

An Analytical Solution for Probabilistic Guarantees of Reservation Based Soft Real-Time Systems

Luigi Palopoli¹, Daniele Fontanelli², Luca Abeni¹ Bernardo Villalba Frías¹

¹Dipartimento di Scienza e Ingegneria dell'Informazione

²Dipartimento di Ingegneria Industriale

University of Trento, Trento, Italy

{luigi.palopoli,daniele.fontanelli,luca.abeni,br.villalbafrías}@unitn.it

Abstract—We show a methodology for the computation of the probability of deadline miss for a periodic real-time task scheduled by a resource reservation algorithm. We propose a modelling technique for the system that reduces the computation of such a probability to that of the steady state probability of an infinite state Discrete Time Markov Chain with a periodic structure. This structure is exploited to develop an efficient numeric solution where different accuracy/computation time trade-offs can be obtained by operating on the granularity of the model. More importantly we offer a closed form conservative bound for the probability of a deadline miss. Our experiments reveal that the bound remains reasonably close to the experimental probability in one real-time application of practical interest. When this bound is used for the optimisation of the overall Quality of Service for a set of tasks sharing the CPU, it produces a good sub-optimal solution in a small amount of time.

Index Terms—Real-time systems, Scheduling, Probabilistic Guarantees

I. INTRODUCTION

The term *soft real-time* is used for a class of real-time applications that are resilient to occasional and controlled timing faults. Significant examples include multimedia streaming [1], computer vision and real-time control [2], [3].

An effective method to express the timing requirements for a soft real-time application is by associating each deadline with a probability that it will be met: the notion of *probabilistic deadlines* [4]. Probabilistic deadlines can be related to the Quality of Service (QoS) delivered by the application [5], [2] and, more generally, enable the expression of a wide range of performance requirements, where classic hard real-time systems can be regarded as a special case.

In traditional hard real-time applications, the use of fixed or dynamic scheduling priorities has gained an undisputed prominence. Part of the reasons of this success is in the presence of efficient numeric techniques that make for the provision of tight conditions for temporal guarantees [6]. At least as important is a group of approximate analytical results. The most famous is the utilisation bound [7], which offers

clear guidelines on how to tweak periods and computation times in order to meet the deadlines of all tasks in the system.

The use of scheduling priorities allows the designer to define a partial order between all the tasks in a set and inevitably couples their timing behaviour. This is acceptable if the purpose is to offer guarantees for the set as a whole. On the contrary, if the designer requires specific QoS levels for each task, scheduling priorities can be too coarse a tool. For this reason an intense research work has produced alternative scheduling solutions for soft real-time systems. One of the most popular is the *Resource Reservations* scheduling (RR) [8], [1], which enables a fine grained control on the fraction of computing power (bandwidth) that each task receives. A key property of RR scheduling is *temporal isolation*: the ability for a task to meet its deadlines solely depends on its computation requirement and on its scheduling parameters. This property enables the provision of specific temporal guarantees to each task and simplifies system design. RR scheduling is now available in the mainstream Linux Kernel¹.

When the probability distribution of inter-arrival time and of computation time are known independent identically distributed (i.i.d.) stochastic processes, temporal isolation allows modelling the evolution of a task scheduled through a RR as a Discrete-Time Markov Chain (DTMC) with an infinite number of states [4], [9]. In this paper, we restrict the focus to the analysis of periodic tasks. For this case, we can see that the DTMC describing the system takes the form of a Quasi-Birth-Death Process (QBDP) [10]. We introduce a granularity parameter that allows us to reduce the complexity of the model at the expense of a conservative approximation in the computation of the probability. We show a novel analysis that exploits the specific structure of the transition matrix of this QBDP. The outcome is an expression for the steady state probability of meeting the deadline, which can be used in different ways. The first one is for the construction of a numeric algorithm for probabilistic guarantees, with a performance comparable to the best state of the art techniques for numeric solutions of QBDP. The second one, the most important, is for the computation of an analytical conservative bound for the probability of meeting the deadline. This bound proves itself

The research leading to these results has received funding from the European Union FP7 Programme (FP7/2007-2013) under grant agreement n° ICT-2011-288917 “DALi - Devices for Assisted Living” and under grant agreement n° FP7-ICT-257462 “HYCON2 NoE”, and from the European Union H2020 programme under grant agreement n° 643544 “ACANTO”

¹<https://www.kernel.org/doc/Documentation/scheduler/sched-deadline.txt>

reasonably accurate for a large set of synthetic test cases. We have also performed a large collection of experimental data for a real-life application, in which the presence of several non-idealities (OS overhead, correlation in the computation times, etc.) challenges the assumptions the method relies on. The small approximation error that we observed in the experiments suggests the practical applicability of the method at least in the considered scenario. The application of the bound is very convenient when solving QoS optimisation problems that require to efficiently identify the minimum bandwidth required for a desired probability of deadline miss. We show a realistic example of this kind where the application of the analytic bound produces a good sub-optimal solution in a tiny fraction of the time required by a numeric approach.

The paper is organised as follows. In Section II, we offer a brief survey of the related work. In Section III we formally describe the problem addressed in the paper. In Section IV, we show how a resource reservation can be conservatively modelled as a QBDP. The computation of our analytical bound is reported in Section V. In Section VI, we prove the validity of the bound in a large set of experiments. In Section VII, we show the concrete application of the method to a QoS optimisation problem. Finally, in Section VIII we offer our conclusions and announce the future work directions.

II. RELATED WORK

The stochastic analysis of performance of soft real-time tasks started two decades ago. The same task model presented in this paper (a triple of period, probability distribution of the task computation time and requested probability of deadline miss in the long run) has been also adopted in the statistical rate monotonic approach [11]. More recently, an important number of research papers has concentrated on the computation of the response time of systems with fixed or dynamic priority when tasks have stochastic variability in computation times [12], [13], [14], in the inter-arrival time [15] or in both [16]. Similar techniques have recently been applied to multiprocessor systems [17]. An obvious point of differentiation between our technique and the ones describes so far is that while these papers propose numeric techniques, we offer an analytic bound that is satisfactorily tight in many cases of interest. A very interesting connection can be established with the work of Diaz et al. [12], where the authors propose the exact solution for a specific numeric example. Our computation, on the contrary, applies to general cases. What is more, all the approaches mentioned above analyse the task set as a whole, since real-time schedulers do not enjoy temporal isolation. This makes QoS optimisation much more difficult than in our case.

Other authors have analysed scheduling approaches other than “traditional” fixed or dynamic priorities. Dong-In et al. [18] have analysed Time Division Multiple Access (TDMA) approaches, Haman et al. [19] have focused on a model where tasks are split in mandatory and optional parts. This paper is based on reservation-based scheduling [1], [8], which allows us to exploit temporal isolation and analyse each

task separately. Abeni and Buttazzo proposed a model for RR scheduling based on queueing theory [4], [9]. The computation of the deadline miss probability requires to numerically solve an eigenvector problem for an infinitely large matrix. Recently, approximated solution techniques have been proposed for efficient numeric computation of a bound for the probability of meeting the deadline [20].

In this paper, we show how the adoption of the reservation scheduler and the restriction to periodic tasks produces a model that is a particular instance of a QBDP. Efficient numeric solutions for QBDP and for M/G/1 queue can be found in the work of Latouche and Ramaswami [21] and of Neuts [22], who pioneered the application of matrix geometric methods for the solution of infinite M/G/1 queues. The literature in the field is rich of optimised methods derived using specific properties of the transition matrix. The most remarkable achievements are summarised in a comprehensive book [10]. In this paper, we consider numeric methods as a basis for comparison but our main focus is on analytical closed form solutions.

Mills and Anderson [23] have recently considered the problem of stochastic analysis for resource reservations on multiprocessor systems. The authors main focus is on the computation of tardiness and response time bounds for the average case. The authors also offer a very conservative result on the probabilistic deadlines, which is applicable only if deadlines much larger than the period are considered.

A customary assumption made in the literature on queueing networks is that inter-arrival times and service times are i.i.d. processes. In this paper, we stick to the same assumption. Different authors have recently questioned on the applicability of the i.i.d. assumption in the area of real-time applications [24]. Remarkable is the so called notion of probabilistic worst case execution time [25], which essentially corresponds to associating a worst case to several execution scenarios that take place within a given probability. A possible evolution of this concept could lead to finding an i.i.d. overapproximation for a computation process that is not i.i.d. A similar idea underpins a recent work by Liu et al. [26], where the authors tackle the correlation problem decomposing the process into a deterministic and an i.i.d. component. In a similar context our results could be used to study the evolution of the system under the action of the i.i.d. component or of the i.i.d. overapproximation of the process.

A complementary issue to our work is how to derive statistically sound estimates for the probability distribution of the computation time. A useful inspiration could come from the application of the Extreme Value Theory [27], but the matter is reserved for future investigations.

The results shown in this paper take to its natural completion a line of work started a few years ago that has produced a number of intermediate results. The relation with our prior achievements is detailed in Section VI-C.

III. PROBLEM DESCRIPTION

A. Task Model

We consider a set of real-time tasks $\{\tau_i\}$ sharing a *processing unit* (CPU). A real-time task τ_i consists of a stream of jobs $J_{i,k}$. Each job $J_{i,k}$ arrives (becomes eligible for execution) at time $r_{i,k}$, and finishes at time $f_{i,k}$ after executing for a time $c_{i,k}$. We restrict to periodic tasks, meaning that two adjacent arrivals are spaced out by a fixed amount of time T_i : $r_{i,k+1} = r_{i,k} + T_i$.

The computation time of each job $c_{i,k}$ is assumed to be an i.i.d. stochastic process \mathcal{U}_i . For each k the computation time is a random variable described by the Probability Mass Function (PMF) $U_i(c) = \Pr\{c_{i,k} = c\}$.

Job $J_{i,k}$ is associated with a deadline $d_{i,k} = r_{i,k} + D_i$ (where D_i is said relative deadline), that is respected if $f_{i,k} \leq d_{i,k}$, and is missed if $f_{i,k} > d_{i,k}$. In this work, *probabilistic deadlines* [4] are used instead of traditional hard deadlines $d_{i,k}$. A probabilistic deadline (D_i, p_i) is respected if $\Pr\{f_{i,k} > r_{i,k} + D_i\} \leq p_i$. If $p_i = 0$ the deadline is hard.

B. The scheduling algorithm

As multiple real-time tasks may be concurrently active, we use a RR scheduler. Each task τ_i is associated with a reservation (Q_i^s, T_i^s) , meaning that τ_i is allowed to execute for Q_i^s (*budget*) time units in every interval of length T_i^s (*reservation period*). The fraction of CPU allocated to the task is said bandwidth B_i and is defined as $B_i = Q_i^s/T_i^s$. The particular implementation of the RR approach that we consider is the Constant Bandwidth Server (CBS) [1]. In the CBS, reservations are implemented by means of an Earliest Deadline First (EDF) scheduler. The EDF schedules tasks $\{\tau_i\}$ based on their *scheduling deadlines* $d_{i,k}^s$, which are dynamically managed by the CBS algorithm. When a new job $J_{i,k}$ arrives, the server checks whether it can be scheduled using the last assigned scheduling deadline $d_{i,k-1}^s$. In the affirmative case, the scheduling deadline of the job is initially set to current deadline $d_{i,k}^s = d_{i,k-1}^s$. Otherwise, the initial deadline $d_{i,k}^s$ is set equal to $r_{i,k} + T_i^s$. Every time the job executes for Q_i^s time units (i.e., its budget is depleted), its scheduling deadline is postponed by T_i^s : $d_{i,k}^s = d_{i,k}^s + T_i^s$. This way, the task is prevented from executing for more than Q_i^s units with the same deadline. As a consequence, each task is reserved an amount of computation time Q_i^s in each server period T_i^s regardless of the behaviour of the other tasks. This property is called *temporal isolation* and it holds as long as the system satisfies the following *schedulability condition*:

$$\sum_i B_i = \sum_i \frac{Q_i^s}{T_i^s} \leq 1. \quad (1)$$

The scheduling deadline $d_{i,k}^s$ has, in general, nothing to do with the deadline $d_{i,k}$ of the job: it is simply instrumental to the implementation of the CBS (see [1] for more details).

C. Problem Statement

In view of the temporal isolation property, each task is guaranteed a minimum share of the processor Q_i^s/T_i^s independently of the behaviour of the other tasks. As a consequence, it is possible to carry out a conservative analysis leading to the computation of a lower bound of the probability of respecting a deadline assuming that the task always receives this minimum (as long as Condition (1) is respected). The advantage is that the behaviour of each task can be studied in isolation. Therefore, we can remove the subscript i meaning that the analysis refers to one specific task.

In this setting, our problem is formulated as follows.

Problem 1: Given a periodic real-time task with a stochastic computation time characterised by a PMF $U(c)$, find conditions on the reservation parameters (Q^s, T^s) such that the task respects the probabilistic deadline (D, p) .

A few remarks are in order. First of all, we look for analytical conditions, which can be inverted and offer easy solution for the problem of system design. Second, in order to be safely utilisable, such conditions have to be *sufficient* (although necessity is certainly a desirable additional requirement).

IV. STOCHASTIC MODEL

In this section, we first recall some basic definitions on Markov chains and in particular on QBDP. Then, we show how a task scheduled by a resource reservation is conveniently modelled as a QBDP (Theorem 1). Finally, we show how to derive a conservative approximation of this model, which has a parametric accuracy and which retains the structure of a QBDP.

A. Background on Markov Chains

A *Discrete-Time Markov Process* (DTMP) $\{X_n\}$ is a discrete-time stochastic process such that its future development only depends on the current state and not on the past history. This can be stated in formal terms on the conditional PMF: $\Pr\{X_n = x_n | X_1 = x_1, X_2 = x_2, \dots, X_{n-1} = x_{n-1}\} = \Pr\{X_n = x_n | X_{n-1} = x_{n-1}\}$. A DTMP defined over a discrete state space is said *Discrete-Time Markov chain* (DTMC). Given a DTMC, let $\pi_n^{(j)}$ represent the probability $\pi_n^{(j)} = \Pr\{X_n = j\}$, π_n be the vector $\pi_n = [\pi_n^{(0)}, \pi_n^{(1)}, \dots]$, $P = [p_{i,j}]$ be a matrix whose generic element $p_{i,j}$ is given by the conditional probability $p_{i,j} = \Pr\{X_n = j | X_{n-1} = i\}$. Starting from an initial probability distribution π_0 , the application of the Bayes theorem and of the properties of the Markov Processes allow us to express the evolution of the distribution by the matrix equation $\pi_{n+1} = \pi_n P$. The matrix P is said probability transition matrix. An *equilibrium point* for this dynamic equation is a vector $\tilde{\pi}$ such that $\tilde{\pi} = \tilde{\pi} P$.

Consider a state i of a DTMC. Let the random variable $\mathcal{T}_i = \min\{n > 1 \text{ s.t. } X_n = i | X_0 = i\}$ denote the first return time to state i . The state i is transient if $\Pr\{\mathcal{T}_i < \infty\} < 1$, i.e., if there is some probability that starting from i the state will never return to i . The state i is *transient* if it is not

recurrent. The *period* d_i of a recurrent state i is defined as the greatest common divisor of the set of all numbers, n , for which $\Pr\{X_m = i \wedge X_{m+n} = i\} > 0, \forall m$. A state is said *aperiodic* if its period $d_i = 1$. A DTMC is said aperiodic, if all of its states are aperiodic.

The mean recurrence time of a state i is the expected value of \mathcal{T}_i : $M_i = \mathbf{E}\{\mathcal{T}_i\}$. The state i is positive recurrent if M_i is finite, and the DTMC is positive recurrent if all its states are positive recurrent.

A DTMC is said *irreducible*, if every state can be reached from any other state in a finite number of steps. It can be shown that in an irreducible DTMC all states are of the same type. So, if one state is aperiodic, so is the DTMC.

A very important property of irreducible and positive recurrent DTMC is *the existence of a single equilibrium* $\tilde{\pi} = \tilde{\pi}P$ where the limiting distributions $\lim_{n \rightarrow \infty} \pi_n$ converge starting from any initial probability distribution π_0 . This equilibrium is called *steady state distribution*.

A DTMC is called a Quasi-Birth-Death Process (QBDP) if its probability transition matrix P has the following block structure:

$$P = \begin{bmatrix} C & A_0 & 0 & 0 & 0 & \cdots \\ A_2 & A_1 & A_0 & 0 & 0 & \cdots \\ 0 & A_2 & A_1 & A_0 & 0 & \cdots \\ 0 & 0 & A_2 & A_1 & A_0 & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \end{bmatrix} \quad (2)$$

When the matrices are scalars, this structure reduces to the standard Birth-Death Process (BDP).

B. A resource reservation as a Markov Chain

We will denote by $F_U(c) = \sum_{h=c_{min}}^c U(h)$ the Cumulative Distribution Function (CDF) of the execution time. For simplicity, we will assume that the server period T^s is chosen as an integer sub-multiple of the activation period T : $T = NT^s$. Other choices are possible but make little practical sense.

Let d_k^s denote the latest scheduling deadline used for job J_k and introduce the symbol $\delta_k = d_k^s - r_k$. The latest scheduling deadline d_k^s is an upper bound for the finishing time of the job (if Equation (1) is respected, then $f_k \leq d_k^s$). Hence, δ_k is an upper bound for the job response time.

Example 1: Consider the schedule in Figure 1. The schedule in the figure considers two adjacent jobs starting at r_k and r_{k+1} and the reservation period is chosen as one third of the task period. Job J_k , in this case finishes beyond the deadline (which in our periodic model is r_{k+1}). More precisely, the last reservation period that it uses (in which its finishing time lies) is upper-limited by the scheduling deadline d_k^s .

The quantity δ_k takes on values in a discrete set: the integer multiples of T^s and the probability p of meeting the deadline is lower bounded by $\Pr\{\delta_k \leq D\}$.

The evolution of δ_k is described as follows [9]:

$$\begin{aligned} v_0 &= c_0 \\ v_{k+1} &= \max\{0, v_k - NQ^s\} + c_{k+1} \\ \delta_k &= \left\lceil \frac{v_k}{Q^s} \right\rceil T^s \end{aligned} \quad (3)$$

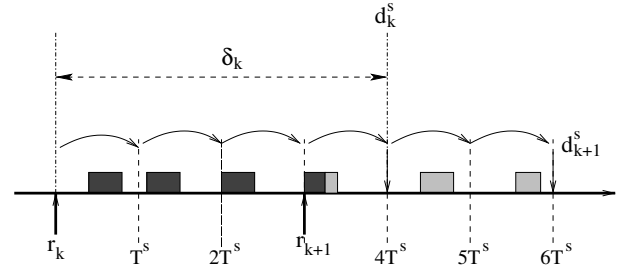


Figure 1. Example schedule of a task by a CBS. The two colours denote different jobs.

The variable v_k cannot be measured directly and it represents the amount of backlogged execution time that has to be served by the CBS scheduler when a new job arrives.

Since the process \mathcal{U} modelling the sequence c_k of the computation time is assumed a discrete valued and i.i.d. random process, the model in Equation (3) represents a Discrete-Time Markov Chain (DTMC) that we define \mathcal{M}_0 , where the states are determined by the possible values of v_k and the transition probabilities by the PMF of the computation time $U(c)$.

This model permits a fine-grained modelling of the behaviour of the reservation, which can be difficult to treat. One possible simplification is to collapse into a single state all the states for which $\delta_k \leq D = NT^s$, which correspond to the values of v_k such that $v_k \leq NQ^s$. In the modified DTMC \mathcal{M} , the state S is defined as

$$S = \begin{cases} 0 & \text{if } v_k \leq NQ^s \\ i & \text{if } v_k = NQ^s + i \end{cases}.$$

By using Equation (3), the transition probabilities for this DTMC can be written as follows:

$$\begin{aligned} p_{i,j} &= \begin{cases} \Pr\{v_{k+1} \leq NQ^s | v_k = i + NQ^s\}, & \text{if } j = 0 \\ \Pr\{v_{k+1} = j + NQ^s | v_k \leq NQ^s\}, & \text{if } i = 0, j \neq 0 \\ \Pr\{v_{k+1} = NQ^s + j | v_k = i + NQ^s\}, & \text{if } i \neq 0, j \neq 0 \end{cases} \\ &= \begin{cases} \Pr\{c_k \leq NQ^s - i\} = F_U(NQ^s - i), & \text{if } j = 0 \\ \Pr\{c_k = j + NQ^s\} = U(j + NQ^s), & \text{if } i = 0, j \neq 0 \\ \Pr\{c_k = NQ^s + j - i\} = U(j - i + NQ^s), & \text{if } i \neq 0, j \neq 0. \end{cases} \end{aligned}$$

Let $\tilde{\pi}_k$ be the (infinite) vector where the i^{th} element represent the probability associated with the i^{th} state of the DTMC \mathcal{M} after k step of evolution starting from an initial probability vector $\tilde{\pi}_0$. The recursive equation for the evolution of $\tilde{\pi}_k$ is $\tilde{\pi}_{k+1} = \tilde{\pi}_k P$. The objective of our analysis can now be stated as *the computation of a lower bound for the first element of the steady state probability vector* $\tilde{\pi} = \lim_{k \rightarrow \infty} \pi_k$. As long as we are not interested in the distribution of δ_k inside the region $\delta_k \leq NQ^s$, collapsing into one state all the values of v_k smaller than NQ^s does not introduce any error because such states do not have influence on the next state ($\max\{0, v_k - NQ^s\} = 0$ in Equation (3)).

The probability matrix P resulting from the computation above has the structure shown in Figure 2, where

$$\begin{aligned} a_{H+h} &= p_{i, i+h} = U(h + NQ^s) \\ b_{H-i} &= p_{i, 0} = F_U(NQ^s - i), \end{aligned}$$

$$\begin{bmatrix}
b_H & a_{H+1} & \dots & a_n & 0 & \dots \\
b_{H-1} & a_H & a_{H+1} & \dots & a_n & \dots \\
b_{H-2} & a_{H-1} & a_H & a_{H+1} & \dots & a_n & \dots \\
\dots & \dots & \dots & \dots & \dots & \dots & \dots \\
b_1 & a_2 & \dots & a_H & a_{H+1} & \dots & a_n & \dots \\
a_0 & a_1 & a_2 & \dots & a_H & a_{H+1} & \dots & \dots \\
0 & a_0 & a_1 & a_{H+H-4} & \dots & a_H & a_{H+1} & \dots \\
0 & 0 & a_0 & a_1 & a_2 & \dots & a_H & \dots \\
\vdots & \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \ddots
\end{bmatrix},$$

Figure 2. Structure of the transition matrix P

and H is the minimum integer such that $U(NQ^s + h) = 0$ for all $h < H$. This structure is recursive: from row H onward, each row is obtained by shifting the previous one to the right and inserting a 0 in the first position. Furthermore, the first element greater than zero of such recursive rows is dubbed a_0 , while the last with a_n : $n = \max\{i | a_i > 0\}$. We now introduce a useful notation for sub-matrices.

Definition 1: Let $P = (p_{i,j})$ be a matrix whose elements are $p_{i,j}$. Let $\alpha = \{i_1, i_2, \dots, i_n\}$ $\beta = \{j_1, j_2, \dots, j_m\}$ two ordered set of indexes. The sub-matrix $P_{[\alpha, \beta]}$ is a matrix whose elements are p_{i_h, j_t} for all $h \in [1, n]$ $t \in [1, m]$. Likewise, if π is a vector, we denote $\pi_{[\alpha]}$ the sub-vector whose elements are π_{i_h} for all $h \in [1, n]$.

From the properties of our transition matrix we can prove the following result [28].

Theorem 1: Let H be the minimum integer such that $U(NQ^s + h) = 0$ for all $h < H$. Let F be defined as $\max\{n - H, H\}$. Define $\alpha(i, F)$ the set $\{i, \dots, i + F - 1\}$ and $\beta(j, F)$ the set $\{j, \dots, j + F - 1\}$. The transition matrix P is block-tri-diagonal with the structure in Equation 2, where $A_0 = P_{[\alpha(F, F), \beta(0, H)]}$, $A_2 = P_{[\alpha(0, F), \beta(F, F)]}$, $A_1 = P_{[\alpha(F, F), \beta(F, F)]}$, $C = P_{[\alpha(0, F), \beta(0, F)]}$, are square matrices of order H . This qualifies the process as a QBDP.

The structure of the QBDP exposed in Theorem 1 enables the application of efficient numeric solutions for the steady state probability [10], as discussed in Section VI.

C. A conservative approximation

In order to make the model tractable from the numeric point of view, it is useful to introduce a conservative approximation. The notion of conservative approximation that we shall adopt here relies on the concept of *first order stochastic dominance* (defining an order relation between probability distributions):

Definition 2: Given two random variables X and Y , with CDFs $F_x(x)$ and $F_y(y)$, X has a first order stochastic dominance over Y ($X \succeq Y$) iff $\forall x F_x(x) \leq F_y(x)$.

Based on this definition, a stochastic real-time task can be seen as a conservative approximation of another one if its probabilistic deadlines are stochastically dominated by the probabilistic deadlines of the original task: considering δ_k in Equation (3), this plainly means that in the modified system the low values of the δ_k will have a greater probability and so will be the probability of the first element of the probability vector (associated with the deadline satisfaction).

As shown by Diaz et al. [13], if \mathcal{U}' stochastically dominates \mathcal{U} , then a system having the execution times distributed according to \mathcal{U}' is a conservative approximation of the original system (with the execution times distributed according to \mathcal{U}).

A simple way to build \mathcal{U}' to obtain such a conservative approximation is to replace c_k with a new variable c'_k whose distribution is given by:

$$U_\Delta(c') = \begin{cases} 0 & \text{if } c' \bmod \Delta \neq 0 \\ \sum_{c=(k-1)\Delta+1}^{k\Delta} U(c') & \text{otherwise,} \end{cases} \quad (4)$$

where Δ is a scaling factor chosen as an integer sub-multiple of Q^s . The transition matrix of the new DTMC has again the structure in Fig. 2, where the different elements of the matrix are functions of the parameter Δ . Large values of Δ correspond to a smaller size for matrices A_2 , A_1 , A_0 in Equation 2. This reduces the time required for the computation of the steady state probability paying the price of a coarser approximation for the computed probability.

V. AN ANALYTICAL BOUND

This section presents an analytic solution for a QBDP described by the transition matrix reported in Fig. 2. In the discussion, we assume that the conservative approximation discussed in Section IV-C for some Δ .

The first key result of the Section is Theorem 2, which shows a general expression for the steady state probability of respecting the deadline. After introducing an additional simplification in the model, this leads to the analytical bound in Theorem 6 and in Corollary 7, which represent the core theoretical results of the paper.

A. A solution for generic QBDP processes

Before going into the theoretic details, let us define the following function $\gamma : \mathbf{N} \times \mathbf{R} \rightarrow \mathbf{R}$ as

$$\gamma_{k,l} = \sum_{j=0}^k \alpha_j l^{k-j},$$

where $\alpha_j = a_j/a_0$. Using this function and the structure of the QBDP, it is possible to write the equation expressing the steady state equilibrium $\tilde{\pi}_k = \tilde{\pi}_k P$, (where $\tilde{\pi}_k = [\tilde{\pi}_k^{(0)}, \tilde{\pi}_k^{(1)}, \dots]$) by expressing the probabilities $\tilde{\pi}_k^{(i)}$, $i > H$, at time k as a function of $\tilde{\pi}_k^{(j)}$, $0 \leq j \leq H$, in the following way:

$$\begin{aligned} \tilde{\pi}_k^{(H)} &= \sum_{j=H+1}^n \alpha_j \tilde{\pi}_k^{(0)} - \sum_{j=1}^{H-1} \gamma_{j,1} \tilde{\pi}_k^{(H-j)}, \\ \tilde{\pi}_k^{(H+l)} &= \left(\gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j \right) \tilde{\pi}_k^{(l)} - \sum_{\substack{j=1 \\ j \neq H}}^{\min(n, l+H)} \alpha_j \tilde{\pi}_k^{(l+H-j)}, \end{aligned} \quad (5)$$

holding for $\forall l > 1$.

The steady state solution for generic $n > H > 0$ is given by the following theorem:

Theorem 2: Consider a QBDP described by the transition probability matrix P given in Fig. 2, in which both a_0 and a_n differ from zero.

Assume that the matrix

$$W = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 0 & \cdots & 1 \\ -\alpha_n & -\alpha_{n-1} & -\alpha_{n-2} & \cdots & w & -\alpha_{H-1} & \cdots & -\alpha_0 \end{bmatrix} \quad (6)$$

where $w = \gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j$, has distinct eigenvalues.

Let $\pi^{(j)} = \lim_{k \rightarrow +\infty} \tilde{\pi}_k^{(j)}$ be the steady state distribution of the state. One of the two following cases apply:

I) if $\sum_{j=0}^{H-1} \gamma_{j,1} \leq \sum_{j=H+1}^n (j-H)\alpha_j$ then the limiting distribution is given by:

$$\tilde{\pi}^{(j)} = \lim_{k \rightarrow +\infty} \tilde{\pi}_k^{(j)} = 0, \quad \forall j, \quad (7)$$

II) if $\sum_{j=0}^{H-1} \gamma_{j,1} > \sum_{j=H+1}^n (j-H)\alpha_j$ then:

$$\tilde{\pi}^{(0)} = \prod_{\beta \in \mathcal{B}_s} (1 - \beta). \quad (8)$$

In the second case, \mathcal{B}_s is the set of stable eigenvalues of W (in this context an eigenvalue β is said stable if $|\beta| < 1$), and the terms $\tilde{\pi}^{(j)}$ with $0 < j < H$ are known linear functions of $\tilde{\pi}^{(0)}$, while the terms $\tilde{\pi}^{(j)}$ with $j \geq H$ are given by (5).

Before showing the proof, we make two important remarks.

Remark 1: The assumption on the eigenvalues of the matrix W is merely technical (it simplifies the proof of the result) and it is not restrictive. In all our examples (both synthetically generated and using data from real applications), it is respected. Artificial examples that violate it could probably be constructed but they are not relevant in practice.

Remark 2: As well as paving the way for Theorem 6, Theorem 2 contains an implicit numeric algorithm for the computation of $\tilde{\pi}^{(0)}$, based on the computation of the eigenvalues of the matrix W . Since the latter is in companion form, in the following we refer to this algorithm as companion.

B. Proof of Theorem 2

This section is devoted to the proof of the fundamental Theorem 2, which will require several definitions and auxiliary results. The section can be skipped over if the reader is only interested in the applications of the Theorem.

The rationale behind the proof is the following. First, the equilibrium point of the QBDP is expressed as an iterative system. The evolution in the iteration step represents the connection between the probabilities of the different states. Using this representation and some property of convergence of the Markov chain, we can express all the steady-state probabilities as a function of $\tilde{\pi}^{(0)}$, which can eventually be found as a solution of a linear system of equations.

We start noticing that having a_0 and a_n different from zero implies that the Markov chain of the QBDP is irreducible and aperiodic. Therefore, it is guaranteed that the probability of the different states converge to a value [29]. Notice, however, that this does not necessarily imply the existence of a steady-state distribution (the distribution could shift toward increasing values of the state without ever reaching the equilibrium, with the probability of each state going to 0).

a) *The case of Positive Recurrent QBDP:* If the QBDP is positive recurrent, it admits indeed a unique steady state distribution. The first step of the proof is then to introduce the following vector: $\Pi_j = [\tilde{\pi}^{(j)}, \dots, \tilde{\pi}^{(j+n-1)}]^T$, whose dimension is equal to n . It is possible to exploit (5) and (6) to derive the equilibrium of the QBDP by the following iterative equation for the vector Π_j :

$$\Pi_1 = \begin{bmatrix} \tilde{\pi}^{(1)} \\ \tilde{\pi}^{(2)} \\ \vdots \\ \tilde{\pi}^{(n)} \end{bmatrix} = W \Pi_0 \Rightarrow \Pi_j = \begin{bmatrix} \tilde{\pi}^{(j)} \\ \tilde{\pi}^{(j+1)} \\ \vdots \\ \tilde{\pi}^{(n-1+j)} \end{bmatrix} = W^j \Pi_0.$$

Using this notation the normalisation constraint $\sum_{h=0}^{\infty} \tilde{\pi}^{(h)} = 1$ can be expressed as

$$\sum_{h=0}^{\infty} \tilde{\pi}^{(h)} = [1 \quad 0 \quad 0 \quad \dots \quad 0] \sum_{i=0}^{+\infty} \Pi_i = 1. \quad (9)$$

The characteristic polynomial of the lower-left companion form matrix W reported in (6) is simply given by

$$P(\lambda) = \lambda^n - \left(\gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j \right) \lambda^{n-H} + \sum_{\substack{j=1 \\ j \neq H}}^n \alpha_j \lambda^{n-j}, \quad (10)$$

from which it is trivially derived that the matrix W has one simple eigenvalue in $\beta_1 = 1$ and additional $n-1$ eigenvalues β_i . Therefore

$$P(\lambda) = (\lambda - 1) \prod_{i=2}^n (\lambda - \beta_i). \quad (11)$$

Since each β_i verifies $P(\beta_i) = 0$, the following relation holds

$$\begin{aligned} \beta_i^n - \left(\gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j \right) \beta_i^{n-H} + \sum_{\substack{j=1 \\ j \neq H}}^n \alpha_j \beta_i^{n-j} &= 0 \Rightarrow \\ \gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j &= \beta_i \gamma_{H-1,\beta_i} + \frac{\sum_{j=H+1}^n \alpha_j \beta_i^{n-j}}{\beta_i^{n-H}}. \end{aligned} \quad (12)$$

Since all the eigenvalues are assumed simple, we can use of the *spectral decomposition* of the matrix W : $W = \sum_{i=0}^{n-1} \beta_i G_i$, where the *spectral projectors* G_i are given by $G_i = \frac{V_i L_i}{L_i V_i} = N_i V_i L_i$, and L_i and V_i are respectively the left and right eigenvectors associated with the i -th eigenvalue β_i . N_i is the normalisation constant needed to satisfy the spectral projectors basic properties, i.e., $G_i G_j = 0$ for $i \neq j$ and $G_i G_i = G_i$. As a consequence, $\Pi_1 = W \Pi_0 = \sum_{i=1}^n \beta_i G_i \Pi_0$, and, in general,

$$\Pi_j = W^j \Pi_0 = \sum_{i=1}^n \beta_i^j G_i \Pi_0 = \sum_{i=1}^n \beta_i^j N_i V_i L_i \Pi_0. \quad (13)$$

Therefore, by combining (13) and (9), one gets:

$$\sum_{i=1}^n \sum_{k=0}^{+\infty} \beta_i^k v_i^{(0)} N_i L_i \Pi_0 = 1, \quad (14)$$

where $v_i^{(0)}$ is the first element of the right eigenvector. Given the expression of the matrix W , the left L_i and right V_i can

be easily found as a function of β_i . From the expression of the eigenvectors, it follows immediately that

$$N_i = \frac{1}{L_i V_i} = \frac{\beta_i^n}{\sum_{j=0}^{H-1} \gamma_{j,\beta_i} \beta_i^{n-j} - \sum_{j=H+1}^n (j-H) \alpha_j \beta_i^{n-j}}. \quad (15)$$

We now state some auxiliary propositions on vector Π_0 .

Proposition 1: The product between the left eigenvector L_i and the initial vector of the iteration Π_0 is given by

$$L_i \Pi_0 = \beta_i^{n-H-1} (\beta_i - 1) \left(\sum_{k=0}^{H-1} \sum_{j=k}^{H-1} \gamma_{H-1-j,\beta_i} \tilde{\pi}^{(k)} \right).$$

Proof: The proof of the proposition follows by first computing the explicit computation of the product $L_i \Pi_0$, in which each term is substituted with the recursive Equations (5) and the constraint given in (12), and then noticing that

$$\beta_i^n - 1 = (\beta_i - 1) \sum_{j=0}^{n-1} \beta_i^j.$$

See [30] for more details. ■

Proposition 2: The initial vector Π_0 is orthogonal to the left eigenvector associated to $\beta_1 = 1$.

Proof: The proof follows from Proposition 1. ■

Proposition 3: For any unstable eigenvalue β_i (i.e., such that $|\beta_i| > 1$) of W it holds that $L_i \Pi_0 = 0$.

Proof: If the QBDP has an equilibrium then (14) holds true. The unitary eigenvalue $\beta_1 = 1$ does not play any role in the summation of (14) in view of Proposition 2. Next, suppose that there exists one or more $|\beta_i| > 1$. From Equation (14) it follows that it may be $L_i \Pi_0 = 0$, $N_i = 0$ or $\Pi_0 = 0$. Since the normalisation factor cannot be null, let us first consider $\Pi_0 = 0$. Using (13) it follows that $\Pi_0 = 0 \Rightarrow \Pi_j = 0, \forall j$. Therefore,

$$\tilde{\pi}^{(j)} = \lim_{k \rightarrow +\infty} \tilde{\pi}^{(j)}(k) = 0, \quad \forall j,$$

and, since the Markov chain is irreducible and aperiodic, the QBDP does not have a unique stationary distribution [29], which contradicts the hypothesis.

It then follows that for any unstable eigenvalue $L_i \Pi_0 = 0$. ■

From Rouche's theorem [31] we have that the number of eigenvalues β_i such that $|\beta_i| \geq 1$ of the matrix W is exactly equal to H , where $H-1$ have $|\beta_i| > 1$. The consequences of Proposition 3 are twofold. First, it states that Proposition 1 defines $H-1$ linear equations

$$\sum_{k=0}^{H-1} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,\beta_i} \tilde{\pi}^{(k)} = 0, \forall \beta_i \in \mathcal{B}_s^*, \quad (16)$$

where \mathcal{B}_s^* is the set of $H-1$ unstable eigenvalues except $\beta_1 = 1$ (the unstable eigenvalue β_1 does not play any role by Proposition 2). The H unknown probabilities $\tilde{\pi}^{(0)}$ to $\tilde{\pi}^{(H-1)}$ of (16) are also the unknowns of the recursion formulae (5). The second consequence is that

$$\sum_{\beta_i \in \mathcal{B}_s} \frac{v_i^{(0)} N_i}{1 - \beta_i} L_i \Pi_0 = 1, \quad (17)$$

where \mathcal{B}_s is the set of stable eigenvalues. By substituting in (17) the result given in Proposition 1 and the expression of the right eigenvector L_i , we get

$$- \sum_{\beta_i \in \mathcal{B}_s} \frac{N_i}{\beta_i^H} \sum_{k=0}^{H-1} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,\beta_i} \tilde{\pi}^{(k)} = 1. \quad (18)$$

By means of Proposition 3, the summation can be extended to the unstable eigenvalues, except for the first eigenvalue $\beta_1 = 1$, which instead induces indefiniteness of (18). The solution to (18) is derived exploiting the spectral projectors property $\sum_{i=1}^n G_i = I_n$. Indeed, summing the elements in position $(n-H, n-j)$, for $1 \leq j \leq H-1$, we have for each j

$$- \sum_{i=1}^n N_i v_i^{(n-H)} l_i^{(n-j)} = - \sum_{i=1}^n \frac{N_i}{\beta_i^H} \gamma_{j,\beta_i} = 0,$$

and hence

$$- \sum_{i=2}^n \frac{N_i}{\beta_i^H} \gamma_{j,\beta_i} = N_1 \gamma_{j,1},$$

where N_1 is easily obtained by (15) for $\beta_1 = 1$, i.e.,

$$N_1 = \frac{1}{\sum_{j=0}^{H-1} \gamma_{j,1} - \sum_{j=H+1}^n (j-H) \alpha_j} = \frac{1}{D_1}.$$

Moreover, for the elements in position $(n-H+1, 1)$, we get

$$- \sum_{i=1}^n N_i v_i^{(n-H+1)} l_i^{(1)} = \sum_{i=1}^n \frac{N_i}{\beta_i^{H-1}} \frac{\alpha_n}{\beta_i} = 0 \Rightarrow - \sum_{i=2}^n \frac{N_i}{\beta_i^H} = N_1.$$

Substituting these relations in (18) produces the equation

$$\sum_{k=0}^{H-1} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,1} \tilde{\pi}^{(k)} = D_1, \quad (19)$$

which, used in conjunction with the $H-1$ equations of (16), determines the set of unknown probabilities.

In order to have an analytic solution of this linear system of H equations in H unknowns, we start by collecting the probability with the highest index, i.e.,

$$\begin{aligned} \tilde{\pi}^{(H-1)} + \sum_{k=0}^{H-2} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,1} \tilde{\pi}^{(k)} &= D_1 \\ \tilde{\pi}^{(H-1)} + \sum_{k=0}^{H-2} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,\beta_i} \tilde{\pi}^{(k)} &= 0, \beta_i \in \mathcal{B}_s^*, \end{aligned}$$

from which it is possible to immediately have the solution

$$\tilde{\pi}^{(H-1)} = - \sum_{k=0}^{H-2} \sum_{q_1=0}^{H-1-k} \gamma_{q_1,\beta_H} \tilde{\pi}^{(k)}$$

and the $H-1$ new linear equations in $H-1$ unknowns

$$\sum_{k=0}^{H-2} \sum_{q_1=0}^{H-1-k} (\gamma_{q_1,1} - \gamma_{q_1,\beta_i}) \tilde{\pi}^{(k)} = D_1, \beta_i \in \mathcal{B}_s^*,$$

that, by simple algebraic manipulations, leads to

$$\sum_{k=0}^{H-2} \sum_{q_1=0}^{H-1-k} \sum_{q_2=0}^{q_1-1} \gamma_{q_2,\beta_i} \tilde{\pi}^{(k)} = \frac{D_1}{1 - \beta_i}, \beta_i \in \mathcal{B}_s^*.$$

From the new set of $H-1$ equations the element $\tilde{\pi}^{(H-2)}$ can be collected, thus leading to a recursive solution formula.

The recursion can be executed for H steps until the following final equation is obtained

$$\tilde{\pi}^{(0)} = \frac{D_1}{\prod_{\beta_i \in \mathcal{B}_s^*} (1 - \beta_i)} = \frac{\sum_{j=0}^{H-1} \gamma_{j,1} - \sum_{j=H+1}^n (j-H)\alpha_j}{\prod_{\beta_i \in \mathcal{B}_s^*} (1 - \beta_i)}. \quad (20)$$

The result in (20) can be suitably rewritten in a more useful way. To this end, we first rewrite the characteristic polynomial (11) as follows

$$P(\lambda) = (\lambda - 1) \prod_{i=2}^{n-1} (\lambda - \beta_i) = \lambda^{n-1} + \sum_{j=1}^{n-1} S_j(\beta) \lambda^{j-1}, \quad (21)$$

where

$$S_j(\beta) = (-1)^{n-j+1} \left(\sum_{J \in \mathcal{C}_1} \prod \beta_J + \sum_{J \in \mathcal{C}_2} \prod \beta_J \right), \quad (22)$$

and where \mathcal{C}_1 and \mathcal{C}_2 are proper sets of indices coming from the explicit computation of the characteristic polynomial. Since the product of all the eigenvalues, except for the first one, is given by

$$\prod_{i=2}^n (1 - \beta_i) = 1 + \sum_{j=1}^{n-1} (-1)^{n-j} \sum_{J \in \mathcal{C}_{n-j}} \prod \beta_J = 1 + \sum_{j=1}^{n-1} \mathcal{W}_j(\beta),$$

where, by means of (22), $\mathcal{W}_k(\beta) = -\sum_{j=1}^k S_j(\beta)$, one gets

$$\prod_{i=2}^n (1 - \beta_i) = 1 - \sum_{j=1}^{n-H} \sum_{k=1}^j S_k(\beta) - \sum_{j=n-H+1}^{n-1} \sum_{k=1}^j S_k(\beta). \quad (23)$$

From (21) and (10), $S_k(\beta) = \alpha_{n-k+1}$, for $1 \leq k \leq n$, and $S_k(\beta) = \gamma_{H-1,1} + \sum_{j=H+1}^n \alpha_j$, for $k = n - H + 1$. Substituting these relations in the last two terms of (23), one gets

$$\begin{aligned} & - \sum_{j=1}^{n-H} \sum_{k=1}^j S_k(\beta) = - \sum_{j=H+1}^n (j-H)\alpha_j, \\ & - \sum_{j=n-H+1}^{n-1} \sum_{k=1}^j S_k(\beta) = (H-1)\gamma_{H-1,1} - \sum_{j=1}^{H-1} (j-1)\alpha_j. \end{aligned}$$

Since

$$1 + (H-1)\gamma_{H-1,1} - \sum_{j=1}^{H-1} (j-1)\alpha_j = \sum_{j=0}^{H-1} \gamma_{j,1},$$

Equation (23) is rewritten as

$$\prod_{i=2}^n (1 - \beta_i) = \sum_{j=0}^{H-1} \gamma_{j,1} - \sum_{j=H+1}^n (j-H)\alpha_j = D_1, \quad (24)$$

that substituted in (20) finally yields Equation (8).

At this point we have proved that *if the QBDP has an equilibrium*, this is given by (8), by the recursive solution of the linear system of equations (19) and (16), and by the recursion formula (5).

b) The case of non-positive recurrent QBDP: If the QBDP is not positive recurrent we can re-write matrix P using its block-tridiagonal representation in (2). We can immediately apply the following theorems.

Theorem 3: [29] An irreducible Markov chain has a stationary distribution if and only if all its states are positive recurrent.

Definition 3: Assume $A = A_0 + A_1 + A_2$ is irreducible. Then, by the Perron–Frobenius Theorem, there exists a unique vector $\mu > 0$ with $\mathbf{1}^T \mu = 1$ and $A\mu = \mu$. The vector μ is called the stationary probability vector of A , while $\mathbf{1}$ is a column vector whose elements are all equal to one.

Theorem 4: [21] The QBDP is transient if $\mathbf{1}^T A_0 \mu < \mathbf{1}^T A_2 \mu$, null recurrent if $\mathbf{1}^T A_0 \mu = \mathbf{1}^T A_2 \mu$ and positive recurrent if $\mathbf{1}^T A_0 \mu > \mathbf{1}^T A_2 \mu$.

By Theorem 3, the QBDP does not have an equilibrium if and only if it has at least one state that is transient or null recurrent. Without loss of generality, assume that $n \leq 2H$ (the case $n > 2H$ can be equivalently derived), which implies $A \in \mathbf{R}^{H+1 \times H+1}$. Since A is irreducible, one immediately has that $\mu = \frac{1}{H+1} \mathbf{1}$, from which it is possible to explicitly compute

$$\begin{aligned} \mathbf{1}^T A_0 \mu &= \frac{1}{H+1} \sum_{j=0}^{H-1} (H-j)a_j \\ \mathbf{1}^T A_2 \mu &= \frac{1}{H+1} \sum_{j=H+1}^n (j-H)a_j. \end{aligned}$$

From Theorem 4, the QBDP does not have an equilibrium if and only if $\mathbf{1}^T A_0 \mu \leq \mathbf{1}^T A_2 \mu$ or, equivalently,

$$\sum_{j=0}^{H-1} (H-j)a_j \leq \sum_{j=H+1}^n (j-H)a_j,$$

that, dividing both terms by a_0 leads to

$$\sum_{j=0}^{H-1} \gamma_{j,1} \leq \sum_{j=H+1}^n (j-H)\alpha_j. \quad (25)$$

This condition is exactly the one that we formulated in the case I of the Theorem, and has just been shown to be equivalent to the process being transient or null recurrent. However, since the QBDP is still irreducible and aperiodic, a limiting probability exists, which is given, as in Equation (7), by:

$$\tilde{\pi}^{(j)} = \lim_{k \rightarrow +\infty} \tilde{\pi}^{(j)}(k) = 0, \quad \forall j,$$

And this ends the proof of Theorem 2.

Remark 3: When condition (25) strictly applies, the numerator of Equation (20) is negative. Since Equation (8) still holds true, the denominator of (20) will be negative too. It follows that in the case of absence of an equilibrium for the QBDP, both (8) and (20) return a coincident value $\tilde{\pi}^{(0)} > 1$, clearly unfeasible.

C. Computation of the bound

As discussed earlier, the steady state probability of meeting the deadline can be found by computing the first element $\tilde{\pi}^{(0)}$ of the $\tilde{\pi}$ that solves the equation $\tilde{\pi} = \tilde{\pi}P$, where P is the infinite transition matrix in Fig. 2 associated with the DTMC

\mathcal{M} . Let us consider a new DTMC whose transition matrix is given by:

$$P' = \begin{bmatrix} b_H & a_{H+1} & a_{H+2} & \dots & a_{n-1} & a_n & 0 & \dots \\ b_{H-1} & a_H & a_{H+1} & \dots & a_{n-2} & a_{n-1} & a_n & \dots \\ 0 & a'_{H-1} & a_H & \dots & a_{n-3} & a_{n-2} & \ddots & \dots \\ 0 & 0 & a'_{H-1} & a_H & \dots & a_{n-3} & \ddots & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \ddots \end{bmatrix}, \quad (26)$$

and $a'_{H-1} = b_{H-1} = a_{H-1} + a_{H-2} + \dots + a_0$.

Remark 4: The underlying idea is very simple. Consider the DTMC associated with matrix P . The terms on the left of the diagonal are transition probabilities toward states with a smaller delay than the current one. By using P' we lump together all these transitions to the state immediately on the left of the current one. For instance, if the current state corresponds to 4 server periods of delay, its only enabled transition to the left will be to the state associated with delay 3. The effect of deleting the transition toward states associated with smaller delays is to slow down the convergence toward small delays, thus decreasing the steady state probability of these states.

Let π represent the steady state probability of this system. We can easily show the following:

Lemma 5: Let Γ be a random variable representing the state of the DTMC evolving with transition matrix P and Γ' be a random variable describing the state of the DTMC associated with the transition matrix P' . If both DTMC are irreducible and aperiodic, then at the steady state Γ' has a first order stochastic dominance over Γ : $\Gamma' \succeq \Gamma$, according to Definition 2. Therefore, for the first element of the steady state probability, we have $\tilde{\pi}^{(0)} \geq \pi^{(0)}$.

Proof: The proof is omitted for the sake of brevity (see [30]).

In view of this Lemma, we can concentrate on the system associated to the transition matrix P' . In such a case, we immediately derive that the equilibrium condition $\pi = \pi P'$ produces the following recursion:

$$\begin{aligned} \pi^{(1)} &= \sum_{j=2}^n \alpha_j \pi^0, \\ \pi^{(l)} &= \left(1 + \sum_{j=2}^n \alpha_j\right) \pi^{(l-1)} - \sum_{j=2}^{\min(n, H+l-1)} \alpha_j \pi^{(l-j)}, \end{aligned} \quad (27)$$

where the equalities hold for $\forall l > 1$. These equations, as well as P' , have been respectively derived from (5) and P by imposing $H = 1$. In such a situation, the following theorem holds.

Theorem 6: Consider a QBDP described by the transition probability matrix (26), in which both a_n and a'_{H-1} differ from zero. Assume that the matrix W in (6) has distinct eigenvalues after imposing $H = 1$. Then, there exists a limiting probability distribution given by

$$\begin{aligned} \pi^{(0)} &= \lim_{k \rightarrow +\infty} \pi^{(0)}(k) = \max\{1 - \sum_{j=2}^n (j-1)\alpha_j, 0\} = \\ &= \max\{1 - \sum_{j=2}^n (j-1)\frac{a_j}{a_0}, 0\}, \end{aligned} \quad (28)$$

while the generic terms $\pi^{(j)}$, with $j > 0$, are given by (27).

Proof: The proof follows immediately from the fact that $H = 1$ implies that $\beta_1 = 1$ is the only unstable eigenvalue if the QBDP has an equilibrium, i.e., \mathcal{B}_s of Theorem 2 comprises all the eigenvalues except $\beta_1 = 1$. Hence, by considering (24) for $H = 1$, the proof follows immediately. ■

We complete the section with a remark. The first one is on the intuitive meaning of the result just proposed. Consider a DTMC with transition matrix as in Fig. 2 and assume for simplicity $n = 4$ and $H = 1$. The analytical bound in Theorem 6 is given by:

$$\pi^{(0)} = 1 - 3\alpha_4 - 2\alpha_3 - \alpha_2 = 1 - 3\frac{a_4}{a_0} - 2\frac{a_3}{a_0} - \frac{a_2}{a_0}$$

In the computation of the steady state probability $\pi^{(0)}$ we have to consider every possible transition to the right (i.e., increasing the delay) that the system can make. For each of them, we compute the ratio between the probability of taking the transition and the aggregate probability of moving to the left (decreasing the delay). In the final computation each of this ratio has a state proportional to the delay introduced. In our example, a_4 corresponds to three steps to the right and is weighted by the factor 3.

The application of this result to our context can be formalised in the following:

Corollary 7: Consider a resource reservation used to schedule a periodic task and suppose that the QBDP produced respects the assumption in Theorem 2. Then the probability of respecting the deadline is greater than or equal to:

$$\pi^{(0)} = 1 - \sum_{j=2}^n (j-1) \frac{U'_\Delta(N+j-1)Q^s}{\sum_{h=0}^{N-1} U'_\Delta(hQ^s)} \quad (29)$$

This corollary descends from the following facts: 1) the DTMC described by the matrix P in Fig. 2 is a conservative approximation of the system, 2) Lemma 5 provides an analytically tractable approximation of the DTMC with transition matrix P' , 3) Theorem 2 and Theorem 6 contain the analytical bounds.

VI. EXPERIMENTAL VALIDATION

We have validated the presented approach in two different ways. First, we have computed the probabilistic deadline using synthetic distributions, to compare accuracy and efficiency of the analytic bound against several other methods and to assess the impact of the scaling factor Δ (Eq. (4)) and of the bandwidth. This set of experiment reveals a very good performance of the bound for appropriate choices of the scaling factor Δ . Its very low computation time allows one to select the best choice of Δ by testing a number of alternative choices. The tightness of the bound improves when the bandwidth is sufficient to achieve an acceptable real-time behaviour for the application.

In a second set of experiments, we have evaluated the method on a real robotic application, for which the mathematical assumptions underlying the model do not apply strictly. The results produced are obviously approximate. Still, the good quality of the approximation makes an interesting case for the practical applicability of the methodology.

A. Synthetic Distributions

We report the results of the comparison between the numeric solution resulting from Theorem 2 and discussed in Remark 2 (*companion*), the analytic approximated bound in Corollary 7 (*analytic*) the Cyclic Reduction algorithm [10] (CR) and the bound developed by Abeni et al. [32] (*gamma*). We have chosen CR after a selection process in which several algorithms for the solution of general QBDP problems and implemented in the SMC Solver tool-suite [33] were tested on a set of example QBDPs derived from our application. The *gamma* algorithm is an approximate bound specifically tailored to the analysis of probabilistic guarantees for resource reservations, so it was considered as a perfect match for our *analytic* bound. The different algorithms have been implemented in C++ in the PROSIT [34] tool. PROSIT can be used for analysis and for synthesis purposes (as shown in Section VII). When the tool is used for analysis, the user specifies activation period and deadline, parameters of the RR (Q^s and T^s), distribution of computation and inter-arrival times and solution algorithm. When the tool is queried in this way, it computes the distribution of the task response times and hence the probability of meeting the deadline.

As a representative sample of our findings, we report below the results obtained for a periodic task with period $T = 100ms$ and random execution time. The computation time was distributed according to a beta distribution: $P\{C = c\} = f_U(c) = J(\alpha, \beta)c^{\alpha-1}(1-c)^{\beta-1}$, with support (i.e., the validity range for the random variable) $c \in [0, 99500] \mu s$, with $\alpha = 2$ and $\beta = 7$ ($J(\alpha, \beta)$ is a normalisation constant). The beta distribution is interesting because it is unimodal and has a finite support, which make it a good fit to approximate the behaviour of a large number of real-time applications.

Effect of Δ . A first set of experiments was to evaluate the impact of the Δ scaling factor. We considered two possible values for the reservation period: $T^s = \frac{1}{4}P = 25ms$ and $T^s = \frac{1}{2}P = 50ms$. The budget was chosen equal to $Q^s = 0.45T^s$ with a bandwidth $B = 45\%$. Figure 3 shows the results for the probability $\pi^{(0)}$ of respecting the deadline achieved for different values of Δ (chosen as a sub-multiple of Q^s). In accordance with our expectations, CR and *companion* produce almost the same result in term of probability (differences are from the 6th digit) and the probability changes monotonically with Δ . For example, for $T^s = 50ms$ the value of the probability is 0.89 for $\Delta = Q^s$ (the coarsest possible granularity), while it is 0.93 for $\Delta = Q^s/45$. The reason for this decrease is obvious since re-sampling introduces a conservative approximation and the error is larger for increasing granularity. For both CR and *companion*, the computation time changes with Δ in a substantial way. For example, for CR and for $T^s = 50ms$, it is 182ms at $\Delta = Q^s$ and 56.179ms at $\Delta = Q^s/45$. In this run of experiments, the computation time of the *companion* algorithm is slightly smaller than the one reported using CR, but the results are too close to claim a clear dominance.

For the *analytic* bound the computed probability is not

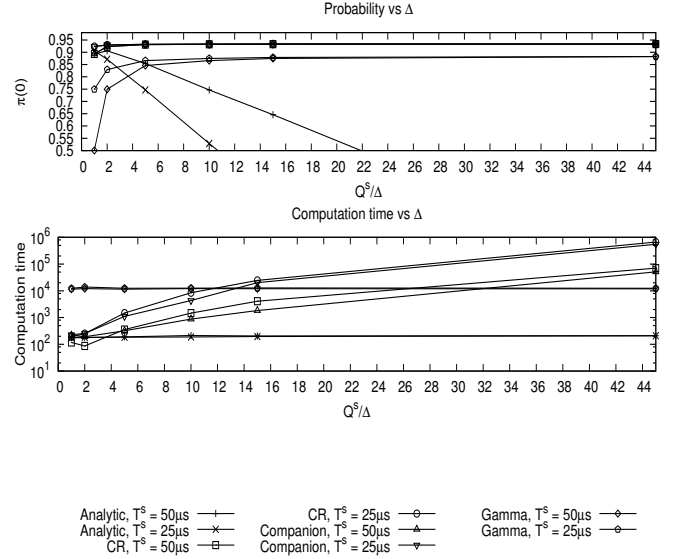


Figure 3. Impact of the scaling factor Δ on the accuracy of the computed probability and on the computation time

always monotonic with Δ . In our example, for $T^s = 50ms$ the probability grows moving from 0.892 at Q^s to 0.906 at $Q^s/2$, and then decreases, finally becoming 0.012 at $Q^s/45$. Sharper changes can be observed for other distributions. The reason is that in the *analytic* bound we have two distinct effects (which play in opposite directions). On the one hand, if we reduce Q^s we have the same conservative approximation effect as for CR or for any other numeric method. On the other, as explained in Remark 4, lumping together all backward transitions reduce the recovery of the error when the computation demand is smaller than the allocated bandwidth. In this example, the first effect determines the growth of the probability when going from $\Delta = Q^s$ to $\Delta = Q^s/2$; the second effect determines the decrease of the probability from $Q^s/2$ onward. The probability computed by *analytic* is very close to the one of the numeric algorithm it derives from (*companion*) for $\Delta = Q^s/2$, while the computation time is several orders of magnitude below. In our experience with different distributions (both synthetic and experimental) the choice of $\Delta = Q^s/2$ has consistently produced an acceptable performance. The *gamma* bound shows an intermediate performance between numeric methods and the *analytic* bound both for the accuracy and for the computation time.

Behaviour with changing bandwidth. In order to compare the accuracy of the *analytic* method against the numeric solutions (CR) for different bandwidths, we considered a task with the activation and scheduling parameters as in the experiments reported above. The budget Q^s was changed so that the resulting bandwidth ranged in $[35\%, 60\%]$. The granularity Δ was fixed for CR to a small value ($50\mu s$) to achieve a good approximation and to $\Delta = Q^s/2$ for the *analytic* solution.

The results reported in Table I show an important gap between *analytic* and CR for small values of the band-

Table I
PROBABILITY FOR DIFFERENT BANDWIDTH AND $\Delta = 50\mu s$

Bandwidth	35%	40%	45%	50%	60%
Analytic Bound	0.602	0.809	0.906	0.956	0.991
Cyclic Reduction	0.773	0.878	0.929	0.965	0.992

width. The gap is significantly reduced for bandwidth greater than 45%/50%. Smaller values of the bandwidth produce a probability level below 0.8, which is not acceptable for most real-time applications. The reason for the improvement of the analytic bound when the bandwidth increases is probably due to the fact that the system recovers more easily from large delays and this alleviates the impact of the conservative simplifications that underlie the analytic model.

B. Real application

As a real test case, we have considered a robotic vision programme that identifies the boundaries of the lane and estimates the position of a mobile robot using a web-camera mounted on the chassis of the robot [35]. The computation was carried out using a Beagle Board (www.beagleboard.org) running Ubuntu. The version of the Kernel used (3.16) supports RR scheduling (under the name of SCHED_DEADLINE policy) alongside the standard POSIX real-time fixed priority policies (SCHED_FIFO and SCHED_RR).

The robot executed 30 different paths across an area delimited by a black line. For each run, we have captured a video stream containing the line. The data sets roughly consisted of 2500 frames each and were later used for multiple off-line execution of the vision algorithm. A first group of ten executions for each data set was with the algorithm executed in a task running alone and scheduled with the the maximum real-time priority (99 for SCHED_FIFO). This allowed us to collect statistics of the computation time associated with the data set. In a second group of executions, we have replicated a real-life condition. The vision algorithm was in this case executed in a periodic task processing a frame every $T = 40ms$. The task was scheduled using SCHED_DEADLINE, with server period $T^s = 20ms$ and with different choices of the bandwidth in the range [35%, 60%]. For each data set and for each choice of the bandwidth, we repeated ten executions recording the probability of deadline miss. The probability averaged through the 10 execution was compared with the one that found using the PROSIT tool, executed with different solution methods and with the distribution estimated from the data set as input. In Figure 4, we report the CDF distributions of the difference between the two probabilities for three representative choices of the bandwidth. The symbol $\Delta_{Analytic}$ denotes the difference obtained using the analytic method (with different choices of the scaling factor Δ), while Δ_{CR} denotes the difference obtained using the cyclic reduction QBDP solver, with Δ set to $50\mu s$. The three levels of bandwidth shown in the three sub-plots produced different probability of meeting the deadline. For bandwidth equal to 40%, this probability ranged in [75%, 97%]. The range was [90.5%, 99%] for bandwidth

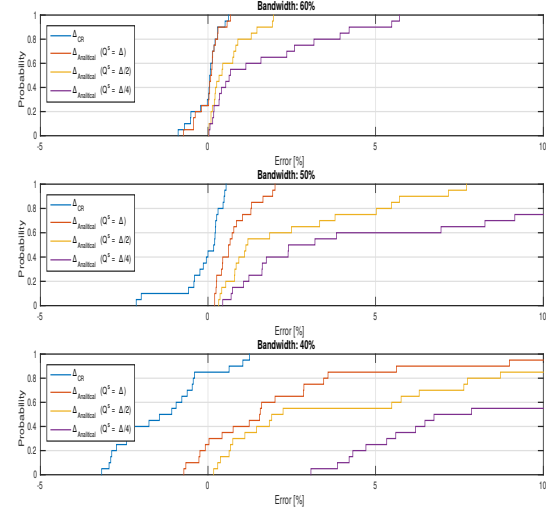


Figure 4. Distribution of the difference between the experimental probability and the one found with PROSIT tool.

equal to 50% and it was [95.2%, 100%] for bandwidth equal to 60%.

As we observe in the plot, the numeric algorithm (CR) produces an error between -3% and 1% for all the three values of the bandwidth. For the analytic bound, in this specific case, the most convenient choice was to set the scaling factor Δ to Q^s (in other cases we found a better performance factor for smaller values). The bound is evidently less accurate, but: 1. it remains below 5% at least 85% of the times even in the most challenging scenario (small bandwidth), 2. is reduced to below 2% for higher values of the bandwidth.

We observe that the vision algorithm iteratively builds upon previous results to produce the estimate. This introduces a strong correlation structure in the process that disrupts the assumptions required for an exact application of the method. In addition, the execution on a “real” operating system comes along with an inevitable amount of un-modelled overhead. Still, the level of approximation that we have reported could be acceptable in most cases. Similar software applications (video-encoding and decoding) were analysed in a previous work [36] with similar conclusions. Clearly, we are not claiming any generality for this fact. We are aware that for other applications dropping the time dependency and the correlation structure of the computation time process could produce very large errors in the estimation of the probability. As reported in the related work, this is a very active research area that is likely to attract the attention of different researchers in the forthcoming years.

C. Discussion

In our first conference paper [28], we derived a model for the evolution of a RR scheduled real-time task. The model was shown to be a QBDP and was solved using the simple numeric algorithm proposed by Latouche and Ramaswami [21]. An important limitation of the model was its pessimism due to

the fact that it neglected the budget shared between adjacent jobs. For instance, in the example in Figure 1, the model would ignore the budget used by the second job in the fourth reservation period. In a later work [36], the same model was instantiated to the sub-case of periodic tasks, it was further simplified in a conservative direction and then used for the computation of an analytic bound.

In the present paper, we start from the more accurate model introduced by Abeni and Buttazzo back in 1998 [4], and we instantiate it to the case of periodic tasks (Section IV-B). We introduce the scaling factor Δ (Section IV-C) obtaining, once again, a QBDP. When the model is used for numeric computations, the Δ parameter allows us to decide the degree of pessimism introduced in the analysis. If we set $\Delta = 1$, we obtain a close approximation of the actual behaviour of the task. If we set $\Delta = Q^s$, we recover the conservative model used in our previous work [28]. As shown in Figure 3, very different trade-offs between computation time and accuracy of the probability result from different choices of Δ .

The key contribution of this paper is found by applying the same type of analytic reasoning as in [36], but with a few substantial differences in the final result. Indeed, Theorem 2 contains an exact formula for the computation of the steady state probability of meeting the deadline, which is used as a basis for a novel numeric algorithm with competitive performance with respect to the state of the art. On the contrary, the key result of [36] is an analytic bound which can sometimes be very conservative. The same bound is rediscovered in this paper specialising Theorem 2 to a conservative approximation of the model (see Theorem 6). Once again, we can take advantage of the configuration options offered by Δ to refine the precision of the result. As shown in Figure 4, the choice $\Delta = Q^s$ (which applies the model proposed in [36]) is not guaranteed to be the best one in all cases. Therefore, the generalisation shown in this paper is relevant both from the theoretical and from the practical point of view.

VII. PROBABILISTIC QUALITY OPTIMISATION

In order to show a practical application of our approach, we have considered a situation where a single computing board (e.g., a video server, or a set-top box) is used to process (in real-time) multiple videos at the same time. This example is based on two different videos (encoded with a bit-rate of 600Kb/s): the first one, “Bridge-Close”, displays a bridge with occasional people coming through (so, it is characterised by a single, almost static scene with slow movements) and comes from a public data base (<http://trace.eas.asu.edu/yuv/index.html>); the second video (“ufo”), instead, is a movie trailer (freely available at <http://www.theufo.net> - trailer 1) characterised by frequent scene changes and rapid movements.

One of the best known ways to evaluate the quality of a video is the Peak Signal to Noise Ratio (PSNR), which is computed comparing pairwise the frames of the original raw video and of the one obtained after encoding and decoding it [37], [5]. This metric can be evaluated considering a video

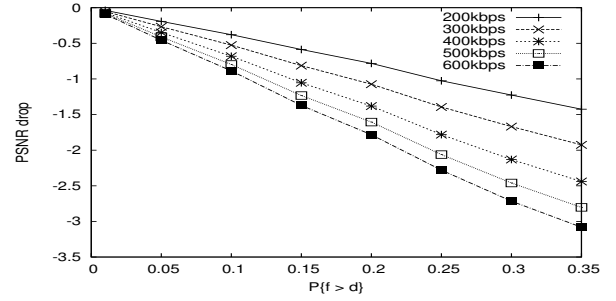


Figure 5. PSNR degradation as a function of the deadline miss probability for “BridgeClose” video.

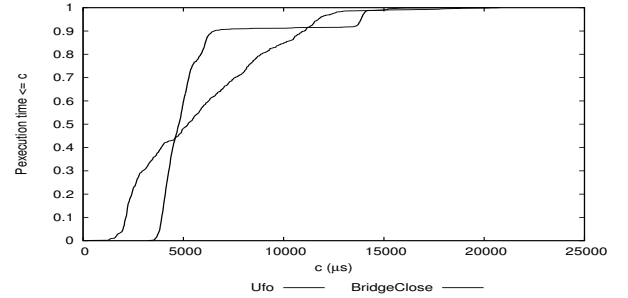


Figure 6. Cumulative Distribution Functions for the execution of the decode for the two streams.

player implemented as a periodic real-time task. If a job misses its deadline, the video frame is not played back but it is decoded (to allow the incremental decoding of the frames that follow). In this case, the behaviour of most players is to fill-in the “hole” by simply repeating the last decoded frame. This is perceived by the user as a reduction in quality, which is well reflected in a degradation of the PSNR. This is visible in Fig. 5, where we show the quality as a function of the probability of deadline miss for the first video. This plot has been created using the PSNR-TOOL software [5].

The PSNR was interpolated by a line with slope 8.9 for “BridgeClose” and 42.051 for “ufo”. This difference is explained by the different nature of the movies (static the former, and dynamic the latter). Both movies were decoded using a player executed by a periodic task and scheduled by the SCHED_DEADLINE policy. The distributions of the execution times were recorded on a notebook powered by an Intel Atom Processor, and the resulting CDFs are shown in Fig. 6.

The problem considered here was to find an optimal allocation of bandwidth between the different tasks. To this end, we have used the synthesis abilities of PROSIT. When PROSIT is used for synthesis, the user specifies for each task: 1) activation period and deadline, 2) reservation period, 3) distribution of the computation time 4) solution algorithm for the probabilistic guarantees, 5) quality as a function of the probability of meeting the deadline and 6) constraints on the minimal value of the quality. The quality of the different tasks can be combined into global quality metrics. In this

Table II
RESULTS OF PROBABILISTIC OPTIMISATION

Cyclic Reduction – Computation time:753801758 μ s				
Task	Opt. Budget	Estim. Prob.	Exact Prob.	Quality
BridgeClose	3000us	0.7427	0.743592	39.65
Ufo	6449us	0.9995	0.9995	41.58
Analytic Bound – Computation time:114524 μ s				
Task	Opt. Budget	Estim. Prob.	Exact Prob.	Quality
BridgeClose	3462us	0.7392	0.8292	40.50
Ufo	3997us	0.8732	0.9138	37.98

particular example, we have used the infinity norm metric: assuming f_i as the quality of the i^{th} task, the cost function to maximise over the budget Q_1^s and Q_2^s is $\max_i \min f_i$. For each candidate choice of Q_i^s the tool evaluates the steady state probability using different solvers for probabilistic guarantees. The optimal solution is found by a bisection algorithm, which uses repeated calls to the algorithm for the computation of the probability. As a solver for the probability computation we have implemented *analytic* (with $\Delta = Q^s/2$) and *CR* (with $\Delta = 50 \mu$ s).

Choosing 30 ms for the activation period (corresponding to 33 fps), setting the server period to 10 ms, and restricting the total bandwidth available to 95% (to leave some room for other applications), the tool produces the results in Table II. We identified empirically the minimum acceptable PSNR as 39 for “Ufo” and 31 for “BridgeClose”. These values were codified as constraints in the optimisation problem. In both cases, the algorithm identified a sub-optimal solution, because the probability evaluated by the solvers is only a lower bound. We re-evaluated the exact probability for each of the sub-optimal assignment of budgets using the CR solver with $\Delta = 1$ (which produces the exact computation of the probability, within the limits of numeric errors). This allowed us to compare the actual quality attained by the optimisation algorithm in the two different configurations. Because the optimiser maximises the worst performance of the two tasks, the algorithm tends to equalise the QoS achieved by the tasks for the optimal budget. For both solvers, the optimal solution assigns a larger bandwidth (almost 64% for the CR and almost 40% for the *analytic*) to the “Ufo” stream; this is because its quality degrades more quickly with the probability of meeting the deadline for “Ufo” than for “BridgeClose”. In this example, the use of the analytic bound produces an optimal value 37.98 which is only 4% away from the value obtained with cyclic reduction, but the computation time (evaluated on an Intel Core i7 with 16GB of RAM) is four orders of magnitude below.

VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we have considered the problem of probabilistic guarantees for RR scheduled soft real-time periodic tasks. We have shown that the evolution of the system can be modelled as a QBDP. The probability of meeting the deadline amounts to the computation of the steady state probability of this process. We have shown how this is possible by numeric means with different performance/accuracy tradeoffs. We have

also shown an analytical bound and offered a comprehensive validation of these results by experiments and simulations.

The gap between the analytic bound and precise numeric solution narrows down when the task is required to meet the deadline with a high probability (e.g., more than 80%). For this reason, the analytic bound appears as a very promising option to solve QoS optimisation problems involving multiple tasks, when the QoS is a function of the probability for the task to meet its deadline and an acceptable level of performance is required to all tasks. In these cases, the frequent calls to the solver to identify the optimal allocation of resources, such as are required by branch and bound or dichotomic search optimisation, can lead to substantial reduction of the computation time when the analytic bound is used in the face of an acceptable distance from the optimal solution.

Future work In our future work, we will investigate further on the connection between QoS and probabilistic deadlines in several application domains, we will extend our analysis and the application of our methods to the case of applications based on multiple tasks and to the case of computation time that is not i.i.d.

REFERENCES

- [1] L. Abeni and G. Buttazzo, “Integrating multimedia applications in hard real-time systems,” in *Proceedings of the IEEE Real-Time Systems Symposium*, Madrid, Spain, December 1998.
- [2] D. Fontanelli, L. Greco, and L. Palopoli, “Soft RealTime Scheduling for Embedded Control Systems,” *Automatica*, vol. 49, pp. 2330–2338, July 2013.
- [3] A. Cervin, B. Lincoln, J. Eker, K. Arzen, and G. Buttazzo, “The jitter margin and its application in the design of real-time control systems,” in *Proceedings of the IEEE International Conference on Real-Time and Embedded Computing Systems and Applications*. Gothenburg, Sweden, 2004.
- [4] L. Abeni and G. Buttazzo, “Qos guarantee using probabilistic dealines,” in *Proceedings of the Euromicro Conference on Real-Time Systems*, York, England, June 1999.
- [5] C. Kiraly, L. Abeni, and R. L. Cigno, “Effects of p2p streaming on video quality,” in *Proceedings of the IEEE International Conference on Communications*. IEEE, 2010.
- [6] M. Joseph and P. Pandya, “Finding response times in a real-time system,” *The Computer Journal*, vol. 29, no. 5, p. 390, 1986.
- [7] C. L. Liu and J. Layland, “Scheduling algorithms for multiprogramming in a hard real-time environment,” *Journal of the ACM*, vol. 20, no. 1, 1973.
- [8] R. Rajkumar, K. Juvva, A. Molano, and S. Oikawa, “Resource kernels: A resource-centric approach to real-time and multimedia systems,” in *Proceedings of the SPIE/ACM Conference on Multimedia Computing and Networking*, January 1998.
- [9] L. Abeni and G. Buttazzo, “Stochastic analysis of a reservation-based system,” in *Proceedings of the IEEE International Parallel and Distributed Processing Symposium*, San Francisco, California, April 2001.
- [10] D. Bini, G. Latouche, and B. Meini, *Numerical methods for structured Markov chains*. Oxford University Press, 2005.
- [11] A. K. Atlas and A. Bestavros, “Statistical rate monotonic scheduling,” in *Proceedings of the IEEE Real-Time Systems Symposium*, Madrid, Spain, December 1998.
- [12] J. L. Diaz, D. F. Garcia, K. Kim, C. G. Lee, L. Lo Bello, J. M. López, S. L. Min, and O. Mirabella, “Stochastic analysis of periodic real-time systems,” in *Proceedings of the IEEE Real-Time Systems Symposium*. IEEE, 2002.
- [13] J. L. Diaz, J. M. López, M. Garcia, A. M. Campos, K. Kim, and L. Lo Bello, “Pessimism in the stochastic analysis of real-time systems: Concept and applications,” in *Proceedings of the IEEE Real-Time Systems Symposium*. IEEE, 2004.

- [14] D. Maxim and L. Cucu-Grosjean, "Response time analysis for fixed-priority tasks with multiple probabilistic parameters," in *Proceedings of the IEEE Real-Time Systems Symposium*, Vancouver, British Columbia, Canada, December 2013.
- [15] L. Cucu and E. Tovar, "A framework for the response time analysis of fixed-priority tasks with stochastic inter-arrival times," *ACM SIGBED Review - Special issue: The work-in-progress (WIP) session of the RTSS 2005*, vol. 3, no. 1, pp. 7–12, January 2006.
- [16] G. A. Kaczynski, L. Lo Bello, and T. Nolte, "Deriving exact stochastic response times of periodic tasks in hybrid priority-driven soft real-time systems," in *Proceedings of the IEEE Conference on Emerging Technologies and Factory Automation*, Patras, Greece, September 2007.
- [17] A. Mills and J. Anderson, "A stochastic framework for multiprocessor soft real-time scheduling," in *Proceedings of the IEEE Real-Time and Embedded Technology and Applications Symposium*. Stockholm, Sweden: IEEE, April 2010.
- [18] D.-I. Kang, R. Gerber, and M. Sakena, "Performance-based design of distributed real-time systems," in *Proceedings of the IEEE Real-Time Technology and Applications Symposium*, Montreal, Quebec, Canada, June 1997.
- [19] C.-J. Hamann, L. Reuther, J. Wolter, H. Haertig, J. Loser, and S. Schonberg, "Quality-assuring scheduling-using stochastic behavior to improve resource utilization," in *Proceedings of the IEEE Real-Time Systems Symposium*, London, December 2001.
- [20] K. S. Refaat and P.-E. Hladik, "Efficient stochastic analysis of real-time systems via random sampling," in *Proceedings of the Euromicro Conference on Real-Time Systems*, Brussels, Belgium, July 2010.
- [21] G. Latouche and V. Ramaswami, *Introduction to matrix analytic methods in stochastic modeling*. Society for Industrial Mathematics, 1987, vol. 5.
- [22] M. F. Neuts, *Matrix-geometric solutions in stochastic models: an algorithmic approach*. Dover publications, 1995.
- [23] A. F. Mills and J. H. Anderson, "A multiprocessor server-based scheduler for soft real-time tasks with stochastic execution demand," in *Proceedings of the IEEE International Conference on Embedded and Real-Time Computing Systems and Applications*, Toyama, Japan, August 2011.
- [24] M. Santos, B. Lisper, G. Lima, and V. Lima, "Sequential composition of execution time distributions by convolution," in *Proceedings of the Workshop on Compositional Theory and Technology for Real-Time Embedded Systems*, R. Davis and L. T. X. Phan, Eds., November 2011, best paper award. [Online]. Available: <http://www.es.mdh.se/publications/2215->
- [25] G. Bernat, A. Burns, and M. Newby, "Probabilistic timing analysis: An approach using copulas," *Journal of Embedded Computing*, vol. 1, no. 2, pp. 179–194, 2005.
- [26] R. Liu, A. Mills, and J. Anderson, "Independence thresholds: Balancing tractability and practicality in soft real-time stochastic analysis," in *Proceedings of the IEEE Real-Time Systems Symposium*, Rome, Italy, December 2014.
- [27] L. Cucu-Grosjean, L. Santinelli, M. Houston, C. Lo, T. Vardanega, L. Kosmidis, J. Abella, E. Mezzetti, E. Quinones, and F. Cazorla, "Measurement-based probabilistic timing analysis for multi-path programs," in *Proceedings of the Euromicro Conference on Real-Time Systems*, Pisa, Italy, July 2012.
- [28] N. Manica, L. Palopoli, and L. Abeni, "Numerically efficient probabilistic guarantees for resource reservations," in *Proceedings of the IEEE International Conference of Emerging Technologies and Factory Automation*, Krakow, Poland, September 2012.
- [29] C. G. Cassandras and S. LaFortune, *Introduction to Discrete Event Systems*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [30] L. Abeni, D. Fontanelli, and L. Palopoli, "Application of the Quasi-Birth-Death Processes techniques to probabilistic guarantees of soft realtime systems scheduled by resource reservations," DISI - Università di Trento, Tech. Rep., 2015, <http://disi.unitn.it/~palopoli/publications/QBDP-TR.pdf>.
- [31] N. G. Lloyd, "Remarks on Generalising Rouch's Theorem," *Journal of the London Mathematical Society*, vol. s2-20, no. 2, pp. 259–272, 1979.
- [32] L. Abeni, N. Manica, and L. Palopoli, "Efficient and robust probabilistic guarantees for real-time tasks," *Journal of Systems and Software*, vol. 85, no. 5, pp. 1147–1156, May 2012.
- [33] D. Bini, B. Meini, S. Steffé, J. F. Pérez, and B. Van Houdt, "Smcsolver and q-mam: tools for matrix-analytic methods," *ACM SIGMETRICS Performance Evaluation Review*, vol. 39, no. 4, pp. 46–46, 2012.
- [34] L. Palopoli, L. Abeni, and D. Fontanelli, "A tool for the optimal design of soft real-time systems," in *Proceeding of WATERS 2014 workshop*, Madrid Spain, July 2014.
- [35] D. Fontanelli, F. Moro, T. Rizano, and L. Palopoli, "Vision-based robust path reconstruction for robot control," *IEEE T. Instrumentation and Measurement*, vol. 63, no. 4, pp. 826–837, 2014. [Online]. Available: <http://dx.doi.org/10.1109/TIM.2013.2289091>
- [36] L. Palopoli, D. Fontanelli, N. Manica, and L. Abeni, "An analytical bound for probabilistic deadline," in *Proceedings of the Euromicro Conference on Real-Time Systems*. Pisa, Italy: IEEE, September 2012.
- [37] J. Klaue, B. Rathke, and A. Wolisz, "Evalvid - a framework for video transmission and quality evaluation," in *Proceedings of the International Conference on Modelling Techniques and Tools for Computer Performance Evaluation*, 2003.