,D^*}(\mathbf{n}) &= 1 + (\lambda + \sigma_1)z_{1,D^*}(\mathbf{n} + \mathbf{e}_1) + \sigma_2 \cdot z_{1,D^*}(\mathbf{n} + \mathbf{e}_2), \\ \Lambda \cdot z_{2,D^*}(\mathbf{n}) &= 1 + (\lambda + \sigma_1)z_{2,D^*}(\mathbf{n} + \mathbf{e}_1) + \sigma_2 \cdot z_{2,D^*}(\mathbf{n} + \mathbf{e}_2). \end{aligned}$$

Since $\mathbf{n} + \mathbf{e}_i \in D_1^* \Rightarrow z_{1;D^*}(\mathbf{n} + \mathbf{e}_i) \leq z_{2;D^*}(\mathbf{n} + \mathbf{e}_i)$ for $i = 1, 2$, we thus have $z_{1;D^*}(\mathbf{n}) \leq z_{2;D^*}(\mathbf{n})$, with strict inequality if any of the above are strict. By Definition 2.5, this implies that $\mathbf{n} \in D_1^*$.

Condition 2: If $k = 1$ and $l = 2$, that is $\mathbf{n} + \mathbf{e}_1 \in D_0^*$, $\mathbf{n} + \mathbf{e}_2 \in D_1^*$, then:

$$\begin{aligned} \Lambda \cdot z_{1;D^*}(\mathbf{n}) &= 1 + (\tilde{p}\lambda + \sigma_1)z_{1;D^*}(\mathbf{n} + \mathbf{e}_1) + [(1 - \tilde{p})\lambda + \sigma_2]z_{1;D^*}(\mathbf{n} + \mathbf{e}_2), \\ \Lambda \cdot z_{2;D^*}(\mathbf{n}) &= 1 + (\lambda + \sigma_1)z_{2;D^*}(\mathbf{n} + \mathbf{e}_1) + \sigma_2 z_{2;D^*}(\mathbf{n} + \mathbf{e}_2) \\ &= 1 + \underbrace{(\tilde{p}\lambda + \sigma_1)z_{2;D^*}(\mathbf{n} + \mathbf{e}_1) + [(1 - \tilde{p})\lambda + \sigma_2]z_{2;D^*}(\mathbf{n} + \mathbf{e}_2)}_A \\ &\quad + (1 - \tilde{p})\lambda[z_{2;D^*}(\mathbf{n} + \mathbf{e}_1) - z_{2;D^*}(\mathbf{n} + \mathbf{e}_2)] \end{aligned}$$

Since $\mathbf{n} + \mathbf{e}_1 \in D_0^* \Rightarrow z_{2;D^*}(\mathbf{n} + \mathbf{e}_1) = z_{1;D^*}(\mathbf{n} + \mathbf{e}_1)$ and $\mathbf{n} + \mathbf{e}_2 \in D_1^* \Rightarrow z_{2;D^*}(\mathbf{n} + \mathbf{e}_2) > z_{1;D^*}(\mathbf{n} + \mathbf{e}_2)$; and $z_{2;D^*}(\mathbf{n} + \mathbf{e}_1) > z_{2;D^*}(\mathbf{n} + \mathbf{e}_2)$ by Lemma 3.6, we thus have $z_{2;D^*}(\mathbf{n}) > A > z_{1;D^*}(\mathbf{n})$. This implies that $\mathbf{n} \in D_1^*$ by Definition 2.5. By symmetry, it is easy to show that the argument also holds for $k = 2$ and $l = 1$. \blacksquare

We now define the *switching curve* for queue i under policy D , $i = 1, 2$. Switching curves often arise naturally when seeking system optimal policies (see, for instance, Stidham and Weber [27] for more information and references therein) and we will also find them useful here.

Definition 3.8. For $D \in \mathcal{D}$, let

$$\begin{aligned} b_D^{(1)}(n_1) &= \sup\{n_2 : \mathbf{n} \in D_1\}, \quad 0 \leq n_1 \leq N_1 - 1, \\ b_D^{(2)}(n_2) &= \sup\{n_1 : \mathbf{n} \in D_2\}, \quad 0 \leq n_2 \leq N_2 - 1, \end{aligned}$$

and let

$$\mathbf{b}_D^{(i)} = \{b_D^{(i)}(n_i), 0 \leq n_i \leq N_i - 1\}$$

be the switching curve for queue i under policy D , $i = 1, 2$.

The following is an easy consequence of Definition 3.1 and 3.8. A similar result can be derived for $b_D^{(1)}(n_1)$.

Lemma 3.9. A decision policy $D \in \mathcal{D}$ is m -level monotone for queue 2 if and only if $b_D^{(2)}(n_2 - 1) \leq b_D^{(2)}(n_2)$, for all n_2 such that $N_2 - m < n_2 \leq N_2 - 1$; and $\mathbf{n} \in D_2$, $\forall \mathbf{n} \in \mathbf{S}$ such that $N_2 - m \leq n_2 \leq N_2 - 1$, $0 \leq n_1 < b_D^{(2)}(n_2)$.

3.2 Proof of main results

We are now ready to prove the main results of this paper, starting with Theorem 2.8, which says that for fixed parameter values, any user equilibrium is monotone.

Proof of Theorem 2.8. Let D^* be a user equilibrium for (N_1, N_2, Γ) . If $\mu_1 < \mu_2$ then $\forall n_1 < N_1 - 1$, it is trivial that $z_{1;D^*}(n_1, N_2 - 1) > z_{1;D^*}(N_1 - 1, N_2 - 1) = 1/\mu_1 > 1/\mu_2 = z_{2;D^*}(n_1, N_2 - 1)$ (see e.g. (3)). Thus $(n_1, N_2 - 1) \in D_2^*$ by Definition 2.5, and D^* is 1-level monotone for queue 2. Moreover $b_{D^*}^{(2)}(N_2 - 1) = N_1 - 1$. Similarly, if

$\mu_1 > \mu_2$ then $(N_1 - 1, n_2) \in D_1^*$ for $n_2 \leq N_2 - 1$, D^* is 1-level monotone with respect to queue 1, and $b_D^{(1)}(N_1 - 1) = N_2 - 1$. Now if $\mu_1 = \mu_2$ then $(N_1 - 1, N_2 - 1) \in D_0^*$ as $1/\mu_1 = 1/\mu_2$; and $(N_1 - 1, n_2) \in D_1^*$ for $n_2 \leq N_2 - 2$, and $(n_1, N_2 - 1) \in D_2^*$ for $n_1 \leq N_1 - 2$, D^* is thus 1-level monotone with respect to queue 1 and queue 2 respectively. That is $b_D^{(1)}(N_1 - 1) = N_2 - 2$ and $b_D^{(2)}(N_2 - 1) = N_1 - 2$.

Suppose that D^* is m -level monotone for queue 2 for some $m > 0$ (the proof is analogous if we assume that D^* is n -level monotone for queue 1). If $b_{D^*}^{(2)}(N_2 - m) = N_1 - 1$ then trivially $b_{D^*}^{(2)}(N_2 - m - 1) \leq b_{D^*}^{(2)}(N_2 - m)$. If $b_{D^*}^{(2)}(N_2 - m) < N_1 - 1$, by Definition 3.8, $(n_1, N_2 - m) \in D_1^* \cup D_0^*$ for n_1 such that $b_{D^*}^{(2)}(N_2 - m) < n_1 \leq N_1 - 1$. In particular, $(N_1 - 1, N_2 - m) \in D_1^* \cup D_0^*$. Then by Lemma 3.5, $z_1(N_1 - 1, N_2 - m - 1) = z_1(N_1 - 1, N_2 - m) \leq z_2(N_1 - 1, N_2 - m) < z_2(N_1 - 1, N_2 - m - 1) \Rightarrow (N_1 - 1, N_2 - m - 1) \in D_1^*$. Thus $b_{D^*}^{(2)}(N_2 - m - 1) \leq b_{D^*}^{(2)}(N_2 - m)$ follows from Lemma 3.7. Furthermore, since $(b_{D^*}^{(2)}(N_2 - m - 1), N_2 - m - 1) \in D_2^*$, we also have $(n_1, N_2 - m - 1) \in D_2^*$ for $0 \leq n_1 < b_{D^*}^{(2)}(N_2 - m - 1)$ by Lemma 3.7. Thus D^* is $(m + 1)$ -level monotone with respect to queue 2 (Lemma 3.9).

By induction on m , D^* is N_2 -level monotone for queue 2. Observe that we also immediately have $b_{D^*}^{(1)}(n_1 - 1) \leq b_{D^*}^{(1)}(n_1)$ for $0 < n_1 \leq N_1 - 1$, and $\mathbf{n} \in D_1^*$ for \mathbf{n} such that $0 \leq n_2 < b_{D^*}^{(1)}(n_1)$. This implies that D^* is also N_1 -level monotone. ■

Definition 3.10. A policy $D^* \in \mathcal{D}$ is a (ξ, τ) -level user equilibrium policy for (N_1, N_2, Γ) if $\forall \mathbf{n} \in \mathbf{S}$ such that $n_1 \geq N_1 - \xi$ and $n_2 \geq N_2 - \tau$,

$$\mathbf{n} \in \begin{cases} D_1^*, & \iff z_{1;D^*}(\mathbf{n}) < z_{2;D^*}(\mathbf{n}) \\ D_0^*, & \iff z_{1;D^*}(\mathbf{n}) = z_{2;D^*}(\mathbf{n}) \\ D_2^*, & \iff z_{1;D^*}(\mathbf{n}) > z_{2;D^*}(\mathbf{n}). \end{cases}$$

Let $\mathcal{D}_{\xi;\tau}^* = \mathcal{D}_{\xi;\tau}^*(N_1, N_2, \Gamma)$ denote the class of all (ξ, τ) -level user equilibrium policies.

Clearly a policy D^* is a user equilibrium policy (Definition 2.5) for (N_1, N_2, Γ) if and only if it is an (N_1, N_2) -level user equilibrium for Γ . Let $\mathcal{D}^* = \mathcal{D}^*(N_1, N_2, \Gamma)$ be the class of all user equilibrium polices.

Proof of Theorem 2.7. We construct a user equilibrium policy D' , via a finite sequence of elements $\mathbf{s}^{(r)}$ of \mathbf{S} and corresponding subpolicies on subsets $\mathbf{S}^{(r)}$ of \mathbf{S} . Let $\mathbf{s}^{(1)} = (N_1 - 1, N_2 - 1)$. If $\mu_1 \geq \mu_2$, we put $(N_1 - 1, n_2) \in D'_1$ for $0 \leq n_2 \leq N_2 - 2$; if the inequality is strict, we put $s^{(1)} \in D'_1$, otherwise $s^{(1)} \in D'_0$. Now define $\mathbf{S}^{(1)} = \{N_1 - 1\} \times \{0, \dots, N_2 - 1\}$ and $\mathbf{s}^{(2)} = (N_1 - 2, N_2 - 1)$. If $\mu_1 < \mu_2$ and we put $(n_1, N_2 - 1) \in D'_2$ for each n_1 and define $\mathbf{S}^{(1)} = \{0, \dots, N_1 - 1\} \times \{N_2 - 1\}$ and $\mathbf{s}^{(2)} = (N_1 - 1, N_2 - 2)$. We proceed iteratively. Suppose that we have already defined $\mathbf{s}^{(r)} = \mathbf{s}$ and specified D' on

$$\mathbf{S}^{(r)} = \{ \{s_1 + 1, \dots, N_1 - 1\} \times \{0, \dots, N_2 - 1\} \} \cup \{ \{0, \dots, N_1 - 1\} \times \{s_2 + 1, \dots, N_2 - 1\} \}.$$

As in Equations (2), (4) and (5) we can compute $z_{1;D'}(\mathbf{s})$ and $z_{2;D'}(\mathbf{s})$ based on $\{D'(m, n) : (m, n) \in \mathbf{S}^{(r)}\}$. If $z_{1;D'}(\mathbf{s}) \leq z_{2;D'}(\mathbf{s})$ then we put $(s_1, n_2) \in D'_1$ for every $n_2 < s_2$; and if the inequality is strict, we put $\mathbf{s} \in D'_1$, otherwise, put $\mathbf{s} \in D'_0$. Then we define $\mathbf{S}^{(r+1)} = \mathbf{S}^{(r)} \cup \{(s_1, n_2) : n_2 \leq s_2\}$ and $\mathbf{s}^{(r+1)} = (s_1 - 1, s_2)$. If

$z_{1;D'}(\mathbf{s}) > z_{2;D'}(\mathbf{s})$ and we put $(n_1, s_2) \in D'_2$ for every $n_1 \leq s_1$, and define $\mathbf{S}^{(r+1)} = \mathbf{S}^{(r)} \cup \{(n_1, s_2) : n_1 \leq s_1\}$ and $\mathbf{s}^{(r+1)} = (s_1, s_2 - 1)$. We stop as soon as $(0, 0) \in \mathbf{S}^{(r)}$ for some r (at which point $\mathbf{S}^{(r)} = \mathbf{S}$).

Observe that for any (N_1, N_2, Γ) , the construction above gives rise to a single decision policy D' that is (N_1, N_2) -level monotone. Although we have not yet shown that D' is a user equilibrium, any user equilibrium for Γ must be equal to D' (i.e. the user equilibrium must be unique) as follows. If D^* is a user equilibrium, then it must agree with D' at $\mathbf{s}^{(1)}$ by Definition 3.10 and it must also agree with D' on $\mathbf{S}^{(1)}$ by Theorem 2.8. Proceeding iteratively, D^* must agree with D' at each $\mathbf{s}^{(r)}$ by Definition 3.10 as $\mathbf{s}^{(r)} \in D'_1 \Leftrightarrow z_{1;D'}(\mathbf{s}^{(r)}) < z_{2;D'}(\mathbf{s}^{(r)})$ and $\mathbf{s}^{(r)} \in D'_2 \Leftrightarrow z_{1;D'}(\mathbf{s}^{(r)}) > z_{2;D'}(\mathbf{s}^{(r)})$ by construction, and therefore D^* agrees with D' on each $\mathbf{S}^{(r)}$ by Theorem 2.8, whence $D^* = D'$.

It remains to prove that D' is a user equilibrium, i.e. that $\mathbf{n} \in D'_1 \Leftrightarrow z_{1;D'}(\mathbf{n}) < z_{2;D'}(\mathbf{n})$, and $\mathbf{n} \in D'_2 \Leftrightarrow z_{1;D'}(\mathbf{n}) > z_{2;D'}(\mathbf{n})$, $\forall \mathbf{n} \in \mathbf{S}$. We already know this for states $\mathbf{s}^{(r)}$. Let $\mathbf{s} = \mathbf{s}^{(r)}$. If $s_1 = N_1 - 1$ and $\mathbf{s} \in D'_1 \cup D'_0$, then $z_{1;D'}(\mathbf{s} - k\mathbf{e}_2) = z_{1;D'}(\mathbf{s}) \leq z_{2;D'}(\mathbf{s}) < z_{2;D'}(\mathbf{s} - k\mathbf{e}_2)$ for $k \geq 1$ by Lemma 3.5. Similarly if $s_1 = N_1 - 1$ and $\mathbf{s} \in D'_2$, $z_{1;D'}(\mathbf{s} - k\mathbf{e}_1) > z_{2;D'}(\mathbf{s} - k\mathbf{e}_1)$.

Assume that D' is an $(N_1 - s_1 - 1, N_2 - s_2 - 1)$ -level user equilibrium. If $s_1 < N_1 - 1$ and $\mathbf{s} \in D'_1 \cup D'_0$ then we also have $\mathbf{s} - \mathbf{e}_2 + \mathbf{e}_1 \in D'_1$, since D' is monotone by construction. We can then apply the argument of Lemma 3.7 to obtain $z_{1;D'}(\mathbf{s} - \mathbf{e}_2) < z_{2;D'}(\mathbf{s} - \mathbf{e}_2)$. Continuing inductively we obtain $z_{1;D'}(\mathbf{s} - k\mathbf{e}_2) < z_{2;D'}(\mathbf{s} - k\mathbf{e}_2)$ for $k = 0, 1, \dots, s_2$. Thus $\mathbf{s} - k\mathbf{e}_2 \in D'_1 \Leftrightarrow z_{1;D'}(\mathbf{s} - k\mathbf{e}_2) < z_{2;D'}(\mathbf{s} - k\mathbf{e}_2)$, i.e. D' is also an $(N_1 - s_1, N_2 - s_2 - 1) = (N_1 - s_1^{(r+1)} - 1, N_2 - s_2^{(r+1)} - 1)$ -level user equilibrium, where $\mathbf{s}^{(r+1)} = (s_1 - 1, s_2)$ since $\mathbf{s} \in D'_1$. A similar argument also holds for $\mathbf{s} \in D'_2$. By induction on r , D' satisfies Definition 2.5 and is therefore a user equilibrium for the model. \blacksquare

These two theorems show that a unique user equilibrium exists, and that it is monotone. These results depend very much on the assumed service discipline. In [2] similar results were obtained for a system with one batch service queue, and one $M|M|1$ queue. However, in [14], the authors study a similar system with parallel processor sharing queues, where existence and uniqueness issues become more delicate. In that case even when a unique user equilibrium exists, a direct construction is not possible, and an iterative algorithm starting from an initial policy is needed to find the user equilibrium.

3.3 Monotonicity of D^* with respect to parameters: an iterative approach

In this section, we adapt the coupling technique of Lemma 3.4 to show that, under a fixed decision policy D , the expected transit time $z_{i;D}(\mathbf{n})$ via queue $i = 1, 2$, is monotone with respect to changes in intrinsic arrival rates and service rates of the queues. It is also of interest to consider monotonicity properties of the user equilibria with respect to parameters. In Theorem 3.15, we use an iterative approach to establish such a property with respect to service rates.

Given a decision policy $D \in \mathcal{D}$. We write $Z_D^{(i)} = \{Z_D^{(i)}(t)\}_{t \geq 0}$ for the process operating under decision policy D with parameters $\Gamma^{(i)}$, $i = 1, 2$, and let $z_2^{(i)}(\mathbf{n}) =$

$z_{2;D}(\mathbf{n}, \Gamma^{(i)})$.

Lemma 3.11. *Let $\Gamma^{(i)} = \{\lambda, \sigma_1, \mu_1, \sigma_2^{(i)}, \mu_2^{(i)}, \tilde{p}\}$ for $i = 1, 2$, with $\sigma_2^{(1)} \geq \sigma_2^{(2)}$ and $\mu_2^{(1)} \geq \mu_2^{(2)}$. If D is $\tau = N_2 - (n_2 + 1)$ level monotone for queue 2, then $z_2^{(1)}(\mathbf{n}) \leq z_2^{(2)}(\mathbf{n})$.*

Proof. We construct a joint process $\{(Y, W)(t)\}_{t \geq 0}$ on the state space $\mathbf{S} \times \mathbf{S}$ such that Y and W follow the laws of $Z_D^{(1)}$ and $Z_D^{(2)}$ respectively. Let $(Y, W)(0) = ((Y_1, Y_2), (W_1, W_2))(0) = (\mathbf{u}, \mathbf{v})$, where \mathbf{u}, \mathbf{v} satisfies either of conditions (I) and (II) in Lemma 3.4. The transitions of the joint process are as follows: $\forall (\mathbf{a}, \mathbf{b}) \in \mathbf{S} \times \mathbf{S}$,

(i) At rate $\sigma_1 + \lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_0\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

(ii) At rate $\sigma_2^{(2)} + \lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_0\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a}, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

(iii) At rate $\sigma_2^{(1)} - \sigma_2^{(2)}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b}).$$

(iv) At rate $\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_2\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_2\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_1, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

(v) At rate $\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_1\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0, \mathbf{b} \in D_1\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_2, \mathbf{b} \in D_0\}}$,

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

Denote by $T_2^X(\mathbf{a}, \mathbf{b}; 0)$ the marginal hitting time to the set H_2 for the system $X \in \{Y, W\}$, starting from state (\mathbf{a}, \mathbf{b}) . Let $(\mathbf{a}', \mathbf{b}')$ denote the state of the system immediately after a transition out of state $(\mathbf{a}, \mathbf{b}) \in \mathbf{S} \times \mathbf{S}$.

For cases (i), (ii), (iv) and (v), the discussion in Lemma 3.4 has shown that if (\mathbf{a}, \mathbf{b}) satisfies either condition (I) or (II), so does the resulting state $(\mathbf{a}', \mathbf{b}')$. For case (iii), if $a_2 = N_2 - 1$, $\mathbf{a}' = (a_1, 0) \in H_2$, so that $T_2^Y(\mathbf{a}, \mathbf{b}; 0) \leq T_2^W(\mathbf{a}, \mathbf{b}; 0)$ holds. If $a_2 < N_2 - 1$, then $(\mathbf{a}', \mathbf{b}') = ((a_1, a_2 + 1), \mathbf{b})$ and $(\mathbf{a}', \mathbf{b}')$ preserves whichever condition was satisfied by (\mathbf{a}, \mathbf{b}) .

Therefore, if (\mathbf{a}, \mathbf{b}) satisfies either of conditions (I) and (II), so does $(\mathbf{a}', \mathbf{b}')$. In our case, since $u_2 \leq v_2$, then $Y_2(t) \geq W_2(t)$ for all $t \in [0, T_2^Y(\mathbf{u}; 0)]$. Thus $T_2^Y(\mathbf{u}, \mathbf{v}; 0) \leq T_2^W(\mathbf{u}, \mathbf{v}; 0)$ almost surely, so $T_2^{(1)}(\mathbf{u}) \stackrel{st}{\leq} T_2^{(2)}(\mathbf{v})$. If $n_2 = N_2 - 1$, by (3), $z_2^{(1)}(\mathbf{n}) = 1/\mu_2^{(1)} \leq 1/\mu_2^{(2)} = z_2^{(2)}(\mathbf{n})$. If $n_2 < N_2 - 1$, let $\mathbf{u} = \mathbf{v} = (\mathbf{n} + \mathbf{e}_2)$ so that (\mathbf{u}, \mathbf{v}) satisfies condition (I) in Lemma 3.4, then $T_2^{(1)}(\mathbf{n} + \mathbf{e}_2) \stackrel{st}{\leq} T_2^{(2)}(\mathbf{n} + \mathbf{e}_2)$. Hence $z_2^{(1)}(\mathbf{n}) = \mathbb{E}[T_2^{(1)}(\mathbf{n} + \mathbf{e}_2)] + 1/\mu_2^{(1)} \leq \mathbb{E}[T_2^{(2)}(\mathbf{n} + \mathbf{e}_2)] + 1/\mu_2^{(2)} = z_2^{(2)}(\mathbf{n})$. \blacksquare

Lemma 3.11 compared expected transit times under a fixed policy with different parameters. The following result compares expected transit times for two different policies with the parameters fixed.

Lemma 3.12. *Let $D^{(i)} \in \mathcal{D}$, $i = 1, 2$, be two monotone decision policies for (N_1, N_2, Γ) such that $D_2^{(2)} \subseteq D_2^{(1)}$ and $D_1^{(2)} \supseteq D_1^{(1)}$. Then $z_{1;D^{(1)}}(\mathbf{n}) \geq z_{1;D^{(2)}}(\mathbf{n})$ and $z_{2;D^{(1)}}(\mathbf{n}) \leq z_{2;D^{(2)}}(\mathbf{n})$ for $\mathbf{n} \in \mathbf{S}$.*

Proof. We consider the joint process of the two systems $\{(Y, W)(t)\}_{t \geq 0}$ on $\mathbf{S} \times \mathbf{S}$ where Y and W follow the laws of $Z_{D^{(1)}}$ and $Z_{D^{(2)}}$ respectively. Let $(Y, W)(0) = (\mathbf{u}, \mathbf{v})$ where (\mathbf{u}, \mathbf{v}) satisfies either condition (I) or (II) in Lemma 3.4. The transitions out of a state $(\mathbf{a}, \mathbf{b}) \in \mathbf{S} \times \mathbf{S}$ are:

$$(i) \text{ At rate } \sigma_1 + \lambda \cdot I_{\{\mathbf{a} \in D_1^{(1)}, \mathbf{b} \in D_1^{(2)}\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_1^{(2)}\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_1^{(1)}, \mathbf{b} \in D_0^{(2)}\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_0^{(2)}\}},$$

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

$$(ii) \text{ At rate } \sigma_2 + \lambda \cdot I_{\{\mathbf{a} \in D_2^{(1)}, \mathbf{b} \in D_2^{(2)}\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_2^{(2)}\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_2^{(1)}, \mathbf{b} \in D_0^{(2)}\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_0^{(2)}\}},$$

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

$$(iii) \text{ At rate } \lambda \cdot I_{\{\mathbf{a} \in D_1^{(1)}, \mathbf{b} \in D_2^{(2)}\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_2^{(2)}\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_1^{(1)}, \mathbf{b} \in D_0^{(2)}\}},$$

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_1, \mathbf{b} + \tilde{\mathbf{e}}_2).$$

$$(iv) \text{ At rate } \lambda \cdot I_{\{\mathbf{a} \in D_2^{(1)}, \mathbf{b} \in D_1^{(2)}\}} + (1 - \tilde{p})\lambda \cdot I_{\{\mathbf{a} \in D_0^{(1)}, \mathbf{b} \in D_1^{(2)}\}} + \tilde{p}\lambda \cdot I_{\{\mathbf{a} \in D_2^{(1)}, \mathbf{b} \in D_0^{(2)}\}},$$

$$(\mathbf{a}, \mathbf{b}) \longrightarrow (\mathbf{a} + \tilde{\mathbf{e}}_2, \mathbf{b} + \tilde{\mathbf{e}}_1).$$

Let $T_2^{(i)}(\mathbf{n}) = T_{2;D^{(i)}}^Z(\mathbf{n})$ be the waiting time via queue 2 for a general commuter who upon arrival sees system $Z_{D^{(i)}}$ in state \mathbf{n} . Since the possible transitions are the same as in Lemma 3.4, we can easily show that $T_2^{(1)}(\mathbf{u}) \stackrel{st}{\leq} T_2^{(2)}(\mathbf{v})$ if (\mathbf{u}, \mathbf{v}) satisfies either condition (I) or (II) in Lemma 3.4. If $n_2 = N_2 - 1$, then $z_{2;D^{(1)}}(\mathbf{n}) = z_{2;D^{(2)}}(\mathbf{n}) = 1/\mu_2$. If $n_2 < N_2 - 1$, let $\mathbf{u} = \mathbf{v} = \mathbf{n} + \mathbf{e}_2 \in \mathbf{S}$ and then, since condition (I) is satisfied, $T_2^{(1)}(\mathbf{n} + \mathbf{e}_2) \stackrel{st}{\leq} T_2^{(2)}(\mathbf{n} + \mathbf{e}_2)$, and then applying (3), $z_{2;D^{(1)}}(\mathbf{n}) = \mathbb{E}[T_2^{(1)}(\mathbf{n} + \mathbf{e}_2)] + 1/\mu_2 \leq \mathbb{E}[T_2^{(2)}(\mathbf{n} + \mathbf{e}_2)] + 1/\mu_2 = z_{2;D^{(2)}}(\mathbf{n})$. The first claim follows by interchanging the labels of the queues. \blacksquare

The following result follows immediately from Lemmas 3.11 and 3.12.

Lemma 3.13. *For $i = 1, 2$, let $\Gamma^{(i)} = \{\lambda, \sigma_1, \mu_1, \sigma_2^{(i)}, \mu_2^{(i)}\}$ be such that $\sigma_2^{(1)} \geq \sigma_2^{(2)}$ and $\mu_2^{(1)} \geq \mu_2^{(2)}$, and let $D^{(i)} = (D_1^{(i)}, D_2^{(i)}) \in \mathcal{D}$ be monotone decision policies such that $D_2^{(2)} \subseteq D_2^{(1)}$ and $D_1^{(2)} \supseteq D_1^{(1)}$. Then $z_{2;D^{(1)}}(\mathbf{n}; \Gamma^{(1)}) \leq z_{2;D^{(2)}}(\mathbf{n}; \Gamma^{(2)})$.*

The following lemma shows that any decision policy $D \in \mathcal{D}$ can be “updated” to a monotone decision policy.

Lemma 3.14. *Given a decision policy $D \in \mathcal{D}$, and $i, k, l \in \{1, 2\}$ such that $k \neq l$, let D' be the policy such that $\mathbf{n} \in D'_1 \Leftrightarrow z_{1;D}(\mathbf{n}) < z_{2;D}(\mathbf{n})$ and $\mathbf{n} \in D'_0 \Leftrightarrow z_{1;D}(\mathbf{n}) = z_{2;D}(\mathbf{n})$, $\forall \mathbf{n} \in \mathbf{S}$. Then each of the following implies that $\mathbf{n} \in D'_i$.*

$$A. \mathbf{n} + \mathbf{e}_k, \mathbf{n} + \mathbf{e}_l \in D'_i,$$

B. $\mathbf{n} + \mathbf{e}_k \in D'_0$, $\mathbf{n} + \mathbf{e}_l \in D'_i$.

Furthermore, D' is (N_1, N_2) -level monotone.

Proof. Without loss of generality, let $i = 1$ (a similar argument will also hold for $i = 2$).

There are nine cases to consider, depending on whether $\mathbf{n} + \mathbf{e}_r \in D_s$, $r = 1, 2$ and $s = 0, 1, 2$. We show the proof for the case $\mathbf{n} + \mathbf{e}_1 \in D_2$, $\mathbf{n} + \mathbf{e}_2 \in D_1$, and omit the details for the rest of the eight cases as they can be proved in a very similar manner.

By (4) and (5),

$$\begin{aligned} \Lambda \cdot z_{1;D}(\mathbf{n}) &= 1 + \sigma_1 z_{1;D}(\mathbf{n} + \mathbf{e}_1) + (\lambda + \sigma_2) z_{1;D}(\mathbf{n} + \mathbf{e}_2). \\ &= 1 + \underbrace{(\lambda + \sigma_1) z_{1;D}(\mathbf{n} + \mathbf{e}_1) + \sigma_2 z_{1;D}(\mathbf{n} + \mathbf{e}_2)}_B + \underbrace{\lambda [z_{1;D}(\mathbf{n} + \mathbf{e}_2) - z_{1;D}(\mathbf{n} + \mathbf{e}_1)]}_{>0 \text{ by Lemma 3.6}}. \end{aligned}$$

$$\begin{aligned} \Lambda \cdot z_{2;D}(\mathbf{n}) &= 1 + (\lambda + \sigma_1) z_{2;D}(\mathbf{n} + \mathbf{e}_1) + \sigma_2 z_{2;D}(\mathbf{n} + \mathbf{e}_2). \\ &= 1 + \underbrace{\sigma_1 z_{2;D}(\mathbf{n} + \mathbf{e}_1) + (\lambda + \sigma_2) z_{2;D}(\mathbf{n} + \mathbf{e}_2)}_C + \underbrace{\lambda [z_{2;D}(\mathbf{n} + \mathbf{e}_1) - z_{2;D}(\mathbf{n} + \mathbf{e}_2)]}_{>0 \text{ by Lemma 3.6}}. \end{aligned}$$

Suppose that $k = 1$ and $l = 2$. Then $\mathbf{n} + \mathbf{e}_1 \in D'_1$ (condition A) and $\mathbf{n} + \mathbf{e}_1 \in D'_0$ (condition B) imply that $z_{1;D}(\mathbf{n} + \mathbf{e}_1) \leq z_{2;D}(\mathbf{n} + \mathbf{e}_1)$. Also under both conditions, $\mathbf{n} + \mathbf{e}_2 \in D'_1 \Leftrightarrow z_{1;D}(\mathbf{n} + \mathbf{e}_2) < z_{2;D}(\mathbf{n} + \mathbf{e}_2)$. Thus $\Lambda z_{1;D}(\mathbf{n}) < C < \Lambda z_{2;D}(\mathbf{n})$, hence $\mathbf{n} \in D'_1$. Similarly, if $k = 2$ and $l = 1$, it is easy to show that $z_{2;D}(\mathbf{n}) > B > z_{1;D}(\mathbf{n})$, which also implies $\mathbf{n} \in D'_1$.

The proof of the monotonicity of D' is very similar to that of the Theorem 2.8. Thus, we have completed the proof. \blacksquare

We now adapt the policy updating idea of Lemma 3.14 to show that D^* is monotone with respect to the service rates.

Theorem 3.15. For $i = 1, 2$, let $\Gamma^{(i)} = \{\lambda, \sigma_1, \mu_1, \sigma_2, \mu_2^{(i)}, \tilde{p}\}$ be such that $\mu_2^{(1)} \geq \mu_2^{(2)}$. Suppose that $D^{*(i)} = \{D_0^{*(i)}, D_1^{*(i)}, D_2^{*(i)}\} \in \mathcal{D}^*$ is the user equilibrium policy for the process operating with parameter $\Gamma^{(i)}$, then $D_1^{*(1)} \subseteq D_1^{*(2)}$ and $D_2^{*(1)} \supseteq D_2^{*(2)}$.

Proof. Let $G^{(0)} = \{G_0^{(0)}, G_1^{(0)}, G_2^{(0)}\} = D^{*(1)}$. We first construct a sequence of decision policies $\{G^{(k)}, k \in \mathbb{N}\}$ as follows:

$$\begin{aligned} z_{1;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) < z_{2;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) &\Leftrightarrow \mathbf{n} \in G_1^{(k)}, \text{ and} \\ z_{1;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) > z_{2;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) &\Leftrightarrow \mathbf{n} \in G_2^{(k)}, \forall \mathbf{n} \in \mathbf{S}. \end{aligned}$$

By Lemma 3.14, the sequence $\{G^{(k)}, k \in \mathbb{N}\}$ exists and each decision policy is (N_1, N_2) -level monotone.

We claim that $G_1^{(k-1)} \subseteq G_1^{(k)}$ and $G_2^{(k-1)} \supseteq G_2^{(k)}$ for each $k \in \mathbb{N}$. Assuming this for the moment, since the number of decision policies $|\mathcal{D}|$ is finite, there must exist an $s \in \mathbb{N}$ such that $G^{(s-1)} = G^{(s)}$. Therefore, $\mathbf{n} \in G_1^{(s)} \Leftrightarrow z_{1;G^{(s)}}(\mathbf{n}; \Gamma^{(2)}) = z_{1;G^{(s-1)}}(\mathbf{n}; \Gamma^{(2)}) < z_{2;G^{(s-1)}}(\mathbf{n}; \Gamma^{(2)}) = z_{2;G^{(s)}}(\mathbf{n}; \Gamma^{(2)})$. Similarly, $\mathbf{n} \in G_2^{(s)} \Leftrightarrow z_{1;G^{(s)}}(\mathbf{n}; \Gamma^{(2)}) > z_{2;G^{(s)}}(\mathbf{n}; \Gamma^{(2)})$. Thus $G^{(s)}$ is a user equilibrium of the system by Definition 2.5.

Theorem 2.7 says that the user equilibrium of a given system is unique. Hence $D^{*(2)} = G^{(s)}$ and $D_1^{*(1)} \subseteq D_1^{*(2)}$ and $D_2^{*(1)} \supseteq D_2^{*(2)}$.

It therefore remains to verify that $G_1^{(k-1)} \subseteq G_1^{(k)}$ and $G_2^{(k-1)} \supseteq G_2^{(k)}$. Notice that for a given decision policy $G^{(0)}$, the service rate of queue 2 does not affect the expected time via queue 1. Thus $z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(1)}) = z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(2)})$. On the other hand, $z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(1)}) \leq z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(2)})$ by Lemma 3.11. Now if $\mathbf{n} \in G_1^{(0)}$, $z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(1)}) < z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(1)})$ as $G^{(0)}$ is the user equilibrium for the system with parameter $\Gamma^{(1)}$ (Definition 2.5). Thus $z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(2)}) < z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(2)}) \Rightarrow \mathbf{n} \in G_1^{(1)}$. Similarly, $\forall \mathbf{n} \in G_2^{(1)}$, $z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(1)}) = z_{1;G^{(0)}}(\mathbf{n}; \Gamma^{(2)}) > z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(2)}) \geq z_{2;G^{(0)}}(\mathbf{n}; \Gamma^{(1)})$. Hence, by Definition 2.5, this implies that $\mathbf{n} \in G_2^{(0)}$. This establishes the claim with $k = 1$. Proceeding by induction, let $k \geq 1$ and suppose that $G_1^{(k-1)} \subseteq G_1^{(k)}$, $G_2^{(k-1)} \supseteq G_2^{(k)}$. If $\mathbf{n} \in G_1^{(k)}$, then by Lemma 3.12 and the definition of $G^{(k)}$, $z_{1;G^{(k)}}(\mathbf{n}; \Gamma^{(2)}) \leq z_{1;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) < z_{2;G^{(k-1)}}(\mathbf{n}; \Gamma^{(2)}) \leq z_{2;G^{(k)}}(\mathbf{n}; \Gamma^{(2)})$. Hence $\mathbf{n} \in G_1^{(k+1)}$. So $G_1^{(k)} \subseteq G_1^{(k+1)}$, and by symmetry we also have $G_2^{(k)} \supseteq G_2^{(k+1)}$, as claimed. \blacksquare

Theorem 3.15 has proved that the user equilibrium policy is monotone with respect to the service rates. One would also expect a natural monotone property with respect to the intrinsic arrival rate, as per the following conjecture.

Conjecture. For $i = 1, 2$, let $\Gamma^{(i)} = \{\lambda, \sigma_1, \mu_1, \sigma_2^{(i)}, \mu_2, \tilde{p}\}$ be such that $\sigma_2^{(1)} \geq \sigma_2^{(2)}$. Suppose that $D^{*(i)} = \{D_0^{*(i)}, D_1^{*(i)}, D_2^{*(i)}\} \in \mathcal{D}^*$ be the user equilibrium policy for the process operating with parameter $\Gamma^{(i)}$, then $D_1^{*(1)} \subseteq D_1^{*(2)}$ and $D_2^{*(1)} \supseteq D_2^{*(2)}$.

4 Examples and discussion

We have shown that the state dependent user equilibrium policy possesses several natural monotonicity properties for the system under consideration. In this section we give some numerical examples illustrating these results, and comparing the expected transit time for a general customer under the user equilibrium policy with that obtained under probabilistic routing. The numerical examples below illustrate that performance under state dependent routing may be considerably better than under probabilistic routing, although it is not always the case that the expected delay under state dependent routing is lower than under probabilistic routing. This section concludes with some discussion and questions for further research.

Under state dependent routing the expected transit time in the stationary regime under the user equilibrium policy can be easily calculated. Let $D^* = (D_0^*, D_1^*, D_2^*) \in \mathcal{D}^*$ be the user equilibrium policy for the network with parameters (N_1, N_2, Γ) and $\vec{\pi}_{D^*} = \{\pi_{D^*}(\mathbf{n}), \mathbf{n} \in \mathbf{S}\}$ be the stationary distribution of the system under D^* . Since the state space is finite (and irreducible if all parameters are positive), $\vec{\pi}_{D^*}$ exists. Thus the expected transit time for a general customer in

the stationary regime under D^* is

$$\begin{aligned}
W_{\text{dep}}(D^*) &= \sum_{\mathbf{n} \in \mathbf{S}} \pi_{D^*}(\mathbf{n}) \left(z_{1;D^*}(\mathbf{n}) \cdot I_{\mathbf{n} \in D_1^*} + z_{2;D^*}(\mathbf{n}) \cdot I_{\mathbf{n} \in D_2^*} + \right. \\
&\quad \left. (\tilde{p} z_{1;D^*}(\mathbf{n}) + (1 - \tilde{p}) z_{2;D^*}(\mathbf{n})) \cdot I_{\mathbf{n} \in D_0^*} \right); \\
&= \sum_{\mathbf{n} \in \mathbf{S}} \pi_{D^*}(\mathbf{n}) \left(z_{1;D^*}(\mathbf{n}) \cdot I_{\mathbf{n} \in D_1^* \cup D_0^*} + z_{2;D^*}(\mathbf{n}) \cdot I_{\mathbf{n} \in D_2^*} \right), \\
&\quad \text{as } \mathbf{n} \in D_0^* \Leftrightarrow z_{1;D^*}(\mathbf{n}) = z_{2;D^*}(\mathbf{n}).
\end{aligned}$$

Figure 2 plots $W_{\text{dep}}(D^*)$ against μ_1 for $\lambda = 4, \sigma_1 = 3, \sigma_2 = 1, \mu_2 = 2, \tilde{p} = 1, N_1 = N_2 = 5$, with μ_1 varying from 0.9 to 1.6 in increments of 0.005. We see that $W_{\text{dep}}(D^*)$ is not, in general, decreasing in μ_1 . The increases in $W_{\text{dep}}(D^*)$ occur when D^* changes. An interesting observation here is that although for a fixed policy the individual expected transit time via queue i decreases as the service rate to that queue increases (Lemma 3.11) and the user equilibrium changes monotonically with respect to the service rates (Theorem 3.15), the overall expected transit for the system $W_{\text{dep}}(D^*)$ is not monotone here.

The example here exhibits similar behaviour to the well-known Downs-Thompson paradox (see for instance Downs [16], Thomson [28], Calvert [12], and Afimeimounga et al. [1, 2]). In a parallel system with a batch-service queue and a $M/M/1$ queue, the paradox arises when increasing the capacity of the $M/M/1$ queue leads to an increase in the delay for the overall system as commuters shift from one queue to the other. Similarly, in our case, as general customers decide to shift their choice of queue to reduce their individual delay, the system performance may get worse.

Now consider the same system of queues under probabilistic routing. Let $W_i(p)$ be the expected transit time in equilibrium for a general customer who travels via queue i , $i = 1, 2$. Then it is easily seen (see, e.g. [1]) that

$$W_1(p) = \frac{1}{\mu_1} + \frac{N_1 - 1}{2(\sigma_1 + p\lambda)}, \quad W_2(p) = \frac{1}{\mu_2} + \frac{N_2 - 1}{2(\sigma_2 + (1 - p)\lambda)}.$$

and the expected transit time for a general customer under probabilistic routing in the stationary regime is

$$W_{\text{ind}}(p) = pW_1(p) + (1 - p)W_2(p), \quad p \in [0, 1].$$

Thus the Wardrop principle [31] here equates to the following definition.

Definition 4.1. *Under probabilistic routing, a user equilibrium $p^* \in [0, 1]$ satisfies one of the following conditions:*

- a) $W_2(0) \leq W_1(0)$ with $p^* = 0$,
- b) $W_1(1) \leq W_2(1)$ with $p^* = 1$,
- c) $W_1(p^*) = W_2(p^*)$ with $p^* \in [0, 1]$.

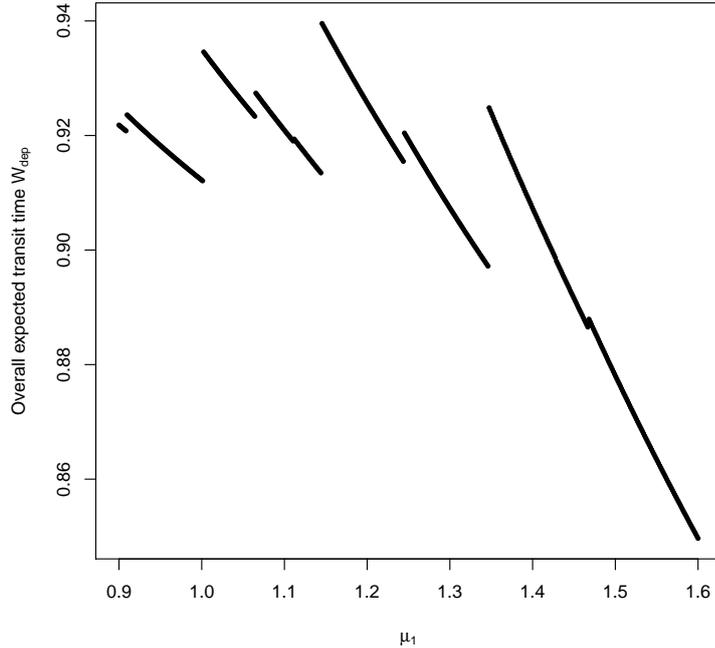


Figure 2: An example showing non-monotonicity of $W_{\text{dep}}(D^*)$ when $\lambda = 4, \sigma_1 = 3, \sigma_2 = 1, \mu_2 = 2, N_1 = N_2 = 5$.

Since $W_1(p)$ is a decreasing function in p whereas $W_2(p)$ is increasing in p , at least one of the conditions in Definition 4.1 holds and hence under probabilistic routing there exists at least one user equilibrium. However, even though the user equilibrium is unique in the state dependent case (Theorem 2.7), this is not necessarily the case under probabilistic routing. Consider an example with $N_1 = N_2 = 5$, $\lambda = 3$, $\sigma_1 = 3$, $\sigma_2 = 1$, $\mu_2 = 2$ and μ_1 equal to 0.2, 1.5 and 3. Figure 3 plots the expected transit times via each queue, $W_1(p)$ and $W_2(p)$, as well as the overall expected transit time through the system under probabilistic routing, $W_{\text{ind}}(p)$. For case A) with $\mu_1 = 0.2$ there is one user equilibrium at $p^* = 0$; in case B) with $\mu_1 = 1.5$ there are three user equilibria, $p^* = 0$, $p^* = 1$ and $p^* = 1/3$; whereas if $\mu_1 = 3$ (case C), there are two user equilibria at $p^* = 0, 1$. The model studied in Afimeimonga et. al. [1] is another example of a system where multiple user equilibria are possible under probabilistic routing.

We are also interested in the stability of user equilibria, that is, whether a user equilibrium p^* is attracting or not. The idea here is as follows. Suppose a proportion $p \in (0, 1)$ of general arrivals choose queue 1. If $W_1(p) < W_2(p)$ then since the expected transit time via queue 1 is lower than via queue 2, over time this proportion will increase until the expected delay via both queues is equal. Similarly, if $W_1(p) > W_2(p)$ then the proportion will decrease until the expected delays are equal. Now suppose $p = p^* + \epsilon$, where p^* is a user equilibrium. Then, assuming p changes continuously, if $W_1(p) > W_2(p)$ for all sufficiently small ϵ then p^* is attracting from above. We need to consider $p < p^*$ as well, to determine whether p^* is attracting. A similar idea has also been discussed in Afimeimonga et al. [1], and

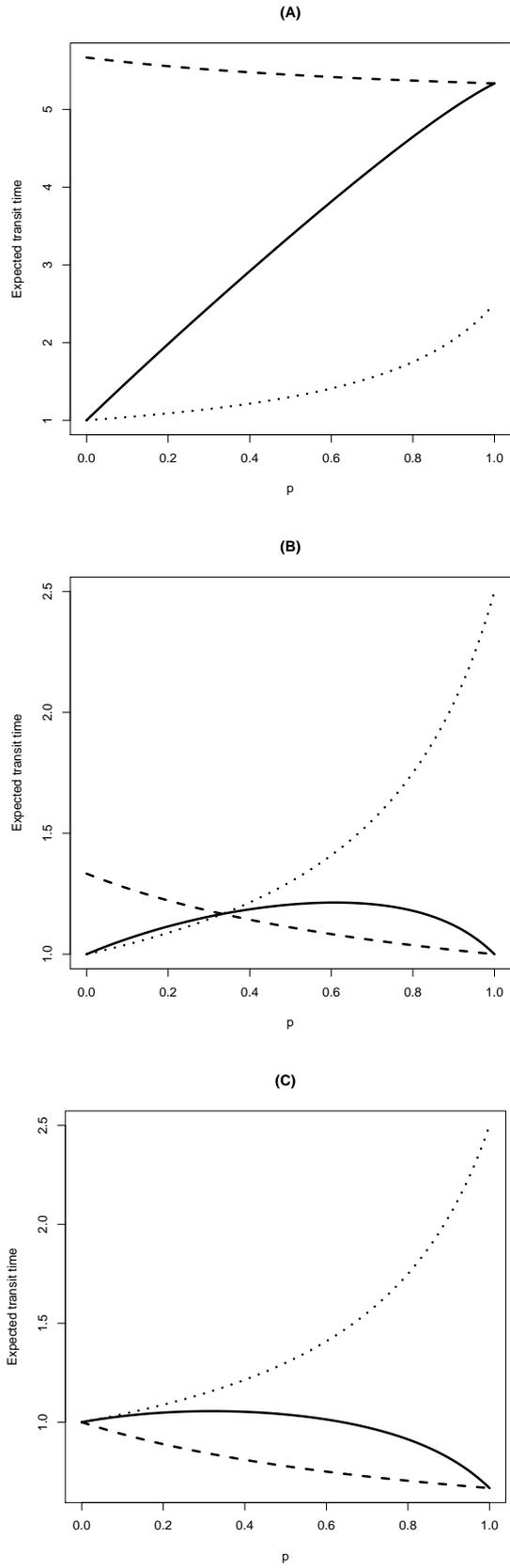


Figure 3: Three cases of user equilibrium p^* ($W_1(p)$ ---, $W_2(p)$ and $W(p)$ —), where $N_1 = N_2 = 5$, $\lambda = 3$, $\sigma_1 = 3$, $\sigma_2 = 1$, $\mu_2 = 2$ and μ_1 is: A) 0.2; B) 1.5; C) 3.

we adopt their definition of stability here.

Definition 4.2. *Given a state independent user equilibrium of the system p^* , we say that it is stable if there exists $\epsilon > 0$ such that*

- 1) $W_1(p) > W_2(p)$ for $p \in (p^*, \min(p^* + \epsilon, 1))$, and
- 2) $W_1(p) < W_2(p)$ for $p \in (\max(p^* - \epsilon, 0), p^*)$.

If the system has a user equilibrium at $p^* = 0$ and $W_2(0) < W_1(0)$ then this is a stable user equilibrium. Similarly, a user equilibrium $p^* = 1$ is stable if $W_1(1) \neq W_2(1)$. However, any user equilibrium $p^* \in (0, 1)$ will be unstable.

One key feature of the batch-service queue is that the expected delay decreases as more customers join the queue. It is not surprising that customers tend to *follow the crowd* (FTC) to obtain a better payoff (in this case, a shorter delay). The model here is an example of the general phenomenon discussed in Hassin and Haviv [19], that under FTC multiple user equilibria may exist, not all of which will be stable.

We give as our final example a comparison of the expected transit time in the stationary regime for the user equilibrium under probabilistic routing and in the state dependent setting. Figure 4 plots $W_{\text{ind}}(p^*)$ and $W_{\text{dep}}(D^*)$ when $N_1 = N_2 = 5$, $\lambda = 4, \mu_2 = 2, \sigma_1 = 3, \sigma_2 = 1$ for μ_1 varying from 0 to 6. For this example, when $\mu_1 \in [0, 0.494)$, there is a single user equilibrium under probabilistic routing at $p^* = 0$. When $\mu_1 \in [4.245, \infty)$ there is also a single user equilibrium under probabilistic routing, but at $p^* = 1$. For $\mu \in (0.494, 4.245)$, three possible user equilibria coexist – one at $p^* = 0$, another at $p^* = 1$, and the third at an intermediate point. Note that although the delay under the state-dependent user equilibrium is, in general, lower than under probabilistic routing, and occasionally substantially lower, there is nevertheless a region where the stable equilibrium at $p^* = 0$ gives lower expected delay than under the state dependent equilibrium. Unlike system optimal policies, state dependent user equilibrium policies do not always have lower expected delays than under probabilistic routing, but in this example we see that under the state dependent user equilibrium policy the worst effects of probabilistic routing are no longer apparent.

There are several possible extensions to the model under consideration. The first natural extension, which we have not considered, is to generalize the times at which service may begin. As mentioned earlier, it may be more natural in many applications, to permit service commencements before a batch is completed. One way of doing this, as discussed in [2], would be to allow service to begin with probability $s(i, j)$ when an arrival finds the system in state $(i, j), 0 \leq i \leq N_1, 0 \leq j \leq N_2$, with the $s(i, j)$ increasing in both i and j , and $s(i, N_2 - 1) = 1, s(N_1 - 1, j) = 1$. This is left for further development.

We have also not considered here a comparison with socially optimal policies – do these possess similar structural properties, and how much does performance of the system improve by using a system rather than user optimal policy.

Similar routing problems can also be considered for queues with other service disciplines where decisions made by later arrivals affect delays experienced by users. A queueing discipline of particular interest here is processor-sharing. In this case, it is not possible to find user equilibrium policies by direct construction, and it is necessary to use an iterative approach to infer properties of user equilibrium policies.

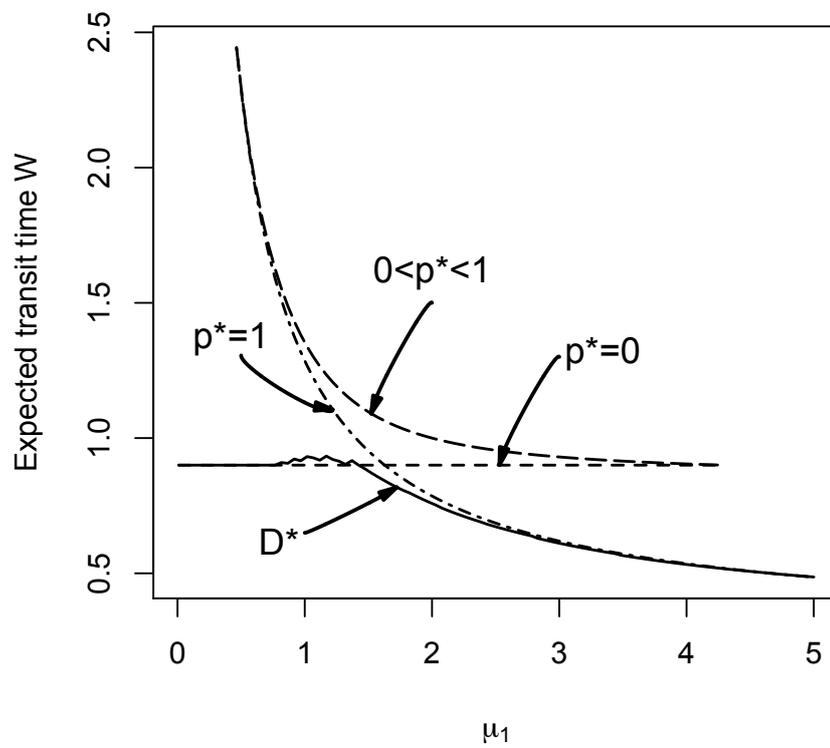


Figure 4: Comparison of the expected transit time for probabilistic routing ($W_{\text{ind}}(p^*)$) vs. state dependent routing ($W_{\text{dep}}(D^*)$). $N_1 = N_2 = 5$, $\sigma_2 = 1$, $\mu_2 = 2$.

We address this in a forthcoming paper [14] where we will also give examples of user optimal policies that are not monotone.

Finally, an even more general question is the extent to which uniqueness, monotonicity and other structural properties hold for broader classes of systems. The simplest extension is to larger systems of parallel queues, but even for these the coupling arguments used here are not straightforward, since general customers now have more alternatives from which to choose.

Acknowledgments

Yizheng Chen would like to thank the TEC for scholarship support of this project, and to EURANDOM and the Isaac Newton Institute, where parts of this research work was carried out. She also thanks Heti Afimeimonga for helpful discussions. Ilze Ziedins is also grateful to the Isaac Newton Institute for hosting a visit there while some of this work was completed.

References

- [1] Afimeimonga, H., Solomon, W. and Ziedins, I. (2004) The Downs-Thomson paradox: Existence, uniqueness and stability of user equilibria. *Queueing Systems* **49**, 321-334.
- [2] Afimeimonga, H., Solomon, W. and Ziedins, I. (2010) User equilibria for a parallel queueing system with state dependent routing. *Queueing Systems* **66**, 169-193.
- [3] Altman, E. (2005) Applications of dynamic games in queues. In: Nowak, A.S. and Szajowski, K. (eds.) *Advances in Dynamic Games: Applications to Economics, Finance, Optimization, and Stochastic Control*. Springer. 309-342.
- [4] Altman, E., Boulogne, T., Ei-Azouzi, R., Jimenez, T. and Wynter, L. (2006) A survey on networking games in telecommunications. *Computers and Operations Research* **33**, 283-311.
- [5] Altman, E. and Shimkin, N. (1998) individual equilibrium and learning in processor sharing systems. *Operations Research* **46**, 776-784.
- [6] Bell, C.E. and Stidham, S. (1983) Individual versus social optimization in the allocation of customers to alternative serves. *Management Science* **29**, 831-839.
- [7] Ben-Shahar, I., Orda, A. and Shimkin, N. (2000) Dynamic service sharing with heterogeneous preferences. *Queueing Systems* **35**, 83-103.
- [8] Boxma, O., Koole, G. and Liu, Z. (1996) Queueing-theoretic solution methods for models of parallel and distributed systems. *CWI/INRIA Tract* **105 & 106**.
- [9] Braess, D. (1969) Über ein Paradoxon aus der Verkehrsplanung. *Unternehmensforschung* **12**, 258-268.

- [10] Braess, D., Nagurney, A. and Wakolbinger, T. (2005) On a paradox of traffic planning. *Transportation Science* **39**, 446-450.
- [11] Brooms, A.C. (2005) On the Nash equilibria for the FCFS queueing system with load-increasing service rate. *Advances in Applied Probability* **37**, 461–481.
- [12] Calvert, B. (1997) The Downs-Thomson effect in a Markov process. *Probability in the Engineering and Information Sciences* **11**, 327-340.
- [13] Calvert, B., Solomon, W. and Ziedins, I. (1997) Braess's paradox in a queueing network with state-dependent routing. *Journal of Applied Probability* **34**, 134-154.
- [14] Chen, Y.Z., Holmes, M. and Ziedins, I. *User equilibria for parallel processor sharing networks*. In preparation (2010).
- [15] Cohen, J. E. and Kelly, F. P. (1990) A paradox of congestion in a queueing network. *Journal of Applied Probability* **27**, 730-734.
- [16] Downs, A. (1962) The law of peak-hour expressway congestion. *Traffic Quarterly* **16**, 393-409.
- [17] El-Taha, M and Stidham, S (1999) *Sample-path Analysis of Queueing Systems*. Kluwer.
- [18] Gelenbe, E and Pujolle, G (1998) *Introduction to Queueing Networks, 2nd Edition*. Wiley.
- [19] Hassin, R. and Haviv, M. (2003) *To Queue or Not to Queue: Equilibrium Behaviour in Queueing Systems*. Kluwer.
- [20] Jia, J. (2005) Optimal routing policies for two batch service queues with partial information. BSc(Hon) Project. Department of Statistics, The University of Auckland.
- [21] Koole, G., Sparaggis, P.D. and Towsley, I. (1999) Minimizing response times and queue lengths in systems of parallel queues. *Journal of Applied Probability* **36**, 1185-1193.
- [22] Lindvall, T. (1992) *Lectures on the Coupling Method*. Wiley.
- [23] Nash, J. (1950) Equilibrium points in n -person games. *Proceedings of the National Academy of Sciences* **36**, 48-49.
- [24] Patriksson, M. (1994) *The Traffic Assignment Problem: Models and Methods*, VSP, The Netherlands.
- [25] Roughgarden, T. and Tardos, E. (2002) How bad is selfish routing? *Journal of the ACM* **49**, 236-259.
- [26] Spicer, S. and Ziedins, I. (2006) User-optimal state-dependent routing in parallel tandem queues with loss. *Journal of Applied Probability* **43**, 274-281.

- [27] Stidham, S. and Weber, R. (1993) A survey of Markov decision models for control of networks of queues. *Queueing Systems*, **13**, 291-314.
- [28] Thomson, J.M. (1977) *Great Cities and Their Traffic*. Gollancz.
- [29] Thorisson, H. (2000) *Coupling, Stationarity and Regeneration*. Springer.
- [30] Walrand, J. (1988) *An Introduction to Queueing Networks*. Prentice-Hall.
- [31] Wardrop, J.G. (1952) Some theoretical aspects of road traffic research, *Proceedings, Institution of Civil Engineers* **1**, 325-378.
- [32] Weber, R.R. (1978) On the optimal assignment of customers to parallel servers. *Journal of Applied Probability* **15**, 406-413.
- [33] Whitt, W. (1986) Deciding which queue to join: Some counterexamples. *Operations Research* **34**, 55-62.
- [34] Winston, W. (1977) Optimality of the shortest line discipline. *Journal of Applied Probability* **14**, 181-189.