

## Modeling and performance analysis of priority queueing systems

Dariusz Strzeciwiłk<sup>†</sup> and Wlodek M. Zuberek<sup>‡</sup>

<sup>†</sup> Department of Applied Informatics, University of Life Sciences,  
ul. Nowoursynowska 159, 02-757 Warszawa, Poland  
email: [dariusz\\_strzeciwiłk@sggw.pl](mailto:dariusz_strzeciwiłk@sggw.pl)

<sup>‡</sup> Department of Computer Science, Memorial University,  
St. John's, NL, Canada A1B 3X5  
email: [wlodek@mun.ca](mailto:wlodek@mun.ca)

**Abstract.** The paper presents results of modeling and performance analysis of systems that support QoS (Quality of Service). Timed Petri nets are used to model the behavior of analyzed system and to obtain their performance characteristics. Traffic shaping mechanisms based on PQS (Priority Queueing Systems) are studied with queueing techniques used in implementations of QoS mechanisms. It is concluded that Petri net models provide a simple, convenient and effective approach to modeling and analysis of the performance of computer and communication systems.

**Keywords:** priority queueing systems, timed Petri nets, performance analysis, quality of service.

### 1 Introduction

Queueing theory is one of popular tools used in modeling and analyzing the quality of transmission in computer networks [1-3]. This theory uses a mathematical apparatus associated with the theory of stochastic processes, and in particular Markov processes [4]. In queueing systems, there are arrivals which require some service from the system's server. If the service requirements exceed the capacity of the server, some request are saved in the queue for later service. A queueing system may be characterized by regulations of queues, i.e., the way one determines the order of service applications in the system [5]. The most common queueing systems are FIFO (First In First Out), LIFO (Last In First Out), SIRO (Select In Random Order), PQ (Priority Queueing). The basic mechanism that supports the transfer of packages is FIFO scheduling that is easy to implement and treats all packets equally. FIFO scheduling is not suitable to provide for a good quality of service transmission, as when the packets come from different traffic flows, one of them can easily disrupt the flow of the other remaining streams. Packet processing in the order of flow means that an aggressive stream can appropriate the higher capacity of a router queue. This can result in poor transmission causing, for example, sudden increase in delays of transmitted packets.

Many packet scheduling algorithms were developed to provide better insulation between the streams [6]. In the case of the priority scheduling algorithms,

some application can be handled before others, regardless of when they occurred in the system. Priority queueing systems form a large class of queueing systems where the incoming requests are to be distinguished by their importance [7]. A typical example of the use of such algorithms are routers, to which are flowing subsequent packets. Core routers classify incoming packets to classes of traffic, and then handle packets belonging to the aggregated streams. Packets are handled in accordance with an implemented queueing mechanism and specific support and traffic shaping policies in order to provide services with the agreed QoS [8]. QoS is one of the most important challenges arising during the design and maintenance of both modern computer networks and next generation networks [9]. Guaranteeing adequate quality of service is of particular importance in the case of real-time applications such as Voice over IP [10] and video - IPTV [11]. These services are particularly sensitive to delay and require a guaranteed bandwidth [12]. To provide the desired QoS packages for the entire route from the sender to the recipient has been the subject of research for many years [13-15]. Research in this area can be divided into two groups. In the analytical methods the authors sought solutions of algebraic or differential equations which bind together the probability of events in the system. In the simulation methods were used most often implantation queueing algorithms, which were then subjected to statistical analysis. Based on the analysis of available research, it can be stated that the analytical methods most commonly include relatively simple queueing systems, which require implementation of many of the assumptions of the stochastic nature of the traffic flow. Complex systems are very difficult in the analysis and their functioning can be effectively examined by simulation methods. It should be noted that, to construct a sufficiently accurate model is not simple, and the waiting time for results could be discouragingly long. Hence, the aim of this study was to use models of Petri nets to assess the efficiency and to study the effectiveness of queueing mechanisms of PQS. Such an assessment may also be useful in the design and analysis of data in computer networks, distributed systems and multiprocessor systems. Constructed queueing models allowed estimation of significant features and parameters of the system under test.

## 2 Petri Nets and Network Models

Petri nets are a graphical and mathematical tool used in many fields. They are seen as a mathematical tool for modeling of concurrent systems [16, 17]. Although there are many varieties of Petri nets [18, 19] their common feature is the structure based on a bipartite directed graph, i.e. graph with two types of vertices, alternately connected by directed edges (or arcs). These two types of vertices represent, in general terms, conditions and events occurring in a modeled system, but each event can occur only when fulfilled are all the conditions associated with it. Formally Petri net is defined as a system  $N = (P, T, A)$  composed of a finite set  $P$  of p-elements (representing conditions), a finite set  $T$  of t-elements (representing events) and a set  $A$  of arcs connecting the p-components with t-elements and t-elements with p-elements,  $A \subseteq P \times T \cup T \times P$ ,  $A$  is also called the flow relation. p-elements connected by arcs directed to an t-element are called its input elements, while p-elements connected by arcs facing away from a t-element are called its output elements.

The overlapping of events is represented here by the so-called tokens assigned to the  $p$ -network elements, typically  $p$ -element, with which is associated at least one marker indicates that the condition represented by the element is fulfilled. Location of markers in the  $p$ -components can be described by marking function,  $m: P \rightarrow \{0, 1, 2 \dots\}$  or presented as a vector specifying the number of tags assigned to the further  $p$  elements of network  $m = [m(p_1), m(p_2), \dots]$ . The network  $N$  together with the (initial) marking function  $m_0$  is called the labeled  $M$  network,  $M = (N, m_0) = (P, T, A, m_0)$ . To evaluate the performance of the modeled system, i.e. to determine how quickly certain events may follow each other, in a Petri networks one must also take into account the duration of the modeled events [20]. Extension of networks by definitions of time allows for their use in modeling of real-time systems [21]. Formally temporal models are an extension of token models, with the additional elements of the description defining the times of overlapping of events and the probability or frequency of occurrence of random events. The temporal  $T$  network is thus defined as the system  $T = (M, c, f)$  where:  $M$  is a labeled network,  $M = (P, T, A, m_0)$ ,  $c$  is a function of the resolution of conflicts  $c: T \rightarrow [0, 1]$ , which for each class of decision gives the probability of particular events belonging to this class, and for other conflict events gives their relative frequency used for random conflict resolution, and  $f$  - defines the times of occurrence of the event,  $f: T \rightarrow \mathbb{R}$ , where  $\mathbb{R}$  - denotes the set of non-negative real numbers. The duration time of events can be deterministic, specified by the  $f(t)$  value or stochastic described by a corresponding function of probability of density with the  $f(t)$  parameter. In the case of distributions described by more parameters the value of the  $f$  function should be treated as vectors of the appropriate values. Performance evaluation of the model using a temporal network is well described in the literature [17, 22].

### 3 The Model

Several versions of Petri nets have been used as models of systems which exhibit concurrent and parallel activities. Stochastic Petri nets and timed Petri nets have many similarities but deal with temporal properties of models in a different way, so sometimes the similarities may be misleading. The basic model used here is known as an inhibitor Petri net  $N$ , which is a bipartite directed graph  $N = (P, T, A, H)$  where  $P$  and  $T$  are two disjoint sets of vertices  $P \cap T = \emptyset$ , called places and transitions, respectively,  $A$  is a set of directed arcs connecting places with transitions and transitions with places,  $A \subseteq P \times T \cup T \times P$ , and  $H$  is a set of inhibitor arcs connecting places with transitions,  $H \subset P \times T$ . Normally,  $A \cap H = \emptyset$ . For performance analysis of Petri net models, temporal characteristics of occurring transitions must be taken into account. This can be done in several ways assigning occurrence time to places, or to transitions, or even to arcs of the Petri net models. In timed nets, the occurrence times of some transitions may be equal to zero, which means that such occurrences are instantaneous, all such transitions are called immediate, while the others are called timed. It should be noted that such a convention effectively introduces the priority of immediate transitions over the timed ones, so the conflicts of immediate and timed transitions are not allowed in timed nets. To evaluate the performance of the modeled systems used was a Petri net model shown

schematically in Fig. 1. Based on the prepared models studied were mechanisms of traffic shaping in systems based on priority queues. In the studied models, have been made assumptions that the PQ model will consist of three priority queues (places  $p_1, p_2, p_3$ ). The three places,  $p_1, p_2$  and  $p_3$ , are queuing for class-1, class-2 and class-3 packets, respectively. It is assumed that all queues have infinite capacities. This will allow the display of traffic data with high, medium and low priority. In the studied models, data has been denoted as class-1 (high priority), class-2 (medium priority), class-3 (low priority). Timed transitions  $t_1, t_2$  and  $t_3$  represent a transmission channel for class-1 packets ( $t_1$ ), class-2 packets ( $t_2$ ) and class-3 packets ( $t_3$ ). In the studied model, for simplicity, has been assumed that the transmission time is deterministic and is 1 time unit (for all three classes). In the studied model, place  $p_5$  is shared by the transitions  $t_1, t_2$  and  $t_3$ . In order to ensure that only one packet is transmitted at the moment, place  $p_5$  has been marked with a single token (i.e.,  $m_0(p_5) = 1$ ). Furthermore, inhibitor arc ( $p_1, t_2$ )

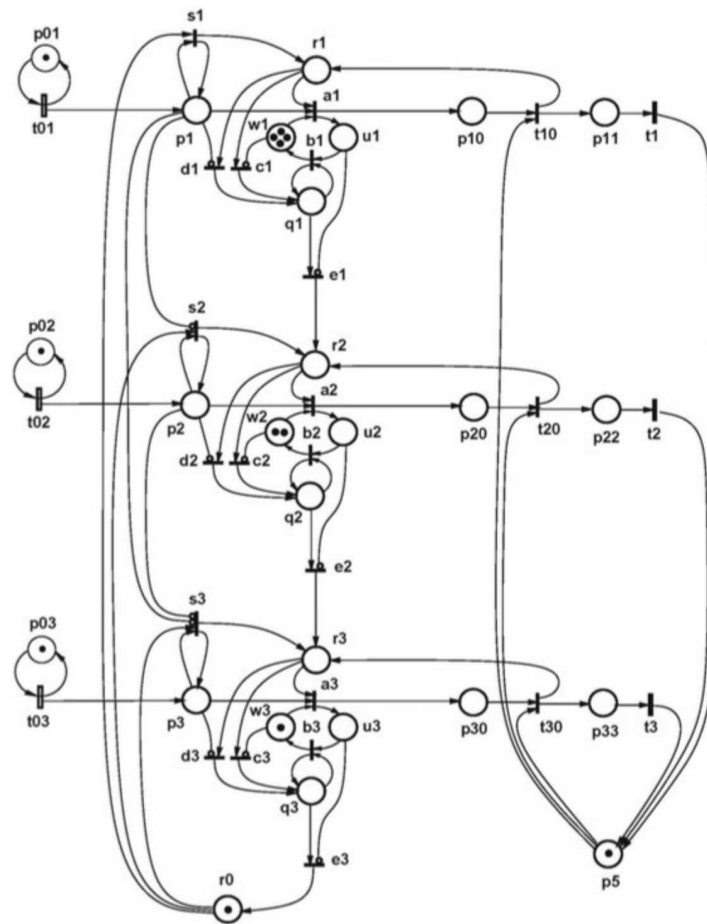


Fig. 1. Petri net model of priority queueing system

does not allow transmission of data with the transit  $t_2$  when the package is ready to be transmitted from class-1. Similarly, inhibitor arc  $(p_1, t_3$  and  $p_2, t_3)$  allows transmissions only when there is no data type class-1 and class-2 waiting for transmission. A characteristic feature of priority queues is the fact that if any packets are in the queue with a higher priority packets waiting in queues with lower priority cannot be sent. Data transfer only the queues with higher priority can lead to “starvation” of the data classified into lower priority queues. The phenomenon of starvation data (*starvation problem*) can be eliminated by introducing configurable queuing CQ (*Custom Queuing*). This mechanism allows to handle the queues on the basis of round-robin by downloading the first  $x$  bytes from the first queue, then  $y$  bytes from the second queue and  $z$  bytes from the third queue. This design prevents the extinction of less privileged streams. In the studied models, assumptions have been made that the model X421 will support successively 4 packages of type class-1, then 2 packages of type class-2, and finally 1 package of type class-3 and the model X432 will support successively 4 packages of type class-1, then 3 packages of type class-2, and 2 package of type class-3. Based on so prepared models were studied mechanisms of traffic shaping in systems based on Priority Queuing with the priorities of 4-2-1 and 4-3-2. As a source of data generation in all three queues, used were M-timed Petri nets (Markovian nets).

#### 4 Results and Discussion

Performance results, obtained by discrete-event simulation of the timed Petri net models, are compared for two cases to show the influence of different weights on the performance characteristics. Both cases use three classes of traffic with weights equal to 4, 2 and 1 in one case (denoted X421) and equal to 4, 3 and 2 in the second case (denoted X432). Fig. 2 shows average waiting times of packets when the packets arrive with the same rates in all three classes. When the traffic intensity approaches 1, the waiting times become significant, especially for class-3 packets as this class receives the least amount of service (according to the weights). Also, the differences of waiting times are less pronounced for the case X432 as the amounts of service for different classes are more similar in the X432 case than in X421.

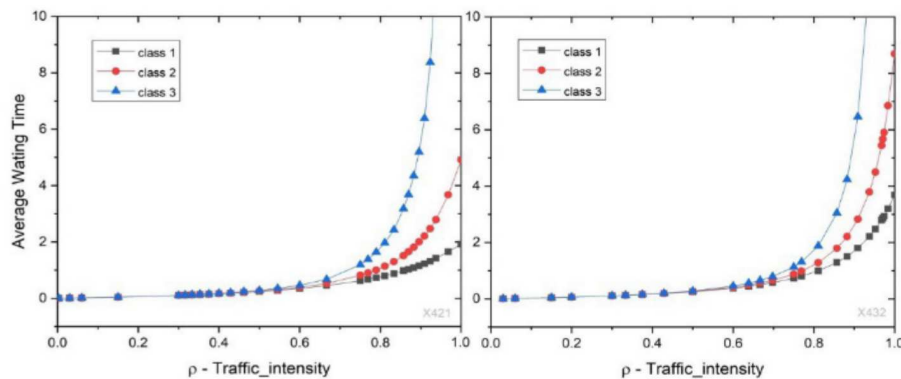
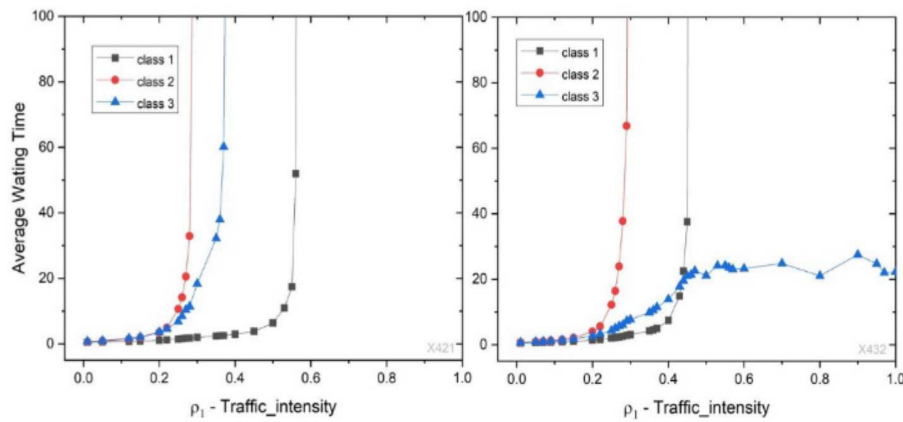


Fig. 2. Average waiting times for X421 model (left) and X432 model (right)

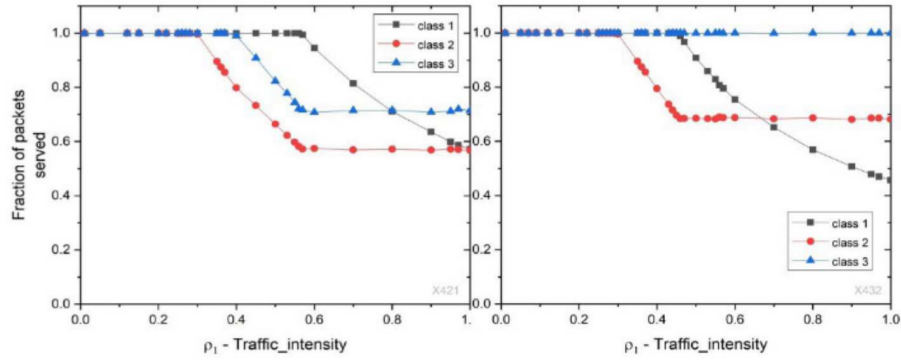


A different characterization of waiting times is shown in Fig. 3, in which the traffic in classes-2 and 3 are constant at the levels of traffic intensity equal to 0.5 for class-2 ( $\rho_2 = 0.5$ ) and 0.2 for class-3 ( $\rho_3 = 0.2$ ), and with traffic intensity changing from 0 to 1 for class-1. It should be observed that the constant traffic in classes 2 and 3 for X421 are at the levels above the levels guaranteed by the weights (these levels are  $2/7 \approx 0.287$  for class-2 and  $1/7 \approx 0.143$  for class-3) while for X432 the guaranteed levels are  $1/3$  for class-2 and  $2/9 \approx 0.222$  for class-3, so the constant traffic is below the guaranteed level for class-3.

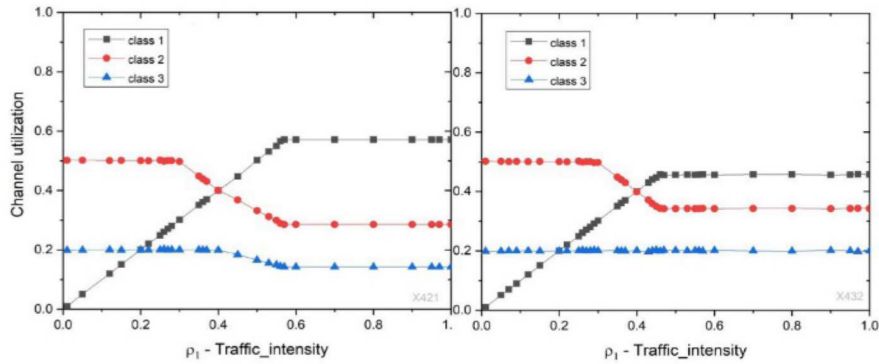


**Fig. 3.** Average waiting times for X421 model (left) and X432 model (right) as a function of  $\rho_1$  with  $\rho_2 = 0.5$  and  $\rho_3 = 0.2$

In Fig.3, as the traffic intensity in class-1 increases, initially the increasing traffic in class-1 uses the remaining capacity of the channel (the constant traffic of classes-2 and 3 use 70 % of this capacity). For  $\rho_1 = 0.3$  the channel is used in 100 %, so any further increase in class-1 traffic is possible only at the expense of traffic in the other classes. For  $\rho_1 > 0.3$ , traffic in class-2 becomes partially blocked and some incoming class-2 packets cannot be forward through the channel (i.e., class-2 becomes nonstationary when the queue is infinite). For X421, when  $\rho_1 > 0.4$ , traffic in class-3 also becomes partially blocked and class-3 becomes nonstationary. Traffic in classes-2 and 3 is gradually reduced to the guaranteed levels, and for  $\rho_1 > 4/7$ , class-1 also becomes nonstationary. For X432, class-2 becomes nonstationary also for  $\rho_1 > 0.3$ , however, class-3 is stationary in the whole range of  $\rho_1$  values. Class-1 becomes nonstationary for  $\rho_1 > 1 - 0.2 - 1/3 \approx 0.467$ . Fig.4 is a complementary illustration to Fig.3. It shows, for each class, the fraction of the total number of incoming packets that are forwarded through the channel. For X421, for  $\rho_1 > 0.3$ , some class-2 packets are blocked (by class-1 traffic), and for  $\rho_1 > 0.4$ , some of class-3 packets cannot be forwarded either. Finally, for  $\rho_1 > 4/7$ , class-1 becomes nonstationary.



**Fig. 4.** Fraction of packet served for X421 model (left) and X432 model (right) as a function of  $\rho_1$  with  $\rho_2 = 0.5$  and  $\rho_3 = 0.2$



**Fig. 5.** Channel utilization for X421 model (left) and X432 model (right) as a function of  $\rho_1$  with  $\rho_2 = 0.5$  and  $\rho_3 = 0.2$

For X432, class-3 traffic is not blocked by traffic in other classes. Class-2 traffic is partially blocked by class-1 traffic for  $\rho_1 > 0.3$ . Class-1 becomes nonstationary for  $\rho_1 > 0.467$ . Fig.5. is yet another illustration of the interplay of different classes of traffic. It shows how the channel is used by the classes of traffic as a function of traffic intensity of class-1. The blocking of traffic in classes-2 and 3 for X421 and in class-2 for X432 is manifested by the reduction of channel utilization to the levels implied by the weights.

## 5 Summary

Through telecommunications networks, including the Internet pass huge amounts of information. We may relatively easily calculate the transfer time from node to node, knowing the distance between the nodes and the transmission speed of the link. Unknown is, however, time and nature of expectancy in the node.

On the one hand, it is important to make the best use of network resources, and on the other an increase in network usage worsens the quality of service (growing queue of nodes, increasing likelihood of overflow). Hence the planning and development of computer networks can be supported using mathematical modeling and computer simulation. The study demonstrated that the use of models of temporal Petri nets can be used to evaluate the performance and effectiveness of the Priority Queuing Systems. It should be noted that the models can often be very complex, and analyze them without a flexible and robust software tools is almost impossible. Although there are many tools for modeling and data analysis with Petri nets, only a few of them can be used in the study of complex systems. It was shown that hierarchical modeling in which we consider some parts of the model at a very detailed level and the other portions on a more general level can be easily implemented based on the Petri net models. Model or its parts can easily undergo further and detailed analysis.

## References

1. Buchholz, P., Kriege, J., Felko, I.: Input modeling with phase-type distributions and Markov models: theory and applications. Springer(2014)
2. Czachórski, T., Domański, A., Domańska, J., Rataj, A.: A study of IP router queues with the use of Markov models. In: International Conference on Computer Networks. Springer International Publishing (2016)
3. Donthi, R., Renikunta, R., Dasari, R., Perati, M.: Study of delay and loss behavior of internet switch-markovian modelling using Circulant Markov Modulated Poisson Process (CMMPP). Appl. Math. **5**(3), 512–519 (2014)
4. Puterman, M., Markov, L.: Decision Processes: Discrete Stochastic Dynamic Programming. Wiley, New York (2014)
5. Bose, S.K.: An Introduction to Queueing Systems. Springer Science & Business Media, New York (2013)
6. Bhatti, S.N., Crowcroft, J.: QoS-sensitive flows: issues in IP packet handling. IEEE Internet Comput. **4**(4), 48–57 (2000)
7. Mishkoy, G., et al.: Priority queueing systems with switchover times: generalized models for QoS and CoS network technologies and analysis. In: The XIV Conference on Applied and Industrial Mathematics, Satellite Conference of ICM, Chisinau (2006)
8. Aurrecochea, C., Campbell, A.T., Hauw, L.: A survey of QoS architectures. Multimedia Syst. **6**(3), 138–151 (1998)
9. Carmona-Murillo, J., et al.: QoS in next generation mobile networks: an analytical study. In: Resource Management in Mobile Computing Environments, pp. 25–41. Springer International Publishing (2014)
10. Strzeciwiłk, D.: Examination of transmission quality in the IP multi-protocol label switching corporate networks. Int. J. Electron. Telecommun. **58**(3), 267–272 (2012)
11. Kim, H.J., Choi, S.G.: A study on a QoS/QoE correlation model for QoE evaluation on IPTV service. In: The 12th International Conference on Advanced Communication Technology (ICACT), vol. 2. IEEE (2010)
12. Li, M.: Queueing Analysis of Unicast IPTV with User Mobility and Adaptive Modulation and Coding in Wireless Cellular Networks. arXiv preprint [arXiv:1511.01794](https://arxiv.org/abs/1511.01794) (2015)
13. Zhang, Q., et al.: Early drop scheme for providing absolute QoS differentiation in optical burst-switched networks. In: Workshop on High Performance Switching and Routing, HPSR. IEEE (2003)
14. Zuberek, W., Strzeciwiłk, D.: Modeling Quality of Service Techniques for Packet-Switched Networks. Dependability Engineering (2018). ISBN 978-953-51-5592-8



15. Tarasiuk, H., et al.: Performance evaluation of signaling in the IP QoS system. *J. Telecommun. Inf. Technol.*, 12–20 (2011)
16. Menth, M., Briscoe, B., Tsou, T.: Precongestion notification: new QoS support for differentiated services IP networks. *IEEE Commun. Mag.* **50**(3), 94–103 (2012)
17. Zuberek, W.M.: Timed Petri nets definitions, properties, and applications. *Microelectron. Reliab.* **31**(4), 627–644 (1991)
18. Gianfranco, B.: Introduction to stochastic Petri nets. *Lectures on Formal Methods and Performance Analysis*, pp. 84–155. Springer, Heidelberg (2001)
19. Roux, O.H., Déplanche, A.M.: A t-time Petri net extension for real time-task scheduling modeling. *Eur. J. Autom.* **36**(7), 973–987 (2002)
20. Coolahan, J.E., Roussopoulos, N.: Timing requirements for time-driven systems using augmented Petri nets. *IEEE Trans. Softw. Eng.* **5**, 603–616 (1983)
21. Cheng, A.: *Real-Time Systems: Scheduling, Analysis, and Verification*. Wiley (2003)
22. Strzeciwlk, D., Zuberek, W.M.: Modeling and performance analysis of QoS data. In: *Photonics Applications in Astronomy, Communications, Industry, and High-Energy Physics Experiments 2016*. International Society for Optics and Photonics (2016)