

Heavy traffic analysis of roving server networks

M. A. A. Boon, R. D. van der Mei & E. M. M. Winands

To cite this article: M. A. A. Boon, R. D. van der Mei & E. M. M. Winands (2016): Heavy traffic analysis of roving server networks, Stochastic Models, DOI: [10.1080/15326349.2016.1226142](https://doi.org/10.1080/15326349.2016.1226142)

To link to this article: <http://dx.doi.org/10.1080/15326349.2016.1226142>



Published online: 30 Sep 2016.



Submit your article to this journal [↗](#)



Article views: 57



View related articles [↗](#)



View Crossmark data [↗](#)

Heavy traffic analysis of roving server networks

M. A. A. Boon^a, R. D. van der Mei^{b,c}, and E. M. M. Winands^d

^aEurandom and Department of Mathematics and Computer Science, Eindhoven University of Technology, Eindhoven, The Netherlands; ^bDepartment of Mathematics, Section Stochastics, VU University, Amsterdam, The Netherlands; ^cCentre for Mathematics and Computer Science (CWI), Amsterdam, The Netherlands; ^dKorteweg-de Vries Institute for Mathematics, University of Amsterdam, Amsterdam, The Netherlands

ABSTRACT

This article studies the heavy-traffic (HT) behavior of queueing networks with a single roving server. External customers arrive at the queues according to independent renewal processes and after completing service, a customer either leaves the system or is routed to another queue. This type of customer routing in queueing networks arises very naturally in many application areas (in production systems, computer- and communication networks, maintenance, etc.). In these networks, the single most important characteristic of the system performance is oftentimes the path time, i.e., the total time spent in the system by an arbitrary customer traversing a specific path. The current article presents the first HT asymptotic for the path-time distribution in queueing networks with a roving server under general renewal arrivals. In particular, we provide a strong conjecture for the system's behavior under HT extending the conjecture of Coffman et al.^[8,9] to the roving server setting of the current article. By combining this result with novel light-traffic asymptotics, we derive an approximation of the mean path time for arbitrary values of the load and renewal arrivals. This approximation is not only highly accurate for a wide range of parameter settings, but is also exact in various limiting cases.

ARTICLE HISTORY

Received July 2013
Accepted August 2016

KEYWORDS



Approximation; heavy traffic; path times; queueing network; waiting times

MATHEMATICS SUBJECT CLASSIFICATION

60K25; 90B22

1. Introduction

This article considers heavy-traffic (HT) limits for queueing networks with a single roving server that visits the queues in a cyclic order according to the gated and exhaustive service. Customers from the outside arrive at the queues according to general renewal processes, and the service time and switch-over time distributions are general as well. After receiving service at queue i , a customer is either routed to queue j with probability $p_{i,j}$, or leaves the system with probability $p_{i,0}$. This model can be seen as an extension of the classical polling model (in which customers always leave the system upon completion of their service) by customer routing.

CONTACT M. A. A. Boon  m.a.boon@tue.nl  Eurandom and Department of Mathematics and Computer Science, Eindhoven University of Technology, P.O. Box 513, 5600MB Eindhoven, The Netherlands.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/istm.

© 2016 Taylor & Francis

The vast majority of the polling literature assumes that the arrival process at each queue follows a Poisson process (see, for example, Refs. [3,13,17,33] for overviews of polling systems and their applications). In many applications, however, the interarrival times are not exponentially distributed. Therefore, in this article, we study a network in which the arrival process at each of the queues is a general renewal process. For open networks an important characteristic of the system performance is the path time, defined as the total time spent in the system by an arbitrary customer traversing a specific path. That is, oftentimes service level agreements are made on the total time required for all the service requests of a customer in a system to be finished. Moreover, in many computer-communication and production-inventory systems the single most important performance measure is often not an aggregate measure like the mean path time, rather the probability that this path time exceeds a pre-defined threshold. In view of dimensioning such systems, the importance of the path-time distribution as a performance measure of interest is evident. Due to the routing of customers, which leads to non-renewal arrival processes at the queues and to strong interdependence of the waiting times within a path time, simulation appears to be the only practical recourse at the present time. In these circumstances, one naturally resorts to asymptotic estimates. In particular, in this article, we study the path-time distribution in a queueing network with customer routing under HT conditions.

The motivation for studying the HT regime – which is also the most challenging regime from a scheduling point of view – is two-fold. First, an attractive feature of HT asymptotics is that in many cases they lead to strikingly simple expressions for the performance measures of interest. This remarkable simplicity of the HT asymptotics leads to structural insights into the dependence of the performance measures on the system parameters and gives fundamental insights in the behavior of the system in general. A second appealing feature of HT asymptotics is that they form an excellent basis for developing simple accurate approximations for the performance measures for stable systems.

The introduced queueing network is very general, which is illustrated by the fact that many special cases have been studied in the past. Some special case configurations are standard polling systems [33], tandem queues [20,35], multi-stage queueing models with parallel queues [15], feedback vacation queues [7,34], symmetric feedback polling systems [32,34], systems with a waiting room [1,31], and many others. Due to the intrinsic complexity of the model, previous studies on the network in its full generality were restricted to queue lengths and waiting-time distributions for stable systems under the assumption of Poisson arrivals (see Refs. [6,28,29,30]). Although the results in these articles are exact they lack an explicit analysis with simple expressions, leading to the need of using numerical techniques to determine performance measures of interest. Moreover, these results are limited to the waiting-time distribution instead of the practically most interesting and theoretically most challenging path-time distribution.

Besides that we have a theoretical interest in the proposed queueing network, the present work is motivated by the fact that customer routing in polling systems

arises naturally in a host of application areas (in production systems, computer- and communication networks, maintenance, etc.). Some examples are a manufacturing system where products undergo service in a number of stages or in the context of rework^[12], a Ferry-based Wireless Local Area Network (FWLAN) in which nodes can communicate with each other or with the outer world via a message ferry^[16], a dynamic order picking system where the order picker drops off the picked items at the depot where sorting of the items is performed^[11], and an internal mail delivery system where a clerk continuously makes rounds within the offices to pick up, sort and deliver mail^[28].

Motivated by the attractiveness of HT asymptotics, several approaches have been proposed to obtain HT limits for polling systems. HT limits have been rigorously proven for systems with Poisson arrivals (cf. Refs.^[21,37]) and renewal arrivals (cf. Refs.^[8,9,14,38]). A central role in these articles is the Heavy-Traffic Averaging Principle (HTAP) which means that the total scaled workload may be considered as a constant during a cycle, whereas the workload of the individual queues change much faster according to deterministic trajectories, or a fluid model. In this article, we also follow this well-established path of deriving HT limits for systems with general renewal arrivals, see also Refs.^[18,19,24,25,26].

One of the main contributions of the current article is that we present the first HT asymptotic for the path-time distribution in queueing networks with a roving server under general renewal arrivals. The main building blocks of this path-time distribution are inevitably the waiting-time distributions at the individual queues, for which the current article also presents the first exact HT asymptotics¹. In particular, we provide a strong conjecture for the systems behavior under HT extending the polling conjecture of Coffman et al.^[8,9] to the roving server setting of the current article. Our conjecture is validated in three ways. First, as stated before we follow a well established and accepted line of thinking. Secondly, our conjecture corresponds to the aforementioned rigorously proven distributional limits in special cases of our network. Thirdly, we give some numerical examples that illustrate that the correct limiting behavior has been derived. As an important by-product of our analytical framework, we obtain exact asymptotics in the large switch-over time regime as well.

The second main contribution concerns the derivation of a simple approximation of the mean path time for arbitrary values of the load and renewal arrivals by combining the HT results with newly derived light-traffic (LT) asymptotics. This approximation technique is shown to be exact in various limiting cases, and is known to be highly accurate for a wide range of parameter settings (see Refs.^[5,10]). Moreover, it satisfies the Pseudo Conservation Law (PCL), and consequently it leads to exact closed-form results of the mean waiting time for symmetric systems with Poisson arrivals. The resulting expressions are very insightful, simple to implement, and suitable for optimization purposes. In particular, the approximation shows explicitly how the path times depend on the system parameters such as the routing probabilities $p_{i,j}$.

¹ Some preliminary results restricted to mean waiting times were derived in Boon et al.^[4].

Our presentation of the analysis, and the analysis itself, may be considered to be somewhat informal throughout. Providing a rigorous presentation of our results, however, would be an extremely interesting, but notoriously difficult area for further research. Furthermore, it would take us far afield from our main goal, i.e., to obtain fundamental insights into customer routing in systems. Lastly, we would like to stress that the current article concerns a continuous-time cyclic system with gated or exhaustive service in each queue, but that all results can be extended to discrete time, to periodic polling, to batch arrivals, or to systems with different branching-type service disciplines such as globally gated service.

The structure of the present article is as follows. In [Section 2](#), we introduce the model and the required notation. In [Section 3](#), we use the HTAP to derive exact queue-length, waiting-time, and path-time asymptotics for networks with gated and/or exhaustive service. Based on these results, we develop novel approximations in [Section 4](#). In the penultimate section, we present some practical cases that illustrate the versatility of the queueing network and the importance of the path-time distribution in practice. Finally, we present some conclusions and points for discussion in the last section.

2. Model and notation

In this article, we consider a queueing network consisting of $N \geq 2$ infinite buffer queues Q_1, \dots, Q_N . External customers arrive at Q_i according to a general renewal arrival process with rate λ_i , and have a generally distributed service requirement B_i at Q_i , with mean value $b_i := \mathbb{E}[B_i]$, and second moment $b_i^{(2)} := \mathbb{E}[B_i^2]$. Throughout this article, we assume that all random variables have finite second moments. The queues are served by a single server in cyclic order. Whenever the server switches from Q_i to Q_{i+1} , a switch-over time R_i is incurred, with mean r_i . The cycle time C_i is the time between successive moments when the server arrives at Q_i . The total switch-over time in a cycle is denoted by $R = \sum_{i=1}^N R_i$ and its first two moments are $r := \mathbb{E}[R]$ and $r^{(2)} := \mathbb{E}[R^2]$. Indices throughout the article are modulo N , so Q_{N+1} actually refers to Q_1 . All service times and switch-over times are mutually independent. Each queue receives exhaustive or gated service. Exhaustive service means that each queue is served until no customers are present anymore, whereas gated service means that only those customers present at the server's arrival at Q_i will be served before the server switches to the next queue. This queueing network can be modeled as a *polling system* with the specific feature that it allows for routing of the customers: Upon completion of service at Q_i , a customer is either routed to Q_j with probability $p_{i,j}$, or leaves the system with probability $p_{i,0}$. Note that $\sum_{j=0}^N p_{i,j} = 1$ for all i , and that the transition of a customer from Q_i to Q_j takes no time. The model under consideration has a branching structure, which is discussed in more detail by Resing^[27]. The total arrival rate at Q_i is denoted by γ_i , which is the unique solution of the following set of linear equations: $\gamma_i = \lambda_i + \sum_{j=1}^N \gamma_j p_{j,i}$, $i = 1, \dots, N$. The offered load to Q_i is $\rho_i := \gamma_i b_i$ and the total utilization is $\rho := \sum_{i=1}^N \rho_i$. We assume that the system

is stable, which means that ρ should be less than one (see Sidi et al.^[30]). The total service time of a customer is the total amount of service given during the presence of the customer in the network, denoted by \tilde{B}_i , and its first two moments by $b_i := \mathbb{E}[\tilde{B}_i]$ and $\tilde{b}_i^{(2)} := \mathbb{E}[\tilde{B}_i^2]$. The first two moments are uniquely determined by the following set of linear equations: For $i = 1, \dots, N$,

$$\tilde{b}_i = b_i + \sum_{j=1}^N \tilde{b}_j p_{i,j}, \quad (1)$$

$$\tilde{b}_i^{(2)} = b_i^{(2)} + 2b_i \sum_{j=1}^N \tilde{b}_j p_{i,j} + \sum_{j=1}^N \tilde{b}_j^{(2)} p_{i,j}. \quad (2)$$

We study this model under heavy-traffic conditions, i.e., we increase the load of the system until it reaches the point of saturation, $\rho \uparrow 1$. As the total load of the system increases, the visit times, cycle times, and waiting times become larger and will eventually grow to infinity. For this reason, we scale them appropriately and consider the scaled versions. We consider several variables as a function of the load ρ in the system. For each variable x that is a function of ρ , we denote its value *evaluated at* $\rho = 1$ by \hat{x} . Scaling is done by varying the interarrival times of the external customers. To be precise, the limit is taken such that the external arrival rates $\lambda_1, \dots, \lambda_N$ are increased, while keeping the service and switch-over time distributions, the routing probabilities, and the *ratios* between these arrival rates fixed. For $\rho = 1$, the generic interarrival time of the stream in Q_i is denoted by \hat{A}_i . Reducing the load ρ is done by scaling the interarrival times, i.e., taking the random variable $A_i := \hat{A}_i/\rho$ as generic interarrival time at Q_i . Hence, the rate of the arrival stream at Q_i satisfies $\lambda_i = 1/\mathbb{E}[A_i]$. Furthermore, we define arrival rates $\hat{\lambda}_i = 1/\mathbb{E}[\hat{A}_i]$, and proportional load at Q_i , $\hat{\rho}_i = \rho_i/\rho$ (“proportional” because $\sum_{i=1}^N \hat{\rho}_i = 1$).

3. HT asymptotics

To obtain HT-results for the waiting-time distributions, we use HT results for polling systems without customer routing, which are obtained by Refs.^[8,9,14]. The key observation in these articles is the occurrence of a so-called HTAP. When a polling system becomes saturated, two limiting processes take place. Let V denote the total workload of the system, i.e., the total service requirement of all customers present in the system including the possible residual service time of a customer being served. As the load offered to the system, ρ , tends to 1, the scaled total workload $(1 - \rho)V$ tends to a Bessel-type diffusion. However, the work *in each queue* is emptied and refilled at a faster rate than the rate at which the total workload is changing. This implies that during the course of a cycle, the total workload can be considered as constant, while the workloads of the individual queues fluctuate according to a fluid model. The HTAP relates these two limiting processes.

Although rigorous proofs have only been presented for standard polling models (see Refs.^[8,9,14]), results in Refs.^[18,19,24,25,26] support the conjecture that the HTAP

holds for a much wider class of systems. As in these articles, we make the crucial assumption that the HTAP holds without providing a rigorous proof of convergence. That is, the HTAP occurs due to a time-scale decomposition that is inherent in the heavy-traffic scaling. Therefore, in HT the multi-dimensional individual workload processes move along a path in the constant workload hyperplane. For systems without routing and switch-over times, this path is described in detail by Jennings^[14]. Using the results from this section, this can easily be adapted to our model with customer routing. Furthermore, it is known that the scaled total workload tends to a Bessel-type limit as the system becomes saturated. More colloquially, the routing of customers impacts the individual workloads, which impels us to considerably modify and extend the HT analysis of polling models, but does not affect the time scale decomposition which directly implies that the HTAP should also hold in the current setting.

We provide justification for this conjecture in four ways. First, we follow a well-founded line of thinking. Secondly, our conjecture corresponds to the rigorously proven distributional limit in the special case of a two-queue polling model and branching-type polling models with Poisson arrivals. Thirdly, in the next section, we present numerical examples supporting the conjecture that the correct limiting behavior has been derived. Finally, in the appendix, we provide a theorem and proof for the limiting queue-length distributions under the assumption of Poisson arrivals. In the next subsection, we start by discussing the fluid model and subsequently discuss the limiting distribution of the scaled total workload. We use these results to obtain the HT limit of the scaled waiting-time and path-time distributions.

3.1. Fluid model: Gated service

3.1.1. Workload

In this section, we consider the fluid version of the queueing network with a single shared server, where the work travels as fluid from one station to another. This fluid model is carefully constructed in such a way that its behavior corresponds to the multi-dimensional individual workload processes, moving along a path in the constant workload hyperplane as discussed in the introduction of this section. An important aspect of this “corresponding fluid model” is the absence of switch-over times, as they become negligible under HT conditions. We start by studying the fluid limit of the per-queue workload, which is obtained by multiplying by $(1 - \rho)$ and letting $\rho \uparrow 1$. For our model, the fluid limit of the workload at Q_i is a piecewise linear function. During the visit time of Q_k , denoted by V_k , $k = 1, \dots, N$, external fluid particles, corresponding to the customers in the original model, flow into Q_i at rate $\hat{\lambda}_i$. Each of these fluid particles brings along \tilde{b}_i units of work into the system. Simultaneously, work is being processed in Q_k at rate one. Since $\sum_{i=1}^N \hat{\lambda}_i \tilde{b}_i = 1$, the total workload remains constant throughout the course of a cycle. In steady state, the length of one cycle is fixed and denoted by c . In this paragraph, we will show that there is a simple linear relation between c and the total workload in the system,

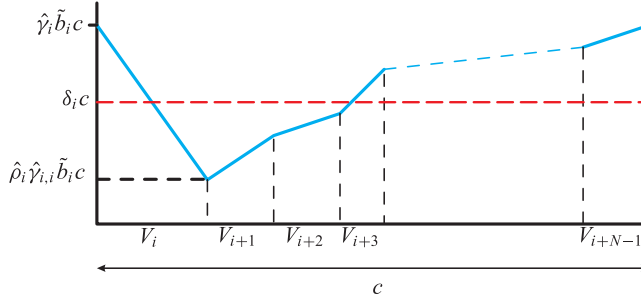


Figure 1. Mean amount of work in Q_i in the fluid limit that arises when the system is in heavy traffic. The length of one cycle is c .

denoted by v . For this reason, we may regard c as a system parameter in the fluid model, instead of v .

Although work is processed at rate one, due to the internal routing work is flowing out of Q_k at rate

$$1 + \frac{1}{b_k} \sum_{i=1}^N p_{k,i} \tilde{b}_i = \frac{\tilde{b}_k}{b_k},$$

which is greater than (or equal to) one. The reason for this anomaly is that work decreases in Q_k either because of the service of fluid particles (customers) in this queue, or because work is shifted due to internal routing of fluid. Work *including* rerouted fluid particles is flowing into Q_i , during V_k , at rate $\hat{\gamma}_{i,k} \tilde{b}_i$, where $\hat{\gamma}_{i,k} := \hat{\lambda}_i + p_{k,i}/b_k$, for $i, k = 1, \dots, N$. It is straightforward to verify that $\tilde{b}_k/b_k = \sum_{i=1}^N \hat{\gamma}_{i,k} \tilde{b}_i$. Figure 1 depicts a graphical representation of the mean amount of work in Q_i in the fluid limit throughout the course of a cycle, the length of which is a constant, denoted by c . Using the fact that fluid particles (external and internal) flow into Q_i during V_k at rate $\hat{\gamma}_{i,k}$ and the fact that the length of V_k in the fluid model is equal to $\hat{\rho}_k c$, one can show that the fluid limit of the mean amount of work in Q_i at the beginning of a visit to Q_j is

$$\sum_{k=i}^{j-1} \hat{\rho}_k \hat{\gamma}_{i,k} \tilde{b}_i c \quad \text{for } j = i+1, \dots, i+N. \quad (3)$$

This reduces to $\hat{\gamma}_i \tilde{b}_i c$ for $j = i+N$. We have used that in the fluid limit the fraction of time that the server is visiting Q_j is $\hat{\rho}_j$ ($j = 1, \dots, N$). Combining these observations, one can obtain the following expression for δ_i , defined as the ratio of the fluid limit of the average amount of work at Q_i and the length of a cycle (see Figure 1).

Definition 3.1.1.1. For $i = 1, \dots, N$,

$$\delta_i = \frac{1}{2} \hat{\rho}_i \tilde{b}_i (\hat{\gamma}_i + \hat{\rho}_i \hat{\gamma}_{i,i}) + \sum_{j=i+1}^{i+N-1} \hat{\rho}_j \left(\frac{1}{2} \hat{\rho}_j \tilde{b}_i \hat{\gamma}_{i,j} + \sum_{k=i}^{j-1} \hat{\rho}_k \tilde{b}_i \hat{\gamma}_{i,k} \right). \quad (4)$$

We introduce the notation v_i for the average amount of work in Q_i ,

$$v_i = \delta_i c.$$

As the *total* inflow in all queues is equal to the total outflow per time unit, the total amount of work during a cycle remains constant at level

$$v = \sum_{i=1}^N v_i = \sum_{i=1}^N \delta_i c = \delta c,$$

where $\delta = \sum_{i=1}^N \delta_i$.

3.1.2. Amount of fluid

We first study the number of fluid particles in each queue at various epochs during the cycle. We denote by $X_{i,k}^{\text{fluid}}$ the number of fluid particles in Q_i at the *beginning* of a visit to Q_k . Using (3) and the fact that the amount of fluid in Q_i is equal to the amount of work in Q_i divided by \tilde{b}_i , we obtain

$$X_{i,k}^{\text{fluid}} = \sum_{j=i}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j} c \quad \text{for } k = i+1, \dots, i+N.$$

Again, note that $X_{i,i+N}^{\text{fluid}}$ can also be written as

$$X_{i,i}^{\text{fluid}} = \hat{\gamma}_i c.$$

As a consequence, the amount of fluid at an arbitrary moment during V_k , denoted by $L_{i,k}^{\text{fluid}}$, is uniformly distributed on the interval

$$\begin{cases} \left[\sum_{j=i}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j} c, \sum_{j=i}^k \hat{\rho}_j \hat{\gamma}_{i,j} c \right] & \text{for } k = i+1, \dots, i+N-1, \\ [\hat{\rho}_i \hat{\gamma}_{i,i} c, \hat{\gamma}_i c] & \text{for } k = i. \end{cases} \quad (5)$$

We can now obtain an expression for the distribution of the amount of fluid in Q_i at an arbitrary epoch, by conditioning on the visit period:

$$L_i^{\text{fluid}} \stackrel{d}{=} L_{i,k}^{\text{fluid}} \quad \text{w.p. } \hat{\rho}_k, \quad (6)$$

where $L_{i,k}^{\text{fluid}}$ is uniformly distributed as in (5). For notational reasons, we introduce the notation $\mathcal{L}_i^{\text{fluid}} := L_i^{\text{fluid}}/c$ and $\mathcal{L}_{i,k}^{\text{fluid}} := L_{i,k}^{\text{fluid}}/c$. One can consider $\mathcal{L}_i^{\text{fluid}}$ as a standardized version of L_i^{fluid} , not depending on the cycle length c .

3.1.3. Waiting times

For the fluid model under consideration, we are interested in the waiting-time distribution of an arbitrary fluid particle, internal or external. We define the waiting time as the time between the arrival in a queue, and the moment of departure from this queue (even if the particle is routed to another, or even the same queue). If we condition on the event that an arbitrary fluid particle arrives in Q_i during V_k ,

its conditional waiting time consists of the residual part of V_k , the visit periods V_{k+1}, \dots, V_{i-1} , and the processing of the amount of fluid that has arrived in Q_i during the elapsed part of the cycle, i.e., V_i, \dots, V_{k-1} plus the elapsed part of V_k . Let u_k be the fraction of V_k that has elapsed at the arrival epoch of a fluid particle in Q_i . We denote by $W_{i,k}^{\text{fluid}}(u_k)$, for $0 \leq u_k \leq 1$, the conditional waiting time of an arbitrary fluid particle arriving in Q_i during V_k , at the moment that a fraction u_k of this visit period has elapsed. We have

$$W_{i,k}^{\text{fluid}}(u_k) = \underbrace{\sum_{j=i-N}^{k-1} \hat{\rho}_j c \hat{\gamma}_{i,j} b_i}_{\mathcal{P}_{i,k}(u_k)} + \underbrace{u_k \hat{\rho}_k c \hat{\gamma}_{i,k} b_i + (1 - u_k) \hat{\rho}_k c + \sum_{j=k+1}^{i-1} \hat{\rho}_j c}_{\mathcal{R}_{i,k}(u_k)}, \quad (7)$$

for $i = 1, \dots, N$, $k = i - N, \dots, i - 1$, and $0 \leq u_k \leq 1$. As (7) indicates, we split $W_{i,k}^{\text{fluid}}(u_k)$ into two parts, namely $\mathcal{P}_{i,k}(u_k)$ representing the waiting time due to the customers that arrived during the elapsed (Past) part of the cycle c , and $\mathcal{R}_{i,k}(u_k)$, which is the time until the next visit of the server to Q_i (Residual cycle). This division into two parts turns out to be useful in the next paragraph, for computing the path-time distributions.

The waiting-time distribution follows after unconditioning. During V_k fluid flows into Q_i at rate $\hat{\gamma}_{i,k}$. Hence, the probability that an arbitrary fluid particle arrives during V_k , given that it arrives in Q_i , is $\pi_{i,k} := \hat{\gamma}_{i,k} \hat{\rho}_k / \hat{\gamma}_i$. The elapsed fraction of the visit period V_k is uniformly distributed on $[0, 1]$. These two results yield the following expression for the waiting-time distribution of an arbitrary fluid particle in Q_i :

$$\begin{aligned} W_i^{\text{fluid}} &\stackrel{d}{=} \mathcal{P}_{i,k}(U_k) + \mathcal{R}_{i,k}(U_k) \quad \text{w.p. } \pi_{i,k} \\ &= c \left(1 + \sum_{j=i-N}^{k-1} \hat{\rho}_j (\hat{\gamma}_{i,j} b_i - 1) + U_k \hat{\rho}_k (\hat{\gamma}_{i,k} b_i - 1) \right) \quad \text{w.p. } \pi_{i,k}, \end{aligned} \quad (8)$$

for $i = 1, \dots, N$ and $k = i - N, \dots, i - 1$. The random variables U_1, \dots, U_N are independent and Uniform $[0, 1]$ distributed. As before, we introduce the notation $\mathcal{W}_i^{\text{fluid}} := W_i^{\text{fluid}} / c$ to represent a standardized version of W_i^{fluid} , not depending on the cycle length c .

3.1.4. Path times

In this paragraph, we derive the path-time (the total time spent in the system) distribution of customers – or, in this case, fluid particles – traversing a specific path through the network. We denote the time spent in the system by a fluid particle traversing the path $Q_{i_1}, Q_{i_2}, \dots, Q_{i_M}$ by $W_{i_1, i_2, \dots, i_M}^{\text{fluid}}$, where $i_k \in \{1, 2, \dots, N\}$ for $k = 1, 2, \dots, M$ and $M \geq 1$. When considering path times, each fluid particle enters the system as an *external* fluid particle in Q_{i_1} . Just like in the previous paragraph, we condition on the visit period and the length of the elapsed part of this visit period at its arrival epoch. Assume that a tagged fluid particle arrives in Q_{i_1} at the moment that the server is visiting Q_k and a fraction u_k of V_k has elapsed ($0 \leq u_k \leq 1$). We denote its conditional path time as $W_{i_1, i_2, \dots, i_M; k}^{\text{fluid}}(u_k)$. Its time spent in Q_{i_1} is exactly

$W_{i_1,k}^{\text{fluid}}(u_k)$. At the moment that the tagged fluid particle leaves Q_{i_1} and is routed to Q_{i_2} , the elapsed time of V_{i_1} is $\mathcal{P}_{i_1,k}(u_k)$. As a consequence, the elapsed *fraction* of V_{i_1} is $\mathcal{P}_{i_1,k}(u_k)/(\hat{\rho}_{i_1}c)$. This result leads to the following recursive equation for the conditional path time,

$$W_{i_1,i_2,\dots,i_M;k}^{\text{fluid}}(u_k) = W_{i_1,k}^{\text{fluid}}(u_k) + W_{i_2,\dots,i_M;i_1}^{\text{fluid}}\left(\frac{\mathcal{P}_{i_1,k}(u_k)}{\hat{\rho}_{i_1}c}\right).$$

We define $W_{i_M;i_{M-1}}^{\text{fluid}}(\cdot) := W_{i_M,i_{M-1}}^{\text{fluid}}(\cdot)$ to ensure that the recursion always ends.

The path-time distribution follows after unconditioning, noting that the probability that an *external* fluid particle enters the system during V_k is $\hat{\rho}_k$:

$$W_{i_1,i_2,\dots,i_M}^{\text{fluid}} \stackrel{\text{d}}{=} W_{i_1,i_2,\dots,i_M;k}^{\text{fluid}}(U_k) \quad \text{w.p. } \hat{\rho}_k, \quad (9)$$

for $i_1, i_2, \dots, i_M \in \{1, 2, \dots, N\}$, and $k = 1, 2, \dots, N$. The random variables U_1, \dots, U_N are independent and Uniform $[0, 1]$ distributed. Note that $W_{i_1,i_2,\dots,i_M}^{\text{fluid}}/c$ does not depend on the cycle time c anymore. Similar to the standardized waiting time, we now introduce the notation $\mathcal{W}_{i_1,i_2,\dots,i_M}^{\text{fluid}} := W_{i_1,i_2,\dots,i_M}^{\text{fluid}}/c$, which turns out to be useful later.

3.2. Fluid model: Exhaustive service

In this subsection, we briefly discuss the fluid model when some of the queues are served exhaustively. Rather than repeating all of the analysis of the previous subsection, we mainly focus on the differences compared to a model with gated service. Note that changing the service discipline of a particular queue has no effects on any of the other queues.

3.2.1. Workload

Assume that the service discipline at a certain queue, say Q_e with $e = 1, \dots, N$, is exhaustive service. When comparing the fluid trajectory of the workload of Q_e during the course of a cycle to the case with gated service (as depicted in Figure 1), the only change is that the trajectory needs to be moved downward such that the queue is empty at the end of the visit period. More precisely, if Q_e receives exhaustive service, then we define

$$\delta_e = \delta_e^{\text{gated}} - \hat{\rho}_e \hat{\gamma}_{e,e} \tilde{b}_e, \quad (10)$$

where δ_e^{gated} is the value of δ_e given by (4) for the case that Q_e would have received gated service.

3.2.2. Amount of fluid

Since the amount of fluid in Q_i is equal to the amount of work in Q_i divided by \tilde{b}_i , and since we have just established that the amount of work for an exhaustively served queue Q_e is equal to the amount of work this queue would have under gated service, minus $\hat{\rho}_e \hat{\gamma}_{e,e} \tilde{b}_e$, we directly obtain that the amount of fluid at an arbitrary

moment during V_k , denoted by $L_{e,k}^{\text{fluid}}$, is uniformly distributed on the interval

$$\begin{aligned} & \left[\sum_{j=e+1}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j} c, \sum_{j=e+1}^k \hat{\rho}_j \hat{\gamma}_{i,j} c \right] && \text{for } k = e+1, \dots, e+N-1, \\ & [0, (\hat{\gamma}_e - \hat{\rho}_e \hat{\gamma}_{e,e}) c] && \text{for } k = e. \end{aligned} \quad (11)$$

3.2.3. Waiting times

Determining the waiting-time distribution of an arbitrary fluid particle also involves the same steps as in the gated case, except for fluid particles arriving in Q_e during V_e , obviously.

$$\begin{aligned} W_{e,k}^{\text{fluid}}(u_k) &= \begin{cases} (1 - u_k) \hat{\rho}_k c + \underbrace{\sum_{j=k+1}^{e-1} \hat{\rho}_j c + \sum_{j=e-N+1}^{k-1} \hat{\rho}_j c \hat{\gamma}_{e,j} b_e + u_k \hat{\rho}_k c \hat{\gamma}_{e,k} b_e}_{\mathcal{P}_{e,k}(u_k)} & (k \neq e), \\ \underbrace{(1 - u_e)(\hat{\gamma}_e c - \hat{\gamma}_{e,e} \hat{\rho}_e c) b_e}_{\mathcal{P}_{e,e}(u_e)} & (k = e). \end{cases} \end{aligned}$$

Note that $\mathcal{R}_{e,e}(u_e) = 0$ when the service in Q_e is exhaustive.

3.2.4. Path times

We consider the conditional path time $W_{e,i_2,\dots,i_M;k}^{\text{fluid}}(u_k)$, where Q_e receives exhaustive service. As in the previous paragraph, we need to treat the case $e = k$ separately. A fluid particle arriving in Q_e during V_e , given that a fraction u_e of this visit period has elapsed, will wait for $W_{e,e}^{\text{fluid}}(u_e)$ before leaving this queue. This waiting time, which can also be denoted as $\mathcal{P}_{e,e}(u_e)$, is a fraction $\frac{\mathcal{P}_{e,e}(u_e)}{\hat{\rho}_e c}$ of the visit period V_e . This gives the following expression for the conditional path time:

$$W_{e,i_2,\dots,i_M;k}^{\text{fluid}}(u_k) = \begin{cases} W_{e,k}^{\text{fluid}}(u_k) + W_{i_2,\dots,i_M;e}^{\text{fluid}}\left(\frac{\mathcal{P}_{e,k}(u_k)}{\hat{\rho}_e c}\right), & (e \neq k), \\ W_{e,e}^{\text{fluid}}(u_e) + W_{i_2,\dots,i_M;e}^{\text{fluid}}\left(u_e + \frac{\mathcal{P}_{e,e}(u_e)}{\hat{\rho}_e c}\right), & (e = k). \end{cases}$$

3.3. Original model: Gated service

In this subsection, we expand upon the HTAP for polling models by relating the limit processes of the total workload process and the workload of the individual queues. The main difference with polling models is that in the roving server network, (the direction of) the shifting of individual workloads is not only determined by which queue is being served but also by the internal routing of the customers. To this end, we return to the original model under HT conditions.

3.3.1. Workload

We denote by V the total amount of work in the system at the start of a cycle starting, without loss of generality, at the beginning of visit period V_1 . As far as the total amount of work is concerned, the system behaves like a polling system in heavy traffic with external customers bringing in an amount of work \tilde{B}_i in Q_i , but with work shifting from one queue to another upon the service completion of a customer. For polling systems with general renewal arrivals the HT limit of the scaled total amount of work at the beginning of a cycle is conjectured by Olsen and Van der Mei^[22]. An adaptation of the conjecture in Olsen and van der Mei^[22], in accordance with the proof for the case of Poisson arrivals in Olsen and van der Mei^[21], to our model leads to the following result.

Conjecture 3.3.1.1. *Define*

$$\sigma^2 = \sum_{i=1}^N \hat{\lambda}_i \left(\text{Var}[\tilde{B}_i] + (\hat{\lambda}_i \tilde{b}_i)^2 \text{Var}[\hat{A}_i] \right), \quad \alpha = 2r\delta/\sigma^2, \quad \mu = 2/\sigma^2,$$

with δ as defined in [Definition 3.1.1.1](#). Then, for $\rho \uparrow 1$, $(1 - \rho)V$ has a Gamma distribution with shape parameter α and rate parameter μ .

For more details, we refer to Olsen and van der Mei^[22] (who, in turn, refer to a result from Coffman et al.^[9]).

3.3.2. Cycle time

Subsequently, the diffusion limit of the *total* workload process and the workload in the individual queues can be related using the HTAP. To this end, we start with the cycle-time distribution under HT scalings, which follows from Conjecture 3.3.1.1 and the fluid analysis carried out in the first part of this section. The length of a cycle depends on the amount of work at the beginning of that cycle (which may be any arbitrarily chosen moment). Denote by $C(x)$ the length of a cycle, given that a total amount of x work is present at its beginning. In steady state, we have the following relation:

$$\delta C(x) = x, \tag{12}$$

where δ is defined as in [Definition 3.1.1.1](#). Hence, given an amount of work x , the cycle time is $C(x) = x/\delta$. In the fluid model, all cycles have the same length. However, in the stochastic model, the cycle lengths are random. Note that, although the cycle-time distribution depends on the service disciplines and the chosen starting point of the cycle, the *mean* cycle time is always $\mathbb{E}[C_i] = \mathbb{E}[C] = r/(1 - \rho)$ (cf. Sidi et al.^[30]). We are now ready to formulate the second conjecture, concerning the limiting distribution of the scaled length-biased cycle time.

Conjecture 3.3.2.1. *For $\rho \uparrow 1$, we find that $(1 - \rho)C_i$ converges in distribution to a random variable having a Gamma distribution with shape parameter α and rate parameter $\delta\mu$.*

3.3.3. Queue lengths

Given the cycle-time distribution, we can finally find the scaled queue-length distributions under HT conditions. We use the fluid analysis, in combination with the conjectured cycle-time distribution, to find the limiting distribution of the scaled queue lengths. In the fluid analysis, the cycle time had a fixed length c . Due to the HTAP, we can replace the constant cycle time from the fluid analysis by the random variable C_i , the scaled *length-biased* cycle time. If a random variable X has a Gamma distribution with parameters a and b , its length-biased equivalent \bar{X} has a Gamma distribution with parameters $a + 1$ and b . Obviously, the replacement of c by C_i can only be carried out because of the independence between the length of the cycle time and the uniformly distributed random variables appearing in (6). The following conjecture summarizes this result. A theorem and proof for the scaled queue length distribution under the assumption of *Poisson arrivals* can be found in the appendix.

Conjecture 3.3.3.1. *As $\rho \uparrow 1$, the scaled queue length $(1 - \rho)L_i$ converges in distribution to the product of two independent random variables. The first has the same distribution as $\mathcal{L}_i^{\text{fluid}}$, and the second random variable Γ has the same distribution as the limiting distribution of the scaled length-biased cycle time, $(1 - \rho)C_i$. For $i = 1, \dots, N$; $k = i, \dots, i + N - 1$, and $\rho \uparrow 1$,*

$$(1 - \rho)L_i \xrightarrow{d} \Gamma \times \mathcal{L}_{i,k}^{\text{fluid}} \quad \text{w.p. } \hat{\rho}_k, \quad (13)$$

where Γ is a random variable having a Gamma distribution with parameters $\alpha + 1$ and $\delta\mu$, and $\mathcal{L}_{i,k}^{\text{fluid}}$ is the “standardized” number of a type- i particles in the fluid model during V_k , introduced below (6). The random variables Γ and $\mathcal{L}_{i,k}^{\text{fluid}}$ are independent.

3.3.4. Waiting times

The distributions of the scaled waiting times are obtained in a similar manner, exploiting the HTAP to combine the results from the fluid analysis and the scaled, length-biased cycle-time distribution, which leads to the following conjecture.

Conjecture 3.3.4.1. *As $\rho \uparrow 1$, the scaled waiting time $(1 - \rho)W_i$ converges in distribution to the product of two independent random variables. The first has the same distribution as $\mathcal{W}_i^{\text{fluid}}$, and the second random variable Γ has the same distribution as the limiting distribution of the scaled length-biased cycle time, $(1 - \rho)C_i$. For $i = 1, \dots, N$; $k = i - N, \dots, i - 1$, and $\rho \uparrow 1$,*

$$(1 - \rho)W_i \xrightarrow{d} \Gamma \times \mathcal{W}_i^{\text{fluid}}, \quad (14)$$

where Γ is a random variable having a Gamma distribution with parameters $\alpha + 1$ and $\delta\mu$, and $\mathcal{W}_i^{\text{fluid}}$ is the “standardized” waiting time of a type- i particle in the fluid model, introduced below (8). The two random variables Γ and $\mathcal{W}_i^{\text{fluid}}$ are independent.

The (HT limit of the) *mean* waiting time of an arbitrary customer in Q_i obviously follows from (14), but an easier way to find it, is by application of Little's Law to the mean queue length at Q_i , which is simply the mean amount of work in Q_i divided by the mean total service time.

Corollary 3.3.4.1. *For $i = 1, \dots, N$,*

$$(1 - \rho)\mathbb{E}[W_i] \rightarrow \left(r + \frac{\sigma^2}{2\delta}\right) \frac{\delta_i}{\hat{\gamma}_i \tilde{b}_i}, \quad (\rho \uparrow 1). \quad (15)$$

3.3.5. Path times

The derivation of the limiting distribution of the scaled path times proceeds along the exact same lines, starting with the fluid model. Due to the HTAP, we may again replace the constant cycle time c from the fluid analysis by the random variable C_i to obtain the limiting distribution of the scaled path times. The conjecture below summarizes this result, which is consistent with the heavy-traffic snapshot principle^[39].

Conjecture 3.3.5.1. *As $\rho \uparrow 1$, the scaled path time $(1 - \rho)W_{i_1, i_2, \dots, i_M}$ converges in distribution to the product of a random variable having the same distribution as $\mathcal{W}_{i_1, i_2, \dots, i_M}^{\text{fluid}}$ and a random variable Γ having the same distribution as the limiting distribution of the scaled length-biased cycle time, $(1 - \rho)C_i$. For $i, k = 1, \dots, N$, and $\rho \uparrow 1$,*

$$(1 - \rho)W_{i_1, i_2, \dots, i_M} \xrightarrow{d} \Gamma \times \mathcal{W}_{i_1, i_2, \dots, i_M}^{\text{fluid}}, \quad (16)$$

where Γ is a random variable having a Gamma distribution with parameters $\alpha + 1$ and $\delta\mu$, and the distribution of $\mathcal{W}_{i_1, i_2, \dots, i_M}^{\text{fluid}}$ is given by Eq. (9).

3.4. Original model: Exhaustive service

In Section 3.3, we have used the HTAP to easily derive the limiting distributions of the scaled workload at cycle beginnings, the cycle times, the waiting times, and path times. This principle can be used in the exact same manner when some of the queues receive exhaustive service. Therefore, we will not repeat all these steps to summarize the results for systems with exhaustive service. Basically, one only has to take the expressions from the fluid model, and replace the constant cycle time c by a random variable with the same limiting distribution as $(1 - \rho)C_i$, i.e., a Gamma distribution with parameters $\alpha + 1$ and $\delta\mu$ (as defined earlier in this section, taking $\delta := \delta_e$ as defined in (10) for exhaustive service).

3.5. Poisson arrivals

In the current article, we have derived the system behavior under heavy traffic for systems with general renewal arrival processes based on the partially conjectured HTAP. Van der Mei^[37] has developed a unifying framework to derive rigorous

proofs of the heavy-traffic behavior of branching-type polling models with (compound) *Poisson* arrivals. By applying this stepwise approach in conjunction with the results of the previous section to the model under consideration, one can rigorously prove the HT asymptotics in queueing networks served by a single shared server under the assumption of *Poisson* arrivals. In the appendix, we provide such a proof for the queue-length distributions.

3.6. Increasing setup times

In HT, the system reaches saturation due to an increase in the total utilization ρ . However, the system might also get saturated due to an increase of the total switch-over time r . These two asymptotic regimes show, however, significantly different behavior. In Refs.^[40,41], it was shown for polling systems that the scaled cycle and intervisit times converge in probability to deterministic quantities in the case that the (deterministic) switch-over times tend to infinity. One has to compare this with the Gamma distribution which is prevalent in the scaled cycle time in the diffusion limit of the present section. The results for polling systems with increasing switch-over times of Refs.^[40,41] can be extended to the setting of the current article. That is, as a consequence of the scaled cycle time converging to a constant, a fluid limit is obtained implying that the scaled delay converges in distribution to a mixture of uniform distributions (cf. Formula (8)).

4. Approximations for general traffic conditions

The HT distributions derived in the preceding section may be used directly as an approximation for the waiting-time and path-time distributions in non-heavy-traffic systems. However, they tend to perform poorly under low or moderate traffic. Therefore, in this section, we combine these HT asymptotics with newly developed LT results leading to highly accurate approximations for the mean performance measures for the whole range of load values. We will assess the accuracy of the resulting approximations in Section 5.

4.1. Waiting-time approximations

In order to derive an approximation for the mean waiting time, we study the LT limit of W_i which can be found by conditioning on the customer type (external or internally routed).

Theorem 4.1.1. *For $i = 1, \dots, N$,*

$$W_i \xrightarrow{d} \begin{cases} R^{\text{res}} & w.p. \lambda_i/\gamma_i, \\ R_{j,i} & w.p. \gamma_j p_{j,i}/\gamma_i, \end{cases} \quad (\rho \downarrow 0), \quad (17)$$

where R^{res} is a residual total switch-over time, with probability density function

$$f_{R^{\text{res}}}(t) = \frac{1 - \mathbb{P}(R \leq t)}{\mathbb{E}[R]},$$

and $R_{j,i} = R_j + R_{j+1} + \dots + R_{i-1}$ is the sum of the switch-over times between Q_j and Q_i , which is a cyclic sum if $i \leq j$. Only if $i = j$ and service in Q_i is exhaustive, we have $R_{j,i} = 0$.

Proof. In LT, we ignore all $O(\rho)$ terms, which implies that we can consider a customer as being alone in the system. Equation (17) can be interpreted as follows. An arbitrary customer in Q_i has arrived from outside the network with probability λ_i/γ_i . In this case, he has to wait for a residual total switch-over time R^{res} . If a customer in Q_i arrives after being served in another queue, say Q_j (with probability $\gamma_j p_{j,i}/\gamma_i$), he has to wait for the mean switch-over times R_j, \dots, R_{i-1} . \square

The LT limit of the *mean* waiting time directly follows from (17):

$$\mathbb{E}[W_i] \rightarrow \frac{\lambda_i}{\gamma_i} \frac{r^{(2)}}{2r} + \sum_{j=i^*}^{i-1} \frac{\gamma_j p_{j,i}}{\gamma_i} \sum_{k=j}^{i-1} r_k, \quad (\rho \downarrow 0), \quad (18)$$

where $i^* = i - N$ if Q_i has gated service, and $i^* = i - N + 1$ if Q_i has exhaustive service. Subsequently, we construct an interpolation between the LT and HT limits that can be used as an approximation for the mean waiting time for arbitrary ρ . For $i = 1, \dots, N$,

$$\mathbb{E}[W_i^{\text{approx}}] = \frac{w_i^{LT} + (w_i^{HT} - w_i^{LT})\rho}{1 - \rho}, \quad (19)$$

where w_i^{LT} and w_i^{HT} are the LT and HT limits of the mean waiting time respectively, as given in (18) and (15). Because of the way $\mathbb{E}[W_i^{\text{approx}}]$ is constructed, it has the nice properties that it is exact as $\rho \downarrow 0$ and $\rho \uparrow 1$. Furthermore, if we have Poisson arrivals, it satisfies a so-called pseudo-conservation law for the mean waiting times, which is derived in Ref.^[30]. This implies that the $\mathbb{E}[W_i^{\text{approx}}]$ yields exact results for symmetric (and, hence, single-queue) systems. Finally, it can be shown that this approximation is exact in the limiting case of deterministic set up times that tend to infinity (see Refs.^[40,41]).

The astute reader has already noticed that the LT result (18) is a first-order Taylor expansion of the mean waiting time at $\rho = 0$, which can be naturally extended with the m th derivatives of the mean waiting time with respect to ρ at $\rho = 0$. Together with the HT limit one has $m + 1$ pieces of information, which can be used to construct an $(m + 1)$ th degree polynomial interpolation (cf. Boon et al.^[5]). Therefore, it is not inconceivable that the approximation can be refined further, but since the primary goal of this article has been the derivation of the path-time distributions under HT conditions such refinements are beyond the scope of the article. Moreover, the presented first-order polynomial interpolation is already quite accurate as can be seen in the numerical evaluation.

4.2. Path-time approximations

The principle used to determine waiting-time approximations, can be applied to path times in the exact same manner. For reasons of compactness, we will not present the details in this article. Almost all of the required ingredients to develop a path-time approximation have been discussed. The only missing piece of information is the LT limit of the mean path time. However, it can easily be found given the LT limit of the mean waiting time (18):

$$\mathbb{E}[W_{i_1, i_2, \dots, i_M}] \rightarrow \frac{r^{(2)}}{2r} + b_{i_1} + \sum_{j=2}^M (r_{i_{j-1}, i_j} + b_{i_j}), \quad (\rho \downarrow 0), \quad (20)$$

where $r_{j,i} = \mathbb{E}[R_{j,i}]$, with $R_{j,i}$ as defined in (17).

In the current section, we have focused on an interpolation of the *mean* waiting time and path times for the ease of presentation. It goes without saying that following the same recipe one could derive an identical interpolation approximation for higher moments as well and, subsequently, fit a phase-type distribution for the complete distribution. Alternatively, one could follow the idea of Dorsman et al.^[10], in which a refined HT distribution is derived as an approximation of the complete distribution for arbitrary loads. However, due to the strong correlation between the waiting times of a customer within a specific path, we have observed in numerical tests that the accuracy of this approximation technique is not as high as for standard polling models. We would like to end with noting that for the mean path times these correlations obviously do not play a role, whence the developed approximation for the mean figures gives excellent results as illustrated in the next section.

5. Numerical evaluation

In the current section, we present four practical cases from completely different application areas, indicating the versatility of the studied roving server network and showing the practical usage of the developed asymptotics and approximations. Moreover, these cases clearly illustrate the fact that the path time is the most important performance measures in most practical applications. Special cases of the network are, for example, standard polling systems^[33], tandem queues^[20,35], multi-stage queueing models with parallel queues^[15], feedback vacation queues^[7,34], symmetric feedback polling systems^[32,34], and systems with a waiting room^[1,31]. We would like to remind the reader that the only practical alternative to our expressions for a complete performance analysis of the studied network is simulation.

5.1. Example 1: A production system with rework

The first application is a multiproduct system with random yields introduced by Grasman et al.^[12], which is a stylized practical case of a producer of plastic bumpers

for automobiles. Every item is produced in a production queue and is defect with probability p ; a defect item has to be reproduced in the same production queue. However, defect items are first routed to a temporary storage queue, which is served immediately after the current queue, before they are routed back to the original queue. The combination of exhaustively served production and storage queues implies that newly arriving items are served during the current cycle, whereas defect items will be reproduced in the next cycle. There are no switch-over time and service time required for these storage queues, implying that the storage queues do not contribute to the utilization of the system.

We consider a system consisting of five production queues with relative loads $(0.1, 0.2, 0.2, 0.2, 0.3)$, Poisson arrivals, exponentially distributed service times with mean 1, and switch-over times with constant value 5. A typical feature of these type of production systems is that the switch-over times are relatively large compared to the processing times, in contrast to the next numerical example that we discuss. We have additional five storage queues to which the defect items are routed. Each item is defect with probability 0.25. We run simulations for various values of ρ , and we compare the results with the limiting HT distributions. For each value of ρ we run 100 simulations, each of length 10^8 . These settings yield extremely accurate estimates for the probability distribution functions.

5.1.1. Convergence to the HT limit

In this example, mostly due to the relatively large switch-over times, the (scaled) path-time distributions converge relatively fast to their HT limits. This is illustrated in Table 1 and in Figure 2. Table 1 depicts the means, standard deviations and tail probabilities of the (scaled) path-time distributions for nine paths. Only production queues are included in the path names. For example, path $2 \rightarrow 2 \rightarrow 2$ actually includes two visits to the storage queue of production queue 2. The tail probabilities are given for path lengths of 20, 50, 80 if the corresponding path consists of respectively 1, 2, and 3 visits to a queue. The results for $\rho < 1$ are obtained using simulation, and the results for $\rho = 1$ are obtained using the theory presented in Section 3. In Table 1 and Figure 2, it can be seen that already for $\rho = 0.8$ the scaled path-time distributions are close to the limiting distributions.

5.1.2. Accuracy of the approximation

Table 2 shows the approximated means using the path-time approximation discussed in Section 4. For nine paths, the means of the (scaled) path-time approximations are given. Only production queues are included in the path names. For example, path $2 \rightarrow 2 \rightarrow 2$ actually includes two visits to the storage queue of production queue 2. Note that we have chosen to display the values for the *scaled* path times $(1 - \rho)W_{i_1, i_2, \dots, i_M}$, because now they can be compared directly to the simulated values in Table 1. The results show that the accuracy is very high for all values of the load and for all paths.

5.2. Example 2: Tandem queues with parallel queues in the first stage

Secondly, we use an example that was introduced by Katayama^[15], who studies call processing in packet switching systems composed of multi-processors. In this application, each processor has two kinds of buffers in parallel for receiving data packets from switching links and from lines and/or subscribers. During the first stage, called input processing, the processor polls both buffers for data packets and moves these packets to a queue in the second stage. During this second stage, called packet processing, the jobs are actually executed. This system can be modeled as a network consisting of three queues. Customers arrive at Q_1 and Q_2 , and are routed

Table 1. Numerical results for Example 1.

Scaled path times: Means									
ρ	1	1→1	1→ 1 →1	2	2→2	2→ 2 →2	5	5→5	5→ 5 →5
0.1	13.40	39.29	65.27	13.32	39.08	65.07	13.21	38.88	64.76
0.3	13.22	38.89	64.84	12.92	38.26	64.14	12.64	37.62	63.37
0.5	13.00	38.58	64.47	12.53	37.50	63.22	12.03	36.41	62.04
0.7	12.79	38.30	64.19	12.10	36.77	62.51	11.40	35.24	60.82
0.8	12.67	38.18	64.10	11.88	36.42	62.18	11.08	34.65	60.23
0.9	12.55	38.08	64.03	11.65	36.07	61.86	10.74	34.06	59.67
0.95	12.48	38.03	64.02	11.53	35.90	61.73	10.57	33.76	59.40
1.00	12.42	38.00	64.53	11.41	35.74	61.95	10.40	33.47	59.37
Scaled path times: Standard deviations									
ρ	1	1→1	1→ 1 →1	2	2→2	2→ 2 →2	5	5→5	5→ 5 →5
0.1	7.39	7.86	8.40	7.33	7.83	8.34	7.25	7.78	8.31
0.3	7.59	8.91	10.29	7.39	8.76	10.16	7.19	8.62	10.05
0.5	7.77	9.94	12.24	7.44	9.70	12.02	7.09	9.45	11.80
0.7	7.93	11.02	14.22	7.46	10.66	13.98	6.98	10.29	13.71
0.8	7.99	11.57	15.27	7.46	11.15	15.01	6.92	10.71	14.70
0.9	8.05	12.12	16.33	7.45	11.63	16.06	6.85	11.14	15.72
0.95	8.07	12.41	16.89	7.45	11.89	16.61	6.82	11.37	16.27
1.00	8.09	12.69	18.70	7.44	12.13	18.07	6.78	11.58	17.45
Scaled path times: Tail probabilities									
ρ	1	1→1	1→ 1 →1	2	2→2	2→ 2 →2	5	5→5	5→ 5 →5
ρ	20	50	80	20	50	80	20	50	80
0.8	0.191	0.154	0.159	0.152	0.120	0.130	0.113	0.089	0.103
0.9	0.187	0.161	0.175	0.144	0.123	0.142	0.103	0.090	0.113
0.95	0.186	0.165	0.184	0.141	0.125	0.149	0.098	0.090	0.118
1.00	0.184	0.168	0.192	0.138	0.126	0.156	0.093	0.091	0.123

Table 2. Accuracy of the approximation in Example 1.

Approximated scaled path times: Means									
ρ	1	1→1	1→ 1 →1	2	2→2	2→ 2 →2	5	5→5	5→ 5 →5
0.1	13.39	39.35	65.40	13.29	39.12	65.14	13.19	38.90	64.89
0.3	13.17	39.05	65.21	12.87	38.37	64.43	12.57	37.69	63.66
0.5	12.96	38.75	65.01	12.45	37.62	63.72	11.95	36.49	62.43
0.7	12.74	38.45	64.82	12.04	36.86	63.01	11.33	35.28	61.21
0.8	12.63	38.30	64.72	11.83	36.49	62.66	11.02	34.68	60.60
0.9	12.52	38.15	64.63	11.62	36.11	62.30	10.71	34.07	59.98
0.95	12.47	38.08	64.58	11.51	35.92	62.13	10.56	33.77	59.68
1.00	12.42	38.00	64.53	11.41	35.74	61.95	10.40	33.47	59.37

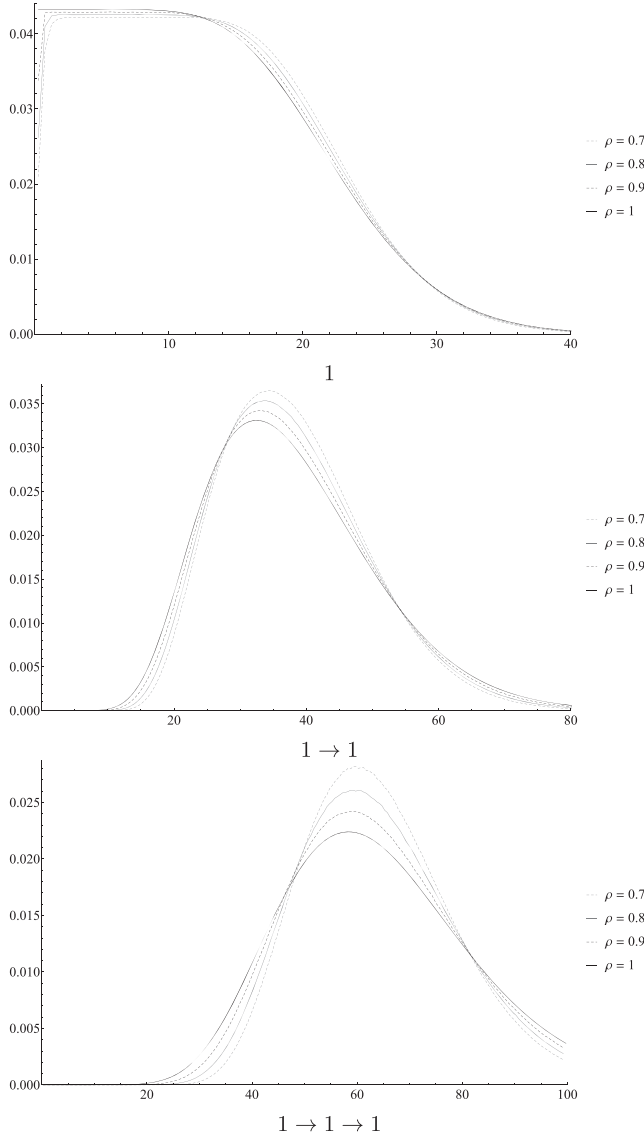


Figure 2. The probability densities for three selected paths of the first numerical example. The results for $\rho < 1$ are obtained using simulation. The results for $\rho = 1$ are obtained using the analytical results from [Section 3](#).

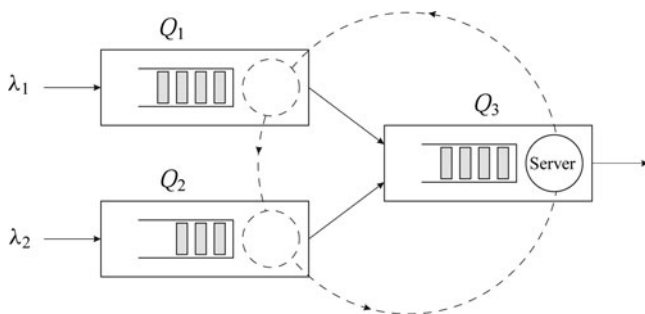
to Q_3 after being served (see [Figure 3](#)). This tandem queueing model with parallel queues in the first stage is a special case of a roving server network. We simply put $p_{1,3} = p_{2,3} = p_{3,0} = 1$ and all other $p_{i,j}$ are zero. We use the same values as in Katayama^[15]: Poisson arrivals with $\lambda_1 = \lambda_2/10$, service times are deterministic with $b_1 = b_2 = 1$, and $b_3 = 5$. The server serves the queues exhaustively, in cyclic order: 1, 2, 3, 1, The only difference with the model discussed in Katayama^[15] is that we introduce (deterministic) switch-over times $r_2 = r_3 = 2$. We assume that no time is required to switch between the two queues in the first stage, so $r_1 = 0$.

Table 3. Numerical results for Example 2.

Scaled path times: Means		
ρ	1 \rightarrow 3	2 \rightarrow 3
0.1	9.58	9.70
0.3	8.69	9.12
0.5	7.74	8.56
0.7	6.75	8.01
0.8	6.24	7.73
0.9	5.72	7.47
0.95	5.46	7.34
1.00	5.20	7.20
Scaled path times: Standard deviations		
ρ	1 \rightarrow 3	2 \rightarrow 3
0.1	1.91	2.02
0.3	2.75	3.04
0.5	3.30	3.75
0.7	3.74	4.30
0.8	3.94	4.54
0.9	4.13	4.76
0.95	4.23	4.88
1.00	4.32	4.97
Scaled path times: Tail probabilities		
ρ	1 \rightarrow 3	2 \rightarrow 3
0.1	0.002	0.002
0.3	0.006	0.008
0.5	0.008	0.013
0.7	0.010	0.018
0.8	0.010	0.020
0.9	0.011	0.022
0.95	0.011	0.023
1.00	0.011	0.023

5.2.1. Convergence to the HT limit

For the two possible paths, the means, standard deviations, and tail probabilities of the (scaled) path-time distributions are given in Table 3 and in Figure 4. The tail probabilities in Table 3 are given for path lengths of 20. The results for $\rho < 1$ are obtained using simulation, and the results for $\rho = 1$ are obtained using the theory presented in Section 3. It can be clearly observed that for loads larger than 0.8 the scaled path-time distributions are almost identical to the limiting distributions.

**Figure 3.** Tandem queues with parallel queues in the first stage, as discussed in Example 2.

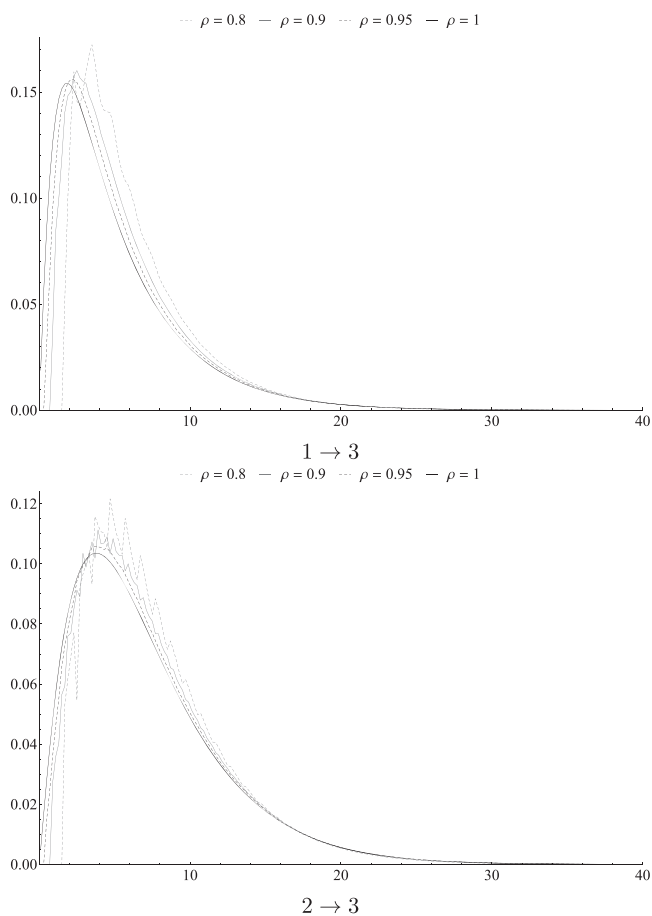


Figure 4. The probability densities for the two paths of the second numerical example. The results for $\rho < 1$ are obtained using simulation. The results for $\rho = 1$ are obtained using the analytical results from Section 3.

5.2.2. Accuracy of the approximation

Table 4 shows the approximated means using the path-time approximation discussed in Section 4 for both paths. Again the quality of the approximation is excellent

Table 4. Approximation results for Example 2. For the two paths, the means of the (scaled) path-time approximations are given.

ρ	Approximated mean scaled path times	
	1→3	2→3
0.1	9.52	9.72
0.3	8.56	9.16
0.5	7.60	8.60
0.7	6.64	8.04
0.8	6.16	7.76
0.9	5.68	7.48
0.95	5.44	7.34
1.00	5.20	7.20

for all possible loads and all paths: The relative approximation error is less than 2% for all values of ρ .

We have also tested the accuracy of the approximation for different interarrival-time distributions, with squared coefficient of variation (SCV) equal to respectively $\frac{1}{2}$ and 2. In the first case, we have fitted a mixed Erlang distribution, and in the second case a hyperexponential distribution. In these cases, the approximation still gives good results but slightly less accurate than for the case with Poisson arrivals: When the SCV equals 2, the relative approximation error is less than 4% for all values of ρ . When the SCV is equal to $\frac{1}{2}$, the relative error is less than 9%. As expected, the relative difference between the simulated and approximated values is biggest when ρ lies between 0.3 and 0.7.

5.3. Example 3: A file-server application

This application is taken from Sidi et al.^[30], who consider a token ring network with one file-server and K workstations, that transmit file requests to the file-server, which in turn replies to the different stations by sending back the files they requested. The performance measure of interest is the response time of a file request, which is the time from the request generation by the workstation (being queued at that time at the output queue of the station, awaiting its turn to be transmitted) until the file arrives back at the station. We can model this system as a roving server network with $K + 1$ queues. External customers arrive at Q_1, \dots, Q_K and are routed to Q_{K+1} after their service completion. Once served at Q_{K+1} the customer leaves the system. We are interested in the path time $W_{i,K+1}$ for $i = 1, \dots, K$. In our numerical example, we take $N = K + 1 = 11$, identical arrival rates at the K workstations and no external arrivals at the file server, i.e., $\lambda_{K+1} = 0$. Since arrival processes in this application area tend to exhibit high variation, we assume that the interarrival times follow a hyperexponential distribution with balanced means (see, e.g., Tijms^[36]) and squared coefficient of variation equal to 4. The routing probabilities are all 0 except $p_{i,K+1} = 1$ for $i = 1, \dots, K$. The service times B_1, \dots, B_K are all deterministic with value 0.1, and the service times at Q_{K+1} are exponentially distributed with mean 1. The switch-over times in this kind of system are typically very small compared to the service times, so we take exponentially distributed switch-over times with mean 0.01. The service discipline is gated at all queues.

We have chosen to highlight only a few typical results for this example, in [Table 5](#) and [Figure 5](#), where the means and standard deviations of the scaled path times are shown. The values for $\rho = 0, 9, 0.95, 0.98$ have been obtained by simulation, whereas the values for $\rho = 1$ are obtained using the analytical results from [Section 3](#). The most typical feature is that the densities of the waiting times and path times exhibit oscillating behavior for smaller values of ρ (see [Figure 5](#)). Most interestingly, this is not caused by simulation inaccuracy, but by the fact that the path times densities are (almost completely) concentrated on a discrete set of values for smaller loads due to the combination of small switch-over times, *deterministic* service times and *deterministic* routing. This effect disappears when the loads increase,

but as a consequence the convergence to the limiting HT distribution is not as fast as in the previous example. Nevertheless, despite the anomalous behavior of the system's performance for small loads, the derived asymptotic turns out to be exact again.

5.4. Example 4: A production system with a joint packaging queue and random yield

In the fourth example, we look at a production system with two types of products and a joint packaging queue and random yield. A set-up which is typically observed in many practical applications. Product type *A* requires three operations, at queue 1, queue 3, and queue 5, before joining the packaging queue 7, whereas the other product type *B* visits queue 2, queue 4, and queue 6 before it is sent to the joint packaging queue. Every production step fails with probability $p = 0.01$, after which the product is sent to the beginning of the production process, i.e., queue 1 for product *A* and queue 2 for product *B*. Packaging is, however, always successful. We assume that all queues are served exhaustively.

We study a system with external arrivals at queue 1 and queue 2 with intensity 0.12 and 0.04, respectively. Furthermore, we assume Erlang distributed arrivals with SCV 0.25, Erlang distributed service times with means 1 (for the regular queues) and 2 (for the packaging queue) and SCV 0.5, and switch-over times with constant value 5.

For two possible paths, the simulated scaled path-time distributions are depicted in Figure 6. The analytically computed limiting distribution at $\rho = 1$ is included as well. As in Example 1, due to the relatively large switch-over times, we see that for loads larger than 0.8 the scaled path-time distributions almost coincide with the limiting distributions.

Lastly, we show an illustration how the performance of the system depends on the rework probability p . Therefore, we vary p while keeping all the other parameters fixed. Note that increasing p will also increase the total workload of the system. When considering the load of the system as a function of p , the stability condition

Table 5. Numerical results for Example 3.

Scaled path times: Means										
ρ	1→11	2→11	3→11	4→11	5→11	6→11	7→11	8→11	9→11	10→11
0.9	1.25	1.38	1.51	1.64	1.77	1.89	2.02	2.14	2.26	2.37
0.95	1.39	1.55	1.71	1.86	2.02	2.17	2.31	2.46	2.60	2.73
0.98	1.47	1.65	1.82	2.00	2.17	2.33	2.50	2.66	2.81	2.96
1.00	1.53	1.73	1.92	2.12	2.32	2.51	2.71	2.91	3.10	3.30
Scaled path times: Standard deviations										
ρ	1→11	2→11	3→11	4→11	5→11	6→11	7→11	8→11	9→11	10→11
0.9	1.44	1.57	1.70	1.83	1.95	2.06	2.17	2.28	2.38	2.47
0.95	1.58	1.72	1.85	1.98	2.11	2.23	2.35	2.45	2.56	2.65
0.98	1.64	1.78	1.92	2.05	2.18	2.30	2.42	2.53	2.63	2.72
1.00	1.74	1.90	2.08	2.25	2.43	2.61	2.79	2.97	3.16	3.34

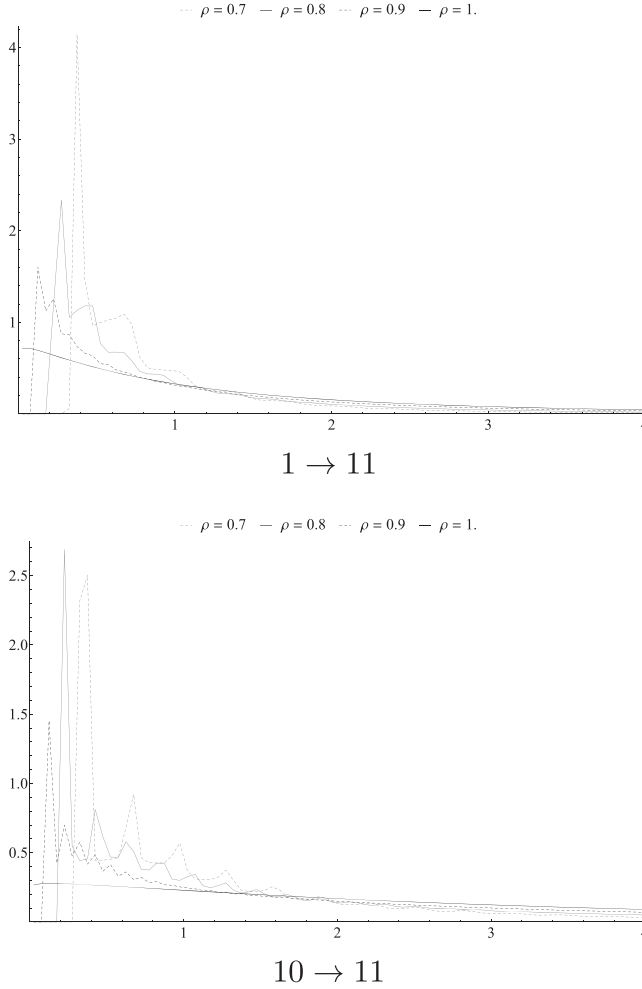


Figure 5. Numerical results for Example 3. The two figures show the densities of the scaled path-time distributions of the paths $1 \rightarrow 11$ and $10 \rightarrow 11$.

can be rewritten to $p < 0.1558$. Figure 7 shows the mean waiting time at queue 1 as function of this probability; this figure clearly indicates that the control of the rework probability is crucial for obtaining satisfactory system performance. Moreover, it shows that the approximation for the mean waiting time in Q_1 is extremely accurate. This is as expected, since $p = 0$ corresponds to a system with total load $\rho = 0.8$, and we have concluded before that the approximation is very accurate for $\rho > 0.8$. We emphasize that the approximations given in this article are closed-form expressions, once all parameter values have been substituted. For example, the approximation for $\mathbb{E}[W_1]$ simplifies to (21). As a consequence, the approximation, is very suitable for optimization purposes.

$$\mathbb{E}[W_1^{\text{approx}}] = (1088p^6 - 6374p^5 + 15814p^4 - 21412p^3 + 16492p^2 - 6406p + 670)^{-1} (192p^{11} - 2168p^{10} + 16940p^9 - 99647p^8$$

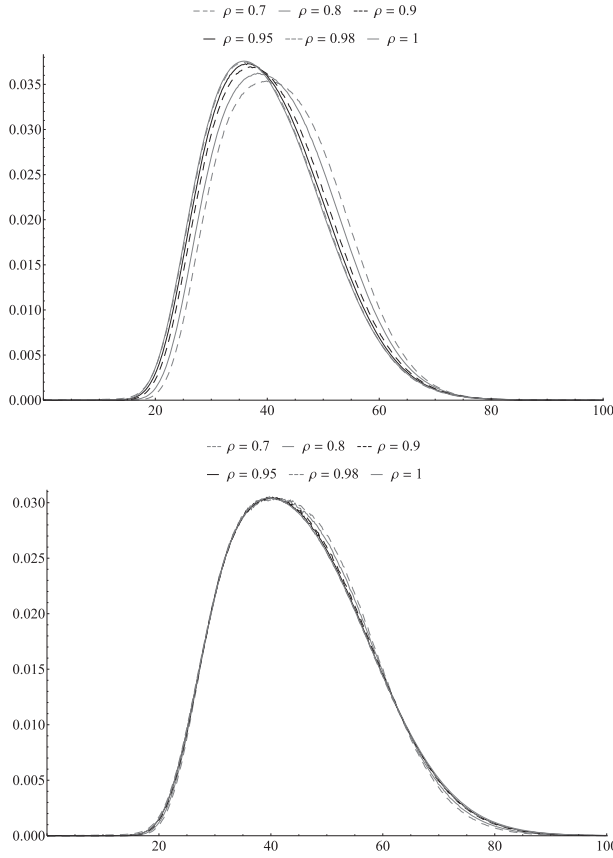


Figure 6. Numerical results for Example 4. The two figures show the densities of the scaled path-time distributions of the paths $1 \rightarrow 3 \rightarrow 5 \rightarrow 7$ and $2 \rightarrow 4 \rightarrow 6 \rightarrow 7$.

$$+378503p^7 - 880840p^6 + 1234922p^5 - 952526p^4 + 202040p^3 + 289849p^2 - 244631p + 52916). \quad (21)$$

We can therefore conclude that, although the four analyzed cases come from completely different fields of application, have dissimilar parameter settings and differ significantly in system's behavior, the derived HT asymptotics are shown to accurately describe the performance in all of these cases. Moreover, in all of these models the key performance metric is obviously the total time spent in the system by an arbitrary customer, i.e., the path time. We hope that this observation illustrates the general nature of our framework, which extends and unifies the HT analysis of many queueing models.

6. Conclusions and suggestions for further research

In the current article, we have analyzed a queueing network with customer routing, where a single shared server serves the queues in a cyclic order. We have not only studied the waiting time of a customer at a certain queue, but also the path time

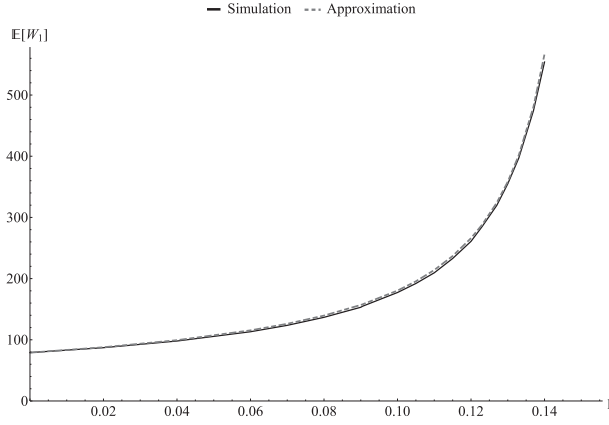


Figure 7. The mean waiting time in the first queue as a function of the rework probability p , for the model in Numerical Example 4.

(the total time spent in the system by an arbitrary customer traversing a specific path). The main complicating factor in this analysis is the routing of customers, which leads to non-renewal arrival processes at the queues and to strong interdependence of the waiting times at the queues. These factors prohibit an exact explicit analysis with closed-form expressions and, therefore, it is natural to resort to asymptotic estimates. That is, we have obtained easily computable expressions for both the waiting-time and path-time distribution in HT. Combining these HT asymptotics with newly developed LT limits leads to highly accurate approximations for the mean performance measures for the whole range of load values. The strength of this refined heavy-traffic approximation lies in its simplicity, which opens up interesting possibilities for optimization of the system performance with respect to the routes of the customers.

Appendix

A.1. Proof for Poisson arrivals

In this appendix, we present a proof for Conjecture 3.3.3.1 for the special case of Poisson arrivals, reformulated in the theorem below. We conclude the appendix with a short discussion.

Theorem A.1.1. *Assume that external customers arrive at Q_i according to a Poisson process with rate λ_i . As $\rho \uparrow 1$, the scaled queue length $(1 - \rho)L_i$ converges in distribution to the product of two independent random variables. The first has the same distribution as $\mathcal{L}_i^{\text{fluid}}$, and the second random variable Γ has the same distribution as the limiting distribution of the scaled length-biased cycle time, $(1 - \rho)C_i$. For $i = 1, \dots, N$; $k = i, \dots, i + N - 1$, and $\rho \uparrow 1$,*

$$(1 - \rho)L_i \xrightarrow{d} \Gamma \times \mathcal{L}_{i,k}^{\text{fluid}} \quad w.p. \ \hat{\rho}_k, \quad (\text{A.1})$$

where Γ is a random variable having a Gamma distribution with parameters $\alpha + 1$ and $\delta\mu$, and $\mathcal{L}_{i,k}^{\text{fluid}}$ is the “standardized” number of a type- i particles in the fluid model during V_k , introduced below (6), independent of Γ . The parameters of the Gamma distribution, in the case of Poisson arrivals, are defined as

$$\alpha = r\delta/\tilde{b}^{\text{res}}, \mu = \tilde{b}^{\text{res}},$$

with δ as defined in [Definition 3.1.1.1](#).

The proof of [Theorem A.1.1](#) is structured in the form of 9 small steps. Steps 1–7 involve finding the limiting scaled joint queue-length distribution at visit beginnings and completions. We rely heavily on existing limiting results for polling systems and branching processes, exploiting the similarities between these models and roving server networks, closely following the framework developed in van der Mei^[37]. The remaining steps follow the approach in Refs.^[6] and^[30], expressing queue lengths at arbitrary moments in terms of the queue-length distributions at visit beginnings, in heavy traffic. We start by giving some preliminary results, required for the remainder of the proof.

A.1.1. Preliminaries

As shown in Refs.^[6,30] the joint queue-length process at visit beginnings constitutes an N -dimensional multi-type branching process (MTBP) with immigration. As preliminaries, we give a brief overview of existing limiting results on critical MTBPs, required for the remainder of this proof. We start by introducing some additional notation. An N -dimensional vector \underline{v} has components (v_1, \dots, v_N) , and we denote $|\underline{v}| := \sum_{i=1}^N v_i$. Finally, I_E is the indicator function on the event E .

Let $\mathbf{Z} = \{\underline{Z}_n, n = 0, 1, \dots\}$ be an N -dimensional multi-type branching process, where $\underline{Z}_n = (Z_n^{(1)}, \dots, Z_n^{(N)})$ is an N -dimensional vector denoting the state of the process in the n th generation. The MTBP is defined by its offspring generating functions and its immigration functions. The one-step offspring generating function is denoted by $f(\underline{z}) = (f^{(1)}(\underline{z}), \dots, f^{(N)}(\underline{z}))$, with $\underline{z} = (z_1, \dots, z_N)$, and

$$f^{(i)}(\underline{z}) = \sum_{j_1, \dots, j_N \geq 0} p^{(i)}(j_1, \dots, j_N) z_1^{j_1} \cdots z_N^{j_N}, \quad (\text{A.2})$$

where $p^{(i)}(j_1, \dots, j_N)$ is the probability that a type- i particle produces j_k particles of type k ($i, k = 1, \dots, N$). The immigration function is denoted as follows:

$$g(\underline{z}) = \sum_{j_1, \dots, j_N \geq 0} q(j_1, \dots, j_N) z_1^{j_1} \cdots z_N^{j_N}, \quad (\text{A.3})$$

where $q(j_1, \dots, j_N)$ is the probability that a group of immigrant consists of j_k particles of type k ($i, k = 1, \dots, N$). Denote the *mean immigration vector* by

$$\underline{g} := (g_1, \dots, g_N), \text{ where } g_i := \left. \frac{\partial g(\underline{z})}{\partial z_i} \right|_{\underline{z}=\underline{1}} \quad (i = 1, \dots, N), \quad (\text{A.4})$$

and the *mean offspring matrix*, or simply *mean matrix*, by

$$\mathbf{M} = (m_{i,j}), \text{ with } m_{i,j} := \frac{\partial f^{(i)}(\underline{z})}{\partial z_j} \Big|_{\underline{z}=\underline{1}} \quad (i, j = 1, \dots, N). \quad (\text{A.5})$$

Thus, for a given type- i particle, $m_{i,j}$ is the mean number of type- j “children” it has in the next generation. Similarly, for a type- i particle, the second-order derivatives are denoted by the matrix

$$\mathbf{K}^{(i)} = (k_{j,k}^{(i)}), \text{ with } k_{j,k}^{(i)} := \frac{\partial^2 f^{(i)}(\underline{z})}{\partial z_j \partial z_k} \Big|_{\underline{z}=\underline{1}}, \quad (i, j, k = 1, \dots, N). \quad (\text{A.6})$$

Denote by $\underline{v} = (v_1, \dots, v_N)$ and $\underline{w} = (w_1, \dots, w_N)$ the left and right eigenvectors corresponding to the largest real-valued, positive eigenvalue ξ of \mathbf{M} , commonly referred to as the *maximum eigenvalue* (cf., e.g., Athreya and Ney^[2]), normalized such that

$$\underline{v}^\top \underline{1} = \underline{v}^\top \underline{w} = 1. \quad (\text{A.7})$$

The following conditions are necessary and sufficient conditions for the ergodicity of the process \mathbf{Z} (cf. Resing^[27]): $\xi < 1$ and

$$\sum_{j_1 + \dots + j_N > 0} q(j_1, \dots, j_N) \log(j_1 + \dots + j_N) < \infty. \quad (\text{A.8})$$

Note that ξ plays a role similar to ρ in the roving server network. The hat-notation introduced in Section 2 is also used here to indicate that x is *evaluated at* $\xi = 1$. Moreover, for $\xi \geq 0$ let

$$\pi_0(\xi) := 0, \text{ and } \pi_n(\xi) := \sum_{r=1}^n \xi^{r-2}, \quad n = 1, 2, \dots \quad (\text{A.9})$$

Quine [Ref.²³, theorem 4] derives the following property for critical MTBPs.

Property A.1.1.1. *Assume that all derivatives of $f(\underline{z})$ through order two exist at $\underline{z} = \underline{1}$ and that $0 < g_i < \infty$ ($i = 1, \dots, N$). Then,*

$$\frac{1}{\pi_n(\xi)} \begin{pmatrix} Z_n^{(1)} \\ \vdots \\ Z_n^{(N)} \end{pmatrix} \xrightarrow{d} A \begin{pmatrix} \hat{v}_1 \\ \vdots \\ \hat{v}_N \end{pmatrix} \Gamma(\alpha, 1) \text{ as } (\xi, n) \rightarrow (1, \infty), \quad (\text{A.10})$$

where $\underline{\hat{v}} = (\hat{v}_1, \dots, \hat{v}_N)$ is the normalized the left eigenvector of $\hat{\mathbf{M}}$, and where $\Gamma(\alpha, 1)$ is a gamma-distributed random variable with scale parameter 1 and shape parameter

$$\alpha := \frac{1}{A} \underline{\hat{g}}^\top \underline{\hat{w}} = \frac{1}{A} \sum_{i=1}^N \hat{g}_i \hat{w}_i, \text{ with } A := \sum_{i=1}^N \hat{v}_i \left(\underline{\hat{w}}^\top \hat{\mathbf{K}}^{(i)} \underline{\hat{w}} \right) > 0. \quad (\text{A.11})$$

A.1.2. Step 1: Characterize the embedded MTBP

We now return to the roving server network. In this step, we show how the joint queue-length process at successive polling instants at Q_1 can be described as an

MTBP with immigration in each state. To this end, let $\underline{X} := (X_1, \dots, X_N)$ be the N -dimensional vector that describes the joint queue length at an arbitrary polling instant at Q_1 . Let $X_{i,n}$ be the number of type- i customers in the system at the n th polling instant of the server at Q_1 , for $i = 1, \dots, N$ and $n = 0, 1, \dots$, and let

$$\underline{X}_n := (X_{1,n}, \dots, X_{N,n}) \quad (\text{A.12})$$

be the joint queue-length vector at the n th polling instant at Q_1 . In this appendix, we will provide the results for *gated service* only. The results for exhaustive service can be obtained along the exact same lines, but require significantly more complex notations (see, for example, Boon et al.^[6]) and are omitted here.

The following result describes the process $\{\underline{X}_n, n = 0, 1, \dots\}$ as an MTBP.

Theorem A.1.2.1. *The discrete-time process $\{\underline{X}_n, n = 0, 1, \dots\}$ constitutes a N -dimensional MTBP with immigration in each state. Denote by $B_i^*(s)$ and $R_i^*(s)$ the Laplace–Stieltjes transforms (LSTs) of, respectively, the service times B_i and the switch-over times R_i . The probability generating function (PGF) of the offspring function is given by the following expression: For $|z_i| \leq 1$ ($i = 1, \dots, N$),*

$$f(\underline{z}) := (f^{(1)}(\underline{z}), \dots, f^{(N)}(\underline{z})), \quad (\text{A.13})$$

where for $i = 1, \dots, N$ the function $f^{(i)}(\underline{z})$ is the unique solution of

$$f^{(i)}(\underline{z}) := B_i^* \left(\sum_{j=1}^N \lambda_j (1 - z_j) \right) \cdot \sum_{j=1}^N p_{i,j} f^{(j)}(\underline{z}), \quad (\text{A.14})$$

and where the PGF of the immigration function is given by

$$g(\underline{z}) := \prod_{i=1}^N R_i^* \left(\sum_{j=1}^i \lambda_j (1 - z_j) + \sum_{j=i+1}^N \lambda_j (1 - f^{(j)}(\underline{z})) \right). \quad (\text{A.15})$$

Proof. Relations (A.13)–(A.15) can be obtained along the lines of Resing^[27] for the case of classical gated service, using simple generating-function manipulations. The only difference is that the offspring function $f^{(i)}(\underline{z})$ has changed in the sense that each customer served at Q_i is effectively replaced *not only* by all customers that arrive at the different queues during its services time (leading to PGF $B_i^*(\sum_{j=1}^N \lambda_j (1 - z_j))$ in Equation (A.14), but *in addition* by a fresh customer at Q_j (which creates an additional offspring $f^{(j)}(\underline{z})$), with probability $p_{i,j}$. Next, Equation (A.15) stems from the fact that the immigration consists of the contributions of newly arriving customers that arrive during the switch-over times R_i , $i = 1, \dots, N$. \square

A.1.3. Step 2: Compute the mean offspring matrix for the roving server network

The following result gives a characterization of the mean offspring matrix \mathbf{M} defined in (A.5).

Lemma A.1.3.1. The mean offspring matrix $\mathbf{M} = (m_{i,j})$ can be expressed by

$$\mathbf{M} = \mathbf{M}_1 \cdots \mathbf{M}_N, \quad (\text{A.16})$$

where for $k = 1, \dots, N$, the elements of the matrix $\mathbf{M}_k = (m_{i,j}^{(k)})$ are given by: For $i, j = 1, \dots, N, i \neq k$,

$$m_{i,j}^{(k)} = I_{\{i=j\}}, \quad (\text{A.17})$$

and for $i = k$,

$$m_{i,j}^{(k)} = \lambda_j b_i + p_{i,j} \text{ for } j = 1, \dots, N. \quad (\text{A.18})$$

Proof. The result can be obtained directly from [Theorem A.1.2.1](#) by taking the partial derivatives of the offspring function defined in [\(A.13\)](#) and [\(A.14\)](#). \square

A.1.4. Step 3: Determine the left and right eigenvectors of \mathbf{M} at $\rho = 1$

The following result gives the left and right eigenvectors (normalized according to [\(A.7\)](#)) of the mean offspring matrix \mathbf{M} , defined in [\(A.16\)](#)–[\(A.18\)](#), evaluated at $\rho = 1$.

Lemma A.1.4.1. The right eigenvector $\underline{\hat{w}}$ of the mean matrix $\hat{\mathbf{M}}$, normalized such that $\underline{\hat{w}}^\top \underline{\hat{w}} = 1$, corresponding with maximum eigenvalue 1, is given by

$$\underline{\hat{w}} = \begin{pmatrix} \hat{w}_1 \\ \vdots \\ \hat{w}_N \end{pmatrix} := |\underline{\tilde{b}}|^{-1} \underline{\tilde{b}}, \text{ with } \underline{\tilde{b}} := \begin{pmatrix} \tilde{b}_1 \\ \vdots \\ \tilde{b}_N \end{pmatrix}. \quad (\text{A.19})$$

The corresponding left eigenvector $\underline{\hat{v}}$, normalized such that $\underline{\hat{v}}^\top \underline{\hat{w}} = 1$, corresponding with maximum eigenvalue 1, is given by

$$\underline{\hat{v}} = \begin{pmatrix} \hat{v}_1 \\ \vdots \\ \hat{v}_N \end{pmatrix} := \frac{|\underline{\tilde{b}}|}{\delta} \underline{\hat{u}}, \text{ with } \underline{\hat{u}} := \begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix}, \text{ where } u_i := \lambda_i \sum_{j=i}^N \rho_j + \sum_{j=i}^N \gamma_j p_{j,i}, \quad (\text{A.20})$$

and where

$$\delta := \underline{\hat{u}}^\top \underline{\tilde{b}} = \sum_{i=1}^N \sum_{j=i+1}^N \hat{\rho}_i \hat{\rho}_j + \sum_{i=1}^N \tilde{b}_i \sum_{j=i}^N \hat{\gamma}_j p_{j,i}. \quad (\text{A.21})$$

Proof. First, it is readily seen by using equations [\(A.17\)](#) and [\(A.18\)](#) and [\(1\)](#) that for $k = 1, \dots, N$, we have

$$\sum_{j=1}^N \hat{m}_{k,j}^{(k)} \tilde{b}_j = \sum_{j=1}^N \left(\hat{\lambda}_j b_k + p_{k,j} \right) \tilde{b}_j = b_k \sum_{j=1}^N \hat{\lambda}_j \tilde{b}_j + \sum_{j=1}^N p_{k,j} \tilde{b}_j = b_k + \sum_{j=1}^N p_{k,j} \tilde{b}_j = \tilde{b}_k. \quad (\text{A.22})$$

This immediately implies that $\hat{\mathbf{M}}_k \hat{\mathbf{w}} = \hat{\mathbf{w}}$ for $k = 1, \dots, N$, and hence from (A.16) that $\hat{\mathbf{M}} \hat{\mathbf{w}} = \hat{\mathbf{M}}_1 \cdots \hat{\mathbf{M}}_N \hat{\mathbf{w}} = \hat{\mathbf{w}}$, which shows that $\hat{\mathbf{w}}$ indeed is a right eigenvector of $\hat{\mathbf{M}}$ with eigenvalue 1. Similar arguments can be used to show that $\hat{\mathbf{M}}^\top \hat{\mathbf{v}} = \hat{\mathbf{v}}$. \square

Remark A.1.4.1. Although δ defined in (A.21) may appear different, at first sight, than δ defined in Definition 3.1.1.1, it can be shown that they are in fact identical.

A.1.5. Step 4: Mean immigration function at $\rho = 1$

We now proceed to specify the mean immigration vector \underline{g} , defined in (A.7), for the model under consideration. Considering the evolution of the N -dimensional state vector as a discrete-time Markov chain $\{\underline{X}_n, n = 0, 1, \dots\}$ at successive polling instants at Q_1 , the “immigrants” in the n th generation are the customers present a time n that are not children of any of the customers present at time $n - 1$. Denote the mean immigration vector by $\underline{g} = (g_1, \dots, g_N)$, where g_i stands for the mean number of type- i immigrants.

Lemma A.1.5.1. The mean immigration function is given by $\underline{g} = (g_1, \dots, g_N)$, where for $j = 1, \dots, N$,

$$g_j = \sum_{i=1}^N r_i \left(\lambda_j I_{\{j \leq i\}} + \sum_{k=i+1}^N \lambda_k m_{k,j} \right). \quad (\text{A.23})$$

Moreover,

$$\hat{\underline{g}}^\top \hat{\underline{w}} = |\underline{b}|^{-1} r. \quad (\text{A.24})$$

Proof. Equation (A.23) follows directly from Theorem A.1.2.1. by differentiating once with respect to s_j and substituting $\underline{s} = (1, \dots, 1)$. Next, to prove (A.24), assume $\rho = 1$. We first observe that it follows from (A.23) that the mean number of type- j customers that immigrate during a cycle is given by

$$\hat{g}_j = \sum_{i=1}^N r_i \left(\hat{\lambda}_j I_{\{j \leq i\}} + \sum_{k=i+1}^N \hat{\lambda}_k \hat{m}_{k,j} \right). \quad (\text{A.25})$$

This implies

$$\hat{\underline{g}}^\top \hat{\underline{w}} := |\tilde{\underline{b}}|^{-1} \sum_{j=1}^N \hat{g}_j \tilde{b}_j = |\tilde{\underline{b}}|^{-1} \sum_{i=j}^N r_i \left(\sum_{j=1}^i \tilde{b}_j \hat{\lambda}_j + \sum_{k=i+1}^N \hat{\lambda}_k \sum_{j=1}^N \hat{m}_{k,j} \tilde{b}_j \right) \quad (\text{A.26})$$

$$= |\tilde{\underline{b}}|^{-1} \sum_{i=1}^N r_i \left(\sum_{j=1}^i \tilde{b}_j \hat{\lambda}_j + \sum_{k=i+1}^N \hat{\lambda}_k \tilde{b}_k \right) = |\tilde{\underline{b}}|^{-1} r \sum_{i=1}^N \hat{\rho}_i = |\tilde{\underline{b}}|^{-1} r, \quad (\text{A.27})$$

by using the definition in (A.23) and the results in Lemma A.1.4.1. \square

A.1.6. Step 5: Determine the parameter A

The following result gives an expression for the scaling parameter A , defined in (A.11).

Lemma A.1.6.1.

$$A = |\underline{b}|^{-1} \delta^{-1} \cdot \tilde{b}^{\text{res}}, \quad (\text{A.28})$$

where $\tilde{b}^{\text{res}} = \frac{\mathbb{E}[\tilde{B}^2]}{2\mathbb{E}[\tilde{B}]}$ denotes the expected residual extended service time \tilde{B} of an *arbitrary* customer in the system. More precisely,

$$\mathbb{E}[\tilde{B}^k] = \sum_{i=1}^N \lambda_i \mathbb{E}[\tilde{B}_i^k] / \sum_{j=1}^N \lambda_j.$$

Proof. This result can be proven following the same lines as in van der Mei^[37] (Remark 5) for the case of branching-type service policies for the case without customer routing (i.e., $p_{i,j} = 0$ for all i, j). \square

A.1.7. Step 6: Asymptotic properties for the maximum eigenvalue of M

The following result describes the limiting behavior of the maximum eigenvalue $\xi(\rho)$ of the matrix M defined in Lemma A.1.3.1, considered as a function of ρ , as ρ goes to 1.

Lemma A.1.7.1. The maximum eigenvalue $\xi = \xi(\rho)$ satisfies the following properties:

- (1) $\xi < 1$ if and only if $\rho < 1$, $\xi = 1$ if and only if $\rho = 1$ and $\xi > 1$ if and only if $\rho > 1$;
- (2) $\xi(\rho)$ is a continuous function of ρ ;
- (3) $\lim_{\rho \uparrow 1} \xi(\rho) = \xi(1) = 1$;
- (4) the derivative of $\xi(\rho)$ at $\rho = 1$ is given by

$$\xi'(1) = \lim_{\rho \uparrow 1} \frac{1 - \xi(\rho)}{1 - \rho} = \frac{1}{\delta}, \quad (\text{A.29})$$

where δ is defined in (A.21).

Proof. The proof can be obtained by following the approach in van der Mei^[37] (Section 3.2), which proves for the case without customer routing (i.e., $p_{i,j} = 0$ for all i, j). \square

A.1.8. Step 7: The joint queue-length vector at visit beginnings and completions

We are now ready to present the HT result for the state vector at polling instants. Without loss of generality, we focus on the evolution of the state vector at embedded polling instants at Q_1 .

Theorem A.1.8.1. *The joint queue-length vector at polling instants at Q_1 has the following asymptotic behavior:*

$$(1 - \rho) \begin{pmatrix} X_1 \\ \vdots \\ X_N \end{pmatrix} \xrightarrow{d} \tilde{b}^{\text{res}} \frac{1}{\delta} \begin{pmatrix} \hat{u}_1 \\ \vdots \\ \hat{u}_N \end{pmatrix} \Gamma(\alpha, 1) \quad (\rho \uparrow 1), \quad (\text{A.30})$$

where

$$\alpha = r\delta/\tilde{b}^{\text{res}}, \quad (\text{A.31})$$

and where \hat{u}_i ($i = 1, \dots, N$) and δ are defined in (A.20) and (A.21).

Proof. First, note that the process that describes the evolution of the state vector $\{\underline{X}_n, n = 0, 1, \dots\}$ at successive polling instants at Q_1 constitutes an N -dimensional MTBP with offspring function $f(s)$ and immigration function $g(z)$ defined in Theorem A.1.2.1., and with mean matrix \mathbf{M} defined in Lemma A.1.3.1. Moreover, from Theorem A.1.2.1. it is readily verified that the assumptions of Property 1 on the finiteness of the second-order derivatives of $f(z)$ and the mean immigration function g are satisfied. Then, using Property A.1.1.1. it follows that

$$\frac{1}{\pi_n(\xi)} \cdot \underline{X}_n^\top \xrightarrow{d} A \cdot \hat{\underline{v}} \cdot \Gamma(\alpha, 1) \text{ as } (\xi, n) \rightarrow (1, \infty), \quad (\text{A.32})$$

where A , $\hat{\underline{v}}$, and α are given in (A.11). Hence, translating this to the polling model and using Lemmas A.1.3.1–A.1.7.1, it readily follows from (A.32) that

$$(1 - \rho) \underline{X}_n^\top \xrightarrow{d} \delta \cdot A \cdot \hat{\underline{v}} \cdot \Gamma(\alpha, 1) \text{ as } (\rho, n) \rightarrow (1, \infty), \quad (\text{A.33})$$

where expressions for δ , A , $\hat{\underline{v}}$, and α are given in (A.21), (A.28), (A.20), and (A.31). Combining these expressions leads to the result. \square

Corollary A.1.8.1. *Let $X_{i,k}^{\text{scaled}} := \lim_{\rho \uparrow 1} (1 - \rho) X_{i,k}$ denote the scaled number of customers in Q_i at the beginning of a visit to Q_k . Its LST is equal to*

$$\mathbb{E} \left[e^{-\omega X_{i,k}^{\text{scaled}}} \right] = \begin{cases} \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \sum_{j=i}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j}} \right)^{r\delta/\tilde{b}^{\text{res}}} & k = i+1, \dots, i+N-1, \\ \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \hat{\gamma}_i} \right)^{r\delta/\tilde{b}^{\text{res}}} & k = i. \end{cases} \quad (\text{A.34})$$

A.1.9. Step 8: The marginal queue-length distribution

Let $L_{i,k}^{\text{scaled}} := \lim_{\rho \uparrow 1} (1 - \rho) L_{i,k}$ denote the scaled number of customers in Q_i at an arbitrary moment during a visit to Q_k , and let $L_i^{\text{scaled}} := \lim_{\rho \uparrow 1} (1 - \rho) L_i$ denote the scaled number of customers in Q_i at an arbitrary moment.

Theorem A.1.9.1.

$$\mathbb{E} \left[e^{-\omega L_i^{\text{scaled}}} \right] = \sum_{k=1}^N \hat{\rho}_k \mathbb{E} \left[e^{-\omega L_{i,k}^{\text{scaled}}} \right], \quad (\text{A.35})$$

where

$$\mathbb{E} \left[e^{-\omega L_{i,i}^{\text{scaled}}} \right] = \frac{\left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \hat{\rho}_i \hat{\gamma}_{i,i}} \right)^{r\delta/\tilde{b}^{\text{res}}} - \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \hat{\gamma}_i} \right)^{r\delta/\tilde{b}^{\text{res}}}}{(\hat{\gamma}_i - \hat{\gamma}_{i,i} \hat{\rho}_i) r \omega}, \quad (\text{A.36})$$

$$\mathbb{E} \left[e^{-\omega L_{i,k}^{\text{scaled}}} \right] = \frac{\left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \sum_{j=1}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j}} \right)^{r\delta/\tilde{b}^{\text{res}}} - \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \sum_{j=1}^k \hat{\rho}_j \hat{\gamma}_{i,j}} \right)^{r\delta/\tilde{b}^{\text{res}}}}{\hat{\rho}_k \hat{\gamma}_{i,k} r \omega}, \quad i \neq k. \quad (\text{A.37})$$

Proof. To prove [Theorem A.1.9.1](#), we need the following results, summarized in [\(A.38\)–\(A.41\)](#), which directly follow from Equations (3.6)–(3.10) in Sidi et al.^[30].

$$\mathbb{E}[z^{L_i}] = \sum_{k=1}^N \left(\hat{\rho}_k \mathbb{E}[z^{L_i^{(V_k)}}] + (1 - \rho) \frac{r_j}{r} \mathbb{E}[z^{L_i^{(R_k)}}] \right), \quad (\text{A.38})$$

where $L_i^{(V_k)}$ and $L_i^{(R_k)}$ denote the number of customers in Q_i at an arbitrary moment during V_k and R_k , respectively.

Denote by $X_{i,k}^*(z)$ and $Y_{i,k}^*(z)$ the PGF of the number of customers in Q_i at the beginning and end of V_k , respectively.

$$L_i^{(V_i)}(z) = (X_{i,i}^*(z) - Y_{i,i}^*(z)) \frac{(1 - \rho)(1 - B_i^*(\lambda_i(1 - z)))}{\lambda_i(1 - z) \rho_i r (1 - B_i^*(\lambda_i(1 - z))(1 - p_{i,i} + p_{i,i}z)/z)}, \quad (\text{A.39})$$

$$L_i^{(V_k)}(z) = (X_{i,k}^*(z) - Y_{i,k}^*(z)) \frac{(1 - \rho)(1 - B_k^*(\lambda_i(1 - z)))}{\lambda_i(1 - z) \rho_k r (1 - B_k^*(\lambda_i(1 - z))(1 - p_{k,i} + p_{k,i}z))}, \quad \text{for } k \neq i. \quad (\text{A.40})$$

Note that

$$Y_{i,k}^*(z) = \frac{X_{i,k+1}^*(z)}{R_k^*(\lambda_i(1 - z))}. \quad (\text{A.41})$$

Our goal is to find the limiting distribution of $(1 - \rho)L_i$ as $\rho \uparrow 1$. In the limit, when substituting $z = e^{-\omega(1-\rho)}$ in [\(A.38\)](#) and letting $\rho \uparrow 1$, the terms that correspond to the queue lengths during switch-over times vanish, caused by the $(1 - \rho)$ term. Intuitively this is exactly what one would expect, as switch-over times become negligible in heavy traffic. In order to prove [\(A.36\)](#) and [\(A.37\)](#), we substitute $z = e^{-\omega(1-\rho)}$, $\lambda_i = \hat{\lambda}_i \rho$, and $\rho_i = \hat{\rho}_i \rho$ in [\(A.38\)](#) and [\(A.40\)](#), respectively, and evaluate the Taylor series of the resulting functions near $\rho = 1$.

For the case $k = i$ this results in

$$\lim_{\rho \uparrow 1} \mathbb{E}[e^{-\omega(1-\rho)L_i^{(V_i)}}] = \frac{\mathbb{E}[e^{-\omega X_{i,i+1}^{\text{scaled}}}] - \mathbb{E}[e^{-\omega X_{i,i}^{\text{scaled}}}]}{\omega r \hat{\rho}_i (1 - \hat{\lambda}_i b_i - p_{i,i}) / b_i}, \quad (\text{A.42})$$

and for the case $k \neq i$, we obtain

$$\lim_{\rho \uparrow 1} \mathbb{E}[e^{-\omega(1-\rho)L_i^{(V_k)}}] = \frac{\mathbb{E}[e^{-\omega X_{i,k}^{\text{scaled}}}] - \mathbb{E}[e^{-\omega X_{i,k+1}^{\text{scaled}}}]}{\omega r \hat{\rho}_k (\hat{\lambda}_i b_i + p_{k,i}) / b_k}. \quad (\text{A.43})$$

After substitution of (A.34) this leads to the following result:

$$\lim_{\rho \uparrow 1} \mathbb{E}[e^{-\omega(1-\rho)L_i^{(V_i)}}] = \frac{\left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \hat{\rho}_i \hat{\gamma}_{i,i} \omega}\right)^{r\delta/\tilde{b}^{\text{res}}} - \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \hat{\gamma}_i \omega}\right)^{r\delta/\tilde{b}^{\text{res}}}}{r(\hat{\gamma}_i - \hat{\rho}_i \hat{\gamma}_{i,i})\omega}, \quad (\text{A.44})$$

$$\lim_{\rho \uparrow 1} \mathbb{E}[e^{-\omega(1-\rho)L_i^{(V_k)}}] = \frac{\left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \sum_{j=i-N}^{k-1} \hat{\rho}_j \hat{\gamma}_{i,j}}\right)^{r\delta/\tilde{b}^{\text{res}}} - \left(\frac{\delta/\tilde{b}^{\text{res}}}{\delta/\tilde{b}^{\text{res}} + \omega \sum_{j=i-N}^k \hat{\rho}_j \hat{\gamma}_{i,j}}\right)^{r\delta/\tilde{b}^{\text{res}}}}{r \hat{\rho}_k \hat{\gamma}_{i,k} \omega}. \quad (\text{A.45})$$

□

A.1.10. Step 9: Proving Theorem A.1.1

To prove that Theorem A.1.1 agrees with the results obtained in Step 8, we need the following, well-known property that is frequently used in HT limits of polling systems.

Property A.1.10.1. Let Γ be a random variable with a Gamma distribution with shape parameter α^* (to avoid confusion since we already have defined a parameter α) and rate parameter β . Let U be a uniform random variable on the interval $[a, b]$, independent of Γ . The LST of the product $U \times I$ is equal to

$$\mathbb{E}[e^{-\omega U I}] = \frac{\left(\frac{\beta/a}{\beta/a + \omega}\right)^{\alpha^* - 1} - \left(\frac{\beta/b}{\beta/b + \omega}\right)^{\alpha^* - 1}}{(b - a)(\alpha^* - 1)\omega/\beta}. \quad (\text{A.46})$$

Theorem A.1.1 states that the scaled queue length L_i^{scaled} is distributed as $L_{i,k}^{\text{scaled}}$ with probability $\hat{\rho}_k$. This is in accordance with what one would obtain when substituting $z = e^{-\omega(1-\rho)}$ in (A.38) and letting $\rho \uparrow 1$.

The distribution of $L_{i,k}^{\text{scaled}}$ is the product of two random variables. The first is a $\text{Gamma}(\alpha + 1, \delta\mu)$ distribution, where $\alpha + 1$ and $\delta\mu$ are the parameters of the Gamma distribution of the scaled, length-biased cycle time, discussed in Theorem A.1.1. When the external arrival process is a Poisson process, the parameter σ^2 is equal to $\tilde{b}^{(2)}/\tilde{b}$, which means that

$$\alpha = r\delta/\tilde{b}^{\text{res}}, \mu = \tilde{b}^{\text{res}}, \quad (\text{A.47})$$

with δ as defined in Definition 3.1.1.1.

The second random variable in the product is $\mathcal{L}_{i,k}^{\text{fluid}}$, which is uniformly distributed. Remember that $\mathcal{L}_{i,k}^{\text{fluid}} = L_{i,k}^{\text{fluid}}/c$, so taking $c = 1$ in (5) gives the parameters of the uniform distribution for each of the $\mathcal{L}_{i,k}^{\text{fluid}}$. Substituting the following values in (A.46) leads to the LST of the limiting distribution of the scaled number of customers in Q_i at an arbitrary epoch during V_k ,

$$\begin{aligned} \alpha^* &= r\delta/\tilde{b}^{\text{res}} + 1, & \beta &= \delta\tilde{b}^{\text{res}}, \\ a &= \hat{\rho}_j\hat{\gamma}_{i,j}, & b &= \hat{\gamma}_i, & (i = k) \\ a &= \sum_{j=i}^{k-1} \hat{\rho}_j\hat{\gamma}_{i,j}, & b &= \sum_{j=i}^k \hat{\rho}_j\hat{\gamma}_{i,j}, & (i \neq k) \end{aligned}$$

It is quickly verified that these substitutions result in an expression completely equivalent to (A.44) for the case $i = k$, and to (A.45) for $i \neq k$. This concludes the proof of Theorem A.1.1.

A.2. Discussion

In this appendix, we have provided a proof for the scaled queue-length distributions in heavy traffic, relying on the framework developed in van der Mei^[37]. In his article, Van der Mei uses the distributional form of Little's Law to obtain the distributions of the scaled waiting times in polling models with Poisson arrivals. We note that the distributional form of Little's Law cannot be applied to our model because of the internal routing, as discussed extensively in Boon et al.^[6]. In Boon et al.^[6], nevertheless a mathematical framework has been developed to derive the LSTs of the steady-state waiting-time distributions. As such, this framework provides an excellent basis to prove the HT limits for waiting times without resorting to the distributional form of Little's Law. The derived LSTs for steady-state waiting-time distributions are given in the form of recursive expressions, which will simplify to the elegant closed-form expressions presented in this article, after taking the HT limit.

For *path times*, there are currently no steady-state results available at all, due to the complex dependencies between successive visit times. In order to prove the HT limit of the scaled path times, one would first have to apply the techniques in Boon et al.^[6] to find path-time LSTs in steady state, which is a separate study by itself.

We conclude this discussion by noting that the *mean* waiting times can easily be obtained by applying Little's Law, which does not require the assumption of Poisson arrivals.

Acknowledgments

The authors are grateful to the anonymous referees, who have provided valuable comments resulting in improved readability of the current manuscript.

References

- [1] Ali, O. M. E.; Neuts, M. F. A service system with two stages of waiting and feedback of customers. *J. Appl. Probab.* **1984**, *21*, 404–413.
- [2] Athreya, K. B.; Ney, P. E. *Branching Processes*; Springer-Verlag: Berlin, 1972.
- [3] Boon, M. A. A.; van der Mei, R. D.; Winands, E. M. M., Applications of polling systems. *Sur. Oper. Res. Manag. Sci.* **2011**, *16*, 67–82.
- [4] Boon, M. A. A.; van der Mei, R. D.; Winands, E. M. M. Queueing networks with a single shared server: light and heavy traffic. *SIGMETRICS Perfo. Eval. Rev.* **2011**, *39*(2), 44–46.
- [5] Boon, M. A. A.; Winands, E. M. M.; Adan, I. J. B. F.; van Wijk, A. C. C. Closed-form waiting time approximations for polling systems. *Perfor. Eval.* **2011**, *68*, 290–306.
- [6] Boon, M. A. A.; van der Mei, R. D.; Winands, E. M. M. Waiting times in queueing networks with a single shared server. *Que. Syst.* **2013**, *74*(4), 403–429.
- [7] Boxma, O. J.; Yechiali, U., An $M/G/1$ queue with multiple types of feedback and gated vacations. *J. App. Prob.* **1997**, *34*, 773–784.
- [8] Coffman, Jr, E. G.; Puhalskii, A. A.; Reiman, M. I., Polling systems with zero switchover times: A heavy-traffic averaging principle. *Ann. App. Prob.* **1995**, *5*(3), 681–719.
- [9] Coffman, Jr, E. G.; Puhalskii, A. A.; Reiman, M. I., Polling systems in heavy-traffic: A Bessel process limit. *Math. Oper. Res.* **1998**, *23*, 257–304.
- [10] Dorsman, J. L.; van der Mei, R. D.; Winands, E. M. M., A new method for deriving waiting-time approximations in polling systems with renewal arrivals. *Stoc. Models* **2011**, *27*(2), 318–332.
- [11] Gong, Y.; de Koster, R., A polling-based dynamic order picking system for online retailers. *IIE Trans.* **2008**, *40*, 1070–1082.
- [12] Grasman, S. E.; Olsen, T. L.; Birge, J. R. Setting basestock levels in multiproduct systems with setups and random yield. *IIE Trans.* **2008**, *40*(12), 1158–1170.
- [13] Grillo, D. Polling mechanism models in communication systems – some application examples. In *Stochastic Analysis of Computer and Communication Systems*; Takagi, H., Eds.; North-Holland: Amsterdam, 1990; 659–699.
- [14] Jennings, O. B. Averaging principles for a diffusion-scaled, heavy-traffic polling station with K job classes. *Math. Oper. Res.* **2010**, *35*(3), 669–703.
- [15] Katayama, T. A cyclic service tandem queueing model with parallel queues in the first stage. *Stoc. Models* **1988**, *4*, 421–443.
- [16] Kavitha, V.; Altman, E. Queueing in space: design of message ferry routes in static ad hoc networks. In *21st International Teletraffic Congress, ITC 21 2009*. Paris, France, 2009; 1–8.
- [17] Levy, H.; Sidi, M. Polling systems: applications, modeling, and optimization. *IEEE Trans. Comm.* **1990**, *38*, 1750–1760.
- [18] Markowitz, D.; Wein, L. Heavy traffic analysis of dynamic cyclic policies: A unified treatment of the single machine scheduling problem. *Oper. Res.* **2001**, *49*(2), 246–270.
- [19] Markowitz, D.; Reiman, M.; Wein, L. The stochastic economic lot scheduling problem: Heavy traffic analysis of dynamic cyclic policies. *Oper. Res.* **2000**, *48*(2), 136–154.
- [20] Nair, S. S. A single server tandem queue. *J. App. Prob.* **1971**, *8*(1), 95–109.
- [21] Olsen, T. L.; van der Mei, R. D. Polling systems with periodic server routeing in heavy traffic: distribution of the delay. *J. App. Prob.* **2003**, *40*, 305–326.
- [22] Olsen, T. L.; van der Mei, R. D. Polling systems with periodic server routing in heavy traffic: renewal arrivals. *Oper. Res. Lett.* **2005**, *33*, 17–25.
- [23] Quine, M. P. The multi-type Galton-Watson process with immigration. *J. App. Prob.* **1970**, *7*(2), 411–422.
- [24] Reiman, M.; Wein, L. Dynamic scheduling of a two-class queue with setups. *Oper. Res.* **1998**, *46*(4), 532–547.

- [25] Reiman, M.; Wein, L. Heavy traffic analysis of polling systems in tandem. *Oper. Res.* **1999**, 47(4), 524–534.
- [26] Reiman, M.; Rubio, R.; Wein, L. Heavy traffic analysis of the dynamic stochastic inventory-routing problem. *Trans. Sci.* **1999**, 33(4), 361–380.
- [27] Resing, J. A. C. Polling systems and multitype branching processes. *Que. Syst.* **1993**, 13, 409–426.
- [28] Sarkar, D.; Zangwill, W. I. File and work transfers in cyclic queue systems. *Manag. Sci.* **1992**, 38(10), 1510–1523.
- [29] Sidi, M.; Levy, H. Customer routing in polling systems. In *Proceedings Performance '90*; King, P., Mitrani, I., Pooley, R., Eds.; North-Holland: Amsterdam, 1990; 319–331.
- [30] Sidi, M.; Levy, H.; Fuhrmann, S. W. A queueing network with a single cyclically roving server. *Que. Syst.* **1992**, 11, 121–144.
- [31] Takács, L. A queueing model with feedback. *Revue française d'automatique, d'informatique et de recherche opérationnelle. Recherche opérationnelle* **1977**, 11(4), 345–354.
- [32] Takagi, H. Analysis and applications of a multiqueue cyclic service system with feedback. *IEEE Trans. Comm. - TCOM* **1987**, 35(2), 248–250.
- [33] Takagi, H. Analysis and application of polling models. In *Performance Evaluation: Origins and Directions volume 1769 of Lecture Notes in Computer Science* Haring, G., Lindemann, C., Reiser, M., Eds.; Springer Verlag: Berlin, 2000; 424–442.
- [34] Takine, T.; Takagi, H.; Hasegawa, T. Sojourn times in vacation and polling systems with Bernoulli feedback. *J. App. Prob.* **1991**, 28(2), 422–432.
- [35] Taube-Netto, M. Two queues in tandem attended by a single server. *Oper. Res.* **1977**, 25(1), 140–147.
- [36] Tijms, H. C. *Stochastic models: an algorithmic approach*; Wiley: Chichester, 1994.
- [37] van der Mei, R. D. Towards a unifying theory on branching-type polling models in heavy traffic. *Que. Syst.* **2007**, 57, 29–46.
- [38] van der Mei, R. D.; Winands, E. M. M. A note on polling models with renewal arrivals and nonzero switch-over times. *Oper. Res. Lett.* **2008**, 36, 500–505.
- [39] Whitt, W. *Stochastic-Process Limits: An Introduction to Stochastic-Process Limits and Their Application to Queues*. Springer Series in Operations Research and Financial Engineering; Springer: Berlin, 2002.
- [40] Winands, E. M. M. On polling systems with large setups. *Oper. Res. Lett.* **2007**, 35, 584–590.
- [41] Winands, E. M. M. Branching-type polling systems with large setups. *OR Spectrum* **2011**, 33(1), 77–97.