



## UvA-DARE (Digital Academic Repository)

### Scaling Limits for Infinite-server Systems in a Random Environment

Heemskerk, M.; van Leeuwen, J.; Mandjes, M.

**DOI**

[10.1214/16-SSY214](https://doi.org/10.1214/16-SSY214)

**Publication date**

2017

**Document Version**

Final published version

**Published in**

Stochastic Systems

**License**

CC BY

[Link to publication](#)

**Citation for published version (APA):**

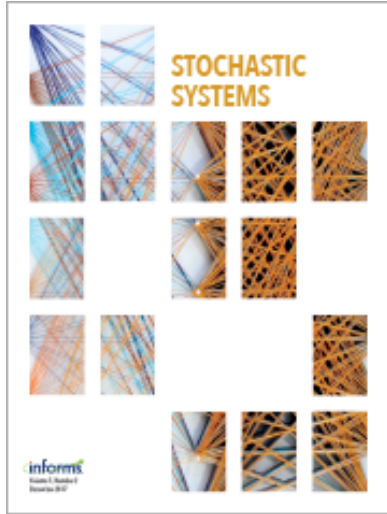
Heemskerk, M., van Leeuwen, J., & Mandjes, M. (2017). Scaling Limits for Infinite-server Systems in a Random Environment. *Stochastic Systems*, 7(1), 1-31.  
<https://doi.org/10.1214/16-SSY214>

**General rights**

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

**Disclaimer/Complaints regulations**

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.



## Stochastic Systems

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### Scaling Limits for Infinite-server Systems in a Random Environment

Mariska Heemskerk, Johan van Leeuwaarden, Michel Mandjes

To cite this article:

Mariska Heemskerk, Johan van Leeuwaarden, Michel Mandjes (2017) Scaling Limits for Infinite-server Systems in a Random Environment. *Stochastic Systems* 7(1):1-31. <https://doi.org/10.1287/16-SSY214>

Full terms and conditions of use: <https://pubsonline.informs.org/page/terms-and-conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2017, The author(s)

Please scroll down for article—it is on subsequent pages

INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

## SCALING LIMITS FOR INFINITE-SERVER SYSTEMS IN A RANDOM ENVIRONMENT

BY MARISKA HEEMSKERK\*, JOHAN VAN LEEUWAARDEN<sup>†</sup>  
AND MICHEL MANDJES\*

*University of Amsterdam\** and *Eindhoven University of Technology<sup>†</sup>*

This paper studies the effect of an overdispersed arrival process on the performance of an infinite-server system. In our setup, a random environment is modeled by drawing an arrival rate  $\Lambda$  from a given distribution every  $\Delta$  time units, yielding an i.i.d. sequence of arrival rates  $\Lambda_1, \Lambda_2, \dots$ . Applying a martingale central limit theorem, we obtain a functional central limit theorem for the scaled queue length process. We proceed to large deviations and derive the logarithmic asymptotics of the queue length's tail probabilities. As it turns out, in a rapidly changing environment (i.e.,  $\Delta$  is small relative to  $\Lambda$ ) the overdispersion of the arrival process hardly affects system behavior, whereas in a slowly changing random environment it is fundamentally different; this general finding applies to both the central limit and the large deviations regime. We extend our results to the setting where each arrival creates a job in multiple infinite-server queues.

**1. Introduction.** Empirical studies show that the number of arrivals in customer contact centers, hospital emergency departments and cloud computing systems typically varies strongly over time [8, 17]. This motivates modeling such arrival processes by a *non-homogeneous Poisson process* (NHPP) with time-dependent arrival rate  $\lambda(t)$ , see e.g. [9]. At the same time, various studies show that in a broad variety of real-life systems the intensity of the fluctuations in the arrival rate is so severe that the Poisson assumption ceases to hold [2, 8]. The observed level of *overdispersion* urges the need to develop stochastic models that can capture such persistent fluctuations.

Starting from the classical Poisson process, it is common practice to increase dispersion by using a *mixed* Poisson process [2, 13], to that end replacing a deterministic parameter  $\lambda$  by a random parameter  $\Lambda$ . This leads to the idea of modeling overdispersed arrival processes by a mixed version

---

Received January 2016.

*AMS 2000 subject classifications:* 60K25, 60F05, 60F10, 60F17, 60H20, 60K37, 97M40, 90B15

*Keywords and phrases:* Scaling limits, overdispersion, non-Poisson arrival processes, Cox processes, infinite-server queues, central limit theorem, large deviations.

of NHPPs, so-called *Cox processes* [5], where the time-dependent rate  $\lambda(t)$  of the classical NHPP is replaced by a stochastic process  $\Lambda(t)$ . For instance, one could use Markov-modulated Poisson processes (MMPPs) in which the arrival rate  $\Lambda(t) = \lambda_{J(t)}$  is a function of a continuous-time Markov chain  $J(\cdot)$  on a finite state space  $S$  and non-negative rates  $\lambda_i$  for  $i \in S$  (see e.g. [1, 3]). Although the MMPP is versatile and has various attractive properties, it has considerable drawbacks as well. First, while a substantial body of results for single-server queues with Markov modulation has been established, considerably less is known about their many-server and infinite-server counterparts; see e.g. an account of this issue for the infinite-server system in [4]. Second, due to the fact that the process  $J(\cdot)$  is not observed, estimating the parameters of an MMPP from data is a non-trivial task [15].

The main objective of this paper is to develop an arrival process simpler than a Markov-modulated Poisson process – arguably the simplest in terms of analysis – that fits the overdispersed and time-dependent setting, and to assess the impact of these characteristics on a corresponding system’s performance. The model we propose is a *mixed Poisson* arrival process in a *random environment*. It is defined as follows. Let  $\Lambda$  a non-negative random variable with finite first two moments and density  $f_\Lambda(\cdot)$ . Introduce a *sampling frequency*  $\frac{1}{\Delta}$ ; then the arrival rate at time  $t$  is given by  $\Lambda_j$  when  $t \in [j\Delta, (j+1)\Delta)$ , where the  $\Lambda_j$  are independent and distributed as a non-negative random variable  $\Lambda$ , for  $j \in \mathbb{Z}$ . In other words, this arrival process is a special case of a stationary Cox process where the arrival rate at time  $t$  is given by

$$(1.1) \quad \Lambda(t) = \sum_j \Lambda_j 1_{[j\Delta, (j+1)\Delta)}(t).$$

To add nonstationarity in the arrivals, one could include a deterministic component  $\bar{\lambda}(t)$  without intrinsically complicating the analysis; for ease of presentation we omit the extra component here. The resulting process can be viewed as an extension of the classical mixed Poisson setting, which is enriched by (independently) resampling the arrival rate after every time slot of length  $\Delta > 0$ . The intuition is that the arrival rate changes every  $\Delta$  time units so that, when observed for a number of consecutive slots, the time between arrivals and hence the number of arrivals per time unit fluctuates more severely than one would expect in a standard Poisson setting. This can be made explicit via an elementary computation. Let the number of arrivals up to time  $t$  be given by  $N_t \sim \text{Pois}(\int_0^t \Lambda_s ds)$  and for simplicity, let  $t$  be an integer multiple of  $\Delta$ . Then  $\mathbb{E}N_t = t\mathbb{E}\Lambda$ , whereas

$$\begin{aligned} \text{Var}(N_t) &= \sum_{j=1}^{t/\Delta} \text{Var}(N_\Delta) = t\Delta^{-1}(\mathbb{E}[\text{Var}(N_\Delta|\Lambda)] + \text{Var}(\mathbb{E}[N_\Delta|\Lambda])) \\ &= t(\mathbb{E}\Lambda + \Delta\text{Var}(\Lambda)). \end{aligned}$$

Conclude that, as desired, the variance-to-mean ratio is strictly larger than 1 for non-deterministic  $\Lambda$ , i.e.,

$$(1.2) \quad \frac{\text{Var}(N_t)}{\mathbb{E}N_t} = 1 + \Delta \frac{\text{Var}(\Lambda)}{\mathbb{E}\Lambda}.$$

Observe that the level of overdispersion is determined by the slot size  $\Delta$  and the level of overdispersion in  $\Lambda$  (through its variance-to-mean ratio).

Given this model for the arrival process, various queueing models can be studied; in this paper we focus on single-class infinite-server systems with exponential service times. The proposed arrival process being overdispersed, the main objective of this paper is to reveal, in a compact manner, the impact of overdispersion on system performance. Infinite-server systems are a natural choice when the system at hand is designed to (almost) immediately serve all customers [16], but it may also serve as a tractable proxy for the more complicated multi-server systems, which is for instance exploited in the modified offered-load (MOL) and pointwise stationary approximation (PSA) methods for staffing large-scale service systems in a time-varying setting [10, 17].

*Contributions.* Infinite-server systems with overdispersed arrivals are, as described above, very tractable. As shown in Section 2, it is fairly straightforward to compute the probability generating function (PGF) of the stationary and time-dependent queue length processes (which in an infinite-server setting refers to the number of jobs in the system) in terms of transforms. This is due to the fact that customers are served immediately upon arrival, independently of each other; as a result, when analyzing the queue length at a given point in time, we can separately consider the individual (independent!) contributions that correspond to each of the preceding intervals of length  $\Delta$ .

The queue length distribution can be characterized in terms of its PGF, which effectively means that evaluation of the accompanying performance measures requires numerical inversion. However, by imposing a scaling on both the time and scale parameters,  $\Delta$  and  $\Lambda$ , we succeed in identifying an asymptotic regime in which the distribution *can* be explicitly given. We inflate the arrival rate and sampling frequency in the following way:

$$(1.3) \quad \Lambda \mapsto N\Lambda \quad \Delta^{-1} \mapsto N^\alpha \Delta^{-1},$$

where we let  $N \rightarrow \infty$ . The scaled counterpart of the variance-to-mean ratio in (1.2) is  $1 + N^{1-\alpha} \Delta \text{Var}(\Lambda) / \mathbb{E}\Lambda$ . Due to the possibility of having different

growth rates for  $\Lambda$  and  $\Delta^{-1}$  under scaling (1.3) this ratio will not always converge to 1. That is, the value of  $\alpha$  determines the nature of the asymptotic behavior of the arrival process, giving rise to a *trichotomy*. In this paper we prove that the queue length process inherits this behavior from the input process. For  $\alpha > 1$ , in which case the arrival rate is resampled relatively frequently, we find that the system behaves as a standard infinite-server queue (no overdispersion), whereas for  $\alpha < 1$  the overdispersion remains present in the asymptotic regime. The case  $\alpha = 1$  essentially reflects a superposition of the two distinct types of behavior.

For preparatory purposes, we show in Section 2 that the centered and normalized stationary queue length is asymptotically normal under the scaling in (1.3). Next, in Section 3 we consider a multidimensional setting with correlated arrivals: an arrival triggers jobs in multiple queues. Hence, we work with a coupled system in which  $d$  parallel queues are fed by a single arrival process; cf. [11, 12]. With  $\mathbf{U}^{(N)}(\cdot)$  denoting the vector of centered and normalized queue length processes, the asymptotic normality now translates to the corresponding limiting process  $\mathbf{U}(\cdot)$  being Gaussian:  $\mathbf{U}(\cdot)$  is a  $d$ -dimensional process of the Ornstein-Uhlenbeck type with parameters that depend on the scaling regime. Following the approach in [1], we show this by applying a lemma due to Kurtz and a martingale central limit theorem (MCLT) to a suitable stochastic integral equation.

Subsequently, in Section 4 we carry out a large deviations analysis to obtain the logarithmic tail asymptotics corresponding to the queue length distribution. The crucial observation in this analysis is that rare events can essentially be realized in two ways: (i) the random arrival rate attains an exceptionally high value, (ii) the Poisson process generates an unusually large number of arrivals given the (not so rare) value of the random parameter. Again, the value of  $\alpha$  determines what type of tail behavior dominates: for  $\alpha < 1$  this is effect (i), for  $\alpha > 1$  effect (ii), and for  $\alpha = 1$  a combination of effects (i) and (ii). These findings complement similar results that have been established for an infinite-server system with Markov-modulated input, where it is noted that the slow regime ( $\alpha \in (0, 1)$ ) was not covered in that setting [3, 6]. We conclude Section 4 by pointing out how the large deviations results can be extended to the multidimensional setting.

**2. Overdispersion in an infinite-server context.** In this section we present a stationary and transient analysis of the single-class Markovian infinite-server system in a random environment just introduced. A crucial role is played by  $\Lambda(t)$ , the arrival rate at time  $t$  given in (1.1). Remember that we assumed that the arrival rates are i.i.d. and distributed as a random variable  $\Lambda \geq 0$  with finite first two moments and density  $f_\Lambda(\cdot)$ . The corre-

sponding service times are assumed i.i.d. (and in addition independent of the arrival process) exponentially distributed random variables with mean  $1/\mu$ .

First, in Section 2.1, we analyze the stationary system behavior, in terms of its PGF and the corresponding moments, which we then extend to the associated transient behavior. We then study the stationary behavior in a central limit regime under parameter scaling (1.3) in Section 2.2. This exposition serves as an illustration for the reader, and is intended to create intuition as for why the scaled stationary queue length is asymptotically normal and why the three different limiting regimes appear; in addition, in Section 4 we need a result that is proven along the same lines. We note that in Section 3 the normality is generalized in several directions: we establish a *functional* central limit theorem (FCLT) for the (scaled) transient process  $M^{(N)}(\cdot)$  corresponding to the  $d$ -dimensional parallel system as defined in the introduction.

*2.1. Pre-limit results.* This subsection presents ‘pre-limit results’; later we study their counterparts in the limiting regime after imposing a parameter scaling.

*Transform of stationary queue length.* Let  $M$  be the random variable associated with the stationary number of jobs (also sometimes referred to as ‘customers’) in the system. Exploiting ‘thinning’ properties, we can identify the PGF  $\phi(z) := \mathbb{E}z^M$  of  $M$ .

In the sequel we write  $p_t := e^{-\mu t}$  for the probability that a job present at  $kt$  is still present at  $(k+1)t$  and  $q_t := (1 - e^{-\mu t})/(\mu t)$  for the probability that a job arriving at a uniform epoch in  $[kt, (k+1)t)$  is still present at  $(k+1)t$ . Denote  $\bar{p}_t := 1 - p_t$ .

Note that  $M$  can be written as the sum of  $M_0, M_1, M_2, \dots$ , where  $M_k$  represents the number of jobs that arrived in  $[-(k+1)\Delta, -k\Delta)$  and are still present at time 0. Furthermore, observe that these ‘thinned’ random variables  $M_k$  are independent. A job that arbitrarily arrived in  $[-(k+1)\Delta, -k\Delta)$  (i.e., having arrived at a uniform epoch in this interval) is still in the system at time 0 with probability

$$\int_0^\Delta \frac{1}{\Delta} e^{-\mu(k\Delta+s)} ds = q_\Delta p_\Delta^k.$$

As a consequence, with  $r_t := tq_t$ ,

$$\begin{aligned} \phi_k(z) &:= \mathbb{E}z^{M_k} \\ &= \sum_{\ell=0}^{\infty} \int_0^{\infty} f_\Lambda(\lambda) e^{-\lambda\Delta} \frac{(\lambda\Delta)^\ell}{\ell!} \sum_{m=0}^{\ell} z^m \binom{\ell}{m} (q_\Delta p_\Delta^k)^m (1 - q_\Delta p_\Delta^k)^{\ell-m} d\lambda \end{aligned}$$

$$\begin{aligned}
&= \int_0^\infty \exp(-\lambda r_\Delta p_\Delta^k (1-z)) f_\Lambda(\lambda) d\lambda \\
(2.1) \quad &= \mathbb{E} \exp(-\Lambda_k r_\Delta p_\Delta^k (1-z)).
\end{aligned}$$

Observe that  $\phi_k(z)$  is a PGF of ‘mixed Poisson’ type: conditional on  $\Lambda_k = \lambda$  the PGF corresponds with that of a Poisson random variable with mean  $\lambda r_\Delta p_\Delta^k$ . We conclude that  $M_k$  is distributed as a mixed Poisson random variable with random parameter

$$\kappa_k(\Lambda_k) := \Lambda_k r_\Delta p_\Delta^k,$$

with  $\Lambda_k$  the value of the arrival rate in the interval  $[-(k+1)\Delta, -k\Delta)$  (note that, in fact, we should have written  $\Lambda_{-(k+1)}$  rather than  $\Lambda_k$ , but due to the i.i.d. assumption the processes  $\{\Lambda(s)\}_{s \geq 0}$  and  $\{\Lambda(-s)\}_{s \geq 0}$  have the same finite-dimensional distributions). Therefore,  $M$  is mixed Poisson as well and its random parameter is given by

$$(2.2) \quad \sum_{k=0}^{\infty} \kappa_k(\Lambda_k) = \int_0^\infty \Lambda(s) e^{-\mu s} ds =: \kappa(\Lambda).$$

(Note that  $\kappa(\cdot)$  is defined as a functional;  $\kappa(\Lambda)$  should be interpreted as  $\kappa(\Lambda(\cdot))$ .)

There is an alternative way to obtain this result. Indeed, since we observe the system in stationarity,

$$(2.3) \quad \phi(z) = \phi(z p_\Delta + \bar{p}_\Delta) g_{\Lambda, \Delta}(z),$$

where  $g_{\Lambda, t}(z)$  is defined by

$$g_{\Lambda, t}(z) := \int_0^\infty \exp(-\lambda r_t (1-z)) f_\Lambda(\lambda) d\lambda = \mathbb{E} \exp(-\Lambda r_t (1-z)).$$

Applying an iteration argument to (2.3) yields

$$(2.4) \quad \phi(z) = \prod_{k=0}^{\infty} g_{\Lambda, \Delta}(z p_\Delta^k + \bar{p}_\Delta \sum_{j=0}^{k-1} p_\Delta^j) = \prod_{k=0}^{\infty} g_{\Lambda, \Delta}(1 - (1-z) p_\Delta^k).$$

In the factors  $g_{\Lambda, \Delta}(1 - (1-z) p_\Delta^k)$  we recognize the expression for  $\phi_k(z)$  as in (2.1).

*First two moments.* We now evaluate the first two moments of  $M$ . This is an interesting computation in its own right, but it also provides useful results that can be exploited when considering this system under the central limit scaling (as is done in the next subsection).



Differentiating (2.3) and letting  $z \uparrow 1$ , we obtain a fixed-point equation,

$$\phi'(1) = \phi'(1)e^{-\mu\Delta} + g'_{\Lambda,\Delta}(1) = \phi'(1)e^{-\mu\Delta} + r_{\Delta} \mathbb{E}\Lambda.$$

Hence  $\mathbb{E}M = \phi'(1) = r_{\Delta} \mathbb{E}\Lambda / (1 - e^{-\mu\Delta}) = \mathbb{E}\Lambda / \mu$ . This quantity could have been computed more directly as well, using a standard identity for conditional means:

$$(2.5) \quad \mathbb{E}M = \sum_{k=0}^{\infty} \mathbb{E}M_k = \sum_{k=0}^{\infty} \mathbb{E}[\mathbb{E}[M_k | \Lambda_k]].$$

Then observe that  $(M_k | \Lambda_k)$  is Poisson, and hence its mean equals its parameter. As a result, (2.5) equals

$$\mathbb{E}M = \sum_{k=0}^{\infty} \mathbb{E}[\kappa_k(\Lambda_k)] = \mathbb{E}[\Lambda] r_{\Delta} \sum_{k=0}^{\infty} p_{\Delta}^k = \mathbb{E}\Lambda / \mu.$$

For the variance we use that

$$\phi''(1) = \phi''(1)p_{\Delta}^2 + 2\phi'(1)p_{\Delta}g'_{\Lambda,\Delta}(1) + g''_{\Lambda,\Delta}(1),$$

and hence

$$\phi''(1) = \frac{2\phi'(1)p_{\Delta}g'_{\Lambda,\Delta}(1)}{1-p_{\Delta}^2} + \frac{g''_{\Lambda,\Delta}(1)}{1-p_{\Delta}^2} = 2\frac{(\mathbb{E}\Lambda)^2}{\mu^2} \frac{p_{\Delta}}{1+p_{\Delta}} + \frac{\mathbb{E}\Lambda^2}{\mu^2} \frac{1-p_{\Delta}}{1+p_{\Delta}}.$$

It thus follows that, after some algebra,

$$(2.6) \quad \begin{aligned} \text{Var}M &= \phi''(1) + \phi'(1) - (\phi'(1))^2 \\ &= \mathbb{E}\Lambda / \mu + C \text{Var}\Lambda / \mu^2, \end{aligned}$$

where  $C := (1 - p_{\Delta}) / (1 + p_{\Delta})$ .

Alternatively, one could use the ‘law of total variance’ to identify  $\text{Var}M$ :

$$(2.7) \quad \text{Var}M = \sum_{k=0}^{\infty} \text{Var}M_k = \sum_{k=0}^{\infty} \mathbb{E}[\text{Var}(M_k | \Lambda_k)] + \sum_{k=0}^{\infty} \text{Var}(\mathbb{E}[M_k | \Lambda_k]).$$

Observe that, because of the ‘mixed Poisson property’,  $\mathbb{E}[\text{Var}(M_k | \Lambda_k)] = \mathbb{E}[\kappa_k(\Lambda_k)]$ , and as a result the first term at the right-hand side of (2.7) equals  $\mathbb{E}M$ . The second term, which is inherently non-negative, gives rise to ‘overdispersion’, i.e., the effect that the variance of the stationary queue length *exceeds* the corresponding mean. This is a distinguishing feature compared to the analogous system in which the Poissonian arrival rate is deterministic: the stationary queue length in an M/M/ $\infty$  system is Poisson, and

cannot accommodate any overdispersion. In order to evaluate the second term in the right-hand side of (2.7), we note that

$$(2.8) \quad \text{Var}(\mathbb{E}(M_k | \Lambda_k)) = \text{Var}(\kappa_k(\Lambda_k)) = r_{\Delta}^2 p_{\Delta}^{2k} \cdot \text{Var}\Lambda.$$

Substituting (2.8) in the second term in the right-hand side of (2.7), we find that  $\text{Var} M$  equals (2.6), as desired.

Formula (2.6) lends itself to a nice interpretation: the term  $\mathbb{E}\Lambda/\mu$  is the contribution to the variance that one would have if the arrival rate would have had the *deterministic* value  $\mathbb{E}\Lambda$ , whereas the term  $C \text{Var}\Lambda/\mu^2$  needs to be added in order to deal with the non-Poisson variability due to the stochasticity of the arrival rate.

*Transient behavior.* As the analysis of the transient system behavior strongly resembles its stationary counterpart, we restrict ourselves to a short account of this. We let the system start empty (for ease of presentation; a non-empty initial condition can be analyzed without any additional difficulty). Denote by  $M(t)$  the number of jobs present at time  $t$ . Then, for  $n$  the smallest integer such that  $t - n\Delta < \Delta$ ,

$$M(t) = \sum_{j=0}^{n-1} \bar{M}_j + \bar{M}_{[n\Delta, t]},$$

where  $\bar{M}_j$  ( $\bar{M}_{[n\Delta, t]}$ ) represents the number of jobs that have arrived between in  $[j\Delta, (j+1)\Delta)$  ( $[n\Delta, t)$ ) and are still around at  $n\Delta$  ( $t$ ). As before, these have PGFs

$$\begin{aligned} \mathbb{E}_z \bar{M}_j &= \mathbb{E} \exp(-\Lambda r_{\Delta} p_{\Delta}^{n-(j+1)} e^{-\mu(t-n\Delta)} (1-z)); \\ \mathbb{E}_z \bar{M}_{[n\Delta, t]} &= \mathbb{E} \exp(-\Lambda/\mu (1 - e^{-\mu(t-n\Delta)}) (1-z)). \end{aligned}$$

As the individual random variables  $\bar{M}_1, \bar{M}_2, \dots$  and  $\bar{M}_{[n\Delta, t]}$  are independent,  $M(t)$  is mixed Poisson with random parameter

$$(2.9) \quad \kappa_t(\Lambda) := \int_0^t \Lambda(s) e^{-\mu s} ds.$$

**2.2. Limit results.** This section focuses on the central limit regime that results from simultaneously scaling, in a controlled way, both the arrival rate  $\Lambda$  and the sampling frequency  $\frac{1}{\Delta}$  as in (1.3). Let the scaled counterpart of  $\Lambda(t)$  be  $N\Lambda^{(N)}(t)$ , with

$$(2.10) \quad \Lambda^{(N)}(t) := \sum_{j=0}^{\infty} \Lambda_j 1_{[j\Delta N^{-\alpha}, (j+1)\Delta N^{-\alpha})}(t).$$

That is, the sampling frequency and the arrival rates are both inflated as we let  $N$  tend to  $\infty$ , but, importantly, at rates that are not necessarily identical. As mentioned in the introduction, depending on the value of  $\alpha$ , we obtain fundamentally different behavior.

We consider a sequence of systems indexed by  $N$ , where the  $N$ -scaled system uses a mixed Poisson arrival process with time-dependent random rate  $N\Lambda^{(N)}(t)$ . Let  $M^{(N)}$  denote the stationary queue length in the  $N$ -scaled system, with parameter  $N\kappa(\Lambda^{(N)})$  (cf. (2.2)). We start our exposition by a preliminary calculation, in which we compute the mean and variance of  $M^{(N)}$  and study their behavior for large  $N$ , which indeed reveals the announced trichotomy. Then, after centering and normalizing  $M^{(N)}$ , we derive a central limit theorem.

*Qualitative behavior of first two moments: trichotomy in variance.* First, we identify the steady-state mean and variance in our scaling regime, using (2.5), (2.7) and (2.8). We find that

$$(2.11) \quad \mathbb{E}M^{(N)} = N\mathbb{E}\Lambda/\mu;$$

$$(2.12) \quad \mathbb{V}\text{ar}M^{(N)} = N\mathbb{E}\Lambda/\mu + N^2 \frac{1 - e^{-\mu\Delta N^{-\alpha}}}{1 + e^{-\mu\Delta N^{-\alpha}}} \mathbb{V}\text{ar}\Lambda/\mu^2,$$

where it is noted that for large  $N$ , (2.12) behaves approximately as

$$N\mathbb{E}\Lambda/\mu + N^{2-\alpha}\Delta\mathbb{V}\text{ar}\Lambda/(2\mu)$$

(the ratio of the two converges to 1). We thus observe the trichotomy

$$(2.13) \quad \mathbb{V}\text{ar}M^{(N)} \sim \begin{cases} N\mathbb{E}\Lambda/\mu & \text{if } \alpha > 1; \\ N^{2-\alpha}\Delta\mathbb{V}\text{ar}\Lambda/(2\mu) & \text{if } \alpha < 1; \\ N(\mathbb{E}\Lambda/\mu + \Delta\mathbb{V}\text{ar}\Lambda/(2\mu)) & \text{if } \alpha = 1. \end{cases}$$

For  $\alpha > 1$ , the sampling frequency dominates the variability of  $\Lambda$ . Consequently, the model behaves essentially as an M/M/ $\infty$  system, with the variance of  $M^{(N)}$  being linear in  $N$  and equal to  $\mathbb{E}M^{(N)}$ , for large  $N$ . For  $\alpha < 1$ , we find a superlinear relation between  $N$  and  $\mathbb{V}\text{ar}\Lambda$ , and both the sampling frequency (i.e., the reciprocal of the interval length  $\Delta$ ) and the variance of  $\Lambda$  play a role. Hence, the asymptotic variance indeed grows faster than the asymptotic mean for  $\alpha < 1$ ; in this regime the system is overdispersed. For  $\alpha = 1$ , the variance is ‘slightly larger’ than for  $\alpha > 1$ , but it is still linear in  $N$ . In this case the sampling frequency and the variance of  $\Lambda$  grow at the same rate, so that the variance for  $M^{(N)}$  combines the effects observed in the two former cases.

As observed from the above computation, the variance of  $M^{(N)}$  is essentially proportional to  $N^\gamma$  with  $\gamma := \max\{1, 2 - \alpha\}$ . As a consequence, one may expect that, under (1.3),

$$(2.14) \quad \check{M}^{(N)} := N^{-\gamma/2}(M^{(N)} - \mathbb{E}M^{(N)})$$

converges to a (zero-mean) normally distributed random variable. It is this property that we verify now.

*Asymptotic normality.* We show how to establish asymptotic normality for the centered and normalized version of  $M^{(N)}$  in (2.14) via evaluation of the corresponding Laplace transform. Appealing to Lévy's convergence theorem, we establish the desired convergence in distribution. For simplicity, the proof of Thm. 2.1 assumes that all moments of  $\Lambda$  are finite; however, as will appear from the proof of Thm. 3.2 only finiteness of the first two moments is necessary.

**THEOREM 2.1 (CLT).** *As  $N \rightarrow \infty$ ,  $\check{M}^{(N)}$  converges to a zero-mean normally distributed random variable with variance*

$$\sigma^2 := \frac{\mathbb{E}\Lambda}{\mu} 1_{\{\alpha \geq 1\}} + \frac{\Delta \text{Var } \Lambda}{2\mu} 1_{\{\alpha \leq 1\}}.$$

**PROOF.** Let  $\phi^{(N)}(z)$  be the counterpart of (2.4) under scaling as in (1.3); likewise  $g_{\Lambda, \Delta}^{(N)}(z)$  is the counterpart of  $g_{\Lambda, \Delta}(z)$ . Then

$$\phi^{(N)}(z) = \prod_{k=0}^{\infty} g_{\Lambda, \Delta}^{(N)}(1 - (1 - z)e^{-\mu k N^{-\alpha} \Delta}).$$

We are interested in the behavior of  $M^{(N)}$  in the central limit regime, hence we need to analyze the limiting distribution of  $\check{M}^{(N)}$ . To this end, we evaluate the logarithm of the corresponding Laplace transform:

$$(2.15) \quad \log \mathbb{E} \exp(-sN^{-\gamma/2}(M^{(N)} - \mathbb{E}M^{(N)})) = sN^{1-\gamma/2} \mathbb{E}\Lambda/\mu + \log \phi^{(N)}(e^{-sN^{-\gamma/2}}).$$

We now use that  $\log \phi^{(N)}(e^{-sN^{-\gamma/2}})$  equals

$$(2.16) \quad \sum_{k=0}^{\infty} \log \mathbb{E} e^{-N\Lambda/\mu(1 - e^{-\mu N^{-\alpha} \Delta})(1 - e^{-sN^{-\gamma/2}})e^{-\mu k N^{-\alpha} \Delta}}.$$

Observe that (2.16) is the sum of cumulant generating functions (which is again a cumulant generating function), each of them related to the random

variable  $\Lambda$  but evaluated at different arguments. Let  $m_\ell$  denote the  $\ell$ -th cumulant of  $\Lambda$  (for  $\ell \in \mathbb{N}$ ); in particular  $m_1 = \mathbb{E}\Lambda$  and  $m_2 = \mathbb{V}\text{ar}\Lambda$ . In addition, we define

$$\zeta_k^{(N)}(s) := -N/\mu(1 - e^{-\mu N^{-\alpha}\Delta})(1 - e^{-sN^{-\gamma/2}})e^{-\mu k N^{-\alpha}\Delta}.$$

Then it follows that

$$(2.17) \quad \log \phi^{(N)}(e^{-sN^{-\gamma/2}}) = \sum_{\ell=1}^{\infty} \frac{m_\ell}{\ell!} \sum_{k=0}^{\infty} (\zeta_k^{(N)}(s))^\ell.$$

Let us first consider the contribution of the term corresponding to  $\ell = 1$ . Observe that, as  $N \rightarrow \infty$ ,

$$(2.18) \quad \begin{aligned} m_1 \sum_{k=0}^{\infty} \zeta_k^{(N)}(s) &= -N \mathbb{E}\Lambda/\mu(1 - e^{-sN^{-\gamma/2}}) \\ &\sim -sN^{1-\gamma/2} \mathbb{E}\Lambda/\mu + \frac{1}{2}s^2 N^{1-\gamma} \mathbb{E}\Lambda/\mu. \end{aligned}$$

Note that the first term in (2.18) is canceled by the first term in the right-hand side of (2.15), so that we are left with the second term, i.e.,

$$(2.19) \quad \frac{1}{2}s^2 N^{1-\gamma} \mathbb{E}\Lambda/\mu.$$

The second term in (2.17) corresponding to  $\ell = 2$  gives

$$(2.20) \quad \frac{1 - e^{-\mu N^{-\alpha}\Delta}}{1 + e^{-\mu N^{-\alpha}\Delta}} \frac{(1 - e^{-sN^{-\gamma/2}})^2}{2\mu^2} N^2 \mathbb{V}\text{ar}\Lambda \sim \frac{1}{2}s^2 N^{2-\alpha-\gamma} \Delta \mathbb{V}\text{ar}\Lambda/(2\mu).$$

Now compare the asymptotic expansion identified in (2.19) and (2.20). In case  $\alpha > 1$ , we have that  $\gamma = 1$ , so that (2.19) equals  $\frac{1}{2}s^2 \mathbb{E}\Lambda/\mu$ , whereas (2.20) converges to zero. On the other hand, for  $\alpha < 1$  we have  $\gamma = 2 - \alpha$ , and hence (2.19) converges to zero, whereas (2.20) behaves as  $\frac{1}{2}s^2 \Delta \mathbb{V}\text{ar}\Lambda/(2\mu)$ . Finally, if  $\alpha = 1$ , we find that both terms converge to the expected finite positive limit.

We now check that the terms in (2.17) for  $\ell \geq 3$  vanish as  $N \rightarrow \infty$ . For large  $N$  the terms can be approximated as follows,

$$\sum_{k=0}^{\infty} (\zeta_k^{(N)}(s))^\ell \sim N^\ell \frac{(1 - e^{-\mu N^{-\alpha}\Delta})^\ell}{1 - e^{-\ell\mu N^{-\alpha}\Delta}} (1 - e^{-sN^{-\gamma/2}})^\ell \sim N^\ell \frac{\mu^\ell N^{-\alpha\ell} \Delta^\ell}{\ell\mu N^{-\alpha}\Delta} \frac{s^\ell}{N^{\gamma\ell/2}},$$

hence being of order  $N^\delta$  with  $\delta = \delta(\alpha) := \ell(1 - \alpha - \gamma/2) + \alpha$ . In case  $\alpha \geq 1$ ,

we get (bearing in mind that  $\gamma = 1$  and  $\ell \geq 1$ )

$$\delta = \ell\left(\frac{1}{2} - \alpha\right) + \alpha = \frac{1}{2}\ell + \alpha(1 - \ell) \leq \frac{1}{2}\ell + 1 - \ell = 1 - \frac{\ell}{2};$$

on the other hand, in case  $\alpha < 1$  we get  $\delta = -\ell\alpha/2 + \alpha = \alpha(1 - \ell/2)$  (with  $\gamma = 2 - \alpha$ ). We conclude that  $\delta < 0$  for  $\ell \geq 3$  and the corresponding terms in (2.17) can indeed be neglected. We have therefore established that, as  $N \rightarrow \infty$ ,

$$\log \mathbb{E} \exp\left(-sN^{-\gamma/2}(M^{(N)} - \mathbb{E}M^{(N)})\right) \rightarrow \frac{1}{2}\sigma^2 s^2,$$

as claimed.  $\square$

**3. Functional central limit theorem.** In this section we generalize the central limit result of Thm. 2.1 in two ways. First, we establish the functional version: the centered and normalized transient queue length process converges to a limiting process of Ornstein-Uhlenbeck type with parameters that depend on the value of  $\alpha$ . Second, we extend this to a multidimensional setting with correlated arrivals: every arrival triggers jobs in multiple queues. The correlation structure of the resulting multidimensional Gaussian limiting process is explicitly identified.

Let us start by describing the mechanics of the generalized setting. We consider a parallel system in which  $d$  queues are fed by a *single* arrival process that was constructed in the same way as the one in the previous section: a Markovian process with arrival rate  $\Lambda(t)$  as in (1.1). The service times in queue  $i$  are i.i.d. exponential random variables with mean  $\mu_i^{-1}$ ; the service processes of the individual queues are independent, and also independent of the arrival process. We perform the same scaling as before: the sampling frequency is sped up by a factor  $N^\alpha$ , while the (random) arrival rate is blown up by a factor  $N$ . This results in a mixed Poisson arrival process with time-dependent rate  $\Lambda^{(N)}(t)$  as in (2.10). Let

$$\mathbf{M}^{(N)}(t) = (M_1^{(N)}(t), \dots, M_d^{(N)}(t))^T,$$

where  $M_i^{(N)}(t)$  is the queue length at time  $t$  in the  $i$ -th queue of the  $N$ -scaled system, for  $i \in \{1, \dots, d\}$ . Note that the  $M_i^{(N)}(t)$  are mixed Poisson with time-dependent random parameter  $N\kappa_{t,i}(\Lambda^{(N)})$ , with  $\kappa_{t,i}(\Lambda^{(N)})$  as defined in (2.9) but with  $\mu$  replaced by  $\mu_i$ .

We now present an alternative way of writing  $M_i^{(N)}(t)$ , which facilitates the use of a martingale central limit theorem. We introduce the functional

$$\Psi[X](t) := \int_0^t X(s) \, ds,$$

mapping the stochastic process  $\{X(s) : s \in [0, t]\}$  to a real number; then  $\mu_i \Psi[M_i^{(N)}](t)$  is to be interpreted as the ‘cumulative service capacity’ in queue  $i$  over the interval  $[0, t]$ . In addition, for the ‘cumulative arrival rate’ we have  $\Psi[\Lambda](t)$ , with scaled counterpart  $N\Psi[\Lambda^{(N)}](t)$ . By the law of large numbers,  $\Psi[\Lambda](t)/t$  converges a.s. to  $\mathbb{E}\Lambda$  as  $t \rightarrow \infty$  and for fixed  $t$ ,  $\Psi[\Lambda^{(N)}](t)$  converges a.s. to  $t\mathbb{E}\Lambda$  as  $N \rightarrow \infty$ .

With  $Y_0(\cdot), \dots, Y_d(\cdot)$  denoting independent unit-rate Poisson processes,

$$(3.1) \quad M_i^{(N)}(t) \stackrel{d}{=} M_i^{(N)}(0) + Y_0(N\Psi[\Lambda^{(N)}](t)) - Y_i(\mu_i \Psi[M_i^{(N)}](t)).$$

Our objective is to derive a  $d$ -dimensional FCLT for  $\mathbf{M}^{(N)}(\cdot)$ . This result characterizes the time-dependent queue length in the scaled system and makes explicit how the correlated arrivals lead to correlation between the individual queue length processes. It will be stated and proved in subsection 3.2; first we study the stationary behavior by presenting the corresponding first two moments of  $\mathbf{M}^{(N)}$  (including the covariances).

3.1. *Qualitative behavior of first two moments in stationarity.* Note that the individual queue lengths are only coupled through the arrival process, so under (1.3) the mean and variance of  $\mathbf{M}^{(N)}$  are, as in (2.11) and (2.12), given by

$$\begin{aligned} \mathbb{E}M_i^{(N)} &= N\mathbb{E}\Lambda/\mu_i, \\ \text{Var}M_i^{(N)} &= N\mathbb{E}\Lambda/\mu_i + N^2\text{Var}\Lambda/\mu_i^2 \frac{1 - p_i(\Delta N^{-\alpha})}{1 + p_i(\Delta N^{-\alpha})} \\ &\sim N\mathbb{E}\Lambda/\mu_i + N^{2-\alpha}\Delta\text{Var}\Lambda/(2\mu_i), \end{aligned}$$

for  $i = 1, \dots, d$ . Hence, we find the same behavior as in (2.13). Interestingly, the same trichotomy is observed for the covariances between the individual queues, as stated in the next lemma.

LEMMA 3.1 (Covariance in  $\mathbf{M}^{(N)}$ ). *For  $i, k \in \{1, \dots, d\}$  with  $i \neq k$ , and for large  $N$ ,*

$$(3.2) \quad \text{Cov}(M_i^{(N)}, M_k^{(N)}) \sim \begin{cases} N\mathbb{E}\Lambda/(\mu_i + \mu_k) & \text{if } \alpha > 1; \\ N^{2-\alpha}\Delta\text{Var}(\Lambda)/(\mu_i + \mu_k) & \text{if } \alpha < 1; \\ N(\mathbb{E}\Lambda/(\mu_i + \mu_k) + \Delta\text{Var}(\Lambda)/(\mu_i + \mu_k)) & \text{if } \alpha = 1. \end{cases}$$

PROOF. Without loss of generality, we take  $i = 1$  and  $k = 2$ . We first consider the non-scaled model, by studying the joint probability generating

function,

$$\mathbb{E}_{z_1}^{M_1(n\Delta)} z_2^{M_2(n\Delta)} = \prod_{j=0}^{n-1} \xi_{jn}(z_1, z_2),$$

where  $\xi_{jn}(z_1, z_2)$  is the contribution due to the slot between  $j\Delta$  and  $(j+1)\Delta$ ; as  $z_1$  and  $z_2$  are held fixed for the moment, we suppress them. Now we introduce functions (for  $\ell = 1, 2$ )

$$f_\ell(r, n) := e^{-\mu_\ell(n\Delta - r)}, \quad g_j(\mu, n) := \frac{1}{\mu\Delta}(1 - e^{-\mu\Delta})e^{-\mu(n-j)\Delta},$$

where it is noted that  $g_j(\mu, n)$  behaves as  $e^{-\mu(n-j)\Delta}$  for small  $\Delta$ . In addition, we define the quantities

$$\begin{aligned} \zeta_{jn}^{++} &:= \int_{j\Delta}^{(j+1)\Delta} \frac{1}{\Delta} f_1(r, n) f_2(r, n) dr = g_j(\mu_1 + \mu_2, n), \\ \zeta_{jn}^{+-} &:= \int_{j\Delta}^{(j+1)\Delta} \frac{1}{\Delta} f_1(r, n) (1 - f_2(r, n)) dr = g_j(\mu_1, n) - g_j(\mu_1 + \mu_2, n), \\ \zeta_{jn}^{-+} &:= \int_{j\Delta}^{(j+1)\Delta} \frac{1}{\Delta} (1 - f_1(r, n)) f_2(r, n) dr = g_j(\mu_2, n) - g_j(\mu_1 + \mu_2, n), \\ \zeta_{jn}^{--} &:= \int_{j\Delta}^{(j+1)\Delta} \frac{1}{\Delta} (1 - f_1(r, n)) (1 - f_2(r, n)) dr \\ &= 1 - g_j(\mu_1, n) - g_j(\mu_2, n) + g_j(\mu_1 + \mu_2, n). \end{aligned}$$

Using arguments similar to those we have used before,

$$\begin{aligned} \xi_{jn} &:= \mathbb{E} \left( \sum_{m=0}^{\infty} e^{-\lambda\Delta} \frac{(\lambda\Delta)^m}{m!} (\zeta_{jn}^{++} z_1 z_2 + \zeta_{jn}^{+-} z_1 + \zeta_{jn}^{-+} z_2 + \zeta_{jn}^{--})^m \right) \\ &= \mathbb{E} \exp \left( \Lambda\Delta (\zeta_{jn}^{++} z_1 z_2 + \zeta_{jn}^{+-} z_1 + \zeta_{jn}^{-+} z_2 + \zeta_{jn}^{--} - 1) \right) \\ &\sim \mathbb{E} \exp \left( \Lambda\Delta \left( \prod_{i=1}^2 ((z_i - 1)e^{-\mu_i(n-j)\Delta} + 1) - 1 \right) \right) \text{ for small } \Delta. \end{aligned}$$

Because the contributions to  $M_1(n\Delta)$  and  $M_2(n\Delta)$  resulting from different time intervals are independent, we obtain that

$$\begin{aligned} &\text{Cov}(M_1(n\Delta), M_2(n\Delta)) \\ &= \sum_{j=0}^{n-1} \left( \frac{\partial^2}{\partial z_1 \partial z_2} \xi_{jn}(z_1, z_2) - \frac{\partial}{\partial z_1} \xi_{jn}(z_1, z_2) \frac{\partial}{\partial z_2} \xi_{jn}(z_1, z_2) \right) \Big|_{z_1 \uparrow 1, z_2 \uparrow 1}. \end{aligned}$$



Now imposing scaling (1.3) and considering the stationary behavior by letting  $n \rightarrow \infty$ , it is readily derived that (for large  $N$ )

$$\mathbb{E} z_1^{M_1^{(N)}} z_2^{M_2^{(N)}} \sim \prod_{j=0}^{\infty} \mathbb{E} \exp \left( \Lambda \Delta N^{1-\alpha} \left( \prod_{\ell=1}^2 ((z_\ell - 1) e^{-\mu_\ell j \Delta / N^\alpha} + 1) - 1 \right) \right);$$

observe that, for reasons of symmetry, it is allowed to replace  $n - j$  by  $j$  in the definition of the  $g_j(\mu, n)$ . We thus arrive at

$$\begin{aligned} \text{Cov}(M_1^{(N)}, M_2^{(N)}) &\sim \sum_{j=0}^{\infty} \left( (\mathbb{E} \Lambda \Delta N^{1-\alpha} + \mathbb{E}[\Lambda^2] \Delta^2 N^{2-2\alpha}) e^{-(\mu_1 + \mu_2) j \Delta / N^\alpha} \right. \\ &\quad \left. - \prod_{\ell=1}^2 \mathbb{E} \Lambda \Delta N^{1-\alpha} e^{-\mu_\ell j \Delta / N^\alpha} \right) \\ &= \frac{\mathbb{E} \Lambda \Delta N^{1-\alpha} + \text{Var}(\Lambda) \Delta^2 N^{2-2\alpha}}{1 - e^{-(\mu_1 + \mu_2) \Delta N^{-\alpha}}}, \end{aligned}$$

which behaves in accordance with (3.2) for  $N$  large.  $\square$

Recall that  $\gamma = \max\{1, 2 - \alpha\}$ ; the above computation shows that the covariance matrix of  $\mathbf{M}^{(N)}$  is essentially proportional to  $N^\gamma$ . Therefore, we expect that the centered and normalized version of the joint stationary queue length process converges to a (zero-mean)  $d$ -dimensional Gaussian random vector with covariance matrix  $C$  such that

$$C_{ik} = \begin{cases} 1_{\{\alpha \leq 1\}} \mathbb{E} \Lambda / \mu_i + 1_{\{\alpha \geq 1\}} \Delta \text{Var}(\Lambda) / (2\mu_i) & \text{if } i = k, \\ 1_{\{\alpha \leq 1\}} \mathbb{E} \Lambda / (\mu_i + \mu_k) + 1_{\{\alpha \geq 1\}} \Delta \text{Var}(\Lambda) / (\mu_i + \mu_k) & \text{if } i \neq k, \end{cases}$$

for  $i, k \in \{1, \dots, d\}$ . This is verified in the next subsection.

**3.2. Proof of functional central limit theorem based on MCLT.** The main objective of this subsection is to derive a functional limit theorem for  $\mathbf{M}^{(N)}(t)$ , the vector describing the queue lengths of the scaled system at time  $t$ . To this end, we consider the process  $\tilde{M}_i^{(N)}(\cdot) := M_i^{(N)}(\cdot)/N$ , for which we have

$$(3.3) \quad \tilde{M}_i^{(N)}(t) = \tilde{M}_i^{(N)}(0) + N^{-1} Y_0(N \Psi[\Lambda^{(N)}](t)) - N^{-1} Y_i(N \mu_i \Psi[\tilde{M}^{(N)}](t)).$$

We will need the following lemma, which uses the law of large numbers for Poisson processes; see [1].

LEMMA 3.2. *Let  $Y$  be a unit-rate Poisson process. Then for any  $U > 0$ , almost surely*

$$\lim_{N \rightarrow \infty} \sup_{0 \leq u \leq U} \left| \frac{Y(Nu)}{N} - u \right| = 0.$$

The uniform convergence in Lemma 3.2 entails that (3.3) converges almost surely to the solution of the functional equation

$$(3.4) \quad \varrho_i(t) = \varrho_{i,0} + t \mathbb{E}\Lambda - \mu_i \Psi[\varrho_i](t),$$

as  $N \rightarrow \infty$ , under the proviso that  $\tilde{M}_i^{(N)}(0)$  converges a.s. to some value  $\varrho_{i,0}$  for  $i = 1, \dots, d$ . The solution is given by a convex mixture of the initial position  $\varrho_i(0) = \varrho_{i,0}$  and the limiting value  $\mathbb{E}\Lambda/\mu_i$ :

$$(3.5) \quad \varrho_i(t) = \varrho_{i,0} e^{-\mu_i t} + \frac{\mathbb{E}\Lambda}{\mu_i} (1 - e^{-\mu_i t}).$$

Having identified this *fluid limit*, the next objective is to establish an FCLT for the centered and normalized process  $\mathbf{U}^{(N)}(\cdot)$  given by

$$(3.6) \quad U_i^{(N)}(t) := N^{\frac{\beta}{2}} (\tilde{M}_i^{(N)}(t) - \varrho_i(t)),$$

with  $\beta := 2 - \gamma = \min\{1, \alpha\}$ . Here we closely follow the approach in [1], where the idea is to use an MCLT, so as to obtain weak convergence to a (generalized) Ornstein-Uhlenbeck process. The version of the MCLT that we need in our setting is stated below.

THEOREM 3.1 (MCLT, [1]). *Let  $\{\mathbf{R}^{(N)}\}_{N \in \mathbb{N}}$  be a sequence of  $\mathbb{R}^d$ -valued martingales. Assume that the following condition on the jump sizes is met:*

$$(3.7) \quad \lim_{N \rightarrow \infty} \mathbb{E} \left[ \sup_{s \leq t} |\mathbf{R}^{(N)}(s) - \mathbf{R}^{(N)}(s^-)| \right] = 0;$$

*in addition, assume that, as  $N \rightarrow \infty$ ,*

$$[R_i^{(N)}, R_k^{(N)}]_t \rightarrow C_{ik}(t)$$

*for a deterministic function  $C_{ik}(t)$ , continuous in  $t$  for all  $t > 0$  and for  $i, k = 1, \dots, d$ . Then the process  $\mathbf{R}^{(N)}$  converges weakly to a centered Gaussian process  $\mathbf{W}$  with independent increments whose covariance matrix is characterized by*

$$\mathbb{E}[W_i(t) \cdot W_k(t)^T] = C_{ik}(t).$$

Introducing compensated unit-rate Poisson processes  $\tilde{Y}_i(t) := Y_i(t) - t$ , we define

$$(3.8) \quad \check{\mathbf{Y}}_0^{(N)}(t) := N^{\frac{\beta}{2}-1} \begin{pmatrix} \tilde{Y}_0(N\Psi[\Lambda^{(N)}](t)) \\ \vdots \\ \tilde{Y}_0(N\Psi[\Lambda^{(N)}](t)) \end{pmatrix},$$

$$(3.9) \quad \check{\mathbf{Y}}^{(N)}(t) := N^{\frac{\beta}{2}-1} \begin{pmatrix} \tilde{Y}_1(N\mu_1\Psi[\tilde{M}_1^{(N)}](t)) \\ \vdots \\ \tilde{Y}_d(N\mu_d\Psi[\tilde{M}_d^{(N)}](t)) \end{pmatrix}.$$

**LEMMA 3.3.** *Consider the  $d$ -dimensional processes  $\check{\mathbf{Y}}_0^{(N)}(\cdot)$  and  $\check{\mathbf{Y}}^{(N)}(\cdot)$ . If  $\alpha \geq 1$ , then as  $N \rightarrow \infty$  these processes converge weakly to  $d$ -dimensional zero-mean Brownian motions with covariance matrices  $K_0(t) := (t\mathbb{E}\Lambda)\mathbf{1}\mathbf{1}^\top$  and  $K(t) := \text{diag}\{\mu_1\Psi[\varrho_1](t), \dots, \mu_d\Psi[\varrho_d](t)\}$ , respectively; if  $\alpha < 1$  the limiting covariance matrices equal  $\mathbf{0}$ .*

**PROOF.** We start by checking the conditions of Thm. 3.1. First, observe that for each  $N$ ,  $\check{\mathbf{Y}}_0^{(N)}(\cdot)$  and  $\check{\mathbf{Y}}^{(N)}(\cdot)$  are  $d$ -dimensional real-valued martingales. Also, condition (3.7) is met, as both for  $\mathbf{R}^{(N)} = \check{\mathbf{Y}}_0^{(N)}$  and  $\mathbf{R}^{(N)} = \check{\mathbf{Y}}^{(N)}$ ,

$$\lim_{N \rightarrow \infty} \mathbb{E} \left[ \sup_{s \leq t} \left| \mathbf{R}^{(N)} - \mathbf{R}^{(N)}(s^-) \right| \right] < \infty,$$

whereas  $N^{\frac{\beta}{2}-1} \leq 1/\sqrt{N}$  converges to zero.

Note that  $\beta = \min\{1, \alpha\}$ , so that  $\beta - 2 = \min\{-1, \alpha - 2\}$ . Now observe that for  $\alpha \geq 1$  (and hence  $\beta - 2 = -1$ ) the quadratic covariation of  $\check{\mathbf{Y}}_0^{(N)}(\cdot)$ ,

$$\left[ N^{\frac{\beta}{2}-1} \tilde{Y}_0(N\Psi[\Lambda^{(N)}](t)) \right]_t = N^{\beta-2} Y_0(N\Psi[\Lambda^{(N)}](t)),$$

converges to  $t\mathbb{E}\Lambda$  as  $N \rightarrow \infty$  (0 for  $\alpha < 1$ ), by virtue of Lemma 3.2. The covariance matrix for  $\check{\mathbf{Y}}^{(N)}(\cdot)$  is determined in the same way; for  $\alpha \geq 1$  the diagonal entries are given by

$$\begin{aligned} \lim_{N \rightarrow \infty} \left[ N^{\frac{\beta}{2}-1} \tilde{Y}_i(N\mu_i\Psi[\tilde{M}_i^{(N)}](t)) \right]_t &= \lim_{N \rightarrow \infty} N^{\beta-2} Y_i(N\mu_i\Psi[\tilde{M}_i^{(N)}](t)) \\ &= \mu_i\Psi[\varrho_i](t) \end{aligned}$$

(which would equal 0 for  $\alpha < 1$ ), whereas for  $i \neq k$  (then  $\tilde{Y}_i(\cdot)$  and  $\tilde{Y}_k(\cdot)$  are independent)

$$\lim_{N \rightarrow \infty} [N^{\frac{\beta}{2}-1} \tilde{Y}_i(N\mu_i \Psi[\tilde{M}_i^{(N)}](t)), N^{\frac{\beta}{2}-1} \tilde{Y}_k(N\mu_k \Psi[\tilde{M}_k^{(N)}](t))]_t = 0,$$

with  $i, k \in \{1, \dots, d\}$ . For  $\alpha \geq 1$ , Thm. 3.1 yields that the processes converge weakly to  $d$ -dimensional Brownian motions with covariance matrices  $K_0(t)$  and  $K(t)$ . On the other hand, for  $\alpha < 1$  the entries of the covariance matrices all vanish as  $N \rightarrow \infty$ . As a result, both  $\tilde{\mathbf{Y}}_0^{(N)}(\cdot)$  and  $\tilde{\mathbf{Y}}^{(N)}(\cdot)$  converge to a process identical to  $\mathbf{0}$ .  $\square$

Stated below is the main theorem of this section: an FCLT for  $\mathbf{U}^{(N)}(\cdot)$ , the process defined via (3.6). In line with earlier findings, three regimes need to be distinguished:  $\alpha > 1$  (the fast regime),  $\alpha < 1$  (the slow regime) and  $\alpha = 1$  (the intermediate regime).

**THEOREM 3.2 (FCLT).** *As  $N \rightarrow \infty$ ,  $\mathbf{U}^{(N)}(\cdot)$  converges weakly to a zero-mean  $d$ -dimensional Gaussian process with covariance matrix given by*

(3.10)

$$C_{ii}(t) := 1_{\{\alpha \geq 1\}} (\mathbb{E}\Lambda/\mu_i + \varrho_{i,0} e^{-\mu_i t}) (1 - e^{-\mu_i t}) \\ + 1_{\{\alpha \leq 1\}} \Delta \text{Var } \Lambda / (2\mu_i) (1 - e^{-2\mu_i t}),$$

(3.11)

$$C_{ik}(t) := (1_{\{\alpha \geq 1\}} \mathbb{E}\Lambda / (\mu_i + \mu_k) + 1_{\{\alpha \leq 1\}} \Delta \text{Var } \Lambda / (\mu_i + \mu_k)) \cdot (1 - e^{-(\mu_i + \mu_k)t}),$$

for  $i \neq k$  ( $i, k \in \{1, \dots, d\}$ ).

**PROOF.** Using (3.3), we write

(3.12)

$$U_i^{(N)}(t) = N^{\frac{\beta}{2}} (\tilde{M}_i^{(N)}(0) + N^{-1} Y_0(N\Psi[\Lambda^{(N)}](t)) \\ - N^{-1} Y_i(N\mu_i \Psi[\tilde{M}_i^{(N)}](t)) - \varrho_i(t)),$$

for  $i = 1, \dots, d$ . Adding and subtracting  $\varrho_{i,0}$ , (3.12) is equivalent to

$$U_i^{(N)}(t) = N^{\frac{\beta}{2}} (\tilde{M}_i^{(N)}(0) - \varrho_{i,0}) - N^{\frac{\beta}{2}} (\varrho_i(t) - \varrho_{i,0}) \\ + N^{\frac{\beta}{2}-1} (Y_0(N\Psi[\Lambda^{(N)}](t)) - Y_i(N\mu_i \Psi[\tilde{M}_i^{(N)}](t))),$$

which, by filling out the implicit form of  $\varrho_i(t)$  as in (3.4), simplifies to

$$U_i^{(N)}(t) = U_{0,i}^{(N)}(t) + U_1^{(N)}(t) + U_{2,i}^{(N)}(t),$$

with

$$\begin{aligned} U_{0,i}^{(N)}(t) &:= U_i^{(N)}(0) - \mu_i \Psi[U_i^{(N)}](t), \\ U_1^{(N)}(t) &:= N^{\frac{\beta}{2}} (\Psi[\Lambda^{(N)}](t) - t \mathbb{E}\Lambda), \\ U_{2,i}^{(N)}(t) &:= N^{\frac{\beta}{2}-1} (\tilde{Y}_0(N\Psi[\Lambda^{(N)}](t)) - \tilde{Y}_i(N\mu_i \Psi[\tilde{M}_i^{(N)}](t))). \end{aligned}$$

We consider the three individual components separately.

- (i) Component  $U_{0,i}^{(N)}(t)$  consists of the starting value of the process, which is assumed to converge to some value  $U_i(0)$ , minus a reverting term. It is now straightforward that, as  $N \rightarrow \infty$ ,  $U_{0,i}^{(N)}(\cdot)$  converges to  $U_{0,i}(t) = U_i(0) - \mu_i \Psi[U_i](t)$ .
- (ii) Then consider  $U_1^{(N)}(\cdot)$ . For  $\alpha \leq 1$  (and hence  $\frac{\beta}{2} = \frac{\alpha}{2}$ ), the standard functional central limit theorem for partial sums of i.i.d. random variables entails that, as  $N \rightarrow \infty$ ,

$$U_1^{(N)}(\cdot) \rightarrow \sqrt{\Delta \text{Var } \Lambda} \cdot V(\cdot),$$

with  $V(\cdot)$  a standard Brownian motion. On the other hand, for  $\alpha > 1$  the limiting process is identical to  $\mathbf{0}$ , as a consequence of  $\frac{\beta}{2} = \frac{1}{2} < \frac{\alpha}{2}$ .

- (iii) Finally, from Lemma 3.3, we conclude that  $U_2^{(N)}(\cdot)$  converges weakly to a  $d$ -dimensional zero-mean Brownian motion with covariance matrix  $K_0(t) + K(t)$  for  $\alpha \geq 1$ , and to  $\mathbf{0}$  else.

Using the above observations, we can now complete the proof. Each of the three regimes will be considered separately.

1. *Fast regime* ( $\alpha > 1$ ). We have obtained above that  $U^{(N)}(\cdot)$  converges weakly to the solution  $U(\cdot)$  of the  $d$ -dimensional stochastic integral equation given by

$$U_i(t) = U_i(0) - \mu_i \Psi[U_i](t) + W_i(t \mathbb{E}\Lambda + \mu_i \Psi[\varrho_i](t)) \quad \text{for } i = 1, \dots, d$$

with  $W_1(\cdot), \dots, W_d(\cdot)$  standard Brownian motions (but not independent), or equivalently

$$U_i(t) = U_i(0) - \mu_i \Psi[U_i](t) + \tilde{W}_0(t \mathbb{E}\Lambda) + \tilde{W}_i(\mu_i \Psi[\varrho_i](t))$$

with  $\tilde{W}_0(\cdot), \tilde{W}_1(\cdot), \dots, \tilde{W}_d(\cdot)$  independent standard Brownian motions. It takes a routine calculation to derive that

$$U_i(t) = e^{-\mu_i t} (U_i(0) + \int_0^t \sqrt{\mathbb{E}\Lambda + \mu_i \varrho_i(s)} e^{\mu_i s} dW_i(s)).$$

All linear combinations of the  $U_i(\cdot)$  are Gaussian processes, so we conclude that this  $d$ -dimensional process is Gaussian. It is readily seen that  $\mathbb{E} U_i(t) = U_i(0)e^{-\mu_i t}$ . For the variance, an elementary computation gives

$$\text{Var } U_i(t) = e^{-2\mu_i t} \left( \int_0^t (\mathbb{E}\Lambda + \mu_i \varrho_i(s)) e^{2\mu_i s} ds \right) = (\mathbb{E}\Lambda / \mu_i + \varrho_{i,0} e^{-\mu_i t}) (1 - e^{-\mu_i t}).$$

Likewise, for the covariance, with

$$\mathcal{U}_i(t) := \sqrt{\mathbb{E}\Lambda} \int_0^t e^{\mu_i s} d\tilde{W}_0(s) + \int_0^t \sqrt{\mu_i \varrho_i(s)} e^{\mu_i s} d\tilde{W}_i(s),$$

it follows that, for  $i \neq k$ ,

$$\begin{aligned} \text{Cov}(U_i(t), U_k(t)) &= e^{-(\mu_i + \mu_k)t} \mathbb{E} [\mathcal{U}_i(t) \mathcal{U}_k(t)] \\ &= e^{-(\mu_i + \mu_k)t} \mathbb{E}\Lambda \cdot \mathbb{E} \left[ \int_0^t e^{-\mu_i s} d\tilde{W}_0(s) \cdot \int_0^t e^{-\mu_k s} d\tilde{W}_0(s) \right] \\ &= e^{-(\mu_i + \mu_k)t} \mathbb{E}\Lambda \int_0^t e^{(\mu_i + \mu_k)s} ds \\ &= \mathbb{E}\Lambda / (\mu_i + \mu_k) (1 - e^{-(\mu_i + \mu_k)t}). \end{aligned}$$

This shows (3.10) for  $\alpha > 1$ .

*2. Slow regime* ( $\alpha < 1$ ). In the slow regime,  $\mathbf{U}^{(N)}(\cdot)$  converges to the solution of

$$U_i(t) = U_i(0) - \mu_i \Psi[U_i](t) + (\sqrt{\Delta \text{Var } \Lambda}) V(t) \text{ for } i = 1, \dots, d,$$

which can be written as

$$dU_i(t) = -\mu_i U_i(t) dt + (\sqrt{\Delta \text{Var } \Lambda}) dV(t).$$

Therefore the  $U_i(\cdot)$  are Ornstein-Uhlenbeck processes given by:

$$U_i(t) = e^{-\mu_i t} (U_i(0) + \int_0^t \sqrt{\Delta \text{Var}(\Lambda)} e^{\mu_i s} dV(s)).$$

As before, we can conclude that this  $d$ -dimensional process is Gaussian with expectation vector given by  $U_i(0)e^{-\mu_i t}$ . Computations as above reveal that for  $\alpha < 1$ , as claimed in (3.10),

$$\text{Cov}(U_i(t), U_k(t)) = \Delta \text{Var}(\Lambda) / (\mu_i + \mu_k) (1 - e^{-(\mu_i + \mu_k)t}).$$

*3. Intermediate regime* ( $\alpha = 1$ ). In this regime, a combination of the processes from the other cases appears:

$$dU_i(t) = -\mu_i U_i(t) dt + \sqrt{\mathbb{E}\Lambda} d\tilde{W}_0(t) + \sqrt{\mu_i \varrho_i(t)} d\tilde{W}_i(t) + \sqrt{\Delta \text{Var}(\Lambda)} dV(t).$$

The marginal solutions are, for  $i = 1, \dots, d$ , equal to

$$U_i(t) = e^{-\mu_i t} (U_i(0) + \int_0^t \sqrt{\mathbb{E}\Lambda} e^{\mu_i s} d\tilde{W}_0(s) + \int_0^t \sqrt{\mu_i \varrho_i(s)} e^{\mu_i s} d\tilde{W}_i(s) + \int_0^t \sqrt{\Delta \text{Var}(\Lambda)} e^{\mu_i s} dV(s)).$$

Again, we conclude that this  $d$ -dimensional process is Gaussian with expectation vector given by  $U_i(0)e^{-\mu_i t}$ ; routine computations yield the desired covariance matrix, as given in (3.10) and (3.11). This completes the proof.  $\square$

It is interesting to study the impact of the scaling parameter  $\alpha$  on the correlation between the individual queue lengths. Remarkably, it turns out that for  $\alpha \neq 1$  this correlation depends on the service rates only, whereas for  $\alpha = 1$  also the first and second moment of  $\Lambda$  play a role; see the following corollary for a result on the stationary regime.

**COROLLARY 3.1** (Correlation coefficients). *In stationarity, the correlation coefficient for  $i \neq k$  satisfies*

$$(3.13) \quad \lim_{N \rightarrow \infty} \text{Corr}(M_i^{(N)}, M_k^{(N)}) = c_{ik}(\alpha) \cdot \frac{\sqrt{\mu_i \mu_k}}{\mu_i + \mu_k},$$

for some constant  $c_{ik}(\alpha) \in [1, 2]$ . The constant  $c_{ik}(\alpha)$  equals 1 for  $\alpha > 1$  and 2 for  $\alpha < 1$ .

**PROOF.** From Thm. 3.2, as  $t \rightarrow \infty$ ,

$$C_{ik}(t) \rightarrow 1_{\{\alpha \geq 1\}} \frac{\mathbb{E}\Lambda}{\mu_i + \mu_k 1_{\{i \neq k\}}} + 1_{\{\alpha \leq 1\}} \frac{\Delta \text{Var}(\Lambda)}{\mu_i + \mu_k}.$$

We observe that (3.13) holds, with

$$c_{ik}(\alpha) = \frac{\mathbb{E}\Lambda 1_{\{\alpha \geq 1\}} + \Delta \text{Var}(\Lambda) 1_{\{\alpha \leq 1\}}}{\mathbb{E}\Lambda 1_{\{\alpha \geq 1\}} + \frac{1}{2} \Delta \text{Var}(\Lambda) 1_{\{\alpha \leq 1\}}}. \quad \square$$

**4. Large deviations.** Where the previous section studied the random-environment infinite-server system under the *central limit scaling*, we now focus on the *large deviations domain*. As it turns out, the previously observed trichotomy remains valid. The results again translate to the setting with  $d$  coupled queues; for ease we first present (and prove) the results for  $d = 1$ , to return to the coupled model at the end of the section.

4.1. *Univariate large deviations.* Let the arrival rate of the  $N$ -scaled model again be given by  $N\Lambda^{(N)}(t)$  (see (2.10)). An important quantity in our analysis is

$$\begin{aligned}\kappa_t(\Lambda^{(N)}) &= \int_0^t \Lambda^{(N)}(s) e^{-\mu s} ds \\ &= \frac{1}{\mu} (1 - e^{-\mu \Delta N^{-\alpha}}) \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} \Lambda_j e^{-\mu j \Delta N^{-\alpha}} + o_p(1),\end{aligned}$$

as  $N \rightarrow \infty$ . As observed earlier,  $M^{(N)}(t)$  is a mixed Poisson random variable, with random parameter distributed as  $N\kappa_t(\Lambda^{(N)})$ . In the large deviations setting we are interested in the tail probabilities of  $M^{(N)}(t)$  for given  $t$  and  $N$  large. More specifically, our objective is to evaluate the decay rate

$$\lim_{N \rightarrow \infty} N^{-\beta} \log \mathbb{P}(M^{(N)}(t)/N \geq a),$$

for any  $a > \rho(t) = \rho(1 - e^{-\mu t})$  (where  $\rho := \lambda/\mu$ ) and some specific  $\beta > 0$ . Given the results obtained in the central limit regime, we expect that  $\beta = \min\{1, \alpha\}$ .

The main idea behind our analysis is to condition on the value of the random Poisson parameter. In self-evident notation,

$$\begin{aligned}\mathbb{P}(M^{(N)}(t)/N \geq a) &= \mathbb{P}(\text{Pois}(N\kappa_t(\Lambda^{(N)})) \geq Na) \\ (4.1) \quad &= \int_0^\infty \mathbb{P}(\text{Pois}(Nx) \geq Na) \mathbb{P}(\kappa_t(\Lambda^{(N)}) \in dx).\end{aligned}$$

In some parts of the analysis we rely on the following lemma, in which we establish a large deviation result for  $\mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq a)$ .

LEMMA 4.1. *Let  $a > \rho(t)$ . Then, with  $M(\theta) := \mathbb{E} e^{\theta \Lambda}$ ,*

$$(4.2) \quad \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq a) = - \sup_{\theta > 0} \left( \theta a - \int_0^t \log M(\theta e^{-\mu s}) ds \right)$$

PROOF. As a first step, we define a proxy for  $\kappa_t(\Lambda^{(N)})$  that is easier to work with:

$$(4.3) \quad k_t(\Lambda^{(N)}) := \Delta N^{-\alpha} \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} \Lambda_j e^{-\mu j \Delta N^{-\alpha}};$$

later we show that  $\kappa_t(\Lambda^{(N)})$  and  $k_t(\Lambda^{(N)})$  are ‘close enough’. Let  $P_N(a) := \mathbb{P}(k_t(\Lambda^{(N)}) \geq a)$ . Writing, for arbitrary  $\theta > 0$ ,



$$\begin{aligned}
P_N(a) &= \mathbb{P}(e^{\theta k_t(\Lambda^{(N)})/(\Delta N^{-\alpha})} \geq e^{\theta a/(\Delta N^{-\alpha})}) \\
&= \mathbb{P}\left(\prod_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} e^{\theta \Lambda_j e^{-\mu j \Delta N^{-\alpha}}} \geq e^{\theta a/(\Delta N^{-\alpha})}\right),
\end{aligned}$$

Markov's inequality immediately yields the upper bound

$$P_N(a) \leq e^{-\theta a/(\Delta N^{-\alpha})} \prod_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} M(\theta e^{-\mu j \Delta N^{-\alpha}}).$$

Recognizing a Riemann sum, we thus obtain

$$\begin{aligned}
\limsup_{N \rightarrow \infty} \Delta N^{-\alpha} \log P_N(a) &\leq \limsup_{N \rightarrow \infty} \Delta N^{-\alpha} \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} \log M(\theta e^{-\mu j \Delta N^{-\alpha}}) - \theta a \\
&= \int_0^t \log M(\theta e^{-\mu s}) ds - \theta a.
\end{aligned}$$

As the established upper bound holds for any  $\theta > 0$ ,

$$(4.4) \quad \limsup_{N \rightarrow \infty} \Delta N^{-\alpha} \log P_N(a) \leq \inf_{\theta > 0} \left( \int_0^t \log M(\theta e^{-\mu s}) ds - \theta a \right).$$

The next goal is to prove that the above upper bound is tight. We do so by first noting that, due to the convexity of the function involved, the infimum in the right-hand side of (4.4) is attained by  $\theta^*$ , being the unique solution to

$$\frac{\partial}{\partial \theta} \int_0^t \log M(\theta e^{-\mu s}) ds \Big|_{\theta=\theta^*} = a.$$

Now we apply a change-of-measure technique. Define a measure  $\mathbb{Q}$  by exponential twisting; the density of the  $\Lambda_j$  is changed into

$$\mathbb{Q}(\Lambda_j \in dx) := \frac{e^{\theta^* e^{-\mu j \Delta N^{-\alpha}} x}}{M(\theta^* e^{-\mu j \Delta N^{-\alpha}})} \mathbb{P}(\Lambda_j \in dx).$$

Fix an  $\varepsilon > 0$ , and let the event  $\mathcal{E}_a^{(N)} := \{k_t(\Lambda^{(N)}) \in [a, a + \varepsilon]\}$ . Then

$$\begin{aligned}
P_N(a) &= \mathbb{E}_{\mathbb{Q}} \left[ \mathbb{1}_{\mathcal{E}_a^{(N)}} \prod_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} M(\theta^* e^{-\mu j \Delta N^{-\alpha}}) e^{-\theta^* \Lambda_j e^{-\mu j \Delta N^{-\alpha}}} \right] \\
&\geq \mathbb{Q}(k_t(\Lambda^{(N)}) \in [a, a + \varepsilon]) e^{-\theta^*(a+\varepsilon)/(\Delta N^{-\alpha})} \prod_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} M(\theta^* e^{-\mu j \Delta N^{-\alpha}}).
\end{aligned}$$

To obtain that  $\mathbb{Q}(k_t(\Lambda^{(N)}) \in [a, a + \varepsilon]) \rightarrow \frac{1}{2}$  as  $N \rightarrow \infty$ , we now show that  $k_t(\Lambda^{(N)})$  is asymptotically normal. It is verified that  $\mathbb{E}_{\mathbb{Q}} k_t(\Lambda^{(N)}) \rightarrow a$  as  $N \rightarrow \infty$ , due to the specific construction of the measure  $\mathbb{Q}$ . Also,

$$\begin{aligned} \lim_{N \rightarrow \infty} \text{Var}_{\mathbb{Q}}(k_t(\Lambda^{(N)})) &= \lim_{N \rightarrow \infty} \text{Var}_{\mathbb{Q}}(\Delta N^{-\alpha} \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} \Lambda_j e^{-\mu j \Delta N^{-\alpha}}) \\ &= \lim_{N \rightarrow \infty} \Delta N^{-\alpha} (\Delta N^{-\alpha} \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} e^{-2\mu j \Delta N^{-\alpha}}) \text{Var}_{\mathbb{Q}}(\Lambda) \\ &= \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \int_0^t e^{-2\mu s} ds \text{Var}_{\mathbb{Q}}(\Lambda) = 0. \end{aligned}$$

Copying the approach – using cumulant generating functions – underlying the proof of Theorem 2.1, it is readily derived that

$$N^{\frac{\alpha}{2}}(k_t(\Lambda^{(N)}) - a) \xrightarrow{d} \mathcal{N}(0, \sigma^2), \quad \text{where } \sigma^2 := \Delta/(2\mu)(1 - e^{-2\mu t})\text{Var}_{\mathbb{Q}}(\Lambda).$$

Hence,

$$\begin{aligned} \liminf_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}_N(a) &\geq \liminf_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{Q}(k_t(\Lambda^{(N)}) \in [a, a + \varepsilon]) \\ &\quad - \theta^*(a + \varepsilon) + \int_0^t \log M(\theta^* e^{-\mu s}) ds \\ &\geq \int_0^t \log M(\theta^* e^{-\mu s}) ds - \theta^*(a + \varepsilon). \end{aligned}$$

By letting  $\varepsilon \downarrow 0$ , together with the upper bound this leads to

$$(4.5) \quad \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log P_N(a) = -\sup_{\theta > 0} \left( \theta a - \int_0^t \log M(\theta e^{-\mu s}) ds \right).$$

Now it remains to show that  $k_t(\Lambda^{(N)})$  can again be replaced by  $\kappa_t(\Lambda^{(N)})$  (which we abbreviate for compactness to  $k_t$  and  $\kappa_t$ ). Note that, as  $N \rightarrow \infty$ ,

$$|\kappa_t - k_t| = \left| \left( \frac{1 - e^{-\mu \Delta N^{-\alpha}}}{\mu \Delta N^{-\alpha}} - 1 \right) \Delta N^{-\alpha} \sum_{j=0}^{\lfloor t/(\Delta N^{-\alpha}) \rfloor - 1} \Lambda_j e^{-\mu j \Delta N^{-\alpha}} + o_p(1) \right| = o_p(1).$$

Let  $\eta > 0$  small enough to guarantee  $a - \eta > \rho(t)$ . Then  $\mathbb{P}(\kappa_t \in (k_t - \eta, k_t + \eta)) \rightarrow 1$  as  $N \rightarrow \infty$ , hence

$$\begin{aligned} \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}_N(a + \eta) &\leq \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}(\kappa_t \geq a) \\ &\leq \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}_N(a - \eta), \end{aligned}$$

which provides bounds for the decay rate of interest of the form

$$(4.6) \quad -\sup_{\theta > 0} \left( \int_0^t \log M e^{-\mu s} ds - \theta(a \pm \eta) \right).$$

The rate function in (4.6) is continuous in  $\eta$ , so now letting  $\eta \downarrow 0$  yields (4.2).  $\square$

As in the central limit regime, we distinguish between the cases  $\alpha > 1$ ,  $\alpha = 1$ , and  $\alpha < 1$ . For all three cases we derive the logarithmic asymptotics.

1. *Fast regime* ( $\alpha > 1$ ). We can bound (4.1) from below by

$$(4.7) \quad \mathbb{P}(\text{Pois}(N(\rho(t) - \varepsilon)) \geq Na) \cdot \mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq \rho(t) - \varepsilon),$$

for some  $\varepsilon \in (0, a - \rho(t))$ . As  $N$  tends to infinity, it is directly shown that the second factor in (4.7) converges to 1, and hence has exponential decay rate 0. Now an application of Cramér's theorem [7] yields

$$\begin{aligned} \liminf_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(\text{Pois}(N \kappa_t(\Lambda^{(N)})) \geq Na) \\ \geq \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(\text{Pois}(N(\rho(t) - \varepsilon)) \geq Na) \\ = -\sup_{\theta} (\theta a - (\rho(t) - \varepsilon)(e^\theta - 1)) \\ = a \log \left( \frac{\rho(t) - \varepsilon}{a} \right) - (\rho(t) - \varepsilon) + a. \end{aligned}$$

On the other hand, (4.1) is majorized by

$$(4.8) \quad \mathbb{P}(\text{Pois}(N(\rho(t) + \varepsilon)) \geq Na) + \mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq \rho(t) + \varepsilon).$$

By Cramér's theorem, the first term in (4.8) decays exponentially in  $N$ . As a consequence of Lemma 4.1, the second term decays exponentially in  $N^\alpha$ , i.e., superexponentially in  $N$ . This yields

$$\begin{aligned} \limsup_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(\text{Pois}(N \kappa_t(\Lambda^{(N)})) \geq Na) \\ \leq \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(\text{Pois}(N(\rho(t) + \varepsilon)) \geq Na) \\ = a \log \left( \frac{\rho(t) + \varepsilon}{a} \right) - (\rho(t) + \varepsilon) + a. \end{aligned}$$

As this holds for all  $\varepsilon > 0$ , we conclude that

$$\lim_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(M_N(t)/N \geq a) = a \log \left( \frac{\rho(t)}{a} \right) - \rho(t) + a.$$

Recognizing the decay rate of a Poisson distribution with mean  $\rho(t)$ , we observe that the essential behavior in the fast regime is again of M/M/ $\infty$  type.

2. *Slow regime* ( $\alpha < 1$ ). In this regime we need to distinguish between the situation in which the random variable  $\Lambda$  almost surely results in a  $\kappa_t(\Lambda^{(N)})$  below  $a$ , and the situation in which this is not the case. The proof of the following lemma is straightforward hence omitted.

LEMMA 4.2. *Given  $\Lambda$ , let  $y = \inf\{x > 0 : \mathbb{P}(\Lambda \leq x) = 1\}$  and  $u(t) = y/\mu(1 - e^{-\mu t})$ . Then, as  $N \rightarrow \infty$ ,*

$$\mathbb{P}(\kappa_t(\Lambda^{(N)}) \leq u(t)) \rightarrow 1.$$

The cases  $u(t) \geq a$  and  $u(t) < a$  should be treated differently, as follows from the following intuitive explanation that is based on the decomposition (4.1). If  $u(t) \geq a$ , then the random variable  $\Lambda$  can be ‘large’ with respect to  $a$ , which enables  $M^{(N)}(t)$  to reach  $Na$  without the Poisson random variable attaining an unlikely value. If on the contrary  $u(t) < a$ , then  $\Lambda$  is ‘small’ with respect to  $a$ ;  $M^{(N)}(t)$  can only exceed level  $Na$  by the Poisson random variable attaining an extraordinarily large value.

We first consider the case  $u(t) < a$ . For ease we assume that  $\Lambda$  attains values in a discrete set of positive values, of which  $y$  is the largest (occurs with probability  $p \in (0, 1)$ ) and  $y' < y$  the one-but-largest. It is directly seen that, for  $\theta > 0$ ,

$$\mathbb{E}e^{\theta M^{(N)}(t)} \geq p^{\lceil t/(\Delta N^{-\alpha}) \rceil} \mathbb{E} \exp(\theta \text{Pois}(Nu(t))).$$

As  $\alpha < 1$ , this leads to

$$(4.9) \quad \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{E}e^{\theta M^{(N)}(t)} \geq u(t)(e^\theta - 1).$$

In addition,  $\mathbb{E}e^{\theta M^{(N)}(t)}$  is majorized by

$$\begin{aligned} & p \mathbb{E} \exp(\theta \text{Pois}(Nu(t))) + (1-p) \mathbb{E} \exp(\theta \text{Pois}(N(y'/\mu)(1 - e^{-\mu t}))) \\ & = p \exp(Nu(t)(e^\theta - 1)) + (1-p) \exp(N(y'/\mu)(1 - e^{-\mu t})(e^\theta - 1)), \end{aligned}$$

which converges to the right-hand side of (4.9) on an exponential scale (use  $y > y'$ ). Applying ‘Cramér’, we thus find that the probability of interest decays exponentially:

$$\begin{aligned} \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(M^{(N)}(t)/N \geq a) &= - \sup_{\theta > 0} (\theta a - u(t)(e^\theta - 1)) \\ &= a \log\left(\frac{u(t)}{a}\right) + a - u(t). \end{aligned}$$

Now we focus on  $u(t) \geq a$ ; in this case

$$(4.10) \quad \mathbb{P}(\text{Pois}(Na) \geq Na) \mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq a)$$

gives an asymptotically non-trivial lower bound for (4.1). Note that for every  $\delta > 0$ , there is an  $N$  large enough such that

$$\mathbb{P}(\text{Pois}(Na) \geq Na) \geq \left(\frac{1}{2} - \delta\right),$$

so the first factor in (4.10) will not contribute to the decay rate. The tail behavior of the second factor follows from Lemma 4.1. On the other hand, (4.1) is majorized by

$$(4.11) \quad \mathbb{P}(\text{Pois}(N(a - \varepsilon)) \geq Na) + \mathbb{P}(\kappa_t(\Lambda^{(N)}) \geq a - \varepsilon).$$

Again it is observed that only the second term in (4.11) contributes to the decay rate: by ‘Cramér’ the first term in (4.11) decays exponentially, whereas the decay of the second term is subexponential (by Lemma 4.1) for  $\varepsilon > 0$  small enough (we need  $a - \varepsilon > \rho(t)$ ). Letting  $\varepsilon \downarrow 0$  while using that the rate function in (4.2) is continuous in  $a$ , we arrive at

$$\lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}(M^{(N)}(t)/N \geq a) = - \sup_{\theta > 0} \left( \theta a - \int_0^t \log M(\theta e^{-\mu s}) ds \right).$$

Note that the decay rate in this fast regime depends on more detailed information on the distribution of  $\Lambda$  than just the mean.

3. *Intermediate regime* ( $\alpha = 1$ ). In this regime we expect exponential decay. Indeed, it is directly derived that

$$\lim_{N \rightarrow \infty} \Delta N^{-1} \log \mathbb{E} e^{\theta M^{(N)}(t)} = \int_0^t \log M(\Delta(e^\theta - 1)e^{-\mu s}) ds,$$

and hence ‘Gärtner-Ellis’ [7] gives

$$\begin{aligned} \lim_{N \rightarrow \infty} \Delta N^{-1} \log \mathbb{P}(M^{(N)}(t)/N \geq a) \\ = - \sup_{\theta > 0} \left( \theta a - \int_0^t \log M(\Delta(e^{\theta/\Delta} - 1)e^{-\mu s}) ds \right). \end{aligned}$$

For deterministic  $\Lambda$  the above result would equal that of the fast regime; the resemblance with the slow regime on the other hand becomes more pronounced for larger values of  $\Delta$ .

4.2. *Large deviations for the coupled model.* We conclude this section by pointing out how the large deviations for the coupled model (where each arrival generates work in  $d$  queues) can be dealt with. For  $\alpha \geq 1$  we are in the regime of exponential decay. The multivariate version of the Gärtner-Ellis theorem entails that, modulo the validity of mild regularity conditions to be imposed on the set  $A \subset \mathbb{R}_+^d$ ,

$$\begin{aligned} & \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{P}(\mathbf{M}^{(N)}(t)/N \in A) \\ &= - \inf_{\mathbf{a} \in A} \sup_{\boldsymbol{\theta}} \left( \sum_{i=1}^d \theta_i a_i - \lim_{N \rightarrow \infty} \frac{1}{N} \log \mathbb{E} \exp \left[ \sum_{i=1}^d \theta_i M_i^{(N)}(t) \right] \right). \end{aligned}$$

The problem therefore reduces to characterizing the limiting log moment generating function. It takes standard computations to verify that for  $\alpha > 1$ , with an argumentation borrowed from specific intermediate results in the proof of Lemma 3.1,

$$\begin{aligned} & \lim_{N \rightarrow \infty} N^{-1} \log \mathbb{E} \exp \left[ \sum_{i=1}^d \theta_i M_i^{(N)}(t) \right] \\ &= t \mathbb{E} \Lambda \left( \int_0^t \frac{1}{t} \prod_{i=1}^d (e^{-\mu_i s} (e^{\theta_i} - 1) + 1) ds - 1 \right), \end{aligned}$$

whereas for  $\alpha = 1$  it turns out to equal

$$\frac{1}{\Delta} \int_0^t \log \mathbb{E} \exp \left[ \Lambda \Delta \left( \left( \prod_{i=1}^d e^{-\mu_i s} (e^{\theta_i} - 1) + 1 \right) - 1 \right) \right] ds.$$

For  $\alpha < 1$ , as before, the decay is either exponential in  $N$  (if the the multi-dimensional random Poisson parameter cannot attain values that are contained in  $A$ ), or exponential in  $N^\alpha$ . The latter regime being the more complicated one, we here include the corresponding decay rate. The probability of our interest can be rewritten as

$$(4.12) \quad \int_{x_1=0}^{\infty} \cdots \int_{x_d=0}^{\infty} F_A(\mathbf{x}) \cdot \pi(dx_1, \dots, dx_d),$$

where

$$\begin{aligned} F_A(\mathbf{x}) &:= \mathbb{P}((\text{Pois}_1(N \kappa_t(\Lambda^{(N)})), \dots, \text{Pois}_d(N \kappa_t(\Lambda^{(N)}))/N \in A), \\ \pi(dx_1, \dots, dx_d) &:= \mathbb{P}(\kappa_{t,1}(\Lambda^{(N)}) \in dx_1, \dots, \kappa_{t,d}(\Lambda^{(N)}) \in dx_d); \end{aligned}$$

here the  $d$  Poisson random variables are independent. Using the same ideas as above, it can be shown that (4.12) decays exponentially in  $N^\alpha$ , where the decay rate is now given by

$$\begin{aligned} & \lim_{N \rightarrow \infty} \Delta N^{-\alpha} \log \mathbb{P}((\kappa_{t,1}, \dots, \kappa_{t,d})(\Lambda^{(N)}) \in A) \\ &= - \inf_{a \in A} \sup_{\theta} \left( \sum_{i=1}^d \theta_i a_i - \int_0^t \log M \left( \sum_{i=1}^d \theta_i e^{-\mu_i s} \right) ds \right). \end{aligned}$$

**5. Discussion and future research.** In this paper we propose to model an overdispersed arrival process by a mixed Poisson process in a random environment. We assess the impact of overdispersion on system performance when feeding such an arrival process into an infinite-server system. Under a specific scaling, we derive (functional) central limit results and large deviations asymptotics.

Various extensions can be explored before using the model in an operational context; a few of them are mentioned here. To start with, many results seem to carry over to the setting with generally distributed service times. In addition, systematically studying the effect of adding a deterministic trend  $\bar{\lambda}(\cdot)$  (e.g. to account for diurnal patterns) to the random environment  $\Lambda(\cdot)$ , the results could be generalized to a setting with nonstationary Cox arrival processes. Moreover, we could pursue to upgrade our model to not only allow for a deterministic trend, but also for dependence between arrival rates corresponding to subsequent time slots; such properties have been observed in (overdispersed) datasets, and are therefore desirable to incorporate [14]. Another challenge lies in refining the logarithmic asymptotics, as obtained in Section 4, to exact asymptotics.

In all of the results obtained, we revealed a trichotomy in system behavior depending on the imposed scaling on system size and sampling frequency. Here the scaling primarily serves to tweak the level of overdispersion in the system. The combination of tunable sampling and tunable overdispersion provides a rich framework for modeling real-world arrival processes. One could imagine that in a rapidly changing environment, the inherent overdispersion of the arrival process hardly plays a role, whereas in a slowly changing random environment, overdispersion is expected to be more dominant. This interplay between sampling frequency and overdispersion is a convenient feature of our model, which could be used to calibrate the model to real-world data. The latter could be a promising direction for future research, involving challenging statistical issues related to determining which of the three asymptotic regimes that arise in the revealed trichotomy suits best with a given data set for a finite  $N$  sys-

tem, and estimating the model parameters conditioned on the asymptotic regime.

An application of our model would be in the area of dimensioning service systems or staffing, and in particular square-root staffing in many-server systems. The general idea behind square-root staffing is as follows: a finite-server system is modeled as a system in heavy traffic, where the number of servers  $s$  is large and at the same time, the system is critically loaded. Under Markovian assumptions this can be achieved by setting  $s = \lambda + \beta\sqrt{\lambda}$  (denoting the load on the system by  $\lambda$ ) and letting  $\lambda \rightarrow \infty$  while keeping  $\beta > 0$  fixed. The system then reaches the desirable Quality-and-Efficiency-Driven (QED) regime, in which the system load approaches 100% while the delays experienced by customers remain limited. In such large-scale service operations, it is natural to use an infinite-server system as a proxy to the many-server system. Infinite-server models are extremely useful because of their tractability; this can be exploited by translating detailed knowledge of the infinite-server system state to the finite-server setting. This returns rather good estimates of future arrivals, even in situations of time-varying arrival processes [16, 17]. The model developed in this paper provides a new way of modeling such large-scale service systems, with the additional feature of a tunable level of overdispersion, essentially replacing a deterministic  $\lambda$  by a stochastically fluctuating  $\Lambda$ . The possibility of using this model, including the above mentioned extensions, for designing dimensioning schemes is currently being investigated by the authors.

## REFERENCES

- [1] D. Anderson, J. Blom, M. Mandjes, H. Thorsdottir, and K. de Turck. A functional central limit theorem for a Markov-modulated infinite-server queue. *Methodology and Computing in Applied Probability*, 2016. [MR3465473](#)
- [2] A. Bassamboo, R.S. Randhawa, and A. Zeevi. Capacity sizing under parameter uncertainty: safety staffing principles revisited. *Management Science*, 56(10):1668–1686, 2010.
- [3] J. Blom, K. de Turck, O. Kella, and M. Mandjes. Tail asymptotics of a Markov-modulated infinite-server queue. *Queueing Systems*, 78(4):337–357, 2014. [MR3269262](#)
- [4] J. Blom, K. de Turck, and M. Mandjes. Analysis of Markov-modulated infinite-server queues in the central-limit regime. *Probability in the Engineering and Informational Sciences*, 29:433–459, 2015. [MR3355613](#)
- [5] D.R. Cox. Some statistical methods connected with series of events journal of the royal statistical society. *Journal of the Royal Statistical Society, Series B (Methodological)*, 17(2):129–164, 1955. [MR0092301](#)
- [6] K. de Turck and M. Mandjes. Large deviations of an infinite-server system with a linearly scaled background process. *Performance Evaluation*, 75–76:36–49, 2014.
- [7] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Springer, New York, 1998. [MR1619036](#)



- [8] S. Kim and W. Whitt. Are call center and hospital arrivals well modeled by non-homogeneous poisson processes? *Manufacturing & Service Operations Management*, 16(3):464–480, 2014.
- [9] S. Kim and W. Whitt. Choosing arrival process models for service systems: Tests of a nonhomogeneous Poisson process. *Naval Research Logistics*, 61(1):66–90, 2014. [MR3162951](#)
- [10] Y. Liu and W. Whitt. Stabilizing customer abandonment in many-server queues with time-varying arrivals. *Operations Research*, 60(6):1551–1564, 2012. [MR3009182](#)
- [11] H. Lu and G. Pang. Gaussian limits for a fork-join network with non-exchangeable synchronization in heavy traffic. *Mathematics of Operations Research*, 2015.
- [12] H. Lu and G. Pang. Heavy-traffic limits for a fork-join network in a renewal random environment. *Submitted*, 2015. [MR3486808](#)
- [13] S. Maman. Uncertainty in the demand of service: the case of call centers and emergency departments. *M. Sc. Thesis, Technion – Israel Institute of Technology, Haifa, Israel*, 2009.
- [14] Ibrahim. R. H. Ye, P. L’Ecuyer, and H. Shen. Modeling and forecasting call center arrivals: a literature survey. *International Journal of Forecasting*, 32:865–874, 2016.
- [15] T. Rydén. An EM algorithm for estimation in Markov-modulated Poisson processes. *Computational Statistics & Data Analysis*, 21:431–447, 1996. [MR1394060](#)
- [16] W. Whitt. Dynamic staffing in a telephone call center aiming to immediately answer all calls. *Operations Research Letters*, 24(5):205–212, 1999. [MR1719916](#)
- [17] W. Whitt, L. V. Green, and P. J. Kolesar. Coping with time-varying demand when setting staffing requirements for a service system. *Production and Operations Management*, 16(1):13–39, 2007.

MARISKA HEEMSKERK  
MICHEL MANDJES  
UNIVERSITY OF AMSTERDAM  
KORTEWEG-DE VRIES INSTITUTE FOR MATHEMATICS  
SCIENCE PARK 107  
1098 XG AMSTERDAM  
THE NETHERLANDS  
E-MAIL: [j.m.a.heemskerk@uva.nl](mailto:j.m.a.heemskerk@uva.nl)  
[m.r.h.mandjes@uva.nl](mailto:m.r.h.mandjes@uva.nl)

JOHAN VAN LEEUWAARDEN  
EINDHOVEN UNIVERSITY OF TECHNOLOGY  
DEPARTMENT OF MATHEMATICS  
AND COMPUTER SCIENCE  
PO BOX 513  
5600 MB EINDHOVEN  
THE NETHERLANDS  
E-MAIL: [j.s.h.v.leeuwaarden@tue.nl](mailto:j.s.h.v.leeuwaarden@tue.nl)