

Priority queues with Gaussian input: a path-space approach to loss and delay asymptotics

Michel Mandjes¹, Petteri Mannersalo² and Ilkka Norros²

¹ Centrum voor Wiskunde en Informatica
P.O. Box 94079, NL-1090 GB Amsterdam, The Netherlands
Michel.Mandjes@cw.nl

² VTT Technical Research Centre of Finland
P.O. Box 12022, FI-02044 VTT, Finland
{Petteri.Mannersalo, Ilkka.Norros}@vtt.fi

Abstract: Priority queueing is the basic technique for providing real-time Quality of Service to packet-based networking. The mathematical analysis of priority queues with general traffic models is, however, prohibitively difficult, in particular when the traffic is long-range dependent. This paper provides some important steps forward in this direction. Our analysis is the first mathematically rigorous treatment of path-space large deviations of priority queues with class-wise heterogeneous Gaussian input having an arbitrary correlation structure. This includes the computation of the most probable paths that lead to overflow in one of the queues. Compared with earlier work on the same topic, the paper provides three novel contributions: a new representation of the workload in the low-priority queues, an exact characterization of the most probable paths, and an extension of the analysis to virtual waiting times, in addition to queue lengths.

Keywords: priority queue, large deviations, Gaussian process, fractional Brownian motion.

1 INTRODUCTION

Priority queueing is the basic technique for providing Quality of Service (QoS) to packet-based networking, which is crucial for real-time services in particular. Advanced scheduling mechanisms are applied in the context of service differentiation, e.g., based on the *Diffserv* principles supported by the IETF (Internet Engineering Task Force). This motivates the need for tractable, yet accurate performance evaluation methods for priority queues that work for a variety of input processes, including long-range dependent traffic. The mathematical analysis of priority queues with such general traffic models is, however, prohibitively difficult. In this paper we present a novel, rigorous methodology for analyzing path-space large deviations of priority queues with general Gaussian input traffic. Though the approach itself is mathematically quite advanced and not necessarily directly applicable to real engineering problems, one of the main implications is that the more tractable performance formulae heuristically derived in [1, 2] can now be motivated rigorously.

Gaussian processes have found wide-spread use as traffic model when considering performance analysis of communication networks. They are attractive because of several nice features; in particular, a Gaussian process A_t with stationary increments is completely characterized by its mean rate $m = \mathbb{E}\{A_1\}$ and cumulative variance function $v(t) = \text{Var}(A_t)$, which can be reasonably well estimated from measurement data. Moreover, Gaussian models of large traffic aggregates can be justified by the Central Limit Theorem [3]. Measurement studies indicate that the Gaussian model is accurate in many situations, as long as the time-scales considered are not too small, and the aggregation level is sufficiently high, see e.g. [4].

The results on path-space large deviations of Gaussian processes were first applied to a queueing context in [5]. The approach was extended to ordinary queues with general Gaussian input in [6, 2], and then, in the form of heuristic approximations and bounds, to priority and Generalized Processor Sharing (GPS) queues in [7, 2]. An essential tool for obtaining exact identification of most probable paths (MPPs) in tandem and priority queues is what could be called the ‘theory of infinite intersections’, as developed in [8]. It was used recently to identify the most likely path to overflow in the tandem queue [9], and the present paper shows how it can be applied in the setting of priority queues. A particular subtlety in our analysis is that the smoothness of the input process plays a dominant role: non-smooth processes such as the highly irregular fractional Brownian motion (fBm) give rise to an essentially different type of MPP than smooth processes such as the integrated Ornstein-Uhlenbeck (iOU) process. We note that other earlier large deviation studies for tandems, priority queues, and GPS predominantly focused on asymptotics for short-range dependent input, see, among many other papers, [10, 11].

This paper is organized as follows. Section 2 presents our definition of a priority queue, which differs from previous definitions, and works very well in the Gaussian case. Section 3 presents preliminaries on the Gaussian large deviation apparatus. In Section 4, we focus, for the two-class setting, on the most probable paths leading to a large queue for the low-priority traffic. In Section 5 we make a corresponding analysis of the virtual waiting time of the second class traffic. Section 6 illustrates the technique with two qualitatively different traffic models: fBm and iOU. Concluding remarks are made in Section 7.

2 THE PRIORITY QUEUE

We use the following notation throughout. For $s < t$, $A_t - A_s$ presents the amount of input traffic in time interval $(s, t]$, and we set $A_0 \equiv 0$. We always assume that the input process is ergodic and has stationary increments, i.e., for any $t_0 \in \mathbb{R}$, the processes $(A_t)_{t \in \mathbb{R}}$ and $(A_{t+t_0} - A_{t_0})_{t \in \mathbb{R}}$ have the same finite-dimensional distributions.

We consider a single queueing system with multiclass traffic and a strict priority queueing discipline. Let the input traffic consist of k classes, and denote the cumulative arrival process of class $j \in \{1, \dots, k\}$ by $(A_t^j)_{t \in \mathbb{R}}$. We also denote $A^j(s, t) \doteq A_t^j - A_s^j$. Consider first the case of a simple queue, i.e., $k = 1$, and let the server have a constant capacity c . The storage process (queue length process) is then defined as $Q_t^1 = \sup_{s \leq t} (A_t^1 - A_s^1 - c(t - s))$. The process Q^1 is obviously stationary, and a sufficient stability condition is that $\mathbb{E}\{A_1^1\} < c$.

Let us then turn to priority queues. Assume that the traffic of class 1 has highest priority, class 2 the second highest, and so on. Our interpretation of the strict priority system is that for any $j = 1, \dots, k - 1$, traffic in classes $1, \dots, j$ is in no way affected by traffic in classes

$j + 1, \dots, k$. Since we are considering fluid type input, there is no distinction between pre-emptive and non-pre-emptive priority. A common way to define the second class queue would be to consider both the total queue, say Q^{12} , and the first class queue, Q^1 as ‘simple’ queues, and then define $Q^2 = Q^{12} - Q^1$. This approach was also applied in the Gaussian context [7, 2], although it was noted that it does not rule out that Q^2 may obtain negative values (albeit with small probability). This unpleasant feature is avoided by the following novel definition (which coincides to the previous one when the processes $(A_t^i)_{t \in \mathbb{R}}$ are non-decreasing).

The idea (for which the authors are grateful to Venkat Anantharam and Takis Konstantopoulos) is to define first the cumulative capacity available for the second class, say C^2 , as a non-decreasing process. Starting from the natural relation $C_t^2 - C_s^2 = (c(t-s) - Q_s^1 - (A_t^1 - A_s^1))^+$ and requiring that C^2 be non-decreasing, it follows that

$$\begin{aligned} C_t^2 - C_s^2 &= \sup_{u \in (s,t]} (c(u-s) - Q_s^1 - (A_u^1 - A_s^1))^+ \\ &= \left(\sup_{u \in (s,t]} (cu - A_u^1) - \sup_{v \leq s} (cv - A_v^1) \right)^+ = \sup_{u \leq t} (cu - A_u^1) - \sup_{v \leq s} (cv - A_v^1). \end{aligned}$$

Thus, we can define simply $C_t^2 = \sup_{s \leq t} (cs - A_s^1)$, which yields a non-decreasing process with stationary increments. We have $C_0^2 = Q_0^1$, which can be subtracted to make the process go through the origin. (Note also that $A_t^1 - ct = C_t^2 - Q_t^1$, i.e., (C^2, Q^1) is the solution of Skorohod’s reflection problem for the process $A_t^1 - ct$.) Our definition of the second class queue length process Q_t^2 is now $Q_t^2 = \sup_{s \leq t} (A_t^2 - A_s^2 - (C_t^2 - C_s^2))$. In general, we define inductively

$$C_t^j \doteq \sup_{s \leq t} (C_s^{j-1} - A_s^{j-1}), \quad Q_t^j \doteq \sup_{s \leq t} (A_t^j - A_s^j - (C_t^j - C_s^j)), \quad j = 2, \dots, k.$$

It is crucial for our large deviation analysis that the event $\{Q_0^2 > x\}$ can be written in a union-intersection form as follows. Since $C_0^2 = \sup_{v \leq 0} (v - A_v^1) \vee \sup_{u \in [s,0]} (u - A_u^1)$ for $s \leq 0$, we have $C_0^2 - C_s^2 = [\sup_{u \in [s,0]} (u - A_u^1) - C_s^2]^+$. Now,

$$\begin{aligned} \{Q_0^2 \geq x\} &= \bigcup_{s \leq 0} \{[\sup_{u \in [s,0]} (cu - A_u^1) - C_s^2]^+ \leq -A_s^2 - x\} \\ &= \bigcup_{s \leq 0} (\{ \sup_{u \in [s,0]} (cu - A_u^1) - C_s^2 \leq -A_s^2 - x\} \cap \{-A_s^2 - x \geq 0\}) \\ &= \bigcup_{s \leq 0} (\{ \sup_{u \in [s,0]} (cu - A_u^1) \leq -A_s^2 - x + \sup_{t \leq s} (ct - A_t^1) \} \cap \{-A_s^2 - x \geq 0\}) \\ &= \bigcup_{s \leq 0} \bigcup_{t \leq s} (\bigcap_{u \in [s,0]} \{A_u^1 - A_t^1 \geq c(u-t) + A_s^2 + x\} \cap \{-A_s^2 - x \geq 0\}) \\ &= \bigcup_{s \leq 0} \bigcup_{t \leq s} \bigcup_{a \geq x} (\bigcap_{u \in [s,0]} \{A_u^1 - A_t^1 \geq c(u-t) + x - a\} \cap \{-A_s^2 = a\}). \end{aligned} \quad (1)$$

Similarly, an expression for $\{Q_0^3 \geq x\}$, etc., can be found.

In addition to queue lengths, we are also interested in the delays. We give the following intuitive definition of the virtual waiting time (v.w.t.) of class j at time t , denoted as V_t^j : we say that $V_t^j = \tau$ if a ‘fluid molecule’ of class j entering the system at time t , is transmitted at time $t + \tau$. In a fluid model with fixed service rate, the v.w.t. of a simple queue is proportional to the

queue length. Evidently, the class 1 queue is simply $V_t^1 \doteq Q_t^1/c$. For $j > 1$, the natural definition of V_t^j is the time remaining to the next time $t + \tau$ such that all work in the queues $1, \dots, j$ at time t has been served and, additionally, all work from the classes $1, \dots, j - 1$ that arrived during $(t, t + \tau]$ has been served, too:

$$V_t^j \doteq \tau_t^j - t, \quad \tau_t^j = \inf \left\{ u > t : \sum_{i=1}^j Q_t^i + \sum_{i=1}^{j-1} (A_u^i - A_t^i) < c(u - t) \right\}. \tag{2}$$

In the rest of this paper, we restrict ourselves to the case of two priority classes.

3 THE LARGE DEVIATION FRAMEWORK FOR GAUSSIAN SYSTEMS

3.1 The Gaussian traffic model

We assume that the processes A^j are independent, continuous Gaussian processes with stationary increments and denote

$$A_t^j = m_j t + Z_t^j, \quad m = \sum_{i=1}^k m_i, \quad \text{Var} \left(Z_t^j \right) = v_j(t), \quad \Gamma_j(s, t) = \text{Cov} \left(Z_s^j, Z_t^j \right),$$

where the Z^j 's are centered (zero-mean) processes. To exclude pathological cases, we assume that there exist numbers $\alpha_0, \alpha_\infty \in (0, 2]$ such that $v(h)/h^{\alpha_0}$ is bounded for $h \in (0, 1)$, and $\lim_{t \rightarrow \infty} v(t)/t^{\alpha_\infty} = 0$. We also assume that all finite-dimensional distributions of the process (Z^1, \dots, Z^k) are non-singular.

We call a Gaussian process Z *smooth* at t , if it has a Bochner derivative at 0, that is, there exists a random variable $Z_t' \in G$ such that $\lim_{h \rightarrow 0} \mathbb{E} \{ (Z_t' - (Z_{t+h} - Z_t)/h)^2 \} = 0$. It follows from the stationarity of increments that if Z is smooth at 0, then it is smooth at all $t \in \mathbb{R}$. It takes some straightforward analysis to verify that fBm, with variance function t^{2H} for some Hurst parameter $H \in (0, 1)$, is non-smooth, as opposed to iOU, with variance function $t - 1 + e^{-t}$ (in fact, Z_t' is a normally distributed with mean 0 and variance $\frac{1}{2}$).

The reproducing kernel Hilbert space (RKHS) of a Gaussian process plays a crucial role in the large deviation asymptotics. For $i = 1, \dots, k$, the RKHS R_i of Z^i is defined as follows: start with the functions $\Gamma_i(t, \cdot)$, $t \in \mathbb{R}$, define their inner product as

$$\langle \Gamma_i(s, \cdot), \Gamma_i(t, \cdot) \rangle_{R_i} \doteq \Gamma_i(s, t),$$

extend to a linear space (with pointwise operations), and complete the space with respect to the norm $\|f\|_{R_i} \doteq \langle f, f \rangle_{R_i}$. The RKHS of the multivariate process (Z^1, \dots, Z^k) can be defined as space $R \doteq R_1 \times \dots \times R_k$ with inner product $\langle (f_1, \dots, f_k), (g_1, \dots, g_k) \rangle_R \doteq \sum_{i=1}^k \langle f_i, g_i \rangle_{R_i}$. The correspondence $Z_t^i \leftrightarrow \Gamma_i(t, \cdot)$ extends to a Hilbert space isometry between a subspace of L^2 and R_i . If Z^i is smooth, then $\Gamma_i(s, t)$ has partial derivatives, and the isometry counterpart of $Z_t^{i'}$ in R is the function $\Gamma_i'(t, s) \doteq (d/dt)\Gamma(t, s)$.

The ‘message’ of this theorem is that the minimizing path is determined through the points S^* where it touches the curve ζ , plus, in the case that the process is smooth, the infinitesimal environments of those points.

4 MOST PROBABLE PATHS WITH A LARGE VALUE OF Q^2

To make a large-deviations analysis of the overflow event $\{Q_0^2 \geq x\}$ accessible to our methods, we should write it in terms of half-spaces in such a way that unions precede intersections: i.e., if $E = \cup_t \cap_s E_{s,t}$, then $\inf_{f \in E} \|f\| = \inf_t \inf_{f \in \cap_s E_{s,t}} \|f\|$, and it turns out that $\inf_{f \in \cap_s E_{s,t}} \|f\|$ can be tackled by applying Theorem 1. Indeed, the form (1) is exactly of this type:

$$\{Q_0^2 \geq x\} = \bigcup_{s \leq 0} \bigcup_{t \leq s} \bigcup_{a \geq x} (B_{x,s,t,a}^1 \cap B_{s,a}^2), \quad (4)$$

where $B_{x,s,t,a}^1 \doteq \{A_u^1 - A_t^1 \geq c(u-t) + x - a \forall u \in [s, 0]\}$ and $B_{s,a}^2 \doteq \{-A_s^2 = a\}$.

The above form tells us already a lot about the most probable paths in $\{Q_0^2 \geq x\}$. First, if we can characterize the most probable path pairs in the events $B_{x,s,t,a}^1 \cap B_{s,a}^2$ and compute their norms, then the rest is just a minimization of a known numerical function (over s, t, a), as argued above. Second, since each event $B_{x,s,t,a}^1 \cap B_{s,a}^2$ is convex and closed and its intersection with the kernel space R is non-empty, it has a unique MPP pair. Third, the independence of the priority classes yields that the MPP pair consists of the MPP s of each of the sets $B_{x,s,t,a}^1$ and $B_{s,a}^2$ taken separately. In particular, as remarked earlier, the MPP in sets of the type $B_{s,a}^2$ is always a multiple of a single covariance function $\Gamma_2(s, \cdot)$.

Let us consider the event $B_{x,s,t,a}^1$ with fixed parameters $t \leq s \leq 0$, $a \geq x$. By the stationarity of the increments, the event $B_{x,s,t,a}^1$ is stochastically equivalent to

$$\tilde{B}_{x,s,t,a}^1 = \{Z_u^1 \geq (c - m_1)u + x - a \forall u \in [s-t, -t]\}.$$

This in turn is very close to the form appearing in Theorem 1. The only difference is that the function $u \mapsto (c - m_1)u + x - a$ does not go through the origin when $a > x$, but this modification is easily seen to be nonessential. Indeed, it suffices to assume that for each $h \in (0, -t)$ there exists a function $\zeta_h \in R$ such that $\zeta_h(u) = (c - m_1)u + x - a$ for $u \in [h, -t]$, which is more or less always true (for all ‘usual’ processes at least). For the case $s = t$, choose h so small that the most probable path in $\{Z_u \geq (c - m_1)u + x - a \forall u \in [h, -t]\}$ is above the line $u \mapsto (c - m_1)u + x - a$ for all $u \in [0, -t]$.

The case that $\zeta_h(-t) \leq 0$ is not interesting, because the corresponding most probable path of Z^1 would be identically zero. Assume then that $\zeta_h(-t) > 0$, and let β^* be the most probable path in $\tilde{B}_{x,s,t,a}^1$. Now there are two possibilities.

- **Case 1.** The MPP is $\beta^* = \frac{a-x-m_1|t|}{v_1(t)} \Gamma_1(|t|, \cdot)$ if

$$\frac{a-x-m_1|t|}{v_1(t)} \Gamma_1(|t|, u) \geq (c - m_1)u + x - a \quad \text{for } u \in [s-t, -t],$$

- **Case 2.** In the remaining case, we can, in principle, determine the MPP β^* with the methods developed in [8], see Theorem 1.

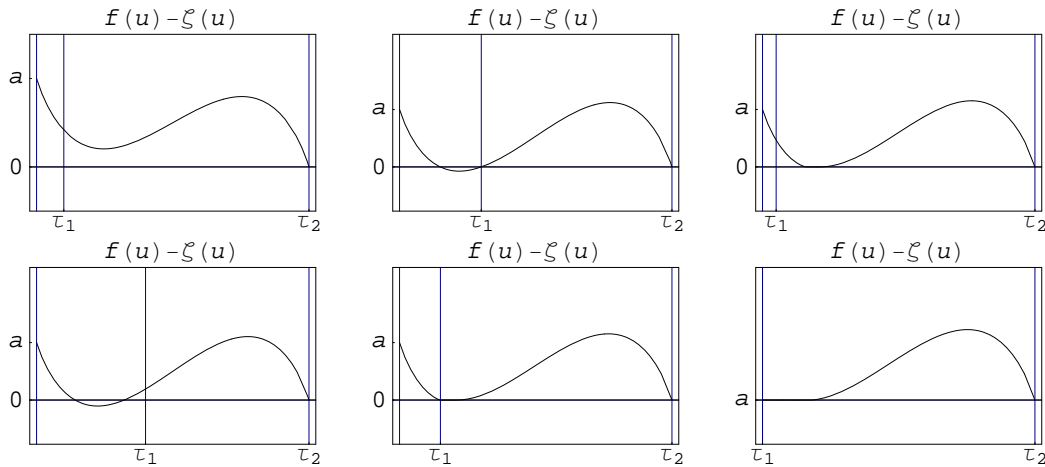


Fig. 1. The shapes of the most probable paths for fractional Brownian motion ($H = 0.8$) in the set $\{Z_u - \zeta(u) \geq 0, \forall u \in [\tau_1, \tau_2]\}$.

where we recall that $A^i(t_1, t_2) = A_{t_2}^i - A_{t_1}^i$. Hence, analogously to (1), we have

$$\{V_0^2 \geq t\} = \bigcup_{v \leq s \leq 0} \bigcup_{a \in \mathbb{R}} (D_{t,s,v,a}^1 \cap D_{s,a}^2), \tag{5}$$

where $D_{t,s,v,a}^1 \doteq \{A^1(v, u) \geq -a + c(u - v) \forall u \in [0, t]\}$ and $D_{s,a}^2 \doteq B_{s,a}^2 = \{A^2(s, 0) = a\}$. After fixing s, t, v, a with $v < s \leq 0 \leq t$, the stochastically equivalent problem is to study the set $\tilde{D}_{t,s,v,a}^1 \doteq \{Z_u^1 \geq (c - m_1)u - a, \forall u \in [-v, -v + t]\}$. Hence, we have again reduced our problem to an event in which the unions precede intersections, and, as a consequence, we can again invoke Theorem 1.

6 EXAMPLES

In our setting, the most probable path of the second class traffic is always a multiple of $\Gamma(s, \cdot)$. Thus we restrict ourselves to finding the most probable path either in $\tilde{B}_{s,t,a}^1$ or in $\tilde{D}_{t,s,v,a}^1$. Define for $0 \leq \tau_1 \leq \tau_2$, $U_{\tau_1, \tau_2} \doteq \{f : f(u) \geq \zeta(u), u \in [\tau_1, \tau_2]\}$, where $\zeta(u) \doteq -a + bu$ with $a \geq 0$ and $b > 0$. It is easy to see that both $\tilde{B}_{s,t,a}$ and $\tilde{D}_{t,s,v,a}$ can be written in this form. (In the delay setting, it is possible that $a < 0$, corresponding to negative A_s^2 . For positively correlated processes, this situation can be ruled out by the ‘waste of energy’ argument.)

With $\beta^* = \operatorname{argmin} \{\|f\| : f \in U_{\tau_1, \tau_2}\}$ and $\beta_S(\cdot) \doteq \mathbb{E}[Z_s | Z_s = \zeta(s) \forall s \in S]$, Theorem 1 states that the MPP for non-smooth processes (such as fBm) is always of the form $\beta_{S^*}^*$ with some $S^* \subset [\tau_1, \tau_2]$. For smooth processes (like iOU), we sometimes need conditions on derivatives.

In the following, we show the MPP s for fBm and iOU. The case $a = \tau_1 = 0$ corresponds to the busy period problem which was studied in [8, 5]. For other parameter values, the proofs of [9] hold with minor changes. Once these paths are known, the corresponding norms are determined; then calculation of the rate function is a straightforward, though computationally involved, minimization problem in 3 dimensions.

For fractional Brownian motion, essentially three different shapes are possible:

- (i) If requiring that $f(\tau_2) = \zeta(\tau_2)$ leads to a feasible path, we are done, i.e., $\beta^* = \beta_{\{\tau_2\}}$. The left column in Fig. 1.

REFERENCES

1. P. Mannersalo and I. Norros, Approximate formulae for Gaussian priority queues, Proceedings of ITC 17, Salvador, Brazil, pp. 991–1002, 2001.
2. P. Mannersalo and I. Norros, A most probable path approach to queueing systems with general Gaussian input, *Computer Networks*, 40 (2002) 399.
3. R. Addie, On weak convergence of long range dependent traffic processes, *Journal of Statistical Planning and Inference*, 80 (1999) 155.
4. J. Kilpi and I. Norros, Testing the Gaussian approximation of aggregate traffic, Proceedings of The 2nd Internet Measurement Workshop, Marseille, France, pp. 49–61, 2002.
5. I. Norros, Busy periods of fractional Brownian storage: a large deviations approach, *Adv. Perf. Anal.*, 2 (1999) 1.
6. R. Addie, P. Mannersalo and I. Norros, Performance formulae for queues with Gaussian input. Proceedings of ITC 16, Edinburgh, UK, pp. 1169–1178, 1999.
7. M. Mandjes and M. van Uitert, Sample path large deviations for tandem and priority queues with Gaussian inputs. *Annals of Applied Probability*, 15 (2005) 1193.
8. M. Mandjes, P. Mannersalo, I. Norros and M. van Uitert, Large deviations of infinite intersections of events in Gaussian processes, Technical Report PNA-E0409, CWI, 2004, <http://www.cwi.nl/>.
9. M. Mandjes, P. Mannersalo and I. Norros, Large deviations of Gaussian tandem queues and resulting performance formulae, Technical Report PNA-E0412, CWI, 2004, <http://www.cwi.nl/>.
10. D. Bertsimas, I. Paschalidis and J. Tsitsiklis, Large deviations analysis of the Generalized Processor Sharing policy, *Queueing Systems*, 32 (1999) 319.
11. C. Chang, P. Heidelberger, S. Juneja and P. Shahabuddin, Effective bandwidth and fast simulation of ATMintree networks, *Performance Evaluation*, 20 (1994) 45.
12. R. Bahadur and S. Zabell, Large deviations of the sample mean in general vector spaces, *Ann. Prob.*, 7 (1979) 587.
13. A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, Springer, New York, 1998.
14. A. Berger and W. Whitt, Effective bandwidths with priorities, *IEEE/ACM Transactions on Networking*, 6 (1998) 447.
15. A. Berger and W. Whitt, Extending the effective bandwidth concept to networks with priority classes, *IEEE Communications Magazine*, August (1998) 78.